

<sup>1</sup> Characterizing the Transcriptome of Turkey Hemorrhagic  
<sup>2</sup> Enteritis Virus in a Turkey B-cell Line

<sup>3</sup>

<sup>4</sup> **Running Title:** Novel Insights into Turkey Hemorrhagic Enteritis Virus Transcriptome

<sup>5</sup> Abraham Quaye<sup>1\*</sup>, Brett E. Pickett<sup>\*</sup>, Joel S. Griffitts<sup>\*</sup>, Bradford K. Berges<sup>\*</sup>, Brian D. Poole<sup>†\*\*</sup>

<sup>6</sup> \*Department of Microbiology and Molecular Biology, Brigham Young University

<sup>7</sup> <sup>1</sup>First-author

<sup>8</sup> <sup>†</sup> Corresponding Author

<sup>9</sup> **Corresponding Author Information**

<sup>10</sup> brian\_poole@byu.edu

<sup>11</sup> Department of Microbiology and Molecular Biology,

<sup>12</sup> 4007 Life Sciences Building (LSB),

<sup>13</sup> Brigham Young University,

<sup>14</sup> Provo, Utah

<sup>15</sup>

16 **ABSTRACT**

17 Characterizing the splice map of turkey hemorrhagic enteritis virus (THEV) is an essential step that would  
18 allow studies of individual genes mediating its immunosuppressive functions. We used an RNA-sequencing  
19 experiment to characterize the transcriptome of THEV for the first time, providing key insight into the THEV  
20 gene expression and mRNA structures. Researchers previously annotated THEV's genome as encoding  
21 23 open reading frames (ORFs). In this work we identified 29 spliced transcripts all of which consisted of  
22 novel exons although some exons matched some previously annotated ORFs. The three annotated splice  
23 junctions were also corroborated by our data. We performed PCR amplification of THEV cDNA, cloned the  
24 PCR products, and used Sanger sequencing to validate all identified splice junctions. During validation we  
25 identified five additional unique transcripts, a subset of which were further validated by 3' rapid amplification of  
26 cDNA ends (3' RACE) experiments. Thus, we report that the genome of THEV contains 34 unique transcripts  
27 with the coding capacity for all annotated ORFs. However, we found six of the previously annotated ORFs  
28 (ORF1, E3, 33K, ORF8, IVa2, and protease) to be truncated ORFs on the basis of the identification of an  
29 in-frame upstream start codon or the detection of additional coding exons. We also identified three of the  
30 annotated ORFs with longer or shorter isoforms, and seven novel unpredicted ORFs that could potentially  
31 be translated; although it is beyond the scope of this manuscript to investigate whether they are translated.  
32 Similar to other adenoviruses (AdVs), THEV also produces multiple distinctly spliced transcripts that code for  
33 the same proteins across its genome. Our data show that all THEV transcripts are spliced and organized  
34 into five transcription units under the control of their cognate promoters like other AdVs. However, our data  
35 suggest that the temporal regulation of THEV may be different from other AdVs.

36 **INTRODUCTION**

37 Adenoviruses (AdVs) are non-enveloped icosahedral-shaped DNA viruses, causing infection in virtually all  
38 types of vertebrates studied to date. Their double-stranded linear DNA genomes range between 26 and  
39 45kb in size, producing a broad repertoire of transcripts via highly complex alternative splicing patterns (1,  
40 2). The AdV genome is one of the most optimally economized; both the forward and reverse DNA strands  
41 harbor protein-coding genes, making it highly gene-dense. There are 16 genes termed “genus-common”  
42 that are homologous in all AdVs; these are thought to be inherited from a common ancestor. All other  
43 genes are termed “genus-specific”. The genus-specific genes tend to be located at the termini of the  
44 genome while genus-common genes are usually towards the center of the genome (1). This pattern is  
45 also observed in *Poxviridae* and *Herpesviridae*, which also have linear DNA genomes (1, 3, 4). The family  
46 *Adenoviridae* consists of five genera: *Mastadenovirus* (MAdV), *Aviadenovirus*, *Atadenovirus*, *Ichtadenovirus*,  
47 and *Siadenovirus* (SiAdV) to which turkey adenovirus 3 also called turkey hemorrhagic enteritis virus (THEV)  
48 belongs (5–10). Members of SiAdV have the smallest genome size (~26 kb) and gene content of all known  
49 AdVs (see **Figure 1**) (1, 2, 6).

50 Virulent THEV strains (THEV-V) and avirulent strains (THEV-A) of THEV both infect turkeys, with THEV-V  
51 causing hemorrhagic enteritis (HE), a debilitating acute disease predominantly affecting turkey pouls charac-  
52 terized by immunosuppression (IS), intestinal lesions leading to bloody diarrhea, and up to 80% mortality  
53 (2, 11–13). While the current vaccine strain (a THEV-A strain isolated from a pheasant; Virginia Avirulent  
54 Strain [VAS]) has proven effective at preventing HE in turkey pouls, it still retains its immunosuppressive  
55 ability. Thus, vaccinated birds are rendered more susceptible to opportunistic infections and death than  
56 unvaccinated birds leading to substantial economic losses (11, 14–16). To eliminate the immunosuppressive  
57 effect of the vaccine strain, a thorough investigation of the culprit viral genes mediating this phenomenon  
58 is essential. However, the transcriptome (splicing and gene expression patterns) of THEV has not been  
59 characterized, making an investigation of specific IS-related viral genes impractical.

60 A myriad of studies have elucidated the AdV transcriptome in fine detail (17, 18). However, a large  
61 preponderance of studies focus on MAdVs – specifically human AdVs. Thus, most of the current knowledge  
62 regarding AdV gene expression and replication is based on MAdV studies, which is generalized for all other  
63 AdVs (10, 19). MAdV transcription is temporally regulated; therefore, genes are categorized into five early  
64 transcription units (E1A, E1B, E2, E3, and E4), two intermediate (IM) units (pIX and IVa2), and one major late  
65 transcription unit (MLTU or major late promoter [MLP] region), which generates five families of late mRNAs  
66 (L1-L5) based on the polyadenylation site. An additional gene (UXP or U exon) is located on the reverse  
67 strand. The early genes encode non-structural proteins such as enzymes or host-cell modulating proteins,

primarily involved in DNA replication, or providing the necessary intracellular niche for optimal replication while late genes encode structural proteins that act as capsid proteins, promote virion assembly, or direct genome packaging. The immediate early genes E1A are expressed first, followed by the delayed early genes, E1B, E2, E3 and E4. Then the intermediate early genes, lVa2 and pIX are expressed followed by the late genes (10, 17, 18). It is noteworthy that the MLP shows basal transcriptional activity during early infection (before DNA replication), with a comparable efficiency to other early viral promoters, but it reaches its maximal activity during late infection (after DNA replication). However, during early infection only a subset of the MLP-derived transcripts are expressed (10). MAdV makes an extensive use of alternative RNA splicing to produce a very complex array of mRNAs. All but the pIX mRNA undergo at least one splicing event. For instance, the MLTU produces over 20 distinct splice variants all of which contain three non-coding exons at the 5'-end (collectively known as the tripartite leader; TPL) (17, 18). There is also an alternate three-exon 5' non-coding leader sequence present in varying amounts on a subset of MLTU mRNAs (known as the x-, y-, and z-leaders). Lastly, there is the i-leader exon, which is infrequently included between the second and third TPL exons, and codes for the i-leader protein (20). Thus, the MLTU produces a complex repertoire of mRNA with diverse 5' untranslated regions (UTRs) spliced onto different 3' coding exons which are grouped into five different 3'-end classes (L1-L5) based on polyadenylation site. Each transcription unit (TU) contains its own promoter driving the expression of the array of mRNA transcripts produced via alternative splicing in the unit (10, 17, 18). The promoters are activated at different phases of the infection by proteins from previously activated TUs. Paradoxically, the early-to-late phase transition during infection requires the L4 gene products, 22K and 33K, which should only be available after the transition. However, a promoter in the L4 region (L4P) that directs the expression of these two proteins independent of the MLP was found, resolving the paradox (10, 17, 21). During translation of AdV mRNA, recent studies using long-read direct RNA sequencing strongly suggest the potential usage of secondary start codons; adding to what was already a highly complex system for gene expression (17, 22).

High throughput sequencing methods have facilitated the discovery of many novel transcribed regions and splicing isoforms. It is also a very powerful tool to study alternative splicing under different conditions at an unparalleled depth (18, 22, 23). In this paper, we use a paired-end deep sequencing experiment to characterize, for the first time, the transcriptome and splicing of THEV (VAS vaccine strain) during different phases of the infection. Our paired-end sequencing allowed for reading 149 bp long high quality (mean Phred Score of 36) sequences from each end of cDNA fragments, which were mapped to the genome of THEV.

99 **RESULTS**

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

101 A previous study by Aboezz *et al* showed that almost all THEV transcripts were detectable beginning at 4  
102 hours, with a full replication cycle of about 18-hours post infection (hpi) (24). Therefore, infected MDTC-RP19  
103 cells were harvested at 4-, 12-, 24-, and 72-hpi to ensure a wide time window to sample all transcripts. Our  
104 paired-end RNA sequencing (RNA-seq) experiment yielded an average of 107.1 million total reads of 149 bp  
105 in length per time-point, which were simultaneously mapped to both the virus (THEV) and host (*Meleagris*  
106 *gallopavo*) genomes using the Hisat2 (25) reference-based aligner. A total of 18.1 million reads from all  
107 time-points mapped to the virus genome; this provided good coverage/depth, leaving no regions unmapped.  
108 The mapped reads to the virus genome increased substantially from a mere 432 reads at 4 hpi to 16.9  
109 million reads at 72 hpi (**Table 1**, **Figure 2A**). From the mapped reads, we identified a total of 2,457 unique  
110 THEV splice junctions from all time-points, with splice junctions from the later time-points being supported by  
111 significantly more sequence reads than earlier time-points. For example, all the 13 unique junctions at 4 hpi  
112 had less than 10 reads supporting each one, averaging a mere 2.8 reads per junction. Conversely, the 2374  
113 unique junctions at 72 hpi averaged 898.4 reads per junction, some junctions having coverage as high as  
114 322,677 reads. The substantial increases in splice junction and mapping reads to the THEV genome over  
115 time denotes an active infection and a successful viral replication, and correlates with our quantitative PCR  
116 (qPCR) assay quantifying the total number of viral genome copies over time (**Figure 2B**).

117 Using StringTie (25), we assembled the data into potential transcripts using the genomic location of the  
118 prior predicted THEV ORFs as a guide. In the consolidated transcriptome, a composite of all non-redundant  
119 transcripts from all time points, we counted a total of 29 novel transcripts. We found that a subset of exons  
120 in the viral transcripts match the predicted ORFs exactly, with the majority of the exons being longer and  
121 spanning multiple predicted ORFs (**Figure 3**).

122 We then validated the splice junctions in all transcripts by PCR amplification of viral cDNA, cloning, and Sanger  
123 sequencing (**Supplementary PCR methods**). During validation, we identified five additional transcripts,  
124 some of which were further validated by 3' Rapid Amplification of cDNA Ends (3' RACE) data. The complete  
125 list of unique splice junctions mapped to the THEV genome has been submitted to the National Center for  
126 Biotechnology Information Gene Expression Omnibus (<http://www.ncbi.nlm.nih.gov/geo>) under accession  
127 number GSE254416.

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

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

130 few pre-mRNAs that undergo alternative splicing to yield a repertoire of mature mRNAs. To evaluate the  
131 activity of each promoter over time, StringTie and Ballgown (a program for statistical analysis of assembled  
132 transcriptomes) (26) were used to estimate the normalized expression levels of all transcripts for each time  
133 point in Fragments Per Kilobase of transcript per Million mapped reads (FPKM) units. Very few unique splice  
134 junctions, reads, and transcripts were counted at 4 hpi; hence, this time point was excluded in this analysis.

135 Considering individual mRNAs, TRXPT\_21 – from the E2 region – was the most significantly expressed at  
136 12 hpi, constituting 33.58% of the total expression of all transcripts. Transcripts in the E3 and E4 regions also  
137 contributed significant proportions, and noticeably, some MLP region transcripts. The later time points were  
138 dominated by the MLP region transcripts – TRXPT\_10 and TRXPT\_14 were the most abundantly expressed  
139 at 24 and 72 hpi, respectively, as expected (**Figure 4A**). When we performed analysis of the FPKM values  
140 of transcripts per region we found a similar pattern: the E2 region was the most abundantly expressed at  
141 12 hpi, after which the MLP region assumes predominance (**Figure 4B**). Secondly, we estimated relative  
142 abundances of all splice junctions at each time point using the raw reads. Only junctions with a read coverage  
143 of at least 1% of the total splice junction reads at the given time point were considered significant and included  
144 in **Tables 2a-2c**. At 12 hpi, 18 junctions meet the 1% threshold, and were comprised of predominantly early  
145 region (E1, E2, E3, and E4) junctions, albeit the MLTU was the single most preponderant region overall,  
146 constituting 38.8% of all the junction reads (**Table 2a** and **Supplementary Table S1a**). The most abundant  
147 junctions at 12 hpi remained the most significantly expressed at 24 hpi also. However, here, the MLP-derived  
148 junctions were unsurprisingly even more preponderant overall, accounting for 45.7% of all the junction reads  
149 counted (**Table 2b** and **Supplementary Table S1b**). At 72 hpi, the trend of increased activity of the MLP  
150 continued as expected; at this time, the MLP region junctions were not only the most abundant overall –  
151 accounting for 67.3% of all junction reads, – but also contained the most significantly expressed individual  
152 junctions (**Table 2c**, **Supplementary Table S1c** and **Figure 4C**). When we limited this analysis to only  
153 junctions in the final transcriptome, we observed the relative abundances of the junctions for each region  
154 over time to be similar to the pattern seen with all the junctions included (**Figure 4D**).

155 We also analyzed splice donor and acceptor site nucleotide usage over time to investigate any peculiarities  
156 that THEV may show, generally or over the course of the infection. We found that most splice donor-acceptor  
157 sequences were unsurprisingly the canonical GU-AG nucleotides. However, the splice acceptor-donor pairing  
158 became less specific over time, such that all combinations of nucleotide pairs were eventually detected  
159 (**Figure 5**).

## 160 **Early Region 1 (E1) transcripts**

161 This region in MAdVs is the first transcribed after successful entry of the viral DNA into the host cell nucleus,

162 albeit at low levels (18). The host transcription machinery solely mediates the transcription of this region.  
163 After their translation, the E1 proteins in concert with a myriad of host transcription factors activate other viral  
164 promoters (10). In MAdVs, this region is subdivided into E1a and E2b units, but the transcripts found in our  
165 data categorized under this region do not appear to be subdivided.

166 Only two ORFs (ORF1 [sialidase] and Hyd) are predicted in this region; however, we discovered four novel  
167 transcripts in this region, which collectively contain 3 unique splice junctions (**Figure 6**). Most of the ORFs of  
168 the novel transcripts are distinct from the predicted ORFs, but they all have the coding potential (CP) for the  
169 predicted Hyd protein as the 3'-most coding sequence (CDS) if secondary start codon usage is considered  
170 as reported for other AdVs (17, 18). The 5'-most CDS of TRXPT\_1 is multi-exonic, encoding a novel 17.9  
171 kilodalton (kDa), 160 residue [amino acids (aa)] protein (ORF9). From its 5'-most start codon (SC), TRXPT\_2  
172 encodes the largest protein in this region – a 64.3 kDa, 580 aa protein (ORF10) with the same SC as ORF9  
173 (position 211 bp). ORF10 spans almost the entire predicted ORF1 and Hyd, coming short in two regards: it  
174 is spliced from 1655 bp to 1964 bp (ORF1's C-terminus, including the stop codon), and its stop codon (STC;  
175 position 2312) is 13 bp short of the Hyd STC. However, it has an SC 102 bp upstream and in-frame with  
176 ORF1's predicted SC. Thus, ORF10 shares substantial protein sequence similarity with ORF1 but not with  
177 Hyd, as the SC of Hyd is not in-frame. Without its splice site removing the ORF1 STC, TRXPT\_2 would  
178 encode a longer variant of ORF1, starting from an upstream SC. TRXPT\_3 is almost identical to TRXPT\_1,  
179 except for the lack of TRXPT\_1's second exon. Our RNA-seq data show that all E1 transcripts share the  
180 same transcription termination site (TTS; at position 2325 bp). However, TRXPT\_3 and TRXPT\_4 seem to  
181 have transcription start sites (TSS) downstream of the TSS of TRXPT\_1 and TRXPT\_2 (E1 TSS; position:  
182 54 bp). Given that studies in MAdVs show that E1 mRNAs share not only a common TTS but also the TSS,  
183 and only differ from each other regarding the internal splicing (18), it is likely that TRXPT\_3 and TRXPT\_4  
184 are incomplete, and their actual TSS just like the TTS are identical for all E1 transcripts. Regardless of  
185 the TSS considered for TRXPT\_3, the CP remains unaffected. Its 5'-most CDS, beginning at 1965 bp and  
186 sharing the same STC as ORF9, produces a 13.1 kDa, 115 residue protein (ORF4). ORF4 was predicted in  
187 an earlier study (27) but was excluded in later studies (1, 12); however, our data suggest it is a genuinely  
188 expressed ORF. Unlike TRXPT\_3, the CP of TRXPT\_4 is affected by the TSS considered; if we consider its  
189 unmodified TSS, then its CP is identical to TRXPT\_3. However, if we assume that TRXPT\_4 uses the E1  
190 TSS, then the 5'-most CDS is a distinct, novel, multi-exonic 15.9 kDa, 143 aa protein (ORF11) with the same  
191 SC as ORF9 and ORF10 but with a unique STC.

192 The splice junctions of all transcripts in this region (except the junction for TRXPT\_4) were validated by  
193 cloning of viral cDNA and Sanger sequencing (**Supplementary PCR methods**). During the validation of

194 TRXPT\_2, we found ORF1 to be present on the agarose gel (an unspliced band size) and Sanger sequencing  
195 results showed it to be a transcribed mRNA (**Supplementary PCR methods**). This was corroborated by  
196 our 3' RACE experiment, which showed a transcript (TRXPT\_2B) spanning the entire ORF1 and Hyd ORFs  
197 without any splicing, with a poly-A tail immediately after the E1 TTS. The 5'-most CDS of this transcript  
198 (TRXPT\_2B) would encode ORF1. However, TRXPT\_2B has an upstream and in-frame SC to the predicted  
199 SC of ORF1, suggesting that the predicted ORF1 CDS is truncated – the actual ORF1 (eORF1) that is  
200 expressed shares the same SC as ORF10, but has a unique STC.

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

202 The E2 TU expressed on the anti-sense strand is subdivided into E2A and E2B and encodes three classical  
203 AdV proteins – pTP and Ad-pol (E2B proteins), and DBP (E2A protein) – essential for genome replication  
204 (17, 18). Unlike MAdV where two promoters (E2-early and E2-late) are known (17), we discovered only a  
205 single TSS (E2 TSS; 18,751 bp) from which both E2A and E2B transcription is initiated. However, similar to  
206 MAdVs, E2A and E2B transcripts have distinct TTSs, and the E2B transcripts share the TTS of the IVa2  
207 transcript of the IM region (17, 18) (**Figure 7**).

208 The E2A ORF, DBP is one of three THEV ORFs predicted to be spliced from two exons. The corresponding  
209 transcript (TRXPT\_21) found in our data matches this predicted splice junction precisely but with a non-coding  
210 additional exon at the 5'-end (E2-5'UTR) at position 18,684-18,751 bp. Thus, TRXPT\_21 is a three-exon  
211 transcript encoding DBP (380 residues, 43.3 kDa) precisely as predicted. This transcript (TRXPT\_21) was  
212 also corroborated in a 3' RACE experiment. Additionally, from the 3' RACE data, we found a splice variant of  
213 TRXPT\_21 which retains the second intron leading to a 2-exon transcript. This new transcript (TRXPT\_21B),  
214 albeit longer due to retaining the second intron and possessing a short 3' UTR, encodes a truncated isoform  
215 of DBP (tDBP) because the SC utilized by TRXPT\_21, is followed shortly by STCs in the retained intron. The  
216 SC 173 bp downstream of the DBP SC yields tDBP (a 346 residue, 39.3 kDa product), which is in-frame  
217 of DBP but entirely contained in the second exon. TRXPT\_21 and TRXPT\_21B share a common TTS but  
218 TRXPT\_21B as seen in our 3'-RACE data, extends 39 bp into an adenine/thymine (A/T)-rich sequence  
219 before the poly-A tail sequence occurs, suggesting this position (16,934 bp) as the true E2A TTS (**Figure 7**).

220 The E2B region transcripts also start with the E2-5'UTR but extend thousands of base pairs downstream to  
221 reach the TTS at 2334 bp in the IM region, which is immediately followed by an A/T-rich sequence (position  
222 2323-2339 bp) where polyadenylation probably occurs. Interestingly, the TTS of the E1 region (position  
223 2,325 bp) on the sense strand is also in the immediate vicinity of this A/T-rich sequence, which is almost  
224 palindromic; hence it likely serves as the polyadenylation signal for both E1 and E2B/IM transcripts. The E2B  
225 transcripts TRXPT\_6 and TRXPT\_7 are almost identical except for an extra splice junction at the 3'-end of

226 TRXPT\_6, making TRXPT\_6 a five-exon transcript and TRXPT\_7, four exons (**Figure 7**). TRXPT\_7 has the  
227 CP for both classical proteins (pTP and Ad-pol) encoded in this region, of which the pTP ORF is predicted to  
228 be spliced from two exons just like in all other AdVs. The predicted splice junction of pTP is corroborated  
229 by our data; however, the full transcript is markedly longer than the predicted ORF: there are two novel  
230 non-coding 5' exons, the third exon (containing the SC of pTP) is significantly longer than predicted, and the  
231 last exon containing the bulk of the CDS is more than triple the predicted size of pTP. The first two exons are  
232 5'-UTRs because the SC here is immediately followed by STCs; thus, the 5'-most SC (position 10,995 bp) of  
233 the third exon which matches the predicted SC of pTP is utilized. The encoded product is identical to the  
234 predicted pTP protein (597 residues; 70.5 kDa). If secondary SC (secSC) usage is considered, with SC at  
235 6768 bp and STC at 3430 bp, the encoded product is identical to the predicted Ad-pol (polymerase) protein  
236 (1112 residues; 129.2 kDa). TRXPT\_6 differs from TRXPT\_7 by containing an extra splice site at 3447-3515  
237 bp. However, the CP remains similar to that of TRXPT\_7 except the Ad-pol encoded from the secSC is a  
238 truncated isoform with a new STC resulting from the splice site.

239 While both TRXPT\_6 and TRXPT\_7 have the CP for Ad-pol with secSC usage, in all AdVs studied, the two  
240 proteins (pTP and Ad-pol) are encoded by separate mRNAs with identical first three 5' exons and TTS, but  
241 the splice junction to the terminal exons are different. We checked for a longer splice junction between the  
242 third and fourth (terminal) exons of TRXPT\_7 with our junction validation method (targeted PCR, cloning,  
243 and Sanger sequencing) and discovered a unique splice junction (10,981-7062 bp) not found in our RNA-seq  
244 data. If initiated from the E2 TSS and terminated at the E2 TTS, this transcript (TRXPT\_31) would encode  
245 Ad-pol exactly as predicted as its 5'-most CDS (**Figure 7**).

246 Our RNA-seq data also showed a novel short transcript (TRXPT\_15) entirely nested within the terminal  
247 exon of TRXPT\_7 but with a unique splice site. This transcript is an incomplete construction from the  
248 mapped reads as it contains a truncated CDS. However, we validated this splice junction to be genuine  
249 (**Supplementary PCR methods**).

250 The IM region is a single-transcript TU, encoding a single classical protein, IVa2. The promoter expressing  
251 this single transcript (TRXPT\_5) is embedded in E2B region and shares a TTS with E2B transcripts (17, 18).  
252 TRXPT\_5 is a two-exon transcript spliced exactly as the last splice junction of TRXPT\_6. The first exon is a  
253 UTR, except the last 2 nucleotides, which connect with the first nucleotide of the second exon to form the  
254 5'-most SC. This first SC is 4 codons upstream and in-frame of the predicted IVa2 SC. Except for the four  
255 extra N-terminus residues, the entire protein sequence is identical to the predicted IVa2.

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

257 The E3 region is wholly contained in the MLTU and encodes proteins involved in modulating and evading  
258 the host immune defenses. In MAdVs, this region contains seven ORFs expressed from several transcripts  
259 which share the same TSS (from the E3 promoter) but have different TTSs (10, 17, 18). However, some E3  
260 transcripts use the TSS of the MLP. Due to sharing the same TSS, in MAdVs, secSC usage is heavily relied  
261 on for gene expression in this region except for 12.5K and transcripts using the MLP TSS, as utilizing only  
262 the first SC cannot produce all the other transcripts in this TU (17).

263 In THEV, only one ORF (E3) was predicted in this region. However, as the E3 TU is nested in the MLTU,  
264 transcripts from the L4P (100K, 22K, 33K, and pVIII) not only overlap the E3 region transcripts entirely as  
265 seen in our RNA-seq results, but also have their TSS and TTS in practically the same locations (**Figure 8**).  
266 Therefore, we have categorized these two groups together as E3 transcripts.

267 We identified seven novel transcripts here (TRXPT\_22, TRXPT\_23, TRXPT\_24, TRXPT\_25, TRXPT\_26,  
268 TRXPT\_27, TRXPT\_29) from our RNA-seq data, all originating from two distinct TSSs – we consider the first  
269 TSS (position 18,230 bp) as corresponding to the L4P and the other at 18,727 bp as corresponding to the  
270 E3 promoter (E3P). These E3 transcripts collectively have the CP for several predicted THEV ORFs: 100K,  
271 22K, 33K, pVIII, and E3, as well as Fiber (IV) and ORF7 belonging to the MLTU. But some of these CDSs  
272 are different than predicted due to either unknown exons or the presence of an in-frame upstream SC. For  
273 instance, 33K is one of the few THEV ORFs predicted to be spliced from two exons; however, we discovered  
274 a significantly longer four-exon ORF (e33K) on TRXPT\_24 that contains it almost entirely. The first two exons  
275 of e33K were not predicted but the last two match the predicted exons and the CDS is in-frame, albeit the  
276 first 20 bp of the predicted 33K (including the SC at 20,142 bp) is spliced out as part of the second intron of  
277 TRXPT\_24. Thus, the bona fide 33K (e33K) is a 19.8 kDa, 171 residue protein spanning four exons instead  
278 of the predicted 120 aa protein. TRXPT\_24 also has the CP for pVIII and E3 if we consider downstream SC  
279 usage. However, the predicted E3 has an upstream in-frame SC; thus, this longer version of E3 (eE3) is likely  
280 the genuinely expressed ORF. TRXPT\_29 is the shortest transcript in this TU. It is a two-exon transcript,  
281 both exons comprising the CDS. The product of TRXPT\_29 is a novel 73 residue protein (8.3KI) sharing  
282 the SC of e33K but with a unique STC. TRXPT\_23 being spliced identically as TRXPT\_29 also encodes  
283 8.3KI from its first SC. Similarly, TRXPT\_22 also encodes a 73 aa novel protein (8.3KII) from its first SC that  
284 shares over 80% similarity with 8.3KI, but it differs from 8.3KI at the C-terminus. Considering downstream  
285 SC usage, both TRXPT\_22 and TRXPT\_23 can encode pVIII and eE3 in that order, but TRXPT\_23 being  
286 longer, has the CP for the Fiber ORF also.

287 As the splice junctions of TRXPT\_22, TRXPT\_23, TRXPT\_24, and TRXPT\_29 essentially share the same  
288 genomic space, their validation was done with a single primer pair, and they were differentiated from each

other by cloning and Sanger sequencing (**Supplementary PCR methods**). In addition to corroborating the splice junctions for the aforementioned transcripts, the Sanger sequencing results also showed another splice variant undetected in our RNA-seq transcriptome. This was a three-exon transcript (TRXPT\_30) with its first and last exons spliced identically as TRXPT\_23, but which also has the second exon of TRXPT\_24 (**Figure 8**). The first CDS on TRXPT\_30 spans all three exons, producing a novel 140 residue, 15.7kDa protein. Interestingly, the last 81 C-terminus residues of this new protein (e22K) are identical to 22K (89 residues), which is a single-exon ORF predicted to use the same SC as 33K (20,142 bp). Just as seen for 33K, all the transcripts in this region exclude the first 20 bp of 22K (including the SC) as part of their introns; therefore, the first 7 residues of 22K are lacking in e22K due to splicing. Hence, we consider e22K as a long variant of the predicted 22K ORF. Albeit the TSS and TTS of TRXPT\_30 was not seen, we presume that they are similar to TRXPT\_23, in which case it would also have the downstream CP of TRXPT\_23.

TRXPT\_25 is the largest transcript in the TU. It also utilizes the L4P TSS but has a distinct TTS. It is a two-exon transcript, encoding a novel protein (t100K; 543 residues), which is a shorter isoform of the predicted 100K ORF. Considering secSC usage on this transcript yields the predicted 22K ORF precisely as predicted. It also has the CP for pVIII and eE3 in that order. Furthermore, during the validation of the TRXPT\_25 splice junction using primers that span its junction (18,350-18,717 bp), we noticed a DNA band that corresponds to the full unspliced sequence (**Supplementary PCR methods**). As TRXPT\_25 only falls short of encoding the complete predicted 100K protein due to its splice junction, this band (which we cloned and validated by Sanger sequencing) suggests that the predicted 100K is indeed expressed. This transcript (TRXPT\_25B) although not seen in full, likely shares the same TSS and TTS as TRXPT\_25. Lastly, TRXPT\_26 and TRXPT\_27 both originate from the E3P but have distinct TTSs. TRXPT\_26 is a three-exon transcript but the first two are UTRs. It encodes pVIII as the 5'-most ORF and has the CP for eE3 and Fiber in that order. TRXPT\_27 on the other hand, is only a two-exon transcript that is similar to TRXPT\_26, only the terminal exon contains the CDSs. It encodes Fiber as the 5'-most ORF, and ORF7 downstream with secSC usage. TRXPT\_13, which is an L4 transcript that uses the MLP TSS is discussed under the MLTU transcripts.

### Early Region 4 (E4) transcripts

This TU is found at the tail-end (3'-end) of the genome and expressed from the anti-sense strand. Based on nucleotide position, ORF7 and ORF8 were predicted in this region (1); however, as ORF7 is neither on the same strand as ORF8 nor transcribed from a promoter in the E4 region, only ORF8 can legitimately be classified as a transcript in this TU. This is corroborated by our RNA-seq data, as only one transcript was identified in this region on the anti-sense strand (**Figure 9**). The transcript (TRXPT\_28) spans 25192-26247

321 bp and is spliced at 25701-26055 bp, making a two-exon transcript. The second exon fully matches the  
322 predicted ORF8 with 12 extra base pairs at the 3'-end. However, there is an SC in the first exon at position  
323 26246 bp (192 bp upstream of the predicted SC). The encoded protein from this SC is in-frame with the  
324 predicted SC found in the second exon; hence, we consider this protein (eORF8 – 26.4 kDa, 229 aa), a  
325 longer isoform of the predicted ORF8, as the genuinely expressed ORF with an identical C-terminus to the  
326 predicted ORF8 protein.

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

328 The MLTU transcripts dominate the late phase (i.e, after DNA replication) of the AdV infectious cy-  
329 cle. The MLP produces all late mRNAs by alternative splicing and alternative polyadenylation of a pri-  
330 mary transcript, grouped into five transcript classes (L1-L5). Most of the coding capacity of THEV falls  
331 within this TU. Specifically, about 13 out of the 23 predicted ORFs were assigned to this TU, some  
332 of which we have categorized under the E3 TU instead. Our RNA-seq data revealed 12 transcripts  
333 (TRXPT\_8, TRXPT\_9, TRXPT\_10, TRXPT\_11, TRXPT\_12, TRXPT\_13, TRXPT\_14, TRXPT\_16, TRXPT\_17,  
334 TRXPT\_18, TRXPT\_19, TRXPT\_20) in this TU, the majority of which have the 5' untranslated TPL sequence  
335 as seen in all AdVs. For three transcripts (TRXPT\_16, TRXPT\_17, TRXPT\_18), a different leader sequence  
336 (sTPL) is used, which differs from the TPL in only one regard: the first TPL exon is substituted for a different  
337 first exon, found between the first and second TPL exons. Also, TRXPT\_20 seems to include only the third  
338 TPL exon (**Figure 10**).

339 We identified five TTSs (10,549 bp, 12,709 bp, 16,870 bp, 17,891 bp, 20,865 bp) in this TU, which we  
340 consider as corresponding to the five late mRNA classes (L1-L5), respectively, as found in all AdVs. L1  
341 mRNAs include TRXPT\_8, which comprises the TPL (non-coding) and the CDS-containing terminal exon.  
342 This transcript encodes the 52K ORF exactly as predicted with the SC beginning from the first nucleotide  
343 of the terminal exon. L2 mRNAs include TRXPT\_16, TRXPT\_17, and TRXPT\_18, all of which consist of  
344 the sTPL (also non-coding) followed by their respective terminal exons. TRXPT\_16 encodes pIIIa exactly  
345 as predicted as the 5'-most ORF, and also has the CP for the ORFs, III and pVII in that order. TRXPT\_17  
346 encodes the ORF, III (penton), and TRXPT\_18 encodes the ORF pVII exactly as predicted. The L3 mRNAs  
347 include TRXPT\_14 and TRXPT\_20, of which TRXPT\_14 utilizes the full TPL whereas TRXPT\_20 uses only  
348 the third TPL exon (TPL3). Both transcripts have the CP for the ORF, hexon (II) but hexon is the only ORF  
349 encoded on TRXPT\_14, whereas the 5'-most ORF on TRXPT\_20 is pX (pre-Mu) followed by pVI and hexon  
350 in that order. L4 mRNAs include TRXPT\_9, TRXPT\_10, TRXPT\_11, and TRXPT\_13 all of which begin  
351 with the TPL followed by three (TRXPT\_9, TRXPT\_10, and TRXPT\_13) or four (TRXPT\_11) coding exons.  
352 These are the largest transcripts found in the transcriptome, each one possessing the CP for several similar

late proteins. Normally, MLTU transcripts encoding particular ORFs splice the TPL onto a splice site just upstream of the ORF to be expressed (17). While this holds true for most MLTU ORFs, several late ORFs (pVI, protease, and ORF7) do not have such close proximity splicing but are contained in larger transcripts such as these L4 mRNAs, strongly suggesting the use of non-standard ribosomal initiation mechanisms such as secSC utility and ribosome shunting found in other AdVs for their translation (17, 28). TRXPT\_9 and TRXPT\_10 are very similar but not identical. The last exon of TRXPT\_9 seems to be truncated and probably shares the same TTS as the other L4 mRNAs. They are both 6-exon transcripts encoding pVII as the 5'-most ORF (fourth exon) and also have the CP for pX, pVI, hexon, a longer variant of protease (eProt) – uses an upstream in-frame SC than predicted, and ORF12 (a novel unpredicted 120 aa protein). TRXPT\_10 (and TRXPT\_9 with the L4 TTS) additionally has the CP for pVIII and eE3. Conversely, TRXPT\_11 is a seven-exon mRNA with hexon as its 5'-most ORF but it also has the CP for eProt, ORF12, e33K, and also pVIII and eE3 in that order. TRXPT\_13 seems to be an E3 ORF utilizing the MLP TSS as it encodes classical L4P genes such as pVIII and eE3 in that order similar to TRXPT\_22 (E3 TU) but it lacks the novel first ORF (8.3KII) of TRXPT\_22.

Lastly, the L5 class includes only TRXPT\_12 which contains the TPL and a coding terminal exon. Its 5'-most ORF is fiber (IV) but it also has the CP for the THEV specific gene, ORF7. TRXPT\_12's CP is identical to TRXPT\_27 of the E3 TU but they differ in their 5'-UTRs.

370 **DISCUSSION/CONCLUSIONS**

371 While the advent of next-generation sequencing has rendered easier the study of large and complex eu-  
372 karyotic transcriptomes, the study of the smaller and compact viral transcriptomes remains unintuitively  
373 challenging, as several transcripts may have significant overlaps due to genome economization. Characteriz-  
374 ing AdV transcriptomes is even more difficult due to the wide array of mRNAs produced via very complex  
375 alternative splicing combined with alternative polyadenylation, all initiated from relatively few promoters. This  
376 makes AdV transcriptomes some of the most intricate for a virus. The challenge is further compounded  
377 by the fact that the standard software programs used in the RNA-seq analysis pipelines are not designed  
378 primarily for such compact, gene-dense, and complex transcriptomes as AdVs. Furthermore, in our case,  
379 there is no prior transcriptomic studies for THEV. Our approach to properly handle this complex data was to  
380 use standard RNA-seq analysis programs coupled with some custom analyses and experimentally validating  
381 all splice junctions with independent methods. Our work provides the first insights into the splicing patterns  
382 of THEV, which is expectedly similar to other MAdVs but with key differences. Our work shows 34 transcripts  
383 in the THEV transcriptome grouped into five TUs, of which the E3 TU shows great complexity of alternative  
384 splicing.

385 An unexpected observation is that the pileup of mapped reads to THEV seems consistently skewed over  
386 similar regions of the genome at all time points. As the gene expression of AdVs is temporally regulated,  
387 we expected to see unambiguous differences in the pileup of reads over different regions of the genome  
388 at different time points, indicating the different stages of infection. While this could simply mean that the  
389 infection was not well synchronized, we speculate that the temporal gene expression regulation of THEV is  
390 probably different from MAdVs. This is supported by a previous study stating the same conclusion with its  
391 finding that almost all THEV transcripts were detectable by at 4h.p.i, and by 8h.p.i, mRNA for all predicted  
392 ORFs (including the late genes) were present (24). Conversely, despite the overall pileup similarity, a close  
393 inspection shows that the relative proportions of reads over some regions show some variation over time.  
394 The breakdown of transcripts detected at different time points in **Figure 3b** seems to support this different  
395 temporal regulation of THEV. Specifically, the MLP of THEV is active significantly earlier in infection – as  
396 early as 4h.p.i and more pronounced at 12h.p.i (**Figure 3b** and **Table 2a**), – whereas the late phase shift in  
397 MAdVs occurs after 24h.p.i. This also lends credence to our speculation. However, generally speaking, the  
398 overall temporal gene expression regulation known in MAdVs – early regions showing their peak expression  
399 at earlier time points followed by predominance of the MLTU at later time points – also holds true for THEV.  
400 Further studies would be necessary to establish the precise temporal regulation of THEV transcription.

401 The use of short read deep sequencing to reconstruct full AdV mRNA structures provides excellent results,

especially for mapping the splice sites. However, due to the substantial overlapping nature of AdV mRNAs coupled with the fragmentation step in the library preparation protocol, mapping the precise TSS and TTS of the assembled transcripts is difficult. Also, similar transcripts with substantial overlaps may be assembled as one longer mRNA since the short reads alone do not provide enough context for the transcript assembler (StringTie) to distinguish them. In our results, we see transcripts in the same TU initiated or terminated in the same approximate area (10-70 bp and 1-300 bp apart for TSS and TTS, respectively) but not precisely at the same position. We consider the most upstream TSS or most downstream TTS for the transcripts involved but we present them unchanged in all the figures shown. Also, by comparison to the more well-studied MAdV transcriptomes, we believe that a few long transcripts in the MLTU (TRXPT\_9, TRXPT\_10, and TRXPT\_11) are probably a result of fusing some L4P-derived transcripts to the terminal exons of the bona fide MLTU transcripts by StringTie, making them significantly longer. These mRNAs do not only have unusually many exons for an AdV, but their last three or four exons are also identical to the L4P-derived mRNAs. Future studies using long read sequencing technologies are necessary to provide conclusive data for precisely mapping the TSS and TTS, as well as teasing apart the bona fide structures of the long MLTU transcripts. Furthermore, it is not unreasonable to presume that several splice variants were undiscovered in our work as evidenced firstly by finding unique transcripts using 3' RACE and during our splice junction validation steps. And secondly, recent studies (17, 18, 22) are still discovering novel mRNA variants for even the best studied MAdVs decades later. Another observation made is that all the TTSs in THEV's transcriptome are in close proximity to A/T-rich sequences which we presume to be polyadenylation signal sequences (PASS). Interestingly, some of these PASSs are located in the immediate vicinity of two closely located TTSs expressed on opposite strands. Namely, the E1 and E2B/IM TTSs have an almost palindromic PASS between them, as do the E4 (anti-sense strand) and the sense strand TRXPT\_12 and TRXPT\_27.

An interesting finding of our analysis is that while most of the predicted ORFs are precisely encoded by the spliced transcripts, we found a few that seem to be truncated predictions, as either an upstream in-frame SC (eORF1, eE3, and eProt) or unknown upstream exons spliced onto them (eIVa2, e33K, and eORF8) were found. Other ORFs were identified that were either shorter (tDBP, t100K) or longer (e22K) isoforms of some predicted ORF but we found evidence to support the predicted ORF itself, making them all possible genuinely translated variants. We also found several novel unpredicted ORFs. Taken together, we surmise that further studies will likely yield even more unpredicted novel ORFs or variants of predicted ORFs.

Eukaryotic mRNAs are typically functionally monocistronic, the 5'-most AUG normally being used to determine the translation reading frame. However, depending on the sequence context, in some organisms, the initiating codon may even be a non-AUG start codon. AdV mRNAs, which mostly span more than one ORF, are known

434 to be functionally polycistronic, employing non-standard mechanisms of translation initiation, namely, secSC  
435 usage and ribosome shunting (10, 22). While there is no reliable method of predicting how efficiently any  
436 given AUG will be used, AdVs use secondary AUGs as initiation codons for most E1b proteins and for some  
437 E3 proteins. In fact, recent studies show that secSC usage is found transcriptome-wide. This is thought to  
438 occur because translation initiation at the first SC is inefficient, allowing downstream SCs to be employed  
439 for initiation (17). The ribosomal shunting or jumping mechanism is utilized for MLTU transcripts that have  
440 the TPL. This mechanism allows the ribosome to translocate to a downstream initiating codon under the  
441 direction of the shunting elements in the TPL, even if a start codon in a good Kozak sequence context is  
442 bypassed. Thus, predicting the protein(s) that are expressed from an AdV mRNA becomes highly uncertain  
443 as any one of the SC may be selected (10, 22). Almost all the THEV transcripts in our data have the CP for  
444 several ORFs, some spanning as many as six ORFs but the majority spanning at least two ORFs. Therefore,  
445 we believe our data support the usage of these special ribosome initiation mechanisms as several predicted  
446 and novel ORFs found on mRNA in our data have no conceivable mechanism of being translated if only the  
447 typical ribosome scanning mechanism is employed. Interestingly, several distinct transcripts have identical  
448 CPs. This is not unique to THEV but is observed in human AdVs in a recent study (17). They proposed that  
449 this may permit protein production to be fine-tuned through alteration in the balance between different mRNA  
450 groups expressing that ORF.

451 It is well established that AdV alternative splicing undergoes a regulated temporal shift in splice site usage.  
452 This was thought to be limited to certain TUs; however, recent studies suggest that AdVs routinely produce  
453 different combinations of splice acceptor–donor pairs and that this is observed in all TUs (10, 17, 22, 29). The  
454 mechanistic details of this phenomenon have been best studied for the E1A and L1 units. The studies show  
455 that AdVs (specifically, late phase AdV-infected nuclear extract) modulate the activities of the splicing factor  
456 U2AF and the cellular SR family of splicing factors (reviewed in reference (29)) and encode several mostly  
457 late phase proteins (E4-ORF3, E4-ORF6, E4-ORF4, L4-33K, and L4-22K) that influence the RNA splice  
458 site used. This phenomenon seems to occur in the THEV transcriptome also, as the stringency of splice  
459 acceptor-donor pairs selected decreases measurably from the onset of the late phase (see **Figure 5**). In  
460 fact, recent studies of some human AdVs show that a virtually unlimited number of combinatorial alternative  
461 splicing events occur in an AdV lytic infection, resulting in menagerie of novel transcripts (17, 22). It is unlikely  
462 that the entire repertoire of mRNA produced via this mechanism will actually be translated. However, it has  
463 been speculated that the plasticity in alternative RNA splicing enables AdVs to fine-tune protein synthesis  
464 by providing different alternatively spliced variants encoding the same protein under changing conditions.  
465 The capacity to produce novel exon combinations offers the virus an evolutionary advantage to adjust the

466 repertoire of mRNA transcription and protein production in the changing environment of the viral replication  
467 cycle (17, 22).

468 Summarizing all the main points above, we see that the THEV transcriptome bears remarkable overall  
469 similarity to the better studied MAdVs. The transcriptome is organized into five TUs, the temporal regulation  
470 is divided into early and late genes, and a broad repertoire of transcripts are produced via virtually unlimited  
471 alternative splicing. However, the THEV transcriptome appears to be less sophisticated (i.e, it encodes  
472 fewer genes) than MAdVs primarily because the MAdV genomes are close to twice as long as that of THEV,  
473 which rationally should encode less genes. The lack of subdivision of the E1 region into E1a and E1b is  
474 one of the most obvious examples. Also, the MAdV E4 region encodes several proteins unlike in THEV  
475 where only one transcript coding for one protein was found. The most conspicuous example is found in  
476 examining the complexity of the MLTU leader sequences. While the majority of the THEV MLTU transcripts  
477 begin with the TPL (267 bp long) just like MAdVs, and also utilize a variant leader sequence (sTPL), it is  
478 well-established that significantly more diverse 5'UTRs are employed for MAdV MLTU transcripts. Namely,  
479 the TPL (used for majority of transcripts), the so-called x, y, and z leaders, and the i-leader are 5' leaders  
480 utilized by MAdV MLTU mRNAs. Granted, the MAdV MLTU transcripts infrequently incorporate the non-TPL  
481 leaders, their absence in our data could mean that the 5'UTR diversity of THEV's MLTU mRNA are indeed  
482 more limited due to its smaller genome size. It is also possible that later studies could uncover more variety  
483 not seen in our results. Specifically, using long read sequencing technologies would be a excellent approach  
484 to discovering and fine-tuning the transcriptome characterized in this work. Also, it is likely that future studies  
485 repeating this work in other cell lines will yield some interesting insights.

486 **MATERIALS AND METHODS**

487 **Cell culture and THEV Infection**

488 The Turkey B-cell line (MDTC-RP19, ATCC CRL-8135) was grown as suspension cultures in 1:1 complete  
489 Leibovitz's L-15/McCoy's 5A medium with 10% fetal bovine serum (FBS), 20% chicken serum (ChS), 5%  
490 tryptose phosphate broth (TPB), and 1% antibiotic solution (100 U/mL Penicillin and 100 $\mu$ g/mL Streptomycin),  
491 at 41°C in a humidified atmosphere with 5% CO<sub>2</sub>. Infected cells were maintained in 1:1 serum-reduced  
492 Leibovitz's L15/McCoy's 5A media (SRLM) with 2.5% FBS, 5% ChS, 1.2% TPB, and 1% antibiotic solution. A  
493 commercially available THEV vaccine was purchased from Hygieia Biological Labs as a source of THEV-A  
494 (VAS strain). The stock virus was titrated using an in-house qPCR assay with titer expressed as genome  
495 copy number (GCN)/mL, similar to Mahshoub *et al* (30) with modifications. Cells were infected in triplicate  
496 at a multiplicity of infection (MOI) of 100 GCN/cell, incubated at 41°C for 1 hour, and washed three times  
497 with phosphate buffered saline (PBS) to get rid of free virus particles. Triplicate samples were harvested at  
498 4-, 12-, 24-, and 72-hpi for total RNA extraction. The infection was repeated but samples in triplicate were  
499 harvested at 12-, 24-, 36-, 48-, and 72-hpi for PCR validation of novel splice sites. Still one more independent  
500 infection was done at time points ranging from 12 to 168-hpi for qPCR quantification of virus titers.

501 **RNA extraction and Sequencing**

502 Total RNA was extracted from infected cells using the Thermo Fisher RNAqueous™-4PCR Total RNA Isolation  
503 Kit (which includes a DNase I digestion step) per manufacturer's instructions. An agarose gel electrophoresis  
504 was performed to check RNA integrity. The RNA quantity and purity was initially assessed using nanodrop,  
505 and RNA was used only if the A260/A280 ratio was 2.0 ± 0.05 and the A260/A230 ratio was >2 and <2.2.  
506 Extracted total RNA samples were sent to LC Sciences, Houston TX for poly-A-tailed mRNA sequencing  
507 where RNA integrity was checked with Agilent Technologies 2100 Bioanalyzer High Sensitivity DNA Chip  
508 and poly(A) RNA-seq library was prepared following Illumina's TruSeq-stranded-mRNA sample preparation  
509 protocol. Paired-end sequencing to generate 150 bp reads was performed on the Illumina NovaSeq 6000  
510 sequencing system.

511 **Validation of Novel Splice Junctions**

512 All splice junctions identified in this work are novel except one predicted splice site each for pTP, DBP,  
513 and 33K, which were corroborated in our work. However, these predicted splice junctions had not been  
514 experimentally validated hitherto, and we identified additional novel exons, giving the complete picture of

515 these transcripts. The novel splice junctions discovered in this work using the StringTie transcript assembler  
516 were validated by PCR, cloning, and Sanger Sequencing (**Supplementary PCR methods**). Briefly, primers  
517 spanning a range of novel exon-exon boundaries for each specific transcript in a transcription unit (TU) were  
518 designed. Universal forward or reverse primers for each respective TU were designed and paired with primers  
519 binding specific positions in each transcript. Each forward primer contained a KpnI restriction site and each  
520 reverse primer, an XbaI site in the primer 5' ends. After first-strand cDNA synthesis of total RNA obtained  
521 from THEV infected MDTC-RP19 cells was done using SuperScript™ IV First-Strand Synthesis System, the  
522 primers were used in a targeted PCR amplification, the products analyzed with agarose gel electrophoresis  
523 to confirm expected band sizes, cloned by traditional restriction enzyme method, and Sanger sequenced  
524 to validate these splice junctions at the sequence level. The total RNA was extracted as described above,  
525 including the DNase I digestion step. We included infected total RNA controls with no reverse transcriptase  
526 (no RT) during the cDNA synthesis step and the parent RNA were digested using RNase H after cDNA  
527 synthesis was complete to ensure that the bands obtained from the targeted PCR amplifications did not  
528 originate from the viral genomic DNA. As seen in the agarose gel images in **Supplementary PCR methods**,  
529 DNA bands were not found in the “no RT” controls, indicating that the DNA bands seen are of cDNA origin.

### 530 **3' Rapid Amplification of cDNA Ends (3' RACE)**

531 A rapid amplification of sequences from the 3' ends of mRNAs (3' RACE) experiment was performed using  
532 a portion of the extracted total RNA of infected MDTC-RP19 cells used for the RNA-seq experiment as  
533 explained above. We followed the protocol described by Green *et al* (31) with modifications. Briefly, 1 $\mu$ g of  
534 total RNA was reverse transcribed to cDNA using SuperScript™ IV First-Strand Synthesis System following  
535 the manufacturing instructions using an adapter-primer with a 3'-end poly(T) and a 5'-end BamHI restriction  
536 site. A gene-specific sense primer with a 5'-end KpnI restriction site paired with an anti-sense adapter-primer  
537 with a 5'-end BamHI site were used to amplify target sections of the cDNA using Invitrogen's Platinum™ Taq  
538 DNA polymerase High Fidelity, following manufacturer's instructions. The PCR amplicons were restriction  
539 digested, cloned, and Sanger sequenced.

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

541 Sequencing reads were analyzed following a well-established protocol described by Pertea *et al* (25), using  
542 Snakemake - version 7.24.0 (32), a popular workflow management system to drive the pipeline. Briefly,  
543 sequencing reads were trimmed with the Trim-galore - version 0.6.6 (33) program to achieve an overall Mean  
544 Sequence Quality (Phred Score) of 36. Trimmed reads were mapped simultaneously to the complete genomic

545 sequence of avirulent turkey hemorrhagic enteritis virus (<https://www.ncbi.nlm.nih.gov/nuccore/AY849321.1/>)  
546 and *Meleagris gallopavo* (<https://www.ncbi.nlm.nih.gov/genome/?term=Meleagris+gallopavo>) using Hisat2  
547 - version 2.2.1 (25) with default settings. The generated binary alignment (BAM) files from each infection  
548 time point were filtered for reads mapping to the THEV genome using Samtools - version 1.16.1 and fed into  
549 StringTie - version 2.2.1 (25) to assemble the transcripts, using a gene transfer format (GTF) annotation  
550 file derived from a gene feature format 3 (GFF3) annotation file obtained from NCBI, which contains the  
551 predicted ORFs of THEV as a guide. GFFCOMPARE - version 0.12.6 was used to merge all transcripts  
552 from all time points without redundancy and using a custom R script, adenovirus transcripts units (regions)  
553 were assigned to each transcript, generating the transcriptome of THEV. StringTie was set to expression  
554 estimation mode to calculate FPKM scores for all transcripts after which Ballgown - version 2.33.0 in R was  
555 used to perform the statistical analysis on the transcript expression levels. Samtools was also used to count  
556 the total sequencing reads for all replicates at each time point and Regtools - version 1.0.0 was used to  
557 count all junctions, the reads supporting them, and extract all other information related to the junction. See  
558 **Supplementary Computational Analysis** for the details of transcript expression level estimations and splice  
559 junction read counts.

## 560 DATA AVAILABILITY

561 The raw sequencing read data (FastQ), transcript expression counts, and total unique junctions have  
562 been deposited at the National Center for Biotechnology Information Gene Expression Omnibus (<http://www.ncbi.nlm.nih.gov/geo>) under accession number GSE254416.  
563  
564 Data is also available on request by contacting the designated corresponding author.

## 565 CODE AVAILABILITY

566 All the code/scripts in the entire analysis pipeline are available on github ([https://github.com/Abraham-Quaye/thev\\_transcriptome](https://github.com/Abraham-Quaye/thev_transcriptome))  
567

## 568 ACKNOWLEDGMENTS

569 We thank the Office of Research Computing at Brigham Young University for granting us access to the high  
570 performance computing systems to perform the memory-intensive steps in the analysis pipeline of this work.

571 **REFERENCES**

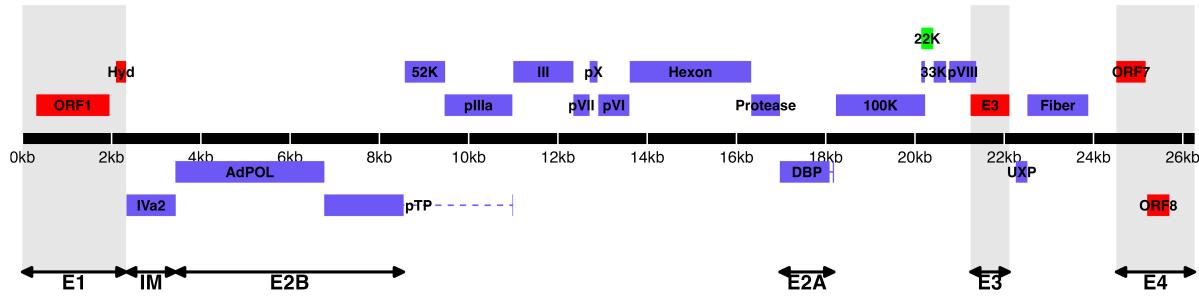
- 572 1. Davison A, Benko M, Harrach B. 2003. Genetic content and evolution of adenoviruses. *The Journal of general virology* 84:2895–908.
- 573 2. Harrach B. 2008. Adenoviruses: General features, p. 1–9. *In* Mahy, BWJ, Van Regenmortel, MHV (eds.), *Encyclopedia of virology* (third edition). Book Section. Academic Press, Oxford.
- 574 3. Upton C, Slack S, Hunter AL, Ehlers A, Roper RL. 2003. Poxvirus orthologous clusters: Toward defining the minimum essential poxvirus genome. *Journal of virology* 77:7590–7600.
- 575 4. McGeoch D, Davison AJ. 1999. Chapter 17 - the molecular evolutionary history of the herpesviruses, p. 441–465. *In* Domingo, E, Webster, R, Holland, J (eds.), *Origin and evolution of viruses*. Book Section. Academic Press, London.
- 576 5. Harrach B, Benko M, Both GW, Brown M, Davison AJ, Echavarría M, Hess M, Jones M, Kajon A, Lehmkühl HD, Mautner V, Mittal S, Wadell G. 2011. Family adenoviridae. *Virus Taxonomy: 9th Report of the International Committee on Taxonomy of Viruses* 125–141.
- 577 6. Kovács ER, Benkő M. 2011. Complete sequence of raptor adenovirus 1 confirms the characteristic genome organization of siadenoviruses. *Infection, Genetics and Evolution* 11:1058–1065.
- 578 7. Davison AJ, Wright KM, Harrach B. 2000. DNA sequence of frog adenovirus. *J Gen Virol* 81:2431–2439.
- 579 8. Kovács ER, Jánoska M, Dán Á, Harrach B, Benkő M. 2010. Recognition and partial genome characterization 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.
- 580 9. Katoh H, Ohya K, Kubo M, Murata K, Yanai T, Fukushi H. 2009. A novel budgerigar-adenovirus belonging to group II avian adenovirus of siadenovirus. *Virus Research* 144:294–297.

- 581 10. Guimet D, Hearing P. 2016. 3 - adenovirus replication, p. 59–84. In Curiel, DT (ed.), Adenoviral  
vectors for gene therapy (second edition). Book Section. Academic Press, San Diego.
- 582 11. Beach NM. 2006. Characterization of avirulent turkey hemorrhagic enteritis virus: A study of the  
molecular basis for variation in virulence and the occurrence of persistent infection. Thesis.
- 583 12. Beach NM, Duncan RB, Larsen CT, Meng XJ, Sriranganathan N, Pierson FW. 2009. Comparison of  
12 turkey hemorrhagic enteritis virus isolates allows prediction of genetic factors affecting virulence. J  
Gen Virol 90:1978–85.
- 584 13. Gross WB, Moore WE. 1967. Hemorrhagic enteritis of turkeys. Avian Dis 11:296–307.
- 585 14. Rautenschlein S, Sharma JM. 2000. Immunopathogenesis of haemorrhagic enteritis virus (HEV) in  
turkeys. Dev Comp Immunol 24:237–46.
- 586 15. Larsen CT, Domermuth CH, Sponenberg DP, Gross WB. 1985. Colibacillosis of turkeys exacerbated  
by hemorrhagic enteritis virus. Laboratory studies. Avian Dis 29:729–32.
- 587 16. Dhami K, Gowthaman V, Karthik K, Tiwari R, Sachan S, Kumar MA, Palanivelu M, Malik YS, Singh  
RK, Munir M. 2017. Haemorrhagic enteritis of turkeys – current knowledge. Veterinary Quarterly  
37:31–42.
- 588 17. Donovan-Banfield I, Turnell AS, Hiscox JA, Leppard KN, Matthews DA. 2020. Deep splicing plasticity  
of the human adenovirus type 5 transcriptome drives virus evolution. Communications Biology 3:124.
- 589 18. Zhao H, Chen M, Pettersson U. 2014. A new look at adenovirus splicing. Virology 456-457:329–341.
- 590 19. Wolfrum N, Greber UF. 2013. Adenovirus signalling in entry. Cell Microbiol 15:53–62.
- 591 20. Falvey E, Ziff E. 1983. Sequence arrangement and protein coding capacity of the adenovirus type 2  
"i" leader. Journal of Virology 45:185–191.

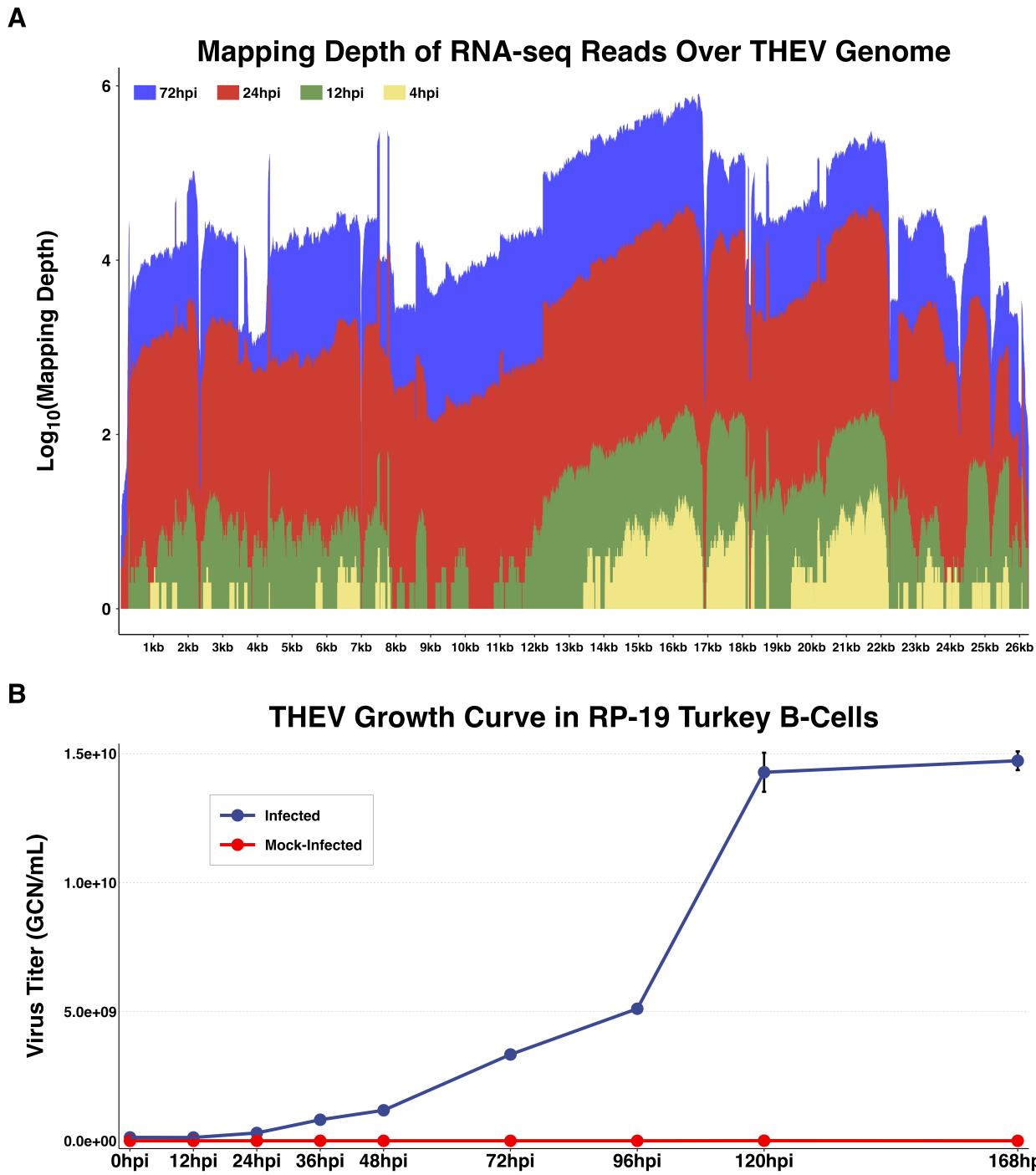
- 592 21. Morris SJ, Scott GE, Leppard KN. 2010. Adenovirus late-phase infection is controlled by a novel L4 promoter. *Journal of Virology* 84:7096–7104.
- 593 22. Westergren Jakobsson A, Segerman B, Wallerman O, Bergström Lind S, Zhao H, Rubin C-J, Petersson U, Akusjärvi G. 2021. The human adenovirus 2 transcriptome: An amazing complexity of alternatively spliced mRNAs. *Journal of Virology* 95.
- 594 23. 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 cells. *Nature* 489:101–108.
- 595 24. Aboezz Z, Mabsoub H, El-Bagoury G, Pierson F. 2019. In vitro growth kinetics and gene expression analysis of the turkey adenovirus 3, a siadenovirus. *Virus Research* 263:47–54.
- 596 25. Pertea M, Kim D, Pertea GM, Leek JT, Salzberg SL. 2016. Transcript-level expression analysis of RNA-seq experiments with HISAT, StringTie and ballgown. *Nature Protocols* 11:1650–1667.
- 597 26. Jack Fu [Aut], Alyssa C. Frazee [Aut, Cre], Leonardo Collado-Torres [Aut], Andrew E. Jaffe [Aut], Jeffrey T. Leek [Aut, Ths]. 2017. Ballgown. Bioconductor.
- 598 27. 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, hemorrhagic enteritis virus. *Virology* 249:307–315.

- 599 28. Yueh A, Schneider RJ. 1996. Selective translation initiation by ribosome jumping in adenovirus-infected  
and heat-shocked cells. *Genes & Development* 10:1557–1567.
- 600 29. Akusjarvi G. 2008. Temporal regulation of adenovirus major late alternative RNA splicing. *Frontiers in  
Bioscience Volume:5006.*
- 601 30. Mabsoub HM, Evans NP, Beach NM, Yuan L, Zimmerman K, Pierson FW. 2017. Real-time PCR-based  
infectivity assay for the titration of turkey hemorrhagic enteritis virus, an adenovirus, in live vaccines.  
*Journal of Virological Methods* 239:42–49.
- 602 31. Green MR, Sambrook J. 2019. Rapid amplification of sequences from the 3' ends of mRNAs: 3'-RACE.  
*Cold Spring Harbor Protocols* 2019:pdb.prot095216.
- 603 32. Mölder F, Jablonski KP, Letcher B, Hall MB, Tomkins-Tinch CH, Sochat V, Forster J, Lee S, Twardziok  
SO, Kanitz A, Wilm A, Holtgrewe M, Rahmann S, Nahnsen S, Köster J. 2021. Sustainable data  
analysis with snakemake. *F1000Research* 10:33.
- 604 33. Krueger F, James F, Ewels P, Afyounian E, Weinstein M, Schuster-Boeckler B, Hulselmans G,  
Sclamons. 2023. FelixKrueger/TrimGalore: v0.6.10 - add default decompression path. Zenodo.

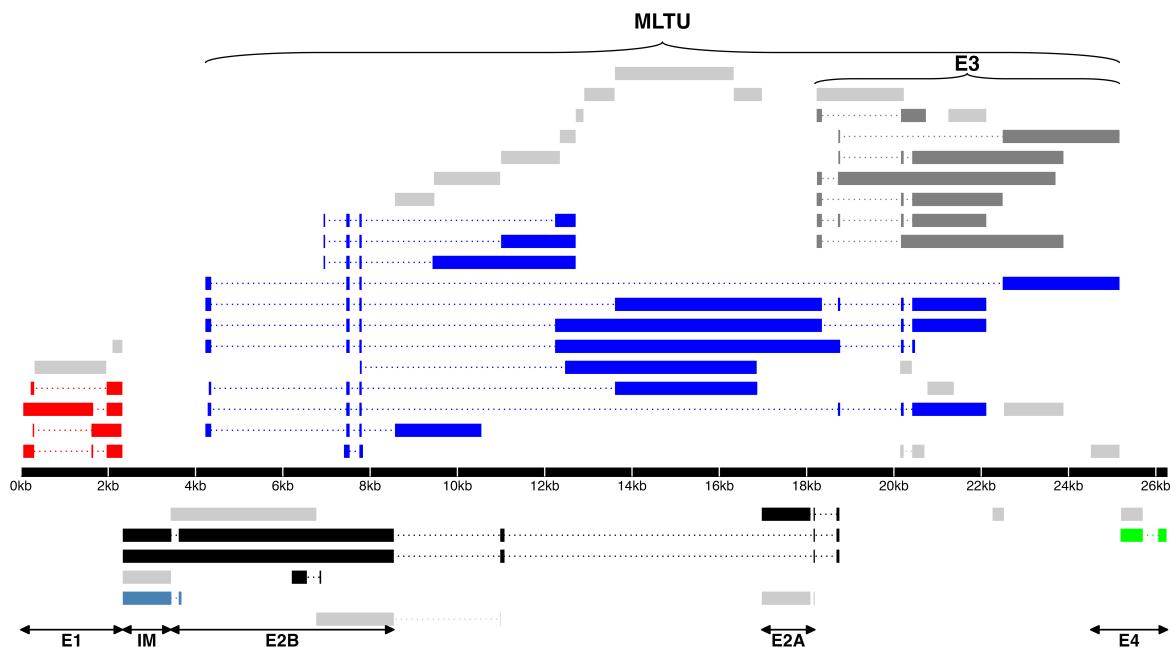
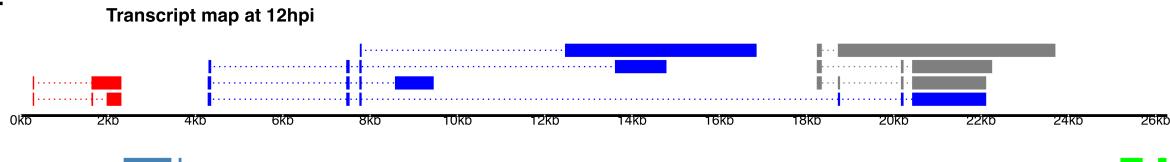
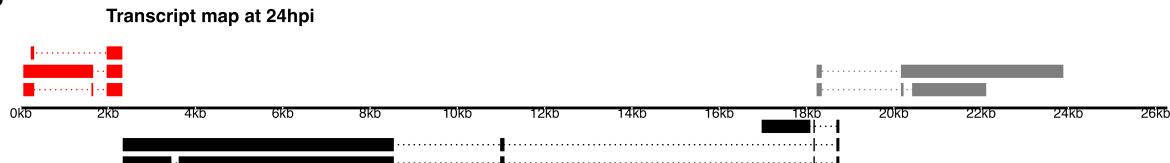
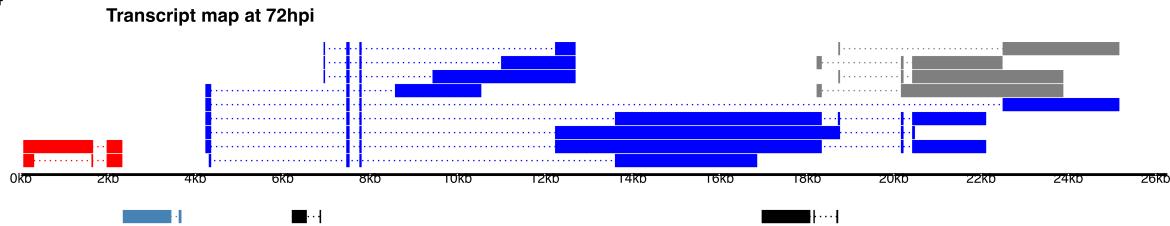
605 **TABLES AND FIGURES**



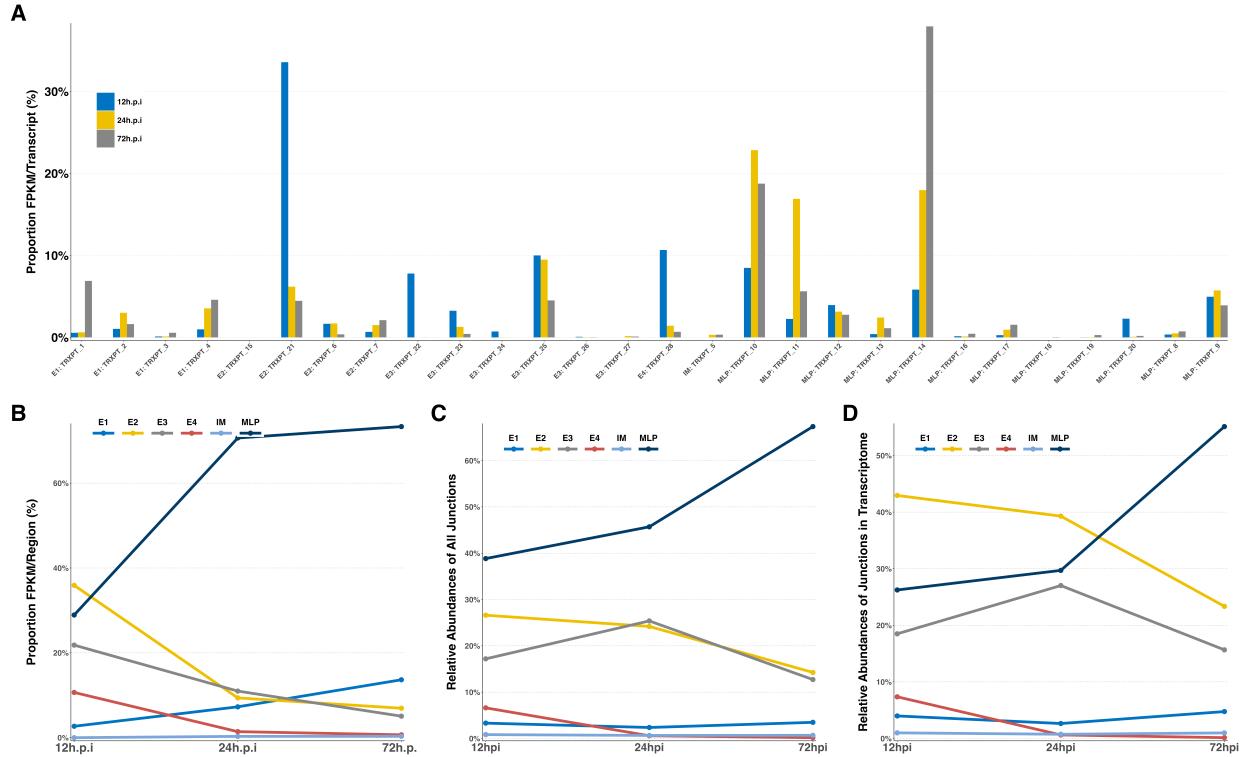
**Figure 1. Predicted ORF map of THEV avirulent strain.** The central horizontal line represents the double-stranded DNA marked at 5kb intervals as white line breaks. Colored blocks represent viral genes. Blocks above the DNA line are transcribed on the sense DNA strand and those below, on the anti-sense strand. pTP, DBP and 33K are predicted to be spliced and are shown as two exons connected with dashed lines. Shaded regions indicate regions containing “genus-specific” genes (colored red). Genes colored in blue are “genus-common”. The gene colored in light green is conserved in all but Atadenoviruses. Regions comprising the different transcription units are labelled at the bottom (E1, E2A, E2B, E3, E4, and IM); the unlabeled regions comprise the MLTU.



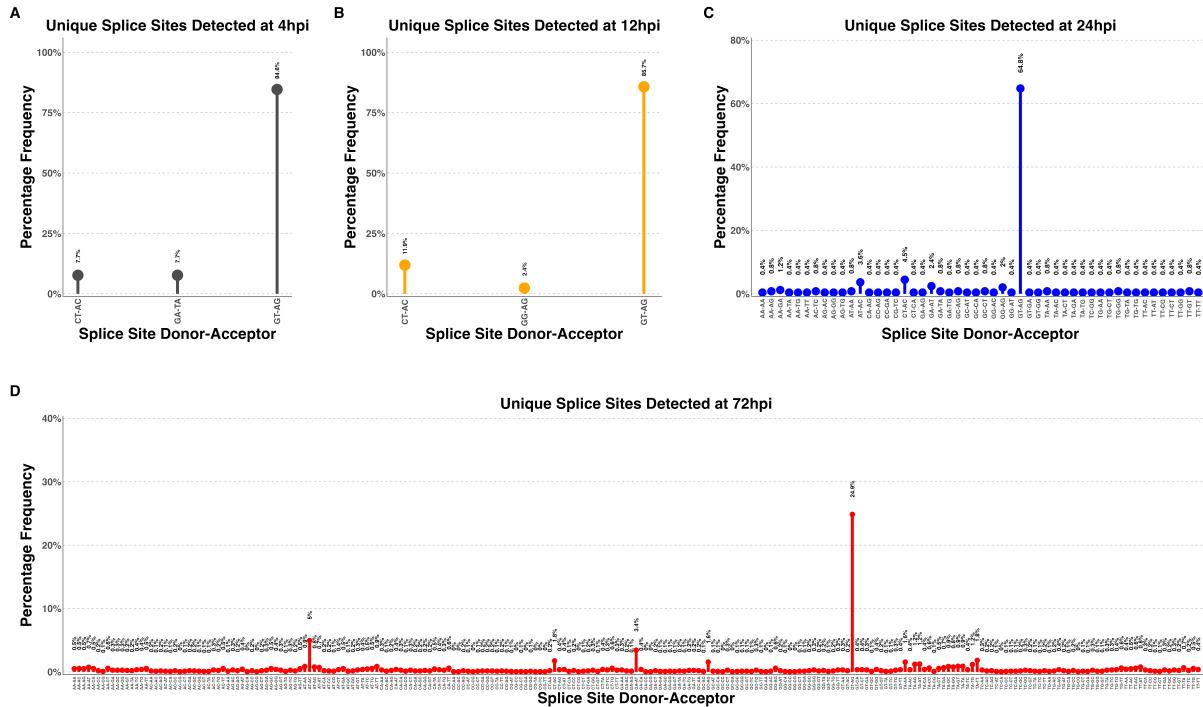
**Figure 2: Increasing levels of THEV over time. A) Per base coverage of sequence reads mapping to THEV genome by time point.** The pileup of mRNA reads mapping to THEV genome at the base-pair level for each indicated time point. **B) Growth curve of THEV (VAS vaccine strain) in MDTC-RP19 cell line.** Virus titers in the freeze-thawed supernatant from infected cells were quantified with a qPCR assay. There is no discernible increase in virus titer up 12 hpi, after which a steady increase in virus titer is measured. The virus titer expands exponentially beginning from 48 hpi, increasing by orders of magnitude before reaching a plateau at 120 hpi. GCN: genome copy number.

**A****1****2****3****4**

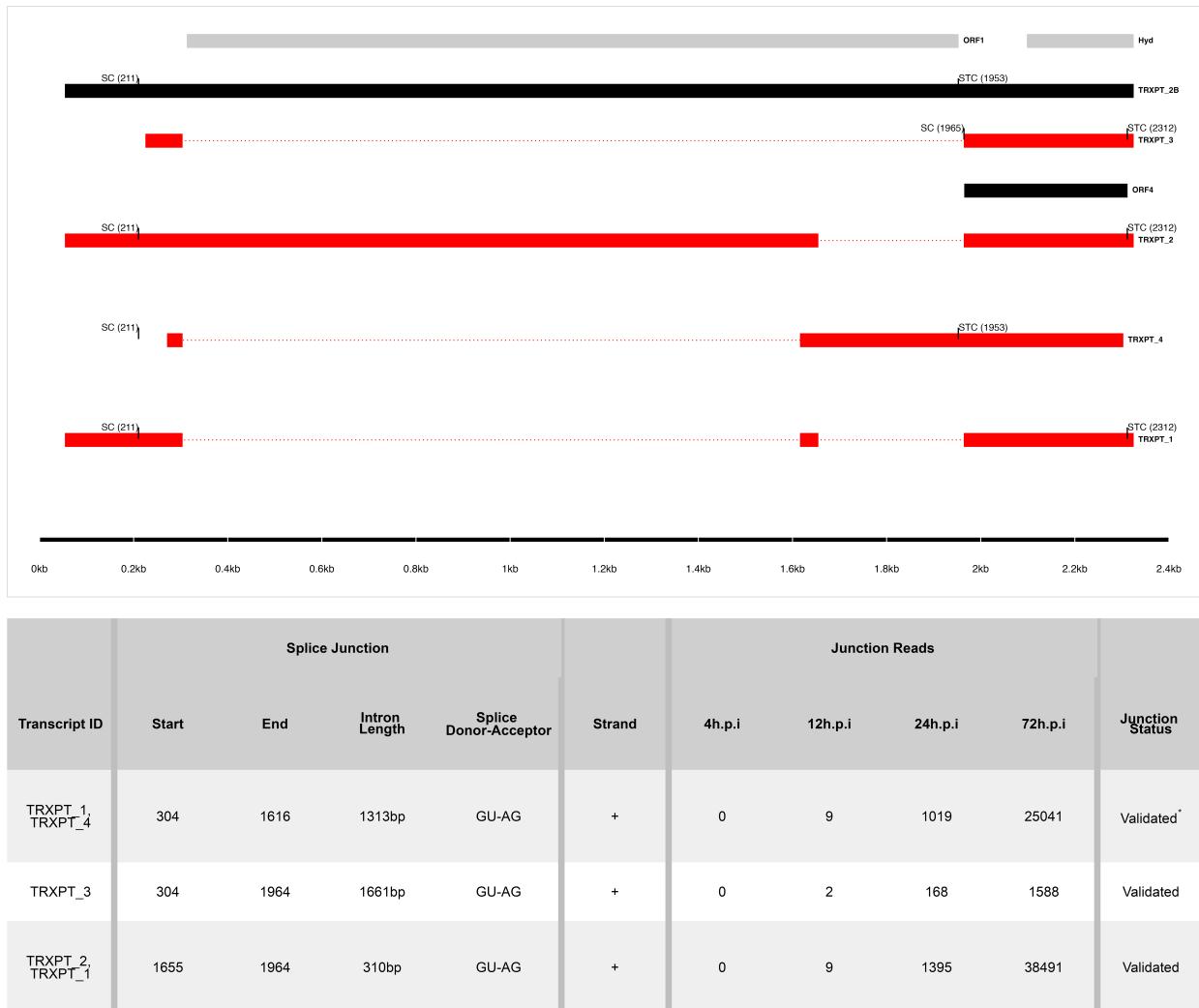
**Figure 3. A) Transcriptome of THEV from RNA-seq.** THEV transcripts assembled from all time points by StringTie are unified forming this final transcriptome (splicing map). Transcripts belonging to the same transcription unit (TU) are located in close proximity on the genome and are color coded and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



**Figure 4: Changes in splicing and expression profile of THEV over time.** **A) Normalized (FPKM) expression levels of transcripts over time.** The expression levels (FPKM) of individual transcripts as a percentage of the total expression of all transcripts at each time point are indicated. Only transcripts from our RNA-seq data are included here. **B) Normalized (FPKM) expression levels of transcripts by region over time.** The expression levels of each region/TU as a percentage of the total expression of all transcripts at each time point are indicated. Region expression levels were calculated by summing up the FPKMs of all transcripts categorized in that region. **C) Relative abundances of all splice junctions grouped by region/TU over time.** After assigning all 2,457 unique junctions to a TU and the total junction reads counted at each time point for each region, the total junction reads for each TU were plotted as percentages of all junction reads at each time point. Note that the junction read counts are not normalized. **D) Relative abundances of junctions in transcriptome grouped by region/TU over time.** This is identical to (C), except that only the junctions found in the full transcriptome obtained from the RNA-seq data were included.

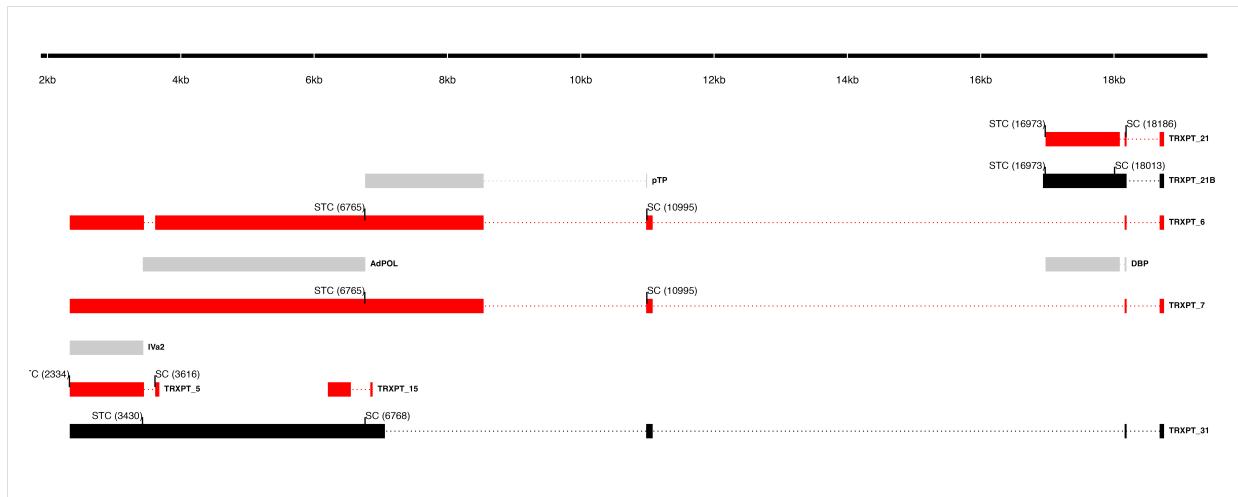


**Figure 5: Changes in splice donor-acceptor nucleotides over time.** The splice donor-acceptor nucleotides of THEV just like other AdVs is mostly the canonical GU-AG. At early time points (4h.p.i and 12h.p.i **(A)** and **(B)**, respectively) the junction nucleotides used appear to be well scrutinized or restricted, utilizing mostly the canonical splice nucleotides. However, as the infection progresses to the late stages (24h.p.i and 72h.p.i **(C)** and **(D)**, respectively), the selectivity of specific splice acceptor-donor pairs seems to degenerate significantly, such that all combinations of nucleotides are utilized.



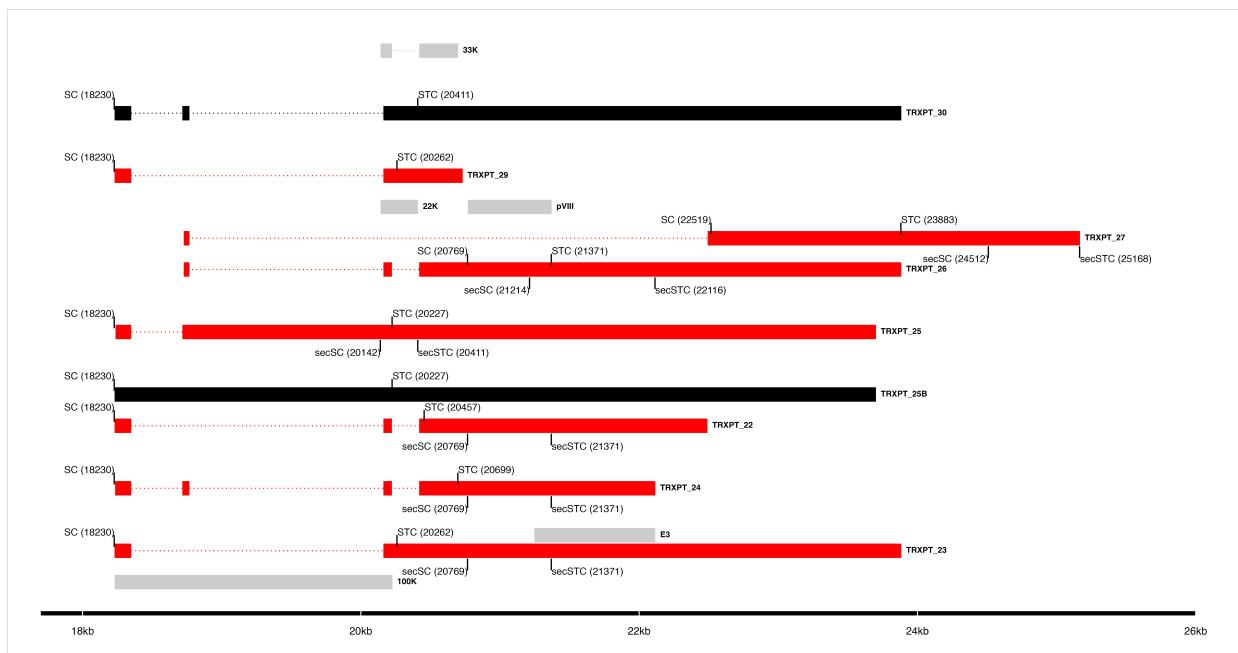
<sup>\*</sup>Not validated for TRXPT\_4

**Figure 6: The splice map of the E1 transcription unit (TU).** Exons are depicted as boxes connected by introns (dotted lines). Transcripts from RNA-seq data are colored red, predicted ORFs are colored grey, and transcripts or ORFs discovered by other means are colored black. Each transcript or ORF is labelled with its name to the right. The start codon (SC) and stop codon (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.



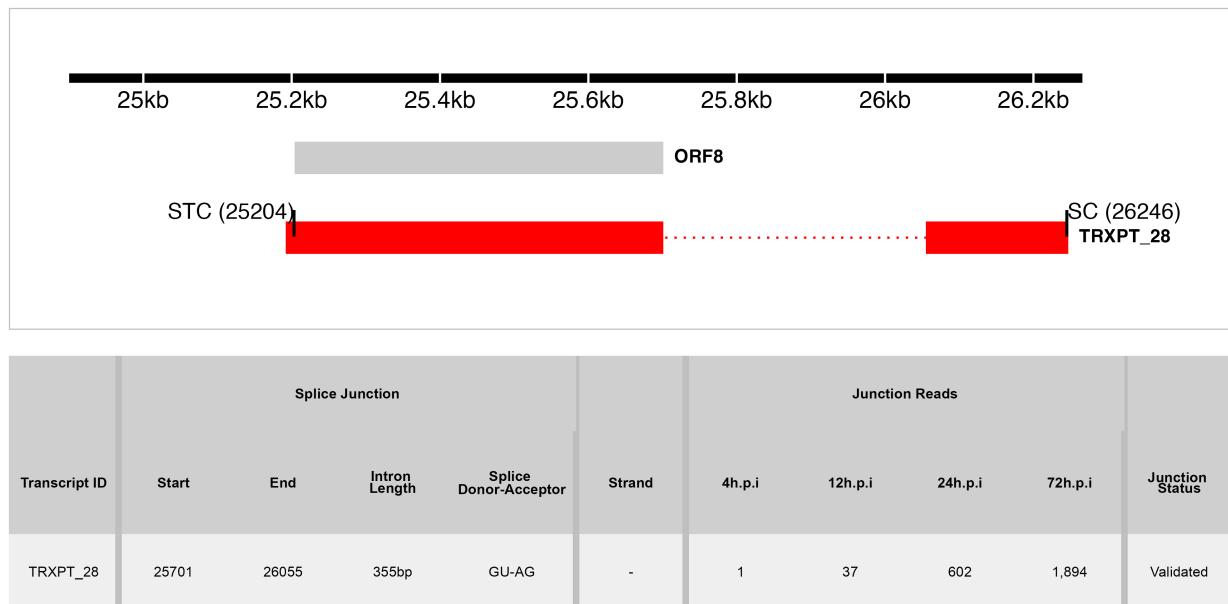
Transcript ID	Splice Junction					Junction Reads					Junction Status
	Start	End	Intron Length	Splice Donor-Acceptor	Strand	4h.p.i	12h.p.i	24h.p.i	72h.p.i		
TRXPT_5, TRXPT_7	3447	3615	169bp	GU-AG	-	1	5	720	13422		Validated
TRXPT_6, TRXPT_7	11079	18159	7081bp	GU-AG	-	0	2	0	0		Validated
TRXPT_21	18087	18159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

**Figure 7: 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 SC and STC of the 5'-most CDS of each transcript are 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				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_25, TRXPT_24, TRXPT_10	18350	18717	368bp	GU-AG	+	4	21	3930	35490	Validated	
TRXPT_23, TRXPT_22, TRXPT_11	18350	20162	1813bp	GU-AG	+	3	18	6619	38841	Validated	
TRXPT_26, TRXPT_24, TRXPT_13, TRXPT_9, TRXPT_10	18768	20162	1395bp	GU-AG	+	2	21	5207	45062	Validated	
TRXPT_26, TRXPT_24, TRXPT_13, TRXPT_11, TRXPT_9, TRXPT_10	20223	20419	197bp	GU-AG	+	3	33	10583	93238	Validated	
TRXPT_27	18768	22492	3725bp	GU-AG	+	0	0	101	1950	Validated	

**Figure 8: The splice map of the E3 TU.** Exons are depicted as boxes connected by introns (dotted lines). Red transcripts are generated from RNA-seq data and predicted ORFs are colored grey. Transcripts discovered by other means are colored black. Each transcript or ORF is labelled with its name to the right. The start codon (SC) and stop codon (STC) of the 5'-most CDS of each transcript are indicated with the nucleotide position in brackets. Similarly, the secondary SC (secSC) and secondary STC (secSTC) are shown. 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.



**Figure 9: The splice map of the E4 TU.** Exons are depicted as boxes connected by introns (dotted lines). The transcript from RNA-seq data is colored red and the predicted ORF, grey. The transcript and ORF are labelled with their names to the right. The start codon (SC) and stop codon (STC) of the 5'-most CDS are 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 junction with its validation status using cloning and Sanger sequencing.



**Figure 10: The splice map of the MLTU.** Exons are depicted as boxes connected by introns (dotted lines). The transcripts from our RNA-seq data are colored red and the predicted ORFs, grey. The transcripts and ORFs are labelled with their names to the right. The start codon (SC) and stop codon (STC) of the 5'-most CDS of each transcript is indicated with the nucleotide position in brackets. Similarly, the secondary SC (secSC) and secondary STC (secSTC) are shown. 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.

Table 1: Overview of sequencing results

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

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

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

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

606 **SUPPLEMENTARY MATERIALS**

607 **Supplementary Table S1A**

Table S1a: Most Transcriptionally Active Regions of THEV at 12h.p.i

Time	Region	Strand	Total Reads	Percentage
12hpi	MLP	+	235	38.8%
12hpi	E2	-	161	26.6%
12hpi	E3	+	104	17.2%
12hpi	E4	-	40	6.6%
12hpi	Unassigned	-,+/-	40	6.6%
12hpi	E1	+	20	3.3%
12hpi	IM	-	5	0.8%

608 **Supplementary Table S1B**

Table S1b: Most Transcriptionally Active Regions of THEV at 24h.p.i

Time	Region	Strand	Total Reads	Percentage
24hpi	MLP	+	52,589	45.7%
24hpi	E3	+	29,208	25.4%
24hpi	E2	-	27,833	24.2%
24hpi	E1	+	2,724	2.4%
24hpi	Unassigned	-,+/-	1,313	1.1%
24hpi	IM	-	744	0.6%
24hpi	E4	-	664	0.6%

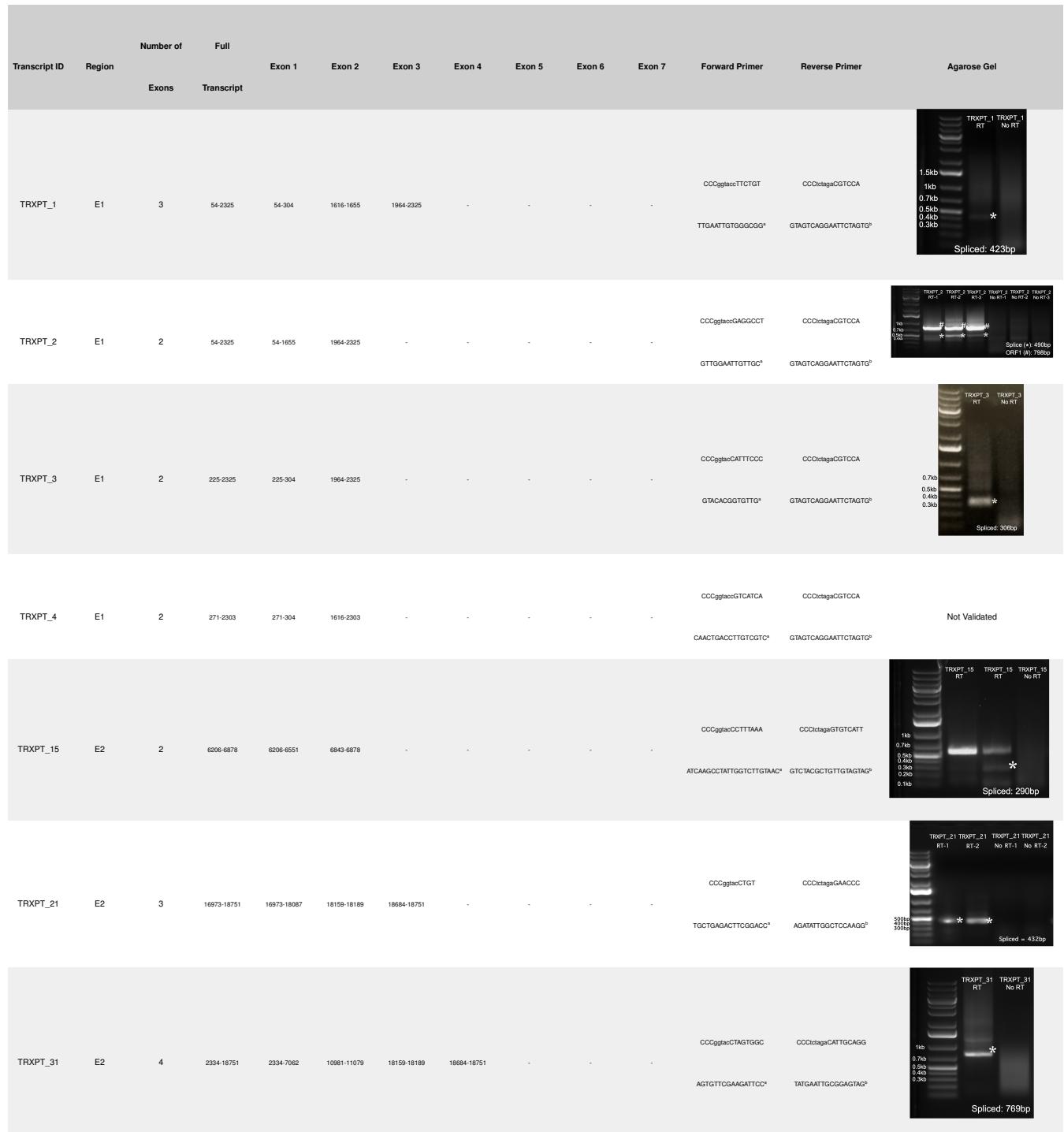
609 **Supplementary Table S1C**

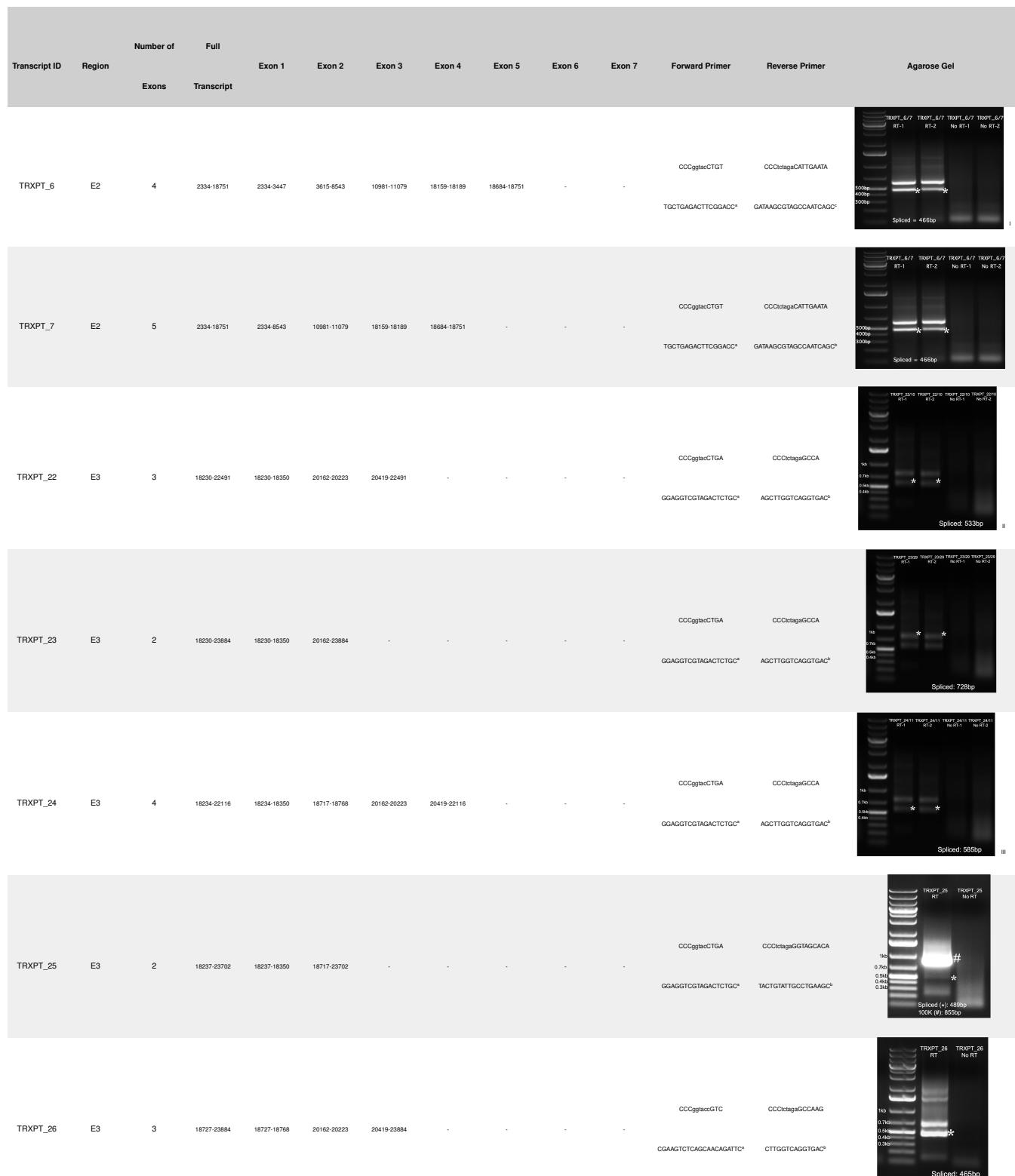
Table S1c: Most Transcriptionally Active Regions of THEV at 72h.p.i

Time	Region	Strand	Total Reads	Percentage
72hpi	MLP	+	1,436,199	67.3%
72hpi	E2	-	304,191	14.3%
72hpi	E3	+	271,310	12.7%
72hpi	E1	+	74,135	3.5%
72hpi	Unassigned	-,+/-	28,921	1.4%
72hpi	IM	-	14,482	0.7%
72hpi	E4	-	3,568	0.2%

610 **Supplementary PCR Methods**

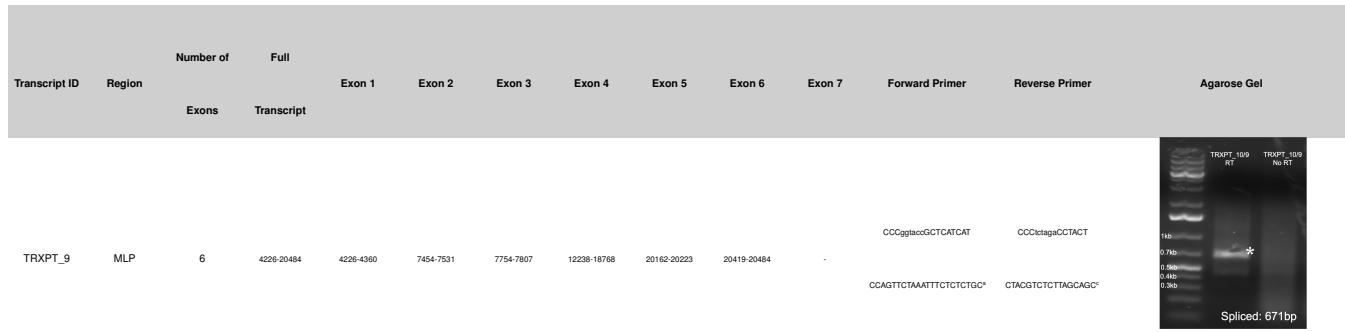
**Agarose Gels Showing PCR Amplification of THEV cDNA With Gene-Specific Primers**





Transcript ID	Region	Number of Exons	Full Transcript							Forward Primer	Reverse Primer	Agarose Gel	
			Exon 1	Exon 2	Exon 3	Exon 4	Exon 5	Exon 6	Exon 7				
TRXPT_27	E3	2	18727-25168	18727-18768	22492-25168	-	-	-	-	CCGgtaccGTC	CCCttagaTGCAAT		
TRXPT_28	E3	2	18230-20732	18230-18350	20162-20732	-	-	-	-	CCGgtacCTGAGGAG	CCCttagaGCCAAG		
TRXPT_28	E4	2	25192-26247	25192-25701	26055-26247	-	-	-	-	CCGgtaccGGACAC	CCCttagaCAGTG		
TRXPT_5	IM	2	2334-3678	2334-3447	3615-3678	-	-	-	-	CCGgtacTCTGGTGAGA	CCCttagaCGCAA		
TRXPT_10	MLP	7	4226-22116	4226-4360	7454-7531	7754-7807	12238-18350	20162-20223	20419-22116	CCGgtaccGCTCATCATC	CCCttagaCCTACTC		
TRXPT_11	MLP	6	4226-22116	4226-4360	7454-7531	7754-7807	13610-18350	18717-18768	20162-20223	20419-22116	CCGgtaccGCTCATCATC	CCCttagaGCTTAG	
TRXPT_12	MLP	4	4226-25168	4226-4360	7454-7531	7754-7807	22492-25168	-	-	CCGgtaccGCTCATCATC	CCCttagaTTCC		

Transcript ID	Region	Number of Exons	Full Transcript							Forward Primer	Reverse Primer	Agarose Gel
			Exon 1	Exon 2	Exon 3	Exon 4	Exon 5	Exon 6	Exon 7			
TRXPT_13	MLP	6	4279-22116	4279-4360	7454-7531	7754-7807	18717-18768	20162-20223	20419-22116	CCGgttaccGCTCATCATC	CCCtctagaGCCAAG	
TRXPT_14	MLP	4	4304-16870	4304-4360	7454-7531	7754-7807	13610-16870	-	-	CCGgttaccGCTCATCATC	CCCtctagaGCTTCAGT	
TRXPT_16	MLP	4	6934-12709	6934-6969	7454-7531	7754-7807	9430-12709	-	-	CCGgttaccGGATCTC	CCCtctagaGCCT	
TRXPT_17	MLP	4	6934-12709	6934-6969	7454-7531	7754-7807	11001-12709	-	-	CCGgttaccGGATCTC	CTCCCCATCTAGAC	
TRXPT_18	MLP	4	6934-12709	6934-6969	7454-7531	7754-7807	12238-12709	-	-	CCGgttaccGGATCTC	CCCtctagaGTTCTC	
TRXPT_19	MLP	2	7401-7836	7401-7531	7754-7836	-	-	-	-	-	-	N/A
TRXPT_20	MLP	2	7765-16856	7765-7807	12466-16856	-	-	-	-	CCGgttaccGAGGATTGA	CCCtctagaCTGAA	
TRXPT_8	MLP	4	4226-10549	4226-4360	7454-7531	7754-7807	8570-10549	-	-	CCGgttaccGCTCATCAT	CCCtctagaCCTATC	



<sup>a</sup>Primer binds inside first exon; <sup>b</sup>Primer binds inside terminal exon; <sup>c</sup>Primer binds inside fourth exon; <sup>d</sup>Agarose gel identical to TRXPT\_7 due to identical splicing; <sup>e</sup>Agarose gel identical to last 3 exons of TRXPT\_10 due to identical splicing; <sup>f</sup>Agarose gel identical to last 4 exons of TRXPT\_11 due to identical splicing; <sup>g</sup>Agarose gel identical to TRXPT\_23 due to identical splicing; <sup>h</sup>Agarose gel identical to TRXPT\_9 due to identical splicing; <sup>i</sup>Agarose gel identical to TRXPT\_14 due to identical splicing;

611 In the table above, the restriction sites in the primer tails are shown in lowercase letters. All the primer melting  
 612 temperatures (TMs) are 58-60°C using a hot start Taq DNA polymerase. The PCR reaction mix was made  
 613 per manufacturer's instructions. The PCR cycling conditions were as follows: Initial denaturation – 95°C for 1  
 614 minute; cyclical denaturation – 95°C for 30 seconds, annealing – variable temperature (53°C-56°C) for 30  
 615 seconds, primer extension – 68°C for variable time, and final elongation – 68°C for 5 minutes. We used 35  
 616 cycles of amplification.

#### 617 **Supplementary Computational Analysis**

618 Snakemake v7.24.0 was used to manage our entire workflow. A graph of the main steps in our pipeline  
 619 generated with Snakemake is shown below. Our trimmed RNA-seq reads were mapped to the genome of  
 620 *M. gallopolo* (with the THEV genome as one of its chromosomes) using Hisat2, to generate the alignment  
 621 (BAM) files and StringTie used to assemble the transcriptome with a GTF annotation file containing the  
 622 predicted THEV ORFs as a guide. The GTF annotation file was derived from a GFF3 annotation file obtained  
 623 from NCBI using Agat - version 1.0.0, a program for converting between many different file formats used in  
 624 bioinformatics. However, the NCBI GFF3 annotation file itself was first modified to remove all unimportant  
 625 features, leaving only the ORFs.

626 StringTie was also used to estimate the normalized expression levels (FPKM) of all the transcripts and  
 627 Ballgown in R was used to perform statistical analysis and comparisons of the transcript expression levels,  
 628 which was instructive in understanding the temporal regulation THEV gene expression.

629 In these steps above, each sample (replicate of each time point) was processed independently and merged  
 630 only in the final transcriptome assembly or during analysis with Ballgown. In the subsequent steps described  
 631 below, all samples for each time point were processed together.

632 We used RegTools to extract and analyze the splice junctions in the BAM files. The command "regtools

633    junctions extract" provides a wealth of information about all the splice sites in the BAM file provided such as:  
634    the start and end positions, the strand, and number of reads supporting the splice junctions. The command  
635    "regtools junctions annotate" gives even more information such as: the splice site donor-acceptor sequences  
636    and transcripts/genes that overlap the junction. This information was the basis for estimating and comparing  
637    the splicing activity of different regions (TUs) of THEV over time. Also, Samtools was also used to count the  
638    total sequencing reads for all replicates at each time point.



A flowchart of the major steps in the computational analysis pipeline (*generated with Snakemake*)