Cell

Genomic Analyses from Non-invasive Prenatal Testing Reveal Genetic Associations, Patterns of Viral Infections, and Chinese Population History

Graphical Abstract



Highlights

- Genome sequencing from low-pass noninvasive prenatal testing samples
- GWAS of 141,431 low-pass genomes reveals 16 unknown genetic associations
- Patterns of clinically relevant viral infection in maternal plasma
- Insights into the genetic structure and history of the Chinese
 population

Siyang Liu, Shujia Huang, Fang Chen, ..., Xin Jin, Rasmus Nielsen, Xun Xu

Authors

Correspondence

wangjian@genomics.cn (J.W.), albrecht@binf.ku.dk (A.A.), jinxin@genomics.cn (X.J.), rasmus_nielsen@berkeley.edu (R.N.), xuxun@genomics.cn (X.X.)

In Brief

Large-scale analysis of genome sequences from non-invasive prenatal testing in Chinese women yields insights into phenotypic trait associations, viral infection patterns, and population history.





Genomic Analyses from Non-invasive Prenatal Testing Reveal Genetic Associations, Patterns of Viral Infections, and Chinese Population History

Siyang Liu,^{1,2,25} Shujia Huang,^{1,3,25} Fang Chen,^{1,23,24,25} Lijian Zhao,^{1,25} Yuying Yuan,^{1,25} Stephen Starko Francis,^{4,5} Lin Fang,¹ Zilong Li,⁶ Long Lin,⁶ Rong Liu,¹ Yong Zhang,¹ Huixin Xu,¹ Shengkang Li,¹ Yuwen Zhou,^{1,6} Robert W. Davies,⁷ Qiang Liu,¹ Robin G. Walters,⁸ Kuang Lin,⁸ Jia Ju,¹ Thorfinn Korneliussen,⁹ Melinda A. Yang,¹⁰ Qiaomei Fu,^{10,11,12} Jun Wang,¹ Lijun Zhou,¹ Anders Krogh,² Hongyun Zhang,¹ Wei Wang,¹ Zhengming Chen,⁸ Zhiming Cai,^{13,14,15} Ye Yin,¹ Huanming Yang,^{1,16} Mao Mao,¹ Jay Shendure,^{17,18,19} Jian Wang,^{1,16,*} Anders Albrechtsen,^{2,*} Xin Jin,^{1,3,*} Rasmus Nielsen,^{20,21,22,*} and Xun Xu^{1,26,*}

¹BGI-Shenzhen, Shenzhen 518083, Guangdong, China

²Bioinformatics Centre, Department of Biology, University of Copenhagen, Copenhagen 2200, Denmark

³School of Medicine, South China University of Technology, Guangzhou 510006, Guangdong, China

⁴Division of Epidemiology, University of Nevada, Reno, NV 89557, USA

⁵Department of Epidemiology and Biostatistics, University of California, San Francisco, San Francisco, CA 94143, USA

⁶BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083, Guangdong, China

⁷The Centre for Applied Genomics, Genetics and Genome Biology, The Hospital for Sick Children, Toronto, ON M6H 3W4, Canada ⁸Clinical Trial Service Unit and Epidemiological Studies Unit, Nuffield Department of Population Health, University of Oxford, Oxford OX3 7LF, UK

⁹Department of Zoology, University of Cambridge, Downing Street, Cambridge CB2 3EJ, UK

¹⁰Key Laboratory of Vertebrate Evolution and Human Origins of Chinese Academy of Sciences, Institute of Vertebrate Paleontology and Paleoanthropology, Chinese Academy of Sciences, Beijing 100044, China

¹¹Center for Excellence in Life and Paleoenvironment, Chinese Academy of Sciences, Beijing 100044, China

¹²University of Chinese Academy of Sciences, Beijing 100049, China

¹³Department of Urological Surgery, The First Affiliated Hospital of Shenzhen University (Shenzhen Second People's Hospital), Shenzhen 518035, Guangdong, China

¹⁴Guangdong Key Laboratory of Systems Biology and Synthetic Biology for Urogenital Tumors, Shenzhen 518035, Guangdong, China

¹⁵Shenzhen University Carson International Cancer Center, Shenzhen 518060, Guangdong, China

¹⁶James D. Watson Institute of Genome Sciences, Hangzhou 310058, Zhejiang, China

¹⁷Department of Genome Sciences, University of Washington, Seattle, WA 98195, USA

¹⁸Howard Hughes Medical Institute, Seattle, WA 98195, USA

¹⁹Brotman Baty Institute for Precision Medicine, Seattle, WA 98195, USA

²⁰Department of Integrative Biology, University of California, Berkeley, Berkeley, CA 94720, USA

²¹Centre for GeoGenetics, Natural History Museum of Denmark, University of Copenhagen, Copenhagen 1350, Denmark

²²Department of Statistics, University of California, Berkeley, Berkeley, CA 942720, USA

²³Laboratory of Genomics and Molecular Biomedicine, Department of Biology, University of Copenhagen, Copenhagen 2100, Denmark
²⁴BGI-Changyuan, Xinxiang 453400, Henan, China

²⁵These authors contributed equally

²⁶Lead Contact

*Correspondence: wangjian@genomics.cn (J.W.), albrecht@binf.ku.dk (A.A.), jinxin@genomics.cn (X.J.),

rasmus_nielsen@berkeley.edu (R.N.), xuxun@genomics.cn (X.X.)

https://doi.org/10.1016/j.cell.2018.08.016

SUMMARY

We analyze whole-genome sequencing data from 141,431 Chinese women generated for non-invasive prenatal testing (NIPT). We use these data to characterize the population genetic structure and to investigate genetic associations with maternal and infectious traits. We show that the present day distribution of alleles is a function of both ancient migration and very recent population movements. We reveal novel phenotype-genotype associations, including several replicated associations with height and BMI, an association between maternal age and EMB, and between twin pregnancy and NRG1. Finally, we identify a unique pattern of circulating viral DNA in plasma with high prevalence of hepatitis B and other clinically relevant maternal infections. A GWAS for viral infections identifies an exceptionally strong association between integrated herpesvirus 6 and MOV10L1, which affects piwi-interacting RNA (piRNA) processing and PIWI protein function. These findings demonstrate the great value and potential of accumulating NIPT data for worldwide medical and genetic analyses.

INTRODUCTION

Sufficient large sample size is of fundamental importance in resolving biological questions in population and medical genetics. Given a fixed budget, sample size tends to play a more essential role compared to sequencing depth (Li et al., 2011). Previous studies have demonstrated that sequencing many individuals at a low depth generally provides a better representation of population genetic variation compared to sequencing a more limited number of individuals at a higher depth (Fumagalli, 2013). Furthermore, when using proper imputation techniques, even sequencing at an average depth of <0.1× in a large enough cohort can be a cost-effective strategy for detecting genetic associations for complex traits (Pasaniuc et al., 2012).

Several large-scale national and international sequencing projects have been carried out in the past decade with sample sizes limited to tens of thousands (Auton et al., 2015; Francioli et al., 2014; Gudbjartsson et al., 2015; Maretty et al., 2017; Walter et al., 2015). Increasing the sample sizes of these studies is a major financial and logistical challenge. However, non-invasive prenatal testing (NIPT) for fetal trisomy-by sequencing of maternal plasma cell-free DNA (cfDNA) (Zhang et al., 2015)-has become the fastest adopted molecular test in history and provide an untapped resource for understanding population genetic variation and its associations with phenotypes. To date, over ten millions of NIPT tests have been carried out globally, among which 70% were conducted on Chinese women. These samples can be leveraged for population genetic investigations of population history, large-scale genetic association studies, and viral screening if the technical issues regarding the use of very large, very low depth $(0.06 \times -0.1 \times)$ samples can be addressed.

Here, we analyze NIPT sequencing data of 141,431 pregnant women with informed consent. We demonstrate that allele frequencies can be estimated with high accuracy, allowing further population genetic analyses. We also show that efficient genotype imputation is feasible and can provide considerable mapping power. We use the data to carry out the hitherto largest analysis of population genetic variation in the Chinese population, perform a genome-wide association study (GWAS) on multiple traits in pregnant Chinese women, and survey the distribution of circulating viral DNA in the maternal plasma.

RESULTS

Study Participants and Chromosomal Coverage

The 141,431 participants were recruited from 31 out of the 34 administrative divisions in China (Figure S1A). Each individual was sequenced using 5-10 million single-end reads (35-49 bp), corresponding to a sequencing depth of 0.06× to 0.1× per individual (Figure S1B). The reads were aligned to the hg19 reference using bwa¹⁰ (see STAR Methods) with a resulting combined read depth distribution that is approximately Poisson and closely follows that expected from the ENCODE mappability track (Figures S1C and S1D). Based on read length and observed depth distribution, we identified regions of the genome accessible to high-confidence mapping in the NIPT data, resulting in a total length of around a 2.13 billion base pair accessible genome (75% of the non-N human reference genome sequence). We also re-sequenced DNA derived from white blood cells of 40 participants to a mean depth of 15x. These data were used in the variant calling and genotype imputation evaluation.

Amount of Genetic Variations and Accuracy of Genotype Imputation

Previous standard methods for allele frequency estimation and joint SNP calling, such as those implemented in GATK (DePristo et al., 2011) and Samtools (Li, 2011), did not scale up to sample sizes over a hundred thousand. Therefore, we developed a new method for fast maximum likelihood estimation of allele frequencies and joint SNP calling using a likelihood ratio test (STAR Methods). Using this method for initial screening, we identified 32.5 million bi-allelic candidate SNPs (Table S1). After recalibration using a Gaussian mixture model, we identified a final call set of 9.04 million single nucleotide variants with a transition/transversion ratio of 2.2 for known variants and 2.4 for novel variants, respectively (Figure S2A), consistent with those obtained for 1000 Genomes Project (1KG) variants (Auton et al., 2015). 81.7% of the variants were called in the 1KG Han Chinese individuals (Auton et al., 2015), ~16% of the variants were in the remainder of 1KG or dbSNP database, while 233,966 (2.6%) were novel variants (Figure 1A). 90% of the variants were found in the union of the gnomAD East Asian (Lek et al., 2016) and 1KG Han Chinese call sets (Figure 1B). Using experimental validation, we estimated an upper bound for the false-positive rate (FPR) of SNP calling of 0.2% (Figures S2B-S2E). However, among novel variants, the FPR was \sim 0.32. These SNPs comprised a small proportion of the total number of SNPs, but did include common variants, likely due to unresolved mapping issues. The squared correlation coefficient (R²) of the frequency of the non-reference allele (i.e., alternative allele estimated in our study and that computed in the 1KG Han Chinese) was 0.98 (Figure 1C).

We subsequently imputed genotypes of 8.9 million known variable DNA sites with allele frequency >0.01 using the 1KG Han Chinese as a reference panel (Davies et al., 2016) (STAR Methods). To estimate imputation accuracy, we compared the squared correlation coefficient (R²) between the genotypes called in the medium coverage whole genome sequencing data of the 40 individuals (MC set, 15×) and the imputed genotypes in the low-coverage data (LC set). 2.13 million of the variants were well imputed with an info score >0.4, and with a p value from a chi-square test of Hardy-Weinberg equilibrium larger than 10^{-6} . The mean imputation accuracy of those variants was 0.89, while it was 0.71 for all variants combined (Figure 1D). The imputation accuracy was negatively correlated to the fraction of fetal DNA present in the plasma but the effect was not very pronounced (Figure S2F).

Population Structure, Recent Population History, and Genetic Adaptations

Even though the Chinese is the world's largest population that comprises 1.4 billion people, it is perhaps surprisingly understudied with respect to population genetic history. We applied information of the 141,431 pregnant women regarding digital geographic location and self-reported ethnicity to the study of the genetic variation in China at multiple timescales. Because of the uncertainty in genotype calling, we conduct all population genetic analyses using methods that either sample a single read per individual or use maximum likelihood estimates of allele frequencies without relying on genotype imputation.



A principal component analysis of all the 141,431 participants suggested that the first three principal components reflected sequencing read length, latitudinal genetic differentiation, and the sequencing error rate (Figures S3A-S3D). After removing participants with 49bp read length and with sequencing error rate >0.00325, a principal component analysis of 45,387 self-reported Han Chinese from the 31 administrative divisions showed that the greatest differentiation of Han Chinese is along a latitudinal gradient (Figures S3E and S3F), consistent with previous studies (Chen et al., 2009; Xu et al., 2009). In contrast, there is, perhaps surprisingly, very little differentiation from East to West. This observation may be explained by the fact that a large proportion of the western Han populations in China are recent immigrants organized by the central government starting from 1949 when the People's Republic of China was founded (Liang and White, 1996). While the Han Chinese were found to be relatively genetically homogeneous, there was greater divergence among the minority ethnic groups for both latitude and longitude (Figures 2A and 2B). The most differentiated ethnic groups are the Turkic speaking Uyghur and Kazakhs, who reside in the Xinjiang province, and the Mongols residing primarily in Inner Mongolia. The Xibe, Tibetans, and Hui from central China, the Yi from southwestern China, and the Zhuang and Buyi minorities from southern China, also differ substantially from the Han Chinese that come from the same area. On the other hand, the Manchu from northeastern China were genetically closest to the Han Chinese in that area, consistent with historical accounts (Rhoads, 2000).

Cell

Figure 1. Allele Frequency Spectrum and Imputation Accuracy of the 141,431 Humans (A) Allele frequency spectrum of known and novel variants.

(B) Sharing of variants among NIPT calls (CMDB, n = 141,431), gnomAD East Asian (gnomAD EAS, n = 811), and the Han Chinese population from the 1000 Genome Project (CHN, n = 301).

(C) Comparison of frequency of the non-reference allele between the NIPT estimations (CMDB) and Han Chinese estimations in the 1000 genomes project (CHN).

(D) Count and imputation accuracy of the known variants as a function of minor allele frequency intervals. The 2.1 million known variants are restricted to those with an imputation info score >0.4 and a test p value of Hardy Weinberg equilibrium frequencies of >10⁻⁶. The error bar in red denotes the 97.5% confidence interval, which is generally very small.

See also Figures S1 and S2 and Table S1.

We further explored the patterns of allele sharing between Han Chinese and major global ethnic groups using private alleles defined from the 1KG populations and using outgroup F3 statistics (Peter, 2016) (STAR Methods). In the northwest and central west, we observed private allele sharing with the 1KG European Central European of Utah (CEU) panel both for individuals selfidentified as Han Chinese and for individ-

uals self-identified as belonging to a minority group. The strongest level of private allele sharing with the CEU was observed for people in the most northwest provinces of Xinjiang and Gansu (Figure 2C), likely reflecting the Turkic speaking ancestry in these minorities. When only the Han Chinese were included, the strongest level of allele sharing with Europeans was observed for people in the Qinghai, Gansu, and Ningxia provinces (Figure 2D). These provinces are located in the Hexi corridor, the most important commercial hub on the Silk Road connecting China to the west since the establishment of the Han Dynasty (206 BC) (Yang et al., 2008). Thus, one potential explanation for the Western ancestry observed in these provinces is gene flow related to their location on the Silk Road. We also observed a pattern of increased allele sharing with the 1KG Indian ITU reference panels in southwestern populations from Xinjiang, Tibet, Yunnan, Guangxi, and Hainan provinces (Figures 2E and 2F), consistent with their geographic proximity to the Indian subcontinent (Yang et al., 2017). Analyses based on the F3 statistic are mostly consistent for the CEU analysis, but for the ITU analysis, we also show high affinity between the Han Chinese in northern provinces and the ITU, likely due to the shared ancestry of the CEU and ITU populations. Furthermore, we applied the F3 statistic to learn patterns of allele sharing between the Chinese provincial populations and 1KGP neighbor populations including three Chinese populations, the Japanese, and the Vietnamese. We observe a pattern of allele sharing among the 33 administrative divisions reflecting the geographical origin of the 1KGP populations (Figures S3G-S3K).



Figure 2. Population Structure and Distribution of Allele Sharing with Related Populations in 1KGP

(A) Geographical distribution of the 36 minorities. Size of the circle reflects the number of minority individuals.
(B) Principal component analysis of the 36 minorities. A random selection of an equal number of Han Chinese matching the same city of each minority are included and the number of the minority and the number of the reflects the number of the number of

and shown as gray colors. English names of the minorities and the number of randomly selected participants from each ethnic and geographical groups from 96,880 participants after QC on error rate and read length are shown in the legend. (C and D) Private allele sharing between each administrative division for all ethnic groups (C) or only the Han (D) and the CEU and ITU reference populations.

(E and F) Private alleles sharing between each administrative division for all ethnic groups (c) of only the Han (F) and the ITU reference populations. (E and F) Private alleles sharing between each administrative division for all ethnic groups (E) or only the Han (F) and the ITU reference populations in the 1KGP. Color corresponds to the private allele frequency defined in the main text (i.e., the frequency of sampling an allele from each division that is private to reference CEU or ITU populations).

See also Figure S3.

Interestingly, we found that the CHB, although annotated as the Han Chinese from Beijing, did not have the closest affinity with Beijing individuals but tended to be closer to populations in the coastal provinces: Shandong, Zhejiang, Jiangsu, Fujian, and Jiangxi (Figure S3G). This likely reflects the recent multiethnic migration into Beijing consistent with the demographic information available for our samples. We also investigated the inter-provincial allele sharing between Han Chinese in the Chinese administrative divisions. The difference in f3 statistic among provinces is very small, but all southern provinces show more genetic affinity with other southern coastal provinces (results not shown). This observation likely reflects a combination of internal migration events organized by the central government since 1949 (Liang and White, 1996) and the country's oriented movement of labor from the interior to the coastal areas since 1979 (Liang and Ma, 2004).

We inferred selection within Han Chinese populations using two approaches. First, we identified variants with significant differentiation along each PC compared to a null distribution expected under a model of genetic drift (STAR Methods). Second, we conducted a scan of the rarer but more important pathogenic variants in the Clinvar database (Landrum et al., 2014) by statistically comparing allele frequency differences of those loci among North, Central and South Han Chinese against a null distribution generated from the genome-wide data. In the PC scan, we identified six loci showing genome-wide significance across latitude: *LILRA3, CR1, FADS2, DOCK9, ABCC11,* and a cluster of *IGH* genes (Figure 3A). The *CR1, DOCK9,* and the *IGH* genes display a higher allele frequency in the south while the *FADS2, ABCC11,*



Figure 3. Genetic Adaptation in Han Chinese Population

(A) Manhattan plot showing the detected selection signals in Han Chinese population across the first principal component. VEP annotated names of the gene loci under selection are displayed.

(B–G) Derived allele frequency per Chinese administrative division for the lead SNP in loci under selection across latitude. Shown is the derived allele frequency distribution of the lead SNP in the *CR1* loci (B), *FADS2* loci (C), *ELK2AP-MIR4507* loci corresponding to the IGH-gene cluster (D), *ABCC11* loci (E), *DOCK9* loci (F), and *LILRA3* loci (G) in (A). The number and the corresponding color in the legend and map indicate derived allele frequency estimated from NIPT data. (H–O) Allele frequency per administrative division for the ClinVar pathogenic variants with a significant difference of allele frequencies across North, Central, and South regions, including the ClinVar pathogenic varients associated with Meckel syndrome type 2 (H), complement component 9 deficiency (I), deafness (J), deficiency of ferroxidase (K), Usher syndrome (L), albinism (M), non-syndromic genetic deafness (N), and G6PD CANTON (O). Number and color in the legend and map represent allele frequency estimated from NIPT data for the risk allele recorded in the ClinVar database. See also Tables S2 and S3.

and *LILRA3* genes display a higher derived allele frequency in the north (Figures 3B–3G; Table S1). Three of these loci are known to be related to immune responses (the *IGH* genes, *LILRA3* and *CR1*). *DOCK9* is associated with bipolar disease and has been

previously shown to under selection in East Asians (Suo et al., 2012). *FADS2* is a well-known target of selection associated with changes in diet to, or from, a diet with a high content of animal fat and has previously been inferred to have been targeted



by selection in Inuit (Fumagalli et al., 2015), South Asians (Kothapalli et al., 2016), Europeans (Buckley et al., 2017), and in Africa (Mathias et al., 2012). Our results suggest more recent selection has also been acting within China. The ABCC11 locus is famously associated with earwax type and has previously been shown to be under selection in Asian, Native American, and European populations (Ohashi et al., 2011), and our results demonstrate that this locus is also under differential selection within China. We also investigated the geographical distribution of possibly pathogenic variants compiled from the ClinVar dataset (Landrum et al., 2014) as candidates for loci under selection. We calculated measures of allele frequency differentiation (Fisher's exact test between northern, central, and southern Han Chinese, comparing against a frequency-matched dataset of 100,000 SNPs chosen at random) (STAR Methods). We identified and reported the nine out of the 42,058 possibly pathogenic variants in eight genes that display the most significant allele frequency difference among the three geographical regions (Fisher exact test with p value $< 10^{-6}$, percentile p value < 5e-3) (Figures 3H-3O; Table S2). Those SNPs include rs72554665, a polymorphic site in G6PD, a gene associated with resistance to malaria (Nkhoma et al., 2009). This variant has a higher frequency in southern China, consistent with historically higher incidence rates of malaria in southern China than in northern and central China.

Phenotype-Genotype Associations of Multiple Complex Traits

In the following, we demonstrate that NIPT data can be used effectively in GWAS. We first investigated associations with Figure 4. Genome-wide Significant Signals for Two Common Quantitative Traits and Two Traits Related to Reproductive Process Known loci, defined as significant variants with a known association with the investigated trait in the GWAS catalog (e90_r2017-10-10) within 1 Mbp region are marked in black. Novel loci are marked in red. For loci where the lead SNP is located in the intergenic region, the most close gene was plotted. Detailed information about the loci can be found in Tables S5 and S6.

See also Figures S4 and S5 and Table S4.

two common traits, height and BMI among 61.7K individuals with both phenotypes recorded (Figures S4A and S4B). We applied a score test (Skotte et al., 2012) to test the association between the traits and the genotype probabilities for each of the previously mentioned ~2 million imputed variants, incorporating covariates such as the first to fifth principal components (Figure S3A), maternal age, gestational age of the fetus, fetus sex, etc. (STAR Methods). The genomic control factor lambda for height and BMI were 1.51 and 1.32, respectively (QQ-plots in Figures S4C

and S4D). Due to the high polygenicity of the traits, we also evaluated inflation of our test statistics using linkage disequilibrium (LD) score regression that did not show severe inflation (intercept 1.03, SE 0.03 and 1.10, SE 0.02, attenuation ratio 0.05, SE 0.04 and 0.24, SE 0.05 for height and BMI, respectively), suggesting that confounding factors, such as population structure, were generally well controlled (Table S4). The estimated SNP heritability obtained from the LD score regression for height and BMI are 0.48 and 0.10, respectively. A comparison of the LD score regression statistics between Giant, UK Biobank, and our study can be found in Table S4. We note that strong inflation was observed if covariates were not applied in the test model (genomic control factor lambda for height and BMI are 9.71 and 2.68, respectively).

In total, 48 and 13 loci reached genome-wide significance for association with height and BMI, respectively, at the classical 5×10^{-8} genome-wide significance level (Figures 4A and 4B; Table 1). Forty-one of the height loci were previously reported, although only 36 of them have previously been found in Asian populations (MacArthur et al., 2017). Seven height loci located in or around the genes *UBQLN2*, *MIR325HG*, *MAST2*, *STRBP/ZBTB26*, *C11orf24-LRP5*, *ARHGEF12*, and *LINC00261* have not been previously reported (Figures S5A–S5F; Table 2). There was one new signal in the intronic region of DNA2, a locus first identified in GIANT (Wood et al., 2014) and another independent signal in the nearby gene MYPN was reported (conditional p value = 4.8e–8). Three BMI-associated loci in the genes PLD5, TRPC6, and CBLN4 were also not previously reported (Figures S5H–S5J; Table 2). We attempted to replicate the novel

Table 1. Replication Status of Height- and BMI-Associated Loci

Number of Variants^b Mean R^c Known Loci^d Infoscore^a Novel Loci Known (Replicated |Notreplicated)^e Novel (Replicated |Notreplicated) Height 0.8 788,385 0.95 24 4 24|0 4|0 0.7 6 1,552,640 0.92 38 38|0 5|1 6 0.6 1,922,780 0.9 40 40|0 5|1 0.5 2,057,331 0.89 41 7 401 6|1 2,104,769 0.89 7 40|1 0.4 41 61 BMI 0.8 788,385 2 6|0 0.95 6 1|1 0.7 8 3 8|0 2|1 1,552,640 0.92 0.6 1,922,780 0.9 10 3 10|0 2|1 0.5 2,057,331 0.89 10 3 10|0 2|1 0.4 2,104,769 0.89 10 3 100 21

^aInfo score is provided by STITCH that measures the ratio of the observed statistical information of the population allele frequency and the complete information (see STAR Methods).

^bNumber of variants refer to number of imputed variants with minor allele frequency >0.01 and p value of hardy Weinberg equilibrium test >10e-6. ^cMean R² refers to the true imputation accuracy comparing the imputed genotype dosage and the true genotypes of the 40 NIPT samples sequenced to 15×.

^dLoci is defined as a 1 Mb window extending 500 kbp at both the 5' and 3' ends centering on the snp with smallest p value in the window. Known refers to the existence of one or more known SNPs in the GWAS catalog found within the 1 Mb window.

^eReplicated refer to the number of loci that have p value <0.05 divided by the number of associated loci and same beta direction in any one of the CKB, Giant, or UK Biobank test sets. Not replicated denotes the number of loci that are not replicated in all three test sets.

and known associations for height and BMI in 32,000 genotyped Chinese participants from the China Kadoorie Biobank (CKB) cohort (Chen et al., 2011) and in the results of the GIANT consortia (Yengo et al., 2018) and the UK Biobank (Ben Neale's website, see STAR Methods). When regressing the effect size of the genome-wide significant variants in the three test sets including CKB, Giant, and UK Biobank on the discovery set (i.e., the NIPT result), we observed a higher regression slopes for the test set of the same Chinese ancestry compared to the test sets of the European ancestry for height (CKB, slope 0.88; Giant, slope 0.614; UK Biobank, slope 0.495) and BMI (CKB, slope 0.709; Giant, slope 0.463; UK Biobank, slope 0.434), respectively (Figures S4E–S4J). When comparing the beta direction and the p value of the lead or proxy SNPs between the discovery and test sets, in almost all cases the p values and the effect sizes of the SNPs are similar (Tables 2, S5, and S6). Only one known and one novel locus for height (FADS2 and LINC00261) and one locus for BMI (TRPC6) failed to be replicated at the significance level when correcting for multiple testing (p < 0.05 divided by the number of test loci) in all the CKB and the GIANT/UK Biobank analyses. However, the variant in FADS2 is nominally significant in the CKB cohort with a p value of 0.002.

The height and BMI results provide a proof of concept for the use of NIPT data in GWAS and suggest that NIPT data can be used to investigate fertility and pregnancy related traits that otherwise would be prohibitively difficult to investigate in such large samples. To illustrate this, we investigated associations for two novel traits: (1) maternal age, which is expected to correlate with female fertility broadly defined, but may also be affected by several other factors, and (2) twin pregnancy. The maternal age distribution follows a bimodal distribution in the NIPT participants (Figure S4K). Similarly to height and BMI, we again do not

observe any severe inflation when including covariates for both maternal age (lambda = 1.29) and twin pregnancy (lambda = 1.06) (QQ plot see Figures S4L and S4M). Strong inflation was also observed for these two traits if covariates were not applied in the test model for maternal age (lambda = 1.88) and twin pregnancy (lambda = 1.36). We find one significant association peak for maternal age located between the HCN1 and EMB loci (rs16828019, p = 1.38E-11) (Figures 4 and S5K). This signal is located near a previously identified association peak for age at first birth, originally reported to be in the HCN1 gene (Barban et al., 2016). Our lead SNP is closer in location to EMB gene which encodes embigin, a transmembrane glycoprotein that is preferentially expressed in the early stages of embryogenesis and enhances integrin-mediated cell-substratum adhesion in mice (Huang et al., 1993). It shows particularly high expression during early post implantation embryogenesis and is therefore a strong candidate gene for further studies of infertility in humans. In European countries, maternal age is associated with education attainment, which is a proxy for intelligence (Barban et al., 2016). We do not have records for education attainment in our study. Whether education has an impact on maternal age in the Chinese population requires further investigation.

Out of 137,646 individuals with ultrasound scans, 476 had more than one fetus (e.g., twins). The lead associated SNP for this trait was located in the gene *NRG1* (rs12056727, p = 5.93E-9, odds ratio = 1.99) (Figures 4 and S5L) and is a very strong expression quantitative trait locus (eQTL) for expression in the thyroid (effect size = 0.5, p value 1.2e-19) in the GTEx database (Carithers and Moore, 2015). The tested allele T increases the twinning probability and the NRG1 expression in the thyroid. Furthermore, the SNP is associated with hyperthyroidism in the UK BioBank (Sudlow et al., 2015) (p = 1.7e-7,

Table 2. Replication Statistics of the New Loci Associated with the Height and BMI Traits Found in NIPT										
	NIPT		СКВ		Giant		UK Biobank			
GENE	SNP	MAF	Р	Beta	Р	Beta	Р	Beta	Р	Beta
UBQLN2-LINC01420	rs7391861	0.34	1.00E-09	0.04	1.85E-07	0.04	NA	NA	NA	NA
MIR325HG-FGF16	rs4892720	0.21	5.67E-13	0.05	4.15E-06	0.04	NA	NA	NA	NA
MAST2	rs7520050	0.35	4.18E-09	-0.04	0.01	-0.02	4.60E-23 ^a	-0.01^{a}	3.49E-07	-0.01
STRBP	rs10818797	0.33	1.41E-09	0.04	0.52	0.01	1.40E-14	0.02	9.08E-06	0.01
C11orf24-LRP5	rs450416	0.49	1.04E-08	0.04	0.09	0.01	1.10E-05 ^a	-0.01^{a}	1.04E-07	0.01
ARHGEF12	rs894839	0.28	3.72E-08	-0.04	6.51E-08	-0.05	6.70E-04 ^a	-0.01^{a}	1.57E-03	-0.01
LINC00261	rs1203887	0.04	4.07E-08	-0.09	0.02	-0.05	1.90E-03 ^a	0.01b	0.90	0.003
PLD5-LINC01347	rs2780797	0.45	2.57E-08	0.04	0.01	0.02	2.10E-06	0.01	3.69E-03	0.01
LOC101054525-TRPC6	rs12803364	0.26	1.25E-09	-0.04	0.37	0.01	0.66	0.00	0.54	-0.003
LINC01441-CBLN4	rs59271815	0.16	3.24E-10	-0.05	2.39E-04	-0.04	3.60E-08 ^a	-0.01^{a}	1.67E-05	-0.01
	2. Replication Statistic GENE UBQLN2-LINC01420 MIR325HG-FGF16 MAST2 STRBP C11orf24-LRP5 ARHGEF12 LINC00261 PLD5-LINC01347 LOC101054525-TRPC6 LINC01441-CBLN4	2. Replication Statistics of the New GENE SNP UBQLN2-LINC01420 rs7391861 MIR325HG-FGF16 rs4892720 MAST2 rs7520050 STRBP rs10818797 C110rf24-LRP5 rs450416 ARHGEF12 rs894839 LINC00261 rs1203887 PLD5-LINC01347 rs2780797 LOC101054525-TRPC6 rs12803364 LINC01441-CBLN4 rs59271815	Asymptotic Statistics Statist	Application Statistics of the New Loci Associated w NIPT GENE SNP MAF P UBQLN2-LINC01420 rs7391861 0.34 1.00E-09 MIR325HG-FGF16 rs4892720 0.21 5.67E-13 MAST2 rs7520050 0.35 4.18E-09 STRBP rs10818797 0.33 1.41E-09 C110rf24-LRP5 rs450416 0.49 1.04E-08 ARHGEF12 rs894839 0.28 3.72E-08 LINC00261 rs1203887 0.04 4.07E-08 PLD5-LINC01347 rs2780797 0.45 2.57E-08 LOC101054525-TRPC6 rs12803364 0.26 1.25E-09 LINC01441-CBLN4 rs59271815 0.16 3.24E-10	Application Statistics of the New Loci Associated with the H NIPT GENE SNP MAF P Beta UBQLN2-LINC01420 rs7391861 0.34 1.00E-09 0.04 MIR325HG-FGF16 rs4892720 0.21 5.67E-13 0.05 MAST2 rs7520050 0.35 4.18E-09 -0.04 STRBP rs10818797 0.33 1.41E-09 0.04 C110rf24-LRP5 rs450416 0.49 1.04E-08 0.04 ARHGEF12 rs894839 0.28 3.72E-08 -0.04 LINC00261 rs1203887 0.04 4.07E-08 -0.09 PLD5-LINC01347 rs2780797 0.45 2.57E-08 0.04 LOC101054525-TRPC6 rs12803364 0.26 1.25E-09 -0.04 LINC01441-CBLN4 rs59271815 0.16 3.24E-10 -0.05	Application Statistics of the New Loci Associated with the Height and B NIPT CKB GENE SNP MAF P Beta P UBQLN2-LINC01420 rs7391861 0.34 1.00E-09 0.04 1.85E-07 MIR325HG-FGF16 rs4892720 0.21 5.67E-13 0.05 4.15E-06 MAST2 rs7520050 0.35 4.18E-09 -0.04 0.01 STRBP rs10818797 0.33 1.41E-09 0.04 0.52 C110rf24-LRP5 rs450416 0.49 1.04E-08 0.04 0.09 ARHGEF12 rs894839 0.28 3.72E-08 -0.04 6.51E-08 LINC00261 rs1203887 0.04 4.07E-08 -0.09 0.02 PLD5-LINC01347 rs2780797 0.45 2.57E-08 0.04 0.01 LOC101054525-TRPC6 rs12803364 0.26 1.25E-09 -0.04 0.37 LINC01441-CBLN4 rs59271815 0.16 3.24E-10 -0.05 2.39E-04	SNP MAF P Beta P Beta UBQLN2-LINC01420 rs7391861 0.34 1.00E-09 0.04 1.85E-07 0.04 MIR325HG-FGF16 rs4892720 0.21 5.67E-13 0.05 4.15E-06 0.04 MAST2 rs7520050 0.35 4.18E-09 -0.04 0.01 -0.02 STRBP rs10818797 0.33 1.41E-09 0.04 0.52 0.01 C110rf24-LRP5 rs450416 0.49 1.04E-08 0.04 0.51E-08 -0.05 LINC00261 rs1203887 0.44 4.07E-08 -0.09 0.02 -0.05 PLD5-LINC01347 rs2780797 0.45 2.57E-08 0.04 0.01 0.02 LOC101054525-TRPC6 rs1280364 0.26 1.25E-09 -0.04 0.37 0.01 LINC01441-CBLN4 rs59271815 0.16 3.24E-10 -0.05 2.39E-04 -0.04	SNP MAF P Beta P Beta P UBQLN2-LINC01420 rs7391861 0.34 1.00E-09 0.04 1.85E-07 0.04 NA MIR325HG-FGF16 rs4892720 0.21 5.67E-13 0.05 4.15E-06 0.04 NA MAST2 rs7520050 0.35 4.18E-09 -0.04 0.01 -0.02 4.60E-23 ^a STRBP rs10818797 0.33 1.41E-09 0.04 0.52 0.01 1.40E-14 C110rf24-LRP5 rs450416 0.49 1.04E-08 0.04 0.52 0.01 1.10E-05 ^a ARHGEF12 rs894839 0.28 3.72E-08 -0.04 6.51E-08 -0.05 6.70E-04 ^a LINC00261 rs1203887 0.04 4.07E-08 -0.09 0.02 -0.05 1.90E-03 ^a PLD5-LINC01347 rs2780797 0.45 2.57E-08 0.04 0.01 0.02 2.10E-06 LOC101054525-TRPC6 rs1280364 0.26 1.25E-09 -0.04	2. Replication Statistics of the New Loci Associated with the Height and BMI Traits Found in NIPT NIPT CKB Giant GENE SNP MAF P Beta P Beta UBQLN2-LINC01420 rs7391861 0.34 1.00E-09 0.04 1.85E-07 0.04 NA NA MIR325HG-FGF16 rs4892720 0.21 5.67E-13 0.05 4.15E-06 0.04 NA NA MAST2 rs7520050 0.35 4.18E-09 -0.04 0.01 -0.02 4.60E-23 ^a -0.01 ^a STRBP rs10818797 0.33 1.41E-09 0.04 0.52 0.01 1.40E-14 0.02 C110rf24-LRP5 rs450416 0.49 1.04E-08 0.04 0.09 0.01 1.10E-05 ^a -0.01 ^a ARHGEF12 rs894839 0.28 3.72E-08 -0.04 6.51E-08 -0.05 6.70E-04 ^a -0.01 ^a LINC00261 rs120387 0.44	2. Replication Statistics of the New Loci Associated with the Height and BMI Traits Found in NIPT NIPT CKB Giant UK Biobani GENE SNP MAF P Beta P Beta P Beta P UBQLN2-LINC01420 rs7391861 0.34 1.00E-09 0.04 1.85E-07 0.04 NA NA NA MIR325HG-FGF16 rs4892720 0.21 5.67E-13 0.05 4.15E-06 0.04 NA NA NA MAST2 rs7520050 0.35 4.18E-09 -0.04 0.01 -0.02 4.60E-23a -0.01a 3.49E-07 STRBP rs10818797 0.33 1.41E-09 0.04 0.52 0.01 1.40E-14 0.02 9.08E-06 C11orf24-LRP5 rs450416 0.49 1.04E-08 0.04 0.09 0.01 1.10E-05a -0.01a 1.04E-07 ARHGEF12 rs894839 0.28 3.72E-08 -0.04 6.51E-08 -0.05 6.70E-04a -0.01a 1.57E-03

^aProxy SNP that have similar p value and effect size in NIPT.

odds ratio = 1.15). Thyroid function has previously been associated with fertility. The *NRG1* gene has mostly been investigated for its effects on behavior (schizophrenia in humans and response to stress and anxiety in rodents), but at least one study notes that matings between knockouts in mice have smaller litter size (Britto et al., 2004). However, the exact reason for this is unknown. More interestingly, twin pregnancies tend to be associated with lower levels of the thyroid-stimulating hormone (TSH) (Soldin, 2006), consistent with the SNP association with hyperthyroidism, which generally involves increased thyroid hormone levels and decreased TSH levels.

Circulating Viral DNA in Maternal Plasma

Despite its importance for public health, few studies have been carried out on the distribution of viral DNA in blood plasma (the virome) at the population level. However, the sequence technologies used in NIPT studies provide an untapped resource for understanding viral epidemiology. We investigated the plasma virome of 138,882 participants by querying reads that do not align to the human genome against the NCBI viral sequence database (Sayers et al., 2009) using BLAST (Altschul et al., 1990) (STAR Methods). The plasma virome, cleaned for phages and virus with low genome coverage (<10%), is represented in Figures 5A and 5B, and all observed viruses, without removing phage, can be found in Figures S6A-S6C. We examined sequencing coverage to individual viruses to determine potential misclassification or contamination, where prevalent viruses with relatively even coverage support true virus identification (Figure S7). Most viruses detected have even sequence coverage, while a few display localized peaks. To understand if these peaks are related to human homology, we aligned viral reference sequences to the human reference genome, hg19 and non-EBV decoys (Table S6). We found that only the localized peak of HCV corresponds to a human homologous sequence, yet these sequences are only found in a small number of our subjects (n = 3). The peaks in coverage may represent highly conserved areas of viruses, misclassification of human sequences, uneven production of viral DNA, select viral integration events, or introduced viral DNA from vaccine. Further investigation and validation is required.

Interestingly, the blood virome in a recent study of Europeans (Moustafa et al., 2017) appeared to have a different viral distribution compared to the pregnant Chinese women participants analyzed in this study (Figure 5A). The use of differing sequencing approaches is a challenge to direct comparisons, but our participants were significantly enriched for hepatitis B virus (HBV) and Parvovirus B19 DNA and showed a lower prevalence of human herpesvirus 7 (HHV-7) DNA compared to Europeans (Moustafa et al., 2017). The prevalence of HBV DNA across Chinese populations estimated in our study is $\sim 2.5\%$ (with high mean abundance 25.6), which is less than the 9.8% prevalence reported in a Chinese 2014 population survey of HBV antibodies (IgM to HBV core antigen or surface antigen [HBsAg⁺]) (Yan et al., 2014). These differences are likely due to varying estimates derived from circulating HBV DNA versus antibody prevalence, the enrichment of our sample for relatively affluent younger women, and the late adoption of the HBV vaccine in China (since 1992) (Yan et al., 2014). We detected 3,421 participants for HBV DNA, yet 1,911 participants self-reported some type of HBV infection. As expected, we have the greatest sensitivity (78.7%) to detect reported active HBV infection (HBsAg⁺). Of 1,032 individuals reporting latent infection (HBsAg⁻), we detected HBV DNA in 53 subjects, where the lower sensitivity (5.1%) is likely due to low levels of circulating HBV DNA during latent infection (Figure S6C). Interestingly, we detected HBV DNA in 2,959 individuals who did not report any HBV infection, suggesting an additional and potentially clinically important use of NIPT where circulating HBV DNA is associated with fetal transmission (Zou et al., 2012).

We detected many hits for human endogenous retrovirus K (HERV-K) (prevalence, \sim 2.1%; mean abundance, 0.22) which was previously shown to be active, capable of expression in humans, and associated with HIV infection (Zwolińska et al., 2013). All humans carry multiple copies of HERV-K in their genomes,



Figure 5. The Viral Spectrum in Maternal Plasma

(A) Prevalence of infection among the investigated population.

(B) Distribution of abundance by each virus. Each dot represents the abundance of one individual.

(C) Manhattan plot showing results from GWAS of carriers of high abundance ciHHV-6A/B versus non-carriers.

(D) Locus Zoom plot denotes lead snp (rs73185306) and the correlated snps around the region of MLC and MOV10L1 genes.

(E-G) Geographic distribution of prevalence for the three most prevalent virus (i.e., HBV, E; HERV-K113, F; and HHV6A/6B, G).

and the recovery of non-aligning HERV-K in a subset of subjects may be the result of insertionally polymorphic HERV-K (Wildschutte et al., 2016) or the result of co-option of HERV-K sequences by other exogenous viruses. Herpesvirus 6A/B (HHV-6A/B), the third most common viral group, has a prevalence of 0.8% (mean abundance, 0.48). HHV-6A/B were grouped together due to high co-occurrence and potential for misclassification due to sequence homology. A distinct bimodal distribution clusters high abundance HHV-6A/B (abundance $>10^{-0.5}$) (Figure 5A), likely separating chromosomally integrated HHV-6A/B (ciHHV-6A/B) from non-integrated circulating HHV-6A/B as noted in previous studies (Moustafa et al., 2017). Human herpesvirus 5 or cytomegalovirus (HHV-5 or CMV) is the fourth most common infection with a prevalence of 0.40% (mean abundance, 0.03). This virus is of particular interest in pregnant women as CMV is one of the leading causes of birth defects (Cheeran et al., 2009). Parvovirus B19 has a prevalence of 0.39% (mean abundance, 1.68) and is of clinical relevance to pregnant women as active infection can cause fetal anemia and death.

To identify germline polymorphisms associated with viral infections, we carried out an association study for each of the major viruses, comparing infected individuals to a control group of 90,531 participants who have no detectable virus in the NIPT sequencing data. We identified an intronic variant, rs73185306, in the MLC1-MOV10L1 region that is significantly associated with the presence of high abundance ciHHV-6A/B (n = 653) but not low abundance HHV-6A/B (n = 1,556) (odds ratio = 3.4, p value = 7.3e-66) (Figures 5B, 5C, S6E, and S6F). rs73185306 is an eQTL for both MCL1 and MOV10L1 genes suggesting a functional role (Lonsdale et al., 2013). To ensure that this strong association was not due to alignment error and homology in the MLC1-MOV10L1 region, we aligned HHV-6A and HHV-6B genomes back to the human genome and found no sequence homology in this genomic region (Table S7). The MCL1 gene is involved in myeloid cell differentiation and has been shown to be upregulated during herpesvirus infection (CMV, EBV, and HHV-8). The MOV10L1 gene is known to be associated with platelet distribution (Astle et al., 2016) that is also correlated with severity of hepatitis B infection (Karagoz et al., 2014), although we found no association with circulating HBV DNA. Intriguingly, MOV10L1 is a PIWI interacting RNA helicase that is active during spermatogenesis and functions as a repressor of retrotransposons (Vourekas et al., 2015). We suspect that the PIWI-interacting RNA represses HHV-6A/B integration, and polymorphisms in this gene allow for more efficient integration of HHV-6A/B during spermatogenesis. We observed no other significant genome-wide SNP associations with other viruses.

Finally, we explored the geographic distribution of detected viruses in the studied population throughout China. We mapped the prevalence of viral sequences among 30 administrative divisions with more than 100 participants. Tibet was excluded due to a small sample size of 13 individuals. We observed different geographic patterns for viruses with occurrence sample size >1,000 (Figures 5D-5F). We observed that HBV DNA in serum has a higher prevalence in southern China compared to central and northern China (Figure 5D), while previous studies have shown higher HBV antibody prevalence in northern China (Yan et al., 2014). We speculate that subgenotype resolution differences in HBV may contribute to circulating DNA levels, a clinically relevant predictor of fetal HBV transmission and tumor progression (Zou et al., 2012). We observed a similar geographic distribution between HBV and HERV-K113 (Figure 5E). HERV-K directly interacts with exogenous viruses such as HIV to reduce the infectivity of the resulting chimeric virions (Zwolińska et al., 2013). HERV-K-HBV co-option may explain the observed geographic co-occurrence. Alternatively, apolio B RNA editing catalytic component (APOBEC) mediated mutation of HERV-K may increase during HBV infection leading to initial alignment errors and subsequent classification by BLAST (Lee et al., 2008; Vartanian et al., 2010).

DISCUSSION

In this study, we develop statistical methods for analysis of NIPT data and illustrate the utility of these data for population genetics, association mapping studies, and studies of the human plasma virome. Despite the low sequencing coverage, we demonstrate that accurate genotype imputation is possible, and we discover novel loci associated with height, BMI, and two pregnancy-related traits. The results illustrate the power and feasibility of association mapping using NIPT data.

We also leverage the data for population genetic inferences and show that the majority Han Chinese have evidence of isolation by distance latitudinally, but not longitudinally, presumably due to recent population movements. In contrast, the genetic diversity in ethnic minorities roughly mirrors geography. We successfully identify known and novel loci that are under selection based on the small allele frequency differences in the Han population. Finally, we identify circulating DNA of viruses with clinical relevance in pregnancy (HBV, CMV, ParvoB19) and reveal a different viral sequencing distribution spectrum compared to Europeans. We analyze genetic and viromic data together and reveal a highly significant association between suspected ciHHV-6A/B and MOV10L1-MCL1 hinting at a possible germline variant affecting the integration of HHV-6A/B.

Our results illustrate the utility of NIPT data for medical genetic studies, particularly for understanding traits related to fertility and pregnancy. Furthermore, the availability of large samples of shotgun DNA sequencing from blood opens up new avenues for investigating hypotheses regarding interactions between viruses and host DNA genetic variability. As NIPT testing expands to millions of individuals globally, obtaining informed consent for patients and effective digital curation of medical records should be prioritized by the medical community.

STAR***METHODS**

Detailed methods are provided in the online version of this paper and include the following:

- KEY RESOURCES TABLE
- CONTACT FOR REAGENT AND RESOURCE SHARING
- EXPERIMENTAL MODEL AND SUBJECT DETAILS
- METHOD DETAILS
 - $\odot\,$ Sequencing and QC
 - Medium depth sequencing of 40 participants
 - Alignment of reads against hg19 and definition of accessible region
 - Variant discovery and allele frequency estimation
 - \odot Annotation
 - Imputation
 - Principal component analysis
 - F3-statistic and private allele frequency analysis
 - Detection of selection across PC coordinates
 - Detection of clinvar pathogenic variants displaying significant allele frequency differentiation
 - Identification of genetic variants significantly associated with a trait

- Viral sequence analysis
- QUANTIFICATION AND STATISTICAL ANALYSIS
- Maximum likelihood estimation of allele frequency
 DATA AND SOFTWARE AVAILABILITY
- DATA AND SOFTWARE AVAILABILITY

SUPPLEMENTAL INFORMATION

Supplemental Information includes seven figures and seven tables and can be found with this article online at https://doi.org/10.1016/j.cell.2018.08.016.

ACKNOWLEDGMENTS

We are grateful to all the participants and BGI colleagues participating in the project. We thank China National GeneBank for computational support, Miaolan Cen from BGI for organising necessary experimental resources in the project, Xiaofeng Wei for support on setting up the CMDB database, Dr. Heng Li for useful discussions on development of the BaseVar to call SNP from NIPT sequencing data, and Dr. Chuan Li for suggestions on structuring the manuscript. Particularly, we would like to thank the professional technical support service provided by Wei Lin, Ruibo Li, Li Tian, Jingren Zhou, Heshan Lin, Changliang Xu, Shaoqing Dai, Qi Zhu, Tiecheng Dengfrom Alibaba Cloud Corperation, and the supercomputing capabilities provided by Alibaba Cloud MaxCompute and BatchCompute products that greatly shortened the research process. We thank the Tianhe Supercomputer Center for computational support. This project was supported by the Natural Science Foundation of Guangdong Province, China (2017A030306026), Funds for Distinguished Young Scholar of South China University of China (2017JQ017), and Funds for Industrial PhD by Innovation Fund Denmark (4135-00130B).

AUTHOR CONTRIBUTIONS

Conceptualization, X.X., R.N., X.J., A.A., S. Liu, J.S., M.M., Y.Y., and Jian Wang; Methodology, S. Liu, A.A., S.H., R.N., S.S.F., R.W.D., T.K., A.K., M.A.Y., and Q.F.; Validation, R.G.W., K.L., Z. Chen, F.C., S. Liu, R.L., and S.H.; Investigation, S. Liu, S.H., L.F., Y. Zhang, H.X., S. Li, Z. Cai, Y. Yuan., and Q.L.; Formal Analysis, S. Liu, S.H., Z.L., L.L., R.L., Y. Zhou, and J.J.; Resources, Y. Yin, L. Zhao, H.Z., and W.W.; Data Curation, Q.L., Y. Yuan, H.Z., L. Zhou, and Jun Wang; Writing – Original Draft, S. Liu; Writing – Review & Editing, S. Liu, R.N., A.A., S.S.F., X.J., J.S., Z. Chen, H.Y., M.A.Y., and Q.F.; Visualization, S. Liu, S.H., Z.L., H.Y., M.A.Y., and Q.F.; Visualization, S. Liu, S.H., Z.L., L.L., R.L., Y. Zhou, and J.J.; Supervision, X.X., R.N., X.J., A.A., and Jian Wang; Project Administration, X.J., S.H., S. Liu, and F.C.; Funding Acquisition, X.J. and S. Liu.

DECLARATION OF INTERESTS

The authors declare no competing interests.

Received: March 20, 2018 Revised: June 12, 2018 Accepted: August 8, 2018 Published: October 4, 2018

REFERENCES

Altschul, S.F., Gish, W., Miller, W., Myers, E.W., and Lipman, D.J. (1990). Basic local alignment search tool. J. Mol. Biol. *215*, 403–410.

Astle, W.J., Elding, H., Jiang, T., Allen, D., Ruklisa, D., Mann, A.L., Mead, D., Bouman, H., Riveros-Mckay, F., Kostadima, M.A., et al. (2016). The allelic landscape of human blood cell trait variation and links to common complex disease. Cell *167*, 1415–1429.

Auton, A., Brooks, L.D., Durbin, R.M., Garrison, E.P., Kang, H.M., Korbel, J.O., Marchini, J.L., McCarthy, S., McVean, G.A., and Abecasis, G.R.; 1000 Genomes Project Consortium (2015). A global reference for human genetic variation. Nature 526, 68–74.

Barban, N., Jansen, R., de Vlaming, R., Vaez, A., Mandemakers, J.J., Tropf, F.C., Shen, X., Wilson, J.F., Chasman, D.I., Nolte, I.M., et al.; BIOS Consortium; LifeLines Cohort Study (2016). Genome-wide analysis identifies 12 loci influencing human reproductive behavior. Nat. Genet. *48*, 1462–1472.

Britto, J.M., Lukehurst, S., Weller, R., Fraser, C., Qiu, Y., Hertzog, P., and Busfield, S.J. (2004). Generation and characterization of neuregulin-2-deficient mice. Mol. Cell. Biol. *24*, 8221–8226.

Buckley, M.T., Racimo, F., Allentoft, M.E., Jensen, M.K., Jonsson, A., Huang, H., Hormozdiari, F., Sikora, M., Marnetto, D., Eskin, E., et al. (2017). Selection in Europeans on fatty acid desaturases associated with dietary changes. Mol. Biol. Evol. *34*, 1307–1318.

Bulik-Sullivan, B.K., Loh, P.-R., Finucane, H.K., Ripke, S., Yang, J., Patterson, N., Daly, M.J., Price, A.L., and Neale, B.M.; Schizophrenia Working Group of the Psychiatric Genomics Consortium (2015). LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. Nat. Genet. *47*, 291–295.

Carithers, L.J., and Moore, H.M. (2015). The Genotype-Tissue Expression (GTEx) Project. Biopreserv. Biobank. *13*, 307–308.

Cheeran, M.C.J., Lokensgard, J.R., and Schleiss, M.R. (2009). Neuropathogenesis of congenital cytomegalovirus infection: disease mechanisms and prospects for intervention. Clin. Microbiol. Rev. *22*, 99–126.

Chen, J., Zheng, H., Bei, J.X., Sun, L., Jia, W.H., Li, T., Zhang, F., Seielstad, M., Zeng, Y.X., Zhang, X., and Liu, J. (2009). Genetic structure of the Han Chinese population revealed by genome-wide SNP variation. Am. J. Hum. Genet. *85*, 775–785.

Chen, Z., Chen, J., Collins, R., Guo, Y., Peto, R., Wu, F., and Li, L.; China Kadoorie Biobank (CKB) collaborative group (2011). China Kadoorie Biobank of 0.5 million people: survey methods, baseline characteristics and long-term follow-up. Int. J. Epidemiol. *40*, 1652–1666.

Chen, Y., Chen, Y., Shi, C., Huang, Z., Zhang, Y., Li, S., Li, Y., Ye, J., Yu, C., Li, Z., et al. (2018). SOAPnuke: A MapReduce acceleration-supported software for integrated quality control and preprocessing of high-throughput sequencing data. Gigascience 7, 1–6.

Davies, R.W., Flint, J., Myers, S., and Mott, R. (2016). Rapid genotype imputation from sequence without reference panels. Nat. Genet. *48*, 965–969.

DePristo, M.A., Banks, E., Poplin, R., Garimella, K.V., Maguire, J.R., Hartl, C., Philippakis, A.A., del Angel, G., Rivas, M.A., Hanna, M., et al. (2011). A framework for variation discovery and genotyping using next-generation DNA sequencing data. Nat. Genet. *43*, 491–498.

Francioli, L.C., Menelaou, A., Pulit, S.L., Van Dijk, F., Palamara, P.F., Elbers, C.C., Neerincx, P.B.T., Ye, K., Guryev, V., Kloosterman, W.P., et al.; Genome of the Netherlands Consortium (2014). Whole-genome sequence variation, population structure and demographic history of the Dutch population. Nat. Genet. *46*, 818–825.

Fumagalli, M. (2013). Assessing the effect of sequencing depth and sample size in population genetics inferences. PLoS ONE 8, e79667.

Fumagalli, M., Moltke, I., Grarup, N., Racimo, F., Bjerregaard, P., Jørgensen, M.E., Korneliussen, T.S., Gerbault, P., Skotte, L., Linneberg, A., et al. (2015). Greenlandic Inuit show genetic signatures of diet and climate adaptation. Science *349*, 1343–1347.

Galinsky, K.J., Bhatia, G., Loh, P.R., Georgiev, S., Mukherjee, S., Patterson, N.J., and Price, A.L. (2016). Fast principal-component analysis reveals convergent evolution of ADH1B in Europe and East Asia. Am. J. Hum. Genet. *98*, 456–472.

Gudbjartsson, D.F., Helgason, H., Gudjonsson, S.A., Zink, F., Oddson, A., Gylfason, A., Besenbacher, S., Magnusson, G., Halldorsson, B.V., Hjartarson, E., et al. (2015). Large-scale whole-genome sequencing of the Icelandic population. Nat. Genet. *47*, 435–444.

Huang, R.P., Ozawa, M., Kadomatsu, K., and Muramatsu, T. (1993). Embigin, a member of the immunoglobulin superfamily expressed in embryonic cells, enhances cell-substratum adhesion. Dev. Biol. *155*, 307–314.

Jiang, F., Ren, J., Chen, F., Zhou, Y., Xie, J., Dan, S., Su, Y., Xie, J., Yin, B., Su, W., et al. (2012). Noninvasive Fetal Trisomy (NIFTY) test: an advanced noninvasive prenatal diagnosis methodology for fetal autosomal and sex chromosomal aneuploidies. BMC Med. Genomics *5*, 57.

Karagoz, E., Ulcay, A., Tanoglu, A., Kara, M., Turhan, V., Erdem, H., Oncul, O., and Gorenek, L. (2014). Clinical usefulness of mean platelet volume and red blood cell distribution width to platelet ratio for predicting the severity of hepatic fibrosis in chronic hepatitis B virus patients. Eur. J. Gastroenterol. Hepatol. *26*, 1320–1324.

Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M., and Haussler, D. (2002). The human genome browser at UCSC. Genome Res. *12*, 996–1006.

Korneliussen, T.S., Albrechtsen, A., and Nielsen, R. (2014). ANGSD: analysis of next generation sequencing data. BMC Bioinformatics 15, 356.

Kothapalli, K.S.D., Ye, K., Gadgil, M.S., Carlson, S.E., O'Brien, K.O., Zhang, J.Y., Park, H.G., Ojukwu, K., Zou, J., Hyon, S.S., et al. (2016). Positive selection on a regulatory insertion-deletion polymorphism in FADS2 influences apparent endogenous synthesis of arachidonic acid. Mol. Biol. Evol. *33*, 1726–1739.

Landrum, M.J., Lee, J.M., Riley, G.R., Jang, W., Rubinstein, W.S., Church, D.M., and Maglott, D.R. (2014). ClinVar: public archive of relationships among sequence variation and human phenotype. Nucleic Acids Res. *42*, D980–D985.

Lee, Y.N., Malim, M.H., and Bieniasz, P.D. (2008). Hypermutation of an ancient human retrovirus by APOBEC3G. J. Virol. 82, 8762–8770.

Lek, M., Karczewski, K.J., Minikel, E.V., Samocha, K.E., Banks, E., Fennell, T., O'Donnell-Luria, A.H., Ware, J.S., Hill, A.J., Cummings, B.B., et al.; Exome Aggregation Consortium (2016). Analysis of protein-coding genetic variation in 60,706 humans. Nature *536*, 285–291.

Li, H. (2011). A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. Bioinformatics *27*, 2987–2993.

Li, H. (2015). FermiKit: assembly-based variant calling for Illumina resequencing data. Bioinformatics *31*, 3694–3696.

Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics 25, 1754–1760.

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., and Durbin, R.; 1000 Genome Project Data Processing Subgroup (2009). The Sequence Alignment/Map format and SAMtools. Bioinformatics *25*, 2078–2079.

Li, Y., Sidore, C., Kang, H.M., Boehnke, M., and Abecasis, G.R. (2011). Low-coverage sequencing: implications for design of complex trait association studies. Genome Res. *21*, 940–951.

Liang, Z., and Ma, Z. (2004). China's floating population: new evidence from the 2000 Census. Popul. Dev. Rev. *30*, 467–488.

Liang, Z., and White, M.J. (1996). Internal migration in China, 1950-1988. Demography *33*, 375–384.

Loh, P.R., Tucker, G., Bulik-Sullivan, B.K., Vilhjálmsson, B.J., Finucane, H.K., Salem, R.M., Chasman, D.I., Ridker, P.M., Neale, B.M., Berger, B., et al. (2015). Efficient Bayesian mixed-model analysis increases association power in large cohorts. Nat. Genet. *47*, 284–290.

Lonsdale, J., Thomas, J., Salvatore, M., Phillips, R., Lo, E., Shad, S., Hasz, R., Walters, G., Garcia, F., Young, N., et al.; GTEx Consortium (2013). The Geno-type-Tissue Expression (GTEx) project. Nat. Genet. *45*, 580–585.

MacArthur, J., Bowler, E., Cerezo, M., Gil, L., Hall, P., Hastings, E., Junkins, H., McMahon, A., Milano, A., Morales, J., et al. (2017). The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog). Nucleic Acids Res. *45* (D1), D896–D901.

Marchini, J., and Howie, B. (2010). Genotype imputation for genome-wide association studies. Nat. Rev. Genet. *11*, 499–511.

Maretty, L., Jensen, J.M., Petersen, B., Sibbesen, J.A., Liu, S., Villesen, P., Skov, L., Belling, K., Theil Have, C., Izarzugaza, J.M.G., et al. (2017).

Sequencing and de novo assembly of 150 genomes from Denmark as a population reference. Nature 548, 87–91.

Mathias, R.A., Fu, W., Akey, J.M., Ainsworth, H.C., Torgerson, D.G., Ruczinski, I., Sergeant, S., Barnes, K.C., and Chilton, F.H. (2012). Adaptive evolution of the FADS gene cluster within Africa. PLoS ONE *7*, e44926.

McLaren, W., Gil, L., Hunt, S.E., Riat, H.S., Ritchie, G.R.S., Thormann, A., Flicek, P., and Cunningham, F. (2016). The Ensembl variant effect predictor. Genome Biol. *17*, 122.

Moustafa, A., Xie, C., Kirkness, E., Biggs, W., Wong, E., Turpaz, Y., Bloom, K., Delwart, E., Nelson, K.E., Venter, J.C., and Telenti, A. (2017). The blood DNA virome in 8,000 humans. PLoS Pathog. *13*, e1006292.

Nkhoma, E.T., Poole, C., Vannappagari, V., Hall, S.A., and Beutler, E. (2009). The global prevalence of glucose-6-phosphate dehydrogenase deficiency: a systematic review and meta-analysis. Blood Cells Mol. Dis. *42*, 267–278.

Ohashi, J., Naka, I., and Tsuchiya, N. (2011). The impact of natural selection on an ABCC11 SNP determining earwax type. Mol. Biol. Evol. *28*, 849–857.

Pasaniuc, B., Rohland, N., McLaren, P.J., Garimella, K., Zaitlen, N., Li, H., Gupta, N., Neale, B.M., Daly, M.J., Sklar, P., et al. (2012). Extremely low-coverage sequencing and imputation increases power for genome-wide association studies. Nat. Genet. 44, 631–635.

Peter, B.M. (2016). Admixture, population structure, and f-statistics. Genetics 202, 1485–1501.

Pruim, R.J., Welch, R.P., Sanna, S., Teslovich, T.M., Chines, P.S., Gliedt, T.P., Boehnke, M., Abecasis, G.R., and Willer, C.J. (2010). LocusZoom: Regional visualization of genome-wide association scan results. Bioinformatics *26*, 2336–2337.

Raghavan, M., Skoglund, P., Graf, K.E., Metspalu, M., Albrechtsen, A., Moltke, I., Rasmussen, S., Stafford, T.W., Jr., Orlando, L., Metspalu, E., et al. (2014). Upper Palaeolithic Siberian genome reveals dual ancestry of Native Americans. Nature *505*, 87–91.

Rhoads, E.J.M. (2000). Manchus and Han: Ethnic Relations and Political Power in Late Qing and Early Republican China, 1861-1928 (University of Washington Press).

Sayers, E.W., Barrett, T., Benson, D.A., Bryant, S.H., Canese, K., Chetvernin, V., Church, D.M., DiCuccio, M., Edgar, R., Federhen, S., et al. (2009). Database resources of the National Center for Biotechnology Information. Nucleic Acids Res. *37*, D5–D15.

Skotte, L., Korneliussen, T.S., and Albrechtsen, A. (2012). Association testing for next-generation sequencing data using score statistics. Genet. Epidemiol. *36*, 430–437.

Soldin, O.P. (2006). Thyroid function testing in pregnancy and thyroid disease: trimester-specific reference intervals. Ther. Drug Monit. 28, 8–11.

Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., Danesh, J., Downey, P., Elliott, P., Green, J., Landray, M., et al. (2015). UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. PLoS Med. *12*, e1001779.

Suo, C., Xu, H., Khor, C.C., Ong, R.T.H., Sim, X., Chen, J., Tay, W.T., Sim, K.S., Zeng, Y.X., Zhang, X., et al. (2012). Natural positive selection and north-south genetic diversity in East Asia. Eur. J. Hum. Genet. 20, 102–110.

Vartanian, J.P., Henry, M., Marchio, A., Suspène, R., Aynaud, M.M., Guétard, D., Cervantes-Gonzalez, M., Battiston, C., Mazzaferro, V., Pineau, P., et al. (2010). Massive APOBEC3 editing of hepatitis B viral DNA in cirrhosis. PLoS Pathog. *6*, e1000928.

Vourekas, A., Zheng, K., Fu, Q., Maragkakis, M., Alexiou, P., Ma, J., Pillai, R.S., Mourelatos, Z., and Wang, P.J. (2015). The RNA helicase MOV10L1 binds piRNA precursors to initiate piRNA processing. Genes Dev. 29, 617–629.

Walter, K., Min, J.L., Huang, J., Crooks, L., Memari, Y., McCarthy, S., Perry, J.R.B., Xu, C., Futema, M., Lawson, D., et al.; UK10K Consortium (2015). The UK10K project identifies rare variants in health and disease. Nature *526*, 82–90.

Welter, D., MacArthur, J., Morales, J., Burdett, T., Hall, P., Junkins, H., Klemm, A., Flicek, P., Manolio, T., Hindorff, L., and Parkinson, H. (2014). The NHGRI

GWAS Catalog, a curated resource of SNP-trait associations. Nucleic Acids Res. 42, D1001–D1006.

Wigginton, J.E., Cutler, D.J., and Abecasis, G.R. (2005). A note on exact tests of Hardy-Weinberg equilibrium. Am. J. Hum. Genet. *76*, 887–893.

Wildschutte, J.H., Williams, Z.H., Montesion, M., Subramanian, R.P., Kidd, J.M., and Coffin, J.M. (2016). Discovery of unfixed endogenous retrovirus insertions in diverse human populations. Proc. Natl. Acad. Sci. USA *113*, E2326–E2334.

Wood, A.R., Esko, T., Yang, J., Vedantam, S., Pers, T.H., Gustafsson, S., Chu, A.Y., Estrada, K., Luan, J., Kutalik, Z., et al.; Electronic Medical Records and Genomics (eMEMERGEGE) Consortium; MIGen Consortium; PAGEGE Consortium; LifeLines Cohort Study (2014). Defining the role of common variation in the genomic and biological architecture of adult human height. Nat. Genet. 46, 1173–1186.

Xu, S., Yin, X., Li, S., Jin, W., Lou, H., Yang, L., Gong, X., Wang, H., Shen, Y., Pan, X., et al. (2009). Genomic dissection of population substructure of Han Chinese and its implication in association studies. Am. J. Hum. Genet. *85*, 762–774.

Yan, Y.P., Su, H.X., Ji, Z.H., Shao, Z.J., and Pu, Z.S. (2014). Epidemiology of hepatitis B virus infection in China: current status and challenges. J. Clin. Transl. Hepatol. 2, 15–22.

Yang, L., Tan, S., Yu, H., Zheng, B., Qiao, E., Dong, Y., Zan, R., and Xiao, C. (2008). Gene admixture in ethnic populations in upper part of Silk Road revealed by mtDNA polymorphism. Sci. China C Life Sci. *51*, 435–444.

Yang, J., Jin, Z.-B., Chen, J., Huang, X.-F., Li, X.-M., Liang, Y.-B., Mao, J.-Y., Chen, X., Zheng, Z., Bakshi, A., et al. (2017). Genetic signatures of high-altitude adaptation in Tibetans. Proc. Natl. Acad. Sci. USA *114*, 4189–4194.

Yengo, L., Sidorenko, J., Kemper, K.E., Zheng, Z., Wood, A.R., Weedon, M.N., Frayling, T.M., Hirschhorn, J., Yang, J., and Visscher, P.M.; GIANT Consortium (2018). Meta-analysis of genome-wide association studies for height and body mass index in ~700,000 individuals of European ancestry. bioRxiv. https://doi. org/10.1101/274654.

Zhang, H., Gao, Y., Jiang, F., Fu, M., Yuan, Y., Guo, Y., Zhu, Z., Lin, M., Liu, Q., Tian, Z., et al. (2015). Non-invasive prenatal testing for trisomies 21, 18 and 13: clinical experience from 146,958 pregnancies. Ultrasound Obstet. Gynecol. *45*, 530–538.

Zou, H., Chen, Y., Duan, Z., Zhang, H., and Pan, C. (2012). Virologic factors associated with failure to passive-active immunoprophylaxis in infants born to HBsAg-positive mothers. J. Viral Hepat. *19*, e18–e25.

Zwolińska, K., Knysz, B., Gąsiorowski, J., Pazgan-Simon, M., Gładysz, A., Sobczyński, M., and Piasecki, E. (2013). Frequency of human endogenous retroviral sequences (HERV) K113 and K115 in the Polish population, and their effect on HIV infection. PLoS ONE *8*, e77820.

STAR*METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Chemicals, Peptides, and Recombinant Proteins		
T4 DNA Polymerase	Enzymatics	Cat#P708L
T4 Polynucleotide Kinase	Enzymatics	Cat#Y904L
Klenow Fragment	Enzymatics	Cat#P706L
dNTP solution set	Enzymatics	Cat#N201L
DATP(2-Deoxyadenosine 5'-Triphosphate)	GE	Cat#28406501
Klenow (3'-5'exo-)	Enzymatics	Cat#P701-LC-L
T4 DNA Ligase (Rapid)	Enzymatics	Cat#L603-HC-L
MgCl2	MILLIPORE	Cat#20-303
PLATINUM PFX DNA POLYMERASE	Invitrogen	Cat#C11708-021
Dimethyl sulfoxide(DMSO)	sigma	Cat#D2650-100ML
Critical Commercial Assays		
TIANamp Micro DNA Kit	TIANGEN	Cat#DP316
TruSeq Rapid SBS Kit HS (200 cycle)	ILLUMINA	Cat#FC-402-4001-2
Hiseq 2500 SR Flowcell	ILLUMINA	Cat#GD-402-4001-PC
TruSeq Rapid SR Cluster Kit HS	ILLUMINA	Cat#GD-402-4001-1
10X T4 PNK Reaction Buffer	Enzymatics	Cat#Y904C
10X Blue Buffer	Enzymatics	Cat#B011L
10X Ligation Buffer	Enzymatics	Cat#B603L
2X Rapid Ligation Buffer	Enzymatics	Cat#B101L
Buffer EB(250ML)	QIAGEN	Cat#19086
Agencourt AMPure XP-Medium	AGENCOURT	Cat#A63882
Agilent DNA 1000 Reagents	Agilent	Cat#5067-1504
TaKaRa Ex Taq Hot Start Version	TaKaRa	Cat#DRR006B
ThermoPol buffer	NEB	Cat#B9004S
TWEEN 20,For molecular biology	sigma	Cat#P9416
EvaGreenTM, 20X	BIOTIUM	Cat#31000
ROX REFERENCE DYE	INVITROGEN	Cat#12223-012
Deposited Data		
hg19 human genome reference	Kent et al.,2002	http://genome.ucsc.edu/
sgdp-263-hs37d5	Li, 2015	https://github.com/lh3/sgdp-fermi/releases/download/v1/ sgdp-263-hs37d5.tgz
wgEncodeDukeMapabilityUniqueness35bp	Kent et al., 2002	http://hgdownload.cse.ucsc.edu/goldenpath/hg19/ encodeDCC/wgEncodeMapability/
Clinvar database	Landrum et al., 2014	https://www.ncbi.nlm.nih.gov/clinvar/
Giant summary statistics for height	Yengo et al., 2018	https://portals.broadinstitute.org/collaboration/giant/ images/0/0f/Meta-analysis_Locke_et_al%2BUKBiobank_ 2018.txt.gz
Giant summary statistics for BMI	Yengo et al., 2018	https://portals.broadinstitute.org/collaboration/giant/ images/6/63/Meta-analysis_Wood_et_al%2BUKBiobank_ 2018.txt.gz
UK biobank summary statistics for height	(Ben Neale's website)	https://www.dropbox.com/s/sbfgb6qd5i4cxku/50. assoc.tsv.gz?dl=0
UK biobank summary statistics for BMI	(Ben Neale's website)	https://www.dropbox.com/s/sweqn7nztyv42zt/21001. assoc.tsv.qz?dl=0

(Continued on next page)

Continued				
REAGENT or RESOURCE	SOURCE	IDENTIFIER		
GWAS catalog database	Welter et al., 2014	ftp://ftp.ebi.ac.uk/pub/databases/gwas/releases/2017/ 02/21/gwas-catalog-associations.tsv		
NCBI viral reference sequence database	Sayers et al., 2009	ftp.ncbi.nih.gov/refseq/release/RefSeq-release84/viral/ viral.2.1.genomic.fna.gz		
Software and Algorithms				
SOAPnuke	Chen et al., 2018	https://github.com/BGI-flexlab/SOAPnuke		
bwa	Li and Durbin, 2009	https://sourceforge.net/projects/bio-bwa/files/		
samtools	Li et al., 2009	http://samtools.sourceforge.net/		
GATK 3.8	DePristo et al., 2011	https://github.com/broadgsa/gatk/releases		
Variant Effect Predictor	McLaren et al., 2016	https://asia.ensembl.org/info/docs/tools/vep/script/ vep_download.html		
STITCH v1.2.7	Davies et al., 2016	http://www.well.ox.ac.uk/~rwdavies/stitch_2017_ 02_14.html		
RSpectral	The R Foundation	https://cran.r-project.org/web/packages/RSpectra/ index.html		
angsd	Korneliussen et al., 2014	http://www.popgen.dk/angsd/index.php/ANGSD		
SNPtest	Marchini and Howie, 2010	https://github.com/ANGSD/angsd		
locuszoom	Pruim et al., 2010	http://locuszoom.sph.umich.edu		
liftover	Kent et al., 2002	http://hgdownload.cse.ucsc.edu/admin/exe/ linux.x86_64/liftOver		
BLAST	Sayers et al., 2009	https://blast.ncbi.nlm.nih.gov/Blast.cgi		

CONTACT FOR REAGENT AND RESOURCE SHARING

Further information and requests for resources and reagents should be directed to and will be fulfilled by the Lead Contact, Xun Xu (xuxun@genomics.cn).

EXPERIMENTAL MODEL AND SUBJECT DETAILS

All participants were recruited via the non-invasive fetal trisomy test at BGI between year 2012 and 2013 (Zhang et al., 2015). They underwent pretest counseling and filled in informed written consent before blood sampling. The study was reviewed and approved by the Institutional Review Board of BGI (BGI-IRB17088) in strict compliance with regulations regarding ethical considerations and personal data protection.

The geographic locations of the participants were provided in an anonymous way via the first six digits of the resident identity card number, indicating birth place information of the participant. Year of birth and the age of the participants were obtained from the written informed consent. The age distribution of 137, 984 participants suggests a bimodal distribution with peak for birthyears of 1977 and 1984 and ages 28 and 35 (Figure S4K-L). The bimodal distribution is likely related to China's first and second child policy and affected by changes to the clinical guidelines for undertaking the NIPT test, where all pregnant women (rather than only high-risk pregnant women) are recommended to take the NIPT test. The height and weight measurements were recorded when the blood sample was taken in the hospital. The BMI is calculated using the standard formula "weight (kg) / height2 (m2)."

Relative fetal fraction can be estimated accurately for a participant if the gender of the child is male based on the proportion of reads that map to Y chromosome relative to the reads that map to the whole genome. Using this proportion, we estimate that the fetal fraction is approximately 3.5% to 30% among all participants, with a median of 8%.

The status of chromosomal aneuploidy was detected using a method we previously developed for screening of fetal chromosomal aneuploidy (Jiang et al., 2012). In addition, both the participant, and if available, the father, reported their karyotype status to the hospital. We removed participants with sequencing error rate greater than 1% (n = 3,278) and with potential abnormal chromosomal aneuploidy detected either via the participant's report or detected using the read coverage method (n = 502) from further analyses, resulting in 141, 431 participants consisting of 118, 576 participants with 35bp read length and 22, 855 participants with 49bp read length.

METHOD DETAILS

Sequencing and QC

Details of the sequencing protocol were published previously in Zhang et al. (2015). In brief, within 8h of blood collection, plasma was extracted from whole blood after two turns of centrifugation. The plasma samples (n = 145, 211) were subsequently subjected to library construction, sample quality control and 36-cycle or 50-cycle single-end multiplex sequencing on Illumina Hiseq 2000 platform. The reads were trimmed to 35bp and 49bp before bioinformatic analysis. Filtering of poor quality reads was carried out using SOAPnuke (https://github.com/BGI-flexlab/SOAPnuke). A read was removed if it contained more than 30% low quality bases (Q \leq 2) or N bases. In general, each participant was whole-genome sequenced to 5-10 million cleaned reads, representing a sequencing depth around 0.06x -0.1x.

Medium depth sequencing of 40 participants

With informed consent, we sequenced 40 participants out of the total 141, 431 participants using the Hiseq X10 system at a medium depth of 15x. We aligned the reads to the same hg19 human reference genome using bwa-mem (Li and Durbin, 2009) and applied the GATK multi-sample best practice (DePristo et al., 2011) to call and genotype SNPs for the 40 participants. The SNP calls and genotype results were used to benchmark the SNP calling performance using the ultra-low depth sequencing data as well as the genotype imputation accuracy.

Alignment of reads against hg19 and definition of accessible region

For each participant, the cleaned reads were aligned against the hg19 human genome reference with bwa, using the single-end read alignment option (Li and Durbin, 2009). Potential PCR duplicates were removed using samtools rmdup (Li et al., 2009). The indel realignment and base quality recalibration modules in GATK were applied to realign the reads around indel candidate loci and to recalibrate the base quality (DePristo et al., 2011). Finally, the alignment files were stored in the standard CRAM format. The alignment of all participants was carried out in the Batch Compute system in Aliyun cloud parallelizing 2000 jobs in a batch (~24 hours per job) (https://github.com/aliyun/aliyun-openapi-java-sdk/tree/master/aliyun-java-sdk-batchcompute). Completion of the whole alignment process took around one week.

After alignment, QC statistics were computed using the stats function implemented in samtools which measures sequencing error rate as the proportion of bases that differ from reference base at base positions, i.e., the mismatch rate. The median sequencing error rate was estimated to be 0.3%. Subsequently, we used the samtools depth to estimate the overall coverage of reads with mapping quality > = 30 and bases with base quality > = 20. Reads or bases with quality lower than this threshold were not included in any of the following analysis. We compared the coverage to the mappability uniqueness (wgEncodeDukeMapabilityUniqueness35bp.bedGraph) as well as gene and repeat density information from the UCSC database (Kent et al., 2002) (Figure S1D).

Considering potential errors in alignment of short single-end reads, we defined an accessible region for variation calling, population genetics and association mapping analysis. We defined the accessible region as follows: 1) regions that are not in the 35-kmer problematic alignment bed file provided by Heng Li (https://github.com/lh3/sgdp-fermi/releases/download/v1/sgdp-263-hs37d5.tgz). Those regions in the hg19 reference genome were detected by Fermi assembler (Li, 2015) as difficult to map uniquely using a 35bp kmer unit sequence. This filter removed 897,319,085 bp hg19 sequence; 2) regions with a mappability uniqueness score equal to one according to wgEncodeDukeMapabilityUniqueness35bp.bedGraph in UCSC. This filter additionally removed 6,459,019 bp hg19 sequence; 3) regions where the sequencing depth was between 3000 and 30,000. An additional 4,340,936 bp are excluded in this step. The final length of accessible region in hg19, including the 22 autosomal chromosomes and the X chromosome is 2,128,184,806 bp.

Variant discovery and allele frequency estimation

We applied a maximum likelihood approach to identify polymorphic sites and infer allele frequencies (S. Liu, unpublished data). We adopted a single-read sampling strategy, where we sampled only one read for each variant candidate site if there were more than one read. This method ensured that, for each site, all reads derived from different individuals and the allele frequency spectrum is therefore not biased due to existence of fetal DNA in the sample. This maximum likelihood framework is much faster than a representation using diploid genotype likelihoods (DePristo et al., 2011; Li, 2011). The mathematical derivation has been detailed in the "QUANTIFICATION AND STATISTICAL ANALYSIS" section in the manuscript.

Annotation

Annotation of the genes mentioned in the manuscript and the annotation of the existence of the variants in database such as dbSNP, GnomeAD, 1KGP was carried out using Variant Effect Predictor (McLaren et al., 2016).

Imputation

We employed STITCH (version 1.2.7) (Davies et al., 2016) to impute genotype probabilities for all 141,431 individuals in a five megabase window with a 250K buffer assuming 10 ancestral haplotypes. The key parameter K (number of ancestral haplotypes) was decided based on tests over a 5Mbp region in chr3 (chr3: 18000000-185000000) via useful discussions with the STITCH author Robert W Davies. Allele frequency information from the Chinese population (CHB+CHS+CDX, n = 301) in the 1KG impute2 reference panel was used for the initial values for the EM optimization of the model parameters. 607 jobs were parallelized in the Tianhe 2 supercomputer in Guangzhou city.

The imputed loci are a target of 8.16 million known polymorphic sites in 22 autosomal chromosomes and chrX with a 1KG East Asian allele frequency > = 0.01. All the loci recorded in the GWAS catalog are also included for imputation.

For each of the imputed site, there is an IMPUTE2-style info score (Marchini and Howie, 2010) and a P value for violation of Hardy Weinberg equilibrium (HWE-pvalue in short) (Wigginton et al., 2005). We used info score greater than 0.4 and HWE-pvalue smaller than 10^{-6} since the remaining variants has greatest power and good replication rate for height and BMI association test.

Principal component analysis

For the PC analysis, we restricted ourselves to known variants with minor allele frequency greater than 0.05 in the data. We continued to use the single read sampling strategy. To compute the covariance between individuals *i* and *j* for *M* loci, we used

$$C_{ij} = \frac{1}{M} \sum_{m=1}^{M} \frac{(h_m^i - f_m) (h_m^j - f_m)}{f_m (1 - f_m)}$$

where f_m is the minor allele frequency and h_m^i is the haploid genotype coded as either 0 or 1 for the major and minor allele, respectively.

Using the above formulas, we parallelized the distance and the covariance matrix computation in 90 nodes from the Aliyun ODPS system and obtained the full matrix of 141,431 individuals in a few days. We applied the Spectra R package (https://cran.r-project. org/web/packages/RSpectra/index.html) to perform the decomposition of the covariance matrix.

Finally, we visualized the top three principal components and colored the points according to the administrative divisions, the ethnic groups, and also the read length and error rate.

We have carried out several principal component analyses to answer different questions using the workflow described above. First, we carried out a PCA for all 141, 431 participants to identify the main principal components (Figures S3A–S3D). After noticing that the first principal component reflects read length and the third reflects the estimated sequencing error rate (measured using the stat function in samtools), we only used the 96,880 participants with read lengths of 35bp and sequencing error rate smaller than 0.00325 for further population genetic analysis (Figure 2). In particular, for the PCA in the Han Chinese population (Figures S3E and S3F), we excluded participants that did not report their ethnicity (although a large majority of them are probably Han). We only used the 45, 387 participants who reported Han ethnicity to understand population structure of the Han.

F3-statistic and private allele frequency analysis

To quantify divergence between populations we use the outgroup F3 statistic, which is a measure of drift time between two populations (Raghavan et al., 2014). The F3 statistic is highly influenced by common alleles, that tend to be older, and we noticed that F3 statistics between CEU and ITU 1KG samples, and samples from each of the Chinese administrative divisions, were highly correlated due to the sharing of ancestry between the CEU and ITU populations after their separation from East Asian populations (Figure not shown). We therefore also measured genetic relatedness using a measure based on alleles that are private to either the African (YRI), east Asian (CHB), South Asian (ITU) or Europeans (CEU) 1KG sample. Private alleles are defined as those that were polymorphic in one group and fixed in the other groups. There were in total 3, 485, 371 and 4, 324,376 private alleles in the CEU and ITU samples respectively. We further applied a filter, using the private alleles that were common in one group with a MAF > 5% and obtained in 66,700 and 45,536 private variants in the CEU and ITU samples, respectively. Those common private variants were used to compute the private allele frequency defined below. We assume that the proportion of allele sharing of these private alleles with any of the administrative division, should be informative regarding genetic exchange at a more recent timescale.

For each administrative division we calculated the fraction of alleles in the NIPT dataset that matched the private allele found in population K in the 1KG. We denote the number of sites that contains an allele that is private to population, K, as M_k . For each site, s, that contains a private allele for population K, we count the number of alleles that match the private allele, n_s , and normalize by the total number of alleles N_s . The private allele frequency for population K is defined as

$$PAF_{k} = \frac{\sum_{s=1}^{M_{k}} n_{s}}{\sum_{s=1}^{M_{k}} N_{s}}$$

Standard errors were estimated using a 5Mb weighted block jackknife where the weights are the number of sites with private alleles within the block.

Detection of selection across PC coordinates

To detect the most extremely differentiated variants, we use a method based on finding deviations from the patterns predicted using the first components of a PCA analysis. We have adapted the FastPCA statistic (Galinsky et al., 2016) to work on covariance matrix

based on PCA (using the NIPT sequencing data). Assuming the eigenvectors obtained by PCA capture the structure in the data in the absence of selection, for each SNP, we use Equation (11) in Galinsky et al. (2016) to calculate a p value associated with deviations from the genomic pattern. The resulting p values were visualized using a Manhattan plot.

Detection of clinvar pathogenic variants displaying significant allele frequency differentiation

We investigated allele frequency differences among the three populations defined by the three latitudinally separated geographic divisions in a total of 3,238 bi-allelic potentially pathogenic variants with a clinical significance level of 5 from the total of 246,385 variants in clinvar database (Landrum et al., 2014) (URL: ftp.ncbi.nlm.nih.gov/pub/clinvar/vcf_GRCh37/clinvar_20170404.vcf.gz). We note that even for the pathogenic variants with clinical significance level of 5, some of them have high allele frequency greater than 1% in both our dataset and the 1KG population. We counted the number of risk and non-risk alleles (B) for each of the North, South and Central populations using the formula

$$C_b = \sum_{i=1}^{N_b} \left(1 - 10^{\frac{-q_b}{10}} \right)$$

$$b \in \{A, C, G, T\}$$

where C_b refers to total allele count of b, q_b refers to the base quality of the observed base b. For each variant, a two-tailed p value using Fisher's exact test was calculated for the 2x3 table with data from the 2 alleles and 3 regions.

For each of the pathogenic variants, as well as for all other variants throughout the genome, we estimated allele frequencies using BaseVar. The risk allele frequency of position *j* is defined as R_{j} . We frequency-matched random SNPs from the whole genome data with SNPs in clinvar by, for SNP *j* in clinvar, randomly selecting 100,000 variants from the non-clinvar variants within a frequency range of $[0.9R_j, 1.1R_j]$. We estimated the rank and percentile of the pathogenic variant comparing to the 100,000 variants. We reported and visualized the top eight loci with a p value retrieved from the Fisher's exact test was less than 10e-6, and the p value from the comparison to non-clinvar variants was less than 5e-3.

Identification of genetic variants significantly associated with a trait

We used the score test (Korneliussen et al., 2014) implemented in Angsd (Korneliussen et al., 2014) to detect the association signal between the imputed genotype probabilities and phenotypes, followed by a linear regression for quantitative trait or logistic regression for qualitative trait to compute the effect size of the top SNPs. For height, we applied the top five principal components, the maternal age and the sex of the fetus as covariates. For BMI, we additionally included the gestational age of the fetus as a covariate. We note that the BMI phenotype in the NIPT cohort is not only related to the mothers' non gestational weight but also related to gestational weight gain including the fetus's weight. We are mostly interested in the genetic effects on the maternal phenotype in this study and therefore, we have used the gestational age and the sex of the fetus as covariates to account for the fetal growth rate and the effect of sex. Even so, there may be some residual variance in the BMI phenotype caused by differences in fetal growth rate. For maternal age, we used the top five principal components and the sex of the fetus as covariates. For the rest of the phenotypes including twins and virus integration and infection, we applied the same covariates as height. Independent loci were defined as significant variants clustered in a 1Mbp window. The lead SNP was defined as the SNP in the 1Mbp window that has most significant, i.e., smallest p value. The conditional test module in SNPtest (Marchini and Howie, 2010) was used to estimate the number of independent signals for each independent loci. Finally, locuszoom (Pruim et al., 2010) was applied to visualize the loci. The reported loci were determined from the conditional test after the single marker analysis using a significance threshold P value $\leq 5 \times 10^{-8}$.

The genomic inflation factor, GC lambda, attenuation ratio, LD score regression intercept and the SNP heritability were estimated using the LD score regression approach (Bulik-Sullivan et al., 2015).

Replication of significant loci

For replication, we compared all the variants reaching the significance threshold to three independent studies- the China Kadoorie Biobank (CKB cohort) (Chen et al., 2011), the recent Giant meta-analysis (Yengo et al., 2018, bioRxiv) and the UK Biobank (Ben Neales's website). The CKB cohort has measurements of height and BMI data and has chip-genotyped 32,000 Chinese participants with imputation into the 1000 genomes Phase reference panel. The GC lambdas from the CKB association test using BOLT-LMM (Loh et al., 2015) were 1.10 and 1.17, respectively, for height and BMI. The GIANT and the UK Biobank summary statistics consist of 2.3 million and 10 million SNP markers, respectively. For some associated loci, the lead SNP is not present in the test. For replication purpose, after ascendingly ranking the SNPs by p value, we chose the first SNP present in the test data as the proxy SNP. In almost all cases the p values and effect sizes of the lead SNPs are similar to the p value of the proxy SNPs. When proxy SNP instead of the lead SNP was used for the replication, we marked it cleary in the result.

We defined a locus as replicated if the lead SNP or the proxy SNP 1) has a p value less than 0.05 divided by the number of loci (for height, n = 48; for BMI, n = 13) 2) has the same directionality of the effect, in at least one of the CKB, the GIANT or the UK Biobank test set. We note that the genomic inflation factor is high in both the GIANT and the UK Biobank (qq plot not shown).

The GWAS catalog database (Welter et al., 2014, e87_r2017-02-20) defines known and novel loci. The b38 coordinates was transferred to the b37 coordinates using the liftover script from UCSC (http://hgdownload.cse.ucsc.edu/admin/exe/linux.x86_64/liftOver).

Viral sequence analysis

We applied BLAST (Altschul et al., 1990) to align the reads that did not map to the human reference genome (hg19+human decoy sequences) to the NCBI viral reference sequence database (Sayers et al., 2009) (ftp.ncbi.nih.gov/refseq/release/RefSeq-release84/viral/viral.2.1.genomic.fna.gz). For each read, we kept the best alignment with smallest e-value. Only reads with an evalue < 1e-5, identity > = 97 and alignment length > = 32bp were counted as a hit. After removing alignments to bacteriophages, we found that out of the 138,882 samples analyzed, 48,298 samples (34.8%) have at least one viral hit and 11,351 samples (8.2%) have at least two significant hits mapped to the viral reference database. In Figure 5, to reduce false positive, we defined individuals with at least two significant hits aligned to the same virus as carriers of that virus. We then carried out prevalence and abundance analysis of each virus using the top viral hit for each read. For prevalence and abundance analysis, we only used the individuals with at least two hits for a specific virus. Virus abundance was calculated by the following equation adapted from Moustafa et al. (2017):

 $Abundance = \frac{2x \frac{number of reads mapped to viral genome}{virus genome size}}{\frac{number of reads mapped to human genome}{human genome size}}$

Viruses with multiple strain entries in RefSeq were aggregated for high homology between entries and ease of graphical display. These viruses include: Anellovirus (TTV), HBV, HHV-6A, HHV-6B, HHV-5, Influenza, etc. For virus sequencing coverage analysis, we aggregated the read depth of all the individuals with mapping quality > 0 for each virus.

QUANTIFICATION AND STATISTICAL ANALYSIS

Maximum likelihood estimation of allele frequency *Likelihood Function for a single site*

For *N* unrelated individuals with a single read covering the position, the likelihood function for the read data D_i , for a single variant candidate site in individual *i*, of the allele frequency $p = (p_A, p_C, p_G, p_T)$, is defined as:

$$L(p) = \prod_{i=1}^{N} P(D_i | p) = \prod_{i=1}^{N} \sum_{b \in \{A,C,G,T\}} p(b | p) p(D_i | b)$$
(1)

where $p(b|p) = p_b$ and the genotype likelihood assuming a haploid model is $p(D_i|b) = \{1 - \varepsilon_i \text{ if } D_i = b \text{ and } \varepsilon_i/3, \text{if } D_i \neq b. \varepsilon_i \text{ corresponds to the GATK corresponds to the GATK-recalibrated error rate converted from the PHRED-scale base quality.$ **Optimization**

We obtain the maximum likelihood estimate $\hat{p} = argmax_p L(p)$ using the EM algorithm with starting value computed by the observed allele frequency:

$$\rho_b = \frac{\sum D_i = b}{N} \tag{2}$$

In the E step, we compute the posterior probability of allele b for individual i at a site j as one of the four A/C/G/T bases:

$$P(b | D_i) = \frac{p(b | p)p(D_i | b)}{\sum_{b' \in \{A, C, G, T\}} p(b' | p)p(D_i | b'))}$$
(3)

We compute the updated allele frequency p' in the M step as

$$p'_{b} = \frac{\sum_{i=1}^{N} P(b \mid D_{i})}{N}$$
(4)

When the change in the maximum likelihood is less than 0.001, we terminate the algorithm.

Decision of allelic type and confidence of SNP calling: Likelihood Ratio Test

Formulae (1) – (4) can be used for estimation of allele frequencies of all four nucleotides simultaneously, and may result in tetra-allelic and tri-allelic variant calls. We will use this formulation for SNP calling and for identifying potential tri- and tetra-allelic loci. Denote the

likelihood value from the four-allelic model in Equation (1) as f_4 . We iteratively set the allele frequency of one of the four nucleotides to zero to obtain models of tri-allelic loci. Let $\hat{f_3}(p_x = 0)$ denote the maximum likelihood value when the frequency of allele *x* is constrained to be zero. We then compute a log likelihood ratio statistic as:

$$LRT_{4vs3} = -2log\left(\frac{\widehat{f}_{3}(p_{x}=0)}{\widehat{f}_{4}}\right)$$
(5)

The tri-allelic model is nested within the tetra-allelic model and, therefore, the distribution of the LRT_{4vs3} statistic asymptotically follows a chi-square distribution with one degree of freedom, under the assumption of a tri-allelic locus. If the p values of one of the four LRT_{4vs3} test are significant (< 10⁻⁶), the variant will be classified as a tetra-allelic loci. If not, we move on to the test a model of a tri-allelic locus versus a bi-allelic locus, where x if the $\hat{f}_3(p_x = 0)$ is the allele with minimum likelihood (which results in maximum p value out of LRT_{4vs3}) was set as the alternative-hypothesis and the reduced hypothesis is $\hat{f}_2(p_x = 0, p_y = 0)$ where p_y is the allele frequency for allele y.

$$LRT_{3vs2} = -2log\left(\frac{\widehat{f}_{2}(p_{x}=0,p_{y}=0)}{\widehat{f}_{3}(p_{x}=0)}\right)$$
(6)

Again, the distribution of LRT_{3vs2} asymptotically follows a chi-square distribution with one degree of freedom under the hypothesis of a bi-allelic locus. If the maximum p value out of the three LRT_{3vs2} is significant, the variant will be classified as a tri-allelic variant. Otherwise, we continue to test the bi-allelic versus mono-allelic assumption, as defined in the equation below, with *y* being the allele with the highest p value

$$LRT_{2vs1} = -2log\left(\frac{\widehat{f_1}(\rho_x = 0, \rho_y = 0, \rho_z = 0)}{\widehat{f_2}(\rho_x = 0, \rho_y = 0)}\right)$$
(7)

Formula (8) is also used to quantify the confidence of the SNP call. We keep variants with p values less than 10⁻⁶.

Note that our method identifies multi-allelic variants. However, since we don't have sufficient validation of the performance for such variants, we focus on reporting results for bi-allelic loci.

Variant quality score recalibration

 \sim 32 million raw variants were obtained using a p value less than 10⁻⁶ based on the maximum likelihood model in the accessible region (Table S1). However, this set of SNPs may contain false positives due to miscalculated quality scores, alignment errors, or other technical issues such as contamination. We therefore applied a Bayesian Gaussian mixture model, similar to the VQSR model in GATK (DePristo et al., 2011) to assign each variant candidate a Phred-scaled probabilistic score (in short, VQSR score) indicating the probability that the variant is a truly polymorphic variant. The higher the VQSR score, the higher probability that the variant candidate is a true polymorphic variant. The Gaussian mixture model was established by learning technical features of a training set that consists of variant likely to be real. In our case, the training set was defined as a subset of the common known variants (n = 50K was randomly chosen). Features include the Fisher exact test statistic for strand bias, sequencing depth, indel density in a 30bp window centered around the variant candidate, and the raw variant quality score using a maximum likelihood model described above. The same likelihood function and expectation and maximization process as that reported in the GATK framework (DePristo et al., 2011) was implemented except that we used prior probability 50% and 50% for all the variants.

The transition versus transversion (Ti/Tv) ratio is high for the raw call set (maximum 3.4 for known variants and 8.9 for novel variants) and decreases as the filtration threshold of VQSR score increases. The final filtration threshold of VQSR score is decided to be 35, which suggests a Ti/Tv ratio of 2.2 for known variants and of 2.4 for novel variants and a sensitivity of 85% for the common known variants used for training (Figure S2A).

DATA AND SOFTWARE AVAILABILITY

Any uploading and sharing of individual genetic data from this project is not allowable according to a review by the Human Genetic Resources Administration of China (HGRAC) based on the regulations documented in the Interim Measures for the Administration of Human Genetic Resources.

However, we have made the maximal efforts to ensure that we can make as detailed as possible summaries of the data available to other researchers, including allele frequencies and GWAS summary statistics. The summary statistics are available at https://db.cngb.org/cmdb/. Researchers who wish to gain access to the data are required to fill in a simple application form and send an email to the bigdata@genomics.cn. After the applicant's identity is verified, they will be given the account and password to access the allele frequency information and other summary statistics data. This process takes no longer than five weekdays. This verification process is necessary in order to adhere to the Chinese regulations.

Supplemental Figures



Figure S1. Geographic Distribution of 140,000 Participants over 34 Administrative Divisions in China and Sequencing Depth, Related to Figure 1, Table S1, and the STAR Methods

(A) Geographic distribution of 140K participants over 34 administrative divisions in China. Color and number in legend indicate the number of participants. Participants from three divisions including Hong Kong, Macau and Taiwan are not involved in the study.

(B) Integrated sequencing depth of all the study participants over accessible and unaccessible regions.

(C) Sequencing depth per individual by chromosomes.

(D) Integrated sequencing depth over the 22 autosomal and the X chromosomes. Sequencing depth for participants with 35bp read length, 49bp read length and a summation of both groups are shown in green, blue and red, respectively.



Figure S2. Quality Measurement of Variation Calling and Imputation Accuracy, Related to Figure 1, Table S1, and the STAR Methods (A) Changes in Ti/Tv ratio and the ratio of the remained variants as a function of the increasing filtration threshold on variant recalibration score. In the legend, "positive" and "negative" refers to the status of variants selected as the positive set and negative set in the training process. "All" refers to the status of all variants.

(B-E) Upper bound false positive counts using variants with > = 2 alleles in 40WGS validation data. "All," "Known" and "Novel" variants refers to all the variants, variants known in dbsnp147 and variants not known in dbsnp147. False calls are those that are called by BaseVar from the 140K NIPT data and have at least two alternative alleles in the low coverage NIPT sequencing data of 40 validation participants but are not identified as variants in the high coverage sequencing data of those same individuals.

(F) Imputation accuracy for all the variants per individual as a function of fetal fraction.



Figure S3. Principal Component Analysis and Sharing of Ancestry with Reference Populations, Related to Figure 2 and STAR Methods (A–D) Principal component analysis for the full sample set (n = 141, 431). (A) Distribution of eigenvalue for the top eight principal components. (B). PCA colored by read length suggests the first principal component is the batch effect caused by differential mapping of different read length. (C). PCA colored by three latitudinal geographical regions (North, Central and South) suggests the second principal component reflects genetic differentiation across latitude. (D). PCA colored by sequencing error rate suggests the third principal component likely corresponds to the error rate.

(E and F) Principal component analysis of the self-reported Han Chinese (n = 43, 387). (E) Geographic locations of the Han. Divisions belong to the "North," "Central" and "South" regions defined by the Chinese government are colored as ""Green," "Red" and "Blue." (F) Visualization of the top two principal components. Colors correspond to their digital geographical information in Panel A. Solid and dashed ellipses line refers to 95% confidence interval of the PC distributions for individuals from the three geographical regions assuming a multivariate t-distribution or a multivariate normal distribution, respectively. (G–K) Allele sharing between Han Chinese in each of the 31 administrative divisions with populations in the 1000 Genomes project by the F3 statistic. CHB, CDX, CHS, JPT and KHV refers to Han Chinese in Beijing, Southern Han Chinese, Chinese Dai in Xishuangbanna, Japanese in Tokyo and Kinh in Ho Chi Minh City, respectively.



Figure S4. Quality Evaluation in Genome-wide Association Studies, Related to Figure 4, Tables 1 and 2, and STAR Methods

(A and B) Distribution of mother's height and body mass index for the 61, 717 participants with height and BMI records.

(C and D) QQ-plot for height and BMI using all 2.1 million variants with Info Score > 0.4, HWE p value < 10-6 are shown in black. Results excluding all loci known to be associated with the trait is shown in gray line.

(E–K) Correlation of effect size of genome-wide significant variants between discovery set (NIPT) and three test sets (Giant, UK Biobank and CKB) for height(E-(G) and BMI(H-K). Linear regression was performed and the fitting line is shown in red.

(L) Distribution of the age of the mother in the year when taking the test. The four number "13," "28," "35," "48" refer to the minimum, two modes of the binomial distribution and maximum of age.

(M) QQ-plot for maternal age.

(N) QQ-plot for twin pregnancy.



Figure S5. Locus Zoom Plot for Novel Association Loci, Related to Figure 4, Tables 1 and 2, and STAR Methods

Locus Zoom plots for height (A-F); BMI (G-I); maternal age (J-K); twin pregnancy (L). The rs number and P value of the most significant primary SNP were labeled on top of the SNP. R2 LD estimated using 1KGP-CHN haplotypes between each SNP and the most significant SNP was color coded.



Figure S6. Extended Information for the Virome Spectrum in Plasma, Related to Figure 5 and STAR Methods

(A–C) Plasma virome not removing phage and low abundance participants. Shown are the coverage of the virus genome (A), prevalence (B) and the distribution of abundance of the viruses per individual (C). Only the top 40 viruses with greater than 10% genome coverage and among the highest rank of prevalence are shown. (D) Medical records for participants with HBV infection (significant hits > = 2). Free/NA means no clinical information are available.

(E) QQ-plot for high abundance ciHHV-6A/B phenotype.

(F) QQ-plot for low abundance ciHHV-6A/B phenotype.



Figure S7. Integrated Coverage from 138,882 Participants toward Virus Genome, Related to Figure 5 and STAR Methods Reads with mapping quality equal to zero were excluded. Viruses with coverage were shown.