

---

**The nucleotide and deduced amino acid sequences of the encephalomyocarditis viral polyprotein coding region**

---

Ann C. Palmenberg\*, Ellen M. Kirby, Michael R. Janda, Neil L. Drake, Gregory M. Duke, Kimberly F. Potratz<sup>+</sup> and Marc S. Collett<sup>+</sup>

---

Biophysics Laboratory of the Graduate School, University of Wisconsin, Madison, WI 53706, and  
<sup>+</sup> Molecular Genetics, Inc., 10320 Bren Road East, Minnetonka, MN 55343, USA

---

Received 19 December 1983; Revised and Accepted 27 February 1984

---

**ABSTRACT**

The nucleotide sequence of 7200 bases of encephalomyocarditis (EMC) viral RNA, including the complete polyprotein-coding region, was determined. The polyprotein is encoded within a unique translational reading frame, 6870 bases in length. Protein synthesis begins with the sequence Met-Ala-Thr, and ends with the sequence Leu-Phe-Trp, 126 bases from the 3' end of the RNA. Viral capsid and noncapsid proteins were aligned with the deduced amino acid sequence of the polyprotein. The proteolytic processing map follows the standard 4-3-4 picornaviral pattern except for a short leader peptide (8 kd), which precedes the capsid proteins. Identification of the proteolytic cleavage sites showed that EMC viral protease, p22, has cleavage specificity for gln-gly or gln-ser sequences with adjacent proline residues. The cleavage specificity of the host-coded protease(s) includes both tyr-pro and gln-gly sequences.

**INTRODUCTION**

The family of picornaviruses contains four major subdivisions: the entero- (polio, coxsackie), cardio- (encephalomyocarditis (EMC), Mengo), rhino- (120 serotypes) and aphthoviruses (foot-and-mouth disease). Although differing greatly in pathogenicity and epidemiology, all picornaviruses share structural and morphogenic properties. Virus particles contain a single-stranded, positive-sense RNA genome ( $2.7 \times 10^6$  MW) enclosed in a protein capsid shell. The capsids are composed of sixty subunits, each of which contains four non-identical polypeptide chains (1).

The 3' ends of viral RNA have poly(A), as is characteristic of most eukaryotic mRNAs (2). However, the 5' ends are not capped in the usual manner with 5'-5' triphosphate linkages. Instead, these viruses have a small, viral-coded, genome-linked protein (VPg) attached by a tyrosine- $O^4$ -phosphodiester bond to the 5' pUp of the RNA (3-5). The cardio- and aphthoviral genomes are distinguished by the presence of a 5'-proximal poly(C) tract, whose length (50-150 bases) and exact location relative to the 5' end

of the RNA (150-500 bases) vary with different serotypes and isolates of virus (6-8).

Translation of picornaviral RNA begins about 700 bases from the 5' end, and produces a large precursor polypeptide which represents most of the coding capacity of the genome (9). The polyprotein is processed in a series of proteolytic cleavage steps to yield mature virion capsid proteins, as well as other viral proteins of a non-structural nature. Although protein molecular weights vary somewhat from virus to virus, the processing scheme shows remarkable similarity among the four subdivisions (1,10).

The first cleavage event occurs while the polyprotein is still nascent on a ribosome and is thought to be catalyzed by a cellular protease, whose identity and specificity remain unclear (11-13). Most subsequent maturation processing is effected by a viral-encoded protease, capable of self-cleaving and autocatalytic reactions within the polyprotein (14-16).

The nucleotide sequences of several different serotypes of polio and portions of foot-and-mouth disease virus (FMDV) have now been completed (17-22). These studies, combined with tryptic peptide mapping and protein end-group determinations, have reconfirmed the major elements of picornaviral homology between aphtho- and enteroviruses. However, the sequences have also highlighted important subdivision differences. For example, translation of FMDV polyprotein begins at two initiation sites, 5' to the capsid region, creating leader peptides attached to the capsid precursor protein (23). Polio translation begins at the start of the capsid proteins, without an attached leader sequence (24). FMDV also contains three tandemly-linked, different protein sequences for VPg, as opposed to the single VPg observed in polio (17,18,25).

There are also differences in the cleavage requirements of the viral proteases. The polio-encoded enzyme appears to have an exclusive specificity for gln-gly peptide linkages. The equivalent FMDV cleavages occur within a broader range of protein sequences, including glu-ser, glu-gly and gln-thr (21,25).

Characterization of picornaviral subdivision differences may help identify those elements responsible for variations in host range, cell specificity, pathogenicity and particle stability. Accordingly, we began to analyze the genome of EMC, a cardiovirus, as example and prototype of this group.

We now present the nucleotide and deduced amino acid sequences of the polyprotein coding region of EMC. Our data show that cardioviruses have structural and biochemical elements which appear to be intermediate between

the entero- and aphthoviruses. EMC polyprotein has a 5' leader sequence like FMDV, but the remainder of the processing map more closely resembles polio. The EMC viral protease has an apparent cleavage specificity for gln-gly and gln-ser sequences which are flanked by proline residues.

## MATERIALS AND METHODS

### Materials

EMC virus and purified viral RNA were obtained from Mark Pallansch and Roland Rueckert (26,27). *E. coli* K12 (strain JM101) and phage M13mp9 were supplied by J. Messing. T4 DNA ligase, DNase I, restriction endonucleases and site-specific primers for di-deoxy M13 sequencing were purchased from New England Biolabs; reverse transcriptase from Life Sciences; large fragment of *E. coli* DNA polymerase I (Klenow) from Promega-Biotec;  $\alpha$  and  $\gamma$   $^{32}\text{P}$ -ATP from Amersham International; oligo(dT) primer and terminal deoxynucleotidyl transferase from P-L Biochemicals; preppacked, mini-C18 columns (Sep-Pac) were from Waters. DNA primer synthesis kits, associated reagents and protocols were obtained from New England Biolabs. *E. coli* strain MC1000 (28) and plasmid pKMG DNA were prepared and provided by Molecular Genetics Inc. pKMG DNA is a linear 3613 base pair derivative of vector pBR322, which has the EcoR I site at one end of the molecule, and a poly(dG)-tailed PstI site at the other.

### Molecular Cloning of EMC

Viral RNA (20  $\mu\text{g}$ ) was reverse transcribed in a 200  $\mu\text{l}$  reaction mixture containing 50 mM Tris-HCl, pH 8.0, 8 mM  $\text{MgCl}_2$ , 75 mM 2-mercaptoethanol, 28 mM KCl, 400  $\mu\text{M}$  each of dATP, dGTP, dCTP and dTTP, 40  $\mu\text{g}/\text{ml}$  oligo(dT)<sub>12-18</sub>, and 40 units of reverse transcriptase. Incubation was for 90 min. at 37°C, then the reaction was terminated by the addition of EDTA to 10 mM and SDS to 0.5%. After phenol extraction and precipitation by ethanol, the material was resuspended in a small volume of 1mM EDTA, adjusted to 0.3 M NaOH, heated at 65°C for 20 min., then layered onto an alkaline 5-30% sucrose gradient containing 0.3 M NaOH, 0.5 M NaCl, and 20 mM EDTA. Centrifugation was performed at 20°C in an SW50.1 rotor at 49,000 rpm for 4 hours. Fractions corresponding to DNA of greater than 5000 nucleotides in length were pooled, neutralized and precipitated with ethanol. A portion of this cDNA (100 ng) was incubated in a 10  $\mu\text{l}$  reaction containing 100 mM sodium cacodylate, pH 7.2, 2mM  $\text{CoCl}_2$ , 0.5 mM EDTA, 500  $\mu\text{M}$  dCTP, and 5 units of terminal deoxynucleotidyl transferase for 5 min. at 37°C. The reaction was terminated by addition of EDTA to 10 mM and heating at 100°C for 1 min. This poly(dC)-tailed cDNA was

annealed to the directional cloning vector, pKMG (50 ng), in 200 mM sodium acetate, for 1-2 hours at 42°C. The annealing mixture was precipitated with ethanol, then resuspended in 20 µl of reverse transcriptase buffer (described above). After addition of reverse transcriptase (5 units) and incubation for 1 hour at 37°C, 80 µl of a solution containing 50 mM Tris-HCl, pH 7.2, 10 mM MgCl<sub>2</sub>, 20 mM dithiothreitol (DTT), 1 mM ATP was added. This mixture was incubated at 20°C for 10 hours with 50-100 units of T4 DNA ligase. A portion of the resultant mixture was used to transform *E. coli* MC1000 cells, which were then plated on nutrient agar containing 20 µg/ml tetracycline. Plasmids from 20 tetracycline-resistant colonies were screened for the size of cloned DNA insert. One clone, pEM3, contained the largest insert DNA (6100 base pairs) and was used for subsequent analyses.

### Sequence of clone pEM3

Plasmid DNA from pEM3 was digested with DNase I under limiting conditions (29) to produce a random assortment of fragments 300-1200 base pairs in length. Additional, defined pieces of pEM3 were prepared by digesting plasmid DNA with PstI plus SalI, isolating the EMC sequence-containing fragments by agarose gel electrophoresis, then redigesting the excised DNA with HpaII. DNA from the DNase I and HpaII reactions was inserted into the vector phage, M13mp9, and used to transfect *E. coli* JM101. White plaques developing from hybrid virus were picked and regrown. The subcloned viral DNA was sequenced by the method of Sanger (29,30). For some reactions, reverse transcriptase was substituted for Klenow fragment in the sequencing reactions (31).

Selected subclones, which gave difficult regions of band compression on sequencing gels, were also analyzed by chemical methods. Double-stranded, replicative-form viral DNA was isolated from infected JM101 cells (30). The inserted EMC fragments were excised by digestion with EcoRI and BamHI, purified by agarose gel electrophoresis, then labeled with polynucleotide kinase and  $\gamma$ -<sup>32</sup>P ATP (32). DNA strands were separated by polyacrylamide gel electrophoresis, purified, then sequenced by the methods of Maxam and Gilbert (32). All sequencing gels were prepared and run as described (33).

### Primer Extension with Viral RNA Template

Oligodeoxynucleotide primers (9-14 bases in length), complementary to specific, selected segments of the EMC genome, were prepared from New England Biolabs synthesis kits. The preparations (1.5 ml) were titrated to pH 7-8 with acetic acid then applied to mini-C18 columns. After washing with H<sub>2</sub>O (5 ml) and 5% methanol (5 ml) the DNA was eluted with 40% methanol (3 ml), then

concentrated (to 200  $\mu$ l) by evaporation under vacuum. Some primer samples were also purified by fractionation on polyacrylamide gels.

For di-deoxy sequencing experiments, primer (0.1-10  $\mu$ g/ml) and EMC RNA (50  $\mu$ g/ml) were annealed in buffer (50 mM Tris-Cl, pH 7.4, 10 mM  $MgCl_2$ , 50 mM NaCl, 10 mM DTT) by heating for 3 minutes at 67°C and cooling (1°C/min) to 42°C. The complexes were then used in standard reactions with reverse transcriptase (31,34).

For chemical sequencing, the primers were treated with polynucleotide kinase and  $\gamma$ -<sup>32</sup>P ATP (32), heated for 3 min. at 100°C, then used in reverse transcriptase reactions as described above (Molecular Cloning of EMC). The resulting cDNA was separated from excess primer by gel filtration on Sephadex G150, concentrated by lyophilization, then analyzed by the methods of Maxam and Gilbert (32).

#### Computer Assisted Assembly of the EMC Sequence

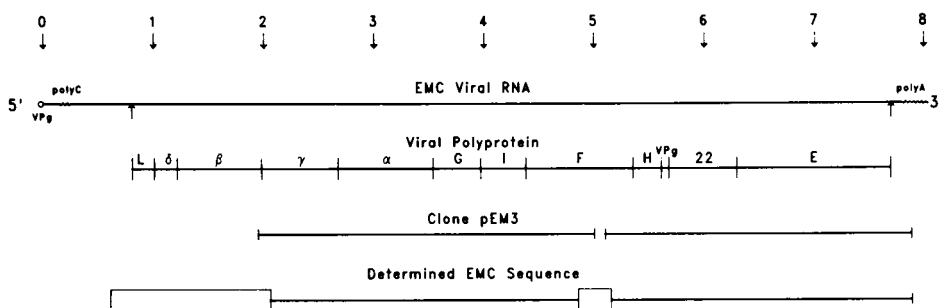
Sequence assembly was aided with programs made available by Rodger Staden (35). Predicted protein translations, sequence manipulation, formatting, and primer searches were carried out with programs developed by John Devereux of the University of Wisconsin Genetics Computer Group (36).

### RESULTS AND DISCUSSION

#### Sequence of Clone pEM3

EMC RNA, extracted from plaque purified virus, was used as template in reverse transcriptase reactions with oligo(dT) primer. The resulting cDNA was elongated at the 5' end with dC residues, purified, then annealed with DNA of the directional cloning vector, pKMG, a derivative of pBR322. After second strand synthesis with reverse transcriptase, the hybrid DNA was treated with DNA ligase and used to transform competent *E. coli* cells. The resulting cloned DNAs were screened for length of the inserted segment, and the largest plasmid, pEM3, chosen for sequence analysis.

The EMC-derived segment within pEM3 was excised, fragmented by digestion with DNase I or HpaII, then subcloned into M13 phage. Di-deoxy and Maxam and Gilbert generated sequence data from 165 DNase I subclones and 38 HpaII subclones was identified, matched, oriented and assembled by computer analysis. A total of 36019 bases, representing an average of 5.9 determinations per nucleotide, were used to compile the 6069 base sequence of the EMC segment within pEM3. The entire lengths of both strands of the inserted piece were covered by these analyses. Figure 1 shows the orientation



**Figure 1.** Schematic representation of EMC viral RNA, polyprotein, clone pEM3 and the 7200 base sequence which is depicted in Figure 2. Those bases sequenced by indirect RNA techniques are represented by boxed areas. The break within pEM3 designates the position of the 89 base deletion. The diagram is drawn to scale. Numbers at the top of the figure represent kilobase distances from the 5' end of the RNA.

of pEM3-determined sequence within the EMC genome. The cloned fragment includes the previously identified 3' terminus of the RNA (37) and extends 5' into the region coding for protein  $\beta$ .

#### The Deletion within pEM3

To encode the large polyprotein (220-260,000 MW) (10,38), EMC RNA must contain an uninterrupted translational reading frame at least 6000 bases in length. However, despite extensive, unambiguous nucleotide analysis, we were unable to find a reading frame of sufficient length within the pEM3 sequence. The longest open frame shifted phase near the middle of the cloned fragment. It was therefore apparent the pEM3 segment contained a discrepant sequence. By comparing homologies between our sequence and that of polio (17,18), we were able to localize the suspect pEM3 region to a 50-100 base sequence centering around the position of the apparent translational phase shift (see Fig 1).

The sequence of this region was then re-examined using viral RNA as template in di-deoxy reactions with reverse transcriptase. A synthetic oligodeoxynucleotide, complementary to an unequivocal EMC sequence region, was used as primer. The resulting gel patterns showed clearly that pEM3 had an 89 base deletion relative to viral RNA. Except for the deleted bases (marked with horizontal arrows on Fig. 2), the indirect RNA sequencing data exactly matched all other nucleotide assignments determined from pEM3.

While the exact molecular event responsible for the deletion within pEM3 is unknown, it is interesting to note that 9 bases at either end of the deleted

segment have almost identical sequences. The 5' sequence is CCAAUCUGC (Fig. 2, bases 4421-4429) and the 3' sequence is CCAGUCUGC (bases 4510-4518). Possibly, during viral replication or cloning, a polymerizing enzyme skipped from one repeated segment to the other, deleting the intervening bases. We consider it unlikely that such a phenomenon would occur more than once within a molecule, creating other deletions which would be undetectable unless they altered the reading frame. Therefore, we assume the sequence of pEM3, with the exception of these 89 bases, accurately reflects the sequence of the viral RNA from which it was derived.

Insertion of the missing bases into the pEM3 sequence established a long open reading frame, ending at UAGUAG codons, 121 bases from the 3' end of the genome (see Fig. 3). Previous sequencing experiments identified this location as the correct termination site of the EMC polyprotein (37). The nucleotide sequence of pEM3 with the deleted bases reinserted, is shown as part of Figure 2. The 5' end of the pEM3 segment is base 1040 of this figure and the 3' end is base 7200. Our sequence corrects several previously misidentified bases near the 3' terminus (37).

#### Primer Extension with Viral RNA

Although clone pEM3 was the largest clone generated in our experiments, it represents only 78% of the EMC RNA. It does not include the 5' terminal 1800 bases of the genome. Sequential primer extension reactions with viral RNA as template were used to sequence the remaining portion of the polyprotein coding region. Starting at a sequence within pEM3, a series of oligodeoxynucleotides complementary to successive segments of EMC RNA were selected and synthesized. As each primer was used in di-deoxy reactions with reverse transcriptase, the newly established bases determined the sequence of the next primer in the series. Overlapping sequences from five primers, spaced an average of 220 bases apart, allowed identification of nearly 1200 EMC nucleotides which were not present in the pEM3 clone (see Fig. 1).

The sequence of this segment was verified by chemical techniques. The synthetic oligonucleotides were radiolabeled with polynucleotide kinase and  $\gamma$ -<sup>32</sup>P ATP, then used to prime reactions with reverse transcriptase and viral RNA. The resulting cDNA products were analyzed by the methods of Maxam and Gilbert. The combined indirect RNA sequencing data, with an average of 4.1 determinations per nucleotide, extended our EMC sequence to within 600 bases of the 5' end of the genome. The sequence of 7200 bases from the 3' end of EMC is shown in Figure 2.

5'... UGGAUAGUUGUGGAAAGAGUCAAAUUGGCCUCUCCUACAAGCGUAUUCACAAAGGGGCGUGAAGGAUGCCCGCAGAGGUA 75

76 CCCCUUGUAUGGGAUCUGAUCUGGGGCCUGGGGCGACAUUGCUUUAACUGUGUUUAGUCGAGGUUAAAAACGUC 150

(start).

151 UAGGCCCCCCGAACCAACGGGGACGUGGUUUUCCUUUGAAAACACGAUGAUUAUUGGCCCAACCAUGGAACAA 225  
M A T T M E Q

226 GAGACUUGCGCGCACUCUCUCACUUUGAGGAAUUGCCAAAUGCUCUGCUCUACAAUACCGUAUUGGAUUUUAC 300  
E T C A A H S L T F E E C P K C S A L Q Y R N G F Y

301 CUGCUAAAGUAUGAAGAAUUGGUACCCAGGAGGUUAUUGCUGAUGGAGGAGGAUGUCUUUUGAACCCGAA 375  
L L K Y D E E W Y P E E L L T D G E D D V F D P C E

(L/δ)

376 UUAGACAUGGAAGUCGUUUUCGAGUUAACGGGCAAUUCCACCUCUCAGACAAGAAUACUCCUCCUGGGAAGGC 450  
L D M E V V F E L Q G N S T S S D K N N S S S E G

451 AAUGAAGGUGUGAUCAUCAAUAACUUUUACUCCAACCAUAUAUCAAACUCUCAAUGACCUCUCUGCUAAUGCAGCC 525  
N E G V Y I I N N F Y S Q Y Q N S I D L S A N A A

526 GGGUCUGACCCACCCAGACUGGCUAAUUUUUGGAAUUCUGGGGCGAGUAAUGCCUUUUUUAUUGCUU 600  
G S D P P P R L R S I F E S L S G A V N A F S N M L

(δ/β)

601 CCAUUGCGACUGAUCAAAUAACAGAAGAAUUGGAGAUCUGUCUGAUGCAGGUCUCAAGACACUGCCGGCAUAU 675  
P L L A D Q N T E E M E N L S D R G L K T L P A I

676 CGGUCACAACACCCAGUCAACAGUGGGCGGUUUGUGUACCGGUUACUGAUGGAGGACAGCCGGCAU 750  
R S Q T P S Q W A V L S V M V P F M M E S I R H

751 CAUGUGCUGACACUGCUUCAGAAAAGAUUCUGGGCGUGGAAAGGUACUACACCUUCAAGGUUAUGAUUGGACAU 825  
H V L T L L Q K R F W R W K G T T P S R L M I G H

826 CAACACAAGAGCCCUUUGAGUACAUCGCGAUUCCCUUCCUACGUCUGCCGUGAAGAUUGGUGGUCUUUG 900  
Q H G K S P L S T S A F P F L T S C P V K M V V S L

901 GUGCCCUCCGGCCACUACUGGAGAAACUGGAGUGGCGGGGUGCAUGUACAGUACAGCCCUCAAAUCCAC 975  
V A L R R R H Y L V K T G W R V Q V C N A S C Q F H

976 CGUGGAGGUUUGCUGGUGUUAUGGCCACAGAGUAUCCAACCUUAAGAUUCUUUUGCCAUUGGACAACCGUUGGUCC 1050  
A G G L L V F M A P E Y P T L D A F A M D N R W S

1051 AAGGAUAACCGCCUUAUGGAACACCAUCAGACAACAAAGAGGACCAUUUGGCCAUUGGACAGACAACUUC 1125  
K D N L P N G T R A T C N K K G P F A M D H Q N F

1126 UGGCAGUGGACCUUGUAUCCCAUCAAUCCUGAUAUCUGAGAACUAACACCACAGUGGAUCUUGAGGUGCCAUAU 1200  
W Q W T L Y P H Q F L N L R T N T T V D L E V P Y

1201 GUAAACAUAAGCCCCCAUUCUCCUGGACACAACAUUGCUUCCUGGACUUGUGGAUUGCAGUGGUUGCUCCCGUC 1275  
V N I A P T S S W T Q H A S W T L V I A V V A P L

1276 ACAUAUCAACCGGGGCUUACCAAGUUUGGAUAUACCGGUUUAUUCAGGCAUUAAGGCGUGUCUUUAUUGGC 1350  
T Y S T G A S T S L D I T A S I Q P V R P V F N G

(β/γ)

1351 CUCCGGCAUGAGACACUUUCUAGACAGCGCCCAUUCGGGUCACAAUAGAGAACAUCUGGUUACCGGUUUAUUCU 1425  
L R H E T L S R Q S P I P V T I R E H A G T W Y S

1426 ACUCUGCCAGACAGUGGCCUUAUUUUGGCAAGACUCCUGGUUUGCUCAUCCAAUUAUUGGAGGCGGAAUAC 1500  
T L P D C S T C A P G P I Y G K T P V A P S N Y M V G E Y

1501 AAGGACUUCUGGAGAUAGCUCAGAUUCCAACCUUUAUUGGAAUAGAGUCCUUAUUGCUGUCCCUACAUUGAG 1575  
K D F L E I A Q I P T F I G N K I P N A V P Y I E

1576 GCAUCACAACAGCCGCUAAGACCCAAACCGGUGGGCACCUAUCAAUGAGACCUUGCUGCUGCUGGCCAAU 1650  
A S N T A V K T Q P L A T Y Q V T L S C S C S C L A N

1651 ACAUUCUUGGGCGGUUUGUUGUAGAAACUUGCUGACUACCGGGAUUAUGGUUUUAUCCUUUGUGUACUGGG 1725  
T F L A A L S R N F A C Q Y R G S L V Y T F V F T G

1726 ACCGGAUGAUGAAGGGCAAGUUCUCAUUGGCUAACCCACCUGGAGCGGGCAAGCCCAUAUGUCGAGACCAA 1800  
T A M M K G K F L I A Y T P P G A G K P T S R D Q

1801 GCCAUCAGGCGCAUUAUGCAUUUGGGAUUUGGGCUAAAUUUCUUAUCCUUAACUGUCCUUUUAUUCU 1875  
A M Q A T Y A I G I W L D L G L N S Y S F T V P F I S

1876 CCCACUCACUUCCGCAUGGUAGGUACUGACCAAGUCAACAUCACUAAUUGCGGAUGGCGUGGUUAACCGUGGGCAG 1950  
P T H F R M V G T D Q V N I T N A D G W V T V W Q

1951 CUCACUCCCUCAUUAACCCACGAGGACGCCGACCUUAGAUUAACAUAUUGGAGGCGAGGCAAGGAAU 2025  
L T P L T Y P P G C P T S A K I L T M V S A G K D



(7/α)

2026 UUCU<sup>˙</sup>CACUAAGAUG<sup>˙</sup>CCUAUCUAC<sup>˙</sup>CUGCC<sup>˙</sup>CCUUGGAG<sup>˙</sup>CCUACGGGAGUAGA<sup>˙</sup>AAACGCUGAAA<sup>˙</sup>AGGGGUCAC<sup>˙</sup> 2100  
F S L K M P I S P A P W S P Q G V E N A E K G V T

2101 GAAACACA<sup>˙</sup>AAACGCAACUG<sup>˙</sup>CUGACUUGUG<sup>˙</sup>GCUAACCA<sup>˙</sup>GUUACUUGC<sup>˙</sup>CUGAGA<sup>˙</sup>ACCA<sup>˙</sup>AAACGAAGUGG<sup>˙</sup>GCUUUC 2175  
E N T N A T A D F V A Q P V Y L P E N Q T K V A F

2176 UUCUA<sup>˙</sup>UAUAGGUC<sup>˙</sup>CAGUCC<sup>˙</sup>UAUUGG<sup>˙</sup>GUCCUAC<sup>˙</sup>CGUGAAGUC<sup>˙</sup>CGCAGUCUAG<sup>˙</sup>AAUCUGGUUU<sup>˙</sup>UGCCCGUUC 2250  
F Y N R S S P I G A F T V K S G S L E S G F A P F

2251 UCUA<sup>˙</sup>AUGGAC<sup>˙</sup>UUGCCG<sup>˙</sup>AAUCAGUGA<sup>˙</sup>UACUGAC<sup>˙</sup>CCUUGG<sup>˙</sup>CCCAAU<sup>˙</sup>UUGACCC<sup>˙</sup>CCCUAUGACCA<sup>˙</sup>ACUCAGG 2325  
S N G T C P N S V I L T P G P Q F D P A Y D Q L R

2326 CCAC<sup>˙</sup>AGCUGACAG<sup>˙</sup>AAAUUGG<sup>˙</sup>GGCAAUGAA<sup>˙</sup>UGAGGAGAC<sup>˙</sup>CUAAA<sup>˙</sup>AGUCUU<sup>˙</sup>UCCGCUUAA<sup>˙</sup>UCCAAC<sup>˙</sup>ACAG 2400  
P Q R L T E I W G N G N E E T S K V F P L K S K Q

2401 GAUUAUCC<sup>˙</sup>UUCUGCC<sup>˙</sup>CUUC<sup>˙</sup>UCCCC<sup>˙</sup>UUUGUGUAUUA<sup>˙</sup>AAUGUGAUUU<sup>˙</sup>UAGAAGUGAC<sup>˙</sup>CUUAGUCU<sup>˙</sup>ACACU 2475  
D Y S F C P N S V I L T P G P Q F D P A Y D Q L R

2476 UCAGG<sup>˙</sup>CAACCAUGG<sup>˙</sup>CGUGUUGG<sup>˙</sup>UGAGGUGG<sup>˙</sup>UGUCC<sup>˙</sup>CUGGUAC<sup>˙</sup>ACCAAC<sup>˙</sup>CAAG<sup>˙</sup>CCACUACCC<sup>˙</sup>AGGUUCUCCAU 2550  
S G N H G L L V R W C P T G T P T K P T T Q V L H

2551 GAAGUAAG<sup>˙</sup>UCCUCUCAG<sup>˙</sup>GAAGG<sup>˙</sup>CAGAACC<sup>˙</sup>CCCC<sup>˙</sup>AGGUUU<sup>˙</sup>UAGUGCCG<sup>˙</sup>ACCUGGCAU<sup>˙</sup>UCAAU<sup>˙</sup>CAGAUUUCC 2625  
E V S S L S E G R T P Q V Y S A G P G I S N Q I S

2626 UUGUA<sup>˙</sup>AGUUCUUA<sup>˙</sup>CAAUUCC<sup>˙</sup>CAUUU<sup>˙</sup>CAGUCC<sup>˙</sup>UAUCAGCUG<sup>˙</sup>UGUUAUAA<sup>˙</sup>UGGACACAAG<sup>˙</sup>AGAUUUGACA<sup>˙</sup>AC 2700  
F V V P Y N S P L S V L S A V W Y N G K R F D N

2701 ACUGG<sup>˙</sup>GAGC<sup>˙</sup>UUGGGCAUUG<sup>˙</sup>CCCC<sup>˙</sup>UAUUUC<sup>˙</sup>GAUUU<sup>˙</sup>CGGC<sup>˙</sup>ACUCUGU<sup>˙</sup>UCUU<sup>˙</sup>UGCUGG<sup>˙</sup>CAC<sup>˙</sup>AAAGCCU<sup>˙</sup>GAC<sup>˙</sup>AUUAAA 2775  
T G S L G I A P N S D F G T L F F A G T K P D I K

2776 UUCAC<sup>˙</sup>AGUCUACU<sup>˙</sup>GAUACAG<sup>˙</sup>AAUAGAGAG<sup>˙</sup>UUUUUG<sup>˙</sup>CCCC<sup>˙</sup>AGUCC<sup>˙</sup>GACUG<sup>˙</sup>UCUUU<sup>˙</sup>UCCCC<sup>˙</sup>UGGCCCC<sup>˙</sup>ACU 2850  
F T V Y L R Y K N K R V F C P R P T V F F P P T

2851 UCCGG<sup>˙</sup>GAGACAAG<sup>˙</sup>AUGUAU<sup>˙</sup>AGCCCC<sup>˙</sup>GAGAG<sup>˙</sup>CGUGGAG<sup>˙</sup>UCUUGAUG<sup>˙</sup>CUAG<sup>˙</sup>AGAGUCCA<sup>˙</sup>AAU<sup>˙</sup>GCCCUAG<sup>˙</sup>CAUUUCA 2925  
S G D K I D M T P R A G V L M L E S P N A L D I S

(α/G)

2926 AGAACAU<sup>˙</sup>ACCC<sup>˙</sup>CAGGUACUGU<sup>˙</sup>CUCAU<sup>˙</sup>CAAU<sup>˙</sup>UACAAC<sup>˙</sup>UAG<sup>˙</sup>AGGUU<sup>˙</sup>UGGAGG<sup>˙</sup>UUAGU<sup>˙</sup>UUUAGACAU<sup>˙</sup>UGG<sup>˙</sup>A 3000  
R T Y P T L H V L I Q F N H R G L E V R L F R H G

3001 CACUUU<sup>˙</sup>UGGG<sup>˙</sup>CUCAC<sup>˙</sup>AGUG<sup>˙</sup>CGGAGUGA<sup>˙</sup>UUCUGA<sup>˙</sup>GAAC<sup>˙</sup>AGCA<sup>˙</sup>AAACAGG<sup>˙</sup>UCUCU<sup>˙</sup>UCCUGAG<sup>˙</sup>CAACGGG 3075  
H F W A E T R A D V I L R S K T K Q V S A F T

3076 AACU<sup>˙</sup>ACCCGUA<sup>˙</sup>AAUG<sup>˙</sup>GACUCUAG<sup>˙</sup>ACUCC<sup>˙</sup>CGGA<sup>˙</sup>AAUCCU<sup>˙</sup>UGGA<sup>˙</sup>AAUACCU<sup>˙</sup>ACC<sup>˙</sup>AGGCGGU<sup>˙</sup>CUU<sup>˙</sup>AAAGAGCAGAA 3150  
N Y P S M D S R A P W N P W K N T Y Q A V L R A E

3151 CCAUG<sup>˙</sup>AGAG<sup>˙</sup>GACCAUGG<sup>˙</sup>AUAUAAUUA<sup>˙</sup>UAGAGAG<sup>˙</sup>UCAGGCC<sup>˙</sup>UUUAG<sup>˙</sup>ACUG<sup>˙</sup>CCCC<sup>˙</sup>UGGU<sup>˙</sup>CAGAG<sup>˙</sup>AAUGG 3225  
P C R V T M D I Y Y K R V R P F R L P L Q K E W

3226 CCCG<sup>˙</sup>UGCGAGAG<sup>˙</sup>GAACGUUU<sup>˙</sup>UGGUU<sup>˙</sup>GUAC<sup>˙</sup>CGGAUCU<sup>˙</sup>CAAU<sup>˙</sup>UGCC<sup>˙</sup>CACUAG<sup>˙</sup>CGUGGU<sup>˙</sup>ACUU<sup>˙</sup>UGCGGACCUA 3300  
P V R E E N V F G L Y R I F N A H Y A G Y F A D L

(G/I)

3301 CUGAUU<sup>˙</sup>CAUGACA<sup>˙</sup>UUGAGACA<sup>˙</sup>AAUCC<sup>˙</sup>AGGG<sup>˙</sup>CCCUCAUGU<sup>˙</sup>UAGAC<sup>˙</sup>CAAG<sup>˙</sup>GAACAGG<sup>˙</sup>UUUCC<sup>˙</sup>AGACCC<sup>˙</sup>AAAGGA 3375  
L I H D I E T N P G P F M F R P R K Q V F Q T Q G

3376 GCGG<sup>˙</sup>CAGUGUCA<sup>˙</sup>UAGG<sup>˙</sup>CUCAA<sup>˙</sup>CCCUACUG<sup>˙</sup>CCGAAC<sup>˙</sup>GACCUU<sup>˙</sup>GCCAGCA<sup>˙</sup>AGCUA<sup>˙</sup>UGGGAUC<sup>˙</sup>AGCUUU<sup>˙</sup>UACG 3450  
A A V S S M A Q T L L P N D L A S K A M G S A F T

3451 GCUU<sup>˙</sup>UGCUC<sup>˙</sup>GAUGCCA<sup>˙</sup>ACGAGG<sup>˙</sup>ACGCCA<sup>˙</sup>AAAGCA<sup>˙</sup>UAG<sup>˙</sup>AGAUUA<sup>˙</sup>AAAGACA<sup>˙</sup>UUAAGU<sup>˙</sup>UCUCU<sup>˙</sup>UACGGAUGCA 3525  
A L L D A N E D A Q K A M K I I K T L S S L S D A

3526 UGGG<sup>˙</sup>AAAAUGUAAA<sup>˙</sup>GAACACUA<sup>˙</sup>AAACACCC<sup>˙</sup>AGAGU<sup>˙</sup>UCUGGA<sup>˙</sup>AGCAGC<sup>˙</sup>CUUGAG<sup>˙</sup>CAGAU<sup>˙</sup>GUGUG<sup>˙</sup>CAGCUGAUU 3600  
W E N V E T L N N P E F W K Q L L S R C V Q L I

3601 GCAGG<sup>˙</sup>GAUGACA<sup>˙</sup>UAG<sup>˙</sup>CAGUAG<sup>˙</sup>CAU<sup>˙</sup>CCGACCC<sup>˙</sup>UUGACUCUG<sup>˙</sup>CUCUG<sup>˙</sup>CUUAG<sup>˙</sup>GAACA<sup>˙</sup>UUGAGCG<sup>˙</sup>CGCCG<sup>˙</sup>CAG 3675  
A G M T I A V M H P D P L T L L C L G T L T A A E

3676 AUUA<sup>˙</sup>CAAGCAGACA<sup>˙</sup>AGUCUG<sup>˙</sup>CGGAAGAA<sup>˙</sup>UAGCAG<sup>˙</sup>CUAAGU<sup>˙</sup>CAAGACA<sup>˙</sup>UUU<sup>˙</sup>UACACUC<sup>˙</sup>CUCCACC<sup>˙</sup>ACGG 3750  
I T S Q T S L C E E I A A K F K T I F I T P P P R

(I/F)

3751 UUUCC<sup>˙</sup>CACA<sup>˙</sup>AUCUCU<sup>˙</sup>UUUCC<sup>˙</sup>AAACA<sup>˙</sup>AUCC<sup>˙</sup>CCCUUG<sup>˙</sup>AAACAGGUA<sup>˙</sup>AAUGAU<sup>˙</sup>UUU<sup>˙</sup>UCCCU<sup>˙</sup>AGCC<sup>˙</sup>AAAGAC 3825  
F P T I S L F Q Q F S P L K Q V N D I F S L N G M A

3826 CUGG<sup>˙</sup>ACUGG<sup>˙</sup>CCCUA<sup>˙</sup>AGACUG<sup>˙</sup>UGG<sup>˙</sup>AAAGGUGG<sup>˙</sup>UUGAU<sup>˙</sup>UGGUU<sup>˙</sup>UGGACA<sup>˙</sup>UGG<sup>˙</sup>AUAGU<sup>˙</sup>ACAGG<sup>˙</sup>AGGAAAGGA<sup>˙</sup> 3900  
L D W A V K T V E K V V D W F G T W I V Q E E K E

3901 CAGAC<sup>˙</sup>CCUAG<sup>˙</sup>AUCAGCUCU<sup>˙</sup>UGCAG<sup>˙</sup>CGUUU<sup>˙</sup>CCCC<sup>˙</sup>GAAC<sup>˙</sup>UAG<sup>˙</sup>CGAAGCGCA<sup>˙</sup>UUUCU<sup>˙</sup>AGU<sup>˙</sup>UCC<sup>˙</sup>GGAAUGGA<sup>˙</sup>AUGGCC 3975  
Q T L D Q L D Q F P E H A K R I S D L F S L N G M A

3976 GCCU<sup>˙</sup>AUGUAGAG<sup>˙</sup>UGCAAG<sup>˙</sup>GAGAGU<sup>˙</sup>UUU<sup>˙</sup>GAUUU<sup>˙</sup>CUU<sup>˙</sup>UAGAA<sup>˙</sup>AGCUGA<sup>˙</sup>CAACU<sup>˙</sup>AGG<sup>˙</sup>CAGUGA<sup>˙</sup>AAAG<sup>˙</sup>AGAGAGA<sup>˙</sup>ACG 4050  
A Y V E C K E S F D F F E K L Y N Q A V K E K R T

4051 GGUAUCGCCGCCGUCUGUGAAAAUUCAGACAGAGCAUGACCACGCCACCCGUCGGUGUGAGCCAGUCUGUAU 4125  
 G I A A V C E K F R Q K H D H A T A R C E P V V I  
 4126 GUGCUCGCCGGAGACCGGGGCAAGGGAUUCUUUAUCAAGUCAGGUUAUUGCCAGGCCGUCUCAAGACCAU 4200  
 V L R G D A G Q G K S L S S Q V I A Q A V S K T I  
 4201 UUCGGCCGGCAUUCUGUGUAUUCUUUCCCCCGAUUCGGAUUUUCUUGAUGGCUAUGAAUACAGUUUGCAGCA 4275  
 F G R Q S V Y S L P P D S D F F D G Y E N Q F A A  
 4276 AUAAUGGAUGAUCUAGGGCAAAUCCUGAUGGCUCUGAUUUACUACGUUCUGUCAGAUUUUGACUACCAU 4350  
 I M D D L G Q N P D G S D F T T F C Q M V S T T N  
 4351 UUUUCUCCCAAUUAGGCUAGUCUAGAGAGAAAGGGCACCCUUUACAUCUACGUUGUGGGGCAACUACCAU 4425  
 F L P N M A S L E R K G T P F T S Q L V V A T T N  
 4426 CUGCCUGAGUUUAGCCUGUCACAUAAGCCCAUUAACCCUGCUGUUGAGAGAAGGAUACUUCGACUUAUUCAGUG 4500  
 L P E F R P V T I A H Y P A V E R R I T F D Y S V  
 4501 UCUCUGGUCCAGUCUGCUCAAAACAGAGGCCGGGUUAUAGGUUUUGGAUUGUUGAAAGCCUUUAGGCCUACC 4575  
 S A G P V C S K T E A G Y K V L D V E R A F R P T  
 4576 GGUAGGGCUCCUUCUUGCUUCCAGAAUAACUGCCUUUUUCCUUGAGAAAGCUGGGCUCCAGUUCAGAGUAUAC 4650  
 G E A P L P C F Q N N C L F L E K A G L Q F R D N  
 4651 CGAACUAAAGAGAUUUUCCUGGUAGAUUGUAUGAGAGCCGUGGCUGAGGAUUGAAGGAAGAAGAAAGUU 4725  
 R T K E I I S L V D V I E R A V A R I E R K K K V  
 4726 CUCACAACCGUGCAGACCCUUGUGGCACAGGUCCAGUAGACAGGUCAGUUUCCAUUCCGUAGUCCAGCAGCU 4800  
 L T T V Q T L V A Q G P V D E V S F H S V V Q Q L  
 4801 AAAGCAAGACAGCAAGGCAGAGUAACAGCUUGAGGAUUGCAAGGCCUUUCGCAAGAUACAGGAGCGUAAC 4875  
 K A R Q A T D E Q L E E L Q E A F A G L Q F R D N  
 4876 UCUGUUUUUCUGAUUGGUUGAAGAUUUUCUGCAUUGUUGUGUGCGGACUUGGCACUUUCCAAUGUUGAACG 4950  
 S V F S D W L K I S A M L C A A T L A L S Q V V K  
 4951 AUGGCCAAGCGGUGAAGCAAGGUGCAAGCCUGAUCUGGUUGUGUGCAUUGGAUGAGCAGGAGCAGGGACCU 5025  
 M A K A P L P C F M V K P D L V R V Q L D E Q E G G P  
 5026 UACAAGAGACAGCGAGAGUUAACCAAAACACUGCAGUUGUUGACAUUCAGGACCAAAACCCUGUAGUGGAC 5100  
 Y N E T A R V K P K T L Q L L D I Q G P N P V M D  
 5101 UUGUAAAAUUGUAGCCAACAUGUAACCGCCCCAUUGGUUUUGUCUACCCACUGGGUGAGCACCCAGACU 5175  
 F E K Y V A K H V T A P I G F V Y P T G V S T Q T  
 5176 UGCCUCCUUGUGAGAGCGCCGACCUUGGUAGUAAUAGACACAUGGCCGAGUCUGACUGGACUUCUAGUAGUG 5250  
 C L L L V R G R T L V N R H M A E S D W T G T F L K A  
 5251 CGUGGAGUCACACAGCCCCGUCUACUGUUAUUUUUGGCCAUAGCUAAAGCAGGCAAGAGACUGACGUUUCU 5325  
 R G V T H A R S T V K I L A I A K A G K E T D V S  
 5326 UUCAUCCGCCUCUCUUCUGGUCCUCUUAUUCAGAGACAAUACAUCCAAUUGUUGAAGGCGUGUGAUGUACUCCU 5400  
 F I R L S S G P L F R D N T S K F V K A G D V L P  
 5401 ACUGGUGCCGCCUCCAGUCACGGGGAUUAAGAACCGGCACUACCAUGAUGUACACAGGACCUUCCUGAAGCU 5475  
 T G A A P V T G I M N T D I P M M Y T G T F L K A  
 5476 GGUUGUGACUCCAGUGGAAACCGGCCAGACCUUUAUACUGUAUUAUUAACAGGCUAACACACGAAGGGC 5550  
 G V S V P V E T G Q T F N H C I H Y K A N T R K G  
 5551 UGGUGUGGACAGCCCUACUGGCGAGAUUCUGGAGGAAGCAAGAAAUCUUGGCAUCCAUUCUGUGGCUCUUG 5625  
 W C G S Q A L L A D L G G S K K I L G I H S A G S M  
 5626 GGAUAGCCGCCGCCUUGAUGUGUCACAGGAGAUUCGGCGGUAUGAUGCCUUUGAGGCCACAGGGGCU 5700  
 G I A A A S I V S Q E M I R A V V N A F E P Q G A  
 5701 CUCGAGAGAUUGCCAGAUUGGCCCGGUUAUUCAGUACCAAGGUAACAGCACUACGCCCCACCGUUGCCCGUCAA 5775  
 L E R L P D G P R I H V P R K T A L R P T V A R Q  
 5776 GUCUCCCAACAGCAUAGCCCCGGCUGUUCUUAUCGAAUUGAAGCCUAGAACAGGCGUGAUGAUGAUG 5850  
 V F Q P A P A V L S K F D P R T E A D G L E V  
 5851 GCUUUCUCCAAACAUAUCCCAACAGGAAGAGCCUCCACAGUGUUUAGAAUGGUAGGCAAGAGUAUGCCAAU 5925  
 A F S K H T S N Q E S L P P V F R M V A K E Y A N  
 5926 AGAGUUUACCUUUGCGGGAAGACAAUGGCCGUCUGACUGUAAGCAGGCUUUGGAAGGACUGGAGGGG 6000  
 R V F T L L G K D N G R L T V K Q A L E G L E G M  
 6001 GACCCCAUGGACAGGAACACCUCCCGGGGUCCAUAUACUGCGCUAGGAUUGCGCAAGACAGAUUGCUGAU 6075  
 D P M D R N T S P G L P Y T A L G M R R T D V V D

```

6076 UGGGAAUACAGCCACCCUGAUCCCGUUGCGGCAGAAAGAUUAGAAAAUAGAAUAGGAGACUUUUCCGAAGUU 6150
      W E S A T L I P F A A E R L R K M N E G D F S E V
6151 GUCUAUCAAACAUCCUCAAGGAUGAGCUUAGACCGAUAGAGAAGGUUCAAGCCGCCAAGACACGGAUUGUAGAU 6225
      V Y Q T F L K D E L R P I E K V Q A A K T R I V D
6226 GUUCCACCAUUUGAGCAUUGCAUUCUGGGUAGACAUAUUGUUGGGAAGUUUGCAUCAAAGUUCAGACCCCAACCG 6300
      V P P F E H C I L G R Q L L G K F A S K F Q T Q P
6301 GGUCUGGAACUAGGAUCAGCAUUGGAUGUAGACCCAGAUUGUACACUGGACUGCCUUCGGUGUGCGCAUUGCAAGGU 6375
      G L E L G S A I G C D P D V H W T A F G V A M Q G
6376 UUGAGCGGUGUCUACGAUGUGGACUACUCCAACUUUGAUUCGACCCAUUCGGUGGCAUUGUCCGUUAUUGGCU 6450
      F E R V Y D V D Y S N F D S T H S V A M F R L L A
6451 GAGGAUUUUUUCACUCCAGAGAUGGUUUUGACCCCGACUAGAGAAUAUCUUGAGCAUUUAGCCAUUUUACCC 6525
      E E F F T P E N G F D P L T R E Y L E S L A I S T
6526 CAUCCGUUUGAGGAGAAGCGCUUUCUGAUAAACCGUGGUCUCCCAUCAGGUUGUGCAGCGACCCUCAAUGCUAAAC 6600
      H A F E E K R F L I T G G L P S G C A A T S M L N
6601 ACUAUAUAGAAUAAUUAUAAUAAUAGGGCGGGUUGUAUCUCACGUAAUAAAAUUUGAAUUUUGAUGAUGAGAAG 6675
      T I M N N I I I R A G L Y L T Y K N F E F D D V K
6676 GUGUUGUGCUACGGAGAUGAUCUCCUUGUGGCCACAAAUAACCAUUGGAUUUUGAUAAGGUGAGAGCAAGCCUC 6750
      V L S Y G D D L L V A T N Y Q L D F D K V R A S L
6751 GCAAGACAGGAUAUAAGAUAACUCCCGUAACACAACUUCUACCUUUCUCUUAUUCGACGCUUGAAGACGUU 6825
      A K T G Y K I I T P A N T T S T F P L N S F L E D V
6826 GUCUUCUAAAAAGAAAGAUUUAAGAAAGAGGCCUUCUGUAUCGGCCUGCAUGAACAGAGAGGCCGUUGGAAGCA 6900
      V F L K R K F K K E G P L Y R P V M N R E A L E A
6901 AUGUUGUCAUACUUCGUCCAGGACUCUAUCUGAGAAAUCACUUCGAUCACUAGUCCUUGCCGUUCAUUCUGGC 6975
      M L S Y Y R P G T L S E K L T S I T M L A V H S G
6976 AAGCAGGAUAUGAUCGGCUCUUCUUGCCCCAUUCCUGAGGUAGGGUUGUGUGGCAUCAUUCGAGAGUGUGGAG 7050
      K Q E Y D R L F A P F R E V G V V V P S F E S V E
7051 UACAGAUGGAGGAGUCUGUUCUGUAGUAGUCACUGGCACAACGCGUACCGGUAAGCCAAUCGGGUUAU 7125
      Y R W R S L F W * *
7126 ACACGGUCGUCAUACUGCAGACAGGUUCUUCUACUUGCAAGAUAAGUCUAGAGUAGUAAAAUAAUAGAUAGAG polyA

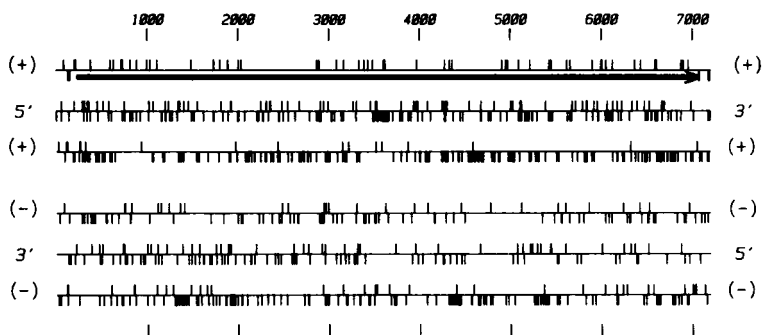
```

Figure 2. The nucleotide and deduced amino acid sequences of 7200 bases from the 3' end of EMC. Polyprotein translation start and stop codons and proteolytic cleavage sites are highlighted with boldface characters. The vertical arrow at position 1040 designates the 5' end of the pEM3 segment. The right and left arrows at positions 4430 and 4518 represent 5' and 3' ends respectively, of the deletion within pEM3.

#### Localization of the Polyprotein

The long open reading frame first detected within pEM3 (designated as frame 1), extends 5' into the region determined by indirect RNA sequencing (see Figure 3). It is bounded on the amino end by a UAG termination codon in position 151. This figure also shows that long open reading frames are not evident in either of the two remaining positive phases. Therefore, viral translation must begin 3' to base 153 within frame 1.

The EMC polyprotein begins with the sequence Met-Ala-Thr followed by a lysine tryptic fragment containing: Glu, Cys, Met, Ser, Pro, Leu, Phe and possibly other (previously unidentified) amino acids (39). The Met-Ala-Thr combination occurs only once within our phase 1 reading frame, beginning at



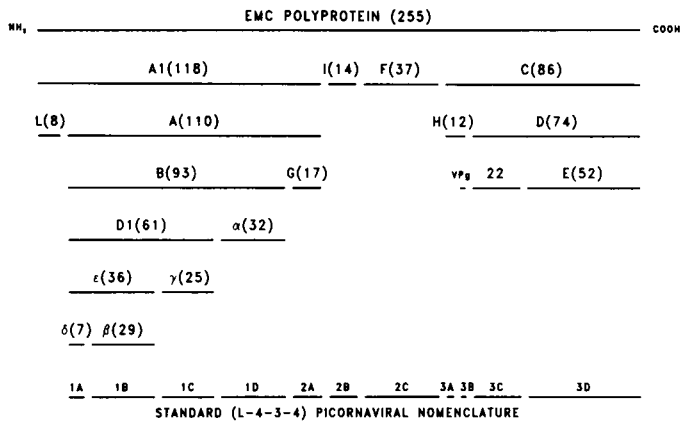
**Figure 3.** Start and stop codons within the EMC sequence. All positive (+) and negative (-) reading frames within the EMC sequence are represented schematically. Tic marks above the lines indicate positions of AUG codons, tic marks below the lines indicate positions of UGA, UAG or UAA codons. The open reading frame encoding the viral polyprotein is marked with a boldface arrow (top line of the figure) and is designated as reading frame 1. The values across the top of the figure orient base position according to the numbering scheme of figure 2.

base position 205. This deduced sequence is predicted to be part of a tryptic fragment ending in lysine, and containing: Glu, Cys, Met, Ser, Pro, Leu, Phe, Gln, Ala and His (see Fig. 2). Since this sequence so clearly matches the established amino end of the polyprotein, we believe EMC translation must initiate at the AUG triplet in position 205. The reading frame extends 2290 codons to the double termination triplets at position 7075-7080 and predicts a polypeptide of 255466 molecular weight.

#### The Proteolytic Cleavage Map

Mature virion capsid and other viral non-structural proteins arise by proteolytic cleavage of the polyprotein. Many carboxy- and amino-termini of these proteins have been identified (1,37,ACP, manuscript in preparation). This information, combined with tryptic peptide analyses (14,40), molecular weight determinations (10,12,38,41,42), amino acid compositions (1,43,44) and homology with FMDV, polio and Mengo (17,21,25,45) have allowed alignment of the viral polypeptides with the deduced amino acid sequence of the open reading frame. The proteolytic processing scheme is summarized in Figure 4.

The viral capsid proteins,  $\delta$ ,  $\beta$ ,  $\gamma$  and  $\alpha$  (7247, 29026, 25141 and 31703 MW, respectively) are preceded by a short leader peptide (7929 MW) and followed by two small proteins designated G (17474 MW) and I (14709 MW). Protein F (36566 MW) lies adjacent to protein H (12166 MW), a VPg-containing precursor (40). The single VPg sequence (2410 MW) is located at the carboxy terminus of H. The



**Figure 4.** Proteolytic processing map of the EMC polyprotein. EMC proteins and proteolytic cleavage precursors are shown in their relative orientation to the viral polyprotein. The map is drawn to scale. Parentheses indicate the molecular weights (in kilodaltons) of each polypeptide. Standard (4-3-4) picornaviral nomenclature designations (46) are shown at the bottom of the diagram.

protease, p22 (21751 MW), and polymerase, E (51934 MW), were previously mapped (14). Precursor proteins, A1, A, B, C, D, D1 and  $\epsilon$ , are shown above their cleavage products.

The processing scheme conforms to the standard 4-3-4 picornaviral polypeptide arrangement (46), and except for the leader sequence is very similar to the primary cleavage map of polio virus (17). Since the leader is most closely analogous to that of FMDV, in some respects the overall structure of EMC can be considered intermediate between the entero- and aphthoviruses. EMC has a leader peptide like FMDV, but the rest of the polyprotein is aligned and processed like polio.

Other reports have suggested the N-terminal position of the EMC cleavage map is occupied by leader peptides (p12/p14) of 12-14000 molecular weight (47). Our protein alignments, which are consistent with the deduced amino acid sequence, indicate a much smaller (7929 MW) leader peptide. However, preliminary nucleotide sequencing data show the possible existence of another open translational reading frame, located between the genomic poly(C) tract and the polyprotein coding region. This segment is potentially long enough to code for protein(s) of at least 14000 molecular weight, and may be the source of p12 and p14. This potential reading frame does not contain a Met-Ala-Thr sequence (ACP, unpublished observations).

<u>VIRAL PROTEASE CATALYZED CLEAVAGE SITES</u>		<u>PROTEINS</u>
Leu-Ser-Arg-GLN /	SER-PRO-Ile-Pro	$\epsilon/\gamma$
Trp-Ser-PRO-GLN /	GLY-Val-Glu-Asn	$\gamma/\alpha$
Phe-Gln-Gln-GLN /	SER-PRO-Leu-Lys	I/F
Leu-Val-Ala-GLN /	GLY-PRO-Val-Asp	F/H
Glu-Gln-Glu-GLN /	GLY-PRO-Tyr-Asn	H/VPg
Leu-Asp-Ile-GLN /	GLY-PRO-Asn-Pro	VPg/22
Phe-Glu-PRO-GLN /	GLY-Ala-Leu-Glu	22/E
<u>HOST (?) PROTEASE CATALYZED CLEAVAGE SITES</u>		
Phe-Glu-Leu-GLN /	GLY-Asn-Ser-Thr	L/ $\delta$
Ser-Arg-Thr-TYR /	PRO-Thr-Leu-His	$\alpha/G$
Phe-Gln-Thr-GLN /	GLY-Ala-Ala-Val	G/I
Pro-Leu-Leu-ALA /	ASP-Gln-Asn-Thr	$\delta/\beta$
<u>CLEAVAGE</u> $\uparrow$ <u>SITE</u>		

**Figure 5.** Proteolytic cleavage sequences. The amino acid sequences surrounding all known proteolytic cleavage sites within the EMC polypeptide are tabulated. The sites are identified by the names of the final, stable processing products. See figure 4 to orient these cleavage sequences within precursor proteins A1, A, B, D1,  $\epsilon$ , C, D and H.

#### Viral Proteolytic (p22) Cleavage Sites

The EMC viral protease, p22, cleaves between proteins  $\epsilon/\gamma$ ,  $\gamma/\alpha$ , H/22, 22/E (14-16). By analogy with the homologous polio enzyme, p22 is also responsible for processing events between I/F, F/H and H/VPg (17). The deduced amino acid sequences surrounding known cleavage sites within the EMC polypeptide are tabulated in Figure 5. All cleavages attributed to p22 are located between gln-gly or gln-ser dipeptides. This result differs from the specificity of the polio protease, which exhibits exclusive preference for gln-gly sequences (17). The equivalent FMDV sites encompass a wider range of sequences, including glu-ser, glu-gly and gln-thr (21,25).

It is not known how the viral proteases recognise or select appropriate cleavage sequences. Both the polio and EMC polypeptides contain gln-gly pairs which do not function in proteolytic reactions (17). However, with EMC, the sequence shows that all sites known to be cleaved by p22 are also flanked by proximal proline residues (see Fig 5). There are six additional gln-gly and gln-ser pairs within the EMC polypeptide. None of these sites lies adjacent to prolines, and none are cleaved by p22. Two of these sequences can serve as processing sites (L/ $\delta$  and G/I), but they are cleaved by host-derived enzymes, not the viral protease. It remains to be determined, whether the proximity of prolines at p22-catalyzed sites indicates a primary sequence recognition requirement specific to this enzyme, or simply reflects a structural preference towards helix-breaking amino acids near the cleavage points.

### Host Proteolytic Cleavage Sites

Capsid precursor protein A1 is released from the nascent EMC polyprotein before the sequences encoding p22 have been translated (12-14). Proteins A and B, as cleavage products of A1, are also evident before the appearance of p22 (12). These observations suggest that release of A1, A or B is catalyzed by an enzyme(s) of cellular origin. Tabulation of amino acids surrounding these sites reveals few common sequences (Fig. 5, L/ $\delta$ ,  $\alpha$ /G and G/I). Two of the cleavages occur at gln-gly pairs (not flanked by prolines), but it is not obvious how these sites are specifically distinguished from others within the polyprotein. Our sequencing data give no indication as to the nature or origin of the amino-terminal blocking group found on mature protein  $\delta$  (1). It is completely unknown whether host-catalyzed proteolysis of L/ $\delta$  has any effect on this reaction.

The tyr-pro protein cleavage ( $\alpha$ /G) is at a site equivalent to the tyr-gly host-catalyzed cleavage within the polio polyprotein (17). Even though these polio and EMC sequences are both preceded by a threonine residue and have histidine in the plus 4 position, it is again not clear whether this observation implies a direct primary recognition requirement, or a selective preference towards helix-breaking amino acids.

The maturation processing of protein  $\epsilon$  to  $\delta$  plus  $\beta$  is the final proteolytic event within the EMC polyprotein. This cleavage is observed only during the last stages of virion morphogenesis and may play a role in stabilizing mature, RNA-containing particles (1). EMC and FMDV share considerable amino acid homology in the region surrounding this cleavage site (P-L-L-A/D-Q-N-T-E-E for EMC, and A-L-L-A/D-K-K-T-E-E for FMDV) (21,22), but neither sequence has as much in common with the equivalent polio site (P-M-L-N/S-P-N-I-E-A) (17). Although the viral or host agents responsible for these cleavages remain to be identified, the homology between EMC and FMDV suggests a similar proteolytic agent may be utilized by these two viruses, if not also by polio.

### Picornaviral Comparisons

Comprehensive comparisons among EMC, polio and FMDV have been undertaken and will be published separately. Base composition, codon preference and dot matrix analyses of protein and nucleic acid sequences will be included (ACP, manuscript in preparation).

### ACKNOWLEDGEMENTS

This research was supported by NIH grant AI 17331. We thank Roland

Rueckert and Paul Kaesberg for helpful discussions and critical reading of the manuscript.

\*To whom reprint requests should be sent

### REFERENCES

1. Rueckert, R. (1976) in *Comprehensive Virology*, Fraenkel-Conrat, H. and Wagner, R. Eds., Vol. 6, pp 131-213, Plenum Publishing Corp., New York.
2. Ahlquist, A. and Kaesberg, P. (1979) *Nucleic Acids Research* 7, 1195-1204.
3. Vartapetian, A., Drygin, Yu., Chumakov, K. and Bogdanova, A. (1980) *Nucleic Acids Research* 8, 3729-3742.
4. Rothberg, P., Harris, T., Nomoto, A. and Wimmer, E. (1978) *Proc. Nat. Acad. Sci. USA* 75, 4868-4872.
5. Ambrose, V. and Baltimore, D. (1978) *J. Biol. Chem.* 253, 5263-5266.
6. Black, D., Stephenson, P., Rowlands, D. and Brown, F. (1979) *Nucleic Acids Research* 6, 2381-2389.
7. Harris, T. and Brown, F. (1976) *J. Gen. Virol.* 33, 493-501.
8. Chumakov, K. and Agol, V. (1976) *Biochem. Biophys. Res. Com.* 71, 551-557.
9. Jacobson, M. and Baltimore, D. (1968) *Proc. Nat. Acad. Sci. USA* 61, 77-84.
10. Butterworth, B. and Korant, B. (1974) *J. Virol.* 14, 282-291.
11. Shih, D., Shih, C., Zimmern, D., Rueckert, R. and Kaesberg, P. (1979) *J. Virol.* 30, 472-480.
12. Shih, D. and Shih, C. (1981) *J. Virol.* 40, 942-945.
13. Pelham, H. (1978) *Eur. J. Biochem.* 85, 457-462.
14. Palmenberg, A., Pallansch, M. and Rueckert, R. (1979) *J. Virol.* 32, 770-778.
15. Palmenberg, A. and Rueckert, R. (1982) *J. Virol.* 41, 244-249.
16. Gorbalenya, A., Svitkin, Yu. and Agol, V. (1981) *Biochem. Biophys. Res. Comm.* 98, 952-960.
17. Kitamura, N., Semler, B., Rothberg, P., Larsen, G., Adler, C., Dorner, A., Emini, E., Hanecak, R., Lee, J., van der Werf, S., Anderson, C. and Wimmer, E. (1981) *Nature* 291, 547-553.
18. Racanelli, V. and Baltimore, D. (1981) *Proc. Nat. Acad. Sci. USA* 78, 4887-4891.
19. Robertson, B., Morgan, D., Moore, D., Grubman, M., Card, J., Fischer, T., Weddell, G., Dowbenko, D. and Yansura, D. (1983) *Virology* 126, 614-623.
20. Harris, T. (1979) *Nucleic Acids Research* 7, 1765-1785.
21. Boothroyd, J., Harris, T., Rowlands, D. and Lowe, P. (1982) *Gene* 17, 153-161.
22. Carroll, A., Rowlands, D. and Clarke, B. (1984) *Nucleic Acids Research*, manuscript submitted.
23. Beck, E., Forss, S., Strebel, K., Cattaneo, R. and Feil, G. (1983) *Nucleic Acids Research* 11, 7873-7885.
24. Dorner, A., Dorner, L., Larsen, G., Wimmer, E. and Anderson, C. (1982) *J. Virol.* 42, 1017-1028.
25. Forss, S. and Schaller, H. (1982) *Nucleic Acids Research* 10, 6441-6450.
26. Hall, L. and Rueckert, R. (1971) *Virology* 43, 152-165.
27. Pallansch, M. and Rueckert, R. (1981) *Methods in Enzymology* 78, 315-325.
28. Casadaban, M. and Cohen, S. (1980) *J. Mol. Biol.* 138, 179-207.
29. Sanger, F., Nicklen, S. and Coulson, A. (1977) *Proc. Nat. Acad. Sci. USA* 74, 5463-5476.
30. Messing, J. (1983) *Methods in Enzymology* 101, 20-78.



31. Smith, A. (1980) *Methods in Enzymology* 65, 560-580.
32. Maxam, A. and Gilbert, W. (1980) *Methods in Enzymology* 65, 499-560.
33. Garroff, H. and Ansorge, W. (1981) *Analyt. Biochem.* 115, 454-457.
34. Zimmer, D. and Kaesberg, P. (1978) *Proc. Nat. Acad. Sci. USA* 75, 4257-4261.
35. Staden, R. (1980) *Nucleic Acids Research* 8, 3673-3694.
36. Devereux, J., Haerberli, P. and Smithies, O. (1984) *Nucleic Acids Research* 12, 378-395.
37. Drake, N., Palmenberg, A., Ghosh, A., Omilianowski, D. and Kaesberg, P. (1982) *J. Virol.* 41, 726-729.
38. Butterworth, B. and Rueckert, R. (1972) *Virology* 50, 535-549.
39. Smith, A. (1973) *Eur. J. Biochem.* 33, 301-313.
40. Pallansch, M., Kew, O., Palmenberg, A., Golini, F., Wimmer, E. and Rueckert, R. (1980) *J. Virol* 35, 414-419.
41. Butterworth, B., Hall, L., Stoltzfus, C. and Rueckert, R. (1971) *Proc. Nat. Acad. Sci. USA* 68, 3083-3087.
42. Butterworth, B. and Rueckert, R. (1972) *J. Virol.* 9, 823-828.
43. Stoltzfus, M. and Rueckert, R. (1972) *J. Virol.* 10, 347-355.
44. Ziola, B. and Scraba, D. (1975) *Virology* 64, 228-235.
45. Ziola, B. and Scraba, D. (1976) *Virology* 71, 111-121.
46. Rueckert, R. and Wimmer, E. (1984) *J. Virol.*, in press.
47. Kazachkov, Yu., Chernovskaya, T., Siyanova, E., Svitkin, Yu., Ugarova, T. and Agol, V. (1982) *FEBS Letters* 141, 153-156.

Note added in proof: Since submission of this manuscript, portions of the 5' end of EMC, which were originally analyzed by indirect RNA sequencing methods, have been cloned into DNA and re-examined. The clones were sequenced (both strands) by Maxam and Gilbert techniques and the newly generated nucleotide assignments agree exactly with those presented in Figure 2 (ACP, manuscript in preparation).