

1

Article

2 **Building a sequence map of the pig pan-genome from multiple *de***
3 ***novo* assemblies and Hi-C data**

4

5 Xiaomeng Tian^{1*}, Ran Li^{1*}, Weiwei Fu^{1*}, Yan Li^{2*}, Xihong Wang¹, Ming Li¹, Duo
6 Du¹, Qianzi Tang², Yudong Cai¹, Yiming Long¹, Yue Zhao¹, Mingzhou Li^{2#}, Yu
7 Jiang^{1#}

8

9 ¹Key Laboratory of Animal Genetics, Breeding and Reproduction of Shaanxi Province, College of
10 Animal Science and Technology, Northwest A&F University, Yangling 712100, China.

11 ²College of Animal Science and Technology, Sichuan Agricultural University, Chengdu 611130,
12 China.

13

14 *These authors contributed equally to this work.

15 #Corresponding author. E-mail: yu.jiang@nwafu.edu.cn (Y.J.) and mingzhou.li@sicau.edu.cn
16 (M.L.).

17 **Abstract**

18 Pigs (*Sus scrofa*) exhibit diverse phenotypes in different breeds shaped by the
19 combined effects of various local adaptation and artificial selection. To
20 comprehensively characterize the genetic diversity of pigs, we construct a pig pan-
21 genome by comparing genome assemblies of 11 representative pig breeds with the
22 reference genome (Sscrofa11.1). Approximately 72.5 Mb non-redundant sequences
23 were identified as pan-sequences which were absent from the Sscrofa11.1. On
24 average, 41.7 kb of spurious heterozygous SNPs per individual are removed and 12.9
25 kb novel SNPs per individual are recovered using pan-genome as the reference for
26 SNP calling, thereby providing enhanced resolution for genetic diversity in pigs.
27 Homolog annotation and analysis using RNA-seq and Hi-C data indicate that these
28 pan-sequences contain protein-coding regions and regulatory elements. These pan-
29 sequences can further improve the interpretation of local 3D structure. The pan-
30 genome as well as the accompanied web-based database will serve as a primary
31 resource for exploration of genetic diversity and promote pig breeding and biomedical
32 research.

33 **Introduction**

34 *Sus scrofa* (i.e., pig or swine) is of enormous agricultural significance and is an
35 attractive biomedical model. Pigs were domesticated from wild boars independently
36 in Anatolia and East Asia approximately 10,000 years ago following long-term gene
37 flow from their local wild counterparts (Larson, et al. 2005; Groenen, et al. 2012;
38 Frantz, et al. 2015). The combined effects of local adaptation and human-driven
39 artificial selection have shaped the genomic diversity of pigs and form the present
40 various phenotypes. However, to date, these variations have been largely interpreted
41 in the context of the annotated representative reference genome by aligning short
42 reads to it. Increasing evidence suggests that a single individual genome is insufficient
43 to capture all genetic diversities within a species since reference genome may have
44 gaps, mis-assigned regions, or unable to provide a repository for all sequences
45 (Golicz, Batley, et al. 2016). Alternatively, comparisons of independently *de novo*-
46 assembled genomes and a reference genome sequence promise a more accurate and
47 comprehensive understanding of genetic variations within a species (Li, et al. 2014;
48 Schatz, et al. 2014).

49 Most recently, the pan-genome, the non-redundant collection of all genomic
50 sequences of a species, has emerged as a fundamental resource for unlocking natural
51 diversity in eukaryotes. Intraspecific comparisons in plants (e.g., soybean (Li, et al.
52 2014), *Brassica oleracea* (Golicz, Bayer, et al. 2016), *Brachypodium distachyon*
53 (Gordon, et al. 2017) and rice (Zhao, et al. 2018)) and in animals (e.g., mosquitoes
54 (Neafsey, et al. 2015), macaques (Yan, et al. 2011) and humans (Li, et al. 2010;
55 Maretty, et al. 2017)) have revealed surprisingly large amounts of variation within a
56 species. To build a high-quality pan-genome, a number of individual genomes are

57 required (Li, et al. 2014; Monat, et al. 2017; Wong, et al. 2018), which remains an
58 obstacle for most mammalian species. The current pig genome (Sscrofa11.1)
59 represents one of the most continuous assemblies in mammalian species
60 **(Supplementary Fig. 1)** and is from one individual (Duroc breed). In addition, our
61 previous studies generated *de novo* assemblies of ten geographically and
62 phenotypically representative pig genomes from Eurasia (Li, et al. 2013; Li, et al.
63 2017). Together with the assembly of Chinese Wuzhishan boar (Fang, et al. 2012), the
64 availability of 12 pig genomes has provided an unprecedented opportunity to
65 investigate their genetic differences at the individual, ethnic/breed or continental
66 level.

67 Here we carried out an in-depth comparison between 11 *de novo* assemblies
68 and the reference genome by analysis of the assembly-versus-assembly alignment.
69 The final pan-genome comprises 39,744 (total length: 72.5 Mb) newly added
70 sequences and of which 607 demonstrate coding potential. Furthermore, the three-
71 dimensional (3D) spatial structure of pan-genome was depicted by revealing the
72 characteristics of pan-genome in A/B compartment (generally euchromatic and
73 heterochromatic regions) and topologically associating domain (TAD). We also build
74 a pig pan-genome database (PIGPAN,
75 <http://animal.nwsuaf.edu.cn/code/index.php/panPig>) which can serve as a
76 fundamental resource for unlocking variations within diverse pig breeds.

77 **Results**

78 **Initial characterization of pan-sequences in the pig genome**

79 To construct the pig pan-genome, we first aligned 11 assemblies from 11 genetically
80 distinct breeds (five from Europe, and six from China) against Sscrofa11.1 using
81 BLASTN to generate the unaligned sequences (**Fig. 1a and Supplementary Fig. 2**).
82 The length of the unaligned sequences in the Chinese pigs was significantly longer
83 than those in the European pigs ($P < 0.01$) since the reference genome is from a
84 European pig (**Fig. 1a**). As expected, the Wuzhishan assembly had the largest number
85 of sequences because this sample is the only male individual among the 11 assemblies
86 and can provide many male-specific sequences (**Fig. 1a and Supplementary Table**
87 **2**). After removing redundant sequences, we obtained 39,744 sequences with a total
88 length of 72.5 Mb (**Fig. 1b**), which were absent from Sscrofa11.1 and thus were
89 defined as pan-sequences. The content of the repetitive elements (45.91%) and GC
90 (44.61%) in these sequences were slightly higher than those in Sscrofa11.1 (45.19%
91 and 41.5%, respectively) (**Fig. 1a and Supplementary Table 3**). Notably, 44% (32
92 Mb) of the pan-sequences can be assigned to a unique assembly, highlighting the
93 limitations of using one single assembly (**Fig. 1b**). All of the pan-sequences were
94 longer than 300 bp. Among them, pan-sequences that are over 5 kb contributed 57%
95 of the total length (**Fig. 1c**).

96 To validate the authenticity of the pan-sequences, we first compared these
97 sequences to ten mammalian genomes and found that the majority (67.5%) of the pan-
98 sequences had homologs in these genomes (E -value $< 1e-5$) (**Fig. 2a and**
99 **Supplementary Table 4**). As expected, Cetacea has the greatest number of best hits
100 in accordance with their close evolutionary relationship with pig (**Fig. 2a**). To explore

101 the potential presence or absence of protein-coding genes within the pan-sequences,
102 we aligned these pan-sequences to Refseq proteins from pig, cattle, goat, human,
103 sheep and sperm whale using TBLASTN (E -value $< 1e-5$). The most significantly
104 overrepresented hits were members of the olfactory receptor family (12.4%), which is
105 the largest gene superfamily in vertebrates (Zhang and Firestein 2002) followed by
106 other highly variable families (typically, O-methyltransferase domain-containing
107 proteins), showing that these additional sequences are more variable than the
108 reference set (**Fig. 2b**). Especially, 18 new and complete olfactory receptor genes
109 were identified (**Supplementary Fig. 3**), implying that our pan-sequences can ensure
110 an enriched repertoire of highly divergent gene families.

111 To explore whether these pan-sequences exhibited population-specific
112 characteristics, we retrieved 87 publicly available pig genomes ($>10\times$ coverage) from
113 China and Europe and aligned them to the whole pan-genome (**Supplementary Table**
114 **5**). The presence of each pan-sequence was determined by calculating their
115 normalized read depth (NRD). The samples were clustered into three distinct groups
116 corresponding to their geographical origin: southern Chinese, northern Chinese and
117 European pigs (**Fig. 2c and Supplementary Table 5**), which was consistent with the
118 previously reported genetic architecture of domestic pigs (Ai, et al. 2015) and showed
119 the presence or absence of pan-sequences can reflect the local adaptation and
120 domestic history of pigs.

121 Furthermore, among the 87 pig genomes, 42 represented males, enabling us to
122 determine the male-specific sequences (**Supplementary Table 5**). We identified
123 10.43 Mb of male-specific sequences that were present in more than 90% of male
124 individuals but absent in females (**Supplementary Table 6**). The distribution pattern

125 of these male-specific pan-sequences revealed an Asian type which is confined to
126 Asia and a Eurasian type which is distributed across Eurasia. This divergent
127 distribution is concordant with the history of the male-biased migration from non-
128 Asian to Asian (Guirao-Rico, et al. 2018) (**Fig. 2d**).

129 To determine the potential genomic positions of these pan-sequences, we aligned
130 these pan-sequences to Sscrofa11.1 using their flanking regions. Only 19.00% (7,554)
131 of the pan-sequences could be anchored in this way (**Supplementary Table 7**). By
132 providing pan-sequences with refined positions, the genomic annotation could be
133 enriched and improved. For instance, a pan-sequence of 14.3 kb containing the
134 complete genic region of *RARRES3* is absent in Sscrofa11.1, which can be acted as a
135 tumour suppressor or growth regulator (Shyu, et al. 2003). We further validated this
136 gene by resequencing and RNA-seq data and found that this gene has high abundance
137 in multiple tissues (FPKM > 1) (**Fig. 2e**). We also found that the expression of this
138 gene is enriched in Chinese pigs, which might be a population-specific gene involved
139 in Chinese pig growth. Another pan-sequence of 18.6 kb included six coding exons of
140 *ZNF622*, which is also missed in Sscrofa11.1 (**Fig. 2e**). *ZNF622* played a role in
141 embryonic development by activating the DNA-bound *MYBL2* transcription factor
142 (Arumemi, et al. 2013). These absent six exons resulted in another spliced transcript
143 isoform, which can be validated using the RNA-seq data.

144 **Constructing a more comprehensive sequence map for genomic and** 145 **transcriptomic analysis**

146 Compared with Sscrofa11.1, the mapping rate of resequencing data in the pan-
147 genome was increased by 0.29-0.43% (**Fig. 3a and Supplementary Table 8**).
148 Meanwhile, the mapping rate of Sscrofa11.1 in the pan-genome was decreased by

149 approximately 1.43%, indicating that many reads had been adjusted to better positions
150 in the pan-sequences accompanied with improved quality (**Fig. 3a, b**). The adjustment
151 of many reads from Sscrofa11.1 to pan-sequences will result in better SNP calling
152 efficacy. An average of 41,729 heterozygous SNPs per sample were depressed and the
153 read depth was also adjusted to the whole-genome average level in these regions
154 where spurious SNPs were removed (**Fig. 3b, c, d, Supplementary Fig. 4 and**
155 **Supplementary Table 9**). Furthermore, 12,888 novel SNPs per individual were
156 recovered using the pan-genome and thus provided enhanced resolution for genetic
157 diversity studies. In addition, the mapping quality and mapping rate of RNA-seq data
158 were also improved based on 92 samples (**Supplementary Figs. 5 and 6**). In total,
159 897 sequences containing 1163 potential transcripts showed appreciable expression
160 ($\text{FPKM} \geq 1$ in at least one sample). To further assess the protein-coding potential of
161 pan-sequences, a total of 607 out of the 897 pan-sequences were predicted to have
162 coding potential by Coding Potential Calculator (CPC) (Kong, et al. 2007). For the
163 gene expression of pan-sequences, more expression differences were found among
164 tissues (Pearson's $r = 0.84$) than within tissues (Pearson's $r = 0.91$), consistent with
165 previous findings (Tang, et al. 2017) (Fig. 3e and Supplementary Fig. 7).

166 **Hi-C based analysis revealing the characteristics of pan-sequences regarding 3D** 167 **structures and their potential function**

168 Adjacent loci generally demonstrated frequent interaction which can be determined by
169 high-throughput Hi-C analysis (Dong, et al. 2017), thus enabling us to anchor these
170 pan-sequences to Sscrofa11.1. Here, we generated 12 Hi-C data from 10 individuals
171 to anchor these pan-sequences to Sscrofa11.1 by inferring their special interactions
172 with adjacent regions (**Supplementary Tables 10 and 11**).

173 To evaluate the robustness and accuracy of Hi-C based localization, we
174 comprehensively investigated the anchored results from five samples digested by the
175 MboI enzyme and another seven samples digested with HindIII enzyme (see
176 methods). The result indicates that the Hi-C-based localization determined by
177 different samples has high consistency (**Supplementary Fig. 8 and 9**). A total of
178 7,554 sequences (19.0%) can be anchored to Sscrofa11.1 by flanking sequences.
179 Using Hi-C based approach, another 3,447 sequences can be further anchored. Based
180 on this result, at 100-kb resolution, we found the pan-sequences are uniformly
181 distributed in A/B compartment which are generally euchromatic and heterochromatic
182 regions, respectively (Dogan and Liu 2018) (**Fig. 4b**). At 20-kb resolution, we found
183 that pan-sequences were more enriched at TAD boundary regions (**Fig. 4c**). Notably,
184 we found that genomic variations (SNPs and CNVs) occurred more frequently at
185 TAD boundary regions than at the TAD interior regions (**Supplementary Figs. 10, 11**
186 **and 12**), indicating that the occurrence of pan-sequences could be associated with
187 genomic variations.

188 Based on the high genome coverage sample (~300×), we identified 201 (5.4%)
189 pan-sequences which were anchored to genomic regions harbouring putative enhancer
190 elements. Furthermore, 47 pan-sequences were shown to contain enhancers by
191 demonstrating enhancer-promoter interactions with high confidence (**Supplementary**
192 **Fig. 13**). These genes which are influenced by remote regulation were significantly
193 enriched in retinol metabolism, olfactory transduction, arachidonic acid metabolism
194 and fatty acid degradation (**Supplementary Table 12**). When the corresponding
195 regions of low interaction were replaced with pan-sequences, we found that 3D spatial
196 structure was greatly improved (**Fig. 4d**). Thus, replacement with pan-sequences will
197 help to depict the 3D structures of the whole genome.

198 **Pig pan-genome database**

199 To facilitate the use of the pig pan-genome by the scientific community, a pig pan-
200 genome database (PIGPAN) was constructed. PIGPAN is a comprehensive repository
201 of integrated genomics, transcriptomics and regulatory data. The system diagram is
202 shown in **Fig. 5a**. In our local UCSC Genome Browser server (Gbrowse), 17 tracks
203 were released against the pig pan-genome (**Fig. 5b**). Users can search the database
204 using a gene symbol or chromosome location to obtain results in terms of four
205 aspects: (i) the reference genome and pan-sequence annotations, (ii) the gene
206 expression in 20 corresponding tissues, seven types of regulation signals
207 (**Supplementary Table 13**) and the conserved elements of a 20-way mammalian
208 alignment, (iii) the chromosome localization of pan-sequences, and (iv) the haploid
209 copy number of 87 pigs. We also provided basic search functions to retrieve basic
210 gene information, GO annotation and KEGG pathways. Here, we present one case
211 using PIGPAN showing the copy number difference of KIT between European and
212 Chinese pigs (**Fig. 5c**). Moreover, users can download data from
213 <http://animal.nwsuaf.edu.cn/panPig/download.php>. As the functions and associated
214 traits of more genes in the pig genome are determined in the future, our browser will
215 be updated regularly to meet the various needs of the scientific community.

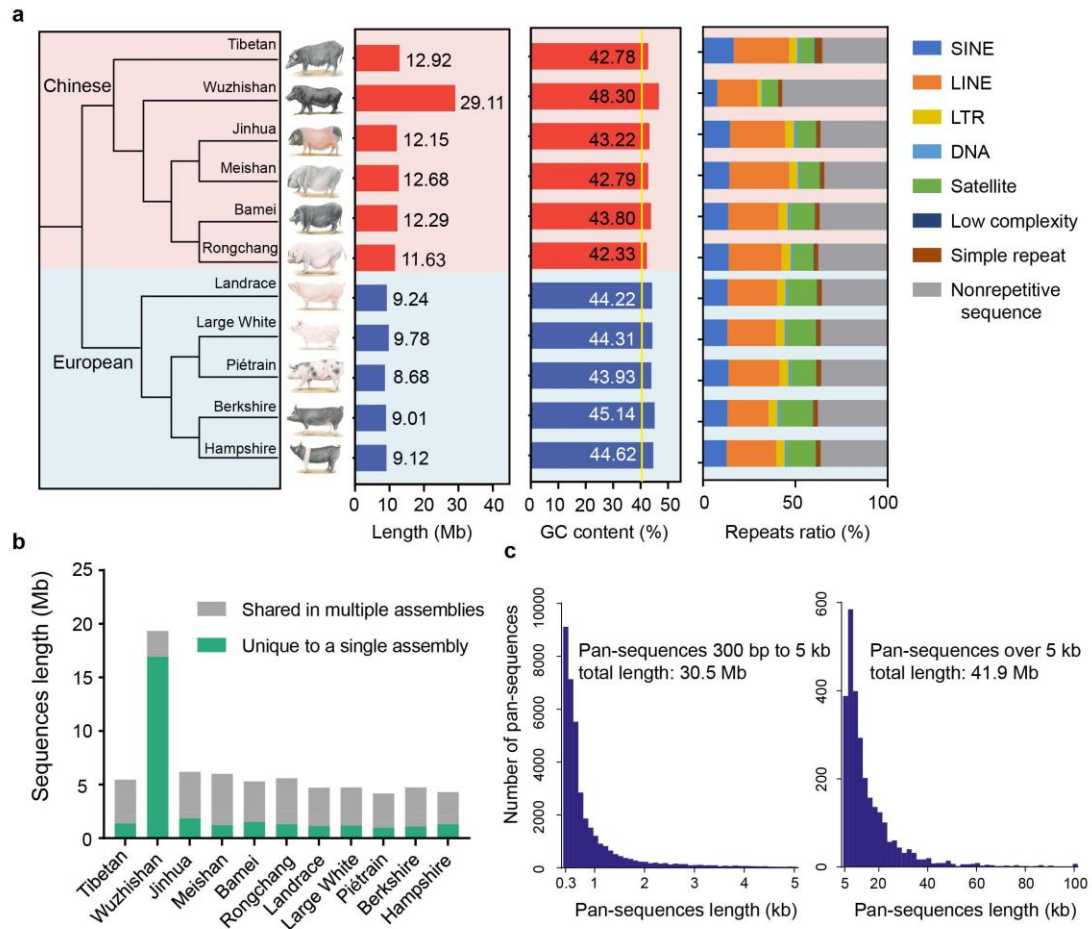
216 **Discussion**

217 In this study, we utilized 12 independent *de novo* assemblies (Fang, et al. 2012; Li, et
218 al. 2013; Li, et al. 2017) and a large amount of whole-genome resequencing data to
219 build a sequence map of the pig pan-genome. The *de novo* assemblies in the present
220 study cover a wide range of diverse breeds across Eurasia and thus ensure a
221 comprehensive discovery of the missing pan-sequences. These pan-sequences as well
222 as the accompanied genomic variation and expression information will be a valuable
223 resource for depicting the complete genetic makeup of porcine phenotypic and
224 genomic diversity.

225 With the rapid decrease in the cost of generating high-quality *de novo* assemblies,
226 the pan-genome strategy is becoming increasingly affordable and will soon become
227 applicable for many other animal species. The importance of pan-genomes has been
228 widely accepted in the field of plant genomics (Hirsch, et al. 2014; Schatz, et al. 2014;
229 Golicz, Bayer, et al. 2016; Sun, et al. 2017; Zhao, et al. 2018). The high genomic
230 plasticity of plants can result in the complete gain/loss of a large number of genes
231 within a species (Golicz, Bayer, et al. 2016; Gordon, et al. 2017; Zhao, et al. 2018). In
232 contrast, animal genomes are much more conserved and have longer genes with
233 complex splicing events, which means that, generally, only intergenic or fragmented
234 genic regions are involved in the gain/loss of genomic sequences in animals.
235 Nonetheless, this difference does not mean that animal pan-genomes are less
236 important. For instance, the pan-sequences that we recovered demonstrated
237 population-specific patterns, indicating that they are potentially associated with
238 adaptations to various environmental conditions (Li, et al. 2010; Gordon, et al. 2017).
239 Furthermore, our research suggests that these pan-sequences may act as enhancers of

240 some genes that regulate metabolic activity in different breeds. We also found a large
241 number of SNPs residing in the pan-sequences, which can lead to an accurate
242 assessment of true variations, thereby providing enhanced resolution for genetic
243 diversity of different pig populations. The enriched genomic sequence repertoire can
244 help in identifying causal mutations that were previously unrecognized by linkage,
245 association and copy-number-variation studies.

246 In conclusion, our study has shown that the pan-genome, when used as a
247 reference, can ensure a more comprehensive repertoire of genomic variations and can
248 facilitate downstream genomic, transcriptomic and even 3D genome analyses.
249 Therefore, we highlight the transition from the current reference genome to the pan-
250 genome.



251

252 **Fig. 1 Construction of the pig pan-genome and the characterization of pan-sequences. a**

253 Maximum likelihood phylogenetic tree, sequence length, GC content and repeat composition

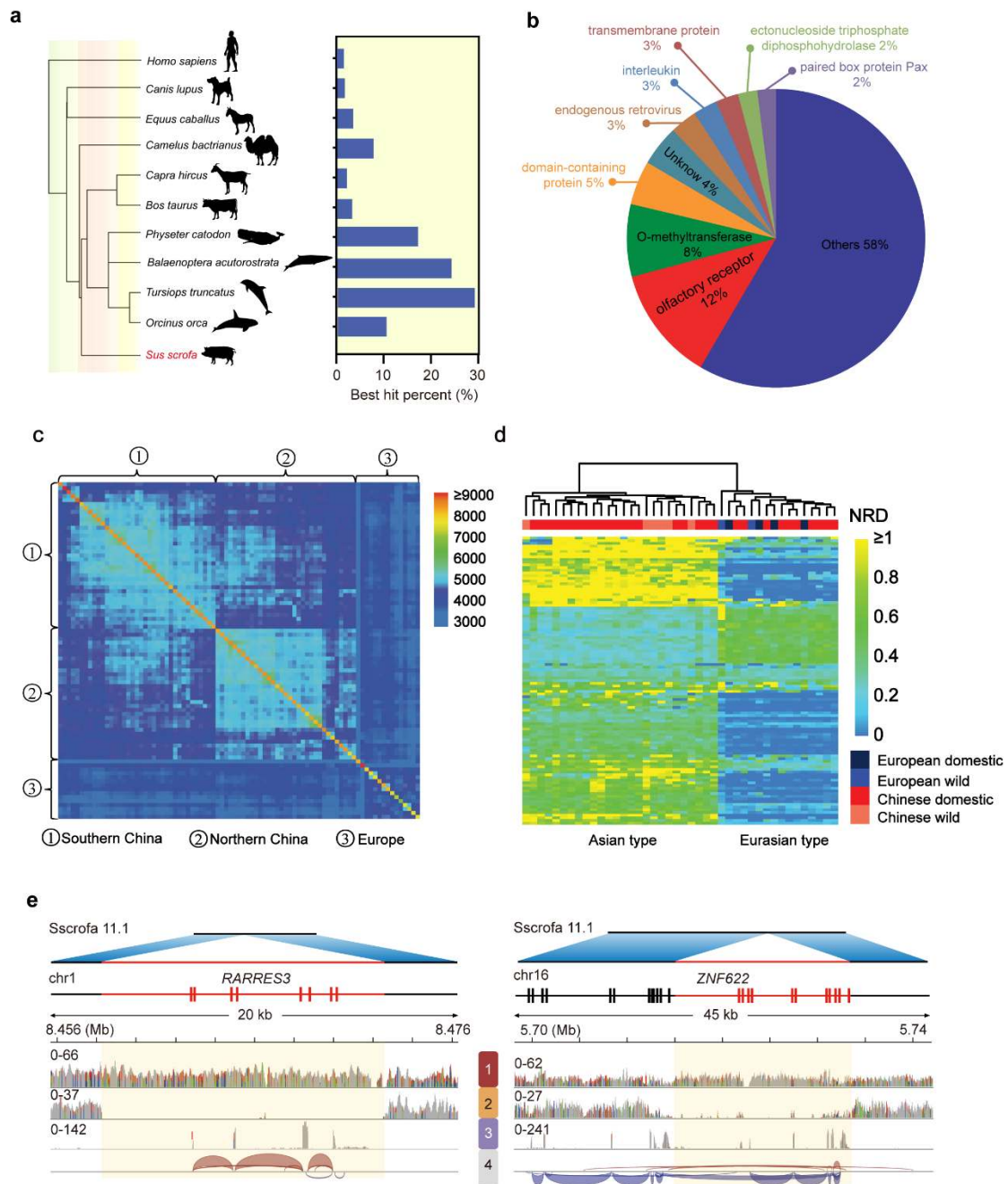
254 of missing sequences identified in each individual assembly of eleven breeds (left to right). **b**

255 The total sequence length and breed-specific sequence length of each breed for non-redundant

256 pan-sequences. **c** Length distribution of all pan-sequences. (Wuzhishan pigs had the largest

257 number of sequences because this individual is the only male among all the 11 assemblies and

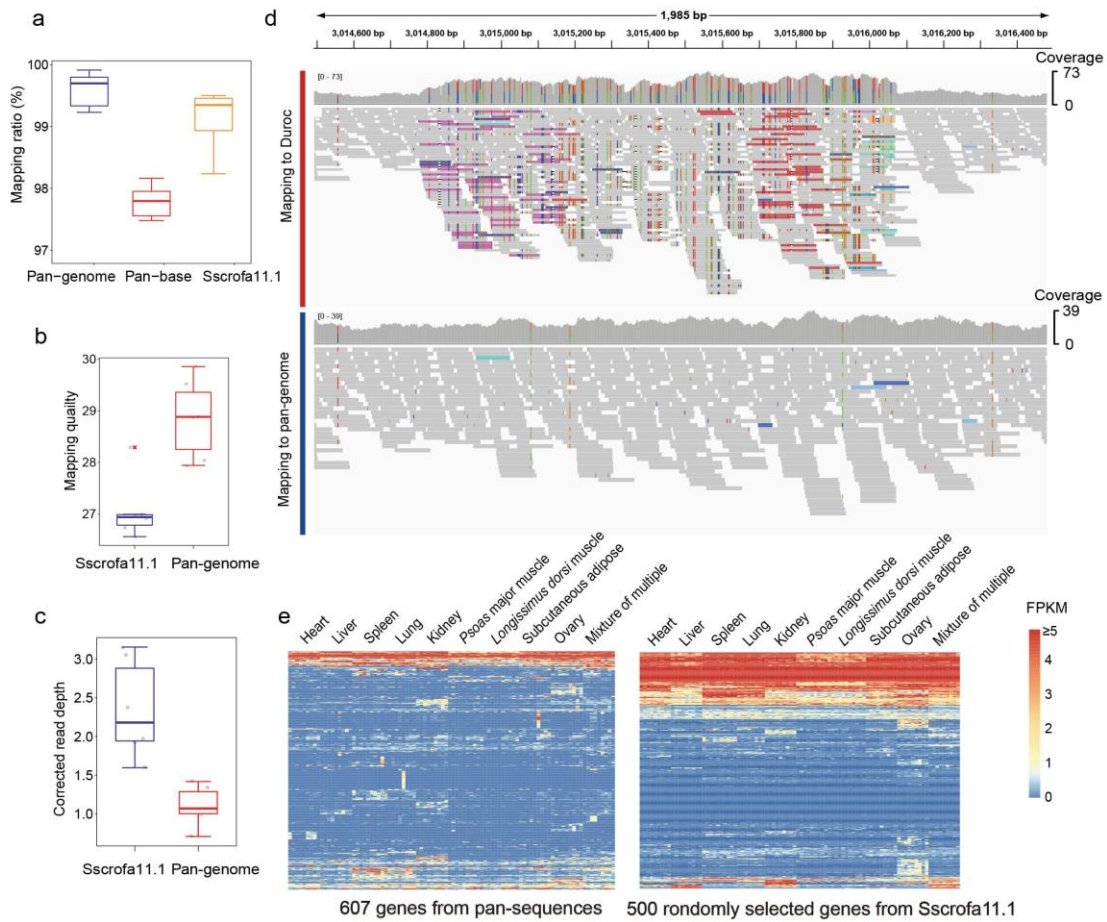
258 the sequencing platform of this individual differed from that used for other samples).



259

260 **Fig 2. Pan-sequences validation and population-specific pattern.** **a** Homolog identification
 261 of pan-sequences in ten mammalian genomes. Only the best hit was retained for each pan-
 262 sequence. **b** Number of hits in pan-sequences to Refseq protein families. **c** An 87×87 matrix
 263 showing the number of shared pan-sequences among all the individuals by pairs. Each cell
 264 represents the number of shared pan-sequences by two individuals. See Supplemental table 8
 265 for the classification of each group. **d** The normalized read depth (NRD) of male-specific pan-
 266 sequences in each male. See Supplemental table 8 for the classification of each group. (Only

267 the sequences with the frequency ranging from 0.5 to 0.9 were shown.) **e** Genes contained in
268 the pan-sequences. One pan-sequence of 14.3 kb harbour the complete genic region of
269 RARRES3, and another covers partial genic regions of ZNF622, representing a new splicing
270 event. The four tracks at the bottom represent the reads mapping of whole-genome
271 resequencing data of two samples (1-2) and the inferred exons as well as their splicing
272 isoforms based on RNA-seq (3-4).



273

274 **Fig. 3 Improvements of genomic and transcriptomic analyses by using the pan-genome.**

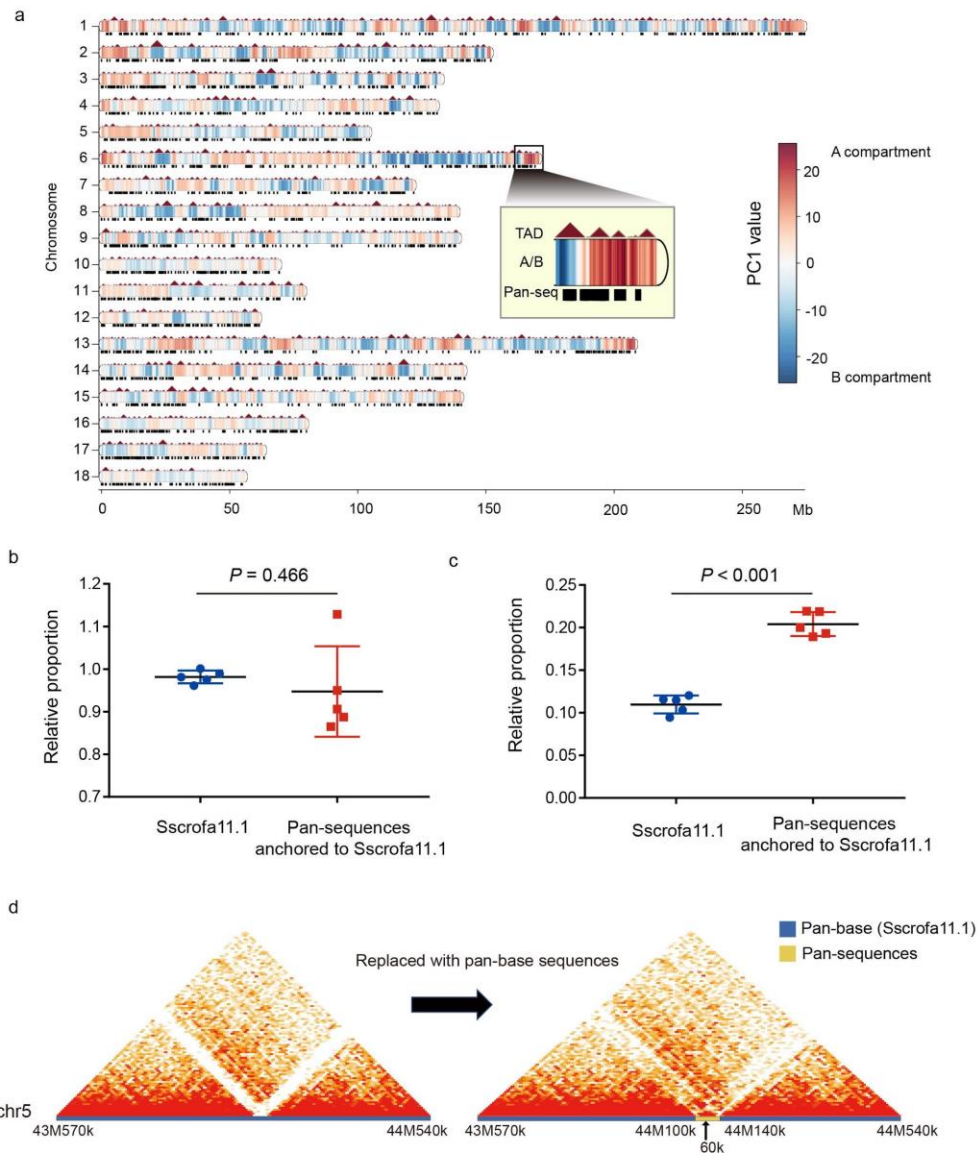
275 **a** Comparison of the mapping ratio of resequencing data using the pan-genome versus

276 Sscrofa11.1. **b** Comparison of read-mapping quality using the pan-genome versus

277 Sscrofa11.1. **c** Comparison of corrected read-mapping depth using the pan-genome versus

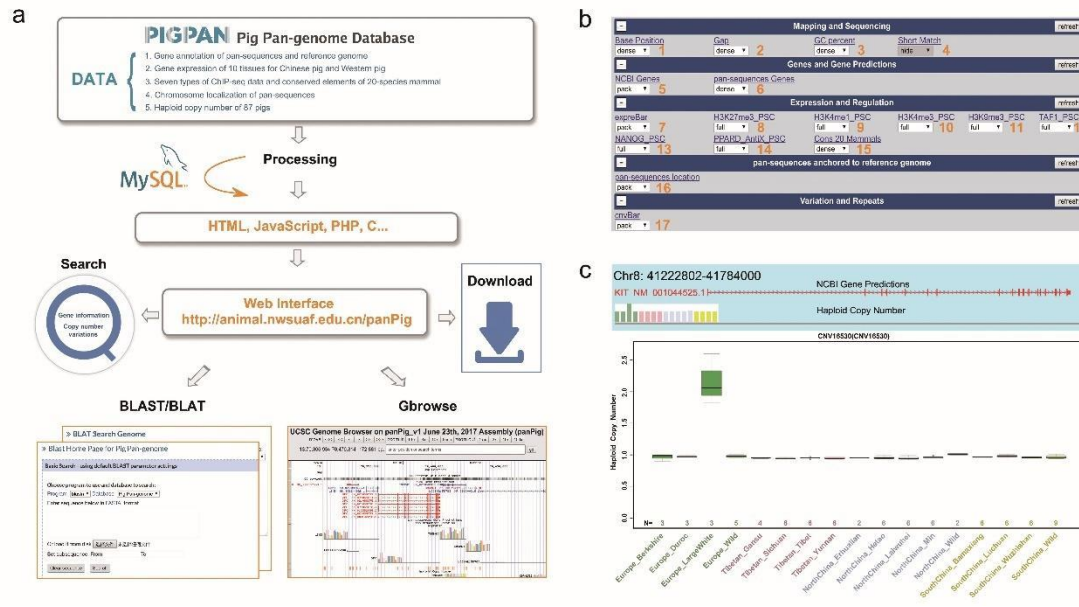
278 Sscrofa11.1. **d** Improved read mapping using the pan-genome versus Sscrofa11.1 as viewed

279 with IGV. **e** Transcriptional potential of the pan-sequences.



280

281 **Fig. 4 The 3D spatial structure of the pan-genome. a** The distributions of the A/B
 282 compartment, TAD and anchored pan-sequences. **b** The relative proportion of A compartment
 283 over B compartment in length in the pig genome (left), and the relative proportion of pan-
 284 sequences located in A compartment over those located in B compartment in length (right).**c**
 285 The relative proportion of TAD boundary regions over TAD interior regions in length in
 286 Sscrofa11.1 (left) and the relative proportion of pan-sequences located in TAD boundary
 287 regions over TAD interior regions in length (right). **d** An example of improving a 3D spatial
 288 structure after replacing the weakly interacting sequences with the non-reference pan-
 289 sequences.



290

291 **Fig. 5. The processing pipeline used to construct PIGPAN.** PIGPAN integrated genomics,
292 transcriptomics and regulatory data. Users can search a gene symbol or a genomic region to
293 obtain results in the form of an interactive table and graph.

294 **Acknowledgements**

295 This work was supported by research grants from the National Natural Science
296 Foundation of China (No. 31572381) to Y.J and the Science & Technology Support
297 Program of Sichuan (2016NYZ0042 and 2017NZDZX0002) to M.Z.L. We thank the
298 High Performance Computing platform of Northwest A&F University for their
299 assistance with the computing.

300 **Author contributions**

301 Y.J. and M.Z. L. conceived the project and designed the research. X.T., Y.L and M.L.
302 analysed the Hi-C data. X.T., R.L., W.F., M.L. and D.D. performed the analysis. X.T.,
303 R.L. and W.F. wrote the manuscript. Y.J, M.Z.L., X.W. revised the manuscript.

304 **Competing interests**

305 The authors declare that no competing interests exist.

306 **Methods**

307 **Construction of the pan-genome**

308 We downloaded the publicly available pig genome assemblies of ten female and one
309 male individuals from 11 diverse breeds (five originated in Europe and six originated
310 in China) (**Supplementary Fig. 2 and Table 2**) (Fang, et al. 2012; Li, et al. 2013; Li,
311 et al. 2017). To identify the sequences which cannot align to the reference genome,
312 we split the 11 assemblies by gap region and iteratively aligned them to the reference
313 pig genome assembly (Sscrofa1.1) using the BLASTN (Camacho, et al. 2009).
314 Sscrofa1.1 was masked by WindowMasker (Morgulis, et al. 2006) before alignment
315 to speed up the alignment process. The sequences with $<90\%$ identity and ≥ 300 bp in
316 length were retained. After that, these low-identity sequences were aligned to each
317 other to remove redundancy. Finally, a non-redundant set of 72.5 Mb of sequences
318 from 11 assemblies was obtained; these sequences were defined as pan-sequences.
319 The Duroc genome (Sscrofa1.1) plus these 72.5 Mb pan-sequences made up the pan-
320 genome.

321 **Determining the characteristics of the pan-sequences**

322 To explore whether the pan-sequences have homologous regions across species and
323 are potential to be functional, we aligned these sequences to ten mammalian reference
324 genomes (i.e., *Homo sapiens*, *Camelus bactrianus*, *Equus caballus*, *Canis lupus*
325 *familiaris*, *Capra hircus*, *Bos Taurus*, *Orcinus orca*, *Physeter catodon*, *Balaenoptera*
326 *acutorostrata scammoni*, *Tursiops truncates*) to search for any matches (E -value $\leq 1e$ -
327 5) using BLASTN (Camacho, et al. 2009) (**Supplementary Table 4**). Only the best
328 hit was remained for each query.

329 To validate the authenticity of these pan-sequences and identify assembly-
330 specific sequences, we aligned all of them to the each of the 11 *de novo* pig
331 assemblies to search for any matches ($\geq 90\%$ coverage and $\geq 95\%$ identity) using
332 BLASTN (Camacho, et al. 2009). If the sequence of an assembly does not have a high
333 similarity with other assemblies, this sequence is considered as the assembly-specific
334 sequences.

335 **Population-based resequencing and CNV calling**

336 We downloaded the whole genome resequencing data for 71 domestic pigs and 16
337 wild boars for population analysis of pan-sequences. The sequences data were
338 retrieved from NCBI under the Bioproject PRJNA213179, PRJNA281548,
339 PRJNA309108 and PRJEB9922 (Ai, et al. 2015; Frantz, et al. 2015; Jeong, et al.
340 2015; Li, et al. 2017) (**Supplementary Table5**). After alignment using BWA (version
341 0.7.15-r1140) (Li and Durbin 2009) with default parameters, we used CNVcaller
342 (Wang, et al. 2017) to calculate the normalized read depth (NRD) of each sequences.
343 The presence and absence of each pan-sequence were then determined by NRD.

344 **ChIP-seq short-read alignment and peak calling**

345 To confirm the content of regulatory elements in pan-sequences, we downloaded
346 seven ChIP-seq data from NCBI Bioproject PRJNA152995, including H3K27me3,
347 H3K4me1, H3K4me3, H3K9me3, NANOG, PPARD AntiX and TAF1 signals (**table**
348 **S13**) (Xiao, et al. 2012). Sequencing reads were aligned to pig pan-genome using
349 BWA (version 0.7.17-r1188) (Li and Durbin 2009) with default parameters. Low-
350 quality and multiple-mapping reads were removed using SAMtools (Li, et al. 2009)
351 with option “-q 20”. Enriched regions (or peaks) were called ($p < 1e-5$; no filtering on
352 fold enrichment or FDR correction) using MACS (version 2.1.1) (Zhang, et al. 2008)

353 with total DNA input as control.

354 **Identification of male-specific sequences**

355 There are 42 males and 45 females in our whole genome resequencing data
356 (Supplementary Table 8). We compared the normalized read depth (NRD) between
357 females (NRD < 0.1, sample size = 42) and males (0.2 < NRD < 0.7, sample size =
358 45) to identify the putative male-specific pan-sequences. Thus, we identified 1,638
359 male-specific scaffolds (**Table S9**) which were present in most all of male individuals
360 (frequency $\geq 50\%$) but absent in females (frequency = 0) with a combined length of
361 10,432,972bp (**Supplementary Table 6**).

362 **Gene annotation and functional enrichment analysis**

363 Homology-based and *de novo* prediction were used to annotate protein-coding genes.
364 For homology-based prediction, pan-sequences were aligned onto the repeat-masked
365 assembly using TblastN (Camacho, et al. 2009) with an *E*-value cutoff of 1e-5.
366 Aligned sequences as well as corresponding query proteins were then filtered and
367 passed to GeneWise to search for accurately spliced alignments (Doerks, et al. 2002).
368 For *de novo* prediction, GenScan (Burge and Karlin 1998), Augustus(Stanke, et al.
369 2006), and geneid(Blanco, et al. 2007) were then used to predict genes.

370 Annotated genes of novel sequences were analysed for Kyoto Encyclopedia of
371 Genes and Genomes (KEGG) terms and pathway enrichment using KOBAS (Xie, et
372 al. 2011).

373 **SNP calling**

374 To verify whether using the pan-genome as reference could improve SNP calling
375 efficacy, we randomly selected six pig samples (ranging from 10 to 30 \times coverage)

376 **(Supplementary Table 8)** and mapped their clean reads to the pan-genome and
377 Sscrofa11.1 for comparison. Duplicate reads were removed using Picard Tools. Then,
378 the Genome Analysis Toolkit (GATK, version 3.6) (McKenna, et al. 2010) was used
379 to detect SNPs. The following criteria were applied to all SNPs: (1) Variant
380 confidence/quality by depth (QD) > 2; (2) RMS mapping quality (MQ) > 30.0.

381 **RNA-seq analysis and noncoding RNA prediction**

382 The 92 strand-specific RNA-seq data (7-10 tissue libraries for each of 10 individuals)
383 were downloaded from the NCBI database (Bioproject: PRJNA311523) (Li, et al.
384 2017). All reads were mapped to the pan-genome by HISAT2 (Kim, et al. 2015).
385 Transcripts including novel splice variants were assembled using StringTie version
386 1.2.2 (Pertea, et al. 2015) and the FPKM (Fragments Per Kilobase per Million
387 mapped reads) values for these transcripts and genes in each sample were determined
388 using Ballgown (Frazee, et al. 2015). Finally, transcripts with FPKM ≥ 1 in at least
389 one sample were retained. After assembling and quantifying all transcripts, the
390 transcripts of pan-sequences were used for identification of high confidence coding
391 RNA by Coding Potential Calculator (CPC) (Kong, et al. 2007) online.

392 **Materials for Hi-C experiment**

393 Liver of BH-33, BH-34, BH-35, and BH-36 were collected from four female 2-years-
394 old Bama minipigs. Liver of F2 were collected from a 90-days-old female fetus of
395 Bama minipig. Ear skin fibroblasts DB-2 and DB-3 were established by using two 12-
396 days-old female Large White pigs. Ear skin fibroblasts XYZ were established by
397 using a 2-years-old female Wild Boar. Embryonic fibroblasts RC-7 and RC-8 were
398 established by using two 40-days-old female fetus of Chinese Rong Chang pig.
399 Mature adipocytes DB-2-Y and DB-3-Y were derived from pre-adipocytes which

400 were established by using the same pigs of Ear skin fibroblasts DB-2 and DB-3, by
401 inducing adipogenic differentiation.

402 All of the fibroblasts were grown in DMEM Dulbecco's Modified Eagle Medium
403 (DMEM, 11995-065, Gibco) containing 10% Fetal Bovine Serum (FBS, 10099-141,
404 Gibco) and 1× penicillin/streptomycin (P/S, 15140-122, Gibco), incubated at 37°C in
405 5% CO₂.

406 Pre-adipocytes were cultured in 10%FBS/DMEM-F12 (11330-032, Gibco) with
407 1×P/S until confluence and induced to differentiation as previously described. Briefly,
408 two days' post-confluence, cells were exposed to differentiation medium containing
409 0.5 mmol/L isobutylmethylxanthine (I5879, Sigma), 1 μmol/L dexamethasone
410 (D2915, Sigma), 850 nmol/L insulin (I6634, Sigma), 1 μmol/L rosiglitazone (R2408,
411 Sigma) and 10% FBS for three days. At the end of day 3, the differentiation medium
412 was replaced into maintenance medium with only 850 nmol/L insulin, 1 μmol/L
413 rosiglitazone and 10% FBS, and replenished every other day. After the differentiation
414 process, at least 90% of the cells had accumulated lipid droplets at day 15, and were
415 used as mature adipocytes (DB-2-Y and DB-3-Y).

416 **Hi-C experimental method**

417 Hi-C experiment on cells were performed according to the previously published
418 Hi-C protocol with some minor modifications (Lieberman-Aiden, et al. 2009).
419 Briefly, 25 million (M) cells were resuspended in 45 ml serum free DMEM, and 37%
420 formaldehyde was added to obtain a final concentration of 2% for chromatin cross-
421 linking. Cells were incubated at room temperature (20–25 °C) for 5 minutes, then
422 glycine was added to obtain a final concentration of 0.25 mol/L to quench the
423 formaldehyde. The mixture was incubated at room temperature for 5 minutes, and

424 subsequently on ice for at least 15 minutes. Fixed cells were lysed using a Dounce
425 homogenizer in the presence of cold lysis buffer (10 mmol/L Tris-HCl, pH 8.0, 10
426 mmol/L NaCl, 0.2% IGEPAL CA-630, and 1× protease inhibitor solution). Chromatin
427 digestion (restriction enzyme HindIII), labelling, and ligation steps were performed
428 according to the original protocol (Lieberman-Aiden, et al. 2009). After
429 deproteinization, removal of biotinylated free-ends, and DNA purification, Hi-C
430 libraries were controlled for quality and sequenced on an Illumina Hiseq X Ten
431 sequencer (paired-end sequencing with 150 bp in read length).

432 Hi-C experiment on liver tissue were performed as previously described using the
433 MboI restriction enzyme (Rao, et al. 2014), with minor modifications pertaining to
434 handling flash frozen primary tissues (Leung, et al. 2015). Briefly, 0.5 g flash frozen
435 liver tissue was pulverized in liquid nitrogen. Then cross-linking by 37%
436 formaldehyde in a final concentration of 4% and incubated at room temperature for 30
437 mins. Glycine was added to obtain a final concentration of 0.25 mol/L to quench the
438 formaldehyde. The mixture was incubated at room temperature for 5 minutes, and
439 subsequently on ice for at least 15 minutes. Cross-linked liver cells were filtered
440 through 70-µm and 40-µm nylon cell strainers and spinning down to collect the liver
441 cells. About 25 mg liver cell precipitate was used for Hi-C library preparation. The
442 Hi-C library preparation procedure was performed as previously described using the
443 MboI restriction enzyme (Rao, et al. 2014).

444 **Hi-C reads mapping, filtering, and generation of contact matrices**

445 Pre-processing paired-end sequencing data, reads mapping as well as filtering of
446 mapped di-tags was performed using the Juicer pipeline (v.1.8.9) (Durand, et al.
447 2016). Briefly, short reads were mapped to pan-genome using BWA (version 0.7.15-

448 r1140) (Li and Durbin 2009). Reads of low mapping quality were filtered using Juicer
449 with default parameters, discarding the invalid self-ligated and un-ligated fragments,
450 as well as PCR artefacts. Filtered di-tags were further processed with Juicer command
451 line tools to bin di-tags (10 kb bins) and to normalize matrices with KR normalization
452 (Knight and Ruiz 2013). Valid Hi-C read pairs should harbour more
453 intrachromosomal (cis) interactions than inter- (trans) (**Supplementary Table 11**). To
454 improve resolution, we combined the Hi-C data from the same tissue of same pig
455 breed after we randomly extracted 20 Gb data for correlation coefficient test. We
456 combined Hi-C data from DB-2 and DB-3 (Pearson's $r = 0.99$); RC-7 and RC-8
457 (Pearson's $r = 0.99$); DB-2-Y and DB-3-Y (Pearson's $r = 0.96$). After combined
458 samples, all processes were done in all the data. Normalized interaction matrices were
459 generated at two resolutions of low (100 kb) and high (20 kb) respectively
460 (**Supplementary Figure 14**).

461 **Identification of compartment A and B**

462 Identification of compartment A/B was performed as previously described using the
463 100-kb interaction matrix (Lieberman-Aiden, et al. 2009). Principal component
464 analysis (PCA) was performed to generate the first principal component (PC1) vectors
465 of each chromosome, and Spearman's correlation between PC1 and genomic
466 characteristics (gene density and GC content) were then calculated. GC content (%)
467 for each bin (100-kb bin sizes) was calculated using SeqKit (v.0.8.0) (Shen, et al.
468 2016). Gene density (number of genes per bin) was calculated based on the number of
469 promoters [from -2,000 to +500 bp of transcription start site (TSS)] located in
470 (namely more than 50% of the region should be overlapped) each bin. Compartment
471 A and B were determined by the PC1 values. Bins with positive Spearman's

472 correlation between PC1 values and genomic features were assigned as compartment
473 A, otherwise B.

474 **Identification of topologically associating domains (TADs) and topological** 475 **boundaries**

476 Higher-resolution TAD calls were generated following the previously described
477 procedure by using the directionality index (DI) metric (Dixon, et al. 2012). DI was
478 calculated using raw interaction counts between 20-kb bins to capture observed
479 upstream or downstream interaction bias of genomic regions. A hidden Markov model
480 (HMM) was then used to predict the states of DI for final TAD generation. The same
481 criteria 400 kb (distance between two adjacent TADs) was used to distinguish
482 unorganized chromatin from topological boundaries. That is, the topological
483 boundaries are less than 400 kb and unorganized chromatin is larger than 400 kb.

484 **Locating pan-sequences on Sscrofa11.1 based on Hi-C**

485 We normalized all Hi-C matrices on the same scale by KR normalization (Knight and
486 Ruiz 2013), ensuring that any differences between Hi-C are not attributable to
487 variation in sequence length. The maximum 100-kb bin of each pan-sequence
488 interacted (Interaction intensity ≥ 5) was collected as the potential location of pan-
489 sequences. Starting with the filtered 100-kb resolution bin of pan-sequences, we get
490 the higher resolution interval of 20 kb by taking the maximum 20-kb bin with each
491 100-kb bin.

492 **Identification of putative promoter and enhancer interactions**

493 We kept the interactions identified by PHYCHIC (Ron, et al. 2017) with FDR < 0.01
494 as high confidence interactions and used them to identify promoter-enhancer

495 interactions (PEI). Promoter segment was determined as a region from $-2,000$ to
496 $+500$ bp of the transcription start site (TSS). When at least half of a promoter segment
497 was in either one of the two bins which involved in a chromatin interaction, this
498 interaction was defined as a putative promoter interaction.

499 The bins which are distal (at least 40 kb upstream or downstream) from the
500 promoter and demonstrate the strongest interaction with the promoter than other
501 regions were determined as the enhancer interacting with the corresponding promoter.
502 This interaction of the two bins corresponding to the promoter and enhancer was
503 defined as a potential PEI. If our pan-sequences were located on a bin harbouring an
504 enhancer of a PEI, the pan-sequences could be potentially involved in the regulatory
505 functions of the enhancer. If the pan-sequences further demonstrate interactions with
506 the promoter of the same PEI, the involvement of the pan-sequences in the regulatory
507 functions of the enhancer would be regarded as highly confident and the pan-
508 sequences could be a potential enhancer itself.

509 **The pig pan-genome web server**

510 The web interface of PIGPAN was built by combining Apache web server, PHP,
511 HTML, JavaScript and relational database MySQL. Users can use all online resources
512 without preregistration. Our browser can be divided into two parts: frontend and
513 backend interfaces. The frontend consists of a home page, a download page and
514 several search pages. The MySQL relational database server stores 16 tables including
515 gap information, GC percent, seven regulatory signals of potential stem cells
516 (H3K27me3, H3K4me1, H3K4me3, H3K9me3, NANOG, PPARD AntiX, TAF1),
517 conserved elements of 20-species mammal, haploid copy number of 87 pigs, gene
518 expression, location of pan-sequences and gene annotation. The appropriate index was

519 built on the corresponding retrieval columns of the table. When a user submits an
520 entry, the backend will respond quickly to execute an SQL statement. PHP and
521 JavaScript manage the data analysis processes and visualize the results. Moreover, we
522 have introduced web-based software such as BLAST (Camacho, et al. 2009), BLAT
523 (Kent 2002) and Gbrowse (Casper, et al. 2018). Accordingly, users can query data
524 with rapid visualization in Gbrowse or enter a query sequence to search for
525 homologous regions in the genome. PIGPAN was tested in all major modern internet
526 browsers, including Firefox, Chrome, Internet Explorer, Safari and Opera. Therefore,
527 PIGPAN is a robust and easy-to-use website to facilitate the search for and
528 visualization of results for pig pan-genome analyses.

529 **Data availability**

530 The sequencing reads of each sequencing library have been deposited at NCBI for Hi-
531 C data (Project ID: PRJNA496307). The assembly of pig pan-genome and subsequent
532 analysis results are available in our PIGPAN website
533 (<http://animal.nwsuaf.edu.cn/code/index.php/panPig>). All other data supporting the
534 findings of this study are available in the article and its supplementary information
535 files or are available from the corresponding author on request.

536 **References and Notes**

- 537 Ai H, Fang X, Yang B, Huang Z, Chen H, Mao L, Zhang F, Zhang L, Cui L, He W, et al. 2015. Adaptation and
538 possible ancient interspecies introgression in pigs identified by whole-genome sequencing. *Nat Genet* 47:217-225.
- 539 Arumemi F, Bayles I, Paul J, Milcarek C. 2013. Shared and discrete interacting partners of ELL1 and ELL2 by
540 yeast two-hybrid assay. *Advances in Bioscience and Biotechnology* 04:774-780.
- 541 Blanco E, Parra G, Guigo R. 2007. Using geneid to identify genes. *Curr Protoc Bioinformatics* Chapter 4:Unit 4 3.
- 542 Burge CB, Karlin S. 1998. Finding the genes in genomic DNA. *Current Opinion in Structural Biology* 8:346-354.
- 543 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+: architecture
544 and applications. *BMC Bioinformatics* 10:421.
- 545 Casper J, Zweig AS, Villarreal C, Tyner C, Speir ML, Rosenbloom KR, Raney BJ, Lee CM, Lee BT, Karolchik D,
546 et al. 2018. The UCSC Genome Browser database: 2018 update. *Nucleic Acids Res* 46:D762-D769.
- 547 Dixon JR, Selvaraj S, Yue F, Kim A, Li Y, Shen Y, Hu M, Liu JS, Ren B. 2012. Topological domains in
548 mammalian genomes identified by analysis of chromatin interactions. *Nature* 485:376-380.
- 549 Doerks T, Copley RR, Schultz J, Ponting CP, Bork P. 2002. Systematic identification of novel protein domain
550 families associated with nuclear functions. *Genome Res* 12:47-56.
- 551 Dong P, Tu X, Chu PY, Lu P, Zhu N, Grierson D, Du B, Li P, Zhong S. 2017. 3D Chromatin Architecture of
552 Large Plant Genomes Determined by Local A/B Compartments. *Mol Plant* 10:1497-1509.
- 553 Durand NC, Shamim MS, Machol I, Rao SS, Huntley MH, Lander ES, Aiden EL. 2016. Juicer Provides a One-
554 Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell Syst* 3:95-98.
- 555 Fang X, Mou Y, Huang Z, Li Y, Han L, Zhang Y, Feng Y, Chen Y, Jiang X, Zhao W, et al. 2012. The sequence
556 and analysis of a Chinese pig genome. *Gigascience* 1:16.
- 557 Frantz LA, Schraiber JG, Madsen O, Megens HJ, Cagan A, Bosse M, Paudel Y, Crooijmans RP, Larson G,
558 Groenen MA. 2015. Evidence of long-term gene flow and selection during domestication from analyses of
559 Eurasian wild and domestic pig genomes. *Nat Genet* 47:1141-1148.
- 560 Frazee AC, Perteza G, Jaffe AE, Langmead B, Salzberg SL, Leek JT. 2015. Ballgown bridges the gap between
561 transcriptome assembly and expression analysis. *Nat Biotechnol* 33:243-246.
- 562 Golicz AA, Batley J, Edwards D. 2016. Towards plant pangenomics. *Plant Biotechnol J* 14:1099-1105.
- 563 Golicz AA, Bayer PE, Barker GC, Edger PP, Kim H, Martinez PA, Chan CK, Severn-Ellis A, McCombie WR,
564 Parkin IA, et al. 2016. The pangenome of an agronomically important crop plant *Brassica oleracea*. *Nat Commun*
565 7:13390.
- 566 Gordon SP, Contreras-Moreira B, Woods DP, Des Marais DL, Burgess D, Shu S, Stritt C, Roulin AC, Schackwitz
567 W, Tyler L, et al. 2017. Extensive gene content variation in the *Brachypodium distachyon* pan-genome correlates
568 with population structure. *Nat Commun* 8:2184.
- 569 Groenen MA, Archibald AL, Uenishi H, Tuggle CK, Takeuchi Y, Rothschild MF, Rogel-Gaillard C, Park C,
570 Milan D, Megens HJ, et al. 2012. Analyses of pig genomes provide insight into porcine demography and
571 evolution. *Nature* 491:393-398.
- 572 Guirao-Rico S, Ramirez O, Ojeda A, Amills M, Ramos-Onsins SE. 2018. Porcine Y-chromosome variation is
573 consistent with the occurrence of paternal gene flow from non-Asian to Asian populations. *Heredity (Edinb)*
574 120:63-76.
- 575 Hirsch CN, Foerster JM, Johnson JM, Sekhon RS, Muttoni G, Vaillancourt B, Penagaricano F, Lindquist E,
576 Pedraza MA, Barry K, et al. 2014. Insights into the maize pan-genome and pan-transcriptome. *Plant Cell* 26:121-
577 135.
- 578 Jeong H, Song KD, Seo M, Caetano-Anolles K, Kim J, Kwak W, Oh JD, Kim E, Jeong DK, Cho S, et al. 2015.
579 Exploring evidence of positive selection reveals genetic basis of meat quality traits in Berkshire pigs through
580 whole genome sequencing. *BMC Genet* 16:104.
- 581 Kent WJ. 2002. BLAT—The BLAST-Like Alignment Tool. *Genome Research* 12:47-56.
- 582 Kim D, Langmead B, Salzberg SL. 2015. HISAT: a fast spliced aligner with low memory requirements. *Nat*
583 *Methods* 12:357-360.
- 584 Knight PA, Ruiz D. 2013. A fast algorithm for matrix balancing. *Ima Journal of Numerical Analysis* 33:1029-
585 1047.

- 586 Kong L, Zhang Y, Ye ZQ, Liu XQ, Zhao SQ, Wei L, Gao G. 2007. CPC: assess the protein-coding potential of
587 transcripts using sequence features and support vector machine. *Nucleic Acids Res* 35:W345-349.
- 588 Larson G, Dobney K, Albarella U, Fang M, Matisoo-Smith E, Robins J, Lowden S, Finlayson H, Brand T,
589 Willerslev E, et al. 2005. Worldwide phylogeography of wild boar reveals multiple centers of pig domestication.
590 *Science* 307:1618-1621.
- 591 Leung D, Jung I, Rajagopal N, Schmitt A, Selvaraj S, Lee AY, Yen CA, Lin S, Lin Y, Qiu Y, et al. 2015.
592 Integrative analysis of haplotype-resolved epigenomes across human tissues. *Nature* 518:350-354.
- 593 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*
594 25:1754-1760.
- 595 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, Genome Project
596 Data Processing S. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25:2078-2079.
- 597 Li M, Chen L, Tian S, Lin Y, Tang Q, Zhou X, Li D, Yeung CKL, Che T, Jin L, et al. 2017. Comprehensive
598 variation discovery and recovery of missing sequence in the pig genome using multiple de novo assemblies.
599 *Genome Res* 27:865-874.
- 600 Li M, Tian S, Jin L, Zhou G, Li Y, Zhang Y, Wang T, Yeung CK, Chen L, Ma J, et al. 2013. Genomic analyses
601 identify distinct patterns of selection in domesticated pigs and Tibetan wild boars. *Nat Genet* 45:1431-1438.
- 602 Li R, Li Y, Zheng H, Luo R, Zhu H, Li Q, Qian W, Ren Y, Tian G, Li J, et al. 2010. Building the sequence map of
603 the human pan-genome. *Nat Biotechnol* 28:57-63.
- 604 Li YH, Zhou G, Ma J, Jiang W, Jin LG, Zhang Z, Guo Y, Zhang J, Sui Y, Zheng L, et al. 2014. De novo assembly
605 of soybean wild relatives for pan-genome analysis of diversity and agronomic traits. *Nat Biotechnol* 32:1045-1052.
- 606 Lieberman-Aiden E, van Berkum NL, Williams L, Imakaev M, Ragoczy T, Telling A, Amit I, Lajoie BR, Sabo PJ,
607 Dorschner MO, et al. 2009. Comprehensive mapping of long-range interactions reveals folding principles of the
608 human genome. *Science* 326:289-293.
- 609 Maretty L, Jensen JM, Petersen B, Sibbesen JA, Liu S, Villesen P, Skov L, Belling K, Theil Have C, Izarugaza
610 JMG, et al. 2017. Sequencing and de novo assembly of 150 genomes from Denmark as a population reference.
611 *Nature* 548:87-91.
- 612 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S,
613 Daly M, et al. 2010. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA
614 sequencing data. *Genome Res* 20:1297-1303.
- 615 Monat C, Pera B, Ndjiondjop MN, Sow M, Tranchant-Dubreuil C, Bastianelli L, Ghesquiere A, Sabot F. 2017. De
616 Novo Assemblies of Three *Oryza glaberrima* Accessions Provide First Insights about Pan-Genome of African
617 Rices. *Genome Biol Evol* 9:1-6.
- 618 Morgulis A, Gertz EM, Schaffer AA, Agarwala R. 2006. WindowMasker: window-based masker for sequenced
619 genomes. *Bioinformatics* 22:134-141.
- 620 Neafsey DE, Waterhouse RM, Abai MR, Aganezov SS, Alekseyev MA, Allen JE, Amon J, Arca B, Arensburger
621 P, Artemov G, et al. 2015. Mosquito genomics. Highly evolvable malaria vectors: the genomes of 16 *Anopheles*
622 mosquitoes. *Science* 347:1258522.
- 623 Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. 2015. StringTie enables improved
624 reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol* 33:290-295.
- 625 Rao SS, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, Sanborn AL, Machol I, Omer AD,
626 Lander ES, et al. 2014. A 3D map of the human genome at kilobase resolution reveals principles of chromatin
627 looping. *Cell* 159:1665-1680.
- 628 Ron G, Globerson Y, Moran D, Kaplan T. 2017. Promoter-enhancer interactions identified from Hi-C data using
629 probabilistic models and hierarchical topological domains. *Nat Commun* 8:2237.
- 630 Schatz MC, Maron LG, Stein JC, Wences AH, Gurtowski J, Biggers E, Lee H, Kramer M, Antoniou E, Ghiban E,
631 et al. 2014. Whole genome de novo assemblies of three divergent strains of rice, *Oryza sativa*, document novel
632 gene space of aus and indica. *Genome Biology* 15.
- 633 Shen W, Le S, Li Y, Hu F. 2016. SeqKit: A Cross-Platform and Ultrafast Toolkit for FASTA/Q File Manipulation.
634 *PLoS One* 11:e0163962.
- 635 Shyu RY, Jiang SY, Chou JM, Shih YL, Lee MS, Yu JC, Chao PC, Hsu YJ, Jao SW. 2003. RARRES3 expression
636 positively correlated to tumour differentiation in tissues of colorectal adenocarcinoma. *Br J Cancer* 89:146-151.
- 637 Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B. 2006. AUGUSTUS: ab initio prediction of
638 alternative transcripts. *Nucleic Acids Res* 34:W435-439.

- 639 Sun C, Hu Z, Zheng T, Lu K, Zhao Y, Wang W, Shi J, Wang C, Lu J, Zhang D, et al. 2017. RPAN: rice pan-
640 genome browser for approximately 3000 rice genomes. *Nucleic Acids Res* 45:597-605.
- 641 Tang Q, Gu Y, Zhou X, Jin L, Guan J, Liu R, Li J, Long K, Tian S, Che T, et al. 2017. Comparative
642 transcriptomics of 5 high-altitude vertebrates and their low-altitude relatives. *Gigascience* 6:1-9.
- 643 Wang X, Zheng Z, Cai Y, Chen T, Li C, Fu W, Jiang Y. 2017. CNVcaller: highly efficient and widely applicable
644 software for detecting copy number variations in large populations. *Gigascience* 6:1-12.
- 645 Wong KHY, Levy-Sakin M, Kwok PY. 2018. De novo human genome assemblies reveal spectrum of alternative
646 haplotypes in diverse populations. *Nat Commun* 9:3040.
- 647 Xie C, Mao X, Huang J, Ding Y, Wu J, Dong S, Kong L, Gao G, Li CY, Wei L. 2011. KOBAS 2.0: a web server
648 for annotation and identification of enriched pathways and diseases. *Nucleic Acids Res* 39:W316-322.
- 649 Yan G, Zhang G, Fang X, Zhang Y, Li C, Ling F, Cooper DN, Li Q, Li Y, van Gool AJ, et al. 2011. Genome
650 sequencing and comparison of two nonhuman primate animal models, the cynomolgus and Chinese rhesus
651 macaques. *Nat Biotechnol* 29:1019-1023.
- 652 Zhang X, Firestein S. 2002. The olfactory receptor gene superfamily of the mouse. *Nat Neurosci* 5:124-133.
- 653 Zhang Y, Liu T, Meyer CA, Eeckhoute J, Johnson DS, Bernstein BE, Nusbaum C, Myers RM, Brown M, Li W, et
654 al. 2008. Model-based analysis of ChIP-Seq (MACS). *Genome Biol* 9:R137.
- 655 Zhao Q, Feng Q, Lu H, Li Y, Wang A, Tian Q, Zhan Q, Lu Y, Zhang L, Huang T, et al. 2018. Pan-genome
656 analysis highlights the extent of genomic variation in cultivated and wild rice. *Nat Genet* 50:278-284.
- 657