

1 **Supplementary Information**

2

3 **Accurate haplotype construction and detection of selection signatures**
4 **enabled by high quality pig genome sequences**

5 **Authors:** Xinkai Tong^{1,2}, Dong Chen¹, Jianchao Hu¹, Shiyao Lin¹, Ziqi Ling¹, Huashui

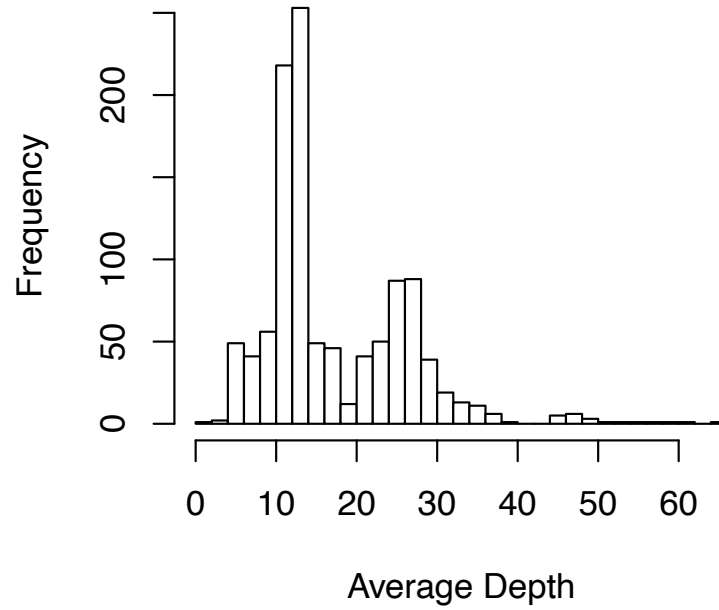
6 Ai¹, Zhiyan Zhang^{1*} & Lusheng Huang^{1*}

7 ¹National Key Laboratory for Swine genetic improvement and production technology,
8 Ministry of Science and Technology of China, Jiangxi Agricultural University,
9 NanChang, Jiangxi Province, PR China.

10 ²College of Life Sciences, Jiangxi Normal University, NanChang, Jiangxi Province, PR
11 China.

12 *Corresponding authors: Prof. Dr. Zhiyan Zhang, Prof. Dr. Lusheng Huang

13



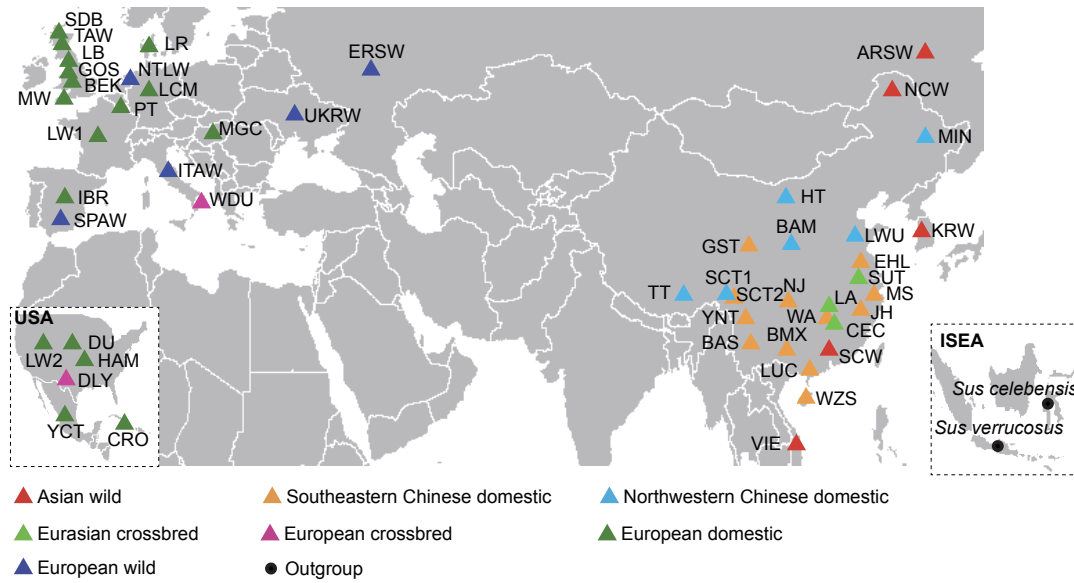
15

16 **Supplementary Fig. 1 The frequency distribution of average sequencing depth of**17 **samples.** A histogram is used to show the sequencing depth of samples. The x-axis18 denotes average sequencing depth (\times) for each sample. The y-axis denotes the number

19 of samples. Each bin is plotted as a bar whose height represents the sample size in that

20 bin.

21



23

24 **Supplementary Fig. 2 Geographic distribution (i.e., sampling sites or country of**25 **origin of breeds) of the samples analyzed in this study.** Color codes of triangles

26 correspond to the seven groups obtained in the phylogenetic tree. ISEA, islands of

27 southeast Asia. ARSW, Asian Russia wild; ERSW, European Russia wild; BAM, Bamei;

28 BAS, Baoshan; BEK, Berkshire; BMX, Bamaxiang; CRO, Creole; DU, Duroc; EHL,

29 Erhualian; CEC or F1/F2/F3, Eurasian Crossbred; GOS, Gloucester old Spot; GST,

30 Gansu Tibetan; HAM, Hampshire; HT, Hetao; IBR, Iberian pig; ITAW, Italian wild; JH,

31 Jinhua; KRW, Korean wild; LA, Lean spotted pig; LB, Large Black; LCM, Leicoma;

32 LR, Landrace; LUC, Luchuan; LW, LargeWhite; LWU, Laiwu; MGC, Mangalica; MS,

33 Meishan; MW, Middle White; NCW, Northern chinese wild pig; NJ, Neijiang; NTLW,

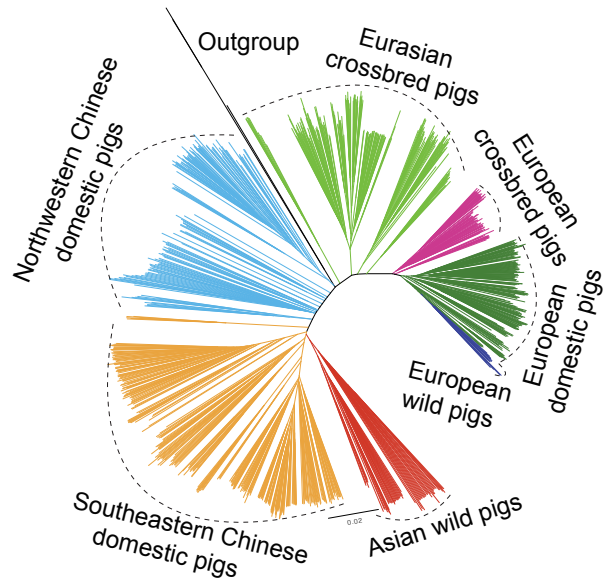
34 Netherland wild; PT, Pietrain; SCT, Sichuan Tibetan; SCW, Sourthern chinese wild;

35 SDB, British Saddleback; SPAW, Spanish wild; SUT, Sutan; TAW, Tamworth; TT,

36 Tibetan Tibetan; UKRW, Ukraine wild; VIE, Vietnam wild; WA, Wanan spotted; WDU,

37 White Duroc; WZS, Wuzhishan; YCT, Yucatan; YNT, Yunnan Tibetan. The map is

38 generated through the ggplot2 package
39 (<https://ggplot2.tidyverse.org/reference/borders.html>) in the R program. The ggplot2
40 package invokes the maps package to generate a world map
41 (<https://rdrr.io/cran/maps/man/world.html>). This world map is imported from the public
42 domain Natural Earth project (<https://www.naturalearthdata.com>).



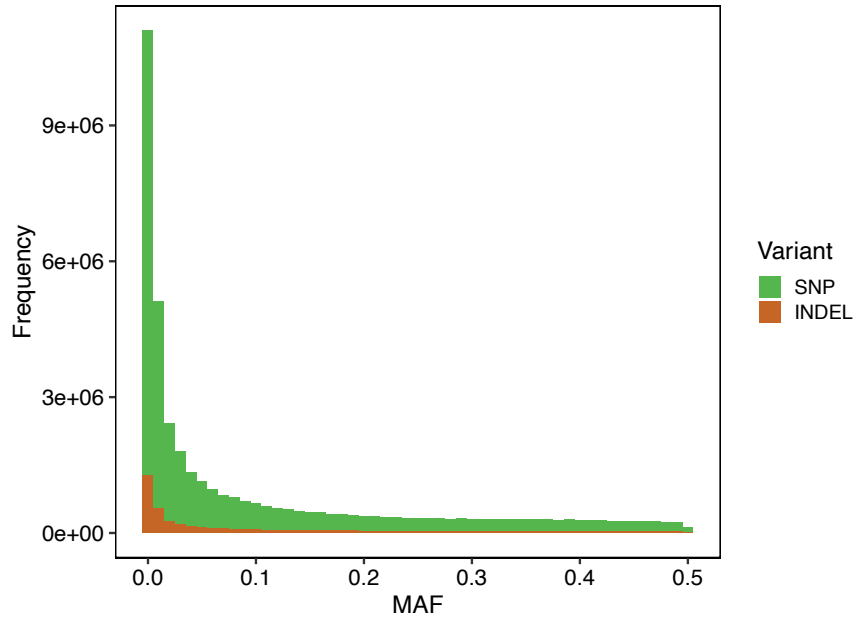
43

44 **Supplementary Fig. 3 Neighbor-Joining phylogenetic tree, using *Sus verrucosus***

45 **(Javan warty pig) and *Sus celebensis* (Celebes warty pig) as outgroup. Each**

46 **branch denotes an individual. All samples are used to construct the phylogenetic tree.**

47



48

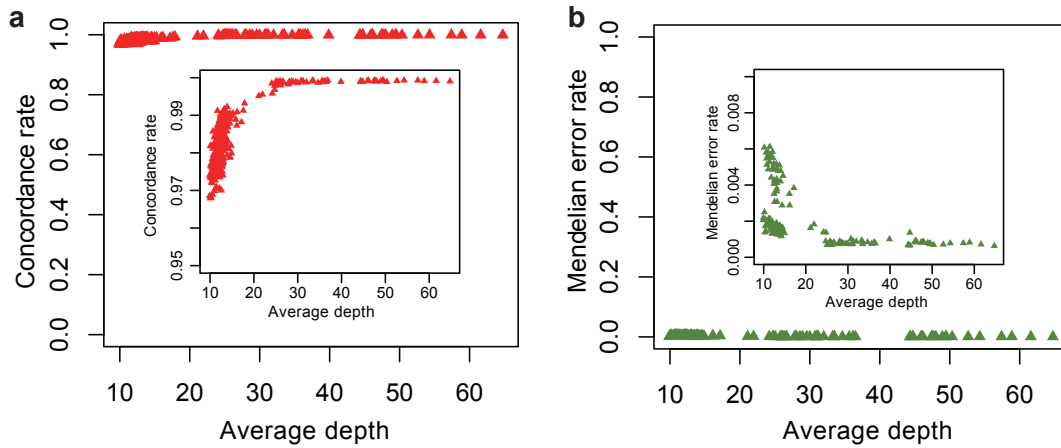
49 **Supplementary Fig. 4 The minor allele frequency (MAF) distribution of SNPs and**

50 **Indels.** A histogram is used to show the frequency of variant under different MAF. The

51 x-axis denotes minor allele frequency. The y-axis denotes the number of variants. Each

52 bin is plotted as a bar whose height represents the number of variants in that bin.

53



54

55

56 **Supplementary Fig. 5 Validation the accuracy of variant genotyping. (a)** The

57 concordance rate between sequence genotypes and array genotypes. The x-axis shows

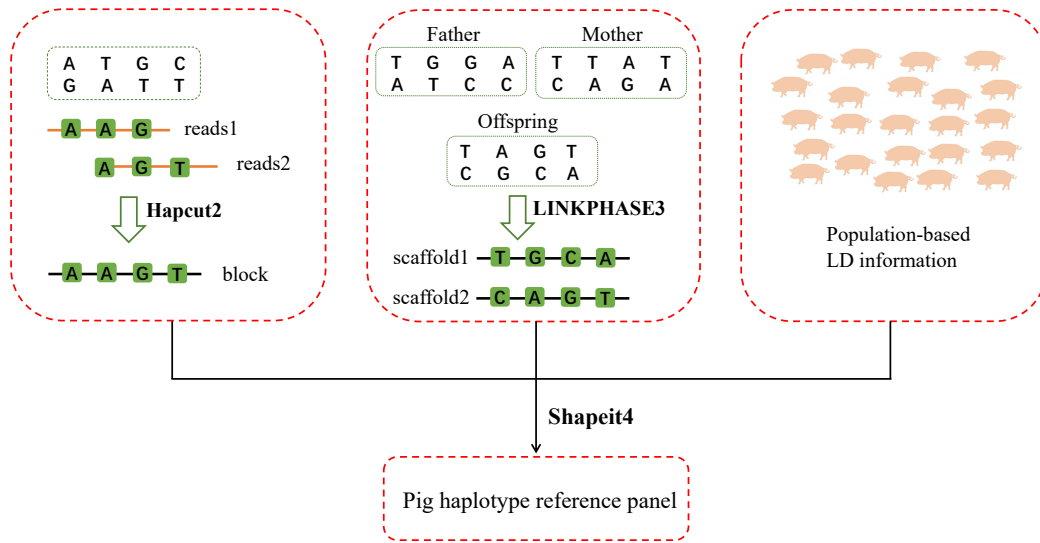
58 average sequencing depth. **(b)** The Mendel Error rate estimated by summarizing error

59 inherited variants in the 179 parent-offspring trios (44) or duos (135). Each triangle

60 represents an individual.

61

62



63

64

65 **Supplementary Fig. 6 The schematic diagram for construction of the pig**

66 **haplotype reference panel.** The left, middle, right boxes represent haplotype blocks

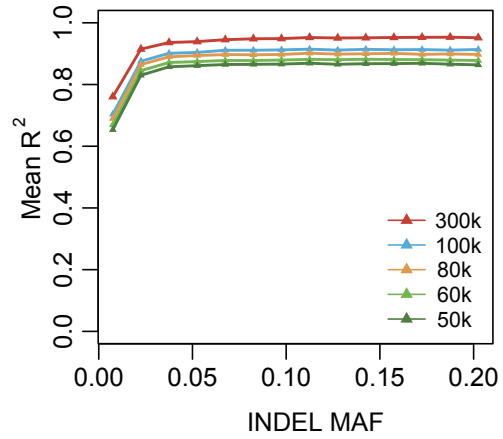
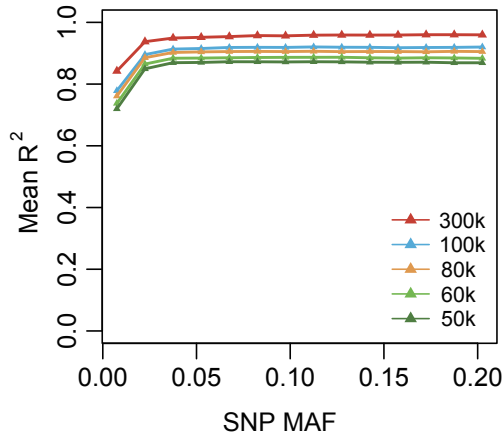
67 establishment, haplotype scaffolds build and Linkage Disequilibrium information

68 extraction from population, respectively. Pig cartoons adapted from Chen, C., Zhou, Y.,

69 Fu, H. et al. Expanded catalog of microbial genes and metagenome-assembled genomes

70 from the pig gut microbiome. Nat Commun 12, 1106 (2021).

71 <https://doi.org/10.1038/s41467-021-21295-0>.



72

73

74 **Supplementary Fig. 7 Imputation accuracy using ten common Chinese indigenous**

75 **pigs (Erhualian, N = 4, Bamaxiang, N = 3, Laiwu, N = 3) as targets and the**

76 **remaining 1854 haplotypes as reference panel.** The x-axis shows the minor allele

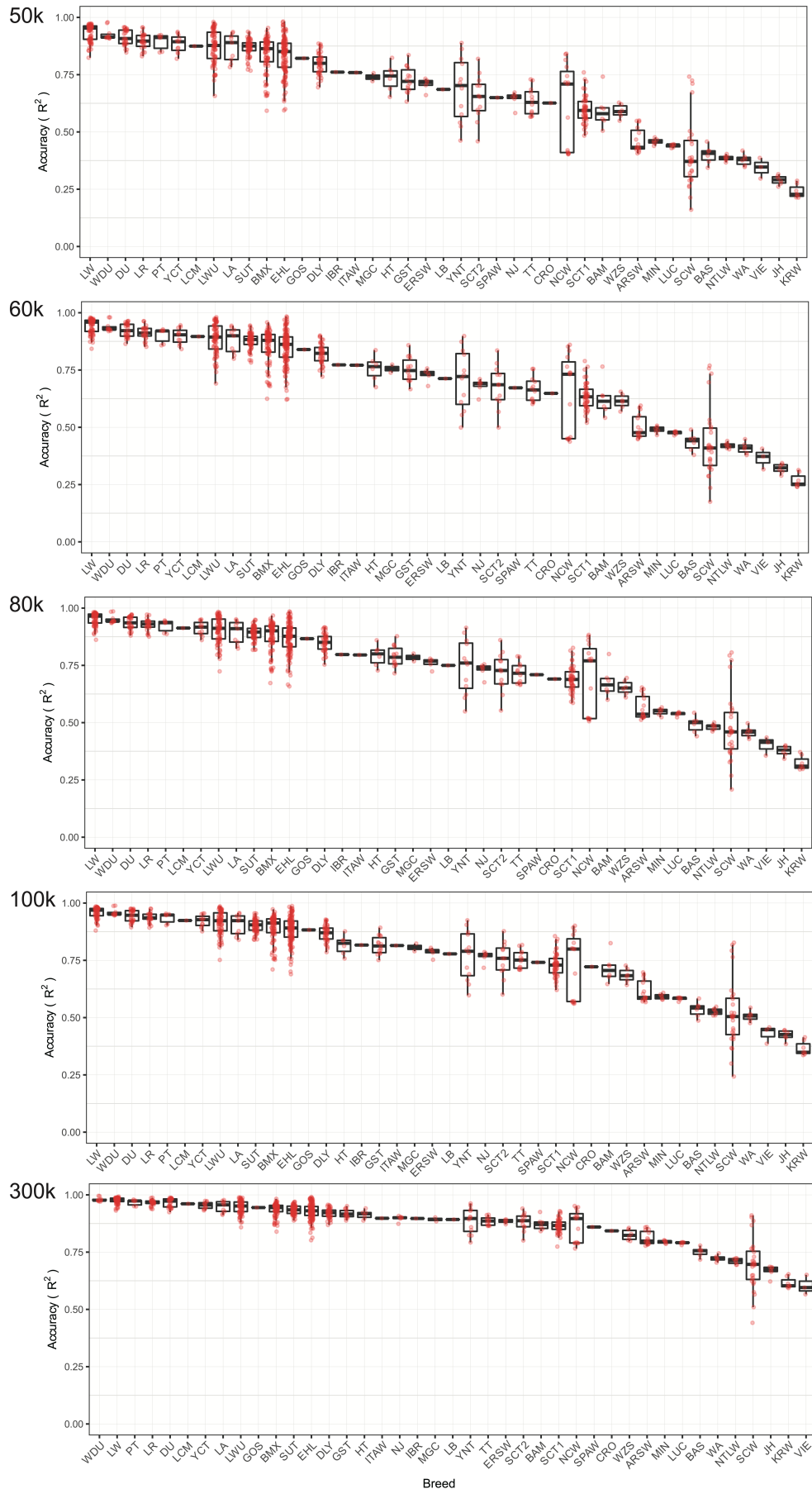
77 frequency of imputed variants. The y-axis shows imputation accuracy measured by

78 average R^2 (squared Pearson correlation). The colored triangles represent the R^2 under

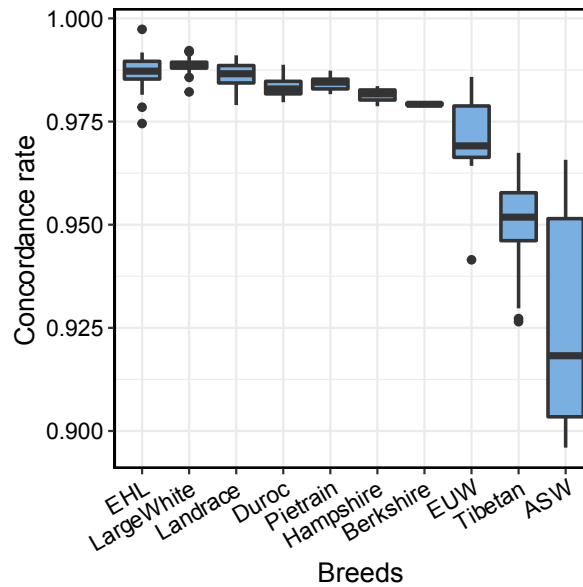
79 different minor allele frequencies. The 50k, 60k, 80k, 100k, and 300k denote the

80 number of randomly selected polymorphism sites being used for imputation.

81



83 **Supplementary Fig. 8 Genotypes imputation accuracy for each breed included in**
84 **the haplotype reference panel.** The accuracy was evaluated as follows: 1) We sampled
85 ten individuals as imputation targets and the remaining 1854 haplotypes as reference
86 panel. The 50k, 60k, 80k, 100k, and 300k variants at autosomes were randomly selected
87 to mimic chips. The genotypes of unselected variants were imputed for sampled ten
88 individuals. 2) Step 1) was repeated until genotypes of all individuals were imputed
89 once. 3) The imputation accuracy was measured by average R^2 (squared Pearson
90 correlation) between sequenced genotypes and imputed genotypes. The letters on the
91 upper left indicates the number of variants selected randomly for imputation. The y-
92 axis represents imputation accuracy measured by R^2 (squared Pearson correlation). The
93 x-axis denotes abbreviation of breeds. Each point represents an individual. Boxplots
94 show median, 25th and 75th percentile, the whiskers indicate the minima and maxima,
95 and the points laying outside the whiskers of boxplots represent the outliers. ARSW,
96 Asian Russia wild; ERSW, European Russia wild; BAM, Bamei; BAS, Baoshan; BMX,
97 Bamaxiang; CRO, Creole; DU, Duroc; EHL, Erhualian; GOS, Gloucester old Spot;
98 GST, Gansu Tibetan; HT, Hetao; IBR, Iberian pig; ITAW, Italian wild; JH, Jinhua; KRW,
99 Korean wild; LA, Lean spotted pig; LB, Large Black; LCM, Leicoma; LR, Landrace;
100 LUC, Luchuan; LW, LargeWhite; LWU, Laiwu; MGC, Mangalica; NCW, Northern
101 chinese wild pig; NJ, Neijiang; NTLW, Netherland wild; PT, Pietrain; SCT, Sichuan
102 Tibetan; SCW, Sourthern chinese wild; SPAW, Spanish wild; SUT, Sutai; TT, Tibetan
103 Tibetan; VIE, Vietnam wild; WA, Wanan spotted; WDU, White Duroc; WZS,
104 Wuzhishan; YCT, Yucatan; YNT, Yunnan Tibetan.



105

106

107 **Supplementary Fig. 9 The concordance rate between sequencing genotypes and**

108 **imputed genotypes for samples being removed for low call rate.** Boxplots show

109 median, 25th and 75th percentile, the whiskers indicate the minima and maxima, and

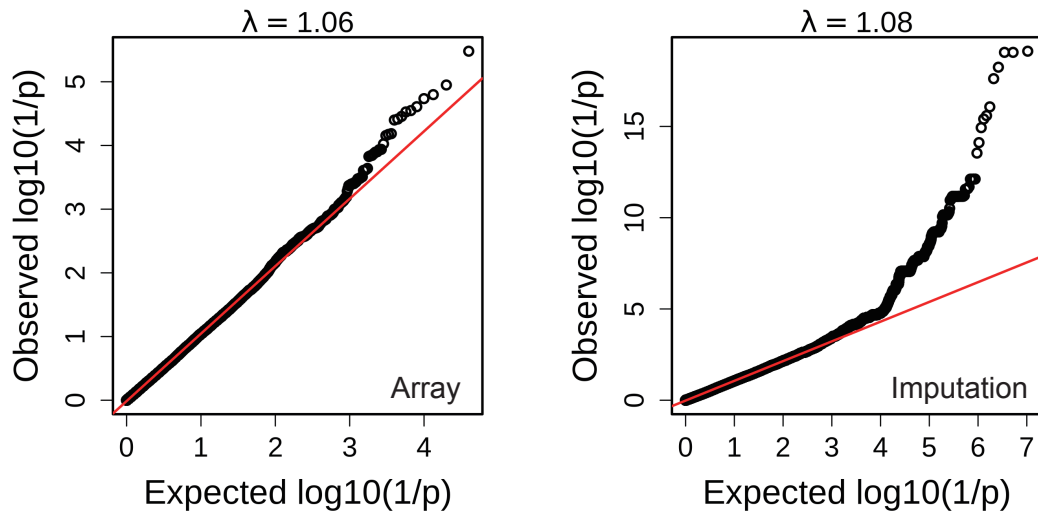
110 the points laying outside the whiskers of boxplots represent the outliers. EHL

111 (Erhualian), n = 22; LargeWhite, n = 20; Landrace, n = 12; Duroc, n = 4; Pietrain, n =

112 5; Hampshire, n = 3; Berkshire, n = 3; EUW (European wild pigs), n = 10; Tibetan, n

113 = 31; ASW (Asian wild pigs), n = 10.

114



115

116

117 **Supplementary Fig. 10 Quantile-quantile plots of array and imputation GWASs.**

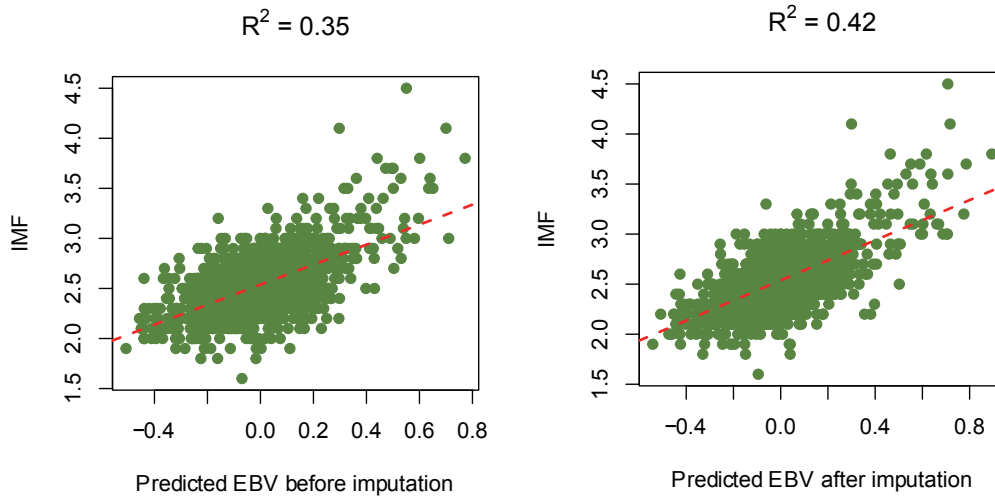
118 Upper lambda denotes inflation factor of association test results. Each dot represents a

119 variant. Red lines are fitted by a simple linear regression of function `lm()` implemented

120 in R software.

121

122



123

124

125 **Supplementary Fig. 11 The phenotypic prediction ability of EBV before and after**

126 **genotype imputation.** Significant variants at the $P < 0.01$ level from GWAS were used

127 for estimating the breeding value of IMF by cross-validated BLUP (leave-one-out). The

128 x-axis shows the predicted EBV. The y-axis shows the phenotypic value of IMF. Each

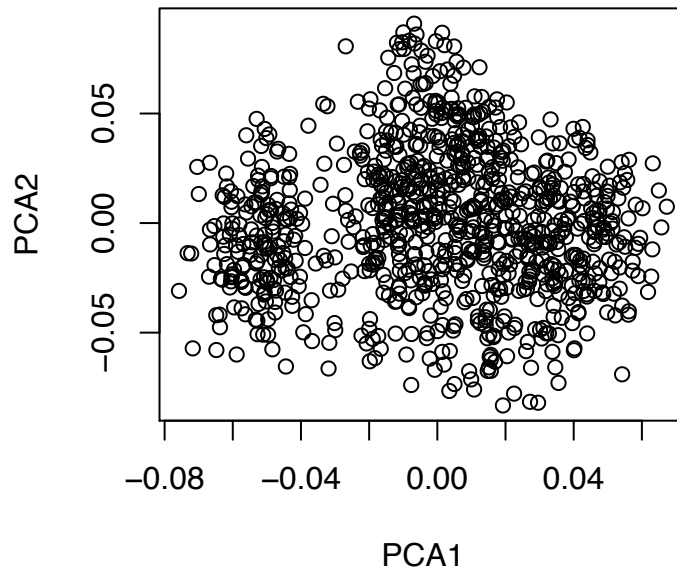
129 dot represents an individual. Pearson's correlation is employed to evaluate the

130 phenotypic prediction ability of EBV, indicated by a R^2 (squared Pearson correlation)

131 between phenotypic value and predicted EBV. The P values of correlation coefficient

132 before and after imputation are 7.3×10^{-144} and 1.3×10^{-177} , respectively.

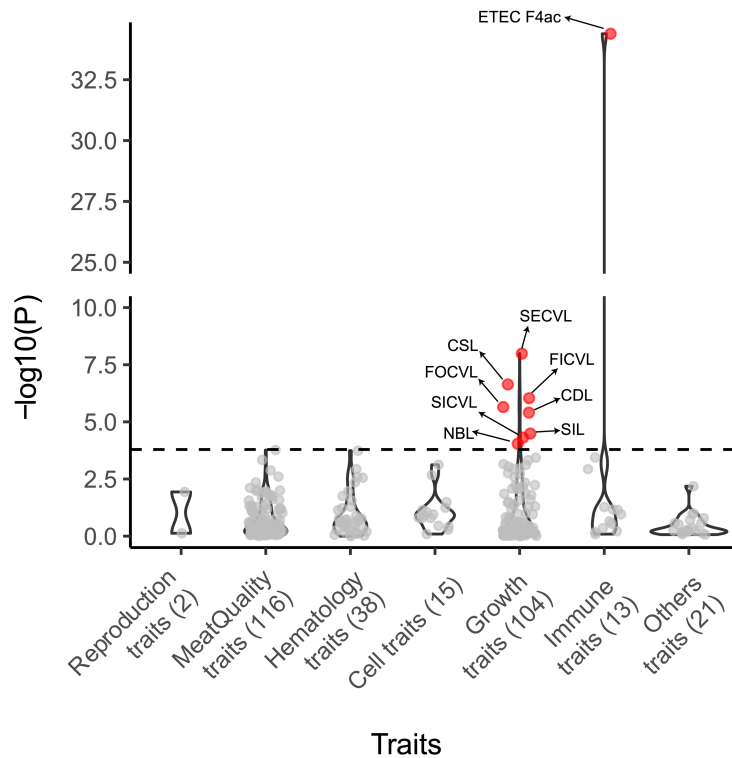
133



134

135 **Supplementary Fig. 12 The principal component analysis using the first two**
136 **principle components based on genotypes of F2 population.** Each dot represents an
137 individual.

138



139

140 **Supplementary Fig. 13 The significance of difference of 309 phenotypes between**

141 **Hap1 and Hap2 of *MUC13*.** Each point represents a phenotype. The y-axis denotes

142 statistical significance tested by Student's test (two-sided) or Pearson's Chi-squared test.

143 Reported log transformed *P* Values are nominal (i.e. not corrected for multiple testing).

144 The numbers in parentheses represent the number of traits. ETEC F4ac,

145 Enterotoxigenic *Escherichia coli* F4ac; CSL, Carcass straight length; CDL, Carcass

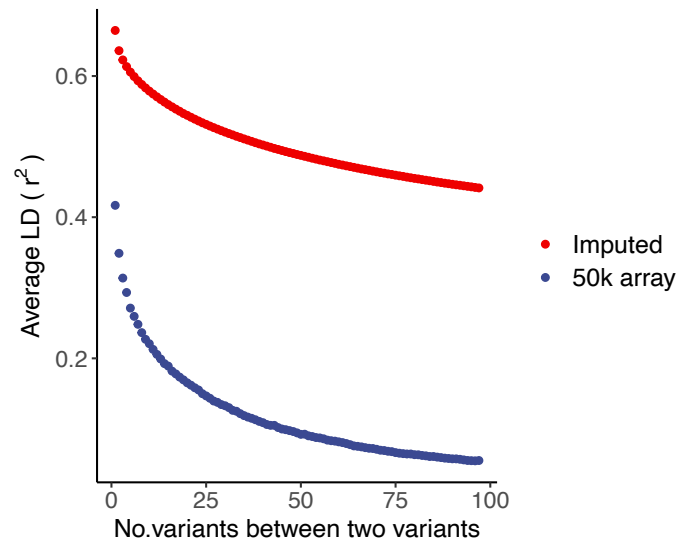
146 diagonal length; SIL, Small intestine length; FOCVL, Fourth cervical vertebra length;

147 FICVL, Fifth cervical vertebra length; SICVL, Sixth cervical vertebra length; SECVL,

148 Seventh cervical vertebra length; NBL, total length of cervical vertebra.

149

150



151

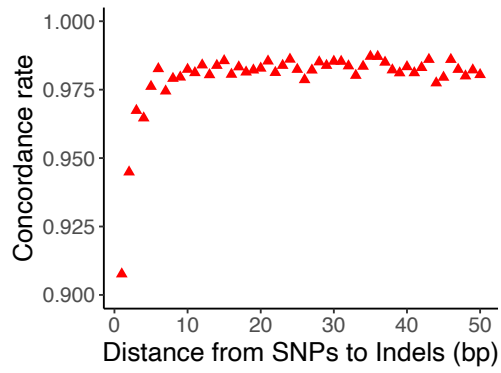
152 **Supplementary Fig. 14 Average Linkage Disequilibrium between all variant pairs.**

153 The y-axis denotes the average LD degree measured by r^2 . The x-axis denotes the

154 number of variants at chromosomes between two variants. Dots represent average LD

155 under different distances measured by the number of variants between two variants.

156



157

158 **Supplementary Fig. 15 The accuracy of SNPs that are within 1-50 bp of INDELS.**

159 Each triangle in red represents a set of SNPs near Indels. The y-axis shows the

160 concordance rate between sequence genotypes and array genotypes.

161

162 **Supplementary Table 1** Chip-based haplotype analysis for top SNPs at novel loci from
 163 imputation GWAS.

Chr	Top SNP	MAF	No.Hap	Begin.Hap	End.Hap	Df	F_value	<i>P</i> value
3	3_106129949	0.025	25	105,756,806	106,275,479	24	3.88	6.9E-10
5	5_62861562	0.012	29	62,753,754	63,543,086	28	11.11	9.0E-47
8	8_66234703	0.319	25	65,489,709	66,846,245	24	9.87	2.6E-35
9	9_10302512	0.012	12	10,180,053	10,403,658	11	5.87	1.7E-09
13	13_162125842	0.027	18	160,745,027	162,777,523	17	7.95	6.1E-20

164 Chr, chromosome; MAF, minor allele frequency; No.Hap, the number of haplotype pattern;
 165 Begin.Hap, start position of haplotype; End.Hap, end position of haplotype; Df, degree of freedom;
 166 F_value, F value in analysis of variance (ANOVA); *P* value, *P* value in ANOVA. The *P* values are
 167 one-sided. Reported *P* values are nominal (i.e. not corrected for multiple testing).

168

169 **Supplementary Table 2** The frequency of Hap1 of *MUC13* across breeds harboring at
 170 least three individuals in populations EUD, NCD, SCD, and Cross. EUD, European
 171 domestic pigs; NCD, Northwestern Chinese domestic pigs; SCD, Southeastern Chinese
 172 domestic pigs; Cross, European and Eurasian crossbred pigs.

Population	Breed	Sample size	Frequency of Hap1	
EUD	DU	29	0.52	0.67
EUD	PT	6	0.67	
EUD	LW1	62	0.69	
EUD	LW2	5	0.70	
EUD	LR	25	0.74	
EUD	YCT	11	0.86	
NCD	MIN	6	0.33	0.77
NCD	SCT1	50	0.56	
NCD	BAM	6	0.58	
NCD	TT	12	0.67	
NCD	LWU	75	0.97	
NCD	HT	6	1.00	
SCD	LUC	6	0.33	0.44
SCD	BMX	84	0.39	
SCD	EHL	132	0.40	
SCD	SCT2	12	0.42	
SCD	GST	14	0.50	
SCD	JH	6	0.50	
SCD	WZS	6	0.50	
SCD	WA	6	0.67	
SCD	YNT	12	0.79	
SCD	NJ	6	0.92	
Cross	F2	52	0.46	0.51
Cross	F1	44	0.47	
Cross	F3	44	0.47	
Cross	SUT	63	0.48	
Cross	DLY	43	0.52	
Cross	WDU	10	0.65	
Cross	LA	9	0.83	

173 BAM, Bamei; BAS, Baoshan; BMX, Bamaxiang; DU, Duroc; EHL, Erhualian; F1/F2/F3, Eurasian
 174 Crossbred; GST, Gansu Tibetan; HT, Hetao; JH, Jinhua; LA, Lean spotted pig; LR, Landrace; LUC,
 175 Luchuan; LW, LargeWhite; LWU, Laiwu; NJ, Neijiang; PT, Pietrain; SCT, Sichuan Tibetan; SUT,
 176 Sutai; TT, Tibetan Tibetan; WA, Wanan spotted; WDU, White Duroc; WZS, Wuzhishan; YCT,
 177 Yucatan; YNT, Yunnan Tibetan.

178