# A new strategy for primary structure determination of proteins: application to bovine $\beta$-casein

Christophe Carles, Jean-Claude Huet° and Bruno Ribadeau-Dumas

*Institut National de la Recherche Agronomique, CRJ, 78350 Jouy-en-Josas and °Institut National de la Recherche Agronomique, CNRA, 78000 Versailles, France*

A new approach has been developed for sequencing proteins. A radioactive label is attached specifically to the C-terminus of the protein. The labelled molecule is subjected to varying proteolysis conditions. From the electrophoretic patterns (SDS-PAGE) of the hydrolysates, appropriate cleavage conditions are selected, giving labelled peptides of different lengths which are purified. The labelled peptides are sequenced in order of increasing size (from 1 to $n$), peptide ($i$) being sequenced until the N-terminal sequence of peptide ($i$–1) is encountered. This approach allows the determination of a complete protein sequence with a minimal number of Edman cycles. The method was successfully applied to bovine $\beta$-casein (209 residues) which was completely resequenced with only 239 Edman cycles.

Protein; Sequence determination; $\beta$-Casein; (Bovine)

## 1. INTRODUCTION

DNA sequencing by the procedures of Sanger [1] or Maxam-Gilbert [2] has led to significant advances in molecular biology. Although direct application of the latter to protein sequence determination appears to be unlikely, a similar approach can provide information concerning the position of particular residues [3,4].

We have developed a new protein sequencing procedure (scheme 1) having some relationships with the Maxam-Gilbert method but more analogous to the ordered deletion methodology used to produce a set of DNA fragments for subsequent sequencing [5].

*Correspondence address:* C. Carles, Institut National de la Recherche Agronomique, CRJ, 78350 Jouy-en-Josas, France

*Abbreviations:* SDS-PAGE, SDS-polyacrylamide gel electrophoresis; TFA, trifluoroacetic acid; IBA, iodosobenzoic acid; TEA, triethylamine; PITC, phenylisothiocyanate; FPLC, fast protein liquid chromatography; RP-HPLC, reverse-phase high-performance liquid chromatography; PMSF, phenylmethanylsulfonyl fluoride; CPD-A, carboxypeptidase A; CPD-Y, carboxypeptidase Y; ATZ, anilinothiazolinone

Following reduction and alkylation, the C-terminus is specifically labelled by a carboxy-peptidase-catalysed transpeptidation between the protein and a tritiated amino acid, as described [6]. Aliquots of the labelled protein are exposed to different cleavage agents. After electrophoresis (SDS-PAGE) of the hydrolysates on the same slab, fluorographic examination indicates the number of labelled peptides, their approximate $M_r$ and gives an idea of their concentration. A large number of digestions can be tested in order to identify the most appropriate ones, which should give a set of peptides with sizes ranging from ~30 residues to the length of the entire molecule and differing by ~45 residues from one another.

The selected digestion conditions are used on larger quantities of the labelled substrate and the selected labelled peptides are purified. Their size is subsequently checked by electrophoresis (SDS-PAGE). Finally these molecules are sequenced automatically in order of increasing length (from 1 to $n$), peptide ($i$) being sequenced until the N-terminal sequence of peptide ($i-1$) is encountered.

Compared to the classical approach, our pro-

cedure presents the advantage of choosing in a reasoned way a limited number of peptides requiring partial sequencing. This allows the determination of a complete protein sequence with a minimal number of Edman cycles.

As each step had to be carefully studied, the proposed methodology has been applied to a protein of known sequence, the bovine $\beta$-casein. Our intention was also to check its primary structure. Indeed there was no complete agreement between the data obtained formerly in our laboratory from direct amino acid sequence determination [7] and those deduced from cDNA analyses carried out by two different groups [8,9]. The 3 sequences bore differences which could not be explained satisfactorily by allelic modifications.

## 2. MATERIALS AND METHODS

The following products were used: acetonitrile HPLC grade (Baker); Amplify (Amersham); Lichrosolv isopropanol, CPD-Y (Merck); IBA (Pierce); p-cresol, S. aureus V8 protease, $\alpha$-chymotrypsin, chymosin (Sigma); endoprotease Lys-C, plasmin, thermolysin (Boehringer); X-Omat films (Kodak).

Crude $\beta$-casein A[2] was prepared as described in [10] from the milk of a cow homozygous at the four casein loci. This fraction was then chromatographed on DEAE-cellulose [11]. The homogeneity of the final product was checked by FPLC [12].

The labelling reaction was performed on 3.75 mg (150 nmol) bovine $\beta$-casein with tritiated Phe amide as described [6]. The homogeneity of the transpeptidated molecule was checked by RP-HPLC and electrophoresis (SDS-PAGE). The radioactivity of the labelled $\beta$-casein was approx. 60000 dpm for 1 $\mu$g (40 pmol).

Cleavage reactions were performed on 50 $\mu$g of labelled $\beta$-casein with IBA [13], endoprotease Lys-C (50 $\mu$l of 50 mM NH$_4$HCO$_3$, pH 8.0, 5 $\mu$l of 0.1 mg/ml enzyme solution, 37°C, 2.5 h), S. aureus protease (25 $\mu$l of 50 mM NH$_4$HCO$_3$, 2 mM EDTA, pH 7.8, 1.7 $\mu$l of 1 mg/ml protease solution, 37°C, reaction stopped after 15 min by PMSF, 3 mM final), plasmin (50 $\mu$l of NH$_4$HCO$_3$, pH 8.2, 10$^{-4}$ U of plasmin, 35°C, 1 h, reaction stopped by heating at 100°C for 15 min).

After electrophoresis of these hydrolysates and fluorographic examination, the same procedures were subsequently applied to 500 $\mu$g samples of the labelled substrate.

To study the influence of ionic strength on the extent of proteolysis by chymotrypsin, 50 $\mu$g aliquots of labelled protein were each solubilized in 50 $\mu$l of 50 mM ammonium acetate, pH 7.8, to which NaCl was added (0–5 M final). 2.5 $\mu$l of a chymotrypsin solution (0.1 mg/ml) were added to the six reaction vessels. Proteolysis conditions were as follows: 30 min, 20°C, reaction stopped by PMSF (1 mM final). After electrophoretic examination, a 500 $\mu$g sample of the labelled substrate was hydrolysed as described above in the selected reaction conditions (5 M NaCl).

Similarly, 50 $\mu$g of labelled protein were dissolved in 12.5 mM Tris, pH 8.5, 25 mM CaCl$_2$. NaCl (0–5 M final) and 1 $\mu$l of thermolysin solution (0.2 mg/ml) were added. Reaction conditions were as follows: 0°C (to prevent precipitation of $\beta$-casein in the presence of Ca$^{2+}$ [14]), 30 min, reaction stopped by adding EDTA (50 mM final) and freezing. For the preparation of selected peptides 3 and 4 (see section 3), the same procedure was applied to 500 $\mu$g of labelled substrate in the presence of 5 M NaCl.

Approx. 10 nmol of peptide 7 were digested by 0.25 $\mu$g of chymosin as follows: 200 $\mu$l of 50 mM sodium citrate, pH 6.2, 30°C, 4 h.

From 1 to 3 $\mu$g of the various digests and the whole labelled protein (1 $\mu$g) were analysed by SDS-PAGE [15]. The gel was then treated with Amplify/glycerol (9:1) for 30 min and dried at 60°C under vacuum over 2 h. Dried gels were examined by fluorography after an exposure time of 48 h.

Transformation of phosphoserine to S-ethylcysteine in presence of ethanethiol was performed according to [16].

The labelled peptides were purified by RP-HPLC. Absorbance of the eluate was monitored at 214 nm. The column was kept at 40°C. Flow rate was 1 ml/min. The following solvent systems were used for the purification of the labelled peptides, with linear gradients between A and B. Solvent system I: A, 0.12% TFA; B, 0.10% TFA in 75% acetonitrile. Solvent system II: A, 15 mM ammonium acetate, pH 7.2; B, 25% of 50 mM ammonium acetate, pH 7.2, and 75% acetonitrile. Solvent system III: A, 0.05% TFA; B, 0.04% TFA in 75% isopropanol/acetonitrile (1:1). Solvent system IV: A, 25 mM ammonium acetate, pH 7.2; B, 25% of 50 mM ammonium acetate, pH 7.2, and 75% isopropanol/acetonitrile (1:1). Solvent system V: A, 0.05% TFA, pH adjusted to 7.5 with TEA; B, 25% A and 75% isopropanol/acetonitrile (1:1). Solvent system VI: A, 10 mM potassium phosphate, pH 7.2; B, 25% A and 75% isopropanol/acetonitrile (1:1). 1 ml fractions were collected along each run and radioactivity was measured on a Beckman LS 1800 scintillation counter. Fractions of interest were pooled, dried and resolubilized in 200 $\mu$l of 4 M Gu-HCl before injection onto HPLC whenever a second chromatographic step was necessary. Peptide 4' (action of plasmin) was purified on a C18 column ($\mu$-Bondapak, 100 Å pore size, Waters). System I was employed for the elution.

Peptide 7 (action of endoprotease Lys-C) was purified on the $\mu$-Bondapak column by using successively systems I and II. The same column was used for the purification of peptide 8 (action of chymosin on peptide 7) with system I. Purification of peptide 6 (action of IBA) was performed on a C8 pre-column, 300 Å pore size (Brownlee). Systems III and IV were employed successively. Peptide 5 (action of S. aureus V8 protease) was purified on a C18 $\mu$-Bondapak column, by using system I, then system V. Peptides 2, 3 and 4 (partial digestions by chymotrypsin and thermolysin) were all purified on a C8 pre-column, 300 Å pore size. For each of them, systems VI and III were employed. Finally, peptide 1 (labelled $\beta$-casein with modified phosphoserine residues) was directly purified on the above pre-column with system III.

Automated Edman degradation of peptides 7,6,5,4',4,3,2 and 1 was performed using an Applied Biosystems 470 A sequenator connected to a PTH HPLC-analyzer (model 120 A) using the manufacturer's reagents and methods. Peptide 8 was manually sequenced [17].

## 3. APPLICATION TO BOVINE $\beta$-CASEIN

Enzymatic labelling does not necessarily occur at the original C-terminal residue. This is the case for $\beta$-casein (209 residues) [6]: labelling occurs after proteolytic removal of the last 2 residues, Ile-207 being replaced afterwards by the exogenous labelled amino acid. However, the rates of release of the C-terminal residues by CPD-Y and CPD-A from the native protein (not shown) clearly indicated that the original C-terminal sequence was F-P-I-I-V.

The C-terminally labelled protein was subjected to the cleavage procedures described above. Fluorographic examination of SDS-PAGE showed one single labelled proteolysis product in the lanes corresponding to IBA and endoproteinase Lys-C hydrolysates. Two labelled peptides were present after digestion by $S.$ $aureus$ V8 protease and plasmin.

With $\alpha$-chymotrypsin and thermolysin, it appeared that the proteolysis level decreased with increasing ionic strength (fig.1). In 0 and 1 M NaCl, there was no residual labelled $\beta$-casein left and only low molecular mass labelled peptides were present. From 2 to 5 M NaCl, well resolved bands appeared corresponding to several proteolysis products differing only at their N-terminal part. Chymotryptic digestion at 20°C and the highest salt concentration (5 M NaCl) was selected since it yielded a major labelled peptide (peptide 2), the molecular mass of which was quite different from those detected in the previously mentioned hydrolysates. Similarly, proteolysis by thermolysin in 5 M NaCl at 0°C produced two main labelled peptides (peptides 3 and 4), completing the set of peptides of different lengths chosen previously. The proteolysis conditions selected for each cleavage agent were then used on 500 $\mu$g $\beta$-casein for the purification of the labelled peptides. Electrophoresis demonstrated that carrying out the reactions on a 10-fold higher quantity of substrate did not change the hydrolysis patterns (at least for the molecules detected by fluorography).

Details of the purification steps used for the different labelled peptides were given above. As far as the stationary phase is concerned, several hydrocarbon chains of varying lengths were tested (C1, C4, C8 and C18). From C4 to C18, no significant difference was found for the separation of the
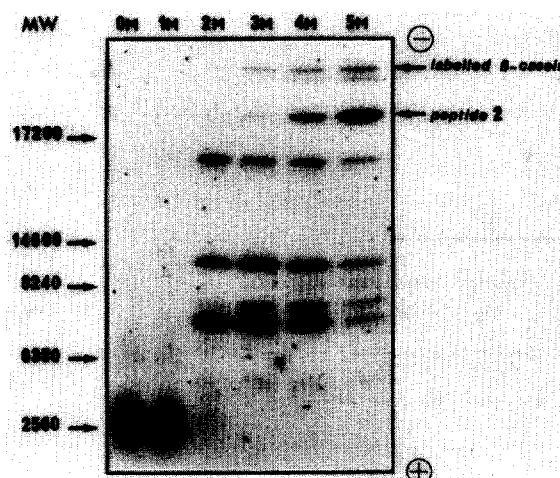


Fig.1. Influence of the ionic strength on the cleavage of labelled $\beta$-casein by chymotrypsin. The six lanes from left to right correspond to increasing concentrations of NaCl (0–5 M). See text for other details.

hydrolysates. With a C1 column, a decrease in the retention times was observed. Resolution was not improved, even for the separation of highly hydrophobic molecules. Moreover, the results obtained with a large pore size (300 Å) C18 column were very similar to those obtained with a classical pore size (100 Å) C18 column from the same manufacturer. Similar results have been observed [18].
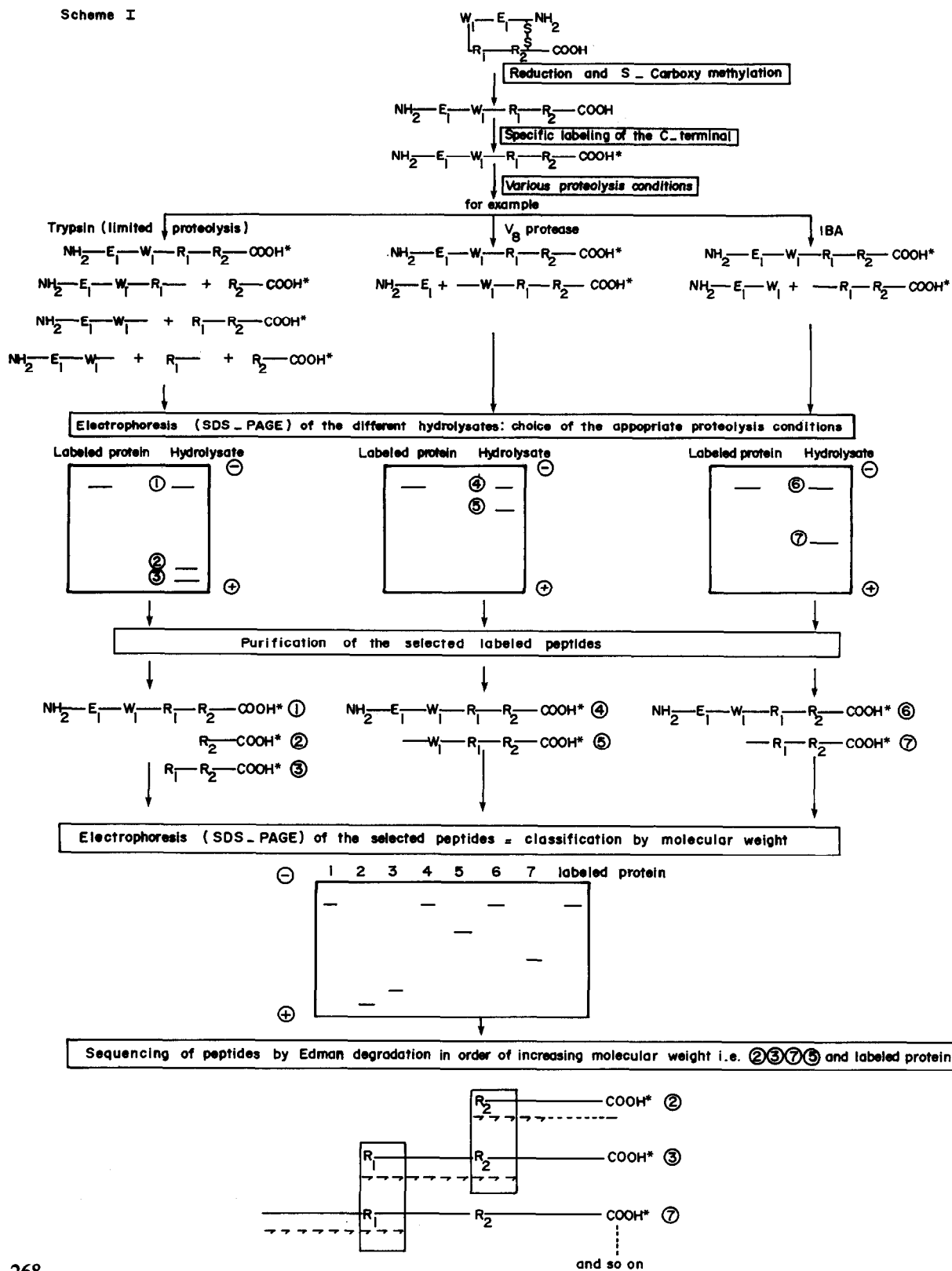
The length of the column seemed to be an important parameter. If a 250 × 4.6 mm column gave satisfactory results for the separation of some hydrolysates (chymosin, endoproteinase Lys-C, $S.$ $aureus$ V8 protease), a shorter column (30 × 4.6 mm) allowed much better separation of very hydrophobic peptides (for example those obtained after the action of chymotrypsin, thermolysin or IBA).

Neutral or acidic mobile phases were used successfully one after another. A volatile mobile phase was always employed for the last chromatographic step prior to sequencing (e.g. ammonium acetate for neutral pH and TFA for acidic pH).

Finally, for the separation of hydrophobic peptides, better results were obtained by using a mixture of acetonitrile and isopropanol in solvent B.

In order to confirm the $M_r$ of the selected pep-

Scheme I



Reduction and S _ Carboxy methylation

Specific labeling of the C_terminal

Various proteolysis conditions

for example

Trypsin (limited proteolysis)

$V_8$ protease

IBA

Electrophoresis (SDS _ PAGE) of the different hydrolysates: choice of the appopriate proteolysis conditions

Labeled protein   Hydrolysate

Purification of the selected labeled peptides

Electrophoresis (SDS _ PAGE) of the selected peptides = classification by molecular weight

⊖    I   2   3   4   5   6   7    labeled protein

Sequencing of peptides by Edman degradation in order of increasing molecular weight i.e. ②③⑦⑤ and labeled protein

and so on

tides following purification by HPLC, it was necessary to perform a final electrophoresis (SDS-PAGE) in which their mobilities were compared with those observed earlier (fig.2). Peptides 7, 6 and 4′ appeared free of any contamination by another labelled peptide. Peptides 5 and 2 were both slightly contaminated by another radioactive peptide. Contamination was more important for peptides 3 and 4.

All these peptides were sequenced in order of increasing size (i.e. 7, 6, 5, 4, 3, 2 and labelled protein; fig.2). Peptide 7 could not be sequenced fully automatically. Starting with 1.5 nmol of this material, 20 residues were determined. A smaller labelled peptide (peptide 8) was generated by hydrolysing peptide 7 with chymosin. Following a one-step purification by RP-HPLC, peptide 8 was satisfactorily sequenced manually. Automatic sequencing was not possible. Determination of radioactivity in the effluents from the reaction chamber clearly indicated that this molecule was rapidly washed out from the filter, particularly during the extraction of the ATZ by chlorobutane. It was verified that only the PTH-amino acid obtained at the last cycle (PTH-Phe at cycle 15) was radioactive. The overlap of 3 residues between peptides 7 and 8 (see fig.3) allowed the determination of the C-terminal sequence of the labelled protein.

Modifications of Met residues were observed on peptide 6. This resulted from the action of IBA on the protein (cleavage at Trp). There was no intact Met present. Part (~30%) was transformed into Met sulfone and at least two other non-identified species were present, as described [19].

Two sequences were simultaneously obtained when sequencing peptides 3 and 4. The low proportion of the contaminating species (~15%) did not interfere with the determination of the sequence of peptide 3. For peptide 4 the contaminating peptide was in larger proportion (~30%) and, as a result, the sequence could not be determined far enough. 17 Edman cycles were performed. To establish the overlap with peptide 5, a peptide obtained by the action of plasmin on labelled β-casein (peptide 4′) and whose size was between those of peptides 4 and 5 (fig.2) was purified and sequenced.

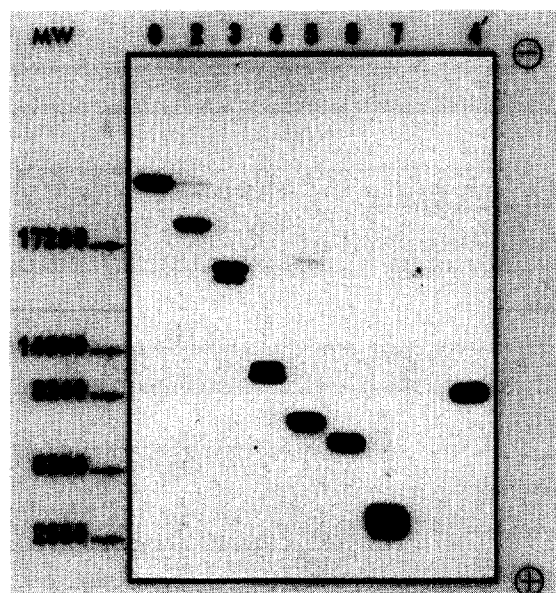Finally, no PTH-amino acid was detected at cycles 15, 17, 18 and 19 during the sequencing of



Fig.2. Electrophoresis of the set of purified peptides selected for the determination of the sequence of the labelled β-casein. The following cleavage agents were used: 7, endoproteinase Lys-C; 6, IBA; 5, *S. aureus* V8 protease; 4 and 3, thermolysin; 2, chymotrypsin; 4′, chymosin.

the original labelled protein. The yield decreased considerably after cycle 19 and, even using a large amount of protein (2 nmol), it was not possible to perform more than 22 cycles. This was probably due to the presence of phosphoamino acids at these positions. The labelled protein was treated with ethanethiol [16] (see above) and the reaction mixture was analyzed by RP-HPLC. In addition to the original protein, a second fraction, named peptide 1, was observed and recovered. Its N-terminal sequence was identical to that of the initial labelled protein for the first 14 residues. At cycles 15, 17, 18, 19 and 35, PTH-*S*-ethylcysteine was released, indicating the original presence of phosphoserine at these positions, since β-casein contains phosphorus but no carbohydrates.

Although the overlap with peptide 2 started at residue 29, the sequence of peptide 1 was conducted to residue 35, which was not determined as previously mentioned.

The whole sequence obtained with peptides 1 to 8 is given on fig.4. Some differences from the initial sequence published [7] were noted (position 117: E instead of Q; 137: L instead of P; 138: P in-
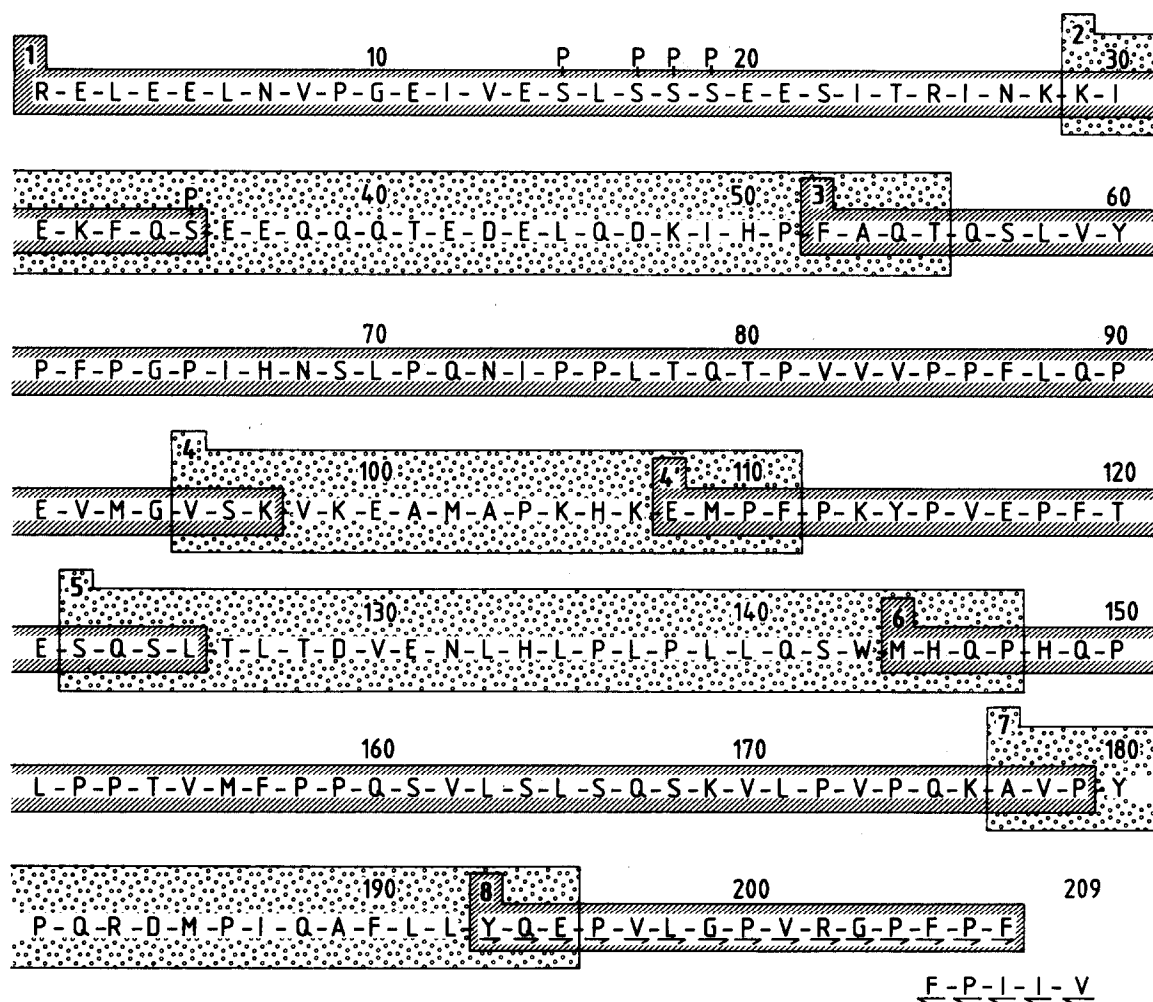
Fig.3. Complete sequence of bovine β-casein, determined as described in the text. For peptides 1–7, residues contained in the frames were automatically sequenced. Peptide 8 was sequenced manually (→ →). The C-terminal sequence of the non-labelled protein was determined by using carboxypeptidases A and Y (← ←) (see text for other details).

stead of L; 176: Q instead of E and 195: E instead of Q). These corrections are in agreement with a recently published cDNA sequence [9], but not with that published in [8].

The whole sequence of the labelled protein was determined with only 239 Edman degradations (for 207 residues).

## 4. DISCUSSION

So far as we know, no rational fragment overlapping strategy has been previously proposed. The approach which has been used here is based on a reasoned choice of a limited number of peptides which require only partial sequencing, with the exception of the C-terminal one which must be fully sequenced. The selection of these peptides on the basis of their approximate $M_r$ confers a considerable degree of flexibility to this approach. If the overlap between two consecutive peptides is not encountered (such an example is given with peptides 4 and 5), another peptide, from a different hydrolysate, whose size is intermediate (in our case peptide 4′), can be used to overlap the two peptides in question. The excess of Edman cycles used to obtain the complete sequence over

the theoretical number (i.e. the number of residues) corresponds to the minimum overlap required.

As far as the adaptation of this method to large proteins is concerned, three likely difficulties can be envisaged. Firstly, the electrophoretic separation of high molecular mass fragments differing only by ~45 residues could be difficult. Secondly, purification by RP-HPLC of large protein fragments of close $M_r$ can be problematic. Thirdly, the higher the $M_r$ of a protein, the more likely an increase in the frequency will be of regions which are either more resistant or more sensitive to proteolysis. The existence of large resistant areas can make it impossible to obtain peptides of the required size, whereas the presence of proteolysis-sensitive regions may make it very difficult to obtain labelled peptides larger than the fragments resulting from cleavage in the sensitivity zones.

Cleavage reactions carried out to completion under classical conditions should be performed with cleaving agents having uncommon target residues so that the likelihood of obtaining large C-terminal peptides is higher, e.g. hydroxylamine for N-G bonds [20]; acidic conditions for D-P [21]; IBA for cleavage at Trp; BrCN at Met [22]; collagenases; plasmin.... For limited proteolysis, it appears that enzymes as different as chymotrypsin and thermolysin cleave with a high yield only a few peptide bonds of $\beta$-casein in the presence of high concentrations of NaCl (2–5 M) (fig.1). Such conditions have been applied with success to other enzymes (*S. aureus* V8 protease, endoproteinase Lys-C) (not shown; to be published). It is feasible to consider that such reaction conditions could be used with any other protein substrates.

The sensitivity of this method should be taken into consideration. Only small amounts of labelled protein are necessary for the selection of proteolysis conditions. The reactions were initially carried out with 50 $\mu$g protein (2 nmol) because the protein was available in large quantities, but this amount could have easily been reduced. For example, some reactions were performed on 0.5 nmol. Thus, even if 20 different hydrolysis conditions were examined, this would consume only 10 nmol of protein.

To obtain ~1 nmol of each labelled peptide, it should be enough to hydrolyse an average quantity of 15 to 20 nmol for each selected condition.

150 nmol of labelled $\beta$-casein were used for the purification of the 8 peptides necessary for the determination of the whole primary structure.

The method described here may have some limitations. As previously mentioned, areas of a protein can be resistant or very sensitive to proteolysis and then make difficult the obtention of some peptides of suitable $M_r$. Should it not be possible to obtain the entire sequence, our approach will still provide valuable structural details.

Another limiting factor is the possible contamination of the labelled peptides by non-labelled ones before sequencing. In this instance, fluorographic examination is ineffective. In order to decrease the risk of contamination, the labelled peptides are chromatographed in at least two different solvent systems. Anyway, it should be noticed that a rate of contamination less than 15–20% does not prevent a correct sequence determination (e.g. peptide 3 in fig.3 was correctly sequenced in spite of contamination by another labelled peptide).

## 5. CONCLUSION

Thus far, our strategy has been tested in only one case, on a known protein available in large amounts. So, this paper is, in part, a 'theoretical' suggestion which requires further study. Nevertheless, we think it can be a suitable basis for the development of further ideas and concepts.

## REFERENCES

[1] Sanger, F., Nicklen, S. and Coulson, A.R. (1977) Proc. Natl. Acad. Sci. USA 74, 5463–5467.
[2] Maxam, A.L. and Gilbert, W. (1977) Proc. Natl. Acad. Sci. USA 74, 560–564.
[3] Jay, D. (1984) J. Biol. Chem. 259, 15572–15578.
[4] Jue, R. and Doolittle, R. (1985) Biochemistry 24, 162–170.
[5] Dale, R.M.K., McClure, B.A. and Houchins, J.P. (1985) Plasmid 13, 31–40.
[6] Carles, C., Gueguen, P. and Ribadeau-Dumas, B. (1987) FEBS Lett. 212, 163–167.
[7] Ribadeau-Dumas, B., Brignon, G., Grosclaude, F. and Mercier, J.C. (1972) Eur. J. Biochem. 25, 505–514.

[8] Jimenez-Flores, R., Yang, Y.C. and Richardson, L. (1987) Biochem. Biophys. Res. Commun. 142, 617–621.

[9] Bayev, A.A., Smirnov, I.K. and Gorodestsky, S. (1987) Mol. Biol. 21, 255–265.

[10] Zittle, C.A. and Custer, J.H. (1963) J. Dairy Sci. 46, 1183–1188.

[11] Mercier, J.C., Maubois, J.L., Pozananski, S. and Ribadeau-Dumas, B. (1968) Bull. Soc. Chim. Biol. 50, 521–530.

[12] Guillou, H., Miranda, G. and Pelissier, J.P. (1987) Le Lait 67, 135–148.

[13] Mahoney, W.C., Smith, P.K. and Hermodson, M.A. (1981) Biochemistry 20, 443–448.

[14] Schmidt, D.G. (1982) in: Developments in Dairy Chemistry (Fox, P.F. ed.) pp.61–86, Applied Science Publishers, London.

[15] Fling, S.P. and Gregerson, D.S. (1986) Anal. Biochem. 155, 83–88.

[16] Meyer, H., Hoffman-Posorske, E., Korte, H. and Heilmeyer, L. (1986) FEBS Lett. 204, 61–66.

[17] Tarr, G.E. (1982) in: Methods in Protein Sequence Analysis (Elzinga, M. ed.) pp.223–232, Humana Press, Preston, NJ.

[18] Van der Zee, R., Hoekzema, T., Welling-Wester, S. and Welling, G. (1986) J. Chromatogr. 368, 283–289.

[19] Fontana, A., Dalzoppo, D., Grandi, C. and Zambonin, M. (1982) in: Methods in Protein Sequence Analysis (Elzinga, M. ed.) pp.325–334, Humana Press, Preston, NJ.

[20] Bornstein, P. (1969) Biochem. Biophys. Res. Commun. 36, 957–964.

[21] Piszkiewicz, D., Landon, M. and Smith, E. (1970) Biochem. Biophys. Res. Commun. 40, 1173–1178.

[22] Gross, E. and Witkop, B. (1962) J. Biol. Chem. 237, 1856–1860.