

# 1 The African Swine Fever Virus Transcriptome

2

3 Gwenny Cackett<sup>1</sup>, Dorota Matelska<sup>1</sup>, Michal Sýkora<sup>1,3</sup>, Raquel Portugal<sup>2</sup>, Michal Malecki<sup>4</sup>, Jürg4 Bähler<sup>1,5</sup>, Linda Dixon<sup>2</sup> and Finn Werner<sup>\*1</sup>5 <sup>1</sup> Institute for Structural and Molecular Biology, Darwin Building, University College London, Gower

6 Street, London WC1E 6BT, United Kingdom

7 <sup>2</sup> Pirbright Institute, Ash Road, Pirbright, Surrey, GU24 0NF, UK8 <sup>3</sup> Institute of Molecular Genetics, Czech Academy of Sciences, Prague, CZ9 <sup>4</sup> Current affiliation: Institute of Genetics and Biotechnology, Faculty of Biology, University of

10 Warsaw, Warsaw, Poland

11 <sup>5</sup> Institute of Healthy Ageing, Department of Genetics, Evolution and Environment, University College

12 London, Gower Street, London WC1E 6BT, United Kingdom

13 \*correspondence should be addressed to linda.dixon@pirbright.ac.uk and f.werner@ucl.ac.uk

14 List of abbreviations:

Abbreviation	Definition
ASFV	African swine fever virus
NCLDV	Nucleocytoplasmic Large DNA Virus
ORF	Open reading frame
MGF	Multigene family
VACV	Vaccinia virus
RNAP	RNA polymerase
Pol II	RNA polymerase II
TBP	TATA-binding protein
TFIIB	transcription initiation factor II B

ETF	early transcription factor
TSS	transcription start site
p/np TSS / TTS	Primary/non-primary TSS/TTS
TTS	transcription termination site
CAGE-seq	cap analysis gene expression sequencing
UTR	untranslated region
NGS	Next generation sequencing
pNG	Putative novel gene
TU	Transcription unit
Inr	Initiator
EPM	Early promoter motif
UCE	Upstream Control Element
LPM	Late promoter motif
sncRNA	Small non-coding RNA
BRE	B-recognition element
ChIP	chromatin immunoprecipitation
CTSS	CAGE-seq TSS
SRA	Sequence Read Archive

15

## 16 Abstract

17 African Swine Fever Virus (ASFV) causes haemorrhagic fever in domestic pigs, presenting the biggest  
18 global threat to animal farming in recorded history. Despite its importance, little is known about the  
19 mechanisms and regulation of ASFV transcription. Using RNA sequencing methods, we have  
20 determined total RNA abundance, transcription start sites and transcription termination sites at  
21 single nucleotide-resolution. This allowed us to characterise DNA consensus motifs of early and late  
22 ASFV core promoters, as well as a poly-thymidylate sequence determinant for transcription  
23 termination. Our results demonstrate that ASFV utilises alternative transcription start sites between  
24 early and late stages of infection, and that ASFV-RNAP undergoes promoter-proximal transcript  
25 slippage at 5' ends of transcription units, adding quasi templated AU- and AUAU-5' extensions to  
26 mRNAs. Here we present the first much-needed genome-wide transcriptome study that provides  
27 unique insight into ASFV transcription and serves as a resource to aid future functional analyses of  
28 ASFV genes which are essential to combat this devastating disease.

## 29 Importance

30 African swine fever virus (ASFV) causes incurable and often lethal haemorrhagic fever in domestic  
31 pigs. In 2019, ASF presents an acute and global animal health emergency that has the potential to  
32 devastate entire national economies as effective vaccines or antiviral drugs are not currently  
33 available (Food and Agriculture Organization of the UN). With major outbreaks ongoing in Eastern  
34 Europe and Asia urgent action is needed to advance our knowledge about the fundamental biology  
35 of ASFV, including the mechanisms and temporal control of gene expression. A thorough  
36 understanding of RNAP and transcription factor function, and the sequence context of their  
37 promoter motifs, as well as accurate knowledge of which genes are expressed when and the amino  
38 acid sequence of the encoded proteins, is direly needed for the development of antiviral drugs and  
39 vaccines.

## 40 Introduction

41 ASFV is the sole characterised member of *Asfarviridae* (1), a family resembling others in the  
42 Nucleocytoplasmic Large DNA Viruses (NCLDV) and *Megavirales* order (2, 3). *Asfarviridae* also include  
43 the uncharacterised Abalone asfarvirus (NCBI:txid2654827), while the Faustoviruses show similarity to  
44 ASFV but have larger genomes and infect amoeba (*Vermamoeba vermiformis*) (4). ASFV originated in  
45 East Sub-Saharan Africa where it remains endemic, it crossed continents to Georgia in 2007 (5) and  
46 its subsequent spread in Europe and to Asia 2018 (6) has resulted in the current emergency situation.  
47 ASFV has a linear double-stranded DNA (dsDNA) genome of ~170–194 kbp encoding ~150–170 open  
48 reading frames (ORFs). Genomic variation between strains predominantly originates from loss or gain  
49 of genes at the genome termini among members of multigene families (MGFs) (7). Despite its global  
50 economic importance, little is known about ASFV transcription, but it is believed to be related to the  
51 vaccinia virus (VACV) system (8–10), a distantly-related NCLDV and *Poxviridae* family member (11).  
52 We have focused our analysis on the BA71V strain (170,101 bp genome, with 153 annotated ORFs  
53 (12, 13), because this is the most well-studied ASFV strain regarding viral molecular biology including  
54 gene expression and mRNA modification (10, 14). Based on a paradigm of the vaccinia virus, several  
55 stages of ASFV gene expression have been hypothesised in the literature including immediate early-,  
56 early-, intermediate- and late genes (10, 15–17). However, the experimental evidence for four  
57 discrete gene expression stages in ASFV leaves room for improvement, though the presence of two  
58 alternative subsets of transcription initiation factors strongly supports the notion of at least two  
59 discrete stages, early and late, likely at pre- and post-replicative stages of the virus life cycle. Previous  
60 individual gene expression studies have made use of chemical inhibitors to inhibit replication or  
61 protein synthesis (10, 15, 16). While valid tools when used with care (18), the application of these  
62 chemicals is not unproblematic due to the possibility of indirect pleiotropic effects. E.g. the  
63 nucleotide analogue cytosine arabinoside (AraC) can be incorporated into DNA and while at low  
64 concentrations mostly inhibiting replication, it can interfere with the action of many DNA-binding

enzymes including RNA polymerases, transcription factors as well as topoisomerase (19). In light of this, in this study we chose to characterise transcription unadulterated by chemical inhibitors. ASFV inhabits the eukaryotic cytoplasm and appears to be self-sufficient in terms of transcription and modification of viral mRNA. It encodes an RNA polymerase (RNAP), a poly-A polymerase and an mRNA capping enzyme, importantly, extracts obtained from mature virus particles are fully transcription competent (10, 20, 21). The basal ASFV transcription machinery resembles the eukaryotic RNAPII system encompassing an (8-subunit) ASFV-RNAP and distant relatives of the TATA-binding protein (TBP), the transcription initiation factor II B (TFIIB) and the elongation factor TFIIS (8, 9, 13). ASFV also encodes a histone-like DNA binding protein pA104R and ASFV topoisomerase II (pP1192R) which collaborate to generate DNA-binding and supercoiling activity (22). Of particular interest is the possibility that the ASFV-RNAP gains promoter-specificity in terms of temporal (early or late) gene expression dependent on the association with either TBP/TFIIB-like or virus-specific factors including those encoded by ASFV BA71V genes D1133L and G1340L, which are homologous to the D6 and A7 (respectively) early transcription factor (ETF) heterodimer (23, 24) from VACV. Promoter consensus motifs for early and late ASFV genes have not been characterised on a genome-wide scale, or in great detail, with the exception of an AT-rich sequence motif upstream of the p72 gene transcription start site (TSS) and some other late genes, as well as a consistently AT-rich region overlapping the TSS (25). Importantly, information about the temporal ASFV gene expression, TSS and transcription termination sites (TTS) is not available (10, 11).

We have applied a combination of NGS techniques including RNA-seq, RNA 5'-end (cap analysis gene expression sequencing or 'CAGE-seq') and RNA 3'-end (3' RNA-seq) determination. We report (i) the ASFV transcriptome map showing differences in gene expression between early and late infection, (ii) a genome-wide TSS map that has allowed us to define early and late ASFV promoter consensus motifs as well as 5'-mRNA leaders, and (iii) a genome-wide TTS map that provides novel insights into the mechanism of transcription termination in ASFV. Figure 1 is a genome-wide map visualising our results from TSS-mapping and differential gene expression in ASFV.

## 91 Results

### 92 Overview of the ASFV transcriptome

93 A transcriptome is defined by the overall expression levels of transcripts, and their 5' and 3' termini.  
94 We carried out RNA-seq, CAGE-seq and 3' RNA-seq in order to characterise these parameters during  
95 early and late ASFV infection, when combined they inform about the ASFV transcriptome and DNA  
96 sequence signatures associated with transcription initiation and termination. The processed data are  
97 compiled in an assembly hub and can be publicly accessed in the UCSC Genome Browser by the  
98 following address: <https://bit.ly/2TazQxK>.

99 *Vero* cells were infected with BA71V, and viral RNA was extracted at 5h and 16h post-infection. These  
100 time points were chosen based on a previous report of a small subset of genes that were  
101 experimentally characterised using nuclease S1 mapping and primer-extension analysis (10, 26).  
102 Bowtie 2 (27) mapping of the RNA-seq, CAGE-seq and 3' RNA-seq reads (summarised in  
103 Supplementary Table 1) showed a strong correlation between replicates (Pearson correlation  
104 coefficient  $r \geq 0.9$ ), with one exception of RNA-seq from 16h ( $r$  of 0.74 and 0.84 for two strands, data  
105 not shown). Figure 2a provides a whole-genome view of mapped reads from all three Next  
106 Generation Sequencing (NGS) approaches, while a selection of individual examples of TSSs and TTSSs  
107 at single-nucleotide resolution is shown in Figure 2 b-e. The sequencing depth of the RNA-seq  
108 approach was more than sufficient to determine significant changes in ASFV transcription (i.e. reads)  
109 at early and late infection due the small genome size (170 kb). The majority of CAGE-seq reads (i.e.  
110 TSSs) were located upstream and proximal to ORF start codons. A subset of late infection TSSs  
111 mapped to more distant locations between ORFs or within ORFs, these are caused by pervasive  
112 transcription, mRNA de-capping and -degradation followed by re-capping, or BA71V genome mis-  
113 annotations (28–31). The increased background of TSSs were more noticeable during late infection  
114 (Figure 2a, 'CAGE-seq 16h') and likely due to pervasive transcription, a phenomenon that has been  
115 observed in humans (32) and in VACV (28). The cause of this low-level and genome-spanning

transcription is unclear but has been attributed to an open chromatin structure in cellular organisms (33). In viral genomes it may reflect differences between nascent, newly replicated genomic DNA during late infection and genomic DNA still associated with histone-like proteins (such as A104R) just released from the virus particle during early infection.

#### Mapping of ASFV Primary Transcription Start Sites

Following mapping of CAGE-seq reads to the ASFV-BA71V genome, we located regions with an enrichment of reads corresponding to the 5' ends of transcripts and thereby the TSS. We detected a 779 clusters of CAGE-seq signals, and CAGE-seq clusters upstream annotated ORFs were manually investigated to confirm that they represent 'primary' TSSs (pTSSs) based on peak height, proximity to the ORF initiation codon, and coverage from our complementing RNA-seq data. We identified pTSSs fulfilling these criteria upstream of 151 BA71V ORFs, thus only two genes, E66L and C62L, were not found associated with a pTSS. Overall, our data showed good agreement with previously individually mapped TSSs of 44 ORFs (Supplementary Table 2). Not all of the ~780 clusters were located within 500 bp upstream of ASFV ORFs, but within, or in the antisense orientation relative to ORF coding sequences (Figure 3a). We reannotated eleven ORFs based on gene-internal TSS and RNA-seq reads (Table 1, and I177L example in Figure 3b); we provide a novel gene feature file based on our revised annotations, (Supplementary GFF).

Several genes have a *bona fide* pTSS upstream of the annotated start codon and an alternative TSS residing within the including J64R (Figure 2d) and B169L (Figure 3b). The alternative downstream TSS of J64R is weaker and specific to 16h p.i., compared to the upstream pTSS. Our genome-wide CAGE results are confirmed by previous analysis of individual genes such as I243L (26), which was shown to have distinct TSSs for different stages of infection (Figure 4a). I243L encodes a homologue of the Pol II transcript cleavage factor TFIIS, that is highly conserved between archaea, eukaryotes and among NCLDV members albeit with limited domain conservation (34). TFIIS has dual functions, it stabilises transcription initiation complexes, and reactivates stalled elongation complexes by transcript cleavage (35, 36). The late TSS is located downstream of the I243L start codon, and the utilisation of

142 the next Methionine codon would result in a TFIIS variant lacking 52 N-terminal amino acid residues  
143 (Figure 4b). While the early and long transcripts encode the fully functional three-domain TFIIS  
144 factor, the late and short transcripts encode a truncation variant lacking the N-terminal domain that  
145 is responsible for initiation functions of TFIIS. In essence, the TFIIS variants expressed during early  
146 and late infection would have a different functionality. We identified seven further genes with  
147 alternative pTSSs during early and late infection (Table 2). In most cases, the re-annotated (single  
148 pTSS downstream of start codon) or alternative pTSSs (multiple pTSSs, some downstream of start  
149 codon) did not substantially alter the ORF protein products, except for re-annotated I177L and  
150 alternative pTSSs of B169L, two putative transmembrane proteins (Figure 3b)(13, 20).

#### 151 Novel Genes Supported by Sequencing Data

152 28 TSSs in our CAGE-seq data set were not associated with annotated ORFs (Supplementary Table 3)  
153 and seven of these pTSSs were associated with transcripts that encode short ORFs, which we call  
154 putative novel genes (pNGs). These encode polypeptides of 25–56 AA length that were missed in the  
155 initial BA71V ORF prediction as only ORFs  $\geq 60$  AA were annotated (13). Five pNG ORFs showed  
156 limited similarity to short ORF-encoding genes from other ASFV strains, while pNG5 showed no clear  
157 similarity (Table 3). Interestingly, pNG6 was homologous to KP93L which is already encoded by  
158 BA71V, but barely expressed according to our data. In contrast, pNG6 was highly expressed at 5h  
159 (Supplementary Table 4). Figure 3c illustrates the features of pNG1 and pNG3, with distinct TSS and  
160 TTS, and robust RNA-seq read coverage across the entire gene. All pNGs had the same orientation as  
161 neighbouring downstream genes (Figure 1), and five of the seven pNGs transcripts terminated  
162 promptly, i.e. were associated with a drop of reads following a 5–8 nucleotide thymidylate sequence  
163 (Figure 3c and (10, 16)). All these observations support the notion that these transcription units are  
164 new *bona fide* genes.



## 165 Highly expressed ASFV genes during Early and Late Infection

166 In order to gain insights into expression of individual genes, we quantified mRNA levels obtained by  
167 CAGE-seq and compared the most abundant mRNAs at early and late time points (Figure 5a).  
168 Supplementary Table 4 summarises expression of all detected ASFV-BA71V genes including the newly  
169 annotated pNGs. For this purpose, we temporarily re-defined ASFV gene transcription units (TUs) as  
170 regions spanning from pTSS to stop codon (as proxy for TTS, see below), and quantified TU  
171 expression based on RNA-seq data (Figure 5b, Supplementary Table 5), which closely reflected the  
172 CAGE-seq analysis. The highly expressed genes matched those identified in the viral proteome of  
173 infected tissue cultures determined by mass spectrometry (highlighted in Figure 5a-b) (37). Six genes  
174 in the top-20 highly expressed genes were common during early and late infection (CP312R, A151R,  
175 K205R, Y118L, pNG1, I73R). While their expression decreases from early to late infection (see below),  
176 these genes are clearly expressed throughout, suggestive of a multistage expression pattern.  
177 Considering their high levels of expression, they are likely important throughout infection which  
178 makes them interesting candidates as potential drug- or vaccine target. However, four (out of six)  
179 have an unknown function (Figure 5a) and await functional investigation.

## 180 Differential Expression of early and late ASFV Genes

181 We characterised differential expression of ASFV genes between early and late infection by  
182 comparing separate DESeq2 analyses of CAGE-seq and RNA-seq datasets (Figure 5c and d,  
183 respectively). Based on RNA-seq data, 103 ASFV TUs showed significant differential expression  
184 (adjusted  $p$ -value < 0.05), with 47 genes down- and 56 genes up-regulated during the progression  
185 from early to late infection. Henceforth, we focused on the CAGE-seq dataset because the reads are  
186 associated with the nascent transcription start sites and thus cannot arise from transcription  
187 readthrough from upstream genes (unlike mRNA quantification using RNA-seq) which would  
188 complicate the analyses. RNA-seq also had the disadvantage of a lower sequencing depth and thus  
189 lower apparent sensitivity compared to CAGE-seq. Indeed, the CAGE-seq identified 149 genes as  
190 significantly differentially expressed with 65 downregulated genes and 84 upregulated genes (Figure

191 5c). Naturally this is not a binary classification i.e. genes that are upregulated during late infection do  
192 not have zero reads during early infection and *vice versa*. Interestingly, the relative expression levels  
193 of early genes at 5 h p.i. appeared significantly higher than late genes at 16 h p.i. (Figure 6a). This is  
194 due to normalisation of the reads and the increase of steady state levels of all transcripts during late  
195 infection, which can be seen from the sequence alignment rates (Supplementary Table 1). While the  
196 number of reads mapping to early genes during early infection is lower than the reads mapping to  
197 late genes during late infection, the total number of reads mapping to *all* ASFV genes is higher during  
198 late infection. The per-gene FPM values and differential expression analyses are normalised for ASFV-  
199 mapped sequencing depth, which therefore reduces this background and emphasises highly  
200 expressed genes during early infection. Overall, we did observe a greater and cleaner contrast in  
201 expression of the genes during early compared to late infection. The expression of the least  
202 expressed genes at 5h p.i. was more consistent and closer to zero than those at 16h p.i. (Figure 6b).  
203 The most highly expressed genes at both time points were more similar, though relative expression  
204 of the most expressed genes at 5h p.i. was higher than at 16h p.i (Figure 6c). In summary, it appears  
205 ASFV maintains a tighter control of gene expression during early infection compared to late, in as  
206 much as early genes are highly expressed and late genes show low or no expression, while during  
207 late infection the total mRNA levels increase, which results in a greater change of absolute late  
208 mRNA levels but lower relative levels of late mRNAs.

209 In order to stringently analyse differential expression in ASFV we identified the genes which showed  
210 the same pattern of differential expression according to separate DESeq2 analyses of the CAGE-seq  
211 and RNA-seq datasets. This minimises any potential biases from each of these complementing  
212 techniques. 101 genes showed significant differential expression according to both independent  
213 techniques, and the changes in expression were significantly correlated between these genes  
214 (Spearman's rank correlation coefficient  $\rho = 0.73$ , Figure 6d). Only a small number of genes, ten out  
215 of 101, showed a discrepancy between the two methods (DP63R, I329L, NP419L, B66L, A224L, E248R,  
216 O174L, D345L, C315R and NP1450L), leaving 91 genes confidently classified as early (36) and late (55)

217 genes. Supplementary Table 6 provides details of these 91 genes, their functions, and whether  
218 previously detected in viral particles (20). The 91 genes with correlated differential expression were  
219 assigned with functional categories based on their annotation in the VOCS database (38)  
220 complemented with ASFVdb (39) (Figure 6e). Around one fifth of early and late genes were classified  
221 as 'uncharacterised' without any functional predictions. The transition between 5 h and 16 h post  
222 infection is characterised by a significant up-regulation of genes important for viral morphology and  
223 structure, but also the overall diversity of differentially expressed genes changed. A significant  
224 difference was seen in the multigene family members; they constitute nearly a half of the early  
225 genes, but only one (MGF 505-2R) among late genes. ORFs annotated as having a 'transmembrane  
226 region' (TR) or a 'putative signal peptide' (PSP) were also overrepresented in late infection (Fisher  
227 Test:  $p < 0.05$ ); they remain poorly characterised beyond a domain prediction and 9 proteins (out of  
228 12) of these ORFs could be detected in BA71V virions by mass spectrometry (20).

## 229 Architecture of ASFV Gene Promoters and Consensus Elements

230 The genome-wide TSS map combined with information about their differential temporal utilisation  
231 allowed us to analyse the sequence context of TSSs and thereby characterise the consensus motifs  
232 and promoter architecture of our clearly defined 36 early and 55 late genes. Eukaryotic RNA pol II  
233 core promoters are characterised by a plethora of motifs, including TATA boxes and BRE elements,  
234 and the Initiator (Inr). The former two interact with initiation factors TBP and TFIIB, while the latter  
235 interacts with RNA pol II (40). Alignment of regions immediately surrounding pTSSs in the BA71V  
236 genome revealed several interesting ASFV promoter signatures: the Inr element overlapping the TSS  
237 is a feature that distinguishes between early and late gene promoters (Figure 7a and b, respectively).  
238 The early gene Inr is a TA(+1)NA tetranucleotide motif (where N has no nucleotide preference, Figure  
239 7c), while the late gene Inr shows a strong preference for the sequence TA(+1)TA (Figure 7d), that is  
240 not to be confused with the TBP-binding TATA box. Our late Inr consensus motif is in good agreement  
241 with those of 20 previously characterised late gene TSSs (10, 25). To search for additional promoter  
242 elements that likely interact with transcription initiation factors, we extended our search to include

243 sequences up to 40 bp upstream of the TSS. Analysis with MEME and FIMO software (41, 42)  
244 identified and located a significant 19-nt motif (Figure 7e) located ~10 bp upstream of pTSSs for 36  
245 (out of 36) early gene promoter sequences (Figure 7f), which we have called the Early Promoter  
246 Motif (EPM). Our EPM is related to the VACV early gene promoter motif ('Upstream Control Element'  
247 or UCE) (43, 44) as well as the yeast Virus-Like Element (VLE) promoters (45). However, the EPM is  
248 not limited to the 36 early genes, since a FIMO software (42) motif search identified the EPM within  
249 60 bp upstream of a much larger subset of 81 TSS/TUs including pNGs and alternative pTSSs, four of  
250 which were the early alternative pTSS for I243L, B169L, J154L and CP80R. Importantly, the limited  
251 distance distribution between the EPM and TSS is indicative of constraints defined by distinct  
252 protein-DNA interactions, e.g. by transcription initiation factors binding upstream of the TSS and  
253 ASFV-RNAP engaging with promoter DNA and TSS (Figure 7f). Figure 7g illustrates expression profiles  
254 of all genes with an EPM upstream according to FIMO, the majority showing a negative log<sub>2</sub> fold  
255 change between 5h and 16h. Since MGF members were overrepresented as early genes (Figure 6e),  
256 we searched directly for the EPM among the FIMO hits. 23 of the 29 MGF members with mapped  
257 pTSSs were associated with the EPM element including a consistent early expression and spacing  
258 relative to their TSS (Figure 7h-i), which suggests that MGF genes are under the control of their own  
259 promoters.

260 Using the same approach, we searched for promoter sequence motifs associated with late genes.  
261 MEME identified a conserved motif upstream of only 17 (out of 55) late genes, which we called Late  
262 Promoter Motif (LPM, Figure 8a). The spacing (4–12 bp) between the LPM and TSS shows a much  
263 greater diversity compared to the EPM (Figure 8b), though genes with the LPM were consistently  
264 upregulated (Figure 8c). A Tomtom (46) search identified the LPM motif as a match for 28 distinct  
265 motifs including the canonical TATA-box (*p*-value: 2.85e-03, *E*-value: 5.16e+00, Figure 8d). However,  
266 this was not a strong hit and these motifs only bear a limited resemblance to each other except for  
267 their AT-rich bias.

## 268 ASFV mRNAs have 5' leader regions

269 Early and late genes in ASFV vary with regard to the length of 5' UTRs i.e. the distance between the 5'  
270 mRNA end and the translation start codon. The 5' UTRs of late genes are significantly shorter and  
271 have a higher AT-content compared to early genes ( $p$ -value < 0.05, Figure 8e-f). Surprisingly, a subset  
272 of late gene CAGE-seq reads extended upstream of the assigned TSSs and were not complementary  
273 to the DNA template strand sequence. In order to rule out any mapping artefacts, we trimmed the  
274 CAGE-seq reads by removing the upstream 25 nt and aligned them to the genome at the 5' boundary  
275 of the reads. This did not significantly impair the mapping statistics but highlighted that nearly half of  
276 the annotated TSSs (74/158) among both early and late genes are associated with mRNAs that have  
277 short 5' extensions (or leaders), including seven genes with multiple TSSs (Supplementary Table 7).  
278 Most 5' leaders consist of two- or four nucleotides (Figure 9a) and the presence of the 5' leaders was  
279 not correlated with early or late expression (Figure 9b). The most common sequence motif in  
280 sequencing reads is AT (33% and 71% of early and late genes, respectively) and ATAT (7% in late  
281 genes, Figure 9c). In order to investigate any potential sequence-dependency of the mRNAs  
282 associated with AU- and AUAU-5' leaders, we scrutinised the template DNA sequence downstream of  
283 the TSS and found that all TUs, contained the motif ATA at positions +1 to +3 (Figure 9d). This  
284 suggests that the formation of AU-leaders is generated by RNA polymerase slippage on the first two  
285 nucleotides of the initial A(+1)TANNN template sequence, generating AUA(+1)UANNN or  
286 AUAUA(+1)UANNN mRNAs. A different but related slippage has been observed in the VACV  
287 transcription system, where all post-replicative mRNAs contain short polyA leaders which are  
288 associated with consensus Inr TAAAT motif (28).

## 289 Transcription termination of ASFV-RNAP

290 Previous mapping of mRNA 3' ends has revealed a conserved sequence motif consisting of  $\geq 7$   
291 thymidylate residues in the template, which is consistent with 3' end formation *via* transcription  
292 termination like the RNA polymerase III paradigm (16, 47). To investigate the genome-wide sequence  
293 context of ASFV transcription termination, we used 3' RNA-seq sequencing to obtain the sequences

immediately preceding ASFV mRNA poly(A) tails, generating a complete map of mRNA 3' end peaks (Figure 2a). Using a similar approach as pTSS mapping, CAGEfightR detected a total of 657 termination site clusters, 212 TTSs within 1000 bp downstream of 1–3 ORFs. Because multiple ORFs had more than one cluster within that region (Supplementary Table 8), we defined 114 primary TTSs (pTTS) as the TTS with the highest CAGEfightR-score in closest proximity to a stop codon; we classified the 98 remaining peaks as non-primary TTSs (npTTS). We identified a highly conserved poly-T signal within 10 bp upstream of 126 TTSs (83 pTTSs, 43 npTTSs) that was characterised by  $\geq 4$  consecutive T residues (Figure 10a), with the ultimate residue located on or 2 bp upstream of the ultimate T residue in the motif (Figure 10b). The remaining 86 TTSs were not associated with any recognisable sequence motif besides a single T residue 1 bp upstream of the TTS. Our results are in good agreement with a previous S1 nuclease mapping of 6 coding mRNAs, but less so with 17 proposed TTSs which were predicted based on transcript length estimates relative to upstream transcription start sites (Supplementary Table 2). This may be because only  $\geq 7$  consecutive Ts in the template were included to serve as terminators. Our results demonstrate that the total number of consecutive Ts of the poly T motif can vary, with poly T tracts of CAGE-early genes being longer than those of late genes (Figure 10c). Finally, we observed differences between CAGE-early and CAGE-late gene termination, in as much as poly T terminators were overrepresented in CAGE-early and underrepresented in CAGE-late genes (Figure 10d). The 3' UTRs (i.e. nt length from translation stop codon to pTTS) of CAGE-late genes were significantly longer compared to CAGE-early genes (Figure 10e), in good agreement with previous studies on a small number of mRNAs which showed ASFV transcripts tended to be longer and more variable in length during late infection (Supplementary Table 2). ORFs are spaced closely in the ASFV genome, and scrutiny of RNA-seq reads reveal a limited extent of transcription readthrough from upstream ORFs into downstream ORFs likely due to leaky termination (Cackett and Werner, unpublished observations). However, any additional downstream ORFs generated aberrantly by transcription readthrough would not be able to be translated since

319 there is no evidence of ASFV utilising internal ribosome entry sites (IRES) that would be required to  
320 enable cap-independent translation (7).

## 321 Discussion

322 Here we report the first comprehensive ASFV transcriptome study at single-nucleotide resolution.  
323 The mapping of 158 TSS and 114 TTS for 159 ASFV genes allowed us to reannotate the BA71V  
324 genome. Our results provide detailed information about differential gene expression during early and  
325 late infection, the sequence motifs for early and late gene promoters (EPM and LPM, and Inr  
326 elements) and terminators (poly-T motif), and evidence quasi-templated 'AU' RNA-5' tailing by the  
327 ASFV-RNAP.

328 We have discovered seven novel putative genes, some of which are highly conserved with the  
329 aggressively virulent strains (Georgia 2007/1 and Belgium 2018/1) that have caused the current  
330 outbreak in Europe (Table 3). This suggests that BA71V has more genes in common with its virulent  
331 cousins than initially thought.

332 Our results demonstrate that the majority of ASFV genes show some degree of differential  
333 expression from early to late infection (Figure 1). Interestingly, our CAGE-seq results demonstrate  
334 that early genes are expressed at relatively higher levels during early infection, than late genes  
335 during late infection (Figure 6a-c). Future experiments including spike-in controls are needed to  
336 confidently quantify the absolute mRNA levels during early- and late infection (48). The RNA  
337 sequencing methods used here quantify the steady-state RNA levels and not RNA synthesis rates,  
338 and without information about ASFV mRNA stability it is not possible to distinguish between early  
339 mRNAs retained until late infection and early genes being transcribed at later stages. Nascent ASFV  
340 mRNA synthesis rates and half-lives could be determined using techniques including TT-seq (49) or by  
341 using transcription inhibitors including actinomycin D (50). Frustratingly, many of the highly  
342 expressed genes are uncharacterised (Figure 5a). These gene products are important candidates for  
343 further functional characterisation and may emerge as promising targets for vaccine development.

344 We have shown that MGFs show a distinct downregulation from early to late infection, while genes  
345 annotated as transmembrane region or putative signal peptides (though poorly characterised beyond  
346 this), along with structural or viral morphology genes, are overrepresented in late infection (Figure  
347 6e). Our CAGE-analysis also identified TSS signals unlikely to serve as primary TSS for annotated  
348 genes (Figure 3a and Supplementary Table 9); these could provide a rich hunting ground for small  
349 non-coding (snc)RNAs. One TSS cluster associated with an sncRNA gene (at 71,302 on the BA71V  
350 genome) was previously reported by Dunn et al. (51) as ASFVsRNA2, that is encoded in the antisense  
351 orientation relative to the ASFV RNA polymerase subunit RPB6-encoding gene. Further investigation  
352 of antisense sncRNAs in the BA71V transcriptome may uncover further examples of riboregulation,  
353 i.e. a more complex method of modulating its own or host gene expression beyond the protein level.  
354 While eukaryotic Pol II and archaeal RNAP critically rely on initiation factors TBP and TFIIB for  
355 transcription initiation on all mRNA genes, bacterial RNAP obtains specificity for subsets of gene  
356 promoters by associating with distinct sigma factors (52). ASFV-RNAP is related to archaeal and  
357 eukaryotic RNA polymerases, detailed phylogenetic analyses reveal that the RPB1 subunit is most  
358 closely related to the RNA polymerase I homologue (3, 45, 53). However, transcription initiation of  
359 early and late genes appears to be directed by two distinct sets of general initiation factors and their  
360 cognate DNA recognition motifs, as our TSS mapping demonstrates. The first feature of all ASFV  
361 promoters is the Inr element, a tetranucleotide motif overlapping the TSS with an A-residue serving  
362 as initiating nucleotides similar to most RNAP systems. The similarity of early and late gene Inr  
363 sequences, is likely because the Inr makes sequence-specific contacts with amino acid sidechains of  
364 the two largest RNAP subunits (RPB1 and 2). The EPM and LPM are located upstream of the TSS, both  
365 are AT-rich, though distinct in sequence (Figure 7e and 8a). The distance distribution of EPM is  
366 narrow (located 9–10 bp upstream of the TSS) while the distance between the LPM and TSS shows  
367 greater variation and is located closer (4–6 bp) to the TSS. The high sequence and distance  
368 conservation of the EPM, especially exemplified for early expressed MGFs (Figure 7h-i), emphasises  
369 the EPM's role in tight control of transcription during early infection. Considering the close



relationship between ASFV and VACV, we posit that the EPM is recognised by heterodimeric ASFV-BA71V D1133L/G1340L initiation factor (VACV D6/A7) (11) consistent with the late expression of these genes (Figure 6d, also ref (54)). Presence of D1133L/G1340L gene-products along with RNAP in viral particles (20) provides a system that is primed to initiate ASFV transcription of early genes.

ASFV-TBP (B263R) is an early gene and ASFV-TFIIB (C315R) is expressed throughout infection. We propose the LPM is utilised by ASFV-TBP and -TFIIB homologues, neither of which were detected in virions (20). A functional comparison of the LPM to the classical Pol II core promoter elements BRE/TATA-box is compelling. However, the tight spacing between the LPM and TSS is incompatible with the overall topology of a classical eukaryotic and archaeal TATA-TBP-TFIIB-RNA pol II preinitiation complex (PIC), where the BRE/TATA promoter elements are located ~ 24 bp upstream of the TSS (55). Considering low sequence conservation between cellular and ASFV-TBP (8) and unusual spacing of LPM and Inr, the structure of ASFV LPM-TBP-TFIIB-RNAP PIC is likely very different from canonical RNA pol II PICs. Additionally, factors including ASFV B175L and B385R may contribute to the PIC, as was proposed for VACV-A1 and A2 (56, 57). At this stage, we cannot rule out a limited overlap between early and late genes without additional information including insights into pre- and post-replicative gene expression pattern, mRNA stability of early and late genes, and knowledge about all regulatory factors that enable the temporal regulation of ASFV transcription. To unequivocally attribute factors to their cognate binding motifs genome-wide, a chromatin immunoprecipitation (ChIP) approach is required; the results may be full of surprises and have the potential to shed light on multistage gene expression pattern including the possibility of a more complex promoter architecture where some genes are under the control of more than one promoter.

An in-depth characterisation of the global gene regulation in ASFV with a higher temporal resolution is essential to assess how closely ASFV follows the cascade-like patterns of VACV (11). While two genes have been proposed to be intermediate genes in ASFV, both of them are also expressed during

intermediate and late (I226R), and during early, intermediate and late stages (I243L). Thus, there is no hard evidence of genes that are specifically expressed during the intermediate stage (26). A combination of a reversible replication inhibitor and a conditionally regulated late transcription factor has been successfully used to characterise intermediate gene expression in VACV (58). Such an approach might also be useful to identify intermediate ASFV genes - and help us refine the LPM that in our current analysis could reflect a combination of late- and intermediate gene promoters'.

We found several examples of alternative, gene-internal, TSS utilisation with the potential to increase the complexity of the viral proteome; protein variants which may provide the means to generate distinct functionalities, which has also been described in VACV by Yang et al. (28). Our TSS mapping uncovered a form of transcript slippage by the ASFV-RNAP occurring on promoters that start with an A(+1)TA motif, where mRNAs are extended by one or two copies of the dinucleotide AU. This is reminiscent of VACV, where late gene transcripts containing a poly-A 5' UTR (28) are associated with improved translation efficiency and reduced reliance on cap-dependent translation initiation (59, 60); likewise, distinct functional attributes of polyA leaders in translation have been documented in eukaryotes (61). Whether the 5' AU- and AUAU-tailing is a peculiarity of the ASFV-RNAP initiation, or whether these mRNA 5' leaders have any functional implications, remains to be investigated. The structural determinants underlying RNAP slippage are interactions between the template DNA sequence and the RNAP and/or transcription initiation factors; the differential use of distinct initiation factors for the transcription of early and late ASFV genes may account for difference in leader sequences.

The mechanisms underlying transcription termination of multisubunit RNAP are diverse (62, 63). Our analyses of genome-wide ASFV RNA-3' ends allowed the mapping of the ASFV 'terminome'. Over half of mRNA 3' ends are characterised by a stretch of seven U residues, with the TTS mostly coinciding with the last T residue in the template DNA motif - in good agreement with ASFV terminators that have been individually mapped (15, 16). In contrast, VACV appears to utilise a motif ~ 40 nt upstream of the mRNA 3' ends (64, 65). In essence, the ASFV-RNAP is akin to archaeal RNAPs and RNA pol III,

421 where a poly-U stretch is the sole *cis*-acting motif without any RNA secondary structures  
422 characteristic of bacterial intrinsic terminators (63). The pTTs without any association with poly-U  
423 motifs are still likely to represent *bona fide* termination sites, since RNA-seq reads were decreasing  
424 towards these termination sites, despite no clear conserved sequence motif. However, ASFV does  
425 encode several (VACV-related) RNA helicases that have been speculated to facilitate transcription  
426 termination and/or mRNA release (10, 66). Future functional studies will address the molecular  
427 mechanisms of termination including the role of putative termination factors.

428 Understanding the molecular mechanisms of the ASFV transcription system is not only of academic  
429 interest. Unless effective vaccines in conjunction with antiviral treatments against ASFV are  
430 developed, a large proportion of the global pig population is projected to die in the context of this  
431 terrible disease (OIE, <https://www.oie.int>). The rational design of drugs that target the gene  
432 expression machinery is crucially reliant on our knowledge about the ASFV-RNAP, the basal factors  
433 that govern its function, and the DNA sequences they interact with, while vaccine development  
434 benefits from the intricate knowledge about gene expression patterns. Our results directly contribute  
435 to these burning issues for animal husbandry.

## 436 Methods

### 437 RNA Sample Extraction from Vero Cells infected with BA71V

438 *Vero* cells (Sigma-Aldrich, cat #84113001) were grown in 6-well plates, plates and were infected in 2  
439 replicate wells for 5h or 16h with a multiplicity of infection of 5 of the ASFV BA71V strain, collected in  
440 Trizol Lysis Reagent (Thermo Fisher Scientific) separately, after growth medium was removed.  
441 Infected cells were collected at 5h post-infection (samples for RNA-seq: S3-5h and S4-5h, CAGE-seq:  
442 S1-5h and S2-5h and 3' RNA-seq: E-5h\_1 and E-5h\_1), and at 16h post-infection (RNA-seq: S5-16h  
443 and S6-16h, CAGE-seq: S3-16h and S4-16h, and 3' RNA-seq: L-16h\_1, L-16h\_1). RNA was extracted  
444 according to manufacturer's instructions for Trizol extraction and the subsequent RNA-pellets were  
445 resuspended in 50µl RNase-free water and DNase-treated (Turbo DNase kit, Invitrogen). RNA  
446 quality was assessed *via* Bioanalyzer (Agilent 2100), before ethanol precipitation. For CAGE-seq and  
447 3' RNA-seq, samples were sent to *CAGE-seq* (Kabushiki Kaisha DNAFORM, Japan) and Cambridge  
448 Genomic Services (Department of Pathology, University of Cambridge, Cambridge, UK), respectively.

### 449 RNA-seq, CAGE-seq and 3' RNA-seq Library Preparations and Sequencing

450 For RNA-seq, samples were resuspended in 100µl RNase-free water, and polyA-enriched using the  
451 BIOO SCIENTIFIC NEXTflex™ Poly(A) Beads kit according to manufacturer's instructions and quality  
452 was assessed *via* Bioanalyzer. NEXTflex™ Rapid Directional qRNA-Seq™ Kit was utilised to produce  
453 paired-end indexed cDNA libraries from the polyA-enriched RNA samples, according to the  
454 manufacturer's instructions. Per-sample cDNA library concentrations were calculated *via* Bioanalyzer,  
455 and Qubit Fluorometric Quantitation (Thermo Fisher Scientific). Sample S3-5h, S4-5h, S5-16h and S6-  
456 16h cDNA libraries were twice separately sequenced on Illumina MiSeq generating 75 bp reads  
457 (Supplementary Table 1) and 12 FASTQ files.

458 Library preparation and CAGE-sequencing of RNA samples S1-5h, S2-5h, S3-16h and S4-16h was  
459 carried out by *CAGE-seq* (Kabushiki Kaisha DNAFORM, Japan). Library preparation produce single-end  
460 indexed cDNA libraries for sequencing: in brief, this included reverse transcription with random

461 primers, oxidation and biotinylation of 5' mRNA cap, followed by RNase ONE treatment removing  
462 RNA not protected in a cDNA-RNA hybrid. Two rounds of cap-trapping using Streptavidin beads,  
463 washing away uncapped RNA-cDNA hybrids. Next, RNase ONE and RNase H treatment degraded any  
464 remaining RNA, and cDNA strands were subsequently released from the Streptavidin beads and  
465 quality-assessed *via* Bioanalyzer. Single strand index linker and 3' linker was ligated to released cDNA  
466 strands, and primer containing Illumina Sequencer Priming site was used for second strand synthesis.  
467 Samples were sequenced using the Illumina NextSeq 500 platform producing 76 bp reads  
468 (Supplementary Table 1).

469 3' RNA-seq was carried out with samples E-5h\_1, E-5h\_2, L-16h\_1 and L-16h\_2 using the Lexogen  
470 QuantSeq 3' mRNA-Seq Library Prep Kit FWD for Illumina according to manufacturer's instructions.  
471 Library preparation and sequencing were carried out Cambridge Genomic Services (Department of  
472 Pathology, University of Cambridge, Cambridge, UK) on a single NextSeq flowcell producing 150 bp  
473 (Supplementary Table 1).

#### 474 Sequencing Quality Checks and Mapping to ASFV and Vero Genomes

475 FastQC (67) analysis was carried out on all FASTQ files: for RNA-seq FASTQ files were uploaded to the  
476 web-platform Galaxy ([www.usegalaxy.org/](http://www.usegalaxy.org/)) (68, 69) and all reads were trimmed by the first 10 and  
477 last 1 nt using FASTQ Trimmer (70). After read-trimming, FASTQ files originating from the same RNA  
478 samples were concatenated. RNA-seq reads were mapped to the ASFV-BA71V (NC\_001659.2) and  
479 *Vero* (GCF\_000409795.2) genomes using Bowtie 2 directly after trimming (27), with alignments  
480 output in SAM file format. FASTQ analysed CAGE-seq reads showed consistent read quality across  
481 the 76 bp reads, except for the nucleotide 1. This was an indicator of the 5' mRNA methylguanosine  
482 due to the reverse transcriptase used in library preparation (71), therefore, the reads were mapped  
483 in their entirety to the ASFV-BA71V (U18466.2) and *Vero* (GCF\_000409795.2) genomes.

484 FASTQC analysed 3' RNA-seq reads showed relatively varying and poorer quality after nucleotide 65.  
485 Cutadapt (72) was utilised to extract only fastq reads with 18 consecutive A's at the 3' end followed  
486 by the sample i7 Illumina adapter, selecting only for reads containing the 3' mRNA end and the polyA

487 tail. The 18A-adaptor sequences were then trimmed and FASTQC-analysed reads were mapped *via*  
488 Bowtie2 to ASFV-BA71V (U18466.2) and *Vero* (GCF\_000409795.2) genomes.

#### 489 CAGE Analysis, TSS-Mapping

490 CAGE-seq mapped sample BAM files were converted to BigWig (BW) format with BEDtools (73)  
491 genomecov, to produce per-strand BW files of 5' read ends. Stranded BW files were input for TSS-  
492 prediction in RStudio (74) with Bioconductor (75) package CAGEfightR (76). Genomic feature  
493 locations were imported as a TxDb object from U18466.2 genome gene feature file (GFF3), modified  
494 to include C44L (12). CAGEfightR was used to quantify the CAGE tag transcripts mapping at base pair  
495 resolution to the ASFV-BA71V genome - at CAGE TSSs (CTSSs). CTSS values were normalized by tags-  
496 per-million for each sample, pooled and only CTSSs supported by presence in  $\geq 2$  samples were kept.  
497 CTSSs were assigned to clusters, merging CTSSs within 50 bp of one another, filtering out pooled,  
498 TPM-normalized CTSS counts below 25, and assigned a 'thick' value as the highest CTSS peak within  
499 that cluster. CTSS clusters were assigned to annotated U18466.2 ORFs (if clusters were between 300  
500 bp upstream and 200 bp downstream of an ORF). Clusters were classified 'tssUpstream' if located  
501 within 300 bp upstream of an ORF, 'proximal' if located within 500 bp of an ORF, 'CDS' if within the  
502 ORF, 'NA' if no annotated ORF was within these regions (excepting pNG), and antisense if within  
503 these regions but antisense relative to the ORF.  
504 Cluster classification was not successful in all cases, therefore, manual adjustment was necessary.  
505 Integrative Genomics Viewer (IGV) (77) was used to visualise BW files relative to the BA71V ORFs,  
506 and incorrectly classified clusters were corrected. Clusters with the 'tssUpstream' classification were  
507 split into subsets for each ORF. 'Primary' cluster subset contained either the highest scoring  
508 CAGEfightR cluster or the highest scoring manually-annotated peak, and the highest peak coordinate  
509 was defined as the primary TSS (pTSS) for an ORF. Further clusters associated these ORFs were  
510 classified as 'non-primary', highest peak as a non-primary TSS (npTSS).  
511 If the strongest CTSS location was intra-ORF and corroborated with RNA-seq coverage, then the ORF  
512 was re-defined as starting from the next ATG downstream. For the 28 intergenic CTSSs, IGV was used

513 to visualise if CAGE BW peaks were followed by RNA-seq coverage downstream, and whether the  
514 transcribed region encode a putative ORF using NCBI Open Reading Frame Finder (78).

#### 515 TTS-Mapping

516 TTSs were mapped in a similar manner to TSSs and CAGEfightR was utilised as above to locate  
517 clusters of 3' RNA-seq peaks, though differed in some respects: input BigWig files contained the 3'  
518 read-end coverage extracted from BAM files using BEDtools genomecov. Clusters were detected for  
519 the 3' RNA-seq peaks in the same manner as before, except merging clusters < 25 nt apart, which  
520 detected a total of 567 clusters. BEDtools was used to check whether the highest point of each  
521 cluster (TTS) was within 500 bp or 1000 bp downstream of annotated ORFs and pNGs. TTSs were  
522 then filtered out if 10 nt downstream of the 3' end had  $\geq 50\%$  As, to exclude clusters potentially  
523 originated from miss-priming. TTS clusters for pNG3 and pNG4 were initially filtered out but included  
524 in final 212 TTSs due to their strong RNA-seq agreement. In cases of multiple TTS clusters per gene  
525 we defined the highest CAGEfightR-scored one within 1000 bp downstream of ORFs as primary  
526 (pTTS) unless no clear RNA-seq coverage was shown, or manually annotated from the literature for  
527 O61R (15).

#### 528 DESeq2 Differential Expression Analysis of ASFV Genes

529 A new GFF was produced for investigating differential expression of ASFV genes across the genome  
530 with changes from the original U18466.2.gff: for all 151 ASFV ORFs which had identified pTSSs, we  
531 defined their transcription unit as beginning from the pTSS coordinate to ORF end. Since no pTSS was  
532 identified for ORFs E66L and C62L these entries were left as ORFs within the GFF, while the 7 putative  
533 pNGs were defined as their pTSS down to the genome coordinate at which the RNA-seq coverage  
534 ends. In 8 cases where genes had alternative pTSSs for the different time-points the TUs were  
535 defined as the most upstream pTSS down to the ORF end. For analysing differential expression with  
536 the CAGE-seq dataset a GFF was created with BEDtools extending from the pTSS coordinate, 25 bp  
537 upstream and 75 bp downstream, however, in cases of alternating pTSSs this TU was defined as 25

538 bp upstream of the most upstream pTSS and 75 bp downstream of the most downstream pTSS.  
539 HTSeq-count (79) was used to count reads mapping to genomic regions described above for both the  
540 RNA- and CAGE-seq sample datasets. The raw read counts were then used to analyse differential  
541 expression across these regions between the time-points using DESeq2 (default normalisation  
542 described by Love et al. (80)) and those regions showing changes with an adjusted  $p$ -value (padj) of  
543  $<0.05$  were considered significant. Further analysis of ASFV genes used their characterised or  
544 predicted functions as found in the VOCS tool database (<https://4virology.net/>) (38, 81) or ASFVdb  
545 (39) entries for the ASFV-BA71V genome.

#### 546 Early and Late Promoter Analysis

547 DESeq2 results were used to categorise ASFV genes into two simple sub-classes: early; genes  
548 downregulated from early to late infection and late; those upregulated from early to late infection.  
549 For those with newly annotated pTSSs (151 including 7 pNGs but excluding 15 alternative pTSSs),  
550 sequences 30 bp upstream and 5 bp downstream were extracted from the ASFV-BA71V genome in  
551 FASTA format using BEDtools. The 36 Early, 55 Late and all 166 pTSSs (including alternative ones) at  
552 once were analysed using MEME software (<http://meme-suite.org>) (82), searching for 5 motifs with a  
553 width of 10–25 nt, other settings at default. Significant motifs (E-value  $< 0.05$ ) detected *via* MEME  
554 were submitted to a following FIMO (42) search ( $p$ -value cut-off  $< 0.0001$ ) of 60 nt upstream of the  
555 total 166 pTSS sequences (including pNGs and alternative pTSSs), and Tomtom software (46) search  
556 (UP00029\_1, Database: uniprobe\_mouse) to find similar known motifs.

#### 557 Data Availability

558 Sequencing data from RNA-seq, CAGE-seq and 3' RNA-seq are available on Sequence Read Archive  
559 (SRA), BioProject: PRJNA590857  
560 (<https://dataview.ncbi.nlm.nih.gov/object/PRJNA590857?reviewer=e597nf6o3r2hk5r45a5hqr9d>).  
561 The processed data for two replicates are visualized in an UCSC Genome Browser [pmid: 24227676]  
562 and can be accessed at <https://bit.ly/2TazQxK>. The tracks include corrected gene annotations  
563 (primary TSSs, primary TTSs, and ORF coordinates), raw coverage of 5' ends (CAGE-seq) and 3' ends



(3'-RNA-seq), and RPKM values for the RNA-seq data. Coverage for the forward and reverse strands are shown in blue and red, respectively. Results from differential gene expression analysis with DESeq2 of CAGE-seq and RNA-seq are found in Supplementary Tables 4 and 5, respectively. The 91 genes showing the same pattern of differential expression according to both of these NGS techniques are found in Supplementary Table 6. Details of non-templated extensions detected from CAGE-seq are in Supplementary Table 7. CAGEfightR-detected cluster peaks from 3' RNA-seq after removal of those arriving from polyA miss-priming are described in Supplementary Table 8. All 779 CAGEfightR-detected cluster peaks from CAGE-seq are listed in Supplementary Table 9.

## Acknowledgements and Funding

Research in the RNAP laboratory at UCL is funded by a Wellcome Investigator Award in Science 'Mechanisms and Regulation of RNAP transcription' to FW (WT 207446/Z/17/Z) and to JB [WT 095598/Z/11/Z]. GC is funded by the Wellcome Trust ISMB 4-year PhD programme 'Macromolecular machines: interdisciplinary training grounds for structural, computational and chemical biology' (WT 108877/B/15/Z). The authors are grateful to all members of the RNAP lab and Tine Arnvig for critical reading of the manuscript.

## Competing Interests

The authors declare that no competing interests exist.

## 584 References

- 585 1. Alonso C, Borca M, Dixon L, Revilla Y, Rodriguez F, Escribano JM, Consortium IR. 2018. ICTV  
586 Virus Taxonomy Profile: Asfarviridae. *J Gen Virol* 99:613–614.
- 587 2. Koonin E V, Yutin N. 2010. Origin and evolution of eukaryotic large nucleo-cytoplasmic DNA  
588 viruses. *Intervirology* 53:284–92.
- 589 3. Yutin N, Koonin E V. 2012. Hidden evolutionary complexity of Nucleo-Cytoplasmic Large DNA  
590 viruses of eukaryotes. *Virol J* 9:161.
- 591 4. Reteno DG, Benamar S, Khalil JB, Andreani J, Armstrong N, Klose T, Rossmann M, Colson P,  
592 Raoult D, La Scola B. 2015. Faustovirus, an asfarvirus-related new lineage of giant viruses  
593 infecting amoebae. *J Virol* 89:6585–94.
- 594 5. Gogin A, Gerasimov V, Malogolovkin A, Kolbasov D. 2013. African swine fever in the North  
595 Caucasus region and the Russian Federation in years 2007-2012. *Virus Res* 173:198–203.
- 596 6. Zhou X, Li N, Luo Y, Liu Y, Miao F, Chen T, Zhang S, Cao P, Li X, Tian K, Qiu H-J, Hu R. 2018.  
597 Emergence of African Swine Fever in China, 2018. *Transbound Emerg Dis* 65:1482–1484.
- 598 7. Dixon LK, Chapman DAG, Netherton CL, Upton C. 2013. African swine fever virus replication  
599 and genomics. *Virus Res* 173:3–14.
- 600 8. Kinyanyi D, Obiero G, Obiero GFO, Amwayi P, Mwaniki S, Wamalwa M. 2018. In silico  
601 structural and functional prediction of African swine fever virus protein-B263R reveals  
602 features of a TATA-binding protein. *PeerJ* 6:e4396.
- 603 9. Rodriguez JM, Salas ML, Viñuela E. 1992. Genes homologous to ubiquitin-conjugating proteins  
604 and eukaryotic transcription factor SII in African swine fever virus. *Virology* 186:40–52.
- 605 10. Rodríguez JM, Salas ML. 2013. African swine fever virus transcription. *Virus Res* 173:15–28.
- 606 11. Broyles SS. 2003. Vaccinia virus transcription. *J Gen Virol* 84:2293–2303.
- 607 12. Kollnberger SD, Gutierrez-Castañeda B, Foster-Cuevas M, Corteyn A, Parkhouse RME. 2002.  
608 Identification of the principal serological immunodeterminants of African swine fever virus by  
609 screening a virus cDNA library with antibody. *J Gen Virol* 83:1331–1342.

- 610 13. Yáñez RJ, Rodríguez JM, Nogal ML, Yuste L, Enríquez C, Rodríguez JF, Viñuela E. 1995. Analysis  
611 of the complete nucleotide sequence of African swine fever virus. *Virology* 208:249–78.
- 612 14. Rodríguez JM, Moreno LT, Alejo A, Lacasta A, Rodríguez F, Salas ML. 2015. Genome Sequence  
613 of African Swine Fever Virus BA71, the Virulent Parental Strain of the Nonpathogenic and  
614 Tissue-Culture Adapted BA71V. *PLoS One* 10:e0142889.
- 615 15. Almazán F, Rodríguez JM, Angulo A, Viñuela E, Rodríguez JF. 1993. Transcriptional mapping of  
616 a late gene coding for the p12 attachment protein of African swine fever virus. *J Virol* 67:553–  
617 6.
- 618 16. Almazán F, Rodríguez JM, Andrés G, Pérez R, Viñuela E, Rodríguez JF. 1992. Transcriptional  
619 analysis of multigene family 110 of African swine fever virus. 66.
- 620 17. Breese SS, DeBoer CJ. 1966. Electron microscope observations of African swine fever virus in  
621 tissue culture cells. *Virology* 28:420–428.
- 622 18. Oda KI, Joklik WK. 1967. Hybridization and sedimentation studies on “early” and “late”  
623 vaccinia messenger RNA. *J Mol Biol* 27:395–419.
- 624 19. Zhang X, Kiechle FL. 2004. Cytosine arabinoside substitution decreases transcription factor-  
625 DNA binding element complex formation. *Arch Pathol Lab Med* 128:1364–1371.
- 626 20. Alejo A, Matamoras T, Guerra M, Andrés G. 2018. A proteomic atlas of the African swine fever  
627 virus particle. *J Virol JVI.01293-18*.
- 628 21. Salas ML, Kuznar J, Viñuela E. 1981. Polyadenylation, methylation, and capping of the RNA  
629 synthesized in vitro by African swine fever virus. *Virology* 113:484–491.
- 630 22. Frouco G, Freitas FB, Coelho J, Leitão A, Martins C, Ferreira F. 2017. DNA-Binding Properties of  
631 African Swine Fever Virus pA104R, a Histone-Like Protein Involved in Viral Replication and  
632 Transcription. *J Virol* 91.
- 633 23. Iyer LM, Balaji S, Koonin E V., Aravind L. 2006. Evolutionary genomics of nucleo-cytoplasmic  
634 large DNA viruses. *Virus Res* 117:156–84.
- 635 24. Yutin N, Wolf YI, Raoult D, Koonin E V. 2009. Eukaryotic large nucleo-cytoplasmic DNA viruses:

- 636 clusters of orthologous genes and reconstruction of viral genome evolution. *Viol J* 6:223.
- 637 25. García-Escudero R, Viñuela E. 2000. Structure of African swine fever virus late promoters:  
638 requirement of a TATA sequence at the initiation region. *J Virol* 74:8176–82.
- 639 26. Rodríguez JM, Salas ML, Viñuela E. 1996. Intermediate class of mRNAs in African swine fever  
640 virus. *J Virol* 70:8584–9.
- 641 27. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods*  
642 9:357–359.
- 643 28. Yang Z, Martens CA, Bruno DP, Porcella SF, Moss B. 2012. Pervasive initiation and 3'-end  
644 formation of poxvirus postreplicative RNAs. *J Biol Chem* 287:31050–60.
- 645 29. Yang Z, Bruno DP, Martens CA, Porcella SF, Moss B. 2011. Genome-wide analysis of the 5' and  
646 3' ends of vaccinia virus early mRNAs delineates regulatory sequences of annotated and  
647 anomalous transcripts. *J Virol* 85:5897–909.
- 648 30. Schoenberg DR, Maquat LE. 2009. Re-capping the message. *Trends Biochem Sci* 34:435–42.
- 649 31. Mukherjee C, Patil DP, Kennedy BA, Bakthavachalu B, Bundschuh R, Schoenberg DR. 2012.  
650 Identification of cytoplasmic capping targets reveals a role for cap homeostasis in translation  
651 and mRNA stability. *Cell Rep* 2:674–84.
- 652 32. Clark MB, Amaral PP, Schlesinger FJ, Dinger ME, Taft RJ, Rinn JL, Ponting CP, Stadler PF, Morris  
653 K V., Morillon A, Rozowsky JS, Gerstein MB, Wahlestedt C, Hayashizaki Y, Carninci P, Gingeras  
654 TR, Mattick JS. 2011. The reality of pervasive transcription. *PLoS Biol* 9:e1000625.
- 655 33. Castelnovo M, Stutz F. 2015. Role of chromatin, environmental changes and single cell  
656 heterogeneity in non-coding transcription and gene regulation. *Curr Opin Cell Biol*. Elsevier  
657 Ltd.
- 658 34. Mirzakhanyan Y, Gershon PD. 2017. Multisubunit DNA-Dependent RNA Polymerases from  
659 Vaccinia Virus and Other Nucleocytoplasmic Large-DNA Viruses: Impressions from the Age of  
660 Structure. *Microbiol Mol Biol Rev* 81:e00010-17.
- 661 35. Kim B, Nesvizhskii AI, Rani PG, Hahn S, Aebersold R, Ranish JA. 2007. The transcription

- 662 elongation factor TFIIS is a component of RNA polymerase II preinitiation complexes. *Proc*  
663 *Natl Acad Sci* 104:16068–16073.
- 664 36. Awrey DE, Shimasaki N, Koth C, Weilbaeher R, Olmsted V, Kazanis S, Shan X, Arellano J,  
665 Arrowsmith CH, Kane CM, Edwards AM. 1998. Yeast transcript elongation factor (TFIIS),  
666 structure and function. II: RNA polymerase binding, transcript cleavage, and read-through. *J*  
667 *Biol Chem* 273:22595–605.
- 668 37. Keßler C, Forth JH, Keil GM, Mettenleiter TC, Blome S, Karger A. 2018. The intracellular  
669 proteome of African swine fever virus. *Sci Rep* 8:14714.
- 670 38. Upton C, Slack S, Hunter AL, Ehlers A, Roper RL, Rock DL. 2003. Poxvirus orthologous clusters:  
671 toward defining the minimum essential poxvirus genome. *J Virol* 77:7590–600.
- 672 39. Zhu Z, Meng G. 2019. ASFVdb: An integrative resource for genomics and proteomics analyses  
673 of African swine fever. *bioRxiv* 670109.
- 674 40. Butler JEF, Kadonaga JT. 2002. The RNA polymerase II core promoter: a key component in the  
675 regulation of gene expression. *Genes Dev* 16:2583–92.
- 676 41. Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, Ren J, Li WW, Noble WS. 2009.  
677 MEME SUITE: tools for motif discovery and searching. *Nucleic Acids Res* 37:W202–W208.
- 678 42. Grant CE, Bailey TL, Noble WS. 2011. FIMO: scanning for occurrences of a given motif.  
679 *Bioinformatics* 27:1017–8.
- 680 43. Davison AJ, Moss B. 1989. Structure of vaccinia virus early promoters. *J Mol Biol* 210:749–769.
- 681 44. Yang Z, Bruno DP, Martens CA, Porcella SF, Moss B. 2010. Simultaneous high-resolution  
682 analysis of vaccinia virus and host cell transcriptomes by deep RNA sequencing. *Proc Natl*  
683 *Acad Sci U S A* 107:11513–8.
- 684 45. Sýkora M, Pospíšek M, Novák J, Mrvová S, Krásný L, Vopálenský V. 2018. Transcription  
685 apparatus of the yeast virus-like elements: Architecture, function, and evolutionary origin.  
686 *PLOS Pathog* 14:e1007377.
- 687 46. Gupta S, Stamatoyannopoulos JA, Bailey TL, Noble W. 2007. Quantifying similarity between

- 688 motifs. *Genome Biol* 8:R24.
- 689 47. Arimbasseri AG, Rijal K, Maraia RJ. 2013. Comparative overview of RNA polymerase II and III  
690 transcription cycles, with focus on RNA polymerase III termination and reinitiation.  
691 Transcription. Taylor and Francis Inc.
- 692 48. Chen K, Hu Z, Xia Z, Zhao D, Li W, Tyler JK. 2016. The Overlooked Fact: Fundamental Need for  
693 Spike-In Control for Virtually All Genome-Wide Analyses.
- 694 49. Schwalb B, Michel M, Zacher B, Hauf KF, Demel C, Tresch A, Gagneur J, Cramer P. 2016. TT-seq  
695 maps the human transient transcriptome. *Science* (80- ) 352:1225–1228.
- 696 50. Kuznar J, Salas ML, Viñuela E. 1980. DNA-dependent RNA polymerase in African swine fever  
697 virus. *Virology* 101:169–75.
- 698 51. Dunn LEM, Ivens A, Netherton CL, Chapman DAG, Beard PM. 2019. Identification of a  
699 functional small non-coding RNA encoded by African swine fever virus. *bioRxiv* 865147.
- 700 52. Kazmierczak MJ, Wiedmann M, Boor KJ. 2005. Alternative Sigma Factors and Their Roles in  
701 Bacterial Virulence. *Microbiol Mol Biol Rev* 69:527–543.
- 702 53. Guglielmini J, Woo AC, Krupovic M, Forterre P, Gaia M. 2019. Diversification of giant and large  
703 eukaryotic dsDNA viruses predated the origin of modern eukaryotes. *Proc Natl Acad Sci U S A*  
704 116:19585–19592.
- 705 54. Yáñez RJ, Rodríguez JM, Bournsnel M, Rodríguez J, Viñuela E. 1993. Two putative African swine  
706 fever virus helicases similar to yeast ‘DEAH’ pre-mRNA processing proteins and vaccinia virus  
707 ATPases D11L and D6R. *Gene* 134:161–174.
- 708 55. Hahn S. 2004. Structure and mechanism of the RNA polymerase II transcription machinery.  
709 *Nat Struct Mol Biol* 11:394–403.
- 710 56. Knutson BA, Liu X, Oh J, Broyles SS. 2006. Vaccinia virus intermediate and late promoter  
711 elements are targeted by the TATA-binding protein. *J Virol* 80:6784–93.
- 712 57. Broyles SS, Knutson BA. 2010. Poxvirus transcription. *Future Virol* 5:639–650.
- 713 58. Yang Z, Reynolds SE, Martens CA, Bruno DP, Porcella SF, Moss B. 2011. Expression Profiling of

- 714 the Intermediate and Late Stages of Poxvirus Replication † Downloaded from. J Virol 85:9899–  
715 9908.
- 716 59. Dhungel P, Cao S, Yang Z. 2017. The 5'-poly(A) leader of poxvirus mRNA confers a translational  
717 advantage that can be achieved in cells with impaired cap-dependent translation. PLoS Pathog  
718 13:e1006602.
- 719 60. Mulder J, Robertson ME, Seamons RA, Belsham GJ. 1998. Vaccinia virus protein synthesis has  
720 a low requirement for the intact translation initiation factor eIF4F, the cap-binding complex,  
721 within infected cells. J Virol 72:8813–9.
- 722 61. Shirokikh NE, Spirin AS. 2008. Poly(A) leader of eukaryotic mRNA bypasses the dependence of  
723 translation on initiation factors. Proc Natl Acad Sci U S A 105:10738.
- 724 62. Kuehner JN, Pearson EL, Moore C. 2011. Unravelling the means to an end: RNA polymerase II  
725 transcription termination. Nat Rev Mol Cell Biol.
- 726 63. Santangelo TJ, Cubonová L, Skinner KM, Reeve JN. 2009. Archaeal intrinsic transcription  
727 termination in vivo. J Bacteriol 191:7102–7108.
- 728 64. Howard ST, Ray CA, Patel DD, Antczak JB, Pickup DJ. 1999. A 43-nucleotide RNA cis-acting  
729 element governs the site-specific formation of the 3' end of a poxvirus late mRNA. Virology  
730 255:190–204.
- 731 65. Shuman S, Moss B. 1988. Factor-dependent transcription termination by vaccinia virus RNA  
732 polymerase. Evidence that the cis-acting termination signal is in nascent RNA. J Biol Chem  
733 263:6220–5.
- 734 66. Freitas FB, Frouco G, Martins C, Ferreira F. 2019. The QP509L and Q706L superfamily II RNA  
735 helicases of African swine fever virus are required for viral replication, having non-redundant  
736 activities. Emerg Microbes Infect 8:291–302.
- 737 67. Andrews S. FastQC A Quality Control tool for High Throughput Sequence Data. 0.11.7.  
738 Babraham Bioinformatics, Babraham, England.
- 739 68. Afgan E, Baker D, van den Beek M, Blankenberg D, Bouvier D, Čech M, Chilton J, Clements D,

- 740 Coraor N, Eberhard C, Grüning B, Guerler A, Hillman-Jackson J, Von Kuster G, Rasche E,  
741 Soranzo N, Turaga N, Taylor J, Nekrutenko A, Goecks J. 2016. The Galaxy platform for  
742 accessible, reproducible and collaborative biomedical analyses: 2016 update. *Nucleic Acids*  
743 *Res* 44:W3–W10.
- 744 69. Gruening BA. 2104. Galaxy wrapper.
- 745 70. Blankenberg D, Gordon A, Von Kuster G, Coraor N, Taylor J, Nekrutenko A, Galaxy Team. 2010.  
746 Manipulation of FASTQ data with Galaxy. *Bioinformatics* 26:1783–1785.
- 747 71. Potter J, Zheng W, Lee J. 2003. Thermal stability and cDNA synthesis capability of Superscript  
748 III reverse transcriptase. *Focus (Madison)* 25:19–24.
- 749 72. Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing  
750 reads. *EMBnet.journal* 17:10.
- 751 73. Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic  
752 features. *Bioinformatics* 26:841–842.
- 753 74. RStudio Team. 2016. RStudio: Integrated Development for R. 3.4.3. RStudio, Inc, Boston, MA.
- 754 75. Huber W, Carey VJ, Gentleman R, Anders S, Carlson M, Carvalho BS, Bravo HC, Davis S, Gatto  
755 L, Girke T, Gottardo R, Hahne F, Hansen KD, Irizarry RA, Lawrence M, Love MI, MacDonald J,  
756 Obenchain V, Oleś AK, Pagès H, Reyes A, Shannon P, Smyth GK, Tenenbaum D, Waldron L,  
757 Morgan M. 2015. Orchestrating high-throughput genomic analysis with Bioconductor. *Nat*  
758 *Methods* 12:115–121.
- 759 76. Thodberg M, Thieffry A, Vitting-Seerup K, Andersson R, Sandelin A. 2018. CAGEfightR: Cap  
760 Analysis of Gene Expression (CAGE) in R/Bioconductor. *bioRxiv* 310623.
- 761 77. Thorvaldsdottir H, Robinson JT, Mesirov JP. 2013. Integrative Genomics Viewer (IGV): high-  
762 performance genomics data visualization and exploration. *Brief Bioinform* 14:178–192.
- 763 78. Wheeler DL, Church DM, Federhen S, Lash AE, Madden TL, Pontius JU, Schuler GD, Schriml  
764 LM, Sequeira E, Tatusova TA, Wagner L. 2003. Database resources of the National Center for  
765 Biotechnology. *Nucleic Acids Res* 31:28–33.



- 766 79. Anders S, Pyl PT, Huber W. 2015. HTSeq--a Python framework to work with high-throughput  
767 sequencing data. *Bioinformatics* 31:166–169.
- 768 80. Love MI, Huber W, Anders S. 2014. Moderated estimation of fold change and dispersion for  
769 RNA-seq data with DESeq2. *Genome Biol* 15:550.
- 770 81. Tu SL, Upton C. 2019. Bioinformatics for Analysis of Poxvirus Genomes, p. 29–62. *In* *Methods*  
771 *in Molecular Biology*. Humana Press Inc.
- 772 82. Bailey TL, Elkan C. 1994. Fitting a mixture model by expectation maximization to discover  
773 motifs in biopolymers. *Proceedings Int Conf Intell Syst Mol Biol* 2:28–36.
- 774 83. Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. 2009. Jalview Version 2-A  
775 multiple sequence alignment editor and analysis workbench. *Bioinformatics* 25:1189–1191.
- 776 84. Kettenberger H, Armache K-J, Cramer P. 2003. Architecture of the RNA Polymerase II-TFIIS  
777 Complex and Implications for mRNA Cleavage. *Cell* 114:347–357.
- 778 85. Crooks GE, Hon G, Chandonia J-M, Brenner SE. 2004. WebLogo: A Sequence Logo Generator.  
779 *Genome Res* 14:1188–1190.
- 780 86. Schneider TD, Stephens RM. 1990. Sequence logos: a new way to display consensus  
781 sequences. *Nucleic Acids Res* 18:6097–6100.
- 782 87. Andrés G, García-Escudero R, Viñuela E, Salas ML, Rodríguez JM. 2001. African swine fever  
783 virus structural protein pE120R is essential for virus transport from assembly sites to plasma  
784 membrane but not for infectivity. *J Virol* 75:6758–68.
- 785 88. NCBI Resource Coordinators NR. 2016. Database resources of the National Center for  
786 Biotechnology Information. *Nucleic Acids Res* 44:D7-19.
- 787
- 788
- 789
- 790

791 **Figure 1. Annotated genome of ASFV-BA71V indicating transcription start sites (TSS) and early and late**  
792 **genes.** The map includes 153 previously annotated as well as novel genes identified in this study and their  
793 differential expression pattern from early to late infection from DESeq2 (80) analysis. Early genes (upregulated,  
794 highlighted in dark blue) and late genes (upregulated, dark red) were differentially expressed according to both  
795 RNA-seq and CAGE-seq approaches. The pale blue and pale red marking indicates a negative (early,  
796 downregulated) or positive (late, upregulated) log2 fold change in expression according to both CAGE-seq and  
797 RNA-seq, but is only statistically significant (adjusted  $p$ -value  $< 0.05$ ) from CAGE-seq, due to its higher  
798 sequencing depth and unlike RNA-seq, is not affected by transcription readthrough. Colour coding in white  
799 suggests ambivalency of early and late expression patterns, i.e. not statistically significant according to either of  
800 the methods, or only according to RNA-seq. These also include ten genes with reversed differential expression  
801 between CAGE-seq and RNA-seq results. The map was visualised with the R package gggenes.

802

803 **Figure 2. The ASFV transcriptome including transcription start sites and termination sites.** (a) Whole genome  
804 view of normalized coverage counts per million (CPM) of RNA-seq, 5' CAGE-seq and 3' RNA-seq reads. The  
805 coverage was capped at 16000 CPM. 153 BA71V annotated ORFs are represented as arrows and coloured  
806 according to strand. Peak cluster shape example from F1055L 5' CAGE-seq ends (b) and 3' RNA-seq ends (c)  
807 showing a wide multi-peaked distribution, and J64R 5' CAGE-seq (d) and 3' RNA-seq (e) showing a narrow peak  
808 distribution.

809

810 **Figure 3. Transcriptome mapping aids the reannotation of the ASFV BA71V genome.** (a-left) Summary bar  
811 graph of CAGEfightR TSS clusters and their locations relative to the 153 annotated BA71V ORFs. (a-right) Types  
812 of CAGEfightR clusters detected and the distribution of their respective CAGEfightR scores. (b) Two examples of  
813 ORFs requiring re-annotation following pTSS identification downstream of annotated start codon, encoding  
814 shorter ORFs from the pTSS (I177L, above) or during one expression stage (B169L, below). (c) Examples of two  
815 putative novel genes (pNG3, left and pNG1 right) annotated with the normalized RNA-seq and CAGE-seq read  
816 coverage (CPM) and their genome neighbourhood.

817 **Figure 4. Analysis of alternative pTSS usage in I243L.** (a) Close up of TSSs (CAGE-seq alignments) on the minus  
818 strand at the start of the I243L ORF. Symbols indicate the TSS sites for early (▼), intermediate (●) and late (▽)  
819 gene expression according to Rodríguez *et al.* (26), while E, I and L indicate their respective pTSS positions  
820 concluded from our data. The first 21 AA residues of the annotated I243L ORF are shown, in yellow is the re-  
821 annotated ORF which could be encoded in transcripts initiating from both our annotated Early pTSS. (b)  
822 ClustalW multiple sequence alignment coloured by percentage identity between sequences, illustrated with  
823 Jalview (83), of TFIIS homologues from ASFV (I243L, UniProt: P27948), *A. thaliana* (Q9ZVH8), *D. melanogaster*  
824 (P20232), human (P23193), mouse (P10711) and *S. cerevisiae* (P07273). *S. cerevisiae* TFIIS domain locations  
825 according to Kettenberger *et al.* (84) are shown below the alignment and acidic (DE) catalytic residues are  
826 indicated with ★. ASFV-TFIIS start codons encoded from alternative transcription start sites are labelled as in  
827 (b).

828

829 **Figure 5. Gene expression of ASFV genes during early and late infection.** (a) FPKM values for 20 most highly  
830 expressed ASFV TUs according to CAGE-seq at 5h (left) and 16h (right) post-infection. Genes highlighted in  
831 maroon indicate those encoding proteins which were also found in the 20 most-abundantly expressed ASFV  
832 proteins during infection of either WSL-HP, HEK293 or Vero cells according to proteome analysis done by  
833 Keßler *et al.* (37). Gene functions are shown after their name with TR and PSP referring to predicted  
834 transmembrane region and putative signal peptide, respectively. (b) 20 most-expressed genes during early  
835 (green) and late (blue) infection according to RNA-seq data over gene TU, defined from TSS to ORF stop codon.  
836 (c) MAplot from DESeq2 analysis of CAGE-seq representing the DESeq2 base mean of transcript levels versus  
837 their log2 fold change, with significantly differentially expressed genes in purple (adjusted *p*-value < 0.05). (d)  
838 MAplot representing expression of ASFV TUs including pNGs from DESeq2 analysis of RNA-seq data.

839

840 **Figure 6. Relative expression during infection stages and defining early and late genes.** (a) Boxplot mean FPM  
841 values for the early and late genes at early and late infection, respectively. Outliers are labelled with their gene  
842 names and Wilcoxon rank sum test showed the mean FPM of early genes during early infection was  
843 significantly greater than that of late genes during late infection (*p*-value: 1.865e-06). (b) Distribution of gene  
844 expressed for the least and most expressed genes during early and late infection. Genes in the 15th percentile

845 for their mean FPM values from each time-point, being below an early FPM threshold of 7.56 (blue) and late  
846 FPM of 199.64 (red). (c) Genes in the 85th percentile for their mean FPM values from each time-point, being  
847 above an early FPM threshold of 8148.91 (blue) and late FPM of 4706.27 (red). In dark blue and dark red are  
848 medians for the plotted expression values for early and late infection respectively. (d) Scatter plot comparing  
849 log2fold changes of the 101 significantly differentially expressed genes in common between RNA-seq and  
850 CAGE-seq. Labels were coloured according to their significant upregulation or downregulation from RNA-seq.  
851 (e) Pie chart of gene functional categories downregulated from 5 h to 16 h (36 early genes) and upregulated  
852 from 5 h to 16 h (55 late genes). Fisher test carried out on gene counts for functional groups between early and  
853 late infection, for this all MGF members were pooled into the 'MGFs' functional group.

854

855 **Figure 7. Initiator and promoter sequence signatures of ASFV genes.** WebLogo 3 (85, 86) of aligned early (a)  
856 and late (b) sequences surrounding the Inr (+1) from -35 to +10, with gradients representing the basepair  
857 conservation of the EPM (blue-white), Inr (purple-white) and LPM (peach-white). WebLogo 3 consensus motif  
858 with error-bars, of the 36 early (c) and 55 late (d) gene sequences surrounding their respective pTSSs (5 nt up-  
859 and downstream), i.e. initiator (Inr) motif. (e) EPM located upstream of all 36 of our classified early genes  
860 according to MEME motif search (E-value: 8.2e-021), FIMO with a threshold of  $p$ -value < 1.0 E-4 then identified  
861 at least one iteration of this motif upstream of 81 ASFV genes. (f) Distances of the EPM motif 3' end (nt 19)  
862 relative to the 78 pTSSs (alternative pTSSs excluded). (4). (g) Expression profiles from DESeq2 analysis (log2fold  
863 change vs. base mean expression) of genes with only an EPM from the FIMO search of 60 bp upstream of  
864 pTSSs. Genes for which FIMO detected both EPM and LPM upstream of pTSSs were excluded. Genes in blue  
865 showed a negative log2 fold change (early genes) and in red a positive log2fold change (regardless of  
866 significance). (h) Expression profiles as in c. for the 26 MGFs where an EPM was detected upstream. (i)  
867 Distances of the EPM motif 3' end (nt 19) relative to the MGF pTSSs.

868

869 **Figure 8. Promoter motif upstream of ASFV late genes.** (a) The LPM detected upstream of 17 of our classified  
870 late genes from MEME motif search (E-value: 1.6e-003). (b) Distances from a FIMO search (threshold  $p$ -value <  
871 1.0E-4) identified the LPM upstream of 53 ASFV genes (excluding those with alternative pTSSs), motif distances  
872 from pTSSs are represented as a bar chart. (c) Expression profiles as in c. of genes with only an LPM from the  
873 FIMO search of 60 bp upstream of pTSSs. (d) The eukaryotic TATA-box motif which was one of 28 hits in a

874 Tomtom search of the LPM. (e) 5' UTR lengths in nt of the 91 early (mean: 39, median: 14) or late (mean: 25,  
875 median: 9) classified ASFV genes, starting from the most upstream pTSS (in the case of alternating pTSSs) until  
876 the first ATG start codon nt, represented as a notched boxplot. 9 genes with 5' UTR's above 80 nt were  
877 excluded from the boxplot: QP509L (92 nt long), pNG2 (105 nt), I267L (110 nt), B318L (118 nt), C44L (131 nt),  
878 DP141L (165 nt), pNG1 (223 nt), EP402R (242 nt) and A118R (332 nt). (f) Percentage AT content of early (mean:  
879 69.0, median: 70.9) and late (mean: 81.7, median: 83.3) 5' UTRs, omitting those of 0 length.

880

881 **Figure 9. Investigating ASFV-RNAP slippage.** (a) Frequency of different lengths of template-free extensions in  
882 early and late stage samples. (b) Relationship between the length of templated 5' UTRs and fraction of  
883 template-free extensions. Gene 5' UTRs split into 36 early (blue), 55 late (orange) and not-classified ('NC',  
884 green). (c) Frequency of most common template-free extensions in the early and late stage samples. (d)  
885 Sequence logo of region surrounding TSS of 'AU' and 'AUAU'-extended transcripts.

886

887 **Figure 10. ASFV transcription termination.** (a) WebLogo 3 motif of 10 nt upstream and 10 downstream of all  
888 pTTS and npTTSs with a polyT upstream with  $\geq 4$  consecutive Ts based on 126 TTSSs. (b) Distance from 3' terminal  
889 T in polyT motif to the TTS (median). (c) The distribution of polyT lengths among 126 polyT TTSSs (median: 7),  
890 split into expression stages according to CAGE-seq differential expression analysis (NC: not-classified), showing  
891 late gene polyTs are shorter in length (Wilcoxon rank sum test,  $p$ -value: 0.0216). (d) Distribution of gene  
892 expression types among the 83 polyT pTTSs and 31 non-polyT pTTSs. Dotted lines labels indicate Fisher test  $p$ -  
893 values of gene types between the two pTTS-types, classified from CAGE-seq. (e) 55 Early and 53 late gene 3'  
894 UTR lengths from stop codon to pTTS (Wilcoxon rank sum test,  $p$ -value: 0.003).

## 895 Tables

ORF	Strand	pTSS Coordinate	Corrected Start Codon Coordinate	ORF Length	Comment
K93L	-	2131	2122	83	Alternative ATG codon 30 nt downstream. Another strong TSS was detected at 2037- whose transcripts would encode a 36 AA protein.
F165R	+	42354	42359	136	Alternative ATG codon 63 nt downstream.
C84L	-	64618	64492   64616	38   76	38 AA ORF was in-frame with original C84L start codon. 76 AA ORF encoded from first ATG after pTSS.
G1211R	+	96370	96377	1207	Alternative ATG codon 12 nt downstream.
CP204L	-	108573	108567	196	Alternative ATG codon 24 nt downstream.
CP312R	+	110491	110501	307	Alternative ATG codon 15 nt downstream.
I177L	-	L: 157857	157849	66	Strongest pTSS only detected in late time- point.
DP93R	+	167971	167980	83	Alternative ATG codon 30 nt downstream.
EP402R	+	56862	57104   56991	115   148	Encodes 115 AA in-frame with original EP402R start codon. 148 AA alternative ORF encoded from first ATG after pTSS.
B169L	-	E: 80983   L: 81025	81018   80745	169   78	Late pTSS can produce full-length B169L and early pTSS: 78 AA.
I243L	-	E: 155122 L: 155124	E/I: 155119 L: 154969	243   191	Late pTSS produces shorter transcript with closest downstream ATG encoding a

		L: 155115			shorter 191 AA protein.
--	--	-----------	--	--	-------------------------

896

897 **Table 1. Summary of ASFV genes where pTSS locations guided the re-annotation of ORFs.** For B169L and

898 I243L, the letters E, I and L refer to alternative pTSSs from early, intermediate and late infection, respectively

899 reported by Rodríguez *et al.* (26).

900

Gene	Early pTSS	Late pTSS	Function
X69R	11315	11280	Uncharacterised
J154R	14174	14150	MGF 300-2R
EP1242L	53125	53135	ASFV-RPB2
C315R	70137	70131	ASFV-TFIIB
CP80R	110208	110191	ASFV-RPB10
D345L	129357	129257	Lambda-like exonuclease (7)
E120R	150949	150911	Structural protein (87)

901

902 **Table 2. Alternative pTSS usage during early and late ASFV infection.** List of ASFV genes with alternative pTSSs

903 used in early and late infection.

Putative Gene	Strand	Transcription Start*	Transcription End**	Putative protein length (AA)	Similarity according to NCBI Blast	Gene-End oligoT (nt)
pNG1	+	13053	13435	25	13 residues had 92% identity to ASFV-G-ACD-00350 (AZP54308.1), E-value: 0.11	6
pNG2	-	30091	29827	50	50 residues had 100% identity with ASFV26544 00600 (AKM05534.1)	8
pNG3	+	12664	12896	44	38 residues had 59% identity to ASFV-G-ACD-00290 (AZP54130.1), E-value: 0.13	6
pNG4	+	10583	10835	44	42 residues had 65% identity with ASFV-G-ACD-00290 (AZP54130.1), E-value: 1e-09.	6
pNG5	+	29817	<u>30080</u>	31	No significant similarity.	None
pNG6	+	167005	167336	56	56 residues aligned with 40% identity to pKP93L (AIY22188.1), E-value: 6e-07	5
pNG7	+	10484	<u>10616</u>	32	32 residues aligned with a 31 AA hypothetical protein with from ASFV Belgium 2018/1 (BioProject: PRJEB31287) 87% identity (VFV47940.1), E-value: 8e-10.	3

904

905 **Table 3. Details of seven novel ASFV candidate genes.** NCBI ORFfinder and BLAST were used to predict the  
906 putative encoded ORFs and subsequently analysed for putative homologous sequences (78, 88). \*: defined as  
907 pTSS from CAGE-seq. \*\*: defined from 3' RNA-seq, underlined transcription ends defined from only RNA-seq.  
908 pNG5 is in the antisense orientation relative to pNG2, and the RNA-3' end of pNG6 is fuzzy according to RNA-  
909 seq and may overlap with DP42R. pNG7 is overlapping pNG4 on the same strand.























