

# N-Terminal pI Determines the Solubility of a Recombinant Protein Lacking an Internal Transmembrane-like Domain in *E. coli*

Sang Jun Lee\*, Yun Hee Han, Young Ok Kim, Bo Hye Nam, Hee Jeong Kong, and Kyung Kil Kim

We examined whether the isoelectric point (pI) of the N-terminal region of the recombinant protein 7xMefp1 acts as a universal index for expression of the protein in soluble form in *E. coli*. Expression analysis of 7xMefp1 fused to various N-terminal sequences with pI values ranging from 2.73 to 13.35 yielded three pI range-specific curves (acidic, neutral, and alkaline curves at pI 2.73–3.25, 4.61–9.58, and 9.90–13.35, respectively) for soluble expression (by facilitated diffusion) as a proportion of total protein. For neutral N-termini (pI 4.61–9.58), the total amount of rMefp1 expressed was minimally affected by  $\Delta G_{\text{RNA}}$  for unfolding the mRNA secondary structure. The highly hydrophilic nature of longer N-terminal sequences with strongly acidic and alkaline pI values reduced the translation of rMefp1-encoding transcripts, thereby reducing the amount of soluble rMefp1 produced. After characterizing both feedback and non-feedback regulation in the acidic, alkaline, and neutral pI ranges, we suggest that three different pI range-specific soluble expression curves exist for the recombinant protein, each defined by specific ranges of the leader sequence pI values.

## INTRODUCTION

Many researchers believe that the expression of recombinant proteins in soluble form in *E. coli* is inherently dependent upon the physical properties (e.g., isoelectric point (pI), hydrophobicity, stability, and molecular weight) of the entire amino acid sequence. Expression may also be affected by the presence of positive charges in the amino (N)-terminal region of the signal peptide (Inouye et al., 1982), by mRNA secondary structure (Mukund et al., 1999; Ramesh et al., 1994), and by codon usage (Ikemura, 1981). Although each of these factors has individually been found to affect the expression of specific proteins, they have not been applicable to the overall problem.

In attempting to develop a more general approach to solving the problem, we previously used hydropathy profile analysis to demonstrate that the presence of an internal, positively charged transmembrane (TM)-like domain in the olive flounder protein hepcidin I inhibits its soluble expression (Lee et al., 2008b).

Using a novel secretion enhancer, we were able to overcome the obstacle posed by the internal TM-like domain and successfully expressed hepcidin I in soluble form. We also investigated Mefp1, an adhesive protein of the marine mussel *Mytilus edulis* (Waite, 1983). This large protein consists primarily of repeated decapeptide units lacking a TM-like domain (Lee et al., 2008b). We found that the addition of a Met residue to the N-terminus allowed low-level expression of recombinant 7xMefp1, a polypeptide comprising seven copies of the decapeptide repeat unit, in soluble form (Lee et al., 2008a) in *E. coli*. When we fused 7xMefp1 to a truncated OmpA signal peptide (OmpASP<sub>tr</sub>) in an attempt to increase its expression in soluble form, we found that as the pI value of the N-terminus of the recombinant 7xMefp1 fusion protein (rMefp1) increased from 9.90 to 10.82, periplasmic expression of soluble rMefp1 also increased (Lee et al., 2008a).

However, the importance of the pI value of N-terminal regions in soluble expression of proteins lacking TM-like domains has not yet been properly investigated. Therefore, in the present study, we investigated the influence of the pI, hydrophobicity, charge,  $\Delta G_{\text{RNA}}$  for mRNA unfolding ( $\Delta G_{\text{RNA}}$ ), codon usage, codon repetition, and molecular weight of the N-terminus on soluble rMefp1 expression in *E. coli*. We found that the N-terminal pI value can be used, over a wide range (2.73–13.35), as a comprehensive biological index to predict the level of soluble rMefp1 expression.

## MATERIALS AND METHODS

### Bacterial strains and plasmids

The *E. coli* strains XL-1 blue (Stratagene) and TOP10 (Invitrogen) were used for cloning, and BL21 (DE3) (Novagen) for direct expression of the fusion protein. The plasmid pBluescript- $\text{II SK}^{+}$  (Stratagene) and the TA cloning vector (Promega) were used for cloning, and pET-22b(+) (Novagen) for protein expression.

### Reagents and molecular techniques

Restriction endonucleases (Roche) and a His-tag purification kit (Qiagen) were used. All other chemicals were of analytical grades. All molecular techniques were conducted as described

Biotechnology Research Division, National Fisheries Research and Development Institute, Pusan 619-902, Korea

\*Correspondence: sangji@nfrdi.go.kr

Received February 10, 2010; revised April 16, 2010; accepted April 26, 2010; published online July 23, 2010

**Keywords:**  $\Delta G_{\text{RNA}}$  value, feedback regulatory mechanism, hydrophilicity, inner membrane channel, pI value trigger (facilitated diffusion)

by Sambrook et al. (1989). Nucleotide sequencing using the dideoxy chain-termination method (Sanger et al., 1977) was performed using the Sequenase 2.0 kit (United States Biochemical). The computer program DNASIS™ (Hitachi, Japan, 1997) was used to analyze the pI characteristics and hydrophobicity of polypeptide sequences.

### Construction of expression vectors

To construct vectors for the expression of the 7xMefp1 target protein fused to N-terminal sequences with various pI values, we designed primer pairs specific for the individual N-termini (Supplementary Table 1). These primers were used to amplify DNA cassettes using the control N-terminal clone pET-22b(+)(*ompASP<sub>1</sub>(Met)-7xmefp1\**) as a template (Lee et al., 2008a). The amplified DNA was cloned into a TA cloning vector, the entire *NdeI-XhoI* fragment of which was then subcloned into pET-22b(+) by replacing the *pel* signal sequence and the poly-linker as described previously (Lee et al., 2008a). The resulting constructs are listed in Table 1. To facilitate Western blot analysis, a C-terminal His-tag was added to each of the fusion proteins in the pET-22b(+) subcloning step.

### Protein expression

*E. coli* BL21 (DE3) cells were transformed with the plasmid constructs listed in Table 1, and transformants were cultured in LB medium overnight at 30°C in the presence of 100 µg/ml ampicillin. The culture was then diluted 1:100 in LB medium and grown until it reached an optical density of 0.6 at 600 nm. Then, isopropyl-β-D-thiogalactopyranoside (IPTG) was added to a final concentration of 1 mM, and the culture was grown for another 3 h to allow expression of the recombinant protein. A 1-ml aliquot was then removed from each culture and centrifuged. Cell pellets were resuspended in 100-200 µl of PBS. Cells were disrupted by sonication, in which 15 pulses at 30% power output were applied in 2-s cycles, and then centrifuged at 16,000 rpm for 30 min at 4°C. (The resulting supernatant fractions contained the soluble protein fraction.) The insoluble pellets were resuspended in an equal volume to that of the supernatants. To prepare periplasmic fractions, IPTG-induced cells were subjected to osmotic shock as described by Nossal and Heppel (1966). Protein fractions were quantified by a Bradford assay (Bradford, 1976), separated by SDS-PAGE on 15% acrylamide gels (Laemmli, 1970), and visualized using Coomassie brilliant blue stain. Band intensities were determined densitometrically using Quantity One 1-D image analysis software (Bio-Rad).

### Western blot analysis

After electrophoresis, His-tagged rMefp1 proteins were transferred to a Hybond-P membrane (GE Healthcare) and detected using (in sequence) an anti-His tag (C-term) primary antibody, an alkaline phosphatase-conjugated anti-mouse secondary antibody, and a chromogenic Western blotting kit (Invitrogen), according to the manufacturers' protocols. Molecular weight markers (Benchmark, His-tagged) were used. Band intensities were quantified densitometrically using Quantity One 1-D image analysis software.

### Determination of N-terminal $\Delta G_{\text{RNA}}$ values

The  $\Delta G_{\text{RNA}}$  values for unfolding the RNA secondary structures of the N-termini were calculated using the program MFOLD 3 (www.bioinfo.rpi.edu/applications/mfold) (Zuker, 2003) (Table 1 and Supplementary Table 1). To investigate the effects of  $\Delta G_{\text{RNA}}$  on soluble expression, we compared longer N-termini (which generally had smaller  $\Delta G_{\text{RNA}}$  values) to shorter N-termini

(which generally had larger  $\Delta G_{\text{RNA}}$  values). To alter  $\Delta G_{\text{RNA}}$ , we used synonymous codons. For example, in the (Glu)<sub>n</sub>-containing N-termini ME<sub>8</sub>-I and ME<sub>6</sub>-I, the Glu<sup>GAA</sup> and Glu<sup>GAG</sup> codons are used with equal frequency, resulting in  $\Delta G_{\text{RNA}}$  values of -8.30 and -8.50, respectively. In contrast, in ME<sub>8</sub>-II and ME<sub>6</sub>-II, the Glu<sup>GAA</sup> codon is used for all but one Glu codon, reducing the  $\Delta G_{\text{RNA}}$  values to -4.00 and -4.30, respectively. Similarly, in the (Lys)<sub>n</sub>-containing N-termini, we used Lys<sup>AAA</sup> codons to reduce  $\Delta G_{\text{RNA}}$  values. In the (Arg)<sub>n</sub>-containing N-termini, we used the Arg<sup>AGA</sup> codon to reduce  $\Delta G_{\text{RNA}}$  values, but it reduced translational efficiency when used as the sole Arg codon in MR<sub>2</sub>AK (Supplementary Fig. S2) and caused translational arrest when used as the sole Arg codon in MR<sub>4</sub>AK, MR<sub>6</sub>AK, and MR<sub>8</sub>AK (Supplementary Fig. S2). Therefore, we instead used alternating Arg<sup>CGT</sup> and Arg<sup>CGC</sup> codons, as detailed in Supplementary Table 1.

## RESULTS AND DISCUSSION

### Effects of N-termini with high pI values (9.90-13.35) on the expression of soluble rMefp1

We previously showed that changing the pI value of the N-terminus of the 7xMefp1 fusion protein rMefp1 (by altering the number of Lys residues) also modulated levels of periplasmic soluble expression (Lee et al., 2008a). Here, to assess the general relationship between N-terminal pI value and soluble rMefp1 expression, using the control clone pET-22b(+)(*ompASP<sub>1</sub>(Met)-7xmefp1\**) (Lee et al., 2008a), we constructed a series of rMefp1 fusion protein clones in which various numbers of Lys or Arg residues were inserted between the *OmpASP<sub>1</sub>(Met)* truncated signal sequence and the 7xMefp1\* sequence. The clones pET-22b(+)(*ompASP<sub>1</sub>(Lys)<sub>n</sub>-7xmefp1\**) (where *n* = 0 (control), 1, 2, 3, 4, 5, 6, or 8) and pET-22b(+)(*ompASP<sub>1</sub>(Arg)<sub>n</sub>-7xmefp1\**) (where *n* = 1, 2, 4, 6, or 8) yielded proteins with N-terminal pI values of 9.90, 10.55, 10.82, 10.99, 11.11, 11.21, 11.28, and 11.41 (for (Lys)<sub>0</sub> (control) to (Lys)<sub>8</sub>, respectively) and of 11.52, 12.51, 12.98, 13.20, and 13.35 (for (Arg)<sub>1</sub>-(Arg)<sub>8</sub>, respectively). We analyzed the pI, hydrophobicity, charge, and  $\Delta G_{\text{RNA}}$  of these N-termini in the *OmpASP<sub>1</sub>(Met)-(Lys/Arg)<sub>n</sub>-Ala<sup>1</sup>-Lys<sup>2</sup>* clones (where Ala<sup>1</sup>-Lys<sup>2</sup> represents the beginning of the native Mefp1 sequence) as previously described (Lee et al., 2008a) (Table 1 and Supplementary Table 1).

Soluble rMefp1 was expressed at higher levels with (Lys)<sub>n</sub>-containing N-termini with pI values ranging from 10.55 to 11.28 (but not the one with a pI value 11.41) than with the control N-terminal sequence *OmpASP<sub>1</sub>(Met)-Ala<sup>1</sup>-Lys<sup>2</sup>* (where Ala is encoded by GCT; sequence name MAK-I; pI 9.90). The highest level of soluble rMefp1 expression was obtained with the N-termini with pI values of 10.82, 10.99, and 11.21 (expression decreased slightly at pI 11.28 and then substantially at pI 11.41) (Fig. 1A, Table 1, and Supplementary Fig. S1D). Therefore, decreased expression with longer N-termini cannot alone be explained by the high number of positive charges.

In the case of the (Arg)<sub>n</sub>-containing N-termini, the highest expression of soluble rMefp1 was obtained using an N-terminus containing a single inserted Arg<sup>AGA</sup> (Met-Arg<sup>AGA</sup>-Ala<sup>1</sup>-Lys<sup>2</sup>; sequence name MRAK-I; pI 11.52). As the number of Arg residues continued to increase (up to an N-terminal pI of 13.25), expression gradually decreased (Fig. 1A, Table 1, and Supplementary Fig. S1E). Again, the decreased expression of soluble rMefp1 with longer N-termini could not be readily explained by the high positive charge.

At present, we do not know the basis for these observed changes in soluble rMefp1 expression. The (Lys)<sub>n</sub>-containing N-termini cover a narrow pI range (9.90-11.41) that encompasses

**Table 1.** Characteristics of the various N-terminal sequences used to create recombinant 7xMefp1<sup>a</sup> fusion proteins (rMefp1) and relative levels of the soluble and insoluble forms of the corresponding fusion proteins synthesized from derivatives of the expression plasmid pET-22b(+)(*ompASP<sub>1</sub>*(Met)-7x*mefp1*)<sup>b</sup>.

Seq. No.	N-terminal sequence	pI <sup>c</sup>	Hydrophobicity index <sup>d</sup>	Charge	$\Delta G_{RNA}$ <sup>e</sup>	Soluble rMefp1 level <sup>f</sup>	Insoluble rMefp1 level <sup>f</sup>	Figure
1	MD <sub>5</sub> AA	2.73	1.09	-5	≤ -9.90	0.50	0.83	1A, S1A
2	MD <sub>3</sub> AA	2.87	0.56	-3	≤ -8.80	0.91	1.20	1A, S1A
3	MDA	3.00	N/A <sup>h</sup>	-1	≤ -6.70	1.40	1.59	1A, S1A
4	ME <sub>8</sub> -I	2.75	2.08	-8	≤ -8.30	0.49	0.42	1A, 2C, S1A
5	ME <sub>8</sub> -II	2.75	2.08	-8	≤ -4.00	N/D <sup>i</sup>	N/D <sup>i</sup>	2C
6	ME <sub>6</sub> -I	2.82	1.82	-6	-8.50	0.65	0.62	1A, 2C, S1A
7	ME <sub>6</sub> -II	2.82	1.82	-6	≤ -4.30	N/D <sup>i</sup>	N/D <sup>i</sup>	2C
8	ME <sub>4</sub>	2.92	N/A <sup>h</sup>	-4	≤ -5.80	0.79	0.66	1A, S1A
9	MEE	3.09	N/A <sup>h</sup>	-2	-5.80	1.42	1.81	1A, S1A
10	MAE	3.25	N/A <sup>h</sup>	-1	≤ -5.60	1.72	1.92	1A, S1A
11	MC <sub>6</sub>	4.61	-0.64	-	≤ -6.80	1.65	2.04	1A, S1B
12	MC <sub>3</sub>	4.75	N/A <sup>h</sup>	-	≤ -7.10	1.93	2.95	1A, S1B
13	MAC	4.83	N/A <sup>h</sup>	-	≤ -5.20	1.96	1.84	1A, S1B
14	MAY	5.16	N/A <sup>h</sup>	-	-5.20	1.74	1.73	1A, S1B
15	MAA	5.60	N/A <sup>h</sup>	-	-5.60	2.25	2.23	1A, S1B
16	MNN	5.70	N/A <sup>h</sup>	-	≤ -2.50	N/D <sup>i</sup>	N/D <sup>i</sup>	2B
17	MTT	5.70	N/A <sup>h</sup>	-	≤ -3.10	N/D <sup>i</sup>	N/D <sup>i</sup>	2B
18	MWW	5.85	N/A <sup>h</sup>	-	≤ -7.50	N/D <sup>i</sup>	N/D <sup>i</sup>	2B
19	MGG	5.85	N/A <sup>h</sup>	-	≤ -7.80	1.93	2.10	1A, 2B, S1B
20	MAKD	6.59	N/A <sup>h</sup>	0	-7.00	2.30	2.63	1A, S1B
21	MAKE	6.79	N/A <sup>h</sup>	0	≤ -5.80	2.05	2.99	1A, S1B
22	MCH	7.13	N/A <sup>h</sup>	-	-3.60	1.83	3.11	1A, S1C
23	MAH	7.65	N/A <sup>h</sup>	-	-5.20	1.81	2.72	1A, S1C
24	MAH <sub>3</sub>	7.89	N/A <sup>h</sup>	-	≤ -5.40	1.54	3.75	1A, S1C
25	MAH <sub>5</sub>	8.01	-0.33	-	-6.70	1.37	4.36	1A, S1C
26	MAKC	8.78	N/A <sup>h</sup>	+1	-7.80	1.73	2.98	1A, S1C
27	MKY	9.58	N/A <sup>h</sup>	+1	≤ -2.50	1.51	4.04	1A, 2B, S1C
28	MAKY	9.58	N/A <sup>h</sup>	+1	-5.20	N/D <sup>i</sup>	N/D <sup>i</sup>	2B
29 <sup>j</sup>	MAK-I (Ala: GCT) (control) <sup>g</sup>	9.90	N/A <sup>h</sup>	+1	-7.80	1.00	1.00	1A, 2A-C, S1A-E
30	MAK-II (Ala: GCA)	9.90	N/A <sup>h</sup>	+1	-7.80	N/D <sup>i</sup>	N/D <sup>i</sup>	2A
31	MAK-III (Ala: GCC)	9.90	N/A <sup>h</sup>	+1	-7.80	N/D <sup>i</sup>	N/D <sup>i</sup>	2A
32	MAK-IV (Ala: GCG)	9.90	N/A <sup>h</sup>	+1	-7.80	N/D <sup>i</sup>	N/D <sup>i</sup>	2A
33 <sup>j</sup>	MKAK	10.55	N/A <sup>h</sup>	+2	≤ -2.40	1.57	3.14	1A, S1D
34	MK <sub>2</sub> AK-I	10.82	N/A <sup>h</sup>	+3	≤ -2.10	1.69	3.43	1A, 2C, S1D
35 <sup>j</sup>	MK <sub>2</sub> AK-II	10.82	N/A <sup>h</sup>	+3	≤ -4.70	N/D <sup>i</sup>	N/D <sup>i</sup>	2C
36	MK <sub>3</sub> AK	10.99	1.14	+4	≤ -2.00	1.80	2.96	1A, S1D
37	MK <sub>4</sub> AK	11.11	1.32	+5	≤ -1.80	1.72	2.59	1A, S1D
38	MK <sub>5</sub> AK	11.21	1.53	+6	≤ -1.70	1.93	2.34	1A, S1D
39	MK <sub>6</sub> AK	11.28	1.69	+7	≤ -1.56	1.39	1.67	1A, S1D
40	MK <sub>8</sub> AK	11.41	1.93	+9	≤ -1.39	0.44	0.39	1A, S1D
41	MRAK-I	11.52	N/A <sup>h</sup>	+2	≤ -2.90	1.69	2.92	1A, 2C, S1E
42	MRAK-II	11.52	N/A <sup>h</sup>	+2	≤ -5.00	N/D <sup>i</sup>	N/D <sup>i</sup>	2C
43	MR <sub>2</sub> AK	12.51	N/A <sup>h</sup>	+3	≤ -5.10	1.26	1.27	1A, S1E
44	MR <sub>4</sub> AK	12.98	1.32	+5	-7.70	1.07	1.14	1A, S1E
45	MR <sub>6</sub> AK	13.20	1.69	+7	-10.20	0.93	0.88	1A, S1E
46	MR <sub>8</sub> AK	13.35	1.93	+9	≤ -11.70	0.55	0.63	1A, S1E

<sup>a</sup>The Mefp1 decapeptide sequence is Ala-Lys-Pro-Ser-Tyr-Pro-Pro-Thr-Tyr-Lys (Waite, 1983).

<sup>b</sup>The plasmid pET-22b(+)(*ompASP<sub>1</sub>*(Met)-7x*mefp1*) (Lee et al., 2008a) encodes the control N-terminal sequence (MAK-I) and was used as the basis for cloning all of the N-terminal variants.

<sup>c</sup>pI values were calculated using DNASIS software.

<sup>d</sup>Hydrophobicity index was calculated using DNASIS software by the Hopp-Woods method, with a window size of 6 and a threshold line of 0.00. On the Hopp-Woods scale, hydrophilic regions are given a positive value, and hydrophobic regions a negative one.

<sup>e</sup>The DNA sequences used to calculate the  $\Delta G_{RNA}$  values for unfolding of the mRNA species encoding the fusion proteins begin near the translation initiation codon region of pET-22b(+) (5'-AAG AAG GAG ATA TA-3') and include the forward primer sequences listed in Supplementary Table 1. MFOLD 3 software (Zuker, 2003) was used to calculate  $\Delta G_{RNA}$  values. In cases in which the program yielded multiple  $\Delta G_{RNA}$  values for a particular sequence (due to the existence of multiple possible folded structures), the higher/highest  $\Delta G_{RNA}$  value is shown.

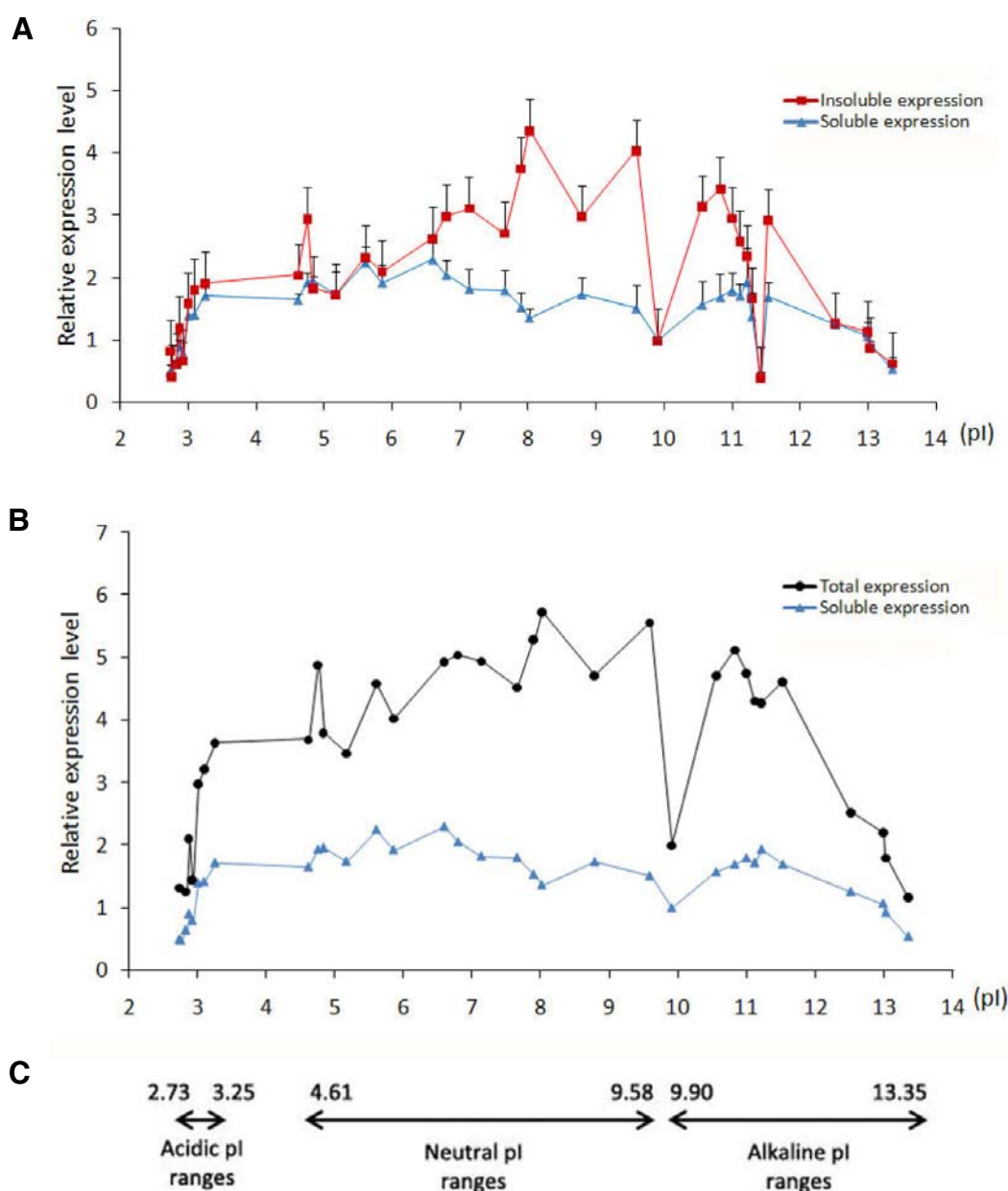
<sup>f</sup>Values indicate mean levels of soluble and insoluble fusion proteins, expressed relative to the levels of soluble and insoluble MAK-I control fusion protein, respectively, which were given values of 1.00 (Fig. 1).

<sup>g</sup>The control sequence MAK-I consists of the N-terminus of Met (OmpASP<sub>1</sub>) + Ala-Lys (the first two amino acids of Mefp1) (Lee et al., 2008a).

<sup>h</sup>Not applicable because the peptide is shorter than the minimum window size of 6 required for Hopp-Woods hydrophobicity calculations.

<sup>i</sup>Not detected.

<sup>j</sup>Primer constructed previously (Lee et al., 2008a).



**Fig. 1.** Western blot analysis of insoluble (■), soluble (▲), and total (●) rMefp1 expression induced by N-terminal peptide sequences with various pI values. (Representative Western blot images showing insoluble and soluble rMefp1 expression induced by N-terminal peptide sequences with various pI values are presented in Supplementary Figs. S1A-S1E.) Gels were loaded with protein (approximately 20  $\mu$ g per well) from the soluble fraction obtained from each clone and with an equal volume of the insoluble fraction. Anti-His-tag antiserum was used to detect the rMefp1 synthesized from a derivative of the *E. coli* expression vector pET-22b(+), which attaches a His-tag sequence to the C-terminal end of the fusion protein. (A) Expression levels were measured by densitometry relative to the levels of soluble and insoluble MAK-I (control protein; pI 9.90) (Lee et al., 2008a), which were given values of 1.00 (without correcting for the slightly higher amount of insoluble rMefp1). Data represent the mean relative levels of rMefp1 in three different samples analyzed by Western blotting (Supplementary Figs. S1A-S1E), each obtained from a different colony (as shown in Table 1). Satellite bands were excluded from the densitometric analysis. (B) The data in (A) were re-plotted as a Hubbert curve (Hubbert, 1956) to show mean total and soluble rMefp1 expression in presenting pI range-specific total and soluble expression for N-termini with pI values ranging from 9.90 to 13.35. (C) The boundaries of the acidic, neutral, and alkaline pI ranges were deduced from the soluble rMefp1 expression curve (B).

the optimal pI—the expression curve peaks within this range. This hyperbolic change in expression may be related to a pI-dependent inner membrane translocation mechanism that secretes the target precursor protein from the cytoplasm to the periplasm. Soluble expression of rMefp1 fusion proteins with

(Arg) $_n$ -containing N-termini (pI 11.52-13.35) gradually decreased as the value of  $n$  increased, suggesting the involvement of a similar inner membrane translocation mechanism.

We next examined hydrophilicity as an alternative index for the charge of the N-termini described above. Those N-termini

with pI values of 11.41 (Met-(Lys)<sub>8</sub>-AK; charge + 9), 12.98 (Met-(Arg)<sub>4</sub>-AK; charge + 5), 13.20 (Met-(Arg)<sub>6</sub>-AK; charge + 7), and 13.35 (Met-(Arg)<sub>8</sub>-AK; charge + 9), which produced only very weak expression of soluble rMefp1, are highly hydrophilic (Hopp-Woods hydrophobicity indices (*h*) (see "Materials and Methods") 1.93, 1.32, 1.69, and 1.93, respectively). Thus, highly hydrophilic N-termini appeared to reduce soluble rMefp1 expression.

However, the (Lys)<sub>*n*</sub>-containing N-termini for which *n* = 4, 5, or 6 (Met-(Lys)<sub>4</sub>-AK (charge +5, pI 11.11), Met-(Lys)<sub>5</sub>-AK (charge + 6, pI 11.21), and Met-(Arg)<sub>6</sub>-AK (charge + 7, pI 11.28)), which are also highly hydrophilic (*h* = 1.32, 1.53, and 1.69, respectively), yielded higher levels of soluble rMefp1 expression than Met-(Arg)<sub>4</sub>-AK (charge +5, pI 12.98), for which *n* = 4 and *h* = 1.32 (Fig. 1A and Table 1). The hydrophilicity of each (Lys)<sub>*n*</sub>-containing N-terminus is equal to or greater than that of the corresponding (Arg)<sub>*n*</sub>-containing N-terminus. Therefore, the (Lys)<sub>*n*</sub>-containing N-termini appear to be superior to the (Arg)<sub>*n*</sub>-containing N-termini in terms of their ability to promote soluble rMefp1 expression.

The highly hydrophilic (Lys)<sub>8</sub>-, (Arg)<sub>4</sub>-, (Arg)<sub>6</sub>-, and (Arg)<sub>8</sub>-containing N-termini failed to increase the expression of soluble rMefp1 (Fig. 1A), which lacks a TM-like domain, whereas we previously found that the insertion of (Lys)<sub>6</sub>-, (Arg)<sub>6</sub>-, (Arg)<sub>8</sub>-, or (Arg)<sub>10</sub> sequences enhanced the secretion of a target protein containing a TM-like domain (Lee et al., 2008b). Therefore, the potential of (Lys)<sub>*n*</sub>- and (Arg)<sub>*n*</sub>-containing N-termini to enhance secretion (Lee et al., 2008b) or reduce soluble expression, by increasing pI and *h*, appears to depend on whether the target protein has a TM-like domain.

#### Effects of N-termini with lower pI values (2.73-9.58) on the expression of soluble rMefp1

To assess the influence of N-termini with lower pI values on soluble rMefp1 expression, we created additional N-terminal variants of the OmpASP<sub>1</sub>(Met)-7xMefp1\* clones with pI values ranging from 2.73 to 9.58. In these clones, 2-8 amino acid residues were inserted after the OmpASP<sub>1</sub> sequence, and the first two residues of the 7xMefp1 polypeptide (Ala<sup>1</sup>-Lys<sup>2</sup>-) were replaced. (These clones were of the form pET-22b(+)(ompASP<sub>1</sub>-(X)-mefp1<sup>3-10</sup>-6xmefp1\*), where X is the insert.) These N-termini were compared with the control MAK-I sequence (Met-Ala<sup>1</sup>-Lys<sup>2</sup>; pI 9.90) in terms of their ability to induce soluble expression of rMefp1 fusion proteins. The sequences of the N-termini were Met-(Asp)<sub>5</sub>-(Ala)<sub>2</sub> (pI 2.73); Met-(Asp)<sub>3</sub>-(Ala)<sub>2</sub> (pI 2.87); Met-Asp-Ala (pI 3.00); Met-(Glu<sup>GAA</sup>Glu<sup>GAG</sup>)<sub>4</sub> (sequence name ME<sub>8</sub>-I; pI 2.75); Met-(Glu<sup>GAA</sup>Glu<sup>GAG</sup>)<sub>3</sub> (sequence name ME<sub>6</sub>-I; pI 2.82); Met-(Glu)<sub>2</sub> (pI 3.09); Met-Ala-Glu (pI 3.25); Met-Cys<sub>6</sub> (pI 4.61); Met-(Cys)<sub>3</sub> (pI 4.75); Met-Ala-Cys (pI 4.83); Met-(Ala)<sub>2</sub> (pI 5.60); Met-(Gly)<sub>2</sub> (pI 5.85); Met-Ala-Lys-Asp (pI 6.59); Met-Ala-Lys-Glu (pI 6.79); Met-Cys-His (pI 7.13); Met-Ala-His (pI 7.65); Met-Ala-(His)<sub>3</sub> (pI 7.89); Met-(His)<sub>5</sub> (pI 8.01); Met-Ala-Lys-Cys (pI 8.78); and Met-Lys-Tyr (pI 9.58) (Table 1 and Supplementary Table 1). Soluble expression generally increased as the N-terminal pI increased from 2.73 to 3.25 and was higher with all N-termini with pI values ranging between 4.61 and 9.58 than with the control N-terminus (Fig. 1A, Table 1, and Supplementary Figs. S1A-S1C).

Analysis of the curves showing soluble expression of rMefp1 fusion proteins with (Asp)<sub>*n*</sub>- and (Glu)<sub>*n*</sub>-containing N-termini (pI 2.73-3.25) revealed that their slopes decreased as the value of *n* increased (Fig. 1A and Table 1). Thus, as hydrophilicity increased, the expression of soluble rMefp1 decreased (as was the case with the (Lys)<sub>*n*</sub>- and (Arg)<sub>*n*</sub>-containing N-termini).

The above N-terminal sequences with pI values in the range

4.75 to 9.58 are shorter than the minimum six residues required for calculation of hydrophobicity indices on the Hopp-Woods scale (the hydrophobicity of the Cys-containing N-terminus Met-(Cys)<sub>6</sub> (pI 4.61) is -0.64). They all contain Cys, Ala, Gly, Lys, Asp, Glu, and/or His residues. Only Met-Ala-Lys-Cys and Met-Lys-Tyr are charged (Table 1), so the size of the positive charge cannot be used as a useful index for soluble rMefp1 expression. Therefore, we examined  $\Delta G_{\text{RNA}}$  as a potential determinant of the expression results; mRNA transcripts with strong secondary structures in their 5' regions are translated less efficiently (Mukund et al., 1999; Ramesh et al., 1994). The calculated  $\Delta G_{\text{RNA}}$  values of the N-termini with pI values in the range 4.61-9.58 varied widely, from -2.50 to -7.80 (Table 1 and Supplementary Table 1), but did not seem to correlate with the level of expression of soluble rMefp1 (Fig. 1A and Table 1).

Overall, analysis of N-termini with pI values ranging from 2.73 to 13.35 yielded at least four different curves describing the relationship between pI and expression: expression increased from pI 2.73 to 3.25, remained relatively high from pI 4.61 to 9.58, increased hyperbolically from pI 9.90 to 11.41, and decreased from pI 11.52 to 13.35 (Fig. 1A). These curves clearly show that the pI value of the N-terminus can be used, over a wide range (2.73-13.35), to predict the level of expression of soluble rMefp1. Both a (Lys)<sub>*n*</sub>-specific curve (pI 9.90-11.41) and an (Arg)<sub>*n*</sub>-specific curve (pI 11.52-13.35) fell within the alkaline range (pI 9.90-13.35) (Fig. 1A). Because soluble expression depends on membrane transport mechanisms (as discussed above), we surmise that the number of specific soluble expression curves at alkaline pI values directly relates to the number of membrane channels involved in the transport of rMefp1.

To test the presumptive number of curves in the alkaline pI range (pI 9.90-13.35), we next re-plotted the data as a Hubbert curve (Hubbert, 1956), which encompassed the (Lys)<sub>*n*</sub>-specific (pI 9.90-11.41) and (Arg)<sub>*n*</sub>-specific (pI 11.52-13.35) curves. Removal of the data points relating to Lys-containing N-termini with pI values of 11.28 and 11.41, which yielded especially low expression of soluble rMefp1, produced a smooth curve in the pI range 9.90-13.35 (Fig. 1B). Three soluble rMefp1 expression curves were apparent: an acidic curve from pI 2.73 to 3.25, a neutral curve from pI 4.61 to 9.58, and an alkaline curve from pI 9.90 to 13.35. These observations again suggest the presence of distinct inner membrane channels that transport rMefp1 fusion proteins with acidic, neutral, and alkaline N-termini.

#### Effect of the N-terminal pI value on the expression of insoluble rMefp1

We next assessed the influence of the N-terminal pI on the expression of insoluble rMefp1 fusion protein using the same N-terminal sequences used to assess soluble expression (Fig. 1A, Table 1, and Supplementary Figs. S1A-S1E). The combined amount of soluble and insoluble rMefp1 was considered to be an indirect index of translation (Fig. 1B). With (Asp)<sub>*n*</sub>- and (Glu)<sub>*n*</sub>-containing N-termini (pI 2.73-3.25,  $\Delta G_{\text{RNA}}$  -5.60 to -9.90), the amount of insoluble rMefp1 decreased gradually as *n* increased and pI decreased. With N-termini of pI 4.61-9.58 ( $\Delta G_{\text{RNA}}$  -2.50 to -7.80), more insoluble rMefp1 was produced than MAK-I control (pI 9.90,  $\Delta G_{\text{RNA}}$  -7.80). Insoluble rMefp1 expression induced by the (Lys)<sub>*n*</sub>-containing N-termini ( $\Delta G_{\text{RNA}}$  -7.80 to -1.39) increased between pI 9.90 and 10.82 and then decreased markedly at pI 11.41, while that induced by the (Arg)<sub>*n*</sub>-containing N-termini ( $\Delta G_{\text{RNA}}$  -2.90 to -11.70) decreased gradually, from a peak at pI 11.52 to a much lower level than that of the control at pI 13.35 (Fig. 1A and Table 1). As with the soluble expression curves (described above), the separately generated insoluble expression curves for (Lys)<sub>*n*</sub>- and (Arg)<sub>*n*</sub>-

containing N-termini (pI ranges 9.90-11.41 and 11.52-13.35, respectively) were combined to yield a single total expression curve encompassing the pI range 9.90-13.35 (Fig. 1B).

Although  $\Delta G_{\text{RNA}}$  influences translational efficiency, we believe that it is more accurately reflected by the total amount of target protein expressed than by the amount of target protein expressed in soluble form, which is also determined by the efficiency of its secretion into the periplasm through an N-terminal pI range-specific inner membrane channel. Therefore, we also investigated the influence of N-terminal pI,  $\Delta G_{\text{RNA}}$ , and hydrophilicity on total and soluble expression of rMefp1 fusion proteins with N-termini with pI 2.73-3.25 (Fig. 1B and Table 1). With the N-termini Met-(Asp)<sub>5</sub>-(Ala)<sub>2</sub> (sequence name MD<sub>5</sub>AA; pI 2.73,  $\Delta G_{\text{RNA}}$  -9.90,  $h$  1.09), Met-(Asp)<sub>3</sub>-(Ala)<sub>2</sub> (MD<sub>3</sub>AA; pI 2.87,  $\Delta G_{\text{RNA}}$  -8.80,  $h$  0.56), and Met-Asp-Ala (MDA; pI 3.00,  $\Delta G_{\text{RNA}}$  -6.70), total rMefp1 expression decreased as  $\Delta G_{\text{RNA}}$  and hydrophilicity increased. Moreover, decreased total rMefp1 expression correlated with decreased soluble rMefp1 expression. Soluble rMefp1 expression changed markedly as a factor of both the N-terminal pI value and overall level of protein expression.

With the (Glu)<sub>*n*</sub>-containing N-terminal sequences ME<sub>8</sub>-I (pI 2.75,  $\Delta G_{\text{RNA}}$  -8.30,  $h$  2.08), ME<sub>6</sub>-I (pI 2.82,  $\Delta G_{\text{RNA}}$  -8.50;  $h$  1.82), ME<sub>4</sub> (pI 2.92;  $\Delta G_{\text{RNA}}$  -5.80), ME<sub>2</sub> (pI 3.09,  $\Delta G_{\text{RNA}}$  -5.80), and MAE (pI 3.25,  $\Delta G_{\text{RNA}}$  -5.60), total rMefp1 expression decreased as  $n$  and  $h$  increased (Fig. 1B and Table 1), consistent with the observed decrease in soluble rMefp1 expression. However, the small difference in the  $\Delta G_{\text{RNA}}$  values of ME<sub>8</sub>-I ( $\Delta G_{\text{RNA}}$  -8.30) and ME<sub>6</sub>-I ( $\Delta G_{\text{RNA}}$  -8.50) would not be expected to greatly influence overall rMefp1 expression, which was seemingly determined primarily by the hydrophilicity of the N-terminus. With (Glu)<sub>*n*</sub>-containing N-termini, pI also greatly influences soluble rMefp1 expression, which (with such N-termini) is thus apparently also determined by the overall level of rMefp1 expression and the activity of a pI range-specific inner membrane channel.

Of the N-termini with pI 2.73-3.25, the majority of rMefp1 was expressed in insoluble form with five, but not with the remaining three, which yielded slightly less insoluble protein than soluble protein (Table 1, seq. nos. 1-4, 6, 8-10), suggesting that no active transport (Overton, 1895) of the soluble form into the periplasm occurred. As all rMefp1 fusion proteins larger than the control Met-TxMefp1\* (9.8 kDa) were secreted into the periplasm in soluble form, we suggest that rMefp1 fusion proteins with acidic N-termini are transported via the acid-specific inner membrane channels of *E. coli* by facilitated diffusion.

Based on these results, we postulate that soluble expression of with these N-termini is regulated by the relationship between the N-terminal pI values of all the cytoplasmic rMefp1 fusion proteins and the pI range-specific inner membrane channel. Notably, the pI range-specific inner membrane channel is irreversibly attached to the inner membrane, while the pI values of the N-termini of all of the target proteins are variable, so that the level of soluble expression would be determined by the various pI values of the N-termini of all of the target proteins. Therefore, we introduced the term "pI value trigger" to describe the proportion of total protein expressed in soluble form in the periplasm, as determined by translocational efficiency or the diffusion coefficient for transit through the pI range-specific inner membrane channel, which is in turn controlled by the various pI values of the N-termini of all of the target proteins.

In general, greater total rMefp1 expression was induced by N-termini with pI values in the range 4.61-9.58 ( $\Delta G_{\text{RNA}}$  -2.50 to -7.80) than by the control MAK-I N-terminus. This expression did not correlate closely with  $\Delta G_{\text{RNA}}$  but did, for MAH<sub>5</sub> (pI 8.01,  $\Delta G_{\text{RNA}}$  -6.70;  $h$  -0.33), appear to be determined by the N-

terminal hydrophobicity (MAH<sub>3</sub> (pI 7.89,  $\Delta G_{\text{RNA}}$  -5.40) and MAH (pI 7.65,  $\Delta G_{\text{RNA}}$  -5.20), do not have calculable hydrophobicities) (Table 1). However, it is possible to represent the pI value as an index for the soluble/total rMefp1 ratio in terms of a pI value trigger or diffusion coefficient. This ratio varies substantially: at high pI value triggers (pI 4.83, 5.16, 5.60), it was almost 1:2, while at low pI value triggers (pI 8.01), it was less than 1:4 (Table 1, seq. nos. 11-15, 19-27). However, soluble expression of rMefp1 was higher with all N-termini with pI values in the range 4.61-9.58 than with the control N-terminus (Fig. 1B). The neutral curve to which it relates shows a non-gradual change in slope, in contrast to the other two curves (pI 2.73-3.25 and 9.90-13.35), which have large positive and/or negative slopes, respectively. We surmise that soluble rMefp1 fusion proteins with N-termini pI values in the range 4.61-9.58 pass from the total rMefp1 pool through the neutral inner membrane channel to the periplasm by facilitated diffusion (as suggested above for fusion proteins with N-termini of pI 2.73-3.25).

In examining total rMefp1 expression induced by the (Lys)<sub>*n*</sub>-containing N-termini (pI 9.90-11.41,  $\Delta G_{\text{RNA}}$  -7.80 to -1.39), we found that total expression was lower with MAK-I (pI 9.90,  $\Delta G_{\text{RNA}}$  -7.80) than with MK<sub>2</sub>AK-I (pI 10.82,  $\Delta G_{\text{RNA}}$  -2.10), consistent with the favorable change in  $\Delta G_{\text{RNA}}$ . However, total expression induced by MK<sub>3</sub>AK (pI 10.99,  $\Delta G_{\text{RNA}}$  -2.00,  $h$  1.14) was much higher than that induced by MK<sub>8</sub>AK (pI 11.41,  $\Delta G_{\text{RNA}}$  -1.39,  $h$  1.93), in spite of the slightly more favorable  $\Delta G_{\text{RNA}}$  value of the latter (Fig. 1A and Table 1). With longer (Lys)<sub>*n*</sub>-containing N-termini, decreases in total rMefp1 expression were not consistent with the reductions in  $\Delta G_{\text{RNA}}$ , but were instead related to increases in hydrophilicity. Here, to decrease  $\Delta G_{\text{RNA}}$ , we used the Lys<sup>AAA</sup> codon (Supplementary Table 1). As the number of these Lys<sup>AAA</sup> codons was increased, total rMefp1 expression rose and then gradually fell (without translational arrest) (Fig. 1A and Table 1), indicating that the use of the Lys<sup>AAA</sup> codon prevented translational inhibition. The soluble and insoluble rMefp1 expression curves for these (Lys)<sub>*n*</sub>-containing N-termini show the typical hyperbolic form (Fig. 1A), leading us to speculate that a feedback mechanism regulates total and soluble rMefp1 protein levels, and that the trigger for this feedback is the increased hydrophilicity of the N-terminus.

Among the (Lys)<sub>*n*</sub>-containing N-termini, total rMefp1 expression was highest with MK<sub>2</sub>AK-I (pI 10.82,  $\Delta G_{\text{RNA}}$  -2.10), but soluble rMefp1 expression was higher with MK<sub>3</sub>AK (pI 10.99,  $\Delta G_{\text{RNA}}$  -2.00), MK<sub>4</sub>AK (pI 11.11,  $\Delta G_{\text{RNA}}$  -1.80), and MK<sub>5</sub>AK (pI 11.21,  $\Delta G_{\text{RNA}}$  -1.70) than with MK<sub>2</sub>AK-I, despite their slightly reduced production of total rMefp1 (Fig. 1B and Table 1). This observation suggests that total protein synthesis and soluble expression are controlled independently. The regulation of total protein synthesis appears to be complicated, and the inner membrane protein very sensitive to the N-terminal pI value.

With (Lys)<sub>*n*</sub>-containing N-termini of pI 9.90-11.41, the most rMefp1 was expressed in insoluble form (with the exception of one N-terminus, which yielded slightly less insoluble protein than soluble protein) (Table 1, seq. nos. 29, 33, 34, 36-40), suggesting that fusion proteins with these termini were translocated through the alkaline pI-specific inner membrane channel by facilitated diffusion (similar to the (Asp)<sub>*n*</sub>- and (Glu)<sub>*n*</sub>-containing N-termini).

Among the (Arg)<sub>*n*</sub>-containing N-termini with pI values in the range 11.52-13.35, highest total expression of rMefp1 was obtained with MRAK-I (pI 11.52,  $\Delta G_{\text{RNA}}$  -2.90), which uses the Arg<sup>AGA</sup> codon. Total expression gradually decreased as the number of Arg-encoding codons was increased: MR<sub>2</sub>AK (R<sub>2</sub> encoded by Arg<sup>CGT</sup> Arg<sup>CGC</sup>, pI 12.51,  $\Delta G_{\text{RNA}}$  -5.10) < MR<sub>4</sub>AK (R<sub>4</sub> = (Arg<sup>CGT</sup> Arg<sup>CGC</sup>)<sub>2</sub>, pI 12.98,  $\Delta G_{\text{RNA}}$  -7.70,  $h$  1.32) < MR<sub>6</sub>AK (R<sub>6</sub>

= (Arg<sup>CGT</sup>Arg<sup>CGC</sup>)<sub>3</sub>, pI 13.20,  $\Delta G_{\text{RNA}}$  -10.20,  $h$  1.69) < MR<sub>8</sub>AK ( $R_8$  = (Arg<sup>CGT</sup>Arg<sup>CGC</sup>)<sub>4</sub>, pI 13.35,  $\Delta G_{\text{RNA}}$  -11.70,  $h$  1.93) (Fig. 1B, Table 1, and Supplementary Table 1). With these N-termini, rMefp1 expression decreased as  $\Delta G_{\text{RNA}}$  and  $h$  increased. While effective in MR<sub>2</sub>AK-I, the Arg<sup>AGA</sup> codon could not be used in MR<sub>2</sub>AK due to a reduction in translational efficiency (Supplementary Fig. S2), or in MR<sub>4</sub>AK, MR<sub>6</sub>AK, or MR<sub>8</sub>AK, due to translational arrest (Supplementary Fig. S2). These results show that repetitive use of a particular codon can slow or arrest translation, regardless of the value of  $\Delta G_{\text{RNA}}$ .

As with the (Lys)<sub>*n*</sub>-containing N-termini, the (Arg)<sub>*n*</sub>-containing N-termini yielded higher insoluble expression of rMefp1 fusion proteins than soluble expression (with the exception of one N-terminus, which yielded slightly less insoluble protein than soluble protein) (Table 1, seq. nos. 41, 43-46). These differences in the expression of soluble rMefp1 appeared to be caused by different pI triggers for inner membrane secretion of the total rMefp1 (the soluble expression of these fusion rMefp1 occurring as a result of facilitated diffusion through the alkaline pI-specific inner membrane channel). The curves displaying soluble and insoluble expression of these rMefp1 fusion proteins showed similar trends (Fig. 1A).

As discussed above, greater amounts of soluble rMefp1 were expressed with the (Lys)<sub>*n*</sub>-containing N-termini with hydrophilicities of 1.32 (MK<sub>4</sub>AK; pI 11.11,  $\Delta G_{\text{RNA}}$  -1.80, pI value trigger 0.40), 1.53 (MK<sub>5</sub>AK; pI 11.21,  $\Delta G_{\text{RNA}}$  -1.70, pI value trigger 0.45), and 1.69 (MK<sub>6</sub>AK; pI 11.28,  $\Delta G_{\text{RNA}}$  -1.56, pI value trigger 0.45) than with the (Arg)<sub>*n*</sub>-containing N-terminus with a hydrophilicity of 1.32 (MR<sub>4</sub>AK; pI 12.98,  $\Delta G_{\text{RNA}}$  -7.70, pI value trigger 0.48) (Table 1). This phenomenon results from the complicated nature of the influence of  $\Delta G_{\text{RNA}}$  and  $h$  on total protein expression and of the pI value trigger on soluble expression. The higher soluble expression with (Lys)<sub>*n*</sub>-containing sequences than with (Arg)<sub>*n*</sub>-containing sequences may be better explained by  $\Delta G_{\text{RNA}}$  than by hydrophilicity or pI trigger values.

The least effective inducers of soluble rMefp1 expression were the N-terminal sequences MD<sub>5</sub>AA (pI 2.73,  $\Delta G_{\text{RNA}}$  -9.90,  $h$  1.09), MK<sub>8</sub>AK (pI 11.41,  $\Delta G_{\text{RNA}}$  -1.39,  $h$  1.93), and MR<sub>8</sub>AK (pI 13.35,  $\Delta G_{\text{RNA}}$  -11.70,  $h$  1.93). However, the fusion proteins containing these N-termini were all expressed in the periplasm, as described in "Materials and Methods" (data not shown). Both soluble and insoluble expression levels were observed (Fig. 1A, Table 1, and Supplementary Figs. S1A, S1D, and S1E), indicating that a proportion of the rMefp1 produced was secreted into the periplasm in soluble form through the acidic or alkaline pI-specific inner membrane channels. The remaining rMefp1 protein in the cytoplasm was collected in insoluble form, as described in "Materials and Methods". Therefore, it seems that a mechanism exists that converts even the smallest amounts of unsecreted protein into an insoluble form in the cytoplasm.

We consider combining the insoluble expression curves for the (Lys)<sub>*n*</sub>- and (Arg)<sub>*n*</sub>-containing N-termini (pI 9.90-13.35) into a single curve to be logical. Removing the data points relating to N-termini with pI values of 11.28 and 11.41, responsible for a sharp but transitory decrease in insoluble expression, yielded a relatively smooth total expression hyperbolic curve (Fig. 1B). The alkaline curves for total and soluble rMefp1 expression (Fig. 1B) showed similar trends.

We have suggested the existence of a hydrophilicity-based feedback regulatory mechanism controlling the expression of soluble and insoluble rMefp1 fused to (Lys)<sub>*n*</sub>-containing N-termini (pI 9.90-11.41). We suggest also that the similar patterns of soluble and insoluble expression obtained with the (Asp)<sub>*n*</sub>- and (Glu)<sub>*n*</sub>-containing N-termini (pI 2.73-3.25) and from the (Arg)<sub>*n*</sub>-containing N-termini (pI 11.52-13.35) (Fig. 1A) are

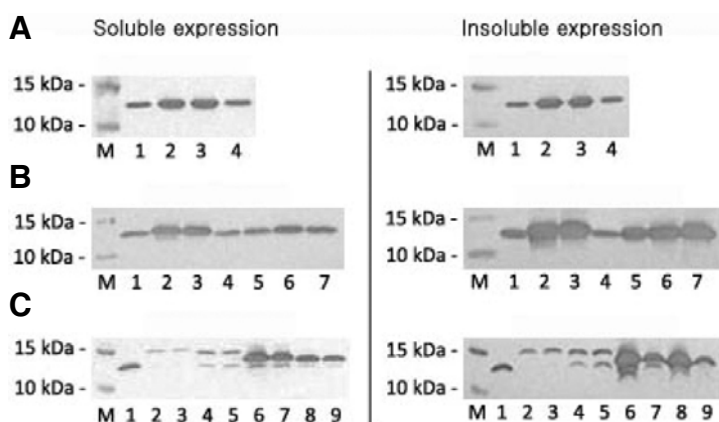
the result of hydrophilicity-based feedback regulation. However, the same is not true of the rMefp1 fusion proteins with (Asp)<sub>*n*</sub>- or (Arg)<sub>*n*</sub>-containing N-termini because their expression is also influenced by  $\Delta G_{\text{RNA}}$ . Designing longer (Asp)<sub>*n*</sub>- or (Arg)<sub>*n*</sub>-containing N-termini with lower  $\Delta G_{\text{RNA}}$  values than those of shorter N-termini (to provide a broader range of  $\Delta G_{\text{RNA}}$  values for closer examination of this issue) was problematic. Instead, we decided to examine the relationship between  $\Delta G_{\text{RNA}}$  and hydrophilicity in feedback regulation of soluble and insoluble rMefp1 expression using two pairs of (Glu)<sub>*n*</sub>-containing N-termini with differing codon usages. These were ME<sub>8</sub>-I (pI 2.75,  $\Delta G_{\text{RNA}}$  -8.30) and ME<sub>8</sub>-II (pI 2.75,  $\Delta G_{\text{RNA}}$  -4.00), which both have a hydrophilicity of 2.08; and ME<sub>6</sub>-I (pI 2.82,  $\Delta G_{\text{RNA}}$  -8.50) and ME<sub>6</sub>-II (pI 2.82,  $\Delta G_{\text{RNA}}$  -4.30), which both have a hydrophilicity of 1.82 (Table 1 and Supplementary Table 1). Analysis of these N-termini showed a modest dependence on  $\Delta G_{\text{RNA}}$ , but a much greater dependence on hydrophilicity (Fig. 2C). Overall, the results from the longer N-termini with pI values in the ranges 2.73-3.25, 9.90-11.41, and 11.52-13.35 show that hydrophilicity is much more important than  $\Delta G_{\text{RNA}}$  in inhibitory feedback regulation of total rMefp1 synthesis and soluble rMefp1 expression. Therefore, this feedback regulation seems to depend primarily on the N-terminal hydrophilicity.

The hydrophilicity of the acid (pI 2.73-3.25) and alkaline (9.901-13.35) N-termini appears to have a dual function: it influences 1) protein synthesis, thus controlling overall rMefp1 expression, and 2) the pI value trigger, thus controlling soluble rMefp1 expression. Therefore, we believe that the reduction in total rMefp1 expression linked to increased N-terminal hydrophilicity results from decreased translational efficiency. This presumptive regulation of the translational efficiency and total and soluble protein expression by N-terminal hydrophilicity suggests that it is important for feedback regulation in *E. coli*. Differences between the hydrophilic amino acids present in different N-termini (Asp, Glu, Lys, and Arg) suggest that hydrophilicity is a more important factor for feedback regulation than is the presence of a specific recognition site.

We also investigated the influence of charge on total rMefp1 expression. Analysis of the (Asp)<sub>*n*</sub>-, (Glu)<sub>*n*</sub>-, (Arg)<sub>*n*</sub>-, and (Lys)<sub>*n*</sub>-containing N-termini showed that, as the positive or negative charge increased, total rMefp1 expression level decreased (Fig. 1A and Table 1). This effect correlated more closely with hydrophilicity than with the value of  $n$ . However, pI values in the ranges 2.73-3.25 and 9.90-13.35 can be used as an index to predict levels of total and soluble rMefp1 expression, regardless of hydrophilicity.

Our analysis of the relationship between N-terminal hydrophilicity and overall rMefp1 expression revealed both feedback and non-feedback regulatory control of soluble rMefp1 expression. In non-feedback regulation, which we observed with N-termini of pI 4.61-9.58 and no calculated hydrophilicity,  $\Delta G_{\text{RNA}}$  had little effect on total rMefp1 expression, while the rate at which the cytoplasmic rMefp1 was secreted into the periplasm in soluble form (by facilitated diffusion through inner membrane channels specific for N-termini with neutral pI values) was controlled by the pI value of the N-terminus.

In feedback regulation, which we observed with the N-termini with values in the ranges 2.73-3.25 and 9.90-13.35, total rMefp1 expression was controlled by  $\Delta G_{\text{RNA}}$  when the N-termini were short and had no calculated hydrophilicity and by hydrophilicity when the N-termini contained multiple hydrophilic amino acids. For both short and long N-termini, the rate at which cytoplasmic rMefp1 was secreted into the periplasm in soluble form was controlled by the specific pI value trigger for facilitated diffusion through inner membrane channels specific



**Fig. 2.** Comparative Western blot analysis of soluble and insoluble rMefp1 expression directed by N-termini with identical pI values. Western blots analysis was performed as described in Fig. 1. (A) Comparison of four Met-Ala-Lys N-termini with different Ala codons. All have a pI value of 9.90 and a  $\Delta G_{\text{RNA}}$  value of -7.80. MAK-I uses Ala<sup>GCT</sup>, MAK-II uses Ala<sup>GCA</sup>, MAK-III uses Ala<sup>GCC</sup>, and MAK-IV uses Ala<sup>GCG</sup>. Lanes: M, molecular markers; 1, MAK-I; 2, MAK-II; 3, MAK-III; 4, MAK-IV. (B) N-termini with different amino acid compositions. Lanes: M, molecular markers; 1, MAK-I (pI 9.90); 2, MNN (pI 5.70); 3, MTT (pI 5.70); 4, MWW (pI 5.85); 5, MGG (pI 5.85); 6, MAKY (pI 9.58); 7, MKY (pI 9.58). (C) N-termini with different  $\Delta G_{\text{RNA}}$  values with or without calculated hydrophobicity indices. Lanes: M, molecular marker; 1, MAK-I (pI 9.90; Ala<sup>GCT</sup>,  $\Delta G_{\text{RNA}}$  -7.80); 2, ME<sub>6</sub>-I (pI 2.75,  $\Delta G_{\text{RNA}}$  -8.30); 3, ME<sub>6</sub>-II (pI 2.75,  $\Delta G_{\text{RNA}}$  -4.00); 4, ME<sub>6</sub>-I (pI 2.82,  $\Delta G_{\text{RNA}}$  -8.50); 5, ME<sub>6</sub>-II (pI 2.82,  $\Delta G_{\text{RNA}}$  -4.30); 6, MK<sub>2</sub>AK-I (pI 10.82,  $\Delta G_{\text{RNA}}$  -2.10); 7, MK<sub>2</sub>AK-II (pI 10.82,  $\Delta G_{\text{RNA}}$  -4.70); 8, MRAK-I (pI 11.52,  $\Delta G_{\text{RNA}}$  -2.90); 9, MRAK-II (pI 11.52,  $\Delta G_{\text{RNA}}$  -7.80).

for N-termini with acidic or alkaline pI values.

N-terminal pI values affect soluble rMefp1 expression levels, even when increased hydrophilicity reduces overall protein synthesis through feedback regulation. Therefore, as an index for soluble rMefp1 expression over a wide pI range, N-terminal pI is a more important than hydrophilicity, charge, or  $\Delta G_{\text{RNA}}$ , irrespective of whether soluble rMefp1 expression is controlled by pI-based feedback, as shown in Figs. 1A, 1B, and Table 1.

#### Effects of codon selection and duplication, and amino acid sequence, on expression of soluble and insoluble rMefp1 fusion proteins with N-termini with identical pI values

We showed that with N-termini of pI 2.73-3.25 and 9.90-13.35, the N-terminal  $\Delta G_{\text{RNA}}$  and hydrophilicity influence overall rMefp1 expression, and that for all N-termini of pI 2.73-13.35, the pI value controls soluble rMefp1 expression via the pI value trigger. Therefore, we hypothesized that in N-termini with identical pI values, differences in codon usage or amino acid sequence might influence expression of soluble and insoluble rMefp1. We assessed the expression of soluble and insoluble rMefp1 induced by three sets of N-termini with identical pI values. These sets possessed (i) identical amino acid sequences and  $\Delta G_{\text{RNA}}$  values but different codon selections (Met-Ala<sup>GCT</sup>-Lys (MAK-I), Met-Ala<sup>GCA</sup>-Lys (MAK-II), Met-Ala<sup>GCC</sup>-Lys (MAK-III), and Met-Ala<sup>GCG</sup>-Lys (MAK-IV) (all pI 9.90,  $\Delta G_{\text{RNA}}$  -7.80)); (ii) identical  $\Delta G_{\text{RNA}}$  values but different amino acid sequences (MNN and MTT (both pI 5.70); MWW and MGG (both pI 5.85); and MKY and MAKY (both pI 9.58)); or (iii) identical amino acid sequences but different  $\Delta G_{\text{RNA}}$  values (Met-(Glu<sup>GAA</sup>Glu<sup>GAG</sup>)<sub>4</sub> (ME<sub>6</sub>-I;  $\Delta G_{\text{RNA}}$  -8.30) and Met-(Glu<sup>GAA</sup>)<sub>7</sub>Glu<sup>GAG</sup> (ME<sub>6</sub>-II;  $\Delta G_{\text{RNA}}$  -4.00); Met-(Glu<sup>GAA</sup>Glu<sup>GAG</sup>)<sub>3</sub> (ME<sub>6</sub>-I;  $\Delta G_{\text{RNA}}$  -8.50) and Met-(Glu<sup>GAA</sup>)<sub>5</sub>Glu<sup>GAG</sup> (ME<sub>6</sub>-II;  $\Delta G_{\text{RNA}}$  -4.30); Met-(Lys<sup>AAA</sup>)<sub>2</sub>-Ala-Lys (MK<sub>2</sub>AK-I;  $\Delta G_{\text{RNA}}$  -2.10) and Met-Lys<sup>AAA</sup>-Lys<sup>AAG</sup>-Ala-Lys (MK<sub>2</sub>AK-II;  $\Delta G_{\text{RNA}}$  -4.70); and Met-Arg<sup>AGA</sup>-Ala-Lys (MRAK-I;  $\Delta G_{\text{RNA}}$  -2.90) and Met-Arg<sup>GCT</sup>-Ala-Lys (MRAK-II;  $\Delta G_{\text{RNA}}$  -5.00)) (Table 1 and Supplementary Table 1).

Levels of soluble and insoluble rMefp1 obtained using MAK-I (using Ala<sup>GCT</sup>) were similar to those for MAK-IV (using Ala<sup>GCG</sup>) but very different from those for MAK-II (Ala<sup>GCA</sup>) and MAK-III (Ala<sup>GCC</sup>), in spite of their identical  $\Delta G_{\text{RNA}}$  values (Fig. 2A). MNN and MTT, which have identical pI values, similar molecular weights (377.41 and 351.41, respectively) and slightly different  $\Delta G_{\text{RNA}}$  values (-2.50 and -3.10, respectively), were quite similar in terms of soluble and insoluble rMefp1 expression (Fig. 2B). MWW and MGG, which have identical pI values, very different

molecular weights (521.61 and 263.31, respectively) and similar  $\Delta G_{\text{RNA}}$  values (-7.50 and -7.80), differed in terms of insoluble rMefp1 expression (the smaller MGG yielded higher expression) but not soluble rMefp1 expression (Fig. 2B). MKY and MAKY, which have identical pI values and similar amino acid sequences and molecular weights (440.55 and 511.62, respectively), but different  $\Delta G_{\text{RNA}}$  values (-2.50 and -5.20, respectively), yielded similar amounts of soluble and insoluble rMefp1 (Fig. 2B).

The longer ME<sub>6</sub>-I/II and ME<sub>6</sub>-I/II pairs, which have identical pI values and hydrophilicity, but different codon selections and  $\Delta G_{\text{RNA}}$  values (ME<sub>6</sub>-I and -II:  $\Delta G_{\text{RNA}}$  -8.30 and -4.00, respectively; ME<sub>6</sub>-I and -II:  $\Delta G_{\text{RNA}}$  -8.50 and -4.30, respectively), yielded similar levels of soluble and insoluble rMefp1 (Fig. 2C). The shorter MK<sub>2</sub>AK-I/II and MRAK-I/II pairs, which also have identical pI values, but different codon selections and  $\Delta G_{\text{RNA}}$  values and no calculated hydrophobicity indices (MK<sub>2</sub>AK-I and -II:  $\Delta G_{\text{RNA}}$  -2.10 and -4.70, respectively; MRAK-I and -II:  $\Delta G_{\text{RNA}}$  -2.90 and -5.00, respectively), yielded very different levels of insoluble rMefp1 according to their  $\Delta G_{\text{RNA}}$  values, but similar amounts of soluble rMefp1 (Fig. 2C). Of course, higher overall rMefp1 expression might be expected to translate into increased levels of soluble rMefp1.

In the case of the MAK N-termini, all of which have the same  $\Delta G_{\text{RNA}}$ , single-base differences at the third "wobble" position of the Ala-encoding codons resulted in large changes in the expression of soluble and insoluble rMefp1. The Ala<sup>GCC</sup> codon yielded higher soluble and insoluble rMefp1 expression than Ala<sup>GCT</sup>, even though both codons use the same anticodon (CGG). Similarly, Ala<sup>GCA</sup> yielded more soluble and insoluble rMefp1 than Ala<sup>GCG</sup>, even though both codons use a CGU anticodon. Thus, MAK-I and MAK-IV (which contain Ala<sup>GCT</sup> and Ala<sup>GCG</sup>, respectively), which use different anticodons (CGG and CGU, respectively), yielded similar amounts of soluble and insoluble rMefp1. This observation suggests that a single base change can alter total rMefp1 expression by affecting the efficiency of translation.

Small differences in N-terminal molecular weight between MNN and MTT, and MKY and MAKY, had little effect on soluble and insoluble rMefp1 expression. The large difference in the molecular weights of the MWW and MGG N-termini had a significant effect on the amount of insoluble protein produced, but not on soluble rMefp1 expression.

To examine the effects of N-terminal  $\Delta G_{\text{RNA}}$  and hydrophilicity on soluble and insoluble rMefp1 expression, we compared N-termini with identical pI values but different  $\Delta G_{\text{RNA}}$  values and/or



hydrophilicity. To reduce the  $\Delta G_{\text{RNA}}$  values of the ME<sub>8</sub>-I (Met-(Glu<sup>GAA</sup>Glu<sup>GAG</sup>)<sub>4</sub>) and ME<sub>6</sub>-I (Met-(Glu<sup>GAA</sup>Glu<sup>GAG</sup>)<sub>3</sub>) N-termini (which have  $\Delta G_{\text{RNA}}$  values of -8.30 and -8.50, respectively), we altered Glu codon selection and thus created ME<sub>8</sub>-II (Met-(Glu<sup>GAA</sup>)<sub>7</sub>Glu<sup>GAG</sup>) and ME<sub>6</sub>-II (Met-(Glu<sup>GAA</sup>)<sub>5</sub>Glu<sup>GAG</sup>), which have  $\Delta G_{\text{RNA}}$  values of -4.00 and -4.30, respectively (Table 1 and Supplementary Table 1). However, these changes in  $\Delta G_{\text{RNA}}$  did not affect the amount of soluble and insoluble rMefp1 produced (Fig. 2C) and there was no translational arrest or reduction in translational efficiency such as we observed with the (Arg)<sub>n</sub>-containing N-termini with repetitive Arg codon usage. Thus, altering the  $\Delta G_{\text{RNA}}$  values of longer N-termini did not affect soluble or insoluble rMefp1 expression. Increased hydrophilicity, meanwhile, greatly reduces the production of soluble and insoluble rMefp1 fusion protein.

By contrast, with shorter N-termini lacking calculated hydrophobicity indices, a change in  $\Delta G_{\text{RNA}}$  had a substantial effect on soluble and insoluble rMefp1 expression. (Compare MK<sub>2</sub>AK-I and -II (pI 10.82;  $\Delta G_{\text{RNA}}$  -2.10 and -4.70, respectively) and MRAK-I and -II (pI 11.52;  $\Delta G_{\text{RNA}}$  -2.90 and -5.00, respectively) in Fig. 2C.) Therefore,  $\Delta G_{\text{RNA}}$  appeared to control the synthesis of rMefp1 fusion protein with shorter N-termini. The reduction in total rMefp1 expression seen with longer N-termini was caused by their increased hydrophilicity.

These attempts to increase soluble expression of a target protein lacking a TM-like domain showed that N-termini with small  $\Delta G_{\text{RNA}}$  values, low hydrophilicity, and codon usage that prevents translational arrest or a reduction in translational efficiency are needed to maximize the rate of translation, and thus increase total protein synthesis. Furthermore, the amount of soluble product obtained increases as the amount of total protein increases, and when the N-terminus has a pI value that can trigger translocation through inner membrane channels.

## CONCLUSION

In this study, we investigated the effect of the pI (in the range 2.73-13.35) of the N-terminal region of 7xMefp1, a protein lacking a TM-like domain (Lee et al., 2008b), on soluble and insoluble expression. We also analyzed the relationships between N-terminal pI value and hydrophilicity, charge,  $\Delta G_{\text{RNA}}$  value, molecular weight, amino acid sequence, and codon selection and repetition. We found that the N-terminal pI value can be used as a comprehensive biological index to represent the level of soluble rMefp1 expression over a wide range because it takes account of all of the factors affecting overall and soluble expression.

In plotting soluble expression against N-terminal pI, we identified three curves, one for each pI range studied (acidic, neutral, and alkaline). We surmise that each such curve is derived from a different periplasmic secretion pathway involving an inner membrane channel that is specific range of N-terminal pI values (acidic, neutral, or alkaline). Each pI range-specific soluble expression curve has a clear form and boundary and does not interfere with the other curves (Fig. 1B). N-terminal pI values may influence membrane channel permeability by guiding appropriate channel selection in a wide pI range, and by modulating the kinetics of channel translocation in the narrow pI ranges that allow soluble expression. Therefore, we suggest that there

exist in the *E. coli* inner membrane three types of membrane channels, each specific for target protein N-termini whose pI values fall in certain ranges (acidic, neutral, and alkaline). Collectively, they may constitute a useful index for defining the periplasmic secretion pathways that underpin soluble expression.

*Note: Supplementary information is available on the Molecules and Cells website (www.molcells.org).*

## ACKNOWLEDGMENT

This work was supported by a grant from the National Fisheries Research and Development Institute (RD-10-BT-007).

## REFERENCES

- Bradford, M. (1976). A rapid and sensitive method for the quantitation of microgram quantities of protein utilizing the principle of protein-dye binding. *Anal. Biochem.* 72, 248-254.
- Hubbert, M.K. (1956). Nuclear Energy and the Fossil Fuels. Drilling and Production Practice. (American Petroleum Institute & Shell Development Co. Publication No. 95), pp. 9-11 and 21-22.
- Ikemura, T. (1981). Correlation between the abundance of *Escherichia coli* transfer RNAs and the occurrence of the respective codons in its protein genes: a proposal for a synonymous codon choice that is optimal for the *E. coli* translational system. *J. Mol. Biol.* 151, 389-409.
- Inouye, S., Soberon, X., Franceschini, T., Nakamura, K., Itakura, K., and Inouye, M. (1982). Role of positive charge on the amino-terminal region of the signal peptide in protein secretion across the membrane. *Proc. Natl. Acad. Sci. USA* 79, 3438-3441.
- Laemmli, U.K. (1970). Cleavage of structural proteins during the assembly of the head of bacteriophage T4. *Nature* 227, 680-685.
- Lee, S.J., Han, Y.H., Nam, B.H., Kim, Y.O., and Reeves, P.R. (2008a). A novel expression system for recombinant marine mussel adhesive protein Mefp1 using a truncated OmpA signal peptide. *Mol. Cells* 26, 34-40.
- Lee, S.J., Park, I.S., Han, Y.H., Kim, Y.O., and Reeves, P.R. (2008b). Soluble expression of recombinant olive flounder hepcidin I with a novel secretion enhancer. *Mol. Cells* 26, 140-145.
- Mukund, M.A., Bannerjee, T., Ghosh, I., and Datta, S. (1999). Effect of mRNA secondary structure in the regulation of gene expression: unfolding of stable loop causes the expression of Taq polymerase in *E. coli*. *Curr. Sci.* 76, 1486-1490.
- Nossal, N.G., and Heppel, L.A. (1966). The release of enzymes by osmotic shock from *Escherichia coli* in exponential phase. *J. Biol. Chem.* 241, 3055-3062.
- Overton, E. (1895). Über die osmotischen Eigenschaften der lebenden Pflanzen und Tierzelle. *Vierteljahreschr. Naturforsch. Ges. Zürich* 40, 159.
- Ramesh, V., De, A., and Nagaraja, V. (1994). Engineering hyper-expression of bacteriophage  $\mu$ C protein by removal of secondary structure at the translation initiation region. *Protein Eng.* 7, 1053-1057.
- Sambrook, J., Fritsch, E.F., and Maniatis, T. (1989). Molecular Cloning: A Laboratory Manual, 2nd eds. (Cold Spring Harbor, NY, Cold Spring Harbor Laboratory Press).
- Sanger, F., Nicklen, S., and Coulson, A.R. (1977). DNA sequencing with chain-terminating inhibitors. *Proc. Natl. Acad. Sci. USA* 74, 5463-5467.
- Waite, J.H. (1983). Evidence for a repeating 3,4-dihydroxyphenylalanine- and hydroxyproline-containing decapeptide in the adhesive protein of the mussel, *Mytilus edulis* L. *J. Biol. Chem.* 258, 2911-2915.
- Zuker, M. (2003). Mfold web server for nucleic acid folding and hybridization prediction. *Nucleic Acids Res.* 31, 3406-3415.