

## Supplementary Information for

1

2 **Genetic study of intrahepatic cholestasis of pregnancy in 101,023 Chinese women**  
3 **unveils East Asian-specific etiology linked to historic HBV infection**

4

5 Yanhong Liu<sup>1\*</sup>, Yuandan Wei<sup>1,2\*</sup>, Xiaohang Chen<sup>3\*</sup>, Shujia Huang<sup>4\*</sup>, Yuqin Gu<sup>1</sup>,  
6 Zijing Yang<sup>1,3</sup>, Liang Hu<sup>3</sup>, Xinxin Guo<sup>1</sup>, Hao Zheng<sup>1</sup>, Mingxi Huang<sup>4</sup>, Shangliang  
7 Chen<sup>5</sup>, Tiantian Xiao<sup>6</sup>, Yang Zhang<sup>1</sup>, Guo-Bo Chen<sup>7</sup>, Likuan Xiong<sup>2,8</sup>, Xiu Qiu<sup>4#</sup>,  
8 Fengxiang Wei<sup>3,9#</sup>, Jianxin Zhen<sup>2#</sup>, Siyang Liu<sup>1#</sup>

9

- 10 1. School of Public Health (Shenzhen), Sun Yat-sen University, Shenzhen,  
11 Guangdong 510006, China
- 12 2. Central Laboratory, Shenzhen Baoan Women's and Children's Hospital,  
13 Shenzhen, Guangdong 518102, China
- 14 3. The Genetics Laboratory, Longgang District Maternity & Child Healthcare  
15 Hospital of Shenzhen City, Shenzhen, Guangdong, 518172, China
- 16 4. Division of Birth Cohort Study, Guangzhou Women and Children's Medical  
17 Center, Guangzhou Medical University, Guangzhou, 510623, China
- 18 5. Department of transfusion, Shenzhen Baoan Women's and Children's Hospital,  
19 Shenzhen, Guangdong 518102, China

- 20 6. Xiangya School of Medicine, Central South University, Changsha, Hunan  
21 410078, China
- 22 7. Center for Productive Medicine, Department of Genetic and Genomic  
23 Medicine, Clinical Research Institute, Zhejiang Provincial People's Hospital,  
24 People's Hospital of Hangzhou Medical College, Hangzhou 310014, Zhejiang,  
25 China
- 26 8. Shenzhen Key Laboratory of Birth Defects Research, Shenzhen, Guangdong  
27 518102, China.
- 28 9. Longgang Maternity and Child Institute of Shantou University Medical  
29 College, Shenzhen, Guangdong, 518172, China

30 \*: Those authors contribute equally as co-first authors

31 #: Correspondence can be addressed to

32 Siyang Liu [liusy99@mail.sysu.edu.cn](mailto:liusy99@mail.sysu.edu.cn)

33 Jianxin Zhen [jxzhen@qq.com](mailto:jxzhen@qq.com)

34 Fengxiang Wei [haowei727499@163.com](mailto:haowei727499@163.com)

35 Xiu Qiu [xiu.qiu@bigcs.org](mailto:xiu.qiu@bigcs.org)

36

37	Table of contents	
38	<b><i>Supplementary Methods</i></b> .....	<b>6</b>
39	<b>GWAS study design and cohort description</b> .....	<b>6</b>
40	<b>Phenotype definition</b> .....	<b>7</b>
41	<b>Sequencing, alignment and imputation</b> .....	<b>8</b>
42	<b>Genome-wide association analysis</b> .....	<b>9</b>
43	<b>Variants Annotation</b> .....	<b>10</b>
44	<b>Replication and comparation</b> .....	<b>10</b>
45	<b>Identification of novel locus and signals</b> .....	<b>12</b>
46	<b>Colocalization analysis</b> .....	<b>12</b>
47	<b>Modern DNA-based selection test</b> .....	<b>13</b>
48	<b><i>Supplementary Reference</i></b> .....	<b>16</b>
49	<b><i>Supplementary Figures</i></b> .....	<b>19</b>
50	<b>Supplementary Fig. 1. Distribution of total bile acid concentrations in the</b>	
51	<b>Baoan and Longgang cohorts.</b> .....	<b>19</b>
52	<b>Supplementary Fig. 2. Power calculations based on the size of the cohort for</b>	
53	<b>the combined meta-analysis.</b> .....	<b>20</b>

54	<b>Supplementary Fig. 3. QQ plots of TBA and ICP GWAS meta-analyses during</b>	
55	<b>gestational weeks 13-42. ....</b>	<b>21</b>
56	<b>Supplementary Fig. 4. Internal replication of TBA and ICP GWAS. ....</b>	<b>22</b>
57	<b>Supplementary Fig. 5. External replication of TBA and ICP meta-GWAS</b>	
58	<b>using independent Chinese population cohorts. ....</b>	<b>23</b>
59	<b>Supplementary Fig. 6. LocusZoom plots of genome-wide significant loci</b>	
60	<b>associated with the TBA trait investigated in the study. ....</b>	<b>25</b>
61	<b>Supplementary Fig. 7. LocusZoom plot of genome-wide significant loci</b>	
62	<b>associated with ICP. ....</b>	<b>27</b>
63	<b>Supplementary Fig. 8. Comparison of effect sizes between TBA and ICP</b>	
64	<b>GWAS with European ICP GWAS results. ....</b>	<b>28</b>
65	<b>Supplementary Fig. 9. Genome-wide association study of HBV and its antigens</b>	
66	<b>&amp; antibodies during pregnancy. ....</b>	<b>29</b>
67	<b>Supplementary Fig. 10. Stacked LocusZoom plot of TBA, ICP and HBV</b>	
68	<b>related traits. ....</b>	<b>31</b>
69	<b>Supplementary Fig. 11. LocusCompare plot of GWAS-GWAS colocalization</b>	
70	<b>between TBA &amp; ICP and HBV. ....</b>	<b>33</b>
71	<b>Supplementary Fig. 12. Geographical allele frequency distribution and linkage</b>	
72	<b>disequilibrium (LD) of rs137983251-G, rs138089855-C, rs47525203 and</b>	
73	<b>rs2296651-A. ....</b>	<b>35</b>

74	<b>Supplementary Fig. 13. Site Frequency Spectrum-based selection tests with</b>	
75	<b>modern DNA samples from 1000 Genomes Project across three populations.</b>	<b>36</b>
76	<b>Supplementary Fig. 14. The plot of LSBL, iHS and XP-EHH for the</b>	
77	<b>chromosome 14 region based on data from the 1000GP.</b>	<b>38</b>
78	<b>Supplementary Fig. 15. Haplotype structure of the ICP risk 14q24.1 locus in</b>	
79	<b>CHS, CEU and YRI populations.</b>	<b>40</b>
80	<b><i>Supplementary Table Legends</i></b>	<b>41</b>
81		
82		

83 **Supplementary Methods**

84 **GWAS study design and cohort description**

85 Previous study has pointed that the more accuracy for presenting population genetic  
86 variation and inferencing population structure is achieved with very large sample  
87 sizes, even when utilizing extremely low sequencing depth<sup>1</sup>. Furthermore, previous  
88 study has demonstrated the substantial value of non-invasive prenatal testing (NIPT)  
89 data for GWAS analysis<sup>2</sup>.

90

91 A total of 121,556 Chinese pregnancies who underwent non-invasive prenatal testing  
92 (NIPT) at two hospital cohorts in Shenzhen, China, were enrolled in our study  
93 between 2017 and 2022. The Baoan study cohort included 70,608 pregnant women  
94 recruited at the Shenzhen Baoan Women's and Children's Healthcare Hospital  
95 (Shenzhen, China). Additionally, 50,948 pregnant women who visited Shenzhen  
96 Longgang District Maternity & Child Healthcare Hospital (Shenzhen, China) were  
97 recruited for the Longgang study cohort. Here, we excluded 5,625 pregnancies  
98 potentially involving multiple gestations, along with 14,908 pregnancies for which  
99 TBA concentration results were unavailable during this period.

100

101 Using the method described in the previous study<sup>2</sup>, maternal genotypes were inferred  
102 from NIPT ultra-low-depth whole-genome sequencing data. Phenotypic data were  
103 obtained from the hospital's electronic medical system during routine pregnancy

104 screening program. Subsequently, we combined genotype data with clinical  
105 phenotypic data to explore the genetic molecular basis of TBA and ICP.

106

107 This study was reviewed and approved by Ethics Committee of School of Public  
108 Health (Shenzhen), Sun Yat-Sen University (2021. No.8), as well as the Institutional  
109 Board of Shenzhen Baoan Women's and Children's Hospital  
110 (LLSC2021-04-01-10-KS) and Longgang District Maternity and Child Healthcare  
111 Hospital of Shenzhen City (LGFYYXLLL-2022-024). The study strictly adhered to  
112 regulations governing ethical considerations and personal data protection. Data  
113 collection was approved by the Human Genetic Resources Administration of China  
114 (HGRAC). Written informed consent of all participants were obtained.

115

#### 116 **Phenotype definition**

117 During the GWAS analysis, Hepatic biochemistry test data, including total bile acid  
118 concentrations, were obtained during second and third trimester pregnancy (13-42  
119 gestational weeks) as part of routine pregnant screening in two hospital cohorts. The  
120 peak level for each individual during this period was defined as TBA quantitative trait  
121 (N = 94,360). ICP cases were identified as pregnancy with TBA concentrations  $\geq$   
122 10 $\mu$ mol/L during the same period. In total, 4,703 cases and 96,320 controls for ICP  
123 were included in the study. The characteristics of the participants are presented in  
124 **Table S1.**

125

126 We also obtained the infection markers, including hepatitis B surface antigen  
127 (HBsAg), hepatitis B surface antibody (HBsAb), hepatitis B e antigen (HBeAg),  
128 hepatitis B e antibody (HBeAb) and hepatitis B core antibody (HBcAb), from  
129 pregnant screening. The sample size for HBsAg is 44,432 (7,575 cases versus 86866  
130 controls), while the sample sizes for the other four traits are all 94,441 each (cases  
131 versus controls are 63627/30834, 1817/92644, 16117/78344, and 23618/70843,  
132 respectively). HBV persistent carriers and spontaneously recovered subjects were  
133 categorized as the HBV cases. HBV persistent carriers were defined as individuals  
134 positive for both HBsAg and HBcAb but negative for hepatitis C virus (HCV)  
135 antibody (anti-HCV). Spontaneously recovered subjects were those who were  
136 negative for HBsAg and anti-HCV but positive for both HBsAb and HBcAb. HBV  
137 controls were individuals negative for all HBV markers, including HBsAg, HBsAb,  
138 HBeAg, HBeAb, and HBcAb, as well as negative for HCV. In total, 21,770 cases and  
139 22,662 controls of HBV were recruited in the study.

140

#### 141 **Sequencing, alignment and imputation**

142 We collected sequencing data from non-invasive prenatal testing in both the Baoan  
143 and Longgang cohorts. The sequencing protocol details were outlined in Zhang et al<sup>3</sup>  
144 & Liu et al<sup>2</sup>. In summary, each participant underwent whole-genome sequencing,  
145 resulting in 9.9-21.9 million cleaned reads, corresponding to sequencing depths of



146 approximately 0.11x-0.25x, with an average of 0.17x. Next, the single-end read  
147 alignment option in BWA was used to align the cleaned reads to the hg38 human  
148 genome reference<sup>4</sup>. The rmdup option in samtools was used to remove potential PCR  
149 duplicates<sup>5</sup>. The realign and base quality recalibration method in GATK were used to  
150 realign the reads and to recalibrate base quality score<sup>6</sup>. Finally, the alignment files  
151 were stored as bam files. After alignment, we employed GLIMPSE<sup>7</sup> (version 1.1.1) to  
152 impute genotype probabilities for all 121,556 individuals with a 10k Chinese  
153 reference panel. All alignment and imputation process were conducted at the National  
154 Supercomputing Center in GuangZhou.

155

#### 156 **Genome-wide association analysis**

157 The GWAS analysis employed the multiple linear regression model for TBA and the  
158 logistic regression model for ICP to examine the association of SNPs using PLINK2.0  
159 (<https://www.cog-genomics.org/plink/2.0/>). Gestational week, maternal age and the  
160 top ten principal components accounting for population stratification were included as  
161 covariates for TBA and ICP. Otherwise, we also conducted GWAS using logistic  
162 regression model for HBV and its antigen & antibody, with maternal age and the top  
163 ten principal components included as covariates. The additive genetic model of SNP  
164 dosage was utilized for genetic-phenotypic association. Principal component analysis  
165 (PCA) was conducted using PLINK2.0<sup>8</sup> on the dataset of 121,556 individuals.

166

167 The results of meta-GWAS were visualized with Manhattan plot using the R package  
168 ggplot2 (version 3.4.2) (<https://cran.r-project.org/web/packages/ggplot2/index.html>)  
169 and ggbreak (version 0.1.2) (<https://cran.r-project.org/web/packages/ggbreak/index.html>)<sup>9</sup>. Quantile-quantile (QQ) plots were generated using the observed and expected  
170  $-\log_{10}(P \text{ value})$  with R package ggplot2. The genomic inflation factor ( $\lambda$ ),  
171 calculated based on the 50th percentile, was 1.069 for TBA and 1.018 for ICP  
172 separately, indicating no significant population stratification. Regional high-resolution  
173 association plots showing the LD between markers in the lead loci were generated  
174 using LocusZoom (version 1.4) (<http://locuszoom.org/>).

176

### 177 **Variants Annotation**

178 Gene annotation was conducted using the Ensembl Variant Effect Predictor<sup>10</sup> (VEP,  
179 version 101), with indexed GRCh38 cache files (version 109). All the data utilized for  
180 annotation were obtained from the Ensembl FTP server (<https://ftp.ensembl.org/pub/>).  
181 Based on a set of VEP default criteria, the "--pick" option was used to assign a single  
182 consequence block to each variant. For variants in the intergenic region, "--nearest"  
183 option was used to identify the nearest gene with a protein-coding transcription start  
184 site (TSS) for variants.

185

### 186 **Replication and comparison**

187 As for external replication, we examined the genetic influence on TBA and ICP with  
 188 two independent study cohorts (Baoan NIPT Plus cohort and BIGCS cohort). SNPs  
 189 meeting the following criteria were regarded as replicated: 1) they exhibited a  
 190 consistent direction of effect for lead SNPs with that of cohorts, and 2) they reached  
 191 Bonferroni-corrected  $P$  values or passed a two-sided two-sample t-test.

192

193 The two-sided two-sample t-test was conducted to evaluate the equivalence of genetic  
 194 effects on the same traits between two independent cohorts, with the following  
 195 hypotheses:

$$\text{Null hypothesis } H_0: \beta_m = \beta_i$$

196

$$\text{Alternative hypothesis } H_1: \beta_m \neq \beta_i$$

197 The  $T$  statistic was computed as follows:

$$T = \frac{\beta_m - \beta_i}{\sqrt{\frac{S_m^2}{n_m} + \frac{S_i^2}{n_i}}} = \frac{\beta_m - \beta_i}{\sqrt{SE_m^2 + SE_i^2}} \sim t(v') \quad (1)$$

198 The degrees of freedom  $v'$  was determined by the formula:

$$v' = \frac{(SE_m^2 + SE_i^2)^2}{\frac{SE_m^4}{n_m - 1} + \frac{SE_i^4}{n_i - 1}} \quad (2)$$

199 Herein,  $\beta_m$  and  $\beta_i$  represent the genetic effects associated with the same TBA and  
 200 ICP traits for two cohorts, respectively.  $S_m^2$  and  $S_i^2$  denote sample variance, whereas  
 201  $SE_m$  and  $SE_i$  stand for estimated standard errors. It is established that the  $T$  statistic  
 202 in equation (1) follows a  $t$ -distribution with a degree of freedom  $v'$ . To address

203 potential inequality between  $S_m^2$  and  $S_i^2$ , the adjusted  $v'$  was employed, computed  
204 according to formula (2).

205

206 If the lead SNP did not exist external GWAS summary datasets, the proxy SNP of the  
207 lead SNP with LD  $R^2$  greater than 0.8 and existed in both data was chosen as a  
208 substitute. Proxy SNPs were queried using the LD proxy Tool  
209 (<https://ldlink.nih.gov/?tab=ldproxy>) through LDlink<sup>11</sup> (version 5.5.1) based on  
210 GRCh38 1000 Genomes Project (1000GP) genome build in East Asian (EAS)  
211 populations.

212

### 213 **Identification of novel locus and signals**

214 We identified SNPs as novel locus if no SNPs in GWAS Catalog within 1Mb block of  
215 the SNP were reported to be associated with our results. If there were SNPs has been  
216 reported before associated with the SNP in 1Mb region, LDpair  
217 (<https://ldlink.nih.gov/?tab=ldpair>) was used to calculate the linkage disequilibrium  $R^2$   
218 with lead SNP. If lead SNPs identified in our study were novel SNP with  $R^2 < 0.2$   
219 seem as a novel signal in a known locus.

220

### 221 **Colocalization analysis**

222 Colocalization evaluates the posterior probabilities of five mutually exclusive  
223 hypotheses: 1)  $H_0$ : no association of any variant in the region with either trait; 2)  $H_1$ :

224 association with first trait but not the second; 3) H<sub>2</sub>: association with second trait but  
225 not the first; 4) H<sub>3</sub>: associated with both traits but have two independent causal  
226 variants and 5) H<sub>4</sub>: associated with both traits and shared one causal variant<sup>12</sup>.  
227 Colocalization analysis has been originally designed for testing two sets of  
228 associations measured on different individuals. While, previous study has confirmed  
229 by simulation that the results running it on the same individuals appear robust to  
230 correlated errors<sup>13</sup>.

231

232 Here, we utilized PP4 (posterior probability that there exists a single causal variant  
233 common to both traits)  $\geq 0.75$  and  $PP4/PP3 \geq 3$  to identify colocalization between the  
234 GWAS and GWAS signals<sup>12</sup>.

235

### 236 **Modern DNA-based selection test**

237 The high-coverage 1000 Genomes Project phased whole-genome sequencing (WGS)  
238 panel  
239 ([http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data\\_collections/1000G\\_2504\\_high\\_cover  
240 age/working/20201028\\_3202\\_phased/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/1000G_2504_high_coverage/working/20201028_3202_phased/)), comprising 3,202 individuals from 26  
241 worldwide populations, was employed for evolutionary analysis. Ancestral allele  
242 information was obtained from 1000 Genomes Project  
243 (<http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/>). Python package  
244 “CrossMap.py” (<https://crossmap.readthedocs.io/en/latest/>)<sup>14</sup> was utilized to change

245 chromosome positions from hg37 to hg38. The sample HG01783 was excluded as it  
246 belongs to both European and African ancestry. SNPs with minor allele frequencies  
247 (MAF) < 0.01 in all five super-populations and genetic variants without ancestral  
248 information were excluded. In total, 18,221,282 SNPs were used for the natural  
249 analysis.

250

251 For SFS-based tests, genetic diversity ( $\pi$ ), Tajima's  $D^{15}$ , and Fay and Wu's  $H^{16}$  were  
252 calculated using Perl scripts from a previously published paper<sup>17</sup>. These three  
253 statistics were computed with a sliding-window approach (window size = 5 kb and  
254 moving step = 1 kb). Statistical significance for these three statistics were evaluated  
255 using the genome-wide empirical distribution. Based on allele frequency  
256 differentiation and extended haplotype homozygosity, we calculated the  
257 locus-specific branch length (LSBL)<sup>18</sup> for the CHS population with the European  
258 population (i.e., CEU) and the African population (i.e., YRI) as reference populations.  
259 The formula used to calculate LSBL was:  $LSBL = (F_{ST\ CHS\_CEU} + F_{ST\ CHS\_YRI} - F_{ST\ CEU\_YRI})/2^{19}$ . The 1% threshold for the whole genome was set at 0.40, and the 0.1%  
260 threshold was set at 0.58. SNPs within the top 1% were considered highly  
261 differentiated, and those within the top 0.1% were considered extremely highly  
262 differentiated.  
263 differentiated.

264

265 We further integrated Haplotype Score (iHS) and the Cross Population Extended  
266 Haplotype Homozygosity (XP-EHH) test using R package rehh 2.0<sup>20</sup>. For iHS, we  
267 calculated the iHS value for each locus in the CHS population, considering SNPs with  
268  $|iHS| > 2$  as exhibiting a signal of positive selection<sup>21</sup>. As for XP-EHH, the score for  
269 each locus in the CHS population was calculated with CEU as the reference  
270 population, and SNPs with  $XP-EHH > 2$  were considered to exhibit a signal of  
271 positive selection<sup>22</sup>.  
272

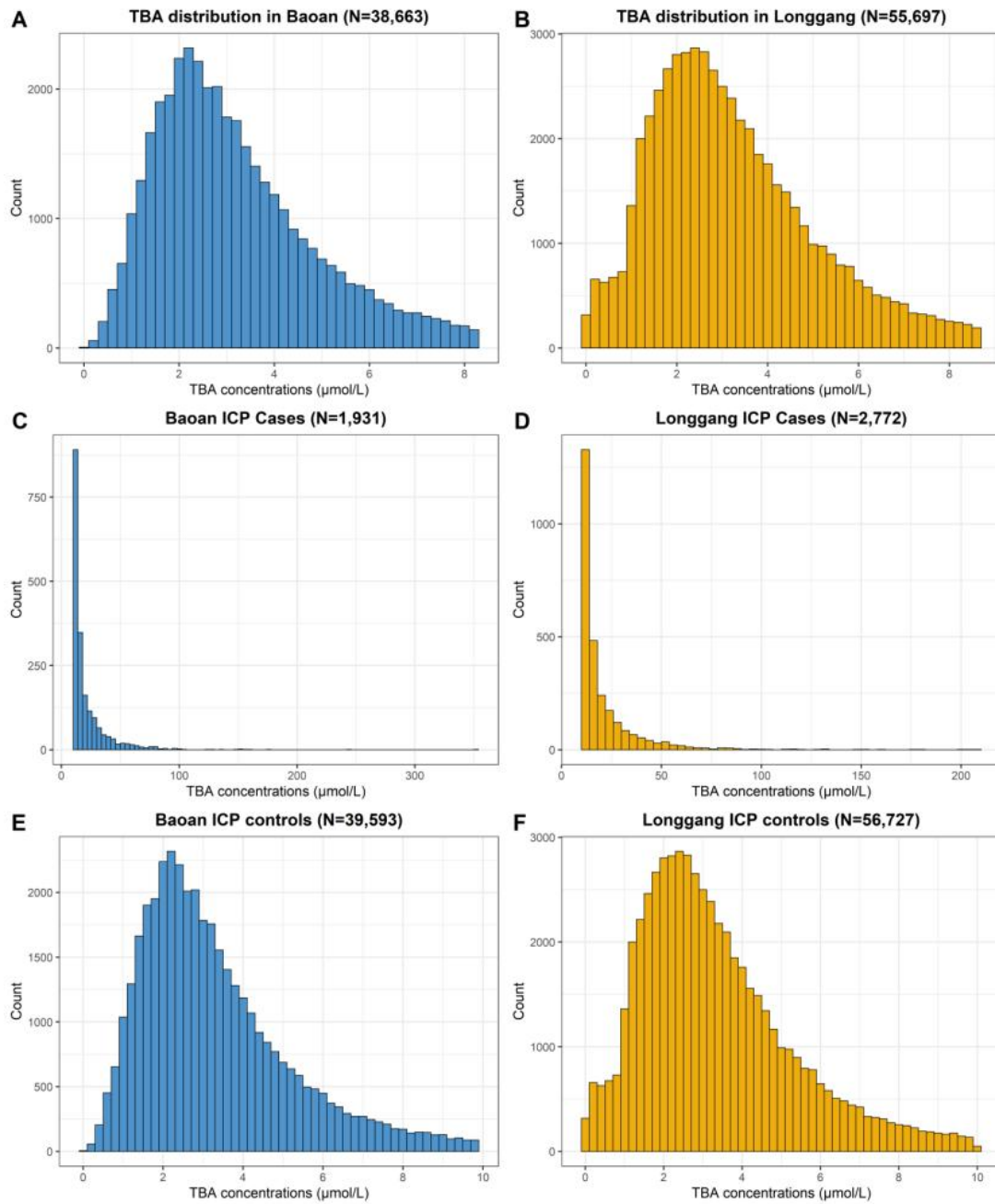
273 **Supplementary Reference**

- 274 1. Fumagalli, M. Assessing the effect of sequencing depth and sample size in  
275 population genetics inferences. *PLoS One* **8**, e79667 (2013).
- 276 2. Liu, S. *et al.* Genomic Analyses from Non-invasive Prenatal Testing Reveal  
277 Genetic Associations, Patterns of Viral Infections, and Chinese Population  
278 History. *Cell* **175**, 347-359.e14 (2018).
- 279 3. Zhang, H. *et al.* Non-invasive prenatal testing for trisomies 21, 18 and 13:  
280 Clinical experience from 146 958 pregnancies. *Ultrasound Obstet. Gynecol.* **45**,  
281 530–538 (2015).
- 282 4. Li, H. & Durbin, R. Fast and accurate short read alignment with  
283 Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
- 284 5. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools.  
285 *Bioinformatics* **25**, 2078–2079 (2009).
- 286 6. DePristo, M. A. *et al.* A framework for variation discovery and genotyping  
287 using next-generation DNA sequencing data. *Physiol. Behav.* **43**, 491–498  
288 (2011).
- 289 7. Rubinacci, S., Ribeiro, D. M., Hofmeister, R. J. & Delaneau, O. Efficient  
290 phasing and imputation of low-coverage sequencing data using large reference  
291 panels. *Nat. Genet.* **53**, 120–126 (2021).
- 292 8. Purcell, S. *et al.* PLINK: A tool set for whole-genome association and  
293 population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).



- 294 9. Xu, S. *et al.* Use ggbreak to Effectively Utilize Plotting Space to Deal With  
295 Large Datasets and Outliers. *Front. Genet.* **12**, 774846 (2021).
- 296 10. Cunningham, F. *et al.* Ensembl 2022. *Nucleic Acids Res.* **50**, D988–D995  
297 (2022).
- 298 11. Machiela, M. J. & Chanock, S. J. LDlink: A web-based application for  
299 exploring population-specific haplotype structure and linking correlated alleles  
300 of possible functional variants. *Bioinformatics* **31**, 3555–3557 (2015).
- 301 12. Giambartolomei, C. *et al.* Bayesian Test for Colocalisation between Pairs of  
302 Genetic Association Studies Using Summary Statistics. *PLoS Genet.* **10**,  
303 e1004383 (2014).
- 304 13. Mitchelmore, J., Grinberg, N. F., Wallace, C. & Spivakov, M. Functional  
305 effects of variation in transcription factor binding highlight long-range gene  
306 regulation by epromoters. *Nucleic Acids Res.* **48**, 2866–2879 (2020).
- 307 14. Zhao, H. *et al.* CrossMap: A versatile tool for coordinate conversion between  
308 genome assemblies. *Bioinformatics* **30**, 1006–1007 (2014).
- 309 15. Tajima, F. Statistical Method for Testing the Neutral Mutation Hypothesis by  
310 DNA Polymorphism. *Pharmatherapeutica* **123**, 585–595 (1989).
- 311 16. Fay, J. C. & Wu, C. I. Hitchhiking under positive Darwinian selection.  
312 *Genetics* **155**, 1405–1413 (2000).

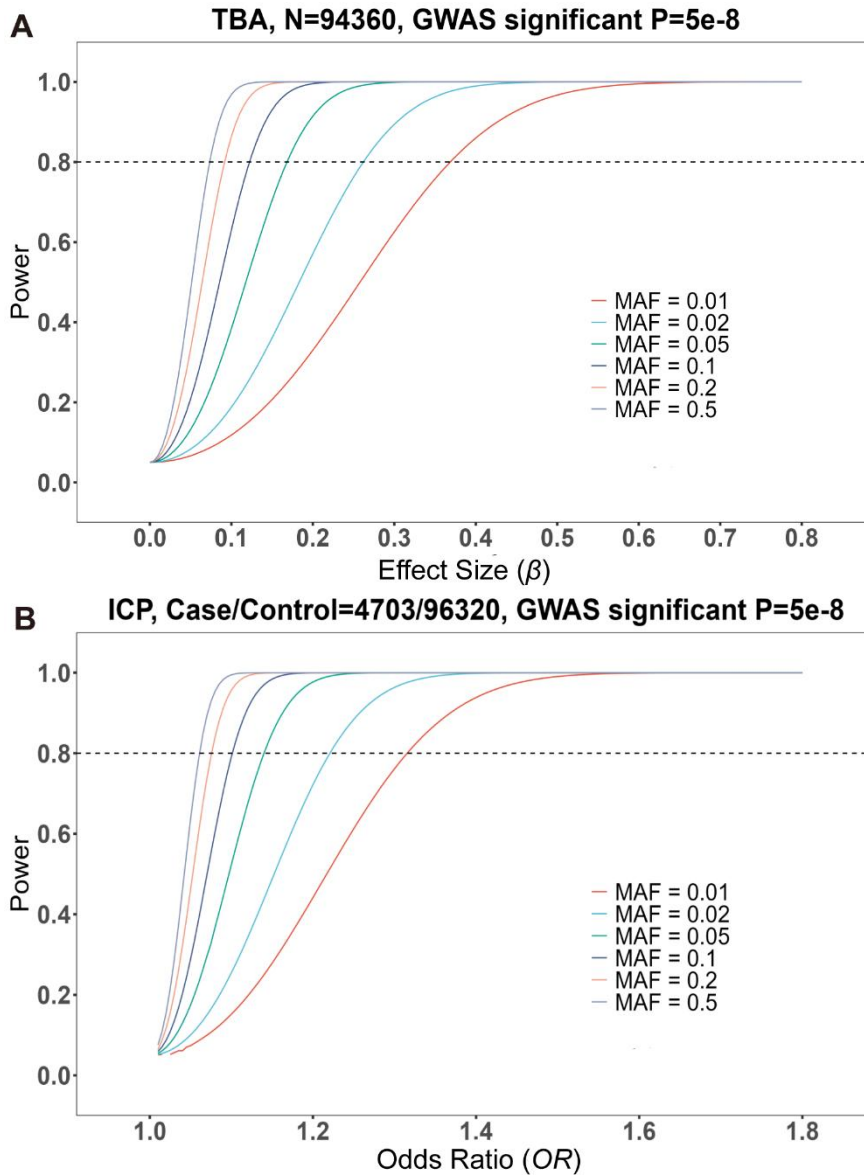
- 313 17. Ye, K., Gao, F., Wang, D., Bar-Yosef, O. & Keinan, A. Dietary adaptation of  
314 FADS genes in Europe varied across time and geography. *Nat. Ecol. Evol.* **1**,  
315 167 (2017).
- 316 18. Shriver, M. D. *et al.* The genomic distribution of population substructure in  
317 four populations using 8,525 autosomal SNPs. *Hum. Genomics* **1**, 274–286  
318 (2004).
- 319 19. Ma, X. & Xu, S. Archaic introgression contributed to the pre-agriculture  
320 adaptation of vitamin B1 metabolism in East Asia. *iScience* **25**, 105614 (2022).
- 321 20. Gautier, M. & Vitalis, R. Rehh An R package to detect footprints of selection  
322 in genome-wide SNP data from haplotype structure. *Bioinformatics* **28**,  
323 1176–1177 (2012).
- 324 21. Voight, B. F., Kudravalli, S., Wen, X. & Pritchard, J. K. A map of recent  
325 positive selection in the human genome. *PLoS Biol.* **4**, e72 (2006).
- 326 22. Sabeti, P. C. *et al.* Genome-wide detection and characterization of positive  
327 selection in human populations. *Nature* **449**, 913–918 (2007).



329

330 **Supplementary Fig. 1. Distribution of total bile acid concentrations in the Baoan**  
 331 **and Longgang cohorts.**

332 Panel (A) and (B) depict TAB levels among normal pregnancies (TBA concentration  
 333  $< 10\mu\text{mol/L}$ ) for each of the two hospitals. Panel (C) and (D) depict TAB levels  
 334 among cases for each of the two hospitals. Panel (E) to (F) illustrate the TBA  
 335 concentrations of cases and controls within the two hospitals.



336

337 **Supplementary Fig. 2. Power calculations based on the size of the cohort for the**  
 338 **combined meta-analysis.**

339 (A) Power analysis for the meta-GWAS of TBA with a sample size of 94,360. (B)

340 Power analysis for the meta-GWAS of ICP with 4,703 cases versus 96,320 controls.

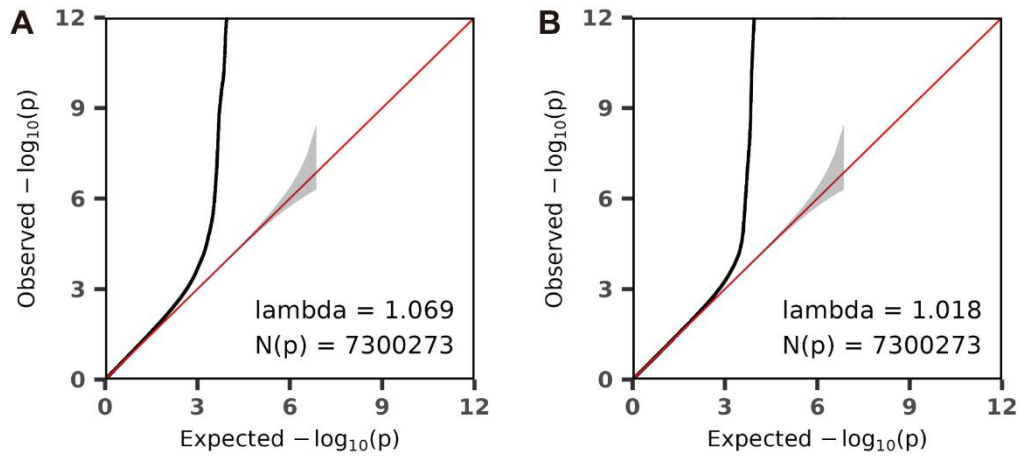
341 The graph illustrates the ability to detect genome-wide associations at a significance

342 threshold of  $P$  value  $< 5 \times 10^{-8}$  for varying odds ratio (x-axis) and minor allele

343 frequencies (MAF). Power calculations were performed using a linear or logistic

344 model under genetic additivity.

345

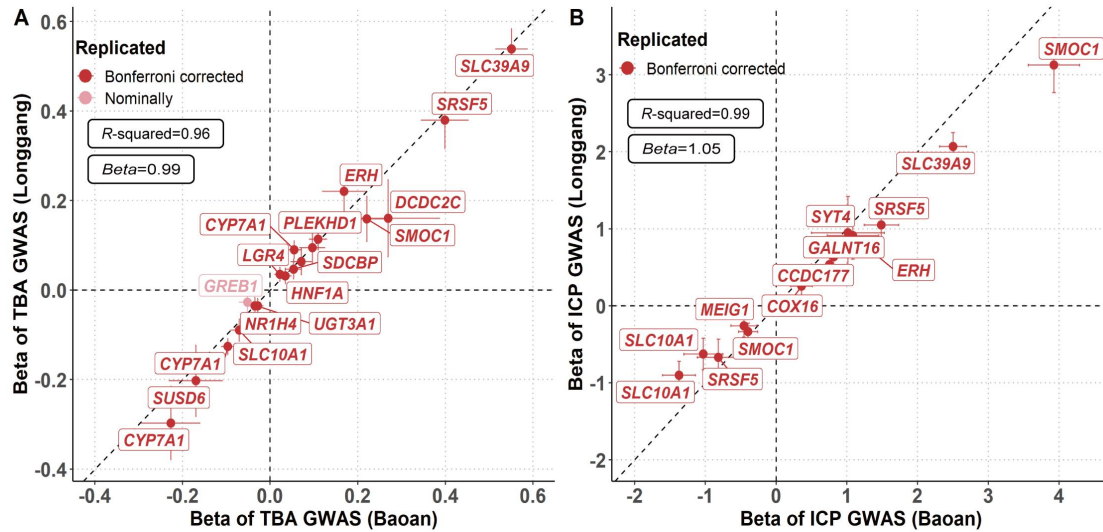


346

347 **Supplementary Fig. 3. QQ plots of TBA and ICP GWAS meta-analyses during**  
348 **gestational weeks 13-42.**

349 QQ plots for GWAS meta-analyses of (A) TBA and (B) ICP depict the relationship  
350 between observed and expected  $-\log_{10}(P)$  values. The red line represents the  
351 distribution of  $P$  values under the null hypothesis, and the gray shaded area indicates  
352 stander errors. The genomic inflation (lambda) is shown in the QQ plots, indicating  
353 no significant population stratification.

354

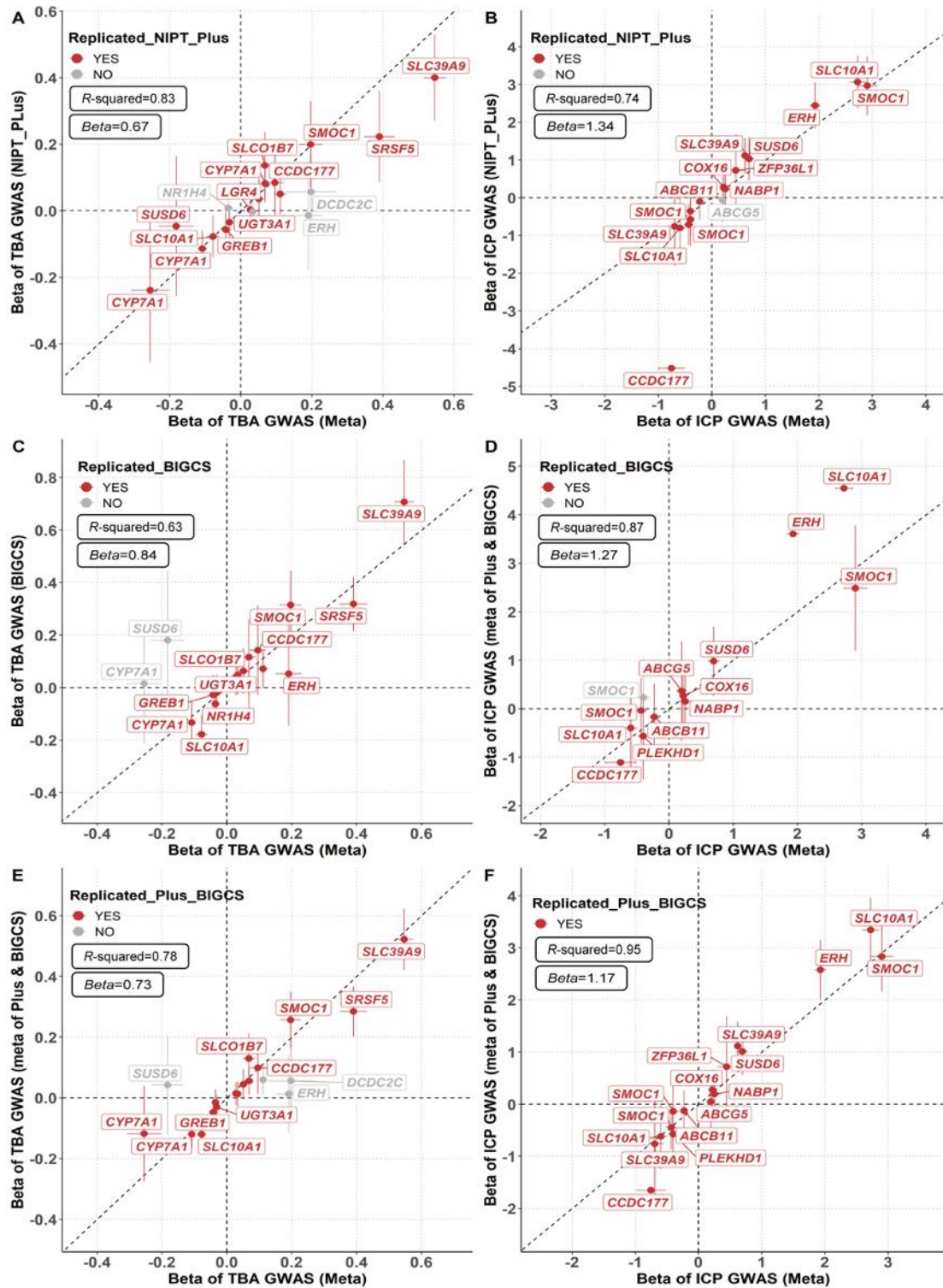


355

356 **Supplementary Fig. 4. Internal replication of TBA and ICP GWAS.**

357 In (A) and (B), the beta and  $P$  values for TBA and ICP traits were replicated in two  
 358 internal cohorts (Baoan and Longgang). The x-axis shows the beta values for TBA or  
 359 ICP GWAS in the Baoan cohort, while the y-axis represents the beta values for the  
 360 Longgang cohort. The error bars indicate the 95% confidence interval of beta. Red  
 361 points denote SNPs selected from GCTA analysis with consistent directions in beta  
 362 and achieved Bonferroni corrected significant  $P$  values. Pink points indicate SNPs  
 363 with the same beta and nominally significant  $P$  values. The Bonferroni significant  
 364 threshold was calculated as 0.05 divided by the number of independent loci for traits.  
 365 Detail data can be found in **Supplementary Table 2**.

366



367

368

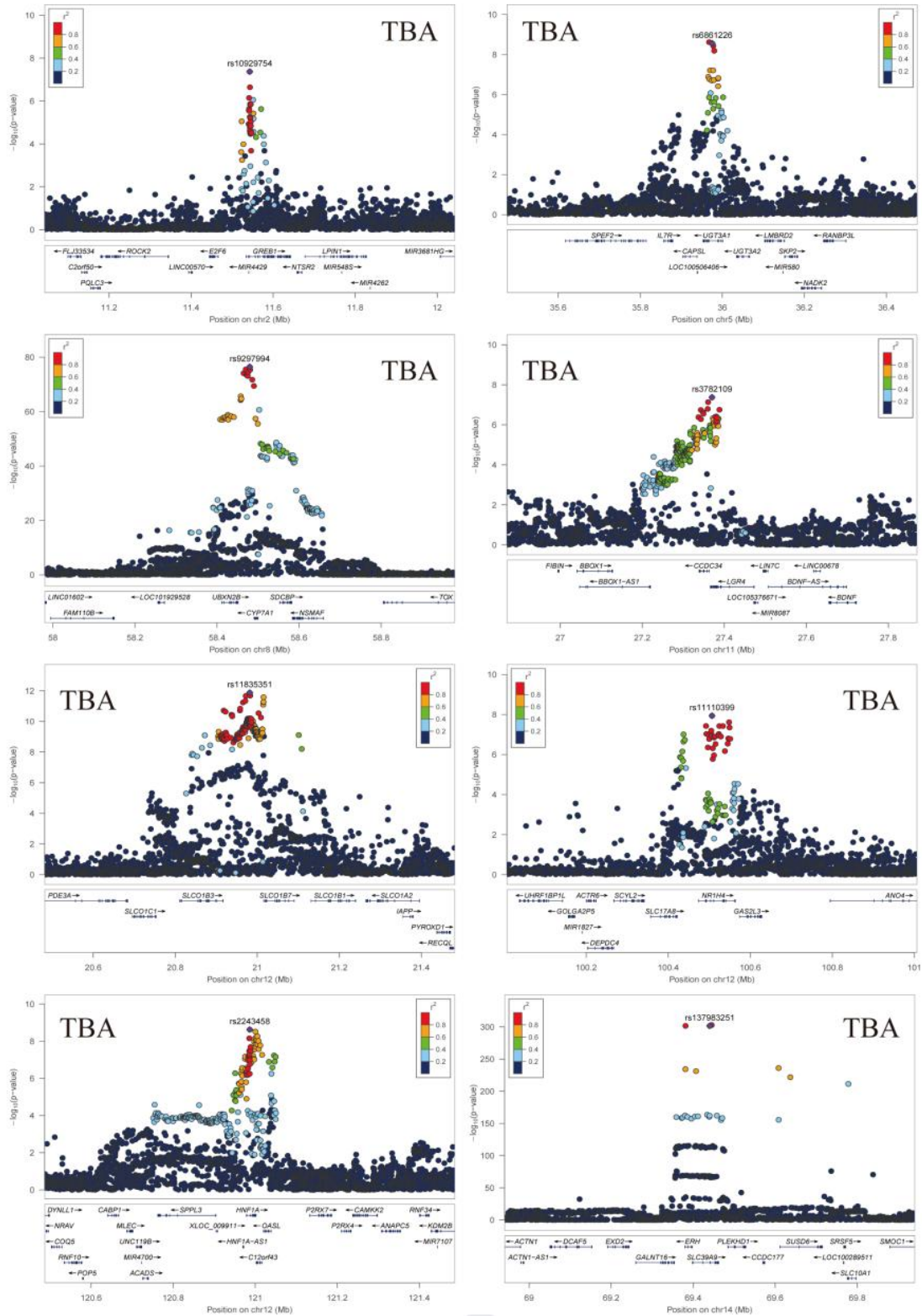
369 **Supplementary Fig. 5. External replication of TBA and ICP meta-GWAS using**  
 370 **independent Chinese population cohorts.**

371 Panel (A) - (F) compared the meta-GWAS results of TBA and ICP against Baoan

372 NIPT PLUS, BIGCS cohorts, and the GWAS-meta of Baoan NIPT PLUS and BIGCS

373 cohorts. The x-axis depicts the beta values of TBA or ICP meta-GWAS, while the  
374 y-axis represents the beta values of two independent cohorts or the result of  
375 meta-analysis of these two independent cohorts. The error bars denote the 95%  
376 confidence interval of beta. Red points denote SNPs selected from GCTA that meet  
377 the following criteria: 1) exhibit consistent direction in beta, and 2) attain a significant  
378 GWAS  $P$  value in an external dataset after Bonferroni correction or pass the T-test  
379 with  $P$  value  $> 0.05$ . Grey points indicate SNPs that do not meet any of these criteria.  
380 Detailed data, excluding empty and proxy loci, are available in **Supplementary**  
381 **Table 2.**





382

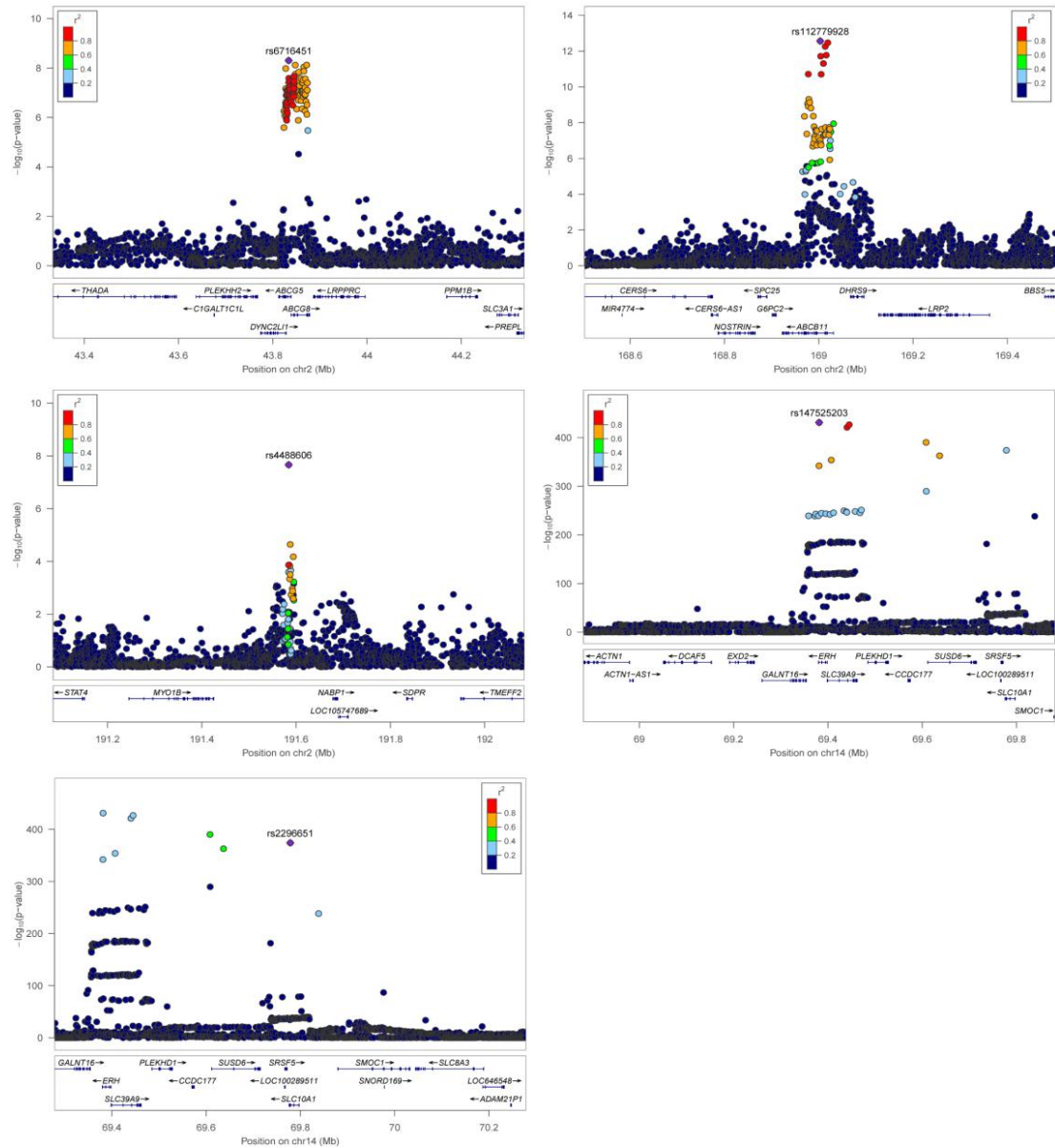
383 **Supplementary Fig. 6. LocusZoom plots of genome-wide significant loci**

384 **associated with the TBA trait investigated in the study.**

385 For the eight lead SNPs associated with TBA level (excluding two SNPs, rs59789496

386 and rs71423384, which appear not credible), regional association and linkage

387 disequilibrium (LD) plots were presented. These plots encompass the upstream and  
388 downstream 500kb flanking region of each lead SNP. The x-axis shows chromosome  
389 positions with respect to GRCh38, and the y-axis indicates  $-\log_{10}(P)$  values for the  
390 associated tests. The purple diamond represents the lead SNP of each locus, while  
391 other SNPs are color-coded based on their LD ( $r^2$ ) with the lead SNP. The plots were  
392 generated using LocusZoom software.  
393



394

395 **Supplementary Fig. 7. LocusZoom plot of genome-wide significant loci**

396 **associated with ICP.**

397 For the five lead SNPs associated with ICP (exclude two SNP, rs74573797 and

398 rs1951510, which appear not credible), regional association and linkage

399 disequilibrium (LD) plots were presented. These plots encompass the upstream and

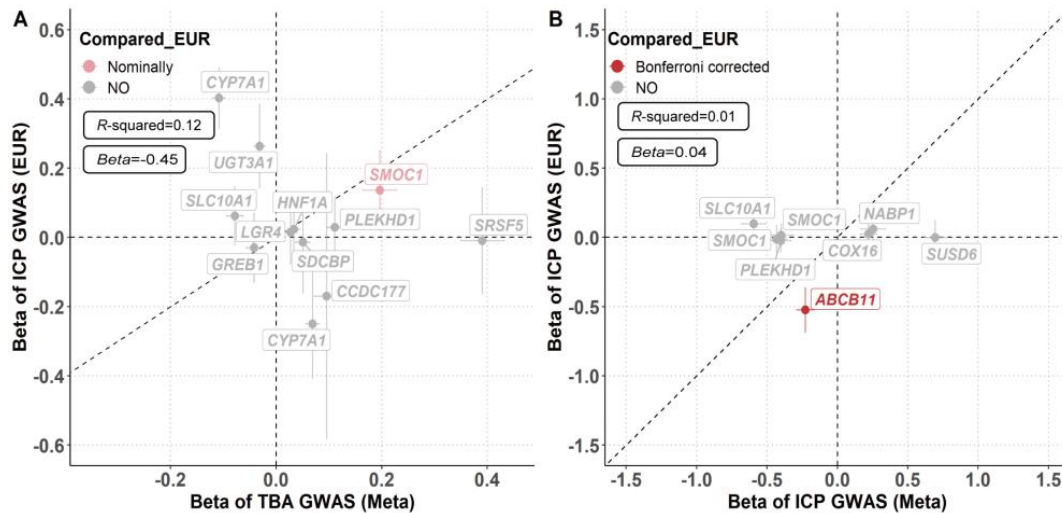
400 downstream 500kb flanking region of each lead SNP. The x-axis displays

401 chromosome positions with respect to GRCh38, while the y-axis indicates  $-\log_{10}(P)$

402 values for the associated tests. The purple diamond represents the lead SNP for each

403 locus, and other SNPs are color-coded based on their LD  $r^2$  with the lead SNP. Plots

404 were generated using LocusZoom software.

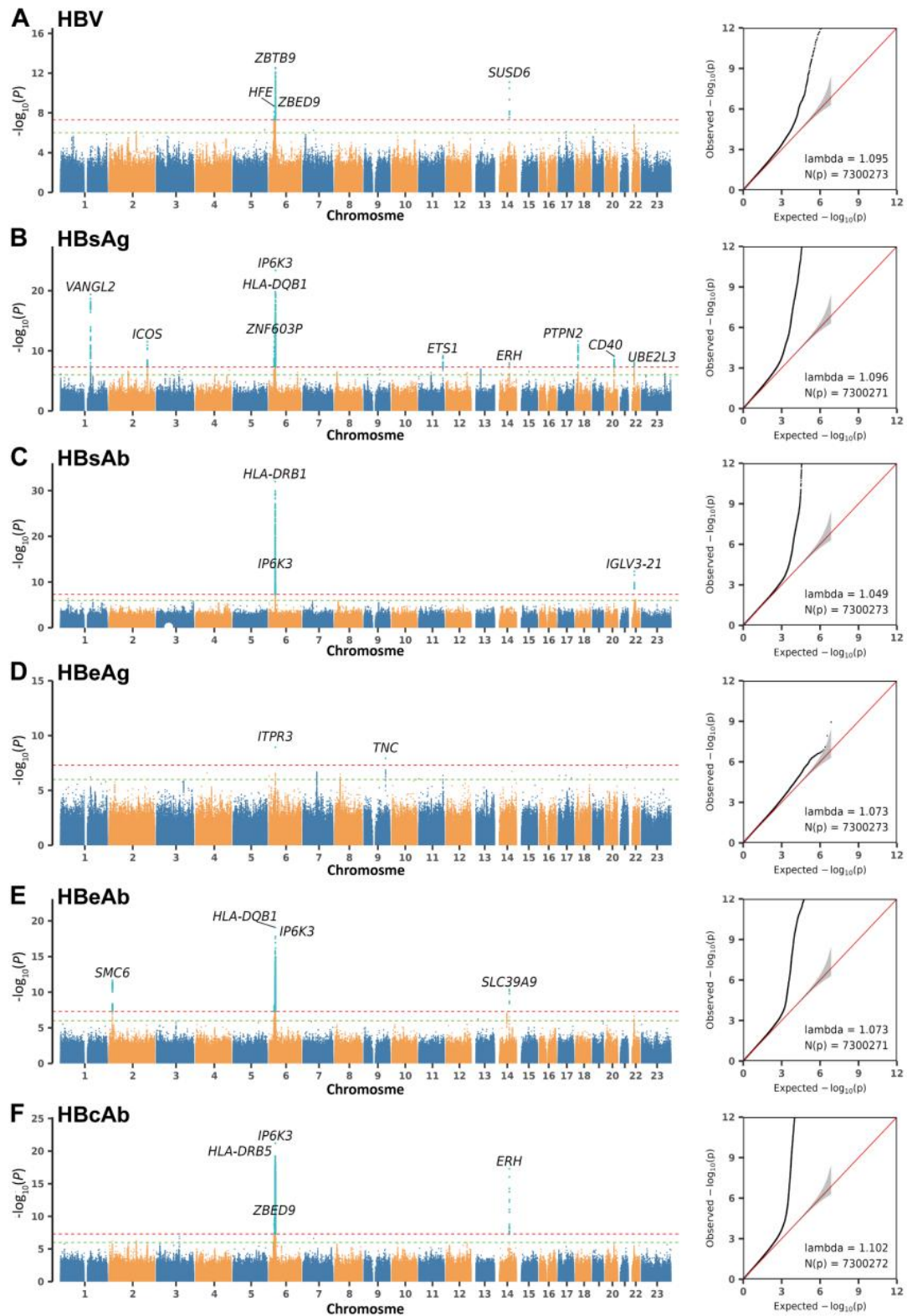


405

406 **Supplementary Fig. 8. Comparison of effect sizes between TBA and ICP GWAS**  
 407 **with European ICP GWAS results.**

408 Panel (A) and (B) compared the beta and *P* values of TBA and ICP with the  
 409 previously published meta-GWAS results of European population. The x-axis shows  
 410 the beta values of TBA or ICP meta-GWAS, while the y-axis illustrates the beta  
 411 values of the European ICP meta-GWAS result. The error bars represent the 95%  
 412 confidence interval of beta. Grey points highlight instances with different directions  
 413 of beta and/or *P* values > 0.05. Detail data, excluding empty and proxy loci, are  
 414 available in **Supplementary Table 2.**

415



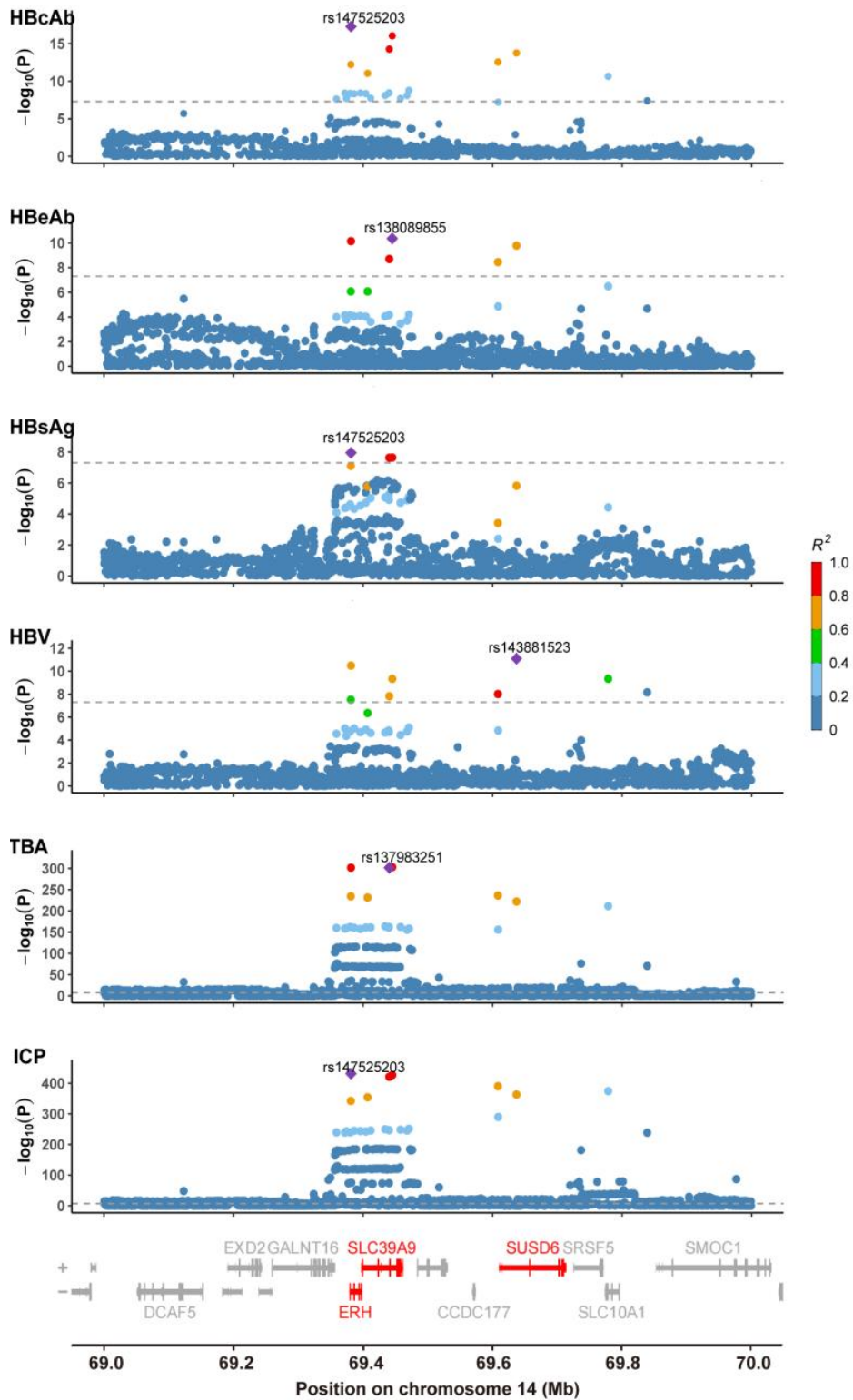
416

417 **Supplementary Fig. 9. Genome-wide association study of HBV and its antigens &**  
 418 **antibodies during pregnancy.**

419 Panel (A) to (F) present the GWAS results and QQ plots of Hepatitis B virus (HBV),

420 Hepatitis B surface antigen (HBsAg), Hepatitis B surface antibody (HBsAb),

421 Hepatitis B e antigen (HBeAg), Hepatitis B e antibody (HBeAb), and Hepatitis B core  
422 antibody (HBcAb) respectively. Horizontal lines delineate the genome-wide  
423 significance ( $P < 5 \times 10^{-8}$ , red line) and suggestive significance ( $P < 5 \times 10^{-6}$ , green line)  
424 thresholds. QQ plots elucidate the relationship between observed and expected  
425  $-\log_{10}(P)$  values from the GWAS meta-analysis of traits, which indicate the absence of  
426 significant population stratification. The sample size for HBV infection, as per  
427 medical records, is 44,432, while for HBsAg, HBsAb, HBeAg, HBeAb and HBcAb, it  
428 ranges from 94,441 to 94,462.  
429



430

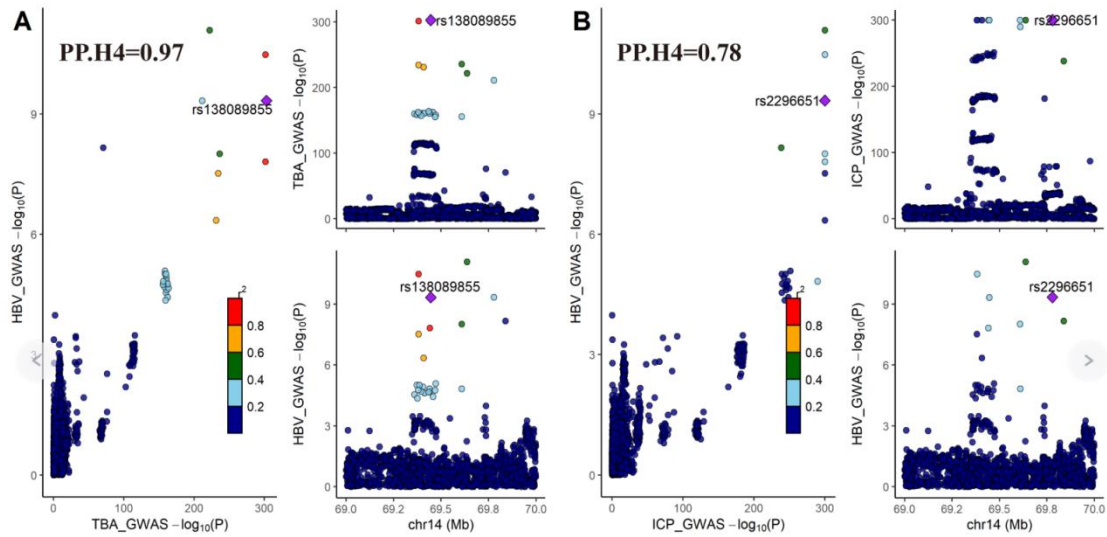
431 **Supplementary Fig. 10. Stacked LocusZoom plot of TBA, ICP and HBV related**  
 432 **traits.**

433 Shown is a Stacked LocusZoom plot depicting TBA, ICP and HBV related

434 phenotypes with genome-wide associations in the chromosome14 spanning 69.0-70.0

435 Mb region. The x-axis shows the chromosome position based on GRCh38, while the  
436 y-axis shows  $-\log_{10}(P)$  values for the associated tests. Genes linked to the lead SNPs  
437 are highlighted in red, and others are marked in grey. The SNP with lowest  $P$  value in  
438 the locus is indicated by a purple diamond. The remaining SNPs in the region are  
439 color-coded based on their  $R^2$  with the lead SNP.  $R^2$  were calculated using the BIGCS  
440 reference panel. Notably, all four lead SNPs exhibit high linkage disequilibrium.  
441





442

443 **Supplementary Fig. 11. LocusCompare plot of GWAS-GWAS colocalization**

444 **between TBA & ICP and HBV.**

445 The x-axis shows the chromosome position based on GRCh38 and the y-axis shows

446  $-\log_{10}(P)$  values for the associated tests. The purple diamond represents the shared

447 SNP of (A) TBA and HBV, (B) ICP and HBV. The remaining SNPs in the region are

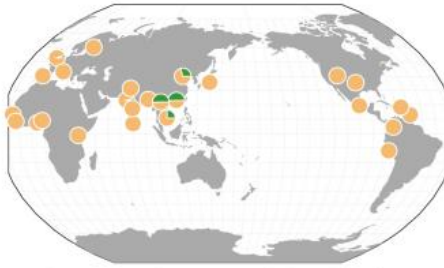
448 color-coded based on  $R^2$  with the lead SNP.  $R^2$  were calculated with 1000GP EAS

449 population as reference panel.

450

**A**

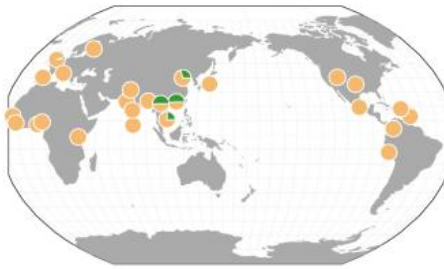
rs137983251  
chr14:69907145 G/A



Frequency Scale = Proportion out of 0.1  
The pie below represents a minor allele frequency of 0.025



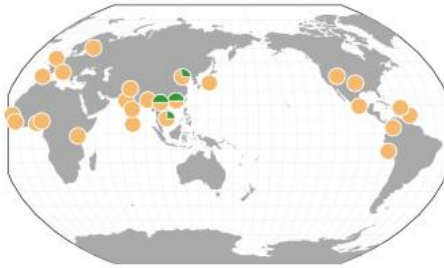
rs138089855  
chr14:69911799 C/T



Frequency Scale = Proportion out of 0.1  
The pie below represents a minor allele frequency of 0.025



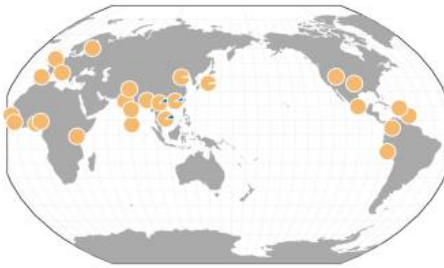
rs147525203  
chr14:69848124 C/T



Frequency Scale = Proportion out of 0.1  
The pie below represents a minor allele frequency of 0.025



rs2296651  
chr14:70245193 A/G



Frequency Scale = Proportion out of 1  
The pie below represents a minor allele frequency of 0.25



**B**

rs147525203  
chr14:69848124

		C	T	
rs137983251 chr14:69907145	A	0	979	979 (0.971)
	G	29	0	29 (0.029)
		29	979	1008
		(0.029)	(0.971)	

Haplotypes	Statistics
A_T: 979 (0.971)	D': 1.0
G_C: 29 (0.029)	R <sup>2</sup> : 1.0
A_C: 0 (0.0)	Chi-sq: 1008.0
G_T: 0 (0.0)	p-value: <0.0001

rs137983251(A) allele is correlated with rs147525203(T) allele  
rs137983251(G) allele is correlated with rs147525203(C) allele

rs138089855  
chr14:69445082

		C	T	
rs137983251 chr14:69440428	A	0	979	979 (0.971)
	G	29	0	29 (0.029)
		29	979	1008
		(0.029)	(0.971)	

Haplotypes	Statistics
A_T: 979 (0.971)	D': 1.0
G_C: 29 (0.029)	R <sup>2</sup> : 1.0
A_C: 0 (0.0)	Chi-sq: 1008.0
G_T: 0 (0.0)	p-value: <0.0001

rs137983251(A) allele is correlated with rs138089855(T) allele  
rs137983251(G) allele is correlated with rs138089855(C) allele

rs2296651  
chr14:70245193

		A	G	
rs147525203 chr14:69848124	C	22	7	29 (0.029)
	T	50	929	979 (0.971)
		72	936	1008
		(0.071)	(0.929)	

Haplotypes	Statistics
T_G: 929 (0.922)	D': 0.7401
T_A: 50 (0.05)	R <sup>2</sup> : 0.2109
C_A: 22 (0.022)	Chi-sq: 212.591
C_G: 7 (0.007)	p-value: <0.0001

rs147525203(C) allele is correlated with rs2296651(A) allele  
rs147525203(T) allele is correlated with rs2296651(G) allele

rs2296651  
chr14:69778476

		A	G	
rs138089855 chr14:69445082	C	22	7	29 (0.029)
	T	50	929	979 (0.971)
		72	936	1008
		(0.071)	(0.929)	

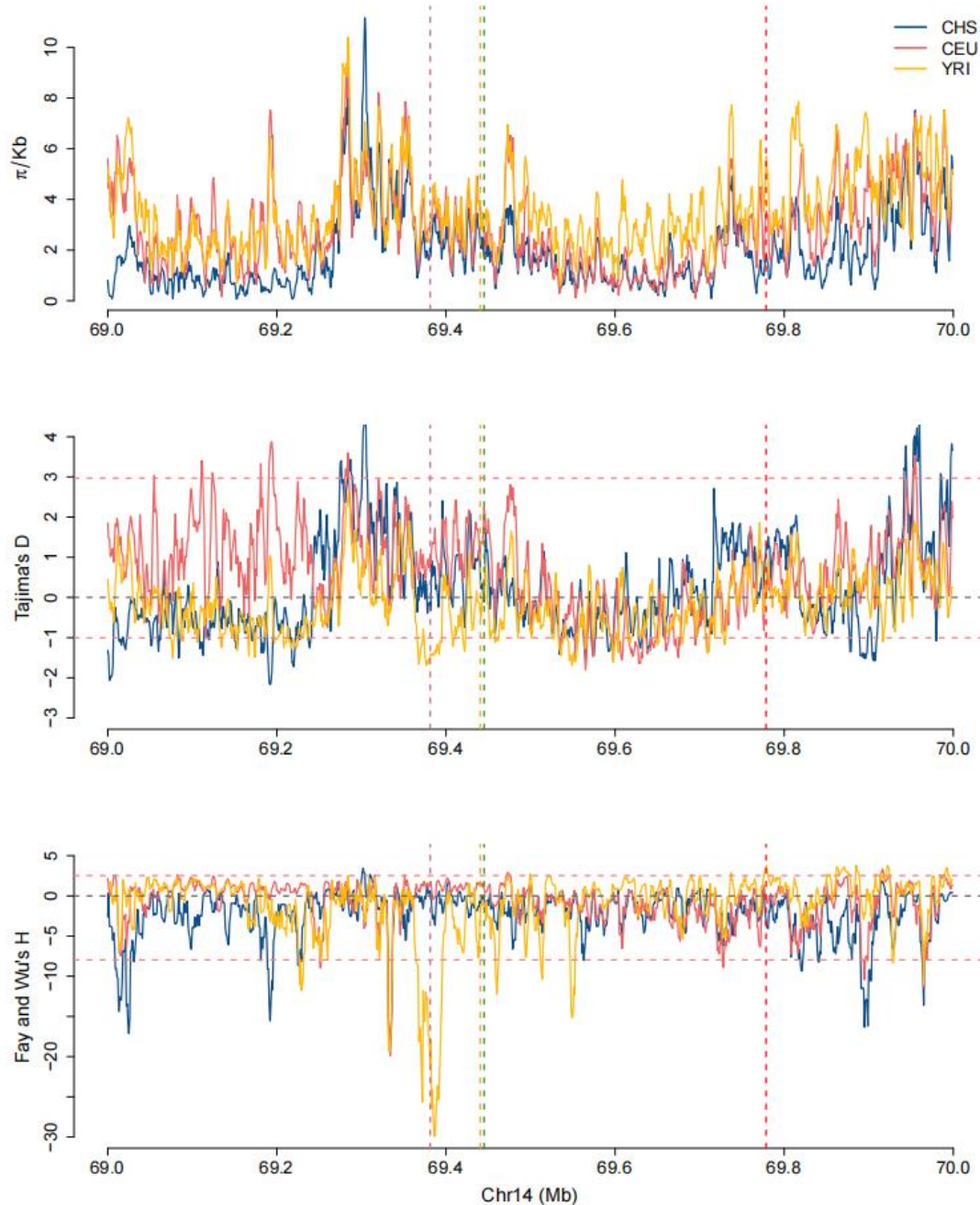
Haplotypes	Statistics
T_G: 929 (0.922)	D': 0.7401
T_A: 50 (0.05)	R <sup>2</sup> : 0.2109
C_A: 22 (0.022)	Chi-sq: 212.591
C_G: 7 (0.007)	p-value: <0.0001

rs138089855(C) allele is correlated with rs2296651(A) allele  
rs138089855(T) allele is correlated with rs2296651(G) allele

452 **Supplementary Fig. 12. Geographical allele frequency distribution and linkage**  
453 **disequilibrium (LD) of rs137983251-G, rs138089855-C, rs47525203 and**  
454 **rs2296651-A.**

455 (A) The variations in all four loci are exclusive to East Asia, displaying a distinct  
456 pattern of higher frequencies in the southern region and lower frequencies in the  
457 northern region. Geography plots were generated using the Geography of Genetic  
458 Variants Browser (<https://popgen.uchicago.edu/ggv/>). (B) All three loci exhibit high  
459 linkage disequilibrium (LD), with significant *P* values less than 0.0001. Plot (B) was  
460 extracted from LDpair (<https://ldlink.nih.gov/?tab=ldpair>).

461

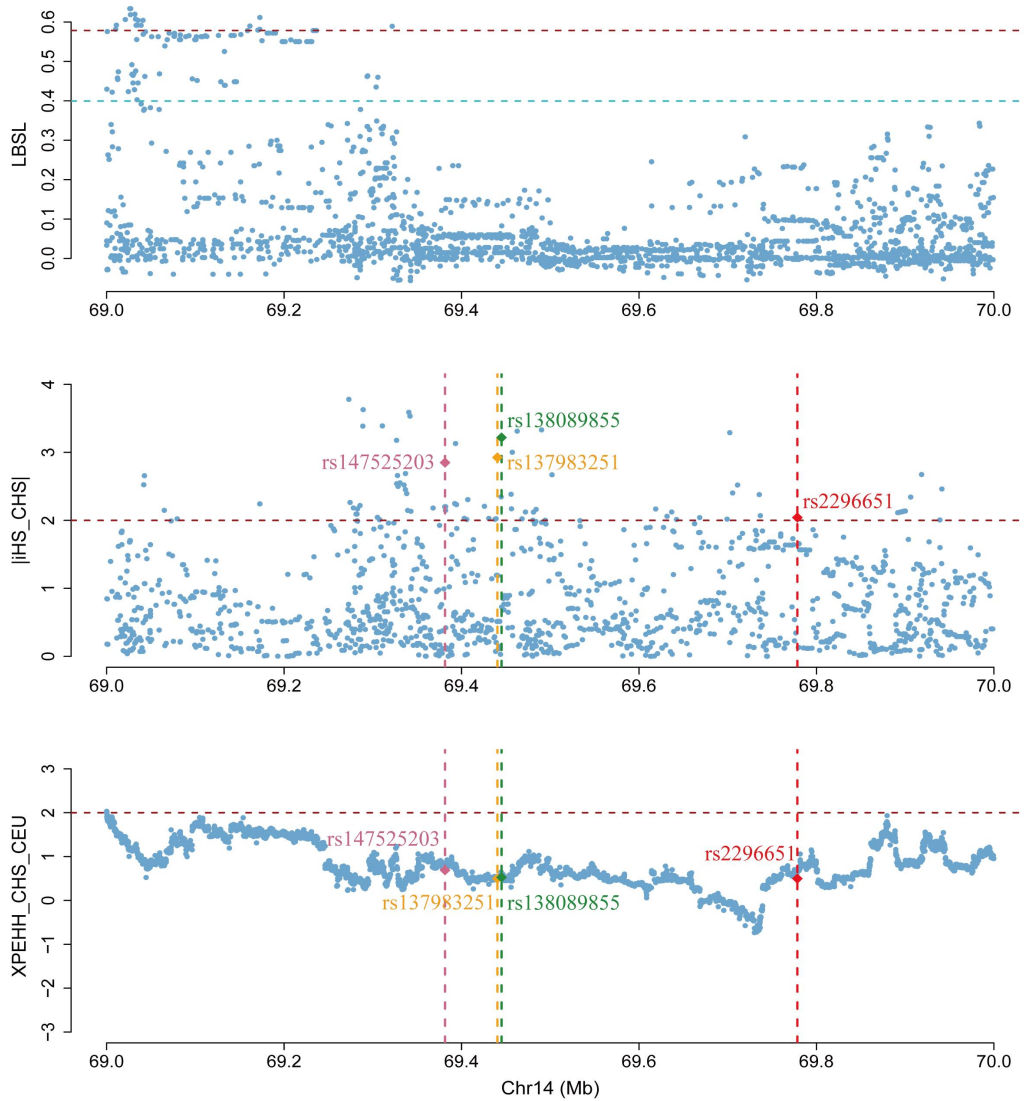


462

463 **Supplementary Fig. 13. Site Frequency Spectrum-based selection tests with**  
 464 **modern DNA samples from 1000 Genomes Project across three populations.**

465 Vertical dashed lines of distinct colors highlight the positions of four target Single  
 466 Nucleotide Polymorphisms (SNPs): rs147525203 (palevioletred), rs137983251  
 467 (orange), rs138089855 (green) and rs2296651 (red). In the bottom two panels, a pink  
 468 dashed line signifies the lowest 5% and 95% cutoffs among the three populations.  
 469 Specifically, genomic regions falling below this pink dashed line achieve statistical  
 470 significance across all populations. It is important to note that some genomic regions

471 above the pink dashed line may still be significant in certain populations, as the 5%  
472 cutoffs in these populations are higher (less extreme). The populations represented are  
473 CHS (Southern Han Chinese), CEU (Utah Residents with Northern and Western  
474 European Ancestry), and YRI (Yoruba in Ibadan, Nigeria).  
475



476

477 **Supplementary Fig. 14. The plot of LSBL, iHS and XP-EHH for the chromosome**  
 478 **14 region based on data from the 1000GP.**

479 The three panels represent the result of LSBL, iHS and XP-EHH, respectively.

480 LSBL was computed for CHS with CEU using YRI as reference population. XP-EHH

481 was calculated for CHS with CEU as the reference population. Each point on the plot

482 represents a genetic variant. In the LSBL panel, the red and blue horizontal dashed

483 lines denote the 0.1% and 1% threshold of the empirical distribution of the whole

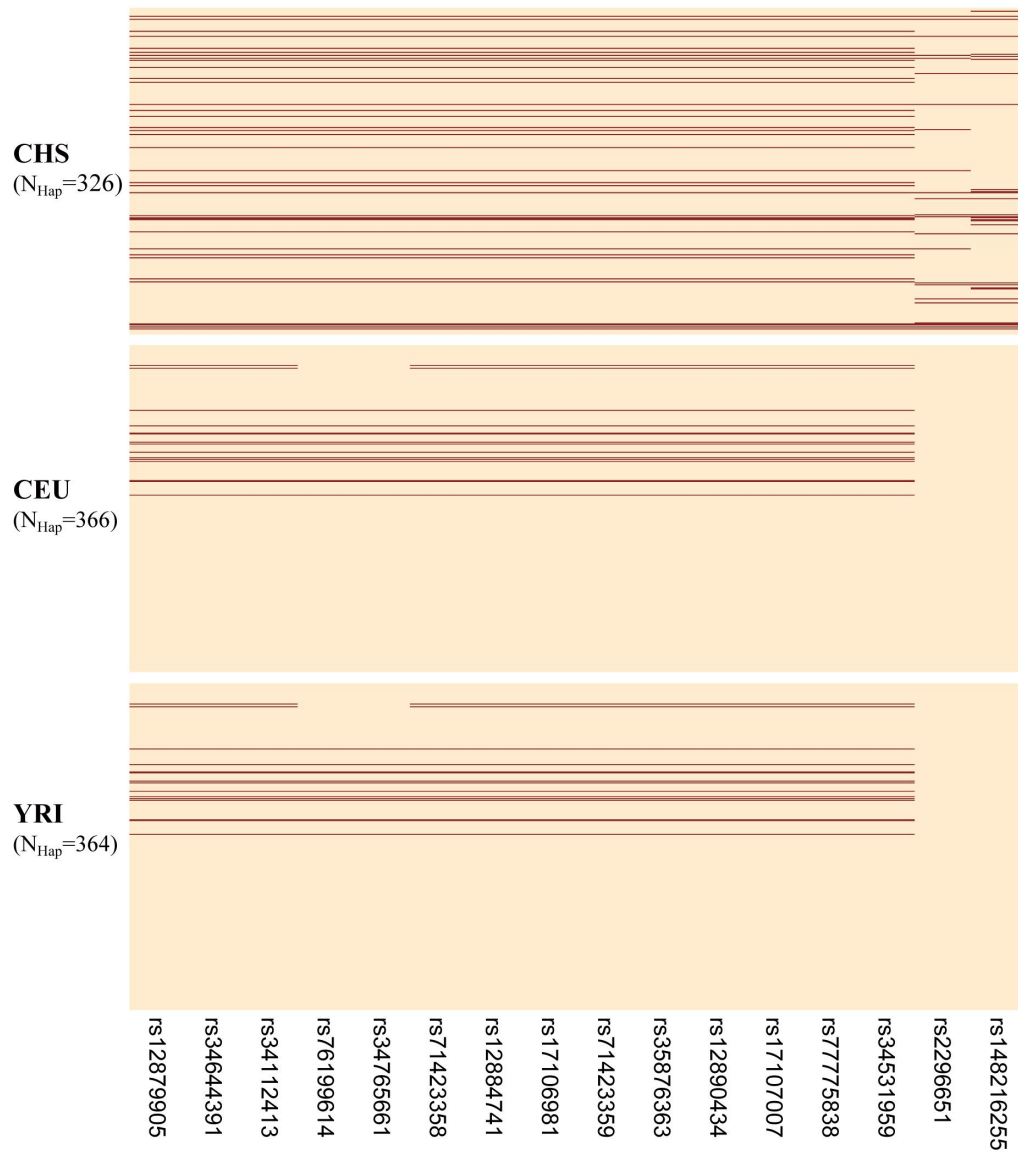
484 genome. In the iHS and XPEHH panels, the red horizontal dashed line signifies the

485 threshold of significance. Positions for the four target SNPs are indicated with vertical

486 dashed lines of different colors: rs147525203 (pale violet), rs137983251 (orange),

487 rs138089855 (green), rs2296651 (red). CHS: Southern Han Chinese; CEU: Utah

488 Residents (CEPH) with Northern and Western European Ancestry; YRI: Yoruba in  
489 Ibadan, Nigeria. Note that genomic portion from chr14:69.4-69.8M is missing in  
490 LSBL results due to high frequency of zero loci in the CEU population in this region.  
491



492

493 **Supplementary Fig. 15. Haplotype structure of the ICP risk 14q24.1 locus in**

494 **CHS, CEU and YRI populations.**

495 Each horizontal line in the plot represents a distinct haplotype, and the composition of

496 each haplotype, consisting of 16 SNPs, is presented below. The organization is

497 organized by population (CHS, CEU, and YRI). In the plot, the blancheted color

498 is used to signify the allele in its ancestral state, while the brown color indicates a

499 derived state.

500



501 **Supplementary Table Legends**

502 **Supplementary Table 1.** Baseline characteristics and TBA concentrations of  
503 participants.

504 **Supplementary Table 2.** Genome-wide significant signals for TBA & ICP, statistics  
505 from internal and external replication and comparison with the European population

506 **Supplementary Table 3.** Genome-wide significant loci of TBA & ICP and internal  
507 replication.

508 **Supplementary Table 4.** Genome-wide significant loci of TBA & ICP and external  
509 replication in two independent cohorts in China.

510 **Supplementary Table 5.** Pathway and gene ontology enrichment analyses for the  
511 TBA associated gene loci (A) and the ICP associated gene loci (B)

512 **Supplementary Table 6.** Meta-GWAS discoveries compared with European effect  
513 estimates

514 **Supplementary Table 7.** GWAS-GWAS colocalization of TBA & ICP with HBV.

515 **Supplementary Table 8.** Frequencies of rs2296651 in the Holocene (>10,000BP) age,  
516 visualized in Fig. 3A.

517 **Supplementary Table 9.** Frequencies of rs2296651 in the Neolithic  
518 (10,000~3,000BP) age, visualized in Fig. 3B.

519 **Supplementary Table 10.** Frequencies of rs2296651 in the Historic (<3,000) age,  
520 visualized in Fig. 3C.

521 **Supplementary Table 11.** Frequencies of rs2296651 in the Present-day populations,  
522 visualized in Fig. 3D.

523 **Supplementary Table 12.** Haplotype of the lead SNPs with high Linkage  
524 disequilibrium (LD) on chromosome 14 in the 1000GP.

525