

# Representations of $E_8$ and other algebras

Richard H. Capps

Department of Physics, Purdue University, West Lafayette, Indiana 47907

(Received 14 August 1986; accepted for publication 18 February 1987)

A simple procedure is given for determining whether or not an arbitrary weight of a simple Lie algebra is contained in an arbitrary irreducible representation. The procedure involves two steps, determining the Weyl class of the weight and determining the classes contained in the representation. The second step applies without alteration to all simple algebras, while the first step is given here only for  $E_8$ . The weights of the shortest 31 Weyl classes of  $E_8$  are listed in a convenient, orthogonal basis.

## I. INTRODUCTION

In the application of finite, simple Lie algebras a common problem is to find whether or not a weight  $M$  is contained in an irrep (irreducible representation)  $\Lambda$ . Dynkin developed a constructive method that may be used to solve this problem for any  $M$  and  $\Lambda$  (Ref. 1). Descriptions of the procedure are given in various reviews and texts.<sup>2,3</sup> Unfortunately, the method is cumbersome for all but the smallest representations, and the number of steps involved increases without limit as one considers very large irreps. One purpose of this paper is to present an alternate method involving a number of steps that is small for all irreps.

The new method is applied here to the algebra  $E_8$ . Recent developments in string theory suggest that this algebra is important in particle physics.<sup>4</sup> The properties of the smallest nontrivial irrep (the adjoint) have been discussed in detail in the literature. However, the other irreps are still unfamiliar. The main reason for this is that the irreps are very large; for example, the sixth smallest irrep has dimension 147 250. A second purpose of this paper is to provide a tractable treatment of  $E_8$  irreps. The weights of the shortest 31 irreps are given in a convenient basis. (The length of an irrep is defined to be the length of the longest weights in the irrep.)

The method involves the Weyl reflection group. The weights that may be obtained from a given weight by a series of Weyl reflections are said to be in the same Weyl class (or Weyl orbit). The Weyl class is a convenient concept, for three well-known reasons. First, each weight is in one and only one Weyl class. Second, the multiplicity of a weight in an irrep is the same for all weights in a class. Third, the sizes of all Weyl classes for an algebra are bounded by the finite size of the Weyl group.

The method involves two separate algorithms. The first, the "weight algorithm," is for finding the Weyl class of an arbitrary weight. The second, the "representation algorithm" is for determining whether or not a given class is contained in a given irrep. The representation algorithm is discussed first, because the method used is general and may be applied to any finite simple Lie algebra. The weight algorithm involves the choice of a convenient basis. Although the method for selecting a basis is general, the details of application depend on the basis and thus are different for different algebras. Identifying the class of a weight is nontrivial only for the algebras  $G_2$ ,  $F_4$ ,  $E_6$ ,  $E_7$ , and  $E_8$ . The algebra  $E_8$  is treated here.

If a Weyl class is in an irrep, the multiplicity may be obtained from the Freudenthal recursion formula,<sup>2,3,5</sup> or from published tables.<sup>6</sup> Multiplicities are not given here.

The basic group-theoretical concepts and formulas that are used in the paper are listed in Sec. II. The representation algorithm is derived and applied in Sec. III. Section IV contains the weight algorithm for  $E_8$ , an iterative procedure involving either zero, one, or two iterations. The basic reasons, and a proof, that two iterations are sufficient are given in Sec. V.

## II. BASIC CONCEPTS AND FORMULAS

The standard Cartan–Weyl construction is used. If  $n$  is the rank of a simple algebra,  $n$  commuting generators, denoted by  $H_1$  to  $H_n$ , are diagonalized in each irrep. The weight vector  $M$  of a state in an irrep is a vector in an  $n$ -dimensional Euclidean space, with real components  $f_i$  given by  $H_i M = f_i M$ .

The roots are the weights of the adjoint representation. Following the Dynkin method, I label the orthogonal axes 1 through  $n$ , and define a weight as positive (or negative) if its first nonzero component is positive (or negative). A simple root is a positive root that cannot be written as a sum of two positive roots; there are  $n$  simple roots. If  $M$  is a weight in an irrep, and  $\alpha$  is a nonzero root, a fundamental equation is,<sup>7</sup>

$$\langle M, \alpha \rangle (2/\alpha^2) = p_{\alpha-} - p_{\alpha+}, \quad (2.1)$$

where the non-negative integers  $p_{\alpha-}$  and  $p_{\alpha+}$  are the maximum numbers of times the root  $\alpha$  may be subtracted from  $M$ , and may be added to  $M$ , to obtain other weights in the irrep. If  $\alpha$  is a simple root  $R_j$ , the integer  $p_{\alpha-} - p_{\alpha+}$  is the Dynkin component  $m_j$ , i.e.,

$$m_j = \langle M, R_j \rangle (2/R_j^2). \quad (2.2)$$

A dominant weight, denoted by  $M^+$ , is one for which all the Dynkin components are non-negative. Most irreps contain more than one dominant weight. However, the most positive weight in the irrep must be dominant, and is used to characterize the irrep. Each dominant weight is the most positive weight of a unique irrep.

The root-basis components and Dynkin components of a weight are denoted by capital and small letters, respectively. The root-basis components are the coefficients in an expansion in the simple roots, i.e.,  $M = \sum_i M_i R_i$ . They are related to the Dynkin component  $m_j$  by<sup>8</sup>

$$m_j = \sum_i M_i A_{ij}, \quad (2.3a)$$

$$M_i = \sum_j m_j (A^{-1})_{ji}, \quad (2.3b)$$

where the elements of the Cartan matrix  $A$  are given by

$$A_{ij} = \langle R_i, R_j \rangle (2/R_j^2). \quad (2.4)$$

A scalar product of two vectors is expressed simply if one vector is in the Dynkin basis and the other is in the root basis. The expression is

$$\langle M, W \rangle = \sum_i m_i W_i (R_i^2/2). \quad (2.5)$$

The length squared of a vector  $M$  may be determined from the formula,

$$M^2 = \sum_{ij} m_i m_j G_{ij}, \quad (2.6)$$

where  $G$  is the symmetric metric tensor, related to  $A^{-1}$  by

$$G_{ij} = (A^{-1})_{ij} (R_j^2/2). \quad (2.7)$$

For every nonzero root  $\alpha$  there is a Weyl reflection operator  $S_\alpha$ , which permutes the weights in an irrep. The action of  $S_\alpha$  on a weight  $M$  is given by<sup>3</sup>

$$S_\alpha(M) = M - (2/\alpha^2) \langle M, \alpha \rangle \alpha. \quad (2.8)$$

All the weights in a Weyl class (related by one or more Weyl reflections) are of the same length, called here the length of the class. The most positive weight in a class is dominant; all other weights are not dominant.

### III. THE REPRESENTATION ALGORITHM

The set of weights for an algebra is the set of vectors with integral Dynkin components. The representation algorithm makes use of the following "containment criterion": the weight  $m$  is in the irrep with most positive weight  $\Lambda$  if and only if the root-basis components of  $\Lambda - M^+$  are all non-negative integers, where  $M^+$  is the dominant weight of the Weyl class of  $M$ . The validity of the criterion is proved at the end of this section.<sup>9</sup>

Equation (2.3b) may be used to base a simple algorithm on this criterion. I denote  $\Lambda - M^+$  by  $\Delta$ . The rule is that the root-basis components  $\Delta_i$  must be non-negative integers, where

$$\Delta_i = \sum_j (\lambda_j - m_j^+) (A^{-1})_{ji}. \quad (3.1)$$

To my knowledge, this simple rule is not contained in the literature.

I illustrate the algorithm with an example from the algebra  $F_4$ . The  $A^{-1}$  matrix for  $F_4$  is given by

$$A^{-1} = \begin{bmatrix} 2 & 3 & 4 & 2 \\ 3 & 6 & 8 & 4 \\ 2 & 4 & 6 & 3 \\ 1 & 2 & 3 & 2 \end{bmatrix}, \quad (3.2)$$

where the root numbering convention is that of Fig. 1. We consider the two dominant weights  $B$  and  $C$ , with Dynkin components  $B = (2000)$  and  $C = (0011)$ , and ask whether or not either of the corresponding Weyl classes is the irrep of

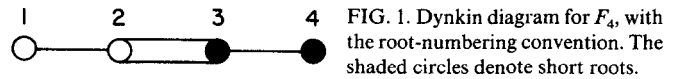


FIG. 1. Dynkin diagram for  $F_4$ , with the root-numbering convention. The shaded circles denote short roots.

the other weight. From Eqs. (3.1) and (3.2), the root basis components of  $B$  and  $C$  are

$$B = [4684], \quad C = [3695].$$

Since the difference  $B - C$  contains both positive and negative root-basis components, neither is in the irrep of the other. [The lengths of these Weyl classes, determined from Eqs. (2.6), (2.7), and (3.2), and the convention that the short and long root lengths are 1 and  $\sqrt{2}$ , are  $B^2 = 8$  and  $C^2 = 7$ .]

For some algebras the weights are in two or more congruence classes. All weights in any Weyl class, or in any irrep, are in the same congruence class. If  $M^+$  is in a different congruence class from  $\Lambda$ , the  $\Delta_i$  of Eq. (3.1) are not all integers. This follows because every nontrivial congruence relation may be associated with a column of the  $A^{-1}$  matrix in which at least one of the elements is nonintegral.<sup>10</sup> Let such a column be the  $k$ th column. Then, if  $C$  is the smallest positive integer such that the product  $C(A^{-1})_{ik}$  is an integer for each  $i$ , the congruence class of a weight with Dynkin indices  $a_i$  is

$$\sum_i a_i C(A^{-1})_{ik} \pmod{C}.$$

Clearly,  $\sum_i a_i (A^{-1})_{ik}$  is an integer only for class  $C$  (class 0). The containment criterion is valid when  $M$  and  $\Lambda$  are in different congruence classes, but is not useful in such a case, since one recognizes the incompatibility of  $M$  and  $\Lambda$  from the congruence relation.

Another useful tool is the well-known rule: if  $M^+$  and  $\Lambda$  are different dominant weights, and if the Weyl class of  $M^+$  is in the irrep  $\Lambda$ , then,

$$\Lambda^2 > (M^+)^2. \quad (3.3)$$

This rule may be proved by using Eq. (2.5), i.e.,

$$\Lambda^2 - (M^+)^2 = 2\langle M^+, \Delta \rangle + \Delta^2 = \sum_i m_i^+ \Delta_i R_i^2 + \Delta^2. \quad (3.4)$$

Since the  $m_i^+$  and  $\Delta_i$  are all non-negative, the right-hand term of Eq. (3.4) is positive, proving the rule.

I will use the algorithm to list the Weyl classes in all irreps of  $E_8$  that are not longer than  $(32)^{1/2}$ . The normalization convention is that the nonzero roots are of length  $\sqrt{2}$ . As seen from Eq. (2.7), the  $A^{-1}$  and  $G$  matrices are identical. The matrix is<sup>11</sup>

$$A^{-1} = \begin{bmatrix} 4 & 7 & 10 & 8 & 6 & 4 & 2 & 5 \\ 7 & 14 & 20 & 16 & 12 & 8 & 4 & 10 \\ 10 & 20 & 30 & 24 & 18 & 12 & 6 & 15 \\ 8 & 16 & 24 & 20 & 15 & 10 & 5 & 12 \\ 6 & 12 & 18 & 15 & 12 & 8 & 4 & 9 \\ 4 & 8 & 12 & 10 & 8 & 6 & 3 & 6 \\ 2 & 4 & 6 & 5 & 4 & 3 & 2 & 3 \\ 5 & 10 & 15 & 12 & 9 & 6 & 3 & 8 \end{bmatrix}. \quad (3.5)$$

The root-numbering convention may be obtained by delet-

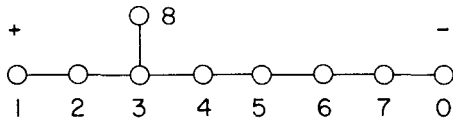


FIG. 2. Extended Dynkin diagram illustrating the  $D_8 \rightarrow E_8$  basis. The  $E_8$  simple roots are numbered 1 to 8.

ing the root 0 and its connecting line from Fig. 2. The elements of  $A^{-1}$  are all integral, so all irreps are in the same congruence class.

It is easy to use Eqs. (2.6) and (3.5) and the requirement that the Dynkin components of dominant weights are non-negative integers to determine all the Weyl classes of  $E_8$  such that  $(M^+)^2 < 32$ . These classes are listed in order of increasing length as the boldface entries in Table I. The subscripts  $a, b, c$ , and  $d$  are used to label different classes of the same length, i.e., **14a** and **14b** are the two classes of length  $(14)^{1/2}$ . The symbols in curly brackets are an abbreviation for the Dynkin components; e.g.,  $\{1^{27}\}$  denotes the Dynkin components (2000 0010).

Because of the rule of Eq. (3.3) each irrep contains its own Weyl class (class of the most positive weight), and some other shorter classes. If  $B$  and  $C$  are dominant weights in the same congruence class, and  $B^2 > C^2$ , usually  $B \supset C$  (irrep  $B$  contains class  $C$ ). Thus I define an anomalous pair of Weyl classes as two classes of the same congruence of unequal length, such that the shorter is not contained in the irrep of the longer. In Table I the classes in parentheses are the shorter anomalous partners of the preceding boldface entries. An  $x$  indicates that there are no shorter anomalous partners.

I discuss briefly the calculation of the table entries. The easiest procedure is to consider the irreps in order of increasing length, checking each for shorter anomalous partners. One can use the fact, obvious from the contents criterion, that if  $B \supset C$  and  $C \supset D$ , then  $B \supset D$ . I illustrate by considering the irrep **16a**, with Dynkin components  $\{1^2\}$ . I assume that the shorter irreps have all been examined and found to contain no shorter anomalous partners. The root-basis components of  $\{1^2\}$  are twice the elements in the first row of  $A^{-1}$ , i.e.,

$$\mathbf{16a}\{1^2\} = [8 \ 14 \ 20 \ 16 \ 12 \ 8 \ 4 \ 10]. \quad (3.6)$$

One uses Eq. (2.3b) to find the corresponding root-basis components of the dominant weights immediately preceding **16a** in Table I. These are

$$\mathbf{14a}\{2\} = [7 \ 14 \ 20 \ 16 \ 12 \ 8 \ 4 \ 10], \quad (3.7a)$$

$$\mathbf{14b}\{67\} = [6 \ 12 \ 18 \ 15 \ 12 \ 9 \ 5 \ 9]. \quad (3.7b)$$

TABLE I. Anomalous pairs of  $E_8$  Weyl classes with  $M_2 < 32$ .

<b>0</b>	$\{0\}x$ ;	<b>2</b>	$\{7\}x$ ;	<b>4</b>	$\{1\}x$ ;	<b>6</b>	$\{6\}x$ ;	<b>8a</b>	$\{8\}x$ ;	<b>8b</b>	$\{7^2\}x$ ;
<b>10</b>	$\{17\}x$ ;	<b>12</b>	$\{5\}x$ ;	<b>14a</b>	$\{2\}x$ ;	<b>14b</b>	$\{67\}x$ ;	<b>16a</b>	$\{1^2\}(14b)$ ;		
<b>16b</b>	$\{78\}x$ ;	<b>18a</b>	$\{16\}x$ ;	<b>18b</b>	$\{7^3\}(14a,16a,16b)$ ;	<b>20a</b>	$\{17^2\}x$ ;				
<b>20b</b>	$\{4\}(18b)$ ;	<b>22a</b>	$\{18\}(18b,20a)$ ;	<b>22b</b>	$\{57\}x$ ;	<b>24a</b>	$\{27\}x$ ;				
<b>24b</b>	$\{6^2\}(22a)$ ;	<b>26a</b>	$\{1^27\}(24b)$ ;	<b>26b</b>	$\{68\}x$ ;	<b>26c</b>	$\{67^2\}(22a,24a)$ ;				
<b>28a</b>	$\{15\}(26c)$ ;	<b>28b</b>	$\{7^28\}(26a)$ ;	<b>30a</b>	$\{167\}x$ ;	<b>30b</b>	$\{3\}(26c,28b)$ ;				
<b>32a</b>	$\{12\}(26c,28b,30a)$ ;	<b>32b</b>	$\{47\}x$ ;	<b>32c</b>	$\{8^2\}(26c,28b,30a)$ ;						
<b>32d</b>	$\{7^4\}(22a,24a,26a,26b,28b,30a,30b)$ ;										

Comparing Eq. (3.6) with Eqs. (3.7a) and (3.7b), one sees that the pair **16a** and **14b** is anomalous, but **14a** is contained in **16a**. Since all shorter classes are contained in **14a**, one concludes that **14b** is the only shorter anomalous partner of **16a**. Thus the irrep **16a** contains the class **16a** and all shorter classes except **14b**. We note that the root-basis difference between Eqs. (3.7a) and (3.7b) contains elements of both signs, but this follows automatically from the fact that the two classes are of the same length.

Finally, I will prove the validity of the containment criterion, stated at the beginning of the section. One-half of the proof is obvious. Since the standard procedure for finding weights in an irrep involves subtracting simple roots from  $\Lambda$ , it is clear that if the root-basis components of  $\Delta = \Lambda - M^+$  are not all non-negative integers,  $M^+$  (and hence  $M$ ) is not in  $\Lambda$ . Hence we turn to the case where the  $\Delta_i$  are all non-negative integers.

We neglect the trivial case  $\Delta = 0$ . Then  $\Delta$  is a positive weight. It follows that at least one of its Dynkin components must be positive, since weights with all nonpositive Dynkin components are the most negative weights of irreps, and so are not positive. Let  $\delta_k$  be a positive Dynkin component. Since the  $\Delta_j$  are non-negative and the off-diagonal elements of the Cartan matrix are nonpositive, it follows from the transformation rule of Eq. (2.3a) that  $\Delta_k > 0$ . Since  $\Delta_k$  is integral,

$$\Delta_k \geq 1. \quad (3.8)$$

The Dynkin component  $\lambda_k$  of  $\Lambda$  satisfies the equation,  $\lambda_k = \delta_k + m_k^+$ . Since  $m_k^+$  is non-negative and  $\delta_k$  is positive,  $\lambda_k \geq 1$ . It follows from Eq. (2.1) that one may subtract the simple root  $R_k$  from  $\Lambda$ , i.e., the weight  $\Lambda' = \Lambda - R_k$  is in the irrep.

One next considers the pair of weights  $\Lambda'$  and  $M^+$ . The new difference weight  $\Delta' = \Lambda' - M^+$  is related to  $\Delta$  by the equation,

$$\Delta'_i = \Delta_i - \delta_{ik}.$$

It follows from Eq. (3.8) that the root-basis components  $\Delta'_i$  are all non-negative integers, so that process may be repeated, until a path is traced from  $\Lambda$  to  $M^+$ , proving  $M^+$  is in  $\Lambda$ . This completes the proof.

## IV. THE WEIGHT ALGORITHM

### A. The problem

In this section I discuss the problem of finding the  $E_8$  Weyl class of an arbitrary weight, expressed in Dynkin components. There exists a standard, finite procedure for finding the answer. One makes a series of Weyl reflections  $S_R$ , always choosing  $R$  to be a simple root corresponding to a negative Dynkin component, so that the reflection leads to a more positive weight. Such a series is called a "Dynkin reflection series" in this paper. It is easy to make each reflection; if the  $j$ th Dynkin component is  $-b$  (where  $b$  is positive) the  $S_j$  reflection adds  $b$  times the simple root  $R_j$ . Unfortunately, the number of reflections necessary to obtain a dominant weight may be as large as the number of positive roots, 120 in the case of  $E_8$ . Clearly, a shorter procedure is needed.

## B. The $D_8$ basis

For each of the classical algebras a basis exists in which the Weyl class of a weight is obvious. Therefore it is convenient to express the  $E_8$  weights in a classical-algebra basis. Such a basis has the additional advantage of giving a clearer picture of the  $E_8$  weights. A convenient basis for  $E_8$  is based on  $D_8$ . I will define this basis here, and apply it to  $E_8$  in part C of this section.

In a standard, orthogonal  $D_8$  basis each nonzero root has two components of magnitude 1, the other components being zero.<sup>12</sup> A shorthand notation is used for the roots in this basis: i.e.,  $3_+5_-$  denotes [0010 - 1000]. The simple roots are  $1_+2_-$ ,  $2_+3_-$ ,  $3_+4_-$ ,  $4_+5_-$ ,  $5_+6_-$ ,  $6_+7_-$ ,  $7_+8_-$ , and  $7_+8_+$ . The Dynkin diagram is obtained by deleting the root labeled  $+$  and its connecting line from Fig. 2. The simple roots listed above correspond to the circles numbered 0, 7, 6, 5, 4, 3, 2, and 8, respectively.

The components of a weight in the orthogonal basis are denoted by  $f_i$ . It is seen from Eq. (2.2) that the set of weights for  $D_8$  is the set of vectors that have integral scalar products with all roots. The weights are of two types. The vector weights are all vectors such that each component  $f_i$  is integral. The spinor weights are all vectors such that each component is half-odd-integral. Both the vector and spinor weights may be classified as even or odd, according to whether the component sum  $\sum_{i=1}^8 f_i$  is even or odd. These four classes are the four congruence classes of  $D_8$ . It is clear that adding a root to a weight leads to another weight of the same congruence class.

Next we examine the structure of the Weyl classes in the orthogonal basis. It follows from Eq. (2.8) that the Weyl reflection corresponding to the root  $j_+k_-$  (or  $j_-k_+$ ) interchanges the  $j$  and  $k$  components of a weight. The reflection corresponding to the root  $j_+k_+$  (or  $j_-k_-$ ) interchanges the  $j$  and  $k$  components and changes both their signs. It is convenient to define the signature of a weight as zero, positive, or negative, corresponding to the component product  $\prod_i f_i$  being zero, positive, or negative. Weyl reflections preserve the signature.

It is clear from the above discussion that if one is given a  $D_8$  weight in the orthogonal basis, he may determine the dominant weight of the Weyl class immediately. The components are arranged in order of decreasing magnitude, and the signs of the first seven components are chosen to be positive. If the eighth component is nonzero, its sign is chosen to be the signature of the weight. A weight is dominant if and only if it satisfies the conditions

$$|f_i| \geq |f_{i+1}|, \quad (4.1a)$$

$$f_i \geq 0 \quad \text{for } i < 8. \quad (4.1b)$$

## C. Application to $E_8$

An  $E_8$  basis based on  $D_8$  may be found by the replacement prescription of Ref. 10. One writes the simple roots of  $D_8$  and adds to this set the most negative weight of a  $D_8$  irrep, such that the weight length is  $\sqrt{2}$ . If the irrep is the odd-signature fundamental spinor, the new root is [ - - - - - - - + ] (where the magnitudes are all  $\frac{1}{2}$ ). One then discards the  $D_8$  root  $1_+2_-$  and reflects the first axis, so that

the eight remaining roots are positive. The new root set is the simple root set of  $E_8$ . Ordered according to the numbers 1 to 8 of Fig. 2, the  $E_8$  simple roots are

$$\begin{aligned} & [ + - - - - - - + ], \\ & 7_+8_-, \quad 6_+7_-, \quad 5_+6_-, \\ & 4_+5_-, \quad 3_+4_-, \quad 2_+3_-, \quad \text{and } 7_+8_+. \end{aligned} \quad (4.2)$$

The replacement procedure is useful because it allows one to consider all bases of a certain type. However, it has been known for years that a convenient  $E_8$  basis may be obtained by taking as  $E_8$  roots the  $D_8$  roots plus the states of the even-signature fundamental spinor.<sup>13</sup> With the usual positivity definition (Sec. II) the simple roots of this set are those of Eq. (4.2).

Since only one of the  $E_8$  simple roots (the spinor) is not a  $D_8$  root, one may determine the  $E_8$  weight set and Weyl classes by considering all the  $D_8$  roots and one spinor. It is simplest to use the spinor

$$\chi^+ = [ + + + + + + + + ],$$

obtainable by combining [ + - - - - - - + ] with  $D_8$  roots. The weight set of  $E_8$  includes all  $D_8$  weights  $W$  such that  $\langle W, \chi^+ \rangle$  is an integer. This includes only two of the four  $D_8$  congruence classes, the vectors of even  $\sum_i f_i$  and the spinors of even  $\sum_i f_i$ . Since the spinor  $E_8$  root connects these classes, there is only one  $E_8$  congruence class.

It follows from Eq. (2.2) and the root set of Eq. (4.2) that the  $E_8$  Dynkin components of a weight are given in terms of their components in the orthogonal basis by

$$a_1 = \frac{1}{2}(f_1 - f_2 - f_3 - f_4 - f_5 - f_6 - f_7 + f_8), \quad (4.3a)$$

$$a_2 = f_7 - f_8, \quad (4.3b)$$

$$a_3 = f_6 - f_7, \quad a_4 = f_5 - f_6, \quad (4.3c)$$

$$a_5 = f_4 - f_5, \quad a_6 = f_3 - f_4, \quad (4.3d)$$

$$a_7 = f_2 - f_3, \quad a_8 = f_7 + f_8. \quad (4.3e)$$

Frequently, one needs to determine the  $f$ 's from the  $a$ 's. The easiest procedure to follow is to use the inverse equations for  $f_1, f_7$ , and  $f_8$ ,

$$f_1 = 2a_1 + \frac{3}{2}a_2 + 5a_3 + 4a_4 + 3a_5 + 2a_6 + a_7 + \frac{5}{2}a_8, \quad (4.4a)$$

$$f_8 = \frac{1}{2}(a_8 - a_2), \quad f_7 = \frac{1}{2}(a_8 + a_2), \quad (4.4b)$$

and then use Eqs. (4.3c)-(4.3e) to determine in order  $f_6$  through  $f_2$  (i.e.,  $f_6 = f_7 + a_3$ , etc.) In order to simplify the determination of dominant weights in the orthogonal basis, I list below the eight fundamental  $E_8$  weights in this basis;  $\{k\}$  denotes the fundamental weight with Dynkin components  $a_i = \delta_{ik}$

$$\begin{aligned} \{1\} & 2, & \{5\} & 3111, \\ \{2\} & \frac{1}{2}(71111111 - 1), & \{6\} & 211, \\ \{3\} & 511111, & \{7\} & 11, \\ \{4\} & 41111, & \{8\} & \frac{1}{2}(51111111), \end{aligned} \quad (4.5)$$

where two or more zeros at the end of a weight are omitted, and the constant  $\frac{1}{2}$  is factored out of the spinor weights.

The  $E_8$  root set includes all the  $D_8$  roots, so each  $E_8$  Weyl class is a sum of complete  $D_8$  Weyl classes. The  $D_8$  classes in an  $E_8$  class are those connected by the reflection  $S_{\chi^+}$  associated with the spinor  $\chi^+$ . It is seen from Eq. (2.8) that



$S_{\chi^+}(f) = f'$ , where

$$f'_i = f_i - \frac{1}{4} \sum_{j=1}^8 f_j. \quad (4.6)$$

When one uses Eq. (4.6) to determine the  $D_8$  classes associated with a dominant class (class of the dominant  $E_8$  weight) he does not have to worry about the order of the components, but only has to consider different combinations of the signs of the  $f_i$ .

I will illustrate the technique by determining the  $D_8$  classes in the  $E_8$  Weyl class  $\{2\}$ , with  $M^2 = 14$ . From Eq.

(4.5) the dominant weight is  $[\frac{1}{2}(7111\ 111 - 1)]$ . From Eq. (4.6), the result of applying the reflection  $S_{\chi^+}$  to this weight is  $[2 - 1 - 1 - 1 - 1 - 1 - 1 - 2]$ . The dominant weight of this  $D_8$  Weyl class is  $[2211\ 111 - 1]$ . One next changes two signs in the original dominant weight, obtaining  $[\frac{1}{2}(7111\ 1 - 1 - 1 - 1)]$ . Reflecting by  $S_{\chi^+}$ , one obtains a weight in the  $D_8$  class with dominant weight  $[\frac{1}{2}(5333\ 111 - 1)]$ . Applying Eq. (4.6) to the weight  $[\frac{1}{2}(711 - 1 - 1 - 1 - 1 - 1)]$  leads to the  $D_8$  Weyl class  $[3111\ 11]$ .

No other  $D_8$  classes may be obtained from one  $S_{\chi^+}$  re-

TABLE II. The  $D_8$  Weyl classes in the shorter  $E_8$  classes.

<b>0</b> $\{0\}$ 1 0d	<b>16b</b> $\{78\}$ $2^{10}3^5$ 32111p 3111 1111p 2221 1110p $\frac{1}{2}(7311\ 1111)d$ $\frac{1}{2}(5531\ 111 - 1)p$ $\frac{1}{2}(5333\ 3111)p$ $\frac{1}{2}(3333\ 333 - 1)p$	<b>24a</b> $\{27\}$ $2^93^5 \cdot 7$ 4211 11p 33211s 3311 111 - 1p 3222 1110p $\frac{1}{2}(9311\ 111 - 1)d$ $\frac{1}{2}(7533\ 111 - 1)p$ $\frac{1}{2}(7333\ 3311)p$ $\frac{1}{2}(5553\ 3111)s$ $\frac{1}{2}(5533\ 333 - 1)p$	<b>28b</b> $\{7^28\}$ $2^{10}3^5$ 43111p 3222 2111p 2222 2220p $\frac{1}{2}(9511\ 1111)d$ $\frac{1}{2}(7731\ 111 - 1)p$ $\frac{1}{2}(7333\ 3333)p$ $\frac{1}{2}(5553\ 3331)p$
<b>2</b> $\{7\}$ $2^43 \cdot 5$ 11d $\frac{1}{2}(1111\ 1111)p$	<b>18a</b> $\{16\}$ $2^63^45 \cdot 7$ 411d 3221s 3211 1110p 2222 11s $\frac{1}{2}(7331\ 111 - 1)p$ $\frac{1}{2}(5533\ 1111)s$ $\frac{1}{2}(5333\ 333 - 1)p$	<b>26a</b> $\{1^27\}$ $2^53^5 \cdot 7$ 51d 4211 111 - 1p 3322s $\frac{1}{2}(7333\ 333 - 1)p$ $\frac{1}{2}(5555\ 1111)s$	<b>30a</b> $\{167\}$ $2^73^45 \cdot 7$ 521d 4321s 4311 1110p 3322 1111s 3222 2210p $\frac{1}{2}(9531\ 111 - 1)p$ $\frac{1}{2}(7733\ 1111)s$ $\frac{1}{2}(7533\ 3331)p$ $\frac{1}{2}(5555\ 3311)s$
<b>4</b> $\{1\}$ $2^43^5$ 2d 1111s $\frac{1}{2}(3111\ 111 - 1)p$	<b>18b</b> $\{7^3\}$ $2^43 \cdot 5$ 33d $\frac{1}{2}(3333\ 3333)p$	<b>26b</b> $\{6^2\}$ $2^63 \cdot 5 \cdot 7$ 422d 3311 1111p 2222 22p	<b>30b</b> $\{3\}$ $2^93^5 \cdot 7$ 5111 11d 4222 11p 3331 11s 33222s 3322 111 - 1p $\frac{1}{2}(9333\ 3111)p$ $\frac{1}{2}(7553\ 311 - 1)p$ $\frac{1}{2}(5553\ 333 - 3)p$
<b>6</b> $\{6\}$ $2^63 \cdot 5 \cdot 7$ 211d 1111 11p $\frac{1}{2}(3311\ 1111)p$	<b>20a</b> $\{17^2\}$ $2^53^5 \cdot 7$ 42d 3311s 2222 1111s $\frac{1}{2}(7511\ 111 - 1)p$ $\frac{1}{2}(5333\ 3331)p$	<b>26c</b> $\{67^2\}$ $2^73 \cdot 5 \cdot 7$ 431d 2222 2211p $\frac{1}{2}(7711\ 1111)p$ $\frac{1}{2}(5533\ 3333)p$	<b>32a</b> $\{12\}$ $2^{10}3^5$ 5111 111 - 1p 4222 111 - 1p 33321s $\frac{1}{2}([11]\ 111\ 111 - 1)d$ $\frac{1}{2}(9333\ 331 - 1)p$ $\frac{1}{2}(7555\ 111 - 1)s$ $\frac{1}{2}(7533\ 333 - 3)p$
<b>8a</b> $\{8\}$ $2^73^5$ 21111p 1111 111 - 1p $\frac{1}{2}(5111\ 1111)d$ $\frac{1}{2}(3331\ 111 - 1)p$	<b>20b</b> $\{4\}$ $2^83^5 \cdot 7$ 41111d 3221 11p 22222s 2222 111 - 1p $\frac{1}{2}(7333\ 1111)p$ $\frac{1}{2}(5551\ 1111)s$ $\frac{1}{2}(5533\ 311 - 1)p$	<b>28a</b> $\{15\}$ $2^73^55^27$ 5111d 4222s 4221 1110p 3331s 3322 11s 3222 211 - 1p $\frac{1}{2}(9333\ 111 - 1)p$ $\frac{1}{2}(7553\ 1111)s$ $\frac{1}{2}(7533\ 331 - 1)p$ $\frac{1}{2}(5555\ 311 - 1)s$	<b>32b</b> $\{47\}$ $2^{10}3^5 \cdot 7$ 52111d 4321 11p 4222 1111p 3331 1111s 3322 2110p $\frac{1}{2}(9533\ 1111)p$ $\frac{1}{2}(7751\ 1111)s$ $\frac{1}{2}(7733\ 311 - 1)p$ $\frac{1}{2}(7553\ 3311)p$ $\frac{1}{2}(5555\ 5111)s$ $\frac{1}{2}(5555\ 333 - 1)p$
<b>8b</b> $\{7^2\}$ $2^43 \cdot 5$ 22d 1111 1111p	<b>22a</b> $\{18\}$ $2^{10}3^5$ 4111 1110p 32221s 3221 111 - 1p $\frac{1}{2}(9111\ 1111)d$ $\frac{1}{2}(7333\ 311 - 1)p$ $\frac{1}{2}(5553\ 111 - 1)s$ $\frac{1}{2}(5333\ 333 - 3)p$		
<b>10</b> $\{17\}$ $2^53^5 \cdot 7$ 31d 2211s 2111 1110p $\frac{1}{2}(5311\ 111 - 1)p$ $\frac{1}{2}(3333\ 1111)s$	<b>22b</b> $\{57\}$ $2^63^45 \cdot 7$ 4211d 332s 3311 11p 3221 1111p 2222 2110p $\frac{1}{2}(7531\ 1111)p$ $\frac{1}{2}(5533\ 3311)p$		
<b>12</b> $\{5\}$ $2^63^5 \cdot 7$ 3111d 222s 2211 11p $\frac{1}{2}(5331\ 1111)p$ $\frac{1}{2}(3333\ 311 - 1)p$			
<b>14a</b> $\{2\}$ $2^93^5$ 3111 11p 22211s 2211 111 - 1p $\frac{1}{2}(7111\ 111 - 1)d$ $\frac{1}{2}(5333\ 111 - 1)p$			
<b>14b</b> $\{67\}$ $2^73 \cdot 5 \cdot 7$ 321d 2211 1111p $\frac{1}{2}(5511\ 1111)p$ $\frac{1}{2}(3333\ 3311)p$			
<b>16a</b> $\{1^2\}$ $2^43^5$ 4d 3111 111 - 1p 2222s			<b>32c</b> $\{8^2\}$ $2^73^5$ 5111 1111d 42222p 3331 111 - 1p 2222 222 - 2p
			<b>32d</b> $\{7^4\}$ $2^43 \cdot 5$ 44d 2222 2222p

flection of a member of the  $\frac{1}{2}(7111\ 111 - 1)$  class. However, if  $S_{\chi^+}$  is applied to all members of the new  $D_8$  classes, one other  $D_8$  class is found, with dominant weight [22211]. (For example, the weight [1100 0 - 2 - 2 - 2] results from applying  $S_{\chi^+}$  to [2211 1 - 1 - 1 - 1].) There are five  $D_8$  classes in the  $E_8$  class {2}.

All the  $D_8$  classes in the  $E_8$  classes with  $M^2 \leq 32$  are listed in Table II. The boldface symbols and curly bracket Dynkin labels (e.g., **14a** {2}) denote the  $E_8$  classes with the same notation used in Table I. The figures after the curly brackets (i.e.,  $2^9 3^3 5$ ) are the numbers of the weights in the  $E_8$  classes, expressed in prime factors. The contained  $D_8$  classes are listed below each  $E_8$  class, with the vector classes preceding the spinor classes. The letter  $d$  denotes the  $D_8$  class containing the dominant  $E_8$  weight, while  $p$  and  $s$  (primary and secondary) label classes that require one and two  $S_{\chi^+}$  reflections, respectively, from a member of the dominant class. The fact that two  $S_{\chi^+}$  reflections are sufficient to obtain all contained classes is proved in Sec. V.

We now return to the basic problem of finding the Weyl class of an arbitrary  $E_8$  weight. We consider first a weight of length no greater than  $(32)^{1/2}$ , the weight  $\mathscr{W}$  with Dynkin components,

$$\mathscr{W} = (10 - 12 - 33 - 41). \quad (4.7)$$

One finds from Eqs. (4.4a) and (4.4b) and (4.3c)–(4.3e) that the components in the orthogonal basis are  $[\frac{1}{2}(1 - 53 - 33 - 111)]$ . The length squared is 14 and the dominant weight of the  $D_8$  class is  $[\frac{1}{2}(5333\ 111 - 1)]$ . From Table II this is the  $E_8$  class **14a** {2}. [In this case if one used a Dynkin reflection series (Sec. IV A), 47 reflections would be required.]

#### D. The long-weight procedure

If the weight length exceeds  $(32)^{1/2}$ , one cannot use Table II, so an extension of the method is needed in order to find the  $E_8$  Weyl class of the weight. One begins as before, expressing the weight in the orthogonal basis and considering the corresponding  $D$ -dominant weight (dominant weight of the  $D_8$  Weyl class). Since all  $E_8$  simple roots except  $R_1$  are  $D_8$  roots, the  $E_8$  Dynkin components  $a_i$  of the  $D$ -dominant weight are non-negative for  $i \geq 2$ . One computes  $a_1$  from Eq. (4.3a). If  $a_1$  is non-negative the weight is  $E$  dominant (a dominant  $E_8$  weight). If  $a_1$  is negative one performs the Weyl reflection  $S_1$  associated with  $R_1$ . The easiest way to make this reflection is to change the signs of the interior components (components 2 through 7), apply Eq. (4.6), and change the signs of the interior components again. Since one is interested only in the  $D_8$  Weyl class of the reflected weight, the final sign-change operation may be omitted. One considers the dominant weight of the new  $D_8$  Weyl class. If the new  $a_1$  is negative one makes a second  $S_1$  reflection. It is shown in Sec. V that two is the maximum number of  $S_1$  reflections needed to produce an  $E$ -dominant weight. One may use Eqs. (4.3a)–(4.3e) to transform the  $e$ -dominant weight to the Dynkin basis.

We consider the example of a weight in the orthogonal basis that is a member of the  $D_8$  Weyl class [7654 3210]. The length squared is 140. From Eq. (4.3a), the value of  $a_1$  is

$-7$ , so an  $S_1$  reflection is needed. One changes the signs of the interior components, yielding [7 - 6 - 5 - 4 - 3 - 2 - 1 0]. Application of Eq. (4.6) yields the weight  $[\frac{1}{2}(21 - 5 - 3 - 1\ 1\ 3\ 5\ 7)]$ . The dominant weight of this  $D_8$  class is  $[\frac{1}{2}(21\ 7\ 5\ 5\ 3\ 3\ 1 - 1)]$ . For this weight  $a_1 = -1$ , so a second  $S_1$  reflection is needed. The second reflection leads to the  $D_8$  class with dominant weight [11 3 2 2 1 1 0 0], which must be  $E$  dominant. From Eqs. (4.3a)–(4.3e), the Dynkin components are (1010 1010).

We consider one more example, a weight  $M$  of the  $D_8$  class [12 11 10 9 8 7 6 - 5]. The value of  $M^2$  is 620. The value of  $a_1$  is  $-22$ , so an  $S_1$  reflection is required. Changing the signs of the interior components and applying Eq. (4.6) yields a weight of the  $D_8$  class,

$$[23\ 6\ 5\ 4\ 3\ 2\ 1\ 0]. \quad (4.8)$$

For this weight,  $a_1 = 1$ , so the weight is  $E$  dominant. The Dynkin indices are (1111 1111). This is the shortest class with a dimension equal to that of the full Weyl group (696, 729, 600) (Ref. 14). The dimensions of some classes are discussed in Sec. V.

#### V. CONVERGENCE RATES AND ASYMPTOTIC LIMITS

The long-weight procedure of Sec. IV D is related to a Dynkin reflection series (Sec. IV A) with all reflections except  $S_1$  made automatically. However, in such a series the total number of reflections is as large as 120 in some cases. Therefore one might expect that the maximum number of  $S_1$  reflections in the long-weight procedure would be on the order of 15 (120/8), rather than two. There are three reasons that the procedure converges so fast. These are discussed below.

(1)  $R_1$  is less active than most roots: Since  $R_1$  is the end root of the second shortest branch of the  $E_8$  diagram of Fig. 2, this root is the second least active simple root of  $E_8$ . This is evidenced by the fact that the first root-basis components of most dominant weights [such as those of Eqs. (3.7a) and (3.7b)] are the second smallest components.

(2) *The  $S_1$  reflections in the long-weight procedure are made with "bottom priority"*: In order to understand this point let us construct a Dynkin reflection series for an  $SU(3)$  weight with Dynkin components  $(-a, -b)$ , where  $a$  and  $b$  are positive. The two simple roots  $R_1$  and  $R_2$  have Dynkin components  $(2, -1)$  and  $(-1, 2)$ , as seen from Eq. (2.3a) and the Cartan matrix for  $SU(3)$ . One can use either  $S_1$  or  $S_2$  for the first reflection. If  $S_1$  is chosen, the Dynkin components of the reflected weight are  $(a, -b - a)$ . The next two reflections must then be  $S_2$  and  $S_1$ , respectively, yielding successively,  $(-b, b + a)$  and  $(b, a)$ , which is dominant. Two  $S_1$  and one  $S_2$  reflections were used. Clearly, if one had applied  $S_2$  first, he would have arrived at the same place in the same number of steps, but with only one  $S_1$  reflection. In order to generalize this phenomenon I define a bottom priority, Dynkin reflection series as one in which a particular reflection  $S_k$  is made only if all Dynkin indices other than  $a_k$  are non-negative. One may minimize the number of  $S_k$  reflections by using a bottom-priority series. The long-weight procedure is related to a bottom-priority series, since the  $S_1$  reflection is made only when all Dynkin components other than  $a_1$  are non-negative.

(3) *The eighth  $D_8$  simple root is used implicitly:* In the bottom-priority series discussed above, one optimizes (reflects until the Dynkin indices are non-negative) with respect to the seven  $D_8$  roots of the  $E_8$  simple root set before and between  $S_1$  reflections. This optimization does not affect the first component in the orthogonal basis, but leads to the optimal signs and ordering of the last seven components. However, in the procedure of Sec. IV D, one orders all eight orthogonal components. This is equivalent to optimizing with respect to all eight  $D_8$  simple roots, the seven of the  $E_8$  set and also  $1_+ 2_-$ . This forces one closer to the top in positivity. The long-weight procedure is equivalent to starting fairly close to the top, and then using a bottom-priority Dynkin reflection series, with all reflections except  $S_1$  made automatically.

I illustrate the above points by considering the weight  $\mathscr{W}$  of Eq. (4.7). If one constructs a Dynkin reflection series choosing always for  $S_i$  the smallest  $i$  such that  $a_i$  is negative (top priority for  $S_1$ ) 5 of the 47 reflections needed to obtain the dominant  $E_8$  weight are  $S_1$  reflections. If one uses  $S_1$  with bottom priority, 47 reflections are still required, but only two are  $S_1$  reflections. The weight preceding the first  $S_1$  reflection is  $[\frac{1}{2}(1533\ 311 - 1)]$ . If one uses the procedure of Sec. IV D, the weight preceding the first  $S_1$  reflection is  $[\frac{1}{2}(5333\ 111 - 1)]$ , and only one  $S_1$  reflection is needed.

Next, I give the proof that the long-root procedure never requires more than two  $S_1$  reflections. It is convenient to define indices  $g_i$  by the equations,

$$g_1 = f_1, \quad g_8 = f_8, \quad g_i = -f_i \quad (2 \leq i \leq 7). \quad (5.1)$$

As seen from Eq. (4.3a) the value of  $a_i$  is given simply as

$$a_i = \frac{1}{2} \sum_{j=1}^8 g_j. \quad (5.2)$$

If  $S_1(g) = g'$ , the transformation equations are

$$g'_i = g_i - \frac{1}{2}a_i. \quad (5.3)$$

If a weight is  $D$  dominant, the conditions of Eqs. (4.1a) and (4.1b) may be written

$$|g_i| \geq |g_{i+1}|, \quad (5.4a)$$

$$g_i \geq 0, \quad g_i \leq 0, \quad \text{for } 2 \leq i \leq 7. \quad (5.4b)$$

If  $g_8 \neq 0$ , the sign of  $g_8$  is the signature of the class.

We may start the iteration procedure with a  $D$ -dominant weight. In the following argument,  $F_i$  and  $G_i$  denote the  $f_i$  and  $g_i$  values of the original  $D$ -dominant weight, while  $G'_i$  and  $G''_i$  denote the  $g_i$  values immediately after the first and second  $S_1$  reflections, respectively.

If the original  $a_1$ , determined from Eq. (5.2) with  $g_i = G_i$ , is non-negative, no  $S_1$  reflections are required, so we assume  $a_1 < 0$ . One performs an  $S_1$  reflection. If the new weight ( $G'$ ) is not  $D$  dominant, one may make it  $D$  dominant by rearranging the coefficients 2 through 8 and changing pairs of signs. (At all stages of the procedure  $g_i \geq |g_j|$ , for  $i > 1$ .) These changes correspond to the Weyl reflections  $S_2$  through  $S_8$ , associated with the  $D_8$  roots. Since the value of the new  $a_1$  does not depend on the order of the  $g_i$ , I will not reorder the  $g_i$ 's, but consider only the sign changes. A weight  $g$  (with  $g_i \geq |g_j|$ ) may be considered  $D$  dominant if either all  $g_i$  are nonpositive (for  $i > 1$ ), or exactly one of these  $g_i$  is

positive, and this  $g_i$  has a magnitude no greater than that of any other  $g_i$ .

We return to the  $S_1$ -reflected weight,  $G'$ . Because of Eqs. (5.4a), (5.4b), and (5.3),

$$G'_i \leq G_{i+1} \quad \text{for } (i > 1). \quad (5.5)$$

One may make  $G'$   $D$  dominant by one of the four following sign-change procedures: (a) *no sign changes*,  $G'$  is already  $D$  dominant; (b) *two sign changes*, one changes the signs of  $G'_7$  and  $G'_8$ ; (c) *four sign changes*, one changes the signs of  $G'_5$ ,  $G'_6$ ,  $G'_7$ , and  $G'_8$ ; or (d) *six sign changes*, one changes the signs of all  $G'_i$  except  $G'_1$  and  $G'_2$ . For example, suppose that  $G'_2, G'_3, G'_4$ , and  $G'_5$  are negative, while  $G'_6, G'_7$ , and  $G'_8$  are positive. One then changes the signs of  $G'_7$  and  $G'_8$  and, if  $G'_6 > |G'_5|$ , one also changes the signs of  $G'_5$  and  $G'_6$ .

We consider the four cases separately, and denote by  $a'_i$  the value after the appropriate sign changes are made. In case (a) it is easy to see that  $a'_1 = -a_1$ . This is positive, so  $G'$  is  $E$  dominant. The one  $S_1$  reflection is sufficient. In case (b) a short calculation yields

$$a'_1 = F_7 - F_8.$$

This is non-negative, so the  $D$ -dominant weight is also  $E$  dominant. The one  $S_1$  reflection is sufficient. In case (d), after the six sign changes are made, a calculation shows that

$$a'_1 = F_1 - F_2.$$

This is non-negative, so the one  $S_1$  reflection is sufficient.

In case (c), after the four sign changes are made,  $a'_1$  may be negative. If this occurs one performs a second  $S_1$  reflection, applying the rule of Eq. (5.3) after the sign changes are made. Because the last four  $G'_i$  have been reflected, the new  $g_2$ - $g_7$  indices form an asymmetric pyramid pattern, i.e.,

$$G''_2 \leq G''_3 \leq G''_4, \quad G''_8 \leq G''_7 \leq G''_6 \leq G''_5. \quad (5.6)$$

Again I do not reorder the columns. A calculation shows that the new  $g_3$  and  $g_6$  indices are

$$G''_3 = \frac{1}{2}[(F_2 - F_1) + (F_4 - F_3)],$$

$$G''_6 = \frac{1}{2}[(F_6 - F_5) + (F_8 - F_7)].$$

These are both nonpositive. Because of this and Eq. (5.6), the only components  $G''_i$  (besides  $G''_1$ ) that might be positive are  $G''_4$  and  $G''_5$ . It follows that if  $G''$  is not  $D$  dominant, it may be made  $D$  dominant with only one pair of sign changes. However, it was shown in case (b) above that when one pair of sign changes leads to  $D$  dominance, the resulting value of  $a_1$  is non-negative. Therefore in case (c) the second  $S_1$  reflection is sufficient. This concludes the proof.

The long-weight procedure is related closely to the procedure used in Sec. IV C to compile Table II. It is easy to show that weights requiring 0, 1, and 2  $S_1$  reflections in the long-weight procedure are those labeled dominant, primary, and secondary, respectively, in Table II.

In order to illuminate some of the procedures discussed here, I will give a short discussion of the rates at which different quantities approach their asymptotic limits as one considers Weyl classes with longer and longer weights. A few definitions are useful. If  $\mathscr{C}$  is a Weyl class of an algebra  $\mathscr{A}$ ,  $\mathscr{A}_0$  is defined as the algebra obtained by writing the Dynkin diagram for  $\mathscr{A}$  and deleting each circle (with its connecting

lines) that corresponds to a positive Dynkin component of  $\mathcal{C}$ . Thus for the class {38} of  $E_8$ ,  $\mathcal{A}_0$  is  $A_2 \times A_4$ . It is well known that the number of weights in the Weyl class  $\mathcal{C}$  is given by<sup>14</sup>

$$N(\mathcal{C}) = N(\mathcal{A})/N(\mathcal{A}_0), \quad (5.7)$$

where  $N(\mathcal{A})$  is the number of reflections in the full Weyl group for the algebra  $\mathcal{A}$ . [If the diagram for  $\mathcal{A}_0$  contains no circles,  $N(\mathcal{A}_0) = 1$ .] For convenience I list below the sizes of the Weyl groups for the  $A$ ,  $D$ , and  $E$  algebras:

$$\begin{aligned} N(A_n) &= (n+1)!, \\ N(D_n) &= n!2^{n-1}, \\ N(E_6) &= 2^7 \cdot 3^4 \cdot 5, \quad N(E_7) = 2^{10} \cdot 3^4 \cdot 5 \cdot 7, \\ N(E_8) &= 2^{14} \cdot 3^5 \cdot 5^2 \cdot 7. \end{aligned}$$

The size of an  $E_8$  Weyl class approaches its limit slowly. The sizes of all the classes in Table II are many times smaller than  $N(E_8)$ . The shortest class such that  $N(\mathcal{C}) = N(E_8)$  is the  $M^2 = 620$  class with dominant weight shown in Eq. (4.8).

Similarly, the number of  $D_8$  classes in an  $E_8$  class approaches its limit slowly. When all the  $E_8$  Dynkin components are positive, this number is  $N(E_8)/N(D_8) = 135$ . It is seen that all the classes of Table II are far from this limit.

We next consider the Dynkin reflection series of Sec. IV A. Let  $N_R(\mathcal{C})$  be the number of reflections necessary to proceed from the most negative weight of class  $\mathcal{C}$  to the dominant weight. This number is given by

$$N_R(\mathcal{C}) = P(\mathcal{A}) - P(\mathcal{A}_0), \quad (5.8)$$

where  $P(\mathcal{A})$  is the number of positive roots in the algebra  $\mathcal{A}$ , and  $\mathcal{A}_0$  is defined as before. It is seen from Eqs. (5.7) and (5.8) that  $N_R(\mathcal{C})$  is equal to its maximum value (120) only when  $N(\mathcal{C})$  is equal to its maximum value. However, the nature of Eq. (5.8) is such that even for short classes  $N_R$  is on the order of 120. Consider, for example, the  $M^2 = 6$  class {6};  $N_R\{6\} = 120 - P(E_6 \times A_1) = 83$ .

Finally, we consider the asymptotic behavior of the classification procedure of Table II and related long-root procedure.

I denote by  $F_d$ ,  $F_p$ , and  $F_s$  the fractions of the weights in a class that are dominant, primary, and secondary, in the sense of Table II. An analysis shows that if  $a_1 > 0$ ,  $F_d = \frac{1}{135}$ ,  $F_p = \frac{64}{135}$ ,  $F_s = \frac{70}{135}$ ; if  $a_1 = 0$ ,  $F_d > \frac{2}{135}$ ,  $F_s < \frac{70}{135}$ . Thus the average number of  $S_1$  reflections necessary in the long-root procedure ( $F_p + 2F_s$ ) is equal to its asymptotic value of 1.51 whenever  $a_1 > 0$ . This occurs for many short classes, such as the class with  $M^2 = 4$ . This argument shows why the long-root procedure does not get more difficult or time consuming as the weight lengths become large.

## ACKNOWLEDGMENT

This paper was supported in part by the U.S. Department of Energy.

<sup>1</sup>E. B. Dynkin, Am. Math. Soc. Transl., Ser. 1, 9, 328 (1962); Ser. 2, 6, 111 (1957).

<sup>2</sup>R. Slansky, Phys. Rep. 79, 1 (1981), Sec. 5.

<sup>3</sup>R. N. Cahn, *Semi-Simple Lie Algebras and Their Representations* (Benjamin/Cummings, Reading, MA, 1984), Chaps. X–XII.

<sup>4</sup>See, for example, D. Gross, J. Harvey, E. Martinec, and R. Rohm, Nucl. Phys. B 256, 253 (1985).

<sup>5</sup>James E. Humphreys, *Introduction to Lie Algebras and Representation Theory* (Springer, New York, 1972).

<sup>6</sup>M. R. Bremner, R. V. Moody, and J. Patera, *Tables of Dominant Weight Multiplicities for Representations of Simple Lie Algebras* (Dekker, New York, 1985).

<sup>7</sup>Howard Georgi, *Lie Algebras in Particle Physics* (Benjamin/Cummings, Reading, MA, 1982), Eq. (VI. 19).

<sup>8</sup>The formulas of Eqs. (2.3a)–(2.7) may be obtained from Ref. 2, pp. 27 and 28.

<sup>9</sup>This theorem may be obtained by combining Secs. 13.4 and 21.3 of Ref. 5 with the well-known fact that all weights in a Weyl class may be obtained by subtracting simple roots from the dominant weight.

<sup>10</sup>R. H. Capps, J. Math. Phys. 27, 914 (1986). The results concerning congruence are given on p. 920.

<sup>11</sup>See Ref. 2, Table 7.

<sup>12</sup>R. Gilmore, *Lie Groups, Lie Algebras, and Some of their Applications* (Wiley, New York, 1974), see Chap. 8.

<sup>13</sup>B. L. van der Waerden, Math. Z. 37, 446 (1933). See also Ref. 12.

<sup>14</sup>The formula for the size of Weyl class is given in Ref. 6, p. 5.

# Generalized Euler angles as intrinsic coordinates for nonlinear spinors

P. Furlan

*International School for Advanced Studies (ISAS), Dipartimento di Fisica Teorica, Università di Trieste, Trieste, Italy and Istituto Nazionale di Fisica Nucleare (INFN), Sezione di Trieste, Trieste, Italy*

(Received 24 June 1986; accepted for publication 31 December 1986)

Among the nonlinear spinor representations of the pseudo-orthogonal rotation groups presented in a preceding work, those ones are considered, whose carrier spaces are isomorphic with group spaces. The general element of such groups is written as a product of one-parameter subgroups, allowing, for the  $SO(\nu, \nu)$  and  $SO(\nu + 1, \nu)$  cases, to write explicitly the nonlinear spinor components in terms of independent intrinsic coordinates satisfying the Cartan's quadratic constraints.

## I. INTRODUCTION

In a recent work,<sup>1</sup> done in collaboration with Rączka, we have introduced a whole class of nonlinear spinor representations of complex and pseudo-orthogonal groups, showing also that they give rise to a natural generalization of the concept of pure spinors introduced by Cartan.<sup>2</sup> The distinguishing characteristic of pure spinors is that their components are not all independent but satisfy some quadratic constraints. Exploiting this fact we have been able, in a subsequent work,<sup>3</sup> to write a set of nonlinear spinor wave equations for the  $SO(\nu, \nu)$  and  $SO(\nu + 1, \nu)$  groups, which simplify in a drastic way when expressed in terms of the intrinsic spinor coordinates (which resemble a kind of "generalized Euler angles" for the considered groups). The key point in getting this result has been the possibility of explicitly writing the generic element of the complementary set  $C$ , by which the pure spinor space can be parametrized (see Refs. 1 and 3 for more details), as a product of one-parameter subgroups.

In this paper we want to show how this result has been achieved, allowing us to write explicitly all the pure spinor components in terms of the independent intrinsic ones in such a way that the quadratic Cartan relations are recovered.<sup>2</sup> We have limited our analysis to the pseudo-orthogonal groups  $SO(\nu, \nu)$ ,  $SO(\nu + 1, \nu)$ ,  $SO(\nu + 1, \nu - 1)$ , and  $SO(\nu + 2, \nu - 1)$ . In fact, from Ref. 1 it follows that these are the only (real) cases in which the carrier space of the nonlinear spinor representations may be represented—up to a set of Haar measure zero—not simply as a homogeneous space but as a group space. This carrier space is obtained by the action of the considered pseudo-orthogonal group  $G$  on the standard spinor  $\psi_m$  (chosen to be the highest-weight eigenspinor). If we denote by  $g$  the Lie algebra of the group  $G$ , by  $h$  the Lie algebra of the stability subgroup  $H$  of  $\psi_m$ , and by  $c$  the set of generators complementary to  $h$  in  $g$ , we see that, when  $G$  coincides with one of the above-mentioned groups,  $c$  is also a Lie algebra and the complementary set  $C$  is a group, obtained—up to a set of Haar measure zero—by exponentiating  $c$  (see Ref. 1 for a more detailed analysis). Therefore the basic technical problem we are confronted with here is to write explicitly  $\exp(c)$  as a product of the one-parameter subgroups obtained by exponentiating the operators of  $c$ .

In Sec. II we give a detailed analysis of the  $SO(\nu + 1, \nu)$

case, taken as a prototype, while we simply sketch the remaining cases in Sec. III. In Sec. IV for the  $SO(\nu, \nu)$  and  $SO(\nu + 1, \nu)$  cases we write the pure spinor components explicitly in terms of the intrinsic coordinates, leaving apart the remaining two cases which present major technical complications but no new interesting features. Finally in Appendix A we give the main features of the Zassenhaus formula,<sup>4</sup> on which all our calculations are based, and in Appendix B we present some cumbersome relations giving the "generalized Euler angles" of the  $SO(\nu + 2, \nu - 1)$  case. Notations and conventions are the ones introduced in Refs. 1 and 3.

## II. THE $SO(\nu+1, \nu)$ CASE

The complementary subalgebra  $c$ , which we are going to exponentiate in the form of a product of one-parameter subgroups, is a  $[1 + \binom{\nu}{2}]$ -dimensional solvable Lie algebra with

$$\binom{\nu}{2} \text{ generators } \tilde{Q}^{kl} = -\tilde{Q}^{lk}, \quad k \neq l = 1, \dots, \nu, \quad (2.1a)$$

$$\nu \text{ generators } \tilde{F}^k, \quad k = 1, \dots, \nu, \quad (2.1b)$$

$$1 \text{ generator } \tilde{D}, \quad (2.1c)$$

satisfying the following commutation relations<sup>1</sup>:

$$[\tilde{Q}^{kl}, \tilde{Q}^{mn}] = 0, \quad k, l, m, n = 1, \dots, \nu, \quad (2.2a)$$

$$[\tilde{D}, \tilde{Q}^{rs}] = 0, \quad r \neq s = 1, \dots, \nu - 1, \quad (2.2b)$$

$$[\tilde{D}, \tilde{F}^r] = 0, \quad r = 1, \dots, \nu - 1, \quad (2.2c)$$

$$[\tilde{F}^k, \tilde{Q}^{lm}] = 0, \quad k, l, m = 1, \dots, \nu, \quad (2.2d)$$

$$[\tilde{D}, \tilde{Q}^{r\nu}] = -\tilde{Q}^{r\nu}, \quad r = 1, \dots, \nu - 1, \quad (2.2e)$$

$$[\tilde{D}, \tilde{F}^\nu] = -\tilde{F}^\nu, \quad (2.2f)$$

$$[\tilde{F}^k, \tilde{F}^{l'}] = 2\tilde{Q}^{kl}, \quad k \neq l = 1, \dots, \nu. \quad (2.2g)$$

In order to get an explicit realization of this algebra we can express the generators (2.1a)–(2.1c) in terms of the basis elements  $\{H_0, H_j, H_{j'}\}, j = 1, \dots, \nu$ , of the Clifford algebra  $\mathbb{R}_{\nu+1, \nu}$ , satisfying the anticommutation relations

$$\{H_r, H_s\} = 2g_{rs} \mathbf{1}, \quad r, s = 0, 1, \dots, \nu, 1', \dots, \nu', \quad (2.3a)$$

where



$$U_2 := \sum_{r=1}^{\nu-1} f_r \tilde{F}^r; \quad (2.20b)$$

then, since from Eqs. (2.2c) and (2.2d)

$$[U_1, U_2] = 0, \quad (2.21)$$

we have

$$e^U = e^{U_1} e^{U_2} = e^{U_2} e^{U_1}. \quad (2.22)$$

Taking into account Eqs. (2.2a) and (2.2b) we get

$$e^{U_1} = \left( \prod_{r < s = 1}^{\nu-1} e^{c_{rs} \tilde{Q}^{rs}} \right) e^{\alpha \tilde{D}}, \quad (2.23a)$$

while, using Eq. (2.2g), we have

$$e^{U_2} = \left( \prod_{r=1}^{\nu-1} e^{f_r \tilde{F}^r} \right) \left( \prod_{s < t = 1}^{\nu-1} e^{-f_s f_t \tilde{Q}^{st}} \right). \quad (2.23b)$$

Inserting Eqs. (2.23a) and (2.23b) into Eq. (2.22), and taking into account Eqs. (2.2b) and (2.2c), the proposition is proved. ▼

Now we are ready to give the final form of  $\exp \tilde{X}$  in the following theorem.

**Theorem 2.1:** The general group element obtained by exponentiating the complementary subalgebra  $c$  of the Lie algebra  $\mathfrak{so}(\nu+1, \nu)$  is expressed in terms of one-parameter subgroups as follows:

$$e^{\tilde{X}} = \left[ \prod_{k=1}^{\nu} \exp(g_k \tilde{F}^k) \right] \left[ \prod_{l < m = 1}^{\nu} \exp(\tilde{d}_{lm} \tilde{Q}^{lm}) \right] e^{\alpha \tilde{D}}, \quad (2.24)$$

where

$$g_r := f_r, \quad r = 1, \dots, \nu-1, \quad (2.25a)$$

$$g_\nu := [(1 - e^{-\alpha})/\alpha] f_\nu, \quad (2.25b)$$

$$\tilde{d}_{rs} = -\tilde{d}_{sr} := c^{rs} - \epsilon(s-r) f^r f^s, \quad r \neq s = 1, \dots, \nu-1, \quad (2.25c)$$

$$\begin{aligned} \tilde{d}_{r\nu} = -\tilde{d}_{\nu r} := & [(1 - e^{-\alpha})/\alpha] c_{r\nu} \\ & + 2[(1 - \alpha - e^{-\alpha})/\alpha^2] f_r f_\nu, \quad r = 1, \dots, \nu-1. \end{aligned} \quad (2.25d)$$

*Proof:* Inserting Eqs. (2.16) and (2.18a) into Eq. (2.10) and using the relation

$$e^{\alpha \tilde{D}} \left( \frac{\tilde{Q}^{r\nu}}{\tilde{F}^\nu} \right) e^{-\alpha \tilde{D}} = e^{-\alpha} \left( \frac{\tilde{Q}^{r\nu}}{\tilde{F}^\nu} \right), \quad r = 1, \dots, \nu-1, \quad (2.26)$$

obtained from Eqs. (2.2e) and (2.2f), the theorem is proved. ▼

*Remark:* The relations (2.25a)–(2.25d) connecting the parameters  $(g_k, \tilde{d}_{lm})$  to the parameters  $(f_k, c_{lm})$  are analogs to the ones expressing the Euler angles in terms of the standard parameters of the  $\text{SO}(3)$  group elements taken as exponentials of the  $\mathfrak{so}(3)$  Lie algebra (keeping in mind the qualitative differences between the two cases).

### III. THE $\text{SO}(\nu, \nu)$ , $\text{SO}(\nu+1, \nu-1)$ , AND $\text{SO}(\nu+2, \nu-1)$ CASES

In this section we want just to sketch the main features of the remaining pseudo-orthogonal cases (in which the complementary set  $C$  has a group structure), leaving apart as much as possible all those technical details which can be recovered by analogy in the preceding section.

#### A. The $\text{SO}(\nu, \nu)$ case

This case can be considered as a subcase of the  $\text{SO}(\nu+1, \nu)$  case in which all  $\tilde{F}^k$  generators are deleted. Let us denote the general element of the complementary subalgebra  $c$  as

$$\tilde{Y} := \alpha \tilde{D} + \sum_{k < l = 1}^{\nu} c_{kl} \tilde{Q}^{kl}. \quad (3.1)$$

Then we have the following theorem.

**Theorem 3.1:** The general group element obtained by exponentiating the complementary subalgebra  $c$  of the Lie algebra  $\mathfrak{so}(\nu, \nu)$  is expressed in terms of one-parameter subgroups as follows:

$$e^{\tilde{Y}} = \left[ \prod_{k < l = 1}^{\nu} \exp(\tilde{c}_{kl} \tilde{Q}^{kl}) \right] e^{\alpha \tilde{D}}, \quad (3.2)$$

where

$$\tilde{c}_{rs} = -\tilde{c}_{sr} := c_{rs}, \quad r \neq s = 1, \dots, \nu-1, \quad (3.3a)$$

$$\tilde{c}_{r\nu} = -\tilde{c}_{\nu r} := [(1 - e^{-\alpha})/\alpha] c_{r\nu}, \quad r = 1, \dots, \nu-1. \quad (3.3b)$$

*Proof:* Put all  $f_r \equiv 0$ ,  $r = 1, \dots, \nu$ , in Eqs. (2.24), (2.25c), and (2.25d) of Theorem 2.1. ▼

#### B. The $\text{SO}(\nu+1, \nu-1)$ case

In this case the complementary subalgebra  $c$  is a  $[1 + (\nu+1)]$ -dimensional solvable Lie algebra with

$$\begin{aligned} \binom{\nu-1}{2} \text{ generators } \tilde{Q}^{kl} = -\tilde{Q}^{lk} = -H_k H_{l'}, \quad k \neq l = 2, \dots, \nu, \end{aligned} \quad (3.4a)$$

$$\begin{aligned} \nu-1 \text{ generators } \tilde{C}^{ll} = -\frac{1}{2}(H_1 + H_{1'}) H_{l'}, \\ \nu-1 \text{ generators } \tilde{D}^{ll} = \frac{1}{2}(H_1 - H_{1'}) H_{l'}, \end{aligned} \quad l = 2, \dots, \nu, \quad (3.4b)$$

$$1 \text{ generator } \tilde{B} = -\tilde{D}^{11}, \quad (3.4c)$$

$$1 \text{ generator } \tilde{B} = -(i/2)[H_1, H_{1'}], \quad (3.4d)$$

$$1 \text{ generator } \tilde{D} = \frac{1}{2}[H_\nu, H_{\nu'}], \quad (3.4e)$$

satisfying the following commutation relations:

$$[\tilde{D}, \tilde{Q}^{r\nu}] = -\tilde{Q}^{r\nu}, \quad r = 2, \dots, \nu-1, \quad (3.5a)$$

$$[\tilde{D}, \tilde{C}^{1\nu}] = -\tilde{C}^{1\nu}, \quad (3.5b)$$

$$[\tilde{D}, \tilde{D}^{1\nu}] = -\tilde{D}^{1\nu}, \quad (3.5c)$$

$$[\tilde{B}, \tilde{C}^{ll}] = \tilde{D}^{ll}, \quad (3.5d)$$

$$[\tilde{B}, \tilde{D}^{ll}] = -\tilde{C}^{ll}, \quad (3.5e)$$

$$[\tilde{C}^{lk}, \tilde{C}^{ll}] = \frac{1}{2} \tilde{Q}^{kl}, \quad k, l = 2, \dots, \nu, \quad (3.5f)$$

$$[\tilde{D}^{lk}, \tilde{D}^{ll}] = \frac{1}{2} \tilde{Q}^{kl}, \quad (3.5g)$$

all other commutators being zero.

Let us define the general element of the complementary subalgebra  $c$  as

$$\begin{aligned} \tilde{W} := & \alpha \tilde{D} + \beta \tilde{B} + \sum_{k < l = 2}^{\nu} c_{kl} \tilde{Q}^{kl} \\ & + \sum_{m=2}^{\nu} c_m \tilde{C}^{1m} + \sum_{n=2}^{\nu} d_n \tilde{D}^{1n}. \end{aligned} \quad (3.6)$$

Then we have the following theorem.

**Theorem 3.2:** The general group element, obtained by exponentiating the complementary subalgebra  $c$  of the Lie algebra  $\mathfrak{so}(\nu + 1, \nu - 1)$ , is expressed in terms of one-parameter subgroups as follows:

$$e^{\tilde{w}} = e^{\beta \tilde{B}} \left[ \prod_{k=2}^{\nu} \exp(\hat{c}_k \tilde{C}^{1k}) \right] \left[ \prod_{l=2}^{\nu} \exp(\hat{d}_l \tilde{D}^{1l}) \right] \times \left[ \prod_{m < n=2}^{\nu} \exp(\hat{c}_{mn} \tilde{Q}^{mn}) \right] e^{\alpha \tilde{D}}, \quad (3.7)$$

where

$$\hat{c}_r := \frac{1}{2} \left( \frac{\sin \beta}{\beta} c_r + \frac{1 - \cos \beta}{\beta} d_r \right), \quad r = 2, \dots, \nu - 1, \quad (3.8a)$$

$$\hat{c}_\nu := \frac{1}{\alpha^2 + \beta^2} \{ [\alpha \cos \beta + \beta \sin \beta - \alpha e^{-\alpha}] c_\nu + [\alpha \sin \beta - \beta \cos \beta + \beta e^{-\alpha}] d_\nu \}, \quad (3.8b)$$

$$\hat{d}_s := \frac{1}{2} \left( \frac{\sin \beta}{\beta} d_s - \frac{1 - \cos \beta}{\beta} c_s \right), \quad s = 2, \dots, \nu - 1, \quad (3.8c)$$

$$\hat{d}_\nu := \frac{1}{\alpha^2 + \beta^2} \{ [\alpha \cos \beta + \beta \sin \beta - \alpha e^{-\alpha}] d_\nu - [\alpha \sin \beta - \beta \cos \beta + \beta e^{-\alpha}] c_\nu \}, \quad (3.8d)$$

$$\hat{c}_{rs} = -\hat{c}_{sr} := c_{rs} + \frac{\beta - \sin \beta}{2\beta^2} (d_r c_s - d_s c_r) - \frac{1 - \cos \beta}{2\beta^2} \epsilon(s-r)(c_r c_s + d_r d_s), \quad r \neq s = 2, \dots, \nu - 1, \quad (3.8e)$$

$$\hat{c}_{r\nu} = -\hat{c}_{\nu r} := \frac{1 - e^{-\alpha}}{\alpha} \left[ c_{r\nu} - \frac{1}{2\beta} (c_r d_\nu - c_\nu d_r) \right] + \frac{1}{2\beta(\alpha^2 + \beta^2)} \{ -[\alpha \sin \beta - \beta \cos \beta + \beta e^{-\alpha}] \times (c_r c_\nu + d_r d_\nu) + [\alpha \cos \beta + \beta \sin \beta - \alpha e^{-\alpha}] \times (c_r d_\nu - c_\nu d_r) \}, \quad r = 2, \dots, \nu - 1, \quad (3.8f)$$

with  $\epsilon(x) = \pm 1$  for  $x \gtrless 0$ .

*Proof:* Following the same procedure used in Sec. II and by repeated use of the Zassenhaus formula (A1) the theorem is proved.  $\blacktriangledown$

### C. The $\mathfrak{SO}(\nu+2, \nu-1)$ case

This case differs from the preceding ones since the complementary subalgebra  $c$  is no more a solvable Lie algebra but a semidirect sum of the  $\mathfrak{so}(3)$  simple Lie algebra with a solvable Lie algebra. Its  $[1 + \binom{\nu+2}{2}]$  generators are those given by Eqs. (3.4a)–(3.4e) plus

$$(\nu - 1) \text{ generators } \tilde{F}^k = -H_k \cdot H_0, \quad k = 2, \dots, \nu, \quad (3.9a)$$

$$1 \text{ generator } \tilde{H}^1 = -\frac{1}{2}(H_1 + H_{1'}) \cdot H_0, \quad (3.9b)$$

$$1 \text{ generator } \tilde{K}^1 = i/2(H_1 - H_{1'}) \cdot H_0. \quad (3.9c)$$

They satisfy the commutation relations (3.5a)–(3.5g) plus

$$[\tilde{D}, \tilde{F}^\nu] = -\tilde{F}^\nu, \quad (3.10a)$$

$$[\tilde{B}, \tilde{H}^1] = \tilde{K}^1, \quad (3.10b)$$

$$[\tilde{H}^1, \tilde{K}^1] = \tilde{B}, \quad (3.10c)$$

$$[\tilde{K}^1, B] = \tilde{H}^1, \quad (3.10d)$$

$$[\tilde{H}^1, \tilde{F}^k] = 2\tilde{C}^{1k}, \quad (3.10e)$$

$$[\tilde{K}^1, \tilde{F}^k] = 2\tilde{D}^{1k}, \quad (3.10f)$$

$$[\tilde{F}^k, \tilde{F}^l] = 2\tilde{Q}^{kl}, \quad (3.10g)$$

$$[\tilde{H}^1, \tilde{C}^{1l}] = -\frac{1}{2}\tilde{F}^l, \quad (3.10h)$$

$$[\tilde{K}^1, \tilde{D}^{1l}] = -\frac{1}{2}\tilde{F}^l, \quad (3.10i)$$

all other commutators being zero.

Let us define the general element of the complementary subalgebra  $c$  as

$$\tilde{Z} := \alpha \tilde{D} + \beta \tilde{B} + \sum_{k < l=2}^{\nu} c_{kl} \tilde{Q}^{kl} + \sum_{k=2}^{\nu} c_k \tilde{C}^{1k} + \sum_{l=2}^{\nu} d_l \tilde{D}^{1l} + \sum_{m=2}^{\nu} f_m \tilde{F}^m + hH^1 + kK^1. \quad (3.11)$$

Then we have the following theorem.

**Theorem 3.3:** The general group element, obtained by exponentiating the complementary subalgebra  $c$  of the Lie algebra  $\mathfrak{so}(\nu + 2, \nu - 1)$ , is expressed in terms of one-parameter subgroups as follows:

$$e^{\tilde{Z}} = e^J \left[ \prod_{k=2}^{\nu} \exp(\hat{c}_k \tilde{C}^{1k}) \right] \left[ \prod_{l=2}^{\nu} \exp(\hat{d}_l \tilde{D}^{1l}) \right] \times \left[ \prod_{m=2}^{\nu} \exp(\hat{f}_m \tilde{F}^m) \right] \left[ \prod_{k < l=2}^{\nu} \exp(\hat{c}_{kl} \tilde{Q}^{kl}) \right] e^{\alpha \tilde{D}}, \quad (3.12a)$$

where  $J$  denotes the general element of the  $\mathfrak{so}(3)$  simple Lie subalgebra

$$J := \beta \tilde{B} + h \tilde{H}^1 + k \tilde{K}^1, \quad (3.12b)$$

and the coefficients  $\hat{c}_k$ ,  $\hat{d}_l$ ,  $\hat{f}_m$ ,  $\hat{c}_{kl}$  are given explicitly in terms of the coefficients of Eq. (3.11) in Appendix B.

*Proof:* Following the same procedure used in Sec. II and by repeated use of the Zassenhaus formula (A1) the theorem is proved.  $\blacktriangledown$

*Remark:* The group element  $e^J$  is easily factorized into one-parameter subgroups by means of the usual expression involving the standard Euler angles.

## IV. INTRINSIC SPINOR COMPONENTS FOR PURE SPINORS

In this section we give the pure spinor components for the  $\mathfrak{SO}(\nu, \nu)$  and  $\mathfrak{SO}(\nu + 1, \nu)$  cases in terms of the “generalized Euler angles” appearing in the complementary set  $C$ , once it is expressed as a product of one-parameter subgroups as in Eqs. (2.24) and (3.2). We prefer to start with the  $\mathfrak{SO}(\nu + 1, \nu)$  case, since the  $\mathfrak{SO}(\nu, \nu)$  case can be easily deduced from it, taking into account Sec. III A.

### A. The $\mathfrak{SO}(\nu+1, \nu)$ case

Choosing the same spinor representation appearing in Ref. 1, we see that the generic  $\mathfrak{SO}(\nu + 1, \nu)$  pure spinor is given by

$$\psi := e^{\tilde{x}} \psi_m, \quad (4.1)$$



where  $e^{\tilde{X}}$  is given by Eq. (2.24) and  $\psi_m$  is a  $2^\nu$ -component spinor with all components zero but the first one equal to unity. Following Refs. 1 and 2, we label the  $2^\nu$  components of  $\psi$  with a set of completely antisymmetric indices, i.e.,  $\psi_{i_1, \dots, i_r}$ , with  $0 \leq r \leq \nu$  and  $1 \leq i_1, \dots, i_r \leq \nu$ . Then we have the following lemmas.

**Lemma 4.1:** The group element  $e^{\tilde{X}}$  of Eq. (2.10) can be rewritten as

$$e^{\tilde{X}} = \left( 1 + \sum_{k=1}^{\nu} g_k \tilde{F}^k + \sum_{r=1}^{\nu-1} g_r g_\nu \tilde{Q}^{r\nu} \right) \times \left[ \prod_{k,l=1}^{\nu} (1 + \omega_{kl} \tilde{Q}^{kl}) \right] e^{\alpha \tilde{D}}, \quad (4.2a)$$

where

$$\omega_{rs} = -\omega_{sr} = \frac{1}{2}(\tilde{d}_{rs} + g_r g_s), \quad r \neq s = 1, \dots, \nu - 1, \quad (4.2b)$$

$$\omega_{r\nu} = -\omega_{\nu r} = \frac{1}{2}\tilde{d}_{r\nu}, \quad r = 1, \dots, \nu - 1. \quad (4.2c)$$

*Proof:* From Eqs. (2.2g) and (2.20b) we have that

$$e^{U_2} = 1 + \sum_{r=1}^{\nu-1} f_r \tilde{F}^r, \quad (4.3a)$$

and from Eqs. (2.2a), (2.2b), and (2.20a)

$$e^{U_1} = \left[ \prod_{r,s=1}^{\nu-1} \left( 1 + \frac{1}{2} c_{rs} \tilde{Q}^{rs} \right) \right] e^{\alpha \tilde{D}}. \quad (4.3b)$$

Furthermore, from Eqs. (2.16) and (2.26)

$$e^{\alpha \tilde{D}} (e^{\nu} e^{c_2} e^{c_3} \dots) e^{-\alpha \tilde{D}} = (1 + g_\nu \tilde{F}^\nu) \left[ \prod_{r=1}^{\nu-1} \exp(\omega_{r\nu} \tilde{Q}^{r\nu}) \right] \left[ \prod_{r=1}^{\nu-1} \exp(\omega_{rr} \tilde{Q}^{rr}) \right]. \quad (4.4)$$

Inserting Eqs. (4.3a) and (4.3b) into Eq. (2.22) and then, together with Eq. (4.4), into Eq. (2.10), with the help of the relation

$$\tilde{F}^k \tilde{F}^l = \tilde{Q}^{kl}, \quad k \neq l = 1, \dots, \nu, \quad (4.5)$$

and of Eqs. (2.25a)–(2.25d), the lemma is proved.  $\blacktriangledown$

**Lemma 4.2:** The even components of  $\psi$  can be written as

$$\psi_{i_1, \dots, i_{2p}} = e^{\alpha/2} \left\{ \left[ \prod_{k,l=1}^{\nu} (1 + \eta_{kl} \tilde{Q}^{kl}) \right] \psi_m \right\}_{i_1, \dots, i_{2p}}, \quad (4.6a)$$

where

$$\eta_{rs} = -\eta_{sr} = \omega_{rs}, \quad r \neq s = 1, \dots, \nu - 1, \quad (4.6b)$$

$$\eta_{r\nu} = -\eta_{\nu r} = \omega_{r\nu} + \frac{1}{2} g_r g_\nu, \quad r = 1, \dots, \nu - 1, \quad (4.6c)$$

while the odd components of  $\psi$  can be written as

$$\psi_{i_1, \dots, i_{2q+1}} = e^{\alpha/2} \left\{ \left( \sum_{k=1}^{\nu} g_k \tilde{F}^k \right) \left[ \prod_{k,l=1}^{\nu} (1 + \omega_{kl} \tilde{Q}^{kl}) \right] \psi_m \right\}_{i_1, \dots, i_{2q+1}}. \quad (4.7)$$

*Proof:* From Refs. 1 and 2 we have that

$$\tilde{D} \psi_m = \frac{1}{2} \psi_m. \quad (4.8)$$

Furthermore, from Refs. 1 and 2 we have that

$$(H_{l'})^2 = 0, \quad \forall l' = 1, \dots, \nu, \quad (4.9a)$$

and then, from Eqs. (2.1a) and (2.1b),

$$\tilde{Q}^{kl} \tilde{Q}^{ml} = 0 = \tilde{Q}^{kl} \tilde{F}^l = 0 = (\tilde{F}^l)^2, \quad \forall k, l, m = 1, \dots, \nu. \quad (4.9b)$$

Since we know from Refs. 1 and 2 that the only nonzero matrix elements of  $H_j$  have the form

$$(H_j)_{k_1, \dots, k_r, j}^{k_1, \dots, k_r} = 1, \quad j \notin \{k_1, \dots, k_r\}; \quad j, k_1, \dots, k_r = 1, \dots, \nu, \quad r = 0, \dots, \nu - 1, \quad (4.10)$$

we see from Eqs. (2.1a) and (2.1b) that each  $\tilde{Q}^{kl}$  in Eq. (4.2a) “creates” a pair of distinct indices for  $\psi$ , while each  $\tilde{F}^k$  (since  $H_0$  is a diagonal matrix) “creates” only one index for  $\psi$ . After separation of the even and odd pieces in Eq. (4.2a), the lemma is proved.  $\blacktriangledown$

Let us introduce the notation for totally antisymmetric tensors

$$T_{[i_1, \dots, i_r]} := \frac{1}{r!} \sum_{\pi(i_1, \dots, i_r)} (-1)^\pi T_{i_{\pi(1)}, \dots, i_{\pi(r)}}, \quad (4.11)$$

where the summation is extended over all permutations of the indices  $i_1, \dots, i_r$  and  $(-1)^\pi$  denotes the signature of the given permutation. Then we are ready to present the main proposition of this subsection.

**Proposition 4.1:** The generic non-null component of the spinor  $\psi$  of Eq. (4.1) with an even number of indices is given by

$$\psi_{i_1, \dots, i_{2p}} = (2p - 1)! e^{\alpha/2} d_{[i_1 i_2} d_{i_3 i_4} \dots d_{i_{2p-1} i_{2p}]}, \quad (4.12a)$$

while that one with an odd number of indices is given by

$$\psi_{i_1, \dots, i_{2q+1}} = -(2q + 1)! e^{\alpha/2} g_{[i_1} d_{i_2 i_3} \dots d_{i_{2q} i_{2q+1]}}, \quad (4.12b)$$

where

$$d_{kl} = -d_{lk} = \tilde{d}_{kl} + g_k g_l \epsilon(l - k), \quad k \neq l = 1, \dots, \nu, \quad (4.12c)$$

with  $\tilde{d}_{kl}$  and  $g_k$  defined by Eqs. (2.25a)–(2.5d) [and  $(-1)!! = 1$ ].

*Proof:* Let us begin by proving Eq. (4.12a). Taking into account Eq. (4.9b) we can write

$$\begin{aligned} & \left[ \prod_{k,l=1}^{\nu} (1 + \eta_{kl} \tilde{Q}^{kl}) \right] \psi_m \\ &= \left\{ 1 + 2 \sum_{\substack{k,l=1 \\ k < l}}^{\nu} \eta_{kl} \tilde{Q}^{kl} + 2^2 \right. \\ & \quad \times \sum_{\substack{k_1 \neq l_1 \neq k_2 \neq l_2 = 1 \\ k_1 < l_1, k_2 < l_2 \\ l_1 < l_2}}^{\nu} \eta_{k_1 l_1} \eta_{k_2 l_2} \tilde{Q}^{k_1 l_1} \tilde{Q}^{k_2 l_2} + \dots \\ & \quad + 2^{[\nu/2]} \sum_{\substack{k_1 \neq l_1 \neq \dots \neq k_{[\nu/2]} \neq l_{[\nu/2]} = 1 \\ k_1 < l_1, \dots, k_{[\nu/2]} < l_{[\nu/2]} \\ l_1 < l_2 < \dots < l_{[\nu/2]}}}^{\nu} \eta_{k_1 l_1} \dots \eta_{k_{[\nu/2]} l_{[\nu/2]}} \\ & \quad \left. \times \tilde{Q}^{k_1 l_1} \dots \tilde{Q}^{k_{[\nu/2]} l_{[\nu/2]}} \right\} \psi_m. \quad (4.13) \end{aligned}$$

Inserting Eq. (4.13) into Eq. (4.6a) and taking into account Eqs. (2.1a) and (4.10) we see that  $\psi_{i_1, \dots, i_{2p}}$  gets a contribution only from those terms of Eq. (4.13) having  $p$   $\tilde{Q}$ 's with the indices being a permutation of  $(i_1 \dots i_{2p})$ , i.e., only from the terms of the kind

$$\tilde{Q}^{k_1 l_1} \dots \tilde{Q}^{k_p l_p} = (-1)^p \pi \left( \begin{matrix} i_1 i_2 \dots i_{2p} \\ k_1 l_1 \dots k_p l_p \end{matrix} \right) H_{i_1} H_{i_2} \dots H_{i_{2p}}, \quad (4.14a)$$

where  $\pi \left( \begin{matrix} i_1 \dots i_{2p} \\ k_1 l_1 \dots k_p l_p \end{matrix} \right)$  denotes the parity of the permutation transforming the sequence  $(i_1 \dots i_{2p})$  into the sequence  $(k_1 l_1 \dots k_p l_p)$ . Taking into account Eq. (4.10), we have then

$$(\tilde{Q}^{k_1 l_1} \dots \tilde{Q}^{k_p l_p} \psi_m)_{i_1 \dots i_{2p}} = \pi \left( \begin{matrix} i_1 i_2 \dots i_{2p} \\ k_1 l_1 \dots k_p l_p \end{matrix} \right), \quad (4.14b)$$

and, using Eqs. (4.6a), (4.13), (4.14b), and the definition (4.11), the first part of the proposition is proved, with the help of the relation

$$\sum_{\substack{(k_1, l_1, \dots, k_p, l_p) \in (i_1, \dots, i_{2p}) \\ k_1 \neq l_1 \neq \dots \neq k_p \neq l_p \\ k_1 < l_1, \dots, k_p < l_p \\ l_1 < l_2 < \dots < l_p}} \pi \left( \begin{matrix} i_1 i_2 \dots i_{2p} \\ k_1 l_1 \dots k_p l_p \end{matrix} \right) \eta_{k_1 l_1} \dots \eta_{k_p l_p} \\ = (2p-1)!! \eta_{[i_1 i_2 \dots i_{2p-1} i_{2p}]}, \quad (4.14c)$$

and of the definitions (4.6b), (4.6c), (4.2b), (4.2c), and (4.12c).

Let us now rewrite Eq. (4.7) as follows:

$$\psi_{i_1 \dots i_{2q+1}} = \sum_{r=0}^{\nu} \sum_{j_1, \dots, j_r=1}^{\nu} \left( \sum_{k=1}^{\nu} g_k \tilde{F}^k \right)_{i_1 \dots i_{2q+1}}^{j_1 \dots j_r} \tilde{\psi}_{j_1 \dots j_r}, \quad (4.15a)$$

where

$$\tilde{\psi}_{j_1 \dots j_r} = e^{\alpha/2} \left\{ \left[ \prod_{k,l=1}^{\nu} (1 + \omega_{kl} \tilde{Q}^{kl}) \right] \psi_m \right\}_{j_1 \dots j_r}. \quad (4.15b)$$

We see that  $\tilde{\psi}_{j_1 \dots j_r}$  is identical to  $\psi_{i_1 \dots i_{2p}}$  of Eq. (4.6a) (non-null contributions come for  $r$  even only), once we substitute  $\eta_{kl}$  with  $\omega_{kl}$ . Its explicit form is given by

$$\tilde{\psi}_{j_1 \dots j_r} = 2^{r/2} (r-1)!! e^{\alpha/2} \omega_{[j_1 j_2 \dots j_{r-1} j_r]}. \quad (4.16)$$

Since  $H_0$  is a diagonal matrix with

$$(H_0)_{i_1 \dots i_s}^{i_1 \dots i_s} = (-1)^s, \quad i_1, \dots, i_s = 1, \dots, \nu, \quad s = 0, \dots, \nu, \quad (4.17)$$

taking into account the definition (2.1b) and Eqs. (4.10), (4.16), and (4.17), we can write

$$\psi_{i_1 \dots i_{2q+1}} = \sum_{k=1}^{2q+1} (-1)^k g_k \tilde{\psi}_{i_1 \dots \hat{i}_k \dots i_{2q+1}}, \quad (4.18)$$

where a  $\hat{\phantom{x}}$  sign over an index means that that particular index is missing from the sequence in which it appears. If we now express the  $\tilde{\psi}$  components by means of Eq. (4.16) and we take into account definitions (4.11), (4.2b), (4.2c), and (4.12c), we see that Eq. (4.12b) is also proved.  $\blacktriangledown$

*Remark:* From the explicit expressions (4.12a) and (4.12b) it is easy to check that the pure spinor's components satisfy the quadratic relations presented in Ref. 2.

## B. The SO( $\nu, \nu$ ) case

In this case we define the generic pure spinor as

$$\psi = e^{\tilde{Y}} \psi_{m^+}, \quad (4.19)$$

where  $e^{\tilde{Y}}$  is given by Eq. (3.2) and  $\psi_{m^+}$  is identical to  $\psi_m$  of the preceding subsection. From Refs. 1 and 2 we see that  $\psi$  is

effectively a semispinor of the first kind, having all components with an odd number of indices equal to zero. [Any second kind semispinor can be obtained by reflections of first kind semispinors (see Ref. 2).]

As already mentioned in Sec. III A, this case can be derived from the SO( $\nu+1, \nu$ ) case, deleting all  $\tilde{F}^k$  generators. We have then the following proposition.

*Proposition 4.2:* The generic non-null component of the semispinor  $\psi$  of Eq. (4.19) is given by

$$\psi_{i_1 \dots i_{2p}} = (2p-1)!! e^{\alpha/2} \tilde{c}_{[i_1 i_2 \dots i_{2p-1} i_{2p}]}, \quad (4.20)$$

with  $\tilde{c}_{kl}$  given by Eqs. (3.3a) and (3.3b).

*Proof:* From Eqs. (2.3) and (3.1) we see that, if we let  $f_k \equiv 0$ ,  $k = 1, \dots, \nu$ , we have

$$e^{\tilde{X}}|_{\{f_k \equiv 0\}} = e^{\tilde{Y}}, \quad (4.21a)$$

then, since  $\psi_{m^+} \equiv \psi_m$ , from Eqs. (4.1), (4.19), and (4.21a) we have that

$$\psi_{\text{SO}(\nu, \nu)} = \psi_{\text{SO}(\nu+1, \nu)}|_{\{f_k \equiv 0\}}. \quad (4.21b)$$

Then the only nonzero components of  $\psi_{\text{SO}(\nu, \nu)}$  are given by Eq. (4.12a), with the parameters  $d_{kl}$  evaluated with  $f_k \equiv 0$ . Taking into account Eqs. (4.12c), (2.25a)–(2.25d) and (3.3a), (3.3b), we see that the proposition is proved.  $\blacktriangledown$

*Remark:* We want to stress that in Eqs. (4.12a), (4.12b), and (4.20) we have been able to write the nonlinear (pure) spinor components in terms of the independent intrinsic coordinates, just because we succeeded in writing both  $e^{\tilde{X}}$  and  $e^{\tilde{Y}}$  as a product of one-parameter subgroups. It can be easily seen also that such a “factorization” of the group elements of the complementary set  $C$  is absolutely necessary for the realization of minimal covariant nonlinear spinor wave equations associated with our nonlinear spinor representations (see Ref. 3).

## APPENDIX A: THE ZASSENHAUS FORMULA

In order to make the paper as self-contained as possible, we give here a brief account on what is presented in Ref. 4 on the Zassenhaus formula, and derive some straightforward consequences to be used in our analysis in the main text.

Let  $R$  be a free ring with two generators  $x, y$  and with rational coefficients. Then we have the following theorem.

**Theorem A1 (Zassenhaus):** There exist uniquely determined Lie elements  $C_n$  ( $n = 2, 3, 4, \dots$ ) in  $R$ , which are exactly of degree  $n$  in  $x, y$ , such that

$$e^{x+y} = e^x e^y e^{C_2} e^{C_3} \dots e^{C_n} \dots \quad (A1)$$

Here a rough definition of  $C_n$  is that it is any linear combination of multiple commutators of  $x$  and  $y$ , in which the total power of  $x$  and  $y$  is  $n$  (see Ref. 4 for a more precise definition).

In order to get the  $C_n$ 's explicitly we have to introduce the “curly bracket operator”  $\{ \}$  having the following properties.

(i) It is a linear operator, i.e., for any two elements  $F_1, F_2 \in R$  and for any two elements  $c_1, c_2$  in the field of coefficients

$$\{c_1 F_1 + c_2 F_2\} = c_1 \{F_1\} + c_2 \{F_2\}. \quad (A2)$$

(ii) Let  $x_\nu$  for  $\nu = 1, 2, \dots, n$  be any one of the generators. Then for any monomial  $x_1 x_2 \cdots x_n$  we define

$$\{x_1 x_2 \cdots x_n\} := [[\cdots [x_1, x_2], x_3], \dots], x_n], \quad (\text{A3})$$

$$\{x_\nu\} = x_\nu, \quad (\text{A4})$$

$$\{1\} = 0. \quad (\text{A5})$$

Furthermore, the two following propositions are needed.

*Proposition A1:* Let  $G$  be a homogeneous Lie element in  $R$  which is of degree  $n$ ; then

$$\{G\} = nG. \quad (\text{A6})$$

*Proposition A2:* If  $G$  is a homogeneous Lie element and  $F$  is any element of  $R$  then

$$\{G^2 F\} = 0. \quad (\text{A7})$$

Then, applying the operator  $\{\}$  to both sides of Eq. (A1) we get

$$\{e^{x+y}\} = x + y \quad (\text{A8})$$

and

$$\begin{aligned} \{e^x e^y e^{c_2} e^{c_3} \dots\} &= x + y + \{xy\} + \{C_2\} + \{xy^2\}/2! \\ &+ \{xC_2\} + \{yC_2\} + \{C_3\} + \dots \end{aligned} \quad (\text{A9})$$

By comparing terms of the same degree in Eqs. (A8) and (A9), we get

$$\{C_2\} + \{xy\} = 0, \quad (\text{A10})$$

$$\{C_3\} + \{xC_2\} + \{yC_2\} + \frac{1}{2}\{xy^2\} = 0, \quad (\text{A11})$$

and so on, giving

$$C_2 = -\frac{1}{2}[x, y], \quad (\text{A12})$$

$$C_3 = -\frac{1}{6}[[x, y], y] - \frac{1}{6}[[x, y], x], \quad (\text{A13})$$

and so on.

If we now look at the relations (A10), (A11), and to the analogous ones of higher degree, we see that

$$\{C_n\} = -\{xC_{n-1}\} + (\text{terms at least of degree 2 in } y), \quad (\text{A14})$$

i.e., from Eq. (A3) and (A6),

$$nC_n = [C_{n-1}, x] + (\text{terms at least of degree 2 in } y). \quad (\text{A15})$$

Iterating Eq. (A15) and using Eq. (A12), we finally get

$$\begin{aligned} C_n &= (1/n!) \underbrace{[[\cdots [y, x], x], \dots], x]}_{n-1} \\ &+ (\text{multiple commutators with } y \\ &\text{appearing at least twice}). \end{aligned} \quad (\text{A16})$$

## APPENDIX B: EXPLICIT FORM OF THE COEFFICIENTS IN THE $SO(\nu+2, \nu-1)$ CASE

The relations connecting the coefficients  $\dot{c}_k, \dot{d}_l, \dot{f}_m, \dot{c}_{kl}$  appearing in Eq. (3.12a) to the coefficients  $\alpha, \beta, c_{kl}, c_k, d_l$ , of Eq. (3.11) are rather cumbersome; therefore we prefer to write them explicitly in this appendix as follows:

$$\begin{aligned} \dot{c}_r &:= \left(2h \frac{\cos \rho - 1}{\rho^2} - \beta k \frac{\sin \rho - \rho}{\rho^3}\right) f_r + \left[1 + \frac{\sin \rho - \rho}{\rho^3} (h^2 + \beta^2)\right] c_r \\ &+ \left(hk \frac{\sin \rho - \rho}{\rho^3} - \beta \frac{\cos \rho - 1 - \frac{1}{2}\rho^2}{\rho^2}\right) d_r, \quad r = 2, \dots, \nu - 1, \end{aligned} \quad (\text{B1a})$$

$$\begin{aligned} \dot{c}_\nu &:= \frac{1 - e^{-\alpha}}{\alpha} c_\nu + \frac{1}{\alpha^2 + \rho^2} \left(\alpha \frac{\sin \rho}{\rho} - \cos \rho + e^{-\alpha}\right) (\beta d_\nu - 2h f_\nu) \\ &- \left[\frac{1}{\alpha \rho^2} - \frac{1}{\alpha^2 + \rho^2} \left(\alpha \frac{\cos \rho}{\rho^2} + \frac{\sin \rho}{\rho} + \frac{e^{-\alpha}}{\alpha}\right)\right] [(\beta^2 + h^2) c_\nu + h k d_\nu + 2k \beta f_\nu], \end{aligned} \quad (\text{B1b})$$

$$\begin{aligned} \dot{d}_r &:= \left(2k \frac{\cos \rho - 1}{\rho^2} + \beta h \frac{\sin \rho - \rho}{\rho^3}\right) f_r + \left(hk \frac{\sin \rho - \rho}{\rho^3} + \beta \frac{\cos \rho - 1 - \frac{1}{2}\rho^2}{\rho^2}\right) c_r \\ &+ \left[1 + (\beta^2 + k^2) \frac{\sin \rho - \rho}{\rho^3}\right] d_r, \quad r = 2, \dots, \nu - 1, \end{aligned} \quad (\text{B1c})$$

$$\begin{aligned} \dot{d}_\nu &:= \frac{1 - e^{-\alpha}}{\alpha} d_\nu - \frac{1}{\alpha^2 + \rho^2} \left(\alpha \frac{\sin \rho}{\rho} - \cos \rho + e^{-\alpha}\right) (\beta c_\nu + 2k f_\nu) \\ &- \left[\frac{1}{\alpha \rho^2} - \frac{1}{\alpha^2 + \rho^2} \left(\alpha \frac{\cos \rho}{\rho^2} + \frac{\sin \rho}{\rho} + \frac{e^{-\alpha}}{\alpha}\right)\right] [(\beta^2 + k^2) d_\nu + h k c_\nu - 2h \beta f_\nu], \end{aligned} \quad (\text{B1d})$$

$$\begin{aligned} \dot{f}_r &:= \left[1 - \frac{1}{2} (h^2 + k^2) \frac{\sin \rho - \rho}{\rho^3}\right] f_r + \frac{1}{2} \left(\beta k \frac{\sin \rho - \rho}{\rho^3} - h \frac{\cos \rho - 1 - \frac{1}{2}\rho^2}{\rho^2}\right) c_r \\ &- \frac{1}{2} \left(\beta h \frac{\sin \rho - \rho}{\rho^3} + k \frac{\cos \rho - 1 - \frac{1}{2}\rho^2}{\rho^3}\right) d_r, \quad r = 2, \dots, \nu - 1, \end{aligned} \quad (\text{B1e})$$

$$\dot{f}_\nu := \frac{1 - e^{-\alpha}}{\alpha} f_\nu + \frac{1}{2} \frac{1}{\alpha^2 + \rho^2} \left(\alpha \frac{\sin \rho}{\rho} - \cos \rho + e^{-\alpha}\right) (h c_\nu + k d_\nu)$$

$$+ \left[ \frac{1}{\alpha \rho^2} - \frac{1}{\alpha^2 + \rho^2} \left( \alpha \frac{\cos \rho}{\rho^2} + \frac{\sin \rho}{\rho} + \frac{e^{-\alpha}}{\alpha} \right) \right] \left[ \frac{1}{2} \beta (h d_v - k c_v) - (h^2 + k^2) f_v \right], \quad (\text{B1f})$$

$$\begin{aligned} \dot{c}_{rs} = -\dot{c}_{sr} = c_{rs} - 2 \left( h \frac{\sin \rho - \rho}{\rho^3} + k \beta \frac{\cos \rho - 1 + \frac{1}{2} \rho^2}{\rho^4} \right) c_{[r} f_{s]} + 2 d_{[r} \left[ h \frac{\cos \rho - 1 + \frac{1}{2} \rho^2}{\rho^4} \left( \beta f_{s]} - \frac{1}{2} k c_{s]} \right) \right. \\ \left. - \frac{\sin \rho - \rho}{\rho^3} \left( \frac{1}{2} \beta c_{s]} + k f_{s]} \right) \right] - \left[ \frac{1}{4} (\dot{c}_r \dot{c}_s + \dot{d}_r \dot{d}_s) + \dot{f}_r \dot{f}_s \right] \epsilon(s-r), \quad r \neq s = 2, \dots, \nu-1, \end{aligned} \quad (\text{B1g})$$

$$\begin{aligned} \dot{c}_{rv} = -\dot{c}_{vr} = \frac{1 - e^{-\alpha}}{\alpha} c_{rv} - \frac{1}{2} \frac{1}{\alpha^2 + \rho^2} \left( \alpha \frac{\sin \rho}{\rho} - \cos \rho + e^{-\alpha} \right) (c_r c_v + d_r d_v + 4 f_r f_v) \\ + \left[ \frac{1}{\alpha \rho^2} - \frac{1}{\alpha^2 + \rho^2} \left( \alpha \frac{\cos \rho}{\rho^2} + \frac{\sin \rho}{\rho} + \frac{e^{-\alpha}}{\alpha} \right) \right] (2 h c_{[r} f_{v]} + 2 k d_{[r} f_{v]} - \beta c_{[r} d_{v]}) \\ - \frac{1}{2 \rho^2} \left[ \frac{1}{\alpha^2 + \rho^2} \left( \alpha \frac{\sin \rho}{\rho} - \cos \rho + e^{-\alpha} \right) + \frac{e^{-\alpha} - 1 + \alpha}{\alpha^2} \right] [h k (c_r d_v + c_v d_r) \\ + 2 \beta k (c_r f_v + c_v f_r) - 2 \beta h (d_r f_v + d_v f_r) - k^2 (c_r c_v + d_r d_v) - 4 \beta^2 f_r f_v], \\ r = 2, \dots, \nu-1, \end{aligned} \quad (\text{B1h})$$

where

$$\rho^2 := \beta^2 + h^2 + k^2 \quad (\text{B2a})$$

and

$$c_{[r} f_{s]} := \frac{1}{2} (c_r f_s - c_s f_r). \quad (\text{B2b})$$

<sup>1</sup>P. Furlan and R. Rączka, "Nonlinear spinor representations," *J. Math. Phys.* **26**, 3021 (1985).

<sup>2</sup>E. Cartan, *The Theory of Spinors* (Hermann, Paris, 1966).

<sup>3</sup>P. Furlan and R. Rączka, "Intrinsic nonlinear spinor wave equations associated with nonlinear spinor representations," *J. Math. Phys.* **27**, 1883 (1986).

<sup>4</sup>W. Magnus, "On the exponential solution of differential equations for a linear operator," *Commun. Pure Appl. Math.* **7**, 649 (1954).

# Decomposition of the $SO^*(8)$ enveloping algebra under $U(4) \supset U(3)$

R. Le Blanc and D. J. Rowe

Department of Physics, University of Toronto, Toronto, Ontario, Canada M5S 1A7

(Received 29 September 1986; accepted for publication 7 January 1987)

The existence of a complete set of  $SU(3)$  tensor operators in the enveloping algebra of  $SO^*(8)$  is demonstrated. The analysis recasts a parallel analysis by Biedenharn and Flath [Commun. Math. Phys. **93**, 143 (1984)] concerning an  $SO(6,2)$  model for  $SU(3)$  in the isomorphic but simpler framework of an  $SO^*(8)$  model.

## I. INTRODUCTION

It has been shown independently by Biedenharn and Flath<sup>1</sup> and Bracken and MacGibbon<sup>2</sup> that the fundamental unirrep of  $SO(6,2)$  defines a model for  $SU(3)$  in the sense of Bernštein *et al.*<sup>3</sup> According to their definition, a representation space is called a model for a group  $G$  if it contains precisely one irrep from every equivalence class of irreps of  $G$ .

Exploiting the local isomorphism  $SO(6,2) \sim SO^*(8)$ , Le Blanc and Rowe<sup>4</sup> gave a Bargmann representation of the  $SO(6,2)$  model for  $SU(3)$  which is simple, naturally unitary with respect to the Bargmann measure, and readily expressed in a Cartan basis. Furthermore, taking advantage of the natural embedding of the  $su(3)$  algebra in the  $so^*(8)$  Lie algebra through the subgroup chain

$$SO^*(8) \supset U(4) \supset U(3) \supset SU(3),$$

where  $U(4) \sim SO(6) \times SO(2)$  is the maximal compact subgroup of  $SO^*(8)$ , it was realized therein that the Cartan generators of the  $so^*(8)$  Lie algebra are the fundamental Wigner operators as defined by Biedenharn and Louck (see, e.g., Louck<sup>5</sup>) in contrast with the  $SO(6,2)$  model which call for more complicated linear combinations of the generators of its Lie algebra for their realization.

Since the fundamental Wigner operators of  $SU(3)$  are the elementary building blocks for a complete set of basic  $SU(3)$  tensor operators, it is appropriate to try to identify the elements of this set in the enveloping algebra of  $SO^*(8)$ . Such an analysis has been carried out by Biedenharn and Flath<sup>1</sup> for the  $SO(6,2)$  model. The purpose of this paper is to recast their analysis in the much simpler  $SO^*(8)$  model framework. Using only the elegant properties of the Gel'fand patterns required for the classification of the  $U(4) \supset U(3)$  tensors arising in the enveloping algebra of  $SO^*(8)$ , we rederive all their results in a simple deductive and therefore more transparent way.

## II. REVIEW OF THE $SO^*(8)$ MODEL FOR $SU(3)$

A basis for the complexification of the  $so^*(8)$  Lie algebra is given (Le Blanc and Rowe<sup>4</sup>) in terms of two four-dimensional Bargmann vectors ( $g_{\alpha\mu}$ ;  $\alpha = 1,2$ ;  $\mu = 1,\dots,4$ ) and with summation over repeated indices by

$$C_{\mu\nu} = \frac{1}{2} \left( g_{\alpha\mu} \frac{\partial}{\partial g_{\alpha\nu}} + \frac{\partial}{\partial g_{\alpha\nu}} g_{\alpha\mu} \right), \quad \mu, \nu = 1,4, \\ = g_{\alpha\mu} \frac{\partial}{\partial g_{\alpha\nu}} + \delta_{\mu\nu}, \quad (2.1a)$$

$$A_{\mu\nu} = -A_{\nu\mu} = \begin{vmatrix} g_{1\mu} & g_{1\nu} \\ g_{2\mu} & g_{2\nu} \end{vmatrix}, \quad (2.1b)$$

$$B_{\mu\nu} = -B_{\nu\mu} = A_{\mu\nu}^\dagger = \begin{vmatrix} \frac{\partial}{\partial g_{1\mu}} & \frac{\partial}{\partial g_{1\nu}} \\ \frac{\partial}{\partial g_{2\mu}} & \frac{\partial}{\partial g_{2\nu}} \end{vmatrix}, \quad (2.1c)$$

where  $(C_{\mu\nu})$  span the maximal compact subalgebra  $u(4)$  of  $so^*(8)$  and  $(A_{\mu\nu})$  and  $(B_{\mu\nu})$  are, respectively, Cartan raising and lowering operators. The  $u(3)$  subalgebra  $(C_{ij})$  is given by the restriction of the indices  $\mu, \nu$  in (2.1a) to  $i, j = 1,2,3$ .

We have the following commutation relations:

$$[C_{\mu\nu}, A_{\gamma\delta}] = \delta_{\nu\gamma} A_{\mu\delta} + \delta_{\nu\delta} A_{\gamma\mu}, \\ [C_{\mu\nu}, B_{\gamma\delta}] = -\delta_{\mu\gamma} B_{\nu\delta} - \delta_{\mu\delta} B_{\gamma\nu}, \quad (2.2) \\ [B_{\mu\nu}, A_{\gamma\delta}] = \delta_{\delta\nu} C_{\gamma\mu} + \delta_{\gamma\mu} C_{\delta\nu} - \delta_{\nu\gamma} C_{\delta\mu} - \delta_{\mu\delta} C_{\gamma\nu}.$$

Under  $u(4)$ ,  $A$  is a  $\{1100\}$  tensor,  $B$  is a  $\{00-1-1\}$  tensor while the elements of  $su(4)$  are the components of a  $\{100-1\}$  self-conjugate tensor.

The lowest weight state for this fundamental representation of  $so^*(8)$  is clearly given by the Bargmann vacuum

$$\langle g|0\rangle = 1, \quad (2.3)$$

which carries an unirrep  $\{1111\}$  of  $u(4)$ .

Since the raising operators  $(A_{\mu\nu})$  transform under the adjoint action of  $u(4)$  as the components of a  $\{1100\}$  tensor, the polynomials of degree  $h_1$  in  $(A_{\mu\nu})$  transform as the components of a  $\{h_1, h_1, 00\}$  tensor. These tensors then reduce under  $u(3)$  according to the branching rule

$$u(4) \downarrow u(3): \{h_1, h_1, 00\} \downarrow \sum_{h_2=0}^{h_1} \{h_1, h_2, 0\} \quad (2.4)$$

given by the usual betweenness conditions of the corresponding Gel'fand patterns. Since there is no multiplicity involved in the  $SO^*(8) \downarrow U(4)$  reduction, it follows that a  $U(3)$  invariant subspace of this  $so^*(8)$  representation space can be uniquely labeled by its  $u(4) \supset u(3)$  quantum numbers

$$\left| \begin{array}{cccc} h_1 & h_1 & 0 & 0 \\ & h_1 & h_2 & 0 \\ & & \eta & \end{array} \right\rangle = P_\eta^{\{h_1, h_2\}}(A) |0\rangle, \quad (2.5a)$$

where  $P_\eta^{\{h_1, h_2\}}(A)$  is a polynomial of degree  $h_1$  in the raising operators  $(A_{\mu\nu})$  and  $\eta$  stands for any appropriate scheme to

label components of the corresponding  $su(3)$  unirrep  $\{h_1, h_2\}$ . For example,  $\eta$  can stand for either the usual lower Gel'fand pattern or for the basis labels  $(\delta)LM$  (Le Blanc and Rowe<sup>6</sup>) corresponding to a canonical  $SU(3) \downarrow SO(3)$  reduction. When  $\eta$  stands for an  $SU(3)$  lowest weight, we have, from Eqs. (4.7) and Eq. (6.13) of Ref. 4,

$$\begin{aligned} & \left| \begin{array}{cccc} h_1 & & 0 & 0 \\ & h_1 & & \\ & h_1 & h_2 & 0 \\ & & lw & \end{array} \right\rangle \\ &= \left[ \frac{h_2!}{h_1!(h_1 - h_2)!} \right]^{1/2} C_{42}^{(h_1, -h_2)} \left\{ \frac{1}{h_1! \sqrt{h_1 + 1}} A_{12}^{h_1} \right\} |0\rangle \\ &= \frac{1}{[(h_1 + 1)! h_2! (h_1 - h_2)!]^{1/2}} A_{12}^{h_2} A_{14}^{h_1 - h_2} |0\rangle. \quad (2.5b) \end{aligned}$$

It does follow that the  $so^*(8)$  unirrep  $\{1111\}$  decomposes under the successive restrictions

$$so^*(8) \downarrow u(4) \downarrow u(3): \{1111\} \downarrow \sum_{h_1} \{h_1, h_1, 00\} \downarrow \sum_{h_1, h_2} \{h_1, h_2, 0\}. \quad (2.6)$$

Thus the  $\{1111\}$  representation space for  $so^*(8)$  contains precisely one representative of every equivalence of  $SU(3)$  irreducible representations. It is therefore, by definition, an  $SU(3)$  model space (Bernstein *et al.*<sup>3</sup>).

In the following sections, we identify in the enveloping algebra of  $so^*(8)$  a complete set of  $SU(3)$  tensor operators which act on this model space and which have well-defined shift properties. In classifying these tensors, we shall make use of the convenient and insightful labeling scheme by upper Gel'fand (operator) patterns introduced by Biedenharn and Louck (see, e.g., Louck<sup>5</sup>) who showed that a complete set of basic  $SU(3)$  tensors can be classified by means of operator patterns of the type

$$\left[ \begin{array}{ccc} & \delta_1 & \\ \gamma_1 & & \gamma_2 \\ h_{13} & h_{23} & h_{33} \end{array} \right] \quad (2.7)$$

$$\left[ \begin{array}{ccc} & \delta_1 & \\ \gamma_1 & & \gamma_2 \\ h_{13} & h_{23} & h_{33} \end{array} \right] = \left[ \begin{array}{ccc} \Gamma_{11} + h_{33} & & \\ \Gamma_{12} + h_{33} & \Gamma_{22} + h_{33} & \\ h_2 + h_{33} & & h_{33} \end{array} \right], \quad (2.12)$$

with

$$h_{33} = \Gamma_{12} + \Gamma_{22} - h_1 - h_2 \leq 0. \quad (2.13)$$

Note that, from (2.10) and (2.11),

$$\gamma_1 + \gamma_2 = \delta_1 + \delta_2 = \sum_i h_{i3}. \quad (2.14)$$

Note also that for given values of  $h_1 = h_{13} - h_{33}$ ,  $h_2 = h_{23} - h_{33}$ ,  $\delta_1$ , and  $\delta_2$ , there generally corresponds a multiplicity set of tensors

$$\left[ \begin{array}{ccc} & \delta_1 & \\ \gamma_1(\rho) & & \gamma_2(\rho) \\ h_{13} & h_{23} & h_{33} \end{array} \right], \quad (2.15)$$

which can be indexed by an integer  $0 \leq \rho \leq \rho_{\max}$  with

$$\gamma_1(\rho) = \gamma_1^s - \rho, \quad \gamma_2(\rho) = \gamma_2^s + \rho, \quad (2.16)$$

with the usual betweenness conditions for  $\gamma_1$ ,  $\gamma_2$ , and  $\delta_1$ . Such a pattern indicates that the corresponding tensor maps the states of a  $U(3)$  representation  $\{\lambda_1, \lambda_2, \lambda_3\}$  into a representation  $\{\lambda_1 + \Delta_1, \lambda_2 + \Delta_2, \lambda_3 + \Delta_3\}$ , where

$$\begin{aligned} \Delta_1 &= \delta_1, & \Delta_2 &= \gamma_1 + \gamma_2 - \delta_1, \\ \Delta_3 &= h_{13} + h_{23} + h_{33} - \gamma_1 - \gamma_2. \end{aligned} \quad (2.8)$$

Thus it maps an  $SU(3)$  representation  $\{\Lambda_1, \Lambda_2\}$  into a representation  $\{\Lambda_1 + \Delta_1 - \Delta_3, \Lambda_2 + \Delta_2 - \Delta_3\}$ .

To label a set of unit tensors for  $SU(3)$ , Biedenharn and Flath restricted the  $U(3)$  patterns (2.7) to the subset of the type

$$\left[ \begin{array}{ccc} & \Gamma_{11} & \\ \Gamma_{12} & & \Gamma_{22} \\ h_1 & h_2 & 0 \end{array} \right], \quad (2.9)$$

and declared the equivalence of the  $U(3)$  unirreps

$$\{\lambda_1, \lambda_2, \lambda_3\} \equiv \{\lambda_1 + \Delta_3, \lambda_2 + \Delta_3, \lambda_3 + \Delta_3\}$$

in their model space. Thus they restricted to the subset of patterns (2.7) with  $h_{33} = 0$ . For our purposes, however, it is more appropriate to restrict to the subset with  $\Delta_3 = 0$ . This is because the  $SU(3)$  states of our model space carry  $U(3)$  representations strictly of the type  $\{\lambda_1, \lambda_2, 0\}$ ; cf. Eq. (2.6). It follows that any tensor operator acting on the model space must have shift  $\Delta_3 = 0$ . We therefore classify tensors acting within the  $so^*(8)$  representation space by patterns of the type (2.7) with the constraint

$$\sum_i h_{i3} = \gamma_1 + \gamma_2. \quad (2.10)$$

The corresponding shifts (2.8) are then given by

$$\begin{aligned} \Delta_1 &= \delta_1, & \Delta_2 &= \gamma_1 + \gamma_2 - \delta_1 \equiv \delta_2, \\ \Delta_3 &= 0. \end{aligned} \quad (2.11)$$

Our labeling is clearly related to that of Biedenharn and Flath by

$$\rho_{\max} = \min(\gamma_1^s - \delta_1, \gamma_1^s - h_{23}, h_{23} - \gamma_2^s, \delta_1 - \gamma_2^s), \quad (2.17)$$

where  $\gamma_1^s$  and  $\gamma_2^s$ , respectively, denote the maximum and minimum (stretched) values of  $\gamma_1$  and  $\gamma_2$  allowed by the Gel'fand betweenness conditions. Thus a basic  $SU(3)$  tensor can be identified by either of the two sets of labels

$$(h_{13}, h_{23}, h_{33}, \delta_1, \rho), \quad (2.18a)$$

$$(h_1, h_2, \delta_1, \delta_2, \rho), \quad (2.18b)$$

from which one readily reconstructs the corresponding operator pattern.

It can be easily verified (cf. Proposition 4.1) that the generators of the  $so^*(8)$  Lie algebra have the following  $U(3)$  tensorial and shift properties:

$$\begin{aligned}
A_{jk}^{\{110\}}: \delta &= (11), & B_{jk}^{\{0-1-1\}}: \delta &= (-1-1), \\
A_{i4}^{\{100\}}: \delta &= (10), & B_{i4}^{\{00-1\}}: \delta &= (-10), \\
C_{i4}^{\{100\}}: \delta &= (01), & C_{4i}^{\{00-1\}}: \delta &= (0-1).
\end{aligned} \tag{2.19}$$

Therefore the expressions for the fundamental Wigner operators are seen to be extremely simple in this model. Their Hermiticity relations are also clearly apparent.

### III. DECOMPOSITION OF $A$ UNDER $SO^*(8) \supset U(4) \supset U(3)$

We now proceed to the decomposition of the  $SO^*(8)$  enveloping algebra  $A$  under  $SU(3)$ .

*Proposition 3.1:* Under  $SO^*(8)$ ,  $A$  decomposes as

$$\{A\} \downarrow \sum_{p=0}^{\infty} \{00-p-p\}, \tag{3.1}$$

where  $\{00-p-p\}$  is a nonunitary finite-dimensional lowest weight irrep of  $SO^*(8)$ .

*Proof:*  $B_{34}^p$  is easily verified to be the unique lowest weight polynomial of degree  $p$  in the generators of the  $so^*(8)$  Lie algebra. It satisfies the equations

$$\begin{aligned}
[B_{\mu\nu}, B_{34}^p] &= 0, \\
[C_{\mu\nu}, B_{34}^p] &= 0, \quad \mu < \nu, \\
[C_{11}, B_{34}^p] &= [C_{22}, B_{34}^p] = 0, \\
[C_{33}, B_{34}^p] &= [C_{44}, B_{34}^p] = -pB_{34}^p.
\end{aligned} \tag{3.2} \quad \text{Q.E.D.}$$

(See Biedenharn and Flath,<sup>1</sup> Proposition 7.3, for an equivalent statement.)

*Proposition 3.3:* The decomposition of  $\{00-p-p\}$  under  $U(4)$  is given by the branching rule  $SO^*(8) \downarrow U(4)$ :

$$\begin{aligned}
\{00-p-p\} \downarrow \sum_{p_1+p_2+p_3+p_4=p} \{p_1+p_3, p_1, -p_2, \\
-p_2-p_3\}.
\end{aligned} \tag{3.3}$$

*Proof:* The  $SO^*(8)$  irreducible tensor  $\{00-p-p\}$  decomposes as a sum of  $U(4)$  Racah tensors  $T^{\{h_{\mu\alpha}\}}$  for which the  $U(4)$  lowest weight component is given by the product of commuting factors

$$T_{lw}^{\{h_{\mu\alpha}\}} = A_{12}^{p_1} B_{34}^{p_2} C_{14}^{p_3} C_{\mu\mu}^{p_4}, \quad \sum_{\nu=1}^4 p_{\nu} = p. \tag{3.4}$$

Here  $A_{12}$ ,  $B_{34}$ ,  $C_{14}$ , and  $C_{\mu\mu}$  are, respectively, lowest weight components of  $U(4)$  tensors of rank  $\{1100\}$ ,  $\{00-1-1\}$ ,  $\{100-1\}$ , and  $\{0000\}$ . It follows that the coupling is stretched and the resulting tensor  $T^{\{h_{\mu\alpha}\}}$  is of rank  $\{h_{\mu\alpha}\}$  with

$$\begin{aligned}
h_{14} &= p_1 + p_3, & h_{24} &= p_1, & h_{34} &= -p_2, \\
h_{44} &= -p_2 - p_3.
\end{aligned} \tag{3.5}$$

Q.E.D.

The decomposition of a  $U(4)$  tensor into a sum of  $U(3)$  tensors is given by the betweenness conditions for the associated  $U(4) \supset U(3)$  Gel'fand patterns:

$$U(4) \downarrow U(3): \{h_{\mu\alpha}\} \downarrow \sum \{h_{i3}\}, \quad h_{i-1,4} < h_{i3} < h_{i+1,4}. \tag{3.6}$$

Observe, however, that from any given  $p$  level  $U(3)$  tensor in  $A$ , we can construct other ("old") tensors at a higher  $p$  level and of the same  $U(3)$  rank by multiplying it by arbitrary polynomials in (a) the  $U(1) \subset U(1) \times SU(4) \subset U(4)$  scalar

$$C_{\mu\mu} \tag{3.7a}$$

and (b) the  $U(1) \subset U(1) \times SU(3) \subset SU(4)$  scalar

$$C_{ii}^{tr} = C_{ii} - \frac{3}{4} C_{\mu\mu}. \tag{3.7b}$$

It is therefore sufficient to restrict consideration to the subset of "new"  $p$  level  $U(3)$  tensors having no such factors. Let  $I$  denote the set of polynomials in the above two  $U(1)$  scalars. Evidently  $I$  is the  $U(1) \times U(1)$  enveloping algebra which is an Abelian subalgebra of  $A$ . The set of new  $U(3)$  tensors in  $A$  is therefore isomorphic to the factor space  $A/I$  and we can identify  $A/I$  with the space of  $SU(3)$  basic tensors.

Factors of the type (3.7a) are easily suppressed by setting

$$p_4 = 0 \tag{3.8}$$

in Eq. (3.4); thus we restrict the set of  $p$  ( $= \sum_{\nu=1}^4 p_{\nu}$ ) level tensors to be subset with  $p = \sum_{i=1}^3 p_i$ .

The suppression of factors of the type (3.7b) can be carried out as follows. Let

$$\begin{Bmatrix} h_{v4} \\ h_{i3} \end{Bmatrix} \tag{3.9}$$

denote a  $U(3) \subset U(4)$  tensor in  $A$ . Since  $C_{ii}^{tr}$  is the  $U(3)$  scalar component of the  $su(4)$  Lie algebra of (self-conjugate) rank  $\{100-1\}$ , the tensor (3.9) is an old  $U(3)$  tensor if it can be expressed as a product

$$\begin{aligned}
&\begin{Bmatrix} 1 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 \end{Bmatrix} \\
&\times \begin{Bmatrix} h_{14} - 1 & & h_{24} & h_{34} & h_{44} + 1 \\ & h_{13} & h_{23} & h_{33} & \end{Bmatrix}.
\end{aligned} \tag{3.10}$$

Thus only the  $U(3)$  tensors  $\{h_{i3}\} \subset \{h_{v4}\}$  for which

$$h_{13} = h_{14} \quad \text{and/or} \quad h_{33} = h_{44} \tag{3.11}$$

are new.

Finally, each  $U(3)$  tensor is also an  $SU(3)$  tensor of rank

$$\{h_1 h_2\} \equiv \{h_{13} - h_{33}, h_{23} - h_{33}\}. \tag{3.12}$$

As an example, we trivially have for the  $p=0$  level the identity operator

$$T^{\{0000\}} \sim 1. \tag{3.13}$$

For the  $p=1$  level, we have the following three cases:

$$(p_i) = (100), \quad T_{lw}^{\{1100\}} \sim A_{12}, \tag{3.14}$$

$$(p_i) = (010), \quad T_{lw}^{\{00-1-1\}} \sim B_{34},$$

$$U(4) \downarrow U(3): \{00-1-1\} \downarrow \{0-1-1\} \oplus \{00-1\}; \tag{3.15}$$

$$(p_i) = (001), \quad T_{lw}^{\{100-1\}} \sim C_{14},$$

$$U(4) \downarrow U(3): \{100-1\} \downarrow \{100\} \oplus \{10-1\} \oplus \{000\} \oplus \{00-1\}. \tag{3.16}$$

Therefore

$$\begin{aligned} &\{1100\}_{\text{new}} \downarrow \{110\} \oplus \{100\}, \\ &\{00-1-1\}_{\text{new}} \downarrow \{0-1-1\} \oplus \{00-1\}, \end{aligned} \quad (3.17)$$

$$\{100-1\}_{\text{new}} \downarrow \{100\} \oplus \{10-1\} \oplus \{00-1\}.$$

The identification of all new SU(3) basic tensors for  $p = 1$  is unambiguous and their lowest weight components are given by

$$\begin{aligned} &\left\{ \begin{array}{cccc} 1 & 1 & 0 & 0 \\ & 1 & 1 & 0 \\ & & lw & \end{array} \right\} \sim A_{12}, & \left\{ \begin{array}{cccc} 0 & 0 & -1 & -1 \\ & 0 & -1 & -1 \\ & & lw & \end{array} \right\} \sim B_{23}, \\ &\left\{ \begin{array}{cccc} 1 & 1 & 0 & 0 \\ & 1 & 0 & 0 \\ & & lw & \end{array} \right\} \sim A_{14}, & \left\{ \begin{array}{cccc} 0 & 0 & -1 & -1 \\ & 0 & 0 & -1 \\ & & lw & \end{array} \right\} \sim B_{34}, \\ &\left\{ \begin{array}{cccc} 1 & 0 & 0 & -1 \\ & 1 & 0 & 0 \\ & & lw & \end{array} \right\} \sim C_{14}, & \left\{ \begin{array}{cccc} 1 & 0 & 0 & -1 \\ & 0 & 0 & -1 \\ & & lw & \end{array} \right\} \sim C_{43}, \\ &\left\{ \begin{array}{cccc} 1 & 0 & 0 & -1 \\ & 1 & 0 & -1 \\ & & lw & \end{array} \right\} \sim C_{13}, \end{aligned} \quad (3.18)$$

which, incidentally, close upon an  $hw(3)$  algebra.

Observe that altogether, there are nine SU(3) lowest weight operators, old and new, at level  $p = 1$ ;

$$C_{\mu\mu}, C_{ii}^{\text{tr}}, A_{12}, A_{14}, B_{34}, B_{23}, C_{14}, C_{43}, C_{13}. \quad (3.19)$$

In Biedenharn and Flath's  $SO(6,2) \sim SO^*(8)$  realization, the two old  $p = 1$  SU(3) [also SU(3) invariant] tensors are given, to within additive constants, by their "quark" and "antiquark" number operators (see also Bracken and McGibbon<sup>2</sup>). Ignoring  $C_{\mu\mu}$  and  $C_{ii}^{\text{tr}}$ , there are seven new SU(3) operators at the  $p = 1$  level.

Consider now the  $p = 2$  level. Since there are a total of nine  $p = 1$  SU(3) tensors, we expect, at first sight, a total of  $9 \times 10/2 = 45$  (old and new) symmetrical  $p = 2$  SU(3) tensors. In fact, listing the complete set of  $U(4) \supset U(3)$   $p = 2$  tensors using Proposition 3.3 and the betweenness conditions of the corresponding Gel'fand patterns, we find only 43 tensors. This implies the existence of two linear relationships between the SU(3) quadratic tensors. Fortunately, there is no need to seek these relationships as the 43 linearly independent tensors are obtained unambiguously from their corresponding U(4) parents [see Eqs. (4.7) and (4.9) below]. This obviates the necessity of imposing constraints, to restrict a linearly independent set, as in Biedenharn and Flath's analysis. Since our analysis is based on the  $U(4) \downarrow U(3)$  branching rules, we automatically obtain linearly independent tensors and thus confirm a previous conjecture (Le Blanc and Rowe<sup>4</sup>) that a primary classification of  $SO^*(8)$  tensors under U(4) would eliminate any ambiguities in the identification of the SU(3) tensors in  $A$ .

Applying the restrictions (3.8) and (3.11), we get the following new SU(3) tensors for the  $p = 2$  level:

$$\begin{aligned} (p_i) = (200): & \{2200\}_{\text{new}} \downarrow \{220\} \oplus \{210\} \oplus \{200\}, \\ (p_i) = (110): & \{11-1-1\}_{\text{new}} \downarrow \{11-1\} \oplus \{10-1\} \oplus \{1-1-1\}, \\ (p_i) = (101): & \{210-1\}_{\text{new}} \downarrow \{210\} \oplus \{21-1\} \oplus \{200\} \oplus \{20-1\} \oplus \{11-1\} \oplus \{10-1\}, \\ (p_i) = (020): & \{00-2-2\}_{\text{new}} \downarrow \{00-2\} \oplus \{0-1-2\} \oplus \{0-2-2\}, \\ (p_i) = (011): & \{10-1-2\}_{\text{new}} \downarrow \{10-1\} \oplus \{10-2\} \oplus \{1-1-1\} \oplus \{1-1-2\} \oplus \{00-2\} \oplus \{0-1-2\}, \\ (p_i) = (002): & \{200-2\}_{\text{new}} \downarrow \{200\} \oplus \{20-1\} \oplus \{20-2\} \oplus \{10-2\} \oplus \{00-2\}. \end{aligned}$$

There are thus 26 new  $p = 2$  SU(3) tensors (as in Biedenharn and Flath<sup>1</sup>).

The classification of SU(3) tensors at any  $p$  level is thus very straightforward.

#### IV. SU(3) SHIFT PROPERTIES OF THE $SO^*(8) \supset U(4) \supset U(3)$ TENSORS

*Proposition 4.1:* The  $SO^*(8) \supset U(4) \supset U(3)$  tensors (3.9) of the factor space  $A/I$  have well-defined U(3) shift properties given by

$$\delta_1 = p_1 - p_2 = h_{14} + h_{44} = h_{24} + h_{34}, \quad (4.1a)$$

$$\delta_2 = \sum_i h_{i3} - \delta_1. \quad (4.1b)$$

*Proof:* First, note that when acting on a state (2.5) of the model space, the tensor (3.9)

$$\left\{ \begin{array}{c} h_{\mu^4} \\ h_{i3} \end{array} \right\} = \left\{ \begin{array}{cccc} p_1 + p_3 & p_1 & -p_2 & -p_2 - p_3 \\ & h_{13} & h_{23} & h_{33} \end{array} \right\}$$

has a  $\delta_1$  shift given by the difference between its number of Cartan raising operators [the number of  $A$ 's in (3.4)] and its number of Cartan lowering operators [the number of  $B$ 's in (3.4)],

$$\delta_1 = p_1 - p_2.$$

Then, note that the U(4) weight operator  $C_{44}$  has eigenvalue  $h_1 - h_2$ ,



$$C_{44} \left| \begin{array}{cccc} h_1 & h_1 & 0 & 0 \\ & h_1 & h_2 & 0 \\ & & \eta & 0 \\ & & & \eta \end{array} \right\rangle$$

$$= (h_1 - h_2) \left| \begin{array}{cccc} h_1 & h_1 & 0 & 0 \\ & h_1 & h_2 & 0 \\ & & \eta & 0 \\ & & & \eta \end{array} \right\rangle,$$

on a model space state. There is therefore no  $C_{44}$  weight multiplicity in the model space as a  $U(3)$  invariant subspace is unambiguously identified by the  $U(4)$  label  $\{h_1, h_1, 00\}$  and the  $C_{44}$  weight  $(h_1 - h_2)$ .

Now the tensor (3.9) has a  $C_{44}$  weight given by

$$\sum_{\mu} h_{\mu 4} - \sum_i h_{i3} = 2(p_1 - p_2) - \sum_i h_{i3}.$$

The tensor will therefore take a  $U(3)$  irreducible subspace  $\{h_1, h_1, 00\} \supset \{h_1, h_2, 0\}$  to a new  $U(3)$  irreducible subspace with  $U(4) \supset U(3)$  labels given by

$$\{h_1 + \delta_1, h_1 + \delta_1, 0, 0\} \supset \{h_1 + \delta_1, h_2 + \delta_2, 0\},$$

where, from the additivity of the  $C_{44}$  weight,

$$h'_1 - h'_2 = h_1 - h_2 + \sum_{\mu} h_{\mu 4} - \sum_i h_{i3},$$

which leads to the result

$$\delta_2 = \sum_i h_{i3} - \delta_1$$

[cf. Eq. (2.14)]. We thus conclude that the tensor (3.9) has unambiguous  $U(3)$  shift properties.

Q.E.D.

We now proceed to prove that the tensors (3.9), classified in the last section according to their  $SO^*(8) \supset U(4) \supset U(3)$  properties, can be put in one-to-one correspondence with the abstract set (2.15) of  $U(3)$  shift tensors.

We first consider, as in Biedenharn and Flath,<sup>1</sup> the subset of  $p$  level tensors of  $SU(3)$  rank  $\{h_1 = p, h_2\}$  and shifts  $(\delta_1, \delta_2)$ . [Recall from (2.18) that an  $SU(3)$  tensor can be identified by either  $(h_{13}, h_{23}, h_{33}, \delta_1, \rho)$  or  $(h_1, h_2, \delta_1, \delta_2, \rho)$  with

$$h_1 = h_{13} - h_{33}, \quad h_2 = h_{23} - h_{33}, \quad \sum_i h_{i3} = \delta_1 + \delta_2.]$$

From (4.1a) and condition  $h_1 = p$ , we have two equations for  $(p_1, p_2, p_3)$ , namely,

$$p_1 + p_2 + p_3 = h_{13} - h_{33}, \quad p_1 - p_2 = \delta_1. \quad (4.2)$$

The third equation is given by the constraint (3.11). We have either

$$h_{13} = h_{14} = p_1 + p_3 \Rightarrow h_{33} = h_{34} = -p_2 \quad (4.3a)$$

or

$$h_{33} = h_{44} = -p_2 - p_3 \Rightarrow h_{13} = h_{24} = p_1. \quad (4.3b)$$

Solving for the  $p_i$ 's, we obtain a unique solution for a given  $\{h_1, h_2, \delta_1, \delta_2, \rho = 0\}$  which we identify with the stretched ( $\rho = 0$ ) tensors of the abstract set (2.15) of shift tensors by the following proposition.

*Proposition 4.4:* To every  $SU(3)$  tensor of rank  $\{h_1, h_2\}$ , of given shift  $(\delta_1, \delta_2)$  and with stretched operator pattern, we can associate a  $p^s = h_1$  level tensor

$$\left\{ \begin{array}{cccc} h_{\mu 4} \\ h_{i3} \end{array} \right\} = \left\{ \begin{array}{cccc} p_1 + p_2 & p_1 & -p_2 & -p_2 - p_3 \\ & h_{13} & h_{23} & h_{33} \end{array} \right\},$$

with  $U(4) \supset U(3)$  labels given by

$$\left\{ \begin{array}{cccc} h_{13} & \delta_1 - h_{33} & h_{33} & -(h_{13} - \delta_1) \\ & h_{13} & h_{23} & h_{33} \end{array} \right\}, \quad (4.4a)$$

if  $\delta_1 - h_{33} \leq h_{13}$  and by

$$\left\{ \begin{array}{cccc} \delta_1 - h_{33} & h_{13} & -(h_{13} - \delta_1) & h_{33} \\ & h_{13} & h_{23} & h_{33} \end{array} \right\}, \quad (4.4b)$$

if  $\delta_1 - h_{33} > h_{13}$ .

Now, for a  $(p = h_1 - \rho)$  level  $SO^*(8) \supset U(4) \supset U(3)$  tensor (3.9) with fixed  $SU(3)$  quantum labels  $\{h_1, h_2\}$  to also have fixed  $U(3)$  shift properties  $(\delta_i)$ , we must have, from (3.11) and (4.1),

$$p_1(\rho) = p_1^s - \rho, \quad p_2(\rho) = p_2^s - \rho, \quad p_3(\rho) = p_3^s + \rho. \quad (4.5)$$

We thus deduce the following proposition.

*Proposition 4.6:* A multiplicity set of  $SU(3)$  tensors (2.5) will have a unique representative on the  $p$  ( $= p_1 + p_2 + p_3$ ) levels starting with the first (stretched) tensor on the level  $p = p^s = h_1$  and the following tensors, indexed by the integer  $\rho$  [Eq. (2.17)], on the consecutive  $p = p^s - \rho$  levels. The corresponding  $U(4) \supset U(3)$  ( $\{h_{\mu 4}\} \supset \{h_{i3}\}$ ) labels for these representatives are given by

$$\left\{ \begin{array}{cccc} h_{13} & \delta_1 - h_{33} - \rho & h_{33} + \rho & -(h_{13} - \delta_1) \\ & h_{13} & h_{23} & h_{33} \end{array} \right\}, \quad (4.6a)$$

if  $\delta_1 - h_{33} \leq h_{13}$  and by

$$\left\{ \begin{array}{cccc} \delta_1 - h_{33} & h_{13} - \rho & -(h_{13} - \delta_1) + \rho & h_{33} \\ & h_{13} & h_{23} & h_{33} \end{array} \right\}, \quad (4.6b)$$

if  $\delta_1 - h_{33} \geq h_{13}$ .

It can be verified that for  $\rho \leq \rho_{\max}$ , all labels are fully consistent with the betweenness conditions of all [upper and lower, U(4) and U(3)] Gel'fand patterns involved. Propositions 4.4 and 4.6 are equivalent to Theorem 8.12 of Biedenharn and Flath.<sup>1</sup>

Acting with the raising operator  $C_{42}$  on the expression (3.4) for the U(4)  $lw$  tensor component, we find that the lowest weight component of the shift tensor (4.4a) is given by

$$C_{42}^{(\delta_1 - h_{33} - h_{23})} \{ A_{12}^{\delta_1 - h_{33}} B_{34}^{-h_{33}} \} C_{14}^{h_{13} - \delta_1 + h_{33}}, \quad (4.7)$$

where

$$\begin{aligned} C_{42}^{(0)} \{ X \} &= X, \\ C_{42}^{(1)} \{ X \} &= [C_{42}, X], \\ C_{42}^{(2)} \{ X \} &= [C_{42}, [C_{42}, X]], \quad \text{etc.} \end{aligned} \quad (4.8)$$

Similarly, acting with  $C_{42}$  and  $C_{43}$ , we find that the lowest weight component of the shift tensor (4.6a) is given by

$$C_{42}^{(\delta_1 - h_{33} - h_{23} - \rho)} \{ A_{12}^{\delta_1 - h_{33} - \rho} B_{34}^{-h_{33} - \rho} \} C_{14}^{h_{13} - \delta_1 + h_{33}} C_{13}^{\rho}. \quad (4.9)$$

Similar expressions can be obtained for the tensors (4.4b) and (4.6b). Equations (4.4), (4.6), (4.7), and (4.9) are also

equivalent to Eqs. (8.21) and (8.22) of Biedenharn and Flath. The group theoretical structure of the tensors (4.9) guarantees their linear independence (Draayer and Akiyama,<sup>7</sup> Le Blanc<sup>8</sup>).

Finally, note that although extremely elegant and simple, the  $SO^*(8) \sim SO(6,2)$  model for SU(3) does not offer (as yet) a group theoretical interpretation for the U(2) aspects of the upper pattern. Such an interpretation has been given elsewhere by the authors (Le Blanc and Rowe<sup>9</sup>) using an alternative model for SU(3) defined in a  $U(2) \times SU(3)$  Bargmann space.

<sup>1</sup>L. C. Biedenharn and D. E. Flath, *Commun. Math. Phys.* **93**, 143 (1984).

<sup>2</sup>A. J. Bracken and J. H. MacGibbon, *J. Phys. A: Math. Gen.* **17**, 2581 (1984); A. J. Bracken, *Commun. Math. Phys.* **94**, 377 (1984).

<sup>3</sup>I. N. Bernštein, I. M. Gel'fand, and S. I. Gel'fand, *Funct. Anal. Appl.* **9**, 332 (1975).

<sup>4</sup>R. Le Blanc and D. J. Rowe, *J. Phys. A: Math. Gen.* **19**, 1111 (1986).

<sup>5</sup>J. D. Louck, *Am. J. Phys.* **38**, 3 (1970).

<sup>6</sup>R. Le Blanc and D. J. Rowe, *J. Phys. A: Math. Gen.* **18**, 1891, 1905 (1985).

<sup>7</sup>J. P. Draayer and Y. Akiyama, *J. Math. Phys.* **14**, 1904 (1973).

<sup>8</sup>R. Le Blanc, Ph.D. thesis, University of Toronto, 1985.

<sup>9</sup>R. Le Blanc and D. J. Rowe, *J. Phys. A: Math. Gen.* **19**, 1093, 2913 (1986).

# Indecomposable modules of the Poincaré algebra in an energy-cyclic angular momentum basis

Romuald Lenczewski

*Department of Mathematics, Southern Illinois University, Carbondale, Illinois 62901*

(Received 7 August 1986; accepted for publication 11 February 1987)

The universal enveloping algebra  $\mathcal{U} = \mathcal{U}_+ \mathcal{U}_- \mathcal{H}$  of the Poincaré algebra  $\text{iso}(3,1)$  is considered. Infinite-dimensional induced modules  $\mathcal{U}_+(\Gamma)$  are studied. Explicit formulas are obtained for  $\mathcal{U}_+(\Gamma)$  in the Poincaré–Birkhoff–Witt basis. Then a change of basis is performed to an energy-cyclic angular momentum basis. For any  $\Gamma \in \mathbb{C}^2$ ,  $\mathcal{U}_+(\Gamma)$  turns out to be indecomposable. It can be represented as an infinite family of interacting Verma modules of the Lorentz algebra  $\text{so}(3,1)$ . Finite-dimensional modules can be obtained as quotient modules of  $\mathcal{U}_+(\Gamma)$  if  $2\Gamma \in \mathbb{Z}^2$ .

## I. INTRODUCTION

Indecomposable modules of physically relevant Lie algebras have been suggested<sup>1-3</sup> for the description of unstable particles. In particular, indecomposable modules of the Poincaré algebra  $\mathcal{P} = \text{iso}(3,1)$  may be useful in modeling unstable particles of arbitrary spin and their interactions. However, the Poincaré algebra, being a nonsemisimple Lie algebra, does not lend itself as easily to methods of the representation theory as semisimple Lie algebras (even in the case of finite-dimensional modules). On the other hand, its maximal solvable subalgebra is Abelian, hence it can be viewed as mathematically relatively easy to work with.

In fact,  $\mathcal{P}$  is a semidirect product of the algebra of translations  $\mathcal{H} = \mathfrak{t}(4)$  and the Lorentz algebra  $\mathcal{L} = \text{so}(3,1)$ . The Lie products satisfy  $[\mathcal{L}, \mathcal{L}] = \mathcal{L}$ ,  $[\mathcal{L}, \mathcal{H}] = \mathcal{H}$ ,  $[\mathcal{H}, \mathcal{H}] = 0$ . Therefore  $\mathcal{H}$  is an Abelian ideal in  $\mathcal{P}$  and this makes the analysis simpler than for other nonsemisimple Lie algebras (except, maybe, some low-dimensional ones).

Hence with both mathematical and physical motivations we resume<sup>4</sup> the study of the indecomposable modules of  $\mathcal{P}$  (infinite dimensional as well as finite dimensional). This work is a natural extension of the study conducted by Gruber and Lenczewski.<sup>4,5</sup> Thus the results contained therein will be basic for the analysis that is carried out in this paper. Moreover, the modules considered previously<sup>4,5</sup> turn out to be incorporated into the present study as special cases. Therefore to a great extent we tried to use the same notation, although in several places it has been changed to avoid confusion.

In Sec. II we define the Poincaré algebra and introduce the raising and lowering operators.

In Sec. III we define the universal enveloping algebra of  $\mathcal{P}$ ,  $\mathcal{U}(\mathcal{P}) = \mathcal{U}_+ \mathcal{U}_- \mathcal{H}$ , as a regular left  $\mathcal{P}$  module. The induced modules  $\mathcal{U}_+(\Gamma)$  (the analogs of Verma<sup>6</sup> modules of semisimple Lie algebras) are explicitly derived in the Poincaré–Birkhoff–Witt (PBW) basis. Although the PBW basis seems to be very natural, a basis more useful in physical applications, especially in the case of the chain  $\text{so}(3) \subset \text{so}(3,1) \subset \text{iso}(3,1)$ , is an angular momentum basis<sup>5</sup> (AMB). Therefore we need to introduce an AMB to the formalism.

In Sec. IV a certain AMB is defined. It turns out that it is

not unique (in opposition to the Lorentz algebra case<sup>5</sup> or some less general modules of the Poincaré algebra<sup>4</sup>). We choose a certain energy-cyclic AMB which allows us to obtain the modules  $\mathcal{U}_+(\Gamma)$  in a closed form.

In Sec. V we give explicit formulas for  $\mathcal{U}_+(\Gamma)$  in the energy-cyclic AMB mentioned above. They are derived by using the induction method. Modules  $\mathcal{U}_+(\Gamma)$  turn out to be indecomposable for all  $\Gamma \in \mathbb{C}^2$ . Each  $\mathcal{U}_+(\Gamma)$  forms, in fact, a composition series. Moreover, when treated as a vector space, it can be written as

$$\mathcal{U}_+(\Gamma) = \sum_{\substack{M, m \in \mathbb{N} \\ m < M}} \mathcal{U}_M^m$$

such that each  $\mathcal{U}_M^m$  is a Verma  $\mathcal{L}$  module. The module  $\mathcal{U}_+(\Gamma)$  can be viewed as a sum of interacting Verma  $\mathcal{L}$  modules.

In Sec. VI we give some examples of finite-dimensional indecomposable modules as quotient modules of  $\mathcal{U}_+(\Gamma)$  [or its extension  $\tilde{\mathcal{U}}_+(\Gamma)$ ] for certain  $2\Gamma \in \mathbb{Z}^2$ . Bases for the modules of dimensions 4, 5, 7, 8, 8' (two different eight dimensional ones) are given as well as their  $\text{so}(3) \subset \text{so}(3,1)$  decomposition. Thus we identify all but one of the finite-dimensional indecomposable modules of dimensions less than 9 according to their recent classification for low dimensions.<sup>7</sup>

The main result of this paper (contained in Sec. V) is of significance in its own right taking into account that the analogous result for the Lorentz algebra<sup>5</sup> reproduced the Gel'fand–Naimark<sup>8,9</sup> representations as a special case. At the same time it provides an important step forward in our quest for a unified approach to finite-dimensional indecomposable modules of  $\mathcal{P}$ .

In the sequel we use the abbreviations introduced above as well as the following: inf(fin)-in(ir)-mod [infinite(finite) dimensional indecomposable (irreducible) module]. All algebras are considered to be over  $\mathbb{C}$ . For general references one is referred to Humphreys<sup>10</sup> and Dixmier.<sup>11</sup>

## II. PRELIMINARIES

The defining relations for the Poincaré algebra  $\mathcal{P}$  are usually given as the following set of Lie products:

$$\begin{aligned}
[P_\mu, P_\nu] &= 0, \\
[M_{\mu\nu}, P_\alpha] &= i(g_{\mu\alpha}P_\nu - g_{\nu\alpha}P_\mu), \\
[M_{\mu\nu}, M_{\alpha\beta}] &= i(g_{\mu\alpha}M_{\nu\beta} - g_{\nu\alpha}M_{\mu\beta} + g_{\mu\beta}M_{\alpha\nu} - g_{\nu\beta}M_{\alpha\mu}),
\end{aligned} \tag{1}$$

where  $M_{\mu\nu}, P_\alpha$  represent infinitesimal generators of rotations and translations, respectively, in the four-dimensional Minkowski space  $[\mu, \nu, \alpha \in \{0, 1, 2, 3\}]$  and  $g_{\mu\nu} = \text{diag}(-1, 1, 1, 1)$ . In our formalism it appears to be necessary to introduce another basis for  $\mathcal{P}$ , consisting of root vectors, vectors corresponding to positive roots (raising operators), and vectors corresponding to negative roots (lowering operators). The root vectors span the Cartan subalgebra of  $\mathcal{P}$ , i.e., the nilpotent subalgebra that equals its normalizer in (it is equal to the Cartan subalgebra of  $\mathcal{L}$ ). They can be chosen as  $h_3 = M_{12}$  and  $f_3 = -iM_{03}$ . Raising operators (with a subscript  $+$ ) and lowering operators (with a subscript  $-$ ) can be chosen as

$$\begin{aligned}
h_+ &= M_{23} + iM_{13}, & h_- &= M_{23} - iM_{13}, \\
f_+ &= -M_{02} - iM_{01}, & f_- &= M_{02} - iM_{01}, \\
k_+ &= P_1 - iP_2, & k_- &= P_1 + iP_2, \\
r_+ &= P_3 + P_0, & r_- &= P_3 - P_0.
\end{aligned} \tag{2}$$

The advantage of using this basis stems from the fact that all structure constants lie in  $\mathbb{Z}$  as can be seen below. In the modules investigated by Lenczewski and Gruber<sup>4</sup> the following basis was used:  $p_+ = if_+, p_- = if_-, p_3 = if_3$  and  $k_3 = P_3, k_0 = -iP_0$ .

Introducing  $r_+$  and  $r_-$  allows us to extend the modules obtained therein as will be seen throughout the paper. For notational convenience we are going to use  $\mathcal{R}$  to represent  $\text{so}(3)$  and  $\mathcal{F}$  to represent the  $\mathbb{C}$ -linear subspace spanned by  $f_+, f_-, f_3$  (boosts). Clearly,  $\mathcal{F}$  is not a subalgebra of  $\mathcal{P}$  contrary to  $\mathcal{X}$  and  $\mathcal{L}$ .

Using the above notation one can derive the following commutation relations in the new basis:

$$\begin{aligned}
[\mathcal{R}, \mathcal{R}] & \begin{cases} [h_3, h_\pm] = \pm h_\pm, \\ [h_+, h_-] = 2h_3, \end{cases} \\
[\mathcal{F}, \mathcal{F}] & \begin{cases} [f_+, f_-] = 2h_3, \\ [f_3, f_\pm] = \pm h_\pm, \end{cases} \\
[\mathcal{R}, \mathcal{F}] & \begin{cases} [h_3, f_\pm] = \pm f_\pm, \\ [f_\pm, h_\mp] = \pm 2f_3, \\ [f_3, h_\pm] = \pm f_\pm, \end{cases} \\
[\mathcal{F}, \mathcal{X}] & \begin{cases} [f_-, r_\pm] = \pm k_-, \\ [f_+, r_\pm] = \pm k_+, \\ [f_\mp, k_\pm] = -r_- + r_+, \\ [f_3, r_\pm] = \pm r_\pm, \end{cases} \\
[\mathcal{R}, \mathcal{X}] & \begin{cases} [h_3, k_\pm] = \pm k_\pm, \\ [h_\pm, k_\mp] = \pm (r_- + r_+), \\ [h_+, r_\pm] = -k_+, \\ [h_-, r_\pm] = k_-. \end{cases}
\end{aligned} \tag{3}$$

It is understood that either all upper or all lower signs (whether as subscripts or arithmetic symbols) are taken simultaneously.

### III. INDUCED MODULES $\mathcal{U}_+(\Gamma)$ IN THE PBW BASIS

We come to the main object of interest: the universal enveloping algebra (UEA) of the algebra  $\mathcal{P}$ ,  $\mathcal{U}(\mathcal{P})$ . It is defined as the quotient algebra  $\mathcal{U}(\mathcal{P}) = \mathcal{T}(\mathcal{P})/\mathcal{I}$ , where  $\mathcal{T}(\mathcal{P})$  is the tensor algebra of  $\mathcal{P}$  and  $\mathcal{I}$  is the ideal generated by  $x \otimes y - y \otimes x - [x, y]$ ,  $x, y \in \mathcal{P}$ . It is well known<sup>10,11</sup> that there exists a basis for  $\mathcal{U}(\mathcal{P})$  that consists of standard monomials. This is the Poincaré–Birkhoff–Witt (PBW) basis.

The left tensor multiplication by  $\mathcal{P}$  makes out of  $\mathcal{U}(\mathcal{P})$  a regular left module. Various quotients of  $\mathcal{U}(\mathcal{P})$  can be considered. For example, for given  $\Gamma \in \mathbb{C}^2$  one can define a module  $\mathcal{U}_+(\Gamma)$  as a quotient module  $\mathcal{U}(\mathcal{P})/\mathcal{J}$ , where the ideal  $\mathcal{J}$  is generated by lowering operators and  $h_3 - \Gamma_1, f_3 - \Gamma_2$ . Similarly,  $\mathcal{U}_-(\Gamma)$  is a quotient module  $\mathcal{U}(\mathcal{P})/\mathcal{G}$  where the ideal  $\mathcal{G}$  is generated by raising operators and  $h_3 - \Gamma_1, f_3 - \Gamma_2$ . We will restrict our attention to modules  $\mathcal{U}_+(\Gamma)$  later on giving a connection between  $\mathcal{U}_+(\Gamma)$  and  $\mathcal{U}_-(\Gamma)$  through an automorphism of  $\mathcal{P}$ . The following relations are obtained on  $\mathcal{U}_+(\Gamma)$ :

$$\begin{aligned}
h_3 X &= (\Gamma_1 + n + s + p)X, & h_+ X &= X(n + 1), \\
f_+ X &= X(s + 1), & k_+ X &= X(p + 1), \\
h_- X &= n(-2\Gamma_1 - 2p - 2s - n + 1)X(n - 1) \\
&\quad - 2s(v + \Gamma_2)X(s - 1) \\
&\quad - s(s - 1)X(s + 1, n - 2) + pX(p - 1, v + 1), \\
f_- X &= s(-2\Gamma_1 - 2p - s + 1 - 2n)X(s - 1) \\
&\quad + 2n(-v - \Gamma_2)X(n - 1) \\
&\quad - n(n - 1)X(s + 1, n - 2) + pX(p - 1, v + 1),
\end{aligned} \tag{4}$$

$$\begin{aligned}
f_3 X &= (\Gamma_2 + v)X + nX(s + 1, n - 1) + sX(n + 1, s - 1), \\
k_- X &= -sX(s - 1, v + 1) - nX(n - 1, v + 1) \\
&\quad + s(s - 1)X(s - 2, p + 1) \\
&\quad - n(n - 1)X(n - 2, p + 1), \\
r_- X &= sX(s - 1, p + 1) + nX(n - 1, p + 1), \\
r_+ X &= X(v + 1) - sX(s - 1, p + 1) + nX(n - 1, p + 1),
\end{aligned}$$

where  $X = X(n, s, p, v) = h_+^n f_+^s k_+^p r_+^v$ ,  $n, s, p, v \in \mathbb{N}$ , and only the parameters that are altered by the action of the operator on the left-hand side of each equation are indicated. The basis elements represent the natural basis induced by the PBW basis of the UEA, hence we refer to it also as a PBW basis of  $\mathcal{U}_+(\Gamma)$ .

The above equations show that  $\mathcal{U}_+(\Gamma)$  is an inf-in-mod for all  $\Gamma \in \mathbb{C}^2$ . The value of  $v$  does not decrease. Hence, on quotients  $\mathcal{U}_+(\Gamma)/\mathcal{O}_i$ , where the ideal  $\mathcal{O}_i$  is generated by  $r_+^i$ , one obtains induced relations similar to those exhibited by Eq. (4). In particular, if we consider the ideal  $\mathcal{O}_1$  we obtain the modules  $\Omega_+(\Gamma)$  (Ref. 4) recently examined. Presently, however, the structure under consideration is richer and will give new results. The next step consists in performing a change of basis from the PBW basis into an AMB.

#### IV. AN ENERGY-CYCLIC ANGULAR MOMENTUM BASIS

In this section we are going to concentrate on finding an AMB for  $\mathcal{U}_+(\Gamma)$ . We will use the fact that  $[h_-, k_0] = 0$  to generate an AMB that is energy cyclic.

To find an AMB we define an  $\mathcal{R}$ -maximal ( $\mathcal{R}$ -minimal) vector to be an eigenvector  $Y$  of  $h_3$  that satisfies  $h_+ Y = 0$  ( $h_- Y = 0$ ). They generate cyclic  $\mathcal{R}$  modules, whose bases will span  $\mathcal{P}$  modules  $\mathcal{U}_+(\Gamma)$  [ $\mathcal{U}_-(\Gamma)$ ].

Let  $\mathcal{Y} = \{Y_j, j \in J\}$  be the set of all  $\mathcal{R}$ -minimal vectors, where  $J$  is a certain index set (or a multi-index set). Then,  $\text{AMB} \subset \{h_+^n \mathcal{Y}\} = \{h_+^n Y_j, j \in J\}$ . In the case of the Lorentz algebra<sup>5</sup> or the modules  $\Omega_+(\Gamma)$  of the Poincaré algebra,<sup>4</sup>  $\text{AMB} = \{h_+^n \mathcal{Y}\}$ . However, in  $\mathcal{U}_+(\Gamma)$  we have  $\text{AMB} \neq \{h_+^n \mathcal{Y}\}$ . In fact,  $\text{card}(\text{AMB}) = \aleph_0$ , whereas  $\text{card} \{h_+^n \mathcal{Y}\} = c$ . This leads to the nonuniqueness of the AMB. In this paper we shall choose an AMB which is energy cyclic.

One arrives at the following form of minimal vectors:

$$\{Y_{NM}\} = \left\{ \sum_{\substack{s+p < N \\ p < M}} c_{sp}^{NM} X(N-s-p, s, p, M-p) \right\},$$

where  $N > 0$ ,  $M > 0$ , and  $h_3 Y_{NM} = (\Gamma_1 + N) Y_{NM}$ . The constants  $c_{sp}^{NM}$  satisfy the following recurrence relation:

$$\begin{aligned} (N-s-p)(-2\Gamma_1-p-s-N+1)c_{sp}^{NM} \\ - 2(M+\Gamma_2-p)(s+1)c_{s+1,p}^{NM} \\ - (s+2)(s+1)c_{s+2,p}^{NM} - (p+1)c_{s,p-1}^{NM} = 0. \end{aligned} \quad (5)$$

Because of the abundance of solutions of the recurrence relation we need to choose certain representatives of the sets of  $\mathcal{R}$ -minimal vectors in the form given. We are going to do it in a fairly natural manner. We will try to reduce the problem to the Lorentz algebra case (where the solutions are unique and generate the bases of  $\mathcal{L}$  modules) in a sufficient number of cases to generate the whole set of  $\mathcal{R}$ -minimal vectors that will be needed by applying the energy operator.

Let us set  $M = 0$  first. This is the prototypical Lorentz algebra case. We obtain the following recurrence relation:

$$\begin{aligned} (N-s)(-2\Gamma_1-s-N+1)c_{s0}^{N0} - 2\Gamma_2(s+1)c_{s+1,0}^{N0} \\ - (s+2)(s+1)c_{s+2,0}^{N0} = 0. \end{aligned} \quad (6)$$

Then  $Y_{N0}$  can be given as

$$Y_{N0} = \sum_{s < N} c_{s0}^{N0} X(N-s, s, 0, 0) \quad (7)$$

and the constants  $c_{s0}^{N0}$  can be uniquely determined. Their explicit form is fairly complicated and will not be needed in the sequel since our derivation of the modules  $\mathcal{U}_+(\Gamma)$  in the new basis proceeds by induction.

Let us set  $N = 0$  now. Then the only possible  $\mathcal{R}$ -minimal vector is  $X(0, 0, 0, M)$  (up to a constant factor, of course). Thus we let

$$Y_{0M} = X(0, 0, 0, M). \quad (8)$$

Now, if  $N \neq 0$  and  $M \neq 0$  then the solutions for  $\mathcal{R}$ -minimal vectors span at least a two-dimensional vector space. We choose the representatives of the sets  $\{Y_{NM}\}$  in the following way:

$$Y_{NM} = \sum_{s < N} c_{s0}^{NM} X(N-s, s, 0, M), \quad (9)$$

with  $c_{s0}^{NM}$  satisfying the recurrence relation:

$$\begin{aligned} (N-s)(-2\Gamma_1-s-N+1)c_{s0}^{NM} \\ - 2(M+\Gamma_2)(s+1)c_{s+1,0}^{NM} \\ - (s+2)(s+1)c_{s+2,0}^{NM} = 0. \end{aligned} \quad (10)$$

Thus  $Y_{NM}$  have exactly the same form as  $Y_{N0}$  except that  $\Gamma_2$  is replaced by  $\Gamma_2 + M$ . Hence we will obtain in this manner  $\mathcal{R}$ -minimal vectors for  $\mathcal{L}$  modules.

Other  $\mathcal{R}$ -minimal vectors are defined inductively using the energy operator as the generator. Thus we obtain the following set of  $\mathcal{R}$ -minimal vectors:

$$Y_{NM}^0 = Y_{NM}, \quad Y_{NM}^m = (2ik_0)^m Y_{N,M-m}^0, \quad (11)$$

where  $m \leq \min(N, M)$ . Now,

$$\begin{aligned} \text{AMB} = \{h_+^n Y_{NM}^m = Y_{NM}^{nm}, n, m, N, M \in \mathbb{N}, \\ m \leq \min(N, M)\}. \end{aligned}$$

We call the above four parameter basis an energy-cyclic AMB since it is generated by the energy operator  $k_0$ . This basis will allow us to obtain the modules  $\mathcal{U}_+(\Gamma)$  in a closed form.

#### V. MODULES $\mathcal{U}_+(\Gamma)$ IN AN ENERGY CYCLIC ANGULAR MOMENTUM BASIS

In this section we shall present the main result of this article. It will consist of the formulas for the modules  $\mathcal{U}_+(\Gamma)$  in the basis defined in Sec. IV. They are obtained by the induction method and their derivation involves fairly extensive calculations. One obtains

$$h_+ Y_{NM}^{nm} = Y_{NM}^{n+1,m}, \quad (12)$$

$$h_- Y_{NM}^{nm} = n F_N^n Y_{NM}^{n-1,m}, \quad (13)$$

$$h_3 Y_{NM}^{nm} = G_N^n Y_{NM}^{nm}, \quad (14)$$

$$\begin{aligned} f_+ Y_{NM}^{nm} = \left( \frac{N-m+1}{N+1} \right) Y_{N+1,M}^{nm} + \left( \frac{m}{N+1} \right) Y_{N+1,M}^{n,m-1} + \beta_{NMm} Y_{NM}^{n+1,m} \\ + \beta_{NMm}^1 Y_{NM}^{n+1,m-1} + \alpha_{NMm} Y_{N-1,M}^{n+2,m} + \alpha_{NMm}^1 Y_{N-1,M}^{n+2,m-1}, \end{aligned} \quad (15)$$

$$f_- Y_{NM}^{nm} = - \left( \frac{n(n-1)(N-m+1)}{N+1} \right) Y_{N+1,M}^{n-2,m} + \left( - \frac{n(n-1)m}{N+1} \right) Y_{N+1,M}^{n-2,m-1} + n F_N^n \beta_{NMm} Y_{NM}^{n-1,m}$$

$$+ nF_N^n \beta_{NMm}^{\downarrow} Y_{NM}^{n-1,m-1} - F_N^n (F_N^n + 1) \alpha_{NMm} Y_{N-1,M}^{nm} - F_N^n (F_N^n + 1) \alpha_{NMm}^{\downarrow} Y_{N-1,M}^{n,m-1}, \quad (16)$$

$$f_3 Y_{NM}^{nm} = \left( \frac{n(N-m+1)}{N+1} \right) Y_{N+1,M}^{n-1,m} + \left( \frac{nm}{N+1} \right) Y_{N+1,M}^{n-1,m-1} + G_N^n \beta_{NMm} Y_{NM}^{nm} \\ + G_N^n \beta_{NMm}^{\downarrow} Y_{NM}^{nm} - F_N^n \alpha_{NMm} Y_{N-1,M}^{n+1,m} - F_N^n \alpha_{NMm}^{\downarrow} Y_{N-1,M}^{n+1,m-1}, \quad (17)$$

$$k_0 Y_{NM}^{nm} = (-i/2) Y_{N,M+1}^{n,m+1}, \quad (18)$$

$$k_3 Y_{NM}^{nm} = \left( -\frac{n}{2(N+1)} \right) Y_{N+1,M+1}^{n-1,m+1} + \left( \frac{n}{2(N+1)} \right) Y_{N+1,M+1}^{n-1,m} + \frac{1}{2} G_N^n \beta_{NMm}^{\uparrow\uparrow} Y_{N,M+1}^{n,m+1} \\ + \frac{1}{2} G_N^n \beta_{NMm}^{\downarrow} Y_{N,M+1}^{nm} - \frac{1}{2} F_N^n \alpha_{NMm}^{\uparrow\uparrow} Y_{N-1,M+1}^{n+1,m+1} - \frac{1}{2} F_N^n \alpha_{NMm}^{\downarrow} Y_{N-1,M+1}^{n+1,m}, \quad (19)$$

$$k_+ Y_{NM}^{nm} = \left( -\frac{1}{2(N+1)} \right) Y_{N+1,M+1}^{n,m+1} + \left( \frac{1}{2(N+1)} \right) Y_{N+1,M+1}^{nm} + \frac{1}{2} \beta_{NMm}^{\uparrow\uparrow} Y_{N,M+1}^{n+1,m+1} \\ + \frac{1}{2} \beta_{NMm}^{\downarrow} Y_{N,M+1}^{n+1,m} + \frac{1}{2} \alpha_{NMm}^{\uparrow\uparrow} Y_{N-1,M+1}^{n+2,m+1} + \frac{1}{2} \alpha_{NMm}^{\downarrow} Y_{N-1,M+1}^{n+2,m}, \quad (20)$$

$$k_- Y_{NM}^{nm} = \left( \frac{n(n-1)}{2(N+1)} \right) Y_{N+1,M+1}^{n-2,m+1} + \left( -\frac{n(n-1)}{2(N+1)} \right) Y_{N+1,M+1}^{n-2,m} + \frac{1}{2} n F_N^n \beta_{NMm}^{\uparrow\uparrow} Y_{N,M+1}^{n-1,m+1} \\ + \frac{1}{2} n F_N^n \beta_{NMm}^{\downarrow} Y_{N,M+1}^{n-1,m} - \frac{1}{2} F_N^n (F_N^n + 1) \alpha_{NMm}^{\uparrow\uparrow} Y_{N-1,M+1}^{n,m+1} - \frac{1}{2} F_N^n (F_N^n + 1) \alpha_{NMm}^{\downarrow} Y_{N-1,M+1}^{n,m}, \quad (21)$$

where the coefficients  $G_N^n, F_N^n, \beta_{NMm}, \beta_{NMm}^{\downarrow}, \beta_{NMm}^{\uparrow}, \beta_{NMm}^{\uparrow\uparrow}, \alpha_{NMm}, \alpha_{NMm}^{\downarrow}, \alpha_{NMm}^{\uparrow}, \alpha_{NMm}^{\uparrow\uparrow}$  are given below as functions of  $\Gamma$ :

$$\alpha_{NMm} = \frac{N(2\Gamma_1 + N + m - 2)((\Gamma_1 + N - 1)^2 - (\Gamma_2 + M - m)^2)}{(\Gamma_1 + N - 1)^2(2\Gamma_1 + 2N - 3)(2\Gamma_1 + 2N - 1)}, \\ \alpha_{NMm}^{\downarrow} = \frac{N((\Gamma_1 + N - 2) - (\Gamma_2 + M - m))((\Gamma_1 + N - 1) - (\Gamma_2 + M - m))}{(\Gamma_1 + N - 1)^2(2\Gamma_1 + 2N - 3)(2\Gamma_1 + 2N - 1)}, \\ \alpha_{NMm}^{\uparrow} = \frac{mN((\Gamma_1 + n - 2) - (\Gamma_2 + M - m))((\Gamma_1 + n - 1) - (\Gamma_2 + M - m))}{(\Gamma_1 + N - 1)^2(2\Gamma_1 + 2N - 3)(2\Gamma_1 + 2N - 1)}, \\ \alpha_{NMm}^{\uparrow\uparrow} = \frac{N((\Gamma_1 + N - 1)^2 - (\Gamma_2 + M - m)^2)}{(\Gamma_1 + N - 1)^2(2\Gamma_1 + 2N - 3)(2\Gamma_1 + 2N - 1)}, \\ \beta_{NMm} = \frac{(\Gamma_2 + M - m)(\Gamma_1 + m - 1)}{(\Gamma_1 + N)(\Gamma_1 + N - 1)}, \\ \beta_{NMm}^{\downarrow} = \frac{(\Gamma_1 + N - 1) - (\Gamma_2 + M - m)}{(\Gamma_1 + N)(\Gamma_1 + N - 1)}, \\ \beta_{NMm}^{\uparrow} = \frac{m((\Gamma_1 + N - 1) - (\Gamma_2 + M - m))}{(\Gamma_1 + N)(\Gamma_1 + N - 1)}, \\ \beta_{NMm}^{\uparrow\uparrow} = \frac{\Gamma_2 + M - m}{(\Gamma_1 + N)(\Gamma_1 + N - 1)},$$

and  $G_N^n = \Gamma_1 + N + n, F_N^n = -2\Gamma_1 - 2N - n + 1$ . The notation that we used in the above formulas is not accidental. Thus the arrows correspond to the action of the operators that increases or decreases values of  $m$  and/or  $M$ . Those are the indices that are implemented by the Poincaré algebra. Hence their interpretation is given below:

$$\downarrow - m \rightarrow m - 1, \\ \uparrow - M \rightarrow M + 1, \\ \uparrow\uparrow - m \rightarrow m + 1, \quad M \rightarrow M + 1.$$

One has to add one more condition on Eqs. (12)–(21). Namely, in the right-hand sides of all equations we define

$$Y_{OM}^{n1} = Y_{OM}^{n0}. \quad (22)$$

Let us recall that the basis  $Y_{NM}^{nm}$  has been defined for  $n, m, N, M \in \mathbb{N}$  and  $m \leq \min(N, M)$ . The artificial condition just given can be treated as a “boundary condition” that enables us to write all the equations in a closed form. The way in

which the matrix elements in Eqs. (12)–(21) are factored follows from the embedding chain  $so(3) \subset so(3,1) \subset iso(3,1)$ .

Let us discuss now the most important features of the modules  $\mathcal{U}_+(\Gamma)$  exhibited by Eqs. (12)–(21).

(A) Define  $\mathcal{U}_M^m = \{Y_{NM}^{nm}, n, N \geq 0, N \geq m\}$ , where  $M \geq m$ . Then each  $\mathcal{U}_M^m$  is a Verma  $\mathcal{L}$  module that is either an inf-in-mod or an inf-ir-mod depending on the value of  $\Gamma$ . Let us treat  $\mathcal{U}_0^0$  as a prototypical  $\mathcal{L}$  module. Here  $\mathcal{U}_0^0$  is the module that would be obtained if one considered the enveloping algebra of  $\mathcal{L}$ . It is essentially the same module that was considered<sup>5</sup> in the case of the Lorentz algebra. Then,  $\mathcal{U}_M^m$  becomes  $\mathcal{U}_0^0$  defined in terms of starred objects:

$$\Gamma_1^* = \Gamma_1 + m, \quad \Gamma_2^* = \Gamma_2 + M - m, \quad N^* = N - m, \\ Y_{NM}^{*nm} = \begin{cases} Y_{NM}^{nm}, & \text{if } N = 0, \\ \frac{1}{(N-m+1) \cdots (N-1)N} Y_{NM}^{nm}, & \text{if } N \neq 0. \end{cases}$$

(B) For each  $\Gamma \in \mathbb{C}^2$ ,  $\mathcal{U}_+(\Gamma)$  is indecomposable since the value of  $M$  never decreases. It is quite clear from Fig. 1 how submodules, quotient modules, and subquotient modules can be defined.

(C) The sum  $\mathcal{U}_0^0 + \mathcal{U}_1^1 + \mathcal{U}_2^2 \cdots$  is a quotient module and is essentially the same module as  $\Omega_+(\Gamma)$  recently<sup>4</sup> examined. Thus the study conducted previously<sup>4,5</sup> is a special (and very illuminating) case of the results contained in this paper.

(D) Equations (12)–(21) can be extended to all  $n \in \mathbb{Z}$ . The extended vector spaces and modules will be denoted with a tilde, i.e.,  $\tilde{\mathcal{U}}_+(\Gamma)$ ,  $\tilde{\mathcal{U}}_M^m$ , etc.

(E) If  $2\Gamma \in \mathbb{Z}^2$  one may obtain fin-in-mods (or fin-ir-mods when trivial on  $\mathcal{K}$ ) of  $\mathcal{P}$  by passing to the quotients. The essential features of the procedure are the same as previously<sup>4</sup> studied. The fin-in-mods obtained on quotients of the composition series discussed in (C) coincide with fin-in-mods therein<sup>4</sup> contained. New fin-in-mods can be obtained as other quotient modules (see Sec. VI).

(F) The module  $\mathcal{U}_+(\Gamma)$  goes over into the module  $\mathcal{U}_-(-\Gamma)$  under the Lie algebra automorphism  $a$ :

$$\begin{aligned} a(h_3) &= -h_3, & a(p_3) &= -p_3, & a(h_+) &= h_-, \\ a(h_-) &= h_+, & a(p_+) &= p_-, \\ a(p_-) &= p_+, & a(k_+) &= k_-, & a(k_-) &= k_+, \\ a(r_-) &= -r_+, & a(r_+) &= -r_-. \end{aligned}$$

(G) It can be seen that the indecomposability of  $\mathcal{U}_+(\Gamma)$  essentially follows from the fact that  $\mathcal{P}$  is a semidirect product of  $\mathcal{K}$  and  $\mathcal{L}$ . Hence the situation is quite different than in the case of  $\mathcal{L}$  itself where passing to the quotients gives<sup>5</sup> fin-ir-mods. This observation might be of a more general nature and apply to analogous modules of other non-

semisimple Lie algebras that are semidirect products of an Abelian ideal and a Levi subalgebra.

## VI. FINITE-DIMENSIONAL INDECOMPOSABLE MODULES

In this section we are going to discuss fin-in-mods of  $\mathcal{P}$ . They are obtained as quotient modules of  $\mathcal{U}_+(\Gamma)$  [ $\tilde{\mathcal{U}}_+(\Gamma)$ ]. One obtains infinitely many such modules. We are going to content ourselves with giving only a few examples. A more systematic treatment goes beyond the scope of this paper.

In general, fin-in-mods are obtained when  $2\Gamma \in \mathbb{Z}^2$ . Then some of the coefficients in Eqs. (12)–(21) vanish giving rise to submodules and hence to quotient modules (finite dimensional). Thus one has to examine the zeros of  $\alpha$ ,  $\alpha^+$ ,  $\alpha^+$ ,  $\beta$ ,  $\beta^+$ ,  $\beta^+$ ,  $\beta^+$ , and  $F_N^n$ ,  $G_N^n$  as well as the singularities of the former ones for specific values of  $\Gamma$ . Clearly, singularities cannot occur on the quotient modules if one wants to obtain meaningful results. The nature of separation of submodules from their complements (that is, what is worth mentioning, hidden in the PBW basis) is for each  $\mathcal{U}_M^m$  ( $\tilde{\mathcal{U}}_M^m$ ) the same and follows the pattern closely examined<sup>5</sup> in the case of the Lorentz algebra (or, simply,  $\mathcal{U}_0^0$ ). The procedure is somewhat geometrical although it is conceivable to look at Eqs. (12)–(21) from the algebraic point of view. Thus if we put the basis of  $\mathcal{U}_M^m$  on a two-dimensional plane with  $N$  on the horizontal axis and  $n$  on the vertical axis we obtain a subset of a  $\mathbb{Z}^2$  lattice. Two kinds of separation lines can be introduced: broken defined by the equation  $F_N^n = -2\Gamma_1 - 2N - n + 1$  and horizontal given by  $N^+ = \Gamma_2 + M - m - \Gamma_1 + 1$ ,  $N^- = -\Gamma_2 - M + m - \Gamma_1 + 1$ . One can see that the action of all operators of  $\mathcal{P}$  is within each  $\mathcal{U}_M^m$  the same and if it amounts to an interaction with an adjacent  $\mathcal{U}_{M+1}^{m+1}$ ,  $\mathcal{U}_{M+1}^m$ , or  $\mathcal{U}_M^{m-1}$  then it is only modified by a translation [in other words, under projection onto the canonical image of  $\mathcal{U}_+(\mathcal{L})$  in each  $\mathcal{U}_M^m$  it is identical for all  $\mathcal{U}_M^m$ ].

Let us now present several examples of fin-in-mods obtained as quotient modules of  $\mathcal{U}_+(\Gamma)$  [ $\tilde{\mathcal{U}}_+(\Gamma)$ ]. We are going to characterize them by giving the value of  $\Gamma$ , number of interacting  $\mathcal{L}$ -modules (those are  $\mathcal{L}$ -fin-ir-mods), their  $\mathcal{R} \subset \mathcal{L}$  decomposition as well as their bases in the explicit form.

(A) A four-dimensional  $\mathcal{P}$ -fin-in-mod is induced on a quotient of  $\mathcal{U}_0^0 \cup \mathcal{U}_1^0$  when  $\Gamma = (-\frac{1}{2}, -\frac{1}{2})$ . The number of interacting  $\mathcal{L}$ -fin-in-mods is 2 and the  $\mathcal{R} \subset \mathcal{L}$  decomposition is  $2 + 2$  (two fin-ir-mods, both two-dimensional). The basis is given by  $\{Y_{00}^{00}, Y_{00}^{10}\} \cup \{Y_{01}^{00}, Y_{01}^{10}\}$  (we separate the bases of interacting modules). The matrix representation is of the block form:

$$\begin{pmatrix} \sigma_1 & 0 \\ \xi_{12} & \sigma_2 \end{pmatrix},$$

where

$$\begin{aligned} \sigma_1(h_3) &= \begin{pmatrix} -\frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{pmatrix}, & \sigma_1(h_+) &= \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}, \\ \sigma_1(h_-) &= \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}, & \sigma_1(f_+) &= \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}, \end{aligned}$$

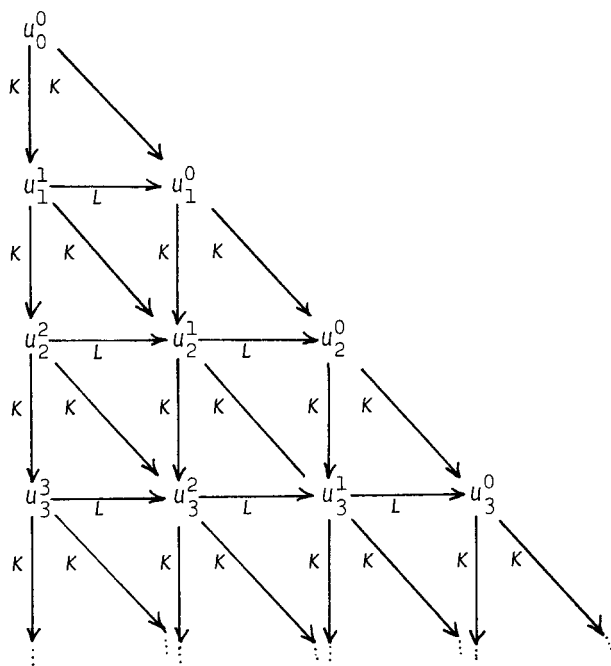


FIG. 1. The module  $\mathcal{U}_+(\Gamma)$ . The action between subquotient modules  $\mathcal{U}_M^m$  is carried by  $\mathcal{K}$  or  $\mathcal{L}$  as shown above.

$$\begin{aligned} \sigma_1(f_-) &= \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}, & \sigma_1(f_3) &= \begin{pmatrix} -\frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{pmatrix}, \\ \sigma_2(h_3) &= \begin{pmatrix} -\frac{1}{2} & 0 \\ 0 & \frac{1}{2} \end{pmatrix}, & \sigma_2(h_+) &= \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix}, \\ \sigma_2(h_-) &= \begin{pmatrix} 0 & 1 \\ 0 & 0 \end{pmatrix}, & \sigma_2(f_+) &= \begin{pmatrix} 0 & 0 \\ -1 & 0 \end{pmatrix}, \\ \sigma_2(f_-) &= \begin{pmatrix} 0 & -1 \\ 0 & 0 \end{pmatrix}, & \sigma_2(f_3) &= \begin{pmatrix} \frac{1}{2} & 0 \\ 0 & -\frac{1}{2} \end{pmatrix}, \\ \xi_{12}(k_+) &= \begin{pmatrix} 0 & 0 \\ -1 & 0 \end{pmatrix}, & \xi_{12}(k_-) &= \begin{pmatrix} 0 & -1 \\ 0 & 0 \end{pmatrix}, \\ \xi_{12}(r_+) &= \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}, & \xi_{12}(r_-) &= \begin{pmatrix} 0 & 0 \\ 0 & -1 \end{pmatrix}. \end{aligned}$$

(B) A five-dimensional  $\mathcal{P}$ -fin-in-mod is induced on a quotient of  $\mathcal{U}_0^0 \cup \mathcal{U}_1^1$  when  $\Gamma = (-1, 0)$ . The number of interacting  $\mathcal{L}$ -fin-ir-mods is two and the  $\mathcal{R} \subset \mathcal{L}$  decomposition is  $4 + 1 = (3 + 1) + 1$ . The basis is given by  $\{Y_{00}^{00}, Y_{00}^{10}, Y_{00}^{20}, Y_{00}^{30}\} \cup \{Y_{11}^{01}\}$ . The matrix representation can be obtained in a straightforward manner.

(C) A seven-dimensional  $\mathcal{P}$ -fin-in-mod is induced on a quotient of  $\tilde{\mathcal{U}}_0^0 \cup \tilde{\mathcal{U}}_1^1$  when  $\Gamma = (1, 2)$ . The number of interacting  $\mathcal{L}$ -fin-ir-mods is two and the  $\mathcal{R} \subset \mathcal{L}$  decomposition is  $4 + 3 = (1 + 3) + 3$ . The basis is given by  $\{\tilde{Y}_{00}^{-10}, \tilde{Y}_{10}^{-10}, \tilde{Y}_{10}^{-20}, \tilde{Y}_{10}^{-30}\} \cup \{\tilde{Y}_{11}^{-11}, \tilde{Y}_{11}^{-21}, \tilde{Y}_{11}^{-31}\}$ .

(D) An eight-dimensional  $\mathcal{P}$ -fin-in-mod is induced on a quotient of  $\mathcal{U}_0^0 \cup \mathcal{U}_1^1$  when  $\Gamma = (-\frac{3}{2}, -\frac{1}{2})$ . The number of interacting  $\mathcal{L}$ -fin-ir-mods is 2 and the  $\mathcal{R} \subset \mathcal{L}$  decomposition is  $6 + 2 = (4 + 2) + 2$ . The basis is given by  $\{Y_{00}^{00}, Y_{00}^{10}, Y_{00}^{20}, Y_{00}^{30}, Y_{10}^{00}, Y_{10}^{10}\} \cup \{Y_{11}^{01}, Y_{11}^{11}\}$ .

(E) An eight dimensional  $\mathcal{P}$ -fin-in-mod is induced on a quotient of  $\mathcal{U}_0^0 \cup \mathcal{U}_1^0 \cup \mathcal{U}_1^1$  when  $\Gamma = (-1, 0)$ . The number of interacting  $\mathcal{L}$ -fin-ir-mods is 3 and the  $\mathcal{R} \subset \mathcal{L}$  decomposition is  $4 + 3 + 1 = (3 + 1) + 3 + 1$ . The basis is given by  $\{Y_{00}^{00}, Y_{00}^{10}, Y_{00}^{20}, Y_{00}^{30}\} \cup \{Y_{01}^{00}, Y_{01}^{10}, Y_{01}^{20}\} \cup \{Y_{11}^{00}\}$ .

The matrix representation of the above  $\mathcal{P}$ -fin-in-mods is of the following block form:

$$\begin{pmatrix} \sigma_1 & 0 \\ \xi_{12} & \sigma_2 \end{pmatrix}$$

in the case of two-interacting  $\mathcal{L}$ -fin-ir-mods  $(A, B, C, D)$ , or

$$\begin{pmatrix} \sigma_1 & 0 & 0 \\ \xi_{12} & \sigma_2 & 0 \\ \xi_{13} & 0 & \sigma_3 \end{pmatrix}$$

in the case of three interacting  $\mathcal{L}$ -fin-ir-mods  $(E)$ . Other  $\mathcal{P}$ -fin-in-mods can be obtained in the manner outlined above. It is interesting to note that from among all the  $\mathcal{P}$ -fin-in-mods up to dimension 8, that were classified by Paneitz,<sup>7</sup> all but one have been identified.

## VII. CONCLUDING REMARKS

We derived the formulas for the induced modules  $\mathcal{U}_+(\Gamma)$  of the Poincaré algebra  $\mathcal{P}$  in an energy-cyclic AMB. It turns out that  $\mathcal{U}_+(\Gamma)$  is indecomposable for each  $\Gamma \in \mathbb{C}^2$ . The subquotient modules of  $\mathcal{U}_+(\Gamma)$  are Verma  $\mathcal{L}$  modules. A large family of new  $\mathcal{P}$ -fin-in-mods can be obtained on quotients of  $\mathcal{U}_+(\Gamma)$  when  $2\Gamma \in \mathbb{Z}^2$ . Several examples were given.

It would be of interest to determine whether other energy-cyclic AMB's give more suitable bases for physical applications. The corresponding formulas could be easily obtained through a linear transformation. Also, a more systematic approach to  $\mathcal{P}$ -fin-in-mods can be attempted.

The extensions of  $\mathcal{U}_+(\Gamma)$  to the scale-invariant algebra ( $\mathcal{P}$  extended by dilations) can be obtained without much effort. However, it would be interesting to study the extension formulas for the conformal algebra.

All of the above matters are currently under investigation.

<sup>1</sup>A. O. Barut and R. Raczka, *Theory of Group Representations and Applications* (Polish Scientific, Warsaw, 1977).

<sup>2</sup>L. Hlavaty and J. Niederle, *Czech. J. Phys. B* t.29 (3), 283 (1979).

<sup>3</sup>P. A. M. Dirac, *Int. J. Theor. Phys.* 23 (8), 677 (1984).

<sup>4</sup>R. Lenczewski and B. Gruber, *J. Phys. A* 19, 1 (1986).

<sup>5</sup>B. Gruber and R. Lenczewski, *J. Phys. A* 16, 3703 (1983).

<sup>6</sup>D. N. Verma, *Bull. Am. Math. Soc.* 74, 160 (1968).

<sup>7</sup>S. M. Paneitz, *Ann. Inst. H. Poincaré* 40 (1), 35 (1984).

<sup>8</sup>I. M. Gel'fand, R. A. Minlos, and Z. Y. Shapiro, *Representations of the Rotation and Lorentz Groups and Their Applications* (Pergamon, Oxford, 1963).

<sup>9</sup>M. A. Naimark, *Linear Representations of the Lorentz Group* (Pergamon, Oxford, 1964).

<sup>10</sup>J. E. Humphreys, *Introduction to Lie Algebras and Representation Theory* (Springer, New York, 1972).

<sup>11</sup>J. Dixmier, *Enveloping Algebras* (North-Holland, Amsterdam, 1977).



# Classes of exactly solvable nonlinear evolution equations for Grassmann variables: The normal form method

Christian Elphick

Laboratoire de Physique Théorique,<sup>a)</sup> Université de Nice, Parc Valrose, 06034 Nice Cedex, France

(Received 4 March 1986; accepted for publication 4 February 1987)

A systematic procedure is presented to solve analytically differential equations for Grassmann variables with the most general nonlinearity. The method consists in the reduction of the original equation to its simplest form (normal form). The classes of solvable normal forms are determined only by the structure of the linear part of the original equation and are parametrized in terms of the number of critical eigenvalues.

## I. INTRODUCTION

Differential equations involving anticommuting (Grassmann) variables (GDE) are often encountered in theoretical physics<sup>1-4</sup> due to the fact that the Grassmann variables  $\theta_i$ ,  $[\theta_i, \theta_j]_+ = 0$  are the classical counterparts ( $\hbar \rightarrow 0$ ) of quantum fermionic operators  $a_i$ ,  $a_i^\dagger$ ,  $[a_i, a_j]_+ = 0$ ,  $[a_i, a_j^\dagger]_+ = \hbar \delta_{ij}$ . For instance, they have been considered by Casalbuoni<sup>4</sup> in the study of the classical mechanics for a Bose-Fermi system; by Berezin and Marinov<sup>1</sup> in the context of the supersymmetric treatment of a classical relativistic free particle which, after quantization, becomes a free Dirac particle; and by Olshanetsky,<sup>5</sup> who considered the supersymmetric version of integrable models by the inverse scattering method. Let us note that GDE also naturally arise when we study a quantum Fermi system via the path integral in the Bargmann-Fock coherent state representation. In this formulation the integral kernel of the evolution operator is given by an integral over a Grassmann algebra whose explicit calculation leads to the evaluation of the action on the classical trajectory satisfying the Grassmann Hamiltonian equations (see Appendix C).

As another example of GDE we mention the widely studied massive Thirring model<sup>6</sup> defined by

$$i \frac{\partial}{\partial x} \phi_1 = m \phi_2 + g \bar{\phi}_2 \phi_2 \phi_1, \quad (1.1)$$

$$i \frac{\partial}{\partial t} \phi_2 = m \phi_1 + g \bar{\phi}_1 \phi_1 \phi_2, \quad (1.2)$$

where  $\phi_1$ ,  $\phi_2$  are Grassmann fields. This model possesses an infinite number of conserved quantities<sup>7</sup> and presents soliton-type behavior. Moreover, Morris<sup>8</sup> has been able to determine a Bäcklund transformation for the anticommuting massive Thirring model by generalizing the prolongation structure method of Wahlquist and Estabrook<sup>9</sup> to Grassmann algebra valued differential forms.

It is the purpose of this paper to study the most general GDE in the sense that they are not necessarily derived from the variation of an action (as in Refs. 1 and 4) and they contain the most general nonlinearity compatible with the Grassmann algebra.

The plan of this paper is as follows: In Sec. II the differential equations to be studied are presented, together with some notation to be used in the remaining sections. For the

sake of clarity I will consider here the case of a discrete number of degrees of freedom (eventually a countable infinity of degrees of freedom), with the generalization to the continuous case (Grassmann fields) being straightforward. Section III is devoted to developing the normal form method for the equations of Sec. II, that is, a method that reduces these equations to their simplest form. Finally, in Sec. IV we comment on other classes of normal forms that can be easily integrated and on the relation of Thirring-type equations with our normal form method.

## II. NONLINEAR GRASSMANN DIFFERENTIAL EQUATIONS

To begin with let  $\mathcal{G}$  be an infinite-dimensional complex Grassmann algebra. Let us consider in  $\mathcal{G}$  a family  $\mathcal{F}$  of  $n$  real odd elements depending on a real parameter  $t$ :

$$\mathcal{F} = \{\theta_i(t); \theta_i(t) = \theta_i^*(t),$$

$$[\theta_i(t), \theta_j(t)]_+ = 0,$$

$$\forall t \in [t_0, +\infty), i, j = 1, \dots, n\},$$

where the asterisk denotes the involution operation<sup>10</sup> in  $\mathcal{G}$  and  $[\ , ]_+$  is the anticommutator.

We assume that the elements of  $\mathcal{F}$  vary with  $t$  according to

$$\frac{\partial}{\partial t} \theta = L \theta + \mathbf{N}(\theta), \quad (2.1)$$

$$\theta(t = t_0) = \theta_0, \quad (2.2)$$

where  $\theta$  is an  $n$ -dimensional vector, with  $\theta = \theta_i e^i$  (sum over repeated indices),  $e^i$ ,  $i = 1, \dots, n$  the canonical basis in a vectorial space that we call  $\mathcal{H}_1$ . Here,  $L$  is a real linear operator acting on  $\mathcal{H}_1$  and  $\mathbf{N}(\theta)$  stands for an arbitrary nonlinearity which preserves the odd character of (2.1), and therefore can be written as

$$\sum_{j=1}^{(1/2)(m-1)} \varphi_i^{(2j+1)} \epsilon_{i_1 i_2 \dots i_{2j+1}} \theta_{i_1} \theta_{i_2} \dots \theta_{i_{2j+1}}, \quad (2.3)$$

where  $\varphi_i^{(2j+1)}$ ,  $i = 1, \dots, n$ ;  $j = 1, \dots, \frac{1}{2}(m-1)$  belong to  $\mathcal{H}_1$ ,  $m = n$  (resp.  $n-1$ ) if  $n$  is odd (resp. even), and  $\epsilon_{i_1 i_2 \dots i_{2j+1}}$  is totally antisymmetric in the indices  $i_k$ ,  $k = 1, \dots, 2j+1$ . Note here that if  $n = 2$ , there is no nonlinear term since the only term available,  $\theta_1 \theta_2$ , is an even element of  $\mathcal{G}$  which is not in  $\mathcal{F}$ .

Let us remark here that the Thirring model (1.1) and

<sup>a)</sup> Unité Associée au C.N.R.S.

(1.2) can be cast into the form (2.1) by writing  $\phi_1 = \varphi_1 + i\varphi_2$ ,  $\phi_2 = \psi_1 + i\psi_2$ , where  $\varphi_1, \varphi_2, \psi_1, \psi_2$  are new independent real Grassmann fields, and by assuming that these fields only depend on a single variable  $\xi = \lambda x - \lambda^{-1}t$ ,  $\lambda \in \mathbb{R}$  (Grassmann solitons).

### III. THE NORMAL FORM METHOD

Let us assume that the spectrum  $S$  of  $L$  is the union of two disjoint parts  $S = S_- \cup S_0$ , where  $\gamma \in S_-$  (resp.  $S_0$ ) if  $\text{Re } \gamma < 0$  (resp.  $\text{Re } \gamma = 0$ ). We suppose that the linear operator  $L$  has  $n_c$  eigenvalues belonging to  $S_0$  which we call critical [since they will determine the solvable classes of (2.1) by this method], while the remaining eigenvalues  $\gamma_\alpha$ ,  $\alpha = 1, \dots, M = n - n_c$  ( $M$  eventually infinite) are all different and have negative real parts. Since  $L$  is a real linear operator the critical eigenvalues can be zero or purely imaginary pairs  $(i\omega_k, -i\omega_k)$ ,  $k = 1, \dots, s$ . Then  $n_c = 2s + l$ , where  $l$  is the algebraic multiplicity of the zero eigenvalue.

In the following we will consider the class of nonlinear equations parametrized by  $s = 1, l = 1$  and it will be proved that it can be exactly solved. (Section IV and Appendices B and C deal with other classes that can also be treated by the method developed in this section.)

The method is as follows. According to our assumptions  $\mathcal{H}_1$  is the direct sum of two subspaces invariant by  $L$ ,  $\mathcal{H}_1 = \mathcal{H}_1^c \oplus \mathcal{H}_1^s$ , where  $\mathcal{H}_1^s$  is spanned by the vectors  $\psi_\alpha$ ,  $L\psi_\alpha = \gamma_\alpha \psi_\alpha$  and  $\mathcal{H}_1^c$  is the critical space spanned by  $\{\phi_1, \phi_2 = \phi, \phi_3 = \bar{\phi}\}$  such that

$$L\phi_i = J_{ji}\phi_j, \quad J = \begin{pmatrix} 0 & 0 & 0 \\ 0 & i\omega & 0 \\ 0 & 0 & -i\omega \end{pmatrix}. \quad (3.1)$$

Let us introduce into  $\mathcal{G}$  a new family of odd variables  $\{A_1, A_2 = A, A_3 = A^*, B_\alpha, \alpha = 1, \dots, M\}$  defined through the nonlinear change of variables

$$\theta(t) = \sum_{i=1}^3 A_i(t)\phi_i + \sum_{\alpha=1}^M B_\alpha(t)\psi_\alpha + \sum_{j=1}^{(m-1)/2} \theta^{[2j+1]}(C_p(t)), \quad (3.2)$$

where  $m = n$  for odd  $n$ ,  $m = n - 1$  for even  $n$ , and  $\theta^{[2j+1]}$  is homogeneous of degree  $2j + 1$  in the variables  $\{C_p\}$ , where  $C_p = A_p, p = 1, 2, 3, C_{3+\alpha} = B_\alpha, \alpha \geq 1$ , and therefore is of the form

$$\theta^{[2j+1]} = \sum (C_{i_1})^{r_1} (C_{i_2})^{r_2} \dots (C_{i_p})^{r_p} \theta_{i_1 i_2 \dots i_p}^{[2j+1] r_1 r_2 \dots r_p}, \quad (3.3)$$

with

$$\sum_{i=1}^p r_i = 2j + 1, \quad r_i \in \{0, 1\}, \quad i_j \neq i_k, \quad \text{if } j \neq k, \quad (3.4)$$

and  $\theta_{i_1 \dots i_p}^{[2j+1] r_1 \dots r_p} \in \mathcal{H}_1$  totally antisymmetric in  $i_1 \dots i_p$  and time independent.

We look for equations for  $\{A_i, B_\alpha\}$  of the form

$$\begin{aligned} \frac{\partial A_i}{\partial t} &= J_{ij} A_j + f_i^{[3]} + f_i^{[5]} + \dots + f_i^{[m]}, \\ \frac{\partial B_\alpha}{\partial t} &= \gamma_\alpha B_\alpha + g_\alpha^{[3]} + g_\alpha^{[5]} + \dots + g_\alpha^{[m]}, \end{aligned} \quad (3.5)$$

(no sum over  $\alpha$ ),

where  $f_i^{[r]}, g_\alpha^{[r]}$  are homogeneous of degree  $r$  in  $\{C_p\}$ .

From (3.2), we have

$$\frac{\partial \theta}{\partial t} = \frac{\partial C_p}{\partial t} \frac{\partial^l}{\partial C_p} \theta, \quad (3.6)$$

where the superscript  $l$  means that in differentiating with respect to  $C_p$  we must displace the variable  $C_p$  to the left before dropping it.

By replacing (3.2) in (2.1) and using (3.5) and (3.6), we obtain, after an identification of each order in  $\{C_p\}$ , a hierarchy of equations for  $\theta^{[r]}$ ,  $r = 3, 5, \dots, m$ . For  $r = 1$ , we obtain

$$J_{ij} A_j \phi_i + \gamma_\beta B_\beta \psi_\beta = L(A_i \phi_i + B_\beta \psi_\beta), \quad (3.7)$$

which is an identity due to (3.1) and  $L\psi_\beta = \gamma_\beta \psi_\beta$ .

At order  $r \geq 3$  (odd  $r$ ), we obtain the following homological equation<sup>11</sup>:

$$\begin{aligned} \mathcal{L}\theta^{[r]} &= (\mathcal{A} + \mathcal{B} - L)\theta^{[r]} \\ &= \mathbf{I}^{[r]} - f_i^{[r]} \phi_i - g_\alpha^{[r]} \psi_\alpha \equiv \mathbf{K}^{[r]}, \end{aligned} \quad r = 3, 5, \dots, m, \quad (3.8)$$

where

$$\mathcal{A} = J_{ij} A_j \frac{\partial^l}{\partial A_i}, \quad \mathcal{B} = \gamma_\alpha B_\alpha \frac{\partial^l}{\partial B_\alpha} \quad (3.9)$$

and

$$\mathbf{I}^{[r]} = \mathbf{N}^{[r]}(\theta) - \sum_{s=3}^{r-2} \left( \frac{\partial}{\partial t} C_p \right)^{[s]} \frac{\partial^l}{\partial C_p} \theta^{[r-s+1]}. \quad (3.10)$$

If  $r = 3$  the last term in (3.10) is absent.

Let us look now at the structure of the space in which (3.8) is written: It is the tensor product  $\mathcal{H} = \mathcal{H}_1 \otimes \mathcal{H}_2 \otimes \mathcal{H}_3$ , where  $\mathcal{H}_2$  (resp.  $\mathcal{H}_3$ ) is generated by monomials constructed from the variables  $A_i$  (resp.  $B_\alpha$ ). Therefore  $\mathcal{H}$  is an orthogonal sum of  $\mathcal{L}$ -invariant subspaces.

$$\mathcal{H} = \bigoplus_{r=3,5,\dots,m} \bigoplus_{N=0}^{N=r} \mathcal{H}_3^{[r-N]} \otimes \mathcal{H}_2^{[N]} \otimes (\mathcal{H}_1^c \oplus \mathcal{H}_1^s), \quad (3.11)$$

where  $\mathcal{H}_2^{[N]}$  is generated by  $A_1^{r_1} A_2^{r_2} A_3^{r_3}$ ,  $\sum_{i=1}^3 r_i = N \leq 3$ ,  $r_i \in \{0, 1\}$ , if  $N > 3$   $\mathcal{H}_2^{[N]}$  is empty, and  $\mathcal{H}_3^{[r-N]}$  is generated by  $B_1^{r_1} \dots B_M^{r_M}$ ,  $\sum_{i=1}^M r_i = r - N$ ,  $r_i \in \{0, 1\}$ .

We now introduce a scalar product in  $\mathcal{H}_2, \mathcal{H}_3$  defined as

$$\begin{aligned} \langle f_1, f_2 \rangle &= \int \exp\left(-\sum_{i=p'}^{p''} C_i^* C_i\right) f_1(C_j) f_2^*(C_j) \\ &\quad \times dC_{p'}^* dC_{p'} \dots dC_{p''}^* dC_{p''}, \end{aligned} \quad (3.12)$$

where  $p' = 3, p'' = 1$  for  $\mathcal{H}_2$ , and  $p' = M, p'' = 4$  for  $\mathcal{H}_3$  [note that (3.12) is a multiple Grassmann integral defined in the Berezin<sup>10</sup> sense]. We define the scalar product in  $\mathcal{H}$  as

$$\langle \cdot, \cdot \rangle_{\mathcal{H}} = \langle \cdot, \cdot \rangle_{\mathcal{H}_1} \langle \cdot, \cdot \rangle_{\mathcal{H}_2} \langle \cdot, \cdot \rangle_{\mathcal{H}_3}, \quad (3.13)$$

where  $\langle \cdot, \cdot \rangle_{\mathcal{H}_1}$  is such that the eigenbasis is orthonormal.

It is easy to see now that  $\mathcal{L}$  is a linear noninvertible operator acting on  $\mathcal{H}$  and therefore (3.8) will have no solutions for  $\theta^{[r]}$  unless  $\mathbf{K}^{[r]}$  is orthogonal to the null space of the adjoint of  $\mathcal{L}$  in the scalar product (3.13). This solvability condition will determine the unknown functions  $f_i^{[r]}, g_\alpha^{[r]}$  in (3.5) in such a way that Eq. (3.8) can be solved for  $\theta^{[r]}$ . We note that  $f_i^{[r]}, g_\alpha^{[r]}$  (resp.  $\theta^{[r]}$ ) are determined modulo terms in  $\text{Ran } \mathcal{L}$  (resp.  $\text{Ker } \mathcal{L}$ ). We choose to make a minimal prescription such that all these gauge terms are taken to be zero. [Obviously the original solution  $\theta(t, t_0)$  of (2.1) will not depend on this choice.]

At this point it is worthwhile remarking that the operators  $A_i, \partial^i/\partial A_i$  (resp.  $B_\alpha, \partial^i/\partial B_\alpha$ ) are adjoint to each other in the scalar product (3.12) and that they satisfy ( $a_i^\dagger = A_i, a_i = \partial^i/\partial A_i, b_\alpha^\dagger = B_\alpha, b_\alpha = \partial^i/\partial B_\alpha$ )

$$[a_i, a_j^\dagger]_+ = \delta_{ij}, \quad [b_\alpha, b_\beta^\dagger]_+ = \delta_{\alpha\beta}, \quad (3.14)$$

i.e., they are Fermionic creation and annihilation operators in the Fock space  $\mathcal{H}_2$  (resp.  $\mathcal{H}_3$ ) [Eqs. (3.14) show that  $\mathcal{L}$  is the Fermionic version of the homological operator introduced by Arnold<sup>[11]</sup> and, for instance, we can write the generators of  $\mathcal{H}_2^{[N]}, \mathcal{H}_3^{[r-N]}$  as

$$\begin{aligned} |r_1, r_2, r_3\rangle &= (a_1^\dagger)^{r_1} (a_2^\dagger)^{r_2} (a_3^\dagger)^{r_3} |0\rangle, \quad |0\rangle = 1, \\ \sum_{i=1}^3 r_i &= N, \quad r_i \in \{0, 1\}, \\ |r_1, \dots, r_M\rangle &= (b_1^\dagger)^{r_1} \dots (b_M^\dagger)^{r_M} |0\rangle, \\ \sum_{i=1}^M r_i &= r - N, \quad r_i \in \{0, 1\}, \end{aligned} \quad (3.15)$$

which are obviously orthonormal.

Therefore the adjoint of  $\mathcal{L}$  in the scalar product (3.12) reads as

$$\mathcal{L}^\dagger = \bar{J}_{ij} A_i \frac{\partial^i}{\partial A_j} + \bar{\gamma}_\alpha B_\alpha \frac{\partial^i}{\partial B_\alpha} - L^\dagger, \quad (3.16)$$

where  $L^\dagger$  is the conjugate transpose of  $L: L^\dagger = \bar{L}^t$ .

It is not difficult to check that the null space of  $\mathcal{L}^\dagger$  in  $\mathcal{H}$  is generated by (provided the set  $\{\gamma_\alpha\}$  satisfies the nonresonant condition  $\gamma_\alpha \neq \sum_{\alpha' \in \Delta} \sigma_{\alpha'}, \forall \Delta \subset \{-2, -1, 1, \dots, M\}, \sigma_{\alpha'} = \{\sigma_{-2} = i\omega, \sigma_{-1} = -i\omega, \gamma_{\alpha'}\}$ )

$$\{A_1 A_2 A_3 \phi_1, B_\alpha A_2 A_3 \psi_\alpha, \alpha = 1, \dots, M\}. \quad (3.17)$$

It is worth noting that although vectors of the form  $X_\alpha = B_\alpha A_1 A_2 A_3 \psi_\alpha$  do satisfy  $\mathcal{L}^\dagger X_\alpha = 0$ , they are not included in (3.17) since they do not belong to  $\mathcal{H}$  (they are even elements in  $\mathcal{G}$ ).

Therefore we have the following solvability conditions:

$$\begin{aligned} \langle f_1^{[r]}, A_1 A_2 A_3 \rangle_{\mathcal{H}_2 \otimes \mathcal{H}_3} &= \langle \mathbf{I}^{[r]}, A_1 A_2 A_3 \phi_1 \rangle_{\mathcal{H}}, \\ \langle g_\alpha^{[r]}, B_\alpha A_2 A_3 \rangle_{\mathcal{H}_2 \otimes \mathcal{H}_3} &= \langle \mathbf{I}^{[r]}, B_\alpha A_2 A_3 \psi_\alpha \rangle_{\mathcal{H}}, \quad \alpha = 1, \dots, M. \end{aligned} \quad (3.18)$$

Consequently we conclude that we can take  $f_2^{[r]} = f_3^{[r]} = 0, f_1^{[r]} = g_\alpha^{[r]} = 0$ , for  $r > 3$ , and

$$\begin{aligned} f_1^{[3]} &= \langle \mathbf{N}^{[3]}(\theta), A_1 A_2 A_3 \phi_1 \rangle_{\mathcal{H}} A_1 A_2 A_3 = k A_1 A_2 A_3, \\ g_\alpha^{[3]} &= \langle \mathbf{N}^{[3]}(\theta), B_\alpha A_2 A_3 \psi_\alpha \rangle_{\mathcal{H}} B_\alpha A_2 A_3 = k_\alpha B_\alpha A_2 A_3. \end{aligned} \quad (3.19)$$

We remark that  $k, k_\alpha$  in (3.19) are totally determined since  $\mathbf{N}^{[3]}(\theta)$  only depends on  $\theta^{[1]} = A_i \phi_i + B_\alpha \psi_\alpha$ .

Having completely determined the functions  $f_i^{[r]}, g_\alpha^{[r]}$  in such a way that (3.8) can be solved for  $\theta^{[r]}$ , and noting that  $\mathbf{I}^{[r]}$  only depends on  $\theta^{[s]}$  for  $s < r$ , we see that (3.8) can be solved by recursion in  $r$ . [In Appendix A it is shown how to solve Eq. (3.8).]

Thus we have finally derived the normal form equations for  $\{A_i, B_\alpha\}$ . By using (3.18) and (3.19) in (3.5), we obtain

$$\begin{aligned} \frac{\partial A_1}{\partial t} &= k A_1 A A^*, \quad \frac{\partial A}{\partial t} = i\omega A, \quad \frac{\partial A^*}{\partial t} = -i\omega A^*, \\ \frac{\partial B_\alpha}{\partial t} &= B_\alpha (\gamma_\alpha + k_\alpha A A^*), \quad \alpha = 1, \dots, M. \end{aligned} \quad (3.20)$$

[Note that  $k, k_\alpha$  are complex numbers since  $A_i^\dagger = A_i, B_\alpha^\dagger = B_\alpha, (A A^*)^* = -A A^*$  (since  $A, A^*$  anticommute).]

Once Eq. (3.8) has been solved for  $\theta^{[r]}$ , the original vector  $\theta(t)$  is obtained through (3.2); if one evaluates this expression at  $t = t_0$ , one obtains a relation between  $\theta_0$  and  $\{A_i(t_0), B_\alpha(t_0)\}$ , which can be solved for  $\{A_i(t_0), B_\alpha(t_0)\}$ , determining in this way the initial conditions for (3.20) as functions of the original initial conditions. In short,  $\{A_i(t_0), B_\alpha(t_0)\}$  can be completely determined from the knowledge of  $\theta_0$ .

Equations (3.20) are easily solved and we obtain

$$\begin{aligned} A(t) &= e^{i\omega(t-t_0)} A(t_0), \quad A^*(t) = e^{-i\omega(t-t_0)} A^*(t_0), \\ A_1(t) &= (1 + k(t-t_0) A(t_0) A^*(t_0)) A_1(t_0), \\ B_\alpha(t) &= e^{\gamma_\alpha(t-t_0)} (1 + k_\alpha(t-t_0) A(t_0) A^*(t_0)) B_\alpha(t_0). \end{aligned} \quad (3.21)$$

By replacing (3.21) back into (3.2) we obtain the solution of (2.1) and (2.2). Let us remark that one can formally obtain from (3.20) the case when  $S_0$  is empty by putting  $A = A^* = A_1 = 0$ ; then the original problem is reduced to the linear equations,

$$\frac{\partial B_\alpha}{\partial t} = \gamma_\alpha B_\alpha, \quad \alpha = 1, \dots, M. \quad (3.22)$$

We finally note that the solvability conditions have a simple interpretation in terms of particle physics:  $\text{Ker } \mathcal{L}^\dagger$  is generated by vectors  $\mathbf{X} = X_1 \phi_1 + X_2 \phi_2 + X_3 \phi_3$  and  $\mathbf{Y} = B_\alpha Y_\alpha \psi_\alpha$  (no sum over  $\alpha$ ) such that  $\mathbf{X} \in \text{Ker } (\mathcal{A}^\dagger - \mathcal{L}^\dagger)$  and  $\mathbf{Y} \in \text{Ker } \mathcal{A}^\dagger$ , which means that they satisfy ( $\mathcal{A}^\dagger = -\mathcal{A}$ )

$$\begin{aligned} \frac{1}{2} \mathcal{A} X_1 &= 0, \quad \frac{1}{2} \mathcal{A} X_2 = \frac{1}{2} X_2, \quad \frac{1}{2} \mathcal{A} X_3 = -\frac{1}{2} X_3, \\ \frac{1}{2} \mathcal{A} Y_\alpha &= 0, \quad \alpha = 1, \dots, M. \end{aligned} \quad (3.23)$$

By writing  $\frac{1}{2} \mathcal{A}$  as

$$\frac{1}{2} a_i^\dagger (\lambda_3)_{ij} a_j, \quad (3.24)$$

where  $\lambda_3$  is the diagonal Gell-Mann matrix  $\text{diag}(0, 1, -1)$ , we easily see that  $\frac{1}{2} \mathcal{A}$  is one of the diagonal generators of the Cartan subalgebra of color  $\text{SU}(3)$  [ $\frac{1}{2} \mathcal{A} = I_3$  is also one of the generators of the isospin subgroup of  $\text{SU}(3)$ ]. Therefore  $\mathbf{X}$  can be regarded as a quark triplet with components  $(s, u, d)$ . It follows that the solvability conditions project Eq. (3.8) on the  $\frac{1}{2}$ -isospin plane  $(u, d)$ .

#### IV. GENERAL COMMENTS

Let us finally mention that besides the class of equations studied here [ $s = 1, l = 1$ ] and the trivial case when  $S_0$  is empty, there are other classes that can also be treated and solved by the method presented here in a straightforward way. These are the following:

$$[s = 0, l = 1], \quad \mathbb{J} = (0); \quad (4.1)$$

$$[s = 0, l = 2], \quad \mathbb{J}_1 = \begin{Bmatrix} 0 & 0 \\ 0 & 0 \end{Bmatrix} \quad \text{or} \quad \mathbb{J}_2 = \begin{Bmatrix} 0 & 1 \\ 0 & 0 \end{Bmatrix}; \quad (4.2)$$

$$[s = 1, l = 0], \quad \mathbb{J}_3 = \begin{Bmatrix} i\omega & 0 \\ 0 & -i\omega \end{Bmatrix}; \quad (4.3)$$

$$[s = 0, l = 3], \quad \mathbb{J} = \begin{Bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{Bmatrix}$$

$$\text{or} \quad \mathbb{J} = \begin{Bmatrix} \mathbb{J}_1 & 0 \\ 0 & 0 \end{Bmatrix} \quad \text{or} \quad \mathbb{J} = \begin{Bmatrix} \mathbb{J}_2 & 0 \\ 0 & 0 \end{Bmatrix}; \quad (4.4)$$

$$[s = 1, l = 2], \quad \mathbb{J} = \begin{Bmatrix} \mathbb{J}_3 & \mathbb{J}_1 \\ \mathbb{J}_1 & \mathbb{J}_1 \end{Bmatrix} \quad \text{or} \quad \mathbb{J} = \begin{Bmatrix} \mathbb{J}_3 & \mathbb{J}_1 \\ \mathbb{J}_1 & \mathbb{J}_2 \end{Bmatrix}. \quad (4.5)$$

For other classes with more complicated matrices  $\mathbb{J}$ , although the method does not give directly the solution to (2.1), it allows us to write (2.1) in its normal form, which means in its simplest form (with the least number of nonlinear terms) the solution can be found either by inspection or by developing new methods such as Bäcklund transformations or inverse scattering methods applied to GDE.

It is also interesting to remark that if we consider the case when

$$\mathbb{J} = \begin{Bmatrix} \mathbb{J}_3 & \mathbb{J}_1 \\ \mathbb{J}_1 & \mathbb{J}_3 \end{Bmatrix}$$

and we apply the method presented in Sec. III, we obtain a normal form equation for the critical variables, which is exactly the anticommuting massive Thirring model when written in terms of the variable  $\xi = \lambda x - \lambda^{-1}t$ , i.e., the normal form equations read as

$$\begin{aligned} i \frac{\partial}{\partial t} \phi_1 &= \omega \phi_2 + g \phi_2^* \phi_2 \phi_1, \\ i \frac{\partial}{\partial t} \phi_2 &= \omega \phi_1 + g \phi_1^* \phi_1 \phi_2, \\ \frac{\partial B_\alpha}{\partial t} &= B_\alpha (\gamma_\alpha + \delta_\alpha (\phi_2^* \phi_2 + \phi_1^* \phi_1) \\ &\quad + \beta_\alpha \phi_1^* \phi_1 \phi_2^* \phi_2), \quad \alpha = 1, \dots, M. \end{aligned} \quad (4.6)$$

Therefore we conclude that the Thirring model is a universal equation for all the nonlinear Grassmann differential equations parametrized by [ $s = 2, l = 0$ ], in the sense that for times  $t \gg \text{Sup}_\alpha \text{Re}|\gamma_\alpha|^{-1}$ , all the dynamics is governed by this model (see Appendix B).

Let us note that rigorously speaking the different classes of GDE should also be parametrized by a third integer  $c^{11}$  (the codimension in the language of bifurcation theory) which counts the minimum number of parameters needed to unfold the critical situation. For instance, the matrix  $\mathbb{J}$  studied here has  $c = 2$ . With the last remark in mind, it is clear

that in a zero-parameter class of Eq. (2.1) the generic situation to be found is the case when  $S_0$  is empty.

Finally, let us comment on possible generalizations of this work. One possibility is to extend the present method to a supersymmetric classical system described by the action

$$S = \int_{t_i}^{t_f} dt d\theta_1 d\theta_2 (-i\Phi_j D_1 D_2 \Phi_j - U(\Phi)), \quad (4.7)$$

where  $D_1, D_2$  are supersymmetric covariant derivatives:

$$D_1 = -\frac{\partial'}{\partial \theta_2} + i\theta_2 \frac{\partial}{\partial t}, \quad D_2 = \frac{\partial'}{\partial \theta_1} + i\theta_1 \frac{\partial}{\partial t}, \quad (4.8)$$

where  $U(\Phi)$  is the superpotential and  $\Phi_j, j = 1, 2, \dots$ , is a set of supervariables,

$$\begin{aligned} \Phi_j &= \varphi_j(t) + i\bar{\theta}\psi_j(t) + \frac{1}{2}\bar{\theta}\theta\varphi_j'(t), \\ \bar{\theta} &= \theta^t \gamma_0, \quad \gamma_0 = \begin{pmatrix} 0 & 1 \\ -1 & 0 \end{pmatrix}, \end{aligned} \quad (4.9)$$

where  $\varphi_j, \varphi_j'$  are commuting variables and  $\psi_j$  is a two-component Majorana spinor.

The classical mechanics of our system is given by an action obtained from (4.7) after integration over  $\theta_1, \theta_2$ :

$$S' = \int_{t_i}^{t_f} dt (-i\varphi_j \ddot{\varphi}_j + \psi_j' \dot{\psi}_j + U(\varphi_j, \psi_j)), \quad (4.10)$$

where an overdot stands for the time derivative.

Variation of  $S'$  leads to nonlinear differential equations coupling commuting and anticommuting variables. In particular, by using the method of Sec. III, we arrive at an equation like (3.8), but where the operator  $\mathcal{L}$  is replaced by a super homological operator  $\hat{\mathcal{L}}$  which contains both bosonic and fermionic creation and annihilation operators. The structure of  $\text{Ker } \hat{\mathcal{L}}^\dagger$  is much more complicated than in the case considered here and therefore the classification of solvable classes or classes reducible to simple normal forms becomes much more complex (details of the normal form method for a supersymmetric classical mechanics will be given elsewhere).

#### ACKNOWLEDGMENT

It is a pleasure for the author to thank Professor E. Tirapegui for helpful discussions and much encouragement during the course of this work.

#### APPENDIX A: SOLUTION OF THE HOMOLOGICAL EQUATION

Here we show how to solve Eq. (3.8) for the Grassmann vector  $\theta^{[r]}$ ,  $r = 3, 5, \dots, m$ . For simplicity we consider the case  $M = 1$  which will suffice for understanding how the solution mechanism works for  $M > 1$ .

We start by decomposing  $\mathbf{N}(\theta)$  in Eq. (2.1) as

$$\mathbf{N}(\theta) = N_1(\theta)\phi_1 + N_2(\theta)\phi_2 + N_3(\theta)\phi_3 + N_4(\theta)\psi, \quad (A1)$$

where

$$N_i(\theta) = k_{i_1 i_2 i_3} \theta_{i_1} \theta_{i_2} \theta_{i_3}, \quad i = 1, \dots, 4, \quad (A2)$$

and  $k_{i_1 i_2 i_3}$  is totally antisymmetric in  $i_1, i_2, i_3$ .

For  $M = 1$ , Eq. (3.8) reduces to the equation for  $\theta^{[3]}$ :

$$\begin{aligned} & \left( i\omega \left( A_2 \frac{\partial^l}{\partial A_2} - A_3 \frac{\partial^l}{\partial A_3} \right) + \gamma B \frac{\partial^l}{\partial B} - L \right) \theta^{[3]} \\ & = \mathbf{N}^{[3]}(\theta) - kA_1A_2A_3\phi_1 - k_1BA_2A_3\psi, \end{aligned} \quad (\text{A3})$$

where use has been made of (3.10) and (3.19).

We note that  $\mathbf{N}^{[3]}(\theta)$  only depends on  $\theta^{[1]}$ , which is known by construction [Eq. (3.2)] and is given by (A1) when we replace  $\theta_i$  by  $\rho_i A_i$ ,  $\rho_i = \langle \mathbf{e}_i, \phi_i \rangle$ ,  $i = 1, 3$ , and  $\theta_4$  by  $\rho_4 B$ ,  $\rho_4 = \langle \mathbf{e}_4, \psi \rangle$ . We obtain

$$\mathbf{N}^{[3]}(\theta) = N_1^{[3]}\phi_1 + N_2^{[3]}\phi_2 + N_3^{[3]}\phi_3 + N_4^{[3]}\psi, \quad (\text{A4})$$

where

$$N_i^{[3]} = a_i A_1 A_2 A_3 + b_i B A_1 A_2 + c_i B A_1 A_3 + d_i B A_2 A_3 \quad (\text{A5})$$

and the constants  $a_i, \dots, d_i$  are functions of  $\rho_i$ ,  $k_{i_1 i_2 i_3}$  [ $a_1 = k$ ,  $d_4 = k_1$  by the solvability conditions (3.18)].

Now writing  $\theta^{[3]}$  as

$$\theta^{[3]} = \theta_1^{[3]}\phi_1 + \theta_2^{[3]}\phi_2 + \theta_3^{[3]}\phi_3 + \theta_4^{[3]}\psi \quad (\text{A6})$$

and using (A3), we arrive at the following:

$$\mathcal{D}\theta_1^{[3]} = B(b_1 A_1 A_2 + c_1 A_1 A_3 + d_1 A_2 A_3), \quad (\text{A7})$$

$$(\mathcal{D} - i\omega)\theta_2^{[3]} = N_2^{[3]}, \quad (\text{A8})$$

$$(\mathcal{D} - \gamma)\theta_4^{[3]} = a_4 A_1 A_2 A_3 + b_4 B A_1 A_2 + c_4 B A_1 A_3, \quad (\text{A9})$$

where

$$\mathcal{D} = i\omega \left( A_2 \frac{\partial^l}{\partial A_2} - A_3 \frac{\partial^l}{\partial A_3} \right) + \gamma B \frac{\partial^l}{\partial B}.$$

The solutions of Eqs. (A7)–(A9) are obtained by expanding  $\theta_1^{[3]}$ ,  $\theta_2^{[3]}$ ,  $\theta_3^{[3]}$  in terms of monomials of degree 3 in  $(A_1, A_2, A_3, B)$  and comparing the coefficients of each monomial in both sides of (A7)–(A9). We obtain

$$\theta_1^{[3]} = B \left( \frac{b_1}{\gamma + i\omega} A_1 A_2 + \frac{c_1}{\gamma - i\omega} A_1 A_3 + \frac{d_1}{\gamma} A_2 A_3 \right), \quad (\text{A10})$$

$$\begin{aligned} \theta_2^{[3]} &= -\frac{a_2}{i\omega} A_1 A_2 A_3 \\ &+ B \left( \frac{b_2}{\gamma} A_1 A_2 + \frac{c_2}{\gamma - 2i\omega} A_1 A_3 + \frac{d_2}{\gamma - i\omega} A_2 A_3 \right), \end{aligned} \quad (\text{A11})$$

$$\theta_4^{[3]} = -(a_4/\gamma)A_1A_2A_3 + (B/i\omega)A_1(b_4A_2 - c_4A_3). \quad (\text{A12})$$

The above expressions completely determine the nonlinear change of variables (3.2) which, together with (3.21), give the complete solution of Eq. (2.1).

## APPENDIX B: THE THIRRING MODEL AS THE NORMAL FORM [ $s=2, l=0$ ]

Here we prove that the Thirring model is the normal form for the class parametrized by [ $s=2, l=0$ ].

We assume that  $t \gg \text{Sup}_\alpha \text{Re}|\gamma_\alpha|^{-1}$ , which means that we are in the asymptotic regime, where all the variables  $B_\alpha$ ,  $\alpha = 1, \dots, M$ , have relaxed to zero and the dynamics is solely described in terms of the critical variables  $(A_1, A_1^*, A_2, A_2^*)$ .

In this regime, Eq. (3.8) reads as

$$\begin{aligned} \mathcal{L}\theta^{[r]} &= \left( \mathbb{J}_{ij} A_i \frac{\partial^l}{\partial A_j} - \mathbb{J} \right) \theta^{[r]} \\ &= I^{[r]} - f_i^{[r]} \phi_i = \mathbf{K}^{[r]}, \quad r = 3, 5, \dots, m. \end{aligned} \quad (\text{B1})$$

To satisfy the solvability condition  $\mathbf{K}^{[r]} \in \text{Ran } \mathcal{L}$ , it is sufficient to choose  $f_i^{[r]} \phi_i$  in  $\text{Ker } \mathcal{L}^\dagger \subset \mathcal{H}_1 \otimes \mathcal{H}_2^{[r]}$ , or equivalently,  $\mathbf{F} = \sum_{r=3,5,\dots,m} f_i^{[r]} \phi_i$  in  $\text{Ker } \mathcal{L}^\dagger \subset \mathcal{H}_1 \otimes \mathcal{H}_2$ :

$$\left( (\mathbb{J}^\dagger)_{ij} A_j \frac{\partial^l}{\partial A_i} - \mathbb{J}^\dagger \right) \mathbf{F} = 0. \quad (\text{B2})$$

By introducing an auxiliary parameter  $\tau$ , Eq. (B2) can be cast into the form

$$\frac{d}{d\tau} e^{-\mathbb{J}^\dagger \tau} \mathbf{F}(e^{\mathbb{J}^\dagger \tau} \mathbf{A}) = 0, \quad \mathbf{A} = A_i \phi_i, \quad (\text{B3})$$

from which we deduce that

$$\mathbf{F}(e^{\mathbb{J}^\dagger \tau} \mathbf{A}) = e^{\mathbb{J}^\dagger \tau} \mathbf{F}(\mathbf{A}). \quad (\text{B4})$$

For the class [ $s=2, l=0$ ],  $\mathbb{J}^\dagger$  is given by

$$\mathbb{J}^\dagger = \begin{pmatrix} -i\omega & 0 & 0 & 0 \\ 0 & +i\omega & 0 & 0 \\ 0 & 0 & -i\omega & 0 \\ 0 & 0 & 0 & i\omega \end{pmatrix}, \quad (\text{B5})$$

and therefore (B4) tells us that  $\mathbf{F}$  is invariant under two independent rotations by  $\omega\tau$  in the planes  $(A_1, A_1^*)$ ,  $(A_2, A_2^*)$ . Employing the above invariance and assuming that the original system is  $T$  invariant [which means that if  $\theta(t)$  is a solution of (2.1) then  $\theta^*(-t)$  is also a solution of (2.1)], we obtain

$$\begin{aligned} F_1 &= igA_2^*A_2A_1, & F_2 &= F_1^*, \\ F_3 &= igA_1^*A_1A_2, & F_4 &= F_3^*, \end{aligned} \quad (\text{B6})$$

where  $g$  is a real constant. Making now the change of variables

$$\varphi_1 = A_1 + iA_2^*, \quad \varphi_2 = A_1 - iA_2^*, \quad (\text{B7})$$

we arrive at the following normal form for the class [ $s=2, l=0$ ]:

$$\begin{aligned} \frac{1}{i} \frac{\partial \varphi_1}{\partial t} &= \omega \varphi_2 + \frac{g}{2} \varphi_2^* \varphi_2 \varphi_1, \\ \frac{1}{i} \frac{\partial \varphi_2}{\partial t} &= \omega \varphi_1 + \frac{g}{2} \varphi_1^* \varphi_1 \varphi_2, \end{aligned} \quad (\text{B8})$$

which is nothing but the Thirring model for Grassmann solitons<sup>8</sup> if we identify  $\omega$  with the soliton mass and  $t$  with an appropriate variable defined in a comoving system with the soliton.

## APPENDIX C: SOME PHYSICAL APPLICATIONS

We show here the application of the method presented in Sec. III to some physical examples. Let us consider a quantum Fermi system described by a normal ordered Hamiltonian  $H(a_j^\dagger, a_j, t)$ . By working in the fermionic coherent state representation of Bargmann–Fock, we obtain that the integral kernel of the evolution operator is given by the

Grassmann integral,

$$U(\theta_{j(f)}^*, t_f; \theta_{j(i)}^*, t_i) = \int \mathcal{D}(\theta_j^*, \theta_j) \times \exp\left(\frac{1}{2} (\theta_{j(f)}^*, \theta_{j(f)} + \theta_{j(i)}^*, \theta_{j(i)}) + iS\right), \quad (C1)$$

where

$$S = \int_{t_i}^{t_f} dt \left( \frac{1}{2i} (\theta_j^* \dot{\theta}_j - \dot{\theta}_j^* \theta_j) - h(\theta_j^*, \theta_j, t) \right) \quad (C2)$$

and  $h$  is the classical value of  $H$  when we replace the operators  $a_j^\dagger, a_j$  by the Grassmann variables  $\theta_j^*, \theta_j$ . Explicit evaluation of the functional integral in (C1) leads to the evaluation of the action  $S$  on the extremal trajectory, satisfying the classical equations of motion (Hamilton's equations):

$$\frac{1}{i} \dot{\theta}_j = \frac{\partial^l h}{\partial \theta_j^*}, \quad (C3)$$

$$\frac{1}{i} \dot{\theta}_j^* = \frac{\partial^l h}{\partial \theta_j}. \quad (C4)$$

Next, we illustrate the normal form method in three examples involving Hamiltonian GDE.

*Example 1:* Let us consider the class parametrized by  $[s = 1, l = 0]$ . Its Jordan matrix reads as

$$\mathbb{J} = \begin{pmatrix} i\omega & \\ & -i\omega \end{pmatrix}. \quad (C5)$$

For times  $t \gg \sup_\alpha \text{Re}|\gamma_\alpha|^{-1}$  the asymptotic dynamics is described only by two variables  $A_1, A_1^*$ . Since the only available nonlinearity is even, we readily conclude that the normal form is linear:

$$\dot{A}_1 = i\omega A_1, \quad (C6)$$

and corresponds to the classical equation of motion of a particle with spin- $\frac{1}{2}$  and magnetic moment  $\mu$  submitted to the action of a constant magnetic field  $B_Z$  ( $\omega = \mu B_Z$ ). The quantum and classical Hamiltonians are

$$H = B_Z S_z = \frac{1}{2} \mu B_Z (a^\dagger a - a a^\dagger), \quad (C7)$$

$$h = \mu B_Z (A_1^* A_1 - \frac{1}{2}). \quad (C8)$$

*Example 2:* We consider two deuterons with an isospin-isospin type interaction. The quantum Hamiltonian reads ( $\hbar = 1$ )

$$H = \omega_1 (p_1^\dagger p_1 + n_1^\dagger n_1) + \omega_2 (p_2^\dagger p_2 + n_2^\dagger n_2) + g |l_1\rangle \langle l_2|, \quad (C9)$$

where  $p_i^\dagger, p_i$  ( $n_i^\dagger, n_i$ ) are Fermionic creation and annihilation operators of proton (neutron) states and the isospin operators  $l_1, l_2$  are given by

$$l_1 = a^\dagger (\sigma_1)_{ij} a_j, \quad l_2 = b^\dagger (\sigma_2)_{ij} b_j \quad (C10)$$

( $\sigma$  are the Pauli matrices and  $a_1 = p_1, a_2 = n_1, b_1 = p_2, b_2 = n_2$ ). Here,  $H$  in (C9) commutes with the charge  $Q = p_1^\dagger p_1 + p_2^\dagger p_2$  and the baryon number  $B = p_1^\dagger p_1 + p_2^\dagger p_2 + n_1^\dagger n_1 + n_2^\dagger n_2$ . We will prove that any term added to  $H$  that breaks the  $Q$ - $B$  conservation laws is an irrelevant or gauge term since its classical counterpart can be completely eliminated from the classical equations of motion (C3) and

(C4) and therefore the evolution operator (C1) will remain unaffected by the nonphysical term.

For simplicity we consider the Hamiltonian  $H + H'$ , where

$$H' = \lambda (p_1^\dagger p_1 p_2 n_2 + n_2^\dagger p_2^\dagger p_1^\dagger n_1^\dagger). \quad (C11)$$

The classical Hamiltonian equations read as ( $p_1 \rightarrow \theta_1, n_1 \rightarrow \theta_2, p_2 \rightarrow \eta_1, n_2 \rightarrow \eta_2$ )

$$(1/i) \dot{\theta}_1 = \omega_1 \theta_1 + g(2\theta_2 \eta_2^* \eta_1 + \theta_1 (\eta_1^* \eta_1 - \eta_2^* \eta_2)) + \lambda \theta_1 (\eta_1 \eta_2 + \eta_2^* \eta_1^*),$$

$$(1/i) \dot{\theta}_2 = \omega_1 \theta_2 + g(2\theta_1 \eta_1^* \eta_2 + \theta_2 (\eta_2^* \eta_2 - \eta_1^* \eta_1)), \quad (C12)$$

$$(1/i) \dot{\eta}_1 = \omega_2 \eta_1 + g(2\eta_2 \theta_2^* \theta_1 + \eta_1 (\theta_1^* \theta_1 - \theta_2^* \theta_2)) - \lambda \theta_1^* \theta_1 \eta_2^*,$$

$$(1/i) \dot{\eta}_2 = \omega_2 \eta_2 + g(2\eta_1 \theta_1^* \theta_2 + \eta_2 (\theta_2^* \theta_2 - \theta_1^* \theta_1)) + \lambda \theta_1^* \theta_1 \eta_1^*.$$

We note that Eqs. (C12) correspond to the class  $[s = 4, l = 0]$  with Jordan matrix  $\mathbb{J} = \text{diag}(\Omega_1, \Omega_1, \Omega_2, \Omega_2)$ , with  $\Omega_j = \text{diag}(i\omega_j, -i\omega_j), j = 1, 2$ . Following Sec. III and using the invariance property (B4) we conclude that there exists a nonlinear change of variables  $(\theta_i, \eta_i) \rightarrow (A_i)$  such that the equations of motion in the new variables contain only nonlinear terms equivariant under the one-parameter Lie group generated by  $\mathbb{J}^\dagger: e^{-\mathbb{J}^\dagger \tau}$ . We easily see from (C12) that the nonlinear terms with  $g$  in factor respect this symmetry while those with  $\lambda$  in factor break it and therefore can be eliminated.

*Example 3:* We consider two interacting fermionic oscillators with Hamiltonian

$$H = \sum_{k=1}^3 \left( \frac{\omega_1}{2} (a_k^\dagger a_k - a_k a_k^\dagger) + \frac{\omega_2}{2} (b_k^\dagger b_k - b_k b_k^\dagger) + g \sum_{l=1}^3 a_k^\dagger b_l^\dagger a_l b_l \right). \quad (C13)$$

The interaction term has been constructed in such a way that  $H$  commutes with the  $SU(3)$  generators:

$$Q_a = \frac{1}{2} (a_k^\dagger (\lambda_a)_{kl} a_l - b_k^\dagger (\lambda_a^*)_{kl} b_l), \quad (C14)$$

where  $\lambda_a$  are the Gell-Mann matrices.

The classical Hamiltonian  $h$  reads as (up to a constant;  $a_i \rightarrow \theta_i, b_i \rightarrow \eta_i$ )

$$h = \omega_1 \theta_k^* \theta_k + \omega_2 \eta_k^* \eta_k + g \theta_k^* \eta_k^* \theta_l \eta_l \quad (C15)$$

and the equations of motion are

$$(1/i) \dot{\theta}_j = \omega_1 \theta_j + g \eta_j^* \theta_l \eta_l, \quad (C16)$$

$$(1/i) \dot{\eta}_j = \omega_2 \eta_j + g \theta_j^* \theta_l \eta_l.$$

The associated Jordan matrix reads  $\mathbb{J} = \text{diag}(\Omega_1, \Omega_1, \Omega_1, \Omega_2, \Omega_2, \Omega_2)$ , with  $\Omega_j = (i\omega_j, -i\omega_j), j = 1, 2$ , and it is easy to see that the Lie group generated by  $\mathbb{J}^\dagger$  leaves Eqs. (C16) invariant; therefore they are automatically written in normal form. In other words, we have proved that the quantum  $SU(3)$  symmetry induces the classical equivariance under  $e^{\mathbb{J}^\dagger \tau}$  implying the normal form of Hamilton's equations.

<sup>1</sup>F. A. Berezin and M. S. Marinov, Ann. Phys. (NY) **104**, 336 (1977).

<sup>2</sup>F. A. Berezin, Theor. Math. Phys. **6**, 194 (1971).

- <sup>3</sup>Y. Oknuki and S. Kamefuchi, *J. Math. Phys.* **21**, 601 (1980).  
<sup>4</sup>R. Casalbuoni, *Nuovo Cimento A* **33**, 389 (1976).  
<sup>5</sup>M. A. Olshanetsky, *Commun. Math. Phys.* **88**, 63 (1983).  
<sup>6</sup>S. Coleman, *Phys. Rev. D* **11**, 2088 (1975).  
<sup>7</sup>D. Kulish and E. Nissimov, *JETP Lett.* **24**, 247 (1976).  
<sup>8</sup>H. C. Morris, *J. Math. Phys.* **19**, 85 (1978).  
<sup>9</sup>H. D. Wahlquist and F. B. Estabrook, *J. Math. Phys.* **16**, 1 (1975).  
<sup>10</sup>F. A. Berezin, *The Method of Second Quantization* (Academic, New York, 1966).  
<sup>11</sup>V. I. Arnol'd, *Geometrical Methods in the Theory of Ordinary Differential Equations* (Springer, New York, 1977).

# Lie transformations, similarity reduction, and solutions for the nonlinear Madelung fluid equations with external potential

G. Baumann and T. F. Nonnenmacher

*Department of Mathematical Physics, University of Ulm, 7900 Ulm, West Germany*

(Received 20 February 1986; accepted for publication 10 December 1986)

The application of Lie-group methods to a system of coupled nonlinear partial differential equations representing what is usually called a Madelung fluid is shown. The generating operators of the transformation group that depends on five arbitrary group constants will be constructed, and all subclasses of systems of ordinary differential equations derived by similarity reduction will be presented in tabular form. Two subclasses of physical interest are investigated in detail and the similarity solutions are compared with solutions found earlier by the application of inverse scattering transform techniques to the cubic nonlinear Schrödinger equation. Similarity solutions for the Madelung equations with linear external potential  $\Gamma(x) = -f_0x$  are presented.

## I. INTRODUCTION

In recent years considerable interest has been focused on nonlinear evolution equations and their methods of solution. Among the most powerful methods for solving nonlinear partial differential equations are the inverse scattering transform technique (IST) and the Lie-group-based similarity method (LSM) originally initiated by Lie<sup>1</sup> in his classical integration theory. The basic idea of Lie's approach is to study the invariance properties of given differential equations under continuous groups of transformations. If the most extended Lie group of transformations of a given system  $S$  is known, then it is possible to construct classes of particular solutions, called the similarity solutions of  $S$ .

The LSM seems to be applicable to a broader class of nonlinear evolution equations than the IST, which has been applied intensively to nondissipative, nonlinear partial differential equations (NPDE) representing completely integrable Hamiltonian systems.<sup>2-6</sup> Solitons are the most prominent examples of nonlinear solutions constructed by the IST. The class of IST-solvable NPDE possesses a number of common properties such as infinite sequences of conservation laws, Bäcklund transformations with associated geometric and group theoretical properties,<sup>7</sup> and Painlevé transcendental equations characterized by no movable critical points. Relations between solitons and Painlevé-type equations has been pointed out,<sup>8,9</sup> and a unified approach to transformations and elementary solutions of Painlevé equations has been developed by Fokas and Ablowitz.<sup>10</sup>

When discussing dissipative nonlinear evolution equations such as diffusion-type equations<sup>11</sup> or nonlinear kinetic (Boltzmann) equations,<sup>12-19</sup> it is seen that IST-like techniques do not exist thus far. However, in all these cases,<sup>11-16</sup> the LSM has been applied successfully in order to construct similarity solutions. The most prominent examples of these types of solutions are the BKW-mode solutions discovered by Bobylev<sup>18</sup> and Krook and Wu<sup>19</sup> and classified in group-theoretic terms<sup>12,13</sup> by making use of the LSM.

Motivated by a successful application of the LSM on dissipative NPDE, one may ask to which extent the LSM is also applicable to IST-solvable NPDE. Lakshmanan and

Kaliappan<sup>20</sup> have already investigated relations between Lie transformations, nonlinear evolution equations, and Painlevé forms of some pertinent examples of NPDE. They applied the LSM and obtained by a similarity reduction of KdV, sine-Gordon, nonlinear Schrödinger equations, etc., a list of nonlinear ordinary differential equations (NODE), which could be classified in part as Painlevé equations.

The central motivation of this article is to apply the LSM to the nonlinear system of Madelung's quantum fluid equations, which represent a nondissipative system of NPDE for the probability density  $\rho(x,t) = \psi^*\psi$  and the phase function  $S = S(x,t)$ . These are related by the Madelung transformation<sup>21</sup>

$$\psi(x,t) = \sqrt{\rho} \exp(-iS/\hbar), \quad (1.1)$$

when  $h = 2\pi\hbar$  is Planck's constant and  $\psi(x,t)$  and  $\psi^*(x,t)$  satisfy the cubic nonlinear Schrödinger equations

$$i\hbar\psi_t = -(\hbar^2/2m)\psi_{xx} + V(x)\psi + \kappa|\psi|^2\psi \quad (1.2)$$

and their complex conjugate. Here,  $V(x)$  is an external potential and  $\kappa$  is a real-valued constant.

It is well known that the system of NPDE (1.2) is IST integrable.<sup>5,6</sup> We note that (1.2) can be written as a Hamiltonian system by making use of canonical Poisson brackets<sup>6,22,23</sup> and we mention that the canonical transformation (1.1) transforms<sup>22,23</sup> the system of NPDE (1.2) for the complex-valued functions  $\psi$  and  $\psi^*$  into a system of NPDE for real-valued functions  $\rho$  and  $S = m\phi$ , which represents what is usually called a Madelung fluid<sup>21-24</sup>:

$$\phi_t + \frac{1}{2}\phi_x^2 + \alpha\rho + \Gamma(x) = \beta \frac{1}{\sqrt{\rho}} \frac{\partial^2 \sqrt{\rho}}{\partial x^2}, \quad (1.3a)$$

$$\rho_t + \partial_x(\rho\phi_x) = 0, \quad (1.3b)$$

with

$$\alpha = \kappa/m, \quad \beta = \hbar^2/2m^2, \quad \Gamma(x) = (1/m)V(x).$$

We shall apply the LSM to the system ( $S$ ) of NPDE (1.3) in order to study the invariance of  $S$  under continuous groups of transformations depending on one infinitesimal parameter  $\epsilon$ . The most extended ( $\epsilon$ ) Lie group of transformations, admitted by  $S$ , will be shown to depend on five



arbitrary group constants, and the general class of similarity solutions will be seen to separate into different subclasses according to the number of nonzero group constants. It will be analyzed to what extent these subclasses of similarity solutions differ from (or coincide with) classes of solutions already obtained by means of quite different analytical methods.

Most of the previous attention in the literature has been focused on the nonlinear Schrödinger equation (1.2) without an external potential [ $V(x) = 0$ ]. Including a potential of type  $V(x) = -mf_0x$ , Alonso<sup>25</sup> studied the invariance of Eq. (1.2) under the Galilei group from the viewpoint of the IST. In recent times, however, some interest in the Madelung fluid equations (1.3) came from hydrodynamics, in order to describe weakly interacting Bose condensates.<sup>26-28</sup> Analytical solutions of Eq. (1.3) have been found for  $\alpha = 0$  and  $\Gamma(x) = 0$  by applying<sup>29</sup> the Lie-Bäcklund transformation, and for  $\alpha \neq 0$  [ $\Gamma(x) = 0$ ] by transforming the known soliton solutions of the nonlinear Schrödinger equation (1.2) to the fluid dynamical field variables<sup>30,31</sup>  $\rho(x,t)$  and  $\phi(x,t)$ . A systematic application of the IST or the LSM on the system of Madelung fluid equations (1.3) has not been given thus far. Hence we shall apply the LSM on (1.3) in order to construct by similarity reduction the most extended class of ordinary differential equations for the similarity functions. Some of them can be classified as Painlevé-type equations and some exact heretofore undiscovered classes of similarity solutions will be presented.

The present article is organized as follows. First, in Sec. II, we define the symbols and notation for the group generators, similarity forms, etc., by giving a brief reformulation of Lie's basic concept of a one-parameter ( $\epsilon$ ) transformation group in a form that is applicable to the system (1.3) of partial differential equations. In Sec. III we construct the generators of the group, the similarity variables for the most extended class, and some subclasses of physical interest, and derive by reduction the corresponding classes of ordinary differential equations, which are presented in tabular forms. In Sec. IV we construct exact classes of similarity solutions for the density  $\rho(x,t)$  and the phase variable  $S(x,t) = m\phi(x,t)$ , taking into account two choices for the external potential [ $\Gamma(x) = 0$  and  $\Gamma(x) = -f_0x$ ].

## II. LIE GROUP OF TRANSFORMATIONS, GENERATORS OF THE GROUP, AND REDUCTION SCHEME

Consider a system of partial differential equations with two dependent variables  $\rho$  and  $\phi$  and two independent variables  $x$  and  $t$  [as in Eq. (1.3)]:

$$H_i(x,t,\phi, \rho, \phi_x, \rho_x, \phi_t, \rho_t, \phi_{xx}, \rho_{xx}, \dots) = 0, \quad i = 1, 2, \quad (2.1)$$

where subscripts denote partial differentiations. Consider further a one-parameter ( $\epsilon$ ) Lie group of transformations

$$\begin{aligned} x' &= f(x,t, \rho, \phi; \epsilon), & t' &= g(x,t, \rho, \phi; \epsilon), \\ \rho' &= h(x,t, \rho, \phi; \epsilon), & \phi' &= j(x,t, \rho, \phi; \epsilon). \end{aligned} \quad (2.2)$$

Let  $\rho = \theta(x,t)$  and  $\phi = \Xi(x,t)$  be solutions of (2.1). If we replace the variables  $\rho, \phi, x$ , and  $t$  in Eq. (2.1) by  $v, w$ , and  $x' = f(x,t, \theta, \Xi; \epsilon)$ ,  $t' = g(x,t, \theta, \Xi; \epsilon)$ , Eq. (2.1) becomes

$$H_i(x', t', v, w, v_{x'}, w_{x'}, v_{t'}, w_{t'}, v_{x'x'}, w_{x'x'}, \dots) = 0, \quad i = 1, 2. \quad (2.3)$$

Then  $v = \theta(x', t')$  and  $w = \Xi(x', t')$  are solutions of (2.3). We say that the transformations (2.2) leave Eq. (2.1) invariant if  $v = h(x, t, \theta, \Xi; \epsilon)$  and  $w = j(x, t, \theta, \Xi; \epsilon)$  are solutions to (2.3) whenever  $\rho = \theta(x, t)$  and  $\phi = \Xi(x, t)$  are solutions to (2.1). This condition implies that if Eqs. (2.1) and (2.3) have a unique solution, then

$$\begin{aligned} \theta(x', t') &= h(x, t, \theta(x, t), \Xi(x, t); \epsilon), \\ \Xi(x', t') &= j(x, t, \theta(x, t), \Xi(x, t); \epsilon). \end{aligned} \quad (2.4)$$

Hence  $\theta(x, t)$  and  $\Xi(x, t)$  satisfy the one-parameter functional equations

$$\begin{aligned} \theta(f(x, t, \theta, \Xi; \epsilon), g(x, t, \theta, \Xi; \epsilon)) &= h(x, t, \theta, \Xi; \epsilon), \\ \Xi(f(x, t, \theta, \Xi; \epsilon), g(x, t, \theta, \Xi; \epsilon)) &= j(x, t, \theta, \Xi; \epsilon). \end{aligned} \quad (2.5)$$

Expanding (2.2) about the identity  $\epsilon = 0$ , one can generate the following infinitesimal transformations:

$$\begin{aligned} x' &= x + \epsilon \xi_1(x, t, \rho, \phi) + O(\epsilon^2), \\ t' &= t + \epsilon \xi_2(x, t, \rho, \phi) + O(\epsilon^2), \\ \phi' &= \phi + \epsilon \eta^1(x, t, \rho, \phi) + O(\epsilon^2), \\ \rho' &= \rho + \epsilon \eta^2(x, t, \rho, \phi) + O(\epsilon^2). \end{aligned} \quad (2.6)$$

The functions  $\xi_1, \xi_2, \eta^1$ , and  $\eta^2$  are the infinitesimals of the transformations for the variables  $x, t, \phi$ , and  $\rho$ , respectively. In order to find the infinitesimals we need to extend the group to calculate how derivative terms transform. The transformations (2.6), together with the transformations for the first, second, ... derivatives, are called first, second, ... extensions. We denote the infinitesimals for  $\rho_x, \rho_t, \rho_{xx}, \phi_x, \phi_t$ , and  $\phi_{xx}$  by  $\eta_x^2, \eta_t^2, \eta_{11}^2, \eta_1^1, \eta_2^1, \eta_{11}^1$ , respectively. As an example, we give

$$\begin{aligned} \eta_1^1 &= \frac{\partial \eta^1}{\partial x} + \left\{ \frac{\partial \eta^1}{\partial \phi} \phi_x + \frac{\partial \eta^1}{\partial \rho} \rho_x \right\} - \frac{\partial \xi_1}{\partial x} \phi_x - \frac{\partial \xi_2}{\partial x} \phi_t \\ &\quad - \frac{\partial \xi_1}{\partial \phi} \phi_x \phi_x - \frac{\partial \xi_1}{\partial \rho} \rho_x \phi_x - \frac{\partial \xi_2}{\partial \phi} \phi_x \phi_t - \frac{\partial \xi_2}{\partial \rho} \rho_x \phi_t. \end{aligned} \quad (2.7)$$

Similar explicit expressions for higher extensions can be given.<sup>11</sup> Using these various extensions, the infinitesimal criteria for the invariance of (2.1) under the group (2.2) is given by

$$\hat{X}H_i|_{H_i=0} = 0, \quad i = 1, 2, \quad (2.8)$$

where the tangent vector field  $\hat{X}$  is given by

$$\begin{aligned} \hat{X} &= \hat{X}_U + \eta_1^1 \frac{\partial}{\partial \phi_x} + \eta_2^1 \frac{\partial}{\partial \phi_t} + \eta_1^2 \frac{\partial}{\partial \rho_x} + \eta_2^2 \frac{\partial}{\partial \rho_t} \\ &\quad + \eta_{11}^1 \frac{\partial}{\partial \phi_{xx}} + \eta_{11}^2 \frac{\partial}{\partial \rho_{xx}} + \dots, \end{aligned} \quad (2.9a)$$

$$\hat{X}_U = \xi_1 \partial_x + \xi_2 \partial_t + \eta^1 \frac{\partial}{\partial \phi} + \eta^2 \frac{\partial}{\partial \rho}. \quad (2.9b)$$

Here,  $\hat{X}_U$  is the unextended operator. Condition (2.8) provides an algorithm for finding  $\xi_1, \xi_2, \eta^1$ , and  $\eta^2$ . For any solutions  $\rho = \theta(x, t)$  and  $\phi = \Xi(x, t)$  of (2.1), Eq. (2.8) may be treated as a form in derivatives of  $\theta$  and  $\Xi$  whose coeffi-

cients depend on  $(\theta, \Xi, x, t)$  and the unknowns  $(\xi_1, \xi_2, \eta^1, \eta^2)$ . Collecting together the coefficients of like derivative terms in  $\theta$  and  $\Xi$  and setting all of them equal to zero we get a system of linear partial differential equations. In practice these equations are solvable and thus  $\xi_1, \xi_2, \eta^1$ , and  $\eta^2$  are determined.

Our objective now is to find  $\theta(x, t)$  and  $\Xi(x, t)$  given that (2.2) leaves (2.1) invariant. Expanding (2.2) about the identity  $\epsilon = 0$  we generate Eq. (2.6). We now use (2.6) to expand the functional equation (2.5) about  $\epsilon = 0$ . This leads to the following first-order partial differential equations:

$$\xi_1 \phi_x + \xi_2 \phi_t = \eta^1, \quad \xi_1 \rho_x + \xi_2 \rho_t = \eta^2, \quad (2.10)$$

provided that  $\xi_1, \xi_2, \eta^1$ , and  $\eta^2$  are known functions of  $x, t, \phi$ , and  $\rho$ . Equation (2.10) is called the invariant surface condition. The solutions of (2.10) are obtained by solving the following characteristic equations:

$$\frac{dx}{\xi_1} = \frac{dt}{\xi_2} = \frac{d\phi}{\eta^1} = \frac{d\rho}{\eta^2}. \quad (2.11)$$

The general solution of these equations will involve three arbitrary constants, of which one constant takes the role of similarity variable  $\zeta$  and the other constants, say  $F_1(\zeta)$  and  $F_2(\zeta)$ , play the role of dependent variables (usually called similarity functions). Thus we finally obtain the similarity forms

$$v = E(x, t, F_1(\zeta))$$

and

$$w = G(x, t, F_2(\zeta)). \quad (2.12)$$

Substitution of (2.12) into (2.1) results in a system of ordinary differential equations for  $F_1(\zeta)$  and  $F_2(\zeta)$ . The results mentioned above for two dependent and independent variables can be extended to any number of dependent and independent variables. In the following sections we will give an application of the above procedure to the nonlinear Madelung fluid equations (1.3).

### III. APPLICATION TO THE MADELUNG FLUID EQUATIONS

#### A. The infinitesimal elements

In applying the infinitesimal Lie-group methods, a straightforward calculation yields the following infinitesimal elements of the  $\epsilon$ -Lie group [for the external potential we take the choice  $\Gamma(x) = -f_0 x$ ]:

$$\begin{aligned} \xi_1 &= Bx + Ft + H + \frac{3}{2} B f_0 t^2, \\ \xi_2 &= 2Bt + C, \\ \eta^1 &= (F + 3Bf_0 t)x + G + f_0 \left[ \frac{1}{2} B f_0 t^3 + (F/2)t^2 + Ht \right], \\ \eta^2 &= -2B\rho. \end{aligned} \quad (3.1)$$

Thus we obtained a one-parameter group of transformations depending on five arbitrary group constants  $(B, C, F, G, H)$ . For the special case  $\Gamma(x) = 0$ , i.e.,  $f_0 = 0$ , it follows that

$$\begin{aligned} \xi_1 &= Bx + Ft + H, \quad \xi_2 = 2Bt + C, \\ \eta^1 &= Fx + G, \quad \eta^2 = -2B\rho. \end{aligned} \quad (3.2)$$

From the most extended group (3.1) one finds 22 nontrivial subgroups listed in Table I, where taking  $f_0$  equal to zero one

obtains the equations for the Madelung fluid without external potential.

The similarity variables  $\zeta$  as well as the similarity solutions  $\theta(x, t)$  and  $\Xi(x, t)$  are found by solving the characteristic equations (2.11). The results of these integrations depend crucially on the number of vanishing group constants. Here, we will focus our attention on the classes of similarity solutions characterized by the choices  $B = 0$  and by  $B = F = 0$ . It will be shown that these subclasses of similarity solutions lead to solutions of the Painlevé II type and a classical soliton solution.

#### B. Lie algebra constructed from the infinitesimal operators

The knowledge of the infinitesimal elements  $\xi_1, \xi_2, \eta^1$ , and  $\eta^2$  given in Eqs. (3.1) enables us to construct, from the unextended operator  $\hat{X}_\nu$  [Eq. (2.9b)], five operators  $\hat{X}_i$  ( $i = 1, \dots, 5$ ) according to the existence of five group constants  $(B, C, F, G, H)$ . Taking the group constants equal to 0 one obtains from (2.9b) via (3.1) the null operator  $\hat{X}_0 = 0$ . The five generating operators  $\hat{X}_i$  can be constructed by taking one of the group constants equal to 1 and the remaining four constants equal to zero: for  $(B, C, F, G, H) = (1, 0, 0, 0, 0)$  one obtains

$$\begin{aligned} \hat{X}_1 &= \left( x + \frac{3}{2} f_0 t^2 \right) \frac{\partial}{\partial x} + 2t \frac{\partial}{\partial t} \\ &\quad + \left( 3f_0 t x + \frac{f_0^2}{2} t^3 \right) \frac{\partial}{\partial \phi} - 2\rho \frac{\partial}{\partial \rho}; \end{aligned}$$

for  $(0, 1, 0, 0, 0)$  one obtains

$$\hat{X}_2 = \frac{\partial}{\partial t};$$

for  $(0, 0, 1, 0, 0)$  one obtains

$$\hat{X}_3 = t \frac{\partial}{\partial x} + \left( x + \frac{f_0}{2} t^2 \right) \frac{\partial}{\partial \phi};$$

for  $(0, 0, 0, 1, 0)$  one obtains

$$\hat{X}_4 = \frac{\partial}{\partial \phi};$$

and for  $(0, 0, 0, 0, 1)$  one obtains

$$\hat{X}_5 = \frac{\partial}{\partial x} + f_0 t \frac{\partial}{\partial \phi}.$$

These operators must form a Lie algebra. Thus, in general, we have to prove that the commutator of any two operators is a linear combination of these same operators with constant coefficients  $C_{ijm}$  (structure constants), i.e., the commutator relation (closure property)

$$[\hat{X}_i, \hat{X}_j] = C_{ijm} \hat{X}_m$$

must hold and it must be shown that the following properties for the commutators are being satisfied: (i) antisymmetry

$$[\hat{X}_i, \hat{X}_j] = -[\hat{X}_j, \hat{X}_i],$$

and (ii) the Jacobi identity

$$[\hat{X}_i, [\hat{X}_j, \hat{X}_k]] + [\hat{X}_j, [\hat{X}_k, \hat{X}_i]] + [\hat{X}_k, [\hat{X}_i, \hat{X}_j]] = 0.$$

Proving property (i) is equivalent to showing that

TABLE I. Similarity variables and ordinary differential equations for the Madelung fluid with  $\Gamma(x) = -f_0x$ .

Case	Similarity variables	Ordinary differential equations
$B, C, G, F, H \neq 0$	$\zeta = \frac{x - (F/B^2)(Bt + C) + H/B - (f_0/2B^2)(B^2t^2 - 2BCt - 2C^2)}{(2Bt + C)^{1/2}}$ $\phi = (2Bt + C)^{1/2} \left\{ \frac{F}{B} + f_0t - \frac{C}{B} f_0 \right\} \zeta$ $+ \frac{1}{2} \ln(2Bt + C) \left\{ \frac{G}{B} + \frac{CHf_0 - FH}{B^2} + \frac{C^3f_0^2 - 2C^2Ff_0 + CF^2}{2B^3} \right\}$ $+ \frac{f_0^2}{3} t^3 + t^2 \frac{f_0}{B} (F - Cf_0) + t \frac{1}{2B^2} (F^2 - C^2f_0^2 - 2HBf_0) + F_1(\zeta)$ $\rho = F_2(\zeta)/(2Bt + C)$	$(\ddot{F}_1 - 2B)F_2 + \dot{F}_2(\dot{F}_1 - B\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2 \left\{ \frac{B}{\beta}\zeta\dot{F}_1 \right.$ $- \frac{1}{2\beta}\dot{F}_1^2 + \frac{FH}{\beta B} + \frac{C^2Ff_0}{\beta B^2} - \frac{CF^2}{2\beta B^2}$ $\left. - \frac{C^3f_0}{2\beta B^2} - \frac{f_0CH}{\beta B} - \frac{G}{\beta} \right\} = 0$
$B = 0$	$\zeta = x - (F/2C)t^2 - (H/C)t$ $\phi = \frac{F^2}{6C^2} t^3 + \frac{FH}{2C^2} t^2 + \frac{F}{C} \zeta t + \frac{G}{C} t$ $+ f_0 \left( \frac{F}{6C} t^3 + \frac{H}{2C} t^2 \right) + F_1(\zeta)$ $\rho = F_2(\zeta)$	$\ddot{F}_1 F_2 + \dot{F}_2(\dot{F}_1 - H/C) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2 \left\{ \left( \frac{f_0}{\beta} - \frac{F}{\beta C} \right) \zeta \right.$ $\left. + \frac{H}{\beta C} \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 - \frac{G}{\beta C} \right\} = 0$
$C = 0$	$\zeta = [x - (F/B)t + H/B - (f_0/2)t^2] / \sqrt{t}$ $\phi = \sqrt{t} \left\{ \frac{F}{B} + f_0t \right\} \zeta + \frac{1}{2} \ln t \left\{ \frac{G}{B} - \frac{FH}{B^2} \right\}$ $+ \frac{f_0^2}{3} t^3 + \frac{F}{B} f_0 t^2 + t \left\{ \frac{F^2}{2B^2} - \frac{Hf_0}{B} \right\} + F_1(\zeta)$ $\rho = F_2(\zeta)/t$	$(\ddot{F}_1 - 1)F_2 + \dot{F}_2(\dot{F}_1 - \frac{1}{2}\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $\times \left\{ \frac{1}{2\beta}\zeta\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 + \frac{FH}{2\beta B^2} - \frac{G}{2\beta B} \right\} = 0$
$G = 0$	$\zeta = \frac{x - (F/B^2)(Bt + C) + H/B - (f_0/2B^2)(B^2t^2 - 2BCt - 2C^2)}{(2Bt + C)^{1/2}}$ $\phi = (2Bt + C)^{1/2} \left\{ \frac{F}{B} + f_0t - \frac{C}{B} f_0 \right\} \zeta$ $+ \frac{1}{2} \ln(2Bt + C) \left\{ \frac{CHf_0 - FH}{B^2} + \frac{C^3f_0 - 2C^2Ff_0 + CF^2}{2B^3} \right\}$ $+ \frac{f_0^2}{3} t^3 + t^2 \frac{f_0}{B} (F - Cf_0) + t \frac{1}{2B^2} (F^2 - C^2f_0^2 - 2HBf_0) + F_1(\zeta)$ $\rho = F_2(\zeta)/(2Bt + C)$	$(\ddot{F}_1 - 2B)F_2 + \dot{F}_2(\dot{F}_1 - B\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3$ $+ F_2^2 \left\{ \frac{B}{\beta}\zeta\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 + \frac{FH}{\beta B} - \frac{CF^2}{2\beta B^2} \right.$ $\left. + \frac{C^2Ff_0}{\beta B^2} - \frac{C^3f_0}{2\beta B^2} - \frac{HCf_0}{\beta B} \right\} = 0$
$F = 0$	$\zeta = \frac{x + H/B - (f_0/2B^2)(B^2t^2 - 2BCt - 2C^2)}{(2Bt + C)^{1/2}}$ $\phi = (2Bt + C)^{1/2} \left\{ f_0t - \frac{C}{B} f_0 \right\} \zeta$ $+ \frac{1}{2} \ln(2Bt + C) \left\{ \frac{G}{B} + \frac{CHf_0}{B^2} + \frac{C^3f_0^2}{2B^3} \right\}$ $+ \frac{f_0^2}{3} t^3 - \frac{Cf_0^2}{B} t^2 - t \frac{1}{2B^2} (C^2f_0^2 + 2HBf_0) + F_1(\zeta)$ $\rho = F_2(\zeta)/(2Bt + C)$	$(\ddot{F}_1 - 2B)F_2 + \dot{F}_2(\dot{F}_1 - B\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $\times \left\{ \frac{1}{\beta}\zeta\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 \right.$ $\left. - \frac{G}{\beta} - \frac{CHf_0}{\beta B} - \frac{C^3f_0^2}{2\beta B^2} \right\} = 0$
$H = 0$	$\zeta = \frac{x - (F/B^2)(Bt + C) - (f_0/2B^2)(B^2t^2 - 2BCt - 2C^2)}{(2Bt + C)^{1/2}}$ $\phi = (2Bt + C)^{1/2} \left\{ \frac{F}{B} + f_0t - \frac{C}{B} f_0 \right\} \zeta$ $+ \frac{1}{2} \ln(2Bt + C) \left\{ \frac{G}{B} + \frac{C^3f_0^2 - 2C^2Ff_0 + CF^2}{2B^3} \right\}$ $+ \frac{f_0^2}{3} t^3 + t^2 \frac{f_0}{B} (F - Cf_0) + t \frac{1}{2B^2} (F^2 - C^2f_0^2) + F_1(\zeta)$ $\rho = F_2(\zeta)/(2Bt + C)$	$(\ddot{F}_1 - 2B)F_2 + \dot{F}_2(\dot{F}_1 - B\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3$ $+ F_2^2 \left\{ \frac{B}{\beta}\zeta\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 + \frac{C^2Ff_0}{\beta B^2} \right.$ $\left. - \frac{CF^2}{2\beta B^2} - \frac{C^3f_0}{2\beta B^2} - \frac{G}{\beta} \right\} = 0$
$C = 0$ and $G = 0$	$\zeta = [x - (F/B)t + H/B - (f_0/2)t^2] / \sqrt{t}$ $\phi = \sqrt{t} \left\{ \frac{F}{B} + f_0t \right\} \zeta - \frac{1}{2} \frac{FH}{B^2} \ln t + \frac{f_0^2}{3} t^3 + \frac{F}{B} f_0 t^2$ $+ t \left\{ \frac{F^2}{2B^2} - \frac{Hf_0}{B} \right\} + F_1(\zeta)$ $\rho = F_2(\zeta)/t$	$(\ddot{F}_1 - 1)F_2 + \dot{F}_2(\dot{F}_1 - \frac{1}{2}\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $\times \left\{ \frac{1}{2\beta}\zeta\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 + \frac{FH}{2\beta B^2} \right\} = 0$
$G = 0$ and $F = 0$	$\zeta = \frac{x + H/B - (f_0/2B^2)(B^2t^2 - 2BCt - 2C^2)}{(2Bt + C)^{1/2}}$ $\phi = (2Bt + C)^{1/2} \left\{ f_0t - \frac{C}{B} f_0 \right\} \zeta$ $+ \frac{1}{2} \ln(2Bt + C) \left\{ \frac{CHf_0}{B^2} + \frac{C^3f_0^2}{2B^3} \right\}$ $+ \frac{f_0^2}{3} t^3 - \frac{Cf_0^2}{B} t^2 - t \frac{1}{2B^2} (C^2f_0^2 + 2HBf_0) + F_1(\zeta)$	$(\ddot{F}_1 - 2B)F_2 + \dot{F}_2(\dot{F}_1 - B\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $\times \left\{ \frac{B}{\beta}\zeta\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 - \frac{CHf_0}{\beta B} - \frac{C^3f_0^2}{2\beta B^2} \right\} = 0$

TABLE I. (Continued.)

Case	Similarity variables	Ordinary differential equations
$F=0$ and $H=0$	$\rho = F_2(\zeta)/(2Bt + C)$ $\zeta = \frac{x - (f_0/2B^2)(B^2t^2 - 2BCt - 2C^2)}{(2Bt + C)^{1/2}}$ $\phi = (2Bt + C)^{1/2} \left\{ f_0 t - \frac{C}{B} f_0 \right\} \zeta$ $+ \frac{1}{2} \ln(2Bt + C) \left\{ \frac{G}{B} + \frac{C^3 f_0^2}{2B^3} \right\}$ $+ \frac{f_0^2}{3} t^3 - \frac{Cf_0^2}{B} t^2 - \frac{C^2 f_0^2}{2B^2} t + F_1(\zeta)$ $\rho = F_2(\zeta)/(2Bt + C)$	$(\ddot{F}_1 - 2B)F_2 + \dot{F}_2(\dot{F}_1 - B\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $\times \left\{ \frac{1}{\beta}\zeta\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 - \frac{G}{\beta} - \frac{C^3 f_0^2}{2\beta B^2} \right\} = 0$
$B=0$ and $G=0$	$\zeta = x - (F/2C)t^2 + (H/C)t$ $\phi = \frac{F^2}{6C}t^3 + \frac{FH}{2C^2}t^2 + \frac{F}{C}\zeta t + \frac{f_0}{C} \left( \frac{F}{6}t^3 + \frac{H}{2}t^2 \right) + F_1(\zeta)$ $\rho = F_2(\zeta)$	$\ddot{F}_1 F_2 + \dot{F}_2(\dot{F}_1 - H/C) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $\times \left\{ \left( \frac{f_0}{\beta} - \frac{F}{\beta C} \right) \zeta + \frac{H}{\beta C} \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 \right\} = 0$
$C=0$ and $F=0$	$\zeta = [x + H/B - (f_0/2)t^2]/\sqrt{t}$ $\phi = \sqrt{t} f_0 t \zeta + \frac{G}{2B} \ln t + \frac{f_0^2}{3} t^3 - \frac{Hf_0}{B} t + F_1(\zeta)$ $\rho = F_2(\zeta)/t$	$(\ddot{F}_1 - 1)F_2 + \dot{F}_2(\dot{F}_1 - \frac{1}{2}\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $\times \left\{ \frac{1}{2\beta} \zeta \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 - \frac{G}{2\beta B} \right\} = 0$
$G=0$ and $H=0$	$\zeta = \frac{x - (F/B^2)(Bt + C) - (f_0/2B^2)(B^2t^2 - 2BCt - 2C^2)}{(2Bt + C)^{1/2}}$ $\phi = (2Bt + C)^{1/2} \left\{ \frac{F}{B} + f_0 t - \frac{C}{B} f_0 \right\} \zeta$ $+ \frac{1}{2} \ln(2Bt + C) \left\{ \frac{C^3 f_0 - 2C^2 F f_0 + C F^2}{2B^3} \right\}$ $+ \frac{f_0^2}{3} t^3 + t^2 \frac{f_0}{B} (F - C f_0) + t \frac{F^2 - C^2 f_0^2}{2B^2} + F_1(\zeta)$ $\rho = F_2(\zeta)/(2Bt + C)$	$(\ddot{F}_1 - 2B)F_2 + \dot{F}_2(\dot{F}_1 - B\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3$ $+ F_2^2 \left\{ \frac{B}{\beta} \zeta \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 - \frac{CF^2}{2\beta B^2} + \frac{C^2 F f_0}{\beta B^2} - \frac{C^3 f_0^2}{2\beta B^2} \right\} = 0$
$B=0$ and $F=0$	$\zeta = x - (H/C)t$ $\phi = (G/C)t + (f_0 H/2C)t^2 + F_1(\zeta)$ $\rho = F_2(\zeta)$	$\ddot{F}_1 F_2 + \dot{F}_2(\dot{F}_1 - H/C) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $\times \left\{ \frac{f_0}{\beta} \zeta + \frac{H}{\beta C} \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 - \frac{G}{\beta C} \right\} = 0$
$C=0$ and $H=0$	$\zeta = [x - (F/B)t - (f_0/2)t^2]/\sqrt{t}$ $\phi = \sqrt{t} \left\{ \frac{F}{B} + f_0 t \right\} \zeta + \frac{G}{2B} \ln t + \frac{f_0^2}{3} t^3$ $+ \frac{F}{B} f_0 t^2 + \frac{F^2}{2B^2} t + F_1(\zeta)$ $\rho = F_2(\zeta)/t$	$(\ddot{F}_1 - 1)F_2 + \dot{F}_2(\dot{F}_1 - \frac{1}{2}\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $\times \left\{ \frac{1}{2\beta} \zeta \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 - \frac{G}{2\beta B} \right\} = 0$
$B=0$ and $H=0$	$\zeta = x - (F/2C)t^2$ $\phi = \frac{F^2}{6C^2}t^3 + \frac{F}{C}\zeta t + \frac{G}{C}t + \frac{f_0 F}{6C}t^3 + F_1(\zeta)$ $\rho = F_2(\zeta)$	$\ddot{F}_1 F_2 + \dot{F}_2 \dot{F}_1 = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $\times \left\{ \left( \frac{f_0}{\beta} - \frac{F}{\beta C} \right) \zeta - \frac{1}{2\beta} \dot{F}_1^2 - \frac{G}{\beta C} \right\} = 0$
$C=0, G=0,$ and $F=0$	$\zeta = [x + H/B - (f_0/2)t^2]/\sqrt{t}$ $\phi = \sqrt{t} f_0 t \zeta + (f_0^2/3)t^3 - (Hf_0/B)t + F_1(\zeta)$ $\rho = F_2(\zeta)/t$	$(\ddot{F}_1 - 1)F_2 + \dot{F}_2(\dot{F}_1 - \frac{1}{2}\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $[(1/2\beta)\zeta\dot{F}_1 - (1/2\beta)\dot{F}_1^2] = 0$
$G=0, F=0,$ and $H=0$	$\zeta = \frac{x - (f_0/2B^2)(B^2t^2 - 2BCt - 2C^2)}{(2Bt + C)^{1/2}}$ $\phi = (2Bt + C)^{1/2} \left\{ f_0 t - \frac{C}{B} f_0 \right\} \zeta$ $+ \frac{1}{2} \frac{C^3 f_0^2}{2B^3} \ln(2Bt + C)$ $+ \frac{f_0^2}{3} t^3 - \frac{Cf_0^2}{B} t^2 - \frac{C^2 f_0^2}{2B^2} t + F_1(\zeta)$ $\rho = F_2(\zeta)/(2Bt + C)$	$(\ddot{F}_1 - 2B)F_2 + \dot{F}_2(\dot{F}_1 - B\zeta) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$ $\times \left\{ \frac{1}{\beta} \zeta \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 - \frac{C^3 f_0^2}{2\beta B^2} \right\} = 0$
$B=0, G=0,$ and $F=0$	$\zeta = x - (H/C)t$ $\phi = (f_0 H/2C)t^2 + F_1(\zeta)$	$\ddot{F}_1 F_2 + \dot{F}_2(\dot{F}_1 - H/C) = 0$ $F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2$

TABLE I. (Continued.)

Case	Similarity variables	Ordinary differential equations
	$\rho = F_2(\zeta)$	$\times \left\{ \frac{f_0}{\beta} \zeta + \frac{H}{\beta C} \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 \right\} = 0$
$C = 0, F = 0,$ and $H = 0$	$\zeta = [x - (f_0/2)t^2]/\sqrt{t}$ $\phi = \sqrt{t} f_0 t \zeta + (G/2B) \ln t + (f_0^3/3)t^3 + F_1(\zeta)$	$(\ddot{F}_1 - 1)F_2 + \dot{F}_2(\dot{F}_1 - \frac{1}{2}\zeta) = 0$ $F_2 \ddot{F}_2 - \frac{1}{2} \dot{F}_2^2 - \frac{\alpha}{\beta} F_2^3 + F_2^2$ $\times \left\{ \frac{1}{2\beta} \zeta \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 - \frac{G}{2\beta B} \right\} = 0$
$C = 0, G = 0,$ and $H = 0$	$\zeta = [x - (F/B)t - (f_0/2)t^2]/\sqrt{t}$ $\phi = \sqrt{t} \left\{ \frac{F}{B} + f_0 t \right\} \zeta + \frac{f_0^2}{3} t^3 + \frac{F}{B} f_0 t^2 + \frac{F^2}{2B^2} t + F_1(\zeta)$	$(\ddot{F}_1 - 1)F_2 + \dot{F}_2(\dot{F}_1 - \frac{1}{2}\zeta) = 0$ $F_2 \ddot{F}_2 - \frac{1}{2} \dot{F}_2^2 - \frac{\alpha}{\beta} F_2^3 + F_2^2$ $\times \left\{ \frac{1}{2\beta} \zeta \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 \right\} = 0$
$B = 0, G = 0,$ and $H = 0$	$\zeta = x - (F/2C)t^2$ $\phi = (F^2/6C^2)t^3 + (F/C)\zeta t + f_0(F/6C)t^3 + F_1(\zeta)$	$\ddot{F}_1 F_2 + \dot{F}_2 F_1 = 0$ $F_2 \ddot{F}_2 - \frac{1}{2} \dot{F}_2^2 - \frac{\alpha}{\beta} F_2^3 + F_2^2$ $\times \left\{ \left( \frac{f_0}{\beta} - \frac{F}{\beta C} \right) \zeta - \frac{1}{2\beta} \dot{F}_1^2 \right\} = 0$
$C = 0, G = 0, F = 0,$ and $H = 0$	$\zeta = [x - (f_0/2)t^2]/\sqrt{t}$ $\phi = \sqrt{t} f_0 t \zeta + (f_0^3/3)t^3 + F_1(\zeta)$	$(\ddot{F}_1 - 1)F_2 + \dot{F}_2(\dot{F}_1 - \frac{1}{2}\zeta) = 0$ $F_2 \ddot{F}_2 - \frac{1}{2} \dot{F}_2^2 - \frac{\alpha}{\beta} F_2^3 + F_2^2$ $\times \left\{ \frac{1}{2\beta} \zeta \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 \right\} = 0$

$C_{ijk} = -C_{jik}$ ,  
and the Jacobi identity (ii) is equivalent to showing that  
 $C_{ijm} C_{mkn} + C_{jkm} C_{min} + C_{kim} C_{mjn} = 0$ .  
All these properties (closure relation, the property of anti-symmetry, and the Jacobi identity) for the operators  $\hat{X}_i$  follow immediately from Table II (commutator table), which shows, in addition, that  $\hat{X}_4$  is a Casimir operator, i.e., it commutes with all the operators  $\hat{X}_i$ .

Finally we notice that the special choice  $f_0 = 0$  in the operators  $\hat{X}_i$  leads to the Lie algebra for the Madelung fluid equations without external potential.

**C. Similarity variables**

(i) The most extended class of group constants  $(B, C, F, G, H)$ , all unequal to zero, leads to the following similarity variable  $\zeta$  and to the similarity solutions  $\phi$  and  $\rho$ :

$$\zeta = \frac{x - (F/B^2)(Bt + C) + (H/B) - (f_0/2B^2)(B^2t^2 - 2BCt - 2C^2)}{(2Bt + C)^{1/2}}, \tag{3.3a}$$

$$\phi = (2Bt + C)^{1/2} \zeta \left\{ f_0 t - \frac{Cf_0}{B} + \frac{F}{B} \right\} + \frac{1}{2} \ln(2Bt + C) \tag{3.3b}$$

$$\times \left\{ \frac{G}{B} + \frac{CHf_0 - HF}{B^2} + \frac{C^3 f_0^2 - 2C^2 Ff_0 + CF^2}{2B^3} \right\}$$

$$+ \frac{f_0^2}{3} t^3 + \frac{f_0 F - f_0^2 C}{B} t^2 + \frac{F^2 - C^2 f_0^2 - 2BHf_0}{2B^2} t$$

$$+ F_1(\zeta), \tag{3.3b}$$

$$\rho = F_2(\zeta)/(2Bt + C). \tag{3.3c}$$

The similarity functions  $F_1(\zeta)$  and  $F_2(\zeta)$  are to be determined later. The special choice  $\Gamma(x) = -f_0 x = 0$ , i.e.,  $f_0 = 0$ , yields

$$\zeta = [x - (F/B)t - FC/B^2 + H/B]/[(2Bt + C)^{1/2}], \tag{3.4a}$$

$$\phi = \frac{F}{B} (2Bt + C)^{1/2} \zeta + \frac{1}{2} \ln(2Bt + C)$$

$$\times \left\{ \frac{G}{B} - \frac{FH}{B^2} + \frac{F^2 C}{2B^3} \right\} + \frac{F^2}{2B^2} t + F_1(\zeta), \tag{3.4b}$$

$$\rho = F_2(\zeta)/(2Bt + C). \tag{3.4c}$$

(ii) Subclasses

(a)  $(B, C, F, G, H) = (0, C, F, G, H)$

$$\zeta = x - (F/2C)t^2 - (H/C)t, \tag{3.5a}$$

TABLE II. Commutator table for the Madelung fluid with external potential  $\Gamma(x) = -f_0 x$ .

[ , ]	$\hat{X}_0$	$\hat{X}_1$	$\hat{X}_2$	$\hat{X}_3$	$\hat{X}_4$	$\hat{X}_5$
$\hat{X}_0$	0	0	0	0	0	0
$\hat{X}_1$	0	0	$-3f_0 \hat{X}_3 - 2\hat{X}_2$	$\hat{X}_3$	0	$-\hat{X}_5$
$\hat{X}_2$	0	$3f_0 \hat{X}_3 + 2\hat{X}_2$	0	$\hat{X}_5$	0	$f_0 \hat{X}_4$
$\hat{X}_3$	0	$-\hat{X}_3$	$-\hat{X}_5$	0	0	$-\hat{X}_4$
$\hat{X}_4$	0	0	0	0	0	0
$\hat{X}_5$	0	$\hat{X}_5$	$-f_0 \hat{X}_4$	$\hat{X}_4$	0	0

$$\phi = \frac{F^2 + f_0 FC}{6C^2} t^3 + \frac{HF + f_0 HC}{2C^2} t^2 + \frac{F\xi + G}{C} t + F_1(\xi), \quad (3.5b)$$

$$\rho = F_2(\xi). \quad (3.5c)$$

For  $\Gamma(x) = 0$ , this result reduces to

$$\xi = x - (F/2C)t^2 - (H/C)t, \quad (3.6a)$$

$$\phi = \frac{F^2}{6C^2} t^3 + \frac{HF}{2C^2} t^2 + \frac{F\xi + G}{C} t + F_1(\xi), \quad (3.6b)$$

$$\rho = F_2(\xi). \quad (3.6c)$$

(b)  $(B, C, F, G, H) = (0, C, 0, G, H)$ . The choices  $B = 0$  and  $F = 0$  produce the results

$$\xi = x - (H/C)t, \quad (3.7a)$$

$$\phi = (f_0 H/2C)t^2 + (G/C)t + F_1(\xi), \quad (3.7b)$$

$$\rho = F_2(\xi). \quad (3.7c)$$

For  $\Gamma(x) = 0$  one obtains

$$\xi = x - (H/C)t, \quad (3.8a)$$

$$\phi = (G/C)t + F_1(\xi), \quad (3.8b)$$

$$\rho = F_2(\xi). \quad (3.8c)$$

We note that the similarity functions  $F_1(\xi)$  and  $F_2(\xi)$  obey ordinary differential equations, which are obtained by insertion of the corresponding expressions for  $\phi$  and  $\rho$  into the original partial differential equations (1.3).

#### D. Reduction to ordinary differential equations

(i) The most extended class [all group constants  $(B, C, F, G, H)$  unequal to zero] leads—by making use of the results given in (3.3)—to the following system of ordinary differential equations (ODE) for  $F_1(\xi)$  and  $F_2(\xi)$ :

$$(\ddot{F}_1 - 2B)F_2 + \dot{F}_2(\dot{F}_1 - B\xi) = 0, \quad (3.9a)$$

$$F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2\left\{\frac{B}{\beta}\xi\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 + \frac{2C^2Ff_0 - CF^2 - C^3f_0^2}{2\beta B^2} + \frac{FH - f_0CH}{\beta B} - \frac{G}{\beta}\right\} = 0. \quad (3.9b)$$

Taking  $\Gamma(x) = 0$  we obtain

$$(\ddot{F}_1 - 2B)F_2 + \dot{F}_2(\dot{F}_1 - B\xi) = 0, \quad (3.10a)$$

$$F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2\left\{\frac{B}{\beta}\xi\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 - \frac{2B}{\beta}\left(\frac{F^2C}{4B^3} + \frac{G}{2B} - \frac{FH}{2B^2}\right)\right\} = 0. \quad (3.10b)$$

(ii) Subclasses

(a)  $(B, C, F, G, H) = (0, C, F, G, H)$  The results (3.5) lead to

$$\ddot{F}_1 F_2 + \dot{F}_2(\dot{F}_1 - H/C) = 0, \quad (3.11a)$$

$$F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2\left\{\frac{H}{\beta C}\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 - \left(\frac{1}{\beta}\left(\frac{F}{C} - f_0\right)\xi + \frac{G}{\beta C}\right)\right\} = 0, \quad (3.11b)$$

and for  $\Gamma(x) = 0$  one obtains

$$\ddot{F}_1 F_2 + \dot{F}_2(\dot{F}_1 - H/C) = 0, \quad (3.12a)$$

$$F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2\left\{\frac{H}{\beta C}\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 - \frac{1}{\beta}\left(\frac{F}{C}\xi + \frac{G}{C}\right)\right\} = 0. \quad (3.12b)$$

(b)  $(B, C, F, G, H) = (0, C, 0, G, H)$  The results (3.7) yield

$$\ddot{F}_1 F_2 + \dot{F}_2(\dot{F}_1 - H/C) = 0, \quad (3.13a)$$

$$F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2\left\{\frac{H}{\beta C}\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 + \frac{f_0}{\beta}\xi - \frac{G}{\beta C}\right\} = 0. \quad (3.13b)$$

Taking  $\Gamma(x) = 0$  one obtains

$$\ddot{F}_1 F_2 + \dot{F}_2(\dot{F}_1 - H/C) = 0, \quad (3.14a)$$

$$F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2\left\{\frac{H}{\beta C}\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 - \frac{G}{\beta C}\right\} = 0. \quad (3.14b)$$

All other ODE are listed in Table I.

#### IV. EXACT CLASSES OF SIMILARITY SOLUTIONS

In Sec. III we reduced the Madelung fluid equations represented by a set of nonlinear partial differential equations (NPDE) (1.3) to various systems of nonlinear ordinary differential equations (NODE) for different choices of the set of group constants. Here we will construct analytical solutions of some special classes of NODE.

The first objective in the study of such NODE is to ascertain whether or not a solution can be obtained either explicitly or implicitly in terms of classical functions. The purpose in such a study is to discover a transformation which will reduce the equation to some type that is known to have a solution of the desired kind. Failing this, one seeks a transformation which will reduce the equation to one that is asymptotic to a form solvable by known functions.

##### A. Solutions of elliptic form

Following along this line we will give, as a first example, the solution of the coupled nonlinear equations:

$$\ddot{F}_1 F_2 + \dot{F}_2(\dot{F}_1 - H/C) = 0 \quad (4.1)$$

and

$$F_2\ddot{F}_2 - \frac{1}{2}\dot{F}_2^2 - \frac{\alpha}{\beta}F_2^3 + F_2^2 \times \left\{\frac{H}{\beta C}\dot{F}_1 - \frac{1}{2\beta}\dot{F}_1^2 - \frac{G}{\beta C}\right\} = 0, \quad (4.2)$$

where overdots denote differentiation with respect to  $\xi$ , the similarity variable. This system is just subclass (ii) of case (b) from Sec. III [Eq. (3.14)] without external potential.

The first of these two equations can be integrated to give

$$F_1(\xi) = \int_{\xi_0}^{\xi} \left(\frac{I_1}{F_2(\xi')} + \frac{H}{C}\right) d\xi', \quad (4.3)$$

where  $I_1$  is a constant of integration. Substituting Eq. (4.3)

into the second of the coupled system, Eq. (4.2), it reduces to the form

$$F_2 \ddot{F}_2 - \frac{1}{2} \dot{F}_2^2 - \frac{\alpha}{\beta} F_2^3 - F_2^2 \times \left\{ \frac{G}{\beta C} - \frac{1}{2\beta} \left( \frac{H}{C} \right)^2 \right\} - \frac{I_1^2}{2\beta} = 0, \quad (4.4)$$

which is an equation for  $F_2$  alone. By means of the transformation

$$\chi^2 = F_2, \quad (4.5)$$

Eq. (4.4) takes the form

$$\ddot{\chi} - \frac{\alpha}{\beta} \chi^3 - \frac{1}{2\beta} \omega \chi - \frac{I_1^2}{4\beta} \frac{1}{\chi^3} = 0, \quad (4.6)$$

with

$$\omega = \{G/C - \frac{1}{2}(H/C)^2\}.$$

Introducing a "potential"  $V(\chi)$  by

$$V(\chi) = -\frac{\alpha}{2\beta} \frac{1}{4} \chi^4 - \frac{1}{2\beta} \omega \frac{1}{2} \chi^2 + \frac{I_1^2}{8\beta} \frac{1}{\chi^2}, \quad (4.7)$$

one can write Eq. (4.6) as

$$\frac{d}{d\xi} \left( \frac{1}{2} \dot{\chi}^2 + V(\chi) \right) = 0. \quad (4.8)$$

By integrating Eq. (4.8) twice it follows that

$$\int^\xi d\xi' = \int^x \frac{d\chi'}{\sqrt{2(E - V(\chi'))}}, \quad (4.9)$$

where  $E$  is a constant of first integration. The potential  $V(\chi)$  given by (4.7) determines this solution in a characteristic way. To see this the following transformation of both the dependent and the independent variables

$$\tilde{\xi} = \sqrt{\alpha/\beta} \xi \quad (4.10)$$

and

$$y = \chi^2 + (1/6\alpha)\omega$$

leads to a standard form of Weierstrass elliptic functions

$$p^{-1}(y) = \int \frac{dy}{\sqrt{4y^3 - g_2 y - g_3}}, \quad (4.11)$$

where  $g_2$  and  $g_3$  are given by

$$g_2 = (3/16\alpha^2)\omega^2 - 8(E\beta/\alpha),$$

$$g_3 = \frac{4}{3} \frac{E\beta}{\alpha^2} \omega - \frac{I_1^2}{\alpha\beta} - \frac{1}{16\alpha^3} \omega^3.$$

If we now invert all our previously used transformations applied to  $\rho$  and  $\phi$ , and to  $\xi$  as well we can write the solution for the density  $\rho$  and phase  $\phi$  as

$$\rho(x,t) = p\{\sqrt{\alpha/\beta} [x - (H/C)t] - I_4\} - (1/6\alpha)\omega \quad (4.12)$$

and

$$\phi(x,t) = \frac{G}{C} + \frac{H}{C} \left( x - \frac{H}{C} t \right) + \int \frac{I_1 d\xi}{p(\sqrt{\alpha/\beta} \xi - I_4) - (1/6\alpha)(\omega/6\alpha)}, \quad (4.13)$$

where  $I_1$  and  $I_4$  are arbitrary integration constants and  $\xi$  has the form of a "moving wave variable"  $\xi = x - (H/C)t$ , i.e., the solutions represent traveling nonlinear waves.

It is well known that there exists a close connection between Weierstrass elliptic functions and Jacobi elliptic functions. To see how this relation appears in this system of equations we turn our attention back to Eq. (4.3). If we set there the arbitrary constant  $I_1$  equal to zero, one can write a first integral of (4.1) as

$$F_1(\xi) = (H/C)\xi + \tilde{I}_1, \quad (4.14)$$

where  $\tilde{I}_1$  is an arbitrary constant of integration. This form of  $F_1(\xi)$  solves (4.1) for any  $F_2(\xi)$ . However, there exists on the other hand a strong coupling between  $F_1$  and  $F_2$  via (4.2). Inserting the solution (4.14) into Eq. (4.2) one can discard this coupling. The result is

$$F_2 \ddot{F}_2 - \frac{1}{2} \dot{F}_2^2 - (\alpha/\beta) F_2^3 - (\omega/\beta) F_2^2 = 0, \quad (4.15)$$

which is a nonlinear ordinary differential equation for  $F_2$  with constant coefficients. By means of the transformation (4.5), Eq. (4.15) is reduced to the form

$$\ddot{\chi} - (\alpha/2\beta)\chi^3 - (1/2\beta)\omega\chi = 0. \quad (4.16)$$

Introduction of the "potential"  $V_1(\chi)$  results in the reduction to

$$\ddot{\chi} = -\frac{\partial V_1(\chi)}{\partial \chi}, \quad (4.17)$$

where  $V_1(\chi)$  is a polynomial in  $\chi$  of fourth order:

$$V_1(\chi) = -\frac{\alpha}{2\beta} \frac{1}{4} \chi^4 - \frac{1}{2\beta} \omega \frac{1}{2} \chi^2. \quad (4.18)$$

Here we dropped an arbitrary constant in  $V_1$ . We note that this potential is a special case of  $V$  resulting from (4.7) by taking  $I_1 = 0$ . An implicit solution for  $\chi$  follows by integration of (4.17). Because  $V_1(\chi)$  is a polynomial of fourth order this solution belongs to the class of Jacobi elliptic functions. The formal solution is

$$\xi = \int^x \frac{d\chi'}{\sqrt{2(E - V(\chi'))}}, \quad (4.19)$$

where  $E$  is an integration parameter which strongly determines the behavior of this solution. For  $\alpha < 0$ ,  $\beta > 0$ , and  $\omega > 0$  integration of (4.19) can be carried out, with the result

$$\chi(\xi) = b_1 \operatorname{cn}(\sqrt{(|\alpha|/4\beta)(a_1^2 + b_1^2)}(\xi - \xi_0), m), \quad (4.20)$$

where  $\operatorname{cn}$  is the cnoidal elliptic function. Here,  $a_1$ ,  $b_1$ , and the modulus  $m$  are given by

$$a_1^2 = -\omega/|\alpha| + \sqrt{(\omega/|\alpha|)^2 + 8\beta E/|\alpha|}, \quad (4.21)$$

$$b_1^2 = \omega/|\alpha| + \sqrt{(\omega/|\alpha|)^2 + 8\beta E/|\alpha|},$$

$$m = b_1/\sqrt{a_1^2 + b_1^2}. \quad (4.22)$$

The initial phase  $\xi_0$  can be expressed by an initial condition for  $\chi$ :

$$\xi_0 = \frac{1}{\sqrt{(|\alpha|/4\beta)(a_1^2 + b_1^2)}} \operatorname{cn}^{-1}(\chi(0)/b_1, m),$$

where  $\chi(0)$  is the value of  $\chi$  at  $\xi = 0$ .

Elliptic functions such as those given in (4.20) belong to

the class of doubly periodic functions with  $2K$ , and  $K$  is the complete elliptic integral of the first kind. Another interesting property of this cnoidal function is the fact that  $\text{cn}$  tends to the hyperbolic function  $\text{sech}$  if the modulus  $m$  is equal to unity, i.e.,  $\text{cn} \rightarrow \text{sech}$  if  $m \rightarrow 1$ .

If we suppose that our group constants  $C, G, H$  are fixed and  $\alpha, \beta$  are finite-valued constants then the modulus  $m$  is a function of  $E$ . Setting this parameter equal to zero  $m$  is equal to 1. In this way solution (4.20) skips to the hyperbolic function  $\text{sech}$ , as noted above. In detail, this solution is given by

$$\chi(\xi) = (2\omega/|\alpha|)^{1/2} \text{sech}(\sqrt{\omega/2\beta}(\xi - \xi_0)). \quad (4.23)$$

Inverting the previously used transformations we can give an explicit solution for the density and phase of the Schrödinger field  $\psi$ :

$$\rho(x,t) = \frac{2\omega}{|\alpha|} \text{sech}^2\left(\sqrt{\frac{\omega}{2\beta}}\left(x - \frac{H}{C}t - \left(x_0 - \frac{H}{C}t_0\right)\right)\right) \quad (4.24)$$

and

$$\phi(x,t) = (H/C)[x - (H/C)t] + G/C + I_3, \quad (4.25a)$$

where

$$\omega = G/C - \frac{1}{2}(H/C)^2. \quad (4.25b)$$

A graphical representation of these solutions in space-time is given in Fig. 1 for  $\rho$ . One notes that the pulselike solution moves with constant velocity through space.

As stated in the Introduction, the Schrödinger field  $\psi$  is related to the density  $\rho$  and phase  $\phi$  via the Madelung transformation (1.1), where  $S(x,t) = m\phi(x,t)$ . Inserting the density and phase into this transformation one finds

$$\psi(x,t) = \sqrt{\frac{2\omega}{|\alpha|}} \text{sech}\left(\sqrt{\frac{\omega}{2\beta}}\left(x - \frac{H}{C}t - \left(x_0 - \frac{H}{C}t_0\right)\right)\right) \times \exp\left(-\frac{i}{\hbar}m\left(\frac{H}{C}\left(x - \frac{H}{C}t\right) + \frac{G}{C} + I_3\right)\right).$$

This representation of  $\psi$  has just the same form as Alonso's<sup>25</sup> result constructed by using IST techniques.

## B. Connection to the Painlevé II function $P_{II}$

The second system of coupled nonlinear ordinary differential equations, which we investigate, includes the equations of subclasses (ii) (a) [Eqs. (3.11) and (3.12)] and the

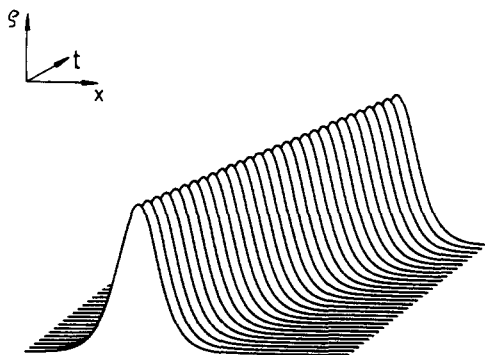


FIG. 1. Classical soliton solution (4.24) of the Madelung fluid for the group constants  $(0, C, 0, G, H)$  and for  $f_0 = 0$ .

first equation of (ii) (b) [Eq. (3.13)] of Sec. III. This system has [in contrast to the previously discussed equations (4.1) and (4.2)] analytical coefficients. It reads as

$$\ddot{F}_1 F_2 + \dot{F}_2 (\dot{F}_1 - H/C) = 0 \quad (4.26a)$$

and

$$F_2 \ddot{F}_2 - \frac{1}{2} \dot{F}_2^2 - \frac{\alpha}{\beta} F_2^3 + F_2^2 \times \left\{ \frac{H}{\beta C} \dot{F}_1 - \frac{1}{2\beta} \dot{F}_1^2 - \frac{1}{\beta} \left[ \gamma \xi + \frac{G}{C} \right] \right\} = 0, \quad (4.26b)$$

where  $\gamma$  is given by the relation

$$\gamma = \begin{cases} \gamma_1 = F/C, & \text{for } f_0 = 0 \text{ and } (0, C, F, G, H), \\ \gamma_2 = F/C - f_0, & \text{for } f_0 \neq 0 \text{ and } (0, C, F, G, H), \\ \gamma_3 = -f_0, & \text{for } f_0 \neq 0 \text{ and } (0, C, 0, G, H). \end{cases} \quad (4.26c)$$

The independent variable  $\xi$  takes the form

$$\xi = x - (F/2C)t^2 - (H/C)t, \quad \text{for } \gamma = \gamma_1,$$

$$\xi = x - (F/2C)t^2 - (H/C)t, \quad \text{for } \gamma = \gamma_2, \quad (4.26d)$$

$$\xi = x - (H/C)t, \quad \text{for } \gamma = \gamma_3.$$

Taking  $\gamma$  equal to zero, i.e., either  $F$  or  $f_0$  equal to zero, it follows the previously discussed ODE system. The main difference between the two systems is a  $\xi$ -dependent term connected with  $\gamma$ . Thus we have to solve a coupled system of equations that possesses a dependence on independent variables. To solve these equations we proceed in the same way as in Sec. IV A. First we integrate (4.26a) and obtain the function  $F_1$  as an implicit solution of  $F_2$ :

$$F_1(\xi) = \int_{\xi_0}^{\xi} \left( \frac{I_1}{F_2(\xi')} + \frac{H}{C} \right) d\xi'. \quad (4.27)$$

For  $I_1 = 0$  it follows that

$$F_1(\xi) = (H/C)\xi + I_2. \quad (4.28)$$

Substituting this result in (4.26b) leads to

$$F_2 \ddot{F}_2 - \frac{1}{2} \dot{F}_2^2 - (\alpha/\beta) F_2^3 - (1/\beta) F_2^2 \{ \gamma \xi + \omega \} = 0, \quad (4.29)$$

where  $\omega$  is defined in (4.25b). Transforming both the dependent and independent variables by

$$z = (2\beta/\gamma)^{1/2} (1/2\beta) \{ \gamma \xi + \omega \} \quad (4.30a)$$

and

$$\chi^2 = (\alpha/4\beta) (2\beta/\gamma)^{2/3} F_2 \quad (4.30b)$$

one obtains

$$\chi'' - 2\chi^3 - z\chi(z) = 0, \quad (4.31)$$

where primes denote differentiation with respect to  $z$ . This second-order nonlinear differential equation is a special case of the so-called Painlevé II type equation. As shown by Painlevé,<sup>32</sup> this type of equation is irreducible and possesses neither movable branch points nor essential singularities. On the other hand, this equation defines a new class of transcendental functions. We see (as a by-product) that—following Ablowitz *et al.*<sup>9</sup>—the Madelung fluid is of the IST type, which is not surprising since the corresponding system of nonlinear Schrödinger equations (1.2) is also IST integrable.



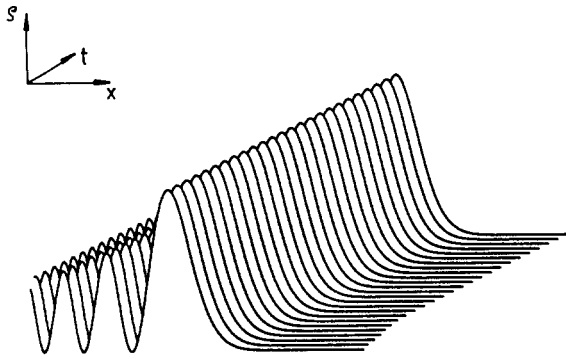


FIG. 2. Representation of the asymptotic solution (4.33) for small  $\rho$  with group constants  $(0, C, 0, G, H)$  and external potential.

To follow our introductory proposals we now approximate (4.31) in such a way that at least an asymptotic analytical solution for this equation can be given. We require that this solution decay rapidly enough as  $|x| \rightarrow \infty$  (say) that the integral of  $\rho$  is defined, i.e.,  $\rho \rightarrow 0$  for  $|x| \rightarrow \infty$ . As stated above,  $\rho$  is connected with  $\chi$  by

$$\rho \propto \chi^2.$$

If we take the positive square root of  $\rho$  we can discuss the whole solution in terms of  $\chi$ . Assuming now for the asymptotic solution that  $\rho$  is a small quantity, we can neglect products of  $\chi$  in Eq. (4.31), which is then reduced to

$$\chi'' - z\chi(z) = 0. \quad (4.32)$$

Equation (4.32) is the definition equation of Airy functions in differential form.<sup>33</sup> The approximate solution for small  $\rho$  is therefore given by

$$\chi \propto \text{Ai}(z). \quad (4.33)$$

The asymptotic solutions of Airy functions are<sup>33</sup>

$$\begin{aligned} \chi &\propto (1/2\sqrt{\pi})(1/z^{1/4})\exp(-\frac{2}{3}(z)^{3/2}), \quad \text{for } z \rightarrow \infty, \\ \chi &\propto (1/z^{1/4})\sin(\frac{2}{3}|z|^{3/2} + \phi_0), \quad \text{for } z \rightarrow -\infty. \end{aligned}$$

Up to now we have solved a coupled nonlinear system of equations with an arbitrary independent variable  $\zeta$ . These solutions are

$$\rho \propto P_{\text{II}}^2(\zeta),$$

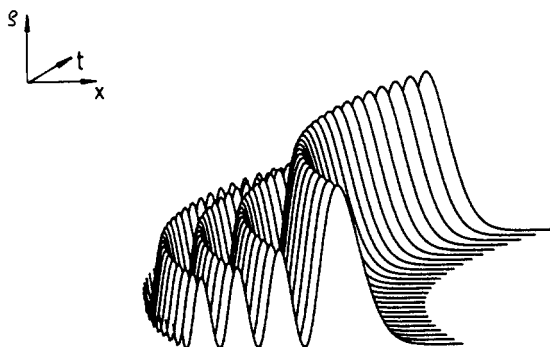


FIG. 3. Asymptotic solution for small  $\rho$  with  $(0, C, F, G, H)$  and  $\gamma = \gamma_1$  or  $\gamma = \gamma_2$ .

$$\phi = \begin{cases} (3.5b) \\ (3.6b), & \text{with } F_1 = (H/C)\zeta + I_2, \\ (3.7b) \end{cases}$$

and  $\zeta$  as in (4.26d). If we now specify the explicit combination of  $x$  and  $t$  in  $\zeta$  we have to distinguish two cases strongly dependent on the values of the group constant  $F$ . Taking the choice  $B = 0$  and  $F = 0$  one gets an ordinary traveling wave with

$$\zeta = x - (H/C)t. \quad (4.34a)$$

The solution for the density  $\rho(\zeta) = \rho(x, t)$  is shown in Fig. 2. The choice  $F \neq 0$ , i.e., regarding the more general subgroup  $(0, C, F, G, H)$ , one obtains the similarity variable

$$\zeta = x - (F/2C)t^2 - (H/C)t. \quad (4.34b)$$

This is somewhat like an "accelerated wave variable." The corresponding solution  $\rho(x, t)$  is plotted in Fig. 3. In contrast to the soliton- or solitary wave-type solutions based on  $\zeta$  given in (4.34a), the solutions based on the similarity variable (4.34b) are called boomerons in the literature (see Ref. 34). Finally, we note that the soliton-type solution based on (4.34a) is, in addition to its Lie-group properties, Galilei invariant as well, while the boomeron solution based on (4.34b) is not invariant with respect to a Galilei transformation.

<sup>1</sup>S. Lie, *Theorie der Transformationsgruppen* (Teubner, Leipzig, 1980), 2nd ed. (reprinted by Chelsea, New York, 1970).

<sup>2</sup>C. S. Gardner, J. M. Green, M. D. Kruskal, and R. M. Miura, *Commun. Pure Appl. Math.* **27**, 97 (1974).

<sup>3</sup>P. D. Lax, *Commun. Pure Appl. Math.* **21**, 467 (1968).

<sup>4</sup>M. J. Ablowitz, D. J. Kaup, A. C. Newell, and H. Segur, *Stud. Appl. Math.* **LIII**, 240 (1974).

<sup>5</sup>V. E. Zakharov and A. B. Shabat, *Sov. Phys. JETP* **34**, 62 (1972).

<sup>6</sup>M. D. Arthur, K. M. Case, *J. Math. Phys.* **23**, 1771 (1982).

<sup>7</sup>R. K. Bullough and P. J. Caudry, *Solitons* (Springer, New York, 1980).

<sup>8</sup>M. J. Ablowitz, A. Ramani, and H. Segur, *J. Math. Phys.* **21**, 715 (1980).

<sup>9</sup>M. J. Ablowitz, A. Ramani, and H. Segur, *J. Math. Phys.* **21**, 1014 (1980).

<sup>10</sup>A. S. Fokas and M. J. Ablowitz, *J. Math. Phys.* **23**, 2033 (1982).

<sup>11</sup>G. W. Bluman and J. D. Cole, *Appl. Math. Sci.* **13**, 1 (1974).

<sup>12</sup>G. Tenti and W. H. Hui, *J. Math. Phys.* **19**, 744 (1978).

<sup>13</sup>T. F. Nonnenmacher, *J. Appl. Math. Phys.* **35**, 680 (1984).

<sup>14</sup>T. F. Nonnenmacher and G. Dukek, *Physica* **135A**, 167 (1986).

<sup>15</sup>G. Dukek and T. F. Nonnenmacher, in *Applications of Mathematics in Technology*, edited by V. C. Boffi and H. Neunzert (Teubner, Stuttgart, 1984).

<sup>16</sup>V. C. Boffi and T. F. Nonnenmacher, *Nuovo Cimento B* **85**, 165 (1985).

<sup>17</sup>G. Spiga, T. F. Nonnenmacher and V. C. Boffi, *Physica* **131A**, 431 (1985).

<sup>18</sup>A. V. Bobylev, *Sov. Phys. Dokl.* **20**, 820 (1976).

<sup>19</sup>M. Krook and T. T. Wu, *Phys. Rev. Lett.* **36**, 1107 (1976).

<sup>20</sup>M. Lakshmanan and P. Kaliappan, *J. Math. Phys.* **24**, 795 (1983).

<sup>21</sup>E. Madelung, *Z. Phys.* **40**, 322 (1926).

<sup>22</sup>T. F. Nonnenmacher, *Lect. Notes Phys.* **253**, 149 (1986).

<sup>23</sup>T. F. Nonnenmacher, G. Dukek and G. Baumann, *Lett. Nuovo Cimento* **36**, 453 (1983).

<sup>24</sup>F. Guerra and R. Marra, *Phys. Rev. D* **28**, 1916 (1983).

<sup>25</sup>L. M. Alonso, *J. Math. Phys.* **23**, 1518 (1982).

<sup>26</sup>L. J. F. Broer, *Physica* **76**, 364 (1974).

<sup>27</sup>C. A. Jones and P. H. Roberts, *J. Phys. A: Math. Gen.* **15**, 2599 (1982).

<sup>28</sup>S. J. Putterman and P. H. Roberts, *Physica* **117A**, 369 (1983).

<sup>29</sup>A. J. Purcell, *Phys. Rev. D* **30**, 2128 (1984).

<sup>30</sup>T. F. Nonnenmacher and J. D. F. Nonnenmacher, *Lett. Nuovo Cimento* **37**, 241 (1983).

<sup>31</sup>T. A. Minelli and A. Pascolini, *Nuovo Cimento* **85B**, 1 (1985).

<sup>32</sup>C. R. Painlevé, *Acta Math.* **25**, 13 (1902).

<sup>33</sup>M. Abramowitz and I. A. Stegun, *Handbook of Mathematical Functions*

(Dover, New York, 1972).

<sup>34</sup>A. Degasperis, in *Nonlinear Evolution Equations Solvable by the Spectral Transform*, edited by F. Calogero (Pitman, London, 1978).

# Some aspects of the isogroup of the self-dual Yang–Mills system

C. J. Papachristou and B. Kent Harrison

*Department of Physics and Astronomy, Brigham Young University, Provo, Utah 84602*

(Received 8 October 1986; accepted for publication 4 February 1987)

A generalized isovector formalism is used to derive the isovectors and isogroup of the self-dual Yang–Mills (SDYM) equation in the so-called  $J$  formulation. In particular, the infinitesimal “hidden symmetry” transformation, a linear system, and a well-known Bäcklund transformation of the SDYM equation are derived in the process. Thus symmetry and integrability aspects of the SDYM system appear in natural relationship to each other within the framework of the isovector approach.

## I. INTRODUCTION

In a recent paper<sup>1</sup> the authors discussed the application of isovector techniques<sup>2,3</sup> to systems of partial differential equations corresponding to exterior equations for vector-valued (and, in particular, matrix-valued) differential forms. It was seen that the application of the Lie derivative operator on vector-valued one-forms presents some technical difficulties, and for this reason an *internal exterior derivative* (i.e., an exterior derivative that acts on the fields but not on the variables of the solution manifold) was introduced by the formula

$$\bar{d}F(x^\mu, \psi^i) \equiv dF - \partial_\mu F dx^\mu, \quad (1.1)$$

where  $F$  is any function of the scalar variables  $x^\mu$  of the solution manifold and the vector-valued fields  $\psi^i$ . If the system of partial differential equations (PDE's) is of order 2 or higher, the variables  $\psi^i$  will comprise the dependent variables  $u^a$  of the PDE's and the derivatives, up to a certain degree, of the  $u^a$  with respect to the  $x^\mu$ . Given that, in the absence of specific restrictions on the exterior differential forms that represent the system, the variables  $\psi^i$  are considered independent of each other (and of the  $x^\mu$ ), we conclude that the problem can be naturally formulated on a jet space with “mixed” (i.e., both scalar- and vector-valued) coordinates.

In the present paper the formalism developed in Ref. 1 is applied to the self-dual Yang–Mills (SDYM) equation in the so-called  $J$  formulation.<sup>4</sup> It is seen that the isovector method provides a natural framework for the unification of such distinct concepts as symmetry and integrability. The independence of the coordinates of the underlying jetlike space is important in this context, as the reader will realize. In Sec. II we calculate the isovectors of the SDYM system. These vector fields can be used to construct infinitesimal symmetries (both geometrical and internal) of the system, as discussed in Ref. 1. The above-mentioned independence of coordinates is used in Sec. III to rewrite certain symmetries in a form equivalent to the parametric infinitesimal transformation introduced in Ref. 5. (This transformation is related to the so-called hidden symmetry of the SDYM field.<sup>6</sup>) Remarkably, the process also yields a pair of linear “inverse scattering” equations, the integrability of which is equivalent to the SDYM equation, and the parameter of which is identical to that of the infinitesimal transformation mentioned above. Finally, the results of Secs. I–III are used

in Sec. IV to derive Bäcklund transformations for the SDYM system. In particular, the process gives the parametric Bäcklund transformation proposed in Ref. 7.

## II. ISOVECTORS OF THE SDYM SYSTEM

The SDYM equation in the  $J$  formulation is written as<sup>4–7</sup>

$$\partial_{\bar{y}}(J^{-1} \partial_y J) + \partial_{\bar{z}}(J^{-1} \partial_z J) = 0. \quad (2.1)$$

The complex coordinates  $y, z, \bar{y}$ , and  $\bar{z}$  are related to the coordinates  $x_1, x_2, x_3$ , and  $x_4$  of complexified Euclidean space by

$$\begin{aligned} 2^{1/2}y &= x_1 + ix_2, & 2^{1/2}z &= x_3 - ix_4, \\ 2^{1/2}\bar{y} &= x_1 - ix_2, & 2^{1/2}\bar{z} &= x_3 + ix_4. \end{aligned} \quad (2.2)$$

[Note that the pairs  $(y, \bar{y})$  and  $(z, \bar{z})$  involve elements that are complex-conjugately related in *real* Euclidean space.] For our purposes,  $J$  is assumed to be a nonsingular element of the algebra  $\text{gl}(N, C)$  in its defining representation.

Equation (2.1) can be rewritten as a set of first-order PDE's:

$$B_{\bar{y}}^1 + B_{\bar{z}}^2 = 0, \quad B^1 = J^{-1} J_y, \quad B^2 = J^{-1} J_z, \quad (2.3)$$

where a standard notation for partial derivatives has been used. We are thus led, in the spirit of Ref. 1, to define the following set of four-forms in seven variables:

$$\begin{aligned} \gamma_1 &= dy dz dB^1 d\bar{z} + dy dz d\bar{y} dB^2, \\ \gamma_2 &= dJ dz d\bar{y} d\bar{z} - JB^1 dy dz d\bar{y} d\bar{z}, \\ \gamma_3 &= dy dJ d\bar{y} d\bar{z} - JB^2 dy dz d\bar{y} d\bar{z}. \end{aligned} \quad (2.4)$$

It is easily seen that the  $d\gamma_k$  are in the ideal of the  $\gamma_k$ ; thus this ideal is closed.

We now proceed to find the isovectors of the system. For this purpose we must expand the Lie derivative of each  $\gamma_i$  into a “linear” combination of all three  $\gamma_k$ . The expansion must be made consistently with the requirement that the Lie derivative preserve the tensorial character of each  $\gamma_i$  separately.

Now, from Eqs. (2.4) it can be seen that the four-forms  $\gamma_k$  have values in  $\text{gl}(N, C)$ , which is closed under both addition and multiplication. This observation suggests the following expansion:

$$\xi\gamma_i = b_i^k \gamma_k + \Lambda_i^k \gamma_k + \gamma_k M_i^k, \quad (2.5)$$

where the  $b_i^k$  are scalars, whereas the zero-forms  $\Lambda_i^k$  and  $M_i^k$  have values in  $\text{gl}(N, C)$ .

The vector field  $V$  is defined on a jetlike space with "coordinates"  $y, z, \bar{y}, \bar{z}, J, B^1$ , and  $B^2$ . As argued in Ref. 1,  $V$  will have a formal representation,

$$V = D^1 \frac{\partial}{\partial y} + D^2 \frac{\partial}{\partial z} + D^3 \frac{\partial}{\partial \bar{y}} + D^4 \frac{\partial}{\partial \bar{z}} + G \frac{\partial}{\partial J} + A^1 \frac{\partial}{\partial B^1} + A^2 \frac{\partial}{\partial B^2}, \quad (2.6)$$

where the  $D^1, \dots, D^4$  are assumed to be scalar functions of  $y, z, \bar{y}, \bar{z}$ , while the  $G, A^1, A^2$  are  $\text{gl}(N, C)$ -valued functions of the above four variables and  $J, B^1$ , and  $B^2$ . As in Ref. 1, we seek vector fields  $V$  for which the coefficients of expansion in Eq. (2.5) depend only on  $y, z, \bar{y}$ , and  $\bar{z}$ .

Substituting Eqs. (2.4) and (2.6) into Eq. (2.5), and using Eq. (1.1) to write

$$\begin{aligned} \xi dJ &= dG = G_{,\mu} dy^\mu + \bar{d}G, \\ \xi dB^k &= dA^k = A^k_{,\mu} dy^\mu + \bar{d}A^k \quad (k=1,2), \end{aligned}$$

where the  $y^\mu$  ( $\mu=1, \dots, 4$ ) denote the  $y, \dots, \bar{z}$ , we obtain a set of three exterior equations for four-forms. By equating the coefficients of  $dy dz d\bar{y} d\bar{z}$  on both sides of each exterior equation, the following set of PDE's is derived:

$$\begin{aligned} A^1_{\bar{y}} + A^2_{\bar{z}} &= -(b_1^2 + \Lambda_1^2)JB^1 - (b_1^3 + \Lambda_1^3)JB^2 \\ &\quad - JB^1 M_1^2 - JB^2 M_1^3, \\ G_y - GB^1 - JA^1 - D^1_{,\mu} JB^1 &= -(b_2^2 + \Lambda_2^2)JB^1 - (b_2^3 + \Lambda_2^3)JB^2 \\ &\quad - JB^1 M_2^2 - JB^2 M_2^3, \\ G_z - GB^2 - JA^2 - D^2_{,\mu} JB^2 &= -(b_3^2 + \Lambda_3^2)JB^1 - (b_3^3 + \Lambda_3^3)JB^2 \\ &\quad - JB^1 M_3^2 - JB^2 M_3^3, \end{aligned} \quad (2.7)$$

where  $D^\mu \equiv D^1, \dots, D^4$ .

We now put

$$\begin{aligned} A^i &= \alpha^{ik}(y^\mu)B^k + \beta^i(y^\mu)J + \bar{A}^i(y^\mu, B^k, J), \\ G &= \delta^k(y^\mu)B^k + \epsilon(y^\mu)J + \bar{G}(y^\mu, B^k, J), \end{aligned} \quad (2.8)$$

where the  $\alpha^{ik}$ ,  $\beta^i$ ,  $\delta^k$ , and  $\epsilon$  are scalars. Then

$$\begin{aligned} \bar{d}A^i &= \alpha^{ik} dB^k + \beta^i dJ + \bar{d}\bar{A}^i, \\ \bar{d}G &= \delta^k dB^k + \epsilon dJ + \bar{d}\bar{G}. \end{aligned}$$

We substitute these expressions into the expansion of Eq. (2.5) and equate coefficients of terms that are *scalar* multiples of similar  $\text{gl}(N, C)$ -valued basis four-forms. There are 12 such basis four-forms; therefore we obtain a set of 36 equations (eight of which are trivial identities). These results can be summarized as follows:

$$\begin{aligned} \beta^1 &= \beta^2 = \delta^1 = \delta^2 = 0; \quad \alpha^{12} = D^3_{\bar{z}}, \quad \alpha^{21} = D^4_{\bar{y}}; \\ D^1_{\bar{y}} &= D^1_{\bar{z}} = 0, \quad D^2_{\bar{y}} = D^2_{\bar{z}} = 0, \\ D^3_y &= D^3_z = 0, \quad D^4_y = D^4_z = 0; \\ b^1_1 &= D^1_y + D^2_z + D^3_{\bar{y}} + \alpha^{22} = D^1_y + D^2_z + D^4_{\bar{z}} + \alpha^{11}, \end{aligned}$$

$$\begin{aligned} b^2_2 &= D^2_z + D^3_{\bar{y}} + D^4_{\bar{z}} + \epsilon, \quad b^3_3 = D^1_y + D^3_{\bar{y}} + D^4_{\bar{z}} + \epsilon, \\ b^2_2 &= -D^2_y, \quad b^2_3 = -D^1_z, \quad b^1_1 = b^3_1 = b^2_2 = b^1_3 = 0. \end{aligned}$$

We notice, in particular, that the  $D^1$  and  $D^2$  depend only on  $y$  and  $z$ , while the  $D^3$  and  $D^4$  depend only on  $\bar{y}$  and  $\bar{z}$ .

The remaining terms in the expansion of Eq. (2.5) are those that cannot be expressed as *scalar* multiples of basis four-forms [in the sense that the coefficients in these terms do not commute with the  $\text{gl}(N, C)$ -valued basis four-forms]. Terms of this type can be divided into four kinds according to their dependence on the basis three-forms  $dy dz d\bar{y}$ ,  $dy dz d\bar{z}$ ,  $dy d\bar{y} d\bar{z}$ , or  $dz d\bar{y} d\bar{z}$ . The  $\text{gl}(N, C)$ -valued coefficients of each of these basis three-forms must be equated in each of the three exterior equations; this process yields a set of 12 equations which can be divided into two general types:

$$\Lambda_i^k dY + (dY)M_i^k = 0, \quad i \neq k, \quad (2.9)$$

and

$$\bar{d}\bar{H} = \Lambda_k^k dY + (dY)M_k^k, \quad (2.10)$$

where  $Y \equiv B^1, B^2, J$  and  $H \equiv A^1, A^2, G$ . The variable  $Y$ , by assumption, does not commute with  $\Lambda_i^k$  and  $M_i^k$ . Thus Eq. (2.9) is satisfied only if  $\Lambda_i^k = M_i^k = 0, i \neq k$ . Also, given that, by definition of the internal exterior derivative and by assumption about the  $\Lambda_i^k$  and  $M_i^k$ ,  $dY = \bar{d}Y$ ,  $\bar{d}\Lambda_k^k(y^\mu) = 0$ ,  $\bar{d}M_k^k(y^\mu) = 0$ , Eq. (2.10) can be integrated immediately:

$$\bar{H} = \Lambda_k^k Y + YM_k^k + h(y^\mu),$$

where  $h(y^\mu)$  is an arbitrary function. Our results are explicitly stated as follows:

$$\begin{aligned} \Lambda_1^1 &\equiv \Lambda^1(y^\mu), \quad M_1^1 \equiv M^1(y^\mu), \\ \Lambda_2^2 &\equiv \Lambda^2(y^\mu), \quad M_2^2 \equiv M^2(y^\mu), \\ \Lambda_i^k &= M_i^k = 0, \quad \text{for } i \neq k; \\ \bar{A}^1 &= \Lambda^1 B^1 + B^1 M^1 + h^1(y^\mu), \\ \bar{A}^2 &= \Lambda^2 B^2 + B^2 M^2 + h^2(y^\mu), \\ \bar{G} &= \Lambda^2 J + JM^2 + g(y^\mu), \end{aligned}$$

where the  $h^1, h^2$ , and  $g$  are arbitrary  $\text{gl}(N, C)$ -valued functions.

Appropriate substitutions into Eqs. (2.8) will now give expressions for  $A^i$  and  $G$ , which can be substituted back into Eqs. (2.7). By using previous results, the coefficients  $b_i^k$  can be eliminated in favor of other quantities, while certain replacements can also be made with regard to the  $\Lambda_i^k$  and  $M_i^k$ . The result is a set of equalities between some kind of generalized "polynomial" expressions in the variables  $B^1, B^2$ , and  $J$ , with  $y^\mu$ -dependent coefficients. The "constant" term in such a "polynomial" is a matrix function  $F(y^\mu)$ , while the other terms are of the following kinds:  $qB^k, qJ, qJB^k, QB^k, B^k Q, QJ, JQ, QJB^k, JQB^k$ , and  $JB^k Q$ , where  $q(y^\mu)$  is a scalar function and  $Q(y^\mu)$  is a  $\text{gl}(N, C)$ -valued function. Equating coefficients of similar terms we obtain a set of partial differential and algebraic equations, which are not hard to solve. In particular, we find

$$\begin{aligned} -\Lambda^1 &= M^1 = M^2 \equiv M(y, z), \quad \Lambda^2 \equiv \Lambda(\bar{y}, \bar{z}), \\ h^1(y, z) &= M_y, \quad h^2(y, z) = M_z, \quad g(y^\mu) = 0. \end{aligned}$$

Equations (2.11) give the complete solution for the components of the isovector field  $V$ :

$$\begin{aligned} D^1 &= c_1 y + k_1 z + \alpha_1, & D^2 &= k_2 y + c_2 z + \alpha_2, \\ D^3 &= (c_2 - c) \bar{y} - k_2 \bar{z} + \alpha_3, \\ D^4 &= -k_1 \bar{y} + (c_1 - c) \bar{z} + \alpha_4, & (2.11) \\ A^1 &= -c_1 B^1 - k_2 B^2 - [M(y,z), B^1] + M_y, \\ A^2 &= -k_1 B^1 - c_2 B^2 - [M(y,z), B^2] + M_z, \\ G &= \epsilon(\bar{y}, \bar{z}) J + \Lambda(\bar{y}, \bar{z}) J + JM(y,z), \end{aligned}$$

where  $c_1, c_2, k_1, k_2, c, \alpha_1, \dots, \alpha_4$  are nine complex parameters,  $\epsilon(\bar{y}, \bar{z})$  is a scalar function, and  $M(y,z)$  and  $\Lambda(\bar{y}, \bar{z})$  are  $\text{gl}(N, \mathbb{C})$ -valued functions. From Eqs. (2.11) we can read off the infinitesimal operators  $P_k$  corresponding to the nine complex parameters (cf. Ref. 1) and we can show that they form the basis of a Lie algebra. In particular, the operators

$$P_{\alpha_\mu} = \frac{\partial}{\partial y^\mu}$$

and

$$P_{c_i} + P_{\alpha_i} = y^\mu \frac{\partial}{\partial y^\mu} - B^k \frac{\partial}{\partial B^k}$$

represent translations and dilatations, respectively.

Following the discussion in Ref. 1, from Eq. (2.11) we can construct the following infinitesimal internal symmetry transformations:

$$\begin{aligned} B^{1'} &\simeq B^1 + [M(y,z), B^1] - M_y, \\ B^{2'} &\simeq B^2 + [M(y,z), B^2] - M_z, & (2.12) \\ J' &\simeq J - \epsilon(\bar{y}, \bar{z}) J - \Lambda(\bar{y}, \bar{z}) J - JM(y,z), \end{aligned}$$

where the  $\epsilon, M$ , and  $\Lambda$  are infinitesimal. The corresponding finite transformations are

$$\begin{aligned} B^{1'} &= UB^1 U^{-1} + U \partial_y U^{-1}, \\ B^{2'} &= UB^2 U^{-1} + U \partial_z U^{-1}, & (2.13) \\ J' &= \beta \bar{U} J U, \end{aligned}$$

where

$$\begin{aligned} U(y,z) &= \exp\{-M(y,z)\}, \\ \bar{U}(\bar{y}, \bar{z}) &= \exp\{-\Lambda(\bar{y}, \bar{z})\}, \end{aligned}$$

and

$$\beta(\bar{y}, \bar{z}) = \exp\{-\epsilon(\bar{y}, \bar{z})\}.$$

These are, of course, familiar symmetries of the SDYM system.

### III. PARAMETRIC INFINITESIMAL TRANSFORMATION AND LINEAR SYSTEM

If we define a new function

$$\xi(y,z, \bar{y}, \bar{z}) \equiv M(y,z) + \epsilon(\bar{y}, \bar{z}) 1_N, \quad (3.1)$$

where  $1_N$  denotes the  $N$ -dimensional unit matrix, then the infinitesimal transformations of Eq. (2.12) with  $\Lambda(\bar{y}, \bar{z}) = 0$  can be rewritten as

$$\begin{aligned} \delta B^1 &= [\xi(y^\mu), B^1] - \xi_y, \\ \delta B^2 &= [\xi(y^\mu), B^2] - \xi_z, & \delta J &= -J\xi(y^\mu), \end{aligned} \quad (3.2)$$

where  $\delta B^k \simeq B^{k'} - B^k$  and  $\delta J \simeq J' - J$ . We wish to rewrite

these symmetries without the restriction (3.1). It turns out that this is possible due to the independence of the coordinates of the underlying jetlike space. Of course, there is a price to be paid for such an adjustment. But this "price" is a most welcome one: Restriction (3.1) is replaced by a set of linear PDE's which, in the case of actual SDYM fields, lead to a linear system for the SDYM equation.

From Eq. (3.1) it is seen that  $\xi(y^\mu)$  satisfies the PDE,

$$[\xi_{\bar{y}}, B^1] + [\xi_{\bar{z}}, B^2] - \xi_{y\bar{y}} - \xi_{z\bar{z}} = 0.$$

Given the independence of the  $y^\mu$  and the  $B^k$  (this is the case as long as no restriction on the solution manifold is imposed), the above equation may be written as

$$\partial_{\bar{y}}([\xi, B^1] - \xi_y) + \partial_{\bar{z}}([\xi, B^2] - \xi_z) = 0. \quad (3.3)$$

This is satisfied if there exists a "potential"  $\psi(y^\mu, B^k)$  such that

$$[\xi, B^1] - \xi_y = \lambda \psi_z, \quad [\xi, B^2] - \xi_z = -\lambda \psi_{\bar{y}}, \quad (3.4)$$

where  $\lambda$  is an arbitrary complex parameter. We thus replace the system of Eqs. (3.1) and (3.2) by the following alternate one:

$$\delta B^1 = \lambda \psi_z, \quad \delta B^2 = -\lambda \psi_{\bar{y}}, \quad \delta J = -J\xi(y^\mu), \quad (3.5)$$

where  $\psi$  and  $\xi$  satisfy the linear system (3.4). Note that Eqs. (3.4) and (3.5) become independent of Eqs. (3.1) and (3.2) upon restriction to the solution manifold, i.e., for actual SDYM fields.

Let us explore further the significance of Eqs. (3.4) for actual SDYM fields (in which case the  $B^k$  are dependent upon the  $y^\mu$ ). In particular, let us examine the ansatz  $\psi(y^\mu) = \xi(y^\mu)$ , all  $y^\mu$ :

$$[\xi, B^1] - \xi_y = \lambda \xi_z, \quad [\xi, B^2] - \xi_z = -\lambda \xi_{\bar{y}}. \quad (3.6)$$

The integrability criterion  $\xi_{y\bar{z}} - \xi_{z\bar{y}} = 0$  yields Eq. (3.3), which, in combination with Eq. (3.6), gives

$$[\xi, B_y^2 - B_z^1 + [B^1, B^2] + \lambda(B_y^1 + B_z^2)] = 0.$$

We seek conditions for  $B^1$  and  $B^2$  in order that the above equality holds for all  $\lambda$  and independently of  $\xi$ . The following pair of PDE's must therefore be satisfied:

$$\partial_y B^2 - \partial_z B^1 + [B^1, B^2] = 0, \quad (3.7)$$

$$\partial_{\bar{y}} B^1 + \partial_{\bar{z}} B^2 = 0. \quad (3.8)$$

Equation (3.7) is a condition for zero curvature and implies that the  $B^1$  and  $B^2$  are pure gauges:

$$B^1 = J^{-1} \partial_y J, \quad B^2 = J^{-1} \partial_z J, \quad (3.9)$$

where  $J$  is a nonsingular  $\text{gl}(N, \mathbb{C})$  matrix. Then Eq. (3.8) becomes identical to the SDYM equation (2.1), of which Eq. (3.6) is seen to be a linear system.

We remark that our results are in agreement with those of Ref. 5 (although they are given in a slightly different form). The thing to notice is that these results were actually *derived* here, in a rather straightforward manner, by using the isovector technique.

### IV. CONNECTION WITH BÄCKLUND TRANSFORMATIONS

By using the original definitions of  $B^1$  and  $B^2$  as given in Eqs. (2.3), the infinitesimal transformations of these quanti-

ties may be written, according to Eqs. (3.5) as

$$J'^{-1}J'_y - J^{-1}J_y = \lambda\psi_z, \quad (4.1a)$$

$$J'^{-1}J'_z - J^{-1}J_z = -\lambda\psi_{\bar{y}}. \quad (4.1b)$$

Clearly, as  $J'$  approaches  $J$ , the  $\psi_{\bar{y}}$  and  $\psi_z$  must approach zero. One way to achieve this is to put

$$\psi = \xi = 1 - J^{-1}J'. \quad (4.2)$$

Now, if the left-hand sides of Eqs. (4.1a) and (4.1b) are considered as *finite*, rather than infinitesimal differences, then Eqs. (4.1) and (4.2) constitute one possible form of the Bäcklund transformation (BT) proposed in Ref. 7. Alternatively, the infinitesimal parametric transformation (4.1) and (4.2) is also an infinitesimal BT. This was observed in Ref. 5, but we include it in the present discussion due to its direct (and quite interesting) relevance to the isovector method.

Incidentally, the transformation (3.1) and (3.2) is also an infinitesimal BT, with Eq. (3.1) being a sort of algebraic constraint. Indeed, putting  $\xi = 1 - J^{-1}J'$  and introducing an arbitrary complex parameter  $\mu$ , we write

$$J'^{-1}J'_y - J^{-1}J_y = \mu\{[J^{-1}J', J^{-1}J_y] - \partial_y(J^{-1}J')\}, \quad (4.3a)$$

$$J'^{-1}J'_z - J^{-1}J_z = \mu\{[J^{-1}J', J^{-1}J_z] - \partial_z(J^{-1}J')\}, \quad (4.3b)$$

$$J^{-1}J' = M(y, z) + \epsilon(\bar{y}, \bar{z})1_N, \quad (4.3c)$$

where  $M(y, z)$  is  $\text{gl}(N, C)$  valued and  $\epsilon(\bar{y}, \bar{z})$  is a scalar. Tak-

ing  $\partial_{\bar{y}}$ (4.3a) +  $\partial_z$ (4.3b) and using (4.3c), we find

$$\begin{aligned} & \{\partial_{\bar{y}}(J'^{-1}J'_y) + \partial_z(J'^{-1}J'_z)\} \\ & - \{\partial_{\bar{y}}(J^{-1}J_y) + \partial_z(J^{-1}J_z)\} \\ & = \mu[J^{-1}J', \partial_{\bar{y}}(J^{-1}J_y) + \partial_z(J^{-1}J_z)], \end{aligned}$$

according to which  $J'$  satisfies the SDYM equation (2.1) if  $J$  does. Note that the BT was constructed so as to yield the trivial solution  $J' = J$  as a particular solution [this corresponds to  $M = 0$  and  $\epsilon = 1$  in the algebraic constraint (4.3c)].

## ACKNOWLEDGMENT

We are indebted to Professor Yong-Shi Wu for very enlightening discussions and for providing us with useful references.

<sup>1</sup>C. J. Papachristou and B. K. Harrison, "Isogroups of differential ideals of vector-valued differential forms: Application to partial differential equations," in *Proceedings of the XV International Colloquium on Group Theoretical Methods in Physics*, edited by R. Gilmore (World Scientific, Singapore, 1987).

<sup>2</sup>B. K. Harrison and F. B. Estabrook, *J. Math. Phys.* **12**, 653 (1971).

<sup>3</sup>D. G. B. Edelen, *Applied Exterior Calculus* (Wiley, New York, 1985), Chap. 6.

<sup>4</sup>C. N. Yang, *Phys. Rev. Lett.* **38**, 1377 (1977); Y. Brihaye, D. B. Fairlie, J. Nuyts, and R. G. Yates, *J. Math. Phys.* **19**, 2528 (1978).

<sup>5</sup>L.-L. Chau, M. L. Ge, and Y. S. Wu, *Phys. Rev. D* **25**, 1086 (1982).

<sup>6</sup>L.-L. Chau and Y. S. Wu, *Phys. Rev. D* **26**, 3581 (1982); L.-L. Chau, M. L. Ge, A. Sinha, and Y. S. Wu, *Phys. Lett.* **121B**, 391 (1983).

<sup>7</sup>M. K. Prasad, A. Sinha, and L.-L. Chau Wang, *Phys. Rev. Lett.* **43**, 750 (1979).

# The invariant density for a class of discrete-time maps involving an arbitrary monotonic function operator and an integer parameter

C. C. Grosjean

Seminarie voor Wiskundige Natuurkunde, Rijksuniversiteit-Gent, Krijgslaan 281, Gebouw S9, B-9000 Gent, Belgium

(Received 29 April 1986; accepted for publication 17 December 1986)

When  $x_t$  and  $x_{t+1}$  represent two random variables, each belonging to a real interval  $[0,1]$  and being related by a first-order difference equation  $x_{t+1} = F(x_t)$ , called a discrete-time map, the probability density distribution connected with  $x_t$  can be translated into that associated with  $x_{t+1}$ . This yields an evolution equation by means of which one can construct an infinite sequence  $\{w_t(x) | t \in \mathbb{N}, x \in [0,1]\}$  starting from an integrable function  $w_0(x)$  normalized to unity on  $[0,1]$ . The question of the convergence of the sequence toward a so-called invariant density function  $w(x)$  as  $t \rightarrow +\infty$  and the problem of finding this limit were examined by a number of authors, mostly studying isolated cases. The present paper solves the problem for a class of discrete-time maps characterized by  $x_{t+1} = f(|\text{sn}[l \text{sn}^{-1} f^{-1}(x_t)]|)$ ,  $l \in \{2,3,\dots\}$ , whereby  $f$  is a real, continuous, monotonically increasing function mapping  $[0,1]$  onto itself and  $\text{sn}$  is the usual symbol for the sinelike Jacobian elliptic function with modulus  $k \in [0,1]$  (including the sine function). Convergence is proven under very general conditions on  $w_0(x)$  and an explicit formula to calculate  $w(x)$  is established. Some properties of  $w(x)$  are discussed. A necessary and sufficient condition for the symmetry of  $w(x)$  about  $x = \frac{1}{2}$  is obtained and attention is also devoted to the inverse problem, leading to a reformulation of the discrete-time map of the type cited above which corresponds to a given invariant density. The examples of practical application considered here cover almost all special cases which were treated in the literature thus far, as well as new cases.

## I. INTRODUCTION

In a number of articles,<sup>1-9</sup> most of them of recent date, the evolution of a normalized single-valued real function (usually a probability density distribution) defined on a finite real interval toward an invariant limit function under a given discrete-time mapping has been studied. Let  $x_t$  and  $x_{t+1}$  be two continuously varying random variables, each belonging, for instance, to  $[0,1]$  and related to one another by a first-order difference equation

$$x_{t+1} = F(x_t). \quad (1.1)$$

If  $w_t(x_t)$  represents the probability density distribution associated with  $x_t$ , so that  $w_t(x_t) |dx_t|$  is the elementary probability that the random variable takes a value between  $x_t$  and  $x_t + dx_t$  (whereby  $dx_t$  may be either positive or negative), one can ask for the corresponding elementary probability  $w_{t+1}(x_{t+1}) |dx_{t+1}|$  associated with the variable  $x_{t+1}$  and its differential. The expression of  $w_{t+1}$  in terms of  $w_t$  and  $F$  ultimately leads to a transformation by means of which one can construct iteratively an infinite sequence of functions  $\{w_t(x) | t \in \mathbb{N}, x \in [0,1]\}$  starting from a given initial function  $w_0(x)$ , with all functions normalized to unity when

$$\int_0^1 w_0(x) dx = 1. \quad (1.2)$$

Thus  $w_0(x)$  is assumed to stem from a single-valued real function defined on  $[0,1]$ , being continuous or maybe only piecewise continuous, but satisfying conditions such that its definite integral from 0 to 1 (possibly an improper integral) be convergent. Then, using a suitable proportionality factor, normalization to unity may be achieved, giving rise to  $w_0(x)$ . The question is, does the infinite sequence of  $w$  func-

tions tend to a limit distribution when  $t \rightarrow +\infty$  and if so, find

$$\lim_{t \rightarrow +\infty} w_t(x), \quad 0 \leq x \leq 1.$$

In Ref. 1, one finds the following particular example:

$$x_{t+1} = F(x_t) = \begin{cases} 2x_t, & 0 \leq x_t \leq \frac{1}{2}, \\ 2(1-x_t), & \frac{1}{2} \leq x_t \leq 1, \end{cases} \quad (1.3)$$

in which case the limit distribution is

$$w(x) = \lim_{t \rightarrow +\infty} w_t(x) = 1, \quad 0 \leq x \leq 1. \quad (1.4)$$

In Ref. 5, it is shown that to the cusp-shaped return map

$$x_{t+1} = 1 - 2\sqrt{|x_t|}, \quad -1 \leq x_t \leq 1, \quad (1.5)$$

by means of which the interval  $[-1,1]$  is mapped onto itself, there corresponds as a limit distribution

$$w(x) = \frac{1}{2}(1-x), \quad -1 \leq x \leq 1. \quad (1.6)$$

We shall reconsider these examples in Sec. III.

Ulam and von Neumann's paper<sup>2</sup> comprises still another example which has received considerable attention as of late. Under the discrete-time quadratic map (a special case of the so-called logistic map),

$$x_{t+1} = 4x_t(1-x_t), \quad (1.7)$$

of the interval  $[0,1]$  onto itself, any normalizable initial function  $w_0(x)$  leads to the invariant density

$$w(x) = 1/\pi [x(1-x)]^{1/2}, \quad 0 \leq x \leq 1. \quad (1.8)$$

Indeed, Falk<sup>6</sup> has calculated the general element  $w_t(x)$  of the sequence  $\{w_t(x) | t \in \mathbb{N}, x \in [0,1]\}$  generated under (1.7) starting from the uniform distribution  $w_0(x) = 1$  and obtained (1.8) by direct transition to the limit. Falk notes that

the symmetry of  $w_0(x)$  about  $x = \frac{1}{2}$  is broken at  $t = 1, 2, 3, \dots$  and returns only in the limit  $t \rightarrow +\infty$ . Nandakumaran<sup>7</sup> has extended Falk's results to all probability density distributions of the form

$$w_0(x) = [x^n(1-x)^n]/[B(n+1, n+1)], \quad n \in \mathbb{N}, \\ 0 \leq x \leq 1,$$

in particular confirming (1.8). Very recently, Grosjean<sup>8,9</sup> generalized Nandakumaran's work to a broad class of initial functions, normalized to unity but not necessarily positive-semidefinite, first under the quadratic mapping (1.7) and later under the polynomial discrete-time maps of any degree which generalize (1.7). Such mappings are first-order difference equations of the form

$$x_{t+1} = p_l(x_t), \quad l \in \{3, 4, \dots\}, \quad (1.9)$$

in which  $p_l$  represents a real polynomial of degree  $l$  and are uniquely defined by the following requirements.

(1) For every  $x_t \in [0, 1]$ , the corresponding  $x_{t+1}$  also belongs to  $[0, 1]$  where, in particular,  $x_t = 0$  entails  $x_{t+1} = 0$ .

(2) To every  $x_{t+1} \in ]0, 1[$ , there correspond  $l$  distinct real values of  $x_t$  belonging to  $]0, 1[$ .

In this way, the explicit expression of (1.9), with (1.7) included, is

$$x_{t+1} = x_t U_{l-1}^2(\sqrt{1-x_t}), \quad l \in \{2, 3, \dots\}, \quad (1.10)$$

where  $U$  symbolizes a Chebyshev polynomial of the second kind, or in parametric form,

$$x_t = (\sin \theta)^2, \quad x_{t+1} = (\sin l\theta)^2, \quad 0 \leq \theta \leq \pi/2, \\ l \in \{2, 3, \dots\}. \quad (1.11)$$

Thinking of initial functions which can possibly be only piecewise continuous by the appearance of a number of finite jumps in  $[0, 1]$ , Grosjean<sup>8,9</sup> proposed to represent  $w_0(x)$  by a convergent series of the type

$$w_0(x) = \frac{a_0}{2} + \sum_{n=1}^{+\infty} [a_n T_{2n}(1-2x) \\ + 2b_n x^{1/2}(1-x)^{1/2} U_{2n-1}(1-2x)], \quad (1.12)$$

deduced from the Fourier series

$$\frac{a_0}{2} + \sum_{n=1}^{+\infty} (a_n \cos 4n\theta + b_n \sin 4n\theta), \quad 0 \leq \theta \leq \pi/2, \quad (1.13)$$

by means of the substitution  $x = (\sin \theta)^2$  and normalized to unity [see (1.2)] when

$$\frac{a_0}{2} - \sum_{n=1}^{+\infty} \frac{a_n}{4n^2 - 1} = 1. \quad (1.14)$$

Grosjean<sup>8,9</sup> established the following (provisional) theorem.

If in  $[0, 1]$  the initial function  $w_0(x)$  is the sum of a convergent series of the type (1.12) whereby the  $a$  coefficients satisfy (1.14) as well as the condition

$$\frac{|a_0|}{2} + \sum_{n=1}^{+\infty} |a_n|: \text{convergent}, \quad (1.15)$$

then under any polynomial discrete-time map comprised in (1.10), the sequence  $\{w_t(x) | t \in \mathbb{N}, x \in [0, 1]\}$  converges to-

ward the limit function

$$w(x) = 1/\pi [x(1-x)]^{1/2}, \quad 0 \leq x \leq 1.$$

As in Refs. 6 and 7, the method used consisted of calculating  $w_t(x)$  explicitly for any  $t \in \mathbb{N}_0$  and letting  $t$  approach infinity. But in the final step of the proof, in the case that the "a part" in (1.12) involves infinitely many terms, there is need for uniform convergence with respect to  $t$  and thus (1.15) is required as a *sufficient* condition. Unfortunately, when (1.15) holds, it deprives the part of  $w_0(x)$  which is symmetric about  $x = \frac{1}{2}$  of the possibility of exhibiting finite jumps between 0 and 1, on account of the criterion of Weierstrass for the absolute and uniform convergence of the "a part" in a series expansion such as (1.13) and *a fortiori* (1.12). In a subsequent note, Barbour<sup>10</sup> succeeded in eliminating this drawback.

The purpose of the present paper is twofold: We wish to make the conditions imposed upon the initial function  $w_0(x)$  much less restrictive and extend the theory to a much wider class of discrete-time maps [rational and irrational functions  $F$  in (1.1)]. The proof of the convergence of the sequence  $\{w_t(x) | t \in \mathbb{N}, x \in [0, 1]\}$  toward a limit density  $w(x)$  as  $t \rightarrow +\infty$  and finding this density will be based upon a procedure differing entirely from the methods applied in the various articles cited in this Introduction, including my own articles.<sup>8,9</sup>

## II. THEORETICAL DEVELOPMENT

The parametric transformation  $x_t = (\sin \theta)^2$  applied to the logistic map (1.7) is implicit in most studies of the map. As was shown in Sec. I, it has led me to the polynomial generalization of arbitrary degree (1.10) via (1.11). As a next step, inspired by the parametric representation (1.11), one can consider the much broader generalization

$$x_t = f(\sin \theta), \quad x_{t+1} = f(|\sin l\theta|), \quad 0 \leq \theta \leq \pi/2, \\ l \in \{2, 3, \dots\}, \quad (2.1)$$

where  $f(y)$  is a single-valued, real, continuous function of  $y$  defined on  $[0, 1]$ , having a continuous or piecewise continuous derivative in that interval and increasing monotonically (*sensu stricto*) from 0 to 1 as  $y$  increases from 0 to 1. The inverse function, which as a function of  $y \in [0, 1]$  we denote by  $f^{-1}(y)$ , is endowed with the same properties and therefore the discrete-time map corresponding to the propounded parametric representation (2.1) is

$$x_{t+1} = f(|\sin[l \arcsin f^{-1}(x_t)]|), \quad 0 \leq x_t \leq 1, \\ l \in \{2, 3, \dots\}. \quad (2.2)$$

But, the sine function is known to be a special case comprised in the Jacobian elliptic sn function. Although one encounters the Jacobian elliptic functions much less frequently than the circular functions in practice, another degree of generality is gained when (2.1) is replaced by

$$x_t = f(\text{sn}(u, k)), \quad x_{t+1} = f(|\text{sn}(lu, k)|), \quad 0 \leq u \leq K(k), \\ l \in \{2, 3, \dots\}, \quad (2.3)$$

where the modulus  $k \in [0, 1[$  and  $K(k)$  is the complete elliptic



integral of the first kind, i.e.,

$$K(k) = \int_0^1 \frac{dx}{(1-x^2)^{1/2}(1-k^2x^2)^{1/2}},$$

whose value is real, finite, and greater than  $\pi/2$  for  $k \in ]0,1[$ . In the special case  $k = 0$ , the Jacobian elliptic functions reduce to circular functions. From here onward, the modulus will not be written explicitly as an argument of the sn function and the complete elliptic integral except where it is desirable to include it in the notation. The counterpart of (2.2) is

$$x_{t+1} = f(|\text{sn}[l \text{sn}^{-1} f^{-1}(x_t)]|), \quad 0 \leq x_t \leq 1, \quad l \in \{2,3,\dots\}, \quad (2.4)$$

with  $\text{sn}^{-1}$  the usual symbol for the inverse of the Jacobian sn function, defined on  $[0,1]$  by

$$\text{sn}^{-1} v = \int_0^v \frac{dx}{(1-x^2)^{1/2}(1-k^2x^2)^{1/2}}, \quad 0 \leq v \leq 1. \quad (2.5)$$

Here, (2.4) is the definitive form of the discrete-time mappings upon which the remaining calculations in the present article will be based. It is sufficiently general to include a variety of interesting special cases, but it is not yet the most general form which can be considered.

Note that the formulas (2.2) and (2.4), each one connecting  $x_t$  and  $x_{t+1}$ , are closely related to the relevant concept of conjugate functions, i.e.,

$$b(x) = h(a[h^{-1}(x)]),$$

which has played an important role in a number of articles. One finds it, for instance, in articles by Halmos<sup>11</sup> and by Ulam,<sup>12</sup> and it was discussed by Grossmann and Thomae as well.<sup>13</sup> Also, Györgyi and Szépfalussy<sup>14</sup> made extensive use of the concept of conjugation to determine invariant measures for various one-dimensional maps. In the important special case of  $l = 2$ , the absolute value stripes may be deleted in (2.2) and (2.4) under the assumptions made for the function  $f(y)$  in (2.1). If a conjugating function  $h(y)$  is such that  $h[f(y)]$  is still a single-valued, real, continuous function of  $y$  on  $[0,1]$ , having a continuous or piecewise continuous derivative in that interval and increasing monotonically (*sensu stricto*) from 0 to 1 as  $y$  increases from 0 to 1, then one can associate with (2.2), where one puts  $l = 2$ , the discrete-time map

$$x_{t+1} = hf(\sin[2 \arcsin f^{-1} h^{-1}(x_t)]); \quad (2.2')$$

the analysis which will follow is also applicable to this map, with the function operator  $f$  simply replaced by  $hf$ . As is well

known, the invariant densities corresponding to the two conjugate dynamical laws (2.2) with  $l = 2$  and (2.2') are closely related, which will be confirmed by the final result for  $w(x)$ . The same remark holds for the two conjugate dynamical laws (2.4) with  $l = 2$  and

$$x_{t+1} = hf(\text{sn}[2 \text{sn}^{-1} f^{-1} h^{-1}(x_t)]). \quad (2.4')$$

To every  $x_t \in ]0,1[$ , there corresponds one real value of  $x_{t+1}$  and it also belongs to  $[0,1]$ . But, conversely, to every  $x_{t+1} \in ]0,1[$ , there correspond  $l$  distinct real values of  $x_t$  in the interval  $[0,1]$ , i.e.,

$$\begin{aligned} x_t^{(1)} &= r_1(x_{t+1}) \\ &= f[\text{sn}((1/l)\text{sn}^{-1} f^{-1}(x_{t+1}))], \\ x_t^{(2)} &= r_2(x_{t+1}) \\ &= f[\text{sn}((2K/l) - (1/l)\text{sn}^{-1} f^{-1}(x_{t+1}))], \\ x_t^{(3)} &= r_3(x_{t+1}) \\ &= f[\text{sn}((2K/l) + (1/l)\text{sn}^{-1} f^{-1}(x_{t+1}))], \\ &\vdots \\ x_t^{(l)} &= r_l(x_{t+1}) = f[\text{sn}((2[l/2]K/l) \\ &\quad - (-1)^l(1/l)\text{sn}^{-1} f^{-1}(x_{t+1}))], \end{aligned} \quad (2.6)$$

with  $[l/2]$  the largest integer smaller than or equal to  $l/2$ . These expressions remain valid as  $x_{t+1} \rightarrow 0$  and  $x_{t+1} \rightarrow 1$  from inside  $[0,1]$ . From them, one easily deduces that  $x_{t+1} = 0$  gives rise to  $x_t^{(1)} = 0$ ,  $x_t^{(2)} = x_t^{(3)}$ ,  $x_t^{(4)} = x_t^{(5)}$ , etc., whereas  $x_{t+1} = 1$  yields  $x_t^{(1)} = x_t^{(2)}$ ,  $x_t^{(3)} = x_t^{(4)}$ , and so on. In these cases, the number of distinct real  $x_t$  values is reduced to either  $[l/2]$  or  $[l/2] + 1$ . As stated in the Introduction, the transformation leading from  $w_t(x)$  to  $w_{t+1}(x)$  stems from probability theory. If  $w_t(x_t) |dx_t|$  represents the elementary probability that the random variable at the integer time instant  $t$  take on a value between  $x_t$  and  $x_t + dx_t$  (where  $dx_t$  can be either positive or negative), then we have, for the analogous elementary probability regarding the random variable at the time instant  $t + 1$ ,

$$\begin{aligned} w_{t+1}(x_{t+1}) |dx_{t+1}| &= \sum_{m=1}^l w_t(x_t^{(m)}) |dx_t^{(m)}| \\ &= \sum_{m=1}^l w_t(r_m(x_{t+1})) |dr_m(x_{t+1})|. \end{aligned}$$

Dividing by  $|dx_{t+1}|$  and dropping the subscript in  $x_{t+1}$ , we obtain the evolution equation yielding  $w_{t+1}(x)$  in terms of  $w_t(x)$ :

$$\begin{aligned} w_{t+1}(x) &= \frac{df^{-1}(x)/dx}{l[1-(f^{-1}(x))^2]^{1/2}[1-k^2(f^{-1}(x))^2]^{1/2}} \left\{ \text{cn} \frac{v}{l} \text{dn} \frac{v}{l} f' \left( \text{sn} \frac{v}{l} \right) w_t \left[ f \left( \text{sn} \frac{v}{l} \right) \right] \right. \\ &\quad + \text{cn} \frac{2K-v}{l} \text{dn} \frac{2K-v}{l} f' \left( \text{sn} \frac{2K-v}{l} \right) w_t \left[ f \left( \text{sn} \frac{2K-v}{l} \right) \right] \\ &\quad + \text{cn} \frac{2K+v}{l} \text{dn} \frac{2K+v}{l} f' \left( \text{sn} \frac{2K+v}{l} \right) w_t \left[ f \left( \text{sn} \frac{2K+v}{l} \right) \right] + \dots \\ &\quad \left. + \text{cn} \frac{2l'K - (-1)^l v}{l} \text{dn} \frac{2l'K - (-1)^l v}{l} f' \left( \text{sn} \frac{2l'K - (-1)^l v}{l} \right) w_t \left[ f \left( \text{sn} \frac{2l'K - (-1)^l v}{l} \right) \right] \right\}, \quad \forall t \in \mathbb{N}, \end{aligned} \quad (2.7)$$

in which  $v = \text{sn}^{-1} f^{-1}(x)$ ,  $l' = [l/2]$ , and  $\text{cn}, \text{dn}$  are the two well-known cosinelike Jacobian elliptic functions associated with the  $\text{sn}$  function. This being by definition the evolution equation connecting any two neighboring functions in the sequence  $\{w_t(x) | t \in \mathbb{N}, x \in [0, 1]\}$ , it is not required after all that the  $w$  functions on which (2.7) is applied be positive semidefinite on  $[0, 1]$ .

When (2.7) is used to transform a given initial function  $w_0(x)$  into  $w_1(x)$ , we obtain

$$w_1(x) = \frac{df^{-1}(x)/dx}{l[1 - (f^{-1}(x))^2]^{1/2}[1 - k^2(f^{-1}(x))^2]^{1/2}} \times \sum_u \text{cn } u \text{ dn } u f'(sn u) w_0[f(sn u)], \quad (2.8)$$

in which the summation with respect to  $u$  runs over the values

$$\frac{v}{l}, \frac{2K-v}{l}, \frac{2K+v}{l}, \dots, \frac{2l'K - (-1)^l v}{l}, \quad (2.9)$$

where  $v = \text{sn}^{-1} f^{-1}(x)$ . These are  $l$  distinct real values when  $x \in ]0, 1[$ , one in each of the open real intervals

$$]m(K/l), (m+1)(K/l)[, \quad \forall m \in \{0, 1, \dots, l-1\}.$$

When  $x \rightarrow 0$  or  $x \rightarrow 1$  from inside  $[0, 1]$ , adjacent pairs of values in (2.9) converge to each other since  $v \rightarrow 0$  or  $v \rightarrow K$ , respectively. The peculiarity of the evolution equation (2.7) is that when it is applied a second time in order to express  $w_2(x)$  in terms of  $w_0(x)$ , one obtains

$$w_2(x) = \frac{df^{-1}(x)/dx}{l^2[1 - (f^{-1}(x))^2]^{1/2}[1 - k^2(f^{-1}(x))^2]^{1/2}} \times \sum_u \text{cn } u \text{ dn } u f'(sn u) w_0[f(sn u)], \quad (2.10)$$

where the summation with respect to  $u$  now runs over the values

$$\frac{v}{l^2}, \frac{2K-v}{l^2}, \frac{2K+v}{l^2}, \dots, \frac{2l'K - (-1)^l v}{l^2}, \\ \frac{2(l-l')K + (-1)^l v}{l^2}, \dots, \frac{2l^2 - 1}{2l^2} K + (-1)^l \\ \times \left( \frac{K}{2l^2} - \frac{v}{l^2} \right), \quad \text{with } v = \text{sn}^{-1} f^{-1}(x), \quad (2.11)$$

which result from the replacement of  $v$  successively in each member of the sequence (2.9) by the entire set of  $l$  values contained in (2.9). When  $x \in ]0, 1[$ , this procedure yields  $l^2$  distinct real values, one belonging to each of the open real intervals

$$]m(K/l^2), (m+1)(K/l^2)[, \quad \forall m \in \{0, 1, \dots, l^2 - 1\}.$$

The complete verification of the results (2.10) and (2.11) would require a considerable amount of space; therefore, it is left to the reader. Only as an illustrative example of how the combinations of arguments simplify, let us calculate the argument with which the function  $w_0$  appears in the first term of the sum in (2.10). To do this, (2.7) is applied for  $t = 1$ ; hence we find in the first term, among other factors,

$$w_1[f(\text{sn}(\text{sn}^{-1} f^{-1}(x)/l))],$$

and according to (2.8), the argument of this  $w_1$  function

should replace  $x$  in

$$w_0[f(\text{sn}(\text{sn}^{-1} f^{-1}(x)/l))].$$

This yields

$$w_0\{f[\text{sn}(1/l)(\text{sn}^{-1} f^{-1}[f(\text{sn}(\text{sn}^{-1} f^{-1}(x)/l))])]\},$$

which clearly reduces to

$$w_0[f(\text{sn}(\text{sn}^{-1} f^{-1}(x)/l^2))].$$

In the same manner, after  $t$  applications of Eq. (2.7), it follows that

$$w_t(x) = \frac{df^{-1}(x)/dx}{l^t[1 - (f^{-1}(x))^2]^{1/2}[1 - k^2(f^{-1}(x))^2]^{1/2}} \times \sum_u \text{cn } u \text{ dn } u f'(sn u) w_0[f(sn u)], \quad (2.12)$$

in which, when  $x \in ]0, 1[$ , the summation with respect to  $u$  runs over  $l^t$  distinct real values resulting from repeated application of the rule described above:

$$\frac{v}{l^t}, \frac{2K-v}{l^t}, \frac{2K+v}{l^t}, \dots, \frac{2l^t - 1}{2l^t} K \\ + (-1)^l \left( \frac{K}{2l^t} - \frac{v}{l^t} \right),$$

still with  $v = \text{sn}^{-1} f^{-1}(x)$ . When the interval  $[0, K]$  is subdivided into  $l^t$  equal subintervals, one finds inside each of these subintervals one  $u$  value taking part in the summation appearing in (2.12). Consequently, when  $x \in ]0, 1[$ , the expression

$$\sum_u \text{cn}(u) \text{dn}(u) f'(sn u) w_0[f(sn u)] \times \frac{K}{l^t} \quad (2.13)$$

recalls the way in which a Riemann definite integral is defined. Indeed, the real interval  $[0, K]$  over which  $u$  can vary is subdivided into  $l^t$  equal subintervals playing the role of  $\Delta u$ , and inside each of these intervals is located one abscissa at which the value of the integrand is taken. With regard to  $w_0(x)$ , the simplest possible case is that this normalized initial function is continuous in  $[0, 1]$ . Then, it is immediately clear that for  $x \in ]0, 1[$  the sequence  $\{w_t(x) | t \in \mathbb{N}\}$  converges toward the following limit as  $t \rightarrow +\infty$ :

$$\frac{df^{-1}(x)/dx}{K[1 - (f^{-1}(x))^2]^{1/2}[1 - k^2(f^{-1}(x))^2]^{1/2}} \times \int_0^K \text{cn } u \text{ dn } u f'(sn u) w_0[f(sn u)] du. \quad (2.14)$$

But,

$$\int_0^K \text{cn } u \text{ dn } u f'(sn u) w_0[f(sn u)] du \\ = \int_0^1 f'(y) w_0[f(y)] dy \\ = \int_0^1 w_0(x) dx = 1$$

in virtue of the normalization of  $w_0(x)$ , and hence  $\{w_t(x) | t \in \mathbb{N}\}$  approaches

$$w(x) = \frac{df^{-1}(x)/dx}{K[1 - (f^{-1}(x))^2]^{1/2}[1 - k^2(f^{-1}(x))^2]^{1/2}}. \quad (2.15)$$

When  $x = 0$ , the  $u$  values over which the summation in (2.12) runs are zero and all integer multiples of  $2K/l'$  smaller than or equal to  $K$ , because, as one lets  $x \rightarrow +0$ , all or almost all pairs of adjacent  $u$  values from the second  $u$  value onward tend to a common limit. When  $l$  is even, one can subdivide  $[0, K]$  as follows:

- $[0, K/l']$  associated with  $u = 0$ ,
- $[(2m - 1)(K/l'), (2m + 1)(K/l')]$   
associated with  $u = 2mK/l'$ ,
- $\forall m \in \{1, 2, \dots, (l' - 2)/2\}$ ,
- $[(1 - 1/l')K, K]$  associated with  $u = K$ .

Here, every  $u$  value between 0 and  $K$  is located in the middle of a segment of length  $2K/l'$ . The sum appearing in (2.12) becomes

$$2 \left\{ \frac{1}{2} f'(0) w_0(0) + \sum_{m=1}^{(l'-2)/2} \operatorname{cn} \frac{2mK}{l'} \operatorname{dn} \frac{2mK}{l'} \times f' \left( \operatorname{sn} \frac{2mK}{l'} \right) w_0 \left[ f \left( \operatorname{sn} \frac{2mK}{l'} \right) \right] + \frac{1}{2} (1 - k^2)^{1/2} [\operatorname{cn} u f'(u)]_{u=K} w_0(1) \right\},$$

and with  $w_0(x)$  still assumed to be continuous in  $[0, 1]$ , this expression, when multiplied by  $K/l'$ , again converges toward the Riemann definite integral appearing in (2.14) when  $t \rightarrow +\infty$ . Hence, when  $l$  is even, (2.15) remains valid in the limit  $x = 0$ . A similar reasoning can be formulated when  $l$  is odd, and after that also for  $x = 1$ . Therefore, (2.15) holds good for all  $x \in [0, 1]$  under the assumption that the initial function  $w_0(x)$  is continuous in the interval  $[0, 1]$ .

Our argument can easily be extended to other cases, where the given initial function  $w_0(x)$  involves discontinuities in  $[0, 1]$ . When  $w_0(x)$  is piecewise continuous, having a limited number of finite jumps in  $[0, 1]$ , there corresponds to every abscissa where such a jump occurs a  $u$  value belonging to  $[0, K]$ , obtained by the transformation  $u = \operatorname{sn}^{-1} f^{-1}(x)$ . These  $u$  values subdivide  $[0, K]$  into a finite number of subintervals in each of which  $w_0[f(\operatorname{sn} u)]$  is continuous, and for sufficiently large  $t$ , when  $x \in ]0, 1[$ , every one of these subintervals is itself subdivided into a number of segments of equal length  $K/l'$ , except perhaps near the boundary points. Letting  $t \rightarrow +\infty$ , the above reasoning can be applied to each of the subintervals; thus one obtains, instead of one Riemann integral like the one in (2.14), a finite sum of Riemann integrals whose intervals of integration are strung together to form  $[0, K]$ . This sum of Riemann integrals defines the value of

$$\int_0^1 w_0(x) dx \tag{2.16}$$

in the considered case and is therefore equal to unity in virtue of the normalization. Once more, this argument can be extended without difficulty to  $x = 0$  and  $x = 1$ . Consequently, (2.15) remains valid when  $w_0(x)$  is piecewise continuous, exhibiting a limited number of finite jumps in  $[0, 1]$ . Analogous considerations also hold in cases where there occur stronger singularities of  $w_0(x)$  in  $[0, 1]$ , making (2.16) an improper integral on the condition, however, that the integral be convergent so that  $w_0(x)$  is normalizable in the strict

sense. An example of such a case is

$$w_0(x) = 1/(2|1 - 2x|^{1/2}), \quad 0 \leq x \leq 1. \tag{2.17}$$

Here, use can be made of suitable cutoffs and one can show that  $w_t(x)$ , as given by (2.12), converges toward the right-hand side of (2.15), multiplied by the definition of (2.16) as a convergent improper integral; thus, e.g., in the case of (2.17),

$$\int_0^1 w_0(x) dx = \lim_{\epsilon \rightarrow +0} \frac{1}{2} \int_0^{1-\epsilon} \frac{dx}{(1 - 2x)^{1/2}} + \lim_{\epsilon \rightarrow +0} \frac{1}{2} \int_{\frac{1}{2}+\epsilon}^1 \frac{dx}{(2x - 1)^{1/2}} = 1.$$

In conclusion, the following theorem can be stated.

Under a discrete-time mapping of the type (2.4) with  $k \in ]0, 1[$ , subjected iteratively to the associated transformation (2.7), a single-valued, real, continuous, or piecewise continuous function  $w_0(x)$ , defined on  $[0, 1]$  and normalized to unity, where

$$\int_0^1 w_0(x) dx$$

may be an improper, yet convergent integral on account of the singularities of  $w_0(x)$  in  $[0, 1]$ , leads to an infinite sequence of functions  $\{w_t(x) | t \in \mathbb{N}, x \in [0, 1]\}$  which converges toward a limit function  $w(x)$  when  $t \rightarrow +\infty$ , with

$$w(x) = \frac{df^{-1}(x)/dx}{K(k) [1 - (f^{-1}(x))^2]^{1/2} [1 - k^2 (f^{-1}(x))^2]^{1/2}}, \quad 0 \leq x \leq 1. \tag{2.18}$$

Note that the limit function is independent of the integer parameter  $l$  appearing in (2.4), but fully determined by the modulus  $k$  and the inverse of the function  $f$  appearing in the parametric representation (2.3). For  $k = 0$ , which means using (2.1) and (2.2), (2.18) reduces to

$$w(x) = \frac{2 df^{-1}(x)/dx}{\pi [1 - (f^{-1}(x))^2]^{1/2}}, \quad 0 \leq x \leq 1. \tag{2.19}$$

Since (2.7) conserves normalization,  $w(x)$  given by (2.18) or by (2.19) is also normalized to unity, as can be verified by inspection.

In regard to the preceding theorem, let me emphasize that the initial function  $w_0(x)$  must no longer be continuous on  $[0, 1]$  in order to ensure convergence of the sequence  $\{w_t(x) | t \in \mathbb{N}, x \in [0, 1]\}$  toward the limit (2.18), in contrast to the continuity imposed on  $w_0(x)$  by the sufficient condition (1.15) which was part of the theorems comprised in my earlier articles.<sup>8,9</sup> In those theorems, leaving out the condition (1.15) would not exclude the possibility of convergence toward the limit function  $\pi^{-1}[x(1-x)]^{-1/2}$ , but convergence would be unproved when  $w_0(x)$  involves singularities in  $[0, 1]$  within the framework of my earlier proofs. The present analysis clearly removes this drawback. It shows that in virtue of the more powerful method used in the present article, condition (1.15) can be deleted after all, in agreement with Barbour's result<sup>10</sup> in the case of finite jumps of  $w_0(x)$  in  $[0, 1]$ .

Reconsidering the two conjugate dynamical laws (2.2)

with  $l = 2$  and (2.2') where the conjugating function  $h$  is such that both  $f$  and  $hf$  have the properties required in the beginning of this section, we clearly have as a relation between the corresponding invariant densities  $w(x)$  and  $w(x;h)$ ,

$$w(x;h) = w[h^{-1}(x)] \frac{dh^{-1}(x)}{dx},$$

a relation which remains valid for (2.4) with  $l = 2$  and (2.4'). The examples given in Ref. 12 (cf. Sec. IV and Figs. 2 and 3) are comprised in the present theory. Indeed, the different conjugating functions  $h(\theta)$  which are coupled with the special case called  $N_2(\theta)$  have the properties specified in the beginning of our Sec. II; therefore this is also the case with  $h[f(y)]$ . Here,  $N_2(\theta)$  is such that

$$x_{t+1} = N_2(x_t)$$

is nothing but the dynamical law (1.3). As is shown in example (5) of Sec. III, this map is included in (2.2), where one puts  $l = 2$  and the functions which I called  $f(y)$  and  $f^{-1}(y)$  are

$$f(y) = (2/\pi)\arcsin y, \quad f^{-1}(y) = \sin(\pi y/2), \quad 0 \leq y \leq 1.$$

In certain applications, as we shall see in Sec. III,  $w(x)$  turns out to be symmetric about  $x = \frac{1}{2}$ . This is not necessarily a consequence of symmetry about  $x_t = \frac{1}{2}$  in the right-hand side of a discrete-time map written in the form (1.1). Let us, for instance, recall the case of the quadratic map (1.7) where, as was pointed out by Falk,<sup>6</sup>  $w_0(x) = 1$  gives rise to the asymmetric functions  $w_t(x)$ ,  $t = 1, 2, \dots$ , despite the symmetry of the right-hand side in (1.7). However, the symmetry about  $x = \frac{1}{2}$  is restored in the limit function  $w(x)$  given by (1.8). In the case of  $l = 3$  in (1.10), being the cubic map

$$x_{t+1} = x_t(3 - 4x_t)^2, \quad 0 \leq x_t \leq 1, \quad (2.20)$$

which does not have the symmetry about  $x_t = \frac{1}{2}$ , starting from either a symmetric or an asymmetric initial function  $w_0(x)$  leads, at any rate, to (1.8), having the symmetry. On the contrary, when one maps linearly  $[-1, 1]$  onto  $[0, 1]$  both for  $x_t$  and  $x_{t+1}$  in (1.5), so as to obtain

$$x_{t+1} = 1 - |1 - 2x_t|^{1/2}, \quad 0 \leq x_t \leq 1, \quad (2.21)$$

which is, of course, also cusp shaped with the singular point at  $x_t = \frac{1}{2}$ , the discrete-time map is symmetric about that point. The evolution corresponding to (2.21) is

$$w_{t+1}(x) = (1-x)[w_t(x-x^2/2) + w_t(1-x+x^2/2)], \quad 0 \leq x \leq 1.$$

Starting from the normalized uniform distribution  $w_0(x) = 1$  which also has the symmetry, one finds, strangely enough,

$$w_t(x) = 2(1-x), \quad \forall t \in \mathbb{N}_0, \quad 0 \leq x \leq 1,$$

and hence

$$w(x) = \lim_{t \rightarrow +\infty} w_t(x) = 2(1-x), \quad (2.22)$$

which is the invariant density, no matter from which initial distribution one starts since it is, under the requirement of normalization to unity, the unique solution of

$$w(x) = (1-x)[w(x-x^2/2) + w(1-x+x^2/2)], \quad 0 \leq x \leq 1,$$

according to Hemmer.<sup>5</sup> Yet, it does not possess the symmetry about  $x = \frac{1}{2}$ . Irrespective of these examples, the explicit form of  $w(x)$ , namely (2.18), associated with a discrete-time map of the type (2.4), enables us to find a necessary and sufficient condition that  $w(x)$  be symmetric about  $x = \frac{1}{2}$ . Indeed, if  $w(x) = w(1-x)$ ,  $\forall x \in [0, 1]$ , then

$$\int_0^y w(x) dx = \int_{1-y}^1 w(x') dx', \quad \forall y \in [0, 1],$$

or

$$\text{sn}^{-1} f^{-1}(y) = K(k) - \text{sn}^{-1} f^{-1}(1-y), \quad \forall y \in [0, 1].$$

If we put

$$\text{sn}^{-1} f^{-1}(y) = u, \quad 0 \leq u \leq K(k),$$

then

$$\text{sn} u = f^{-1}(y), \quad \text{sn}[K(k) - u] = f^{-1}(1-y).$$

But

$$\text{sn}[K(k) - u] = (\text{cn} u) / (\text{dn} u),$$

and therefore, we find as a necessary and sufficient condition for symmetry,

$$f^{-1}(1-x) = \left( \frac{1 - (f^{-1}(x))^2}{1 - k^2 (f^{-1}(x))^2} \right)^{1/2}, \quad \forall x \in [0, 1], \quad (2.23)$$

in which we returned to  $x$  as the independent variable. If one keeps in mind that  $f^{-1}(x)$  increases steadily from 0 to 1 as  $x$  increases from 0 to 1, (2.23) can be rewritten rationally as  $(f^{-1}(1-x))^2 = [1 - (f^{-1}(x))^2] / [1 - k^2 (f^{-1}(x))^2]$ .

$$(2.24)$$

For  $k = 0$ , which corresponds to (2.2), the condition is

$$(f^{-1}(x))^2 + (f^{-1}(1-x))^2 = 1. \quad (2.25)$$

Expressed in terms of  $f$  itself, (2.23) becomes

$$f(x) + f(\sqrt{(1-x^2)/(1-k^2x^2)}) = 1, \quad \forall x \in [0, 1], \quad k \in [0, 1], \quad (2.26)$$

with, as a special case for  $k = 0$ ,

$$f(x) + f(\sqrt{1-x^2}) = 1, \quad \forall x \in [0, 1]. \quad (2.27)$$

Examples of solutions of Eq. (2.27) which increase steadily from 0 to 1 as  $x$  increases from 0 to 1, are

$$f(x) = x^2 \quad \text{and} \quad f(x) = (2/\pi)\arcsin x. \quad (2.28)$$

They will be brought in relation with certain special cases comprised in Sec. III. Note that in regard to possible symmetry of  $w(x)$  about  $x = \frac{1}{2}$ , the integer parameter  $l$  in (2.2) or (2.4) is totally unimportant.

The explicit formula (2.18) enables us to discuss some characteristic features of the limit function  $w(x)$ . In virtue of the properties  $f(x)$  is assumed to have for  $x \in [0, 1]$ , we have

$$\frac{df^{-1}(x)}{dx} \geq 0,$$

entailing that

$$w(x) \geq 0, \quad 0 \leq x \leq 1, \quad (2.29)$$

irrespective of whether the initial function  $w_0(x)$  is positive

semidefinite or not. Furthermore, if

$$\left(\frac{df^{-1}(x)}{dx}\right)_{x=1} > 0,$$

which means  $0 \leq f'(1) < +\infty$ , then one certainly has  $w(1) = +\infty$ . If at the same time  $f^{-1}(x)$  satisfies (2.24) or (2.25), one obviously also has  $w(0) = +\infty$  in virtue of the symmetry of  $w(x)$  about  $x = \frac{1}{2}$ . When  $x \in [0, 1[$ , the denominators in (2.18) and (2.19) are finite and positive. To any abscissa  $\alpha \in [0, 1[$ , where  $f'(\alpha)$  happens to be zero, there corresponds  $w(\beta) = +\infty$  on account of  $(df^{-1}(x)/dx)_{x=\beta} = +\infty$ , with  $\beta = f(\alpha) \in [0, 1[$ . This can solely occur at isolated abscissas  $\beta$ . Furthermore, also for  $x \in [0, 1[$ , the equality sign in (2.29) can only hold at points where  $df^{-1}(x)/dx = 0$ , which is associated with  $f'(x)$  becoming infinite. Thus if  $f(x)$  is a function with a continuous derivative in  $[0, 1[$ , then  $w(x) > 0, x \in [0, 1[$ . It could happen, however, that  $f'(x)$  is only piecewise continuous in  $[0, 1[$ . Then, for an abscissa  $\kappa \in [0, 1[$ , where  $f'(\kappa) = +\infty$ , there comes  $w(\lambda) = 0$  on account of  $(df^{-1}(x)/dx)_{x=\lambda} = 0$ , with  $\lambda = f(\kappa)$ . This can occur solely at isolated abscissas  $\lambda$  because  $f(x)$  is assumed to be single valued. It could also happen, still when  $f'(x)$  is only piecewise continuous in  $[0, 1[$ , that  $f'(x)$  is double valued at one or several abscissas  $\mu$ , which would mean that  $f(x)$  has a left derivative and a right derivative, unequal at such an abscissa. Then  $df^{-1}(x)/dx$  and therefore  $w(x)$  are also double valued at  $x = \nu, \nu = f(\mu)$ , with the discontinuity a finite or an infinite jump.

The particular form of the right-hand side in (2.18) makes it possible to establish a way to solve the inverse problem: Given a positive-semidefinite invariant density, find the discrete-time maps of the type (2.4) or (2.2) from which it originates. Indeed, integrating on both sides between 0 and  $y$  in (2.18), one obtains

$$\int_0^y w(x) dx = \frac{1}{K(k)} \operatorname{sn}^{-1} f^{-1}(y), \quad 0 \leq y \leq 1,$$

in virtue of (2.5), and therefore

$$f^{-1}(y) = \operatorname{sn}\left(K(k) \int_0^y w(x) dx\right), \quad 0 \leq y \leq 1. \quad (2.30)$$

Then  $f$  clearly follows from the inversion of  $f^{-1}$ . This procedure will be applied in two of the examples considered in Sec. III. Equation (2.30) even permits the discrete-time map (2.4), which yields  $w(x)$ , to be rewritten solely in terms of  $w(x)$ :

$$\operatorname{sn}\left(K(k) \int_0^{x_{i+1}} w(x) dx\right) = \left| \operatorname{sn}\left(lK(k) \int_0^{x_i} w(x) dx\right) \right|, \quad 0 \leq x_i \leq 1, \quad 0 \leq x_{i+1} \leq 1, \quad l \in \{2, 3, \dots\}, \quad k \in [0, 1]. \quad (2.31)$$

In the special case  $k = 0$ , the corresponding formulas are

$$f^{-1}(y) = \sin\left(\frac{\pi}{2} \int_0^y w(x) dx\right), \quad (2.32)$$

$$\sin\left(\frac{\pi}{2} \int_0^{x_{i+1}} w(x) dx\right) = \left| \sin\left(\frac{l\pi}{2} \int_0^{x_i} w(x) dx\right) \right|. \quad (2.33)$$

Expressed in this manner, a discrete-time map is in an im-

plicit form, i.e., solved neither with respect to  $x_{i+1}$  nor with respect to  $x_i$ .

### III. SOME EXAMPLES OF PRACTICAL APPLICATION

(1) In the case of a polynomial discrete-time map as described by (1.10) or (1.11), we have  $k = 0$  and  $f(x) = x^2$  with  $x \in [0, 1]$ . Hence,  $K = \pi/2$  and  $f^{-1}(x) = \sqrt{x}$ , and consequently there is convergence of  $\{w_i(x) | i \in \mathbb{N}\}$  toward

$$w(x) = \frac{2 d\sqrt{x}/dx}{\pi(1-x)^{1/2}} = \frac{1}{\pi[x(1-x)]^{1/2}}, \quad 0 \leq x \leq 1, \quad (3.1)$$

which confirms the result (1.8) obtained by several authors in the special case of the quadratic map. On the basis of Sec. II, one understands why all the polynomial discrete-time mappings comprised in (1.10) give rise to the same invariant density. Since  $f'(1) = 2$ , thus positive,  $w(1) = +\infty$ . Since  $f'(0) = 0$ , we have  $(df^{-1}(x)/dx)_{x=0} = +\infty$  and consequently  $w(0) = +\infty$ . Because of the continuity of  $f'(x)$  in  $[0, 1]$ ,  $w(x) > 0$  for all  $x \in [0, 1]$ . The function  $f(x) = x^2$  is a solution of (2.27), which explains the symmetry of (3.1) about  $x = \frac{1}{2}$ .

(2) The most direct generalization of (1.11) is

$$x_i = (\operatorname{sn}(u, k))^2, \quad x_{i+1} = (\operatorname{sn}(lu, k))^2, \quad 0 \leq u \leq K(k), \quad l \in \{2, 3, \dots\}, \quad 0 < k < 1. \quad (3.2)$$

In this case, the discrete-time map is rational. When  $l = 2$ , it is

$$x_{i+1} = [4x_i(1-x_i)(1-k^2x_i)] / (1-k^2x_i^2)^2, \quad (3.3)$$

being a generalization of the quadratic map (1.7). As expected, its inverse is expressed by two formulas:

$$\begin{aligned} x_i^{(1)} = r_1(x_{i+1}) &= [1 - (1 - x_{i+1})^{1/2}] \\ &\times [1 + (1 - k^2x_{i+1})^{1/2}]^{-1}, \\ x_i^{(2)} = r_2(x_{i+1}) &= [1 + (1 - x_{i+1})^{1/2}] \\ &\times [1 + (1 - k^2x_{i+1})^{1/2}]^{-1}. \end{aligned}$$

Here again, we have  $f(x) = x^2, f^{-1}(x) = \sqrt{x}$ , and consequently, irrespective of the value of  $l \in \{2, 3, \dots\}$ , the invariant density associated with (3.2) is

$$w(x) = \{2K(k)[x(1-x)(1-k^2x)]^{1/2}\}^{-1}, \quad 0 < k < 1, \quad (3.4)$$

which generalizes (3.1). It is not symmetric about  $x = \frac{1}{2}$  because  $x^2$  is not a solution of (2.26).

(3) With the discrete-time mapping

$$x_{i+1} = x_i \left| U_{l-1}(\sqrt{1-x_i^2}) \right| \quad (3.5)$$

stemming from

$$x_i = \sin \theta, \quad x_{i+1} = |\sin l\theta|, \quad 0 \leq \theta \leq \pi/2,$$

there is associated as invariant density,

$$w(x) = 2/[\pi(1-x^2)^{1/2}], \quad 0 \leq x \leq 1, \quad (3.6)$$

since  $f(x) = x$  and  $f^{-1}(x) = x$ . As  $x$  is not a solution of (2.27), there is no symmetry about  $x = \frac{1}{2}$ , but on account of  $f'(x) = 1$ , which entails  $df^{-1}(x)/dx = 1, \forall x \in [0, 1]$ ,  $w(x) > 0$ , and  $w(1) = +\infty$ .

(4) Intermediate between  $f(x) = x^2$  and  $f(x) = x$ , we

can consider

$$f(x) = cx^2 + (1 - c)x, \quad 0 \leq x \leq 1, \quad (3.7)$$

increasing steadily from 0 to 1 as  $x$  increases from 0 to 1 when  $-1 \leq c \leq 1$ . For  $f^{-1}(x)$ , we find

$$f^{-1}(x) = 2x / \{1 - c + [(1 - c)^2 + 4cx]^{1/2}\}.$$

Putting  $k = 0$ , the invariant density corresponding to the discrete-time mappings

$$\begin{aligned} x_t &= c(\sin \theta)^2 + (1 - c)\sin \theta, \\ x_{t+1} &= c(\sin l\theta)^2 + (1 - c)|\sin l\theta|, \\ &0 \leq \theta \leq \pi/2, \quad l \in \{2, 3, \dots\}, \end{aligned}$$

or

$$\begin{aligned} x_{t+1} &= c(\sin \theta U_{l-1}(\cos \theta))^2 \\ &+ (1 - c)\sin \theta |U_{l-1}(\cos \theta)|, \end{aligned} \quad (3.8)$$

in which

$$\begin{aligned} \sin \theta &= 2x_t / \{1 - c + [(1 - c)^2 + 4cx_t]^{1/2}\}, \\ \cos \theta &= \left\{ 1 - \frac{4x_t^2}{\{1 - c + [(1 - c)^2 + 4cx_t]^{1/2}\}^2} \right\}^{1/2}, \end{aligned}$$

is

$$w(x) = 2 \left\{ \pi [(1 - c)^2 + 4cx]^{1/2} \left[ 1 - \frac{4x^2}{\{1 - c + [(1 - c)^2 + 4cx]^{1/2}\}^2} \right]^{1/2} \right\}^{-1}. \quad (3.9)$$

(5) In Ref. 4, it was stated that the discrete-time mapping

$$x_{t+1} = \begin{cases} 2x_t, & 0 \leq x_t \leq \frac{1}{2}, \\ 2(1 - x_t), & \frac{1}{2} \leq x_t \leq 1 \end{cases} \quad (3.10)$$

(rewritten in our notation) gives rise to the constant limit density

$$w(x) = 1, \quad 0 \leq x \leq 1, \quad (3.11)$$

according to Kac.<sup>1</sup> It is amusing to verify this result, proceeding in the backward direction starting from  $w(x) = 1$ . Let us assume that the function called  $f$  throughout this article exists and let us at first put the modulus  $k$  equal to zero, for simplicity. Then, according to (2.19), we must have

$$\frac{2 df^{-1}(x)/dx}{\pi [1 - (f^{-1}(x))^2]^{1/2}} = 1, \quad 0 \leq x \leq 1.$$

Integration from 0 to  $y$  on both sides, where  $y \in [0, 1]$ , yields

$$(2/\pi) \arcsin f^{-1}(y) = y,$$

and so

$$f^{-1}(y) = \sin(\pi y/2), \quad 0 \leq y \leq 1. \quad (3.12)$$

In turn, this gives

$$f(x) = (2/\pi) \arcsin x, \quad 0 \leq x \leq 1, \quad (3.13)$$

which satisfies the condition (2.27) in virtue of the symmetry of  $w(x) = 1$  about  $x = \frac{1}{2}$ . Since  $f(x)$  has all the necessary properties to be inserted into (2.1), we know at once that for every  $l \in \{2, 3, \dots\}$ ,

$$\begin{aligned} x_t &= (2/\pi) \arcsin(\sin \theta) = (2/\pi)\theta, \\ x_{t+1} &= (2/\pi) \arcsin(|\sin l\theta|), \quad 0 \leq \theta \leq \pi/2, \end{aligned}$$

is the parametric representation of a discrete-time mapping to which there corresponds  $w(x) = 1$ . Eliminating  $\theta$ , the result is

$$x_{t+1} = (2/\pi) \arcsin(|\sin(l\pi x_t/2)|), \quad (3.14)$$

or

$$x_{t+1} = \begin{cases} lx_t, & 0 \leq x_t \leq 1/l, \\ 2 - lx_t, & 1/l \leq x_t \leq 2/l, \\ -2 + lx_t, & 2/l \leq x_t \leq 3/l, \\ \vdots \\ (-1)^l(2[l/2] - lx_t), & (l-1)/l \leq x_t \leq 1. \end{cases} \quad (3.15)$$

The case (3.10) is simply that of  $l = 2$ . It is somewhat astonishing that repeating the calculations with  $0 < k < 1$ , hence with elliptic functions, does not yield a different result compared to what we just obtained. The finiteness of  $w(x)$  at  $x = 1$  implies that  $(df^{-1}(x)/dx)_{x=1}$  be equal to zero or, equivalently, that  $f'(1) = +\infty$ . According to (3.12) and (3.13), such is indeed the case.

(6) Let us return to the subject of Ref. 5, namely, the cusp-shaped map (1.5) and the associated limit distribution (1.6), which, after linear transformation so that both  $x_t$  and  $x_{t+1}$  belong to  $[0, 1]$ , give rise to (2.21) and (2.22), respectively. Again, if we assume the existence of  $f$  and its inverse  $f^{-1}$ , as well as  $k = 0$ , we can integrate from 0 to  $y$  in (2.19) with  $w(x) = 2(1 - x)$ . The result is

$$2y - y^2 = (2/\pi) \arcsin f^{-1}(y)$$

and so

$$f^{-1}(y) = \sin(\pi y - (\pi/2)y^2), \quad 0 \leq y \leq 1.$$

This, in turn, yields

$$f(x) = 1 - (1 - (2/\pi) \arcsin x)^{1/2}, \quad 0 \leq x \leq 1. \quad (3.16)$$

Here,  $f(x)$  and its inverse satisfy all requirements to be utilized for the construction of a discrete-time map of the type (2.1) and (2.2):

$$\begin{aligned} x_t &= 1 - [1 - (2/\pi) \arcsin(\sin \theta)]^{1/2} \\ &= 1 - (1 - 2\theta/\pi)^{1/2}, \\ x_{t+1} &= 1 - [1 - (2/\pi) \arcsin(|\sin l\theta|)]^{1/2}, \quad 0 \leq \theta \leq \pi/2, \end{aligned} \quad (3.17)$$

or

$$x_{t+1} = 1 - \{1 - (2/\pi)\arcsin [|\sin l\pi(x_t - x_t^2/2)|]\}^{1/2},$$

$$0 \leq x_t \leq 1, \quad l \in \{2, 3, \dots\}. \quad (3.18)$$

We know that solving this discrete-time map with respect to  $x_t$  results in  $l$  real expressions [see (2.6) with  $k = 0$ ]. Hence, in order to compare our result to (2.21), namely,

$$x_{t+1} = 1 - |1 - 2x_t|^{1/2}, \quad 0 \leq x_t \leq 1, \quad (3.19)$$

we have to put  $l = 2$ . In this case, we find

$$x_t = 1 - (1 - 2\theta/\pi)^{1/2},$$

$$x_{t+1} = 1 - [1 - (2/\pi)\arcsin(\sin 2\theta)]^{1/2}$$

$$= 1 - |1 - 4\theta/\pi|^{1/2}, \quad 0 \leq \theta \leq \pi/2, \quad (3.20)$$

or

$$x_{t+1} = 1 - |1 - 4x_t + 2x_t^2|^{1/2}, \quad 0 \leq x_t \leq 1. \quad (3.21)$$

This is a cusp-shaped return map, but it is asymmetric with respect to  $x = \frac{1}{2}$ , and therefore it differs quantitatively from (3.19). The Cartesian graph of (3.21) starts at the origin with slope 2 and increases steadily toward the point  $(1 - (\sqrt{2}/2), 1)$ , where the slope is  $+\infty$ . The slope jumps to  $-\infty$  as the curve starts to decrease, tending toward  $(1, 0)$ , where the slope becomes 0. Both under (3.19) and (3.21) there is convergence of  $\{w_t(x) | t \in \mathbb{N}, x \in [0, 1]\}$  toward  $2(1 - x)$  as  $t \rightarrow +\infty$ , but this is not paradoxical. Different discrete-time maps may give rise to the same invariant density, as is, for instance, the case with the maps comprised in (2.2) or (2.4) for  $l = 2, 3, \dots$ . But here, the difference lies deeper: (3.19) is simply not a discrete-time map of the kind (2.2). Parametrically, (3.19) can be represented as

$$x_t = 2\theta/\pi, \quad x_{t+1} = 1 - |1 - 4\theta/\pi|^{1/2}, \quad 0 \leq \theta \leq \pi/2.$$

Still making use of  $f(x)$  as defined by (3.16), this parametric representation of (3.19) becomes

$$x_t = 2f(\sin \theta) - [f(\sin \theta)]^2,$$

$$x_{t+1} = f(\sin 2\theta), \quad 0 \leq \theta \leq \pi/2, \quad (3.22)$$

confirming our statement made near the beginning of Sec. II, namely, that the type of discrete-time maps represented by (2.4) does not comprise all possible maps. It is practically certain that to all discrete-time maps with parametric form

$$x_t = 2f(\sin \theta) - [f(\sin \theta)]^2, \quad x_{t+1} = f(|\sin l\theta|),$$

$$0 \leq \theta \leq \pi/2, \quad l \in \{2, 3, \dots\}, \quad (3.23)$$

or even more generally,

$$x_t = 2f(\text{sn}(u, k)) - [f(\text{sn}(u, k))]^2,$$

$$x_{t+1} = f(|\text{sn}(lu, k)|), \quad 0 \leq u \leq K(k), \quad l \in \{2, 3, \dots\}, \quad (3.24)$$

there corresponds a formula to calculate the limit function  $w(x)$ , being the counterpart of (2.19) or (2.18), respectively. In the special case of (3.16), this formula, where  $k = 0$ , should yield the same result as (2.19), i.e.,  $w(x) = 2(1 - x)$ .

(7) In Ref. 13, one finds two examples of piecewise linear mappings of  $[0, 1]$  onto itself, symmetric with respect to  $x = \frac{1}{2}$  and giving rise to limit functions which are piecewise constant, with a finite jump at  $x = \frac{1}{2}$  (cf. Fig. 2 of Ref. 14). It appears of interest to calculate with which discrete-time

maps of the type (2.2), where  $l = 2$ , these limit functions are associated. Györgyi and Szépfalussy<sup>14</sup> consider

$$w(x) = \begin{cases} 1 + c, & 0 \leq x < \frac{1}{2}, \\ 1 - c, & \frac{1}{2} \leq x \leq 1, \end{cases} \quad c \in ]-1, 1[,$$

and their first numerical example corresponds to  $c = -0.6$ , in which case they obtain the symmetric map

$$x_{t+1} = \begin{cases} 5x_t, & 0 \leq x_t < \frac{1}{10}, \\ \frac{5}{2}x_t + \frac{3}{8}, & \frac{1}{10} \leq x_t < \frac{1}{2}, \\ \frac{13}{8} - \frac{5}{2}x_t, & \frac{1}{2} \leq x_t < \frac{9}{10}, \\ 5(1 - x_t), & \frac{9}{10} \leq x_t \leq 1. \end{cases}$$

This discrete-time map is not comprised in (2.2) with  $l = 2$ , since the present formalism yields another dynamical law, nonsymmetric with respect to  $x_t = \frac{1}{2}$ . Indeed, inserting

$$w(x) = \begin{cases} 0.4, & 0 \leq x < \frac{1}{2}, \\ 1.6, & \frac{1}{2} \leq x \leq 1, \end{cases}$$

into (2.19) and integrating from 0 to  $y$ , we obtain

$$(2/\pi)\arcsin f^{-1}(y) = \begin{cases} 2y/5, & 0 \leq y < \frac{1}{2}, \\ 8y/5 - \frac{3}{5}, & \frac{1}{2} \leq y \leq 1, \end{cases}$$

and consequently,

$$f^{-1}(y) = \begin{cases} \sin(\pi y/5), & 0 \leq y < \frac{1}{2}, \\ \sin(4\pi y/5 - 3\pi/10), & \frac{1}{2} \leq y \leq 1. \end{cases}$$

In turn, this leads to

$$f(y) = \begin{cases} (5/\pi)\arcsin y, & 0 \leq y \leq \sin(\pi/10) = (\sqrt{5} - 1)/4, \\ \frac{3}{8} + (5/4\pi)\arcsin y, & (\sqrt{5} - 1)/4 \leq y \leq 1. \end{cases}$$

Equation (2.2) with  $l = 2$  yields

$$x_{t+1} = \begin{cases} f(\sin(2\pi x_t/5)), & 0 \leq x_t < \frac{1}{2}, \\ f(\sin[(8x_t - 3)\pi/5]), & \frac{1}{2} \leq x_t \leq 1. \end{cases}$$

Taking into account the various intervals of validity, we obtain as final result:

$$x_{t+1} = \begin{cases} 2x_t, & 0 \leq x_t < \frac{1}{4}, \\ x_t/2 + \frac{3}{8}, & \frac{1}{4} \leq x_t < \frac{1}{2}, \\ 2x_t - \frac{3}{8}, & \frac{1}{2} \leq x_t < \frac{11}{16}, \\ \frac{19}{8} - 2x_t, & \frac{11}{16} \leq x_t < \frac{15}{16}, \\ 8(1 - x_t), & \frac{15}{16} \leq x_t \leq 1, \end{cases}$$

which is a piecewise linear, nonsymmetric discrete-time map. The absence of symmetry with respect to  $x_t = \frac{1}{2}$  could be foreseen since  $f^{-1}$  does not satisfy (2.25) in this case. The situation is similar to that of our example (6).

In the case  $c = 0.6$ , Györgyi and Szépfalussy obtain

$$x_{t+1} = \begin{cases} \frac{5}{2}x_t, & 0 \leq x_t < \frac{2}{3}, \\ 5x_t - \frac{3}{2}, & \frac{2}{3} \leq x_t < \frac{1}{2}, \\ \frac{7}{2} - 5x_t, & \frac{1}{2} \leq x_t < \frac{2}{3}, \\ \frac{5}{2}(1 - x_t), & \frac{2}{3} \leq x_t \leq 1, \end{cases}$$

whereas analogous calculations as we carried out above yield

$$x_{t+1} = \begin{cases} 2x_t, & 0 \leq x_t \leq \frac{1}{4}, \\ 8x_t - \frac{3}{2}, & \frac{1}{4} \leq x_t \leq \frac{5}{16}, \\ \frac{7}{2} - 8x_t, & \frac{5}{16} \leq x_t \leq \frac{3}{8}, \\ \frac{5}{4} - 2x_t, & \frac{3}{8} \leq x_t \leq \frac{1}{2}, \\ \frac{1}{2}(1 - x_t), & \frac{1}{2} \leq x_t \leq 1. \end{cases}$$

The calculations may also be worked out for arbitrary  $c \in ] - 1, 1[$  and the result compared to Eq. (3.9) of Ref. 13, but several subcases have to be distinguished.

(8) Finally, let  $f(x)$  be of the form

$$f(x) = a + (b + b'x^2)^{1/2}, \quad 0 \leq x \leq 1.$$

This should be a single-valued, real, continuous function of  $x$  on  $[0, 1]$ , increasing steadily from 0 to 1 as  $x$  increases from 0 to 1. Such is the case when

$$a = -b^{1/2}, \quad b' = 1 + 2b^{1/2}, \quad b \geq 0.$$

Putting  $b^{1/2} = c$ , we obtain

$$x_{t+1} = -c + \left\{ c^2 + \frac{4(2c+1)^2 x_t (1-x_t)(2c+x_t)(2c+1+x_t)[2c+1-k^2 x_t(2c+x_t)]}{[(2c+1)^2 - k^2 x_t^2 (2c+x_t)^2]^2} \right\}^{1/2}.$$

Even for  $l=3$ , the formula relating  $x_{t+1}$  to  $x_t$  is already considerably more complicated and yet, for any  $l \in \{2, 3, \dots\}$ , the sequence  $\{w_t(x) | t \in \mathbb{N}, x \in [0, 1]\}$  starting from an arbitrary normalized initial  $w_0(x)$  converges toward

$$w(x) = \frac{(2c+1)^{1/2}(c+x)}{K(k)\{x(1-x)(2c+x)(2c+1+x)[2c+1-k^2x(2c+x)]\}^{1/2}}.$$

For  $c > 0$ , we indeed have  $w(0) = w(1) = +\infty$ . For  $c = 0$ , wherefore  $f(x) = x, f'(x) = 1$ , the singularity at  $x = 0$  disappears as  $w(x)$  reduces to

$$w(x) = \{K(k)[(1-x^2)(1-k^2x^2)]^{1/2}\}^{-1},$$

generalizing (3.6).

## ACKNOWLEDGMENT

I wish to thank Dr. H. De Meyer, Senior research associate N.F.W.O. (Belgium) for having drawn my attention to a number of articles concerning this subject and also for stimulating discussions resulting from his keen interest in my research work.

<sup>1</sup>M. Kac, *Ann. Math.* **47**, 33 (1946).

<sup>2</sup>S. M. Ulam and J. von Neumann, *Bull. Am. Math. Soc.* **53**, 1120 (1947).

$$f(x) = -c + [c^2 + (2c+1)x^2]^{1/2}, \quad 0 \leq x \leq 1, \quad c \geq 0, \quad (3.25)$$

satisfying all requirements and having its derivative function  $f'(x)$  continuous in  $[0, 1]$ . It constitutes a way of generalizing  $f(x) = x$ , which is the special case  $c = 0$ . When  $c > 0$ ,  $f'(0) = 0$  and  $f'(1) = (2c+1)/(c+1) > 0$ , and therefore,  $w(0) = w(1) = +\infty$  can be expected for any  $k \in [0, 1[$ . Using the Jacobian sn function, an acceptable discrete-time map is obtained in parametric form:

$$\begin{aligned} x_t &= -c + [c^2 + (2c+1)(\text{sn}(u, k))^2]^{1/2}, \\ x_{t+1} &= -c + [c^2 + (2c+1)(\text{sn}(lu, k))^2]^{1/2}, \quad (3.26) \\ 0 &\leq u \leq K(k), \quad c \geq 0, \end{aligned}$$

for every value of  $l$  belonging to  $\{2, 3, \dots\}$ . This example is given in order to show that a function  $f$  which is still relatively simple may give rise to discrete-time maps of high degree of complexity, especially as the integer value of  $l$  increases. The simplest map contained in (3.26) is that corresponding to  $l = 2$ . Its explicit form is

<sup>3</sup>A. Lasota and J. A. Yorke, *Trans. Am. Math. Soc.* **186**, 481 (1973).

<sup>4</sup>J. B. McGuire and C. J. Thompson, *Bull. Aus. Math. Soc.* **22**, 133 (1980).

<sup>5</sup>P. C. Hemmer, *J. Phys. A: Math. Gen.* **17**, L247 (1984).

<sup>6</sup>H. Falk, *Phys. Lett. A* **105**, (3), 101 (1984).

<sup>7</sup>V. M. Nandakumaran, *J. Phys. A: Math. Gen.* **18**, L1021 (1985).

<sup>8</sup>C. C. Grosjean, *J. Phys. A: Math. Gen.* **19**, 3535 (1986).

<sup>9</sup>C. C. Grosjean, *Acad. Analecta (Proc. R. Ac. Sc., Litt., Fine Arts, Brussels, Belgium)* **48**(5), 95 (1986).

<sup>10</sup>A. D. Barbour, *J. Phys. A: Math. Gen.* **19**, 3921 (1986).

<sup>11</sup>P. R. Halmos, *Lectures on Ergodic Theory* (Chelsea, New York, 1956).

<sup>12</sup>S. M. Ulam, *Intersci. Tracts Pure Appl. Math.* **8**, 73 (1960).

<sup>13</sup>S. Grossmann and S. Thomae, *Z. Naturforsch.* **32a**, 1353 (1977).

<sup>14</sup>G. Györgyi and P. Szépfalussy, *Z. Phys. B* **55**, 179 (1984).



# On unitary $SU(N)$ ordered exponentials in a strong coupling limit

H. M. Fried<sup>a)</sup>

*Physique Théorique, Université de Nice,<sup>b)</sup> 06034 Nice Cedex, France*

(Received 1 December 1986; accepted for publication 4 February 1987)

A construction is given of the leading, averaged output dependence of a unitary ordered exponential in the strong coupling limit of rapidly fluctuating input. While valid for any  $SU(N)$ , the method does not provide "fine structure" corrections to the leading output behavior. Numerical illustrations are given using a simple  $SU(3)$  example.

## I. INTRODUCTION

This paper is meant as a first, and partial generalization to  $SU(N)$  of the techniques and results of two previous papers<sup>1,2</sup> dealing with the approximation of ordered  $SU(2)$  exponentials in the stochastic limit. Such ordered exponentials (OE's) are found in almost every field of mathematical physics which deals with more than one degree of freedom. Typically, one finds that analytic, continuum calculations for quantities of physical interest are rendered impossible by the presence of an OE, with perturbation expansions remaining the only straightforward method of approach. For strong coupling (SC) problems, however, that avenue is blocked; with the exception of machine estimates written for specific problems, it has not been possible to proceed in any systematic way.

The goal of the present approach, for which this paper is, hopefully, a useful first step, is to exhibit for arbitrary  $N$  the functional dependence of an OE on its input dependence, in the sense that the expected, rapid fluctuations of the output in the stochastic regime are suppressed, and only the "average," or relatively slowly varying, amplitudes are reproduced. (This slowly varying dependence is, however, quite dependent upon the frequency of the rapid fluctuations.) Certain technical limitations of the present work restrict the validity of this analysis, rendering incomplete the "fine structure" (FS) analysis of the previous papers; nevertheless, the construction given below should provide a description of the leading behavior of an OE, as a functional of the defining input dependence, in the stochastic or rapidly fluctuating limit.

The general OE of interest may be described as the solution to the differential equation

$$\frac{\partial U}{\partial t}(t;E) = i\lambda_a E_a(t) U(t;E), \quad (1.1)$$

with the initial condition  $U(0;E) = 1$ ; that is,

$$U(t;E) = \left( \exp \left( i \int_0^t dt' \lambda \cdot E(t') \right) \right)_+, \quad (1.2)$$

with  $\lambda_a$  the  $N^2 - 1$  independent, Hermitian, fundamental representation matrices (the Gell-Mann matrices) of  $SU(N)$ , satisfying the properties

$$\text{tr}[\lambda_a] = 0, \quad \text{tr}[\lambda_a \lambda_b] = 2\delta_{ab},$$

$$\frac{1}{2} [\lambda_a, \lambda_b] = if_{abc} \lambda_c,$$

$$\frac{1}{2} \{\lambda_a, \lambda_b\} = (2/N)\delta_{ab} + d_{abc} \lambda_c.$$

The arguments of the OE  $U(t;E)$  have been written to emphasize its functional dependence on the input vector  $E_a(t')$ , as well as on the specific  $t$  dependence of its output form. Frequently, the OE's which appear in various problems are intended to be used as part of the integrand of a subsequent functional integration, weighted with an appropriate probability distribution. For example, if  $E_a(t) = g^* v_\mu(t) \lambda_a^* A_\mu^a(x - \int_0^t dt' v(t'))$ , where  $A_\mu^a$  is a gauge field and the four-velocity  $v_\mu(t)$  is represented by  $dx_\mu(t)/dt$ , the trace of the corresponding OE will define—before quantum fluctuations are attempted—a QCD Wilson loop.<sup>3</sup> If, on the other hand,  $E_a$  is proportional to the symmetric, spatial gradients of a fluid velocity, one can formulate<sup>4</sup> the "vortex stretching" term of three-dimensional fluid dynamics. When the probability weightings used are, or are close to, white-noise Gaussian, one may expect very rapid fluctuations of the input  $E$  vectors; and this behavior has previously been termed "stochastic." Because this paper deals only with the "deterministic" behavior of  $U(t;E)$  under rapid, nonrandom fluctuations of the unit vector  $\hat{E} = E/E$ ,  $E = (E^2)^{1/2}$ , we will henceforth use the phrase "rapidly fluctuating input" (RFI) to denote that portion of the SC regime studied here.

The weak coupling regime, defined by  $\int_0^t dt' E(t') \ll 1$ , is essentially perturbative and poses no problem. The SC regime, for which  $\int_0^t dt' E(t') > 1$ , has two natural and opposite limits defined most simply in terms of the dimensionless ratio  $\rho = |dE/dt|/E$ . For  $\rho = 0$ , that is for a constant unit vector  $\hat{E}$ , the OE of (1.2) reduces to an ordinary exponential, if one merely rotates the coordinate axes such that any one of them points in the  $\hat{E}$  direction. Hence, for  $\rho \ll 1$ , one may expect a set of "adiabatic" approximations to  $U$ , defined as corrections to the  $\rho = 0$  limit; this has been discussed in some detail for  $SU(2)$  in Ref. 1, and can be extended to  $SU(N)$ , although it will not be treated here. The other, opposite limit  $\rho \gg 1$  defines the much more interesting RFI situation, and is the subject of this paper.

As in the  $SU(2)$  case, one finds that the complex coefficient functions  $F_0, F_a$  in the representation

$$U = F_0 + \sum_a \lambda_a F_a \quad (1.3)$$

are given as rapid fluctuations superimposed on a slowly

<sup>a)</sup> Permanent address: Physics Department, Brown University, Providence, Rhode Island 02912.

<sup>b)</sup> Unité Associée 767 au Centre National de la Recherche Scientifique, Physique Théorique, Parc Valrose, 06034 Nice Cedex, France.

varying, or averaged background; and the time variations of the latter depend in a nontrivial way on the rapid fluctuations of the input. In contrast to the SU(2) case, where the four coefficient functions  $F_0, -i^*F_a$  are real, the corresponding functions of SU( $N$ ) are complex, and one must pay special attention to the unitarity property of the  $U$ , following from (1.1) or (1.2),

$$U^\dagger * U = U * U^\dagger = 1, \quad (1.4)$$

which places certain restrictions on the coefficient functions.

Because of the need to maintain unitarity, the theoretical work in this paper will make use of the representation

$$U = \exp\left(i \sum_a \lambda_a G_a\right), \quad (1.5)$$

with  $\hat{G} = \mathbf{G}/G$ ,  $G = (G^2)^{1/2}$ , and  $G_a$  a real vector with  $L = N^2 - 1$  components. Approximate forms will be found for  $G$  and  $\hat{G}$  in the RFI limit, and a simple conversion made to the coefficient functions of (1.3). For that conversion, as well as for use in other steps of the analysis, a version of an ancient representation due to Lagrange and Sylvester<sup>5</sup> will be used, and will be denoted by the phrase "normal form,"

$$\mathcal{F}(\lambda \cdot \mathbf{v}) = \sum_T \mathcal{F}(\xi_l) \left[ \frac{1}{N} + \frac{1}{2} \lambda_a \frac{\partial \xi_l}{\partial v_a} \right], \quad (1.6)$$

where the  $\xi_l$  denote the  $N$  eigenvalues of  $\lambda \cdot \mathbf{v}$ , with  $v_a$  an arbitrary,  $L = N^2 - 1$  component vector. The derivation of (1.6) is elementary, holds for any  $\mathcal{F}(z)$  whose Fourier transform exists, and has been relegated to the Appendix.

Because it is appropriate to compare the results of the RFI approximation with the "correct," numerically integrated coefficient functions, the latter have been computed using an algorithm<sup>6</sup> that guarantees unitarity. There, for a fixed time interval  $\Delta t$  one first replaces (1.1) by the difference equation,

$$U(t + \Delta t) = U(t) + i\lambda \cdot \mathbf{E}(t) \Delta t U(t); \quad (1.7)$$

and then, in order to insure the unitarity of the numerically integrated solutions, one replaces (1.7) by the equation

$$U(t + \Delta t) = U(t) + (i/2)\lambda \cdot \mathbf{E}(t) \Delta t [U(t + \Delta t) + U(t)], \quad (1.8)$$

which is the same as (1.7) and (1.1) only in the limit  $\Delta t = 0$ . The "solution" to (1.8) may be written as

$$U(t + \Delta t) = \left( \frac{1 + (i/2)\lambda \cdot \mathbf{E}(t) \Delta t}{1 - (i/2)\lambda \cdot \mathbf{E}(t) \Delta t} \right) U(t), \quad (1.9)$$

and is rigorously unitary, by inspection; that is, if  $U(t)$  is unitary, then  $U(t + \Delta t)$  is also. Equation (1.9) has been used to compute the "exact" amplitudes with which our approximate forms are compared.

As in Ref. 1, the essential idea is to extract, in the large- $\rho$  limit, the slowly varying envelope upon which the rapid fluctuations ride, for in any physical context it is surely only the slowly varying behavior that is of interest. Simple (unitarity) considerations would, as for the SU(2) case, suggest the improvement of the leading, "averaged" dependence, by including FS corrections expressible in powers of  $1/\rho$ . However, for general SU( $N$ ), and in particular as  $N$  increases, there will appear a certain difficulty in calculating FS terms, and for reasons quite distinct from the "technical" difficul-

ties alluded to above. Very simply, most of the coefficient functionals  $F_a$  will be expressed as increasingly nonlinear functions of the  $\hat{G}_b$ ; and if one is interested in the slowly varying, or averaged behavior of the  $F_a$ , one will have to perform averages over the corresponding combinations of the  $\hat{G}_b$ . Our analysis, however, is really simple only in its extraction of average values of the  $\hat{G}_b$ ; and since the average of products is not equal to the product of averages—in particular, some slowly varying, FS dependence is always missed—one may expect that, for certain  $F_a$ , the FS dependence estimated in the manner of Ref. 1 is bound to be incomplete. For this reason, as well as for the "technical" reasons described below, the results of this paper are restricted to the leading-order estimates of  $F_0$  and the  $F_a$ , which are constructed from those parts of the averaged  $\hat{G}_b$  independent of  $\rho$  (and, as in Ref. 1, from the  $\rho$ -dependent  $G$ ).

Finally, a word must be said about the level of rigor—or its absence—in this paper. When the analysis becomes too difficult to permit a brute force extraction of a desired result, an appeal will be made to intuition, to an argument labeled by the phrase "...what else can it be?" And it may well be that some subtlety not foreseen by the author will answer this question in a manner different from that found here. On the basis of the numerical comparisons noted below this would not appear too likely; but it must be understood that, without honest mathematical rigor, such possibilities exist. Nevertheless, the predictions made in this paper are, at the very least, interesting; and they may even be true.

## II. ESTIMATING $\hat{G}$

### A. The differential equation

One begins with the representation (1.5), differentiating  $U$  and substituting into the original Eq. (1.1), to obtain the nonlinear relation

$$\lambda \cdot \mathbf{E} = \int_0^1 d\mu e^{i\mu\lambda \cdot \mathbf{G}} \lambda \cdot \frac{d\mathbf{G}}{dt} e^{-i\mu\lambda \cdot \mathbf{G}},$$

or

$$E_a = \frac{1}{2} \int_0^1 d\mu \operatorname{tr} \left[ \lambda_a e^{i\mu\lambda \cdot \mathbf{G}} \lambda \cdot \frac{d\mathbf{G}}{dt} e^{-i\mu\lambda \cdot \mathbf{G}} \right]. \quad (2.1)$$

In all of the following we will assume that the time variation of the components  $\hat{E}_a$  follows the rule

$$\frac{d\hat{E}_a}{dt} = f_{abc} \rho_b E_c, \quad (2.2)$$

where the quantities  $\rho_b$  are constants. It then follows that  $\mathbf{p}$  and  $\hat{E}$  are orthogonal, and that the norm of  $\hat{E}$  is unchanged. For simplicity, we will also assume that the magnitude  $E$  does not depend on time. From the experience of Ref. 1, validated by the SU(3) example of the next section, one may expect that the form of our leading dependence will be essentially unchanged even if  $\mathbf{p}$  and  $\mathbf{E}$  are time dependent, as long as the condition  $\rho \gg 1$  is maintained.

Separating  $\mathbf{G}$  into the product  $\hat{G}G$ , and calculating the trace of all the multicommutators of (2.1) (which are obtained by expanding and resumming in powers of  $\mu$ ), one builds

$$E_a = \hat{G}_a(\mathbf{E} \cdot \hat{G}) + \frac{i}{2} \sum_b \frac{d\hat{G}_b}{dt} \left[ (1 - e^{2i\mathbf{A} \cdot \hat{G}G}) \cdot \frac{1}{\mathbf{A} \cdot \hat{G}} \right]_{ba}, \quad (2.3)$$

where the  $(A^i)_{jk} = i f_{ijk}$  form the adjoint representation matrices of  $SU(N)$ . Solving for  $d\hat{G}/dt$ , with the aid of the relation

$$(1 - e^{2iz})^{-1} = \frac{1}{2} [1 + i \cot(z)],$$

one can rewrite (2.3) in the form

$$\frac{d\hat{G}_a}{dt} = f_{abc} \hat{G}_b E_c + [\mathbf{A} \cdot \hat{G} \cot(\mathbf{A} \cdot \hat{G}G)]_{ab} (E_b - \hat{G}_b(\mathbf{E} \cdot \hat{G})). \quad (2.4)$$

As in the  $SU(2)$  case, the magnitude  $G$  is given in terms of  $\hat{G}$  in the following way. Multiplication of (2.1) by  $\Sigma_a \hat{G}_a$  yields the relation

$$\hat{G} \cdot \mathbf{E} = \frac{dG}{dt},$$

which, together with the initial condition  $G(0) = 1$ , provides

$$G(t) = \int_0^t dt' \mathbf{E}(t') \cdot \hat{G}(t'), \quad (2.5)$$

showing that  $G$  is completely specified by knowledge of  $\hat{G}$ . It will be convenient to introduce a new variable,  $\tau$ , defined by  $d\tau = dt * E(t)$ ,  $\tau = \int_0^t dt' * E(t')$ , in order to rewrite (2.5) and (2.4) as

$$G(\tau) = \int_0^\tau d\tau' \hat{E}(\tau') \cdot \hat{G}(\tau') \quad (2.6)$$

and

$$\frac{d\hat{G}_a}{d\tau} = -f_{abc} \hat{E}_b \hat{G}_c + [\mathbf{A} \cdot \hat{G} \cot(\mathbf{A} \cdot \hat{G}G)]_{ab} (\hat{E}_b - \hat{G}_b(\hat{E} \cdot \hat{G})). \quad (2.7)$$

For  $SU(2)$ , where the  $f_{abc} = \epsilon_{abc}$ , it is easy to see that  $[(\mathbf{A} \cdot \hat{G})^2]_{ab} = \delta_{ab} - \hat{G}_a \hat{G}_b$ ,

and that (2.7) reduces to

$$\frac{d\hat{G}}{d\tau} = -\hat{E} \times \hat{G} + (\hat{E} - \hat{G}(\hat{E} \cdot \hat{G})) \cot(G), \quad (2.9)$$

as quoted in Ref. 2. For  $SU(N)$ , with  $N > 2$ , there is no reason to believe that the equation corresponding to (2.9) will depend only on the "group invariant"  $G^2$ , for there are other invariants, depending on  $N$ , such as  $d_{abc} G_a G_b G_c$ , which could appear in the corresponding equation. It is here that the first stumbling block to a straightforward analysis arises, for the averaging method to be used requires a certain closure property, which is difficult to see directly from (2.6). What we shall do, instead, is to consider first a model problem, in which one pretends that the internal structure of (2.7) is analogous to that of  $SU(2)$ . Then, on the basis of that model result, we guess the form of the more realistic  $\hat{G}$  dependence in the large- $\rho$  limit. Finally, with the aid of the exact (2.7), we can estimate the size of the coefficients of the various terms, and find, for large  $\rho$  at any finite  $\tau$ , that the leading  $\hat{G}$  dependence is just that of the model calculation.

Needless to say, one would prefer a more straightforward argument; but indirect as the present estimates are, they do agree with numerical computations in at least one simple, but decidedly nontrivial case of  $SU(3)$ , as shown in the following section.

## B. A simplified model

For  $N > 2$ , the relation generalizing (2.8) can be written in various ways, but each of them apparently too complicated to permit the estimation of multiple products of  $(\mathbf{A} \cdot \hat{G})^2$ , leading to a closed form for the quantity  $(\mathbf{A} \cdot \hat{G}) \cot(\mathbf{A} \cdot \hat{G})$ . To obtain a simple form for the latter, we introduce the model approximation

$$(\mathbf{A} \cdot \hat{G})_{ab}^2 = P \delta_{ab} + Q \hat{G}_a \hat{G}_b, \quad (2.10)$$

which, effectively, pretends that the internal group structure here is that of  $SU(2)$ . There are other terms which could possibly be added to the rhs of (2.10), such as  $d_{asr} \hat{G}_r * d_{bst} \hat{G}_t$ , and  $d_{abr} * d_{rst} \hat{G}_s \hat{G}_t$ , which can appear automatically for  $N > 2$ . A compact method of writing the exact rhs of (2.8) uses the "normal form" to express the combination in terms of Gell-Mann matrices and the corresponding  $d_{abc}$  of  $SU(L)$ , with  $L = N^2 - 1$ , and is noted in the Appendix.

With the model approximation of (2.10), the only reflection of arbitrary  $N$  lies in the coefficients  $P$  and  $Q$ , which are determined by insisting that the combination continue to be orthogonal to either  $\hat{G}_a$  or  $\hat{G}_b$ ; and that the normalization be appropriate to  $\text{tr}((\mathbf{A} \cdot \hat{G})^2)$ . In this way, one finds that  $-Q = P = C_A(N)/(L - 1) = N/(N^2 - 2)$ , where  $C_A$  is the eigenvalue of the first Casimir operator in the adjoint representation. Repeated products of (2.10) can then be summed to yield a differential equation for  $\hat{G}$ ,

$$\frac{d\hat{G}_a}{d\tau} = -f_{abc} \hat{E}_b \hat{G}_c + \sqrt{P} \cot(\sqrt{P}G) (\hat{E}_a - \hat{G}_a(\hat{E} \cdot \hat{G})). \quad (2.11)$$

With the model (2.11), it will be possible to formulate an "averaged" approximation to the coefficient functions of  $U(t;E)$ ; the technique to be used resembles that of the  $SU(2)$  calculation, but is different in detail because of the need to work with a relatively large number of  $\hat{G}_a$  components, growing as  $N^2$  for large  $N$ . If an equation, similar to (2.11), could be obtained without recourse to the model approximation of (2.10), one would be in somewhat better shape as far as the FS is concerned; for the leading dependence of  $\hat{G}$  in the large- $\rho$  limit, however, it will be argued below that the model estimates are sufficient.

We now form the following three quantities,  $I = \Sigma_{abc} \hat{G}_a \rho_b \hat{E}_c$ ,  $J = \Sigma_a \hat{E}_a \hat{G}_a$ , and  $K = \Sigma_a \rho_a \hat{G}_a$ , and ask for the equations satisfied by these objects in an averaged sense, retaining only bilinear products of rapidly oscillating dependence, such as  $\hat{E}_a \hat{E}_b$  and  $J$  itself. As in the  $SU(2)$  case, a cursory examination of any numerical output suggests that  $J$  is approximately a constant, of order  $1/\rho$ , which means that  $\hat{G}$  will have small but rapidly oscillating components in phase with those of  $\hat{E}$ . If it turns out that  $K$  has a slowly varying part, then that will also be true of  $\hat{G}$ .

With the aid of (2.11), one builds the relations

$$\frac{dI}{d\tau} = \frac{N}{L} (K - \rho^2 J) - JI\sqrt{P} \cot(\sqrt{P}G), \quad (2.12)$$

$$\frac{dJ}{d\tau} = I + \sqrt{P} \cos(\sqrt{P}G)(1 - J^2), \quad (2.13)$$

$$\frac{dK}{d\tau} = -I - \sqrt{P}JK \cot(\sqrt{P}G), \quad (2.14)$$

where the last two of these equations follow from (2.11) without approximation, while the first has used the averaging replacements

$$\langle \hat{E}_a \hat{E}_b \rangle = \delta_{ab}/L, \quad (2.15)$$

and

$$\langle \hat{E}_a \hat{G}_b \rangle = \delta_{ab} \langle J \rangle / L + f_{abc} \rho_c \langle I \rangle / N \rho^2. \quad (2.16)$$

In writing (2.15) and (2.16), the tacit assumption has been made that there are in fact  $L = N^2 - 1$  rapidly oscillating components  $\hat{E}_a$ ; and that their time average over rapid fluctuations is of the same form as is their stochastic average. For other situations, such as the SU(3) example of the next section, the normalization factor ( $L/N$ ) of (2.12) will be changed; but this change will effect only the FS dependence. Equation (2.16) represents a statement of internal consistency, which is compatible with the definitions of both  $I$  and  $J$ .

As in the SU(2) case, we use the "experimental" fact that  $J$  can be considered to have a small constant, or averaged value, on which is superimposed rapid fluctuations. Neglecting, or averaging over such fluctuations, one then sets  $dJ/dt = 0$  to obtain from (2.13)

$$I \sim -\sqrt{P}(1 - J^2) \cot(\sqrt{P}G), \quad (2.17)$$

whose derivative then yields

$$\frac{dI}{d\tau} \sim -\sqrt{P}(1 - J^2)(-J\sqrt{P}) \cdot \frac{1}{\sin^2(\sqrt{P}G)}. \quad (2.18)$$

Comparing (2.18) with (2.12), one finds that  $K$  must, on the average, be considered as a constant, of value

$$K = J\rho^2 + (L/N)PJ(1 - J). \quad (2.19)$$

Since the lhs of (2.14) is to vanish, one then concludes that

$$I = -JK\sqrt{P} \cot(\sqrt{P}G). \quad (2.20)$$

Finally, a comparison of (2.20) with (2.17) yields

$$JK = (1 - J^2), \quad (2.21)$$

which, together with (2.19), serves to determine  $J$  as a function of  $\rho$ ,

$$(L/N)PJ(1 - J^2) + J\rho^2 = (1 - J^2)/J,$$

or

$$J^2 = \frac{(1 + (L/N)P + \rho^2)}{2P(L/N)} \times \left[ 1 - \left( 1 - \left[ \frac{2\sqrt{PL/N}}{(1 + (L/N)P + \rho^2)} \right]^2 \right)^{1/2} \right]. \quad (2.22)$$

The negative sign of the square root has been chosen so that  $J$  tends to zero as  $\rho$  increases. For large  $\rho$ , one has

$$J^2 \sim (1 + (L/N)P + \rho^2)^{-1} + \dots, \quad (2.23)$$

that is,  $J \sim 1/\rho$ . Factors depending on  $N$  and  $L$  appear only in higher-order corrections to the leading RPI behavior, and

even though this is just a model computation, we argue below that the same feature will be true in an exact rendering of the leading, large- $\rho$  behavior.

From (2.21) and (2.23) we now infer that  $K$  should be assigned an average value,  $K = (1 - J^2)/J$ , and from the definition of  $K$  one concludes that, to leading order, there must be assigned an averaged, nonzero value:  $\hat{G}_a = \hat{\rho}_a$ . Our model solution, denoted by the quantity  $U$ , is then

$$U = e^{i\lambda_a \hat{\rho}_a G}, \quad G \sim \int_0^t dt' \frac{E}{\rho}. \quad (2.24)$$

This result is equivalent to that of the previous (and somewhat different) SU(2) calculation in the limit of very large  $\rho$ . Since the  $\rho_a$  are specified input constants, one can proceed from (2.24) to the coefficient functions  $F_0, F_a$  by the use of the normal form, (1.6). If the differences between products of the averages of  $\hat{G}$  components and the averages of the same products is a FS effect, as is true in all the SU(2) work and in the SU(3) example of the next section, then the leading, large- $\rho$  behavior of the coefficient functions is obtained from (2.24).

### C. Guessing the solution

How could the use of the exact (2.7) change the model result (2.24) that followed from (2.11)? The only (reasonable) difference would be that  $\hat{G}_a$  would have a constant (averaged over the rapid fluctuations) part that now depends on unit vectors  $\hat{W}_a^{(i)}(\rho)$  more general and more complicated than the  $\hat{W}_a^{(1)} = \hat{\rho}_a$  alone. Such unit vectors, orthogonal to the  $\rho_a$ , can be constructed out of the gradients of the  $N - 1$  independent invariants  $\text{tr}([\lambda \cdot \hat{\rho}]^i)$ ; e.g., if  $d_{abc} \hat{\rho}_b \hat{\rho}_c = \frac{1}{3}(d/d\hat{\rho}_a)(d_{rst} \hat{\rho}_r \hat{\rho}_s \hat{\rho}_t)$ , then one choice for  $\hat{W}_a^{(2)}$  is

$$(d_{abc} \hat{\rho}_b \hat{\rho}_c - \hat{\rho}_a d_{\alpha\beta\gamma} \hat{\rho}_\alpha \hat{\rho}_\beta \hat{\rho}_\gamma) \cdot \left( \sum d_{abc}^2 \hat{\rho}_b \hat{\rho}_c - [d_{rst} \hat{\rho}_r \hat{\rho}_s \hat{\rho}_t]^2 \right)^{-1}.$$

One would then expect to find a representation for the averaged  $\hat{G}$  in the form of a sum of such terms,

$$\hat{G}_a = \sum_i c_i \hat{W}_a^{(i)}(\hat{\rho}), \quad (2.25)$$

where the coefficients  $c_i$  remain to be determined, but where the overall normalization is chosen to satisfy  $\hat{G}^2 = 1$ . Under a set of reasonable assumptions, we now obtain the leading behavior of these coefficients.

To see how this goes, return to the exact equation (2.7) for  $\hat{G}$ , and calculate the time variation of the same  $I, J, K$  quantities as before. One can no longer carry through the process of writing the quantity

$$Q_{ab} = [A \cdot \hat{G} \cot(A \cdot \hat{G} G)]_{ab}$$

in closed form; but since  $Q_{ab}$  will only be used to estimate leading orders of magnitude when averaged, or contracted with an  $\hat{E}_c$ , this will not be too important. Based on the SU(2) work, on the model of the previous section, and on the SU(3) example of the next section, we now introduce two assumptions which form the basis of the argument to follow.

(1) Contraction with any  $\hat{E}_c(t)$  produces for an arbitrary  $\mathcal{F}(\hat{G})$  a result of one nominal order lower in  $\rho$ ,

$$\langle \hat{E}_a \mathcal{F}(\hat{G}) \rangle = \langle \hat{E}_a \hat{G}_c \rangle \frac{\delta}{\delta \hat{G}_c} \mathcal{F}(\hat{G})|_{\hat{G}=\hat{G}^{(0)}}, \quad (2.26)$$

where  $\hat{G}^{(0)}$  denotes the constant dependence of  $\hat{G}$  of  $O(1)$ , after averaging over the rapid fluctuations. Again, we adopt the internally consistent relation (2.16), and assume an averaging over all  $L = N^2 - 1$  components  $\hat{E}_a$ . By this assumption,  $\langle I \rangle \sim O(1)$ ,  $\langle J \rangle \sim O(1/\rho)$ , and  $\langle K \rangle \sim O(\rho)$ .

(2) The time average  $\langle q \rangle$  over rapid fluctuations of any quantity  $q(\tau)$  may still depend on a slowly varying  $\tau$  dependence; and  $\langle dq/d\tau \rangle = d\langle q \rangle/d\tau$ . All slowly varying dependence—such as that of  $G = \int d\tau' J$ —will be treated as constants during the averaging over rapid fluctuations. [Were this not the case, then some of the  $O(1/\rho)$  could conceivably be associated with slow variation of the  $G$  dependence.]

Since no approximation is now considered for the  $Q_{ab}$ , and since  $\hat{G}_a Q_{ab} = Q_{ab} \hat{G} = 0$  by virtue of the definition of  $A_{ab}$ , we can drop the corresponding term on the rhs of the exact (2.7). The counterparts of Eqs. (2.12) through (2.14) then become

$$\frac{dI}{d\tau} = \frac{N}{L} (K - J\rho^2) + O\left(\frac{1}{\rho^2}\right), \quad (2.27)$$

$$\frac{dJ}{d\tau} = I - \frac{1}{L} \text{tr} Q, \quad (2.28)$$

$$\frac{dK}{d\tau} = -I + \rho_a \langle \hat{E}_b \hat{G}_c \rangle \frac{\partial}{\partial \hat{G}_c} Q_{ab}|_{\hat{G}=\hat{G}^{(0)}}, \quad (2.29)$$

where the averaging symbols  $\langle \rangle$  have been suppressed, as in Sec. II B, but where an averaging has been performed to obtain each of these three equations. If, again, we call upon the “experimental” knowledge that  $J$  is a constant, in the large- $\rho$  limit, we may infer from (2.28) that

$$I \approx (1/L) \text{tr} Q, \quad (2.30)$$

and then, from (2.27), that

$$K \approx J\rho^2 + \frac{1}{N} \frac{d}{d\tau} \text{tr} [ \ ]. \quad (2.31)$$

Substituting into (2.29), and comparing with (2.30), one can see that  $\text{tr} Q(\tau)$  satisfies a forced harmonic oscillator equation, with frequency  $(1/2\pi) * (N/L)^{1/2}$  and driving term given by  $\rho_a \langle \hat{E}_b \hat{G}_c \rangle (\partial/\partial \hat{G}_c) Q_{ab}|_{\hat{G}^{(0)}}$ . More detailed information is difficult to extract; but for our purposes essentially all that is needed is (2.31), which says that to leading order,  $K = J\rho^2$ . In terms of  $\hat{G}$ , this means that  $\sum_a \hat{\rho}_a \hat{G}_a = \rho J$ , and one can conclude that  $\hat{G}$  is given by a series of terms of the form of (2.25), of which the coefficient multiplying  $\hat{\rho}$ , say  $c_1$ , is given by  $\rho * J$ . In the previous model calculation, we knew that  $c_1 = 1$ , but that information is missing here. In general,  $\hat{G}$  may very well have projections in (group) directions orthogonal to  $\hat{\rho}$ , as explained just before (2.25).

Suppose that  $\hat{G}$  does have a projection in the direction of a unit vector  $\hat{v}(\hat{\rho})$  perpendicular to  $\hat{\rho}$ . Analogous to the quantities  $K$  and  $I$ , one may define the sums  $M = \sum_a v_a \hat{G}_a$  and  $P = \sum_a f_{abc} \hat{G}_a v_b \hat{E}_c$ , where  $M$  measures the projection of  $\hat{G}$  in the  $\hat{v}$  direction, and the averaging of  $P$  must generate a term of  $O(1/\rho)$ . Following the same averaging techniques, one then finds that

$$\frac{dM}{d\tau} = -P + v_a \langle \hat{E}_b \hat{G}_c \rangle \frac{\partial}{\partial \hat{G}_c} Q_{ab}|_{\hat{G}^{(0)}} \quad (2.32)$$

and

$$\frac{dP}{d\tau} = -\frac{N}{L} J \mathbf{v} \cdot \mathbf{p} + \frac{N}{L} M, \quad (2.33)$$

where the first term on the rhs of (2.33) has been included, for clarity, even if the value of this term is zero; another pair of terms, not shown, have vanished by symmetry. That first term is zero because  $\mathbf{v}$  has been chosen orthogonal to  $\mathbf{p}$ ; but if it did not vanish it would contribute a term of order unity to (2.33). Because all rapid fluctuations of  $P$  have been averaged away, so that  $dP/d\tau$  is at least of  $O(1/\rho)$ , the inescapable conclusion is that  $M$  must also be of  $O(1/\rho)$ . If this arbitrarily chosen projection of  $\hat{G}$  is  $O(1/\rho)$ , then to  $O(1)$ , the properly normalized  $\hat{G}$  can only be given by  $\hat{G} = \hat{\rho}$ ; then,  $c_1 = 1$ , and  $J \sim 1/\rho$ .

By this construction, the model result (2.24) is correct for the non-FS terms of  $O(1)$  in the coefficient functions  $F_0, F_a$  built from (2.24) with the aid of the normal form representation. This statement will next be verified with the aid of a simple but nontrivial  $SU(3)$  calculation. It must be remembered, however, that in the passage from the  $\hat{G}$  description to that of the coefficient functions, we are going to assume that the difference between the average of a product of  $\hat{G}$  components and the product of the corresponding averages is of higher order in  $1/\rho$ . In fact, in the discussion of Sec. II B, we have cavalierly ignored the difference between the average of products of the  $I, J, K$  and the corresponding products of the averages; this is obviously wrong, since it is easy to see that

$$\langle KJ \rangle \sim (1 + 1/L) \langle K \rangle \langle J \rangle,$$

$$\langle IJ \rangle \sim (1 + 2/L) \langle I \rangle \langle J \rangle, \quad \langle J^2 \rangle \sim \langle J \rangle^2 + 1/L.$$

In spite of these omissions, the  $O(1)$  results of the model are unchanged. For the  $\hat{G}_a$ , the neglect of such correlations does seem to be a viable assumption, as is born out by the calculations (and figures) of the next section.

### III. AN $SU(3)$ EXAMPLE

The simplest, nontrivial example of an  $SU(3)$ , OE has but a single component of  $\rho$ , driving two components of  $\hat{E}_a$  each in a different “sector” of  $SU(3)$ . That is, we suppose that the only nonzero components are  $\hat{E}_2 = \cos(\omega t)$  and  $\hat{E}_5 = \sin(\omega t)$ , so that  $\rho_a = \delta_{a7} \cdot \rho$  with  $\rho = 2 * (\omega/E)$ . Hence

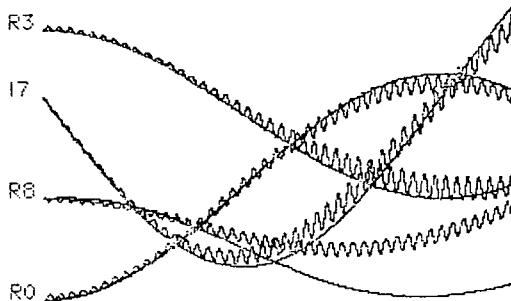


FIG. 1. A superposition of  $O(1)$  averaged approximations and four numerically integrated functions, for  $E = 10$  and  $\omega = 60$ .

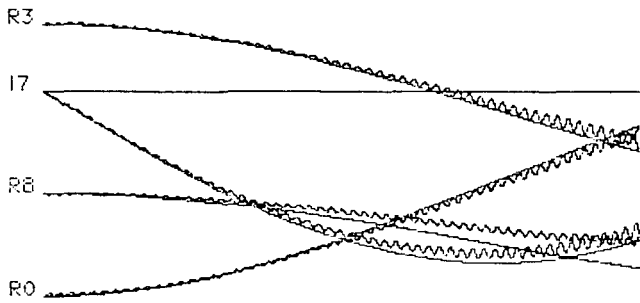


FIG. 2. The same as in Fig. 1, except that  $\omega = 90$ .

the leading  $\hat{G}_a^{(0)} = \delta_{a7}$ , and the nonzero coefficient functions  $F_0, F_a$  are easily worked out to be

$$F_0 = \frac{1}{3}[1 + 2 \cos G], \quad (3.1)$$

and

$$\begin{aligned} F_7 &= i \sin G, & F_3 &= \frac{1}{2}[1 - \cos G], \\ F_8 &= (1/2\sqrt{3})[1 - \cos G]. \end{aligned} \quad (3.2)$$

A superposition of these four predictions with those of the numerically integrated solutions is given in Figs. 1 and 2, in the simplest possible situation of constant  $\omega$  and  $E$ , for two different values of  $\rho$ . (Here and subsequently we use the notation  $R_a$  and  $I_a$  to specify, for this special input, the nonzero real and imaginary parts, respectively, of the  $F_a$ .) Only for  $R_8$  is the agreement less than satisfactory; and this will be discussed separately, below.

At larger  $t$  values, however, there is another source of error in the graphical presentation of these results, as the differences between the numerical integrations and the true solutions become evident. This is to be expected since (1.8) and (1.1) become the same only in the limit  $\Delta t \rightarrow 0$ ; for finite  $\Delta t$  the numerically integrated functions start to separate from the exact solutions, and this separation becomes greater the larger the value of  $t$ . In Figs. 3 and 4, produced for the same  $t$ -dependent  $\rho(t)$  but for different values of  $\Delta t$ , this " $\Delta t$ " effect can clearly be seen. One infers that, in the limit of very small  $\Delta t$ , the smoothly varying part of the exact functions are well described by (3.1) and (3.2), even when  $\rho$  depends upon  $t$ ; the only requirement is that  $\rho \gg 1$ .

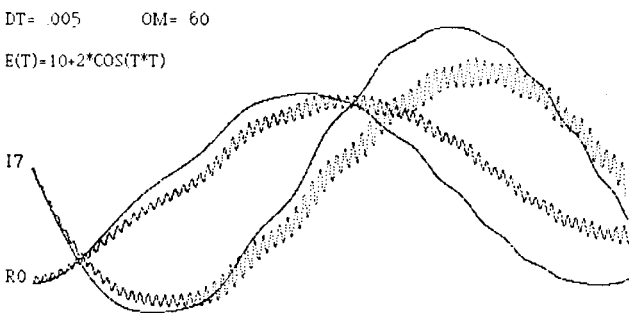


FIG. 3. A superposition of approximate and numerically integrated functions for  $R_0$  and  $I_7$ , with  $\omega = 60$  and  $E(t) = 10 + 2^* \cos(t^2)$ , for  $\Delta t = 0.005$ .

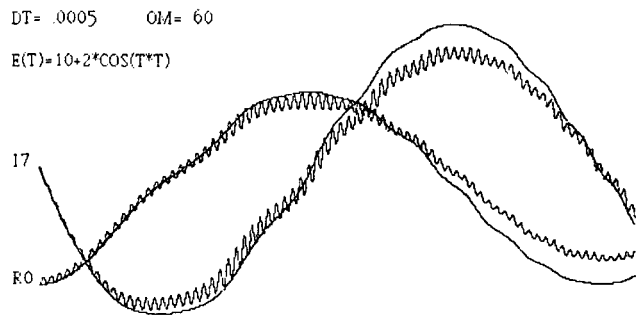


FIG. 4. The same as in Fig. 3, except that  $\Delta t = 0.0005$ .

Finally, one comes to the remaining five functions whose  $O(1)$  dependence is predicted to be zero by this analysis,  $I_2, I_5, R_1, R_4$ , and  $R_6$ . Figure 5 clearly shows that these functions decrease at least as fast as  $1/\rho$ ; and it also shows the existence of higher harmonics present in those curves, which becomes especially clear for larger values of  $\rho$ .

Although no statement of FS has been attempted in this paper, it is perhaps worth remarking that an attempt to write down appropriate corrections to all of these functions can be organized by using the model forms of Sec. II B, for this simplest case of  $SU(3)$ . This is not completely correct for reasons mentioned above: neglect of the  $O(1/\rho)$  terms due to the difference of averages of products and products of averages, and due to the uncertainties of the model itself. But it is at least interesting to see how closely such FS can approximate the structure of the numerically integrated functions. For a set of  $O(1/\rho)$  corrections to the  $\hat{G}_a$ , computed in a straightforward way and converted to the coefficient functions, this is shown in Fig. 6. Clearly, one is on the right track, even if other, neglected effects become important at later times.

#### IV. SUMMARY AND CRITIQUE

A partial, first step towards the SC, RFI approximation of  $SU(N)$  OE's has been suggested in this paper, in which the leading terms of that OE are specified. A comparison with the numerically integrated functions of a simple but nontrivial  $SU(3)$  example illustrates the sort of agreement that may be expected.

While interesting, and of probable future use in estimating functional integrals over distributions close to white-noise Gaussian, the analysis is incomplete in that it is unable to extract the  $O(1/\rho)$  corrections. The method of analysis

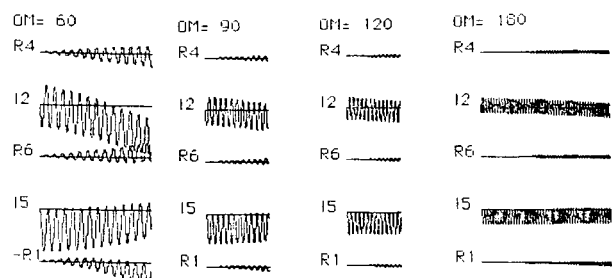


FIG. 5. Numerically integrated functions which are (at least) of  $O(1/\rho)$  illustrated for  $E = 10$  and four choices of  $\rho$ .

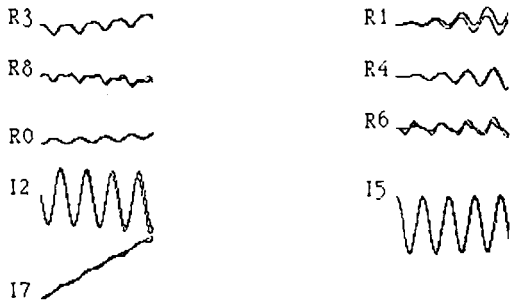


FIG. 6. An example of FS for the set of nine functions, in the SU(3) example described in the text.

also leaves unproved an assertion on the sources of  $O(1/\rho)$  dependence. It will therefore be worthwhile to list those inadequacies of this treatment; if such difficulties could be overcome, one would have a much more satisfying approach to this entire problem.

(1) The difference between an average of products of the components  $\hat{G}_a$  and the product of corresponding averages has been treated as an effect of relative  $O(1/\rho)$ . This could be proved if one were sure that slowly varying  $G(\tau)$  dependence does not contribute  $O(1/\rho)$  terms which have here been associated with the  $\hat{G}_a$  dependence. In this paper, averages over rapid fluctuations have been taken under the assumption that the  $G(\tau)$  factors may be treated as constants, for which case the  $O(1/\rho)$  behavior is as stated. With one possible exception, the SU(3) example seems to bear this out; but no proof has been given.

(2) The averaging process used has been the simplest possible, with every contraction of an  $\hat{E}$  component with a  $\hat{G}$  component assumed to be of  $O(1/\rho)$ , and all noncontracted  $\hat{G}$  components replaced by their  $O(1)$  average value. In effect, no higher harmonics of the RFI were considered, even though one can see them appearing in Fig. 5. Further, the averaging process depends on the number of nonzero  $\hat{E}$  components, with different normalization factors appearing in different situations.

With one exception—the discrepancy associated with  $R_8$  in Figs. 1 and 2—these assumptions appear to work quite well. Within the context of this approach, one can try to understand possible reasons for such poor accuracy. Since  $R_8$  really depends on  $\hat{G}_7^2$ , one might expect that this is an indication that the difference between  $\langle \hat{G}_7^2 \rangle$  and  $\langle \hat{G}_7 \rangle^2$  can be of  $O(1)$ . However,  $R_3$  also depends upon  $\hat{G}_7^2$ , and there is no corresponding effect there. [ $R_8$  has been numerically calculated for a variety of large  $\rho$  values, and one finds no appreciable difference between the vertical separation of the approximate and numerically integrated curves, measured at the same geometric points (e.g., at that value of  $t$  corresponding to the intersection of  $R_3$  and  $R_0$ ).] One might next ask if this is a “ $\Delta t$ ” effect, due to the inaccuracies of the numerical integration; but the answer is, again, no. The only possibility left is an error in the numerical integration, although that has also been rechecked. As it stands, the author has no explanation for the discrepancy of the  $R_8$  curves.

Many of these difficulties would disappear if a useful representation of  $(A * G)^2$  could be found, similar to that of

the model of Sec. II B. This is important because it surely holds the key to a development of the FS corrections. In the previous SU(2) analysis, it was possible to calculate the  $O(1/\rho)$  corrections, which can always be relevant to some physical situation. For example, in the problem of neutrino–antineutrino oscillations in the presence of a constant plus-time-periodic magnetic field,<sup>7</sup> the probability  $|\bar{N}|^2$  of finding an antineutrino at time  $t$ , when there was none at time  $t = 0$ , is just one of finding the FS to the OE representing the solution to a pair of coupled, first-order differential equations for the neutrino and antineutrino amplitudes. In the notation of Ref. 7, the answer to that problem is

$$|\bar{N}|^2 = \left[ \frac{\left( \frac{\omega_m}{\omega_B} \right) \sin((2\omega_B r/\omega) \sin^2(\omega t/2) - \omega_B t)}{[1 - r \sin(\omega t)]} \right]^2$$

for fixed  $\omega t$ ,  $r \neq \pm 1$ , and  $\omega_m \ll \omega_B$ . It corresponds to the square of certain  $O(1/\rho)$  terms ( $F_1^2 + F_2^2$ , in the notation of Ref. 2), where  $\rho = 2(\omega_B/\omega_m)|1 - r \sin(\omega t)|$ .

In spite of the limitations described above, it is hoped that this paper will point the way towards a systematic method for the approximation of OE's in the RFI limit.

## ACKNOWLEDGMENTS

It is a pleasure to acknowledge the kind hospitality of colleagues at Physique Théorique, Université de Nice, and at the Observatoire de Nice, during the academic year 1985–1986, where most of this work was performed.

An NSF–CNRS travel grant was useful in the preliminary stages of this project. This work was supported in part by the U. S. Department of Energy Contract No. DE-AC0276A03130.A009-Task A.

## APPENDIX: DERIVATION OF THE “NORMAL FORM”

The derivation of (1.6) may be sketched as follows. One writes an arbitrary  $\mathcal{F}(\lambda \cdot \mathbf{G})$  in terms of its Fourier transform,

$$\mathcal{F}(\lambda \cdot \mathbf{G}) = \int_{-\infty}^{+\infty} d\omega \tilde{\mathcal{F}}(\omega) e^{i\omega \lambda \cdot \mathbf{G}},$$

and considers the unitary exponential written in terms of its coefficient functions,

$$e^{i\omega \lambda \cdot \mathbf{G}} = F_0 + i\lambda \cdot \mathbf{F}.$$

Let the  $N$  eigenvalues of  $\lambda \cdot \mathbf{G}$  be denoted by  $\xi_i(\mathbf{G})$ . Then

$$F_0 = \frac{1}{N} \text{tr} [e^{i\omega \lambda \cdot \mathbf{G}}] = \frac{1}{N} \sum_{i=1}^N e^{i\omega \xi_i},$$

and

$$F_i = -\frac{i}{2} \text{tr} [\lambda_i e^{i\omega \lambda \cdot \mathbf{G}}] = -\frac{N}{2\omega} \frac{\partial}{\partial G_i} F_0.$$

Hence

$$\begin{aligned} \mathcal{F}(\lambda \cdot \mathbf{G}) &= \frac{1}{2\pi} \int_{-\infty}^{+\infty} dz \mathcal{F}(z) \int_{-\infty}^{+\infty} d\omega e^{-i\omega z} \\ &\times \sum_{i=1}^N e^{i\omega \xi_i} \left\{ \frac{1}{N} + \frac{\lambda_i}{2} \frac{\partial \xi_i}{\partial G_i} \right\}, \end{aligned}$$

or

$$\mathcal{F}(\lambda \cdot \mathbf{G}) = \sum_{i=1}^N \mathcal{F}(\xi_i) \left\{ \frac{1}{N} + \frac{1}{2} \lambda_i \frac{\partial \xi_i}{\partial G_i} \right\},$$

which is the “normal form” quoted in the text. A similar representation may be given for any  $\mathcal{F}(\mathbf{A} \cdot \mathbf{G})$  in terms of the  $SU(L)$  defining representation matrices  $\Lambda_\alpha$ , with  $L = N^2 - 1$  and  $\mathbf{A} \cdot \mathbf{G} = \psi_\alpha \cdot \Lambda_\alpha$ ,

$$\mathcal{F}(\mathbf{A} \cdot \mathbf{G}) = \sum_{i=1}^L \mathcal{F}(\xi_i[\psi]) \left\{ \frac{1}{L} + \frac{1}{2} \Lambda_\alpha \frac{\partial \xi_i[\psi]}{\partial \psi_\beta} \right\}.$$

<sup>1</sup>M. E. Brachet and H. M. Fried, *Phys. Lett. A* **103**, 309 (1984).

<sup>2</sup>M. E. Brachet and H. M. Fried, *J. Math. Phys.* **28**, 15 (1987).

<sup>3</sup>Another approach to Wilson loops, involving a strong coupling estimate by the extraction of relevant infrared effects, directly in the continuum, may be found in H-T. Cho, Ph. D. thesis, Brown University, 1987.

<sup>4</sup>H. M. Fried and J. Tesselndorf, *J. Math. Phys.* **25**, 1144 (1984). The  $SU(3)$  OE discussed and approximated in the adiabatic limit there is not unitary, in contrast to that of the present paper.

<sup>5</sup>The author is indebted to P. Sorba for pointing out the origin of this “normal form.”

<sup>6</sup>This is a variant of the Runge-Kutta method. The use of this sort of algorithm was kindly suggested to the author by D. Weingarten.

<sup>7</sup>W. P. Trower and N. Zovko, *Phys. Rev. D* **25**, 3088 (1982). The author thanks N. Zovko for an interesting discussion.



# Existence and approximation of the solutions of some nonlinear problems

S. R. Vatsya

Atomic Energy of Canada Limited, Whiteshell Nuclear Research Establishment, Pinawa, Manitoba, Canada ROE 1L0

(Received 11 September 1986; accepted for publication 7 January 1987)

Nonlinear equations defined by a positive definite, elliptic operator and nonlinear functions are considered. It is proved that a unique solution exists for a wider class of problems than previously determined, which may be approximated by an iterative method. These results are shown to hold for an equation arising in some reaction-diffusion phenomena.

## I. INTRODUCTION

Keller<sup>1</sup> and Pennline<sup>2</sup> have considered two-point boundary value problems of the form

$$-\frac{d^2u}{dx^2} = f(x;u), \quad 0 < x < 1, \quad u(0) = u(1) = 0, \quad (1)$$

where  $f$  is a function of  $u$ . Equation (1) is equivalent to the integral equation

$$u(x) = \int_0^1 d\xi g_k(x,\xi) [ku(\xi) + f(\xi;u(\xi))],$$

where  $g_k(x,\xi)$  is the Green's function defining the inverse of the operator  $(-d^2/dx^2 + k)$  with the stated boundary conditions, whenever it exists.

Keller<sup>1</sup> has shown that if, for all  $u$  and all  $x \in [0,1]$ ,  $\partial f / \partial u = f'$  is continuous and  $0 \geq f' \geq -N$  with some  $N \geq 0$ , then Eq. (1) has a unique solution  $u$  given by

$$u = \lim_{n \rightarrow \infty} u_n,$$

where  $u_0(x) = 0$ ,

$$u_{n+1}(x) = \int_0^1 d\xi g_k(x,\xi) [ku_n(\xi) + f(\xi;u_n(\xi))], \quad n = 0,1,2,\dots,$$

and  $k \geq N$ . The convergence is uniform with respect to  $x \in [0,1]$ . Pennline<sup>2</sup> points out that it is possible to take  $k \geq N/2$ , which improves the rate of convergence. Also an attempt was made at relaxing the conditions further.

The work of Ref. 2 was motivated by the equation

$$\frac{d^2v}{dx^2} = \phi^2 v^\nu, \quad \nu \geq 1, \quad \phi^2 \geq 0, \quad 0 < x < 1, \\ v'(0) = 0, \quad v(1) = 1, \quad (2)$$

which arises in the problem of steady-state, isothermal, reaction-diffusion of a substance involving  $n$ th-order kinetics.<sup>3</sup> By setting  $u = 1 - v$ , Eq. (2) reduces to an equivalent equation,

$$-\frac{d^2u}{dx^2} = \phi^2(1-u)^\nu, \quad 0 < x < 1, \\ u'(0) = 0, \quad u(1) = 0. \quad (3)$$

While the results of Refs. 1 and 2 are not applicable in this case, a similar but independent treatment was used in Ref. 2 to conclude that, with a restriction on the values of  $\phi$ , Eq.

(2) has a unique solution which may be approximated by the iterative method.

Let  $D$  be an open, bounded subset of  $R^l, l \geq 1$ , and let  $\partial D$ , the boundary of  $D$ , be piecewise continuous. Also let  $H$  be the Hilbert space of square integrable functions of  $x$  on  $\bar{D} = D \cup \partial D$  with the norm denoted by  $\|\cdot\|$ . In the present note, we consider the equation

$$Lu = f(u), \quad (4)$$

where  $L$  is a positive definite operator from  $H$  to  $H$ . Let  $\lambda$  be the greatest lower bound of  $L$ . It will be assumed that  $(L + k)^{-1}$  for each  $k > -\lambda$  is an integral operator with its kernel, still denoted by  $g_k(x,\xi)$ , being a non-negative, bounded function on  $\bar{D} \times \bar{D}$ . As a result of the boundedness of  $g_k(x,\xi)$ ,  $(L + k)^{-1}$  is a Hilbert-Schmidt operator and hence  $\lambda$  is an eigenvalue. These conditions are satisfied, in particular, for the case when  $L$  is the elliptic operator defined by

$$(Lu)(x) = - \sum_{i,j=1}^l \frac{\partial}{\partial x_i} \left( a_{ij}(x) \frac{\partial u(x)}{\partial x_j} \right) + a_0(x)u(x), \quad x \in D;$$

$$\alpha(x)u(x) + \beta(x) \frac{\partial u(x)}{\partial \eta} = 0, \quad x \in \partial D;$$

$$\alpha(x) \text{ and } \beta(x) \geq 0, \quad \alpha(x) \text{ or } \beta(x) > 0, \quad \alpha(x) \neq 0;$$

where  $\partial / \partial \eta$  is the conormal derivative, and the coefficients,  $a_{ij}(x)$ , are the elements of a continuously differentiable matrix-valued function on  $D$  that is bounded below by a positive constant and  $a_0(x) \geq 0$  is continuous.<sup>4</sup> It is clear that the operators appearing in Eqs. (1) and (3) are special cases of the elliptic operator considered here.

In addition to the generalization described above, we extend the results of Refs. 1 and 2 to include  $f$  with  $f'$  bounded above by a positive constant instead of zero. We show further that this weaker condition on  $f'$  need be satisfied on a conveniently characterized narrower set of functions only. As an application, it is deduced that Eq. (3) has a unique solution for all values of  $\phi$  which may be approximated by the iterative method. Present results complement, also, those of Ref. 5 where Eq. (4) with  $f$  belonging to a different class was studied.

## II. EXISTENCE AND APPROXIMATION OF $u$

Let  $S$  be a given set of functions and let  $|||\cdot|||$  denote the supremum norm, i.e.,  $|||v||| = \text{Sup}_{x \in \bar{D}} |v(x)|$  whenever it exists. Condition (C1) will be said to be satisfied on  $S$  if

(C1)  $S \subseteq H$ ,  $S$  is convex and it is closed with respect to  $|||\cdot|||$ .

A set  $S$  is convex iff  $v, w \in S$  implies that  $[tv + (1-t)w] \in S$  for each  $t \in [0, 1]$ , and  $S$  is closed in the stated sense iff, with a given sequence  $\{v_m\} \subset S$  of bounded functions  $|||v_n - v_m||| \rightarrow_{n,m \rightarrow \infty} 0$  implies that  $v = \lim_{n \rightarrow \infty} v_n \in S$ . The set  $S$  may or may not be closed with respect to  $||\cdot||$ .

With (C1) satisfied, let  $f$  be such that

(C2) there exist constants  $\gamma$  and  $N$  such that for all  $v \in S$  and all  $x \in \bar{D}$ ,  $\lambda > \gamma \geq f'(v) \geq -N$ ,  $N \geq 0$ ;

(C3) there exists a  $v \in S$  such that  $f(v) \in H$ .

If (C1) and (C2) are satisfied, then (C3) implies that  $f(v) \in H$  for each  $v \in S$  for the following. With  $v, w \in S$ ,

$$f(w) = f(v) + \int_0^1 dt f'[tv + (1-t)w](v-w).$$

From the convexity of  $S$  and (C2), it follows that

$$\gamma \geq \int_0^1 dt f'[tv + (1-t)w] \geq -N, \quad (5)$$

implying that  $|f(w)| \leq |f(v)| + \max(\gamma, N)|v-w|$ . Since  $|v-w| \in H$ ,  $f(v) \in H$  implies that  $f(w) \in H$ .

In view of the above,

$$A_k v = (L+k)^{-1}(kv + f(v)), \quad v \in S,$$

defines a one-parameter family of operators  $A_k$  from  $S$  to  $H$  and  $k \in (-\lambda, \infty)$ . Furthermore, the fixed points of  $A_k$  and the solutions of Eq. (4) are in a one-to-one correspondence. The parameter  $k$  will be assumed to be restricted by

$$(C4) \quad k \geq \frac{1}{2}[\max(\gamma, N) - \gamma] = \bar{k}.$$

We shall also assume

$$(C5) \quad \text{for some } k \text{ satisfying (C4), } A_k S \subseteq S.$$

With  $u_0 \in S$ , let  $u_{n+1} = A_k u_n$ ,  $n = 0, 1, 2, \dots$ ;  $\{u_n\}$  will be called the iterative sequence generated by  $u_0$ . It follows, by induction, that  $\{u_n\} \subset S$ , for all  $k$  in conformity with (C4) and (C5).

**Theorem 1:** Let  $S$ ,  $f$ , and  $k$  be such that the conditions (C1)–(C5) are satisfied, and let  $\{u_n\}$  be the iterative sequence generated by a  $u_0 \in S$ . Then  $A_k$  has a unique fixed point  $u$  in  $S$  and

$$|||u_n - u||| \rightarrow_{n \rightarrow \infty} 0.$$

*Proof:* We divide the proof in the following five steps.

*Step 1:*  $A_k$  is a contraction of  $S$  in  $H$ .

*Proof:* Since  $A_k S \subseteq S$ ,  $A_k$  is a map from  $S$  to  $S$ . With  $v, w \in S$ , we have

$$\begin{aligned} A_k v - A_k w &= (L+k)^{-1}[k(v-w) + f(v) - f(w)] \\ &= (L+k)^{-1} B_k(v, w)(v-w), \end{aligned}$$

where  $B_k(v, w)$  is the operation of multiplication defined by

$$B_k(v, w)h = \left(k + \int_0^1 dt f'[tv + (1-t)w]\right)h, \quad h \in H.$$

Since  $S$  is convex and (C2) holds, the inequality given by Eq. (5) is valid, which in view of (C4) yields

$$\begin{aligned} k + \gamma &\geq k + \int_0^1 dt f'[tv + (1-t)w] \\ &\geq -(k + \gamma), \quad k + \gamma \geq 0. \end{aligned}$$

Hence,  $\|B_k(v, w)\| \leq (k + \gamma)$ . Also,  $\|(L+k)^{-1}\| \leq 1/(k + \lambda)$ . It follows that

$$\|A_k v - A_k w\| \leq \mu_k \|v - w\|,$$

where

$$\mu_k = (k + \gamma)/(k + \lambda) < 1.$$

*Step 2:* For  $v \in S$ ,  $A_k v$  is bounded.

*Proof:* For each  $x \in \bar{D}$ ,

$$\begin{aligned} |(A_k v)(x)| &= \left| \int_{\bar{D}} d\xi g_k(x, \xi) [kv(\xi) + f(\xi; v(\xi))] \right| \\ &\leq \|kv + f(v)\| \left[ \int_{\bar{D}} d\xi g_k^2(x, \xi) \right]^{1/2} \end{aligned}$$

by the Schwarz inequality. Since  $g_k(x, \xi)$  is bounded and the measure of  $\bar{D}$  is finite, the bracketed term is bounded, say by  $M$ . Thus

$$|(A_k v)(x)| \leq M \|kv + f(v)\|.$$

*Step 3:* For  $v, w \in S$ ,  $|||A_k v - A_k w||| \leq M(k + \gamma) \times \|v - w\|$ .

*Proof:* For each  $x \in \bar{D}$ ,

$$\begin{aligned} |(A_k v - A_k w)(x)| &= \left| \int_{\bar{D}} d\xi g_k(x, \xi) \right. \\ &\quad \left. \times [B_k(v, w)(v-w)](\xi) \right| \\ &\leq M \|B_k(v, w)\| \|v - w\|, \end{aligned}$$

as in Step 2. The result follows by observing that  $\|B_k(v, w)\| \leq (k + \gamma)$ , as in Step 1.

*Step 4:*  $\{u_n\}$  is a Cauchy sequence with respect to  $|||\cdot|||$ .

*Proof:* A proof of the fact that  $\{u_n\} \subset S$  is a Cauchy sequence in  $H$ , i.e.,

$$\|u_n - u_m\| \rightarrow_{n,m \rightarrow \infty} 0,$$

follows from Step 1. Since  $\{u_n\}_0^\infty \subset S$ , from Step 2,  $\{u_n\}_1^\infty$  is bounded. From Step 3,

$$\begin{aligned} |||u_{n+1} - u_{m+1}||| &= |||A_k u_n - A_k u_m||| \\ &\leq M(k + \gamma) \|u_n - u_m\| \rightarrow_{n,m \rightarrow \infty} 0 \end{aligned}$$

for  $\{u_n\}$  is a Cauchy sequence in  $H$ .

*Step 5:* Result of Theorem 1.

*Proof:* Since  $S$  is closed with respect to  $|||\cdot|||$ , it follows from Step 4 that there exists a  $u \in S$  such that

$$\lim_{n \rightarrow \infty} |||u_n - u||| = \lim_{n \rightarrow \infty} |||A_k u_n - u||| = 0.$$

Since  $A_k$  is defined on  $S$ ,  $A_k u$  is well defined. From Step 3,

$$\begin{aligned} \|A_k u_n - A_k u\| &\leq M(k + \gamma) \|u_n - u\| \\ &\leq \bar{M}(k + \gamma) \|u_n - u\| \\ &\rightarrow 0, \end{aligned}$$

where  $\bar{M}$  is a constant. This implies that  $u = A_k u$ . To show the uniqueness, let  $A_k \tilde{u} = \tilde{u} \in S$ . Then from Step 1,

$$\begin{aligned} \|u - \tilde{u}\| &= \|A_k u - A_k \tilde{u}\| \\ &\leq \mu_k \|u - \tilde{u}\| \\ &= 0 \end{aligned}$$

for  $\mu_k < 1$ . This implies that  $u = \tilde{u}$  almost everywhere. However,  $u = \tilde{u}$  everywhere, for, using Step 3,

$$\begin{aligned} \|u - \tilde{u}\| &= \|A_k u - A_k \tilde{u}\| \\ &\leq M(k + \gamma) \|u - \tilde{u}\| \\ &= 0. \end{aligned}$$

Since  $\mu_k = (k + \gamma)/(k + \lambda)$  of Theorem 1, Step 1, is an increasing function of  $k$  on  $(-\lambda, \infty)$ , the rate of convergence is expected to improve with decreasing values of  $k$

$$\begin{aligned} (u_{n+1} - u_n)(x) &= \int_{\bar{D}} d\xi g_k(x, \xi) [B_k(u_n, u_{n-1})(u_n - u_{n-1})](\xi) \\ &\geq 0 \end{aligned}$$

for

$$k + \int_0^1 dt f'[tu_n + (1-t)u_{n-1}] \geq 0$$

implying that

$$[B_k(u_n, u_{n-1})(u_n - u_{n-1})](\xi) \geq 0.$$

The result follows by the induction principle.

We have used the non-negativity of  $g_k(x, \xi)$  in Proposition 1, which was not needed in Theorem 1.

If we take  $S = H$  in Theorem 1, then (C1) is clearly satisfied. If (C2) and (C3) hold and  $k$  is chosen according to (C4), then (C5) is also satisfied. Therefore the result of Theorem 1 is valid with  $S = H$ . This generalizes the result of Ref. 1. However, the requirement that (C2) be valid for all  $v \in H$  is too stringent to be satisfied by a large number of problems of interest. An attempt at relaxing this requirement was made in Ref. 2, but the new conditions are too limiting to be useful. In particular, it was assumed that, for all non-negative  $v$  bounded by a fixed constant,  $0 \geq f(v) \geq -\frac{1}{2}(\delta + N)v$ , with some  $\delta, N \geq 0$ . By setting  $v = 0$  one has that  $f(0) = 0$ , which is serious limitation.

In the following, we give some useful characterizations of sets that may replace  $S$  in Theorem 1 and Proposition 1.

Let  $f(0) \in H$  and

$$\bar{u}_\xi = (L + \xi)^{-1} |f(0)|, \quad \xi > -\lambda.$$

Non-negativity of  $g_\xi(x, \xi)$  implies that  $\bar{u}_\xi \geq 0$ . It is also bounded as in Step 2 of Theorem 1. Since

$$(L + k)\bar{u}_\xi = |f(0)| + (k - \xi)\bar{u}_\xi$$

it follows that

within the restriction imposed by (C4). Thus, for computational purposes, it is desirable to take  $k = \bar{k}$ . Further, it is to one's advantage to obtain tight bounds on  $f'$ , which can be checked by considering the behavior of  $\bar{\mu} = (\bar{k} + \gamma)/(\bar{k} + \lambda)$  with respect to the variations of  $\gamma$  and  $N$ . This point was discussed also in Ref. 2.

For  $f$  in a smaller class of functions, a stronger restriction on  $k$  enables one to obtain a monotonically convergent  $\{u_n\}$ , which we show in Proposition 1. This result may prove useful when a slower rate of convergence may be tolerated in favor of the monotonicity.

*Proposition 1:* In addition to the conditions of Theorem 1, let  $u_0 = 0 \in S, f(0) \geq 0$  for all  $x \in \bar{D}$ , and let  $k \geq N$ . Then, in addition to the result of Theorem 1,  $\{u_n\}$  converges monotonically to  $u$  from below.

*Proof:* We need show only that  $\{u_n\}$  is a nondecreasing sequence. Since  $g_k(x, \xi) \geq 0$ ,

$$u_1(x) = \int_{\bar{D}} d\xi g_k(x, \xi) f(\xi; 0) \geq 0 = u_0(x).$$

Let  $u_n \geq u_{n-1}$ ; we have that

$$\bar{u}_\xi = (L + k)^{-1} [|f(0)| + (k - \xi)\bar{u}_\xi], \quad k > -\lambda.$$

Let

$$Q_\xi = \{v: |v| \leq \bar{u}_\xi\}$$

and let

$$P_\xi = \{v: 0 \leq v \leq \bar{u}_\xi\}.$$

It is straightforward to check that (C1) and (C3) hold on  $Q_\xi$  and  $P_\xi$ . Assuming that (C2) is satisfied, one may restrict  $k$  according to (C4). In Proposition 2, we determine conditions that imply (C5) with  $S = Q_\xi$  and  $S = P_\xi$ .

*Proposition 2:* Let the symbols be as above. If there exist  $\xi$  and  $k$  such that, for all  $x \in \bar{D}$ ,

- (i) for each  $v \in P_\xi, 0 \leq kv + f(v) \leq |f(0)| + (k - \xi)\bar{u}_\xi$ ,
- (ii) for each  $v \in Q_\xi, 0 \leq |kv + f(v)| \leq |f(0)| + (k - \xi)\bar{u}_\xi$ ,

then (i)  $A_k P_\xi \subseteq P_\xi$  and (ii)  $A_k Q_\xi \subseteq Q_\xi$ .

*Proof:* We give a proof of (i); a proof of (ii) follows by similar arguments by estimating  $|A_k v|$  instead of  $A_k v$ .

For each  $v \in P_\xi$ , we have

$$\begin{aligned} 0 \leq (A_k v)(x) &= \int_{\bar{D}} d\xi g_k(x, \xi) [kv(\xi) + f(\xi; v(\xi))] \\ &\leq \int_{\bar{D}} d\xi g_k(x, \xi) [|f(\xi; 0)| + (k - \xi)\bar{u}_\xi(\xi)] \\ &= \bar{u}_\xi(x). \end{aligned}$$

It is clear that, if (C2) holds on  $P_\xi, Q_\xi$  and with  $k$  in conformity with (C4), the corresponding condition of Proposition 2 also holds; then (C1)-(C5) are satisfied with  $S = P_\xi, S = Q_\xi$ , implying the validity of the result of

Theorem 1. In case of  $P_\zeta$ , if  $u_0 = 0$  and  $k \geq N$ , then, from Proposition 1, a monotonic convergence results.

### III. APPLICATION

It follows from the results of Sec. III of Ref. 2 that Eq. (3) has a unique solution  $u$ , satisfying

$$0 \leq u(x) \leq 1 - (\cosh \phi x) / (\cosh \phi) = \chi(x)$$

for each  $\phi$  such that

$$(\cosh \phi - 1) / (\cosh \hat{k} - 1) \geq \phi^2 / (\hat{k}^2 \cosh \phi),$$

where

$$\hat{k} = \frac{1}{2} \nu \phi^2 [1 + (1/\cosh \phi)^{\nu-1}].$$

Also the iterative sequence generated by  $\chi(x)$ , with  $k = \hat{k}$ , was shown to converge to  $u(x)$ . The author considered Eq. (2); the results given here for Eq. (3) follow by setting  $u = 1 - v$ . In the following we use the results of Sec. II to show that the existence and convergence results hold for Eq. (3) for all real values of  $\phi$ , with any  $k \geq \hat{k}$  and any  $u_0(x)$  such that  $0 \leq u_0(x) \leq \chi(x)$ .

First we note that, in this case,  $l = 1$ ,  $\bar{D} = [0, 1]$ , and  $\lambda = \pi^2/4$ . For  $\nu = 1$ , Eq. (3) is solved easily, yielding  $u = \chi$ .

Let  $\nu \geq 2$ . With  $\zeta = \phi^2$ ,  $\bar{u}_\zeta = \chi \leq 1$ , and  $0 \leq f(0) = \phi^2 \in H$ . Let  $P_\zeta$  be defined as in Sec. II, i.e.,  $v \in P_\zeta$  implies that  $0 \leq v \leq \bar{u}_\zeta \leq 1$ . Now for  $v \in P_\zeta$ ,

$$-\delta = -\nu \phi^2 (1/\cosh \phi)^{\nu-1} \geq f' \geq -\nu \phi^2.$$

Thus (C2) is satisfied with  $\gamma = -\delta$  and  $N = \nu \phi^2$ . Let

$$k \geq \frac{1}{2} (\nu \phi^2 + \delta),$$

which satisfies (C4). For each  $v \in P_\zeta$ ,

$$\phi^2 + (k - \phi^2)v = kv + \phi^2(1 - v) \geq kv + f(v) \geq 0.$$

For  $\nu \geq 2$ ,  $k \geq \phi^2$ ; consequently,

$$\phi^2 + (k - \phi^2)v \leq \phi^2 + (k - \phi^2)\bar{u}_\zeta = f(0) + (k - \zeta)\bar{u}_\zeta,$$

and the condition (i) of Proposition 2 is satisfied.

Since the conditions (C1)–(C5) of Theorem 1 are satisfied, Eq. (4) has a unique solution,  $u \in P_\zeta$ , for each value of  $\phi$  and the iterative sequence generated by any  $u_0 \in P_\zeta$  converges uniformly to  $u$ . If we take  $u_0 = 0$  and  $k \geq \nu \phi^2$ , then the convergence is also monotonic from below.

It is not necessary to restrict  $\nu \geq 2$ . In fact,  $\nu \geq 1$  may be allowed to be a continuous parameter. For the above arguments to hold, all one has to do is to take

$$k \geq \max(\phi^2, \frac{1}{2}(\nu \phi^2 + \delta)).$$

<sup>1</sup>H. B. Keller, *Numerical Methods for Two-Point Boundary Value Problems* (Blaisdell, Waltham, MA, 1968), pp. 108–109.

<sup>2</sup>J. A. Pennline, *Math. Comp.* **37**, 127 (1981).

<sup>3</sup>R. Aris, *The Mathematical Theory of Diffusion and Reaction in Permeable Catalysis* (Clarendon, Oxford, 1975), Vol. 1, pp. 101–239.

<sup>4</sup>N. Aronszajn and K. Smith, *Am. J. Math.* **79**, 611 (1957); see also, H. B. Keller, in *Bifurcation Theory and Nonlinear Eigenvalue Problems*, edited by J. B. Keller and S. Antman (Benjamin, New York, 1969), pp. 217–255.

<sup>5</sup>S. R. Vatsya, *J. Math. Phys.* **22**, 2977 (1981).

# On the existence of infinitely many resonances in the scattering problem based on moment conditions and entire functions

J. S. Hwang

*Institute of Mathematics, Academia Sinica, Taipei, Taiwan, China 11529*

(Received 7 August 1986; accepted for publication 30 January 1987)

It has been well proved that there are infinitely many resonances in the scattering problem provided the potential  $V(t)$  satisfies the moment conditions  $\int_0^\infty |V(t)|t^n dt = O(n^{(1-\epsilon)n})$ ,  $\epsilon > 0$ ,  $n = 0, 1, \dots$ . In particular, if  $V(t)$  has a compact support then the result of Rollnik [H. Rollnik, *Z. Phys.* **145**, 654 (1956)] and Regge [T. Regge, *Nuovo Cimento*, **8**, 671 (1958)] is obtained and if  $|V(t)| \sim \exp(-\tau t^{1+\epsilon})$ ,  $\tau, \epsilon > 0$ , we have the results of Sartori [L. Sartori, *J. Math. Phys.* **4**, 1408 (1963)].

## I. INTRODUCTION

It is well known (see Agranovich and Marchenko,<sup>1</sup> p. 20) that the Schrödinger equation

$$y'' + z^2 y = V(x)y, \quad 0 < x < \infty,$$

has a solution  $E(x, z)$  given by

$$E(x, z) = e^{izx} + \int_x^\infty K(x, t) e^{izt} dt, \quad (1)$$

where the kernel  $K(x, t)$  satisfies

$$|K(x, t)| \leq \frac{1}{2} e^{\sigma(x)} \sigma[(x+t)/2], \quad (2)$$

where

$$\sigma(x) = \int_x^\infty |V(t)| dt < \infty$$

and

$$\sigma_1(x) = \int_x^\infty |V(t)| t dt < \infty.$$

For  $x = 0$ , the function  $f(z) = E(0, z)$  is analytic in  $\text{Im } z > 0$  (the Jost function, see Newton,<sup>2</sup> p. 340) and as usual, a zero of  $f(z)$  is called a resonance in the scattering problem (excluding the bound state, see Ref. 2, p. 360). Naturally, we may ask under what condition on the potential  $V(x)$  can there be infinitely many resonances? In this paper, we answer this question by the following.

**Theorem 1:** There are infinitely many resonances if the Jost function  $f(z)$  is entire and the potential satisfies the moment conditions

$$\int_x^\infty |V(t)| t^n dt = O(n^{(1-\epsilon)n}), \quad x \geq 0, \quad \epsilon > 0, \quad n = 0, 1, \dots \quad (3)$$

As a consequence of Theorem 1, we obtain the following known result due to Rollnik<sup>3</sup> and Regge.<sup>4</sup>

*Corollary 2:* There are infinitely many resonances if the potential has a compact support.

To prove Corollary 2, we need only observe that  $|V(t)| = 0$  for all  $t \geq T$ , where  $T > 0$  is fixed, so that

$$\int_x^\infty |V(t)| t^n dt \leq CT^{n+1}, \quad \text{for some } C > 0.$$

This satisfies the moment conditions (3). Clearly, the Jost function is entire and hence the assertion follows from Theorem 1.

Although the regular solution is always an entire function, the Jost solution and the scattering operator are only entire functions for potentials which decrease faster than an exponential at infinity, see Nussenzweig (Ref. 5, p. 201).

## II. ENTIRE FUNCTIONS

As we saw in the Introduction, the resonances in the scattering problem are the same as the zeros of a certain class of entire functions. To study such a class, we call an entire function  $f \in B$  if  $f(z)$  is bounded in a half-plane  $H$  and  $f(z)$  has a radial limit  $l$ , as  $z \rightarrow \infty$ , along the boundary of  $H$ , where we require  $f$  be nonconstant.

We note that the function  $f(z) = E(0, z)$  defined in (1) belongs to class  $B$  but not the function  $e^{izx}$ ,  $x > 0$ . Functions in class  $B$  can possibly have no zeros as will be seen from the following.

**Theorem 3:** If  $f \in B$ , then  $e^f \in B$  and the function  $f(z)$  tends to  $l$  uniformly in  $H$ , as  $z \rightarrow \infty$ , so that  $f(z)$  has at most finitely many zeros in  $H$  provided  $l \neq 0$ . If  $f(z)$  is of exponential type  $\tau$  and belongs to  $L^2$  on the real axis, then  $f(z)$  can be represented by

$$f(z) = \int_{-\tau}^{\tau} e^{izt} V(t) dt, \quad V \in L^2(-\tau, \tau).$$

*Proof:* We first observe that both of the boundedness and the radial limit of  $f \in B$  are preserved by the exponential function and hence  $e^f \in B$ .

Next, in view of Montel's theorem (see Boas,<sup>6</sup> p. 5), we see that the function  $f(z)$  tends to  $l$  uniformly in  $H$ , as  $z \rightarrow \infty$ , so that it cannot have infinitely many zeros in  $H$  provided  $l \neq 0$ .

Finally, the representation of  $f(z)$  follows from the Paley-Wiener theorem (see Ref. 6, p. 103) and the proof is complete.

Note that the hypothesis  $l \neq 0$  in the above result is necessary. For instance, the function

$$f(z) = (e^{iz} - e^{-1}) / (z - i)$$

belongs to  $B$  in the upper half-plane with  $l = 0$  and has infinitely many zeros there at  $i + 2n\pi$ ,  $n = 1, 2, \dots$

Note that the function  $e^f$  is of infinite order and therefore the condition on the finite order is what we need to have infinitely many zeros as will be seen from the following.

**Theorem 4:** If  $f \in B$  and if  $f$  is of finite order, then  $f$  has infinitely many zeros.

*Proof:* Suppose on the contrary that  $f$  has only finitely many zeros. Then  $f$  can be represented by

$$f(z) = e^{P_n(z)} Q_m(z),$$

where  $P_n$  and  $Q_m$  are polynomials of degree  $n$  and  $m$ , respectively. Without loss of generality, we may assume that the half-plane in question is the upper half-plane  $U = \{z: \text{Im } z > 0\}$ . We set

$$P_n(z) = a_0 + a_1 z + \dots + a_n z^n, \quad a_n = r e^{i\alpha} \neq 0, \quad 0 \leq \alpha < 2\pi.$$

Then it is easy to see that  $P_n(z) \rightarrow \infty$ , as  $z \rightarrow \infty$  along the ray  $\Gamma = \{\text{Re}^{i(2\pi - \alpha)/n}; 0 \leq R < \infty\}$ . Clearly, if  $n > 1$  then the ray  $\Gamma$  lies on  $U$  which contradicts the boundedness of  $f$  in  $U$ . Hence we must have  $n = 1$  and

$$f(z) = e^{irz} Q_m(z).$$

This in turn implies that  $f(z) \rightarrow \infty$ , as  $z \rightarrow \infty$  along the real axis, a contradiction again. We thus conclude that  $f$  has infinitely many zeros.

We note that any entire function which is not a polynomial takes every value in its range with one possible exception infinitely often due to Picard's theorem (see Titchmarsh,<sup>7</sup> p. 277). The function  $e^f$ ,  $f \in B$ , has the exceptional value 0. However, any function considered in Theorem 4 has no exceptional values in its range as will be seen from the following extension.

*Corollary 5:* Under the hypothesis of Theorem 4, the function  $f$  assumes every value infinitely often.

*Proof:* For any value  $v$ , the function  $f - v \in B$  and of finite order. Hence the assertion follows from Theorem 4.

**Theorem 6:** Let  $f(z) = c + \int_x^\infty K(x,t) e^{izt} dt$ , where  $c$  is a constant,  $x \geq 0$  is fixed, and  $K(x,t)$  satisfies

$$\int_x^\infty |K(x,t)| t^n dt = O(n^{(1-\epsilon)n}), \quad \epsilon > 0, \quad n = 0, 1, \dots \quad (4)$$

If  $f(z)$  is an entire function, then it is of order  $\leq 1/\epsilon$  and has infinitely many zeros.

*Proof:* Since  $f(z)$  is entire, it can be expanded as

$$f(z) = \sum_{n=0}^{\infty} a_n z^n,$$

where

$$a_n = \frac{(-i)^n}{n!} \int_x^\infty K(x,t) t^n dt.$$

Using (4) and Stirling's formula, we see that  $f(z)$  is an entire function of order (see Ref. 6, p. 9)

$$\rho = \limsup_{n \rightarrow \infty} \frac{n \log n}{\log(1/|a_n|)} \leq \lim_{n \rightarrow \infty} \frac{n \log n}{n \epsilon (\log n - 1)} = \frac{1}{\epsilon}.$$

Since the function  $f(z)$  is bounded in  $\text{Im } z \geq 0$  and  $f(z) \rightarrow c$ , as  $z \rightarrow \pm \infty$  along the real axis, it follows from Theorem 4 that  $f$  has infinitely many zeros. This completes the proof.

Note that from Corollary 5 we see that the function

$$F(z) = \int_x^\infty K(x,t) e^{izt} dt$$

assumes any value infinitely often in the lower half-plane.

With the help of Theorem 6, we are now ready to prove Theorem 1.

### III. PROOF OF THEOREM 1

According to (2), we have the inequality

$$|K(x,t)| \leq \frac{1}{2} e^{\sigma_1(x)} \sigma(x+t)/2.$$

It follows from (3) that

$$\begin{aligned} \int_x^\infty |K(x,t)| t^n dt &\leq \frac{1}{2} e^{\sigma_1(x)} \int_x^\infty \left( \int_{(x+t)/2}^\infty |V(s)| ds \right) t^n dt \\ &= \frac{1}{2} e^{\sigma_1(x)} \int_x^\infty |V(s)| ds \int_x^{2s-x} t^n dt \\ &< \frac{1}{n+1} e^{\sigma_1(x)} \int_x^\infty |V(s)| (2s-x)^{n+1} ds \\ &= \frac{2^{n+1}}{n+1} e^{\sigma_1(x)} O((n+1)^{(1-\epsilon)(n+1)}) \\ &= O(n^{(1-\epsilon)n}). \end{aligned}$$

Hence the kernel  $K(x,t)$  satisfies condition (4) in Theorem 6 and the assertion follows from that theorem.

As a consequence of Theorem 1, we obtain Corollary 2 and the following result of Sartori.<sup>8</sup>

*Corollary 7:* There are infinitely many resonances if the potential satisfies

$$|V(t)| \sim \exp(-\tau t^{1+\epsilon}), \quad \tau, \epsilon > 0. \quad (5)$$

*Proof:* Clearly condition (5) gives

$$\int_0^\infty \exp(-\tau t^{1+\epsilon}) t^n dt = ((1+\epsilon)\tau^N)^{-1} \Gamma(N+1),$$

where  $N = (1+n)/(1+\epsilon)$  and  $\Gamma(N+1) = N! < N^N < n^{(1-\epsilon/2)n}$  for all  $n > 2(\epsilon(1-\epsilon))^{-1}$ ,  $0 < \epsilon < 1$  (it suffices to prove the case  $\epsilon < 1$ ). It follows that

$$\int_0^\infty |V(t)| t^n dt = O(n^{(1-\delta)n}), \quad 0 < \delta < \frac{\epsilon}{2}.$$

This together with Theorem 1 yields the assertion.

### IV. CONCLUSIONS AND OPEN PROBLEMS

Note that the method used does not generalize to potentials with exponential or slower decrease at infinity.

In view of the function defined in Theorem 3,

$$f(z) = (e^{iz} - e^{-1})/(z-i),$$

which has infinitely many zeros in  $U$ , we may ask whether there is a potential  $V(t)$  such that the function

$$f_V(z) = \int_0^\infty e^{izt} V(t) dt$$

is entire and has infinitely many zeros in  $U$ . In general, we do not know the answer. However, if we replace  $U$  by the closure  $\bar{U}$ , then such a potential does exist. For instance, we may take  $V(t) = 1$ ,  $0 \leq t \leq 2\pi$ , and  $V(t) = 0$ ,  $t > 2\pi$ . Then

$$f_V(z) = (e^{i2\pi z} - 1)/(iz),$$

which has infinitely many zeros at  $z = 1, 2, \dots$ .

Naturally, we may ask whether there is a potential such that

$$(e^{iz} - e^{-1})/(z - i) = \int_0^\infty e^{izt} V(t) dt. \quad (6)$$

The answer should be negative. To see this, we expand both sides of (6) and we then have

$$\int_0^\infty t^n e^{-t} V(t) dt = e^{-1}, \quad n = 0, 1, \dots \quad (7)$$

We conjecture that no summable function  $f(t)$  can satisfy the following condition of constant moment:

$$\int_0^\infty t^n f(t) dt = c \neq 0, \quad n = 0, 1, \dots \quad (8)$$

Clearly, the nonexistence of the potential  $V(t)$  in (7) is a particular case in (8). However, if  $c = 0$ , then there does exist a summable function satisfying (8). For instance, the function (see Shohat and Tamarkin,<sup>9</sup> p. 22)

$$f(t) = (\sin t^{1/4}) \exp(-t^{1/4}), \quad 0 \leq t < \infty,$$

is summable and satisfies (8) when  $c = 0$ .

Of course, if we consider the distribution  $d\mu(t)$  in place of the density  $f(t)dt$ , then there does exist a distribution satisfying

$$\int_0^\infty t^n d\mu(t) = c, \quad n = 0, 1, \dots \quad (9)$$

In fact, we can define the point mass:  $\mu(1) = c$ ,  $t \geq 1$ , and  $\mu(t) = 0$ ,  $t < 1$ , then (9) obviously holds.

We now explain the reason for our conjecture with regard to (8). More precisely, we shall prove that no potentials with compact supports can satisfy (8). Suppose on the contrary that there is such a potential  $V(t) = 0$ ,  $t \geq T$ , satisfying (8). Then we have

$$\int_0^T t^n V(t) dt = c \neq 0, \quad n = 0, 1, \dots \quad (10)$$

Subtracting two consecutive terms in (10) gives

$$\int_0^T t^n (1-t) V(t) dt = 0, \quad n = 0, 1, \dots$$

It follows from the Müntz-Szász theorem (see Rudin,<sup>10</sup> p. 304) that  $(1-t)V(t) = 0$  or  $V(t) = 0$  holds almost everywhere on  $(0, T)$ . This in turn implies that  $c = 0$ , a contradiction. We thus conclude that no potentials with compact supports can satisfy (8).

We note that the above assertion can also be proved via a different approach (Ref. 9, p. 5). Also note that some extensions of the Müntz-Szász theorem with regard to the completeness have been done by the author.<sup>11</sup>

In closing this note, we finally pose the following moment problem: What is a necessary and sufficient condition on a sequence  $\{\mu_n\}$  of real numbers such that there is a potential  $V(t)$  satisfying

$$\int_0^\infty t^n V(t) dt = \mu_n, \quad n = 0, 1, \dots ?$$

## ACKNOWLEDGMENT

I am indebted to the referee for his valuable comments.

<sup>1</sup>Z. S. Agranovich and V. A. Marchenko, *The Inverse Problem of Scattering Theory* (Gordon and Breach, New York, 1963).

<sup>2</sup>R. G. Newton, *Scattering Theory of Waves and Particles* (McGraw-Hill, New York, 1982), 2nd ed.

<sup>3</sup>H. Rollnik, "Zur theorie der zerfallenden zustände," *Z. Phys.* **145**, 654 (1956).

<sup>4</sup>T. Regge, "Analytic properties of the scattering matrix," *Nuovo Cimento* **8**, 671 (1958).

<sup>5</sup>H. M. Nussenzweig, *Causality and Dispersion Relations* (Academic, New York, 1972).

<sup>6</sup>R. P. Boas, *Entire Functions* (Academic, New York, 1954).

<sup>7</sup>E. C. Titchmarsh, *The Theory of Functions* (Oxford U. P., New York, 1939), 2nd ed.

<sup>8</sup>L. Sartori, "Asymptotic behavior of Schrödinger scattering amplitudes," *J. Math. Phys.* **4**, 1408 (1963).

<sup>9</sup>J. A. Shohat and J. D. Tamarkin, *The Problem of Moments* (Am. Math. Soc., Providence, RI, 1943).

<sup>10</sup>W. Rudin, *Real and Complex Analysis* (McGraw-Hill, New York, 1966).

<sup>11</sup>J. S. Hwang, "A note on Bernstein and Müntz-Szász theorems with applications to the order statistics," *Ann. Inst. Statist. Math.* **30**, 167 (1978), Part A.

# Do trilinear commutation relations in quantum mechanics admit coordinate space realization in three dimensions?

Ranjan Bhattacharya and Siddhartha Bhowmick  
 Department of Physics, Jadaupur University, Calcutta 700032, India

(Received 7 October 1985; accepted for publication 14 January 1987)

The trilinear commutation relations involving coordinates and momenta introduced by Wigner [E. P. Wigner, *Phys. Rev.* **77**, 711 (1950)] are generalized to three dimensions. It is shown that the only realizable coordinate space representation of the momenta implies the usual bilinear commutation relations.

## I. INTRODUCTION

It is well known that in quantum mechanics trilinear commutation relations involving coordinates and momenta, rather than the usual bilinear ones, can be introduced for a one-dimensional harmonic oscillator.<sup>1</sup> A nontrivial coordinate space representation of the momentum operator necessitates the use of wave functions that are not analytic in the usual sense. It seems worthwhile to investigate whether similar conclusions are imperative for the same problem in three dimensions.

The trilinear commutation relations for the one-dimensional oscillator follow by requiring that the equations of motion obtained from the Heisenberg equations are the same as the classical ones. The trilinear relations thus obtained are

$$[q, \{q, p\}] = 2iq, \quad [p, \{q, p\}] = -2ip, \quad (1)$$

where the brackets  $\{ \}$  and  $[ \ ]$  refer, respectively, to an anti-commutator and a commutator. Yang<sup>2</sup> found the coordinate representation of the momentum operator  $p$ ,

$$p = -i \frac{d}{dq} + i \frac{C}{q} R, \quad (2)$$

where  $C$  is a real constant, and  $R$  the inversion operator, i.e.,

$$RqR^{-1} = -q, \quad R \frac{d}{dq} R^{-1} = -\frac{d}{dq}. \quad (3)$$

In trying to solve the oscillator problem with the representation (2), Yang concluded that  $C$  must be zero if one requires the wave functions to be analytic. Ohnuki and Kamefuchi,<sup>3,4</sup> by introducing generalized functions, or "hyperfunctions" as wave functions, solved the oscillator problem with  $C \neq 0$ .

## II. TRILINEAR COMMUTATION RELATIONS IN THREE DIMENSIONS

We now consider what generalizations can be made when three degrees of freedom, rather than a single one, are involved. We, therefore, consider a three-dimensional oscillator with the Hamiltonian

$$H = \frac{1}{2} \sum_k (p_k^2 + x_k^2). \quad (4)$$

Here, and subsequently, all the indices run over the values 1, 2, and 3. The demand that the classical equations of motion

follow from the Heisenberg equations leads to the trilinear relations

$$[x_k, \{p_l, x_k\}] = 2i\delta_{kl}x_k, \quad [p_k, \{x_l, p_k\}] = -2i\delta_{kl}p_k. \quad (5)$$

We assume relations more general than the ones above and write

$$\begin{aligned} \text{(i)} \quad & [x_k, \{p_l, x_m\}] = 2i\delta_{kl}x_m, \\ \text{(ii)} \quad & [p_k, \{x_l, p_m\}] = -2i\delta_{kl}p_m, \\ \text{(iii)} \quad & [p_k, \{x_l, x_m\}] = -2i\delta_{kl}x_m - 2i\delta_{km}x_l, \\ \text{(iv)} \quad & [x_k, \{p_l, p_m\}] = 2i\delta_{kl}p_m + 2i\delta_{km}p_l. \end{aligned} \quad (6)$$

Not all of the above relations are independent. In fact, (iii) and (iv) follow from (i) and (ii), respectively, by using the (generalized) Jacobi identity. If we now introduce the operators  $a_k$  and  $a_k^\dagger$  defined by

$$a_k \equiv (1/\sqrt{2})(x_k + ip_k), \quad a_k^\dagger \equiv (1/\sqrt{2})(x_k - ip_k), \quad (7)$$

the relations (6) become in terms of these operators,

$$\begin{aligned} [a_k, \{a_l^\dagger, a_m\}] &= 2\delta_{kl}a_m, \\ [a_k, \{a_l^\dagger, a_m^\dagger\}] &= 2\delta_{kl}a_m^\dagger + 2\delta_{km}a_l^\dagger, \\ [a_k, \{a_l, a_m\}] &= 0. \end{aligned} \quad (8)$$

These are the well-known trilinear relations for the creation and annihilation operators for a system of para-Bose oscillators.<sup>5</sup>

## III. COORDINATE SPACE REALIZATION

To find the coordinate representation of the momenta  $p_i$  satisfying the trilinear relations (6), consider the operator

$$S_{ij} \equiv [x_i, p_j] - i\delta_{ij}. \quad (9)$$

Therefore

$$S_{ij}^\dagger = -S_{ij}. \quad (10)$$

Then using the relations (6), we have

$$\{p_k, S_{ij}\} = \{x_k, S_{ij}\} = 0. \quad (11)$$

Taking the matrix element of the second relation of (11) between the states  $|\bar{x}'\rangle$  and  $|\bar{x}''\rangle$ , we get

$$\langle \bar{x}' | \{x_k, S_{ij}\} | \bar{x}'' \rangle = (x_k' + x_k'') \langle \bar{x}' | S_{ij} | \bar{x}'' \rangle = 0. \quad (12)$$

Or,



$$\langle \bar{x}' | S_{ij} | \bar{x}'' \rangle = 2iC_{ij}(\bar{x}')\delta(\bar{x}' + \bar{x}''), \quad (13)$$

where the  $2i$  has been introduced for convenience. Thus

$$S_{ij} = 2iC_{ij}(\bar{x})R. \quad (14)$$

Here  $R$  is the inversion operator, i.e.,

$$R\Psi(\bar{x}) = \Psi(-\bar{x}).$$

We also note, using (10),

$$C_{ij}^*(\bar{x}) = C_{ij}(-\bar{x}). \quad (15)$$

The matrix element of  $S_{ij}$  is also given by

$$\begin{aligned} \langle \bar{x}' | S_{ij} | \bar{x}'' \rangle &= (x'_i - x''_i) \langle \bar{x}' | p_j | \bar{x}'' \rangle - i\delta_{ij}\delta(\bar{x}' - \bar{x}'') \\ &= 2iC_{ij}(\bar{x}')\delta(\bar{x}' + \bar{x}''), \end{aligned}$$

using (9) and (13).

We therefore have

$$\begin{aligned} \langle \bar{x}' | p_j | \bar{x}'' \rangle &= -i\frac{\partial}{\partial x_j}\delta(\bar{x}' - \bar{x}'') + \frac{iC_{ij}(\bar{x}')\delta(\bar{x}' + \bar{x}'')}{x'_i} \\ &\quad + B_j(\bar{x}')\delta(\bar{x}' - \bar{x}''), \end{aligned} \quad (16)$$

where we have used

$$x_i \frac{\partial}{\partial x_j} \delta(\bar{x}) = -\delta_{ij}\delta(\bar{x}).$$

Here  $B_j(\bar{x}')$  is an arbitrary real function as  $p_j$  is Hermitian. Since the right-hand side of (16) cannot depend on the index  $i$ ,  $C_{ij}(\bar{x}')$  must be proportional to  $x'_i$ . Now, from the requirement of the correct transformation properties of  $p_j$  under rotations, we must have

$$C_{ij}(\bar{x}') = x_i x_j f(|\bar{x}|), \quad B_j(\bar{x}) = x_j g(|\bar{x}|). \quad (17)$$

Hence we have from (16),

$$\begin{aligned} \langle \bar{x}' | p_j | \Psi \rangle &= -i\frac{\partial}{\partial x_j}\Psi(\bar{x}) + ix_j f(|\bar{x}|)\Psi(-\bar{x}) \\ &\quad + x_j g(|\bar{x}|)\Psi(\bar{x}), \end{aligned} \quad (18)$$

whence,

$$p_j = -i\frac{\partial}{\partial x_j} + ix_j f(|\bar{x}|)R + x_j g(|\bar{x}|). \quad (19)$$

We know, however, that the arbitrary function  $B_j(\bar{x})$  in (16) can be gauged away with a unitary transformation if the integrability condition  $\partial B_j/\partial x_i = \partial B_i/\partial x_j$  holds. We find from (17) that this condition is obviously satisfied. It is easy to see that under this transformation  $C_{ij}(\bar{x}) \rightarrow C'_{ij}(\bar{x})$  which also satisfies (15). The argument leading to the explicit form of  $C_{ij}$  in (17) also holds for  $C'_{ij}$ . Therefore we have

$$p_j = -i\frac{\partial}{\partial x_j} + ix_j f(|\bar{x}|)R. \quad (20)$$

The form of  $f$  can now be determined from the first of the relations (10). We have, using (14),

$$\{p_k, S_{ij}\} = \left\{ -i\frac{\partial}{\partial x_k} + ix_k f(|\bar{x}|)R, 2iC_{ij}(\bar{x})R \right\} = 0. \quad (21)$$

This gives, finally,

$$\delta_{ik}x_j f(|\bar{x}|)R + \delta_{jk}x_i f(|\bar{x}|)R + x_i x_j \frac{\partial}{\partial x_k} f(|\bar{x}|)R = 0. \quad (22)$$

Or, since  $R^{-1} = R$  exists,

$$\delta_{ik}x_j f(|\bar{x}|) + \delta_{jk}x_i f(|\bar{x}|) + \frac{x_i x_j x_k}{|\bar{x}|} \frac{d}{d|\bar{x}|} f(|\bar{x}|) = 0. \quad (23)$$

Equation (23) has the solution

$$f(|\bar{x}|) = c\delta(|\bar{x}|). \quad (24)$$

This does not contribute to  $p_j$  as  $x_j\delta(|\bar{x}|) = 0$ . Therefore we have for the coordinate representation of the momentum operator,

$$p_i = -i\frac{\partial}{\partial x_i}. \quad (25)$$

It is clear that this representation necessarily implies the canonical bilinear commutation relations between  $x$  and  $p$ . This, of course, does not rule out the possible existence of other inequivalent representations. In fact, the usual parastatistics does provide a representation of trilinear commutation relations in terms of the so called generalized Bose numbers (see the Appendix). We therefore conclude that the generalization of Wigner's trilinear commutation relations for coordinates and momenta does not admit a nontrivial coordinate representation in three dimensions.

#### IV. CONCLUSION

To summarize, in one dimension a nontrivial coordinate representation of the momentum operator satisfying trilinear commutation relations is possible at the expense of introducing generalized functions as wave functions. A generalization to three dimensions shows that the only realizable coordinate space representation of the momentum operator implies the usual bilinear commutation relations. In Schrödinger quantum mechanics, therefore, there is no place for trilinear commutation relations in three dimensions.

#### ACKNOWLEDGMENTS

We thank the referee for pointing out the general solution (24) of Eq. (23). We are grateful to Dr. A. Chatterjee and Professor B. Dutta Roy for giving more help than they admit.

#### APPENDIX: GENERALIZED BOSE NUMBER REPRESENTATION

Representations of trilinear commutation relations exist in terms of generalized Bose numbers<sup>4</sup> which are the analogs of Green components.<sup>5</sup> Let

$$x_k = \sum_{\alpha=1}^p x_k^\alpha \quad \text{and} \quad p_k = -i \sum_{\alpha=1}^p \frac{\partial}{\partial x_k^\alpha},$$

where

$$\{x_k^\alpha, x_l^\beta\} = 0, \quad \alpha \neq \beta, \quad [x_k^\alpha, x_l^\alpha] = 0,$$

and

$$\left\{ \frac{\partial}{\partial x_k^\alpha}, \frac{\partial}{\partial x_l^\beta} \right\} = 0, \quad \alpha \neq \beta, \quad \left[ \frac{\partial}{\partial x_k^\alpha}, \frac{\partial}{\partial x_l^\alpha} \right] = 0,$$

with

$$\left\{ \frac{\partial}{\partial x_k^\alpha}, x_l^\beta \right\} = 0, \quad \alpha \neq \beta, \quad \left[ \frac{\partial}{\partial x_k^\alpha}, x_l^\alpha \right] = \delta_{kl}.$$

Then the above  $x$  and  $p$  satisfy the trilinear commutation relations (6).

<sup>1</sup>E. P. Wigner, Phys. Rev. **77**, 711 (1950).

<sup>2</sup>L. M. Yang, Phys. Rev. **84**, 788 (1951).

<sup>3</sup>Y. Ohnuki and S. Kamefuchi, J. Math. Phys. **19**, 67 (1978).

<sup>4</sup>Y. Ohnuki and S. Kamefuchi, *Quantum Field Theory and Parastatistics* (Springer, Berlin, 1982).

<sup>5</sup>H. S. Green, Phys. Rev. **90**, 270 (1953); O. W. Greenberg and A. Messiah, *ibid.* **138**, B1155 (1965).

# Electromagnetic scattering of an arbitrary plane wave from a spherical shell with a circular aperture

Richard W. Ziolkowski

*Electronics Engineering Department, Lawrence Livermore National Laboratory, P.O. Box 5504, L-156, Livermore, California 94550*

William A. Johnson<sup>a)</sup>

*Division 7553, Sandia National Laboratories, Albuquerque, New Mexico 87185*

(Received 23 July 1985; accepted for publication 30 January 1987)

The problem of the scattering of an electromagnetic plane wave with arbitrary polarization and angle of incidence from a perfectly conducting spherical shell with a circular aperture is solved with a generalized dual series approach. This canonical problem encompasses coupling to an open spherical cavity and scattering from a spherical reflector. In contrast to the closed sphere problem, the electromagnetic boundary conditions couple the TE and TM modes. A pseudodecoupling of the resultant dual series equations system into dual series problems for the TE and TM modal coefficients is accomplished by introducing terms that are proportional to the associated Legendre functions  $P_0^{-m}$ . The solutions of the TE and TM dual series problems require the further introduction of terms proportional to  $P_n^{-m}$ , where  $0 \leq n < m$ . These functions effectively complete the standard spherical harmonic basis set when an aperture is present and guarantee the satisfaction of Meixner's edge conditions. Having generated the modal coefficients, all desired electromagnetic quantities follow immediately. Numerical results for the currents induced on the open spherical shell and for the energy density of the field at its center are presented for the case of normal incidence.

## I. INTRODUCTION

The number of electromagnetic boundary value problems that can be solved exactly is rather small, especially in three dimensions. The desire and the need for these canonical problems, however, is very strong. They reveal the basic physics underlying the phenomena and help establish insights that can usually be extrapolated to more general situations. Moreover, they act as valuable test cases for general numerical approaches to related problems.

The scattering of an electromagnetic plane wave from a perfectly conducting closed sphere is probably the best known three-dimensional canonical scattering problem. Its generalization, the scattering of a plane wave from a perfectly conducting spherical shell with a circular aperture, is important from both theoretical and practical points of view. In particular, when the circular hole has a relatively small angular extent, this problem allows one to study the coupling of a wave from an external source through an aperture into an enclosed region. On the other hand, when the shell has a relatively small angular extent, the problem describes the scattering of a plane wave from a spherical reflector. A complete solution to this canonical mixed boundary value problem is given in this paper.

A Debye potential formulation is employed, but in contrast to standard treatments in spherical geometries, the associated Legendre polynomials of negative order ( $P_n^{-m}$ ,  $n \geq m$ ) are utilized for the modal expansions. Enforcement of the electromagnetic boundary conditions leads to a coupled set of dual series equations for the TE and TM modal coefficients of each azimuthal mode. A pseudodecoupling ansatz is developed to allow separate treatment of the TE and TM

dual series systems. It requires the introduction of terms proportional to the associated Legendre polynomials  $P_0^{-m}$  ( $m$  being the azimuthal mode number) which are homogeneous solutions of the boundary condition equations. Solutions of the resulting "uncoupled" TE and TM dual series systems are given. They require the further introduction of terms proportional to the associated Legendre polynomials whose degree is less than its order:  $P_n^{-m}$ , where  $0 \leq n < m$ . These terms guarantee satisfaction of Meixner's edge conditions and effectively complete the spherical harmonic basis set in the presence of the aperture. Infinite systems of Fredholm equations of the second kind for the modal coefficients are obtained. A rigorous truncation procedure is given that leads to a straightforward numerical evaluation of those coefficients. Results for the currents induced on the open spherical shell and for an energy density ratio as a function of the fundamental parameter  $ka$  ( $2\pi \times$  radius/wavelength) are presented for the case of normal incidence. It is shown analytically that the behavior of the currents near the edge of the aperture are in agreement with Meixner's edge conditions; the graphical results further confirm this. The energy density scans highlight the resonance features of the coupling physics.

This paper is organized as follows. In Sec. II the coupled dual series systems are derived for the scattering of a general plane wave from an open spherical shell. The decoupling ansatz is presented in Sec. III, and the resulting TE and TM dual series systems are solved in Sec. IV. The results are then restricted to the normal incidence case in Sec. V. In Sec. VI the currents induced on the open spherical shell are given for various values of  $ka$ , aperture size, and the two allowed angles of incidence. Their modal structure is exhibited with a set of three-dimensional color figures. The energy density scans are discussed in Sec. VII.

<sup>a)</sup> Present address: Division 1265, Sandia National Laboratories, Albuquerque, New Mexico 87185.

There have been several reports of solutions to the normally incident case of the open spherical shell problem from both analytical<sup>1-10</sup> and numerical<sup>11-15</sup> points of view. In the numerical papers, various convergence problems and erroneous results are encountered. Of the analytical papers only Ref. 9 seems to lead to correct results for the scattering problem. Unfortunately, direct comparisons for that case are difficult because the dual series systems and their solutions (which were checked with our validation scheme) differ from those obtained here and no current or field values were calculated there. Moreover, the basic tenets of Ref. 9 appear to be restricted to the normal incidence case. The errors in Refs. 1-8 and 10 are either that the dual series systems were solved incorrectly or, more fundamentally, that the wrong dual series systems were solved. The latter stems from the erroneous assumption that the TE and TM dual series are completely decoupled. This error is identical to the one made by Meixner in his original Debye potential solution to the scattering of a plane wave from a circular hole in a perfectly conducting ground plane.<sup>16</sup> In analogy with our approach, Meixner corrected that error in Ref. 17 by introducing additional potentials that were homogeneous solutions of the equations resulting from enforcement of the electromagnetic boundary conditions. The coefficients of these potentials were chosen to guarantee that the fields satisfy the correct edge behavior, hence accounting for the presence of the aperture. The pseudodecoupling ansatz can be shown to be equivalent to a gauge transformation, which in analogy with Dirac string analyses, involves discontinuous potentials, the gauge conditions being identical to the pseudodecoupling constraint conditions.<sup>18</sup>

The results for normal incidence were closely compared with those generated with a general, numerical surface patch scattering code; and these comparisons were reported in Ref. 19. The agreement is excellent; and since that code has been validated with a variety of different scattering problems, this lends further credence to the validity of the solution presented below. The present work represents a generalization of related aperture coupling work<sup>20-24</sup> to three dimensions. A more detailed presentation is available.<sup>25</sup> It includes many complementary results that were omitted here simply because of length considerations.

## II. REDUCTION TO COUPLED DUAL SERIES PROBLEM

Consider the problem configuration shown in Fig. 1. A perfectly conducting open thin spherical shell is represented by the surface  $r = a$ ,  $0 \leq \theta < \theta_0$  in the spherical coordinate system  $(r, \theta, \phi)$  erected at the shell's center. The negative  $z$  axis of that system passes through the center of the aperture, the latter being defined as  $\{(r, \theta, \phi) | r = a \text{ and } \theta_0 < \theta \leq \pi\}$ . The opening angle of the aperture,  $\theta_{ap}$ , is defined simply as  $\theta_{ap} = \pi - \theta_0$ . The medium inside and outside the shell is free space. The unit vectors  $(\hat{r}, \hat{\theta}, \hat{\phi})$  are defined in the standard manner in the directions of positively increasing coordinate values.

Mathematically, we are seeking, for an arbitrary incident plane wave, the field scattered by the open spherical shell. This scattered field must satisfy the Sommerfeld radiation condition as  $r \rightarrow \infty$ . The total field (incident + scat-

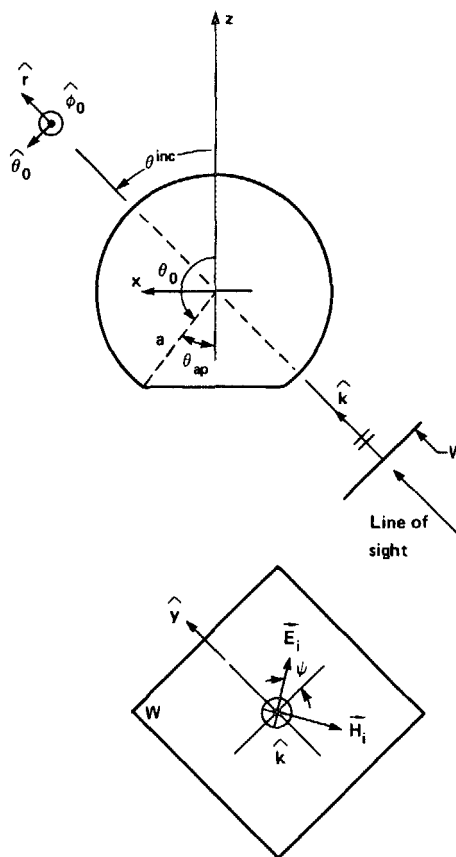


FIG. 1. Configuration of the scattering of an arbitrary plane wave from a spherical shell having a circular aperture.

tered) must satisfy (1) the electromagnetic conditions,  $E_{tan} = 0$  on the metallic shell and  $H_{tan}$  continuous over the aperture; and (2) Meixner's edge conditions,<sup>17,26</sup> i.e., the total energy of the field must be finite near the aperture rim.

### A. Debye expansions

A plane wave with electric field strength  $E_0$  is incident on the open spherical shell. It is characterized by a wave vector  $\mathbf{k}$ , which for convenience is assumed to lie in the  $xz$  plane; an incident angle  $\theta^{inc}$  with respect to the  $z$  axis so that  $\mathbf{k} \cdot \hat{z} = \cos \theta^{inc}$ ; and a polarization angle  $\psi$  between  $\mathbf{E}$  and the projection of the positive  $z$  axis on the incident wave front. The incident field has the form

$$\begin{bmatrix} \mathbf{E}^{inc} \\ Z_0 \mathbf{H}^{inc} \end{bmatrix} = -E_0 e^{i\mathbf{k} \cdot \mathbf{r}} \begin{bmatrix} (\cos \psi) \hat{\theta}_0 - (\sin \psi) \hat{\phi}_0 \\ (\sin \psi) \hat{\theta}_0 + (\cos \psi) \hat{\phi}_0 \end{bmatrix}, \quad (2.1)$$

where  $\hat{\theta}_0$  and  $\hat{\phi}_0$  are the incident polarization vectors and where, as throughout this paper, an  $e^{-i\omega t}$  time dependence has been assumed and suppressed. The free-space impedance  $Z_0$  is related to the free-space admittance  $Y_0 = (\epsilon/\mu)^{1/2}$  as  $Z_0 = Y_0^{-1}$  and the wave number  $k = \omega(\epsilon\mu)^{1/2}$ , where  $\epsilon$  and  $\mu$  denote, respectively, the permittivity and permeability of free space. The incident field parameters are indicated in Fig. 1. The incident electric field is polarized perpendicular to the edge of the aperture when  $\psi = 0$  and is polarized parallel to it when  $\psi = \pi/2$ . Since any incident plane wave can be reduced to a linear superposition of these waves, only they

TABLE I. The electric and magnetic field components in spherical coordinates in terms of Debye potentials.

$E_r = -\frac{1}{i\omega\epsilon} \left\{ \frac{\partial^2}{\partial r^2} + k^2 \right\} (r\Psi)$
$E_\theta = -\frac{1}{r \sin \theta} \frac{\partial}{\partial \phi} (r\Phi) - \frac{1}{(i\omega\epsilon)r} \frac{\partial^2}{\partial r \partial \theta} (r\Psi)$
$E_\phi = \frac{1}{r} \frac{\partial}{\partial \theta} (r\Phi) - \frac{1}{(i\omega\epsilon)r \sin \theta} \frac{\partial^2}{\partial r \partial \phi} (r\Psi)$
$H_r = -\frac{1}{i\omega\mu} \left\{ \frac{\partial^2}{\partial r^2} + k^2 \right\} (r\Phi)$
$H_\theta = \frac{1}{r \sin \theta} \frac{\partial}{\partial \phi} (r\Psi) - \frac{1}{(i\omega\mu)r} \frac{\partial^2}{\partial r \partial \theta} (r\Phi)$
$H_\phi = -\frac{1}{r} \frac{\partial}{\partial \theta} (r\Psi) - \frac{1}{(i\omega\mu)r \sin \theta} \frac{\partial^2}{\partial r \partial \phi} (r\Phi)$

need to be addressed explicitly. Moreover, with the intrinsic symmetry of the field components in Maxwell's equations, the  $\psi = \pi/2$  case is readily obtained from the  $\psi = 0$  case. Consequently, we restrict our considerations to the  $\psi = 0$  case with no loss in generality.

Following standard analyses of problems in a spherical-symmetric geometry, we employ a Debye potential formalism.<sup>27-29</sup> In particular, the electric and magnetic fields are expressed in terms of the two vector potentials  $\Phi\mathbf{r}$  and  $\Psi\mathbf{r}$  as

$$\mathbf{E} = -\text{curl}(\Phi\mathbf{r}) - (i\omega\epsilon)^{-1} \text{curl} \text{curl}(\Psi\mathbf{r}), \quad (2.2)$$

$$\mathbf{H} = +\text{curl}(\Psi\mathbf{r}) - (i\omega\mu)^{-1} \text{curl} \text{curl}(\Phi\mathbf{r}), \quad (2.3)$$

where the radial vector  $\mathbf{r} = r\hat{r}$ . Their components are given explicitly in Table I. The scalar functions  $\Phi$  and  $\Psi$  may represent any combination of the incident and scattered fields. The function  $\Phi$  defines the field TE with respect to  $r$ ,  $\Psi$  the field TM with respect to  $r$ . The descriptor "with respect to  $r$ " is assumed and suppressed throughout the rest of this paper.

The spherical wave expansion of the incident field (2.1) with  $\psi = 0$  given, for instance, in Ref. 30 or Ref. 31, can be generated with the Debye scalar potentials,  $\Phi^{\text{inc}}$  and  $\Psi^{\text{inc}}$ , defined below. Since the scattered potentials,  $\Phi^s$  and  $\Psi^s$ , assume an analogous form, we have

$$\begin{pmatrix} \Phi^{\text{inc}} \\ \Phi^s \end{pmatrix} = -E_0 \sum_{m=0}^{\infty} \begin{pmatrix} \Phi_m^{\text{inc}} \\ \Phi_m^s \end{pmatrix} \sin m\phi, \quad (2.4)$$

$$\begin{pmatrix} \Psi^{\text{inc}} \\ \Psi^s \end{pmatrix} = Y_0 E_0 \sum_{m=0}^{\infty} \begin{pmatrix} \Psi_m^{\text{inc}} \\ \Psi_m^s \end{pmatrix} \cos m\phi, \quad (2.5)$$

where the azimuthal modal coefficients

$$\Phi_m^{\text{inc}} = \sum_{n=m}^{\infty} \gamma_{mn} \left[ \frac{m P_n^m(\cos \theta^{\text{inc}})}{\sin \theta^{\text{inc}}} \right] j_n(kr) P_n^{-m}(\cos \theta), \quad (2.6)$$

$$\Psi_m^{\text{inc}} = \sum_{n=m}^{\infty} \gamma_{mn} \left[ -\left( \frac{\partial}{\partial \theta} P_n^m \right) (\cos \theta^{\text{inc}}) \right] \times j_n(kr) P_n^{-m}(\cos \theta), \quad (2.7)$$

$$\Phi_m^s = \sum_{n=m}^{\infty} A_{mn} P_n^{-m}(\cos \theta) \begin{cases} j_n(kr) h_n(ka) & (r < a), \\ j_n(ka) h_n(kr) & (r > a), \end{cases} \quad (2.8)$$

$$\Psi_m^s = \sum_{n=m}^{\infty} B_{mn} P_n^{-m}(\cos \theta) \times \begin{cases} j_n(kr) [k a h_n(ka)]' & (r < a), \\ [k a j_n(ka)]' h_n(kr) & (r > a), \end{cases} \quad (2.9)$$

where

$$\gamma_{mn} = (-1)^{m+1} i^n [(2n+1)/n(n+1)] \epsilon_m (1 - \delta_{0n}). \quad (2.10)$$

The terms  $j_n$  and  $h_n$  are, respectively, the spherical Bessel and Hankel (of the first kind) functions of order  $n$ . The associated Legendre polynomials of degree  $n$ , order  $\pm m$ , are denoted by  $P_n^{\pm m}$ . The prime in an expression  $[x f_n(x)]'$  denotes the derivative with respect to  $x$ . The term  $\epsilon_m = 2$  for  $m \neq 0$  and  $\epsilon_m = 1$  for  $m = 0$ . Kronecker's delta  $\delta_{ij} = 0$  for  $i \neq j$  and  $\delta_{ij} = 1$  for any  $j$ . Because the term corresponding to both  $m = 0$  and  $n = 0$  is identically zero in the incident field, we set the corresponding scattered potential coefficients identically to zero:  $A_{00} \equiv B_{00} \equiv 0$ .

These representations of the interior and exterior scattered potentials have been chosen so that  $\Phi^s$  and  $\partial_r(r\Psi^s)$  are continuous at  $r = a$ , thereby ensuring the continuity of the tangential scattered electric field components  $E_\theta$  and  $E_\phi$  across that surface. The resultant fields satisfy Sommerfeld's radiation condition; the dependence of the scattered fields on  $h_n(kr)$  for  $r > a$  ensures their decay to zero as  $r \rightarrow \infty$ . The modal coefficients  $A_{mn}$  and  $B_{mn}$  are the quantities that must be determined.

## B. Electromagnetic boundary conditions

The electromagnetic boundary conditions:  $E_{\text{tan}} = 0$  on the metal and  $H_{\text{tan}}$  continuous in the aperture, are now enforced. Referring to Table I,  $E_\theta = 0$  on the metal if

$$\sum_{m=0}^{\infty} \cos m\phi \left\{ \frac{m}{\sin \theta} (\Phi_m^{\text{inc}} + \Phi_m^s) - \frac{Y_0}{i\omega\epsilon} \frac{1}{r} \frac{\partial^2}{\partial r \partial \theta} [r(\Psi_m^{\text{inc}} + \Psi_m^s)] \right\}_{r=a} = 0,$$

$E_\phi = 0$  on the metal if

$$\sum_{m=0}^{\infty} \sin m\phi \left\{ -\frac{\partial}{\partial \theta} (\Phi_m^{\text{inc}} + \Phi_m^s) + \frac{Y_0}{i\omega\epsilon} \frac{m}{r \sin \theta} \frac{\partial}{\partial r} [r(\Psi_m^{\text{inc}} + \Psi_m^s)] \right\}_{r=a} = 0,$$

$H_\theta$  is continuous across the aperture if for  $\epsilon \rightarrow 0$

$$\sum_{m=0}^{\infty} \sin m\phi \left\{ \frac{-m Y_0}{\sin \theta} (\Psi_m^{\text{inc}} + \Psi_m^s) + \frac{1}{i\omega\mu} \frac{1}{r} \frac{\partial^2}{\partial r \partial \theta} [r(\Phi_m^{\text{inc}} + \Phi_m^s)] \right\}_{r=a-\epsilon}^{r=a+\epsilon} = 0,$$

and  $H_\phi$  is continuous across the aperture if for  $\epsilon \rightarrow 0$

$$\sum_{m=0}^{\infty} \cos m\phi \left\{ -Y_0 \frac{\partial}{\partial \theta} (\Psi_m^{\text{inc}} + \Psi_m^s) + \frac{1}{i\omega\mu} \frac{m}{r \sin \theta} \frac{\partial}{\partial r} [r(\Phi_m^{\text{inc}} + \Phi_m^s)] \right\}_{r=a-\epsilon}^{r=a+\epsilon} = 0.$$

Because the azimuthal eigenfunctions  $\sin m\phi$  and  $\cos m\phi$  form an orthogonal set over  $[0, 2\pi]$ , these conditions must be satisfied on a mode by mode basis. They require satisfac-

tion of the following coupled set of dual series equations for the modal coefficients  $A_{mn}$  and  $B_{mn}$ :

$$ika \sum_{n=m}^{\infty} \{A_{mn} j_n(ka) h_n(ka) - f_{mn}\} m P_n^{-m}(\cos \theta) \\ = \sin \theta \partial_{\theta} \sum_{n=m}^{\infty} \{B_{mn} [kaj_n(ka)]' [kah_n(ka)]' - g_{mn}\} P_n^{-m}(\cos \theta) \quad (0 \leq \theta < \theta_0), \quad (2.11a)$$

$$\sum_{n=m}^{\infty} A_{mn} m P_n^{-m}(\cos \theta) = -ika \sin \theta \partial_{\theta} \sum_{n=m}^{\infty} B_{mn} P_n^{-m}(\cos \theta) \quad (\theta_0 < \theta \leq \pi), \quad (2.11b)$$

$$ika \sin \theta \partial_{\theta} \sum_{n=m}^{\infty} \{A_{mn} j_n(ka) h_n(ka) - f_{mn}\} P_n^{-m}(\cos \theta) \\ = \sum_{n=m}^{\infty} \{B_{mn} [kaj_n(ka)]' [kah_n(ka)]' - g_{mn}\} m P_n^{-m}(\cos \theta) \quad (0 \leq \theta < \theta_0), \quad (2.12a)$$

$$\sin \theta \partial_{\theta} \sum_{n=m}^{\infty} A_{mn} P_n^{-m}(\cos \theta) = -ika \sum_{n=m}^{\infty} B_{mn} m P_n^{-m}(\cos \theta) \quad (\theta_0 < \theta \leq \pi), \quad (2.12b)$$

where

$$f_{mn} = \gamma_{mn} [m P_n^m(\cos \theta^{\text{inc}}) / \sin \theta^{\text{inc}}] j_n(ka), \quad (2.13)$$

$$g_{mn} = \gamma_{mn} \left[ \left( \frac{\partial}{\partial \theta} P_n^m \right) (\cos \theta^{\text{inc}}) \right] [kaj_n(ka)]'. \quad (2.14)$$

Equations (2.11) result from the  $E_{\theta}$  and  $H_{\phi}$  boundary conditions and are naturally paired because their  $\theta$  and  $\phi$  dependencies are the same; Eqs. (2.12) result from the  $E_{\phi}$  and  $H_{\theta}$  boundary conditions. The absence of any spherical Bessel or Hankel function in (2.11b) and (2.12b) results from application of the modified Wronskian relation

$$j_n(x) [xh_n(x)]' - h_n(x) [xj_n(x)]' = i/x. \quad (2.15)$$

Note that  $\partial_{\theta} \equiv \partial / \partial \theta$ .

### III. PSEUDODECOUPLING ANSATZ

In the problem of plane wave scattering from a solid sphere it is known<sup>27</sup> that the TE and TM portions of the problem may be decoupled. Satisfaction of independent boundary conditions applied directly to the TE and TM Debye potentials leads to series defined over the entire  $\theta$  interval,  $[0, \pi]$ ; and orthogonality arguments then produce a complete decoupling. Introducing the hole results in mixed boundary conditions over partial  $\theta$  intervals and a coupling of the TE and TM modes. Nonetheless, one might anticipate some form of TE/TM decoupling even in this case if the proper set of basis functions were employed.

Consider the associated Legendre functions of negative order  $P_n^{-m}(\cos \theta)$ . For all  $n$  and  $m$  they are known independent solutions to Legendre's equation<sup>32-35</sup>:

$$\mathcal{L}_{\theta} P_n^{-m}(\cos \theta) = -n(n+1)(\sin^2 \theta) P_n^{-m}(\cos \theta), \quad (3.1)$$

where the operator

$$\mathcal{L}_{\theta} \equiv (\sin \theta \partial_{\theta})(\sin \theta \partial_{\theta}) - m^2. \quad (3.2)$$

They have the integral representations<sup>32</sup>

$$P_n^{-m}(\cos \theta) = (-1)^m \left( \frac{2}{\pi} \right)^{1/2} \frac{\csc^m \theta}{\Gamma(m + \frac{1}{2})} \\ \times \int_0^{\theta} \frac{\cos[(n + \frac{1}{2})t] dt}{[\cos t - \cos \theta]^{1/2-m}}. \quad (3.3)$$

Note that our definition of  $P_n^{-m}$  differs from that in Ref. 32 by the factor  $(-1)^m$ . The related functions

$$\bar{P}_n^{-m}(\cos \theta) \\ \equiv (-1)^{n+m} P_n^{-m}(\cos(\pi - \theta)) \\ = (-1)^m \left( \frac{2}{\pi} \right)^{1/2} \frac{\csc^m \theta}{\Gamma(m + \frac{1}{2})} \\ \times \int_{\theta}^{\pi} \frac{\sin[(n + \frac{1}{2})t] dt}{[\cos \theta - \cos t]^{1/2-m}} \quad (3.4)$$

also satisfy (3.1). For  $n \geq m$ , these functions are identical by the standard symmetry relation

$$\bar{P}_n^{-m}(\cos \theta) \equiv (-1)^{n+m} P_n^{-m}(\cos(\pi - \theta)) \\ \equiv P_n^{-m}(\cos \theta). \quad (3.5)$$

However, for  $0 \leq n < m$ , this relation no longer holds true. In particular,  $P_n^{-m}(\cos \theta)$  is finite at  $\theta = 0$  but infinite at  $\theta = \pi$ . On the other hand,  $\bar{P}_n^{-m}(\cos \theta)$  is finite at  $\theta = \pi$  but infinite at  $\theta = 0$ . This behavior is immediately apparent for  $n = 0$  where

$$P_0^{-m}(\cos \theta) = [(-1)^m / m!] \tan^m(\theta/2), \quad (3.6a)$$

$$\bar{P}_0^{-m}(\cos \theta) = (1/m!) \cot^m(\theta/2). \quad (3.6b)$$

Return now to the dual series systems (2.11) and (2.12). They are self-consistent if

$$\mathcal{L}_{\theta} \sum_{n=m}^{\infty} \{A_{mn} j_n(ka) h_n(ka) - f_{mn}\} P_n^{-m}(\cos \theta) = 0 \\ \mathcal{L}_{\theta} \sum_{n=m}^{\infty} \{B_{mn} [kaj_n(ka)]' [kah_n(ka)]' - g_{mn}\} P_n^{-m}(\cos \theta) = 0 \\ (0 \leq \theta < \theta_0);$$

$$\begin{aligned} \mathcal{L}_\theta \sum_{n=m}^{\infty} A_{mn} P_n^{-m}(\cos \theta) &= 0 \\ \mathcal{L}_\theta \sum_{n=m}^{\infty} B_{mn} P_n^{-m}(\cos \theta) &= 0 \end{aligned} \quad (\theta_0 < \theta \leq \pi).$$

Bounded homogeneous solutions of these equations are admissible and are proportional to  $P_0^{-m}(\cos \theta)$  and  $\bar{P}_0^{-m}(\cos \theta)$  over their respective intervals. Solutions to the TE dual series system

$$\sum_{n=m}^{\infty} \{A_{mn} j_n(ka) h_n(ka) - f_{mn}\} P_n^{-m}(\cos \theta) = \alpha_m P_0^{-m}(\cos \theta) \quad (0 \leq \theta < \theta_0), \quad (3.7a)$$

$$\sum_{n=m}^{\infty} A_{mn} P_n^{-m}(\cos \theta) = \bar{\alpha}_m \bar{P}_0^{-m}(\cos \theta) \quad (\theta_0 < \theta \leq \pi), \quad (3.7b)$$

for  $m \geq 1$  and to the TM dual series system

$$\sum_{n=m}^{\infty} \{B_{mn} [kaj_n(ka)]' [kah_n(ka)]' - g_{mn}\} P_n^{-m}(\cos \theta) = \beta_m P_0^{-m}(\cos \theta) \quad (0 \leq \theta < \theta_0), \quad (3.8a)$$

$$\sum_{n=m}^{\infty} B_{mn} P_n^{-m}(\cos \theta) = \bar{\beta}_m \bar{P}_0^{-m}(\cos \theta) \quad (\theta_0 < \theta \leq \pi), \quad (3.8b)$$

for  $m \geq 0$  are therefore solutions to (2.11) and (2.12) provided that the "decoupling" constants  $\alpha_m$ ,  $\beta_m$ ,  $\bar{\alpha}_m$ , and  $\bar{\beta}_m$  are constrained by those coupled dual series equations. Since

$$\sin \theta \partial_\theta \begin{bmatrix} P_0^{-m}(\cos \theta) \\ \bar{P}_0^{-m}(\cos \theta) \end{bmatrix} = m \begin{bmatrix} +P_0^{-m}(\cos \theta) \\ -\bar{P}_0^{-m}(\cos \theta) \end{bmatrix}, \quad (3.9)$$

the required constraint relations for  $m \geq 1$  are simply

$$\beta_m = ika\alpha_m, \quad (3.10)$$

$$\bar{\alpha}_m = ika\bar{\beta}_m. \quad (3.11)$$

There is no  $m = 0$  constraint relation because there is no  $m = 0$  TE dual series equation.

Consequently, although the TE and TM portions of the dual series systems (2.11) and (2.12) have been decoupled, the TE and TM modal coefficients are still coupled through these "decoupling" constant constraint relations. This explains the connotation "pseudodecoupling ansatz." The solutions to the TE and TM dual series systems (3.7) and (3.8) subject to the constraints (3.10) and (3.11) comprise the desired result.

#### IV. TE AND TM DUAL SERIES SOLUTIONS

The TE and TM dual series systems can be reduced to more manageable and physically revealing forms with several manipulations. First, by introducing for  $n \geq 1$  the functions  $\chi_n^\phi$  and  $\chi_n^\psi$  so that

$$j_n(ka) h_n(ka) = (1 + \chi_n^\phi) / ika(2n + 1), \quad (4.1)$$

$$\begin{aligned} [kaj_n(ka)]' [kah_n(ka)]' \\ = -[n(n + 1) / ika(2n + 1)] (1 + \chi_n^\psi), \end{aligned} \quad (4.2)$$

the dual series systems (3.7) and (3.8) can be rewritten as

$$\begin{aligned} \sum_{n=m}^{\infty} \frac{A_{mn}}{n + \frac{1}{2}} (1 + \chi_n^\phi) P_n^{-m} \\ = 2ika\alpha_m P_0^{-m} + 2ika \sum_{n=m}^{\infty} f_{mn} P_n^{-m} \quad (0 \leq \theta < \theta_0), \end{aligned}$$

$$\sum_{n=m}^{\infty} A_{mn} P_n^{-m} = \bar{\alpha}_m \bar{P}_0^{-m} \quad (\theta_0 < \theta \leq \pi);$$

$$\begin{aligned} \sum_{n=m}^{\infty} B_{mn} \frac{n(n + 1)}{n + \frac{1}{2}} (1 + \chi_n^\psi) P_n^{-m} \\ = -2ika\beta_m P_0^{-m} - 2ika \\ \times \sum_{n=m}^{\infty} g_{mn} P_n^{-m} \quad (0 \leq \theta < \theta_0), \end{aligned}$$

$$\sum_{n=m}^{\infty} B_{mn} P_n^{-m} = \bar{\beta}_m \bar{P}_0^{-m} \quad (\theta_0 < \theta \leq \pi).$$

As shown in the Appendix, in the quasistatic limit

$$\lim_{ka \rightarrow 0} \chi_n^\phi = 0 + \mathcal{O}((ka)^2), \quad (4.3a)$$

$$\lim_{ka \rightarrow 0} \chi_n^\psi = 0 + \mathcal{O}((ka)^2). \quad (4.3b)$$

Thus, in analogy with the dual series treatments of the two-dimensional slit cylinder coupling problems given in Refs. 20–23, the static terms have been extracted. These TE and TM dual series must now be solved subject to Meixner's edge conditions; i.e., one must account for the singular behavior near the rim of the aperture required by the finite energy condition. The large  $n$  behavior of the solution coefficients is responsible for this edge behavior. Since, as shown in the Appendix, for large values of this index

$$\lim_{n \rightarrow \infty} \chi_n^\phi \sim \mathcal{O}(n^{-2}), \quad (4.4a)$$

$$\lim_{n \rightarrow \infty} \chi_n^\psi \sim \mathcal{O}(n^{-2}), \quad (4.4b)$$

the terms proportional to  $\chi_n^\phi$  and  $\chi_n^\psi$  are of order  $n^{-2}$  smaller than the static pieces. To enhance the isolation of the large index behavior in the TM systems, we introduce the additional functions

$$\begin{aligned} \tilde{\chi}_n^\psi &= n(n + 1)(1 + \chi_n^\psi) / (n + \frac{1}{2})^2 - 1 \\ &= -\{1 + [4ika / (2n + 1)] [kaj_n(ka)]' \\ &\quad \times [kah_n(ka)]'\}, \end{aligned} \quad (4.5)$$

which exhibit the limiting behaviors

$$\lim_{n \rightarrow \infty} \tilde{\chi}_n^\psi = \mathcal{O}(n^{-2}) \quad \text{and} \quad \lim_{ka \rightarrow 0} \tilde{\chi}_n^\psi = -(2n + 1)^{-2}, \quad (4.6)$$

and rewrite the TM dual series systems as

$$\begin{aligned} \sum_{n=m}^{\infty} B_{mn} \left(n + \frac{1}{2}\right) (1 + \tilde{\chi}_n^\psi) P_n^{-m} \\ = -2ika\beta_m P_0^{-m} - 2ika \end{aligned} \quad (4.7)$$

$$\times \sum_{n=m}^{\infty} g_{mn} P_n^{-m} \quad (0 \leq \theta < \theta_0),$$

$$\sum_{n=m}^{\infty} B_{mn} P_n^{-m} = \bar{\beta}_m \bar{P}_0^{-m} \quad (\theta_0 < \theta \leq \pi).$$

We then treat the terms proportional to  $\chi_n^\phi$  in (4.2) and  $\tilde{\chi}_n^\psi$  in (4.7) as forcing terms by moving them to the right-hand sides. This isolates the pieces responsible for the singularities near the rim of the aperture on the left-hand sides. Defining the forcing terms

$$F_{mn} = 2ikaf_{mn} - [A_{mn}/(n + \frac{1}{2})]\chi_n^\phi, \quad (4.8a)$$

$$G_{mn} = -2ikag_{mn} - \tilde{\chi}_n^\psi(n + \frac{1}{2})B_{mn}, \quad (4.8b)$$

the TE and TM dual series systems for  $m \geq 1$  become

$$\begin{aligned} \sum_{n=m}^{\infty} \frac{A_{mn}}{n + \frac{1}{2}} P_n^{-m} \\ = 2ika\alpha_m P_0^{-m} + \sum_{n=m}^{\infty} F_{mn} P_n^{-m} \quad (0 \leq \theta < \theta_0), \end{aligned} \quad (4.9a)$$

$$\sum_{n=m}^{\infty} A_{mn} \bar{P}_n^{-m} = \bar{\alpha}_m \bar{P}_0^{-m} \quad (\theta_0 < \theta \leq \pi); \quad (4.9b)$$

$$\begin{aligned} \sum_{n=m}^{\infty} B_{mn} \left(n + \frac{1}{2}\right) P_n^{-m} \\ = -2ika\beta_m P_0^{-m} + \sum_{n=m}^{\infty} G_{mn} P_n^{-m} \quad (0 \leq \theta < \theta_0), \end{aligned} \quad (4.10a)$$

$$\sum_{n=m}^{\infty} B_{mn} \bar{P}_n^{-m} = \bar{\beta}_m \bar{P}_0^{-m} \quad (\theta_0 < \theta \leq \pi). \quad (4.10b)$$

Equation (3.5) has been invoked to convert the  $P_n^{-m}$  to their duals  $\bar{P}_n^{-m}$  over the aperture interval. This form of the dual series systems strongly suggests the solution process we introduce below.

For  $m = 0$ , the TM dual series becomes

$$\begin{aligned} \sum_{n=1}^{\infty} B_{0n} \left(n + \frac{1}{2}\right) P_n \\ = -2ika\beta_0 + \sum_{n=1}^{\infty} G_{0n} P_n \quad (0 \leq \theta < \theta_0), \end{aligned} \quad (4.10a')$$

$$\sum_{n=1}^{\infty} B_{0n} P_n = \bar{\beta}_0 \quad (\theta_0 < \theta \leq \pi), \quad (4.10b')$$

since  $P_n^0 \equiv P_n$ , Legendre's polynomial, and  $B_{00} \equiv g_{00} \equiv 0$ .

The singular behavior of the fields near the aperture rim ( $\theta = \theta_0$ ) is reflected in the corresponding behavior of the current components

$$\begin{aligned} J_\theta(\theta, \phi) &= H_\phi^<(a, \theta, \phi) - H_\phi^>(a, \theta, \phi) \\ &= \left[ \frac{-Y_0 E_0}{(ka)^2} \right] \sum_{m=0}^{\infty} \cos m\phi \sum_{n=m}^{\infty} \left\{ A_{mn} \left[ \frac{m P_n^{-m}(\cos \theta)}{\sin \theta} \right] + ika B_{mn} \left[ \frac{\partial}{\partial \theta} P_n^{-m}(\cos \theta) \right] \right\}, \end{aligned} \quad (4.11a)$$

$$\begin{aligned} J_\phi(\theta, \phi) &= H_\theta^<(a, \theta, \phi) - H_\theta^>(a, \theta, \phi) \\ &= \left[ \frac{+Y_0 E_0}{(ka)^2} \right] \sum_{m=0}^{\infty} \sin m\phi \sum_{n=m}^{\infty} \left\{ A_{mn} \left[ \frac{\partial}{\partial \theta} P_n^{-m}(\cos \theta) \right] + ika B_{mn} \left[ \frac{m P_n^{-m}(\cos \theta)}{\sin \theta} \right] \right\}, \end{aligned} \quad (4.11b)$$

where, for instance,  $H_\phi^<(a, \theta, \phi)$  [ $H_\phi^>(a, \theta, \phi)$ ] is the  $\phi$  component of the magnetic field for  $r < a$  ( $r > a$ ) evaluated at  $r = a$ . Because the angular dependence of the terms depending on  $A_{mn}$  and  $B_{mn}$  in these expressions is distinct, we use it to guide our constructions of the physically correct TE and TM solutions. These solutions require several summation formulas:

$$\sum_{n=0}^{\infty} P_n^{-m}(\cos \theta) \partial_{\theta_0}^j \cos\left(n + \frac{1}{2}\right)\theta_0 = 0 \quad (0 \leq \theta < \theta_0), \quad (4.12a)$$

$$\sum_{n=0}^{\infty} \bar{P}_n^{-m}(\cos \theta) \partial_{\theta_0}^j \sin\left(n + \frac{1}{2}\right)\theta_0 = 0 \quad (\theta_0 < \theta \leq \pi), \quad (4.12b)$$

$$\sum_{n=0}^{\infty} P_n^{-m}(\cos \theta) \frac{\sin(n + \frac{1}{2})\theta_0}{n + \frac{1}{2}} = \sum_{n=0}^{\infty} \frac{(-1)^n P_n^{-m}(\cos \theta)}{n + \frac{1}{2}} \quad (0 \leq \theta < \theta_0), \quad (4.12c)$$

$$\sum_{n=0}^{\infty} \bar{P}_n^{-m}(\cos \theta) \frac{\cos(n + \frac{1}{2})\theta_0}{n + \frac{1}{2}} = \sum_{n=0}^{\infty} \frac{\bar{P}_n^{-m}(\cos \theta)}{n + \frac{1}{2}} \quad (\theta_0 < \theta \leq \pi), \quad (4.12d)$$

$$\sum_{n=0}^{\infty} P_n^{-m}(\cos \theta) \frac{\cos(n + \frac{1}{2})\theta_0}{(n + \frac{1}{2})^2} = (\pi - \theta_0) \sum_{n=0}^{\infty} (-1)^n \frac{P_n^{-m}(\cos \theta)}{n + \frac{1}{2}} \quad (0 \leq \theta < \theta_0), \quad (4.12e)$$

which are derived from the basic expressions [see Ref. 36, Eq. (3.71)]

$$\sum_{n=0}^{\infty} P_n^{-m}(\cos \theta) \cos\left(n + \frac{1}{2}\right)\psi = \begin{cases} \frac{(-1)^m (\pi/2)^{1/2} [\cos \psi - \cos \theta]^{m-1/2}}{\Gamma(m + \frac{1}{2}) \sin^m \theta} & (0 < \psi < \theta), \\ 0 & (\theta < \psi < \pi) \end{cases} \quad (4.13a)$$

(modified to our sign convention) and its dual

$$\sum_{n=0}^{\infty} \bar{P}_n^{-m}(\cos \theta) \sin\left(n + \frac{1}{2}\right)\psi = \begin{cases} 0 & (0 < \psi < \theta), \\ \frac{(\pi/2)^{1/2} [\cos \theta - \cos \psi]^{m-1/2}}{\Gamma(m + 1/2) \sin^m \theta} & (\theta < \psi < \pi), \end{cases} \quad (4.13b)$$



and the identities (see Ref. 36, 1.19 and 1.16)

$$\sum_{n=0}^{\infty} (-1)^n \frac{\cos(n + \frac{1}{2})t}{n + \frac{1}{2}} = \frac{\pi}{2} \quad (0 \leq t < \pi), \quad (4.14a)$$

$$\sum_{n=0}^{\infty} \frac{\sin(n + \frac{1}{2})t}{n + \frac{1}{2}} = \frac{\pi}{2} \quad (0 \leq t < \pi). \quad (4.14b)$$

### A. TE dual series solution

We would like to reduce the associated Legendre function dual series system (4.9) to one in sines and cosines. This conversion would appear to be straightforward with the representations (3.3) and (3.4) and with an interchange of the summations and integrations. However, consider Meixner's edge conditions (see Ref. 26, Sec. 9.2), which, when applied to the field generated by the TE Debye potential, imply that as the edge is approached along the surface  $r = a$ ,

$$H_{\phi}^s(a, \theta, \phi) \sim \partial_r (r \Phi_m^s)|_{r=a} \sim (\theta_0 - \theta)^{+1/2},$$

$$H_{\theta}^s(a, \theta, \phi) \sim \partial_r \partial_{\theta} (r \Phi_m^s)|_{r=a} \sim (\theta_0 - \theta)^{-1/2}.$$

The corresponding portions of  $J_{\theta}$  and  $J_{\phi}$  near  $\theta = \theta_0$  must behave, respectively, as

$$\sum_{n=m}^{\infty} A_{mn} P_n^{-m}(\cos \theta) \sim (\theta_0 - \theta)^{+1/2}, \quad (4.15a)$$

$$\sum_{n=m}^{\infty} A_{mn} \partial_{\theta} P_n^{-m}(\cos \theta) \sim (\theta_0 - \theta)^{-1/2}, \quad (4.15b)$$

where  $m \geq 1$ . Analogously, the  $\theta$  dependency of Eq. (4.9a) near  $\theta = \theta_0$  differs from that of (4.9b) by  $(\theta_0 - \theta)^{+1}$ . The factor  $(n + \frac{1}{2})^{-1}$  in (4.9a) is responsible for this difference. Thus, with [see Ref. 37, (8.10.7) and (6.1.37)]

$$\lim_{n \rightarrow \infty} P_n^{-m}(\cos \theta) \sim n^{-m+1/2}$$

$$\times \frac{\cos[(n + \frac{1}{2})\theta - m\pi/2 - \pi/4]}{(\pi \sin \theta/2)^{1/2}}$$

and (4.13), Meixner's conditions are satisfied if

$$\lim_{n \rightarrow \infty} A_{mn} \sim n^{m-1}. \quad (4.16)$$

The simple summation-integration interchange is then not directly permitted because it will introduce terms that are proportional to delta functions and their derivatives; i.e., with (1.135) from Ref. 36 one finds, for instance, that near  $\theta = \theta_0$

$$\sum_{n=0}^{\infty} n^j \cos\left(n + \frac{1}{2}\right)\theta \cos\left(n + \frac{1}{2}\right)\theta_0 \sim \delta^{(j)}(\theta_0 - \theta),$$

the  $j$ th derivative of the Dirac distribution. Interchange can be accomplished by preconditioning the dual series as follows.

We need to introduce terms into (4.9) that will cancel the potential delta function contributions. This is accomplished with Eqs. (4.12)–(4.14). In particular, we define the modified solution coefficients

$$\tilde{A}_{mn} = \sum_{j=0}^{m-1} a_{mj} \partial_{\theta_0}^j \sin\left(n + \frac{1}{2}\right)\theta_0 + \begin{cases} A_{mn} & (n \geq m \geq 1), \\ 0 & (0 \leq n < m), \end{cases} \quad (4.17)$$

and the modified forcing term coefficients

$$\tilde{F}_{mn} = F_{mn} \quad (n \geq m = 1), \quad (4.18a)$$

$$\tilde{F}_{mn} = \frac{-(ka)^2}{2} \sum_{j=0}^{m-2} a_{m(j+1)} \partial_{\theta_0}^j \frac{\cos(n + \frac{1}{2})\theta_0}{(n + \frac{1}{2})^2}$$

$$+ \begin{cases} F_{mn} & (n \geq m \geq 2), \\ 0 & (0 \leq n < m), \end{cases} \quad (4.18b)$$

so that the TE dual series systems for  $m \geq 1$  become

$$\sum_{n=0}^{\infty} \frac{\tilde{A}_{mn}}{n + \frac{1}{2}} P_n^{-m}$$

$$= 2ika\alpha_m P_0^{-m} + (a_{m0} - \kappa_m^E) \sum_{n=0}^{\infty} \frac{(-1)^n P_n^{-m}}{n + \frac{1}{2}}$$

$$+ \sum_{n=0}^{\infty} \tilde{F}_{mn} P_n^{-m} \quad (0 \leq \theta < \theta_0), \quad (4.19a)$$

$$\sum_{n=0}^{\infty} \tilde{A}_{mn} \bar{P}_n^{-m} = \bar{\alpha}_m \bar{P}_0^{-m} \quad (\theta_0 < \theta \leq \pi). \quad (4.19b)$$

The constants

$$\kappa_m^E = \frac{-(ka)^2}{2} \begin{cases} 0, & \text{for } m = 1, \\ (\pi - \theta_0)a_{m1}, & \text{for } m = 2, \\ (\pi - \theta_0)a_{m1} - a_{m2}, & \text{for } m \geq 3. \end{cases} \quad (4.20)$$

The additional unknown coefficients  $a_{mj}$  ( $j = 0, 1, \dots, m-1$ ) provide the extra degrees of freedom needed to remove the unphysical singularities and permit the desired summation-integration interchange. In particular, their values will be fixed by our solution process so that for any  $m \geq 1$

$$\lim_{n \rightarrow \infty} \tilde{A}_{mn} \sim \mathcal{O}(n^{-1}), \quad (4.17')$$

$$\lim_{n \rightarrow \infty} \tilde{F}_{mn} \sim \mathcal{O}(n^{-3}). \quad (4.18')$$

Note that it can be inferred from the form of (4.19) that we have completed the basis function set for this open geometry by including the associated Legendre polynomials  $P_n^{-m}$  and  $\bar{P}_n^{-m}$  for  $0 \leq n < m$ .

Inserting (3.3) and (3.4) into (4.19) and interchanging the summations and integrations, the desired TE dual series are generated:

$$\sum_{n=0}^{\infty} \frac{\tilde{A}_{mn}}{n + \frac{1}{2}} \cos\left(n + \frac{1}{2}\right)t$$

$$= \frac{\pi}{2} (a_{m0} - \kappa_m^E) + 2ika\alpha_m \cos \frac{t}{2}$$

$$+ \sum_{n=0}^{\infty} \tilde{F}_{mn} \cos\left(n + \frac{1}{2}\right)t \quad (0 \leq t < \theta_0), \quad (4.21a)$$

$$\sum_{n=0}^{\infty} \tilde{A}_{mn} \sin\left(n + \frac{1}{2}\right)t = \bar{\alpha}_m \sin \frac{t}{2} \quad (\theta_0 < t \leq \pi). \quad (4.21b)$$

A solution of (4.21) is constructed as in Refs. 20–24 by first making the metal and aperture equations display the same  $t$  dependence. Two possibilities exist: integrating (4.21b) or differentiating (4.21a). Only the former guarantees satisfaction of (4.17'). Applying  $\int_0^t dt$  to (4.21b) leads to the dual series system

$$\sum_{n=0}^{\infty} \frac{\tilde{A}_{mn}}{n + \frac{1}{2}} \cos\left(n + \frac{1}{2}\right)t = \begin{cases} \frac{\pi}{2} (a_{m0} - \kappa_m^E) + 2ika\alpha_m \cos \frac{t}{2} + \sum_{n=0}^{\infty} \tilde{F}_{mn} \cos\left(n + \frac{1}{2}\right)t & (0 \leq t < \theta_0), \\ 2\tilde{\alpha}_m \cos \frac{t}{2} & (\theta_0 < t \leq \pi). \end{cases} \quad (4.22)$$

Since the left-hand side of (4.22) is now defined over the entire  $[0, \pi]$  interval, Fourier inversion then yields the coefficients

$$\frac{\tilde{A}_{ml}}{l + \frac{1}{2}} = (a_{m0} - \kappa_m^E) \frac{\sin(l + \frac{1}{2})\theta_0}{l + \frac{1}{2}} + (2ika\alpha_m - 2\tilde{\alpha}_m)\Lambda_{0l}^E + 2\tilde{\alpha}_m\delta_{0l} + \sum_{n=0}^{\infty} \tilde{F}_{mn}\Lambda_{nl}^E \quad (l = 0, 1, \dots), \quad (4.23)$$

where the inversion terms

$$\Lambda_{nl}^E = \frac{2}{\pi} \int_0^{\theta_0} \cos\left[\left(n + \frac{1}{2}\right)\psi\right] \cos\left[\left(l + \frac{1}{2}\right)\psi\right] d\psi = \begin{cases} \frac{1}{\pi} \left[ \frac{\sin(n-l)\theta_0}{n-l} + \frac{\sin(n+l+1)\theta_0}{n+l+1} \right] & (n \neq l), \\ \frac{1}{\pi} \left[ \theta_0 + \frac{\sin(2l+1)\theta_0}{2l+1} \right] & (n = l). \end{cases} \quad (4.24)$$

Explicitly, (4.23) means

$$\begin{aligned} \sum_{j=1}^{m-1} a_{mj} \partial_{\theta_n}^j \frac{\sin(l + \frac{1}{2})\theta_0}{l + \frac{1}{2}} \\ = (2ika\alpha_m - 2\tilde{\alpha}_m)\Lambda_{0l}^E + 2\tilde{\alpha}_m\delta_{0l} - \kappa_m^E \frac{\sin(l + \frac{1}{2})\theta_0}{l + \frac{1}{2}} \\ + \sum_{n=0}^{\infty} \tilde{F}_{mn}\Lambda_{nl}^E \quad (l = 0, 1, \dots, m-1), \end{aligned} \quad (4.23')$$

$$\begin{aligned} \frac{A_{ml}}{l + \frac{1}{2}} = - \sum_{j=1}^{m-1} a_{mj} \partial_{\theta_n}^j \frac{\sin(l + \frac{1}{2})\theta_0}{l + \frac{1}{2}} \\ + (2ika\alpha_m - 2\tilde{\alpha}_m)\Lambda_{0l}^E + 2\tilde{\alpha}_m\delta_{0l} \\ - \kappa_m^E \frac{\sin(l + \frac{1}{2})\theta_0}{l + \frac{1}{2}} + \sum_{n=0}^{\infty} \tilde{F}_{mn}\Lambda_{nl}^E \\ (l = m, m+1, \dots). \end{aligned} \quad (4.23'')$$

Furthermore, inversion requires continuity of the right-hand side of (4.22) across  $t = \theta_0$ . This yields

$$\begin{aligned} a_{m0} = \kappa_m^E - \frac{2}{\pi} \left[ (2ika\alpha_m - 2\tilde{\alpha}_m) \cos \frac{\theta_0}{2} \right. \\ \left. + \sum_{n=0}^{\infty} \tilde{F}_{mn} \cos\left(n + \frac{1}{2}\right)\theta_0 \right]. \end{aligned} \quad (4.25)$$

The system (4.23) and (4.25) is an infinite system of linear equations for the TE solution coefficients. The first  $(m+1)$  of these, (4.23') and Eq. (4.25), would not have appeared without the introduction of the terms  $a_{mj}$ ,  $\alpha_n$ , and  $\tilde{\alpha}_n$ ; hence, the terms  $\cos(n + \frac{1}{2})\theta$  and  $\sin(n + \frac{1}{2})\theta$  for  $n = 0, 1, \dots, m-1$  into (4.21). The  $\Lambda_{0l}^E$  ( $l = 0, 1, \dots$ ) terms originate in this completion of the expansion. Moreover, since there are no solution coefficients  $A_{mn}$  ( $0 \leq n < m$ ) in those first  $(m+1)$  equations, we may view them as orthogonality relations. They determine the interchange coefficients  $a_{mj}$  ( $j = 1, \dots, m-1$ ) and a relation between the decoupling coefficients  $\alpha_m$  and  $\tilde{\alpha}_m$ . In a similar fashion, the TM case generates a relation between  $\beta_m$  and  $\tilde{\beta}_m$ . The remaining two degrees of freedom are determined by the constraint relations (3.10) and (3.11). The coefficient  $a_{m0}$  is defined by (4.25) and is coupled to all of the other  $a_{mj}$  ( $j = 1, \dots, m-1$ ) through the relation for  $(2ika\alpha_m - 2\tilde{\alpha}_m)$ . However, it does not contribute directly to the solution coefficients  $A_{ml}$  ( $l = m, m+1, \dots$ ). It only

provides that degree of freedom needed to insure continuity across the boundary between the metal and aperture intervals.

Equations (4.23) are solutions of the original dual series (4.9) subject to Meixner's edge conditions (4.15a). They are general solutions if these results are independent of the decoupling and interchange constants. This has been confirmed numerically for  $m = 1, 2, 3$ . A rich set of new associated Legendre polynomial identities is obtained from this validation process.<sup>25</sup> We consider explicitly only the  $m = 1$  relations since they are employed for the normal incidence case discussed below.

For  $m = 1$ , the solution system

$$\frac{A_{1l}}{l + \frac{1}{2}} = (2ika\alpha_1 - 2\tilde{\alpha}_1)\Lambda_{0l}^E + \sum_{n=1}^{\infty} F_{1n}\Lambda_{nl}^E \quad (l = 1, 2, \dots), \quad (4.26a)$$

$$0 = (2ika\alpha_1 - 2\tilde{\alpha}_1)\Lambda_{00}^E + 2\tilde{\alpha}_1 + \sum_{n=1}^{\infty} F_{1n}\Lambda_{n0}^E, \quad (4.26b)$$

gives the coefficients

$$\frac{A_{1l}}{l + \frac{1}{2}} = \sum_{n=1}^{\infty} F_{1n}\Gamma_{1,nl}^E + 2\tilde{\alpha}_1 L_{1l}^E \quad (l = 1, 2, \dots), \quad (4.27a)$$

where

$$\Gamma_{1,nl}^E = \Lambda_{nl}^E - \Lambda_{n0}^E \Lambda_{0l}^E / \Lambda_{00}^E, \quad (4.27b)$$

$$L_{1l}^E = -\Lambda_{0l}^E / \Lambda_{00}^E. \quad (4.27c)$$

Substituting these expressions into the  $m = 1$  versions of (4.9a) and (4.9b), the original dual series system is satisfied since on the metal ( $0 \leq \theta < \theta_0$ )

$$\sum_{l=1}^{\infty} \Gamma_{1,nl}^E P_l^{-1} - P_n^{-1} - L_{1n}^E P_0^{-1} = 0 \quad (n = 1, 2, \dots), \quad (4.28a)$$

$$\sum_{l=1}^{\infty} L_{1l}^E P_l^{-1} - \frac{\Lambda_{00}^E - 1}{\Lambda_{00}^E} P_0^{-1} = 0, \quad (4.28b)$$

and in the aperture ( $\theta_0 < \theta < \pi$ )

$$\sum_{l=1}^{\infty} (2l+1)\Gamma_{1,nl}^E P_l^{-1} = 0 \quad (n = 1, 2, \dots), \quad (4.29a)$$

$$\sum_{l=1}^{\infty} (2l+1)L_{1l}^E - \bar{P}_0^{-1} = 0. \quad (4.29b)$$

When evaluated over the metal interval ( $0 \leq \theta < \theta_0$ ), the left-hand sides of Eqs. (4.29) yield Eqs. (4.15).

### B. TM dual series solution

We proceed as in the TE case. Consider the dual series systems (4.10) and (4.10'). Meixner's edge conditions applied to the fields generated by the TM Debye potential imply that near  $\theta = \theta_0$

$$H_\theta^s(a, \theta, \phi) \sim (r\Psi_m^s)|_{r=a} \sim (\theta_0 - \theta)^{+3/2},$$

$$H_\phi^s(a, \theta, \phi) \sim \partial_\theta(r\Psi_m^s)|_{r=a} \sim (\theta_0 - \theta)^{+1/2}.$$

Thus, from (4.11a) and (4.11b), the portions of  $J_\theta$  and  $J_\phi$  generated by  $\Psi_m^s$  near the aperture edge behave, respectively, as

$$\sum_{n=m}^{\infty} B_{mn} P_n^{-m}(\cos \theta) \sim (\theta_0 - \theta)^{+3/2}, \quad (4.30a)$$

$$\sum_{n=m}^{\infty} B_{mn} \partial_\theta P_n^{-m}(\cos \theta) \sim (\theta_0 - \theta)^{+1/2}, \quad (4.30b)$$

where  $m \geq 0$ . Analogously, the  $\theta$  dependency of the metal equations near  $\theta = \theta_0$  differs from that of the aperture equations by  $(\theta_0 - \theta)^{-1}$ . The factor  $(n + \frac{1}{2})$  in the metal equations is responsible for this difference.

The requisite edge behavior (4.30a) is obtained if

$$\lim_{n \rightarrow \infty} B_{mn} \sim n^{m-2}. \quad (4.31)$$

Consider first the cases with  $m > 0$ . Anticipating the effects of the operator interchange, we introduce the modified coefficients

$$\tilde{B}_{mn} = \sum_{j=0}^{m-1} b_{mj} \partial_{\theta_0}^j \frac{\cos(n + \frac{1}{2})\theta_0}{n + \frac{1}{2}} + \begin{cases} B_{mn} & (n \geq m \geq 1), \\ 0 & (0 \leq n < m), \end{cases} \quad (4.32)$$

$$\sum_{n=0}^{\infty} \tilde{B}_{mn} \sin\left(n + \frac{1}{2}\right)t = \begin{cases} -4ika\beta_m \sin \frac{t}{2} + \sum_{n=0}^{\infty} \frac{\tilde{G}_{mn}}{n + \frac{1}{2}} \sin\left(n + \frac{1}{2}\right)t - \kappa_m^H \frac{\pi}{2} t & (0 \leq t < \theta_0), \\ \bar{\beta}_m \sin \frac{t}{2} + b_{m0} \frac{\pi}{2} & (\theta_0 < t \leq \pi). \end{cases} \quad (4.35)$$

Introducing the terms

$$\mu_l(\theta_0) = \int_0^{\theta_0} t \sin\left(l + \frac{1}{2}\right)t = \frac{\sin(l + \frac{1}{2})\theta_0}{(l + \frac{1}{2})^2} - \theta_0 \frac{\cos(l + \frac{1}{2})\theta_0}{l + \frac{1}{2}}, \quad (4.36)$$

$$\Lambda_{nl}^H = \frac{2}{\pi} \int_0^{\theta_0} \sin\left[\left(n + \frac{1}{2}\right)\psi\right] \sin\left[\left(l + \frac{1}{2}\right)\psi\right] d\psi = \begin{cases} \frac{1}{\pi} \left[ \frac{\sin(n-l)\theta_0}{n-l} - \frac{\sin(n+l+1)\theta_0}{n+l+1} \right] & (n \neq l), \\ \frac{1}{\pi} \left[ \theta_0 - \frac{\sin(2l+1)\theta_0}{2l+1} \right] & (n = l), \end{cases} \quad (4.37)$$

Fourier inversion leads to the coefficient expressions for  $m \geq 1$

$$\begin{aligned} \tilde{B}_{ml} = & -(4ika\beta_m + \bar{\beta}_m) \Lambda_{0l}^H + \bar{\beta}_m \delta_{0l} \\ & + b_{m0} \frac{\cos(l + \frac{1}{2})\theta_0}{l + \frac{1}{2}} - \kappa_m^H \mu_l(\theta_0) \\ & + \sum_{n=0}^{\infty} \frac{\tilde{G}_{mn}}{n + \frac{1}{2}} \Lambda_{nl}^H \quad (l = 0, 1, \dots). \end{aligned} \quad (4.38a)$$

the modified forcing term coefficients

$$\tilde{G}_{mn} = G_{mn} \quad (n \geq m \geq 1), \quad (4.33a)$$

$$\begin{aligned} \tilde{G}_{mn} = & \frac{(ka)^2}{2} \sum_{j=0}^{m-2} b_{m(j+1)} \partial_{\theta_0}^j \frac{\sin(n + \frac{1}{2})\theta_0}{n + \frac{1}{2}} \\ & + \begin{cases} G_{mn} & (n \geq m \geq 2), \\ 0 & (0 \leq n < m), \end{cases} \end{aligned} \quad (4.33b)$$

and the constants

$$\kappa_m^H = \begin{cases} 0, & \text{for } m = 1, \\ [(ka)^2/2]b_{m1}, & \text{for } m \geq 2. \end{cases} \quad (4.33c)$$

The interchange constants  $b_{mj}$  ( $j = 0, 1, \dots, m-1$ ) will be adjusted so that for all  $m \geq 1$

$$\lim_{n \rightarrow \infty} \tilde{B}_{mn} \sim \mathcal{O}(n^{-2}), \quad (4.32')$$

$$\lim_{n \rightarrow \infty} \tilde{G}_{mn} \sim \mathcal{O}(n^{-3}). \quad (4.33')$$

The dual series systems (4.10) become

$$\begin{aligned} & \sum_{n=0}^{\infty} \left(n + \frac{1}{2}\right) \tilde{B}_{mn} P_n^{-m} \\ & = -2ika\beta_m P_0^{-m} + \sum_{n=0}^{\infty} \tilde{G}_{mn} P_n^{-m} - \kappa_m^H \\ & \quad \times \sum_{n=0}^{\infty} \frac{(-1)^n P_n^{-m}}{n + \frac{1}{2}} \quad (0 \leq \theta < \theta_0), \end{aligned} \quad (4.34a)$$

$$\sum_{n=0}^{\infty} \tilde{B}_{mn} \bar{P}_n^{-m} = \bar{\beta}_m \bar{P}_0^{-m} + b_{m0} \sum_{n=0}^{\infty} \frac{\bar{P}_n^{-m}}{n + \frac{1}{2}} \quad (\theta_0 < \theta \leq \pi). \quad (4.34b)$$

Introducing (3.3) and (3.4), interchanging summations and integrations, and applying the operator  $\int_0^t dt$  to the resulting metal equation to attain similar  $t$  dependencies for both equations, the TM dual series prior to inversion are

These contain explicitly  $m$  orthogonality relations

$$\begin{aligned} & \sum_{j=0}^{m-1} b_{mj} \partial_{\theta_0}^j \frac{\cos(l + \frac{1}{2})\theta_0}{l + \frac{1}{2}} \\ & = -(4ika\beta_m + \bar{\beta}_m) \Lambda_{0l}^H + \bar{\beta}_m \delta_{0l} \\ & \quad + b_{m0} \frac{\cos(l + \frac{1}{2})\theta_0}{l + \frac{1}{2}} - \kappa_m^H \mu_l(\theta_0) + \sum_{n=0}^{\infty} \frac{\tilde{G}_{mn}}{n + \frac{1}{2}} \Lambda_{nl}^H \\ & \quad (l = 0, \dots, m-1), \end{aligned} \quad (4.38b)$$

and the coefficients

$$B_{ml} = - (4ika\beta_m + \bar{\beta}_m) \Lambda_{0l}^H + \bar{\beta}_m \delta_{0l} - \kappa_m^H \mu_l(\theta_0) - \sum_{j=1}^{m-1} b_{mj} \partial_{\theta_0}^j \frac{\cos(l + \frac{1}{2})\theta_0}{l + \frac{1}{2}} + \sum_{n=0}^{\infty} \frac{\tilde{G}_{mn}}{n + \frac{1}{2}} \Lambda_{nl}^H \quad (l = m, m+1, \dots). \quad (4.38c)$$

Continuity across  $t = \theta_0$  of the right-hand side of (4.35) gives

$$b_{m0} = -\kappa_m^H \theta_0 - \frac{2}{\pi} \left[ (4ika\beta_m + \bar{\beta}_m) \sin \frac{\theta_0}{2} - \sum_{n=0}^{\infty} \frac{\tilde{G}_{mn}}{n + \frac{1}{2}} \sin \left( n + \frac{1}{2} \right) \theta_0 \right]. \quad (4.39)$$

Note that  $b_{m0}$  does not contribute directly to the solution coefficients  $B_{ml}$  ( $l = m, m+1, \dots$ ).

Equations (4.38) are general solutions of (4.10) subject to Meixner's edge conditions (4.30b). As with the TE case, this has been confirmed numerically for  $m = 1, 2, 3$ . Explicitly for  $m = 1$  the resulting solution system

$$B_{1l} = - (4ika\beta_1 + \bar{\beta}_1) \Lambda_{0l}^H + \sum_{n=1}^{\infty} \frac{G_{1n}}{n + \frac{1}{2}} \Lambda_{nl}^H \quad (l = 1, 2, \dots), \quad (4.40a)$$

$$0 = \bar{\beta}_1 - (4ika\beta_1 + \bar{\beta}_1) \Lambda_{00}^H + \sum_{n=1}^{\infty} \frac{G_{1n}}{n + \frac{1}{2}} \Lambda_{n0}^H \quad (4.40b)$$

yields the coefficients

$$B_{1l} = \sum_{n=1}^{\infty} \frac{G_{1n}}{n + \frac{1}{2}} \Gamma_{1,nl}^H + \bar{\beta}_1 L_{1l}^H \quad (l = 1, 2, \dots), \quad (4.41a)$$

where

$$\Gamma_{1,nl}^H = \Lambda_{nl}^H - \Lambda_{n0}^H \Lambda_{0l}^H / \Lambda_{00}^H, \quad (4.41b)$$

$$L_{1l}^H = -\Lambda_{0l}^H / \Lambda_{00}^H. \quad (4.41c)$$

Satisfaction of the  $m = 1$  versions of the original dual series system (4.10) is guaranteed since *on the metal* ( $0 < \theta < \theta_0$ ),

$$\sum_{l=1}^{\infty} (2l+1) \Gamma_{1,nl}^H P_l^{-1} - (2n+1) P_n^{-1} - L_{1n}^H P_0^{-1} = 0 \quad (n = 1, 2, \dots), \quad (4.42a)$$

$$\sum_{l=1}^{\infty} (2l+1) L_{1l}^H P_l^{-1} - \frac{\Lambda_{00}^H - 1}{\Lambda_{00}^H} P_0^{-1} = 0, \quad (4.42b)$$

and since *in the aperture* ( $\theta_0 < \theta < \pi$ ),

$$\sum_{l=1}^{\infty} \Gamma_{1,nl}^H P_l^{-1} = 0, \quad (4.43a)$$

$$\sum_{l=1}^{\infty} L_{1l}^H P_l^{-1} - \bar{P}_0^{-1} = 0. \quad (4.43b)$$

Finally, consider the  $m = 0$  case. Introducing the modified coefficients

$$\bar{B}_{0n} = \begin{cases} B_{00} & (n = 0), \\ B_{0n} & (n \geq 1), \end{cases} \quad (4.44)$$

and the modified constants

$$\beta'_0 = \beta_0 - B_{00}/4ika, \quad (4.45a)$$

$$\bar{\beta}'_0 = \bar{\beta}_0 + B_{00}, \quad (4.45b)$$

the dual series system (4.8') can be rewritten as

$$\sum_{n=0}^{\infty} \left( n + \frac{1}{2} \right) \tilde{B}_{0n} P_n = -2ika\beta'_0 P_0 + \sum_{n=1}^{\infty} G_{0n} P_n \quad (0 < \theta < \theta_0), \quad (4.46a)$$

$$\sum_{n=0}^{\infty} \tilde{B}_{0n} P_n = \tilde{\beta}'_0 P_0 \quad (\theta_0 < \theta < \pi). \quad (4.46b)$$

Our TM solution process leads to the solution coefficients

$$B_{0l} = \sum_{n=1}^{\infty} \frac{G_{0n}}{n + \frac{1}{2}} \Gamma_{0,nl}^H \quad (l = 1, 2, \dots), \quad (4.47a)$$

where

$$\Gamma_{0,nl}^H = \Lambda_{nl}^H - \frac{\sin(n + \frac{1}{2})\theta_0}{\sin \theta_0/2} \Lambda_{0l}^H \quad (4.47b)$$

and the constant

$$\bar{\beta}'_0 = - \sum_{n=1}^{\infty} \frac{G_{0n}}{n + \frac{1}{2}} \Gamma_{0,n0}^H. \quad (4.47c)$$

The remaining constant  $\beta_0$  follows immediately from the continuity condition:

$$(4ika\beta'_0 + \bar{\beta}'_0) \equiv 4ika\beta_0 + \bar{\beta}_0 = \sum_{n=1}^{\infty} \frac{G_{0n}}{n + \frac{1}{2}} \frac{\sin(n + \frac{1}{2})\theta_0}{\sin \theta_0/2}. \quad (4.48)$$

### C. Coupled TE and TM solution systems

The modal coefficients of the original electromagnetics problem can now be constructed from the TE and TM dual series results. The TE solution systems for  $m \geq 1$ , (4.23), are still coupled to the corresponding TM solution systems (4.38) through the constraint relations (3.10) and (3.11). For  $m = 0$  only TM coefficients exist, and they are generated from (4.47). For each  $m$  an infinite linear system of the form (an invertible Fredholm system of the second kind)

$$V_{ml} + \sum_{n=0}^{\infty} \mathcal{M}_{m,nl} V_{mn} = \sum_{n=0}^{\infty} \mathcal{N}_{m,nl} W_{mn} \quad (l = 0, 1, 2, \dots) \quad (4.49)$$

is obtained and must be solved. A solution process analogous to the one developed in Ref. 20 can then be applied.

The infinite linear system (4.49) is reduced to a finite one by recognizing that as  $n \rightarrow \infty$  several terms rapidly go to zero. In particular, the TE and TM solutions have been constructed so that for all  $m$

$$\lim_{n \rightarrow \infty} \mathcal{M}_{m,nl} \sim \mathcal{O}(n^{-3}), \quad (4.50a)$$

$$\lim_{n \rightarrow \infty} W_{mn} \sim \mathcal{O}(n^{m-3/2-n}). \quad (4.50b)$$

Let us assume that  $N$  unknown coefficients are desired:  $V_{m1}, \dots, V_{mN}$ . Truncation then occurs in (4.49) after the  $N$ th term and the following square system results:

$$V_{ml} + \sum_{n=0}^N \mathcal{M}_{m,nl} V_{mn} = \sum_{n=0}^N \mathcal{N}_{m,nl} W_{mn} \quad (l = 0, 1, \dots, N). \quad (4.51)$$

This system can be solved numerically, for instance, by Gauss elimination. Any additional coefficients can then be

generated recursively from (4.51) by setting  $l = N + 1, N + 2, \dots, L$ .

To illustrate this procedure, consider the  $m = 1$  case. Introducing the terms  $\xi = -\bar{\beta}_1, \eta = \alpha_1, \bar{f}_{1n} = 2ikaf_{1n}, \bar{g}_{1n} = -2ikag_{1n}/(n + \frac{1}{2})$ , and  $\bar{A}_{1n} = A_{1n}/(n + \frac{1}{2})$  and combining the constraint conditions (3.10) and (3.11), the orthogonality relations (4.26b) and (4.40b), and the coefficient expressions (4.27) and (4.41), one obtains the solution system ( $l = 1, 2, \dots$ )

$$(2ikaL_{1l}^E)\xi + \bar{A}_{1l} + \sum_{n=1}^{\infty} (\chi_n^\phi \Gamma_{1,nl}^E) \bar{A}_{1n} = \sum_{n=1}^{\infty} \Gamma_{1,nl}^E \bar{f}_{1n}, \quad (4.52a)$$

$$(L_{1l}^H)\xi + B_{1l} + \sum_{n=1}^{\infty} (\tilde{\chi}_n^\psi \Gamma_{1,nl}^H) B_{1n} = \sum_{n=1}^{\infty} \Gamma_{1,nl}^H \bar{g}_{1n}, \quad (4.52b)$$

$$[2ika(1 - \Lambda_{00}^E)]\xi + (-2ika\Lambda_{00}^E)\eta + \sum_{n=1}^{\infty} (\chi_n^\phi \Lambda_{n0}^E) \bar{A}_{1n} = \sum_{n=1}^{\infty} \Lambda_{n0}^E \bar{f}_{1n}, \quad (4.52c)$$

$$(1 - \Lambda_{00}^H)\xi + [-4(ka)^2 \Lambda_{00}^H]\eta + \sum_{n=1}^{\infty} (\tilde{\chi}_n^\psi \Lambda_{n0}^H) B_{1n} = \sum_{n=1}^{\infty} \Lambda_{n0}^H \bar{g}_{1n}. \quad (4.52d)$$

These equations clearly are coupled and take the form of (4.49). The infinite system (4.52) is reduced to a finite one by noticing that

$$\lim_{n \rightarrow \infty} \chi_n^\phi \Gamma_{1,nl}^E \sim \lim_{n \rightarrow \infty} \chi_n^\psi \Gamma_{1,nl}^H \sim \mathcal{O}(n^{-3}), \quad (4.53a)$$

$$\lim_{n \rightarrow \infty} \bar{f}_{1n} \sim \lim_{n \rightarrow \infty} \bar{g}_{1n} \sim \mathcal{O}(n^{-(n+1/2)}). \quad (4.53b)$$

Assuming that the coefficients  $\bar{A}_{1n}$  and  $B_{1n}$  are desired for  $n = 1, \dots, N$ , the truncated solution system is

$$(2ikaL_{1l}^E)\xi + \bar{A}_{1l} + \sum_{n=1}^N (\chi_n^\phi \Gamma_{1,nl}^E) \bar{A}_{1n} = \sum_{n=1}^N \Gamma_{1,nl}^E \bar{f}_{1n} \quad (l = 1, 2, \dots, N), \quad (4.54a)$$

$$(L_{1l}^E)\xi + B_{1l} + \sum_{n=1}^N (\tilde{\chi}_n^\psi \Gamma_{1,nl}^H) B_{1n} = \sum_{n=1}^N \Gamma_{1,nl}^H \bar{g}_{1n} \quad (l = 1, 2, \dots, N), \quad (4.54b)$$

$$[2ika(1 - \Lambda_{00}^E)]\xi + (-2ika\Lambda_{00}^E)\eta + \sum_{n=1}^N (\chi_n^\phi \Lambda_{n0}^E) \bar{A}_{1n} = \sum_{n=1}^N \Lambda_{n0}^E \bar{f}_{1n}, \quad (4.54c)$$

$$(1 - \Lambda_{00}^H)\xi + [-4(ka)^2 \Lambda_{00}^H]\eta + \sum_{n=1}^N (\tilde{\chi}_n^\psi \Lambda_{n0}^H) B_{1n} = \sum_{n=1}^N \Lambda_{n0}^H \bar{g}_{1n}. \quad (4.54d)$$

Numerical results generated from this system will be presented below.

## V. NORMAL INCIDENCE CASE

The plane wave is normally incident when  $\theta^{\text{inc}} = 0$  or  $\theta^{\text{inc}} = \pi$ . The  $\theta^{\text{inc}} = 0$  geometry is illustrated in Fig. 2. The

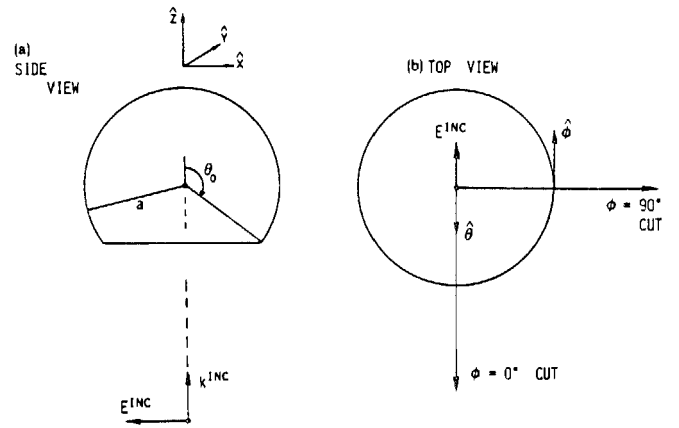


FIG. 2. Configuration of the scattering of a normally incident plane wave from a spherical shell having a circular aperture: (a) side view, (b) top view.

restriction to normal incidence provides a great simplification because

$$\begin{aligned} & \left[ \frac{mP_n^m(\cos \theta)}{\sin \theta} \right] (\theta = \theta^{\text{inc}} = 0) \\ &= \left[ + \frac{\partial}{\partial \theta} P_n^m(\cos \theta) \right] (\theta = \theta^{\text{inc}} = 0) \\ &= \frac{n(n+1)}{2} \delta_{m1}, \end{aligned} \quad (5.1)$$

$$\begin{aligned} & \left[ \frac{mP_n^m(\cos \theta)}{\sin \theta} \right] (\theta = \theta^{\text{inc}} = \pi) \\ &= \left[ - \frac{\partial}{\partial \theta} P_n^m(\cos \theta) \right] (\theta = \theta^{\text{inc}} = \pi) \\ &= (-1)^{n+1} \frac{n(n+1)}{2} \delta_{m1}. \end{aligned} \quad (5.2)$$

As a result, the potentials reduce to single sums involving only the  $m = 1$  azimuthal mode:

$$\begin{pmatrix} \Phi^{\text{inc}} \\ \Phi^s \end{pmatrix} = -E_0 \begin{pmatrix} \Phi_1^{\text{inc}} \\ \Phi_1^s \end{pmatrix} \sin \phi, \quad (5.3a)$$

$$\begin{pmatrix} \Psi^{\text{inc}} \\ \Psi^s \end{pmatrix} = Y_0 E_0 \begin{pmatrix} \Psi_1^{\text{inc}} \\ \Psi_1^s \end{pmatrix} \cos \phi. \quad (5.3b)$$

Consequently, the coupled dual series systems for normal incidence coincide with the  $m = 1$  case treated in Sec. IV; the modal coefficients  $A_{1n}$  and  $B_{1n}$  are numerically generated from the solution system (4.54) with  $\theta^{\text{inc}} = 0$  or  $\pi$ . The field components for normal incidence in terms of these coefficients are listed in Table II for convenient reference. These expressions isolate the coefficients, the  $r$ , the  $\theta$ , and the  $\phi$  dependencies; hence they are very useful for current, energy density, and cross section calculations.

An important analytical property of the dual series solution is simply revealed by the normal incidence case results. This is its trivial recovery of the scattering coefficients for the closed sphere case when  $\theta_0 = \pi$ . Let  $\theta^{\text{inc}} = 0$ . The terms

$$\Lambda_{nl}^{E,H}(\theta_0 = \pi) = \delta_{nl}, \quad (5.4)$$

TABLE II. Electric and magnetic field components for normal incidence.

Field components	
$E_r = E_0 \sum_{n=1}^{\infty} in(n+1)\tau_{1n} \frac{Z_n(kr)}{kr} P_n^{-1}(\cos \theta) \cos \phi$	
$E_\theta = E_0 \sum_{n=1}^{\infty} \left\{ \sigma_{1n} Z_n(kr) \bar{v}_{1n}(\theta) - i\tau_{1n} \frac{[krZ_n(kr)]'}{kr} \bar{w}_{1n}(\theta) \right\} \cos \phi$	
$E_\phi = E_0 \sum_{n=1}^{\infty} \left\{ \sigma_{1n} Z_n(kr) \bar{w}_{1n}(\theta) - i\tau_{1n} \frac{[krZ_n(kr)]'}{kr} \bar{v}_{1n}(\theta) \right\} \sin \phi$	
$H_r = -Y_0 E_0 \sum_{n=1}^{\infty} in(n+1)\sigma_{1n} \frac{Z_n(kr)}{kr} P_n^{-1}(\cos \theta) \sin \phi$	
$H_\theta = -Y_0 E_0 \sum_{n=1}^{\infty} \left\{ \tau_{1n} Z_n(kr) \bar{v}_{1n}(\theta) - i\sigma_{1n} \frac{[krZ_n(kr)]'}{kr} \bar{w}_{1n}(\theta) \right\} \sin \phi$	
$H_\phi = +Y_0 E_0 \sum_{n=1}^{\infty} \left\{ \tau_{1n} Z_n(kr) \bar{w}_{1n}(\theta) - i\sigma_{1n} \frac{[krZ_n(kr)]'}{kr} \bar{v}_{1n}(\theta) \right\} \cos \phi$	
Incident field	
$Z_n(kr) = j_n(kr)$	
$\sigma_{1n}^{\text{inc}} = i^n(2n+1) \begin{cases} 1, & \theta^{\text{inc}} = 0 \\ (-1)^{n+1}, & \theta^{\text{inc}} = \pi \end{cases}$	
$\tau_{1n}^{\text{inc}} = i^n(2n+1) \begin{cases} -1, & \theta^{\text{inc}} = 0 \\ (-1)^{n+1}, & \theta^{\text{inc}} = \pi \end{cases}$	
Scattered field for $r < a$	
$Z_n(kr) = j_n(kr)$	
$\sigma_{1n}^< = A_{1n} h_n(ka)$	
$\tau_{1n}^< = B_{1n} [kah_n(ka)]'$	
Scattered field for $r > a$	
$Z_n(kr) = h_n(kr)$	
$\sigma_{1n}^> = A_{1n} j_n(ka)$	
$\tau_{1n}^> = B_{1n} [kaj_n(ka)]'$	
Terms	
$\bar{v}_{1n}(\theta) = \frac{P_n^{-1}(\cos \theta)}{\sin \theta} = \frac{-1}{n(n+1)} \frac{P_n^1(\cos \theta)}{\sin \theta}$	
$\bar{w}_{1n}(\theta) = -\partial_\theta P_n^{-1}(\cos \theta) = \frac{1}{n(n+1)} \partial_\theta P_n^1(\cos \theta)$	

so that the modal coefficients

$$A_{1l} = \frac{ika(2l+1)f_{1l}}{1 + \chi_l^\phi} = \frac{-i^l(2l+1)}{h_l(ka)}, \quad (5.5a)$$

$$B_{1l} = \frac{-4ikag_{1l}}{(2l+1)(1 + \chi_l^\psi)} = \frac{i^l(2l+1)}{[kah_l(ka)]'}. \quad (5.5b)$$

Referring to Eqs. (2.4)–(2.14) and to Table II, this means that (1) for  $r \leq a$ ,  $\Phi_{1<}^s = -\Phi_{1<}^{\text{inc}}$  and  $\Psi_{1<}^s = -\Psi_{1<}^{\text{inc}}$  so that the total potentials, hence the fields, are identically zero there and the boundary conditions  $E_{\text{tan}}^s(r=a) = -E_{\text{tan}}^{\text{inc}}(r=a)$  are satisfied and (2) for  $r > a$ , the standard results for the scattered potentials—fields given, for instance, in Ref. 38, Sec. 6.9, and Ref. 39, Sec. 16.9, are recovered.

## VI. CURRENTS ON THE SPHERICAL SHELL

The most stringent test of the dual series solution is the calculation of the currents  $J_\theta$  and  $J_\phi$  on the open spherical

shell. Verification of the required current behavior near the aperture edge is immediately apparent from graphical results. Moreover, the vanishing of the current in the aperture is an excellent test of the results and reflects the satisfaction of the corresponding TE and TM dual series equations in that region. For normal incidence the current expressions (4.11) simply become

$$J_\theta(\theta, \phi) = \left[ \frac{-Y_0 E_0}{(ka)^2} \cos \phi \sum_{n=1}^{\infty} \left[ A_{1n} \frac{P_n^{-1}(\cos \theta)}{\sin \theta} + ikaB_{1n} \partial_\theta P_n^{-1}(\cos \theta) \right] \right], \quad (6.1)$$

$$J_\phi(\theta, \phi) = \left[ \frac{+Y_0 E_0}{(ka)^2} \sin \phi \sum_{n=1}^{\infty} \left[ A_{1n} \partial_\theta P_n^{-1}(\cos \theta) + ikaB_{1n} \frac{P_n^{-1}(\cos \theta)}{\sin \theta} \right] \right]. \quad (6.2)$$

### A. Analytical preconditioning

Consider first the quasistatic case where  $ka = 0.01$ ,  $\theta_0 = 120^\circ$ . In all of the examples  $a = 1.0$ . Simply performing the sums in (6.2) with the solution coefficients generated from Eqs. (4.54), we find that the number of terms required to track the square root singularity in  $J_\phi$  is large. The polynomial sum  $\sum_{n=1}^L A_{1n} \partial_\theta P_n^{-1}$  is the cause of this difficulty. However, the truncation number  $N$  (Sec. IV C) need not be large; and the remaining coefficients  $n = N + 1, \dots, L$  are recursively defined from (4.54a) and (4.54b). This is demonstrated in Fig. 3 where the real part of  $J_\phi(\theta, \pi/2) [(ka)^2/Y_0 E_0]$  is given for various truncation numbers. In Figs. 3(a), 3(b), and 3(c) the truncation numbers  $N = 5$  and  $L = 50, 500$ , and  $5000$ , respectively. However, the results may be improved by treating the singularity analytically as follows.

Inserting the coefficient expressions (4.27a) and (4.41) into (6.1) and (6.2) and referring to the definitions given in Table III, the current components

$$J_\theta = \frac{-Y_0 E_0}{2(ka)^2} \cos \phi \left\{ \sum_{n=1}^N \left[ F_{1n} \frac{K_n^E(\theta)}{\sin \theta} + 4ika \frac{G_{1n}}{2n+1} \partial_\theta K_n^H(\theta) \right] - 2ka\xi \left[ \frac{K_0^E(\theta)}{\sin \theta} + \partial_\theta K_0^H(\theta) \right] \right\}, \quad (6.3)$$

$$J_\phi = \frac{+Y_0 E_0}{2(ka)^2} \sin \phi \left\{ \sum_{n=1}^N \left[ F_{1n} \partial_\theta K_n^E(\theta) + 4ika \frac{G_{1n}}{2n+1} \frac{K_n^H(\theta)}{\sin \theta} \right] - 2ka\xi \left[ \partial_\theta K_0^E(\theta) + \frac{K_0^H(\theta)}{\sin \theta} \right] \right\} \quad (6.4)$$

result. Two advantages of these expressions are immediate. First, the coefficients obtained from the matrix inversion can be used directly without calculating any additional coefficients by recursion. Second, the currents vanish analytically in the aperture. The terms proportional to  $F_{1n}$  and  $G_{1n}$  give

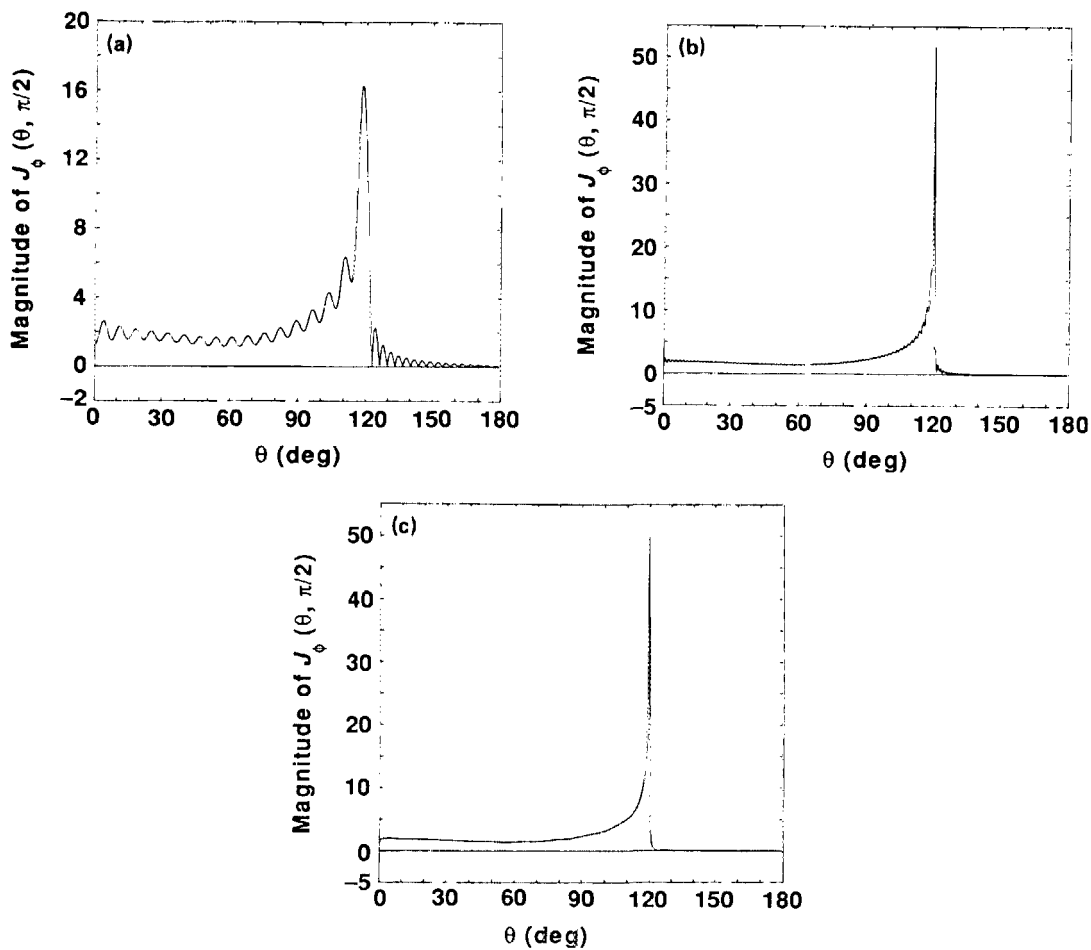


FIG. 3. Brute force summation of the  $J_\phi(\theta, \pi/2)$  current expression requires a large number of terms to eliminate the numerical Gibbs' phenomena: (a)  $N = 5, L = 50$ ; (b)  $N = 5, L = 500$ ; (c)  $N = 5, L = 5000$ .

TABLE III. Special functions and relations for the current expressions.

$$\begin{aligned}
 K_0^H(\theta) &= \sum_{l=1}^{\infty} L_{1,l}^H(\theta_0) P_l^{-1}(\cos \theta) = \begin{cases} +\bar{P}_0^{-1}(\theta) & (\theta_0 < \theta \leq \pi) \\ +\bar{P}_0^{-1}(\theta) + s_H(\theta) & (0 \leq \theta < \theta_0) \end{cases} \\
 K_0^E(\theta) &= \sum_{l=1}^{\infty} (2l+1) L_{1,l}^E(\theta_0) P_l^{-1}(\cos \theta) = \begin{cases} +\bar{P}_0^{-1}(\theta) & (\theta_0 < \theta \leq \pi) \\ +\bar{P}_0^{-1}(\theta) + s_E(\theta) & (0 \leq \theta < \theta_0) \end{cases} \\
 K_n^H(\theta) &= \sum_{l=1}^{\infty} \Gamma_{1,nl}^H(\theta_0) P_l^{-1}(\cos \theta) = \begin{cases} 0 & (\theta_0 < \theta \leq \pi) \\ S_n^H(\theta) + \Lambda_{n0}^H s_H(\theta) & (0 \leq \theta < \theta_0) \end{cases} \\
 K_n^E(\theta) &= \sum_{l=1}^{\infty} (2l+1) \Gamma_{1,nl}^E(\theta_0) P_l^{-1}(\cos \theta) = \begin{cases} 0 & (\theta_0 < \theta \leq \pi) \\ S_n^E(\theta) + \Lambda_{n0}^E s_E(\theta) & (0 \leq \theta < \theta_0) \end{cases}
 \end{aligned}$$

where

$$\begin{aligned}
 s_H(\theta) &= \frac{+4}{\pi \Lambda_{00}^H \sin \theta} \left\{ \cos \frac{\theta_0}{2} \left( \cos^2 \frac{\theta}{2} - \cos^2 \frac{\theta_0}{2} \right)^{1/2} - \cos^2 \frac{\theta}{2} \left[ \arccos \left( \frac{\cos(\theta_0/2)}{\cos(\theta/2)} \right) \right] \right\} \\
 s_E(\theta) &= \frac{-4}{\pi \Lambda_{00}^E \sin \theta} \left\{ \cos \frac{\theta_0}{2} \left( \cos^2 \frac{\theta}{2} - \cos^2 \frac{\theta_0}{2} \right)^{1/2} + \cos^2 \frac{\theta}{2} \left[ \arccos \left( \frac{\cos(\theta_0/2)}{\cos(\theta/2)} \right) \right] \right\} \\
 u_n(\theta) &= \frac{4}{\pi} \cos \left( n + \frac{1}{2} \right) \theta_0 \frac{[2(\cos \theta - \cos \theta_0)]^{1/2}}{\sin \theta} \\
 S_n^H(\theta) &= \sum_{l=1}^{\infty} \Lambda_{nl}^H P_l^{-1}(\cos \theta) + \Lambda_{n0}^H \bar{P}_0^{-1}(\theta) \\
 S_n^E(\theta) &= u_n(\theta) + (2n+1) S_n^H(\theta)
 \end{aligned}$$

TABLE IV. Derivative relations employed in the current calculations.

$$\begin{aligned} \partial_\theta s_H(\theta) &= \frac{-4}{\pi \Lambda_{00}^H} \left\{ \frac{-\cos \theta}{\sin^2 \theta} \left[ \cos \frac{\theta_0}{2} \left( \cos^2 \frac{\theta}{2} - \cos^2 \frac{\theta_0}{2} \right)^{1/2} - \cos^2 \frac{\theta}{2} \arccos \left( \frac{\cos(\theta_0/2)}{\cos(\theta/2)} \right) \right] + \frac{1}{2} \arccos \left( \frac{\cos(\theta_0/2)}{\cos(\theta/2)} \right) \right\} \\ \partial_\theta s_E(\theta) &= \frac{+4}{\pi \Lambda_{00}^E} \left\{ \frac{-\cos \theta}{\sin^2 \theta} \left[ \cos \frac{\theta_0}{2} \left( \cos^2 \frac{\theta}{2} - \cos^2 \frac{\theta_0}{2} \right)^{1/2} \right. \right. \\ &\quad \left. \left. + \cos^2 \frac{\theta}{2} \arccos \left( \frac{\cos(\theta_0/2)}{\cos(\theta/2)} \right) \right] - \frac{1}{2} \arccos \left( \frac{\cos(\theta_0/2)}{\cos(\theta/2)} \right) - \frac{\cos(\theta_0/2)}{2(\cos^2(\theta/2) - \cos^2(\theta_0/2))^{1/2}} \right\} \\ \partial_\theta u_n(\theta) &= \frac{-8 \cos(n+1)\theta_0}{\pi \sin^2 \theta (\cos^2(\theta/2) - \cos^2(\theta_0/2))^{1/2}} \left[ \cos^2 \frac{\theta}{2} \left( \cos^2 \frac{\theta}{2} - \cos^2 \frac{\theta_0}{2} \right) + \sin^2 \frac{\theta}{2} \cos^2 \frac{\theta_0}{2} \right] \\ \partial_\theta S_n^H(\theta) &= \sum_{i=1}^{\infty} \Lambda_{ni}^H \partial_\theta P_i^{-1} - \Lambda_{n0}^H \frac{\bar{P}_0^{-1}}{\sin \theta} \\ \partial_\theta S_n^E(\theta) &= \partial_\theta u_n(\theta) + (2n+1) \partial_\theta S_n^H(\theta) \end{aligned}$$

no contributions in the aperture because  $K_n^E$  and  $K_n^H$  are zero there. The terms proportional to  $\xi$  also reduce to zero there by (3.9).

Restricting now our attention to the behavior of the current on the metal, (6.3) and (6.4) yield

$$\begin{aligned} J_\theta &= -\frac{Y_0 E_0}{2(ka)^2} \cos \phi \left\{ \sum_{n=1}^N \left[ F_{1n} \frac{K_n^E(\theta)}{\sin \theta} \right. \right. \\ &\quad \left. \left. + 4ika \frac{G_{1n}}{2n+1} \partial_\theta K_n^H(\theta) \right] \right. \\ &\quad \left. - 2ka\xi \left[ \frac{s_E(\theta)}{\sin \theta} + \partial_\theta s_H(\theta) \right] \right\} \quad (0 \leq \theta \leq \theta_0), \end{aligned} \tag{6.3'}$$

$$\begin{aligned} J_\phi &= +\frac{Y_0 E_0}{2(ka)^2} \sin \phi \left\{ \sum_{n=1}^N \left[ F_{1n} \partial_\theta K_n^E(\theta) \right. \right. \\ &\quad \left. \left. + 4ika \frac{G_{1n}}{2n+1} \frac{K_n^H(\theta)}{\sin \theta} \right] \right. \\ &\quad \left. - 2ka\xi \left[ \partial_\theta s_E(\theta) + \frac{s_H(\theta)}{\sin \theta} \right] \right\} \quad (0 \leq \theta \leq \theta_0). \end{aligned} \tag{6.4'}$$

Near the aperture edge  $\theta = \theta_0$  the terms  $K_n^E(\theta)$  and  $K_n^H(\theta)$  behave, respectively, as  $(\theta_0 - \theta)^{1/2}$  and  $(\theta_0 - \theta)^{3/2}$ . Consequently, the square root singularity in  $J_\phi$  is generated by the term  $\partial_\theta K_n^E(\theta)$  and, referring to Table IV, by the term  $\partial_\theta s_E(\theta)$ . If the former is generated numerically, a large number of coefficients are required. However, referring to

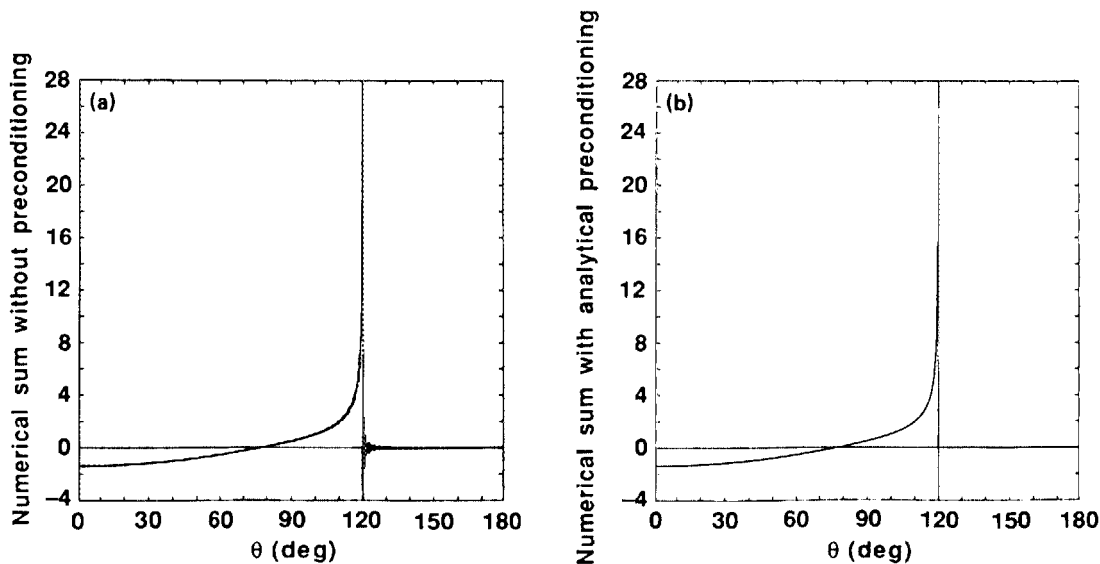


FIG. 4. The Gibbs' phenomena is removed by handling the edge singularity analytically. The dominant sum in the  $J_\phi(\theta, \pi/2)$  expression (a) without analytical preconditioning, and (b) with analytical preconditioning.



Tables III and IV for  $0 \leq \theta < \theta_0$ , the relations

$$K_n^E(\theta) = S_n^E(\theta) + \Lambda_{n0}^E S_E^E(\theta) \\ = (2n + 1)S_n^H(\theta) + u_n(\theta) + \Lambda_{n0}^E S_E^E(\theta), \quad (6.5)$$

$$\partial_\theta K_n^E(\theta) = (2n + 1)\partial_\theta S_n^H(\theta) + \partial_\theta u_n(\theta) + \Lambda_{n0}^E \partial_\theta S_E^E(\theta) \quad (6.6)$$

indicate that one only needs to evaluate  $S_n^H(\theta)$  and  $\partial_\theta S_n^H(\theta)$  numerically. Since  $S_n^H(\theta) = K_n^H(\theta) - \Lambda_{n0}^H S_H^H(\theta)$  over the metal and near the aperture edge  $\partial_\theta K_n^H(\theta) \sim (\theta_0 - \theta)^{1/2}$  and  $\partial_\theta S_H^H \sim (\theta_0 - \theta)^{1/2}$ , the term  $\partial_\theta S_n^H(\theta) \sim (\theta_0 - \theta)^{1/2}$  near  $\theta = \theta_0$ , which circumvents the numerical difficulties. The square root singularity is handled analytically through the terms  $\partial_\theta u_n$  and  $\partial_\theta S_E^E$ . A comparison of  $\partial_\theta K_1^E(\theta)$  evaluated directly and with (6.6) is given in Fig. 4. Each sum included 800 terms. As desired, the (numerical) oscillations were removed by the analytical preconditioning.

### B. Numerical results

Because the current components

$$J_\theta(\theta, \phi) = J_\theta(\theta, 0) \cos \phi, \quad (6.7)$$

$$J_\phi(\theta, \phi) = J_\phi(\theta, \pi/2) \sin \phi, \quad (6.8)$$

their important features are illustrated succinctly by considering  $J_\theta(\theta, 0)$  and  $J_\phi(\theta, \pi/2)$ . Examples of the scaled current terms  $\mathcal{J}_\theta = J_\theta(\theta, 0)[-2(ka)^2/Y_0 E_0]$  and  $\mathcal{J}_\phi = J_\phi(\theta, \pi/2)[2(ka)^2/Y_0 E_0]$  are given in Figs. 5–10 for various  $ka$ , aperture sizes, and angles of incidence.

Values of  $|\mathcal{J}_\theta|$  and  $|\mathcal{J}_\phi|$  are given in Figs. 5 and 6 for the quasistatic limit ( $ka = 0.01$ ), the angle of incidence  $\theta^{\text{inc}} = 0.0$ , and, respectively, the aperture angles  $\theta_0 = 120^\circ$  and  $\theta_0 = 170^\circ$ . For both cases the truncation number  $N = 10$ . Essentially the same results were generated with  $N = 3$ . This low truncation number is typical for quasistatic cases because the  $n = 1$  term dominates the behavior. The term  $|\mathcal{J}_\theta|$  is given in Fig. 7 for  $\theta^{\text{inc}} = 0.0$ ,  $\theta_0 = 120^\circ$ , and the  $ka$  values 1.0, 3.0, 5.0, and 10.0. The corresponding graphs of  $|\mathcal{J}_\phi|$  are given in Fig. 8. For all of these cases the truncation number was taken to be  $N = 10(ka)$ . This choice yields convergent results. The plots in Fig. 7 clearly demonstrate that our solution reproduces the required  $(\theta_0 - \theta)^{1/2}$  behavior of  $J_\theta$  near  $\theta = \theta_0$ ; Fig. 8 demonstrates that the required square root singularity of  $J_\phi$  near  $\theta = \theta_0$  is present. In Fig. 9 the terms  $\text{Re}(\mathcal{J}_\theta)$  and  $|\mathcal{J}_\theta|$  are plotted for  $\theta_0 = 120^\circ$ ,  $ka = 1.0$ ,  $\theta^{\text{inc}} = 0^\circ$ , and  $\theta^{\text{inc}} = 180^\circ$ . Very different behaviors are obtained. When the wave is incident on the shell

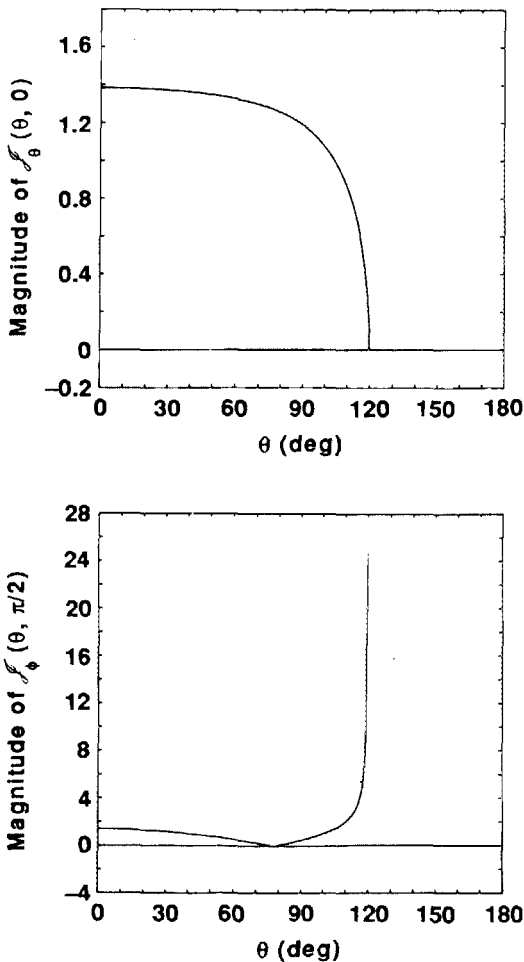


FIG. 5. The magnitudes of the current terms  $\mathcal{J}_\theta(\theta, 0)$  and  $\mathcal{J}_\phi(\theta, \pi/2)$  induced on an open spherical shell with  $\theta_0 = 120^\circ$  when  $ka = 0.01$  and  $\theta^{\text{inc}} = 0.0^\circ$ .

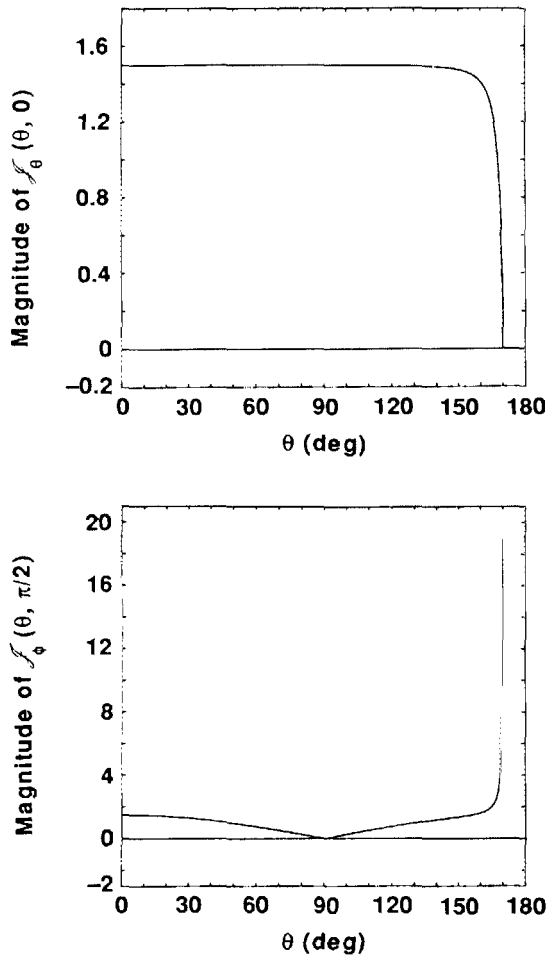


FIG. 6. The magnitudes of the current terms  $\mathcal{J}_\theta(\theta, 0)$  and  $\mathcal{J}_\phi(\theta, \pi/2)$  induced on an open spherical shell with  $\theta_0 = 170^\circ$  when  $ka = 0.01$  and  $\theta^{\text{inc}} = 0.0^\circ$ .

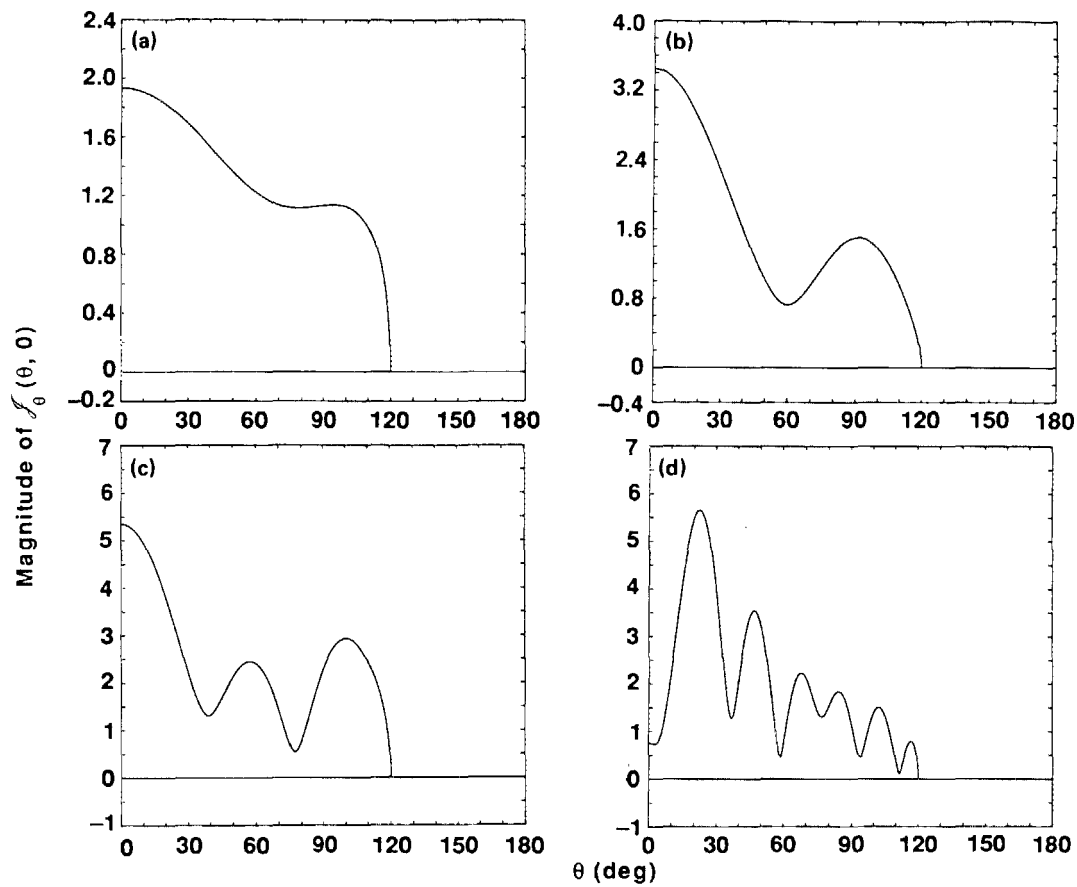


FIG. 7. The current term  $\mathcal{J}_\theta(\theta, 0)$  induced on an open spherical shell with  $\theta_0 = 120^\circ$  when  $\theta^{\text{inc}} = 0^\circ$  and the  $ka$  of the incident plane wave is (a) 1.0, (b) 3.0, (c) 5.0, and (d) 10.0.

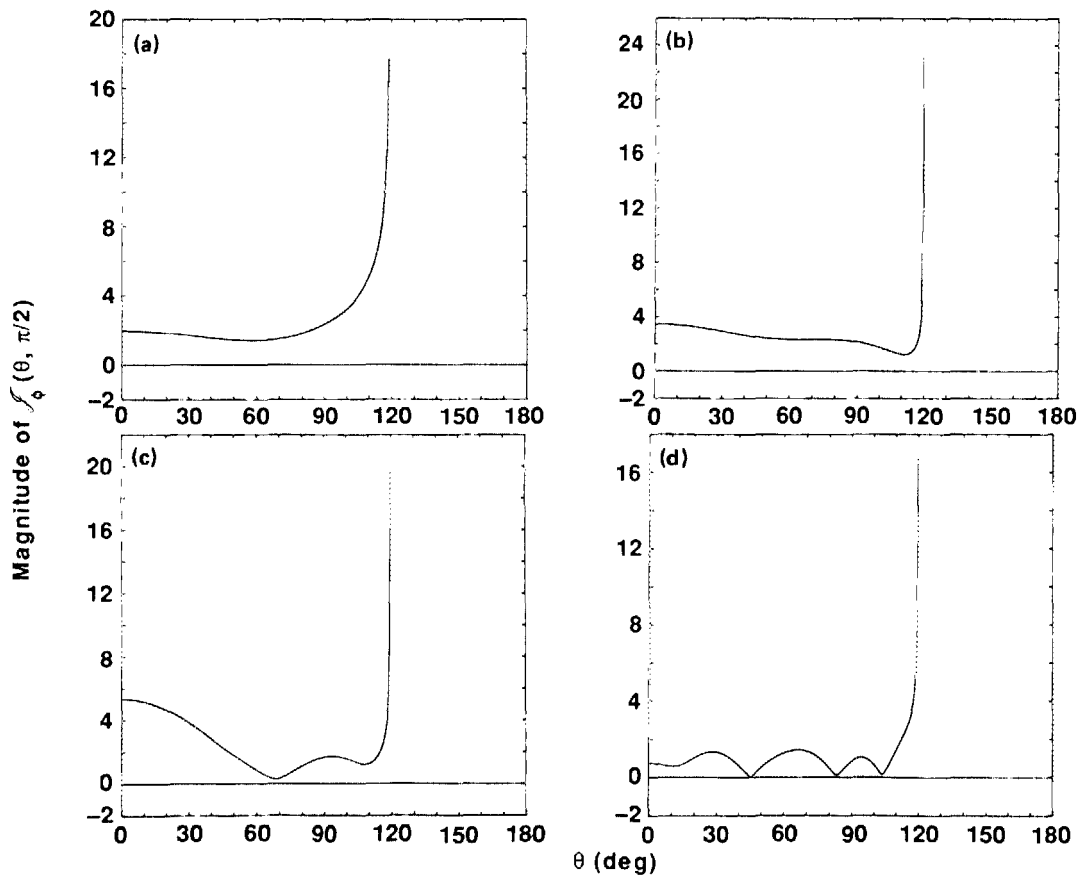


FIG. 8. The current term  $\mathcal{J}_\phi(\theta, \pi/2)$  induced on an open spherical shell with  $\theta_0 = 120^\circ$  when  $\theta^{\text{inc}} = 0^\circ$  and the  $ka$  of the incident plane wave is (a) 1.0, (b) 3.0, (c) 5.0, and (d) 10.0.

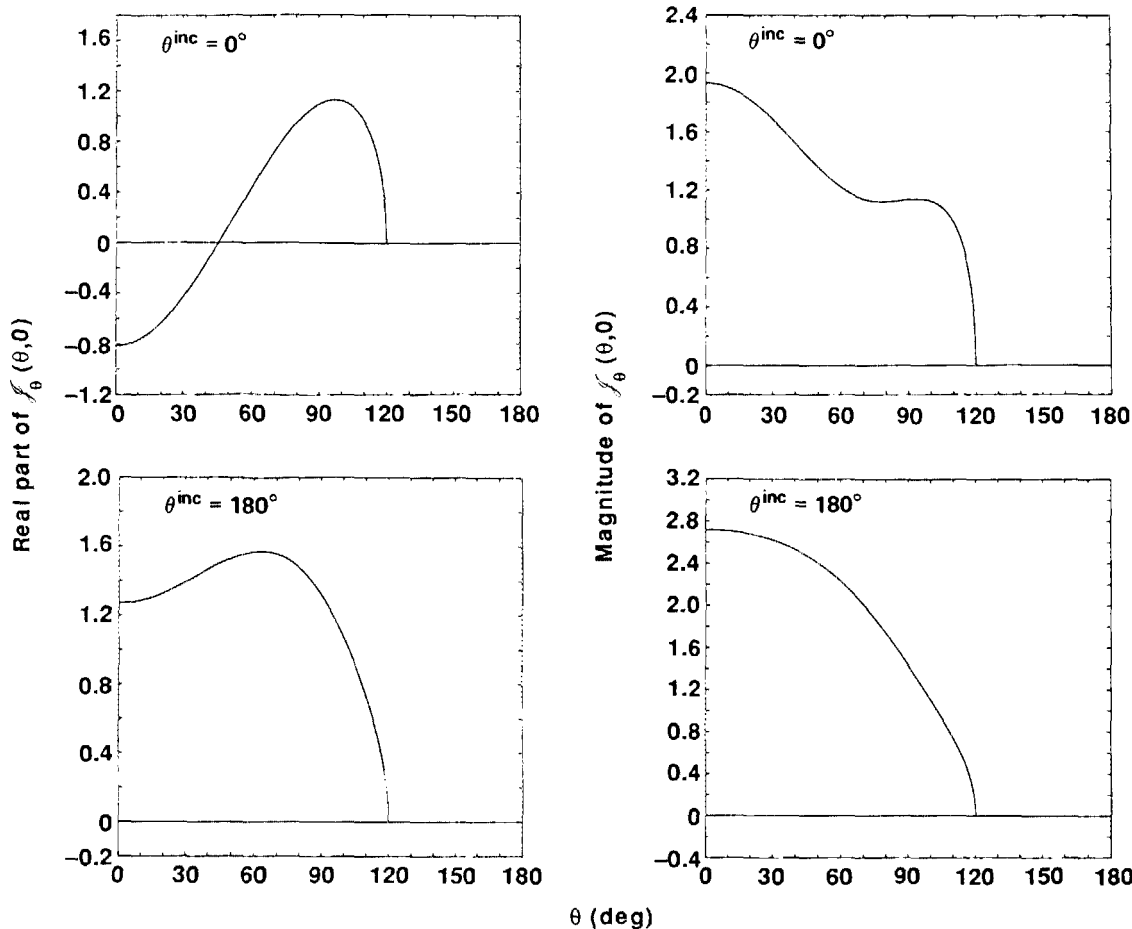


FIG. 9. A comparison of the real part and the magnitude of the current term  $\mathcal{J}_\theta(\theta, 0)$  induced on an open spherical shell with  $\theta_0 = 120^\circ$  when  $ka = 1.0$  and the angle of incidence  $\theta^{\text{inc}} = 0^\circ$  and  $\theta^{\text{inc}} = 180^\circ$ .

( $\theta^{\text{inc}} = 180^\circ$ ),  $|\mathcal{J}_\theta|$  is much more uniformly distributed over the shell. The hump near the aperture edge which appeared when  $\theta^{\text{inc}} = 0^\circ$  is no longer present. In Fig. 10,  $|\mathcal{J}_\theta|$  is given for  $\theta^{\text{inc}} = 0^\circ$  and  $ka = 1.0$  when  $\theta_0 = 120^\circ$  and  $\theta_0 = 170^\circ$ . The latter case exhibits a more pronounced hump near the aperture edge.

The distributions of  $J_\theta$  and  $J_\phi$  over the entire spherical shell for  $ka = 3$ ,  $\theta^{\text{inc}} = 0^\circ$ , and  $\theta_0 = 120^\circ$  are shown in Figs. 11 and 12. In Fig. 11 the values of  $|\mathcal{J}_\theta|$  and  $|\mathcal{J}_\phi|$  are replotted in more detail to provide a reference for Figs. 12. In Figs. 12(a) and 12(b)  $\mathcal{J}_\theta$  is viewed from the directions ( $\theta = 0$ ,  $\phi = 0$ ) and ( $\theta = 76^\circ$ ,  $\phi = 125^\circ$ ). Dark red represents the largest values; dark blue the smallest ones. The characteristic cosine pattern and null at the aperture edge are very apparent. The corresponding views of  $\mathcal{J}_\phi$  are given in Figs. 12(c) and 12(d). The associated sine pattern and edge singularity are nicely reproduced.

The current results have been validated with a totally independent method<sup>19</sup>: a completely numerical solution based upon a method of moments (MoM) analysis of the problem. It has been demonstrated that the MoM solution converges to the dual series results when the former is applicable.

## VII. ENERGY DENSITIES

To provide some measure of the degree of coupling of the incident field into the spherical cavity, the energy density at the center of the shell normalized to the incident field energy density there was calculated. This also allows a direct comparison with Senior-Desjardins results.<sup>12,13</sup>

Consider the normal incidence field expressions given in Table II for  $r = 0$ . With the small argument relations in the Appendix, one obtains ( $n \neq 0$ )

$$j_n(0) \equiv 0, \quad \{[xj_n(x)]'/x\}_{x=0} = \frac{2}{3}\delta_{n1},$$

$$[j_n(x)/x]_{x=0} = \frac{1}{3}\delta_{n1}.$$

Moreover,  $P_1^{-1}(\cos \theta)/\sin \theta = -\frac{1}{2}$  and  $-\partial_\theta P_1^{-1}(\cos \theta) = \cos(\theta/2)$ . Therefore, the general electric and magnetic field vectors at the origin are

$$(E_r, E_\theta, E_\phi)(r=0)$$

$$= (i/3)E_0\tau_{11}(\sin \theta \cos \phi, \cos \theta \cos \phi, -\sin \phi), \quad (7.1)$$

$$(H_r, H_\theta, H_\phi)(r=0)$$

$$= (i/3)Y_0E_0\sigma_{11}(\sin \theta \sin \phi, \cos \theta \sin \phi, +\cos \phi), \quad (7.2)$$

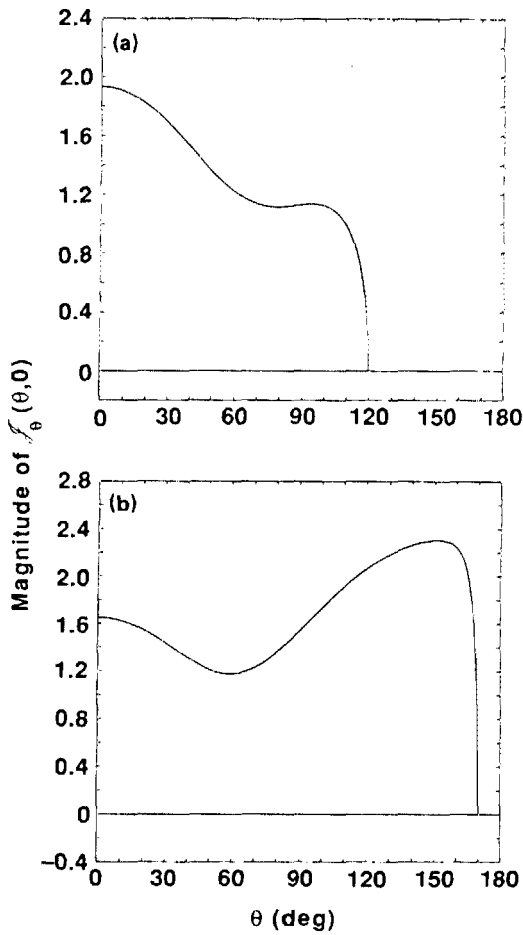


FIG. 10. A comparison of the magnitudes of the current term  $f_{\theta}(\theta, 0)$  induced on open spherical shells with (a)  $\theta_0 = 120^\circ$  and (b)  $\theta_0 = 170^\circ$ , when  $ka = 1.0$  and the angle of incidence  $\theta^{\text{inc}} = 0^\circ$ .

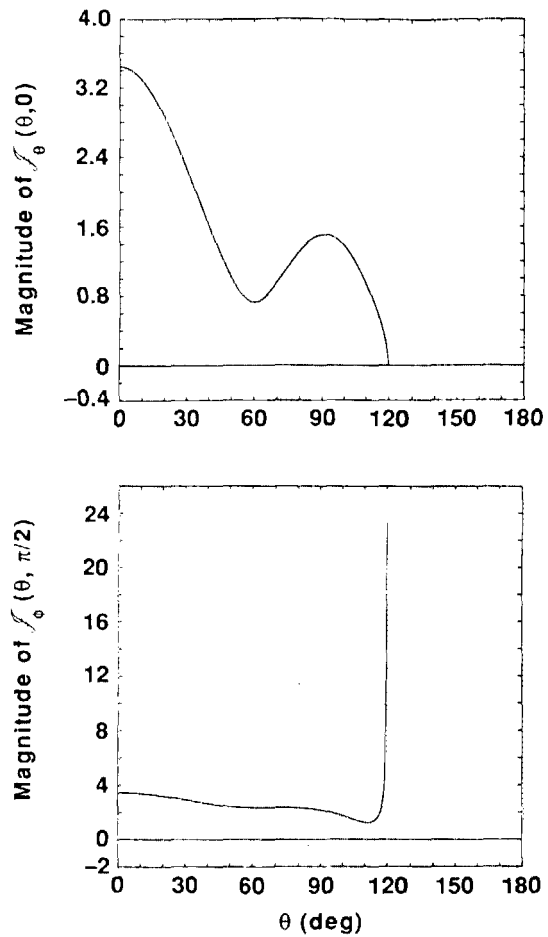


FIG. 11. The magnitudes of the current terms  $f_{\theta}(\theta, 0)$  and  $f_{\theta}(\theta, \pi/2)$  induced on an open spherical shell with  $\theta_0 = 120^\circ$  when  $ka = 3.0$  and the angle of incidence  $\theta^{\text{inc}} = 0^\circ$ .

where

$$\sigma_{11} = \sigma_{11}^{\text{inc}} + \sigma_{11}^{\leq} = \sigma_{11}^{\text{inc}} + A_{11}h_1(ka), \quad (7.3)$$

$$\tau_{11} = \tau_{11}^{\text{inc}} + \tau_{11}^{\leq} = \tau_{11}^{\text{inc}} + B_{11}[kah_1(ka)]', \quad (7.4)$$

the incident field terms being

$$\sigma_{11}^{\text{inc}} = \begin{cases} 3i, & \theta^{\text{inc}} = 0, \\ 3i, & \theta^{\text{inc}} = \pi, \end{cases} \quad (7.5)$$

$$\tau_{11}^{\text{inc}} = \begin{cases} -3i, & \theta^{\text{inc}} = 0, \\ 3i, & \theta^{\text{inc}} = \pi. \end{cases} \quad (7.6)$$

Consequently, the associated energy density relation is simply

$$\begin{aligned} U(r=0) &= \frac{1}{2}(|\mathbf{E}|^2 + |Z_0\mathbf{H}|^2)(r=0) \\ &= E_0^2(|\sigma_{11}|^2 + |\tau_{11}|^2)/18, \end{aligned} \quad (7.7)$$

which leads to the desired energy density ratio

$$\begin{aligned} \frac{U_{\text{tot}}(r=0)}{U_{\text{inc}}(r=0)} &= \frac{|\sigma_{11}^{\text{inc}} + \sigma_{11}^{\leq}|^2 + |\tau_{11}^{\text{inc}} + \tau_{11}^{\leq}|^2}{|\sigma_{11}^{\text{inc}}|^2 + |\tau_{11}^{\text{inc}}|^2} \\ &= \frac{1}{18} \{ |3i + A_{11}h_1(ka)|^2 + |3i \mp B_{11}[kah_1(ka)]'|^2 \} \\ &= \frac{1}{2} \{ |1 - (i/3)A_{11}h_1(ka)|^2 + |1 \pm (i/3)B_{11}[kah_1(ka)]'|^2 \}. \end{aligned} \quad (7.8)$$

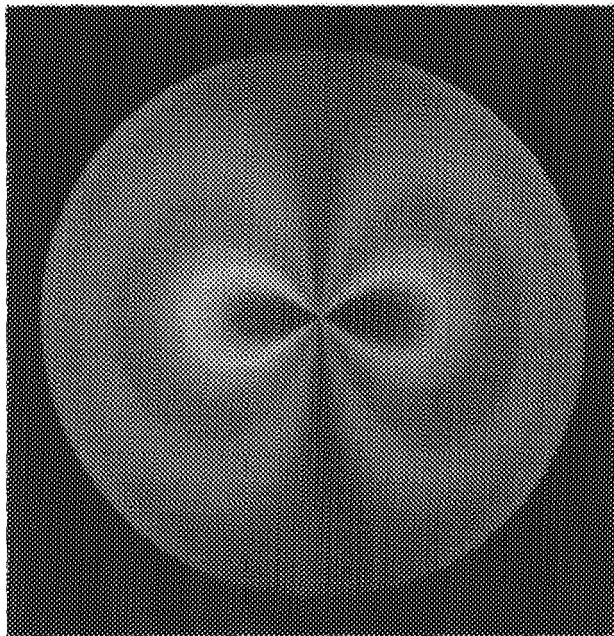
The upper sign is appropriate for  $\theta^{\text{inc}} = 0$ ; the lower sign for  $\theta^{\text{inc}} = \pi$ .

The quantity  $10 \log_{10} [U_{\text{tot}}(r=0)/U_{\text{inc}}(r=0)]$  is plotted in Figs. 13–15 for the aperture angles  $10^\circ$ ,  $30^\circ$ , and  $60^\circ$  (i.e., for  $\theta_0 = 170^\circ$ ,  $150^\circ$ , and  $120^\circ$ ) and for the angle of incidence  $\theta^{\text{inc}} = 0^\circ$ . Several interesting features are apparent immediately. For  $\theta_0 = 170^\circ$  the open spherical shell acts very similarly to a spherical cavity of the same size. The peaks in the data at  $ka = 4.49$ ,  $2.74$ ,  $3.87$ , and  $4.97$  closely correspond, respectively, to the lowest TE and TM modes of a closed cavity (see Ref. 36, pp. 268–271); i.e., to the lowest-order zero  $x_{11}$  of  $\{x_{j_1}(x)\}$  and to the zeros  $x'_{11}$ ,  $x'_{12}$ , and  $x'_{13}$  of  $\{x_{j_1}(x)\}'$ . They are slightly offset (detuned) from the closed cavity values because of the presence of the aperture. Extensions of the discussion in Sec. V for  $\theta_0$  near  $\pi$  and for  $ka$  small lead to the approximate coefficient expressions in this  $ka$  region,

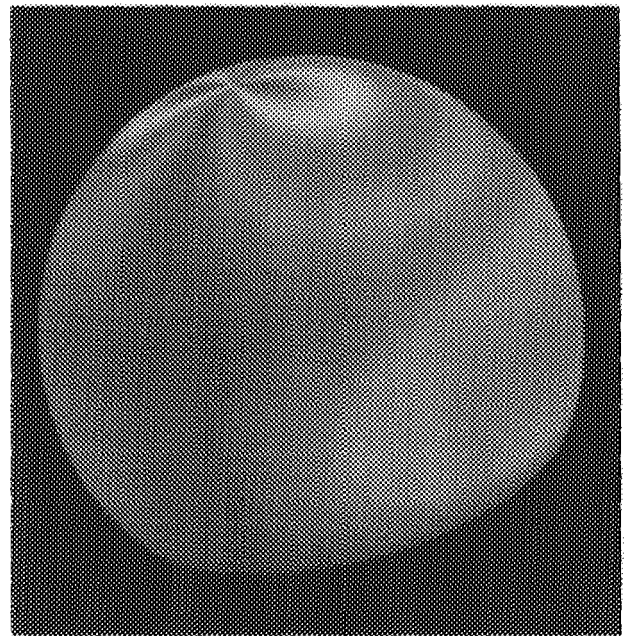
$$A_{11}h_1(ka) \sim \frac{\sum_{n=1}^N f_{1n} \Gamma_{n1}^E}{j_1(ka)}, \quad (7.9a)$$

$$B_{11}[kah_1(ka)]' \sim \frac{\sum_{n=1}^N (3g_{1n} \Gamma_{n1}^H / 2n + 1)}{[kaj_1(ka)]'}, \quad (7.9b)$$

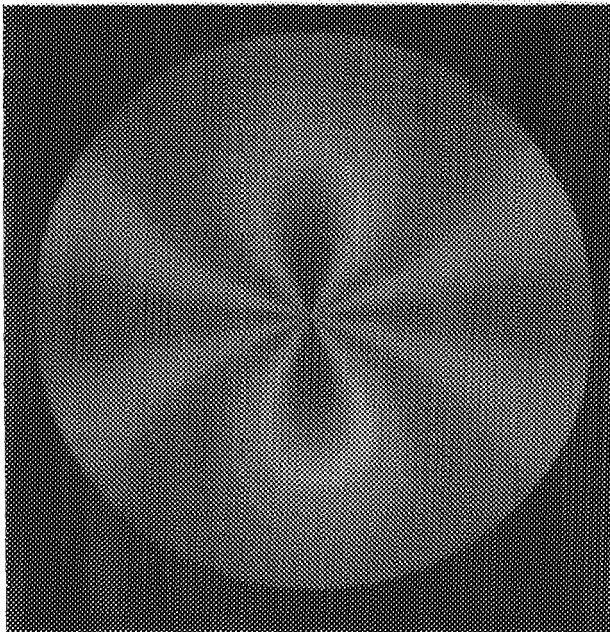
which readily explain the locations of the observed features. At higher  $ka$  peaks corresponding to the roots  $x_{np}$  of



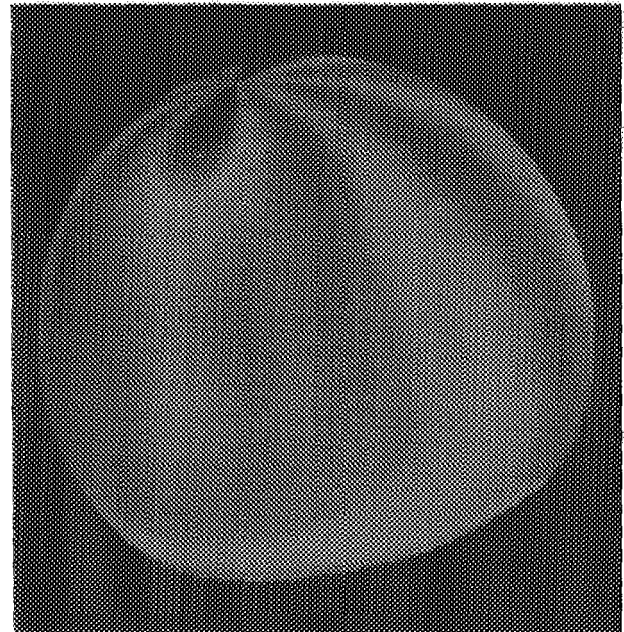
(a)



(b)



(c)



(d)

FIG. 12. The modal structure of the currents induced on an open spherical shell with  $\theta_0 = 120^\circ$  when  $ka = 3.0$  and the angle of incidence  $\theta^{inc} = 0^\circ$  is revealed with three dimensional graphics: (a) a top view of the magnitude of the current term  $\mathcal{J}_\theta(\theta, \phi)$ ; (b) a side view of the magnitude of the current term  $\mathcal{J}_\theta(\theta, \phi)$ ; (c) a top view of the magnitude of the current term  $\mathcal{J}_\phi(\theta, \phi)$ ; (d) a side view of the magnitude of the current term  $\mathcal{J}_\phi(\theta, \phi)$ .

$[xj_n(x)] = 0$  and  $x'_{np}$  of  $[xj_n(x)]' = 0$  appear. The antiresonance form of the peaks at  $ka = 3.87$  and  $4.97$  was not anticipated. In fact, only the  $TE_{n1}$  and  $TM_{n1}$  modes ( $n = 1, 2, \dots$ ) develop the resonance form of the peaks; all others have the antiresonant form. This behavior is a result of (1) the modal patterns induced in the open cavity—all  $TE_{np}$  and  $TM_{np}$  ( $p \neq 1$ ) modes have nulls at  $r = 0$ , and (2) a reradiation effect that occurs because the aperture is backed

by a resonant cavity.<sup>40</sup> The important features in the  $ka$  scans of (7.8) for larger apertures are associated with the modes effecting the antiresonant behavior.

Detuning of the cavity by the larger aperture is noticeable in the  $\theta_0 = 150^\circ$  data. The resonance peaks are broadened and the antiresonance peaks have become broad depressions. The resonance locations are downshifted to lower  $ka$  values (lower frequency); the antiresonance locations are

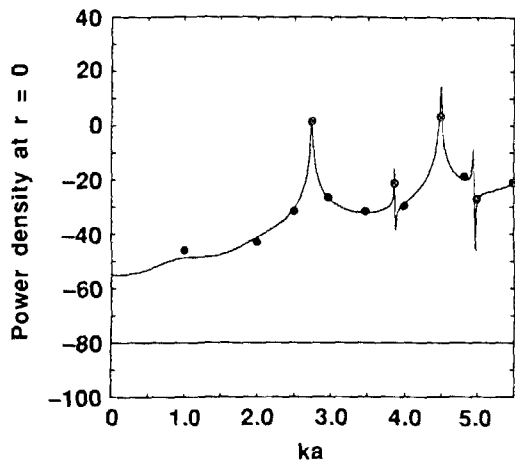


FIG. 13. A scan in  $ka$  of the total energy density of the field at the origin normalized to the energy density of the incident field there for an open spherical shell having an aperture angle of  $170^\circ$  and a plane wave incident at  $\theta^{inc} = 0^\circ$ . The solid line is generated by the dual series solution, the dots by a MoM surface patch code.

upshifted to higher  $ka$  values (higher frequency). The depressions at the antiresonance locations indicate that this slightly open cavity may have poor energy storage characteristics, hence a large scattering cross section at those points. Energy storage and cross-section calculations are also in progress to study this (for instance see Ref. 40).

The largest aperture ( $\theta_0 = 120^\circ$ ) data shows nearly a complete detuning of the cavity. The observed depressions are shallower and broadened. They correspond to the original antiresonance locations  $ka = 3.87$  and  $4.97$ , thus demonstrating the considerable upshift in  $ka$  of their locations as the aperture size increases. The data also indicates a focusing of the energy near  $r = 0$  over a large range of  $ka$ . This is expected since the shell is beginning to look largely like a spherical reflector when  $\theta_0 = 120^\circ$ .

Comparing these results with those of Senior and Desjardins, very distinct dissimilarities are evident. Although the resonance peaks at  $ka = 2.74$  and  $4.49$  are present in their results, the antiresonance peaks at  $ka = 3.87$  and  $4.97$

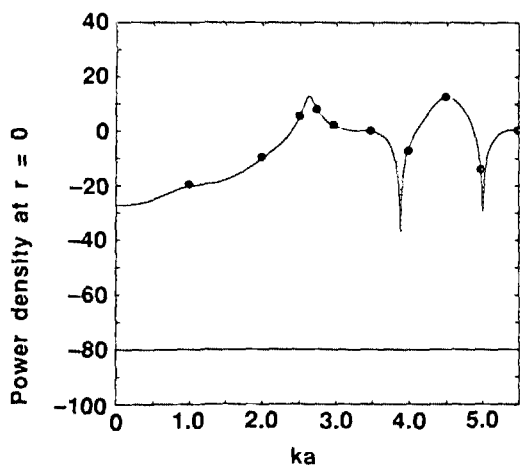


FIG. 14. A scan in  $ka$  of the total energy density of the field at the origin normalized to the energy density of the incident field there for an open spherical shell having an aperture angle of  $150^\circ$  and a plane wave incident at  $\theta^{inc} = 0^\circ$ . The solid line is generated by the dual series solution, the dots by a MoM surface patch code.

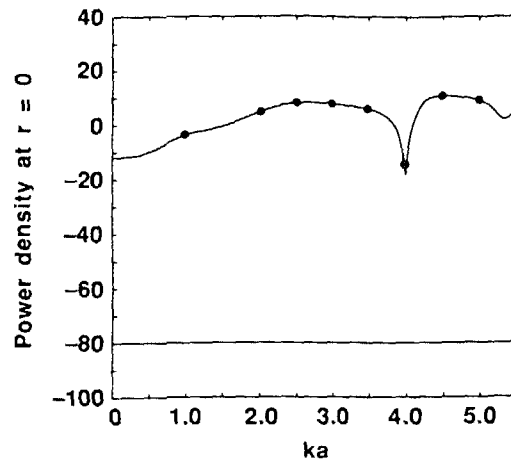


FIG. 15. A scan in  $ka$  of the total energy density of the field at the origin normalized to the energy density of the incident field there for an open spherical shell having an aperture angle of  $120^\circ$  and a plane wave incident at  $\theta^{inc} = 0^\circ$ . The solid line is generated by the dual series solution, the dots by a MoM surface patch code.

are not. Similarly, the antiresonance depressions in the  $\theta_0 = 150^\circ$  data are absent. Also, the levels they predicted for small  $ka$  are close to 50 dB smaller in the  $\theta_0 = 170^\circ$  data and 20 dB smaller in the  $\theta_0 = 150^\circ$  data than ours. Analogously, their resonance peak levels are higher than those we predict. Our results again have been validated with a method of moments calculation.<sup>19</sup> Sample data points from those checks have been included in Figs. 13–15. Agreement is very good.

## VIII. SUMMARY

A complete solution of the scattering of a plane wave from a spherical shell having a circular aperture was developed in this paper. The angle of incidence and the polarization of the plane wave were arbitrary. This solution was constructed with a dual series equations approach and was validated in several different ways. Numerical results were given for the case of normal incidence. Induced currents on the open spherical shell were presented and it was demonstrated that they satisfy Meixner's edge conditions. Energy density scans in  $ka$  were also given; they were dominated by resonance features characteristic of the open spherical cavity.

Several new concepts and techniques were reported. The associated Legendre functions  $P_n^{-m}$  for  $0 \leq n < m$  and their duals  $\bar{P}_n^{-m}$  were introduced to produce a system of pseudodecoupled TE and TM dual series equations and to insure satisfaction of Meixner's edge conditions. Procedures were described in detail that generated an analytical solution of the resulting, previously untreated dual series systems and a numerical solution of the resulting infinite system of linear equations for the modal coefficients. Analytical preconditioning of the current sums led to results free of any Gibbs oscillations. The resonance features in the  $ka$  scans of an energy density ratio at the origin were observed to be predominantly of an antiresonant form.

Cross-section and stored energy calculations are currently in progress. Preliminary cross-section results are also dominated by antiresonance features and suggest that they

are characteristic of a cavity-backed aperture. These studies are summarized in Ref. 40.

## ACKNOWLEDGMENTS

The authors would like to extend special thanks to Art Gautesen of the Department of Mathematics at Iowa State University and to Larry Warne at Sandia National Laboratories for a variety of thoughts and discussions that led to the successful conclusion of our work. The authors would also like to thank Brian Cabral of the Computer Systems Research Group at Lawrence Livermore National Laboratory for developing the graphics package that created the three-dimensional color figures.

This work was performed by the Lawrence Livermore National Laboratory and Sandia National Laboratories under the auspices of the U. S. Department of Energy under Contracts No. W-7405-ENG-48 and DE-AC04-76DP00789, respectively.

## APPENDIX: ASYMPTOTIC BEHAVIOR OF THE TERMS $\chi_n^\phi$ AND $\chi_n^\psi$

The small argument and large index behavior of the terms  $\chi_n^\phi$  and  $\chi_n^\psi$  will be developed for  $n > 0$ . The spherical Bessel function expansions (see Ref. 37, Eqs. 10.1.2 and 10.1.3)

$$j_n(x) = \frac{(\pi/4)^{1/2}(x/2)^n}{\Gamma(n + \frac{1}{2})} \times \left[ 1 - \frac{x^2}{2(2n+3)} + \mathcal{O}\left(\frac{x^4}{n^2}\right) \right], \quad (A1)$$

$$h_n(x) = \frac{-i}{(4\pi)^{1/2}} \left(\frac{x}{2}\right)^{-(n+1)} \Gamma\left(n + \frac{1}{2}\right) \times \left[ 1 + \frac{x^2}{2(2n-1)} + \mathcal{O}\left(\frac{x^4}{n^2}\right) \right], \quad (A2)$$

and the identity

$$\Gamma\left(n + \frac{1}{2}\right) = \{[(2n-1) \cdot \dots \cdot 5 \cdot 3 \cdot 1]/2^n\} \pi^{1/2} \quad (A3)$$

provide the necessary expressions. Equations (A1) and (A2) yield

$$[xj_n(x)]' = \left(\frac{\pi}{4}\right)^{1/2} \frac{(n+1)}{\Gamma(n + \frac{1}{2})} \left(\frac{x}{2}\right)^n \times \left[ 1 - \frac{(n+3)x^2}{2(n+1)(2n+3)} + \mathcal{O}\left(\frac{x^4}{n^2}\right) \right], \quad (A4)$$

$$[xh_n(x)]' = \frac{i}{(4\pi)^{1/2}} n \Gamma\left(n + \frac{1}{2}\right) \left(\frac{x}{2}\right)^{-(n+1)} \times \left[ 1 + \frac{(n-2)x^2}{2n(2n-1)} + \mathcal{O}\left(\frac{x^4}{n^2}\right) \right]. \quad (A5)$$

Combining (A1)–(A3) gives

$$j_n(x)h_n(x) = \frac{1}{i(2n+1)x} \times \left[ 1 + \frac{2x^2}{(2n-1)(2n+3)} + \mathcal{O}\left(\frac{x^4}{n^2}\right) \right]; \quad (A6)$$

combining (A3)–(A5) gives

$$[xj_n(x)]' [xh_n(x)]' = \frac{n(n+1)}{-i(2n+1)x} \times \left[ 1 - \frac{(2n^2+2n+3)x^2}{n(n+1)(2n-1)(2n+3)} + \mathcal{O}\left(\frac{x^4}{n^2}\right) \right]. \quad (A7)$$

Consequently, for small arguments

$$\lim_{x \rightarrow 0} \chi_n^\phi = \lim_{x \rightarrow 0} \{ [i(2n+1)xj_n(x)h_n(x)] - 1 \} \sim 0, \quad (A8)$$

$$\lim_{x \rightarrow 0} \chi_n^\psi = \lim_{x \rightarrow 0} \left\{ \left[ \frac{-i(2n+1)x}{n(n+1)} \times [xj_n(x)]' [xh_n(x)]' \right] - 1 \right\} \sim 0, \quad (A9)$$

and for indices larger than the argument

$$\lim_{n \rightarrow \infty} \chi_n^\phi(x) = \lim_{n \rightarrow \infty} \mathcal{O}(x^2/n^2) \sim 0, \quad (A10)$$

$$\lim_{n \rightarrow \infty} \chi_n^\psi(x) = \lim_{n \rightarrow \infty} \mathcal{O}(x^2/n^2) \sim 0. \quad (A11)$$

Note that this limiting behavior is responsible for the number of terms required for convergence of the solution. In particular, for large enough  $N$ , the terms

$$\chi_N^\phi(ka) \sim \chi_N^\psi(ka) \sim (ka/N)^2$$

and the elements of the matrix  $\mathcal{M}_{ij}$  in (4.49) are small. This explains the choice  $N = 10ka$  for the examples.

- <sup>1</sup>A. M. Radin and V. P. Shestopalov, USSR Comp. Math. Math. Phys. **14**(5), 137 (1974).
- <sup>2</sup>A. M. Radin and V. P. Shestopalov, Sov. Phys. Dokl. **18**, 642 (1974).
- <sup>3</sup>A. M. Radin, V. A. Rezunencko, and V. P. Shestopalov, USSR Comp. Math. Math. Phys. **17**(2), 104 (1977).
- <sup>4</sup>S. S. Vinogradov and V. P. Shestopalov, Sov. Phys. Dokl. **22**, 638 (1977).
- <sup>5</sup>S. S. Vinogradov, Yu A. Tuchkin, and V. P. Shestopalov, Sov. Phys. Dokl. **23**, 650 (1978).
- <sup>6</sup>S. S. Vinogradov, Yu A. Tuchkin, and V. P. Shestopalov, Sov. Phys. Dokl. **25**, 531 (1980).
- <sup>7</sup>S. S. Vinogradov, Yu A. Tuchkin, and V. P. Shestopalov, Sov. Phys. Dokl. **26**, 169 (1981).
- <sup>8</sup>S. S. Vinogradov and V. P. Shestopalov, Sov. Phys. Dokl. **26**, 314 (1981).
- <sup>9</sup>S. S. Vinogradov, Radiophys. Quantum Electron. **26**, 78 (1983).
- <sup>10</sup>K. F. Casey, Proceedings of the National Radio Science Meeting, June 1981, p. 69.
- <sup>11</sup>S. Chang and T. B. A. Senior, U. S. Air Force Weapons Lab. Interaction Note 141, Albuquerque, NM, April 1969.
- <sup>12</sup>T. B. A. Senior and G. A. Desjardins, U. S. Air Force Weapons Lab. Interaction Note 142, Albuquerque, NM, August 1973.
- <sup>13</sup>T. B. A. Senior and G. A. Desjardins, IEEE Trans. Electromag. Compat. **EMC-16**, 205 (1974).
- <sup>14</sup>T. B. A. Senior, U. S. Air Force Weapons Lab. Interaction Note 220, Albuquerque, NM, January 1975.
- <sup>15</sup>R. K. Jones and T. H. Shumpert, IEEE Trans. Antennas Propagat. **AP-28**, 128 (1980).
- <sup>16</sup>J. Meixner, Nachr. Akad. Wiss. Gottingen **2**, 74 (1946).
- <sup>17</sup>J. Meixner, A. Naturforsch. **3a**, 506 (1948).
- <sup>18</sup>R. W. Ziolkowski and W. A. Johnson, Radio Sci. **22**, 169 (1987).
- <sup>19</sup>W. A. Johnson and R. W. Ziolkowski, Proceedings of the National Radio Science Meeting, June 1984, p. 174.
- <sup>20</sup>W. A. Johnson and R. W. Ziolkowski, Radio Sci. **19**, 275 (1984).
- <sup>21</sup>R. W. Ziolkowski, W. A. Johnson, and K. F. Casey, Radio Sci. **19**, 1425 (1984).

- <sup>22</sup>R. W. Ziolkowski, *SIAM J. Math. Anal.* **16**, 358 (1985).
- <sup>23</sup>R. W. Ziolkowski and J. B. Grant, to be published in *IEEE Trans. Antennas Propag.*
- <sup>24</sup>R. W. Ziolkowski and J. B. Grant, *IEEE Trans. Microwave Theor. Tech.* **MTT-34**, 1164 (1986).
- <sup>25</sup>R. W. Ziolkowski and W. A. Johnson, UCRL-92940, Lawrence Livermore National Laboratory, Livermore, CA, July 1985.
- <sup>26</sup>D. S. Jones, *The Theory of Electromagnetism* (Pergamon, New York, 1964).
- <sup>27</sup>M. Born and E. Wolf, *Principles of Optics* (Macmillan, New York, 1964), 2nd ed.
- <sup>28</sup>C. H. Wilcox, *J. Math. Mech.* **6**, 167 (1957).
- <sup>29</sup>C. J. Bouwkamp and H. B. G. Casimir, *Physica* **20**, 539 (1954).
- <sup>30</sup>J. H. Bruning and Y. T. Lo, *IEEE Trans. Antennas Propagat.* **AP-19**, 378 (1971); **AP-19**, 391 (1971).
- <sup>31</sup>C. Liang and Y. T. Lo, *Radio Sci.* **2**, 1481 (1967).
- <sup>32</sup>I. S. Gradshteyn and I. M. Ryzhik, *Table of Integrals, Series, and Products* (Academic, New York, 1965).
- <sup>33</sup>T. M. MacRobert, *Spherical Harmonics* (Pergamon, New York, 1967), 3rd ed.
- <sup>34</sup>E. W. Hobson, *The Theory of Spherical and Ellipsoidal Harmonics* (Chelsea, New York, 1955).
- <sup>35</sup>L. Robin, *Fonctions Spheriques de Legendre et Fonctions Spheroidales* (Gauthier-Villars, Paris, 1958).
- <sup>36</sup>F. Oberhettinger, *Fourier Expansions: A Collection of Formulas* (Academic, New York, 1973).
- <sup>37</sup>*Handbook of Mathematical Functions*, edited by M. Abramowitz and I. A. Stegun (Dover, New York, 1970).
- <sup>38</sup>R. F. Harrington, *Time-Harmonic Electromagnetic Fields* (McGraw-Hill, New York, 1977).
- <sup>39</sup>J. D. Jackson, *Classical Electrodynamics* (Wiley, New York, 1962).
- <sup>40</sup>R. W. Ziolkowski, to be published in *Radio Sci.*



# A geometric setting for internal motions of the quantum three-body system

Toshihiro Iwai

*Department of Applied Mathematics and Physics, Kyoto University, Kyoto 606, Japan*

(Received 15 May 1986; accepted for publication 7 January 1987)

Quantum mechanics for internal motions of the three-body system is set up on the basis of the complex vector bundle theory. The three-body system is called a triatomic molecule in the Born–Oppenheimer approximation. The internal states of the molecule are described as cross sections in the complex vector bundle assigned by an eigenvalue of the square of the total angular momentum operator. This bundle is equipped with a linear connection, which is a natural consequence of a geometric interpretation of the so-called Eckart condition. The coupling of the internal motion with the rotation is understood naturally in terms of this connection. The internal Hamiltonian operator is obtained which includes the internal motion–rotation coupling and a centrifugal potential. The complex vector bundle for the triatomic molecule proves to be a trivial bundle, though the geometric setting for the internal motion is independent of whether the bundle is trivial or not.

## I. INTRODUCTION

Quantum mechanics of a few-body system has been drawing increasing interest in quantum chemistry.<sup>1</sup> A few-body system is called a molecule in the Born–Oppenheimer approximation. Theoretical treatment of nonrigid molecules has been related more or less with the Eckart frame.<sup>2</sup> Tachibana and the author<sup>3</sup> have discussed quantum mechanics of nonrigid molecules without using the Eckart frame, but they did not refer to the geometric setting in the large. This paper is a continuation of the previous paper<sup>3</sup> with particular interest in an application of the theory of connections in complex vector bundles.<sup>4</sup>

It is Guichardet<sup>5</sup> who observed that the center-of-mass system is made into a principal fiber bundle with a rotation group as structure group, and is endowed with a connection by the Eckart condition. Using the holonomy theorem,<sup>4</sup> he proved that the vibration cannot, in general, be separated from the rotation.

On the basis of Guichardet's observation, the present author<sup>6</sup> showed that a moving frame, called the Eckart frame, exists relative to which the molecule moves without rotation, but it depends on a choice of the molecular motion and is not unique for any molecular configuration. For this reason the Eckart frame is not suited for a description of quantum mechanics of nonrigid molecules.

The organization of this article is outlined in the following way: Section II deals with the center-of-mass system. It is shown that the center-of-mass system for a triatomic molecule is made into a principal  $SO(3)$  bundle  $\hat{Q}$  and that the base manifold  $M: = \hat{Q}/SO(3)$ , called the internal space, is diffeomorphic to  $\mathbb{R}_+^3 = \{(x, y, z) \in \mathbb{R}^3; z > 0\}$ . Since the base manifold is contractible to a point, this  $SO(3)$  bundle becomes a trivial bundle<sup>7</sup>;  $\hat{Q} \cong \mathbb{R}_+^3 \times SO(3)$ . Any point of  $\hat{Q}$  can then be assigned uniquely by the internal coordinates in  $\mathbb{R}_+^3$  and the Euler angles in  $SO(3)$ . A coordinate system, called the Dragt coordinate system, is introduced on  $\hat{Q}$  so that the Euler angles and the internal coordinates may assign, respectively, the principal axes of moment of inertia and the molecular configuration relative to the principal axes.

Section III is concerned with the connection due to Guichardet. The connection form and the curvature form are defined on the  $SO(3)$  bundle  $\hat{Q}$  and expressed in the Dragt coordinate system. It is to be noted that even though the  $\hat{Q}$  is a trivial bundle, the connection has nonvanishing curvature, i.e., the  $\hat{Q}$  is not flat with respect to this connection.

Section IV is a review of the total angular momentum operator. The angular momentum operator is obtained as the infinitesimal generator of the  $SO(3)$  action on  $\hat{Q}$ . The right and left actions of  $SO(3)$  are strictly distinguished. The left and right actions give rise to the angular momentum operator with respect to a fixed frame and to the principal axis frame, respectively.

In Sec. V, the connection is discussed again in terms of vector fields. The connection is by definition an assignment of the vibrational (or horizontal) subspace of the tangent space at every point of  $\hat{Q}$ . Rotational and vibrational vectors are discussed to demonstrate that the connection theory naturally fits mechanics of several-particle systems. The connection form is in fact dual to the angular momentum.

Section VI discusses the metric and the volume element induced on the internal space. The results will be used for constructing the internal Hamiltonian operator in Sec. VIII.

In Sec. VII, the complex vector bundles  $V_l$ ,  $l$  being non-negative integers, are introduced in association with the  $SO(3)$  bundle  $\pi: \hat{Q} \rightarrow M$ . The linear connection and the curvature are defined in  $V_l$  and expressed in the Dragt coordinate system. Like the  $SO(3)$  bundle  $\hat{Q}$ , the  $V_l$  is a trivial bundle.

In Sec. VIII, quantum mechanics for internal states is established on the vector bundle  $V_l$  together with the assigned internal Hamiltonian operator  $H_l$  which acts on cross sections in  $V_l$ . The energy functional is used to derive the operator  $H_l$  from the usual Hamiltonian operator acting on wave functions on the center-of-mass system. The  $H_l$  is a matrix-valued second-order differential operator which contains both covariant derivation operators with respect to the linear connection introduced in Sec. VII and a matrix-valued centrifugal potential arising from the conservation of the total angular momentum. This section contains also remarks about internal Hamiltonian operators.

The result obtained can be interpreted in terms of gauge theory as follows: The rotation of the molecule induces a gauge field (the curvature introduced in Sec. VII). The internal motion is coupled with the gauge field through the gauge potential (the connection introduced in Sec. VII). The gauge field plays the role of a magnetic field on the internal space, and may be called the Coriolis field. The internal Hamiltonian describes the internal motion in coupling with the Coriolis field under the presence of the centrifugal potential.

The three-body problem in  $\mathbb{R}^3$  is in marked contrast with the same problem in  $\mathbb{R}^2$ . In fact, in the former case, the  $\text{SO}(3)$  bundle and the associated complex vector bundles are both trivial, but in the latter case, the  $\text{SO}(2)$  bundle and the associated complex line bundles are neither trivial.<sup>8</sup>

## II. SETTINGS IN THE CENTER-OF-MASS SYSTEM

### A. An orthonormal system

Consider a three-body system in  $\mathbb{R}^3$ , which we call a triatomic molecule (or a molecule in short) in the Born–Oppenheimer approximation. Let  $y_k$  and  $m_k$ ,  $k = 1, 2, 3$ , denote the position vector and the mass of each particle, respectively. Then the configuration space  $Q_0$  of the molecule is the linear space of all the triples  $y = (y_1, y_2, y_3)$ . The center-of-mass system  $Q$  is defined as the linear subspace of  $Q_0$  defined by

$$\sum_{k=1}^3 m_k y_k = 0. \quad (2.1)$$

We define the inner product  $K$  on  $Q_0$  by

$$K(x, y) = \sum m_k (x_k | y_k), \quad (2.2)$$

where the round brackets denote the standard inner product in  $\mathbb{R}^3$ . The induced inner product in  $Q$  will be also denoted by  $K$  without any confusion.

The rotation group  $\text{SO}(3)$  has a natural action on the configuration space  $Q_0$ ;

$$y = (y_1, y_2, y_3) \rightarrow gy = (gy_1, gy_2, gy_3), \quad (2.3)$$

where  $g \in \text{SO}(3)$ . Since Eq. (2.1) is invariant under the  $\text{SO}(3)$  action, the  $\text{SO}(3)$  also acts on the center-of-mass system  $Q$ .

We start with the following proposition.

*Proposition 1:* The following system  $\{c_k, f_k, f_{3+k}\}$ ,  $k = 1, 2, 3$ , is an orthonormal system in  $Q_0$  with respect to the inner product  $K$ :

$$c_k = N_0(e_k, e_k, e_k), \quad k = 1, 2, 3, \\ f_k = N_1(m_3 e_k, 0, -m_1 e_k), \quad k = 1, 2, 3, \quad (2.4)$$

$$f_{3+k} = N_2(-m_2 e_k, (m_1 + m_3)e_k, -m_2 e_k), \quad k = 1, 2, 3,$$

where  $N_j$ ,  $j = 0, 1, 2$ , are normalization constants given by

$$N_0 = (m_1 + m_2 + m_3)^{-1/2}, \\ N_1 = (m_1 m_3 (m_1 + m_3))^{-1/2}, \\ N_2 = (m_2 (m_1 + m_3) (m_1 + m_2 + m_3))^{-1/2}. \quad (2.5)$$

It is a matter of calculation to verify that  $\{c_k, f_k, f_{3+k}\}$ ,  $k = 1, 2, 3$ , form an orthonormal system. The vectors  $\{f_j\}$ ,  $j = 1, 2, \dots, 6$ , span the center-of-mass system  $Q$ , and the  $\{c_k\}$ ,

$k = 1, 2, 3$ , are the basis of  $Q^\perp$ , the orthogonal complement of  $Q$ . The space  $Q^\perp$  is identified with the space of center-of-mass vectors. In fact, let  $B$  denote the center-of-mass vector

$$B = \sum_{k=1}^3 m_k y_k \left( \sum_{k=1}^3 m_k \right)^{-1} = \sum_{k=1}^3 B^k e_k. \quad (2.6)$$

Then any triple  $(y_1, y_2, y_3)$  is as usual broken up into

$$(y_1, y_2, y_3) = (B, B, B) + (x_1, x_2, x_3) \quad (2.7)$$

with  $\sum m_k x_k = 0$ . It is now easy to see that

$$(B, B, B) = \sum_{k=1}^3 \frac{B^k c_k}{N_0}, \quad (2.8)$$

which assures the above assertion. The triple  $x = (x_1, x_2, x_3)$  is of course expressed in terms of  $\{f_j\}$ ,  $j = 1, \dots, 6$ ;

$$x = \sum_{j=1}^6 q^j f_j, \quad q^j = K(x, f_j). \quad (2.9)$$

Thus we may consider  $(B^k / N_0, q^j)$ ,  $k = 1, 2, 3$ ,  $j = 1, \dots, 6$ , as the Cartesian coordinates in  $Q_0$ ;

$$y = \sum_{k=1}^3 \frac{B^k c_k}{N_0} + \sum_{j=1}^6 q^j f_j. \quad (2.10)$$

It is now an easy matter to calculate  $K(gf_i, f_j)$ ,  $i, j = 1, 2, \dots, 6$ , etc., for  $g \in \text{SO}(3)$  in order to express the  $\text{SO}(3)$  action in the block diagonal form

$$\begin{pmatrix} g & & \\ & g & \\ & & g \end{pmatrix}, \quad g \in \text{SO}(3), \quad (2.11)$$

with respect to the basis  $\{c_k, f_k, f_{3+k}\}$ ,  $k = 1, 2, 3$ , where missing matrix entries are all zero.

In view of (2.11), we are allowed to think of the center-of-mass system  $Q \simeq \mathbb{R}^6$  as a product space  $\mathbb{R}^3 \times \mathbb{R}^3$ , the first factor spanned by  $f_k$ , and the second by  $f_{3+k}$ ,  $k = 1, 2, 3$ . The  $\text{SO}(3)$  action on  $Q \simeq \mathbb{R}^3 \times \mathbb{R}^3$  is then a diagonal one. Moreover, the vectors  $\sum q^k f_k$  and  $\sum q^{3+k} f_{3+k}$  can be represented in the original space  $\mathbb{R}^3$ . In effect, on setting  $r^k = q^k$  and  $s^k = q^{3+k}$ ,  $k = 1, 2, 3$ , to define the vectors  $r = \sum r^k e_k$  and  $s = \sum s^k e_k$  in  $\mathbb{R}^3$ , we obtain, after a straightforward calculation,

$$r = ((m_1 m_3) / (m_1 + m_3))^{1/2} (x_1 - x_3), \\ s = \left( \frac{m_2 (m_1 + m_3)}{m_1 + m_2 + m_3} \right)^{1/2} \left( x_2 - \frac{m_1 x_1 + m_3 x_3}{m_1 + m_3} \right). \quad (2.12)$$

These are Jacobi vectors used often in the three-body problem.<sup>9,10</sup> Thus the points  $x$  of the center-of-mass system can be thought of as the pairs of vectors  $(r, s)$ .

### B. The internal space

Following Guichardet,<sup>5</sup> we make the center-of-mass system  $Q$  into a principal fiber bundle with structure group  $\text{SO}(3)$ . To this end, we have to investigate whether the  $\text{SO}(3)$  action on  $Q$  is free or not. Suppose that  $gr = r$  and  $gs = s$  for some vectors  $r$  and  $s$ . If  $r$  and  $s$  are linearly independent,  $g$  must be the identity. If they are linearly dependent,  $g$  need not be the identity. Therefore setting

$$D = \{(r, s) \in \mathbb{R}^3 \times \mathbb{R}^3; \lambda r + \mu s = 0 \text{ with } (\lambda, \mu) \neq (0, 0)\},$$

we find that the  $SO(3)$  acts freely on

$$\dot{Q} = (\mathbb{R}^3 \times \mathbb{R}^3) - D, \quad (2.13)$$

so that the quotient space  $M := \dot{Q}/SO(3)$  is a manifold, called the internal space. Thus we obtain a principal fiber bundle  $\pi: \dot{Q} \rightarrow M$  with structure group  $SO(3)$ . We note here that the excluded set  $D$  corresponds to the configurations in which three particles lie in the same line, and two or three of them may happen to collide.

Now we wish to study the topology of the internal space. For this purpose we consider the mapping

$$(r, s) \rightarrow (|r|^2 - |s|^2, 2(r|s), 2|r \times s|). \quad (2.14)$$

Note that the quantities in the mapping image are invariant under the  $SO(3)$  action. Since the set  $D$  is excluded,  $|r \times s|$  must be positive. The quantities  $|r|^2 - |s|^2$  and  $(r|s)$  range over all the real numbers. Hence, by (2.14),  $\dot{Q}$  is mapped onto  $\mathbb{R}_+^3$ , the half-space of  $\mathbb{R}^3$ . Furthermore, it is easy to verify that

$$(|r|^2 - |s|^2)^2 + 4(r|s)^2 + 4|r \times s|^2 = (|r|^2 + |s|^2)^2. \quad (2.15)$$

Conversely, when given a point  $(w^k)$  of  $\mathbb{R}_+^3$ , the inverse image of the mapping (2.14) is the set of solutions to the coupled equations  $|r|^2 - |s|^2 = w^1$ ,  $2(r|s) = w^2$ ,  $2|r \times s| = w^3 > 0$ . In solving these equations, we see from (2.15) that  $|r|^2 + |s|^2$  is determined in terms of  $w^k$ ,  $k = 1, 2, 3$ , so that  $|r|$  and  $|s|$  become known. Moreover, the angle made by two vectors  $r$  and  $s$  is determined by  $2(r|s) = w^2$ . Thus we can obtain a triangle formed by a certain pair of linearly independent vectors  $r_0$  and  $s_0$  in  $\mathbb{R}^3$ . The solutions are then the  $SO(3)$  orbit of  $r_0$  and  $s_0$ . This shows that the inverse image of  $(w^k) \in \mathbb{R}_+^3$  under the mapping (2.14) is topologically  $SO(3)$ . Therefore we obtain the following theorem.

**Theorem 2:** The internal space  $M = \dot{Q}/SO(3)$  is topologically  $\mathbb{R}_+^3$ . Since  $\mathbb{R}_+^3$  is contractible to a point, the principal fiber bundle  $\pi: \dot{Q} \rightarrow M$  becomes a trivial bundle,<sup>7</sup> where  $\pi$  is equivalently given by (2.14).

The triviality of the bundle  $\pi: \dot{Q} \rightarrow M$  implies the existence of a cross section<sup>7</sup>  $\sigma_0: M \rightarrow \dot{Q}$ ,  $\pi \circ \sigma_0 = \text{id}$ . This fact was known tacitly by Dragt,<sup>11</sup> and Levy-Leblond and Levy-Nahas,<sup>12</sup> and used effectively in analyzing the states of noninteracting three particles. According to them,  $r$  and  $s$  can be expressed in the form

$$\begin{aligned} r &= \rho(\cos(\psi/2)\cos(\chi/2)u_1 - \sin(\psi/2)\sin(\chi/2)u_2), \\ s &= \rho(\sin(\psi/2)\cos(\chi/2)u_1 + \cos(\psi/2)\sin(\chi/2)u_2), \end{aligned} \quad (2.16)$$

where  $(\rho, \chi, \psi)$  are internal coordinates subject to

$$0 < \rho < +\infty, \quad 0 < \chi < \pi/2, \quad 0 \leq \psi < 2\pi,$$

and  $u_k = ge_k$ ,  $k = 1, 2, 3$ , are a moving frame such that the molecule is set on the plane spanned by  $u_1$  and  $u_2$ . We will refer to these coordinates as Dragt's coordinates. In the succeeding section we will see that the  $u_k$  are chosen so as to lie in the direction of the principal axes of moment of inertia for the molecule. It is also to be noted that the determinant of the coefficient matrix of (2.16) is  $\frac{1}{2} \sin \chi$ , so that  $\chi \neq 0$ . Now Eq. (2.16) with  $u_k = e_k$ , i.e.,  $g = \text{id}$ , provides a cross section  $\sigma_0: M \rightarrow \dot{Q}$ . Every element of  $\dot{Q}$  is then given in the form  $g\sigma_0(w)$ ,

$g \in SO(3)$ ,  $w \in M$ . Here we remark that the structure group acts on  $\dot{Q}$  to the left, contrary to the usual convention according to which the structure group acts on the principal fiber bundle to the right.

From (2.16), the projection  $\pi$  given by (2.14) takes the form

$$\begin{aligned} w^1 &= |r|^2 - |s|^2 = \rho^2 \cos \psi \cos \chi, \\ w^2 &= 2(r|s) = \rho^2 \sin \psi \cos \chi, \\ w^3 &= 2|r \times s| = \rho^2 \sin \chi > 0, \end{aligned} \quad (2.17)$$

so that  $(\rho^2, \psi, \pi/2 - \chi)$  can be thought of as the spherical coordinates in  $\mathbb{R}_+^3$ .

We conclude this section by fixing the Euler angles as follows:

$$\begin{aligned} u_1 &= e_1(\cos \beta \cos \alpha \cos \gamma - \sin \alpha \sin \gamma) \\ &\quad + e_2(\cos \beta \sin \alpha \cos \gamma + \cos \alpha \sin \gamma) \\ &\quad + e_3(-\sin \beta \cos \gamma), \\ u_2 &= e_1(-\cos \beta \cos \alpha \sin \gamma - \sin \alpha \cos \gamma) \\ &\quad + e_2(-\cos \beta \sin \alpha \sin \gamma + \cos \alpha \cos \gamma) \\ &\quad + e_3(\sin \beta \sin \gamma), \\ u_3 &= e_1 \sin \beta \cos \alpha + e_2 \sin \beta \sin \alpha + e_3 \cos \beta, \end{aligned} \quad (2.18)$$

where  $0 \leq \alpha < 2\pi$ ,  $0 \leq \beta < \pi$ ,  $0 \leq \gamma < 2\pi$ .

### III. THE CONNECTION DUE TO GUICHARDET

#### A. A review of the connection

In Ref. 6, using the connection form due to Guichardet,<sup>5</sup> we have shown that there exists the Eckart frame, a frame relative to which the molecule moves without rotation, along any curve  $x(t)$  in the center-of-mass system, but it depends inevitably on the choice of  $x(t)$ . Hence the Eckart frame is not suitable for description of molecular motions in quantum as well as classical mechanics.

The connection form is defined as follows<sup>5,6</sup>: Let  $\Lambda^2 \mathbb{R}^d$  denote the space of antisymmetric tensors of order 2 on  $\mathbb{R}^d$ , and  $\mathfrak{so}(d)$  the Lie algebra of  $SO(d)$ . A linear isomorphism  $R$  of  $\Lambda^2 \mathbb{R}^d$  to  $\mathfrak{so}(d)$  is defined for  $\xi = \sum_{i < j} \xi_{ij} e_i \wedge e_j$  and  $x = \sum x^i e_i$  by

$$R_\xi(x) = \sum \left( \sum \xi_{ij} x^j \right) e_i. \quad (3.1)$$

The inertia operator  $A_x: \Lambda^2 \mathbb{R}^d \rightarrow \Lambda^2 \mathbb{R}^d$  is defined on the center-of-mass system by

$$A_x(\xi) = - \sum m_k x_k \wedge R_\xi(x_k), \quad (3.2)$$

which is symmetric and positive definite, and  $k$  ranges from 1 to  $N$ , the number of particles. The connection form  $\omega$  on the center-of-mass system is then given by

$$\omega = R \left( -A_x^{-1} \sum m_k x_k \wedge dx_k \right), \quad (3.3)$$

where  $R(\xi)$  stands for  $R_\xi$  for notational convenience.

In our case,  $d = 3$ . If we set  $\xi_{12} = \phi^3$ ,  $\xi_{23} = \phi^1$ ,  $\xi_{31} = \phi^2$ , the two-vector  $\xi = \sum_{i < j} \xi_{ij} e_i \wedge e_j$  is identified with  $\phi = \sum \phi^i e_i$ . Put another way,  $\Lambda^2 \mathbb{R}^3$  is identified with  $\mathbb{R}^3$  by

$e_1 \wedge e_2 \rightarrow e_3$  and the cyclic permutations. Accordingly,  $R$  becomes a linear isomorphism of  $\mathbb{R}^3$  to  $\mathfrak{so}(3)$ ;

$$R_{\xi}(x) = R_{\phi}(x) = -\phi \times x, \quad \text{for } x \in \mathbb{R}^3.$$

Alternatively,  $R(e_1)$  is the matrix  $(\xi_{ij})$  with nonzero elements  $\xi_{23} = -\xi_{32} = 1$  only, and so on. It should be noted that  $R$  is Ad-equivariant in the sense that

$$R(g\phi) = \text{Ad}_g R(\phi) = gR(\phi)g^{-1}. \quad (3.4)$$

With the above identification in mind, we treat  $A_x$  and  $\omega$  in the form

$$A_x(\phi) = \sum m_k x_k \times (\phi \times x_k), \quad (3.5)$$

$$\omega = R \left( -A_x^{-1} \sum m_k x_k \times dx_k \right). \quad (3.6)$$

We note again that  $A_x$  is symmetric and positive definite, because

$$(\phi' | A_x(\phi)) = \sum m_k (\phi' \times x_k | \phi \times x_k). \quad (3.7)$$

Moreover, using (3.4), we can verify that

$$A_{gx}(\phi) = gA_x(g^{-1}\phi), \quad (3.8)$$

$$\lambda_g^* \omega = \text{Ad}_g \omega = g\omega g^{-1}, \quad (3.9)$$

where  $\lambda_g^*$  denotes the pullback by the  $g$  action;  $\lambda_g(x) = gx$ .

A vector field  $v = (v_k)$ ,  $k = 1, 2, 3$ , on  $\dot{Q}$ , which is subject to  $\sum m_k v_k = 0$ , is called vibrational if it satisfies  $\omega(v) = 0$ . From (3.6) this condition becomes equivalent to

$$\sum m_k x_k \times v_k = 0, \quad (3.10)$$

because  $dx_k(v) = v_k$ , and  $R$  and  $A_x$  are nonsingular. Equation (3.10) means that the total angular momentum vanishes for the vibrational vector fields. In this sense Eq. (3.10) is viewed as a generalization of the Eckart condition to any configuration of the molecule. Rotational vector fields are defined as infinitesimal generators of the  $\text{SO}(3)$  action on  $\dot{Q}$ , which are given by  $R_{\phi}(x) := (R_{\phi}(x_k))$ ,  $k = 1, 2, 3$ . We notice that for rotational vector fields  $R_{\phi}(x)$  one has

$$\begin{aligned} \omega(R_{\phi}(x)) &= R \left( -A_x^{-1} \sum m_k x_k \times (-\phi \times x_k) \right) \\ &= R(A_x^{-1} A_x(\phi)) = R_{\phi}. \end{aligned} \quad (3.11)$$

Equations (3.9) and (3.11) are characteristic of the connection form.<sup>4</sup> We notice here that because of the left action Eq. (3.9) is expressed in a form different from the usual one. In what follows we will express  $A_x$  and  $\omega$  in terms of  $r$  and  $s$  and of Dragt's coordinates.

## B. The inertia operator

We start by calculating  $\sum m_k x_k \times dx_k$ . Since  $r$  and  $s$  are given by (2.12), the vectors  $x_k$ ,  $k = 1, 2, 3$ , with  $\sum m_k x_k = 0$ , are expressed as linear combinations of  $r$  and  $s$ . Consequently, a straightforward calculation gives

$$\sum m_k x_k \times dx_k = r \times dr + s \times ds. \quad (3.12)$$

Applying (3.12) to a rotational vector  $R_{\phi}(x)$ , one obtains

$$\sum m_k x_k \times R_{\phi}(x_k) = r \times R_{\phi}(r) + s \times R_{\phi}(s). \quad (3.13)$$

We note here that since  $\text{SO}(3)$  acts on  $\dot{Q}$  in the form  $(gr, gs)$ , rotational vector fields are expressed as  $(R_{\phi}(r), R_{\phi}(s))$ .

From (3.5) and (3.13) it follows that

$$A_x(\phi) = r \times (\phi \times r) + s \times (\phi \times s). \quad (3.14)$$

The matrix elements  $(e_j | A_x(e_k))$  of  $A_x$  are then easy to calculate in terms of  $r$  and  $s$ :

$$\begin{aligned} (e_j | A_x(e_k)) &= -(r^j r^k + s^j s^k) \quad (j \neq k), \\ (e_k | A_x(e_k)) &= |r|^2 + |s|^2 - (r^k)^2 - (s^k)^2. \end{aligned} \quad (3.15)$$

Applying (3.14) to  $u_k = ge_k$ ,  $k = 1, 2, 3$ , together with (2.16), we obtain

$$\begin{aligned} A_x(u_1) &= \rho^2 \sin^2(\chi/2) u_1, \\ A_x(u_2) &= \rho^2 \cos^2(\chi/2) u_2, \\ A_x(u_3) &= \rho^2 u_3. \end{aligned} \quad (3.16)$$

These equations show that the  $u_k$  lie in the direction of the principal axes, and at the same time give the principal moments of inertia.

## C. The connection form

The connection form (3.6) is now given from (3.12) as

$$\omega = R(-A_x^{-1}(r \times dr + s \times ds)), \quad (3.17)$$

where  $A_x^{-1}$  is the inverse of the matrix given in (3.15).

In what follows we are going to calculate  $\omega$  in Dragt's coordinates. For this purpose, the following formula, easy to prove, is of great help.

*Proposition 3:* The frame  $(u_k)$ ,  $k = 1, 2, 3$ , satisfies the relation

$$(du_1, du_2, du_3) = (u_1, u_2, u_3) \begin{pmatrix} 0 & -\Theta^3 & \Theta^2 \\ \Theta^3 & 0 & -\Theta^1 \\ -\Theta^2 & \Theta^1 & 0 \end{pmatrix}, \quad (3.18)$$

where  $\Theta^k$ ,  $k = 1, 2, 3$ , are left-invariant one-forms on  $\text{SO}(3)$  which are expressed as

$$\begin{aligned} \Theta^1 &= \sin \gamma d\beta - \sin \beta \cos \gamma d\alpha, \\ \Theta^2 &= \cos \gamma d\beta + \sin \beta \sin \gamma d\alpha, \\ \Theta^3 &= \cos \beta d\alpha + d\gamma. \end{aligned} \quad (3.19)$$

*Remark:* When we set  $g = (u_k)$  and  $\Theta = -\sum R(e_k)\Theta^k$ , a matrix-valued one-form, Eq. (3.18) is written as  $dg = g\Theta$ , or  $\Theta = g^{-1}dg$ .

*Proposition 4:* The right-invariant one-forms  $\Psi^j$ ,  $j = 1, 2, 3$ , are given by  $\Psi^j = \sum g_{jk}\Theta^k$ ,  $g_{jk}$  being the components of  $g$ . From Eqs. (2.18) and (3.19) one has

$$\begin{aligned} \Psi^1 &= -\sin \alpha d\beta + \sin \beta \cos \alpha d\gamma, \\ \Psi^2 &= \cos \alpha d\beta + \sin \beta \sin \alpha d\gamma, \\ \Psi^3 &= d\alpha + \cos \beta d\gamma. \end{aligned} \quad (3.20)$$

For the proof we use (3.4). Since the right-invariant matrix one-form is defined as  $\Psi = dg g^{-1} = g\Theta g^{-1}$ , one obtains

$$\begin{aligned} g\Theta g^{-1} &= -g \sum R(e_k \Theta^k) g^{-1} \\ &= -\sum R(ge_k \Theta^k) = -\sum R \left( \sum g_{jk} e_j \Theta^k \right). \end{aligned}$$

Thus  $\Psi = -\sum R(e_j)\Psi^j$  with  $\Psi^j = \sum g_{jk}\Theta^k$ . In the course of the calculation, we have also obtained

$$\Psi = -\sum R(e_j)\Psi^j = -\sum R(u_k)\Theta^k. \quad (3.21)$$

Now we proceed to the calculation of  $\omega$ . Using (2.16) and (3.18), we are ready to compute  $r \times dr + s \times ds$  in a straightforward manner to obtain

$$r \times dr + s \times ds = \rho^2 \sin^2(\chi/2)u_1\Theta^1 + \rho^2 \cos^2(\chi/2)u_2\Theta^2 + \rho^2 u_3(\Theta^3 - \frac{1}{2} \sin \chi d\psi). \quad (3.22)$$

Since  $A_x^{-1}$  is given from (3.16), the connection form (3.17) turns out to have the form

$$\omega = \Psi + R(u_3)\frac{1}{2} \sin \chi d\psi. \quad (3.23)$$

We denote the components of  $\omega$  as follows:

$$\omega = \sum R(u_k)\sigma^k = \sum R(e_k)\omega^k. \quad (3.24)$$

On account of Eqs. (3.21) and (3.23), one has

$$\begin{aligned} \sigma^1 &= -\Theta^1, \\ \sigma^2 &= -\Theta^2, \\ \sigma^3 &= -\Theta^3 + \frac{1}{2} \sin \chi d\psi, \\ \omega^1 &= -\Psi^1 + \frac{1}{2} \sin \beta \cos \alpha \sin \chi d\psi, \\ \omega^2 &= -\Psi^2 + \frac{1}{2} \sin \beta \sin \alpha \sin \chi d\psi, \\ \omega^3 &= -\Psi^3 + \frac{1}{2} \cos \beta \sin \chi d\psi. \end{aligned} \quad (3.26)$$

In summary, we have the following.

**Theorem 5:** The connection form  $\omega$  takes the form (3.23) or (3.24) with (3.25) and (3.26) in Dragt's coordinates (2.16).

Since the principal fiber bundle  $\pi: \dot{Q} \rightarrow M$  is trivial (Theorem 2), the connection form  $\omega$  is pulled back on the internal space  $M$  through the cross section  $\sigma_0: M \rightarrow \dot{Q}$ . In fact, setting  $\alpha = \beta = \gamma = 0$  in (3.23), we obtain

$$\sigma_0^*\omega = R(e_3)\frac{1}{2} \sin \chi d\psi. \quad (3.27)$$

#### D. The curvature form

The curvature form  $\Omega$  is given by the structure equation

$$\Omega = d\omega - \omega \wedge \omega. \quad (3.28)$$

Note here that  $\text{SO}(3)$  acts on  $\dot{Q}$  to the left, so that the structure equation takes a different form from the usual one.<sup>4</sup> We wish to express the  $\Omega$  in Dragt's coordinates, as we have done for  $\omega$ . To this end, the following propositions are of great use.

**Proposition 6:** The left- and right-invariant one-forms  $\Theta$  and  $\Psi$  on  $\text{SO}(3)$  are subject to the following Maurer–Cartan equation,<sup>13</sup> respectively:

$$d\Theta = -\Theta \wedge \Theta, \quad d\Psi = \Psi \wedge \Psi. \quad (3.29)$$

**Proposition 7:** The matrices  $R(u_k)$ ,  $k = 1, 2, 3$ , satisfy the commutation relations

$$[R(u_1), R(u_2)] = -R(u_3) \quad (\text{cycl.}). \quad (3.30)$$

To prove these equations, we note that  $[R(e_1), R(e_2)] = -R(e_3)$  (cycl.), which are easy to see. On using (3.4) with  $ge_k = u_k$ , Eq. (3.30) proves to be a consequence of these commutation relations.

We are now in a position to write out (3.28) in Dragt's coordinates. From (3.23) we obtain

$$d\omega = d\Psi + R(u_3)\frac{1}{2} \cos \chi d\chi \wedge d\psi + (R(u_1)\Theta^2 - R(u_2)\Theta^1) \wedge \frac{1}{2} \sin \chi d\psi,$$

where we have used (3.18) to get  $R(du_3) = R(u_1)\Theta^2 - R(u_2)\Theta^1$ . On the other hand, using (3.21) and (3.30) to get

$$R(u_3)\Psi = \Psi R(u_3) + R(u_2)\Theta^1 - R(u_1)\Theta^2,$$

one obtains

$$\omega \wedge \omega = \Psi \wedge \Psi + (R(u_1)\Theta^2 - R(u_2)\Theta^1) \wedge \frac{1}{2} \sin \chi d\psi.$$

Thus the curvature form (3.28) turns out to be

$$\Omega = R(u_3)\frac{1}{2} \cos \chi d\chi \wedge d\psi. \quad (3.31)$$

**Theorem 8:** The curvature form  $\Omega$  is expressed as (3.31) in Dragt's coordinates (2.16).

The  $\Omega$  is pulled back on the internal space through the cross section  $\sigma_0$  to give

$$\sigma_0^*\Omega = R(e_3)\frac{1}{2} \cos \chi d\chi \wedge d\psi. \quad (3.32)$$

From (3.27) and (3.32) we have  $\sigma_0^*\Omega = d\sigma_0^*\omega$ .

#### IV. THE TOTAL ANGULAR MOMENTUM OPERATOR

The total angular momentum operator is defined through the infinitesimal generator of the  $\text{SO}(3)$  action. The operator  $\hat{J}_j$  will be defined as  $i$  times the infinitesimal generator  $J_j$  of the action of  $\exp tR(e_j)$ ,  $j = 1, 2, 3$ . While we have treated the left action of  $\text{SO}(3)$ , the right action is also of great importance in dealing with the total angular momentum operator. Recall Eq. (2.16);  $r$  and  $s$  are assigned by the moving frame  $u_k$  and the internal coordinates. The left action of  $h \in \text{SO}(3)$  maps  $u_k = \sum u_k^i e_i$  to  $u_k' = hu_k = \sum (\sum h_{ji} u_k^i) e_j$  with the Euler angles  $(\alpha', \beta', \gamma')$ , keeping the internal coordinates fixed. This means that we are observing the rotation of the molecule in the fixed frame  $\{e_k\}$ . Contrary to this, under the right action of  $h$ , the  $u_k$  are mapped to  $\sum u_j h_{jk}$ , while the internal coordinates are left invariant. This implies that we are observing the rotation of the molecule in the moving frame  $\{u_k\}$ . In effect, the right action is related to a left action as follows: Let the matrix  $g$  be viewed as a set of the column vectors  $u_k = ge_k$ . Then under the left action of  $\exp tR(u_j)$ , the moving frame  $g = (u_k)$  is mapped as a whole to  $(\exp tR(u_j))g$ , which are put into, by using (3.4),

$$\begin{aligned} (\exp tR(u_j))g &= g(\exp tR(e_j))g^{-1}g \\ &= g(\exp tR(e_j)). \end{aligned} \quad (4.1)$$

Thus the left action of  $\exp tR(u_j)$ , the rotation of the molecule around the  $u_j$  axis, becomes equivalent to the right action of  $\exp tR(e_j)$  on the frame  $g = (u_k)$ .

Though the total angular momentum operator is already well known, we wish to derive them in order to realize the difference between right and left actions.

#### A. The left action

We start with the left  $\text{SO}(3)$  action on a triple  $x \in Q$ . Let  $h(t) = \exp tR(\phi)$  be a one-parameter subgroup of  $\text{SO}(3)$ .

Then its infinitesimal generator is given by

$$\begin{aligned} \sum \left( R(\phi) x_k \left| \frac{\partial}{\partial x_k} \right. \right) &= \sum \left( -\phi \times x_k \left| \frac{\partial}{\partial x_k} \right. \right) \\ &= \left( -\phi \left| \sum x_k \times \frac{\partial}{\partial x_k} \right. \right), \end{aligned} \quad (4.2)$$

where  $\partial/\partial x_k$ ,  $k = 1, 2, 3$ , denote the gradient vectors. Let

$$J = \sum x_k \times \frac{\partial}{\partial x_k} = -\sum e_j J_j. \quad (4.3)$$

Then the infinitesimal generator (4.2) is written as  $(-\phi|J)$ . Since  $J_j = (-e_j|J)$ , each  $J_j$  turns out to be the infinitesimal generator of  $\exp tR(e_j)$ ,  $j = 1, 2, 3$ .

We turn to the expression of  $J$  in  $r$  and  $s$ . Since the action of  $h(t)$  has the form  $(h(t)r, h(t)s)$ , in the same manner as (4.2) we have

$$J = r \times \frac{\partial}{\partial r} + s \times \frac{\partial}{\partial s} = -\sum e_j J_j, \quad (4.4)$$

where  $J_j = (-e_j|J)$  is expressed in  $r$  and  $s$ .

Following the same procedure as the above, we can express  $J_j$ ,  $j = 1, 2, 3$ , in the Euler angles. For any  $r$  and  $s$ , we set  $r(t) = h(t)r$ ,  $s(t) = h(t)s$ , and  $u_k(t) = h(t)u_k$ , where  $u_k(t)$  has the Euler angles  $(\alpha(t), \beta(t), \gamma(t))$  with  $\alpha(0) = \alpha$ ,  $\beta(0) = \beta$ ,  $\gamma(0) = \gamma$ . Then, making use of (3.21), we have

$$\frac{dr(t)}{dt} = \left( \sum \Theta^k \left( \frac{d}{dt} u_k \right) \right) \times r(t), \quad (4.5)$$

and the same equation for  $s(t)$ , where  $d/dt$  stands for the tangent vector to the curve  $r(t)$ . In fact

$$\begin{aligned} \frac{dr(t)}{dt} &= \frac{dh(t)}{dt} h(t)^{-1} r(t) = \Psi \left( \frac{d}{dt} \right) r(t) \\ &= -\sum \Theta^k \left( \frac{d}{dt} \right) R(u_k) r(t) \\ &= \left( \sum \Theta^k \left( \frac{d}{dt} u_k \right) \right) \times r(t). \end{aligned}$$

The same calculation is good for  $s(t)$ .

Let  $h(t) = \exp tR(e_j)$  for a fixed  $j$ . Then  $(dr/dt)(0) = R(e_j)r = -e_j \times r$ , and  $(ds/dt)(0) = R(e_j)s = -e_j \times s$ , so that we have

$$-e_j = \sum \Theta^k \left( \frac{d}{dt} \right) u_k(0).$$

For each  $j$ , this is the equation which determines  $(\dot{\alpha}(0), \dot{\beta}(0), \dot{\gamma}(0))$  and hence  $J_j$ . A straightforward calculation then results in

$$\begin{aligned} J_1 &= \cos \alpha \cot \beta \frac{\partial}{\partial \alpha} + \sin \alpha \frac{\partial}{\partial \beta} - \frac{\cos \alpha}{\sin \beta} \frac{\partial}{\partial \gamma}, \\ J_2 &= \sin \alpha \cot \beta \frac{\partial}{\partial \alpha} - \cos \alpha \frac{\partial}{\partial \beta} - \frac{\sin \alpha}{\sin \beta} \frac{\partial}{\partial \gamma}, \\ J_3 &= -\frac{\partial}{\partial \alpha}. \end{aligned} \quad (4.6)$$

The  $J_j$ ,  $j = 1, 2, 3$ , are known as the right-invariant vector fields on  $SO(3)$ , which satisfy, for the right-invariant one-forms (3.20),

$$\Psi^j(J_k) = -\delta^j_k. \quad (4.7)$$

The total angular momentum operators are defined as  $\hat{J}_k = iJ_k$  to satisfy

$$[\hat{J}_1, \hat{J}_2] = i\hat{J}_3 \quad (\text{cycl.}). \quad (4.8)$$

## B. The right action

The right action of  $h(t)$  is expressed as  $u_k(t) = \sum u_j h_{jk}(t)$ . Differentiation gives  $\dot{u}_k(0) = \sum u_j \dot{h}_{jk}(0)$ . On the other hand, Eq. (3.18) provides

$$(\dot{u}_1(0), \dot{u}_2(0), \dot{u}_3(0)) = (u_1, u_2, u_3) \Theta \left( \frac{d}{dt} \right) \Big|_{t=0}.$$

Hence, for each  $h(t) = \exp tR(e_j)$ , this equation reads

$$\left( \sum u_l \dot{h}_{lk}(0) \right) = (u_k) \Theta \left( \frac{d}{dt} \right) \Big|_{t=0}.$$

This determines  $(\dot{\alpha}(0), \dot{\beta}(0), \dot{\gamma}(0))$ , and therefore  $L_j$ , the infinitesimal generator of the right action of  $\exp tR(e_j)$ . To obtain  $L_j$  is now a matter of calculation:

$$\begin{aligned} L_1 &= \frac{\cos \gamma}{\sin \beta} \frac{\partial}{\partial \alpha} - \sin \gamma \frac{\partial}{\partial \beta} - \cos \gamma \cot \beta \frac{\partial}{\partial \gamma}, \\ L_2 &= -\frac{\sin \gamma}{\sin \beta} \frac{\partial}{\partial \alpha} - \cos \gamma \frac{\partial}{\partial \beta} + \sin \gamma \cot \beta \frac{\partial}{\partial \gamma}, \\ L_3 &= -\frac{\partial}{\partial \gamma}. \end{aligned} \quad (4.9)$$

These are known as the left-invariant vector fields on  $SO(3)$ , which satisfy, for the left-invariant one-forms (3.17),

$$\Theta^k(L_j) = -\delta^k_j. \quad (4.10)$$

The total angular momentum operators are defined as  $\hat{L}_k = -iL_k$  to satisfy

$$[\hat{L}_1, \hat{L}_2] = i\hat{L}_3 \quad (\text{cycl.}). \quad (4.11)$$

It is worth mentioning that  $J_k$  and  $L_k$  are related by

$$\sum_{k=1}^3 g_{jk} L_k = J_j, \quad j = 1, 2, 3, \quad (4.12)$$

as the right- and left-invariant one-forms are related in the same fashion (Proposition 4). From (4.3) and (4.12) it follows that

$$J = -\sum e_j J_j = -\sum u_k L_k. \quad (4.13)$$

This equation is consistent with Eq. (4.1). In fact, from Eq. (4.1) the infinitesimal generator  $L_k$  of the right action of  $\exp tR(e_j)$  must be the infinitesimal generator of the left action of  $\exp tR(u_k)$ . Hence the  $L_k$  must equal  $(-u_k|J)$  on account of (4.2) and (4.3). However, this is also a consequence of (4.13). We note in conclusion that in a different manner the operators  $J_k$  and  $L_k$  are derived in Ref. 14 in the coordinates  $(\psi, \theta, \phi) = (-\alpha, -\beta, -\gamma)$ .

In summary we have the following.

*Proposition 9:* The total angular momentum operator  $\hat{J}$  is expressed as

$$\begin{aligned} \hat{J} &= -i \sum x_k \times \frac{\partial}{\partial x_k} = -i \left( r \times \frac{\partial}{\partial r} + s \times \frac{\partial}{\partial s} \right) \\ &= \sum e_j \hat{J}_j = -\sum u_k \hat{L}_k, \end{aligned} \quad (4.14)$$

where  $\hat{J}_j = iJ_j$  and  $\hat{L}_k = -iL_k$  are given by (4.6) and (4.9), respectively, in the Euler angles.

## V. ROTATIONS AND VIBRATIONS

### A. Rotational and vibrational vectors

We have introduced the connection in Sec. III in terms of differential one-forms, and defined vibrational vector fields as those vector fields for which the connection form vanishes. Rotational vector fields are the infinitesimal generators of the SO(3) action. In this section, we review again the connection in terms of vector fields. Let  $K_x$  denote the inner product induced from  $K$  into the tangent space  $T_x(\dot{Q})$ ;

$$K_x(u, v) = \sum m_k (u_k | v_k). \quad (5.1)$$

We show that the vibrational vector fields are orthogonal to the rotational vector fields with respect to  $K_x$ . In fact, for a vibrational vector field  $v$ , subject to (3.10), and a rotational vector field  $R_\phi(x)$ , we have

$$\begin{aligned} K_x(v, R_\phi(x)) &= \sum m_k (v_k | -\phi \times x_k) \\ &= \left( \sum m_k v_k \times x_k | \phi \right) = 0. \end{aligned} \quad (5.2)$$

Thus every tangent space to  $\dot{Q}$  is decomposed into an orthogonal direct sum of the rotational and vibrational subspaces. This is indeed the alternative definition of the connection that Guichardet<sup>5</sup> first gave.

Let  $P_x$  denote the projection such that  $P_x(v)$ ,  $v \in T_x(\dot{Q})$ , is the rotational component of  $v$ . Then we have, for  $v = \sum (v_k | \partial / \partial x_k)$ ,

$$P_x(v)_j = \left( A_x^{-1} \sum m_k x_k \times v_k \right) \times x_j, \quad j = 1, 2, 3. \quad (5.3)$$

Because the complement  $v - P_x(v)$  is vibrational;

$$\sum m_j x_j \times \left( v_j - \left( A_x^{-1} \sum m_k x_k \times v_k \right) \times x_j \right) = 0.$$

The rotational vector fields  $J_j$  and  $L_j$  have of course the rotational components only. Indeed, for

$$\begin{aligned} J_j &= \sum \left( -e_j \times x_k \left| \frac{\partial}{\partial x_k} \right. \right), \\ L_j &= \sum \left( -u_j \times x_k \left| \frac{\partial}{\partial x_k} \right. \right), \end{aligned} \quad (5.4)$$

we have, from (5.3) and (5.4),  $P_x(J_j)_k = -e_j \times x_k$  and  $P_x(L_j)_k = -u_j \times x_k$ .

It is of great importance to recognize that the rotational vector field is dual to the connection form. Ineed, from (3.5), (3.6), (3.24), and (5.4), it follows that

$$\omega^k(J_j) = \delta^k_j, \quad \sigma^k(L_j) = \delta^k_j. \quad (5.5)$$

Each of these equations is a specialization of Eq. (3.11). These are also verified by (3.25), (3.26), (4.7), and (4.10).

We now deal with vibrational vector fields. For a vector field  $X$  on the internal space  $M$ , its vibrational (or horizontal) lift<sup>4,5</sup> is defined as the vibrational vector field on  $\dot{Q}$  which projects to  $X$ ;  $\pi_* X_x^* = X_{\pi(x)}$ , where  $\pi_*$  is the tangent map

of  $\pi: \dot{Q} \rightarrow M$ . Alternatively, the  $X^*$  is determined by  $\omega(X^*) = 0$  and  $\pi_* X^* = X$ . We are going to obtain vibrational lifts in Dragt's coordinates. For  $\partial / \partial \xi^j$ ,  $(\xi^j) = (\rho, \chi, \psi)$ , their vibrational lifts  $(\partial / \partial \xi^j)^*$  are found, by using (5.5), to be

$$\begin{aligned} \left( \frac{\partial}{\partial \rho} \right)^* &= \frac{\partial}{\partial \rho}, \quad \left( \frac{\partial}{\partial \chi} \right)^* = \frac{\partial}{\partial \chi}, \\ \left( \frac{\partial}{\partial \psi} \right)^* &= \frac{\partial}{\partial \psi} - \frac{1}{2} \sin \chi L_3. \end{aligned} \quad (5.6)$$

In summary, the  $J_j$  (or  $L_j$ ) and  $(\partial / \partial \xi^j)^*$  constitute a basis of the space of vector fields on  $\dot{Q}$ .

### B. Rotational and vibrational covectors

In Sec. V A we have seen that the connection is an assignment of the vibrational subspace of the tangent space at every point of  $\dot{Q}$ . In this section we deal again with the connection in the cotangent space. Let  $T_x^*(\dot{Q})$  denote the cotangent space to  $\dot{Q}$  at  $x$ . Then the inner product  $K_x$  provides an isomorphism  $K_x^b$  of  $T_x(\dot{Q})$  to  $T_x^*(\dot{Q})$ ;

$$K_x^b(u) \cdot v := K_x(u, v), \quad u, v \in T_x(\dot{Q}). \quad (5.7)$$

Let  $K_x^b(u) = p = (p_k)$ . Then from (5.7) one has  $p_k = m_k u_k$ . Further, every cotangent space is equipped with the inner product  $K_x^*$ ; for  $p, q \in T_x^*(\dot{Q})$ ,  $K_x^*$  is defined through  $K_x^* = (K_x^b)^{-1}$  as

$$K_x^*(p, q) := K_x(K_x^b(p), K_x^b(q)) = \sum \frac{(p_k | q_k)}{m_k}. \quad (5.8)$$

An assignment of the vibrational subspace of the cotangent space is made as follows: With the infinitesimal generator  $R_\phi(x)$  of the SO(3) action, we can associate the rotational covector  $K_x^b(R_\phi(x)) = (m_k R_\phi(x_k))$ , a triple. Then, like vibrational vectors in the tangent space, vibrational covectors are defined as those which are orthogonal to rotational covectors with respect to  $K_x^*$ . For a covector  $p = (p_k)$ , the orthogonality condition takes the form

$$\sum x_k \times p_k = 0, \quad (5.9)$$

which is similar to Eq. (3.10), and interpreted as a generalization of the Eckart condition.

Now, in a dual manner to (5.3), we can obtain the rotational component  $P_x^*(p)$  of a covector  $p$  as follows:

$$P_x^*(p)_j = \left( A_x^{-1} \sum x_k \times p_k \right) \times m_j x_j. \quad (5.10)$$

For the proof, we have only to show that  $p - P_x^*(p)$  is vibrational;

$$\sum x_k \times \left( p_k - \left( A_x^{-1} \sum x_j \times p_j \right) \times m_k x_k \right) = 0.$$

This is, however, easy to verify.

The components of the connection form,  $\omega^j$  or  $\sigma^j$ , have indeed the rotational components only. To see this, we note that  $\omega^j$  and  $\sigma^j$  are expressed as

$$\omega^j = \left( -e_j | A_x^{-1} \sum m_k x_k \times dx_k \right), \quad j = 1, 2, 3, \quad (5.11)$$

$$\sigma^j = \left( -u_j | A_x^{-1} \sum m_k x_k \times dx_k \right), \quad j = 1, 2, 3. \quad (5.12)$$

These are easy consequences of (3.6) and (3.24). Because of these expressions of  $\omega^j$  and  $\sigma^j$ , we have dealt with the connection form  $\omega$  in the vector-valued form  $\omega = A_x^{-1} \Sigma m_k x_k \times dx_k$  in Ref. 3. Equations (5.11) and (5.12) are put into the form

$$\omega^j = \sum (-A_x^{-1}(e_j) \times m_k x_k | dx_k), \quad (5.13)$$

$$\sigma^j = \sum (-A_x^{-1}(u_j) \times m_k x_k | dx_k), \quad (5.14)$$

respectively. Thinking of  $\omega^j$  and  $\sigma^j$  as covectors and applying (5.10), we have  $P_x^*(\omega^j)_k = -A_x^{-1}(e_j) \times m_k x_k$  and  $P_x^*(\sigma^j) = -A_x^{-1}(u_j) \times m_k x_k$ .

We turn to vibrational one-forms. On account of (5.9), a one-form  $\nu = \Sigma(p_k | dx_k)$  is vibrational if and only if it vanishes for any rotational vector  $R_\phi(x)$ ;

$$\nu(R_\phi(x)) = \sum(p_k | -\phi \times x_k) = -\left(\sum x_k \times p_k | \phi\right) = 0.$$

For  $d\xi^j$ ,  $(\xi^j) = (\rho, \chi, \psi)$ , the pullbacks  $\pi^* d\xi^j$  are vibrational, as is easily seen. In summary, we have found the basis  $\omega^j$  (or  $\sigma^j$ ) and  $d\xi^j$ , viewed as one-forms on  $\dot{Q}$ , in the space of one-forms on  $\dot{Q}$ . This basis is dual to the basis  $J_j$  (or  $L_j$ ) and  $(\partial/\partial\xi^j)^*$  obtained in the last section.

**Proposition 10:** The one-forms  $\omega^j$  (or  $\sigma^j$ ) and vector fields  $J_j$  (or  $L_j$ ) and  $(\partial/\partial\xi^j)^*$  constitute dual bases on  $\dot{Q}$ , where  $\omega^j$  and  $\sigma^j$  are the components of the connection form  $\omega$  [see (3.24)], and  $J_j$  and  $L_j$  the components of the infinitesimal rotational vector  $J$  [see (4.13)], respectively, and  $(\partial/\partial\xi^j)^*$  are the horizontal lifts of  $\partial/\partial\xi^j$ ,  $(\xi^j) = (\rho, \chi, \psi)$ .

**Remark:** In a local sense, this proposition holds true for any internal coordinates.

## VI. THE METRIC AND THE VOLUME ELEMENT ON THE INTERNAL SPACE

### A. The metric

We have defined the inner product  $K$  on  $Q$ , which induces the metric tensor  $K_x$  given in (5.1). In a symbolic way, we can write  $K_x$  as

$$\begin{aligned} K_x &= \sum m_k (dx_k | dx_k) \\ &= (dr|dr) + (ds|ds). \end{aligned} \quad (6.1)$$

Here the last equality holds because of the fact the  $(r^k)$  and  $(s^k)$  are the components of  $x \in Q$  with respect to the orthonormal system  $\{f_k, f_{3+k}\}$ ,  $k = 1, 2, 3$ . Our purpose in this section is to express  $K_x$  in terms of  $\omega^j$  (or  $\sigma^j$ ) and  $d\xi^j$ , the basis obtained in the last section.

Using (5.4) and (3.7), we obtain, for  $J_j$  and  $L_j$ ,

$$\begin{aligned} K_x(J_j, J_l) &= \sum m_k (e_j \times x_k | e_l \times x_k) \\ &= (e_j | A_x(e_l)), \end{aligned} \quad (6.2)$$

and similarly

$$K_x(L_j, L_l) = (u_j | A_x(u_l)). \quad (6.3)$$

We turn to  $K_x((\partial/\partial\xi^i)^*, (\partial/\partial\xi^j)^*)$ ,  $(\xi^j) = (\rho, \chi, \psi)$ . A

metric tensor  $(b_{ij})$  is defined on the internal space  $M$  by

$$b_{ij} = B_{\pi(x)}\left(\frac{\partial}{\partial\xi^i}, \frac{\partial}{\partial\xi^j}\right) = K_x\left(\left(\frac{\partial}{\partial\xi^i}\right)^*, \left(\frac{\partial}{\partial\xi^j}\right)^*\right). \quad (6.4)$$

The  $b_{ij}$  are well defined, that is, the right-hand side of (6.4) is independent of the  $x$  chosen, because every vibrational subspace at  $x$  in the same  $SO(3)$  orbit is isomorphic with the tangent space to  $M$  at  $\pi(x)$ , and because  $K_x$  is  $SO(3)$  invariant. Thus on account of the orthogonality of  $J_j$  (or  $L_j$ ) and  $(\partial/\partial\xi^j)^*$ , and of the duality (Proposition 10), we have

$$\begin{aligned} K_x &= \sum (e_i | A_x(e_j)) \omega^i \omega^j + \sum b_{ij} d\xi^i d\xi^j \\ &= \sum (u_i | A_x(u_j)) \sigma^i \sigma^j + \sum b_{ij} d\xi^i d\xi^j. \end{aligned} \quad (6.5)$$

In order to get  $(b_{ij})$  in an explicit form, we refer to Eq. (6.1). Using (2.16), (3.18), and (3.25), we can work out  $K_x$  to get

$$\begin{aligned} (dr|dr) + (ds|ds) &= d\rho^2 + \frac{1}{4}\rho^2(d\chi^2 + \cos^2\chi d\psi^2) + \rho^2 \sin^2(\chi/2)(\sigma^1)^2 \\ &\quad + \rho^2 \cos^2(\chi/2)(\sigma^2)^2 + \rho^2(\sigma^3)^2. \end{aligned} \quad (6.6)$$

From (6.5) and (6.6) together with (3.16), the induced metric on the internal space turns out to be

$$B_{\pi(x)} = d\rho^2 + \frac{1}{4}\rho^2(d\chi^2 + \cos^2\chi d\psi^2). \quad (6.7)$$

If we set  $\rho^2 = \tau$  and  $\chi = \pi/2 - \theta$ , Eq. (6.7) goes over into

$$B_{\pi(x)} = (1/4\tau)(d\tau^2 + \tau^2(d\theta^2 + \sin^2\theta d\psi^2)). \quad (6.8)$$

This shows that  $B_{\pi(x)}$  is a conformally flat metric, because  $4\tau B_{\pi(x)}$  is the standard flat metric expressed in the spherical coordinates  $(\tau, \theta, \psi)$  in  $\mathbb{R}^3_+$ .

**Theorem 11:** The internal space is endowed with the conformally flat metric which is expressed as Eq. (6.7) in Dragt's coordinates (2.16).

For the sake of use in Sec. VII, we discuss also  $K_x^*$  in a dual manner to  $K_x$ . From the definition (5.8) of  $K_x^*$  together with (6.2), (6.3), and (6.4), we obtain

$$K_x^*(\omega^i, \omega^j) = (e_i | A_x^{-1}(e_j)), \quad (6.9)$$

$$K_x^*(\sigma^i, \sigma^j) = (u_i | A_x^{-1}(u_j)), \quad (6.10)$$

$$K_x^*(d\xi^i, d\xi^j) = b^{ij}, \quad (6.11)$$

where  $(b^{ij})$  is the inverse of  $(b_{ij})$ . Equation (6.11) can serve as a generalized definition of Wilson's  $G$  matrix.<sup>15</sup>

Using (6.9)–(6.11), we obtain, for a real-valued function  $f$  on  $Q$ ,

$$\begin{aligned} K_x^*(df, df) &= \sum \frac{1}{m_k} \left( \frac{\partial f}{\partial x_k} \middle| \frac{\partial f}{\partial x_k} \right) \\ &= \sum (e_i | A_x^{-1}(e_j)) J_i f J_j f + \sum b^{ij} \left( \frac{\partial f}{\partial \xi^i} \right)^* f \left( \frac{\partial f}{\partial \xi^j} \right)^* f \\ &= \sum (u_i | A_x^{-1}(u_j)) L_i f L_j f + \sum b^{ij} \left( \frac{\partial f}{\partial \xi^i} \right)^* f \left( \frac{\partial f}{\partial \xi^j} \right)^* f. \end{aligned} \quad (6.12)$$

This equation is dual to Eq. (6.5), and the coefficients are



ready to be known from (6.6) in Dragt's coordinates;

$$K_x^*(df, df) = \frac{1}{\rho^2 \sin^2(\chi/2)} (L_1 f)^2 + \frac{1}{\rho^2 \cos^2(\chi/2)} (L_2 f)^2 + \frac{1}{\rho^2} (L_3 f)^2 + \left(\frac{\partial f}{\partial \rho}\right)^2 + \frac{4}{\rho^2} \left(\frac{\partial f}{\partial \chi}\right)^2 + \frac{4}{\rho^2 \cos^2 \chi} \left(\left(\frac{\partial}{\partial \psi}\right)^* f\right)^2. \quad (6.13)$$

## B. The volume element

In order to perform integration on the internal space, we have to fix the volume element which is to be deduced from that on the configuration space  $Q_0$ . The volume element on  $Q_0$  is of course  $dQ_0 = dy_1 \wedge dy_2 \wedge dy_3$ , where  $dy_k$  stands for  $dy_k^1 \wedge dy_k^2 \wedge dy_k^3$ ,  $k = 1, 2, 3$ . The volume element  $dV_0$  defined by the inner product  $K$  is related to  $dQ_0$  by

$$dV_0 = (m_1 m_2 m_3)^{3/2} dQ_0. \quad (6.14)$$

We recall that  $(B^k/N_0, q^j)$ ,  $k = 1, 2, 3$ ,  $j = 1, \dots, 6$ , serve as the Cartesian coordinates in  $Q_0$ , so that we have

$$dV_0 = N_0^{-3} dB^1 \wedge \dots \wedge dB^3 \wedge dq^1 \wedge \dots \wedge dq^6. \quad (6.15)$$

Separating off the center-of-mass coordinates  $(B^k)$  from  $dQ_0$ , we obtain, from (6.14) and (6.15), the volume element on  $Q$

$$dQ = \mu dq^1 \wedge \dots \wedge dq^6, \\ = \mu dV, \quad \mu = \left(\frac{\sum m_k}{\prod m_k}\right)^{3/2}, \quad (6.16)$$

where  $dV = dq^1 \wedge \dots \wedge dq^6$  equals the volume element defined by the metric  $K_x$  on  $Q$ .

Expressing  $dV$  in Dragt's coordinates is straightforward. From (6.5) the volume element  $dV$  takes the form

$$dV = (\det A_x \det(b_{ij}))^{1/2} \omega^1 \wedge \omega^2 \wedge \omega^3 \wedge d\xi^1 \wedge d\xi^2 \wedge d\xi^3 \\ = (\det A_{g^{-1}x} \det(b_{ij}))^{1/2} \sigma^1 \wedge \sigma^2 \wedge \sigma^3 \wedge d\xi^1 \wedge d\xi^2 \wedge d\xi^3. \quad (6.17)$$

Here we have used the equality  $(u_i | A_x (u_j)) = (e_i | A_{g^{-1}x} (e_j))$ . We note that  $\det A_x = \det A_{g^{-1}x}$  because of Eq. (3.8). Put another way,  $\det A_x$  depends only on  $\pi(x)$ . Further, inserting (3.25) and (3.26) into (6.17), we obtain

$$dQ = \mu dV = dG \wedge dM, \quad (6.18)$$

where

$$dG = -\Theta^1 \wedge \Theta^2 \wedge \Theta^3 = -\Psi^1 \wedge \Psi^2 \wedge \Psi^3, \quad (6.19)$$

$$dM = \mu (\det A_x \det(b_{ij}))^{1/2} d\xi^1 \wedge d\xi^2 \wedge d\xi^3. \quad (6.20)$$

We remark that  $dM$  is  $\mu (\det(A_x))^{1/2}$  times the volume element defined by the metric  $(b_{ij})$  on  $M$ . Expressing  $dG$  and  $dM$  in Dragt's coordinates is an easy matter. From (3.19) and (3.20), and from (3.16) and (6.7), it follows that

$$dG = \sin \beta d\alpha \wedge d\beta \wedge d\gamma, \quad (6.21)$$

$$dM = (\mu/16) \rho^5 \sin 2\chi d\rho \wedge d\chi \wedge d\psi, \quad (6.22)$$

respectively.

*Proposition 12:* The volume element  $dM$  on the internal space  $M$  is given by (6.20) or (6.22), which is  $\mu (\det A_x)^{1/2}$ ,  $\mu = (\sum m_k / \prod m_k)^{3/2}$ , times the volume element defined by the induced metric  $B$  on  $M$ .

## VII. ASSOCIATED COMPLEX VECTOR BUNDLES

### A. The bundles $V_l$

Our goal is to set up quantum mechanics for internal molecular motions. To accomplish the task, we have to answer the question of how the Euler angles should be eliminated from the wave function  $f(x)$  on  $Q$  to give wave functions on  $M$ . To do so, we invoke the vector bundle theory. Recall that the center-of-mass system is made into the principal  $SO(3)$  bundle  $\pi: Q \rightarrow M$ . With this bundle, complex vector bundles are associated as follows: Fix a non-negative integer  $l$ . Let  $D^l$  denote the  $l$ th unitary irreducible representation of  $SO(3)$  and  $\mathbb{C}^{2l+1}$  its representation space. By  $D_{jk}^l(g)$ ,  $g \in SO(3)$ , we mean the matrix elements; for  $z = (z_j) \in \mathbb{C}^{2l+1}$  (Ref. 16), one has

$$(D^l(g)z)_k = \sum_{j=-l}^{-l} D_{kj}^l(g) z_j. \quad (7.1)$$

For a basis  $|lm\rangle$  with  $\hat{J}_3 |lm\rangle = m |lm\rangle$ , this equation is often written as

$$D^l(g) |lm\rangle = \sum |lm'\rangle D_{m'm}^l(g).$$

It is also well known<sup>14</sup> that the matrix  $D^l(g)$  satisfies

$$\hat{J}^2 D^l(g) = l(l+1) D^l(g), \quad \hat{J}^2 = \sum \hat{J}_k^2 = \sum \hat{L}_k^2, \quad (7.2)$$

$$\hat{J}_k D^l(g) = -[\hat{J}_k] D^l(g), \quad k = 1, 2, 3, \quad (7.3)$$

$$\hat{L}_k D^l(g) = D^l(g) [\hat{L}_k], \quad k = 1, 2, 3, \quad (7.4)$$

where  $[\hat{J}_k] = [\hat{L}_k]$  denote the representation matrices of  $\hat{J}_k$  and  $\hat{L}_k$ , respectively. Especially, one has

$$[\hat{J}_3] = [\hat{L}_3] = \text{diag}(l, l-1, \dots, -l). \quad (7.5)$$

We notice that Eqs. (7.3) and (7.4) result from the very definition of  $J_k$  and  $L_k$ . Indeed,  $J_k$  and  $L_k$  are infinitesimal generators of the left and right actions of  $SO(3)$ , respectively.

Define a left action of  $SO(3)$  on the product space  $\dot{Q} \times \mathbb{C}^{2l+1}$  by

$$(x, z) \rightarrow (gx, D^l(g)z). \quad (7.6)$$

This action gives an equivalence relation in  $\dot{Q} \times \mathbb{C}^{2l+1}$ . The quotient manifold, denoted by  $\dot{Q} \times_{SO(3)} \mathbb{C}^{2l+1}$  is made into a complex vector bundle  $V_l = (\dot{Q} \times_{SO(3)} \mathbb{C}^{2l+1}, \pi_l, M)$  via the commutative diagram

$$\begin{array}{ccc} \dot{Q} \times \mathbb{C}^{2l+1} & \xrightarrow{q_l} & \dot{Q} \times_{SO(3)} \mathbb{C}^{2l+1} \\ \text{pr} \downarrow & & \downarrow \pi_l \\ \dot{Q} & \xrightarrow{\pi} & M \end{array}, \quad (7.7)$$

where pr denotes the projection onto the first factor, and  $q_l$  is the natural projection. The  $\pi_l$  is the projection in  $V_l$  such that  $\pi_l \circ q_l = \pi \circ \text{pr}$ . A map  $\sigma: M \rightarrow \dot{Q} \times_{SO(3)} \mathbb{C}^{2l+1}$  such that  $\pi \circ \sigma = \text{id}$  is called a cross section in  $V_l$ . The internal states of

the molecule are then described as the cross sections in the complex vector bundle  $V_l$ . The "quantum" number  $l$  will prove to assign the total angular momentum of the molecule;  $\hat{J}^2 = l(l+1)$ .

A  $\mathbb{C}^{2l+1}$ -valued function  $F$  on  $\hat{Q}$  is said to be  $D^l$ -equivariant if it satisfies

$$F(gx) = D^l(g)F(x). \quad (7.8)$$

To any  $D^l$ -equivariant function, there corresponds a cross section in the complex vector bundle  $V_l$ , and vice versa.<sup>4</sup> We denote by  $q_l^\#$  the one-to-one correspondence from the cross section to the  $D^l$ -equivariant function.

In our case, the correspondence  $q_l^\#$  is quite simple on account of Theorem 2. Let  $\sigma_0$  denote the cross section given in Sec. II B. Then any point  $x$  of  $\hat{Q}$  is of the form  $g\sigma_0(w)$ ,  $w \in M$ . Let  $F$  be a  $D^l$ -equivariant function on  $\hat{Q}$ . Then from (7.8) it follows that

$$F(x) = D^l(g)F(\sigma_0(w)). \quad (7.9)$$

Thus, setting  $\Phi = F \circ \sigma_0$ , one can identify  $\Phi$  with a cross section and has  $q_l^\# \Phi = F$ . This means that any cross section in  $V_l$  becomes a  $\mathbb{C}^{2l+1}$ -valued function on  $M$  because of the triviality of the  $SO(3)$  bundle  $\pi: \hat{Q} \rightarrow M$ . We note also that Eq. (7.9) is also expressed as

$$F_k(x) = \sum_{j=l}^{-l} D^l_{kj}(g) \Phi_j(w)$$

in components. This form of wave functions are frequently used in the three-body problem.<sup>10,12</sup> If  $\Phi_j(w) = \text{const}$ , the  $F_k$ 's are those used by van Winter for the asymmetric rotator.<sup>17</sup> In what follows, we treat  $D^l$ -equivariant functions in the form  $D^l(g)\Phi(w)$ . Now, the quantum number  $l$  is easy to understand. Indeed, from (7.2) and (7.3) the  $D^l$ -equivariant function  $D^l(g)\Phi(w)$  is a set of simultaneous eigenstates of  $\hat{J}^2$  and  $\hat{J}_3$ . The cross section  $\Phi$  in  $V_l$  is therefore thought of as a set of the internal states of the molecule with a specified eigenvalue  $\hat{J}^2 = l(l+1)$ . This understanding goes well irrespective of whether the complex vector bundle  $V_l$  is trivial or not.

**Theorem 13:** The internal states of the triatomic molecule with an eigenvalue  $l(l+1)$  of the square of the total angular momentum operator are described as the cross sections in the complex vector bundle  $V_l$  given in (7.7). The  $V_l$  is, however, a trivial bundle, and hence the cross sections become  $\mathbb{C}^{2l+1}$ -valued functions on the internal space  $M$ .

The inner product for cross sections in  $V_l$  is defined as follows: Let  $\Phi$  and  $\Psi$  be cross sections in  $V_l$ . Then the Hermitian metric in  $V_l$  is defined for  $\Phi$  and  $\Psi$  by

$$\langle \Phi | \Psi \rangle (w) := ((q_l^\# \Phi)(x) | (q_l^\# \Psi)(x)), \quad (7.10)$$

where the round brackets in the right-hand side indicate the Hermitian inner product in  $\mathbb{C}^{2l+1}$ ;

$$(u|v) = \sum_{k=l}^{-l} \bar{u}_k v_k, \quad u, v \in \mathbb{C}^{2l+1}.$$

We note that the right-hand side of (7.10) depends only on  $\pi(x) = w$ . The inner product of  $\Phi$  and  $\Psi$  is then defined as

$$\langle \Phi | \Psi \rangle_M := \int_M (\Phi | \Psi) dM. \quad (7.11)$$

The  $L^2$  space of cross sections with respect to this inner product is considered the state space of the internal molecular motions. From (6.18) and (7.10), the inner product is expressed as

$$\begin{aligned} & \int_{\hat{Q}} ((q_l^\# \Phi)(x) | (q_l^\# \Psi)(x)) dQ \\ &= \int_{SO(3)} dG \int_M (\Phi | \Psi) dM \\ &= 8\pi^2 \langle \Phi | \Psi \rangle_M. \end{aligned} \quad (7.12)$$

## B. The linear connection on $V_l$ .

The vector bundle  $V_l$  is equipped with the linear connection associated with the connection on  $\hat{Q}$  (Ref. 4). Let  $X$  be a vector field on  $M$  and  $X^*$  its horizontal lift. Then for a cross section  $\sigma$  in  $V_l$ , its covariant derivative with respect to  $X$  is defined by

$$\nabla_X \sigma = q_l^{\#-1} X^*(q_l^\# \sigma). \quad (7.13)$$

The operator  $\nabla$  is called the associated linear connection, which is linear in  $X$  and  $\sigma$ , and satisfies for arbitrary functions  $f$  on  $M$  the conditions

$$\nabla_{fX} \sigma = f \nabla_X \sigma, \quad (7.14)$$

$$\nabla_X f \sigma = (Xf) \sigma + f \nabla_X \sigma. \quad (7.15)$$

The curvature of  $\nabla$  is defined by

$$R(X, Y) \sigma = [\nabla_X, \nabla_Y] \sigma - \nabla_{[X, Y]} \sigma. \quad (7.16)$$

From (7.13), the curvature is also written as

$$R(X, Y) \sigma = q_l^{\#-1} ([X^*, Y^*] - [X, Y]^*) q_l^\# \sigma. \quad (7.17)$$

It is easy to express the associated linear connection and its curvature in Dragt's coordinates. From (5.6), (7.4), and (7.13), it follows for a cross section  $\sigma$  with  $q_l^\# \sigma = D^l(g)\Phi(w)$  that

$$\begin{aligned} \nabla_{\partial/\partial\rho} \sigma &= \frac{\partial}{\partial\rho} \otimes I, & \nabla_{\partial/\partial\chi} \sigma &= \frac{\partial}{\partial\chi} \otimes I, \\ \nabla_{\partial/\partial\psi} \sigma &= \frac{\partial}{\partial\psi} \otimes I - \frac{i}{2} \sin \chi [\hat{L}_3], \end{aligned} \quad (7.18)$$

where  $I$  is the  $(2l+1) \times (2l+1)$  identity matrix. We note here that the connection form  $\sigma_0^* \omega$  given in (3.27) is represented in the linear connection  $\nabla$  so as to couple with  $\partial/\partial\xi^j$ ,  $(\xi^j) = (\rho, \chi, \psi)$ . The curvature is computed by using (7.17) to give

$$R\left(\frac{\partial}{\partial\psi}, \frac{\partial}{\partial\chi}\right) = \frac{i}{2} \cos \chi [\hat{L}_3], \quad (7.19)$$

and the other components vanishing. From (3.32) this curvature  $R$  proves to be a representation of  $\sigma_0^* \Omega$ .

**Theorem 14:** The linear connection and its curvature on the complex vector bundle  $V_l$  are expressed as (7.18) and (7.19), respectively, in Dragt's coordinates.

## VIII. QUANTUM MECHANICS FOR INTERNAL STATES

We are in the final stage to set up quantum mechanics for internal states of the triatomic molecule. What we have to do is to obtain the internal Hamiltonian operator acting on cross sections in  $V_l$ .

## A. The Laplacian on $Q$

We start with the kinetic energy of the molecule in  $Q_0$ , which is expressed as

$$\frac{1}{2} \int_{Q_0} \sum \frac{1}{m_k} \left( \frac{\partial f}{\partial x_k} \middle| \frac{\partial f}{\partial x_k} \right) dQ_0, \quad (8.1)$$

where  $\partial/\partial x_k$ ,  $k=1,2,3$ , are gradient vectors. When integrated by parts, this functional yields the Laplacian  $\Delta_0$  with respect to the inner product  $K$ ;

$$\Delta_0 = \sum \frac{1}{m_k} \left( \frac{\partial}{\partial x_k} \right)^2. \quad (8.2)$$

Owing to the fact that  $(B^k/N_0, q^j)$  serve as the Cartesian coordinates in  $Q_0$ , the Laplacian  $\Delta_0$  are put into the form

$$\Delta_0 = \sum_{i=1}^6 \left( \frac{\partial}{\partial q^i} \right)^2 + N_0^2 \sum_{k=1}^3 \left( \frac{\partial}{\partial B^k} \right)^2. \quad (8.3)$$

Separating off the center-of-mass coordinates  $(B^k)$ ,  $k=1,2,3$ , or assuming that the wave functions under consideration are independent of  $(B^k)$ , we obtain the Laplacian  $\Delta = \Sigma(\partial/\partial q^i)^2$  on  $Q$ . This operator can be also derived from the functional

$$\int_Q \sum \left| \frac{\partial f}{\partial q^i} \right|^2 dQ = \int_Q K_x^* (\overline{df}, df) dQ. \quad (8.4)$$

From (6.12), this expression turns into

$$\begin{aligned} & \int_Q \sum \left( e_i \middle| A_x^{-1}(e_j) \right) \overline{J_i f} J_j f dQ \\ & + \int_Q \sum b^{ij} \overline{\left( \frac{\partial}{\partial \xi^i} \right)^* f} \left( \frac{\partial}{\partial \xi^j} \right)^* f dQ \\ & = \int_Q \sum \left( u_i \middle| A_x^{-1}(u_j) \right) \overline{L_i f} L_j f dQ \\ & + \int_Q \sum b^{ij} \overline{\left( \frac{\partial}{\partial \xi^i} \right)^* f} \left( \frac{\partial}{\partial \xi^j} \right)^* f dQ. \end{aligned} \quad (8.5)$$

In both sides, the first term is twice the rotational energy, and the last twice the vibrational energy.

We are ready to express  $\Delta$  in Dragt's coordinates, using (6.13), (6.18), and (6.22). After integration by parts we can obtain the following.

**Theorem 15:** The Laplacian  $\Delta$  on the center-of-mass system  $Q$  with respect to the metric  $K_x$ , given in (6.6), takes the form

$$\begin{aligned} \Delta = & \frac{\partial}{\partial \rho^2} + \frac{5}{\rho} \frac{\partial}{\partial \rho} + \frac{4}{\rho^2} \left[ \frac{\partial^2}{\partial \chi^2} + 2 \cot 2\chi \frac{\partial}{\partial \chi} + \frac{1}{\cos^2 \chi} \right. \\ & \times \left. \left( \frac{\partial}{\partial \psi} - \frac{i}{2} \sin \chi \hat{L}_3 \right)^2 \right] - \left( \frac{1}{\rho^2 \sin^2(\chi/2)} (\hat{L}_1)^2 \right. \\ & \left. + \frac{1}{\rho^2 \cos^2(\chi/2)} (\hat{L}_2)^2 + \frac{1}{\rho^2} (\hat{L}_3)^2 \right), \end{aligned} \quad (8.6)$$

where  $L_k = i\hat{L}_k$ ,  $k=1,2,3$ , are given in (4.9).

*Remark:* The last term including  $(\hat{L}_k)^2$  is derived from the rotational energy and the rest from the vibrational energy. The vibrational part of  $\Delta$  includes the differential operators coming from the horizontal lifts only [see (5.6)]. If the vibrational part is separated off and  $(\rho, \chi, \psi)$  are fixed, the operator  $-\Delta/2$  reduces to the well-known Hamiltonian for

the rigid rotator. The Laplacian  $\Delta$  was given in Ref. 12 in the expanded form.

## B. The internal Hamiltonian operator

We are now in a position to obtain the internal Hamiltonian operator  $H_I$  acting on cross sections in  $V_I$ . Let  $H$  denote the Hamiltonian operator on  $Q$ ;

$$H = -\frac{1}{2}\Delta + U, \quad (8.7)$$

where  $U$  is the potential function invariant under the  $SO(3)$  action. Let  $\Phi$  be a cross section in  $V_I$  and  $F$  the corresponding  $D^1$ -equivariant function on  $Q$ ;  $q_i^\# \Phi = F$ . Then, on account of Eq. (7.12), the internal Hamiltonian operator  $H_I$  is defined through

$$\int_Q (F | HF) dQ = 8\pi^2 \langle \Phi | H_I \Phi \rangle_M. \quad (8.8)$$

We note that the left-hand side of (8.8) is thought of as a sum of the kinetic energy for a wave function  $F_k$ . Written out for  $F(x) = D^1(g)\Phi(w)$ , and integrated on  $SO(3)$ , the left-hand side will yield the operator  $H_I$ . For the sake of convenience we treat the left-hand side in a form similar to (8.5). Then we have

$$\begin{aligned} & \frac{1}{2} \int_Q \sum b^{ij} \left( \left( \frac{\partial}{\partial \xi^i} \right)^* F \middle| \left( \frac{\partial}{\partial \xi^j} \right)^* F \right) dQ \\ & + \frac{1}{2} \int_Q \sum A_{ij}^{-1} (L_i F | L_j F) dQ \\ & = \frac{1}{2} \int_Q \sum b^{ij} \left( \left( \frac{\partial}{\partial \xi^i} \right)^* F \middle| \left( \frac{\partial}{\partial \xi^j} \right)^* F \right) dQ \\ & + \frac{1}{2} \int_Q \sum a_{ij}^{-1} (J_i F | J_j F) dQ, \end{aligned} \quad (8.9)$$

where the round brackets in each integrand indicate the inner product in  $\mathbb{C}^{2l+1}$ , and

$$A_{ij}^{-1} = (u_i | A_x^{-1}(u_j)), \quad a_{ij}^{-1} = (e_i | A_x^{-1}(e_j)).$$

Using the definition of the covariant derivative and the fact that  $D^1(g)$  and  $[\hat{L}_k] = [\hat{J}_k]$  are unitary and Hermitian matrices, respectively, together with (7.3) and (7.4), we obtain, from (8.9),

$$\begin{aligned} & \frac{1}{2} 8\pi^2 \int_M \sum b^{ij} (\nabla_i \Phi | \nabla_j \Phi) dM + \frac{1}{2} 8\pi^2 \\ & \times \int_M \sum A_{ij}^{-1} (\Phi | [\hat{L}_i] [\hat{L}_j] \Phi) dM \\ & = \frac{1}{2} 8\pi^2 \int_M \sum b^{ij} (\nabla_i \Phi | \nabla_j \Phi) dM \\ & + \frac{1}{2} 8\pi^2 \int_M \sum a_{ij}^{-1} (\Phi | D^1(g^{-1}) \\ & \times [\hat{J}_i] [\hat{J}_j] D^1(g) \Phi) dM, \end{aligned} \quad (8.10)$$

where  $\nabla_i$ ,  $i=1,2,3$ , stand for  $\nabla_{\partial/\partial \xi^i}$ . Integration by parts gives the operator

$$\begin{aligned}
T_l &= -\frac{1}{2} \frac{1}{J_M} \sum \nabla_i (J_M b^i \nabla_j) + \frac{1}{2} \sum A_{ij}^{-1} [\hat{L}_i] [\hat{L}_j] \\
&= -\frac{1}{2} \frac{1}{J_M} \sum \nabla_i (J_M b^i \nabla_j) \\
&\quad + \frac{1}{2} \sum a_{ij}^{-1} D^i (g^{-1}) [\hat{J}_i] [\hat{J}_j] D^j (g), \quad (8.11)
\end{aligned}$$

where  $J_M = \mu(\det A_x \det(b_{ij}))^{1/2}$  is the volume density on  $M$  [see (6.20)]. The internal Hamiltonian  $H_l$  is then the sum  $T_l + U \otimes I$ , where  $I$  is  $(2l+1) \times (2l+1)$  identity matrix. Writing out (8.11) in Dragt's coordinates, we obtain

$$\begin{aligned}
H_l &= -\frac{1}{2} \left[ \left( \frac{\partial}{\partial \rho^2} + \frac{5}{\rho} \frac{\partial}{\partial \rho} \right. \right. \\
&\quad \left. \left. + \frac{4}{\rho^2} \left( \frac{\partial^2}{\partial \chi^2} + 2 \cot 2\chi \frac{\partial}{\partial \chi} \right) \right) \otimes I \right. \\
&\quad \left. + \frac{4}{\rho^2 \sin^2 \chi} \left( \frac{\partial}{\partial \psi} \otimes I - \frac{i}{2} \sin \chi [\hat{L}_3] \right)^2 \right] \\
&\quad + \frac{1}{2} \left[ \frac{1}{\rho^2 \sin^2(\chi/2)} [\hat{L}_1]^2 + \frac{1}{\rho^2 \cos^2(\chi/2)} [\hat{L}_2]^2 \right. \\
&\quad \left. + \frac{1}{\rho^2} [\hat{L}_3]^2 \right] + U \otimes I. \quad (8.12)
\end{aligned}$$

**Theorem 16:** The internal Hamiltonian operator acting on cross sections in  $V_l$  is given by  $H_l = T_l + U \otimes I$ , where  $T_l$  is defined in (8.11). In Dragt's coordinates, the  $H_l$  takes the form of (8.12).

We conclude this section with some remarks on the Schrödinger operator  $H_l$ . Zickendraht<sup>10</sup> has already derived the system of coupled equations  $H_l \Phi = E \Phi$  without reference to the internal motion. The quantum three-body problem without interaction was also analyzed in Refs. 11 and 12 by using an  $SU(3)$  action on the center-of-mass system  $Q \simeq \mathbb{R}^3 \times \mathbb{R}^3 \simeq \mathbb{C}^3$ , but the internal motion was received little attention. In Refs. 18 and 19, the same problem was studied through the hyperspherical functions on  $S^5 \subset \mathbb{R}^6 \simeq Q$ . For several-particle systems, the internal Hamiltonian operator will be obtained in the same manner as was done in this paper, though the topology of the internal space is difficult to study.

*Remark:* For the four-body problem, the internal space can be shown to be homeomorphic to  $\mathbb{R} \times (S^5 - P)$ , where  $P$  is a submanifold of  $S^5$  homeomorphic to the projective plane  $\mathbb{R}P^2$ . (For the proof, see Ref. 20, Lemma 6.3.)

Further, we notice that when the angular momentum vanishes, the internal Hamiltonian reduces to the one we have presented in the previous paper.<sup>3</sup> That is, for  $l=0$ , we have

$$H_0 = -\frac{1}{2} J_M^{-1} \sum \partial_i (J_M b^i \partial_j) + U, \quad \partial_i = \frac{\partial}{\partial \xi^i}.$$

It is worth mentioning that the first term is not  $-\frac{1}{2}$  times the Laplacian on the Riemannian space  $M$  endowed with the metric  $B$ , because  $J_M$  is  $(\sum m_k / \prod m_k)^{3/2} (\det A_x)^{1/2}$  times  $(\det(b_{ij}))^{1/2}$ , and  $\det A_x$  is not constant.

The advantage of the use of the connection theory is that in terms of connection theory one can understand how the internal motion is coupled with the rotation. The rotation of the molecule induces on the internal space the matrix-valued

gauge field or the curvature (7.19), with which the internal motion is coupled through the gauge potential or the connection (7.18). The curvature of this connection may be called the Coriolis field in the sense that it is induced by the rotation and plays the role of a matrix-valued magnetic field on the internal space. The internal motion is driven also by the matrix-valued centrifugal potential, each of the second terms in the right-hand side of (8.11). The internal Hamiltonian operator is a second order matrix-valued differential operator including both the internal motion-Coriolis field coupling and the centrifugal potential. The ambiguity in an early paper<sup>21</sup> by Hirschfelder and Wigner is thus cleared up in terms of vector bundle theory.

The last but very important point to make is the fact that for a quantum system one can describe the internal motion of a molecule (Theorems 13 and 16) but for a classical system one cannot. This is explained by Jacobi's celebrated "elimination of nodes." According to Wintner,<sup>22</sup> the dimensions of the classical Hamiltonian system are reduced by  $3+1=4$  by eliminating the angular momentum. However, in order to get a classical Hamiltonian system for the internal motion this reduction of dimensions must be  $3 \times 2 = 6$  in number. If so, the reduced Hamiltonian system would be of dimension  $3 \times 2 = 6$ , twice the dimension of the internal space.

## ACKNOWLEDGMENTS

The author would like to thank Dr. A. Tachibana at Kyoto University for helpful discussions.

This work was supported by a Grant-in-Aid for Scientific Research from the Ministry of Education in Japan.

- <sup>1</sup>J. Tennyson and B. T. Sutcliffe, *J. Chem. Phys.* **77**, 4061 (1982), and references therein.
- <sup>2</sup>B. T. Sutcliffe, *Quantum Dynamics of Molecules*, edited by R. G. Woolley (Plenum, New York, 1980).
- <sup>3</sup>A. Tachibana and T. Iwai, *Phys. Rev. A* **33**, 2262 (1986).
- <sup>4</sup>S. Kobayashi and K. Nomizu, *Foundations of Differential Geometry* (Interscience, New York, 1963), Vol. I.
- <sup>5</sup>A. Guichardet, *Ann. Inst. H. Poincaré* **40**, 329 (1984).
- <sup>6</sup>T. Iwai, "A geometric setting for a classical molecular dynamics," Technical Report #85009, Department of Applied Mathematics and Physics, Kyoto University, Kyoto, Japan.
- <sup>7</sup>N. Steenrod, *The Topology of Fiber Bundles* (Princeton U.P., Princeton, NJ, 1951).
- <sup>8</sup>T. Iwai, *J. Math. Phys.* **28**, 964 (1987).
- <sup>9</sup>F. T. Smith, *Phys. Rev.* **120**, 1058 (1960); *J. Math. Phys.* **3**, 735 (1962).
- <sup>10</sup>W. Zickendraht, *Ann. Phys. (NY)* **35**, 18 (1965).
- <sup>11</sup>A. J. Dragt, *J. Math. Phys.* **6**, 533 (1965).
- <sup>12</sup>J.-M. Levy-Leblond and M. Levy-Nahas, *J. Math. Phys.* **6**, 1571 (1965).
- <sup>13</sup>Y. Matsushima, *Differentiable Manifolds* (Marcel Dekker, New York, 1972).
- <sup>14</sup>J.-M. Normand, *A Lie Group: Rotations in Quantum Mechanics* (North-Holland, Amsterdam, 1980).
- <sup>15</sup>E. B. Wilson Jr., J. C. Decius, and P. C. Cross, *Molecular Vibrations* (McGraw-Hill, New York, 1955).
- <sup>16</sup>As long as  $D$  is concerned, the components of vectors are designated by the subscript indices.
- <sup>17</sup>C. van Winter, *Physica* **20**, 274 (1954).
- <sup>18</sup>H. Mayer, *J. Phys. A* **8**, 1562 (1982).
- <sup>19</sup>F. del Aguila, *J. Math. Phys.* **21**, 2327 (1980).
- <sup>20</sup>M. S. Narasimhan and T. R. Ramadas, *Commun. Math. Phys.* **67**, 121 (1979).
- <sup>21</sup>J. O. Hirschfelder and E. Wigner, *Proc. Natl. Acad. Sci. USA* **21**, 113 (1935).
- <sup>22</sup>A. Wintner, *The Analytical Foundations of Celestial Mechanics* (Princeton U.P., Princeton, NJ, 1941).

# Quantum mechanics of a charged scalar boson with respect to an observer's past light cone

G. H. Derrick

*School of Physics, University of Sydney, Sydney, New South Wales 2006, Australia*

(Received 7 August 1986; accepted for publication 14 January 1987)

An observer whose instantaneous "here-now" has Minkowski coordinates  $z^\lambda (\lambda = 0, 1, 2, 3)$  can only be aware of events within or on the past light cone with vertex at  $z^\lambda$ . In conventional quantum mechanics his current quantum state would refer to some spacelike surface containing  $z^\lambda$ , for example,  $x^0 = z^0$ . This is, however, a region of space-time about which the observer can know nothing except the single event  $z^\lambda$ , his current here-now. The aim of the present paper is to give a version of quantum mechanics in which the intrinsically unknowable "quantum state at the present time" is replaced by the "quantum state on the past light cone." The theory is an extension and adaptation of Dirac's point mechanics [Rev. Mod. Phys. **21**, 392 (1949)].

## I. INTRODUCTION

A previous paper<sup>1</sup> drew attention to the problems that stem from the finiteness of the velocity of light  $c$  in the conventional approaches both to classical and to quantum theory. At a certain time  $t_0$  an observer simply cannot know all the initial data of a system if that data is specified on the spacelike surface  $t = t_0$ . Reference 1 considered an adaptation of Dirac's classical Hamiltonian point mechanics<sup>2</sup> in which dynamical variables are specified by their values on an observer's past light cones, progression from data on an initial light cone to that on the current light cone being achieved by a canonical transformation. The aim of this present paper is to give a corresponding quantum treatment.

We first need to recall the definition of the light cone coordinates<sup>1</sup> belonging to an observer. Let the observer's trajectory in four dimensional Minkowski space be given in parametric form by<sup>3</sup>

$$x^\lambda = z^\lambda(\tau), \quad (1)$$

where the parameter  $\tau$  is taken as the proper interval  $\int (\eta_{\lambda\mu} dz^\lambda dz^\mu)^{1/2}$  measured along the trajectory from some arbitrary event. Thus  $\eta_{\lambda\mu} v^\lambda v^\mu = 1$ ,  $v^0 \geq 1$ , where  $v^\lambda = dz^\lambda(\tau)/d\tau$  is the observer's four-velocity vector. We now define a change of coordinates  $x^\lambda \rightarrow (\tau, y^1, y^2, y^3)$  by

$$x^\lambda = z^\lambda(\tau) + y^\lambda, \quad y^0 = -y, \quad (2)$$

where  $y = |y|$ . Thus the three-surface  $\tau = \tau_0$  is the past light cone with vertex  $z^\lambda(\tau_0)$ , while the past-pointing null vector  $y^\lambda$  serves to parametrize the cone. All compatible dynamical data on this cone can be known by the observer when his personal ideal clock reads  $\tau_0/c$ , i.e., when his "here-now" is  $z^\lambda(\tau_0)$ .

Reference 1 considered in detail the classical motion of a spin-zero particle of mass  $m$ , which was either free or suffered electromagnetic interactions. For this system the evolution of any dynamical variable  $f$  was shown to be governed by

$$\frac{df}{d\tau} = v_\lambda \{f, p^\lambda\} + \frac{\partial f}{\partial \tau}. \quad (3)$$

Here the Poisson bracket is defined with respect to a set of conjugate generalized coordinate-momentum pairs of which  $f$  and  $p^\lambda$  are functions. The second term on the right of (3) arises if  $f$  is additionally an explicit function of  $\tau$ . The four-vector  $p^\lambda$  is the energy-momentum vector of the system. In Ref. 1 it was shown that for a wide class of dynamical variables (3) is equivalent to

$$\frac{\partial f}{\partial z^\lambda} = \{f, p_\lambda\} + \left(\frac{\partial f}{\partial z^\lambda}\right)_{\text{explicit}}, \quad (4)$$

where we now have *four* independent variables  $z^\lambda$ . The second term on the right of (4) arises if  $f$  depends not only on the conjugate variables which define the Poisson brackets but also explicitly on  $z^\lambda$ . Dynamical variables which satisfy (4) depend only on the particle trajectory and where the latter cuts the past light cone with vertex  $z^\lambda$ . The route by which the observer arrived at  $z^\lambda$  is irrelevant.

In the present paper we seek Schrödinger picture quantum analogs of (3) and (4) of the form

$$i\hbar \frac{d}{d\tau} |\Psi(\tau)\rangle = v_\lambda p_{\text{OP}}^\lambda |\Psi(\tau)\rangle, \quad (5)$$

$$i\hbar \frac{\partial}{\partial z^\lambda} |\Psi(z^\lambda)\rangle = p_{\text{OP}\lambda} |\Psi(z^\lambda)\rangle. \quad (6)$$

The problem is to define a suitable Hilbert space  $\mathcal{H}_{\text{phys}}$  of physical states  $|\Psi\rangle$  and to specify the action of the four-momentum operator  $p_{\text{OP}}^\lambda$  on these states. In (5),  $\tau$  is regarded as an external parameter whose interpretation is that an observer following trajectory (1) finds that his quantum state  $|\Psi(\tau)\rangle$  evolves with his proper time according to (5). Likewise in (6),  $|\Psi(z^\lambda)\rangle$  is the state belonging to an observer whose here-now is  $z^\lambda$ , with (6) governing how the state varies with  $z^\lambda$ .

As in Ref. 1 we shall focus our attention on a system of one spin-zero boson of mass  $m$ . For this system a suitable set of classical variables is the pair  $y, \pi$ , where the coordinate  $y$  specifies the particle position on the observer's past light cone according to (2), and the conjugate variable  $\pi$  is defined as in Ref. 1. In the case of a free particle, the appropri-

ate momentum vector  $p^\lambda \equiv p^\lambda(\mathbf{y}, \boldsymbol{\pi})$ , which governs the classical motion via (3) and (4), is given by

$$\begin{aligned} p^0 &= \frac{1}{2}(\mathbf{y} \cdot \boldsymbol{\pi})^{-1}(\boldsymbol{\pi}^2 + m^2 c^2)y, \\ \mathbf{p} &= \boldsymbol{\pi} - \frac{1}{2}(\mathbf{y} \cdot \boldsymbol{\pi})^{-1}(\boldsymbol{\pi}^2 + m^2 c^2)\mathbf{y}. \end{aligned} \quad (7)$$

When an electromagnetic field derived from a potential  $A^\lambda \equiv A^\lambda(x^\kappa) \equiv A^\lambda(z^\kappa + y^\kappa)$  is present, and the particle has a charge  $e$ , (7) is changed to

$$\begin{aligned} p^0 - (e/c)A^0 &= \frac{1}{2}(\mathbf{y} \cdot \boldsymbol{\pi}_E)^{-1}(\boldsymbol{\pi}_E^2 + m^2 c^2)y, \\ \mathbf{p} - (e/c)\mathbf{A} &= \boldsymbol{\pi}_E - \frac{1}{2}(\mathbf{y} \cdot \boldsymbol{\pi}_E)^{-1}(\boldsymbol{\pi}_E^2 + m^2 c^2)\mathbf{y}, \end{aligned} \quad (8)$$

where

$$\begin{aligned} \boldsymbol{\pi}_E &= \boldsymbol{\pi} - (e/c)(\mathbf{A} + \hat{\mathbf{y}}A^0), \\ \hat{\mathbf{y}} &= \mathbf{y}/y. \end{aligned} \quad (9)$$

What we need are quantum analogs of (7) and (8), where  $\mathbf{y}$  and  $\boldsymbol{\pi}$  are replaced by operators in some Hilbert space. An obvious candidate for the latter is  $\mathcal{H}_y = L^2(\mathbb{R}^3, d^3\mathbf{y}/y)$ . This is defined as the Hilbert space of complex scalar functions  $\psi(\mathbf{y})$  for which the norm

$$(\psi, \psi)_y = \int \frac{d^3\mathbf{y}}{y} |\psi(\mathbf{y})|^2$$

exists, and with the scalar product of two elements  $\psi_1(\mathbf{y})$  and  $\psi_2(\mathbf{y})$  defined by

$$(\psi_1, \psi_2)_y = \int \frac{d^3\mathbf{y}}{y} \psi_1^*(\mathbf{y})\psi_2(\mathbf{y}). \quad (10)$$

If  $\psi_1$  and  $\psi_2$  are SO(1,3) scalars then this scalar product is Lorentz invariant on account of the like invariance of  $d^3\mathbf{y}/y$ . This motivates using the measure  $d^3\mathbf{y}/y$  rather than  $d^3\mathbf{y}$ .

In the Hilbert space  $\mathcal{H}_y$  the operator analogs of the classical variables  $\mathbf{y}$  and  $\boldsymbol{\pi}$  act according to

$$\begin{aligned} \mathbf{y}_{\text{OP}} \psi(\mathbf{y}) &= \mathbf{y}\psi(\mathbf{y}), \\ \boldsymbol{\pi}_{\text{OP}} \psi(\mathbf{y}) &= -i\hbar \mathbf{y}^{1/2} \left( \frac{\partial}{\partial \mathbf{y}} \right) [y^{-1/2} \psi(\mathbf{y})]. \end{aligned} \quad (11)$$

The operators so defined are Hermitian with respect to (10) in domains which are dense subspaces of  $\mathcal{H}_y$ . From these operators our aim is to construct an appropriate four-momentum operator  $p_{\text{OP}}^\lambda$  which determines the evolution of quantum states with respect to the external parameters  $\tau$  or  $z^\lambda$  according to (5) or (6).

By analogy with nonrelativistic quantum mechanics, we shall tentatively interpret  $|\psi(\mathbf{y}, \tau)|^2 d^3\mathbf{y}/y$  as the probability that a measurement of the position of the particle on the past light cone with vertex  $z^\lambda(\tau)$  will yield a value in the range  $\mathbf{y}$  to  $\mathbf{y} + d\mathbf{y}$  (normalizing  $\psi$  to unity). Such a measurement, of course, needs the collaboration of a large number of auxiliary observers spread throughout the current past light cone of the central observer, the measurement having been prearranged.

In a relativistic quantum theory we should anticipate the appearance of antiparticle states. This means that we should not expect  $\mathcal{H}_y$  to be the same as  $\mathcal{H}_{\text{phys}}$ , the Hilbert space of physical states  $|\Psi\rangle$ . We shall in fact resolve the elements  $\psi(\mathbf{y}) \in \mathcal{H}_y$  into particle and antiparticle components, and then construct the states  $|\Psi\rangle$  from the particle amplitude and the complex conjugate of the antiparticle amplitude.

In Sec. II we review the conventional theory of a charged scalar boson based on the Klein–Gordon equation. This will be written both in the usual Minkowski coordinates and in the light cone coordinates defined by (2). The latter formulation adds some insight to Sec. III, which addresses the difficult problem of finding energy momentum operators  $p_{\text{OP}}^\lambda$  which parallel the classical expressions (7) and (8).

## II. KLEIN–GORDON THEORY

### A. The Klein–Gordon equation in Minkowski coordinates

In the conventional, first-quantized theory of a noninteracting charged scalar boson we have a complex amplitude  $\Phi \equiv \Phi(x^\lambda)$ , which obeys the Klein–Gordon equation

$$\eta^{\lambda\mu} \frac{\partial^2 \Phi}{\partial x^\lambda \partial x^\mu} + \left( \frac{mc}{\hbar} \right)^2 \Phi = 0. \quad (12)$$

By taking a Fourier transform we obtain the general solution of (12) in the form<sup>3</sup>

$$\Phi = (2\pi)^{-3/2} \int dS_k [a(\mathbf{k})e^{-ik_\lambda x^\lambda} + \{b(\mathbf{k})e^{-ik_\lambda x^\lambda}\}^*]. \quad (13)$$

The integration in (13) is over all future-pointing vectors  $k^\lambda$  lying on the mass shell, i.e.,

$$\begin{aligned} \eta_{\lambda\mu} k^\lambda k^\mu &= \kappa^2, \quad \kappa = mc/\hbar, \\ k^0 &= \epsilon_k = (\mathbf{k}^2 + \kappa^2)^{1/2}, \end{aligned} \quad (14)$$

with the standard Lorentz invariant measure<sup>4</sup>

$$dS_k = d^3\mathbf{k}/\epsilon_k. \quad (15)$$

The particle and antiparticle amplitudes for four-momentum  $\hbar k^\lambda$  are, respectively,  $a(\mathbf{k})$  and  $b(\mathbf{k})$ . A Lorentz invariant scalar product<sup>4</sup> may be introduced between any two solutions  $\Phi$  and  $\Phi'$  of (12):

$$(\Phi, \Phi')_{\text{KG}} = \frac{i}{2} \int_S dS^\mu \left[ \Phi^* \frac{\partial \Phi'}{\partial x^\mu} - \frac{\partial \Phi^*}{\partial x^\mu} \Phi' \right], \quad (16)$$

where  $S$  is any unbounded spacelike three-surface. For sufficiently localized solutions (16) is independent of the choice of  $S$ . However the norm  $(\Phi, \Phi)_{\text{KG}}$  can be negative so that the scalar product (16) is not suitable for the definition of a Hilbert space. The Hilbert space of physical states  $\mathcal{H}_{\text{phys}}$  is constructed instead from the amplitudes  $a(\mathbf{k})$  and  $b(\mathbf{k})$ . Let us write

$$\phi(\mathbf{k}) = \begin{bmatrix} a(\mathbf{k}) \\ b(\mathbf{k}) \end{bmatrix}, \quad \phi^\dagger(\mathbf{k}) = [a^*(\mathbf{k}), b^*(\mathbf{k})], \quad (17)$$

and define the Hilbert space  $\mathcal{H}_{\text{phys}} = \mathcal{H}_k$  to be the linear vector space of all  $\phi(\mathbf{k})$  for which the norm  $(\phi, \phi)$  exists, based on the scalar product

$$\begin{aligned} (\phi, \phi') &= \int dS_k \phi^\dagger(\mathbf{k})\phi'(\mathbf{k}), \\ &= \int dS_k [a^*(\mathbf{k})a'(\mathbf{k}) + b^*(\mathbf{k})b'(\mathbf{k})]. \end{aligned} \quad (18)$$

In  $\mathcal{H}_k$  the momentum four-vector and position three-vector operators are given, respectively, by

$$p_{\text{OP}}^\lambda \phi(\mathbf{k}) = \hbar k^\lambda \phi(\mathbf{k}), \quad (19)$$

$$\mathbf{x}_{\text{OP}} \phi(\mathbf{k}) = i\epsilon_k^{1/2} \frac{\partial}{\partial \mathbf{k}} [\epsilon_k^{-1/2} \phi(\mathbf{k})], \quad (20)$$

while the charge operator is

$$Q = e \begin{pmatrix} 1 & 0 \\ 0 & -1 \end{pmatrix}. \quad (21)$$

The Hermitian operators which generate infinitesimal rotations and boosts are defined, respectively, by

$$(j_k^{23}, j_k^{31}, j_k^{12}) \phi(\mathbf{k}) = -i\hbar \mathbf{k} \times \frac{\partial \phi(\mathbf{k})}{\partial \mathbf{k}}, \quad (22)$$

$$(j_k^{01}, j_k^{02}, j_k^{03}) \phi(\mathbf{k}) = -ih\epsilon_k \frac{\partial \phi(\mathbf{k})}{\partial \mathbf{k}}, \quad (23)$$

with the tensor  $j_k^{\lambda\mu}$  being interpreted as the angular momentum operator in  $\mathcal{H}_k$ . The commutation relations of  $p_{\text{OP}}^\lambda$  and  $j_k^{\lambda\mu}$  are those appropriate for the generators of the Poincaré group:

$$[p_{\text{OP}}^\kappa, p_{\text{OP}}^\lambda] = 0, \quad (24)$$

$$[j_k^{\mu\kappa}, p_{\text{OP}}^\lambda] = i\hbar [\eta^{\kappa\lambda} p_{\text{OP}}^\mu - \eta^{\lambda\mu} p_{\text{OP}}^\kappa], \quad (25)$$

$$[j_k^{\mu\kappa}, j_k^{\lambda\nu}] = i\hbar [\eta^{\mu\nu} j_k^{\kappa\lambda} + \eta^{\kappa\lambda} j_k^{\mu\nu} - \eta^{\lambda\mu} j_k^{\nu\kappa} - \eta^{\nu\kappa} j_k^{\mu\lambda}]. \quad (26)$$

When an external electromagnetic field derived from the vector potential  $A^\lambda$  is present we modify (12) by the ansatz  $\partial/\partial x^\lambda \rightarrow \partial/\partial x^\lambda + [ie/(\hbar c)] A_\lambda$ :

$$\eta^{\lambda\mu} \left( \frac{\partial}{\partial x^\lambda} + \frac{ie}{\hbar c} A_\lambda \right) \left( \frac{\partial}{\partial x^\mu} + \frac{ie}{\hbar c} A_\mu \right) \Phi + \kappa^2 \Phi = 0. \quad (27)$$

However, no longer is there in general any natural separation into particle and antiparticle amplitudes based on the sign of the frequency. Thus it is now difficult to know how to construct the Hilbert space of physical states  $\mathcal{H}_{\text{phys}}$ , and to define therein suitable energy, momentum, and angular momentum operators. A satisfactory resolution of these problems requires second quantization and the use of the interaction picture.

## B. The Klein-Gordon equation in light cone coordinates

It is instructive to rewrite the Klein-Gordon equation in terms of the light cone coordinates of an observer. The formalism so obtained is equivalent to the usual one, and does not yield the past light cone quantum theory that we are seeking. Nevertheless it does provide some useful pointers on how to attain our goal, and in particular leads to the valuable identities (49) and (50) below.

The metric tensor for the coordinates  $(\tau, y^1, y^2, y^3)$  has been given in Ref. 1, and with its aid standard tensor analysis transforms the free Klein-Gordon equation (12) to the form

$$i\Sigma \frac{\partial \Phi}{\partial \tau} = v_\lambda R^\lambda \Phi. \quad (28)$$

In (28)  $\Sigma$  is a scalar operator and  $R^\lambda$  is a vector operator, with the definitions

$$\Sigma = -i \left( \mathbf{y} \cdot \frac{\partial}{\partial \mathbf{y}} + 1 \right), \quad (29)$$

$$R^0 = \frac{1}{2} y \left( -\frac{\partial^2}{\partial y^2} + \kappa^2 \right), \quad (30)$$

$$\mathbf{R} = -\frac{\partial}{\partial \mathbf{y}} \left( \mathbf{y} \cdot \frac{\partial}{\partial \mathbf{y}} \right) - \frac{1}{2} y \left( -\frac{\partial^2}{\partial y^2} + \kappa^2 \right). \quad (31)$$

These operators are Hermitian with respect to the  $\mathcal{H}_y$  scalar product (10), and are subject to the identity

$$\eta_{\lambda\mu} R^\lambda R^\mu = \kappa^2 (\Sigma^2 - 1). \quad (32)$$

The indefinite Klein-Gordon scalar product (16) assumes a simple form in light cone coordinates. Choosing a sequence of spacelike three-surfaces which have the past light cone  $\tau = \text{const}$  as their limit, we find that for sufficiently localized  $\Phi, \Phi'$ ,

$$\begin{aligned} (\Phi, \Phi')_{\text{KG}} &= \frac{i}{2} \int \frac{d^3 \mathbf{y}}{y} \left[ \Phi' \cdot \frac{\partial \Phi^*}{\partial \mathbf{y}} - \Phi^* \cdot \frac{\partial \Phi'}{\partial \mathbf{y}} \right] \\ &= (\Phi, \Sigma \Phi')_y, \end{aligned} \quad (33)$$

where the latter is an  $\mathcal{H}_y$  scalar product [see (10)]. As we shall see later in this section,  $\Sigma$  has positive, negative, and zero eigenvalues, consistent with the indefinite character of the Klein-Gordon scalar product.

In terms of our light cone coordinates the general free particle solution (13) becomes

$$\Phi = \int dS_k \left[ a(\mathbf{k}) u_k(\mathbf{y}) e^{-ik_\lambda z^\lambda} + \{b(\mathbf{k}) u_k(\mathbf{y}) e^{-ik_\lambda z^\lambda}\}^* \right], \quad (34)$$

where

$$u_k(\mathbf{y}) = (2\pi)^{-3/2} e^{-ik_\lambda y^\lambda} = (2\pi)^{-3/2} e^{i(\epsilon_k y^0 + \mathbf{k} \cdot \mathbf{y})}. \quad (35)$$

Substituting (34) into (28) yields the eigenvalue equations

$$R^\lambda u_k(\mathbf{y}) = k^\lambda \Sigma u_k(\mathbf{y}), \quad (36)$$

$$R^\lambda u_k^*(\mathbf{y}) = -k^\lambda \Sigma u_k^*(\mathbf{y}).$$

In Appendix A we prove the following orthogonality and completeness properties:

$$(u_k, \Sigma u_{k'})_y = \epsilon_k \delta(\mathbf{k} - \mathbf{k}'), \quad (37)$$

$$(u_k^*, \Sigma u_{k'}^*)_y = -\epsilon_k \delta(\mathbf{k} - \mathbf{k}'), \quad (38)$$

$$(u_k^*, \Sigma u_{k'})_y = 0, \quad (39)$$

$$2 \text{Re} \left\{ \int dS_k \Sigma u_k(\mathbf{y}) u_k^*(\mathbf{y}') \right\} = y \delta(\mathbf{y} - \mathbf{y}'). \quad (40)$$

Using the orthogonality relations (37)–(39) we can project out the particle and antiparticle amplitudes from (34):

$$a(\mathbf{k}) = e^{ik_\lambda z^\lambda} (u_k, \Sigma \Phi)_y, \quad (41)$$

$$b(\mathbf{k}) = e^{ik_\lambda z^\lambda} (u_k, \Sigma \Phi^*)_y.$$

Thus the value of  $\Phi(\mathbf{y})$  on any initial light cone suffices to determine  $a(\mathbf{k})$  and  $b(\mathbf{k})$ , and hence the value of  $\Phi(\mathbf{y})$  on all later light cones. Contrast this behavior with that of the usual theory, where both  $\Phi$  and  $\partial\Phi/\partial x^0$  are needed as initial data. We must make one proviso, however. Integrating (28) with respect to  $y$  yields

$$i \frac{\partial \Phi}{\partial \tau} = v_\lambda \int_0^1 (R^\lambda \cdot \Phi(\alpha \mathbf{y})) d\alpha + \frac{C(\tau, \hat{\mathbf{y}})}{y}. \quad (42)$$

(The notation in the integrand means that  $R^\lambda \Phi$  is to be evaluated at  $\alpha \mathbf{y}$  before integrating over  $\alpha$ .) An arbitrary function  $C(\tau, \hat{\mathbf{y}})$  appears in (42), indicating that the evolution is not unique if one allows solutions which fall off with

distance like  $y^{-1}$ . Such solutions would involve infinite energies and would not be normalizable.

The significance of the operators  $R^\lambda$  may be seen from the expression for the expectation value of the momentum operator  $p_{OP}^\lambda$  defined by (19). We have

$$\langle p_{OP}^\lambda \rangle = \int dS_k [|a(\mathbf{k})|^2 + |b(\mathbf{k})|^2] \hbar k^\lambda = (\Phi, R^\lambda \Phi)_y, \quad (43)$$

where use has been made of (36) and (40). Despite the appearance of (43),  $R^\lambda$  cannot be regarded as a momentum operator in  $\mathcal{H}_y$  because its components do not mutually commute. Further, the identity (32) shows that the eigenvalues of  $\eta_{\lambda\mu} R^\lambda R^\mu$  can have either sign.

The form of (28)–(31) is quite suggestive. Comparison with the expressions (7) for the classical momentum vector  $p^\lambda$  shows  $\hbar^2 R^\lambda$  is the Hermitian part of the operator obtained from  $(\mathbf{y} \cdot \boldsymbol{\pi}) p^\lambda$  by the replacement  $\boldsymbol{\pi} \rightarrow \boldsymbol{\pi}_{OP}$  [see (11)]. The operator  $\hbar \Sigma$  is likewise obtained from the classical quantity  $D \equiv \mathbf{y} \cdot \boldsymbol{\pi}$ . Consider now a solution  $\Phi_+$  of (28) which has no antiparticle content. We can heuristically write (28) as

$$i\hbar \frac{\partial \Phi_+}{\partial \tau} = v_\lambda \hbar \Sigma^{-1} R^\lambda \Phi_+, \quad (44)$$

ignoring the fact that  $\Sigma$  can have zero eigenvalues. Equation (44) looks like (5), with  $p_{OP}^\lambda = \hbar \Sigma^{-1} R^\lambda$ . Unfortunately this operator is not Hermitian with respect to the scalar product (10), but only with respect to the indefinite product (33), and further, its components do not mutually commute. On these grounds  $\hbar \Sigma^{-1} R^\lambda$  is not acceptable as the momentum operator, even if we manage to give it a rigorous meaning in some domain.

In subsequent sections the eigenvalues and eigenvectors of the operator  $\Sigma$ , defined in (29), will be of importance. A complete, orthonormal set of simultaneous eigenvectors of the commuting Hermitian operators  $\Sigma$  and  $\hat{\mathbf{y}}$ , with respective eigenvalues  $\sigma$  and  $\mathbf{w}$ , is

$$v_{\sigma\mathbf{w}}(\mathbf{y}) = (2\pi)^{-1/2} y^{i\sigma-1} \delta(\hat{\mathbf{y}}, \mathbf{w}). \quad (45)$$

In (45),  $\sigma$  takes any real value in  $(-\infty, \infty)$ ,  $\mathbf{w}$  is any real unit vector, while  $\delta(\hat{\mathbf{y}}, \mathbf{w})$  is the surface Dirac delta function for a unit sphere. Strictly, these functions lie outside  $\mathcal{H}_y$  because they are not normalizable, but nevertheless they comprise a useful expansion set for the elements of  $\mathcal{H}_y$  on account of their orthogonality and completeness:

$$(v_{\sigma\mathbf{w}}, v_{\sigma'\mathbf{w}'})_y = \delta(\sigma - \sigma') \delta(\mathbf{w}, \mathbf{w}'), \quad (46)$$

$$\int d\sigma d^2\mathbf{w} v_{\sigma\mathbf{w}}(\mathbf{y}) v_{\sigma'\mathbf{w}'}^*(\mathbf{y}') = y \delta(\mathbf{y} - \mathbf{y}').$$

In particular, the function  $u_k(\mathbf{y})$  of (35) has the expansion

$$u_k(\mathbf{y}) = \int d\sigma d^2\mathbf{w} u_k(\sigma, \mathbf{w}) v_{\sigma\mathbf{w}}(\mathbf{y}), \quad (47)$$

$$u_k(\sigma, \mathbf{w}) = \frac{ie^{1/2\pi\sigma}}{4\pi^2} \Gamma(1 - i\sigma) (\epsilon_k + \mathbf{k} \cdot \mathbf{w})^{i\sigma-1}.$$

The integration in (47) is over the complete range of the eigenvalues,  $-\infty$  to  $\infty$  for  $\sigma$  and the surface of a unit sphere for  $\mathbf{w}$ . Substituting (47) into the orthogonality relation (37) yields

$$\frac{1}{16\pi^3} \int \sigma^2 d\sigma d^2\mathbf{w} \frac{e^{\pi\sigma}}{\sinh \pi\sigma} (\epsilon_k + \mathbf{k} \cdot \mathbf{w})^{-i\sigma-1} \times (\epsilon_{k'} + \mathbf{k}' \cdot \mathbf{w})^{i\sigma-1} = \epsilon_k \delta(\mathbf{k} - \mathbf{k}'). \quad (48)$$

If we take the complex conjugate of (48) and replace the variable of integration  $\sigma$  by  $-\sigma$ , the left-hand side assumes the same form except that  $e^{\pi\sigma}$  is replaced by  $-e^{-\pi\sigma}$ . Averaging this integral with (48) then yields the important identity

$$\frac{1}{16\pi^3} \int \sigma^2 d\sigma d^2\mathbf{w} (\epsilon_k + \mathbf{k} \cdot \mathbf{w})^{-i\sigma-1} (\epsilon_{k'} + \mathbf{k}' \cdot \mathbf{w})^{i\sigma-1} = \epsilon_k \delta(\mathbf{k} - \mathbf{k}'). \quad (49)$$

In a similar manner (47) substituted into the completeness relation (40) leads to

$$\frac{1}{16\pi^3} \int dS_k \sigma \sigma' (\epsilon_k + \mathbf{k} \cdot \mathbf{w})^{i\sigma-1} (\epsilon_{k'} + \mathbf{k}' \cdot \mathbf{w}')^{-i\sigma'-1} = \frac{1}{2} \delta(\sigma - \sigma') \delta(\mathbf{w}, \mathbf{w}') + \delta(\sigma + \sigma') G(\sigma, \mathbf{w}, \mathbf{w}'). \quad (50)$$

The form of the function  $G$  is left undetermined by this argument, but is proved in Appendix A to be

$$G(\sigma, \mathbf{w}, \mathbf{w}') = - (i\sigma/4\pi) (\frac{1}{2} \kappa^2)^{i\sigma} (1 - \mathbf{w} \cdot \mathbf{w}')^{+i\sigma-1}. \quad (51)$$

The identities (49) and (50) will play an important role in Sec. III.

### III. QUANTIZATION OF THE FREE SYSTEM

#### A. Introduction

Passage from classical to quantum theory is notoriously ambiguous<sup>5</sup> because of the problem of how to order noncommuting factors. The classical system we are dealing with here is that of one free charged scalar boson whose dynamics are described in terms of the light cone coordinates  $(\tau, y^1, y^2, y^3)$  belonging to an observer with trajectory (1). The energy momentum four-vector  $p^\lambda$  of this system is given by (7) as a function of the Hamiltonian conjugate variables  $\mathbf{y}$  and  $\boldsymbol{\pi}$ . Suppose we attempt in (7) the substitutions  $\mathbf{y} \rightarrow \mathbf{y}_{OP}$  and  $\boldsymbol{\pi} \rightarrow \boldsymbol{\pi}_{OP}$  according to (11). How are we to order the various noncommuting factors, and how are we to interpret  $(\mathbf{y} \cdot \boldsymbol{\pi})^{-1}$  when this denominator becomes an operator? Classically this denominator causes no problem, because  $\mathbf{y} \cdot \boldsymbol{\pi}$  is intrinsically positive, and even when  $y \rightarrow 0$  the ratio  $\mathbf{y} \cdot \boldsymbol{\pi}/y$  stays finite.

Because of the above difficulties we adopt a different approach, that of defining the energy-momentum operator  $p_{OP}^\lambda$  by its eigenvectors and eigenvalues. Guided by Klein-Gordon theory, we seek a complete orthonormal set of states  $|\Psi_{kq}\rangle \in \mathcal{H}_{\text{phys}}$ , labeled by a future pointing four-vector  $k^\lambda$  lying on the mass shell (14), and additionally by a charge index  $q = \pm 1$ . Thus

$$\langle \Psi_{kq} | \Psi_{k'q'} \rangle = \epsilon_k \delta(\mathbf{k} - \mathbf{k}') \delta_{qq'}, \quad (52)$$

$$\sum_q \int dS_k |\Psi_{kq}\rangle \langle \Psi_{kq}| = I,$$

where  $I$  is the unit operator in  $\mathcal{H}_{\text{phys}}$ . Having found such a set we define the momentum and charge operators by



$$p_{OP}^\lambda = \sum_q \int dS_k |\Psi_{kq}\rangle \hbar k^\lambda \langle \Psi_{kq}|, \quad (53)$$

$$Q = \sum_q \int dS_k |\Psi_{kq}\rangle e q \langle \Psi_{kq}|. \quad (54)$$

By construction,  $|\Psi_{kq}\rangle$  is a simultaneous eigenstate of these operators with eigenvalues  $\hbar k^\lambda$  and  $eq$ , respectively.

The above program is also subject to ambiguity. We must first decide how the Hilbert space of physical states  $\mathcal{H}_{\text{phys}}$  is related to  $\mathcal{H}_y$ . Then we must decide which of the infinite number of complete orthonormal sets satisfying (52) is the one that corresponds most closely in the classical limit to the classical system described by (7).

Section III B considers the specification of  $\mathcal{H}_{\text{phys}}$ , based on resolution of  $\psi(\mathbf{y}) \in \mathcal{H}_y$  into particle and antiparticle amplitudes. Section III C deals with the construction of complete orthonormal sets and then the remaining parts of Sec. III concern the application of these sets to (53) and (54). The modifications necessary for the incorporation of electromagnetic interactions are considered in Sec. IV.

## B. Particle and antiparticle amplitudes and the specifications of $\mathcal{H}_{\text{phys}}$

Let us first see how to accommodate antiparticles in the classical Hamilton theory. A free classical particle has a future-pointing timelike momentum vector  $p^\lambda$ , so that the quantity  $D = \mathbf{y} \cdot \boldsymbol{\pi} = y p^0 + \mathbf{y} \cdot \boldsymbol{\pi}$  is necessarily positive. The same is true for an antiparticle. Thus classically we could treat particle and antiparticle as two disjoint systems, with conjugate variables  $\mathbf{y}_p, \boldsymbol{\pi}_p$  and  $\mathbf{y}_a, \boldsymbol{\pi}_a$ , respectively. In this approach the phase spaces are restricted by  $\mathbf{y}_p \cdot \boldsymbol{\pi}_p > 0$  and  $\mathbf{y}_a \cdot \boldsymbol{\pi}_a > 0$ , and the same functional form (7) applies to the evolution generators  $p^\lambda(\mathbf{y}_p, \boldsymbol{\pi}_p)$  and  $p^\lambda(\mathbf{y}_a, \boldsymbol{\pi}_a)$ , which are interpreted, respectively, as the particle and antiparticle momentum vectors. However, particle and antiparticle can be treated alternatively as a single system by exploiting Feynman's idea that an antiparticle behaves like a particle traveling backwards in time. In this alternative approach we enlarge the classical phase space to allow all values of  $\mathbf{y} \cdot \boldsymbol{\pi}$ , both positive and negative. The variables  $\mathbf{y}$  and  $\boldsymbol{\pi}$  evolve according to

$$\frac{d\mathbf{y}}{d\tau} = v^\lambda \{\mathbf{y}, P_\lambda\}, \quad \frac{d\boldsymbol{\pi}}{d\tau} = v^\lambda \{\boldsymbol{\pi}, P_\lambda\}, \quad (55)$$

with the function  $P^\lambda(\mathbf{y}, \boldsymbol{\pi})$  defined by

$$P^0 = \frac{1}{2} (\mathbf{y} \cdot \boldsymbol{\pi})^{-1} (\boldsymbol{\pi}^2 + m^2 c^2) y, \quad (56)$$

$$\mathbf{P} = \boldsymbol{\pi} - \frac{1}{2} (\mathbf{y} \cdot \boldsymbol{\pi})^{-1} (\boldsymbol{\pi}^2 + m^2 c^2) \mathbf{y}.$$

Here  $P^\lambda(\mathbf{y}, \boldsymbol{\pi})$  takes the same form as  $p^\lambda(\mathbf{y}, \boldsymbol{\pi})$  in (7), except that it is now defined for both signs of  $\mathbf{y} \cdot \boldsymbol{\pi}$ . If  $\mathbf{y} \cdot \boldsymbol{\pi} > 0$  we interpret  $\mathbf{y}_p = \mathbf{y}$ ,  $\boldsymbol{\pi}_p = \boldsymbol{\pi}$  as the conjugate variables for a particle with momentum  $p_p^\lambda = P^\lambda$ . On the other hand, if  $\mathbf{y} \cdot \boldsymbol{\pi} < 0$  then we have an antiparticle with conjugate variables  $\mathbf{y}_a = \mathbf{y}$ ,  $\boldsymbol{\pi}_a = -\boldsymbol{\pi}$  and momentum  $p_a^\lambda = -P^\lambda$ . In each case the momentum  $p^\lambda = \text{sgn}(\mathbf{y} \cdot \boldsymbol{\pi}) P^\lambda$  is future pointing timelike, but the evolution generator is  $P^\lambda$  rather than  $p^\lambda$ . The charge is  $e \text{sgn}(\mathbf{y} \cdot \boldsymbol{\pi})$  (Ref. 3).

The above suggests that in quantum theory, particles and antiparticles should be associated, respectively, with the

positive and negative eigenvalues of  $\Sigma$ , the operator analog of  $D/\hbar$  defined by (29). As we saw in Sec. II B,  $\Sigma$  has a complete orthonormal set of eigenfunctions  $v_{\sigma\omega}(\mathbf{y})$ , given by (45), the eigenvalue  $\sigma$  taking all values in  $(-\infty, \infty)$ . Any element  $\psi(\mathbf{y}) \in \mathcal{H}_y$  has the decomposition

$$\psi(\mathbf{y}) = \psi_+(\mathbf{y}) + \psi_-(\mathbf{y}), \quad (57)$$

where

$$\psi_+(\mathbf{y}) = \int_0^\infty d\sigma \int d^2\mathbf{w} v_{\sigma\omega}(\mathbf{y}) (v_{\sigma\omega}, \psi)_y,$$

$$\psi_-(\mathbf{y}) = \int_{-\infty}^0 d\sigma \int d^2\mathbf{w} v_{\sigma\omega}(\mathbf{y}) (v_{\sigma\omega}, \psi)_y.$$

The above equations may be written

$$\psi_+(\mathbf{y}) = \Theta(\Sigma)\psi(\mathbf{y}), \quad \psi_-(\mathbf{y}) = \Theta(-\Sigma)\psi(\mathbf{y}),$$

where the step functions  $\Theta(\Sigma)$  and  $\Theta(-\Sigma)$  are projection operators with matrix elements

$$\langle \mathbf{y} | \Theta(\Sigma) | \mathbf{y}' \rangle = \int_0^\infty d\sigma \int d^2\mathbf{w} v_{\sigma\omega}(\mathbf{y}) v_{\sigma\omega}^*(\mathbf{y}'),$$

$$= \frac{i\delta(\hat{\mathbf{y}}, \hat{\mathbf{y}}')}{2\pi y y' [\log(y/y') + i\epsilon]},$$

$$= \frac{1}{2} y \delta(\mathbf{y} - \mathbf{y}') + \frac{i}{2\pi} \mathcal{P} \frac{\delta(\hat{\mathbf{y}}, \hat{\mathbf{y}}')}{y y' \log(y/y')}, \quad (58)$$

$$\langle \mathbf{y} | \Theta(-\Sigma) | \mathbf{y}' \rangle = \langle \mathbf{y} | \Theta(\Sigma) | \mathbf{y}' \rangle^*. \quad (59)$$

In (58)  $\epsilon \rightarrow +0$  and  $\mathcal{P}$  denotes the principal value. An equivalent form is

$$\Theta(\pm \Sigma)\psi(\mathbf{y}) = \frac{1}{2} \psi(\mathbf{y}) \pm \frac{i}{2\pi} \mathcal{P} \int_0^\infty \frac{\psi(\alpha\mathbf{y}) d\alpha}{-\log \alpha}. \quad (60)$$

Corresponding to (57),  $\mathcal{H}_y$  has the decomposition

$$\mathcal{H}_y = \mathcal{H}_y^+ \oplus \mathcal{H}_y^-, \quad (61)$$

where the two subspaces  $\mathcal{H}_y^+$  and  $\mathcal{H}_y^-$  are themselves Hilbert spaces, whose elements satisfy  $\Theta(\Sigma)\psi_+ = \psi_+$  and  $\Theta(-\Sigma)\psi_- = \psi_-$ , respectively.

It remains now to relate  $\psi(\mathbf{y}) \in \mathcal{H}_y$  to physical states  $|\Psi\rangle \in \mathcal{H}_{\text{phys}}$ . By analogy with Klein-Gordon theory let us suppose that physical states  $|\Psi\rangle$  can be represented by

$$|\Psi(\mathbf{y})\rangle = \begin{bmatrix} \psi_p(\mathbf{y}) \\ \psi_a(\mathbf{y}) \end{bmatrix}, \quad (62)$$

with  $\psi_p(\mathbf{y})$  and  $\psi_a(\mathbf{y})$  being the particle and antiparticle amplitudes. The appropriate scalar product is

$$(\Psi, \Psi')_{\text{phys}} = \int \frac{d^3\mathbf{y}}{y} [\psi_p^*(\mathbf{y}) \psi_p'(\mathbf{y}) + \psi_a^*(\mathbf{y}) \psi_a'(\mathbf{y})]. \quad (63)$$

Let us make the identification

$$\psi_p(\mathbf{y}) = \psi_+(\mathbf{y}) = \Theta(\Sigma)\psi(\mathbf{y}),$$

$$\psi_a(\mathbf{y}) = [\psi_-(\mathbf{y})]^* = \Theta(\Sigma)\psi^*(\mathbf{y}), \quad (64)$$

$$\mathcal{H}_{\text{phys}} = \mathcal{H}_y^+ \oplus \mathcal{H}_y^+.$$

Thus both components of the physical state  $\Psi(\mathbf{y})$  belong to the positive eigenspace of  $\Sigma$ :

$$\Theta(\Sigma)\Psi(\mathbf{y}) = \Psi(\mathbf{y}).$$

This is the quantum analog of the classical result that both  $\mathbf{y}_p \cdot \boldsymbol{\pi}_p$  and  $\mathbf{y}_a \cdot \boldsymbol{\pi}_a$  are positive.

The hypothesis given by (62)–(64) is similar to that of Klein–Gordon and Dirac theory in that particle and antiparticle amplitudes must be projected out from the coordinate space wave function in order to define a physical state. We have a one-to-one mapping, albeit nonlinear, between the points  $\Psi(\mathbf{y}) \in \mathcal{H}_{\text{phys}}$  and the points  $\Psi(\mathbf{y}) \in \mathcal{H}_y$ :

$$\Psi(\mathbf{y}) = \begin{bmatrix} \psi_p(\mathbf{y}) \\ \psi_a(\mathbf{y}) \end{bmatrix} = \Theta(\Sigma) \begin{bmatrix} \psi(\mathbf{y}) \\ \psi^*(\mathbf{y}) \end{bmatrix}, \quad (65)$$

$$\psi(\mathbf{y}) = \psi_p(\mathbf{y}) + \psi_a^*(\mathbf{y}).$$

The momentum operator  $p_{\text{OP}}^\lambda$ , whose explicit form we have yet to determine, governs the evolution of the physical states  $\Psi(\mathbf{y})$  according to (5). Analogy with the classical evolution equation (55) suggests that there should also exist an operator  $P_{\text{OP}}^\lambda$  which determines the evolution of coordinate space wave functions  $\psi(\mathbf{y})$  according to

$$i\hbar \frac{\partial \psi}{\partial \tau} = v_\lambda P_{\text{OP}}^\lambda \psi. \quad (66)$$

Note that the operators  $p_{\text{OP}}^\lambda$  and  $P_{\text{OP}}^\lambda$  act in different Hilbert spaces, viz  $\mathcal{H}_{\text{phys}}$  and  $\mathcal{H}_y$ , respectively. We now consider the problem of how to give effect to the ansatz (53) for determining  $p_{\text{OP}}^\lambda$ , and how to specify  $P_{\text{OP}}^\lambda$  in terms of the latter.

### C. Complete and orthonormal sets

We seek complete, orthonormal sets of states  $|\Psi_{kq}\rangle \in \mathcal{H}_{\text{phys}}$  satisfying (52), which states are to be interpreted as energy-momentum-charge eigenstates according to (53) and (54). Let  $\psi_{kq}(\mathbf{y})$  be the element of  $\mathcal{H}_y$  from which  $\Psi_{kq}(\mathbf{y})$  is derived by the prescription (65). Following the ideas of the previous section, we assume that the particle states and the antiparticle states correspond, respectively, to the positive and to the negative eigenspaces of  $\Sigma$ . Thus

$$\Theta(\Sigma)\psi_{k1}(\mathbf{y}) = \psi_{k1}(\mathbf{y}), \quad (67)$$

$$\Theta(-\Sigma)\psi_{k(-1)}(\mathbf{y}) = \psi_{k(-1)}(\mathbf{y}).$$

Motivated by the desire to achieve maximum particle–antiparticle symmetry we make the additional assumption  $\psi_{k(-1)}(\mathbf{y}) = \psi_k^*(\mathbf{y})$ , which is consistent with (67). For brevity we write  $\psi_k(\mathbf{y})$  for  $\psi_{k1}(\mathbf{y})$  so that (67) becomes

$$\psi_{k1}(\mathbf{y}) = \psi_k(\mathbf{y}) = \psi_k^*(\mathbf{y}), \quad (68)$$

$$\Theta(\Sigma)\psi_k(\mathbf{y}) = \psi_k(\mathbf{y}).$$

The orthogonality and completeness relations (52) now reduce to

$$(\psi_k, \psi_{k'})_y = \epsilon_k \delta(\mathbf{k} - \mathbf{k}'), \quad (69)$$

$$\int dS_k \psi_k(\mathbf{y}) \psi_k^*(\mathbf{y}') = \langle \mathbf{y} | \Theta(\Sigma) | \mathbf{y}' \rangle. \quad (70)$$

Since  $\psi_k(\mathbf{y}) \in \mathcal{H}_y^+$ ,  $\Theta(\Sigma)$  acts as the unit operator in the completeness relation (70). [See (58).] Having found a solution of (68)–(70) we can reconstruct  $\Psi_{kq}(\mathbf{y})$  by

$$\Psi_{k1}(\mathbf{y}) = \begin{bmatrix} \psi_k(\mathbf{y}) \\ 0 \end{bmatrix}, \quad (71)$$

$$\Psi_{k(-1)}(\mathbf{y}) = \begin{bmatrix} 0 \\ \psi_k(\mathbf{y}) \end{bmatrix}. \quad (72)$$

[See (62), (64), and (65).]

The functions  $\Psi_{kq}(\mathbf{y})$  define a unitary mapping between the points  $\Psi(\mathbf{y}) \in \mathcal{H}_{\text{phys}}$  and  $\phi(\mathbf{k}) \in \mathcal{H}_k$  [see (17)]:

$$a_k = (\Psi_{k1}, \Psi)_{\text{phys}} = (\psi_k, \psi)_y, \quad (73)$$

$$b_k = (\Psi_{k(-1)}, \Psi)_{\text{phys}} = (\psi_k, \psi^*)_y.$$

The inverse mapping is

$$\Psi(\mathbf{y}) = \int dS_k [a_k \Psi_{k1}(\mathbf{y}) + b_k \Psi_{k(-1)}(\mathbf{y})], \quad (74)$$

$$\psi(\mathbf{y}) = \int dS_k [a_k \psi_k(\mathbf{y}) + b_k^* \psi_k^*(\mathbf{y})].$$

We now show that angular momentum considerations dictate that  $\Psi_{kq}(\mathbf{y})$ , and hence  $\psi_k(\mathbf{y})$ , be taken as a function solely of the SO(1,3) scalar

$$\zeta = -k_\lambda y^\lambda = \epsilon_k y + \mathbf{k} \cdot \mathbf{y}. \quad (75)$$

[This is analogous to the result of nonrelativistic quantum theory that the function which maps coordinate space onto momentum space, namely  $(2\pi)^{-3/2} \exp(i\mathbf{k} \cdot \mathbf{x})$ , depends solely on the SO(3) scalar  $\mathbf{k} \cdot \mathbf{x}$ .] Let us derive the form of the angular momentum tensor operator using Wigner's  $P_R$  prescription.<sup>6</sup> An infinitesimal Lorentz transformation parametrized by the antisymmetric matrix  $\omega_{\nu\mu}$ ,

$$y'^\lambda = y^\lambda + \eta^{\lambda\nu} \omega_{\nu\mu} y^\mu,$$

changes the functional form of  $\Psi(\mathbf{y})$  to

$$\Psi'(\mathbf{y}) = (I - \frac{1}{2} i \omega_{\lambda\mu} j_y^{\lambda\mu}) \Psi(\mathbf{y}),$$

where the components of the angular momentum operator  $j_y^{\lambda\mu}$  are

$$(j_y^{23}, j_y^{31}, j_y^{12}) = -i\hbar \mathbf{y} \times \frac{\partial}{\partial \mathbf{y}}, \quad (76)$$

$$(j_y^{01}, j_y^{02}, j_y^{03}) = i\hbar \mathbf{y} \frac{\partial}{\partial y}.$$

We now demand that the unitary transformation (73) and (74) should transform  $j_k^{\lambda\mu}$  given by (22) and (23) into  $j_y^{\lambda\mu}$ . This will be the case if and only if

$$(j_k^{\lambda\mu} + j_y^{\lambda\mu}) \Psi_{kq}(\mathbf{y}) = 0.$$

The above equation is the condition that  $\Psi_{kq}(\mathbf{y})$  is unchanged by a combined Lorentz transformation of the two vectors  $k^\lambda, y^\lambda$ , which implies that  $\Psi_{kq}(\mathbf{y})$ , and hence  $\psi_k(\mathbf{y})$  are functions only of the scalar  $\zeta$  defined by (75):

$$\psi_k(\mathbf{y}) = f(\zeta). \quad (77)$$

Let us now expand  $\psi_k(\mathbf{y})$  in terms of the complete, orthonormal set  $v_{\sigma\mathbf{w}}(\mathbf{y})$  given in (45). Only terms with  $\sigma \geq 0$  are needed on account of (68). The expansion coefficients are

$$\begin{aligned} \psi_{k\sigma\mathbf{w}} &= (v_{\sigma\mathbf{w}}, \psi_k)_y, \\ &= F(\sigma) (\epsilon_k + \mathbf{k} \cdot \mathbf{w})^{i\sigma - 1}, \end{aligned} \quad (78)$$

where

$$F(\sigma) = (2\pi)^{-1/2} \int_0^\infty f(\zeta) \zeta^{-i\sigma} d\zeta.$$

Hence

$$\begin{aligned}\psi_k(\mathbf{y}) &= \int_0^\infty d\sigma \int d^2\mathbf{w} \psi_{k\sigma\mathbf{w}} v_{\sigma\mathbf{w}}(\mathbf{y}), \\ &= (2\pi)^{-1/2} \int_0^\infty d\sigma F(\sigma) \xi^{i\sigma-1}.\end{aligned}\quad (79)$$

The function  $F(\sigma)$  must be chosen such that the orthogonality and completeness relations (69) and (70) are satisfied. Substituting (79) into these equations yields

$$\begin{aligned}\int_0^\infty |F(\sigma)|^2 d\sigma \int d^2\mathbf{w} (\epsilon_k + \mathbf{k}\cdot\mathbf{w})^{-i\sigma-1} (\epsilon_{k'} + \mathbf{k}'\cdot\mathbf{w})^{i\sigma-1} \\ = \epsilon_k \delta(\mathbf{k} - \mathbf{k}'),\end{aligned}\quad (80)$$

$$\begin{aligned}\int dS_k F(\sigma) F^*(\sigma') (\epsilon_k + \mathbf{k}\cdot\mathbf{w})^{i\sigma-1} (\epsilon_k + \mathbf{k}\cdot\mathbf{w}')^{-i\sigma'-1} \\ = \delta(\sigma - \sigma') \delta(\mathbf{w}, \mathbf{w}').\end{aligned}\quad (81)$$

The form of these equations is similar to that of (49) and (50), except that here only positive values of  $\sigma$  and  $\sigma'$  are allowed.

We now prove a theorem derived from (49) and (50).

**Theorem:** A solution of (80) and (81) is given by any function  $F(\sigma)$  whose modulus is  $(2\pi)^{-3/2}\sigma$ , i.e., we may write

$$F(\sigma) = (2\pi)^{-3/2} i \sigma e^{ig(\sigma)}, \quad \sigma \geq 0, \quad (82)$$

where  $g(\sigma)$  is an arbitrary real function. The proof of this theorem depends on two lemmas.

**Lemma 1:** For any  $F_1(\sigma)$  and  $F_2(\sigma)$  defined in  $-\infty < \sigma < \infty$ , the integral

$$\begin{aligned}J = \int_{-\infty}^{\infty} F_1^*(\sigma) F_2(\sigma) d\sigma \\ \times \int d^2\mathbf{w} (\epsilon_k + \mathbf{k}\cdot\mathbf{w})^{-i\sigma-1} (\epsilon_{k'} + \mathbf{k}'\cdot\mathbf{w})^{i\sigma-1}\end{aligned}$$

is unchanged if  $\mathbf{k}$  and  $\mathbf{k}'$  are interchanged.

**Lemma 2:** The integral  $J$  of Lemma 1 is given by

$$\begin{aligned}J = \int_0^\infty [F_1^*(\sigma) F_2(\sigma) + F_1^*(-\sigma) F_2(-\sigma)] d\sigma \\ \times \int d^2\mathbf{w} (\epsilon_k + \mathbf{k}\cdot\mathbf{w})^{-i\sigma-1} (\epsilon_{k'} + \mathbf{k}'\cdot\mathbf{w})^{i\sigma-1}.\end{aligned}$$

Lemma 1 is proved by observing that  $J$  is the Lorentz invariant scalar product  $(f_1, f_2)_y$ , where

$$\begin{aligned}f_1(\mathbf{y}) &= (2\pi)^{-1/2} \int_{-\infty}^{\infty} d\sigma F_1(\sigma) (-k_\lambda y^\lambda)^{i\sigma-1}, \\ f_2(\mathbf{y}) &= (2\pi)^{-1/2} \int_{-\infty}^{\infty} d\sigma F_2(\sigma) (-k'_\lambda y^\lambda)^{i\sigma-1}.\end{aligned}$$

Since both  $k_\lambda$  and  $k'_\lambda$  lie on the mass shell (14), the only scalars we can form from these vectors are  $k^\lambda k'_\lambda$  and  $\kappa = mc/\hbar$  (Ref. 7). Hence  $J$  is a function only of these scalars, and is consequently symmetric in  $k_\lambda$  and  $k'_\lambda$ , proving the first lemma. To prove Lemma 2, divide the range of integration over  $\sigma$  into the two intervals  $(0, \infty)$  and  $(-\infty, 0)$ . If in the second integral we replace  $\sigma$  by  $-\sigma$  and invoke Lemma 1, then we obtain the form for  $J$  given by Lemma 2.

With the aid of these lemmas we can now prove the theorem. Applying Lemma 2 to (49), with  $F_1(\sigma) = F_2(\sigma) = \sigma$  we obtain

$$\begin{aligned}(2\pi)^{-3} \int_0^\infty \sigma^2 d\sigma \int d^2\mathbf{w} (\epsilon_k + \mathbf{k}\cdot\mathbf{w})^{-i\sigma-1} \\ \times (\epsilon_{k'} + \mathbf{k}'\cdot\mathbf{w})^{i\sigma-1} = \epsilon_k \delta(\mathbf{k} - \mathbf{k}').\end{aligned}\quad (83)$$

With the restriction to positive  $\sigma$  values (50) becomes

$$\begin{aligned}(2\pi)^{-3} \int dS_k \sigma \sigma' (\epsilon_k + \mathbf{k}\cdot\mathbf{w})^{i\sigma-1} \\ \times (\epsilon_k + \mathbf{k}\cdot\mathbf{w}')^{-i\sigma'-1} = \delta(\sigma - \sigma') \delta(\mathbf{w}, \mathbf{w}').\end{aligned}\quad (84)$$

From (83) and (84) we see that any  $F(\sigma)$  of the form (82) will satisfy (80) and (81), which establishes the theorem.

To summarize Sec. III C, a complete orthonormal set of states  $\Psi_{kq}(\mathbf{y})$  is given by (71) and (72), with

$$\begin{aligned}\psi_k(\mathbf{y}) &= f(-k_\lambda y^\lambda) \\ &= (2\pi)^{-2i} \int_0^\infty d\sigma \sigma e^{ig(\sigma)} (-k_\lambda y^\lambda)^{i\sigma-1},\end{aligned}\quad (85)$$

$g(\sigma)$  being any arbitrary real function. Some examples for particular values of  $g(\sigma)$  are given in Appendix B. The question of what function  $g(\sigma)$  should be adopted will now be considered.

#### D. Choice of the phase function $g(\sigma)$

Equations (71), (72), and (85) together define a complete orthonormal set of states  $\Psi_{kq}(\mathbf{y})$ . Consider an arbitrary element  $\Psi(\mathbf{y}) \in \mathcal{H}_{\text{phys}}$  derived from the element  $\psi(\mathbf{y}) \in \mathcal{H}_y$  via the mapping (65). According to (53) the action of the momentum operator in  $\mathcal{H}_{\text{phys}}$  is

$$P_{\text{OP}}^\lambda \Psi(\mathbf{y}) = \int dS_k \psi_k(\mathbf{y}) \hbar k^\lambda \begin{bmatrix} (\psi_k, \psi_p)_y \\ (\psi_k, \psi_a)_y \end{bmatrix}.\quad (86)$$

To be consistent with (66) the evolution operator in  $\mathcal{H}_y$  has to act according to

$$\begin{aligned}P_{\text{OP}}^\lambda \psi(\mathbf{y}) &= \int dS_k \hbar k^\lambda [\psi_k(\mathbf{y}) (\psi_k, \psi)_y \\ &\quad - \psi_k^*(\mathbf{y}) (\psi_k^*, \psi)_y].\end{aligned}\quad (87)$$

Now any choice of the real function  $g(\sigma)$  gives rise to a unitary representation in  $\mathcal{H}_{\text{phys}}$  of the Poincaré group, based on the infinitesimal generators  $p_{\text{OP}}^\lambda$  and  $j_y^{\lambda\mu}$  [see (76)]. This representation is unitarily equivalent via (73) and (74) to the direct sum of two copies of the standard mass  $m$  spin-zero representation,<sup>8</sup> corresponding to particle and antiparticle. A representation is likewise generated in  $\mathcal{H}_y$  by  $P_{\text{OP}}^\lambda$  and  $j_y^{\lambda\mu}$ .

However, we cannot leave  $g(\sigma)$  arbitrary if we wish the quantum theory to correspond in the classical limit to the formalism based on the classical momentum vector of (7). We want  $\mathbf{y}$  to correspond to the light cone coordinate defined by (2), with  $|\psi(\mathbf{y})|^2 d^3\mathbf{y}/y$  representing the probability that a measurement of this coordinate will yield a value in the interval  $(\mathbf{y}, \mathbf{y} + d\mathbf{y})$ . In what follows we present a heuristic argument which makes plausible the particular choice

$$\begin{aligned}g(\sigma) &= \arg[\Gamma(-i\sigma)], \\ \psi_k(\mathbf{y}) &= (2\pi)^{-2i} \int_0^\infty d\sigma \sigma (\pi^{-1} \sigma \sinh \pi \sigma)^{1/2} \\ &\quad \times \Gamma(-i\sigma) (-k_\lambda y^\lambda)^{i\sigma-1},\end{aligned}\quad (88)$$

where  $\Gamma$  denotes the usual gamma function.

Our starting point is the observation that there exists a classical canonical transformation which transforms the classical evolution generator  $P^\lambda$  of (56) into  $-P^\lambda$ . The transformed variables  $\mathbf{\Pi}, Y^\lambda$ , with  $Y^0 = -|\mathbf{Y}|$ , are derived from the generating function  $mc(2y_\lambda Y^\lambda)^{1/2}$ , and take the form

$$\begin{aligned} Y^0 &= -(mc)^{-2}\pi^2 y, \\ \mathbf{Y} &= (mc)^{-2}[\pi^2 \mathbf{y} - 2(\mathbf{y}\cdot\boldsymbol{\pi})\boldsymbol{\pi}], \\ \boldsymbol{\Pi} &= (mc)^2 \boldsymbol{\pi}/\pi^2. \end{aligned} \quad (89)$$

Substitution into (56) leads to

$$P^\lambda(\mathbf{Y}, \boldsymbol{\Pi}) = -P^\lambda(\mathbf{y}, \boldsymbol{\pi}), \quad \mathbf{Y}\cdot\boldsymbol{\Pi} = -\mathbf{y}\cdot\boldsymbol{\pi}. \quad (90)$$

Thus in the  $\mathbf{Y}, \boldsymbol{\Pi}$  system, particle and antiparticle have interchanged roles. Double application of this canonical transformation leads back to the original variables.

These classical results suggest that there ought to be a quantum counterpart to (89) in the form of a unitary operator  $W$  which interchanges the role of particle and antiparticle wave functions in  $\mathcal{H}_y$ . Thus  $W$  should satisfy

$$\begin{aligned} WP_{\text{OP}}^\lambda W^\dagger &= -P_{\text{OP}}^\lambda, \\ W\psi_k &= \psi_k^*, \quad W = W^{-1} = W^\dagger. \end{aligned} \quad (91)$$

An operator equivalent to (89) in  $\mathcal{H}_y$  is

$$\begin{aligned} Y_{\text{OP}}^0 &= \kappa^{-2} y \frac{\partial^2}{\partial \mathbf{y}^2}, \\ \mathbf{Y}_{\text{OP}} &= \kappa^{-2} \left[ -\mathbf{y} \frac{\partial^2}{\partial \mathbf{y}^2} + 2 \frac{\partial}{\partial \mathbf{y}} \left( \mathbf{y} \cdot \frac{\partial}{\partial \mathbf{y}} \right) \right], \end{aligned}$$

where  $\kappa = mc/\hbar$ . These four operators are Hermitian in a dense subspace of  $\mathcal{H}_y$ , mutually commute, and satisfy  $\eta_{\lambda\mu} Y_{\text{OP}}^\lambda Y_{\text{OP}}^\mu = 0$ . A complete orthonormal set of simultaneous eigenfunctions with respective eigenvalues  $Y^\lambda$  ( $Y^0 = -|\mathbf{Y}|$ ) is given by

$$w(y_\lambda Y^\lambda) = (\kappa^2/4\pi) J_0([2\kappa^2 y_\lambda Y^\lambda]^{1/2}). \quad (92)$$

The eigenvalue property may be proved readily by direct differentiation. The orthogonality and completeness relations are

$$\begin{aligned} \int \frac{d^3 \mathbf{y}}{y} w(y_\lambda Y^\lambda) \overline{w(y_\lambda Y'^\lambda)} &= Y \delta(\mathbf{Y} - \mathbf{Y}'), \\ \int \frac{d^3 \mathbf{Y}}{Y} w(y_\lambda Y^\lambda) \overline{w(y'_\mu Y'^\mu)} &= y \delta(\mathbf{y} - \mathbf{y}'), \end{aligned}$$

which may be proved by methods akin to those used in Appendix A.

The appropriate definition of the action of  $W$  on  $\psi(\mathbf{y}) \in \mathcal{H}_y$  is then

$$(W\psi)(\mathbf{y}) = \int \frac{d^3 \mathbf{Y}}{Y} w(y_\lambda Y^\lambda) \psi(\mathbf{Y}). \quad (93)$$

Applying (93) to  $\psi_k(\mathbf{y})$  given by (85) we find (see Appendix C)

$$\begin{aligned} (W\psi_k)(\mathbf{y}) &= -(2\pi)^{-2} i \int_0^\infty d\sigma \sigma e^{ig(\sigma)} \\ &\quad \times \frac{\Gamma(i\sigma)}{\Gamma(-i\sigma)} (-k_\lambda y^\lambda)^{-i\sigma-1}, \end{aligned} \quad (94)$$

which coincides with  $\psi_k^*(\mathbf{y})$  if

$$\begin{aligned} e^{ig(\sigma)} &= [\Gamma(-i\sigma)/\Gamma(i\sigma)]^{1/2}, \\ &= \pm (\pi^{-1} \sigma \sinh \pi \sigma)^{1/2} \Gamma(-i\sigma). \end{aligned}$$

Thus apart from a trivial assignment of sign, (88) is the only choice of the set  $\psi_k(\mathbf{y})$  for which the classical symmetry property (90) goes over to the quantum symmetry property (91) under the unitary mapping  $W$  corresponding to (89).

It does not appear possible to express the momentum eigenfunctions  $\psi_k(\mathbf{y})$  given by (88) in terms of elementary functions. However only the asymptotic behavior for large  $\zeta \equiv -k_\lambda y^\lambda$  will be needed in this paper. The integrand contains the factor  $\exp\{i \arg[\Gamma(-i\sigma)] + i\sigma \log \zeta\}$ , which, for large  $\zeta$ , oscillates rapidly except in the neighborhood of a stationary point of the exponent. The only stationary point occurs at the value of  $\sigma = \zeta - (12\zeta)^{-1} + O(\zeta^{-3})$ . The method of stationary phase then gives the asymptotic behavior for large  $\zeta$  as

$$\psi_k(\mathbf{y}) \sim \frac{i\zeta^{1/2} e^{i\zeta}}{(2\pi)^{3/2}} \left[ 1 + \frac{1}{12\zeta} - \frac{1}{24\zeta^2} + O(\zeta^{-3}) \right]. \quad (95)$$

At the other extreme, when  $\zeta$  is small, the integral is dominated by values of  $\sigma$  of the order of  $-1/\log \zeta$ . The limiting behavior as  $\zeta \rightarrow 0$  is

$$4\pi^2 \zeta \psi_k(\mathbf{y}) \sim (\log \zeta)^{-2} - 2\gamma (\log \zeta)^{-3} + O((\log \zeta)^{-4}), \quad (96)$$

where  $\gamma \equiv 0.57722$  is Euler's constant.

Further creditability will now be given to (88), by showing that this choice of  $g(\sigma)$  leads to the correct nonrelativistic quantum mechanics limit ( $c \rightarrow \infty$ ), and to the correct classical mechanics limit ( $\hbar \rightarrow 0$ ).

## E. The nonrelativistic limit

Consider an observer whose world line passes from  $z^\lambda(\tau_0)$  to  $z^\lambda(\tau)$ . If at proper time  $\tau_0/c$  the observer finds the particle with the momentum eigenfunction  $\psi_k(\mathbf{y})$ , then (66) predicts that at the later proper time  $\tau/c$  his wave function will have evolved to

$$\psi_k(\mathbf{y}, z) = \psi_k(\mathbf{y}) \exp\{-ik_\lambda [z^\lambda(\tau) - z^\lambda(\tau_0)]\}. \quad (97)$$

These functions satisfy the orthogonality relations

$$(\psi_k(\mathbf{y}, z), \psi_{k'}(\mathbf{y}, z))_y = \epsilon_k \delta(\mathbf{k} - \mathbf{k}'),$$

which may be written

$$\int d^3 \mathbf{y} \phi_k^*(\mathbf{y}, z) \phi_{k'}(\mathbf{y}, z) = \delta(\mathbf{k} - \mathbf{k}'), \quad (98)$$

with

$$\phi_k(\mathbf{y}, z) = (\epsilon_k y)^{-1/2} \psi_k(\mathbf{y}, z).$$

We will interpret  $c \rightarrow \infty$  as meaning that both the dimensionless quantities  $\kappa y$  and  $\kappa/k$ , with  $\kappa = mc/\hbar$ , become indefinitely large. In this limit  $\zeta = \epsilon_k y + \mathbf{k}\cdot\mathbf{y}$  becomes large and we may use (95). The leading term is

$$\begin{aligned} \phi_k(\mathbf{y}, z) &\sim i(2\pi)^{-3/2} \exp\{-ik_\lambda [y^\lambda + z^\lambda(\tau) - z^\lambda(\tau_0)]\} \\ &\sim i(2\pi)^{-3/2} \exp\{i\mathbf{k}\cdot\mathbf{x} - i\hbar^{-1}[mc^2 + \hbar^2 k^2/(2m)]t \\ &\quad + ik_\lambda z^\lambda(\tau_0)\}, \end{aligned}$$

where  $x^\lambda = z^\lambda(\tau) + y^\lambda$  [see (2)], and  $x^0 = ct$ . Apart from an irrelevant phase factor, the above is the usual form for the energy-momentum eigenfunction in nonrelativistic quantum mechanics (including the rest energy  $mc^2$ ). Note that in the nonrelativistic limit the integration  $\int d^3\mathbf{y}$  in (98), taken over the past light cone  $\tau = \text{const}$ , reduces to  $\int d^3\mathbf{x}$ , taken over the hyperplane  $t = \text{const}$ .

## F. The classical limit

Following Landau and Lifshitz,<sup>9</sup> we shall seek the classical limit by considering wave functions with large phases and large wavenumbers. A quasiclassical wave function for a particle is of the form

$$\psi(\mathbf{y}, z) = \int dS_k \psi_k(\mathbf{y}) A(\hbar\mathbf{k}) \exp\{i\hbar^{-1}\phi(\hbar\mathbf{k}) - ik_\lambda z^\lambda(\tau)\}.$$

Here  $\phi(\hbar\mathbf{k})$  is real, and we assume that  $A(\hbar\mathbf{k})$  and  $\phi(\hbar\mathbf{k})$  are functions which differ from zero in regions where  $p^\lambda = \hbar k^\lambda$  has a macroscopic value, with  $k$  becoming indefinitely large as we let  $\hbar \rightarrow 0$ . Thus  $\zeta = \epsilon_\lambda y^\lambda + \mathbf{k}\cdot\mathbf{y}$  is likewise large and once again the asymptotic form (95) is applicable. Changing variable to  $\mathbf{p} = \hbar\mathbf{k}$  and using (95) yields

$$\begin{aligned} \psi(\mathbf{y}, z) &= \frac{i}{(2\pi)^{3/2}\hbar^2} \int \frac{d^3\mathbf{p}}{p^0} \zeta^{1/2} A(\mathbf{p}) \\ &\quad \times \exp\{i\hbar^{-1}[\phi(\mathbf{p}) - p_\lambda x^\lambda]\}, \end{aligned} \quad (99)$$

where  $x^\lambda = z^\lambda(\tau) + y^\lambda$  and  $p^0 = (\mathbf{p}^2 + m^2c^2)^{1/2}$ . As  $\hbar \rightarrow 0$  the exponent in (99) oscillates rapidly, and we may use the method of stationary phase. The stationary point is given by

$$\frac{\partial\phi}{\partial\mathbf{p}} = \frac{\mathbf{p}}{p^0} x^0 - \mathbf{x},$$

whose solution  $\mathbf{p} = \mathbf{p}_s(x)$  is a function of  $x^\lambda$ . The stationary phase approximation to (99) then assumes the form

$$\psi = B \exp\{i\hbar^{-1}S(x)\},$$

where  $B$  is a relatively slowly varying function and

$$S(x) = [\phi(\mathbf{p}) - p_\lambda x^\lambda]_{\mathbf{p}=\mathbf{p}_s(x)}. \quad (100)$$

The differential of (100) is

$$dS(x) = -(p_{s\lambda})_\lambda dx^\lambda.$$

Whence  $\partial S(x)/\partial x^\lambda = -(p_{s\lambda})_\lambda$ , which substituted into the identity  $p_s^0 = (\mathbf{p}_s^2 + m^2c^2)^{1/2}$  leads to the Hamilton-Jacobi equation

$$\frac{\partial S}{\partial x^0} + \left[ \left( \frac{\partial S}{\partial \mathbf{x}} \right)^2 + m^2c^2 \right]^{1/2} = 0.$$

Thus a wave packet of quasiclassical form, whose dimensions tend to zero with  $\hbar$ , will follow the classical trajectory of a free particle.

## IV. ELECTROMAGNETIC INTERACTIONS

### A. Classical treatment of particle and antiparticle

When an electromagnetic field derived from a potential  $A^\lambda(x^\mu) = A^\lambda(z^\mu + y^\mu)$  is present and the particle has charge

$e$ , the classical momentum  $p^\lambda(\mathbf{y}, \boldsymbol{\pi})$  is given by (8). This is obtained from the noninteracting form (7) by the prescription  $p^\lambda \rightarrow p^\lambda - (e/c)A^\lambda$ ,  $\boldsymbol{\pi} \rightarrow \boldsymbol{\pi}_E \equiv \boldsymbol{\pi} - (e/c)(\mathbf{A} + \hat{\mathbf{y}}A^0)$ . Note the identity<sup>1</sup>

$$\mathbf{y}\cdot\boldsymbol{\pi}_E = -mcy_\lambda w^\lambda, \quad (101)$$

where  $w^\lambda$  is the four-velocity of the particle at the point where the past light cone with vertex  $z^\lambda$  intersects the particle trajectory. Thus (101) implies that  $\mathbf{y}\cdot\boldsymbol{\pi}_E$  is Lorentz invariant and non-negative. To make the treatment more like that for the noninteracting case, let us effect a canonical transformation  $\mathbf{y}, \boldsymbol{\pi} \rightarrow \mathbf{y}', \boldsymbol{\pi}'$ , which converts the relation  $\mathbf{y}\cdot\boldsymbol{\pi}_E \geq 0$  into  $\mathbf{y}'\cdot\boldsymbol{\pi}' \geq 0$ . Consider the generating function  $F(z^\lambda, \mathbf{y}, \boldsymbol{\pi}') = \mathbf{y}\cdot\boldsymbol{\pi}' + (e/c)\chi(z^\lambda, \mathbf{y})$ , where  $\chi(z^\lambda, \mathbf{y})$  is to be chosen to make  $\mathbf{y}\cdot\boldsymbol{\pi}_E = \mathbf{y}'\cdot\boldsymbol{\pi}'$ . We have

$$\mathbf{y}' = \frac{\partial F}{\partial \boldsymbol{\pi}'} = \mathbf{y}, \quad \boldsymbol{\pi} = \frac{\partial F}{\partial \mathbf{y}} = \boldsymbol{\pi}' + \frac{e}{c} \frac{\partial \chi}{\partial \mathbf{y}}, \quad (102)$$

$$p'_\lambda = p_\lambda + \frac{\partial F}{\partial z^\lambda} = p_\lambda + \frac{e}{c} \frac{\partial \chi}{\partial z^\lambda}.$$

Then  $\mathbf{y}\cdot\boldsymbol{\pi}_E = \mathbf{y}\cdot\boldsymbol{\pi}'$  provided  $\chi$  is any solution of

$$\mathbf{y}\cdot\frac{\partial \chi}{\partial \mathbf{y}} = -y_\kappa A^\kappa. \quad (103)$$

The general solution of (103) is

$$\chi(z^\lambda, \mathbf{y}) = - \int_\gamma^y \frac{y_\kappa}{y} A^\kappa \left( z^\lambda + \frac{ty^\lambda}{y} \right) dt, \quad (104)$$

where  $\gamma$  is an arbitrary function of  $z^\lambda$  and  $\hat{\mathbf{y}}$ . The theory presented in what follows does not depend on the choice of this function  $\gamma$ . [See (114) and (121).] The transformed evolution generator is then given by

$$\begin{aligned} p^0(\mathbf{y}, \boldsymbol{\pi}') &= \frac{e}{c} \left( A^0 + \frac{\partial \chi}{\partial z^0} \right) \\ &\quad + \frac{1}{2} (\mathbf{y}\cdot\boldsymbol{\pi}')^{-1} [(\boldsymbol{\pi}'_E)^2 + m^2c^2] y, \\ \mathbf{p}(\mathbf{y}, \boldsymbol{\pi}') &= \frac{e}{c} \left( \mathbf{A} - \frac{\partial \chi}{\partial \mathbf{z}} \right) + \boldsymbol{\pi}'_E \\ &\quad - \frac{1}{2} (\mathbf{y}\cdot\boldsymbol{\pi}')^{-1} [(\boldsymbol{\pi}'_E)^2 + m^2c^2] \mathbf{y}, \\ \boldsymbol{\pi}'_E &= \boldsymbol{\pi}' - \frac{e}{c} \left( \mathbf{A} + \hat{\mathbf{y}}A^0 - \frac{\partial \chi}{\partial \mathbf{y}} \right). \end{aligned} \quad (105)$$

The above does *not* represent an electromagnetic gauge transformation because  $\chi(z^\lambda, \mathbf{y})$  does not depend on its arguments solely through the combination  $z^\lambda + y^\lambda$ .

A like treatment may be given for an antiparticle, with  $e \rightarrow -e$  in all expressions. Thus we could treat particle and antiparticle as two disjoint systems with conjugate variables  $\mathbf{y}_p, \boldsymbol{\pi}_p$  and  $\mathbf{y}_a, \boldsymbol{\pi}_a$ , respectively, with  $\mathbf{y}_p \cdot \boldsymbol{\pi}_p \geq 0$  and  $\mathbf{y}_a \cdot \boldsymbol{\pi}_a \geq 0$ . The evolution generators  $p_p^\lambda(\mathbf{y}_p, \boldsymbol{\pi}_p)$  and  $p_a^\lambda(\mathbf{y}_a, \boldsymbol{\pi}_a)$  would then be given by (105) with  $\mathbf{y} \rightarrow \mathbf{y}_p$ ,  $\boldsymbol{\pi}' \rightarrow \boldsymbol{\pi}_p$ , and  $\mathbf{y} \rightarrow \mathbf{y}_a$ ,  $\boldsymbol{\pi}' \rightarrow \boldsymbol{\pi}_a$ ,  $e \rightarrow -e$ , respectively. However, just as in the noninteracting case, particle and antiparticle may be treated more elegantly as a single system, with the variables  $\mathbf{y}, \boldsymbol{\pi}'$  now being in an enlarged phase space allowing both signs of  $\mathbf{y}\cdot\boldsymbol{\pi}'$ . Henceforth let us drop the prime and write  $\boldsymbol{\pi}$  instead of  $\boldsymbol{\pi}'$ . The particle-antiparticle system then has the following classical description of its dynamics. The conjugate variables  $\mathbf{y}, \boldsymbol{\pi}$

each vary over the whole of  $\mathbb{R}^3$ . As we change the here-now  $z^\lambda$  from which the past light cone is drawn, the variables  $\mathbf{y}$  and  $\boldsymbol{\pi}$  evolve according to

$$\frac{\partial \mathbf{y}}{\partial z^\lambda} = \{\mathbf{y}, P_\lambda\}, \quad \frac{\partial \boldsymbol{\pi}}{\partial z^\lambda} = \{\boldsymbol{\pi}, P_\lambda\}, \quad (106)$$

where  $P^\lambda$  is given by

$$\begin{aligned} P^0 &= \frac{e}{c} \left( A^0 + \frac{\partial \chi}{\partial z^0} \right) + \frac{1}{2} (\mathbf{y} \cdot \boldsymbol{\pi})^{-1} (\boldsymbol{\pi}_E^2 + m^2 c^2) y, \\ \mathbf{P} &= \frac{e}{c} \left( \mathbf{A} + \frac{\partial \chi}{\partial \mathbf{z}} \right) + \boldsymbol{\pi}_E - \frac{1}{2} (\mathbf{y} \cdot \boldsymbol{\pi})^{-1} (\boldsymbol{\pi}_E^2 + m^2 c^2) \mathbf{y}, \\ \boldsymbol{\pi}_E &= \boldsymbol{\pi} - \frac{e}{c} \left( \mathbf{A} + \hat{\mathbf{y}} A^0 - \frac{\partial \chi}{\partial \mathbf{y}} \right). \end{aligned} \quad (107)$$

If  $\mathbf{y} \cdot \boldsymbol{\pi} > 0$ , we interpret  $\mathbf{y}_p = \mathbf{y}$ ,  $\boldsymbol{\pi}_p = \boldsymbol{\pi}$  as the conjugate variables of a particle with  $p_p^\lambda = P^\lambda$ , and if  $\mathbf{y} \cdot \boldsymbol{\pi} < 0$ ,  $\mathbf{y}_a = \mathbf{y}$ ,  $\boldsymbol{\pi}_a = -\boldsymbol{\pi}$  as those of an antiparticle with  $p_a^\lambda = -P^\lambda$ . In each case  $p^\lambda = \text{sgn}(\mathbf{y} \cdot \boldsymbol{\pi}) P^\lambda$ .

Equation (107) gives  $P^\lambda$  as a quadratic function of  $e/c$ . We may write (107) in the form

$$P^\lambda = P_{(0)}^\lambda + (e/c) P_{(1)}^\lambda + (e/c)^2 P_{(2)}^\lambda, \quad (108)$$

where  $P_{(0)}^\lambda$  is given by (7),  $P_{(1)}^\lambda$  is linear in  $A_\kappa$  and  $\chi$  and  $P_{(2)}^\lambda$  is quadratic in these fields. Explicitly,

$$\begin{aligned} P_{(1)}^\lambda &= \frac{\partial \chi}{\partial z^\lambda} + \{\chi, P_{(0)}^\lambda\} - (\mathbf{y} \cdot \boldsymbol{\pi})^{-1} P_{(0)}^\kappa A_\kappa y^\lambda, \\ P_{(2)}^\lambda &= (2\mathbf{y} \cdot \boldsymbol{\pi})^{-1} \left[ (A^0)^2 - \left( \mathbf{A} - \frac{\partial \chi}{\partial \mathbf{y}} \right)^2 \right] y^\lambda. \end{aligned} \quad (109)$$

[The transition from (107) to (109) is facilitated by use of the identities  $\{P_{(0)}^\lambda, y^\mu\} = g^{\lambda\mu} + (\mathbf{y} \cdot \boldsymbol{\pi})^{-1} y^\lambda P_{(0)}^\mu$  and  $y_\lambda P_{(0)}^\lambda = -\mathbf{y} \cdot \boldsymbol{\pi}$ .]

## B. Invariances of the classical momentum function

The system (106) and (107) has a number of symmetries which one would expect to be preserved upon quantization. These are summarized in (110)–(114) below, and may be verified by direct computation. The evolution generator is given by (107) as a function of  $\mathbf{y}, \boldsymbol{\pi}, e, A_\kappa$ , and  $\chi$ , which dependence will be indicated when necessary by the notation  $P^\lambda(\mathbf{y}, \boldsymbol{\pi}, e, A_\kappa, \chi)$ . Recall that  $A_\kappa$  is a function of  $z^\lambda + y^\lambda$ , and  $\chi$  any function of  $z^\lambda$  and  $\mathbf{y}$  which satisfies (103).

Invariance of the sign  $\mathbf{y} \cdot \boldsymbol{\pi}$  means

$$\{\Theta(\mathbf{y} \cdot \boldsymbol{\pi}), P^\lambda\} = 0. \quad (110)$$

Self-consistency of the evolution equation (106) means

$$\frac{\partial P_\lambda}{\partial z^\mu} - \frac{\partial P_\mu}{\partial z^\lambda} + \{P_\lambda, P_\mu\} = 0. \quad (111)$$

Charge conjugation invariance means

$$P^\lambda(\mathbf{y}, \boldsymbol{\pi}, e, A_\kappa, \chi) = -P^\lambda(\mathbf{y}, -\boldsymbol{\pi}, -e, A_\kappa, \chi). \quad (112)$$

Electromagnetic gauge invariance means, for any  $\Lambda(z^\lambda + y^\lambda)$ ,

$$P^\lambda(\mathbf{y}, \boldsymbol{\pi}, e, A_\kappa, \chi) = P^\lambda \left( \mathbf{y}, \boldsymbol{\pi}, e, A_\kappa + \frac{\partial \Lambda}{\partial z^\kappa}, \chi - \Lambda \right). \quad (113)$$

Invariance under changing  $\chi$  by an arbitrary function  $M(z^\lambda, \hat{\mathbf{y}})$  means

$$\begin{aligned} P_\lambda(\mathbf{y}, \boldsymbol{\pi}, e, A_\kappa, \chi) \\ = P_\lambda \left( \mathbf{y}, \boldsymbol{\pi} + \frac{e}{c} \frac{\partial M}{\partial \mathbf{y}}, e, A_\kappa, \chi - M \right) + \frac{e}{c} \frac{\partial M}{\partial z^\lambda}. \end{aligned} \quad (114)$$

## C. Quantization

Our aim is to find an operator  $P_{\text{OP}}^\lambda$  acting in  $\mathcal{H}_y$  which is a quantum analog of the classical  $P^\lambda$  given by (107), and which governs the evolution of  $\psi(\mathbf{y}) \in \mathcal{H}_y$  according to

$$i\hbar \frac{\partial \psi}{\partial z^\lambda} = P_{\text{OP}^\lambda} \psi. \quad (115)$$

We shall assume that the same relation holds between the Hilbert space  $\mathcal{H}_y$  and  $\mathcal{H}_{\text{phys}}$  as applied without interactions [see (57)–(65)]. Thus particle and antiparticle are again associated, respectively, with the positive and negative eigenspaces  $\mathcal{H}_y^+, \mathcal{H}_y^-$  of the operator  $\Sigma = -i(\mathbf{y} \cdot \partial / \partial \mathbf{y} + 1)$ . This hypothesis makes sense when the classical variables  $\mathbf{y}$  and  $\boldsymbol{\pi}$  [designated  $\boldsymbol{\pi}'$  in (102)] are chosen as in Sec. IV A, because then the sign of  $\mathbf{y} \cdot \boldsymbol{\pi}$  determines the particle's charge, just as was the case without electromagnetic interactions.

Analogy with (108) suggests that  $P_{\text{OP}}^\lambda$  should be quadratic in  $e/c$ :

$$P_{\text{OP}}^\lambda = P_{(0)}^\lambda + (e/c) P_{(1)}^\lambda + (e/c)^2 P_{(2)}^\lambda. \quad (116)$$

Here  $P_{(0)}^\lambda$  is the operator defined by (87) and (88), and the Hermitian operators  $P_{(1)}^\lambda$  and  $P_{(2)}^\lambda$  are functionals of  $A_\kappa$  and  $\chi$  which are, respectively, linear and homogeneous of degree 2. Let us henceforth omit the suffix OP from  $P_{\text{OP}}^\lambda$ , it being understood that  $P^\lambda$  now represents an operator on  $\mathcal{H}_y$ . This operator is a function of  $z^\lambda$  and  $e$ , and a functional of the fields  $A_\kappa$  and  $\chi$ , which dependence will be indicated when needed by the notation  $P^\lambda(z^\lambda, e; A_\kappa, \chi)$ .

In addition to requiring that  $P^\lambda$  should correspond in the classical limit to the classical evolution generator specified by (107), we shall demand that the quantum counterparts of the classical symmetries (110)–(114) shall hold, namely

$$[\Theta(\Sigma), P^\lambda] = 0, \quad (117)$$

$$i\hbar \left( \frac{\partial P_\lambda}{\partial z^\mu} - \frac{\partial P_\mu}{\partial z^\lambda} \right) + [P_\lambda, P_\mu] = 0, \quad (118)$$

$$[P^\lambda(z^\mu, e; A_\kappa; \chi) \psi]^* = -P^\lambda(z^\mu, -e; A_\kappa; \chi) \psi^*, \quad (119)$$

$$P^\lambda(z^\mu, e; A_\kappa; \chi) = P^\lambda \left( z^\mu, e; A_\kappa + \frac{\partial \Lambda}{\partial z^\kappa}; \chi - \Lambda \right), \quad (120)$$

$$\begin{aligned} P_\lambda(z^\mu, e; A_\kappa; \chi) &= \exp \left( -\frac{ie}{\hbar c} M \right) P_\lambda(z^\mu, e; A_\kappa; \chi - M) \\ &\quad \times \exp \left( \frac{ie}{\hbar c} M \right) + \frac{e}{c} \frac{\partial M}{\partial z^\lambda}. \end{aligned} \quad (121)$$

A consequence of (117) is that  $\mathcal{H}_y^+$  and  $\mathcal{H}_y^-$  [see (57)–(61)] are invariant subspaces of  $P^\lambda$ . This implies that if  $\psi \in \mathcal{H}_y$  evolves according to (115), then there exists an operator  $p^\lambda$  on  $\mathcal{H}_{\text{phys}}$  such that  $\Psi$ , the image in  $\mathcal{H}_{\text{phys}}$  of  $\psi$  under the mapping (65), satisfies

$$i\hbar \frac{\partial \Psi}{\partial z^\lambda} = P_\lambda \Psi. \quad (122)$$

Explicitly,

$$P_\lambda \begin{bmatrix} \psi_p \\ \psi_a \end{bmatrix} = \begin{bmatrix} P_\lambda(z^\mu, e; A_\kappa; \chi) \psi_p \\ P_\lambda(z^\mu, -e; A_\kappa; \chi) \psi_a \end{bmatrix}, \quad (123)$$

where use has been made of (119).

The consistency condition (118) implies the existence of a unitary operator  $U(z) \equiv U(z, e; A_\kappa; \chi)$ , a function of  $z^\lambda$  and  $e$  and a functional of  $A_\kappa$  and  $\chi$ , such that

$$P_\lambda = i\hbar \frac{\partial U(z)}{\partial z^\lambda} U^\dagger(z). \quad (124)$$

The wave function  $\psi$  at the current here-now  $z^\lambda$  is related to its value  $\psi_0$  at an initial here-now  $z_0^\lambda$  by

$$\psi = U(z) U^\dagger(z_0) \psi_0. \quad (125)$$

Note that if  $V$  is any unitary operator which does not depend on  $z^\lambda$ , then  $U(z)$  and  $U(z)V$  give the same results in (124) and (125), which equivalence will be denoted

$$U(z)V \cong U(z).$$

The symmetries (119) to (121) now become

$$\begin{aligned} [U(z, e; A_\kappa; \chi) \psi]^* &= U(z, -e; A_\kappa; \chi) \psi^*, \\ U(z, e; A_\kappa; \chi) &\cong U\left(z, e; A_\kappa + \frac{\partial \Lambda}{\partial z^\kappa}; \chi - \Lambda\right), \end{aligned} \quad (126)$$

$$U(z, e; A_\kappa; \chi) \cong \exp\left(\frac{e}{i\hbar c} M\right) U(z, e; A_\kappa; \chi - M).$$

The author has been unable to obtain any solution of (117) to (121), or equivalently of (126), for which  $P^\lambda$  takes the quadratic form (116). However a solution to first order in  $e/c$  may be obtained by writing

$$\begin{aligned} U(z) &= \exp\left[\frac{e}{i\hbar c} \left(H_1 + \frac{e}{c} H_2 + \dots\right)\right] \\ &\times \exp\left[\frac{z_\lambda P^\lambda(0)}{i\hbar}\right], \end{aligned} \quad (127)$$

and expanding the first factor as a power series in  $e/c$ , keeping only the leading terms. Equating terms in  $e/c$  in (124) yields

$$\frac{\partial H_1}{\partial z^\lambda} + \frac{1}{i\hbar} [H_1, P_{(0)\lambda}] = P_{(1)\lambda}. \quad (128)$$

The classical counterpart of (128) is

$$\begin{aligned} \frac{\partial H_{Cl}}{\partial z^\lambda} + \{H_{Cl}, P_{(0)\lambda}\} &= P_{(1)\lambda}, \\ &= \frac{\partial \chi}{\partial z^\lambda} + \{\chi, P_{(0)\lambda}\} - \frac{P_{(0)\kappa} A_\kappa y_\lambda}{\mathbf{y} \cdot \boldsymbol{\pi}}, \end{aligned} \quad (129)$$

upon using (109) [ $P_{(0)\lambda}$  and  $P_{(1)\lambda}$  are now functions of  $\mathbf{y}, \boldsymbol{\pi}$ ]. A solution of (129) is

$$\begin{aligned} H_{Cl} &= \chi + \frac{1}{2} \int_0^\infty P_{(0)\kappa} [A_\kappa(z^\lambda + y^\lambda - tP_{(0)\lambda}^\lambda) \\ &\quad - A_\kappa(z^\lambda + y^\lambda + tP_{(0)\lambda}^\lambda)] dt, \end{aligned} \quad (130)$$

which may be verified by direct substitution. We now seek a

Hermitian operator  $H_1$  on  $\mathcal{H}_y$  which is an analog of (130), and further, is consistent to first order in  $(e/c)$  with the symmetries (117) to (121). Let us specify  $H_1$  by its matrix elements with respect to  $\psi_k$  and  $\psi_k^*$ , the complete orthonormal set of eigenfunctions of  $P_{(0)}^\lambda$ , defined by (88). An operator  $H_1$  which satisfies all the required constraints is given by

$$\begin{aligned} (\psi_k, H_1 \psi_{k'})_y &= \int \frac{d^3 \mathbf{y}}{y} \psi_k^*(\mathbf{y}) \psi_{k'}(\mathbf{y}) \\ &\times \left\{ \chi + \frac{1}{2} \int_0^\infty K^\kappa [A_\kappa(z^\lambda - y^\lambda - tK^\lambda) \right. \\ &\quad \left. - A_\kappa(z^\lambda - y^\lambda + tK^\lambda)] dt \right\}, \end{aligned} \quad (131)$$

$$K_\kappa = k_\kappa + k'_\kappa,$$

$$(\psi_k^*, H_1 \psi_{k'}^*)_y = (\psi_k, H_1 \psi_{k'})^*, \quad (\psi_k^*, H_1 \psi_{k'}) = 0.$$

This form was found by trial and error, guided by the classical result (130). Finally, the matrix elements of  $P_{(1)\lambda}$  may be calculated from (128):

$$\begin{aligned} (\psi_k, P_{(1)\lambda} \psi_{k'})_y &= \left[ \frac{\partial}{\partial z^\lambda} + i(k_\lambda - k'_\lambda) \right] (\psi_k, H_1 \psi_{k'})_y, \\ (\psi_k^*, P_{(1)\lambda} \psi_{k'}^*)_y &= (\psi_k, P_{(1)\lambda} \psi_{k'})_y^*, \\ (\psi_k^*, P_{(1)\lambda} \psi_{k'})_y &= 0. \end{aligned} \quad (132)$$

No satisfactory way of extending these results to second order in  $e/c$  has been found.

## V. DISCUSSION

### A. Summary

In the previous sections we have arrived at the following past light cone quantum picture of a charged boson of spin-zero and nonzero rest mass. An observer at here-now  $z^\lambda$  ascribes to the particle an  $SO(1,3)$  scalar wave function  $\psi(\mathbf{y}, z^\lambda)$ , which belongs to the Hilbert space  $\mathcal{H}_y$  defined by (10). As  $z^\lambda$  varies the wave function evolves according to

$$i\hbar \frac{\partial \psi}{\partial z^\lambda} = P_\lambda \psi, \quad (133)$$

where  $P_\lambda$  is a Hermitian operator in  $\mathcal{H}_y$ . The wave function describes the quantum state on the past light cone with vertex at  $z^\lambda$ , i.e., on the three-surface  $x^\lambda = z^\lambda + y^\lambda$ , with  $|\psi(\mathbf{y}, z^\lambda)| d^3 \mathbf{y}/y$  representing the probability that a position measurement will locate the particle on the light cone between  $\mathbf{y}$  and  $\mathbf{y} + d\mathbf{y}$  ( $\psi$  being normalized to unity).

An important role is played by the operator  $\Sigma = -i(\mathbf{y} \cdot \partial / \partial \mathbf{y} + 1)$  and the associated projection operator  $\Theta(\Sigma)$ . Particle and antiparticle states are eigenfunctions of  $\Theta(\Sigma)$  with eigenvalues 1 and 0, respectively, so that we can resolve any wave function into a particle component and an antiparticle component according to the prescription (57). However, the Hilbert space of physical states  $\mathcal{H}_{phys}$  is not identified with  $\mathcal{H}_y$ , but is rather constructed by combining the particle amplitude with the complex conjugate of the antiparticle amplitude. [See (62)–(64).]

The evolution operator  $P_\lambda$  must satisfy the conditions (118) in order that (133) be self-consistent. For a free particle (88) gives a complete orthonormal set of eigenfunctions

of  $P_\lambda$ , which correspond to energy-momentum eigenstates in  $\mathcal{H}_{\text{phys}}$ . When an external electromagnetic field is present the theory is less complete. In this case (131) and (132) specify  $P_\lambda$  up to first order in the particle charge.

## B. Alternative hypotheses

There are a number of places in the development of the theory where a hypothesis different from the one adopted might have been made. Particular examples follow.

(1) The probability density for finding the particle might be taken as  $(|\psi_p|^2 + \psi_a)^2/y$  rather than as  $|\psi|^2/y$ . [See (65).] Note that these quantities are not the same, though their integrals over all  $\mathbf{y}$  are equal.

(2) Perhaps neither of these expressions for the probability density is valid and one should seek instead to find a set of position eigenfunctions  $\Psi_Y(\mathbf{y})$  which span  $\mathcal{H}_{\text{phys}}$ . These should be labeled by a past light cone vector  $Y^\lambda$  with  $Y^0 = -|\mathbf{Y}|$ , and satisfy the orthogonality and completeness relations

$$(\Psi_Y, \Psi_{Y'})_{\text{phys}} = Y\delta(\mathbf{Y} - \mathbf{Y}'), \quad (134)$$

$$\int \frac{d^3\mathbf{Y}}{Y} \Psi_Y(\mathbf{y}) \Psi_Y^\dagger(\mathbf{y}') = \langle \mathbf{y} | \Theta(\Sigma) | \mathbf{y}' \rangle \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}.$$

We can then define a position operator in  $\mathcal{H}_{\text{phys}}$  by

$$Y^\lambda_{\text{phys}} \Psi(\mathbf{y}) = \int \frac{d^3\mathbf{Y}}{Y} Y^\lambda \Psi_Y(\mathbf{y}) (\Psi_Y, \Psi)_{\text{phys}},$$

and interpret  $|(\Psi_Y, \Psi)_{\text{phys}}|^2 d^3\mathbf{Y}/Y$  as the probability of finding the light cone position three-vector  $\mathbf{y}$  in the range  $\mathbf{Y}$  to  $\mathbf{Y} + d\mathbf{Y}$ . Note that the operator  $y_{\text{OP}}$  of (11) acts in  $\mathcal{H}_y$ , not in  $\mathcal{H}_{\text{phys}}$ . A possible candidate for  $\Psi_Y(\mathbf{y})$ , which satisfies (134), is

$$\Psi_Y(\mathbf{y}) = \int dS_k \psi_k(\mathbf{y}) \begin{bmatrix} \psi_k^*(\mathbf{Y}) \\ \psi_k(\mathbf{Y}) \end{bmatrix},$$

$$\equiv \Theta(\Sigma) \begin{bmatrix} y\delta(\mathbf{y} - \mathbf{Y}) \\ \frac{\kappa^2}{4\pi} J_0([2\kappa^2 y_\lambda Y^\lambda]^{1/2}) \end{bmatrix}.$$

[See (58), (60), (70), (91), (92).] With this choice, the effective coordinate space wave function is

$$\psi_{\text{phys}}(\mathbf{Y}) = (\Psi_Y, \psi)_{\text{phys}} = \psi_p(\mathbf{Y}) + W\psi_a(\mathbf{Y}), \quad (135)$$

where  $W$  is the unitary operator defined by (93).

(3) One might identify  $\mathcal{H}_y$  directly with the space of physical states  $\mathcal{H}_{\text{phys}}$ . In this case the momentum operator would be  $\Theta(\Sigma)P_\lambda$ , the charge operator would be  $\text{sgn}(\Sigma)$ , and  $y_{\text{OP}}$  given by (11) would be acceptable as the light cone position operator. Such a theory would in fact be equivalent with that of (2) above, with the correspondence  $\psi \rightarrow \psi_{\text{phys}}$  [see (135)]. This theory would be more akin to that of the nonrelativistic Schrödinger equation rather than that of the Dirac and Klein-Gordon equations.

(4) We could make a choice of phase function  $g(\sigma)$  other than that of (88). It would be seen from the arguments of Secs. III E and III F that any  $g(\sigma)$  which behaves asymptotically like  $-\sigma \log \sigma + \sigma + O(1)$  for large  $\sigma$  will lead to

the correct nonrelativistic and to the correct classical limits. An example of such a  $g(\sigma)$  is given by (B5) in Appendix B. The price to pay if (88) is not adopted is that the operator  $W$  effecting the symmetry (91) no longer takes a simple form.

## C. Some unsolved problems

This paper leaves many questions unanswered, particularly those relating to the measurement of observables. Suppose the quantum state of an observer at the here-now  $z^\lambda$  is  $|\Psi\rangle$ , and he measures the observable represented by a Hermitian operator  $\mathcal{L}$  which has a complete orthonormal set of eigenstates  $|\Phi_n\rangle$  with corresponding eigenvalues  $\lambda_n$ . One would like to postulate that  $|\langle \Phi_n | \Psi \rangle|^2$  is the probability that the measurement will yield  $\lambda_n$ , and that if this eigenvalue is obtained, then immediately after the measurement the state will be  $|\Phi_n\rangle$ . However, how are we to interpret the phrase "a measurement at here-now  $z^\lambda$ ?" If the measurement occupies negligible extension in space and time so that it can be carried out by a single observer at  $z^\lambda$  then there is no difficulty. Otherwise the observer at  $z^\lambda$  must collaborate with auxiliary observers in other parts of space time. Various possibilities suggest themselves: (1) a prearranged experiment, carried out by a number of observers on the past light cone with vertex at  $z^\lambda$ ; (2) an unmediated experiment, where the observer at  $z^\lambda$  radios instructions to his collaborators to measure some quantity in their locality, in which case the experiment involves a number of observers on the future light cone of  $z^\lambda$ ; and (3) independent measurements made by different observers on a spacelike three-surface through  $z^\lambda$ , the results being communicated to the central observer.

Type (1) experiments comply most closely with the spirit of the present formalism, and have the advantage that all the data from the auxiliary observers reaches the central observer simultaneously at  $z^\lambda$ . However, the experiment can only involve a finite region of the past light cone and the central observer must plan the experiment sufficiently in advance to enable him to instruct his collaborators.

With experiments of type (2) and (3) the experimental results from the auxiliary observers will reach the central observer at different times, all in the future of  $z^\lambda$ . As each new piece of data arrives the central observer should presumably update his quantum state, but we lack a prescription for doing this. It seems likely that we shall need a description in terms of density matrices rather than of pure states to resolve this problem. That a solution to the problem of measurement in relativistic quantum theory may require the introduction of density matrices has been suggested by Houtappel, Van Dam, and Wigner.<sup>10</sup>

## APPENDIX A: ORTHOGONALITY AND COMPLETENESS RELATIONS

*Proof of (37)–(39):* Since  $k^\lambda$  and  $k^{\lambda'}$  both lie on the mass shell (14),  $k^\lambda + k^{\lambda'}$  is future pointing timelike and there exists a "center of momentum" inertial frame in which this sum vector has no spatial component, i.e.,  $k^\lambda = (\epsilon_k, \mathbf{k})$  and  $k^{\lambda'} = (\epsilon_k, -\mathbf{k})$ . In this frame (37) becomes



$$\begin{aligned}
(u_k, \Sigma u_{k'})_y &= \frac{1}{2}(u_k, \Sigma u_{k'})_y + \frac{1}{2}(\Sigma u_k, u_{k'})_y \\
&= (2\pi)^{-3} \int \epsilon_k e^{-2ik \cdot y} d^3y \\
&= \epsilon_k \delta(2\mathbf{k}) = \epsilon_k \delta(\mathbf{k} - \mathbf{k}'). \tag{A1}
\end{aligned}$$

Since each side of (A1) is an SO(1,3) scalar, the truth of this equation in any particular inertial frame implies its truth in all inertial frames, and hence (37) has been proved. Similar proofs may be given for (38) and (39).

*Proof of (40):* The left-hand side of (40) takes the form

$$\begin{aligned}
\Sigma(2\pi)^{-3/2} \int dS_k [e^{-ik^\lambda(y_\lambda - y'_\lambda)} - e^{ik^\lambda(y_\lambda - y'_\lambda)}] \\
= 2i\Sigma\Delta(y_\lambda - y'_\lambda). \tag{A2}
\end{aligned}$$

The function  $\Delta$  occurs in quantum field theory in connection with spin-zero field commutators, and may be expressed in terms of Bessel functions.<sup>11</sup> Putting in the explicit form for  $\Delta$  and applying the differential operator  $\Sigma$  then yields (40).

*Proof of (50) and (51):* The Beltrami operator belonging to the mass shell (14) is proportional to the operator

$$\begin{aligned}
\mathcal{H} &= \frac{1}{2} \tilde{\eta}^{-2} \eta_{\lambda\nu} \eta_{\mu\rho} \hat{J}_k^{\mu\nu} \hat{J}_k^{\rho\nu} \\
&= \left( \mathbf{k} \times \frac{\partial}{\partial \mathbf{k}} \right)^2 - \left( \epsilon_k \frac{\partial}{\partial \mathbf{k}} \right)^2. \tag{A3}
\end{aligned}$$

[See (22) and (23).] Direct computation shows that  $(-k^\lambda y_\lambda)^{i\sigma-1}$  is an eigenfunction of  $\mathcal{H}$  with eigenvalues  $1 + \sigma^2$ , if  $y^\lambda$  is any vector lying on the past light cone. Now  $\mathcal{H}$  is Hermitian with respect to the measure  $dS_k = d^3\mathbf{k}/\epsilon_k$ , which implies orthogonality for eigenfunctions belonging to different values of  $\sigma^2$ . Hence

$$\begin{aligned}
\int dS_k (-k^\lambda y_\lambda)^{i\sigma-1} (-k^\mu y'_\mu)^{-i\sigma'-1} \\
= C_1 \delta(\sigma - \sigma') + C_2 \delta(\sigma + \sigma'), \tag{A4}
\end{aligned}$$

where  $C_1$  and  $C_2$  are functions of  $\sigma$ ,  $y_\lambda$ , and  $y'_\lambda$  which are SO(1,3) scalars. A scalar function must yield zero when operated upon by any component of the appropriate total angular momentum tensor. Thus

$$\begin{aligned}
\left( \mathbf{y} \times \frac{\partial}{\partial \mathbf{y}} + \mathbf{y}' \times \frac{\partial}{\partial \mathbf{y}'} \right) C_s = 0, \\
\left( y \frac{\partial}{\partial y} + y' \frac{\partial}{\partial y'} \right) C_s = 0, \tag{A5}
\end{aligned}$$

with  $s = 1, 2$ . [See (76).] The general solution of (A5) is readily found to be

$$C_s = \alpha_s(\sigma, y^\lambda y'_\lambda) + (yy')^{-1} \beta_s(\sigma, y/y') \delta(\hat{\mathbf{y}}, \hat{\mathbf{y}}'), \tag{A6}$$

where  $\alpha_s$  and  $\beta_s$  are arbitrary functions of their arguments, and  $\delta(\hat{\mathbf{y}}, \hat{\mathbf{y}}')$  is the Dirac delta function for the surface of a unit sphere. Multiplying (A4) by  $y^{-i\sigma+1} (y')^{i\sigma'+1}$  yields

$$\begin{aligned}
\int dS_k (\epsilon_k + \mathbf{k} \cdot \hat{\mathbf{y}})^{i\sigma-1} (\epsilon_k + \mathbf{k} \cdot \hat{\mathbf{y}}')^{-i\sigma'-1} \\
= yy' (y/y')^{-i\sigma} C_1 \delta(\sigma - \sigma') \\
+ (yy')^{-i\sigma+1} C_2 \delta(\sigma + \sigma'). \tag{A7}
\end{aligned}$$

The left-hand side of (A7) is independent of  $y$  and  $y'$ , and hence so must be the right-hand side. The only way in which this independence can be consistent with (A6) is for

$$\begin{aligned}
\alpha_1 = 0, \quad \beta_1 = \beta(y/y')^{i\sigma}, \\
\alpha_2 = \alpha(y^\lambda y'_\lambda)^{i\sigma-1}, \quad \beta_2 = 0, \tag{A8}
\end{aligned}$$

where  $\alpha$  and  $\beta$  are functions of  $\sigma$  only. To find these functions, apply the integration  $\int d^2\hat{\mathbf{y}}'$  over the unit sphere, to both sides of (A7). The left-hand side yields

$$8\pi^3 (\sigma\sigma')^{-1} \kappa^{i\sigma} \kappa^{-i\sigma'} [\delta(\sigma - \sigma') - \delta(\sigma + \sigma')],$$

and the right-hand side

$$\beta \delta(\sigma - \sigma') - 2\pi i \alpha \sigma^{-1} 2^{i\sigma} \delta(\sigma + \sigma').$$

Identification of the coefficients of the two delta functions then gives  $\beta = 8\pi^3 \sigma^{-2}$ ,  $\alpha = 4\pi^2 i \sigma^{-1} (\kappa^2/2)^{i\sigma}$ . Thus

$$\begin{aligned}
\int dS_k (\epsilon_k + \mathbf{k} \cdot \hat{\mathbf{y}})^{i\sigma-1} (\epsilon_k + \mathbf{k} \cdot \hat{\mathbf{y}}')^{-i\sigma'-1} \\
= 8\pi^3 \sigma^{-2} \delta(\hat{\mathbf{y}}, \hat{\mathbf{y}}') \delta(\sigma - \sigma') \\
+ 4\pi^2 i \sigma^{-1} (\kappa^2/2)^{i\sigma} (1 - \hat{\mathbf{y}} \cdot \hat{\mathbf{y}}')^{i\sigma-1} \delta(\sigma + \sigma'). \tag{A9}
\end{aligned}$$

Multiplying (A9) by  $\sigma\sigma'/(16\pi^3)$  and making the identification  $\mathbf{w} = \hat{\mathbf{y}}$ ,  $\mathbf{w}' = \hat{\mathbf{y}}'$  then yields (50) and (51).

## APPENDIX B: COMPLETE ORTHONORMAL SETS IN $\mathcal{H}_y$

In Sec. III C we considered function sets of the form

$$\psi_k(\mathbf{y}) = (2\pi)^{-2} i \int_0^\infty d\sigma \sigma e^{ig(\sigma)} \xi^{i\sigma-1}, \tag{B1}$$

where  $\xi = -k_\lambda y^\lambda$ . These functions belong to the eigenvalue  $+1$  of the operator  $\Theta(\Sigma)$  defined by (58). It was shown that for any choice of the real function  $g(\sigma)$  the functions  $\psi_k(\mathbf{y})$  and  $\psi_k^*(\mathbf{y})$  form a complete orthonormal set in  $\mathcal{H}_y$ :

$$\begin{aligned}
(\psi_k, \psi_{k'})_y &= (\psi_k^*, \psi_{k'}^*)_y = \epsilon_k \delta(\mathbf{k} - \mathbf{k}'), \\
(\psi_k^*, \psi_{k'})_y &= 0, \tag{B2}
\end{aligned}$$

$$\int dS_k [\psi_k(\mathbf{y}) \psi_k^*(\mathbf{y}') + \psi_k^*(\mathbf{y}) \psi_k(\mathbf{y}')] = y \delta(\mathbf{y} - \mathbf{y}').$$

We exhibit here the form taken by  $\psi_k(\mathbf{y})$  for three particular choices for  $g(\sigma)$  (the case  $g(\sigma) = \arg(\Gamma(-i\sigma))$  was discussed in Sec. III D):

$$\begin{aligned}
e^{ig(\sigma)} &= 1, \\
\psi_k(\mathbf{y}) &= [4\pi^2 i \xi (\log \xi + i\epsilon)^2]^{-1}, \quad \epsilon \rightarrow 0, \tag{B3}
\end{aligned}$$

$$= \Theta(\Sigma) (2\pi)^{-1} \frac{d}{d\xi} \delta(\xi - 1);$$

$$\begin{aligned}
e^{ig(\sigma)} &= -\Gamma(-i\sigma)/\Gamma(i\sigma), \\
\psi_k(\mathbf{y}) &= \Theta(\Sigma) (2\pi)^{-1} J_0(2\xi^{1/2}); \tag{B4}
\end{aligned}$$

$$\begin{aligned}
e^{ig(\sigma)} &= \pi^{-1/2} \Gamma(\frac{1}{2} - i\sigma) [\cosh(\frac{1}{2}\pi\sigma) - i \sinh(\frac{1}{2}\pi\sigma)], \\
\psi_k(\mathbf{y}) &= \Theta(\Sigma) (2\pi^3)^{-1/2} \frac{d}{d\xi} (\xi^{1/2} \sin \xi). \tag{B5}
\end{aligned}$$

If we drop the requirement that our functions be eigenfunctions of  $\Theta(\Sigma)$ , then we can find further orthonormal, complete sets from a generalization of the theorem given by (82). The more general result is that any set of functions  $v_{kq}(\mathbf{y})$ ,  $q = 1$  or  $2$ , of the form

$$v_{k_1}(\mathbf{y}) = (2\pi)^{-2} i \int_{-\infty}^{\infty} d\sigma v(\sigma) \xi^{i\sigma-1},$$

$$v_{k_2}(\mathbf{y}) = \left[ (2\pi)^{-2} i \int_{-\infty}^{\infty} d\sigma \operatorname{sgn}(\sigma) v(\sigma) \xi^{i\sigma-1} \right]^*, \quad (\text{B6})$$

where  $v(\sigma)$  satisfies

$$|v(\sigma)|^2 + |v(-\sigma)|^2 = \sigma^2,$$

are orthonormal and complete

$$(v_{k_q}, v_{k'_q})_{\mathbf{y}} = \epsilon_k \delta(\mathbf{k} - \mathbf{k}') \delta_{qq'},$$

$$\sum_q \int dS_k v_{k_q}(\mathbf{y}) v_{k_q}^*(\mathbf{y}') = \mathbf{y} \delta(\mathbf{y} - \mathbf{y}').$$

This theorem is proved by techniques similar to those used to establish (82). Two examples of sets of type (B6) are furnished by (B7) and (B8) below:

$$v(\sigma) = (2\pi)^{-1/2} \sigma \Gamma(\frac{1}{2} - i\sigma) e^{1/2\pi(\sigma + (1/2)i)}, \quad (\text{B7})$$

$$v_{k_1}(\mathbf{y}) = (2\pi)^{-3/2} \frac{d}{d\xi} (\xi^{1/2} e^{i\xi});$$

$$v(\sigma) = (2\pi)^{-1/2} \sigma \frac{\Gamma(-i\sigma)}{\Gamma(i\sigma)} \Gamma(\frac{1}{2} + i\sigma) e^{1/2\pi(\sigma - i)},$$

$$v_{k_1}(\mathbf{y}) = 2^{-5/2} \pi^{-1} M(\frac{3}{2}, 1, i\xi), \quad (\text{B8})$$

$$= 2^{-5/2} \pi^{-1} [(1 + i\xi) J_0(\frac{1}{2}\xi) - \xi J_1(\frac{1}{2}\xi)] e^{(1/2)i\xi}.$$

In (B8)  $M$  denotes Kummer's confluent hypergeometric function. In neither of the above examples is the function  $v_{k_2}(\mathbf{y})$  expressible in elementary form.

### APPENDIX C: PROOF OF (94)

Because (94) has manifest Lorentz invariance, it is sufficient to prove it in the inertial frame in which  $k^\lambda = (\kappa, 0, 0, 0)$ . In this frame

$$W\psi_k(\mathbf{y}) = \int Y dY d^2\hat{\mathbf{Y}} w(y_\lambda Y^\lambda) (2\pi)^{-2} i$$

$$\times \int_0^\infty d\sigma \sigma e^{ig(\sigma)} (\kappa Y)^{i\sigma-1}. \quad (\text{C1})$$

The function  $w(y_\lambda Y^\lambda)$  may be expressed as the integral transform

$$w(y_\lambda Y^\lambda) = \frac{i\kappa^2}{8\pi^2} \int_{-\infty}^{\infty} d\sigma'$$

$$\times \left( \frac{1}{2} \kappa^2 y_\lambda Y^\lambda \right)^{-i\sigma'-1} \sigma' \frac{\Gamma(i\sigma')}{\Gamma(-i\sigma')}. \quad (\text{C2})$$

Upon substituting (C2) into (C1) one obtains the product of the two integrals

$$\int dY Y^{i(\sigma-\sigma')-1} = 2\pi \delta(\sigma - \sigma'),$$

$$\int d^2\hat{\mathbf{Y}} \left[ \frac{1}{2} (1 - \hat{\mathbf{y}} \cdot \hat{\mathbf{Y}}) \right]^{-i\sigma'-1} = 4\pi i \sigma'^{-1}.$$

This leads to

$$W\psi_k(\mathbf{y}) = - (2\pi)^{-2} i \int_0^\infty d\sigma \sigma e^{ig(\sigma)}$$

$$\times \frac{\Gamma(i\sigma)}{\Gamma(-i\sigma)} (\kappa \mathbf{y})^{-i\sigma-1},$$

which is just (94) in this special inertial frame. Invoking Lorentz invariance, (94) holds generally.

<sup>1</sup>G. H. Derrick, *J. Math. Phys.* **28**, 64 (1987).

<sup>2</sup>P. A. M. Dirac, *Rev. Mod. Phys.* **21**, 392 (1949).

<sup>3</sup>*Alphabet conventions*: Greek lowercase letters  $\iota, \kappa, \lambda, \dots = 0, 1, 2, 3$ , with summation over repeated indices. *Metric tensor*:  $\eta_{\lambda\mu} = \text{diag}(1, -1, -1, -1)$ . *Conjugation operations*: A superscript  $*, T, \dagger$  applied to a quantity denotes, respectively, the complex conjugate, transpose, Hermitian conjugate. *Step function*:  $\Theta(x) = 1$  if  $x > 0$  and 0 if  $x < 0$ . *Sign function*:  $\operatorname{sgn}(x) = \Theta(x) - \Theta(-x)$ .

<sup>4</sup>C. Itzykson and J.-B. Zuber, *Quantum Field Theory* (McGraw-Hill, New York, 1980), p. 114.

<sup>5</sup>L. Cohen, *J. Math. Phys.* **11**, 3296 (1970); E. Kerner and W. Sutcliffe, *J. Math. Phys.* **11**, 391 (1970); F. J. Testa, *J. Math. Phys.* **12**, 1471 (1971); I. W. Mayes and J. S. Dowker, *J. Math. Phys.* **14**, 434 (1973); M. M. Mizrahi, *J. Math. Phys.* **16**, 2201 (1975); J. S. Dowker, *J. Math. Phys.* **17**, 1873 (1976); A. C. Hirschfeld, *Phys. Lett. A* **67**, 5 (1978).

<sup>6</sup>E. P. Wigner, *Group Theory and its Application to the Quantum Mechanics of Atomic Spectra* (Academic, New York, 1959), Chap. 11.

<sup>7</sup>Note that if the mass  $m$  is zero, Lemma 1 and the derived theorem are false, because in this case there are an infinite number of unsymmetric  $\text{SO}(1,3)$  scalars derived from  $k^\lambda$  and  $k^{\lambda'}$  of the form  $(kk')^{-1} \beta(k/k') \delta(\hat{\mathbf{k}}, \hat{\mathbf{k}}')$ , with  $\beta$  an arbitrary function, cf. (A6).

<sup>8</sup>E. P. Wigner, *Ann. Math.* **40**, 149 (1939); V. Bargmann and E. P. Wigner, *Proc. Natl. Acad. Sci. USA* **34**, 211 (1948); L. L. Foldy, *Phys. Rev.* **102**, 568 (1956).

<sup>9</sup>L. D. Landau and E. M. Lifshitz, *Quantum Mechanics* (Pergamon, Oxford, 1965), p. 20.

<sup>10</sup>R. M. F. Houtappel, H. Van Dam, and E. P. Wigner, *Rev. Mod. Phys.* **37**, 595 (1965).

<sup>11</sup>S. S. Schweber, *An Introduction to Relativistic Quantum Field Theory* (Harper and Row, New York, 1961), p. 177.

# Sufficient conditions for zero not to be an eigenvalue of the Schrödinger operator

A. G. Ramm

Mathematics Department, Kansas State University, Manhattan, Kansas 66506

(Received 9 September 1986; accepted for publication 4 February 1987)

It is proved that if  $H = -\nabla^2 + q(x) \geq 0$ ,  $\text{Im } q = 0$ ,  $|q(x)| \leq c(1 + |x|)^{-a}$ ,  $c = \text{const} > 0$ ,  $a > 2$ , then zero is not an eigenvalue of  $H$ . An example is given of  $H \geq 0$ , with zero a resonance (half-bound state) and  $q = q(|x|)$  compactly supported and integrable. An example of a potential  $q = O(r^{-2})$  is known, for which  $H \geq 0$  and zero is an eigenvalue. This shows that  $a > 2$  is the optimal condition for zero not to be an eigenvalue of  $H \geq 0$ . If the condition  $H \geq 0$  does not hold and  $H$  is an operator in  $L^2(\mathbb{R}^3)$ , then zero can be an eigenvalue even if  $q \in C_0^\infty$ . If  $H$  is an operator in  $L^2(\mathbb{R}^1)$  or in  $L^2(\mathbb{R}_+^1)$ ,  $\mathbb{R}_+^1 = [0, \infty)$ , then zero cannot be an eigenvalue of  $H$  provided that  $a > 2$ ; here conditions  $H \geq 0$  and  $\text{Im } q = 0$  can be dropped. Global estimates of the Green's function of  $H$  from below and above are given.

## I. INTRODUCTION

In Ref. 1 Newton asked the following question. Let  $H = -\nabla^2 + q(x)$ ,  $\text{Im } q(x) = 0$ ,  $x \in \mathbb{R}^3$ ,  $\int (1 + |x|) \times |q(x)| dx < \infty$ ,  $\int = \int_{\mathbb{R}^3}$ . Assume that

$$H \geq 0. \quad (1)$$

Can 0 be an eigenvalue of  $H$ ?

The assumption about the decay of  $q(x)$  can be relaxed: It is sufficient to assume that for all sufficiently large  $x$  the estimate

$$|q(x)| \leq c(1 + |x|)^{-a}, \quad a > 2, \quad c = \text{const} > 0 \quad (2)$$

holds and that  $q \in L_{\text{loc}}^2$ . By  $c$  we denote below various constants. The assumption (\*)  $\int (1 + |x|) |q(x)| dx < \infty$  means, roughly speaking, that  $a > 4$  in (2). On the other hand, (\*) allows local singularities in a neighborhood of infinity, which are excluded by (2). Since assumption (2) covers most, if not all, of the potentials of interest which decay at infinity faster than  $|x|^{-2}$ , we will use this assumption.

The question raised by Newton was discussed in Ref. 2, where it was proved that if  $a > 3$  then the answer is no. The assumptions on  $q$  in Ref. 2 were given in terms of weighted  $L^p$  spaces. Our argument is different from the one in Ref. 2.

The purpose of this paper is to give the exact value of  $a$  for which the answer is no. We prove that if  $a > 2$ , then the answer is no and if  $a \leq 2$ , then zero can be an eigenvalue of  $H \geq 0$ .

The second part of this statement is known: It is shown in Ref. 3 (p. 375) that for central potentials  $q = q(r)$ ,  $r = |x|$ , if  $H\psi = 0$ ,  $\psi = r^{-1}u(r)Y_l(x^0)$ , where  $Y_l$  is the spherical harmonic  $x^0 = xr^{-1}$ , then  $u = O(r^{-l})$  as  $r \rightarrow \infty$ . Therefore  $\psi \in L^2(\mathbb{R}^3)$  provided that  $l > \frac{1}{2}$ . This conclusion does not use assumption (1). The potential  $q(r)$  can be chosen in  $C_0^\infty$ .

If  $l = 0$ ,  $q = q(r)$ ,  $\text{Im } q = 0$ , is compactly supported and integrable in a neighborhood of the origin and (1) holds, then zero can be a half-bound state (a resonance) although it cannot be a bound state.

In the one-dimensional case, if condition (2) holds and  $u'' - q(r)u = 0$ ,  $r \geq 0$ ,  $u \in L^2[0, \infty)$ , (3)

then  $u = 0$ . This conclusion holds without assumption (1) and without assumption  $\text{Im } q = 0$ .

We prove that the Green's function  $G(x, y)$  of  $H$  under the assumptions  $q(x) \geq 0$  and (2) satisfies the global estimates

$$c|x - y|^{-1} \leq G(x, y) \leq (4\pi|x - y|)^{-1}, \quad c = \text{const} > 0. \quad (4)$$

The low energy scattering has been studied in Refs. 4 and 5.

## II. THE RESULTS

**Theorem 1:** If  $H \geq 0$  and  $a > 2$ , then zero is not an eigenvalue of  $H$ .

*Proof:* Let  $H = H_0 + q$ ,  $H_0 = -\nabla^2$ ,  $H\psi = 0$ ,  $\psi \in L^2(\mathbb{R}^3)$ . Then  $\psi > 0$ , being the ground state of  $H$ . Let us first assume that  $q(x) = -p(x)$  and  $p(x) \geq 0$ . If  $q(x) = q_+(x) - q_-(x)$ , where  $q_+ = \max(0, q(x))$ ,  $q_-(x) = \max(0, -q(x))$ , then we can use a similar argument taking  $p(x) = q_-(x)$  and  $H_0 = -\nabla^2 + q_+(x)$ . The key estimate for the Green's function of the operator  $-\nabla^2 + q_+$  that is needed for the proof is the estimate (4). This estimate is proved in Lemma 2.

*Step 1:* If  $q(x) = -p(x)$ ,  $p(x) \geq 0$ , then  $H\psi = 0$  can be written as

$$\psi = H_0^{-1}p\psi = \int g p \psi dy, \quad g = (4\pi|x - y|)^{-1}. \quad (5)$$

Since  $\psi > 0$  and  $p \geq 0$ , we have  $p\psi \geq 0$  and

$$\psi = (4\pi|x|)^{-1} \int p \psi dy + o(|x|^{-1}), \quad |x| \rightarrow \infty. \quad (6)$$

Estimate (6) is proved in Lemma 1.

If  $\psi \in L^2$  then (6) implies that  $\int p \psi dy = 0$ . Since  $p\psi \geq 0$ , this means that  $p\psi = 0$ . Therefore  $\nabla^2 \psi = 0$  and  $\psi \in L^2(\mathbb{R}^3)$ . Thus  $\int |\nabla \psi|^2 dx = 0$  and  $\psi = 0$ . The same conclusion follows from the equation  $p\psi = 0$  if  $p \neq 0$  on an open set. On this set  $\psi = 0$  and, by the unique continuation property for elliptic equations,  $\psi \equiv 0$ .

*Step 2:* If  $q = q_+ - p(x)$ ,  $H_0 = -\nabla^2 + q_+$ , and  $H_0 G$

$= \delta(x - y)$ , then  $G$  satisfies estimate (4) and (5) holds with  $G$  in place of  $g$ . Therefore

$$c|x|^{-1} \int p\psi dy + o(|x|^{-1}) \leq \psi \leq (4\pi|x|)^{-1} \int p\psi dy + o(|x|^{-1}), \quad |x| \rightarrow \infty. \quad (7)$$

If  $\psi \in L^2$ , then  $\int p\psi dy = 0$ , and the rest of the argument is the same. Theorem 1 is proved.

**Lemma 1:** If  $a > 2$  and  $\psi \in L^2(\mathbb{R}^3)$ , then (5) implies (6).

*Proof:* One has

$$|\psi| \leq 4\pi \left( \int |x-y|^{-2} |p|^2 dy \right)^{1/2} \left( \int |\psi|^2 dy \right)^{1/2} \leq c(1 + |x|)^{-1}. \quad (8)$$

Here and below  $c$  denotes various positive constants, the first integral in (8) was estimated with the help of inequality (2), and the inequality  $p > 0$  was not used. From (5) and (8) one obtains

$$\psi = (4\pi|x|)^{-1} \int_{|y| < \epsilon|x|} dy p\psi(1 + O(\epsilon)) + \int_{|y| > \epsilon|x|} (4\pi|x-y|)^{-1} p\psi dy =: J_1 + J_2. \quad (9)$$

Here  $0 < \epsilon = \epsilon(r) \rightarrow 0$ ,  $\epsilon(r)r \rightarrow \infty$  as  $r = |x| \rightarrow \infty$ . Since  $a > 2$  and (8) holds we have  $\int p\psi dy < \infty$ . Therefore  $J_1 = (4\pi|x|)^{-1} \int dy p\psi + o(|x|^{-1})$ . Let us show that  $J_2 = o(|x|^{-1})$  as  $|x| \rightarrow \infty$ . One has

$$J_2 \leq c \int_{\epsilon r}^{\infty} dt t^2 (1+t)^{-(a+1)} \int_{-1}^1 (r^2 + t^2 - 2rts)^{-1/2} ds \leq c \int_{\epsilon r}^{\infty} dt (1+t)^{-a+1} (2rt)^{-1} [r+t - |r-t|] \leq cr^{-1} \left[ \int_{\epsilon r}^r dt (1+t)^{-a+1} + \int_r^{\infty} dt (1+t)^{-a+1} \right] \leq cr^{-a+1} + o(r^{-1}). \quad (10)$$

Here  $s = \cos \theta$  and we used the spherical coordinates with the  $y_3$  axis along the vector  $x$ ,  $|y| + t$ . Lemma 1 is proved.

**Lemma 2:** If  $q \geq 0$  satisfies estimates (2), with  $a > 2$  for  $|x| \geq R$ , where  $R > 0$  is an arbitrary large number, and  $q \in L^2_{loc}$ , then the Green's function of  $H = -\nabla^2 + q(x)$  satisfies estimate (4).

*Proof:* The right inequality in (4) follows from the maximum principle. Indeed, first note that  $G > 0$ : If  $G(x_0, y_0) \leq 0$ , then, since  $G(x, y_0) \rightarrow +\infty$  as  $x \rightarrow y_0$  and  $G(x, y_0) \rightarrow 0$  as  $|x| \rightarrow \infty$ , the function  $G(x, y_0)$  attains a nonpositive minimum at a certain point  $\xi \neq y_0$ . Since  $G \neq \text{const}$ ,  $q \geq 0$ , and  $\nabla^2 G = qG$  in a neighborhood of  $\xi$ , we have a contradiction which proves that  $G(x, y) > 0$ . Since  $G = g - \int gqG dz \leq g$ , we conclude that  $g \geq G$ .

Let us prove the left inequality (4).

We have

$$G(x, y) = g(x, y) - \int g(x, z)q(z)G(z, y)dz. \quad (11)$$

Suppose that the left inequality (4) does not hold. Then there exist sequences  $x_n$  and  $y_n$  such that

$$|x_n - y_n|G(x_n, y_n) \leq n^{-1}, \quad n \rightarrow \infty. \quad (12)$$

We will show this is impossible, so that (4) holds with some  $c > 0$ . There are three cases to consider.

*Case 1:* There exist numbers  $m$  and  $\delta$  which do not depend on  $n$ , such that  $|x_n| \leq m$ ,  $|y_n| \leq m$ ,  $|x_n - y_n| \geq \delta > 0$ . In this case choose  $x_n \rightarrow x_0$ ,  $y_n \rightarrow y_0$  and pass to the limit in (12) to get  $G(x_0, y_0) = 0$ , which is a contradiction since  $G(x_0, y_0) > 0$ .

*Case 2:* There exists a number  $m$  such that  $|x_n| \leq m$ ,  $|y_n| \leq m$ , and  $|x_n - y_n| \rightarrow 0$ . Then choose  $x_n \rightarrow x_0$ ,  $y_n \rightarrow y_0$ , multiply (11) by  $|x - y|$ , set  $x = x_n$ ,  $y = y_n$ , and pass to the limit  $n \rightarrow \infty$  to get  $0 = (4\pi)^{-1}$ , which is a contradiction. Here we used  $\lim |x_n - y_n| \int g(x_n, z)q(z)G(z, y_n)dz = 0$ , which holds since the integral is bounded.

*Case 3:* Either  $|x_n| \rightarrow \infty$  or  $|y_n| \rightarrow \infty$ , or  $|x_n| \rightarrow \infty$  and  $|y_n| \rightarrow \infty$ . From (11) and (12) one obtains

$$0 = (4\pi)^{-1} - \lim_{n \rightarrow \infty} |x_n - y_n| \int g(x_n, z)q(z)G(z, y_n)dz. \quad (13)$$

If  $|x_n - y_n| \leq m$ , where  $m$  does not depend on  $n$ , then the limit in (13) is zero since  $g(x_n, z) = (4\pi|x_n - z|)^{-1} \rightarrow 0$  as  $|x_n| \rightarrow \infty$ , and  $G(z, y_n) \leq (4\pi|z - y_n|)^{-1} \rightarrow 0$  as  $|y_n| \rightarrow \infty$ . In this case Eq. (13) becomes  $0 = (4\pi)^{-1}$ , which is a contradiction.

If  $|x_n - y_n| \rightarrow \infty$ , then

$$\lim |x_n - y_n| \int g(x_n, z)q(z)G(z, y_n)dz \leq \lim c|x_n - y_n| \left[ \int_{|z| < R} \frac{dz}{|x_n - z||z - y_n|(1 + |z|)^a} + \int_{|z| > R} \frac{dz}{|x_n - z||y_n - z|(1 + |z|)^a} \right] \leq c \lim |x_n - y_n| \cdot O(|x_n - y_n|^{-a}) = 0.$$

So again Eq. (13) leads to a contradiction. Therefore (12) cannot hold and the left inequality (4) is proved (cf. Ref. 6, p. 314).

*Remark:* In the one-dimensional case if  $H\psi = -\psi'' + q(r)\psi = 0$ ,  $r \geq 0$ ,  $a > 2$ , and  $\psi \in L^2[R, \infty)$ , then  $\psi \equiv 0$ . Here  $R > 0$  is an arbitrary (large) fixed number. This conclusion holds without assumption  $H \geq 0$  and without assumption  $\text{Im } q = 0$ . Indeed, the differential equation implies

$$\psi = A + Br + \int_r^{\infty} (t-r)q\psi dt, \quad (14)$$

where  $A$  and  $B$  are constants. If  $\psi \in L^2[R, \infty)$  and  $a > 2$ , then

$$\left| \int_r^{\infty} (t-r)q\psi dt \right| \leq \left( \int_r^{\infty} t^2(1+t)^{-2a} dt \right)^{1/2} \cdot \left( \int_r^{\infty} |\psi|^2 dt \right)^{1/2} = o(r^{-a+1.5}), \quad \text{as } r \rightarrow +\infty.$$

Therefore (14) and  $\psi \in L^2[R, \infty)$  imply that  $A = B = 0$ . If  $A = B = 0$ , then (14) becomes a homogeneous Volterra equation for  $\psi$  and therefore  $\psi = 0$ .

*Remark 2:* (a) Although, as we stated in Introduction, it is shown in Ref. 3 that there are central potentials such that

$\nabla^2\psi - q\psi = 0$ ,  $\psi \in L^2(\mathbb{R}^3)$ , no specific examples are given in Ref. 3. We give such examples and show that  $q(r)$  can be chosen in  $C_0^\infty$ . The corresponding  $\psi = h(r)Y_1(x^0) \in L^2(\mathbb{R}^3)$ ,  $x^0 = xr^{-1}$ . The construction is simple. Take  $h(r) = r^{-2}$  for  $r \geq 1$ ,  $h(r) > 0$  for  $0 \leq r \leq 1$ ,  $h \in C^\infty$ . Define  $q = \psi^{-1}\nabla^2\psi$ . Then  $q \in C^\infty$  and  $q = 0$  for  $r \geq 1$  since  $\nabla^2 r^{-2} Y_1 = 0$  for  $r \geq 1$ . Clearly  $\psi \in L^2(\mathbb{R}^3)$ . A similar example is in Ref. 7.

(b) An explicit example of an integrable near  $r = 0$  compactly supported  $q(r)$ , such that  $H \geq 0$  and zero is a half-bound state can also be constructed. One defines  $\psi = r^{-1}$  for  $r \geq 1$ ,  $\psi = r^{-1}u(r)$  for  $r \leq 1$ , and chooses  $u(r)$  so that  $q(r) := \psi^{-1}\nabla^2\psi$  is integrable and  $H \geq 0$ . Note that  $q(r) = 0$  for  $r \geq 1$ . The desired  $u$  one can choose, for example, in the form  $u = r^\gamma(1 + \gamma - \gamma r)$ , with any  $0 < \gamma < (\sqrt{2} - 1)/2$ .

### III. CONCLUSIONS

In this paper we prove that  $H \geq 0$  does not have zero eigenvalue if (2) with  $a > 2$  holds and may have zero if  $q$  falls off as  $O(r^{-2})$ ,  $r \rightarrow \infty$ . An example of compactly supported

integrable potential is given for which  $H \geq 0$  and zero is a resonance.

### ACKNOWLEDGMENT

The author thanks Professor R. G. Newton and Professor O. L. Weaver for discussions.

This work was supported by ONR. It was written while the author was a research professor at the Mathematical Institute of Academia Sinica. The author is grateful to the Academia for hospitality.

<sup>1</sup>R. C. Newton, *J. Math. Phys.* **18**, 1353 (1977).

<sup>2</sup>M. Klaus and B. Simon, *Ann. Phys. (NY)* **130**, 251 (1980).

<sup>3</sup>R. C. Newton, *Scattering Theory of Waves and Particles* (Springer, New York, 1982).

<sup>4</sup>R. C. Newton, *J. Math. Phys.* **27**, 2720 (1986).

<sup>5</sup>A. G. Ramm, *J. Math. Phys.* **21**, 308 (1980); *Theory and Applications of Some New Classes of Integral Equations* (Springer, New York, 1980), Appendix 4.

<sup>6</sup>A. G. Ramm, *Scattering by Obstacles* (Reidel, Dordrecht, 1986).

<sup>7</sup>C. L. Dolph, B. McLeod, and D. Thoe, *J. Math. Anal. Appl.* **16**, 311 (1966).

# The Korteweg–de Vries hierarchy of isospectral transformations: Towards a general explicit expression

Avia Rosenhouse<sup>a)</sup> and Jacob Katriel

Department of Chemistry, Technion—Israel Institute of Technology, Haifa 32000, Israel

(Received 15 October 1985; accepted for publication 4 February 1987)

The structure of the Korteweg–de Vries hierarchy of evolution equations, generating isospectral transformations, is elucidated by means of a study of its recurrence relations. For the  $m$ th member of the KdV hierarchy, which can be written in the form  $V_t = -2A_{m+1,x}$ , where the  $A_i$  satisfy the recurrence relation  $A_{m+1,x} = VA_{m,x} + \frac{1}{2}A_m V_x - \frac{1}{4}A_{m,xxx}$ , it is shown that  $A_m$  is a homogeneous polynomial in  $\partial^i V / \partial x^i$ . A general combinatorial formula for the coefficients of all the monomials entering  $A_m$ , up to a set of constants determined by means of a recurrence relation, is derived.

## I. INTRODUCTION

The Korteweg–de Vries (KdV) hierarchy of isospectral transformations was introduced by Lax<sup>1</sup> and by Gardner *et al.*<sup>2</sup> It is very intimately related to the question of uniqueness of spectral inversion. The connection with the inverse scattering method has been particularly clearly studied with respect to the original KdV equation.<sup>3</sup>

The classical limit of the KdV hierarchy was recently discussed,<sup>4</sup> and it was shown that in this limit the hierarchy reduces to the first-order equation

$$V_t = f(V) \cdot V_x \quad (1)$$

in which  $f(V)$  is an arbitrary function of  $V$ . This equation generates an isoperiodic transformation of  $V_0(x) = V(x,0)$  into  $V(x,t)$ . Thus, to each  $f(V)$  corresponds some isoperiodic transformation. It was also shown in Ref. 4 that if the two isoperiodic potentials  $V(x,0)$  and  $V(x,1)$  are given, the form of  $f(V)$  generating the transformation between them via Eq. (1) can easily be written down.

The complete KdV hierarchy has so far not been studied extensively at all. In particular, the general form of an arbitrary isospectral transformation has not been explicitly derived.

In the present article we present an attempt to derive the explicit form of the higher-order members of the hierarchy. General explicit forms are obtained, up to some numerical coefficients for which a set of recurrence relations is derived. The form of the results enables the association of an isospectral transformation with any  $F(V)$ , although the inversion problem, concerning the determination of the form of  $F(V)$  generating a particular isospectral transformation, has not been solved.

A related problem, which has been worked out in considerable detail, concerns the determination of the infinite sequence of polynomial conservation laws of the KdV equation.<sup>5,6</sup> In view of the relation between these conservation laws and the members of the KdV hierarchy, one could consider the results, in particular in Ref. 6, as almost providing the explicit form of the KdV hierarchy. However, the eluci-

ation of the structure of the KdV hierarchy achieved in the present article is a prerequisite for the consideration of the general isospectral transformation as presented in the concluding section, i.e., in terms of an arbitrary  $F(V)$ .

## II. PRELIMINARY CONSIDERATIONS

The evolution equation

$$V_t = -2A_{m+1,x}, \quad (2)$$

where

$$A_{m+1,x} = V \cdot A_{m,x} + \frac{1}{2}A_m \cdot V_x - (\hbar^2/4)A_{m,xxx} \quad (3)$$

specifies an isospectral transformation of the one-dimensional Schrödinger equation  $H\psi = E\psi$ , where  $H = -\hbar^2(d^2/dx^2) + V$ . Here  $x$  is the dynamical coordinate,  $t$  is a parameter such that  $V = V(x,t)$ ,  $V_t = \partial V / \partial t$ ,  $V_x = \partial V / \partial x$ , and  $A_{m,x} = \partial A_m / \partial x$ . Here  $m$  is a running index,  $m = 0, 1, 2, \dots$ .

The fact that the transformation is isospectral means that for all eigenvalues we have  $\partial E / \partial t = 0$ , or equivalently  $\int_{-\infty}^{\infty} \psi^* V_t \psi dx = 0$ , where  $\psi$  is any one of the eigenfunctions. That  $\int_{-\infty}^{\infty} \psi^* A_{m+1,x} \psi dx = 0$  was shown in Refs. 1, 2, and 4. It follows immediately that an arbitrary linear combination  $\sum_m \alpha_m A_{m+1,x}$ , where  $\{\alpha_m\}$  is a set of constant coefficients, will also specify an isospectral transformation.

$A_m$  was implicitly presented in Refs. 1 and 2 for  $m = 0, 1, 2, 3$ , i.e.,

$$A_0 = -1, \quad A_1 = -V/2,$$

$$A_2 = -\frac{3}{8}V^2 + (\hbar^2/8)V_{xx},$$

$$A_3 = -\frac{5}{16}V^3 + \frac{5}{16}\hbar^2[V \cdot V_{xx} + \frac{1}{2}V_x^2] - (\hbar^4/32)V_{(4)},$$

resulting in the evolution equations

$$V_t = V_x, \quad V_t = \frac{3}{2}VV_x - (\hbar^2/4)V_{(3)},$$

$$V_t = \frac{15}{8}V^2V_x - \frac{5}{8}\hbar^2[V \cdot V_{(3)} + 2V_x V_{xx}] - (\hbar^4/32)V_{(5)},$$

where  $V_{(i)} = \partial^i V / \partial x^i$ .

To elucidate the explicit form of the general isospectral transformation we first consider the successive terms that can be generated using Eq. (3).

Writing

<sup>a)</sup> Based on a part of a thesis to be submitted by AR to the senate of the Technion—Israel Institute of Technology, in partial fulfillment of the requirements for the M.Sc. degree.

$$A_m = \sum_{i=0}^{m-1} (-1)^{i+1} \hbar^{2i} A_{i,m}, \quad (4)$$

we obtain

$$A_{i+1,m+1,x} = V \cdot A_{i+1,m,x} + \frac{1}{2} A_{i+1,m} V_x + \frac{1}{4} A_{i,m,xxx}. \quad (5)$$

The nature of this two-dimensional recurrence relation is illustrated in Fig. 1. It determines  $A_{i+1,m+1,x}$ ,  $0 \leq i < m-1$ , in terms of the two elements  $A_{i+1,m}$  and  $A_{i,m}$  or, if  $i = m-1$ , in terms of  $A_{i,m}$  only. It is therefore necessary to specify both  $A_{0,0}$  and the constants of integration entering upon evaluation of  $A_{i+1,m+1}$  from  $A_{i+1,m+1,x}$ . We shall set  $A_{0,0} = -1$ , and all the integration constants will be chosen to be zero. Although some further flexibility in the final expression for the general isospectral transformation could be incorporated by allowing the constants of integration to be arbitrary functions of  $t$ , this further flexibility, whose classical analog was discussed in Ref. 4, will not be explicitly retained.

The recurrence relation becomes particularly simple for  $A_{0,m}$ , obtaining the form

$$A_{0,m+1,x} = V \cdot A_{0,m,x} + \frac{1}{2} A_{0,m} \cdot V_x,$$

which is just the classical limit of the original recurrence relation, Eq. (3). This recurrence relation, which involves stepping along the horizontal sequence  $i = 0$  in Fig. 1, was solved in Ref. 4, where it was shown that

$$A_{0,m} = [(2m-1)!!/2^m \cdot m!] V^m. \quad (6)$$

Having obtained all the terms corresponding to  $i = 0$ , it is now a simple matter to proceed along the horizontal line corresponding to  $i = 1$  in Fig. 1. The terms along this line constitute the lowest-order quantal terms. Thus substituting Eq. (6) in Eq. (5) we obtain

$$\begin{aligned} A_{1,m+1,x} = & V \cdot A_{1,m,x} + \frac{1}{2} A_{1,m} \cdot V_x \\ & + [(2m-1)!!/2^m \cdot (m-1)!] \\ & \times [(m-1)(m-2)V^{m-3}V_{(1)}^3 \\ & + 3(m-1)V^{m-2}V_{(2)}V_{(1)} + V^{m-1} \cdot V_{(3)}]. \end{aligned}$$

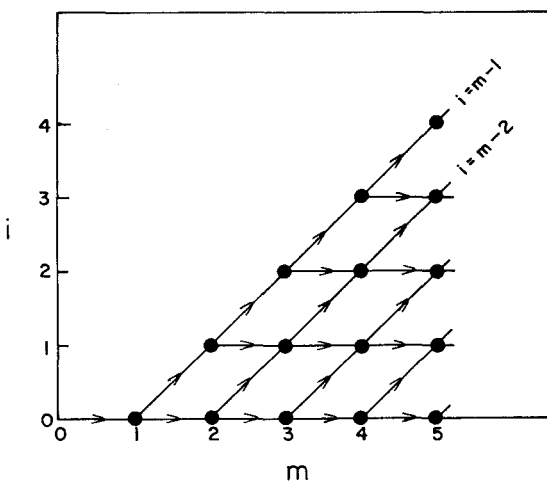


FIG. 1. The two-dimensional recurrence relation for  $A_{m,i}$ .

It follows straightforwardly that

$$A_{1,2} = \frac{1}{8} V_{(2)},$$

$$A_{1,3} = \frac{5}{16} (V \cdot V_{(2)} + \frac{1}{2} V_{(1)}^2),$$

$$A_{1,4} = \frac{35}{64} (V^2 \cdot V_{(2)} + V \cdot V_{(1)}^2),$$

⋮

and it can be shown by induction that

$$A_{1,m} = \frac{(2m-1)!!}{2^m \cdot 6} \left[ \frac{V^{m-2}}{(m-2)!} \cdot V_{(2)} + \frac{V^{m-3} \cdot V_{(1)}^2}{(m-3)! \cdot 2!} \right]. \quad (7)$$

This result can be substituted in Eq. (5) to obtain

$$A_{2,3} = \frac{1}{32} V_{(4)},$$

$$A_{2,4} = \frac{7}{64} V \cdot V_{(4)} + \frac{7}{32} V_{(3)} V_{(1)} + \frac{21}{128} V_{(2)}^2,$$

$$A_{2,5} = \frac{63}{256} V^2 V_{(4)} + \frac{63}{64} V V_{(3)} V_{(1)}$$

$$+ \frac{189}{256} V \cdot V_{(2)}^2 + \frac{231}{256} V_{(2)} V_{(1)}^2,$$

$$A_{2,6} = \frac{231}{512} V^3 V_{(4)} + \frac{693}{256} V^2 V_{(3)} V_{(1)} + \frac{2079}{1024} V^2 \cdot V_{(2)}^2$$

$$+ \frac{2541}{512} V V_{(2)} V_{(1)}^2 + \frac{1155}{2048} V_{(1)}^4,$$

⋮

On the basis of these results one can already make certain observations concerning the form of  $A_{i,m}$  in general. First,  $A_{i,m}$  is homogeneous, of degree  $m-i$ , with respect to  $V$  (including its derivatives). Furthermore,  $A_{i,m}$  is a sum of terms, each one of which contains a product of derivatives of  $V$ , the sum of whose orders is  $2i$ . This sum is called the derivative index in Ref. 6.

Thus the orders of the derivatives in each term of  $A_{i,m}$  constitute a partition of  $2i$ . These two statements can be established by induction, using Eq. (5). The induction has to be carried out over the two indices  $i$  and  $m$ , which can be done by imbedding an induction over  $m$  within an induction over  $i$ , as follows: checking Eq. (6) to establish that the two statements are true for all  $A_{0,m}$  (i.e.,  $i = 0$  and all  $m$ ) we shall assume that they are true for  $i$  and all  $m$  and show that they are, consequently, true for  $i+1$  and all  $m$ .

The last statement will be proved by induction over  $m$ , with fixed  $i$ : From Eq. (5) and the induction hypothesis it follows that the theorem holds for  $A_{i+1,i+2}$ . We assume that it holds for  $A_{i+1,m}$  and show that it holds for  $A_{i+1,m+1}$  by noting that each one of the three terms in the rhs of the recurrence relation, Eq. (5), is of order  $2i+3$  and degree  $m-i$  from which it follows that  $A_{i+1,m+1}$  is of order  $2(i+1)$  and degree  $(m+1) - (i+1) = m-i$ . This completes the proof.

Denoting by  $\xi$  the number of summands in a given partition of  $2i$ , we observe that the ratios between the coefficients

of terms corresponding to partitions of equal  $\xi$  are independent of  $m$ . Thus, for  $i = 2$ ,  $\xi = 2$  the possible partitions of  $2i = 4$  are  $3 + 1$  and  $2 + 2$ . The ratios of the corresponding terms are  $4/3$  in  $A_{2,4}$ ,  $A_{2,5}$  as well as in  $A_{2,6}$ . The general validity of this observation will be established later on. It suggests that the coefficient of each term in  $A_{i,m}$  can be written as a product of two factors, one of which depends on the partition but not on  $m$ , the other one depending on  $m$ ,  $i$ , and  $\xi$  but not on the specific partition.

The observations just made suggest that the general form of  $A_{2,m}$  is

$$A_{2,m} = \frac{(2m-1)!!}{2^m \cdot 180} \left\{ 3 \cdot \frac{V^{m-3} \cdot V_{(4)}}{(m-3)!} + \frac{V^{m-4}}{(m-4)!} \left[ 6V_{(3)} V_{(1)} + 9 \cdot \frac{V_{(2)}^2}{2!} \right] + 11 \cdot \frac{V^{m-5} \cdot V_{(2)} V_{(1)}^2}{(m-5)! \cdot 2!} + 15 \cdot \frac{V^{m-6} \cdot V_{(1)}^4}{(m-6)! \cdot 4!} \right\}. \quad (8)$$

This relation can be established by induction. One can proceed and obtain

$$A_{3,m} = \frac{(2m-1)!!}{2^m \cdot 2520} \left\{ 3 \cdot \frac{V^{m-4}}{(m-4)!} \cdot V_{(6)} + \frac{V^{m-5}}{(m-5)!} \left[ 9V_{(5)} V_{(1)} + 19 \cdot V_{(4)} V_{(2)} + \frac{23}{2!} V_{(3)}^2 \right] + \frac{V^{m-6}}{(m-6)!} \left[ \frac{25}{2!} V_{(4)} V_{(1)}^2 + 43V_{(3)} V_{(2)} V_{(1)} + \frac{61}{3!} V_{(2)}^3 \right] + \frac{V^{m-7}}{(m-7)!} \left[ \frac{60}{3!} V_{(3)} V_{(1)}^3 + \frac{83}{2! \cdot 2!} V_{(2)}^2 V_{(1)}^2 \right] + \frac{V^{m-8}}{(m-8)!} \cdot \frac{119}{4!} V_{(2)} V_{(1)}^4 + \frac{V^{m-9}}{(m-9)!} \cdot \frac{175}{6!} V_{(1)}^6 \right\}. \quad (9)$$

Further progress along these lines becomes rather cumbersome. However, before we attempt to undertake some more general considerations, let us obtain one further simple special case, along the line characterized by  $i = m - 1$ , in Fig. 1.

Along that line the recurrence relation obtains the form

$$A_{m,m+1,x} = \frac{1}{4} A_{m-1,m,xxx}$$

or

$$A_{m,m+1} = \frac{1}{4} A_{m-1,m,xx}$$

It follows immediately that

$$A_{m-1,m} = (1/2^{2m-1}) \cdot V_{(2m-2)}. \quad (10)$$

A comparison of this result with Eqs. (6)–(9) suggests that

$$A_{i,m} = \frac{(2m-1)!!}{3^i \cdot (2i+1)!! \cdot 2^{m+i}} \cdot \sum_{i_0, i_1, \dots, i_{2i}} C[1^{i_1} 2^{i_2} \dots (2i)^{i_{2i}}] \cdot \prod_{j=0}^{2i} \frac{V_{(j)}^{i_j}}{i_j!} \left( \sum_{j=0}^{2i} i_j = m-i, \sum_{j=1}^{2i} j \cdot i_j = 2i \right), \quad (11)$$

where  $C[1^{i_1} 2^{i_2} \dots (2i)^{i_{2i}}]$  are numerical coefficients which depend on the partition of  $2i$  but will be shown to be independent of  $m$ . These coefficients have already been determined, via Eqs. (6)–(9), for  $i = 0, 1, 2, 3$ . They are presented, along with further coefficients, in Table I.

For a given  $i$  and a large enough  $m$ , the number of terms in  $A_{i,m}$  is equal to the number of partitions of  $2i$  into sums of positive integers. Since  $i_j$  is the number of times that  $j$  appears in the partition, it follows that  $\xi = \sum_{j=1}^{2i} i_j$  is the number of summands in the partition considered. The maximum number of summands, obtained by writing  $2i = 1 + 1 + \dots + 1$  (i.e.,  $i_1 = 2i$ ;  $i_j = 0, j > 1$ ) is  $\xi_m = 2i$ . The power of  $V$  in a term corresponding to a partition into  $\xi$  summands is

$$i_0 = m - i - \sum_{j=1}^{2i} i_j = m - i - \xi.$$

Therefore, the minimum value of  $m$  for which all possible partitions of  $2i$  appear in  $A_{i,m}$  is  $m = i + \xi_m = 3i$ .

For  $i = 0$  the minimal value of  $m$  is zero. However, for  $i > 0$  the minimal value of  $\xi$  is 1, corresponding to  $i_j = 0, i \leq j < 2i; i_{2i} = 1$ . Therefore, the minimal value of  $m$ , corresponding to  $i_0 = 0$ , is  $i + 1$ . For  $i + 1 \leq m < 3i$  only partitions into at most  $m - i$  summands are allowed.

### III. THE GENERAL RECURRENCE RELATION FOR THE COEFFICIENTS

Having written  $A_{i,m}$  in Eq. (11), we shall now derive general recurrence relations for the coefficients  $C[1^{i_1} 2^{i_2} \dots (2i)^{i_{2i}}]$  which will enable us to show that these coefficients are indeed independent of  $m$ . The recurrence relations for the coefficients are obtained by substituting Eq. (11) in Eq. (5) and comparing coefficients. A typical term appearing in Eq. (5) after substitution of (11) is

$$g = \prod_{j=0}^n V_{(j)}^{i_j},$$

where

$$\sum_{j=0}^n i_j = m - 1, \quad \sum_{j=1}^n j \cdot i_j = 2i + 3, \quad n \leq 2i + 3.$$

To obtain the coefficient of  $g$  in Eq. (5) we note that it appears in the following ways.

#### A. From $A_{i+1,m+1,x}$

Here  $g$  can only arise upon differentiation of terms of the form  $h = gV_{(k)} / V_{(k+1)}$  which, being of degree  $m - i$  and order  $2(i + 1)$ , are contained in  $A_{i+1,m+1}$ . Such a term is only present if  $g$  contains  $V_{(k+1)}$  (i.e.,  $i_{k+1} > 0$ ). Upon differentiation it will generate many other terms in addition to  $g$ , but all these other terms are not relevant. The coefficient of the corresponding contribution to  $g$  will be

$$\alpha_{i+1,m+1}(h) \cdot (i_k + 1) \{k + 1\}.$$



TABLE I. The coefficients for  $i < 5$ .

$C[0]$	= 1	$C[8]$	= 81	$C[10]$	= 243	$C[2^2 3^2]$	= 2183 517/5
$C[2]$	= 3	$C[1,7]$	= 324	$C[1,9]$	= 1 215	$C[1^4,6]$	= 101 331
$C[1^2]$	= 3	$C[2,6]$	= 891	$C[2,8]$	= 4 131	$C[1^3,2,5]$	= 1223 262/5
$C[4]$	= 9	$C[3,5]$	= 1 539	$C[3,7]$	= 9 234	$C[1^3,3,4]$	= 370 332
$C[1,3]$	= 18	$C[4^2]$	= 1 863	$C[4,6]$	= 14 823	$C[1^2,2^2,4]$	= 2539 593/5
$C[2^2]$	= 27	$C[1^2,6]$	= 1 215	$C[5^2]$	= 17 253	$C[1^2,2,3^2]$	= 3192 777/5
$C[1^2,2]$	= 33	$C[1,2,5]$	= 2 916	$C[1^2,8]$	= 5 751	$C[1,2^3,3]$	= 4386 879/5
$C[1^4]$	= 45	$C[1,3,4]$	= 4 374	$C[1,2,7]$	= 17 658	$C[2^5]$	= 1209 411
$C[6]$	= 27	$C[2^2,4]$	= 6 075	$C[1,3,6]$	= 176 904/5	$C[1^5,5]$	= 363 285
$C[1,5]$	= 81	$C[2,3^2]$	= 7 533	$C[1,4,5]$	= 248 427/5	$C[1^4,2,4]$	= 754 515
$C[2,4]$	= 171	$C[1^3,5]$	= 4 131	$C[2^2,6]$	= 243 081/5	$C[1^4,3^2]$	= 949 887
$C[3^2]$	= 207	$C[1^2,2,4]$	= 42 687/5	$C[2,3,5]$	= 85 293	$C[1^3,2^2,3]$	= 6518 718/5
$C[1^2,4]$	= 225	$C[1^2,3^2]$	= 53 379/5	$C[2,4^2]$	= 514 593/5	$C[1^2,2^4]$	= 8960 139/5
$C[1,2,3]$	= 387	$C[1,2^2,3]$	= 73 548/5	$C[3^2,4]$	= 643 059/5	$C[1^6,4]$	= 1136 025
$C[2^2]$	= 549	$C[2^4]$	= 102 141/5	$C[1^3,7]$	= 25 272	$C[1^5,2,3]$	= 1964 655
$C[1^3,3]$	= 540	$C[1^4,4]$	= 12 393	$C[1^2,2,6]$	= 346 923/5	$C[1^4,2^3]$	= 13495 977/5
$C[1^2,2^2]$	= 747	$C[1^3,2,3]$	= 106 677/5	$C[1^2,3,5]$	= 610 983/5	$C[1^7,3]$	= 2993 760
$C[1^4,2]$	= 1071	$C[1^2,2^3]$	= 146 853/5	$C[1^2,4^2]$	= 147 177	$C[1^6,2^2]$	= 4113 747
$C[1^6]$	= 1575	$C[1^5,3]$	= 31 590	$C[1,2^2,5]$	= 837 621/5	$C[1^8,2]$	= 6330 555
		$C[1^4,2^2]$	= 216 999	$C[1,2,3,4]$	= 1265 949/5	$C[1^{10}]$	= 9823 275
		$C[1^6,2]$	= 65 205	$C[1,3^3]$	= 1588 734/5		
		$C[1^8]$	= 99 225	$C[2^3,4]$	= 1742 553/5		

The factor  $i_k + 1$  is due to the differentiation of  $V_{(k)}^{i_k+1}$  and the factor  $\{k + 1\} \equiv 1 - \delta_{i_k+1,0}$  takes care of the requirement  $i_{k+1} > 0$ . Here  $\alpha_{i+1,m+1}(h)$  is the coefficient of  $h$  in  $A_{i+1,m+1}$ , which, according to Eq. (11), can be written in the form

$$\alpha_{i+1,m+1}(h) = \frac{(2m+1)!!}{3^{i+1} \cdot (2i+3)!! \cdot 2^{m+i+2}} C[h] \left( \prod_j (i_j(h)!) \right)^{-1}.$$

### B. From $A_{i+1,m,x} \cdot V$

Here  $g/V$  is obtained (among other terms) upon differentiation of each one of the functions  $h' = (g/V) \cdot (V_{(k)}/V_{(k+1)})$ ,  $k = 0, 1, \dots, n-1$ . Since all these functions are contained in  $A_{i+1,m}$ , each one of them contributes the quantity

$$\alpha_{i+1,m}(h') \cdot (i_k + 1 - \delta_{k,0}) \{k + 1\}$$

to the coefficient of  $g$  in Eq. (5).

### C. From $\frac{1}{2} A_{i+1,m} \cdot V_{(1)}$

The term in  $A_{i+1,m}$  contributing to  $g$  will be  $h'' = g/V_{(1)}$ , with the coefficient

$$\frac{1}{2} \cdot \alpha_{i+1,m} \cdot \{1\}.$$

As indicated by the last factor, there is a contribution only if  $g$  contains  $V_{(1)}$ , i.e.,  $i_1 > 0$ .

### D. From $\frac{1}{4} A_{i,m,xxx}$

The general form of terms in  $A_{i,m}$  which contribute to  $g$  upon triple differentiation is

$$h''' = g \cdot (V_{(j)} V_{(k)} V_{(l)} / V_{(j+1)} V_{(k+1)} V_{(l+1)}),$$

$$0 \leq j < k < l \leq n-1.$$

The coefficient of  $g$  in the third derivative of  $h'''$  is specified for each one of the following cases.

#### 1. $j+1 < k < l-1$

$$6 \cdot (i_j + 1)(i_k + 1)(i_l + 1)\{j + 1\} \times \{k + 1\}\{l + 1\} \cdot \alpha_{i,m}(h''').$$

#### 2. $j+1 = k < l-1$

In this case  $h'''$  reduces to  $g \cdot V_{(j)} V_{(l)} / (V_{(j+2)} \times V_{(l+1)})$  and the coefficient of  $g$  in the third derivative is

$$3(i_j + 1)(2i_{j+1} + 1)(i_l + 1)\{j + 2\}\{l + 1\} \alpha_{i,m}(h''').$$

#### 3. $j = k < l-1$

Since  $h''' = g \cdot V_{(j)}^2 V_{(l)} / (V_{(j+1)}^2 V_{(l+1)})$  we obtain

$$3 \cdot (i_j + 2)(i_j + 1)(i_l + 1)\{\{j + 1\}\}\{l + 1\} \alpha_{i,m}(h'''),$$

where

$$\{\{\alpha\}\} \equiv 1 - \delta_{i_a,0} - \delta_{i_a,1} = \begin{cases} 0, & i_a = 0, 1, \\ 1, & i_a \geq 2. \end{cases}$$

#### 4. $j+1 < k = l-1$

$$3(i_j + 1)(i_k + 1)(2i_{k+1} + 1)\{j + 1\}\{k + 2\} \alpha_{i,m}(h''').$$

#### 5. $j+1 < k = l$

$$3(i_j + 1)(i_l + 2)(i_l + 1)\{j + 1\}\{\{l + 1\}\} \alpha_{i,m}(h''').$$

#### 6. $j+1 = k = l-1$

Here  $h''' = g \cdot V_{(j)} / V_{(j+3)}$  and the coefficient of  $g$  is

$$(i_j + 1) [6i_{j+1} \cdot i_{j+2} + 3(i_{j+1} + i_{j+2}) + 1] \times \{j + 3\} \alpha_{i,m}(h''').$$

**7.  $j+1=k=l$**

$$3(i_j + 1)(i_{j+1} + 1)^2 \cdot \{ \{ j + 2 \} \} \cdot \alpha_{i,m}(h^m).$$

**8.  $j=k=l-1$**

$$3(i_j + 1)(i_j + 2)i_{j+1} \{ j + 1 \} \{ j + 2 \} \alpha_{i,m}(h^m).$$

**9.  $j=k=l$**

$$(i_j + 3)(i_j + 2)(i_j + 1) \{ \{ j + 1 \} \} \alpha_{i,m}(h^m).$$

Using these results it is a straightforward, though tedious, task to equate the coefficients of any term  $g$  in the recurrence relation, Eq. (5).

The resulting equation is

$$\begin{aligned} & \sum_{k=0}^{n-1} \frac{i_{k+1}}{6(2i+3)} \left[ \left( m + \frac{1}{2} \right) \cdot C \left[ \frac{g(V) V_{(k)}}{V_{(k+1)}} \right] + (i + \xi - m) \cdot C \left[ \frac{g(V) \cdot V_{(k)}}{V \cdot V_{(k+1)}} \right] \right] - \frac{1}{2} \frac{i_1}{6(2i+3)} \cdot C \left[ \frac{g(v)}{V_{(1)}} \right] \\ &= \frac{1}{4} \left\{ \sum_{j=0}^{n-5} \sum_{k=j+2}^{n-3} \sum_{l=k+2}^{n-1} 6i_{j+1} i_{k+1} i_{l+1} \cdot C \left[ g(V) \cdot \frac{V_{(j)} V_{(k)} V_{(l)}}{V_{(j+1)} V_{(k+1)} V_{(l+1)}} \right] \right. \\ &+ \sum_{j=0}^{n-3} \sum_{l=j+2}^{n-1} 3i_{j+1} (i_{j+1} - 1) i_{l+1} \cdot C \left[ g(V) \cdot \frac{V_{(j)}^2 V_{(l)}}{V_{(j+1)}^2 V_{(l+1)}} \right] \\ &+ \sum_{j=0}^{n-3} \sum_{l=j+2}^{n-1} 3i_{j+1} i_{l+1} (i_{l+1} - 1) \cdot C \left[ g(V) \cdot \frac{V_{(j)} V_{(l)}^2}{V_{(j+1)} V_{(l+1)}^2} \right] \\ &+ \sum_{j=0}^{n-1} i_{j+1} (i_{j+1} - 1) (i_{j+1} - 2) \cdot C \left[ g(V) \cdot \frac{V_{(j)}^3}{V_{(j+1)}^3} \right] \\ &+ \sum_{j=0}^{n-4} \sum_{l=j+3}^{n-1} 3i_{j+2} (2i_{j+1} + 1) i_{l+1} \cdot C \left[ g(V) \frac{V_{(j)} V_{(l)}}{V_{(j+2)} V_{(l+1)}} \right] \\ &+ \sum_{j=0}^{n-4} \sum_{k=j+2}^{n-2} 3i_{j+1} i_{k+2} (2i_{k+1} + 1) \cdot C \left[ g(V) \frac{V_{(j)} V_{(k)}}{V_{(j+1)} V_{(k+2)}} \right] \\ &+ \sum_{j=0}^{n-2} 3i_{j+2} (i_{j+2} - 1) (i_{j+1} + 1) \cdot C \left[ g(V) \cdot \frac{V_{(j)} V_{(j+1)}}{V_{(j+2)}^2} \right] \\ &+ \sum_{j=0}^{n-2} 3i_{j+1}^2 i_{j+2} C \left[ g(V) \frac{V_{(j)}^2}{V_{(j+1)} V_{(j+2)}} \right] \\ &+ \left. \sum_{j=0}^{n-3} i_{j+3} (6i_{j+1} i_{j+2} + 3(i_{j+1} + i_{j+2}) + 1) \cdot C \left[ g(V) \frac{V_{(j)}}{V_{(j+3)}} \right] \right\}. \end{aligned} \tag{12}$$

In Appendix B we use Eq. (12) to show the  $m$  independence of the coefficients  $C[\dots]$ . Using this result the left-hand side of Eq. (12) obtains the form

$$\begin{aligned} & \sum_{k=0}^{n-1} \frac{i_{k+1}}{6(2i+3)} \left( i + \xi + \frac{1}{2} \right) C \left[ \frac{g V_{(k)}}{V_{(k+1)}} \right] \\ & - \frac{1}{2} \cdot \frac{i_1}{6(2i+3)} C \left[ \frac{g}{V_{(1)}} \right], \end{aligned} \tag{12'}$$

the right-hand side remaining the same as in Eq. (12).

The number of distinct terms appearing in the recurrence relation is larger than the number of terms in  $A_{i+1, m+1}$ , because upon differentiation of a typical term of the latter several contributions to the former appear. More precisely, the number of terms in  $A_{i+1, m+1}$  is the number of partitions of  $2(i+1)$ , whereas the number of terms in the recurrence relation is the number of partitions of  $2i+3$ . Thus the set of linear equations obtained upon equating coefficients in the recurrence relation is redundant. It is shown in Appendix A that one way of obtaining a nonredundant set of equations consists of considering only those terms in which the highest derivative appears linearly, i.e.,

$$g = V^{i_0} V_{(1)}^{i_1} \cdots V_{(n-1)}^{i_{n-1}} V_{(n)}.$$

**IV. RESULTS**

The general expression for  $A_{i,m}$  which we now write in the form

$$\begin{aligned} A_{i,m} &= \frac{(2m-1)!!}{3^i (2i+1)!! \cdot 2^{m+i}} \sum_{\xi \in \{i\}} \frac{V^{m-i-\xi}}{(m-i-\xi)!} \\ & \cdot \sum_{i_1, i_2, \dots, i_{2i}} C [1^{i_1} 2^{i_2} \cdots (2i)^{i_{2i}}] \prod_{j=1}^{2i} \left( \frac{V_{(j)}^{i_j}}{i_j!} \right) \\ & \left( \sum_{j=1}^{2i} i_j = \xi, \quad \sum_{j=1}^{2i} j \cdot i_j = 2i \right), \end{aligned}$$

where  $\{i\} = 1 - \delta_{i,0}$  involves a set of coefficients which have been shown in Appendix B to be independent of  $m$ . These coefficients satisfy the recurrence relation, Eq. (12'), which was used to obtain the coefficients appearing in Table I. The recurrence relation was further used to derive closed form expressions for certain types of coefficients, which are pre-

TABLE II. The coefficients for certain sets of terms.

$\xi = 1$	$[2i] = 3^i$
$\xi = 2$	$[1, 2i - 1] = 3^i \cdot i$ $[2, 2i - 2] = 3^{i-1} (2i^2 + 1)$ $[3, 2i - 3] = 3^{i-1} (i^3 - i^2/2 + i/2 - 1)$ $[4, 2i - 4] = (3^{i-1}/5)(2i^4 - 3i^3 + 2i^2 - 3i + 5)$
$\xi = 3$	$[1^2, 2i - 2] = 3^{i-1} (3i^2 - i + 1)$ $[1, 2, 2i - 3] = 3^{i-1} (2i^3 - 3i^2/2 + 3i/2 - 2)$ $[1, 3, 2i - 4] = (3^{i-1}/10)(10i^4 - 17i^3 + 13i^2 - 22i + 28)$ $[2^2, 2i - 4] = (3^{i-2}/5)(20i^4 - 32i^3 + 28i^2 - 52i + 63)$
$\xi = 4$	$[1^3, 2i - 3] = 3i(i^3 - i^2 + i - 1)$ $[1^2, 2, 2i - 4] = (3^{i-2}/5)(30i^4 - 55i^3 + 54i^2 - 95i + 99)$
$\xi = 5$	$[1^4, 2i - 4] = 3^{i-1} (3i^4 - 6i^3 + 7i^2 - 12i + 11)$
$\xi = 2i$	$[1^{2i}] = (2i + 1) [(2i - 1)!!]^2$

sented in Table II. These expressions provide some clues to the form of the general expression for an arbitrary coefficient, but the actual derivation of a closed form general expression has not been achieved.

### V. THE GENERAL FORM OF AN ISOSPECTRAL TRANSFORMATION

It was pointed out in Sec. II that the general form of an isospectral transformation is

$$V_t = \sum_m \alpha_m A_{m+1,x}, \quad (13)$$

where  $\{\alpha_m, m = 0, 1, \dots\}$  is an arbitrary set of constants. As a matter of fact, if the  $\alpha_m$  were arbitrary functions of  $t$  the above expression would still be a valid isospectral transformation.

Let

$$F(V) = \sum_m \frac{(2m-1)!!}{2^m} \alpha_m \cdot \frac{V^m}{m!}$$

and note that

$$F_{(k)}(V) \equiv \frac{\partial^k F}{\partial V^k} = \sum_m \frac{(2m-1)!!}{2^m} \alpha_m \cdot \frac{V^{m-k}}{(m-k)!}. \quad (14)$$

To obtain the general expression for an isospectral transformation characterized by an arbitrary  $F(V)$  we substitute Eqs. (4) and (11) in Eq. (13) and use Eq. (14) to express the sums over  $m$ . The resulting equation is

$$V_t = -2 \left[ \sum_{i=0}^{2i} \frac{(-1)^{i+1} \cdot \kappa^{2i}}{2^i \cdot 3^i \cdot (2i+1)!!} \cdot \sum_{\xi=\{i, i_1, i_2, \dots, i_{2i}\}} C [1^{i_1} \dots (2i)^{i_{2i}}] \times \prod_{j=1}^{2i} \left( \frac{V_{(j)}^{i_j}}{i_j!} \right) \cdot F_{(i+\xi)}(V) \right]_{x} \left( \sum_{j=1}^{2i} i_j = \xi, \sum_{j=1}^{2i} j \cdot i_j = 2i \right).$$

Comparison of the classical limit of this expression,  $V_t = 2F_x = 2F_{(1)} V_x$  with Eq. (1), indicates that  $f = 2F_{(1)}$ .

In a publication which appeared after the present article was submitted for publication, Torriani<sup>7</sup> presented a conjec-

ture enabling the combinatorial enumeration of the terms appearing in the KdV densities as well as the partial determination of their numerical coefficients. The relation suggested by these conjectures between the results derived in the present article and the representation theory of the symmetric group seem to deserve further attention.

### ACKNOWLEDGMENTS

Helpful discussions with Professor R. Pauncz are gratefully acknowledged. We are grateful to the referee for his suggestions.

This research was supported by the Fund for the Promotion of Research at Technion and the Technion V. P. R. Fund—Lawrence Deutsch Research Fund.

### APPENDIX A: NONREDUNDANT SET OF EQUATIONS FOR THE COEFFICIENTS

We shall now present a nonredundant set of equations for the coefficients appearing in  $A_{i+1,m}$ . The number of coefficients is equal to  $p_{2(i+1)}$ , the number of partitions of  $2(i+1)$ , and this will also be the number of equations to be presented. These equations will be ordered in such a way that each one of them contains one new coefficient, in addition to coefficients appearing in preceding equations.

The set of equations described will be obtained by equating the coefficient of a particular set of terms in Eq. (5), and in a particular order which we now specify.

Since each term  $g$  appearing in Eq. (5) is specified by a partition of  $2i+3$ , let us arrange these partitions in "increasing order" as in the following example, corresponding to  $2i+3=5$ :

$$\xi = 1 \quad (5)$$

$$\xi = 2 \quad (1) (4)$$

$$(2) (3)$$

$$\xi = 3 \quad (1) (1) (3)$$

$$(1) (2) (2)$$

$$\xi = 4 \quad (1) (1) (1) (2)$$

$$\xi = 5 \quad (1) (1) (1) (1) (1)$$

This ordering can either be specified as increasing in  $\xi$  and arranged dictionarywise for each  $\xi$ , or, if the partition is read as a number (in the basis  $2i+4$ ), the ordering is according to increasing numerical value ( $5 < 14 < 23 < 113 < \dots$  etc.). The same ordering was used in Ref. 6.

Let us first consider the relation obtained by equating the coefficients of

$$g = V^{m-i-\xi} V_{(1)}^{i_1} \dots V_{(n-1)}^{i_{n-1}} V_{(n)} \quad (A1)$$

in Eq. (5). One of the terms in  $A_{i+1,m+1}$  which contribute to  $g$  upon differentiation is

$$f_1(n-1) = V^{m-i-\xi} V_{(1)}^{i_1} \dots V_{(n-1)}^{i_{n-1}+1} = g \cdot V_{(n-1)} / V_{(n)}.$$

The other terms obtained by differentiation of  $f_1(n-1)$  are of the form

$$\tilde{g} = V^{m-i-\xi} \dots V_{(j)}^{i_j-1} V_{(j-1)}^{i_{j-1}+1} \dots V_{(n-1)}^{i_{n-1}+1}.$$

Since  $g$  precedes all possible  $\tilde{g}$  in the ordering specified above, it is obvious that  $f_1(n-1)$  could not have contribut-

ed to any one of the linear equations corresponding to the coefficients of terms preceding  $g$  in the above ordering. Other terms in  $A_{i+1,m+1}$  contributing to  $g$  upon differentiation, such as

$$f_1(k) = V^{m-i-\xi} V_{(1)}^{i_1} \cdots V_{(k)}^{i_k+1} V_{(k+1)}^{i_{k+1}-1} \cdots V_{(n)}$$

also generate additional terms which precede  $g$ , which means that the coefficient of  $f_1(k)$  has already appeared in a preceding equation. Since Eq. (5) is a two-dimensional recurrence relation we have to assume that when we attempt to determine  $A_{i+1,m+1}$  all the terms in the rhs, i.e.,  $A_{i+1,m}$  and  $A_{i,m}$ , are already available. Hence, the only new coefficient appearing is  $C[g \cdot V_{(n-1)} / V_{(n)}]$ . For a given  $1 \leq n \leq 2i+3$  the number of different  $g$ 's of the form (A1) is the number of partitions of  $2i+3-n$  into integers not larger than  $n-1$ ,  $P_{n-1}(2i+3-n)$ . In view of the identity

$$\sum_{n=2}^{2i+3} P_{n-1}(2i+3-n) = \sum_{n^*=1}^{2i+2} P_{n^*}(2i+2-n^*) = p(2i+2)$$

it is obvious that the number of equations generated by all  $g$  of the form (A1) is equal to the number of coefficients in  $A_{i+1,m+1}$ . Since each one of these equations, if ordered as specified above, contains one new coefficient, together they determine all  $p(2i+2)$  coefficients in  $A_{i+1,m+1}$ .

## APPENDIX B: $m$ INDEPENDENCE OF THE COEFFICIENTS

To demonstrate the fact that the coefficients are not  $m$  dependent, we start from the recurrence relation, Eq. (12), which we write in the form

$$\begin{aligned} & \sum_{k=0}^{n-1} \frac{i_{k+1} \cdot m}{6(2i+3)} \left[ C \left[ \frac{g V_{(k)}}{V_{(k+1)}} \right] - C \left[ \frac{g \cdot V_k}{(V \cdot V_{(k+1)})} \right] \right] \\ & + \sum_{k=0}^{n-1} \frac{i_{k+1}}{6(2i+3)} \left[ \frac{1}{2} C \left[ \frac{g \cdot V_{(k)}}{V_{(k+1)}} \right] \right. \\ & \left. + (i + \xi) C \left[ \frac{g \cdot V_{(k)}}{(V \cdot V_{(k+1)})} \right] \right] \\ & - \frac{1}{2} \frac{i_1}{6(2i+3)} C \left[ \frac{g}{V_{(1)}} \right] = \cdots \end{aligned} \quad (\text{B1})$$

All the terms not explicitly written down are neither explicit-

ly  $m$  dependent nor containing coefficients corresponding to partitions of  $2(i+1)$ .

It was shown in Appendix A that each equation introduces one new coefficient. Since for  $i=0,1$  the coefficients were explicitly shown to be  $m$  independent, we can invoke the following inductive argument to establish the  $m$  independence of all coefficients. Assuming that all coefficients preceding the last one are  $m$  independent we have in particular

$$C[g V_{(k)} / V_{(k+1)}] = C[g \cdot V_{(k)} / (V \cdot V_{(k+1)})], \quad k < n-1.$$

Thus, the only remaining  $m$  dependence is in

$$\begin{aligned} & [i_n \cdot m / 6(2i+3)] [C[g \cdot V_{(n-1)} / V_{(n)}] \\ & - C[g \cdot V_{(n-1)} / (V \cdot V_{(n)})]] \\ & + [i_n / 6(2i+3)] \left[ \frac{1}{2} C[g V_{(n-1)} / V_{(n)}] \right. \\ & \left. + (i + \xi) C[g \cdot V_{(n-1)} / (V \cdot V_{(n)})] \right] = \cdots \equiv \mu. \end{aligned} \quad (\text{B2})$$

The lowest value of  $m$  for which  $C[g V_{(n-1)} / V_{(n)}]$  can appear is  $m = i + \xi$ . For this value of  $m$   $C[g V_{(n-1)} / (V \cdot V_{(n)})]$  cannot appear so that Eq. (B2) results in

$$C[g \cdot V_{(n-1)} / V_{(n)}] = \mu \cdot 6(2i+3) / i_n \cdot (i + \xi + \frac{1}{2}).$$

Assuming that up to some  $m^*$   $C[V^{m^*} \cdot g \cdot V_{(n-1)} / V_{(n)}] = C[g V_{(n-1)} / V_{(n)}]$  we show that the same applies to  $C[V^{m^*+1} \cdot g \cdot V_{(n-1)} / V_{(n)}]$ . Note that the latter corresponds to  $m = i + \xi + m^* + 1$ , so that from (B2)

$$\begin{aligned} & C[V^{m^*+1} \cdot g \cdot V_{(n-1)} / V_{(n)}] \\ & = (\mu \cdot [6(2i+3) / i_n] + (m^* + 1) \\ & \cdot C[V^{m^*} \cdot g \cdot V_{(n-1)} / V_{(n)}]) / (m + \frac{1}{2}) \\ & = \mu \cdot 6(2i+3) / i_n (i + \xi + \frac{1}{2}) = C[g V_{(n-1)} / V_{(n)}]. \end{aligned}$$

This concludes the proof of the  $m$  independence of the coefficients.

<sup>1</sup>P. D. Law, Commun. Pure Appl. Math. **21**, 467 (1968).

<sup>2</sup>C. S. Gardner, J. M. Greene, M. D. Kruskal, and R. M. Miura, Commun. Pure Appl. Math. **27**, 97 (1974).

<sup>3</sup>G. L. Lamb, Jr., *Elements of Soliton Theory* (Wiley, New York, 1980).

<sup>4</sup>J. Katriel and A. Rosenhouse, Phys. Rev. D **32**, 884 (1985).

<sup>5</sup>R. M. Miura, C. S. Gardner, and M. D. Kruskal, J. Math. Phys. **9**, 1204 (1968).

<sup>6</sup>M. D. Kruskal, R. M. Miura, C. S. Gardner, and N. J. Zabusky, J. Math. Phys. **11**, 952 (1970).

<sup>7</sup>H. H. Torriani, Phys. Lett. A **113**, 345 (1986).

# Schrödinger-like equation for relativistic particles

I. B. Goldberg

*Racah Institute of Physics, Hebrew University, Jerusalem, 91904 Israel*

R. H. Pratt

*Department of Physics and Astronomy, University of Pittsburgh, Pittsburgh, Pennsylvania 15260*

(Received 17 November 1986; accepted for publication 11 February 1987)

The Dirac equation in a spherically symmetric screened Coulomb potential is transformed to a modified Schrödinger equation of the form  $d^2u/dr^2 + k^2(r)u = 0$ . This transformation is induced by expressing the Dirac function as a linear combination of the function  $u$  and its derivative  $du/dr$ . Various properties of the transformation and of the resulting equations are studied. The close similarity between the modified Schrödinger equation and the Schrödinger equation suggests that methods applied to the Schrödinger equation to derive nonrelativistic relations can be applied to the modified Schrödinger equation to derive the analogous relativistic relations. As an example, this approach is applied to the single channel quantum defect theory to give a new derivation of its relativistic form.

## I. INTRODUCTION

Our objective in this paper is to develop a Schrödinger-equation-like formalism for the Dirac equation in a screened Coulomb potential permitting relativistic calculations utilizing procedures previously developed for the nonrelativistic case. In various atomic processes, such as photoionization, Compton scattering, or bremsstrahlung in an ionic field, the potential seen by the free electron at large distances from the ion is a point Coulomb potential, corresponding to the ionic charge  $Z_{\text{ion}}$ . Closer to the ion, when the free electron is penetrating the charge distribution of the bound electrons, bound electron screening of the nuclear charge is no longer complete and the potential is no longer of the point Coulomb type; in addition exchange and correlation effects also affect the free electron. Following the terminology of the  $R$ -matrix method<sup>1</sup> we thus divide the space around the ion into two regions: the external region, where the potential is point Coulomb, and the internal, non-Coulomb region. (Of course at the center of the internal region is another Coulombic region, but now characterized by the nuclear charge  $Z$  rather than the ionic charge  $Z_{\text{ion}}$ .) Several theoretical models have been applied in order to calculate a local atomic potential in the internal region. The most commonly used are the Hartree-Fock-Slater<sup>2</sup> potential in the nonrelativistic case and the Dirac-Fock-Slater<sup>3</sup> potential in the relativistic case. Other methods (local-density approximation,<sup>4</sup> random phase approximation,<sup>5</sup> and Fock approximation<sup>6</sup>), which may go beyond a local potential description for the interior, are used as well.

Except for very special cases, there is no analytic solution to the Dirac equation for a screened central potential in the internal region, and numerical methods must be applied. On the other hand, in the external region (referred to hereafter also as the tail region) the solutions are well known.<sup>7</sup> The complete solution to the Dirac equation, describing both internal and external regions, may be obtained by matching numerically integrated functions in the internal region to an appropriate linear combination of the regular and irregular solutions in the point Coulomb potential. With this matching, carried out at some point in the tail region, one is able to

normalize the continuum wave function, and to calculate the phase shift caused by the presence of the potential in the interior region.

The analytic solution of the Dirac equation in the point Coulomb potential may be expressed in terms of confluent hypergeometric functions. These functions appear also in the Coulomb functions which provide the solution of the Schrödinger equation in the point Coulomb potential, and it turns out that the Dirac functions in the point Coulomb potential are given as linear combinations of such functions and their derivatives, albeit for modified arguments. Utilizing this linear combination of solutions, the Dirac equation is transformed into a modified Schrödinger equation<sup>8</sup> in the Coulomb potential, which we may call a modified Coulomb equation.

The question now arises as to whether there exists a similar transformation for the Dirac equation in a screened Coulomb potential, which would result in a modified Schrödinger equation of the form  $d^2u/dr^2 + k^2(r)u = 0$ , where the effective potential  $k^2(r)$  is a function to be determined. It turns out that there is such a transformation. The Dirac functions can be expressed as a linear combination of the function  $u$  and its derivative, where  $u$  is a solution of the modified Schrödinger equation.

The close similarity between the modified Schrödinger equation and the Schrödinger equation suggests that various relativistic quantities can be obtained with the same methods used to obtain the corresponding nonrelativistic quantities. This approach is applied here to the quantum defect theory, and application of the formalism to a relativistic WKB approximation is in progress.<sup>9</sup> We expect that other features and consequences of the Dirac equation can be obtained by applying nonrelativistic methods to the modified Schrödinger equation.

In Sec. II we will begin by reviewing the Dirac equation for the point Coulomb potential and how it may be written in terms of a modified Schrödinger equation (modified Coulomb equation). In Sec. III we demonstrate that the Dirac equation in a screened Coulomb potential may also be written in terms of a modified Schrödinger equation. In Sec. IV,

as an application of this formalism, we rederive the relativistic quantum defect theory relations.

## II. THE DIRAC EQUATION FOR THE POINT COULOMB POTENTIAL

### A. The Dirac equation

The Dirac equation for an electron moving in a spherically symmetric potential  $V(r)$  is

$$H\psi = (c\alpha \cdot p + \beta mc^2 + V)\psi = E\psi. \quad (2.1)$$

We choose the usual representation in which

$$\alpha = \begin{pmatrix} 0 & \sigma \\ \sigma & 0 \end{pmatrix}, \quad \beta = \begin{pmatrix} I & 0 \\ 0 & -I \end{pmatrix}. \quad (2.2)$$

$I$  is the  $2 \times 2$  unit matrix and  $\sigma$  are the usual  $2 \times 2$  Pauli matrices. Here  $E = mc^2 + T$  is the total energy of the electron,  $T$  being the energy relative to the rest mass.

The eigenfunction  $\psi$  is a four-rank spinor which can be written as

$$\psi_{\kappa m} = \frac{1}{r} \begin{pmatrix} g_{\kappa} \cdot \Omega_{\kappa m} \\ i f_{\kappa} \cdot \Omega_{-\kappa m} \end{pmatrix}, \quad (2.3)$$

where  $\kappa$  is a quantum number which combines angular momentum  $j$  and parity  $l$ ,

$$\kappa = \mp (j + \frac{1}{2}) \quad \text{as } j = l \pm \frac{1}{2}. \quad (2.4)$$

The spherical spin orbit spinor of rank 2 is given by

$$\Omega_{\kappa m} = \sum_{s=\pm 1/2} \left( l \frac{1}{2} m - s s |jm \right) Y_{l, m-s} X^s; \quad (2.5)$$

$(l \frac{1}{2} m - s s |jm)$  are the Clebsch-Gordan coefficients,  $Y_{l, m-s}$  are the spherical harmonics, and the spinors  $X^s$  are defined by

$$X^{1/2} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \quad X^{-1/2} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}. \quad (2.6)$$

With this form for the eigenfunction, the Dirac equation reduces to

$$\frac{dg}{dr} = -\frac{\kappa}{r} g + \left( \frac{1+\epsilon}{\lambda_c} - U \right) f, \quad (2.7a)$$

$$\frac{df}{dr} = \frac{\kappa}{r} f + \left( \frac{1-\epsilon}{\lambda_c} + U \right) g, \quad (2.7b)$$

where  $\lambda_c = \hbar/mc$  is the Compton wavelength,  $\epsilon = E/mc^2 = 1 + T/mc^2$ , and  $U = V/\hbar c$ .

### B. The modified Coulomb equation

When the potential is given by  $V = -Ze^2/r$ , the Dirac equation (2.7) takes the form

$$\frac{dg}{dr} = \frac{\kappa}{r} g + \left( \frac{1+\epsilon}{\lambda_c} + \frac{a}{r} \right) f, \quad (2.8a)$$

$$\frac{df}{dr} = \frac{\kappa}{r} f + \left( \frac{1-\epsilon}{\lambda_c} - \frac{a}{r} \right) g, \quad (2.8b)$$

where  $a = [(Ze^2)/(\hbar c)] = Z\alpha$ ,  $\alpha$  being the fine structure constant.

The solutions of (2.8) are well known<sup>7</sup>; they can be writ-

ten in terms of solutions of the Schrödinger equation in the same potential, but with modified parameters. Here we will deduce a form for this connection by only considering boundary conditions, to prepare for a generalization of the connection to other potentials. When the energy  $T$  is positive ( $\epsilon > 1$ ), the regular solution of (2.8) has the following properties.

(i) As  $r \rightarrow 0$  we have

$$g = Nr^\gamma, \quad f = [(\kappa + \gamma)/a]g. \quad (2.9)$$

Here  $\gamma = (\kappa^2 - a^2)^{1/2}$  and  $N$  is a constant which depends on the normalization of the continuum solution of the Dirac equation at large distances.

(ii) As  $r \rightarrow \infty$  we have

$$\begin{aligned} g &\sim A \sin(pr - \frac{1}{2}l\pi + \delta_c + n \ln|2pr|), \\ f &\sim AQ \cos(pr - \frac{1}{2}l\pi + \delta_c + n \ln|2pr|), \end{aligned} \quad (2.10)$$

where

$$\delta_c = -\arg \Gamma(\gamma + i\eta) + (l+1-\gamma)(\pi/2) + \xi, \quad (2.11a)$$

and

$$\begin{aligned} p &= (1/\lambda_c)(\epsilon^2 - 1)^{1/2}, \quad \eta = a\epsilon/(\epsilon^2 - 1)^{1/2}, \\ Q &= ((\epsilon - 1)/(\epsilon + 1))^{1/2}, \quad \tan \xi = aQ/(\kappa - \gamma). \end{aligned} \quad (2.11b)$$

The constant  $A$  in (2.10), like  $N$  in (2.9), is determined by the normalization of the continuum wave functions. When the normalization is on the energy scale, i.e.,

$$\int_0^\infty dr [g(r, E)g(r, E') + f(r, E)f(r, E')] = \delta(E - E'), \quad (2.12)$$

we have

$$A = [(mc^2 \cdot \pi \lambda_c \cdot Q)^{1/2}]^{-1}, \quad (2.13a)$$

$$\begin{aligned} N &= A \cdot [|\Gamma(\gamma + i\eta)| e^{\pi\eta/2} / \Gamma(2\gamma + 1)] \\ &\quad \times (2p\lambda_c)^\gamma (\gamma \cos \xi - \eta \sin \xi). \end{aligned} \quad (2.13b)$$

The boundary conditions (2.9) and (2.10) for the large component function  $g$  are similar to the boundary conditions of the regular Coulomb function  $u_c$ , which provide the continuum solution to the Schrödinger equation in the point Coulomb potential, but corresponding to a noninteger angular momentum parameter  $\gamma - 1$ . Namely,

$$u_c \underset{r \rightarrow 0}{\propto} r^\gamma, \quad (2.14)$$

$$u_c \underset{r \rightarrow \infty}{\propto} \sin(pr - \frac{1}{2}l\pi + \sigma_c + \eta \ln(|2pr|)),$$

where

$$\sigma_c = -\arg \Gamma(\gamma + i\eta) + (l+1-\gamma)(\pi/2) = \delta_c - \xi. \quad (2.15)$$

The combination  $v_c = (1/p)(u'_c - (\gamma/r)u_c)$  behaves like  $r^\gamma$  near the origin and like

$$\cos(pr - \frac{1}{2}l\pi + \sigma_c + \eta \ln|2pr|)$$

as  $r \rightarrow \infty$ . Thus the boundary conditions suggest the following form for the functions  $g$  and  $f$ :

$$g = \cos \xi \cdot u + \sin \xi \cdot (1/p) \cdot (u' - (\gamma/r)u), \quad (2.16a)$$

$$f = Q [ -\sin \xi \cdot u + \cos \xi \cdot (1/p) \cdot (u' - (\gamma/r)u) ]. \quad (2.16b)$$

Indeed, on substituting (2.16) in the Dirac equations we find that both equations, (2.8a) and (2.8b), reduce to the same modified Schrödinger equation in the point Coulomb potential,

$$\frac{d^2u}{dr^2} + \left( p^2 + \frac{2a\epsilon}{\lambda_c r} - \frac{\gamma(\gamma-1)}{r^2} \right) u = 0, \quad (2.17)$$

which we shall call the *modified Coulomb equation*. The coefficients of the first derivative  $du/dr$  in (2.17) vanish,

$$\begin{aligned} -(\kappa - \gamma) \tan \xi + aQ &= 0, \\ -(\kappa + \gamma) / (\tan \xi) + a/Q &= 0, \end{aligned} \quad (2.18)$$

since  $\gamma = (\kappa^2 - a^2)^{1/2}$  and  $\tan \xi = aQ / (\kappa - \gamma)$ , according to our previous definitions.

However, if we tried to solve (2.18) to determine  $\gamma$  and  $\xi$ , we would find a second solution which also leads to the modified Coulomb equation (2.17), namely  $\gamma = -(\kappa^2 - a^2)^{1/2}$ . The regular solution of (2.17) in this case behaves near the origin like  $r^{1-\gamma}$  and the phase  $\sigma_c$  would be obtained by replacing  $\gamma$  with  $1 - \gamma$  in Eq. (2.15). Either solution leads to the same solution of the Dirac equation.

It follows then that  $\gamma$  and  $\xi$ , which we have related initially to the properties of the Dirac functions, may alternately be considered as free parameters, to be determined by the requirements that (a)  $g$  and  $f$  given by (2.16) solve the Dirac equation (2.8) and (b)  $u$  is a solution of the modified Coulomb equation (2.17).

The arguments which led us to the form (2.16) for the Dirac wave functions were based on the boundary condition of the regular continuum functions. However, since we have shown that the Dirac equation (2.8) is equivalent to the modified Coulomb equation (2.17) when  $g$  and  $f$  have the form (2.16), it follows that this relation is not restricted to the regular continuum case. Thus, when  $u$  is an irregular solution of the modified Coulomb equation, we get from (2.16) an irregular solution of the Dirac equation. The relation also holds for the negative energy ( $\epsilon < 1$ ) case. In this situation the momentum  $p$  and the functions  $\eta$ ,  $Q$ , and  $\sin \xi$  become purely imaginary. Equations (2.16) remain, however, real, and so does the modified Coulomb equation (2.17).

### C. The nonrelativistic limit

The choice of the appropriate sign of  $\gamma$  in Eq. (2.16) will now be made by considering the behavior of the modified Coulomb equation (2.17) in the nonrelativistic limit. This limit is obtained by letting the speed of light tend to infinity. Thus we get

$$P \rightarrow \sqrt{2mT}/\hbar, \quad a/\lambda_c \rightarrow Z_{\text{nuc}} e^2/\hbar^2, \quad |\gamma| \rightarrow |\kappa|. \quad (2.19)$$

The modified Coulomb equation reduces to the Schrödinger equation for the point Coulomb potential,

$$\frac{d^2u}{dr^2} + \left( \frac{2mT}{\hbar^2} + \frac{2m}{\hbar^2} \cdot \frac{Z_{\text{nuc}} e^2}{r} - \frac{L(L+1)}{r^2} \right) u = 0, \quad (2.20)$$

where

$$L = \begin{cases} |\kappa| - 1, & \gamma > 0, \\ |\kappa|, & \gamma < 0. \end{cases} \quad (2.21)$$

We require that in the nonrelativistic limit, the solution of (2.20) coincides (up to a sign) with the large component function  $g$ . Therefore  $L(L+1) = \kappa(\kappa+1)$  and  $\text{sgn}(\gamma) = -\text{sgn}(\kappa)$ .<sup>10</sup>

With this choice of the sign  $\gamma$  we get in the nonrelativistic limit

$$g = \text{sgn}(\kappa) \cdot u, \quad f = \frac{1}{c} \frac{\hbar}{2m} r^{-\kappa} \frac{d}{dr} (r^\kappa \cdot g). \quad (2.22)$$

(Note that other conventions are often used; for example in Rose<sup>7</sup> the  $\text{sgn} \kappa$  factor is omitted.)

### D. Free particle case

The Dirac equation for a free particle in spherical coordinates is given by Eq. (2.8) with  $a = 0$ . The expressions we have derived are valid for this case, for which we have

$$\gamma = -\kappa, \quad \eta = 0, \quad \tan \xi = 0, \quad (2.23)$$

and the modified Coulomb equation reduces to the Bessel equation

$$\frac{d^2u}{dr^2} + \left( p^2 - \frac{\kappa(\kappa+1)}{r^2} \right) u = 0. \quad (2.24)$$

The solution of (2.24), regular at the origin, is

$$u = \tilde{A} r j_l(pr), \quad (2.25)$$

where  $\tilde{A}$  is a constant of normalization. The corresponding Dirac functions are obtained from Eq. (2.16),

$$\begin{aligned} g &= u = \tilde{A} r j_l(pr), \\ f &= (Q/p)(u' + (\kappa/r)u) = (\kappa/|\kappa|) \tilde{A} Q r j_l'(pr), \end{aligned} \quad (2.26)$$

where  $\bar{l} = l - \kappa/|\kappa|$ . The solution of the Dirac equation which is irregular at the origin is obtained by replacing in Eq. (2.26) the Bessel functions  $j_l$  and  $j_l'$  with the Neumann functions  $n_l$  and  $n_l'$ , respectively.

## III. SCREENED COULOMB POTENTIAL AND THE MODIFIED SCHRÖDINGER EQUATION

We consider now the case of an electron in a screened Coulomb potential. We assume that the potential is of the form

$$U(r) = -(a_0/r)s(r), \quad (3.1)$$

where  $a_0 = \alpha Z_{\text{nuc}}$  and  $s(r)$  is a smooth and monotonic function such that

$$s(r) = \begin{cases} 1, & r = 0, \\ s_t, & r \geq r_t. \end{cases} \quad (3.2)$$

The tail radius  $r_t$  is the position at which screening has its full effect. Beyond this point the potential is point Coulomb ( $s_t > 0$ ) or it vanishes ( $s_t = 0$ ),

$$U(r) = -a_t/r, \quad r \geq r_t, \quad (3.3)$$

where  $a_t = a_0 \cdot s_t$ .

### A. The modified Schrödinger equation

We look for a solution of the Dirac equation of the form (2.16) in the screened potential, with  $u$  again the solution of

a Schrödinger-like equation. Since the effective charge seen by the electron is now position dependent, we now must allow  $\xi$  and  $\gamma$  to be also position dependent. Thus we now have three functions  $u, \gamma, \xi$ ; we shall use our freedom in the choice of  $\gamma$  and  $\xi$  to obtain a convenient result for  $u$ . (Other choices might also be useful.)

When we substitute Eqs. (2.16) in the Dirac equation, we find that the function  $u$  must satisfy two (second-order) differential equations, resulting from the two equations (2.7a) and (2.7b). It turns out that these two equations for  $u$  are the same if we require that the coefficients of the first derivative  $du/dr$  vanish. Thus we get two equations connecting the "rotation angle"  $\xi$  and the "angular momentum parameter"  $\gamma$  with the screening function  $s(r)$

$$r \frac{d\xi}{dr} = -(\kappa - \gamma(r)) \tan \xi(r) + a_0 s(r) Q, \quad (3.4a)$$

$$r \frac{d\gamma}{dr} = -\frac{\kappa + \gamma(r)}{\tan \xi(r)} + \frac{a_0 s(r)}{Q}. \quad (3.4b)$$

When these equations are satisfied,  $u$  is a solution of the equation

$$\frac{d^2 u}{dr^2} + \left( p^2 + \frac{2a_0 \epsilon}{\lambda_c r} w(r) - \frac{\gamma(r)[\gamma(r) - 1]}{r^2} \right) u = 0, \quad (3.5)$$

where

$$w(r) = \frac{\lambda_c}{2a_0 \epsilon} \left[ -\frac{d\gamma}{dr} + p \left( \frac{\kappa + \gamma(r)}{\tan \xi(r)} + (\kappa - \gamma(r)) \tan \xi(r) \right) \right]. \quad (3.6)$$

Equation (3.5) will be referred to as the modified Schrödinger equation. It should be emphasized that Eqs. (3.4) are energy dependent through the parameter  $Q = \sqrt{(\epsilon - 1)/(\epsilon + 1)}$ . Therefore, both  $\tan \xi(r)$  and  $\gamma(r)$  are energy dependent and so is the effective screening function  $w(r)$ . However, as we will show below, we can obtain both  $\xi$  and  $\gamma$  by solving for one energy-independent function.

Before solving the modified Schrödinger equation (3.5) we have to solve Eqs. (3.4), which are the generalization of Eqs. (2.18) for the case of a screened Coulomb potential. We eliminate  $\gamma(r)$  between (3.4a) and (3.4b) and get

$$r \frac{d \tan \xi}{dr} = \frac{a_0 s}{Q} \tan^2 \xi - 2\kappa \tan \xi + a_0 s Q, \quad (3.7a)$$

$$\gamma = \{ \kappa (\tan^2 \xi - 1) + a_0 s (1 - Q^2) (\tan \xi / Q) \} \times \{ \tan^2 \xi + 1 \}^{-1}. \quad (3.7b)$$

There are two "natural" initial values for the solution of Eqs. (3.7), which correspond to the values of  $\xi$  and  $\gamma$  in a point Coulomb potential of an appropriate nuclear charge, which depends on the point  $r_{\text{init}}$  at which the integration of Eq. (3.7a) is started. Namely,

$$\gamma_0 = -\text{sgn}(\kappa) \sqrt{\kappa^2 - a_0^2}, \quad r_{\text{init}} = 0, \quad (3.8a)$$

$$\tan \xi_0 = a_0 Q / (\kappa - \gamma_0),$$

and

$$\gamma_t = -\text{sgn}(\kappa) \sqrt{\kappa^2 - a_t^2}, \quad r_{\text{init}} = r_t. \quad (3.8b)$$

$$\tan \xi_t = a_t Q / (\kappa - \gamma_t),$$

We will choose that initial value for which the solution of Eqs. (3.7) is smooth and well behaved, as will be discussed in the next section.

## B. The subsidiary equation

We transform now Eq. (3.7a) to an energy-independent form. We define the subsidiary function

$$\theta(r) = [\tan \xi(r)] / Q, \quad (3.9)$$

and we get

$$r \frac{d\theta}{dr} = a_0 s \theta^2 - 2\kappa \theta + a_0 s. \quad (3.10)$$

Equation (3.10), which will be referred to as the subsidiary equation, is energy independent, and so are the natural Coulombic initial values which correspond to Eqs. (3.8)

$$\theta_0 = a_0 / (\kappa - \gamma_0), \quad r_{\text{init}} = 0, \quad (3.11a)$$

$$\theta_t = a_t / (\kappa - \gamma_t), \quad r_{\text{init}} = r_t. \quad (3.11b)$$

It follows, then, that when one of these initial values is chosen,  $\theta(r)$  is energy independent everywhere. This property of  $\theta(r)$  makes the numerical solution of the modified Schrödinger equation (3.5) more tractable than would have appeared from Eqs. (3.4).

The question now arises as to whether there is a solution  $\theta(r)$  which varies smoothly between these two initial values as  $r$  moves from the origin to the tail region or vice versa. Such a solution, if it exists, has a Coulombic feature in the two regions where the potential is point Coulombic. It is usually impossible for a solution of a first-order differential equation to satisfy two boundary conditions. However, the solution  $\theta(r)$  of the subsidiary equation (3.10) has the property that its limiting values as  $r \rightarrow 0$  and  $r \rightarrow \infty$  are independent of the initial values (see Secs. 1 and 2 in the Appendix). In particular we have

$$\lim_{r \rightarrow 0} \theta(r) = \theta_0, \quad \text{for } \kappa < 0, \quad r_{\text{init}} > 0,$$

$$\lim_{r \rightarrow \infty} \theta(r) = \theta_t, \quad \text{for } \kappa > 0, \quad r_{\text{init}} < \infty.$$

Therefore, when  $\kappa < 0$  we start the integration at  $r_{\text{init}} = r_t$  with the initial value  $\theta_t$  and integrate inward. The solution  $\theta(r)$  is a monotonically increasing function (see Appendix C) satisfying the two boundary values (3.11). When  $\kappa > 0$  we start the integration at  $r_{\text{init}} = 0$  with the initial value  $\theta_0$  and integrate outward. The solution  $\theta(r)$  in this case is a monotonically decreasing function which approaches  $\theta_t$  as  $r$  tends to infinity (see Appendix D). At  $r = r_t$ , however, we have (for  $\kappa > 0$ )  $\theta(r_t) > \theta_t$ , and the boundary condition at this point is not satisfied. As far as the solution of the Dirac equation is concerned, this poses no special difficulty, as we shall see later.



### C. Some properties of the modified Schrödinger equation and its solution

With the function  $\theta(r)$  given, we calculate  $\tan \xi(r)$  and  $\gamma(r)$  by (3.9) and (3.7b),

$$\begin{aligned} \tan \xi(r) &= Q \cdot \theta(r), \\ \gamma(r) &= \{\kappa [Q^2 \theta^2(r) - 1] + a_0 s(r) (1 - Q^2) \theta(r)\} \\ &\quad \times \{Q^2 \theta^2(r) + 1\}^{-1}. \end{aligned} \quad (3.12)$$

Then we integrate the modified Schrödinger equation (3.5) to obtain  $u$ , and using Eq. (2.16) we obtain the solution  $g$  and  $f$  to the Dirac equation. Let us look at some properties of Eq. (3.5) and of its solution  $u$ .

(i) *Near the origin* ( $r \rightarrow 0$ ). Since  $\theta(0) = \theta_0$ , we get from (3.12) and (3.10) that  $\gamma(0) = \gamma_0$  and  $\lim_{r \rightarrow 0} r(d\gamma/dr) = 0$ . Therefore the dominant term of the modified Schrödinger equation near the origin is the centrifugal potential  $\gamma_0(\gamma_0 - 1)/r^2$ . This assures  $r^{\gamma_0}(r^{1-\gamma_0})$  behavior of  $u$  near the origin for the regular (irregular) solution for  $\kappa < 0$ , or for the irregular (regular) solution when  $\kappa > 0$ . This, in turn, assures  $r^{|\gamma_0|}(r^{-|\gamma_0|})$  behavior for the regular (irregular) solution (2.14) of the Dirac equation.

(ii) *The tail region* ( $r \gg r_t$ ). When  $\kappa < 0$  we have  $\theta(r \gg r_t) = \theta_t$  and by (3.6) and (3.12) we get in this region  $\gamma(r) = \gamma_t$ ,  $\tan \xi(r) = \tan \xi_t$ ,  $w(r) = s_t$  ( $r \gg r_t$ ). (3.13)

Thus the modified Schrödinger equation (3.5) reduces to the modified Coulomb equation (2.17) for charge  $Z = Z_{\text{nuc}} \cdot s_t$  (ionic potential,  $s_t > 0$ ) or to the Bessel equation (2.24) (neutral atom potential,  $s_t = 0$ ).

When  $\kappa > 0$  the boundary value (3.8b) is not satisfied and the modified Schrödinger equation (3.5) does not reduce to either Eq. (2.17) or to Eq. (2.24). However, since  $\lim_{r \rightarrow \infty} \theta(r) = \theta_t$ ,  $\xi(r)$  and  $\gamma(r)$  approach the Coulombic values  $\xi_t$  and  $\gamma_t$  as  $r \rightarrow \infty$ . The asymptotic behavior of the solution  $u(r)$  of Eq. (3.5) is the same as the asymptotic behavior of a particular solution  $u_t$  of the modified Coulomb equation (2.17) with  $\gamma = \gamma_t$  and  $a = a_t$ , as we will now show.

The solution ( $g, f$ ) of the Dirac equation in the tail region can be obtained in two ways: (i) by solving the modified Schrödinger equation (3.5) and substituting its solution in Eqs. (2.16); and (ii) by solving the modified Coulomb equation (2.17) with  $\gamma = \gamma_t$  and  $\xi = \xi_t$ , and substituting its solution in Eqs. (2.16). We denote by  $u_t$ , that particular solution of (2.17) for which the two solutions of the Dirac equation coincide at  $r = r_t$  (and therefore for  $r \gg r_t$ ),

$$\begin{aligned} g_t(r) &= \cos \xi_t \cdot u_t(r) + \sin \xi_t \cdot v_t(r), \\ f_t(r) &= Q [-\sin \xi_t \cdot u_t(r) + \cos \xi_t \cdot v_t(r)], \end{aligned} \quad (3.14)$$

where  $v_t = (1/p)[u_t'(r) - (\gamma_t/r)u_t(r)]$ .

By equating the two solutions we get

$$\begin{aligned} u(r) &= \cos(\xi_t - \xi) \cdot u_t(r) + \sin(\xi_t - \xi) \cdot v_t(r), \\ v(r) &= -\sin(\xi_t - \xi) u_t(r) + \cos(\xi_t - \xi) \cdot v_t(r), \end{aligned} \quad (3.15)$$

where  $v(r) = (1/p)[u'(r) - (\gamma(r)/r)u(r)]$  and  $\xi = \arctan[Q \cdot \theta(r)]$ . Since  $\xi(r) \rightarrow_{r \rightarrow \infty} \xi_t$  we see that the asymptotic behavior of  $u(r)$  and its derivative  $u'(r)$  is the

same as the asymptotic behavior of the Coulomb function  $u_t(r)$  and its derivative  $u_t'(r)$ , respectively.

(iii) *The nonrelativistic limit*. The boundary values  $\theta_0$  and  $\theta_t$  [Eq. (3.11)] are given in this limit by

$$\theta_0 \approx (1/c) \cdot (e^2 Z_{\text{nuc}} / 2\hbar\kappa), \quad \theta_t \approx \theta_0 \cdot s_t.$$

Therefore from Eq. (3.10) we find that to the lowest order in  $1/c$ ,  $\theta(r)$  is given by  $\theta(r) = (1/c)\hat{\Theta}(r)$ , where  $\hat{\Theta}(r)$  does not depend on  $c$ . On expanding Eqs. (3.6) and (3.12) in  $1/c$  we get

$$\begin{aligned} \gamma(r) &\approx -\kappa + (1/c^2)(e^2 Z_{\text{nuc}} / \hbar) \cdot s(r) \cdot \hat{\Theta}(r), \\ \tan \xi(r) &\approx (1/c^2) \sqrt{(T/2m)} \hat{\Theta}(r), \end{aligned} \quad (3.16)$$

$$w(r) \approx s(r),$$

and the modified Schrödinger equation reduces to the Schrödinger equation for a screened Coulomb potential,

$$\frac{d^2 u}{dr^2} + \left( \frac{2mT}{\hbar^2} + \frac{2me^2 Z_{\text{nuc}}}{\hbar^2} \cdot \frac{s(r)}{r} - \frac{l(l+1)}{r^2} \right) u = 0, \quad (3.17)$$

and again,  $g$  and  $f$  are given by (2.22).

### IV. APPLICATION: DERIVATION OF RELATIVISTIC QUANTUM DEFECT THEORY

Here we illustrate the application of these ideas, using as our example a rederivation of single channel relativistic quantum defect theory<sup>11,12</sup> with the methods of the nonrelativistic derivation. Applying a nonrelativistic derivation based on the modified Schrödinger equation (3.5), we obtain the corresponding relativistic theory from Eqs. (2.16).

It was shown in Sec. III that the asymptotic behavior of the solution  $u$  of Eq. (3.5) is given by the function  $u_t$ , which is a solution of the modified Coulomb equation

$$\frac{d^2 u}{dr^2} + \left( p^2 + \frac{2a\epsilon}{\lambda_c r} + \frac{\gamma(\gamma-1)}{r^2} \right) u = 0, \quad (4.1)$$

where

$$a = \alpha Z_{\text{nuc}} s_t, \quad (4.2a)$$

and

$$\gamma = \begin{cases} \gamma_t, & \kappa < 0, \\ 1 - \gamma_t, & \kappa > 0. \end{cases} \quad (4.2b)$$

We consider the case of an ionic potential ( $a > 0$ ) and we apply Seaton's quantum defect theory<sup>13</sup> to the modified Schrödinger equation. We follow the QDT method as closely as possible, bearing in mind that the angular momentum parameter is not an integer in our case.

#### A. Solutions of the modified Coulomb equation

We first change the variable  $r$  in (4.1) to  $x = a\epsilon r / \lambda_c$  and we obtain the modified Coulomb equation in its dimensionless form

$$\frac{d^2u}{dx^2} + \left( -\frac{1}{\sigma^2} + \frac{2}{x} - \frac{\gamma(\gamma-1)}{x^2} \right) u = 0, \quad (4.3)$$

where

$$\sigma = \begin{cases} i\eta = i a \epsilon / \sqrt{\epsilon^2 - 1}, & \epsilon > 1, \\ \nu = a \epsilon / \sqrt{1 - \epsilon^2}, & \epsilon < 1. \end{cases} \quad (4.4)$$

We consider the following solutions to Eq. (4.3) (Ref. 14):

$$y_1(\sigma, \gamma; x) = [(2x)^\gamma / \Gamma(2\gamma)] e^{-x/\sigma} M(\gamma - \sigma, 2\gamma; 2x/\sigma),$$

$$y_2(\sigma, \gamma; x) = y_1(\sigma, 1 - \gamma; x), \quad (4.5)$$

where  $M(\alpha, \beta; Z)$  is the confluent hypergeometric function as defined in Ref. 15. It was shown by Seaton<sup>13</sup> that these functions are analytic functions of  $1/\sigma^2$  and therefore they are analytic functions of the energy for  $\epsilon \gg \epsilon_0 > 0$ . Moreover, since  $\gamma$  is not an integer or half-integer, these functions are algebraically independent.

## B. Asymptotic behavior

From the asymptotic behavior of the confluent hypergeometric function<sup>16</sup> we get the following expression for  $y_1$ :

$$y_1 \underset{x \rightarrow \infty}{\sim} \sigma^\gamma \left\{ \frac{e^{x/\sigma}}{\Gamma(\gamma - \sigma)} \left( \frac{2x}{\sigma} \right)^{-\sigma} + \frac{e^{-x/\sigma}}{\Gamma(\gamma + \sigma)} \left( \frac{2x}{\sigma} \right)^\sigma e^{-i\pi(\gamma - \sigma)} \right\}. \quad (4.6)$$

(i) Below threshold ( $\epsilon < 1$ ). In this case  $\sigma = \nu$  and the function grows exponentially as  $x \rightarrow \infty$ . The function  $y_1$  can represent a bound state only when the exponentially growing part of the function vanishes. This is the case when  $\gamma - \sigma$  is a nonpositive integer. The conditions for the energy levels are therefore

$$\nu(\epsilon_{n,\kappa}) = a \epsilon_{n,\kappa} / \sqrt{1 - \epsilon_{n,\kappa}^2} = \gamma - |\kappa| + n$$

$$(n = |\kappa|, |\kappa| + 1, |\kappa| + 2, \dots). \quad (4.7)$$

This relation is equivalent to the expression for the energy levels of hydrogenlike atoms in the Dirac theory.

(ii) Above threshold ( $\epsilon > 1$ ). In this case  $\sigma = i\eta$ , the function remains bounded as  $x \rightarrow \infty$ , and we have

$$y_1 \underset{x \rightarrow \infty}{\sim} [2n^\gamma e^{-\pi\eta/2} / |\Gamma(\gamma + i\eta)|] \sin \omega, \quad (4.8)$$

where

$$\omega = x/\eta + \eta \ln(2x/\eta) + (\pi/2)(1 - \gamma) - \arg \Gamma(\gamma + i\eta). \quad (4.9)$$

The asymptotic behavior of  $y_2(\sigma, \gamma; x)$  is obtained from Eq. (4.6) by replacing  $\gamma$  with  $1 - \gamma$ ,

$$y_2 \underset{x \rightarrow \infty}{\sim} \sigma^{1-\gamma} \left\{ \frac{e^{x/\sigma}}{\Gamma(1 - \gamma - \sigma)} \left( \frac{2x}{\sigma} \right)^{-\sigma} + \frac{e^{-x/\sigma}}{\Gamma(1 - \gamma + \sigma)} \left( \frac{2x}{\sigma} \right)^\sigma e^{-i\pi(1 - \gamma - \sigma)} \right\}. \quad (4.10)$$

(i) Below threshold ( $\epsilon < 1$ ). Here this function also grows exponentially as  $x \rightarrow \infty$ . Since the function is irregular at the origin, bound states of the point Coulomb potential are not given by the function  $y_2$  for hydrogenlike atoms. However, in the case of a screened potential with a Coulomb tail, bound states are represented in the exterior region by the

linear combination of  $y_1$  and  $y_2$  which decays exponentially as  $x \rightarrow \infty$ . Therefore we have to look for the linear combination of  $y_1$  and  $y_2$  in which the exponentially growing part vanishes.

Using the relation  $\Gamma(Z)\Gamma(1-Z) = \pi/\sin \pi Z$ , we find that the exponentially growing part of  $y_1$  is given by  $K \sin \pi(\gamma - \nu)$ , where

$$K = (1/\pi) \nu^\gamma \Gamma(1 - \gamma + \nu) (2x/\nu)^{-\nu} e^{x/\nu}. \quad (4.11)$$

To describe bound states (for the screened case) in a convenient form, we look for a linear combination of  $y_1$  and  $y_2$  in which the exponentially growing part is given by  $-K \cos \pi(\gamma - \nu)$ . The coefficients of the linear combination are given by

$$K_1 = -\cot 2\pi\gamma,$$

$$K_2 = -\nu^{2\gamma-1} [\Gamma(1 - \gamma + \nu) / \Gamma(\gamma + \nu)] \cdot (\sin 2\pi\gamma)^{-1}, \quad (4.12)$$

and we get

$$y_0 = K_1 y_1 + K_2 y_2 = -K \cos \pi(\gamma - \nu)$$

$$- \nu^\gamma \cdot [e^{-x/\nu} / \Gamma(\gamma + \nu)] (2x/\nu)^\nu [\cot 2\pi\gamma e^{-i\pi(\gamma - \nu)}$$

$$+ (\sin 2\pi\gamma)^{-1} e^{-i\pi(1 - \gamma - \nu)}]. \quad (4.13)$$

The coefficient  $K_2$  in Eq. (4.12) can be analytically continued to energies above threshold<sup>17</sup>

$$K_2(\epsilon > 1) = -\sigma^{2\gamma-1} \frac{\Gamma(1 - \gamma + \sigma)}{\Gamma(\gamma + \sigma)} \cdot \frac{1}{\sin 2\pi\gamma}$$

$$= -\eta^{2\gamma-1} \left| \frac{\Gamma(1 - \gamma + i\eta)}{\Gamma(\gamma + i\eta)} \right| \frac{1}{\sin 2\pi\gamma}. \quad (4.14)$$

Therefore the function  $y_0$  defined in Eq. (4.13) is an analytic function of the energy near the threshold.

(ii) Above threshold ( $\epsilon > 1$ ). Here we have

$$y_2 \underset{x \rightarrow \infty}{\sim} [2\eta^{1-\gamma} e^{-\pi\eta/2} / |\Gamma(1 - \gamma + i\eta)|] \sin(\tilde{\omega}), \quad (4.15)$$

where

$$\tilde{\omega} = x/\eta + \eta \ln(2x/\eta) + \pi\gamma/2 - \arg \Gamma(1 - \gamma + i\eta). \quad (4.16)$$

In the case of a screened potential, we characterize continuum wave functions by phase shifts. It is thus more convenient to have two independent solutions to the modified Coulomb equation (4.3); the regular solution behaves asymptotically like  $\sin(\omega)$  and the irregular solution like  $\cos(\omega)$ . We define therefore

$$L_1 = -\cot(\tilde{\omega} - \omega),$$

$$L_2 = \eta^{2\gamma-1} \left| \frac{\Gamma(1 - \gamma + i\eta)}{\Gamma(\gamma + i\eta)} \right| \frac{1}{\sin(\tilde{\omega} - \omega)} \quad (4.17)$$

$$= -K_2 \frac{\sin 2\pi\gamma}{\sin(\tilde{\omega} - \omega)},$$

and we get

$$\tilde{y}_0 = L_1 y_1 + L_2 y_2 \underset{x \rightarrow \infty}{\sim} [2\eta^\gamma e^{-\pi\eta/2} / |\Gamma(\gamma + i\eta)|] \cos(\omega). \quad (4.18)$$

At threshold we have

$$K_1(\epsilon = 1) = L_1(\epsilon = 1) = -\cot 2\pi\gamma, \quad (4.19)$$

$$K_2(\epsilon = 1) = L_2(\epsilon = 1) = -(\sin 2\pi\gamma)^{-1},$$

### C. Quantum defect and phase shift of the solution of the modified Schrödinger equation

Let  $u(r)$  be a solution of Eq. (3.5), regular at the origin. Let  $u_i(r)$  be a solution of Eq. (4.1) such that in the tail region ( $r \gg r_i$ ); for  $\kappa < 0$  we have  $u_i(r) = u(r)$ , and for  $\kappa > 0$  Eqs. (3.15) are satisfied. We will use the basis functions defined above and we change therefore the variable  $r$  to  $x = a\epsilon r/\lambda_c$ ,

$$y(x) = u(r), \quad y_i(x) = u_i(r). \quad (4.20)$$

(i) Below threshold ( $\epsilon < 1$ ). We write  $y_i(x)$  as

$$y_i(x) = c(\cos \pi\mu \cdot y_1(x) + \sin \pi\mu \cdot y_0(x)), \quad (4.21)$$

where  $c$  is a constant and  $\mu$  is an analytic function of the energy.<sup>18</sup>

For the function  $y(x)$  to remain bounded, so that corresponding Dirac functions  $g$  and  $f$  represent a bound state, the exponentially growing part of  $y_i(x)$  must vanish. This condition is satisfied when  $\sin \pi(\gamma - \nu - \mu) = 0$ , or

$$\gamma - \nu(\epsilon) - \mu(\epsilon) = |\kappa| - n = \text{integer}, \quad (4.22)$$

from which the eigenvalue  $\epsilon_{n,\kappa}$  is obtained. We write Eq. (4.22) as  $\nu = \gamma - |k| + (n - \mu)$  and compare it with (4.7). It follows that  $\mu(\epsilon_{n,\kappa})$  is the quantum defect caused by the non-Coulomb part of the modified Schrödinger equation (3.5).

(ii) Above threshold ( $\epsilon > 1$ ). We write  $y_i(x)$  as

$$y_i(x) = c(\cos \delta \cdot y_1(x) + \sin \delta \cdot \tilde{y}_0(x))$$

$$\sim c [2\eta^\gamma e^{-\pi\eta/2} / |\Gamma(\gamma + i\eta)|] \sin(\omega + \delta). \quad (4.23)$$

Thus,  $\delta(\epsilon)$  is the phase shift caused by the non-Coulomb part of Eq. (3.5). We may also express  $y_i(x)$  in the form of Eq. (4.21) for energies above threshold. Comparing coefficients of  $y_1$  and  $y_2$  we get

$$\cot \delta(\epsilon) = \cot(\tilde{\omega} - \omega) - [(\sin 2\pi\gamma) / \sin(\tilde{\omega} - \omega)] \times (\cot \pi\mu(\epsilon) - \cot 2\pi\gamma). \quad (4.24)$$

Near threshold, for  $\eta \gg 1$ , we have

$$\cot \delta(\epsilon) = (1 - \cos 2\pi\gamma e^{-2\pi\eta}) \cot \pi\mu(\epsilon) - \sin 2\pi\gamma e^{-2\pi\eta}, \quad (4.25)$$

and at threshold

$$\cot \delta(\epsilon = 1) = \cot \pi\mu(\epsilon = 1). \quad (4.26)$$

Equations (4.25) and (4.26) are the quantum defect relations of the modified Schrödinger equation. As we will see below,  $\mu(\epsilon)$  and  $\delta(\epsilon)$  are also the quantum defect and phase shift of the Dirac functions, and therefore Eqs. (4.24)–(4.26) are the relativistic quantum defect relations. They are identical with the relations obtained by Johnson and Cheng.<sup>11</sup>

### D. Quantum defect and phase shift of the Dirac functions

We substitute now an unnormalized regular solution  $u(r)$  of the modified Schrödinger equation (3.5) in the expressions (2.16) and we get an unnormalized regular solution of the Dirac equation,

$$\begin{pmatrix} g \\ f \end{pmatrix} = \begin{pmatrix} \cos \xi \\ -Q \sin \xi \end{pmatrix} u + \frac{1}{p} \begin{pmatrix} \sin \xi \\ Q \cos \xi \end{pmatrix} \left( u' - \frac{\gamma}{r} u \right), \quad (4.27)$$

where  $\xi$  and  $\gamma$  are the solution of Eqs. (3.10) and (3.12).

In the tail region ( $r \gg r_i$ ) the Dirac functions  $g$  and  $f$  are the same as the functions defined in Eq. (3.14) in terms of the Coulomb function  $u_i(r)$ . The function  $u_i(r)$  can be expressed as a linear combination of the basis functions  $u_1(r)$  and  $u_0(r)$ ,

$$u_i(r) = y_1(x), \quad u_0(r) = y_0(x), \quad (4.28)$$

and therefore  $g$  and  $f$  are linear combinations, with the same coefficients, of the two solutions ( $g_1 f_1$ ) and ( $g_0 f_0$ ) obtained by substituting  $u_1$  and  $u_0$  in Eqs. (2.16), respectively.

(i) Below threshold ( $\epsilon < 1$ ). From Eq. (4.21) we get

$$\begin{pmatrix} g \\ f \end{pmatrix} = c \cdot \left\{ \cos \pi\mu \begin{pmatrix} g_1 \\ f_1 \end{pmatrix} + \sin \pi\mu \begin{pmatrix} g_0 \\ f_0 \end{pmatrix} \right\}, \quad r \gg r_i. \quad (4.29)$$

When  $\mu$  satisfies Eq. (4.22),  $g$  and  $f$  decay exponentially as  $r \rightarrow \infty$ , and thus represent a bound state. Therefore, the quantum defect function  $\mu(\epsilon)$  defined for the solutions of the modified Schrödinger equation is identical with the quantum defect function of the Dirac functions.

(ii) Above threshold ( $\epsilon > 1$ ). Let  $(\tilde{g}_0, \tilde{f}_0)$  be the Dirac functions obtained by substituting  $\tilde{u}_0(r) = \tilde{y}_0(x)$  in Eqs. (2.16). Then, from Eq. (4.23) we get

$$\begin{pmatrix} g \\ f \end{pmatrix} = c \left\{ \cos \delta \begin{pmatrix} g_1 \\ f_1 \end{pmatrix} + \sin \delta \begin{pmatrix} \tilde{g}_0 \\ \tilde{f}_0 \end{pmatrix} \right\}, \quad r \gg r_i. \quad (4.30)$$

The asymptotic behavior of  $g_1$  and  $f_1$  is given by (2.10). The asymptotic behavior of  $\tilde{g}_0$  and  $\tilde{f}_0$  is obtained from (4.18),

$$\begin{aligned} \tilde{g}_0 &\sim A \cos(pr - \frac{1}{2}l\pi + \delta_c + \eta \ln(2pr)), \\ \tilde{f}_0 &\sim -AQ \sin(pr - \frac{1}{2}l\pi + \delta_c + \eta \ln(2pr)), \end{aligned} \quad (4.31)$$

where the relativistic Coulomb phase shift is given by Eq. (2.11a). It follows from (4.30) that  $\delta(\epsilon)$  is indeed the phase shift, with respect to the regular solution of the Dirac equation in the point Coulomb potential, caused by the non-Coulomb part of the potential. As we already saw,  $\delta(\epsilon)$  is also the phase shift, with respect to the Coulomb function  $u_1(r)$ , caused by the non-Coulomb part in the modified Schrödinger equation (3.5).

Thus, Eqs. (4.24)–(4.26), which connect the phase shift  $\delta(\epsilon)$  and the quantum defect  $\mu(\epsilon)$  of the modified Schrödinger equation, are also the desired single channel relativistic quantum defect relations.

### V. CONCLUSION

We have shown that the Dirac equation in a screened Coulomb potential can be transformed to a modified Schrödinger equation [Eq. (3.5)] with an effective potential. The close similarity between the modified Schrödinger equation and the Schrödinger equation suggests that various proper-

ties of the Dirac equation can be derived from the corresponding properties of the Schrödinger equation. We have demonstrated the usefulness of this approach by obtaining the relativistic quantum defect relations from the corresponding relations of the modified Schrödinger equation. We expect this approach to be useful in obtaining other properties associated with the Dirac equation.

## ACKNOWLEDGMENT

Helpful discussions with Mr. J. C. Parker, III and Dr. J. Stein are gratefully acknowledged.

## APPENDIX: PROPERTIES OF THE SOLUTION OF THE SUBSIDIARY EQUATION (3.10)

### 1. Solution near the origin

We expand  $\theta(r)$  and  $s(r)$  in power series,

$$\theta(r) = \sum_{l=0}^{\infty} \tilde{\theta}_l r^l, \quad (A1)$$

$$s(r) = \sum_{l=0}^{\infty} s_l r^l = 1 + \sum_{l=1}^{\infty} s_l r^l.$$

On substituting in Eq. (3.10) and equating coefficients of the same order in  $r$  we get the following set of equations:

$$a_0 \tilde{\theta}_0^2 - 2\kappa \tilde{\theta}_0 + a_0 = 0, \quad (A2a)$$

$$\tilde{\theta}_l = a_0 [(c_l + s_l)/(2\kappa + l - 2a_0 \tilde{\theta}_0)], \quad l = 1, 2, \dots, \quad (A2b)$$

where

$$C_l = \sum_{\substack{m,n=0 \\ m+n < l}}^{l-1} S_{l-m-n} \tilde{\theta}_m \tilde{\theta}_n. \quad (A3)$$

The point  $r = 0$  is a singular point of the subsidiary equation (3.10). The behavior of its solution near the origin depends on the sign of the quantity

$$S = a_0 \tilde{\theta}_0 - \kappa,$$

where  $\tilde{\theta}_0$  is one of the two solutions of Eq. (A2a),<sup>19</sup>  $\tilde{\theta}_{0,1} = \theta_0$ ;  $\tilde{\theta}_{0,2} = 1/\theta_0$  [see (3.11a)]. When  $S > 0$ , the point ( $r = 0, \theta = \tilde{\theta}_0$ ) is a nodal point, and all the curves  $\theta(r)$  reach the point  $\tilde{\theta}_0$  as  $r \rightarrow 0$ . On the other hand, when  $S < 0$ , the point ( $r = 0, \theta = \theta_0$ ) is a saddle point, and apart from two lines (the principal lines), no other curve  $\theta(r)$  reaches the point  $\tilde{\theta}_0$ . For the solution  $\tilde{\theta}_0 = \theta_0$ , which corresponds to the initial value (3.11a), we have

$$\text{sgn}(S) = -\text{sgn}(k),$$

and therefore  $(0, \theta_0)$  is a nodal point for  $\kappa < 0$  and a saddle point for  $\kappa > 0$ . On the other hand, the point  $(0, 1/\theta_0)$  is a saddle point for  $\kappa < 0$  and a nodal point for  $\kappa > 0$ . It follows that for  $\kappa < 0$  we get  $\lim_{r \rightarrow 0} \theta(r) = \theta_0$  even when the integration of (3.10) begins at  $r_{\text{init}} > 0$ , with arbitrary initial value, while for  $\kappa > 0$  we get  $\lim_{r \rightarrow 0} \theta(r) = \theta_0$  only if we start the integration at  $r_{\text{init}} = 0$  with the initial value  $\theta_0$ , and integrate outward (or integrate along the principal line, which is very unlikely when the integration is carried out numerically).

### 2. Solution in the tail region

When  $r \gg r_t$ , the screening is constant and Eq. (3.10) can be solved analytically.

(i) Ionic potential,  $s_t > 0$ . Let  $\theta(r_t)$  be the initial value at  $r = r_t$ . Then

$$\theta(r) = \theta_t + (2\gamma_t/a_t) \cdot [cr^{2\gamma_t}/(1 - cr^{2\gamma_t})] \quad (r \gg r_t), \quad (A4)$$

where the constant  $c$  is given by

$$c = (1/r_t^{2\gamma_t}) \cdot \{[\theta(r_t) - \theta_t]/[\theta(r_t) - 1/\theta_t]\}. \quad (A5)$$

(ii) Neutral atom potential,  $s_t = 0$ . Here we have  $\theta_t = 0$  and

$$\theta(r) = \theta(r_t) (r_t/r)^{2k} \quad (r \gg r_t). \quad (A6)$$

For  $k < 0$  we choose  $\theta(r_t) = \theta_t$  and get

$$\theta(r) = \theta_t = 0 \quad \text{for } r \gg r_t. \quad (A7)$$

For  $\kappa > 0$   $\theta(r_t)$  is determined by integrating from the origin outward, and we get

$$\lim_{r \rightarrow \infty} \theta(r) = \theta_t = 0. \quad (A8)$$

### 3. Solution for $\kappa < 0, r \leq r_t$

We will show that in this case  $\theta(r)$  is a monotonically increasing function. We start the integration of Eq. (3.10) at  $r_{\text{init}} = r_t$  with the initial value  $\theta(r_t) = \theta_t$ . Let  $n$  be the order of the first nonvanishing, left side, derivative of the screening function  $s(r)$  at  $r = r_t$ . Then by induction the first nonvanishing, left side, derivative of  $\theta(r)$  at  $r = r_t$  is of the order  $n + 1$ , and we have

$$r_t \left( \frac{d^{n+1}\theta}{dr^{n+1}} \right)_{r_t} = a_0 (\theta_t^2 + 1) \left( \frac{d^n s}{dr^n} \right)_{r_t}. \quad (A9)$$

Since  $s(r)$  is assumed to be a monotonically decreasing function we get

$$\left( \frac{d^{n+1}\theta}{dr^{n+1}} \right)_{r_t} \begin{cases} > 0, & n \text{ even,} \\ < 0, & n \text{ odd,} \end{cases} \quad (A10)$$

and therefore, in the vicinity of  $r_t$  ( $r \leq r_t$ )  $\theta(r)$  is an increasing function of  $r$  and  $d\theta/dr > 0$  near  $r_t$ .

Suppose that while integrating Eq. (3.10) inward we encounter a point  $\bar{r}$  for which  $(d\theta/dr)_{\bar{r}} = 0$ . Then, at this point we have [ $s(r)$  is a monotonically decreasing function]

$$\left( \frac{d^2\theta}{dr^2} \right)_{\bar{r}} = a_0 (\theta^2(\bar{r}) + 1) \left( \frac{ds}{dr} \right)_{\bar{r}} < 0. \quad (A11)$$

Thus  $\theta(r)$  has a maximum at  $r = \bar{r}$ , which is impossible since  $d\theta/dr > 0$  for  $\bar{r} < r < r_t$ . Therefore  $d\theta/dr > 0$  for the whole range  $(0, r_t)$ .

### 4. Solution for $\kappa > 0$

We will show that  $\theta(r)$  is a monotonically decreasing function for  $r > 0$ . Let  $s_n$  ( $n \geq 1$ ) be the first nonvanishing coefficient in (A1). Then, the first nonvanishing coefficient (for  $n \geq 1$ ) in the expansion of  $\theta(r)$  is

$$\theta_n = 2\theta_0 \kappa S_n / (n - 2\gamma_0) < 0, \quad (A12)$$

and therefore  $d\theta/dr < 0$  for  $r \geq 0$ .

Suppose that while integrating Eq. (3.10) outward we encounter a point  $\bar{r}$  for which  $(d\theta/dr)_{\bar{r}} = 0$ . Then  $\theta(r)$  has a maximum at this point, which is, however, impossible since  $d\theta/dr < 0$  for  $r < \bar{r}$ . Therefore  $d\theta/dr < 0$  for the whole range  $r > 0$ .

<sup>1</sup>P. G. Burke and W. D. Robb, *Adv. At. Mol. Phys.* **11**, 153 (1975).

<sup>2</sup>F. Herman and S. Skillman, *Atomic Structure Calculations* (Prentice-Hall, Englewood Cliffs, NJ, 1963).

<sup>3</sup>D. Liberman, J. T. Waber, and D. T. Cromer, *Phys. Rev.* **137**, A27 (1965).

<sup>4</sup>Nonrelativistic: A Zangwill and Paul Soven, *Phys. Rev. A* **21**, 1561 (1980); relativistic: F. A. Parpia and W. R. Johnson, *J. Phys. B* **17**, 531 (1984).

<sup>5</sup>Nonrelativistic: M. Ya Amusia and N. A. Cerepkov, *Case Stud. At. Phys.* **5**, 47 (1975); relativistic: W. R. Johnson and C. D. Lin, *Phys. Rev. A* **20**, 964 (1979).

<sup>6</sup>Nonrelativistic: A. Dalgarno and G. A. Victor, *Proc. R. Soc. London Ser. A* **291**, 291 (1966); relativistic: I. P. Grant, *Adv. Phys.* **19**, 747 (1970).

<sup>7</sup>See, e.g., H. A. Bethe and E. E. Salpeter, *Quantum Mechanics of One and Two Electron Atoms* (Academic, New York, 1957); M. E. Rose, *Relativistic Electron Theory* (Wiley, New York, 1961).

<sup>8</sup>For the point Coulomb potential this is essentially the transformation studied by L. C. Biedenharn, *Phys. Rev.* **126**, 845 (1962), which also has recently been discussed by J. Y. Su, *Phys. Rev. A* **32**, 3251 (1985). In this paper we are concerned with the possibility of generalizing the transformation to the screened potential case.

<sup>9</sup>I. B. Goldberg, J. Stein, Akiva Ron, and R. H. Pratt, Fourteenth International Conference on the Physics of Electronic and Atomic Collisions, Abstracts of Contributed Papers, Palo Alto, CA, 1985, p. 135.

<sup>10</sup>This choice of the sign of  $\gamma$  differs from the choice  $\text{sign}(\gamma) = \text{sign}(\kappa)$  made in Ref. 8.

<sup>11</sup>W. R. Johnson and K. T. Cheng, *J. Phys. B* **12**, 863 (1979).

<sup>12</sup>V. A. Zilitis, *Opt. Spectrosc. (USSR)* **50**, 227 (1981).

<sup>13</sup>M. J. Seaton, *Mon. Not. R. Astron. Soc.* **118**, 504 (1958).

<sup>14</sup>Reference 13, Eqs. (8) and (9).

<sup>15</sup>M. Abramowitz and I. A. Stegun, *Handbook of Mathematical Functions* (National Bureau of Standards, Washington, 1964), Eq. (13.1.2).

<sup>16</sup>Bateman Manuscript Project, *Higher Transcendental Functions* (McGraw-Hill, New York, 1953), Vol. 1, p. 278.

<sup>17</sup>Reference 16, p. 48, Eq. (12), and p. 37, Eq. (12).

<sup>18</sup>F. S. Ham, *Solid State Physics*, edited by F. Seitz and D. Turnbull (Academic, New York, 1955), Vol. 1, p. 127.

<sup>19</sup>E. L. Ince, *Integration of Ordinary Differential Equations* (Interscience, New York, 1952), pp. 35–39.

# A geometric approach to quantum scattering with group symmetry

Yihren Wu

Department of Mathematics, Hofstra University, Hempstead, New York 11550

(Received 16 October 1986; accepted for publication 11 February 1987)

A geometric theory for quantum scattering when the symmetry group is semisimple is presented. This theory is seen as a generalization of the partial wave analysis. As an application of this theory, the  $S$ -matrix elements for scattering in the Pöschl–Teller potential with symmetry group  $SO(1,2)$ , Coulomb potential with  $SO(1,3)$ , and a perturbed Coulomb potential with  $SO(2,3)$  are calculated. The last example may be considered as a model for heavy-ion scattering.

## I. INTRODUCTION

In a series of papers,<sup>1</sup> Alhassid *et al.* have studied the quantum scattering problems having an  $SU(1,1)$  symmetry group. Some examples are scattering with the Morse potential and Pöschl–Teller potential. Their work involves a procedure called the “Euclidean connection” in which they show that the shift operators acting on the set of asymptotic states can be written in terms of operators from the Euclidean algebra, thus giving rise to recursive relations for the transmission and reflection coefficients in the adjacent states, and the  $S$ -matrix elements are determined. This procedure is completely group theoretical, in a sense that no physical aspect of the problem enters the calculation, so the result will hold [modulo the relation between the Casimir operator of  $SU(1,1)$  and the energy Hamiltonian] for all situations having such a group symmetry. This procedure is extended to study the  $SO(2,n)$  groups.

Recently, we have shown<sup>2</sup> that the scattering problem is closely related to the Radon transform if the symmetry group is semisimple. Based on the assumption that the  $S$ -matrix elements should look the same, we consider the “standard” scattering problem, namely the geodesic flow problem on the symmetric space associated with the group. We quantize (via the techniques of geometric quantization<sup>3</sup>) this classical system and show that the wave operators are the dual Radon transform on the symmetric space corresponding to different choices of the Weyl chambers. If the group is of rank 1, there are only two Weyl chambers, we can label the resulting wave operators as incoming and outgoing, and the scattering operator can be calculated accordingly. The calculation is the same in the higher ranking cases, their physical significance still remains to be explored.

The dual Radon transform (wave operator) operates from the set of functions on the dual space (free states) to functions on the symmetric space (interacting states). This seems to suggest that the free states are expressed more conveniently as functions on the dual space, rather than the Euclidean space. This observation leads us to consider an alternative to the Euclidean connection procedure. Our objective here is to show that the  $S$ -matrix elements can be found by comparing the asymptotic states on the symmetric space and the free states on the dual space, a method similar to the one employed by Frank and Wolf.<sup>4</sup> Moreover, the comparison has been performed for the spherically symmetric functions

[the states  $|k,0\rangle$  in the  $SU(1,1)$  case or, in the Coulomb case, the states that have only radial dependence]. The results are the celebrated Harish-Chandra  $c$ -functions, and have been computed for all semisimple Lie groups by Gindikin and Karpelevic.<sup>5</sup>

This paper is organized as follows. In Sec. II we present the formal setup for this comparison, and we show that this is a generalization of the partial wave analysis. If we adhere to this formalism, there is a constraint on the dimension of the phase space relative to the dimension of the various groups and subgroups in question. In the cases when the dimensional constraints are satisfied, the calculation is straightforward. Examples for the nice dimensions are the Pöschl–Teller potential and the Coulomb potential with symmetry group  $SO(1,3)$ . These will be done in Sec. III. In the other cases, we need to know more about the geometry of the symmetry group so as to choose the state spaces with the right dimensions. In Sec. IV we give an example for the group  $SO(2,3)$  as the symmetry group for a perturbed Coulomb potential. The procedure is analogous, we hope to report on the formal principle in the future.

## II. PARTIAL WAVE ANALYSIS AND RADON TRANSFORM

The relevant details on Lie group theory can be found in Helgason.<sup>6</sup>

Let  $G$  be a semisimple Lie group of noncompact type with finite center,  $K$  a maximal compact subgroup. Fix a Cartan decomposition  $\mathfrak{g} = \mathfrak{k} + \mathfrak{p}$  of the Lie algebra  $\mathfrak{g}$ , where the analytic subgroup of  $\mathfrak{k}$  is  $K$ . Let  $\mathfrak{a}$  be the maximal Abelian subalgebra of  $\mathfrak{p}$ ,  $\mathfrak{a}^+$  the Weyl chamber. Let  $A = \exp \mathfrak{a}$ ,  $A^+ = \exp \mathfrak{a}^+$  be the analytic subgroups of  $\mathfrak{a}$  and  $\mathfrak{a}^+$  in  $G$ , respectively. Let  $M, M'$  denote the centralizer and normalizer of  $A$  in  $K$ , then  $B = K/M$  can be viewed as the boundary of the symmetric space  $X = G/K$ . Let  $G = KAN$  be the Iwasawa decomposition of  $G$  corresponding to our previous choices. Define a function  $H: X \times B \rightarrow \mathfrak{a}$ ,  $H(x,b) = \alpha$  if  $X = b \exp \mathfrak{a} n K$  for some  $n \in N$ . Let  $\sigma$  be one-half the sum of positive roots,  $\mu \in \mathfrak{a}^*$ , the dual vector space, then  $\exp(i\mu - \sigma, H(x,b))$  are the plane waves. We get the plane wave decomposition of functions on  $X$ ,

$$F^{\sim}(b, \mu) = \int_x f(x) \exp(i\mu - \sigma, H(x,b)) dx. \quad (1)$$

Here  $F^-(\mu)$ , which is considered as a function on  $B$ , can be decomposed further once we have a representation of the compact group  $K$ . This allows us to label the states as  $(\mu_1, \mu_2, \dots, m_1, \dots, m_2, \dots)$ , where  $m_i$ 's are discrete corresponding to the state labeling of  $K$ , and  $\mu_i$ 's are continuous as components of  $\mu \in \mathfrak{a}^*$ .

Closely related to the plane wave decomposition of functions on  $X$ , we have the Radon transform, first discussed by Newton<sup>7</sup> in the Euclidean case and was shown to have the same significance in all symmetric spaces of the form  $G/K$  (see Ref. 2). There the integration is over horocycles  $C(b, \alpha) = \{x \in X | H(x, b) = \alpha\}$  or

$$F(b, \alpha) = \int_x f(x) \delta(H(x, b) - \alpha) dx. \quad (2)$$

The set of horocycles are parameterized by  $\Sigma = G/MN$ , where  $C(b, \alpha)$  is given as  $b \exp \alpha MN \in \Sigma$ , where  $\Sigma$  is known as the dual space of  $X$ . We can rewrite (2) as

$$F(b \exp \alpha MN) = \int f(b \exp \alpha nK) dn. \quad (3)$$

This is the Radon transform.

Physically,  $b$  is the outgoing asymptotic "direction,"  $\alpha$  is the dual coordinate to the "energy"  $\mu$ . So  $\Sigma$  is the space of asymptotic data, and (3) can be interpreted as the inverse of the scattering (outgoing) wave operator. In the incoming direction, the calculation is the same as above except we choose the opposite Weyl chamber  $\mathfrak{a}^- = -\mathfrak{a}^+$ .

Let  $L^2(\Sigma)$  and  $L^2(X)$  be carrier spaces of representations of  $G$ . Denote  $|\pm \mu_i, m_i\rangle_\Sigma$  and  $|\mu_i, m_i\rangle_X$  the states in the two spaces with the given label. We can consider them as the free and interacting states of given asymptotic data, since  $|\pm \mu_i, m_i\rangle_\Sigma$  are the Radon transforms of  $|\mu_i, m_i\rangle_X$  with respect to  $\mathfrak{a}^+$  and  $\mathfrak{a}^-$ . The transmission and reflection coefficients are given by comparing these states. This comparison is well defined via the polar coordinates on  $X$  and  $\Sigma$ :

$$X = K/M \times A^+, \quad \Sigma = K/M \times A. \quad (4)$$

We now have asymptotically

$$|\mu, m\rangle_X \approx A |\mu, m\rangle_\Sigma + B |-\mu, m\rangle_\Sigma. \quad (5)$$

For the  $K$  invariant states,  $A = c(\mu)$  is the Harish-Chandra  $c$  function. In this sense, the transmission and reflection coefficients are extensions to the  $c$  functions, corresponding to choosing  $\mathfrak{a}^+$  and  $\mathfrak{a}^-$ , respectively. Furthermore, this comparison is just a generalization of the partial wave analysis. In the standard case,  $G$  is of rank 1 (rank  $G = \dim \mathfrak{a}$ ),  $\mu$  is interpreted as the energy,  $K \approx \text{SO}(3)$ ,  $M \approx \text{SO}(2)$ , so  $B \approx S^2$ . The explicit formulas calculated by Gindikin and Karpelevic will suggest a functional form of the matrix elements.

If the symmetry group is of rank 1, then there is an essentially unique Casimir operator. This operator plays the role of the quantum Hamiltonian for the dynamical system. In the case where the group is of higher rank, we have more than one Casimir operator, and the correspondence between the quantum observables and the Lie algebra representation is unclear. Different correspondence will result in completely different dynamical systems. For instance, we know that the group  $\text{SU}(1,1)$  is the symmetry for a particle in the Pöschl-Teller potential, so if we have two noninteracting

particles moving in this potential, the symmetry group for elastic scattering is  $\text{SU}(1,1) \times \text{SU}(1,1) \approx \text{SO}(2,2)$ . This is of course different from the case considered by Alhassid *et al.* in which one particle is moving in a two-dimensional space, with a modified Pöschl-Teller potential.

If the space  $G/K$  is appropriate, then the procedure as described will give correct  $S$ -matrix elements. Examples include the Pöschl-Teller potential with  $\text{SU}(1,1) \approx \text{SO}(1,2)$ , Coulomb potential with  $\text{SO}(1,3)$ , and the two noninteracting particles with  $\text{SO}(2,2)$ . We will also consider the group  $\text{SO}(2,3)$  as the symmetry group of heavy-ion scattering, here the isotropy group is not the maximal compact subgroup. The procedure is analogous, it turns out that  $X$  is the coset space  $\text{SO}(2,3)/\text{SO}(1,3)$ . We note here that Radon transform has been generalized to these situations (known as double fibrations in Helgason<sup>6</sup>), although not as complete.

### III. APPLICATIONS TO SCATTERING

#### A. Pöschl-Teller potential

Here we assume that the dynamical symmetry decomposes according to the chain  $\text{SO}(1,2) \supset \text{SO}(2)$ . The symmetric space  $X = \text{SO}(1,2)/\text{SO}(2)$  can be realized as a hyperboloid. Let us take the upper sheet of the two-sheeted hyperboloid

$$X = \{x \in \mathbb{R}^3 | [x, x] = 1, x_0 \geq 0\}, \quad (6)$$

where  $x = (x_0, x_1, x_2)$ ,  $[ \ , \ ]$  denote the pseudometric  $(+ \ - \ -)$ . The dual space is the cone

$$\Sigma = \{y \in \mathbb{R}^3 | [y, y] = 0, y_0 \geq 0\}. \quad (7)$$

Since  $\text{SO}(1,2)$  acts on  $X$  and  $\Sigma$ , we define interacting states and free states as eigenfunctions to the Casimir operator of  $\text{SO}(1,2)$  on  $X$  and  $\Sigma$ , respectively. As indicated earlier, we will perform a comparison between the states on  $X$  and  $\Sigma$  corresponding to the same eigenvalues for the chain  $\text{SO}(1,2) \supset \text{SO}(2)$ . Denote  $|\mu, m\rangle_X$  and  $|\mu, m\rangle_\Sigma$  as the states on  $X$  and  $\Sigma$ , respectively, we want to determine the transmission and reflection coefficients  $A$  and  $B$  such that

$$|\mu, m\rangle_X \approx A |+\mu, m\rangle_\Sigma + B |-\mu, m\rangle_\Sigma. \quad (8)$$

Choosing coordinate  $x_0 = \cosh \alpha$ ,  $x_1 = \sinh \alpha \cos \theta$ ,  $x_2 = \sinh \alpha \sin \theta$ , the Casimir operator takes the form

$$C_x^2 = -\frac{1}{\sinh \alpha} \frac{\partial}{\partial \alpha} \sinh \alpha \frac{\partial}{\partial \alpha} - \frac{1}{\sinh^2 \alpha} \frac{\partial^2}{\partial \theta^2}, \quad (9)$$

and

$$\begin{aligned} |\mu, m\rangle_X &= u_m^\mu(\alpha) e^{im\theta}, \\ u_m^\mu(\alpha) &= P_{i\mu - 1/2}^m(\cosh \alpha), \\ C_x^2 |\mu, m\rangle_X &= (-\mu^2 - 1/4) |\mu, m\rangle_X, \end{aligned} \quad (10)$$

where  $P$  denotes the conical functions.<sup>8</sup> Similarly, choose coordinate on  $\Sigma$  as  $y_0 = e^\alpha$ ,  $y_1 = e^\alpha \cos \theta$ ,  $y_2 = e^\alpha \sin \theta$ , then the Casimir operator is

$$C_\Sigma^2 = \frac{d^2}{d\alpha^2} + \frac{d}{d\alpha}. \quad (11)$$

The eigenfunction are of the form  $e^\sigma f(\theta)$  with eigenvalue  $\sigma^2 + \sigma$ . If  $\sigma = \pm i\mu - \frac{1}{2}$ ,  $\sigma^2 + \sigma = \mu^2 - \frac{1}{4}$ , we set

$$|\pm \mu, m\rangle_{\Sigma} = e^{(\pm i\mu - 1/2)\alpha} e^{im\theta} \quad (12)$$

The coordinates are chosen so that the comparison is performed under the limit as  $\alpha$  approaches infinity. One need only to get the asymptotic expansion of  $P^m_{i\mu-1/2}(\cosh \alpha)$ . We have<sup>8</sup>

$$(2\pi \sinh \alpha)^{1/2} P^m_{i\mu-1/2}(\cosh \alpha) \approx e^{i\mu\alpha} \frac{\Gamma(i\mu)}{\Gamma(\frac{1}{2} + m + i\mu)} + e^{-i\mu\alpha} \frac{\Gamma(-i\mu)}{\Gamma(\frac{1}{2} + m - i\mu)} \quad (13)$$

The  $S$ -matrix elements are

$$\frac{\Gamma(\frac{1}{2} + m - i\mu)}{\Gamma(\frac{1}{2} + m + i\mu)} \cdot \frac{\Gamma(i\mu)}{\Gamma(-i\mu)} \quad (14)$$

## B. Coulomb potential

Let  $X = \{x \in R^4 | [x, x] = 1, x_0 \geq 1\}$  be the upper sheet of the unit hyperboloid on which the symmetry group of the Coulomb potential  $SO(1,3)$  acts naturally,  $[ , ]$  is the metric  $+ - - -$ ,  $X = SO(1,3)/SO(3)$ . Choose coordinates

$$\begin{aligned} x_0 &= \cosh \alpha, \\ (x_1, x_2, x_3) &= \sinh \alpha \Omega, \quad \Omega \in S^2 \text{ in } R^3. \end{aligned} \quad (15)$$

Via a similarity transform  $1/\sinh \alpha$ , the Casimir operator takes the form

$$-\frac{\partial^2}{\partial \alpha^2} + 1 - \frac{L}{\sinh^2 \alpha} \quad (16)$$

(where  $L$  is the Laplace operator on  $S^2$ ), the eigenvalues are  $\mu^2 + 1$ . So the eigenfunctions satisfy the equation

$$\left(-\frac{\partial^2}{\partial \alpha^2} + \frac{l(l+1)}{\sinh^2 \alpha}\right) \phi = \mu^2 \cdot \phi, \quad (17)$$

where  $\phi(\alpha, \Omega) = \sinh \alpha \cdot P^{-l-1/2}_{i\mu-1/2}(\cosh \alpha) \cdot Y_{lm}(\Omega)$ , the  $P$  are the conical function and the  $Y$  are the spherical harmonics.<sup>9</sup> The work of Ref. 9 can now be reinterpreted as a partial wave analysis described above. Notice that in both these cases, via a similarity transform, the Casimir operator on  $\Sigma$  becomes simply  $\partial^2/\partial \alpha^2$ . We suspect that this is true, in general, thus we always compare the interacting states with  $e^{i\mu\alpha}$ .

## IV. HEAVY-ION SCATTERING

It was reported<sup>1</sup> that heavy-ion scattering possesses an  $SO(2,3)$  symmetry. Here we will show that this problem can be recognized as a short range perturbation on the Coulomb interaction, and that the group  $SO(2,3)$  enters naturally. The comparison procedure described above is applied and the results are consistent with those obtained by Alhassid *et al.*

Let  $X = \{x \in R^5 | [x, x] = 1\}$ ,  $[ , ]$  is the metric  $+ + - - -$ ,  $X = SO(2,3)/SO(1,3)$ . We parametrize  $X$  as

$$\begin{aligned} (x_1, x_2) &= \cosh \alpha \cdot w, \\ (x_3, x_4, x_5) &= \sinh \alpha \cdot \Omega, \end{aligned} \quad (18)$$

with  $w \in S^1$  in  $R^2$ ,  $\Omega \in S^2$  as before. Via the similarity transform  $1/\sinh \alpha \cdot \cosh^{1/2} \alpha$ , using the fact that the eigenvalue is  $\mu^2 + \frac{9}{8}$ , the Laplace equation on  $X$  reads

$$C_{l,k\Phi} = \left(-\frac{\partial^2}{\partial \alpha^2} - \frac{l(l+1)}{\sinh^2 \alpha} - \frac{k^2 - \frac{1}{4}}{\cosh^2 \alpha}\right) \Phi = \mu^2 \cdot \Phi, \quad (19)$$

where  $-k^2$  is the eigenvalue of the operator  $\partial^2/\partial w^2$ . One sees that this equation is a perturbation of (17), the last term corresponds to a potential of the form  $e^{-r}/r^2$  (see the Appendix). Notice that if  $k = \frac{1}{2}$ , Eq. (19) becomes that of the Coulomb problem. Our procedure suggests that we should compare the eigenfunctions with eigenfunctions on the asymptotic cone  $\Sigma = \{y \in R^5 | [y, y] = 0\}$ .

From Ref. 1, the  $S$ -matrix elements are

$$S(\mu, l, k) = \frac{\Gamma(\frac{1}{2}(l+k+\frac{3}{2}+i\mu))\Gamma(\frac{1}{2}(l-k+\frac{3}{2}+i\mu))}{\Gamma(\frac{1}{2}(l+k+\frac{3}{2}-i\mu))\Gamma(\frac{1}{2}(l-k+\frac{3}{2}-i\mu))} \quad (20)$$

Note that when  $k = \frac{1}{2}$ , this reduces to the Coulomb  $S$  matrix. So our comparison trick will give the same result in these cases. Equation (20) yields the following recursive relation:

$$\begin{aligned} S(\mu, l-1, k-1) &= [(l+k-\frac{1}{2}+i\mu)/(l+k-\frac{1}{2}-i\mu)] S(\mu, l, k). \end{aligned} \quad (21)$$

So we need to show that our comparison gives the same relations.

Let

$$J = \frac{\partial}{\partial \alpha} + l \cdot \coth \alpha + (k - \frac{1}{2}) \cdot \tanh \alpha, \quad (22)$$

then we have,  $C$  as in (19),

$$C_{(l-1), (k-1)} J |\mu, l, k\rangle_X = J C_{l,k} |\mu, l, k\rangle_X. \quad (23)$$

Thus  $J$  plays the role of a shift operator, whether  $J$  belongs to the  $SO(2,3)$  algebra is immaterial. So asymptotically,  $\text{const} \cdot |\mu, l-1, k-1\rangle$

$$\begin{aligned} &= J |\mu, l, k\rangle \approx \left(\frac{\partial}{\partial \alpha} + l + k - \frac{1}{2}\right) \cdot (A \cdot e^{i\mu\alpha} + B \cdot e^{-i\mu\alpha}) \\ &= (i\mu + l + k - \frac{1}{2}) A \cdot e^{i\mu\alpha} \\ &\quad + (-i\mu + l + k - \frac{1}{2}) B \cdot e^{-i\mu\alpha}, \end{aligned} \quad (24)$$

and we do get the same recursive relations.

In conclusion, we have embarked on a geometric theory, which complements the algebraic theory of Alhassid *et al.*, for scattering in the presence of a group symmetry. We see that this theory is an extension to the partial wave analysis, and that the geometric consideration simplifies the calculation.

## APPENDIX: CANONICAL TRANSFORM FOR POSITIVE COULOMB SPACE AND THE HYPERBOLOID

Here we give explicitly the  $SO(1,3)$  equivariant canonical transform between the positive energy phase space of the Coulomb problem and the phase space of the geodesic flow problem on the hyperboloid. This allows us to estimate the size of the perturbation term in (19). The construction here is similar to the one given in Souriau and Onofri<sup>10</sup> for the negative energy case.

Let  $T^*_+R^3$  denote the subspace of the cotangent bundle on  $R^3$  on which the energy  $E = p^2/2 - 1/q$  is positive. Denote  $T^*H$  the cotangent bundle on the hyperboloid



$H = \{y \in \mathbb{R}^4 \mid [y, y] = 1, y_0 \geq 1\}$ . For convenience, we can (via a reduction<sup>11</sup>) identify  $T^*H$  as the subspace  $\{(y, \eta) \in \mathbb{R}^4 \times \mathbb{R}^4 \mid [y, y] = 1, [y, \eta] = 0\}$ . The symplectic form is  $d\eta_0 \wedge dy_0 - \sum d\eta_i \wedge dy_i$ , so that the  $SO(1,3)$  actions are canonical. Then the canonical transform  $T^*H \rightarrow T^*\mathbb{R}^3$  is

$$\begin{aligned} y &= A \cosh t + (B/\sqrt{-[B, B]}) \sinh t, \\ \eta &= \sqrt{-[B, B]} \cdot A \cdot \sinh t + B \cosh t, \end{aligned} \tag{A1}$$

where

$$\begin{aligned} t &= \langle q, p \rangle \sqrt{2E}, \\ A &= \begin{bmatrix} p^2 q - 1 \\ \sqrt{2E} \cdot q \cdot p_i \end{bmatrix}, \quad B = \left[ \frac{\langle q, p \rangle}{(\langle q, p \rangle p_i - q_i/q) / \sqrt{2E}} \right]. \end{aligned} \tag{A2}$$

The perturbation term  $(k^2 - \frac{1}{4})/y_0^2$ , both near and far away from the scattering center, assumes the form

$$(k^2 - \frac{1}{4}) \exp(-4Eq) / (2Eq)^2. \tag{A3}$$

Thus we see that this term is a short range potential which grows as  $1/q^2$  at the origin.

<sup>1</sup>Y. Alhassid, F. Iachello, and J. Wu, *Phys. Rev. Lett.* **56**, 271 (1986); Y. Alhassid, F. Gürsey, and F. Iachello, *ibid.* **50**, 873 (1983); *Ann. Phys. (NY)* **148**, 346 (1983); Yale preprint No. 3074-837.

<sup>2</sup>Y. Wu, *J. Phys. A* **20**, 429 (1987).

<sup>3</sup>B. Kostant, *Quantization and Unitary Representation*, in *Lecture Notes in Mathematics*, Vol. 170 (Springer, New York, 1970); J. Sniatycki, *Geometric Quantization and Quantum Mechanics* (Springer, New York, 1980).

<sup>4</sup>A. Frank and K. B. Wolf, *Phys. Rev. Lett.* **52**, 1737 (1984).

<sup>5</sup>S. G. Gindikin and F. I. Karpelevic, *Dok. Akad. Nauk SSSR* **145**, 252 (1962).

<sup>6</sup>S. Helgason, *Groups and Geometric Analysis* (Academic, New York, 1984).

<sup>7</sup>R. G. Newton, *Scattering Theory in the Mixed Representation*, *Lecture Notes in Physics*, Vol. 130 (Springer, New York, 1980).

<sup>8</sup>W. Magnus, F. Oberhettinger, and R. P. Soni, *Formulas and Theorems for the Special Functions of Mathematical Physics* (Springer, New York, 1966).

<sup>9</sup>Y. Wu, *J. Phys. A* **18**, L499 (1985).

<sup>10</sup>J. M. Souiau, *Sur la Variete de Kepler*, *Symposia Math* (Academic, London, 1974); E. Onofri, *J. Math. Phys.* **17**, 401 (1976); Y. Wu, *Hadronic J.* **8**, 101 (1985).

<sup>11</sup>V. W. Guillemin and S. Sternberg, *Am. J. Math.* **101**, 915 (1979).

# Existence and observability of spinor structure

Anne M. R. Magnon<sup>a)</sup>

*Département de Mathématiques, Université de Clermont-Ferrand, 63170, Aubière, France; Physics Department, University of Syracuse, Syracuse, New York 13244-1130; and Institute for Theoretical Physics, University of California, Santa Barbara, California 93106*

(Received 3 March 1986; accepted for publication 7 January 1987)

A mechanism by which space-time topological modifications could have been controlled, in the early universe or at the Planck length, to enable onset of spinor structure is investigated. This mechanism (based on a reshuffling of topological charges and related modification of characteristic classes) could provide a gravitational analog of the Aharonov–Susskind *Gedankenexperiment* proposed to detect relative rotation in the universe, spinor behavior, or to keep track of the two homotopy classes of the Lorentz Lie group. The space-time topology [and in particular the trivial (nontrivial) bundle structure at conformal null infinity] provide a labeling of the asymptotic Lorentz homotopy classes which originates in the first Chern class (enclosed magnetic mass) or in the parametrization of the second homology group, and gives rise to a necessary (and sufficient) condition for the existence of spinor structure. This underlines the intertwined roles of topology and curvature. The mechanism could also be viewed as an “unwinding” of gravitational magnetic monopoles with one asymptotic region into electric mass (black-hole) solutions with two asymptotic regions. In such situations a discrete PT symmetry could emerge from a continuous transformation. Possible implications on the CPT theorem are mentioned.

## I. INTRODUCTION

Since the discovery, in 1956, of parity nonconservation in weak interactions, it is believed that the concepts of right and left physical systems can be unambiguously defined provided space is orientable. Indeed such a statement could not be made if a closed circuit within a disoriented laboratory could transform a particle into an antiparticle: in this case the combined transformations, charge conjugation (C) and space inversion (P), would not enable us to define the conservation of right–left symmetry in space. Equivalently, in a disoriented three-space, space inversion is not a discrete transformation but a continuous one<sup>1</sup> and CP invariance does not enable us to decide whether two remote particles are identical or are a particle and an antiparticle since the conclusion would depend on the path along which the two particles are brought to the same point. However, the absolute difference between matter and antimatter has been experimentally confirmed by the CP violation in the  $K^0$ -meson decay, and nonorientability of space seems to be excluded by the same token.<sup>2,3</sup> It is to be underlined here that such conclusions are drawn under the requirement that causality violation is excluded, i.e., that space-time does not contain closed timelike world lines.

Since one might question the availability of causality, e.g., in the early universe or on a cosmic scale (and we shall do so here), it is of interest to search for mechanisms that could explain the onset of orientability or even bring support to a further question: Why should there be orientability in the universe or PCT invariance?

It has often been underlined that this issue is related to that of observability of a relative  $2\pi$  rotation of two systems (certainly a  $2\pi$  rotation of the entire universe should not be observable), and to the existence of spin structure.<sup>4–7</sup>

Underlining a key point—a spinor changes sign when a basis completes a  $2\pi$  rotation, returning to its original position after  $4\pi$  rotation—Penrose has suggested<sup>5</sup> a mechanism (based on a *Gedankenexperiment* proposed by Aharonov and Susskind<sup>8</sup>) to test spin structure. The Aharonov–Susskind apparatus is a box with perfectly reflecting walls which is divided into two identical compartments (by an impenetrable partition which may be open or closed off by a shutter) and which contains an electron. A relative (quasistatic)  $2\pi$  rotation between the two compartments can be detected if the electron wave function in box 1 is allowed to produce an interference pattern with the wave function in box 2 after uniform equal magnetic fields have been applied to the disconnected compartments in the direction of the spin of the electron. The relation with spinor structure emerges from the following argument. The physically measurable properties of a spinor  $\psi^A$  can be determined from its corresponding null bivector (null flag). If  $\psi^A$  undergoes a (phase) rotation  $\psi^A \rightarrow e^{i\theta} \psi^A$ , it is clear that  $\psi^A$  changes sign; it takes a  $2\pi$  rotation to bring it back to its initial value, in this rotation the null flag undergoes a  $4\pi$  rotation around its attached null direction  $k^a = \psi^A \psi^{A'}$ . This is related to the existence of two homotopy classes in the Lorentz Lie group  $\mathcal{L}$ —every closed path in this group being either homotopic to the  $2\pi$  rotation (element of the nontrivial homotopy class) or the  $4\pi$  rotation (the trivial class)—and to the fact that the  $SL(2, C)$  Lie group is the universal covering group of  $\mathcal{L}$ , two elements  $L_B^A$  and  $M_B^A$  of  $SL(2C)$  giving rise to the same Lorentz group element if and only if  $L_B^A = \pm M_B^A$ .

<sup>a)</sup> Détachée du Ministère des Relations Extérieures, Paris, France.

Since the definition of a spinor bundle  $\mathcal{S}$  over a curved space-time  $M$  involves the intertwining of the fundamental group of  $M$  with that of the Lorentz group, the role of topology is crucial and has been investigated by various authors.<sup>6,7,9</sup> The basic idea is to start from the principal fiber bundle  $\mathcal{P}$  of oriented-time oriented orthonormal bases, to unwrap each fiber into a principal  $SL(2, C)$  bundle  $B$  over  $M$ , a spinor structure being defined as a second principal fiber bundle  $\mathcal{S}$  over  $B$  with a 2-1 mapping  $\phi: \mathcal{S} \rightarrow B$ . If  $\mathcal{P}$  is simply connected, closed curves contained within a fiber and corresponding to the nontrivial  $2\pi$ -rotation homotopy class can be deformed outside the fiber into the trivial curve and the notion of a spinor field cannot be defined. It is known that  $\mathcal{P}$  will fail to be simply connected if the second Stiefel-Whitney class of  $M$  vanishes. If  $M$  is noncompact, it is also known<sup>6</sup> that a spinor structure exists if and only if  $M$  is parallelizable.

In presence of intricate space-time topologies,  $M$  may fail to be orientable, or time oriented. In such a case  $\mathcal{P}$  does not exist and the usual notion of spinors cannot be defined. If  $M$  is orientable and time oriented,  $\mathcal{P}$  can be unwrapped into  $B$  without unwrapping  $M$  provided  $\pi_1(\mathcal{P}) = \pi_1(M) \times \mathbb{Z}_2$ ; if this condition is not satisfied, the notion of spinor structure cannot be defined. It is also of interest to underline that the above features of the  $SL(2, C)$  representation are related to the presence at the quantum level of two irreducible representations<sup>10</sup> which are conveniently labeled by the values of the Casimir operators  $m^2$  (squared mass) and  $S^2$  (squared angular momentum) of the  $SL(2, C)$  Lie algebra. The results presented here might suggest that a canonical generalization in the context of quantum gravity could be provided by the squared mass and the squared angular momentum monopole.

These considerations are intended to motivate the viewpoint we would like to develop, i.e., the existence of spinor structure and orientability could have been controlled in the early universe or at the Planck length through an appropriate reshuffling of topological charges and related modification in the homology classes of the space-time manifold.

We shall take advantage of the availability of various Maxwellian features of gravity in presence of suitable nontrivial space-time topologies to propose a mechanism rather reminiscent of that designed by Aharonov and Susskind for the detection of spinor structure in presence of a Maxwellian magnetic field. Our analysis will be based on previous results<sup>11-13</sup> concerning the structure of magnetic mass source-free solutions to Einstein's equation.

A simplified picture of the gravitational apparatus to be investigated could be the following: the two compartments of the Aharonov-Susskind boxes are provided by the two asymptotic regions (or mirror universes) of the Kruskal-Schwarzschild black-hole solution. A reshuffling of the (electric) mass monopole into the magnetic mass monopole (which could be supported by the following geometrical interpretation: decreasing the area of the event horizon, shrinking the region of trapped surfaces, and smooth absorption into a causality violation and single asymptotic region) has been investigated in Refs. 11-13. As a result of the transformation, a nontrivial bundle structure ( $S^1$  Hopf fibering

and transition functions) is induced at conformal null infinity  $\mathcal{I}$ . One could thus draw an analogy between the above setup and the Aharonov-Susskind apparatus, where the opening of the shutter would correspond to the absorption of the "entropy reservoir" into new topological features, and the interference pattern to nontrivial bundle structure at  $\mathcal{I}$ ,  $S^1$  Hopf fibering and transition functions measuring the enclosed magnetic charge.

In the main body of this paper we shall prove that the parametrization of characteristic classes involved in the definition of magnetic mass (second homology class and first Chern class) induces a (0-1) labeling of the homotopy classes of the (asymptotic) Lorentz Lie group as well as a mechanism to keep track of their possible mingling (as one moves from fiber to fiber within the principal frame bundle) along closed space-time paths, thus providing a necessary (and in many cases sufficient) condition for the existence of spinor structure. This is strongly reminiscent of a setup proposed by Geroch<sup>6</sup> to keep track of spinor structure in the universe by moving Maxwellian-Aharonov-Susskind boxes along space-time tracks. However, the control here is purely gravitational, global, and uses topological charges (characteristic classes) which originate in the bundle structure at null infinity.

Since the above criterion is purely global and since we rely on the space-time asymptotic structure, we would like to propose the following viewpoint. If gravitational magnetic monopoles could have "unwinded" into black holes (electric mass monopoles) in the early universe, a mechanism to control onset of spinor structure would be provided.

Finally the above transformation can be shown<sup>12</sup> to be associated to an invariant (gravitational) charge  $C$ , and to the onset of a PT discrete symmetry (between two-mirror-asymptotic regions) emerging from a continuous transformation (acting on the—single—asymptotic region of a gravitational magnetic monopole solution). This leads us to comment, in our concluding remarks, on a possible (theoretical or experimental) link between the CPT theorem and the onset of spinor structure.

## II. PRELIMINARY REMARKS

Let us briefly summarize some relevant properties of the spinor representation. Recall that a spinor space is a pair  $(W, \epsilon_{AB})$ , where  $W$  is a two-dimensional vector space over  $C$  and  $\epsilon_{AB} = -\epsilon_{BA}$ , a skew-symmetric two-index tensor inducing a mapping from  $W$  into its dual  $W^*$ , hence playing the role of a nondegenerate metric. The corresponding group of transformations is that of linear mappings  $L: W \rightarrow W$  which are metric preserving:

$$\epsilon_{CD} = L^A_C L^B_D \epsilon_{AB}, \quad (1)$$

i.e., elements of  $SL(2, C)$ .

A four-dimensional complex vector space  $Y$  can be associated to  $W$ , which admits a real section  $V$  with Lorentzian metric (signature  $+, -, -, -$ ) given by

$$g_{AA'BB'} = \epsilon_{AB} \bar{\epsilon}_{A'B'}. \quad (2)$$

It can be checked that

$$\Lambda^{AA'}_{BB'} = L^A_B \bar{L}^{A'}_{B'}. \quad (3)$$

is metric preserving (a proper Lorentz transformation):

$$g_{CC'DD'} = \Lambda_{CC'}^{AA'} \Lambda_{DD'}^{BB'} g_{AA'BB'} \quad (4)$$

Furthermore  $L_B^A$  and  $M_B^A$  give rise to the same Lorentz element iff  $L_B^A = \pm M_B^A$ , reflecting the fact that  $SL(2, C)$  is the universal covering group of the Lorentz group (two-valuedness of the spin representation). It is the condition that this two-valuedness can be consistently eliminated over an entire space-time which enables the definition of global spinor fields, reflecting the presence of a spinor structure.

Since the Lorentz Lie group  $\mathcal{L}$  is doubly connected—a path representing a  $2\pi$  rotation cannot be continuously deformed to the identity whereas a  $4\pi$  rotation can be—a spin structure must be able to keep track of the presence of these two homotopy classes, providing a homomorphism from  $\pi_1(\mathcal{L})$  into  $Z_2$ , taking the value 0 on the trivial class (containing the  $4\pi$  rotation), unity on the other class. Following a geometrical argument proposed by Penrose,<sup>5</sup> we suggest that such a homomorphism could be associated to the onset of topological charges or equivalently specific homology classes. Recall that to every one-spinor  $\xi^A \in \mathcal{W}$  one can associate a bivector  $F_{AA'BB'}$  defined via

$$F_{[ab]} \leftrightarrow F_{AA'BB'} = \epsilon_{AB} \bar{\phi}_{A'B'} + \bar{\epsilon}_{A'B'} \phi_{AB}, \quad (5)$$

where  $\phi_{AB} = \xi_A \xi_B$ . This bivector (or null flag) is attached to a pole represented by the direction of the null vector  $\xi^A \xi^{A'}$ . Let us denote by  $f_\theta$  the (phase) rotation defined via

$$f_\theta(\xi^A) = e^{i\theta} \xi^A. \quad (6)$$

The induced transformation on  $F_{ab}$  is

$$f_\theta(F_{ab}) = F_{ab} \cos 2\theta - *F_{ab} \sin 2\theta. \quad (7)$$

On another hand, a basis  $(\xi^A, \eta^A, \xi \cdot \eta = 1)$  being chosen in  $\mathcal{W}$ ,  $F_{ab}$  can be expressed as

$$F_{ab} = 2\xi_{[a} w_{b]} \quad (\xi \cdot w = 0), \quad (8a)$$

the (phase) rotation inducing a rotation of the two-flat defined by

$$f_\theta w_a = w_a \cos 2\theta + v_a \sin 2\theta, \quad (8b)$$

where  $(\xi^a, w^a, v^a)$  is an orthonormal triad in Minkowski space. Hence  $f_\theta$  induces a reversal of the spinor sign ( $f_\pi \xi^A = -f_0 \xi^A$ ), and a winding of the null flag around its pole  $\xi^A$ , a  $4\pi$  rotation of this flag being required to bring back  $\xi^A$  to its original value. Thus the flag and its null pole have the availability to keep track of the double valuedness of the spin representation. Let  $S_2$  denote a two-sphere surrounding the point at spatial infinity in Minkowski space, and let us assume that  $F_{ab}$  is a (plane-wave) source-free Maxwell field. Denote by

$$C_1 = \int_{S_2} *F_{ab} dS^{ab} \quad (9)$$

its electric charge, and by

$$C_2 = \int_{S_2} F_{ab} dS^{ab} \quad (10)$$

its magnetic charge. It is straightforward to check that

$$f_\theta F_{ab} = F_{ab} \cos 2\theta - *F_{ab} \sin 2\theta, \quad (11)$$

$$f_\theta *F_{ab} = F_{ab} \sin 2\theta + *F_{ab} \cos 2\theta. \quad (12)$$

The quantity

$$C = C_1^2 + C_2^2 \quad (13)$$

is an invariant of the transformation, its vanishing (nonvanishing) being governed by the second homology group of the space-time manifold.

A gravitational analog of the transformation  $f_\theta$  and related reshuffling of charges is available. Its topological implications have been mentioned in previous papers,<sup>11-13</sup> and will be crucial when we shall derive criteria for the existence and possible observability of spinor structure. The role of spinor structure has been carefully underlined<sup>14</sup> in the investigation of the asymptotic behavior of zero rest-mass fields. Let us briefly summarize the situation.

From now onwards,  $(M, g_{ab}, \xi^a)$  will denote a solution of Einstein's equation  $R_{ab} = 0$ , with Killing vector field  $\xi^a$ . The norm and twist of  $\xi^a$  are, respectively,  $-\lambda = \xi^a \xi_a$  and  $\omega_a = \epsilon_{abcd} \xi^b \nabla^c \xi^d = \text{grad } \omega$ . Let  $h_{ab} = g_{ab} + \lambda^{-1} \xi_a \xi_b$  be the induced metric on  $T$ , the manifold of orbits of  $\xi^a$ . Assuming  $\lambda \neq 0$ , a rescaled metric  $\tilde{h}_{ab}$  and complex potential  $\tau$  can be defined<sup>15</sup> on  $T$ :

$$\tilde{h}_{ab} = \lambda h_{ab}, \quad (14)$$

$$\tau = \omega + i\lambda. \quad (15)$$

A transformation (which was initially introduced<sup>15-18</sup> to generate circle families of explicit, exact, source-free solutions to Einstein's equation with one Killing vector field, starting from one of them) is given by

$$G_\theta \tau = (\tau \cos \theta + \sin \theta) / (-\tau \sin \theta + \cos \theta). \quad (16)$$

The expression of the resulting solution  $g_{ab}(\theta)$  is given explicitly in Ref. 16. The topological charges that can emerge from the action of  $G_\theta$  are obtained from three real divergence-free vector fields  $V_i^a$  ( $i = 1, 2, 3$ ) on  $T$ , and their associated curl-free two-forms  $\epsilon_{abc} V_i^c \equiv \tilde{F}_{ab}^i$  ( $i = 1, 2, 3$ ) on  $T$ . These forms admit the following pullbacks on  $(M, g_{ab})$ :

$$F_{ab}^1 = \nabla_{[a} \lambda^{-1} \xi_{b]}, \quad (17)$$

$$F_{ab}^2 = \nabla_{[a} \omega \lambda^{-1} \xi_{b]} - \frac{1}{2} \epsilon_{abcd} \nabla^c \xi^d, \quad (18)$$

$$F_{ab}^3 = \nabla_{[a} (\lambda^{-1} (\omega^2 + \lambda^2) \xi_{b]} - 2\lambda \nabla_a \xi_b - \omega \epsilon_{abcd} \nabla^c \xi^d. \quad (19)$$

Integrating  $\tilde{F}_{ab}^i$  on a two-sphere  $S_2^\infty$  surrounding the point at spacelike infinity on  $T$  leads to various conserved quantities:

$$Q_i = \int_{S_2^\infty} \tilde{F}_{ab}^i dS^{ab}, \quad i = 1, 2, 3. \quad (20)$$

Furthermore, under the action of the circle group, these charges are reshuffled according to

$$G_\theta V_1^a = V_1^a \cos^2 \theta - V_2^a \sin 2\theta + V_3^a \sin^2 \theta, \quad (21)$$

$$G_\theta V_2^a = \frac{1}{2} V_1^a \sin 2\theta + V_2^a \cos 2\theta - \frac{1}{2} V_3^a \sin 2\theta, \quad (22)$$

$$G_\theta V_3^a = V_1^a \sin^2 \theta + V_2^a \sin 2\theta + V_3^a \cos^2 \theta. \quad (23)$$

It is easily checked that  $Q = Q_1 Q_3 - Q_2^2$  is an invariant of the transformation  $G_\theta$ , a gravitational analog of the charge  $C$ .

The reshuffling of  $Q_i$ 's under the action of the cyclic group (phase rotation) can be related to a modification in the structure of characteristic classes and consequently decide on the existence of spinor structure. This will be investigated in the following sections.

### III. MAGNETIC MASS AND CHARACTERISTIC CLASSES

Let us briefly review a definition of the magnetic mass based on various results<sup>11-13</sup> concerning the asymptotic structure of the NUT (gravitational magnetic monopole) source free solution to Einstein's equation. For the exact NUT solution it has been shown elsewhere<sup>11</sup> that this conserved quantity reduces to the NUT parameter (angular momentum monopole).

An asymptotically NUT gravitational magnetic monopole source-free solution to Einstein's equation will be a space-time with conformal null boundary  $\mathcal{S}$  exhibiting the topology of an  $S^1$  (Hopf fibering) bundle over  $S^\infty$ , the two-sphere of orbits of  $n^a$  (the null normal to  $\mathcal{S}$ ). The structure of this bundle has been analyzed in Refs. 11-13. The key point here is the availability (and expression) of the bundle curvature two-form  $\Omega_{ab}$  defined by the pullback to  $\mathcal{S}$  of

$$\hat{F}_{ab}^1 = \Omega^{-1} * C^m_{c^n d} n^c n^d l_m \epsilon_{nab}, \quad (24)$$

where  $\Omega^{-1} C_{abcd}$  denotes the rescaled Weyl tensor,  $(l^a, n^a, m^a, \bar{m}^a)$  is the usual Newman-Penrose null tetrad in the neighborhood of  $\mathcal{S}$  ( $l \cdot n = -1$ ), and  $\epsilon_{abc} = \epsilon_{abcd} l^d$ ,  $\epsilon^{abc} = \epsilon^{abcd} n_d$ . We know also from Ref. 13 that  $\Omega_{ab}$  is related to the unphysical Riemann curvature via

$$\Omega_{ab} = D_{[a} S_{b]} l_c, \quad (25)$$

where  $S_{ab} = R_{ab} - \frac{1}{6} R g_{ab}$ .

If a timelike Killing vector field is available on the physical space-time, we also know<sup>13</sup> that  $\Omega_{ab}$  is (essentially) the pullback to  $\mathcal{S}$  of  $F_{ab}^1 = \nabla_{[a} \lambda^{-1} \xi_{b]}$ .

Finally  $\Omega_{ab}$  is the lift to  $\mathcal{S}$  of a closed two-form  $\tilde{\Omega}_{ab}$  on  $S^\infty$ , an element of  $H^2_\infty(M)$ , second homology group (in the asymptotic region). The fact that this  $\tilde{\Omega}_{ab}$  is not globally exact on  $S^\infty$  (discontinuities in its potential) gives rise to the magnetic mass.

**Definition:** The magnetic mass of an asymptotically NUT (source-free) gravitational magnetic monopole solution  $M$ , is defined as the value of an element  $\tilde{\Omega}_{ab}$  of  $H^2_\infty(M)$ —second homology group<sup>19</sup> of the asymptotic region of  $M$ —on the two-chain  $S^\infty$ .

**Corollary 1:** The magnetic mass can be conveniently labeled by

$$C_1 = \int_{S^\infty} \tilde{\Omega}_{ab} dS^{ab},$$

the first Chern class of the  $\mathcal{S}$  bundle. The parametrization of  $H^2(M, \mathbb{R})$  provides the (integer) number of twists of  $\mathcal{S}$  around its  $S^1$  Hopf fiber.

**Corollary 2:** If  $H^2(M, \mathbb{R})$  is trivial,  $\tilde{\Omega}_{ab}$  is a closed and globally exact two-form: the magnetic mass vanishes, implying that  $\mathcal{S}$  is a trivial  $S^2 \times R$  bundle (zero twist).

As we shall see, in the next section, the value of the magnetic mass provides a homomorphism from  $H^2(M)$  into the first homotopy group of the asymptotic Lorentz Lie

group which enables us to keep track of the homotopy classes as one moves within a (possibly multiply connected) space-time. This is precisely what is expected from a criterion for the existence of spinor structure.

### IV. CONFORMAL NULL BOUNDARY AND AHARONOV-SUSSKIND APPARATUS

We want to show that the available structure at conformal null infinity (for the source-free nonradiative solutions considered in the previous section) could be used to design a device for the detection of spinor structure. As was already mentioned in the Introduction, the resulting setup could be viewed as a gravitational analog of that proposed by Aharonov-Susskind.

Recall the existence of an exact sequence of homotopy groups<sup>20-23</sup>:

$$\begin{aligned} \pi_n(M, x) &\xrightarrow{h_n} \pi_{n-1}(G_x) \xrightarrow{i_{n-1}} \pi_{n-1}(G, u) \\ &\rightarrow \pi_{n-1}(M, x) \rightarrow \cdots \rightarrow \pi_0(M, x) \rightarrow 0, \end{aligned}$$

where  $M$  is the base space of a principal bundle with fiber  $G$ ,  $x$  a point of  $M$ ,  $G_x$  the corresponding fiber, and each homomorphism in the sequence is the kernel of the next. [This sequence reflects in particular the intertwining of the fundamental group of  $G, (\pi_1(G))$  with that of the base space  $(\pi_1(M))$ .] If  $M$  is the space-time manifold and  $G = \mathcal{L}$  (the general Lorentz Lie group),  $h_2 \equiv 0$  is a necessary and sufficient condition for  $i_1$  to be a one to one mapping, i.e., for the existence<sup>6</sup> on (a noncompact)  $M$  of a global system of (orthonormal) tetrads. This in turn provides a necessary and, in many cases,<sup>9</sup> sufficient condition for the existence of spinor structure. Since a spinor structure is associated to a homomorphism from  $\pi_1(L)$  into  $Z_2$ , we shall search for a mechanism to detect a charge taking the value zero or unity after a physical object (e.g., an orthonormal tetrad) has been transported around closed paths in  $\mathcal{L}$  or  $M$ . If a spinor structure is to be available, one expects such a mechanism to be associated to a mapping from  $\pi_2(M, x)$  to the trivial homotopy class of  $\mathcal{L}_x - x$  chosen in the asymptotic region. Along these lines, an argument has already been proposed,<sup>7</sup> which relates the space-time curvature in the neighborhood of a point  $p$  to an upper bound on the length of those closed paths in  $\mathcal{L}_p$ , which should be contractible to a point. Since the topology of the underlying space-time manifold determines the existence of spinor structure, sufficiently "twisted" manifolds should be such that any metric that could be defined on them would contain a minimum amount of curvature that would prevent the existence of spinor structure. We shall show that the space-time topology (second homology class) or equivalently the bundle structure available at conformal null infinity, not only controls this amount of curvature but provides a necessary (and in many cases sufficient) condition for the existence of spinor structure. Our result will be the following: the nonvanishing of the first Chern class—a measurement of the magnetic mass enclosed within an asymptotic region, as related to a parametrization of the second homology class of the space-time manifold—provides a nonzero mapping  $h_2: \pi_2(M, x) \rightarrow \pi_1(\mathcal{L}_x)$ , thus preventing the existence of a one-to-one mapping between the

homotopy classes of the fibers of the frame bundle, and consequently preventing the existence of a spinor structure. If the magnetic mass vanishes (trivial bundle structure at conformal null infinity),  $h_2 = 0$ ,  $i_2$  is one to one, and, with the exception of situations involving a simply connected principal frame bundle, a spinor structure will be available.

Let  $p$  be a point in the neighborhood of infinity,  $\mathcal{L}_p$  the Lorentz Lie group at  $p$ , and  $\mathcal{C} \subset \mathcal{L}_p$ , a closed path through  $p$ , i.e., a one-parameter family  $R_b^a(s)$ ,  $0 \leq s \leq 1$  of Lorentz rotations at  $p$ . We want to evaluate the length of  $\mathcal{C}$ ; thus we need to introduce a metric on the rotation group. Choose in  $M$  a spacelike two-sphere  $S$  through  $p$ , which surrounds the point at spacelike infinity. Denote by  $\gamma_s(u)$ ,  $0 \leq s \leq 1$ ,  $\gamma_s(0) = \gamma_s(1) = p$ , a one-parameter family of loops on  $S$ , originating and ending at  $p$ , and spanning  $S$  (all points of  $S$ , except  $p$ , lie on exactly one curve). Let  $t^a$  denote an arbitrary unit timelike vector field on  $M$  ( $t \cdot t = -1$ ), adjust  $S$  so that  $t^a$  and  $S$  are orthogonal at  $p$ , and introduce at  $p$  a triad  $x_\lambda^a$ ,  $\lambda = 1, 2, 3$ , which, together with  $t^a$ , defines an orthonormal tetrad. For each value of  $s$ ,  $0 < s < 1$ , we shall (Fermi) transport the tetrad around  $\gamma_s$  according to

$$u^m \nabla_m x_\lambda^a = + t^a (t_c u^m \nabla_m x_\lambda^c), \quad \lambda = 1, 2, 3, \quad (26)$$

where  $u^m$  denotes the tangent to  $\gamma_s(u)$  ( $u^m \nabla_m u = 1$ ). In this transport the tetrad remains an orthonormal tetrad although (in general) after a complete tour on  $\gamma_s$ ,  $x_\lambda^a$  ( $\lambda = 1, 2, 3$ ) might not be tangential to  $S$ . As a result of this transport, we define at  $p$  an element  $R_b^a(s)$  of the Lorentz Lie group  $\mathcal{L}_p$ :  $x_\lambda^a|_{u=1} = R_b^a(s) (x_\lambda^b|_{u=0})$ . As  $s$  varies from 0 to 1,  $\gamma_s$  spans the two sphere  $S$  and  $R_b^a(s)$  ( $0 \leq s \leq 1$ ) describes a closed path  $\mathcal{C}$  in  $\mathcal{L}_p$  [ $R_b^a(0) = R_b^a(1) = \text{identity}$ ]. We want to evaluate the length of  $\mathcal{C}$  and decide on its homotopy class. If a spinor structure is to exist, this class should be the trivial one (connected to the identity), (i.e.,  $\mathcal{C}$  should be contractible to a point). Here  $\tau_b^a = (d/ds)R_b^a(s)$ , the tangent to  $\mathcal{C}$ , is associated to an infinitesimal rotation (a generator  $F_{[ab]}$  of the Lorentz Lie algebra) and its length is proportional to the amount of curvature enclosed on  $S$ , within the two-loops  $\gamma_s$  and  $\gamma_{s+ds}$ . Since  $S$  is (in the asymptotic region) surrounding the point at spacelike infinity, the leading term is provided by

$$\int_\sigma \Omega_{ab} dS^{ab}, \quad (27)$$

where  $\Omega_{ab}$  is the two-form, introduced in Sec. III, whose pullback to  $\mathcal{S}$  provides the bundle curvature two-form, and where  $\sigma$  is, on  $S^\infty$ , the area enclosed within  $\gamma_s^+$  and  $\gamma_{s+ds}^+$ . (Here  $S^\infty$  is the two-sphere of null generators of  $\mathcal{S}$ .) Thus the total length of  $\mathcal{C}$  is given by

$$C_1 = \left| \int_{S^\infty} \Omega_{ab} dS^{ab} \right|, \quad (28)$$

i.e., the first Chern class of the  $\mathcal{S}$  bundle (or equivalently the enclosed magnetic mass). A convenient labeling of  $C_1$  is provided by the second homology class of  $M$ , hence the following theorems.

**Theorem 1:** If a space-time with one asymptotic region is enclosing a nonvanishing magnetic mass, a global system of

orthonormal tetrads cannot be defined, which rules out spinor structure.

**Theorem 2:** The nonvanishing of the first Chern class or, equivalently, the existence of  $S_1$  (Hopf fibering) nontrivial bundle structure over  $S^2$  at conformal null infinity, prevents the existence of global spinor fields on the corresponding (space-time) asymptotic region.

**Theorem 3:** The vanishing of the first Chern class provides a necessary (and in many cases sufficient) condition for the existence of spinor structure.

For an example where the space-time topology prevents the condition from being sufficient, see, e.g., the Appendix in Ref. 9.

## V. CONCLUDING REMARKS

(a) It was shown in Ref. 7 that a necessary and sufficient condition for a noncompact space-time  $M$  to have spinor structure is that  $M$  may be given a global system of orthonormal tetrads. It was also underlined that this condition might not be the most convenient one to decide on this issue. In Sec. IV we related the criterion to the global topological structure of the space-time manifold, reflected by the topology of its conformal null boundary, and reduced the test for existence of spinor structure to that of detection of a vanishing or nonvanishing topological charge (the magnetic mass) conveniently labeled by the value of suitable homology classes. Thus, although curvature plays a role [e.g., in the labeling of the homotopy classes of the (asymptotic) Lorentz Lie group], it is rather its intertwining with the underlying topology that is crucial. The existence of spinor structure reflects the global structure of the space-time manifold, and originates in the topology.

(b) The existence of spinor structure has also been characterized<sup>7</sup> in terms of the type of the Weyl tensor  $C_{abcd}$ , a sufficient condition being that this tensor be *everywhere* type  $[1, 1, 1, 1]$ ,  $[2, 1, 1]$ ,  $[3, 1]$ , or  $[4]$ . Here again the criterion involves the existence of continuous transport of orthonormal tetrads related to the principal null directions of  $C_{abcd}$ .

The above results suggest that a relation might exist between the algebraic type of solution to the source-free Einstein equation and its global (topological) structure as described by topological charges, or by the parametrization of suitable homology (cohomology) classes. Since these characteristic classes are conveniently used in the classification of fiber bundles, one might hope to be able to relate them to the algebraic type of the ( $\mathcal{S}$ -bundle) curvature two-form or to the algebraic type<sup>24</sup> of the Weyl curvature.

(c) We have proved elsewhere that gravitational magnetic monopoles exhibit causality violations—continuous time reversal (e.g., if an everywhere timelike and complete Killing vector field is available, its orbits must be closed, or in absence of isometry, nontrivial bundle structure with  $S^1$  Hopf fibering over  $S^2$  at conformal null infinity). Such solutions thus could illustrate a situation where parity reversal becomes a continuous transformation, PT invariance being a reasonable statement, e.g., when the principle frame bundle over the space-time manifold is simply connected. The “unwinding” of gravitational magnetic monopoles (i.e., NUT solution) into (black-hole) electric mass monopoles (i.e.,

Schwarzschild black hole) could be viewed as an onset of two asymptotic regions, of a discrete PT symmetry and of unambiguous notions of parity and time reversal. As was already mentioned, an invariant (gravitational charge C) is also involved in the transformation.

Since nontrivial topologies and multiply connected space-times could provide observational evidence of their existence<sup>25</sup> at some stage of the evolution of the universe (periodicity in the distribution of quasar red shifts?) we hope that the above considerations might consequently shed light on a possible (experimental) confirmation of a link between the CPT theorem and the existence of spinor structure.<sup>26,27</sup>

## ACKNOWLEDGMENTS

I would like to thank Professor D. Boulware, Professor E. T. Newman, and Professor R. Wald for useful conversations. Special thanks are due to Professor J. Hartle for his very kind help. I also want to thank J. Reynolds for his prompt processing of this manuscript at the Institute for Theoretical Physics, Santa Barbara.

This research was supported in part by the National Science Foundation under Grant No. PHY82-17853, supplemented by funds from the National Aeronautics and Space Administration, at the University of California at Santa Barbara.

<sup>1</sup>Ya.B. Zel'dovich and I. D. Novikov, *JETP Lett.* **1967**, 2136; I. D. Novikov, A. G. Polnarev, A. A. Starobinsky, and Ya.B. Zel'dovich, *Astron. Astrophys.* **80**, 104 (1979); Ya.B. Zel'dovich and I. D. Novikov, *Structure and Evolution of the Universe* (U. Chicago P., Chicago, 1983), revised edition.

<sup>2</sup>M. Suveges, *Acta Phys. Hung.* **20**, 273 (1966).

<sup>3</sup>A. S. Eddington, *Space-Time and Gravitation: An Outline of the Theory of Gravitation* (Cambridge U.P., Cambridge, 1959).

<sup>4</sup>R. Penrose, *Structure of Space-Time, Battelle Rencontres 1967* (Benjamin, New York, 1968).

<sup>5</sup>R. Penrose, "Null hypersurface initial data for classical fields of arbitrary spin and for General Relativity," in report ARL 63-56 USAF, 1963.

<sup>6</sup>R. Geroch, *J. Math. Phys.* **9**, 1739 (1968).

<sup>7</sup>R. Geroch, *J. Math. Phys.* **11**, 343 (1970).

<sup>8</sup>Y. Aharonov and L. Susskind, *Phys. Rev.* **158**, 1237 (1967).

<sup>9</sup>C.T.S. Clarke, *Gen. Relativ. Gravit.* **2**, 43 (1971).

<sup>10</sup>E. P. Wigner, *Ann. Math.* **40**, 149 (1939).

<sup>11</sup>A. Magnon, *J. Math. Phys.* **27**, 1059 (1986).

<sup>12</sup>A. Magnon, "Mass, dual mass and gravitational entropy," *J. Math. Phys.*, submitted for publication.

<sup>13</sup>A. Magnon, *J. Math. Phys.* **27**, 1066 (1986).

<sup>14</sup>R. Penrose, *Proc. R. Soc. London Ser. A* **284**, 159 (1965).

<sup>15</sup>R. Geroch, *J. Math. Phys.* **12**, 6 (1971).

<sup>16</sup>M. Buchdahl, *Quart. J. Math.* **5** (1954).

<sup>17</sup>J. Ehlers, in "Les Théories relativistes de la gravitation," CNRS, Paris, 1959.

<sup>18</sup>B. K. Harrison, *J. Math. Phys.* **9**, 1744 (1968).

<sup>19</sup>C. Nash and S. Sen, *Topology and Geometry for Physicists* (Academic, New York, 1983).

<sup>20</sup>N. Steenrod, *The topology of Fiber Bundles* (Princeton U.P., Princeton, NJ, 1951).

<sup>21</sup>J. Milnor, "Lectures on characteristic classes," Princeton University, 1957 (mimeographed notes).

<sup>22</sup>Sz. T. Hu, *Homotopy Theory* (Academic, New York, 1959).

<sup>23</sup>F. Hirzebruch, *Neue Topologische Methoden in der Algebraischen Geometrie* (Springer, Berlin, 1956).

<sup>24</sup>R. Penrose, *Ann. Phys. (NY)* **10**, 171 (1960).

<sup>25</sup>L. Z. Fang and H. Sato, "Is the periodicity in the distribution of quasar redshifts an evidence of a multiply connected universe?," 1985 first award-winning essay, Gravity Research Foundation, Gloucester, MA 01930.

<sup>26</sup>R. Penrose and W. Rindler, *Spinors and Space-Time* (Cambridge U.P., Cambridge, 1984), Vol. 1.

<sup>27</sup>R. Penrose, W. Rindler, *Spinors and Space-Time* (Cambridge U.P., Cambridge, 1986), Vol. 2.

# Einstein–Maxwell equations and the conformal Ricci collineations

Abbas M. Faridi

*Department of Physics, University of California, Santa Barbara, California 93206*

(Received 1 October 1986; accepted for publication 11 February 1987)

The Einstein–Maxwell field equations for non-null electromagnetic fields are studied under the assumption of admitting a conformal Ricci collineation. It is shown that a non-null electromagnetic field does not admit any conformal Ricci collineation, unless the generators of the symmetry groups are Killing vector fields. Furthermore, it is shown that the energy-momentum tensor of a non-null electromagnetic field can admit a conformal Ricci collineation, if and only if the collineation is homothetic. The restrictions on non-null Maxwell field, its sources, and its invariants implied by the symmetry condition are calculated. An example of a space-time satisfying the Einstein–Maxwell equations, and admitting a homothetic conformal vector field is also given.

## I. INTRODUCTION

The connection between the inherent symmetries of a system, and their corresponding conservation laws has been well established since the early observation of Nöther.<sup>1</sup> In particular, the importance of groups of motions, generated by Killing vector fields of a space-time, and their relation to the conservation laws of energy, momentum, and angular momentum is well known.<sup>2,3</sup> Collineations other than groups of motions have also been studied to a limited extent. It has been shown, for example, that for space-times with zero Ricci tensor, the more familiar symmetries such as motions and conformal and homothetic motions are subcases of a more general symmetry requirement known as curvature collineations.<sup>4</sup> Moreover, it has been shown that the allowable conformal symmetries admitted by a nonflat empty space-time (except in the case of specific type- $N$  metrics) are the homothetic motions.<sup>5</sup> The significance of homothetic motions in general relativity is yet to be elaborated; though they have been used by a number of authors. For example, homothetic motions have been utilized in the analysis of spherically symmetric, and plane symmetric self-similar space-times<sup>6,7</sup>; homothetic Weyl space-times,<sup>8</sup> as well as in certain self-similar cosmologies.<sup>9</sup>

In this paper we are concerned with Einstein–Maxwell field equations

$$R_{\mu\nu} = f_{\mu\sigma} f^{\sigma}_{\nu} - \frac{1}{4} g_{\mu\nu} f_{\alpha\beta} f^{\alpha\beta}, \quad (1.1)$$

$$f^{\mu\nu};_{\nu} = j^{\mu}, \quad (1.2)$$

$$*f^{\mu\nu};_{\nu} = 0, \quad (1.3)$$

admitting a conformal motion

$$L_{\xi} g_{\mu\nu} = 2\phi g_{\mu\nu}, \quad (1.4)$$

where  $L$  denotes the Lie derivative,  $g_{\mu\nu}$  is the metric of the space-time, and  $\phi = \frac{1}{4} \xi^{\mu};_{\mu}$  is a scalar function. The electromagnetic field  $f_{\mu\nu}$ , and its dual  $*f_{\mu\nu}$  satisfy the identities

$$f_{\mu\sigma} f^{\sigma\nu} - *f_{\mu\sigma} *f^{\sigma\nu} = \frac{1}{2} (f_{\alpha\beta} f^{\beta\alpha}) \delta_{\mu}^{\nu}, \quad (1.5)$$

$$f_{\mu\sigma} *f^{\sigma\nu} = *f_{\mu\sigma} f^{\sigma\nu} = \frac{1}{4} (f_{\alpha\beta} *f^{\beta\alpha}) \delta_{\mu}^{\nu}. \quad (1.6)$$

Throughout this paper, the Greek indices are tensor indices and range over the values 1, 2, 3, 4; and the lightface Latin  $a, b, c, \dots$  denote tetrad indices.

It has been shown (see, e.g., Ref. 4) that every curvature collineation

$$L_{\xi} R^{\mu\nu\rho\sigma} = 0, \quad (1.7)$$

generated by a conformal vector field  $\xi^{\mu}$  [Eq. (1.4)] is also a Ricci collineation,

$$L_{\xi} R_{\mu\nu} = 0, \quad (1.8)$$

provided the scalar function  $\phi$  satisfies

$$\phi;_{\mu\nu} = 0. \quad (1.9)$$

The problem that we wish to investigate may be posed in the following context. Suppose that the field equations (1.1)–(1.3) are satisfied for a non-null electromagnetic field, and the space-time admits a conformal Ricci collineation. What restrictions are then imposed on the Maxwell field  $f_{\mu\nu}$ , and its source  $j^{\mu}$ ? In particular, we want to determine whether or not these limitations depend on the form of the scalar function  $\phi$ , as defined by Eq. (1.4).

## II. SOME GEOMETRIC RELATIONS DUE TO RICCI CONFORMAL SYMMETRY CONDITIONS

In this section we make use of the proposed symmetry condition, and derive relations to be fulfilled by the field equations. These results will be used in the subsequent section to simplify the computations.

It is well known that a non-null electromagnetic field  $f_{\mu\nu}$ , and its dual  $*f_{\mu\nu}$ , can be expressed in the form

$$f_{\mu\nu} = 2\phi_0 (n_{\mu} l_{\nu} - n_{\nu} l_{\mu}) + 2i\psi_0 (m_{\mu} \bar{m}_{\nu} - \bar{m}_{\mu} m_{\nu}), \quad (2.1)$$

$$*f_{\mu\nu} = 2\psi_0 (n_{\mu} l_{\nu} - n_{\nu} l_{\mu}) - 2i\phi_0 (m_{\mu} \bar{m}_{\nu} - \bar{m}_{\mu} m_{\nu}), \quad (2.2)$$

where  $\phi_0$  and  $\psi_0$  are the real and imaginary parts of the scalar

$$\Phi = \frac{1}{2} f_{\mu\nu} (l^{\mu} n^{\nu} + \bar{m}^{\mu} m^{\nu}) = \phi_0 + i\psi_0. \quad (2.3)$$

The null vectors  $l^{\mu}$ ,  $n^{\mu}$  are real,  $m^{\mu}$  complex, and satisfy

$$l_{\mu} n^{\mu} = -m_{\mu} \bar{m}^{\mu} = 1, \quad (2.4)$$

with all other contractions being zero.

The two invariants of the electromagnetic fields are

$$f_{\mu\nu} f^{\nu\mu} = 4(\Phi^2 + \bar{\Phi}^2) = 8(\phi_0^2 - \psi_0^2), \quad (2.5)$$



$$f_{\mu\nu} *f^{\nu\mu} = 4i(\Phi^2 - \bar{\Phi}^2) = 16\phi_0\psi_0. \quad (2.6)$$

Einstein–Maxwell equations may now be written in terms of the above invariants and the tetrad of null vectors  $l^\mu$ ,  $n^\mu$ , and  $m^\mu$  and  $\bar{m}^\mu$ . We have<sup>10–12</sup>

$$\begin{aligned} R_{\mu\nu} &= \Lambda(l_\mu n_\nu + n_\mu l_\nu + m_\mu \bar{m}_\nu + \bar{m}_\mu m_\nu) \\ &= 2\Lambda(l_\mu n_\nu + n_\mu l_\nu) - \Lambda g_{\mu\nu}, \end{aligned} \quad (2.7)$$

where  $\Lambda$  is defined by the relation

$$\Lambda = 4\Phi\bar{\Phi} = 4(\phi_0^2 + \psi_0^2), \quad (2.8)$$

and the metric tensor has the form

$$g_{\mu\nu} = l_\mu n_\nu + n_\mu l_\nu - \bar{m}_\mu m_\nu - m_\mu \bar{m}_\nu. \quad (2.9)$$

We may express  $\Lambda$  in terms of the invariants by forming the expression

$$R_{\mu\sigma} R^{\sigma\nu} = \Lambda^2 \delta_\mu^\nu = \frac{1}{4}[(f_{\alpha\beta} f^{\beta\alpha})^2 + (f_{\alpha\beta} *f^{\beta\alpha})^2] \delta_\mu^\nu, \quad (2.10)$$

where use has been made of Eqs. (1.5) and (1.6). We therefore have

$$\Lambda^2 = \frac{1}{4}(R_{\mu\sigma} R^{\sigma\mu}) = \frac{1}{4}[(f_{\mu\nu} f^{\nu\mu})^2 + (f_{\mu\nu} *f^{\nu\mu})^2]. \quad (2.11)$$

From Eqs. (1.8) and (2.11) we obtain the restriction imposed on the invariant of the electromagnetic field, and may be written in the form

$$L_\xi \Lambda = -2\phi\Lambda. \quad (2.12)$$

Lie differentiation of Eq. (2.7) yields

$$\begin{aligned} l_\mu L_\xi(n_\nu) + n_\nu L_\xi(l_\mu) + n_\mu L_\xi(l_\nu) + l_\nu L_\xi(n_\mu) \\ - 2\phi(l_\mu n_\nu - n_\mu l_\nu) = 0. \end{aligned} \quad (2.13)$$

The contraction of Eq. (2.13) with  $l^\mu$ , and  $n^\mu$ , gives, respectively,

$$L_\xi(l_\mu) = 2\phi l_\mu - l_\mu l^\alpha L_\xi(n_\alpha), \quad (2.14)$$

$$L_\xi(n_\mu) = 2\phi n_\mu - n_\mu n^\alpha L_\xi(l_\alpha). \quad (2.15)$$

Equations (2.14) and (2.15) may be simplified if use is made of the normalization condition Eq. (2.4), which after differentiation takes the form

$$n^\mu L_\xi(l_\mu) + l^\mu L_\xi(n_\mu) = 2\phi. \quad (2.16)$$

Eliminating the quantity  $n^\alpha L_\xi(l_\alpha)$  by means of Eq. (2.16), we may write

$$L_\xi(l_\mu) = (\phi + \lambda)l_\mu, \quad (2.17)$$

$$L_\xi(n_\mu) = (\phi - \lambda)n_\mu, \quad \lambda = \phi - l^\alpha L_\xi(n_\alpha). \quad (2.18)$$

Similarly, making use of Eqs. (1.4), (2.9), (2.17), and (2.18), we obtain

$$m^\mu L_\xi(\bar{m}_\mu) + \bar{m}^\mu L_\xi(m_\mu) = -2\phi, \quad (2.19)$$

$$L_\xi(m_\alpha) = (\phi - \delta)m_\alpha, \quad (2.20)$$

$$L_\xi(\bar{m}_\alpha) = (\phi + \delta)\bar{m}_\alpha, \quad \delta = \phi + \bar{m}^\sigma L_\xi(m_\sigma). \quad (2.21)$$

We are now in a position to write the restrictions implied on Maxwell field, due to conformal Ricci collineation. We obtain, after some algebra

$$\begin{aligned} L_\xi f_{\mu\nu} &= 2L_\xi \phi_0(n_\mu l_\nu - n_\nu l_\mu) \\ &\quad + 2iL_\xi \psi_0(m_\mu \bar{m}_\nu - \bar{m}_\mu m_\nu) + 2\phi f_{\mu\nu}, \end{aligned} \quad (2.22)$$

$$\begin{aligned} L_\xi *f_{\mu\nu} &= 2L_\xi \psi_0(n_\mu l_\nu - n_\nu l_\mu) \\ &\quad - 2iL_\xi \phi_0(m_\mu \bar{m}_\nu - m_\nu \bar{m}_\mu) + 2\phi *f_{\mu\nu}. \end{aligned} \quad (2.23)$$

Equations (2.22) and (2.23) may be transformed into a more useful form. We may write these equations in the form

$$L_\xi f_{\mu\nu} = A f_{\mu\nu} + B *f_{\mu\nu}, \quad (2.24)$$

$$L_\xi *f_{\mu\nu} = A *f_{\mu\nu} - B f_{\mu\nu}, \quad (2.25)$$

where

$$A = (\phi_0 L_\xi \phi_0 + \psi_0 L_\xi \psi_0)/(\phi_0^2 + \psi_0^2) + 2\phi, \quad (2.26)$$

$$B = (\psi_0 L_\xi \phi_0 - \phi_0 L_\xi \psi_0)/(\phi_0^2 + \psi_0^2). \quad (2.27)$$

Taking the covariant divergence of Eqs. (2.24) and (2.25), and making use of Maxwell equations, we obtain for the right-hand side of these equations, the relations

$$\gamma^{\mu\nu};\nu = A_\nu f^{\mu\nu} + B_\nu *f^{\mu\nu} + A j^\mu, \quad (2.28)$$

$$\gamma^{*\mu\nu};\nu = A_\nu *f^{\mu\nu} - B_\nu f^{\mu\nu} - B j^\mu, \quad (2.29)$$

where

$$\gamma_{\mu\nu} = L_\xi f_{\mu\nu}. \quad (2.30)$$

Finally, applying the identities Eqs. (A3) and (A5) and the compatibility condition Eq. (A4) for the vector field  $\xi^\mu$ , we obtain for the left-hand side of Eqs. (2.28) and (2.29) the expressions

$$\gamma^{\mu\nu};\nu = L_\xi(j^\mu) + 4\phi j^\mu, \quad (2.31)$$

$$\gamma^{*\mu\nu};\nu = 0. \quad (2.32)$$

The implication of the equations derived in this section on the functions  $A$ ,  $B$ , and  $\phi$  as defined by Eqs. (2.26) and (2.27) are sought in the next section.

### III. DETERMINATION OF THE FUNCTIONS $A$ , $B$ , AND $\phi$

It will be shown here that in order for the system of equations (1.1)–(1.3) to be consistent for a nontrivial, non-null electromagnetic field under Ricci collineation, Eq. (1.8), the allowable conformal vector field must be homothetic,

$$\phi = \text{const.} \quad (3.1)$$

There are two cases to be treated separately; namely we can either have a conformally invariant non-null electromagnetic field

$$\text{case (a): } L_\xi f_{\mu\nu} = 0 \quad (3.2)$$

or

$$\text{case (b): } L_\xi f_{\mu\nu} = \gamma_{\mu\nu} \neq 0. \quad (3.3)$$

Case (a): In this case we obtain from Eqs. (1.1), (2.24), and (2.25), the simple relation

$$(A - 2\phi)R_{\mu\nu} = 0, \quad (3.4)$$

which implies either (i)  $A - 2\phi \neq 0$ ,  $R_{\mu\nu} = 0$ , or (ii)  $A = 2\phi$ ,  $R_{\mu\nu} \neq 0$ . The results for the case of vanishing Ricci tensor have been obtained (see, e.g., Ref. 5). On the other hand, if  $R_{\mu\nu} \neq 0$ , Eqs. (2.8) and (2.26) yield

$$L_\xi \Lambda = 0, \quad \Lambda = 4(\phi_0^2 + \psi_0^2). \quad (3.5)$$

According to Eq. (2.12), a non-null electromagnetic field satisfies Eq. (3.5) if and only if

$$\phi = 0. \quad (3.6)$$

The condition  $\phi = 0$ , for a nonvanishing energy-momentum tensor of a non-null electromagnetic field, implies by virtue of Eq. (1.4), that the vector field  $\zeta^\mu$  must be the generator of an isometry group. The details of this case have also been worked out by Wooley,<sup>13,14</sup> in his analysis of Einstein–Maxwell equations, under the assumptions  $L_\zeta R_{\mu\nu} = 0$ ,  $\zeta_{(\mu;\nu)} = 0$ . The results for case (a) may be summarized in the following theorem.

**Theorem 1:** In the combined Einstein–Maxwell theory, a non-null electromagnetic field, and the Ricci tensor, cannot simultaneously be conformally invariant.

An immediate consequence of Theorem 1 is that the conformal invariance of Ricci tensor may be shared by the energy-momentum tensor of a non-null electromagnetic field rather than the field itself. This in fact turns out to be the case, and embodies the main content of the case (b), specified by Eq. (3.3).

Case (b): In this case Eqs. (2.28)–(2.32), and the vanishing of the covariant divergence of the energy-momentum tensor yield

$$B_\nu *f^{\mu\nu} = L_\zeta(j^\mu) + (4\phi - A)j^\mu - A_\nu f^{\mu\nu}, \quad (3.7)$$

$$B_\nu f^{\mu\nu} = A_\nu *f^{\mu\nu} - B j^\mu, \quad (3.8)$$

$$f^\mu \nu j^\nu = 0. \quad (3.9)$$

The analysis of these equations depends on whether the electromagnetic field is source-free,  $j^\mu = 0$ , which we label it as case (1b), or otherwise,  $j^\mu \neq 0$ , which we denote it as case (2b). Consequently, the behavior of the functions  $A$ ,  $B$ , and  $\phi$ , if they exist at all, may depend on quite different conditions. Similarly, by virtue of Eqs. (2.24)–(2.27), the behavior of a nontrivial Maxwell field may likewise be different.

Case (1b): Let us consider the source-free case. In this case, Eqs. (3.7)–(3.9) become

$$B_\nu f^{\mu\nu} = A_\nu *f^{\mu\nu}, \quad j^\mu = 0, \quad (3.10)$$

$$B_\nu *f^{\mu\nu} = -A_\nu f^{\mu\nu}. \quad (3.11)$$

Inserting in these equations, the expressions for the field and its dual from Eqs. (2.1) and (2.2) and contracting the resulting two equations with the null vectors  $l^\mu$ ,  $n^\mu$ ,  $m^\mu$ , and  $\bar{m}^\mu$ , we find that for a non-null electromagnetic field, we must have

$$A_\mu = 0, \quad B_\mu = 0. \quad (3.12)$$

The constant value of  $A$  may be obtained from Eqs. (2.12) and (2.26). In fact we have

$$A = \phi = \text{const.} \quad (3.13)$$

The constant  $B$  is then restricted by either one of the following relations, obtained from Eqs. (2.26) and (2.27):

$$L_\zeta \phi_0 = -\phi \phi_0 + B \psi_0, \quad j^\mu = 0, \quad (3.14)$$

$$L_\zeta \psi_0 = -\phi \psi_0 - B \phi_0, \quad j^\mu = 0. \quad (3.15)$$

The corresponding relations for the Maxwell field, given by Eqs. (2.24) and (2.25), assume the form

$$L_\zeta f_{\mu\nu} = \phi f_{\mu\nu} + B *f_{\mu\nu}, \quad j^\mu = 0, \quad (3.16)$$

$$L_\zeta *f_{\mu\nu} = \phi *f_{\mu\nu} - B f_{\mu\nu}. \quad (3.17)$$

Case (2b): In this case, where  $j^\mu \neq 0$ , Eq. (3.9) implies that for a nontrivial  $j^\mu$  to exist, the determinant of the admis-

sible non-null electromagnetic field must vanish. The vanishing of the determinant is equivalent to the vanishing of the invariant<sup>15</sup>

$$f_{\mu\nu} *f^{\nu\mu} = 0. \quad (3.18)$$

The physical implication of Eq. (3.18) is the orthogonality of the electric and the magnetic fields in a local Minkowskian frame. From Eqs. (2.6) and (3.18) we have either

$$\psi_0 = 0, \quad \phi_0 \neq 0, \quad j^\mu \neq 0, \quad (3.19)$$

or

$$\phi_0 = 0, \quad \psi_0 \neq 0. \quad (3.20)$$

In either case, the function  $B$ , as given by Eq. (2.27) is equal to zero. Equations (3.7) and (3.8) and the Lie differentiation of (3.9) give us

$$L_\zeta j^\mu + (4\phi - A)j^\mu - A_\nu f^{\mu\nu} = 0, \quad (3.21)$$

$$A_\nu f^{*\mu\nu} = 0, \quad (3.22)$$

$$f_{\nu\mu} L_\zeta j^\mu = 0. \quad (3.23)$$

Multiplying Eq. (3.21) by  $f_{\sigma\mu}$ , and making use of Eqs. (1.5), (1.6), and (3.22), we obtain

$$(f \cdot f) A_\mu = 0, \quad (f \cdot f) = f_{\mu\nu} f^{\nu\mu}. \quad (3.24)$$

Since  $(f \cdot f)$  is the only nonvanishing invariant of the electromagnetic field in this case, we must have  $A_\mu = 0$ , or  $A$  is at most a constant. With  $A$  being a constant, Eqs. (2.12) and (2.26) give the results

$$A = \phi = \text{const.}, \quad (3.25)$$

$$L_\zeta(\phi_0) = -\phi \phi_0, \quad \phi_0 \neq 0, \quad \psi_0 = 0, \quad j^\mu \neq 0, \quad (3.26)$$

$$L_\zeta(\psi_0) = -\phi \psi_0, \quad \psi_0 \neq 0, \quad \phi_0 = 0. \quad (3.27)$$

The corresponding restrictions on the electromagnetic field and its source are obtained from Eqs. (2.24), (2.25), and (3.21). They will take the form

$$L_\zeta f_{\mu\nu} = \phi f_{\mu\nu}, \quad (3.28)$$

$$L_\zeta *f_{\mu\nu} = \phi *f_{\mu\nu}, \quad j^\mu \neq 0, \quad (3.29)$$

$$L_\zeta j^\mu = -3\phi j^\mu. \quad (3.30)$$

The results obtained for case (b) may be stated in the following theorem.

**Theorem 2:** In the combined Einstein–Maxwell theory, the energy momentum tensor of a non-null electromagnetic field, with or without a source, and the Ricci tensor, are simultaneously conformally invariant, if and only if the conformal vector field is homothetic.

#### IV. A PARTICULAR SOLUTION

In this section we construct a special source-free ( $j^\mu = 0$ ) solution, satisfying the conditions (a) the principal null directions  $l^\mu$ , and  $n^\mu$  are geodesics, (b) the null tetrad of vector  $l^\mu$ ,  $n^\mu$ ,  $m^\mu$ , and  $\bar{m}^\mu$  are parallelly propagated along  $l^\mu$  and  $n^\mu$ , and (c) the null geodesics with tangent vectors  $l^\mu$  and  $n^\mu$  are twist-free.

The particular solution to Einstein–Maxwell equations satisfying the above requirements has been found, using the Cartan's equations of structure. We will not, however, present the details of the calculations here, since it has also been

worked out by Tariq and Tupper<sup>16,17</sup> in Newman–Penrose formalism (Debney and Zund<sup>18,19</sup> also give some useful results shared by the principal null geodesics, and when the tetrad vectors are parallelly propagated along them). Our problem therefore reduces to determine whether or not Eq. (1.4) admits a nontrivial homothetic solution for the given metric obtained from the field equations.

The solution to Einstein–Maxwell equations under the assumed conditions (a), (b), and (c), may be specified by the relations

$$l^\mu = \delta_r^\mu, \quad n^\mu = \delta_u^\mu, \quad m^\mu = \lambda_1 \delta_y^\mu + \lambda_2 \delta_z^\mu, \quad (4.1)$$

where the coordinates are  $x^\mu$ : ( $x^1 = u$ ,  $x^2 = r$ ,  $x^3 = y$ ,  $x^4 = z$ ), and the functions  $\lambda_1$ , and  $\lambda_2$  are, respectively,

$$\lambda_1 = (1/\sqrt{2})u^n r^m, \quad \lambda_2 = (i/\sqrt{2})u^m r^n, \quad (4.2)$$

with  $m = (\sqrt{3} - 1)/4$ , and  $n = -(\sqrt{3} + 1)/4$ . The metric is defined by Eq. (2.9) and the Maxwell field by Eqs. (2.1)–(2.3), with

$$\Phi \bar{\Phi} = \phi_0^2 + \psi_0^2 = 1/8ur. \quad (4.3)$$

The only nonvanishing spin coefficients are

$$\rho = \bar{\rho} = -1/4r, \quad \mu = \bar{\mu} = 1/4u, \quad (4.4)$$

$$\sigma = \bar{\sigma} = \sqrt{3}/4r, \quad \lambda = \bar{\lambda} = \sqrt{3}/4u, \quad (4.5)$$

having the intrinsic derivatives

$$D\rho = 4\rho^2, \quad D\sigma = 4\rho\sigma, \quad \Delta\mu = -4\mu^2, \quad \Delta\lambda = -4\lambda\mu, \quad (4.6)$$

with all other intrinsic derivatives being zero.

With these spin coefficients, the commutation relations Eqs. (B10a)–(B10d), the tetrad components of the homothetic vector field  $\zeta^a$  ( $a = 1, 2, 3, 4$ ), Eqs. (B5a)–(B5j), and its corresponding compatibility conditions, Eqs. (B13)–(B17), assume the form, respectively;

$$(\Delta D - D\Delta)\psi = 0, \quad (4.7a)$$

$$(\delta\bar{\delta} - \bar{\delta}\delta)\psi = 0, \quad (4.7b)$$

$$(\delta D - D\delta)\psi = -\rho\delta\psi - \sigma\bar{\delta}\psi, \quad (4.7c)$$

$$(\delta\Delta - \Delta\delta)\psi = \mu\delta\psi + \lambda\bar{\delta}\psi, \quad (4.7d)$$

$$D\zeta_1 = 0, \quad (4.8a)$$

$$\Delta\zeta_2 = 0, \quad (4.8b)$$

$$\delta\zeta_3 = \lambda\zeta_1 - \sigma\zeta_2, \quad (4.8c)$$

$$\Delta\zeta_1 + D\zeta_2 = 2\phi, \quad (4.8d)$$

$$\delta\zeta_1 + D\zeta_3 = -\rho\zeta_3 - \sigma\zeta_4, \quad (4.8e)$$

$$\delta\zeta_2 + \Delta\zeta_3 = \mu\zeta_3 + \lambda\zeta_4, \quad (4.8f)$$

$$\bar{\delta}\zeta_3 + \delta\zeta_4 = -2\phi + 2\mu\zeta_1 - 2\rho\zeta_2, \quad (4.8g)$$

$$\bar{\delta}\zeta_4 = \lambda\zeta_1 - \sigma\zeta_2, \quad (4.8h)$$

$$\bar{\delta}\zeta_1 + D\zeta_4 = -\rho\zeta_4 - \sigma\zeta_3, \quad (4.8i)$$

$$\bar{\delta}\zeta_2 + \Delta\zeta_4 = \mu\zeta_4 + \lambda\zeta_3, \quad (4.8j)$$

$$\zeta_1 \Delta\psi_0 + \zeta_2 D\psi_0 - \zeta_3 \bar{\delta}\psi_0 - \zeta_4 \delta\psi_0 + 2\psi_0 (D\zeta_2 - \delta\zeta_4) = 2\psi_0 (\phi + \rho\zeta_2 - \mu\zeta_1), \quad (4.9a)$$

$$\psi_0 \Delta\zeta_4 + \psi_2 (2D\zeta_3 - \delta\zeta_1) = 2\psi_0 (\lambda\zeta_3 + \mu\zeta_4) - 2\psi_2 (\rho\zeta_3 + \sigma\zeta_4), \quad (4.9b)$$

$$\zeta_1 \Delta\psi_2 + \zeta_2 D\psi_2 - \zeta_3 \bar{\delta}\psi_2 - \zeta_4 \delta\psi_2 = -2\phi\psi_2, \quad (4.9c)$$

$$\psi_4 D\zeta_3 - \psi_2 (\bar{\delta}\zeta_2 - 2\Delta\zeta_4) = 2\psi_2 (\lambda\zeta_3 + \mu\zeta_4) - \psi_4 (\rho\zeta_3 + \sigma\zeta_4), \quad (4.9d)$$

$$\zeta_1 \Delta\psi_4 + \zeta_2 D\psi_4 - \zeta_3 \bar{\delta}\psi_4 - \zeta_4 \delta\psi_4 + 2\psi_4 (\Delta\zeta_1 - \bar{\delta}\zeta_3) = 2\psi_4 (\phi - \mu\zeta_1 + \rho\zeta_2), \quad (4.9e)$$

where

$$\psi_0 = 2\rho\sigma, \quad \psi_2 = 2\rho\mu, \quad \psi_4 = 2\lambda\mu, \quad (4.10)$$

are the nonvanishing scalars associated with the Weyl tensor.

The compatibility conditions, Eqs. (4.9a)–(4.9e), and the commutation relations Eqs. (4.7a)–(4.7d), when use is made of Eq. (4.10), give us the results

$$\delta\zeta_1 = \bar{\delta}\zeta_1 = \delta\zeta_2 = \bar{\delta}\zeta_2 = 0, \quad (4.11)$$

$$\Delta\zeta_1 = 4\mu\zeta_1, \quad (4.12)$$

$$D\zeta_2 = -4\rho\zeta_2, \quad (4.13)$$

$$\Delta\zeta_3 = \mu\zeta_3 + \lambda\zeta_4, \quad (4.14)$$

$$\delta\zeta_3 = \lambda\zeta_1 - \sigma\zeta_2, \quad (4.15)$$

$$D\zeta_4 = -\rho\zeta_4 - \sigma\zeta_3, \quad (4.16)$$

$$\Delta\zeta_4 = \mu\zeta_4 + \lambda\zeta_3. \quad (4.17)$$

The integration of Eqs. (4.8a)–(4.8j) gives for the tetrad components of  $\zeta^a$ ,

$$\zeta_1 = cu, \quad (4.18)$$

$$\zeta_2 = (2\phi - c)r, \quad (4.19)$$

$$\zeta_3 = d_1 y u^{-n} r^{-m} - id_2 z u^{-m} r^{-n} + F_1(u, r), \quad (4.20)$$

$$\zeta_4 = d_1 y u^{-n} r^{-m} + id_2 z u^{-m} r^{-n} + F_2(u, r), \quad (4.21)$$

where  $c$ ,  $d_1$ , and  $d_2$  are constants, and  $F_1, F_2$  are two arbitrary integration constants. The constants  $d_1$  and  $d_2$ , may be expressed in terms of the constant  $c$ , and the homothetic constant  $\phi$ ,

$$d_1 = (\alpha - \sqrt{2}\phi)/2, \quad d_2 = \alpha/2, \quad (4.22)$$

$$\alpha = \sqrt{6}c/2 + \sqrt{2}(1 - \sqrt{3})\phi/2.$$

The arbitrary constants of integrations  $F_1, F_2$  must satisfy

$$F_{1,r} - (1/4r)F_1 + (\sqrt{3}/4r)F_2 = 0, \quad (4.23)$$

$$F_{1,u} - (1/4u)F_1 - (\sqrt{3}/4u)F_2 = 0,$$

$$F_{2,r} - (1/4r)F_2 - (\sqrt{3}/4r)F_1 = 0,$$

$$F_{2,u} - (1/4u)F_2 - (\sqrt{3}/4u)F_1 = 0.$$

The only admissible solution to the system of equations in (4.23) is the trivial solution

$$F_1(u, r) = F_2(u, r) = 0. \quad (4.24)$$

The complete particular solution to the Eqs. (1.1)–(1.9) with respect to the system of coordinates  $(u, r, y, z)$  may be expressed in the following form.

For the metric we have

$$g_{\mu\nu} = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & -u^{-2n}r^{-2m} & 0 \\ 0 & 0 & 0 & -u^{-2m}r^{-2n} \end{pmatrix}, \quad (4.25)$$

$$g^{\mu\nu} = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & -u^{2n}r^{2m} & 0 \\ 0 & 0 & 0 & -u^{2m}r^{2n} \end{pmatrix},$$

where  $m = (\sqrt{3} - 1)/4$ , and  $n = -(\sqrt{3} + 1)/4$ . The Maxwell field becomes

$$f_{\mu\nu} = \begin{pmatrix} 0 & -[(\cos \epsilon)/\sqrt{2}](ur)^{-1/2} & 0 & 0 \\ [(\cos \epsilon)/\sqrt{2}](ur)^{-1/2} & 0 & 0 & 0 \\ 0 & 0 & 0 & (\sin \epsilon)/\sqrt{2} \\ 0 & 0 & -(\sin \epsilon)/\sqrt{2} & 0 \end{pmatrix}, \quad (4.26)$$

where  $\epsilon$  is an arbitrary real constant. The homothetic vector field  $\zeta^\mu$  takes the form

$$\begin{aligned} \zeta^u &= cu, & \zeta^r &= (2\phi - c)r, \\ \zeta^y &= -c_1 y, & \zeta^z &= c_2 z, \end{aligned} \quad (4.27)$$

where the constants  $c_1$  and  $c_2$  are

$$c_1 = \sqrt{3}(c - \phi)/2, \quad c_2 = \sqrt{3}c/2 + (1 - \sqrt{3})\phi/2. \quad (4.28)$$

From Eq. (4.26), we see that both electric and magnetic fields are in the radial direction with respect to the adopted coordinate system  $(u, r, y, z)$ . In addition to the homothetic motion, the solution admits a three-parameter group of motions,

$$\begin{aligned} \eta_1^\mu &= \delta_y^\mu, & \eta_2^\mu &= \delta_z^\mu, \\ \eta_3^\mu &= u\delta_u^\mu - r\delta_r^\mu - dy\delta_y^\mu + dz\delta_z^\mu, \end{aligned}$$

where the constant  $d = \sqrt{3}c/2$  may be obtained from Eq. (4.28), by letting  $\phi = 0$ .

## APPENDIX A: CONFORMAL MOTION

We summarize here some of the relevant relations concerning a conformal motion. A space-time is said to admit a conformal motion if there exists a vector field  $\zeta^\mu$ , such that

$$L_\zeta g_{\mu\nu} = \zeta_{\mu;\nu} + \zeta_{\nu;\mu} = 2\phi g_{\mu\nu}, \quad (A1)$$

where the symbol  $L_\zeta$  denotes Lie differentiation with respect to the vector field  $\zeta^\mu$ , and  $\phi$  is a scalar function satisfying

$$\phi = \frac{1}{4} \zeta^\mu{}_{;\mu}. \quad (A2)$$

Every conformal motion must satisfy

$$L_\zeta \Gamma_{\nu\sigma}^\mu = \phi;_\sigma \delta_\nu^\mu + \phi;_\nu \delta_\sigma^\mu - g_{\nu\sigma} g^{\mu\alpha} \phi;_\alpha. \quad (A3)$$

The integrability condition for (A1) can be shown to have the form

$$\zeta_{\mu;\nu;\sigma} = \zeta^\alpha R_{\alpha\nu\mu} + \phi;_\sigma g_{\mu\nu} + \phi;_\nu g_{\mu\sigma} - \phi;_\mu g_{\nu\sigma}. \quad (A4)$$

For an arbitrary tensor  $K^{\mu\nu}$ , we have the identity<sup>20</sup>

$$L_\zeta (K^{\mu\nu};_\alpha) - (L_\zeta K^{\mu\nu});_\alpha = (L_\zeta \Gamma_{\alpha\beta}^\mu) K^{\beta\nu} + (L_\zeta \Gamma_{\alpha\beta}^\nu) K^{\mu\beta}. \quad (A5)$$

The Lie differentiation of the Riemann tensor can be expressed in the form

$$L_\zeta R_{\mu\nu\sigma}^\alpha = (L_\zeta \Gamma_{\sigma\mu}^\alpha);_\nu - (L_\zeta \Gamma_{\nu\mu}^\alpha);_\sigma. \quad (A6)$$

The curvature collineation is defined by the relation

$$L_\zeta R_{\mu\nu\sigma}^\alpha = 0, \quad (A7)$$

which by virtue of Eqs. (A2) and (A3) gives the restriction on  $\phi$

$$\phi;_{\mu\nu} = 0. \quad (A8)$$

The contraction of Eq. (A7) yields the Ricci collineation

$$L_\zeta R_{\mu\nu} = 0. \quad (A9)$$

The integrability condition for (A1) is generally expressed in terms of the vanishing of the Lie derivative of the conformal Weyl tensor

$$L_\zeta C_{\mu\nu\sigma}^\alpha = 0. \quad (A10)$$

## APPENDIX B: CONFORMAL MOTION AND ITS COMPATIBILITY CONDITIONS IN TETRAD REPRESENTATION

In the following we write the tetrad components equations of Eqs. (A1) and (A10). We choose a tetrad of null vector  $l^\mu, n^\mu, m^\mu, \bar{m}^\mu$ , with  $l^\mu$  and  $n^\mu$  real and  $m^\mu$  complex. The only nonvanishing contractions are

$$l_\mu n^\mu = -m_\mu \bar{m}^\mu = 1. \quad (B1)$$

A frame defined by the inner product

$$e_a^\mu e_\mu^b = \delta_a^b, \quad e_a^\mu e_\nu^a = \delta_\nu^\mu, \quad (B2)$$

where  $e_a^\mu (l^\mu, n^\mu, m^\mu, \bar{m}^\mu)$ , induces a metric of the form

$$\eta_{ab} = \eta^{ab} = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1 \\ 0 & 0 & -1 & 0 \end{pmatrix}, \quad a, b = 1, 2, 3, 4. \quad (B3)$$

Equation (A1) in tetrad representation takes the form

$$\zeta_{a;b} + \zeta_{b;a} = 2\phi \eta_{ab} + \zeta^c (\gamma_{acb} + \gamma_{bca}), \quad (B4)$$

where  $\gamma_{abc}$  are the Ricci rotation coefficients. Equation (B4) is equivalent to the set of scalar equations

$$D\xi_1 = (\epsilon + \bar{\epsilon})\xi_1 - \bar{\kappa}\xi_3 - \kappa\xi_4, \quad (B5a)$$

$$\Delta\xi_2 = -(\gamma + \bar{\gamma})\xi_2 + \nu\xi_3 + \bar{\nu}\xi_4, \quad (B5b)$$

$$\delta\xi_3 = \bar{\lambda}\xi_1 - \sigma\xi_2 - (\bar{\alpha} - \beta)\xi_3, \quad (B5c)$$

$$\bar{\delta}\xi_4 = \lambda\xi_1 - \bar{\sigma}\xi_2 - (\alpha - \bar{\beta})\xi_4, \quad (B5d)$$

$$\Delta\xi_1 + D\xi_2 = 2\phi + (\gamma + \bar{\gamma})\xi_1 - (\epsilon + \bar{\epsilon})\xi_2 + (\pi - \bar{\pi})\xi_3 + (\bar{\pi} - \tau)\xi_4, \quad (B5e)$$

$$\delta\xi_1 + D\xi_3 = (\bar{\alpha} + \beta + \bar{\pi})\xi_1 - \kappa\xi_2 + (\epsilon - \bar{\epsilon} - \bar{\rho})\xi_3 - \sigma\xi_4, \quad (B5f)$$

$$\bar{\delta}\xi_1 + D\xi_4 = (\alpha + \bar{\beta} + \pi)\xi_1 - \bar{\kappa}\xi_2 + (\bar{\epsilon} - \epsilon - \rho)\xi_4 - \bar{\sigma}\xi_3, \quad (B5g)$$

$$\delta\xi_2 + \Delta\xi_3 = \bar{\nu}\xi_1 - (\bar{\alpha} + \beta + \tau)\xi_2 + (\mu + \gamma - \bar{\gamma})\xi_3 + \bar{\lambda}\xi_4, \quad (B5h)$$

$$\bar{\delta}\xi_2 + \Delta\xi_4 = \nu\xi_1 - (\alpha + \bar{\beta} + \bar{\tau})\xi_2 + (\bar{\mu} + \bar{\gamma} - \gamma)\xi_4 + \lambda\xi_3, \quad (B5i)$$

$$\bar{\delta}\xi_3 + \delta\xi_4 = -2\phi + (\mu + \bar{\mu})\xi_1 - (\rho + \bar{\rho})\xi_2 + (\alpha - \bar{\beta})\xi_3 + (\bar{\alpha} - \beta)\xi_4, \quad (B5j)$$

where the spin coefficients  $\alpha\beta\gamma\cdots$  are related to the Ricci rotations and may be expressed in the form

$$\gamma_{131} = \kappa, \quad \gamma_{132} = \tau, \quad \gamma_{133} = \sigma, \quad \gamma_{134} = \rho, \quad (B6)$$

$$\gamma_{241} = -\pi, \quad \gamma_{242} = -\nu, \quad (B7)$$

$$\gamma_{243} = -\mu, \quad \gamma_{244} = -\lambda,$$

$$\gamma_{121} = \epsilon + \bar{\epsilon}, \quad \gamma_{122} = \gamma + \bar{\gamma}, \quad (B8)$$

$$\gamma_{123} = \beta + \bar{\alpha}, \quad \gamma_{124} = \alpha + \bar{\beta}.$$

The intrinsic derivatives are defined according to the relations

$$\begin{aligned} &\xi_1\Delta\psi_0 + \xi_2D\psi_0 - \xi_3\bar{\delta}\psi_0 - \xi_4\delta\psi_0 + 2\psi_0(D\xi_2 - \delta\xi_4) + 2\psi_1(\delta\xi_1 - D\xi_3) \\ &= 2\phi\psi_0 + 2\psi_0[(2\gamma - \mu)\xi_1 + (\epsilon - \bar{\epsilon} + \bar{\rho})\xi_2 + (\pi - 2\alpha)\xi_3 + (\bar{\pi} - \beta - \bar{\alpha})\xi_4] \\ &+ 2\psi_1[(\beta + \bar{\alpha} - 2\tau - \bar{\pi})\xi_1 - \kappa\xi_2 + (\bar{\epsilon} - \epsilon + 2\rho - \bar{\rho})\xi_3 + \sigma\xi_4], \end{aligned} \quad (B13)$$

$$\begin{aligned} &\xi_1\Delta\psi_1 + \xi_2D\psi_1 - \xi_3\bar{\delta}\psi_1 - \xi_4\delta\psi_1 - \psi_0\Delta\xi_4 + \psi_1(D\xi_2 - \delta\xi_4) + \psi_2(\delta\xi_1 - 2D\xi_3) \\ &= \psi_0[(\pi + \bar{\tau})\xi_2 - \lambda\xi_3 + (\gamma - \bar{\gamma} - \mu)\xi_4] + \psi_1[(2\gamma - \mu)\xi_1 + (\epsilon - \bar{\epsilon} + \bar{\rho})\xi_2 + (\pi - 2\alpha)\xi_3 + (\bar{\pi} - \bar{\alpha} - \beta)\xi_4] \\ &+ \psi_2[(\bar{\alpha} + \beta - 2\bar{\pi} - 3\tau)\xi_1 - \kappa\xi_2 + (3\rho - \bar{\rho} - 2\epsilon + 2\bar{\epsilon})\xi_3 + 2\sigma\xi_4], \end{aligned} \quad (B14)$$

$$\begin{aligned} &\xi_1\Delta\psi_2 + \xi_2D\psi_2 - \xi_3\bar{\delta}\psi_2 - \xi_4\delta\psi_2 + \psi_1(\bar{\delta}\xi_2 - \Delta\xi_4) + \psi_3(\delta\xi_1 - D\xi_3) \\ &= -2\phi\psi_2 + \psi_1[\nu\xi_1 + (2\pi + \bar{\tau} - \alpha - \bar{\beta})\xi_2 - \lambda\xi_3 + (\gamma - \bar{\gamma} + \bar{\mu} - 2\mu)\xi_4] \\ &+ \psi_3[(\bar{\alpha} + \beta - 2\tau - \bar{\pi})\xi_1 - \kappa\xi_2 + (2\rho - \bar{\rho} - \epsilon + \bar{\epsilon})\xi_3 + \sigma\xi_4], \end{aligned} \quad (B15)$$

$$\begin{aligned} &\xi_1\Delta\psi_3 + \xi_2D\psi_3 - \xi_3\bar{\delta}\psi_3 - \xi_4\delta\psi_3 - \psi_4D\xi_3 + \psi_3(\Delta\xi_1 - \bar{\delta}\xi_3) + \psi_2(\bar{\delta}\xi_2 - 2\Delta\xi_4) \\ &= \psi_2[\nu\xi_1 + (3\pi + 2\bar{\tau} - \alpha - \bar{\beta})\xi_2 - 2\lambda\xi_3 + (2\gamma - 2\bar{\gamma} + \bar{\mu} - 3\mu)\xi_4] \\ &+ \psi_3[(\bar{\gamma} - \gamma - \bar{\mu})\xi_1 + (\rho - 2\epsilon)\xi_2 + (\alpha + \bar{\beta} - \bar{\tau})\xi_3 + (2\beta - \tau)\xi_4] + \psi_4[-(\tau + \bar{\pi})\xi_1 + (\bar{\epsilon} - \epsilon + \rho)\xi_3 + \sigma\xi_4], \end{aligned} \quad (B16)$$

$$\begin{aligned} &\xi_1\Delta\psi_4 + \xi_2D\psi_4 - \xi_3\bar{\delta}\psi_4 - \xi_4\delta\psi_4 + 2\psi_4(\Delta\xi_1 - \bar{\delta}\xi_3) + 2\psi_3(\bar{\delta}\xi_2 - \Delta\xi_4) \\ &= 2\phi\psi_4 + 2\psi_3[\nu\xi_1 + (2\pi + \bar{\tau} - \alpha - \bar{\beta})\xi_2 - \lambda\xi_3 + (\gamma - \bar{\gamma} + \bar{\mu} - 2\mu)\xi_4] \\ &+ 2\psi_4[(\bar{\gamma} - \gamma - \bar{\mu})\xi_1 + (\rho - 2\epsilon)\xi_2 + (\alpha + \bar{\beta} - \bar{\tau})\xi_3 + (2\beta - \tau)\xi_4]. \end{aligned} \quad (B17)$$

$$D\phi = \phi;\mu l^\mu, \quad \Delta\phi = \phi;\mu n^\mu, \quad (B9)$$

$$\delta\phi = \phi;\mu m^\mu, \quad \bar{\delta}\phi = \phi;\mu \bar{m}^\mu,$$

and are satisfied by the commutation relation

$$\begin{aligned} \Delta D - D\Delta &= (\gamma + \bar{\gamma})D + (\epsilon + \bar{\epsilon})\Delta \\ &- (\tau + \bar{\pi})\bar{\delta} - (\bar{\tau} + \pi)\delta, \end{aligned} \quad (B10a)$$

$$\delta D - D\delta = (\bar{\alpha} + \beta - \bar{\pi})D + \kappa\Delta - \sigma\bar{\delta} - (\bar{\rho} + \epsilon - \bar{\epsilon})\delta, \quad (B10b)$$

$$\begin{aligned} \delta\Delta - \Delta\delta &= -\bar{\nu}D + (\tau - \bar{\alpha} - \beta)\Delta \\ &+ \bar{\lambda}\bar{\delta} + (\mu - \gamma + \bar{\gamma})\delta, \end{aligned} \quad (B10c)$$

$$\begin{aligned} \bar{\delta}\delta - \delta\bar{\delta} &= (\bar{\mu} - \mu)D + (\bar{\rho} - \rho)\Delta \\ &- (\bar{\alpha} - \beta)\bar{\delta} - (\bar{\beta} - \alpha)\delta. \end{aligned} \quad (B10d)$$

Equation (A10) can be written in the form

$$\begin{aligned} C_{abcd;p}\xi^p + C_{pbcd}\xi^p;a + C_{apcd}\xi^p;b + C_{abpd}\xi^p;c + C_{abc p}\xi^p;d \\ = 2\phi C_{abcd} + \xi^p [C_{rbcd}(\gamma_a{}^r p + \gamma^r{}_{pa}) \\ + C_{arcd}(\gamma_b{}^r p + \gamma^r{}_{pb}) + C_{abrd}(\gamma_c{}^r p + \gamma^r{}_{pc}) \\ + C_{abcr}(\gamma_d{}^r p + \gamma^r{}_{pd})]. \end{aligned} \quad (B11)$$

The independent components of conformal curvature tensor may be expressed in terms of the five complex scalars  $\psi_0, \psi_1, \dots, \psi_4$ . We have

$$\begin{aligned} C_{1212} &= \psi_2 + \bar{\psi}_2, \quad C_{1213} = \psi_1, \quad C_{1223} = -\bar{\psi}_3, \\ C_{1234} &= -\psi_2 + \bar{\psi}_2, \quad C_{1313} = \psi_0, \quad C_{1324} = -\psi_2, \\ C_{1334} &= -\psi_1, \quad C_{2323} = \bar{\psi}_4, \quad C_{2334} = -\bar{\psi}_3, \\ C_{3434} &= \psi_2 + \bar{\psi}_2, \quad C_{1314} = C_{1323} = C_{2324} = 0. \end{aligned} \quad (B12)$$

Making use of Eqs. (B5a)–(B10d), and (B12), Eq. (B11) becomes

- <sup>1</sup>E. Nöther, *Nachr. Akad. Wiss. Goettingen II, Math. Phys. Kl.* **1918**, 235.
- <sup>2</sup>W. R. Davis and M. K. Moss, *Nuovo Cimento* **33**, 1558 (1965).
- <sup>3</sup>G. H. Katzin and J. Levine, *J. Math. Phys.* **9**, 8 (1968).
- <sup>4</sup>G. H. Katzin, J. Levine, and W. R. Davis, *J. Math. Phys.* **10**, 617 (1969).
- <sup>5</sup>C. D. Collinson and C. D. Frech, *J. Math. Phys.* **8**, 701 (1967).
- <sup>6</sup>M. E. Cahill and A. H. Taub, *Commun. Math. Phys.* **21**, 1 (1971).
- <sup>7</sup>A. H. Taub, *General Relativity: Papers in Honour of J. L. Synge*, edited by L. O'Raifeiraigh (Oxford U. P., London, 1972), Chap. VIII, p. 133.
- <sup>8</sup>B. B. Godfrey, *Gen. Relativ. Gravit.* **3**, 3 (1972).
- <sup>9</sup>D. M. Eardley, *Commun. Math. Phys.* **37**, 287 (1974).
- <sup>10</sup>C. W. Misner and J. A. Wheeler, *Ann. Phys. (NY)* **2**, 525 (1957).
- <sup>11</sup>P. A. Goodinson and R. A. Newing, *J. Int. Math. Appl.* **6**, 212 (1970).
- <sup>12</sup>P. A. Goodinson and R. A. Newing, *J. Int. Math. Appl.* **5**, 72 (1969).
- <sup>13</sup>M. L. Woolley, *Commun. Math. Phys.* **31**, 75 (1973).
- <sup>14</sup>M. L. Woolley, *Commun. Math. Phys.* **33**, 135 (1973).
- <sup>15</sup>R. K. Sachs, *Proc. R. Soc. London Ser. A* **264**, 309 (1961).
- <sup>16</sup>N. Tariq and B. O. J. Tupper, *Gen. Relativ. Gravit.* **6**, 345 (1975).
- <sup>17</sup>N. Tariq and B. O. J. Tupper, *Tensor (N.S.)* **28**, 83 (1974).
- <sup>18</sup>G. C. Debney and J. Zund, *Tensor (N.S.)* **22**, 333 (1971).
- <sup>19</sup>G. C. Debney and J. Zund, *Tensor (N.S.)* **25**, 53 (1972).
- <sup>20</sup>K. Yano, *The Theory of Lie Derivatives and its Applications* (North-Holland, Amsterdam, 1957).

# Exact solutions for space-times with local rotational symmetry in which the Dirac equation separates

B. R. Iyer and C. V. Vishveshwara  
Raman Research Institute, Bangalore 560080, India

(Received 5 February 1986; accepted for publication 18 February 1987)

The field equations for the class of perfect fluid space-times with local rotational symmetry in which the authors had earlier shown the Dirac equation separates are studied. For the vacuum and dust cases all possible solutions are exhibited. Other solutions correspond to radiation, a stiff fluid, and a fluid with negative pressure.

## I. INTRODUCTION

In an earlier paper<sup>1</sup> (hereafter referred to as I) we investigated the problem of separability of the Dirac equation in perfect fluid space-times with local rotational symmetry and showed that separation was possible only in a certain subclass of the whole family. The geometrical properties of these space-times were also obtained but the question of the specific space-times in this subclass was left unanswered. In this paper we study the field equations for these particular space-times and attempt to isolate those exact solutions which fall in this category. For the vacuum ( $p = \rho = 0$ ) and dust ( $p = 0$ ) cases all the possible solutions are exhibited. Some exact solutions for other interesting sources like radiation ( $p = \frac{1}{3}\rho$ ), a stiff fluid ( $p = \rho$ ), and fluid with negative pressure ( $p + \rho = 0$ ) are also obtained. Though most of these solutions were known earlier we present a unified and systematic treatment of the different cases of particular interest as background metrics wherein our earlier separation of variables for the Dirac equation is applicable.

In the next section we set up the field equations for the relevant solutions. In Secs. III and IV we obtain all the vacuum and dust solutions, respectively. Section V contains some solutions corresponding to radiation, a stiff fluid, and a fluid with negative pressure.

## II. SPACE-TIMES WITH LOCAL ROTATIONAL SYMMETRY WHEREIN THE DIRAC EQUATION SEPARATES

As demonstrated in I the space-times with local rotational symmetry in which the Dirac equation separates are of the following four types.

Case I:

$$ds^2 = (1/F^2)d\bar{x}^{0^2} - dx^{1^2} - Y^2(dx^{2^2} + t^2 dx^{3^2}), \quad (1a)$$

where

$$F = F(x^1), \quad Y = Y(x^1). \quad (1b)$$

Case III:

$$ds^2 = dx^{0^2} - X^2 dx^{1^2} - Y^2(dx^{2^2} + t^2 dx^{3^2}), \quad (2a)$$

with

$$X = X(x^0), \quad Y = Y(x^0). \quad (2b)$$

Case II a:

$$ds^2 = (1/F^2)d\bar{x}^{0^2} - X^2 d\bar{x}^{1^2} - Y^2(dx^{2^2} + t^2 dx^{3^2}), \quad (3a)$$

where

$$F = F(\bar{x}^0), \quad X = X(\bar{x}^1), \quad Y = Y(\bar{x}^1). \quad (3b)$$

Case II b:

$$ds^2 = (1/F^2)d\bar{x}^{0^2} - X^2 d\bar{x}^{1^2} - Y^2(d\bar{x}^{2^2} + t^2 dx^{3^2}), \quad (4a)$$

where

$$F = F(\bar{x}^0), \quad X = X(\bar{x}^1), \quad Y = Y(\bar{x}^0). \quad (4b)$$

In the above equations  $t$  is one of the four functions

$$(i) \ t = \text{const}, \quad (ii) \ t = x^2, \\ (iii) \ t = \sin(x^2), \quad (iv) \ t = \sinh(x^2). \quad (5)$$

It is clear that the solutions corresponding to  $t = \text{const}$  and  $t = x^2$  are related trivially by transformations from Cartesian to cylindrical coordinates in the  $x^2$ - $x^3$  plane, i.e.,

$$x^{2'} = x^2 \cos(x^3), \quad x^{3'} = x^2 \sin(x^3). \quad (6)$$

Consequently, these two cases can be treated together. Further, in Cases II a and II b, by the following transformation of coordinates

$$x^0 = \int \frac{d\bar{x}^0}{F(\bar{x}^0)}, \quad x^1 = \int X(\bar{x}^1)d\bar{x}^1, \quad (7)$$

the line elements become the following.

Case II a:

$$ds^2 = dx^{0^2} - dx^{1^2} - Y^2(dx^{2^2} + t^2 dx^{3^2}), \quad (8a)$$

where

$$Y = Y(x^1). \quad (8b)$$

Case II b:

$$ds^2 = dx^{0^2} - dx^{1^2} - Y^2(dx^{2^2} + t^2 dx^{3^2}), \quad (9a)$$

with

$$Y = Y(x^0). \quad (9b)$$

In this form Eq. (8) is a special case of (1) with  $F = \text{const}$  while Eq. (9) is a special case of Eq. (2) with  $X = \text{const}$ .

Choosing units  $c = 8\pi G = 1$  and signature  $(+, -, -, -)$  the field equations are

$$G_{ab} = T_{ab}, \quad (10a)$$

where for a perfect fluid

$$T_{ab} = (\rho + p)U_a U_b - p g_{ab}. \quad (10b)$$

Introducing  $\epsilon$  such that

$$\epsilon = \begin{cases} 0, & \text{for } t = \text{const}, \\ +1, & \text{for } t = \sin(x^2), \\ -1, & \text{for } t = \sinh(x^2). \end{cases} \quad (11)$$

Equation (10) for Case I becomes

$$2 \frac{Y_{,11}}{Y} + \frac{Y_{,1}^2}{Y^2} - \frac{\epsilon}{Y^2} = -\rho, \quad (12a)$$

$$2 \frac{Y_{,1} F_{,1}}{YF} - \frac{Y_{,1}^2}{Y^2} + \frac{\epsilon}{Y^2} = -\rho, \quad (12b)$$

$$\frac{F_{,11}}{F} - \frac{Y_{,11}}{Y} - \frac{F_{,1}}{F} \left( 2 \frac{F_{,1}}{F} - \frac{Y_{,1}}{Y} \right) = -\rho. \quad (12c)$$

For Case III one obtains

$$2 \frac{X_{,0} Y_{,0}}{XY} + \frac{Y_{,0}^2}{Y^2} + \frac{\epsilon}{Y^2} = \rho, \quad (13a)$$

$$2 \frac{Y_{,00}}{Y} + \frac{Y_{,0}^2}{Y^2} + \frac{\epsilon}{Y^2} = -\rho, \quad (13b)$$

$$\frac{Y_{,00}}{Y} + \frac{X_{,0} Y_{,0}}{XY} + \frac{X_{,00}}{X} = -\rho. \quad (13c)$$

As mentioned earlier, Case II a corresponds to  $F = \text{const}$  in Eqs. (12) while Case II b corresponds to  $X = \text{const}$  in Eqs. (13).

The field equations should be supplemented by the equation of state for the perfect fluid which we prescribe to be of the form

$$p = (\gamma - 1)\rho. \quad (14)$$

The conservation equation for  $T^{ab}$  gives

$$T^{ab}_{;b} = 0. \quad (15)$$

For Case I Eq. (15) gives

$$\rho_{,0} = p_{,2} = p_{,3} = 0 \quad (16a)$$

and

$$\rho Y^4 F^{\gamma/(\gamma-1)} = \text{const}, \quad (16b)$$

while for Case III we have

$$\rho (XY^2)^\gamma = \text{const} = \rho_0, \quad (17a)$$

$$p_{,1} = p_{,2} = p_{,3} = 0. \quad (17b)$$

Though not useful for the vacuum and dust cases the above "first integrals" are useful in the other cases.

### III. VACUUM SPACE-TIMES ( $\gamma = 1; p = \rho = 0$ )

Case I: If  $F = \text{const}$ , Eqs. (12) become

$$Y_{,1}^2 = \epsilon, \quad Y_{,11} = 0. \quad (18)$$

Thus for  $\epsilon = 0$ , one obtains a flat space-time in Cartesian coordinates while for  $\epsilon = +1$  one finds  $Y^2 = (x^1)^2$ , which is a flat space-time in spherical polar coordinates. There is no solution for  $\epsilon = -1$ .

From Eqs. (12a) and (12b)  $Y = \text{const}$  is possible only if  $\epsilon = 0$ . In this case Eq. (12c) gives

$$F_{,11}/F - 2F_{,1}^2/F^2 = 0, \quad (19)$$

which on integration yields

$$1/F^2 = (x^1)^2 \quad (20)$$

(here and in later parts all trivial integration constants are

transformed away by a suitable translation or scaling). This is just a Minkowski space-time in Rindler coordinates

$$\bar{x}^0 = x^1 \sinh(x^0), \quad \bar{x}^1 = x^1 \cosh(x^0). \quad (21)$$

In general (i.e., if  $F_{,1} \neq 0, Y_{,1} \neq 0$ ) by adding Eqs. (12a) and (12b) and integrating we obtain

$$FY_{,1} = C_1. \quad (22)$$

Equation (12a) can be rewritten as

$$2Y_{,11}/(Y_{,1}^2 - \epsilon) + 1/Y = 0, \quad (23)$$

which on integration gives

$$Y(Y_{,1}^2 - \epsilon) = C_2. \quad (24)$$

Solutions of (22) and (24) satisfy (12c) identically. Hence a solution of Eq. (24) yields a solution of the field equation. From Eq. (24) for  $\epsilon = 0$ , we obtain

$$Y^2 = (x^1)^{4/3}, \quad F^2 = (x^1)^{2/3}, \quad (25)$$

which is the plane symmetric Taub solution<sup>2</sup>

$$ds^2 = z^{-1/2}(dT^2 - dz^2) - z(dx^2 + dy^2); \quad z > 0, \quad (26a)$$

as follows by the transformations

$$T = \left(\frac{3}{4}\right)^{1/3} x^0, \quad z = \left(\frac{3}{4}\right)^{4/3} (x^1)^{4/3}, \quad (26b)$$

$$x = \left(\frac{3}{4}\right)^{-2/3} x^2, \quad y = \left(\frac{3}{4}\right)^{-2/3} x^3.$$

For  $\epsilon = 1$ , the solution may be implicitly given as

$$Y = -(c_2/2)(1 + \cosh(2\rho)),$$

$$F = \pm c_1 \coth \rho, \quad (27)$$

$$\pm x^1 + c_3 = -(c_2/2)(\sinh(2\rho) + 2\rho).$$

On transforming to coordinates  $(x^0, \rho, x^2, x^3)$ , one obtains

$$ds^2 = \tanh^2 \rho dx^{0^2} - 4c_2^2 \cosh^4 \rho d\rho^2$$

$$- c_2^2 \cosh^4 \rho (dx^{2^2} + \sin^2 x^2 dx^{3^2}). \quad (28)$$

This is just a Schwarzschild solution of mass  $c_2/2$  as can be seen by transforming to coordinate  $r$ :

$$r = c_2 \cosh^2(\rho). \quad (29)$$

It is also one of the Levi-Civita degenerate static vacuum solution type AI (Ref. 3).

For  $\epsilon = -1$  one obtains

$$Y = c_2 \sin^2(\rho), \quad F = \pm c_1 \tan(\rho),$$

$$\pm x^1 + c_3 = -(c_2/2)(\sin(2\rho) - 2\rho), \quad (30)$$

which in coordinates  $(x^0, \rho, x^2, x^3)$  give

$$ds^2 = \cot^2 \rho dx^{0^2} - c_2^2 \sin^4 \rho (4d\rho^2 + dx^{2^2} + \sinh^2 x^2 dx^{3^2}). \quad (31)$$

This is the degenerate static vacuum solution due to Levi-Civita<sup>3</sup> which in the classification of Ehlers and Kundt<sup>3</sup> is class AII. In terms of coordinates

$$z = c_2 \sin^2 \rho, \quad (32)$$

$$ds^2 = (c_2/z - 1)dx^{0^2} - ((c_2/z) - 1)^{-1} dz^2$$

$$- z^2(dx^{2^2} + \sinh^2 x^2 dx^{3^2}). \quad (33)$$

Case III: Let us now turn to Eqs. (13). If  $X = \text{const}$  they become

$$Y_{,0}^2 = -\epsilon, \quad Y_{,00} = 0. \quad (34)$$



As for Eqs. (12) for  $\epsilon = 0$  one has a flat space-time in Cartesian coordinates while for  $\epsilon = -1$  one obtains  $Y^2 = x^{0^2}$ . This is just a Milne universe: a flat space-time in Rindler-like coordinates, as can be seen by the transformations

$$\begin{aligned} \bar{x}^0 &= x^0 \cosh x^2, & \bar{x}^1 &= x^1, \\ \bar{x}^2 &= x^0 \sinh x^2 \cos x^3, & \bar{x}^3 &= x^0 \sinh x^2 \sin x^3. \end{aligned} \quad (35)$$

From Eqs. (13a) and (13b)  $Y = \text{const}$  is possible only if  $\epsilon = 0$ . In this case Eq. (13c) gives

$$X_{,00} = 0, \quad \text{i.e., } X = x^0. \quad (36)$$

This again is a flat space-time in Rindler-like coordinates

$$\bar{x}^0 = x^0 \cosh x^1, \quad \bar{x}^1 = x^0 \sinh x^1. \quad (37)$$

We now consider cases when neither  $X$  nor  $Y$  is constant. As before, taking the difference of Eqs. (13a) and (13b) and integrating we obtain

$$X = c_1 Y_{,0}. \quad (38)$$

Equation (13b) on integration gives

$$Y(Y_{,0}^2 + \epsilon) = c_2. \quad (39)$$

Solutions of Eqs. (38) and (39) satisfy Eq. (13c) identically. For  $\epsilon = 0$  one obtains

$$Y = x^{0^{2/3}}, \quad X = (x^0)^{-1/3}, \quad (40)$$

which is a Kasner space-time with local rotational symmetry. The Dirac equation in this case is treated in more detail elsewhere.<sup>4</sup>

For  $\epsilon = +1$  we obtain

$$Y = c_2 \sin^2 T, \quad X = \pm c_1 \cot T, \quad (41)$$

$$\pm x^0 + c_3 = (c_2/2)(2T - \sin 2T),$$

which in terms of coordinates  $(T, x^1, x^2, x^3)$  gives

$$\begin{aligned} ds^2 &= 4c_2^2 \sin^4 T dT^2 - c_1^2 \cot^2 T dx^{1^2} \\ &\quad - c_2^2 \sin^4 T(dx^{2^2} + \sin^2 x^2 dx^{3^2}). \end{aligned} \quad (42)$$

Transforming to

$$\bar{T} = c_2 \sin^2 T, \quad r = c_1 x^1 \quad (43)$$

yields the "inner" sector of the Schwarzschild solution, i.e. ( $r < c_2$ )

$$\begin{aligned} ds^2 &= (c_2/\bar{T} - 1)^{-1} d\bar{T}^2 - (c_2/\bar{T} - 1) dr^2 \\ &\quad - \bar{T}^2(dx^{2^2} + \sin^2 x^2 dx^{3^2}). \end{aligned} \quad (44)$$

For  $\epsilon = -1$ , on the other hand,

$$\begin{aligned} Y &= -c_2 \cosh^2 T, \quad X = \mp c_1 \tanh T, \\ \pm x^0 + c_3 &= (c_2/2)(\sinh 2T + 2T), \end{aligned} \quad (45)$$

which in terms of  $(T, x^1, x^2, x^3)$  yields

$$\begin{aligned} ds^2 &= 4c_2^2 \cosh^4 T dT^2 - c_1^2 \tanh^2 T dx^{1^2} \\ &\quad - c_2^2 \cosh^4 T(dx^{2^2} + \sinh^2 x^2 dx^{3^2}). \end{aligned} \quad (46)$$

Once again going over to

$$\bar{T} = c_2 \cosh^2 T, \quad r = c_1 x^1, \quad (47)$$

we obtain

$$\begin{aligned} ds^2 &= (1 - c_2/\bar{T})^{-1} d\bar{T}^2 - (1 - c_2/\bar{T}) dr^2 \\ &\quad - \bar{T}^2(dx^{2^2} + \sinh^2 x^2 dx^{3^2}). \end{aligned} \quad (48)$$

This solution is to the Levi-Civita static solution AII (Ref. 3), the analog of the  $R < 2M$  region of the Schwarzschild solution.

#### IV. THE DUST SOLUTIONS ( $\gamma = 1, \rho = 0$ )

In Eq. (12) corresponding to Case I for dust, if  $F_{,1} = 0$ , then Eqs. (12b) and (12c) become

$$Y_{,1}^2 - \epsilon = 0, \quad Y_{,11} = 0, \quad (49)$$

which when compared with Eq. (12a), gives  $\rho = 0$ . Thus no dust solutions are possible in this case. Similarly, for  $Y_{,1}^2 = \epsilon$  no dust solutions exist.

In general, however, Eq. (12b) gives

$$F_{,1}/F = (Y_{,1}^2 - \epsilon)/2YY_{,1}. \quad (50)$$

Differentiating (50) and substituting in Eq. (12c) one obtains

$$2Y_{,11}/Y + (Y_{,1}^2 - \epsilon)/Y^2 = 0, \quad (51)$$

which employing (12a) gives  $\rho = 0$ . Thus no dust solution is possible for Eq. (12). They seem to be possible only in metrics of subclass III corresponding to Eq. (13).

If  $X_{,0} = 0$ , Eqs. (13) yield

$$\begin{aligned} Y_{,0}^2 + \epsilon &= \rho, \\ Y_{,00}/Y + (Y_{,0}^2 + \epsilon)/Y^2 &= 0, \\ Y_{,00}/Y &= 0, \end{aligned} \quad (52)$$

which are consistent only for  $\rho = 0$ . Thus one does not have dust solutions with  $X = \text{const}$ . From Eq. (13b)  $Y = \text{const}$  solutions are only possible for  $\epsilon = 0$ , which from (13a) implies  $\rho = 0$ . Thus one does not have such dust solutions either. If  $X_{,0} \neq 0$ ,  $Y_{,0} \neq 0$ , Eq. (13b) can be rewritten as

$$2Y_{,0}Y_{,00}/(Y_{,0}^2 + \epsilon) + Y_{,0}/Y = 0, \quad (53)$$

which gives

$$Y(Y_{,0}^2 + \epsilon) = \text{const}. \quad (54)$$

For  $\epsilon = 0$ , Eq. (54) is solved by

$$Y = (c_1 x^0 + c_2)^{2/3}. \quad (55)$$

Replacing  $Y$  in (13c) from Eq. (55) one obtains for  $X$ , the differential equation

$$X_{,TT} - \frac{1}{3}X_{,T} - \frac{2}{3}X = 0, \quad (56a)$$

where

$$T = \log(c_1 x^0 + c_2). \quad (56b)$$

Consequently, the general solutions for  $X$  is

$$X = [c_3(c_1 x^0 + c_2) + c_4]/(c_1 x^0 + c_2)^{1/3}. \quad (57)$$

Substituting for  $X$  and  $Y$  from Eqs. (55) and (57) in Eq. (13a) yields  $\rho$ :

$$\rho = \frac{4}{3}c_3 c_1^2 / (c_1 x^0 + c_2) [c_3(c_1 x^0 + c_2) + c_4]. \quad (58)$$

By a simple translation and scaling, the metric becomes

$$\begin{aligned} ds^2 &= dx^{0^2} - ((x^0 + \alpha)/x^{0^{1/3}})^2 dx^{1^2} \\ &\quad - x^{0^{4/3}}(dx^{2^2} + dx^{3^2}), \end{aligned} \quad (59a)$$

where

$$\alpha = c_4/c_3 \text{ and } \rho = \frac{4}{3}(x^{0^2} + \alpha x^0)^{-1}. \quad (59b)$$

For the special choice of  $\alpha = 0$ , Eqs. (59) yield

$$ds^2 = dx^{0^2} - x^{0^{4/3}}(dx^{1^2} + dx^{2^2} + dx^{3^2}), \quad (60a)$$

$$\rho = 4/(3x^{0^2}). \quad (60b)$$

This is the Einstein–de Sitter solution for dust which has homogeneous and isotropic spatial sections. However, if  $\alpha \neq 0$  we obtain a more general solution which does not seem obviously equivalent to the  $\alpha = 0$  case. For  $\epsilon = 1$ , the solution may be written in the implicit form

$$Y = c_1 \sin^2 T, \quad (61)$$

$$\pm x^0 + c_2 = (c_1/2)[2T - \sin 2T].$$

Substituting (61) in Eq. (13c) then gives

$$X_{,TT}/X - 2/\sin^2 T = 0. \quad (62)$$

By inspection  $X = \cot(T)$  is a solution to the above equation. To find the other solution let

$$X = V \cot(T) \quad (63)$$

in Eq. (62) so that  $V$  satisfies

$$V_{,TT}/V_{,T} = 2 \csc^2 T / \cot T. \quad (64)$$

The above equation is integrated and finally one has

$$X = c_3(1 - T \cot T) + c_4 \cot T. \quad (65)$$

Substituting (61) and (65) in Eq. (12a) we have

$$\rho = c_3/c_1^2 \sin^4 T [c_3 - \cot T(c_3 T - c_4)]. \quad (66)$$

In terms of  $(T, x^1, x^2, x^3)$  one has

$$ds^2 = 4 \sin^4 T dT^2 - [1 - \cot T(T - c)]^2 dx^{1^2} - \sin^4 T(dx^{2^2} + \sin^2 x^2 dx^{3^2}), \quad (67a)$$

$$\rho = 1/\sin^4 T [1 - \cot T(T - C)]. \quad (67b)$$

Similarly, for  $\epsilon = -1$ ,

$$Y = c_1 \sinh^2 T, \quad (68)$$

$$\pm x^0 + c_2 = (c_1/2)(\sinh 2T - 2T).$$

Substituting into Eq. (13c) gives

$$X_{,TT}/X - 2/\sinh^2 T = 0. \quad (69)$$

As before, since  $X = \coth T$  is a solution of (69) we write

$$X = V \coth T, \quad (70a)$$

$V$  is then a solution of

$$V_{,TT}/V_{,T} = (2 \operatorname{csch}^2 T)/(\coth T), \quad (70b)$$

and consequently

$$X = c_3(T \coth T - 1) + c_4 \coth T. \quad (71)$$

For this case

$$\rho = c_3/c_1^2 \sinh^4 T [c_3(T \coth T - 1) + c_4 \coth T]. \quad (72)$$

In terms of  $(T, x^1, x^2, x^3)$  we thus have

$$ds^2 = 4 \sinh^4 T dT^2 - [(T + \beta) \coth T - 1]^2 dx^{1^2} - \sinh^4 T(dx^{2^2} + \sinh^2 x^2 dx^{3^2}), \quad (73a)$$

$$\rho = 1/\sinh^4 T [(T + \beta) \coth T - 1]. \quad (73b)$$

The dust solutions given by Eqs. (67) and (73) for  $\epsilon = \pm 1$  are those obtained by Kantowski and Sachs.<sup>5</sup>

## V. OTHER SOLUTIONS

If  $F = \text{const}$ , adding twice Eq. (12c) to (12b) one finds that the equations are consistent with (12a) only if  $\rho + 3p = 0$ . This case is not of physical interest. Similarly if  $X = \text{const}$  adding two times (13c) to (13a) one finds that there is consistency with (13b) only for  $\rho = p$ , i.e.,  $\gamma = 2$ . In this case, choosing  $X = 1$  we obtain from Eq. (17)

$$\rho = c_1^2/4Y^4. \quad (74)$$

Substituting in Eq. (13a) and integrating we get

$$x^0 = \int \frac{2Y dY}{\sqrt{c_1^2 - 4\epsilon Y^2}}. \quad (75)$$

For  $\epsilon = 0$ , the solution after suitable scalings give

$$ds^2 = dx^{0^2} - dx^{1^2} - x^0(dx^{2^2} + x^{2^2} dx^{3^2}), \quad (76a)$$

$$\rho = c_1^2/4x^{0^2}. \quad (76b)$$

For  $\epsilon = 1$ , similarly,

$$ds^2 = dx^{0^2} - dx^{1^2} - (c_1^2/4 - x^{0^2})(dx^{2^2} + \sin^2 x^2 dx^{3^2}), \quad (77a)$$

$$\rho = (c_1^2/4)(c_1^2/4 - x^{0^2})^{-2}, \quad (77b)$$

whereas for  $\epsilon = -1$ ,

$$ds^2 = dx^{0^2} - dx^{1^2} - (x^{0^2} - c_1^2/4)(dx^{2^2} + \sinh^2 x^2 dx^{3^2}), \quad (78a)$$

$$\rho = c_1^2/4/(x^{0^2} - c_1^2/4)^2. \quad (78b)$$

The solutions given by Eqs. (76), (77), (78), for a  $\gamma = 2$  fluid is to our knowledge new.

Let us consider Eq. (13) for  $\epsilon = 0$ . Adding  $(\gamma - 1)$  times Eq. (13a) to Eq. (13b) and integrating one gets

$$X = c_1(Y_{,0}^2 Y^\gamma)^{1/2(\gamma-1)}. \quad (79)$$

Since

$$\rho = c(XY^2)^{-\gamma}, \quad (80)$$

one thus gets

$$\rho = c_\gamma = (Y_{,0})^{\gamma/(\gamma-1)} Y^{\gamma(4-3\gamma)/2(\gamma-1)}, \quad (81a)$$

where

$$c_\gamma = c/c_1^\gamma. \quad (81b)$$

Equation (13b) thus becomes

$$2Y_{,00}/Y + Y_{,0}^2/Y^2 = -(\gamma - 1)c_\gamma Y_{,0}^{\gamma/(\gamma-1)} Y^{\gamma(4-3\gamma)/2(\gamma-1)}. \quad (82)$$

The above equation will now be solved for the following interesting physical cases.

(a)  $\gamma = 2$  ( $p = \rho$ ).

For this value the right-hand side of Eq. (82) is proportional to  $Y_{,0}^2/Y^2$ . Thus integrating (82) yields

$$Y = (c_2 x^0 + c_3)^{1/(1+\alpha)}, \quad (83)$$

where

$$2\alpha = 1 + c/c_1^2.$$

Then Eq. (79) gives  $X$  as

$$X = (c_1/c_2)(1 + \alpha)(c_2x^0 + c_3)^{(\alpha-1)/(\alpha+1)}. \quad (84)$$

After the usual scalings one thus has

$$ds^2 = dx^{0^2} - (x^0)^{2(\alpha-1)/(\alpha+1)} dx^{1^2} - (x^0)^{2/(1+\alpha)}(dx^{2^2} + dx^{3^2}), \quad (85a)$$

$$\rho = c/c_1^2(1 + \alpha)^2x^{0^2}. \quad (85b)$$

This solution is identical to one of the solutions in Vajk and Eltgroth.<sup>6</sup>

$$(b) \gamma = \frac{4}{3} \quad (\rho = \frac{1}{3}\rho).$$

In this case  $\rho$  is proportional to  $Y_{,0}^4$  and hence Eq. (82) becomes

$$2Y_{,00}/Y + Y_{,0}^2/Y^2 = -\beta Y_{,0}^4, \quad (86a)$$

where

$$\beta = \frac{1}{3}c_1^{-4/3}. \quad (86b)$$

Substituting  $YY_{,0}^2 = u$  into the above equation and integrating one obtains

$$u = YY_{,0}^2 = (c_2 + \beta Y)^{-1}, \quad (87)$$

whose solution may be written as

$$Y = (c_2/\beta)\sinh^2 T, \\ \pm x^0 + c_3 = (c_2^2/16\sqrt{\beta^3})[\sinh 4T - 4T], \quad (88)$$

$$X = c_1c_2\sqrt{\beta}(\cosh^3 T)/(\sinh T).$$

In terms of  $(T, x^1, x^2, x^3)$  the space-time is described by

$$ds^2 = (c_2^4/4\beta^3)\sinh^4 2T dT^2 - \cosh^4 T \coth^2 T dx^{1^2} - \sinh^4 T(dx^{2^2} + dx^{3^2}), \quad (89a)$$

$$\rho = 16\beta^2 c_1/c_2^4 \sinh^4 2T. \quad (89b)$$

Like the earlier case, this is also a particular solution from Vajk and Eltgroth.<sup>6</sup>

$$(c) \gamma = 0 \quad (\rho + \rho = 0).$$

For this value of  $\gamma$ ,  $\rho = \text{const} = \rho_0$  and  $X = c_1 Y_{,0}$ . Thus Eq. (82) yields

$$2Y_{,00}/Y + Y_{,0}^2/Y^2 = \rho_0. \quad (90)$$

The above equation can be integrated by letting  $Y_{,0}/Y = u$ . We get

$$Y = c_3 [\cosh [\sqrt{3\rho_0}(x^0 + c_2)/2]]^{2/3}, \\ X = c_1c_3\sqrt{\rho_0/3} [\cosh [\sqrt{3\rho_0}(x^0 + c_2)/2]]^{-1/3} \\ \times (\sinh [\sqrt{3\rho_0}(x^0 + c_2)/2]). \quad (91)$$

Thus the metric may be written as

$$ds^2 = dx^{0^2} - (\cosh(\sqrt{3\rho_0}x^0/2))^{-2/3} \sinh^2(\sqrt{3\rho_0}x^0/2) dx^{1^2} - (\cosh(\sqrt{3\rho_0}x^0/2))^{4/3} (dx^{2^2} + dx^{3^2}). \quad (92)$$

To the best of our knowledge Eq. (92) is a new solution. For the various values of  $\gamma$  dealt with above we have not been able to obtain solutions of Eq. (13) for  $\epsilon = \pm 1$  or of Eq. (12) for  $\epsilon = 0, +1$ .

In the foregoing we have systematically obtained the various exact solutions with local rotational symmetry in which the Dirac equation is separable. As was mentioned at the outset many of them turn out to be already known solutions sometimes in terms of unconventional coordinates. Other solutions, given by Eqs. (59), (76), (77), (78), and (92), are new as far as we know. Our results, while incorporating a regular classification of these space-times would also facilitate the study of the Dirac equation in backgrounds exhibiting local rotational symmetry.

## ACKNOWLEDGMENT

It is a pleasure to thank Professor M. A. H. MacCallum for making available an algebraic computational program with the help of which we have checked our calculations.

<sup>1</sup>B. R. Iyer and C. V. Vishveshwara, *J. Math. Phys.* **26**, 1034 (1985).

<sup>2</sup>D. Kramer, H. Stephani, M. MacCallum, and E. Herlt, *Exact Solutions of Einstein's Field Equations* (Cambridge U. P., Cambridge, 1980), p. 159.

<sup>3</sup>Reference 2, p. 188.

<sup>4</sup>B. R. Iyer, *Phys. Lett. A* **112**, 313 (1985).

<sup>5</sup>R. Kantowski and R. K. Sachs, *J. Math. Phys.* **7**, 443 (1966).

<sup>6</sup>J. P. Vajk and P. G. Eltgroth, *J. Math. Phys.* **11**, 2212 (1970); see, also, Ref. 2, p. 146.

# Comment on the two “new” classes of Bianchi type II solutions

Dieter Lorenz-Petzold

*Fakultät für Physik, Universität Konstanz, D-7750 Konstanz, Federal Republic of Germany*

(Received 24 June 1986; accepted for publication 11 February 1987)

Hajj-Boutros's claim [J. Math. Phys. **27**, 1592 (1986)] that two new classes of Bianchi type II solutions can be generated from Lorenz's solution [Phys. Lett. A **79**, 19 (1980)] is shown to be wrong.

In a recent paper by Hajj-Boutros<sup>1</sup> the locally rotationally symmetric (LRS) Bianchi type II stiff matter solution derived by Lorenz<sup>2</sup> (and independently by Ruban<sup>3</sup>) has been reconsidered. Our solution is the unique stiff matter solution of Eqs. (2.6)–(2.8).<sup>1</sup> This solution includes the Taub<sup>4</sup> vacuum solution as a special case. It is an easy matter of calculation to derive the corresponding non-LRS solution.<sup>5</sup> The LRS case is given by Eqs. (3.1)–(3.5).<sup>1</sup>

We make the following comments. The crucial equation in the paper of Hajj-Boutros<sup>1</sup> is given by Eq. (2.9). By taking  $R = R(t)$ ,  $S = S(t)$ , and  $r = \dot{R}/R$  (or  $\dot{S}/S$ ) this equation can be considered as a Riccati equation in  $r$ . Hajj-Boutros finds that from a known solution  $r_0 = \dot{R}_0/R_0$  ( $\dot{S}_0/S$ ) some new solutions (of Bianchi type II) are given by Eqs. (2.14)–(2.17). However, this idea is entirely wrong. First of all the Bianchi type II stiff matter solution  $R = R(\tau)$ ,  $S = S(\tau)$  [see Eqs. (3.1)–(3.4)], where the temporal variable  $\tau$  is related to  $t$  by the relation  $dt = SR^2 d\tau$ , is not a “particular” solution of the field equations (2.4)–(2.8): it is the most general (LRS) solution of Eq. (2.9)! Equation (2.9) can be rewritten in the simple form

$$(R'/R + S'/S)' = 0, \quad (1)$$

where  $( )' = d( )/d\tau$ . The most general solution is given by

$$(SR)^2 = \exp 2(q\tau + \phi), \quad q, \phi = \text{const.} \quad (2)$$

Moreover, our solution is given by  $R = R(\tau)$  and  $S = S(\tau)$  and not by  $R = R(t)$  and  $S = S(t)$ . Thus Eqs. (3.6) and (3.7) of Ref. 1 are meaningless! This is the main error made by Hajj-Boutros. By using  $dt = SR^2 d\tau$  the solutions given can be reexpressed (at least in principle) in  $t$  time.

For the sake of completeness we also mention the recent critical remarks of MacCallum<sup>6</sup> concerning various “new” and incorrect Bianchi type II solutions.

<sup>1</sup>J. Hajj-Boutros, J. Math. Phys. **27**, 1592 (1986).

<sup>2</sup>D. Lorenz, Phys. Lett. A **79**, 19 (1980).

<sup>3</sup>V. A. Ruban, Preprint 412, Leningrad Institut of Nuclear Physics B. P. Konstantinova, 1978.

<sup>4</sup>A. H. Taub, Ann. Math. **53**, 472 (1951).

<sup>5</sup>D. Lorenz, Phys. Lett. A **80**, 235 (1980).

<sup>6</sup>M. A. H. MacCallum, Gen. Relativ. Gravit. **17**, 659 (1985).

# Nonexistence of static conformally flat solutions in a new scalar-tensor theory

Tarkeshwar Singh

Department of Mathematics, Shree Ramdeobaba Kamala Nehru Engineering College, Gittikhadan, Katol Road, Nagpur 440013, India

(Received 28 October 1986; accepted for publication 11 February 1987)

For the special case when the scalar field is massless and conformally invariant, it is shown that there do not exist spherically symmetric conformally flat solutions in a new scalar-tensor theory proposed by Schmidt *et al.* [Phys. Rev. D **24**, 1484 (1981)] representing disordered radiation and in the presence of a source-free electromagnetic field except for the trivial empty flat space-time of Einstein's theory. The solution in the vacuum case is also only a flat space-time of Einstein's theory.

## I. INTRODUCTION

Schmidt *et al.*<sup>1</sup> proposed a new scalar-tensor theory of gravitation, where the gravitational constant depends on a scalar field which itself couples to the surrounding masses through the curvature scalar. The idea was to obtain a possible stable configuration as a final situation in the history of a collapsing object due to the generation of a strong scalar field. The result was contrary to what was expected. In the new theory a mass term was added and an arbitrary coupling constant  $\beta$  between the scalar field  $\phi$  and the curvature invariant  $R$  was also allowed. The theory was subsequently applied to a Friedmann-Robertson-Walker universe by Banerjee and Santos.<sup>2</sup> Further, Singh and Singh<sup>3</sup> have shown that the spatially homogeneous stationary perfect fluid cosmological model in this theory cannot include the radiation-filled universe or the empty universe at the limit in the presence of a massive scalar field. Banerjee *et al.*<sup>4</sup> have discussed a stiff fluid Bianchi type I cosmological model in this theory by considering the cosmological constant  $\Lambda$  and the mass term both being equal to zero. Finally, they have considered some special cases for  $\beta \geq 1$  and have shown that solutions for matter-free space include the one previously found by Accioly *et al.*<sup>5</sup> for the conformally invariant scalar field ( $\beta = 1$ ). Very recently, the author<sup>6</sup> has shown that an analog of the Birkhoff theorem in general relativity exists in this new scalar-tensor theory for the special case when the scalar field is massless and independent of time.

In the present case we apply this general theory to the static spherically symmetric conformally flat space-time for the special case when the scalar field  $\phi$  is massless and conformally invariant ( $\beta = 1$ ) (see Refs. 7 and 8).

## II. FIELD EQUATIONS AND THEIR SOLUTIONS

The gravitational field equations in the scalar-tensor theory proposed by Schmidt *et al.*<sup>1</sup> are given by

$$(\gamma - (\beta/12)\phi^2)G_{ij} = -\frac{1}{2}T_{ij} - \frac{1}{2}[\phi_{,i}\phi_{,j} - \frac{1}{2}g_{ij}(\phi_{,k}\phi^{,k} - \mu^2\phi^2)] + (\beta/12)[(\phi^2)_{;ij} - g_{ij}(\phi^2)_{;k}{}^k], \quad (1)$$

and the wave equation is

$$\square\phi + [\mu^2 + (\beta/6)R]\phi = 0. \quad (2)$$

Here  $\mu$  is the mass of the scalar field,  $\beta$  is an arbitrary coupling constant, and  $\gamma = c^2/16\pi G$  is half of the inverse gravitational constant. The effective inverse gravitational coupling in this theory becomes

$$\gamma_{\text{eff}} = \gamma - (\beta/12)\phi^2$$

and the effective mass of the scalar field is now

$$\mu_{\text{eff}} = [\mu^2 + (\beta/6)R]^{1/2}.$$

We consider the static spherically symmetric conformally flat metric in the form

$$ds^2 = e^\alpha(-dr^2 - r^2 d\theta^2 - r^2 \sin^2 \theta d\Phi^2 + dt^2), \quad (3)$$

where  $\alpha$  is the function of  $r$  alone. We assume that the scalar field  $\phi$  also has spherical symmetry.

Taking  $\phi$  as a function of  $r$  only and using (3) in (1) and (2), the explicit field equations for the new scalar-tensor theory can be written as

$$\frac{1}{2}e^\alpha T_1^1 = -(\gamma - (\beta/12)\phi^2)(3\alpha'^2/4 + 2\alpha'/r) + \frac{1}{4}(\phi')^2 - \frac{1}{4}e^\alpha \mu^2 \phi^2 + (\beta/12)[(2/r)(\phi^2)' + 3\alpha'(\phi^2)'/2], \quad (4)$$

$$\frac{1}{2}e^\alpha T_2^2 = -(\gamma - (\beta/12)\phi^2)(\alpha'' + \alpha'^2/4 + \alpha'/r) - \frac{1}{4}(\phi')^2 - \frac{1}{4}\mu^2 \phi^2 e^\alpha + (\beta/12)[(\phi^2)'' + (\phi^2)'(\alpha' + 1/r)] = \frac{1}{2}e^\alpha T_3^3, \quad (5)$$

$$\frac{1}{2}e^\alpha T_4^4 = -(\gamma - (\beta/12)\phi^2)(\alpha'' + \alpha'^2/4 + 2\alpha'/r) - \frac{1}{4}(\phi')^2 - \frac{1}{4}\mu^2 \phi^2 e^\alpha + (\beta/12)[(\phi^2)'' + 2(\phi^2)'/r + (\phi^2)'\alpha'/2], \quad (6)$$

$$e^{-\alpha}[(\phi^2)'' + (\phi^2)'(\alpha' + 2/r)] = \frac{(\beta/6)\phi}{\gamma + \beta(\beta-1)\phi^2/12} \left\{ \frac{6}{\beta} \mu^2 \gamma + \frac{1}{2} T + \frac{1}{2} \mu^2 \phi^2 - \frac{(1-\beta)}{2} e^{-\alpha} (\phi')^2 \right\}. \quad (7)$$

Here a prime indicates differentiation with respect to  $r$ .

## III. SOLUTIONS OF THE FIELD EQUATIONS

As the field equations are nonlinear, the problem becomes difficult, in the general case with the nonvanishing

mass term  $\mu$ . So the solutions in a special case for a massless conformally invariant scalar field ( $\mu = 0, \beta = 1$ ) are presented. The question of overdeterminacy is settled by satisfying all the field equations by the actual substitution of the solution obtained. Before proceeding further we simply assume  $\mu = 0$  and  $\beta = 1$ .

### A. Vacuum solutions

It can be easily verified that when the scalar field  $\phi$  is a constant and  $T_{ij} = 0$ , the field equations (4)–(7) yield a solution which describes an empty flat space-time of Einstein's theory, and when  $T_{ij} = 0$ , the field equations (4)–(7) reduce to the vacuum case.

From Eqs. (5)–(7) we can easily obtain

$$e^\alpha = k / (12\gamma - \phi^2), \quad (8a)$$

$$\phi^2 = 12\gamma - \exp(k_1/kr + k_2), \quad (8b)$$

where  $k, k_1$ , and  $k_2$  are integration constants. Solution (8) satisfies Eq. (4) only when  $k_1 = 0$ . This, in view of (8), gives

$$\phi^2 = \text{const}, \quad \alpha = \text{const}.$$

Thus the solution of Eqs. (4)–(7) is

$$\alpha = \text{const}, \quad \phi = \text{const}. \quad (9)$$

Hence the only spherically symmetric static conformally flat solution of the new scalar–tensor theory<sup>1</sup> is simply the empty flat space-time of Einstein's theory, when the scalar field is massless and conformally invariant.

### B. Electrovac solution

Here we consider the energy momentum tensor for a trace-free electromagnetic field in the form

$$T_{ij} = F_{ij}F_j^i - \frac{1}{2}g_{ij}F_{\prime\prime\prime}F^{\prime\prime\prime}, \quad (10)$$

where  $F_{ij}$  is the electromagnetic field tensor satisfying

$$F_{\prime j}^{\prime j} = 0 \quad (11)$$

and

$$F_{\prime j, k} + F_{\prime k, i} + F_{\prime k, j} = 0. \quad (12)$$

For a static charged particle the only nonzero component of the electromagnetic field  $F_{ij}$  is  $F_{14}$ . Equations (11) and (12) now lead to  $F_{14} = q/r^2$ , where  $q$  is a constant which can be identified with the electric charge of the particle.

With metric (3) the nonvanishing components of the energy momentum are

$$-T_1^1 = T_2^2 = T_3^3 = -T_4^4 = -\frac{1}{2}(q^2/r^4)e^{-2\alpha}. \quad (13)$$

Using (13), the field equations (5)–(7) admit the solution

$$e^\alpha = (c_1r^2 + c_3)/(\phi^2 - 12\gamma)r^2, \quad (14a)$$

$$\phi^2 = 12\gamma - \exp\left\{\frac{c_2}{(-c_1c_3)^{1/2}} \tanh^{-1} \frac{r(-c_1c_2)^{1/2}}{c_3} + \phi_0\right\}, \quad (14b)$$

where  $\phi_0, c_1$ , and  $c_2$  are constants of integration and  $c_3$  is set equal to  $-3q^2$ . On actual verification it was found that solution (14) satisfies each of the field equations only when  $c_2 = 0$  and  $c_3 = 0$ , which, in turn reduces the field equations (4)–(7) to Einstein's vacuum case.

Hence the only spherically symmetric static conformally flat solution of the new scalar–tensor theory in the pres-

ence of a source-free electromagnetic field is the empty flat space-time of Einstein's theory.

### C. Disordered radiation

Here we consider the energy momentum tensor due to that of a perfect fluid distribution in the form

$$T_{ij} = (\rho + p)u_i u_j - pg_{ij}, \quad (15)$$

with equation of state

$$\rho = 3p, \quad (16)$$

where  $\rho$  is energy density and  $p$  is the pressure of the fluid. From (3), (15), and (16) the components of  $T_j^i$  in comoving coordinates are

$$T_j^i = \text{diag}(-p, -p, -p, 3p). \quad (17)$$

The conservation equation  $T_{j,i}^i = 0$  leads to

$$\frac{dp}{dr} + (\rho + p) \frac{\alpha'}{2} = 0. \quad (18)$$

Using Eqs. (15)–(17) in field equations (4)–(7), one gets the field equations of a new scalar–tensor theory with disordered radiation.

Now using (8) in the difference of Eqs. (5) and (6) and using (17) in the sum of Eqs. (5) and (6) one easily gets  $p = 0$  and therefore  $\rho = 0$ . This leads to the vacuum field equations in which case, as shown in Sec. III A, the only solution is the flat space-time of Einstein's theory.

Thus there are no spherically symmetric conformally flat solutions of the new scalar–tensor theory<sup>1</sup> representing disordered radiation in the presence of massless conformally invariant scalar field.

### IV. CONCLUSION

In order to understand fully the scalar–tensor theories of gravitation it is useful to have a knowledge of some exact solutions of these equations. The search for an analytic solution is important due to the fact that once such a solution is obtained one can study all of its physical properties. Exact static spherically symmetric conformally flat solutions in vacuum, in the presence of an electromagnetic field and for disordered radiation, are considered in a new scalar–tensor theory proposed by Schmidt *et al.*<sup>1</sup> for the special case when the scalar field is massless and conformally invariant. It is observed that the only spherically symmetric conformally flat solution in this new scalar–tensor theory is the flat space-time of Einstein's theory.

### ACKNOWLEDGMENTS

The author is thankful to Professor A. A. Kayande for his encouragement. The author is also thankful to Dr. T. Singh, Department of Applied Mathematics, Banaras Hindu University, Varanasi for his valuable suggestions.

<sup>1</sup>G. Schmidt, W. Greiner, U. Heinz, and B. Müller, Phys. Rev. D **24**, 1484 (1981).

<sup>2</sup>A. Banerjee and N. O. Santos, Phys. Rev. D **26**, 3747 (1982).

<sup>3</sup>T. Singh and T. Singh, Indian J. Phys. **58** B, 41 (1984).

<sup>4</sup>A. Banerjee, A. K. G. De Oliveira, and N. O. Santo, Gen. Relativ. Gravit. **17**, 371 (1985).

<sup>5</sup>A. J. Accioli, A. N. Vaidya, and M. M. Som, Phys. Rev. D **27**, 2282 (1983).

<sup>6</sup>T. Singh, Phys. Rev. D **34**, 646 (1986).

<sup>7</sup>F. Gürsey, Ann. Phys. (NY) **24**, 211 (1963).

<sup>8</sup>R. Penrose, Proc. R. Soc. London Ser. A **284**, 159 (1965).

# Series representations for calculations in quantum statistics. II

William A. Barker

Department of Physics, Santa Clara University, Santa Clara, California 95053

(Received 12 August 1986; accepted for publication 30 January 1987)

The methods developed in the first paper of this series (I) [J. Math. Phys. **27**, 302 (1986)] are used to calculate the pressure in several cases of physical interest: conduction electrons inside and outside a heated tungsten cathode; helium atoms in liquid helium II and He vapor; electrons in a white dwarf star, and photons in a cavity. The latter two cases involve ultrarelativistic particles whose energy  $\epsilon = pc$ . The formulation in I is extended to include particle speeds comparable to  $c$ . The grand potential is used to establish two useful, exact relationships between pressure and energy density:  $P = 2u/3$  and  $P = u/3$ , valid for  $v \ll c$  and  $v \approx c$ , respectively.

## I. INTRODUCTION

In an earlier paper,<sup>1</sup> hereafter referred to as I, the three distribution functions of statistical physics are treated in a similar manner. In each case a single integral suffices. The quantity  $C\epsilon^p$  is averaged over the Maxwell-Boltzmann (MB), Bose-Einstein (BE), and Fermi-Dirac (FD) distributions. Here  $C$  is a constant and  $\epsilon$  is the particle energy raised to any power  $p$ . Reference 1 develops simple schemes for evaluating these integrals. The MB integral [Eq. (2)] is written in closed form in terms of the gamma function of  $p + 1$ . The BE integral [Eq. (8)] is the product of the gamma function and the Riemann zeta function of  $p + 1$ . The FD integral [Eq. (27)] is written as a rapidly converging series of derivatives for any value of  $p$ . In the literature, some calculations, using the quantum statistics of highly degenerate gases, are quite formidable. Now it is really quite easy. The physical properties of various systems of interest can be investigated with very little mathematical effort.

The purpose of this paper is to illustrate this point for seven cases of physical interest. We focus on a single thermodynamic coordinate, the pressure, and we calculate this quantity for the following: (1) the conduction electrons inside a heated tungsten cathode, (2) the conduction electrons in the evacuated region between a heated tungsten cathode and an anode, (3) the helium atoms in liquid helium II at 1 K, (4) the helium atoms in helium vapor at 1 K, (5) the ultrarelativistic electrons in a white dwarf star, (6) a photon gas in a cavity, and (7) a laser beam.

## II. GENERAL RELATIONSHIPS BETWEEN PRESSURE AND ENERGY. THE GRAND POTENTIAL

The grand potential<sup>2</sup>

$$\Omega = \sum_k \Omega_k = \mp kT \sum_k \ln [1 \pm e^{(\mu - \epsilon_k)/kT}] \quad (1)$$

is applicable in determining various thermodynamic properties of a Fermi, Bose, or classical gas. Here  $\Omega_k$  is the grand potential for the energy  $\epsilon_k$ . The familiar distribution function of a particle in the  $k$ th state is obtained by

$$n_k = - \frac{\partial \Omega_k}{\partial \mu} = \frac{1}{e^{(\epsilon_k - \mu)/kT} \pm 1}. \quad (2)$$

If  $\mu/kT$  is large and negative, as it is in the classical case,

$$n_k = 1/e^{(\epsilon_k - \mu)/kT}. \quad (3)$$

Formally the energy  $U$  for each of these three cases, in the nonrelativistic limit, is

$$U = C \int_0^\infty \frac{\epsilon^{3/2} d\epsilon}{e^{(\epsilon - \mu)/kT} \pm 1,0}, \quad (4)$$

where  $C = (2S + 1)(V/4\pi^2)(2m/\hbar^2)^{3/2}$ . Here  $S = \text{spin}$ ,  $V = \text{volume}$ ,  $m = \text{particle mass}$ , and  $\hbar = \text{Planck's constant divided by } 2\pi$ .

Now convert (1) from a sum to an integral, integrate by parts, and compare with (4). We find that

$$\Omega = -2U/3. \quad (5)$$

Using<sup>3</sup>  $\Omega = -PV$ , it follows that

$$PV = 2U/3 \quad (6)$$

for nonrelativistic fermions, bosons, and classical particles. If, on the other hand, the particles are ultrarelativistic, then  $\epsilon = pc$  replaces  $\epsilon = p^2/2m$ . The constant  $C$  in Eqs. (2), (6), and (16) of I should be replaced by

$$C' = (2S + 1)V/2\pi^2\hbar^3c^3. \quad (7)$$

This follows from the changes in the volume element  $4\pi p^2 dp$  in momentum space. There are corresponding changes in the integrands which are taken care of by the proper choice for the power  $p$ . The expression for the grand potential now becomes

$$\Omega = - \frac{C'}{3} \int_0^\infty \frac{\epsilon^3 d\epsilon}{e^{(\epsilon - \mu)/kT} \pm 1,0} = - \frac{U}{3}. \quad (8)$$

Therefore,

$$PV = U/3 \quad (9)$$

for ultrarelativistic particles.

## III. THE PRESSURE OF A GAS OF CONDUCTION ELECTRONS INSIDE AND OUTSIDE A HEATED TUNGSTEN CATHODE. $S = \frac{1}{2}$

*Case 1 (Inside):* The conduction electrons collide frequently with phonons. Consequently their ambient speed is very small indeed compared to  $c$ .

It is easy to show that the conduction electrons inside a tungsten cathode constitute a degenerate Fermi gas: Compare  $n_Q$ , the quantum concentration,<sup>4</sup> with  $n$ , the actual concentration of electrons in tungsten,

$$n_Q = (mkT/2\pi\hbar^2)^{3/2} = 2.16 \times 10^{20}/\text{cm}^3. \quad (10)$$

Here  $m$  is the electron mass,  $k$  is Boltzmann's constant, and  $T$  is the temperature of the heated cathode, taken here to be 2000 K. On the other hand

$$n = fN_0\rho/M = 1.27 \times 10^{23}/\text{cm}^3, \quad (11)$$

with  $f$  = the number of free electrons per atom = 2,  $\rho$  = density of tungsten = 19.35 g/cm<sup>3</sup>,  $M$  = atomic weight = 183.85, and  $N_0$  = Avogadro's number. We find, as expected, that  $n/n_Q = 588 \gg 1$ .

We seek, therefore, the pressure of a nonrelativistic, degenerate gas of electrons of concentration  $n$ . We use  $PV = 2U/3$  from (6) and the FD expression Eq. (30) for  $U$  from I,

$$P = (3\pi^2)^{2/3} \frac{\hbar^2}{5m} n^{5/3} \left( 1 + \frac{5\pi^2}{12} \left( \frac{T}{T_F} \right)^2 + \dots \right). \quad (12)$$

The contribution of the first term in (16) is  $P_0$ , the pressure at absolute zero. Numerically  $P_0 = 7.40 \times 10^5$  atm, with  $n$  taken from (11). The increase,  $\Delta P$ , in the pressure due to the second term in (12) is  $1.05 \times 10^3$  atm. Here  $T_F$  = the Fermi temperature =  $\mu_0/k = 1.072 \times 10^5$  K for tungsten. The very large value for the pressure is due to the high concentration,  $n$ , of conduction electrons, which is raised to the 5/3 power in (12).

*Case 2 [Outside (halfway between the cathode and the anode)]:* The conduction electrons which are boiled off of the cathode are clearly nonrelativistic since  $e\phi_b \ll mc^2$ , where  $\phi_b$  is the anode voltage.

It is easy to show that the conduction electrons in the evacuated region between the cathode and the anode constitute a classical MB gas. Again we need  $n_Q$  and  $n$ . The quantum concentration  $n_Q = 1.25 \times 10^{19}/\text{cm}^3$  where we have used a nominal  $T = 300$  K in (10) instead of  $T = 2000$  K. This is justified because a vacuum is a good thermal insulator.

For simplicity, we use parallel plate geometry and the assumption of space charge limitation. The solution is very well known.<sup>5</sup> It leads to the Child-Langmuir law<sup>6</sup> and enables us to calculate  $n$ .

The current density

$$J = nev = 2.33 \times 10^{-6} \phi_b^{3/2}/d^2 = 2.16 \times 10^{20}/\text{cm}^3 \quad (13)$$

if we choose  $\phi_b$  = plate voltage = 50 V with  $d$  = cathode-plate separation =  $2 \times 10^{-2}$  m.

The potential

$$K_1 x^{4/3} = 19.8 V \quad (14)$$

for  $K_1 = 9.18 \times 10^3 \text{ V/m}^{4/3}$ , if we take  $x = d/2$ . The electron speed follows from

$$mv^2/2 = e\phi. \quad (15)$$

We find  $v = 2.64 \times 10^6$  m/sec<sup>-1</sup> and  $n = 4.88 \times 10^6$  cm<sup>-3</sup>. We see that  $v/c = 8.8 \times 10^{-3} \ll 1$  and that  $n/n_Q$

=  $3.90 \times 10^{-13} \ll 1$ . The electron gas is nonrelativistic and nondegenerate as expected.

The pressure of this gas follows readily from  $PV = 2U/3$  and  $U = 3NkT/2$  [I, Eq. (5)]. We obtain the familiar equation of state for an ideal classical gas

$$P = nkT = 2 \times 10^{-13} \text{ atm}. \quad (16)$$

Comparing Cases 1 and 2, it is interesting that conduction electrons boiled off from the cathode represent a minute fraction,  $3.8 \times 10^{-17}$ , of the conduction electrons inside a hot cathode. The outside electron gas pressure is  $2.7 \times 10^{-19}$  times the inside gas pressure. It is even 16 orders of magnitude less than the pressure due to the small  $T$ -dependent term in (12).

#### IV. THE PRESSURE OF A GAS OF HELIUM ATOMS. $S=0$

*Case 3 (Liquid helium at 1 K):* We expect that this system will be a nonrelativistic Bose gas. The Einstein condensation temperature<sup>7</sup> in He<sup>4</sup> occurs at  $T = 2.174$  K. This is the temperature at which the number of atoms in excited states equals the total number of atoms. As the temperature is lowered, there is a condensation of atoms into states of zero momentum. At 1 K only 20% of the normal component of He<sup>4</sup> remains,<sup>8</sup> 80% is superfluid with no velocity. The normal atoms have speeds<sup>9</sup>  $\sim 50$  m sec<sup>-1</sup>  $\ll c$ .

The density of liquid helium at 1 K is 0.145 g/cm<sup>3</sup>. To find the concentration  $n$ , we use (11) with  $f = 1$ ,  $M = 4$ . Thus  $n = 2.18 \times 10^{22}/\text{cm}^3$ . To calculate the quantum concentration  $n_Q$ , we use (10) with  $m$  = helium atomic mass =  $6.64 \times 10^{-24}$  g and  $T = 1$  K. Thus  $n_Q = 1.52 \times 10^{21}/\text{cm}^3$  and  $n/n_Q = 14.3 > 1$ .

The helium gas is a nonrelativistic, degenerate boson gas, as expected. To find the pressure we use  $PV = 2U/3$  with  $U$  taken from I, Eq. (14):

$$P = 1.341 n_Q kT. \quad (17)$$

The numerical factor in (17) is the Riemann zeta function  $\zeta(5/2) = \sum_{m=1}^{\infty} m^{-5/2}$ . Numerically  $P = 0.278$  atm. We note that the pressure of the degenerate He gas is independent of the concentration of helium atoms, varies as  $T^{5/2}$  since  $n_Q \sim T^{3/2}$ , and is slightly greater than  $\frac{1}{4}$  of one atmosphere. Inasmuch as we are using an ideal gas model to describe a liquid, the results can only be considered qualitatively correct. However, if we use Eq. (11) of I to calculate the Einstein temperature on this model it turns out to be 3.15 K as compared to the experimental value of 2.174 K. Accordingly, we might expect the pressure results to be correct within a factor of 2 or so.

*Case 4 (The vapor pressure of helium atoms at 1 K):* We expect that this gas will be nonrelativistic and nondegenerate. The vapor pressure of helium is important in measuring temperature. Consequently it has received a lot of attention and can be calculated in several ways. We chose the empirical equation of Clement, Logan, and Gaffney,<sup>10</sup>

$$\ln P = I - A/T + B \ln T + CT^2/2 - D((\alpha\beta/\beta^2 - 1) - T^{-1})\tan^{-1}(\alpha T - \beta). \quad (18)$$

Here  $P$  is measured in mm of Hg. The constants are  $I = 4.6202$ ,  $A = 6.399$ ,  $B = 2.541$ ,  $C = 0.00612$ ,



$D = 0.5197$ ,  $\alpha = 7.00$ , and  $\beta = 14.14$ . This expression, which lacks the sixth term given by Clement *et al.*, is of sufficient accuracy for our purposes. For a given  $T$ , we can compute the pressure which we use in the following statistical calculation.

At the lambda point,  $T = 2.174$  K and  $P = 38.47$  mm of Hg. At  $T = 1$  K,  $P = 0.116$  mm of Hg, a decrease by a factor of 332 from the vapor pressure at the condensation temperature.

At 1 K, using (18) and assuming the gas is classical,  
 $P = nkT = 1.53 \times 10^{-4}$  atm (0.116 mm of Hg). (19)

From (19),  $n = 1.20 \times 10^{18}/\text{cm}^3$ . We calculated the quantum concentration in Case 3 and accordingly we find that  $n/n_Q = 7.73 \times 10^{-4}$ , well within the MB regime. Using

$$mv^2/2 = 3kT/2 \quad (20)$$

we calculate  $v = 7.90 \times 10^3$  cm sec $^{-1} \ll c$ .

This justifies the assumption that the helium gas, in this case, is nonrelativistic and classical. It is interesting that the ratio of helium atoms in the vapor to helium atoms in the liquid at 1 K is  $5.51 \times 10^{-5}$ . The ratio of the corresponding pressures is  $5.50 \times 10^{-4}$ . Both these values are several orders of magnitude larger than the corresponding ratios of electrons outside and inside a hot metal.

## V. ULTRARELATIVISTIC PARTICLES

*Case 5 (Electrons.  $S = \frac{1}{2}$ ):* Consider a Fermi gas of ultrarelativistic electrons. This problem is important in the theory of white dwarfs.<sup>11</sup> Relativistic effects occur when a gas is compressed, the average energy of the electrons rises, and the Fermi energy becomes comparable to the rest energy.

It is no longer correct to take

$$\mu_0 = (\hbar^2/2m)(3\pi^2 n)^{2/3} \quad (21)$$

to solve for  $n$ , because (21) is obtained by considering particle energy nonrelativistically as  $\epsilon = p^2/2m$ . Instead by using  $\epsilon = pc$ , the correct expression for  $\mu_0$  is obtained from the integral

$$N = C' \int_0^\infty \frac{\epsilon^2 d\epsilon}{e^{(\epsilon - \mu)/kT} + 1}, \quad (22)$$

where  $C'$  is given by (7). It follows from I, Eqs. (16) and (27), that

$$\mu_0 = (3n\pi^2)^{1/3} \hbar c. \quad (23)$$

Now take  $\mu_0 = m_0 c^2$  and solve for  $n$ ,

$$n = (m_0 c / \hbar)^3 / 3\pi^2 = 5.86 \times 10^{29} \text{ cm}^{-3}. \quad (24)$$

This is six orders of magnitude larger than the concentration of conduction electrons in a metal. This suggests, but does not establish, high degeneracy. The reason is that it is not correct to compare this  $n$  with the quantum concentration  $n_Q$  given by (10). We note that there  $n_Q \sim m^{3/2}$  and that for the ultrarelativistic case  $m_0 = 0$ . Equation (10) is derived for  $\epsilon = p^2/2m$  and like the Fermi energy must be rederived for  $\epsilon = pc$ . With the aid of Eq. (2) from I,  $p = 2$ ,  $C \rightarrow C'$ , it follows that

$$\mu = kT \ln(n/n_Q) (2S + 1)^{-1}. \quad (25)$$

This is the same form as Eq. (3) in I, but now

$$n_Q = (kT/\hbar c)^3/\pi^2. \quad (26)$$

The temperature in the interior of a white dwarf<sup>12</sup> is  $\sim 10^7$  K. Thus  $n/n_Q = (m_0 c^2/kT)^3/3 = 6.98 \times 10^7$ , which proves this gas to be highly degenerate.

In the ultrarelativistic case  $PV = U/3$  [see (9)], where  $U$  can be found from Eqs. (16) and (27) in I with  $C \rightarrow C'$  and  $p = 3$ . It follows directly that

$$P = (3\pi^2)^{1/3} (n^{4/3}) \hbar c / 4 = 1.19 \times 10^{17} \text{ atm}, \quad (27)$$

where  $n$  is given by (24).

*Case 6 [Photons:  $v = c$ ,  $m_0 = 0$ ,  $S = 1$  (the quantum degeneracy factor is 2 not 3)]:* Consider a gas of photons in a cavity. The fact that this gas is degenerate is very well known historically. The correct explanation of cavity radiation by Planck started quantum theory.<sup>13</sup> We establish degeneracy here by using the formalism of I,

$$N = C' (kT)^3 \int_0^\infty \frac{\epsilon^2 d\epsilon}{e^{(\epsilon - \mu)/kT} - 1}, \quad (28)$$

where  $C'$  is given by (7). The chemical potential  $\mu$  for a Bose (or Fermi) gas has a finite nonzero value if the number of particles  $N$  of the system is constant. But photons may be emitted or absorbed by the walls of the cavity. Thus  $N$  is not conserved. Therefore  $\mu = 0$ . Using Eqs. (6) and (8) in I,

$$n = 0.244 (kT/\hbar c)^3. \quad (29)$$

From (26),  $n/n_Q = 2.408$ , which establishes the degeneracy for any  $T$ . To obtain (29), we use  $\zeta(3) = \sum_{m=1}^\infty m^{-3} = 1.202$ .

In the ultrarelativistic regime,  $PV = U/3$ . The average energy  $U$  may be calculated from Eqs. (6) and (8) in I with  $C \rightarrow C'$  and  $p = 3$ . The pressure follows immediately,

$$P = 2(kT)^4 \zeta(4) / \pi^2 \hbar^3 c^3, \quad (30)$$

where the Riemann zeta function  $\zeta(4) = 1.082$ . The photon gas pressure may be written in an equivalent, but more familiar form,<sup>13</sup>

$$P = 4\sigma T^4 / 3c, \quad (31)$$

where  $\sigma = \pi^2 k^4 / 60 \hbar^3 c^2 \text{ W/m}^2 (\text{K})^4$  is the Stefan-Boltzmann constant. Numerically  $P = 3.99 \times 10^{-8}$  atm for  $T = 2000$  K and  $\sigma = 5.67 \times 10^{-8}$  in SI units.

*Case 7 (The pressure due to a laser beam):* In the previous example, the cavity is maintained at a temperature = 2000 K. The radiation emitted, through a small hole in the cavity, is characteristic of a thermal equilibrium distribution. Photons of all wavelengths leave and enter this hole.

In the case of a laser beam, photons of a single wavelength emerge from a "hole" at one end of a Fabry-Perot cavity. Here, we ask, what steady state power in the He-Ne laser at 632.8 nm is required to generate a beam pressure which is equal to the black-body radiation pressure of Case 6? We note that this physical situation does not constitute a statistical problem. But it does provide an interesting comparison.

Suppose the steady state power  $L$  of this laser is measured in watts. The beam intensity is the power divided by the cross sectional area of the beam,

$$I = L / \pi r^2. \quad (32)$$

TABLE I. The pressure of various physical systems.

Physical system	Spin	Statistics	Temperature (K)	Pressure (atm)
Free electrons in tungsten	$\frac{1}{2}$	Nonrelativistic FD	2000	$7.4 \times 10^5$
Space charge limited electron flow	$\frac{1}{2}$	Nonrelativistic MB	300	$2 \times 10^{-13}$
Liquid helium atoms	0	Nonrelativistic BE	1	$2.78 \times 10^{-1}$
Helium vapor atoms	0	Nonrelativistic MB	1	$1.53 \times 10^{-4}$
Electrons in a white dwarf star	$\frac{1}{2}$	Ultrarelativistic FD	$10^7$	$1.19 \times 10^{17}$
Cavity radiation	1	Ultrarelativistic BE	2000	$3.98 \times 10^{-8}$

The radiation pressure is

$$P = I/c = L/\pi r^2 c. \tag{33}$$

We take  $P = 4.03 \times 10^{-3} \text{ N/m}^2$  from the black-body radiation problem and solve for  $L$ , with  $r = 1 \text{ mm}$ .

Here  $L = 3.80 \text{ W}$ . This corresponds to  $N = L\lambda/hc = 1.21 \times 10^{19}$  red He-Ne laser photons/sec.

The beam pressure is independent of wavelength, but of course the associated number of photons/sec depends on the energy of monochromatic photons.

## VI. CONCLUSION

The results of this paper are summarized in Table I.

The methods developed in I make it quite easy to investigate the physical properties of diverse systems of interest. This paper focuses on the pressure of each system. This thermodynamic coordinate ranges from  $10^{17}$  atm in a white dwarf to  $10^{-13}$  atm in the evacuated region of a diode.

Each physical system in Table I is investigated to determine whether the particle speed is nonrelativistic or ultrarelativistic and whether the concentration is smaller or larger than the quantum concentration. This must be established, along with the particle's spin, in order to identify the appropriate statistical formulation, as given in Table I, column 3.

Ultrarelativistic particles, characterized by  $\epsilon = pc$ , lead to an extension of the formalism of I. There the nonrelativistic

expression  $\epsilon = p^2/2m$  is used. The essential change occurs in the volume element in momentum space. In the nonrelativistic case  $4\pi p^2 dp = 4\pi 2^{1/2} m^{3/2} \epsilon^{1/2} d\epsilon$ . In the ultrarelativistic case this becomes  $4\pi \epsilon^2 d\epsilon/c^3$ . All of the familiar quantities like chemical potential, quantum concentration, and internal energy are quite different in the ultrarelativistic regime.

<sup>1</sup>W. A. Barker, *J. Math. Phys.* **27**, 302 (1986).

<sup>2</sup>C. Kittel, *Elementary Statistical Physics* (Wiley, New York, 1958), p. 63.

<sup>3</sup>L. D. Landau and E. M. Lifshitz, *Statistical Physics* (Addison-Wesley, Reading, MA, 1958), p. 70.

<sup>4</sup>See, for example, C. Kittel and H. Kroemer, *Thermal Physics* (Freeman, San Francisco, 1980), p. 73.

<sup>5</sup>See, for example, *Applied Electronics*, edited by T. S. Gray (Wiley, New York, 1954), pp. 125-132.

<sup>6</sup>C. D. Child, *Phys. Rev.* **1**, 32, 498 (1911); I. Langmuir, *ibid.* **2**, 450 (1913).

<sup>7</sup>*Handbook of Physics and Chemistry* (CRC, Cleveland, 1977), 57th ed., p. B-25.

<sup>8</sup>Reference 4, p. 206.

<sup>9</sup>Reference 4, p. 215.

<sup>10</sup>J. R. Clement, J. K. Logan, and J. Gaffney, *Bull. Inst. Int. Froid Suppl. Annexe 3*, 601 (1955). See also *Handbuch der Physik, Low Temperature Physics II* (Springer, Berlin, 1956), p. 405.

<sup>11</sup>M. Harwit, *Astrophysical Concepts* (Wiley, New York, 1973), p. 163.

<sup>12</sup>See Ref. 7, p. 198.

<sup>13</sup>See Ref. 3, p. 175.

# Series representations for calculations in quantum statistics. III

William A. Barker

Department of Physics, Santa Clara University, Santa Clara, California 95053

(Received 12 August 1986; accepted for publication 30 January 1987)

In two previous papers [J. Math. Phys. **27**, 1 (1986) (I); **28**, 1385 (1987) (II)] simple mathematical procedures are developed for averaging the particle energy  $\epsilon$ , raised to any power  $p \geq 1$  over Maxwell-Boltzmann, Bose-Einstein, and Fermi-Dirac distributions. In I, the particle energy is treated nonrelativistically:  $\epsilon = p^2/2m$ . In II, the particle energy is treated ultrarelativistically:  $\epsilon = pc$ . Here, the particle energy is treated relativistically  $\epsilon^2 = p^2c^2 + m_0^2c^4$ . For each distribution function, expressions are obtained for the chemical potential, internal energy, and heat capacity for two cases:  $(kT/m_0c^2)^2 \ll 1$  and  $(m_0c^2/kT)^2 \ll 1$ .

## I. INTRODUCTION

In the literature of classical and quantum statistics, most calculations are made for nonrelativistic particles.<sup>1</sup> There are some which use the ultrarelativistic case,<sup>2</sup> but there are very few which treat the particle energy exactly.<sup>3</sup> From a mathematical point of view, this is understandable. To average the relativistic density of states expression over any of the distributions

$$f(\epsilon) = 1/(e^{(\epsilon - \mu)/kT} \pm 1, 0) \quad (1)$$

is indeed formidable. Furthermore, most of the cases of physical interest are nonrelativistic. However, the mathematical simplifications introduced in I make it quite easy to use two approximate forms: either  $(kT/m_0c^2)^2 \ll 1$  or  $(m_0c^2/kT)^2 \ll 1$ . The purpose of this paper is to use these two approximations to the relativistic density of states expression and extend I and II to derive the corresponding expressions for the chemical potential, energy, and heat capacity for Maxwell-Boltzmann (MB), Bose-Einstein (BE), and Fermi-Dirac (FD) statistics.

## II. DENSITY OF STATES

The density of states takes three different forms in calculations using MB, BE, and FD statistics.

$$\text{In the nonrelativistic approximation,}^4 \epsilon = p^2/2m \text{ and} \\ D(\epsilon)d\epsilon = (2S + 1)2^{1/2}Vm^{3/2}\epsilon^{1/2}d\epsilon/(2\pi^2\hbar^3). \quad (2)$$

In the ultrarelativistic limit,<sup>5</sup>  $\epsilon = pc$  and

$$D(\epsilon)d\epsilon = (2S + 1)V\epsilon^2d\epsilon/(2\pi^2\hbar^3c^3). \quad (3)$$

In the relativistic case, which includes the foregoing approximations,

$$\epsilon^2 = p^2c^2 + m_0^2c^4 \quad (4)$$

and

$$D(\epsilon)d\epsilon = (2S + 1)V\epsilon(\epsilon^2 - m_0^2c^4)^{1/2}d\epsilon/(2\pi^2\hbar^3c^3). \quad (5)$$

Each of these expressions is easily found from

$$D(\mathbf{p}, \mathbf{r})d\tau = (2S + 1)dp_x dp_y dp_z dx dy dz/(2\pi\hbar)^3. \quad (6)$$

First, transform to polar coordinates in momentum space,  $dp_x dp_y dp_z = 4\pi p^2 dp$ , and then use each of the three relationships connecting particle momentum and energy.

## III. CHEMICAL POTENTIALS; THE FERMI, BOSE, AND CLASSICAL ENERGY

The range of energy in the three types of integrals in I and II is from 0 to  $\infty$ . This makes the mathematics tractable. However, in the relativistic expression (4), the lower limit on  $\epsilon = m_0c^2$ , corresponding to  $p = 0$ . This suggests a change in variables to  $x = (\epsilon - m_0c^2)/kT$ . The general form of the integral is now

$$I = \int_0^\infty \frac{f(x)dx}{e^xe^{-(\mu - m_0c^2)/kT} \pm 1, 0}. \quad (7)$$

The quantity  $(\mu - m_0c^2)$  in the denominator of (7) is a crucial parameter in this theory. In the FD case, the Fermi energy  $E_F = \mu - m_0c^2$ . By analogy, in the MB and BE cases, we define the Bose energy  $E_B$  and the classical energy  $E_C$  to be  $\mu - m_0c^2$ . The chemical potential  $\mu$  is different for each of the distribution functions. It depends on particle concentration or number, rest energy, and temperature. The corresponding energy parameter reflects this dependence.

## IV. EXACT AND APPROXIMATE SOLUTIONS

At this point, we make two exact calculations using (5), with a limiting value for the FD distribution. This will be used to provide an internal consistency check on the approximate FD calculations and as a guide for interpreting results obtained from the other statistical distributions.

Define  $W \equiv \epsilon - m_0c^2$  and consider the FD distribution function

$$f(W) = 1/(e^{(W - E_F)/kT} + 1) \quad (8)$$

in the limit as  $T \rightarrow 0$ . For  $0 \leq W \leq E_F$ ,  $f(W) = 1$ ;  $W$  has an upper limit, namely,  $E_F(0, m_0)$ , the Fermi energy at  $T = 0$ .

There is a very simple argument which leads to an exact expression for  $E_F(0, m_0)$ . Start with the density of states expression (6), take  $S = \frac{1}{2}$ , write out the integral for the total number of particles  $N$ , using the Fermi function  $f(p) = 1$ ,

$$N = \frac{V}{\pi^2\hbar^3} \int_0^{p_0} p^2 dp, \quad (9)$$

where  $p_0$  is the limiting value of the momentum corresponding to  $E_F(0, m_0)$ . Use the relativistic formula (4) and  $W = \epsilon - m_0c^2$  to rewrite  $N$  as an integral over  $W$ ,

$$N = \frac{V}{\pi^2 \hbar^3 c^3} \int_0^{E_F(0, m_0)} (W + m_0 c^2) \times (W^2 + m_0 c^2)(W^2 + 2Wm_0 c^2)^{1/2} dW. \quad (10)$$

This yields

$$n = N/V = (E_F^2 + 2E_F m_0 c^2)^{3/2} / 3\pi^2 \hbar^3 c^3. \quad (11)$$

In the limit as  $m_0 \rightarrow 0$ ,

$$E_F(0, 0) = (3\pi^2 n)^{1/3} \hbar c = p_0 c, \quad (12)$$

a value which may be used in other expressions for  $E_F(T, m_0)$ . Using (11), the exact relation for

$$E_F(0, m_0) = (m_0^2 c^4 + E_F^2(0, 0))^{1/2} - m_0 c^2. \quad (13)$$

If  $E_F/m_0 c^2 \ll 1$ , it follows from (13) that

$$E_F(0, \infty) = E_F^2(0, 0) / 2m_0 c^2 = (3\pi^2 n)^{2/3} \hbar^2 / 2m_0. \quad (14)$$

If  $m_0 c^2 / E_F \ll 1$ , it follows from (13) that

$$E_F(0, m_0) = E_F(0, 0)(1 - m_0 c^2 / E_F(0, 0)). \quad (15)$$

The physical significance of these four expressions for the Fermi energy will be discussed later in the paper.

To find the internal energy  $U = U(0, m_0)$ , multiply the integrand in (9) by  $\epsilon$  from (4)

$$U = \frac{Vc}{\pi^2 \hbar^3} \int_0^{p_0} p^2 (p^2 + m_0^2 c^2)^{1/2} dp. \quad (16)$$

The result, expressed in terms of  $E_F = E_F(0, 0)$ , is

$$U = V/8\pi^2 \hbar^3 c^3 (E_F (2E_F^2 + m_0^2 c^4) (E_F^2 + m_0^2 c^4)^{1/2} - m_0^4 c^8 \sinh^{-1}(E_F / m_0 c^2)). \quad (17)$$

Equation (17) is equivalent to the formula for the energy given by Landau and Lifschitz.<sup>6</sup> Its form is more convenient than an alternate expression for  $U$  written in terms of  $E_F(0, m_0)$ , which may be obtained by multiplying the integrand of (10) by  $\epsilon = W + m_0 c^2$ . The foregoing formula (17) will be used later in the paper to compare with various approximate forms for  $U(T, m_0)$ .

## V. THE MAXWELL-BOLTZMANN DISTRIBUTION

*Case A* [ $(kT/m_0 c^2)^2 \ll 1$ ]: In this approximation, the density of states (5) is

$$D(x) dx \cong [2^{1/2} (2S + 1) V (m_0 kT)^{3/2} / 2\pi^2 \hbar^3] \times (x^{1/2} + 5ax^{3/2}/4 + 7a^2 x^{5/2}/32 + \dots), \quad (18)$$

where  $a = (kT/m_0 c^2)$ . The total number of particles

$$N = e^{E_c/kT} \int_0^\infty D(x) e^{-x} dx = A (kT)^{3/2} e^{E_c/kT} (1 + 1.88a + 0.82a^2), \quad (19)$$

where  $A = (2S + 1) V m_0^{3/2} / (2\pi \hbar^2)^{3/2}$ . The integral (19) is the sum of three gamma functions of the form  $\Gamma(p + 1)$  where  $p$  is the exponent of  $x$ . Solving for the chemical energy yields

$$E_c = kT \ln(n/n_Q (2S + 1)) (1 - 1.88a + 2.70a^2), \quad (20)$$

where

$$n_Q = (m_0 kT / 2\pi \hbar^2)^{3/2}. \quad (21)$$

Compare with Eqs. (3) and (4) of I, which apply to the nonrelativistic approximation. Here,  $E_c = \mu - m_0 c^2$  replaces the chemical potential. It has the standard form, with two correction terms. The quantum concentration is identical with the nonrelativistic value in I.

The energy and the heat capacity follow directly. Multiply  $dN$  by  $\epsilon = m_0 c^2 + kTx$  and integrate,

$$U = Nm_0 c^2 (1 + 1.25a - 1.25a^2), \quad (22)$$

$$C_V = \frac{3}{2} Nk (1 + 2.50a - 3.75a^2). \quad (23)$$

Certainly the rest energy term in (22) was to be expected. The second term in (22) and the first term in (23) agree with Eqs. (5) and (6) in I. The remaining terms are thermal corrections to the nonrelativistic approximation.

*Case B* [ $(m_0 c^2 / kT)^2 \ll 1$ ]: In this approximation, the density of states is given by

$$D(x) dx = [(2S + 1) V (kT)^3 / 2\pi^2 \hbar^3 c^3] \times (x^2 + 2bx + 7b^2/8) dx, \quad (24)$$

where  $b = m_0 c^2 / kT$ .

The total number of particles

$$N = e^{E_c/kT} \int_0^\infty D(x) e^{-x} dx = B (kT)^3 e^{E_c/kT} (1 + b + 0.438b^2), \quad (25)$$

where  $B = (2S + 1) V / 2\pi^2 \hbar^3 c^3$ . Solving for the classical energy

$$E_c = kT \ln(n/n_Q (2S + 1)) (1 - b + 0.562b^2), \quad (26)$$

where

$$n_Q = (kT / \hbar c)^3 \pi^{-2} \quad (27)$$

is the ultrarelativistic expression for the quantum concentration. [See Eq. (26) of II.]  $E_c$  has the expected form plus both linear and quadratic rest energy correction terms.

To obtain the energy  $U$ , multiply the integrand in (25) by  $kTx + m_0 c^2$  and integrate,

$$U = e^{E_c/kT} \int_0^\infty (kTx + m_0 c^2) D(x) e^{-x} dx. \quad (28)$$

With the aid of (25), the energy and heat capacity are written

$$U = 3NkT (1 + 1.92b^2), \quad (29)$$

$$C_V = \frac{dU}{dT} = 3Nk (1 - 1.92b^2). \quad (30)$$

The ultrarelativistic results<sup>7</sup> follow from (29) and (30) in the limit as  $m_0 \rightarrow 0$ . Here  $U$  and  $C_V$  both have a quadratic rest energy correction term. It is interesting to note that (22), in the nonrelativistic limit, and (29) in the ultrarelativistic limit, lead to precisely the same familiar equation of state  $PV = NkT$ . This follows because  $PV = 2U/3$  and  $U/3$  in these two limits, respectively. However, neither of these relationships are consequences of the relativistic expression connecting particle energy and momentum.

## VI. THE BOSE-EINSTEIN DISTRIBUTION

The BE integrals have the form

$$I_{\text{BE}} \propto \int_0^\infty \frac{f(x)dx}{e^x e^{-E_B/kT} - 1}, \quad (31)$$

where  $E_B = (\mu - m_0c^2)$  is the Bose energy. Let  $z = E_B/kT$ . The quantity  $e^{-z}$  can be shown to be  $\approx 1$  below the Einstein condensation temperature. The argument parallels the development in I [See Eqs. (10)–(12)] inclusive. If  $N_0$  is the occupation number of the ground state when  $\epsilon = m_0c^2$  and  $x = 0$ , then

$$N_0 = (e^{-z} - 1)^{-1}, \quad (32)$$

$$e^{-z} = 1 + 1/N_0, \quad (33)$$

$$E_B \cong -kT/N_0. \quad (34)$$

*Case A* [ $(kT/m_0c^2)^2 \ll 1$ ]: The density of states is given by (18). The total number of bosons in excited energy states

$$N_e = \int_0^\infty \frac{D(x)dx}{e^x - 1}. \quad (35)$$

The integral involves the sum of three gamma-Riemann zeta function products. As in I, this integrates immediately to

$$N_e = 2.612(2S + 1)n_Q V(1 + 0.962a + 0.355a^2). \quad (36)$$

This agrees with the nonrelativistic expression (10) in I, with two thermal correction terms. To find  $U$ , multiply the integrand in (35) by  $(m_0c^2 + kTx)$ ,

$$U = m_0c^2 \int_0^\infty \frac{(1 + ax)D(x)dx}{e^x - 1}. \quad (37)$$

This integrates to

$$U = N_e m_0c^2 + 1.341(2S + 1)(3NkT/2) \times (1 + 2.62a + 1.51a^2)n_Q/n. \quad (38)$$

There is an important distinction between  $N_e$  and  $N$  in Eq. (38).  $N_e$  is the number of excited bosons whose  $p > 0$ . Here  $N$  is the total number of bosons, excited and condensed. It enters (38) via  $V = N/n$ . This distinction can be kept in mind by realizing that the Einstein condensation is a condensation in momentum space, but not in coordinate space.

The heat capacity  $C_V$  follows by differentiating  $U$  with respect to  $T$ ,

$$C_V = 3.35(2S + 1)(3NkT/2) \times (1 + 9.19a + 6.77a^2)n_Q/n. \quad (39)$$

Equations (38) and (39) agree with the nonrelativistic Eqs. (14) and (15) in I in the appropriate limit. Both  $U$  and  $C_V$  have linear and quadratic correction terms.

*Case B* [ $(m_0c^2/kT)^2 \ll 1$ ]: In this approximation, the density of states is given by (24). The total number of bosons in excited states

$$N_e = \int_0^\infty \frac{D(x)dx}{e^x - 1}. \quad (40)$$

This integrates to

$$N_e = 1.202(2S + 1)N(1 + 1.31b)n_Q/n, \quad (41)$$

where the quantum concentration is the ultrarelativistic value given in (27). The third term in the density of states

expression (24) is not carried in this calculation, as it leads to a divergent term in (40). The energy

$$U = \int_0^\infty \frac{(kTx + m_0c^2)D(x)dx}{e^x - 1} \quad (42)$$

integrates to

$$U = N_e mc^2 + 1.082(2S + 1)(3NkT)(1 + 0.741b)n_Q/n. \quad (43)$$

The corresponding heat capacity is

$$C_V = 13(2S + 1)Nk(1 + .556b)n_Q/n. \quad (44)$$

The correction terms in (43) and (44) are linear terms in  $b = m_0c^2/kT$ .

The Einstein condensation temperature, in the ultrarelativistic case, is an additional quantity of interest in BE statistics. The Einstein temperature  $T_E$  is that temperature for which the number of bosons  $N_e$  in excited states equals the total number of bosons  $N$ . Using (41) with  $m_0 = 0$

$$T_E = 2.017(n/2S + 1)^{1/3}/\hbar c/k. \quad (45)$$

The fraction of particles in excited states is

$$N_e/N = (T/T_E)^3, \quad (46)$$

and the number of particles in the ground state is

$$N_0 = N - N_e = N(1 - (T/T_E)^3). \quad (47)$$

Consider a photon gas at  $T = 2000$  K with a concentration of  $n = 10n_Q = 6.74 \times 10^{11}/\text{cm}^3$ . The Einstein condensation temperature from (45) is  $3.22 \times 10^3$  K. The fraction in excited states is about 0.24. Compare this with the Einstein condensation for liquid helium. The experimental value for  $T_E = 2.174$  K. Using  $N_e/N = (T/T_E)^{3/2}$ , the fraction in excited states at 1 K is about 0.31. Actually, there is no photon condensation because the number of photons is not a constant.<sup>8</sup>

## VII. THE FERMI-DIRAC DISTRIBUTION

Fermi-Dirac integrals, treated relativistically, are represented by a superposition of integrals of the form

$$I_{\text{FD}} \propto \int_0^\infty \frac{x^p dx}{e^x e^{-E_F/kT} + 1}, \quad (48)$$

where  $E_F = \mu - m_0c^2$  and  $p = 0, \frac{1}{2}, 1, \frac{3}{2}, 2, \dots$ . With  $z = E_F/kT$ , (48) has the same structure as Eq. (16) in I. The Blankenbecler<sup>9</sup> method applies and we can write

$$I_{\text{FD}} \propto \left(1 + \frac{\pi^2}{6} \frac{\partial^2}{\partial z^2} + \dots\right) \frac{z^{p+1}}{p+1}. \quad (49)$$

*Case A* [ $(kT/m_0c^2)^2 \ll 1$ ]: The approximate form for the density of states in this case is given by (18). The total number of particles

$$N = \int_0^\infty \frac{D(x)dx}{e^x e^{-z} + 1}. \quad (50)$$

Using (49) and solving for  $N$  with  $S = \frac{1}{2}$ ,

$$N = \frac{4}{3} A E_F^{3/2} (1 + (3E_F/4m_0c^2) + (\pi^2/8)(kT/E_F)^2), \quad (51)$$

where  $A = Vm_0^{3/2}/(2\pi\hbar^2)^{3/2}$ . In the limit as  $T \rightarrow 0$ ,  $m_0 \rightarrow \infty$ , we find that

$$E_F(0, \infty) = \frac{E_F^2(0,0)}{2m_0c^2} = (3\pi^2n)^{1/3}/2m_0, \quad (52)$$

in agreement with (14). Retaining  $T = 0$ , we can solve for

$$E_F(0, m_0) = \frac{E_F^2(0,0)}{2m_0c^2} \left( 1 - \frac{E_F^2(0,0)}{4m_0^2c^4} \right), \quad (53)$$

in agreement with an expansion of the exact expression (13). Finally

$$E_F(T, m_0) = E_F(0, m_0) \left( 1 - (\pi^2/12)(kT/E_F)^2 \right), \quad (54)$$

which is the standard result,<sup>10</sup> subject to the new meaning for  $E_F(0, m_0)$ , as given by Eq. (53).

To solve for  $U$  multiply  $dN$  in (50) by  $\epsilon = m_0c^2 + kTx$ ,

$$U = \int_0^\infty \frac{(m_0c^2 + kTx)D(x)dx}{e^xe^{-z} + 1}. \quad (55)$$

The result is

$$U = Nm_0c^2 + (\frac{3}{2})NE_F(0, m_0) \left( 1 + (5\pi^2/12)(kT/E_F)^2 \right). \quad (56)$$

The algebra is tedious, but the result is of the expected form, subject to a new meaning for  $E_F(0, m_0)$ .

It is a simple matter to demonstrate that  $\lim_{T \rightarrow 0} U$  is an approximation for  $U$  given in (17). To show this, note that  $V = N/n$ ,  $n\hbar^3c^3 = E_F^3(0,0)/3\pi^2$ , and  $E_F(0, m_0) \simeq E_F^2(0,0)/2m_0c^2$ . Expand (17) to the second power in  $E_F(0,0)$ . The values for  $U(0, m_0)$  from (56) and (17) agree,

$$U(0, m_0) = Nm_0c^2 \left( 1 + (3E_F^2(0,0)/10m_0^2c^4) \right). \quad (57)$$

The heat capacity

$$C_V = \pi^2 Nk^2 k^2 T / 2E_F(0, m_0) \quad (58)$$

is of standard form, but there is a correction for  $E_F(0, m_0)$ , as shown in (53).

*Case B* [ $(m_0c^2/kT)^2 \ll 1$ ]: The density of states is given by (24). The total number of particles

$$N = \int_0^\infty \frac{D(x)dx}{e^xe^{-z} + 1} = \frac{VE_F^3}{3\pi^2\hbar^3c^3} \left( 1 + (3m_0c^2/E_F) \right) \times \left( 1 + m_0c^2/2E_F + \pi^2(kT/E_F)^2 \right). \quad (59)$$

From (59)

$$E_F(0, m_0) = E_F(0,0) \left( 1 - m_0c^2/E_F + m_0^2c^4/2E_F^2 \right) \quad (60)$$

in agreement with an expansion of the exact expression (13). Again from (59) and (60),

$$E_F(T, m_0) = E_F(0,0) \left( 1 - m_0c^2/E_F + m_0^2c^4/2E_F^2 - \pi^2(kT/E_F)^2 \right). \quad (61)$$

To find  $U$ , multiply the integrand in (59) by  $\epsilon = kTx + m_0c^2$  and integrate. After some tedious algebra, the result is

$$U(T, m_0) = \frac{3}{2}NE_F(0,0) \times \left( 1 + \frac{4}{3}(m_0^2c^4/E_F^2) + \pi^2(kT)^2/E_F^2 \right). \quad (62)$$

This expression, at  $T = 0$ , agrees with an expansion of the exact formula (17). It is easy to see that the second term in (17) does not contribute to this order. Write

$$0.5(m_0c^2/E_F)^4 \sinh^{-1}(E_F/m_0c^2) = -0.5(m_0c^2/E_F)^4 \times \log \left[ m_0c^2 / ((E_F^2 + m_0^2c^4)^{1/2} + E_F) + 1 \right]. \quad (63)$$

An expansion generates terms  $O(m_0c^2/E_F)^5$  and smaller. The heat capacity

$$C_V = (3\pi^2/4)Nk(kT/E_F(0,0)). \quad (64)$$

This has the same structure as the Case A result [Eq. (58)]. However, the heat capacity in Case A is larger than in Case B because  $E_F(0, m_0)$  is smaller than  $E_F(0,0)$ . See Eq. (52).

## VIII. DISCUSSION

### A. Three tests distinguish among the 12 statistical formulations

For the individual who is about to make some calculations in statistical mechanics, there may appear to be a bewildering array of choices. For each of the three distributions, MB, BE, and FD, there are four possible expressions for the density of states: nonrelativistic, relativistic A, relativistic B, and ultrarelativistic.

The choice is clear once three tests are made.

In the first test, a comparison of  $kT$  with  $m_0c^2$  distinguishes amongst the density of states formulas. If  $kT/m_0c^2 \leq 10^{-2}$ , use the nonrelativistic formula [Eq. (2)]. If  $kT/m_0c^2 > 10^{-2}$  and  $(kT/m_0c^2)^2 \leq 10^{-2}$ , use the relativistic A formula [Eq. (18)]. If  $m_0c^2/kT > 10^{-2}$  and  $(m_0c^2/kT)^2 \leq 10^{-2}$ , use the relativistic B formula [Eq. (24)]. If  $m_0c^2/kT \leq 10^{-2}$ , use the ultrarelativistic formula [Eq. (3)].

In the second test, the particle density is compared with the quantum concentration  $n_Q$ . This determines whether the problem is classical or degenerate.

If the first test establishes the density of states to be either nonrelativistic or relativistic A, then if  $n \ll n_Q = (mkT/2\pi\hbar^2)^{3/2}$ , the problem is nondegenerate. Use MB statistics. If  $n \gtrsim n_Q$ , the problem is degenerate. Use BE or FD statistics. On the other hand, if the first test establishes the density of states to be either relativistic B or ultrarelativistic, then if  $n \ll n_Q = (kT/\hbar c)^3/\pi^2$ , the problem is nondegenerate and calls for MB statistics. If  $n \gtrsim n_Q$ , the problem is degenerate, requiring either BE or FD statistics.

In the third test, the particle spin is used to distinguish bosons from fermions. Let  $S_0 = n\hbar/2$ , where  $n = 0, 1, 2, \dots$ . Then for  $n$  odd (even), the particle is a fermion (boson).

If tests 1 and 2 establish the problem as degenerate and either nonrelativistic, relativistic A, relativistic B, or ultrarelativistic, then  $n$  even (odd) requires BE (FD) statistics.

### B. The classical, Bose, and Fermi energies

In MB statistics, the leading term for the classical energy  $E_C$  has the same form for both the relativistic A and B regimes [Eqs. (20) and (26)]. However, the quantum concentration  $n_Q$ , as described in test 2, has two quite different meanings [Eqs. (21) and (27)]. The first (second) is appropriate for relativistic A (B).

Suppose that  $S = \frac{1}{2}$  and that  $n/n_Q \leq 10^{-2}$ , then  $E_C/kT = \ln(n/n_Q) \leq -4.6$ . The classical energy  $E_C$  will always be negative, as required by  $n \ll n_Q$ . In this situation, MB statistics are used.

The Bose energy  $E_B$  has a small negative value  $\sim -kT/N$ , if the particle number  $N$  is conserved. However,  $E_B = 0$  if the particle number is not conserved. The former value is valid if the actual temperature is less than the Einstein condensation temperature.

The Fermi energy takes several forms. In the case of a problem which is nonrelativistic or relativistic A, the leading term is  $E_F(0, m_0) = (3\pi^2 n)^{2/3} \hbar^2 / 2m$ . However, if the problem is relativistic B or ultrarelativistic, the leading term is  $E_F(0, 0) = (3\pi^2 n)^{1/3} \hbar c$ . These two expressions differ substantially in structure and magnitude. Consider, for example, conduction electrons of concentration  $n = 10^{22}/\text{cm}^3$ . Then  $E_F(0, m_0) = 1.70$  eV and  $E_F(0, 0) = 1.32 \times 10^3$  eV. All the other contributions to  $E_F$  are due to rest energy and thermal effects and are less than these leading terms.

### C. The third law of thermodynamics

The third law of thermodynamics requires that  $\lim_{T \rightarrow 0} C_V = 0$ . This law is violated by the MB heat capacity expressions (23) and (30), as is well known. However, the BE and FD expressions for  $C_V$ , as given in Eqs. (39), (44), (58), and (64), are in agreement with the third law. In verifying this for BE statistics, note that  $n_Q \sim T^{3/2}$  or  $T^3$ . This temperature dependence of the quantum concentration guarantees that the correction terms, as well as the leading terms, go to zero as  $T \rightarrow 0$ .

### IX. CONCLUSION

The methods developed in I are used, in this paper, to give a systematic discussion of MB, BE, and FD distributions when the density of states expression is based on the

exact relativistic energy-momentum relationship. Two approximations are used. When  $(kT/m_0c^2)^2 \ll 1$ , the results are the same as for the nonrelativistic approximation, but with correction terms in  $(kT/m_0c^2)$ . When  $(m_0c^2/kT)^2 \ll 1$ , the results agree with the ultrarelativistic approximation but with correction terms in  $m_0c^2/kT$ . The principal focus of the paper is on the quantity  $\mu - m_0c^2$ , where  $\mu$  is the chemical potential. This is the Fermi energy in FD statistics, and it is given a comparable definition for BE and MB statistics. As in I, calculations are made for the chemical potential, internal energy, and heat capacity for each of the three distributions. A feature of the paper are two exact solutions. The relativistic density of states formula is used, without approximation, with the FD distribution function at  $T = 0$ ,  $\epsilon \ll E_F$  to obtain exact expressions for the Fermi energy and the internal energy. When these formulas are expanded in a series, they agree to the same order with the results obtained using the series representation in I.

<sup>1</sup>C. Kittel, *Elementary Statistical Physics* (Wiley, New York, 1958), pp. 86–96; see also C. Kittel and H. Kroemer, *Thermal Physics* (Freeman, San Francisco, 1980), pp. 183–198.

<sup>2</sup>L. D. Landau and E.M. Lifschitz, *Statistical Physics* (Addison-Wesley, Reading, MA, 1968), pp. 165–167.

<sup>3</sup>See Ref. 2, pp. 167 and 168.

<sup>4</sup>William A. Barker, *J. Math. Phys.* **27**, 1 (1986).

<sup>5</sup>William A. Barker, *J. Math. Phys.* **28**, 1385 (1987).

<sup>6</sup>See Ref. 2, p. 168.

<sup>7</sup>See C. Kittell, Ref. 1, p. 60.

<sup>8</sup>See Ref. 1, Kittel and Kroemer, p. 202.

<sup>9</sup>R. Blankenbecler, *Am. J. Phys.* **25**, 279 (1957).

<sup>10</sup>See Ref. 1, Kittel, p. 94. [Eq. (20.24)].

# Stable thermodynamic states

J. S. Cohen

*Philips Research Laboratories, Eindhoven, The Netherlands*

M. Winnink

*Institute for Theoretical Physics, University of Groningen, The Netherlands*

(Received 2 December 1986; accepted for publication 11 February 1987)

A general class of perturbations of the dynamics for thermodynamic quantum systems is discussed. Without making use of weak asymptotic Abelianess, stability of a state for these perturbations is shown to lead to the  $\phi$ -KMS condition and to the KMS condition in particular cases. Conversely,  $\phi$ -KMS states satisfy the stability property introduced here.

## I. INTRODUCTION

The derivation of equilibrium properties for states of thermodynamic systems has already been of interest for some time. It is well known that so-called KMS states may be obtained from stability for perturbation of the dynamics. This has been discussed initially by Haag, Kastler, and Trych-Pohlmeyer, by Kastler and Bratteli, and by Hoekman; for a review cf. Ref. 1, Chap. 5.4.2. However, owing to the assumed rapid decay of the time correlation functions these KMS states can only describe pure thermodynamic phases. A further restriction of the method is that it can be applied only to states of dynamical systems that are weakly asymptotically Abelian; viz.  $\int dt \omega([A, \alpha_t B]) = 0$ . Here  $A, B$  denote elements of the  $C^*$  algebra  $\mathfrak{A}$ ;  $\omega \in E_{\mathfrak{A}}$  is a state over  $\mathfrak{A}$ , and  $\alpha_t \in \text{aut } \mathfrak{A}$  describes the time evolution.

In this paper we shall discuss a stability property which leads to states that satisfy the  $\phi$ -KMS condition introduced recently.<sup>2</sup> The main advantage of the stability criterion put forward here is that neither are assumptions made on the decay of the correlation functions nor is the dynamics assumed to act weakly asymptotically Abelian.

Depending on the details of the perturbation for which stability is imposed, the  $\phi$ -KMS states in some cases are KMS states. For an infinite quantum lattice system the  $\phi$ -KMS condition and KMS condition are equivalent.<sup>2</sup> Consequently, in this instance either our stability condition leads to a KMS state or the system does not admit states that are stable for the particular perturbation. For any finite system or for continuous quantum systems, however, the only states that fulfill the stability criterium are  $\phi$ -KMS states.

For a finite system the presently proposed stability property is stronger than the condition imposed by Lebowitz *et al.*<sup>3</sup> For thermodynamic systems our conditions are weaker than those introduced by Kastler<sup>4</sup> (cf. Ref. 5).

## II. A GENERALIZED PERTURBED DYNAMICS

In the Heisenberg picture the equation of motion for the unperturbed evolution reads

$$\frac{d}{dt} \alpha_t(A) = i \alpha_t(\delta(A)), \quad A \in D(\delta) \subset \mathfrak{A}, \quad (2.1)$$

where the derivation  $\delta$  is the infinitesimal generator of the group of  $*$  automorphisms  $\{\alpha_t\}$ . A perturbed dynamics can be considered as the solution of the differential equation

$$\frac{d}{dt} \tilde{\alpha}_t^h(A) = i \tilde{\alpha}_t^h(\delta(A)) + i \tilde{\alpha}_t^h([h_t, A]), \quad (2.2)$$

cf., e.g., Ref. 6, with  $h_t = h^* \in \mathfrak{A}$ . One may choose the particular form  $h_t = f(t)h$  with  $h = h^* \in \mathfrak{A}$  and  $f: \mathcal{R} \rightarrow \mathcal{R}$ , so that the perturbation becomes localized in time if, e.g.,  $\text{supp } f$  is compact or  $f \in L_1(\mathcal{R})$ . The family  $\{h_t\}$  can be interpreted as the action of some external agent on the system. Owing to the time dependence of  $h_t$ , the mappings  $\{\tilde{\alpha}_t^h\}$  do not form a group.

A further generalization of the perturbed dynamics is obtained from the following equation of motion:

$$\frac{d}{dt} \tilde{\alpha}_t^h(A) = i \tilde{\alpha}_t^h(\delta(A)) + i f_1(t) \tilde{\alpha}_t^h(hA) - i f_2(t) \tilde{\alpha}_t^h(Ah). \quad (2.3)$$

The solution to this equation is the family of mappings  $\tilde{\alpha}_t^h: \mathfrak{A} \rightarrow \mathfrak{A}$  given by

$$\tilde{\alpha}_t^h(A) = \tilde{\gamma}_t^h(\alpha_t(A)), \quad (2.4a)$$

$$\tilde{\gamma}_t^h(A) = \tilde{u}_t^h(t) A \tilde{u}_t^h(t)^*, \quad (2.4b)$$

$$\tilde{u}_t^h(t) = \sum_{n=0}^{\infty} \left\{ i^n \int_0^t ds_1 \cdots \int_0^{s_{n-1}} ds_n \prod_{k=1}^n [f_j(s_k) \alpha_{s_k}(h)] \right\}, \quad (2.4c)$$

with  $f_j \in L_1(\mathcal{R})$  and  $h = h^* \in \mathfrak{A}_0$ . The unperturbed dynamics  $\alpha_t$  is assumed to be strongly continuous on a  $\sigma(\mathfrak{A}, N)$ -dense subalgebra  $\mathfrak{A}_0 \subset \mathfrak{A}$ . Here  $N$  denotes the set of locally normal states on  $\mathfrak{A}$  (cf. Ref. 7). In general, the mappings  $\tilde{\alpha}_t^h$  and  $\tilde{\gamma}_t^h$  will not be positivity preserving. The operators  $\tilde{u}_t^h(t)$  are easily seen to be unitary. The integrals in (2.4c) exist as Bochner integrals. We now give some useful properties of the generalized perturbed dynamics in the following.

*Proposition 2.1:* For  $A \in \mathfrak{A}$  and  $h = h^* \in \mathfrak{A}_0$ ,

$$\tilde{u}_t^h(t) = \mathbb{1} + i \int_0^t ds f_j(s) \alpha_s(h) \tilde{u}_t^h(s), \quad (2.5a)$$

$$\frac{d}{dt} \tilde{u}_t^h(t) = i f_j(t) \alpha_t(h) \tilde{u}_t^h(t); \quad (2.5b)$$

$$\tilde{\gamma}_t^h(A) = A + i \int_0^t ds [f_1(s) \tilde{\gamma}_s^h(\alpha_s(h)A) - f_2(s) \tilde{\gamma}_s^h(A \alpha_s(h))], \quad (2.5c)$$



$$\begin{aligned} \tilde{\alpha}_t^h(A) = & \alpha_t(A) + i \int_0^t ds [f_1(s)\alpha_s(h)\alpha_t(A) \\ & - f_2(s)\alpha_t(A)\alpha_s(h)] + \dots \end{aligned} \quad (2.5d)$$

The omitted terms in (2.5d) are  $O(h^2)$ .

*Proof:* By iteration of (2.5a) we obtain (2.4c). The equivalence of (2.5c) and (2.5b) then follows from the initial condition  $\tilde{u}_t^h(0) = 1$ . Finally, (2.5c) and (2.5d) are obtained with the use of (2.4b) and (2.4c) along similar lines as in the discussion of the cocycle property (cf., e.g., Ref. 1).

We now turn to the introduction of the notion of stability for perturbations from the unperturbed dynamics  $\alpha_t$ , as described by (2.5d). Succinctly, one assumes that close to the original state  $\omega \in \mathbb{U}^*$  there exists a bounded linear functional  $\omega^h \in \mathbb{U}^*$  that is *almost* invariant for the perturbed evolution  $\tilde{\alpha}_t^h$ . At this point we shall impose some restrictions on the functions  $f_j$ .

**Definition 2.2:** For a pair of functions  $f_1$  and  $f_2$  such that

- (1)  $f_j \in L_1(\mathcal{R}) \cap C^1(\mathcal{R})$ ;
- (2)  $\hat{f}_j \in C^\infty(\mathcal{R})$  and invertible on  $\text{sp } \alpha$ , i.e.,  
 $\hat{f}_j(\lambda) \neq 0 \quad \forall \lambda \in \text{sp } \alpha = \{\lambda \in \mathcal{R} \mid \hat{g}(\lambda) \neq 0\}$

$$\forall g: \int dt g(t)\alpha_t(A) = 0$$

$$\forall A \in \mathbb{U}_0 \};$$

- (3)  $\hat{f}_1(\lambda) = \hat{f}_2(\lambda) \quad \text{iff } \lambda = 0$ ;

[ $\hat{g}(\lambda) = \int dt e^{-i\lambda t} g(t)$  is the Fourier transform] we say that a state  $\omega$  is  $(f_1, f_2)$ -stable if there exists a bounded linear functional  $\omega^{\mu h} \in \mathbb{U}^*$  such that for  $\mu$  in a neighborhood of the origin

$$\lim_{t \rightarrow \pm \infty} \omega^{\mu h}(\tilde{\gamma}_t^{\mu h} A) = \omega_{\pm}(A); \quad (2.6)$$

$$\omega_+^{\mu h}(A) - \omega_-^{\mu h}(A) = o(\mu); \quad (2.7)$$

and

$$\lim_{\mu \rightarrow 0} \omega^{\mu h}(\alpha_t A) = \omega(\alpha_t A), \quad \text{uniformly in } t \quad (2.8)$$

for all  $A \in \mathbb{U}$ . With the use of (2.5c) a simple estimate shows that  $\omega_+^{\mu h}(A) - \omega_-^{\mu h}(A) = O(\mu)$  so that (2.7) does not seem to be a very severe assumption.

The conditions (2.6) and (2.7) are in fact the same as the ones introduced by Kastler<sup>4</sup> and Hoekman<sup>5</sup> because there  $\omega^{\mu h}$  is a perturbed state which is invariant for the perturbed dynamics. In Ref. 5 a perturbed dynamics  $\alpha_t^h$  is considered that is an Abelian group of transformations. As a consequence, the perturbed state could be explicitly constructed, viz.  $\omega^{\mu h}(A) = \mathfrak{M}_t \omega(\alpha_t^{\mu h} A)$ , where  $\mathfrak{M}$  is an invariant mean over the additive group of the real numbers and  $t$  is a dummy variable.<sup>8</sup> If, in addition, one has that  $(\mathbb{U}_0, \alpha_t)$  is  $L_1$ -asymptotically Abelian, then the convergence (2.8) can be derived.

We shall now proceed with the demonstration that without loss of generality the perturbed state  $\omega^{\mu h}$  may be assumed to be *approximately invariant* for the perturbed dynamics.

**Lemma 2.3:** Let  $\mathfrak{M}$  be an invariant mean over the addi-

tive group of the real numbers.<sup>8</sup> Then the time-averaged perturbed state  $\mathfrak{M}_t \omega^{\mu h}(\alpha_t A) \equiv \mathfrak{M} \omega^{\mu h}(A)$ , where  $t$  is to be considered as a dummy variable, satisfies

$$\lim_{\mu \rightarrow 0} \mathfrak{M} \omega^{\mu h}(\alpha_t A) = \mathfrak{M} \omega(\alpha_t A), \quad \text{uniformly in } t \quad (2.9a)$$

and

$$\lim_{\mu \rightarrow 0} [\mathfrak{M} \omega^{\mu h}(\tilde{\alpha}_t^h A) - \mathfrak{M} \omega^{\mu h}(A)] = 0, \quad \text{uniformly in } t. \quad (2.9b)$$

*Proof:* Because  $\mathfrak{M}$  is an invariant mean we have<sup>8</sup>

$$|\mathfrak{M}[\omega^{\mu h}(\alpha_t A) - \omega(\alpha_t A)]|$$

$$\leq \sup_{t \in \mathcal{R}} |\omega^{\mu h}(\alpha_t A) - \omega(\alpha_t A)| < \epsilon,$$

for  $\mu < \mu_0(\epsilon)$ . This establishes the continuity property of the time averaging. Similarly, with the use of (2.5c) we obtain

$$\begin{aligned} & |\mathfrak{M}[\omega^{\mu h}(\tilde{\alpha}_t^{\mu h} A) - \omega^{\mu h}(A)]| \\ &= |\mathfrak{M}[\omega^{\mu h}(\tilde{\alpha}_t^{\mu h} A) - \omega^{\mu h}(\alpha_t A)]| \\ &\leq \|\omega^{\mu h}\| \|\tilde{\alpha}_t^{\mu h}(A) - \alpha_t(A)\| \\ &\leq \|\omega^{\mu h}\| |\mu| \|h\| \|A\| (\|f_1\|_1 + \|f_2\|_1), \end{aligned}$$

so that (2.9b) follows from this estimate.

We shall denote the set of  $(\hat{f}_1, \hat{f}_2)$ -stable states by  $I_{1,2}$ . In order to study the consequences of  $(f_1, f_2)$  stability we shall derive a condition which involves only the unperturbed entities  $\omega$  and  $\alpha_t$  and the functions  $f_1, f_2$ . Here  $\check{f}(x)$  denotes  $f(-x)$ .

**Proposition 2.4:** Let  $\omega \in I_{1,2}$  be continuous in the  $\sigma(\mathbb{U}, N)$  topology. Then

$$\begin{aligned} & \int_{-\infty}^{\infty} dt f_1(t) \omega(A \alpha_t B) \\ &= \int_{-\infty}^{\infty} dt f_2(t) \omega(\alpha_t(B) A), \quad \text{for } A, B \in \mathbb{U}. \end{aligned} \quad (2.10)$$

*Proof:* From Lemma 2.3 it follows that without loss of generality one may assume  $\omega^{\mu h}$  to be approximately invariant for  $\tilde{\alpha}_t^{\mu h}$ , in the sense of (2.9b). For  $h = h^* \in \mathbb{U}$ ,  $A \in \mathbb{U}$  we write

$$\int_{T_1}^{T_2} dt \frac{d}{dt} [\omega^{\mu h}(\tilde{\gamma}_t^{\mu h} A)] = \omega^{\mu h}(\tilde{\gamma}_{T_2}^{\mu h} A) - \omega^{\mu h}(\tilde{\gamma}_{T_1}^{\mu h} A).$$

With the use of (2.5c) and (2.6) we find

$$\begin{aligned} & \int_{-\infty}^{\infty} dt [\check{f}_1(t) \omega^{\mu h}(\tilde{\alpha}_t^{\mu h}(h \alpha_{-t} A)) \\ & - \check{f}_2(t) \omega^{\mu h}(\tilde{\alpha}_t^{\mu h}(\alpha_{-t}(A) h))] \\ &= (i/\mu) [\omega_+^{\mu h}(A) - \omega_-^{\mu h}(A)]. \end{aligned}$$

The right-hand side vanishes as  $\mu \rightarrow 0$  due to (2.7). Because  $f_i \in L_1(\mathcal{R})$  the Lebesgue dominated convergence theorem yields for  $\mu \rightarrow 0$ ,

$$\int_{-\infty}^{\infty} dt f_1(t) \omega(h \alpha_t A) = \int_{-\infty}^{\infty} dt f_2(t) \omega(\alpha_t(A) h), \quad (*)$$

for  $h = h^* \in \mathbb{U}$  and  $A \in \mathbb{U}$ . Now consider the GNS representation  $(\mathfrak{h}_\omega, \pi_\omega, \Omega_\omega)$  associated with the state  $\omega$ . Owing to Kaplansky's density theorem  $\pi_\omega(h) = \pi_\omega(h)^*$  can be approximated strongly by a net  $h_\alpha \in \pi_\omega(\mathbb{U}_0)$  of self-adjoint elements.

With the help of a three- $\epsilon$  argument and the polarization method we can now extend (\*) to all  $h \in \mathbb{U}$ .

Throughout this paper we shall adopt the  $\sigma(\mathbb{U}, N)$  continuity which we assumed in the preceding proposition.

### III. INVARIANCE, SEPARATING CHARACTER, AND THE MODULAR GROUP

From the stability condition (2.11) we now explore the ensuing properties of a state  $\omega \in I_{1,2}$ . For a state to be stable it should at least be time invariant; i.e., invariant for  $\alpha_t$ . To deal with this problem we formulate the following.

**Lemma 3.1:** (See Ref. 1.) Let  $F$  be a bounded function of two variables and  $h \in L_1(\mathcal{R}_2)$ . If  $F(h) = \int ds dt F(s,t)h(s,t)$  vanishes for all  $h$  with  $\hat{h}(p,q)$  having compact support not containing  $q = 0$  and  $h(s,t)$  is differentiable with respect to  $t$ , with  $\partial h(s,t)/\partial t \in L_1(\mathcal{R}_2)$ , then

$$F(s,t) = G(s),$$

for some bounded function  $G$ . Now we are able to prove the desired invariance.

**Proposition 3.2:** If  $\omega \in I_{1,2}$  then  $\omega$  is invariant for the unperturbed dynamics  $\alpha_t$ .

**Proof:** Let  $A = 1$  and  $B = C_g = \int dt g(t)\alpha_t(C)$ , then (2.11) yields

$$\omega(C_{h_1} - C_{h_2}) = \omega(C_{h_1 - h_2}) = 0, \quad (3.1)$$

for  $C \in \mathbb{U}_0, \hat{g} \in D$ , and  $h_j = f_j * g$ . From Lemma 3.1 we now conclude that  $\omega(\alpha_t(C))$  is a constant for all  $C \in \mathbb{U}_0$ . Invoking the continuity of  $\omega$  yields invariance, viz.  $\omega \circ \alpha_t = \omega$ .

To proceed further it is now convenient to write the stability condition (2.11) in the GNS representation. Let  $(\mathfrak{h}, \pi, \Omega)$  be the GNS triple associated with  $\omega \in I_{1,2}$ . Since  $\omega$  is invariant, the group of \* automorphisms  $\{\alpha_t\}$  can be implemented by a strongly continuous group of unitaries on  $\mathfrak{h}$ . To this end we must also assume that the correlation functions  $t \rightarrow \omega(A\alpha_t B)$  are continuous. Explicitly, we then have  $\pi(\alpha_t(A)) = U_t \pi(A) U_t^{-1}$  and  $U_t \Omega = \Omega$ . The stability criterion (2.11) can now be written as

$$\int dt f_1(t) (\Omega, A U_t B \Omega) = \int dt f_2(t) (\Omega, B U_t^{-1} A \Omega), \quad (3.2)$$

for  $A, B \in \pi(\mathbb{U})$ .

The infinitesimal generator of  $U_t$ , i.e., the Liouville operator, will be denoted by  $L$ , with the spectral representation  $L = \int dE_\lambda \lambda$ .

**Proposition 3.3:** If  $\omega \in I_{1,2}$  then  $\Omega$  is separating.

**Proof:** From (3.2) we have

$$(\Omega, A \hat{f}_1(L) B \Omega) = (\Omega, B \hat{f}_2(L) A \Omega), \quad (3.3)$$

where  $\check{f}(\lambda) = f(-\lambda)$ . Let  $A \Omega = 0$  then

$$(\Omega, A \hat{f}_1(L) B \Omega) = 0,$$

so that

$$(\hat{f}_1(L) A * \Omega, B \Omega) = 0, \quad (3.4)$$

for all  $B \in \pi(\mathbb{U})$ . Because  $\hat{f}_1$  is invertible on  $\text{sp } \alpha = \{\lambda \mid \hat{g}(\lambda) = 0 \forall g: \int g(t)\alpha_t(A) dt = 0 \forall A \in \mathbb{U}_0\} \supset \{\lambda \mid \hat{g}(\lambda) = 0 \forall g: \int g(t)U_t A \Omega dt = 0 \forall A \in \pi(\mathbb{U}_0)\} = \text{sp } L$ ,  $f_1$  is also invertible on  $\text{sp } L$ . Since  $\Omega$  is cyclic it follows now from (3.4) that  $A * \Omega = 0$  and therefore  $A = 0$  by standard arguments.<sup>9</sup>

As a consequence we have that the state  $(\Omega, \cdot \Omega)$  is a Tomita state on  $\pi(\mathbb{U})$  with the modular automorphism  $\sigma_t(\cdot) = \Delta^{it} \cdot \Delta^{-it}$ , where  $\Delta$  is the modular operator.

The preceding results already show great similarity of the  $(f_1, f_2)$ -stable states with thermodynamic equilibrium states. We shall make this connection more explicit in the following.

**Theorem 3.4:** If  $\omega \in I_{1,2}$  with  $f_{1,2}$  such that  $\phi(\lambda) = \hat{f}_1(\lambda)/\hat{f}_2(\lambda)$  satisfies

$$(i) \phi(\lambda)\check{\phi}(\lambda) = 1, \quad \text{with } \check{\phi}(\lambda) = \phi(-\lambda);$$

$$(ii) \phi(\lambda) > 0 \quad \text{for } \lambda \in \text{sp } L;$$

then  $\omega \in K_\phi$ , i.e.,  $\omega$  is a  $\phi$ -KMS state.

Whenever (i) and (ii) are not fulfilled  $I_{1,2} = \emptyset$ . Conversely, if  $\omega \in K_\phi$  then  $\omega \in I_{1,2}$  for some nonunique  $f_{1,2}$  such that  $\hat{f}_1/\hat{f}_2 = \phi$ .

**Proof:** Suppose  $\omega \in I_{1,2}$ , then from (3.3) it follows that

$$(A * \Omega, \check{f}_1(L) B \Omega) = (B * \Omega, \hat{f}_2(L) A \Omega).$$

Now choose  $A = A_g$  with  $\hat{g} = \check{f}_1$  then

$$\begin{aligned} & (\check{f}_1(L) A * \Omega, \check{f}_1(L) B \Omega) \\ &= (B * \Omega, \hat{f}_2(L) \check{f}_1(L) A \Omega) \\ &= (\check{f}_1(L) B * \Omega, [\hat{f}_2(L)] / [\check{f}_1(L)] \check{f}_1(L) A \Omega). \end{aligned}$$

Furthermore, we may let  $B = A *$  and since we assumed  $\phi = \hat{f}_1/\hat{f}_2 > 0$  on  $\text{sp } L$  we have

$$\| \check{f}_1(L) A * \Omega \| = \| \check{\phi}(L)^{1/2} \check{f}_1(L) A \Omega \|. \quad (3.5)$$

Since  $\check{f}_1(L)$  is invertible we can use the same reasoning as in Ref. 5 to conclude from (3.5) that the modular operator  $\Delta$  can be written

$$\Delta^2 = \check{\phi}(L) = \left[ \frac{\hat{f}_1(L)}{\hat{f}_2(L)} \right]^\vee = \frac{\hat{f}_2(L)}{\hat{f}_1(L)}. \quad (3.6)$$

It was shown in Ref. 2 that (3.6) is equivalent with  $\omega \in K_\phi$ .

The proof of the converse is quite easy. If  $\omega \in K_\phi$  then

$$\int dt f_\phi(t) \omega(A \alpha_t B) = \int dt f(t) \omega(\alpha_t(B) A), \quad (3.7)$$

for all  $A, B \in \mathbb{U}$ ,  $\hat{f} \in D$ , and  $\hat{f}_\phi = \phi \hat{f}$ , with  $\phi \in C^\infty(\mathcal{R})$ . Now choose  $\hat{F} \neq 0$  on  $\text{sp } L$  and a sequence  $(\hat{g}_n)_{n=1}^\infty$  in  $D$ , such that  $\hat{g}_n \rightarrow \hat{F}$  and  $\phi \hat{g}_n \rightarrow \hat{F}_\phi$  in  $S$ . As  $\phi$  may have an essential singularity at infinity, owing to a theorem of Weierstrass, we can write  $\phi = \exp(g)$ . Here  $g$  is odd and finite in the finite complex plane. Now choose a function  $h$  with a Laurent expansion such that  $\psi = \exp(h) \in S$  and  $\phi \psi \in S$ . Then we have for any  $\hat{G} \in S$  that  $\hat{F} = \psi \hat{G} \in S$  and  $\hat{F}_\phi = (\phi \psi) \hat{G} \in S$ . Then it follows that  $\omega$  satisfies (2.11), with  $f_1 = F_\phi$  and  $f_2 = F$ , where  $F, F_\phi \in L_1(\mathcal{R}) \cap C^1(\mathcal{R})$ . Obviously we have  $\hat{f}_1/\hat{f}_2 = \phi$ ; and since  $\check{\phi}(L) = \Delta^2$ , (3.7) can only be satisfied if on  $\text{sp } L$   $\phi > 0$  and  $\phi \check{\phi} = 1$  (Ref. 2, Lemma 4).

We conclude with a further remark which can now be made regarding the set  $I_{1,2}$ .

**Remark 3.5:** From (3.6) it follows that the modular operator  $\Delta$  commutes with the Liouville operator  $L$ . Then one may follow the line of reasoning given in Ref. 10 to establish that  $I_{1,2}$  is a lattice in its own order. In general  $I_{1,2}$  will not be closed and hence *a fortiori* not compact. If one assumes in addition the compactness of  $I_{1,2}$  in the  $w^*$  topology, then it follows that  $I_{1,2}$  is a Choquet simplex.

## ACKNOWLEDGMENT

Part of this work was done while one of us (J.S.C.) was at the Institute for Theoretical Physics of the University of Amsterdam as a member of the scientific staff of the "Stichting Fundamenteel Onderzoek der Materie" (FOM), which is financially supported by the "Nederlandse Organisatie voor Zuiver Wetenschappelijk Onderzoek" (ZWO).

<sup>1</sup>O. Bratteli and D. W. Robinson, *Operator Algebras and Quantum Statistical Mechanics* (Springer, Berlin, 1979, 1981).

<sup>2</sup>J. S. Cohen, H. A. M. Daniëls, and M. Winnink, *Commun. Math. Phys.* **84**, 449 (1982).

<sup>3</sup>J. L. Lebowitz, M. Aizenman, and S. Goldstein, *J. Math. Phys.* **16**, 1284 (1975).

<sup>4</sup>D. Kastler, *Symp. Math.* **20**, 49 (1976).

<sup>5</sup>F. Hoekman, Ph.D. thesis, University of Groningen, 1977.

<sup>6</sup>W. Pusz and S. L. Woronowicz, *Commun. Math. Phys.* **58**, 273 (1978).

<sup>7</sup>G. ten Brinke and M. Winnink, *Commun. Math. Phys.* **51**, 135 (1976).

<sup>8</sup>F. P. Greenleaf, *Invariant Means on Topological Groups* (Van Nostrand, New York, 1969).

<sup>9</sup>M. Winnink, "Some general properties of thermodynamic states in an algebraic approach," in *Statistical Mechanics and Field Theory*, edited by R. N. Sen and C. Weil (Keter, Jerusalem, 1971).

<sup>10</sup>M. Takesaki and M. Winnink, *Commun. Math. Phys.* **30**, 129 (1973).

# Nonequilibrium entropy in classical and quantum field theory

Henry E. Kandrup<sup>a)</sup>

Center for Theoretical Physics, University of Maryland, College Park, Maryland 20742

(Received 31 July 1986; accepted for publication 18 February 1987)

This paper proposes a definition of nonequilibrium entropy appropriate for a bosonic classical or quantum field, viewed as a collection of oscillators with equations of motion which satisfy a Liouville theorem (as is guaranteed for a Hamiltonian system). This entropy  $S$  is constructed explicitly to provide a measure of correlations and, as such, is conserved absolutely in the absence of couplings between degrees of freedom. This means, e.g., that there can be no entropy generation for a source-free linear field in flat space, but that  $S$  need no longer be conserved in the presence of couplings induced by nonlinearities, material sources, or a nontrivial dynamical background space-time. Moreover, through the introduction of a "subdynamics," it is proved that, in the presence of such couplings, the entropy will satisfy an  $H$ -theorem inequality, at least in one particular limit. Specifically, if at some initial time  $t_0$  the field is free of any correlations, it then follows rigorously that, at time  $t_0 + \Delta t$ , the entropy will be increasing:  $dS/dt > 0$ . Similar arguments demonstrate that this  $S$  is the only measure of "entropy" consistent mathematically with the subdynamics. It is argued that this entropy possesses an intrinsic physical meaning, this meaning being especially clear in the context of a quantum theory, where a direct connection exists between entropy generation and particle creation. Reasonable conjectures regarding the more general time dependence of the entropy, which parallel closely the conventional wisdom of particle mechanics, lead to an interpretation of  $S$  which corroborates one's naive intuition as to the behavior of an "entropy."

## I. INTRODUCTION

Conventional wisdom holds that the "entropy" associated with some system should be interpreted probabilistically as a measure of how generic its state really is. A state which is comparatively random, and which could be realized in many different ways, has associated with it a large entropy; a state which is somehow improbable, requiring, e.g., a special preparation, is considered to have low entropy. "Equilibrium" is, in this context, interpreted as a state of "maximum randomness" and "maximum entropy." The content of Boltzmann's<sup>1</sup> classic  $H$ -theorem is that a system will evolve towards this state of maximum entropy.

There is, however, a well-known difficulty in implementing this general picture. Consider, for example, a collection of  $N$  classical point particles. Suppose in the usual way that this system is characterized by an  $N$ -particle distribution function  $\mu$  and that the evolution of this  $\mu$  is governed by an  $N$ -particle Liouville equation which expresses probability conservation. The standard paradigm then implies (i) that the entropy  $S = -\text{Tr} \mu \log \mu$ , where  $\text{Tr}$  denotes a trace over the degrees of freedom of the  $N$  particles, and (ii) that the unique equilibrium corresponds to a state  $\mu \propto \exp(-\beta H)$ , where  $\beta$  is a constant and  $H$  the  $N$ -particle Hamiltonian. The problem, however, is that the Liouville equation guarantees that  $dS/dt \equiv 0$ ! This  $S$  does not change with time and, consequently, there can be no systematic evolution towards an equilibrium state of maximum entropy.

This well-known difficulty led historically to the idea of

a "coarse-grained averaging," namely the notion that an  $H$ -theorem does not hold on a truly microscopic level, but, instead, holds only on a quasimacroscopic level for an appropriately averaged  $\mu$ . This is, at least superficially, an extremely attractive idea, but it suffers from two related drawbacks: (1) it seems difficult to generate a general algorithm to effect the desired coarse graining; and (2) even if one were to construct a working algorithm, one would be faced with the additional problem of demonstrating that it is "canonical" in some natural sense. There would remain, e.g., the task of either ascertaining the scale on which the averaging is to be implemented, or, alternatively, of demonstrating the scale invariance of the averaging.

To the extent that no canonical prescription exists, one seems forced ultimately to the viewpoint adopted by Jaynes, namely that "Entropy is a property, not of the physical system, but of the particular experiments you or I choose to perform on it."<sup>2</sup> It is, therefore, natural to ask what it is that one typically measures when one probes the state of the system. And the answer to that would appear quite clear. One seeks typically to measure the one-particle distribution function  $f$  (i.e., the probability density for finding a particle at a given point  $\mathbf{x}$  with momentum  $\mathbf{p}$  at time  $t$ ), or perhaps the pair or three-body correlation functions. But one does not even try to measure the detailed correlations amongst the particles buried in the full  $N$ -particle  $\mu$ . One might, therefore, argue that, as a practical matter, the physically relevant notion of coarse graining does not involve a macroscopic averaging but, instead, entails a loss of information about higher-order interparticle correlations.

From this point of view, it would seem natural to conjecture that the entropy of the system should be defined in terms of the reduced one-particle  $f$ , rather than the full  $N$ -particle  $\mu$ . And, as such, it would be natural to propose an entropy

<sup>a)</sup> Present address: Department of Physics, Syracuse University, Syracuse, New York 13244.

$$S(t) \equiv - \sum_i \int d\mathbf{x}_i d\mathbf{p}_i f(i) \log f(i), \quad (1.1)$$

where

$$f(i) \equiv f(\mathbf{x}_i, \mathbf{p}_i, t) = \int \prod_{j \neq i} d\mathbf{x}_j d\mathbf{p}_j \mu(\mathbf{x}_1, \mathbf{p}_1, \dots, \mathbf{x}_N, \mathbf{p}_N; t). \quad (1.2)$$

There are solid reasons to believe that this  $S$  provides a reasonable notion of entropy. Most obvious is the fact that this  $S$  will in general be time dependent. Thus, for example, allowing for particle interactions derived from two-body potentials  $V_{ij}$ , one concludes that, for a system of  $N$  identical particles,

$$\begin{aligned} \frac{dS}{dt} &= N(N-1) \int d\mathbf{x}_i d\mathbf{p}_i \int d\mathbf{x}_j d\mathbf{p}_j f_2(i, j) \\ &\times \frac{\partial V_{ij}}{\partial \mathbf{x}_i} \cdot \frac{\partial}{\partial \mathbf{p}_i} \log f(i), \end{aligned} \quad (1.3)$$

where  $f_2(i, j)$  is the reduced two-particle distribution function for particles  $i$  and  $j$ .

Also significant is the fact that this  $S$  is truly canonical, being constructed in a systematic and unambiguous fashion from the one-particle  $f$ , an object of obvious physical significance. Of particular relevance in this regard is the fact that this  $f$  satisfies a "subdynamics,"<sup>3</sup> decoupled from the higher-order correlations. Specifically, by means of projection operator techniques, one can derive for the evolution of  $f$  an exact, closed (albeit nonlocal and nonlinear) equation which contains no explicit reference to the higher-order correlations buried in such quantities as  $f_2(i, j)$ . It is the nonlinearity of this equation which leads to a nonconserved entropy; it is the linearity of the fundamental Liouville equation that guarantees that  $\text{Tr } \mu \log \mu$  is a constant of the motion.

This  $S$  coincides, moreover, with the entropy entering into the standard  $H$ -theorem; and, consequently, one anticipates that this  $S$  really will increase monotonically at least in some approximate limit. It is in fact well known that the exact equation satisfied by  $f$  reduces to the standard Landau equation<sup>3</sup> if one assumes (i) that the interactions are weak and comparatively short range (dilute gas approximation), (ii) that initial conditions were specified at a time in the past long compared with the duration of a typical interaction, and (iii) that one can neglect the effects of nontrivial initial conditions and suppose that, at some initial time  $t_0$ ,  $\mu = \Pi_i f(i)$ .

The first two of these requirements seem reasonable physically; and indeed, they can be relaxed, at least in principle, in the context of a systematic perturbation expansion. The third requirement needs some further justification. For "ordinary," reasonably well behaved interactions, one can argue convincingly, and in certain cases prove,<sup>3</sup> that generic nontrivial initial conditions will in fact decay as time goes by. The physical content of this statement is that any correlations present at the outset will eventually become irrelevant compared with the systematically evolving correlations generated by the subsequent dynamics. For more perverse, long-range interactions, like Newtonian gravity, this argument is most likely invalid: numerical simulations suggest that self-

gravitating systems really do "remember" their initial conditions for a very long time.<sup>4</sup> For these perverse interactions, a systematic statistical analysis really requires the ansatz that no significant initial correlations be present. Whether this is reasonable is of course exceedingly difficult to say. Such a viewpoint is, however, if nothing else, consistent with the speculation that the Universe originated from a state of "maximum simplicity." After all, it is only on very large scales that self-gravity becomes important.

It should also be stressed that, even neglecting "special" initial conditions, not every physical system of interacting entities will satisfy an  $H$ -theorem for all times. Consider, for example, a pair of harmonic oscillators with natural frequencies  $\omega_1$  and  $\omega_2$ , connected by a linear coupling derived from a pair potential  $V_{12} = \lambda x_1 x_2$ , where  $\lambda$  measures the strength of the interaction. In this simple case, one can solve explicitly for the evolution of the system in terms of arbitrary initial conditions; and, for appropriate choices of  $\omega_1$ ,  $\omega_2$ , and  $\lambda$ , it is easy to see that the motion will be periodic. After some time  $\tau$ , the system will return to its initial state: there can be no progression towards equilibrium. This means that, if the entropy increases at one point of time, it must at some other time decrease.

This does not, however, imply that the entropy (1.1) is devoid of physical meaning. Even if this  $S$  does not satisfy an  $H$ -theorem, it can provide a useful measure of the degree of correlations in the system. Thus, e.g., if one supposes that, at some initial time  $t_0$ , the system was free of correlations, the interactions between the oscillators will result in an initial generation of correlations and a concomitant increase in entropy!

This initial entropy increase is in fact a very general result. Consider a collection of  $N$  objects which interact via arbitrary two- and higher-body forces characterized by a coupling constant  $\lambda$ . Assume then that, at some initial time  $t_0$ , the system was completely free of correlations, so that  $\mu = \Pi_i f(i)$ . It follows that, an instant  $\Delta t$  later,

$$\frac{dS(t_0 + \Delta t)}{dt} = \lambda^2 |a|^2 \Delta t > 0, \quad (1.4)$$

where  $|a|^2$ , which is intrinsically positive, reflects the form of the initial state. A proof of this claim for a specific model interaction was provided in Ref. 5. A more general proof will be provided below. The important point to note here is that, even if an  $H$ -theorem does not hold for all times, the entropy (2.1) can still provide a useful measure of correlations. To show that  $S$  reflects correlations, and to show that  $S$  always increases, are two distinct and separate issues.

To the extent that a system is truly periodic and no systematic evolution towards a more "random" state is obtained, one might perhaps argue that  $S$  does not warrant the appellation entropy. What does, however, appear to be an empirical fact is that, for realistic complicated systems, couplings between degrees of freedom lead to a progression towards a more random sort of state.

The object of this paper is to formulate a notion of non-equilibrium entropy for a classical or quantum field which parallels as closely as possible the particle entropy described in this Introduction. This notion of entropy is such that it is conserved absolutely for a source-free linear field in Min-

kowski space, but assumes a time dependence in the presence of couplings induced, e.g., by nonlinearities, sources, or a nontrivial background space-time.

The program of this paper is as follows: Section II sketches out the basic picture for a classical or quantum field in Minkowski space, emphasizing in particular (i) the question of physical interpretation and (ii) the formal similarities between the classical and quantum theories. Section III then addresses the additional complications, both conceptual and practical, which arise for fields propagating in a nontrivial dynamical space-time. Section IV introduces the notion of a subdynamics to show explicitly that the field entropy  $S$  defined in Secs. II and III does in fact provide a measure of correlations, demonstrating that, in the presence of couplings between degrees of freedom, an initially uncorrelated state leads necessarily to an initial increase in entropy. Section V turns to the special question of physical interpretation in the context of a quantum theory, focusing upon a fundamental connection between changes in the field entropy and the phenomenon of particle creation. The principal conclusion here is that the mechanism which gives rise to an initial increase in entropy will also cause an initial enhancement in the rate of particle creation. Section VI demonstrates that, in a precise and well-defined mathematical sense, the  $S$  defined in Secs. II and III is the only measure of entropy consistent with the notion of subdynamics defined in Sec. IV. Finally, Sec. VII summarizes the principal results, speculates upon the more generic time dependence of  $S$ , and concludes by reflecting upon the connection between the entropy defined here and the "geometric" entropy associated, e.g., with the event horizon of a black hole in general relativity.

It should, perhaps, be noted that, although the analysis presented here is comparatively abstract, the key ideas were motivated originally by a desire to understand nonequilibrium processes in the early Universe. The sense in which these ideas are relevant there is considered in a companion paper by Hu and Kandrup.<sup>6</sup>

Units are chosen throughout such that Planck's constant  $\hbar/2\pi$  and the speed of light  $c$  are equal to unity.

## II. NONEQUILIBRIUM ENTROPY IN MINKOWSKI SPACE

### A. Classical field theory

As noted already, the object of this paper is to construct a notion of entropy for a classical or quantum field which parallels as closely as possible the notion of particle entropy posited in Sec. I. This implies, in particular, three specific requirements.

(1) This entropy  $S$  must provide a measure of the degree of correlations in the system. If there are no couplings between degrees of freedom, so that no correlations can be generated,  $S$  must be conserved.

(2) This  $S$  must satisfy an  $H$ -theorem inequality, at least in some appropriate limit. One demands that an "initially uncorrelated" system leads to an initial increase in  $S$ ; and one anticipates (or at least hopes) that  $S$  will increase monotonically for all times, this corresponding to an approach towards equilibrium.

(3) This  $S$  must in some sense be physically meaningful.

The key idea underlying the analysis is that the field entropy  $S$  should be constructed from an appropriate analog of the one-particle distribution functions  $f(i)$ . One knows that, at least in flat space, a classical field theory is equivalent mathematically to an infinite set of oscillators. The statistical description of such a field takes, therefore, as its fundamental object a many-oscillator distribution function  $\mu$ , the evolution of which is governed by a Liouville equation which expresses conservation of probability. Given this fundamental  $\mu$ , one can then define reduced distribution functions  $g(q_A, \pi_A, t) \equiv g(A)$  in the obvious way for each oscillator by integrating over the "position" and "momentum,"  $q_B$  and  $\pi_B$ , of each of the remaining oscillators:

$$g(A) \equiv g(q_A, \pi_A, t) = \int \prod_{B \neq A} dq_B d\pi_B \mu(\{q_C, \pi_C\}; t). \quad (2.1)$$

Given these one-oscillator  $g(A)$ 's, one is then instructed to define an entropy

$$S(t) = - \sum_A \int dq_A d\pi_A g(A) \log g(A). \quad (2.2)$$

Because the field is equivalent to an infinite set of oscillators, the sum in this expression is an infinite one, so that the entropy so defined might well prove infinite, at least formally. This, however, is not especially relevant for the present discussion. What is relevant is whether this  $S$  can be shown to increase in the presence of evolving correlations.

The first obvious point to note is that, although, in the absence of material sources,  $\text{Tr} \mu \log \mu$  is a constant of the motion, the entropy of Eq. (2.2) will in fact be time dependent if there exist couplings between degrees of freedom. It follows trivially that, for a linear free field,  $dS/dt \equiv 0$ , but, in the presence of nonlinearities such as those arising in a  $\lambda \Phi^p$  field theory, one concludes instead that, in general,  $dS/dt \neq 0$ .<sup>5</sup> This is no different from the statement that a system of  $N$  noninteracting particles must conserve its entropy, but that, once one allows for particle interactions, the entropy will in fact change with time.

One can, moreover, demonstrate that, in the presence of such couplings induced by nonlinearities,  $S$  will satisfy for short times the same sort of  $H$ -theorem as did the particle entropy of Sec. I. Specifically, if one supposes that, at some initial time  $t_0$ , the system was free of correlations, so that  $\mu = \prod_A g(A)$ , it follows rigorously that  $dS(t_0 + \Delta t)/dt > 0$ . The proof of this statement, provided in Sec. IV, is the same for particle and field theories. The crucial point is simply that, for early times,  $dS/dt$  will be quadratic in the interaction Liouvillian which generates the evolving correlations! For a system of particles interacting via two-body forces, the interaction Liouvillian for a pair of particles  $i$  and  $j$  takes the form

$$L_{ij}^I = \lambda \frac{\partial V_{ij}}{\partial \mathbf{x}_i} \cdot \frac{\partial}{\partial \mathbf{p}_i}, \quad (2.3)$$

where  $\lambda$  is a coupling constant. For a simple  $\lambda \Phi^3$  field theory, the analogous object is

$$\mathcal{L}_{AB}^I = \lambda q_B q_{A-B} \frac{\partial}{\partial \pi_A}. \quad (2.4)$$

In each case, one concludes that, to first order in  $\Delta t$ ,

$$\frac{dS(t_0 + \Delta t)}{dt} = \lambda^2 |a|^2 \Delta t > 0, \quad (2.5)$$

where  $|a|^2 > 0$  depends upon the form of the initial state.

In the presence of material sources (e.g., when considering a collection of charged particles interacting via an electromagnetic field) the analysis becomes more complicated mathematically, but the physical picture remains unchanged. The fundamental object for the composite system of particles and fields is now an enlarged distribution function  $\nu$  depending upon both particle and field variables which satisfies a Liouville equation expressing probability conservation in an enlarged particle-plus-field phase space. It follows trivially from the linearity of the Liouville equation that  $\text{Tr } \nu \log \nu$  is conserved absolutely, but that the total entropy

$$S \equiv - \sum_I \int dx_i, d\mathbf{p}_i f(i) \log f(i) - \sum_A \int dq_A d\pi_A g(A) \log g(A) \quad (2.6)$$

will in general exhibit a nontrivial time dependence. Local equations for  $\partial f(i)/\partial t$  or  $\partial g(A)/\partial t$  analogous to Eq. (1.3) now involve the particle-oscillator correlation function

$$h(i,A) = \int \prod_{j \neq i} dx_j d\mathbf{p}_j \int \prod_{B \neq A} dq_B d\pi_B \times \nu(\mathbf{x}_1, \mathbf{p}_1, \dots, \mathbf{x}_N, \mathbf{p}_N, \{q_C, \pi_C\}; t). \quad (2.7)$$

And the couplings buried in this  $h(i,A)$  induce a nontrivial  $dS/dt$  just as surely as do the couplings associated with the direct particle-particle interactions of ordinary Newtonian dynamics or the nonlinearities discussed above.

Since the physical picture here is not different from that arising in a nonlinear field theory, but significantly messier mathematically, this situation will not be considered in any detail in this paper. One may, however, note that it is completely straightforward to use the techniques developed in Ref. 7 to parallel the discussion here and obtain a theory of nonequilibrium entropy for an interacting system of particles and fields.

It remains here to at least address the issue of whether the field entropy (2.2) is physically meaningful. After all, the oscillators that one is considering are only mathematical constructs, and one would not expect to measure any  $g(A)$ 's in a realistic experiment. One reason to believe that the  $g(A)$ 's, and hence  $S$ , are meaningful is that the  $g(A)$ 's serve at least to define physically measurable average values. Thus, for example, if one considers a scalar field

$$\Phi(\mathbf{x}, t) \equiv \sum_k q_k(t) \exp(-i\mathbf{k} \cdot \mathbf{x}), \quad (2.8)$$

it follows immediately that the statistical average value

$$\langle \Phi(\mathbf{x}, t) \rangle = \sum_k \langle q_k(t) \rangle \exp(-i\mathbf{k} \cdot \mathbf{x}), \quad (2.9)$$

where

$$\langle q_k(t) \rangle = \text{Tr } \mu q_k(t) = \int dq_k d\pi_k g(k) q_k(t). \quad (2.10)$$

This is completely analogous to the situation in a particle theory. Here, for example, the true mass density

$$\rho(\mathbf{x}, t) = \sum_I m_i \delta_D[\mathbf{x} - \mathbf{x}_i(t)], \quad (2.11)$$

whereas the average density

$$\langle \rho(\mathbf{x}, t) \rangle = \sum_I m_i \langle n_i(\mathbf{x}, t) \rangle, \quad (2.12)$$

with

$$\begin{aligned} \langle n_i(\mathbf{x}, t) \rangle &= \text{Tr } \mu \delta_D(\mathbf{x} - \mathbf{x}_i) \\ &= \int d\mathbf{x}_i d\mathbf{p}_i f(i) \delta_D[\mathbf{x} - \mathbf{x}_i(t)]. \end{aligned} \quad (2.13)$$

One should note, moreover, that if an  $H$ -theorem holds,  $S$  can provide a useful diagnostic for the evolution of physical observables. For a free, linear field, each mode will of course evolve independently, so that the entropy, which is after all conserved, will not be of particular significance. Thus a scalar field satisfying the Klein-Gordon equation

$$-\partial_t^2 \Phi + \Delta \Phi = 0 \quad (2.14)$$

leads trivially to the averaged equation

$$-\partial_t^2 \langle \Phi \rangle + \Delta \langle \Phi \rangle = 0. \quad (2.15)$$

If, however, there exist correlations generated by couplings between the modes, one acquires an effective "source"  $\Sigma$ , and the evolving  $S$  can provide useful information about the effects of this  $\Sigma$ . Thus, for example, the true nonlinear equation

$$-\partial_t^2 \Phi + \Delta \Phi + \lambda \Phi^2 = 0 \quad (2.16)$$

leads to a statistically averaged equation

$$-\partial_t^2 \langle \Phi \rangle + \Delta \langle \Phi \rangle + \lambda \langle \Phi \rangle^2 = \Sigma, \quad (2.17)$$

where the source  $\Sigma$  involves couplings between each mode  $A$  and the modes  $B$  and  $A - B$ .<sup>5</sup>

## B. Quantum field theory

The discussion hitherto has focused exclusively upon a classical field theory. It is, however, completely straightforward to formulate a corresponding quantum theory for a bosonic system with integral spin, provided only that the fundamental dynamics can be cast into a Hamiltonian form<sup>8</sup> (the generalization to a spin- $\frac{1}{2}$  fermionic system is currently under investigation<sup>9</sup>). All that one need do is implement the standard formalism of canonical quantization, assuming that the Poisson bracket is to be replaced by a commutator. Indeed, one could equally well construct a quantum statistical theory for a non-Hamiltonian system provided only that one is willing to accept *some* (perhaps quite *ad hoc*) quantization prescription.

It is, however, important to emphasize that the interpretation of the field entropy (2.2) as providing a measure of correlations *does* rely upon one important feature guaranteed for Hamiltonian systems which will not, however, hold in general, namely the notion of conservation of phase expressed by the Liouville theorem. Consider, e.g., a collection

of  $N$  identical classical, noninteracting particles which feel the influence of an external force  $F_a$  which depends upon the coordinates, momenta, and time in a perverse fashion which cannot be formulated in a Hamiltonian framework. In this case, one still has a well-defined notion of probability conservation in an appropriate  $6N$ -dimensional  $(\mathbf{x}, \mathbf{p})$  space, and it follows at once that the  $N$ -particle  $\mu$  will satisfy the relation

$$\frac{\partial \mu}{\partial t} + \sum_{i=1}^N \frac{\partial}{\partial x_i^a} \left( \frac{dx_i^a}{dt} \mu \right) + \sum_{i=1}^N \frac{\partial}{\partial p_i^a} \left( \frac{dp_i^a}{dt} \mu \right) = 0, \quad (2.18)$$

so that, since the particles are noninteracting,

$$\frac{\partial f}{\partial t} + \frac{p^a}{m} \frac{\partial f}{\partial x^a} + \frac{\partial}{\partial p_a} (F_a f) = 0. \quad (2.19)$$

It does not, however, follow in the usual way that the Boltzmann entropy will be conserved! Rather, one calculates explicitly that

$$\frac{dS}{dt} = N \int dx dp \frac{\partial F_a}{\partial p_a} f(1 + \log f), \quad (2.20)$$

which certainly need not vanish. In the absence of conservation of phase, which would be guaranteed if

$$\frac{\partial}{\partial x^a} \left( \frac{dx^a}{dt} \right) + \frac{\partial}{\partial p_a} \left( \frac{dp_a}{dt} \right) = 0, \quad (2.21)$$

even a noninteracting theory entails a time-dependent entropy.

It may also be emphasized that this is not a purely technical observation with no physical relevance. A similar problem arises in general relativity if one wishes to reexpress geodesic flows in a space-time with metric  $g'_{\alpha\beta}$  in terms of a slightly different metric  $g_{\alpha\beta}$  (as would, e.g., be required for a phase space description of linearized perturbations away from some static background). In this case, the prescription of Israel and Kandrup<sup>10</sup> leads to a four-force

$$F_\alpha = \delta \Gamma_{\mu\nu}^\lambda \Delta_{\alpha\lambda} p^\mu p^\nu / m,$$

where  $\Delta_{\alpha\lambda}$  is the spatial projection tensor constructed from  $p_\alpha$  and  $g_{\alpha\lambda}$ , and  $\delta \Gamma_{\mu\nu}^\lambda$  is the difference between the Christoffel symbols associated with  $g'_{\alpha\beta}$  and  $g_{\alpha\beta}$ . It then follows that  $\partial F_\alpha / \partial p_\alpha \neq 0$ , and this implies that the entropy flux  $s^\mu$ , as defined, e.g., by Israel,<sup>11</sup> will not be divergence-free.

In any case, at least for Hamiltonian systems, the only really new feature of a quantum statistical description is that the distribution function  $\mu$  must be reinterpreted as a density matrix with an evolution governed by the usual quantum Liouville equation. Reduced distribution functions, such as  $g(A)$ , realized as reduced density matrices, are obtained by partial traces over the degrees of freedom of some subset of the oscillators. The field entropy  $S$  is then defined in terms of the reduced  $g(A)$ 's by the obvious prescription

$$S = - \sum_A \text{Tr}_A g(A) \log g(A), \quad (2.22)$$

where  $\text{Tr}_A$  denotes a trace over the degrees of freedom of the  $A$  th oscillator, realized in an arbitrary (e.g., coordinate or momentum) representation.

As illustrated in Sec. IV, one can, either classically or in the framework of a quantum description, formulate a sub-

dynamics for the evolution of the  $g(A)$ 's which contains no explicit reference to the higher-order correlations. And, by using this subdynamics, one can again prove an initial entropy increase (2.5) for a system which evidences no initial correlations. There remains the well-known problem of obtaining a true probabilistic interpretation of the quantum  $\mu$  or  $g(A)$ 's, but this has no immediate bearing upon the interpretation of  $S$  as a measure of correlations.

Given the quantum or classical equations of motion, and the existence of a density matrix or distribution function  $\mu$  which satisfies a linear Liouville equation, the quantum and classical statistical theories are essentially identical!

The quantum theory does, however, admit to one enrichment of interpretation. Specifically, in the context of a quantum description, it is conventional to interpret a change in the field as representing the creation or destruction of particles. Thus, in particular, if the  $A$  th oscillator is characterized by a natural frequency  $\omega_A$ , it is customary to interpret the statistical average

$$\langle N_A \rangle \equiv \langle \frac{1}{2} \pi_A^2 + \frac{1}{2} \omega_A^2 q_A^2 \rangle / \omega_A - \frac{1}{2} \quad (2.23)$$

as representing the "number of quanta in the  $A$  th mode." One might, therefore, seek to establish a connection between changes in the entropy (2.22) and changes in the average  $\langle N_A \rangle$ . The gratifying fact, discussed in great detail by Hu and Kandrup,<sup>6</sup> and considered briefly in Sec. V, is that such a connection does exist. Evolving correlations lead not only to an increase in the entropy, but to an increase in particle number.

That this is the case is illustrated by the following statement, a proof of which is presented in Sec. V. Consider a source-free nonlinear field theory in flat space realized as a collection of oscillators, and suppose that the  $\langle N_A \rangle$  of Eq. (2.23) represents the number of particles in the  $A$  th mode. Write the total density matrix  $\mu$  in the form

$$\mu = \prod_A g(A) + \mu_I, \quad (2.24)$$

where  $\mu_I \equiv \mu - \prod_A g(A)$  reflects the "piece" of the total  $\mu$  which contains information about correlations amongst the degrees of freedom. Suppose then that, at some initial time  $t_0$ ,  $\mu_I$  vanishes identically. It follows rigorously that, at time  $t_0 + \Delta t$ , the contribution to  $d \langle N_A \rangle / dt$  which involves the generated  $\mu_I(t_0 + \Delta t)$  is intrinsically positive.

### III. NONEQUILIBRIUM ENTROPY IN CURVED SPACETIMES

#### A. Conceptual issues

The further generalization of this basic picture to a quantum field theory in a fixed curved background space-time is again comparatively straightforward, at least formally, provided only that one can implement a preferred  $3 + 1$  splitting into space and time, and that one specifies a preferred set of spatial functions at each instant of time to generalize the standard decomposition into plane waves. Provided that the space-time is not too perverse, the Cauchy problem will be well defined, and, given a consistent notion of dynamics, a statistical description should be possible. The only new significant technical complications arise from the fact (i)



that, in general, one cannot expand in spatial plane waves, and (ii) that, in a dynamical space-time, the three-Hamiltonian  $H$  derived from the fundamental action will manifest an explicit time dependence.

The fact that  $H$  is time dependent does not, in itself, generate any serious conceptual or practical problem.<sup>12</sup> A time-dependent  $H$  most likely precludes the possibility of an equilibrium,<sup>11</sup> but a nonequilibrium statistical description is still quite possible. At a conceptual level, field theory in a classical time-dependent gravitational field is not very different, in many respects, from quantum electrodynamics in a time-dependent, externally imposed, classical electromagnetic field.

Strictly speaking, the introduction of a specific 3 + 1 splitting, or, equivalently, the specification of a family of preferred observers, arises already in Minkowski space. Thus ordinary quantum field theory is usually formulated from the point of view of inertial observers; and it is well known that, even for the simplest case of a source-free linear field theory, an accelerated observer will interpret the physics very differently. Thus, for example, the state which, to an inertial observer, corresponds to a true vacuum will, to a uniformly accelerated observer, correspond instead to a thermal state with a temperature proportional to its acceleration.<sup>13</sup>

One might, therefore, conjecture that in curved, as in flat, space one ought simply to restrict attention to freely falling observers. This, however, does not suffice to solve the problem. The symmetries of Minkowski space imply that any inertial observer will see "space" as homogeneous and isotropic, free of event horizons and other global complications, so that there is a natural decomposition of any field into plane waves  $\exp(-i\mathbf{k}\cdot\mathbf{x})$ . In curved spaces, however, there is in general no obvious analog of these plane waves, and two different freely falling observers might find it convenient to expand in two very different sets of "spatial" functions.

Given that the natural spatial functions are no longer necessarily plane waves, it is not especially convenient to introduce the concept of a Wigner function except in a type of "quasilocal" approximation. Indeed, the conventional noncovariant construction of such Wigner functions by means of the Weyl prescription exploits, in a deep and fundamental way, the space translational symmetry of Minkowski space.<sup>14</sup> And, similarly, the covariant Wigner functions of flat space exploit the time translational symmetry as well.<sup>15</sup>

This does not, however, imply that Wigner functions are totally useless in curved space-times. By assuming that the space-time is comparatively smooth, i.e., that the gravitational "field" does not change too quickly in space and time, one can use a Riemann normal coordinate construction to obtain a notion of the Wigner function which does provide useful insights into the overall dynamics. Thus, for example, one can show that a free scalar field may be characterized by a Wigner function  $f_w(x^\alpha, p_\alpha)$  which, in a first approximation, satisfies the "collisionless Boltzmann" equation<sup>16</sup>

$$\frac{df_w}{d\tau} = p^\alpha \frac{\partial f}{\partial x^\alpha} + \Gamma_{\alpha\mu}^\lambda p_\lambda p^\mu \frac{\partial f}{\partial p_\alpha} = 0, \quad (3.1)$$

the same equation which would be satisfied by a collection of noninteracting classical point masses. This is a very intuitive sort of relation, implying simply that  $f_w$  is conserved under Lie transport along the classical geodesics

$$\frac{dx^\alpha}{d\tau} = p^\alpha \quad \text{and} \quad \frac{dp_\alpha}{d\tau} = \Gamma_{\alpha\mu}^\lambda p_\lambda p^\mu. \quad (3.2)$$

The important point to observe, however, is that, whereas this equation is exact for classical point masses, it is only an approximation for the quantum field. A more careful analysis shows that<sup>16</sup>

$$\frac{df_w}{d\tau} = C[f_w], \quad (3.3)$$

where  $C[f_w]$ , which can be calculated perturbatively, involves the Riemann curvature tensor and derivatives thereof.

This is qualitatively similar to what happens to a charged particle in flat space. Here, classically, the one-particle  $f(x^\alpha, p_\alpha)$  will satisfy exactly the equation

$$\frac{df}{d\tau} = p^\alpha \frac{\partial f}{\partial x^\alpha} + eF_{\alpha\beta}^\mu p^\mu \frac{\partial f}{\partial p^\alpha} = 0, \quad (3.4)$$

where  $e$  is the charge and  $F_{\alpha\beta}$  the Maxwell tensor. However, the Wigner function  $f_w(x^\alpha, p_\alpha)$  appropriate for a Dirac field will instead satisfy an equation of the form

$$\frac{df_w}{d\tau} = \mathcal{C}[f_w], \quad (3.5)$$

where  $\mathcal{C}[f_w]$  involves the derivative  $\nabla_\alpha F_{\mu\nu}$ . In each case, one acquires corrections proportional to the "tidal forces" acting between nearby points which may be interpreted as reflecting the effects of "virtual particles."

For better or ill, the basic point of view adopted in this paper circumvents completely the notion of the Wigner function by working directly with the density matrix  $\mu$ . One advantage therein is the fact that one has not built into the basic formalism any underlying assumptions that reflect the special symmetries of flat space. The density matrix  $\mu$  has a natural meaning in and of itself in terms of the 3 + 1 splitting and an arbitrary choice of spatial functions.

Another related feature of this description is that it is intrinsically nonlocal. The  $g(A)$ 's are defined in an abstract  $(q, \pi)$  space which has no direct connection with the space-time manifold or the associated cotangent bundle. This means that there is no natural sense in which some piece of the total entropy  $S$  can be associated with some piece of the space-time. The  $g(A)$ 's are dependent upon the choice of spatial functions which contain important information about the global properties of the space-time. As will be discussed in Sec. VII, this seems a desirable feature if one wishes to establish a connection between the ideas presented here and the geometric entropy associated, e.g., with a black hole in a static space-time. Here, following Bekenstein and Hawking,<sup>17</sup> one is wont to associate with the space-time an entropy proportional to the area of the event horizon, but it is clear that the existence of such an event horizon is manifest only in a global description of the physics.

This nonlocality does, however, have one perhaps undesirable implication, namely that it would seem difficult, if

not impossible, to construct a natural covariant generalization. The situation is, e.g., very different in relativistic kinetic theory.<sup>12</sup> There one is wont to associate with the distribution of matter in the space-time an entropy flux  $s^\mu(x^\alpha)$ , a vector field defined in the space-time manifold. The content of an  $H$ -theorem is then encapsulated in the covariant statement that the covariant divergence of  $s^\mu$  is intrinsically non-negative:  $\nabla_\mu s^\mu \geq 0$ . Only by breaking manifest covariance and integrating over some arbitrary spacelike hypersurface does one reach the noncovariant (but observer-independent!) conclusion that  $dS/dt \geq 0$ .

In the absence of a covariant theory, it is not at all clear that the monotonic increase of the field entropy (2.22) is a statement with which all observers would agree. Indeed, one is confronted with two even more fundamental questions: (1) Will all observers agree that correlations do, or do not, exist between degrees of freedom for the field?; and (2) will they even agree whether couplings exist between these degrees of freedom? These, however, are issues which arise already in ordinary quantum field theory in curved space-time, and, especially, in analyses of particle creation, so that it is probably fair to say: In this regard, at least, the notion of entropy presented here is as reasonable—or unreasonable—as the standard discussions of particle creation.

In any case, to recapitulate: If one is willing to implement a preferred  $3 + 1$  splitting and a preferred decomposition into spatial functions, the formal structure of the statistical description is not very different from that arising in Minkowski space. What is different are the new sorts of physical implications induced by a nontrivial dynamical background space-time.

## B. Three examples

### 1. A static space-time

In this case, there exists a natural  $3 + 1$  splitting,  $t$  being the coordinate associated with the time translational invariance. The metric can then be written in the form

$$ds^2 = g_{\mu\nu} dx^\mu dx^\nu = g_{tt}(dt)^2 + g_{ab} dx^a dx^b \quad (a, b = 1, 2, 3), \quad (3.6)$$

where the metric functions  $g_{tt}$  and  $g_{ab}$  depend only upon the spatial coordinates  $x^a$ . As a simple illustration, consider the minimally coupled massive Klein–Gordon field satisfying

$$\nabla_\mu \nabla^\mu \Phi - m^2 \Phi = 0, \quad (3.7)$$

where  $\nabla_\mu$  denotes a covariant derivative. For the metric of Eq. (3.6),  $\Phi$  satisfies the equation

$$g^{tt} \frac{\partial^2 \Phi}{\partial t^2} + (-g)^{-1/2} \frac{\partial}{\partial x^a} (-g)^{1/2} g^{ab} \frac{\partial}{\partial x^b} \Phi - m^2 \Phi = 0, \quad (3.8)$$

and, consequently, it is natural to expand in terms of the eigenfunctions of the operator

$$\Delta \equiv (-g^{tt})^{-1} \left\{ (-g)^{-1/2} \frac{\partial}{\partial x^a} (-g)^{1/2} g^{ab} \frac{\partial}{\partial x^b} - m^2 \right\}, \quad (3.9)$$

the natural generalization of the flat space Laplacian. One

discovers thereby that the modes effectively decouple, each  $q_A$  satisfying a distinct equation; and, consequently, it follows that the entropy is conserved. For a source-free linear field in a static space-time,  $dS/dt \equiv 0$ .

This result manifests an obvious connection with the phenomenon of particle creation, it being well known that, in the absence of perverse global effects like event horizons, a source-free linear field in a static space-time admits no particle creation.<sup>18</sup> If, however, one allows for nonlinearities in the basic field equation, as would be induced, e.g., by a  $\lambda \Phi^p$  coupling, and if, as would seem natural, one still expands in terms of the eigenfunctions of  $\Delta$ , the  $q_A$ 's will be coupled with one another. And this coupling guarantees that  $dS/dt$  will no longer be conserved. Correlations will evolve and these will lead to changes in the entropy. This change in  $S$  is again related to particle creation. As will be shown in Sec. V, if the system is free of correlation at some initial time  $t_0$ , the evolving correlations induced by the dynamics will lead both to an initial increase in entropy and to an enhancement in the rate of particle creation.

### 2. A conformally static space-time

A net overall expansion or contraction has comparatively little effect upon the evolution of the entropy. It is natural now to write the metric in the form

$$ds^2 = \Omega^2(t) [g_{tt}(dt)^2 + g_{ab} dx^a dx^b], \quad (3.10)$$

which differs from Eq. (3.6) only by the presence of the conformal factor  $\Omega(t)$ . And, with such a metric, it is again natural to expand the Klein–Gordon field (3.7) in terms of the eigenfunctions of the  $\Delta$  of Eq. (3.9). The equation satisfied by each  $q_A$  now becomes more complicated,<sup>19</sup> but, nevertheless, one discovers thereby that, in the absence of nonlinearities, the modes remain decoupled so that  $dS/dt \equiv 0$ . If, alternatively, one allows for nonlinearities, one again induces couplings between the modes which lead to a nonconserved entropy and a simple  $H$ -theorem inequality.

For the case of a conformally static space-time, the connection between entropy generation and particle creation becomes more subtle, it being well known that even a source-free linear field theory can lead to the creation of particles. The point to observe, however, is that particle creation can be triggered by two quite distinct mechanisms.

(1) The fact that the space-time is dynamic leads to a mixing of positive and negative frequency states, and hence a change in the number of particles in any given mode.<sup>20</sup> This may be interpreted as a manifestation of the type of “parametric amplification” well known in quantum optics.<sup>6</sup>

(2) Couplings between modes will lead to the creation or destruction of particles whether or not the space-time is dynamic.

The entropy  $S$  defined in this paper says absolutely nothing about the effects of “parametric amplification.” The particle creation associated with “mode–mode coupling” does, however, correspond to an increase in the entropy. If one starts with an initially uncorrelated state at some time  $t_0$ , the generated correlations reflected in the  $\mu_t$  of Eq. (2.24) necessarily induce a net creation of particles at time  $t_0 + \Delta t$ .

### 3. A Mixmaster Universe (Ref. 21)

These models correspond to a class of cosmological solutions to the Einstein equation which admit a foliation into a family of spacelike hypersurfaces, each of which is homogeneous but characterized by a nontrivial and dynamic spatial curvature. The existence of the foliation leads to a natural 3 + 1 splitting, and the theory of harmonic analysis leads to a natural selection of spatial eigenfunctions in terms of which to expand, so that the choice of mode decomposition is very nearly canonical.

The crucial new effect in such a universe is the fact that even a linear field theory leads to a collection of *coupled* oscillators! The ordinary flat space equations

$$\frac{d^2 q_A}{dt^2} + \omega_A^2 q_A = 0, \quad \omega_A^2 = k_A^2 + m^2, \quad (3.11)$$

for a Klein-Gordon field are now replaced by more complicated equations<sup>22</sup>

$$\frac{d^2 q_A}{dt^2} + \sum_B b_{AB} \frac{dq_B}{dt} + \sum_B c_{AB} q_B = 0, \quad (3.12)$$

where  $b_{AB}$  and  $c_{AB}$  are time-dependent coefficients which reflect the effects of the dynamical curvature. For the special case of a conformally static Friedmann Universe, the matrices  $b_{AB}$  and  $c_{AB}$  are diagonal and the modes decouple,<sup>19</sup> but, in general, one infers a linear coupling between the modes which induces a time dependent  $S(t)$ . Formally, this may be understood by observing that, in the presence of nontrivial dynamics, the appropriate choice of spatial eigenfunctions will in fact exhibit a parametric time dependence.

This illustrates the important fact that, in a realistic dynamical space-time, even the simplest source-free linear field theory can lead to a nontrivial generation of entropy. The notion of entropy becomes, if anything, even more important in a dynamical space-time manifold.

## IV. THE DERIVATION OF A SUBDYNAMICS AND AN $H$ -THEOREM INEQUALITY

### A. The derivation of a subdynamics (Ref. 23)

The object of this section is to derive an exact closed equation for the evolution of  $S(t)$  involving only the reduced one-oscillator  $g(A)$ 's, and to use that exact equation to prove a short time  $H$ -theorem for a state evidencing no initial correlations. Most of the work entails the construction of a subdynamics for the  $g(A)$ 's and a derivation of an exact equation for their evolution which involves no explicit reference to the higher-order correlations. Once such an equation has been obtained, it is straightforward to evaluate  $dS/dt$ , and it is completely trivial to show that an initially uncorrelated state leads to an initial entropy generation.

The starting point for the analysis is a collection of "objects" that interact in a fashion described by the classical equations of motion

$$\frac{dq_A}{dt} = \alpha_A(q_A, \pi_A; t) + \xi_A(\{q_B, \pi_B\}; t)$$

and

$$\frac{d\pi_A}{dt} = \gamma_A(q_A, \pi_A; t) + \eta_A(\{q_B, \pi_B\}; t), \quad A = 1, 2, 3, \dots, \quad (4.1)$$

where  $\alpha$ ,  $\gamma$ ,  $\xi$ , and  $\eta$  are arbitrary time-dependent functions of the  $q$ 's and  $\pi$ 's, required only to satisfy the conditions

$$\sum_A \left( \frac{\partial \alpha_A}{\partial q_A} + \frac{\partial \gamma_A}{\partial \pi_A} \right) = 0$$

and

$$\sum_A \left( \frac{\partial \xi_A}{\partial q_A} + \frac{\partial \eta_A}{\partial \pi_A} \right) = 0. \quad (4.2)$$

These equations could have been derived from a Hamiltonian—as would be expected for most realistic systems (recall, however, the caveat in Sec. II B!)—but that is by no means necessary. The conditions (4.2) ensure an analog of the standard Liouville theorem, i.e., conservation of phase, so that the fundamental Liouville equation is truly linear. The construction of a subdynamics *per se* does not require the imposition of these conditions, but, if they are not imposed, one cannot infer a one-to-one correspondence between changes in the entropy (2.22) and the evolution of correlations in the system.

Note also that the equations of motion (4.1) are significantly more general than one might reasonably expect for any realistic field theory. This generality should, however, serve to emphasize that the proof presented here is of very broad generality, including as a special case a collection of  $N$  point masses interacting via two-body forces.

The fundamental object for the statistical description of such a classical system is a many-object distribution function  $\mu$  defined in an infinite-dimensional phase space. The details of this phase space are not particularly relevant here; and, as discussed by Kandrup,<sup>19</sup> one can assume simply that it is constructed in the obvious way as the direct product of an infinite number of two-dimensional, flat, one-object phase spaces. The only important point is that there exist a notion of probability conservation, so that

$$\frac{\partial \mu}{\partial t} + \sum_A \left( \frac{dq_A}{dt} \mu \right) + \sum_A \left( \frac{d\pi_A}{dt} \mu \right) = 0. \quad (4.3)$$

Given this expression of conservation of probability, it follows immediately from Eqs. (4.1) and (4.2) that the evolution of  $\mu$  is governed by the linear Liouville equation

$$\begin{aligned} \frac{\partial \mu}{\partial t} + \sum_A (\alpha_A + \xi_A) \frac{\partial \mu}{\partial q_A} + \sum_A (\gamma_A + \eta_A) \frac{\partial \mu}{\partial \pi_A} \\ \equiv \frac{\partial \mu}{\partial t} + Lu = 0. \end{aligned} \quad (4.4)$$

It is useful to view this  $L$  as a sum of free and interaction Liouvillians, denoted, respectively,  $L^0$  and  $L^I$ ,

$$L^0 \equiv \sum_A \alpha_A \frac{\partial}{\partial q_A} + \sum_A \gamma_A \frac{\partial}{\partial \pi_A} \equiv \sum_A L^0_A \quad (4.5)$$

and

$$L^I \equiv \sum_A \xi_A \frac{\partial}{\partial q_A} + \sum_A \eta_A \frac{\partial}{\partial \pi_A} \equiv \sum_A L^I_A. \quad (4.6)$$

The obvious point is that it is the  $L^I_A$ 's which induce evolu-

ing correlations amongst the  $q_A$ 's. If  $L^I_A \equiv 0$  for all the  $A$ 's, the system is noninteracting and  $dS/dt \equiv 0$ .

Two important identities should be noted.

(1) Let  $\text{Tr}_A$  denote a trace over the degrees of freedom of the  $A$ th object, i.e., an integration  $\int dQ_A d\pi_A$ . It then follows trivially that the operators  $L^0_A$  and  $L^I_A$  are antisymmetric in the sense that, for any two functions  $\psi$  and  $\chi$ ,

$$\text{Tr}_A \psi L^0_A \chi = -\text{Tr}_A \chi L^0_A \psi$$

and (4.7)

$$\text{Tr}_A \psi L^I_A \chi = -\text{Tr}_A \chi L^I_A \psi.$$

This implies, in particular, that

$$\text{Tr}_A L^0_A \chi \equiv 0 \equiv \text{Tr}_A L^I_A \chi. \quad (4.8)$$

(2) Similarly, it follows that, for any function  $\psi$  of the  $g(A)$ 's,

$$L^0_A \psi[g(A)] = \frac{d\psi}{dg(A)} L^0_A g(A)$$

and (4.9)

$$L^I_A \psi[g(A)] = \frac{d\psi}{dg(A)} L^I_A g(A).$$

The operators  $L^0_A$  and  $L^I_A$  satisfy the Leibnitz rule. It is this identity which would fail in the absence of conservation of phase.

The generalization to a quantum description introduces no serious mathematical complications. If the equations of motion (4.1) derive from a Hamiltonian, one can simply implement the standard prescription of canonical quantization. In this event, the classical Liouvillian  $L$  is nothing other than the Poisson bracket  $\{H, \dots\}$ , and the corresponding quantum  $L$  is a commutator  $[H, \dots]_-$ . If a Hamiltonian  $H$  does not exist, one requires some other rule to generate the quantum analog of the equations of motion. In either case, the distribution function  $\mu$  is simply reinterpreted as a density matrix and the classical equation (4.4) as a quantum Liouville equation

$$\frac{\partial \mu}{\partial t} + L\mu = 0, \quad (4.10)$$

where  $L$  is an abstract operator. The only important requirements are (i) that  $L$  remain linear, (ii) that  $L$  still admit a decomposition into linear contributions  $L^0_A$  and  $L^I_A$ , and (iii) that the quantum  $L^0_A$  and  $L^I_A$  still satisfy the identities (4.7) and (4.9), where  $\text{Tr}_A$  is interpreted now as an abstract trace. These requirements are, e.g., guaranteed in the framework of canonical quantization. The analysis henceforth in this section will be formulated abstractly in a fashion applicable to either a classical or a quantum description.

The idea underlying a subdynamics is that one can view the total distribution function or density matrix  $\mu$  as being a sum of "relevant" and "irrelevant" contributions,  $\mu_R$  and  $\mu_I$ , and that one can extract from the full Liouville equation an exact equation for the evolution of  $\mu_R$  which contains no explicit reference to  $\mu_I$ . The desired entropy  $S(t)$  is to be

constructed in terms of the reduced one-object  $g(A)$ 's defined by the prescription

$$g(A) \equiv \prod_{B \neq A} \text{Tr}_B \mu, \quad (4.11)$$

and, consequently, it is natural to introduce the decomposition (2.24), setting

$$\mu = \mu_R + \mu_I,$$

where

$$\mu_R \equiv \prod_A g(A). \quad (4.12)$$

The object now is to implement this decomposition in a canonical fashion by means of projection operators.<sup>24</sup> Specifically, what is required is a linear operator  $P(t)$  defined to satisfy the three requirements

$$P(t)\mu(t) = \mu_R(t), \quad (4.13)$$

$$P(t_2)P(t_1) = P(t_2) \quad \text{for } t_2 \geq t_1, \quad (4.14)$$

and

$$\left[ P(t), \frac{\partial}{\partial t} \right] \mu(t) = 0. \quad (4.15)$$

The first of these ensures that  $P$  projects out the desired  $\mu_R$ . The second implies that  $P$  is in fact idempotent. The third guarantees that the operations of projection and time evolution commute, at least when restricted to the fundamental  $\mu(t)$ .<sup>25</sup>

Given such a  $P$ , it is straightforward to obtain the desired closed equation for  $\mu_R$ . By acting upon the fundamental Liouville equation with the operators  $P$  and  $(1 - P)$ , one is led to the coupled system

$$\frac{\partial \mu_R}{\partial t} + PL\mu_R = -PL\mu_I \quad (4.16)$$

and

$$\frac{\partial \mu_I}{\partial t} + (1 - P)L\mu_I = -(1 - P)L\mu_R. \quad (4.17)$$

It is then trivial to write down a formal solution to Eq. (4.17), yielding  $\mu_I(t)$  in terms of the retarded  $\mu_R(t - \tau)$ ; and, by substituting that solution back into Eq. (4.16), one obtains the desired closed equation for  $\mu_R$ . Thus, in terms of an initial condition  $\mu_I(t_0)$ , one concludes that

$$\begin{aligned} & \frac{\partial \mu_R(t)}{\partial t} + P(t)L(t)\mu_R(t) \\ &= -P(t)L(t)\mathcal{G}(t, t_0)\mu_I(t_0) \\ &+ \int_0^{t-t_0} d\tau P(t)L(t)\mathcal{G}(t, t-\tau) \\ &\times [1 - P(t-\tau)]L(t-\tau)\mu_R(t-\tau), \end{aligned} \quad (4.18)$$

where

$$\mathcal{G}(t_2, t_1) = T \exp \left\{ - \int_{t_1}^{t_2} d\tau [1 - P(\tau)]L(\tau) \right\} \quad (4.19)$$

and  $T$  denotes a time ordering operator.

By integrating over the degrees of freedom of all but one of the objects in the system, one then obtains equations for each  $\partial g(A)/\partial t$  in terms of  $g(A)$  and the remaining  $g(B)$ 's.

And, given these equations, it is straightforward to evaluate  $dS/dt$ . Alternatively, one could evaluate that quantity directly by observing that the field entropy (2.22) can also be written in the form

$$S = -\text{Tr} \mu_R \log \mu_R, \quad (4.20)$$

where  $\text{Tr}$  denotes a trace over all the degrees of freedom. It remains, however, to prove that the desired  $P$  actually exists and to consider in greater detail its action on such quantities of interest as  $L\mu_R$ . These issues could be addressed at purely formal level by endowing  $P$  with additional properties and then constructing a rigorous existence proof.<sup>26</sup>

It is, however, instructive instead to proceed constructively by exhibiting a specific  $P$  that satisfies Eqs. (4.13)–(4.15) and then evaluating directly its effects upon the relevant quantities. One convenient choice, considered in the past,<sup>5,7</sup> is obtained by generalizing the approach of Willis and Picard<sup>2</sup> for a system of  $N$  interacting point masses. Specifically, one is instructed to view the infinite collection of objects as the  $N \rightarrow \infty$  limit of a finite system, and then define

$$P(t) = \lim_{N \rightarrow \infty} \sum_{A=1}^N \prod_{B \neq A} g(B) \text{Tr}_B - \lim_{N \rightarrow \infty} (N-1) \prod_{B=1}^N g(B) \text{Tr}_B. \quad (4.21)$$

It is simple to verify that this  $P$  satisfies Eq. (4.13), and a proof of Eq. (4.15) is also not hard. The proof that  $P$  is in fact idempotent is somewhat more difficult, but can be constructed straightforwardly by observing that any function  $\psi(q_A, \pi_A, \dots; t)$  can be approximated with arbitrary precision as a sum of terms

$$\psi_i \equiv \prod_A \psi_i^A(q_A, \pi_A). \quad (4.22)$$

By exploiting the explicit form of  $P$  and the first identity (4.8), it is easy to see that, when acting upon  $\mu_R$ , the operators  $P$  and  $L^0$  commute, so that

$$PL^0 \mu_R = L^0 P \mu_R = L^0 \mu_R. \quad (4.23)$$

Similarly, one can prove that, when acting on  $\mu_R$ ,  $PL^1$  serves to define an “average” value. Thus, one verifies explicitly that

$$PL^1 \mu_R = \{L^1_A\} \mu_R, \quad (4.24)$$

where  $\{L^1_A\}$  denotes an “average interaction Liouvillian” involving only the variables  $q_A$  and  $\pi_A$ . The form of this average will of course depend upon the form of  $L^1_A$ , and, particularly, whether the fundamental interactions are two-, three-, or higher-body. If, for example,  $L^1_A$  can be written as a sum of contributions  $L^{B_1 \dots B_n}_A$  involving the interaction of  $A$  with objects  $B_1, \dots, B_n$ , each of which satisfies Eq. (4.7), one can conclude that

$$\{L^1_A\} \mu_R = \sum_{B_1} \dots \sum_{B_n} \{L^{B_1 \dots B_n}_A\} \mu_R, \quad (4.25)$$

where, explicitly,

$$\{L^{B_1 \dots B_n}_A\} = \text{Tr}_{B_1} \dots \text{Tr}_{B_n} L^{B_1 \dots B_n}_A g(B_1) \dots g(B_n). \quad (4.26)$$

Note in particular that  $\{L^1_A\}$ , like  $L^1_A$ , satisfies the identities (4.7) and (4.9).

Now introduce the operator

$$\Delta_A \equiv L^1_A - \{L^1_A\}. \quad (4.27)$$

One concludes then that, in the absence of initial correlations, so that  $\mu_I(t_0) \equiv 0$ ,

$$\begin{aligned} \frac{\partial \mu_R(t)}{\partial t} + \sum_A L^0_A(t) \mu_R(t) + \sum_A \{L^1_A(t)\} \mu_R(t) \\ = \int_0^{t-\tau} d\tau P(t) L(t) \mathcal{G}(t, t-\tau) \\ \times \sum_B \Delta_B(t-\tau) \mu_R(t-\tau). \end{aligned} \quad (4.28)$$

This implies that each  $g(A)$  will satisfy an equation of the form

$$\begin{aligned} \frac{\partial g(A, t)}{\partial t} + L^0_A(t) g(A, t) + \{L^1_A(t)\} g(A, t) \\ = \mathcal{S}[g(A, t)] \\ \equiv \prod_{C \neq A} \text{Tr}_C \int_0^{t-\tau} d\tau P(t) L(t) \mathcal{G}(t, t-\tau) \\ \times \sum_B \Delta_B(t-\tau) \prod_D g(D, t-\tau). \end{aligned} \quad (4.29)$$

If  $\mathcal{S}[g]$  were assumed to vanish identically, one would be reduced to a type of “mean field” description appropriate in the limit that  $\mu \simeq \prod_A g(A)$ . It is the nontrivial  $\mathcal{S}$  which leads to a nonconserved entropy.

Equation (4.29) simplifies further. Thus, by virtue of Eq. (4.7) and the definition of  $P$ , it follows that, for any function  $\psi$ ,

$$\prod_{B \neq A} \text{Tr}_B (1-P) \psi = 0, \quad (4.30)$$

where the trace extends over the degrees of freedom of all but one of the objects in the system. Given, moreover, that  $\mathcal{G}(t, t-\tau)$  is constructed as a sum of terms involving  $(1-P)$ , it follows that

$$\prod_{B \neq A} \text{Tr}_B \mathcal{G}(t, t-\tau) \psi = 0. \quad (4.31)$$

This means that, in Eq. (4.29), the operator  $P(t)$  can be omitted. Equation (4.7) implies further that  $L(t)$  may be replaced by  $L_A(t)$ , and it is also found easy to see that, since  $L^0_A$  is independent of the variables  $B \neq A$ , its contribution vanishes identically. Finally, note that Eq. (4.31) implies also that  $L^1_A$  can be replaced by  $\Delta_A$ . One concludes, therefore, that

$$\begin{aligned} \mathcal{S}[g] = \prod_{C \neq A} \text{Tr}_C \int_0^{t-\tau} dt \Delta_A(t) \mathcal{G}(t, t-\tau) \\ \times \sum_B \Delta_B(t-\tau) \mu_R(t-\tau). \end{aligned} \quad (4.32)$$

## B. The time dependence of the entropy

Given Eq. (4.32), it becomes straightforward to evaluate  $dS/dt$ . It follows at once that

$$\frac{dS}{dt} = - \sum_A \text{Tr}_A [1 + \log g(A)] \frac{\partial g(A)}{\partial t}. \quad (4.33)$$

Equations (4.9) and (4.29) then guarantee that, in the absence of initial correlations,

$$\frac{dS}{dt} = \sum_A \text{Tr} [L^0_A g(A) \log g(A) + \{L^I_A\} g(A) \log g(A)] - \sum_A \text{Tr} [1 + \log g(A)] \mathcal{S}[g(A)]. \quad (4.34)$$

Equation (4.7), and the analog satisfied by  $\{L^I_A\}$ , imply further that the first two terms in Eq. (4.34) vanish identically, so that

$$\begin{aligned} \frac{dS(t)}{dt} &= -\text{Tr} \int_0^{t-\tau} d\tau \sum_A [1 + \log g(A)] \Delta_A(t) \mathcal{S}(t, t-\tau) \\ &\quad \times \sum_B \Delta_B(t-\tau) \mu_R(t-\tau). \end{aligned} \quad (4.35)$$

The antisymmetry of  $\Delta_A$  means that

$$\begin{aligned} \frac{dS(t)}{dt} &= \text{Tr} \int_0^{t-\tau} dt \left[ \sum_A \Delta_A(t) \log g(A, t) \right] \mathcal{S}(t, t-\tau) \\ &\quad \times \sum_B \Delta_B(t-\tau) \mu_R(t-\tau), \end{aligned} \quad (4.36)$$

and the fact that it satisfies the Leibnitz rule guarantees further that

$$\Delta_A \log g(A) = g^{-1}(A) \Delta_A g(A) = \mu_R^{-1} \Delta_A \mu_R, \quad (4.37)$$

so that, finally, one concludes that

$$\frac{dS(t)}{dt} = \text{Tr} \int_0^{t-\tau} d\tau \mu_R^{-1}(t) \xi(t) \mathcal{S}(t, t-\tau) \xi(t-\tau), \quad (4.38)$$

where

$$\xi = \sum_A \Delta_A \mu_R. \quad (4.39)$$

Given Eq. (4.38), the “short-time”  $H$ -theorem follows immediately. In the limit that  $\tau \rightarrow 0$ ,  $\mathcal{S}(t, t-\tau) \rightarrow 1$ , so that for small intervals  $\Delta t$ ,

$$\frac{dS(t_0 + \Delta t)}{dt} = \text{Tr} \mu_R^{-1}(t) |\xi(t)|^2 \Delta t > 0. \quad (4.40)$$

An absence of initial correlations guarantees an initial increase in entropy. Note further that this quantity is quadratic in the interaction Liouvillian, and, as such, if the typical interaction is characterized by a coupling constant  $\lambda$ , the spontaneous entropy generation induced by  $\mu_I$  will, at least for short times, scale as  $\lambda^2$ .

For the case of a classical system, where  $\mu$  is realized as a distribution function, the final inequality in (4.40) follows immediately from the fact that  $\mu$  can never be negative. In the quantum case, however, the density matrix  $\mu$  need not be positive definite, so that the positivity of  $dS(t_0 + \Delta t)/dt$  is less transparent. The basic conclusion does, however, follow immediately by working in a representation in which  $\mu_R$  is diagonal. That  $\mu_R$  may be diagonalized follows in turn from the fact that the density matrix must be symmetric (i.e., Hermitian), and, once  $\mu_R$  has been diagonalized, the inequality obtains trivially from the fact that the diagonal ele-

ments of  $\mu_R$ , and hence  $\mu_R^{-1}$ , must be non-negative.<sup>27</sup>

If, initially,  $\mu_I(t)$  does not vanish identically, the expression (4.38) for  $dS/dt$  will contain an additional contribution of the form

$$\text{Tr} [1 + \log \mu_R(t)] P(t) L(t) \mathcal{S}(t, t_0) \mu_I(t_0). \quad (4.41)$$

By exploiting the identity (4.30) and the fact that  $L$  is anti-symmetric and satisfies the Leibnitz rule, one concludes then that, quite generally,

$$\begin{aligned} \frac{dS}{dt} &= \text{Tr} \int_0^{t-\tau} d\tau \mu_R^{-1}(t) \xi(t) \mathcal{S}(t, t-\tau) \xi(t-\tau) \\ &\quad - \text{Tr} \mu_R^{-1}(t) \sigma(t) \mathcal{S}(t, t_0) \mu_I(t_0), \end{aligned} \quad (4.42)$$

where

$$\sigma = L \mu_R. \quad (4.43)$$

This means that, in the limit that  $t = t_0$ , the derivative  $dS/dt$  need not vanish identically. Rather, one sees that

$$\frac{dS(t_0)}{dt} = -\text{Tr} \mu_R^{-1}(t_0) \sigma(t_0) \mu_I(t_0), \quad (4.44)$$

which, depending upon the form of the initial  $\mu_I$  and  $\mu_R$ , could be either positive or negative! By judicious choice of initial conditions, one can, at least in principle, induce an initial decrease in the entropy.

This might, naively, seem a disturbing result, but further reflection demonstrates that such a possibility should have been expected. If one is allowed to choose the initial conditions arbitrarily, one *should* be able to stimulate either an initial increase or an initial decrease in  $S$ . As discussed in Ref. 6, the situation is qualitatively similar to that arising in the analysis of particle creation in a time-dependent electromagnetic or gravitational field. Here again the choice of an especially simple state, such as the vacuum or an eigenstate of the number operator, leads to a spontaneous generation of particles; but an allowance for more generic initial conditions, in this case reflecting nontrivial “phase” information, leads to the possibility of stimulated creation or destruction of quanta.

Given these observations, it becomes clear that, as was intimated in the Introduction, an  $H$ -theorem could only be expected to hold in full generality at late times, after which any nontrivial initial correlations have either “decayed” or else have been completely “dwarfed” by the systematic correlations induced spontaneously by the evolving dynamics.

## V. ENTROPY GENERATION AND PARTICLE CREATION

### A. A simple example

The purpose of this section is to provide some general insights into the phenomenon of particle creation induced by mode-mode coupling and by parametric amplification, and, specifically, to demonstrate that the former mechanism is connected intimately with the generation of entropy. Specifically, it is also possible to prove a short-term “ $H$ -theorem” for particle creation which states that correlations induced from an initially uncorrelated state lead to an overall enhancement in the rate at which quanta are created. This theorem, and other related phenomena, will be discussed first for a system described by a specific model Hamiltonian

$H(t)$ , and thereafter, a general proof of the particle creation  $H$ -theorem will be presented for a more general system described by an arbitrary interaction Hamiltonian  $H^I$  constructed from the "coordinates"  $q_A$ .

The specific model Hamiltonian is chosen to take the form

$$H = \sum_A \frac{1}{2} (\pi_A^2 + \omega_A^2 q_A^2) + \sum_{B \neq A} \frac{1}{2} c_{AB} q_A q_B \equiv \sum_A H^0_A + \sum_{B \neq A} H^I_{AB}, \quad (5.1)$$

where  $\omega_A^2 > 0$  and  $c_{AB}$  are arbitrary real functions of time. This leads to classical equations of motion

$$\frac{dq_A}{dt} = \pi_A$$

and

$$\frac{d\pi_A}{dt} = -\omega_A^2 q_A - \sum_B c_{AB} q_B, \quad (5.2)$$

or, equivalently,

$$\frac{d^2 q_A}{dt^2} + \omega_A^2 q_A + \sum_B c_{AB} q_B = 0. \quad (5.3)$$

This  $H$  is reasonable for several reasons.

(1) If  $c_{AB} \equiv 0$  and the  $\omega_A$ 's are assumed time-independent constants, one recovers the massless Klein-Gordon equation in flat space. Thus, if one interprets  $A$  as labeling a  $\mathbf{k}$  vector and sets  $\omega_A^2 = k_A^2$ , the decomposition (2.8) leads immediately to Eq. (2.14).

(2) If  $c_{AB} \equiv 0$  but the  $\omega_A$ 's are allowed to be functions of time, Eq. (5.1) includes as a special example the Klein-Gordon equation in a  $k = 0$  Friedmann cosmology. Thus the identification

$$\omega_k^2 = k^2 - \ddot{\Omega}/\Omega, \quad (5.4)$$

where  $\Omega$  is the conformal factor and an overdot is a conformal time derivative, together with the decomposition

$$\Phi(\mathbf{x}, t) = \Omega^{-1} \sum_k q_k(t) \exp(-i\mathbf{k} \cdot \mathbf{x}), \quad (5.5)$$

leads to a field equation

$$-\frac{\partial^2 \Phi}{\partial t^2} - \frac{2\dot{\Omega}}{\Omega} \frac{\partial \Phi}{\partial t} + \delta^{ab} \frac{\partial}{\partial x^a} \frac{\partial}{\partial x^b} \Phi = 0, \quad (5.6)$$

which is nothing other than the Klein-Gordon equation appropriate for a conformally flat space-time with  $ds^2 = \Omega^2(t) \eta_{\mu\nu} dx^\mu dx^\nu$ .

(3) In the case that  $c_{AB}$  and  $\omega_A$  are both nonvanishing and exhibit a nontrivial time dependence, the Hamiltonian (5.1) mocks quite reasonably the behavior of a scalar field in a Mixmaster Universe. One might like also to allow for more complicated nonlinear couplings, as induced, e.g., by a  $\lambda \Phi^p$  field theory, but, as will be seen eventually, more complicated contributions involving  $c_{ABC}$ ,  $c_{ABCD}$ , etc. do not alter the existence of the basic  $H$ -theorem for particle creation.

The passage to a quantum theory can now be effected trivially by canonical quantization, the Poisson bracket structure being replaced by a canonical commutation relation

$$[q_A, \pi_B]_- = i\delta_{AB}. \quad (5.7)$$

The quantum system is then characterized by a many-oscillator density matrix  $\mu$ , the evolution of which is governed by the quantum Liouville equation

$$\frac{\partial \mu}{\partial t} = -L\mu = -i[H, \mu]_-. \quad (5.8)$$

Note for future reference that the free Hamiltonian  $L^0_A$  satisfies

$$L^0_A \xi = i[H^0_A, \xi]_-, \quad (5.9)$$

and that the average interaction Liouvillian defined in Sec. IV takes the form

$$\{L^I_A\} \xi = i[\{H^I_A\}, \xi]_-, \quad (5.10)$$

where

$$\{H^I_A\} = \sum_{B \neq A} \text{Tr}_B H^I_{AB} g(B) = \sum_B \frac{1}{2} c_{AB} \langle q_B \rangle q_A, \quad (5.11)$$

the angular brackets  $\langle \rangle$  denoting an expectation value with respect to the full many-oscillator  $\mu$ . This implies further that the  $\Delta_A$  of Sec. IV can be realized in the form

$$\sum_A \Delta_A \xi = \sum_{B \neq A} \sum i[h_{AB}, \xi]_-, \quad (5.12)$$

where

$$h_{AB} \equiv H^I_{AB} - \text{Tr}_B H^I_{AB} g(B) - \text{Tr}_A H^I_{AB} g(A). \quad (5.13)$$

At this stage, it is customary to define creation and annihilation operators

$$a_A = (\sqrt{2\omega_A})^{-1} (\omega_A q_A + i\pi_A) \quad \text{and} \quad (5.14)$$

$$a^\dagger_A = (\sqrt{2\omega_A})^{-1} (\omega_A q_A - i\pi_A),$$

in terms of which the free Hamiltonian

$$\begin{aligned} H^0 &= \sum_A (\omega_A^2 q_A^2 + \pi_A^2) \\ &= \sum_A \frac{1}{2} \omega_A (a^\dagger_A a_A + a_A a^\dagger_A) \\ &= \sum_A \omega_A \left( a^\dagger_A a_A + \frac{1}{2} \right). \end{aligned} \quad (5.15)$$

In the absence of any interaction  $H^I$ , the operator

$$N_A \equiv a^\dagger_A a_A \quad (5.16)$$

would admit to an unambiguous interpretation as "the number of quanta in mode  $A$ ." In the presence of nontrivial interactions, this interpretation is less clear-cut, but, at least in the limit of weak couplings, where, typically,  $|c| \ll \omega^2$ , it is still conventional to think of  $N_A$  as reflecting some measure of "particle number." Whether this generalization is reasonable is, perhaps, unclear, but for the remainder of this section, changes in the normal-ordered statistical expectation value

$$\begin{aligned} \langle N_A \rangle &\equiv \text{Tr} a^\dagger_A a_A \mu = \text{Tr} (2\omega_A)^{-1} (\pi_A^2 + \omega_A^2 q_A^2) \mu \\ &= (\omega_A)^{-1} \langle H^0_A \rangle \end{aligned} \quad (5.17)$$

will be spoken of as corresponding to "changes in particle number." The statements in the Introduction and in the remainder of this section regarding the net creation of quanta refer specifically to the time dependence of  $\langle N_A \rangle$ .

The object now is to evaluate  $\partial \langle N_A \rangle / \partial t$ , allowing for contributions arising both from a changing  $\mu$  and from the time dependence of the  $\omega_A$ 's. Thus, one sees immediately that

$$\frac{\partial \langle N_A \rangle}{\partial t} = \text{Tr} \frac{1}{\omega_A} H^0_A \frac{\partial \mu}{\partial t} + \text{Tr} \mu \frac{\partial}{\partial t} \left( \frac{1}{\omega_A} H^0_A \right). \quad (5.18)$$

It is then trivial to verify that

$$\begin{aligned} \text{Tr} H^0_A \frac{\partial \mu}{\partial t} &= -i \text{Tr} H^0_A [H, \mu]_- \\ &= -i \text{Tr} \mu [H^0_A, H]_- \\ &= -i \text{Tr} \mu [H^0_A, H^I]_-. \end{aligned} \quad (5.19)$$

And, therefore, one concludes that

$$\begin{aligned} \frac{\partial \langle N_A \rangle}{\partial t} &= -\frac{\dot{\omega}_A}{2\omega_A^2} \text{Tr} \mu (\pi_A^2 - \omega_A^2 q_A^2) \\ &\quad - \sum_{B \neq A} \frac{c_{AB}}{\omega_A} \text{Tr} \mu \pi_A q_B. \end{aligned} \quad (5.20)$$

The first term in Eq. (5.20) reduces at once to a partial trace  $\text{Tr}_A$  weighted by the reduced  $g(A)$ . The second entails couplings between modes  $A$  and  $B$ , so that one obtains instead a double trace  $\text{Tr}_A \text{Tr}_B$  weighted by the two-oscillator  $g_2(A, B)$ . Thus, if one writes

$$g_2(A, B) = g(A)g(B) + \nu(A, B), \quad (5.21)$$

where  $\nu(A, B)$  denotes the pair correlation function, one concludes that

$$\begin{aligned} \frac{\partial \langle N_A \rangle}{\partial t} &= -\frac{\dot{\omega}_A}{2\omega_A^2} \langle \pi_A^2 - \omega_A^2 q_A^2 \rangle \\ &\quad - \sum_{B \neq A} \frac{c_{AB}}{\omega_A} \langle \pi_A \rangle \langle q_B \rangle + \frac{\partial \langle N_A \rangle^c}{\partial t}, \end{aligned} \quad (5.22)$$

where

$$\frac{\partial \langle N_A \rangle^c}{\partial t} \equiv - \sum_{B \neq A} \text{Tr}_A \text{Tr}_B \nu(A, B) \pi_A q_B \frac{c_{AB}}{\omega_A}. \quad (5.23)$$

Alternatively, in terms of the  $\mu_I$  defined in Sec. IV,

$$\frac{\partial \langle N_A \rangle^c}{\partial t} = - \sum_{B \neq A} \text{Tr} \mu_I \pi_A q_B \frac{c_{AB}}{\omega_A}. \quad (5.24)$$

The first two terms in Eq. (5.22) can be interpreted as representing a mean field particle creation  $\partial \langle N_A \rangle^m / \partial t$  which will be present even if, in the spirit of the collisionless Boltzmann equation, one were to pretend that  $\mu_I \approx 0$ . The final term represents instead a "correlational" particle creation induced specifically by the correlations amongst the modes. The field entropy (5.22) has no direct connection with the mean field contributions, but it does connect intimately with  $\partial \langle N_A \rangle^c / \partial t$ . Specifically, the same arguments which led to the short-time  $H$ -theorem (5.40) imply also that  $\partial \langle N_A(t_0 + \Delta t) \rangle^c / \partial t > 0$ .

## B. Correlational particle creation for the simple example

Turn, therefore, to an explicit evaluation of this quantity for an initially uncorrelated state. The analysis of Sec. IV indicates that, in the absence of any initial correlations,

$$\mu_I(t) = \int_0^{t-t_0} d\tau \mathcal{G}(t, t-\tau) \sum_A \Delta_A(t-\tau) \mu_R(t-\tau). \quad (5.25)$$

And, thus, in terms of the  $\Delta_A$  of Eq. (4.27), one concludes immediately that

$$\begin{aligned} \frac{\partial \langle N_A(t) \rangle^c}{\partial t} &= i \sum_B \text{Tr} \frac{c_{AB}}{\omega_A} \pi_A q_B \int_0^{t-t_0} d\tau \mathcal{G}(t, t-\tau) \\ &\quad \times \sum_{C \neq D} [h_{CD}(t-\tau), \mu_R(t-\tau)]_-. \end{aligned} \quad (5.26)$$

A simple application of the cyclic trace rule

$$\text{Tr} ABC = \text{Tr} CAB \quad (5.27)$$

then implies that

$$\begin{aligned} \frac{\partial \langle N_A(t) \rangle^c}{\partial t} &= -i \sum_B \text{Tr} \frac{c_{AB}}{\omega_A} \int_0^{t-t_0} d\tau \\ &\quad \times \sum_{C \neq D} [h_{CD}(t-\tau), q_B \pi_A \mathcal{G}(t, t-\tau)]_- \\ &\quad \times \mu_R(t-\tau). \end{aligned} \quad (5.28)$$

For small  $\tau$ ,  $\mathcal{G}(t, t-\tau) \rightarrow 1$  and all the quantities in Eq. (5.28) can be approximated by their values at time  $t_0$ , so that

$$\begin{aligned} \frac{\partial \langle N_A(t_0 + \Delta t) \rangle^c}{\partial t} &= -i \sum_B \text{Tr} \frac{c_{AB}}{\omega_A} \sum_{C \neq D} [h_{CD}, q_B \pi_A]_- \mu_R(t_0) \Delta t. \end{aligned} \quad (5.29)$$

One then verifies immediately that

$$\begin{aligned} [h_{CD}, q_B \pi_A]_- &= (i/2) q_B (c_{AC} (q_C - \langle q_C \rangle) \delta_{AD} \\ &\quad + c_{AD} (q_D - \langle q_D \rangle) \delta_{AC}), \end{aligned} \quad (5.30)$$

so that

$$\begin{aligned} \frac{\partial \langle N_A(t_0 + \Delta t) \rangle^c}{\partial t} &= \sum_B \text{Tr} \frac{c_{AB}}{\omega_A} \sum_C q_B c_{AC} (q_C - \langle q_C \rangle) \mu_R(t_0) \Delta t. \end{aligned} \quad (5.31)$$

The trace in Eq. (5.31) implies that the only nonvanishing contributions arise when  $C = B$ , so that

$$\frac{\partial \langle N_A(t_0 + \Delta t) \rangle^c}{\partial t} = \sum_B \frac{(c_{AB})^2}{\omega_A} (\langle q_B^2 \rangle - \langle q_B \rangle^2) \Delta t > 0! \quad (5.32)$$

One obtains a positive particle creation at a rate proportional to the square of the coupling constant  $c$ .

## C. General proof of spontaneous particle creation induced by correlations

The object now is to demonstrate explicitly that an analog of Eq. (5.32) will hold for an arbitrary interaction Ham-



iltonian  $H^I$  constructed from the  $q$ 's. For such a more general  $H^I$ , Eq. (5.24) will be replaced by the relation

$$\frac{\partial \langle N_A \rangle^c}{\partial t} = -i \text{Tr} \frac{1}{\omega_A} [H^0_A, H^I] - \mu_I, \quad (5.33)$$

which can be written in the form

$$\frac{\partial \langle N_A \rangle^c}{\partial t} = -\frac{1}{2\omega_A} \text{Tr} \left( \pi_A \frac{\partial H^I}{\partial q_A} + \frac{\partial H^I}{\partial q_A} \pi_A \right) \mu_I, \quad (5.34)$$

where  $\partial/\partial q_A$  is interpreted as an operator derivative.

The irrelevant contribution  $\mu_I$  now may be written as

$$\mu_I(t) = i \int_0^{t-t_0} d\tau \mathcal{G}(t, t-\tau) [h^I(t-\tau), \mu_R(t-\tau)]_-, \quad (5.35)$$

where  $h^I$  denotes an appropriate "fluctuating" interaction Hamiltonian. For short times  $\Delta t$ , one then infers that

$$\mu_I(t_0 + \Delta t) = i [h^I, \mu_R(t_0)]_- \Delta t, \quad (5.36)$$

so that

$$\frac{\partial \langle N_A(t_0 + \Delta t) \rangle^c}{\partial t} = \frac{1}{\omega_A} \text{Tr} \frac{\partial H^I}{\partial q_A} \frac{\partial h^I}{\partial q_A} \mu_R(t_0) \Delta t. \quad (5.37)$$

Quite generally, the fluctuating  $h^I$  can be written as a difference  $H^I - \{H^I\}$ , where  $\{H^I\}$  denotes an average interaction Hamiltonian. And thus, one sees that the right hand side of Eq. (5.37) reduces to

$$\frac{1}{\omega_A} \left( \left\langle \left( \frac{\partial H^I}{\partial q_A} \right)^2 \right\rangle - \left\langle \frac{\partial H^I}{\partial q_A} \right\rangle^2 \right) \Delta t > 0, \quad (5.38)$$

which establishes the inequality.

As a concrete example, one can consider an interaction  $H^I$  again involving only pairs of oscillators, but now quadratic in the  $q$ 's,

$$H^I_A = \sum_{B \neq A} \frac{1}{2} c_{AB} q_A^2 q_B^2. \quad (5.39)$$

In this case,

$$h^I_A = \sum_{B \neq A} \frac{1}{2} c_{AB} q_A^2 (q_B^2 - \langle q_B^2 \rangle), \quad (5.40)$$

and one concludes that

$$\frac{\partial \langle N_A(t_0 + \Delta t) \rangle^c}{\partial t} = \sum_B \frac{(c_{AB})^2}{\omega_A} \langle q_A^2 \rangle (\langle q_B^4 \rangle - \langle q_B^2 \rangle^2). \quad (5.41)$$

Similarly, one can suppose instead that each oscillator  $A$  interacts linearly with  $n$  other oscillators  $B_1, \dots, B_n$  via a coupling  $c_{AB_1 \dots B_n}$ . The net result then is that

$$\begin{aligned} \frac{\partial \langle N_A(t_0 + \Delta t) \rangle^c}{\partial t} &= \sum_{B_1 \neq \dots \neq B_n} \frac{(c_{AB_1 \dots B_n})^2}{\omega_A} \left\{ \prod_{B_i} \langle q_{B_i}^2 \rangle - \prod_{B_i} \langle q_{B_i} \rangle^2 \right\}. \end{aligned} \quad (5.42)$$

## VI. THE UNIQUENESS OF THE ENTROPY

Another important question remains to be asked. To what extent does the  $S(t)$  defined by Eq. (2.22) constitute the unique field entropy? Even if one accepts the notion that any measure of entropy should be constructed solely from the  $g(A)$ 's, and even if one could prove rigorously a very general  $H$ -theorem for  $S(t)$ , there would remain the issue of whether some other entropy  $Z(t)$  could be shown to satisfy the same properties.

Broadly speaking, given the general philosophy developed in this paper, there are only three mathematical requirements which must be imposed upon such a candidate  $Z(t)$ .

(1) This  $Z(t)$  must be constructed solely from the  $g(A)$ 's and contain no information about the interoscillator correlations.

(2) In the absence of couplings between degrees of freedom, it must follow that  $dZ/dt \equiv 0$ .

(3) At least in some appropriate limit, one must be able to show that  $dZ/dt \geq 0$ .

To the extent, however, that one takes very seriously the notion of a subdynamics introduced in Sec. IV, these demands can be refined somewhat. Thus, in particular, it would be reasonable to demand not simply that  $Z(t)$  be constructed only from the  $g(A)$ 's, but that it be realizable as a functional of the relevant  $\mu_R$  at the same instant of time  $t$ , i.e., that

$$Z(t) = \text{Tr} \Psi[\mu_R(t)], \quad (6.1)$$

where  $\Psi$  is some arbitrary real function of  $\mu_R$ . Similarly, in the same spirit, one might demand further that, in the absence of initial correlations,  $dZ/dt$  be realizable quite generally as the simple trace of some object  $\mathcal{A}$  involving only the relevant  $\mu_R(t-\tau)$ , the Greenian  $\mathcal{G}(t, t-\tau)$ , and the fluctuating Liouvillian  $\Delta \equiv \Sigma_A \Delta_A$ ,

$$\frac{dZ}{dt} = \text{Tr} \mathcal{A}[\mu_R, \Delta, \mathcal{G}]. \quad (6.2)$$

An explicit dependence on the individual  $g(A)$ 's or any functionals thereof would be considered inappropriate. As illustrated by Eq. (4.38), this requirement is satisfied by  $S(t)$ . And finally, it would be reasonable to demand at least that, if  $\mu_I(t_0) = 0$ ,

$$\frac{dZ(t_0 + \Delta t)}{dt} > 0. \quad (6.3)$$

As will be demonstrated below, it is in fact easy to construct  $Z$ 's which will satisfy the criteria (6.1) and (6.3), but it is far more difficult to satisfy the condition (6.2). Thus, for example, Eqs. (6.1) and (6.3) will hold for any

$$Z(t) = -\text{Tr} \mu_R^p(t), \quad (6.4)$$

where  $p$  is a real number greater than unity, whereas Eq. (6.2) will not be satisfied. Indeed, it is possible to prove that, modulo an overall multiplicative factor and the addition of conserved quantities, the field entropy  $S(t)$  is the unique object to satisfy the conditions (6.1)–(6.3).

Consider first the proof that the  $Z(t)$  defined by Eq. (6.4) will satisfy the short-time  $H$ -theorem (6.3). A simple application of the chain rule shows that

$$\frac{dZ}{dt} = -p \sum_A \prod_{B \neq A} \sigma(B) \text{Tr}_A g^{p-1}(A) \frac{\partial g(A)}{\partial t}, \quad (6.5)$$

where

$$\sigma(A) \equiv \text{Tr}_A g^p(A). \quad (6.6)$$

One then verifies immediately that the mean field contributions to  $\partial g(A)/\partial t$  have a vanishing contribution to  $dZ/dt$  by virtue of Eqs. (4.7) and (4.9), so that

$$\frac{dZ}{dt} = -p \sum_A \prod_{B \neq A} \sigma(B) \text{Tr}_A g^{p-1}(A) \mathcal{S}[g(A)]. \quad (6.7)$$

By inserting into Eq. (6.7) the explicit form of  $\mathcal{S}[g(A)]$  given by Eq. (4.32), one sees that

$$\begin{aligned} \frac{dZ}{dt} &= -p \sum_A \prod_{B \neq A} \sigma(B, t) \\ &\times \text{Tr} \int_0^{t-\tau} d\tau g^{p-1}(A, t) \Delta_A(t) \mathcal{G}(t, t-\tau) \\ &\times \sum_C \Delta_C(t-\tau) \mu_R(t-\tau). \end{aligned} \quad (6.8)$$

The facts that  $\Delta_A(t)$  is antisymmetric and that it satisfies the Leibnitz rule then imply that

$$\begin{aligned} \frac{dZ}{dt} &= p(p-1) \sum_A \prod_{B \neq A} \sigma(B, t) \\ &\times \text{Tr} \int_0^{t-\tau} d\tau g^{p-2}(A, t) [\Delta_A(t) g_A(t)] \mathcal{G}(t, t-\tau) \\ &\times \sum_C \Delta_C(t-\tau) \mu_R(t-\tau), \end{aligned} \quad (6.9)$$

or, equivalently, in terms of the quantity

$$\theta_A = \Delta_A \mu_R, \quad (6.10)$$

that

$$\begin{aligned} \frac{dZ}{dt} &= p(p-1) \sum_A \prod_{B \neq A} \sigma(B, t) \\ &\times \text{Tr} \int_0^{t-\tau} d\tau g^{p-1}(A, t) \mu_R^{-1}(t) \theta_A(t) \mathcal{G}(t, t-\tau) \\ &\times \sum_C \theta_C(t-\tau). \end{aligned} \quad (6.11)$$

The presence of the contributions involving  $\sigma(B)$  and  $g^{p-1}(A)$  demonstrate explicitly that this candidate entropy cannot satisfy the criterion (6.2). At this point, however, it is easy to see that Eq. (6.3) will in fact hold. For small  $\Delta t$ , it follows immediately that

$$\begin{aligned} \frac{dZ(t_0 + \Delta t)}{dt} &= p(p-1) \sum_A \prod_{B \neq A} \sigma(B) \\ &\times \text{Tr} g^{p-1}(A) \mu_R^{-1} \theta_A \sum_C \theta_C \Delta t. \end{aligned} \quad (6.12)$$

And thus, since

$$\text{Tr}_A \theta_A(t) \equiv 0, \quad (6.13)$$

one can conclude that

$$\begin{aligned} \frac{dZ(t_0 + \Delta t)}{dt} &= p(p-1) \sum_A \prod_{B \neq A} \sigma(B, t_0) \\ &\times \text{Tr} g^{p-1}(A, t_0) \mu_R^{-1}(t_0) |\theta_A(t_0)|^2 \Delta t. \end{aligned} \quad (6.14)$$

For the special cases  $p=0$  and  $p=1$ ,  $dZ/dt$  vanishes identically, as must of course be the case. However, for  $p>1$ , one concludes instead that  $dZ(t_0 + \Delta t)/dt > 0$ .

If the system under consideration were a collection of  $N$  identical point masses, rather than an infinite set of different oscillators, Eqs. (6.11) and (6.14) would assume a more palatable form. Thus, if one considers  $N$  particles characterized by a distribution function or density matrix symmetric under particle interchange, one would conclude that

$$Z = \sigma^N, \quad (6.15)$$

where, now,

$$\sigma(t) \equiv \text{Tr}_i f^p(i; t) \quad (6.16)$$

is the same for each particle. This means that Eq. (6.11) will be replaced by a relation of the form

$$\begin{aligned} \frac{d\sigma}{dt} &= p(p-1) \text{Tr} f^{p-1}(i, t) \mu_R^{-1}(t) \theta_i(t) \\ &\times \mathcal{G}(t, t-\tau) \sum_{j=1}^N \theta_j(t-\tau), \end{aligned} \quad (6.17)$$

so that, for short times,

$$\begin{aligned} \frac{d\sigma(t_0 + \Delta t)}{dt} &= p(p-1) \text{Tr} f^{p-1}(i, t_0) \\ &\times \mu_R^{-1}(t_0) |\theta_i(t_0)|^2 \Delta t, \end{aligned} \quad (6.18)$$

where  $\mu_R(t_0)$  is now a product of  $N$  identical quantities.

Turn finally to the general uniqueness proof. For an arbitrary function of the form (6.1), one concludes again that a nontrivial  $dZ/dt$  is induced only by the  $\mathcal{S}[g(A)]$ 's, so that, in the absence of initial correlations,

$$\begin{aligned} \frac{dZ}{dt} &= -\text{Tr} \sum_A \frac{d\Psi}{d\mu_R} g^{-1}(A) \mu_R \prod_{C \neq A} \text{Tr}_C \Delta_A^{(C)} \\ &\times \int_0^{t-\tau} d\tau \mathcal{G}(t, t-\tau) \sum_D \Delta_D(t-\tau) \mu_R(t-\tau), \end{aligned} \quad (6.19)$$

where the superscript  $(C)$  is a reminder that  $\Delta_A^{(C)}$  depends also upon the variables  $C \neq A$ . The antisymmetry of  $\Delta_A$  and the fact that it satisfies the Leibnitz rule then imply that

$$\begin{aligned} \frac{dZ}{dt} &= \text{Tr}_A \prod_{B \neq A} \text{Tr}_B \frac{d^2\Psi}{d\mu_R^2} g^{-1}(A) \mu_R \\ &\times \prod_{C \neq A} \text{Tr}_C [\Delta_A^{(C)} \mu_R(A, B)] \int_0^{t-\tau} d\tau \mathcal{G}(t, t-\tau) \\ &\times \sum_D \Delta_D(t-\tau) \mu_R(t-\tau). \end{aligned} \quad (6.20)$$

Here  $\Delta_A^{(C)} \mu_R(A, B)$  means that  $\Delta_A^{(C)}$ , an operator involving  $A$  and the dummy variables  $C \neq A$ , acts on  $\mu_R$ , viewed as a function of  $A$  and different dummy variables  $B \neq A$ .

If  $dZ/dt$  is to depend only on  $\mu_R$ , rather than upon the individual  $g(B)$ 's, one must demand that

$$\prod_{B \neq A} \text{Tr}_B \frac{d^2 \Psi}{d\mu_R^2} \mu_R(A, B) \Delta_A^{(C)} \mu_R(A, B) \quad (6.21)$$

be independent of the  $g(B)$ 's. Since, however, the only relevant universal property of  $g(B)$  is that it have unit trace, this implies further that  $(d^2 \Psi / d\mu_R^2) \mu_R$  must itself be independent of  $g(B)$ . Given, moreover, that  $d^2 \Psi / d\mu_R^2$  must be a function only of  $\mu_R$ , this can only be true if

$$\frac{d^2 \Psi}{d\mu_R^2} = a \mu_R^{-1}, \quad (6.22)$$

where  $a$  is an arbitrary constant. In this case,

$$\begin{aligned} \frac{dZ(t)}{dt} &= a \text{Tr} \sum_A \Delta_A \mu_R \int_0^{t-\tau} d\tau \mathcal{G}(t, t-\tau) \\ &\quad \times \sum_B \Delta_B(t-\tau) \mu_R(t-\tau) \\ &= a \text{Tr} \int_0^{t-\tau} d\tau \zeta(t) \mathcal{G}(t, t-\tau) \zeta(t-\tau), \end{aligned} \quad (6.23)$$

where  $\zeta$  was defined in Eq. (4.39).

The right-hand side of Eq. (6.23) is, modulo the constant factor  $a$ , of precisely the form satisfied by the field entropy  $S(t)$  of Eq. (2.22). Thus, one concludes (i) that  $Z(t)$  is consistent with the demand (6.2) if and only if  $\Psi[\mu_R]$  satisfies Eq. (6.22), and (ii) that  $Z(t)$  will evidence a monotonic evolution if and only if  $S(t)$  increases monotonically. The most general solution to Eq. (6.22) is of course

$$\Psi[\mu_R] = a \mu_R \log \mu_R + b \mu_R + c, \quad (6.24)$$

where  $b$  and  $c$  are arbitrary constants. And, since

$$\frac{d}{dt} \text{Tr}(b \mu_R + c) \equiv 0, \quad (6.25)$$

one can set  $b = c = 0$  without any loss of generality. The most general entropy  $Z(t)$  consistent with Eq. (6.2) is, therefore,

$$Z(t) = -a \text{Tr} \mu_R \log \mu_R, \quad (6.26)$$

where  $a$  is a constant, and the demand that  $a$  be real and positive guarantees that  $Z(t)$  will satisfy Eq. (6.3).

*Net conclusion:* Modulo the addition of conserved quantities and multiplication by a positive constant, the  $S(t)$  defined by Eq. (2.22) is the unique entropy guaranteed to satisfy a short-time  $H$ -theorem (6.3) and consistent with the general subdynamics in the sense that it satisfies Eq. (6.2). Thus  $S(t)$  may, or may not, satisfy an  $H$ -theorem in some broad generality. It is, however, the only candidate entropy consistent with the demands outlined at the outset of this section.

## VII. DISCUSSION

### A. The late time evolution of $S$

Section I of this paper introduced a notion of particle entropy constructed explicitly as a measure of correlations which, at least in one special limit, is guaranteed to satisfy a short-time  $H$ -theorem inequality. Sections II and III then generalized that definition to a classical or quantum field, emphasizing in particular the question of "physical significance." The entropy so constructed is conserved absolutely

in the absence of couplings between degrees of freedom, and, as such, there can be no entropy generation for a source-free linear field in Minkowski space. If, however, couplings are induced by sources, by nonlinearities, or by a nontrivial dynamical space-time, the entropy need no longer be conserved. And, as was illustrated in Sec. IV, one can prove a short-time  $H$ -theorem which guarantees that an initially uncorrelated state leads (at least) to an initial increase in entropy. Section V demonstrated that there also exists a direct connection between entropy and particle creation, proving that evolving correlations which induce a net entropy generation lead also to an enhancement in the rate at which quanta are produced. Finally, Sec. VI showed that, in a well defined mathematical sense, the entropy  $S(t)$  defined in Secs. II and III is the only measure of entropy consistent with the notion of subdynamics considered in Sec. IV.  $S(t)$  may—or may not—satisfy a sufficiently general  $H$ -theorem to be physically useful; it is, however, the only viable candidate for an entropy compatible with the underlying symmetries inherent in the subdynamics.

These are encouraging results, but it remains to consider in greater detail the more generic time dependence of  $S$ . After all, the principal utility of the intuitive notion of entropy is that  $S$  satisfies an  $H$ -theorem quite generally. It is, therefore, important to ask whether the  $S$  of Eq. (2.22) necessarily increases monotonically for all times and whether, or under what circumstances, there exists an approach of the field towards a unique equilibrium state.

The basic desired properties of  $S$  can be summarized by three conjectures, the validity of which will be motivated and discussed in the remainder of this subsection.

(1) Neglecting "initial transients," the entropy  $S$  is guaranteed to increase monotonically for all times.

(2) If there exists a well defined static equilibrium state  $\mu_{\text{eq}} \propto \exp(-\beta H)$ , the increase in  $S$  coincides with an approach of  $\mu$  towards its equilibrium form.

(3) The equilibrium  $\mu_{\text{eq}}$  maximizes  $S$ , at least locally, with respect to variations  $\delta\mu$  which satisfy appropriate constraints.

It is important to stress once again that  $S$  cannot increase monotonically for every possible system. If, e.g., the system evidences a periodic evolution,  $S$  must also be periodic; and thus, if the entropy increases at one point of time, it must decrease at another. A universal  $H$ -theorem can hold only for systems which are, in some suitable sense, "complex" or "ergodic." It is also important to emphasize that the neglect of initial transients is an essential caveat. In the absence of initial correlations, the immediate response of the system is to increase its entropy. But, as discussed in Sec. IV, if there do exist such initial correlations, they could in principle induce a net decrease in the entropy. Only after these nontrivial initial correlations have "died away" could one expect a universal  $H$ -theorem to hold.

The existence of such a universal  $H$ -theorem, and the existence of a systematic evolution, is especially difficult to address for a system in which no static end state can exist. Some concrete results are, however, known for the special case of systems characterized by a time-independent Hamiltonian. Most important, perhaps, is the fact that an arbitrary

initial  $\mu_{in}$  cannot converge pointwise towards the equilibrium  $\mu_{eq}$ ! It may well tend towards  $\mu_{eq}$  in some appropriate norm or in some suitable “time-averaged” sense, but a true pointwise convergence is simply impossible.

Indeed, that this is true is very easy to see. Given that  $H$  is independent of time, there exists a conserved energy and, as such, if some  $\mu_{in}$  were to converge towards a  $\mu_{eq}$ , the final inverse temperature  $\beta$  would be determined uniquely by the demand that  $\langle H \rangle_{in} = \langle H \rangle_{eq}$ . The problem, however, is that this is not the only constraint that  $\mu_{eq}$  must satisfy. The linearity of the fundamental Liouville equation implies that any functional  $\text{Tr} \Lambda[\mu]$  must also be conserved, and it is easy to see that, in general, the unique  $\mu_{eq}$  determined by energy conservation cannot satisfy this infinite number of additional constraints. Only if one were to violate the Liouville equation by allowing interactions with some external “bath” could one hope to obtain a true pointwise convergence towards the canonical  $\mu_{eq}$ . A completely isolated system cannot converge pointwise towards the true equilibrium.

However, the fact that  $\mu$  cannot converge pointwise towards  $\mu_{eq}$  does not imply that an approximate  $H$ -theorem cannot exist. Indeed, a naive perturbation expansion appropriate in a “dilute gas” ( $|H^I| \ll |H^0|$ ) approximation shows that an homogeneous system of  $N$  identical point masses will in fact evidence a systematic increase in entropy. Thus, in the absence of initial correlations, one concludes that the one-particle distribution function  $f(i)$  will satisfy the Landau equation.<sup>3,28</sup> And it is well known that the Landau equation implies that  $dS/dt \geq 0$ , with equality holding if and only if  $f(i) \propto \exp(-\beta H^0_i)$ , where, in terms of the particle mass  $m$ ,  $H^0_i = \mathbf{p}_i^2/2m$ .<sup>3</sup>

For these and related reasons, it is natural to ask whether the entropy (2.22) is in fact maximized by the equilibrium  $g_{eq}(A)$ 's. More specifically, one would like to ascertain whether the equilibrium  $S_{eq}$  is a maximum with respect to perturbations  $\delta\mu$  which preserve probability [so that  $\text{Tr}_A g(A) = 1$ ], energy  $\langle H \rangle = \text{Tr} \mu H$ , and any other conserved quantities. This is difficult to determine in general because  $\langle H \rangle$  involves not only the  $g(A)$ 's, but the higher-order correlations buried in such quantities as  $\nu(A,B)$ . This problem of entropy maximization would, however, seem tractable in a dilute gas approximation, where the total energy  $\langle H \rangle$  can be approximated by the mean field value  $\langle H \rangle^m \equiv \text{Tr} \mu_R H$ .

Indeed, by using Lagrange multiplier techniques, Lynden-Bell and Wood<sup>29</sup> have proved that a self-gravitating system of  $N$  identical particles confined within a sufficiently small spherical box will in fact maximize  $S$  at least locally with respect to infinitesimal changes in the  $f$ 's which conserve probability and mean field energy. Their analysis does, however, reveal two interesting (and related?) points: (1) if the box is too large,  $S$  will still be extremized by the equilibrium  $f_{eq}(i)$ 's, but it will not be a local maximum; and (2) although the  $f_{eq}(i)$ 's can maximize  $S$  locally, they never constitute a global entropy maximum. These conclusions appear related to the fact that, for a self-gravitating system, the Hamiltonian is not bounded from below. For more well behaved interactions, one might expect that the equilibrium state  $\mu_{eq}$  does indeed correspond to a global entropy maxi-

mum.

As a concrete test of some interest, discussed in Ref. 6, one can examine the Hamiltonian system of Eq. (5.1) in the limit that the  $\omega_A$ 's and  $c_{AB}$ 's are independent of time, contrasting the equilibrium entropy  $S_{eq}$  with the entropy  $S_{in}$  associated with some plausible initial state of the same energy  $\langle H \rangle$ . Consider, for example, an initial state corresponding to a “pseudothermal” density matrix  $\mu_{in} \propto \exp(-\beta H^0)$ . One concludes then that (i) if the  $\omega_A^{2s}$  are positive and (ii) if the  $c_{AB}$ 's are sufficiently small, so that  $H$  is a positive quadratic form and  $\mu_{eq}$  is well defined, the entropy  $S_{eq}$  associated with  $\mu_{eq}$  will be greater than the initial  $S_{in}$ . It is, however, also easy to show that  $\mu_{in}$  could not evolve pointwise towards the equilibrium  $\mu_{eq}$  since  $\text{Tr} \mu_{in} \log \mu_{in} \neq \text{Tr} \mu_{eq} \log \mu_{eq}$ .

## B. Connection with black hole entropy

Superficially, at least, it might appear that the field entropy  $S(t)$  defined in this paper has precious little to do with the entropy  $S_{BH}$  attributed to a black hole in general relativity. The former provides a measure of correlations induced by the evolving dynamics, whereas the latter, being defined in terms of an event horizon, would appear to be an intrinsically geometric quantity.<sup>17</sup>

There is, however, one significant feature which these two entropies share in common. The existence of a nontrivial entropy reflects in each instance the fact that one considers only an incomplete description of the system. If, as a practical matter, an observer is unable to measure complicated correlations between degrees of freedom, his knowledge of the system is only partial. And similarly, given that an observer cannot probe the insides of a black hole without being lost forever to the external Universe, the existence of such a black hole in his system imposes an inherent limitation upon what he can discover about its state.

If one starts with a field characterized initially by comparatively weak and insignificant correlations amongst degrees of freedom, a knowledge of the  $g(A)$ 's and the associated relevant  $\mu_R$  constitutes a quite reasonable approximation to a complete characterization of the system, so that the field entropy  $S$  will be quite small. However, as the evolving dynamics leads to the generation of more significant and detailed correlations, the  $g(A)$ 's provide an increasingly incomplete characterization of the state of the system. One knows less and less about the true state of the system, and this loss of information is reflected by the fact that  $S(t)$  will increase.

The scattering of a scalar field  $\Phi$  by a Schwarzschild black hole admits to a similar interpretation. An initial state of the field, defined with support only outside the hole, can, at least in principle, be measured with arbitrary precision (modulo, perhaps, uncertainties in the form of the initial correlations). But there is no way that an observer can measure directly the complete final state: the portion of the field captured by the hole will remain, at least classically, forever inaccessible. This scattering experiment entails an intrinsic loss of information, and, as such, should entail also an increase in the entropy. Indeed, this is entirely consistent with the geometric notion of black hole entropy. The piece of the

field captured by the hole will presumably have a non-negative total energy, so that, as a consequence of the scattering process, the mass  $M$  of the hole will increase. This, however, corresponds to an increase in the area  $A = 16\pi M^2$  of the hole and a concomitant increase in the entropy  $S_{\text{BH}} = A/4$ . Indeed, Zurek and Thorne<sup>30</sup> have very much adopted this point of view in their “derivation” of the generalized second law of thermodynamics for a system containing a black hole.

It is important to stress that, even in the presence of a black hole, the fundamental equations characterizing the evolution of the field remain deterministic. Given the specification of initial data, an observer could, at least in principle, predict uniquely the subsequent evolution of the system if he knew already what was inside the hole. Formally, at least, one can define a field theory in a Schwarzschild space-time, even inside the event horizon. It then makes sense to introduce a density matrix  $\mu$  which will satisfy a linear Liouville equation, and, given this Liouville equation, one concludes that the “conventional” entropy  $-\text{Tr} \mu \log \mu$  will be a constant of the motion. The point, however, is that the total  $\mu$  is not accessible to an observer outside the hole. What is relevant to him is instead the “piece” of the density matrix  $\mu$ , which could be measured directly by someone outside the hole. And, as discussed, e.g., by Sorkin and his co-workers,<sup>31</sup> the entropy  $-\text{Tr} \mu$ ,  $\log \mu$ , need not be conserved.

#### ACKNOWLEDGMENTS

It is a pleasure to acknowledge fruitful collaborations with Salman Habib and Bei Lok Hu. I am also grateful to Chris Stevens for finding the patience to explain the subtleties of quantum field theory in curved spaces to a rank novice.

Limited financial support was provided by the National Science Foundation and by the Center for Theoretical Physics at the University of Maryland.

<sup>1</sup>L. Boltzmann, *Wien. Ber.* **66**, 275 (1872).

<sup>2</sup>E. T. Jaynes, *Am. J. Phys.* **33**, 391 (1965).

<sup>3</sup>The notion of a “subdynamics,” at least in its most elegant form, may be attributed primarily to the so-called Brussels school of statistical mechanics. The classic paper on the subject is R. Balescu and J. Wallenborn, *Physica (Utrecht)* **54**, 477 (1971). A more systematic and leisurely discourse may be found in R. Balescu, *Equilibrium and Non-Equilibrium Statistical Mechanics* (Wiley, New York, 1975). A more general, but less elegant, approach has been considered by such authors as C. R. Willis and R. H. Picard, *Phys. Rev. A* **9**, 1343 (1974) or H. E. Kandrup, *Astrophys. J.* **244**, 316 (1981).

<sup>4</sup>See, e.g., M. Haggerty and G. Severne, *Adv. Chem. Phys.* **35**, 119 (1977), and references contained therein.

<sup>5</sup>H. E. Kandrup, *Class. Quantum Grav.* **3**, L 55 (1986).

<sup>6</sup>B. L. Hu and H. E. Kandrup, *Phys. Rev. D* **35**, 1776 (1987).

<sup>7</sup>H. E. Kandrup, *J. Math. Phys.* **26**, 2850 (1985).

<sup>8</sup>Consistent with the point of view adopted by such authors as R. Hakim, *J. Math. Phys.* **8**, 1315 (1967) or W. Israel and H. E. Kandrup, *Ann. Phys. (NY)* **152**, 30 (1984), the point of view adopted in this paper is that one can define a classical theory of nonequilibrium statistical mechanics for any dynamical system with well-defined equations of motion, whether or not these equations derive from a Hamiltonian. The notion of an “equilibrium” theory will presumably require the existence of a time-independent  $H$ , but an analog of a Liouville equation exists quite generally, provided only that one can implement the notion of conservation of probability.

<sup>9</sup>S. Habib, preprint.

<sup>10</sup>W. Israel and H. E. Kandrup, *Ann. Phys. (NY)* **152**, 30 (1984).

<sup>11</sup>See, e.g., W. Israel, in *General Relativity, Papers in Honour of J. L. Synge*, edited by L. O’Raifeartaigh (Clarendon, Oxford, 1972), pp. 201–241.

<sup>12</sup>The existence of a time-independent Hamiltonian is, however, necessary if one is to implement the especially elegant approach to nonequilibrium statistical mechanics developed by the Brussels school.

<sup>13</sup>D. W. Sciama, P. Candelas, and D. Deutsch, *Adv. Phys.* **30**, 327 (1981).

<sup>14</sup>See, e.g., the book by Balescu quoted in Ref. 3.

<sup>15</sup>S. R. de Groot, W. A. van Leeuwen, and Ch. G. van Weert, *Relativistic Kinetic Theory, Principles and Applications* (North-Holland, Amsterdam, 1980).

<sup>16</sup>E. A. Calzetta, S. Habib, and B. L. Hu, submitted to *Phys. Rev. D*.

<sup>17</sup>J. D. Bekenstein, *Phys. Rev. D* **7**, 2333 (1973); S. W. Hawking, *Commun. Math. Phys.* **43**, 199 (1975).

<sup>18</sup>N. D. Birrell and P. C. W. Davies, *Quantum Fields in Curved Spaces* (Cambridge U. P., Cambridge, 1982).

<sup>19</sup>H. E. Kandrup, *J. Math. Phys.* **25**, 3286 (1984).

<sup>20</sup>L. Parker, in *Asymptotic Structure of Spacetime*, edited by F. P. Esposito and L. Witten (Plenum, New York, 1977), pp. 107–226; Ya. B. Zel’dovich, in *Physics of the Expanding Universe*, edited by M. Demianski (Springer, New York, 1979), pp. 60–80.

<sup>21</sup>C. W. Misner, *Phys. Rev. Lett.* **22**, 1071 (1969).

<sup>22</sup>B. L. Hu, *Phys. Rev. D* **9**, 3263 (1974).

<sup>23</sup>The “spirit” of the exposition here is due to Balescu and his co-workers, although the actual manipulations described in this section parallel more closely the approach in the paper of Willis and Picard quoted in Ref. 3.

<sup>24</sup>The idea of using projection operators in statistical mechanics was first systematized by R. Zwanzig, in *Lectures in Theoretical Physics*, edited by W. E. Brittin, B. W. Downes, and J. Downes (Interscience, New York, 1961), pp. 106–141.

<sup>25</sup>As discussed, e.g., by the authors cited in Ref. 3, this third criterion is not just convenient practically, but in fact constitutes a deep and fundamental consistency check on the entire formalism. The conventional Brussels school would demand further that  $P$  be explicitly time independent, but, given that  $L$  itself may well be a function of  $t$ , this will not in general be possible. The application of these ideas to a time-dependent system has been considered in a Newtonian sense by H. E. Kandrup and S. Hill Kandrup, *Astrophys. J.* **277**, 1 (1984), where they were used to drive an exact closed equation for the evolution of the galaxy–galaxy pair correlation function. That successful application shows that these ideas can, at least in principle, be used to construct very different types of subdynamics.

<sup>26</sup>This is, e.g., the approach adopted by the Brussels school.

<sup>27</sup>See, e.g., R. P. Feynman, *Statistical Mechanics, A Set of Lectures* (Benjamin, Reading, MA, 1972).

<sup>28</sup>L. D. Landau, *Phys. Z. Sowjetunion* **10**, 154 (1936).

<sup>29</sup>D. Lynden-Bell and R. Wood, *Mon. Not. R. Astron. Soc.* **138**, 495 (1968).

<sup>30</sup>W. H. Zurek and K. S. Thorne, *Phys. Rev. Lett.* **54**, 2171 (1985).

<sup>31</sup>L. Bombelli, R. K. Koul, J. Lee, and R. D. Sorkin, *Phys. Rev. D* **34**, 373 (1986).

# A phase cell approach to Yang–Mills theory. II. Analysis of a mode

Paul Federbush

Department of Mathematics, University of Michigan, Ann Arbor, Michigan 48109

Calvin Williamson

Department of Mathematics, University of Missouri—Columbia, Columbia, Missouri 65211

(Received 18 November 1986; accepted for publication 30 January 1987)

Properties of a mode announced in a previous paper are proved. This involves some complicated calculations in linear algebra, and observation of the structure of a function of several complex variables.

## I. $A'_1$ AND $A'_2$ , LINEAR ALGEBRA

$A'_i(p)$  is given by (3.12) of Ref. 1. We will consider a mode at level zero, specified by bond assignments at level zero, having only one bond assignment different from zero. We choose the bond to be the bond in the  $+1$  direction from the origin, and the assignment to it to be unity. By invariance under certain interchanges of coordinate axes it is sufficient to find  $A'_1$  and  $A'_2$  to determine  $A'_i$ .

With  $i$  fixed, (3.12) of Ref. 1 is a matrix product,  $ABC$ , with  $A, B, C$ , respectively, a  $(1 \times 6)$  matrix,  $(6 \times 6)$  matrix, and  $(6 \times 1)$  matrix. The coordinates of the six-dimensional space are labeled by oriented plaquette directions, for which we choose the ordering

$$(1,2), (2,3), (3,4), (1,3), (1,4), (2,4). \quad (1.1)$$

We find it convenient to introduce certain special notation and matrices. We let

$$f_i \equiv (e^{-ip_i} - 1) \quad (1.2)$$

and  $D$  be the  $6 \times 6$  diagonal matrix with

$$D_{(ij),(ij)} \equiv f_i f_j \quad (1.3)$$

(here labeling rows and columns of  $D$  by the associated oriented plaquette directions). We let

$$\langle h(p) \rangle \equiv \sum_n \frac{1}{(p + 2\pi n)^2} \prod_{i=1}^4 \left| \frac{e^{ip_i} - 1}{(p_i + 2\pi n_i)} \right|^2 h(p + 2\pi n), \quad (1.4)$$

where  $h$  is any function of  $p = (p_1, p_2, p_3, p_4)$ .

We now detail the matrices occurring in (3.12) of Ref. 1. We first view the  $(6 \times 1)$  matrix

$1/p_1^2 + 1/p_2^2$	$-1/p_2^2$	0	$1/p_1^2$	$1/p_1^2$	$-1/p_2^2$
$-1/p_2^2$	$1/p_2^2 + 1/p_3^2$	$-1/p_3^2$	$1/p_3^2$	0	$1/p_2^2$
0	$-1/p_3^2$	$1/p_3^2 + 1/p_4^2$	$-1/p_3^2$	$1/p_4^2$	$1/p_4^2$
$1/p_1^2$	$1/p_3^2$	$-1/p_3^2$	$1/p_1^2 + 1/p_3^2$	$1/p_1^2$	0
$1/p_1^2$	0	$1/p_4^2$	$1/p_1^2$	$1/p_1^2 + 1/p_4^2$	$1/p_4^2$
$-1/p_2^2$	$1/p_2^2$	$1/p_4^2$	0	$1/p_4^2$	$1/p_2^2 + 1/p_4^2$

Here  $\langle M_0 \rangle$  is a positive Hermitian matrix, and  $w_1$  is in its range. We deal with the limit

$$\lim_{\epsilon \rightarrow 0^+} (\langle M_0 \rangle + \epsilon)^{-1} w_1. \quad (1.14)$$

$$\beta_\gamma e^{i\gamma \cdot p} \equiv v_1 \equiv \bar{f}_1^{-1} \bar{D} w_1, \quad (1.5)$$

$$w_1 \equiv (1,0,0,1,1,0)^T, \quad (1.6)$$

where the bars indicate complex conjugates. We next view the left  $(1 \times 6)$  matrix, two different vectors depending on whether one is studying  $A'_1$  or  $A'_2$ . With  $r_L$  defined by

$$r_L \equiv -\frac{i}{(2\pi)^2} \prod_k \left( \frac{e^{-ip_k} - 1}{p_k} \right), \quad (1.7)$$

we have

$$\tilde{P}_1 = r_L (1/p_1) w_1^T D \quad (1.8)$$

and

$$\tilde{P}_2 = r_L (1/p_2) w_2^T D \quad (1.9)$$

with

$$w_2 \equiv (-1,1,0,0,0,1)^T. \quad (1.10)$$

At this point we may write (3.12) of Ref. 1 as

$$A'_i(p) = (1/p^2) \tilde{P}_i M^{-1} v_1. \quad (1.11)$$

Recall  $M$  is singular; and  $M^{-1}$  is defined on suitable vectors as  $\lim_{\epsilon \rightarrow 0^+} (\epsilon + M)^{-1}$ .

Substituting (1.5) and (1.8) or (1.9) into (1.11) we get

$$A'_i = (-1/p^2) r_L \bar{f}_1^{-1} (1/p_i) w_i^T \langle M_0 \rangle^{-1} w_1, \quad i = 1, 2, \quad (1.12)$$

where

$$M = \bar{D} \langle M_0 \rangle D \quad (1.13)$$

and the inverse of  $\langle M_0 \rangle$  taken in the same sense as the inverse of  $M$ . We now display  $M_0$ , itself, in its full glory:

This limit may be taken by the following procedure. Let  $P$  be the (orthogonal) projection onto the range of  $\langle M_0 \rangle$ ,  $S$ . Let  $m = P \langle M_0 \rangle P$ , a strictly positive Hermitian matrix. We invert  $m$ , in  $S$ , writing it as  $m^{-1}$ . One has for the result of (1.14)

$$m^{-1}w_1. \quad (1.15)$$

The observations of this paragraph reduce the computations of (1.12) to calculations in a three-dimensional vector space. We note that  $S$  is spanned by  $w_1, w_2,$  and  $w_3 = (0,0,1,0,1,1)^T$ .

We first present the results of the computation in (1.15):

$$\begin{aligned} \langle M_0 \rangle^{-1}w_1 = & (1/\mathcal{D})(4bc + 4bd + 8cd, 4bd - 4cd, \\ & 4bc - 4bd, 4bc + 8bd + 4cd, \\ & 8bc + 4bd + 4cd, 4bc - 4cd)^T, \end{aligned} \quad (1.16)$$

where

$$\begin{aligned} a = \langle 1/p_1^2 \rangle, \quad b = \langle 1/p_2^2 \rangle, \\ c = \langle 1/p_3^2 \rangle, \quad d = \langle 1/p_4^2 \rangle, \end{aligned} \quad (1.17)$$

and

$$\mathcal{D} = 16 \sum_k \prod_{j \neq k} \langle 1/p_j^2 \rangle. \quad (1.18)$$

Here  $\mathcal{D}$  has arisen as the determinant of the  $3 \times 3$  matrix involved in the computation of  $m^{-1}$ .

Collecting our results and substituting into (1.12) we easily find

$$\begin{aligned} A'_1 = & \left( -r_L \bar{f}_1^{-1} \frac{1}{p^2} \right) \frac{1}{p_1} \frac{16}{\mathcal{D}} \left( \left\langle \frac{1}{p_2^2} \right\rangle \left\langle \frac{1}{p_3^2} \right\rangle \right. \\ & \left. + \left\langle \frac{1}{p_2^2} \right\rangle \left\langle \frac{1}{p_4^2} \right\rangle + \left\langle \frac{1}{p_3^2} \right\rangle \left\langle \frac{1}{p_4^2} \right\rangle \right), \end{aligned} \quad (1.19)$$

$$A'_2 = \left( -r_L \bar{f}_1^{-1} \frac{1}{p^2} \right) \frac{1}{p_2} \frac{16}{\mathcal{D}} \left( - \left\langle \frac{1}{p_3^2} \right\rangle \left\langle \frac{1}{p_4^2} \right\rangle \right). \quad (1.20)$$

## II. $A'_i$ AND $X(p)$

From (1.19) and (1.20), and the invariance under interchange of coordinate directions—except for the special one-direction—we find our expression for  $A'_i(p)$ ,

$$A'_i = \left( -r_L \bar{f}_1^{-1} \frac{1}{p^2} \right) \frac{1}{p_i} \frac{16}{\mathcal{D}} l_i, \quad (2.1)$$

$$l_1 = \left\langle \frac{1}{p_2^2} \right\rangle \left\langle \frac{1}{p_3^2} \right\rangle + \left\langle \frac{1}{p_2^2} \right\rangle \left\langle \frac{1}{p_4^2} \right\rangle + \left\langle \frac{1}{p_3^2} \right\rangle \left\langle \frac{1}{p_4^2} \right\rangle, \quad (2.2)$$

$$l_i = - \prod_{j \neq 1, i} \left\langle \frac{1}{p_j^2} \right\rangle, \quad i \neq 1. \quad (2.3)$$

A little study of these expressions for the region near  $p_1 \sim p_2 \sim p_3 \sim p_4 \sim 0$  shows that the  $A'_i$  have a singularity in the complex four-dimensional region on the surface  $p^2 = 0$ , of course hitting the real axis. We define  $A_i^N(p)$ , a gauge transformation of  $A'_i(p)$ , by

$$A_i^N(p) = A'_i(p) + p_i X(p), \quad (2.4)$$

with

$$X(p) = \left( -r_L \bar{f}_1^{-1} \frac{1}{(p^2)^2} \right) \frac{1}{\langle 1/p_1^2 \rangle} \frac{1}{(1+p^2)^s}, \quad (2.5)$$

where  $s$  is an arbitrary integer, sufficiently large. The analyt-

ic properties of  $A_i^N(p)$  will be studied in the next section.

## III. ANALYTIC PROPERTIES OF $A_i^N(p)$

The analysis we follow herein is analogous to the similar analysis due to Gawedzki and Kupiainen in the appendix of Ref. 2. We take the analytic extensions of the expressions in Sec. II from real  $p$  to complex  $p$ . For example,

$$\bar{f}_i \rightarrow e^{ip_i} - 1$$

and

$$\left| \frac{e^{-ip_i} - 1}{p_i + 2\pi n_i} \right|^2 \rightarrow \frac{e^{-ip_i} - 1}{p_i + 2\pi n_i} \frac{e^{ip_i} - 1}{p_i + 2\pi n_i}.$$

We make the preliminary observation that  $f_i, D, \mathcal{D}$ , and  $\langle \cdot \rangle$  are periodic functions of  $p$ , invariant under

$$p \rightarrow p + 2\pi n \quad (3.1)$$

( $n$  is a four-vector with integer components). We note that  $r_L, p^2$ , and  $p_i$  are not periodic functions of  $p$ .

Because of the periodicity mentioned in the last paragraph it is natural to divide our results into the following three theorems.

**Theorem 3.1 (Local analyticity):** There is an  $\epsilon_0 > 0$  such that in the domain,  $\mathcal{D}_L$ , specified by

$$-\pi < \text{Re } p_j < \pi, \quad |\text{Im } p_j| < \epsilon_0, \quad (3.2)$$

$A_i^N(p)$  is analytic.

**Theorem 3.2 (Global analyticity):**  $A_i^N(p)$  is analytic in the domain,  $\mathcal{D}_G$ , specified by

$$|\text{Im } p_j| < \epsilon_0. \quad (3.3)$$

**Theorem 3.3 (Boundedness):** Within the domain,  $\mathcal{D}_B$ , specified by

$$|\text{Im } p_j| < \epsilon_0/2, \quad (3.4)$$

$A_i^N(p)$  satisfies bounds of the form

$$|A_i^N(p)| < c \prod_j \frac{1}{(|p_j| + 1)} \frac{1}{|p^2| + 1}. \quad (3.5)$$

In Theorems 3.2 and 3.3 the  $\epsilon_0$  is as defined in Theorem 3.1.

Theorems 3.2 and 3.3 are easy consequences of the form of  $A_i^N(p)$  and Theorem 3.1. We will devote our attention entirely to the proof of Theorem 3.1, which is carried out in Sec. V.

## IV. CONCLUSIONS

Equations (3.13)–(3.15) of Ref. 1 follow directly from Theorems 3.2 and 3.3 of the last section by standard techniques.

## V. LOCAL ANALYTICITY

We must study  $A_i^N(p)$  on  $\mathcal{D}_L$  [of (3.2)]. We write  $A_i^N(p)$  from (2.1)–(2.5),

$$A_i^N(p) = -r_L \bar{f}_1^{-1} \frac{1}{p^2} \left[ \frac{1}{p_i} \frac{16}{\mathcal{D}} l_i + \frac{p_i}{p^2} \frac{1}{\langle 1/p_1^2 \rangle} \frac{1}{(1+p^2)^s} \right]. \quad (5.1)$$

In studying analyticity we may remove the factor  $-r_L$  and replace  $\bar{f}_1^{-1}$  by  $p_1^{-1}$ . We thus study

$$a_i(p) = \frac{1}{p_1} \frac{1}{p^2} \left[ \frac{1}{p_i} \frac{16}{\mathcal{D}} l_i + \frac{p_i}{p^2} \frac{1}{\langle 1/p_i^2 \rangle} \frac{1}{(1+p^2)^s} \right]. \quad (5.2)$$

We introduce some notation

$$p_i^2 \langle 1/p_i^2 \rangle \equiv r_0(p) e_j(p) (1/p^2), \quad (5.3)$$

where  $e_j(0) = 1$  and

$$r_0(p) = \prod_{i=1}^4 \left| \frac{e^{-ip_i} - 1}{p_i} \right|^2. \quad (5.4)$$

It will be sufficient to study  $a_1(p)$  and  $a_2(p)$ . With the notation of (5.3) and (5.4) we have

$$r_0(p) a_1(p) = \frac{1}{\sum_k p_k^2 \prod_{j \neq k} e_j} (p_4^2 e_2 e_3 + p_3^2 e_2 e_4 + p_2^2 e_3 e_4) + \frac{1}{p^2} \frac{p_1^2}{e_1} \frac{1}{(1+p^2)^s}, \quad (5.5)$$

$$r_0(p) a_2(p) = \frac{-p_2 p_1}{\sum_k p_k^2 \prod_{j \neq k} e_j} e_3 e_4 + \frac{p_2 p_1}{p^2} \frac{1}{e_1} \frac{1}{(1+p^2)^s}. \quad (5.6)$$

We now write  $e_i(p)$  as

$$e_i(p) = 1 + \delta_i(p), \quad (5.7)$$

near  $p = 0$ ,  $\delta_i$  is small. In (5.5) we substitute (5.7) and set  $\delta_i = 0$  to get

$$r_0(p) a_1^0(p) \equiv 1 + \frac{p_1^2}{p^2} \left( \frac{1}{(1+p^2)^s} - 1 \right) \quad (5.8)$$

and

$$r_0(p) a_2^0(p) \equiv \frac{p_2 p_1}{p^2} \left( \frac{1}{(1+p^2)^s} - 1 \right). \quad (5.9)$$

Writing

$$a_i(p) = a_i^0(p) + R_i(p), \quad (5.10)$$

we see from (5.8) and (5.9) that  $a_i^0(p)$  are analytic in  $\mathcal{D}_L$ . This is the main "algebraic miracle" involved in showing local analyticity of  $A_i^N(p)$ . This algebraic miracle motivated, of course, our choice of  $X_i(p)$ , which, however, was far from unique. We turn our attention to  $R_i(p)$ . By showing the analyticity in  $\mathcal{D}_L$  of  $R_i$  we will have completed the proof. (We have already used the analyticity in  $\mathcal{D}_L$  of  $r_L$ ,  $p_1 \tilde{f}_1^{-1}$ , and  $r_0^{-1}$ .)

All hinges on certain properties of  $e_i$  and  $\delta_i$  in  $\mathcal{D}_L$ , which we now pursue.

## VI. PROPERTIES OF THE $e_i(p)$ IN $\mathcal{D}_L$

$$(1) \quad e_i(p) \text{ is analytic}, \quad (6.1)$$

$$(2) \quad \delta_i(p) = p_i^2 p^2 h_i(p) \text{ with } h_i(p) \text{ analytic}, \quad (6.2)$$

$$(3) \quad e_i^{-1}(p) \text{ is analytic}, \quad (6.3)$$

$$(4) \quad \left[ \frac{1}{\sum_k p_k^2 \prod_{j \neq k} e_j} - \frac{1}{p^2} \right] \text{ is analytic}. \quad (6.4)$$

It is quite immediate that (6.1)–(6.4) imply the analyticity of  $R_i(p)$ . We are faced with our final task, the proof of these properties. (1) and (2) follow from the form of  $e_i(p)$  upon inspection. It is only (3) and (4) that must be studied.

We write  $e_i(p)$  in elaborate fashion [ $e_1(p)$  and  $e_i(p)$ ,  $i \neq 1$ , have the same properties],

$$e_1(p) = 1 + p^2 \sum_{\rho} \left[ \chi_{\rho}(1) \prod_{i \in \rho} p_i^2 K_{\rho}(p) + \chi_{\rho^c}(1) p_1^2 \prod_{i \in \rho} p_i^2 K_{\rho}^1(p) \right], \quad (6.5)$$

$$K_{\rho}(p) = \sum_{n \sim \rho} \frac{1}{(p + 2\pi n)^2} \prod_{i \in \rho} \frac{1}{(p_i + 2\pi n_i)^2}, \quad (6.6)$$

$$K_{\rho}^1(p) = \sum_{n \sim \rho} \frac{1}{(p + 2\pi n)^2} \prod_{i \in \rho} \frac{1}{(p_i + 2\pi n_i)^2} \times \frac{1}{(p_1 + 2\pi n_1)^2}, \quad (6.7)$$

where

(a)  $\rho$  is a proper subset of  $(1,2,3,4)$ ;

$$(b) \quad \chi_{\rho}(1) = \begin{cases} 1 & 1 \in \rho, \\ 0 & 1 \notin \rho, \end{cases}$$

$$\chi_{\rho^c}(1) = \begin{cases} 0 & 1 \in \rho, \\ 1 & 1 \notin \rho; \end{cases}$$

(c)  $n \sim \rho$  means that  $\rho = \{i, n_i = 0\}$ .

We easily see, for  $p$  in  $\mathcal{D}_L$ , that

$$|K_{\rho}(p)| < m, \quad |K_{\rho}^1(p)| < m, \quad (6.8)$$

for some fixed  $m$ . And that

$$|\text{Arg } K_{\rho}(p)| < \epsilon', \quad |\text{Arg } K_{\rho}^1(p)| < \epsilon', \quad (6.9)$$

where  $\epsilon'$  can be made arbitrarily small by choosing  $\epsilon_0$  of (3.2), suitably small. Again we have

$$\begin{aligned} |\text{Arg } p_i^2| < \epsilon'' & \quad \text{if } |p_i| > \bar{\epsilon}, \\ |\text{Arg } p^2| < \epsilon'' & \quad \text{if } |p|^2 > \bar{\epsilon}, \end{aligned} \quad (6.10)$$

where  $\epsilon''$  and  $\bar{\epsilon}$  can be fixed arbitrarily small, choosing  $\epsilon_0$  small enough. Picking  $\epsilon'$ ,  $\epsilon''$ , and  $\bar{\epsilon}$  small enough, and using (6.8), we see  $e_1(p)$  is invertible in  $\mathcal{D}_L$  and so analytic implying (3), (6.3). [The terms on the right side of (6.5), individually, will either be small, or have small argument.]

We turn to the study of (4), (6.4). We write the brackets in (6.4) as

$$\frac{1}{p^2(1+g)} - \frac{1}{p^2} = \frac{1}{p^2} \frac{g}{1+g}, \quad (6.11)$$

this relation defining  $g$ . Property (4) will be proved by showing

$$(5) \quad g/p^2 \text{ is analytic}, \quad (6.12)$$

$$(6) \quad (1+g)^{-1} \text{ is analytic}. \quad (6.13)$$

For  $g$  we find the formula

$$g = \sum_k \frac{p_k^2}{p^2} \left[ \prod_{j \neq k} (1 + \delta_j) - 1 \right]. \quad (6.14)$$

From (6.2) we see that (5) [(6.12)] holds. The proof that  $(1+g)$  is invertible, and so analytic, follows quite immediately from the proof above that  $e_1(p)$  is invertible.



## ACKNOWLEDGMENTS

One of the authors (C.W.) was supported in part by the University of Missouri Research Council. This work was supported in part by the National Science Foundation under

Grant No. PHY-85-02074.

<sup>1</sup>P. Federbush, "A phase cell approach to Yang-Mills theory. I. Modes, lattice-continuum duality," *Commun. Math. Phys.* **107**, 319 (1986).

<sup>2</sup>K. Gawedzki and A. Kupiainen, "A rigorous block spin approach to massless lattice theories," *Commun. Math. Phys.* **77**, 31 (1980).

# Quantum noncompact $\sigma$ models<sup>a)</sup>

J. W. van Holten

NIKHEF-H, P.O. Box 41882, 1009 DB Amsterdam, The Netherlands

(Received 9 December 1986; accepted for publication 4 February 1987)

Two possible quantization procedures for field theories with noncompact symmetry are discussed: one in a positive-definite Hilbert space with negative-energy states and one in a Hilbert space with indefinite metric and positive-definite energy. The physical interpretation of these alternative procedures is explained in terms of the realization of the noncompact symmetry, which is broken in one case and not in the other. It is shown that in exactly solvable quantum  $\sigma$  models the noncompact symmetry is broken in the vacuum. Some arguments supporting the general validity of this result are given. It is concluded that in contrast to theories like quantum electrodynamics, the noncompact  $\sigma$  models are to be quantized in the positive-definite Hilbert space with negative-energy modes.

## I. INTRODUCTION

Positivity of the energy of physical systems is one of the fundamental postulates of quantum theory. It is introduced to explain the existence of a stable ground state of the system. In addition, it is fundamental in establishing the second law of thermodynamics (i.e., the increase of entropy with energy). In particular, the kinetic energy of physical systems is usually assumed to be positive and therefore takes the form (after proper normalization of the variables):

$$T(\dot{q}_i) = \sum_i \frac{1}{2} \dot{q}_i^2. \quad (1.1)$$

If the total energy is

$$H = T(\dot{q}_i) + \lambda V(q_i), \quad (1.2)$$

and if  $T(\dot{q}_i) \gg \lambda V(q_i)$ , the system has an approximate  $O(N)$  symmetry:

$$q'_i = R_{ij} q_j, \quad R_{ij} = R_{ji}^{-1}, \quad (1.3)$$

which becomes exact if either  $V(q_i)$  is an invariant itself,  $V(q_i) = f(q_i^2)$  (central forces), or if  $\lambda \rightarrow 0$ .

In systems with infinitely many degrees of freedom, like field theories, similar observations can be made. In particular, for fields  $\phi_i$  described by a positive-definite Hamiltonian density

$$H = \frac{1}{2} \sum_i (\dot{\phi}_i^2 + (\nabla \phi_i)^2) + \lambda V(\phi_i), \quad (1.4)$$

the kinetic term and gradient terms are  $O(N)$  invariant, a consequence of the positivity of the theory. However, there are physically relevant theories that admit an (exact or approximate) noncompact symmetry, such as  $O(N, M)$ . Field theories of this type have Hamiltonian densities of the form

$$H = \frac{1}{2} \sum_i (\dot{\pi}_i^2 + (\nabla \pi_i)^2) - \frac{1}{2} \sum_r (\dot{\sigma}_r^2 + (\nabla \sigma_r)^2) + \lambda V(\sigma, \pi), \quad (1.5)$$

where  $\{\pi_i\}$  denotes a set of fields with positive energy and  $\{\sigma_r\}$  a set of fields with negative energy. Of course, the appearance of degrees of freedom with negative energy contradicts the physical principles mentioned above. Although in some cases one may question the absolute validity of these principles, we study in this paper systems for which the degrees of freedom with negative energy are unphysical and can be eliminated by imposing suitable constraints, resulting in Hamiltonians which are bounded below. [For example, no instability arises in the absence of perturbations or interactions coupling to both the positive and the negative-energy modes; therefore, under special circumstances, it is possible to create systems with negative temperatures (Hamiltonians bounded above) in the laboratory.<sup>1</sup>]

Examples of theories of this kind are quite easy to find, and some are actually very familiar. As a first example, consider the  $O(N, 1)$  nonlinear  $\sigma$  model.<sup>2-12</sup> This model can be defined by the  $O(N, 1)$  invariant Lagrangian [our metric has signature  $(+, +, +, -)$ ],

$$\mathcal{L} = -\frac{1}{2} \sum_{i=1}^N (\partial_\mu \pi_i)^2 + \frac{1}{2} (\partial_\mu \sigma)^2, \quad (1.6)$$

with the additional constraint

$$\sigma^2 - \pi^2 = 1/g. \quad (1.7)$$

To prove, that the energy of the system is bounded below, we solve the constraint in terms of  $N$  scalars  $\varphi_i$ :

$$\pi = \frac{\varphi}{\sqrt{1-g\varphi^2}}, \quad \sigma = \frac{1}{\sqrt{g}} \frac{1}{\sqrt{1-g\varphi^2}}. \quad (1.8)$$

Then the Lagrangian becomes

$$\mathcal{L} = -\frac{1}{2} \frac{(\partial_\mu \varphi)^2}{(1-g\varphi^2)^2}. \quad (1.9)$$

The conjugate momentum to  $\varphi$  is then

$$\Pi_\varphi = \frac{\dot{\varphi}}{(1-g\varphi^2)^2}, \quad (1.10)$$

and the corresponding Hamiltonian is

$$H = \frac{1}{2} (1-g\varphi^2)^2 \Pi_\varphi^2 + \frac{1}{2} \frac{(\nabla \varphi)^2}{(1-g\varphi^2)^2}. \quad (1.11)$$

<sup>a)</sup> This is a combined and revised version of two earlier papers, which appeared as preprints WUB 84-9 (unpublished) and NIKHEF 85-12 (unpublished).

It is manifestly positive definite. Note, that the Lagrangian (1.9) still possesses the noncompact  $O(N,1)$  invariance, but now realized nonlinearly:

$$\delta\varphi_i = a_{ij}\varphi_j + (1/\sqrt{g})(b_i + gb_i\varphi^2 - 2g\mathbf{b}\cdot\boldsymbol{\varphi}\varphi_i). \quad (1.12)$$

Hence, the physical scalars  $\pi_i$  or  $\varphi_i$  may be interpreted as Goldstone bosons describing the spontaneous breaking of an internal  $O(N,1)$  symmetry to its compact subgroup  $O(N)$  (Refs. 2,4,5). This example can easily be generalized to other nonlinear  $\sigma$  models on  $G/H$ , where  $G$  is noncompact and  $H$  its maximal compact subgroup. Such noncompact  $\sigma$  models form an intrinsic part of  $N$ -extended supergravity theories ( $N \geq 4$ ), where they describe the scalar partners of the graviton.<sup>13</sup>

A second example of a constrained system with noncompact symmetry is provided by Maxwell's theory of electromagnetism. In this theory, the fundamental fields may be represented by a vector potential  $A_\mu$ , from which one can compute electric and magnetic field strengths  $F_{\mu\nu} = \partial_\mu A_\nu - \partial_\nu A_\mu$ . Since the physically observable field strengths are invariant under a gauge transformation  $A_\mu \rightarrow A'_\mu = A_\mu + \partial_\mu \Lambda$ , we can impose the Lorentz condition  $\partial \cdot A = 0$ . The Maxwell equations which describe the dynamics of the electromagnetic fields can then be obtained from the Lagrangian density

$$\mathcal{L} = -\frac{1}{2}(\partial_\mu \mathbf{A})^2 + \frac{1}{2}(\partial_\mu A_0)^2, \quad (1.13)$$

with the fields constrained by

$$\partial \cdot A = \nabla \cdot \mathbf{A} - \dot{A}_0 = 0. \quad (1.14)$$

This Lagrangian is manifestly invariant under  $O(3,1)$  rotations on the fields  $(\mathbf{A}, A_0)$ . Of course, these transformations are just the standard Lorentz transformations acting on the four-vector potential  $A_\mu$ . The correct Hamiltonian density for the electric and magnetic fields is

$$H = \frac{1}{2}(\mathbf{E}^2 + \mathbf{B}^2), \quad (1.15)$$

where  $\mathbf{E} = \nabla A_0 - \dot{\mathbf{A}}$  and  $\mathbf{B} = \nabla \times \mathbf{A}$ , as usual. Again the Hamiltonian is positive definite.

A third example worth mentioning here is the theory of the relativistic string.<sup>14</sup> The free string may be described by the action<sup>15</sup>

$$S = \frac{1}{2} \int d^2\xi \sqrt{g} g^{ab} \partial_a X^\mu \partial_b X^\nu \eta_{\mu\nu}. \quad (1.16)$$

The variables  $X^\mu$  represent the string coordinates, whilst  $g_{ab}$  is the metric of the two-dimensional world sheet of the string;  $\eta_{\mu\nu}$  is the Lorentzian metric of the  $D$ -dimensional space-time in which the string moves. Clearly, the action is invariant under the (noncompact)  $D$ -dimensional Lorentz transformations. The 2-D metric  $g$  is not a dynamical degree of freedom, but may be interpreted as a kind of Lagrange multiplier imposing constraints

$$T_{ab} = \partial_a X^\mu \partial_b X_\mu - \frac{1}{2} g_{ab} g^{cd} \partial_c X^\mu \partial_d X_\mu = 0. \quad (1.17)$$

Again, these constraints guarantee the positivity of the energy in this theory.<sup>16</sup>

Thus we see that constrained systems with noncompact symmetries are quite common. Note, however, that there is an important difference between the first example and the

other two examples: in the nonlinear  $\sigma$  model we could interpret the noncompact symmetry as a spontaneously broken internal symmetry, whereas the noncompact symmetry of the electromagnetic field and the relativistic string is Lorentz invariance (a space-time symmetry) which normally we do not expect to be broken.

In this paper we address the problem of quantizing systems with noncompact symmetries, and we show that the physical realization of the noncompact symmetry (spontaneously broken or not broken) determines the correct quantization scheme. We illustrate the general principles only for a few simple cases like the  $O(N,1)$  model and Maxwell's theory, but the results carry over directly to more complicated models. Applications have been given elsewhere,<sup>11</sup> but a brief discussion of the consistency of the quantized theories is presented in the last section of this paper.

## II. CANONICAL QUANTIZATION OF THE NEGATIVE-ENERGY MODES

Free fields may be considered as collections of harmonic oscillators. Therefore the study of harmonic oscillators is directly applicable to the case of relativistic fields. In this section we discuss two different quantization procedures for oscillators with negative kinetic energy.<sup>17-20</sup> Consider a classical harmonic oscillator with negative energy and Lagrangian

$$L = -\frac{1}{2}\dot{q}^2 + \frac{1}{2}\omega^2 q^2. \quad (2.1)$$

Note, that owing to the time-reversal invariance of  $L$ , the action  $\int L dt$  is obtained from that of the positive-energy oscillator by reversing time:  $t \rightarrow -t$ . The first method to construct a quantized version of this system uses the standard correspondence between Poisson brackets and commutators. For the momentum and Hamiltonian of the classical theory defined by (2.1) we obtain

$$p = \frac{\partial L}{\partial \dot{q}} = -\dot{q}, \quad H = p\dot{q} - L = -\frac{1}{2}(p^2 + \omega^2 q^2). \quad (2.2)$$

Clearly the classical Hamiltonian is negative definite (hence bounded above). The Hamiltonian equations of motion are

$$\dot{q} = \{H, q\}, \quad \dot{p} = \{H, p\}, \quad (2.3)$$

where the Poisson brackets are defined by

$$\{A, B\} = \frac{\partial A}{\partial p} \frac{\partial B}{\partial q} - \frac{\partial B}{\partial p} \frac{\partial A}{\partial q}. \quad (2.4)$$

It follows that

$$\{p, q\} = 1. \quad (2.5)$$

According to the correspondence principle, one obtains a quantum theory by replacing the coordinates and momenta by operators  $\hat{p}$  and  $\hat{q}$  and the Poisson bracket  $\{, \}$  by a commutator  $(i/\hbar) [, ]$ . Thus

$$[\hat{p}, \hat{q}] = -i\hbar. \quad (2.6)$$

As for the standard harmonic oscillator, we introduce the creation and annihilation operators defined by

$$\hat{a} = (1/\sqrt{2\hbar\omega})(\omega\hat{q} + i\hat{p}), \quad \hat{a}^\dagger = (1/\sqrt{2\hbar\omega})(\omega\hat{q} - i\hat{p}). \quad (2.7)$$

They satisfy the commutation relation

$$[\hat{a}, \hat{a}^\dagger] = 1. \quad (2.8)$$

With these operators one can construct a complete Hilbert space of orthonormal states for the negative-energy oscillator, starting from the ground state defined by

$$\hat{a}|0\rangle = 0, \quad \langle 0|0\rangle = 1. \quad (2.9)$$

This state is well defined, as may be checked in the coordinate representation, in which

$$|0\rangle \rightarrow \psi_0 = c e^{- (\omega/2\hbar) q^2}. \quad (2.10)$$

The excited states are

$$|n\rangle = ((\hat{a}^\dagger)^n / \sqrt{n!})|0\rangle, \quad \langle n|m\rangle = \delta_{n,m}. \quad (2.11)$$

Thus all states have positive norm. The Hamiltonian reads

$$\hat{H} = -\hbar\omega(\hat{a}^\dagger\hat{a} + \frac{1}{2}). \quad (2.12)$$

Its eigenvalues are negative definite:

$$E_n = - (n + \frac{1}{2})\hbar\omega, \quad n = 0, 1, 2, \dots \quad (2.13)$$

The time dependence of the operators in the Heisenberg representation follows from the Schrödinger equation

$$i\hbar \frac{d\hat{a}}{dt} = [\hat{a}, \hat{H}]. \quad (2.14)$$

This gives

$$\hat{a}(t) = \hat{a}(0)e^{i\omega t}, \quad \hat{a}^\dagger(t) = \hat{a}^\dagger(0)e^{-i\omega t}. \quad (2.15)$$

Since  $\omega$  is taken positive, the annihilation operator  $\hat{a}$  here has the same time dependence as the creation operator  $\hat{a}^\dagger$  in the case of the positive-energy oscillator, and vice versa. From Eqs. (2.7) and (2.15) one obtains

$$\hat{q}(t) = (\hbar/2\omega)^{1/2}(\hat{a}(0)e^{i\omega t} + \hat{a}^\dagger(0)e^{-i\omega t}), \quad (2.16)$$

and for the Green's function:

$$\begin{aligned} \langle 0|T\hat{q}(t)\hat{q}(0)|0\rangle &= \frac{\hbar}{2\omega}(\theta(t)e^{i\omega t} + \theta(-t)e^{-i\omega t}) \\ &= -\frac{i\hbar}{2\pi} \int_{-\infty}^{\infty} dk \frac{e^{-ikt}}{k^2 - \omega^2 - i\epsilon}. \end{aligned} \quad (2.17)$$

This Green's function differs from that of the standard positive-energy oscillator in two respects: (1) the  $i\epsilon$  prescription has the opposite sign; and (2) the residue has an extra minus sign as well.

The above results are easily understandable from the point of view of time reversal, since negative energies and positive time direction are equivalent to positive energies and negative time direction. The notions of advanced and retarded Green's functions have therefore been interchanged. We now contrast this treatment with an alternative method of quantization involving negative-norm states. This method results in positive energy expectation values. As noted already below Eq. (2.15) positive energies require interchanging the role of creation and annihilation operators.<sup>17</sup> Thus if we define

$$\begin{aligned} \hat{b}(t) &= \hat{b}(0)e^{-i\omega t} = \hat{a}^\dagger(t), \\ \hat{b}^\dagger(t) &= \hat{b}^\dagger(0)e^{i\omega t} = \hat{a}(t), \end{aligned} \quad (2.18)$$

then

$$\hat{q}(t) = (\hbar/2\omega)^{1/2}(\hat{b}(0)e^{-i\omega t} + \hat{b}^\dagger(0)e^{i\omega t}). \quad (2.19)$$

The Hamiltonian becomes

$$\hat{H} = \hbar\omega(-\hat{b}^\dagger\hat{b} + \frac{1}{2}). \quad (2.20)$$

In spite of its somewhat unusual appearance, this represents an oscillator with positive energy. The difference with the ordinary harmonic oscillator is that  $\hat{b}, \hat{b}^\dagger$  satisfy unphysical equal-time commutation relations<sup>17,18</sup>:

$$[\hat{b}, \hat{b}^\dagger] = -1. \quad (2.21)$$

If we define the energy eigenstates by

$$\hat{b}|0\rangle = 0, \quad \langle 0|0\rangle = 1, \quad (2.22)$$

for the ground state, and

$$|n\rangle = ((b^\dagger)^n / \sqrt{n!})|0\rangle, \quad (2.23)$$

for the excited states, then the energy expectation values are  $\langle E_n \rangle = \hbar\omega(n + \frac{1}{2})$ , while the normalization of the states is

$$\langle n|m\rangle = (-)^n \delta_{n,m}. \quad (2.24)$$

Thus all odd eigenstates have negative norm. Using (2.19), (2.21), and (2.22), we can compute the two-point function:

$$\begin{aligned} \langle 0|T\hat{q}(t)\hat{q}(0)|0\rangle &= -(\hbar/2\omega)(\theta(t)e^{-i\omega t} + \theta(-t)e^{i\omega t}) \\ &= -\frac{i\hbar}{2\pi} \int_{-\infty}^{\infty} dk \frac{e^{-ikt}}{k^2 - \omega^2 + i\epsilon}. \end{aligned} \quad (2.25)$$

Again, the residue is negative, but the  $i\epsilon$  prescription is now the same as for the positive-energy harmonic oscillator.

Summarizing and comparing the two different quantum theories of the negative-energy oscillator, we note that they are related formally by the correspondence (2.18). However they differ in the choice of groundstates used; namely, the state  $|0\rangle$  is annihilated by the operator  $\hat{a}$ , while the state  $|0\rangle$  is annihilated by  $\hat{b} = \hat{a}^\dagger$ . Because  $|0\rangle \neq |0\rangle$ , the two procedures correspond to different dynamical realizations of the system. In this respect one may think of the relation (2.18) as a kind of Bogoliubov transformation between the groundstates  $|0\rangle, |0\rangle$ . An order parameter can be introduced as  $\langle \hat{b}^\dagger \hat{b} \rangle$ , which is positive in the state  $|0\rangle$ , and vanishes in the other one,  $|0\rangle$ .

### III. REALIZATION OF NONCOMPACT SYMMETRIES IN QUANTUM SYSTEMS

In order to study the interpretation of the two quantization schemes in the context of systems with noncompact symmetry, it suffices to consider the simple example of two oscillators with opposite energy and an  $O(1,1)$  symmetry:

$$L = \frac{1}{2}(\dot{x}_1^2 - \dot{x}_2^2) - \frac{1}{2}\omega^2(x_1^2 - x_2^2). \quad (3.1)$$

This is just the linear  $O(1,1)$  sigma model in  $(0+1)$  dimensions, invariant under the infinitesimal  $O(1,1)$  transformations

$$\delta x_1 = \alpha x_2, \quad \delta x_2 = \alpha x_1. \quad (3.2)$$

If the quantum system is defined using the positive-definite Hilbert space with negative-energy states, the operators describing the two degrees of freedom of the system are

$$\begin{aligned} \hat{x}_1(t) &= (\hbar/2\omega)^{1/2}(\hat{a}_1 e^{-i\omega t} + \hat{a}_1^\dagger e^{i\omega t}), \\ \hat{x}_2(t) &= (\hbar/2\omega)^{1/2}(\hat{a}_2 e^{i\omega t} + \hat{a}_2^\dagger e^{-i\omega t}). \end{aligned} \quad (3.3)$$

The raising and lowering operators  $\hat{a}_i^\dagger, \hat{a}_i$  satisfy the equal-time commutation relations

$$[\hat{a}_i, \hat{a}_j^\dagger] = \delta_{ij}, \quad (i, j) = (1, 2). \quad (3.4)$$

The Hamiltonian is

$$\hat{H} = \hbar\omega(\hat{a}_1^\dagger \hat{a}_1 - \hat{a}_2^\dagger \hat{a}_2). \quad (3.5)$$

The noncompact symmetry is realized on the creation and annihilation operators by

$$\delta\hat{a}_1 = \alpha\hat{a}_2^\dagger, \quad \delta\hat{a}_2 = \alpha\hat{a}_1^\dagger. \quad (3.6)$$

One can show without difficulty, that this leaves the commutator (3.4) invariant. The transformations (3.6) are generated by a charge

$$\hat{Q} = \hat{a}_1 \hat{a}_2 - \hat{a}_1^\dagger \hat{a}_2^\dagger = -\hat{Q}^\dagger. \quad (3.7)$$

Its commutators with the  $\hat{a}_i, \hat{a}_i^\dagger$  are

$$[\hat{Q}, \hat{a}_1] = \hat{a}_2^\dagger, \quad [\hat{Q}, \hat{a}_2] = \hat{a}_1^\dagger. \quad (3.8)$$

The  $O(1,1)$  charge  $\hat{Q}$  commutes with the Hamiltonian:

$$[\hat{Q}, \hat{H}] = 0. \quad (3.9)$$

Hence the Hamiltonian is invariant and the charge conserved.

We define the groundstate as the state with occupation numbers (0,0) (it is more appropriate to call this a vacuum state, since it is not the state of lowest energy):

$$\hat{a}_i |0\rangle = 0, \quad \langle 0|0\rangle = 1. \quad (3.10)$$

Since the two oscillators do not interact, no energy can be exchanged between the positive and negative-energy subsystems. Therefore the vacuum (3.10) is stable under perturbations that couple only to the positive-energy degree of freedom (as one expects when the negative-energy degree of freedom is unphysical and can be eliminated by imposing suitable constraints).

It is immediately evident, that the vacuum (3.10) is not  $O(1,1)$  invariant:

$$\hat{Q} |0\rangle = -|1;1\rangle. \quad (3.11)$$

Here  $|n;m\rangle$  denotes the state with occupation numbers  $(n,m)$  for the oscillators (1,2). By continuing to apply  $\hat{Q}$ , we can actually construct an infinite multiplet of states with the same (vanishing) energy. This is expected from the theory of unitary representations of noncompact groups.<sup>19,20</sup> It is easily explained here by observing that the commutator of  $\hat{Q}$  with the occupation number operators  $\hat{n}_i = \hat{a}_i^\dagger \hat{a}_i$  is nonvanishing, even when acting on the vacuum. We wish to stress that the noninvariance of the vacuum (3.10) is not a result of diagonalizing a wrong set of observables in searching for a ground state. Attempts to construct an  $O(1,1)$  invariant state which does not have well-defined occupation numbers,

$$\Psi = \sum_{n,m} c_{n,m} |n;m\rangle \quad \text{with} \quad \hat{Q}\Psi = 0, \quad (3.12)$$

fail, because such states turn out to have infinite norm. Thus we have a conserved charge  $\hat{Q}$  (commuting with the Hamiltonian), but no invariant states. In this sense the noncompact symmetry is broken in this realization of the quantum  $O(1,1)$  model.

The observations made for this simple system can be extended to the case of noncompact field theories as well.

Thus a quantum noncompact field theory with a positive-definite Hilbert space does not have an invariant vacuum, although the charges connected with the noncompact symmetry may be conserved.

When the other quantization procedure is chosen, using a Hilbert space with indefinite metric, the dynamical variables of the theory can be expanded as

$$\begin{aligned} \hat{x}_1(t) &= (\hbar/2\omega)^{1/2} (\hat{a}_1 e^{-i\omega t} + \hat{a}_1^\dagger e^{i\omega t}), \\ \hat{x}_2(t) &= (\hbar/2\omega)^{1/2} (\hat{a}_2 e^{-i\omega t} + \hat{a}_2^\dagger e^{i\omega t}), \end{aligned} \quad (3.13)$$

where the previous  $\hat{b}$  operators are now denoted by  $\hat{a}_2, \hat{a}_2^\dagger$ , and the fundamental commutators are

$$[\hat{a}_i, \hat{a}_j^\dagger] = (\tau_3)_{ij}. \quad (3.14)$$

The Hamiltonian is

$$\hat{H} = \hbar\omega(\hat{a}_1^\dagger \hat{a}_1 - \hat{a}_2^\dagger \hat{a}_2 + 1). \quad (3.15)$$

In this case the noncompact  $O(1,1)$  symmetry is realized by

$$\delta\hat{a}_1 = \alpha\hat{a}_2, \quad \delta\hat{a}_2 = \alpha\hat{a}_1. \quad (3.16)$$

Again, the commutator (3.14) is invariant. The generator of the transformations is

$$\hat{Q} = \hat{a}_1^\dagger \hat{a}_2 - \hat{a}_2^\dagger \hat{a}_1. \quad (3.17)$$

It commutes with the occupation number operators only on states with occupation numbers (0,0). There is a unique state with this property; thus the ground state defined by

$$\hat{a}_i |\underline{0}\rangle = 0, \quad \langle \underline{0}|\underline{0}\rangle = 1, \quad (3.18)$$

is  $O(1,1)$  invariant:  $\hat{Q} |\underline{0}\rangle = 0$ .

This result holds more generally for systems with noncompact symmetries quantized in a Hilbert space with indefinite metric, which has only positive-energy states.<sup>18</sup>

Summarizing, the two quantization schemes for systems with noncompact symmetry correspond to a realization of the system in a positive-definite Hilbert space with negative-energy states and spontaneously broken noncompact symmetry, in the sense defined below Eq. (3.11); or alternatively a realization of the system in a Hilbert space with indefinite metric, but with strictly positive energy and an invariant vacuum.

#### IV. PATH INTEGRAL QUANTIZATION

As might be expected, the two different quantum realizations of theories with negative-energy modes also have different path-integral descriptions. The correct form of the path integral for the two quantum theories discussed in Secs. II and III can be established in various ways, for example, by discretizing time, from the holomorphic representation, or using the completeness relation for the eigenfunctions to compute the integral kernel for the Schrödinger equation. Here we do not give a derivation from first principles, but merely state the results. However, we compute the two-point functions and show that they reproduce those of Eqs. (2.17) and (2.25). For the case of noninteracting oscillators this amounts in fact to a complete proof of the results.

For the case of the negative-energy harmonic oscillator quantized in a positive definite Hilbert space the correct form of the path integral is

$$\begin{aligned}
Z[j] &= \langle 0, T \rightarrow \infty | 0, T \rightarrow -\infty \rangle_j \\
&= \lim_{\epsilon \rightarrow 0} \int Dq(t) \exp \frac{i}{\hbar} \int_{-\infty}^{\infty} \left( L + \frac{i}{2} \epsilon q^2 + jq \right) dt
\end{aligned} \tag{4.1}$$

with  $L$  as in Eq. (2.1). The integration over the  $q(t)$  runs from  $-\infty$  to  $+\infty$  along the real axis. Strictly speaking, the  $i\epsilon$  term is not necessary for convergence of the integral, but it is included to indicate the proper analytic continuation of the frequency  $\omega$  to complex values. Such a continuation allows a smooth limit to the case of zero mass ( $\text{Re } \omega = 0$ ), which needs extra regularization (the determinant develops a zero mode). In fact the integral is well defined for all complex values of  $\omega$  such that  $\text{Im } \omega \geq 0$ ; however it does not exist for  $\text{Im } \omega < 0$ . [Of course, the determinant does exist for  $\text{Im } \omega < 0$ ; however, for such values it cannot be represented by the integral (4.1).] The analytic continuation made in (4.1) also defines the correct time ordering, as may be seen by computing the two-point function and comparing it to (2.17). This computation is facilitated by the inclusion of the external source term  $j(t)q(t)$ : the propagator is now obtained in standard fashion from the quadratic term in the expansion of  $\log Z[j]$  in terms of the sources  $j(t)$ :

$$\begin{aligned}
\langle q(t)q(0) \rangle &= \frac{-\hbar^2}{Z[0]} \frac{\delta^2 Z[j]}{\delta j(t)\delta j(0)} \Big|_{j=0} \\
&= \frac{-i\hbar}{2\pi} \int_{-\infty}^{\infty} dk \frac{e^{-ikt}}{k^2 - \omega^2 - i\epsilon}.
\end{aligned} \tag{4.2}$$

This is indeed the same as in Eq. (2.17), including the  $i\epsilon$  prescription. From the derivation of (2.17) it can be seen that the  $i\epsilon$  prescription determines the  $\theta(t)$  functions necessary for obtaining the correct time ordering of the product  $\hat{q}(t)\hat{q}(0)$ .

In contradistinction to this, the description of the system in a negative-norm Hilbert space is recovered from the path integral

$$\begin{aligned}
\bar{Z}[j] &= \langle \bar{0}, T \rightarrow \infty | \bar{0}, T \rightarrow -\infty \rangle \\
&= \lim_{\epsilon \rightarrow 0} \int Dq(t) \exp \frac{i}{\hbar} \int_{-\infty}^{\infty} \left( L - \frac{i\epsilon}{2} q^2 + jq \right) dt,
\end{aligned} \tag{4.3}$$

where the domain of integration of  $q(t)$  is now from  $-i\infty$  to  $+i\infty$ . The reason for this is, that in the negative-norm space the operators  $\hat{q}$  have imaginary eigenvalues.<sup>21</sup> Here the  $i\epsilon$  prescription indicates convergence of the integrals for  $\text{Im } \omega \leq 0$ . It agrees with the time ordering in the canonical treatment (2.25), as follows from the two-point function:

$$\begin{aligned}
\langle q(t)q(0) \rangle &= - \frac{\hbar^2}{Z[0]} \frac{\delta^2 Z[j]}{\delta j(t)\delta j(0)} \Big|_{j=0} \\
&= - \frac{i\hbar}{2\pi} \int_{-\infty}^{\infty} dk \frac{e^{-ikt}}{k^2 - \omega^2 + i\epsilon}.
\end{aligned} \tag{4.4}$$

Comparison with the expression (2.25) shows this to be correct indeed.

## V. NONCOMPACT $\sigma$ MODELS

In the previous sections we have shown that there are two alternative methods of quantizing theories with non-

compact symmetries, differing in the way the noncompact symmetry is realized in the spectrum. Now we must determine which quantization procedure to follow in actual physical theories. For Maxwell's theory of the electromagnetic field or the theory of the relativistic string the noncompact symmetry involved is Lorentz invariance. Because this symmetry is not broken, these theories must be quantized in a Hilbert space with indefinite metric. This leads to the standard Gupta-Bleuler formulation of quantum electrodynamics.<sup>22</sup>

The situation is different in the case of noncompact  $\sigma$  models. Below we present a number of arguments showing that in these theories the noncompact symmetry is broken (in the sense of Sec. III). Therefore we are forced to conclude that the quantum description requires one to work with an extended Hilbert space having positive-definite metric and negative-energy states.

An explicit proof of this statement can be given for a number of exactly solvable models. The first of these is the nonlinear  $O(1,1)$  model. This model is a nonlinear counterpart of the model that was briefly discussed in its  $(0+1)$ -dimensional version in Sec. III. In line with the general discussion of nonlinear  $O(N,1)$  models in Eqs. (1.6)–(1.12) it is defined by the Lagrangian

$$\mathcal{L} = \frac{1}{2}(\partial_\mu \sigma)^2 - \frac{1}{2}(\partial_\mu \pi)^2, \tag{5.1}$$

with the constraint

$$\sigma^2 - \pi^2 = 1/g > 0. \tag{5.2}$$

One may now proceed to solve the constraint as in Eq. (1.8), in terms of a single scalar field  $\varphi$ :

$$\begin{aligned}
\pi &= \frac{\varphi}{\sqrt{1-g\varphi^2}}, \quad \sigma = \frac{1}{\sqrt{g}} \frac{1}{\sqrt{1-g\varphi^2}}, \\
\mathcal{L} &= -\frac{1}{2}(\partial_\mu \varphi)^2 / (1-g\varphi^2)^2.
\end{aligned} \tag{5.3}$$

This is invariant under the nonlinear Abelian transformations

$$\delta\varphi = (1/\sqrt{g})b(1-g\varphi^2). \tag{5.4}$$

However, in this case there exists a more convenient parametrization of the model, using a scalar field  $\vartheta$  defined by

$$\pi = (1/\sqrt{g})\sinh(\sqrt{g}\vartheta), \quad \sigma = (1/\sqrt{g})\cosh(\sqrt{g}\vartheta). \tag{5.5}$$

In terms of this field  $\vartheta$  the Lagrangian becomes

$$\mathcal{L} = -\frac{1}{2}(\partial_\mu \vartheta)^2. \tag{5.6}$$

Thus the theory describes a single real, massless scalar field, and is invariant under nonlinear  $O(1,1)$  transformations of the form

$$\vartheta \rightarrow \vartheta' = \vartheta + \eta, \tag{5.7}$$

where  $\eta$  is a constant.

The importance of this result is, that the spectrum of a massless field is well known to be continuous. In particular, there is no normalizable ground state. This may be checked directly, by enclosing the system in a box of finite dimension  $L$  and decomposing the field into plane waves:

$$\vartheta(x) = \sum_{\mathbf{k} \in \mathbb{Z}^d} \vartheta_{\mathbf{k}}(t) e^{2\pi i \mathbf{k} \cdot \mathbf{x} / L}, \quad \vartheta_{\mathbf{k}}^* = \vartheta_{-\mathbf{k}}. \tag{5.8}$$

The Hamiltonian separates into a direct sum of Hamilto-

nians for the modes (5.8). Correspondingly, the wave function factorizes into an infinite product of wave functions, one for each mode. For all the modes with space momentum  $\mathbf{k} \neq 0$ , the Hamiltonian is that of a harmonic oscillator with frequency  $\omega_{\mathbf{k}} = (2\pi/L)|\mathbf{k}|$ . Thus these modes have normalizable wave functions  $\psi_{\mathbf{k}}$ . For the zero mode  $\vartheta_0$  the Hamiltonian is that of a free, nonrelativistic particle in one dimension. Its eigenstates are plane waves, which are not normalizable. To get normalizable states, we have to form wave packets by linear superposition of the zero modes. Therefore the wave function for the free, massless scalar field has the form

$$\Psi[\vartheta_{\mathbf{k}}] = \int_{-\infty}^{\infty} d\alpha c(\alpha) e^{i\alpha\vartheta_0} \prod_{\mathbf{k} \neq 0} \psi_{\mathbf{k}}[\vartheta_{\mathbf{k}}], \quad (5.9)$$

where

$$\frac{1}{2\pi} \int_{-\infty}^{\infty} d\alpha |c(\alpha)|^2 = 1.$$

The nonlinear  $O(1,1)$  transformation (5.7) shifts the zero mode by  $\eta$ , but does not affect the other modes:

$$\delta\vartheta_0 = \eta; \quad \delta\vartheta_{\mathbf{k}} = 0, \quad \mathbf{k} \neq 0. \quad (5.10)$$

The result of this transformation on the wave function (5.9) is that all  $c(\alpha)$  are multiplied by a phase factor  $\exp(i\alpha\eta)$ . This is a unitary transformation, preserving the normalization, but changing the form of the superposition and hence the state of the system. In fact, the plane wave basis for the zero modes is an infinite-dimensional unitary representation of the Abelian  $O(1,1)$  group. Therefore we have no normalizable  $O(1,1)$ -invariant ground state for the nonlinear  $O(1,1)$  model.

The conclusion arrived at for the Abelian  $O(1,1)$  model also holds for the non-Abelian  $O(N,1)$  models. This has been proved in detail for the  $O(3,1)$  model in  $(0+1)$  dimensions by Velo and Wess,<sup>6</sup> and for the supersymmetric  $O(2,1)$  model in  $(0+1)$  dimensions by Davis *et al.*<sup>3</sup> In this last model there are in fact normalizable ground states with energy zero, separated by a finite energy gap from the rest of the spectrum, which is purely continuous. However, these zero energy states form two infinite-dimensional discrete series, which are unitary representations of  $O(2,1)$  (Refs. 23, 24). Therefore there is still no  $O(2,1)$  invariant vacuum. Such discrete series arise in fact for all supersymmetric  $O(N,1)$  models with  $N$  even.<sup>4</sup> Therefore none of these models possesses an  $O(N,1)$  invariant vacuum. That this situation generalizes to the case of the supersymmetric  $O(N,1)$  field theories in  $(d+1)$  dimensions ( $d \geq 1$ ) is shown by using Wittens index theorem.<sup>23</sup>

Before proceeding to discuss more general arguments indicating the breaking of the noncompact symmetry in the nonlinear  $O(N,1)$  models, we would like to return once more to the  $O(1,1)$  model and discuss the construction of its Hilbert space of states in the formalism with constraints, as in (5.1) and (5.2). We consider the model in  $(0+1)$  dimensions for simplicity, but this restriction is not essential. Before imposing the constraint, the wave functions of the system in the extended Hilbert space with positive-definite metric (and negative energies) are

$$\psi(\sigma, \pi) = \int_{-\infty}^{\infty} \frac{d\alpha}{2\pi} \int_{-\infty}^{\infty} \frac{d\eta}{2\pi} c(\alpha, \eta) e^{i(\alpha\pi - \eta\sigma)}. \quad (5.11)$$

We can now impose the constraint (5.2) as a constraint on the wave functions

$$(\sigma^2 - \pi^2 - (1/g))\psi(\sigma, \pi) = 0. \quad (5.12)$$

This constraint is satisfied if the Fourier coefficients  $c(\alpha, \eta)$  satisfy the  $(1+1)$ -dimensional Klein-Gordon equation:

$$(-\partial_{\eta}^2 + \partial_{\alpha}^2 - 1/g)c(\alpha, \eta) = 0. \quad (5.13)$$

The solutions of this equation are of the form

$$c(\alpha, \eta) = \int_{-\infty}^{\infty} d\pi \int_{-\infty}^{\infty} d\sigma \delta(\sigma^2 - \pi^2 - (1/g)) \times \varphi(\sigma, \pi) e^{-i(\alpha\pi - \eta\sigma)}. \quad (5.14)$$

Taking  $\eta = 0$  this reduces to the plane wave basis (5.9). Thus we can indeed obtain the correct physical states by imposing the constraint (5.2) on the states of the extended Hilbert space with positive-definite metric and negative-energy modes.

With the nonlinear transformation law of the physical fields under the  $O(N,1)$  transformations (1.12), it is actually not surprising, that the  $O(N,1)$  symmetry is broken to its compact  $O(N)$  subgroup in the quantum nonlinear  $\sigma$  models. Assuming the quantum theory to exist, it follows from the Goldstone theorem. Denoting the compact  $O(N)$  generators by  $\hat{R}$ , and the remaining noncompact ones by  $\hat{N}$ , we can show that the vacuum expectation value of the massless physical fields  $\langle \varphi_i \rangle$  vanishes because of the  $O(N)$  invariance of the vacuum. Consider an infinitesimal  $O(N)$  transformation of  $\langle \varphi_i \rangle$ :

$$\delta(R)\langle \varphi_i \rangle = a_{ij} \langle \varphi_j \rangle = \langle 0 | [\hat{R}(a), \varphi_i] | 0 \rangle = 0. \quad (5.15)$$

The vanishing of this expression follows, since the compact  $O(N)$  symmetry of the vacuum may safely be assumed in the absence of a symmetry breaking potential. However, for a transformation in a noncompact direction we obtain

$$\begin{aligned} \delta(N)\langle \varphi_i \rangle &= \langle 0 | [\hat{N}(b), \varphi_i] | 0 \rangle \\ &= (1/\sqrt{g})b_j (\delta_{ij} + g[\delta_{ij} \langle \varphi_k^2 \rangle - 2\langle \varphi_i \varphi_j \rangle]). \end{aligned} \quad (5.16)$$

For  $N \geq 2$  the right-hand side cannot vanish, proving that the vacuum cannot be invariant under these transformations. For  $N = 1$  we have already proved the nonexistence of an invariant vacuum. We conclude that, provided the quantum theory is well defined, the vacuum of the nonlinear  $O(N,1)$   $\sigma$  model is not invariant under the noncompact symmetries. Finally we briefly consider the constraint (1.7). Taking vacuum expectation values on both sides of the equation, we get after proper renormalization:

$$\langle \sigma^2 \rangle_{\text{ren}} - \langle \pi^2 \rangle_{\text{ren}} = 1/g_{\text{ren}} > 0. \quad (5.17)$$

If the renormalized coupling constant is positive (and finite), we are forced to conclude that

$$\langle \sigma^2 \rangle_{\text{ren}} \neq \langle \pi^2 \rangle_{\text{ren}} \quad \text{and} \quad \langle \sigma^2 \rangle_{\text{ren}} > 0. \quad (5.18)$$

The first condition can be avoided only if both vacuum expectation values become infinite:

$$(\langle \sigma^2 \rangle_{\text{ren}}, \langle \pi^2 \rangle_{\text{ren}}) \rightarrow \infty.$$

This possibility is presumably realized in the  $(1+1)$ -dimensional models,<sup>4,10</sup> since for  $d=1$  no properly massless scalars exist.<sup>24</sup> However, for  $d \geq 2$ , model calculations and the analysis of the  $O(1,1)$  model above indicate that the vacuum expectation values are finite<sup>4,9,11</sup> and the conditions (5.18) must be satisfied. Again we are led to conclude that the  $O(N,1)$  symmetry is broken.

The arguments presented above indicate, that for the noncompact nonlinear  $\sigma$  models the quantization in an extended Hilbert space for constrained fields  $(\sigma, \pi)$  should proceed with positive norm, but negative-energy states for the unphysical degrees of freedom  $\sigma$ . The energy of the full theory is nevertheless bounded below, because the constraint (1.7) [(5.2)] guarantees, that the excitation of negative-energy states can only take place, if simultaneously sufficiently many positive-energy modes of the  $\pi$  fields are excited as well.

## VI. A MODEL COMPUTATION

The conclusion of Sec. V, that in the formulation with constraints noncompact  $\sigma$  models are to be quantized using a positive-definite Hilbert space with negative-energy states for the unphysical degrees of freedom, implies that the free propagators of the fields  $(\sigma, \pi)$  used in perturbation theory take the form

$$\begin{aligned} \langle \sigma(x) \sigma(0) \rangle &= - \int \frac{d^n k}{(2\pi)^n} \frac{e^{ik \cdot x}}{k^2 + m^2 + i\epsilon}, \\ \langle \pi_i(x) \pi(0) \rangle &= \delta_{ij} \int \frac{d^n k}{(2\pi)^n} \frac{e^{ik \cdot x}}{k^2 + m^2 - i\epsilon}, \end{aligned} \quad (6.1)$$

where  $n$  denotes the dimensionality of space-time; cf. Eqs. (2.17) and (4.2). In this section we show that this prescription leads to correct results in a simple toy-model computation. Although this model is strongly simplified, we believe it has all the necessary features to test the analytic continuation of  $m^2$  implicit in Eqs. (6.1), as it appears for example in path integrals of the type (4.1) and (4.3). Instead of these infinite-dimensional integrals, we consider the ordinary double integral

$$\begin{aligned} Z_-(m_1^2, m_2^2) &= \int_{-\infty}^{\infty} d\pi \int_{-\infty}^{\infty} d\sigma \delta\left(\pi^2 - \sigma^2 + \frac{1}{g}\right) \\ &\times e^{-i/2(m_1^2 \pi^2 - m_2^2 \sigma^2)}, \end{aligned} \quad (6.2)$$

where for definiteness we choose  $m_2^2 \geq m_1^2 \geq 0$  and  $g > 0$ . This is a zero-dimensional version of the nonlinear  $O(1,1)$  model, with an additional  $O(1,1)$  breaking for  $m_1^2 \neq m_2^2$ . Inspection of the integral shows it to be well defined for  $\text{Im } m_1^2 \leq \text{Im } m_2^2$ . The integral can be evaluated directly by first carrying out the integration over  $\sigma$ ; this eliminates the  $\delta$  function:

$$\begin{aligned} Z_-(m_1^2, m_2^2) &= \int_{-\infty}^{\infty} d\pi \int_{-\infty}^{\infty} d\sigma \int_{-\infty}^{\infty} \frac{d\alpha}{4\pi} \exp\left\{-\frac{i}{2}\left[m_1^2 \pi^2 - m_2^2 \sigma^2 + \alpha\left(\pi^2 - \sigma^2 + \frac{1}{g}\right)\right]\right\} \\ &= \int_{-\infty}^{\infty} d\pi \int_{-\infty}^{\infty} d\sigma \int_{-\infty}^{\infty} \frac{d\alpha}{4\pi} \exp\left\{-\frac{i}{2}\left[(m_1^2 + \alpha)\pi^2 - (m_2^2 + \alpha)\sigma^2 + \frac{\alpha}{g}\right]\right\}. \end{aligned} \quad (6.10)$$

The perturbation theory treatment is obtained by interchanging the order of integration, allowing the integrations over  $(\sigma, \pi)$

$$Z_-(m_1^2, m_2^2) = 2\sqrt{g} e^{im_2^2/2g} \int_0^{\infty} d\pi \frac{e^{(i/2)(m_2^2 - m_1^2)\pi^2}}{\sqrt{1 + g\pi^2}}. \quad (6.3)$$

It is convenient to redefine the parameters as follows:

$$\begin{aligned} \mu^2 &= (1/4g)(m_1^2 + m_2^2), \\ \omega &= (1/4g)(m_2^2 - m_1^2), \quad s = 1 + 2g\pi^2. \end{aligned} \quad (6.4)$$

Then

$$Z_-(m_1^2, m_2^2) = e^{i\mu^2} \int_1^{\infty} ds \frac{e^{i\omega s}}{\sqrt{s^2 - 1}}. \quad (6.5)$$

This integral is defined for  $\text{Im } \omega \geq 0$  ( $\omega \neq 0$ ) and all complex values of  $\mu^2$ , in agreement with the complex domains of  $m_1^2, m_2^2$ . The result of (6.5) is a Hankel function of the first kind:

$$Z_-(m_1^2, m_2^2) = (i\pi/2) e^{i\mu^2} H_0^{(1)}(\omega). \quad (6.6)$$

Of course, the Hankel function can be defined in the whole complex  $\omega$  plane, but the integral representation (6.6) exists only for  $\text{Im } \omega \geq 0$  (Ref. 24).

It is instructive to compare this result with the compact  $O(2)$  version of this integral:

$$\begin{aligned} Z_+(m_1^2, m_2^2) &= \int_{-\infty}^{\infty} d\pi \int_{-\infty}^{\infty} d\sigma \delta\left(\sigma^2 + \pi^2 - \frac{1}{g}\right) \\ &\times \exp\left(-\frac{i}{2}(m_1^2 \pi^2 + m_2^2 \sigma^2)\right) \\ &= 2\sqrt{g} e^{-im_2^2/2g} \int_0^{1/\sqrt{g}} d\pi \frac{\exp[(i/2)(m_2^2 - m_1^2)\pi^2]}{\sqrt{1 - g\pi^2}}, \end{aligned} \quad (6.7)$$

which exists for all complex values of  $m_1^2, m_2^2$ . With  $\mu^2$  and  $\omega$  as in (6.4) and taking

$$s = -1 + 2g\pi^2,$$

the integral (6.7) becomes

$$Z_+(m_1^2, m_2^2) = e^{-i\mu^2} \int_{-1}^{+1} ds \frac{e^{i\omega s}}{\sqrt{1 - s^2}}. \quad (6.8)$$

In contrast to (6.5), this integral is indeed well defined for all complex values of  $\omega$ . It is an integral representation of the Bessel function of order zero:

$$Z_+(m_1^2, m_2^2) = \pi e^{-i\mu^2} J_0(\omega). \quad (6.9)$$

In actual quantum-field theoretical models, the path integrals usually cannot be evaluated directly and one must take recourse to perturbation theory. The constraint for  $(\sigma, \pi)$  is then imposed using a Lagrange multiplier  $\alpha$ . For the integral (6.2) such a formulation is obtained by substitution of the Fourier representation of the  $\delta$  function:



to be performed first and that over  $\alpha$  afterwards. For the integral (6.10) this interchange is allowed for  $\text{Im } m_1^2 \leq 0$  and  $\text{Im } m_2^2 \geq 0$ , and gives

$$Z_-(m_1^2, m_2^2) = \int_{-\infty}^{\infty} \frac{d\alpha}{4\pi} \exp\left(-\frac{i\alpha}{2g}\right) \left\{ \int_{-\infty}^{\infty} d\pi \int_{-\infty}^{\infty} d\sigma \exp\left(-\frac{i}{2}[(m_1^2 + \alpha)\pi^2 - (m_2^2 + \alpha)\sigma^2]\right) \right\} \\ = \frac{1}{2} \int_{-\infty}^{\infty} d\alpha \frac{\exp(-i\alpha/2g)}{\sqrt{(m_1^2 + \alpha)(m_2^2 + \alpha)}}. \quad (6.11)$$

Now define

$$t = -((\alpha/2g) + \mu^2).$$

Then

$$Z_-(m_1^2, m_2^2) = \frac{1}{2} e^{i\mu^2} \int_{-\infty}^{\infty} dt \frac{e^{it}}{\sqrt{t^2 - \omega^2}}. \quad (6.12)$$

A similar treatment of the compact model (6.7) requires  $\text{Im } m_1^2 \leq 0$ ,  $\text{Im } m_2^2 \leq 0$ . In order to show that this reduces to the form (6.5), we must take care of the multivaluedness of the square root, and deal with the singularities at  $t = \pm \omega$ . The square root becomes single valued if we make cuts along the half-lines  $[\omega, \infty)$  and  $(-\infty, -\omega]$ . For  $\text{Im } \omega > 0$  the integration does not encounter the cuts (Fig. 1). In the limit  $\text{Im } \omega \rightarrow 0$ , the path of integration passes along the upper edge of the cut  $(-\infty, -\omega]$  and the lower edge of the cut  $[\omega, \infty)$ . Now consider the contour  $C$  of Fig. 2. Since there are no singularities inside the contour, we have

$$\oint_C dz \frac{e^{iz}}{\sqrt{z^2 - \omega^2}} = 0. \quad (6.13)$$

In the limit  $R \rightarrow \infty$  the part of the integral along the semicircular arc vanishes. Therefore

$$\int_{-\infty}^{\text{Re } \omega} dt \frac{e^{it}}{\sqrt{t^2 - \omega^2}} = - \int_{\text{Re } \omega = \omega - i \text{Im } \omega}^{\omega + i \text{Im } \omega} dt \frac{e^{it}}{\sqrt{t^2 - \omega^2}} \\ - \int_{\omega + i \text{Im } \omega}^{\infty + 2i \text{Im } \omega} dt \frac{e^{it}}{\sqrt{t^2 - \omega^2}}. \quad (6.14)$$

Hence the contour of Fig. 1 may be deformed to pass along the upper and lower edges of the cut  $[\omega, \infty)$ , as in Fig. 3. Finally, the two parts above and below the cut contribute equal amounts, because the change in sign of the square root is balanced by the reversal of the direction of the path. The final result is therefore

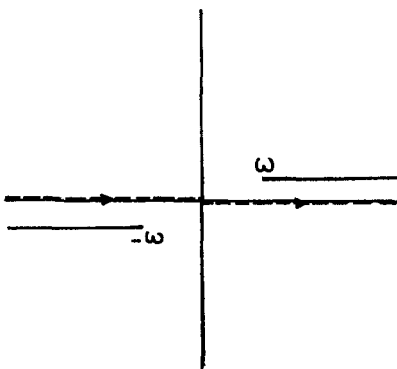


FIG. 1. Path of integration in complex  $t$  plane for Eq. (6.12).

$$Z_-(m_1^2, m_2^2) = e^{i\mu^2} \int_{\omega}^{\infty} dt \frac{e^{it}}{\sqrt{t^2 - \omega^2}} \\ = e^{i\mu^2} \int_1^{\infty} ds \frac{e^{i\omega s}}{\sqrt{s^2 - 1}}. \quad (6.15)$$

Thus the result (6.5) is reproduced and the representation of the integral by (6.11) is correct with  $\text{Im } m_1^2 \leq 0$ ,  $\text{Im } m_2^2 \geq 0$ . Writing explicitly the real and imaginary parts as  $m_{1,2}^2 \rightarrow m_{1,2}^2 \pm i\epsilon_{1,2}$ , we can replace (6.11) in the limit  $\epsilon_{1,2} \downarrow 0$  (such that the logarithm is single valued) by

$$Z_-(m_1^2, m_2^2) = \int_{-\infty}^{\infty} \frac{d\alpha}{4\pi} \exp\left(-\frac{i\alpha}{2g}\right) \\ \times \left\{ \exp\left(-\frac{1}{2} \ln(m_1^2 + \alpha - i\epsilon_1)\right) \right. \\ \left. - \frac{1}{2} \ln(m_2^2 + \alpha + i\epsilon_2) \right\}. \quad (6.16)$$

This is the form that, generalized to higher-dimensional models, is the starting point for a perturbative analysis. It is now straightforward to establish the results

$$\langle \sigma^2 \rangle = -2i \frac{\partial}{\partial m_2^2} Z_-(m_1^2, m_2^2) = \left\langle \frac{i}{m_2^2 + \alpha + i\epsilon_2} \right\rangle_{\alpha}, \\ \langle \pi^2 \rangle = 2i \frac{\partial}{\partial m_1^2} Z_-(m_1^2, m_2^2) = \left\langle \frac{-i}{m_1^2 + \alpha - i\epsilon_1} \right\rangle_{\alpha}, \quad (6.17)$$

which are identical to Eqs. (6.1) for  $n = 0$ . It is obvious that the replacement of  $(\pi, \sigma)$  by fields  $(\pi(x), \sigma(x))$  and  $m_{1,2}^2$  by differential operators  $\partial_{\mu}^2 + m_{1,2}^2$  does not alter these results. Hence the analytic continuation chosen for the propagators in Eq. (6.1) is consistent for the noncompact models and leads to a correct perturbative treatment of the quantum field theory,<sup>11</sup> provided that this theory is well defined.

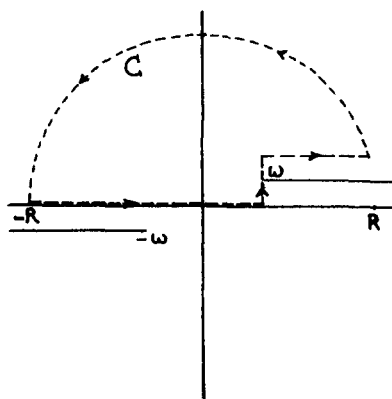


FIG. 2. Closed contour in  $z$  plane, Eq. (6.13).

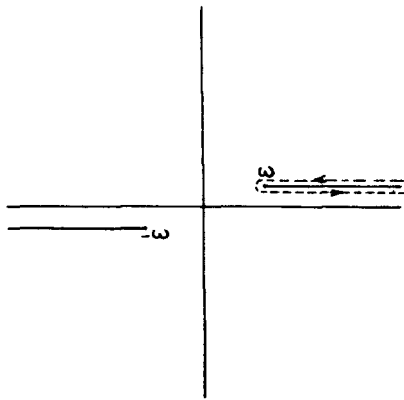


FIG. 3. Equivalent path of integration in  $t$  plane for  $Z_-(m_1^2, m_2^2)$ .

## VII. DISCUSSION

The argument presented in this paper may be summarized as follows. There exist two different quantization schemes for field theories with noncompact symmetry of the kinetic energy operator: one with the noncompact symmetry broken, but a positive-definite Hilbert space and one with the noncompact symmetry not broken, but being realized in a Hilbert space of indefinite metric. It is argued from the examples of some explicitly solvable models that in nonlinear  $\sigma$  models with noncompact symmetry the noncompact symmetry is broken in the sense that there is no invariant vacuum state, even though the charges of the noncompact symmetry are conserved. The general validity of this result in nonlinear  $\sigma$  models is supported by other, more model-independent arguments. From this it is inferred that the quantization of the linearized formulation of such theories in an extended Hilbert space with constraints requires the use of a positive-definite Hilbert space with negative-energy modes for the unphysical degrees of freedom.

The question of whether such a quantization is consistent with causality and stability is far from obvious. As we have argued in Secs. I and V, the constraint (1.7) [(5.2)] is crucial in establishing the boundedness of the Hamiltonian and therefore it must necessarily be conserved in the full quantum theory. This requires the quantum theory including the constraint to be renormalizable, since otherwise counterterms might modify the constraint in such a way that the negative-energy modes can no longer be removed from the theory. Obviously, the renormalization must be carried out very carefully; for example, mixing bare and renormalized quantities in the analysis may easily lead to inconsistencies.<sup>12</sup>

Most studies have been done on  $(0+1)$ - and  $(1+1)$ -dimensional models. The  $d=0$  models can be solved exactly and confirm the analysis presented here.<sup>3,4,6</sup> In  $d=1$  exact results can be obtained in the large- $N$  limit of  $O(N,1)$  and noncompact  $CP(N,1)$  models.<sup>4,5,7-12</sup> However, in this case the precise quantization of the negative-energy modes is not so important, because most of its contributions disappear in the limit  $N \rightarrow \infty$ . Nevertheless, negative-metric quantization has been shown to give rise to certain inconsistencies.<sup>12</sup> Also, as explained in Sec. V, the behavior of the noncompact mod-

els in  $d=1$  is somewhat peculiar<sup>4,10</sup> because of the problems with infrared divergences for massless scalars.<sup>25</sup>

In  $(3+1)$  dimensions a study has been made of the  $SU(1,1)/U(1)$  nonlinear  $\sigma$  model,<sup>11</sup> which describes the scalar sector of  $N=4$  supergravity.<sup>13</sup> The model has been analyzed at the one-loop level using the positive metric quantization. At this level it was shown to be renormalizable, with the energy bounded below.<sup>11</sup> These results can in fact be extended to the case of  $SU(N,N)/SU(N) \times SU(N) \times U(1)$  models. However, it is not known what happens when one includes higher loop contributions. Some arguments against consistency of the models at higher loops were coined by the authors of Ref. 12, but we believe their arguments to be inconclusive (Ref. 26). In particular, the derivation of cutting rules and the conditions of unitarity are modified by the use of negative-energy modes, since their propagator has a non-standard analytic continuation. Also, in the context of  $N$ -extended supergravity, it is known that the models are finite to two loops in perturbation theory.<sup>27</sup> Thus it seems that any inconsistencies are suppressed, at least till the next order of computation. This may be particularly relevant for gauged  $N=8$  supergravity,<sup>28</sup> where the dynamical generation of composite  $SU(8)$  gauge bosons may be possible by the mechanism of Ref. 11. This would greatly improve the phenomenological prospects for this theory as a unified model of electro-weak, strong, and gravitational interactions.

<sup>1</sup>E. M. Purcell and R. V. Pound, Phys. Rev. **81**, 279 (1951).

<sup>2</sup>K. Cahill, Phys. Rev. D **18**, 2930 (1978); B. Julia and J. F. Luciani, Phys. Lett. B **90**, 270 (1980).

<sup>3</sup>A. C. Davis, A. J. Macfarlane, and J. W. van Holten, Nucl. Phys. B **216**, 493 (1983).

<sup>4</sup>A. C. Davis, A. J. Macfarlane, and J. W. van Holten, Nucl. Phys. B **232**, 473 (1984); J. W. van Holten, Z. Phys. C **27**, 57 (1985).

<sup>5</sup>Y. Cohen and E. Rabinovici, Phys. Lett. B **124**, 371 (1983); Y. Ohnuki, preprint DPNU-86-21.

<sup>6</sup>G. Velo and J. Wess, Nuovo Cimento **1**, 177 (1971).

<sup>7</sup>M. Gomes and Y. K. Ha, Phys. Lett. B **145**, 235 (1984); Y. K. Ha B **256**, 687 (1985).

<sup>8</sup>N. Ohta, Phys. Lett. B **134**, 75 (1984).

<sup>9</sup>A. C. Davis, A. J. Macfarlane and J. W. van Holten, Phys. Lett. B **125**, 151 (1983).

<sup>10</sup>D. J. Amit and A. C. Davis, Nucl. Phys. B **225** (FS9), 221 (1983).

<sup>11</sup>J. W. van Holten, Nucl. Phys. B **242**, 307 (1984); Phys. Lett. B **135**, 427 (1984).

<sup>12</sup>A. C. Davis, M. D. Freeman, and A. J. Macfarlane, Nucl. Phys. B **258**, 393 (1985); T. Morozumi and S. Nojiri, Prog. Theor. Phys. **75** (1985).

<sup>13</sup>E. Cremmer, S. Ferrara, and J. Scherk, Phys. Lett. B **74**, 61 (1978); E. Cremmer and B. Julia, Nucl. Phys. B **159**, 141 (1979).

<sup>14</sup>Y. Nambu, *Proceedings of the International Conference on Symmetries and Quark Models* (Wayne State Univ., Detroit, MI, 1969); H. B. Nielsen, *15th Conference on High Energy Physics* (Naukova Dumka, Kiev, 1970); L. Susskind, Nuovo Cimento **69**, 457 (1970).

<sup>15</sup>L. Brink, P. Di Vecchia, and P. S. Howe, Phys. Lett. B **65**, 471 (1976); S. Deser and B. Zumino, Phys. Lett. B **65**, 369 (1976).

<sup>16</sup>R. C. Brower, Phys. Rev. D **6**, 1655 (1972); P. Goddard and C. B. Thorn, Phys. Lett. B **40**, 235 (1972).

<sup>17</sup>W. Heisenberg, Z. Phys. **90**, 209 (1934); J. R. Oppenheimer and W. Furry, Phys. Rev. **45**, 245 (1934).

<sup>18</sup>P. A. M. Dirac, Proc. R. Soc. London Ser. A **180**, 1 (1942); W. Pauli, Rev. Mod. Phys. **15**, 175 (1943).

<sup>19</sup>V. Bargmann, Ann. Math. **48**, 568 (1947); A. O. Barut and C. Fronsdal, Proc. R. Soc. London Ser. A: **287**, 532 (1965); W. J. Holman and L. C. Biedenharn, Ann. Phys. (NY) **39**, 1 (1966).

- <sup>20</sup>B. G. Wybourne, *Classical Groups for Physicists* (Wiley, New York, 1974), Chaps. 17 and 18.
- <sup>21</sup>H. Arisue, T. Fujiwara, T. Inoue, and K. Ogawa, *J. Math. Phys.* **22**, 2055 (1981).
- <sup>22</sup>S. N. Gupta, *Proc. Phys. Soc.* **63**, 68 (1950); K. Bleuler, *Helv. Phys. Acta* **23**, 567 (1950).
- <sup>23</sup>E. Witten, *Nucl. Phys. B* **202**, 253 (1982).
- <sup>24</sup>R. Courant and D. Hilbert, *Methoden der Mathematischen Physik I* (Springer, Berlin, 1968), Chap. VII.
- <sup>25</sup>S. Coleman, *Commun. Math. Phys.* **31**, 264 (1973).
- <sup>26</sup>I thank T. Morozumi and S. Nojiri for correspondence concerning this point.
- <sup>27</sup>M. T. Grisaru, *Phys. Lett. B* **66**, 75 (1977); M. T. Grisaru and P. van Nieuwenhuizen, *Deeper Pathways in High Energy Physics*, edited by A. Perlmutter and L. F. Scott (Plenum, New York, 1977).
- <sup>28</sup>B. de Wit and H. Nicolai, *Nucl. Phys. B* **208**, 323 (1982).

# Definition and general properties of a class of non-Markovian collision processes

M. Moreau

*Laboratoire de Physique Theorique des Liquides, Université Pierre et Marie Curie, 4, place Jussieu, 75252-Paris Cedex 05, France*

B. Gaveau

*UER de Mathématiques, Université Pierre et Marie Curie, 4, place Jussieu, 75252-Paris Cedex 05, France*

(Received 16 July 1986; accepted for publication 31 December 1986)

A class of non-Markovian stochastic processes is defined; it generalizes several anterior models, and allows a trajectorial description of collisions in dense fluids when the system loses all memory at each collision. The abstract formalism is studied here: an integral evolution equation is derived, as well as an integrodifferential equation that generalizes the master equation; some asymptotic properties of these processes are established. Applications to specific models will be treated in other articles.

## I. INTRODUCTION

Several authors<sup>1,2</sup> have recently developed stochastic models to describe the molecular relaxation and chemical reactions in a fluid. In most of these models the intermolecular potentials are supposed to have a very short range, compared to the average distance between two molecules (long range potentials may be simulated, in many cases, by mean field potential<sup>2,3</sup> and are not considered explicitly here). Then the interaction of a molecule, or of a reacting complex, with the other particles of the fluid may be represented as a collision localized in time and space; between two collisions, the molecule undergoes a determinist evolution under the influence of the external field, or effective external field.

Naturally, an exact probabilistic description of the fluid is obtained formally by applying the Liouville theorem to the complete density function, and writing the well-known hierarchy of equations for the reduced probability densities. However, this method does not yield explicit results easily, except in the case of the Boltzmann approximation, where multiple collisions are neglected.

Another method is to consider the evolution of a test molecule as a perturbation of its deterministic evolution under the action of collisions with other particles, represented as stochastic events; if the distribution laws in time and phase space of these collisions may be determined *a priori* then it is generally not difficult to study the stochastic evolution of the molecule. Although this approach is less satisfying from a theoretical point of view, since it implies a heuristic evaluation of the collision laws, it yields a simple description of the fluid and permits us to integrate experimental data or semiempirical deductions; furthermore, it appears to be very convenient for numerical simulation,<sup>3,4</sup> and has given valuable results in the theory of chemical reactivity in liquid phase. Thus it seems interesting to extend the validity of this method. In the works published on the subject,<sup>3,5</sup> three main assumptions are generally made: the collisions are instantaneous; the collision times of the test molecule are exponentially distributed: the collision rate at any time does not depend on the past collisions, nor on the state of the molecule; and in a collision the transition probabilities only depend on the initial and final states.

Of course these assumptions are related to the fundamental hypothesis that the evolution of the test molecule may be represented by a Markov process. However, it is clear that they can only be justified approximately and in special conditions; in particular, they are not appropriate in the case of a liquid medium: then, and generally, the evolution of the molecule is non-Markovian and cannot be treated by a standard formalism.

The difficulty is to define non-Markovian processes which are powerful enough to represent the physical phenomena, at least schematically, and not too general in order to be of practical use. For this purpose it is often assumed that the processes obey a generalized Langevin equation,<sup>6,7</sup> as it may be shown in the linear response approximation: but the memory function is hardly obtained from the first principles, and it must generally be approximated by a truncated continued fraction expansion, where the few generalized friction coefficients are chosen in order to fit experimental data, which is not very satisfying from a conceptual point of view. Among the various non-Markovian processes introduced by previous authors one may quote the semi-Markov processes defined by Feller<sup>8</sup> or the generalized random-walk processes used by Kenkre, Montroll, and Shlesinger<sup>9</sup>; the most interesting for physical purposes seems to be the composite stochastic processes studied by van Kampen<sup>10</sup>; however, their definition is not wide enough to be applied to the modelization of chemical reactions. Thus, using the mathematical formalism of stochastic processes, we introduced<sup>11</sup> the so-called collision processes which give a convenient generalization of the anterior works. Here we present this class of processes in a more physical way and study their main properties. Applications to reactional microdynamics are given in another paper.<sup>12</sup>

## II. NON-MARKOVIAN COLLISION PROCESSES

### A. The model

We consider a system, which may be a test molecule or group of molecules (a reacting complex, for instance) evolving under the action of internal and external forces. We assume that this evolution suffers different phases, labeled by

the index  $\epsilon$  running from 0 to  $m$ , which will be called "regimes"; for instance, regime 0 may be the free evolution of the system, or the evolution in an external field between two collisions; regime 1 may be the evolution during a collision with a particle of kind 1, etc. More generally  $\epsilon$  could be an element of a countable set  $\Sigma$  (or even of an uncountable set; but we will not consider such a general case in order to avoid technical difficulties which are unnecessary for application to concrete models).

During regime  $\epsilon$ , the instantaneous state of the system is represented by an element of a probability space  $E_\epsilon$  which may depend on  $\epsilon$ . Generally  $E_\epsilon$  is a  $n$ -dimensional Euclidean space  $R^n$  or a discrete set, or a manifold, and its dimension or its structure may change with  $\epsilon$ ; then the system is described at time  $t$  by a  $n$ -dimensional random vector  $X$ , and its probability law is defined by its probability density  $p_\epsilon(x, t)$ ,  $x \in R^n$ , or more generally by the probability measure  $P_\epsilon(A, t)$  which gives the probability to find the system at time  $t$  in some subset  $A$  of  $E_\epsilon$ . We have given in Ref. 12 three examples modeling chemical reactions in the presence of a solvent. In one of these examples, the state space  $E_\epsilon$  can be either one point space (for example a bound state), or a spatial interval with two velocities to take into account the diffusion on the top of a barrier potential.

Whatever may be the physical interpretation of regimes and state spaces, our fundamental hypothesis will be that *the process keeps no memory of the events prior to the beginning of a regime*, specified by its initial time and initial state. This is of course a kind of partial Markov property and does not correspond to general non-Markov processes; however, it allows for nonexponential waiting times and is sufficient in many problems. The same hypothesis was made by van Kampen in defining the composite stochastic process,<sup>10</sup> with more restrictive conditions. Now it is necessary to specify the description of this class of processes.

## B. Evolution during a regime

Under regime  $\epsilon$  the evolution of the system is a Markov process with an infinitesimal operator  $L_\epsilon(x)$ ,

$$\frac{\partial}{\partial t} p_\epsilon(x, t) = L_\epsilon(x) p_\epsilon(x, t), \quad (1)$$

where  $L_\epsilon(x)$  is a linear operator acting on the  $x$  variables only, satisfying the condition

$$\int dx L_\epsilon(x) p_\epsilon(x, t) = 0, \quad (2)$$

which can be written in the integral form

$$L_\epsilon(x) p_\epsilon(x, t) = \int dx' K_\epsilon(x|x'; t) p_\epsilon(x', t), \quad (3)$$

with

$$\int dx K_\epsilon(x|x'; t) = 0. \quad (4)$$

For a stationary continuous process with transition rate  $W_\epsilon(x'|x)$  from state  $x$  to state  $x'$ ,

$$L_\epsilon(x) p_\epsilon(x, t) = \int dx' W_\epsilon(x|x') p_\epsilon(x', t) - p_\epsilon(x, t) \int dx' W_\epsilon(x'|x), \quad (5)$$

which can be written in the integral form (3) with the kernel

$$K_\epsilon(x|x') = W_\epsilon(x|x') - \delta(x - x') \int dx'' W_\epsilon(x''|x'). \quad (6)$$

In the case of a deterministic process with conjugated coordinates  $x = \{q_i, p_i\}$  and Hamiltonian  $H_\epsilon(x)$ ,  $L_\epsilon(x)$  is of course the Lagrangian operator

$$L_\epsilon(x) p_\epsilon(x, t) = \sum_i \left( \frac{\partial H_\epsilon}{\partial q_i} \frac{\partial}{\partial p_i} - \frac{\partial H_\epsilon}{\partial p_i} \frac{\partial}{\partial q_i} \right) p_\epsilon(x, t) \quad (7)$$

corresponding to the integral kernel,

$$K_\epsilon(x|x') = \sum_i \left( \frac{\partial H_\epsilon}{\partial p_i} \delta(p_i - p'_i) \frac{\partial}{\partial q'_i} \delta(q_i - q'_i) - \frac{\partial H_\epsilon}{\partial q_i} \delta(q_i - q'_i) \frac{\partial}{\partial p'_i} \delta(p_i - p'_i) \right). \quad (8)$$

## C. Evolution of the regime $\epsilon$

This evolution contains all the (possibly) non-Markovian behavior of the system. It is well known<sup>12</sup> that when the successive regimes  $\epsilon_0, \epsilon_1, \dots$  constitute a Markov chain with an exponentially distributed pausing time between two changes, the process  $\epsilon(t)$  is Markovian, since for  $t > t_0$ ,  $\epsilon(t)$  is independent of the events preceding  $t$  once  $\epsilon(t_0)$  is known. Conformally with the fundamental hypothesis, we assume that once the state  $x_0$  is known at the beginning  $t_0$  of a regime  $\epsilon_0$ , the process  $\{\epsilon(t), X(t)\}$  at further times  $t > t_0$  is independent of the events preceding  $t_0$ ; this is not true if  $t_0$  is not the beginning of the regime  $\epsilon_0$  (unless the pausing time in  $\epsilon_0$  is exponentially distributed) so that the process is generally non-Markovian.

In our formalism the changes of regime are instantaneous; let  $C(\epsilon_0)$  be the first regime following  $\epsilon_0$ , and  $T_{\epsilon_0, \epsilon_1}$  the duration of regime  $\epsilon_0$  when  $C(\epsilon_0) = \epsilon_1$ ; if  $\epsilon_0$  begins at time  $t_0$  and is changed into  $\epsilon_1$  at time  $t_1$ , then  $T_{\epsilon_0, \epsilon_1} = t_1 - t_0$ ; then a stochastic variable  $T_{\epsilon_0, \epsilon_1}$  may be defined for any regime  $\epsilon$  by taking  $T_{\epsilon_0, \epsilon} = \infty$  if  $\epsilon \neq C(\epsilon_0) = \epsilon_1$ . We leave the possibility of having  $C(\epsilon_0) = \epsilon_0$ , since the transition  $\epsilon_0 \rightarrow \epsilon_0$  may be defined (for instance it may account for instantaneous elastic collisions); in other cases we take  $T_{\epsilon_0, \epsilon_0} = \infty$ . Since the evolution of  $X(t)$  under regime  $\epsilon_0$  is specified by the infinitesimal operator  $L_{\epsilon_0}$ , the whole process  $\{\epsilon(t), X(t)\}$  is completely determined after the beginning  $t_0$  of regime  $\epsilon_0$  by giving the laws of  $T_{\epsilon_0, \epsilon}$  for any  $\epsilon$ , and the law of change of the state in the transition  $\epsilon_0 \rightarrow \epsilon_1$ . Of course, these laws are conditioned by the initial values  $\epsilon(t_0) = \epsilon_0, x(t_0) = x_0$ , and can be interdependent.

We now assume the following.

(i) The laws of  $T_{\epsilon_0, \epsilon}$  and  $C(\epsilon_0)$  only depend on  $\epsilon$ , on  $\epsilon_0$ , and on the invariants of the evolution under the action of  $L_{\epsilon_0}$  [the fact that these laws depend only on these quantities will

be explicitly used in the derivation of (Eq. 29)]. These laws are specified by

$$\Lambda_{\epsilon_0}(t|x_0, t_0) = \text{Prob}(\inf_{\epsilon} T_{\epsilon_0} > t - t_0 | x_0, t_0, \epsilon_0) \quad (9)$$

(which is the conditional probability of having no transition in the time interval  $]t_0, t]$ ), and by

$$\begin{aligned} \Gamma_{\epsilon_1, \epsilon_0}(]t, t'] | x_0, t_0) \\ = \text{Prob}(T_{\epsilon_1, \epsilon_0} = t_1 - t_0 \in ]t - t_0, t' - t_0], \\ C(\epsilon_0) = \epsilon_1 | x_0, t_0, \epsilon_0) \end{aligned} \quad (10)$$

(which is the conditional probability that the first transition after  $t_0$  occurs in the time interval  $]t, t']$  ( $t_0 \leq t < t'$ ) and yields the regime  $\epsilon_1$ ). We shall also use the corresponding density  $\gamma_{\epsilon_1, \epsilon_0}(t|x_0, t_0)$ ,

$$\begin{aligned} \Gamma_{\epsilon_1, \epsilon_0}(dt|x_0, t_0) \equiv \Gamma_{\epsilon_1, \epsilon_0}(]t, t+dt] | x_0, t) \\ = \gamma_{\epsilon_1, \epsilon_0}(t|x_0, t_0) dt. \end{aligned} \quad (10')$$

(ii) If the state of the system is  $x_1^- \in E_{\epsilon_0}$  at the end  $t_1$  of regime  $\epsilon_0$  [ $X(t_1 - 0) = x_1^-$ ], the state  $X(t_1)$  at the beginning of the next regime  $\epsilon_1$  is an element of  $E_{\epsilon_1}$ , which is independent of all other events before  $t_1$ ; thus the effect of the transition is specified by the transition probability

$$\begin{aligned} Y_{\epsilon_1, \epsilon_0}(A|x_1^-, t_1) = \text{Prob}(X(t_1) = x_1 \in A | X(t_1 - 0) = x_1^-; \\ \text{a transition } \epsilon_0 \rightarrow \epsilon_1 \text{ occurs at } t_1) \end{aligned} \quad (11)$$

(where  $A$  is a subset of the phase space  $E_{\epsilon_1}$ ) or by the corresponding density  $y_{\epsilon_1, \epsilon_0}(x_1|x_1^-, t_1)$ ,

$$Y_{\epsilon_1, \epsilon_0}(dx_1|x_1^-, t_1) = y_{\epsilon_1, \epsilon_0}(x_1|x_1^-, t_1) dx_1. \quad (11')$$

In some cases  $X(t_1)$  may be determined unambiguously from  $x_1^-$  by a mapping  $f$  of  $E_{\epsilon_0}$  into  $E_{\epsilon_1}$ ,

$$X(t_1) = f(x_1^-) \text{ with probability } 1;$$

then  $y_{\epsilon_1, \epsilon_0}$  is obviously a  $\delta$  function. In particular  $E_{\epsilon_1}$  may be identical to  $E_{\epsilon_0}$ , with  $X(t_1) = x_1^-$  if the transition does not change the state of the system.

#### D. Compound transition probabilities

Our main purpose is to find the conditional probability  $P_{\epsilon_0}(A, t|x_0, t_0)$  that the system be at time  $t$  under regime  $\epsilon$  and in some state  $x$  of the subset  $A$  of the phase space  $E_{\epsilon}$ , knowing that a transition has occurred at time  $t_0$ , yielding regime  $\epsilon_0$  and state  $x_0 \in E_{\epsilon_0}$ ,

$$\begin{aligned} P_{\epsilon_0}(A, t|x_0, t_0) = \text{Prob}(X(t) = x \in A, \\ \epsilon(t) = \epsilon | X(t_0) = x_0, \\ \epsilon(t_0) = \epsilon_0; \text{ a transition occurs at } t_0). \end{aligned} \quad (12)$$

Here  $t$  is any time posterior to  $t_0$  and not necessarily a collision time.

This conditional probability may be computed from the following two basic quantities, which we call *compound transition probabilities*.

(a) The first is the conditional probability  $\Theta_{\epsilon_0}(A, t|x_0, t_0)$  that no transition occurs in the time interval  $]t_0, t]$  and that the state  $X(t) = x$  at time  $t$  belongs to the subset  $A$  of the phase space  $E_{\epsilon_0}$ ,

$$\begin{aligned} \Theta_{\epsilon_0}(A, t|x_0, t_0) \\ = \text{Prob}(x(t) = x \in A, \inf_{\epsilon} T_{\epsilon_0} > t - t_0 | x_0, t_0, \epsilon_0). \end{aligned} \quad (13)$$

(b) The other is the conditional probability  $\Phi_{\epsilon_1, \epsilon_0}(A, ]t, t'] | x_0, t_0)$  that the first transition after  $t_0$  occurs at some epoch  $t_1$  of the time interval  $]t, t']$  ( $t_0 \leq t_1 < t'$ ), leading to the regime  $\epsilon$ , and to some state  $x_1$  belonging to the subset  $A$  of the phase space  $E_{\epsilon_1}$ ,

$$\begin{aligned} \Phi_{\epsilon_1, \epsilon_0}(A, ]t, t'] | x_0, t_0) \\ = \text{Prob}(X(t_1) = x_1 \in A; \\ T_{\epsilon_1, \epsilon_0} = t_1 - t_0 \in ]t - t_0, t' - t_0]; \\ C(\epsilon_0) = \epsilon_1 | x_0, t_0, \epsilon_0). \end{aligned} \quad (14)$$

Although these two quantities  $\Theta_{\epsilon_0}$  and  $\Phi_{\epsilon_1, \epsilon_0}$  may possibly be defined in more general situations, we shall restrict ourselves to the particular case described by the hypotheses (i) and (ii) settled in 2.3: then the pausing time in regime  $\epsilon_0$  does not depend on the state evolution during  $\epsilon_0$ , and we have

$$\Theta_{\epsilon_0}(A, t|x_0, t_0) = \Lambda_{\epsilon_0}(t|x_0, t_0) P_{\epsilon_0}(A, t|x_0, t_0), \quad (15)$$

where  $\Lambda_{\epsilon_0}(t|x_0, t_0)$ , defined by (9), specifies the pausing time in  $\epsilon_0$  and  $P_{\epsilon_0}(A, t|x_0, t_0)$  is the transition probability from state  $x_0$  at time  $t_0$  under regime  $\epsilon_0$ . Furthermore,

$$\begin{aligned} \Phi_{\epsilon_1, \epsilon_0}(A, dt_1|x_0, t_0) \\ = \int Y_{\epsilon_1, \epsilon_0}(A|x_1^-, t_1) \Gamma(dt_1|x_0, t_0) P_{\epsilon_0}(dx_1^-, t_1|x_0, t_0), \end{aligned} \quad (16)$$

where  $\Phi_{\epsilon_1, \epsilon_0}$  and  $Y_{\epsilon_1, \epsilon_0}$  are defined by (10) and (11).

We shall also use the densities  $\theta_{\epsilon_0}$  and  $\varphi_{\epsilon_1, \epsilon_0}$  corresponding to  $\Theta_{\epsilon_0}$  and  $\Phi_{\epsilon_1, \epsilon_0}$  when they exist,

$$\Theta_{\epsilon_0}(x, t|x_0, t_0) = \Lambda_{\epsilon_0}(t|x_0, t_0) p_{\epsilon_0}(x, t|x_0, t_0), \quad (17)$$

$$\begin{aligned} \varphi_{\epsilon_1, \epsilon_0}(x_1, t_1|x_0, t_0) = \int dx_1^- y_{\epsilon_1, \epsilon_0}(x_1|x_1^-, t_1) \gamma_{\epsilon_1, \epsilon_0}(t_1|x_0, t_0) \\ \times p_{\epsilon_0}(x_1^-, t_1|x_0, t_0). \end{aligned} \quad (18)$$

Clearly the quantities  $\Theta$  and  $\Phi$  are not independent. As a matter of fact it results from definitions (13) and (14) that

$$\Theta_{\epsilon_0}(E_{\epsilon_0}, t|x_0, t_0) + \sum_{\epsilon_1} \Phi_{\epsilon_1, \epsilon_0}(E_{\epsilon_1}, ]t_0, t] | x_0, t_0) = 1, \quad (19)$$

since the second term on the left-hand side is the conditional probability that a first transition after  $t_0$  occurs before or at  $t$ . With the expressions (15) and (16) for  $\Theta$  and  $\Phi$ , this relation becomes

$$\Lambda_{\epsilon_0}(t|x_0, t_0) + \sum_{\epsilon_1} \Gamma_{\epsilon_1, \epsilon_0}(]t_0, t] | x_0, t_0) = 1 \quad (19')$$

with the same interpretation.

Clearly  $\Theta(E_{\epsilon_0}, t|x_0, t_0)$  or  $\Lambda_{\epsilon_0}(t|x_0, t_0)$  cannot increase with  $t$ . From now on we shall assume that these quantities tend to 0 when  $t \rightarrow \infty$ , that is to say, *for any initial regime and state, at least one transition will almost surely occur*; in this case

$$\Theta_{\epsilon_0}(E, \infty|x_0, t_0) = \Lambda_{\epsilon_0}(\infty|x_0, t_0) = 0. \quad (20)$$

Although the statement of the previous definitions and hypotheses has been somewhat lengthy, it is now easy to obtain

the evolution equation of the conditional probability  $P_{\epsilon\epsilon_0}(A, t | x_0, t_0)$ .

### III. EVOLUTION EQUATIONS OF THE SYSTEM

#### A. Fundamental evolution equation

We want to obtain an evolution equation for the probabilities  $P_{\epsilon\epsilon_0}(A, t | x_0, t_0)$ . We divide the set of trajectories starting from  $x_0$  in regime  $\epsilon_0$  at time  $t_0$  (just after a collision has occurred) and arriving in the subset  $A$  of  $E_\epsilon$  in regime  $\epsilon$  at time  $t$ , into two subclasses.

(1) The first subset is the set of trajectories that do not suffer any collision and thus do not change the regime in the time interval  $]t_0, t]$  so that in particular  $\epsilon(t) = \epsilon_0$ ; this means that  $\text{Inf}_\epsilon T_{\epsilon\epsilon_0} > t - t_0$  and by Ref. 9, this subset of trajectories has the probability

$$\Theta_{\epsilon_0}(A, t | x_0, t_0) \delta_{\epsilon\epsilon_0}. \quad (21)$$

(2) The second subset is the set of trajectories suffering at least one collision in the time interval  $]t_0, t]$ ; if  $t_1 \in ]t_0, t]$  denotes the first collision time occurring in this interval, then the probability that the trajectories change their regime from  $\epsilon_0$  to  $\epsilon_1$  at a first collision in the time interval  $]t_1, t_1 + dt_1]$ , and are found in the volume element  $dx_1$  in  $E_{\epsilon_1}$  just after the collision, is  $\Phi_{\epsilon_1\epsilon_0}(dx_1, dt_1 | x_0, t_0)$ . After that collision, the process follows independently and has probability  $P_{\epsilon\epsilon_1}(A, t | x_1, t_1)$  to go from  $(x_1, t_1)$  to  $(A, t)$  (because, precisely,  $t_1$  is a collision time). The probability of this subset is then the sum over all these possibilities,

$$\int_{t_1 \in ]t_0, t]} \sum_{\epsilon_1} \int_{x_1^- \in E_{\epsilon_1}} P_{\epsilon\epsilon_1}(A, t | x_1, t_1) \Phi_{\epsilon_1\epsilon_0}(dx_1, dt_1 | x_0, t_0). \quad (21')$$

Thus  $P$  satisfies the integral equation obtained by adding the contribution (21) and (21'),

$$P_{\epsilon\epsilon_0}(A, t | x_0, t_0) = \Theta_{\epsilon_0}(A, t | x_0, t_0) \delta_{\epsilon\epsilon_0} + \sum_{\epsilon_1} \int_{t_1 \in ]t_0, t]} \int_{x_1 \in E_{\epsilon_1}} P_{\epsilon\epsilon_1}(A, t | x_1, t_1) \times \Phi_{\epsilon_1\epsilon_0}(dx_1, dt_1 | x_0, t_0), \quad (22)$$

and the conditional density  $p$  satisfies

$$p_{\epsilon\epsilon_0}(A, t | x_0, t_0) = \theta_{\epsilon_0}(x, t | x_0, t_0) \delta_{\epsilon\epsilon_0} + \sum_{\epsilon_1} \int_{t_0}^t dt_1 \int dx_1 \times p_{\epsilon\epsilon_1}(x, t | x_1, t_1) \varphi_{\epsilon_1\epsilon_0}(x_1, t_1 | x_0, t_0). \quad (23)$$

Equations (22) and (23) are the fundamental evolution

equations under their integral form. It is seen that they are "backward" equations, since the operator of the right-hand side acts on the initial values  $\epsilon_0, x_0, t_0$  only.

*Remark:* This is a generalized version of a renewal equation [see (18)] but with memory and nonidentical laws.

#### B. Integrodifferential evolution equation

In order to compare the evolution equation of our processes with the usual master equation of Markov processes, we derive Eq. (23) with respect to the initial time  $t_0$ , and obtain

$$\begin{aligned} \frac{\partial}{\partial t_0} p_{\epsilon\epsilon_0}(x, t | x_0, t_0) &= \frac{\partial}{\partial t_0} \theta_{\epsilon_0}(x, t | x_0, t_0) \delta_{\epsilon\epsilon_0} + \sum_{\epsilon_1} \int_{t_0}^t dt_1 \int dx_1 \\ &\times p_{\epsilon\epsilon_1}(x, t | x_1, t_1) \frac{\partial}{\partial t_0} \varphi_{\epsilon_1\epsilon_0}(x_1, t_1 | x_0, t_0) \\ &- \sum_{\epsilon_1} \int dx_1 p_{\epsilon\epsilon_1}(x, t | x_1, t_0) \varphi_{\epsilon_1\epsilon_0}(x_1, t_0 | x_0, t_0). \end{aligned} \quad (24)$$

Now we use the expressions (17) and (18) of  $\theta$  and  $\varphi$ ,

$$\theta_{\epsilon_0}(x, t | x_0, t_0) = \Lambda_{\epsilon_0}(t | x_0, t_0) p_{\epsilon_0}(x, t | x_0, t_0), \quad (25)$$

$$\begin{aligned} \varphi_{\epsilon_1\epsilon_0}(x_1, t_1 | x_0, t_0) &= \int dx_1^- y_{\epsilon_1\epsilon_0}(x_1 | x_1^-, t_1) \delta_{\epsilon_1\epsilon_0} \\ &\times (t_1 | x_0, t_0) p_{\epsilon_0}(x_1^-, t_1 | x_0, t_0). \end{aligned} \quad (26)$$

We notice that under regime  $\epsilon_0$ , the evolution equation (1) implies the forward equation,

$$\frac{\partial}{\partial t} p_{\epsilon_0}(x, t | x_0, t_0) = L_{\epsilon_0}(x) p_{\epsilon_0}(x, t | x_0, t_0), \quad (27)$$

whereas the backward equation (which is more convenient mathematically, see Ref. 13) reads

$$\frac{\partial}{\partial t_0} p_{\epsilon_0}(x, t | x_0, t_0) = p_{\epsilon_0}(x, t | x_0, t_0) \bar{L}_{\epsilon_0}(x_0), \quad (27')$$

where  $\bar{L}_{\epsilon_0}(x_0)$  operates on the left, on variable  $x_0$ ;  $\bar{L}_{\epsilon_0}$  is the opposite of the adjoint  $L_{\epsilon_0}^*$  of  $L_{\epsilon_0}$ ,

$$\bar{L}_{\epsilon_0} = -L_{\epsilon_0}^*. \quad (28)$$

Since  $\Lambda_0(t | x_0, t_0)$  (the probability that no transition occurs before  $t$ ) and  $\gamma_{\epsilon_1\epsilon_0}(t_1 | x_0, t_0) dt_1$  (the probability that the first transition occurs in  $dt_1$  and leads to  $\epsilon_1$ ) only depend on  $x_0$  through the evolution constants, we have, for instance,

$$(\theta_{\epsilon_0} p_{\epsilon_0}) \bar{L}_{\epsilon_0} = \theta_{\epsilon_0} (p_{\epsilon_0} \bar{L}_{\epsilon_0}).$$

Then Eq. (24) becomes

$$\begin{aligned} \frac{\partial}{\partial t_0} p_{\epsilon\epsilon_0}(x, t | x_0, t_0) &= \Lambda_{\epsilon_0}(t | x_0, t_0) p_{\epsilon_0}(x, t | x_0, t_0) \bar{L}_{\epsilon_0}(x_0) \delta_{\epsilon\epsilon_0} + \frac{\partial}{\partial t_0} \Lambda_{\epsilon_0}(t | x_0, t_0) p_{\epsilon_0}(x, t | x_0, t_0) \delta_{\epsilon\epsilon_0} \\ &+ \left\{ \sum_{\epsilon_1} \int_{t_0}^t dt_1 \int dx_1 dx_1^- p_{\epsilon\epsilon_1}(x, t | x_1, t_1) y_{\epsilon_1\epsilon_0}(x_1 | x_1^-, t_1) \delta_{\epsilon_1\epsilon_0}(t_1 | x_0, t_0) p_{\epsilon_0}(x_1^-, t_1 | x_0, t_0) \right\} \bar{L}_{\epsilon_0}(x_0) \\ &+ \sum_{\epsilon_1} \int_{t_0}^t dt_1 \int dx_1 dx_1^- p_{\epsilon\epsilon_1}(x, t | x_1, t_1) y_{\epsilon_1\epsilon_0}(x_1 | x_1^-, t_1) \frac{\partial}{\partial t_0} \gamma_{\epsilon_1\epsilon_0}(t_1 | x_0, t_0) p_{\epsilon_0}(x_1^-, t_1 | x_0, t_0) \\ &- \sum_{\epsilon_1} \int dx_1 p_{\epsilon\epsilon_1}(x, t | x_1, t_0) \varphi_{\epsilon_1\epsilon_0}(x_1, t_0 | x_0, t_0). \end{aligned} \quad (29)$$

Using the integral equation (23), it is seen that the first and third term on the right-hand side sum up to give  $p_{\epsilon_0} \bar{L}_{\epsilon_0}$ . Furthermore, the second term may be expressed in function of  $p_{\epsilon_0}$  by using (23) again, if  $\Lambda_{\epsilon_0}^{-1}(\partial/\partial t)\Lambda_{\epsilon_0}$  is finite. Then (29) can be written

$$\begin{aligned} \frac{\partial}{\partial t_0} p_{\epsilon_0}(x, t | x_0, t_0) &= p_{\epsilon_0}(x, t | x_0, t_0) \bar{L}_{\epsilon_0}(x_0) - p_{\epsilon_0}(x, t | x_0, t_0) \Lambda_{\epsilon_0}^{-1} \frac{\partial}{\partial t_0} \Lambda_{\epsilon_0}(x, t | x_0, t_0) \\ &\quad - \sum_{\epsilon_1} \int dx_1 p_{\epsilon_1}(x, t | x_1, t_0) y_{\epsilon_1, \epsilon_0}(x_1 | x_0, t_0) \gamma_{\epsilon_1, \epsilon_0}(t_0 | x_0, t_0) + \sum_{\epsilon_1} \int_{t_0}^t dt_1 \int dx_1 dx_1^- p_{\epsilon_1}(x, t | x_1, t_1) \\ &\quad \times y_{\epsilon_1, \epsilon_0}(x_1 | x_1^-, t_1) N_{\epsilon_1, \epsilon_0}(t, t_1 | x_0, t_0) p_{\epsilon_0}(x_1^-, t_1 | x_0, t_0), \end{aligned} \quad (30)$$

with

$$\begin{aligned} N_{\epsilon_1, \epsilon_0}(t, t_1 | x_0, t_0) &= \frac{\partial}{\partial t_0} \gamma_{\epsilon_1, \epsilon_0}(t_1 | x_0, t_0) - \gamma_{\epsilon_1, \epsilon_0}(t_1 | x_0, t_0) \Lambda_{\epsilon_0}^{-1} \\ &\quad \times \frac{\partial}{\partial t_0} \Lambda_{\epsilon_0}(t | x_0, t_0). \end{aligned} \quad (31)$$

This is the generalized evolution equation of the process. It is not a master equation in the usual sense since the  $p_{\epsilon_0}(x, t | x_0, t_0)$  are not transition probabilities (remember that  $t_0$  is a time of "collision," or of a change of regime, and not any time). Only the fourth term on the right-hand side of (30) contains a memory kernel, because of the memory factor  $N_{\epsilon_1, \epsilon_0}$ . This memory vanishes if

$$\begin{aligned} -\gamma_{\epsilon_1, \epsilon_0}^{-1} \frac{\partial}{\partial t_0} \gamma_{\epsilon_1, \epsilon_0}(t_1 | x_0, t_0) \\ = -\Lambda_{\epsilon_0}^{-1} \frac{\partial}{\partial t_0} \Lambda_{\epsilon_0}(t | x_0, t_0) \equiv \lambda_{\epsilon_0}(x_0, t_0), \end{aligned} \quad (32)$$

which implies that this common positive value  $\lambda_{\epsilon_0}(x_0, t_0)$  is independent of  $\epsilon_1$ ,  $t_1$ , and  $t$ . Then we see from (32) that

$$\Lambda_{\epsilon_0}(t | x_0, t_0) = \exp - \int_{t_0}^t dt' \lambda_{\epsilon_0}(x_0, t'), \quad (33)$$

$$\gamma_{\epsilon_1, \epsilon_0}(t_1 | x_0, t_0) = \gamma_{\epsilon_1, \epsilon_0}(t_1 | x_0, t_1) \exp - \int_{t_0}^{t_1} dt' \lambda_{\epsilon_0}(x_0, t'). \quad (33')$$

It results from these relations that Eq. (30) reduces to an ordinary evolution equation without memory if

$$\Lambda_{\epsilon_0}^{-1}(t | x_0, t_0) \frac{\partial}{\partial t_0} \Lambda_{\epsilon_0}(t | x_0, t_0) = -\lambda_{\epsilon_0}(x_0, t), \quad (34)$$

$$\Lambda_{\epsilon_0}^{-1}(t | x_0, t_0) \gamma_{\epsilon_1, \epsilon_0}(t | x_0, t_0) = \gamma_{\epsilon_1, \epsilon_0}(t | x_0, t). \quad (34')$$

These quantities are, respectively, the conditional probability rate of having any transition at time  $t$ , and of having a

transition towards regime  $\epsilon_1$  at time  $t$ , knowing that there has been no transition from  $t_0$  to  $t$ . Equation (34) and (34') show that if the memory vanishes they may depend on  $t$  but not on the time  $t_0$  of the previous transition. These conditions are necessary to obtain a Markov process.

*Remark:* In the case of Markov processes, the transition probabilities specify the Markov process in an unique way (i.e., up to an isomorphism of probability spaces). In our case, we can study by the generalized evolution equation only some special transition probabilities. In general, they will not be sufficient to determine the process completely, but for physical applications they will give information on the two-time correlation functions, the equilibrium value distribution, and the rate constants (see Ref. 12 and subsequent publications).

On the other hand, we have described in Ref. 11 the trajectories of our processes; this work stresses the analytical aspect and we refer to<sup>11</sup> for a detailed description of the trajectories.

### C. Time-homogeneous processes

It is conceptually interesting to distinguish the initial and final times in the conditional probabilities, as has been done previously, in order to point out the backward nature of the equations. However, many of the physical processes are time homogeneous: the waiting time  $T_{\epsilon_0 \rightarrow \epsilon_1}$  of a transition  $\epsilon_0 \rightarrow \epsilon_1$  does not depend on the beginning  $t_0$  of regime  $\epsilon_0$ , and all two-time quantities only depend on the difference between these times; then the quantities  $\lambda_{\epsilon_0}$  and  $\gamma_{\epsilon_1, \epsilon_0}$  defined by (34) and (34') are independent of  $t$ , as well as the transition rates  $y_{\epsilon_1, \epsilon_0}(x_1 | x_1^-)$ . In this case we write

$$p_{\epsilon_0}(x, t | x_0, t_0) = p_{\epsilon_0}(x, t - t_0 | x_0),$$

and the generalized master equation becomes

$$\begin{aligned} \frac{\partial}{\partial t} p_{\epsilon_0}(x, t | x_0) \\ = -p_{\epsilon_0}(x, t | x_0) \bar{L}_{\epsilon_0}(x_0) - p_{\epsilon_0}(x, t | x_0) \Lambda_{\epsilon_0}^{-1} \frac{\partial \Lambda}{\partial t}(x, t | x_0) + \sum_{\epsilon_1} \int dx_1 p_{\epsilon_1}(x, t | x_1) y_{\epsilon_1, \epsilon_0}(x_1 | x_0) \gamma_{\epsilon_1, \epsilon_0}(x_0) \\ + \sum_{\epsilon_1} \int_{t_0}^t d\tau \int dx_1 dx_1^- p_{\epsilon_1}(x, t - \tau | x_1) y_{\epsilon_1, \epsilon_0}(x_1 | x_1^-) N_{\epsilon_1, \epsilon_0}(t, \tau | x_0) p_{\epsilon_0}(x_1^-, \tau | x_0). \end{aligned} \quad (35)$$



The memory term  $N_{\epsilon_1, \epsilon_0}$  is easily deduced from (31),

$$N_{\epsilon_1, \epsilon_0}(t, \tau | x_0) = \frac{\partial}{\partial \tau} \gamma_{\epsilon_1, \epsilon_0}(\tau | x_0) - \gamma_{\epsilon_1, \epsilon_0}(\tau | x_0) \Lambda_{\epsilon_0}^{-1}(t | x_0) \frac{\partial}{\partial t} \Lambda_{\epsilon_0}(t | x_0). \quad (36)$$

From (33) and (33'),  $N_{\epsilon_1, \epsilon_0}$  vanishes if

$$\Lambda_{\epsilon_0}(t | x_0) = \exp(-\lambda_{\epsilon_0}(x_0)t), \quad (37)$$

$$\gamma_{\epsilon_1, \epsilon_0}(t | x_0) = \gamma_{\epsilon_1, \epsilon_0}(x_0) \exp(-\lambda_{\epsilon_0}(x_0)t), \quad (37')$$

thus exponential waiting times are necessary and sufficient conditions to recover a process without memory, as was expected; in this case (35) is exactly the backward master equation.

From now on we will treat time-homogeneous processes, unless otherwise specified.

#### D. A particular case

It may occur that the conditional probability of having a transition  $\epsilon_0 \rightarrow \epsilon_1$  between  $t$  and  $t + dt$ , knowing that there is a transition from  $\epsilon_0$  in this time interval, is independent of the pausing time in  $\epsilon_0$ ; then one may write

$$\frac{\gamma_{\epsilon_1, \epsilon_0}(t | x_0)}{-(\partial/\partial t)\Lambda_{\epsilon_0}(t | x_0)} = a_{\epsilon_1, \epsilon_0}(x_0) = a_{\epsilon_1, \epsilon_0}(x_1^-) \quad (38)$$

(since the dependence on  $x_0$  is only through constant of motions, which are the same in state  $x_1^-$  at the end of regime  $\epsilon_0$ ).

Physically, this property can be true if all transitions have the same cause which determines the end of  $\epsilon_0$ , the choice of the next regime  $\epsilon_1$  being independent of the pausing time in  $\epsilon_0$ . In this case the generalized master equation may be written in a simple form by defining the matrices

$$\begin{aligned} \mathbf{p} &= (p_{\epsilon\epsilon_0}), \quad \mathbf{p}^0 = (\delta_{\epsilon\epsilon_0} p_{\epsilon_0}), \\ \Lambda &= (\delta_{\epsilon\epsilon_0} \Lambda_{\epsilon_0}), \quad \boldsymbol{\mu} = (y_{\epsilon\epsilon_0} a_{\epsilon\epsilon_0}), \\ \mathbf{N}(t, \tau | x_0) &= \frac{\partial^2}{\partial \tau^2} \Lambda(\tau | x_0) - \frac{\partial}{\partial \tau} \Lambda(\tau | x_0) \Lambda^{-1}(t | x_0) \frac{\partial}{\partial t} \Lambda(\tau | x_0), \end{aligned} \quad (39)$$

and the (diagonal) matrix of operators,

$$\bar{\mathbf{L}}(x_0) = (\delta_{\epsilon\epsilon_0} \bar{L}_{\epsilon_0}(x_0)).$$

Then, using matrices products, Eq. (35) becomes

$$\begin{aligned} \frac{\partial}{\partial t} \mathbf{p}(x, t | x_0) &= -\mathbf{p}(xt | x_0) \bar{\mathbf{L}}(x_0) + \mathbf{p}(xt | x_0) \Lambda^{-1} \frac{\partial}{\partial t} \Lambda(t | x_0) \\ &\quad - \int dx_1 \mathbf{p}(xt | x_1) \boldsymbol{\mu}(x_1 | x_0) \frac{\partial}{\partial t} \Lambda(0 | x_0) \\ &\quad - \int_0^t d\tau \int dx_1 dx_1^- \mathbf{p}(x, t - \tau | x_1) \\ &\quad \times \boldsymbol{\mu}(x_1 | x_1^-) \mathbf{p}^0(x_1^-, \tau | x_0) \mathbf{N}(t, \tau | x_0), \end{aligned} \quad (40)$$

which is the final equation given in Ref. 11.

It is clear from (37) and (37') that Markov processes must be of this kind. In this case we can write

$$\Lambda(t | x_0) = \exp(-\lambda t), \quad (41)$$

and Eq. (40) becomes a standard master equation,

$$\begin{aligned} \frac{\partial}{\partial t} \mathbf{p}(xt | x_0) &= -\mathbf{p}(xt | x_0) \bar{\mathbf{L}}(x_0) + \int dx_1 \{ \mathbf{p}(xt | x_1) \\ &\quad - \mathbf{p}(xt | x_0) \} \boldsymbol{\mu}(x_1 | x_0) \lambda. \end{aligned} \quad (42)$$

If the Markov process during regime  $\epsilon$  is a jump process with an exponential pausing time with time constant  $\eta_\epsilon$  and transition rate  $\nu_\epsilon(x' | x)$  from  $x$  to  $x'$ , it is seen that the total Markov process is a jump process with time constant  $\lambda_\epsilon + \eta_\epsilon$  and transition rates from  $\epsilon$  to  $\epsilon'$ ,

$$\frac{\mu_{\epsilon\epsilon'}(x' | x) \lambda_\epsilon + \delta_{\epsilon\epsilon'} \nu_\epsilon(x' | x) \eta_\epsilon}{\lambda_\epsilon + \eta_\epsilon}.$$

#### E. A remark on non-Markovian regimes

It can be interesting for some applications to study the more general situation where the evolution in the regimes is not Markovian. Then it can be seen that, if quantities  $\theta$  and  $\phi$  defined by (13) and (14) can be computed, the integral equation (22) applies, but the integrodifferential equations such as (29) obviously do not make sense.

### IV. GENERAL PROPERTIES

#### A. Condensed evolution equation

Much works on non-Markov processes use a generalized Langevin equation: for this reason we have carefully studied how such an equation may be deduced from the fundamental evolution equation (23). However, the general properties of the process are more easily derived from the integral equation (22) for the conditional probability  $P_{\epsilon_0}(A, t | x_0, t_0)$  to be at time  $t$  in regime  $\epsilon$  and in subset  $A$  of the state space, knowing that the system was in regime  $\epsilon_0$  and state  $x_0$  at time  $t_0$ , just after a transition.

We again use the compound transition probabilities  $\Theta_{\epsilon_0}$  and  $\Phi_{\epsilon_0}$  defined in Sec. II D with the same notations.

In order to write the evolution equations concisely we introduce the matrices

$$\begin{aligned} \mathbf{P} &= (P_{\epsilon\epsilon_0}), \\ \boldsymbol{\Theta} &= (\Theta_{\epsilon\epsilon_0}) \equiv (\delta_{\epsilon\epsilon'} \Theta_{\epsilon_0}), \\ \boldsymbol{\Phi} &= (\Phi_{\epsilon\epsilon_0}), \end{aligned} \quad (43)$$

and we define the  $*$  product as the matricial product followed by the integrations on  $x_1$  and on  $t_1$ . Then Eq. (22) may be written

$$\mathbf{P}(A, t | x_0, t_0) = (\boldsymbol{\Theta} + \mathbf{P} * \boldsymbol{\Phi})(A, t | x_0, t_0). \quad (44)$$

Here  $A$  is a subset of the Cartesian product  $E$  of all the state spaces  $E_\epsilon$ ,

$$E = \prod_{\epsilon} E_{\epsilon}.$$

In these condensed notations, relation (19) between  $\Theta$  and  $\Phi$  becomes

$$\mathbf{1} = (\mathbf{U} + \mathbf{1} * \Phi)(E, t | x_0, t_0), \quad (45)$$

where  $\mathbf{1}$  denotes the  $(m + 1)$  line vector  $(1, 1, \dots, 1)$  and  $\mathbf{U}$  the  $(m + 1)$  line vector defined by

$$\mathbf{U} = (\Theta_{\epsilon_0}), \quad \epsilon_0 = 0, 1, 2, \dots \quad (46)$$

This condition is clearly necessary in order that the solution of (22) be normalized, since it results from (22) by writing  $\mathbf{P}(E, t | x_0, t_0) = 1$ . We shall see that it is not always sufficient.

Before studying the solution of Eqs. (22) and (23), we notice that since  $\mathbf{U}(E, t | x_0, t_0) \rightarrow 0$  as  $t \rightarrow \infty$ , Eq. (45) shows that  $\Phi$  is a solution of the equation

$$\mathbf{v} = (\mathbf{v} * \Phi)(E, \infty | x_0, t_0), \quad (47)$$

which reads explicitly

$$v_{\epsilon_0}(x_0, t_0) = \sum_{\epsilon_1} \int_{t_1 > t_0} \int_{x_1 \in E_1} v_{\epsilon_1}(x_1, t_1) \times \Phi_{\epsilon_1 \epsilon_0}(dx_1, dt_1 | x_0, t_0), \quad (48)$$

where  $v_{\epsilon_0}(E, \infty | x_0, t_0) \equiv v_{\epsilon_0}(x_0, t_0)$  is a function  $\epsilon_0, x_0, t_0$ .

On the other hand, if we assume that  $P_{\epsilon \epsilon_0}(A, t | x_0, t_0)$  has a limit  $Q_{\epsilon \epsilon_0}(A | x_0, t_0) = P_{\epsilon \epsilon_0}(A, \infty | x_0, t_0)$  when  $t \rightarrow \infty$ , Eq. (22) shows that the  $(p + 1)$  vector  $\mathbf{Q}_\epsilon = (Q_{\epsilon \epsilon_0})$ ,  $\epsilon_0 = 0, 1, 2, \dots$ , is a solution of (47) for any  $\epsilon$  or  $A$ . Thus if this equation has a unique solution (up to a multiplicative constant) we conclude that

$$\mathbf{Q}_\epsilon \propto \mathbf{1},$$

which means that  $P_{\epsilon \epsilon_0}(A, \infty | x_0, t_0)$  do not depend on  $\epsilon_0, x_0, t_0$ .

The previous assumptions must be verified on particular models; nevertheless it can be said that the asymptotic probability law, if defined, is generally independent of the initial conditions, as could be expected.

## B. Formal solution of the evolution equation

Iterating Eq. (22) yields the formal solution of Eqs. (22) and (23),

$$\mathbf{P} = \Theta + \Theta * \Phi + \Theta * \Phi * \Phi + \dots + \Theta * \Phi^{*n} + \dots, \quad (49)$$

where the right-hand side displays the successive contributions of the evolutions with no change of regime between  $t_0$  and  $t$ , with one change, two changes, etc.

This is the equation written by van Kampen for the composite stochastic processes.<sup>10</sup> In Appendix A, it is shown that the formal series that appears in Eq. (49) is indeed convergent.

In most practical cases its sum is normalized (unless an infinite number of transitions can occur during a finite time interval) and then it yields the unique physically significative solution of the evolution equation. However, this solution, under the form (49), is of little practical use; for this reason we now turn to the solution of the evolution equation by Laplace transforms.

## V. STUDY BY LAPLACE TRANSFORM

### A. General solution

The Laplace transform of a matrix function  $\mathbf{f}(x, t | x_0)$  will be denoted  $\hat{\mathbf{f}}(x, s | x_0)$ ,

$$\hat{\mathbf{f}}(x, s | x_0) = \int_0^\infty dt e^{-st} \mathbf{f}(x, t | x_0). \quad (50)$$

We now explicitly restrict to time homogeneous processes, represented by their probability density matrix  $\mathbf{p}(x, t | x_0, 0) \equiv \mathbf{p}(x, t | x_0)$ . Then the fundamental equation (23) for  $\mathbf{p}$  yields

$$\hat{\mathbf{p}}(x, s | x_0) = \hat{\boldsymbol{\theta}}(x, s | x_0) + \int dx_1 \hat{\mathbf{p}}(x, s | x_1) \hat{\boldsymbol{\varphi}}(x_1, s | x_0), \quad (51)$$

or in condensed notations

$$\hat{\mathbf{p}} = \hat{\boldsymbol{\theta}} + \hat{\mathbf{p}} \hat{\boldsymbol{\varphi}}, \quad (52)$$

where the product  $\hat{\mathbf{p}} \hat{\boldsymbol{\varphi}}$  implies the matricial product and the summation on the intermediate coordinate  $x_1$ .

By iteration it follows from (52) that

$$\hat{\mathbf{p}} = \hat{\boldsymbol{\theta}} (\hat{\mathbf{I}} + \hat{\boldsymbol{\varphi}} + \hat{\boldsymbol{\varphi}} \hat{\boldsymbol{\varphi}} + \dots + \hat{\boldsymbol{\varphi}}^n + \dots), \quad (53)$$

where  $\hat{\mathbf{I}} = (\hat{I}_{\epsilon \epsilon_0})$  and

$$\hat{I}_{\epsilon \epsilon_0}(x | x_0) = \delta_{\epsilon \epsilon_0} \delta(x - x_0).$$

The results of Sec. IV B show that the series (52) is convergent [or at least, the analogous series for  $\hat{P}(A, s | x_0)$ , which is obtained by integrating (52) on the subset  $A$  of the  $x$  space].

If the matrix  $\hat{\mathbf{I}} - \hat{\boldsymbol{\varphi}}$  is invertible in the sense of the product used in (52), the solution of (52) may also be written

$$\hat{\mathbf{p}} = \hat{\boldsymbol{\theta}} (\hat{\mathbf{I}} - \hat{\boldsymbol{\varphi}})^{-1}. \quad (54)$$

In general no condensed expression of  $(\hat{\mathbf{I}} - \hat{\boldsymbol{\varphi}})^{-1}$  is known, and the expansion (53) must be used. However, in the next subsection we consider a particular case where this expansion may be avoided.

### B. System "without phase space"

Under this name we consider systems the description of which includes no state  $x$  (or a single, definite  $x$  for each regime  $\epsilon$ ). Then the product used in (52) is only a matricial product and the solution (54) may be easily studied.

#### 1. Solution

We notice that the conservation relation (45),

$$1 = \theta_{\epsilon_0}(t) + \sum_{\epsilon'} \int_0^t dt' \varphi_{\epsilon' \epsilon_0}(t')$$

implies

$$\hat{\boldsymbol{\theta}} = s^{-1} (\hat{\mathbf{I}} - \hat{\mathbf{F}}(s)) \quad (55)$$

with

$$\hat{F}_{\epsilon \epsilon_0}(s) = \delta_{\epsilon \epsilon_0} \sum_{\epsilon'} \hat{\varphi}_{\epsilon' \epsilon_0}(s). \quad (56)$$

Except for special values of  $s$ , the matrices  $\mathbf{A}(s) = \hat{\mathbf{I}} - \hat{\mathbf{F}}(s)$  and  $\hat{\mathbf{I}} - \hat{\mathbf{F}}(s)$  can be inverted, so that the solution (54) is

$$\hat{\mathbf{P}}(s) = \hat{\boldsymbol{\theta}}(s)\hat{\mathbf{A}}(s) = s^{-1}(\hat{\mathbf{I}} - \hat{\mathbf{F}}(s))(\hat{\mathbf{I}} - \hat{\boldsymbol{\phi}}(s)) = s^{-1}\mathbf{B}^{-1}. \quad (57)$$

It is clear from (55) and (56) that the vector  $\mathbf{1} = (1, 1, \dots, 1)$  is a left eigenvector of  $\mathbf{B} = (\hat{\mathbf{I}} - \hat{\boldsymbol{\phi}})(\hat{\mathbf{I}} - \hat{\mathbf{F}})^{-1}$ , with eigenvalue 1; thus the same property is true for  $\mathbf{B}^{-1}$ , and it results from (57) that

$$\sum_{\epsilon} \hat{p}_{\epsilon\epsilon_0}(s) = s^{-1},$$

which implies that the conditional probabilities  $p_{\epsilon\epsilon_0}(t)$  are indeed normalized.

## 2. Stationary probability

The asymptotic value  $\mathbf{p}(\infty)$  of  $\mathbf{p}(t)$  when  $t \rightarrow \infty$  is obtained by taking the limit of  $s\hat{\mathbf{p}}(s)$  when  $s \rightarrow 0$ :  $s\hat{\mathbf{p}}(s) \rightarrow \mathbf{p}(\infty)$ , and it is shown in Appendix B that

$$p_{\epsilon}(\infty) = \lim_{s \rightarrow 0} s\hat{p}_{\epsilon}(s) = \frac{\bar{t}_{\epsilon}q_{\epsilon}}{\sum_{\epsilon'} \bar{t}_{\epsilon'}q_{\epsilon'}}, \quad (58)$$

where

$$\bar{t}_{\epsilon} = \int_0^{\infty} dt t \sum_{\epsilon'} \varphi_{\epsilon'\epsilon}(t) \quad (59)$$

is the mean waiting time in regime  $\epsilon$ , and  $q = (q_{\epsilon})$  is the only right eigenvector of  $\hat{\mathbf{A}}(0)$  corresponding to the eigenvalue 0.

Thus the stationary probability  $p_{\epsilon}(\infty)$  is independent of the initial regime, as it should be.

## 3. Relaxation towards the stationary state

If the required conditions are fulfilled<sup>13</sup> the conditional probability density  $p(t)$  can be obtained from its Laplace transform  $\tilde{\mathbf{p}}(s)$  by using the Laplace inversion theorem,<sup>14</sup>

$$\mathbf{p}(t) = \lim_{s \rightarrow 0} s\hat{\mathbf{p}}(s) + \sum_i \text{res}_{s_i}(e^{st}\hat{\mathbf{p}}(s)), \quad (60)$$

where  $\{s_i\}$  denotes the poles of  $e^{st}\mathbf{p}(s)$ , other than  $s = 0$ . These poles are the zeros of  $\det \mathbf{A}(s) = \det(\mathbf{I} - \boldsymbol{\varphi}(s))$ , except in exceptional cases, which will be discarded in this general discussion.

It is shown in Appendix C that  $\det \mathbf{A}(s)$  has no zero for  $\text{Re } s > 0$ . Thus the poles of  $\tilde{\mathbf{p}}(s)$  have negative real parts, and the corresponding time exponentials of  $p(t) - p(\infty)$  are time decreasing, which shows that  $p(t)$  indeed tends to  $p(\infty)$ .

## VI. CONCLUSION

It has been seen that the present model of non-Markovian processes leads to a rather simple formalism, which may be treated by generalized master equations, or by integral equations; in some cases it is possible to obtain the most important properties of the process analytically.

Many points remain to be studied, such as the possible approximation schemes [the simplest one being implied by the iterative solution (49)] or, more fundamentally, the connection of the conditional probabilities used here with the theory of fluctuations in thermodynamic equilibrium.

We postpone such discussions to further papers. Applications of this formalism are given in other articles.<sup>12</sup>

## APPENDIX A: CONVERGENCE OF THE FORMAL SOLUTION OF THE EVOLUTION EQUATION

Using the notation of Sec. IV let us consider the formal solution (49) of the evolution equation (22),

$$\mathbf{P} = \Theta + \Theta * \Phi + \Theta * \Phi * \Phi + \dots + \Theta * \Phi^{*n} + \dots \quad (A1)$$

We shall see that this formal series is indeed convergent.

As a matter of fact let us define the conditional probability  $\mathbf{P}^{\bar{n}}$ ,

$$\mathbf{P}^{\bar{0}} = \Theta, \quad (A2)$$

$$\mathbf{P}^{\bar{n}} = \Theta * \sum_{k=0}^n \Phi^{*k}. \quad (A3)$$

Here  $P_{\epsilon\epsilon_0}^{\bar{n}}(A, t | x_0, t_0)$  is the conditional probability of being in regime  $\epsilon$  and subset  $A$  of  $E_{\epsilon}$  at time  $t$  after at most  $n$  transitions since time  $t_0$ . Clearly

$$\mathbf{P}^{\bar{n+1}} = \Theta + \mathbf{P}^{\bar{n}} * \Phi = \mathbf{P}^{\bar{n}} + \Theta * \Phi^{*(n+1)} \quad (A4)$$

and  $P_{\epsilon\epsilon_0}^{\bar{n}}(A, t | x_0, t_0)$  increases with  $n$ .

On the other hand,

$$P_{\epsilon\epsilon_0}^{\bar{n}}(A, t | x_0, t_0) \leq P_{\epsilon\epsilon_0}^{\bar{n}}(E_{\epsilon}, t | x_0, t_0).$$

But if

$$\sum_{\epsilon} P_{\epsilon\epsilon_0}^{\bar{n}}(E_{\epsilon}, t | x_0, t_0) \leq 1, \quad (A5)$$

then by (A4) and (A5),

$$\sum_{\epsilon} P_{\epsilon\epsilon_0}^{\bar{n+1}}(E_{\epsilon}, t | x_0, t_0) \leq (\mathbf{U} + \mathbf{1} * \Phi)(E_{\epsilon}, t | x_0, t_0) = 1,$$

so that (A5) is satisfied for all  $n$ , since it is true for  $n = 0$ . Thus  $P_{\epsilon\epsilon_0}^{\bar{n}}(A, t | x_0, t_0)$  converges to a limit as  $n \rightarrow \infty$ ,

$$P_{\epsilon\epsilon_0}^{\bar{n}}(A, t | x_0, t_0) \rightarrow P_{\epsilon\epsilon_0}(A, t | x_0, t_0) \leq P_{\epsilon\epsilon_0}(E_{\epsilon}, t | x_0, t_0) \leq 1 \quad (A6)$$

and by (A4) the limit matrix  $\mathbf{P}(A, t | x_0, t_0)$  satisfies the fundamental equation (44).

However, it may happen that the limit  $\sum_{\epsilon} P_{\epsilon\epsilon_0}(E_{\epsilon}, t | x_0, t_0)$  is inferior to 1. But from definition (A6) it is the probability that a finite number of transitions occurs between  $t$  and  $t'$ , including 0 transition; thus  $1 - \sum_{\epsilon} P_{\epsilon\epsilon_0}(E_{\epsilon}, t | x_0, t_0)$  is the probability of having an infinite number of transitions in this time interval: it may differ from 0 in some special cases which will be excluded from this study.

In particular, if  $\Theta_{\epsilon_0}(E_{\epsilon_0}, t | x_0, t_0)$  is superior to some decreasing positive function of  $t$ ,  $1 - a(t) > 0$ , independent of the initial conditions,

$$\Theta_{\epsilon_0}(E_{\epsilon_0}, t | x_0, t_0) \geq 1 - a(t) > 0, \quad (A7)$$

then it is easily shown by recurrence that

$$P_{\epsilon_0}^{\bar{n}}(t | x_0, t_0) \equiv \sum_{\epsilon} P_{\epsilon\epsilon_0}^{\bar{n}}(E_{\epsilon}, t | x_0, t_0) \geq 1 - (a(t))^{n+1}, \quad (A8)$$

which implies

$$P_{\epsilon_0}^{\bar{n}}(t | x_0, t_0) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

Condition (A7) is satisfied if, for instance, the pausing times are exponential,

$$\Theta_{\epsilon_0}(E_{\epsilon_0}, t | x_0, t_0) \propto \exp(-\lambda_{\epsilon_0}(x_0)(t - t_0))$$

with  $\lambda_{\epsilon_0}(x_0)$  inferior to some constant  $\lambda$  independent of  $\epsilon_0$  and  $x_0$ .

On the contrary (A7) is not true in the case of periodic changes of regime with a definite period  $\tau$  (which is considered lower) since  $\Theta_{\epsilon_0} = 0$  for  $t > \tau$ . But in this case  $P_{\epsilon_0}^n(t | x_0, t_0)$ , which is the probability of having at most  $n$  transitions during  $]t_0, t]$ , is obviously 1 for  $n > t/\tau$ , so that  $P$ , given by formula (A6), is again normalized.

Finally, it should be observed that if (A6) is normalized it is the only acceptable solution of Eq. (22), since another acceptable solution should satisfy, for any subset  $A$  of  $E = \Pi_{\epsilon} E_{\epsilon}$ ,

$$Q_{\epsilon \in \epsilon_0}(A, t | x_0, t_0) \geq \delta_{\epsilon \in \epsilon_0} \Theta_{\epsilon_0}(A, t | x_0, t_0).$$

In the same way

$$Q = \Theta + Q * \Phi \geq \Theta + \Theta * \Phi = P^2,$$

and by recurrence

$$Q \geq P^n,$$

which implies

$$Q_{\epsilon \in \epsilon_0}(A, t | x_0, t_0) \geq P_{\epsilon \in \epsilon_0}(A, t | x_0, t_0), \quad (\text{A9})$$

$$\sum_{\epsilon} Q_{\epsilon \in \epsilon_0}(E_{\epsilon}, t | x_0, t_0) \geq \sum_{\epsilon} P_{\epsilon \in \epsilon_0}(E_{\epsilon}, t | x_0, t_0) = 1. \quad (\text{A10})$$

Thus  $Q$  cannot be normalized to 1, unless the equality holds in (A10), so that it also holds in (A9) for any and any  $A$ .

## APPENDIX B: STATIONARY SOLUTION OF THE EVOLUTION EQUATION

According to Sec. V B, the conditional probabilities  $P_{\epsilon \in \epsilon_0}(t)$  in systems without phase space are determined by their Laplace transforms under matricial form

$$\hat{\mathbf{P}}(s) = \hat{\boldsymbol{\theta}}(s)(\hat{\mathbf{I}} - \hat{\boldsymbol{\phi}}(s))^{-1} = s^{-1} \mathbf{B}^{-1}(s), \quad (\text{B1})$$

and the asymptotic value  $\mathbf{P}(\infty)$  of  $\mathbf{P}(t)$  is the limit of  $s\mathbf{P}(s)$  when  $s \rightarrow 0$ .

Although the matrix  $\mathbf{A}(s) = \hat{\mathbf{I}} - \hat{\boldsymbol{\phi}}(s)$  is singular for  $s = 0$  since

$$\sum_{\epsilon} \varphi_{\epsilon \in \epsilon_0}(0) = \sum_{\epsilon} \int_0^{\infty} dt \varphi_{\epsilon \in \epsilon_0}(t) = 1, \quad (\text{B2})$$

$\mathbf{B}^{-1}(s)$  generally has a finite limit as  $s \rightarrow 0$ . As a matter of fact, if  $\bar{A}_{\epsilon' \epsilon}(s)$  is the minor of the element  $A_{\epsilon' \epsilon}(s) = \delta_{\epsilon' \epsilon} - \hat{\varphi}_{\epsilon' \epsilon}(s)$  of  $\mathbf{A}(s)$ , the determinant of  $\mathbf{A}(s)$  is, for any  $\epsilon$ ,

$$\det \mathbf{A}(s) = \sum_{\epsilon'} \left( \sum_{\epsilon''} A_{\epsilon' \epsilon''}(s) \right) \bar{A}_{\epsilon \epsilon'}(s), \quad (\text{B3})$$

whereas

$$A_{\epsilon \epsilon_0}^{-1}(s) = \bar{A}_{\epsilon_0 \epsilon}(s) (\det \mathbf{A}(s))^{-1}. \quad (\text{B4})$$

Generally  $\bar{A}_{\epsilon \epsilon_0}(s)$  tends to a finite limit  $A_{\epsilon \epsilon_0}(0)$  when  $s \rightarrow 0$ . Furthermore,

$$s\hat{\boldsymbol{\theta}}_{\epsilon_0}(s) = \sum_{\epsilon} A_{\epsilon \epsilon_0}(s) = s\bar{t}_{\epsilon_0} + O(s), \quad (\text{B5})$$

where

$$\bar{t}_{\epsilon_0} = \int_0^{\infty} dt t \sum_{\epsilon} \varphi_{\epsilon \in \epsilon_0}(t) \quad (\text{B6})$$

is the mean waiting time in regime  $\epsilon_0$ . Then it results from (B5) that if  $s \rightarrow 0$ ,

$$\begin{aligned} s\hat{\mathbf{P}}_{\epsilon \in \epsilon_0}(s) &= s\hat{\boldsymbol{\theta}}_{\epsilon_0}(s) A_{\epsilon \in \epsilon_0}^{-1}(s) \\ &\rightarrow \frac{\bar{t}_{\epsilon_0} \bar{A}_{\epsilon_0 \epsilon}(0)}{\sum_{\epsilon'} \bar{t}_{\epsilon'} \bar{A}_{\epsilon_0 \epsilon'}(0)} = P_{\epsilon \in \epsilon_0}(\infty). \end{aligned} \quad (\text{B7})$$

Now for any  $\epsilon_0$ ,

$$\sum_{\epsilon_1} A_{\epsilon_0 \epsilon_1}(0) \bar{A}_{\epsilon_0 \epsilon_1}(0) = \det \mathbf{A}(0) = 0. \quad (\text{B8})$$

But the matrix  $\hat{\boldsymbol{\phi}}(0)$  is a stochastic matrix: its elements are non-negative and by (B2) they add up to one in each column; thus by the theorem of Frobenius one is a simple eigenvalue of  $\hat{\mathbf{A}}(0)$ , which corresponds to a unique right eigenvector  $\mathbf{q} = (q_{\epsilon_0})$  with non-negative components which add up to 1: it is the stationary probability of a Markov chain on the regimes  $\epsilon$ , with transition probability from  $\epsilon_0$  to  $\epsilon$  given by  $\hat{\varphi}_{\epsilon \in \epsilon_0}(0) = \int_0^t dt \varphi_{\epsilon \in \epsilon_0}(t)$ . This Markov chain is constructed from the primitive process by deleting the influence of time.

Thus (B8) shows that there exists constant  $\lambda_{\epsilon_0}$  such that

$$\bar{A}_{\epsilon_0 \epsilon_1}(0) = \lambda_{\epsilon_0} q_{\epsilon_1}$$

and (B7) yields the stationary probability,

$$P_{\epsilon \in \epsilon_0}(\infty) = \frac{\bar{t}_{\epsilon} q_{\epsilon}}{\sum_{\epsilon'} \bar{t}_{\epsilon'} q_{\epsilon'}}, \quad (\text{B9})$$

which is independent of the initial regime.

## APPENDIX C: ASYMPTOTIC BEHAVIOR OF SYSTEMS WITHOUT PHASE SPACE

It is known<sup>15</sup> that if a function  $f(t)$  satisfies

$$f(t) \propto e^{-\lambda t} (1 + \alpha(t)),$$

where  $\alpha(t)$  is bounded and tends to 0 as  $t \rightarrow \infty$ , then its Laplace transform  $\hat{f}(s)$  admits a pole of first order for  $s = -\lambda$ , and no other pole in the region  $R_e s \geq -\lambda$ . Thus the asymptotical behavior of the conditional probability  $\mathbf{P}(t)$  is found by calculating the poles of  $\hat{\mathbf{P}}(s)$ , which are the zeros of  $\det \mathbf{A}(s)$  (excepted in particular cases); but  $\det \mathbf{A}(s)$  has no zero for  $R_e s \geq 0$ . As a matter of fact,

$$\det \mathbf{A}(s) = \det(\hat{\mathbf{I}} - \hat{\boldsymbol{\phi}}(s)) = \prod_i (1 - \lambda_i(s)), \quad (\text{C1})$$

if  $\{\lambda_i(s)\}$  is the spectrum of matrix  $\hat{\boldsymbol{\phi}}(s)$ .

We may write

$$\begin{aligned} \max_i |\lambda_i(s)| &\leq \max_{\epsilon} \sum_{\epsilon'} |\hat{\varphi}_{\epsilon' \epsilon}(s)| \\ &\leq \max_{\epsilon} \sum_{\epsilon'} \int_0^{\infty} dt e^{-(R_e t)} \varphi_{\epsilon' \epsilon}(t), \end{aligned} \quad (\text{C2})$$

and thus, if  $\hat{\boldsymbol{\phi}}_{\epsilon}(s) = \sum_{\epsilon'} \hat{\varphi}_{\epsilon' \epsilon}(s)$ ,

$$\max_i |\lambda_i(s)| \leq \max_{\epsilon} \hat{\varphi}_{\epsilon}(R_e s). \quad (\text{C3})$$

(i) When  $R_{\epsilon}s > 0$ ,  $\hat{\varphi}_{\epsilon}(R_{\epsilon}s) < \hat{\varphi}_{\epsilon}(0) = 1$ : by (C1) it is seen that  $\det \mathbf{A}(s)$  cannot vanish.

(ii) When  $R_{\epsilon}s = 0$  and  $s = iy \neq 0$ , in (C2) the equality holds if and only if the waiting time  $T_{\epsilon'\epsilon}$  of a transition  $\epsilon \rightarrow \epsilon'$  takes discrete values  $\tau_{\epsilon'\epsilon} + 2k\pi/y$ ,  $k = 0, 1, 2, \dots$ . Then it may be shown that  $\det \mathbf{A}(iy)$  cannot vanish, unless the different regimes are initiated at discrete, periodic times, with period  $2\pi/y$ . Such a special case, which should be treated as a discrete time process, is out of the scope of the present work.

*Proof of (ii):* Let  $p_{\epsilon'\epsilon}$  be the probability that the first transition from  $\epsilon$  leads to  $\epsilon'$ ,

$$p_{\epsilon\epsilon'} = \int_0^{\infty} dt \varphi_{\epsilon'\epsilon}(t) = \hat{\varphi}_{\epsilon'\epsilon}(0). \quad (\text{C4})$$

Thus  $\varphi_{\epsilon'\epsilon}(t)/p_{\epsilon'\epsilon}$  is a probability density on positive times, and the complex number

$$\int_0^{\infty} dt e^{-iyt} \varphi_{\epsilon'\epsilon}(t)/p_{\epsilon'\epsilon} = z_{\epsilon'\epsilon} \quad (\text{C5})$$

necessarily lies *inside* the complex circle  $|z| < 1$ , unless the probability density  $\varphi_{\epsilon'\epsilon}(t)/p_{\epsilon'\epsilon}$  is concentrated on times for which  $e^{-iyt}$  has the same value; this means that there exists a deterministic time  $\tau_{\epsilon'\epsilon} > 0$  such that

$$\text{Prob}(e^{-iyT_{\epsilon'\epsilon}} = e^{-iy\tau_{\epsilon'\epsilon}}) = 1;$$

in this case the only possible values of  $T_{\epsilon'\epsilon}$  are

$$t_{\epsilon'\epsilon}^k = \tau_{\epsilon'\epsilon} + 2k\pi/y, \quad k = 0, 1, 2, \dots \quad (\text{C6})$$

As a result we may write

$$|\hat{\varphi}_{\epsilon'\epsilon}(iy)| = p_{\epsilon'\epsilon} |z_{\epsilon'\epsilon}| \leq p_{\epsilon'\epsilon},$$

$$\max_{\epsilon} \sum_{\epsilon'} |\hat{\varphi}_{\epsilon'\epsilon}(iy)| \leq 1,$$

the equality being possible only if  $|z_{\epsilon'\epsilon}| = 1$ , which implies (C6) for at least one  $\epsilon$ , as shown previously.

More precisely it is seen from (C1) that  $\det \mathbf{A}(iy)$  vanishes only if 1 is an eigenvalue of  $\hat{\varphi}(iy)$ , which implies the existence of a left eigenvector  $\mathbf{u} = \{\hat{u}_{\epsilon}\}$  such that

$$u_{\epsilon} = \sum_{\epsilon'} u_{\epsilon'} \varphi_{\epsilon'\epsilon}(iy). \quad (\text{C7})$$

Let the maximum value of  $|u_{\epsilon}|$  be realized for some regime  $\epsilon_0$ ; then using definitions (C5) and (C7) for  $\epsilon = \epsilon_0$ ,

$$1 = \sum_{\epsilon'} \frac{u_{\epsilon'}}{u_{\epsilon_0}} z_{\epsilon'\epsilon_0} p_{\epsilon'\epsilon_0}. \quad (\text{C8})$$

Now it should be noticed that, since  $|z_{\epsilon'\epsilon_0}| < 1$ ,

$$|(u_{\epsilon'}/u_{\epsilon_0})z_{\epsilon'\epsilon_0}| \leq 1.$$

Thus the equality (C8) can hold only if

$$(u_{\epsilon'}/u_{\epsilon_0})z_{\epsilon'\epsilon_0} = 1 \quad (\text{C9})$$

for every regime  $\epsilon'$  such that  $p_{\epsilon'\epsilon_0} \neq 0$ . Equation (C9) implies  $|u_{\epsilon'}/u_{\epsilon_0}| = 1$  and  $|z_{\epsilon'\epsilon_0}| = 1$ , which in turn implies condition (C6), and by (C5)

$$z_{\epsilon'\epsilon_0} = e^{-iy\tau_{\epsilon'\epsilon_0}}. \quad (\text{C10})$$

Defining  $\tau_{\epsilon}$  up to an additive constant by

$$u_{\epsilon} = |u_{\epsilon}| e^{-iy\tau_{\epsilon}}, \quad (\text{C11})$$

it results from (C9)–(C11) that

$$\tau_{\epsilon'\epsilon_0} = \tau_{\epsilon'} - \tau_{\epsilon_0} + 2k\pi/y. \quad (\text{C12})$$

But (C8) may be used in the same way when  $\epsilon_0$  is replaced by any regime  $\epsilon$  accessible from  $\epsilon_0$  ( $p_{\epsilon\epsilon_0} \neq 0$ ) since then  $|u_{\epsilon}| = |u_{\epsilon_0}|$ . In the simplest case  $p_{\epsilon'\epsilon} \neq 0$  for any two regimes  $\epsilon$  and  $\epsilon'$ ; as a consequence the waiting time  $T_{\epsilon'\epsilon}$  for a first transition  $\epsilon \rightarrow \epsilon'$  can only take the values

$$\tau_{\epsilon'} - \tau_{\epsilon} + 2k\pi/y, \quad k = 0, 1, 2, \dots,$$

which implies the assertions (ii).

<sup>1</sup>R. Kapral, *Adv. Chem. Phys.* **48**, 71 (1981).

<sup>2</sup>D. Chandler and L. R. Pratt, *J. Chem. Phys.* **65**, 2925 (1976).

<sup>3</sup>D. Chandler, *J. Chem. Phys.* **68**, 2959 (1978).

<sup>4</sup>J. A. Montgomery, Jr., D. Chandler, and B. J. Berne, *J. Chem. Phys.* **70**, 4056 (1979).

<sup>5</sup>S. Adelman, *J. Chem. Phys.* **73**, 3145 (1980).

<sup>6</sup>B. Carmeli and A. Nitzan, *J. Chem. Phys.* **80**, 3596 (1984).

<sup>7</sup>S. H. Northrup and J. T. Hynes, *J. Chem. Phys.* **71**, 871, 884 (1979).

<sup>8</sup>W. Feller, *Proc. Natl. Acad. Sci. USA* **51**, 653 (1964).

<sup>9</sup>W. M. Kenkre, E. W. Montroll, and M. F. Shlesinger, *J. Stat. Phys.* **9**, 45 (1979).

<sup>10</sup>N. G. Van Kampen, *Physica A* **96**, 435 (1979).

<sup>11</sup>B. Gaveau and M. Moreau, *Lett. Math. Phys.* **9**, 213 (1985).

<sup>12</sup>D. Borgis, B. Gaveau, and M. Moreau, "A class of collision processes with memory and application to simple chemical reactions in a solvent," to be published in *J. Stat. Phys.*

<sup>13</sup>W. Feller, *An Introduction to Probability Theory and its Applications* (Wiley, New York, 1971), Vol. II.

<sup>14</sup>M. R. Spiegel, *Laplace Transforms*, Schaum series (McGraw-Hill, New York, 1973).

<sup>15</sup>J. L. Doob, *Stochastic Processes* (Wiley, New York, 1967).