

¹ Characterizing the Splice Map of Turkey Hemorrhagic Enteritis
² Virus

³

⁴ Abraham Quaye^{†,a}, Brett E. Pickett^a, Joel S. Griffitts^a, Bradford K. Berges^a, Brian D. Poole^{a,*}

⁵ ^aDepartment of Microbiology and Molecular Biology, Brigham Young University

⁶ [†]First-author

⁷ ^{*}Corresponding Author

⁸ **Corresponding Author Information**

⁹ brian_poole@byu.edu

¹⁰ Department of Microbiology and Molecular Biology,

¹¹ 4007 Life Sciences Building (LSB),

¹² Brigham Young University,

¹³ Provo, Utah

¹⁴

15 **ABSTRACT**

16 **Background:** Hemorrhagic enteritis, caused by *Turkey Hemorrhagic Enteritis Virus (THEV)*, is a disease
17 affecting turkey pouls characterized by immunosuppression and bloody diarrhea. An avirulent THEV strain
18 that retains the immunosuppressive ability is used as a live vaccine. Characterizing the splice map of THEV
19 is an essential step that would allow studies of individual genes mediating its immunosuppressive functions.
20 We used RNA sequencing to characterize the splice map of THEV for the first time, providing key insights
21 into the THEV gene expression and mRNA structures.

22 **Methods:** After infecting a turkey B-cell line with the vaccine strain, samples in triplicates were collected
23 at 4-, 12-, 24-, and 72-hours post-infection. Total RNA was extracted, and poly-A-tailed mRNA sequenced.
24 Reads were mapped to the THEV genome after trimming and transcripts assembled with StringTie. We
25 performed PCR of THEV cDNA, cloned the PCR products, and used Sanger sequencing to validate all
26 identified splice junctions.

27 **Results:** Researchers previously annotated the THEV genome as encoding 23 open reading frames (ORFs).
28 We identified 29 spliced transcripts from our RNA sequencing data, all containing novel exons although
29 some exons matched some previously annotated ORFs. The three annotated splice junctions were also
30 corroborated by our data. During validation we identified five additional unique transcripts, a subset of which
31 were further validated by 3' rapid amplification of cDNA ends (3' RACE). Thus, we report that the genome of
32 THEV contains 34 transcripts with the coding capacity for all annotated ORFs. However, we found six of the
33 previously annotated ORFs to be truncated ORFs on the basis of the identification of an in-frame upstream
34 start codon or the detection of additional coding exons. We also identified three of the annotated ORFs with
35 longer or shorter isoforms, and seven novel unannotated ORFs that could potentially be translated; although
36 it is beyond the scope of this manuscript to investigate whether they are translated.

37 **Conclusions:** Similar to human adenoviruses, all THEV transcripts are spliced and organized into five
38 transcription units under the control of their cognate promoters. The genes are expressed under temporal
39 regulation and THEV also produces multiple distinctly spliced transcripts that code for the same protein.
40 **Studies of the newly identified potential proteins should be urgently performed as these proteins may have**
41 **roles in THEV-induced immunosuppression. Also, knowing the splicing of THEV genes should be invaluable**
42 **to future research focusing of studying THEV genes, as this will allow accurate cloning of the mRNAs.**

43 **KEY WORDS**

44 Alternative splicing, Turkey hemorrhagic enteritis virus, Adenovirus, Transcriptome, RNA sequencing.

45 **BACKGROUND**

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

59 Virulent THEV strains (THEV-V) and avirulent strains (THEV-A) of THEV both infect turkeys, with THEV-
60 V causing hemorrhagic enteritis (HE), a debilitating acute disease predominantly affecting turkey pouls
61 characterized by immunosuppression, intestinal lesions leading to bloody diarrhea, and up to 80% mortality
62 (2, 11–13). While the current vaccine strain (a THEV-A called Virginia Avirulent Strain [VAS]) has proven
63 effective at preventing HE in turkey pouls, it still retains its immunosuppressive ability. Thus, vaccinated
64 birds are rendered more susceptible to opportunistic infections and death than unvaccinated birds leading to
65 substantial economic losses (11, 14–16). To eliminate the immunosuppressive immunosuppressive effect of
66 the vaccine strain, a thorough investigation of the culprit viral genes mediating this phenomenon is essential.
67 However, the transcriptome (splicing and gene expression patterns) of THEV has not been characterized,
68 making an investigation of specific immunosuppressive viral genes impractical.

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

77 in DNA replication, or providing the necessary intracellular niche for optimal replication while late genes
78 encode structural proteins that act as capsid proteins, promote virion assembly, or direct genome packaging.
79 The immediate early genes E1A are expressed first, followed by the delayed early genes, E1B, E2, E3 and
80 E4. Then the intermediate early genes, IVa2 and pIX are expressed followed by the late genes (10, 17,
81 18). It is noteworthy that the MLP shows basal transcriptional activity during early infection (before DNA
82 replication), with a comparable efficiency to other early viral promoters, but it reaches its maximal activity
83 during late infection (after DNA replication). However, during early infection only a subset of the MLP-derived
84 transcripts are expressed (10). ~~MAdV makes an~~ MAdVs make extensive use of alternative RNA splicing
85 and polyadenylation to produce a very complex array of mRNAs. All but the pIX mRNA undergo at least
86 one splicing event. For instance, the MLTU produces over 20 distinct splice variants all containing three
87 non-coding exons at the 5'-end (collectively known as the tripartite leader; TPL) (17, 18). There is also
88 an alternate three-exon 5' non-coding leader sequence present in varying amounts on a subset of MLTU
89 mRNAs (known as the x-, y-, and z-leaders). Lastly, there is the i-leader exon, which is infrequently included
90 between the second and third TPL exons, and codes for the i-leader protein (20). Thus, the MLTU produces
91 a complex repertoire of mRNA with diverse 5' untranslated regions (UTRs) spliced onto different 3' coding
92 exons which are grouped into five different 3'-end classes (L1-L5) based on polyadenylation site. Each
93 transcription unit (TU) contains its own promoter driving the expression of the array of mRNA transcripts
94 produced via alternative splicing in the unit (10, 17, 18). The promoters are activated at different phases of
95 infection by proteins from previously activated TUs. Paradoxically, the early-to-late phase transition during
96 infection requires the L4 gene products, 22K and 33K, which should only be available after the transition.
97 However, a promoter in the L4 region (L4P) that directs the expression of these two proteins independent of
98 the MLP was found, resolving the paradox (10, 17, 21). During translation of AdV mRNA, recent studies
99 using long-read direct RNA sequencing strongly suggest the potential usage of secondary start codons;
100 adding to what was already a highly complex system for gene expression (17, 22).

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

108 **RESULTS**

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

110 A prior study by Aboeazz *et al.* demonstrated that nearly all THEV transcripts became detectable starting
111 at 4 hours post-infection (hpi), with one replication cycle concluding around 18 hpi (24). Consequently, we
112 harvested infected MDTC-RP19 cells (**MOI**
[multiplicity of infection \(MOI\)](#) of 100 genome copy numbers/cell)
113 at 4-, 12-, 24-, and 72-hpi to capture all transcripts within a broad time window. Our paired-end RNA
114 sequencing (RNA-seq) experiment generated an average of 107.1 million total reads of 149 bp length per
115 time-point. These reads were concurrently mapped to both the virus (THEV) and host (*Meleagris gallopavo*)
116 genomes using the Hisat2 (25) reference-based aligner. A total of 18.1 million reads from all time-points
117 mapped to the virus genome, providing comprehensive coverage and leaving no regions unmapped. The
118 mapped reads to the virus genome increased significantly from a scant 432 reads at 4 hpi to 16.9 million
119 reads at 72 hpi (**Table 1, Figure 2A**). From these mapped reads, we identified 2,457 unique THEV splice
120 junctions across all time-points, with later time-points exhibiting significantly more sequence reads supporting
121 the splice junctions than earlier time-points. For instance, all 13 unique junctions at 4 hpi had fewer than 10
122 supporting reads each, averaging only 2.8 reads per junction. In contrast, the 2374 unique junctions at 72
123 hpi averaged 898.4 reads per junction, with some junctions reaching as high as 322,677 reads. The marked
124 increase in splice junction and mapping reads to the THEV genome over time indicates an active infection
125 and successful viral replication, which is corroborated by our quantitative PCR (qPCR) assay that quantified
126 the total number of viral genome copies over time (**Figure 2B**).

127 Using StringTie (25), we assembled the data into potential transcripts, guided by the genomic locations of
128 the previously predicted THEV ORFs. In the consolidated transcriptome, a composite of all non-redundant
129 transcripts across all time points, we identified a total of 29 novel transcripts. We found that a subset of
130 exons in the viral transcripts match some predicted ORFs exactly, with the majority of the exons being longer
131 and spanning multiple predicted ORFs (**Figure 3**).

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

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

139 AdV gene expression is subject to meticulous temporal regulation, with each promoter typically generating
140 one or a few pre-mRNAs. These pre-mRNAs undergo alternative splicing to produce a diverse array of
141 mature mRNAs. To assess the temporal activity of each promoter, we utilized StringTie and Ballgown (a tool
142 for statistical analysis of assembled transcriptomes) (26). These tools estimated the normalized expression
143 levels of all transcripts at each time point, measured in Fragments Per Kilobase of transcript per Million
144 mapped reads (FPKM) units. At 4 hpi, we counted very few unique splice junctions, reads, and transcripts;
145 hence, this time-point was excluded from this analysis.

146 Examining individual mRNAs, TRXPT_21 – from the E2 region – was the most significantly expressed at
147 12 hpi, constituting 33.58% of the total expression of all transcripts. Transcripts in the E3 and E4 regions
148 also contributed substantial proportions, along with some MLP region transcripts. The later time points were
149 dominated by the MLP region transcripts — TRXPT_10 and TRXPT_14 were the most abundantly expressed
150 at 24 and 72 hpi, respectively (**Figure 4A**). Our analysis of the FPKM values of transcripts per region/TU
151 revealed a similar pattern: the E2 region was the most abundantly expressed at 12 hpi, after which the MLP
152 region assumed dominance (**Figure 4B**).

153 Next, we estimated the relative abundances of all splice junctions at each time point using the raw reads.
154 Only junctions with a read coverage of at least 1% of the total splice junction reads at the given time point
155 were considered significant and included in **Tables 2a-2c**. At 12 hpi, 18 junctions met the 1% threshold,
156 predominantly from early regions (E1, E2, E3, and E4), although the MLTU was the single most predominant
157 region overall, constituting 38.8% of all the junction reads (**Table 2a** and **Supplementary Table S1a**). The
158 most abundant junctions at 12 hpi remained the most significantly expressed at 24 hpi. However, here, the
159 MLP-derived junctions unsurprisingly became even more predominant overall, accounting for 45.7% of all
160 the junction reads counted (**Table 2b** and **Supplementary Table S1b**). At 72 hpi, the trend of increased
161 activity of the MLP continued as expected; at this time, the MLP region junctions were not only the most
162 abundant overall — accounting for 67.3% of all junction reads, — but also contained the most significantly
163 expressed individual junctions (**Table 2c**, **Supplementary Table S1c** and **Figure 4C**). When we limited this
164 analysis to only junctions in the final transcriptome, we observed the relative abundances of the junctions for
165 each region over time to be similar to the pattern seen with all the junctions included (**Figure 4D**).

166 Finally, we analyzed splice donor and acceptor site nucleotide usage over time to investigate any peculiarities
167 that THEV may exhibit, generally or over the course of the infection. We found that most splice donor-acceptor
168 sequences were, unsurprisingly, the canonical GU-AG nucleotides. However, the splice acceptor-donor
169 pairing became less specific over time, such that all combinations of nucleotide pairs were eventually detected
170 (**Figure 5**).

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

172 In MAdVs, E1 is the first region transcribed post-viral DNA entry into the host cell nucleus, mediated solely
173 by host transcription machinery (18). Translated E1 proteins subsequently activate other viral promoters in
174 conjunction with host transcription factors (10). Despite subdivision into E1a and E1b units in MAdVs, our
175 THEV data does not reflect this. This region is predicted to encode only two ORFs: ORF1 (sialidase) and
176 Hyd (a hydrophobic product with unknown function) in THEV.

177 We identified four novel transcripts in this region, containing 3 unique splice junctions (**Figure 6**), and
178 encoding four distinct novel ORFs in addition to Hyd. All transcripts have coding potential (CP) for the Hyd
179 protein as the 3'-most coding sequence if secondary start codon usage is considered (17, 18). Also, all
180 the transcripts have a common transcription termination site (TTS; at position 2325 bp), but TRXPT_1 and
181 TRXPT_2 have an upstream transcription start sites (TSS) to TRXPT_3 and TRXPT_4. Given that E1
182 mRNAs in MAdVs share a common TTS and TSS, differing only in the internal splicing (18), we consider
183 the upstream TSS (position 54 bp) as the actual TSS for all E1 region transcripts. We also identified the
184 canonical polyadenylation signal (PAS; AAUAAA) in the immediate context of the TTS at position 2323 bp
185 (location of the “U” in the PAS sequence); see **Supplementary Table S2**.

186 From the 5'-most start codon (SC), TRXPT_1 encodes a multi-exonic novel 17.9 kDa, 160 residue protein
187 (ORF9). TRXPT_2 encodes a two-exon 66.4 kDa, 597 residue novel protein (ORF10), spanning almost the
188 entire predicted ORF1 and Hyd. The intron of TRXPT_2 excludes the C-terminus of ORF1 (including its
189 stop codon) from ORF10 but the SC of ORF10 is 102 bp upstream and in-frame with the predicted SC of
190 ORF1. TRXPT_3, similar to TRXPT_1 but lacking the second exon, encodes a 13.1 kDa, 115 residue protein
191 (ORF4), previously predicted (27) but excluded in later annotations (1, 12). Our data suggest it is genuinely
192 expressed. Lastly, TRXPT_4 encodes a distinct novel 15.9 kDa, 143 residue protein (ORF11).

193 The splice junctions of all transcripts in this region, except for TRXPT_4, were validated by cloning of viral
194 cDNA and Sanger sequencing (see **Supplementary PCR methods**). During TRXPT_2 validation, ORF1
195 was found on the agarose gel (an unspliced band size) and Sanger sequencing results showed it to be a
196 transcribed mRNA (**Supplementary PCR methods**). This was corroborated by our 3' RACE experiment,
197 which showed a transcript (TRXPT_2B) spanning the entire ORF1 and Hyd ORFs without splicing, with a
198 poly-A tail immediately after the E1 TTS. The 5'-most coding sequence (CDS) of this transcript (TRXPT_2B)
199 encodes ORF1. However, TRXPT_2B has an upstream and in-frame SC to the predicted SC of ORF1,
200 suggesting that the predicted ORF1 CDS is truncated – the expressed ORF1 (eORF1) shares the same SC
201 as ORF10, but has a unique stop codon (STC). See **Supplementary Table S3** for all transcripts and their
202 encoded proteins.

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

204 The E2 TU expressed on the anti-sense strand, is subdivided into E2A and E2B and encodes three classical
205 AdV proteins – pTP and Ad-pol (E2B proteins), and DBP (E2A protein) – essential for genome replication (17,
206 18). Unlike MAdV where two promoters are known (17), we discovered only a single TSS (E2 TSS; 18,751
207 bp) for both E2A and E2B transcripts in THEV. However, E2A and E2B transcripts have distinct TTSs, with
208 E2B transcripts sharing the TTS of the IVa2 transcript of the IM region similar to MAdVs (17, 18) (**Figure 7**).
209 The E2A ORF, DBP, is one of three THEV ORFs predicted to be spliced from two exons. The corresponding
210 transcript (TRXPT_21) in our data matches the predicted splice junction but includes an additional non-coding
211 exon at the 5'-end (E2-5'UTR). Thus, TRXPT_21 is a three-exon transcript encoding DBP (380 residues,
212 43.3 kDa) precisely. TRXPT_21 was also corroborated in a 3' RACE experiment. Additionally, from the 3'
213 RACE data, we found a splice variant of TRXPT_21 which retains the second intron, leading to a 2-exon
214 transcript (TRXPT_21B). Although longer, TRXPT_21B encodes a truncated isoform of DBP (tDBP; a 346
215 residue, 39.3 kDa product) using a downstream in-frame SC but the same STC as DBP. Both TRXPT_21 and
216 TRXPT_21B share a common TTS, seen in our 3' RACE data, located 39 bp downstream of the CDS in an
217 adenine/thymine (A/T)-rich sequence followed by the poly-A tail sequence, suggesting this position (16,934
218 bp) as the true E2A TTS. There are two canonical PASs (AAUAAA; 16,964 and 16950 bp) immediately after
219 the CDS any of which can serve as the PAS without affecting the encoded proteins (**Supplementary Table**
220 **S2**).

221 The E2B region transcripts also start with the E2-5'UTR but extend downstream to reach the TTS at
222 2334 bp in the IM region, which is in the immediate context of a canonical PAS (position 2333 bp) where
223 polyadenylation likely occurs. The E2B transcripts, TRXPT_6 and TRXPT_7, are almost identical except
224 for an extra splice junction at the 3'-end of TRXPT_6 (**Figure 7**). TRXPT_7 has the CP for both classical
225 proteins (pTP and Ad-pol) encoded in this region, with the pTP ORF predicted to be spliced from two exons.
226 The predicted splice junction of pTP is corroborated by our data but the full transcript is markedly longer
227 than the predicted ORF, although the encoded product (pTP) remains unchanged. Ad-pol (polymerase) is
228 encoded downstream of pTP with secondary start codon SC (secSC) usage. The CP of TRXPT_6 slightly
229 differs from TRXPT_7 because a new STC resulting from the extra splice site forms a minimal truncation of
230 the Ad-pol encoded from its secSC.

231 While both TRXPT_6 and TRXPT_7 have the CP for Ad-pol with secSC usage, in all AdVs studied, the
232 two proteins (pTP and Ad-pol) are encoded by separate mRNAs with identical first three 5' exons and TTS,
233 but different splice junctions to the terminal coding exons. Hence, we checked for a longer splice junction
234 between the third and fourth (terminal) exons of TRXPT_7 with our junction validation method (targeted PCR,

235 cloning, and Sanger sequencing). We discovered a unique splice junction (10,981-7062 bp) not present in
236 our RNA-seq data. If initiated from the E2 TSS and terminated at the E2B TTS, this transcript (TRXPT_31)
237 would encode Ad-pol in its 5'-most CDS (**Figure 7**).

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

242 The IM region is a single-transcript TU, encoding a single classical protein, IVa2. The promoter expressing
243 this single transcript (TRXPT_5) is embedded in the E2B region and shares a TTS with E2B transcripts (17,
244 18). TRXPT_5 is a two-exon transcript with a non-coding first exon, except the last 2 nucleotides, which
245 connect with the first nucleotide of the second exon to form the 5'-most SC. This new SC is four codons
246 upstream and in-frame of the predicted IVa2 SC. Beside the four additional N-terminus residues, the protein
247 sequence is unchanged.

248 **Early Region 3 (E3) transcripts**

249 The E3 region, nested within the MLTU, encodes proteins that modulate and evade host immune defenses.
250 In MAdVs, this region contains seven ORFs expressed from multiple transcripts sharing the same TSS (from
251 the E3 promoter) but having different TTSs (10, 17, 18). However, some E3 transcripts use the TSS of the
252 MLP. Due to sharing the same TSS, in MAdVs, secSC usage is heavily relied on for gene expression in this
253 region as utilizing only the first SC cannot produce the downstream proteins in this TU. The 12.5K ORF and
254 transcripts using the MLP TSS are exceptions (17).

255 In THEV, only one ORF (E3) was predicted in this region. However, as the E3 TU is nested in the MLTU,
256 transcripts from the L4 promoter (100K, 22K, 33K, and pVIII) overlap the E3 region transcripts entirely and
257 share similar TSS and TTS locations (**Figure 8**). Therefore, we have categorized these two groups together
258 as E3 transcripts.

259 We identified seven novel transcripts (TRXPT_22, TRXPT_23, TRXPT_24, TRXPT_25, TRXPT_26,
260 TRXPT_27, TRXPT_29) from our RNA-seq data, all originating from two distinct TSSs. We consider the
261 first TSS (position 18,230 bp) as corresponding to the L4 promoter (L4P) and the other at 18,727 bp as
262 corresponding to the E3 promoter (E3P). We also identified the canonical or other known PAs (28) near
263 the TTS of the transcripts (see **Supplementary Table S2**). These E3 transcripts collectively have the CP
264 for several predicted THEV ORFs: 100K, 22K, 33K, pVIII, and E3, as well as Fiber (IV) and ORF7 of the
265 MLTU (see **Supplementary Table S3**). However, some of these CDSs differ from the predictions due to

either unknown exons or the presence of an in-frame upstream SC. For instance, we discovered that 33K, predicted to be spliced from two exons, is actually a significantly longer four-exon ORF (e33K; 19.8 kDa, 171 residues) encoded on TRXPT_24. Its first two exons were unknown but the last two match the predicted exons and the CDS is in-frame, albeit the first 20 bp of the predicted 33K (including the SC) is spliced out as part of the second intron of TRXPT_24. TRXPT_24 also has the CP for pVIII and E3 if we consider downstream SC usage. However, we found an upstream in-frame SC for the predicted E3; thus, this longer version of E3 (eE3) is likely the genuinely expressed ORF. TRXPT_29, the shortest transcript in this TU, encodes a novel 73 residue protein (8.3KI) across its two exons using the SC of e33K with a unique STC. TRXPT_23, spliced identically as TRXPT_29, also encodes 8.3KI from its first SC. Similarly, TRXPT_22 encodes a 73 residue novel protein (8.3KII) from its first SC that shares over 80% similarity with 8.3KI, but they differ at the C-terminus. Considering downstream SC usage, both TRXPT_22 and TRXPT_23 can encode pVIII and eE3 in that order, but TRXPT_23 being longer, also has the CP for the Fiber ORF.

As the splice junctions of TRXPT_22, TRXPT_23, TRXPT_24, and TRXPT_29 share the same genomic space, their validation was done with a single primer pair, and they were differentiated from each other by cloning the cDNA and Sanger sequencing (**Supplementary PCR methods**). In addition to corroborating the splice junctions for the aforementioned transcripts, the Sanger sequencing results also showed a distinct splice variant undetected in our RNA-seq transcriptome. This was a three-exon transcript (TRXPT_30) with identical first and last exons as TRXPT_23, which also contained the second exon of TRXPT_24 (**Figure 8**). TRXPT_30 encodes a novel 140 residue, 15.7kDa protein (e22K), spanning all three exons. Interestingly, the last 81 C-terminus residues of e22K are identical to 22K (89 residues), a single-exon ORF predicted to use the same SC as 33K. Just as seen for 33K, the first 20 bp of 22K is intronic, excluding the first 7 residues of 22K from e22K. We consider e22K as a long variant of the predicted 22K ORF. Assuming TRXPT_30 shares the same TSS and TTS as TRXPT_23, it would also have the downstream CP of TRXPT_23.

TRXPT_25, the largest transcript in the TU, is a two-exon transcript, encoding a novel protein (t100K; 543 residues), which is a shorter isoform of the predicted 100K ORF. secSC usage on this transcript yields the predicted 22K ORF. It also has the CP for pVIII and eE3 downstream. 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 corresponding 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. We assume that this transcript (TRXPT_25B) 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 encodes

298 pVIII as the 5'-most ORF and has the CP for eE3 and Fiber in that order. TRXPT_27, a two-exon transcript,
299 encodes Fiber as the 5'-most ORF, and ORF7 downstream with secSC usage. TRXPT_13 is an L4P
300 transcript that uses the MLP TSS; it is discussed under the MLTU transcripts.

301 **Early Region 4 (E4) transcripts** This TU is found at the 3'-end of the genome and expressed on the
302 anti-sense strand. Based on nucleotide position, ORF7 and ORF8 were predicted in this region (1); however,
303 as ORF7 is neither on the anti-sense strand nor transcribed from a promoter in the E4 region, we only
304 classify ORF8 in this TU. This is corroborated by our RNA-seq data, showing only one transcript in this
305 region on the anti-sense strand (**Figure 9**). The transcript (TRXPT_28) spans 25192-26247 bp and is spliced
306 at 25701-26055 bp, forming a two-exon transcript. The second exon fully matches the predicted ORF8 with
307 12 extra base pairs at the 3'-end. However, we identified a SC 192 bp upstream of the predicted SC in the
308 first exon from which an in-frame protein is encoded. We consider this longer isoform (eORF8 – 26.4 kDa,
309 229 residues) as the genuinely expressed ORF. We also identified a canonical PAS 11 bp upstream of TTS
310 (**Supplementary Table S2**).

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

312 The MLTU transcripts, dominant in the late phase of the AdV infectious cycle, are produced by alternative
313 polyadenylation and splicing of a primary transcript and grouped into five transcript classes (L1-L5). About
314 13 out of the 23 predicted ORFs in THEV fall within this TU, some of which we have categorized under the
315 E3 TU instead. Our RNA-seq data revealed 12 transcripts (TRXPT_8, TRXPT_9, TRXPT_10, TRXPT_11,
316 TRXPT_12, TRXPT_13, TRXPT_14, TRXPT_16, TRXPT_17, TRXPT_18, TRXPT_19, TRXPT_20) in this
317 TU, most of which have the 5' TPL sequence as in all AdVs. However, three transcripts (TRXPT_16,
318 TRXPT_17, TRXPT_18) use a different leader sequence (sTPL), where a different first exon is used instead
319 of the first TPL exon, and TRXPT_20 uses only the third TPL exon (TPL3); see **Figure 10**.

320 We identified five TTSs (10,549, 12,709, 16,870, 22,116, 25,168 bp) in this TU, which we consider as
321 corresponding to the five late mRNA classes (L1-L5), respectively. L1 mRNAs include TRXPT_8, encoding
322 the 52K ORF as predicted. L2 mRNAs include TRXPT_16, TRXPT_17, and TRXPT_18, all containing the
323 sTPL with their respective coding exons. They encode pIIa, III (penton), and pVII, respectively. The L3
324 mRNAs, TRXPT_14 and TRXPT_20, both encode the hexon (II) ORF but hexon is the only ORF encoded on
325 TRXPT_14, whereas TRXPT_20 encodes pX (pre-Mu), pVI, and hexon in that order. L4 mRNAs, TRXPT_9,
326 TRXPT_10, TRXPT_11, and TRXPT_13 are the largest transcripts in the transcriptome and encode several
327 similar late proteins. TRXPT_9 and TRXPT_10 are very similar but not identical. The last exon of TRXPT_9
328 seems to be truncated and likely shares the same TTS as TRXPT_10. They both encode pVII as the 5'-most
329 ORF and also have the CP for pX, pVI, hexon, a longer variant of protease (eProt) from an upstream in-frame

330 SC, and ORF12 (a novel 120 residue protein). Additionally, they have the CP for pVIII and eE3. TRXPT_11
331 encodes hexon as its 5'-most ORF and also has the CP for eProt, ORF12, e33K, pVIII and eE3. Typically,
332 MLTU transcripts splice the TPL onto a splice site just upstream of the ORF to be expressed (17). While
333 this holds true for most MLTU ORFs, several late ORFs (pVI, protease, and ORF7) do not have such close
334 proximity splicing but are contained in larger transcripts such as these L4 mRNAs, strongly suggesting the
335 use of non-standard ribosomal initiation mechanisms such as secSC usage or ribosome shunting described
336 in other AdVs for their translation (17, 29). TRXPT_13, an E3 ORF utilizing the MLP TSS, encodes the
337 classical L4P genes, pVIII and eE3. Lastly, the L5 class transcript, TRXPT_12, encodes Fiber as its 5'-most
338 ORF but also has the CP for ORF7. Interestingly, the CP of TRXPT_12 and TRXPT_27 of the E3 TU are
339 identical but are initiated from different TSSs.

340 **DISCUSSION**

341 While the advent of next-generation sequencing has rendered easier the study of large and complex
342 eukaryotic transcriptomes, studying the smaller, compact viral transcriptomes is counterintuitively more
343 challenging, as the transcripts typically have significant overlaps due to genome economization. AdV
344 transcriptomes escalate the difficulty due to the wide array of mRNAs produced via very complex alternative
345 splicing and polyadenylation, all initiated from relatively few promoters. Standard RNA-seq analysis programs,
346 not primarily designed for such compact, gene-dense, and complex transcriptomes, further compound this
347 challenge. Furthermore, in our case, no prior transcriptomic studies for THEV exist; hence, assembling the
348 transcripts without any prior experimentally-derived annotation of THEV splicing using only short illumina
349 reads proved difficult. Lastly, we had initially planned to sequence RNA from another time point (8 hpi);
350 however, all the RNA samples from 8 hpi and one replicate sample from 12 hpi got too degraded during
351 the library preparation steps to be yield any useful data. We believe that these would have contributed to
352 better insights into the temporal expression levels and splicing. Our approach combines standard RNA-seq
353 analysis programs with custom analyses and experimentally validating all splice junctions with independent
354 methods. The transcript map for THEV produced from our analysis is strikingly similar to that of the MAdVs.

355 Our work provides the first insights into THEV splicing, revealing 34 transcripts grouped into five transcription
356 units (TUs). The general temporal gene expression regulation observed in MAdVs, with early regions
357 peaking at earlier time points followed by MLTU predominance at later time points, seems to also apply
358 to THEV. An unexpected observation is that the pileup of mapped reads to THEV seems consistently
359 skewed over similar regions of the genome at all time points. Given the temporal regulation of AdVs gene
360 expression, we anticipated distinct differences in read pileups over the genome at different time points,
361 indicating the different stages of infection. This may be due to an unsynchronized infection, leading to
362 transcripts overlapping the time points. Further research is required to determine the precise temporal
363 regulation of THEV.

364 Short read deep sequencing effectively reconstructs full AdV mRNA structures, particularly mapping splice
365 sites (18). However, the substantial overlapping nature of AdV mRNAs and fragmentation during library prepa-
366 ration make it challenging to map the exact TSS, TTS, and PASs of assembled transcripts. Also, As AdVs
367 make heavy use of alternative polyadenylation, short read RNA-seq is ill-equipped to discriminate mRNA
368 variants of the same gene produced via alternative polyadenylation. Thus, shorter variants of alternatively
369 polyadenylated mRNAs may potentially be incorporated into the longer variants during transcript assembly,
370 significantly diminishing the diversity of mRNA in the transcriptome. Also, independent transcripts with
371 significant overlaps may be assembled as a single, longer mRNA, since the short reads alone do not

372 provide enough context for the transcript assembler (StringTie) to distinguish them. Such fusions may affect
373 the transcript expression level estimations by inflating or deflating the expression levels of the transcripts
374 involved, affecting the proper understanding of the temporal gene expression regulation and also the
375 diversity of the transcriptome. Transcripts that have reads mistakenly fused with them would have inflated
376 expression levels while those whose reads are counted elsewhere would show false lower expression levels.
377 In our case, since we used other independent methods to validate the splice junctions, we believe these
378 drawbacks to be minimized. Our results show transcripts in the same TU initiating or terminating in similar
379 areas, but not at the exact same position. We consider the most upstream TSS or most downstream TTS for
380 for the transcripts involved but we present them unchanged in all the figures shown (see **Supplementary**
381 **Table S2**). Also, comparing our results to the better-studied MAdV transcriptomes, we believe some long
382 transcripts in the MLTU (TRXPT_9, TRXPT_10, and TRXPT_11) are likely due to fusing some E3 transcripts
383 to the terminal exons of the MLTU transcripts by StringTie, making them significantly longer. These mRNAs
384 have unusually many exons and their last few exons are identical to some E3 mRNAs. Future studies using
385 long read sequencing technologies will provide more precise mapping of the TSS~~and TTS and~~, TTS, and
386 ~~PAs and~~ clarify the structures of the long MLTU transcripts.

387 While most predicted ORFs are encoded by the spliced transcripts, we found some that seem to be truncated
388 predictions, as either an upstream in-frame SC or unknown upstream exons were found. Other ORFs
389 were identified that were either shorter or longer isoforms of some predicted ORF. We also found several
390 novel unpredicted ORFs (**Supplementary Table S3**). On this basis, we anticipate that further studies
391 will likely reveal more unpredicted novel ORFs or new variants of predicted ORFs. Furthermore, it is not
392 unreasonable to presume that several splice variants will likely be found as evidenced firstly by finding
393 unique transcripts using 3' RACE and during our splice junction validation steps. And secondly, recent
394 studies (17, 18, 22) are still discovering novel mRNA variants for even the best studied MAdVs decades
395 later. These new potential proteins and isoforms significantly extend the repertoire of THEV gene products
396 than predicted, adding a hefty number of proteins of undefined function to the previously predicted five (See
397 **Figure 1**). These new potential proteins and isoforms of unknown functions may mediate or contribute to
398 viral replication efficiency, or the immunosuppression associated with THEV. Hence, further studies of these
399 potential proteins is urgently needed.

400 Eukaryotic mRNAs are typically functionally monocistronic, with the 5'-most AUG determining the translation
401 reading frame. However, AdV mRNAs, which span more than one ORF, are functionally polycistronic,
402 employing non-standard mechanisms of translation initiation such as secSC usage and ribosome shunting
403 (10, 22). AdVs use secondary AUGs as initiation codons for most E1b proteins and some E3 proteins. In

fact, recent studies show that secSC usage is found transcriptome-wide. This is thought to occur because translation initiation at the first SC is inefficient, allowing downstream SCs to be employed (17). Ribosomal shunting or jumping mechanism is utilized for MLTU transcripts that have the TPL. This mechanism allows the ribosome to translocate to a downstream AUG, under the direction of the shunting elements in the TPL, even if a **start codon SC** in a good Kozak sequence context is bypassed. Thus, predicting the protein(s) expressed from an AdV mRNA is uncertain as any one of the AUGs may be selected (10, 22). Almost all the THEV transcripts in our data have the CP for several ORFs, some spanning as many as six ORFs. This supports the usage of these special ribosome initiation mechanisms as several predicted and novel ORFs found on mRNA in our data could not be translated using only the typical ribosome scanning mechanism. Interestingly, several distinct THEV mRNAs have identical CPs. This is also observed in human AdVs in a recent study (17). They proposed that this may permit protein production to be fine-tuned through alteration in the balance between different mRNA groups expressing that ORF.

AdV alternative splicing undergoes a regulated temporal shift in splice site usage, previously thought to be limited to certain TUs. However, recent studies suggest that AdVs routinely produce different combinations of splice acceptor–donor pairs across all TUs (10, 17, 22, 30). The details of this phenomenon have been best studied for the E1A and L1 units. AdVs modulate the activities of the splicing factor U2AF and the cellular SR family of splicing factors (reviewed here (30)), and encode several proteins that influence the RNA splice site used. This phenomenon appears to occur in the THEV transcriptome, as the stringency of splice acceptor-donor pairs selected decreases from the onset of the late phase (**Figure 5**). Recent studies show that a virtually unlimited number of combinatorial alternative splicing events occur in an AdV lytic infection, resulting in a variety of novel transcripts (17, 22). It is unlikely that the entire repertoire of mRNA produced via this mechanism will actually be translated. However, it has been speculated that the plasticity in alternative RNA splicing enables AdVs to fine-tune protein synthesis by providing different alternatively spliced variants encoding the same protein under changing conditions, conferring an evolutionary advantage (17, 22).

CONCLUSIONS

The THEV transcriptome bears remarkable similarity to the better-studied MAdVs. The transcriptome is organized into five TUs, with temporal regulation divided into early and late genes, and a broad repertoire of transcripts are produced via virtually unlimited alternative splicing. However, the THEV transcriptome appears less sophisticated (i.e, it encodes fewer genes) than MAdVs, primarily because the MAdV genomes are close to twice as long as that of THEV. The lack of subdivision of the E1 region into E1a and E1b is one of the most obvious examples. Also, the MAdV E4 region encodes several proteins unlike in THEV

435 where only one transcript encoding one protein was found. The complexity of the MLTU leader sequences is
436 another example. While the majority of the THEV MLTU transcripts begin with the TPL just like MAdVs with
437 a small subset using a variant leader sequence (sTPL), significantly more diverse 5'-UTRs are employed for
438 MAdV MLTU transcripts. Namely, the TPL, the so-called x, y, and z leaders, and the i-leader are 5' leaders
439 utilized by MAdV MLTU mRNAs. The absence of these non-TPL leaders in our data could mean that the
440 5'-UTR diversity of THEV's MLTU mRNAs is more limited due to its smaller genome size or future studies
441 could uncover more variety not seen in our results. We also note that although THEV genomic sequences
442 show minimal differences between strains (12), the transcriptomes may have significant variations; hence,
443 our results may vary from other THEV strains. Also, performing the study *in vivo* or with primary turkey
444 cells may show different results. The potential new proteins identified in our work adds to the number
445 of proteins with undefined functions in THEV; these may have roles in viral replication efficiency, or the
446 immunosuppression associated with THEV. Hence, further studies of these proteins and other predicted
447 proteins of unknown function should be useful in elucidating THEV-induced immunosuppression. Being the
448 first transcriptomic characterization of THEV, this work should serve as useful resource to future THEV
449 gene expression and transcriptomic studies, especially, mapping the mRNA splice sites. More importantly,
450 this work characterizing the splicing of THEV mRNAs will allow researchers to accurately clone any THEV
451 gene(s) of interest to study its potential role in inducing immunosuppression or other functions.

452 **METHODS**

453 **Cell culture and THEV Infection**

454 The Turkey B-cell line (MDTC-RP19, ATCC CRL-8135) was grown as suspension cultures in 1:1 complete
455 Leibovitz's L-15/McCoy's 5A medium with 10% fetal bovine serum (FBS), 20% chicken serum (ChS), 5%
456 tryptose phosphate broth (TPB), and 1% antibiotic solution (100 U/mL Penicillin and 100 μ g/mL Streptomycin),
457 at 41°C in a humidified atmosphere with 5% CO₂. Infected cells were maintained in 1:1 serum-reduced
458 Leibovitz's L15/McCoy's 5A media (SRLM) with 2.5% FBS, 5% ChS, 1.2% TPB, and 1% antibiotic solution. A
459 commercially available THEV vaccine was purchased from Hygieia Biological Labs as a source of THEV-A
460 (VAS strain). The stock virus was titrated using an in-house qPCR assay with titer expressed as genome
461 copy number (GCN)/mL, similar to Mahshoub *et al* (31) with modifications. Cells were infected in triplicate at
462 ~~a multiplicity of infection (MOI) MOI~~ of 100 GCN/cell, incubated at 41°C for 1 hour, and washed three times
463 with phosphate buffered saline (PBS) to get rid of free virus particles. Triplicate samples were harvested at
464 4-, 12-, 24-, and 72-hpi for total RNA extraction. The infection was repeated but samples in triplicate were
465 harvested at 12-, 24-, 36-, 48-, and 72-hpi for PCR validation of novel splice sites. Still one more independent
466 infection was done at time points ranging from 12 to 168-hpi for qPCR quantification of virus titers.

467 **RNA extraction and Sequencing**

468 Total RNA was extracted from infected cells using the Thermo Fisher RNAqueous™-4PCR Total RNA Isolation
469 Kit (which includes a DNase I digestion step) per manufacturer's instructions. An agarose gel electrophoresis
470 was performed to check RNA integrity. The RNA quantity and purity was initially assessed using nanodrop,
471 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.
472 Extracted total RNA samples were sent to LC Sciences, Houston TX for poly-A-tailed mRNA sequencing
473 where RNA integrity was checked with Agilent Technologies 2100 Bioanalyzer High Sensitivity DNA Chip
474 and poly(A) RNA-seq library was prepared following Illumina's TruSeq-stranded-mRNA sample preparation
475 protocol. Paired-end sequencing to generate 150 bp reads was performed on the Illumina NovaSeq 6000
476 sequencing system.

477 **Validation of Novel Splice Junctions**

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

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

496 **3' Rapid Amplification of cDNA Ends (3' RACE)**

497 A rapid amplification of sequences from the 3' ends of mRNAs (3' RACE) experiment was performed using
498 a portion of the extracted total RNA of infected MDTC-RP19 cells used for the RNA-seq experiment as
499 explained above. We followed the protocol described by Green *et al* (32) with modifications. Briefly, 1 μ g of
500 total RNA was reverse transcribed to cDNA using SuperScript™ IV First-Strand Synthesis System following
501 the manufacturing instructions using an adapter-primer with a 3'-end poly(T) and a 5'-end BamHI restriction
502 site. A gene-specific sense primer with a 5'-end KpnI restriction site paired with an anti-sense adapter-primer
503 with a 5'-end BamHI site were used to amplify target sections of the cDNA using Invitrogen's Platinum™ Taq
504 DNA polymerase High Fidelity, following manufacturer's instructions. The PCR amplicons were restriction
505 digested, cloned, and Sanger sequenced.

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

507 Sequencing reads were analyzed following a well-established protocol described by Pertea *et al* (25), using
508 Snakemake - version 7.24.0 (33), a popular workflow management system to drive the pipeline. Briefly,
509 sequencing reads were trimmed with the Trim-galore - version 0.6.6 (34) program to achieve an overall Mean
510 Sequence Quality (Phred Score) of 36. Trimmed reads were mapped simultaneously to the complete genomic

511 sequence of avirulent turkey hemorrhagic enteritis virus THEV (<https://www.ncbi.nlm.nih.gov/nuccore/AY849321.1/>) and *Meleagris gallopavo* (<https://www.ncbi.nlm.nih.gov/genome/?term=Meleagris+gallopavo>)
512 using Hisat2 - version 2.2.1 (25) with default settings. The generated binary alignment (BAM) files from
513 each infection time point were filtered for reads mapping to the THEV genome using Samtools - version
514 1.16.1 and fed into StringTie - version 2.2.1 (25) to assemble the transcripts, using a gene transfer format
515 (GTF) annotation file derived from a gene feature format 3 (GFF3) annotation file obtained from NCBI, which
516 contains the predicted ORFs of THEV as a guide. GFFCOMPARE - version 0.12.6 was used to merge all
517 transcripts from all time points without redundancy and using a custom R script, adenovirus transcripts
518 units (regions) were assigned to each transcript, generating the transcriptome of THEV. StringTie was set
519 to expression estimation mode to calculate FPKM scores for all transcripts after which Ballgown - version
520 2.33.0 in R was used to perform the statistical analysis on the transcript expression levels. Samtools was
521 also used to count the total sequencing reads for all replicates at each time point and Regtools - version
522 1.0.0 was used to count all junctions, the reads supporting them, and extract all other information related to
523 the junction. See **Supplementary Computational Analysis** for the details of transcript expression level
524 estimations and splice junction read counts.

526 **LIST OF ABBREVIATIONS**

Abbreviation	Definition
AdV	Adenovirus
MAdV	Mastadenovirus
SiAdV	Siadenovirus
THEV	Turkey Hemorrhagic Enteritis Virus
HE	Hemorrhagic Enteritis
UTR	Untranslated Region
TU	Transcription Unit
L4P	L4 Promoter
MLP	Major Late Promter
E3P	E3 Promoter
hpi	Hours Post-infection
qPCR	Quantitative Polymerase Chain Reaction
FPKM	Fragments Per Kilobase of transcript per Million mapped reads
CP	Coding Potential
TSS	Transcription Start Site
TTS	Transcription Termination Site
SC	Start Codon

Abbreviation	Definition
STC	Stop Codon
secSC	Secondary Start Codon
secSTC	Secondary Stop Codon
ORF	Open Reading Frame
CDS	Coding Sequence
MLTU	Major Late Transcription Unit
TPL	Tripartite Leader
sTPL	Short Tripartite Leader
TPL3	Third exon of Tripartite Leader
GCN	Genome Copy Number

527 **DATA AVAILABILITY**

528 The raw sequencing read data (FastQ), transcript expression counts, and total unique junctions have
 529 been deposited at the National Center for Biotechnology Information Gene Expression Omnibus (<http://www.ncbi.nlm.nih.gov/geo>) under accession number GSE254416.
 530
 531 Data is also available on request by contacting the designated corresponding author.

532 **CODE AVAILABILITY**

533 All the code/scripts in the entire analysis pipeline are available on github (https://github.com/Abraham-Quaye/thev_transcriptome)

535 **ACKNOWLEDGMENTS**

536 We thank the Office of Research Computing at Brigham Young University for granting us access to the **high**
 537 **performance** **high-performance** computing systems to perform the memory-intensive steps in the analysis

⁵³⁸ pipeline of this work.

539 **REFERENCES**

- 540 1. Davison A, Benko M, Harrach B. 2003. Genetic content and evolution of adenoviruses. *The Journal of general virology* 84:2895–908.
- 541 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.
- 542 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.
- 543 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.
- 544 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.
- 545 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.
- 546 7. Davison AJ, Wright KM, Harrach B. 2000. DNA sequence of frog adenovirus. *J Gen Virol* 81:2431–2439.
- 547 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.
- 548 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.

- 549 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.
- 550 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.
- 551 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.
- 552 13. Gross WB, Moore WE. 1967. Hemorrhagic enteritis of turkeys. Avian Dis 11:296–307.
- 553 14. Rautenschlein S, Sharma JM. 2000. Immunopathogenesis of haemorrhagic enteritis virus (HEV) in
turkeys. Dev Comp Immunol 24:237–46.
- 554 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.
- 555 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.
- 556 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.
- 557 18. Zhao H, Chen M, Pettersson U. 2014. A new look at adenovirus splicing. Virology 456-457:329–341.
- 558 19. Wolfrum N, Greber UF. 2013. Adenovirus signalling in entry. Cell Microbiol 15:53–62.
- 559 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.

- 560 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.
- 561 22. Westergren Jakobsson A, Segerman B, Wallerman O, Bergström Lind S, Zhao H, Rubin C-J, Pettersson U, Akusjärvi G. 2021. The human adenovirus 2 transcriptome: An amazing complexity of alternatively spliced mRNAs. *Journal of Virology* 95.
- 562 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.
- 563 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.
- 564 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.
- 565 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.
- 566 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.

- 567 28. Beaudoin E, Freier S, Wyatt JR, Claverie J-M, Gautheret D. 2000. Patterns of variant polyadenylation
signal usage in human genes. *Genome Research* 10:1001–1010.
- 568 29. Yueh A, Schneider RJ. 1996. Selective translation initiation by ribosome jumping in adenovirus-infected
and heat-shocked cells. *Genes & Development* 10:1557–1567.
- 569 30. Akusjarvi G. 2008. Temporal regulation of adenovirus major late alternative RNA splicing. *Frontiers in
Bioscience Volume:5006*.
- 570 31. 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.
- 571 32. Green MR, Sambrook J. 2019. Rapid amplification of sequences from the 3' ends of mRNAs: 3'-RACE.
Cold Spring Harbor Protocols 2019:pdb.prot095216.
- 572 33. 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.
- 573 34. 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.

574 **TABLES AND FIGURES**

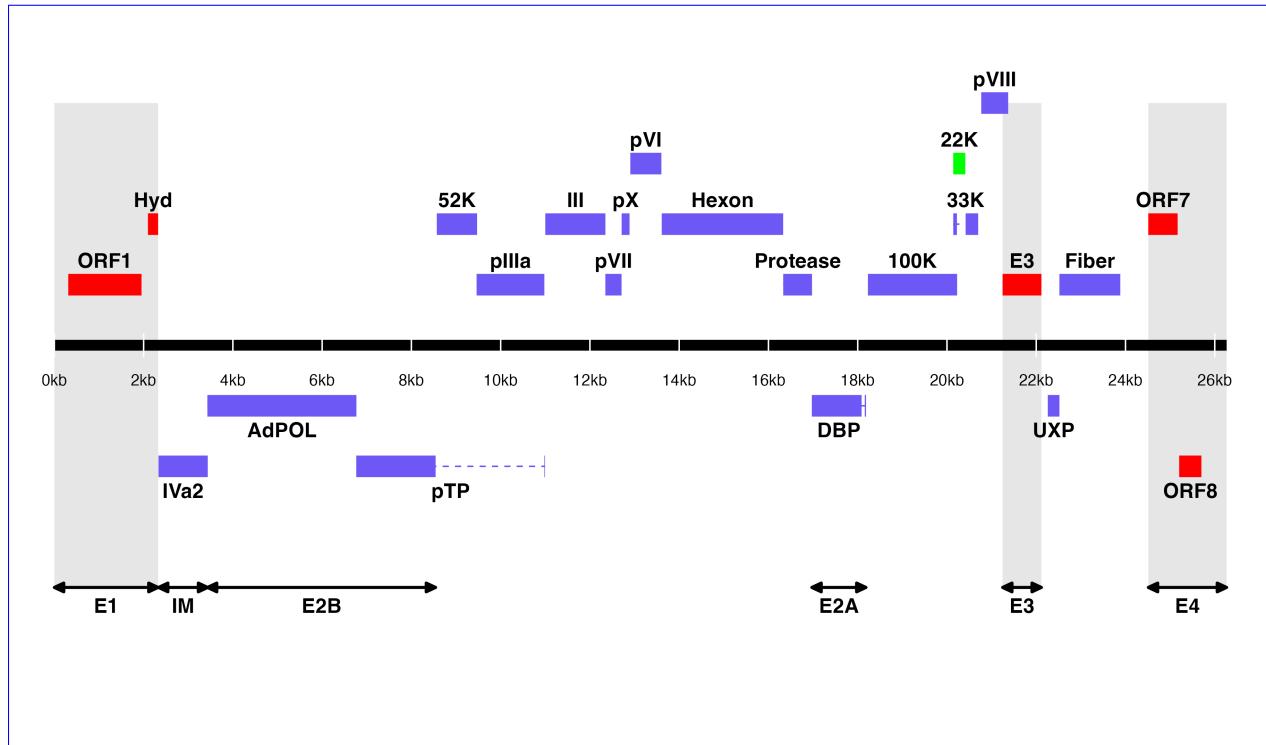


Figure 1. Predicted ORF map of THEV virulent strain. The central horizontal line represents the double-stranded DNA marked at 2kb intervals as white line breaks. Colored blocks represent viral genes. Blocks above the DNA line are transcribed on the sense DNA strand and those below, on the anti-sense 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 “the five genus-specific” genes of undefined functions (colored red). Genes colored in blue are “genus-common”. The gene colored in light green is conserved in all but Adenoviruses. Regions comprising the different transcription units-TUs are labelled at the bottom (E1, E2A, E2B, E3, E4, and IM); the unlabeled regions comprise the MLTU.

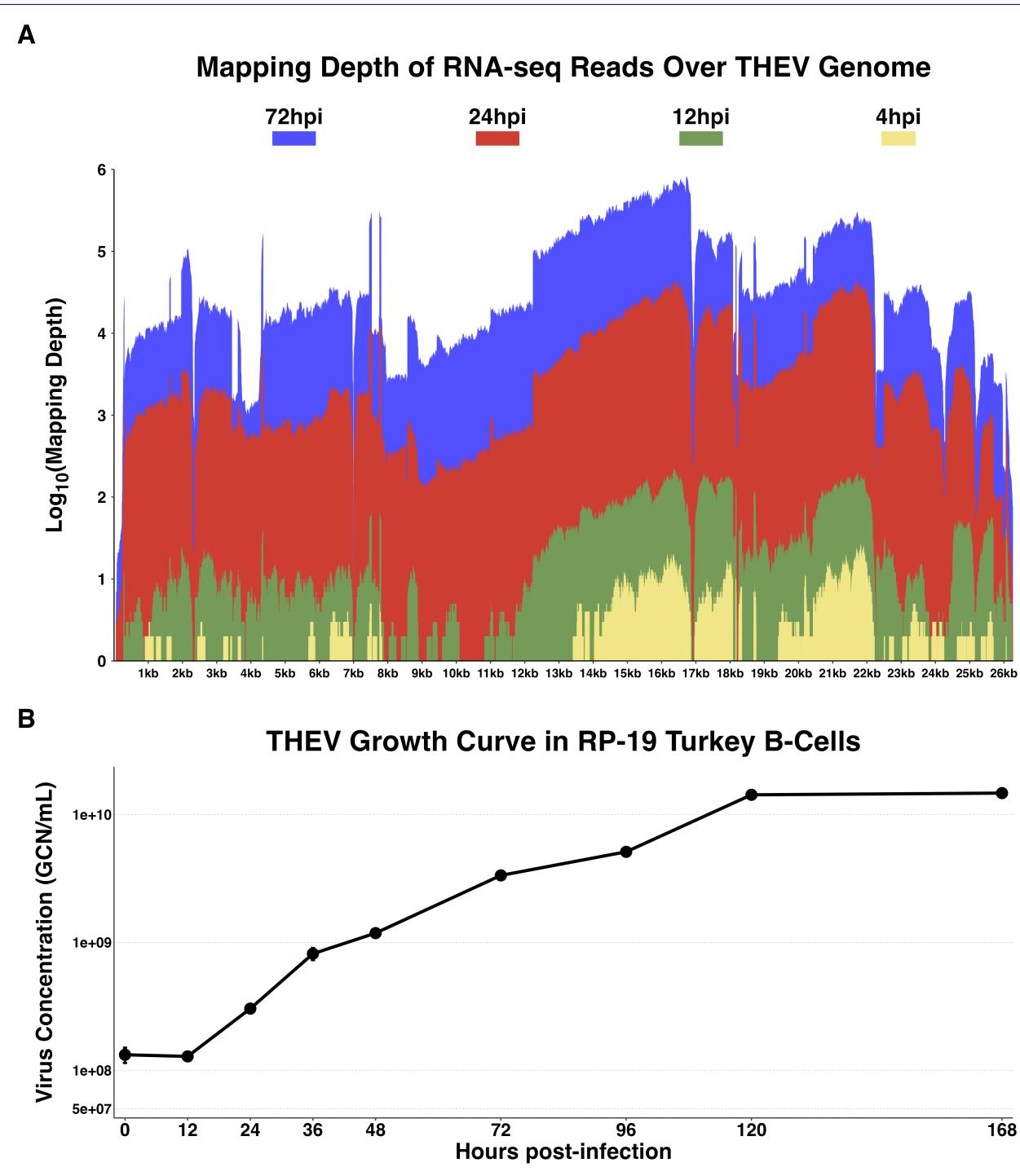


Figure 2: Increasing levels of THEV over time. **A) Per base coverage of sequence reads mapping to THEV genome by time point.** The pileup of mRNA reads mapping to THEV genome at the base-pair level for each indicated time point. **B) Growth curve of THEV (VAS vaccine strain) in MDTC-RP19 cell line.** Virus quantities 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.

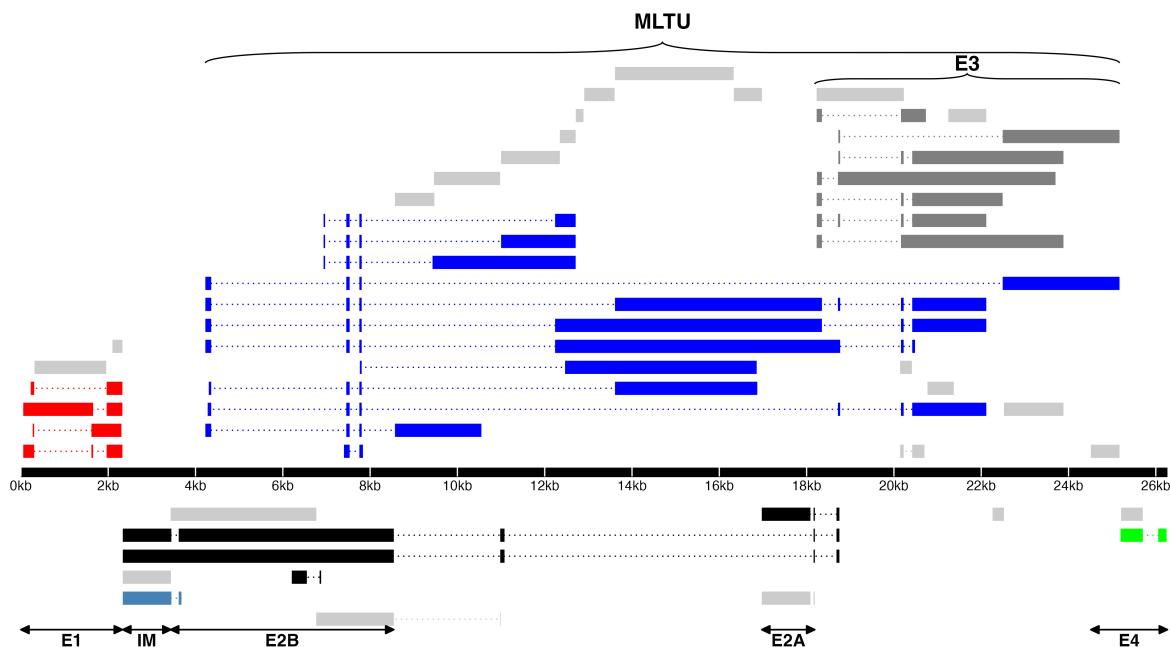
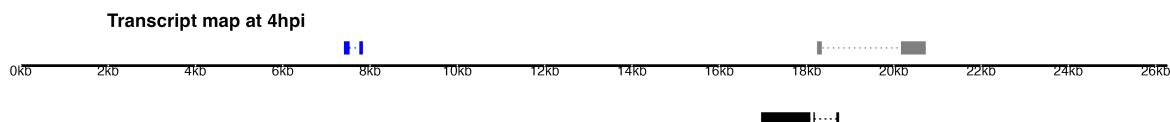
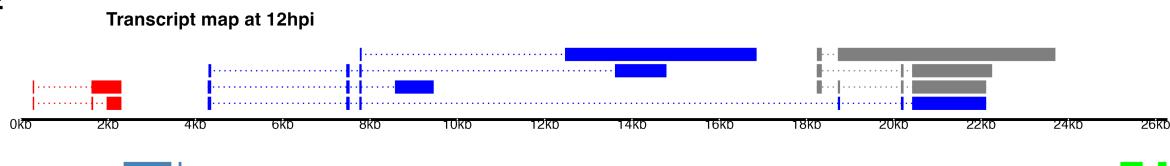
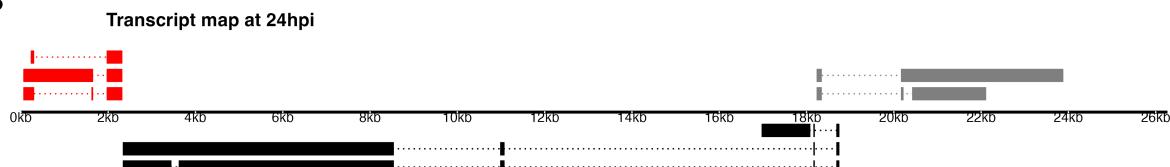
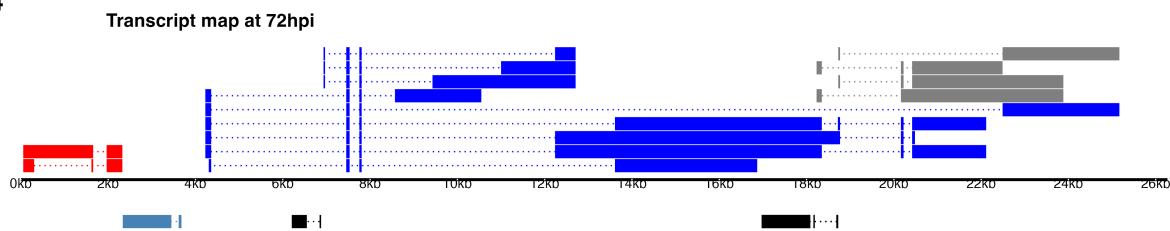
A**1****2****3****4**

Figure 3. A) Transcriptome of THEV from RNA-seq. THEV transcripts assembled from all time points by StringTie are unified forming this 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).

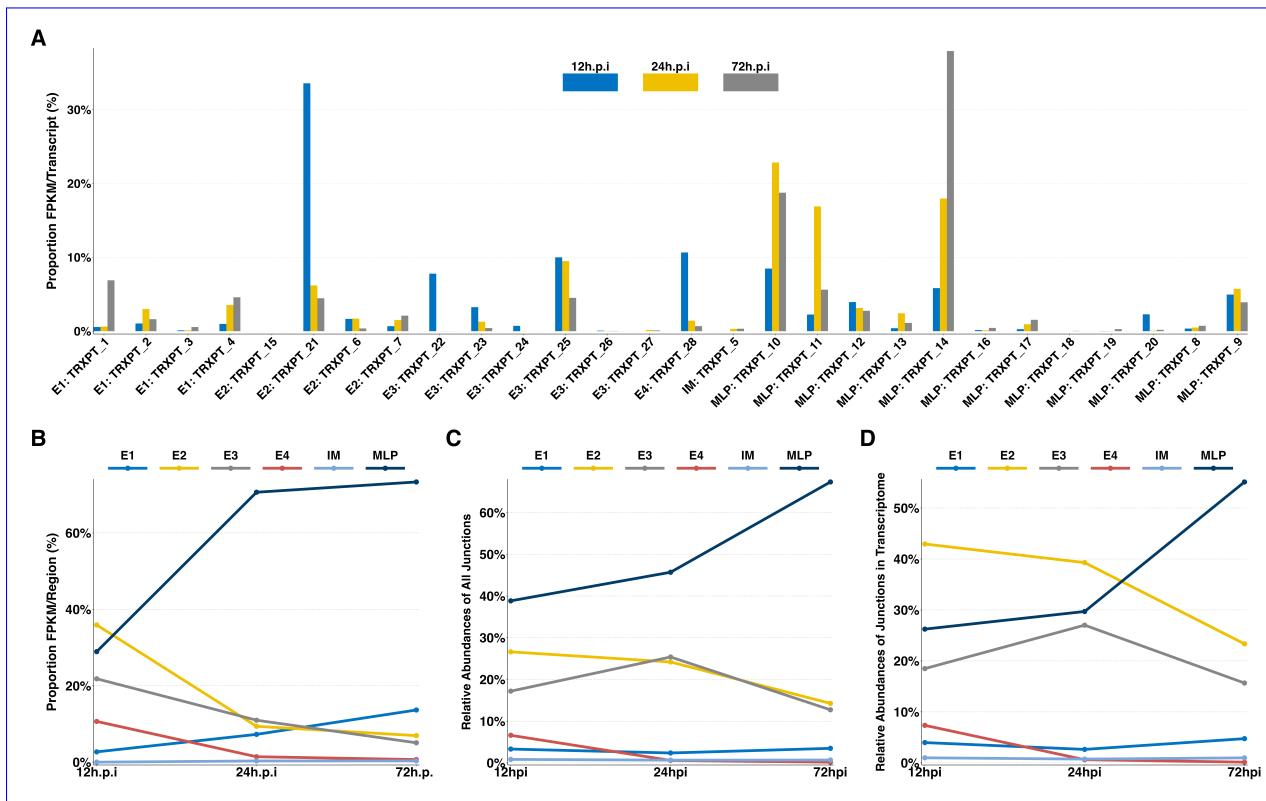


Figure 4: Changes in splicing and expression profile of THEV over time. **A) Normalized (FPKM) expression levels of transcripts over time.** The expression levels (FPKM) of individual transcripts as a percentage of the total expression of all transcripts at each time point are indicated. Only transcripts 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.

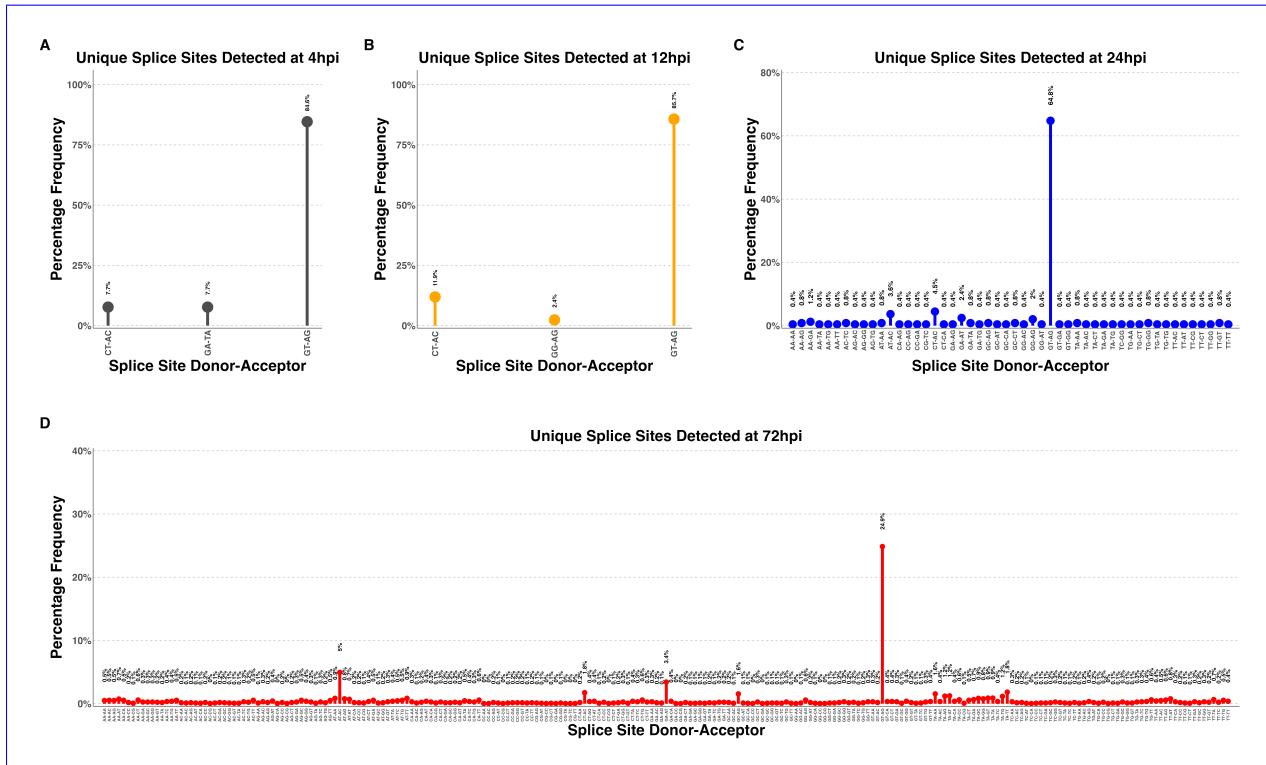


Figure 5: Changes in splice donor-acceptor nucleotides over time. The splice donor-acceptor nucleotides of THEV just like other AdVs is mostly the canonical GU-AG. At early time points (4h.p.i and 12h.p.i **(A)** and **(B)**, respectively) the junction nucleotides used appear to be well scrutinized or 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.

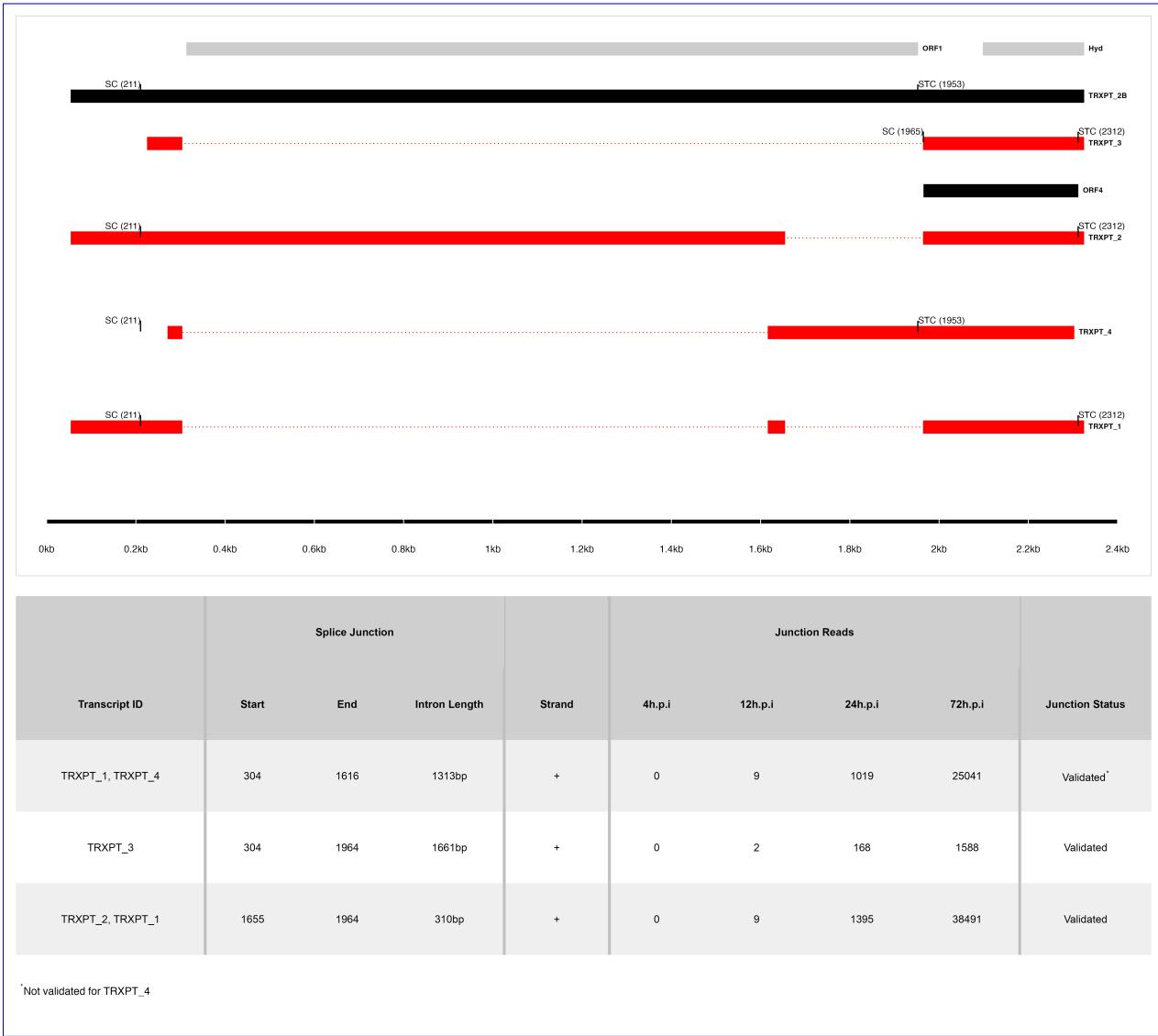


Figure 6: The splice map of the E1 transcription unit (TU). Exons are depicted as boxes connected by introns (dotted lines). Transcripts from RNA-seq data are colored red, predicted ORFs are colored grey, and transcripts or ORFs discovered by other means are colored black. Each transcript or ORF is labelled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

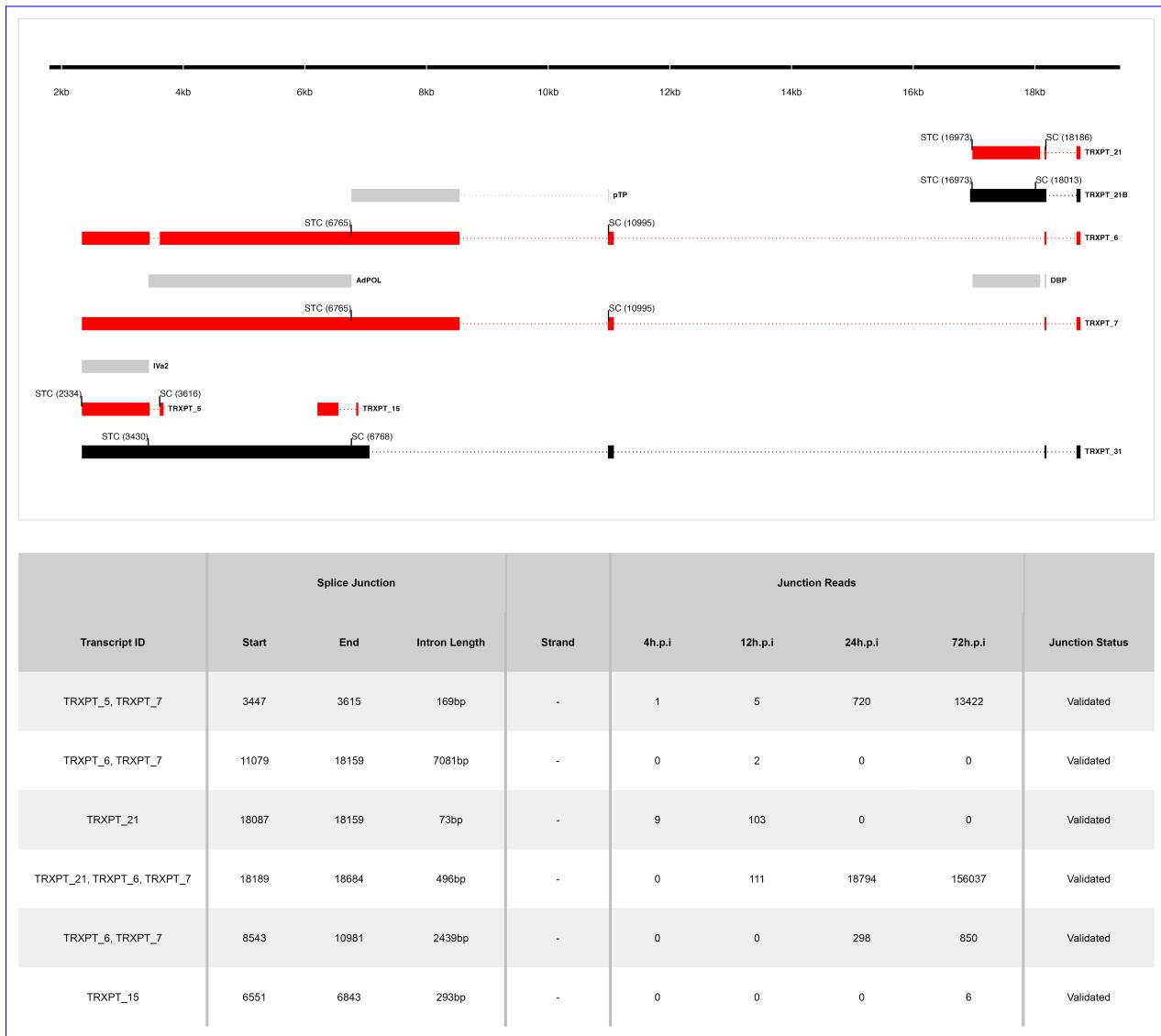


Figure 7: The splice map of the E2 and IM TUs. Exons are depicted as boxes connected by introns (dotted lines). Red transcripts are generated from RNA-seq data and predicted ORFs are colored grey. TRXPT_21B discovered by 3'RACE is colored black. Each transcript or ORF is labelled with its name to the 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.

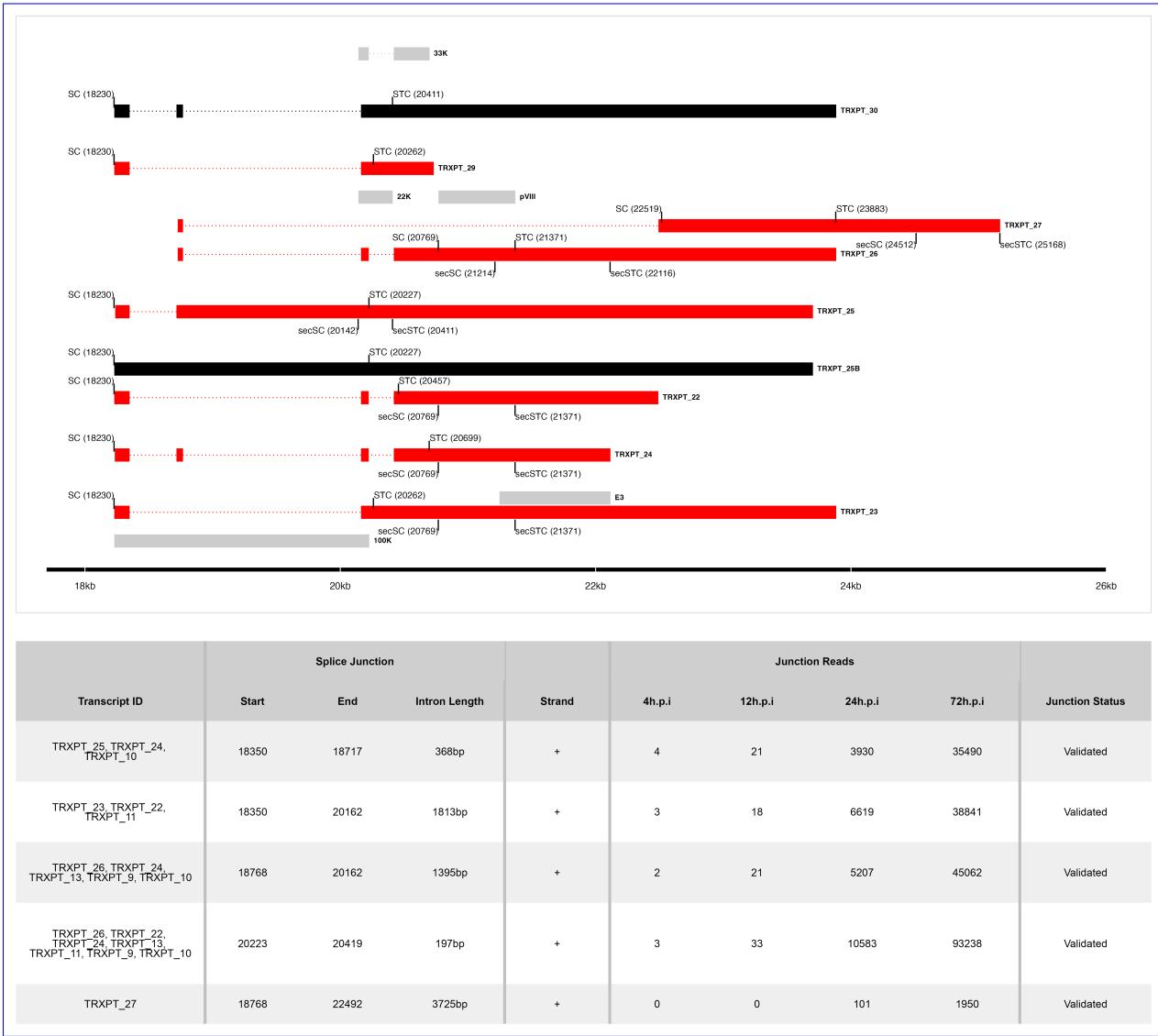


Figure 8: The splice map of the E3 TU. Exons are depicted as boxes connected by introns (dotted lines). Red transcripts are generated from RNA-seq data and predicted ORFs are colored grey. Transcripts discovered by other means are colored black. Each transcript or ORF is labelled with its name to the 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.

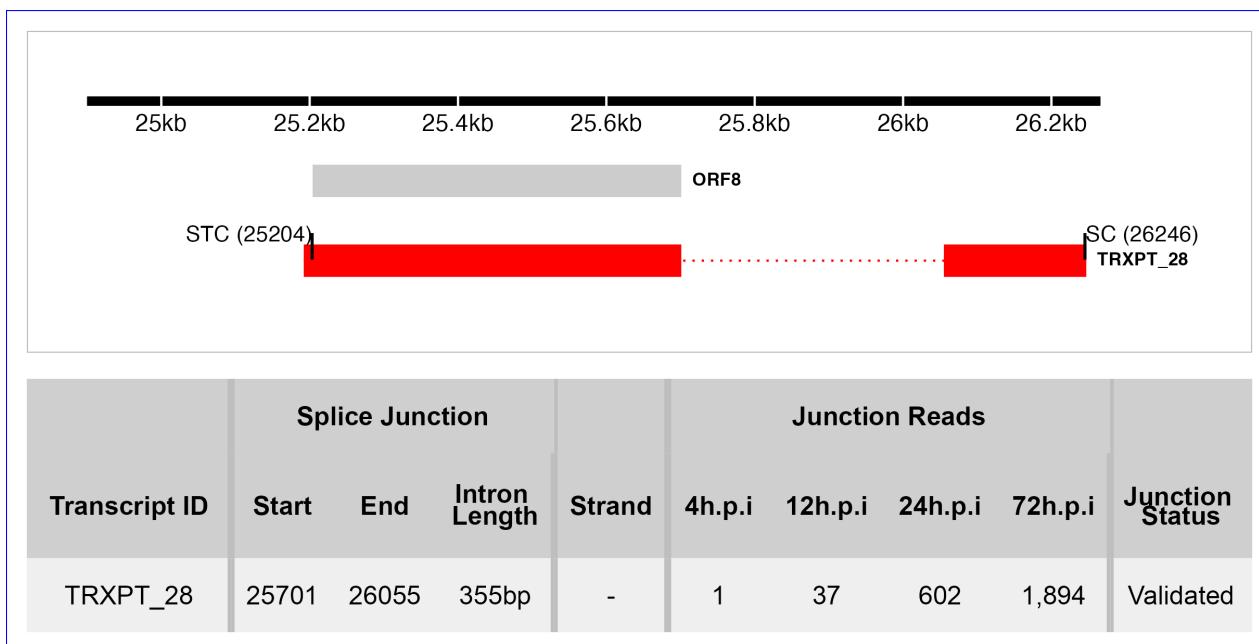


Figure 9: The splice map of the E4 TU. Exons are depicted as boxes connected by introns (dotted lines). The transcript from RNA-seq data is colored red and the predicted ORF, grey. The transcript and ORF are labelled with their names to the right. The **start codon** (SC) and **stop codon** (STC) of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

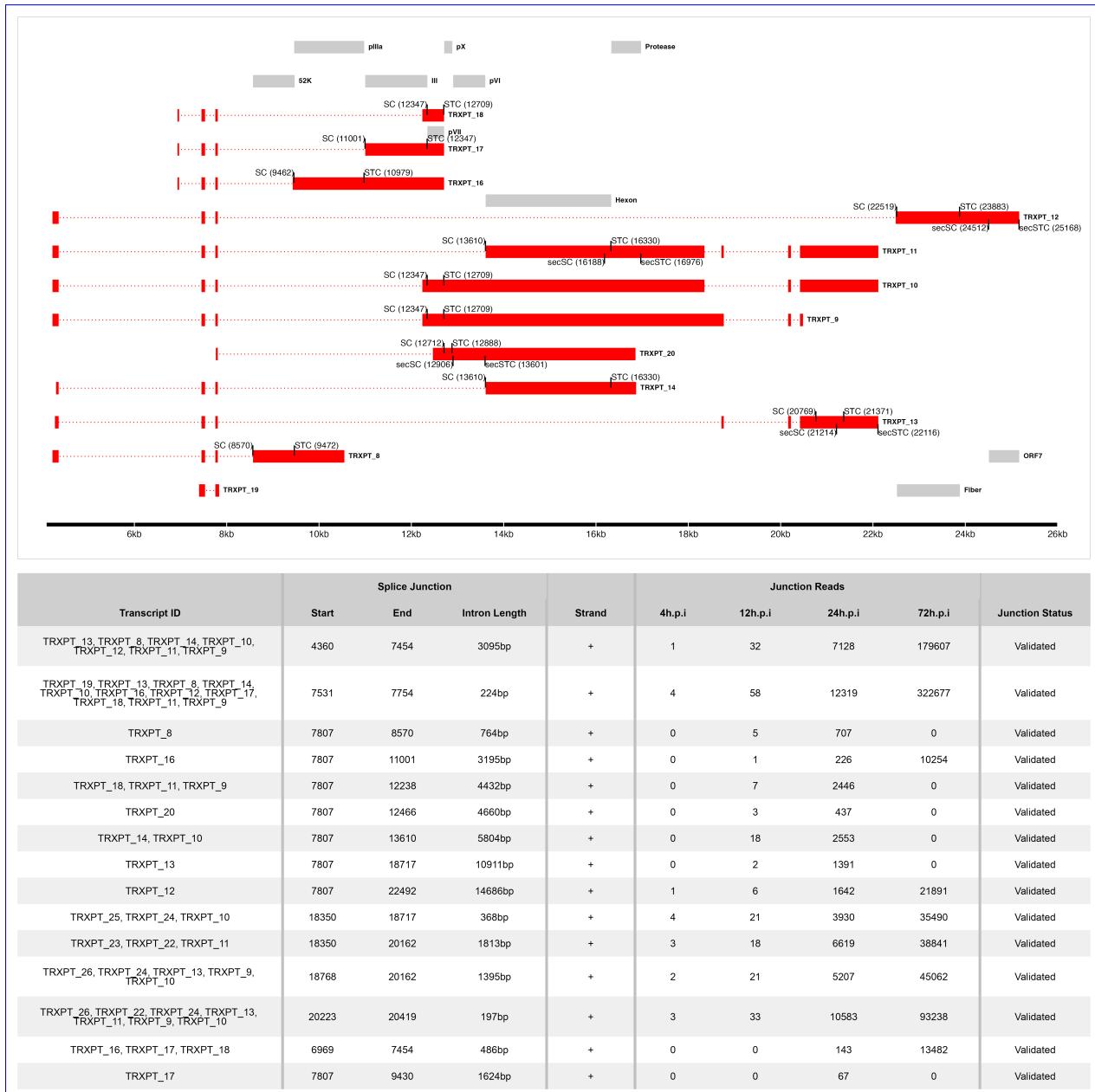


Figure 10: The splice map of the MLTU. Exons are depicted as boxes connected by introns (dotted lines). The transcripts from our RNA-seq data are colored red and the predicted ORFs, grey. The transcripts and ORFs are labelled with their names to the right. The **start codon** (SC) and **stop codon** (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
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 2a: Most abundant splice junctions at 12h.p.i

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

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

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

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

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

575 **SUPPLEMENTARY MATERIALS**

576 **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%

577 **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%

⁵⁷⁸ **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%

⁵⁷⁹ **Supplementary Table S2**

Table S2: Mapping Transcript TTS with Closest Polyadenylation Signal

Region	Sub Region	Transcripts	Nearest 'AAUAAA' TTS	
			TTS	signal (location refers to 'U' in the sequence)
E1	E1	TRXPT_1, TRXPT_2, TRXPT_3, TRXPT_4	2,325	2,323
E2	E2A	TRXPT_21, TRXPT_21B	16,934	16,950
E2/IM	E2B/IM	TRXPT_5, TRXPT_6, TRXPT_7, TRXPT_31	2,334	2,333
E3	E3A	TRXPT_22, TRXPT_24	22,491	22,458 (AGUAAA)

Table S2: Mapping Transcript TTS with Closest Polyadenylation Signal

Region	Sub Region	Transcripts	TTS	Nearest 'AAUAAA' TTS signal (location refers to 'U' in the sequence)
E3	E3B	TRXPT_23, TRXPT_25, TRXPT_25B, TRXPT_26, TRXPT_30	23,884	23,889
E3	E3C	TRXPT_27	25,168	25,166
E4	E4	TRXPT_28	25,192	25,203
MLTU	L1	TRXPT_8	10,549	10,537
MLTU	L2	TRXPT_16, TRXPT_17	12,709	-
MLTU	L3	TRXPT_14, TRXPT_20	16,870	16,903
MLTU	L4	TRXPT_9, TRXPT_10, TRXPT_11, TRXPT_13	22,116	22,098 (AAGAAA)
MLTU	L5	TRXPT_12	25,168	25,166

580 **Supplementary Table S3**

Table S3: The Encoded Proteins and Coding Potential of THEV Transcripts

mRNA	Primary Product	Size (kDa)	Number of Residues	Start Codon Position	Stop Codon Position	Downstream Coding Potential
TRXPT_1	ORF9 ^{new}	17.9	160	211	2312	Hyd, ORF4
TRXPT_2	ORF10 ^{new}	66.4	597	211	2312	Hyd, ORF4
TRXPT_2B	eORF1 ^L	64.3	580	1965	1953	Hyd, ORF4
TRXPT_3	ORF4	13.1	115	211	2312	Hyd

Table S3: The Encoded Proteins and Coding Potential of THEV Transcripts

mRNA	Primary Product	Size (kDa)	Number of Residues	Start Codon Position	Stop Codon Position	Downstream Coding Potential
TRXPT_4	ORF11 ^{new}	15.9	143	211	1953	Hyd, ORF4
TRXPT_5	IVa2	42.3	371	3616	2334	-
TRXPT_6	pTP	70.5	597	10995	6765	Ad-pol
TRXPT_7	pTP	70.5	597	10995	6765	Ad-pol
TRXPT_8	52K	33.8	300	8570	9472	-
TRXPT_9	pVII	13.2	120	12347	12709	pX, pVI, hexon, eProt, ORF12, pVIII, eE3
TRXPT_10	pVII	13.2	120	12347	12709	pX, pVI, hexon, eProt, ORF12, pVIII, eE3
TRXPT_11	Hexon (II)	101.1	906	13610	16330	eProt, ORF12, e33K, pVIII, eE3
TRXPT_12	Fiber (IV)	49	454	22519	23883	ORF7
TRXPT_13	pVIII	21.8	200	20769	21371	eE3
TRXPT_14	Hexon (II)	101.1	906	13610	16330	-
TRXPT_15	-	-	-	-	-	-
TRXPT_16	pIIIa	57.6	505	9462	10979	Penton (III), pVII
TRXPT_17	Penton (III)	50.8	448	11001	12347	-
TRXPT_18	pVII	13.2	120	12347	12709	-
TRXPT_19	-	-	-	-	-	-
TRXPT_20	pX	6.1	58	12712	12888	pVI, hexon

Table S3: The Encoded Proteins and Coding Potential of THEV Transcripts

mRNA	Primary Product	Size (kDa)	Number of Residues	Start Codon Position	Stop Codon Position	Downstream Coding Potential
TRXPT_21	DBP	43.3	380	18186	16973	-
TRXPT_21B	tDBP ^S	39.3	346	18013	16973	-
TRXPT_22	8.3KII ^{new}	8.6	73	18230	20457	pVIII, eE3
TRXPT_23	8.3KI ^{new}	8.3	73	18230	20262	pVIII, eE3, Fiber
TRXPT_24	e33K ^L	19.8	171	18230	20699	pVIII, E3
TRXPT_25	t100K	62.4	543	18230	20227	22K, pVIII, eE3
TRXPT_25B	100K	76.5	665	18230	20227	22K, pVIII, eE3
TRXPT_26	pVIII	21.8	200	20769	21371	eE3, Fiber
TRXPT_27	Fiber (IV)	49	454	22519	23883	ORF7
TRXPT_28	eORF8 ^L	26.4	229	26246	25204	-
TRXPT_29	8.3KI	8.3	73	18230	20262	-
TRXPT_30	e22K ^L	15.7	140	18230	20411	pVIII, eE3, Fiber
TRXPT_31	Ad-pol	129.2	1112	6768	3430	-

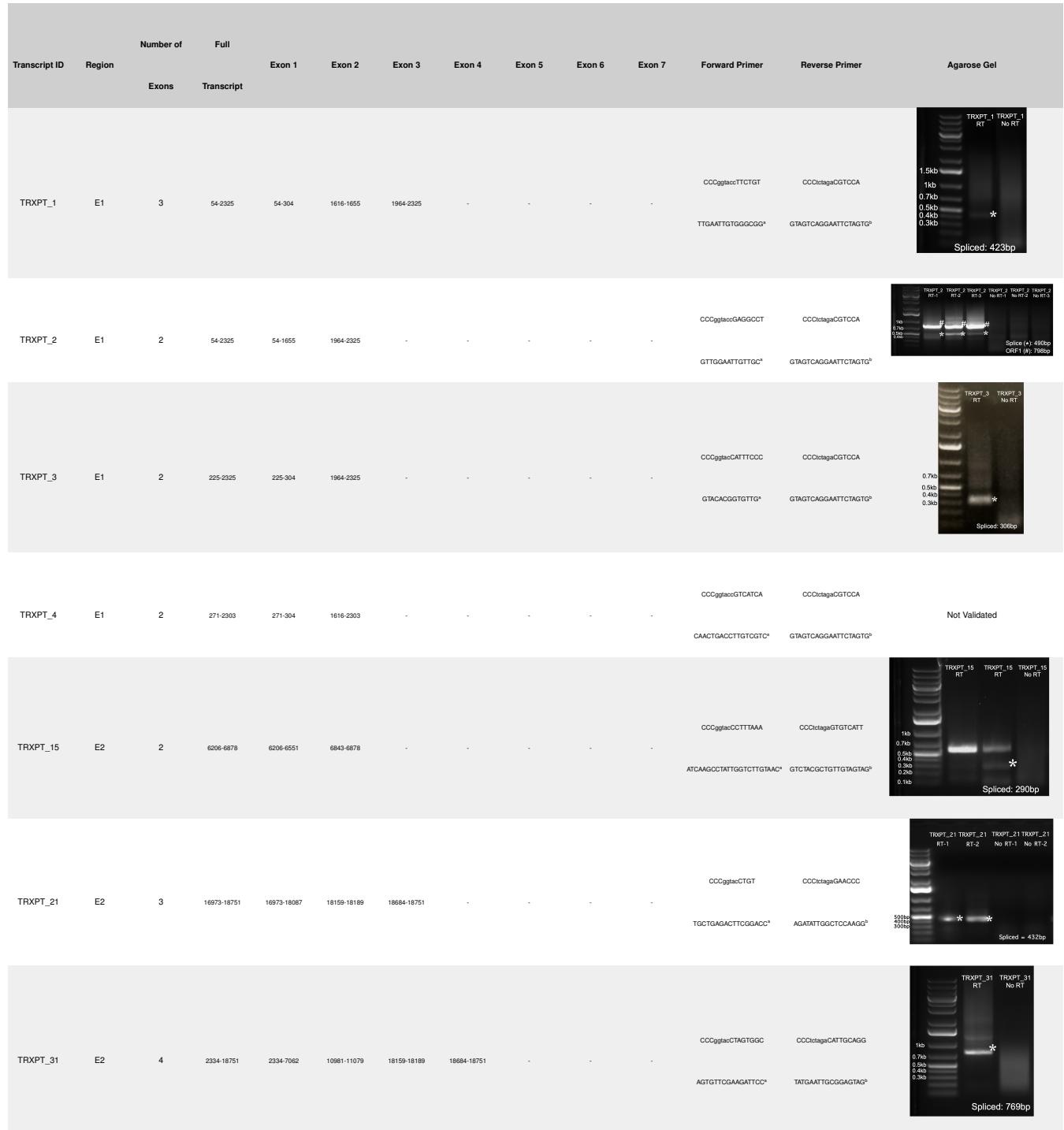
^{new}Novel product identified in this work;

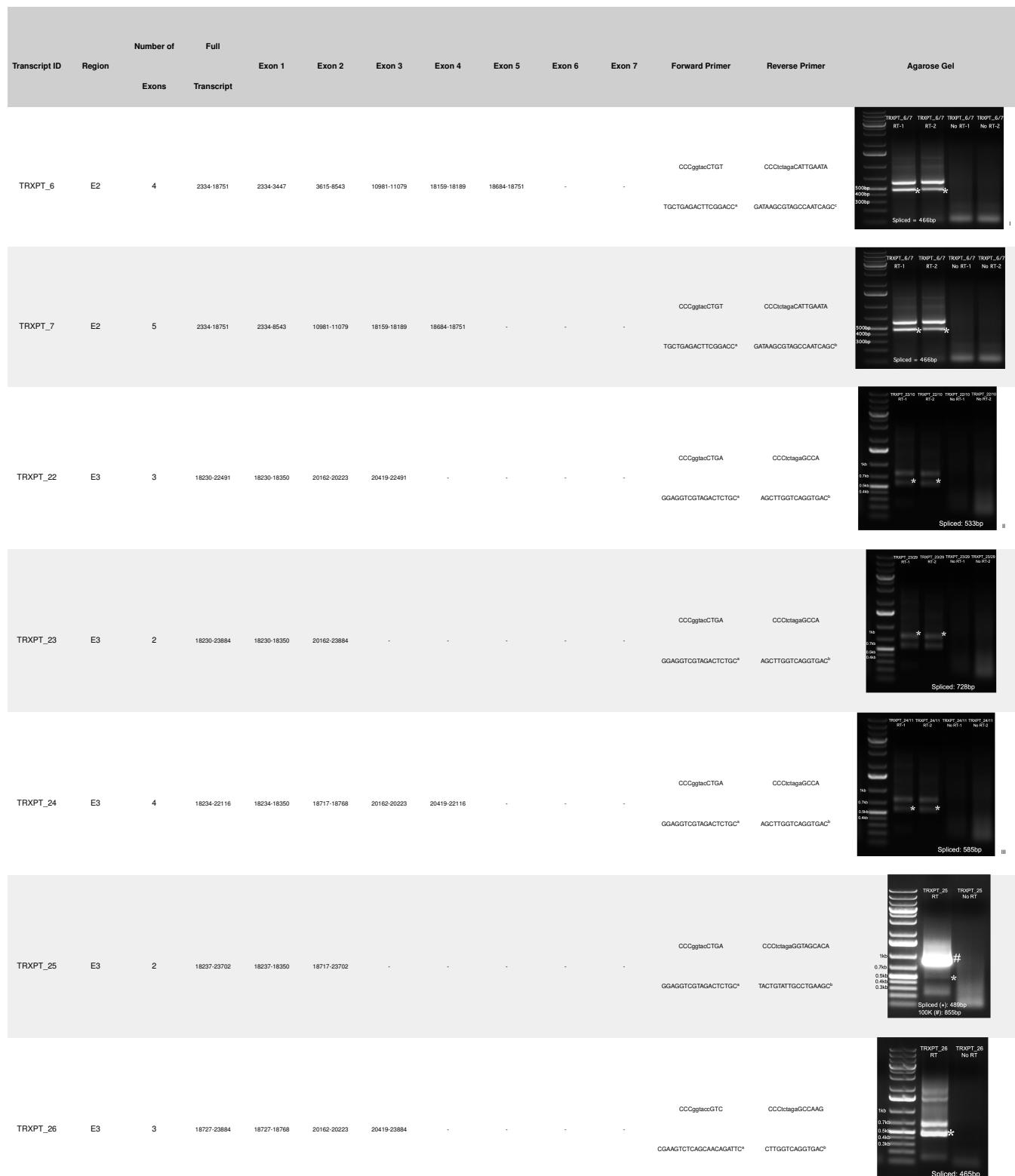
^LLonger isoform identified in this work

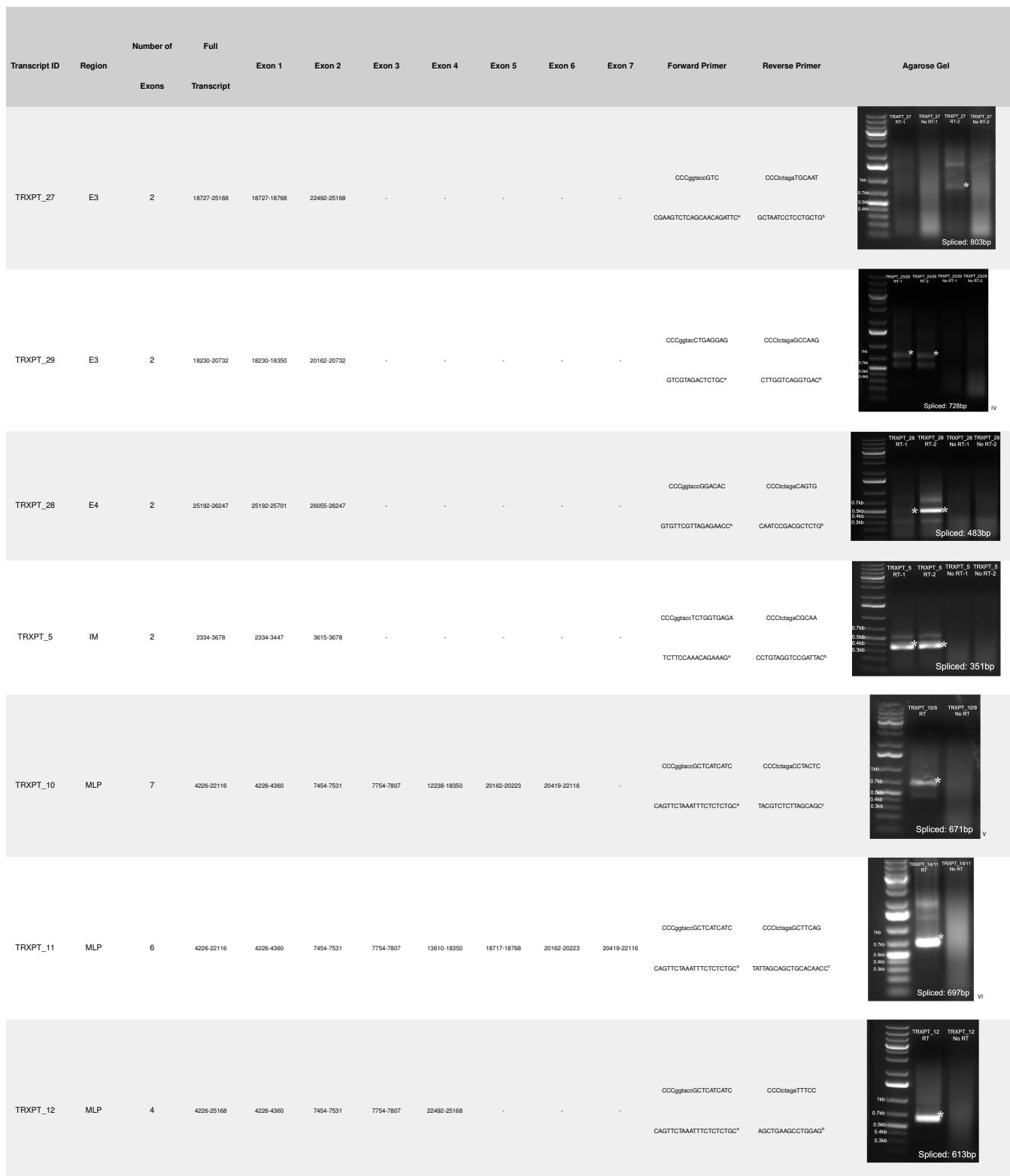
^SShorter isoform identified in this work

581 **Supplementary PCR Methods**

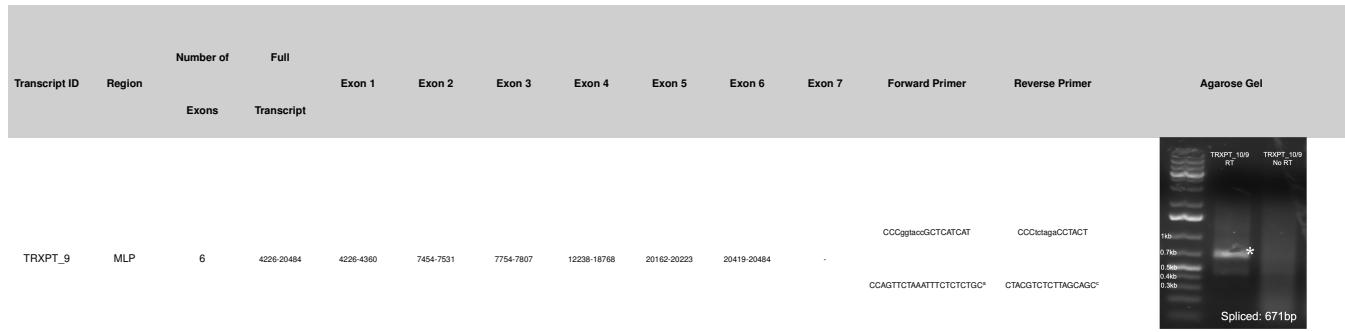
Agarose Gels Showing PCR Amplification of THEV cDNA With Gene-Specific Primers







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_13	MLP	6	4279-22116	4279-4360	7454-7531	7754-7807	18717-18768	20162-20223	20419-22116	CCGgttaccGCTCATCATC	CCCtctagaGCCAAAG	
TRXPT_14	MLP	4	4304-16870	4304-4360	7454-7531	7754-7807	13610-16870	-	-	CCGgttaccGCTCATCATC	CCCtctagaGCCCTCAGT	
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	



^aPrimer binds inside first exon; ^bPrimer binds inside terminal exon; ^cPrimer binds inside fourth exon; ^dAgarose gel identical to TRXPT_7 due to identical splicing; ^eAgarose gel identical to last 3 exons of TRXPT_10 due to identical splicing; ^fAgarose gel identical to last 4 exons of TRXPT_11 due to identical splicing; ^gAgarose gel identical to TRXPT_23 due to identical splicing; ^hAgarose gel identical to TRXPT_9 due to identical splicing; ⁱAgarose gel identical to TRXPT_14 due to identical splicing;

582 In the table above, the restriction sites in the primer tails are shown in lowercase letters. All the primer melting
 583 temperatures (TMs) are 58-60°C using a hot start Taq DNA polymerase. The PCR reaction mix was made
 584 per manufacturer's instructions. The PCR cycling conditions were as follows: Initial denaturation – 95°C for
 585 1 minute; cyclical denaturation – 95°C for 30 seconds, annealing – variable temperature (53°C-56°C) for
 586 30 seconds, primer extension – 68°C for variable time, and final elongation – 68°C for 5 minutes. We used
 587 35 cycles of amplification. The PCR products of THEV cDNA shown in the gels in the table above were
 588 cloned and Sanger sequenced to confirm the splice junctions. We included "No reverse transcriptase (No
 589 RT)" template controls to ensure the PCR products are of RNA origin.

590 Supplementary Computational Analysis

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

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

602 In these steps above, each sample (replicate of each time point) was processed independently and merged
 603 only in the final transcriptome assembly or during analysis with Ballgown. In the subsequent steps described

604 below, all samples for each time point were processed together.

605 We used RegTools to extract and analyze the splice junctions in the BAM files. The command “regtools
606 junctions extract” provides a wealth of information about all the splice sites in the BAM file provided such as:
607 the start and end positions, the strand, and number of reads supporting the splice junctions. The command
608 “regtools junctions annotate” gives even more information such as: the splice site donor-acceptor sequences
609 and transcripts/genes that overlap the junction. This information was the basis for estimating and comparing
610 the splicing activity of different regions (TUs) of THEV over time. Also, Samtools was also used to count the
611 total sequencing reads for all replicates at each time point.



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