

Supplementary information for

Whole genome sequencing and analysis of 4,053 individuals in trios and mother-infant duos from the Born in Guangzhou Cohort Study

Shujia Huang^{1*}, Mingxi Huang^{1*}, Siyang Liu^{2*}, Chengrui Wang¹, Jianrong He^{1,3,4}, Yashu Kuang¹, Jinhua Lu¹, Yuqin Gu², Xiaoyan Xia¹, Shanshan Lin¹, Huimin Xia^{3,5,6#}, Xiu Qiu^{1,3,7#}

Affiliations:

¹ Division of Birth Cohort Study, Guangzhou Women and Children's Medical Center, Guangzhou Medical University, Guangzhou, 510623, China.

² School of Public Health (Shenzhen), Sun Yat-sen University, Shenzhen, Guangdong 510006, China.

³ Provincial Clinical Research Center for Child Health, Guangzhou, 510623, China.

⁴ Department of Obstetrics and Gynecology, Guangzhou Women and Children's Medical Center, Guangzhou Medical University, Guangzhou, 510623, China

⁵ Provincial Key Laboratory of Research in Structure Birth Defect Disease, Guangzhou Women and Children's Medical Center, Guangzhou Medical University, Guangzhou, 510623, China.

⁶ Department of Pediatric Surgery, Guangzhou Women and Children's Medical Center, Guangzhou Medical University, Guangzhou, 510623, China

⁷ Department of Women's Health, Provincial Key Clinical Specialty of Woman and Child Health, Guangzhou Women and Children's Medical Center, Guangzhou Medical University, Guangzhou, 510623, China

* These authors contributed equally to this work: Shujia Huang, Mingxi Huang, Siyang Liu

These authors supervised this work: Xiu Qiu and Huimin Xia. Correspondence should be addressed to xiu.qiu@bigcs.org and huimin.xia@bigcs.org

Table of Contents

<i>Supplementary Figures</i>.....	3
Figure S1. Ethnic distribution of the 4,053 samples of BIGCS.....	3
Figure S2. Quality and overview of the sequencing samples.....	4
Figure S3. The length distribution of the Indels.	5
Figure S4. The number of the variants along geographical regions and ethnicities.....	6
Figure S5. The genotype validation.....	7
Figure S6. PCA analysis for the BIGCS and the 1KGP3 samples.	8
Figure S7. PCA of the BIGCS samples from seven geographical groups of China.	9
Figure S8. The distribution of cross-validation error of ADMIXTURE analysis.....	10
Figure S9. The LocusZoom plots of 15 lead variants for the ten traits with GWAS significant loci revealed in the study.....	13
Figure S10. RNA tissue specificity expression of gene <i>TTC28</i> , <i>SLC10A1</i> and <i>SOAT2</i>	14

Supplementary Figures

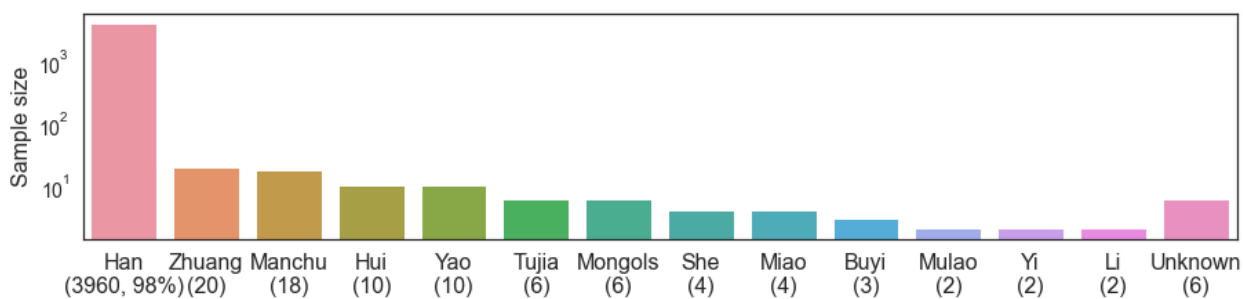


Figure S1. Ethnic distribution of the 4,053 samples of BIGCS. The y-axis represents the sample size in log-scale and each color bar represents each of the ethnic groups. The number of individuals is indicated in the parenthesis. Related to Figure 1.

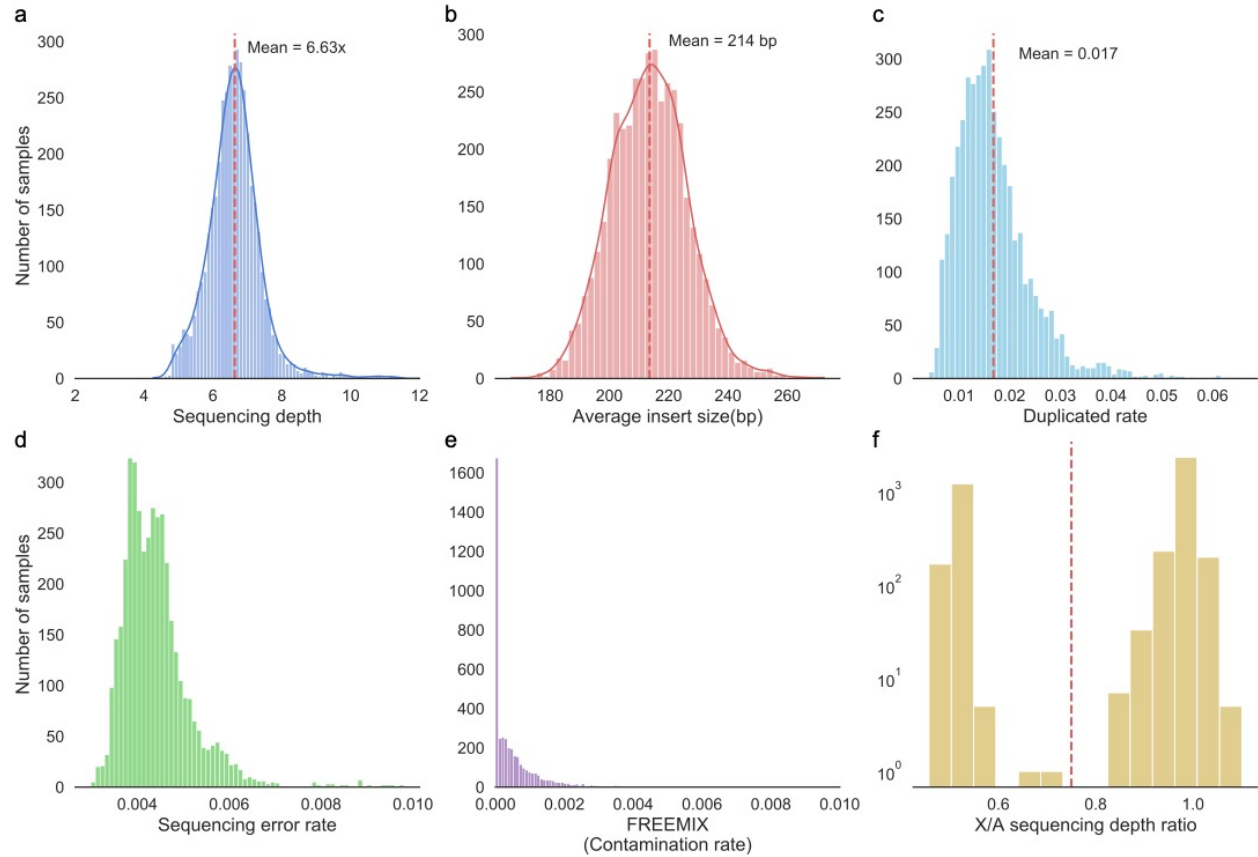


Figure S2. Quality and overview of the sequencing samples.

(a) Distribution of sequencing depth in the WGS data of BIGCS. (b) Distribution of average insert size. (c) Distribution of the duplication rate of the sequencing data. (d) Distribution of sequencing error rate provided by BGI-Shenzhen Co., Ltd. We used 0.01 as the strict cutoff to filter high sequencing error samples. (e) Distribution of contamination rates estimated by verifyBamID2. A strict cutoff of 0.01 was used to filter contaminated samples. (f) Distribution of average sequencing depth between chromosome X and autosomes (X/A ratio). A cutoff of 0.75 was used to infer the gender of the BIGCS samples.

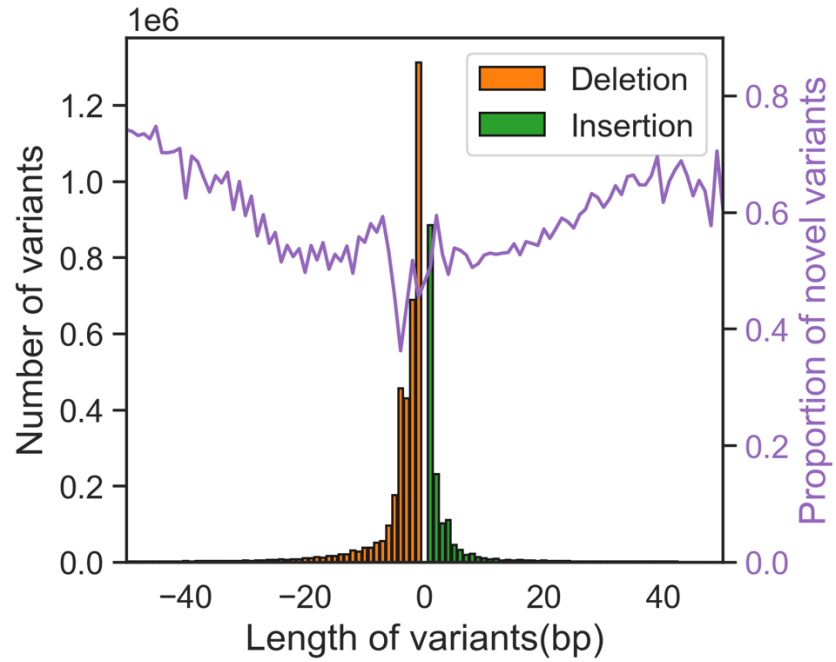


Figure S3. The length distribution of the Indels.

The x-axis represents the length of Indels variants from -50 bp (Deletion) to +50 bp (Insertion). The left y-axis is the count of Indels. The purple line and right y-axis represent the proportion of novel variants. Related to Figure 1.

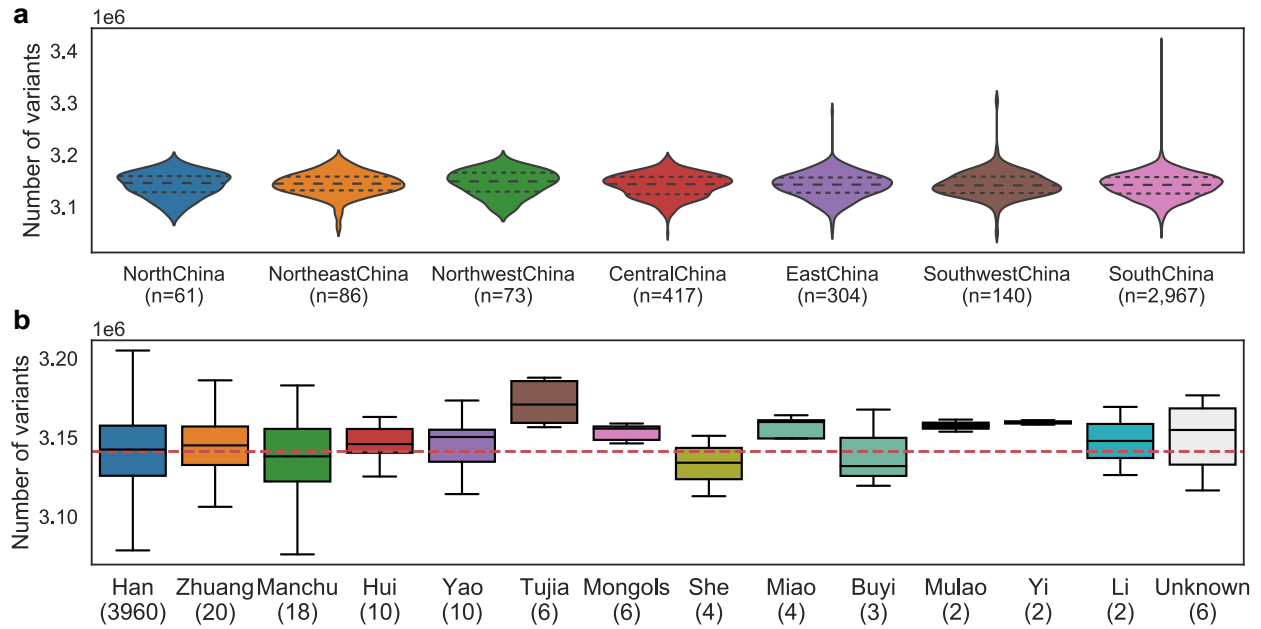


Figure S4. The number of the variants along geographical regions and ethnicities.

The distribution of number of variants (SNPs and Indels) grouped by the seven large geographical regions of China (a) and 13 ethnic groups in BIGCS (b) with sample size indicated in parenthesis. The red dash horizontal line demonstrates the mean of individuals' variants.

Related to Figure 1.

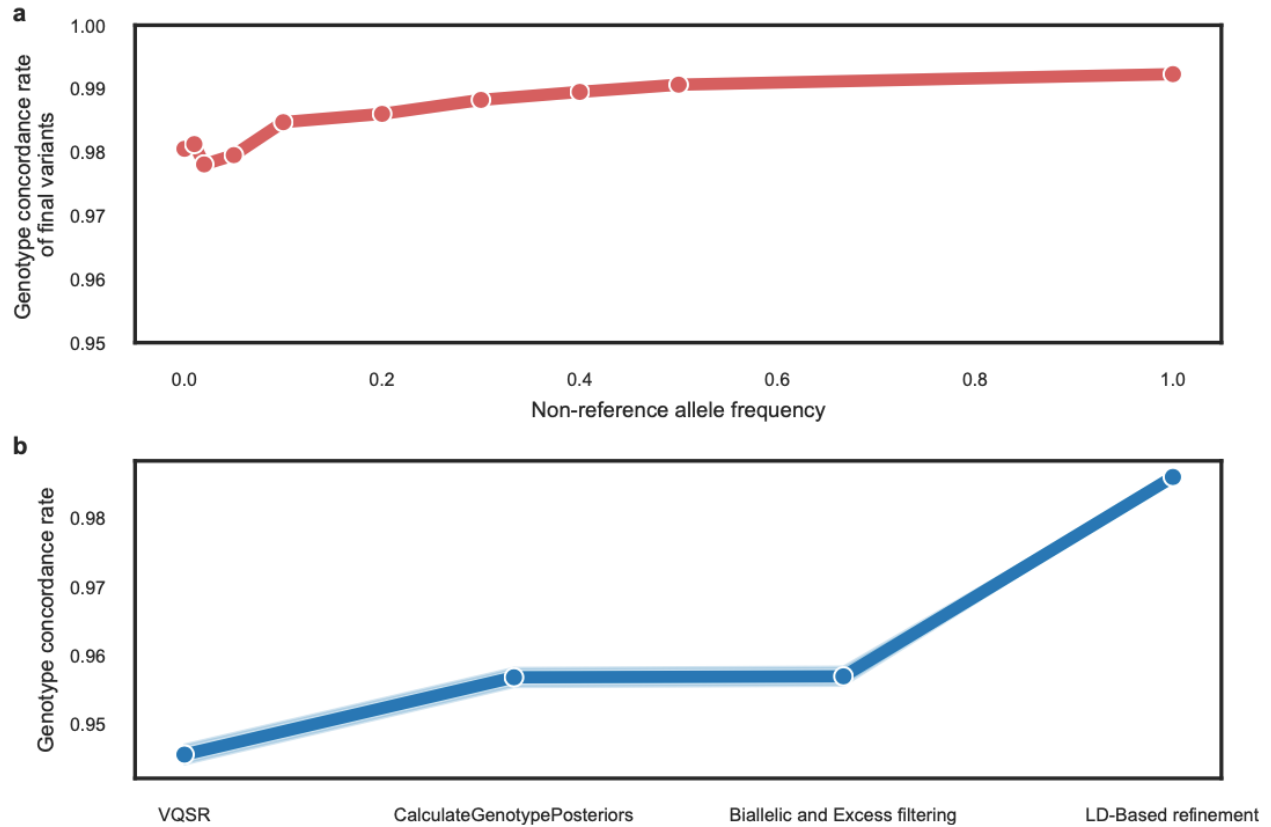


Figure S5. The genotype validation.

Genotype concordance rate of the final variants call set (after LD-based refinement) compare with the genotype data of 240 SNP-array samples in different non-reference allele frequency bins (a) and genotype concordance rate in different variants filtration processes (b). Related to Figure 1.

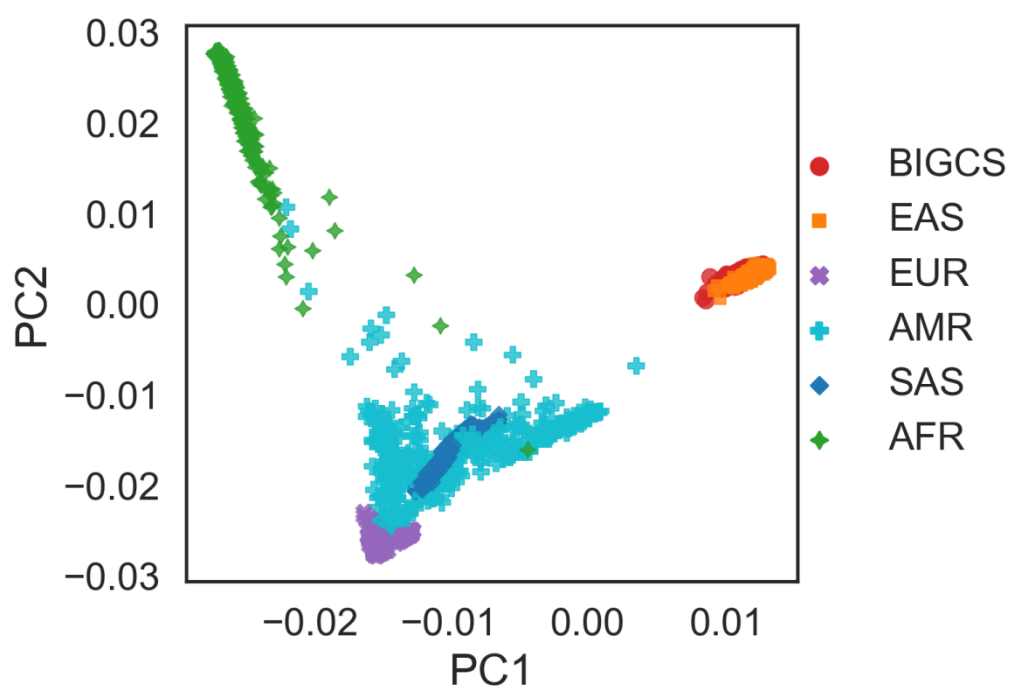


Figure S6. PCA analysis for the BIGCS and the 1KGP3 samples.

Each point represents one participant and is placed according to their PC scores. Related to Figure 2.

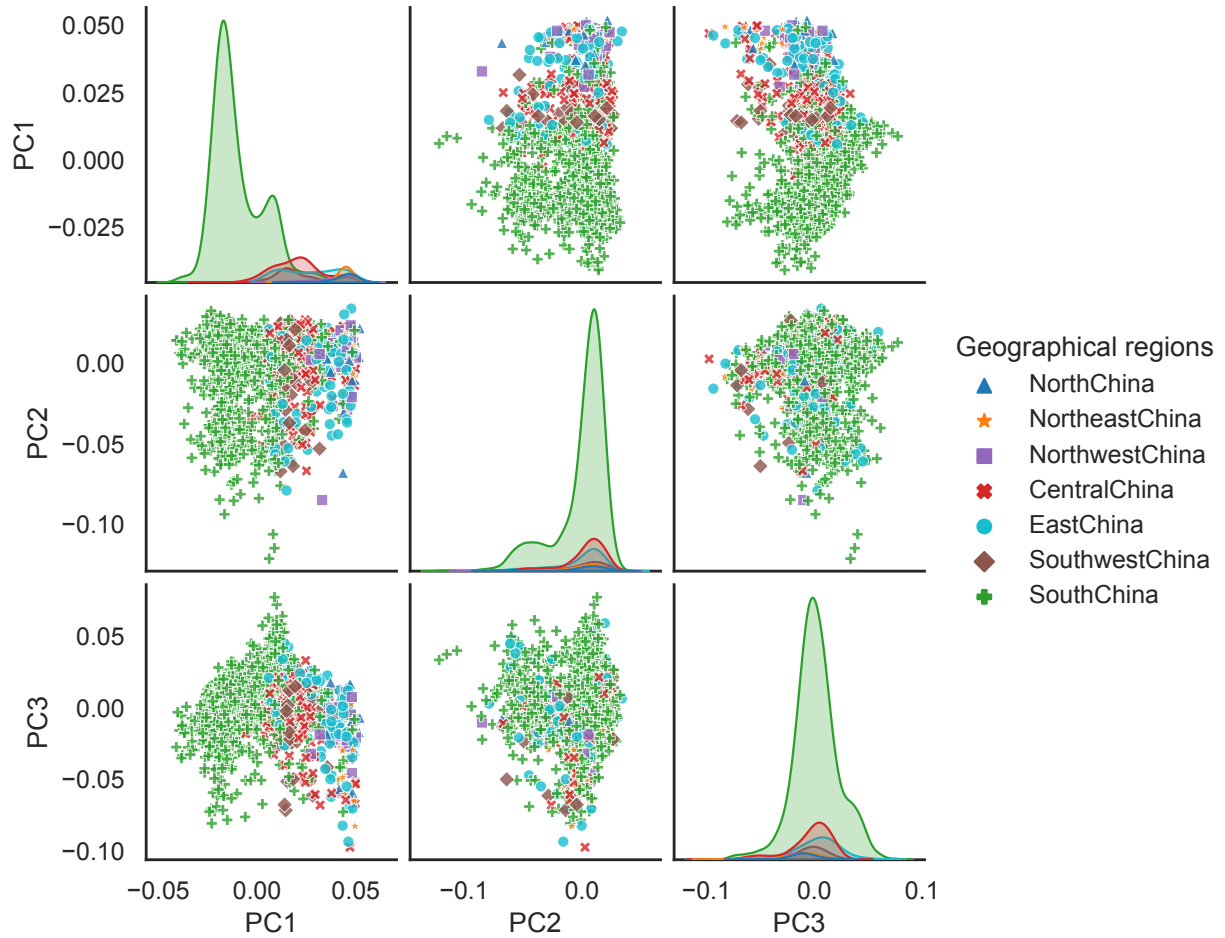


Figure S7. PCA of the BIGCS samples from seven geographical groups of China. Pair plot for the first three components of PCA analysis. Each point represents one participant and is placed according to their PC scores. Related to Figure 2.

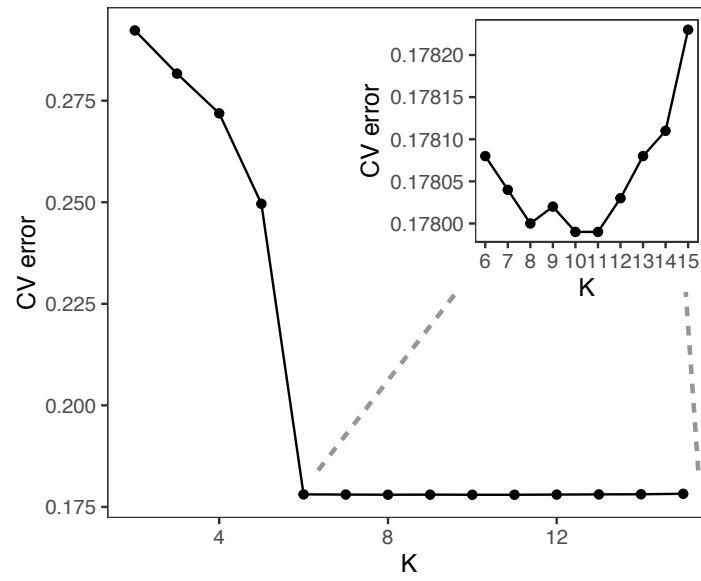
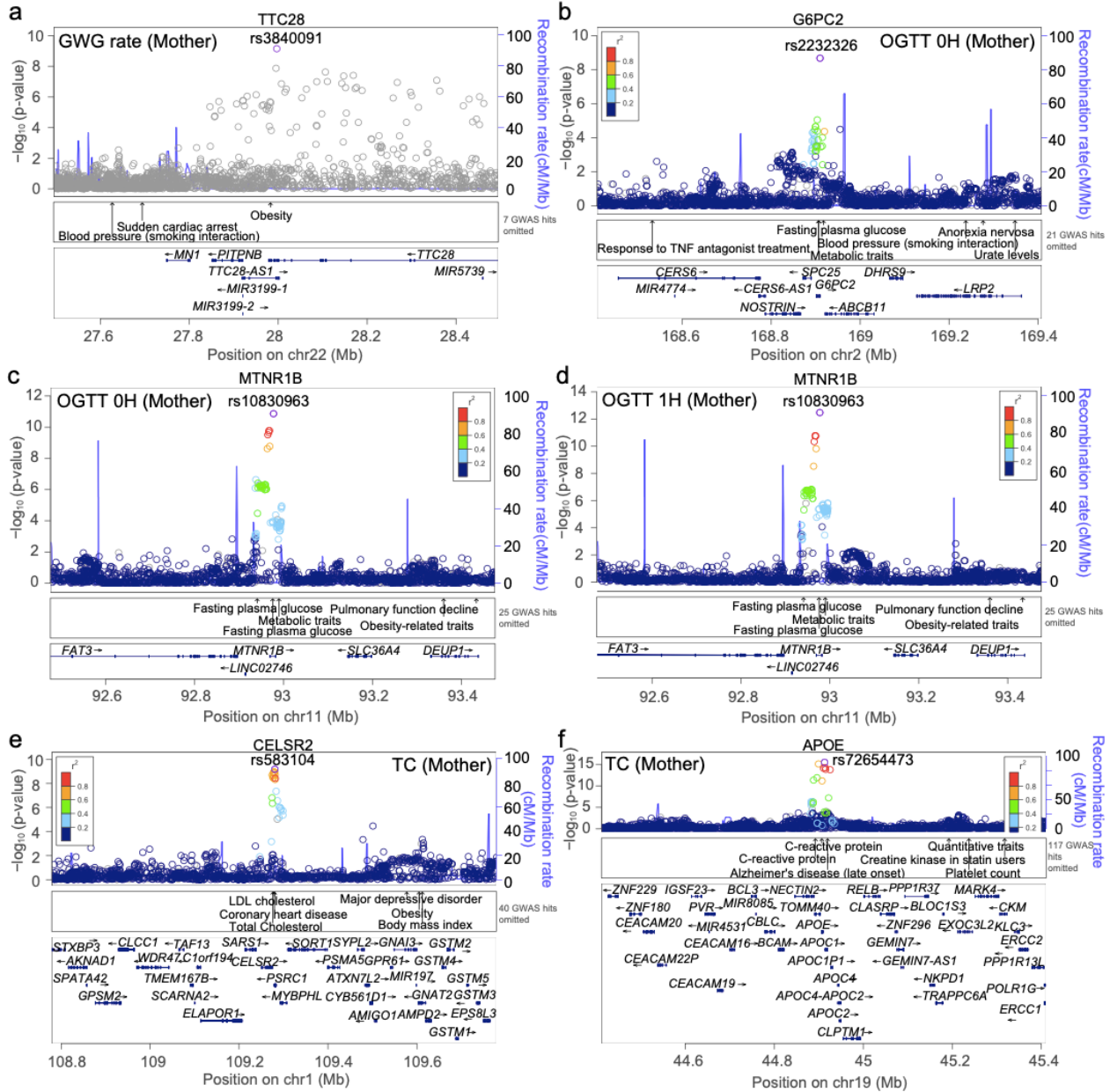
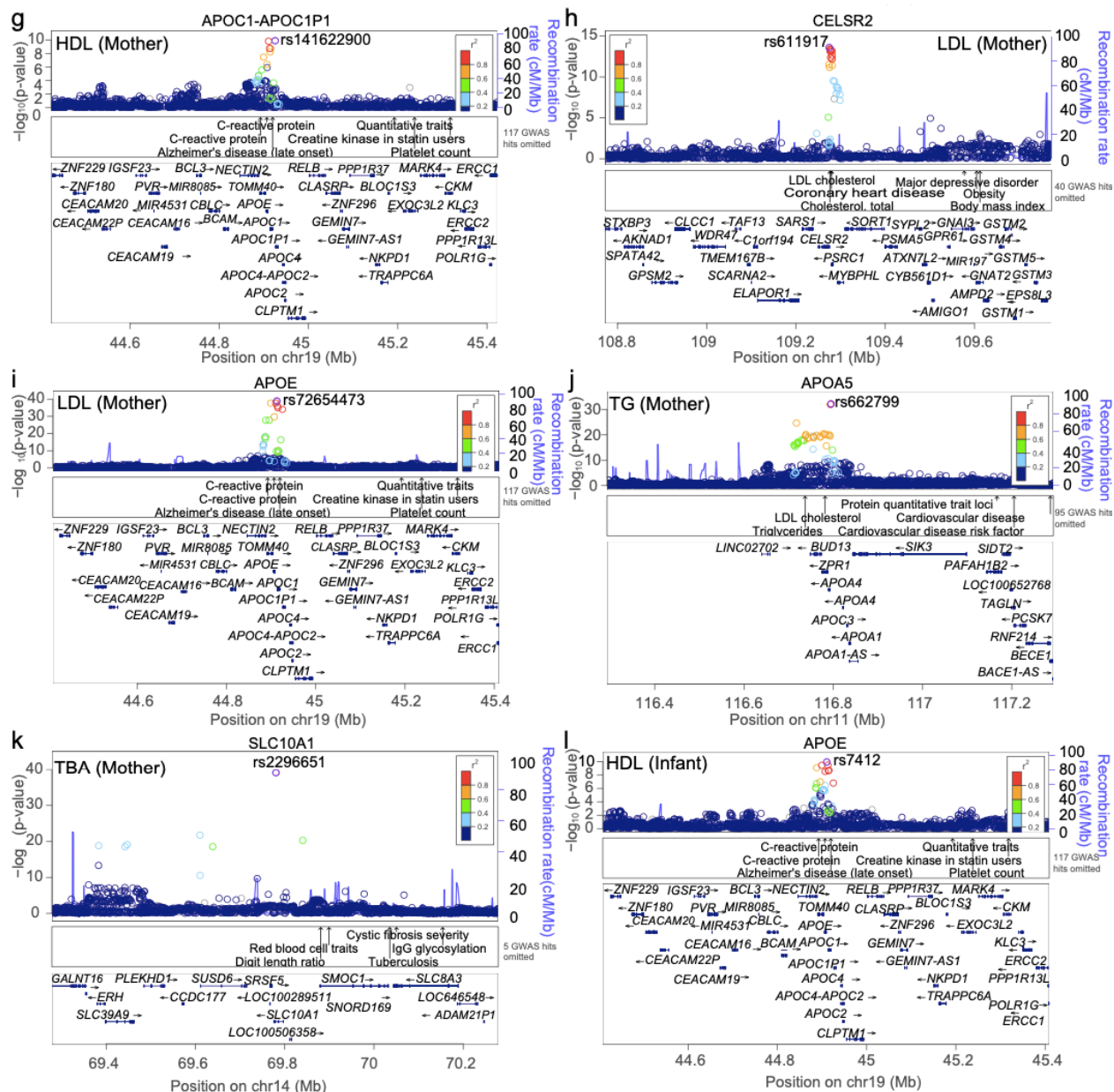


Figure S8. The distribution of cross-validation error of ADMIXTURE analysis. Related to Figure 2.





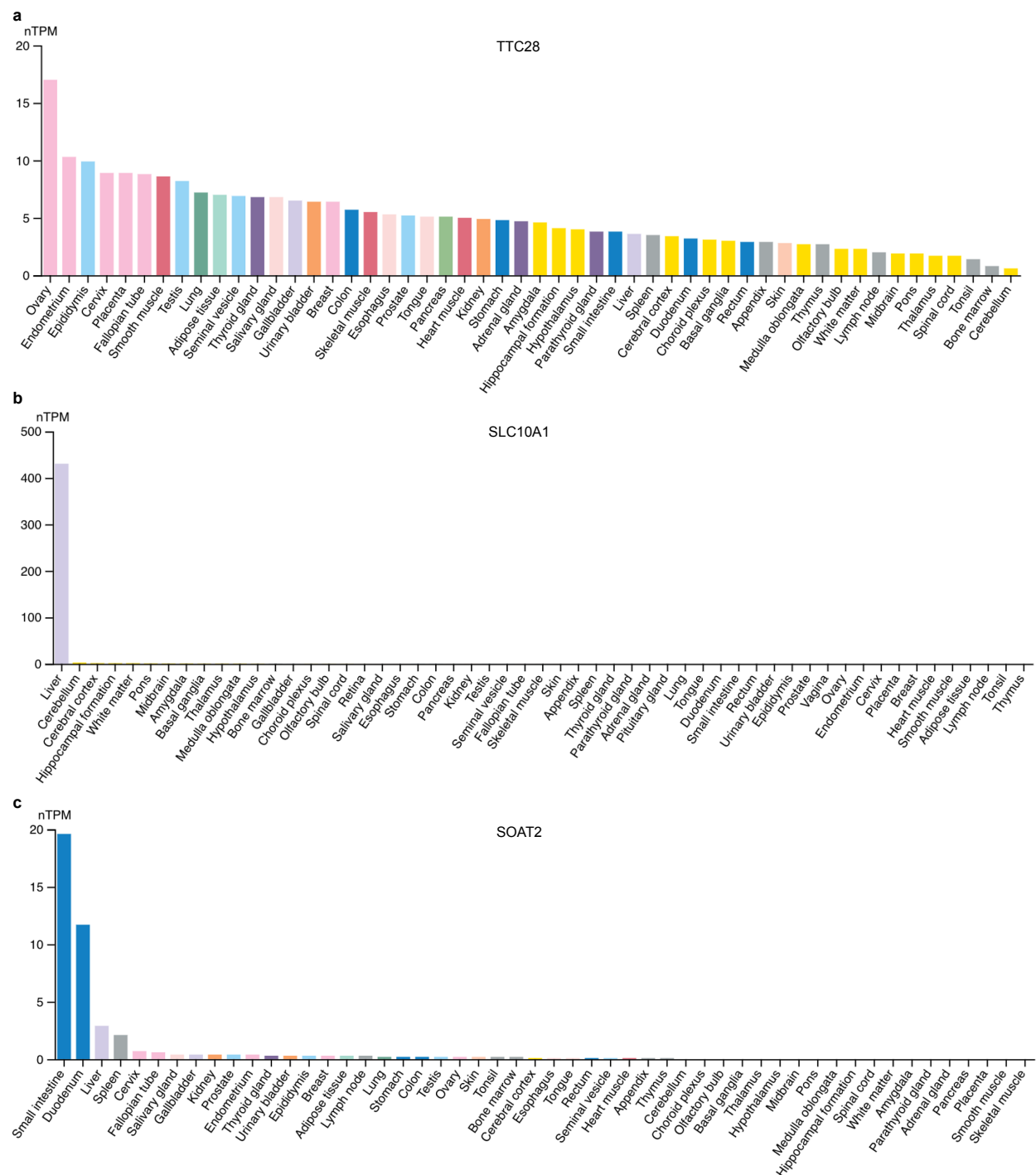


Figure S10. RNA tissue specificity expression of gene *TTC28*, *SLC10A1* and *SOAT2*.
 RNA tissue specificity expression of gene *TTC28* (a), *SLC10A1* (b) and *SOAT2* (c) according to the Human Protein Atlas dataset (<https://www.proteinatlas.org/>). Related to Figure 3.