1 African Swine Fever Virus and host response - transcriptome

2 profiling of the Georgia 2007/1 strain and porcine macrophages

3

4 Gwenny Cackett[a§], Raquel Portugal[b§] , Dorota Matelska[a], Linda Dixon[b#] and Finn Werner[a#]

5 [a]Institute for Structural and Molecular Biology, Darwin Building, University College London, Gower

6 Street, London WC1E 6BT, United Kingdom

7 [b]Pirbright Institute, Ash Road, Pirbright, Surrey, GU24 0NF, United Kingdom

8 [§] have contributed equally to this work

9

10 Short Title:

11 The ASFV Georgia 2007/1 Strain Transcriptome

12 #Address correspondence to Linda Dixon, linda.dixon@pirbright.ac.uk, or Finn Werner,

13 f.werner@ucl.ac.uk.

1

## Abstract [222 words]

African swine fever virus (ASFV) has a major global economic impact. With a case fatality in domestic pigs approaching 100%, it currently presents the largest threat to animal farming. Although genomic differences between attenuated and highly virulent ASFV strains have been identified, the molecular determinants for virulence at the level of gene expression have remained opaque. Here we characterise the transcriptome of ASFV genotype II Georgia 2007/1 (GRG) during infection of the physiologically relevant host cells, porcine macrophages. In this study we applied Cap Analysis Gene Expression sequencing (CAGE-seq) to map the 5' ends of viral mRNAs at 5 and 16 hours post-infection. A bioinformatics analysis of the sequence context surrounding the transcription start sites (TSSs) enabled us to characterise the global early and late promoter landscape of GRG. We compared transcriptome maps of the GRG isolate and the lab-attenuated BA71V strain that highlighted GRG virulent-specific transcripts belonging to multigene families, including two predicted MGF 100 genes I7L and I8L. In parallel, we monitored transcriptome changes in the infected host macrophage cells. Of the 9,384 macrophage genes studied, transcripts for 652 host genes were differentially regulated between 5 and 16 hours-post-infection compared with only 25 between uninfected cells and 5 hours post-infection. NF-kB activated genes and lysosome components like S100 were upregulated, and chemokines such as CCL24, CXCL2, CXCL5 and CXCL8 downregulated.

## Importance [183 words]

African swine fever virus (ASFV) causes haemorrhagic fever in domestic pigs with case fatality rates approaching 100%, and no approved vaccines or antivirals. The highly-virulent ASFV Georgia 2007/1 strain (GRG) was the first isolated when ASFV spread from Africa to the Caucasus region in 2007. Then spreading through Eastern Europe, and more recently across Asia. We used an RNA-based next generation sequencing technique called CAGE-seq to map the starts of viral genes across the GRG DNA genome. This has allowed us to investigate which viral genes are expressed during early or late stages of infection and how this is controlled, comparing their expression to the non-virulent ASFV-BA71V strain to identify key genes that play a role in virulence. In parallel we investigated how host cells respond to infection, which revealed how the ASFV suppresses components of the host immune response to ultimately win the arms race against its porcine host.

## Introduction [1,317 words]

ASFV originated in Sub-Saharan Africa where it remains endemic. However, following the introduction in 2007, of a genotype II isolate to Georgia (1), and subsequent spread in Russia and Europe. The virus

2

45  was then introduced to China in 2018 (2), from here it spread rapidly across Asia, strongly emphasizing

46  this disease as a severe threat to global food security. ASFV is the only characterised member of the

47  Asfarviridae family (3) in the recently classified Nucleocytoviricota (ICTV Master Species List 2019.v1)

48  phylum (4,5). ASFV has a linear double-stranded DNA (dsDNA) genome of ~170–193 kbp encoding

49  ~150–~200 open reading frames (ORFs). Until recently, little was known about either the transcripts

50  expressed from the ASFV genome or the mechanisms of ASFV transcription. Much of what is known

51  about transcription is extrapolated from vaccinia virus (VACV), a distantly-related Nucleocytoviricota

52  member, from the Poxviridae family (6). ASFV encodes a eukaryotic-like 8-subunit RNA polymerase

53  (RNAP), an mRNA capping enzyme and poly-A polymerase, all of which are carried within mature virus

54  particles (7). These virions are transcription competent upon solubilisation in vitro (8) and support

55  mRNA modification by including a 5'-methylated cap and a 3' poly-adenylated (polyA) tail of ~33

56  nucleotide-length (8,9).

57  Viral genes are typically classified according to their temporal expression patterns - ASFV genes have

58  historically been categorised as 'immediate early' when expressed immediately following infection, as

59  'early genes' following the onset of viral protein synthesis, as 'intermediate genes' after the onset of

60  viral DNA replication, or as 'late genes' thereafter. The temporal regulation of transcription is likely

61  enabled by different sets of general transcription initiation factors that recognise distinct early or late

62  promoter motifs (EPM or LPM, respectively), as we previously investigated in the ASFV-BA71V strain

63  (10), and address further in this study. EPM recognition is likely enabled by the ASFV homologue of

64  heterodimeric VACV early transcription factor (VETF), consisting of D1133L (D6) and G1340L (A7) gene

65  products, which bind the Poxvirus early gene promoter motif (11–13), which the ASFV EPM strongly

66  resembles. Both ASFV-D6 and ASFV-A7 are late genes, i.e. synthesised late during infection (10) and

67  packaged into virus particles (7). The ASFV LPM is less well defined than the EPM, but a possible

68  initiation factor involved in its recognition is the ASFV-encoded viral homolog of the eukaryotic TATA-

69  binding protein (TBP), expressed during early infection (10). By analogy with the VACV system,

70  additional factors including homologs of A1, A2 and G8 may also contribute to late transcription

71  initiation (6)

72  We have recently carried out a detailed and comprehensive ASFV whole genome expression analysis

73  using complimentary next-generation sequencing (NGS) results and computational approaches to

74  characterise the ASFV transcriptome following BA71V infection of Vero cells at 5 hpi and 16 hpi post-

75  infection (hpi) (10). Most of our knowledge about the molecular biology of ASFV, including gene

76  expression, has been derived from cell culture-adapted, attenuated virus strains, such as BA71V

77  infecting Vero tissue culture cells (9,10). These model systems provide convenient models to study

3

78   the replication cycle but have deletions of many genes that are not essential for replication, but have

79   important roles in virulence within its natural porcine hosts. (14–16). To date 24 ASFV genotypes have

80   been identified in Africa (16–23), while all strains spreading across Asia and Europe belong to the Type

81   II genotype. Most of these are highly virulent in domestic pigs and wild boar, including the ASFV

82   Georgia 2007/1 (GRG) (24), and the Chinese ASFV Heilongjiang, 2018 (Pig/HLJ/18) (25) isolates.

83   Though a number of less virulent isolates have been identified in wild boar in the Baltic States and

84   domestic pigs in China (26–29). It is crucial to understand the similarities and commonalities between

85   ASFV strains, and to characterise the host response to these in order to understand the molecular

86   determinants for ASFV pathogenicity. Information about the gene content and genome organisation

87   can be gained from comparing virus genome sequences. However, only functional genomics such as

88   transcriptome or proteome analyses can provide information about the differences in gene expression

89   programmes and the host responses to infection.

90   On the genome level, most differences between virulent (e.g. GRG) and attenuated (e.g. lab-

91   attenuated BA71V) ASFV strains reside towards the genome termini. Figure 1a shows a whole genome

92   comparison of GRG (left) and BA71V (right) strains with the sequence conservation colour coded in

93   different shades of blue. The regions towards the ends of the genome are more dynamic compared to

94   the central region which is highly conserved, as genes at the termini are prone to deletion, duplication,

95   insertion, and fusion (17,30). Most of the GRG-specific genes are expressed early during infection

96   (early genes are colour coded blue in the outer arch of Figure 1a) and many belong to Multi-Gene

97   Families (MGFs, purple in the inner arch). The functions of many MGF members remain poorly

98   understood, though variation among MGFs is linked to virulence (31) and deleting members of MGF

99   360 and 505 families has been shown to reduce virulence (32,33). Deletion of MGF 505-7R or MGF

100  110-9L also partially attenuated the virus in pigs (34,35). In contrast deletion of MGF 110-1L and MGF

101  100-1R did not reduce virus virulence (17). Members of MGF 110 are highly expressed both on the

102  mRNA and protein level in infections with the BA71V isolate or OURT88/3 (10,36), suggesting MGF

103  110 holds importance during infection. Overall, the functions of MGF 360 and 505 members are better

104  characterised than other MGFs, playing a role in evading the host type I interferon (IFN) response

105  (15,32,37–40). In summary, comparing the expression of ASFV genes, especially MGFs, between the

106  virulent GRG- and the lab adapted BA71V strains, is fundamental in identification of virulence factors

107  and better MGF characterisation.

108  Macrophages are the primary target cells for ASFV, they are important immune effector cells that

109  display remarkable plasticity allowing efficient response to environmental signals (41).  There are

110  some studies which have investigated how host macrophages respond to infection, including a

4

111     microarray analysis of primary swine macrophage cells infected with virulent GRG (42). There are two

112     RNA-seq studies of whole blood or tissues isolated from pigs post-mortem, which were infected with

113     either a low pathogenic ASFV-OURT 88/3 or ASFV-GRG (43), or infected with a pathogenic Chinese

114     isolate ASFV-SY18 (44). Recently, two reports have been published about the transcriptomic response

115     of porcine macrophages to infection with a virulent Chinese genotype II isolate using a low multiplicity

116     of infection (MOI: 1) and classical RNA-seq (45,46), but due to different experimental conditions the

117     varying results are somewhat challenging to compare with other studies. It must also be remembered

118     that neither these classical RNA-seq nor microarray analyses, have sufficient resolution to accurately

119     capture viral gene expression in the compact ASFV genome alongside that of the host.

120     Here we applied CAGE-seq to characterise the transcriptome of the highly virulent GRG isolate (24),

121     in primary porcine macrophages, the biologically relevant target cells for ASFV infection. In this study

122     we used a high multiplicity of infection (MOI: 5), so that transcripts expressed during a single cycle

123     time course could be measured without the complication of variable proportions of uninfected cells

124     being present. We investigated the differential gene expression patterns of viral mRNAs at early and

125     late time points of 5- and 16 hpi, and mapped the viral promoter motifs. Importantly, we have

126     compared the expression levels and temporal regulation of genes conserved in both the virulent GRG

127     isolate, and attenuated tissue-culture adapted BA71V strain. With a few exceptions, both mRNA

128     expression levels and temporal regulation of the conserved genes are surprisingly similar. This

129     confirms that it is not deregulation of their conserved genes, but the virulent isolate-specific genes,

130     which are the key determinants for ASFV virulence. Most of these genes are MGF members, likely

131     involved in suppression of the host immune-response. Indeed, transcriptome analysis of the porcine

132     macrophages upon GRG infection reflects a modulation of host immune response genes, although the

133     bulk of the ~ 9000 genes studied did not significantly change expression levels during infection.

## Results [4,633 words]

### Genome-wide Transcription Start Site-Mapping

136     We infected primary porcine alveolar macrophages with ASFV GRG at a high multiplicity of infection

137     (MOI 5.0), isolated total RNA at 5 hpi and 16 hpi and sequenced using CAGE-seq (Supplementary Table

138     1a). The resulting mRNA 5' ends were mapped to the GRG genome (Figure 1b) resulting in the

139     annotation of 229 and 786 TSSs at 5 and 16 hpi, respectively (Figure 1c and d, from Supplementary

140     Table 1b and c, respectively). The majority of TSSs were identified within 500 bp upstream of the start

141     codon of a given ORF, a probable location for a *bona fide* gene TSS. The strongest and closest TSSs

142     upstream of ORFs were annotated as 'primary' TSS (pTSS, listed in Supplementary Table 1d) and in this

5

143  manner we could account TSS for 177 out of 189 GRG ORFs annotated in the FR682468.1 genome.

144  TSSs signals below the threshold for detection included MGF_110-11L, C62L, and E66L, the remainder

145  being short ORFs designated as 'ASFV_G_ACD', predicted solely from the FR682468 genome sequence

146  (24). The E66L ORF was originally predicted from only the BA71V genome sequence, but likewise was

147  undetectable with CAGE-seq (10), making its expression unlikely. Our TSS mapping identified novel

148  ORFs (nORFs) downstream of the TSS, which were included in the curated GRG genome map

149  (Supplementary Table 1d includes pTSSs of annotated ORFs and nORFs in gene feature file or 'GFF'

150  format). In addition to ORF-associated TSSs, some were located within ORFs (intra-ORF or ioTSS), or

151  in between them (inter-ORF TSS), and all detected TSSs are listed in Supplementary Table 1b-c.

152  ## Expression of GRG genes during Early and Late Infection

153  Having annotated TSSs across the GRG genome, we quantified the viral mRNAs originating from pTSSs

154  from CAGE-seq data, normalising against the total number of reads mapping to the ASFV genome (i.

155  e. RPM or reads per million mapped reads per sample). We compared gene expression between early

156  and late infection, and simplistically defined genes as 'early' or 'late' if they are significantly down- or

157  upregulated (respectively), using DESeq2 (47). In summary, 165 of the 177 detectable genes were

158  differentially expressed (adjusted p-value or padj < 0.05, Supplementary Table 1e). Those showing no

159  significant change were D345L, DP79L, I8L, MGF_100-1R, A859L, QP383R, B475L, E301R, DP63R,

160  C147L, and I177L. 87 of those 165 differentially expressed genes were significantly downregulated,

161  thus representing the 'early genes', while 78 of the 165 genes were upregulated or 'late genes'. The

162  majority of MGFs were early genes, apart from MGF 505-2R, MGF 360-2L and MGF 100-1L (Figure 2a).

163  Figure 2b shows the expression patterns of GRG-exclusively expressed genes, which we defined as

164  only having a detectable CAGE-seq TSS in GRG, and not in BA71V (regardless of presence in the BA71V

165  genome). These unsurprisingly, consist of many MGFs (19), all of which were early genes (Figure 2b),

166  barring MGF 100-1L. In addition genes I9R, I10L and I11L and several of the newly annotated short

167  ORFs were specific to GRG.

168  We extracted the top twenty most highly expressed genes of GRG (as RPM) during 5 hpi (Figure 2c)

169  and 16 hpi (Figure 2d) post-infection. Ten genes are shared between both top 20 lists: MGF 110-3L,

170  A151R, MGF 110-7L, MGF 110-5L-6L, I73R, 285L, CP312R, ASFV_G_ACD_00600, MGF 110-4L, and

171  CP204L. It is important to note that the relative expression values (RPM) for genes at 5 hpi are

172  significantly higher than those at 16 hpi. This is consistent with our observations in the BA71V strain

173  (10) and due to the increase in global viral transcript levels during late infection discussed below.

174  Supplementary Table 1f includes all the GRG annotated ORFs, their TSS locations during early and late

6

175    infection, their relative distances if these TSS locations differ, and their respective 5' Untranslated

176    Region (UTR) lengths.

### GRG and BA71V Share Strong Similarity between Conserved Gene Expression

178    Next we carried out a direct comparison of mRNA levels from 132 conserved genes between the

179    virulent GRG and attenuated BA71V (10) strain making use of our previously published CAGE-seq data.

180    The relative transcript levels (RPM) of the genes conserved between the two strains showed a

181    significant correlation at 5 hpi (Figure 3a) and 16 hpi (Figure 3b), supported by the heatmap in

182    Supplementary Figure 1, the RPM for each gene, across both time-points and replicates, showing a

183    strong congruence between the two strains. Of the 132 conserved genes, 125 showed significant

184    differential expression in both strains. 119 of these 125 showed the same down- or up-regulated

185    patterns of significant differential expression from 5 hpi to 16 hpi (Figure 3c, early genes in blue, late

186    genes in red). The exceptions are D205R, CP80R, C315R, NP419L, F165R, and DP148R (MGF 360-18R),

187    encoding RNA polymerase subunits RPB5 and RPB10 (15), Transcription Factor IIB (TFIIB) (15), DNA

188    ligase (48), a putative signal peptide-containing protein, and a virulence factor (49), respectively. The

189    ASFV-TFIIB homolog (C315R) is classified as an early gene in GRG but not in BA71V, in line with the

190    predominantly early-expressed TBP (B263R), its predicted interaction partner. It is worth noting

191    however, that D205R, CP80R, and C315R are close to the threshold of significance, with transcripts

192    being detected at both 5 hpi and 16 hpi (Supplementary Table 1e).

### Increased and pervasive transcription during late infection

194    During late infection of BA71V (10), we noted an increase in genome-wide mRNA abundance, as well

195    as an increasing number of TSSs and transcription termination sites, reminiscent of pervasive

196    transcription observed during late infection of Vaccinia virus (50). To quantify and compare the global

197    mRNA increase both in BA71V and GRG, we calculated the ratio of read coverage per nucleotide, at

198    16 hpi versus 5 hpi (log2 transformed ratio of RPM), across the viral genome (Figure 4a, increase shown

199    above- and decrease below the x-axis). This dramatic increase is due to the overall increase of virus

200    mRNAs present, which is visible in both strains (Figure 4b), with a ~2 fold increase in GRG from 5 hpi

201    to 16 hpi, versus ~8 fold in BA71V (Figure 4c).

202    This observation can at least in part be attributed to the larger number of viral genomes during late

203    infection, with increased levels of viral RNAP and associated factors available for transcription,

204    following viral protein synthesis. Viral DNA-binding proteins, such as histone-like A104R (51), may

205    remain associated with the genome originating from the virus particle in early infection. This could

206    suppress spurious transcription initiation, compared to freshly replicated nascent genomes that are

7

207   highly abundant in late infection. In order to test whether the increased mRNA levels correlated with

208   the increased number of viral genomes in the cell, we determined the viral genome copy number by

209   using quantitative PCR (qPCR against the p72 capsid gene sequence) using purified total DNA from

210   infected cells isolated at 0 hpi, 5 hpi and 16 hpi, and normalized values to the total amount of input

211   DNA. Using this approach, we observed genome copy levels that were consistent from 0 hpi to 5 hpi,

212   consistent with this being pre-DNA replication, followed by a substantial increase at 16 hpi, which was

213   more pronounced in BA71V infection (Figure 4d). This corresponded to a 15-fold increase in GRG

214   genome copy numbers from late, compared to early times post-infection of porcine macrophages,

215   and a 30-fold increase in BA71V during infection of Vero cells (Figure 4e). In summary, the ASFV

216   transcriptome changes both qualitatively and quantitatively as infection progresses, and the increase

217   of virus mRNAs during late infection is accompanied by the dramatic increase in viral genome copies.

218   Interestingly, the increase in viral transcripts and genome copies was less dramatic in the virulent GRG

219   strain.

### Correcting the bias of temporal expression pattern

221   The standard methods of defining differential gene expression are well established in transcriptomics

222   using programs like DESeq2 (47). This is a very convenient and powerful tool which captures the

223   nuances of differential expression in complex organisms. However, virus transcription is often

224   characterised by more extreme changes, typically ranging from zero to millions of reads. Furthermore,

225   in both BA71V and GRG strains the genome-wide mRNA levels and total ASFV reads increase over the

226   infection time course (Figure 4 and Supplementary Table 1a). As a consequence, such normalisation

227   against the total mapped transcripts per sample (RPM) generates overestimated relative expression

228   values at 5 hpi, and understates those at 16 hpi (10). In order to validate the early-late expression

229   patterns derived from CAGE-seq, we carried out RT-PCR for selected viral genes, as this signal is

230   proportionate to the number of specific mRNAs regardless of the level of other transcripts – with the

231   minor caveat that it can pick up readthrough transcripts from upstream genes. We tested differentially

232   expressed conserved genes including GRG early- (MGF 505-7R, MGF 505-9R, NP419L), and D345L

233   which showed stable relative expression values (RPM values in Figure 1e). All selected genes showed

234   a consistently stronger RT-PCR signal during late infection in both BA71V and GRG (Figure 5a-d). The

235   exception is NP419L whose levels were largely unchanged, and this is an example of how a gene whose

236   transcript levels remain constant would be considered downregulated, when almost all other mRNA

237   levels increase (Figure 5b).

238   The standard normalisation of NGS reads against total mapped reads (RPM) is regularly used as it

239   enables a statistical comparison between samples and conditions, subject to experimental variations

   8

240 (52). Keeping this in mind, we used an additional method of analysing the 'raw' read counts to

241 represent global ASFV transcript levels that are not skewed by the normalisation against total mapped

242 reads. Figure 5 shows a side-by-side comparison of RT-PCR results, and the CAGE-seq data normalised

243 (RPM) or expressed as raw counts, beneath each RT-PCR gel. Unlike CAGE-seq, RT-PCR will detect

244 transcripts originating from read-through of transcripts initiated from upstream TSS including intra-

245 ORF TSS (ioTSSs). To detect such 'contamination' we used multiple primer combinations in upstream

246 and downstream segments of the gene (Figure 5c, cyan and yellow arrows) to capture and account for

247 possible variations. Overall, our comparative analyses shows that the normalised data (RPM) of early

248 genes such as MGF 505-7R and 9R indeed skews and overemphasises their early expression, while the

249 raw counts are in better agreement with the mRNA levels detected by RT-PCR. In contrast, late genes

250 such as NP419L and D345L would be categorised as late using all three quantification methods, in

251 agreement with GRG CAGE-seq but not BA71V from Figure 3c. We validated the expression pattern of

252 the early GRG-specific gene MGF 360-12L (Figure 5e). While the RPM values indicated a very strong

253 decrease in mRNA levels from early to late time points, the decrease in raw counts was less

254 pronounced and more congruent with the RT-PCR analysis, showing a specific signal with nearly equal

255 intensity during early and late infection. Lastly, we used qRT-PCR to quantify C315R transcript levels,

256 as this was close to the early vs late threshold, (a log2fold change of 0 in Figure 3c), which showed

257 again that qRT-PCR better agreed with the raw counts.

258 ## An improved temporal classification of ASFV genes

259 Based on the considerations above, we prepared a revised classification of temporal gene expression

260 of the genes conserved between the two strains based on raw counts. The heatmap in Figure 6a shows

261 the mRNA levels at early and late infection stages of BA71V and GRG strains (all in duplicates) with the

262 genes clustered into five subcategories (1 to 5, Figure 6a) according to their early and late expression

263 pattern, which are shown in Figure 6b. Genes that are expressed at high or intermediate levels during

264 early infection but that also show high or intermediate mRNA levels during late infection are classified

265 as 'early' genes belonging to cluster-1 (8 genes, levels: high to high, H-H), cluster-4 (33 genes, mid to

266 mid, M-M) and cluster-5 (16 genes, low-mid to low-mid, LM-LM). Genes with low or undetectable

267 mRNA levels during early infection, which increase to intermediate or high levels during late infection

268 are classified as 'late' genes and belong to cluster-2 (15 genes, low to high, L-H) and cluster-3 (60

269 genes, low to mid, L-M), respectively. Overall, the clustered heatmap based on raw counts shows a

270 similar but more emphasised pattern compared to the normalised (RPM) data (compare Figure 6 and

271 Supplementary Figure 1). Calculating the percentage of reads per gene, which can be detected at 16

272 hpi compared to 5 hpi, reveals only a small number of genes have most ( ≥70%) of their reads

9

273    originating during early infection: 30 genes in the GRG strain and 5 genes in the BA71V strain. For over

274    half of the BA71V-GRG conserved genes, 90-100 % of reads can be detected during late infection

275    (Figure 6c). For all GRG genes, this generates a significant difference between the raw counts per gene

276    between time-points (Figure 6d).

277    Below we discuss specific examples of genes subcategorised in specific clusters. I73R is among the top

278    twenty most-expressed genes during both early and late infection according to the normalised RPM

279    values (Figure 2c and d) resides in cluster-1 (H-H) (Figure 6a). While I73R is expressed during early

280    infection, the mRNA levels remain high with >1/3 of all reads detected during late infection in both

281    strains when calculated as raw counts (34 % in GRG and 45 % in BA71V). This new analysis firmly

282    locates I73R into cluster-1 (H-H) and is classified confidently as early gene. Notably, our new approach

283    results in biologically meaningful subcategories of genes that are likely to be coregulated, e. g. the

284    eight key genes that encode the ASFV transcription system including RNAP subunits RPB1 (NP1450L),

285    RPB2 (EP1242L), RPB3 (H359L), RPB5 (D205R), RPB7 (D339L) and RPB10 (CP80R), the transcription

286    initiation factor TBP (B263R) and the capping enzyme (NP868R) belong to cluster-4 (M-M), and

287    transcription factors TFIIS (I243L) and TFIIB (C315R) belong to cluster-5 (LM-LM). The overall mRNA

288    levels of cluster-4 and -5 genes are different, but remain largely unchanged during early and late

289    infection, consistent with the transcription machinery being required throughout infection. In

290    contrast, the mRNAs encoding the transcription initiation factors D6 (D1133L) and A7 (G1340L) are

291    only present at low levels during early- but increase during late infection and thus belong to cluster-3

292    (L-M), classifying them as late genes. This is meaningful since the heterodimeric D6-A7 factor is

293    packaged into viral particles (7), presumably during the late stage of the infection cycle. The mRNAs

294    of the major capsid protein p72 (B646L) and the histone-like-protein A104R (51,53) follow a similar

295    late pattern but are present at even higher levels during late infection and therefore belong to cluster-

296    2 (L-H).

### Architecture of ASFV promoter motifs

298    In order to characterise early promoter motifs (EPM) in the GRG strain, we extracted sequences 35 bp

299    upstream of all early gene TSSs and carried out multiple sequence alignments. As expected, this region

300    shows a conserved sequence signature in good agreement with our bioinformatics analyses of EPMs

301    in the BA71V strain, including the correct distance between the EPM and the TSS (9-10 nt from the

302    EPM 3' end) and the 'TA' motif characteristic of the early gene Initiator (Inr) element (Figure 7a) (10).

303    A motif search using MEME (54) identified a core (c)EPM motif with the sequence 5'-AAAATTGAAT-3'

304    (Figure 7b), within the longer EPM. The cEPM is highly conserved and is present in almost all promoters

305    controlling genes belonging to cluster-1, -4 and -5 (Supplementary Table 3). A MEME analysis of

10

306    sequences 35 bp upstream of late genes (Figure 7c), provided a 17-bp AT-rich core late promoter motif

307    (cLPM, Figure 7d), however, this could only be detected in 46 of the late promoters.

308    In an attempt to improve the promoter motif analyses and deconvolute putative sequence elements

309    further, we probed the promoter sequence context of the five clusters (clusters 1-5 in Figure 7e-i,

310    respectively) of temporally expressed genes with MEME (Supplementary Table 3). The early gene

311    promoters of clusters-1 (H-H), -4 (M-M) and -5 (LM-LM) are each associated with different expression

312    levels, and all of them contain the cEPM located 15-16 nt upstream of the TSS with two exceptions

313    that are characterized by relatively low mRNA levels (Figure 7k). Interestingly, cluster-2 (L-H)

314    promoters are characterized by a conserved motif with significant similarity to eukaryotic TATA-box

315    promoter element that binds the TBP-containing TFIID transcription initiation factor (Figure 7f

316    highlighted with red bracket, detected via Tomtom (55) analysis of the MEME motif output). Cluster-

317    3 (L-M) promoters contain a long motif akin to the cLPM, derived from searching all late gene

318    promoter sequences, and which is similar to the LPM identified in BA71V (Figure 7d and g, green

319    bracket). All motifs described in the cluster analysis above could be detected with statistically

320    significance (p-value < 0.05) via MEME, in every gene in each respective cluster with only two

321    exceptions: MGF 110–3L from cluster-1, and MGF 360-19R from cluster-4, for the latter see details

322    below.

### Updating Genome Annotations using Transcriptomics Data

324    TSS-annotation provides a useful tool for re-annotating predicted ORFs in genomes like ASFV (10)

325    where many of the gene products have not been fully characterized and usually rely on prediction

326    from genome sequence alone. We have provided the updated ORF map of the GRG genome in GFF

327    format (Supplementary Table 1f). This analysis identified an MGF 360-19R ortholog (Figure 8),

328    demonstrating how transcriptomics enhances automated annotation of ASFV genomes by predicting

329    ORFs from TSSs. The MGF 360-19R was included in subsequent DESeq2 analysis showing it was not

330    highly nor significantly differentially expressed (Supplementary Table 1e). Another important feature

331    is the identification of intra-ORF TSSs (ioTSSs) within MGF 360-19R that potentially direct the synthesis

332    of N-terminally truncated protein variants expressed either during early or late infection. The presence

333    of EPM and LPM promoter motifs lends further credence to the ioTSSs (Figure 8). Similar truncation

334    variants were previously reported for I243L and I226R (56) and in BA71V (10). In addition, we detected

335    multiple TSSs within MGF 360-19R encoding very short putative novel ORFs (nORF) 5, 7 or 12 aa

336    residues long; since these ioTSSs were present in both early and late infection they are not all likely to

337    be due to pervasive transcription during late infection.

11

338     We investigated the occurrence of ioTSS genome wide and uncovered many TSSs with ORFs

339     downstream that were not annotated in the GRG genome (Supplementary Table 2a). These ORFs

340     could be divided into sub-categories: in-frame truncation variants (Supplementary Table 2b, akin to

341     MGF 360-19R in Figure 8), nORFs (Supplementary Table 2c), and simply mis-annotated ORFs. All

342     updated annotations are found in Supplementary Table 1f. Putative truncation variants generated

343     from ioTSSs were predominantly identified during late infection, suggesting these could be a by-

344     product of pervasive transcription. Therefore, those detected early or throughout infection are

345     perhaps more interesting, they span a variety of protein functional groups, and many gene-products

346     are entirely uncharacterised (Figure 9a). The truncation variants additionally showed a size variation

347     of 5'-UTRs between the ioTSSs and downstream start codon (Figure 9b). An example of a mis-

348     annotation would be CP204L (Phosphoprotein p30, Figure 9c) gene that is predicted to be 201 residues

349     long. The TSS determined by CAGE-seq and validated by Rapid Amplification of cDNA Ends (5'-RACE)

350     is located downstream of the annotated start codon; based on our results we reannotated the start

351     codon of CP204L which results in a shorter ORF of 193 amino acids (Figure 9c).

352     Our GRG TSS map led to the discovery of many short nORFs, which are often overlooked in automated

353     ORF annotations due to a minimum size, e. g. 60 residues in the original BA71V annotation (15). Some

354     short ORFs have been predicted for the GRG genome including those labeled 'ASFV_G_ACD' in the

355     Georgia 2007/1 genome annotation (19). However, their expression was not initially supported by

356     experimental evidence, though we have now demonstrated their expression via CAGE-seq (Figure 2b,

357     Supplementary Table 1e). We have now identified TSSs for most of these short ORFs, indicating at

358     minimum they are transcribed. As described above, we noted that TSSs were found throughout the

359     genome in intergenic regions in addition to those identified upstream of the 190 annotated GRG ORFs

360     (including MGF 360-19R, Supplementary Table 2c). Our systematic, genome-wide approach identified

361     175 novel putative short ORFs. BLASTP (57) alignments showed that 13 were homologous to ORFs

362     predicted in other strains, including DP146L and pNG4 from BA71V . We validated the TSSs for these

363     candidates using 5'-RACE, which demonstrates the presence of these mRNAs and their associated TSSs

364     at both time-points (Figure 9d and e, respectively), compared to our CAGE-seq data (Figure 9f and g,

365     respectively).

366     **Putative single-SH2 domain protein encoding genes in MGF 100**

367     Our understanding of the ASFV genome is hampered by the large number of genes with unknown

368     functions. We approached this problem by searching for conserved domains of uncharacterised MGF

369     members *in silico*. MGF 100 genes form the smallest multigene family and include three short (100–

370     150 aa) paralogs located at both genome ends (right, R and left, L): 1R, 1L (MGF_100-2L or DP141L in

12

371   BA71V), and 3L (DP146L in BA71V). We predicted the two highly similar GRG ORFs I7L and I8L (51%

372   sequence identity) to belong to the MGF 100 family (Figure 10a), as designated in the Malawi LIL20/1

373   strain (58). Both I7L and I8L show similar overall transcript levels to the annotated MGF 100 members

374   -1L and 1R, though newly annotated MGF 100-3L (nORF_180573) was expressed at much higher levels.

375   I7L and I8L are both early genes like MGF 100-3L, while MGF 100-1L and 1R are expressed late and not

376   significantly changing, respectively (Supplementary Table 1e). Several lines of evidence suggest that

377   I7L and I8L play an role during infection. I7L and I8L are expressed early and at high levels, their

378   deletion along with L9R, L10L, and L11L ORFs reduces virulence in swine (59), and their loss is

379   associated with the adaptation of the GRG2007/1 strain to tissue culture infection (60). To gain insight

380   into the function of MGF family members including I7L and I8L, we generated computational

381   homology models of MGF 100-1L -1R, I7L and I8L using Phyre2 (61) (Figure 10b). The structures

382   selected by the algorithm for the modeling of MGF 100 proteins, included suppressor of cytokine

383   signalling proteins 1 and -2, and the PI3-kinase subunit alpha, all of which are characterized by Src

384   Homology 2 (SH2) domains (Figure 10b and Supplementary Table 2d). Canonical SH2 domains bind to

385   phosphorylated Tyrosine residues and are an integral part of signalling cascades involved in the

386   immune response (62). HHpred searches (63) predicted that indeed all MGF 100 members in BA71V

387   and GRG include SH2 domains (Figure 10c).

### The response of the porcine macrophage transcriptome to ASFV infection

389   In order to evaluate the impact of ASFV on the gene expression of the host cell, we analysed

390   transcriptomic changes of infected porcine macrophages using the CAGE-seq data from the control

391   (uninfected cells), 5 hpi, and 16 hpi. We annotated 9,384 macrophage-expressed protein-coding genes

392   with CAGE-defined TSSs (Supplementary Table 4). Although primary macrophages are known to vary

393   largely in their transcription profile, the CAGE-seq reads were highly similar between RNA samples

394   obtained from macrophages from two different animals in this study (Spearman's correlation

395   coefficients ≥ 0.77).

396   As TSSs are not well annotated for the swine genome, we annotated them *de novo* using our CAGE-

397   seq data with the RECLU pipeline. 37,159 peaks could be identified, out of which around half (18,575)

398   matched unique CAGE-derived peaks annotated in Robert et al. (64) i.e. they were located closer than

399   100 nt to the previously described peaks. Mapping CAGE-seq peaks to annotated swine protein-coding

400   genes led to identification of TSSs for 9,384 macrophage-expressed protein-coding genes

401   (Supplementary Table 4). The remaining 11,904 swine protein-coding genes did not have assigned

402   TSSs, and therefore their expression levels were not assessed. The majority of genes were assigned

403   with multiple TSSs, and these TSS-assigned genes, corresponded to many critical functional

13

404    macrophage markers, including genes encoding 56 cytokines and chemokines (including CXCL2, PPBP,

405    CXCL8 and CXCL5 as the most highly expressed), ten S100 calcium binding proteins (S100A12, S100A8,

406    and S100A9 in the top expressed genes), as well as interferon and TNF receptors (IFNGR1, IFNGR2,

407    IFNAR1, IFNAR2, IFNLR1, TNFRSF10B, TNFRSF1B, TNFRSF1A, etc.), and typical M1/M2 marker genes

408    such as TNF, ARG1, CCL24, and NOS2 (Supplementary Table 5

409    The 9,384 genes with annotated promoters were subjected to differential expression analysis using

410    DESeq2 to compare the 5 and 16 hour infected cell time points with control non-infected cells (c, 5

411    and 16) in a pairwise manner i.e. between each condition. Expression of only 25 host genes was

412    significantly deregulated between the control and 5 hpi, compared to 652 genes between 5 hpi and

413    16 hpi, and 1325 genes between mock-infected  and 16 hpi (at FDR of 0.05). Based on the pairwise

414    comparisons, we could distinguish major response profiles of the host genes. Late response genes,

415    whose expression was significantly deregulated both between the uninfected control and 16 hpi and

416    5 and 16 hpi, and early response genes, whose expression was significantly deregulated between the

417    control and 5 hpi, but not 5-16 hpi (Figure 11a). The latter category included only 20 genes, whereas

418    more than 500 genes showed the late differentially regulated response: 344 genes were up-regulated,

419    and 180 genes were down-regulated. The majority of the > 9000 genes analysed therefore were not

420    differentially regulated. Comparison of differences between expression levels in the different samples

421    indicate that macrophage differentially expressed transcription programs change mostly between 5

422    and 16 hpi (Figure 11b and c). The upregulated late response genes with highest expression levels

423    included several S100 calcium binding proteins. In contrast, expression of important cytokines

424    (including CCL24, CXCL2, CXCL5 and CXCL8) significantly decreased from 5 hpi to 16 hpi (Figure 11d).

425    To investigate the transcriptional response pathways and shed light on possible transcription factors

426    involved in the macrophage response to ASFV infection, we searched for DNA motifs enriched in

427    promoters of the four categories of deregulated genes in Figure 11a. Both late response promoter sets

428    were significantly enriched with motifs, some of which contained sub-motifs known to be recognised

429    by human transcription factors (Supplementary Figure 2). The highest-scored motif found in

430    promoters of upregulated genes contained a sub-motif recognised by a family of human interferon

431    regulatory factors (IRF9, IRF8 and IRF8, (Supplementary Figure 2a) that play essential roles in the anti-

432    viral response. Interestingly, both upregulated and downregulated promoters (Supplementary Figure

433    2b and c, respectively) were enriched with extended RELA/p65 motifs. p65 is a Rel-like domain-

434    containing subunit of the NF-kappa-B complex, regulated by I-kappa-B, whose analog is encoded by

435    ASFV. This pathway being a known target for ASFV in controlling host transcription (65–68).

14

436     To understand functional changes in the macrophage transcriptome, we also performed gene set

437     enrichment analysis using annotations of human homologs. The top enriched functional annotations

438     in the upregulated late response genes include glycoproteins and disulfide bonds, transmembrane

439     proteins, innate immunity, as well as positive regulation of inflammatory response (Figure 11e). In

440     contrast, sterol metabolism, rRNA processing, cytokines, TNF signalling pathway, inflammatory

441     response as well as innate immunity were the top enriched functional clusters among the

442     downregulated late response genes. Interestingly, the genes associated with innate immunity appear

443     overrepresented in both up- and downregulated gene subsets, yet cytokines are 8-fold enriched only

444     in the downregulated genes). The mRNA levels of genes of interest were additionally verified using

445     RT-PCR (Figure 11f).

446     Protein expression of selected genes.

447     In order to determine whether the regulation exerted by GRG on host transcription of

448     immunomodulatory genes could also translate to protein levels, we selected representative proteins

449     whose genes showed significant changes. ISG15 expression, part of the antiviral response genes of the

450     type I IFN stimulation pathway, was measured with Western blot (Figure 12a), with ASFV infection

451     being monitored via P30 levels (Figure 12b). Cytokines released from infected PAMS were quantified

452     using ELISA tests for pig cytokines, TNF-α, CXCL8 and CCL2 (Figure 12c, d and e, respectively). As shown

453     in Figure 12, the release/expression for all the tested proteins during GRG infection were similar or

454     decreased in comparison to the control uninfected cells at both 5 hpi and 16 hpi, while the production

455     of viral protein P30 increased, confirming an effective viral infection.

## Discussion [3,009 words]

457     In order to shed light on the gene expression determinants for ASF virulence, we focussed our analyses

458     on the similarities and differences in gene expression between a highly virulent Georgia 2007/1 isolate

459     and a nonvirulent, lab-adapted strain BA71V. Previous annotation identified 125 ASFV ORFs that are

460     conserved between all ASFV strain genomes irrespective of their virulence (16). They represent a 'core'

461     set of genes required for the virus to produce infectious progeny and include gene products like those

462     involved in virus genome replication, virion assembly, RNA transcription and modification. These

463     genes are located in the central region of the genome (Figure 1a). Besides such essential genes, about

464     one third are non-essential for replication, but have roles in evading host defence pathways. Some

465     genes are conserved between isolates, but not necessarily essential core genes, for example apoptosis

466     inhibitors: Bcl-2 family member A179L and IAP family member A224L (69). Other non-essential genes,

467     especially MGF members, vary in number between isolates. Our transcriptomic analysis captured TSS

15

468    signals from 119 genes both shared between the BA71V and GRG genomes, which also matched

469    expression patterns during early and late infection, according to CAGE-seq (Figure 3, Figure 4a-c).

470    Outliers include DP148R, which is unsurprising, given its promoter region is deleted in BA71V, and its

471    coding region is interrupted by a frame shift mutation, therefore functional protein expression

472    unlikely. DP148R is a non-essential, early-expressed virulence factor in the Benin 97/1 strain (49) –

473    consistent with our GRG data. Many additional GRG genes, lost from BA71V are MGFs, which are

474    mostly upregulated during early infection and located at the ends of the linear genome (Figure 1a).

475    MGFs have evolved on the virus genome by gene duplication, and do not share significant similarity

476    to other proteins, though some conserved domains, including ankyrin repeats, are present in some

477    MGF 360 and 505 family members (17,19).

478    Using advanced sequence searches and computational homology modelling we predict the members

479    of the MGF 100 family to encode SH2 domains, including I7L and I8L. Although SH2 domains are

480    primarily specific to eukaryotes, rare cases of horizontally transferred SH2 domains found in viruses,

481    are implicated in hijacking host cell pTyr signalling (70). A large family of 'super-binding' SH2 domains

482    were discovered in Legionella. Its members, including single SH2 domain-proteins are likely effector

483    proteins during infection (71). We also identified a further MGF 100 member in the GRG genome as

484    one of our nORFs, a partial 100-residue copy of DP146L (MGF 100-3L) (Supplementary Table 2c).

485    Unlike its annotated MGF 100-1L and MGF 100-1R cousins it was downregulated from 5 hpi to 16 hpi

486    (Supplementary Table 1e). Together with I7L and I8L, GRG encodes a total of 5 MGF 100 genes (Figure

487    10a). Interestingly, loss of MGF 100 members was observed during the process of adapting a virulent

488    Georgia strain to grow in cultured cell lines (60). Deletion of MGF 100-1R, from a virulent genotype II

489    Chinese strain (72) or of I8L from Georgia 2010 was shown not to reduce virulence of the virus in pigs

490    or reduce virus replication in porcine macrophages (73). However, simultaneous deletion of genes I7L,

491    I8L, I9L, I10L and I11L from a Chinese virulent isolate reduced virulence and surviving pigs were

492    protected against challenge (59). In summary, although deletion of some individual MGF 100 genes

493    does not lead to attenuation, deletion of I7L and I8L, in combination with I9L, I10L, and I11L did have

494    an impact.

495    The Georgia 2007/1 genome was recently re-sequenced which identified a small number of genome

496    changes affecting mapped ORFs and identified new ORFs (18). Adjacent to the covalently cross-linked

497    genome termini, the BA71V genome contains terminal inverted repeats of >2 kbp, in which two short

498    ORFs were identified (DP93R, DP86L). These were not included in previous GRG sequence annotations,

499    however our nORFs included a 55-residue homolog of DP96R, which was a late, but not highly

500    expressed gene. These are yet further examples of how transcriptomics aid in improving ASFV genome

16

501    annotation. Functional data is available for only a few of proteins coded by ORFs not conserved

502    between BA71V and GRG. This includes the p22 protein (KP177R), which is expressed on the cell

503    membrane during early infection, and also incorporated into the virus particle inner envelope. The

504    function of the KP177R-like GRG gene l10L has not been studied, but may provide an antigenically

505    divergent variant of P22, enabling evasion of the host immune response (19). We found KP177R was

506    highly expressed at 16 hpi, while l10L was also expressed late, but at much lower levels. Their function

507    is unknown, though the presence of an SH2 domain indicates possible roles in signalling pathways

508    (7,19,74).

509    MGF 110 members are among the highest expressed genes during early infection both in GRG (this

510    study), and in BA71V (10), suggesting high importance during infection, at least in porcine

511    macrophages and Vero cells, respectively. However, MGF 110 remains poorly characterised, and 13

512    orthologues were identified thus far, with numbers present varying between isolates (30). MGF 110

513    proteins possess cysteine-rich motifs, optimal for an oxidizing environment as found in the

514    endoplasmic reticulum (ER) lumen or outside the cell, and MGF 110-4L (XP124L) contains a KDEL signal

515    for retaining the protein in the ER (75). Since highly virulent isolates have few copies of these genes

516    (for example, only 5 in the Benin 97/1 genome), it was assumed they are not importance for virulence

517    in pigs (17), but their high expression warrant further investigation, which has recently begun in the

518    form of deletion mutants. For example, deletion of MGF 110-9L from a Chinese genotype II virulent

519    strain, reduced virulence (35), whereas deletion of MGF 110-1L from Georgia 2010 (76) did not

520    substantially affect virulence.

521    There is however, good evidence that MGF 360 and 505 carry out important roles in evading the host

522    type I interferon (IFN) response - the main host antiviral defence pathway (37). Evidence for the role

523    of MGF 360 and 505 genes in virulence was obtained from deletions in tissue-culture adapted and

524    field attenuated isolates, as well as targeted gene deletions This correlated with induction of the type

525    I interferon response, which itself is inhibited in macrophages infected with virulent ASFV isolates

526    (32,38,39). Deletions of these MGF 360, and 505 genes also correlated with an increased sensitivity of

527    ASFV replication, to pre-treatment of the macrophage cells with type I IFN (40). Thus, the MGF 360

528    and 505 genes have roles in inhibiting type I IFN induction and increasing sensitivity to type I IFN.

529    However, it remains unknown if these MGF 360 and MGF 505 genes act synergistically or if some have

530    a more important role than others type I IFN suppression. Our DESeq2 analysis did show that members

531    of both these families showed very similar patterns of early expression (Figure 2 and Figure 3),

532    conserved cEPM-containing promoters, and almost exclusive presence in clusters-1 (H-H), -4 (M-M),

17

533    and -5 (LM-LM) (Figure 6 and Figure 7), consistent with ASFV prioritising inhibition of the host immune

534    response during early infection.

535    An interesting pattern which emerged during our CAGE-seq analysis was the clear prevalence of ioTSSs

536    within ORFs, especially in MGFs (Figure 8 and Figure 9). However, it is not clear whether subsequent

537    in-frame truncation variants generate stable proteins, nor what their function could be. Perhaps even

538    more interesting was the discovery of 176 nORFs (including MGF 360-19R), with clear TSSs according

539    to CAGE-seq, highlighting the power of transcriptomics to better annotate sequenced genomes. We

540    were able to detect previously unannotated genes from other strains, and partial duplications of genes

541    already encoded in GRG (Supplementary Table 2).

542    The increase in transcription across the ASFV genome during late infection (10), appears ubiquitous.

543    At least 50 genes have previously been investigated in single gene expression studies using Northern

544    blot or primer extension (for review see references (10,77). Transcripts from over two thirds of these

545    genes were detected during late infection, and a quarter had transcripts detected during both early

546    and late infection. Therefore, clear evidence using several techniques now support this increase in

547    ASFV transcripts at late times post-infection. It is not entirely clear whether it is due to pervasive

548    transcription, high mRNA stability or a combination of factors. However, there is a correlated increase

549    in viral genome copies, potentially available as templates for pervasive transcription. The increase in

550    genome copies is more pronounced in BA71V compared to GRG, which likewise is reflected in the

551    increase in transcripts during late infection (Figure 4).

552    Our transcriptomic analysis of the porcine macrophage host revealed 522 genes whose expression

553    patterns significantly changed between 5 and 16 hrs post-infection (Figure 11a) and only 20 genes

554    were found to change between the control cells and those infected for 5 hpi. In aggregate, this reflects

555    a relatively slow host response to ASFV infection following expression of early ASFV genes. We

556    observed mild downregulation of some genes e.g. ACTB coding for ß-actin, eIF4A, and eIF4E

557    (Supplementary Table 5), resembling patterns previously shown by RT-qPCR (78). The macrophage

558    transcriptome mainly shuts down immunomodulation between 5 hpi to 16 hpi post-infection;

559    cytokines appeared highly expressed at 5 hpi, but downregulated from 5 hpi to 16 hpi. Of the 54

560    cytokine genes we detected, expression of thirteen was decreased: four interleukin genes (IL1A, IL1B,

561    IL19, IL27), four pro-inflammatory chemokines (CCL24, CXCL2, CXCL5, CXCL8), and tumor necrosis

562    factor (TNF) genes. Since inflammatory responses serve as the first line of host defense against viral

563    infections, viruses have developed ways to neutralise host pro-inflammatory pathways. ASFV encodes

564    a structural analog of IκB, A238L, which was proposed to act as a molecular off-switch for NFκB-

18

565    targeted pro-inflammatory cytokines (67). In our study, A238L is one of the most expressed ASFV

566    genes at 5 hpi, but significantly downregulated afterwards (Figure 2c). Accordingly, swine homologs

567    of human NFκB target genes were significantly over-represented (3.8 fold) among downregulated

568    macrophage genes (Fisher's exact p-value < 1e-5, based on human NFκB target genes from

569    https://www.bu.edu/nf-kb/gene-resources/target-genes/). Downregulated genes include

570    interleukins 1A, 1B, and 8, and 27 (IL1A, IL1B, CXCL8, IL27), TNF, as well as a target for common

571    nonsteroidal anti-inflammatory drugs, prostaglandin-endoperoxide synthase 2 (PTGS2 or COX-2)

572    (Supplementary Figure 2). Interestingly, promoters of both up- and downregulated genes contained a

573    motif with the sequence preferentially recognised by the human p65-NFκB complex (79). Expression

574    of TNF, a well-known marker gene for acute immune reaction and M1 polarisation, was recorded at a

575    high level in control samples  and at 5 hpi, but significantly dropped at 16 hpi. It has been already

576    shown that ASFV inhibits transcription of TNF and other proinflammatory cytokines (67). On the other

577    hand, the downregulation of TNF stands in contrast to previous results from ASFV-E75 strain-infected

578    macrophages in vitro, where TNF expression increased significantly after 6 hpi (80). Therefore, the

579    different time courses of TNF expression induced by the moderately virulent E75 and more virulent

580    Georgia strain may reflect different macrophage activation programs (81).

581    We investigated if the modulation of transcription we observed by CAGE-seq during GRG infection of

582    PAMS was also observed at the protein level. We analysed the secretion or expression of different

583    immunomediators (cytokines CCL2, CXCL8, TNF-α and interferon stimulated gene ISG15) at different

584    times following infection of PAMS. We confirmed that that the infection did not lead to an increase of

585    these mediators at either 5h or 16h infection. Secretion or expression of these proteins were similar

586    or slightly decreased in infected cells in comparison to control non-infected cells. The results indicated

587    that the control by virulent Georgia 2007/1 of host cell responses to infection we observed at the

588    transcription level can lead to a control also at the level of the protein production. Interestingly, CCL2

589    transcription was somewhat upregulated at late infection (Supplementary Table 5), whereas its

590    protein release to the supernatant was decreased (Figure 12e). ASFV has been shown to prioritize

591    expression of its encoded proteins by sequestering components of the host translation machinery to

592    viral factories (82). The levels or functions of host proteins may also be modulated by targeting for

593    post-translational modification or degradation (82–84). Therefore, in addition to control at the

594    transcriptional level ASFV may modulate the production of immunomodulatory host proteins at a later

595    step, as seems to occur for CCL2, a known chemoattractant for myeloid and lymphoid cells (85), that

596    could be an important target for regulation by ASFV.

19

597 Four S100 family members are among the host genes that are upregulated after 5 hpi (Figure 11b)

598 including S100A8, S100A11, S100A12, and S100A13. S100A8 and S100A12 are among the most highly

599 expressed genes on average throughout infection. S100 proteins are calcium-binding cytosolic

600 proteins that are released and serve as a danger signal, and stimulate inflammation (86). Once

601 released from the cell, S100A12 and S100A8 function as endogenous agonists to bind TLR4 and induce

602 apoptosis and autophagy in various cell types (86). S100A8 and S100A9 were also found in the RNA-

603 seq whole blood study as the top upregulated upon infection of the pigs with Georgia 2007/1, but not

604 of a low pathogenic ASFV isolate OURT 88/3 (43).

605 Previous studies described global swine transcriptome changes upon ASFV infection using short read

606 sequencing (Illumina): including the RNA-seq described above (43) and a microarray study of primary

607 swine macrophage cell cultures infected with the GRG strain, at six time points post-infection (42).

608 Although these varied in designs and selected methods, results of these works both give some

609 indication into the main host immune responses and ways how ASFV could evade them. The latter

610 microarray study indicated similar suppression of inflammatory response after 16 hpi as we observed

611 in this study, with expression of many cytokines down-regulated relative to non-infected macrophages

612 (42). More-recently, there have been several transcriptomic studies using classical RNA-seq of ASFV

613 infections from Chinese isolates (44–46). Fan et al (44) investigated the transcriptomic and proteomic

614 response within tissues of pigs following ASFV infection and death, though this was not directly

615 comparable to our own analysis in PAMs, due to their observations being of a far later infection stage

616 (post-mortem) than our 16 hrs time-point. The two most-comparable studies to ours were carried out

617 on a Chinese genotype II pathogenic strain during infection of PAMs. Ju et al. (45) investigated 6, 12

618 and 24 hpi, while Yang et al. (46) investigated 12, 24, and 36 hpi. However, comparison of the

619 overlapping time points of 12 hpi and 24 hpi did not yield similar host gene expression changes,

620 possibly due to variation among primary macrophages or due to the low MOI of 1 used in both studies.

621 In summary, these differences highlight that our understanding of the host-virus relationship during

622 ASFV infection is still not well understood, and further work is needed to understand why such

623 substantial variation in host gene expression can arise.

624 A further important note, is that all of the studies described above are using classical RNA-seq-based

625 methods, the nucleotide resolution of which, is not sufficient to investigate differential expression of

626 both the virus and host simultaneously. Investigating the viral transcriptome is especially difficult in a

627 compact genome like that of ASFV, where transcription read-through can undermine results from

628 classical RNA-sequencing techniques (10,87). A recent investigation into ASFV RNA transcripts using

629 long-read based Oxford Nanopore Technologies (ONT) – provides fascinating insight into their length

20

630   and read-through heterogeneity. This new method highlighted how misleading short read sequencing

631   with classical RNA-seq can be when quantifying ASFV gene expression, due to the abundance of

632   readthrough occurring in ASFV, generating transcripts covering multiple viral ORFs. This study did

633   however, unfortunately lack the read coverage for in-depth analysis of host transcripts alongside that

634   of viral transcripts (88,89).

635   Here we have demonstrated that CAGE-seq is an exceptionally powerful tool for quantifying relative

636   expression of viral genes across the ASFV genome, as well as making direct comparison between

637   strains for expression of shared genes, and further highlighting the importance of highly-expressed

638   but still functionally uncharacterised viral genes. CAGE-seq conveniently circumvents the issue in

639   compact viral genomes like those of ASFV and VACV, of transcripts reading through into downstream

640   genes which cannot be distinguished from classical short-read RNA-seq (10,43,90). Furthermore, it

641   enables us to effectively annotate genome-wide, the 5' ends of capped viral transcripts, and thus TSSs

642   of viral genes, and subsequently their temporal promoters. This 5' end resolution in ASFV is still not

643   achievable via ONT long read sequencing (88,89). We have now expanded on promoter motifs we

644   previously described (Figure 7), to identify 5 clusters of genes (Figure 6), with distinct patterns of

645   expression. Three of these clusters (-1: high to high levels, -4: mid to mid, and -5 low-mid to low-mid)

646   have slightly differing promoters, with a highly conserved core EPM. This is akin to the early gene

647   promoter of VACV (87) for VETF recognition and early gene transcription initiation (13,91,92). We have

648   found late genes can be categorised into two types that either increase from low to extremely high

649   expression levels (e. g. p72-encoding B646L) in cluster-2, or from low to medium expression levels in

650   cluster 3 (e. g VETF-encoding genes). The promoters of these genes show resemblance to the

651   eukaryotic TATA-box (93) or the BA71V LPM (10), respectively. Our analysis additionally shows the

652   potential for a variety of non-pTSSs: alternative ones used for different times in infection, ioTSSs which

653   could generate in-frame truncation variants of ORFs, sense or antisense transcripts relative to

654   annotated ORFs, and finally TSSs generating nORFs, which predominantly have no known homologs.

655   In summary, it is becoming increasingly clear that the transcriptomic landscape of ASFV and its host

656   during infection is far more complex than originally anticipated. Much of this raises further questions

657   about the basal mechanisms underlying ASFV transcription and how it is regulated over the infection

658   time course. Which subsets of initiation factors enable the RNAPs to recognise early and late

659   promoters? Does ASFV include intermediate genes, and what factors enables their expression? What

660   is the molecular basis of the pervasive transcription during late infection? The field of ASFV

661   transcription has been understudied and underappreciated and considering the severe threat that ASF

662   poses for the global food system and -food security, we now need to step up and focus our attention

21

663 and resources to study the fundamental biology of ASFV to develop effective antiviral drugs and
664 vaccines.

## Methodology [2871 words]

### GRG-Infection of Macrophages and RNA-extraction

667 Primary porcine alveolar macrophage cells were collected from two animals following approval by the
668 local Animal Welfare and Ethical Review Board at The Pirbright Institute. Cells were seeded in 6-well
669 plates (2x10$^6$ cells/well) with RPMI medium (with GlutaMAX), supplemented with 10% Pig serum and
670 100 IU/ml penicillin, 100 μg/ml streptomycin. They were infected as 2 replicate wells for 5 hpi or 16
671 hpi with a multiplicity of infection (MOI) of 5 of the ASFV Georgia 2007/1 strain, while uninfected cells
672 were seeded in parallel as a control (mock-infection). Total RNA was extracted according to
673 manufacturer's instructions for extraction with Trizol Lysis Reagent (Thermo Fisher Scientific and the
674 subsequent RNAs were resuspended in 50μl RNase-free water and DNase-treated (Turbo DNAfree kit,
675 Invitrogen). RNA quality was assessed via Bioanalyzer (Agilent 2100). 5 μg of each sample was ethanol
676 precipitated before sending to CAGE-seq (Kabushiki Kaisha DNAFORM, Japan). Samples were named
677 as follows: uninfected cells or 'mock' (C1-ctrl and C2-ctrl), at 5 hpi post-infection (samples G1-5h and
678 G2-5h), and at 16 hpi post-infection (G3-16h and G4-16h).

### CAGE-sequencing and Mapping to GRG and *Sus scrofa* Genomes

680 Library preparation and CAGE-sequencing of RNA samples was carried out by CAGE-seq (Kabushiki
681 Kaisha DNAFORM, Japan). Library preparation produced single-end indexed cDNA libraries for
682 sequencing: in brief, this included reverse transcription with random primers, oxidation and
683 biotinylation of 5' mRNA cap, followed by RNase ONE treatment removing RNA not protected in a
684 cDNA-RNA hybrid. Two rounds of cap-trapping using Streptavidin beads, washed away uncapped RNA-
685 cDNA hybrids. Next, RNase ONE and RNase H treatment degraded any remaining RNA, and cDNA
686 strands were subsequently released from the Streptavidin beads and quality assessed via Bioanalyzer.
687 Single strand index linker and 3' linker was ligated to released cDNA strands, and primer containing
688 Illumina Sequencer Priming site was used for second strand synthesis. Samples were sequenced using
689 the Illumina NextSeq 500 platform producing 76 bp reads. FastQC (94) analysis was carried out on all
690 FASTQ files at Kabushiki Kaisha DNAFORM and CAGE-seq reads showed consistent read quality across
691 their read-length, therefore, were mapped in their entirety to the GRG genome (FR682468.1) in our
692 work using Bowtie2 (95), and *Sus scrofa* (GCF_000003025.6) genome with HISAT2 (95,96) by Kabushiki
693 Kaisha DNAFORM.

22

## Transcription Start Site-mapping Across Viral GRG Genome

694

695 CAGE-seq mapped sample BAM files were converted to BigWig (BW) format with BEDtools (97)

696 genomecov, to produce per-strand BW files of 5' read ends. Stranded BW files were input for TSS-

697 prediction in RStudio (98) with Bioconductor (99) package CAGEfightR (100). Genomic feature

698 locations were imported as a TxDb object from FR682468.1 genome gene feature file (GFF3).

699 CAGEfightR was used to quantify the CAGE reads mapping at base pair resolution to the GRG genome

700 - at CAGE TSSs, separately for the 5 hpi and 16 hpi replicates. TSS values were normalized by tags-per-

701 million for each sample, pooled, and only TSSs supported by presence in both replicates were kept.

702 TSSs were assigned to clusters, if within 25 bp of one another, filtering out pooled, RPM-normalized

703 TSS counts below 25 bp for 5 hpi samples, or 50 bp for 16 hpi, and assigned a 'thick' value as the

704 highest TSS peak within that cluster. A higher cut-off for 16 hpi was used to minimise the extra noise

705 of pervasive transcription observed during late infection (10). TSS clusters were assigned to annotated

706 FR682468.1 ORFs using BEDtools intersect, if its highest point ('thick' region) was located within 500

707 bp upstream of an ORF, 'CDS' if within the ORF, 'NA' if no annotated ORF was within these regions.

708 Multiple TSSs located within 500 bp of ORFs were split into subsets: 'Primary' cluster subset contained

709 either the highest scoring CAGEfightR cluster or the highest scoring manually-annotated peak (when

710 manual ORF corrections necessary), and the highest peak coordinate was defined as the primary TSS

711 (pTSS) for an ORF. Further clusters associated with these ORFs were classified as 'non-primary', with

712 their highest peak as a non-primary TSS (npTSS). If the strongest TSS location was intra-ORF, without

713 any TSSs located upstream of the ORF, then the ORF was manually re-defined as starting from the next

714 ATG downstream.

## DESeq2 Differential Expression Analysis of GRG Genes

715

716 For analysing differential expression with the CAGE-seq dataset, a GFF was created with BEDtools

717 extending from the pTSS coordinate, 25 bp upstream and 75 bp downstream, however, in cases of

718 alternating pTSSs this region was defined as 25 bp upstream of the most upstream pTSS and 75 bp

719 downstream of the most downstream pTSS. HTSeq-count (101) was used to count reads mapping to

720 genomic regions described above for both the RNA- and CAGE-seq sample datasets. The raw read

721 counts were then used to analyse differential expression across these regions between the time-

722 points using DESeq2 (default normalisation described by Love et al. (47)) and those regions showing

723 changes with an adjusted p-value (padj) of <0.05 were considered significant. A caveat of this 'early'

724 or 'late' definition is that it is a binary definition of whether a gene is up- or downregulated between

725 conditions (time-points), relative to the background read depth of reads, which map to the genome

23

726    in question. Further analysis of ASFV genes used their characterised or predicted functions, from the

727    VOCS tool database (https://4virology.net/) (102,103) entries for the GRG genome.

### Quantification of viral genome copies at different time points of infection

729    Porcine lung macrophages were seeded and infected as described above. *Vero* cells were similarly

730    cultured in 6-well plates in DMEM medium supplemented with 10% Fetal calf serum, 100 IU/ml

731    penicillin and 100 µg/ml streptomycin, when semi-confluent they were infected with MOI 5 of Ba71V.

732    Immediately after infection (after 1h adsorption period, considered '0 hpi), or at 5 hpi, and 16 hpi, the

733    supernatant was removed and nucleic acids were extracted using the Qiamp viral RNA kit (Qiagen)

734    and quantified using a NanoDrop spectrophotometer (ThermoFisher Scientific). For quantification of

735    viral genome copy equivalents, 50 ng of each nucleic acid sample was used in qPCR with primers and

736    probe targeting the viral capsid gene B646L. As previously described (104), standard curve

737    quantification qPCR was carried out on a Mx3005P system (Agilent Technologies) using the primers

738    CTGCTCATGGTATCAATCTTATCGA    and    GATACCACAAGATC(AG)GCCGT    and    probe    5'-(6-

739    carboxyfluorescein    [FAM])-CCACGGGAGGAATACCAACCCAGTG-3'-(6-carboxytetramethylrhodamine

740    [TAMRA]).

### Analysis of mRNA levels by RT-PCR and quantitative real time PCR (qPCR)

742    RNA from GRG or Ba71V infected macrophages, or *Vero* cells respectively, or from uninfected cell

743    controls, was collected at the different time points post-infection with Trizol, as described above. RNA

744    was reverse transcribed (800 ng RNA per sample) using SuperScript III First-Strand Synthesis System

745    for RT-PCR and random hexamers (Invitrogen). For PCR, cDNAs were diluted 1:20 with nuclease free

746    water and 1 µl each sample was amplified in a total volume of 20 µl using Platinum™ Green Hot Start

747    PCR Master Mix (Invitrogen) and 200 nM of each primer. Annealing temperatures were tested for each

748    primer pair in gradient PCR to determine the one optimal for amplification.

749    Supplementary Table 7a shows the primers used for each gene target, the amplicon size, PCR reaction

750    conditions, and NCBI accession numbers for sequences used primer design. PCRs were then

751    performed with limited cycles of amplification to have a semi-quantitative comparison of transcript

752    abundance between infection timepoints (by not reaching the maximum product amplification

753    plateau). Amplification products were viewed using 1.5% agarose gel electrophoresis.

754    C315R transcript levels were assessed by qPCR, using housekeeping gene glyceraldehyde-3-phosphate

755    dehydrogenase (GAPDH) expression was used for normalisation. Primer details and the qPCR

756    amplification program are shown in

24

757 Supplementary Table 7b (GAPDH primers used for *Vero* cells were previously published by
758 Melchjorsen et al., 2009 (105)). Primers were used at 250 nM concentration with Brilliant III Ultra-Fast
759 SYBR® Green QPCR Master Mix (Agilent 600882), 1 µl cDNA in 20 µl (1:20) total reaction volumes, and
760 qPCRs carried out in Mx3005P system (Agilent Technologies). Similar amplification efficiencies (97-
761 102%) for all primers had been observed upon amplification of serially diluted cDNA samples, and the
762 relative expression at each timepoint of infection was calculated using the formula $2^{\Delta Ct}$ ($2^{Ct\_GAPDH-}$
763 $^{Ct\_C315R}$).

### Preparation of supernatant and cell lysis extracts for ELISA and Western blot detection of host proteins

766 Lung macrophage cultures from two donor outbred pigs (same cells used for CAGEseq) were prepared
767 in 6-well plates. Approximately $1.5x\ 10^6$ cells were seeded per well with 3 ml medium (RPMI with
768 penicillin/streptomycin and 10% pig serum) and incubated at 37 degrees 5% $CO_2$ overnight. Cultures
769 were washed once with culture medium to remove non-adherent cells and inoculated with MOI 5 of
770 ASFV-Georgia 2007/1 (or left uninfected as control) and centrifuged 1h at 600xg 26 degrees
771 (adsorption period). Supernatants from cell cultures were collected immediately after adsorption for
772 obtaining the 0 hpi timepoint and stored at -70 degrees until analysis. Adherent cells were washed
773 twice with cold DPBS (Sigma) and then lysed with 0.12 ml/well cold RIPA buffer (Thermo Scientific)
774 supplemented with protease inhibitors (Halt Protease Inhibitor Cocktail, Thermo Scientific). For 5h
775 and 16h timepoints, the inoculum was removed after adsorption, cells were washed twice in culture
776 medium and returned to the incubator with fresh 3 ml medium per well for the specified times of
777 infection. Supernatants and lysis volumes were collected similarly to the control. Supernatants were
778 analysed for the presence of CCL2 (Porcine CCL2/MCP-1 ELISA Kit, ES2RB Invitrogen), CXCL8
779 (Quantikine® ELISA, Porcine IL-8/CXCL8 Immunoassay, P8000 R&D) and TNF-α (Quantikine® ELISA,
780 Porcine TNF-α Immunoassay, PTA00 R&D) as recommended by the manufacturers. A volume of 25 µl
781 each lysate was analysed in Western Blot for expression of ISG15 (anti-ISG15 antibody ab233071,
782 Abcam; used at 1:1000 dilution), γ-Tubulin (anti-gamma Tubulin antibody ab11321, Abcam; used at
783 1:1000 dilution); and viral ASFV protein P30 (in-house mouse monoclonal antibody used at 1:500
784 dilution). Secondary antibodies used were Goat Anti-Rabbit IgG H&L (HRP) (ab205718, Abcam) and
785 Goat Anti-Mouse Immunoglobulins/HRP (P0447, Dako) both at 1:2000 dilution. Western blot
786 membranes were revealed using Pierce ECL Western Blotting Substrate (32106, Thermo Scientific).
787 Band densities were quantified using ImageJ (Rasband, W.S., ImageJ, U. S. National Institutes of
788 Health, Bethesda, Maryland, USA, https://imagej.nih.gov/ij/, 1997-2018).

25

### ASFV Promoter Motif Analysis

789 

790 DESeq2 results were used to categorise ASFV genes into two simple sub-classes: early; 87 genes

791 downregulated from early to late infection and late; the 78 upregulated from early to late infection.

792 These characterised gene pTSSs were then pooled with the nORF pTSSs, and sequences upstream and

793 downstream of the pTSS were extracted from the GRG genome in FASTA format using BEDtools.

794 Sequences 35 bp upstream of and including the pTSSs were analysed using MEME software

795 (http://meme-suite.org) (106), searching for 5 motifs with a maximum width of 20 nt and 27 nt,

796 respectively (other settings at default). The input for MEME motif searches included sequences

797 upstream of 134 early pTSSs (87 genes and 47 nORFs) for early promoter searching, while 234 late

798 pTSSs (78 genes and 156 nORFs) were used to search for late promoters. For analysis of conserved

799 motifs upstream of the five clusters described in Figure 6a-b, sequences were extracted in the same

800 manner as above, but grouped according to their cluster. MEME motif searches were carried out for

801 sequences in each cluster, searching for 3 motifs, 5-36 bp in length, with zero or one occurrence per

802 sequence ('zoops' mode).

### Identification of TSSs by rapid amplification of cDNA ends - 5'RACE

804 For 5'RACE of GRG genes DP146L, pNG4 and CP204L we designed the gene specific primers (GSP)

805 shown in

806 Supplementary Table 7c, and used the kit: "5′ RACE System for Rapid Amplification of cDNA Ends"

807 (Invitrogen), according to manufacturer instructions. Briefly, 150 ng RNA from either 5 hpi or 16 hpi

808 macrophages (one of the replicate RNA samples used for CAGE-seq) was used for cDNA synthesis with

809 GSP1 primers, followed by degradation of the mRNA template with RNase Mix, and column

810 purification of the cDNA. A homopolymeric tail was added to the cDNA 3'ends with Terminal

811 deoxynucleotidyl transferase, which allowed PCR amplification with an "Abridged Anchor Primer"

812 (AAP) from the 5'RACE kit and a nested GSP2 primer. A second PCR was performed over an aliquot of

813 the previous, with 5'RACE "Abridged Universal amplification Primer" (AUAP), and an additional nested

814 primer GSP3, except for pNG4 where GSP2 was re-used due to the small predicted size of the

815 amplicon. Platinum™ Green Hot Start PCR Master Mix (Invitrogen) was used for PCR and products

816 were run in 2% agarose gel electrophoresis (see

817 Supplementary Table 7c for expected sizes). Efficient recovery of cDNA from the purification column

818 requires a product of at least 200 bases and therefore, due to the small predicted size of pNG4

819 transcripts its GSP1 primer was extended at the 5' end with an irrelevant non-annealing sequence of

820 extra 50 nt in order to create a longer recoverable product.

26

### CAGE-seq Analysis for the *Sus scrofa* Genome

Analyses of TSS-mapping, gene expression and motif searching with CAGE-seq reads mapped to the *Sus scrofa* 11.1 genome were carried out by DNAFORM (Yokohama, Kanagawa, Japan). The 5' ends of CAGE-seq reads were utilised as input for the Reclu pipeline (107) with a cutoff of 0.1 RPM, and irreproducible discovery rate of 0.1. 37,159 total CAGE-seq peaks could be identified, of which around half (16,720) match unique CAGE peaks previously identified by Roberts et al. (64) (i.e. within 100 nt of any of them). TSSs for 9,384 protein-coding genes (out of 21,288) were annotated de novo from the CAGE-defined TSSs (Supplementary Table 4).

Protein-coding genes with annotated TSSs (9,384 out of 21,288) were then subjected to differential expression analysis. CAGE-seq reads were summed up over all TSSs assigned to a gene and compared between two time points using edgeR (108) at maximum false discovery rate of 0.05. The full list of host genes with annotated promoters together with their estimated expression levels is provided in Supplementary Table 5. Gene set enrichment analysis was performed with the DAVID 6.8 Bioinformatics Resources (109), using best BLASTP (110) human hits (from the UniProt (111) reference human proteome). The 9,331 genes with human homologs were used as a background, and functional annotations of the four major expression response groups (late/early up-/down-regulated genes) were clustered in DAVID 6.8 using medium classification stringency. MEME motif searches were conducted for promoters of four differentially regulated subsets of host genes, as defined in Figure 11a. Promoters sequences were extended 1000 bp upstream and 200 bp downstream of TSSs, searched with MEME (max. 10 motifs, max. 100 bp long, on a given strand only, zero or one site per sequence, E < 0.01), and then compared against known vertebrate DNA motifs with Tomtom (p-value < 0.01).

### Data Availability

Raw sequencing data are available on the Sequence Read Archive (SRA) database under BioProject: PRJNA739166. This also includes CAGE-seq data aligned to the ASFV-GRG (FR682468.1 *Sus scrofa* (GCF_000003025.6) genomes (see methods above) in BAM format.

### Acknowledgements

27

## References

855 1. Gogin A, Gerasimov V, Malogolovkin A, Kolbasov D. African swine fever in the North Caucasus

856 region and the Russian Federation in years 2007-2012. Virus Res. 2013 Apr 1;173(1):198–203.

857 2. Zhou X, Li N, Luo Y, Liu Y, Miao F, Chen T, et al. Emergence of African Swine Fever in China,

858 2018. Transbound Emerg Dis. 2018 Dec 1;65(6):1482–4.

859 3. Alonso C, Borca M, Dixon L, Revilla Y, Rodriguez F, Escribano JM, et al. ICTV Virus Taxonomy

860 Profile: Asfarviridae. J Gen Virol. 2018 May 1;99(5):613–4.

861 4. Koonin E V., Yutin N. Origin and evolution of eukaryotic large nucleo-cytoplasmic DNA viruses.

862 Intervirology. 2010;53(5):284–92.

863 5. Yutin N, Koonin E V. Hidden evolutionary complexity of Nucleo-Cytoplasmic Large DNA viruses

864 of eukaryotes. Virol J. 2012 Aug 14;9(1):161.

865 6. Broyles SS. Vaccinia virus transcription. J Gen Virol. 2003 Sep 1;84(9):2293–303.

866 7. Alejo A, Matamoros T, Guerra M, Andrés G. A proteomic atlas of the African swine fever virus

867 particle. J Virol. 2018 Sep 5;JVI.01293-18.

868 8. Salas ML, Kuznar J, Viñuela E. Polyadenylation, methylation, and capping of the RNA

869 synthesized in vitro by African swine fever virus. Virology. 1981 Sep 1;113(2):484–91.

870 9. Rodríguez JM, Salas ML. African swine fever virus transcription. Vol. 173, Virus Research.

871 Elsevier B.V.; 2013. p. 15–28.

872 10. Cackett G, Matelska D, Sýkora M, Portugal R, Malecki M, Bähler J, et al. The African Swine Fever

873 Virus Transcriptome. J Virol. 2020 Feb 19;94(9).

874 11. Iyer LM, Balaji S, Koonin E V., Aravind L. Evolutionary genomics of nucleo-cytoplasmic large

875 DNA viruses. Virus Res. 2006 Apr;117(1):156–84.

876 12. Yutin N, Wolf YI, Raoult D, Koonin E V. Eukaryotic large nucleo-cytoplasmic DNA viruses:

877 clusters of orthologous genes and reconstruction of viral genome evolution. Virol J. 2009 Dec

878 17;6(3):223.

879 13. Fischer U, Grimm C, Bartuli J, Böttcher B, Szalay A. Structural basis of the complete poxvirus

28

880       transcription initiation process. 2021 Apr 28;

881   14.   Rodríguez JM, Moreno LT, Alejo A, Lacasta A, Rodríguez F, Salas ML. Genome Sequence of

882         African Swine Fever Virus BA71, the Virulent Parental Strain of the Nonpathogenic and Tissue-

883         Culture Adapted BA71V. Munderloh UG, editor. PLoS One. 2015 Nov 30;10(11):e0142889.

884   15.   Yáñez RJ, Rodríguez JM, Nogal ML, Yuste L, Enríquez C, Rodriguez JF, et al. Analysis of the

885         complete nucleotide sequence of African swine fever virus. Virology. 1995 Apr 1;208(1):249–

886         78.

887   16.   Dixon LK, Chapman DAG, Netherton CL, Upton C. African swine fever virus replication and

888         genomics. Virus Res. 2013;173(1):3–14.

889   17.   Chapman DAG, Tcherepanov V, Upton C, Dixon LK. Comparison of the genome sequences of

890         non-pathogenic and pathogenic African swine fever virus isolates. J Gen Virol. 2008 Feb

891         1;89(2):397–408.

892   18.   Forth JH, Forth LF, King J, Groza O, Hübner A, Olesen AS, et al. A deep-sequencing workflow for

893         the fast and efficient generation of high-quality African swine fever virus whole-genome

894         sequences. Viruses. 2019;11(9).

895   19.   Chapman DAG, Darby AC, da Silva M, Upton C, Radford AD, Dixon LK. Genomic analysis of highly

896         virulent Georgia 2007/1 isolate of African swine fever virus. Emerg Infect Dis. 2011

897         Apr;17(4):599–605.

898   20.   Farlow J, Donduashvili M, Kokhreidze M, Kotorashvili A, Vepkhvadze NG, Kotaria N, et al. Intra-

899         epidemic genome variation in highly pathogenic African swine fever virus (ASFV) from the

900         country of Georgia. Virol J. 2018 Dec 14;15(1):190.

901   21.   Mazur-Panasiuk N, Woźniakowski G, Niemczuk K. The first complete genomic sequences of

902         African swine fever virus isolated in Poland. Sci Rep. 2019 Dec 1;9(1):3–5.

903   22.   Granberg F, Torresi C, Oggiano A, Malmberg M, Iscaro C, De Mia GM, et al. Complete genome

904         sequence of an African swine fever virus isolate from Sardinia, Italy. Genome Announc.

905         2016;4(6):1220–36.

906   23.   Wang Z, Jia L, Li J, Liu H, Liu D. Pan-Genomic Analysis of African Swine Fever Virus. Virologica

907         Sinica. Science Press; 2019. p. 1–4.

908   24.   Rowlands RJ, Michaud V, Heath L, Hutchings G, Oura C, Vosloo W, et al. African swine fever

29

909         virus isolate, Georgia, 2007. Emerg Infect Dis. 2008 Dec;14(12):1870–4.

910    25.    Zhao D, Liu R, Zhang X, Li F, Wang J, Zhang J, et al. Replication and virulence in pigs of the first
911         African swine fever virus isolated in China. Emerg Microbes Infect. 2019 Jan 1;8(1):438–47.

912    26.    Zani L, Forth JH, Forth L, Nurmoja I, Leidenberger S, Henke J, et al. Deletion at the 5'-end of
913         Estonian ASFV strains associated with an attenuated phenotype. Sci Reports 2018 81. 2018 Apr
914         25;8(1):1–11.

915    27.    Gallardo C, Nurmoja I, Soler A, Delicado V, Simón A, Martin E, et al. Evolution in Europe of
916         African swine fever genotype II viruses from highly to moderately virulent. Vet Microbiol. 2018
917         Jun 1;219:70–9.

918    28.    Pershin A, Shevchenko I, Igolkin A, Zhukov I, Mazloum A, Aronova E, et al. A Long-Term Study
919         of the Biological Properties of ASF Virus Isolates Originating from Various Regions of the
920         Russian Federation in 2013–2018. Vet Sci 2019, Vol 6, Page 99. 2019 Dec 6;6(4):99.

921    29.    Sun E, Zhang Z, Wang Z, He X, Zhang X, Wang L, et al. Emergence and prevalence of naturally
922         occurring lower virulent African swine fever viruses in domestic pigs in China in 2020. Sci China
923         Life Sci. 2021 May 1;64(5):752–65.

924    30.    Imbery J, Upton C. Organization of the multigene families of African Swine Fever Virus. Fine
925         Focus. 2017;3(2):155–70.

926    31.    Netherton CL, Connell S, Benfield CTO, Dixon LK. The Genetics of Life and Death: Virus-Host
927         Interactions Underpinning Resistance to African Swine Fever, a Viral Hemorrhagic Disease.
928         Front Genet. 2019 May 3;10(MAY):402.

929    32.    Reis AL, Abrams CC, Goatley LC, Netherton C, Chapman DG, Sanchez-Cordon P, et al. Deletion
930         of African swine fever virus interferon inhibitors from the genome of a virulent isolate reduces
931         virulence in domestic pigs and induces a protective response. Vaccine. 2016 Sep
932         7;34(39):4698–705.

933    33.    O'Donnell V, Risatti GR, Holinka LG, Krug PW, Carlson J, Velazquez-Salinas L, et al. Simultaneous
934         Deletion of the 9GL and UK Genes from the African Swine Fever Virus Georgia 2007 Isolate
935         Offers Increased Safety and Protection against Homologous Challenge. J Virol. 2017 Jan;91(1).

936    34.    Li D, Zhang J, Yang W, Li P, Ru Y, Kang W, et al. African swine fever virus protein MGF-505-7R
937         promotes virulence and pathogenesis by inhibiting JAK1- and JAK2-mediated signaling. J Biol
938         Chem. 2021 Nov;297(5):101190.

30

939   35.   Li D, Liu YY, Qi X, Wen Y, Li P, Ma Z, et al. African Swine Fever Virus MGF-110-9L-deficient
940        Mutant Has Attenuated Virulence in Pigs. Virol Sin. 2021 Apr 1;36(2):187–95.

941   36.   Keßler C, Forth JH, Keil GM, Mettenleiter TC, Blome S, Karger A. The intracellular proteome of
942        African swine fever virus. Sci Rep. 2018 Oct 2;8(1):14714.

943   37.   Randall RE, Goodbourn S. Interferons and viruses: An interplay between induction, signalling,
944        antiviral responses and virus countermeasures. Vol. 89, Journal of General Virology.
945        Microbiology Society; 2008. p. 1–47.

946   38.   Afonso RCL, Piccone ME, Zaffuto KM, Neilan J, Kutish GF, Lu Z, et al. African swine fever virus
947        multigene family 360 and 530 genes affect host interferon response. J Virol. 2004;78:1858–64.

948   39.   Neilan JG, Zsak L, Lu Z, Kutish GF, Afonso CL, Rock DL. Novel Swine Virulence Determinant in
949        the Left Variable Region of the African Swine Fever Virus Genome. J Virol. 2002 Apr
950        1;76(7):3095–104.

951   40.   Golding JP, Goatley L, Goodbourn S, Dixon LK, Taylor G, Netherton CL. Sensitivity of African
952        swine fever virus to type I interferon is linked to genes within multigene families 360 and 505.
953        Virology. 2016 Jun 1;493:154–61.

954   41.   Mosser DM, Edwards JP. Exploring the full spectrum of macrophage activation. Vol. 8, Nature
955        Reviews Immunology. NIH Public Access; 2008. p. 958–69.

956   42.   Zhu JJ, Ramanathan P, Bishop EA, O'Donnell V, Gladue DP, Borca M V. Mechanisms of African
957        swine fever virus pathogenesis and immune evasion inferred from gene expression changes in
958        infected swine macrophages. PLoS One. 2019;14(11).

959   43.   Jaing C, Rowland RRR, Allen JE, Certoma A, Thissen JB, Bingham J, et al. Gene expression
960        analysis of whole blood RNA from pigs infected with low and high pathogenic African swine
961        fever viruses. Sci Rep. 2017 Dec 31;7(1):10115.

962   44.   Fan W, Cao Y, Jiao P, Yu P, Zhang H, Chen T, et al. Synergistic effect of the responses of different
963        tissues against African swine fever virus. Transbound Emerg Dis. 2021;

964   45.   Ju X, Li F, Li J, Wu C, Xiang G, Zhao X, et al. Genome-wide transcriptomic analysis of highly
965        virulent African swine fever virus infection reveals complex and unique virus host interaction.
966        Vet Microbiol. 2021 Oct 1;261:109211.

967   46.   Yang B, Shen C, Zhang D, Zhang T, Shi X, Yang J, et al. Mechanism of interaction between virus

31

968    and host is inferred from the changes of gene expression in macrophages infected with African

969    swine fever virus CN/GS/2018 strain. Virol J. 2021 Dec 1;18(1).

970    47.    Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq

971    data with DESeq2. Genome Biol. 2014 Dec 5;15(12):550.

972    48.    Hammond JM, Kerr SM, Smith GL, Dixon LK. An African swine fever virus gene with homology

973    to DNA ligases. Nucleic Acids Res. 1992 Jun 11;20(11):2667–71.

974    49.    Reis AL, Goatley LC, Jabbar T, Sanchez-Cordon PJ, Netherton CL, Chapman DAG, et al. Deletion

975    of the African Swine Fever Virus Gene DP148R Does Not Reduce Virus Replication in Culture

976    but Reduces Virus Virulence in Pigs and Induces High Levels of Protection against Challenge. J

977    Virol. 2017 Dec 15;91(24).

978    50.    Yang Z, Martens CA, Bruno DP, Porcella SF, Moss B. Pervasive initiation and 3'-end formation

979    of poxvirus postreplicative RNAs. J Biol Chem. 2012 Sep 7;287(37):31050–60.

980    51.    Frouco G, Freitas FB, Coelho J, Leitão A, Martins C, Ferreira F. DNA-Binding Properties of African

981    Swine Fever Virus pA104R, a Histone-Like Protein Involved in Viral Replication and

982    Transcription. J Virol. 2017;91(12).

983    52.    Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A, McPherson A, et al. A survey

984    of best practices for RNA-seq data analysis. Genome Biol. 2016 Jan 26;17:13.

985    53.    García-Escudero R, Viñuela E. Structure of African Swine Fever Virus Late Promoters:

986    Requirement of a TATA Sequence at the Initiation Region. J Virol. 2000 Sep 1;74(17):8176–82.

987    54.    Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, et al. MEME SUITE: tools for motif

988    discovery and searching. Nucleic Acids Res. 2009 Jul 1;37(Web Server):W202–8.

989    55.    Gupta S, Stamatoyannopoulos JA, Bailey TL, Noble W. Quantifying similarity between motifs.

990    Genome Biol. 2007 Feb 26;8(2):R24.

991    56.    Rodríguez JM, Salas ML, Viñuela E. Intermediate class of mRNAs in African swine fever virus. J

992    Virol. 1996 Dec;70(12):8584–9.

993    57.    Kim DE, Chivian D, Baker D. Protein structure prediction and analysis using the Robetta server.

994    Nucleic Acids Res. 2004;32(WEB SERVER ISS.):W526–31.

995    58.    Vydelingum S, Baylis SA, Bristow C, Smith GL, Dixon LK. Duplicated genes within the variable

996    right end of the genome of a pathogenic isolate of African swine fever virus. J Gen Virol. 1993

32

997        Oct 1;74(10):2125–30.

998    59.   Zhang J, Zhang Y, Chen T, Yang JJ, Yue H, Wang L, et al. Deletion of the L7L-L11L Genes
999          Attenuates ASFV and Induces Protection against Homologous Challenge. Viruses. 2021 Feb
1000         1;13(2).

1001   60.   Krug PW, Holinka LG, O'Donnell V, Reese B, Sanford B, Fernandez-Sainz I, et al. The Progressive
1002         Adaptation of a Georgian Isolate of African Swine Fever Virus to Vero Cells Leads to a Gradual
1003         Attenuation of Virulence in Swine Corresponding to Major Modifications of the Viral Genome.
1004         J Virol. 2015 Feb 15;89(4):2324–32.

1005   61.   Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJE. The Phyre2 web portal for protein
1006         modeling, prediction and analysis. Nat Protoc. 2015 May 7;10(6):845–58.

1007   62.   Liu H, Li L, Voss C, Wang F, Liu J, Li SSC. A Comprehensive Immunoreceptor Phosphotyrosine-
1008         based Signaling Network Revealed by Reciprocal Protein–Peptide Array Screening. Mol Cell
1009         Proteomics. 2015 Jul 1;14(7):1846–58.

1010   63.   Gabler F, Nam SZ, Till S, Mirdita M, Steinegger M, Söding J, et al. Protein Sequence Analysis
1011         Using the MPI Bioinformatics Toolkit. Curr Protoc Bioinforma. 2020 Dec 1;72(1).

1012   64.   Robert C, Kapetanovic R, Beraldi D, Watson M, Archibald AL, Hume DA. Identification and
1013         annotation of conserved promoters and macrophage-expressed genes in the pig genome. BMC
1014         Genomics. 2015 Nov 18;16(1).

1015   65.   Ganchi PA, Sun SC, Greene WC, Ballard DW. A novel NF-kappa B complex containing p65
1016         homodimers: implications for transcriptional control at the level of subunit dimerization. Mol
1017         Cell Biol. 1993 Dec;13(12):7826–35.

1018   66.   Dixon LK, Abrams CC, Bowick G, Goatley LC, Kay-Jackson PC, Chapman D, et al. African swine
1019         fever virus proteins involved in evading host defence systems. In: Veterinary Immunology and
1020         Immunopathology. 2004. p. 117–34.

1021   67.   Powell PP, Dixon LK, Parkhouse RM. An IkappaB homolog encoded by African swine fever virus
1022         provides a novel mechanism for downregulation of proinflammatory cytokine responses in
1023         host macrophages. J Virol. 1996;70(12):8527–33.

1024   68.   Granja AG, Sánchez EG, Sabina P, Fresno M, Revilla Y. African swine fever virus blocks the host
1025         cell antiviral inflammatory response through a direct inhibition of PKC-theta-mediated p300
1026         transactivation. J Virol. 2009 Jan 15;83(2):969–80.

33

1027    69.    Nogal ML, González de Buitrago G, Rodríguez C, Cubelos B, Carrascosa AL, Salas ML, et al.
1028           African swine fever virus IAP homologue inhibits caspase activation and promotes cell survival
1029           in mammalian cells. J Virol. 2001 Mar 15;75(6):2535–43.

1030    70.    Takeya T, Hanafusa H. DNA sequence of the viral and cellular src gene of chickens. II.
1031           Comparison of the src genes of two strains of avian sarcoma virus and of the cellular homolog.
1032           J Virol. 1982;44(1):12–8.

1033    71.    Kaneko T, Stogios PJ, Ruan X, Voss C, Evdokimova E, Skarina T, et al. Identification and
1034           characterization of a large family of superbinding bacterial SH2 domains. Nat Commun. 2018
1035           Dec 1;9(1).

1036    72.    Liu Y, Li Y, Xie Z, Ao Q, Di D, Yu W, et al. Development and in vivo evaluation of MGF100-1R
1037           deletion mutant in an African swine fever virus Chinese strain. Vet Microbiol. 2021 Oct 1;261.

1038    73.    Vuono E, Ramirez-Medina E, Pruitt S, Rai A, Silva E, Espinoza N, et al. Evaluation in swine of a
1039           recombinant georgia 2010 african swine fever virus lacking the i8l gene. Viruses. 2021 Jan
1040           1;13(1).

1041    74.    Camacho A, ViÑuela E. Protein p22 of African swine fever virus: An early structural protein that
1042           is incorporated into the membrane of infected cells. Virology. 1991 Mar 1;181(1):251–7.

1043    75.    Netherton C, Rouiller I, Wileman T. The subcellular distribution of multigene family 110
1044           proteins of African swine fever virus is determined by differences in C-terminal KDEL
1045           endoplasmic reticulum retention motifs. J Virol. 2004 Apr 1;78(7):3710–21.

1046    76.    Ramirez-Medina E, Vuono E, Pruitt S, Rai A, Silva E, Espinoza N, et al. Development and In Vivo
1047           Evaluation of a MGF110-1L Deletion Mutant in African Swine Fever Strain Georgia. Viruses.
1048           2021 Feb 1;13(2).

1049    77.    Cackett G, Sýkora M, Werner F. Transcriptome view of a killer: African swine fever virus. Vol.
1050           48, Biochemical Society Transactions. Portland Press Ltd; 2020. p. 1569–81.

1051    78.    Quintas A, Pérez-Núñez D, Sánchez EG, Nogal ML, Hentze MW, Castelló A, et al.
1052           Characterization of the African Swine Fever Virus Decapping Enzyme during Infection. Jung JU,
1053           editor. J Virol. 2017 Dec 15;91(24):e00990-17.

1054    79.    Kunsch C, Ruben SM, Rosen CA. Selection of optimal kappa B/Rel DNA-binding motifs:
1055           interaction of both subunits of NF-kappa B with DNA is required for transcriptional activation.
1056           Mol Cell Biol. 1992 Oct;12(10):4412–21.

34

1057    80.    Gómez del Moral M, Ortuño E, Fernández-Zapatero P, Alonso F, Alonso C, Ezquerra A, et al.
1058           African Swine Fever Virus Infection Induces Tumor Necrosis Factor Alpha Production:
1059           Implications in Pathogenesis. J Virol. 1999 Mar 1;73(3):2173–80.

1060    81.    Roy S, Schmeier S, Arner E, Alam T, Parihar SP, Ozturk M, et al. Redefining the transcriptional
1061           regulatory dynamics of classically and alternatively activated macrophages by deepCAGE
1062           transcriptomics. Nucleic Acids Res. 2015;43(14):6969–82.

1063    82.    Castelló A, Quintas A, Sánchez EG, Sabina P, Nogal M, Carrasco L, et al. Regulation of host
1064           translational machinery by African swine fever virus. PLoS Pathog. 2009 Aug;5(8):e1000562.

1065    83.    Sánchez EG, Quintas A, Nogal M, Castelló A, Revilla Y. African swine fever virus controls the
1066           host transcription and cellular machinery of protein synthesis. Virus Res. 2013 Apr 1;173(1):58–
1067           75.

1068    84.    Barrado-Gil L, Del Puerto A, Muñoz-Moreno R, Galindo I, Cuesta-Geijo MA, Urquiza J, et al.
1069           African Swine Fever Virus Ubiquitin-Conjugating Enzyme Interacts With Host Translation
1070           Machinery to Regulate the Host Protein Synthesis. Front Microbiol. 2020 Dec 15;11.

1071    85.    Gschwandtner M, Derler R, Midwood KS. More Than Just Attractive: How CCL2 Influences
1072           Myeloid Cell Behavior Beyond Chemotaxis. Front Immunol. 2019 Dec 13;0:2759.

1073    86.    Xia C, Braunstein Z, Toomey AC, Zhong J, Rao X. S100 proteins as an important regulator of
1074           macrophage inflammation. Vol. 8, Frontiers in Immunology. Frontiers Media S.A.; 2018.

1075    87.    Yang Z, Bruno DP, Martens CA, Porcella SF, Moss B. Simultaneous high-resolution analysis of
1076           vaccinia virus and host cell transcriptomes by deep RNA sequencing. Proc Natl Acad Sci U S A.
1077           2010 Jun 22;107(25):11513–8.

1078    88.    Olasz F, Tombácz D, Torma G, Csabai Z, Moldován N, Dörmő Á, et al. Short and Long-Read
1079           Sequencing Survey of the Dynamic Transcriptomes of African Swine Fever Virus and the Host
1080           Cells. Front Genet. 2020 Jul 28;11:2020.02.27.967695.

1081    89.    Torma G, Tombácz D, Csabai Z, Moldován N, Mészáros I, Zádori Z, et al. Combined short and
1082           long-read sequencing reveals a complex transcriptomic architecture of African swine fever
1083           virus. Viruses. 2021 Apr 1;13(4).

1084    90.    Yang Z, Bruno DP, Martens CA, Porcella SF, Moss B. Genome-Wide Analysis of the 5' and 3'
1085           Ends of Vaccinia Virus Early mRNAs Delineates Regulatory Sequences of Annotated and
1086           Anomalous Transcripts. J Virol. 2011 Jun;85(12):5897–909.

35

1087    91.    Gershon PD, Moss B. Early transcription factor subunits are encoded by vaccinia virus late
1088         genes. Proc Natl Acad Sci U S A. 1990 Jun 1;87(11):4401–5.

1089    92.    Li J, Broyles SS. Recruitment of vaccinia virus RNA polymerase to an early gene promoter by
1090         the viral early transcription factor. J Biol Chem. 1993;268(4):2773–80.

1091    93.    Patikoglou GA, Kim JL, Sun L, Yang SH, Kodadek T, Burley SK. TATA element recognition by the
1092         TATA box-binding protein has been conserved throughout evolution. Genes Dev. 1999 Dec
1093         15;13(24):3217–30.

1094    94.    Andrews S. FastQC A Quality Control tool for High Throughput Sequence Data. Babraham,
1095         England: Babraham Bioinformatics;

1096    95.    Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods. 2012 Apr
1097         4;9(4):357–9.

1098    96.    Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory requirements.
1099         Nat Methods. 2015 Apr 9;12(4):357–60.

1100    97.    Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features.
1101         Bioinformatics. 2010 Mar 15;26(6):841–2.

1102    98.    RStudio Team. RStudio: Integrated Development for R. Boston, MA: RStudioe, Inc; 2016.

1103    99.    Huber W, Carey VJ, Gentleman R, Anders S, Carlson M, Carvalho BS, et al. Orchestrating high-
1104         throughput genomic analysis with Bioconductor. Nat Methods. 2015 Feb 1;12(2):115–21.

1105    100.    Thodberg M, Thieffry A, Vitting-Seerup K, Andersson R, Sandelin A. CAGEfightR: Analysis of 5'-
1106         end data using R/Bioconductor. BMC Bioinformatics. 2019 Oct 4;20(1):487.

1107    101.    Anders S, Pyl PT, Huber W. HTSeq--a Python framework to work with high-throughput
1108         sequencing data. Bioinformatics. 2015 Jan 15;31(2):166–9.

1109    102.    Upton C, Slack S, Hunter AL, Ehlers A, Roper RL, Rock DL. Poxvirus orthologous clusters: toward
1110         defining the minimum essential poxvirus genome. J Virol. 2003 Jul 1;77(13):7590–600.

1111    103.    Tu SL, Upton C. Bioinformatics for Analysis of Poxvirus Genomes. In: Methods in Molecular
1112         Biology. Humana Press Inc.; 2019. p. 29–62.

1113    104.    DP K, SM R, GH H, SS G, PJ W, LK D, et al. Development of a TaqMan PCR assay with internal
1114         amplification control for the detection of African swine fever virus. J Virol Methods. 2003
1115         Jan;107(1):53–61.

36

1116    105.    Melchjorsen J, Kristiansen H, Christiansen R, Rintahaka J, Matikainen S, Paludan SR, et al.
1117            Differential regulation of the OASL and OAS1 genes in response to viral infections. J Interf
1118            Cytokine Res. 2009 Apr 1;29(4):199–207.

1119    106.    Bailey TL, Elkan C. Fitting a mixture model by expectation maximization to discover motifs in
1120            biopolymers. Proceedings Int Conf Intell Syst Mol Biol. 1994;2:28–36.

1121    107.    Ohmiya H, Vitezic M, Frith MC, Itoh M, Carninci P, Forrest ARR, et al. RECLU: A pipeline to
1122            discover reproducible transcriptional start sites and their alternative regulation using capped
1123            analysis of gene expression (CAGE). BMC Genomics. 2014 Apr 25;15(1):269.

1124    108.    Robinson MD, McCarthy DJ, Smyth GK. edgeR: A Bioconductor package for differential
1125            expression analysis of digital gene expression data. Bioinformatics. 2009 Nov 11;26(1):139–40.

1126    109.    Huang DW, Sherman BT, Lempicki RA. Systematic and integrative analysis of large gene lists
1127            using DAVID bioinformatics resources. Nat Protoc. 2009;4(1):44–57.

1128    110.    Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. J Mol
1129            Biol. 1990 Oct 5;215(3):403–10.

1130    111.    Bateman A. UniProt: A worldwide hub of protein knowledge. Nucleic Acids Res. 2019 Jan
1131            8;47(D1):D506–15.

1132    112.    Ramírez F, Dündar F, Diehl S, Grüning BA, Manke T. deepTools: a flexible platform for exploring
1133            deep-sequencing data. Nucleic Acids Res. 2014 Jul;42(Web Server issue):W187-91.

1134    113.    Crooks GE, Hon G, Chandonia JM, Brenner SE. WebLogo: A sequence logo generator. Genome
1135            Res. 2004 May 12;14(6):1188–90.

1136    114.    Grant CE, Bailey TL, Noble WS. FIMO: scanning for occurrences of a given motif. Bioinformatics.
1137            2011 Apr 1;27(7):1017–8.

1138    115.    McCarthy DJ, Chen Y, Smyth GK. Differential expression analysis of multifactor RNA-Seq
1139            experiments with respect to biological variation. Nucleic Acids Res. 2012 May 1;40(10):4288–
1140            97.

1141    116.    Benjamini Y, Hochberg Y. Controlling the False Discovery Rate: A Practical and Powerful
1142            Approach to Multiple Testing. J R Stat Soc Ser B. 1995 Jan;57(1):289–300.

1143    117.    Yates AD, Achuthan P, Akanni W, Allen J, Allen J, Alvarez-Jarreta J, et al. Ensembl 2020. Nucleic
1144            Acids Res. 2020 Jan 1;48(D1):D682–8.

37

1145 118. Wheeler DL, Church DM, Federhen S, Lash AE, Madden TL, Pontius JU, et al. Database resources

1146 of the National Center for Biotechnology. Nucleic Acids Res. 2003 Jan 1;31(1):28–33.

1147

## Figures

1149

1150 Figure 1. Functional genome annotation of ASFV GRG. (a) Comparison between the genomes of BA71V

1151 and GRG, generated with Circos (http://circos.ca/). Blue lines represent sequence conservation (Blast

1152 E-values per 100 nt). The Inner ring represents genes defined as MGF members (purple), and all others

1153 (grey). The outer ring shows annotated genes which we have defined as early or late according to

1154 downregulation or upregulation between 5 hpi and 16 hpi from DESeq2 analysis. (b) 189 GRG

1155 annotated ORFs are represented as arrows and coloured according to strand. CAGE-seq peaks across

1156 the GRG genome at 5 hpi (c) and 16 hpi (d), normalized coverage reads per million mapped reads

1157 (RPM) of 5' ends of CAGE-seq reads. The coverage was capped at 20000 RPM for visualisation, though

1158 multiple peaks exceeded this. DeepTools (112), was used to convert bam files to bigwig format and

1159 imported into Rstudio for visual representation via packages ggplot, ggbio, rtracklayer, and gggenes

1160 was used to generate the ORF map in (b).

1161

1162 Figure 2. Summary of GRG gene expression (a) Expression profiles for 164 genes for which we

1163 annotated pTSSs from CAGE-seq and which showed significant differential expression. Log2 fold

1164 change and basemean expression values were from DESeq2 analysis of raw counts (see methods).

1165 Genes are coloured according to their log2 fold change in expression as red (positive: upregulated

1166 from 5 hpi to 16 hpi) or blue (negative: downregulated). MGFs are emphasised with a black outline to

1167 highlight their overrepresentation in the group of downregulated genes. (b) Expression profiles for 41

1168 genes (excluding nORFs) only detected as being expressed in GRG and not BA71V, format as in (a). (c)

1169 Expression (RPM) of 20 highest-expressed genes at 5 hpi, error bars represent standard deviation

1170 between replicates. (d) Expression (RPM) of 20 highest-expressed genes at 16 hpi pi, error bars are

1171 the standard deviation between replicates.

1172

1173 Figure 3. Comparison of gene expression profiles for genes shared between GRG and BA71V. Scatter

1174 plots of mean RPM across replicates for shared genes at 5 hpi (a) and 16 hpi (b), coloured according

1175 to whether genes show significant downregulation (blue), or upregulation (red) according to DESeq2

38

1176    analysis in GRG. In both (b) and (c) genes with RPM values above 40000 RPM in either strain are

1177    labelled. (c) Comparison of log2 fold change in expression values of genes in GRG and BA71V, in blue

1178    are downregulated (early) genes in both strains, red are upregulated (late) genes in both strains, while

1179    the genes which disagree in their differential expression patterns between strains are in black. R

1180    represents the Pearson Correlation coefficient for each individual plot in (a), (b), and (c). Due to

1181    inconsistencies in their genome annotations, two genes were omitted from the BA71V-GRG

1182    transcriptome comparisons in Figures 2b and 3a-d: EP296R in GRG known as E296R in BA71V, and

1183    C122R (GRG) is the old nomenclature for C105R (BA71V), which are now correctly named in

1184    Supplementary Table 1e and Figure 2a. Both genes showed the same early expression patterns in

1185    BA71V (10) and GRG (Supplementary Table 1e) so would strengthen the patterns observed.

1186

1187    Figure 4. Increase in virus genome copy number mRNA levels during late infection. (a) The 'log2

1188    change' represents log2 of the ratio of CAGE-seq reads (normalised per million mapped reads) at 16

1189    hpi vs. 5 hpi per nucleotide across the genome. Alignment comparisons and calculations were done

1190    with deepTools (112). (b) Replicate means of CAGE-seq reads mapped to either the BA71V (green) or

1191    GRG (purple) genomes throughout infection. (c) Fold change in CAGE-seq reads during infection,

1192    calculated via mean value across 2 replicates, but with the assumption number of reads at 0 hpi is 0,

1193    therefore dividing by values from 5 hpi. (d) Change in genome copies from DNA qPCR of B646L gene,

1194    dividing by value at 0 hpi to represent '1 genome copy per infected cell'. (e) Fold change in genome

1195    copies present at 0 hpi , 5 hpi and 16 hpi from qPCR in (d). (d) calculated as for (c), but with actual

1196    vales for 0 hpi.

1197

1198    Figure 5. RT-PCR results of genes for comparison to CAGE-seq data from (a) MGF 505-7R, (b) NP419L,

1199    (c) D345L, (d) MGF 360-12L, (e) MGF 505-9R, and (f) qRT-PCR results of C315R (ASFV-TFIIB). (NT = no

1200    template control). For each panel at the top is a diagrammatic representation of each gene's TSSs

1201    (bent arrow, including both pTSS and ioTSSs), annotated ORF (red arrow), and arrow pairs in cyan or

1202    yellow represent the primers used for PCR (see methods for primer sequences). Beneath each PCR

1203    results are bar charts representing the CAGE-seq results as either normalised (mean RPM) or raw

1204    (mean read counts) data, error bars show the range of values from each replicate.

1205

39

1206    Figure 6. Comparison of the raw read counts for genes shared between BA71V and GRG. (a) clustered

1207    heatmap representation of raw counts for genes shared between BA71V and GRG, generated with

1208    pheatmap. (b) broad patterns represented by genes in the 5 clusters indicated in (a). (c) histogram

1209    showing the percentage of the total raw reads per gene which are detected at 16 hpi vs. 5 hpi post-

1210    infection, and comparing the distribution of percentages between GRG and BA71V. (d) Mean read

1211    counts from GRG at 5 hpi vs 16 hpi replicates, showing a significant increase (T-test, p-value: 0.045)

1212    from 5 hpi to 16 hpi.

1213

1214    Figure 7. Promoter motifs and initiators detected in early and late ASFV GRG TSSs including alternative

1215    TSSs and those for nORFs. (a) Consensus of 30 bp upstream and 5 bp downstream of all 134 early TSSs

1216    including nORFs, with the conserved EPM (10) and Inr annotated. (b) 30 bp upstream and 5 bp

1217    downstream of all 234 late gene and nORFs TSSs, with the LPM and Inr annotated (c) The conserved

1218    EPM detected via MEME motif search of 35 bp upstream for 133 for 134 early TSSs (E-value: 3.1e-

1219    069). The conserved LPM detected via MEME motif search of 35 bp upstream for 46 for 234 late gene

1220    TSSs (E-value: 2.6e-003). The locations of the EPM shown in (b) and LPM shown in (d) are annotated

1221    with brackets in (a) and (b), respectively. Motifs detected via MEME search of 35 bp upstream of genes

1222    in clusters from Figure 6: cluster 1 (7 genes, E-value: 9.1e-012), 2 (15 genes, E-value: 2.6e-048), 3 (60

1223    genes, E-value: 1.0e-167), 4 (32 genes, E-value: 4.7e-105), 5 (16 genes, E-value: 5.7e-036), are shown

1224    in e-i, respectively. For ease of comparison, (e), (g), (i) and (f), (h) are aligned at TSS position. All motifs

1225    were generated using Weblogo 3 (113). (k) shows the distribution of MEME motif-end distances, from

1226    last nt (in coloured bracket), to their respective downstream TSSs.

1227

1228    Figure 8. The TSSs of MGF 360-19R. Panels (a) 5 hpi and (b) 16 hpi show CAGE-seq 5' end data from

1229    these time-points, in red are reads from the plus strand and blue from the minus strand, the RPM

1230    scales are on the right. (c) TSSs are annotated with arrows if they can generate a minimum of 5 residue-

1231    ORF downstream, and grey bars indicate where they are located on the CAGE-seq coverage in (a) and

1232    (b). ORFs identified downstream of TSSs are shown as red arrows (visualized with R package gggenes),

1233    including three short nORFs out of frame with MGF 360-19R. Also shown are three in-frame truncation

1234    variants, from TSSs detected inside the full-length MGF 360-19R 269-residue ORF, downstream of its

1235    pTSS at 185213. Blue or yellow boxes upstream of TSSs indicate whether the EPM or LPM

1236    (respectively) could be detected within 35 nt upstream of the TSS using FIMO searching (114).

1237

40

1238 Figure 9. Summary of intra-ORF TSSs (ioTSSs) and nORFs detected in the GRG genome, further

1239 information in Supplementary Table 2. (a) Summarises the gene types in which ioTSSs were detected,

1240 showing an overrepresentation of MGFs, especially from families 360 and 505, furthermore, the

1241 majority of ioTSSs are detected at 16 hpi. (b) For ioTSSs in-frame with the original, summarised are

1242 the subsequent UTR lengths i.e. distance from TSS to next in frame ATG start codon, which could

1243 generate a truncation variant. (c) Example of a miss-annotation for CP204L, whereby the pTSS is

1244 downstream the predicted start codon. (d) and (e) show the results of 5'RACE for three genes (DP146L,

1245 pNG4, and CP204L, see methods for primers), at 5 hpi and 16 hpi, respectively. Examples of genome

1246 regions around DP146L (f) and pNG4 (g), wherein ioTSSs were detected with capacity for altering ORF

1247 length in subsequent transcripts, and therefore protein output. Primers used for 5'RACE for DP146L

1248 and pNG4 are represented as black arrows in (f) and (g), respectively.

1249

1250 Figure 10. MGF 100 genes likely encode SH2-domain factors. (a) Occurrence of MGF 100 genes in

1251 selected ASFV strains, with genotype and pathogenicity indicated (as yes. Y, or no, N). '1L/2L' refers to

1252 the gene MGF 100-2L (DP141L in BA71V) and MGF 100-1L in the FR682468.1 genome annotation. (b)

1253 The top panel illustrates representative SH2 domain structures (Suppressor of Cytokine Signalling 1

1254 and -2 and the PI3K alpha), and the bottom shows structural homology models of MGF 100 members

1255 1L, 1R, and I7L and I8L superimposed. The PHYRE2 algorithm (56) was used to predict models for MGF

1256 100 members (Supplementary Table 2d), and the structures at the top were detected as the top hits

1257 for each of the MGF 100 models shown in the lower panel. (c) Structure-guided multiple sequence

1258 alignment of selected MGF 100 member models, alongside known SH2 domain structures (annotated
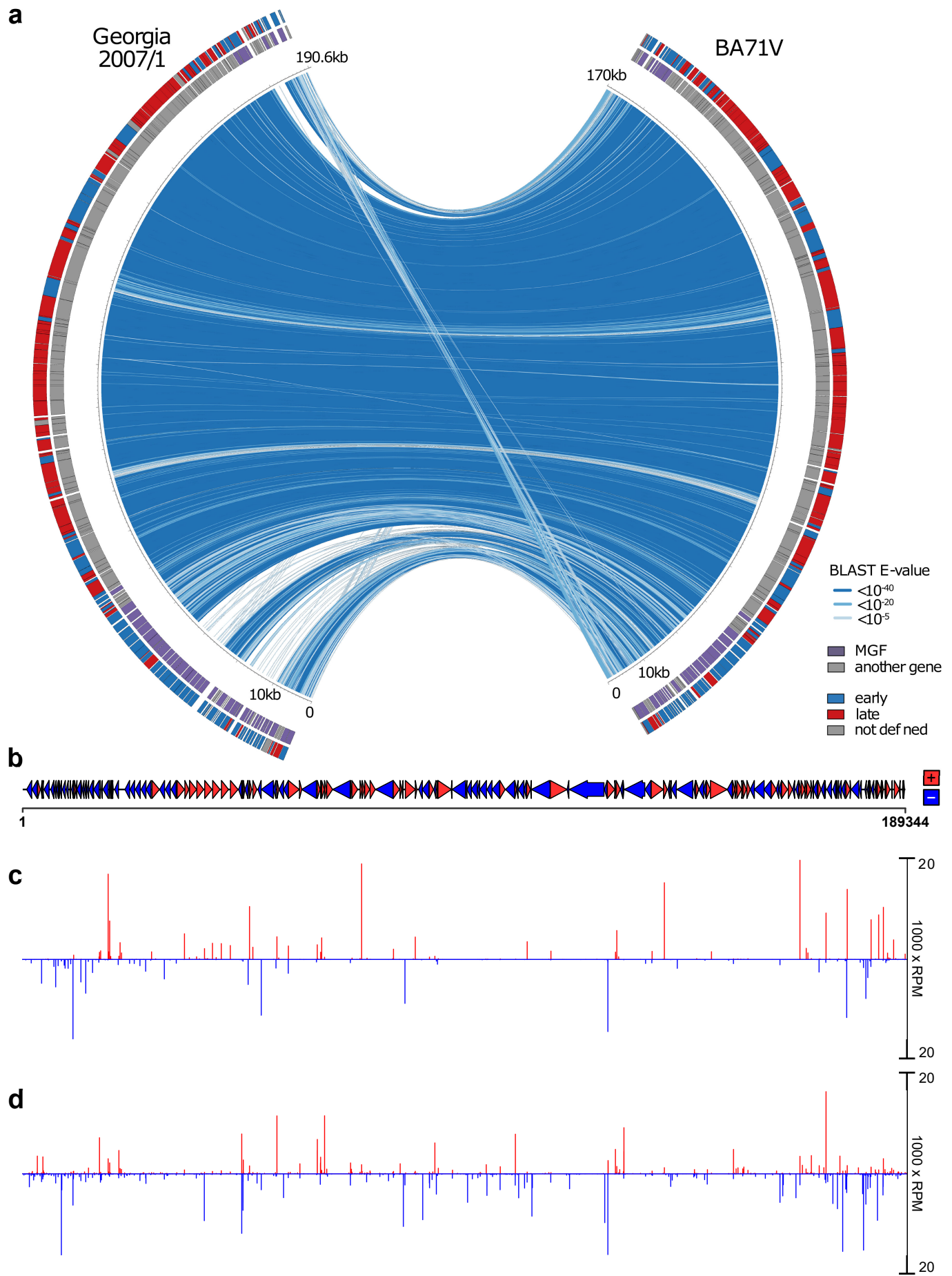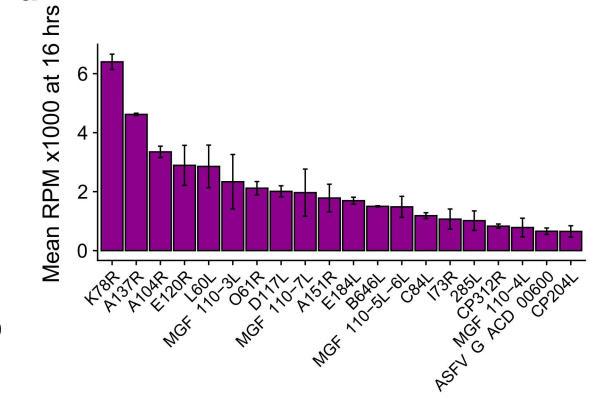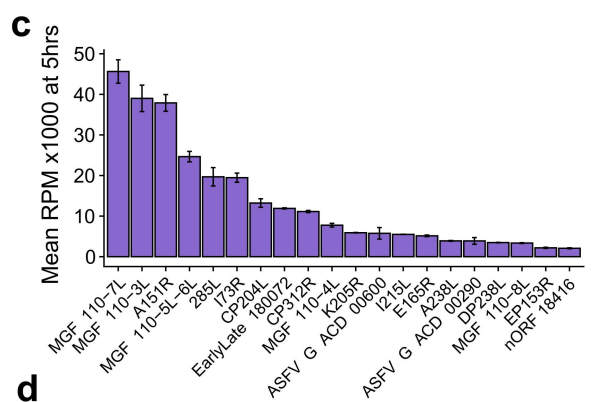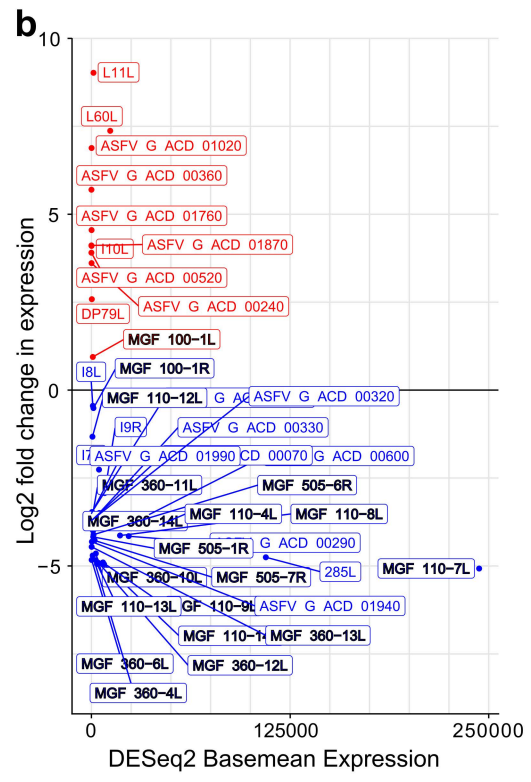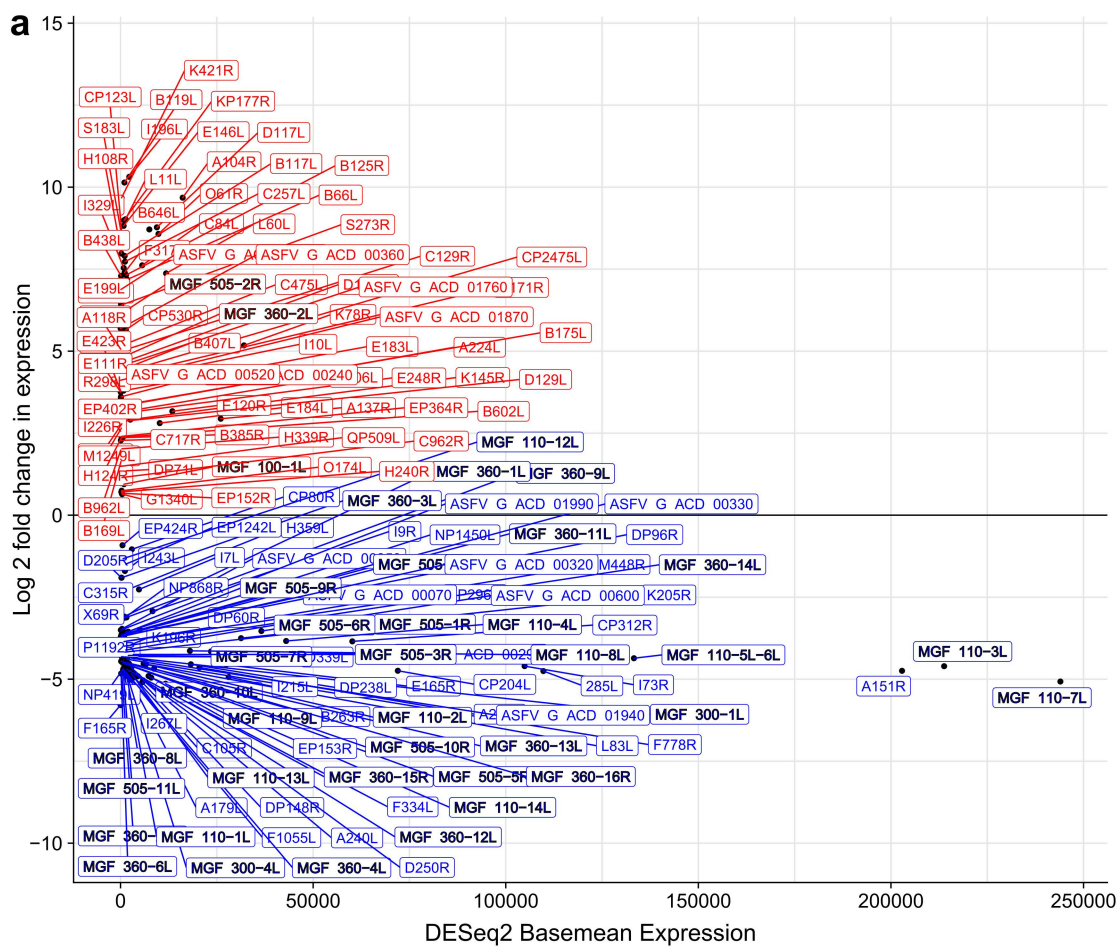
1259 as SH2_name_PDB number).

1260

1261 Figure 11. Changes in the swine macrophage transcriptome upon ASFV GRG infection. (a) Major

1262 expression response profiles of the pig macrophage transcriptome. Late response genes are

1263 significantly deregulated (false discovery rate < 0.05) in one direction both between mock-infected

1264 (ctrl) and 16 hpi as well as between 5 and 16 hpi, but not between mock-infected and 5 hpi. Early

1265 response genes are significantly deregulated in one direction both between ctrl and 5 hpi as well as

1266 ctrl and 16 hpi, but not between 5 and 16 hpi. (b) Relationship of log fold changes (logFC) of TSS-

1267 derived gene expression levels of the total 9,384 swine genes expressed in macrophages between 5–

1268 16 hpi and ctrl–16 hpi. Colors correspond to the response groups from the panel a. (c) Relationship of

1269 log fold changes of TSS-derived gene expression levels of the total 9,384 swine genes expressed in
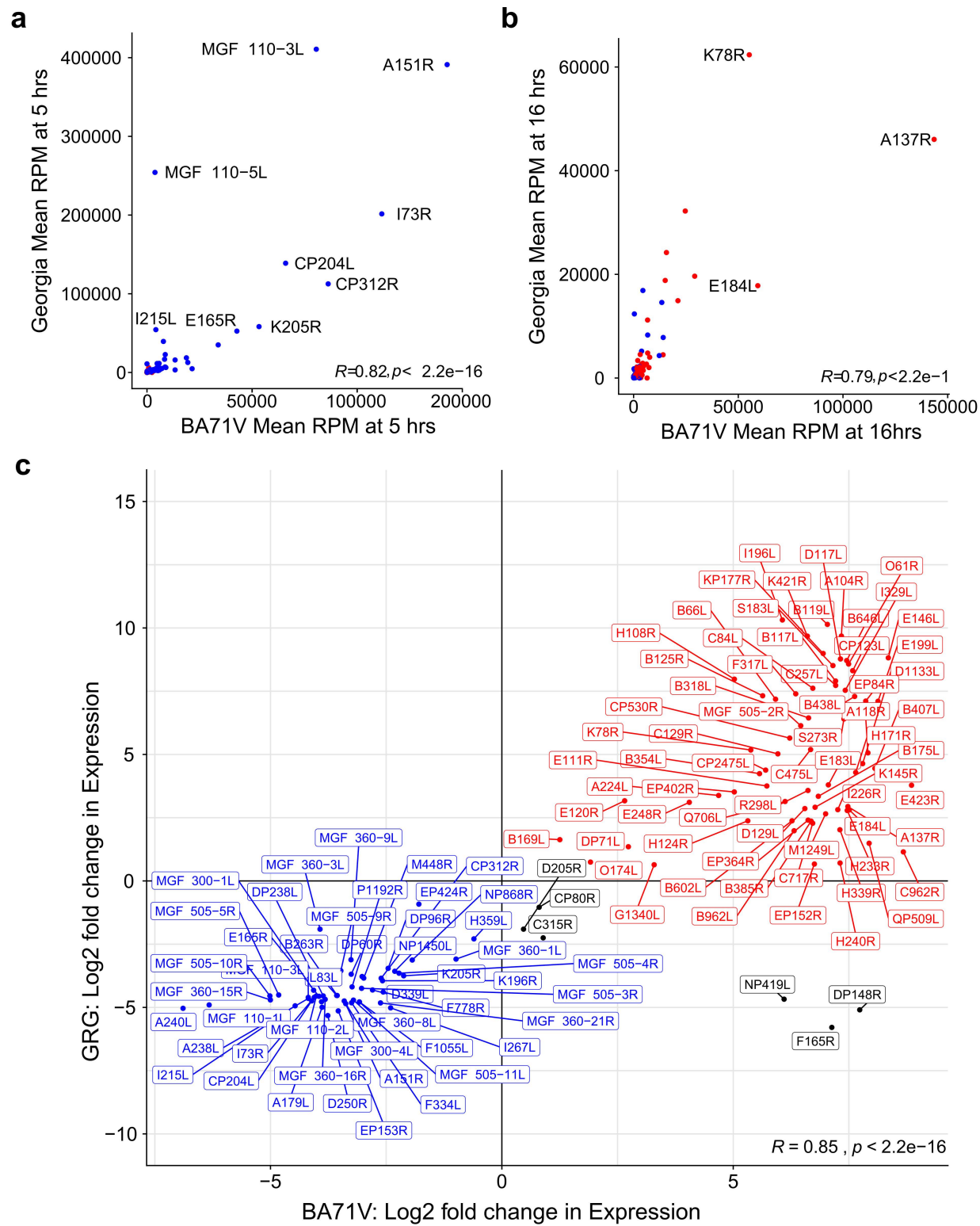
41

1270    macrophages between 5–16 hpi and ctrl–5 hpi. (d) MA plot of the TSS-derived gene expression levels

1271    between 5 and 16 hpi based on differential expression analysis with edgeR (108,115). (e)

1272    Representative overrepresented functional annotations of the upregulated (red) and downregulated

1273    (blue) macrophage genes following late transcription response (Benjamini-corrected p-value lower

1274    than 0.05). Numbers on the right to the bars indicate total number of genes from a given group

1275    annotated with a given annotation. (f) RT-PCR of four genes of interest indicated in (d). 'C' is the

1276    uninfected macrophage control, NTC is the Non Template Control for each PCR, excluding template

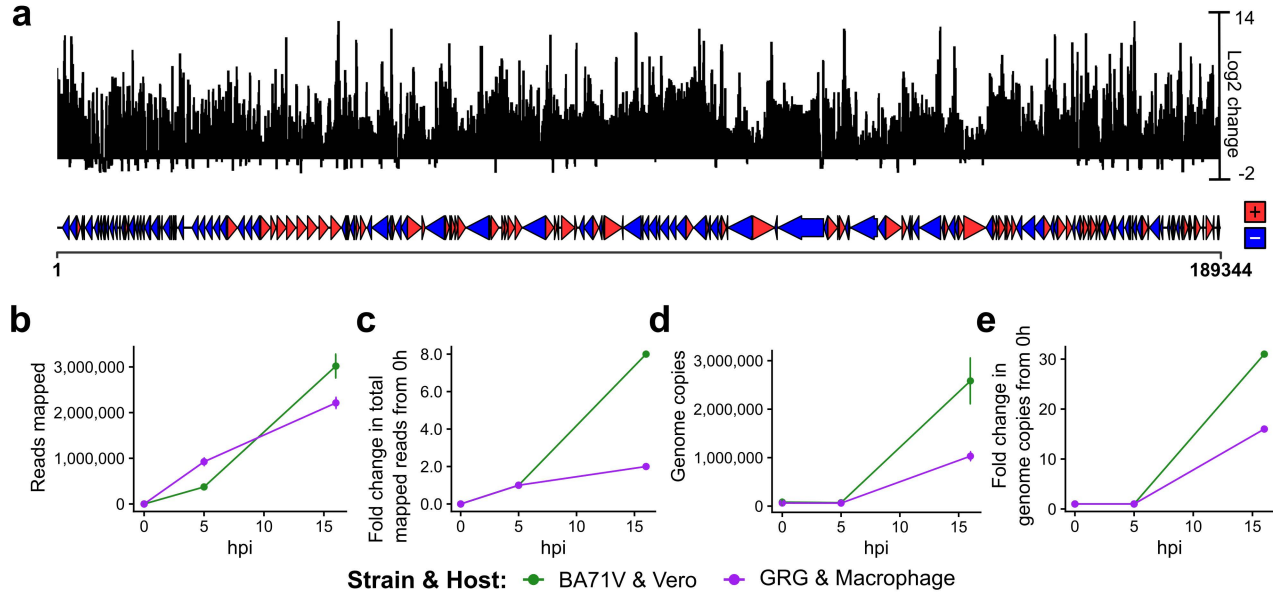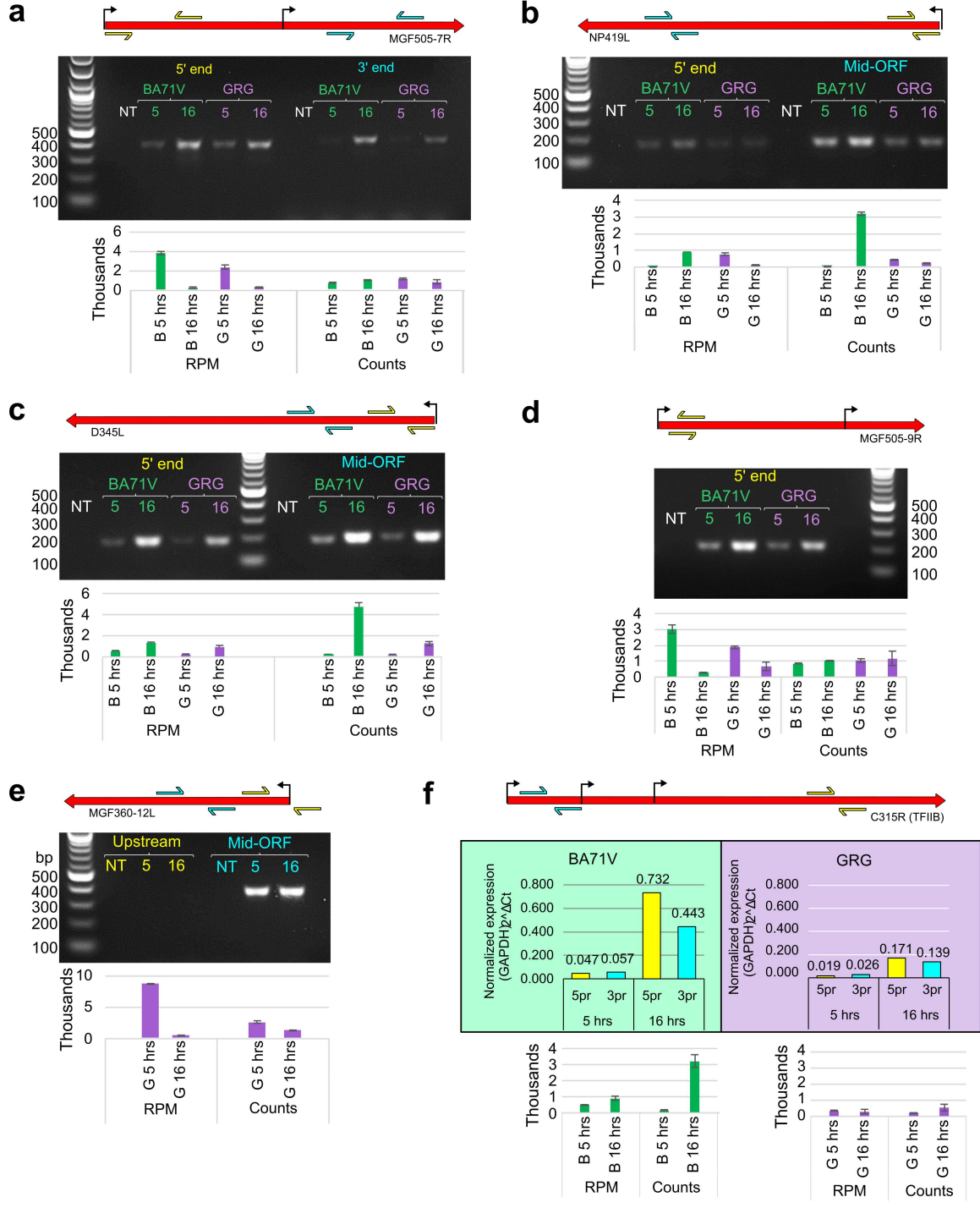1277    DNA. See methods for primers used.

1278

1279    Figure 12. Protein expression at different times during infection of swine macrophages with ASFV-

1280    GRG. Two different batches of macrophages (S1 and S2) were infected with MOI 5 or left uninfected

1281    as a control (Ctrl) and at 0, 5 and 16hpi cellular extracts were collected and analysed via SDS-PAGE

1282    Western blot for the presence of ISG15 and γ-Tubulin as a protein loading control (a) and for the

1283    presence of viral protein P30 as control of ASFV infection (b). (c), (d) and (e) are the results from ELISAs

1284    for detection of porcine TNF-α, IL-8/CXCL8, and CCL2/MCP-1, respectively, in culture supernatants.

1285    Results are presented as 'Relative to control' values (y-axis of c-e) calculated by performing ELISAs in

1286    parallel for control and GRG infection at each timepoint.

42

**a**

Georgia
2007/1

190.6kb

170kb

BA71V

BLAST E-value
$<10^{-40}$
$<10^{-20}$
$<10^{-5}$

MGF
another gene

early
late
not def ned

10kb

0

10kb

0

**b**

1                                                                                                    189344

+
−

**c**

20

1000 × RPM

20

**d**

20

1000 × RPM

20

**a**

Georgia Mean RPM at 5 hrs vs BA71V Mean RPM at 5 hrs

$R=0.82, p< 2.2e-16$

**b**

Georgia Mean RPM at 16 hrs vs BA71V Mean RPM at 16hrs

$R=0.79, p<2.2e-1$

**c**

GRG: Log2 fold change in Expression vs BA71V: Log2 fold change in Expression

$R = 0.85 , p < 2.2e-16$

a

b

| | Counts 5h | Counts 16h | Cluster pattern |
|---|---|---|---|
| 1 | high | high | both high (H-H) |
| 2 | low | high | low to high (L-H) |
| 3 | low | mid | low to mid (L-M) |
| 4 | mid | mid | both mid (M-M) |
| 5 | low-mid | low-mid | both low-mid (LM-LM) |

c

d

**a** Time ioTSS detected: Early (blue) Late (red) Both (purple)

Virus structure / assembly
Uncharacterised
Transcription / RNA modification
PSP / TR
NAm / DNA replication / repair
MGF 505
MGF 360
MGF 300
MGF 100
Infection / immune evasion
Enzymes / other

Number of ioTSS (0, 5, 10, 15)

**b** Number of genes vs UTR length (nt) (0, 7, 14, 21, 28, 35, 42, 49)

**c** 50000 RPM 50000
newly-annotated start codon
original annotated start codon
CP204L
125270 — 125400

**d** DP146L pNG4 CP204L (1000, 500, 400, 300, 200, 100)

**e** DP146L pNG4 (1000, 500, 400, 300, 200, 100)

**f** 5h 8000 RPM 8000 / 16h 8000 RPM 8000
Frame 1 DP146L-like (nORF 180574)
180250 — 180415 — 180580

**g** 5h 38000 RPM 38000 / 16h 7200 RPM 7200
Frame 3 / 2 / 1
nORF 16717
ASFV_G_ACD_00290
pNG4 (nORF 16814)
ASFV_G_ACD_00300
16550 — 16820 — 17090

**a**

| Genotype | ASFV strain | Pathogenic | Source | 1R | 1L/2L (DP141L) | 3L (DP146L) | I7L | I8L |
|---|---|---|---|---|---|---|---|---|
| I | BA71V | N | NC_001659.2 | | ← | ← | | |
| I | BA71 | Y | NC_044942 | // | ← | ← | ← | ← |
| I | Portugal, OURT88/3 | N | AM712240.1 | → | // ← | ← | ← | ← |
| I | Benin 97/1 | Y | AM712239.1 | | ← | ← | ← | ← |
| II | Georgia 2007/1 | Y | FR682468.2 | → | // ← | ← | ← | ← |
| II | Georgia 2007/1-VP110 | N | Krug *et al.* | → | // ← | ← | | |
| II | China/2018/AnhuiXCGQ | Y | MK128995.1 | → | // ← | ← | ← | ← |
| IX | Ken05/Tk1 | Y | KM111294.1 | → | // ← | ← | ← | ← |
| X | Kenya 1950 | Y | AY261360.1 | → | // ← | ← | ← | ← |

**b**



**PDB ID: 2C9W .A**
Suppressor of Cytokine Signalling 1

**PDB ID: 6C5X .D**
Suppressor of Cytokine Signalling 2

**PDB ID: 4L1B .B**
PIK3 regulatory subunit alpha

**Model: MGF 100-1L (GRG)**
Template PDB ID: 2C9W .A

**Model: MGF 100-1R (GRG)**
Template PDB ID: 6C5X .D

**Model: I7L & I8L (GRG)**
Template PDB ID: 4L1B .B

**c**

**a**

0hpi　　5hpi　　16hpi

Ctrl　GRG　Ctrl　GRG　Ctrl　GRG

S1

ISG15

γ-Tub

Density
ISG15/γTub　　0.70　0.85　1.00　0.89　1.26　1.11

S2

ISG15

γ-Tub

Density
ISG15/γTub　　1.01　1.13　1.11　0.65　0.61　0.56

**b**

S1

Ctrl　GRG 0hpi　GRG 5hpi　GRG 16hpi

P30

S2

Ctrl　GRG 0hpi　GRG 5hpi　GRG 16hpi

P30

**c**

■ S1　■ S2

TNF-α relative to Ctrl

GRG-0h　GRG-5h　GRG-16h

**d**

■ S1　■ S2

CXCL8 relative to Ctrl

GRG-0h　GRG-5h　GRG-16h

**e**

■ S1　■ S2

CCL2 relative to Ctrl

GRG-0h　GRG-5h　GRG-16h