

1 Homologous Recombination as an Evolutionary Force in 2 African Swine Fever Viruses

3

4 Zhaozhong Zhu^{1,#}, Chao-Ting Xiao^{1,#}, Yunshi Fan¹, Zena Cai¹, Congyu Lu¹, Gaihua
5 Zhang², Taijiao Jiang^{3,4}, Yongjun Tan¹, Yousong Peng^{1,*}

6

7 ¹ College of Biology, Hunan University, Changsha, China

8 ² College of Life Sciences, Hunan Normal University, Changsha 410081, China

9 ³ Center of System Medicine, Institute of Basic Medical Sciences, Chinese Academy
10 of Medical Sciences & Peking Union Medical College, Beijing, China

11 ⁴ Suzhou Institute of Systems Medicine, Suzhou, China

12 # These authors contributed equally to this work

13 * To whom correspondence should be addressed. Email: pys2013@hnu.edu.cn

14

15

16 **Abstract**

17 Recent outbreaks of African swine fever virus (ASFV) in China severely
18 influenced the swine industry of the country. Currently, there is no
19 effective vaccine or drugs against ASFVs. How to effectively control the
20 virus is challenging. In this study, we have analyzed all the publicly
21 available ASFV genomes and demonstrated that there was a large genetic
22 diversity of ASFV genomes. Interestingly, the genetic diversity was

23 mainly caused by extensive genomic insertions and/or deletions (indels)
24 instead of the point mutations. The genomic diversity of the virus resulted
25 in proteome diversity. Over 250 types of proteins were inferred from the
26 ASFV genomes, among which only 144 were observed in all analyzed
27 viruses. Further analyses showed that the homologous recombination may
28 contribute much to the indels, as supported by significant associations
29 between the occurrence of extensive recombination events and the indels
30 in the ASFV genomes. Repeated elements of dozens of nucleotides in
31 length were observed to widely distribute and cluster in the adjacent
32 positions of ASFV genomes, which may facilitate the occurrence of
33 homologous recombination. Moreover, two enzymes, which were
34 possibly related to the homologous recombination, i.e., a Lambda-like
35 exonuclease with a YqaJ-like viral recombinase domain, and a DNA
36 topoisomerase II, were found to be conservative in all the analyzed
37 ASFVs. This work highlighted the importance of the homologous
38 recombination in the evolution of the ASFVs, and helped with the
39 strategy development of the prevention and control of the virus.

40

41

42 **Introduction**

43 African swine fever virus (ASFV), the causative agent of African swine
44 fever (ASF), is a complex, large, icosahedral multi-enveloped DNA virus.

45 It is classified as the only member in the family *Asfarviridae*^{1,2}. The
46 genome of the virus belongs to double-stranded DNA, with the size
47 ranging from 170 kb to 190 kb³. ASFV mainly infect suids and soft ticks.
48 The suids include domestic pigs and wild boars, and were reported as the
49 natural hosts of the virus^{4,5}. ASFV was firstly discovered in Kenya in
50 1921⁶. It remained restricted in Africa till 1957, when it was reported in
51 Spain and Portugal. Up to now, the virus has caused ASF outbreaks in
52 more than fifty countries in Africa, Europe, Asia, and South America⁴.
53 The latest reports showed that the virus has caused outbreaks in more
54 than fifteen provinces in China^{7,8}. Because of the high lethality of ASFV
55 in domestic pigs, the most commonly used strategies to control the virus
56 were the massive culling campaigns and the restriction of pig movement⁵.
57 Both strategies have resulted in a huge economic loss for pig industry and
58 affected people's livelihoods. Unfortunately, currently there is no
59 available effective vaccine against ASFVs.

60

61 Many efforts have been devoted to developing the vaccine for the ASFV
62 ^{1,5,9-11}, however, most of these attempts failed. One of the most important
63 reasons was the complex composition of the antigenic proteins^{5,12}.
64 Previous reports showed that p72, p30, and p54 were the three important
65 antigenic proteins during the infection of ASFVs, but the immunity
66 against them could only provide a partial protection^{12,13}. Many other

67 proteins or other factors such as phospholipid composition may also
68 influence the antigen of the virus ¹². Therefore, it is necessary to
69 understand the mechanisms of the antigen diversity of the ASFV virus ¹.
70
71 The genetic diversity of ASFVs has been investigated in many studies.
72 The ASFV genome encodes over 150 proteins, including viral enzymes,
73 viral transcription and replication-related proteins, structural proteins,
74 other proteins involved in the virus assembly, the evading of host defense
75 systems and the modulation of host cell function, etc ^{3,14,15}. For example,
76 the transcription of the virus is independent on the host RNA polymerase
77 because the virus contains relevant enzymes and factors ³. The viral
78 genome contains a conservative central region of about 125 kb and two
79 variable ends, which results in the variable size of the genome ^{3,16,17}.
80 There are significant variations among the ASFV genomes due to the
81 genomic insertion or deletion, such as the deletion of the multigene
82 family (MGF) members ³. Although much progress have been made on
83 genetic diversity of the virus, the extent and mechanisms are still not
84 clear. Besides, most of these studies either only investigated the genetic
85 diversity of some common genes, such as p72 and p54 ^{18,19}, or only used
86 one or several isolate genomes ^{3,16,17}. The number of discovered viral
87 genomes has increased rapidly as the development of DNA sequencing
88 technology. Therefore, a comprehensive study on the genetic diversity of

89 ASFVs is necessary.

90

91 Homologous recombination, which has been reported to occur in several
92 groups of viruses²⁰⁻²³, such as herpesvirus, retroviruses, and
93 coronaviruses, has played an important role in viral evolution²¹. A few
94 studies on several ASFV genes have suggested the occurrence of
95 homologous recombination in the evolution of ASFVs^{3,18}. However, a
96 comprehensive study on the homologous recombination in ASFV at the
97 genomic scale is lacking, and the role of the recombination on the genetic
98 diversity and the evolution of the virus is still unknown. In this study, we
99 have systematically investigated the genomic diversity and the
100 homologous recombination of ASFVs based on the analysis on all the
101 publicly available ASFV genomes. The results demonstrated that the
102 homologous recombination contributed much to the genetic diversity of
103 ASFVs. This work would help to understand the evolution of the ASFV
104 and thus facilitate the prevention and control of the virus.

105

106 **Results**

107 *1 ASFV genomes*

108 A total of 36 genome sequences of ASFVs were obtained from the NCBI
109 GenBank database, which were listed in Table S1. They were mainly
110 isolated from Africa and Europe during the years from 1950 to 2017. The

111 size of the ASFV genomes ranged from 170,101 bp to 193,886 bp,
112 averaged at 185,800 bp. The viral isolate Kenya50 had the largest size,
113 while the isolate BA71V had the smallest size. No increasing or
114 decreasing trend in the genome size was observed from 1950 to 2017
115 (Figure 1A), suggesting the dynamic changes of the viral genomes.

116

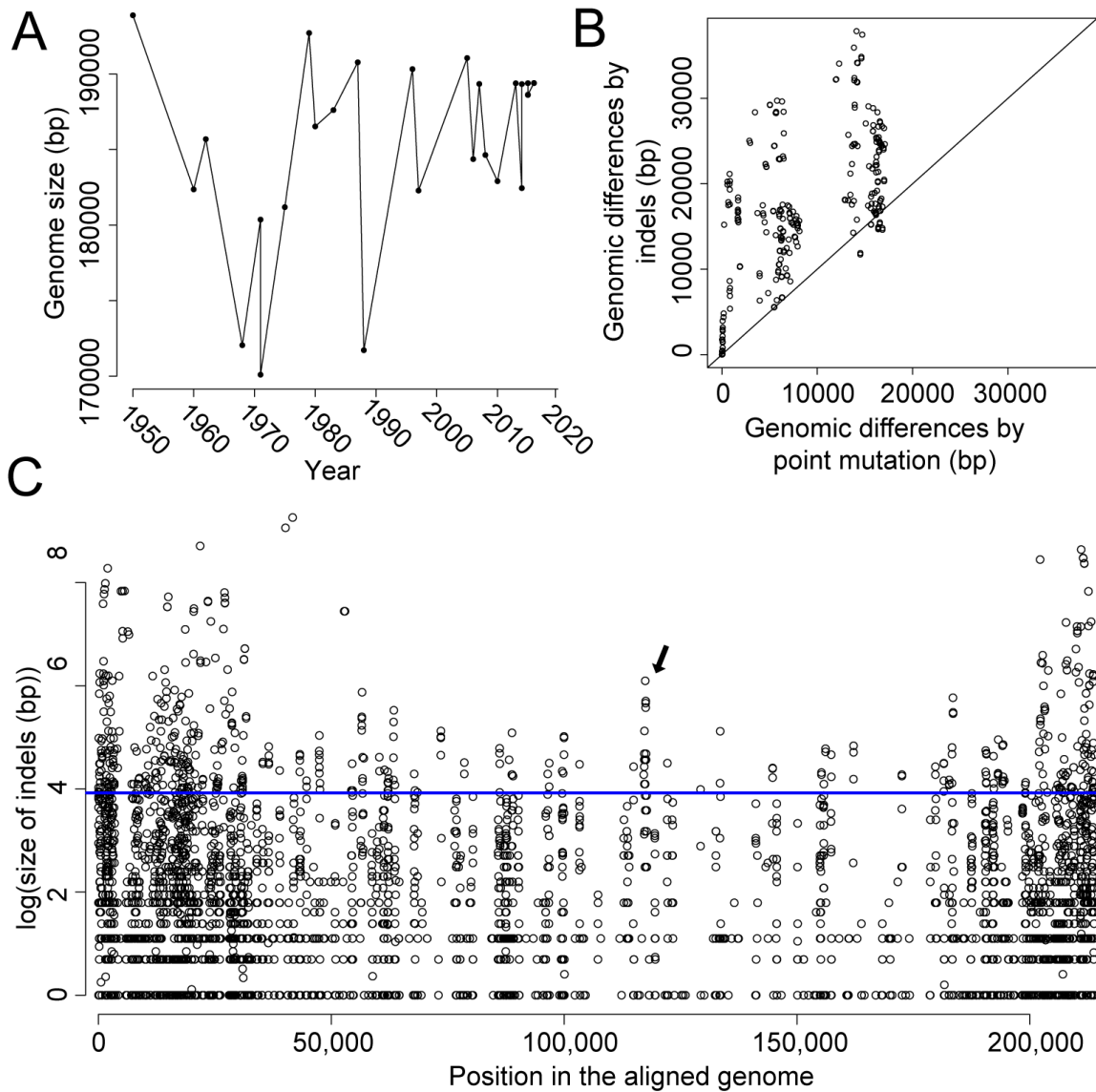
117 ***2 Genomic diversity of ASFVs***

118 Pairwise comparisons between ASFV genomes were conducted after the
119 genome alignment. The average genomic difference between viruses was
120 24,570 bp, which accounts for more than 10% of the genome.

121 Interestingly, the genomic differences caused by the insertions and
122 deletions (indels) were much more significant than those caused by the
123 point mutations (Figure 1B & Figure S1) in most cases. For example,
124 there were 31,833 bp differences between virus Mkuzi79 and BA71V, 78%
125 of which were caused by indels.

126 The size and position of indels in ASFV genomes were also analyzed. 70%
127 of indels were no longer than 10 bp, and about 10% of indels were 50 bp
128 or longer (Figure S2). The occurrence of indels was much more frequent
129 in both ends of the genome, especially in the 5' end (Figure 1C). Besides,
130 the size of indels in both ends was also much larger than that in the
131 middle region. Large indels with over 50 bp (above the blue line in Figure
132 1C) were mostly observed in both ends. It should be noted that the

133 variation to some extent was observed in the middle region (marked in
134 black arrow), which were considered to be conservative in previous
135 reports.



136

137

138 **Figure 1.** Genomic differences between ASFV genomes. (A) The
139 variation of genome size along the isolation time of the virus. (B)
140 Genomic differences caused by point mutations and indels. (C) The size
141 and location of indels along the aligned genomes. For clarity, the natural

142 logarithm of the indel size was used. Position for an indel is defined as
143 the middle position of the indel. The average size was used if more than
144 one indel was found in the position. The blue line refers to indel size of
145 50 bp.

146

147 ***3 Proteome diversity of ASFV***

148 Genomic diversity could lead to proteome diversity. Therefore, the
149 proteome diversity of ASFVs was further analyzed. Firstly, the candidate
150 proteins encoded by ASFV genomes were inferred (Materials and
151 Methods) (Table S2). The plus strand encoded 95-126 proteins, with an
152 average of 109 proteins; the minus strand encoded 106-128 proteins, with
153 an average of 118 proteins. Considering both the plus and minus strands,
154 the ASFV genome encoded 205-254 proteins, with an average of 227
155 proteins. The viral isolate Russia14 encoded most proteins, although the
156 size of this genome was not the largest. The ratio of coding region in each
157 genome ranged from 88% to 91%.

158

159 Furthermore, the ortholog or paralog groups based on sequence homology
160 were identified. A total of 252 protein groups plus 28 singletons were
161 obtained (Table S3), each of which stood for one type of protein encoded
162 by ASFV genomes. The obtained proteins contained almost all the
163 proteins identified in previous experiments (Table S3). Each protein

164 group included 2-99 proteins. Only 144 protein groups were observed in
165 all 36 ASFV viruses, which could be considered as core protein sets of
166 the virus, and were mainly encoded by both plus and minus strands in the
167 middle regions (Figure 2). The protein groups could be further separated
168 into seven classes by function based on previous studies (Figures 2 & S3).
169 Only about 30% of protein groups were observed to have the known
170 functions, including replication and transcription (in red), host cell
171 interactions (in magenta), structure and morphogenesis (in blue), and
172 enzymes (in yellow). Most of the above-described protein groups with the
173 known functions belonged to the core proteins of the virus. In addition,
174 forty protein groups belonged to “Multigene Families (MGF)” (in cyan),
175 most of which had unknown functions. The MGFs were encoded by both
176 ends of the genome. Besides, the remaining 146 protein groups belonged
177 to either the class of “Proteins with unknown function” (in gray) or
178 “Hypothetical proteins” (in black).

179

180 Analysis of the protein conservation showed that except the functional
181 class of MGFs, the proteins in other functional classes had an average of
182 pairwise sequence identities greater than 90% (Figure S4). The proteins
183 in functional classes of “Other enzymes” and “Replication &
184 transcription” were most conservative, with average pairwise sequence
185 identities larger than 95%. Proteins of these two functional classes also

186 had the smallest ratios of dN/dS (Figure S5), suggesting strong negative
187 selection on them. While the proteins in the functional class of MGF and
188 “Hypothetical proteins” had the largest ratio of dN/dS. The hypothetical
189 proteins had a median dN/dS ratio of 0.82, suggesting strong positive
190 selection on these proteins.

191

192 Membrane proteins, which may be located in the inner or outer envelope,
193 were observed to be distributed widely in the proteome of ASFVs
194 (marked with asterisks in Figure 2). A total of 67 protein groups belonged
195 to membrane proteins, including 35 in the core protein groups, such as
196 p54 and EP402R. Among the membrane proteins, only 11 protein groups
197 had the known functions, including 8 protein groups in the functional
198 class of “Structural and morphogenesis”, and 1 protein group in each
199 functional class of “Host cell interactions”, “Replication & transcription”
200 and “Other enzymes”.

201

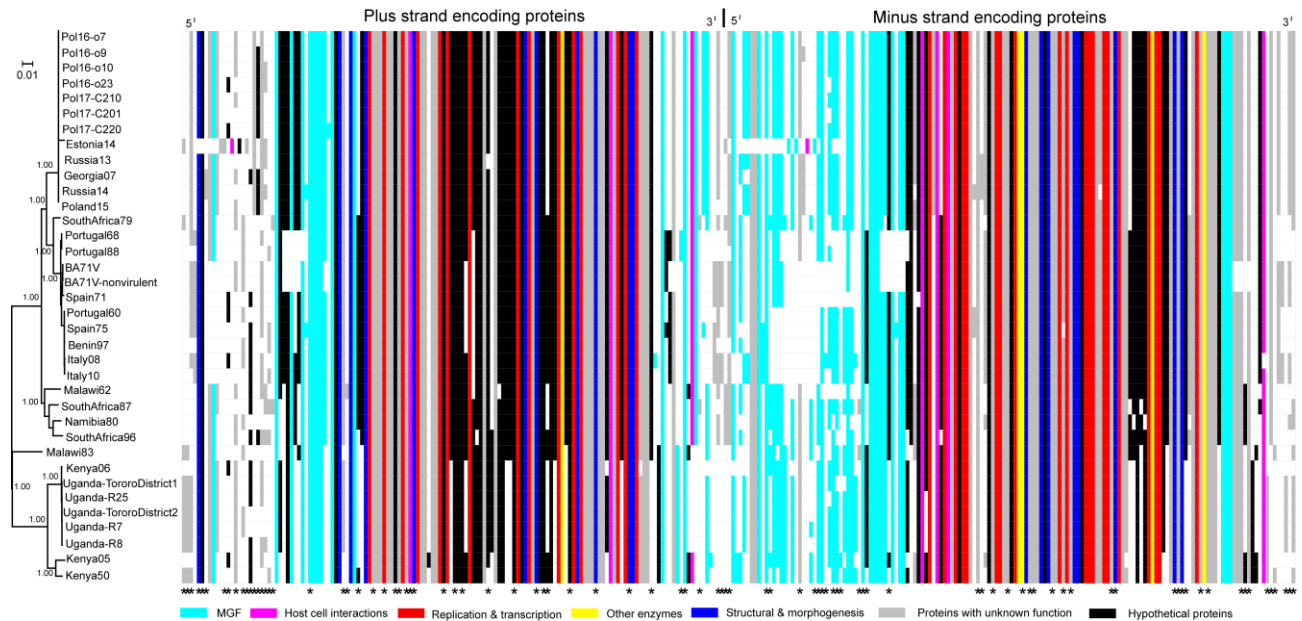
202 Thirty-one protein groups were observed to have paralogs (duplicated
203 proteins) in at least one virus (colored in Figure S6). They were mostly
204 located in both ends of the genome. Thirteen of them belonged to MGFs.
205 In addition, two protein groups, “DP71L” and “DP96R”, belonged to the
206 class of “Host cell interactions”. The rest protein groups belonged to
207 either the class of “Proteins with unknown function” or “Hypothetical

208 proteins”. Most of the paralogs were clustered in adjacent positions.
209 Exceptions were observed for some protein groups which were encoded
210 by the first one to three thousands nucleotides in the plus and minus
211 strands, such as the protein group “p01990-3L” (marked with black
212 arrows in Figure S6). Further analysis showed that a segment of 200-3000
213 bp was exactly the same in the beginning of the plus and minus strands in
214 most viral genomes (Table S4).

215

216 Extensive insertion and deletions of proteins were observed in the
217 proteome of ASFVs after alignment. Viruses in the adjacent positions in
218 the phylogenetic tree tended to have similar proteomes. The number of
219 different proteins between different viruses ranged from 1 to 84, with an
220 average of 43, which was about one-fifth of the viral proteome. The
221 differences of the proteome among the viruses were mainly caused by
222 proteins of the class of “Hypothetical proteins”, “Proteins with unknown
223 function” and “MGF” (Figure 2).

224



225

226 **Figure 2.** The phylogenetic tree of ASFVs and the alignment of their
227 proteomes in plus (left side) and minus strand (right side). Each row
228 refers to the proteome of the virus in the phylogenetic tree; each column
229 refers to one protein group. Protein groups were colored according to
230 their functions. “White” refers to no protein group in the virus. Asterisks
231 in the bottom refer to membrane proteins. For clarity, the singletons were
232 ignored in the alignment.

233

234 ***4 Extensive homologous recombination in ASFV genomes***

235 As numerous indels have been revealed in the ASFV genomes, then, we
236 investigated the mechanism of generating indels. According to the results
237 in previous studies, three factors may contribute to the extensive indels in
238 ASFVs: replication slippage, retrotransposition and recombination²³.

239 Replication slippage mainly produced duplications of short genetic
240 sequences and may cause short indels, but it is unlikely to generate large

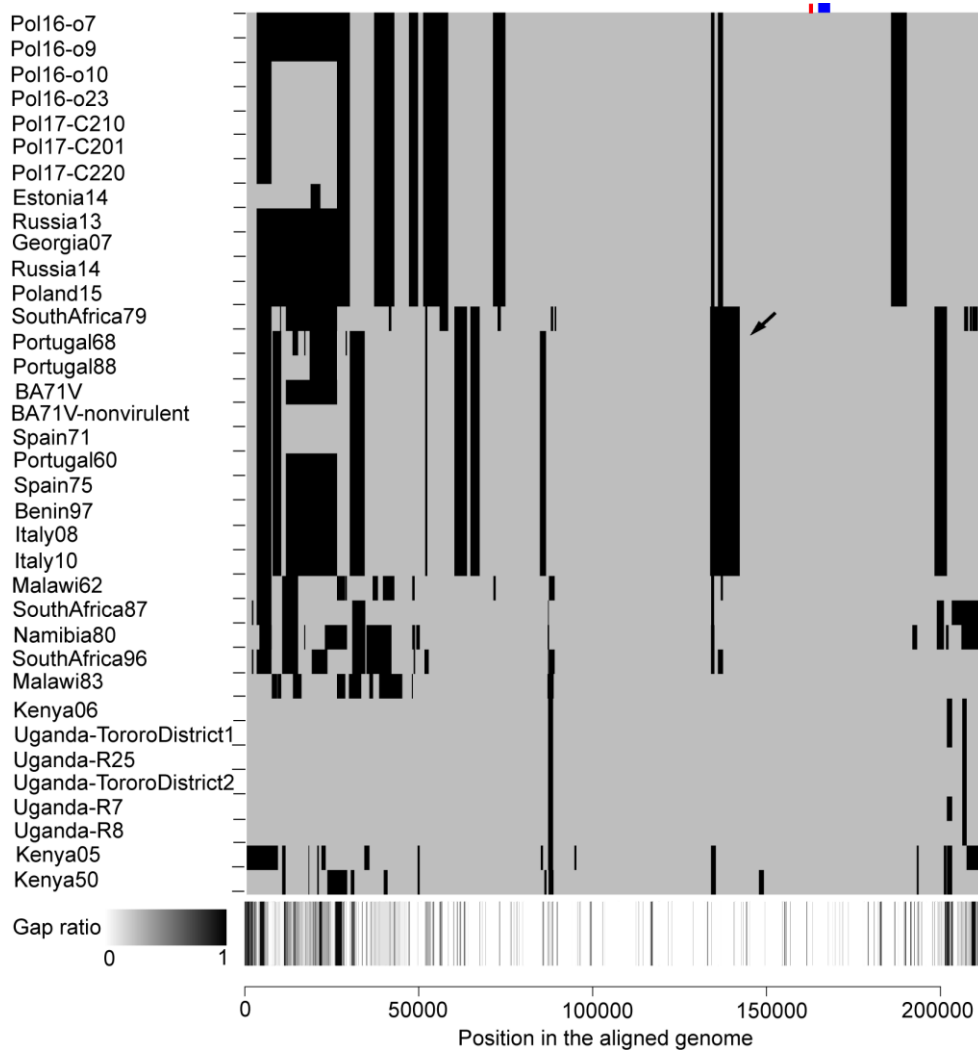
241 indels observed in ASFVs. Retrotransposition can result in duplication of
242 large genetic sequences or genes, but the location of duplicates would be
243 randomly distributed in the genome. However, most paralogs shown in
244 Figure S6 were clustered in adjacent positions, thus these paralogs may
245 be not caused by retrotransposition. Besides, no retrotransposons were
246 observed in the analyzed ASFV genomes (as described in Materials and
247 Methods).

248

249 Finally, we investigated the role of recombination in the generation of
250 indels in the ASFV genomes. The analyses on the recombination showed
251 that there were a total of 103 recombination events, and each ASFV
252 genome had 3-22 recombination events (Figure 3 & Table S5). The virus
253 isolate SouthAfrica79 experienced the largest number of recombination
254 events. On average, each virus experienced 11 recombination events. The
255 sizes of recombination region ranged from 174 to 22,628 bp. The ratio of
256 recombination region in each genome ranged from 2% to 27%. In total,
257 the regions in the ASFV genomes involved in all recombination events
258 covered a total of 101,569 nucleotide positions, accounting for 47% of
259 the aligned genome. Most recombination events happened at both ends,
260 especially at the 5' end. Interestingly, the recombination event in the
261 aligned genomes was observed to be consistent with the ratio of the gap
262 in the genome (the bottom of Figure 3). Almost all the recombination

263 events happened in or close to the gap-rich regions where the indels were
264 observed. The ratios of the gaps in the recombination regions were found
265 to be much larger than those in other regions (Figure S7). Further
266 comparison of the number of indels in the recombination regions and
267 other regions showed that for indels of varying length, such as those
268 greater than 5, 10, or 50 bp, the number of indels in the recombination
269 regions was much larger than those in other regions (Figure S8).

270



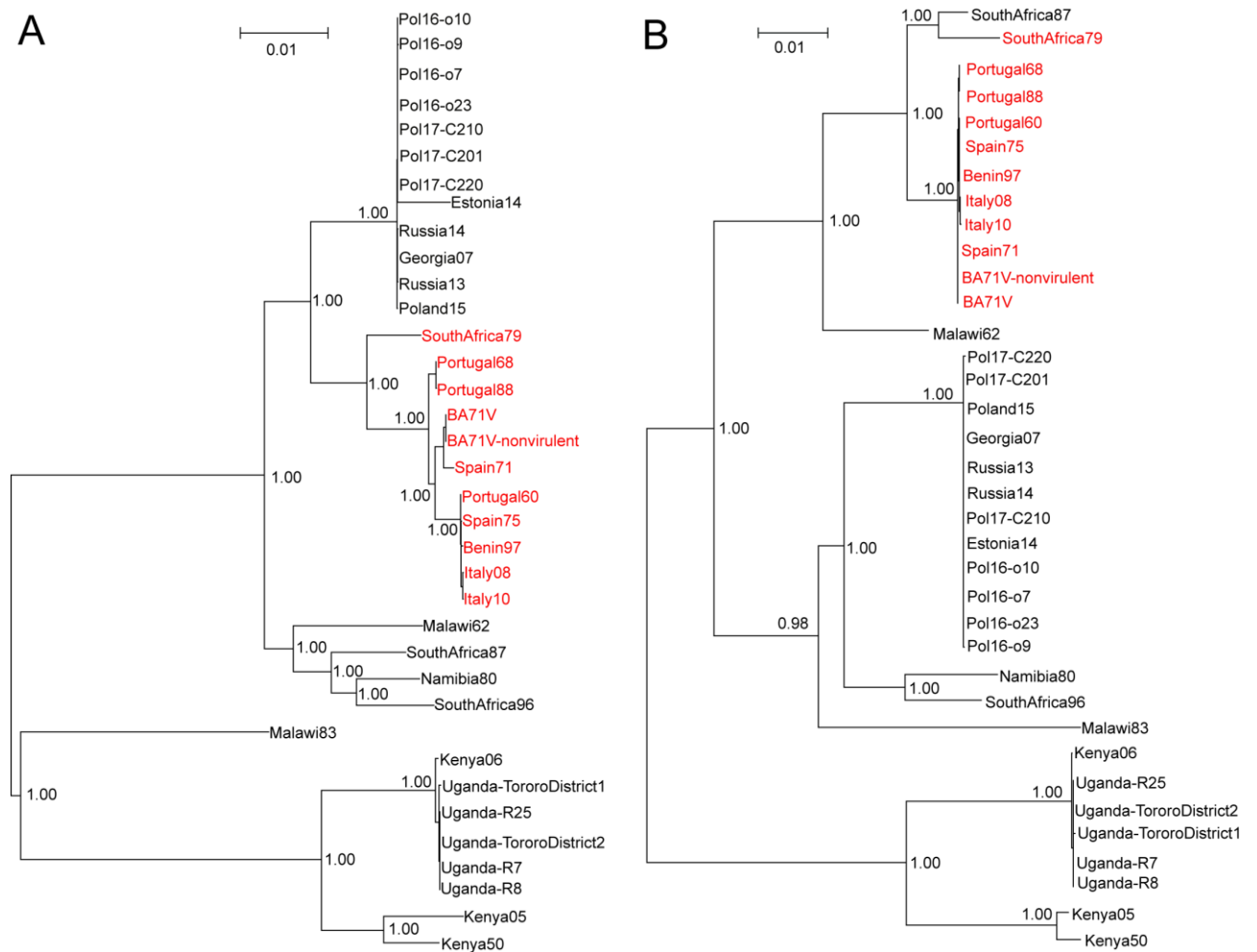
271

272

273 **Figure 3.** Recombination of ASFV genomes. The black areas indicate the
274 recombination region for each genome. The bottom panel shows the ratio
275 of gap in each position of the aligned genome. The panel uses the
276 grayscale color bar at the bottom-left. The red and blue rectangles in the
277 top-right indicate the coding region of pD345L (Lambda-like exonuclease)
278 and P1192R (DNA topoisomerase II), respectively. The black arrow
279 refers to the recombination event displayed in Figure 4.

280

281 Figure 4 illustrates the recombination event in 11 viral isolates (colored in
282 red), including two viral isolates from Africa (SouthAfrica79 and
283 Benin97) and nine viral isolates from Europe. These 11 viral isolates
284 formed a separate lineage in the phylogenetic tree. The recombination
285 region ranged from 133,683 to 142,222 bp, located in the central
286 conservative region of the genome (shown by the black arrow in Figure
287 3). In the phylogenetic tree built with genomic sequences without the
288 recombination regions, the recombinants are the neighbors of a clade
289 containing viruses from Eastern Europe countries (Figure 4A); while in
290 the tree built with genomic sequences of the recombination regions, the
291 recombinants are the descendants of Malawi62 from Africa (Figure 4B).



292

293 **Figure 4.** An example of recombination events in 11 ASFVs (colored in
294 red). Figure (A) refers to the maximum-likelihood phylogenetic tree built
295 with genome sequences without the recombination region
296 (133,683-142,222 in the aligned genome). Figure (B) refers to the
297 phylogenetic tree built with genome sequences of the recombination
298 region. The numbers refer to the bootstrap values of nodes in the
299 bootstrapping test with 100 replicates.

300

301 In addition, the indels introduced directly by the recombination were
302 investigated. The results of comparing the sequence in the recombination

303 regions of recombinants to that in the major parent viruses showed that an
304 average of 37% of differences were caused by the indels in all the
305 recombination events. Further comparison on the proteins showed that in
306 58 of 103 recombination events, there was at least one different protein
307 encoded by the recombination region of the recombinants from that
308 encoded by the major parent viruses.

309

310 Furthermore, the proteins involved in the recombination events were
311 further analyzed. A total of 110 protein groups were involved in the
312 recombination events, including 47 core protein groups and 63 variable
313 protein groups (Table S3). 34 of 110 protein groups belonged to
314 membrane proteins. Four protein groups of “Host cell interactions”
315 (EP153R, A238L, DP96R and DP71L), eight protein groups of “Structure
316 & morphogenesis” and nine protein groups of “Replication &
317 transcription” were involved in the recombination events.

318

319 *5 Identification of possible recombinase and DNA topoisomerase in*

320 *ASFVs*

321 Interestingly, we found a protein, named pD345L, denoted as the
322 Lambda-like exonuclease, was possibly involved in the recombination
323 because it contained the YqaJ-like viral recombinase domain. The protein
324 pD345L has 345 amino acids and is encoded by the minus strand. It was

325 highly conservative in ASFVs, with an average sequence identity of 96.7%
326 between ASFVs. Even higher level of conservation was observed in the
327 recombinase domain of pD345L, with an average sequence identify of
328 97.2%. Although the YqaJ-like recombinase domain was extensively
329 distributed in Bacteria, Virus and Eukaryota, it was considered as the
330 viral origin²⁴. The recombinase domain in ASFV was most similar to that
331 in two giant viruses, Pacmanvirus and Kaumoebavirus (see Materials and
332 Methods) which were possibly distant relatives of ASFVs^{25,26}.

333

334 Besides, topoisomerase was also reported to be related to homologous
335 recombination. We found that a type II DNA topoisomerase, i.e., P1192R,
336 exists in all analyzed ASFV isolates. P1192R was a protein including
337 1192 amino acids, and was encoded by the plus strand. P1192R is also
338 conservative in these viruses, with an average sequence identity of 98.3%
339 between ASFVs. Although pD345L and P1192R were encoded by
340 different strands, they were encoded by the genomic sequences in
341 adjacent regions: the former was encoded by the sequences in the
342 positions of 162,236-163,273 (colored in red in Figure 3), while the latter
343 was encoded by the sequences in the positions of 164,743-168,319
344 (colored in blue in Figure 3). Both enzymes were not involved in any
345 recombination events. The phylogenetic trees for proteins pD345L and
346 P1192R were similar to the tree built with the whole genome (Figure S9).

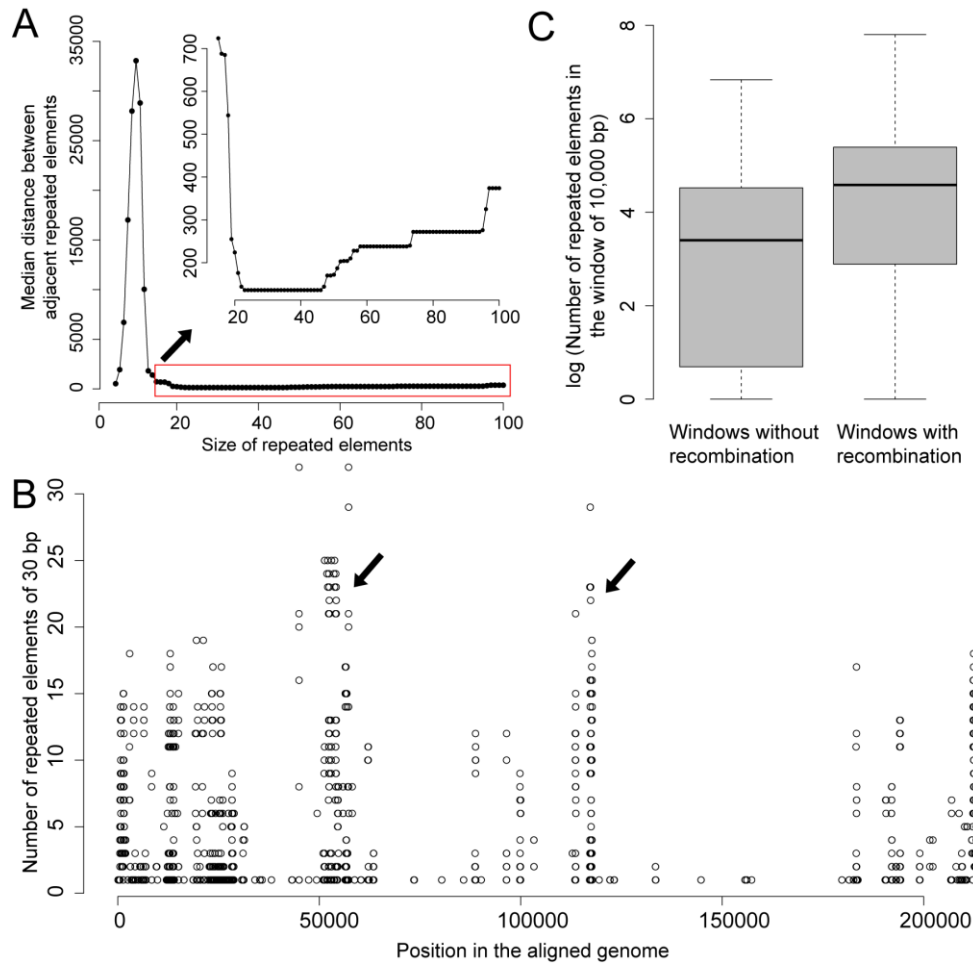
347

348 ***6 An abundance of repeated elements in ASFV genomes***

349 Repeated elements could facilitate the homologous recombination. In this
350 study, lots of repeated elements ranging from 5-100 bp were identified,
351 and then the distribution of the repeated elements in the ASFV genomes
352 was analyzed. As shown in Figure S10, the number of repeated elements
353 in ASFV genomes decreased monotonously as the size of elements
354 increased. Then, the distances between adjacent elements for a given
355 repeated element was investigated (Figure 5A). As the size of the
356 elements increased from 5 to 10, the average distance between the
357 adjacent elements also increased because the number of repeated
358 elements in the genome decreased. Interestingly, the average distance
359 decreased as the size of the elements increased from 11 to 23; it reached
360 to the minimum (136 bp) when the size was 23; then the distance kept
361 unchanged as the size increased from 23 to 46; finally, it increased as the
362 size of repeated element increased from 47 to 100. It should be noted that
363 the average distance was still less than 400 bp even for the repeated
364 elements of 100 bp. These phenomena suggested that the repeated
365 elements of 11 bp or larger tended to cluster in the genome, especially for
366 those of 23-46 bp.

367 For example, when the size of elements was 30 bp, each genome had a
368 median of 427 types of elements which repeated at least two times in the

369 genome. Some elements appeared for over ten times in the genome, such
370 as the element “AGGCGTTAAACATTAAAATTATTACTACTG” in
371 the viral strain BA71V. The region covered by repeated elements
372 accounted for 1%-3% of the genome in ASFVs. The distance between
373 repeated elements was analyzed and demonstrated to have a median
374 distance of 136 bp, suggesting they tend to cluster in adjacent regions.
375 Figure 5B shows the distribution of repeated elements in the aligned
376 genome. Most repeated elements were located at both ends of the genome.
377 Besides, there were two clusters of repeated elements in the positions of
378 around 55,000 bp and 120,000 bp (marked by black arrows), respectively.
379
380 Finally, the contribution of repeated elements to the recombination was
381 investigated. For elements of 10 or more nucleotides, the number of
382 repeated elements in the windows (2000-10,000bp in length) including
383 the recombination was significantly larger than those without the
384 recombination (Table S6). Figure 5C shows the comparison of the
385 number of repeated elements (15 bp in length) in the windows of 10,000
386 bp with and without the recombination in viral genomes. The windows
387 including the recombination had a mean of 194 repeated elements, which
388 was twice of that in the windows without the recombination.



389

390 **Figure 5.** Distribution of the repeated elements. (A) The median distance
391 between adjacent repeated elements versus the size of repeated elements.
392 (B) Number of the repeated elements with the size of 30 bp observed in
393 each genomic position. (C) Comparison of the number of repeated
394 elements (15 bp in length) in the window of 10,000 bp with and without
395 the recombination in viral genomes. For clarity, the natural logarithm of
396 the number of repeated elements was used.

397

398 Discussion

399 This work systematically analyzed the genetic diversity of ASFVs. The

400 large genome of the virus enabled the encoding of an abundance of
401 proteins. Each of the functions of the virus in its life circle could be
402 accomplished by multiple proteins. For example, 35, 19, and 7 protein
403 groups were involved in DNA replication and transcription, structure and
404 morphogenesis, and host cell interactions, respectively. On one hand, this
405 multiple protein mechanism could facilitate the efficient control of the
406 host cell by protein-protein interactions, such as inhibiting the
407 transcriptional activation of host immunomodulatory by A238L³ and
408 inhibiting Toll-like receptor 3 signaling pathways by I329L²⁷; on the
409 other hand, this mechanism could facilitate the precise regulation of the
410 viral activities. For example, the ASFV virus was considered to contain
411 all the enzymes and factors which were required for the transcription and
412 post-treatment of mRNAs³.

413

414 Significant differences were observed among the different proteomes of
415 ASFVs, which may be caused by the following two reasons: i) over 40%
416 of the proteins were non-essential among ASFVs, and ASFVs may have
417 variable number of these proteins; ii) there were 31 genes with
418 replications in ASFV genomes. Diverse proteome among ASFVs may
419 lead to diverse phenotype, such as the diversity in antigen and virulence.
420 The diversity may result in a great challenge for the prevention and
421 control of the virus. For example, the viruses with diverse antigens may

422 need multiple types of vaccines because the effectivity of the
423 cross-protection on viruses may be limited.

424

425 Although lots of efforts have been devoted to developing vaccines against
426 ASFVs^{1,10-12}, unfortunately, most of the attempts have been unsuccessful.

427 The failure could be caused by many factors¹², including the absence of
428 neutralizing antibodies, the diverse antigen-related proteins, the

429 complexity of neutralization, etc. In this study, a total of 65 membrane
430 proteins have been identified. Over eighty percent (80%) of the

431 membrane proteins had unknown functions, many of which may

432 contribute to the antigenic diversity of the virus. Previous studies showed

433 that immunized pigs with the baculovirus expressed hemagglutinin of

434 ASFV were protected against the viral lethal infection²⁸. The results

435 suggested that incorporation of multiple antigens in the vaccine may

436 provide better protection. Therefore, much more efforts are needed to

437 determine the role of membrane proteins in stimulating neutralizing

438 antibodies, and to investigate the neutralization mechanisms and

439 efficiency of the antibodies.

440

441 Indels were found to have larger contribution to the genetic diversity of

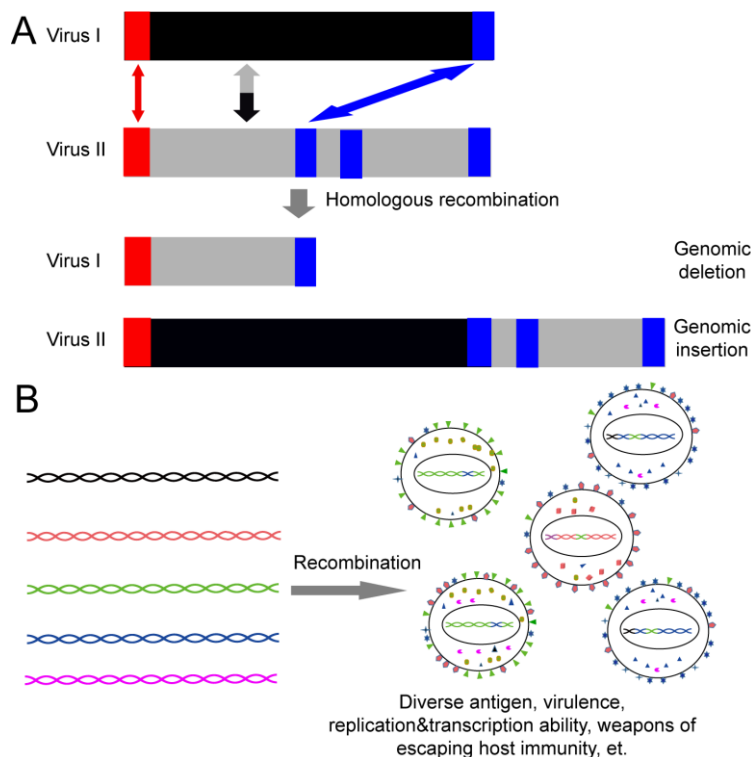
442 ASFVs than the point mutations. Compared to point mutations, indels

443 could introduce a larger variation to the genome, and cause a more severe

444 damage to the genome structures, which may lead to the death of viruses.
445 Therefore, only few indels were observed in viruses with small genomes,
446 such as influenza viruses and hepatitis B viruses (HBV). However, it was
447 more robust for the indels to occur inside the viruses with large genomes,
448 such as ASFVs and poxviruses^{3,29}, because the viruses with large
449 genomes had lots of repeated elements and duplicated proteins (paralogs).
450 Moreover, indels may provide a more efficient way of survival than the
451 point mutations under the natural selection pressure, since the virus with
452 indels could rapidly change its phenotype^{3,29}, such as antigen, virulence,
453 or ability of replication and transcription. For example, the deletion of
454 some MGF genes in ASFV could reduce the viral replication or virulence,
455 which may help with the viral infection of soft ticks^{3,30}.

456
457 Several factors could contribute to the indels and the gene duplications,
458 including replication slippage, retrotransposition, recombination,
459 aneuploidy, polyploidy, etc³¹. The replication slippage may introduce
460 short indels which were widely observed in ASFV genomes, but it is
461 unlikely to cause large indels. This study has demonstrated that the
462 ectopic homologous recombination³², during which the segments with
463 unequal length were exchanged (Figure 6A), may contribute much to the
464 extensive indels observed in ASFV genomes. As a proof, significant
465 associations were observed between the occurrence of extensive

466 recombination events and the indels. Two factors may facilitate the
467 homologous recombination in ASFVs: firstly, a large amount of clustered
468 repeated elements were observed in ASFV genomes (Figure 5); secondly,
469 all the analyzed ASFVs in this study contained a possible recombinase
470 and DNA topoisomerase, both of which were commonly observed
471 enzymes responsible for homologous recombination. Both of enzymes
472 were very conservative and experienced no recombination, suggesting
473 their important roles in ASFVs. Taken together, the homologous
474 recombination should be the effective strategy of ASFVs to generate the
475 genetic diversity, which further leads to the diverse phenotypes, including
476 antigen, virulence, replication and transcription ability, and the “weapons”
477 of escaping from the host immunity (Figure 6B).



478

479 **Figure 6.** Homologous recombination leads to (A) the indels, and (B) the

480 genetic diversity of ASFVs.

481

482 There were some limitations to this study. Firstly, the number of ASFV
483 genomes was limited, which hindered a comprehensive analysis on the
484 evolution of ASFV genomes. Fortunately, the isolates included in this
485 study covered a long time period from 1950 to 2017, and also covered a
486 large area including Africa and Europe, which were the two major areas
487 of the ASFV circulation. Thus the results based on these isolates reflect
488 the genetic diversity of the ASFVs to a large extent. Secondly, the
489 location and size of the indels observed in ASFV genomes may be
490 affected by the sequence alignment algorithm. Two common methods for
491 the alignment of ASFV genomes were used in this study. In both methods,
492 frequent indels were observed, and the indels were demonstrated to be
493 more responsible for the genetic diversity than the point mutations.
494 Thirdly, the proteome of each ASFV was inferred by computational
495 methods. All of the obtained proteins had significant homology to the
496 proteins in the NCBI Protein database, however, most of the proteins in
497 the NCBI Protein database were only predicted without experimental
498 validations. Besides, functions of nearly 70% of ASFV proteins were
499 unknown. Further experimental studies were needed to determine the
500 proteome and the functions in ASFVs^{14,15}. Lastly, the extensively
501 repeated elements in ASFV genomes could facilitate the frequent

502 occurrence of recombination events. However, some of recombination
503 events cannot be detected by the recombination detection method because
504 of the exchange between the genomic segments with small indels. Such
505 kinds of recombination events are difficult to detect. Increasing the
506 sensitivity of the recombination detection method can help detect them,
507 but may also bring false positives. Therefore, the sensitivity and
508 specificity should be balanced in the recombination detection methods.

509

510 Overall, this work provided a systematic view of the genetic diversity of
511 ASFVs. Extensive homologous recombination detected in this study may
512 contribute much to the widespread indels observed in ASFV genomes,
513 which further lead to the large genetic diversity of ASFVs. The results on
514 the causes of the diversity of ASFVs would help with the understanding
515 of the evolution of the virus and thus facilitate the prevention and control
516 of ASFVs.

517

518 **Materials and Methods**

519 *1 ASFV genome and alignment*

520 All the ASFV genomic sequences with over 160, 000 bp were obtained
521 from NCBI GenBank database on October 7, 2018³³. After removing the
522 genomic sequences derived from a patent, a total of 36 ASFV genomes
523 were kept in the analysis. The genomic sequences were aligned by

524 MAFFT (version 7.127b)³⁴. To ensure the robustness of the alignment,
525 the traditional tool of CLUSTAL (version 2.1)³⁵ was also used to align
526 these genome sequences.

527

528 ***2 ORF prediction***

529 To obtain the proteins encoded by the ASFV genomes, each genome
530 sequence was searched against all the ASFV protein sequences obtained
531 from the NCBI protein database on October 7, 2018, with the help of
532 blastx³⁶. All genomic regions with significant hits (e-value < 0.001) were
533 checked using a Perl script: overlapping regions in the same coding frame
534 were merged to obtain open reading frames (ORFs) as long as possible;
535 regions without start codon or stop codon were extended upstream or
536 downstream to search for the start or stop codon. Then, the genomic
537 regions which had i) significant hit, ii) both sequence identity and query
538 coverage percentage greater than 60%, iii) both start and stop codons, and
539 iv) over 120 bps, were defined as the candidate ORFs. The candidate
540 ORFs were then translated into proteins using a Perl script. The proteins,
541 which were either completely embedded within another protein, or
542 contained less than 40 amino acids due to early termination of translation,
543 were removed.

544

545 ***3 Protein grouping***

546 All the inferred proteins of ASFVs were grouped based on sequence
547 homology using OrthoFinder (version 2.2.7)³⁷ with the default
548 parameters. Manual check was conducted to ensure that each protein
549 group contains one type of protein.

550

551 ***4 Calculation of the ratio of dN/dS for proteins***

552 The coding sequences of proteins in each protein group were aligned by
553 codon according to the protein sequence alignment using a Perl script.
554 The ratios of dN/dS between pairwise coding sequences were calculated
555 by yn00 in PAML (version 4.1)³⁸. The average of pairwise dN/dS ratios
556 was calculated as the ratio of dN/dS for the protein.

557

558 ***5 Alignment of ASFV proteome***

559 An ASFV proteome was defined as all the proteins encoded by the ASFV
560 genome. Because both the plus and minus strands could encode proteins,
561 the proteins in a proteome were separated into plus and minus proteome
562 based on the coding strands. Proteome alignment was conducted
563 separately for the plus and the minus proteomes. Firstly, proteins in each
564 proteome were sorted with the order from the 5' end to the 3' end of the
565 genome, based on the coding regions of the proteins. Then, the proteomes
566 were aligned using a dynamic programming algorithm. Manual check
567 was conducted to ensure that there was no mismatch of proteins in the

568 alignment.

569

570 ***6 Function inference and classification of ASFV proteins***

571 The name of each protein group was obtained from the names of BLAST
572 best hit of proteins included in the protein group. To infer the function of
573 each protein group, the longest protein sequence in each protein group
574 was selected as the representative of the protein group. InterproScan
575 (version 5)³⁹ was used to infer the function of the representative protein
576 sequence. The TMHMM Server (version 2.0)⁴⁰ was used to predict
577 whether the representative protein had a trans-membrane helix.

578 Membrane proteins were defined as those who had at least one
579 trans-membrane helices. The functional classification of the proteins was
580 adapted from Dixon's³ and Alejo's¹⁵ work.

581

582 ***7 Detection of homologous recombination events***

583 RDP (version 4)⁴¹ was used to detect the recombination events in the
584 aligned ASFV genomes. Multiple methods in RDP were used. Only the
585 recombination events which were detected by at least two methods were
586 used for further analysis.

587

588 ***8 Evolutionary analysis of YqaJ-like viral recombinase domain***

589 All viral protein sequences of the family of Yqaj (YqaJ-like viral

590 recombinase domain, ID: PF09588) were downloaded from the Pfam
591 database⁴² on November 21, 2018. With the help of blastp, the
592 recombinase domain in ASFV was found to be most similar to that in two
593 giant viruses, Pacmanvirus and Kaumoebavirus, with sequence identities
594 equal to 0.35 and 0.30, respectively.

595

596 ***9 Searching for retrotransposon in ASFV genomes***

597 All retrotransposons in the databases of RepBase (Version 23.10)⁴³ and
598 TREP⁴⁴ were downloaded on November 11, 2018. All ASFV genomes
599 were searched against these retrotransposons using blastn. No hits were
600 obtained under the e-value cutoff of 0.001.

601

602 ***10 Phylogenetic tree inference and visualization***

603 Maximum-likelihood phylogenetic trees were inferred using MEGA
604 (version 5.0)⁴⁵ with the default values of parameters. Bootstrap analysis
605 was conducted with 100 replicates. The phylogenetic tree was visualized
606 using Denscrope (version 2.4)⁴⁶.

607

608 ***11 Statistics analysis***

609 All the statistical analyses were conducted in R (version 3.2.5)⁴⁷.

610

611

612 **Acknowledgements**

613 This work was supported by the National Key Plan for Scientific
614 Research and Development of China (2016YFD0500300,
615 2016YFC1200200 and 2017YFD0500104), the National Natural Science
616 Foundation of China (31500126 and 31671371) and the Chinese
617 Academy of Medical Sciences (2016-I2M-1-005).

618

619 **Competing interests**

620 The authors have declared that no competing interests exist.

621

622 **References**

623

- 624 1 Arias, M., Jurado, C., Gallardo, C., Fernandez-Pinero, J. & Sanchez-Vizcaino, J. M. Gaps in
625 African swine fever: Analysis and priorities. *Transboundary and emerging diseases* **65**,
626 235-247, doi:10.1111/tbed.12695 (2018).
- 627 2 Galindo, I. & Alonso, C. African Swine Fever Virus: A Review. *Viruses* **9**,
628 doi:10.3390/v9050103 (2017).
- 629 3 Dixon, L. K., Chapman, D. A. G., Netherton, C. L. & Upton, C. African swine fever virus
630 replication and genomics. *Virus research* **173**, 3-14, doi:10.1016/j.virusres.2012.10.020
631 (2013).
- 632 4 Costard, S., Mur, L., Lubroth, J., Sanchez-Vizcaino, J. M. & Pfeiffer, D. U. Epidemiology of
633 African swine fever virus. *Virus research* **173**, 191-197, doi:10.1016/j.virusres.2012.10.030
634 (2013).
- 635 5 Sanchez-Cordon, P. J., Montoya, M., Reis, A. L. & Dixon, L. K. African swine fever: A
636 re-emerging viral disease threatening the global pig industry. *Vet J* **233**, 41-48,
637 doi:10.1016/j.tvjl.2017.12.025 (2018).
- 638 6 Arzt, J., White, W. R., Thomsen, B. V. & Brown, C. C. Agricultural Diseases on the Move
639 Early in the Third Millennium. *Veterinary pathology* **47**, 15-27,
640 doi:10.1177/0300985809354350 (2010).
- 641 7 World Organization for animal health. *African Swine Fever (ASF) Report N°4: October 5 - 18,*
642 *2018,*
643 <<http://www.oie.int/en/animal-health-in-the-world/information-on-aquatic-and-terrestrial-ani>

- 644 mal-diseases/african-swine-fever/reports-on-asf/> (2018).
- 645 8 Ge, S. Q. *et al.* Molecular Characterization of African Swine Fever Virus, China, 2018. *Emerg*
646 *Infect Dis* **24**, 2131-2133, doi:10.3201/eid2411.181274 (2018).
- 647 9 SS, S. & WR, H. Antibody response to inactivated preparations of African swine fever virus
648 in pigs. *American journal of veterinary research* **28**, 6 (1967).
- 649 10 King, K. *et al.* Protection of European domestic pigs from virulent African isolates of African
650 swine fever virus by experimental immunisation. *Vaccine* **29**, 4593-4600,
651 doi:10.1016/j.vaccine.2011.04.052 (2011).
- 652 11 Reis, A. L. *et al.* Deletion of the African Swine Fever Virus Gene DP148R Does Not Reduce
653 Virus Replication in Culture but Reduces Virus Virulence in Pigs and Induces High Levels of
654 Protection against Challenge. *Journal of virology* **91**, doi:UNSP
655 e01428-1710.1128/JVI.01428-17 (2017).
- 656 12 Escribano, J. M., Galindo, I. & Alonso, C. Antibody-mediated neutralization of African swine
657 fever virus: Myths and facts. *Virus research* **173**, 101-109, doi:10.1016/j.virusres.2012.10.012
658 (2013).
- 659 13 P, G.-P. *et al.* Neutralizing antibodies to different proteins of African swine fever virus inhibit
660 both virus attachment and internalization. *Journal of virology* **70**, 6 (1996).
- 661 14 Kessler, C. *et al.* The intracellular proteome of African swine fever virus. *Scientific reports* **8**,
662 doi:Artn 1471410.1038/S41598-018-32985-Z (2018).
- 663 15 A, A., T, M., M, G. & G, A. A proteomic atlas of the African swine fever virus particle.
664 *Journal of virology*, doi:10.1128/JVI.01293-18 (2018).
- 665 16 Chapman, D. A., Tcherepanov, V., Upton, C. & Dixon, L. K. Comparison of the genome
666 sequences of non-pathogenic and pathogenic African swine fever virus isolates. *The Journal*
667 *of general virology* **89**, 397-408, doi:10.1099/vir.0.83343-0 (2008).
- 668 17 de Villiers, E. P. *et al.* Phylogenomic analysis of 11 complete African swine fever virus
669 genome sequences. *Virology* **400**, 128-136, doi:10.1016/j.virol.2010.01.019 (2010).
- 670 18 Fraczyk, M. *et al.* Evolution of African swine fever virus genes related to evasion of host
671 immune response. *Veterinary microbiology* **193**, 133-144, doi:10.1016/j.vetmic.2016.08.018
672 (2016).
- 673 19 Michaud, V., Randriamparany, T. & Albina, E. Comprehensive phylogenetic reconstructions
674 of African swine fever virus: proposal for a new classification and molecular dating of the
675 virus. *PloS one* **8**, e69662, doi:10.1371/journal.pone.0069662 (2013).
- 676 20 Wang, Y. *et al.* Origin and Possible Genetic Recombination of the Middle East Respiratory
677 Syndrome Coronavirus from the First Imported Case in China: Phylogenetics and Coalescence
678 Analysis. *mBio* **6**, e01280-01215, doi:10.1128/mBio.01280-15 (2015).
- 679 21 Nagy, P. D. & Bujarski, J. J. Homologous RNA recombination in brome mosaic virus:
680 AU-rich sequences decrease the accuracy of crossovers. *Journal of virology* **70**, 415-426
681 (1996).
- 682 22 Roossinck, M. J. Mechanisms of plant virus evolution. *Annual review of phytopathology* **35**,
683 191-209, doi:10.1146/annurev.phyto.35.1.191 (1997).
- 684 23 Wikipedia. *Homologous recombination*,
685 <https://en.wikipedia.org/wiki/Homologous_recombination#In_viruses> (2018).
- 686 24 *YqaJ protein domain*, <https://en.wikipedia.org/wiki/YqaJ_protein_domain> (2018).
- 687 25 Bajrai, L. H. *et al.* Kaumobavirus, a New Virus That Clusters with Faustoviruses and

- 688 Asfarviridae. *Viruses-Basel* **8**, doi:Artn 27810.3390/V8110278 (2016).
- 689 26 Andreani, J. *et al.* Pacmanvirus, a New Giant Icosahedral Virus at the Crossroads between
690 Asfarviridae and Faustoviruses. *Journal of virology* **91**, doi:UNSP
691 e0021210.1128/JVI.00212-17 (2017).
- 692 27 de Oliveira, V. L. *et al.* A novel TLR3 inhibitor encoded by African swine fever virus (ASFV).
693 *Archives of virology* **156**, 597-609, doi:10.1007/s00705-010-0894-7 (2011).
- 694 28 Ruiz-Gonzalvo, F., Rodriguez, F. & Escribano, J. M. Functional and immunological
695 properties of the baculovirus-expressed hemagglutinin of African swine fever virus. *Virology*
696 **218**, 285-289, doi:10.1006/viro.1996.0193 (1996).
- 697 29 Elde, N. C. *et al.* Poxviruses Deploy Genomic Accordions to Adapt Rapidly against Host
698 Antiviral Defenses. *Cell* **150**, 831-841, doi:10.1016/j.cell.2012.05.049 (2012).
- 699 30 TG, B., Z, L., JG, N., DL, R. & L, Z. African swine fever virus multigene family 360 genes
700 affect virus replication and generalization of infection in *Ornithodoros porcinus* ticks. *Journal*
701 *of virology* **78**, 9 (2004).
- 702 31 Zhang, J. Z. Evolution by gene duplication: an update. *Trends in Ecology & Evolution* **18**,
703 292-298, doi:10.1016/S0169-5347(03)00033-8 (2003).
- 704 32 Freitas-Junior, L. H. *et al.* Frequent ectopic recombination of virulence factor genes in
705 telomeric chromosome clusters of *P-falciparum*. *Nature* **407**, 1018-1022, doi:Doi
706 10.1038/35039531 (2000).
- 707 33 Agarwala, R. *et al.* Database resources of the National Center for Biotechnology Information.
708 *Nucleic Acids Res* **44**, D7-D19, doi:10.1093/nar/gkv1290 (2016).
- 709 34 Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7:
710 improvements in performance and usability. *Molecular biology and evolution* **30**, 772-780,
711 doi:10.1093/molbev/mst010 (2013).
- 712 35 Larkin, M. A. *et al.* Clustal W and Clustal X version 2.0. *Bioinformatics* **23**, 2947-2948,
713 doi:10.1093/bioinformatics/btm404 (2007).
- 714 36 Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment
715 search tool. *Journal of molecular biology* **215**, 403-410 (1990).
- 716 37 Emms, D. M. & Kelly, S. OrthoFinder: solving fundamental biases in whole genome
717 comparisons dramatically improves orthogroup inference accuracy. *Genome biology* **16**, 157,
718 doi:10.1186/s13059-015-0721-2 (2015).
- 719 38 Yang, Z. H. PAML 4: Phylogenetic analysis by maximum likelihood. *Molecular biology and*
720 *evolution* **24**, 1586-1591, doi:10.1093/molbev/msm088 (2007).
- 721 39 Quevillon, E. *et al.* InterProScan: protein domains identifier. *Nucleic Acids Res* **33**, W116-120,
722 doi:10.1093/nar/gki442 (2005).
- 723 40 Krogh, A., Larsson, B., von Heijne, G. & Sonnhammer, E. L. Predicting transmembrane
724 protein topology with a hidden Markov model: application to complete genomes. *Journal of*
725 *molecular biology* **305**, 567-580, doi:10.1006/jmbi.2000.4315 (2001).
- 726 41 Martin, D. & Rybicki, E. RDP: detection of recombination amongst aligned sequences.
727 *Bioinformatics* **16**, 562-563 (2000).
- 728 42 Finn, R. D. *et al.* The Pfam protein families database: towards a more sustainable future.
729 *Nucleic Acids Res* **44**, D279-D285, doi:10.1093/nar/gkv1344 (2016).
- 730 43 *Repbase*, <<https://girinst.org/>> (2018).
- 731 44 Wicker, T. *et al.* A unified classification system for eukaryotic transposable elements. *Nature*

- 732 *Reviews Genetics* **8**, 973-982, doi:10.1038/nrg2165 (2007).
- 733 45 Tamura, K., Stecher, G., Peterson, D., Filipski, A. & Kumar, S. MEGA6: Molecular
734 Evolutionary Genetics Analysis version 6.0. *Molecular biology and evolution* **30**, 2725-2729,
735 doi:10.1093/molbev/mst197 (2013).
- 736 46 Huson, D. H. *et al.* Dendroscope: An interactive viewer for large phylogenetic trees. *Bmc*
737 *Bioinformatics* **8**, 460, doi:10.1186/1471-2105-8-460 (2007).
- 738 47 R Core Team, *R: A language and environment for statistical computing*. *R Foundation for*
739 *Statistical Computing, Vienna, Austria*, <<https://www.R-project.org/>> (2018).
- 740