Key words: IBV/coronavirus/nucleotide sequence

# Completion of the Sequence of the Genome of the Coronavirus Avian Infectious Bronchitis Virus

By M. E. G. BOURSNELL,\* T. D. K. BROWN, I. J. FOULDS, P. F. GREEN, F. M. TOMLEY AND M. M. BINNS

Houghton Poultry Research Station, Houghton, Huntingdon, Cambridgeshire PE17 2DA, U.K.

(Accepted 19 September 1986)

#### SUMMARY

The nucleotide sequence determination of the genome of the Beaudette strain of the coronavirus avian infectious bronchitis virus (IBV) has been completed. The complete sequence has been obtained from 17 overlapping cDNA clones, the 5'-most of which contains the leader sequence (as determined by direct sequencing of the genome) and the 3'-most of which contains the poly(A) tail. Approximately 8 kilobases at the 3' end of this sequence have already been published. These contain the sequences of mRNAs A to E within which are the genes for the spike, the membrane and the nucleocapsid polypeptides: the main structural components of the virion. The remainder of the sequence, equivalent to the 'unique' region of mRNA F, is some 20 kilobases in length and is thought to code for a polymerase or polymerases which are involved in the replication of the genome and the production of the subgenomic messenger RNAs. This sequence contains two large open reading frames, potentially coding for polypeptides of molecular weights 441000 and 300000. Unlike other large open reading frames in the virus, the 300000 open reading frame appears to have no subgenomic RNA associated with it which would allow it to be at the 5' end of an mRNA species. Because of this, and because of the characteristics of the sequence in the region immediately upstream of its start codon, other mechanisms of translation, such as ribosome slippage, must be postulated.

## INTRODUCTION

Avian infectious bronchitis virus (IBV) is the type species of the family Coronaviridae (Siddell et al., 1983a). Coronaviruses are enveloped, pleomorphic particles with a distinctive 'corona' of club-shaped surface projections, and a large single-stranded RNA genome of positive polarity (Siddell et al., 1983b). In infected cells, in addition to genome-sized RNA, a number of subgenomic RNAs can be detected which have a common 3' terminus, but extend for different lengths in the 5' direction, forming a nested set (Stern & Kennedy, 1980a, b; Leibowitz et al., 1981). In the case of IBV these are designated mRNAs A to F, mRNA A being the smallest and mRNA F being of genome length. In vitro translation studies have demonstrated that mRNAs A, C and E code for the nucleocapsid polypeptide, the membrane polypeptide and the precursor polypeptide to the spike or surface projection respectively (Stern & Sefton, 1984). These three polypeptides form the three known structural proteins of coronavirus virions (Cavanagh, 1981). Sequencing of cDNA clones derived from IBV genomic RNA has shown that, in the case of mRNAs A, C and E, only the 5' region of each mRNA which is not present in the next smallest mRNA is translated (Boursnell et al., 1985a, 1984; Binns et al., 1985b). This region is often referred to, for convenience, as the 'unique' region of the particular mRNA. For mRNAs B and D the situation is more complicated in that each mRNA has more than one open reading frame (ORF) and also has ORFs overlapping the next smallest mRNA (Boursnell & Brown, 1984; Boursnell et al., 1985b).

The genome of IBV is infectious (Lomniczi, 1977) indicating that it has a messenger function. There is also no evidence for a virion-associated RNA polymerase (Schochetman *et al.*, 1977).

On entry into the cell therefore the virion RNA probably codes for a polymerase, the gene for which must lie in the large 5' region of the genome, the 'unique' region of mRNA F, which does not contain the genes for the structural polypeptides. This polymerase would then be used to synthesize a negative-stranded template. The negative strand could then be used by another polymerase, or a modified form of the same polymerase, to produce the subgenomic mRNAs and virion RNA. Both the negative strand and two distinct polymerase activities have been detected in cells infected with the coronavirus mouse hepatitis virus (MHV) (Lai *et al.*, 1982; Brayton *et al.*, 1982). Translation of MHV virion RNA in reticulocyte lysates produced three structurally related polypeptides of molecular weights greater than 200000 (200K) (Leibowitz *et al.*, 1982).

In this paper we present the nucleotide sequence, obtained from cDNA clones, of the 'unique' region of mRNA F, the genome-sized mRNA. The sequence of approximately 8 kilobases from the 3' end of the genome, containing the genes for the major structural polypeptides, has already been published (Boursnell & Brown, 1984; Boursnell *et al.*, 1984, 1985*a*, *b*; Binns *et al.*, 1985*b*). The 20500 bases of sequence reported here complete the sequence of the IBV genome, which is, as far as we are aware, the first complete sequence of a coronavirus and the largest RNA virus sequenced to date.

#### METHODS

cDNA cloning. Seventeen cDNA clones covering the 3'-most 27569 kb of the genome have been obtained. These are shown in Fig. 1. They have been derived from RNA isolated from gradient-purified virus of the Beaudette strain (Beaudette & Hudson, 1937; Brown & Boursnell, 1984). cDNA has been obtained by three methods: oligo(dT) priming (Brown & Boursnell, 1984), priming with specific oligonucleotides (Boursnell *et al.*, 1984) and random priming with calf thymus DNA oligonucleotides (Binns *et al.*, 1985*a*). The Southern blotting technique was used to identify overlapping clones (Southern, 1975). Specific cDNA clones were identified using 'prime-cut' probes. These are made by synthesizing labelled DNA from selected M13 clones using the normal sequencing primer, cutting with a restriction enzyme, and eluting the labelled, single-stranded probe from denaturing acrylamide gels (Biggin *et al.*, 1984).

Subcloning for M13 sequencing. Random subclones of each cDNA clone were generated by sonication (Deininger, 1983) and subcloning into SmaI-cut, phosphatase-treated M13mp10 (Amersham). Bacterial colonies containing M13 with inserts were grown, transferred to nitrocellulose filters, and probed with nick-translated purified viral insert DNA from the cDNA clone. Single-stranded templates were prepared from M13 clones identified as viral in this way.

DNA sequencing. Sequencing was carried out by the dideoxy method (Sanger *et al.*, 1977; Bankier & Barrell, 1983). [ $\alpha$ -3<sup>5</sup>S]dATP was used in the sequencing reactions and the products were analysed on buffer gradient gels (Biggin *et al.*, 1983). Additional sequencing information was obtained by reverse sequencing (Hong, 1981). For regions containing compressions due to DNA secondary structure, sequencing samples were run on hot (80 °C) gels or gels containing 42% formamide. For some regions cytosine residues were modified by the method of Ambartsumyan & Mazo (1980) prior to separating on gels, to reduce GC base pairing. Deoxyinosine triphosphate (Bankier & Barrell, 1983) and deoxy-7-deazaguanosine triphosphate (Mizusawa *et al.*, 1986) were used in place of deoxyguanosine triphosphate in some cases, again to reduce GC base pairing. For sequencing directly from the viral RNA the method used was essentially as described by Caton *et al.* (1982).

Computer analysis of the sequence data. Sequence data were read directly into a BBC microcomputer using a sonic digitizer (Graf/Bar, Science Accessories Corporation) and data were analysed on a VAX 11/750 using the programs of Staden (1982a, b, 1984a, b). Comparisons with the National Biomedical Research Foundation (NBRF) protein identification resource was made using the programs SEARCH and FASTP (George et al., 1986; Lipman & Pearson, 1985) and SEQHP (Kanehisa, 1982).

#### RESULTS

#### Selection of cDNA clones

The majority of the cDNA clones which have been used to obtain the sequence of the 'unique' region of mRNA F were produced by a random priming method (Binns *et al.*, 1985*a*). Clone 182 was produced by priming with a specific oligonucleotide from existing sequence at the 5' end of mRNA D. Clone 227 was identified as coming from the 5' end of the genome by probing a random library with leader-specific probes. The randomly primed clones 217, 216, 204, 210, 205, 220 and 249 were mapped by identifying overlaps using Southern blotting. The nine clones were



Fig. 1. Diagram showing the positions of all the cDNA clones used in obtaining the nucleotide sequence. The squares at the end of some of the clones show the positions of oligonucleotide primers used to prime synthesis of cDNA for adjacent clones. Above the clones are shown mRNAs A to F.

not contiguous but formed four blocks. cDNA clones in the region of the three remaining gaps were obtained using specific oligonucleotide primers. Clones spanning the gaps were identified using either 'prime-cut' probes (Biggin *et al.*, 1984) made from M13 subclones of cDNA clones on either side of the gap or by using Southern blotting. Five clones, 256, 263, BP3, BP5 and BP8 were identified in this way and the overlaps confirmed by sequencing. Fig. 1 shows the positions of all the cDNA clones used in obtaining the complete sequence of the virus, and the positions of the oligonucleotide primers.

# DNA sequencing

Fourteen cDNA clones have been sequenced to obtain the complete sequence of the 'unique' region of mRNA F, the genome-sized messenger RNA. The 20500 bases of sequence presented here stretch from the 5' end of the genome to an arbitrary position 190 bases 3'-wards of the end of the body of mRNA E. The 39 nucleotides at the very 5' end of the genome have not been obtained in cDNA clones from the Beaudette strain, and the sequence here is derived from Maxam & Gilbert (1980) sequencing of primer-extended products from Beaudette virion RNA (Brown *et al.*, 1986). Fig. 2 shows the DNA sequence obtained from the cDNA clones, with a translation in single-letter amino acid code of the main ORFs.

#### Sequence analysis

Fig. 3 shows the positions of ORFs in this region. Most of the sequence encodes two very large ORFs which could code for polypeptides of predicted molecular weights 441K and 300K. These two large ORFs have been designated F1 and F2.

The first large ORF, F1, is not the first ORF to occur after the homology region. At position 131 there is an AUG codon followed by a small ORF which could code for a polypeptide of 11 amino acids. This AUG is the first initiation codon to occur on the genome. The second initiation codon is at the start of F1. Both the large ORFs have a codon usage (Staden & McLachlan, 1982) very similar to that of the genes for the structural polypeptides S, M and N. The small ORF also appears to have the same codon usage, insofar as that is significant for such a short sequence. After the end of the small ORF the reading frame is open, in the other two possible frames, for a further 232 or 73 bases but the codon usage of the predicted amino acids for these sections of ORF is not similar to that previously found for IBV. The sequence context around the first AUG codon is not similar to that used by most eukaryotic mRNAs (Kozak, 1983) in that it has a pyrimidine at position -3. The context around the second AUG on the other hand has a purine at -3, in addition to a C at positions -1 and -4, both of which mean that it conforms well to the consensus for functional initiation codons.

| 1    | ACTTAAGATAGATATATATATATATATATATATATATATA   | 100  |
|------|--|------|
| 101  | M A P G H L S G F C Y *<br>GTTCCATTGCAGTGCACTTTAGTGCCCTGGCAGCCTGGCCACCTGTCAGGTTTTTGTTATTAAAATCTTATTGTTGCTGGTATCACTGCTTGTTTTG   | 200  |
| 201  | CCGTGTCTCACTTTATACATCTGTTGCTTGGGCTACCTAGTGTCCAGCGTCCTACGGCGTCGTGGCTGGTTCGAGTGCGAGGAACCTCTGGTTCATCTA  | 300  |
| 301  | GCGGTAGGEGGGTGTGTGGAAGTAGCACTTCAGACGTAECGGTTCTGTTGTGTAAATACGGGGTCACCTECCCCACATACCTCTAAGGGETTTTGAGC   | 400  |
| 401  | CTAGEGTTEGGETAEGTTETEGEATAAGGTEGGETATAEGAEGTTTGTAGGGGGTAGTGECAAACAACCECTGAGGTGAEAGGTTETEGTGGTGTTTAGT   | 500  |
| 501  | M A S S L K Q G V S P K P R D V I L V S K D I P<br>GAGCAGACATAGACAGGGACAACATGGCTTCAAGCCTAAAACAGGGGAGTATCTCCCAAACCACGGGATGTCATTCTTGTGTCCAAAGACATCCCT                          | 600  |
| 601  | E Q L C D A L F F Y T S H N P K D Y A D A F A V R Q K F D R S L Q T<br>GAACAACTITGTGACGCTITGTTTICFATACGTCACATAACCCTAAGGATTACGCTGATGCTTITGCAGTTAGGCAGAAGTTIGACCGTAGTCTCCAGA   | 700  |
| 701  | G K Q F K F E T V C G L F L L K G V D K I T P G V P A K V L K A T<br>CTGGGAAACAGTTCAAATTTGAAACTGTGTGTGTGTCTCTCTC   | 800  |
| 801  | SKLADLEDIFGVSPLARKYRELLKTACQWSLTV<br>TTCTAAGTTGGCAGATTTAGAAGACATCTTTGGTGTCTCTCTTTAGCGCGGAAGTACCGTGAATTGTTGAAAACAGCGTGTCAGTGGTCTCTTACTGTA                                     | 900  |
| 901  | E A L D V R A Q T L D E I F D P T E I L W L Q V A A K I H V S S M A<br>GAAGCACTGGATGTTCGTGCACAAAACTCTCGATGAAATTTTTGACCCCACTGAAATACTTTGGCTTCAGGTGGCTGCAAAAATTCATGTTTCATCTATGG | 1000 |
| 1001 | M R R L V G E V T A K V M D A L G S N L S A L F Q I V K Q Q I A R<br>CAATGCGCAGGCTTGTTGGAGAAGTAACTGCAAAGTCATGGATGCTCTGGGCTCAAACTGGGATGCTCTTGTCAAATTGTTAAACAACAAATAGCCAG      | 1100 |
| 1101 | IFQKALAIFENVNELPQRIAALKMAFAKCARSI<br>AATCTTTCAAAAGGCACTGGCTATTTTGAGAATGTGGATGGA  | 1200 |
| 1201 | T V V V E R T L V V K E F A G T C L A S I N G A V A K F F E E L P<br>ACTGTIGTGGTIGTGAAAGAACTCTAGTIGTTAAAGAGTTGCCAAGAACTTGTCTTGCAAGAGTTGCGCAGAAAATTCTTTGAAGAGTTGC             | 1300 |
| 1301 | N G F M G S K I F T T L A F F K E A A V R V V E N I P N A P R G T<br>CAAACGGCTTCATGGGTTCTAGATTTTCACAACACTTGCCTTCTTTAAAGAGGCAGCTGTGGAGAGTTGTGGGAGAACATACCAAATGCACCGAGAGGTAC   | 1400 |
| 1401 | K G F E V V G N A K G T Q V V V R G M R N D L T L L D Q K A D I P<br>TAAGGGATTTGAAGTTGTTGGCAATGCCAAAAGGCACACAGGTAGTTGTGGCGCGCATGCCGAAATGACTTAACATTGCTTGACCAAAAAGCTGATATTCCT  | 1500 |
| 1501 | V E P E G W S A I L D G H L C Y V F R S G D R F Y A A P L S G N F A GTTGAACCAGAAGGTTGGTCTGCAATTITGGATGGACATCTTGCTATGTCTTTAGGAGGTGGTGATCGCTTTTATGCTGCACCTCTTTCAGGAAATTITG     | 1600 |
| 1601 | L S D V H C C E R V V C L S D G V T P E I N D G L I L A A I Y S S<br>CTITGAGTGATGTTCATTGCTGTGGGCGTGTAGTCTGCTGCTGTGTGTG   | 1700 |
| 1701 | FSVSELVTALKKGEPFKFLGHKFVYAKDAAVSF<br>TTTTAGTGTCTCTGAGCTTGTAACAGGCTCTAAAAAGGGTGAACCATTCAAGTTCTTGGGCCATAAATTCGTGTATGCGAAGGATGCAGCAGTGTCTTTT                                    | 1800 |
| 1801 | T L A K A A T I A D V L R L F Q S A R V I A E D V W S S F T E K S F<br>ACTITAGEGAAGGETGEEACTATTGEAGATGTEITGAGEETGTTEAATEAGETEGTGTGATAGEAGAAGATGTTTGGTEITCATTAETGAAAAGTETT    | 1900 |
| 1901 | E F W K L A Y G K V R N L E E F V K T Y V C K A Q M S I V I L A A<br>TTGAATTCTGGAAGCTTGCATATGGAAAAGTGCGCAACCTTGAAGAATTTGTGAAGACCTATGTTTGTAAGGCTCAAATGTCGATTGTGATTCTAGCAGC    | 2000 |
| 2001 | VLGEDIWHLVSQVIYKLGVLFTKVVDFCDKHWK<br>AGTGCTTGGAGAGGACATTTGGCATCTTGTCTCACAAGTCATCTAAAGTCGTTGACTTGTGACAAACACTGGAAA   | 2100 |
| 2101 | G F C V Q L K R A K L I V T E T F C V L K G V A Q H C F Q L L L D A<br>GGTTTTTGTGTACAGTTGAAAAGAGCTAAGCTCATTGTCACCGAAACCTTCTGTGTTTAAAAGGAGTTGCACAGCATTGTTTCAACTGCTGCTAGATG    | 2200 |
| 2201 | I H S L Y K S F K K C A L G R I H G D L L F W K G G V H K I V Q D<br>CAATACACTCTITGTACAAGAGGTTTTAAGAAGTGTGCACTTGGTAGAATCCATGGGAGATTTGCTCTTGGAAAGGAGGTGTGCATAAAATTGTTCAAGA    | 2300 |
| 2301 | G D E I W F D A I D S V D V E D L G V V Q E K S I D F E V C D D V<br>TGGCGATGAAATATGGTTTGACGCCATTGATAGTGTTGATGTTGAAGATCGAGGTGTGGTGGGGGGGG                                    | 2400 |
| 2401 | T L P E N Q P G H M V Q I E D D G K N Y M F F R F K K D E N I Y Y T<br>ACACTTCCAGAAAACCAACCTGGTCATATGGTTCAAAAAGGGATGATGATGATGATGATGAGAACATTTATTATA                           | 2500 |

|               | Coronavirus IBV sequence completed  | 61   |
|---------------|---|------|
| 2501          | PMSQLGAINVVCKAGGKTVTFGETTVQEIPPPD<br>CACCAATGTCTCAACTTGGTGCTATAATGTGGTTTGCAAAGCAGGCGGTAAGACTGTCACCTTTGGAGAAACTACAGTACAAGAGATACCACCACCTGA  | 2600 |
| 2601          | VVPIKVSIECCGEP₩NTIFKKAYKEPIEVDTDL<br>TGTCGTGCCTATTAAGGTTAGCATAGAATGTTGTGGTGGACCATGGAATACGAACTTCAAGAAGGGCTTATAAAGGGCCTATAGAAGTAGAAGTAGAACAGACCTC                                   | 2700 |
| 2 <b>7</b> 01 | T V E Q L L S V I Y E K M C D D L K L F P E A P E P P F E N V A L<br>ACAGTAGAACAATTGCTCTCTGTGATCTATGAGAAAATGTGTGACGACCTTAAATTGTTTCCAGAGGCACCAGAGCCTCCACCATTGAGAATGTCGCAC          | 2800 |
| 2801          | V D K N G K D L D C I K S C H L I Y R D Y E S D D D I E E E D A E<br>TTGTTGATAAGAACGGTAAAGATTTGGATTGTATAAAATCTTGCCATTTGATCTATCGTGACTATGAGAGCGATGATGACATCGAGGAGGAGATGCTGA          | 2900 |
| 2901          | E C D T D S G E A E E C D T N S E C E E E D E D T K V L A L I Q D<br>GGAGTGTGACACAGACTCAGGTGGAGGTGGGGGGTGGGCACTAATTCAGAAGAGGAGGAGGAGGATGACGATACTAAAGGGTTGGCTCTTATACAAGAC          | 3000 |
| 3001          | PASIKYPLPLDEDYSVYNGCIVHKDALDVVNLPS<br>CCGGCAAGTATTAAATACCCTCTGCCTCTTGATGAAGATTATAGCGTCTATAATGGATGTATGT  | 3100 |
| 3101          | G E E T F V V N N C F E G A V K P L P Q K V V D V L G D W G E A V<br>CTGGTGAAGAAACTTTTGTGTCAATAACTGTTTTGAGGGAGCTGTTAAACCACTTCCACAGAAGGTAGTTGATGTCTTGGTGACTGGGGAGAGGCTGT           | 3200 |
| 3201          | DAQEQLCQQEPLQHTFEEPVENSTGSSKTMTEQ<br>TGATGCGCAAGAACTGTGTCAACAAGAGCCTCTGCAACATACCTTTGAAGAACCAGTCGAAAATTCTACTGGTAGTTCTAAGACAATGACTGAACAA  | 3300 |
| 3301          | VVVEOQELPVVEQDQDVVVYTPTDLEVAKETAEE<br>GTCGTTGTAGAAGATCAAGAACTACCTGTTGTTGAAGAACAAGAAACAAGAAGAAGAAGAAGAAGAAGAAGAA   | 3400 |
| 3401          | V D E F I L I F A V P K E E V V S Q K D G A Q I K Q E P I Q V V K<br>AGGITGATGAGITTATICTCATITTIGCTGTTCCTAAAGAAGAAGAAGTTGTGTCCCAGAAAGATGGGGGCACAGATTAAACAAGAGCCTATTCAAGTIGTTAA     | 3500 |
| 3501          | PQREKKAKKFKVKPATCEKPKFLEYKTCVGDLT<br>ACCACAACGTGAGAAGAAGGCTAAAAAGTTCAAAGTTAAACCAGCCACATGTGGAGAAACCTAAATTTITGGAGTATAAAACATGTGGGGGGATTTGACT   | 3600 |
| 3601          | VVIAKALDEFKEFCIVNAANEHMTHGSGVAKAIA<br>GTTGTAATTGCCAAAGCATTGGATGAGTTTAAAGAGTTCTGCATTGTAATGCTGCAAAGCATATGACTCATGGTAGTGGCGTTGCAAAGGCAATTG  | 3700 |
| 3701          | DFCGLDFVEYCEDYVKKHGPQQRLVTPSFVKGI<br>CAGACTTTTGTGGACTGGATTTTGTGAGGAGGACTATGTTAAGAAACATGGGCCACAACAGAGACTTGTTACACCTTCGTTTGTCAAAGGCAT  | 3800 |
| 3801          | Q C V N N V V G P R H G D N N L H E K L V A A Y K N V L V D G V V<br>TCAATGTGTGAATAATGTTGTAGGACCCCCCCCATGGAGAACAACTTGCATGAGAAGCTTGTTGCTGCCTACAAGAATGTGCTTGTAGATGGCGTAGTC          | 3900 |
| 3901          | NYVVPVLSLGIFGVDFKMSIDAMREAFEGCTIRV<br>AATTATGTTGTGCCAGTTCTTTCATTAGGAATTTTTGGTGTGGAGATTTTGAAGGTTGCACCATACGCG   | 4000 |
| 4001          | LLFSLSQEHIDYFDVTCKQKTIYLTEDGVKYRS<br>TTCTTTGTTTCCTGAGCCAAGAACACATCGATCGTTATTCCGAGGATCGTGTTAAATACCGCTC   | 4100 |
| 4101          | IVLKPGDSLGQFGQVYAKNKIVFTADDVEDKEI<br>CATTGTTCTAAAAACCTGGTGACTTGGGCCAATTTGGACAGGTAGGT  | 4200 |
| -+201         | LYVPTTDKSILEYYGLDAQKYVIYLQTLAQKWNV<br>CTCTACGTCCCCACGACTGATAAAAGCATTCTTGAATACTATGGTTTAGATGCGCAAAAGGTATGTAATATATTTGCAAAAGGCTTGCGCAGAAATGGAATG                                      | 4300 |
| 4301          | QYRDNFLILEWRDGNCWISSAIVLLQAAKIRFK<br>TCCAATATAGGGACAATTTTCTTATACTAGAGTGGCGCGATGGAAATTGTGTGGATAGTTCAGCAATAGTTCTCCTTCAAGCTGCTAAAATTAGGTTTAA   | 4400 |
| 4401          | G F L T E A W A K L L G G D P T D F V A W C Y A S C T A K V G D F<br>AGGITITCTAACAGAAGCGTGGGCTAAACTGTTAGGTGGGGGGAGATCCTACAGGACTITGTTGCCTGGTGTTATGCAAGTTGTACTGCTAAAGTAGGTGGTGATTTC | 4500 |
| 4501          | S D A N W L L A N L A E H F D A D Y T N A F L K K R V S C N C G I K<br>TCAGATGCTAATTGGCTTTTAGCGAATTAGCAGAACATTTTGACGCAGATTACACAAATGCGTTTCTTAAGAAGCGCGTTTCGTGTAACTGTGGTATTA        | 4600 |
| 4601          | SYELRGLEACIQPVRATNLLHFKTQYSNCPTCG<br>AGAGCTATGAGCTTAGAGGCCTTGAAGCTTGTATTCAGCCGGCAACTAATCTGCCCAACCTGTGG  | 4700 |
| 4701          | ANNT DEVIEASLPYLLLFAT DGPAT VDC DEDAV<br>CGCAAATAATACGGATGAAGTAATAGAAGCTTCGTTACCGTACCTGTTGCTTATTGCTACTGATGGTCCTGCTACAGTTGATTGTGATGAAGATGCTGTG                                     | 4800 |
|               |   |      |

G T V V F V G S T N S G H C Y T Q A A G Q A F D N L A K D R K F G K 4801 GGGACTGTCGTGTTGTTGGTTCTACTACTAGTGGGCCATTGTTATACACAAGCTGCGGGCAAGCTTTGATAATCTTGCTAAAGATAGAAAATTTGGAA 4900

| 4901 | K S P Y I T A M Y T R F A F K N E T S L P V A K Q S K G K S K S V<br>AGAAGTEGECETTACATTACTGEAATGTATACGEGATTEGECTITTAAGAATGAAACECTETTTGECTAGTEGGAAGAGEAAGGETAAGTEGGT        | 5000 |
|------|--|------|
| 5001 | K E D V S N L A T S S K A S F D N L T D F E Q W Y D S N I Y E S L<br>AAAGGAAGATGTTTCTAACCTTGCTACTAGTTCTAAGGCCAGTTTGATAATCTTACTGACTTCGAACAGTGGTATGATAGTAACATCTATGAAAGTCTT   | 5100 |
| 5101 | K V Q E S P D N F D K Y V S F T T K E D S K L P L T L K V R G I K S<br>AAAGTGCAGGAATCACCTGATAACTTTGATAAATATGTGTCATTCAAAAGGAAGATTCTAAGTTGCCATTGACACTTAAGTTAGAGGTATTAAAT     | 5200 |
| 5201 | V V D F R S K D G F I Y K L T P D T D E N S K A P V Y Y P V L D A<br>CAGTIGITGACTITAGATCGAAGGATGGITITATTIATAAGITAACACCIGATACTGATAAAAAGCAACCAGTCTACTACCAGTCTTGGACGC         | 5300 |
| 5301 | ISLKAIWVEGNANFVVGHPNYYSKSLHIPTFWE<br>TATTAGTCTTAAGGCAATATGGGTGGAAGGTAATGCTAACTTTGTTGTGGTCATCCAAATTATTATAGTAAGTCTCTTCATATTCCTACTTTTGGGAA                                    | 5400 |
| 5401 | N A E N F V K M G D K I G G V T M G L W R A E H L N K P N L E R I F<br>AATGCTGAGAATTTTGTTAAAATGGGGGGTAAAAATTGGGGGGTGAAGAGCGTGGGGGGGG                                       | 5500 |
| 5501 | NIAKKAIVGSSVVTTQCGKLIGKAATTFIADKVG<br>TCAACATTGCTAAGAAAGCCATTGTTGGGGTGATGTGGTGAAAGCGGGGAAATTAATAGGTAAAGCAAGC   | 5600 |
| 5601 | G G V V R N I T D S I K G L C G I T R G H F E R K M S P Q F L K T<br>TGGTGGTGGTTGGCAATATTACAGATAGGGTTATAGGGGTCTTTGGGAATTACAGGAGGGGCATTTTGGAAAGAAA                          | 5700 |
| 5701 | LMFFLFYFLKASVKSVVASYKTVLCKVVLATLLI<br>CTTATGTTCTTTTATTCTATTCTTGAAGGCTAGTGTTAAGAGGTGGTAGTAGCTAGC  | 5800 |
| 5801 | V W F V Y T S N P V M F T G I R V L D F L F E G S L C G P Y K D Y<br>TAGTITGGTTGCTACACAAGTAACCCAGTAATGTTTACAGGAATACGTGTGTTGTATATAGAGGTTCTTTGTGTGGTCCTTATAAAGACTA           | 5900 |
| 5901 | G K D S F D V L R Y C A D D F I C R V C L H D K D S L H L Y K H A<br>TGGTAAAGATTCTTTTGATGTGTTACGATATTGTGGGAGATGATTTTATTTGTCGTGTGTGT  | 6000 |
| 6001 | YSVEQVYKDAASGFIFNWNWLYLVFLILFVKPVA<br>TATAGTGTAGAGCAGGTCTATAAAGATGCAGCTTCTGGTTTATATGGTATTGGTCTTTCTAAATATTATTTGTTAAACCAGTGG   | 6100 |
| 6101 | G F V I I C Y C V K Y L V L N S T V L Q T G V C F L D W F V Q T V<br>CAGGTTTTGTTATTATTTGCTATTGTGTTAAGTATTTGGTATTGGATTGAATTGAATTGGTTTGTGCTGCTGCAAACTGGTGTTTGTT              | 6200 |
| 6201 | FSHFNFMGAGFYFWLFYKIYIQVHHILYCKDVT<br>TTTTAGTCACTTTAATTITATGGGAGCAGGGTTTTATTTTTTTTTT  | 6300 |
| 6301 | C E V C K R V A R S N R Q E V S V V V G G R K Q I V H V Y T N S G Y<br>TGTGAAGTGTGCAAAAGGGTTGCACGCAACAGGCAAGAGGGTTAGCGTGGTTGTTGGTGGACGCAAGCAGGATAGTGCATGTTTACACTAACTCTGGCT | 6400 |
| 6401 | NFCKRHNWYCRNCDDYGHQNTFMSPEVAGELSE<br>ATAACTTTTGTAAGAGGACATAGTTATGGTATTGTGAGAAATTGGTCACCAAAATACATTTATGTCTCCTGAAGTTGCTGGCGAGGCTCTCTGA  | 6500 |
| 6501 | KLKRHVKPTAYAYHVVDEACLVDDFVNLKYKAA<br>AAAGCTTAAGCGCCATGTTAAACCTACAGCATACCCTTACCACGTTGTGGATGGCGCATGCTTAGTGATTTTGTCAATTTAAAATATAAAGCTGCA                                      | 6600 |
| 6601 | T P G K D S A S S A V K C F S V T D F L K K A V F L K E A L K C E Q<br>ACTCCTGGTAAGGATAGTGCATCTTCAGCTGTTAAGTGTTTCGAGTGTTACAGATAGGAAAGCTGTTTTTCTTAAGGAAGCACTGAAATGTGAAC     | 6700 |
| 6701 | ISNDGFIVCNTQSAHALEEAKNAAIYYAQYLCK<br>AAATATCTAATGATGGTTTTATAGTGTGTAATACACAGAGTGCTCATGCATAGCGCAAGGAAGCAAAGAATGCAGCCATCTATTATGCGCAATATCTGTGTAA                               | 6800 |
| 6801 | PILILDQALYEQLVVEPVSKSVIDKVCSILSSI<br>GCCAATACTTATACTIGACCAGGCACTTATGAGCAATTAGTAGAGAGCCTGTGTCTAAGAGTGTTATAGATAAAGTGTGTAGCATTTTGTCTAGTATA                                    | 6900 |
| 6901 | ISVDTAALNYKAGTLRDALLSITKDEEAVDMAIF<br>ATATCTGTAGATACTGCAGCTTTAAATTATAAGGCAGGCA   | 7000 |
| 7001 | C H N H D V D Y T G D G F T N V I P S Y G I D T G K L T P R D R G<br>TCTGTCATAATCATGATGTGGGATTACACTGGTGATGGTTTACCAATGTGATACGGCCATGTGACAGCTCGGCAAGTTAACACCTCGTGATAGAGG      | 7100 |
| 7101 | FLINADASIANLRVKNAPPVVWKFSELIKLSDS<br>GTTTTTGATAAATGCAGATGCTTCTATTGCTAACTTAAGAGTTAAAAATGCTCCGCCGGTAGTATGGAAGTTTTCTGAGCTTATTAAGTTGTCTGACAGT                                  | 7200 |
| 7201 | CLKYLISA TVKSGVRFFITKSGAKQVVIA CHTQKL<br>TGTCTTAAATATTTAATTTCGGCTACTGTTTAAGGCAGGTGTTCGTTC  | 7300 |

| 7301 | LVEKKAGGIVSGTFKCFKSYFKWLLIFYILFTA<br>TGTTAGTAGAGAAAAAGGCAGGTGGTATTGTTAGCGCGCACCTTTAAGTGTTTTAAGAGTTATTTTAAATGGCTCTTTGCATCTTTTACATACTTTTTACAGC                               | 7400 |
|------|--|------|
| 7401 | CCSGYYYMEVSKSFVHPMYDVNSTLHVEGFKVI<br>ATGTTGTTCGGGTTATTATGGAGGTGAGTAAAAGTTTTGTTCACCCCATGTATGATGTAGACTCCACACTGCATGTTGAAGGTTTTAAAGTTATA                                       | 7500 |
| 7501 | DKGVLREIVPEDTCFSNKFVNFDAFWGRPYDNSR<br>GATAAAGGTGTTCTTAGGGGAAATTGTACCAGAAGATACATGTTCCTAATAAATTTGTAGTGCTTTTTGGGGCAGACCATATGATAATAGTA   | 7600 |
| 7601 | N C P I V T A V I D G D G T V A T G V P G F V S W V M D G V M F I<br>GAAACTGTCCAATTGTCACAGCTGTTATAGATGGTGTGGGGGGGG   | 7700 |
| 7701 | H M T Q T E R K P W Y I P T W F N R E I V G Y T Q D S I I T E G S<br>ACATATGACACAGACTGAGAGAAAACCGTGGTACATTCCTACTGAGGTTAATAGAGAAATTGTCGGTTACACTCAGGATTCAATTATTACTGAGGGTAGT  | 7800 |
| 7801 | FYTSIALFSARCLYLTASNTPQLYCFNG DNDAPG<br>TTTTATACATCTATAGCGTTATTTTCCGCTAGGTGTTTATTTA   | 7900 |
| 7901 | A L P F G S I I P H R V Y F Q P N G V R L I V P Q Q I L H T P Y V<br>GGGCTTTGCCATTTGGTAGTATTATTCCTCATAGAGTTTATTTCCAACCCAATGGTGTTAGGTTCAAGTACAAAAAACTGCACAACAACAACCCTACGT   | 8000 |
| 8001 | VKFVSDSYCRGSVCEYTRPGYCVSLNPQWVLFN<br>AGTAAAGTITGTATCAGACAGCTATTGTAGGGGTAGGGGTGGGGT   | 8100 |
| 8101 | DEYTSKPGVFCGSTVRELMFSMVSTFFTGVNPNI<br>GACGAATACACAAGTAAACCCGGTGTTTTCTGTGGGTCCTACTGTTAGGAACTTATGTTTGGTTAGGTAGG  | 8200 |
| 8201 | Y M Q L A T M F L I L V V V V L I F A M V I K F Q G V F K A Y A T<br>TCTATATGCAATTAGCAACTATGTTTTTAATACTAGTTGTTGTATTAATCTTTGCAATGGTTATAAAGTTTCAAGGTGTTTTTAAAGCTTATGCAAC     | 8300 |
| 8301 | TVFITMLVWVINAFILCVHSYNSVLAVILLVLY<br>CACTGITTTTATAACAATGITAGTITGGGTAAT!AACGCATITATTIGTGTGTGTACAACAGGTGTTTTAGCTGTTATATTACTAGTACTCTAT  | 8400 |
| 8401 | CYASLVTSRNTVIIMHCWLVFTFGLIVPTWLACC<br>TGCTATGCGTCATTGGTTACAAGTCGCAATACTGTTATAATAGTACCCACATGGTTGGCTTGTT<br>TGCTATGCGTCATTGGTTACAAGTCGCAATACTGTTATAATAGTACCCACATGGTTGGCTTGTT | 8500 |
| 8501 | YLGFIIYMYTPLFLWCYGTTKNTRKLYDGNEFV<br>GCTACCTGGGATTTATTATTATGTATACACCGTTGTTTTTATGGTGTTATGGTACTACAAAAAAACACTCGTAAGCTGTATGATGGCAATGAGTTTGT                                    | 8600 |
| 8601 | G N Y D L A A K S T F V I R G S E F V K L T N E I G D K F E A Y L<br>TGGTAATTATGATCITGCTGCGAAGAGCACTTITGTTATTCGCGGCTCTGAATTGTTAAGCTTACTAATGAGATAGGTGATAAATTTGAGGCCTACCTT   | 8700 |
| 8701 | SAYARLKYYSG TG SE QD YLQACRAWLAYALDQYR<br>TCAGCGTATGCTAGATTAAAGTACTATTCAGGCACTGGCAGGCA   | 8800 |
| 8801 | N S G V E I V Y T P P R Y S I G V S R L Q S G F K K L V S P S S A<br>GAAATAGTGGTGTGGAAATTGTTTATACTCCGCCACGTTACTATTGGTGTTAGGAATCTGGTTTTAAGAAACTGGTTTCTCCTAGTAGTGC           | 8900 |
| 8901 | V E K C I V S V S Y R G N N L N G L W L G D T I Y C P R H V L G K<br>TGTTGAAAAGTGCATTGTTAGTGTCTCTTATAGAGGTAATAATCTTAATGGACTGTGGCTAGGTGACACTATCTACTGTCCTCGTCATGTATTGGGTAAG  | 9000 |
| 9001 | FSGDQWNDVLNLANNHEFEVTTQHGVTLNVVSRR<br>TTTTCAGGTGACCAATGGAATGATGTACTTAATCTTGCTAATAATCATGAGTTTGAAGTTACAACATGGTGTTACTTTGAATGTTGTCAGTAGGG                                      | 9100 |
| 9101 | LKGAVLILQTAVANAETPKYKFIKANCGDSFTI<br>GTITAAAAGGTGCAGTITTAATITTACAAACTGCTGTTGCTAATGCTGAAAGTCCAAAGTATAAGTTTATAAAGCTAATTGTGGTGATAGTITCACTAT                                   | 9200 |
| 9201 | A C A Y G G T V V G L Y P V T M R S N G T I R A S F L A G A C G S<br>AGCTTGTGCTTATGGTGGTACAGTTGTAGGACTCTACCCTGTTACTATGGGGTCTAATGGTACTATTAGAGCATCTTTTCTTGEGGGAGCCTGTGGTTCA  | 9300 |
| 9301 | VGFNIEKGVVNFFYMHHLELPNALHTGTDLMGEF<br>GTGGTTTTAATATAGAAAAGGGTGTAGTTAATTTCTTTTATATGCACCATCTGAGTTACCACTGCATTACACACTGGAACTGACCTAATGGGTGAAT                                    | 9400 |
| 9401 | Y G G Y V D E E V A Q R V P P D N L V T N N I V A W L Y A A I I S<br>TCTATGGTGGTTATGTTGATGAAGAGGTTGCACAAAGAGTGCCACCAGATAATTTAGTTACTAACAATATTGTAGCATGGCTCTATGCGGCAATTATTAG  | 9500 |
| 9501 | VKESSFSLPKWLESTTVSVDDYNKWAGDNGFTP<br>TGTTAAGGAGAGTAGTTICTCGCTGCCTAAATGGTGGGAGAGTACTACTGTTAGTGTGTGATGATTATAATAAGTGGGCTGGTGACAATGGTTTTACACCA                                 | 9600 |
| 9601 | FSTSTAITKLSAITGVDVCKLLRTIMVKNSQWGG<br>TTTTCTACTAGTACCGCTATTAACTAGATGCTATAACTGGGAGTGGATGGTGTGTGT  | 9700 |

63

| 9701  | D P I L G Q Y N F E D E L T P E S V F N Q I G G V R L Q S S F V R<br>GTGACCCCATTITAGGGCAATATAATTITGAAGATGAACTGACGGGGTGTTAATTAATCAGATGGGGGTGTTAGATTACAATCYTCTITIGTAAG        | 9800  |
|-------|---|-------|
| 9801  | KATSWFWSRCVLACFLFVLCAIVLFTAVPLKFY<br>ARAAGCTACATCTTGGTTTGGAGTAGATGTGGTGTGGTGCTGTTTGTGTTGTGTTGTGTTGTGTTGGGGGG  | 9900  |
| 9901  | VYAAVILLMAVLFISFTVKHVMAYMDTFLLPTLI<br>GTATATGCAGCTGTTATTTGGTTAATGGCTGTAATGGCATATTGGGATACTTTTGCTAATGCCAACATTGA   | 10000 |
| 10001 | T V I I G V C A E V P F I Y N T L I S Q V V I F L S Q W Y D P V V<br>TTACAGTTATTATTGGAGTTGTGCTGAAGTGCCTTTCATCTACAATACTCTAATTAGTCAAGTTGTTATTTTCTTAAGTCAATGGTATGACCCAGTAGT    | 10100 |
| 10101 | FDTMVPWMFLPLVLYTAFKCVQGCYMNSFNTSL<br>CTITGATACTATGGTACCATGGATGTTCTGCCACTAGTGTTGTATACTGCTTTTAAGTGTGCTACAAGGTTGCTATATGAATTCTTTCAATACTTCTTTG                                   | 10200 |
| 10201 | LMLYQFVKLGFVIYTSSNTLTAYTCGCATACACAGAAGGTAATTGGGAGTTATTCTTCGAGT  | 10300 |
| 10301 | VHTTVLANVSSNSLIGLFVFKCAKWMLYYCNAT<br>TGGTGCACACTACTGTGTTGGTTGGTTGGTTGTTTGGTTTTGAGTGTGCTAAATGGATGTTGTATTATTGTAATGCAAC  | 10400 |
| 10401 | YLNNYVLMAVMVVNCIGULCTCYFGLYWWVNKVF<br>ATACTTAAACAATTATGTACTAATGGCAGTTAACTGCATTGGCTGCACTTGTTACTTTGGGTTGTATTGGTGGGTTAATAAGGTTTTT  | 10500 |
| 10501 | G L T L G K Y N F K V S V D Q Y R Y M C L H K I N P P K T V W E V F<br>GGTTTAACCTTAGGTAAATACAATTTTAAAGTTTCAGTAGATCAATATAGGTATATGTGTTTGCACAAGATAAACCCACCTAAAACTGTGTGGGAAGTCT | 10600 |
| 10601 | S T N I L I Q G I G G D R V L P I A T V Q A K L S D V K C T T V V<br>TTTCGACAAATATACTTATACAAGGAATTGGTGGTGGCGGTGGTGCCTATTGCTACAGTTCAAGCTAAATTGAGTGATAGAAGTGTACAACTGTTGT      | 10700 |
| 10701 | L M Q L L T K L N V E A N S K M H V Y L V E L H N K I L A S D D V<br>TTTAATGCAGCTTTTGACTAAGCTTAATGTTGAAGCAAATTCAAAAATGCATGTTTATCTTGTTGAGTTACACAATAAAATTCTTGCTTCTGATGATGTT   | 10800 |
| 10801 | G E C M D N L L G M L I T L F C I D S T I D L S E Y C D D I L K R S<br>GGAGAGTGCATGGATAATITGTTGGGTATGCTTATAACACTATTTGTATAGAGTTCTACTATTGATTTGAGTGAG                          | 10900 |
| 10901 | T V L Q S V T Q E F S H I P S Y A E Y Ê R A K N L Y E K V L V D S<br>CAACTGTATTACAATCGGTTACTCAAGAATTCTCACATATACCCTCTTATGCTGAATATGAAAGGGCTAAGAATCTTTATGAAAAGGTTTTAGTTGATTC   | 11000 |
| 11001 | K N G G V T Q Q E L A A Y R K A A N I A K S V F D R D L A V Q K K<br>TAAAAATGGTGGTGTTACACAGCAAGAGCTTGCTGCATATCGTAAAGCTGCCAATATTGCAAAGTCAGTTTTTGATAGAGACTTGGCTGTCCAAAAGAAG   | 11100 |
| 11101 | L D S M A E R A M T T M Y K E A R V T D R R A K L V S S L H A L L F<br>TTAGATAGCATGGCAGAGCGTGCTATGACAACAATGTATAAAGAGGGGCGTGTAACAGATGAGCGAGC                                 | 11200 |
| 11201 | S M L K K I D S E K L N V L F D Q A S S G V V P L A T V P I V C S<br>TCTCAATGCTTAAGAAAATAGATTCTGAAAAGCTTAATGTCTTGTTGACCAGGCTAGTGGTGTGTGCCCCTAGCGACTGTTCCAATTGTTTGT          | 11300 |
| 11301 | NKLTLVIPDPETWVKCVEGVHVTYSTVVWNIOT<br>TAATAAGCTTACACTTGTAATACCAGACCCAGAAACGTGGGGGGGG   | 11400 |
| 11401 | VIDADGTELHPTSTGSGLTYCISGANIAWPLKVN<br>GTTATTGATGCCGATGGCACAGGTTACAGGTAGTGGATTGACATACTGTATAAGTGGTGCTAATATAGCATGGCCTTTAAAGGTTA  | 11500 |
| 11501 | LTRNGHNKVDVVLQNNELMPHGVKTKACVAGVD<br>ACTTGACTAGGAATGGGCATAATAAGGTTGATGTTGTTGCAAAATAAGGCTTGCGTAGCAGGTGTAGA   | 11600 |
| 11601 | Q A H C S V E S K C Y Y T N I S G N S V V A A I T S S N P N L K V<br>TCAAGCACATTGTAGCGTAGAGTCTAAATGTTATTACAAATATTAGTGGCAATTCAGTTGTAGCTGCTATTACTTCTTCAAATCCAAATCCTGAAAGTA    | 11700 |
| 11701 | A S F L N E A G N Q I Y V D L D P P C K F G M K V G V K V E V V Y L<br>GCTTCGTTTTTGAATGAGGCAGGCAATCAGATTTATGTAGACTTAGACCCACCATGTAAATTTGGCATGAAAGTGGGTGTCAAGGTTGAGGTTGTTTACT | 11800 |
| 11801 | Y F I K N T R S I V R G M V L G A I S N V V V L Q S K G H E T E E<br>TGTATTITATAAAGAATACAAGGTCGATTGTTAGGGGTATGGTACTTGGTGCTATATGTTGTTGTCTTACAGTCTAAAGGGCATGAAACAGAGGA        | 11900 |
| 11901 | V D A V G I L S L C S F A V D P A D T Y C K Y V A A G N Q P L G N<br>AGTGGATGCTGTTGGCATTCTTTCACTATGTTCATTTGCAGTAGATCCGCGGGCACACATATTGTAAATATGTGGCAGCAGGTAATCAACCTTTAGGTAAC  | 12000 |
| 12001 | CVKMLTVHNGSGFAITSKPSPTPDQDSYGGASVC<br>TGTGTTAAAATGTTGACAGTGCTATATGGTAGTGGTTTTGCTATAACTTCAAAGCCAAGTCCTACTCCTGGACCAGGATTCTTATGGAGGAGCTTCTGTGT                                 | 12100 |

| 12101          | LYCRAHIAHPGSVGNLDGRCQFKGSFVQIPTTE<br>GTCTCTATTGTAGRAGCACATAGCACATCCAGGAAAGTGTAGGAAATTTAGATGGACGTTGTCAATTTAAAGGTTCTTTTGTGCAAATACCTACTACGGA                                  | 12200 |
|----------------|--|-------|
| 12201          | K D P V G F C L R N K V C T V C Q C W I G Y G C Q C D S L R Q P K<br>GAAAGACCCCGTTGGATTCTGTCTACGTAATAAGGTTTGCACTGTTTGCAGTGTGGATTGGATGGA                                    | 12300 |
| 12301          | S S V Q S V A G A S D F D K N Y L N G Y G V A V R L G *<br>TCTTCTGTTCAATCAGTTGCTGGAGCATCTGATTTTGATAAGAATTATTTAAACGGGTACGGGGTAGCAGTGAGGCTCGGCTGATACCCCTTGCTAGTGG            | 12400 |
| 12401          | M F Q N L K R N C A R F Q E<br>ATGTGATCCTGATGTTGTAAAGCGAGCCTTTGATGTTTGTAATAAGGAATCAGCTGGTATGTTTCAAAATTTGAAGCGTAACTGCGCTAGATTCCAGGAA  | 12500 |
| 12501          | L R D T E D G N L E Y L D S Y F V V K Q T T P S N Y E H E K S C Y E<br>CTACGCGATACTGAAGATGGAAATCTTGAGTATCTTGATTCTTACTTTGTAGTAAACAAAC                                       | 12600 |
| 12601          | DLKSEVTADHDFFVFNKNIYNISRQRLTKYTMM<br>AAGACTTAAAGTCAGAAGTAACAGCTGACCATGACTTCTTGTGTTCAATAAGAACATTTACCAATATTAGTAGGCAACGGCTTACTAAATATACTATGAT                                  | 12700 |
| 12701          | DFCYALRHFDPKDCEVLKEILVTYGCIEDYHPK<br>GGACTTCTGCTATGCTTTGAGACATTTCGACCCAAAGGATTGTGAAGAAATACTTGTCACTTATGGTTGTATAGAAGAACTATCACCCTAAG  | 12800 |
| 12801          | WFEENKDWYDPIENSKYYVMLAKMGPIVRRALLN<br>TGGTTTGAGGAGAATAAGGATTGGTACGACCCAATAGAAAACTCAAAATGTGACGACCTATTGTACGACGTGCTTTATTGA  | 12900 |
| 12901          | A I E F G N L M V E K G Y V G V I T L D N Q D L N G K F Y D F G D<br>ATGCTATTGAGTTCGGAAAACCTTATGGTTGAAAAAGGTTATGTTGGTGTTATTACACTCGATAACCAAGACCTTAATGGCAAATTTTATGATTTTGGTGA | 13000 |
| 1 3001         | FQKTAPGAG VPVFDTYYSYMMPIIAMTDALAPE<br>TTITCAGAAGACGGCACCTGGTGCTGGTGTTCCTGTTTTGATACGTATTATTCTTACATGATGCCATCATAGCCATGACGGATGCTTTAGCACCTGAG                                   | 13100 |
| 13101          | RYFEYDVHKGYKSYDLLKYDYTEEKQELFQKYFK<br>AGGTACTTTGAATATGATGTGGGCACAAGGGTTATAAATCTTATGATCTCCGAGGTAGGAAGGA   | 13200 |
| 13201          | YWDQEYHPNCRDCSDDRCLIHCANFNILFSTLI<br>AGTACTGGGATCAAGAGTATCATCCTAACTGCCGTGACTGCGACACTTCAACACTGCGACACTTCAACACTGCTGTTTCTACACTTAT  | 13300 |
| 13301          | PQTSFGNLCRKVFVDGVPFIATCGYHSKELGVI<br>ACCGCAGACTICTITCGGTAATTIGTGTGAGAAAAGTITITTGTTGATGGTGTGTGT   | 13400 |
| 13401          | M N Q D N T M S F S K M G L S Q L M Q F V G D P A L L V G T S N N L<br>ATGAATCAAGATAACACCATGTCTTTTTCAAAAATGGGTTTAAGTCAACTCATGCGGTTTGTTGGGGGAACATCCAATAATT                  | 13500 |
| 13501          | V D L R T S C F S V C A L T S G I T H Q T V K P G H F N K D F Y D<br>TAGTTGATCTTAGAACGTCTTGTTTTAGTGTTTGTGCGTTAACATCTGGTATTACCAAGGGTAACGGTAAAGCCAGGTCACTTTAACAAGGATTTCTATGA | 13600 |
| 13601          | FAEKAGMFKEGSSIPLKHFFYPQTGNAAINDYD<br>TTTTGCAGAAGGCTGGTATGTTTAAGGAGGGTTCGTCATACCACTTAAACCATTTTTCTATCCTCAAACCGGTAATGCTGCTATAAACGATTATGAT                                     | 13700 |
| 1 <b>37</b> 01 | YYRYNRPTMFDICQLLFCLEVTSKYFECYEGGCI<br>TATTATCGTTATAACAGGCCTACCATGTTGACATAGTGCCACCTCTATTTGGTTATGAAGGCGGCTGTA  | 13800 |
| 13801          | PASQVVVNNLDKSAGYPFNKFGKARLYYEMSLE<br>TACCAGCTAGCCAAGTTGTAGTAACAACTTAGAAAAGAGTGCAGGCTATCCATTTAATAAGTTTGGAAAAGCCCGCCTCTATTATGAAATGAGTCTAGA                                   | 13900 |
| 13901          | E Q D Q L F E I T K K N V L P T I T Q M N L K Y A I S A K N R A R<br>GGAACAGGACCAACTCTTCGAGATTACGAAGAATGTCCTACCCACTATAACTCAAATGAATTTAAAATATGCCATATCCGCGAAAAATAGAGCGCGT     | 14000 |
| 14001          | T V A G V S I L S T M T N R Q F H Q K I L K S I V N T R N A S V V I<br>ACAGTGGCAGGTGTGTCTATCCTTTCTCATGAGTAATAGGCAGTTCTCAGAGAGATTCTTAAGTCCATAGTCAACACTAGAAATGCTTCTGTAGTTA   | 14100 |
| 14101          | G T T K F Y G G W D N M L R N L I Q G V E D P I L M G W D Y P K C<br>TTGGAACAACCAAGTTTTATGGGGTTGGGACAACATGTTGAGAAACCTGATTCAGGGTGTTGAAGACCCAATTCTTATGGGTTGGGATTATCCTAAGTG   | 14200 |
| 14201          | DRAMPNLLRIAASLVLARKHTNCCSWSERIYRL<br>TGATAGAGCAATGCCTAATTTGTTGCGTATAGCAGCACTCCTTAGTACTTGCTCGCAAACACCACTAACTGTTGTGGTCGGACGCATTTATAGGTTG                                     | 14300 |
| 14301          | YNECAQVLSETVLATGGIYVKPGGTSSGDATTAY<br>TATAATGAATGCGCCCAGGTCTTATCTGAAACTGTACTGCTACTAGGTGGTACTACTGCTT  | 14400 |
| 14401          | A N S V F N I I Q A T S A N V A R L L S V I T R D I V Y D N I K S<br>ATGCARACAGTGTTTTTAACATAATACAAGCCACATCTGCTAATGTTGCGCGTCTTTTGAGTGTTATAACGCGTGATATTGTCTATGATAATATTAAGAG  | 14500 |

| 14501         | LQYELYQQVYRRVNFDPAFVEKFYSYLCKNFSL<br>CTTGCAGTATGAATTGTATCAGGAGGTCTACAGGGGGGTTAATTTTGACCCAGCCTTTGTTGAAAAGTTTTATTCTTACTTA   | 14600 |
|---------------|---|-------|
| 14601         | MILSDDGVVCYNNTLAKQGLVADISGFREVLYYQ<br>ATGATCTTGTCTGAEGAEGGTGTTGTTTGTAEAACAACAACATTAGCCAAACAAGGTCTTGTAGCAGATATTTCTGGTTTTAGAGAGGGTCTCTACTATC  | 14700 |
| 14701         | NNVFMADSKCWVEPDLEKGPHEFCSQHTMLVEV<br>AGAATAATGTTTTTATGGCTGATTCTAAATGTTGGGTTGAACCAGAATTTAGAAAAAGGCCCACATGAGTTTTGTTCACAACACACAC   | 14800 |
| 14801         | DGEPKYLPYPDPSRILGACVFVDDVDKTEPVAV<br>TGATGGTGACCCTAAGTATTTGCCATACCCAGACCCTCACGCATTTTGGGTGACGTGGATGAGGTGGATAAGACAGAACCTGTGGCTGTT   | 14900 |
| 14901         | MERYIALAIDAYPLVHHENEEYKKVFFVLLAYIR<br>ATGGAGGGTTATATAGGTCTTGCCATAGATGCATCATGAAAATGAAGAAGAGGTATCCTTGTCTCCCTTGCATATATCA   | 15000 |
| 15001         | KLYQELSQNMLMDYSFVMDIDKGSKFWEQEFYE<br>GAAAACTCTATCAAGAGCTTTCTCAGAATATGCTTATGGACTACTCTTTGGAATAGGCTAGGGTGGTAGAATTTTGGGAACAGGGAGTTCTATGA  | 15100 |
| 15101         | NMYRAPTTLQSCGVCVVCNSQTILRCGNCIRKP<br>GAATATGTATAGAGCTCCTACGACTTTGTAGTCGGCGTTGTAGTTGTAATAGTCAAACTATACTACGCTGCGGTAATTGTATTCGTAAACCG   | 15200 |
| 15201         | FLCCKCCYDHVMHTDHKNVLSINPYICSQLGCGE<br>TTTTTGTGTTGTAAGTGTTGCTATGACCAGCATGCGTAGCGACCACAAAAATGTTTTATCTATAAATCCTTATATTTGCTCACAGCTAGGTTGCGGTG  | 15300 |
| 15 <b>301</b> | A D V T K L Y L G G M S Y F C G N H K P K L S I P L V S N G T V F<br>AAGCAGATGTTACTAAAATTGTACCTCGGGGGTATGTCGTACTCTGTGGTAATCATAAACCGAAATTGTCAATACCGTTAGTAGTACTGGTACTGTTTT  | 15400 |
| 15401         | G I Y R A N C A G S E N V D D F N Q L A T T N W S I V E P Y I L A<br>TGGAATTTACAGGGCTAATTGTGCTGGTAGTGAAAATGTTGATGATTTAATCAACTAGCTACTACTAATTGGTCCATTGTCGAACCTTATATTTAGCA   | 15500 |
| 15501         | N R C S D S L R R F A A E T V K A T E E L H K Q Q F A S A E V R E V<br>AATGGCTGTAGTGATTGATGAGAGGTTTTGCTGCAGAGAGCAGGAAGAAGTAGCAACAATTTGCTAGTGCGAGAAGTGGGAGAAG  | 15600 |
| 15601         | FSDRELILSWEPGKTRPPLNRNYVFTGYHFTRT<br>TATTCTCAGATCGTGAATTGATTCATCATGGGAACCAGGGAAAAACCAGGCCGCCATTGAATAGAAATTATGTTTTCACAGGTTATCACTTTACAAGAAC   | 15700 |
| 15701         | SKVQLGDFTFEKGEGKDVVYYKATSTAKLSVGD<br>TAGTAAGGTGCAGCTTGGTGATTTTGAAAAAGGTGAAGGTGAAGGTAAGGTAGTGTCTGTAGGAGACGTCTACTGCTAAATTGTCTGTAGGAGAC  | 15800 |
| 15801         | IFVLTSHNVVSLVAPTLCPQQTFSRFVNLRPNVM<br>ATTITITAACCTCACACAATGTTGTTTCTCCGTAGCGCCAACATTGTGTCACAAACCTTTTCTAGGTTTGTAAATTTAAGACCTAATGTAA   | 15900 |
| 15901         | V P E C F V N N I P L Y H L V G K Q K R T T V Q G P P G S G K S H<br>TGGTACCTGAATGTTTTGTAAATAACATTCCACTTTACCATTTAGTAGGTAAAACAGGAAGCGTACTACAAGGTCCTCCTGGCAGTGGTAAATCCCA  | 16000 |
| 16001         | FAIGLAVYFSSARVVFTACSHAAVDALCEKAFK<br>CTTTGCTATAGGCCTTGCAGTATACTTTAGTGCAGCCGCTGGTGAGAAAAGCTTTTAAG  | 15100 |
| 16101         | FLKVDDCTRIVPQRTTVDCFSKFKANDTGKKYIF<br>TTTCTTAAAGTTGATGATTGCACTCGTATAGTACCCCAAAGGACTACTGCTCGTCGATTGCTCCCAAAATTTAAAGCTAATGACACGACAAGAAAGTACATTT   | 16200 |
| 16201         | S T I N A L P E V S C D I L L V D E V S M L T N Y E L S F I N G K<br>TTAGTACTATTAATGCCTTGCCGGAAGTTAGTTGTGATGATAGTTGACGAGGTTAGTAGTTGACCAATTACGAATTGTCCTTTATTAATGGTAA   | 16300 |
| 16301         | INYQYVVVG DPAQLPAPRTLLNG SLSPKDYNV<br>GATAAATTACCAATATGTTGTGTGTGTGAGGTGATCCGGCTCAATTACCGGCACCCCGCACTTACTT   | 16400 |
| 16401         | V T N L M V C V K P D I F L A K C Y R C P K E I V D T V S T L V Y D<br>GTCACAAACCTTATGGTTTGTGTTAAACCTGATATTTTCCTTGCAAAGTGTTATCGTTGTCTAAGGAAATTGTAGACACTGTGTCTACTCTTGTTTATG  | 16500 |
| 16501         | G K F I A N N P E S R E C F K V I V N N G N S D V G H E S G S A Y<br>ATGGAAAGTTTATTGCAAATAACCCAGAATCACGTGAGTGTTTCAAGGTTATAGTTAATAATGGCAATTCTGATGTAGGACATGAAAGTGGTTCAGCCTA   | 16600 |
| 16601         | NTTQLEFVKDFVCRNKQWREAIFISPYNAMNQR<br>CAACACAACACAATTGGAATTGGGAAGACTTTGTTGTCGCAATAAACAATGGCGGGAAGCAATATTTATT   | 16700 |
| 16701         | A Y R M L G L N V Q T V D S S Q G S E Y D Y V I F C V T A D S Q H A GCTTACCGTATGCTTGGACTTAATGTCCAAACAGTAGATTCTTCTCAAGGTTCAGAGTATGATTATGTCATCTTCTGTGTTACTGCAGATTCGCAGGTTCGGTTCGCGGTTCGCGGTTCGGTTGGTT | 16800 |
| 15801         | L N I N R F N V A L T R A K R G I L V V M R Q R D E L Y S A L K F<br>CACTGAATATTAATAGATTTAATGTGGCGCCTTACAAGAGCTAAGCGTGGTATACTAGTTGTCATGCGCCAGCGTGATGAATTGTATTCTGCTCTTAAGTT  | 16900 |

| 16901          | TELDSETSLQGTGLFKICNKEFSGVHPAYAVTT<br>TACAGAGCTAGATAGTGAAACAAGTCTGCAAGGTACAGGTTTGTTAAAATTTGCAACAAAGAATTTAGTGGTGTCCATCCTGCTTATGCAGTCACAACT                                   | 17000          |
|----------------|--|----------------|
| 1 <b>7</b> 001 | KALAATYKVNDELAALVNVEAGSEITYKHLISLL<br>AAGGETCTTGCTGCAACCTATAAAGTTAATGATGAAACTATAAACATATAAACATCTTATTTTCTCTGT  | 17100          |
| 17101          | G F K M S V N V E G C H N M F I T R D E A I R N V R G W V G F D V<br>TAGGATTCAAGATGAGTGTTAATGTTGAAGGCTGCCACAACATGTTTATAACACGTGATGAGGCAATCCGCAATGTAAGAGGTTGGGTAGGTTTTGATGT  | 17200          |
| 17201          | EATHACGTNIGTNLPFQVGFSTGADFVVTPEGL<br>AGAAGCAACACATGCTTGTGGGCACTAACATTGGTACTAACCTGCCTTTTCAAGTAGGTTTCTCTACTGGTGCAGACTTTGTAGTCACGCCTGAGGGACTT                                 | 17300          |
| 17301          | V D T S I G N N F E P V N S K A P P G E Q F N H L R V L F K S A K P<br>GTAGATACTTCAATAGGCAATAATITTGAGCCTGTGAATTCTAAAGCACCTCCAGGTGAACAATTTAACAGTGTTATTTAAAAGTGCTAAAC        | 17400          |
| 17401          | WHVIRPRIVQMLADNLCNVSDCVVFVTWCHGLE<br>CTTGGCATGTTATAAGACCAAGGATAGTGCGGAGATGTTAGCAGACGATTGTCAGAGTGTTGTCACATGGTGTCATGGCCTAGA  | 17500          |
| 17501          | L T T L R Y F V K I G K E Q V C S C G S R A T T F N S H T Q A Y A<br>ACTAACTACTITGCGCTATITTGTTAAAATAGGCAAGGAACAAGTTTGTTCTTGTGGTTCTAGAGCTACAACTTTAATTCTCATACTCAAGCTTATGCT   | 17600          |
| 17601          | CWKHCLGFDFVYNPLLVDIQQWGYSGNLQFNHDL<br>TGTTGGAAGCATTGTTTGGGTTTTGATTTGATTGGTGGATATTCAACAGTGGGGTTACTCGGGTAACCTACAGTTTAATCATGATT   | 17700          |
| 17701          | HCNVHGHAHVASVDAIMTRCLAINNAFCQDVNW<br>TGCACTGTAATGTGGCACGCCCATGTAGGCTCTGTGGGCGCTATAATGACTCGTTGTCTGCAATTAACAATGCATTTTGTCAAGATGTCAACTG  | 17800          |
| 17801          | DLTYPHIANEDEVNSSCRYLQRMYLNACVDALK<br>GGATTTGACATACCCTCACATTGCAAATGAGGATGAAGTCAATTCTAGTTGTAGATACTCTAAATGCGTGTGTGT   | 17900          |
| 1 <b>7</b> 901 | VNVVYDIGNPKGIKCVRRGDVNFRFYDKNPIVRN<br>GTTAATGTTGTCTATGATATAGGCAACCCTAAAGGTATTAAATGTGTTAGGCGTGGGGATGTTAATTTTAGATTCTATGATAGAACCCAATAGTACGCA                                  | 18000          |
| 18001          | VKQFEYDYNQHKDKFADGLCMFWNCNVDCYPDN<br>ACGTCAAGCAGTTTGGGATAGCATAAATCAGCACAAAGATAGTTGCTGATGGTCTTTGGTATTGTTTGGAATTGTGATGTGATGTTATCCTGATAA                                      | 18100          |
| 18101          | SLVCRYDTRNLSVFNLPGCNGGSLYVNKHAFYT<br>TTCCTTGGTTTGTAGGTATGACACGAAATTTGAGTGTGTTAACCAGGCTGTAAGGTGGTAGGCTGTAACAAACA  | 18200          |
| 18201          | PKFDRISFRNLKAMPFFFYDSSPCETIQVDGVAQ<br>CCTAAATTTGACCGCATTAGCTTCCGCAATTTGAAGCTATGCCATTCTTTTTTATGACTCATCGCCTTGTGAAACCATTCAAGTGGATGGA  | 18300          |
| 18301          | DLVSLATKDCITKCNIGGAVCKKHAQMYAEFVT<br>AAGACCTTGTGTCTCAGGTAGGTAGGTGGTGGTGGTGGTGGTGGTGGTGGTGGTG   | 18400          |
| 18401          | SYNAAVTAGFTFWVTNKLNPYNLWKSFSALQSI<br>TTCTTACAATGCAGCTGTCACGGCTGGCTTTACTTTCGGGTAACTAATAAACTTATAACTTATGGAAAAGTTTTTCAGCTCTCCAGTCTATC  | 18500          |
| 18501          | DNIAYNMYK GGHYDAIAGEMPTVIT GDKVFVIDQ<br>GACAATATIGCTTATAATATGTATAAGGGTGGTCATTATGATGCTATGCTGGGAGAAATGCCCACTGTCATAACTGGAGACAAAGTTITTGTTATTGATC                               | 18600          |
| 18601          | G V E K A V F V N Q T T L P T S V A F E L Y A K R N I R T L P N N<br>AAGGTGTAGAAAAGGCAGTTTTTGTTAATCAAACAACTCTACCTAC  | 18700          |
| 18701          | RILKGLG V D V T NG F V I W D Y A NQ T PLYR N T V K V C<br>CCGTATTTTGAAAGGTTTAGGTGTAGACGTAACCAATGGATTTGTAATTTGGGATTATGCTAACACCAATGTATCGTAATACCGTCAAGGTATGT                  | 18800          |
| 18801          | A Y T D I E P N G L V V L Y D D R Y G D Y Q S F L A A D N A V L V S<br>GCATATACAGATATTGAGCCAAATGGCCTAGTAGTTCTGATGATGATGATAGATA   | 18900          |
| 18901          | T Q C Y K R Y S Y V E I P S N L L V Q N G M P L K D G A N L Y V Y<br>CTACACAGTGTTATAAGCGATATTCATACGTAGAAATACCATCTAATTTGCTCGTTCAGAATGGTATGCCATTAAAAGATGGAGCGAACCTGTATGTTTA  | 19000          |
| 19001          | K R V N G A F V T L P N T I N T Q G R S Y E T F E P R S D I E R D<br>TAAGCGTGTTAATGGTGCGTTTGTTACACTACCTAACACAATAAACACCCAGGGTCGAAGTTATGAAACTTTTGAACCTCGTAGTGACATTGAGCGTGAT  | 19100          |
| 19101          | FLAMSEESFVERYGKDLGLQHILYGEVDKPQLGG<br>TTTCTCGCTATGTCAGAGGGGGGGGTTTGGTAGAAGGCTAGGCCTACAGCACATACTGTATGGTGAAGTTGATAAGCCCCAATTAGGTG  | 19200          |
| 19201          | L H T V I G M Y R L L R A N K L N A K S V T N S D S D V M Q N Y F<br>GTTTACACACTGTTATAGGTATGTACAGACTCTTACGTGCGAAAAGTGCAAGGCCGAAAGTCTGTAACTAATTCGGATTCTGATGTCATGCAAAATTACTT | 19 <b>30</b> 0 |

| 19301 | VLSDNGSYKQVCTVVDLLLDDFLELLRNILKEY<br>TGTATTGTCGGACAATGGTTCTTACAAGCAAGTGTGTGTG  | 19400 |
|-------|--|-------|
| 19401 | G T N K S K V V T V S I D Y H S I N F M T W F E D G S I K T C Y P Q<br>GGTACTAATAAGTCGAAAAGTTGTAACAGTGTCAATTGATTACCATAGCATAGCATAGATTTATGACTTGGTTTGGAAGATGGCAGGTATTAAAACATGTTATCCAC | 19500 |
| 19501 | L Q S A W T C G Y N M P E L Y K V Q N C V M E P C N I P N Y G V G<br>AGCTTCAATCAGCATGGACGTGTGGTTATAAAATATGCCTGAACTTTATAAAGTTCAGAATTGTGTTATGGAACCTTGCAACATTCCTAATTATGGTGTTGG        | 19600 |
| 19601 | ITLPSGILMNVAKYTQLCQYLSKTTICVPHNMR<br>AATAACGTTGCCTAGCGGTATTCTTATGAATGTGGCAAAGTATACACAACTTTGTCAATACCTTTCGAAAACAACTATGTGTACCGCATAACATGCGA  | 19700 |
| 19701 | VMHFGAGSDKGVAPGSTVLKQWLPEGTLLVDNDI<br>GTAATGCATTTCGGAGCAGAGGGGCGAAAAGGAGTGGCGCCAGGTAGTACTGTTCTTAAACAATGGCTCCCAGAAGGGACACTCCTTGTCGATAATGATA   | 19800 |
| 19801 | V D Y V S D A H V S V L S D C N K Y N T E H K F D L V I S D M Y T<br>TTGTAGACTATGTGTCTGATGCACATGTTTCTGTGCTTTCAGATTGCATAAATAA   | 19900 |
| 19901 | D N D S K R K H E G V I A N N G N D D V F I Y L S S F L R N N L A<br>AGATAATGATTCAAAAAAGAAAGCATGAAGGCGTGATAGCCAATAATGGCAATGATGACGTTTTCATATATCTCTCAAGTTTTCTTCGTAACAATTTGGCT         | 20000 |
| 20001 | LGGSFAVKVTETSWHEVLYDIAQDCAWWTMFCTA<br>CTAGGTGGTAGTTTTGCTGTAAAAGTGACAGAGAGACAGTTGGCACGAGGTTTGTGCATGGTGGACAATGTTTTGTACAG   | 20100 |
| 20101 | V N A S S S E A F L I G V N Y L G A S E K V K V S G K T L H A N Y<br>CAGTGAATGCCICTTCTTCAGAAGCATTCTTGATTGGTGTTAATTATTTGGGTGCAAGTGAAAAGGTTAAGGTTGGGGGAAAAAGCGCTGCACGCAAATTA         | 20200 |
| 20201 | IFWRNCNYLQTSAYSIFDVAKFDLRLKATPVVN<br>TATATTTTIGGAGGAATTGTAATTATTTACAAACCTCTGCTTATAGTATATTTTGGCGTTGCTAAGTTTGAATTGAAAGCAACGCCAGTGCTAAT   | 20300 |
|       | L K T E Q K T D L V F N L I K C G K L L V R D V G N T S F T S D S F  |       |
| 20301 | MLVTPLLLVTL<br>TTGAAAACTGAACAAAAGACAGACTTAGTCTTTAATTTAAGTGTGGTAAGTTACTGGTAAGAGATGTTGGTAACACCTCTTTTACTAGTGACCTCT  | 20400 |
| 20401 | V C T M *<br>L C A L C S A V L Y D S S S Y V Y Y Y Q S A F R P P S G W H L Q G<br>TTGTGTGCACTATGTAGTGCTGTTTTGTATGACAGTAGTTCTTACGATTACCAAGGGGCCTTCAGAGCGCCTGGGCATTTACCAAGGGG        | 20500 |

Fig. 2. The sequence of the 'unique' region of mRNA F from the Beaudette strain of IBV. Translations of the ORFs are shown in single-letter amino acid code. The amino acid is shown above the first base of the appropriate codon. The translation starting at position 20368 is the  $NH_2$  terminus of the spike precursor protein.



Fig. 3. Diagram showing the positions of the main ORFs in the 'unique' region of mRNA F. The two large ORFs, designated F1 and F2 are shown, as well as a small ORF at the 5' end of the genome, and the start of the spike precursor gene, which overlaps with F2.

The second large ORF, F2, extends into the 'unique' region of mRNA E and in fact overlaps the coding sequences for the spike protein gene by 16 amino acids.

# Potential sources of error

All the sequence information has been confirmed by sequencing M13 clones obtained from both strands of the DNA. In addition most of it has been sequenced several times from different M13 clones. The 14 cDNA clones used to obtain the sequence of mRNA F contain, including overlaps, 24765 bases. During the shotgun sequencing of these clones 203113 bases have been sequenced, so that each base has, on average, been sequenced 8.2 times. However there are two regions we have checked more carefully. The first is at positions 12340 to 12390 where F1 ends and F2 begins. An error here leading to a frameshift could make the difference between two large ORFs and one very large ORF. The second is at position 167 where the very small 11 amino acid ORF ends. A frameshifting error here could mean that this first ORF can continue for another 77 amino acids until position 397. There are two possible sorts of error. The first is an artefact in the sequencing gels leading to a misreading. The sequence on both strands appears perfectly clear in both these regions. Both regions have been sequenced using formamide gels, high temperature gels, in addition to the use of deoxyinosine triphosphate (Bankier & Barrell, 1983) or deoxy-7-deazaguanosine triphosphate (Mizusawa *et al.*, 1986) to replace deoxyguanosine triphosphate and cytosine-modified sequence reaction products (Ambartsumyan & Mazo, 1980) to avoid gel compressions.

The second potential source of error is either a reverse transcriptase error during the synthesis of the cDNA or the occurrence of a mutant RNA molecule from which the cDNA was copied, both of which would lead to an incorrect cDNA clone. In the case of position 167 the sequence has been obtained from an equivalent clone from the M41 strain of IBV and is identical. In the case of the sequence between F1 and F2 the sequence has been confirmed from two additional independent cDNA clones, by sequencing directly from the double-stranded DNA using an oligonucleotide primer (Korneluk *et al.*, 1985). Fig. 4(*a*) shows the relevant sequence in this region and Fig. 4(*b*) shows a sequencing gel of bases 12333 to 12390 obtained directly from a cDNA clone using an oligonucleotide primer. In addition the sequence has been obtained directly from the virion RNA using specific oligonucleotide primers at both of these points and has confirmed the original gel readings. At positions 12333 to 12390 the sequence in this region is identical.

Gel compressions are thought to be caused by the presence of hairpin loops in the DNA migrating down the gel. Examination of the sequence in these regions shows that there are several possibilities for the formation of fairly large hairpins, including for example, at the position between F1 and F2, the sequence GGGGTA with its exact complement TACCCC 24 bases further on. At this position (12380), in the region where the reading frame changes between F1 and F2, the sequence has been determined from ten separate M13 clones. It is interesting to note that one of these clones gave a different sequence reading in that a CT dinucleotide, which appears in the other nine M13 readings, was not present. This is unusual as normally all independent M13 clones agree. It is possible that the secondary structure in this region has some effect on the fidelity of copying by polymerases.

## Computer analysis

Extensive computer analysis has been carried out in an attempt to identify some salient features on the bleak landscapes of these large ORFs. Searches for homologies with other viral polymerases have been performed using the NBRF protein identification resource (George *et al.*, 1986). Short regions of fairly low homology with several viral polymerases can be identified but in general they do not rise significantly above the background of matches with proteins that are apparently unrelated. One region, between amino acids 1342 and 1350, has a fairly good match (8/9 amino acids) with the nsP2 protein of Sindbis virus, a protein which is known to be involved in RNA replication (Strauss & Strauss, 1983). This region also has a match with the 1a protein of brome mosaic virus. These matches are shown in Fig. 5. One of the most interesting matches is at the 5' end of the first large ORF. The first 300 amino acids have a low-level but extensive homology with the replication initiation protein from *Escherichia coli* (Germino & Bastia, 1982). The homology is statistically significant and it may indicate that this region of the polymerase protein is involved in initiation of replication of either the positive or negative strands.

The predicted amino acid sequences of the large ORFs have been compared against themselves and against each other to see whether there are any repeats which might represent 
 S
 L
 R
 Q
 P
 K
 S
 V
 Q
 S
 V
 A
 G
 A
 S
 D
 F
 N
 Y
 L
 N
 G
 Y
 R
 L
 G
 Y
 P
 L
 F
 T
 T
 K
 I
 F
 C
 S
 I
 S
 F
 T
 F
 T
 F
 T
 T
 K
 I
 F
 C
 S
 I
 S
 F
 T
 T
 K
 I
 F
 C
 S
 I
 S
 F
 T
 T
 F
 T
 T
 F
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T
 T

L V D V I L M L \* S E P L M F V I R N Q L V C F K I \* S V T A L D S R N <u>A S G C D P D V V K R A F D V C N K E S A G M F Q N L K R N C A R F Q E</u> C \* W M \* S \* C C K A S L \* C L \* \* G I S W Y V S K F E A \* L R \* I P G GCTAGTGGATGTGATCCTGATGTTGTAAAGGAACCAGCTGGTATGTTTGAAAATTTGAAGGGTAACTGCGCTAGATTCCAGAAA 12400 1241D 12420 12430 12440 12450 12460 1247D 12480 12490 12500





| BMV | SCHRLLV | DEAGLLI | HYGQLL | VVAAL  | SKCSQ | VLAF- | -GDTEC | )   | ISFK | GRDAGI | FKLLHG | NLQYDR | RDV-\ | HKTYRCF  | QDVIA   | AVNLLI | KRKCGNRDTI | КY  |
|-----|---------|---------|--------|--------|-------|-------|--------|-----|------|--------|--------|--------|-------|----------|---------|--------|------------|-----|
| TDU |         |         | ••     | •••    | ••    | :.    |        |     | •••• | • •    | ••••   | • • •  |       | . : :::: |         | .:.:   | WORKETAN   |     |
| IBV | SCUILLV | DEVSML  | INYELS | F INGK | INYUY |       | -GDPAL |     | LNGS | _SPKU  | YNVVIN |        | PUIFL | AKLYNL   | KEIVU   | 1V51L  |            | .vP |
| SV  | AVEVLYV | DEAFACI | HAGALL | ALIAI  | VRPRK | KVVLI | CGDPMC | )CC | FFNM | NQLKVI | HENHPE | KDICTK | -TFYK | YISRRC   | IQPVTA: | IVSTL  | HYDGKMKTT  | ΝP  |

Fig. 5. Comparison between amino acid sequences of brome mosaic virus (BMV), infectious bronchitis virus (IBV) and Sindbis virus (SV). The BMV sequences are amino acids 748 to 838 of the 1a protein. The SV sequences are amino acids 785 to 878 of the nsP2 protein, The IBV sequences are amino acids 1248 to 1356 of F2. A colon shows identical amino acids and a dot shows similar (Kanehisa, 1982) amino acids. The dashes in the sequences are blank characters inserted to achieve optimal alignment.

two separate but similar polymerases. A dot matrix comparison, such as DIAGON (Staden, 1982*a*), reveals no repeats. However several low homology repeats can be detected using the program FASTP (Lipman & Pearson, 1985). These are shown on Fig. 6(a) beneath a hydrophilicity plot (Kyte & Doolittle, 1982) of the amino acid sequences of F1 and F2. Fig. 6(b to *e*) shows the amino acid matches in these regions. The spacing between the repeats marked A and B is very similar in both cases, 1157 amino acids in F1 and 1183 amino acids in F2. It is possible that these represent residual domains of homology between two polymerases which were at one time more closely related. The areas marked C and D also show regions of homology. The diagram also shows several very hydrophobic regions in the first large ORF which represent potential membrane-spanning domains.

Computer analysis has also detected a homology between the non-coding region at the 5' end of the positive strand, and the 5' end of the negative strand (i.e. the reverse complement of the non-coding region at the 3' end of the positive strand). This is shown in Fig. 7. These sequences, on the positive and negative strands, are approximately the same distance from their 5' ends, 52 bases and 48 bases [excluding the poly(A) tail] respectively, and may play some role in the replication of the positive and negative strands.

## Homology regions

At position 599 the sequence CTGAACAA occurs. This is identical to the sequence which occurs in the 'homology regions' at the 5' ends of the bodies of mRNAs D and E (Boursnell *et al.*, 1985*b*; Binns *et al.*, 1985*b*). These sequences are thought to be recognition sites for binding of the polymerase/leader complex during the synthesis of the subgenomic RNAs (Baric *et al.*, 1983). The same sequence CTGAACAA occurs at position 3293. Neither of these positions are known to be situated at the 5' end of an mRNA species as are all the other homology regions. We have attempted to determine whether there is some feature of the sequence context surrounding these homology regions which sets them apart from homology regions which are known to occur at the 5' end of the bodies of mRNAs. Accordingly, a consensus sequence has been calculated from the sequences surrounding the known homology regions at the ends of mRNAs A to F. This consensus sequence includes six bases to the left of the core homology

Fig. 4. (a) The nucleotide sequence in the region between F1 and F2, with a translation in single-letter amino acid code of three reading frames. The amino acid is shown above the second base of the appropriate codon. Stop codons are marked as asterisks. The frames which are open in F1 and F2 are underlined and the methionine at the start of F2 is boxed in. (b). A DNA sequencing gel obtained by sequencing a double-stranded cDNA clone using an oligonucleotide primer. The sequence shown is from 12333 to 12390, and is the reverse complement of the sequence shown in (a). (c) The same three reading frames as shown in (a), with a graph for each showing the extent to which that reading frame conforms to the codon usage found for the amino acid sequence of F1 and F2. The frame which conforms best to the F1/F2 codon usage is marked with a series of dots and marked F1 or F2. Stop codons are marked as short vertical lines along the centre of each frame, and start codons as bars with filled-in circles on top. The two stop codons at 12339 (TAA) and 12382 (TGA) are marked as is the start codon at 12459. The program used is the 'codon usage' option from ANALYSEQ (Staden, 1984b, 1983c) and uses the method of Staden & McLachlan (1982). The parameters used were a window length of 25 and an output length of 1. (Codon usage analysis from the spike, membrane and nucleocapsid gene data gives a very similar result.)



(b) Repeat A

F1 484 EFVKTYVCKAQMSIVILAAVLGEDIWHLVSQVIYKLGVLFTKVVDFC---DKHWKGFCVQLKRAKLIVTE

F1 TFCVLKGVAQHCFQLLLDAIHSLYKSFKKCALGR---IHGDLLF

F2 RFNVALTRAKRGILVVMRQRDELYSALKFTELDSETSLQGTGLF

(c) Repeat B

| Fl | 1630 | VKMGDKIC | GVTM | GLWRAE | HLNK | PNL | ERIFNI | AKK | AIN  | /GSSI | VVTT | QC  | GKLI | SKAP | ATFIA | DK۱ | /GGG | VVRNITD | ) |
|----|------|----------|------|--------|------|-----|--------|-----|------|-------|------|-----|------|------|-------|-----|------|---------|---|
|    |      |          | •••  | .::.   |      | ••  | .::    | ::  | ••   | ••    |      |     | •    | ••   | .:    | :   | :    | . : :   |   |
| F2 | 2570 | VKVSGKTI | HANY | TEWRNC | NYID | TSA | YSTEDV | AKE | ol F | A KA  | τργι | /NI | KTEQ | (TDI | VENI. | IKC | GKL  | LVRDVGN | L |

(d) Repeat C

| Fl | 3696 | VKTKACVAC | GVDQAH | CSVESKCY | YTN: | ISGNSV | VAAI | TSSN  | IPN     | -LKI | IAS | FLN  | EAC | iN  | -QI |
|----|------|-----------|--------|----------|------|--------|------|-------|---------|------|-----|------|-----|-----|-----|
|    |      | ***       | ::.:   | : :      | ٠    | ••••   | ••   | .:.   |         | ::   | :   | ::.  | ::  | •   | ::  |
| F2 | 1996 | VKPTAYAY  | /VDEA- | CLVDDFVN | ILKY | KAATPG | KDSł | ISSAV | KCFSVTD | FLK  | CAV | FLKI | EAL | KCE | .01 |

(e) Repeat D

Fig. 6. (a) Hydropathicity plots (Kyte & Doolittle, 1982) of the predicted amino acid sequences of ORFs F1 and F2. Values above the line are hydrophobic and values below the line are hydrophilic. The hydropathicity is calculated using a moving window of 41 amino acids, with a value plotted every 21 residues. The pairs of bars marked A, B, C and D show regions of partial homology [see Results and (b) to (e)]. (b to e) Amino acid sequences of the matches depicted by the bars in (a). A colon shows identical amino acids and a dot shows similar (Kanehisa, 1982) amino acids. The dashes in the sequences are padding characters inserted to achieve optimal alignment.

region CT(T/G)AACAA present in all the regions, the eight bases of the core homology itself, and four bases to the right. The consensus has been compared to the complete sequence using the computer program FITCONSENSUS (Devereux *et al.*, 1984). The program successfully identifies the known homology regions with scores ranging from 74.6 to 64.1. The 14 next best fitting regions identified have a range of scores well separated from those of the known 52 TTTAACTTAACAAAACGGACTTAAATACCTACAGCTGGTCCTCATAGGTGTTCCATTGCAGTGCACT 118 11 TTTAACTTAACAAAACGGACTTAAATACCTACAGCTGGTCCTCATAGGTGTTCCATTGCAGTGCACT 118 48 TTAAACTTAACTTAA---ACTAAAATT--TAGCTCTTCCCCTAATGGCGGTCCTAGTGCTGTACCCT 109

Fig. 7. Comparison between (top) the nucleotide sequence of the 5' end of the genome and (bottom) the reverse complement of the 3' end of the genome (i.e. the 5' end of the negative strand). Colons show identical bases. The dashes in the sequences are padding characters inserted to achieve optimal alignment.

G A S D F D K N Y L N G Y G V A V R L G \* GGAGCATCTGATTTTGATAAGAATTATTTAAACGGGTACGGGTAGCAGTGAGGCTCGGCTGATACCCCTTGCTAGTG :::: :: :: :: :: :: :: :: :: :: :: :: GTAGCTATGGTTAGAGGGAGTATCCTAGGAAGAGATTGTCTGCAGGGCCTAGGGCCTCGCTTGACAAATTTATAGGGA V A M V R G S I L G R D C L Q G L G L R L T N L \*

Fig. 8. Nucleotide and predicted amino acid sequences where ribosomal frameshifting may occur. The top sequence is at the F1/F2 junction of IBV, and the bottom sequence is at the gag/pol junction of Rous sarcoma virus. Colons show identical bases.

homology regions, with a tight cluster of scores (53.6 to 58.8). The CTGAACAA sequence at position 599 scores even lower. It seems probable, therefore, that the two CTGAACAA sequences at 599 and 3293 are chance matches with the core sequence, but when surrounding sequences are taken into account the differences are enough to ensure that they are not major sites for the binding of the leader/polymerase complex.

# DISCUSSION

The 20500 bases of sequence presented in this paper complete the sequence of the Beaudette strain of avian infectious bronchitis virus, the type species of the Coronaviridae. The complete sequence, excluding the poly(A) tail at the 3' end, is 27608 residues. This is somewhat larger than the previously estimated size of the viral RNA which had been put at 20 to 24 kilobases (Lomniczi & Kennedy, 1977). The sequence of the 'unique' regions of mRNAs A, B, C, D and E have already been published, covering some 8 kilobases at the 3' end of the genome and including the genes for the major structural proteins of the virus. The 20 kilobases at the 5' end of the viral RNA constitutes the 'unique' region of mRNA F, the genome-sized RNA. This is thought to code for a polymerase or polymerases which carry out all the necessary replication and transcription functions of the virus.

Sequence analysis shows that the main part of the 'unique' region of mRNA F appears to contain two large ORFs. Because of the importance of determining whether there are one or two ORFs, we have considered the possibility that mRNA F in fact contained one very large ORF, and that a sequencing error or a mutant cDNA clone had led to a frameshift. Because of this the sequence in the region between the two ORFs has been checked exceedingly carefully. The relevant sequence is shown, with translations in the three reading frames, in Fig. 4(a). Any frameshift error must occur within 43 bases between positions 12341 and 12383. Two independent cDNA clones and direct RNA sequences from virion RNA give the same result. There are no obvious signs of sequence artefacts such as compressions, and indeed several gel systems and sequencing methods which could resolve compressions (see Methods and Results) do not show any change in the sequence. Fig. 4(b) shows a sequencing gel representing this region, obtained by sequencing a cDNA clone directly using an oligonucleotide primer. It can be seen that the sequence appears clear and unambiguous. Unless, therefore, there is some singular form of unresolvable and undetectable sequencing artefact, we must accept that the sequence here is correct.

The problem now arises as to how translation of the second ORF, F2, is achieved. No mRNA has been detected at this point, and no homology region which might suggest the presence of one can be seen in the RNA sequence (see Results). It is possible that the ribosomes, having completed translation of the first ORF, F1, reinitiate translation at the first AUG of F2, or that internal initiation occurs, as appears to be the case with the phosphoprotein mRNA of vesicular

stomatitis virus (Herman, 1986). There is however one piece of evidence that suggests that neither of these alternatives is the case. If the second ORF is genuinely a separate gene, then the 70 or so bases preceding its initiation codon should be non-coding sequences, comparable to the 5' non-coding sequences preceding other IBV genes. In fact, if translated, they exhibit a heavy codon bias (Staden & McLachlan, 1982; Staden, 1984c) similar to the bias found in other IBV genes. This is shown graphically in Fig. 4(c) where it can be seen that the frame with typical IBV codon bias switches from that of F1 to that of F2 exactly at the point where the ORF changes. This strongly suggests that the sequences before the AUG of F2 have a coding function. One way to resolve this problem is to postulate that on some occasions, during translation of mRNA F, a ribosome slippage occurs, which introduces a frameshift and allows translation to continue unhindered from F1 into F2. Ribosomal frameshifting has been described in bacteriophage (Kastelein et al., 1982), prokaryotic (Atkins et al., 1972) and eukaryotic (Fox & Weiss-Brummer, 1980; Jacks & Varmus, 1985) systems. Such a mechanism could be conceived in the case of IBV as a form of translational control designed to provide coordinated expression of two polymerases, with the protein from the first gene being produced at a higher level than that from the second gene. In the case of Rous sarcoma virus (Jacks & Varmus, 1985) expression of the pol gene requires a frameshift by the ribosome. Some well-controlled work by these authors, using cell-free translation systems, has demonstrated that the frameshifting is sequence-specific. Moreover it occurs ten times more efficiently in a eukaryotic system than in a prokaryotic system, indicating that there are specific eukaryotic signals to which the prokaryotic system responds poorly. The region of sequence responsible for the frameshifting has been narrowed down to 24 nucleotides. Both IBV and Rous sarcoma virus require a shift into the -1 frame to occur, and it may be that similar frameshifting signals are present in both sequences. Accordingly the 24 nucleotides of Rous sarcoma virus sequence have been compared to the 43 nucleotides of IBV sequence within which any frameshift must occur (see Fig. 4a). Interestingly a match of 8/9 nucleotides can be found, both sequences occurring in the same frame and both within 20 bases of the termination codon (see Fig. 8). Further work will be needed to determine whether this sequence forms part of any signals which may promote ribosomal frameshifting.

For each of the other IBV mRNAs, the first AUG to occur after the homology region either is used to initiate synthesis of a protein, as is the case for the spike and membrane proteins (Binns et al., 1985b; Boursnell et al., 1984), or is present at the start of a reasonable sized ORF which could code for a polypeptide of 7K or more. Thus it is surprising to find the first AUG, at position 131, at the start of a small, 11 amino acid, ORF. The sequence context around this first AUG does not conform to Kozak's consensus for functional initiation codons whereas the context round the second AUG does. A similar small ORF of 12 amino acids occurs at the 5' end of RNA 1 of alfalfa mosaic virus (Cornelissen et al., 1983), an RNA species encoding a 115K product thought to be involved in RNA replication. In this case also only the second AUG conforms to the Kozak consensus. Both these cases suggest the possibility that the ribosomes can bypass the first, non-functional, AUG and initiate translation at the second. It is likely that this also occurs in mRNA D of IBV to allow translation of the second and third ORFs (Boursnell et al., 1985b).

It is not known for coronaviruses whether the sequences at the 5' end of the genome produce a polyprotein which is subsequently cleaved into separate proteins, as is the case for alphaviruses (Strauss *et al.*, 1984), or whether the viral polymerase acts as an extremely large multifunctional enzyme. Whether or not it is cleaved post-translationally into separate proteins, such an enzyme would need to perform several functions. First it must synthesize the negative-stranded template. From this template it must synthesize the leader sequence and then the subgenomic mRNAs, for which it needs the ability to recognize highly conserved signal sequences (Baric *et al.*, 1983, 1985; Spaan *et al.*, 1983; Brown & Boursnell, 1984), a capping ability (Lai *et al.*, 1982) and probably the ability to reinitiate transcription at these points (Lai *et al.*, 1985; Makino *et al.*, 1986). If it is cleaved into separate proteins it may encode a protease function to do this. Two polymerase activities, early and late, have been identified in MHV-infected cells (Brayton *et al.*, 1982). These have different ionic requirements and different pH optima. Both polymerase activities are associated with two different membrane fractions, a light fraction which appears

to synthesize positive-stranded genome-size RNA and a heavy fraction which also synthesizes subgenomic RNAs (Brayton *et al.*, 1984). Some evidence for two polymerase-coding genes can be found in the nucleotide sequence of mRNA F, in that there are small regions of residual homology between the predicted amino acid sequences of F1 and F2 (see Results and Fig. 6).

The question of whether the cDNA clones sequenced in this study might derive from mutant, non-viable RNA molecules is an interesting one. The error rate of RNA polymerases is fairly high (Steinhauer & Holland, 1986) and many of the RNA molecules in an infected cell may be different from that in the original infecting virus. If the mutation rate is 1 in 10000 then over the 20 kilobases of sequence presented here, there may be one or two changes each time one strand was copied into another. While the viral RNA is replicating within the cell, it is likely that mutant, and possibly defective, virion RNA molecules will accumulate with little selection against them, and, unless they have gross structural defects, most of them will be packaged into virions. It is these virions, without any further selection for viability, which are used to extract the RNA which is used to synthesize cDNA. In addition the infecting virus will be a mixture of different RNA molecules, even though it has been plaque-purified. However, be that as it may, there is no evidence for very high mutation rates in the cDNA clones which we have sequenced here. For the clones covering the 20 kilobases there are 4659 bases of overlap between separate, independent clones (all made from the same RNA preparation). In the overlap regions there was not one difference, there being 100% agreement between the sequences from adjacent clones.

This is in contrast to results found by Schubert *et al.* (1984) while sequencing the polymerase gene of vesicular stomatitis virus. The gene spans 6380 nucleotides and each region was sequenced from approximately three cDNA clones, giving 19140 nucleotides of overlap. In these 19140 nucleotides they found 20 nucleotide changes, including four insertions or deletions, giving an overall mutation rate of approximately  $10^{-3}$ . In the 9318 (4659 × 2) nucleotides of IBV cDNA clones which can be checked on another clone, there were no changes. Over 9318 nucleotides change; thus, since there were no changes, the overall mutation rate is probably lower than this. Given the number of rounds of replication which will have occurred between the original plaque isolation and the production of the cDNA clones, the mutation rate per base incorporated is likely to be considerably lower than this. It is interesting to speculate on the disparity between the vesicular stomatitis virus and the IBV results in this case, and on whether the (presumably) very large IBV polymerase, or polymerases, has a lower intrinsic error rate than the VSV polymerase.

Sequencing of cDNA clones from the 'unique' region of mRNA F has revealed the rather unexpected presence of two large ORFs. Although the sequence in the region between these has been obtained from three independent cDNA clones and from the virion RNA, the possibility of some bizarre form of sequence artefact cannot be totally discounted. It will be interesting to see if a similar frameshift occurs in an equivalent position in the coronavirus MHV genome. Experiments can now be designed to confirm the reading frame switch by other means. For example *in vitro* translation of SP6 polymerase transcripts from this region can be performed and the sizes of the products determined. Although no mRNA has been detected with a 5' end near the beginning of the second ORF, a search for a low abundance mRNA species can now be carried out by primer extension from mRNA preparations. In addition, the availability of sequence data from the IBV polymerase(s) allows antisera to be raised against products expressed from selected parts of the sequence. These will prove useful in determining the fate of the large polypeptides predicted from the nucleotide sequence, showing whether posttranslational cleavage occurs, and attempting to unravel the relationship between the various polymerase activities which have been detected in coronavirus-infected cells.

We are grateful to Bridgette Britton, Penny Gatter, Neil Macey, Rona Chellew and Steve Laidlaw for excellent technical assistance. We would like to thank Dave Cavanagh and Phil Davis for help with the sequencing of the virion RNA. We would also like to thank Alan Bankier for the gift of some deoxy-7-deazaguanosine triphosphate and for general advice and encouragement during the DNA sequencing.

#### REFERENCES

- AMBARTSUMYAN, N. S. & MAZO, A. M. (1980). Elimination of the secondary structure effect in gel sequencing of nucleic acids. FEBS Letters 114, 265-268.
- ATKINS, J. F., ELSEVIERS, D. & GORINI, L. (1972). Low activity of beta-galactosidase in frameshift mutants of Escherichia coli. Proceedings of the National Academy of Sciences, U.S.A. 69, 1192-1195.
- BANKIER, A. & BARRELL, B. G. (1983). Shotgun DNA sequencing. In *Techniques in the Life Sciences (Biochemistry)*, vol. B5: Techniques in Nucleic Acid Biochemistry, pp. 1-34. Edited by R. A. Flavell. Ireland: Elsevier.
- BARIC, R. S., STOHLMAN, S. A. & LAI, M. M. C. (1983). Characterisation of replicative intermediate RNA of mouse hepatitis virus: presence of leader RNA sequences on nascent chains. Journal of Virology 48, 633-640.
- BARIC, R. S., STOHLMAN, S. A., RAZAVI, M. K. & LAI, M. M. C. (1985). Characterisation of leader-related small RNAs in coronavirus-infected cells: further evidence for leader-primed mechanism of transcription. Virus Research 3, 19–33.
- BEAUDETTE, F. R. & HUDSON, C. B. (1937). Cultivation of the virus of infectious bronchitis. Journal of the American Veterinary Medical Association 90, 51-60.
- BIGGIN, M. D., GIBSON, T. J. & HONG, G. F. (1983). Buffer gradient gels and <sup>35</sup>S label as an aid to rapid DNA sequence determination. Proceedings of the National Academy of Sciences, U.S.A. 80, 3963-3965.
- BIGGIN, M., FARRELL, P. J. & BARRELL, B. G. (1984). Transcription and DNA sequence of the BamHI L fragment of B95-8 Epstein-Barr virus. EMBO Journal 3, 1083-1090.
- BINNS, M. M., BOURSNELL, M. E. G., FOULDS, I. J. & BROWN, T. D. K. (1985a). The use of a random priming procedure to generate cDNA libraries of infectious bronchitis virus, a large RNA virus. Journal of Virological Methods 11, 265–269.
- BINNS, M. M., BOURSNELL, M. E. G., CAVANAGH, D., PAPPIN, D. J. C. & BROWN, T. D. K. (1985b). Cloning and sequencing of the gene encoding the spike protein of the coronavirus IBV. Journal of General Virology 66, 719–726.
- BOURSNELL, M. E. G. & BROWN, T. D. K. (1984). Sequencing of coronavirus IBV genomic RNA: a 195-base open reading frame encoded by mRNA B. Gene 29, 87-92.
- BOURSNELL, M. E. G., BROWN, T. D. K. & BINNS, M. M. (1984). Sequence of the membane protein gene from avian coronavirus IBV. Virus Research 1, 303-313.
- BOURSNELL, M. E. G., BINNS, M. M., FOULDS, I. J. & BROWN, T. D. K. (1985*a*). Sequences of the nucleocapsid genes from two strains of avian infectious bronchitis virus. *Journal of General Virology* **66**, 573-580.
- BOURSNELL, M. E. G., BINNS, M. M. & BROWN, T. D. K. (1985b). Sequencing of coronavirus IBV genomic RNA: three open reading frames in the 5' 'unique' region of mRNA D. Journal of General Virology 66, 2253-2258.
- BRAYTON, P. R., LAI, M. M. C., PATTON, C. D. & STOHLMAN, S. A. (1982). Characterisation of two polymerase activities induced by mouse hepatitis virus. *Journal of Virology* 42, 847–853.
- BRAYTON, P. R., STOHLMAN, S. A. & LAI, M. M. C. (1984). Further characterisation of mouse hepatitis virus RNAdependent RNA polymerases. Virology 133, 197–201.
- BROWN, T. D. K. & BOURSNELL, M. E. G. (1984). Avian infectious bronchitis virus genome RNA contains sequence homologies at the intergenic boundaries. Virus Research 1, 15–24.
- BROWN, T. D. K., BOURSNELL, M. E. G., BINNS, M. M. & TOMLEY, F. M. (1986). Cloning and sequencing of 5' terminal sequences from avian infectious bronchitis virus genomeic RNA. Journal of General Virology 67, 221-228.
- CATON, A. J., BROWNLEE, G. G., YEWDELL, J. W. & GERHARD, W. (1982). The antigenic structure of the influenza virus A/PR/8/34 hemagglutinin (H1 subtype). Cell 31, 417-427.
- CAVANAGH, D. (1981). Structural polypeptides of coronavirus IBV. Journal of General Virology 53, 93-103.
- CORNELISSEN, J. C., BREDERODE, F. T., MOORMANN, R. J. M. & BOL, J. F. (1983). Complete nucleotide sequence of alfalfa mosaic virus RNA 1. Nucleic Acids Research 11, 1253-1265.
- DEININGER, P. L. (1983). Random subcloning of sonicated DNA: application to shotgun DNA sequence analysis. Analytical Biochemistry 129, 216-223.
- DEVEREUX, J., HAEBERLI, P. & SMITHIES, O. (1984). A comprehensive set of sequence analysis programs for the VAX. Nucleic Acids Research 12, 387-395.
- FOX, T. D. & WEISS-BRUMMER, B. (1980). Leaky +1 and -1 frameshift mutations at the same site in a yeast mitochondrial gene. *Nature, London* 288, 60-63.
- GEORGE, D. G., BARKER, W. C. & HUNT, L. T. (1986). The protein identification resource (PIR). Nucleic Acids Research 14, 11-15.
- GERMINO, J. & BASTIA, D. (1982). Primary structure of the replication initiation protein of plasmid R6K. Proceedings of the National Academy of Sciences, U.S.A. 79, 5475-5479.
- HERMAN, R. C. (1986). Internal initiation of translation on the vesicular stomatitis virus phosphoprotein mRNA yields a second protein. Journal of Virology 58, 797–804.
- HONG, G. F. (1981). A method for sequencing single-stranded cloned DNA in both directions. *Bioscience Reports* 1, 243–252.
- JACKS, T. & VARMUS, H. E. (1985). Expression of the Rous sarcoma virus pol gene by ribosomal frameshifting. Science 230, 1237-1242.
- KANEHISA, M. I. (1982). Los Alamos sequence analysis package for nucleic acids and proteins. Nucleic Acids Research 10, 183-196.
- KASTELEIN, R. A., REMAUT, E., FIERS, W. & VAN DUIN, J. (1982). Lysis gene expression of RNA phage MS2 depends on a frameshift during translation of the overlapping coat protein gene. *Nature, London* 295, 35-41.
- KORNELUK, R. G., QUAN, F. & GRAVEL, R. A. (1985). Rapid and reliable dideoxy sequencing of double-stranded DNA. Gene 40, 317-323.

- KOZAK, M. (1983). Comparison of initiation of protein synthesis in procaryotes, eucaryotes, and organelles. Microbiological Reviews 47, 1-45.
- KYTE, J. & DOOLITTLE, R. F. (1982). A simple method for displaying the hydropathic character of a protein. Journal of Molecular Biology 157, 105-132.
- LAI, M. M. C., PATTON, C. D. & STOHLMAN, S. A. (1982). Replication of mouse hepatitis virus: negative-stranded RNA and replicative form RNA are of genome length. *Journal of Virology* 44, 487–492.
- LAI, M. M. C., BARIC, R. S., MAKINO, S., KECK, J. G., EGBERT, J., LEIBOWITZ, J. L. & STOHLMAN, S. A. (1985). Recombination between nonsegmented RNA genomes of murine coronaviruses. *Journal of Virology* 56, 449–456.

LEIBOWITZ, J. L., WILHELMSEN, K. C. & BOND, C. W. (1981). The virus-specific intracellular RNA species of two murine coronaviruses: MHV-A59 and MHV-JHM. Virology 114, 39-51.

LEIBOWITZ, J. L., WEISS, S. R., PAAVOLA, E. & BOND, C. W. (1982). Cell-free translation of murine coronavirus RNA. Journal of Virology 43, 905–913.

LIPMAN, D. J. & PEARSON, W. R. (1985). Rapid and sensitive protein similarity searches. Science 227, 1435-1441.

LOMNICZI, B. (1977). Biological properties of avian coronavirus RNA. Journal of General Virology 36, 531-533.

- LOMNICZI, B. & KENNEDY, I. (1977). Genome of infectious bronchitis virus. *Journal of Virology* 24, 99–107. MAKINO, S., STOHLMAN, S. A. & LAI, M. M. C. (1986). Leader sequences of murine coronavirus mRNAs can be freely
- reassorted: evidence for the role of free leader RNA in transcription. Proceedings of the National Academy of Sciences, U.S.A. 83, 4204-4208.
- MAXAM, A. M. & GILBERT, W. (1980). Sequencing end-labeled DNA with base-specific chemical cleavages. Methods in Enzymology 65, 499-560.
- MIZUSAWA, S., NISHIMURA, S. & SEELA, F. (1986). Improvement of the dideoxy chain termination method of DNA sequencing by use of deoxy-7-deazaguanosine triphosphate in place of dGTP. Nucleic Acids Research 14, 1319–1324.
- SANGER, F., NICKLEN, S. & COULSON, A. R. (1977). DNA sequencing with chain-terminating inhibitors. Proceedings of the National Academy of Sciences, U.S.A. 74, 5463-5467.
- SCHOCHETMAN, G., STEVENS, R. H. & SIMPSON, R. W. (1977). Presence of infectious polyadenylated RNA in the coronavirus avian bronchitis virus. Virology 77, 772–782.
- SCHUBERT, M., HARMISON, G. G. & MEIER, E. (1984). Primary structure of the vesicular stomatitis virus polymerase (L) gene: evidence for a high frequency of mutations. *Journal of Virology* **51**, 505–514.
- SIDDELL, S. G., ANDERSON, R., CAVANAGH, D., FUJIWARA, K., KLENK, H. D., MACNAUGHTON, M. R., PENSAERT, M., STOHLMAN, S. A., STURMAN, L. & VAN DER ZEIST, B. A. M. (1983a). Coronaviridae. Intervirology 20, 181–189.
- SIDDELL, S., WEGE, H. & TER MEULEN, V. (1983b). The biology of coronaviruses. Journal of General Virology 64, 761– 776.
- SOUTHERN, E. M. (1975). Detection of specific sequences among DNA fragments separated by gel electrophoresis. Journal of Molecular Biology 98, 503-517.
- SPAAN, W., DELIUS, H., SKINNER, M., ARMSTRONG, J., ROTTIER, P., SMEEKENS, S., VAN DER ZEIJST, B. A. M. & SIDDELL, s. G. (1983). Coronavirus mRNA synthesis involves fusion of non-contiguous sequences. *EMBO Journal* 2, 1839–1844.
- STADEN, R. (1982a). An interactive graphics program for comparing and aligning nucleic acid and amino acid sequences. Nucleic Acids Research 10, 2951–2961.
- STADEN, R. (1982b). Automation of the computer handling of gel reading data produced by the shotgun method of DNA sequencing. Nucleic Acids Research 10, 4731–4751.
- STADEN, R. (1984a). A computer program to enter DNA gel reading data into a computer. Nucleic Acids Research 12, 499–503.
- STADEN, R. (1984b). Graphic methods to determine the function of nucleic acid sequences. Nucleic Acids Research 12, 521–538.
- STADEN, R. (1984c). Measurements of the effects that coding for a protein has on a DNA sequence and their use for finding genes. Nucleic Acids Research 12, 551-567.
- STADEN, R. & MCLACHLAN, A. D. (1982). Codon preference and its use in identifying protein coding regions in long DNA sequences. Nucleic Acids Research 10, 141–157.
- STEINHAUER, D. A. & HOLLAND, J. J. (1986). Direct method for quantitation of extreme polymerase error frequencies at selected single base sites in viral RNA. Journal of Virology 57, 219-228.
- STERN, D. F. & KENNEDY, S. I. T. (1980a). Coronavirus multiplication strategy. I. Identification and characterisation of virus-specified RNA. Journal of Virology 34, 665–674.
- STERN, D. F. & KENNEDY, S. I. T. (1980b). Coronavirus multiplication strategy. II. Mapping the avian infectious bronchitis virus intracellular RNA species to the genome. *Journal of Virology* 36, 440–449.
- STERN, D. F. & SEFTON, B. M. (1984). Coronavirus multiplication: the locations of genes for the virion proteins on the avian infectious bronchitis virus genome. Journal of Virology 50, 22–29.
- STRAUSS, E. G. & STRAUSS, J. H. (1983). Replication strategies of the single stranded RNA viruses of eukaryotes. Current Topics in Microbiology and Immunology 105, 1–98.
- STRAUSS, E. G., RICE, C. M. & STRAUSS, J. H. (1984). Complete sequence of the genomic RNA of Sindbis virus. Virology 133, 92-110.

(Received 21 August 1986)