1   **Identification of a novel ephemerovirus in a water buffalo (*Bubalus bubalis* [Linnaeus, 1758])**

2

3   Sakiho Imai[1], Mai Kishimoto[1], Masayuki Horie[1,2,3] *

4

5   [1] Laboratory of Veterinary Microbiology, Graduate School of Veterinary Medicine, Osaka

6   Metropolitan University

7   [2] Osaka International Infectious Diseases Research Center, Osaka Metropolitan University

8

9   *Corresponding author

10  Masayuki Horie, DVM, PhD

11  Email: mhorie@omu.ac.jp

12

**Abstract**

Ephemeroviruses, belonging to the genus *Ephemerovirus* within the family *Rhabdoviridae* of the *Mononegavirales*, are non-segmented, negative-strand RNA viruses that infect artiodactyls and blood-sucking arthropods. Although recent advances in sequencing technology have facilitated the identification of novel ephemeroviruses, thereby expanding our understanding of this viral genus, their diversity remains elusive, as evidenced by phylogenetic gaps between currently known ephemeroviruses. In this study, we analyzed publicly available RNA-seq data and identified a novel ephemerovirus, tentatively named Punjab virus (PBV), in a water buffalo (*Bubalus bubalis* [Linnaeus, 1758]). We obtained two separate PBV contigs from the RNA-seq data; the first contig covers the N, P, and M genes, while the second contig covers the G, α, β, γ, and L genes. Together, these PBV contigs represent 99% of the estimated complete viral genome. Mapping analysis revealed a typical transcriptional gradient pattern commonly observed in mononegaviruses, suggesting that the water buffalo is the authentic host for PBV. Sequence comparisons with its closest relatives indicate that the newly identified virus meets the ICTV species demarcation criteria for sequence divergence. Thus, this study contributes to a deeper understanding of the diversity of ephemeroviruses.

## Introduction

The diversity of viruses remains largely unknown, with the viruses we know today representing only a small portion of the entire virosphere [10]. Understanding the diversity of viruses is important from various perspectives, such as the control of infectious diseases, the elucidation of viral evolution, and the development of virotherapy. Identification of novel viruses can help with pandemic preparedness, as emerging infectious diseases can be caused by previously unidentified viruses [17]. Additionally, many phylogenetic gaps are present between known viruses [26], which suggests the existence of undiscovered viruses filling the gaps. Furthermore, viruses can be utilized for various treatments and vaccines. Consequently, understanding the diversity of viruses is crucial from a variety of perspectives.

Ephemeroviruses, members of the genus *Ephemerovirus* within the family *Rhabdoviridae* of the order *Mononegavirales*, possess single-stranded negative-strand RNA as their genomes [30]. These genomes encode structural proteins N, P, M, G, and L, as well as several nonstructural proteins that are encoded between the G and L genes. The genes encoding structural proteins are conserved among all the known ephemeroviruses, whereas the genes encoding nonstructural proteins may vary among the viruses.

Ephemeroviruses have been detected in artiodactyls and blood-sucking arthropods [30, 31]. Bovine ephemeral fever virus (BEFV), which is the type virus of the genus *Ephemerovirus*, is an arbovirus proven to cause a disease called BEF in certain ruminants, such as cattle [31]. Animals affected by BEF exhibit symptoms like acute fever, arthralgia, and dropping milk production, which result in economic loss. On the other hand, the pathogenicity of other ephemeroviruses remains unclear, although some ephemeroviruses have been reported in animals showing symptoms [3].

A diverse range of ephemeroviruses have been detected through various methods to date, and recent advancements in metagenomic analyses have expedited these discoveries. However, there still remain large phylogenetic gaps among ephemeroviruses [30], suggesting the presence of undiscovered ephemeroviruses that fill the gaps. Therefore, further exploration of ephemeroviruses is crucial for a more comprehensive understanding of the diversity and potential pathogenicity of ephemeroviruses.

59    In this study, we investigated the diversity of ephemeroviruses by analyzing publicly available

60    RNA-seq data. We identified novel ephemerovirus sequences, tentatively named Punjab virus (PBV),

61    in RNA-seq data obtained from a water buffalo (*Bubalus bubalis* [Linnaeus, 1758]). Our detailed

62    sequence analysis revealed that the viral genome sequence is divided into two contigs, estimated to

63    cover approximately 99% of the complete genome. Phylogenetic analysis showed that this novel

64    ephemerovirus forms a cluster with Kokolu virus, Puchong virus, and Hayes Yard virus. Moreover,

65    PBV meets the ICTV species demarcation criteria for sequence divergence, suggesting that this virus

66    can be classified as a new species within the genus *Ephemerovirus*.

67

68    **Material Method**

69    ***Detection of ephemerovirus-like contigs***

70    RNA-seq data (accession number SRR8476835) were downloaded from the NCBI SRA [24], which

71    were preprocessed by fastp 0.23.2 [6] using the "-l 35  -x  -y" options. The preprocessed reads were

72    then assembled by Trinity 2.14.0 [12], SKESA 2.5.1 [27], SPAdes v3.15.5 [23], or metaSPAde

73    v3.15.5 [22] using the default setting. Contigs obtained by Trinity that were 100 nucleotides or more

74    were extracted using SeqKit 2.3.0 [25], and then clustered using CD-HIT v4.8.1 [9] with a threshold

75    of  0.98. The clustered contigs were used for a two-step sequence similarity search as follows.

76    In the first step, a sequence similarity search was conducted against a custom database

77    containing protein sequences of viruses belonging to the kingdom *Orthornavirae* by MMseqs2

78    version c48da9d781b81804727b5cccfed7f97cfcc20c9d [28] using the clustered contigs as queries.

79    From the hit contigs, those with $E$-values of less than $10^{-20}$ and whose top hits (hit sequence with the

80    highest score) were viruses were extracted. Among the extracted contigs, only one representative

81    sequence was used in the subsequent analysis for a group of sequences considered to be isoforms

82    based on the contig IDs.

83    The second sequence similarity search was performed against the NCBI nr database [24] by

84    BLASTx 2.13.0 [4] with the options "-evalue 1e-20 -max_target_seqs 10 -word_size 2 -

85    lcase_masking" using the extracted contigs as queries. The contigs whose BLAST best hits were

86    viruses were extracted and used in subsequent analyses as virus-like contigs.

87

*Validation of the virus-like contigs by mapping analysis*

To validate the accuracy of virus-like contigs, a mapping analysis was performed. The original RNA-seq data (accession number SRR8476835) were mapped to the obtained PBV contigs by HISAT2 2.2.1 [18], and the read depth at each position was calculated using SAMtools 1.16.1 [7]. The positions covered by five or more reads were considered as reliable positions.

93

*Annotation of the virus-like contigs*

Open reading frames (ORFs) consisting of more than 256 nucleotides (based on the lengths of ephemeroviral ORFs) were identified in the virus-like contigs using Geneious Prime (https://www.geneious.com). ORFs that spanned transcription signals (see below) were manually corrected. BLASTp searches were conducted against the protein sequences of viruses (taxid:10239) in the NCBI nr database on the BLAST web server using translated sequences of ORFs as queries. The following options were used: Word size: 3; Expect threshold: $10^{-10}$.

To identify transcription signals, conserved motifs were searched by MEME 5.5.0 [1] with the options "-mod oops -maxw 10 -nmotifs 3 -dna" using the sequences of intergenic regions. Each identified motif sequence with its flanking 4 nucleotides was extracted, which were aligned by MAFFT v7.490 using the E-INS-i algorithm [15]. Putative transcription signals were determined based on the alignments.

Signal peptide prediction was performed using the SignalP 6.0 web server [29].

107

*Phylogenetic analyses*

Phylogenetic trees were inferred using the putative amino acid sequences of N, G, or L proteins of PBV, known members of the genus *Ephemerovirus*, and three other rhabdoviruses (outgroups) (Table S1). The sequences were aligned by MAFFT v7.490 using the E-INS-i algorithm, and the ambiguously aligned regions were trimmed by Trimal v1.4.rev22 with the "-strict" option [5]. Phylogenetic trees were reconstructed by the maximum likelihood method using RAxML Next

114    Generation 1.1.0 [20]. LG+I+G4, WAG+I+G4+FC, and LG+I+G4+FC models, chosen by

115    ModelTest-NG v0.2.0 [8], were used for the inference of N, G, and L trees, respectively.

116

117    *Determination of pairwise sequence identities*

118    Amino acid sequences of N, G, or L proteins of PBV and closely related ephemeroviruses (Table S1)

119    were aligned by MAFFT, and then pairwise sequence identities were determined using Sequence

120    Demarcation Tool version 1.2 [21].

121

122    *Mapping analysis to detect Punjab virus infection*

123    To detect PBV infection, a total of 46,244 publicly available RNA-seq data (accession numbers are

124    available in Supplementary Materials) were downloaded and preprocessed by fastp 0.23.2 with the

125    options "-x -y -l 35". The preprocessed reads were then mapped to the PBV contigs using HISAT2

126    2.2.1. The numbers of mapped reads were counted using SAMtools. The mapped reads were also

127    manually analyzed to check the accuracy of the mapping.

128

129    **Results**

130    *Identification of a novel ephemerovirus*

131    In our previous study, we performed a large-scale metaviromic analysis and detected many RNA

132    viruses from publicly available RNA-seq data [16]. However, detailed analyses were conducted only

133    for the viral sequences that were close to the full-length genomes. Consequently, many of the detected

134    partial viral sequences have not yet been analyzed well. Therefore, we reanalyzed the BLAST results

135    obtained in the previous study and found that one of the RNA-seq data sets (accession number

136    SRR8476835) obtained from the blood of a water buffalo (*B. bubalis*) [13] contains ephemerovirus-

137    like sequences. To confirm this result, we again performed *de novo* assembly and a two-step sequence

138    similarity search using the resultant contigs. Consistent with the previous result, we detected two

139    ephemerovirus-like contigs whose respective BLAST best hits were Hayes Yard virus N protein

140    (QEA08650.1; 90.3% identity) and Puchong virus L protein (QEA08648.1; 78.4% identity) (Table 1).

141    To validate the accuracy of obtained ephemerovirus-like contigs, we mapped the original short

142    reads to the contigs and measured the read depths. In this study, we defined positions mapped by five

143    or more reads as "reliable" regions. As a result, we removed some of the extreme terminal sequences

144    of the obtained contigs, resulting in contigs with lengths of 3007 (Contig 1) and 11819 (Contig 2)

145    nucleotides. It is important to note that the mapping pattern showed a typical transcription gradient

146    observed in mononegaviruses, further supporting the assertion that the contigs are derived from an

147    ephemerovirus (Fig. 1).

148

149    *Characterization of the ephemerovirus-like contigs*

150    To determine the genomic structure of  PBV, we extracted ORFs from the contigs and performed

151    BLASTp searches using each of the ORFs as a query. As a result, we identified three and seven ORFs

152    in Contig 1 and Contig 2, showing sequence similarities to N, P, and M and G, Gns, α1, α2, β, γ, L

153    genes of other ephemeroviruses, respectively (Table S2). Because it was initially unclear whether the

154    annotated G gene is full length or not due to its location at the end of contig (Fig. 1b), we

155    characterized the putative G protein *in silico*. The putative G protein sequence was predicted to

156    contain a signal peptide at the N-terminus (Fig. S1a). Furthermore, the putative G protein was

157    alignable with the full-length G proteins of related viruses (Fig. S1b). These results strongly suggest

158    that the annotated G gene is full-length. On the other hand, the stop codon of L gene was not included

159    in the contig (Fig. 1b).

160    We subsequently performed MEME searches to identify putative transcription signal sequences

161    in the intergenic regions. These searches, in combination with manual curation, identified putative

162    transcription initiation signals 5'-AACAGG-3' and termination/polyadenylation signals 5'-

163    ATGAAAAAAA-3' (Fig. 1c).

164

165    *Phylogenetic analysis*

166    To understand the evolutionary relationships between PBV and other ephemeroviruses, we conducted

167    phylogenetic analyses using the amino acid sequences of the N, G, and L proteins (Figs. 2 and S2).

168    All the trees show that PBV forms a well-supported cluster with Kokolu virus, Puchong virus, and

169    Hayes Yard virus, and diverged earlier than these three viruses. Additionally, the clade containing

170    PBV and the aforementioned three viruses is closely related to the clade containing BEFV.

171

172    ***Amino acid sequence divergence between Punjab virus and the closely related viruses***

173    To investigate the amino acid divergence between PBV and the closely related viruses, we determined

174    pairwise identities using SDT with the amino acid sequences of N, G, and L proteins from closely

175    related viruses listed in Table S1. The maximum amino acid identities were 91.2% for the N protein

176    (Puchong virus), 64.5% for the G protein (Kokolu virus), and 77.8% for the L protein (Berrimah

177    virus), respectively (Fig. 3).

178

179    ***Mapping analysis to detect infection from other public RNA-seq data***

180    To gain more insight into PBV infection, we searched for PBV-like sequences in public RNA-seq

181    data by mapping analysis. Given that some ephemeroviruses are known to be arboviruses, we mapped

182    short reads from publicly available RNA-seq data of ticks (subclass Acari), mosquitos (family

183    Culicidae), biting midges (family Ceratopogonidae), and bovines (subfamily Bovinae) to PBV

184    contigs, and then counted the number of mapped reads. We detected a small amount of mapped reads

185    from three RNA-seq data sets belonging to the same BioProject from which we originally detected the

186    viral contigs (Table S3). However, we cannot rule out the possibility that these were due to cross-

187    contamination and/or index hopping, and therefore it is unclear whether these samples really

188    contained the virus.

189

190    **Discussion**

191    To date, 13 species of viruses have been identified in the genus *Ephemerovirus*. However, the

192    divergence of ephemeroviruses remains unclear, as suggested by the presence of phylogenetic gaps

193    [30]. In this study, we identified a novel ephemerovirus in publicly available RNA-seq data obtained

194    from *B. bubalis*. A series of analyses showed that the identified viral sequences possess a typical

195    ephemerovirus genome structure and also exhibit a characteristic transcription pattern of

196    mononegaviruses. Importantly, our analyses demonstrated that the PBV exhibits an amino acid

8

197    sequence divergence of 8.8%, 35.5%, and 35.2% from the most closely related viral N, G, and L

198    proteins, respectively (Fig. 3). This fulfills the current species demarcation criteria of the genus

199    *Ephemerovirus* in terms of sequence divergence. Although we were unable to obtain a single contig of

200    this virus, and the L protein lacks its C-terminal sequence, our data suggest that the PBV can be

201    classified as a new species within the genus *Ephemerovirus*.

202          The pathogenicity of PBV remains uncertain. The RNA-seq data (SRR8476835), used for this

203    study was sourced from a blood sample collected from a water buffalo affected by metritis. In the

204    same BioProject (PRJNA514883), there exists additional RNA-seq data (SRR8476836) obtained from

205    another individual also affected by metritis. However, only a few viral reads were detected in this

206    second individual (SRR8476836), creating ambiguity regarding whether the water buffalo was indeed

207    infected by PBV, especially considering cross-contamination and index hopping [19, 32]. Hayes Yard

208    virus, one of the closely related ephemeroviruses, was isolated from a bull (*Bos indicus* [Linnaeus,

209    1758]) afflicted with a severe ephemeral fever-like illness, but it remains inconclusive whether this

210    virus was the causative agent [3]. Furthermore, while preparing this manuscript, we noted that another

211    study identified a novel ephemerovirus, which can be classified into the same species as PBV, in a

212    febrile cow (Figs. S3) [11]. Further epidemiological studies are essential to improve our

213    understanding of the pathogenicity of ephemeroviruses, including PBV.

214          The mapping pattern provides strong evidence that *B. bubalis* is a legitimate host for PBV. In

215    viral metagenomic analysis, host identification can sometimes be challenging because samples may

216    contain nucleic acids from viruses of various environmental and dietary origins. Our mapping analysis

217    showed the typical transcription gradient from N to L genes observed in mononegaviruses, implying

218    that PBV was actively transcribing in the samples. Since the RNA was extracted from blood, the

219    likelihood of contamination is minimal. Moreover, the genetically related viruses were also detected

220    from bovines  [2, 3, 14]. Considering these points, *B. bubalis* would be an authentic host for PBV. It

221    should be noted that some ephemeroviruses are known to be arboviruses. As PBV was detected in

222    blood samples, it is plausible that the virus could be transmitted by arthropod vectors. Further studies

223    are required to elucidate the transmission route of PBV.

224    In this study, we only obtained two separate contigs of the PBV genome, but not a single one.

225    Besides using Trinity, we performed *de novo* assembly with several assemblers (SKESA, SPAdes,

226    metaSPAdes), yet we consistently obtained two separate contigs (data not shown). This is likely

227    because the mRNA-seq does not adequately cover the intergenic regions (Fig. 1). Unfortunately, only

228    a few viral reads were detected from RNA-seq data other than the initially detected one (Table S3),

229    making co-assembly unavailable. Further accumulation of data or in-depth molecular epidemiological

230    studies are required to determine the complete genome of PBV.

231    Together, we identified a novel ephemerovirus from public RNA-seq data, thereby contributing

232    to a deeper understanding of the diversity of ephemeroviruses. However, the virological

233    characteristics of PBV, such as its pathogenicity and infection route, remain unclear. Further

234    identification of infected individuals and the accumulation of sequence information would contribute

235    to the characterization of PBV.

236

237    **Conflict of interest**

238    The authors declare no conflict of interest in this study.

239

245

246    **References**

247    1.    Bailey, T. L., Johnson, J., Grant, C. E. and Noble, W. S. 2015. The MEME Suite. *Nucleic Acids*

248          *Res.* **43**: W39-49.

249    2.    Balinandi, S., Hayer, J., Cholleti, H., Wille, M., Lutwama, J. J., Malmberg, M. and Mugisha, L.

250          2022. Identification and molecular characterization of highly divergent RNA viruses in cattle,

251          Uganda. *Virus Res.* **313**: 198739.

252    3.    Blasdell, K. R., Davis, S. S., Voysey, R., Bulach, D. M., Middleton, D., Williams, S., Harmsen,

253           M. B., Weir, R. P., Crameri, S., Walsh, S. J., Peck, G. R., Tesh, R. B., Boyle, D. B., Melville, L.

254           F. and Walker, P. J. 2020. Hayes Yard virus: a novel ephemerovirus isolated from a bull with

255           severe clinical signs of bovine ephemeral fever is most closely related to Puchong virus. *Vet.*

256           *Res.* **51**: 58.

257    4.    Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K. and Madden, T.

258           L. 2009. BLAST+: architecture and applications. *BMC Bioinformatics*. **10**: 421.

259    5.    Capella-Gutiérrez, S., Silla-Martínez, J. M. and Gabaldón, T. 2009. trimAl: a tool for automated

260           alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*. **25**: 1972–1973.

261    6.    Chen, S., Zhou, Y., Chen, Y. and Gu, J. 2018. fastp: an ultra-fast all-in-one FASTQ

262           preprocessor. *Bioinformatics*. **34**: i884–i890.

263    7.    Danecek, P., Bonfield, J. K., Liddle, J., Marshall, J., Ohan, V., Pollard, M. O., Whitwham, A.,

264           Keane, T., McCarthy, S. A., Davies, R. M. and Li, H. 2021. Twelve years of SAMtools and

265           BCFtools. *Gigascience*. **10**:.

266    8.    Darriba, D., Posada, D., Kozlov, A. M., Stamatakis, A., Morel, B. and Flouri, T. 2020.

267           ModelTest-NG: A New and Scalable Tool for the Selection of DNA and Protein Evolutionary

268           Models. *Mol. Biol. Evol.* **37**: 291–294.

269    9.    Fu, L., Niu, B., Zhu, Z., Wu, S. and Li, W. 2012. CD-HIT: accelerated for clustering the next-

270           generation sequencing data. *Bioinformatics*. **28**: 3150–3152.

271    10.  Geoghegan, J. L. and Holmes, E. C. 2017. Predicting virus emergence amid evolutionary noise.

272           *Open Biol.* **7**:.

273    11.  Golender, N., Klement, E., Ofer, L., Hoffmann, B., Wernike, K., Beer, M. and Pfaff, F. 2023.

274           Hefer valley virus: a novel ephemerovirus detected in the blood of a cow with severe clinical

275           signs in Israel in 2022. *Arch. Virol.* **168**: 234.

276    12.  Grabherr, M. G., Haas, B. J., Yassour, M., Levin, J. Z., Thompson, D. A., Amit, I., Adiconis, X.,

277           Fan, L., Raychowdhury, R., Zeng, Q., Chen, Z., Mauceli, E., Hacohen, N., Gnirke, A., Rhind, N.,

278           di Palma, F., Birren, B. W., Nusbaum, C., Lindblad-Toh, K., Friedman, N. and Regev, A. 2011.

279           Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat.*

280      *Biotechnol.* **29**: 644–652.

281   13.  Gurao, A., Vasisth, R., Singh, R., Dige, M. S., Vohra, V., Mukesh, M., Kumar, S. and Kataria, R.

282      S. 2022. Identification of differential methylome signatures of white pigmented skin patches in

283      Nili Ravi buffalo of India. *Environ. Mol. Mutagen.* **63**: 408–417.

284   14.  Karabatsos, N. 1978. Supplement to International Catalogue of Arboviruses including certain

285      other viruses of vertebrates. *Am. J. Trop. Med. Hyg.* **27**: 372–440.

286   15.  Katoh, K. and Standley, D. M. 2013. MAFFT multiple sequence alignment software version 7:

287      improvements in performance and usability. *Mol. Biol. Evol.* **30**: 772–780.

288   16.  Kawasaki, J., Kojima, S., Tomonaga, K. and Horie, M. 2021. Hidden Viral Sequences in Public

289      Sequencing Data and Warning for Future Emerging Diseases. *MBio*. **12**: e0163821.

290   17.  Kawasaki, J., Tomonaga, K. and Horie, M. 2023. Large-scale investigation of zoonotic viruses in

291      the era of high-throughput sequencing. *Microbiol. Immunol.* **67**: 1–13.

292   18.  Kim, D., Paggi, J. M., Park, C., Bennett, C. and Salzberg, S. L. 2019. Graph-based genome

293      alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* **37**: 907–915.

294   19.  Kircher, M., Sawyer, S. and Meyer, M. 2012. Double indexing overcomes inaccuracies in

295      multiplex sequencing on the Illumina platform. *Nucleic Acids Res.* **40**: e3.

296   20.  Kozlov, A. M., Darriba, D., Flouri, T., Morel, B. and Stamatakis, A. 2019. RAxML-NG: a fast,

297      scalable and user-friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics*.

298      **35**: 4453–4455.

299   21.  Muhire, B. M., Varsani, A. and Martin, D. P. 2014. SDT: a virus classification tool based on

300      pairwise sequence alignment and identity calculation. *PLoS One*. **9**: e108277.

301   22.  Nurk, S., Meleshko, D., Korobeynikov, A. and Pevzner, P. A. 2017. metaSPAdes: a new

302      versatile metagenomic assembler. *Genome Res.* **27**: 824–834.

303   23.  Prjibelski, A., Antipov, D., Meleshko, D., Lapidus, A. and Korobeynikov, A. 2020. Using

304      SPAdes De Novo Assembler. *Curr. Protoc. Bioinformatics*. **70**: e102.

305   24.  Sayers, E. W., Bolton, E. E., Brister, J. R., Canese, K., Chan, J., Comeau, D. C., Connor, R.,

306      Funk, K., Kelly, C., Kim, S., Madej, T., Marchler-Bauer, A., Lanczycki, C., Lathrop, S., Lu, Z.,

307      Thibaud-Nissen, F., Murphy, T., Phan, L., Skripchenko, Y., Tse, T., Wang, J., Williams, R.,

308          Trawick, B. W., Pruitt, K. D. and Sherry, S. T. 2022. Database resources of the national center

309          for biotechnology information. *Nucleic Acids Res.* **50**: D20–D26.

310   25.  Shen, W., Le, S., Li, Y. and Hu, F. 2016. SeqKit: A Cross-Platform and Ultrafast Toolkit for

311          FASTA/Q File Manipulation. *PLoS One*. **11**: e0163962.

312   26.  Shi, M., Zhang, Y.-Z. and Holmes, E. C. 2018. Meta-transcriptomics and the evolutionary

313          biology of RNA viruses. *Virus Res.* **243**: 83–90.

314   27.  Souvorov, A., Agarwala, R. and Lipman, D. J. 2018. SKESA: strategic k-mer extension for

315          scrupulous assemblies. *Genome Biol.* **19**: 153.

316   28.  Steinegger, M. and Söding, J. 2017. MMseqs2 enables sensitive protein sequence searching for

317          the analysis of massive data sets. *Nat. Biotechnol.* **35**: 1026–1028.

318   29.  Teufel, F., Almagro Armenteros, J. J., Johansen, A. R., Gíslason, M. H., Pihl, S. I., Tsirigos, K.

319          D., Winther, O., Brunak, S., von Heijne, G. and Nielsen, H. 2022. SignalP 6.0 predicts all five

320          types of signal peptides using protein language models. *Nat. Biotechnol.* **40**: 1023–1025.

321   30.  Walker, P. J., Freitas-Astúa, J., Bejerman, N., Blasdell, K. R., Breyta, R., Dietzgen, R. G., Fooks,

322          A. R., Kondo, H., Kurath, G., Kuzmin, I. V., Ramos-González, P. L., Shi, M., Stone, D. M.,

323          Tesh, R. B., Tordo, N., Vasilakis, N., Whitfield, A. E. and Ictv Report Consortium 2022. ICTV

324          Virus Taxonomy Profile: Rhabdoviridae 2022. *J. Gen. Virol.* **103**:.

325   31.  Walker, P. J. and Klement, E. 2015. Epidemiology and control of bovine ephemeral fever. *Vet.*

326          *Res.* **46**: 124.

327   32.  Wright, E. S. and Vetsigian, K. H. 2016. Quality filtering of Illumina index reads mitigates

328          sample cross-talk. *BMC Genomics*. **17**: 876.

329

330

331   **Figure legends**

332   **Figure 1. Genomic organization of Punjab virus. (a)** Genomic organization and transcription

333   profile of Punjab virus. Pink arrow boxes show open reading frames. Short reads from SRR8476835

334   were mapped to the Punjab virus contigs and visualized. **(b)** Putative transcription signal sequences of

335   Punjab virus.

336

337   **Figure 2. Phylogenetic relationship of Punjab virus and ephemeroviruses.** Phylogenetic trees

338   were reconstructed by the maximum likelihood method using amino acid sequences of N **(a)** or  L **(b)**

339   protein of Punjab virus, ephemeroviruses, and outgroup rhabdoviruses. Bootstrap values equal to or

340   more than 70 are shown on each branch. The scale bar indicates the number of amino acid

341   substitutions per site.

342

343   **Figure 3. Sequence divergence of Punjab virus and related ephemeroviruses.**

344   Pairwise amino acid sequence identities of N **(a)**, G **(b)**, or L **(c)** proteins between Punjab virus and

345   related ephemeroviruses were determined using Sequence Demarcation Tool [21]. Punjab virus;

346   PUCV, Puchong virus; KoV, Kokolu virus; HYV, Hayes Yard virus; BRMV, Berrimah virus; BEFV,

347   Bovine ephemeral fever virus.

348

14

**Table 1. The top hits of BLASTx analysis.**

| Query | BLAST best hit | | | | |
| --- | --- | --- | --- | --- | --- |
| | Accession | Virus name | Protein | Identity (%) | Length (aa) |
| Contig 1 | QEA08650.1 | Hayes Yard virus | N protein | 90.3 | 432 |
| Contig 2 | QEA08648.1 | Puchong virus | L protein | 78.4 | 2098 |

# Figure 1

**a**



Contig 1

**b**



Contig 2

**c**

| Region | Initiation | Termination |
|--------|-----------|-------------|
| N-P | AACAGG | ATGAAAAAAA |
| P-M | AACAGG | ATGAAAAAAA |
| G-Gns | AACAGG | ATGAAAAAAA |
| Gns-α | AACAGG | ATGAAAAAAA |
| α-β | AACAGG | ATGAAAAAAA |
| β-γ | AACAGG | ATGAAAAAAA |
| γ-L | AACAGG | ATGAAAAAAA |

# Figure 2

**a**

**N protein**



**b**

**L protein**

## Figure 3



**a** — N protein

|  | Punjab virus | PUCV | KoV | HYV | BRMV | BEFV |
|---|---|---|---|---|---|---|
| **Punjab virus** | 100.0 | | | | | |
| **PUCV** | 91.2 | 100.0 | | | | |
| **KoV** | 90.8 | 96.2 | 100.0 | | | |
| **HYV** | 90.3 | 95.1 | 94.8 | 100.0 | | |
| **BRMV** | 77.0 | 78.6 | 77.7 | 77.9 | 100.0 | |
| **BEFV** | 76.9 | 75.9 | 76.0 | 76.9 | 91.6 | 100.0 |

**b** — G protein

|  | Punjab virus | PUCV | KoV | HYV | BRMV | BEFV |
|---|---|---|---|---|---|---|
| **Punjab virus** | 100.0 | | | | | |
| **KoV** | 64.5 | 100.0 | | | | |
| **PUCV** | 64.5 | 83.8 | 100.0 | | | |
| **HYV** | 64.5 | 74.6 | 76.0 | 100.0 | | |
| **BRMV** | 49.1 | 49.7 | 49.8 | 50.1 | 100.0 | |
| **BEFV** | 50.5 | 50.5 | 52.0 | 49.9 | 75.3 | 100.0 |

**c** — L protein

|  | Punjab virus | PUCV | KoV | HYV | BRMV | BEFV |
|---|---|---|---|---|---|---|
| **Punjab virus** | 100.0 | | | | | |
| **PUCV** | 77.8 | 100.0 | | | | |
| **KoV** | 77.5 | 90.4 | 100.0 | | | |
| **HYV** | 78.4 | 86.7 | 86.7 | 100.0 | | |
| **BRMV** | 64.8 | 64.3 | 64.8 | 65.0 | 100.0 | |
| **BEFV** | 64.0 | 63.9 | 64.3 | 64.2 | 84.0 | 100.0 |

Identity (%)

100

40