# Chromosome-scale *de novo* assembly and phasing of a Chinese indigenous pig genome

Yalan Yang[1,2,#], Jinmin Lian[3,#], Bingkun Xie[4,#], Muya Chen[1,2], Yongchao Niu[3],

Qiaowei Li[1,2], Yuwen Liu[1,2], Guoqiang Yi[1,2], Xinhao Fan[1,2], Yijie Tang[1,2], Jiang Li[3],

Ivan Liachko[5], Shawn T. Sullivan[5], Bradley Nelson[5], Erwei Zuo[1,2], Zhonglin Tang[1,2*]

[1] Innovation Team of Pig Genome Design and Breeding, Agricultural Genome Institute at Shenzhen, Chinese Academy of Agricultural Sciences, Shenzhen, 518124, China;

[2] Genome Analysis Laboratory of the Ministry of Agriculture and Rural Affairs, Agricultural Genome Institute at Shenzhen, Chinese Academy of Agricultural Sciences, Shenzhen, 518124, China;

[3] Biozeron Shenzhen, Inc., Shenzhen, 518081, China;

[4] Guangxi Key Laboratory of Livestock Genetic Improvement, Guangxi Institute of Animal Science, Nanning, 530001, China;

[5] Phase Genomics Inc, Seattle, WA 98109, USA.

[#] These authors contributed equally to this work.

[*] To whom correspondence should be addressed. Tel. 86-0755-23251432, Fax. 86-0755-23251432. Email: tangzhonglin@caas.cn (Z.T.)

## Abstract

Chinese indigenous pigs differ significantly from Western commercial pig breeds in phenotypic and genomic characteristics. Thus, building a high-quality reference genome for Chinese indigenous pigs is pivotal to exploring gene function, genome evolution and improving genetic breeding in pigs. Here, we report an ultrahigh-quality phased chromosome-scale genome assembly for a male Luchuan pig, a representative Chinese domestic breed, by generating and combining data from PacBio Sequel reads, Illumina paired-end reads, high-throughput chromatin conformation capture and BioNano optical map. The primary assembly is ~ 2.58 Gb in size with contig and scaffold N50s of 18.03 Mb and 140.09 Mb, respectively. Comparison between primary assembly and alternative haplotig reveals numerous haplotype-specific alleles, which provide a rich resource to study the allele-specific expression, epigenetic regulation, genome structure and evolution of pigs. Gene enrichment analysis indicates that the Luchuan-specific genes are predominantly enriched in Gene Ontology terms for phosphoprotein phosphatase activity, signaling receptor activity and phosphatidylinositol binding. We provide clear molecular evolutionary evidence that the divergence time between Luchuan and Duroc pigs is dated back to about 1.7 million years ago. Meanwhile, Luchuan exhibits fewer events of gene family expansion and stronger gene family contraction than Duroc. The positively selected genes (PSGs) in Luchuan pig significantly enrich for protein tyrosine kinase activity, microtubule motor activity, GTPase activator activity and ubiquitin-protein transferase activity, whereas the PSGs in Duroc pig enrich for G-protein coupled receptor activity. Overall, our findings not only provide key benchmark data for the pig genetics community, but also pave a new avenue for utilizing porcine biomedical models to study human health and diseases.

## Introduction

*Sus scrofa* (pig) is one of the most important domesticated animals for its enormous value in food supply and biomedical research. Plenty of archaeological and molecular evidence suggests that pigs were independently domesticated in the Near East and China about 9,000 years ago [1-3]. The effects of geographical divergence, local adaptation and artificial selection result in great phenotypic and genomic diversity among pigs from distinct locations and breeds [4, 5]. In China, there are ~ 100 native breeds (China National Commission of Animal Genetic Resources 2011), accounting for about one-third of world breeds. To study pig genetics, the present pig reference genome (Sscrofa11.1) was derived from a Western pig (the Duroc breed) [6, 7]. However, Eastern and Western pigs have different genetic backgrounds. To better explore gene function, genome evolution and improve genetic breeding in pigs, it is of great value to build a reference genome for Chinese indigenous pigs.

Two main challenges for assembling a state-of-the-art high-quality reference genome are chromosome-scale contiguity and diploid phasing. Previous studies reported multiple *de novo* assemblies of Chinese native breeds using whole-genome shotgun-based strategies, and shed light on genomic and phenotype diversities of Chinese domestic pigs [4, 8-10]. Nonetheless, these shotgun-based approaches cannot yield large continuous genome scaffolds, significantly limiting the quality and contiguity of the current Chinese pig genome assemblies. Beyond genome assembly at the chromosome scale, accurate representation of haplotypes is crucial to identifying single-nucleotide polymorphisms (SNPs) and structural variants (SVs), haplotype structure and heterozygosities between two homologous chromosomes. Therefore, a phased genome assembly is essential for studies on intraspecific variation, allele-specific expression, epigenetic regulation, and chromosome evolution, as well as understanding how combinations of variants impact phenotypes [15-17]. Among the new technologies to tackle the two challenges in genome assembly, long-read sequencing, high-throughput chromatin conformation capture (Hi-C) and optical

75    mapping technologies have been developed for ordering and orienting assembly contigs,

76    and thus can create phased chromosome-scale genome assemblies [11]. These

77    technologies have substantially improved genome assemblies for human, goat and

78    gorilla [12-14]. However, a phased genome assembly with chromosome-scale

79    contiguity for pigs is not yet to available, which results in the lack of resolution for pigs

80    inter-haplotype variations, and impedes the dissection of the genetic basis of phenotypic

81    differences in domestication between Eastern and Western pigs.

82    Here, we applied long-read sequencing (Pacbio), short paired-end reads (Illumina),

83    Hi-C and optical map (BioNano) technologies to generate an assembly of the Luchuan

84    pig, an indigenous breed from Guangxi province in South China. As a representative of

85    the native breeds in China, Luchuan pig has many distinguishing phenotypic features

86    comparing with Western domesticated pigs, including low growth rate, high fat content，

87    excellent meat quality, early maturity, high fecundity, good maternal stability, wide

88    adaptability to coarse feeding and strong disease resistance [4, 5]. To study the genetic

89    basis underlying these phenotypic differences, Luchuan is an ideal material for building

90    a high-quality reference genome representing Chinese indigenous pigs. In our study, a

91    high-contiguous, chromosome-scale phased assembly of the Luchuan pig genome was

92    *de novo* assembled. To our knowledge, this is the first published phased chromosome-

93    scale assembly for mammals, providing important genetic resources and

94    methodological references for future studies of animal genomic evolution, molecular

95    breeding and biomedical research.

## Material and Methods

## Sample collection and sequencing

98    A Luchuan boar was obtained from the Institute of Animal Science of Guangxi

99    province, China, for genome assembly. Genomic DNA was extracted from its blood

100    sample. In order to generate a chromosome-scale assembly, four different genome

101    libraries were constructed and sequenced according to the manufacturers' instructions:

102    (i) Whole genome sequencing (WGS) by PacBio Sequel platform (20-kb library); (ii)

103 Hi-C chromosome conformation captured reads sequencing by Phase genomics; (iii)

104 Short reads paired-end sequencing (150bp in length) by Illumina NovaSeq 6000

105 platform; (iv) BioNano optical map data (Nt.BspQI, Nb.BssSI and DLE-1 enzymes).

106 　　To fully assist genome annotation, thirty-seven RNAs from 14 tissues (heart, lung,

107 adipose, kidney, liver, brain, spleen, stomach, leg muscle, dorsal muscles, testis, ovary,

108 large intestine, small intestine) at four developmental stages (Days 0, 14, 50 and adult

109 pigs) of Luchuan pigs (4 individuals) were equally pooled together. Two strand-specific

110 RNA-seq libraries with an insert size of 350 bp using the NEBNext® Ultra™

111 Directional RNA Library Prep Kit for Illumina® (NEB, USA) were prepared and

112 sequenced on an Illumina NovaSeq 6000 platform, to generate 150bp paired-end reads

113 (Berry Genomics Co., Ltd., Tianjin, China). A PacBio full-length transcriptome library

114 was constructed and sequenced on the Pacific Bioscience RS II sequencer (Berry

115 Genomics, Co., Ltd., Beijing, China).

116 　　All animals and samples used in this study were collected according to the

117 guidelines for the care and use of experimental animals established by the Ministry of

118 Agriculture and Rural Affairs of China.

119

## *De novo* genome assembly and scaffolding

121 　　The primary contigs were assembled with the Falcon software packages (v2.0.5)

122 [16] followed by the FALCON-Unzip and Arrow (v2.2.2) polishing, then a Hi-C-based

123 contigs phasing was processed by FALCON-Phase to create phased, diploid contigs.

124 Phase Genomics' Proximo Hi-C genome scaffolding platform was used to establish

125 chromosome-scale scaffolds from the draft assembly using a method similar to that

126 described previously [14]. Following diploid chromosomal scaffolding, a round of

127 polishing using Juicebox (v1.8.8) [18, 19] was performed to correct small errors in

128 chromosome assignment, ordering and orientation. After a draft set of scaffolds was

129 generated, FALCON-Phase was run again for Hi-C based scaffold phasing. The

130 Illumina sequencing data were further used to improve the assembly by Pilon (v1.22)

131 software. Given the availability of a relatively good quality of the Duroc pig

132    (Sscrofa11.1) genome, a reference-assisted scaffolding strategy was conducted to get

133    chromosome-level pseudomolecules with Chromosomer software (v0.1.4a) [20].

134    Quality control on the integrity of the assembly of genic regions was performed by

135    using the independent BUSCO v3 benchmark (http://busco.ezlab.org/) [21].

136

## Assembly quality assessment

138    BioNano optical map data was used to assess the assembly quality, which produces

139    physical maps with unique sequence motifs that can provide long-range structural

140    information of the genome. Briefly, high-molecular weight DNA was extracted from

141    the pig blood sample and digested with nickases Nt.BspQI, Nb.BssSI and Direct

142    Labeling Enzyme 1 (DLE-1), respectively. After labeling and staining, DNA was loaded

143    onto the Saphyr chip for sequencing. Raw data for each enzyme library were collected

144    and converted into a BNX file by AutoDetect software, to obtain basic labeling and

145    DNA length information. The filtered raw DNA molecules in BNX format were aligned,

146    clustered and assembled into the BNG map by using the Bionano Solve pipeline. Two

147    enzyme (Nt.BspQI, Nb.BssSI) hybrid scaffolding was firstly processed to produce a set

148    of initial hybrid scaffold, a second round of hybrid scaffolding with genome map of

149    DEL-1 enzyme was followed.

150

## Repeat annotation

152    There are two main types of repeats in the genome: tandem and interspersed.

153    Tandem repetitive sequences were identified using Tandem Repeats Finder (TRF,

154    version4.07). The interspersed repeat contents were identified using two methods: *de*

155    *novo* repeat identification and known repeat searching against existing databases.

156    RepeatModeler (version 1.0.8, http://www.repeatmasker.org/RepeatModeler/) was

157    used to predict repeat sequences in the genome, and RepeatMasker (version 4.0.7) [22]

158    was then used to search the Luchuan pig genome against the *de novo* transposable

159    elements (TE) library. The homology-based approach involved applying commonly

160    used databases of known repetitive sequences, RepeatMasker (version 4.0.7) and the

161    Repbase database (version 21) [23] were used to identify TEsin the assembled genome.

162    RepeatMasker and Repeat Protein Masker (http://repeatmasker.org) were applied for

163    TEs identification at the DNA and protein levels, respectively.

164

## Gene prediction and annotation

166    Protein-coding region identification and gene prediction were conducted through a

167    combination of three approaches as following:

168    (i) Homology-based prediction. Protein sequences for human and five animal

169    genomes (mouse, cattle, dog, goat and the Duroc pig) were downloaded from Ensembl

170    release-95, and aligned to the Luchuan assembly using the TBLASTN program

171    available in the BLAST v2.2.24 (E-value cutoff 1e-05). Then the SOLAR

172    (version0.9.6), a dynamic program algorithm to link putative exons together, was

173    employed to analyze the TBLASTN results. GeneWise (version 2.4.1) [24] was used to

174    predict the exact gene structure of the corresponding genomic regions on each matched

175    sequences;

176    (ii) *De novo* prediction. Four *ab initio* gene prediction programs including

177    Augustus (version 3.2.1) [25], GlimmerHMM (version 3.0.4) [26], Geneid (version

178    1.4.4) [27] and SNAP (version 2006-07-28) [24], were employed to predict coding

179    regions in the repeat-masked genome;

180    (iii) Transcriptome-based prediction methods. RNA-seq data (26.35 Gb) reads

181    were mapped to the assembly using Hisat2 (version 2.1.0) [28]. Stringtie (version 1.2.2)

182    and TransDecoder (version 3.0.1) were used to assemble the transcripts and identify

183    candidate coding regions into gene models. For PacBio full-length transcriptome data

184    (Iso-Seq), transcripts were identified by IsoSeq3 (version 3.1.0) with default parameters,

185    then the Iso-Seq data were mapped to the reference genome with minimap2 (version

186    2.15-r905). Furthermore, Cupcake ToFU (v5.8) was used to get the final unique, full-

187    length and high-quality isoforms of Pacbio data.

188    All gene models predicted based on the above three approaches were combined by

189 EvidenceModeler (EVM) into a non-redundant set of gene structures, and the produced

190 gene models were finally refined using the Program to Assemble Spliced Alignments

191 (PASA v2.3.3) [29]. Functional annotation of protein-coding genes (PCGs) was

192 achieved using BLASTP (E-value 1e-05) against two integrated protein sequence

193 databases: SwissProt and TrEMBL. Protein domains were annotated by InterProScan

194 (v5.30). The Gene Ontology (GO) terms for each gene were extracted with

195 InterProScan [30]. The pathways in which the genes might be involved were assigned

196 by BLAST against the KEGG databases (release 59.3) [31] with an E-value cutoff of

197 1e-05.

198

## Noncoding RNAs annotation

200 The transfer RNAs (tRNA) genes were predicted by tRNAscan-SE (version 1.3.1)

201 [32] with eukaryote parameters. The ribosomal RNA (rRNA) fragments were predicted

202 by aligning to human template rRNA sequences using BlastN (version 2.2.26) at an E-

203 value of 1e-5. The microRNAs (miRNAs) and small nuclear RNAs (snRNAs) were

204 detected by searching against the Rfam database (release 12.0) [33] with INFERNAL

205 (version 1.1.1) [34]. Long non-coding RNAs (LncRNAs) and Circular RNAs

206 (circRNAs) were predicted by methods described previously [4, 35, 36].

207

## Identification of orthologous gene sets across species

209 A gene family indicates a set of similar genes that descended from a single original

210 gene in the last common ancestor of considered species. Orthologous gene sets of

211 Luchuan pig, Duroc pig, cattle, goat, dog, mouse and human were used for genome

212 comparisons. For a gene with multiple isoforms, we chose the longest transcript ($\geq 50$

213 amino acids) to represent the gene. The TreeFam methodology [37] was used to define

214 a gene family and result in 3,733 single-copy orthologous genes for the six mammalian

215 species. In addition, the one-to-one orthologous between these species were defined

216 using BLASTP based on the Bidirectional Best Hit (BBH) method with a sequence

217    coverage > 80% and identity > 80%, followed by selection of the best match.

218

## Variants calling

220    The primary assembly of Luchuan genome was aligned with the alternative

221    haplotig assembly and the Duroc contigs by MUMmer (version 3.23)[38] with default

222    parameters, and one-to-one genomic alignment results were extracted with the 'delta-

223    filter -1' parameter. SNPs and indels were identified by show-snp from the one-to-one

224    alignment blocks (parameter '-ClrT –x 1'). Structural variations were identified by

225    Assemblytics (v1.0) software [39] base on the alignment blocks from MUMmer.

226

## Phylogenetic tree construction and evolution rate estimation

228    Single-copy gene families were used to construct a phylogenetic tree for Luchuan

229    pig and the other mammalian genomes (Duroc pig, cattle, goat, dog, mouse and human).

230    Four-fold degenerate sites were extracted from each family and concatenated into one

231    supergene for each species. PhyML v3.0 was adopted to reconstruct the phylogenetic

232    tree based on the GTR+gamma substitution model [40]. The divergence time among

233    Luchuan pig, Duroc pig, cattle, goat, dog, mouse and human were estimated using the

234    MCMCtree program (version 4.4) as implemented in the Phylogenetic Analysis of

235    Maximum Likelihood (PAML) package [41], with an independent rates clock and

236    HKY85 nucleotide substitution model. The calibration times (differentiation time

237    between human and mouse, human and goat, cattle and goat, pig and goat) were derived

238    from the TimeTree database [42].

239

## Results

## Assembly and phasing of the Luchuan pig genome

242    To construct a high-quality reference genome for Chinese indigenous pigs, a male

243    Luchuan pig was used for WGS, which generated ~140× Pacbio Sequel long reads

244     (348.71 Gb), ~41× Hi-C reads (102.42 Gb, Phase Genomics), ~86× Illumina paired-

245     end reads (214.48 Gb), and ~351× BioNano optical map data (879.44 Gb, Bionano

246     Genomics).

247         The Pacbio reads were first assembled *de novo*, producing an initial contig

248     assembly with N50 of 18.68 Mb and a total length of 2.52 Gb. Then the assembly was

249     integrated with Hi-C data to create phased diploid chromosome-scale scaffolds

250     (Supplementary Figure 1), generating an alternative haplotype sequence with contig

251     N50 of 18.79Mb, scaffold N50 of 141.24Mb and a total length of 2.55 Gb. After

252     improving the assembly based on Illumina sequencing data, the optical map data were

253     used to validate, correct and merge the scaffolds (Supplementary Table 1). Given the

254     high quality of the present Duroc reference genome (Sscrofa11.1), a reference-assisted

255     scaffolding strategy was used to get chromosome-level pseudomolecules (Figure 1A).

256     Finally, we generated a high-contiguous, chromosome-scale and phased assembly of

257     the Luchuan genome, yielding a 2.58 Gb primary assembly with a contig N50 of 18.03

258     Mb and a scaffold N50 of 140.09 Mb. This assembly is comparable in quality to the

259     Duroc genome [7] and much better than other published pig genomes [4, 8-10] (Table

260     1). Remarkably, the alternative haplotig assembly size is very close to the primary

261     assembly with a contig N50 of 17.77 Mb and a scaffold N50 of 140.08 Mb.

262

263     **Table 1. Comparison of features between the Luchuan pig and other assemblies**.

| | Luchuan | Duroc*[7] | Tibetan wild[8] | Wuzhishan[9] | Bama[10] |
|---|---|---|---|---|---|
| **Sequenced genome size (Gb)** | 2.58 | 2.50 | 2.43 | 2.64 | 2.49 |
| **Contig N50 (Mb)** | 18.03 | 41.89 | 0.0207 | 0.0235 | 1.01 |
| **Scaffold N50 (Mb)** | 140.09 | 138.97 | 1.06 | 5.43 | 140.44 |
| **Percentage of anchoring and ordering** | 96.1% | 97.34% | - | - | 97.49% |
| **Predicted PCGs** | 22,710 | 22,452 | 21,806 | 20,326 | 21,334 |
| **Repeat proportion (%)** | 40.16 | 40.55 | 39.47 | 38.20 | 37.32 |
| **Complete BUSCOs (%)** | 95.1 | 96.0 | 93.1 | 95.2 | 93.9 |

264     * Statistic of Duroc pig genome was based on Sscrofa11.1 (Ensembl release-95).

265    The reference assessment revealed that approximately 96.1% of the 2.58 Gb

266    assembled final Luchuan assembly was assigned to 20 chromosomes (18 autosomes

267    and X/Y chromosome) (Supplementary Table 2-3). The 20 chromosomes were made

268    up of 466 contigs, reflecting the low fragmentation of these assemblies. We further

269    evaluated the genome assembly quality, and found that 95.1% of the 4,104 core genes

270    in the OrthoDB mammalian database were identified in the Luchuan primary assembly,

271    of which 94.4% were single-copy, 0.7% duplicated, 2.9% fragmented and 2.0% missing

272    (Supplementary Table 4).

273

## Validation of the phased diploid assemblies

275    The pseudo-chromosomes of Luchuan pig presented great colinearity with

276    Sscrofa11.1, supporting a high-quality genome assembly (Figure 1B; Supplementary

277    Figure 2). It is worth noting that the alternative haplotig also has highly collinear

278    relationships with Duroc pig assembly (Supplementary Figure 2). Additionally, to

279    assess the scaffolding accuracy of Luchuan assembly, we adopted the nickases

280    Nt.BspQI, Nb.BssSI and DLE-1 for optical map library construction, and got 453 Gb,

281    345 Gb, and 618 Gb raw data using these three enzymes, respectively. After removing

282    molecules in lengths less than 150 kb, we obtained 303 Gb, 268 Gb and 308 Gb high-

283    quality optical molecules, accounting for > 100× coverage of genome size. The N50 of

284    the molecules are 358 kb, 394 kb and 248 kb for Nt.BspQI , Nb.BssSI nickase and

285    DLE-1 enzymes, respectively (Supplementary Table 5). The high concordance between

286    the assembly and the optical map data provides strong support for the robustness of the

287    assembly (Figure 2). By comparison between the contigs/scaffolds and optical maps,

288    74 and 73 conflicts were detected for the primary assembly and alternative haplotig,

289    respectively. After conflict correction, we assembled 63 and 64 hybrid scaffolds based

290    on genome map hybrid assembly for the primary assembly and alternative haplotig,

291    respectively. These results demonstrated the high reliability of the alternate haplotype

292    assembly.

293

## Genetic variations between primary assembly and alternate haplotig

By comparing the primary assembly to the alternate haplotig, we identified numerous haplotype-specific alleles, including 6.83 million SNPs, 1.64 million short indels and 23,539 SVs (Figure 3). Among the SNPs, most (97.54%) were located in intergenic regions (63.01%) and intronic regions (34.51%), only 0.56% were located in coding sequences. Of the SNPs present in coding regions, 24,056 were synonymous and 13,959 were non-synonymous. In addition, 2,479 and 463 indels may result in frameshift and non-frameshift variations, respectively. These variations are valuable to further study the allele-specific expression, epigenetic regulation, genome structure and evolution in pigs.

## Genome annotation

We predicted a total of 22,710 PCGs with strong evidence in Luchuan by combining *ab initio* prediction, homologous protein prediction and transcriptome alignment. Of these PCGs, ~90% gain clear supporting evidence based on transcriptome sequencing data and functional annotation information (Table 1; Supplementary Table 6). The average length of gene, exon and intron were 40,062bp, 177bp and 4,709bp, respectively. We also annotated 2,835 small ncRNAs including 388 miRNAs, 394 rRNAs, 1,076 tRNAs and 977 snRNAs. Additionally, 3,066 novel lncRNAs and 1,019 novel circRNAs were identified in the Luchuan genome (Supplementary Table 7).

Repeat elements accounted for ~40.16% of the Luchuan genome (Supplementary Table 8-9). The two largest repeat classes were long-interspersed elements (LINEs) and short interspersed nuclear elements (SINEs), which comprised 27.83% and 10.86% of the genome, respectively. Tandem repeats constituted 3.88% of the genome. The number and length of genes and the proportion of repeat elements were similar to those present in the pig reference genome and other assemblies [4, 6-9].

## Comparative Genomic and Phylogenetic Analyses

We identified a total of 8,481 homologous gene families that are shared among

323 Luchuan, Duroc, goat and human. Interestingly, there are 163 and 134 gene families

324 specifically identified in Luchuan and Duroc, respectively (Figure 4A). Among those

325 Luchan-specific gene families, 421 genes with supporting evidence of transcription or

326 Interpro functional annotation were considered to be the high-quality Luchuan-specific

327 genes. These genes are significantly (FDR < 0.05) enriched in GO terms for

328 phosphoprotein phosphatase activity, signaling receptor activity and

329 phosphatidylinositol binding. By comparison, the 207 Duroc-specific genes are

330 functionally over-represented in biological processes related to actin filament binding,

331 peptidase inhibitor activity, pheromone receptor activity, microtubule motor activity

332 and epidermis development (Supplementary Table 10).

333 　　A phylogenetic tree was constructed using the pigs (Luchuan and Duroc) and five

334 other mammals (cattle, goat, dog, human and mouse). As shown in Figure 4B, the

335 divergence time between Luchuan and Duroc was estimated to be about 1.7 million

336 years ago (MYA). Compared with Duroc, Luchuan showed fewer events of gene family

337 expansion (63 vs. 433), and more events of gene family contraction (560 vs. 161)

338 (Figure 4B). Notably, expanded genes in Luchuan were closely related to response to

339 oxidative stress and biotic stimulus. In Duroc, the olfactory-related genes were

340 significantly expanded, consistent with a previous study [8]. In addition, expanded

341 genes in Duroc are significantly (P < 0.05) enriched in GO terms for

342 galactosyltransferase activity, antioxidant activity and growth factor activity.

## Bidirectional selection between Luchuan and Duroc pigs

344 　　To study the bidirectional selection between Luchan and Duroc pigs, we further

345 screened out 7,222 one-to-one orthologous gene sets from the seven mammals. We

346 found 272 and 768 positively selected genes (PSGs) in the Luchuan and Duroc pigs (P

347 < 0.05, likelihood ratio test), respectively. It is worth noting that 25 PSGs were shared

348 in both breeds, such as *CACNA1F*, a calcium channel subunit gene, and *RBM46*, an

349 RNA binding motif protein. Enrichment analysis revealed the PSGs detected in

350 Luchuan were especially enriched in GO terms related to protein tyrosine kinase

351 activity (8 PSGs), microtubule motor activity (6 PSGs), GTPase activator activity (6

352 PSGs) and ubiquitin-protein transferase activity (6 PSGs), whereas PSGs in Duroc pigs

353 were significantly enriched in GO terms for G-protein coupled receptor activity, which

354 is closely related with the olfactory receptors (Supplementary Table 11).

355

## Discussion

357     Chinese indigenous and Western pigs are independently domesticated and exhibit

358 a great spectrum of phenotypic and genomic differences [4, 6-9]. A comprehensive

359 exploration of the genetic diversity within and between pig breeds is important for

360 animal breeding and biomedical research. The present pig reference genome

361 (Sscrofa11.1) was derived from a Western breed (Duroc pig) [7] with high continuity

362 and quality. Increased accessibility to short-read sequencing has resulted in a deluge of

363 genome assemblies for Chinese indigenous pigs, although incomplete and fragmented

364 compared with Duroc [4, 8-10]. Until now, no chromosome-level phased assemblies of

365 Chinese indigenous pigs have been built, so accurately investigating the full range of

366 genetic variations and phased diploid architecture is extremely difficult. Recently, rapid

367 progress in high-throughput DNA sequencing and library preparation methods have

368 enabled the generation of phased genome assemblies with chromosome-level quality

369 [15-17]. Built upon these most recent technology breakthroughs, here we present, to

370 our knowledge, the first phased chromosome-scale genome assembly of pigs, which is

371 also the first such type of published assembly for mammals. Our genome assembly

372 yields a 2.58 Gb primary assembly with a contig N50 of 18.03 Mb, with comparable

373 quality to the current reference genome [7].

374     Synteny analysis revealed strong collinearity between the genomes of Luchuan and

375 Duroc pigs, supporting great overall quality of our assembly. Notably, our assembly

376 approach also makes it possible to construct a high-quality alternate haplotig assembly,

377 which is comparable to the primary assembly with a scaffold N50 size of 17.77 Mb.

378 Using the phased diploid assembly, we are able to identify structural variations between

379 two homologous chromosomes [15, 16], which are important for understanding how

380   combinations of variants impact phenotypes. Millions of genetic variations between
381   primary assembly and alternate haplotig of Luchuan genome were identified in our
382   study, which provided an unprecedentedly detailed resource to further study the allele-
383   specific expression, epigenetic regulation, genome structure and evolution of Eastern
384   pigs [15-17]. Moreover, combining our Luchuan genome and the classic Duroc
385   assembly would provide foundational resources to study the genetic basis underlying
386   the phenotypic differences between Eastern and Western pigs.

387   To study the evolution and domestication of Luchuan pigs, we reconstructed the
388   phylogenetic tree among Luchuan, Duroc, cattle, goat, dog, human and mouse. Our
389   analysis revealed that the divergence time between Luchuan and Duroc was about 1.7
390   MYA, which is in close proximity to the split time between Asian and European wild
391   boars (0.8-2 MYA) [43-45]. Gene replication is one of the basic mechanisms for
392   acquiring new functions and physiological features, and accordingly studying gene
393   family expansion and contraction provides unique perspectives on the genetic basis of
394   local domestication and adaptation [46, 47]. The Luchuan genome exhibited fewer
395   events of gene family expansion and stronger gene family contraction compared with
396   Duroc pig, which is in accordance with the comparative analysis between Duroc pig
397   and Tibetan wild boar [8]. Duroc pig was reported to have markedly more olfactory-
398   related genes than Tibetan wild boar [8]. Our results also confirmed that these genes
399   were significantly expanded and positively selected in Duroc pig. The oxidative stress
400   and response to biotic stimulus–related genes were expanded and GTPase activator
401   activity-related genes were positively selected in Luchuan pig, which might confer the
402   remarkable capabilities of Luchuan to adapt to coarse feeding and strong resistance to
403   diseases, which are important features shared by many Chinese indigenous breeds. The
404   PSGs analysis suggested that the Duroc pig had experienced stronger selection
405   pressures during breeding than Luchuan pig. These results provided novel insights into
406   the distinct evolutionary scenarios occurring under different local adaptation and
407   artificial selection between Chinese indigenous and Western pig breeds.

408   Overall, we presented the first phased chromosome-scale genome assembly of a
409   Chinese indigenous breed, which provides great resources for understanding pig

410 evolution and domestication. This Luchuan pig genome assembly would benefit the

411 dissection of the genetic basis and molecular mechanisms underlying phenotypic

412 differences between and within pig breeds, facilitate molecular breeding to improve

413 economical traits, and shed light on the etiology of human traits and diseases.

414

## Acknowledgements

420

## Author contributions

422 Z.L.T conceived, coordinated and managed the project; Y.L.Y, Y.W.L, G.Q.Y, M.Y.C,

423 Y.C.N, J.M.L, J.L, I.L, S.T.S and B.N assembled and annotated the genome sequences,

424 and carried out other computational and bioinformatics analysis; B.K.X provided the

425 Luchuan pigs and helped in samples collection. Z.L.T, X.H.F, Y.L.Y and Y.J.T

426 performed animal experiment and collected biological samples; Y.L.Y, Y.C.N and J.M.L

427 wrote the manuscript; Z.L.T, Y.L.Y, G.Q.Y, and E.W.Z revised the paper. All authors

428 read and approved the final manuscript.

429

## References

431 1.    Giuffra E, Kijas J, Amarger V, Carlborg Ö, Jeon J-T, Andersson L: **The origin of the domestic**
432       **pig: independent domestication and subsequent introgression.** *Genetics* 2000, **154:**1785-
433       1791.
434 2.    Larson G, Dobney K, Albarella U, Fang M, Matisoo-Smith E, Robins J, Lowden S, Finlayson
435       H, Brand T, Willerslev E, et al: **Worldwide phylogeography of wild boar reveals multiple**
436       **centers of pig domestication.** *Science* 2005, **307:**1618-1621.
437 3.    Fang M, Andersson L: **Mitochondrial diversity in European and Chinese pigs is consistent**
438       **with population expansions that occurred prior to domestication.** *Proceedings of the Royal*
439       *Society B: Biological Sciences* 2006, **273:**1803-1810.
440 4.    Li M, Chen L, Tian S, Lin Y, Tang Q, Zhou X, Li D, Yeung CK, Che T, Jin L: **Comprehensive**
441       **variation discovery and recovery of missing sequence in the pig genome using multiple de**

novo assemblies. *Genome research* 2017, **27:**865-874.

5. Ai H, Fang X, Yang B, Huang Z, Chen H, Mao L, Zhang F, Zhang L, Cui L, He W, et al: **Adaptation and possible ancient interspecies introgression in pigs identified by whole-genome sequencing.** *Nature Genetics* 2015, **47:**217-225.

6. Groenen MA, Archibald AL, Uenishi H, Tuggle CK, Takeuchi Y, Rothschild MF, Rogel-Gaillard C, Park C, Milan D, Megens H-J: **Analyses of pig genomes provide insight into porcine demography and evolution.** *Nature* 2012, **491:**393.

7. Warr A, Affara N, Aken B, Beiki H, Bickhart DM, Billis K, Chow W, Eory L, Finlayson HA, Flicek P: **An improved pig reference genome sequence to enable pig genetics and genomics research.** *bioRxiv* 2019**:**668921.

8. Li M, Tian S, Jin L, Zhou G, Li Y, Zhang Y, Wang T, Yeung CK, Chen L, Ma J, et al: **Genomic analyses identify distinct patterns of selection in domesticated pigs and Tibetan wild boars.** *Nat Genet* 2013, **45:**1431-1438.

9. Fang X, Mou Y, Huang Z, Li Y, Han L, Zhang Y, Feng Y, Chen Y, Jiang X, Zhao W: **The sequence and analysis of a Chinese pig genome.** *GigaScience* 2012, **1:**16.

10. Zhang L, Huang Y, Wang M, Guo Y, Liang J, Yang X, Qi W, Wu Y, Si J, Zhu S, et al: **Development and Genome Sequencing of a Laboratory-Inbred Miniature Pig Facilitates Study of Human Diabetic Disease.** *iScience* 2019, **19:**162-176.

11. Kronenberg ZN, Rhie A, Koren S, Concepcion G, Peluso P, Munson K, Hiendleder S, Fedrigo O, Jarvis E, Phillippy A: **Extended haplotype phasing of de novo genome assemblies with FALCON-Phase.** *bioRxiv* 2019**:**327064.

12. Gordon D, Huddleston J, Chaisson MJ, Hill CM, Kronenberg ZN, Munson KM, Malig M, Raja A, Fiddes I, Hillier LW, et al: **Long-read sequence assembly of the gorilla genome.** *Science* 2016, **352:**aae0344.

13. Chaisson MJ, Wilson RK, Eichler EE: **Genetic variation and the de novo assembly of human genomes.** *Nat Rev Genet* 2015, **16:**627-640.

14. Bickhart DM, Rosen BD, Koren S, Sayre BL, Hastie AR, Chan S, Lee J, Lam ET, Liachko I, Sullivan ST, et al: **Single-molecule sequencing and chromatin conformation capture enable de novo reference assembly of the domestic goat genome.** *Nat Genet* 2017, **49:**643-650.

15. Zhou B, Ho SS, Greer SU, Zhu X, Bell JM, Arthur JG, Spies N, Zhang X, Byeon S, Pattni R: **Comprehensive, integrated, and phased whole-genome analysis of the primary ENCODE cell line K562.** *Genome research* 2019, **29:**472-484.

16. Chin CS, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, Dunn C, O'Malley R, Figueroa-Balderas R, Morales-Cruz A, et al: **Phased diploid genome assembly with single-molecule real-time sequencing.** *Nat Methods* 2016, **13:**1050-1054.

17. Seo JS, Rhie A, Kim J, Lee S, Sohn MH, Kim CU, Hastie A, Cao H, Yun JY, Kim J, et al: **De novo assembly and phasing of a Korean human genome.** *Nature* 2016, **538:**243-247.

18. Durand NC, Robinson JT, Shamim MS, Machol I, Mesirov JP, Lander ES, Aiden EL: **Juicebox Provides a Visualization System for Hi-C Contact Maps with Unlimited Zoom.** *Cell Syst* 2016, **3:**99-101.

19. Rao SSP, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, Sanborn AL, Machol I, Omer AD, Lander ES, Aiden EL: **A 3D Map of the Human Genome at Kilobase Resolution Reveals Principles of Chromatin Looping.** *Cell* 2014, **159:**1665-1680.

20. Tamazian G, Dobrynin P, Krasheninnikova K, Komissarov A, Koepfli KP, O'Brien SJ:

486    **Chromosomer: a reference-based genome arrangement tool for producing draft**
487    **chromosome sequences.** *Gigascience* 2016, **5:**38.

488  21.  Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM: **BUSCO: assessing**
489    **genome assembly and annotation completeness with single-copy orthologs.** *Bioinformatics*
490    2015, **31:**3210-3212.

491  22.  Tarailo‐Graovac M, Chen N: **Using RepeatMasker to identify repetitive elements in**
492    **genomic sequences.** *Current protocols in bioinformatics* 2009, **25:**4.10. 11-14.10. 14.

493  23.  Bao W, Kojima KK, Kohany O: **Repbase Update, a database of repetitive elements in**
494    **eukaryotic genomes.** *Mobile Dna* 2015, **6:**11.

495  24.  Korf I: **Gene finding in novel genomes.** *BMC bioinformatics* 2004, **5:**59.

496  25.  Stanke M, Schöffmann O, Morgenstern B, Waack S: **Gene prediction in eukaryotes with a**
497    **generalized hidden Markov model that uses hints from external sources.** *BMC*
498    *bioinformatics* 2006, **7:**62.

499  26.  Majoros WH, Pertea M, Salzberg SL: **TigrScan and GlimmerHMM: two open source ab**
500    **initio eukaryotic gene-finders.** *Bioinformatics* 2004, **20:**2878-2879.

501  27.  Guigo R: **Assembling genes from predicted exons in linear time with dynamic**
502    **programming.** *Journal of Computational Biology* 1998, **5:**681-702.

503  28.  Kim D, Langmead B, Salzberg SL: **HISAT: a fast spliced aligner with low memory**
504    **requirements.** *Nat Methods* 2015, **12:**357-360.

505  29.  Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, White O, Buell CR, Wortman JR:
506    **Automated eukaryotic gene structure annotation using EVidenceModeler and the**
507    **Program to Assemble Spliced Alignments.** *Genome biology* 2008, **9:**R7.

508  30.  Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell
509    A, Nuka G: **InterProScan 5: genome-scale protein function classification.** *Bioinformatics*
510    2014, **30:**1236-1240.

511  31.  Kanehisa M, Goto S, Sato Y, Kawashima M, Furumichi M, Tanabe M: **Data, information,**
512    **knowledge and principle: back to metabolism in KEGG.** *Nucleic Acids Res* 2014, **42:**D199-
513    205.

514  32.  Lowe TM, Eddy SR: **tRNAscan-SE: a program for improved detection of transfer RNA**
515    **genes in genomic sequence.** *Nucleic Acids Res* 1997, **25:**955-964.

516  33.  Kalvari I, Nawrocki EP, Argasinska J, Quinones-Olvera N, Finn RD, Bateman A, Petrov AI:
517    **Non-Coding RNA Analysis Using the Rfam Database.** *Curr Protoc Bioinformatics* 2018,
518    **62:**e51.

519  34.  Nawrocki EP, Eddy SR: **Infernal 1.1: 100-fold faster RNA homology searches.**
520    *Bioinformatics* 2013, **29:**2933-2935.

521  35.  Yang Y, Zhou R, Zhu S, Li X, Li H, Yu H, Li K: **Systematic Identification and Molecular**
522    **Characteristics of Long Noncoding RNAs in Pig Tissues.** *Biomed Res Int* 2017,
523    **2017:**6152582.

524  36.  Liang G, Yang Y, Niu G, Tang Z, Li K: **Genome-wide profiling of Sus scrofa circular RNAs**
525    **across nine organs and three developmental stages.** *DNA research* 2017, **24:**523-535.

526  37.  Li H, Coghlan A, Ruan J, Coin LJ, Heriche JK, Osmotherly L, Li R, Liu T, Zhang Z, Bolund L,
527    et al: **TreeFam: a curated database of phylogenetic trees of animal gene families.** *Nucleic*
528    *Acids Res* 2006, **34:**D572-580.

529  38.  Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, Salzberg SL: **Versatile**

530          **and open software for comparing large genomes.** *Genome Biol* 2004, **5:**R12.

531  39.  Nattestad M, Schatz MC: **Assemblytics: a web analytics tool for the detection of variants**
532          **from an assembly.** *Bioinformatics* 2016, **32:**3021-3023.

533  40.  Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O: **New algorithms and**
534          **methods to estimate maximum-likelihood phylogenies: assessing the performance of**
535          **PhyML 3.0.** *Syst Biol* 2010, **59:**307-321.

536  41.  Yang Z: **PAML 4: phylogenetic analysis by maximum likelihood.** *Mol Biol Evol* 2007,
537          **24:**1586-1591.

538  42.  Kumar S, Stecher G, Suleski M, Hedges SB: **TimeTree: A Resource for Timelines, Timetrees,**
539          **and Divergence Times.** *Mol Biol Evol* 2017, **34:**1812-1819.

540  43.  Frantz LAF, Schraiber JG, Madsen O, Megens HJ, Bosse M, Paudel Y, Semiadi G, Meijaard E,
541          Li N, Crooijmans RPMA, et al: **Genome sequencing reveals fine scale diversification and**
542          **reticulation history during speciation in Sus.** *Genome Biology* 2013, **14.**

543  44.  Frantz LA, Madsen O, Megens HJ, Schraiber JG, Paudel Y, Bosse M, Crooijmans RP, Larson
544          G, Groenen MA: **Evolution of Tibetan wild boars.** *Nat Genet* 2015, **47:**188-189.

545  45.  Groenen MA: **A decade of pig genome sequencing: a window on pig domestication and**
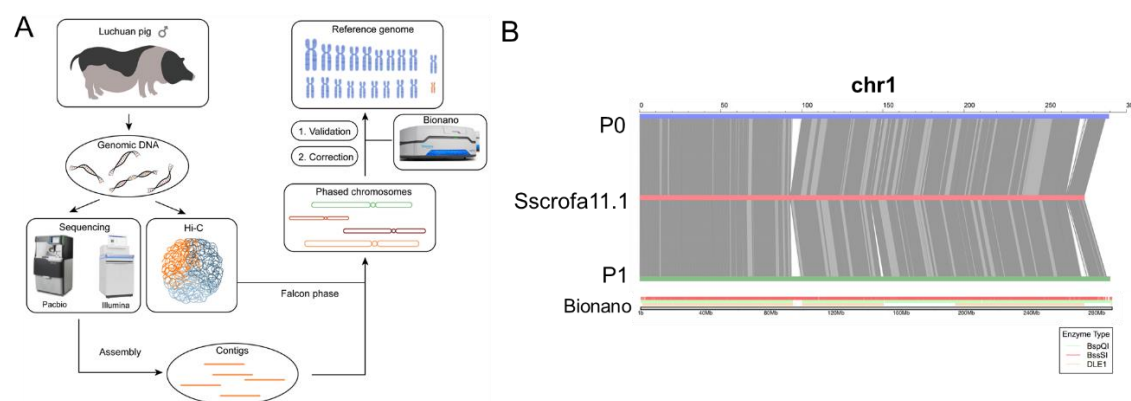546          **evolution.** *Genet Sel Evol* 2016, **48:**23.

547  46.  Nasvall J, Sun L, Roth JR, Andersson DI: **Real-time evolution of new genes by innovation,**
548          **amplification, and divergence.** *Science* 2012, **338:**384-387.

549  47.  Nowoshilow S, Schloissnig S, Fei JF, Dahl A, Pang AWC, Pippel M, Winkler S, Hastie AR,
550          Young G, Roscito JG, et al: **The axolotl genome and the evolution of key tissue formation**
551          **regulators.** *Nature* 2018, **554:**50-55.

552

553

554



555

556 **Figure 1. Genome Assembly.** (A) The flowchart of contig, scaffold and chromosome
557 assembly in this study. (B) Collinearity analysis for Chr1 between Sscrofa11.1 (*Middle*)
558 and primary assembly (P0, *Upper*) and alternate haplotigs (P1, *Lower*) assemblies. Gray
559 lines indicate collinearity between the genomes. Bionano optical map of Chr1 is shown
560 in the bottom. Collinearity Analysis for other chromosomes were shown in
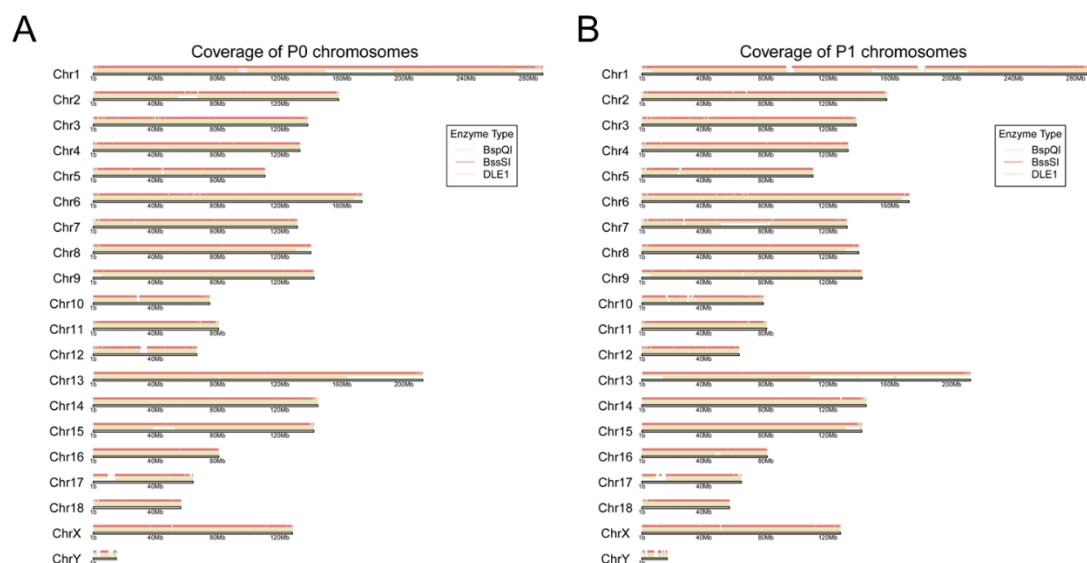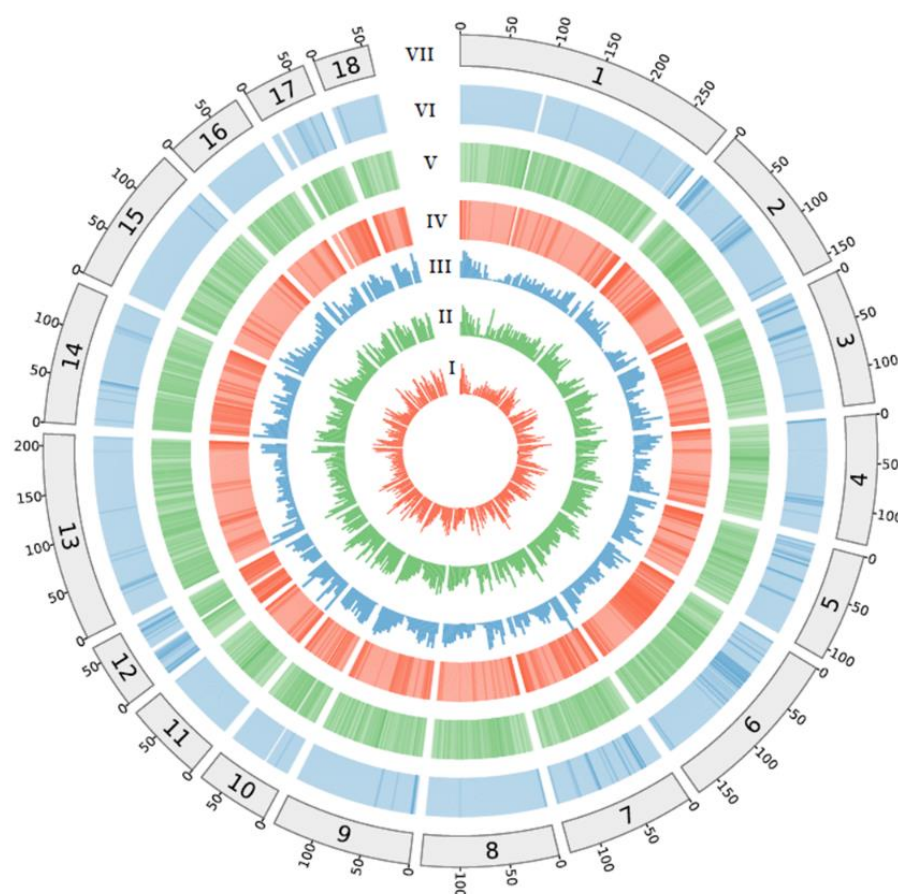561 Supplementary Figure 2.

562

**Figure 2. Assembly quality assessment by BioNano optical map data.** (A) A comparison of Bionano optical maps and primary assembly of Luchuan pig. (B) A comparison of Bionano optical maps and alternate haplotigs of Luchuan pig. The optical genome maps are constructed by three enzymes (Nt.BspQI, Nb.BssSI and DLE-1) and shown in different colors. The black bar corresponds to the pseudo-chromosomes of Luchuan pig.

571



572

**Figure 3. Circos plot showing the characterization of the Luchuan pig.**

I: Number of SNPs between primary assembly (P0) and alternate haplotigs (P1) in non-overlapping 5Mb windows;

II: Number of indels between primary assembly (P0) and alternate haplotigs (P1) in non-overlapping 5Mb windows;

III: Number of structural variants between primary assembly (P0) and alternate haplotigs (P1) in non-overlapping 5Mb windows;

IV: GC content in non-overlapping 1Mb windows;

V: Percent coverage of TEs in non-overlapping 1Mb windows;

VI: Gene density calculated on the basis of the number of genes in non-overlapping 1Mb windows;

VII: The length of pseudo-chromosome in the size of Mb.
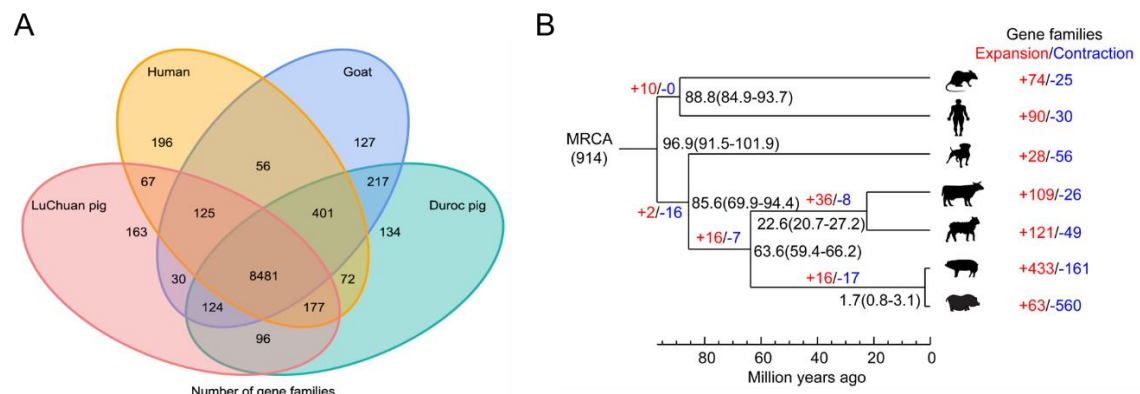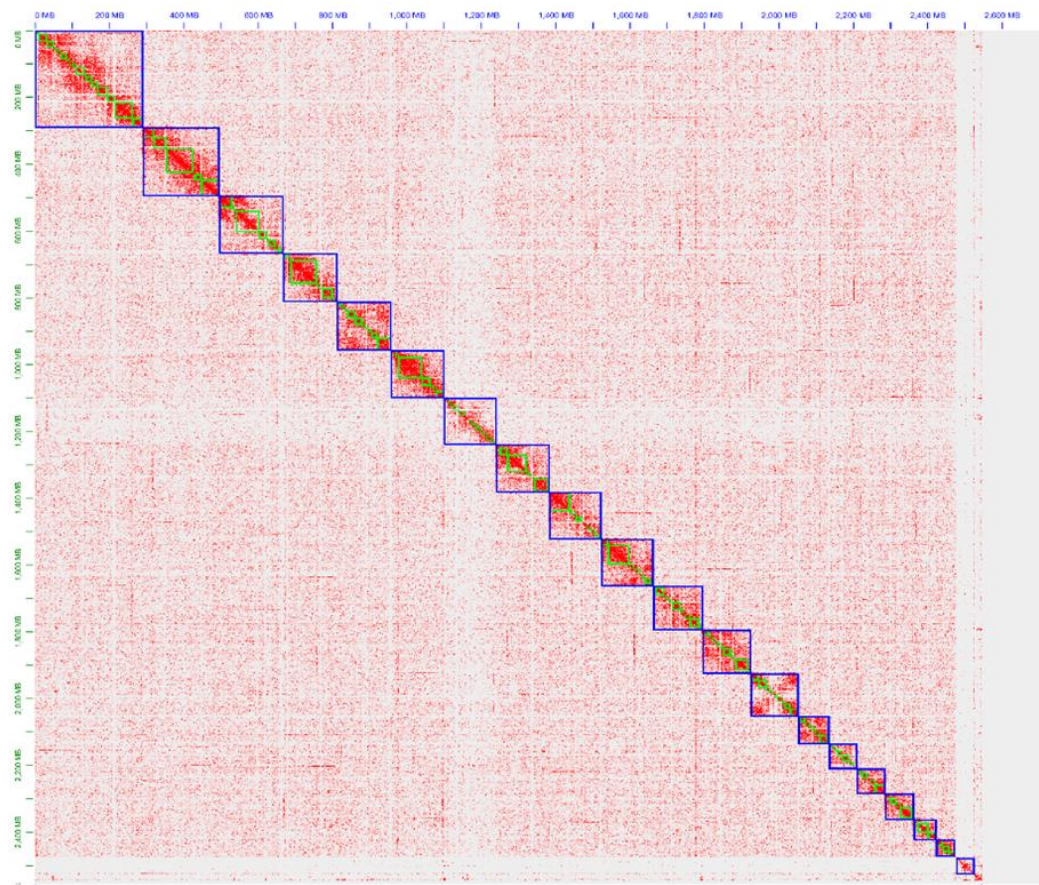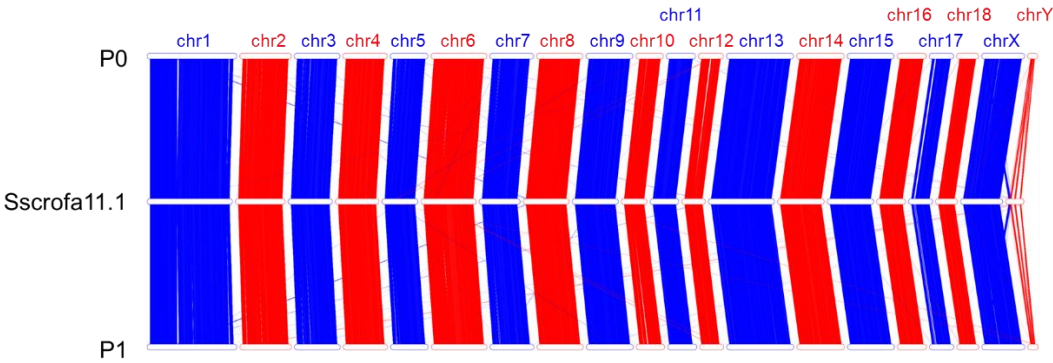
585

586

**Figure 4. Comparative genomic and phylogenetic analyses**. (A) Venn diagram showing shared orthologous gene families among genomes of Luchuan, Duroc, goat and human. (B) Phylogenetic tree with divergence times and history of orthologous gene families. Numbers on the nodes represent divergence times, with the error range shown in parentheses. The numbers of gene families that expanded (red) or contracted (blue) in each lineage after speciation are shown on the corresponding branch. MRCA, most recent common ancestor.

595

**Supplementary Figure 1. HiC contact heatmap.** Genome-wide analysis of chromatin interactions in Luchuan genome.

599

**Supplementary Figure 2. Collinearity analysis between Sscrofa11.1 (*Middle*) and primary assembly (P0, *Upper*) and alternate haplotigs (P1, *Lower*) assemblies.** Red and blue lines indicate collinearity between the genomes.

603