1    **Running title: The transcriptional and translational landscape of equine torovirus**

2    Hazel Stewart[1±], Katherine Brown[1±], Adam M. Dinan[1¥], Nerea Irigoyen[1], Eric J.

3    Snijder[2], Andrew E. Firth[1]*

4    [1]Division of Virology, Department of Pathology, University of Cambridge, Cambridge,

5    United Kingdom.

6    [2]Molecular Virology Laboratory, Department of Medical Microbiology, Leiden

7    University Medical Center, Leiden, The Netherlands.

8    [¥] Current address: Fios Genomics, Edinburgh, United Kingdom.

9    * corresponding author

10    [±] These two authors contributed equally.

11    **Abstract**

12    The genus *Torovirus* (subfamily *Torovirinae*, family *Coronaviridae,* order *Nidovirales*)

13    encompasses a range of species that infect domestic ungulates including cattle,

14    sheep, goats, pigs and horses, causing an acute self-limiting gastroenteritis. Using the

15    prototype species equine torovirus (EToV) we performed parallel RNA sequencing

16    (RNA-seq) and ribosome profiling (Ribo-seq) to analyse the relative expression levels

17    of the known torovirus proteins and transcripts, chimaeric sequences produced via

18    discontinuous RNA synthesis (a characteristic of the nidovirus replication cycle) and

19    changes in host transcription and translation as a result of EToV infection. RNA

20    sequencing confirmed that EToV utilises a unique combination of discontinuous and

21    non-discontinuous RNA synthesis to produce its subgenomic RNAs; indeed, we

22    identified transcripts arising from both mechanisms that would result in sgRNAs

23    encoding the nucleocapsid. Our ribosome profiling analysis revealed that ribosomes

24    efficiently translate two novel CUG-initiated ORFs, located within the so-called 5'

25    UTR. We have termed the resulting proteins U1 and U2. Comparative genomic

26    analysis confirmed that these ORFs are conserved across all available torovirus

27    sequences and the inferred amino acid sequences are subject to purifying selection,

28    indicating that U1 and U2 are functionally relevant. This study provides the first high-

29    resolution analysis of transcription and translation in this neglected group of

30    livestock pathogens.

31    **Importance**

32    Toroviruses infect cattle, goats, pigs and horses worldwide and can cause

33    gastrointestinal disease. There is no treatment or vaccine and their ability to spill

34    over into humans has not been assessed. These viruses are related to important

35    human pathogens including severe acute respiratory syndrome (SARS) coronavirus

36    and they share some common features, however the mechanism that they use to

37    produce subgenomic RNA molecules differs. Here we performed deep sequencing to

38    determine how equine torovirus produces subgenomic RNAs. In doing so, we also

39    identified two previously unknown open reading frames "hidden" within the

40    genome. Together these results highlight the similarities and differences between

41    this domestic animal virus and related pathogens of humans and livestock.

42

**Introduction**

The order *Nidovirales* currently contains four families of positive-sense, single-stranded RNA viruses: the *Coronaviridae*, *Arteriviridae*, *Roniviridae* and *Mesoniviridae* (1). Their grouping into the one taxonomic order is based upon replicase protein conservation, genome organisation and replication strategy. However these viral families are nonetheless very diverse with respect to their virion structure, host range, pathogenic potential and genome size.

The genus *Torovirus* (family *Coronaviridae,* subfamily *Torovirinae*) encompasses a range of species with worldwide distribution that infect domestic ungulates including cattle, goats, sheep, pigs and horses, causing an acute self-limiting gastroenteritis. Approximately 55 % of cattle within the United Kingdom are seropositive for bovine torovirus and this pathogen represents a significant burden to the industry (2, 3). Similarly porcine torovirus is endemic in Europe and causes disease in production herds (4-6). Despite this, limited research has been conducted upon these pathogens and neither specific antiviral treatments nor vaccines are available. The prevalence of toroviruses in non-domestic reservoirs and potential for cross-species transmission has not been assessed, although they are known to undergo recombination events (7). The extensive research conducted upon the related coronaviruses would not necessarily be relevant in the event of an emerging torovirus infection, due to the divergent nature of these viruses.

The genomes of *Nidovirales* are positive-sense, polycistronic RNAs. One of the hallmarks of this virus order is the utilisation of an unusual transcription mechanism to express the genes encoding structural and accessory proteins, which reside downstream of the large replicase open reading frames (ORFs) 1a and 1b (Figure 1). These proteins are typically translated from a nested set of 3' coterminal subgenomic mRNAs (sg mRNAs). Although, with the exception of the smallest species, these sgRNAs are structurally polycistronic, translation is normally limited to the 5' ORF of each mRNA. Studies of coronaviruses and arteriviruses have revealed that they produce negative-sense subgenome-sized RNAs via a mechanism of "discontinuous" extension (8). This process resembles homology-assisted copy-

3

73    choice recombination (9) and requires the presence of multiple copies of a species-

74    specific short motif, the transcription regulatory sequence (TRS). TRS motifs are

75    located immediately upstream of the structural protein ORFs (body TRSs) and within

76    the 5' UTR (leader TRS).

77    Negative strand RNA synthesis initiates at the 3' end of the positive-sense viral

78    genome. When the RNA-dependent RNA polymerase (RdRp) has copied a TRS

79    sequence, a translocation event may occur during which the anti-TRS at the 3' end of

80    the nascent strand basepairs with the leader TRS within the 5' UTR. Transcription

81    reinitiates and continues to the 5' end of the genomic template. The resulting "anti-

82    leader" sequence that is added ranges from 55 – 92 nt in coronaviruses to ~200 nt in

83    arteriviruses. These negative-sense transcripts are therefore 5'- and 3'-coterminal

84    with the full length negative RNA strand and are identifiable as chimaeras with

85    distinct leader-body junctions. The anti-leader sequence in each of the negative-

86    sense templates then functions as a promoter, to drive synthesis of a mirror set of

87    positive-sense sgRNAs that are translated to produce the structural proteins.

88    However not all details of the mechanism outlined above are wholly conserved

89    across the *Nidovirales.* Specifically, the two sg mRNAs of roniviruses (pathogens of

90    shrimp) do not possess conserved 5' leader sequences, indicative of the lack of a

91    discontinuous step during their production (10). Despite the presence of a conserved

92    body TRS in each sg mRNA, an equivalent leader TRS is not readily identifiable in the

93    5' UTR. It may therefore be reasoned that the ronivirus body TRSs stimulate

94    termination of RNA synthesis without RdRp translocation and reinitiation.

95    Mesoniviruses (a branch of *Nidovirales* recently identified in insects) are thought to

96    produce two major sgRNAs possessing leader sequences of different lengths,

97    indicating the nidoviral mechanism for discontinuous RNA synthesis may allow two

98    very different leader/body TRS pairs to be utilised in a single viral species (11).

99    Toroviruses appear to represent a nidovirus subgroup with a remarkably flexible

100    transcription strategy: equine torovirus (EToV) possesses a leader TRS-like sequence

101    (CUUUAGA) but it is only involved in the synthesis of the mRNA used for expression

102    of the spike (S) protein gene (12). Despite similarities to the corona- and arteriviral

4

103     mechanism, the preceding leader sequence incorporated into this mRNA is merely 6

104     nt in length (ACGUAU). Additionally, this case is unusual in that the translocation

105     event is thought to be prompted by an RNA structure - a predicted RNA hairpin

106     upstream of the S protein gene, rather than a body TRS (12). Body TRSs are located

107     upstream of the three remaining structural protein genes, yet a non-discontinuous

108     mechanism is utilised for their production, as is the case for roniviruses. As a result,

109     the sg mRNAs for membrane (M), nucleocapsid (N) and haemagglutinin-esterase

110     (HE) do not normally possess a conserved 5' leader sequence; they each possess a

111     variable and unique extended version of the TRS at their 5' end. It is clear there is

112     significant difference between how the various *Nidovirales* families synthesise their

113     sgRNAs.

114     Here we describe the first high-resolution analysis of viral transcription during

115     infection by EToV, which is one of the few toroviruses that can be propagated in cell

116     culture (13, 14). RNA sequencing (RNA-seq) confirmed previous reports that EToV

117     utilises a unique combination of both discontinuous and non-discontinuous RNA

118     synthesis to generate its repertoire of sgRNAs. Strikingly, we also identified a small

119     proportion of chimaeric transcripts spanning from the leader to the body TRS of the

120     N protein gene, indicating that discontinuous and non-discontinuous mechanisms

121     compete in this location. We also identified numerous locations across the genome

122     where non-canonical RdRp translocation occurs, leading to a vast array of

123     (presumably mostly non-functional) chimaeric transcripts.

124     Ribosome profiling (Ribo-seq) conducted in tandem with the RNA-seq indicated

125     ribosomes were actively translating within the so-called 5' UTR. Further analysis

126     confirmed the existence of two novel ORFs in this region, which are conserved in all

127     torovirus genome sequences analysed to date. The specific function(s) of these

128     proteins will be the topic of future work. Together, these results provide an overview

129     of the transcriptional and translational events that accompany infection by this wide-

130     ranging pathogen.

131     **Results**

5

132 **Tandem RNA-seq and Ribo-seq of EToV infected cells.** We conducted tandem RNA-

133 seq and Ribo-seq of EToV infected equine dermal (ED) cells. Two biological replicates

134 of virus-infected and mock-infected cells were analysed, generating 25 to 53 million

135 reads per sample. For RNA-seq, 77-92 % of reads mapped to the host genome, of

136 which a mean of 1.5 % mapped to rRNA, 19 % to mRNA, 32 % to ncRNA and 47 %

137 elsewhere in the genome. For Ribo-seq, 46-60 % of reads mapped to the host

138 genome, of which a mean of 56 % mapped to rRNA, 13 % to mRNA, 4.9 % to ncRNA

139 and 26 % elsewhere in the genome (Supplementary Table 1). 1.3 % and 2.3 % of

140 reads mapped to the virus genome in the two EToV-infected RNA-seq replicates and

141 0.41 % and 0.21 % in the two virus-infected Ribo-seq replicates.

142 The viral genome was assembled *de novo* from RNA-seq reads and confirmed as

143 EToV, Berne isolate. A single 27694-nt contig was assembled representing almost the

144 entire viral genome. Only 18 nt at the 5' terminus and 300 nt at the 3' terminus of

145 this contig failed to assemble automatically; however these regions were clearly

146 covered by reads consistent with the reference sequence on inspection and so were

147 added manually to the consensus sequence. Four single nucleotide changes were

148 present in all reads but not the reference sequence compiled from previous

149 sequencing data, at positions 18078 (ORF 1b, C > U), 21429 (ORF S, A > U), 21814

150 (ORF S, C > A) and 25596 (ORF S, C > U). The full-length virus sequence has been

151 deposited in GenBank (Accession MG996765).

152 The distribution of reads on the virus genome and the phasing of these reads are

153 shown in Figure 2. There was good coverage across the viral genome for both RNA-

154 seq and Ribo-seq. The Ribo-seq/RNA-seq ratio along the genome was calculated

155 (Figure 2C) to estimate translation efficiency (note that this simple estimate is naive

156 since it does not account for the fact that the genomic RNA and different sgRNA

157 species overlap one another). Ribo-seq density, RNA-seq density and translational

158 efficiency were also calculated separately for each ORF (Figure 3), based on the

159 density of Ribo-seq reads in each ORF divided by the density of the RNA-seq reads

160 for either the same region (for subgenomic RNAs) or the region of the genome which

161 does not overlap the subgenomic RNAs (for genomic RNA). RNA-seq density was

162  adjusted based on the "decumulation" methodology described previously (15) (see

163  Materials and Methods) to account for the fact that not all of the RNA-seq density in

164  the 3' ORFs derives from transcripts from which the ORFs can be expressed. Ribo-seq

165  coverage is much higher towards the 3' end of the genome, particularly across the M

166  and N genes, reflecting the translation of abundant subgenomic RNAs in this region

167  (Figure 2, Figure 3). ORFs 1a and 1b contain a considerably lower density of Ribo-seq

168  reads. The relatively low translation efficiencies calculated for ORFs 1a and 1b may

169  be partly due to some gRNA being packaged (or destined for packaging) and

170  unavailable for translation but still contributing to the estimate of gRNA RNA-seq

171  density. ORF1a has a higher Ribo-seq density and a higher translational efficiency

172  than ORF1b, reflecting the proportion of ribosomes terminating at the ORF1a stop

173  codon and not undergoing the -1 frameshift into ORF1b (Figure 2, Figure 3). As

174  expected, RNA-seq density is similar across ORF1a and ORF1b, as both are present

175  only on the full-length genomic RNA (Figure 2). The region covering the HE ORF also

176  has low ribosomal coverage (Figure 2), which may be due to the fact that the EToV

177  HE gene is nonfunctional due to a large deletion including the canonical AUG (16). HE

178  is not shown in Figure 3 as the HE transcript is much less abundant than the

179  "upstream" M transcript which makes the decumulation procedure susceptible to

180  noise (see Irigoyen et al., 2016). Translational efficiency appears highest for the M

181  and S subgenomic RNAs. The high RNA-seq density in the 5' UTR may be indicative of

182  one or more defective interfering (DI) RNAs in the sample (see below). Ribosome

183  protected fragments (RPFs) were also identified mapping to the second half of the 5'

184  UTR, mostly in the +2/-1 frame with respect to ORF1a (Figure 2A).

185  To calculate the length distributions of host- and virus-mapped RPFs, we used reads

186  mapping within coding regions. After adaptor trimming, the majority (75 %) of Ribo-

187  seq reads were 27 – 29 nt in length, which is consistent with the expected size of

188  mammalian ribosome footprints. As expected, the distribution of read lengths for

189  RNA-seq was much broader, peaking between 60 and 70 nt (Supplementary Figure

190  1). For quality control, histograms of the 5' end positions of host mRNA Ribo-seq and

191  RNA-seq reads relative to initiation and termination codons were constructed

192  (Supplementary Figures 2, 3). This confirmed we had high quality RPFs arising from

193 host transcripts, with strong triplet periodicity ("phasing") and very few reads

194 mapping to 3' UTRs. As in other datasets, a ramp effect of decreased RPF density was

195 seen over a region of ~30 codons following initiation sites; but, unusually, in this

196 dataset we did not observe a density peak at the initiation site itself (cf. Irigoyen et al

197 2016). This may be due to the flash freezing without cycloheximide pretreatment

198 used for these samples, as for a later cycloheximide-treated sample this peak is

199 present (Supplementary Figure 2). Within coding sequences, the 5' ends of the

200 majority of reads from the host (65-81 %) and virus (60-75 %) mapped to the first

201 positions of codons (Supplementary Figure 4).

202 The relative RPF density allowed us to estimate the efficiency of ribosomal

203 frameshifting in the context of virus infection. After translating ORF1a, a proportion

204 of ribosomes undergo a −1 ribosomal frameshift to translate ORF1b (17). This is

205 (presumably) required to produce a specific ratio of pp1a to pp1ab, thereby

206 controlling the ratio of RNA-synthesing enzymes such as RdRp and helicase to other

207 components of the replicase complex, including the proteinases and trans-

208 membrane subunits encoded in ORF1a. The ORF1a/1b −1 ribosomal frameshifting

209 event is stimulated by a pseudoknot structure 3'-adjacent to the U_UUA_AAC

210 slippery heptanucleotide frameshift site. The efficiency of −1 ribosomal frameshifting

211 (measured by dividing the mean RPF density in ORF1b by the mean density in ORF1a)

212 was estimated to be 29.9 % for replicate one and 27.5 % for replicate 2, which is in

213 accordance with the rates measured previously outside of the context of virus

214 infection (20 – 30 %) (17).

215 **RNA sequencing indicates both discontinuous and non-discontinuous mechanisms**

216 **are utilised for N protein gene sgRNA synthesis.** RNA sequencing reads that did not

217 map to either the viral genome or host databases were analysed for containing

218 potential viral chimaeric junctions, indicative of leader-to-body joining during

219 discontinuous sgRNA synthesis (Figure 4). Relative abundances were calculated by

220 normalising read counts to the number of non-chimaeric reads spanning each

221 junction. Between the two replicates combined, 8330 reads were identified as

222 chimaeras, mapping to 2837 putative junction sites. Of these, 213 were considered

223    to be highly supported by the data, either due to being identified in at least 10

224    chimaeric reads or containing the full 5' leader and TRS sequence. Adjacent donor or

225    acceptor sites were then merged (see Materials and Methods), leaving 70 unique

226    junctions (Figure 4).

227    Three chimaeric junctions were identified where the first nucleotide of the

228    corresponding read mapped to the first nucleotide of the viral genome. Of these,

229    one junction was consistent with the previously characterised sgRNA produced via

230    discontinuous RNA synthesis encoding the S gene (280 reads, or 3 % of total

231    chimaeric reads) (12). These reads spanned the entire leader-body junction of the S

232    gene, possessing 14 - 18 nt of the 5' UTR (i.e. the actual 5'-derived sequence is at

233    least 14 nt, ACGUAUCUUUAGAA, comprising the so-called 6-nt leader, the leader TRS

234    CUUUAGA, and an additional A), followed by the stretch of ORF1b just upstream of

235    the S gene. A second set of transcripts containing 5' leader sequence was identified

236    by four unique reads starting with the 5' leader (ACGUAU) and TRS sequence

237    (CUUUAGA), where the remainder of the read mapped to the start of the N gene.

238    This indicates that, contrary to previous reports, low levels of discontinuous RNA

239    synthesis are used during production of the N gene negative-strand RNA. The final

240    chimaera which included the 6 nt leader was represented by three reads. These

241    reads included 44 - 46 nt of the 5' UTR (i.e., significantly more than the predicted

242    leader-TRS) followed by a sequence mapping to position 19987-19989 which is

243    within ORF1b.

244    A substantial number of additional chimaeric reads were identified, indicative of

245    non-TRS-driven cases of discontinuous RNA synthesis, although formally it is possible

246    that some of these are template-switching artefacts introduced during library

247    preparation and/or sequencing. Additionally, a large number of reads spanning from

248    the 5' UTR to either within the N protein gene or the 3' UTR were identified. Indeed,

249    the only junction represented by over 1000 reads spanned nucleotides 673 to 27649;

250    similarly the second most commonly identified junction spanned 687 to 27550 (642

251    reads). If chimaeric reads were predominantly a sequencing artefact, the abundance

252    of any particular chimaera would be approximately proportional to the product of

9

253    the abundances of the sequences from which the 5' and 3' ends of the chimaera are

254    derived (with some variation due to sequence-specific biases), and thus a high

255    density of chimaeras would be expected to fall entirely within the N transcript. In

256    contrast, most of the observed chimaeric reads were between N and the 5' UTR. The

257    relative paucity of reads mapping to generic locations in the ORF1ab region also

258    argues against the majority of chimaeras being simply artefactual. The 5' UTR

259    preference may be due to genome circularisation during negative-sense synthesis as

260    has been proposed for coronaviruses (18). Alternatively these may derive from

261    autonomously replicating defective interfering RNAs, rather than multiple

262    independent RNA translocation and reinitiation events. Such defective interfering

263    RNAs have been extensively analysed previously and are a common complication of

264    EToV studies (19). Consistent with the high level of 5'UTR:N chimaeric sequences,

265    there was high RNA-seq density throughout much of the 5' UTR, with the 3' extent of

266    the region of high density coinciding approximately with the region to which a large

267    number of the chimaeric 5' ends mapped (Figure 2, Figure 4).

268    **Gene expression analysis indicates multiple pathways are perturbed by EToV**

269    **infection.** The RNA-seq data were analysed to identify genes that were differentially

270    expressed between virus-infected and mock-infected ED cells. We identified 61

271    genes that were upregulated in virus-infected cells; amongst which eight gene

272    ontology (GO) terms were overrepresented, mostly related to the nucleosome or

273    immune responses (Figure 5). We found 24 genes that were downregulated in

274    infected cells, amongst which four GO terms were overrepresented, two of which

275    were related to the ribosome. We also analysed differential translational efficiency

276    (based on the RPF to mRNA ratio) between mock- and virus-infected cells. We

277    identified 22 genes that were translated more efficiently in infected cells; GO

278    analysis indicated that these genes tend to encode proteins that are involved in RNA

279    binding. Only two genes were found to be translated less efficiently in infected cells

280    compared to mock (Supplementary Table 2 and Figure 4). Note that these analyses

281    measure changes in individual genes relative to the global mean and do not inform

282    on global changes in host transcription or translation as a result of virus infection.

10

283 **Two additional proteins are translated from 5' CUG-initiated ORFs.** Our initial

284 dataset indicated an excess density of ribosomes translating within the +2/-1 frame

285 upstream of ORF1a and overlapping the 5' end of ORF1a (Figure 2A). To further

286 investigate this, we repeated the ribosome profiling using infected cells treated with

287 translation inhibitors prior to flash freezing (harringtonine, HAR, and/or

288 cycloheximide, CHX). HAR specifically arrests initiating ribosomes whilst allowing

289 "run-off" of elongating ribosomes; conversely CHX stalls elongating ribosomes whilst

290 allowing on-going accumulation at initiation sites. Our quality control analysis

291 confirmed the datasets were of similar quality to our previous experiment

292 (Supplementary Figures 1, 2 and 4) and mapping of the RPFs provided good coverage

293 of the EToV genome (Figure 6).

294 This Ribo-seq data confirmed translation of two ORFs located within the so-called 5'

295 UTR and overlapping the 5' end of ORF1a. We have termed these U1 (80 codons) and

296 U2 (258 codon). We predict that translation of both U1 and U2 is initiated from CUG

297 codons, as a close inspection indicated that ribosomes accumulated at these two

298 sites (Figure 7). It must be noted that pretreatment with CHX or HAR can introduce

299 artefacts into ribosome profiling data: CHX can lead to an excess of RPF density over

300 ~30 codons following initiation sites when cells are stressed (15, 20). It has also been

301 suggested that both drugs can promote upstream initiation due to scanning pre-

302 initiation complexes stacking behind ribosomes paused at canonical initiation sites

303 (21). However, the distance between the U1 CUG, the U2 CUG and the ORF1a

304 initiation site, besides observation of efficient translation of U2 downstream of the

305 ORF1a initiation site makes these artefacts unlikely to be significant confounding

306 factors in the case of U1 and U2.

307 Revisiting our first non-drug-treated dataset, we calculated the RPF densities and

308 translational efficiencies within the U1 and U2 ORFs (Figure 8). U1 has a higher

309 translational efficiency than any of the other ORFs translated from genomic RNA,

310 whereas U2 has a translational efficiency similar to that of ORF1a.

311 To assess the coding potential of U1, we calculated the ratio of non-synonymous to

312 synonymous substitutions (dN/dS), where dN/dS < 1 indicates selection against non-

11

313  synonymous substitutions which is a strong indicator that a sequence encodes a

314  functional protein. Application of codeml (22) to a codon alignment of eight

315  torovirus U1 nucleotide sequences resulted in a dN/dS estimate of 0.31 ± 0.08,

316  indicating that the U1 ORF is likely to encode a functional protein. MLOGD (23) uses

317  a principle similar to the dN/dS statistic but also accounts for conservative amino

318  acid substitutions (i.e. similar physico-chemical properties) being more probable

319  than non-conservative substitutions in biologically functional polypeptides. MLOGD

320  3-frame "sliding window" analysis of a full-genome alignment revealed a strong

321  coding signature in the known protein-coding ORFs (as expected) and also in the U1

322  ORF (Figure 9).

323  We previously predicted the existence of U2 via an analysis of coding potential and

324  synonymous site conservation across the two torovirus genomes available at that

325  time (24). Six additional torovirus genome sequences have now become available.

326  We therefore extended the bioinformatics analysis using all eight currently available

327  torovirus genome sequences (Figure 9). Since the U2 ORF overlaps ORF1a, leading to

328  constraint on dS, the dN/dS analysis is not appropriate for U2. MLOGD analysis

329  indicated that the U2 ORF has a higher coding potential than the corresponding part

330  of ORF1a (Figure 9). Overlapping genes are thought mainly to evolve through

331  "overprinting" of an ancestral gene by the *de novo* gene (25). The *de novo* gene

332  product is often an accessory protein and often disordered (26). Interestingly, the

333  fragment of pp1a encoded by the region of ORF1a that is overlapped by U2 has no

334  tblastn (27) nor HHpred (28) homologues outside of the *Torovirus* genus. Thus, it is

335  unclear which of U2 and the N-terminal domain of pp1a is ancestral. To provide

336  further comparative genomic evidence for the functionality of U2, we used synplot2

337  to assess conservation at synonymous sites in the ORF1a reading frame, since

338  overlapping functional elements are expected to place extra constraints on

339  synonymous site evolution (29). Consistent with the earlier 2-sequence analysis (24),

340  synplot2 revealed greatly enhanced ORF1a-frame synonymous-site conservation in a

341  region coinciding precisely with the conserved absence of stop codons that defines

342  the U2 ORF (Figure 9), with the mean rate of synonymous substitutions in that region

343  being 0.20 of the genome average. Summed over the 230-codon overlap region, the

12

344    probability that the observed level of conservation would occur by chance is $p = 6.5$ x

345    $10^{-40}$.

346    Both U1 and U2 are conserved in all eight torovirus sequences with no variation in

347    length or initiation or termination position (Supplementary Figure 5). In all

348    sequences, U1 and U2 begin with a CUG codon in a strong initiation context ('A' at -3

349    for U1, and 'A' at -3 and 'G' at +4 for U2) (30). The U1 protein is predicted to contain

350    two central transmembrane domains and has a C-terminus containing many charged

351    amino acids. The U2 protein is predicted to form alternating α helix and antiparallel β

352    sheet domains, however no structural homologs were found through searches of

353    public databases (31-33). Their function(s) will be the topic of future work.

354

355    **Discussion**

356    **RNA-seq reveals the complexity of torovirus transcription mechanisms.** The factors

357    influencing which transcriptional mechanism is utilised for the synthesis of each

358    sgRNA during torovirus replication have not been elucidated. The EToV genome

359    contains seven occurrences of the canonical TRS motif (CUUUAGA): within the 5' UTR

360    (leader TRS), the end of U1, central ORF1a, central ORF1b, and immediately before

361    the M, HE and N ORFs (Figure 1). Consistent with experimental evidence (12), we did

362    not identify any chimaeric transcripts encompassing the body TRS of M or HE, or

363    those within ORF1b or ORF1a. It appears that these sites do not stimulate

364    interruption of negative strand RNA synthesis followed by subsequent re-pairing and

365    reinitiation. The nucleotides flanking the N, M and HE TRSs are semi-conserved

366    (Supplementary Figure 6) and it has been suggested previously that the motif

367    definition should be extended to $cACN_{3-4}CUUUAGA$ to reflect this (34). It is likely that

368    these flanking nucleotides contribute to the degree of utilisation.

369    For the S gene, the chimaeric junction occurs within the run of uridines 3'-adjacent

370    to the hairpin (Figure S6I). Our results lend support to the hypothesis suggested

371    previously that a short conserved RNA hairpin, 174 nt upstream of the AUG start

13

372   codon of the EToV S protein gene, mediates discontinuous extension of negative

373   strand RNA synthesis to produce this sgRNA (12) (Supplementary Figure 6). The

374   predicted hairpin structure was not present in S gene chimaeric reads, indicating that

375   translocation may indeed be prompted by the RdRp encountering a physical block

376   after synthesising the reverse complement of the S ORF. This is in contrast to the

377   coronaviral and arteriviral mechanism, wherein RNA structures are insufficient and

378   an accompanying body TRS is required to act as a transcriptional attenuation signal,

379   prompting translocation and re-pairing of the nascent RNA. We cannot

380   unambiguously identify which nucleotides are templated before or after the

381   translocation event, as a GUUU sequence maps to genomic RNA on either side of the

382   breakpoint.

383   The leader-TRS chimaeric reads mapping to the N protein gene initially appear

384   consistent with the coronaviral and arteriviral mechanism of TRS-driven

385   discontinuous RNA synthesis. However close inspection indicated that the

386   homologous motif mediating copy-choice recombination-like translocation and re-

387   pairing of RNA strands was actually a short AGAA sequence, not the true TRS

388   (tetranucleotides underlined in Figures S6A and S6G). This would result in the

389   nascent anti-TRS mispairing with the leader TRS; two nucleotides are "skipped" once

390   reinitiation occurs. This may explain why the discontinuous mechanism is utilised so

391   rarely for this mRNA.

392   This leads to the suggestion that homology between any two sites may be sufficient

393   to induce discontinuous RNA synthesis, i.e. that provided adequate sequence

394   homology exists, the nascent RNA strand may re-pair with upstream sites within the

395   genomic RNA regardless of the presence of a predefined TRS. This is consistent with

396   the 5' UTR-ORF1b chimaeric transcripts, which again revealed a particular sequence

397   that could be templated from either region, in this case AACCUUA rather than the

398   TRS.

399   If TRS sequence-specificity is not required to stimulate EToV discontinuous RNA

400   synthesis, it is presumably constrained by alternative roles. The highly conserved

401   nature of the canonical leader, M, HE and N TRS (CUUUAG[A/U]) across all torovirus

14

402     genomes (Supplementary Figure 6) suggests it is not tolerant to mutations, however

403     this has not been formally confirmed. Lack of conservation of the EToV U1, ORF1a

404     and ORF1b TRS sequences is consistent with them not being functionally relevant.

405     Our results indicate this essential nature is likely due to a role in transcriptional

406     termination, as we did not identify a significant role of this motif in the generation of

407     chimaeric transcripts. Conversely, the upstream region of the "extended" TRS

408     ($cACN_{3-4}CUUUAGA$) is tolerant to modifications, reflecting the variable nature within

409     sequences; even when this spacer is extended to six nucleotides, transcripts are still

410     detectable at 20 % of WT levels (34). Again, this is consistent with a role in

411     termination rather than a requirement for re-pairing with upstream sequences. The

412     canonical TRS sequences also presumably contribute to subgenomic promoter

413     recognition, as the initial CAC is essential though the adenylate is the first nucleotide

414     on all positive-strand subgenomic transcripts (34). Initiation of sgRNA transcription at

415     AC dinucleotides is also found in the roniviruses (10). It may be that in these

416     *Nidovirales* families, the conserved TRS is utilised primarily for signalling

417     transcriptional termination followed by promoter recognition, and any use for

418     discontinuous RNA synthesis is merely a byproduct of RdRp promiscuity.

419     The unique combination of discontinuous and non-discontinuous mechanisms within

420     the one virus so far appears unique to the mammalian toroviruses. The one

421     bafinivirus isolated to date (white bream virus, family *Coronaviridae*, subfamily

422     *Torovirinae*, genus *Bafinivirus*) has an extended TRS sequence (CA[G/A]CACUAC)

423     which is not conserved with the mammalian toroviruses analysed in this study.

424     Bafinivirus replication produces three sgRNAs which share an identical 42-nt leader

425     also found at the far 5' terminus of the genome, indicating this species utilises

426     discontinuous RNA synthesis in a manner similar to the corona- and arteriviruses

427     (35). However there was preliminary evidence that two of the three sgRNAs exhibit

428     diversity in their junction sites, suggesting the anti-TRS may bind to multiple sites

429     within the 5' leader during strand transfer, consistent with our suggestion that whilst

430     a threshold level of homology is required this is not limited to particular primary

431     sequences. This is reflected in the fact that the bafinivirus leader-TRS is not fully

432     identical to the body TRSs.

15

433    It is not known which mechanism was utilised by the last common ancestor of

434    nidovirids, and thus which represents divergence from the original model. It has

435    been suggested that convergent evolution has resulted in the mechanism for

436    discontinuous negative strand synthesis arising multiple times within the *Nidovirales*.

437    Similarly, whether the initial role of the TRS motif was to merely stimulate the

438    attenuation of RNA synthesis or to direct the discontinuous mechanism is not

439    known. Our data suggests that transcription mechanisms in the *Nidovirales* fall into

440    multiple categories, each requiring a distinct role of the TRS: (i) homology-driven

441    reinitiation (canonical discontinuous RNA synthesis, as seen in coronaviruses and

442    arteriviruses and to a low extent, EToV N protein-coding mRNAs); (ii) structure-

443    driven discontinuous transcription (EToV S protein gene); and (iii) transcription

444    termination (EToV M, HE and the majority of N protein-coding transcripts). These

445    mechanisms all require a RdRp which is prone to translocating when even relatively

446    short homologous sequences are present, potentially leading to a large number of

447    irrelevant transcripts being produced (as previously observed in an arterivirus (36))

448    and also facilitating the production of defective interfering RNAs (34) and

449    recombinant strains (7).

450    **Effects upon the host: transcriptional and translational differential expression.** The

451    differential transcription analysis indicated that infection with EToV induces

452    increased transcription of multiple genes, the products of which are significantly

453    more likely than random to be involved in (i) nucleosome function and DNA binding,

454    and (ii) immune responses to infection than genes which were not differentially

455    transcribed. Some of the identified GO categories, including cytokine signalling,

456    innate immune responses and ribosome biogenesis have been identified in previous

457    RNA-seq analyses of various coronaviruses (37, 38). Similarly, although differential

458    translational analyses or proteomic studies have not been conducted upon

459    toroviruses, some of the identified proteins have been recognised as being

460    incorporated into nidovirid virions (for example, TCP-1 and multiple heat shock

461    proteins within arterivirus particles) (39). Others have been identified as being

462    upregulated upon infection with coronaviruses, such as the solute carrier family 25

463    members (40). Notably, both poly(C) and poly(A) binding proteins were

16

464 preferentially translated in infected cells; these have been previously identified as

465 interaction partners of arteriviral non-structural protein 1β and contribute to viral

466 RNA replication (41). It therefore appears that torovirus infection induces a similar

467 host response to many nidovirids.

468 To the best of our knowledge, this is the first analysis of differential gene expression

469 following infection with a torovirus. It would be of interest to repeat this analysis at

470 later time points, as a previous study found that EToV-mediated global inhibition of

471 host protein synthesis was only detectable at 16 h.p.i. (38). The same study found

472 induction of both the intrinsic and extrinsic apoptotic pathways was evident only by

473 24 h.p.i. (42). It is clear that the transcriptional and translational profile of the host

474 cell may differ significantly throughout the course of infection. Additionally, it must

475 be noted that the horse (*Equus caballus*) genome is not highly annotated and thus

476 many Ensembl gene identifications do not possess an annotated orthologue, a

477 limiting factor in our analysis.

478 **What is the function of U1 and U2?** The current lack of a published reverse genetics

479 system to study torovirus replication means we are unable to perform targeted

480 mutagenesis. This would enable definitive experimental confirmation that U1 and U2

481 are translated from their respective CUG codons, followed by phenotypic analysis of

482 knock-out mutants. However the comparative genomic analysis together with the

483 accumulation of ribosomes on both CUG codons is highly suggestive of this being the

484 site of initiation; CUG has previously been reported as the most commonly utilised

485 non-AUG initiation codon in mammalian systems (43). In the case of U1, the coding

486 sequence contains no AUG codons (in any frame), a situation that would facilitate

487 pre-initiation ribosomes to continue scanning to the U2 CUG and the ORF1a AUG

488 initiation sites (44). It remains a possibility that U2 translation initiates at a

489 downstream AUG, however the only in-frame AUG is located 336 nt downstream of

490 our presumed start site and is in a poor initiation context ('C' at −3) and 3' of the

491 ORF1a AUG. We are therefore confident that the CUG codons that were identified in

492 the ribosome profiling data represent the genuine translational start sites.

493    The ORFs of both U1 and U2 are intact in all torovirus genomic sequences that we

494    have analysed to date, including bovine (45, 46), caprine and porcine isolates (47).

495    Most of the U2 ORF is constrained by the fact that the sequence must also retain

496    ORF1a coding capacity in another frame. U1 is not under such limitations, although it

497    is likely that the viral genome must maintain specific 5' UTR structures to facilitate

498    viral replication. Previous investigations utilising defective interfering RNAs have

499    confirmed that no more than the first 604 nt of the 5' UTR and the entirety of the 3'

500    UTR are sufficient to allow both positive and minus strand RNA synthesis (34); it is

501    notable that this region only includes one-third of the U1 ORF (which starts at

502    nucleotide 524) and hence only this subdomain would be constrained by maintaining

503    two distinct functional roles. We suggest that the so-called 5' UTR is actually limited

504    to 523 nt preceding the CUG of U1, and the remainder of U1 and U2 is not under

505    pressure to maintain *cis*-replication elements.

506    Neither ORF could be identified within the white bream virus genome, a bafinivirus

507    that constitutes another genus within the subfamily *Torovirinae* (35), although the

508    lack of multiple bafinivirus sequences makes comparative genomic analysis

509    impossible.

510    The function(s) of the proteins encoded by both U1 and U2 remain to be elucidated.

511    Despite the relatively large size of the U2 protein (~30 kDa), after extensive database

512    searches no structural homologs were identified. By comparison, the U1 protein is

513    small (~10 kDa), highly basic (pI = 10.4) and possesses many of the predicted features

514    of a double-spanning transmembrane protein, including two hydrophobic stretches

515    separated by a 'hinge' and a predicted coiled-coil tertiary topology. Based on

516    structural similarity to known proteins, one potential function might be a virally

517    encoded ion channel (viroporin) embedded in either intracellular or plasma

518    membranes. It is possible that U1 plays a similar role in toroviruses to that of the

519    coronaviral and arteriviral E proteins, which have no known toroviral homologue.

520    The coronavirus E protein is a small transmembrane protein (~10 kDa) which

521    possesses ion channel activity and is required for virion assembly, forming a

522    pentamer that traverses the viral envelope (48). E proteins also possess a

18

523    membrane-proximal palmitoylated cysteine residue, which is a predicted (and

524    conserved) posttranslational modification for U1 (31).

525    Alternatively viroporin activity may be mediated by a small, basic double-

526    transmembrane protein, the ORF of which is embedded within the EToV N gene in

527    the +1 frame (with respect to N). An analogous "N+1" protein has been identified in

528    some group II coronaviruses and is postulated to play a structural role, however it is

529    not essential for replication (49, 50). Neither our ribosome profiling nor comparative

530    genomic analysis provides evidence that this ORF is utilised in toroviruses. We did

531    not observe ribosomes translating in this frame in either the initial dataset or the

532    drug-treated samples (although Ribo-seq may not always detect poorly translated

533    overlapping genes); further, the ORF is not preserved in all torovirus genomes.

534    Our data has revealed that the transcriptional landscape of a prototypic torovirus is

535    complex and driven by many factors beyond the canonical "multi-loci TRS" model of

536    coronaviruses. The development of a torovirus reverse genetics system would allow

537    manipulation of potential translocation-inducing sequences and allow us to elucidate

538    which features of the toroviral TRS cause them to act as terminators of RNA

539    synthesis, rather than consistently inducing homology-assisted recombination. Our

540    accompanying translational analysis has revealed two conserved novel ORFs, and has

541    shortened the EToV 5' UTR to a mere 523 nt. Together these data provide an insight

542    into the molecular biology of the replication cycle of this neglected pathogen and

543    highlight the disparities between the families of the *Nidovirales.*

544    **Materials and Methods**

545    **Virus isolates.** A plaque-purified isolate of equine torovirus, Berne strain (isolate

546    P138/72) (EToV) was kindly provided by Raoul de Groot (Utrecht University) and

547    cultured in equine dermis (ED) cells. This virus was initially isolated from a

548    symptomatic horse in 1972 (13). ED cells were maintained in Dulbecco's modified

549    Eagle's medium (Invitrogen), supplemented with 10 % foetal calf serum, 100 IU/mL

550    penicillin, 100 μg/mL streptomycin, 1 mM non-essential amino acids, 25 mM HEPES

551    and 1 % L-glutamine in a humidified incubator at 37°C with 5% $CO_2$.

552 **RNA sequencing and ribosome profiling.** ED cells were infected with EToV for 1 hour

553 (h) in serum-free media (MOI = 0.1) and flash-frozen in liquid nitrogen at 8 h post

554 infection (h.p.i.) prior to either RNA isolation or ribosome purification for profiling.

555 Cells were either not pretreated or, where stated, were treated with a final

556 concentration of 100 µg/mL cycloheximide (CHX) for 2 minutes (Sigma-Aldrich) or 2

557 µg/mL of harringtonine for 3 minutes (LKT Laboratories) followed by CHX for 2

558 minutes, before flash-freezing. RNA and ribosomes were harvested according to

559 previously published protocols (15, 51) with minor modifications. Following either

560 RPF or RNA isolation, duplex-specific nuclease was not utilised but instead rRNA was

561 depleted with the RiboZero [human/mouse/rat] kit (Illumina). Libraries were

562 prepared and sequenced using the NextSeq500 platform (Illumina).

563 **Bioinformatic analysis of Ribo-seq and RNA-seq data**. Both Ribo-seq and RNA-seq

564 reads were demultiplexed and adaptor sequences trimmed using the FASTX-Toolkit

565 (hannonlab.cshl.edu/fastx_toolkit/). Reads shorter than 25 nt after trimming were

566 discarded. Bowtie (version 1.2.1.1) databases were generated as follows. Horse

567 ribosomal RNA (rRNA) sequences were downloaded from the National Center for

568 Biotechnology Information (NCBI) Entrez Nucleotide database (accessions

569 EU081775.1, NR_046271.1, NR_046309.2, EU554425.1, XM_014728542.1 and

570 FN402126.1) (52). As the full-length virus RNA (vRNA) reference genome was not

571 available for EToV, a reference was constructed from the following overlapping

572 segments available from Entrez Nucleotide: DQ310701.1 (positions 1-14531),

573 X52374.1 (13475-21394), X52506.1 (21250-26086), X52505.1 (26054-26850),

574 X52375.1 (26784-27316) and D00563.1 (27264-279923). Horse messenger RNA

575 (mRNA) sequences from EquCab2.0 (GCF_000002305.2) were downloaded from

576 NCBI RefSeq (53). Horse non-coding RNA (ncRNA) sequences were obtained from

577 Ensembl release 89 (54) and combined with horse transfer RNA (tRNA) sequences

578 from GtRNADB (55). Horse genomic DNA (gDNA) was obtained from Ensembl release

579 89. All horse sequences were from the EquCab2.0 genome build. Trimmed reads

580 were then mapped sequentially to the rRNA, vRNA, mRNA and ncRNA databases

581 using bowtie version 1.2.1.1 (56), with parameters -v 2 --best (i.e. maximum 2

582 mismatches, report best match), with only unmapped reads passed to each following

20

583    stage. Reads that did not align to any of the aforementioned databases were then

584    mapped to the host gDNA using STAR version 2.5.4a (57), again allowing a maximum

585    of 2 mismatches per alignment. Remaining reads were classified as unmapped.

586    Ribo-seq density and RNA-seq density were calculated for each gene in the EToV

587    genome (Figure 3, Figure 8).  To normalise for different library sizes, reads per million

588    mapped reads (RPM) values were calculated using the sum of positive-sense virus

589    RNA reads and host RefSeq mRNA reads as the denominator. In order to standardise

590    the regions used to calculate RNA-seq and Ribo-seq density, the following regions

591    were selected: ORF1a, start codon (position 882) to 5' end of frameshift site

592    (position 14518); ORF1b, 3' end of frameshift site (position 14525) to 5' end of the S

593    gene hairpin (position 21118); all other ORFs, initiation codon to termination codon.

594    For U2, a region overlapping with ORF1a was used because only 46 bases are unique

595    to U2 and, for Figure 8, the ORF1a coordinates were updated to exclude the region

596    which overlaps with U2, giving a range from 1552 to 21394. In addition, for all ORFs,

597    only Ribo-seq reads mapping to the predominant phase (i.e. reads mapping to the

598    first positions of codons) were used, as this should greatly diminish misassignment of

599    ORF1a-translating ribosomes to U2 or *vice versa*. Reads mapping to the first five

600    codons at the 5' end of each region or the last six codons at the 3' end of each region

601    were excluded. For subgenomic RNAs, RNA-seq density was calculated for the same

602    regions as described for Ribo-seq. For the genomic RNA the regions for ORF1a and

603    ORF1b were combined into the interval from the start codon of ORF1a (position 882)

604    to the 5' end of the S gene hairpin (position 21118). Ribo-seq and RNA-seq densities

605    were calculated as the number of reads per million mapped reads for which the 5'

606    end maps to each region, divided by the length of the region in nt, multiplied by

607    1000 (i.e. RPKM). For RNA-seq, a decumulation strategy was used to subtract the

608    estimated RNA-seq density for longer overlapping genomic and subgenomic

609    transcripts that would contribute to the RNA-seq density measured for each of the 3'

610    ORFs: the genomic RNA-seq density was subtracted from all subgenomic densities,

611    and then the RNA-seq densities of overlapping "upstream" subgenomic transcripts

612    were iteratively subtracted from "downstream" regions (e.g. RNA-seq density in the

613    unique region of M was subtracted from HE, and this was subtracted from N).

21

614 Translation efficiency for each gene was calculated as Ribo-seq density /

615 decumulated RNA-seq density. Translational efficiencies for HE could not be

616 accurately estimated as the low expression of the HE transcript made the

617 decumulation procedure for HE susceptible to noise.

618 Read length distributions were calculated for Ribo-seq and RNA-seq reads mapping

619 to positive-sense host mRNA annotated CDSs or to the positive- or negative-sense

620 EToV genome (Supplementary Figure 1). Histograms of host mRNA Ribo-seq and

621 RNA-seq 5' end positions relative to initiation and termination codons

622 (Supplementary Figure 2, Supplementary Figure 3) were derived from reads mapping

623 to mRNAs with annotated CDSs ≥ 450 nt in length and annotated 5' and 3' UTRs ≥ 60

624 nt in length. Host mRNA Ribo-seq and RNA-seq phasing distributions (Supplementary

625 Figure 4) were calculated taking into account interior regions of annotated coding

626 ORFs only (specifically, reads for which the 5' end mapped between the first

627 nucleotide of the initiation codon and 30 nt 5' of the termination codon) in order to

628 exclude reads on or near initiation or termination codons. For viral genome coverage

629 plots, but not for meta-analyses of host RefSeq mRNA coverage, mapping positions

630 of RPF 5' ends were offset + 12 nt to approximate the location of the ribosomal P-

631 site (15).

632 **Analysis of viral transcripts.** The EToV (Berne isolate) genome sequence was

633 confirmed by *de novo* assembly of unmapped and vRNA reads from the infected

634 RNA-seq samples. Assembly was performed using Trinity (58) with the default

635 settings for stranded single ended (--SS_lib type "F") data. Viral contigs were

636 identified using BLASTN (27) against a database of EToV reference sequences based

637 on the NCBI records listed above. The viral contig was aligned to the reference using

638 the MAFFT L-INS-i method (59).

639 Chimaeric reads were classified as reads for which the entire read mapped uniquely

640 to the viral genome, with no mismatches, after adding a single breakpoint, with a

641 minimum of 12 nt mapping on either side of the breakpoint, at least 5 nt apart. To

642 identify such reads, all unmapped reads were split into two sub-reads at every

643 possible position ≥12 nt from either end and these sub-reads were mapped to the

22

644    viral genome using bowtie with no mismatches and no multimapping permitted.

645    Transcription junctions were defined as "donor/acceptor" pairs that were either

646    supported by at least 10 chimaeric reads or contained the entire 5' leader and TRS

647    sequence in the 5' segment of the read. At some positions single nucleotide

648    resolution for the chimaeric break-point could not be established; where reads were

649    found to break at adjacent possible positions these positions were merged to give a

650    short region containing the breakpoint. The number of non-chimaeric reads spanning

651    each donor and acceptor site was calculated as the number of reads which

652    overlapped the site by at least 12 nt in either direction (as chimaeric reads

653    overlapping the site by < 12 nt are not detectable). The proportion of chimaeric

654    reads at each "donor" or "acceptor" site is therefore the number of chimaeric reads

655    with a breakpoint at the site divided by this number plus the number of non-

656    chimaeric reads spanning the site (Figure 4B).

657    To visualise TRS conservation, multiple sequence alignments were generated using

658    Clustal Omega with default parameters (60). RNA structure was predicted using RNA-

659    Alifold (61) and visualised using VARNA (62).

660    **Differential gene expression analysis.** For analysis of host differential expression

661    between non-drug treated infected and mock-infected cells, all reads which did not

662    map to rRNA or vRNA were mapped to the EquCab2.0 reference genome and

663    annotations (Ensembl release 89) using STAR (57) with a maximum of two

664    mismatches and removal of non-canonical, non-annotated splice junctions. Read

665    counts were generated using HTSeq 0.8.0 (63). For differential transcription analysis,

666    gene level counts were generated across the Ensembl release 89 EquCab2.0 gtf file,

667    filtered to include only protein-coding genes. For differential translation efficiency

668    analysis only coding regions (CDS) were considered: both RNA-seq and Ribo-seq

669    counts were generated at CDS level using intersection-strict mode, based on the

670    same annotation set. Multimapping reads were excluded from both analyses.

671    Differential transcript abundance analysis was performed using the standard DESeq2

672    (64) pipeline described in the vignette. Genes to which <10 reads mapped were

673    discarded and shrinkage of $\log_2$ fold changes for lowly expressed genes was

674    performed using the lfcshrink method of DESeq2. All recommended quality control

675    plots were inspected, and no major biases were identified in the data. False

676    discovery rate (FDR) values were calculated using the R fdrtool package (65). Genes

677    with a $\log_2$ fold change >1 and an FDR less than 0.1 were considered to be

678    differentially expressed. Gene ontology (GO) term enrichment analysis (66) was

679    performed against a background of all horse protein-coding genes in the Ensembl gtf

680    using a Fisher Exact Test and corrected for multiple testing with a Bonferroni

681    correction. GO annotations for horse genes were downloaded from BiomaRt

682    (Ensembl release 90) (67). Differential translational efficiency analysis was carried

683    out using the CDS counts table, normalised using the DESeq2 "sizeFactors"

684    technique. Similar to the differential transcription analysis, genes to which <10 reads

685    mapped were discarded. Again all recommended quality control plots for DESeq2

686    were inspected and no major biases were identified in the data. Differential

687    translation efficiency analysis was performed using Xtail (68), following the standard

688    pipeline described in the vignette. *P*-values were adjusted automatically within Xtail

689    using the Benjamini–Hochberg method. Genes with a $\log_2$ fold change >1 and an

690    adjusted *p*-value less than 0.1 were considered to be differentially translated. GO

691    enrichment analysis was performed as described for the differential transcript

692    abundance analysis.

693    **Comparative genomics.** The Genbank accession numbers utilised for comparative

694    genomic analysis were as follows: DQ310701.1 (Berne virus), AY427798.1 (Breda

695    virus) (45), KR527150.1 (goat torovirus), JQ860350.1 (porcine torovirus) (47),

696    KM403390.1 (porcine torovirus) (69), LT900503.1 (porcine torovirus), LC088094.1

697    (bovine torovirus) and LC088095.1 (bovine torovirus) (46). The ratio of

698    nonsynonymous to synonymous substitution rates (dN/dS) was estimated using the

699    codeml program in the PAML package (22). The eight torovirus U1 nucleotide

700    sequences were translated, aligned as amino acids with MUSCLE (70), and the amino

701    acid alignment used to guide a codon-based nucleotide alignment (EMBOSS

702    tranalign) (71). Alignment columns with gap characters in any sequence were

703    removed, resulting in a reduction from 81 to 79 codon positions. PhyML (72) was

704    used to produce a nucleotide phylogenetic tree for the U1 alignment and, using this

24

705     tree topology, dN/dS was calculated with codeml. The standard deviation for the

706     codeml dN/dS value was estimated via a bootstrapping procedure, in which codon

707     columns of the alignment were randomly resampled (with replacement); 100

708     randomized alignments were generated, and their dN/dS values calculated with

709     codeml.

710     Coding potential within each reading frame was analysed using MLOGD (23) and

711     synonymous site conservation was analysed with synplot2 (29). For these analyses

712     we generated a codon-respecting alignment of the eight torovirus full-genome

713     sequences using a procedure described previously (29). In brief, each individual

714     genome sequence was aligned to a reference sequence using code2aln version 1.2

715     (73). Breda virus (GenBank accession AY427798) was used as reference, since unlike

716     Berne virus it contains an intact HE gene. Genomes were then mapped to reference

717     sequence coordinates by removing alignment positions that contained a gap

718     character in the reference sequence, and these pairwise alignments were combined

719     to give the multiple sequence alignment. This was analysed with MLOGD using a 40-

720     codon sliding window and a 5-codon step size. For each of the three reading frames,

721     within each window the null model is that the sequence is non-coding whereas the

722     alternative model is that the sequence is coding in the given reading frame.

723     Positive/negative values indicate that the sequences in the alignment are

724     likely/unlikely to be coding in the given reading frame. To assess conservation at

725     synonymous sites, the concatenated coding regions were extracted from the

726     alignment and analysed with synplot2.

727     **Data availability**

728     The sequencing data reported in this paper have been deposited in ArrayExpress

729     (http://www.ebi.ac.uk/arrayexpress) under the accession number E-MTAB-6656.

730     **Acknowledgements**

25

731 We thank Raoul de Groot and Arno van Vliet (Utrecht University) for providing the

732 virus isolates and helpful advice, and Polly Roy (London School of Hygiene and

733 Tropical Medicine) for ED cells.

738 **Figure Legends**

739 **Figure 1.** Schematic of the equine torovirus genome (EToV). Open reading frames

740 (ORFs) are coloured according to their respective reading frames (pink: phase 0

741 yellow: phase -1; blue: phase +1). Polyproteins pp1a and pp1ab are translated from

742 genomic RNA, with pp1ab generated via -1 programmed ribosomal frameshifting.

743 Structural proteins are translated from a series of subgenomic RNAs. Untranslated

744 regions of subgenomic RNAs are represented by black bars. The leader transcription

745 regulatory sequence (TRS) (green) and putative body TRSs (blue) are displayed below

746 the viral genome. The frameshift site and a putative RNA hairpin involved in S sgRNA

747 synthesis are indicated above the genome.

748 **Figure 2.** Read density of (A) Ribo-seq and (B) RNA-seq reads across the viral genome

749 from EToV infected cells. Red lines represent total reads per million mapped reads at

750 each position; pink: reads in phase 0; yellow: phase -1; blue: phase +1. Densities are

751 smoothed with a 15-nt running mean filter and plotted on a $\log_{10}(1+x)$ scale.

752 Negative-sense reads (grey) are displayed below the x-axis for total reads only. Each

753 line represents a single replicate. For Ribo-seq reads, a +12 nt offset has been

754 applied to read 5' end positions to map approximate P-site positions. (C) The positive

755 sense Ribo-seq/RNA-seq ratio after applying a 100-nt running mean filter to each

756 distribution. Each line represents one of the two paired Ribo-seq and RNA-seq

757 replicates.

758     **Figure 3.** Relative gene expression levels. (A) Ribo-seq density in reads per kilobase

759     per million mapped reads (RPKM) for each ORF in the EToV genome. For each ORF,

760     only reads mapping in the predominant phase (i.e. mapping to first positions of

761     codons) were included. (B) "Decumulated" RNA-seq density in RPKM for each ORF.

762     For subgenomic RNAs, density was calculated across the regions used for Ribo-seq in

763     A; for genomic RNAs the regions for ORF1a and ORF1b were combined, as these

764     ORFs are both translated from gRNA. A decumulation strategy was used to correct

765     for the fact that the measured RNA density in 3' ORFs derives from multiple 3'-

766     coterminal transcripts (see Materials and Methods). (C) Translation efficiency for

767     each gene in the EToV genome, calculated as Ribo-seq density / decumulated RNA-

768     seq density. For each ORF, the two bars represent two repeats.

769     **Figure 4.** Analysis of chimaeric viral reads. (A) Sashimi plot showing junctions in the

770     EToV genome across which chimaeric RNA-seq reads were identified in EToV

771     infected, non-drug treated samples. Chimaeric reads were defined as reads for which

772     the intact read could not be mapped but for which the 5' and 3' ends could be

773     uniquely mapped to non-contiguous regions of the EToV genome. Junctions that

774     were either covered by at least 10 chimaeric reads (grey) and/or for which the 5'

775     section of the read contained the full 5' leader sequence and leader TRS (red) were

776     identified and adjacent positions merged. These junctions are shown as curved lines

777     connecting the position of the 3' end of the 5' mapped segment of the read and the

778     5' end of the 3' mapped segment of the read. The apical height of each curved line

779     shows the absolute number of reads spanning this junction on a $\log_{10}(1+x)$ scale. (B)

780     Inverted bar chart showing, for the 5' (orange) and 3' (blue) breakpoints for each

781     junction, the  number of chimaeric reads as a fraction of the total number of

782     chimaeric and non-chimaeric reads at each site.

783     **Figure 5.** Volcano plots showing the results of (A) differential transcription analysis

784     performed using DESeq2 (64) and (B) differential translation efficiency analysis

785     performed using Xtail (68), between cells infected with EToV (infected) and

786     uninfected cells (mock). Genes which were expressed at significantly higher levels

787     (FDR ≤ 0.05 and absolute $\log_2$(fold change) ≥ 1) in infected cells are highlighted in

788    pink (transcription, A) and blue (translational efficiency, B). Genes which were

789    expressed at significantly higher levels in mock infected cells are highlighted in green

790    (transcription, A) and orange (translational efficiency, B). The five most significant

791    genes in each category are labelled with the gene symbol where available and

792    otherwise with the Ensembl gene ID. (C) Absolute $\log_2$(fold change) for all gene

793    ontology (GO) terms which were significantly overrepresented compared to a

794    background of all horse protein-coding genes for genes significantly more

795    transcribed in infected cells (pink), genes significantly more efficiently translated in

796    infected cells (blue) and genes significantly more transcribed in mock cells (green).

797    No terms were identified for genes significantly more efficiently translated in mock

798    cells.


799    **Figure 6.** Read density of Ribo-seq reads along the viral genome for EToV infected

800    cells pretreated with (A) cycloheximide or (B) harringtonine. Red lines represent total

801    reads per million mapped reads (RPM) at each position. Densities are smoothed with

802    a 15-nt running mean filter and plotted on a $\log_{10}(1+x)$ scale. Negative-sense reads

803    (grey) are displayed below the x-axis. Each line represents a single replicate. A +12 nt

804    offset has been applied to read 5' end positions to map approximate P-site positions.


805    **Figure 7.** Read density of Ribo-seq reads across (A) U1, U2 and ORF1a; and (B) the U1

806    ORF and surrounding regions, for EToV infected cells with no drug treatment or with

807    cycloheximide or harringtonine pretreatment. Pink: reads in phase 0; yellow: phase -

808    1; blue: phase +1. Graphs show total reads per million mapped reads (RPM) at each

809    position. In (A) densities are smoothed with a 15-nt running mean filter while (B)

810    shows the RPM counts at single-nt resolution. Each plot represents a single replicate.

811    A +12 nt offset has been applied to read 5' end positions to map approximate P-site

812    positions.


813    **Figure 8.** Relative translation efficiencies for U1, U2, ORF1a and ORF1b. To reduce

814    misassignment of reads in the U2/ORF1a overlap region, for all ORFs only reads

815    mapping in the predominant phase (i.e. mapping to first positions of codons) were

816    included. Ribo-seq densities were divided by the ORF1ab RNA-seq densities for the

817    corresponding paired sample. For each ORF, the two bars represent two repeats.

818    **Figure 9.** Coding potential statistics for the torovirus genome. A map of the torovirus

819    genome is shown at top. Breda virus (AY427798.1) was used as the reference

820    genome for this analysis since EToV has a deletion in the HE gene. In Breda virus, U1

821    is in-frame with ORF1a due to a 2-nt insertion relative to EToV in the short non-

822    coding region between U1 and U2. The next four panels show an analysis of

823    synonymous site conservation in the concatenated coding ORFs (with the reading

824    frame of the longer ORF being used wherever two ORFs overlap). Red lines show the

825    probability that the degree of conservation within a given window (25- or 65-codons

826    as indicated) could be obtained under a null model of neutral evolution at

827    synonymous sites, whereas brown lines depict the absolute amount of conservation

828    as represented by the ratio of the observed number of substitutions within a given

829    window to the number expected under the null model. Greatly enhanced

830    synonymous site conservation is seen in the region of ORF1a that is overlapped by

831    the U2 ORF. The next three panels show MLOGD coding potential scores and stop

832    codon plots for each of the three reading frames. The positions of stop codons are

833    shown for each of the eight torovirus sequences mapped onto the Breda virus

834    reference sequence coordinates. Note the conserved absence of stop codons in the

835    U1 and U2 ORFs. MLOGD was applied in a 40-codon sliding-window (5-codon step

836    size). Positive scores indicate that the sequence is likely to be coding in the given

837    reading frame. Note the positive scores within the U1 and U2 ORFs besides the

838    previously known ORFs. The bottom panel (green line) indicates the total amount of

839    phylogenetic divergence contributing to the analyses at each alignment position

840    (regions containing alignment gaps have reduced summed divergence leading to

841    reduced statistical power). Pink regions in the stop codon plots (e.g. EToV sequence

842    in the HE region) indicate regions excluded from the analyses due to poor or locally

843    out-of-frame mapping to the Breda reference sequence (see Firth, 2014 for details).

844    **Supplementary Table 1.** Read counts for each sample. Too short, adaptor only, rRNA

845    forward/reverse, mRNA forward/reverse, ncRNA forward/reverse, gDNA

846    forward/reverse, and vRNA forward/reverse reads were summed to give the total

847    mapped read count. Remaining reads were classified as unmapped.

29

848 **Supplementary Table 2.** Gene descriptions (Ensembl gene identifiers and gene
849 symbols) for transcripts which were differentially transcribed or translated, in EToV
850 compared to mock infected cells.

851 **Supplementary Figure 1.** Comparison of read length distribution for reads mapping
852 to EToV in infected cells (orange), host mRNAs in non-infected cells (blue) and host
853 mRNAs in infected cells (red) for (A) Ribo-seq data in non-drug treated cells; (B) RNA-
854 seq data in non-drug treated cells; (C) Ribo-seq data in cycloheximide-treated cells;
855 and (D) Ribo-seq data in harringtonine-treated cells.

856 **Supplementary Figure 2.** Histograms of Ribo-seq read 5' end positions (nt) relative to
857 annotated initiation (left) and termination (right) sites, summed across all host
858 mRNAs. Bars are coloured by phase relative to the first base of the start codon (pink:
859 phase 0; blue: phase +1; yellow: phase -1). Histograms are scaled so that the
860 maximum value is 1. For clarity, the y-axis is cropped at 0.3 for non-drug treated and
861 0.1 for drug treated cells; bars which extended beyond this point are marked with an
862 asterisk (*).

863 **Supplementary Figure 3.** Histograms of RNA-seq read 5' end positions (nt) relative to
864 annotated initiation (left) and termination (right) sites, summed across all host
865 mRNAs. Bars are coloured by phase relative to the first base of the start codon (pink:
866 phase 0; blue: phase +1; yellow: phase -1). Histograms are scaled so that the
867 maximum value is 1.

868 **Supplementary Figure 4.** Phasing of the 5' ends of reads (pink: phase 0; blue: phase
869 +1; yellow: phase -1) for (A) Ribo-seq reads mapping to host mRNA coding regions,
870 (B) RNA-seq reads mapping to host mRNA coding regions, (C) Ribo-seq reads
871 mapping to virus mRNA coding regions and (D) RNA-seq reads mapping to virus
872 mRNA coding regions.

873 **Supplementary Figure 5.** Conservation of uORF1 and uORF2 in the eight publicly
874 available torovirus genomes. Individual amino acid residues are coloured according
875 to their biochemical properties.

876    **Supplementary Figure 6.** Conservation of TRSs and regulatory structures in the eight

877    publicly available torovirus genomes. Regions were selected based on the presence

878    of a putative TRS in the EToV genome. The TRS and six flanking nucleotides are

879    displayed; putative TRS nucleotides are highlighted in red. Nucleotide conservation

880    between all eight sequences is indicated by an asterisk (*). The predicted hairpin

881    structure (I) is based upon nucleotide conservation across all eight genomes. Variant

882    nucleotides are circled in either red (covariance indicates the predicted pairing may

883    occur in all but one genome) or blue (variable). R indicates a purine exists in all

884    genomes.

885

886    **References**

887

888

889    1.    Lauber C, Ziebuhr J, Junglen S, Drosten C, Zirkel F, Nga PT, Morita K,
890        Snijder EJ, Gorbalenya AE. 2012. Mesoniviridae: a proposed new family in
891        the order Nidovirales formed by a single species of mosquito-borne
892        viruses. Arch Virol 157:1623-8.
893    2.    Brown DW, Beards GM, Flewett TH. 1987. Detection of Breda virus
894        antigen and antibody in humans and animals by enzyme immunoassay. J
895        Clin Microbiol 25:637-40.
896    3.    Hoet AE, Saif LJ. 2004. Bovine torovirus (Breda virus) revisited. Anim
897        Health Res Rev 5:157-71.
898    4.    Hanke D, Pohlmann A, Sauter-Louis C, Hoper D, Stadler J, Ritzmann M,
899        Steinrigl A, Schwarz BA, Akimkin V, Fux R, Blome S, Beer M. 2017. Porcine
900        Epidemic Diarrhea in Europe: In-Detail Analyses of Disease Dynamics and
901        Molecular Epidemiology. Viruses 9.
902    5.    Pignatelli J, Grau-Roma L, Jimenez M, Segales J, Rodriguez D. 2010.
903        Longitudinal serological and virological study on porcine torovirus
904        (PToV) in piglets from Spanish farms. Vet Microbiol 146:260-8.
905    6.    Alonso-Padilla J, Pignatelli J, Simon-Grife M, Plazuelo S, Casal J, Rodriguez
906        D. 2012. Seroprevalence of porcine torovirus (PToV) in Spanish farms.
907        BMC Res Notes 5:675.
908    7.    Smits SL, Lavazza A, Matiz K, Horzinek MC, Koopmans MP, de Groot RJ.
909        2003. Phylogenetic and evolutionary relationships among torovirus field
910        variants: evidence for multiple intertypic recombination events. J Virol
911        77:9567-77.
912    8.    Pasternak AO, Spaan WJ, Snijder EJ. 2006. Nidovirus transcription: how to
913        make sense...? J Gen Virol 87:1403-21.
914    9.    van Marle G, Dobbe JC, Gultyaev AP, Luytjes W, Spaan WJ, Snijder EJ. 1999.
915        Arterivirus discontinuous mRNA transcription is guided by base pairing

916     between sense and antisense transcription-regulating sequences. Proc
917     Natl Acad Sci U S A 96:12056-61.
918  10. Cowley JA, Dimmock CM, Walker PJ. 2002. Gill-associated nidovirus of
919     Penaeus monodon prawns transcribes 3'-coterminal subgenomic mRNAs
920     that do not possess 5'-leader sequences. J Gen Virol 83:927-35.
921  11. Zirkel F, Roth H, Kurth A, Drosten C, Ziebuhr J, Junglen S. 2013.
922     Identification and characterization of genetically divergent members of
923     the newly established family Mesoniviridae. J Virol 87:6346-58.
924  12. van Vliet AL, Smits SL, Rottier PJ, de Groot RJ. 2002. Discontinuous and
925     non-discontinuous subgenomic RNA transcription in a nidovirus. EMBO J
926     21:6571-80.
927  13. Weiss M, Steck F, Horzinek MC. 1983. Purification and partial
928     characterization of a new enveloped RNA virus (Berne virus). J Gen Virol
929     64 (Pt 9):1849-58.
930  14. Kuwabara M, Wada K, Maeda Y, Miyazaki A, Tsunemitsu H. 2007. First
931     isolation of cytopathogenic bovine torovirus in cell culture from a calf
932     with diarrhea. Clin Vaccine Immunol 14:998-1004.
933  15. Irigoyen N, Firth AE, Jones JD, Chung BY, Siddell SG, Brierley I. 2016. High-
934     Resolution Analysis of Coronavirus Gene Expression by RNA Sequencing
935     and Ribosome Profiling. PLoS Pathog 12:e1005473.
936  16. Snijder EJ, den Boon JA, Horzinek MC, Spaan WJ. 1991. Comparison of the
937     genome organization of toro- and coronaviruses: evidence for two
938     nonhomologous RNA recombination events during Berne virus evolution.
939     Virology 180:448-52.
940  17. Snijder EJ, den Boon JA, Bredenbeek PJ, Horzinek MC, Rijnbrand R, Spaan
941     WJ. 1990. The carboxyl-terminal part of the putative Berne virus
942     polymerase is expressed by ribosomal frameshifting and contains
943     sequence motifs which indicate that toro- and coronaviruses are
944     evolutionarily related. Nucleic Acids Res 18:4535-42.
945  18. Yang D, Leibowitz JL. 2015. The structure and functions of coronavirus
946     genomic 3' and 5' ends. Virus Res 206:120-33.
947  19. Snijder EJ, den Boon JA, Horzinek MC, Spaan WJ. 1991. Characterization of
948     defective interfering RNAs of Berne virus. J Gen Virol 72 ( Pt 7):1635-43.
949  20. Gerashchenko MV, Gladyshev VN. 2014. Translation inhibitors cause
950     abnormalities in ribosome profiling experiments. Nucleic Acids Res
951     42:e134.
952  21. Andreev DE, O'Connor PB, Loughran G, Dmitriev SE, Baranov PV, Shatsky
953     IN. 2017. Insights into the mechanisms of eukaryotic translation gained
954     with ribosome profiling. Nucleic Acids Res 45:513-526.
955  22. Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. Mol
956     Biol Evol 24:1586-91.
957  23. Firth AE, Brown CM. 2006. Detecting overlapping coding sequences in
958     virus genomes. BMC Bioinformatics 7:75.
959  24. Firth AE, Atkins JF. 2009. A case for a CUG-initiated coding sequence
960     overlapping torovirus ORF1a and encoding a novel 30 kDa product. Virol J
961     6:136.
962  25. Keese PK, Gibbs A. 1992. Origins of genes: "big bang" or continuous
963     creation? Proc Natl Acad Sci U S A 89:9489-93.

26. Rancurel C, Khosravi M, Dunker AK, Romero PR, Karlin D. 2009. Overlapping genes produce proteins with unusual sequence properties and offer insight into de novo protein creation. J Virol 83:10719-36.

27. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. J Mol Biol 215:403-10.

28. Soding J, Biegert A, Lupas AN. 2005. The HHpred interactive server for protein homology detection and structure prediction. Nucleic Acids Res 33:W244-8.

29. Firth AE. 2014. Mapping overlapping functional elements embedded within the protein-coding regions of RNA viruses. Nucleic Acids Res 42:12425-39.

30. Kozak M. 1986. Point mutations define a sequence flanking the AUG initiator codon that modulates translation by eukaryotic ribosomes. Cell 44:283-92.

31. Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJ. 2015. The Phyre2 web portal for protein modeling, prediction and analysis. Nat Protoc 10:845-58.

32. Haas J, Roth S, Arnold K, Kiefer F, Schmidt T, Bordoli L, Schwede T. 2013. The Protein Model Portal--a comprehensive resource for protein structure and model information. Database (Oxford) 2013:bat031.

33. McGuffin LJ, Bryson K, Jones DT. 2000. The PSIPRED protein structure prediction server. Bioinformatics 16:404-5.

34. Smits SL, van Vliet AL, Segeren K, el Azzouzi H, van Essen M, de Groot RJ. 2005. Torovirus non-discontinuous transcription: mutational analysis of a subgenomic mRNA promoter. J Virol 79:8275-81.

35. Schutze H, Ulferts R, Schelle B, Bayer S, Granzow H, Hoffmann B, Mettenleiter TC, Ziebuhr J. 2006. Characterization of White bream virus reveals a novel genetic cluster of nidoviruses. J Virol 80:11598-609.

36. Di H, Madden JC, Jr., Morantz EK, Tang HY, Graham RL, Baric RS, Brinton MA. 2017. Expanded subgenomic mRNA transcriptome and coding capacity of a nidovirus. Proc Natl Acad Sci U S A 114:E8895-E8904.

37. Cong F, Liu X, Han Z, Shao Y, Kong X, Liu S. 2013. Transcriptome analysis of chicken kidney tissues following coronavirus avian infectious bronchitis virus infection. BMC Genomics 14:743.

38. Raaben M, Groot Koerkamp MJ, Rottier PJ, de Haan CA. 2007. Mouse hepatitis coronavirus replication induces host translational shutoff and mRNA decay, with concomitant formation of stress granules and processing bodies. Cell Microbiol 9:2218-29.

39. Zhang C, Xue C, Li Y, Kong Q, Ren X, Li X, Shu D, Bi Y, Cao Y. 2010. Profiling of cellular proteins in porcine reproductive and respiratory syndrome virus virions by proteomics analysis. Virol J 7:242.

40. VanLeuven JT, Ridenhour BJ, Gonzalez AJ, Miller CR, Miura TA. 2017. Lung epithelial cells have virus-specific and shared gene expression responses to infection by diverse respiratory viruses. PLoS One 12:e0178408.

41. Beura LK, Dinh PX, Osorio FA, Pattnaik AK. 2011. Cellular poly(c) binding proteins 1 and 2 interact with porcine reproductive and respiratory syndrome virus nonstructural protein 1beta and support viral replication. J Virol 85:12939-49.

1012 42.  Maestre AM, Garzon A, Rodriguez D. 2011. Equine torovirus (BEV)
1013      induces caspase-mediated apoptosis in infected cells. PLoS One 6:e20972.
1014 43.  Touriol C, Bornes S, Bonnal S, Audigier S, Prats H, Prats AC, Vagner S.
1015      2003. Generation of protein isoform diversity by alternative initiation of
1016      translation at non-AUG codons. Biol Cell 95:169-78.
1017 44.  Firth AE, Brierley I. 2012. Non-canonical translation in RNA viruses. J Gen
1018      Virol 93:1385-409.
1019 45.  Draker R, Roper RL, Petric M, Tellier R. 2006. The complete sequence of
1020      the bovine torovirus genome. Virus Res 115:56-68.
1021 46.  Ito M, Tsuchiaka S, Naoi Y, Otomaru K, Sato M, Masuda T, Haga K, Oka T,
1022      Yamasato H, Omatsu T, Sugimura S, Aoki H, Furuya T, Katayama Y, Oba M,
1023      Shirai J, Katayama K, Mizutani T, Nagai M. 2016. Whole genome analysis of
1024      Japanese bovine toroviruses reveals natural recombination between
1025      porcine and bovine toroviruses. Infect Genet Evol 38:90-95.
1026 47.  Sun H, Lan D, Lu L, Chen M, Wang C, Hua X. 2014. Molecular
1027      characterization and phylogenetic analysis of the genome of porcine
1028      torovirus. Arch Virol 159:773-8.
1029 48.  Ruch TR, Machamer CE. 2012. The coronavirus E protein: assembly and
1030      beyond. Viruses 4:363-82.
1031 49.  Senanayake SD, Hofmann MA, Maki JL, Brian DA. 1992. The nucleocapsid
1032      protein gene of bovine coronavirus is bicistronic. J Virol 66:5277-83.
1033 50.  Fischer F, Peng D, Hingley ST, Weiss SR, Masters PS. 1997. The internal
1034      open reading frame within the nucleocapsid gene of mouse hepatitis virus
1035      encodes a structural protein that is not essential for viral replication. J
1036      Virol 71:996-1003.
1037 51.  Irigoyen N, Dinan AM, Brierley I, Firth AE. 2018. Ribosome profiling of the
1038      retrovirus murine leukemia virus. Retrovirology 15:10.
1039 52.  Coordinators NR. 2016. Database resources of the National Center for
1040      Biotechnology Information. Nucleic Acids Res 44:D7-19.
1041 53.  Pruitt KD, Tatusova T, Brown GR, Maglott DR. 2012. NCBI Reference
1042      Sequences (RefSeq): current status, new features and genome annotation
1043      policy. Nucleic Acids Res 40:D130-5.
1044 54.  Flicek P, Amode MR, Barrell D, Beal K, Brent S, Carvalho-Silva D, Clapham
1045      P, Coates G, Fairley S, Fitzgerald S, Gil L, Gordon L, Hendrix M, Hourlier T,
1046      Johnson N, Kahari AK, Keefe D, Keenan S, Kinsella R, Komorowska M,
1047      Koscielny G, Kulesha E, Larsson P, Longden I, McLaren W, Muffato M,
1048      Overduin B, Pignatelli M, Pritchard B, Riat HS, Ritchie GR, Ruffier M,
1049      Schuster M, Sobral D, Tang YA, Taylor K, Trevanion S, Vandrovcova J,
1050      White S, Wilson M, Wilder SP, Aken BL, Birney E, Cunningham F, Dunham
1051      I, Durbin R, Fernandez-Suarez XM, Harrow J, Herrero J, Hubbard TJ, et al.
1052      2012. Ensembl 2012. Nucleic Acids Res 40:D84-90.
1053 55.  Chan PP, Lowe TM. 2016. GtRNAdb 2.0: an expanded database of transfer
1054      RNA genes identified in complete and draft genomes. Nucleic Acids Res
1055      44:D184-9.
1056 56.  Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-
1057      efficient alignment of short DNA sequences to the human genome.
1058      Genome Biol 10:R25.

1059  57.  Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P,
1060       Chaisson M, Gingeras TR. 2013. STAR: ultrafast universal RNA-seq aligner.
1061       Bioinformatics 29:15-21.
1062  58.  Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis
1063       X, Fan L, Raychowdhury R, Zeng Q, Chen Z, Mauceli E, Hacohen N, Gnirke
1064       A, Rhind N, di Palma F, Birren BW, Nusbaum C, Lindblad-Toh K, Friedman
1065       N, Regev A. 2011. Full-length transcriptome assembly from RNA-Seq data
1066       without a reference genome. Nat Biotechnol 29:644-52.
1067  59.  Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment
1068       software version 7: improvements in performance and usability. Mol Biol
1069       Evol 30:772-80.
1070  60.  Sievers F, Higgins DG. 2014. Clustal omega. Curr Protoc Bioinformatics
1071       48:3 13 1-16.
1072  61.  Bernhart SH, Hofacker IL, Will S, Gruber AR, Stadler PF. 2008. RNAalifold:
1073       improved consensus structure prediction for RNA alignments. BMC
1074       Bioinformatics 9:474.
1075  62.  Darty K, Denise A, Ponty Y. 2009. VARNA: Interactive drawing and editing
1076       of the RNA secondary structure. Bioinformatics 25:1974-5.
1077  63.  Anders S, Pyl PT, Huber W. 2015. HTSeq--a Python framework to work
1078       with high-throughput sequencing data. Bioinformatics 31:166-9.
1079  64.  Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change
1080       and dispersion for RNA-seq data with DESeq2. Genome Biol 15:550.
1081  65.  Strimmer K. 2008. fdrtool: a versatile R package for estimating local and
1082       tail area-based false discovery rates. Bioinformatics 24:1461-2.
1083  66.  Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP,
1084       Dolinski K, Dwight SS, Eppig JT, Harris MA, Hill DP, Issel-Tarver L,
1085       Kasarskis A, Lewis S, Matese JC, Richardson JE, Ringwald M, Rubin GM,
1086       Sherlock G. 2000. Gene ontology: tool for the unification of biology. The
1087       Gene Ontology Consortium. Nat Genet 25:25-9.
1088  67.  Durinck S, Spellman PT, Birney E, Huber W. 2009. Mapping identifiers for
1089       the integration of genomic datasets with the R/Bioconductor package
1090       biomaRt. Nat Protoc 4:1184-91.
1091  68.  Xiao Z, Zou Q, Liu Y, Yang X. 2016. Genome-wide assessment of differential
1092       translations with ribosome profiling data. Nat Commun 7:11194.
1093  69.  Anbalagan S, Peterson J, Wassman B, Elston J, Schwartz K. 2014. Genome
1094       sequence of torovirus identified from a pig with porcine epidemic
1095       diarrhea virus from the United States. Genome Announc 2.
1096  70.  Edgar RC. 2004. MUSCLE: multiple sequence alignment with high
1097       accuracy and high throughput. Nucleic Acids Res 32:1792-7.
1098  71.  Rice P, Longden I, Bleasby A. 2000. EMBOSS: the European Molecular
1099       Biology Open Software Suite. Trends Genet 16:276-7.
1100  72.  Guindon S, Gascuel O. 2003. A simple, fast, and accurate algorithm to
1101       estimate large phylogenies by maximum likelihood. Syst Biol 52:696-704.
1102  73.  Stocsits RR, Hofacker IL, Fried C, Stadler PF. 2005. Multiple sequence
1103       alignments of partially coding nucleic acid sequences. BMC Bioinformatics
1104       6:160.
1105

**Figure 1**

**Figure 2**

**Figure 3**

**Figure 4**

**Figure 5**

**Figure 6**

**Figure 7**

**Figure 8**

**Figure 9**

**Supplementary Figure 1**

## Supplementary Figure 2



Ribo-Seq Infected Untreated Rep 1

Ribo-Seq Infected Untreated Rep 2

Ribo-Seq Mock Untreated Rep 1

Ribo-Seq Mock Untreated Rep 2

Ribo-Seq Infected Cycloheximide Treated

Ribo-Seq Infected Harringtonine Treated

## Supplementary Figure 3

## Supplementary Figure 4

**Supplementary Figure 5**

**uORF1**

```
EToV      MLFYIFKFCGFTLFTYGWQFVWLFAQFFVIPIFTVFLLFTVKAIFYLLRLVELSFSTLIIWLIEVLRKRK
KM403390  MIFFILKFCVYTLFSYGWQFVWLFAQFLLIPIFTVLLVFTAKAVFYLLKLLEAAFTTLVLILIDKLKGAR
LC088094  MILFILKFCVYTLFSYGWQFVWLFAQFLLIPIFTVLLVFTAKAVFYLLKLLEAAFTILVLVLIEKLKGVR
LC088095  MILFILKFCVYTLFSYGWQFVWLFAQFLLIPIFTVLLVFTAKAVFYLLKLLEAAFTILVLVLIEKLKGVR
JQ860350  MILFILKFCVYTLFSYGWQFVWLFAQFLLIPIFTVLLVFTAKAVFYLLKLLEAAFTTLVLISIEKLRGVR
LT900503  MILFILKFCVYTLFSYGWQFVWLFAQFLLIPIFTVLLVFTAKAVFYLLKLLEAAFTTLVLILIEKLKGVR
AY427798  MIFIVTKFCVYTLFSYGWQFVWLFAQFLLIPIFTVLLVFSAKAVFHILSLLEVTFTTSVLFLVSKVKSFR
KR527150  MIFIVVKFCVYTLFSYGWQFVWLFAQFLLIPIFTVLLVYSAKAVFYILSLVEAIFTTFILFLISKVKSLR
          *::  :  *** .:***.*************::*****.*:::.**:*::* *.*  *:  ::  :. :.  .

EToV      QQFRRRSLDV
KM403390  KQRRRPSCGV
LC088094  KQRRRPSCGV
LC088095  KQRRRPSCGV
JQ860350  KQHRRPSCGV
LT900503  KQHRRPSCGV
AY427798  KKHRRPSCGV
KR527150  KHHRRHSCGV
          ::  ** * .*
```

**uORF2**

```
EToV      MANKYQVIDSLWSETYEYQFQYFGHPFKNVQDLKKQHQRNRAAFVLKYLGPNFQVPAFGPVFRYTRNNGI
AY427798  MANKYQVIDSLWSETYEYQFAYFGHPYKNVQDLKRAHQRNRAAFVLKYLGPNFQVPAFGPVFRYTTKSGI
KR527150  MANKYQVIDSLWSETYEYQFAYFGHPYKDVQDLKRAHQRNRAAFVLKYLGPSFQVPAFGPVFRYTTKPGI
JQ860350  MANKYQVIDSLWSETYEYQFAYFGHPYKNVQDLKKAHQRNRAAFVLKYLGPNFQVPVFGPVFRYTTKPGI
LC088094  MANKYQVVDSLWSETYEYQFAYFGHPYKNVQDLKKAHQRNRAAFVLKYLGPNFQVPAFGPVFRYTTKPGI
LC088095  MANKYQVVDSLWSETYEYQFAYFGHPYKNVQDLKKAHQRNRAAFVLKYLGPNFQVPAFGPVFRYTTKPGI
LT900503  MANKYQVIDSLWSETYEYQFAYFGHPYKNVQDLKKAHQRNRAAFVLKYLGPNFQVPAFGPVLRYTTKPGI
KM403390  MANKYQVIDSLWSETYEYQFAYFGHPYKNVQDLKKAHQRNRAAFVLKYLGPNFQVPAFGPVFRYTTKPGI
          *******:.************* *****.*.***** .***************.****.****.***  : **

EToV      AFKNGAIYLGVSELGTQIHINPLQLFTKFTVTCDEHLVHPVQMDYRVYLECEGSVGERIVQGVSAFERYY
AY427798  SFKDGSIYLGVTDFGTQIHINPLQLFTKFAITCPEHLIHPVQMDYRVYLETEGSFGERIVQGVSSFERFY
KR527150  SFKDGSIYLGVTDFGTQIHINPLQLFTKFAVTCPEHLIHPVQMDYRVYLETEGSFGERIVQGVSSFERFY
JQ860350  TFKDGSIYLGVTDFGTQVHINPLQLFTKFAITCPEHLIHPVQMDYRVYLETEGSVGERIVQGVSAFERYY
LC088094  TFKDGSIYLGTTDFGTQVHINPLQLFTKFAVTCPEHLIHPVQMDYRVYLETEGSFGERIVQGVSAFERFY
LC088095  TFKDGSIYLGTTDFGTQVHINPLQLFTKFAVTCPEHLIHPVQMDYRVYLETEGSFGERIVQGVSAFERFY
LT900503  TFKDGSVYLGTDFGTQVHINPLQLFTKFAITCPEHLIHPVQMDYRVYLETEGSFGERIVQGVSAFERFY
KM403390  TFKDGSIYLGITDFGTQVHINPLQLFTKFAVTCPEHLIHPVQMDYRVYLETEGSFGERIVQGVSAFERFY
          .*.**.*:.*** ::.***.************::** ***.*********** *** .*********.***.*

EToV      PKKQLCGAITADPFNFDWERNIHNYYFTRNTLRYGTKYYQLCGKHLIERSSGIERTGILPRILSECQLPI
AY427798  PKRQLCGVIIDDPFSFDWAGNIHNYYFTRNVLRYGTKLYQVNGNRLIERSSGIERSDVLPRILSECQLPI
KR527150  PKRQLCGIIIDDPFSFDWAGNIHNYYFTRNVLRYGTKLYQVNGNRLIERSSGIERSDILPRILSECQLPI
JQ860350  PRKQLCGTIVNDPFTFDWAGNIHNYYFTRNVLRYGTKLYQVNGNKLIERCSGIERSDILPRILSECQLPI
LC088094  PKRQLCGSIVSDPFTFDWAGNIHNYYFTRNVLRYGTKLYQVNGNKLIERCSGIERSDILPRILSECQLPV
LC088095  PKRQLCGSIVSDPFTFDWAGNIHNYYFTRNVLRYGTKLYQVNGNKLIERCSGIERSDILPRILSECQLPI
LT900503  PKRQLCGTIVNDPFTFDWAGNIHNYYFTRNVLRYGTKLYQVNGNKLIERCSGIERSDILPRILSECQLPI
KM403390  PKRQLCGTIVNDPFTFDWAGNIHNYYFTRNVLRYGTKLYQVNGNKLIERCSGIERSDILPRILSECQLPI
          *.:.**** *  ***.*** **********.****** **.:.**:.****.*****.::.***********.

EToV      LDTTASAAEFDEDVICCGFESLDITEHPTLAETQPFPWRHFSQLCNSN
AY427798  LDTTPTPSECDEDVICCGFESLDIREYPALAETQPFPWRHFSQLHLND
KR527150  LDTTPTPSERDEDVICCGFESLDIREYPALAETQPFPWRHFSQLHLND
JQ860350  LDTTPTPAERDEDVICCDFESLDIREYPALAETKPFPWRHFSQLHLND
LC088094  LDTTPTPSERDEDVICCDFESLDIREYPALAETKPFPWRHFSQLHLND
LC088095  LDTTPTPSERDEDVICCDFESLDIREYPALAETKPFPWRHFSQLHLND
LT900503  LDTTPTPAERDEDVICCDFESLDIREYPALAETKPFPWRHFSQLHLND
KM403390  LDTTPTPAERDEDVICCDFESLDIREYPALAETKPFPWRHFSQLHLND
          ****.::.* *******.****** *.:.****.**********   .:
```
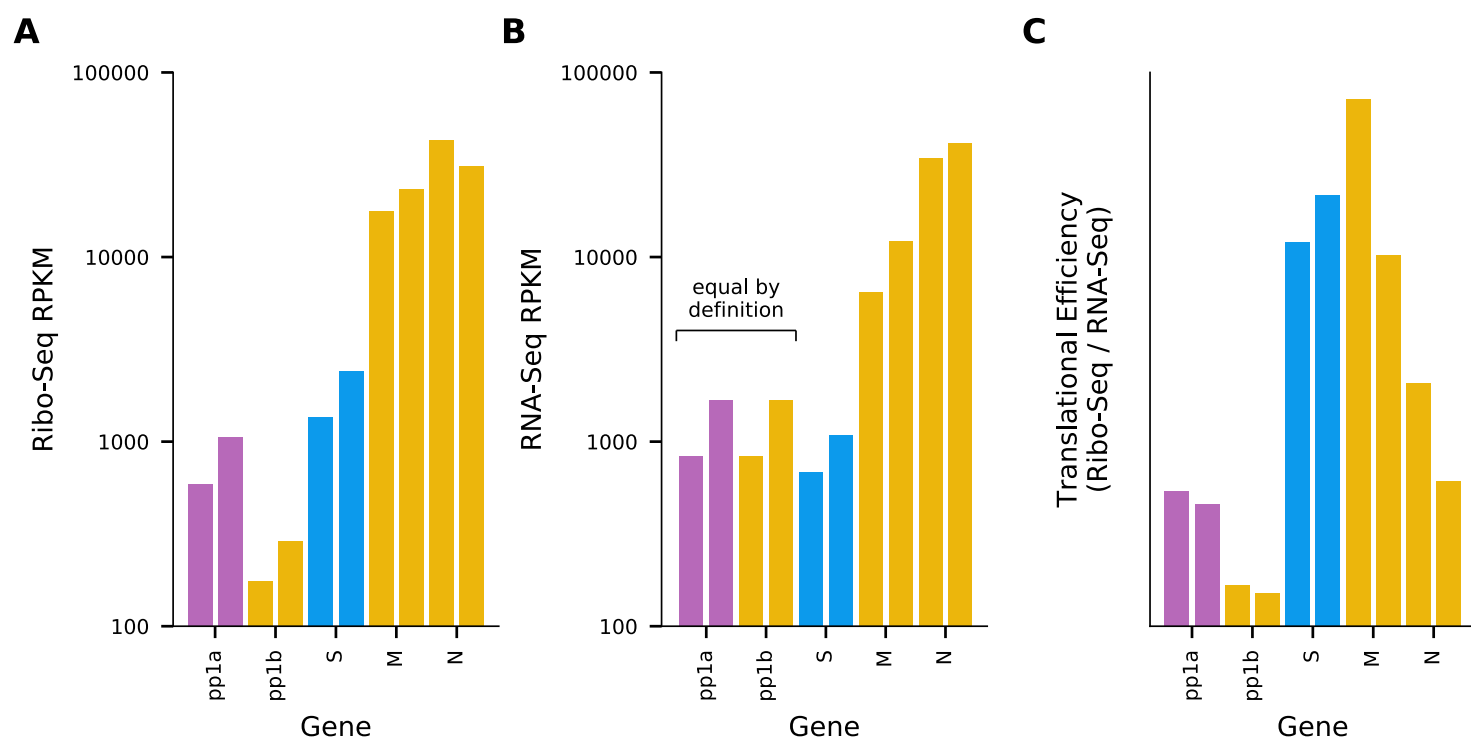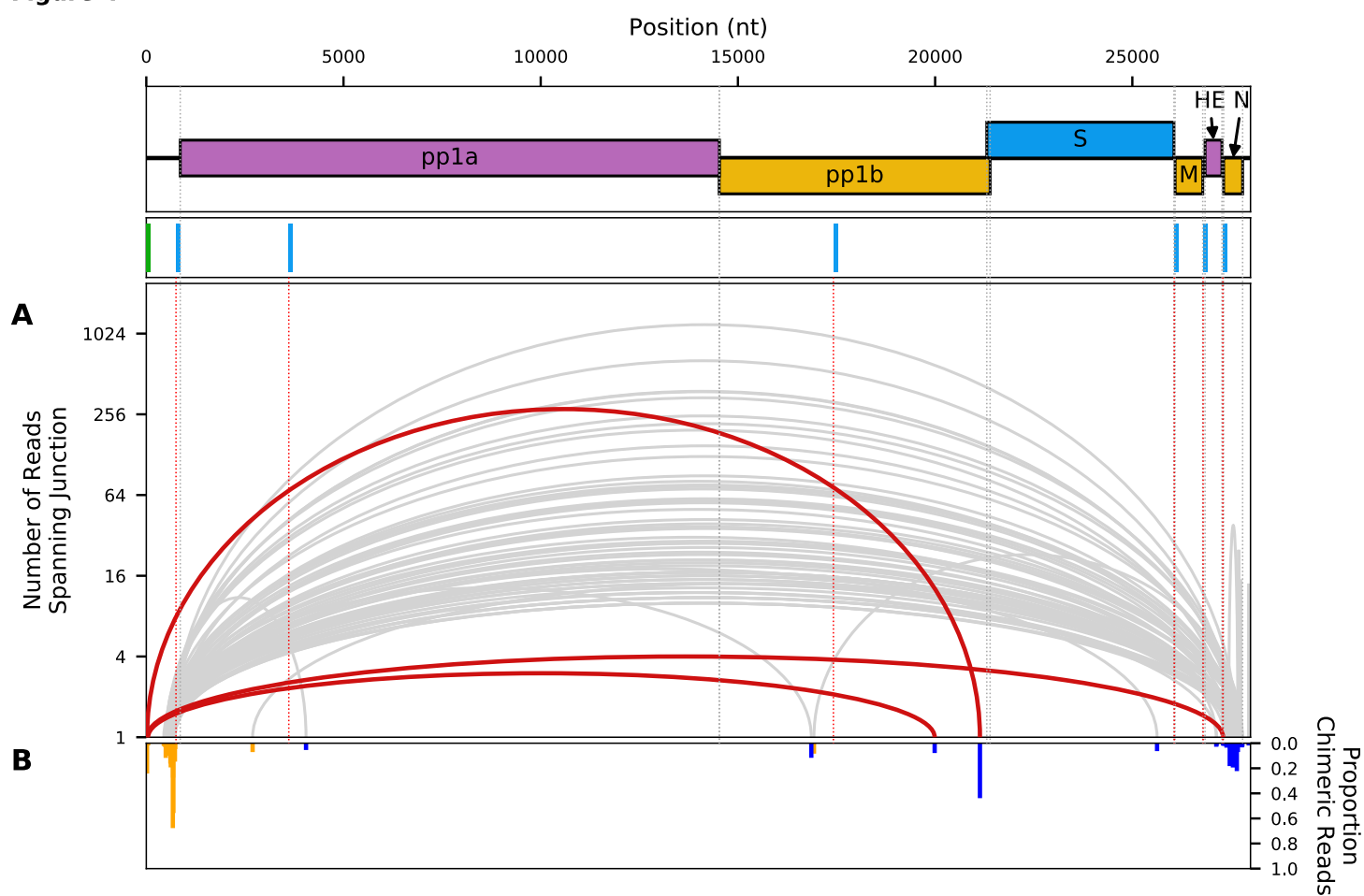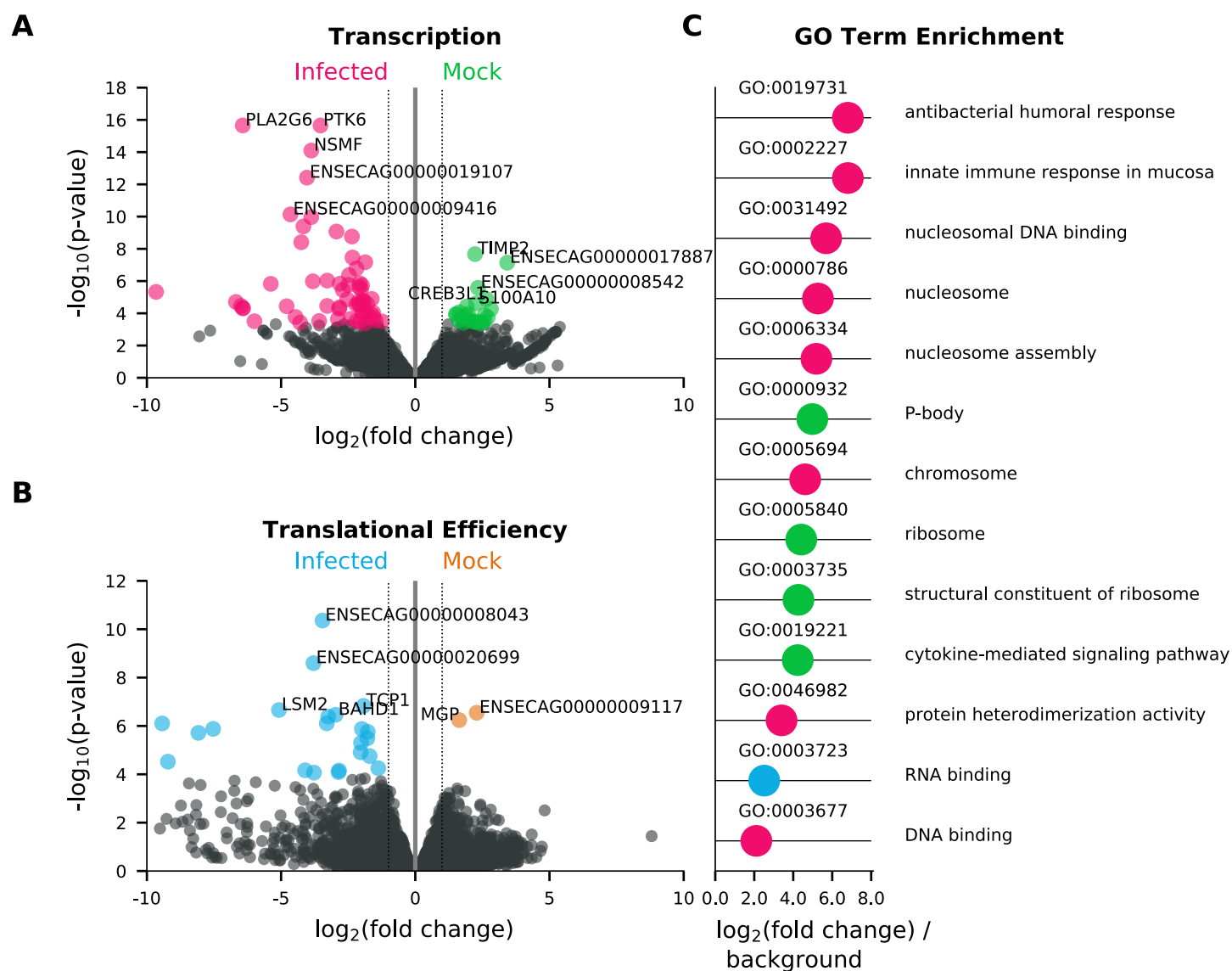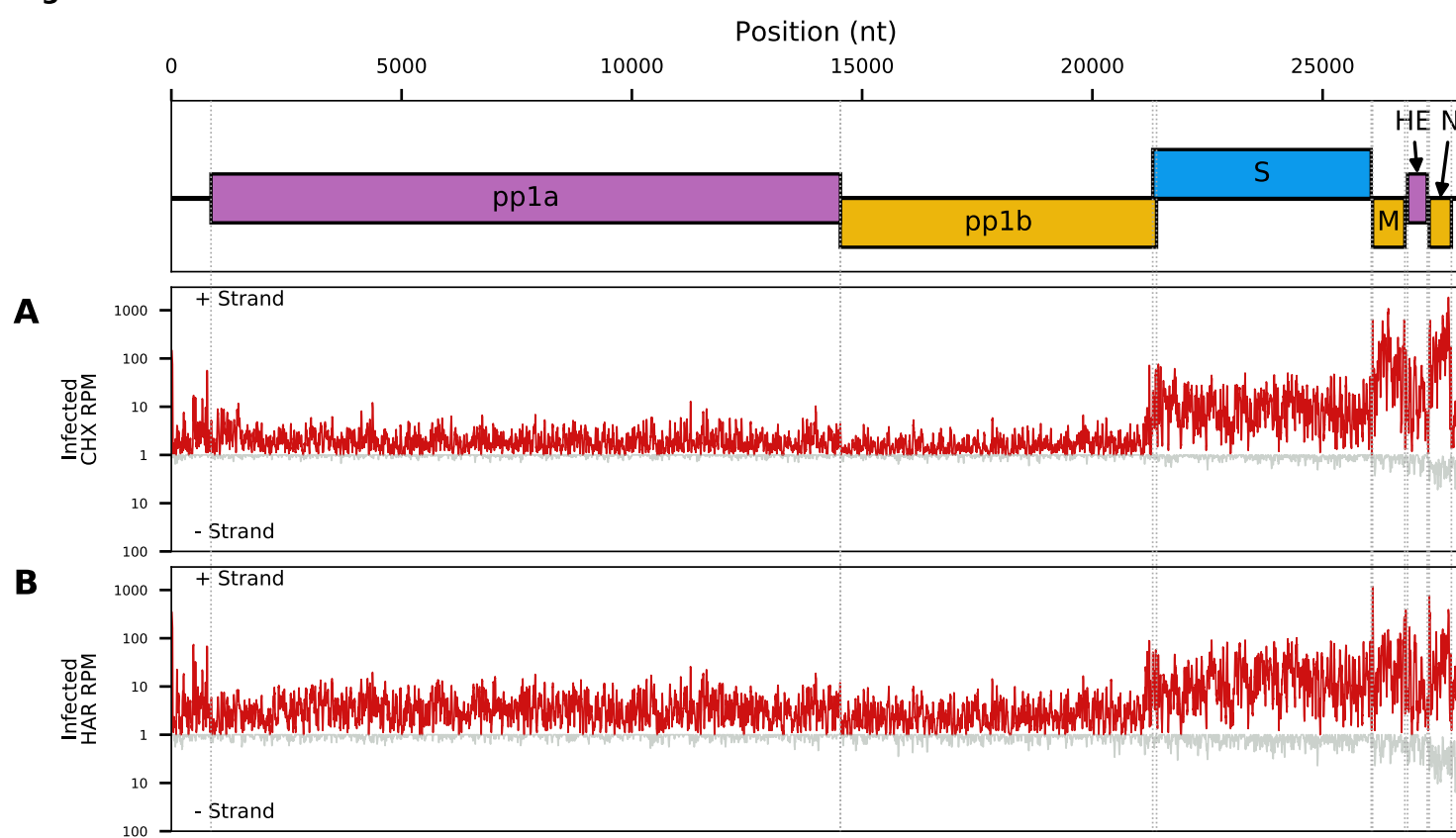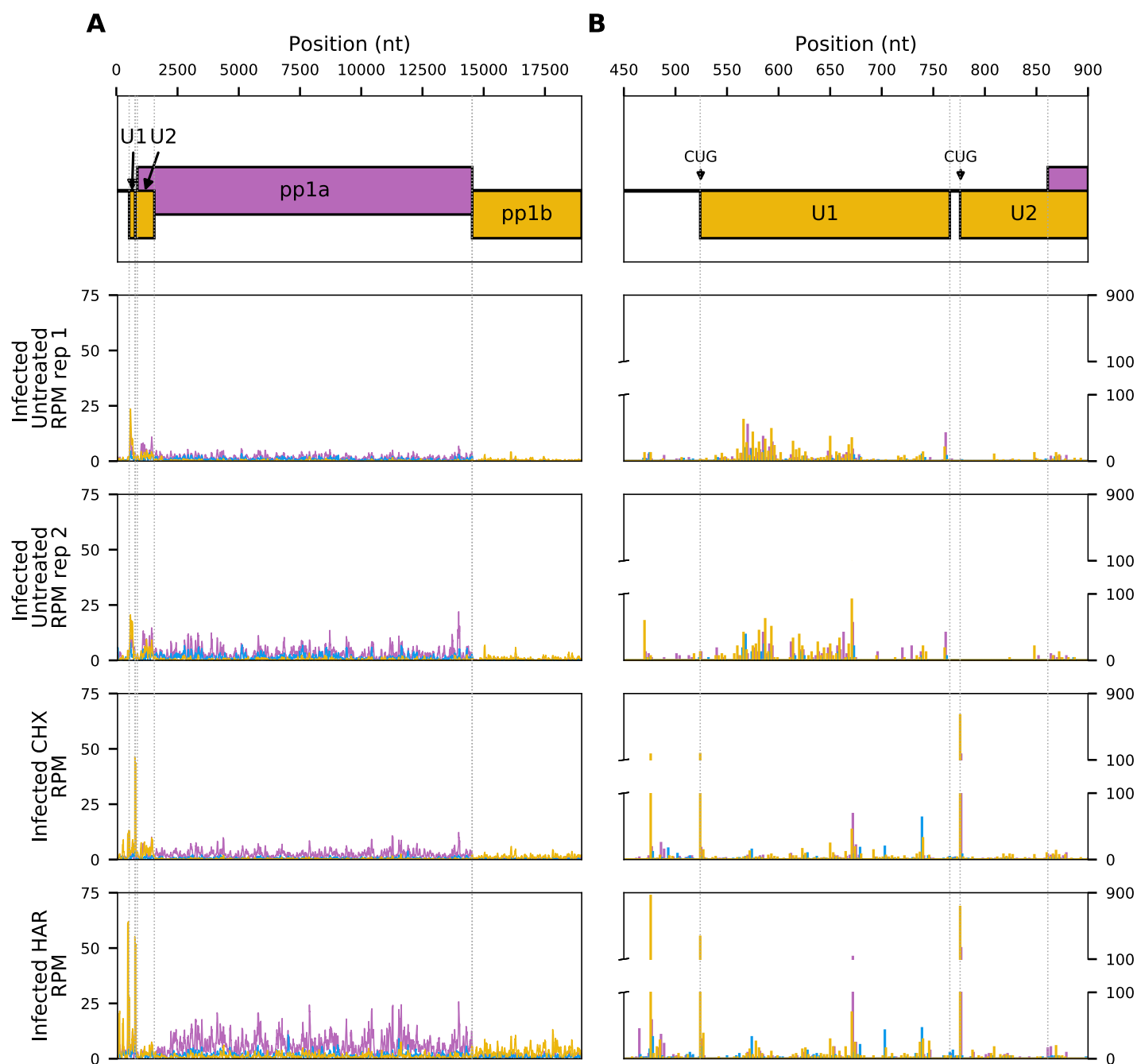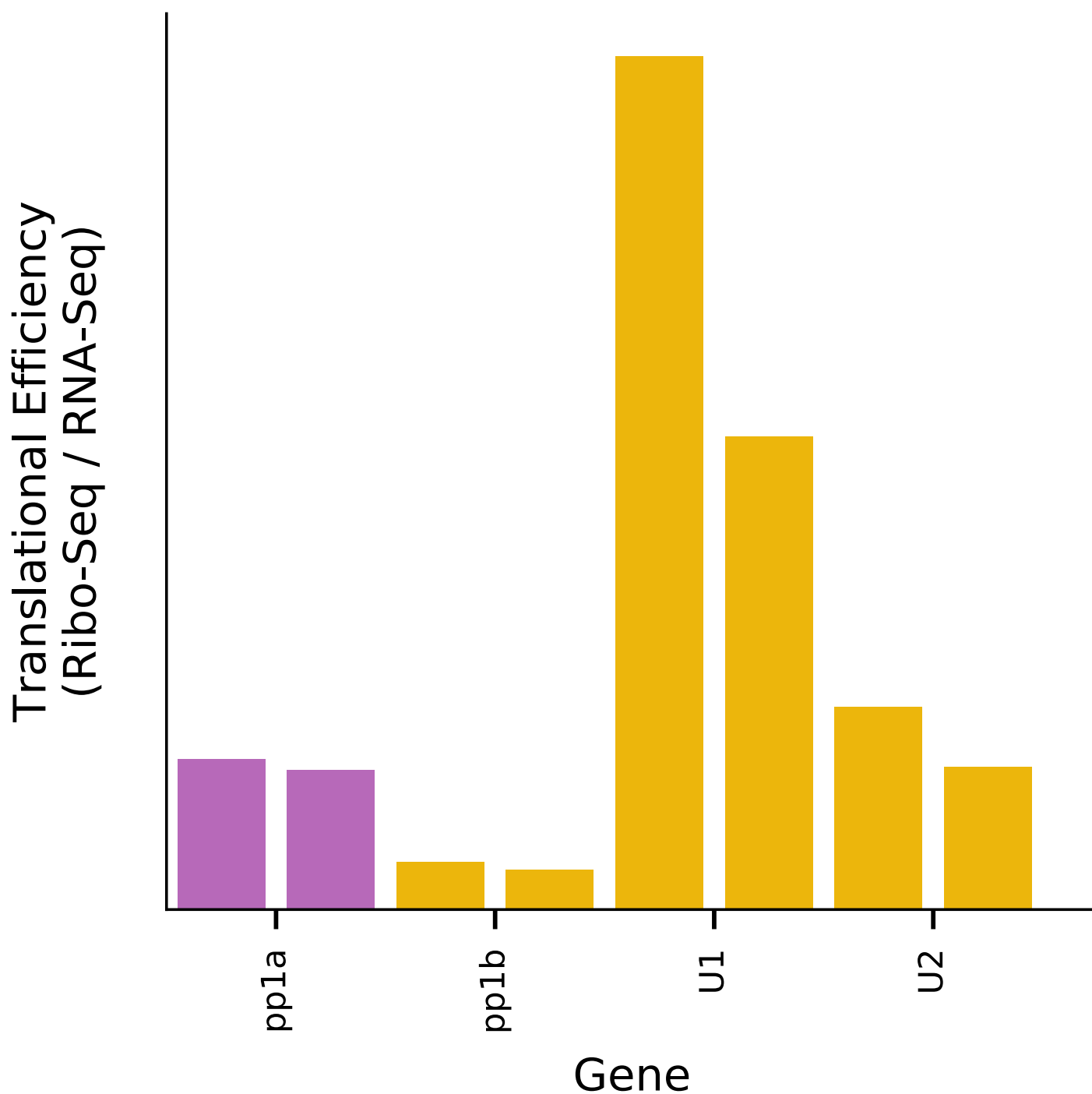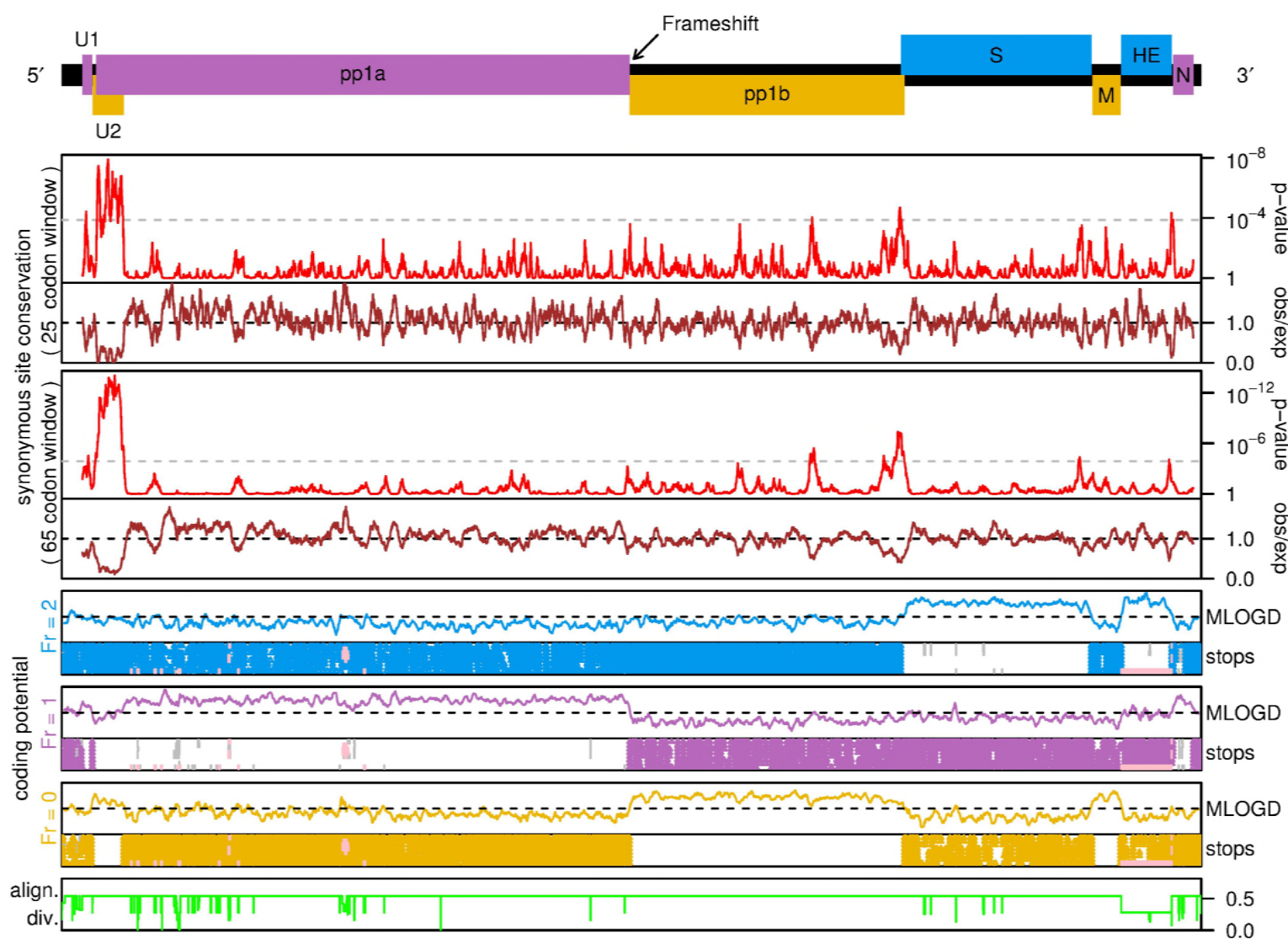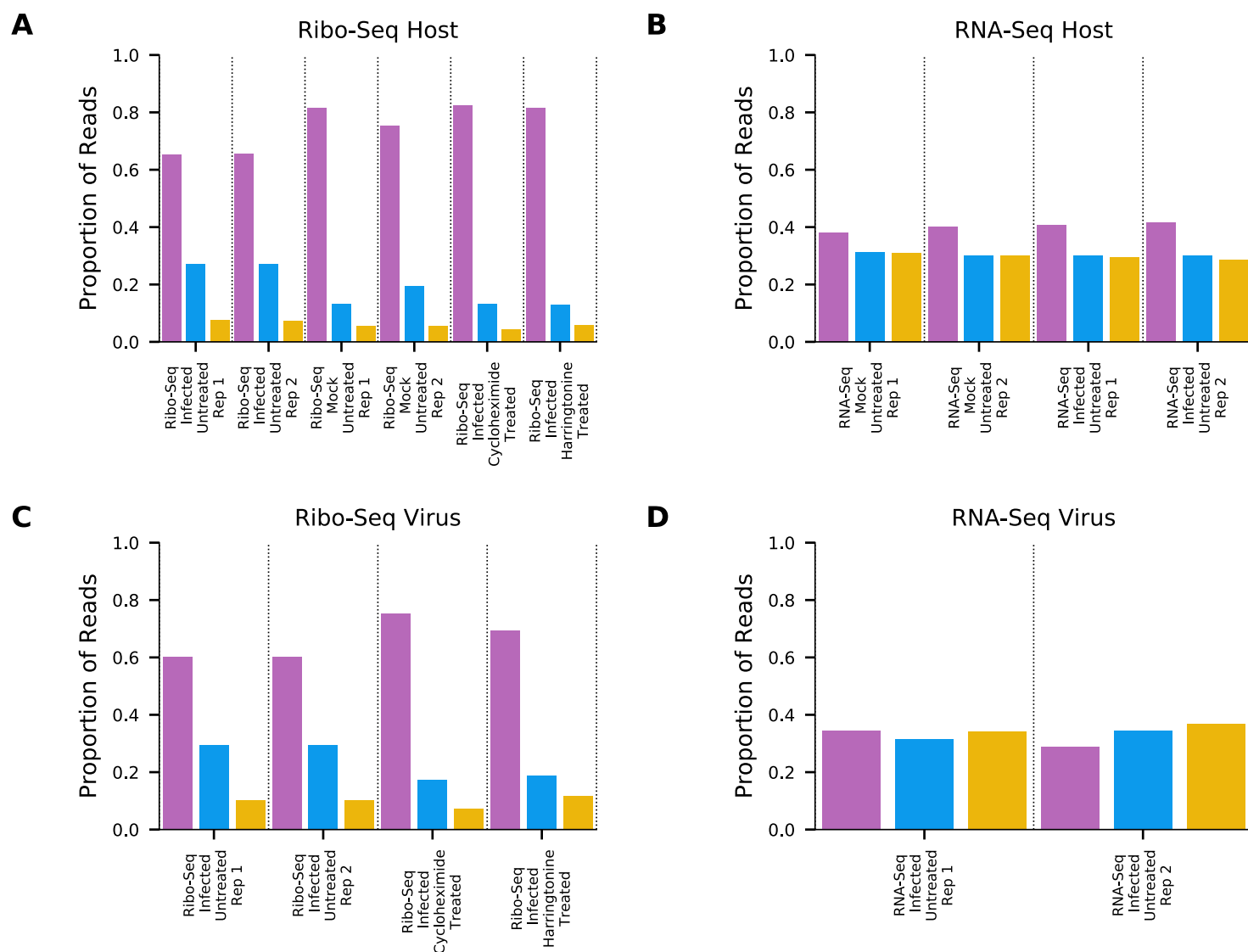
Legend:

| | | | | | | |
|---|---|---|---|---|---|---|
| I | Hydrophobic position | K | Basic position | Y | Tyrosine or Histidine | |
| D | Acidic position | S | Other polar position | P | Proline | |

**Supplementary Figure 6**

**(A)     Leader + TRS + 6 nt**

```
EToV        ACGUAUCUUUAGAAGUUUA
AY427798    ACGUAUCUUUAGUUGAUUU
KR527150    ACGUAUCUUUAGUUGAUUU
LC088094    ACGUAUCUUUAGUUGAUUU
LC088095    ACGUAUCUUUAGUUGAUUU
LT900503    - - - - - - - - - - - - - - - - - - -
JQ0860350   ACGUAUCUUUAGUUGAUUU
KM403390    ACGUAUCUUUAGUUGAUUU
            ************  *  **
```

**(F)     6 nt + HE TRS + 6 nt**

```
EToV        ACUUAUCUUUAGAAGAUGU
AY427798    ACUUAUCUUUAGAAGAUGC
KR527150    ACUUAUCUUUAGAAGAUGC
LC088094    ACGUAUCUUUAGAAGAUGC
LC088095    ACUUAUCUUUAGAAGAUGC
LT900503    ACUUAUCUUUAGAAGAUGU
JQ0860350   ACUUAUCUUUAGAUGAUGU
KM403390    ACUUAUCUUUAGAUGAUGU
            ** ********** ****
```

**(B)     6 nt + U1 TRS + 6 nt**

```
EToV        GUCGUUCUUUAGACGUCUA
AY427798    GCCCAUCUUGUGGGUGUCUA
KR527150    GCCAUUCUUGUGGGUGUCUA
LC088094    GCCCUUCUUGUGGGUGUCUA
LC088095    GCCCUUCUUGUGGGUGUCUA
LT900503    GCCCUUCUUGUGGGUGUCUA
JQ860350    GCCCUUCUUGUGGGUGUCUA
KM403390    GCCCUUCUUGUGGGUGUCUA
            *  *   ****   *  *****
```

**(G)     6 nt + N TRS + 6 nt**

```
EToV        CACUAUCUUUAG-AGAAAGA
AY427798    CACUAUCUUUAG-AGAGAGA
KR527150    CACUAUCUUUAG-AGAGAGA
LC088094    CACUAUCUUUAG-UGAGUGA
LC088095    CACUAUCUUUAGUUGAGUGA
LT900503    CACUAUCUUUAG-UGAGUGA
JQ860350    CACUAUCUUUAG-UGAGUGA
KM403390    CACUAUCUUUAG-UGAGUGA
            ************  **   **
```

**(C)     6 nt + 1a TRS + 6 nt**

```
EToV        GUCGGCCUUUAGAGAAAUU
AY427798    AUUGUCCUAUUGGGAAUUU
KR527150    GUCGUCCUAUUGAGAAUUU
LC088094    GCUGUCCUUUGGAGAAUCU
LC088095    GCUGUCCUUUGGAGAAGCU
LT900503    GCUGCCCUUUAGAGAAGUU
JQ860350    GUUGCCCAUUGGAGAAGUU
KM403390    GUUGUCCAUUAGAGAGUUU
                 *  **   *  *  **    *
```

**(H)     Hairpin**

```
            ((((((.....))))))....
EToV        ACCUCCUCUUCGGAGGUUUUU
AY427798    ACCUCUUCAUCGGAGGUUUUU
KR527150    ACCUCUUCUUCAGAGGUUUUU
LC088094    ACCUCUUCUUCAGAGGUUUUU
LC088095    ACCUCUUCUUCAGAGGUUUUU
LT900503    ACCUCGUCUUCAGAGGUUUUU
JQ860350    ACCUCUUCUUCAGAGGUUUUU
KM403390    ACCUCUUCGUCAGAGGUUUUU
            *****  **  **  *********
```

**(D)     6 nt + 1b TRS + 6 nt**

```
EToV        AUGUAUCUUUAGACUGGAA
AY427798    AUGUGUCUUUGGAUUGGAA
KR527150    AUGUGUCUUUGGAUUGGAA
LC088094    AUAUUUCAUUAGAUUGGAA
LC088095    AUAUUUCAUUGGAUUGGAA
LT900503    ACAUUUCAUUAGACUGGAA
JQ860350    AUAUUUCUUUAGAUUGGAA
KM403390    AUAUUUCAUUAGAUUGGAA
            *   *  **  **  **  *****
```

**(I)**



**(E)     6 nt + M TRS + 6 nt**

```
EToV        CACUUUCUUUAGAAGAAGG
AY427798    CACUAUCUUUAGUUGAAGG
KR527150    CACUAUCUUUAGUUGAAGG
LC088094    CACUAUCUUUAGUUGAAGG
LC088095    CACUAUCUUUAGUUGAAGG
LT900503    CACUAUCUUUAGUUGAAGA
JQ860350    CACUAUCUUUAGUUGAAGA
KM403390    CACUAUCUUUAGUUGAAGA
            **** *******    ****
```