



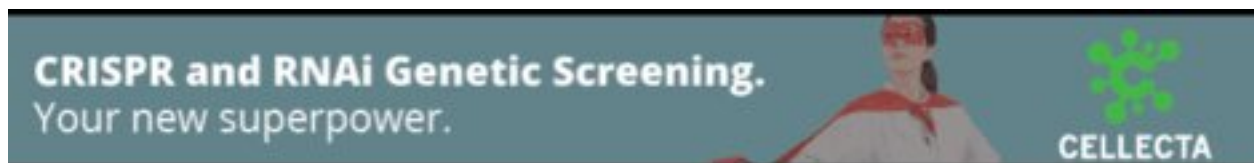
## Comprehensive variation discovery and recovery of missing sequence in the pig genome using multiple de novo assemblies

Mingzhou Li, Lei Chen, Shilin Tian, et al.

*Genome Res.* published online September 19, 2016  
Access the most recent version at doi:[10.1101/gr.207456.116](https://doi.org/10.1101/gr.207456.116)

---

<b>P&lt;P</b>	Published online September 19, 2016 in advance of the print journal.
<b>Accepted Manuscript</b>	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
<b>Open Access</b>	Freely available online through the <i>Genome Research</i> Open Access option.
<b>Creative Commons License</b>	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at <a href="http://creativecommons.org/licenses/by-nc/4.0/">http://creativecommons.org/licenses/by-nc/4.0/</a> .
<b>Email Alerting Service</b>	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or <a href="#">click here</a> .



---

To subscribe to *Genome Research* go to:  
<https://genome.cshlp.org/subscriptions>

---

Published by Cold Spring Harbor Laboratory Press

---

1 **Comprehensive Variation Discovery and Recovery of Missing**  
2 **Sequence in the Pig Genome using Multiple *De Novo***  
3 **Assemblies**

4 Mingzhou Li,<sup>1,9</sup> Lei Chen,<sup>2,9</sup> Shilin Tian,<sup>1,3,9</sup> Yu Lin,<sup>3,9</sup> Qianzi Tang,<sup>1,9</sup> Xuming Zhou,<sup>4,9</sup>  
5 Diyan Li,<sup>1</sup> Carol KL Yeung,<sup>3</sup> Tiandong Che,<sup>1</sup> Long Jin,<sup>1</sup> Yuhua Fu,<sup>1,5</sup> Jideng Ma,<sup>1</sup> Xun  
6 Wang,<sup>1</sup> Anan Jiang,<sup>1</sup> Jing Lan,<sup>2</sup> Qi Pan,<sup>3</sup> Yingkai Liu,<sup>1</sup> Zonggang Luo,<sup>2</sup> Zongyi Guo,<sup>2</sup>  
7 Haifeng Liu,<sup>1</sup> Li Zhu,<sup>1</sup> Surong Shuai,<sup>1</sup> Guoqing Tang,<sup>1</sup> Jiugang Zhao,<sup>2</sup> Yanzhi Jiang,<sup>1</sup> Lin  
8 Bai,<sup>1</sup> Shunhua Zhang,<sup>1</sup> Miaomiao Mai,<sup>1</sup> Changchun Li,<sup>5</sup> Dawei Wang,<sup>3</sup> Yiren Gu,<sup>6</sup>  
9 Guosong Wang,<sup>1,7</sup> Hongfeng Lu,<sup>3</sup> Yan Li,<sup>3</sup> Haihao Zhu,<sup>3</sup> Zongwen Li,<sup>3</sup> Ming Li,<sup>8</sup> Vadim N.  
10 Gladyshev,<sup>4</sup> Zhi Jiang,<sup>3</sup> Shuhong Zhao,<sup>5</sup> Jinyong Wang,<sup>2</sup> Ruiqiang Li,<sup>3</sup> and Xuewei Li<sup>1</sup>

11 <sup>1</sup> Institute of Animal Genetics and Breeding, College of Animal Science and Technology,  
12 Sichuan Agricultural University, Chengdu 611130, China;

13 <sup>2</sup> Key Laboratory of Pig Industry Sciences (Ministry of Agriculture), Chongqing Academy of  
14 Animal Sciences, Chongqing 402460, China;

15 <sup>3</sup> Novogene Bioinformatics Institute, Beijing 100089, China;

16 <sup>4</sup> Division of Genetics, Department of Medicine, Brigham and Women's Hospital, Harvard  
17 Medical School, Boston, Massachusetts, 02115 USA;

18 <sup>5</sup> College of Animal Science and Technology, Huazhong Agricultural University, Wuhan  
19 430070, China;

20 <sup>6</sup> Sichuan Animal Science Academy, Chengdu 610066, China;

21 <sup>7</sup> Department of Animal Science, Texas A & M University, College Station, Texas, 77843  
22 USA;

23 <sup>8</sup> Key Laboratory of Animal Ecology and Conservation Biology, Institute of Zoology,  
24 Chinese Academy of Sciences, Beijing 100101, China.

25 <sup>9</sup>These authors contributed equally to this work.

26 Corresponding authors: [mingzhou.li@sicau.edu.cn](mailto:mingzhou.li@sicau.edu.cn), [kingyou@vip.sina.com](mailto:kingyou@vip.sina.com),  
27 [lirq@novogene.cn](mailto:lirq@novogene.cn), [xuewei.li@sicau.edu.cn](mailto:xuewei.li@sicau.edu.cn)

28

29

30

## Abstract

31 Uncovering genetic variation through resequencing is limited by the fact that  
32 only sequences with similarity to the reference genome are examined.  
33 Reference genomes are often incomplete and cannot represent the full range  
34 of genetic diversity as a result of geographical divergence and independent  
35 demographic events. To more comprehensively characterize genetic variation  
36 of pigs (*Sus scrofa*), we generated *de novo* assemblies of nine geographically  
37 and phenotypically representative pigs from Eurasia. By comparing them to  
38 the reference pig assembly, we uncovered a substantial number of novel  
39 SNPs, structural variations, as well as 137.02 Mb sequences harboring 1,737  
40 protein coding genes that were absent in the reference assembly, revealing  
41 variants left by selection. Our results illustrate the power of whole-genome *de*  
42 *novo* sequencing relative to resequencing, and provide valuable genetic  
43 resources that enable effective use of pigs in both agricultural production and  
44 biomedical research.

45

46

## Introduction

47 *Sus scrofa* (i.e., pig or swine) is of enormous agricultural importance and also  
48 an attractive model for biomedical research and applications. There are over  
49 730 distinct pig breeds worldwide, of which two thirds reside in Europe and  
50 China ([Chen et al., 2007](#)), whose diverse phenotypes are shaped by the  
51 combined effects of local adaptation and artificial selection ([Ai et al., 2015](#)).  
52 Efforts have been made to characterize the genetic variation that underlies this

---

53 phenotypic diversity using resequencing data and the genome of the European  
54 domestic Duroc pig as a reference (Choi et al., 2015; Groenen et al., 2012;  
55 Moon et al., 2015; Rubin et al., 2012). Nonetheless, resequencing is limiting in  
56 terms of capturing genetic variation and assessing gaps and misassigned  
57 regions of the reference genome (Weisenfeld et al., 2014). In contrast, multiple  
58 *de novo* assemblies of pig genomes from different regions and breeds promise  
59 a more accurate and comprehensive understanding of genetic variation within  
60 this species (Besenbacher et al., 2015; Chaisson et al., 2015b). Among  
61 populations of plants (coccolithophores (Read et al., 2013), *Arabidopsis*  
62 *thaliana* (Gan et al., 2011), soybean (Li et al., 2014) and rice (Zhang et al.,  
63 2014)), animals (mosquitoes (Neafsey et al., 2015) and macaques (Yan et al.,  
64 2011)) and even modern humans (Li et al., 2010a), surprisingly large amount  
65 of variation has been uncovered by *de novo* assemblies.

66 To advance the characterization of the genetic diversity of pigs, we  
67 generated *de novo* assemblies of nine geographically and phenotypically  
68 representative individuals from Eurasia. Combining this resource with genome  
69 assembly of the Tibetan wild boar (Li et al., 2013), we carried out in-depth  
70 comparisons between ten *de novo* assemblies and the reference genome. We  
71 uncovered a substantial number of single nucleotide polymorphisms (SNPs)  
72 and structural variations, as well as hundreds of millions of base pairs that are  
73 not present in the reference genome, including thousands of protein-coding  
74 genes that are either missing or fragmented in the reference genome, which  
75 contain potentially important genetic information pertaining to porcine  
76 evolution.

77

78

## Results

### 79 ***De novo* genome assemblies of nine pig breeds**

80 We sequenced the genomes of nine female individuals from nine diverse  
81 breeds (five originated in Europe and four originated in China) to an average of  
82 ~100-fold coverage (~229.5 gigabase (Gb)) using Illumina sequencing  
83 technology and a whole-genome shotgun strategy (**Supplemental Fig. S1** and  
84 **Table S1**). The genomes were independently assembled using SOAPdenovo  
85 ([Li et al., 2010b](#)) (**Supplemental Methods**), which yielded contig N50 sizes of  
86 28.99 to 42.66 kilobases (kb), scaffold N50 sizes of 1.26 to 2.45 megabases  
87 (Mb) and a total of 2.45 to 2.49 Gb of ungapped sequences that exhibited  
88 genomic features similar to those of the reference assembly ([Groenen et al.,](#)  
89 [2012](#)) (**Supplemental Figs. S2-S7** and **Tables S2-S6**). We also improved the  
90 available genome assembly of the Tibetan wild boar ([Li et al., 2013](#)) by  
91 increasing the contig N50 size from 20.69 kb to 22.54 kb and the ungapped  
92 genome assembly size from 2.43 Gb to 2.44 Gb (**Supplemental Table S4**).

### 93 **Discovery and characterization of SNPs**

94 We identified 8.86-15.95 million (M) SNPs in individual pig genomes using an  
95 assembly-versus-assembly approach (**Supplemental Methods**). These SNPs  
96 were consistent with more than 98% SNPs identified from the Illumina's  
97 porcine 60K Genotyping BeadChip (v.2) (**Supplemental Table S7**), and  
98 covered most SNPs identified by resequencing as implemented in SAMtools  
99 ([Li et al., 2009](#)) (98.78%) and GATK tool ([McKenna et al., 2010](#)) (97.65%) and

---

100 3.12-5.40 M SNPs (33.46-35.25%) in divergent regions that failed to be  
101 cataloged by these algorithms, where unassembled short reads are difficult to  
102 be mapped (**Fig. 1** and **Supplemental Figs. S8,9**).

103 Extensive inter-continental genomic divergence was reflected by the  
104 significantly larger amount of variation of Chinese pigs when compared to the  
105 reference Duroc genome of European origin (15.14 to 15.95 M SNPs; the Ts /  
106 Tv ratio: 2.13 to 2.15) than that between European pigs and Duroc (8.86 to  
107 10.14 M SNPs; the Ts / Tv ratio: 1.95 to 1.99) ( $P < 10^{-16}$ , Mann-Whitney *U* test)  
108 (**Figs. 2A,B** and **Supplemental Fig. S10**), attributable to considerable  
109 divergence time between European and Asian lineages (at least 1 million  
110 years) and their independent domestication in multiple locations across  
111 Eurasia in the past ~10,000 years ([Frantz et al., 2013](#); [Groenen et al., 2012](#);  
112 [Larson et al., 2005](#)).

113 We also observed higher genomic diversity of Chinese pigs than European  
114 pigs, reflected by the former's higher heterozygous SNP ratio ( $2.17 \times 10^{-3}$  to  
115  $2.69 \times 10^{-3}$  vs.  $0.94 \times 10^{-3}$  to  $1.63 \times 10^{-3}$ ) and lower homozygosity (382 regions  
116 of homozygosity (ROHs) with a total size of 107.5 Mb vs. 907 ROHs with  
117 totaling 289.9 Mb per assembly) ( $P < 10^{-16}$ , Mann-Whitney *U* test) (**Fig. 2B** and  
118 **Supplemental Figs. S11,12**). Principal component analysis (PCA) and identity  
119 score (IS) analysis of pairwise breed genomes also recapitulated these  
120 findings (**Fig. 2C** and **Supplemental Fig. S13A**). This may be a reflection of  
121 the fact that European origin breeds have undergone intense selection in  
122 inbred commercial lines for economical traits, while Chinese breeds  
123 experienced weaker selection in scattered, individual farms, and exhibited

---

124 relatively weak linkage disequilibrium (LD) ([White, 2011](#)) ([Supplemental Fig.](#)  
125 [S13B](#)). Another possible explanation is that European wild boars (ancestors of  
126 European domestic pigs) may have suffered more pronounced population  
127 bottlenecks during the last glacial maximum (~20,000 years ago) compared to  
128 their Asian counterparts ([Bosse et al., 2012](#); [Groenen et al., 2012](#)).

129 We pooled the SNPs of ten breeds into a non-redundant set of 33.60 M sites  
130 that account for ~81.25 % of the estimated repertoire of SNPs in the pig  
131 ([Supplemental Fig. S14](#), [Table S8](#) and [Supplemental Methods](#)), of which  
132 6.34 M (18.87% of 33.60 M) SNPs were considered to be novel based on their  
133 absence in the pig dbSNP (Build 143) entries ([Supplemental Fig. S15](#)).  
134 Compared with synonymous SNPs (122.44 k), missense SNPs (83.39 k)  
135 exhibited greater diversity among breeds (77.44% of the estimated repertoire  
136 compared to 80.61%), accounted for a larger proportion of breed-specific (and  
137 thus rare) SNPs (32.42% compared to 30.18%) and had a higher ratio of  
138 homozygous to heterozygous SNPs (0.37 compared to 0.32) ([Supplemental](#)  
139 [Figs. S14,16](#)), which may be associated with breed-specific adaptation.

#### 140 **Maps of structural variations**

141 We detected 161.45-279.98 k insertions (15.99-23.07 Mb in length) and  
142 137.89-269.55 k deletions (3.61-5.63 Mb in length) in individual genomes  
143 against the reference genome ([Fig. 2B](#), [Supplemental Table S9](#) and  
144 [Supplemental Methods](#)). More than 80% of the insertions and deletions  
145 (indels) were 1-10 bp in length, and there was also a relatively high number of  
146 indels ~300 bp in length, due to the enrichment of indels of tRNA<sup>Glu</sup>-derived  
147 short interspersed element (SINE/tRNA<sup>Glu</sup>) ([Supplemental Fig. S17](#)) ([Ai et al.,](#)

---

148 [2015](#)). Repetitive elements (38.05% of the genome) comprised ~52.73% of  
149 indels, which are an important source of structural variation in pig genome.  
150 Moreover, the SINE/tRNA<sup>Glu</sup> (290.47 Mb and containing 18.09% of indels)  
151 showed higher incidence of indels than the predominant long interspersed  
152 elements (LINE/L1) (636.50 Mb and containing 15.48% of indels)  
153 ([Supplemental Fig. S18](#)).

154 The indels appeared to be regulated by selection: most indels were located  
155 in intergenic regions (72.20% to 74.14%), the indel ratio was lower in the  
156 coding sequences than in introns ([Supplemental Fig. S19](#)) and more  
157 conserved genes showed fewer structural variations ([Supplemental Fig. S20](#)).  
158 We observed an enrichment of short indels (1-15 bp in length) in coding  
159 sequences (414 of 1,582, or 26.17%) that were multiples of 3 bp, which is  
160 expected to preserve the reading frame, and identified 1,152 frameshift  
161 mutations in 947 genes ([Supplemental Fig. S21](#) and [Table S10](#)), which  
162 mainly represented the cellular functions of 'binding of nucleoside, ATP, and  
163 cation' and 'neuron development' ([Supplemental Table S11](#)). As the SNPs,  
164 distribution of indels across the genome also reflected a deep phylogenetic  
165 split between European and Chinese pigs and higher genetic variability of  
166 Chinese pigs than European pigs ([Bosse et al., 2012](#); [Groenen et al., 2012](#))  
167 ([Supplemental Fig. S22](#)).

### 168 **Signatures of diversifying selection in pig breeds**

169 To uncover genetic variation underlying phenotypic diversity of pigs, we  
170 identified breed-specific signatures left by diversifying selection using a relative  
171 homozygous SNP density (RSD) algorithm ([Atanur et al., 2013](#))

---

172 (**Supplemental Methods**). We identified 493 separate genomic regions of 20  
173 to 150 kb (a total of 20.10 Mb and containing 308 genes) to be under selection  
174 (FDR < 0.05) (**Fig. 3A** and **Supplemental Table S12**). These putative selected  
175 regions also exhibited significantly strong LD and lower negative Tajima's *D*  
176 values ( $P < 10^{-16}$ , Mann-Whitney *U* test) (**Supplemental Fig. S23**), and distinct  
177 phylogenetic relationships compared to genomic background (**Supplemental**  
178 **Fig. S24**).

179 Most homozygous SNPs (88.60%) in the selected regions were unique to a  
180 particular breed (**Fig. 3A**), exhibiting lower degree of haplotype sharing with  
181 other breeds than pairwise between other breeds (**Fig. 3B** and **Supplemental**  
182 **Fig. S25**). These private SNPs were highly concentrated in a small number of  
183 discrete genomic regions (0.79% of the genome), and may be associated with  
184 phenotypes described by standard breed criteria ([Wang et al., 2011](#)): typically,  
185 9 (out of 49, or 18.37%;  $P = 0.004$ ,  $\chi^2$  test) and 6 (out of 59, or 10.17%;  $P =$   
186  $0.491$ ,  $\chi^2$  test) genes within or in the vicinity of the selected regions in the fatty  
187 Rongchang and Jinhua pigs were orthologous to well-established mammalian  
188 fat deposition genes ([Kunej et al., 2013](#)) (**Fig. 3B** and **Supplemental Fig.**  
189 **S25A**), including factors involved in the regulation of feed intake and energy  
190 homeostasis (*CEP120*, *GABRA2*, *NPPA*, *NPY1R* and *NYP5R*), lipid  
191 metabolism (*ABCC4*, *ANGPT2*, *LRPAP1* and *PRKAG2*) and indicators of  
192 obesity-induced hypertension, inflammatory signaling and insulin resistance  
193 (*ADD1*, *HSPD1*, *MMP2*, *PIK3R4*, *RAE1* and *TBCA*) (**Supplemental Table**  
194 **S13**). In contrast to highly inbred European pigs that have undergone selection  
195 for lean growth (high protein and low fat content; lean meat percentage of  
196 carcass ranging between 63-73%) as a response to demands for reduced

---

197 calorie intake in modern society, Chinese pigs have been selected for extreme  
198 fatness all along (typical lean meat percentage is under 45%) (**Supplemental**  
199 **Fig. S1**), driven by demand for energy-rich food in developing countries until  
200 ~10 years ago ([Wang et al., 2011](#)).

201 We also identified 16 (31.37%;  $P = 8.21 \times 10^{-11}$ ,  $\chi^2$  test) out of 51 genes with  
202 strong selective sweep signals in the Tibetan wild boar (**Supplemental Fig.**  
203 **S25B** and **Table S13**) that were likely driven by the harsh and hypoxic  
204 environment of the Tibetan plateau and might have a role in the formation of  
205 characteristic phenotypes, such as an insulating layer formed by hard skin and  
206 long, dense hair, and larger lungs and hearts ([Li et al., 2013](#)).

#### 207 **Identifying missing sequences of the reference pig genome**

208 There are considerable unidentified regions (289.24 Mb of 2.81 Gb, or 10.29%)  
209 in the reference pig assembly (Sscrofa10.2) ([Groenen et al., 2012](#)), of which  
210 266.15 Mb (91.92%) is composed of 5,317 gaps of at least 50 kb long  
211 (**Supplemental Fig. S26**). To recover such missing genetic information, we  
212 retrieved ~9.17 G ‘orphan reads’ for which neither end mapped to the  
213 reference genome (**Supplemental Fig. S27**), and re-localized them to their  
214 respective assemblies of origin. Consequently, we identified 83.8 k sequences  
215 of  $\geq 500$  bp (137.02 Mb in length) that were missing in the reference genome  
216 (**Table 1** and **Supplemental Table S14**). Only a small portion of missing  
217 sequences was considered to be insertions (~0.91Mb) or copy number gains  
218 (~4.16%) (**Supplemental Tables S14,15** and **Supplemental Methods**).  
219 Compared with whole assemblies, these missing sequences exhibited similar  
220 heterozygous SNP ratio ( $2.67 \times 10^{-3}$  vs.  $2.56 \times 10^{-3}$ ;  $P = 0.623$ , Mann-Whitney

---

221 *U* test), but significantly higher GC content (43.07% vs. 41.41%;  $P < 10^{-16}$ ,  
222 Mann-Whitney *U* test) and repeat ratio (47.57% vs. 38.38%;  $P < 10^{-16}$ ,  
223 Mann-Whitney *U* test) (**Supplemental Fig. S28**).

224 Most sequences missing in the reference genome were common between  
225 different assemblies, as most orphan reads (95.04%) could crossly align to  
226 missing sequences of other assemblies with coverage (97.10% with depth  $\geq 4$   
227 per base) comparable to that against their respective assemblies (mapping  
228 ratio = 95.83% and coverage = 99.51%) (**Supplemental Fig. S29**). Pairwise  
229 similarity of the orphan reads and the missing sequences between ten breeds  
230 revealed a clear distinction between European and Asian pigs, as well as  
231 relatively high genetic variability in Chinese pigs than in European pigs ([Bosse  
232 et al., 2012; Groenen et al., 2012](#)) (**Supplemental Fig. S30**), suggesting that  
233 these common sequences, which were absent in the reference genome, may  
234 be important sources of pig diversity and contain biologically meaningful  
235 information.

236 We were also able to fill in 71.37% (3,795 of 5,317) of the gaps in the  
237 reference genome by missing sequences of at least one breed (**Supplemental  
238 Fig. S27, Tables S14,16 and Supplemental Methods**). These filled missing  
239 sequences were highly collinear across ten breeds and exhibited similar  
240 distribution over the reference assembly gaps (average of pairwise Person's  $r$   
241 = 0.89,  $P < 10^{-16}$ ) (**Supplemental Fig. S31**). Typical examples are shown in  
242 **Supplemental Fig. S32**.

### 243 **Recovery of missing genes**

244 Of the average 20,782 protein-coding genes (87.13% were supported by

---

245 evidence of transcription) predicted in each of the ten assemblies  
246 (**Supplemental Figs. S33-S35**, **Table S17** and **Supplemental Methods**), we  
247 found an average of 1,096 (5.27%) genes to be embedded or almost  
248 completely contained (> 50% to > 90% overlap of gene length, respectively) in  
249 the missing sequences of the reference assembly (**Table 1** and **Supplemental**  
250 **Table S18**), which we referred to as ‘missing genes’ ([Genovese et al., 2013](#);  
251 [Kidd et al., 2010](#)).

252 To check whether these predicted missing genes are likely to be functional,  
253 we compared their conservation level across 19 mammalian genomes and  
254 found that they generally exhibited similar identity (81.55% vs. 83.60%) and  
255 coverage (96.32% vs. 97.37%) as annotated genes (**Supplemental Fig. S36**).  
256 Coding sequences of missing genes were enriched at high cross-species  
257 (human, cow and sheep) identity level (> 90%), also consistent with the  
258 sequence identity distribution of well-annotated coding sequences of the  
259 reference genome (**Supplemental Fig. S37**). We then retrieved ~0.59 G  
260 orphan reads against the reference genome from each of 96 pair-end  
261 RNA-seq libraries (7 to 10 libraries for each breed) and mapped them onto the  
262 missing genes in their respective assemblies (**Supplemental Figs. S38A,B**).  
263 Consequently, an average of 91.51% (1,003 of 1,096) missing genes in each  
264 assembly showed  $\log_2$ -transformed FPKM expression values (denoted as  
265 fragments per kb of transcript per Mb orphan reads) greater than 0.3 in at least  
266 one library (**Supplemental Fig. S38C**), suggesting that a considerable number  
267 of missing genes are functional and biologically important.

268 To determine the collinear relationships of missing genes among ten breeds,  
269 we separately aligned the protein sequences of nine assemblies to the

---

270 assembly of the Large White breed, which had the longest scaffold N50 size  
271 (2.45 Mb). Using MCScanX toolkit ([Wang et al., 2012](#)), we found that 10,313 of  
272 10,959 (94.10%) missing genes in all of ten assemblies belonged to 1,091  
273 inter-assembly collinear gene models, of which 871 (79.84%) models were  
274 present in all ten assemblies ([Table 1](#) and [Supplemental Tables S18,19](#)).

275 There were a total of 646 missing genes (14 to 95 per assembly) assembled in  
276 only a single breed, which could be found in other assemblies when orphan  
277 reads from short-insert (180 and 500 bp) libraries were used for mapping  
278 (coverage 94.05% at least 1 × depth), suggesting that the absence of these  
279 singleton genes from other assemblies is likely an artifact of fragmentation or  
280 misassignment in short reads assembly ([Alkan et al., 2011](#); [Chaisson et al.,](#)  
281 [2015b](#)) ([Supplemental Fig. S39](#)).

282 Together with the longest gene model of 1,091 inter-assembly collinear  
283 genes and 646 singleton genes, we obtained 1,737 missing gene models  
284 ([Table 1](#)). Aligning these missing genes to RefSeq proteins of pig, human, cow  
285 and mouse yielded 1,731 (99.65%) hits at least in one species ([Supplemental](#)  
286 [Table S19](#)), of which 359 (20.66%) missing genes could not be aligned to any  
287 known RefSeq proteins of pig, indicating that these genes have not been  
288 characterized in pig. Among hits that matched functionally classified proteins,  
289 the most abundant were members of olfactory receptors (65 hits,  $P = 1.60 \times$   
290  $10^{-12}$ ,  $\chi^2$  test), G-protein coupled receptors (104 hits,  $P = 9.81 \times 10^{-6}$ ,  $\chi^2$  test)  
291 and those involved in neurological system processes (112 hits,  $P = 4.26 \times 10^{-6}$ ,  
292  $\chi^2$  test) ([Supplemental Table S20](#)), which are known to be rapidly evolving  
293 between species ([Mainland et al., 2014](#)). We also recovered genes  
294 corresponding to economically important traits that are valuable for future

---

295 functional analyses and improvements of pig as an important livestock species,  
296 such as genes related to pork production (74 of 1,515 fat deposition genes  
297 (Kunej et al., 2013), or 4.88%) and disease resistance (76 of 1,517 genes  
298 annotated with the GO: 0002376; immune system process, or 5.01%)  
299 (**Supplemental Table S19**). Typical examples are shown in **Supplemental**  
300 **Fig. S40**.

### 301 **Selection in missing genes**

302 To reveal variants left by selection, we measured pairwise the extent of  
303 population differentiation of the coding SNPs in the missing genes between  
304 Chinese wild boars (32.57 k coding SNPs) and seven Chinese domestic  
305 populations (23.02 k coding SNPs per population) using the *FDIST* approach  
306 as implemented in Arlequin (Excoffier and Lischer, 2010) (**Supplemental Figs.**  
307 **S41,42** and **Table S21**). A total of 605 non-redundant coding SNPs embedded  
308 in 328 missing genes were found to be under directional selection in seven  
309 Chinese domestic populations (FDR < 0.05, *FDIST* test) (**Supplemental Fig.**  
310 **S43** and **Table S22**), which also exhibited significantly lower  $q$  values in a  
311 Bayesian test (Foll and Gaggiotti, 2008) and  $F_{ST}$  values in a ‘model-free’ global  
312  $F_{ST}$  test when compared to other unselected loci ( $P < 10^{-16}$ , Mann-Whitney  $U$   
313 test) (**Supplemental Fig. S44**). The missing genes under selection in seven  
314 Chinese domestic populations were commonly enriched for biological  
315 processes related to ‘binding of actin, calcium ion, and cytoskeletal protein’  
316 (**Supplemental Fig. S45A**). Intriguingly, 71 genes harboring 110 selected  
317 coding SNPs of domestic Erhualian pigs (one of the most prolific pig breeds  
318 known) (Wang et al., 2011) belonged predominantly to fertility-related

---

319 categories, such as ‘sexual reproduction’ (7 genes: *ADAM20*, *AKT1*, *GMCL1*,  
320 *MICALCL*, *NOTCH1*, *SPIN4* and *SPTBN4*;  $P = 0.001$ ) and ‘placenta  
321 development’ (3 genes: *AKT1*, *RXRA* and *VWF*;  $P = 0.012$ ), which may  
322 underlie the breed’s markedly larger litters (~3 to 5 more piglets per litter)  
323 (**Supplemental Fig. 45B**).

324 The expression of missing genes under selection also showed remarkably  
325 higher tissue specificity, reflected by the lower Shannon entropy ( $H$ ) values (a  
326 measure of the specificity of gene expression across tissues) (**Schug et al.,**  
327 **2005**) than unselected missing genes (1.98 vs. 2.37 per gene;  $P < 10^{-16}$ ,  
328 Mann-Whitney  $U$  test) (**Supplemental Fig. S46**). As opposed to constitutive  
329 genes that are ubiquitously expressed and essential for basic cellular functions,  
330 tissue-specific genes are usually associated with the development of generally  
331 desirable traits, such as disease resistance, muscle growth, fat deposition and  
332 reproduction, and thus are more likely prone to be shaped by selection.

333 None of the selected coding SNPs was a nonsense mutation (resulting in  
334 premature stop codons in transcribed mRNAs) (**Supplemental Table S22**),  
335 supporting the idea that gene inactivation did not play a prominent role during  
336 pig domestication and consistent with the results from screens in chickens  
337 (**Rubin et al., 2010**), rabbits (**Carneiro et al., 2014**) and pigs (based on  
338 reference genome) (**Rubin et al., 2012**). Compared to synonymous  
339 substitutions, missense substitutions showed significantly lower genetic  
340 differentiation (global  $F_{st}$ , 0.05 compared to 0.10 per locus;  $P < 10^{-16}$ ,  
341 Mann-Whitney  $U$  test) between Chinese wild boars and domestic pigs  
342 (**Supplemental Fig. S47**). Nonetheless, there were still 127 genes harboring

343 selected missense mutations, which were over-represented in the highly  
344 variable olfactory receptor family (12 genes;  $P = 0.02$ ,  $\chi^2$  test)<sup>34</sup>  
345 (**Supplemental Table S22**). Of these, three missense mutations embedded in  
346 two genes related to the development of obesity were of interest: the closely  
347 linked Asn566-His (T1,696-G) and Ser578-Cys (G1,733-C) substitutions ( $D' =$   
348  $1$ ,  $r^2 = 0.975$ ) found in *ALPK3* (alpha kinase 3) (**Fig. 4** and **Supplemental Fig.**  
349 **S48**), and a Thr18-Ile (C53-T) substitution in *PKD1L2* (polycystin 1 like 2  
350 (gene/pseudogene)) (**Supplemental Fig. S49**). These three missense  
351 mutations exhibited significant selection signals ( $FDR < 0.05$ , *FDIST* test)  
352 between Chinese wild boars and one of seven domestic populations (Min and  
353 Erhualian, respectively), but were nearly fixed in the more genetically  
354 homogenous European/North American domestic pigs, possibly as a result of  
355 stronger selective pressure in western societies, although larger sample sizes,  
356 inter-continental genetic discrepancy of pig genomes and functional analyses  
357 are required to validate the non-neutrality of these genes.

358 *ALPK3* plays a role in cardiomyocyte differentiation; knockout of this gene in  
359 mice was associated with marked hypertrophic and dilated forms of  
360 cardiomyopathy ([Van Sligtenhorst et al., 2012](#)). *ALPK3* shows the strongest  
361 evidence of positive selection in the polar bear, which has a lipid-rich diet  
362 throughout life ([Liu et al., 2014](#)). Selection of *ALPK3* in domestic pigs suggests  
363 that potential protection against the chronically deleterious effects of a  
364 'diabetogenic' environment (high calorie, atherogenic diet and little physical  
365 exercise) on the cardiovascular system may be favorable ([Gerstein and](#)  
366 [Waltman, 2006](#); [Koopmans and Schuurman, 2015](#)). *PKD1L2* is primarily  
367 associated with fatty acid synthase in the skeletal muscle fiber; its

---

368 overexpression in mice provokes myofiber atrophy and suppressed  
369 lipogenesis ([Mackenzie et al., 2009](#)). The triglycerides accumulated between  
370 or within myofibers represent a large energy source (contributes up to 20% of  
371 total energy turnover during physical exercise in human) ([Roepstorff et al.,](#)  
372 [2005](#)). Selection of *PKD1L2* is likely related to the relatively weak athletic  
373 performance of domestic pigs compared to wild boars due to limited active  
374 space in pig farms.

375

376

## Discussion

377 We describe an assembly-versus-assembly approach that relies on multiple  
378 independently assembled genomes for improving the power of variant  
379 detection, as opposed to the currently dominant resequencing approach. This  
380 catalogue of variants, including SNPs, indels, and common and rare variants is  
381 a valuable resource for further investigation of the genetic makeup of porcine  
382 phenotypic diversity and adaptive evolution. We show that high-quality *de*  
383 *novo* assembly of individual genomes followed by comparison with the  
384 reference sequence is necessary for identifying novel genetic variation across  
385 geographical ranges and different evolutionary histories. Such experimental  
386 design is increasingly affordable with the advances in sequencing technology  
387 ([Zook and Salit, 2015](#)), especially long-read sequencing ([Chaisson et al.,](#)  
388 [2015a](#)) and single-molecule mapping ([Koren et al., 2012](#)) technologies.

389 Interpretation of the consequences of genetic variation has typically relied  
390 on reference sequences, relative to which genes and variants are annotated  
391 and examined. However, we recovered hundreds of millions of base pairs that

---

392 were not present in the pig reference genome, including thousands of  
393 protein-coding genes that are either missing or fragmented in the reference  
394 genome, which harbor abundant variants associated with economic traits that  
395 are likely subjected to artificial selection. These newly recovered genes can  
396 now be incorporated into genotyping platforms and expression microarrays to  
397 facilitate their functional characterization. Recovered sequences missing from  
398 the reference genome could also be the source of genetic signals that have  
399 been ascertained by linkage, association and copy number variation studies  
400 but not yet been mapped to causal mutations.

401

402

## Methods

### 403 ***De Novo* sequencing and assembly of pig genomes**

404 We sequenced the genomes of nine geographically and phenotypically  
405 representative pig breeds using Illumina sequencing technology and a  
406 whole-genome shotgun strategy ([Fig. 2A](#) and [Supplemental Fig. S1](#)).  
407 Short-insert (180 and 500 bp) and long-insert (2, 5, 6 and 10 kb) DNA libraries  
408 were paired-end sequenced on the Illumina HiSeq 2500 platform  
409 ([Supplemental Fig. S2](#) and [Table S1](#)). We independently assembled nine  
410 genomes using SOAPdenovo ([Li et al., 2010b](#)), which is a *de Bruijn* graph  
411 algorithm-based *de novo* genome assembler ([Supplemental Methods](#)). We  
412 performed repeat annotation for ten breed assemblies and the reference  
413 genome using the same pipeline ([Supplemental Figs. S5,6](#) and  
414 [Supplemental Methods](#)).

### 415 **SNP and indel calling using an assembly-versus-assembly method**

---

416 We took advantage of an assembly-versus-assembly approach to identify  
417 candidate variations and further filtered out spurious variations by aligning  
418 short sequencing reads (**Supplemental Methods**). In brief, we first extracted  
419 candidate SNPs and small- and intermediate-scale indels (1-50 kb) in ten  
420 assemblies by pairwise gapped alignment of the ten assemblies and the  
421 reference genome assembly (Sscrofa10.2) using the LASTZ program. Then,  
422 the paired-end short-insert reads (180 and 500 bp) were separately aligned to  
423 the ten assembled genomes and the reference genome using the BWA  
424 software (v.0.7.12) (Li and Durbin, 2009). We filtered spurious SNPs and  
425 determined the heterozygous or homozygous mutations (depth  $\geq 10$ ) using  
426 SAMtools (v.1.3) (Li et al., 2009). With regard to indels, we eliminated spurious  
427 indel calls based on the calculation of read coverage for each indel locus with  
428 different criteria for indels  $\leq 50$  bp or  $> 50$  bp (Li et al., 2011).

#### 429 **Identification of selected regions using RSD algorithm**

430 To identify signatures of diversifying selection of pig breeds, relative  
431 homozygous SNP density (RSD) in nonoverlapping 10kb windows across the  
432 reference genome were calculated for each individual using a previously  
433 reported methodology (Atanur et al., 2013) (**Supplemental Methods**).

#### 434 **RNA-seq and data processing**

435 The 92 strand-specific RNA libraries (7 to 10 tissue libraries for each of ten  
436 individuals, which were used for *de novo* genome assemblies) were  
437 sequenced on the Illumina HiSeq 2500 platform (**Supplemental Methods**).  
438 High-quality reads were mapped to their respective *de novo* assemblies  
439 (**Supplemental Figs. S35,46**) or the reference genome (**Supplemental Fig.**

---

440 **S38**) using TopHat (v.2.1.0) (Trapnell et al., 2009). Cufflinks (v.2.2.1) (Trapnell  
441 et al., 2012) was used to quantify gene expression.

#### 442 **Discovery of missing sequences and missing genes**

443 We retrieved 'orphan reads' from paired-end DNA libraries with insert sizes of  
444 180 and 500 bp for each of ten breeds, where both ends of the read cannot be  
445 uniquely mapped to the reference genome (**Supplemental Figs. S27**). We  
446 re-localized these orphan reads to their respective assemblies. Sequences ( $\geq$   
447 500 bp in length) absent from the public reference genome assembly that were  
448 mapped by at least four orphan reads per base were considered 'missing  
449 sequences' (Kidd et al., 2010).

450 To identify genes embedded in the missing sequences, we conducted  
451 annotation of protein coding genes in the ten assemblies separately, using a  
452 combination of evidences from reference assembly-guided approach, the *ab*  
453 *initio*- and the homology-based methods, as well as RNA-seq data  
454 (**Supplemental Figs. S33,34, Table S17** and **Supplemental Methods**). We  
455 considered genes that showed  $> 50\%$  overlap of the gene length with missing  
456 sequences to be either missing or fragmented in the reference genome, and  
457 referred to them as 'missing genes'.

#### 458 **Determination of inter-assembly collinear genes**

459 The protein sequences of genes in the nine assemblies were separately  
460 queried against the protein sequences of the Large White assembly, which had  
461 the longest scaffold N50 size ( $\sim 2.45$  Mb) using BLASTp with an E-value cutoff  
462 of  $10^{-5}$  and restricting the output to a maximum of five hits per gene to serve as  
463 input for the MCScanX algorithm (Wang et al., 2012), which was used to detect

---

464 and classify high-confidence collinear blocks of coding genes (**Supplemental**  
465 **Tables S18,19**).

#### 466 **Detecting coding SNPs in missing genes under selection**

467 To test whether the recovered genes missing in the reference genome were  
468 under selection, we retrieved ~365.55 Gb orphan reads against the reference  
469 genome from 117 publicly available pig genomes (in total 6.03 trillion base  
470 resequencing data) ([Ai et al., 2015](#); [Choi et al., 2015](#); [Moon et al., 2015](#)) and  
471 aligned them to the intact scaffolds harboring missing genes across the ten  
472 breed assemblies (~636.38 Mb per assembly) (**Supplemental Fig. S41**). Of  
473 these, 6 wild boars and 41 domestic pigs belonging to 7 populations in China  
474 that have high-coverage depth (27.29 × of the reference genome, 14.43 × of  
475 the missing genes embedded scaffolds by 3.91 Gb orphan reads per individual)  
476 ([Ai et al., 2015](#)) were used to test for differentiation and possibly selection  
477 (**Supplemental Fig. S42**). The remaining 70 individuals (including 10 Korean  
478 wild boars and 60 European/ North American domestic pigs) with  
479 intermediate-coverage (15.87 × of the reference genome, 6.99 × of missing  
480 gene embedded scaffolds by 2.60 Gb orphan reads per individual) ([Choi et al.,](#)  
481 [2015](#); [Moon et al., 2015](#)) were used to investigate the patterns of selected loci  
482 (**Supplemental Fig. S41**).

483 We measured pairwise the extent of population differentiation of the coding  
484 SNPs in the missing genes between Chinese wild boars and seven Chinese  
485 domestic populations using the *FDIST* approach as implemented in Arlequin  
486 (v.3.5.2.2) ([Excoffier and Lischer, 2010](#)) (**Supplemental Fig. S43**,  
487 **Supplemental Table S21** and **Supplemental Methods**). We also measured

---

488 pairwise global  $F_{ST}$  values (**Supplemental Fig. S44A**) and performed a  
489 Bayesian test using the program BayeScan (v.2.1) (**Foll and Gaggiotti, 2008**)  
490 (**Supplemental Fig. S44B**) for every gene to detect highly differentiated SNPs  
491 between populations.

492

### 493 **Data access**

494 The nine pigs and Tibetan wild boar BioProjects are accessible at NCBI  
495 BioProject (<http://www.ncbi.nlm.nih.gov/bioproject>) under accession numbers  
496 PRJNA309108 and PRJNA186497, respectively. The assembled whole  
497 genome sequences have been deposited at NCBI GenBank  
498 (<http://www.ncbi.nlm.nih.gov/genbank>) under accession numbers  
499 LUXQ00000000.1 (Meishan), LUXR00000000.1 (Rongchang),  
500 LUXS00000000.1 (Hampshire), LUXT00000000.1 (Landrace),  
501 LUXU00000000.1 (Piétrain), LUXV00000000.1 (Bamei), LUXW00000000.1  
502 (Berkshire), LUXX00000000.1 (Large White), LUXY00000000.1 (Jinhua) and  
503 AORO00000000.2 (Tibetan wild boar, v.2). The unassembled sequencing  
504 reads of nine pigs and Tibetan wild boar have been deposited in NCBI  
505 Sequence Read Archive (SRA: <http://www.ncbi.nlm.nih.gov/sra>) under  
506 accession numbers SRP068560 and SRA065461, respectively. All RNA-seq  
507 reads and the genotyping data of the Illumina's porcine 60K Genotyping  
508 BeadChip (v.2) have been deposited in NCBI Gene Expression Omnibus  
509 (GEO: <https://www.ncbi.nlm.nih.gov/geo/>) under accession numbers  
510 GSE77776 and GSE83910, respectively.

511

512

## Acknowledgments

513 This work was supported by grants from the National Natural Science Foundation of  
514 China (31530073, 31522055, 31472081, 31372284, 31402046 and 31401073), the  
515 National Special Foundation for Transgenic Species of China (2014ZX0800950B and  
516 2014ZX08006-003), the National Program for Support of Top-notch Young Professionals,  
517 the Program for Innovative Research Team of Sichuan Province (2015TD0012), the  
518 Specialized Research Fund of Ministry of Agriculture of China (NYCYTX-009), the  
519 Program for Changjiang Scholars and Innovative Research Team in University  
520 (IRT13083), the National High Technology Research and Development Program of China  
521 (863 Program) (2013AA102502), the Science & Technology Support Program of Sichuan  
522 (Pig breeding-16ZC2850), the Fund of Fok Ying-Tung Education Foundation (141117), the  
523 National Key Technology R & D Program of China (2011BAD28B01), Chongqing Fund of  
524 Application & Development (CSTC2013YYKFC80003), Modern Agricultural Industry  
525 Technology System (CARS-36) and Chongqing Foundation of Agricultural Development  
526 (12404 and 14409).

527 *Author contributions:* Mingz.L., S.T., J.W., R.L. and X.L. led the experiments and designed  
528 the analytical strategy. L.C., D.L., A.J., Yingk.L, S.S., L.Z., Y.J. and L.B. performed animal  
529 work and prepared biological samples. L.J., J.M., X.W., Zongg.L., S.Z., and Z.J.  
530 constructed the DNA and RNA libraries and performed sequencing. Mingz.L., S.T., Yu.L.  
531 Q.T., Hongf.L. and T.C. designed the bioinformatics analysis process. Yu L., X.Z., Y.F.,  
532 Haif.L., D.W., Zongw.L. and H.Z. performed the genome assembly and annotation. S.T.,  
533 Mingz. L., L.C., Yan L., C.L. and G.W. performed the variation calling. Mingz.L., S.T., Y.G.,  
534 C.L., Z.G., G.T. and J.Z. identified missing sequences and missing genes. S.T., Mingz.L.,  
535 Q.T., X.Z., Q.P., M.M., and C.Y. performed selective sweep analyses. Mingz.L., L.C., R.L.  
536 and X.L. wrote the paper. Ming.L., J.W., V.N.G. and S.Z. revised the paper.

537

---

## 538 **Figure legends**

539

### 540 **Figure 1. Comparison of SNP calling between assembly-versus-assembly method** 541 **and resequencing approaches based on read mapping.**

542 Venn diagram with colors corresponding to bar chart shows the sharing of identified SNPs  
543 among assembly-versus-assembly method and two resequencing algorithms as  
544 implemented in SAMtools and GATK. An average of 4.25 M SNPs per breed were  
545 specifically identified by assembly-versus-assembly method (marked as yellow), while  
546 only 0.24 k SNPs per breed were categorized by resequencing approaches (marked as  
547 red). A significant fraction of the detected SNPs by the SAMtools (8.11 M per individual)  
548 and GATK (7.77 M per individual) was coincident (7.41 M; or 91.24% of SAMtools and  
549 95.34% of GATK) ([Supplemental Fig. S8](#)).

550

### 551 **Figure 2. Genomic variation between Chinese and European pigs.**

552 **(A)** Geographic locations of the original pig breeds. The Duroc (donor of the reference  
553 genome; it is denoted by star) and Hampshire pigs were developed mainly in North  
554 America but originated in Europe.

555 **(B)** Neighbor-joining phylogenetic tree, number of SNPs, transition / transversion ratio (Ts  
556 / Tv), heterozygous SNP ratio, patterns of regions of homozygosity (ROHs), and length  
557 and number of indels in the ten breeds (left to right). Violin plots of the heterozygous SNP  
558 ratio and Ts / Tv ratio were generated using non-overlapping 1 Mb windows (the medians  
559 are shown). For ROH, the circled area indicates the total length of ROHs in each breed.

560 **(C)** Pairwise genomic similarity of Chinese and European pigs by identical score (IS)  
561 values within each 10 kb window across the genome ( $n = 259,511$ ).

562

### 563 **Figure 3. Identification of breed-specific selective sweeps.**

564 **(A)** Number of homozygous SNPs in breed-specific selected regions. Of 74.21 k  
565 homozygous SNPs in 20.10 Mb selected regions, 65.75 k (88.60%) were unique to a  
566 particular breed, which highly concentrated in a small fraction (0.79%) of the genome and  
567 likely contributed to diversifying selection.

568 **(B)** Selective sweep regions identified in the Rongchang pig. Top panels, top half: genes  
569 residing within or in the vicinity ( $\pm 5$  kb) of the selected regions are presented for each

570 chromosome and ordered according to their locations. Top panels, lower half: degree of  
 571 haplotype sharing of selected regions in pairwise comparisons among the ten breeds.  
 572 Homozygous SNP frequencies in individual breeds were used to calculate identity scores  
 573 in 10 kb windows. Boxes (left) indicate pairwise comparison presented on that row (E,  
 574 European pigs; C, Chinese pigs) according to the color assigned to each pig breed (right).  
 575 Heat-map colors indicate identity scores. Second panels: Percentage stacked column  
 576 showing RSD values in the Rongchang-specific selected regions across ten breeds  
 577 sequenced. Rongchang showed predominantly higher RSD values than other breeds,  
 578 indicating that only this breed shows SNPs against the reference genome in this region.  
 579 Third panels: RSD in 10 kb windows for Rongchang plotted along chromosomes. Black  
 580 lines indicate selected regions ( $FDR < 0.05$ ). Nine selected genes orthologous to the  
 581 mammalian fat deposition genes were marked in red.

582

583 **Figure 4. Details of assembled *ALPK3* gene and selected variants.**

584 **(A)** Structure of assembled *ALPK3*. Top panels: the inter-assembly collinear genes  
 585 (colored rectangles) among ten assemblies are linked by gray lines, and the genes not  
 586 present in all ten assemblies are marked in black. *ALPK3* is denoted by a circle. Different  
 587 scaffolds are shown as alternating white and gray backgrounds. Second panels:  
 588 comparison of structure of *ALPK3* among the ten assemblies. Boxes and lines indicate  
 589 exons and introns, respectively.

590 **(B)** Coverage and depth for the longest gene model of *ALPK3* (Gene ID: RCGENE17759)  
 591 by crossly mapping reads from paired-end DNA libraries (insert sizes of 180 and 500 bp)  
 592 of the ten assemblies. The higher coverage depth ( $\geq 30 \times$ ) suggests slightly different  
 593 structures of *ALPK3* which is attributable to limitations of short read assembly; as such the  
 594 longest gene model is considered more reliable and used for subsequent analyses.

595 **(C)** Two selected missense mutations (T1,696-G and G1,733-C) in *ALPK3* between  
 596 Chinese wild boars ( $n = 6$ ) and domestic Min pigs ( $n = 6$ ). Top panels:  $F_{ST}$  and  
 597 Heterozygosity /  $(1 - F_{ST})$ , FDR (Arlequin) and  $q$  values (BayeScan) are plotted for 45  
 598 coding SNPs (18 missenses and 27 synonymous mutations). Second panels: LD pattern  
 599 of 45 SNPs in 101 domestic pigs from China ( $n = 41$ ), North America ( $n = 12$ ) and Europe  
 600 ( $n = 48$ ). Squares shaded in pink or red indicate significant LD between SNP pairs (bright  
 601 red indicates pairwise  $D' = 1$ ), white squares indicate no evidence of significant LD, and  
 602 blue squares indicate pairwise  $D' = 1$  without statistical significance. The adjacent  
 603 T1,696-G and G1,733-C are closely linked ( $D' = 1$ ,  $r^2 = 0.975$ , LOD = 41.6).

604

605 **Table 1. Summary of missing sequences and genes of the reference genome**  
 606 **(Sscrofa10.2)**

Assemblies	Missing sequence		Missing genes			
	Number	Length (Mb)	Number	Assigned by inter-assembly collinearity		
				Singleton	Assembled in 2-9 breeds	Assembled in all 10 breeds
Hampshire	82,824	136.33	1,105	67	167	
Berkshire	82,958	136.40	1,092	63	158	
Landrace	82,741	135.86	1,093	65	157	
Piértrain	82,472	135.75	1,096	60	165	
Large White	82,987	136.04	1,105	14	220	
Bamei	84,336	137.49	1,104	83	150	871
Jinhua	84,031	137.34	1,090	65	154	
Meishan	85,197	138.65	1,116	95	150	
Rongchang	84,062	137.88	1,064	49	144	
Tibetan wild boar	86,592	138.42	1,094	85	138	

607

608

## References

- 609 Ai H, Fang X, Yang B, Huang Z, Chen H, Mao L, Zhang F, Zhang L, Cui L, He W, et al.  
 610 2015. Adaptation and possible ancient interspecies introgression in pigs identified by  
 611 whole-genome sequencing. *Nat Genet* **47**: 217-225.
- 612 Alkan C, Sajjadian S, Eichler EE. 2011. Limitations of next-generation genome sequence  
 613 assembly. *Nat Methods* **8**: 61-65.
- 614 Atanur SS, Diaz AG, Maratou K, Sarkis A, Rotival M, Game L, Tschannen MR, Kaisaki PJ,  
 615 Otto GW, Ma MC, et al. 2013. Genome sequencing reveals loci under artificial selection  
 616 that underlie disease phenotypes in the laboratory rat. *Cell* **154**: 691-703.
- 617 Besenbacher S, Liu S, Izarzugaza JM, Grove J, Belling K, Bork-Jensen J, Huang S, Als  
 618 TD, Li S, Yadav R, et al. 2015. Novel variation and *de novo* mutation rates in  
 619 population-wide *de novo* assembled Danish trios. *Nat Commun* **6**: 5969.

- 
- 620 Bosse M, Megens HJ, Madsen O, Paudel Y, Frantz LA, Schook LB, Crooijmans RP,  
621 Groenen MA. 2012. Regions of homozygosity in the porcine genome: consequence of  
622 demography and the recombination landscape. *PLoS Genet* **8**: e1003100.
- 623 Carneiro M, Rubin CJ, Di Palma F, Albert FW, Alfoldi J, Barrio A.M, Pielberg G, Rafati N,  
624 Sayyab S, Turner-Maier J, et al. 2014. Rabbit genome analysis reveals a polygenic basis  
625 for phenotypic change during domestication. *Science* **345**: 1074-1079.
- 626 Chaisson MJ, Huddleston J, Dennis MY, Sudmant PH, Malig M, Hormozdiari F, Antonacci  
627 F, Surti U, Sandstrom R, Boitano M, et al. 2015a. Resolving the complexity of the human  
628 genome using single-molecule sequencing. *Nature* **517**: 608-611.
- 629 Chaisson MJ, Wilson RK, Eichler EE. 2015b. Genetic variation and the *de novo* assembly  
630 of human genomes. *Nat Rev Genet* **16**: 627-640.
- 631 Chen K, Baxter T, Muir WM, Groenen MA, Schook LB. 2007. Genetic resources, genome  
632 mapping and evolutionary genomics of the pig (*Sus scrofa*). *Int J Biol Sci* **3**: 153-165.
- 633 Choi JW, Chung WH, Lee KT, Cho ES, Lee SW, Choi BH, Lee SH, Lim W, Lim D, Lee YG,  
634 et al. 2015. Whole-genome resequencing analyses of five pig breeds, including Korean  
635 wild and native, and three European origin breeds. *DNA Res* **22**: 259-267.
- 636 Excoffier L, Lischer HE. 2010. Arlequin suite ver 3.5: a new series of programs to perform  
637 population genetics analyses under Linux and Windows. *Mol Ecol Resour* **10**: 564-567.
- 638 Foll M, Gaggiotti O. 2008. A genome-scan method to identify selected loci appropriate for  
639 both dominant and codominant markers: a Bayesian perspective. *Genetics* **180**: 977-993.
- 640 Frantz LA, Schraiber JG, Madsen O, Megens HJ, Bosse M, Paudel Y, Semiadi G,  
641 Meijaard E, Li N, Crooijmans RP, et al. 2013. Genome sequencing reveals fine scale  
642 diversification and reticulation history during speciation in *Sus*. *Genome Biol* **14**: R107.
- 643 Gan X, Stegle O, Behr J, Steffen JG, Drewe P, Hildebrand KL, Lyngsoe R, Schultheiss SJ,  
644 Osborne EJ, Sreedharan VT, et al. 2011. Multiple reference genomes and transcriptomes  
645 for *Arabidopsis thaliana*. *Nature* **477**: 419-423.

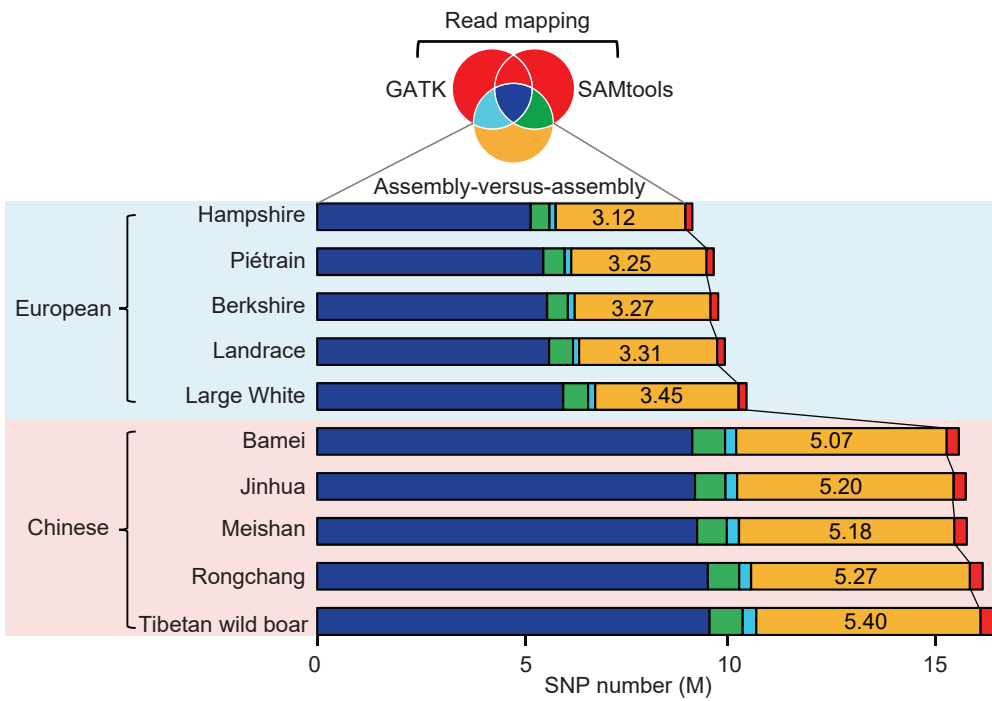
- 
- 646 Genovese G, Handsaker RE, Li H, Altemose N, Lindgren AM, Chambert K, Pasaniuc B,  
647 Price AL, Reich D, Morton CC, et al. 2013. Using population admixture to help complete  
648 maps of the human genome. *Nat Genet* **45**: 406-414.
- 649 Gerstein HC, Waltman L. 2006. Why don't pigs get diabetes? Explanations for variations  
650 in diabetes susceptibility in human populations living in a diabetogenic environment. *Can*  
651 *Med Assoc J* **174**: 25-26.
- 652 Groenen MA, Archibald AL, Uenishi H, Tuggle CK, Takeuchi Y, Rothschild MF,  
653 Rogel-Gaillard C, Park C, Milan D, Megens HJ, et al. 2012. Analyses of pig genomes  
654 provide insight into porcine demography and evolution. *Nature* **491**: 393-398.
- 655 Kidd JM, Sampas N, Antonacci F, Graves T, Fulton R, Hayden HS, Alkan C, Malig M,  
656 Ventura M, Giannuzzi G, et al. 2010. Characterization of missing human genome  
657 sequences and copy-number polymorphic insertions. *Nat Methods* **7**: 365-371.
- 658 Koopmans SJ, Schuurman T. 2015. Considerations on pig models for appetite, metabolic  
659 syndrome and obese type II diabetes: From food intake to metabolic disease. *Eur J*  
660 *Pharmacol* **759**: 231-239.
- 661 Koren S, Schatz MC, Walenz BP, Martin J, Howard JT, Ganapathy G, Wang Z, Rasko DA,  
662 McCombie WR, Jarvis ED, et al. 2012. Hybrid error correction and *de novo* assembly of  
663 single-molecule sequencing reads. *Nat Biotechnol* **30**: 693-700.
- 664 Kunej T, Jevsinek Skok D, Zorc M, Ogrinc A, Michal JJ, Kovac M, Jiang Z. 2013. Obesity  
665 gene atlas in mammals. *J Genomics* **1**: 45-55.
- 666 Larson G, Dobney K, Albarella U, Fang M, Matisoo-Smith E, Robins J, Lowden S,  
667 Finlayson H, Brand T, Willerslev E, et al. 2005. Worldwide phylogeography of wild boar  
668 reveals multiple centers of pig domestication. *Science* **307**: 1618-1621.
- 669 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler  
670 transform. *Bioinformatics* **25**: 1754-1760.
- 671 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin

- 
- 672 R, Genome Project Data Processing Subgroup. 2009. The Sequence Alignment/Map  
673 format and SAMtools. *Bioinformatics* **25**: 2078-2079.
- 674 Li M, Tian S, Jin L, Zhou G, Li Y, Zhang Y, Wang T, Yeung CK, Chen L, Ma J, et al. 2013.  
675 Genomic analyses identify distinct patterns of selection in domesticated pigs and Tibetan  
676 wild boars. *Nat Genet* **45**: 1431-1438.
- 677 Li R, Li Y, Zheng H, Luo R, Zhu H, Li Q, Qian W, Ren Y, Tian G, Li J, et al. 2010a. Building  
678 the sequence map of the human pan-genome. *Nat Biotechnol* **28**: 57-63.
- 679 Li R, Zhu H, Ruan J, Qian W, Fang X, Shi Z, Li Y, Li S, Shan G, Kristiansen K, et al. 2010b.  
680 *De novo* assembly of human genomes with massively parallel short read sequencing.  
681 *Genome Res* **20**: 265-272.
- 682 Li Y, Zheng H, Luo R, Wu H, Zhu H, Li R, Cao H, Wu B, Huang S, Shao H, et al. 2011.  
683 Structural variation in two human genomes mapped at single-nucleotide resolution by  
684 whole genome *de novo* assembly. *Nat Biotechnol* **29**: 723-730.
- 685 Li YH, Zhou G, Ma J, Jiang W, Jin LG, Zhang Z, Guo Y, Zhang J, Sui Y, Zheng L, et al.  
686 2014. *De novo* assembly of soybean wild relatives for pan-genome analysis of diversity  
687 and agronomic traits. *Nat Biotechnol* **32**: 1045-1052.
- 688 Liu S, Lorenzen ED, Fumagalli M, Li B, Harris K, Xiong Z, Zhou L, Korneliussen TS,  
689 Somel M, Babbitt C, et al. 2014. Population genomics reveal recent speciation and rapid  
690 evolutionary adaptation in polar bears. *Cell* **157**: 785-794.
- 691 Mackenzie FE, Romero R, Williams D, Gillingwater T, Hilton H, Dick J, Riddoch-Contreras  
692 J, Wong F, Ireson L, Powles-Glover N, et al. 2009. Upregulation of *PKD1L2* provokes a  
693 complex neuromuscular disease in the mouse. *Hum Mol Genet* **18**:3553-3566.
- 694 Mainland JD, Keller A, Li YR, Zhou T, Trimmer C, Snyder LL, Moberly AH, Adipietro KA,  
695 Liu WL, Zhuang H, et al. 2014. The missense of smell: functional variability in the human  
696 odorant receptor repertoire. *Nat Neurosci* **17**: 114-120.
- 697 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,

- 
- 698 Altshuler D, Gabriel S, Daly M, et al. 2010. The Genome Analysis Toolkit: a MapReduce  
699 framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**:  
700 1297-1303.
- 701 Moon S, Kim TH, Lee KT, Kwak W, Lee T, Lee SW, Kim MJ, Cho K, Kim N, Chung WH, et  
702 al. 2015. A genome-wide scan for signatures of directional selection in domesticated pigs.  
703 *BMC Genomics* **16**: 130.
- 704 Neafsey DE, Waterhouse RM, Abai MR, Aganezov SS, Alekseyev MA, Allen JE, Amon J,  
705 Arca B, Arensburger P, Artemov G, et al. 2015. Highly evolvable malaria vectors: the  
706 genomes of 16 *Anopheles* mosquitoes. *Science* **347**: 1258522.
- 707 Read BA, Kegel J, Klute MJ, Kuo A, Lefebvre SC, Maumus F, Mayer C, Miller J, Monier A,  
708 Salamov A, et al. 2013. Pan genome of the phytoplankton *Emiliana* underpins its global  
709 distribution. *Nature* **499**: 209-213.
- 710 Roepstorff C, Vistisen B, Kiens B. 2005. Intramuscular triacylglycerol in energy  
711 metabolism during exercise in humans. *Exerc Sport Sci Rev* **33**:182-188.
- 712 Rubin CJ, Megens HJ, Martinez Barrio A, Maqbool K, Sayyab S, Schwochow D, Wang C,  
713 Carlborg O, Jern P, Jorgensen CB, et al. 2012. Strong signatures of selection in the  
714 domestic pig genome. *Proc. Natl. Acad. Sci. USA* **109**: 19529-19536.
- 715 Rubin CJ, Zody MC, Eriksson J, Meadows JR, Sherwood E, Webster MT, Jiang L, Ingman  
716 M, Sharpe T, Ka S, et al. 2010. Whole-genome resequencing reveals loci under selection  
717 during chicken domestication. *Nature* **464**: 587-591.
- 718 Schug J, Schuller WP, Kappen C, Salbaum JM, Bucan M, Stoeckert CJ Jr. 2005.  
719 Promoter features related to tissue specificity as measured by Shannon entropy. *Genome*  
720 *Biol* **6**: R33.
- 721 Trapnell C, Pachter L, Salzberg SL. 2009. TopHat: discovering splice junctions with  
722 RNA-Seq. *Bioinformatics* **25**: 1105-1111.
- 723 Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, Pimentel H, Salzberg SL, Rinn

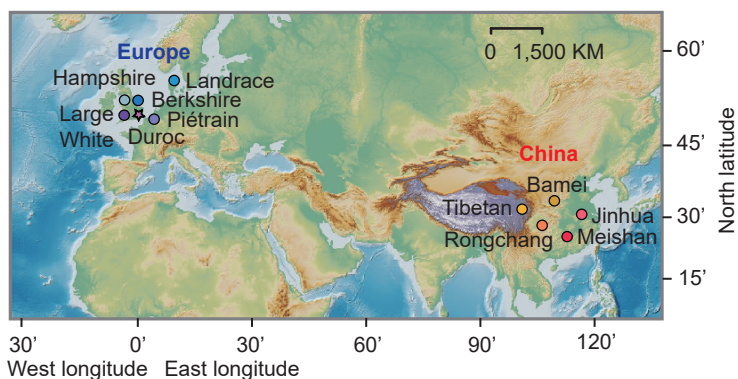
- 
- 724 JL, Pachter L. 2012. Differential gene and transcript expression analysis of RNA-seq  
725 experiments with TopHat and Cufflinks. *Nat Protoc* **7**: 562-578.
- 726 Van Sligtenhorst I, Ding ZM, Shi ZZ, Read RW, Hansen G, Vogel P. 2012.  
727 Cardiomyopathy in  $\alpha$ -kinase 3 (*ALPK3*)-deficient mice. *Vet Pathol* **49**: 131-141.
- 728 Wang LY, Wang AG, Wang LX, Li K, Yang GS, He RG, Qian L, Xu NY, Huang RH, Peng  
729 ZZ, et al. 2011. Animal genetic resources in China: pigs. (ed. China National Commission  
730 of Animal Genetic Resources) , pp. 2-16. China Agricultural Press, Beijing.
- 731 Wang Y, Tang H, Debarry JD, Tan X, Li J, Wang X, Lee TH, Jin H, Marler B, Guo H, et al.  
732 2012. MCScanX: a toolkit for detection and evolutionary analysis of gene synteny and  
733 collinearity. *Nucleic Acids Res* **40**: e49.
- 734 Weisenfeld NI, Yin S, Sharpe T, Lau B, Hegarty R, Holmes L, Sogoloff B, Tabbaa D,  
735 Williams L, Russ C, et al. 2014. Comprehensive variation discovery in single human  
736 genomes. *Nat Genet* **46**: 1350-1355.
- 737 White S. 2011. From globalized pig breeds to capitalist pigs: a study in animal cultures  
738 and evolutionary history. *Environ Hist* **16**: 94-120.
- 739 Yan G, Zhang G, Fang X, Zhang Y, Li C, Ling F, Cooper DN, Li Q, Li Y, van Gool AJ, et al.  
740 2011. Genome sequencing and comparison of two nonhuman primate animal models, the  
741 cynomolgus and Chinese rhesus macaques. *Nat Biotechnol* **29**: 1019-1023.
- 742 Zhang QJ, Zhu T, Xia EH, Shi C, Liu YL, Zhang Y, Liu Y, Jiang WK, Zhao YJ, Mao SY, et al.  
743 2014. Rapid diversification of five *Oryza* AA genomes associated with rice adaptation.  
744 *Proc. Natl. Acad. Sci. USA* **111**: E4954-4962.
- 745 Zook JM, Salit M. 2015. Advancing benchmarks for genome sequencing. *Cell Systems* **1**:  
746 176-177.

**Fig. 1**

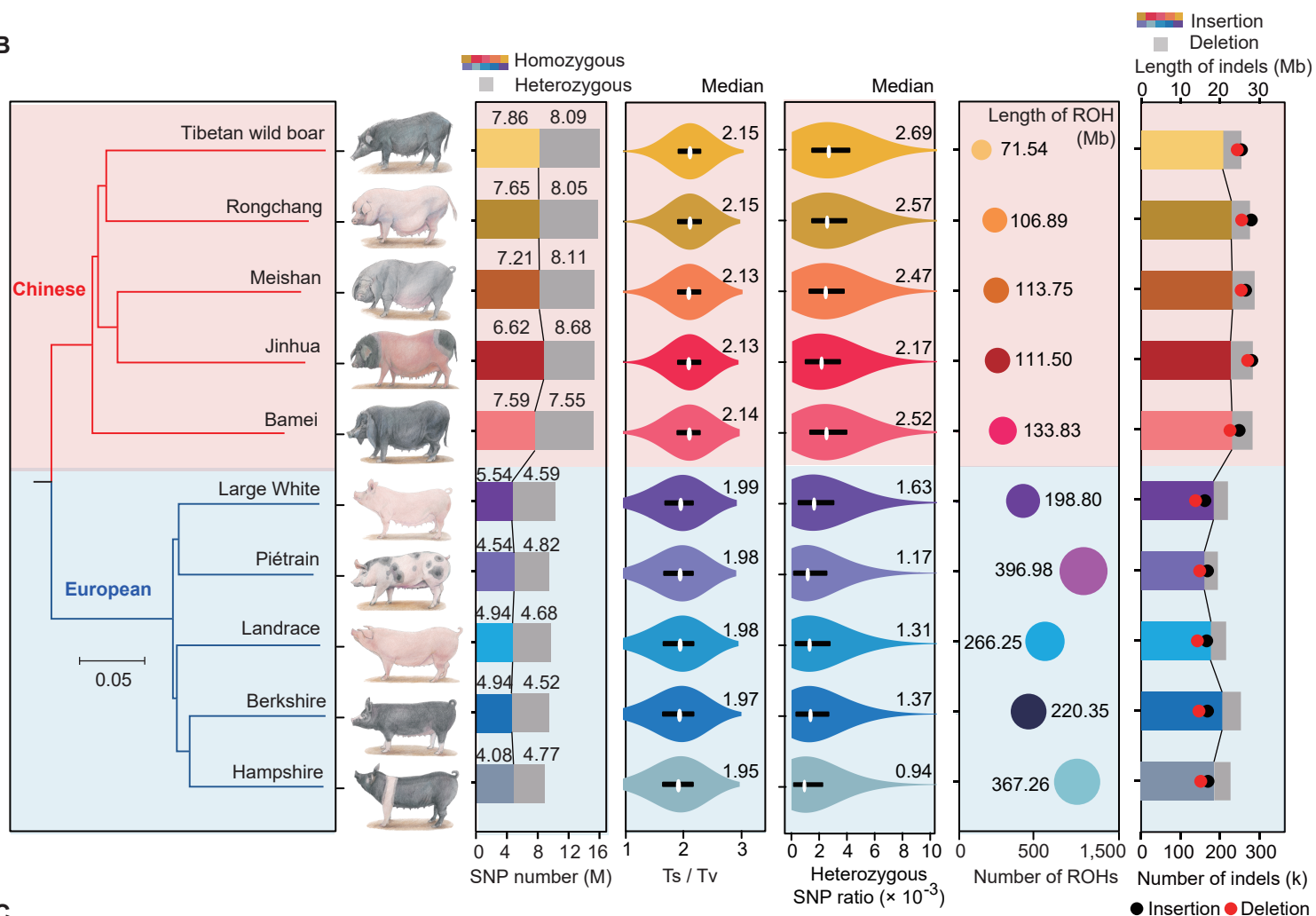


**Fig. 2**

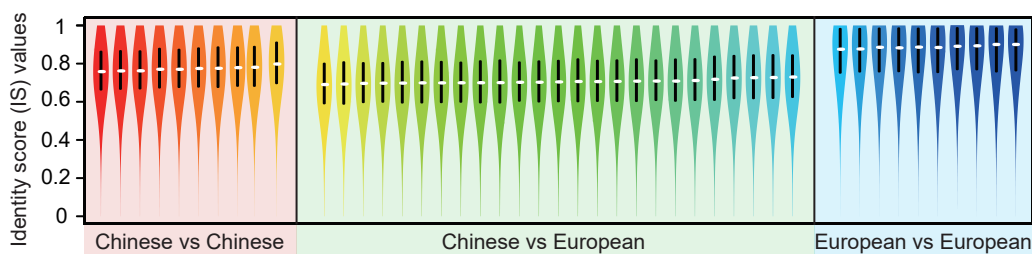
**A**



**B**

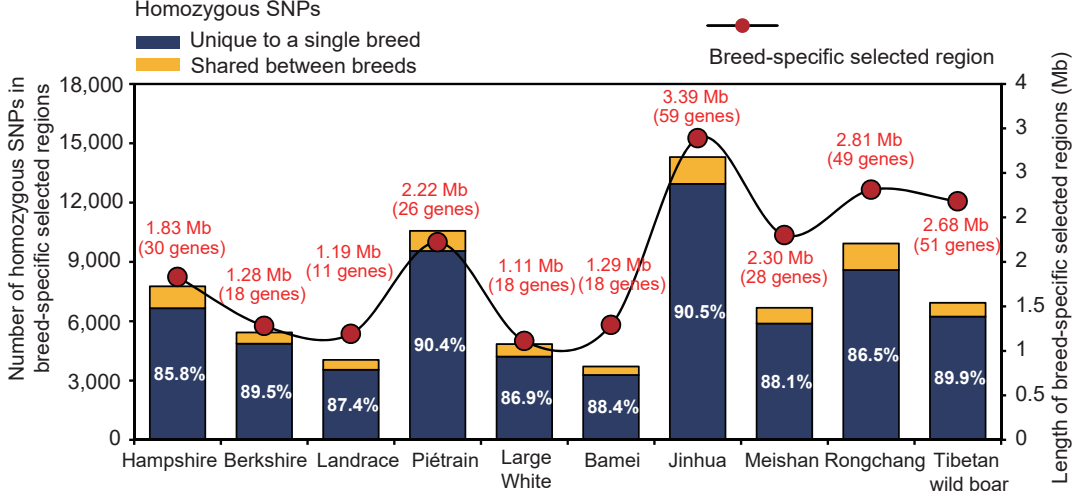


**C**

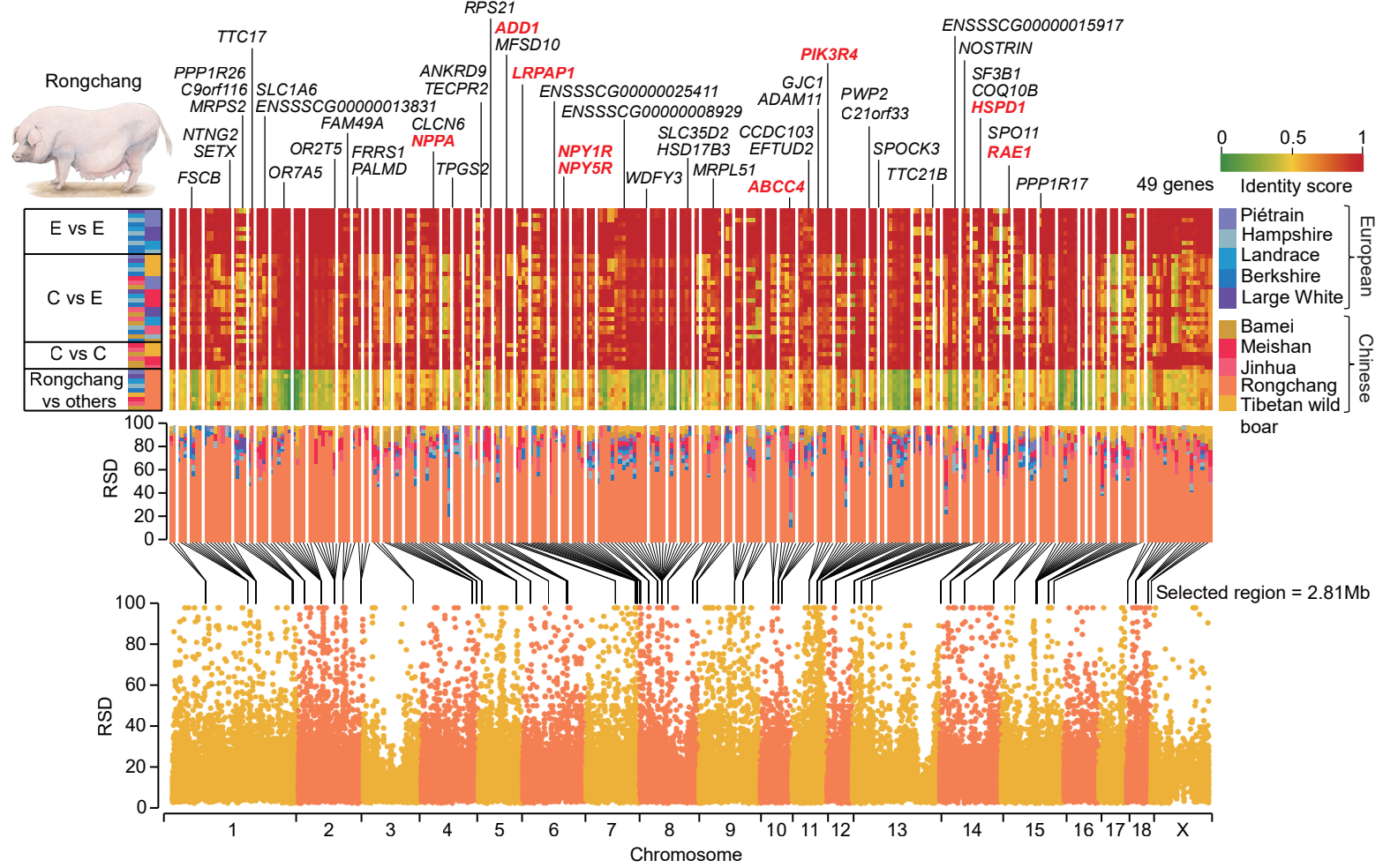


**Fig. 3**

**A**

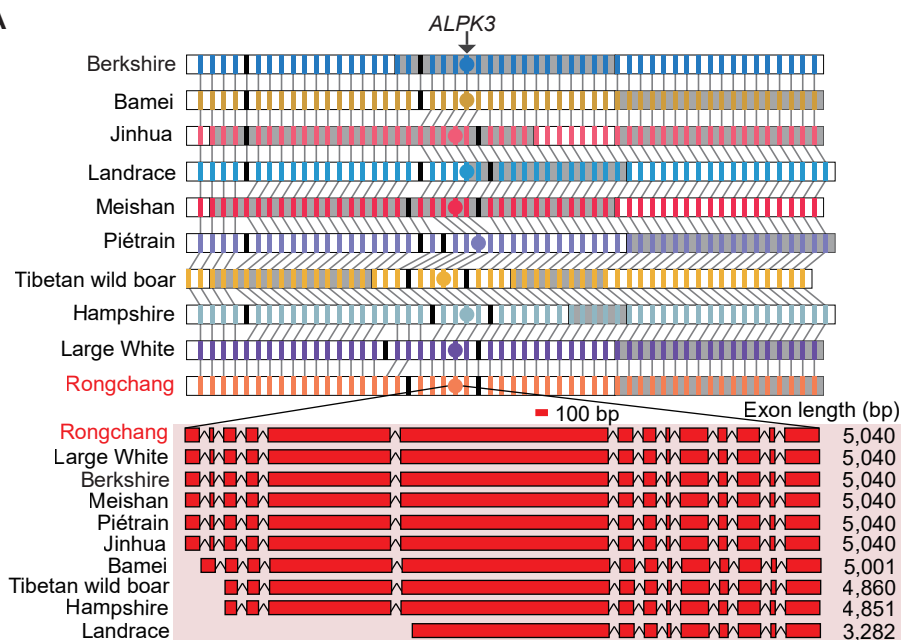


**B**

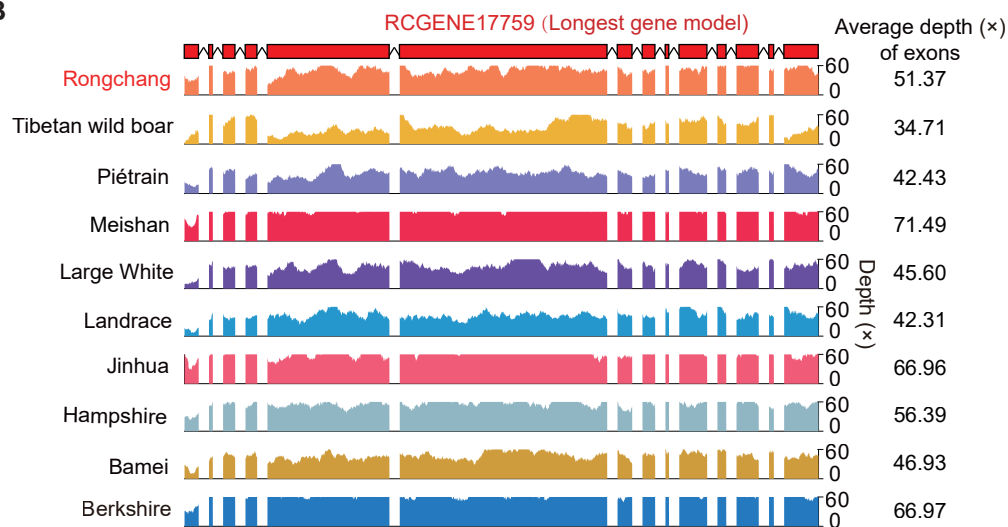


**Fig. 4**

**A**



**B**



**C**

