

Characterizing the Transcriptome of Turkey Hemorrhagic Enteritis Virus

3

4 Running Title: Novel Insights into Turkey Hemorrhagic Enteritis Virus Transcriptome

⁵ Abraham Quaye^{1*}, Brett Pickett^{*}, Joel S. Griffitts^{*}, Bradford K. Berges^{*}, Brian D. Poole^{†*}

⁶*Department of Microbiology and Molecular Biology, Brigham Young University

7 1 First-author

⁸ † Corresponding Author

9 Corresponding Author Information

¹⁰ brian_poole@byu.edu

¹¹ Department of Microbiology and Molecular Biology,

¹² 4007 Life Sciences Building (LSB),

¹³ Brigham Young University,

14 Provo, Utah

15

16 **ABSTRACT**

17 Hemorrhagic enteritis (HE) is a disease affecting 6-12-week-old turkeys characterized by immunosuppres-
18 sion (IS) and bloody diarrhea. This disease is caused by *Turkey Hemorrhagic Enteritis Virus* (THEV) of
19 which avirulent strains (THEV-A) that do not cause HE but retain the immunosuppressive ability have been
20 isolated. The THEV-A Virginia Avirulent Strain (VAS) is still used as a live vaccine despite its immuno-
21 suppressive properties. We have performed the first RNA-sequencing experiment characterizing THEV's
22 transcriptome, yielding the most detailed insight into THEV gene expression, to set the stage for further
23 experimentation with specific viral genes that may mediate IS. After infecting a turkey B-cell line (MDTC-
24 RP19) with the VAS vaccine strain, samples in triplicates were collected at 4-, 12-, 24-, and 72-hours
25 post-infection. Total RNA was subsequently extracted, and poly-A-tailed mRNA sequencing done. After
26 trimming the raw sequencing reads with the Trim-galore, reads were mapped to the THEV genome using
27 Hisat2 and transcripts assembled with StringTie. We identified 29 transcripts from our RNA-seq data all
28 of which consisted of novel exons albeit some exons matched the predicted ORFs. The three predicted
29 splice junctions were also corroborated in our data. We performed PCR amplification of THEV cDNA,
30 cloned the PCR products, and Sanger sequencing was used to validate all identified splice junctions. Dur-
31 ing validation, we identified 5 additional transcripts some of which were further validated by 3'RACE data.
32 Thus, the transcriptome of THEV consists of 34 unique transcripts with the coding capacity for all predicted
33 ORFs. However, we found 8 predicted ORFs to be incomplete as either an upstream, in-frame start codon
34 was identified or additional coding exons were found, making the actual expressed versions of these ORFs
35 longer. We also identified 7 novel unpredicted ORFs that could be encoded by some transcripts; albeit it
36 is beyond the scope of this manuscript to investigate whether they are indeed expressed. In keeping with
37 all Adenoviruses, our data shows that all THEV transcripts are spliced, and organized in transcription units
38 under the control of their cognate promoter.

39 **INTRODUCTION**

40 Adenoviruses (AdVs) are non-enveloped icosahedral-shaped DNA viruses, causing infection in virtually all
41 vertebrates. Their double-stranded linear DNA genomes range between 26 and 45kb in size, producing a
42 broad repertoire of transcripts via highly complex alternative splicing patterns (1, 2). The AdV genome is
43 one of the most optimally economized; both the forward and reverse DNA strands harbor protein-coding
44 genes, making it highly gene-dense. There are 16 genes termed “genus-common” that are homologous in
45 all AdVs; these are thought to be inherited from a common ancestor. All other genes are termed “genus-
46 specific”. “Genus-specific” genes tend to be located at the termini of the genome while “genus-common”
47 genes are usually central (1). This pattern is observed in *Adenoviridae*, *Poxviridae*, and *Herpesviridae* (1,
48 3, 4). The family *Adenoviridae* consists of five genera: *Mastadenovirus* (MAdV), *Aviadenovirus*, *Ataden-
49 ovirus*, *Ichtadenovirus*, and *Siadenovirus* (SiAdV) (5, 6). Currently, there are three recognized members
50 of the genus SiAdV: frog adenovirus 1, raptor adenovirus 1, and turkey adenovirus 3 also called turkey
51 hemorrhagic enteritis virus (THEV) (5, 7–10). Members of SiAdV have the smallest genome size (~26 kb)
52 and gene content (~23 genes) of all known AdVs, and many “genus-specific” putative genes of unknown
53 functions have been annotated (see **Figure 1**) (1, 2, 7).

54 Virulent THEV strains (THEV-V) and avirulent strains (THEV-A) of THEV are serologically indistinguishable,
55 infecting turkeys, chickens, and pheasants, with the THEV-V causing different clinical diseases in these
56 birds (2, 11). In turkeys, the THEV-V cause hemorrhagic enteritis (HE), a debilitating acute disease affect-
57 ing predominantly 6-12-week-old turkeys characterized by immunosuppression (IS), weight loss, intestinal
58 lesions leading to bloody diarrhea, splenomegaly, and up to 80% mortality (11–13). HE is the most econom-
59 ically significant disease caused by any strain of THEV (11). While the current vaccine strain (a THEV-A
60 isolated from a pheasant, Virginia Avirulent Strain [VAS]) has proven effective at preventing HE in young
61 turkey pouls, it still retains the immunosuppressive ability. Thus, vaccinated birds are rendered more sus-
62 ceptible to opportunistic infections and death than unvaccinated cohorts leading to substantial economic
63 losses (11, 14–16). To eliminate this immunosuppressive side-effect of the vaccine, a thorough investiga-
64 tion of the culprit viral factors (genes) mediating this phenomenon is essential. However, the transcriptome
65 (splicing and gene expression patterns) of THEV has not been characterized, making the investigation of
66 specific viral genes for possible roles in causing IS impractical. A well-characterized transcriptome of THEV
67 is required to enable experimentation with specific viral genes that may mediate IS.

68 Myriads of studies have elucidated the AdV transcriptome in fine detail (17, 18). However, a large pre-
69 ponderance of studies focus on MAdVs – specifically human AdVs. Thus, most of the current knowledge

70 regarding AdV gene expression and replication is based on MAdV studies, which is generalized for all other
71 AdVs (6, 19). MAdV genes are transcribed in a temporal manner; therefore, genes are categorized into five
72 early transcription units (E1A, E1B, E2, E3, and E4), two intermediate (IM) units (pIX and IVa2), and one
73 major late unit (MLTU or major late promoter [MLP] region), which generates five families of late mRNAs
74 (L1-L5) based on the polyadenylation site. An additional gene (UXP or U exon) is located on the reverse
75 strand. The early genes encode non-structural proteins such as enzymes or host cell modulating proteins,
76 primarily involved in DNA replication or providing the necessary intracellular niche for optimal replication
77 while late genes encode structural proteins that act as capsid proteins, promote virion assembly, and direct
78 genome packaging. The immediate early gene E1A is expressed first, followed by the delayed early
79 genes, E1B, E2, E3 and E4. Then the intermediate early genes, IVa2 and pIX are expressed followed by
80 the late genes (6, 17, 18). Noteworthily, the MLP shows basal transcriptional activity during early infection
81 (before DNA replication), with a comparable efficiency to other early viral promoters, but reaches its max-
82 imal activity during late infection (after DNA replication). However, during early infection the repertoire of
83 late transcripts from the MLP is restricted until late infection (6). MAdV makes an extensive use of alterna-
84 tive RNA splicing to produce a very complex array of mRNAs. All but the pIX mRNA undergo at least one
85 splicing event. For instance, the MLTU produces over 20 distinct splice variants all of which contain three
86 non-coding exons at the 5'-end (collectively known as the tripartite leader, TPL) (17, 18). There is also
87 an alternate 5' three non-coding exons present in varying amounts on a subset of MLTU mRNAs (known
88 as the x-, y- and z-leaders). Lastly, there is the i-leader exon, which is infrequently included between the
89 second and third TPL exons, and codes for the i-leader protein (20). Thus, the MLTU produces a complex
90 repertoire of mRNA with diverse 5' untranslated regions (UTRs) spliced onto different 3' coding exons which
91 are grouped into five different 3'-end classes (L1-L5) based on polyadenylation site. Each transcription unit
92 (TU) contains its own promoter driving the expression of all the array of mRNA transcripts produced via
93 alternative splicing in the unit (6, 17, 18). The promoters are activated at different phases of the infection by
94 proteins from previously activated TUs. Paradoxically, the early-to-late phase transition during infection re-
95 quires the L4 genes, 22K and 33K, which should only be available after the transition. However, a promoter
96 in the L4 region (L4P) that directs the expression of these two proteins independent of the MLP was found,
97 resolving the paradox (6, 17, 21). During translation of AdV mRNA, recent studies strongly suggest the
98 potential usage of secondary start codons; adding to what was already a highly complex system for gene
99 expression (17).

100 High throughput sequencing methods have facilitated the discovery of many novel transcribed regions and
101 splicing isoforms. It is also a very powerful tool to study alternative splicing under different conditions at

102 an unparalleled depth (18, 22). In this paper, a paired-end deep sequencing experiment was performed to
103 characterize for the first time the transcriptome of THEV (VAS vaccine strain) during different phases of the
104 infection, yielding the first THEV splicing map. Our paired-end sequencing allowed for reading **149** bp long
105 high quality (mean Phred Score of 36) sequences from each end of cDNA fragments, which were mapped
106 to the genome of THEV.

107 **RESULTS**

108 **Overview of sequencing data and analysis pipeline outputs**

109 A previous study by Zeinab *et al* showed that almost all THEV transcripts were detectable beginning at
110 4 hours (23). Therefore, infected MDTC-RP19 cells were harvested at 4-, 12-, 24-, and 72-hours post-
111 infection(h.p.i) to ensure an amply wide time window to sample all transcripts. Our paired-end RNA se-
112 quencing (RNA-seq) experiment yielded an average of **107.1** million total reads of **149bp** in length per
113 time-point, which were simultaneously mapped to both the virus (THEV) and host (*Meleagris gallopavo*)
114 genomes using the Hisat2 (24) alignment program. A total of **18.1** million reads from all time-points mapped
115 to the virus genome; this provided good coverage/depth, leaving no regions unmapped. The mapped reads
116 to the virus genome increased substantially from **432** reads at 4 h.p.i to **16.9** million reads at 72 h.p.i (**Table**
117 **1**, **Figure 2a**). From the mapped reads, we identified a total of **2,457** unique THEV splice junctions from all
118 time-points, with splice junctions from the later time-points being supported by significantly more sequence
119 reads than earlier time-points. For example all the **13** unique junctions at 4 h.p.i had less than 10 reads
120 supporting each one, averaging a mere **2.8** reads/junction. Conversely, the **2374** unique junctions at 72 h.p.i
121 averaged **898.4** reads/junction, some junctions having coverage as high as **322,677** reads. The substantial
122 increases in splice junction and mapping reads to the THEV genome over time denotes an active infection,
123 and correlates with our quantitative PCR (qPCR) assay quantifying the total number of viral genome copies
124 over time (**Figure 2b**).

125 Using StringTie (24), an assembler of RNA-seq alignments into potential transcripts, the mapped reads for
126 each time point were assembled into transcripts using the genomic location of the predicted THEV ORFs as
127 a guide. In the consolidated transcriptome, a composite of all unredudant transcripts from all time points,
128 we counted a total of **30** novel transcripts. Although some exons in some transcripts match the predicted
129 ORFs exactly, most of our identified exons are longer, spanning multiple predicted ORFs (**Figure 3**).

130 We validated the splice junctions in all transcripts by PCR amplification of viral cDNA, cloning, and Sanger
131 sequencing (**Supplementary PCR methods**). During validation, we identified 5 additional transcripts some
132 of which were further validated by 3' Rapid Amplification of cDNA Ends (3'RACE) data. The complete
133 list of unique splice junctions mapped to THEV's genome has been submitted to the National Center for
134 Biotechnology Information Gene Expression Omnibus (<http://www.ncbi.nlm.nih.gov/geo>) under **accession**
135 **no. XXXXXX**.

136 **Changes in THEV splicing profile over time**

137 AdV gene expression occurs under exquisite temporal control with each promoter typically producing one or

138 few pre-mRNAs that undergo alternative splicing to yield the manifold repertoire of complex transcripts. To
139 evaluate the activity of each promoter over time, *StringTie* and *Ballgown* (a program for statistical analysis
140 of assembled transcriptomes) (25) were used to estimate the normalized expression levels of all transcripts
141 for each time point in Fragments Per Kilobase of transcript per Million mapped reads (FPKM) units. Very few
142 unique splice junctions, reads, and transcripts were counted at 4 h.p.i; hence, this time point was excluded
143 in this analysis.

144 Individually, TRXPT_21 (DBP) – from the E2 region – was the most significantly expressed at 12 h.p.i,
145 comprising about **33.58%** of the total transcripts. Transcripts in the E3 and E4 regions also contributed sig-
146 nificant proportions, and noticeably, some MLP region transcripts. The later time points were dominated by
147 the MLP region transcripts – TRXPT_10 and TRXPT_14 were the most abundantly expressed at 24 and 72
148 h.p.i, respectively, as expected (**Figure 4a**). When we performed analysis of the FPKM values of transcripts
149 per region we found a similar pattern: the E2 region was the most abundantly expressed at 12 h.p.i, after
150 which the MLP region assumes predominance (**Figure 4b**). Secondly, we estimated relative abundances
151 of all splice junctions at each time point using the raw reads. For individual junctions, we counted as signif-
152 icantly expressed only junctions with coverage of at least 1% of the total splice junction reads at the given
153 time point. At 12 h.p.i, **18** junctions meet the 1% threshold, and were comprised of predominantly early
154 region (E1, E2, E3, and E4) junctions, albeit the MLTU was the single most preponderant region overall,
155 constituting **38.8%** of all the junction reads (**Table 2a** and **Supplementary Table 1a**). The top most abun-
156 dant junctions at 12 h.p.i remained the most significantly expressed at 24 h.p.i also. However, here, the
157 MLP-derived junctions were unsurprisingly even more preponderant overall, accounting for **45.7%** of all the
158 junction reads counted (**Table 2b** and **Supplementary Table 1b**). At 72 h.p.i, the trend of increased activity
159 of the MLP continued as expected; at this time, the MLP region junctions were not only the most abundant
160 overall – accounting for **67.4%** of all junction reads, – but also contained the most significantly expressed
161 individual junctions (**Table 2c**, **Supplementary Table 1c** and **Figure 4c**). When we limited this analysis to
162 only junctions in the final transcriptome, the relative abundances of the junctions for each region over time
163 was generally similar to the pattern seen with all the junctions included (**Figure 4d**).

164 We also analyzed splice donor and acceptor site nucleotide usage over time to investigate any peculiarities
165 that THEV may show, generally or over the course of the infection. We found that most splice donor-acceptor
166 sequences were unsurprisingly the canonical GU-AG nucleotides.

167 **Early Region 1 (E1) transcripts**

168 This region in MAdVs is the first transcribed after successful entry of the viral DNA into the host cell nucleus,
169 albeit at low levels (18). The host transcription machinery solely mediates the transcription of this region.

170 After their translation, the E1 proteins in concert with a myriad of host transcription factors activate the other
171 viral promoters (6).

172 Only two ORFs (ORF1 [sialidase] and Hyd) are predicted in this region; however, we discovered **four** novel
173 transcripts in this region, which collectively contain **3** unique splice junctions (**Figure 5**). Most of the ORFs
174 of the novel transcripts are distinct from the predicted ORFs, but they all have the coding potential (CP)
175 for the predicted Hyd protein as the 3'-most coding sequence (CDS) if secondary start codon usage is
176 considered as reported for other AdVs (17, 18). The 5'-most CDS of TRXPT_1 is multi-exonic, encoding
177 a novel 17.9 kilodalton (kDa), 160 residue [amino acids (aa)] protein (ORF9). From its 5'-most start codon
178 (SSC), TRXPT_2 encodes the largest protein in this region – a 64.3 kDa, 580 aa protein (ORF10) with the
179 same SSC as ORF9 (position 211bp). ORF10 spans almost the entire predicted ORF1 and Hyd, coming
180 short in two regards: it is spliced from 1655bp to 1964bp (ORF1's C-terminus, including the stop codon), and
181 it's stop codon (STC; position 2312) is 13 bp short of Hyd's STC. However, it has an SSC 102 bp upstream
182 and in-frame with ORF1's predicted SSC. Thus, ORF10 shares substantial protein sequence similarity with
183 ORF1 but not with Hyd, as the SSC of Hyd is not in-frame. Without its splice site removing the ORF1 STC,
184 TRXPT_2 would encode a longer variant of ORF1, starting from an upstream SSC. TRXPT_3 is almost
185 identical to TRXPT_1, except for the lack of TRXPT_1's second exon. Our RNA-seq data shows that all E1
186 transcripts share the same transcription termination site (TTS; at position 2325bp). However, TRXPT_3 and
187 TRXPT_4 seem to have transcription start sites (TSS) downstream of the TSS of TRXPT_1 and TRXPT_2
188 (E1 TSS; position: 54bp). Given that studies in MAdVs show that E1 mRNAs share not only a common
189 TTS but also the TSS, and only differ from each other regarding the internal splicing (18), it is likely that
190 TRXPT_3 and TRXPT_4 are incomplete, and their actual TSS just like the TTS are identical for all E1
191 transcripts. Regardless of the TSS considered for TRXPT_3, the coding potential (CP) remains unaffected.
192 Its 5'-most CDS, beginning at 1965bp and sharing the same STC as ORF9, produces a 13.1 kDa, 115
193 residue protein (ORF4). ORF4 was predicted in an earlier study (26) but was excluded in later studies (1,
194 12); however, our data suggests it is a bona fide ORF. Unlike TRXPT_3, the CP of TRXPT_4 is affected by
195 the TSS considered; if we consider its unmodified TSS, then its CP is the same as TRXPT_3 (ORF4 as the
196 first CDS and Hyd as second CDS using the secondary SSC). However, if we assume that TRXPT_4 uses
197 the E1 TSS, then the 5'-most CDS is a distinct, novel, multi-exonic 15.9 kDa, 143 aa protein (ORF11) with
198 the same SSC as ORF9 and ORF10 but with a unique STC. The splice junctions of all transcripts in this
199 region (except the junction for TRXPT_4) were validated by cloning of viral cDNA and Sanger sequencing
200 (**Supplementary PCR methods**).

201 During the validation of TRXPT_2, ORF1 was present on the agarose gel (an unspliced band size) and

202 Sanger sequencing results as a bona fide transcript (**Supplementary PCR methods**). This was corroborated by our 3'RACE experiment, which showed a transcript (TRXPT_2B) spanning the entire ORF1 and
203 Hyd ORFs without any splicing, with a poly-A tail immediately after the E1 TTS. The 5'-most CDS of this
204 transcript (TRXPT_2B) would encode ORF1. However, TRXPT_2B has an upstream and in-frame SSC
205 to the predicted SSC of ORF1, suggesting that the predicted ORF1 CDS is truncated – the actual ORF1
206 (eORF1) that is expressed shares the same SSC as ORF10, but has a unique STC.

208 **Early Region 2 (E2) and Intermediate Region (IM) transcripts**

209 The E2 TU expressed on the anti-sense strand is subdivided into E2A and E2B and encodes three classical
210 AdV proteins – pTP and Ad-pol (E2B proteins), and DBP (E2A protein) – essential for genome replication
211 (17, 18). Unlike MAdV where two promoters (E2-early and E2-late) are known (17), we discovered only a
212 single TSS (E2 TSS; 18,751bp) from which both E2A and E2B transcription is initiated. However, similar
213 to MAdVs, E2A and E2B transcripts have distinct TTSs, and the E2B transcripts share the TTS of the IVa2
214 transcript of the IM region (17, 18) (**Figure 6**).

215 The E2A ORF, DBP is one of three THEV ORFs predicted to be spliced from two exons. The correspond-
216 ing transcript (TRXPT_21) found in our data matches this predicted splice junction precisely but with a
217 non-coding additional exon at the 5'-end (E2-5'UTR) at position 18,684-18,751 bp. Thus, TRXPT_21 is
218 a three-exon transcript encoding DBP (380 residues, 43.3 kDa) precisely as predicted. This transcript
219 (TRXPT_21) was also corroborated in a 3'RACE experiment. Additionally, from the 3'RACE, a splice vari-
220 ant of TRXPT_21 which retains the second intron leading to a 2-exon transcript was found. This transcript
221 (TRXPT_21B), albeit longer due to retaining the second intron and possessing a short 3' UTR, encodes a
222 truncated isoform of DBP (tDBP) because the SSC utilized by TRXPT_21, is followed shortly by STCs in the
223 retained intron. The SSC 173 bp downstream of DBP's SSC yields tDBP (a 346 residue, 39.3 kDa product),
224 which is in-frame of DBP but entirely contained in the second exon. TRXPT_21 and TRXPT_21B share a
225 common TTS but TRXPT_21B as seen in our 3'-RACE data, extends 39 bp into an adenine-thymine (A-T)
226 rich sequence before the poly-A tail sequence occur, suggesting this position (16,934bp) as the bona fide
227 E2A TTS (**Figure 6**).

228 The E2B region transcripts also start with the E2-5'UTR but extend thousands of base pairs downstream to
229 reach the TTS at 2334bp in the IM region, which is immediately followed by an A-T rich sequence (position
230 2323-2339bp) where polyadenylation probably occurs. Interestingly, the TTS of the E1 region (position
231 2,325bp) on the sense strand is also in the immediate vicinity of this A-T rich sequence, which is almost
232 palindromic; hence it likely serves as the polyadenylation signal for both E1 and E2B/IM transcripts. The
233 E2B transcripts, TRXPT_6 and TRXPT_7 are almost identical except for an extra splice junction at the 3'-

234 end of TRXPT_6, making TRXPT_6 a five-exon transcript and TRXPT_7, four exons (**Figure 6**). TRXPT_7
235 has the CP for both classical proteins (pTP and Ad-pol) encoded in this region, of which the pTP ORF is
236 predicted to be spliced from two exons just like in all other AdVs. The predicted splice junction of pTP
237 is corroborated by our data; however, the full transcript is markedly longer than the predicted ORF: there
238 are two novel non-coding 5' exons, the third exon (containing the SSC of pTP) is significantly longer than
239 predicted, and the last exon containing the bulk of the CDS is more than triple the predicted size of pTP.
240 The first two exons are 5'-UTRs because the SSC here is immediately followed by STCs; hence, the 5'-
241 most SSC (position 10,995bp) of the third exon which matches the predicted SSC of pTP is utilized. The
242 encoded product is identical to the predicted pTP protein (597 residues; 70.5 kDa). If secondary SSC
243 (secSSC) usage is considered, with SSC at 6768bp and STC at 3430bp, the encoded product is identical
244 to the predicted Ad-pol (polymerase) ORF (1112 residues; 129.2 kDa). TRXPT_6 differs from TRXPT_7 by
245 containing an extra splice site at 3447-3515bp. However, the CP remains similar to that of TRXPT_7 except
246 the Ad-pol encoded from the secSSC is a truncated isoform with a new STC resulting from the splice site.

247 While both TRXPT_6 and TRXPT_7 have the CP for Ad-pol with secSSC usage, in all AdVs studied, the two
248 proteins (pTP and Ad-pol) are encoded by separate mRNAs with identical first three 5' exons and TTS, but
249 the splice junction to the terminal exons are different. We checked for a longer splice junction between the
250 third and fourth (terminal) exons of TRXPT_7 with our junction validation method (targeted PCR, cloning,
251 and Sanger sequencing) and discovered a unique splice junction (10,981-7062bp) not found in our RNA-
252 seq data. If initiated from the E2 TSS and terminated at the E2 TTS, this transcript(TRXPT_31) would
253 encode Ad-pol exactly as predicted as its 5'-most CDS (**Figure 6**).

254 Our RNA-seq data also showed a novel short transcript (TRXPT_15) entirely nested within the terminal
255 exon of TRXPT_7 but with a unique splice site. This transcript is an incomplete construction from the
256 mapped reads as it contains a truncated CDS. However, we validated this splice junction to be genuine
257 (**Supplementary PCR methods**).

258 The IM region is a single-transcript TU, encoding a single classical protein, IVa2. The promoter expressing
259 this single transcript (TRXPT_5) is embedded in E2B region and shares a TTS with E2B transcripts (17,
260 18). TRXPT_5 is a two-exon transcript spliced exactly as the last splice junction of TRXPT_6. The first
261 exon is a UTR, except the last 2 nucleotides, which connect with the first nucleotide of the second exon to
262 form the 5'-most SSC. This first SSC is 4 codons upstream and in-frame of the predicted IVa2 SSC. Except
263 for the four extra N-terminus residues, the entire protein sequence is identical to the predicted IVa2.

264 **Early Region 3 (E3) transcripts.**

265 The E3 region is wholly contained in the MLTU and encodes proteins involved in modulating and evading

266 the host immune defenses. In MAdVs, this region contains seven ORFs expressed from several transcripts
267 which share the same TSS (from the E3 promoter) but have different TTSs (6, 17, 18). However, some
268 E3 transcripts use the TSS of the MLP. Due to sharing the same TSS, in MAdVs, secSSC usage is heavily
269 relied on for gene expression in this region except for 12.5K and transcripts using the MLP's TSS, as utilizing
270 only the first SSC cannot produce all the other transcripts in this TU (17).

271 In THEV, only one ORF (E3) was predicted in this region. However, as the E3 TU is nested in the MLTU,
272 transcripts from the L4P (100K, 22K, 33K, and pVIII) not only overlap the E3 region transcripts entirely as
273 seen in our RNA-seq results, but also have their TSS and TTS in practically the same locations (**Figure 7**).
274 Therefore, we have categorized these two groups together as E3 transcripts.

275 We identified seven novel transcripts here (**TRXPT_22, TRXPT_23, TRXPT_24, TRXPT_25, TRXPT_26,**
276 **TRXPT_27, TRXPT_29**) from our RNA-seq data, all originating from two distinct TSSs – one corresponding
277 to the TSS of the L4P (position 18,230bp) and the other at 18,727bp corresponding the E3 promoter (E3P).
278 These E3 transcripts collectively have the CP for several predicted THEV ORFs: 100K, 22K, 33K, pVIII, E3,
279 Fiber (IV), and ORF7 belonging to the MLTU. But some CDSs are nonidentical due to unpredicted splicing
280 or the use of an in-frame upstream SSC. For instance, 33K is one of the few THEV ORFs predicted to be
281 spliced from two exons; however, we discovered a significantly longer four-exon ORF (e33K) on TRXPT_24
282 that contains it almost entirely. The first two exons of e33K were not predicted but the last two match
283 the predicted exons and the CDS is in-frame, but the first 20bp of the predicted 33K (including the SSC at
284 20,142bp) is spliced out as part of the second intron of TRXPT_24. Thus, the bona fide 33K (e33K) is a 19.8
285 kDa, 171 residue protein spanning four exons instead of the predicted 120 aa protein. TRXPT_24 also has
286 the CP for the ORFs, pVIII and a longer variant of E3 (eE3; starting from an in-frame upstream SSC) if we
287 consider downstream SSC usage. TRXPT_29 is the shortest transcript in this TU. It is a two-exon transcript,
288 both exons comprising the CDS. The product of TRXPT_29 is a novel 73 residue protein (8.3KI) sharing the
289 SSC of e33K but with a unique STC. TRXPT_23 being spliced identically as TRXPT_29 also encodes 8.3KI
290 from its first SSC. Similarly, TRXPT_22 also encodes a 73 aa novel protein (8.3KII) from its first SSC that
291 shares over 80% similarity with 8.3KI, but it differs from 8.3KI at the C-terminus. Considering downstream
292 SSC usage, both TRXPT_22 and TRXPT_23 can encode pVIII and eE3 in that order, but TRXPT_23 being
293 longer, has the CP for the Fiber ORF also. As the splice junctions of TRXPT_22, TRXPT_23, TRXPT_24,
294 and TRXPT_29 essentially share the same genomic space, their validation was done with a single primer
295 pair and they were differentiated from each other by cloning and Sanger sequencing (**Supplementary PCR**
296 **methods**).

297 In addition to corroborating the splice junctions for the aforementioned transcripts, the Sanger sequencing

298 results also showed another splice variant undetected in our RNA-seq transcriptome. This was a three-exon
299 transcript (TRXPT_30) with its first and last exons spliced identically as TRXPT_23, but which also has the
300 second exon of TRXPT_24 (**Figure 7**). The first CDS on TRXPT_30 spans all three exons, producing a
301 novel 140 residue, 15.7kDa protein (e22K). Interestingly, the last 81 C-terminus residues of e22K are iden-
302 tical to 22K (89 residues), which is a single-exon ORF predicted to use the same SSC as 33K (20,142bp).
303 Just as seen for 33K, all the transcripts in this region exclude the first 20bp of 22K (including the SSC) as
304 part of their introns; therefore, the first 7 residues of 22K are lacking in e22K due to splicing. Hence, we
305 may consider e22K as a long variant of the predicted 22K ORF. Albeit the TSS and TTS of TRXPT_30 was
306 not seen, we presume that they are similar to TRXPT_23, in which case it would also have the downstream
307 CP of TRXPT_23. TRXPT_25 is the largest transcript in the TU. It also utilizes the L4P TSS but has a
308 distinct TTS. It is a two-exon transcript, encoding a novel protein (t100K; 543 residues), which is a shorter
309 isoform of the predicted 100K ORF. Considering secSSC usage on this transcript yields the predicted 22K
310 ORF precisely. It also has the CP for pVIII and eE3 in that order. Furthermore, during the validation of
311 TRXPT_25's splice junction using primers that span its junction (18350-18717bp), we noticed a DNA band
312 that corresponds to the full unspliced sequence (**Supplementary PCR methods**). As TRXPT_25 only falls
313 short of encoding the complete predicted 100K protein due to its splice junction, this band (which we cloned
314 and validated by Sanger sequencing) suggests that the predicted 100K is indeed expressed. This tran-
315 script (TRXPT_25B) although not seen in full, likely shares the same TSS and TTS as TRXPT_25. Lastly,
316 TRXPT_26 and TRXPT_27 both originate from the E3 TSS but have distinct TTSs. TRXPT_26 is a three-
317 exon transcript but the first two are UTRs. It encodes pVIII as the 5'-most ORF and has the CP for eE3 and
318 Fiber in that order. TRXPT_27 on the other hand, is only a two-exon transcript but similar to TRXPT_26,
319 only the terminal exon contains the CDSs. It encodes Fiber as the 5'-most ORF, and ORF7 downstream
320 with secSSC usage. TRXPT_13, which is an L4 transcript that uses the MLP TSS is discussed under the
321 MLTU transcripts.

322 **Early Region 4 (E4) transcripts**

323 This TU is found at the tail-end (3'-end) of the genome and expressed from the anti-sense strand. Based
324 on nucleotide position, ORF7 and ORF8 were predicted in this region (1); however, as ORF7 is neither on
325 the same strand as ORF8 nor transcribed from a promoter in the E4 region, only ORF8 can legitimately
326 be classified as a transcript in this TU. This is corroborated by our RNA-seq data, as only one transcript
327 was identified in this region on the anti-sense strand (**Figure 8**). The transcript (TRXPT_28) spans 25192-
328 26247bp and is spliced at 25701-26055bp, making a two-exon transcript. The second exon fully matches
329 the predicted ORF8 with 12 extra base pairs at the 3'-end. However, there is a SSC in the first exon at

330 position 26246bp (192bp upstream of the predicted SSC). The encoded protein from this SSC is in-frame
331 with the predicted SSC found in the second exon; hence, we consider this protein (eORF8 – 26.4 kDa, 229
332 aa), a longer isoform of the predicted ORF8, the genuinely expressed ORF with an identical C-terminus to
333 the predicted ORF8 protein.

334 **Major Late Transcription Unit (MLTU) or MLP Region transcripts**

335 The MLTU transcripts dominate the late phase (i.e., after DNA replication) of the AdV infectious cycle.
336 The MLP produces all late mRNAs by alternative splicing and alternative polyadenylation of a primary
337 transcript, grouped into five transcript classes (L1-L5). Most of THEV's coding capacity falls within this
338 TU. Specifically, about 13 out of the 23 predicted ORFs were assigned to this TU, some of which we have
339 categorized under the E3 TU instead. Our RNA-seq data revealed 12 transcripts (**TRXPT_8, TRXPT_9,**
340 **TRXPT_10, TRXPT_11, TRXPT_12, TRXPT_13, TRXPT_14, TRXPT_16, TRXPT_17, TRXPT_18,**
341 **TRXPT_19, TRXPT_20**) in this TU, the majority of which have the 5' untranslated TPL leader sequence as
342 seen in all AdVs. For three transcripts (**TRXPT_16, TRXPT_17, TRXPT_18**), a different leader sequence
343 (sTPL) is used, which differs from the TPL in only one regard: the first TPL exon is substituted for a different
344 first exon, found between the first and second TPL exons. Also, TRXPT_20 seems to include only the third
345 TPL exon (**Figure 9**).

346 We identified five TTSSs (10,549bp, 12,709bp, 16,870bp, 17,891bp, 20,865bp) in this TU, corresponding
347 to the five late mRNA classes (L1-L5), respectively, as found in all AdVs. L1 mRNAs include TRXPT_8,
348 which comprises the TPL (non-coding) and the CDS-containing terminal exon. This transcript encodes
349 the 52K ORF exactly as predicted with the SSC beginning from the first nucleotide of the terminal exon.
350 L2 mRNAs include TRXPT_16, TRXPT_17, and TRXPT_18, all of which consist of the sTPL (also non-
351 coding) followed by their respective terminal exons. TRXPT_16 encodes pIIIa exactly as predicted as the
352 5'-most ORF, and also has the CP for the ORFs, III and pVII in that order. TRXPT_17 encodes the ORF, III
353 (penton), and TRXPT_18 encodes the ORF pVII exactly as predicted. The L3 mRNAs include TRXPT_14
354 and TRXPT_20, of which TRXPT_14 utilizes the full TPL whereas TRXPT_20 uses only the third TPL
355 exon (TPL3). Both transcripts have the CP for the ORF, hexon (II) but hexon is the only ORF encoded
356 on TRXPT_14, whereas the 5'-most ORF on TRXPT_20 is pX (pre-Mu) followed by pVI and hexon in
357 that order. L4 mRNAs include TRXPT_9, TRXPT_10, TRXPT_11, and TRXPT_13 all of which begin with
358 the TPL followed by three (TRXPT_9, TRXPT_10, and TRXPT_13) or four (TRXPT_11) coding exons.
359 These are the largest transcripts found in the transcriptome, each one possessing the CP for several similar
360 late proteins. Normally, MLTU transcripts encoding particular ORFs splice the TPL onto a splice site just
361 upstream of the ORF to be expressed (17). While this holds true for most MLTU ORFs, several late ORFs

362 (pVI, protease, and ORF7) do not have such close proximity splicing but are contained in larger transcripts
363 such as these L4 mRNAs, strongly suggesting the use of non-standard ribosomal initiation mechanisms
364 such as secSSC utility and ribosome shunting found in other AdVs for their translation (17, 27). TRXPT_9
365 and TRXPT_10 are very similar but not identical. The last exon of TRXPT_9 seems to be truncated and
366 probably shares the same TTS as the other L4 mRNAs. They are both 6-exon transcripts encoding pVII
367 as the 5'-most ORF (fourth exon) and also have the CP for pX, pVI, hexon, a longer variant of protease
368 (eProt) – uses an upstream in-frame SSC than predicted, and ORF12 (a novel unpredicted 120 aa protein).
369 TRXPT_10 (and TRXPT_9 with the L4 TTS) also has the CP for pVIII and eE3. Conversely, TRXPT_11 is
370 a seven-exon mRNA with hexon as its 5'-most ORF but it also has the CP for eProt, ORF12, e33K, and
371 also pVIII and eE3 in that order. TRXPT_13 seems to be an E3 ORF utilizing the MLP TSS as it encodes
372 classical L4P genes such as pVIII and eE3 in that order similar to TRXPT_22 (E3 TU) but lacks TRXPT_22's
373 novel first ORF (8.3KII).

374 Lastly, the L5 class includes only TRXPT_12 which contains the TPL and a coding terminal exon. Its 5'-
375 most ORF is fiber (IV) but it also has the CP for the THEV specific gene, ORF7. TRXPT_12's CP is identical
376 to TRXPT_27 of the the E3 TU but they differ in their 5'-UTRs.

377 **DISCUSSION/CONCLUSIONS**

378 expression level changes over time: The pattern seen is similar to other AdVs. the MLTU is significantly
379 expressed relative to other TUs because at 12h.p.i, the infectious cycle is well underway. And the MLP is
380 expressed sub optimally even before full activation of all late genes

381 For fig2a: There is a dramatic increase of mean coverage/depth from **2.42** at 4 h.p.i to **95,042** at 72 h.p.i,
382 strongly demonstrating an active infection. Unexpectedly, the pileup of reads seems consistently skewed
383 over similar regions of the genome. We could speculate that the temporal gene expression regulation
384 of THEV is different from MAdVs or this could simply mean that the infection was not well synchronized.
385 However, the relative proportions over these similar regions shows some variation over time. For fig2b:
386 titer reaching a plateau at 120 h.p.i, probably due to high cell death TRXPT_2 and ORF1 are isoforms
387 Presumably, if the junction reads were normalized, MLTU would not be predominant at 12hpi. The TTSs
388 were all in the context of A-T rich sequences; which presumably serve as polyA signals. All splice junctions
389 were confirmed by cloning and Sanger sequencing of cDNA (**Supplementary PCR methods**). We did not
390 find the x,y,z or i-leaders for MLP transcripts probably because THEV doesn't use it due to its smaller size
391 The E3 ORF has an upstream, in-frame SSC.

392 **MATERIALS AND METHODS**

393 **Cell culture and THEV Infection**

394 The Turkey B-cell line (MDTC-RP19, ATCC CRL-8135) was grown as suspension cultures in 1:1 complete
395 Leibovitz's L-15/McCoy's 5A medium with 10% fetal bovine serum (FBS), 20% chicken serum (ChS), 5%
396 tryptose phosphate broth (TPB), and 1% antibiotics solution (100 U/mL Penicillin and 100ug/mL Strepto-
397 mycin), at 41°C in a humidified atmosphere with 5% CO₂. Infected cells were maintained in 1:1 serum-
398 reduced Leibovitz's L15/McCoy's 5A media (SRLM) with 2.5% FBS, 5% ChS, 1.2% TPB, and 1% antibiotics
399 solution (100 U/mL Penicillin and 100ug/mL Streptomycin). A commercially available HE vaccine was pur-
400 chased from Hygieia Biological Labs as a source of THEV-A (VAS strain). The stock virus was titrated using
401 an in-house qPCR assay with titer expressed as genome copy number(GCN)/mL, similar to Mahshoub *et al*
402 (28) with modifications. Cells were infected in triplicates at a multiplicity of infection (MOI) of 100 GCN/cell,
403 incubate at 41°C for 1 hour, and washed three times to get rid of free virion particles. Samples in tripli-
404 cates were harvested at 4-, 12-, 24-, and 72-h.p.i for total RNA extraction. The infection was repeated but
405 samples in triplicates were harvested at 12-, 24-, 36-, 48-, and 72-h.p.i for PCR validation of novel splice
406 sites. Still one more independent infection was done at time points ranging from 12 to 168-h.p.i for qPCR
407 quantification of virus titers.

408 **RNA extraction and Sequencing**

409 Total RNA was extracted from infected cells using Thermofishers' RNAqueous™-4PCR Total RNA Isolation
410 Kit (#AM1914) per manufacturer's instructions. An agarose gel electrophoresis was performed to check
411 RNA integrity. The RNA quantity and purity was initially assessed using nanodrop, and RNA was used only
412 if the A260/A280 ratio was 2.0 ± 0.05 and the A260/A230 ratio was >2 and <2.2. Extracted total RNA sam-
413 ples were sent to LC Sciences, Houston TX for poly-A-tailed mRNA sequencing where RNA integrity was
414 checked with Agilent Technologies 2100 Bioanalyzer High Sensitivity DNA Chip and poly(A) RNA-
415 seq library was prepared following Illumina's TruSeq-stranded-mRNA sample preparation protocol.
416 Paired-end sequencing was performed on Illumina's NovaSeq 6000 sequencing system.

417 **Validation of Novel Splice Junctions**

418 All splice junctions identified in this work are novel except one predicted splice site each for pTP and DBP,
419 which were corroborated in our work. However, these predicted splice junctions had not been experimen-

420 tally validated hitherto, and we identified additional novel exons, giving the complete picture of these tran-
421 scripts. The novel splice junctions in this work discovered in the assembled transcripts using the StringTie
422 transcript assembler which we validated by PCR and Sanger Sequencing are shown in **Supplementary**
423 **PCR methods.** Briefly, we designed primers that crossed a range of novel exon-exon boundaries for each
424 specific transcript in a transcription unit (TU) paired with their respective universal primers for the TU. Each
425 forward primer contained a KpnI restriction site and reverse primers, an XbaI site. After first-strand cDNA
426 synthesis with SuperScript™ III First-Strand Synthesis System, these primers were used in a targeted
427 PCR amplification, the products analyzed with agarose gel electrophoresis to confirm expected band sizes,
428 cloned by traditional restriction enzyme method, and Sanger sequenced to validate these splice junctions
429 at the sequence level.

430 **3' Rapid Amplification of cDNA Ends (3'-RACE)**

431 We performed a rapid amplification of sequences from the 3' ends of mRNAs (3'-RACE) experiment us-
432 ing a portion of the extracted total RNA of infected MDTC-RP19 cells used for the RNA-seq experiment
433 as explained above. We followed the protocol described by Green *et al* (29) with modifications. Briefly,
434 1ug of total RNA was reverse transcribed to cDNA using SuperScript™ IV First-Strand Synthesis System
435 following the manufacturing instructions using an adapter-primer with a 3'-end poly(T) and a 5'-end BamHI
436 restriction site. A gene-specific sense primer with a 5'-end KpnI restriction site paired with an anti-sense
437 adapter-primer with a 5'-end BamHI site were used to amplify target sections of the cDNA using Invitrogen's
438 Platinum™ Taq DNA polymerase High Fidelity, following manufacturer's instructions. The PCR amplicons
439 were restriction digested, cloned, and Sanger sequenced.

440 **Computational Analysis of RNA Sequencing Data: Mapping and Transcript characterization**

441 Our sequence reads were analyzed following a well established protocol described by Pertea *et*
442 *al* (24), using Snakemake - version 7.24.0 (30), a popular workflow management system to
443 drive the pipeline. Briefly, sequencing reads were trimmed with the Trim-galore - version
444 0.6.6 (31) program to achieve an overall Mean Sequence Quality (Phred Score) of 36. Trimmed
445 reads were mapped simultaneously to the complete genomic sequence of avirulent turkey hemor-
446 rhagic enteritis virus (<https://www.ncbi.nlm.nih.gov/nuccore/AY849321.1/>) and *Meleagris gallopavo*
447 (<https://www.ncbi.nlm.nih.gov/genome/?term=Meleagris+gallopavo>) using Hisat2 - version 2.2.1 (24)
448 with default settings. The generated alignment (BAM) files from each infection time point were filtered

449 for reads mapping to the THEV genome using `Samtools` – version 1.16.1 and fed into `StringTie` –
450 version 2.2.1 (24) to assemble the transcripts, using a GTF annotation file derived from a GFF3 annotation
451 file obtained from NCBI, which contains the predicted ORFs of THEV as a guide. `GFFCOMPARE` – version
452 0.12.6 was used to merge all transcripts from all time points without redundancy and using a custom R
453 script, adenovirus transcripts units (regions) were assigned to each transcript, generating the transcriptome
454 of THEV. `StringTie` set to expression estimation mode was used to calculate FPKM scores for all
455 transcripts after which `Ballgown` – version 2.33.0 in R was used to perform the statistical analysis on the
456 transcript expression levels. `Samtools` was also used to count the total sequencing reads for all replicates
457 at each time point and `Regtools` – version 1.0.0 was used to count all junctions, the reads supporting
458 them, and extract all other information related to the junction. See **Supplementary Computational**
459 **Analysis** for the details of transcript expression level estimations and splice junction read counts.

460 **SCRIPTS AND SUPPLEMENTARY MATERIALS**

461 **DATA AVAILABILITY**

462 **CODE AVAILABILITY**

- 463 All the code/scripts written for analysis of the data are available on github ([https://github.com/Abraham-
464 Quaye/thev_transcriptome](https://github.com/Abraham-Quaye/thev_transcriptome))

465 **ACKNOWLEDGMENTS**

466 LC Sciences - RNA sequencing was done here

467 Eton Bioscience, Inc, San Diego, CA - All Sanger sequencing validations was done here BYU comput-

468 ing systems

469 REFERENCES

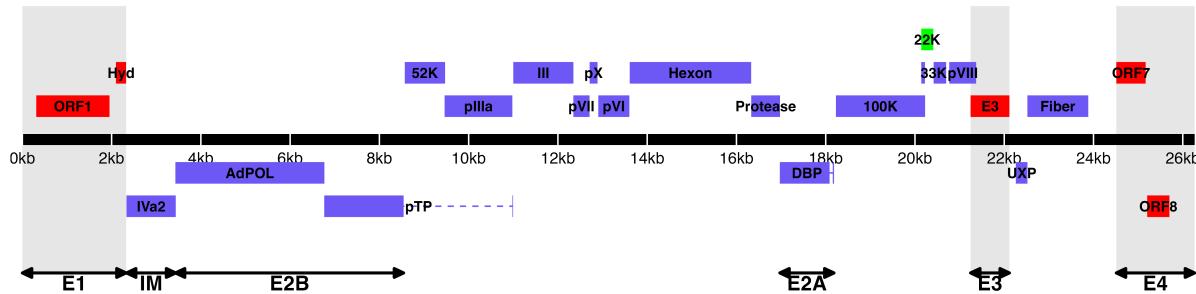
- 470 1. Davison A, Benko M, Harrach B. 2003. Genetic content and evolution of adenoviruses. *The Journal*
471 of general virology
- 472 2. Harrach B. 2008. Adenoviruses: General features, p. 1–9. *In* Mahy, BWJ, Van Regenmortel, MHV
473 (eds.), *Encyclopedia of virology* (third edition). Book Section. Academic Press, Oxford.
- 474 3. Upton C, Slack S, Hunter AL, Ehlers A, Roper RL. 2003. Poxvirus orthologous clusters: Toward
475 defining the minimum essential poxvirus genome. *Journal of virology* 77:7590–7600.
- 476 4. McGeoch D, Davison AJ. 1999. Chapter 17 - the molecular evolutionary history of the herpesviruses,
477 p. 441–465. *In* Domingo, E, Webster, R, Holland, J (eds.), *Origin and evolution of viruses*. Book
Section. Academic Press, London.
- 478 5. Harrach B, Benko M, Both GW, Brown M, Davison AJ, Echavarría M, Hess M, Jones M, Kajon A,
Lehmkuhl HD, Mautner V, Mittal S, Wadell G. 2011. Family adenoviridae. *Virus Taxonomy: 9th*
479 *Report of the International Committee on Taxonomy of Viruses* 125–141.
- 480 6. Guimet D, Hearing P. 2016. 3 - adenovirus replication, p. 59–84. *In* Curiel, DT (ed.), *Adenoviral*
481 *vectors for gene therapy* (second edition). Book Section. Academic Press, San Diego.
- 482 7. Kovács ER, Benkő M. 2011. Complete sequence of raptor adenovirus 1 confirms the characteristic
483 genome organization of siadenoviruses. *Infection, Genetics and Evolution* 11:1058–1065.
- 484 8. Davison AJ, Wright KM, Harrach B. 2000. DNA sequence of frog adenovirus. *J Gen Virol* 81:2431–
485 2439.
- 486 9. Kovács ER, Jánoska M, Dán Á, Harrach B, Benkő M. 2010. Recognition and partial genome char-
487 acterization by non-specific DNA amplification and PCR of a new siadenovirus species in a sample
originating from parus major, a great tit. *Journal of Virological Methods* 163:262–268.
- 488 10. Katoh H, Ohya K, Kubo M, Murata K, Yanai T, Fukushi H. 2009. A novel budgerigar-adenovirus
489 belonging to group II avian adenovirus of siadenovirus. *Virus Research* 144:294–297.
- 490 11. Beach NM. 2006. Characterization of avirulent turkey hemorrhagic enteritis virus: A study of the
491 molecular basis for variation in virulence and the occurrence of persistent infection. Thesis.

- 492 12. Beach NM, Duncan RB, Larsen CT, Meng XJ, Sriranganathan N, Pierson FW. 2009. Comparison of
493 12 turkey hemorrhagic enteritis virus isolates allows prediction of genetic factors affecting virulence.
494 J Gen Virol 90:1978–85.
- 495
- 496 13. Gross WB, Moore WE. 1967. Hemorrhagic enteritis of turkeys. Avian Dis 11:296–307.
- 497
- 498 14. Rautenschlein S, Sharma JM. 2000. Immunopathogenesis of haemorrhagic enteritis virus (HEV) in
499 turkeys. Dev Comp Immunol 24:237–46.
- 500 15. Larsen CT, Domermuth CH, Sponenberg DP, Gross WB. 1985. Colibacillosis of turkeys exacerbated
501 by hemorrhagic enteritis virus. Laboratory studies. Avian Dis 29:729–32.
- 502 16. Dhama K, Gowthaman V, Karthik K, Tiwari R, Sachan S, Kumar MA, Palanivelu M, Malik YS, Singh
503 RK, Munir M. 2017. Haemorrhagic enteritis of turkeys – current knowledge. Veterinary Quarterly
504 37:31–42.
- 505
- 506 17. Donovan-Banfield I, Turnell AS, Hiscox JA, Leppard KN, Matthews DA. 2020. Deep splicing plasticity
507 of the human adenovirus type 5 transcriptome drives virus evolution. Communications Biology 3:124.
- 508 18. Zhao H, Chen M, Pettersson U. 2014. A new look at adenovirus splicing. Virology 456-457:329–341.
- 509
- 510 19. Wolfrum N, Greber UF. 2013. Adenovirus signalling in entry. Cell Microbiol 15:53–62.
- 511
- 508 20. Falvey E, Ziff E. 1983. Sequence arrangement and protein coding capacity of the adenovirus type 2
509 "i" leader. Journal of Virology 45:185–191.
- 510 21. Morris SJ, Scott GE, Leppard KN. 2010. Adenovirus late-phase infection is controlled by a novel L4
511 promoter. Journal of Virology 84:7096–7104.

- 512 22. Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, Tanzer A, Lagarde J, Lin W,
Schlesinger F, Xue C, Marinov GK, Khatun J, Williams BA, Zaleski C, Rozowsky J, Röder M, Kokocinski F, Abdelhamid RF, Alioto T, Antoshechkin I, Baer MT, Bar NS, Batut P, Bell K, Bell I, Chakrabortty S, Chen X, Chrast J, Curado J, Derrien T, Drenkow J, Dumais E, Dumais J, Duttagupta R, Falconnet E, Fastuca M, Fejes-Toth K, Ferreira P, Foissac S, Fullwood MJ, Gao H, Gonzalez D, Gordon A, Gunawardena H, Howald C, Jha S, Johnson R, Kapranov P, King B, Kingswood C, Luo OJ, Park E, Persaud K, Preall JB, Ribeca P, Risk B, Robyr D, Sammeth M, Schaffer L, See L-H, Shahab A, Skancke J, Suzuki AM, Takahashi H, Tilgner H, Trout D, Walters N, Wang H, Wrobel J, Yu Y, Ruan X, Hayashizaki Y, Harrow J, Gerstein M, Hubbard T, Reymond A, Antonarakis SE, Hannon G, Giddings MC, Ruan Y, Wold B, Carninci P, Guigó R, Gingeras TR. 2012. Landscape of transcription in human
513 cells. *Nature* 489:101–108.
- 514 23. Aboeza Z, Mabsoub H, El-Bagoury G, Pierson F. 2019. In vitro growth kinetics and gene expression
515 analysis of the turkey adenovirus 3, a siadenovirus. *Virus Research* 263:47–54.
- 516 24. Pertea M, Kim D, Pertea GM, Leek JT, Salzberg SL. 2016. Transcript-level expression analysis of
517 RNA-seq experiments with HISAT, StringTie and ballgown. *Nature Protocols* 11:1650–1667.
- 518 25. Jack Fu [Aut], Alyssa C. Frazee [Aut, Cre], LeonardoCollado-Torres [Aut], Andrew E. Jaffe [Aut],
519 Jeffrey T. Leek[Aut, Ths]. 2017. Ballgown. Bioconductor.
- 520 26. Pitcovski J, Mualem M, Rei-Koren Z, Krispel S, Shmueli E, Peretz Y, Gutter B, Gallili GE, Michael A,
Goldberg D. 1998. The complete DNA sequence and genome organization of the avian adenovirus,
521 hemorrhagic enteritis virus. *Virology* 249:307–315.
- 522 27. Yueh A, Schneider RJ. 1996. Selective translation initiation by ribosome jumping in adenovirus-
523 infected and heat-shocked cells. *Genes & Development* 10:1557–1567.
- 524 28. Mabsoub HM, Evans NP, Beach NM, Yuan L, Zimmerman K, Pierson FW. 2017. Real-time PCR-
525 based infectivity assay for the titration of turkey hemorrhagic enteritis virus, an adenovirus, in live
vaccines. *Journal of Virological Methods* 239:42–49.
- 526 29. Green MR, Sambrook J. 2019. Rapid amplification of sequences from the 3' ends of mRNAs: 3'-
527 RACE. *Cold Spring Harbor Protocols* 2019:pdb.prot095216.

- 528 30. Mölder F, Jablonski KP, Letcher B, Hall MB, Tomkins-Tinch CH, Sochat V, Forster J, Lee S, Twardziok
529 SO, Kanitz A, Wilm A, Holtgrewe M, Rahmann S, Nahnsen S, Köster J. 2021. Sustainable data
analysis with snakemake. *F1000Research* 10:33.
- 530 31. Krueger F, James F, Ewels P, Afyounian E, Weinstein M, Schuster-Boeckler B, Hulselmans G, Scla-
531 mons. 2023. FelixKrueger/TrimGalore: v0.6.10 - add default decompression path. Zenodo.

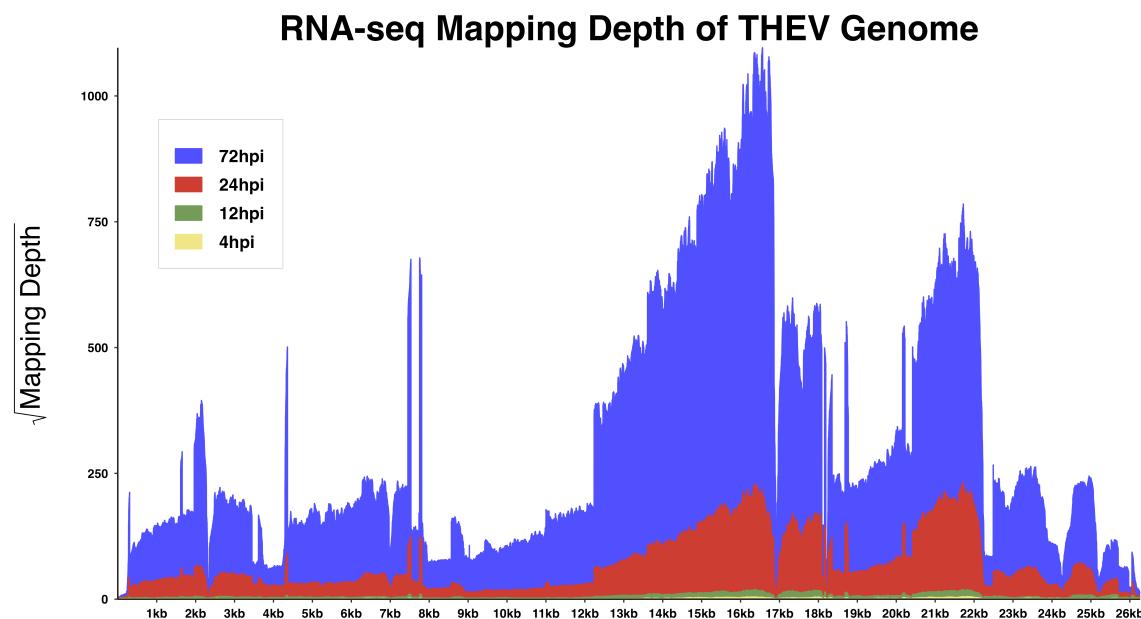
532 **TABLES AND FIGURES**



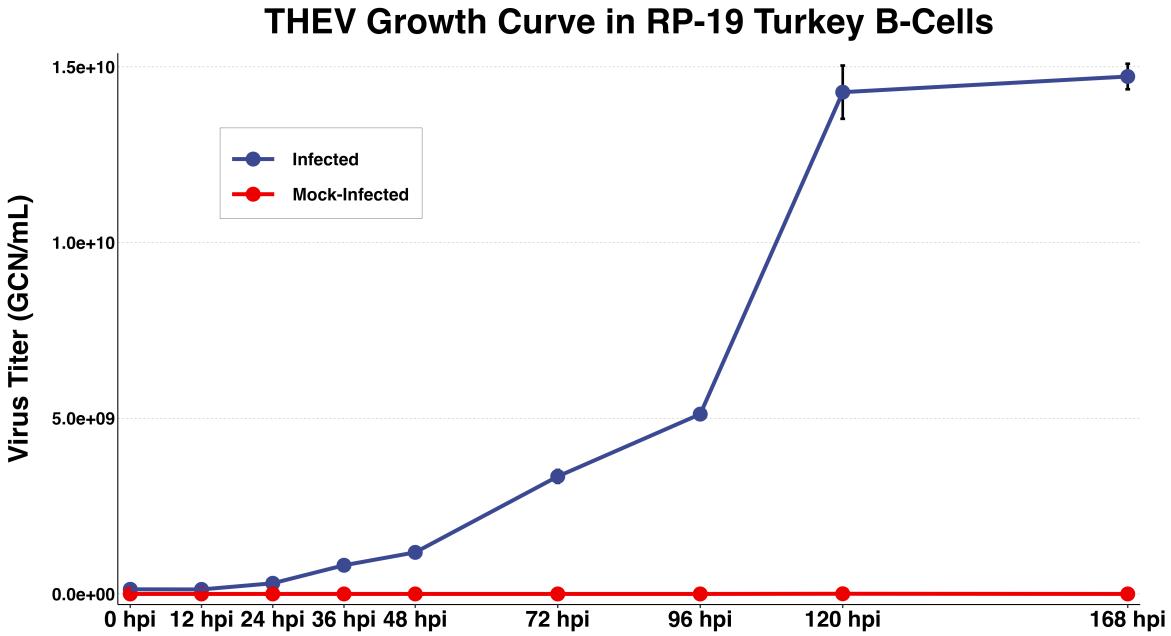
533

534 **Figure 1. Genomic map of THEV avirulent strain annotated ORFs.** The central horizontal line repre-
 535 presents the double-stranded DNA marked at 5kb intervals as white line breaks. Blocks represent viral genes.
 536 Blocks above the DNA line are transcribed rightward, those below are transcribed leftward. pTP, DBP
 537 and 33K predicted to be spliced are shown as having tails. Shaded regions indicate regions containing
 538 "genus-specific" genes (colored red). Genes colored in blue are "genus-common". Gene colored in light
 539 green is conserved in all but Atadenoviruses. The UXP (light blue) is an incomplete gene present in almost
 540 all AdVs. Regions comprising the different transcription units are labelled at the bottom (E1, E2A, E2B, E3,
 541 E4, and IM); the unlabeled regions comprise the MLTU.

A



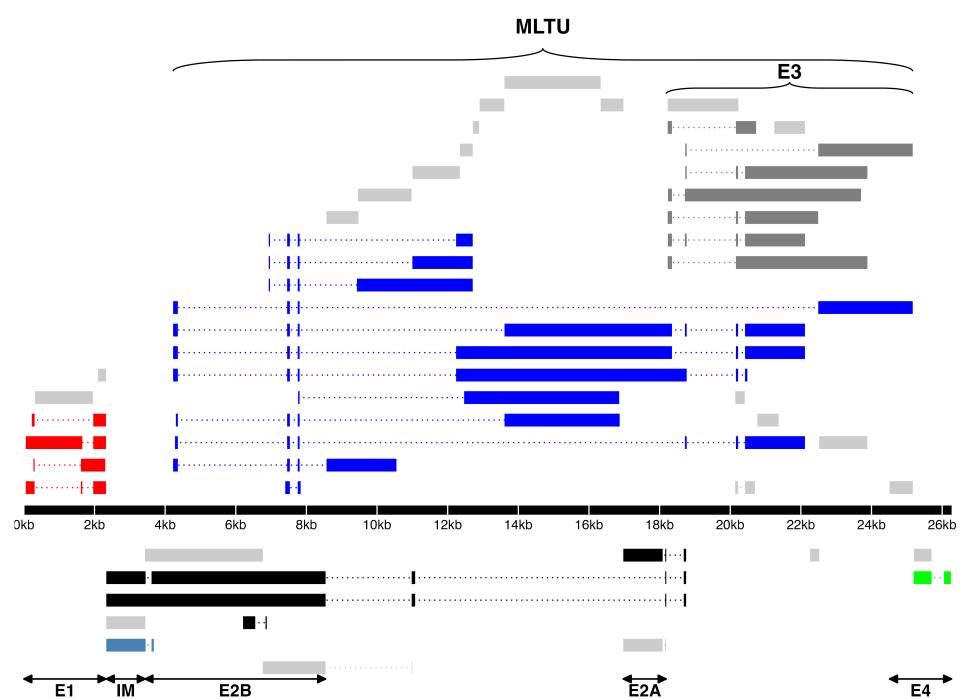
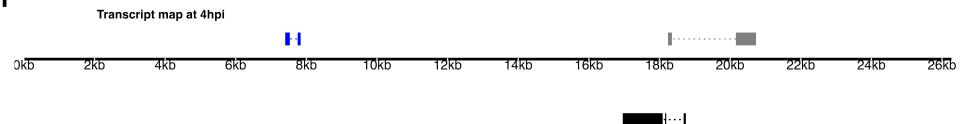
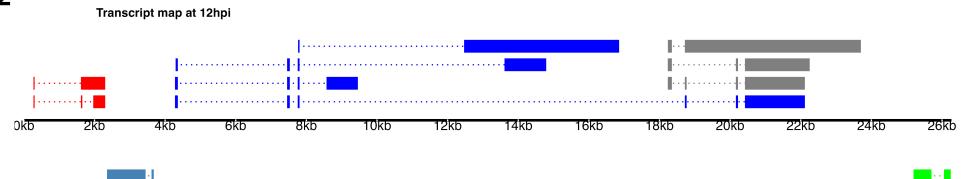
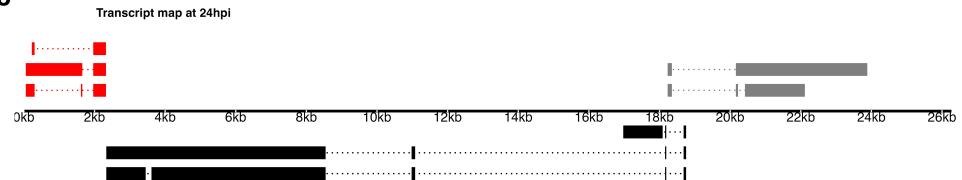
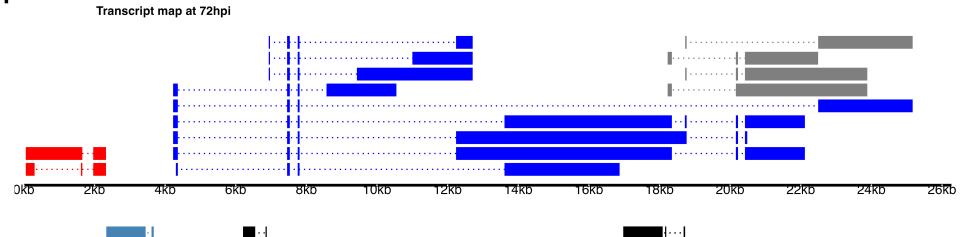
B



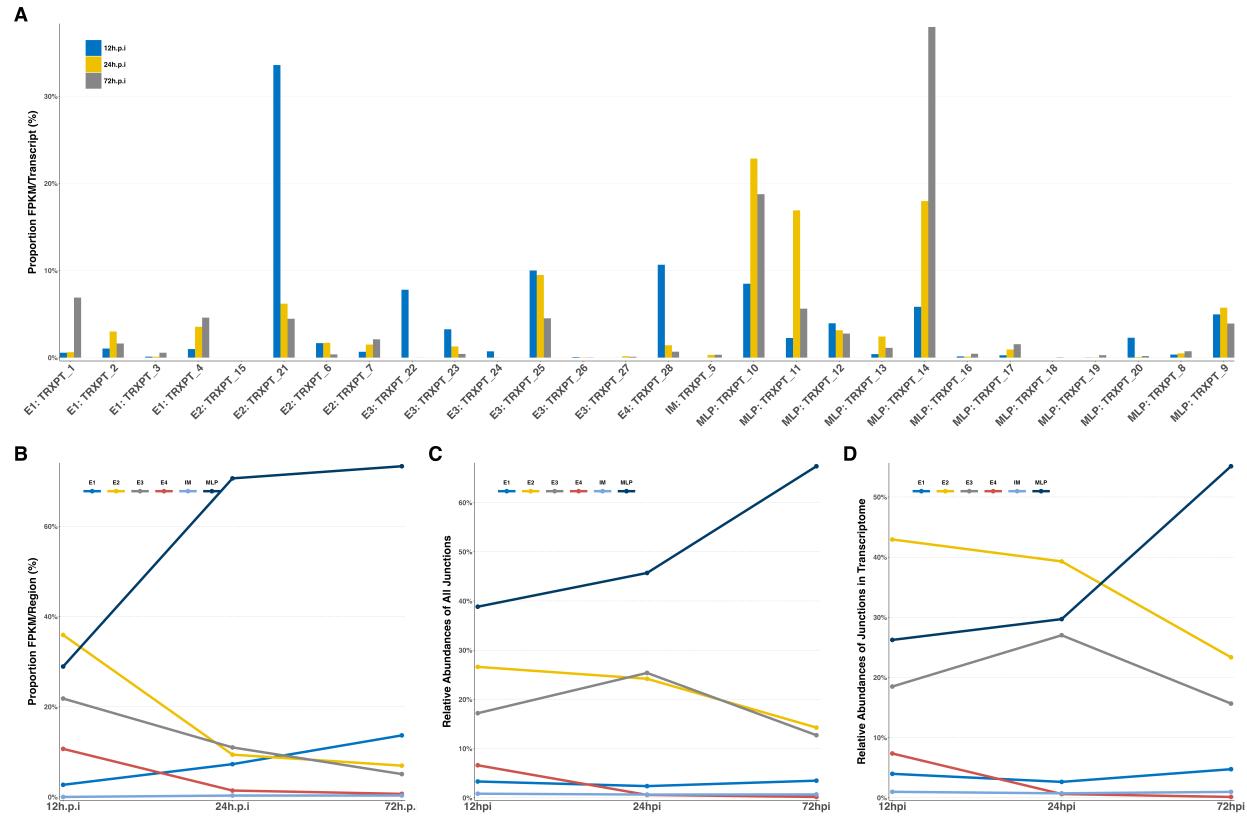
542

543 **Figure 2: Increasing levels of THEV over time. a) Per base coverage of sequence reads mapping to
544 THEV genome by time point.** The pileup of mRNA reads mapping to THEV genome at the base-pair level
545 for each indicated time point. **b) Growth curve of THEV (VAS vaccine strain) in MDTC-RP19 cell line.**
546 Virus titers were quantified with a qPCR assay. There is no discernible increase in virus titer up 12 h.p.i,
547 after which a steady increase in virus titer is measured. The virus titer expands exponentially beginning

⁵⁴⁸ from 48 h.p.i, increasing by orders of magnitude before reaching a plateau at 120 h.p.i. GCN: genome copy
⁵⁴⁹ number.

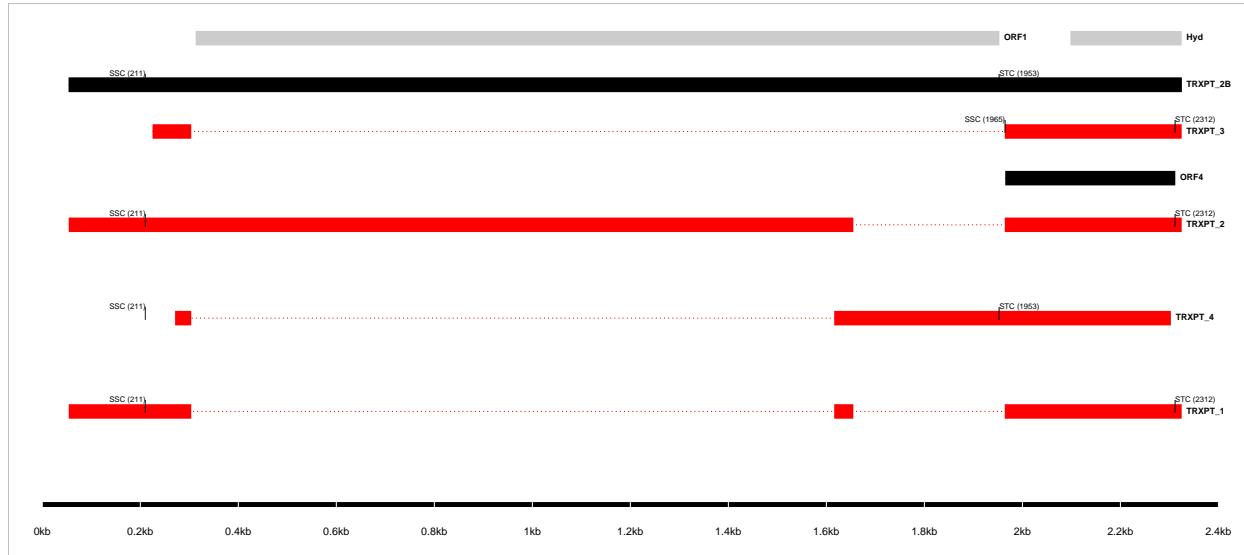
A**B1****B2****B3****B4**

551 **Figure 3. a) Transcriptome of THEV from RNA-seq.** THEV transcripts assembled from all time points
 552 by StringTie are unified forming this final transcriptome (splicing map). Transcripts belonging to the same
 553 transcription unit (TU) are located in close proximity on the genome and are color coded and labeled in this
 554 figure as such. The organization of TUs in the THEV genome is unsurprisingly similar to MAdVs; however,
 555 the MAdV genome shows significantly more transcripts. The TUs are color coded: E1 transcripts - red, E2
 556 - black, E3 - dark grey, E4 - green, MLTU - blue. Predicted ORFs are also indicated here, colored light grey.
 557 **b) THEV transcripts identified at given time points.** Transcripts are color coded as explained in (a).



558 **Figure 4: Changes in splicing and expression profile of THEV over time.** **a)** Normalized (FPKM)
 559 expression levels of transcripts over time. The expression levels (FPKM) of individual transcripts as a
 560 percentage of the total expression of all transcripts at each time point are indicated. Only transcripts from
 561 our RNA-seq data are included here. **b)** Normalized (FPKM) expression levels of transcripts by region over
 562 time. The expression levels of each region/TU as a percentage of the total expression of all transcripts at
 563 each time point are indicated. Region expression levels were calculated by summing up the FPKMs of all
 564 transcripts categorized in that region. **c)** Relative abundances of all splice junctions grouped by region/TU
 565 over time. After assigning all 2,457 unique junctions to a TU and the total junction reads counted at each
 566 time point for each region, the total junction reads for each TU plotted as percentage of all junction reads at
 567 each time point for each region.

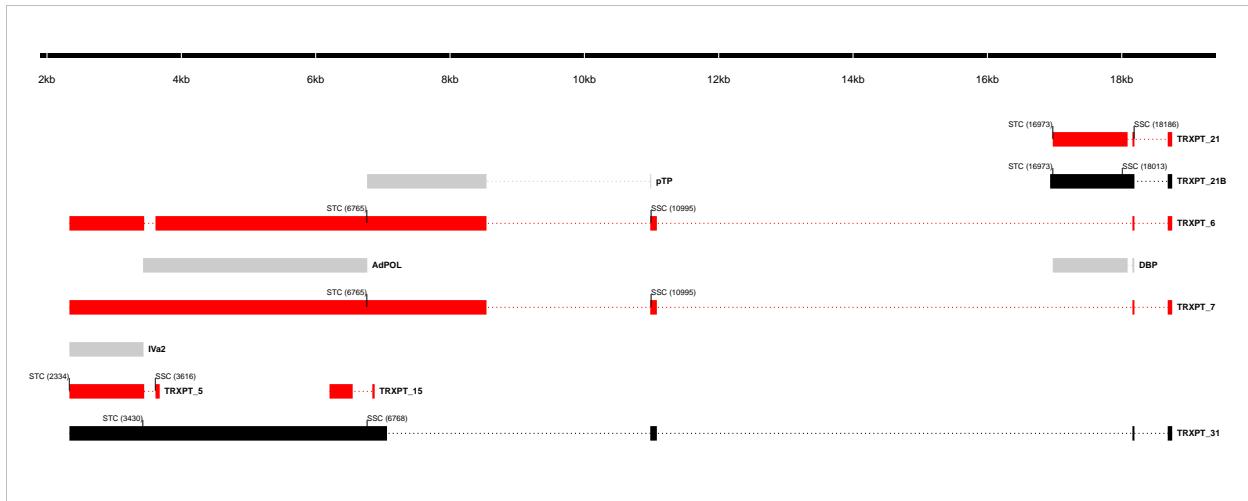
568 each time point is indicated. Note that the junction read counts are not normalized. **d) Relative abundances**
 569 *of junctions in transcriptome grouped by region/TU over time*. This is identical to **(c)**, except that only the
 570 junctions found in the full transcriptome obtained from the RNA-seq data were included.



Transcript ID	Splice Junction					Strand	Junction Reads				Junction Status
	Start	End	Intron Length	Splice Donor-Acceptor			4h.p.i	12h.p.i	24h.p.i	72h.p.i	
TRXPT_1, TRXPT_4	304	1616	1313bp	GT-AG		+	0	9	1019	25041	Validated [*]
TRXPT_3	304	1964	1661bp	GT-AG		+	0	2	168	1588	Validated
TRXPT_2, TRXPT_1	1655	1964	310bp	GT-AG		+	0	9	1395	38491	Validated

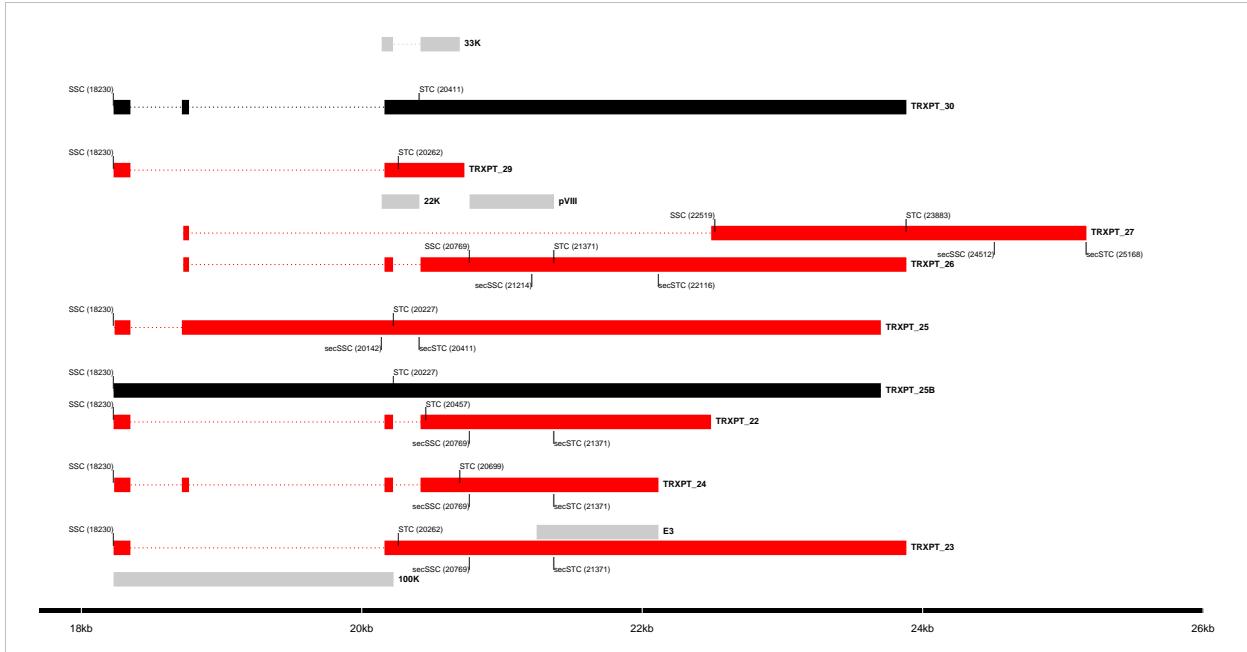
571 ^{*}Not validated for TRXPT_4

572 **Figure 5: The splice map of the E1 transcription unit (TU).** Exons are depicted as boxes connected by
 573 introns (dotted lines). Transcripts from RNA-seq data are colored red, predicted ORFs are colored grey, and
 574 transcripts or ORFs discovered by other means are colored black. Each transcript or ORF is labelled with
 575 its name to the right. The start codon (SSC) and stop codon (STC) of the 5'-most CDS of each transcript
 576 is indicated with the nucleotide position in brackets. The region of the virus is depicted at the bottom as a
 577 black line with labels of the nucleotide positions for reference. The table shows sequence reads covering
 578 the splice junctions with information about their validation status using cloning and Sanger sequencing.



Transcript ID	Splice Junction				Strand	region	Junction Reads				Junction Status
	Start	End	Splice Donor-Acceptor	Intron Length			4h.p.i	12h.p.i	24h.p.i	72h.p.i	
TRXPT_5 TRXPT_7	3447	3615	GT-AG	169bp	-	IM, E2	1	5	720	13422	Validated
TRXPT_6 TRXPT_7	11079	18159	GT-AG	7081bp	-	E2	0	2	0	0	Validated
TRXPT_21	18087	18159	GT-AG	73bp	-	E2	9	103	0	0	Validated
TRXPT_21, TRXPT_6, TRXPT_7	18189	18684	GT-AG	496bp	-	E2	0	111	18794	156037	Validated
TRXPT_6, TRXPT_7	8543	10981	GT-AG	2439bp	-	E2	0	0	298	850	Validated
TRXPT_15	6551	6843	GT-GC	293bp	-	E2	0	0	0	6	Validated

579 **Figure 6: The splice map of the E2 and IM TUs.** Exons are depicted as boxes connected by introns (dotted lines). Red transcripts are generated from RNA-seq data and predicted ORFs are colored grey. TRXPT_21B discovered by 3'RACE is colored black. Each transcript or ORF is labelled with its name to the right. The SSC and STC of the 5'-most CDS of each transcript is indicated with the nucleotide position in brackets. The region of the virus is depicted at the bottom as a black line with labels of the nucleotide positions for reference. The table shows sequence reads covering the splice junctions with information about their validation status using cloning and Sanger sequencing.



Transcript ID	Splice Junction					Junction Reads					Junction Status
	Start	End	Splice Donor-Acceptor	Intron Length	Strand	region	4h.p.i	12h.p.i	24h.p.i	72h.p.i	
TRXPT_25, TRXPT_24, TRXPT_10	18350	18717	GT-AG	368bp	+	E3, MLP	4	21	3930	35490	Validated
TRXPT_23, TRXPT_22, TRXPT_11	18350	20162	GT-AG	1813bp	+	E3, MLP	3	18	6619	38841	Validated
TRXPT_26, TRXPT_24, TRXPT_13, TRXPT_11, TRXPT_10	18768	20162	GT-AG	1395bp	+	E3, MLP	2	21	5207	45062	Validated
TRXPT_26, TRXPT_22, TRXPT_24, TRXPT_13, TRXPT_11, TRXPT_10	20223	20419	GT-AG	197bp	+	E3, MLP	3	33	10583	93238	Validated
587 TRXPT_27	18768	22492	GT-AG	3725bp	+	E3	0	0	101	1950	Validated

588 **Figure 7: The splice map of the E3 TU.** Exons are depicted as boxes connected by introns (dotted
 589 lines). Red transcripts are generated from RNA-seq data and predicted ORFs are colored grey. Transcripts
 590 discovered by other means are colored black. Each transcript or ORF is labelled with its name to the right.
 591 The start codon (SSC) and stop codon (STC) of the 5'-most CDS of each transcript is indicated with the
 592 nucleotide position in brackets. Similarly, the secondary SSC (secSSC) and secondary STC (secSTC)
 593 are shown. The region of the virus is depicted at the bottom as a black line with labels of the nucleotide
 594 positions for reference. The table shows sequence reads covering the splice junctions with information
 595 about their validation status using cloning and Sanger sequencing.



597 **Figure 8: The splice map of the E4 TU.** Exons are depicted as boxes connected by introns (dotted lines).
 598 The transcript from RNA-seq data is colored red and the predicted ORF, grey. The transcript and ORF are
 599 labelled with their names to the right. The start codon (SSC) and stop codon (STC) of the 5'-most CDS
 600 is indicated with the nucleotide position in brackets. The region of the virus is depicted at the bottom as a
 601 black line with labels of the nucleotide positions for reference. The table shows sequence reads covering
 602 the splice junction with its validation status using cloning and Sanger sequencing.



604 **Figure 9: The splice map of the MLTU.** Exons are depicted as boxes connected by introns (dotted lines).
605 The transcripts from our RNA-seq data are colored red and the predicted ORFs, grey. The transcripts and
606 ORFs are labelled with their names to the right. The start codon (SSC) and stop codon (STC) of the 5'-most
607 CDS of each transcript is indicated with the nucleotide position in brackets. Similarly, the secondary SSC
608 (secSSC) and secondary STC (secSTC) are shown. The region of the virus is depicted at the bottom as a
609 black line with labels of the nucleotide positions for reference. The table shows sequence reads covering
610 the splice junctions with information about their validation status using cloning and Sanger sequencing.

Table 1: Table 1: Overview of sequencing results

Metric	4h.p.i	12h.p.i	24h.p.i	72h.p.i	Total
Total reads	1.17e+08	7.63e+07	1.20e+08	1.15e+08	4.28e+08
Mapped (Host)	1.04e+08	6.79e+07	1.06e+08	8.38e+07	3.62e+08
Mapped (THEV)	4.32e+02	6.70e+03	1.18e+06	1.69e+07	1.81e+07
Mean Per Base Coverage/Depth	2.42	37.71	6,666.96	95,041.7	101,749
Total unique splice junctions	13	37	236	2374	2,457
Junction coverage Total (at least 1 read)	37	605	115075	2132806	2.25e+06
Junction coverage Mean reads	2.8	16.4	487.6	898.4	351.3
Junction coverage (at least 10 reads)	0	13	132	1791	1,936
Junction coverage (at least 100 reads)	0	1	53	805	859
Junction coverage (at least 1000 reads)	0	0	18	168	186

Table 2: Table 2a: Most abundant splice junctions at 12h.p.i

Timepoint	Strand	Start	End	Splice_Site	Region	Intron Length	Reads (Percentage)
12hpi	-	18,087	18,159	GT-AG	E2	72 bp	103 (17%)
12hpi	+	18,189	18,684	CT-AC	MLP	495 bp	97 (16%)
12hpi	+	7,531	7,754	GT-AG	MLP	223 bp	58 (9.6%)
12hpi	-	25,701	26,055	GT-AG	E4	354 bp	37 (6.1%)
12hpi	+	20,223	20,419	GT-AG	E3	196 bp	33 (5.5%)
12hpi	+	4,360	7,454	GT-AG	MLP	3,094 bp	32 (5.3%)
12hpi	-	18,751	20,668	GT-AG	E2	1,917 bp	22 (3.6%)
12hpi	+	18,350	18,717	GT-AG	E3	367 bp	21 (3.5%)
12hpi	+	18,768	20,162	GT-AG	E3	1,394 bp	21 (3.5%)
12hpi	+	7,807	13,610	GT-AG	MLP	5,803 bp	18 (3%)
12hpi	+	18,350	20,162	GT-AG	E3	1,812 bp	18 (3%)
12hpi	-	18,189	18,684	GT-AG	E2	495 bp	14 (2.3%)
12hpi	-	18,751	21,682	GT-AG	E2	2,931 bp	10 (1.7%)
12hpi	+	304	1,616	GT-AG	E1	1,312 bp	9 (1.5%)
12hpi	+	1,655	1,964	GT-AG	E1	309 bp	9 (1.5%)
12hpi	-	18,087	18,163	GT-AG	E2	76 bp	8 (1.3%)
12hpi	+	7,807	12,238	GT-AG	MLP	4,431 bp	7 (1.2%)
12hpi	+	7,807	22,492	GT-AG	MLP	14,685 bp	6 (1%)

Table 3: Table 2b: Most abundant splice junctions at 24h.p.i

Timepoint	Strand	Start	End	Splice_Site	Region	Intron Length	Reads (Percentage)
24hpi	-	18,087	18,159	GT-AG	E2	72 bp	18,825 (16.4%)
24hpi	+	18,189	18,684	CT-AC	MLP	495 bp	17,670 (15.4%)
24hpi	+	7,531	7,754	GT-AG	MLP	223 bp	12,319 (10.7%)
24hpi	+	20,223	20,419	GT-AG	E3	196 bp	10,583 (9.2%)
24hpi	+	4,360	7,454	GT-AG	MLP	3,094 bp	7,128 (6.2%)
24hpi	+	18,350	20,162	GT-AG	E3	1,812 bp	6,619 (5.8%)
24hpi	+	18,768	20,162	GT-AG	E3	1,394 bp	5,207 (4.5%)
24hpi	+	18,350	18,717	GT-AG	E3	367 bp	3,930 (3.4%)
24hpi	-	18,751	20,668	GT-AG	E2	1,917 bp	3,870 (3.4%)
24hpi	+	7,807	13,610	GT-AG	MLP	5,803 bp	2,553 (2.2%)
24hpi	+	7,807	12,238	GT-AG	MLP	4,431 bp	2,446 (2.1%)
24hpi	+	7,807	22,492	GT-AG	MLP	14,685 bp	1,642 (1.4%)
24hpi	+	1,655	1,964	GT-AG	E1	309 bp	1,395 (1.2%)
24hpi	+	7,807	18,717	GT-AG	MLP	10,910 bp	1,391 (1.2%)
24hpi	-	18,189	18,684	GT-AG	E2	495 bp	1,124 (1%)
24hpi	-	18,751	21,128	GT-AG	E2	2,377 bp	1,124 (1%)
24hpi	+	20,223	20,894	GT-AG	E3	671 bp	1,208 (1%)

Table 4: Table 2c: Most abundant splice junctions at 72h.p.i

Timepoint	Strand	Start	End	Splice_Site	Region	Intron Length	Reads (Percentage)
72hpi	+	7,531	7,754	GT-AG	MLP	223 bp	322,677 (15.1%)
72hpi	+	4,360	7,454	GT-AG	MLP	3,094 bp	179,607 (8.4%)
72hpi	-	18,087	18,159	GT-AG	E2	72 bp	161,336 (7.6%)
72hpi	+	18,189	18,684	CT-AC	MLP	495 bp	146,425 (6.9%)
72hpi	+	20,223	20,419	GT-AG	E3	196 bp	93,238 (4.4%)
72hpi	+	7,807	13,610	GT-AG	MLP	5,803 bp	81,420 (3.8%)
72hpi	+	7,807	12,238	GT-AG	MLP	4,431 bp	77,616 (3.6%)
72hpi	+	18,768	20,162	GT-AG	E3	1,394 bp	45,062 (2.1%)
72hpi	+	1,655	1,964	GT-AG	E1	309 bp	38,491 (1.8%)
72hpi	+	18,350	20,162	GT-AG	E3	1,812 bp	38,841 (1.8%)
72hpi	+	18,350	18,717	GT-AG	E3	367 bp	35,490 (1.7%)
72hpi	+	304	1,616	GT-AG	E1	1,312 bp	25,041 (1.2%)
72hpi	-	18,751	20,668	GT-AG	E2	1,917 bp	26,338 (1.2%)
72hpi	+	7,807	12,904	GT-AG	MLP	5,097 bp	21,946 (1%)
72hpi	+	7,807	22,492	GT-AG	MLP	14,685 bp	21,891 (1%)