# The Exceptionally Large Genome of Hendra Virus: Support for Creation of a New Genus within the Family *Paramyxoviridae*

LIN-FA WANG,* MENG YU, ERIC HANSSON, L. IAN PRITCHARD, BRIAN SHIELL,
WOJTEK P. MICHALSKI, AND BRYAN T. EATON

*CSIRO Livestock Industries, Australian Animal Health Laboratory, Geelong, Victoria 3220, Australia*

An outbreak of acute respiratory disease in Hendra, a suburb of Brisbane, Australia, in September 1994 resulted in the deaths of 14 racing horses and a horse trainer. The causative agent was a new member of the family *Paramyxoviridae*. The virus was originally called *Equine morbillivirus* but was renamed *Hendra virus* (HeV) when molecular characterization highlighted differences between it and members of the genus *Morbillivirus*. Less than 5 years later, the closely related *Nipah virus* (NiV) emerged in Malaysia, spread rapidly through the pig population, and caused the deaths of over 100 people. We report the characterization of the HeV L gene and protein, the genome termini, and gene boundary sequences, thus completing the HeV genome sequence. In the highly conserved region of the L protein, the HeV sequence GDNE differs from the GDNQ found in almost all other nonsegmented negative-strand (NNS) RNA viruses. HeV has an absolutely conserved intergenic trinucleotide sequence, 3′-GAA-5′, and highly conserved transcription initiation and termination sequences similar to those of respiroviruses and morbilliviruses. The large genome size (18,234 nucleotides), the unique complementary genome terminal sequences of HeV, and the limited homology with other members of the *Paramyxoviridae* suggest that HeV, together with NiV, should be classified in a new genus in this family. The large genome of HeV also fills a gap in the spectrum of genome sizes observed with NNS RNA virus genomes. As such, it provides a further piece in the puzzle of NNS RNA virus evolution.

Although they manifest diverse biological properties, viruses in the families *Filoviridae*, *Paramyxoviridae*, *Rhabdoviridae*, and *Bornaviridae* all contain a nonsegmented negative-strand (NNS) RNA genome and share features of genome organization. These facts, together with similarities in domain structure and sequence of the viral polymerase proteins, suggest a close phylogenetic relationship. The four families are now grouped taxonomically in the order *Mononegavirales*, the first taxon above family level to be recognized in virus taxonomy (23, 25). The genome size of viruses in the order varies significantly, ranging from 8.9 kb in the *Bornaviridae* to 19.1 kb in the *Filoviridae*. Members of the *Rhabdoviridae* and *Paramyxoviridae* have intermediate genome sizes, 10.8 to 14.9 kb and 15.1 to 15.9 kb, respectively. Two interesting observations can be made from the comparison of genome sizes. First, there is no overlap of genome size between virus families. Second, genome size ranges differ significantly between the two families in which multiple genera have been defined, the *Rhabdoviridae* and *Paramyxoviridae*. Within the *Rhabdoviridae*, genome length can vary more than 40%, whereas variation within the *Paramyxoviridae* is no more than 5%. Thus, paramyxoviruses, especially those in the subfamily *Paramyxovirinae*, have traditionally been described as having a "uniform genome size" (23, 27). The universality of this feature is now challenged with the discovery, reported here, of a much larger genome for Hendra virus (HeV).

Members of the family *Paramyxoviridae* include highly contagious human and animal pathogens such as human parainfluenza viruses, *Measles virus*, *Canine distemper virus*, *Rinderpest virus*, *Mumps virus*, *Newcastle disease virus* (NDV), *Human respiratory syncytial virus*, and *Turkey rhinotracheitis virus*. Clas-

sification within the family has undergone major changes in recent years, and the current taxonomy (17, 24, 27) divides the family into two subfamilies, *Paramyxovirinae* and *Pneumovirinae*. The *Paramyxovirinae* include three genera, *Respirovirus* (formerly known as *Paramyxovirus*), *Morbillivirus*, and *Rubulavirus*, whereas the *Pneumovirinae* contains two genera, *Pneumovirus* and *Metapneumovirus*.

HeV was the causative agent of an explosive outbreak of a respiratory disease that resulted in the deaths of 14 horses and one human in a 2-week period in September 1994 in Hendra, a suburb of Brisbane, Australia (19). The virus was also responsible for a fatal human case of encephalitis in 1995, the infection almost certainly being acquired during necropsy of two horses that had died as a result of HeV infection 13 months previously (20). In January 1999, an additional fatal equine case was reported in North Queensland (13). Serological surveys and virus isolation studies indicated that flying foxes (fruit bats) in the genus *Pteropus* are likely to be the natural host of this new virus (11, 16, 34).

In March 1999, a virus closely related to HeV emerged in Malaysia, spread rapidly via the respiratory route through the pig population, and caused the death by encephalitis of over 100 people. Efforts to control the spread of the pathogen, Nipah virus (NiV), included the culling of over 1 million pigs (5, 6). NiV is closely related to HeV, and antibodies raised against one virus can neutralize the other in serum neutralization tests, albeit with reduced efficiency (5, 12). Positive antibody responses to NiV have been recorded in the Malaysian fruit bat population (8).

In addition to these two viruses, several other newly emerged *Mononegavirales* members of bat origin have been identified. These include *Australian bat lyssavirus* (9) and *Menangle virus* (21). Australian bat lyssavirus is closely related to rabies virus and was responsible for the death of a bat handler in 1996 (1). Menangle virus caused fetal death and abortion in pigs and respiratory disease in humans (4, 21). It

---

* Corresponding author. Mailing address: CSIRO Livestock Industries, Australian Animal Health Laboratory, PO Bag 24, Geelong, Victoria 3220, Australia. Phone: 61-3-52275121. Fax: 61-3-52275555. E-mail: linfa.wang@li.csiro.au.

TABLE 1. Mass spectrometry analysis of Lys-C-digested peptides of HeV L protein

| Digest fragment(s) detected[a] | Mass (MH$^+$) (Da) | | Amino acid sequence | Residues[b] |
|---|---|---|---|---|
| | Observed | Calculated | | |
| 1 | 2,770.64 | 2,770.34 | MAHELSISDIIYPECHLDSPIVSGK | 1–25[c] |
| 2 | 2,167.25 | 2,167.13 | LISAIEYAQLRHNQPNGDK | 26–44 |
| 3 | 873.54 | 873.51 | RLTENIK | 45–51 |
| 5 | 2,594.66 | 2,594.43 | RRSVYISRQSRLGNYIRDNIK | 58–78 |
| 8 | 694.42 | 694.41 | SLFSLK | 94–99 |
| 13 | 950.54 | 950.54 | AYNIVSRK | 117–124 |
| 14 | 2,171.33 | 2,171.19 | IIEMLQNITRNLITQDQK | 125–142 |
| 15 | 1,949.11 | 1,949.00 | DEVLGIYEQDRLSNIGK | 143–159 |
| 18 | 701.65 | 701.40 | NSQKPK | 186–191 |
| 31 | 1,065.84 | 1,065.57 | VLEYGPIMK | 386–394[c] |
| 47 | 883.49 | 883.45 | DEHELLK | 587–593 |
| 48 | 3,098.81 | 3,098.51 | TLFQLSISSVPRGNSQGRDSEFSNNTEK | 594–621 |
| 49 | 660.44 | 660.43 | SLISLK | 622–627 |
| 53, 54 | 2,149.17 | 2,148.08 | HNILPNTRNRHK | 659–676 |
| 55 | 1,693.87 | 1,693.76 | TFLDYHMEFSPYK | 677–689[c] |
| 55 | 1,677.77 | 1,677.76 | TFLDYHMEFSPYK | 677–689 |
| 56 | 1,690.05 | 1,690.73 | SDRMDRTETSDFSK | 690–703[c] |
| 63 | 2,106.11 | 2,106.01 | HINLDDTPEDDIFIHSPK | 778–795 |
| 64 | 938.48 | 938.46 | GGIEGYSQK | 796–804 |
| 66 | 967.58 | 967.54 | VHPNLPYK | 842–849 |
| 80 | 1,942.23 | 1,942.10 | NITARTILRTSPNPMLK | 1074–1090[c] |
| 80 | 1,926.35 | 1,926.10 | NITARTILRTSPNPMLK | 1074–1090 |
| 81 | 716.37 | 716.37 | GLFHDK | 1091–1096 |
| 83 | 843.59 | 843.54 | GLIRSGLK | 1143–1150 |
| 85 | 629.35 | 629.36 | SGIQPK | 1152–1157 |
| 86 | 2,757.65 | 2,757.31 | LVSRLSNHDYNQFLILNRLLSNK | 1158–1180 |
| 91 | 2,362.33 | 2,362.16 | EHSSIRVPYVGSSTDERSDIK | 1263–1283 |
| 94 | 2,628.77 | 2,628.40 | AITPVSTSNNLSHRLRDRSTQFK | 1335–1357 |
| 101 | 1,315.80 | 1,315.78 | TVAQTVLEIITK | 1498–1509 |
| 107 | 1,467.91 | 1,467.84 | VLSNALSHPRVFK | 1587–1599 |
| 117 | 578.13 | 578.32 | STSRK | 1790–1794 |
| 118 | 764.47 | 764.43 | VFNLGSK | 1795–1801 |
| 120 | 1,547.89 | 1,546.78 | YRRIGLNSSSCYK | 1809–1821 |
| 128 | 2,493.35 | 2,493.25 | VLDHSYLSDEINDQGITSVIFK | 2032–2053 |
| 133 | 889.48 | 889.52 | LLMQAGLK | 2088–2095[c] |
| 133 | 873.54 | 873.52 | LLMQAGLK | 2088–2095 |
| 136 | 503.22 | 503.33 | TRVK | 2151–2154 |
| 138 | 922.48 | 922.56 | VTVYSLIK | 2164–2171 |
| 141 | 963.54 | 963.53 | SPELYNIK | 2177–2184 |
| 150 | 877.49 | 877.54 | IIGYLSLV | 2237–2244 |

[a] Chronological order of theoretical Lys-C digest fragments of the L protein.
[b] Position of the peptide within the L-protein sequence.
[c] Additional mass of 16 Da is due to an oxidized methionine residue.

appears to be a member of the *Rubulavirus* genus (M. Westenberg, personal communication). A new member of the *Paramyxoviridae* has recently been isolated from bat urine in Malaysia, and it displays some antigenic cross-reactivity with Menangle virus (K. Chua, personal communication).

Tidona et al. (30) reported the isolation and characterization of a novel virus, *Tupaia paramyxovirus* (TPMV), from tree shrews, and Renshaw et al. (26) recently published the molecular characterization of the *Salem virus*, yet another novel paramyxovirus isolated from horses. These two new viruses are phylogenetically related to each other and to HeV and morbilliviruses. The isolation of seven new viruses, at least four of which are zoonotic and five of which appear to have originated from fruit bats, opens a new and exciting era in the investigation of the natural history of *Paramyxoviridae* and NNS RNA viruses in general.

In this paper, we report the molecular characterization of the HeV L gene, which encodes the RNA polymerase, and determine the sequence of the genome termini and gene boundaries and thus complete the sequence of the largest genome in the *Paramyxoviridae* to be described. Important molecular features will be summarized to support the establishment of a new genus for HeV and NiV within the subfamily *Paramyxovirinae*.

## MATERIALS AND METHODS

**Genome sequencing.** The construction and screening of a cDNA library from purified viral genomic RNA have been described previously (33). Briefly, HeV was purified by zonal centrifugation in sucrose gradients, and genomic RNA was isolated using standard methods. The TimeSaver cDNA Synthesis Kit (Pharmacia) was used to make total cDNA by using random priming followed by the addition of an *Eco*RI adaptor and cloning into the pZEr0-1 vector (Invitrogen). A genome walking strategy was used to isolate specific clones by colony hybridization. DNA sequencing was performed by a combination of manual sequencing with the Sequenase Kit (USB) and automatic sequencing using the Big-Dye Dideoxyl Termination Cycle Sequencing kit (Pharmacia) and the ABI 377 Sequencer. Each nucleotide position was sequenced at least twice, either by sequencing overlapping clones or by sequencing the opposite strand of the same clone. Sequence confirmation by direct sequencing of PCR fragments without cloning was also carried out for several regions of importance. Sequences were analyzed and aligned using the software package, including Clone Manager 5 and Align Plus 4, from S & E Software (Durham, N.C.) as described previously (33).

TABLE 2. Comparison of L gene features of HeV and selected *Paramyxovirinae* members

| Virus (accession no.)[a] | Length of gene element or product | | | |
|---|---|---|---|---|
| | mRNA (nt) | 5′ UTR (nt) | 3′ UTR (nt) | Protein (aa) |
| HeV (AF017149) | 6,955 | 153 | 67 | 2,244 |
| TPMV (AF079780) | 6,927 | 34 | 80 | 2,270 |
| *Morbillivirus* | | | | |
| MeV (AB016162) | 6,643 | 22 | 69 | 2,183 |
| CDV (AF014953) | 6,642 | 22 | 65 | 2,184 |
| *Respirovirus* | | | | |
| SeV (M19661) | 6,800 | 28 | 85 | 2,228 |
| PIV3 (Z11575) | 6,795 | 22 | 71 | 2,233 |
| NDV (AF077761) | 6,703 | 11 | 77 | 2,204 |
| *Rubulavirus* | | | | |
| MuV (AB000388) | 6,930 | 8 | 136 | 2,261 |
| SV5 (AF052755) | 6,800 | 8 | 34 | 2,255 |

[a] See the legend to Fig. 1 for abbreviations.

**Cloning and sequencing of genome termini.** Two different methods were used to determine genome end sequences. (i) For inverse PCR, purified viral genomic RNA was denatured at 100°C for 40 s in the presence of deionized formamide and digested with tobacco acid pyrophosphatase (Epicentre Technologies) according to the manufacturer's instructions. The RNA was then phenol-chloroform extracted and ethanol precipitated at −70°C for 60 min. The 5′ and 3′ ends of the genomic RNA were ligated using RNase-free T4 RNA Ligase (Promega) according to the manufacturer's instructions. Genomic RNA was again phenol-chloroform extracted and ethanol precipitated. First-strand cDNA synthesis was performed with genome-specific primers by using the SuperScript preamplification system (Life Technologies). PCR fragments incorporating the 5′ and 3′ ends of the viral genome were amplified and sequenced using genome specific primers. (ii) For rapid amplification of cDNA ends, first-strand cDNA synthesis was performed using SuperScript II (Life Technologies) and virus-specific primer, followed by ligation of anchor (annealed from the top strand, 5′-CTAAT AC GAC TCACT ATAGG GCTCG AGCGC CCGCC CGGGC AGGT-3′, and the 5′-phosphorylated bottom strand, 5′-ACCTG CCC-3′) using T4 RNA ligase (New England Bio-Labs). Subsequent PCR amplifications were carried out using a combination of virus-specific primers and two nested primers annealing to the anchor (AP1, 5′-GGATC CTAAT ACGAC TCACT ATAGG GC; AP2, 5′-AATAG GGCTC GAGCG GC). PCR products were then cloned as blunt-end fragments into the *Eco*RV site of pZEr0-1 vector (Invitrogen) for sequence determination. A total of 24 independent clones were sequenced.

**Production of monospecific antibodies in rabbits.** A recombinant polypeptide corresponding to amino acid residues 245 to 517 of the deduced L protein was expressed in *Escherichia coli* by using the pRSET vector as described previously (32). Approximately 1.5 mg of the recombinant protein was purified using preparative sodium dodecyl sulfate-polyacrylamide gel electrophoresis (SDS-PAGE), and the band of interest was visualized in ice-cold 0.3 M KCl, cut out, and sliced into 2- by 2-mm pieces, followed by passive elution in 50 mM $NH_4HCO_3$ containing 0.1% SDS. Eluted protein solution (0.5 ml at 0.5 mg/ml) was mixed with an equal volume of Freund's incomplete adjuvant and used to immunize rabbits. The same injection was repeated 4 weeks later, followed by a third injection 6 weeks later without adjuvant. Antibody titers and specificity were determined by enzyme-linked immunosorbent assay and Western blotting using both recombinant antigens and purified virus proteins.

**Peptide mapping by in situ enzymatic digestion and mass spectrometry.** An adaptation of the method described by Moritz et al. (18) was employed to obtain internal amino acid sequence and mass spectrometry data. Approximately 200 μg of purified virus was inactivated at 100°C for 5 min with 2% SDS and 3% dithiothreitol. Following electrophoresis, a stained band identified as the L protein (approximately 200 kDa) was excised from the gel and fragments were cut into 2-mm pieces and washed with 2 M $NH_4HCO_3$–50% acetonitrile for 30 min. Gel pieces were dried and rehydrated in 50 mM Tris-HCl–1 mM EDTA–6% acetonitrile (pH 8.5). Endoproteinase Lys-C (sequencing grade; Wako) was added at an enzyme-to-substrate ratio of approximately 1:10, and gels were incubated overnight at 37°C. After centrifugation, supernatant fluids containing peptides were removed. Gel pieces were reextracted with 1% trifluoroacetic acid for 30 min with sonication, and after centrifugation, gel pieces were washed again with 0.05% trifluoroacetic acid–80% acetonitrile for 30 min. Supernatant fluids were combined and concentrated, and peptides were separated using reverse-phase chromatography as previously described (33). Electrospray ionization mass spectrometry (ESI-MS) analysis of Lys-C digests was performed on a Finnigan LCQ instrument with a Hewlett-Packard 1090 HPLC. The Pros-

pector MS-FIT program was set to consider the following modifications: phosphorylation of serine, threonine, and tyrosine; oxidation of methionine; N-terminal acetylation of unmodified cysteine; and 200.00 ppm mass tolerance.

**Nucleotide sequence accession number.** The sequence reported in this paper has been deposited in the GenBank database (accession no. AF017149).

## RESULTS

**Sequence analysis of the L gene and protein.** A total of 21 overlapping cDNA clones covering the entire coding region of L gene were isolated from the cDNA library. Highly conserved transcriptional start and stop signals were identified, and the size of L gene mRNA was predicted to be 6,955 nucleotides (nt). Sequence analysis revealed a single long open reading frame coding for a large protein of 2,244 amino acids (aa) with a calculated molecular mass of 257,280 Da. The identity of the L gene was confirmed using two different approaches. Firstly, a monospecific antiserum raised to a portion of the L protein (aa 245 to 517) expressed in *E. coli* was used in Western blotting to identify the protein in purified virus. The antiserum specifically reacted with a protein whose apparent molecular mass exceeded 200 kDa (data not shown). Secondly, peptides derived from the putative L protein band by Lys-C digestion and separated by reverse-phase chromatography were analyzed by mass spectrometry. A total of 47 masses was identified, and a majority (59%) matched those of the L protein sequence deduced from cDNA clones (Table 1).
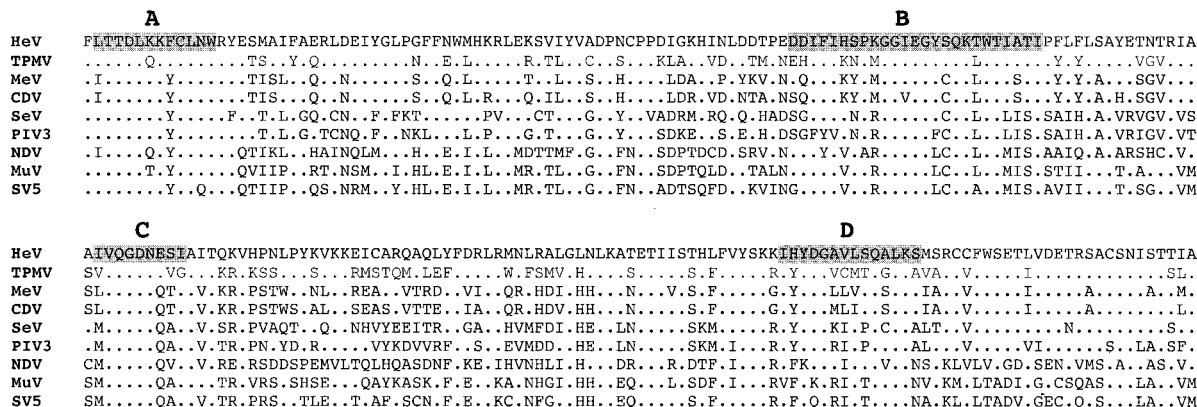
The first AUG codon in the HeV L gene was located at nt 24 in the same reading frame as that encoding the L protein. However, there were eight in-frame stop codons existing between this and the second AUG codon at nt 154. The use of this second AUG codon as the initiation site of HeV L gene translation was supported by the conserved N-terminal amino acid sequence of the deduced L protein in comparison with those of morbilliviruses and respiroviruses (data not shown). The 5′ untranslated region (UTR) of the HeV L gene was therefore 153 nt in length (Table 2). On the other hand, the 3′ UTR of the HeV L gene (67 nt) was very similar in size to those of morbilliviruses (varying from 65 to 69 nt) and respiroviruses (71 to 85 nt).

The HeV L protein was rich in leucine and isoleucine (19.3%, which is more than 1.3 times that of an average protein) and carried a net positive charge of +37 at neutral pH. The linear domain structure of L proteins previously suggested by Poch et al. (22) could also be identified in the HeV L protein. A sequence comparison of the most conserved domain (domain III) and the four conserved motifs therein is shown in Fig. 1A. The C motif in domain III of HeV L had a sequence of GDNE instead of GDNQ, a surprising finding in view of the fact that GDNQ is conserved among all NNS RNA virus L genes sequenced so far, the sole exception being the TPMV L protein (see Fig. 1B for a more detailed and expanded comparison).

Phylogenetic analysis of the full-length L protein sequences or the domain III sequences produced very similar trees, indicating that the HeV L protein is most closely related to, but distinct from, viruses in the *Morbillivirus* and *Respirovirus* genera in the subfamily *Paramyxovirinae* (data not shown).

**Sequences of the 3′ and 5′ termini of the HeV genome.** The genome-end sequences were determined using a combination of two different approaches, and the results are summarized in Fig. 2A. The termini of the HeV RNA genome and antigenome are identical for the first 12 nt and 19 of the first 23 nt. A comparison of the 3′ leader sequences of selected *Paramyxovirinae* viruses is shown in Fig. 2B. Overall, the sequence of HeV is more closely related to those of respiroviruses and morbilliviruses than rubulaviruses. However, HeV is unique in
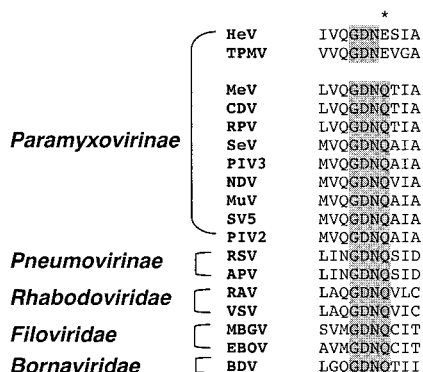
FIG. 1. (A) Alignment of four polymerase motifs (A to D) located within domain III of L proteins from selected *Paramyxovirinae* members. Residues identical in all viruses are shown as dots. Abbreviations: HeV, Hendra virus; TPMV, Tupaia paramyxovirus; MeV, measles virus; CDV, canine distemper virus; SeV, Sendai virus; PIV3, human parainfluenza virus type 3; NDV, Newcastle disease virus; MuV, mumps virus; SV5, simian virus 5. See Table 4 for GenBank accession numbers for the sequences used in the alignment. (B) A more extensive alignment of the motif C sequence to include viruses from all four families of the order *Mononegavirales*. The unique E residue in the GDNE sequence of HeV and TPMV L proteins is indicated by an asterisk. Abbreviations and GenBank accession numbers for viruses not listed above are as follows: RPV, rinderpest virus (Z30697); PIV2, human parainfluenza virus type 2 (X57559); RSV, respiratory syncytial virus (U39661); APV, avian paramyxovirus (U65312); RAV, rabies virus (M31046); VSV, vesicular stomatitis virus (J02428); MBGV, Marburg virus (Z12132); EBOV, Ebola virus (AF086833); BDV, Borna disease virus (U04608).

that it has a G residue at position 4 in lieu of an A residue in all the other viruses. At positions 5, 10, and 12, HeV resembles members of the *Respirovirus* genus (A, G, and G, respectively) and differs from viruses in the *Morbillivirus* genus (G, A, and T, respectively). The converse is true at position 11, where a G residue is shared by HeV and morbilliviruses while respiroviruses have an A residue there.

The HeV 3′ leader sequence was 55 nt, a length identical to those of the 3′ leader sequences of other members of the subfamily *Paramyxovirinae* (Fig. 2B). The 5′ trailer sequence of HeV was 33 nt. Unlike that of the 3′ leader, the size of the 5′ trailer can vary quite substantially within the subfamily (see Table 4).

**Gene start, stop, and intergenic sequences.** For members of the family *Paramyxoviridae*, transcription of individual genes is carried out by a stop-and-reinitiation mechanism that is controlled by conserved transcriptional sequences at the gene borders. Sequence comparison analyses revealed that HeV is similar to members of the *Respirovirus* and *Morbillivirus* genera in having a conserved intergenic trinucleotide sequence, 3′-GAA-5′. The transcriptional start and stop signals of all HeV genes together with the consensus sequences are shown in Table 3.

Comparison of the HeV consensus sequences with those of the respiroviruses and morbilliviruses suggests that while the consensus sequence of HeV gene start signals is similar to the consensus sequences of the start signals of both, the HeV gene stop signals are more like those of respiroviruses.

**Rule of six and other genomic features.** It has been suggested that *Paramyxovirinae* genomes are replicated efficiently only when they are a multiple of 6 nt in length, and this has been dubbed the "rule-of-six" (3). The genome length of HeV (18,234 nt) is a multiple of 6 and does conform to this rule (Table 4). Also listed in Table 4 are the subunit hexamer-phasing positions of the transcription initiation site for each of the six genes and the P gene-editing site from selected representative members of the subfamily *Paramyxovirinae*. Several interesting observations can be made. (i) Members of each of the three existing genera seem to have conserved hexamer-phasing positions for most of the genes. (ii) HeV has a pattern that is significantly different from that of any of the other members of the subfamily. (iii) Although phylogenetic analysis suggested relatedness of TPMV and HeV (30), TPMV seems to be more closely related to morbilliviruses with respect to the hexamer-phasing positions.

**A.**



|  |  | 10 | 20 | 30 | 40 | 50 |
| --- | --- | --- | --- | --- | --- | --- |

Antigenome    5'-ACCGAACAAGGGGCAAATATGGATACGTGTTAAAAAACTGCGTATGTTTAAAACTT
Genome        5'-ACCGAACAAGGGTAAAGAGAGATCGTTATTAAG

**B.**



|  | * | 10 | 20 | 30 | 40 | 50 | ▽ 60 |

HeV   ACCGAACAAGGG GAAATATGGATACGTGTTAAAAAACTGCGTATGTTTAAAACTTAGGAA
TPMV  ACCACAAAAGGG TGGTCACGGGGTCGTAAGTTTCAAAGAAAATTATTTATGGCTTAGGAA

*Morbillivirus* [
MeV   ACCAAACAAAGT TGGGTAAGGATAGATCAATCAATGATCATATTCTAGTACACTTAGGAT
CDV   ACCAGACAAAGT TGGCTAAGGATAGTTAAATTATTGAATATTTTATTAAAAACTTAGGGT
]

*Respirovirus* [
SeV   ACCAAACAAGAG AAGAAACTTGTTTGGAATATATAATGAAGTTAGACAGGATTTTAGGGT
PIV3  ACCAAACAAGAG AAGAAACTTGTCTGGGAATATAAATTTAACTTTTAAATTAACTTAGGAT
]

NDV   ACCAAACAGAGA ATCCGTGAGTTACGATAAAAGGCGAAGGAGCAATTGAAGTCGCACGGG

*Rubulavirus* [
MuV   ACCAAGGGGAAA AAGAAGATGGGATGTTGGTAGAACAAATAGTGTAAGAAACAGTAGCC
SV5   ACCAAGGGGAAA ATGAAGTGGTGACTCAAATCATCGAAGACCCTCGAGATTACATAGGTC
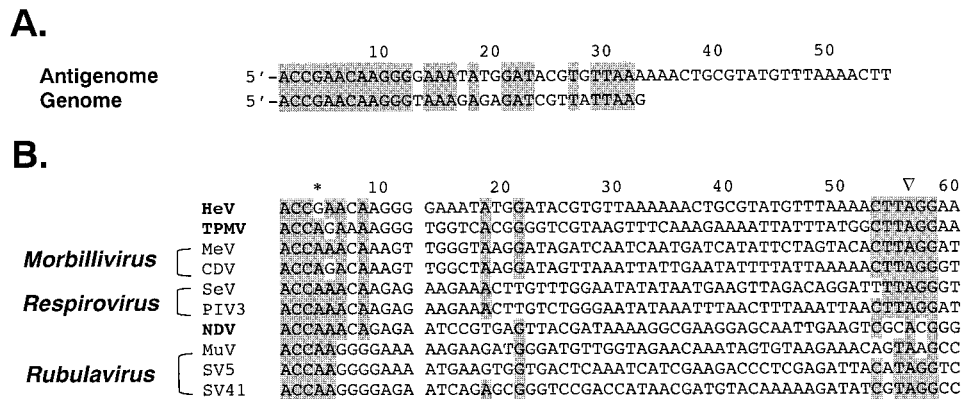SV41  ACCAAGGGGAGA ATCAGAGCGGGTCCGACCATAACGATGTACAAAAAGATATCGTAGGCC
]

FIG. 2. Alignment of genome end sequences. (A) Alignment of the 3' leader and 5' trailer sequences of HeV. (B) Alignment of 3' leader sequences of selected *Paramyxovirinae* members. The unique HeV residue at position 4 is indicated by an asterisk. The downward pointing triangle indicates the N gene transcription start site. Abbreviations are given in the Fig. 1 legend; SV41, simian virus 41 (GenBank accession no. X64275).

While the genome size of HeV is much larger than those of most other *Paramyxoviridae* members, the sizes of its encoded six major proteins are very similar to those of others in the family with the exception of the P protein, which is approximately 100 aa longer (33). The increase in genome size is mainly due to the expansion of UTRs, especially at the 3' end of the mRNA with the exception of the L gene (Fig. 3). This has resulted in a significant drop in the coding percentage from an average of approximately 92% for known *Paramyxovirinae* viruses to 82.1% for HeV. A similar observation has been made for TPMV following the completion of sequencing of the whole genome. Filoviruses are a further example of NNS viruses that have a low coding percentage, 71.6% for Ebola virus and 76.5% for Marburg virus. On the other hand, the small genome of Borna disease virus has a very high coding percentage, 95.2%.

## DISCUSSION

The genome structure of HeV resembles that of other *Paramyxovirinae* members in having six genes in the order 3'-N-P/V/C-M-F-G-L-5'. The size of each gene product is similar to that of other viruses in the subfamily with the exception of the P protein, which is approximately 100 aa longer than the longest P protein previously characterized (33). Analyses of the protein sequences of the first five genes strongly suggest that HeV represents a new evolutionary lineage within the subfamily *Paramyxovirinae* with a closer relationship to members of

the *Respirovirus* and *Morbillivirus* genera than the *Rubulavirus* genus (10, 33, 35, 36). This is corroborated by the present study of the L gene sequence. RNA polymerase (L) genes of NNS RNA viruses have been considered a reliable target for phylogenetic analysis due to their high degree of conservation during evolution (22). The size, sequence, and domain structure of the HeV L protein are very similar to those of other viruses in the subfamily. Extensive phylogenetic analyses of both full-length L protein and the more conserved domain III sequences confirmed that HeV is a member of the subfamily *Paramyxovirinae* and is more closely related to morbilliviruses and respiroviruses (data not presented).

However, the L protein of HeV and TPMV differs from other NNS virus L proteins in one important respect. The 4-aa sequence GDNQ that is absolutely conserved among all known NNS virus L genes, from the smallest genome of Borna disease virus to the largest genome of Marburg virus, is replaced in these new viruses by GDNE (Fig. 1). This sequence resides in the highly conserved C motif of domain III of NNS virus RNA polymerases and is believed to be the core polymerase motif in which the GD dipeptide resides precisely within a β-turn-β structure (22). The functional importance of the GDNQ sequence in the L protein has been examined experimentally for rabies virus and vesicular stomatitis virus (VSV). Changing the rabies virus L protein GDN sequence to GDD (the motif shared by most positive-strand RNA viruses) or SDD (the motif observed among segmented negative-strand RNA viruses) completely abolished RNA polymerase activity (28). Replacement of GDNQ in rabies virus by GDNE generated a mutant L protein that retained less than 1% of wild-type RNA polymerase activity. In the case of the VSV L protein, changing GDN to GDD did not completely inactivate the RNA polymerase but reduced its activity by 73%. However, alteration of GDNQ to GDNN reduced activity by more than 95% (29). These studies suggested that the invariant GDNQ sequence is optimal for RNA polymerase activity in the order *Mononegavirales* and indicated that the effect of mutation in the region varies from virus to virus. Following the establishment of an HeV reverse genetics system, it will be interesting to determine whether mutation of GDNE to the more conserved GDNQ will have any effect on HeV RNA polymerase activity.

In the genus *Morbillivirus*, the length of the 5' UTR for the L gene is 21 nt. In the genus *Respirovirus*, this varies from 22 nt for human parainfluenza virus type 3 to 28 nt for Sendai virus. The HeV L gene has a 5' UTR of 154 nt, which is the longest

TABLE 3. Gene start, stop, and intergenic sequences

| Virus | Gene(s) | Gene boundary sequence |
| --- | --- | --- |
| HeV | /N | CTT AGGAACCAAG |
|  | N/P | ATTAAGAAAAA CTT AGGATCCAAG |
|  | P/M | ATTAAGAAAAA CTT AGGAGACAGG |
|  | M/F | ATTAAGAAAAA CTT AGGAGCCAAG |
|  | F/G | TTTACAAAAA CTT AGGACCCAAG |
|  | G/L | ATTAAGAAAAA CTT AGGACCCAAG |
|  | L/ | ATTAAGAAAAA CTT |
| Consensus sequences |  |  |
| HeV |  | WTTAMRAAAAA CTT AGGAnMCARG |
| Morbilliviruses |  | VHHWHDnAAAA CKT AGGRnRMARG |
| Respiroviruses |  | ADWAHVAAAAA CYY AGGRnnAAHG |

TABLE 4. Comparison of genome properties of HeV and selected *Paramyxovininae* viruses

| Virus[a] | Accession no. | Genome length (nt) | Coding region (%) | 5' trailer (nt) | Rule of six (length/6) | Hexamer-phasing positions[b] | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | N | P(e) | M | F | A | L |
| HeV | AF017149 | 18,234 | 82.1 | 33 | 3,039 | 2 | 3 (5) | 4 | 4 | 4 | 3 |
| TPMV | AF079780 | 17,904 | 82.5 | 590 | 2,984 | 2 | 2 (6) | 1 | 3 | 3 | 2 |
| *Morbillivirus* | | | | | | | | | | | |
| MeV | AB016162 | 15,894 | 89.1 | 41 | 2,649 | 2 | 2 (6) | 4 | 3 | 3 | 2 |
| CDV | AF014953 | 15,690 | 92.2 | 41 | 2,616 | 2 | 2 (2) | 4 | 2 | 3 | 2 |
| *Respirovirus* | | | | | | | | | | | |
| SeV | M19661 | 15,384 | 93.8 | 58 | 2,564 | 2 | 1 (1) | 1 | 1 | 1 | 2 |
| PIV3 | Z11575 | 15,462 | 94.0 | 45 | 2,577 | 2 | 1 (2) | 1 | 1 | 1 | 2 |
| NDV | AF077761 | 15,185 | 90.6 | 114 | 2,531 | 2 | 4 (1) | 4 | 4 | 3 | 6 |
| *Rubulavirus* | | | | | | | | | | | |
| MuV | AB000388 | 15,384 | 93.1 | 25 | 2,564 | 2 | 1 (3) | 1 | 6 | 1 | 6 |
| SV5 | AF052755 | 15,246 | 92.0 | 32 | 2,541 | 2 | 1 (3) | 1 | 6 | 1 | 6 |

[a] Abbreviations are given in the legend to Fig. 1.
[b] The number given is the subunit hexamer-phasing position at each of the six gene start sites and, for P gene, the editing site (in parentheses). N, nucleoprotein; P, phosphoprotein; M, matrix protein; F, fusion protein; A, attachment protein, i.e., HN, H, or G protein depending on the virus; L, L protein or RNA polymerase.

to be found within the subfamily *Paramyxovirinae*. The 3' UTR of the HeV L gene is 67 nt, which is similar to those of morbillivirus L genes (Table 2). It is interesting to note that the L gene is the only gene in the HeV genome that does not have a longer 3' UTR than other *Paramyxovirinae* members. In that regard, it is also interesting to point out that the 5' UTR of the HeV N gene was similar to those of other *Paramyxovirinae* members (36). Considering the vast variation in the lengths of UTRs for internal HeV genes (10, 33, 35, 36), it is tempting to suggest that there might be selective pressure to maintain the sizes of not only the genome leader and trailer sequences in HeV (see below) but also the genome end proximal UTRs, i.e., the 5' UTR of the N gene and the 3' UTR of the L gene.

The genome terminal sequences of paramyxoviruses are highly conserved, and there is complementarity between the 3'- and 5'-terminal sequences. These conserved terminal sequences, especially the first 12 nt, are thought to contain the genome and antigenome promoters and are largely genus specific (2, 15). This feature has been used as an important criterion in classification of the *Paramyxoviridae*. For example, in a recent analysis, the difference between the genome-end sequences of NDV and other members of the genus *Rubulavirus*

has been cited as one of the main reasons for reclassification of NDV in a separate genus within the *Paramyxovirinae* (7). In the case of HeV, the difference between its genome end sequences and those of other subfamily members is obvious and significant, especially the presence of a G residue at position 4 in lieu of the A residue found in all other members of the subfamily. On the other hand, the size of the HeV 3' leader sequence, 55 nt, is identical to those found for all members of the subfamily *Paramyxovirinae*, a fact consistent with the placing of HeV within the subfamily based on phylogenetic analyses. The length of the 5' trailer sequence varies significantly within the subfamily, with TPMV having the longest at 590 nt and simian virus 41 (SV41) the shortest at 20 nt. In general, rubulaviruses have a shorter 5' trailer sequence (20 to 32 nt), whereas morbilliviruses and respiroviruses tend to have slightly longer sequences of 41 to 45 and 45 to 58 nt, respectively. For HeV, the 5' trailer is 33 nt, which is shorter than those of morbilliviruses and respiroviruses.

Members of the family *Paramyxoviridae* not only have highly conserved genome end sequences for replication and transcription but also have highly conserved transcription start and stop sequences at each gene boundary (2, 15). For the *Respirovirus*
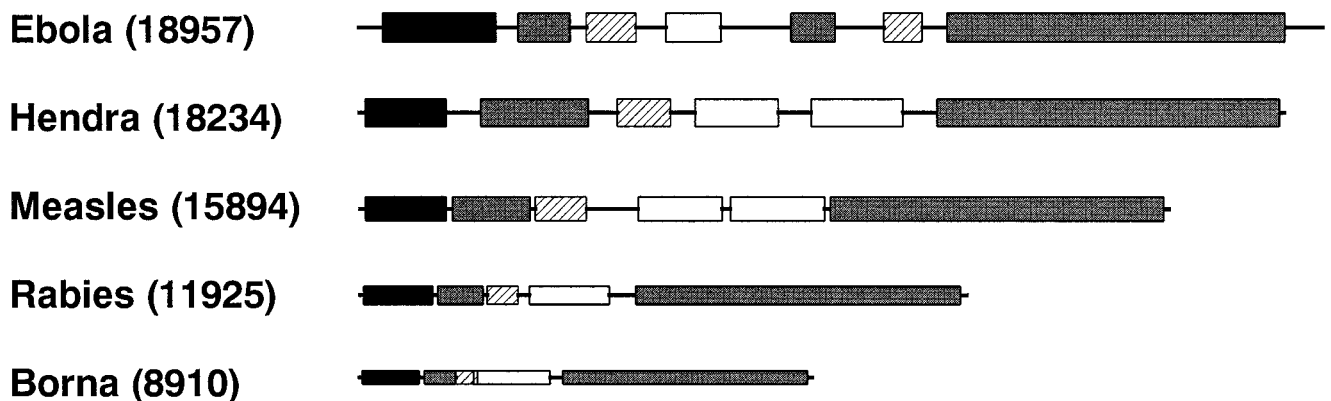


FIG. 3. Genome structure and size comparison of selected members from the four families within the order *Mononegavirales*. Genome size (in nucleotides) is given in parentheses for each virus. The sizes are relative to that of the Ebola virus. Only genes coding for major structural proteins are shown, using the following functional grouping: nucleocapsid proteins (solid box), proteins associated with RNA polymerase and/or ribonucleoprotein complex (shaded box), matrix proteins (hatched box), and membrane proteins (open box). See the Fig. 1 legend for GenBank accession numbers.

and *Morbillivirus* genera, the nontranscribed intergenic sequence that is found not only between genes but also before the N and after the L genes is 3 nt in length. Its sequence, 3′-GAA-5′ (or 5′-CTT-3′ in the antigenome), is conserved for most, but not all, of the six genes of viruses within these genera. HeV is the first virus to be sequenced that has the identical CTT sequence in all seven positions of a six-transcription-unit genome (Table 3). In morbilliviruses and respiroviruses, at least one "imperfect" intergenic sequence has been found in each of the viruses sequenced so far. These include trinucleotide sequences AAA, GTA, GGG, GCA, GTT, and TTT, and they are mainly located toward the 3′ end of the antigenome, i.e., at the end of the H (HN) gene or L gene (14), with the exception of Sendai virus, which has a TTT sequence before the N gene (Fig. 2B). The gene start and stop signals for the six HeV transcription units are also more conserved than their morbillivirus and respirovirus counterparts. Whether this near-perfect conservation of transcription signals represents a more ancient or modern configuration remains to be seen. A close examination of the complete genome sequence of TPMV, which became available only very recently (GenBank accession no. AF079780), indicated that TPMV also has perfect 5′-CTT-3′ intergenic sequences at all of the seven positions. Together with the larger genome sizes (17,904 nt for TPMV and 18,234 nt for HeV) and the unusual sequence GDNE at the putative catalytic site of their L proteins, this makes HeV and TPMV unique among *Paramyxovirinae* genomes that have been completely sequenced. Whether these unique features, i.e., perfect intergenic sequences and the GDNE sequence of L protein, are interrelated and common to all large-genome paramyxoviruses remains to be seen.

The template for paramyxovirus RNA synthesis is not naked RNA but the helical ribonucleoprotein core of the virus, a structure in which nucleotide hexamers are believed to be associated with individual nucleocapsid (N) protein molecules (14). As a consequence, many *Paramyxoviridae* members have a genome length that is a multiple of 6 nt. Moreover, it has been shown that these genomes are effectively replicated only when they are a multiple of 6 nt (14). In addition, it has been found that, within a genus, the transcription start site for each gene tends to be conserved in relation to the hexamer-phasing position as shown in Table 4 (7, 14). While the conservation of genome length in multiples of 6 nt may have a functional advantage, the functional role for a conserved hexamer-phasing position for each gene is less obvious. Nevertheless, this may serve as an additional molecular marker in studying virus evolution and the classification of new viruses. The hexamer-phasing position of the HeV N gene is the same as that of cognate genes in the subfamily as a result of the uniform length of 3′ leader sequences within the subfamily. However, it is intriguing to find that HeV has a hexamer-phasing pattern, "2, 3, 4, 4, 4, 3," that is significantly different from those of other viruses in the subfamily. The hexamer-phasing positions for the HeV P, A (G), and L genes are unique and not used by cognate genes in the subfamily. HeV is also unique in that it uses hexamer-phasing position 5 for the P editing site. This position is not used by any other virus for either a P editing site or transcription start site (14).

Previous studies based on RNA polymerase sequence analysis have revealed that, within the order *Mononegavirales*, the *Filoviridae* are more closely related to viruses in the genera *Respirovirus* and *Morbillivirus* of the family *Paramyxoviridae* than to viruses in the family *Rhabdoviridae* (31). It is interesting that the 3′ leader sequence of Ebola virus is 55 nt, the same length as that of all members of the subfamily *Paramyxovirinae*. The discovery that HeV has a genome approximately 15%

larger than other members of the *Paramyxoviridae* family and that each HeV gene has long UTRs is also consistent with a possible evolutionary relationship between the *Paramyxovirinae* and the *Filoviridae*. As pointed out previously, until the discovery of HeV, there had been a large genome size gap of approximately 3 kb between the largest genome of the *Paramyxoviridae* and the smallest genome of the *Filoviridae*. This gap has now been filled by the 18.2-kb genome of HeV. It remains to be seen whether NiV will have a larger genome size. The recent completion of the 17.9-kb TPMV genome sequence has provided a further reduction in the size gap of paramyxovirus genomes.

The data reported here on the L protein sequence, genome size, genome end sequences, hexamer-phasing positions, and gene start and stop signals, together with sequence and molecular features previously reported for the other five genes (10, 33, 35, 36), suggest that the newly identified HeV is a member of the subfamily *Paramyxovirinae*. The data also demonstrate a clear need for the creation of a separate genus within the subfamily to accommodate the many significant differences between HeV and other existing members of the *Paramyxovirinae*. The name *Henipavirus* is here proposed for a new genus that would include HeV as the type species and NiV as the second member of the genus. Although the complete genome structure of NiV is not yet available, the strong antigenic cross-reactivity between NiV and HeV (5) and high sequence homology for their N, P, M, F, and G genes (12) suggest that HeV and NiV are very closely related viruses.

Identification and characterization of HeV and related new viruses is important not only because a number of them are lethal to humans and/or domestic animals but also because of the insights they may provide into virus evolution. The pattern of genome organization, the correlation between genome size and virus family, and the lack of recombination suggest that virus families in the order *Mononegavirales* have evolved by expansion or deletion of UTRs and by gene duplication, rather than by introduction of genetic information from outside (25). In this regard, the exceptionally large genome of HeV reported here represents an important link in the evolution of *Mononegavirales* viruses. It is interesting that an emerging virus with a high profile due to its lethal infection of humans and a variety of animals has also significantly expanded the viral diversity within the *Paramyxoviridae*, which has raised more questions about the nature of virus evolution.

## REFERENCES

1. **Allworth, A., K. Murray, and J. Morgan.** 1996. A human case of encephalitis due to a lyssavirus recently identified in fruit bats. Commun. Dis. Intellig. **20:**504.
2. **Bellini, W. J., P. A. Rota, and L. J. Anderson.** 1998. Paramyxoviruses, p. 435–461. *In* L. Collier, A. Balows, and M. Sussman (ed.), Microbiology and microbial infections, vol. 1. Arnold, London, England.
3. **Calain, P., and L. Roux.** 1993. The rule of six, a basic feature for efficient replication of Sendai virus defective interfering RNA. J. Virol. **67:**4822–4830.
4. **Chant, K., R. Chan, M. Smith, D. E. Dwyer, and P. Kirkland.** 1998. Probable human infection with a newly described virus in the family Paramyxoviridae. Emerg. Infect. Dis. **4:**273–275.
5. **Chua, K. B., W. J. Bellini, P. A. Rota, B. H. Harcourt, A. Tamin, S. K. Lam, T. G. Ksiazek, P. E. Rollin, S. R. Zaki, W. J. Shieh, C. S. Goldsmith, J. T. Roehrig, B. Eaton, A. R. Gould, J. Olson, H. Field, P. Daniels, A. E. Ling, C. J. Peters, L. J. Anderson, and B. W. J. Mahy.** 2000. Nipah virus, a newly emergent deadly praramyxovirus. Science **288:**1432–1435.
6. **Chua, K. B., K. J. Joh, K. T. Wong, A. Kamarulzaman, P. S. K. Tan, T. G.**

Ksiazek, S. R. Zaki, G. Paul, S. K. Lam, and C. T. Tan. 1999. Fatal encephalitis due to Nipah virus among pig-farmers in Malaysia. Lancet 354:1257–1259.

7. de Leeuw, O., and B. Peeters. 1999. Complete nucleotide sequence of New-castle disease virus: evidence for the existence of a new genus within the subfamily Paramyxoviridae. J. Gen. Virol. 80:131–136.

8. Field, H., J. M. Yob, C. Morrissy, and P. Selleck. 1999. Nipah virus in Malaysia—preliminary wildlife surveillance, p. 187. Proceedings of the XIth International Congress of Virology, Sydney, Australia.

9. Fraser, G. C., P. T. Hooper, R. A. Lunt, A. R. Gould, L. J. Gleeson, A. D. Hyatt, G. M. Russell, and J. A. Kattenbelt. 1996. Encephalitis caused by a lyssavirus in fruit bats in Australia. Emerg. Infect. Dis. 2:327–331.

10. Gould, A. R. 1996. Comparison of the deduced matrix and fusion protein sequences of equine morbillivirus with cognate genes of the Paramyxoviridae. Virus Res. 43:17–31.

11. Halpin, K., P. Young, and H. Field. 1996. Identification of likely natural hosts for equine morbillivirus. Commun. Dis. Intellig. 20:476.

12. Harcourt, B. H., A. Tamin, T. G. Ksiazek, P. E. Rollin, L. J. Anderson, W. J. Bellini, and P. A. Rota. 2000. Molecular characterization of Nipah virus, a newly emergent paramyxovirus. Virology 271:334–349.

13. Hooper, P. T., A. R. Gould, A. D. Hyatt, M. A. Braun, J. A. Kattenbelt, S. G. Hengstberger, and H. A. Westbury. 2000. Identification and molecular characterisation of Hendra virus in a horse in Queensland. Aust. Vet. J. 78:281–282.

14. Kolakofsky, D., T. Pelet, D. Garcin, S. Hausmann, J. Curran, and L. Roux. 1998. Paramyxovirus RNA synthesis and the requirement for hexamer genome length: the rule of six revisited. J. Virol. 72:891–899.

15. Lamb, R. A., and D. Kolakofsky. 1996. Paramyxoviridae: the viruses and their replication, p. 1177–1204. In B. N. Fields, D. M. Knipe, and P. M. Howley (ed.), Fields virology, 3rd ed. Lippincott-Raven, Philadelphia, Pa.

16. Mackenzie, J. S. 1999. Emerging viral diseases: an Australian perspective. Emerg. Infect. Dis. 5:1–8.

17. Mayo, M. A., and C. R. Pringle. 1998. Virus taxonomy—1997. J. Gen. Virol. 79:649–657.

18. Moritz, R. L., J. Eddes, H. Ji, G. E. Reid, and R. J. Simpson. 1995. Rapid separation of proteins and peptides using conventional silica-based supports: identification of 2-D gel proteins following in-gel proteolysis, p. 311–319. In J. W. Crabb (ed.), Techniques in protein chemistry, vol. VI. Academic Press, New York, N.Y.

19. Murray, K., P. Selleck, P. Hooper, A. Hyatt, A. Gould, L. Gleeson, H. Westbury, L. Hiley, L. Selvey, and B. Rodwell. 1995. A morbillivirus that caused fatal disease in horses and humans. Science 268:94–97.

20. O'Sullivan, J. D., A. M. Allworth, D. L. Paterson, T. M. Snow, R. Boots, L. J. Gleeson, A. R. Gould, A. D. Hyatt, and J. Bradfield. 1997. Fatal encephalitis due to novel paramyxovirus transmitted from horses. Lancet 349:93–95.

21. Philbey, A. W., P. D. Kirkland, A. D. Ross, R. Davis, J., A. B. Gleeson, R. J. Love, P. W. Daniels, A. R. Gould, and A. D. Hyatt. 1998. An apparently new virus (family Paramyxoviridae) infectious for pigs, humans, and fruit bats. Emerg. Infect. Dis. 4:269–271.

22. Poch, O., B. M. Blumberg, L. Bougueleret, and N. Tordo. 1990. Sequence comparison of five polymerases (L proteins) of unsegmented negative-strand RNA viruses: theoretical assignment of functional domains. J. Gen. Virol. 71:1153–1162.

23. Pringle, C. R. 1991. The order Mononegavirales. Arch. Virol. 117:137–140.

24. Pringle, C. R. 1998. Virus taxonomy—San Diego 1998. Arch. Virol. 143:1449–1459.

25. Pringle, C. R., and A. J. Easton. 1997. Monopartite negative strand RNA genomes. Semin. Virol. 8:49–57.

26. Renshaw, R. W., A. L. Glaser, H. Van Campen, F. Weiland, and E. J. Dubovi. 2000. Identification and phylogenetic comparison of Salem virus, a novel paramyxovirus of horses. Virology 270:417–429.

27. Rima, B., D. J. Alexander, M. A. Billeter, P. L. Collins, D. W. Kingsbury, M. A. Lipkind, Y. Nagai, C. Orvell, C. R. Pringle, and V. ter Meulen. 1995. Family Paramyxoviridae, p. 268–274. In F. A. Murphy, C. M. Fauquet, D. H. L. Bishop, S. A. Ghabrial, A. W. Jarvis, G. P. Martelli, M. A. Mayo, and M. D. Summers (ed.), Virus taxonomy. Sixth report of the International Committee on Taxonomy of Viruses. Springer-Verlag, Vienna, Austria.

28. Schnell, M. J., and K. K. Conzelmann. 1995. Polymerase activity of in vitro mutated rabies virus L protein. Virology 214:522–530.

29. Sleat, D. E., and A. K. Banerjee. 1993. Transcriptional activity and mutational analysis of recombinant vesicular stomatitis virus RNA polymerase. J. Virol. 67:1334–1339.

30. Tidona, C. A., H. W. Kurz, H. R. Gelderblom, and G. Darai. 1999. Isolation and molecular characterization of a novel cytopathogenic paramyxovirus from tree shrews. Virology 258:425–434.

31. Volchkov, V. E., V. A. Volchkova, A. A. Chepurnov, V. M. Blinov, O. Dolnik, S. V. Netesov, and H. Feldmann. 1999. Characterization of the L gene and 5′ trailer region of Ebola virus. J. Gen. Virol. 80:355–362.

32. Wang, L. F., A. R. Gould, and P. W. Selleck. 1997. Expression of equine morbillivirus (EMV) matrix and fusion proteins and their evaluation as diagnostic reagents. Arch. Virol. 142:2269–2279.

33. Wang, L. F., W. P. Michalski, M. Yu, L. I. Pritchard, G. Crameri, B. Shiell, and B. T. Eaton. 1998. A novel P/V/C gene in a new member of the Paramyxoviridae family, which causes lethal infection in humans, horses, and other animals. J. Virol. 72:1482–1490.

34. Young, P. L., K. Halpin, P. W. Selleck, H. Field, J. L. Gravel, M. A. Kelly, and J. S. Mackenzie. 1996. Serologic evidence for the presence in Pteropus bats of a paramyxovirus related to equine morbillivirus. Emerg. Infect. Dis. 2:239–240.

35. Yu, M., E. Hansson, J. P. Langedijk, B. T. Eaton, and L. F. Wang. 1998. The attachment protein of Hendra virus has high structural similarity but limited primary sequence homology compared with viruses in the genus Paramyxovirus. Virology 251:227–233.

36. Yu, M., E. Hansson, B. Shiell, W. Michalski, B. T. Eaton, and L. F. Wang. 1998. Sequence analysis of the Hendra virus nucleoprotein gene: comparison with other members of the subfamily Paramyxovirinae. J. Gen. Virol. 79:1775–1780.