



OPEN

## Genetic profiling of Vietnamese population from large-scale genomic analysis of non-invasive prenatal testing data

Ngoc Hieu Tran<sup>1,7</sup>, Thanh Binh Vo<sup>1,2</sup>, Van Thong Nguyen<sup>3</sup>, Nhat-Thang Tran<sup>4</sup>, Thu-Huong Nhat Trinh<sup>5</sup>, Hong-Anh Thi Pham<sup>1,2</sup>, Thi Hong Thuy Dao<sup>1,2</sup>, Ngoc Mai Nguyen<sup>1,2</sup>, Yen-Linh Thi Van<sup>1,2</sup>, Vu Uyen Tran<sup>1,2</sup>, Hoang Giang Vu<sup>1,2</sup>, Quynh-Tram Nguyen Bui<sup>1,2</sup>, Phuong-Anh Ngoc Vo<sup>1,2</sup>, Huu Nguyen Nguyen<sup>1,2</sup>, Quynh-Tho Thi Nguyen<sup>2</sup>, Thanh-Thuy Thi Do<sup>2</sup>, Nien Vinh Lam<sup>6</sup>, Phuong Cao Thi Ngoc<sup>2,8</sup>, Dinh Kiet Truong<sup>2</sup>, Hoai-Nghia Nguyen<sup>6</sup>✉, Hoa Giang<sup>1,2</sup>✉ & Minh-Duy Phan<sup>1,9</sup>✉

The under-representation of several ethnic groups in existing genetic databases and studies have undermined our understanding of the genetic variations and associated traits or diseases in many populations. Cost and technology limitations remain the challenges in performing large-scale genome sequencing projects in many developing countries, including Vietnam. As one of the most rapidly adopted genetic tests, non-invasive prenatal testing (NIPT) data offers an alternative untapped resource for genetic studies. Here we performed a large-scale genomic analysis of 2683 pregnant Vietnamese women using their NIPT data and identified a comprehensive set of 8,054,515 single-nucleotide polymorphisms, among which 8.2% were new to the Vietnamese population. Our study also revealed 24,487 disease-associated genetic variants and their allele frequency distribution, especially 5 pathogenic variants for prevalent genetic disorders in Vietnam. We also observed major discrepancies in the allele frequency distribution of disease-associated genetic variants between the Vietnamese and other populations, thus highlighting a need for genome-wide association studies dedicated to the Vietnamese population. The resulted database of Vietnamese genetic variants, their allele frequency distribution, and their associated diseases presents a valuable resource for future genetic studies.

Following the successful initiative of the 1000 genomes project<sup>1</sup>, several large-scale genome and exome sequencing projects have been conducted, either as international collaboration efforts such as ExAC<sup>2</sup>, gnomAD<sup>3</sup>, or for a specific country or population<sup>4–8</sup>. Those projects have provided comprehensive profiles of human genetic variation in some populations, paving the way for unprecedented advance in treatment of common genetic diseases. However, the lack of diversity and the under-representation of several populations in genome sequencing projects and genome-wide association studies (GWAS) have increasingly become a critical problem<sup>9,10</sup>. For instance, Gurdasani et al. found that the representation of ethnic groups in GWAS was significantly biased, with nearly 78% of the participants having European ancestries, whereas the two major populations, Asian and African, only accounted for 11% and 2.4%, respectively<sup>10</sup>. Vietnam has a population of 96.5 million, the 15th highest in the world and the 9th highest in Asia. Yet there was merely one dataset of 99 Vietnamese individuals that had been studied as part of the 1000 genomes project (population code KHV, the Kinh ethnic group in Ho Chi Minh City, Vietnam). A recent study has sequenced genomes and exomes of another 305 individuals to further expand the Vietnamese genetic database<sup>11</sup>. However, costs and technologies to perform large-scale genome sequencing projects still remain a challenge for most developing countries, including Vietnam.

<sup>1</sup>Gene Solutions, Ho Chi Minh City, Vietnam. <sup>2</sup>Medical Genetics Institute, Ho Chi Minh City, Vietnam. <sup>3</sup>Hung Vuong Hospital, Ho Chi Minh City, Vietnam. <sup>4</sup>University Medical Center, Ho Chi Minh City, Vietnam. <sup>5</sup>Tu Du Hospital, Ho Chi Minh City, Vietnam. <sup>6</sup>University of Medicine and Pharmacy, Ho Chi Minh City, Vietnam. <sup>7</sup>David R. Cheriton School of Computer Science, University of Waterloo, Waterloo, Canada. <sup>8</sup>Division of Molecular Hematology, Lund Stem Cell Center, Lund University, Lund, Sweden. <sup>9</sup>School of Chemistry and Molecular Biosciences, University of Queensland, Brisbane, Australia. ✉email: nhnghia81@gmail.com; gianghoa@gmail.com; m.phan1@uq.edu.au

An alternative approach has been proposed recently to re-use the low-coverage genome sequencing data from non-invasive prenatal testing (NIPT) for large-scale population genetics studies<sup>12,13</sup>. NIPT is a method that sequences cell-free DNA from maternal plasma at an ultra-low depth of 0.1–0.2 × to detect fetal aneuploidy<sup>14</sup>. By combining a sufficiently large number of NIPT samples, one could obtain a good representation of the population genetic variation. The benefits of re-using NIPT data for population genetics are manifold. First, the data can be re-used at no extra cost given the approval and consent of the participants. As one of the most rapidly adopted genetic tests, NIPT has been successfully established and become a standard screening procedure with thousands to millions of tests performed world-wide, including many developing countries such as Vietnam<sup>14</sup>. Using NIPT data for population genetic studies may also reduce privacy concerns since the genetic variants can only be analyzed by aggregating a large number of samples and the results can only be interpreted at the population level. A single sample tells little about the genetic information of an individual due to low sequencing depth. Last but not least, previous studies have suggested that sequencing a large number of individuals at a low depth might provide more accurate inferences of the population genetic structure than the traditional approach of sequencing a limited number of individuals at a higher depth, especially when the budget is limited<sup>15,16</sup>.

In this paper, we presented the first study of Vietnamese genetic variations from non-invasive prenatal testing data, and to the best of our knowledge, the third of such kinds in the world<sup>12,13</sup>. We analyzed the NIPT data of 2683 pregnant Vietnamese women to identify genetic variants and their allele frequency distribution in the Vietnamese population. We also studied the relationships between the Vietnamese genetic profile and common genetic disorders, and discovered pathogenic variants related to prevalent diseases in Vietnam. Finally, we highlighted the differences in the distribution of disease-associated genetic variants between the Vietnamese and other populations, thus highlighting a need for genome-wide association studies dedicated to the Vietnamese population. The resulted database of Vietnamese genetic variants from NIPT data is made available to facilitate future research studies in population genetics and associated traits or diseases.

## Results

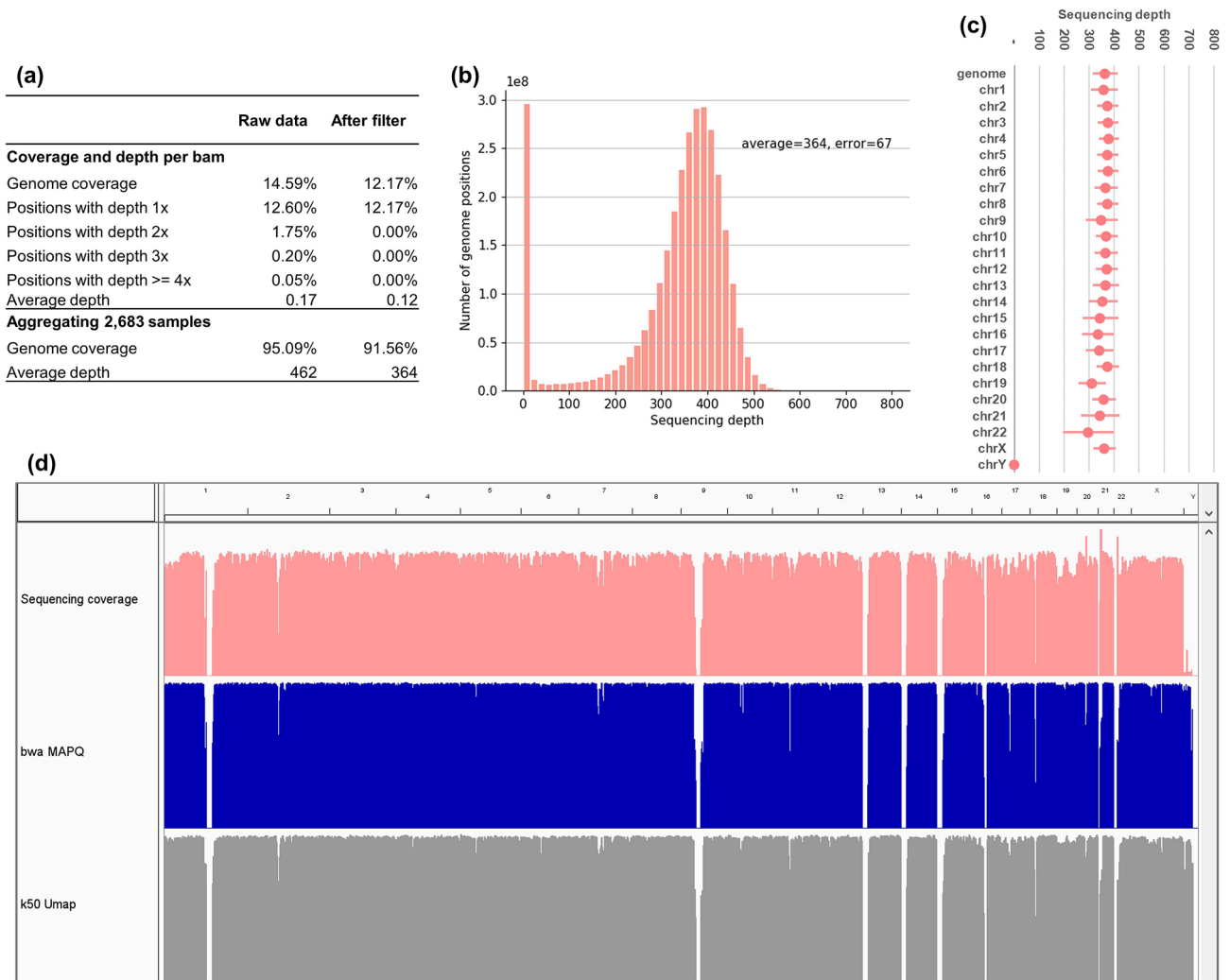
**Data collection.** A total of 2683 pregnant Vietnamese women who performed non-invasive prenatal testing during the period from 2018 to 2019 at the Medical Genetics Institute, Vietnam, were recruited to the study. The participants have approved and given written informed consent to the anonymous re-use of their genomic data for the study. All information of the participants is confidential and not available to the authors, except the records that their NIPT and pregnancy results are normal. The study was approved by the institutional ethics committee of the University of Medicine and Pharmacy, Ho Chi Minh city, Vietnam. The whole genome of each participant was sequenced to an average of 3.6 million paired-end reads of 2 × 75 bp, which corresponds to a sequencing depth of 0.17 × per sample.

**Genome coverage and sequencing depth of the NIPT dataset.** Data pre-processing was first performed on each of 2683 samples and the results were stored in binary alignment map (BAM) format, one BAM file per sample. The data pre-processing steps include: quality control of raw data using FastQC<sup>17</sup>, trimmomatic<sup>18</sup>; alignment of paired-end reads to the human reference genome (build GRCh38) using bwa<sup>19</sup>, samtools<sup>20</sup>, MarkDuplicates<sup>21</sup>; and summary of alignment results using Qualimap<sup>22</sup>, bedtools<sup>23</sup>, IGV<sup>24</sup>. The quality of raw data and alignment results are presented in Supplementary Figs. S1 and S2 for 1,000 samples randomly selected from our dataset. Raw data showed high sequencing quality, no bias was observed. The mapping quality and insert size distributions followed closely what expected across the reference genome. The overall sequencing error rate was estimated to be about 0.3% by Qualimap. More details of the data pre-processing steps can be found in the “Methods” section.

The average genome coverage and depth were 14.59% and 0.17 × per sample, respectively, and aggregated to 95.09% and 462 × across 2683 samples (Fig. 1a). Although the sequencing depth per sample was low, there might be more than one read from the same sample overlapping at a genome position. This problem may affect the estimation of allele frequency because the estimation is based on the assumption that a sample may contribute only 0 or 1 allele (read) at any given genome position<sup>12,13</sup>. For instance, we found that the average percentage of genome positions with depth 2 × (i.e. covered by two overlapping reads) in a sample was 1.75% (Fig. 1a). These overlapping reads occurred randomly across the reference genome and the samples. At any genome position, there were on average 47 out of 2683 samples that each contributed two reads (Supplementary Fig. S3). To address this problem, we followed a filtering strategy from previous studies<sup>12,13</sup> to keep only one read if there were overlapping reads in a sample. Thus, for every genome position, each sample could only contribute up to one read, and when the samples were aggregated, all reads at any position were obtained from different samples. In addition, we also removed alignments with low mapping quality scores (MAPQ < 30).

After the filtering step, the sequencing depth was reduced from 0.17 × to 0.12 × per sample. The aggregated sequencing depth of 2683 samples was 364 × and the genome coverage was 91.56% (Fig. 1a). The distributions of genome coverage and sequencing depth are presented in Fig. 1b–d. The sequencing depth was approximately uniform across the reference genome, except for low-mappability regions and chromosome Y. The distributions of sequencing depth and MAPQ score also closely followed the mappability of the human reference genome obtained from Umap<sup>25</sup>. The average sequencing depth of chromosome Y was about 2.1% that of the whole genome, consistent with the proportion of fetal DNA in NIPT samples (8–10%)<sup>14</sup>.

**Variant calling and validation.** We aggregated 2683 samples into one and used Mutect2 from GATK<sup>26–28</sup> for variant calling and allele frequency estimation. In addition to its main function of somatic calling, Mutect2 can also be used on data that represents a pool of individuals, such as our NIPT dataset, to call multiple variants at a genome site<sup>12,13,27</sup>. The called variants were further checked against strand bias, weak evidence, or contami-



**Figure 1.** Distributions of genome coverage and sequencing depth of the NIPT dataset. **(a)** Average genome coverage and sequencing depth per sample and from all samples combined. **(b)** Summary histogram of sequencing depth over all genome positions. **(c)** Distribution of sequencing depth per chromosome. **(d)** IGV tracks of sequencing depth, bwa MAPQ score, and Umap k50 mappability across the whole genome (the figure was produced using IGV, Integrative Genomics Viewer, version 2.8.9<sup>24</sup>).

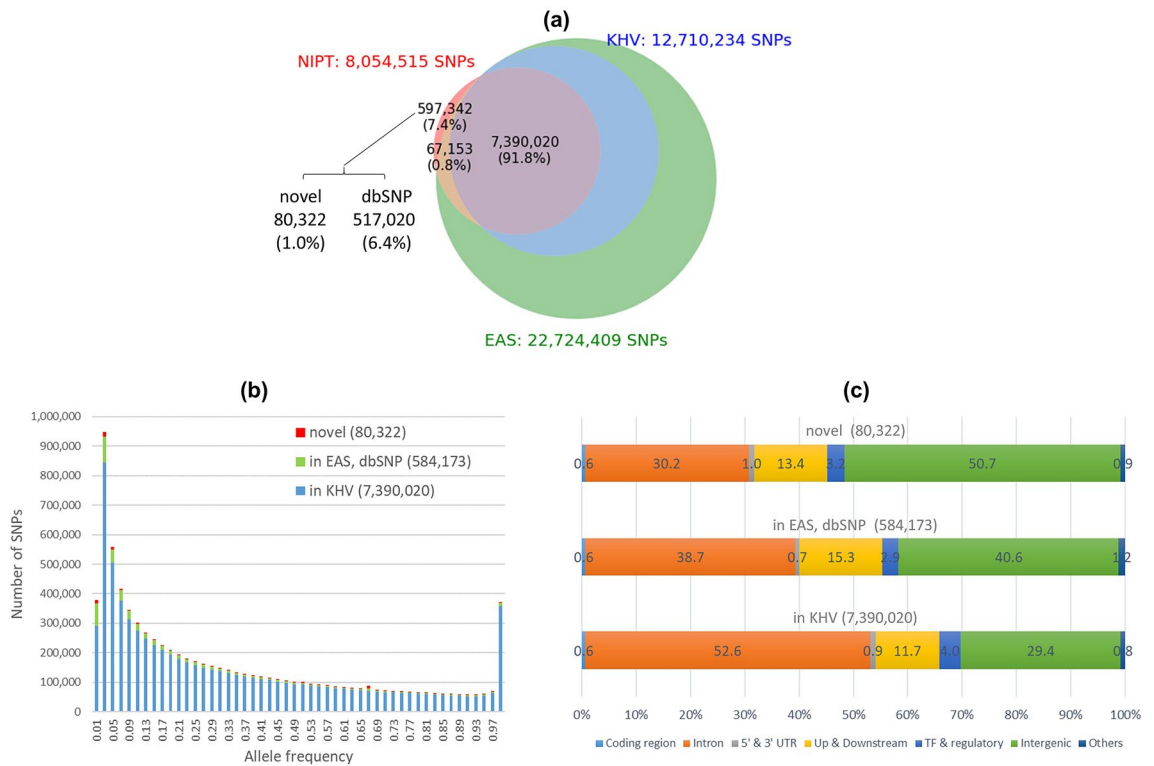
nation using FilterMutectCalls. The allele frequencies were estimated based on the numbers of reads aligned to the reference and the alternate alleles. For validation, we compared our NIPT call set to the KHV (Kinh in Ho Chi Minh City, Vietnam) and EAS (East Asian) populations from the 1000 genomes project<sup>1</sup>, as well as the dbSNP database (version 151<sup>29</sup>).

We identified a total of 8,054,515 SNPs from the NIPT dataset. The transition to transversion ratio was 2.0 over the whole genome and 2.8 over protein coding regions, which was similar to the observed ratios from previous genome or exome sequencing projects. As expected, a majority of these SNPs, 7,390,020 or 91.8%, had been reported earlier in the KHV call set (Fig. 2a). Since the KHV population only had 99 individuals, we further looked into its common SNPs that were shared by at least two individuals. We found that the NIPT call set recovered 90.5% of the KHV common SNPs (Supplementary Fig. S4;  $6,889,016/7,609,526 = 90.5\%$ ). This sensitivity is in line with the genome coverage reported earlier in Fig. 1a.

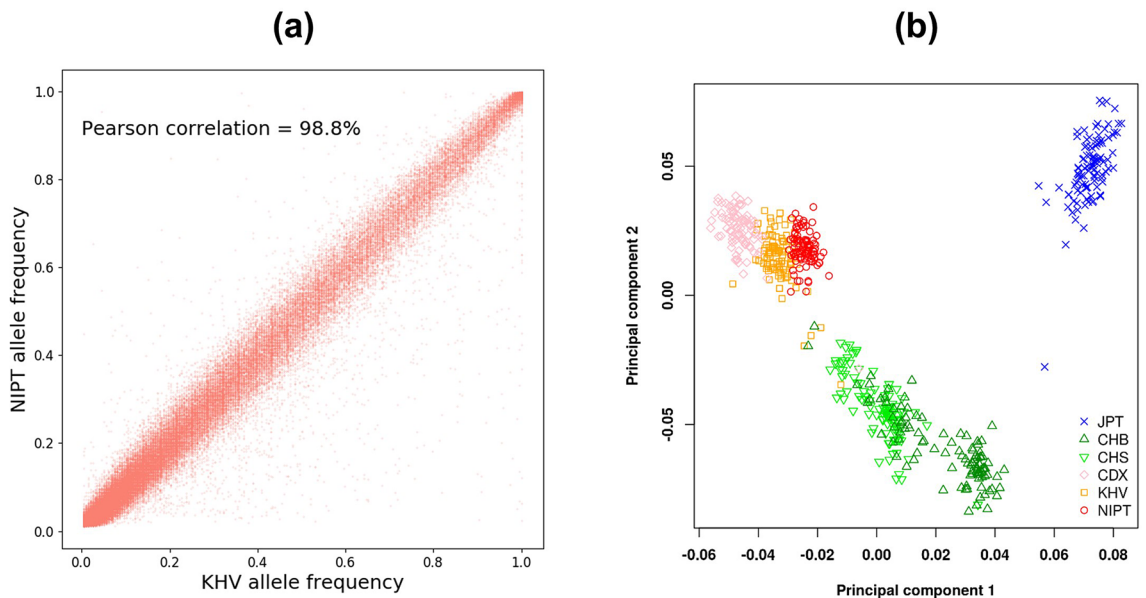
Our NIPT call set included 664,495 (8.2%) SNPs that had not been reported in the KHV call set. Among them, 67,153 (0.8%) were found in the EAS call set, another 517,020 (6.4%) were found in the dbSNP database, and the remaining 80,322 (1.0%) were novel SNPs (Fig. 2a). Majority of those SNPs had allele frequencies less than 10%. The overall allele frequency distribution of our NIPT call set is presented in Fig. 2b.

We used VEP (Variant Effect Predictor<sup>30</sup>) to analyze the effects of 8,054,515 variants in the NIPT call set (Fig. 2c). About 1.5% of the SNPs were located in the coding and UTR regions, 12% in the upstream and downstream regions, and 4% in the TF regulatory regions. More than 80% of the SNPs were located in the intron and intergenic regions. We also noted that the new SNPs and those in KHV had similar proportions of coding, UTR, upstream and downstream, and TF regulatory regions (Fig. 2c).

An important advantage of NIPT data is the ability of sampling a large number of individuals to better represent a population and to accurately estimate the allele frequency. We found a strong Pearson correlation of 98.8%



**Figure 2.** Summary of the NIPT call set. (a) Venn diagram comparison between the NIPT call set, the KHV and EAS call sets from the 1000 genomes project, and the dbSNP database. The percentages were calculated with respect to the NIPT call set. (b) Allele frequency distribution of the NIPT call set. (c) Distribution of locations and effects of variants in the NIPT call set.



**Figure 3.** Principal component analysis of the NIPT call set and other East Asia populations. (a) Scatter plot comparison of allele frequency estimated from the NIPT and the KHV call sets. (b) Principal component analysis. (JPT Japanese in Tokyo, Japan, CHB Han Chinese in Beijing, China, CHS Southern Han Chinese, CDX Chinese Dai in Xishuangbanna, China, KHV Kinh in Ho Chi Minh City, Vietnam, NIPT Non-Invasive Prenatal Testing data).

ClinVar clinical significance	In KHV	Not in KHV
Benign or likely benign	23,414	300
Uncertain significance	353	74
Pathogenic or likely pathogenic	4	1
Drug response	114	3
Others	211	13
Total number of annotations	24,096	391

**Table 1.** Summary of ClinVar annotations for the NIPT call set.

Variant information					ClinVar annotations			Allele frequency		
chr	Position	dbSNP	Ref	Alt	ID	Gene	Conditions	NIPT (%)	gnomAD EAS (%)	gnomAD (%)
chr18	57,571,588	rs2272783	A	G	562	FECH	Erythropoietic protoporphyria	28.10	32.57	11.23
chr13	20,189,473	rs72474224	C	T	17023	GJB2	Nonsyndromic hearing loss and deafness	13.40	8.35	0.76
chr13	72,835,359	rs17089782	G	A	217689	PIBF1	Joubert syndrome	6.80	5.13	1.36
chr6	26,090,951	rs1799945	C	G	10	HFE	Hemochromatosis type 1	5.10	3.41	10.82
chr2	31,529,325	rs9332964	C	T	3351	SRD5A2	5-alpha reductase deficiency	2.90	0.67	0.05

**Table 2.** Pathogenic variants identified from the NIPT call set.

between the allele frequency of the NIPT call set and that of the KHV call set (Fig. 3a). Furthermore, thanks to its larger sample size, the NIPT allele frequency indeed showed better resolution than the KHV one, as evidenced by vertical trails in Fig. 3a or a zoomed-in view in Supplementary Fig. S5. We also performed principal component analysis (PCA) on our NIPT data together with the data of East Asia populations from the 1000 Genome Project by using PLINK (version 1.9<sup>31</sup>). The PCA results in Fig. 3b show that the NIPT group is closely clustered with the KHV as both represent the Vietnamese population. The distribution of the populations in the PCA plot is also consistent with their geographic locations.

There are a few possible reasons as why the NIPT and KHV groups did not entirely overlap in the PCA plot. Firstly, they were obtained from two different types of data, ultra-low sequencing data of cell-free DNA from maternal plasma and whole-genome sequencing data. Secondly, in our variant calling approach for NIPT data, we aggregated all samples into one BAM file and then performed variant calling on that BAM file. Thus, we did not obtain individual genotypes, and instead, the allele frequencies were estimated based on the numbers of reads. Since PCA tools like PLINK require individual genotypes (e.g., 100 individuals for each population in the 1000 Genome Project), we used the NIPT allele frequency distribution to simulate 100 sets of genotypes. That could be another source of difference between the NIPT and the KHV groups in the PCA plot, as the former consisted of simulated genotypes from allele frequency distribution and the latter consisted of genotypes from real individuals.

**Analysis of pathogenic SNPs and their allele frequencies in Vietnamese population.** We searched the NIPT call set against the ClinVar database (version 20191105<sup>32</sup>) to explore the associations between Vietnamese genomic variants and common genetic diseases. We identified 24,487 SNPs with ClinVar annotations that have been reviewed by at least one research group (Table 1). Among them, five SNPs were classified as “Pathogenic” or “Likely pathogenic”, 117 SNPs were found to affect a “drug response”, and a majority of the remaining were classified as “Benign” or “Likely benign”. We also noted that 391 ClinVar-annotated SNPs (1.6%), including 1 pathogenic SNP, had not been reported in the KHV call set.

Table 2 and Supplementary Table S1 present the details of five pathogenic SNPs identified in our NIPT call set. Their associated genetic diseases include: erythropoietic protoporphyria, non-syndromic genetic deafness, Joubert syndrome, hemochromatosis type 1, and 5-alpha reductase deficiency. The SNP rs9332964 C > T in *SRD5A2*, which is associated with 5-alpha reductase deficiency, had not been reported in the KHV call set. 5-alpha reductase deficiency is an autosomal recessive disorder that affects male sexual development. This SNP is rare in the world and East Asia populations, but was found to be more common in the Vietnamese population (allele frequency of 0.05%, 0.67%, and 2.90%, respectively). This SNP had also been reported in a recent study<sup>11</sup> at a very low allele frequency of 1.36%.

We noticed that the allele frequencies of the five pathogenic SNPs varied considerably between the Vietnamese, the East Asia, and the world populations (Table 2). For instance, the SNP rs72474224 C > T in *GJB2* is commonly linked to non-syndromic hearing loss and deafness, which is also the most prevalent genetic disorder in the Vietnamese population. We found that its allele frequency in the Vietnamese population was 60% higher than in the East Asia population, which in turn was an order of magnitude higher than in the world population (13.40%, 8.35%, and 0.76%, respectively). The allele frequency was consistent with the estimated carrier frequency of 1 in 5 in the Vietnamese population. Similarly, the allele frequency of rs2272783 A > G in *FECH*, which is associated with erythropoietic protoporphyria, was nearly three times higher in the Vietnamese and

East Asia populations than in the world population (28.10%, 32.57%, and 11.23%, respectively). The prevalence of this pair of SNP and disease in East and Southeast Asia has been reported previously in<sup>33</sup>. On the other hand, the allele frequency of rs1799945 C>G in *HFE*, which is associated with hemochromatosis type 1, was about two and three times lower in the Vietnamese and East Asia populations than in the world population (5.10%, 3.41%, and 10.82%, respectively). Such discrepancies were also observed for “Benign” variants, e.g., those related to autosomal recessive non-syndromic hearing loss (Supplementary Table S2). The variations strongly suggest that population-specific genome-wide association studies are required to provide a more accurate understanding of the clinical significance of genetic variants and the true disease prevalence in the Vietnamese population.

## Discussion

In this study, we analyzed the genomes of 2683 pregnant Vietnamese women from their non-invasive prenatal testing data. The genomes were originally sequenced at a low depth of approximately  $0.17 \times$  per sample for the purpose of fetal aneuploidy testing<sup>14</sup>. Here we combined the 2683 samples to a total sequencing depth of  $364 \times$  and performed variant calling and analysis for the Vietnamese population. We identified a comprehensive set of 8,054,515 SNPs at a high level of sensitivity and accuracy: 90.5% of Vietnamese common SNPs were recovered; 99% of identified SNPs were confirmed in existing databases; and a strong correlation of 98.8% to known allele frequencies. The results were exciting given that the total sequencing depth of our dataset,  $364 \times$ , was merely equivalent to sequencing 20 individuals at a moderate depth of  $20 \times$ . It also suggests that there is still plenty of room for improvement by increasing the number of NIPT samples. For instance, Liu et al. have demonstrated a large-scale population genetic analysis based on hundreds of thousands of NIPT samples for the Chinese population<sup>12</sup>.

Another benefit of using NIPT data is cost-effective. In our study, the dataset was re-used at no cost with written informed consent from the participants. The whole analysis pipeline was done within a week on a local computer with 32 CPUs, 128 GB memory, and 5 TB hard disk. Thus, the overall cost was negligible compared to that of a typical genome sequencing project. The cost advantage of this approach may play a major role in large-scale genome sequencing projects, especially in developing countries where technologies and resources are still limited.

However, conducting population genetic studies using NIPT data is not without challenges, the most important of which is how to properly design and collect sufficiently large NIPT data across a population. This requires a coordinated, nation-wide effort, and more importantly, the informed consents of thousands of participants for their data to be used for such studies to improve public healthcare in the population. To facilitate this nation-wide data collection, a legal framework, including privacy policy, code of ethics, and standards of practice, is needed to protect confidential information and data of the participants.

Our study revealed 24,487 disease-associated genetic variants, especially five pathogenic variants for prevalent genetic disorders in the Vietnamese population. We also found major discrepancies in the allele frequency distribution of genetic variants between the Vietnamese, the East Asia, and the world populations. Thus, a comprehensive genetic profile and genome-wide association studies dedicated to the Vietnamese population are highly desired. Knowing the distribution of genetic disorders in the population will be useful for public health policy and planning, preventive medicine, early genetic screening strategies, etc.

Some technical and design limitations in our study could be addressed in future research to improve the application of NIPT data in population genetics studies. First, currently there is no variant calling tool that is designed specifically for NIPT data. Here we used Mutect2 and previous studies also used similar somatic calling tools with the purpose of identifying all possible variants at a genome site<sup>12,13</sup>. Thus, we took a conservative approach to consider only SNPs but not indels to ensure a reliable call set. The novel variants reported in Fig. 2a may include false positives, especially those with high allele frequency ( $> 5\%$ ) as it is not likely that they would have been missed in previous studies<sup>1</sup>. The transition to transversion ratio is 2.09 for known variants and 1.56 for new variants. We were not able to experimentally validate those novel variants and they should be interpreted with caution. Another limitation was the exclusion of chromosome Y due to its low coverage as a result of limited amount of fetal DNA in NIPT samples. This problem could be addressed by increasing the number of samples to obtain enough sequencing coverage for reliable variant calling. NIPT data is also biased by sex, with only  $\sim 5\%$  of the data coming from male population (assuming a 10% cell-free fetal DNA fraction with 50% male fetuses).

## Conclusions

We showed that non-invasive prenatal testing data could be reliably used to reconstruct the genetic profile of the Vietnamese population. Our study identified pathogenic variants for prevalent genetic diseases in the Vietnamese population and called for a need for population-specific genome-wide association studies. The resulted database presents a valuable resource for future studies of genetic variations and associated traits or diseases, not only for the Vietnamese population but also for other Southeast-Asia and Asia populations. Our results also demonstrated that non-invasive prenatal testing data provides a valuable and cost-effective resource for large-scale population genetic studies.

## Methods

**Sample preparation.** Cell-free DNA (cfDNA) in maternal plasma was extracted using MagMAX Cell-Free DNA Isolation Kit from Thermo Fisher Scientific (Waltham, MA, USA). Library preparation was done using NEB-Next Ultra II DNA Library Prep Kit from New England BioLabs (Ipswich, MA, USA). The samples were sequenced on the NextSeq 550 platform using paired-end  $2 \times 75$  bp Reagent Kit from Illumina (San Diego, CA, USA).

**Bioinformatics analysis pipeline.** Quality check of raw sequencing data was performed using FastQC (version 0.11.8<sup>17</sup>). Paired-end reads were trimmed to 75 bp, adapters (“TruSeq3-PE-2.fa”) and low-quality bases were removed using trimmomatic (version 0.39<sup>18</sup>). We only kept pairs with both reads surviving the trimming (about 95.7% of the dataset). The reads were then aligned to the human reference genome, build GRCh38 (hg38), using bwa mem (version 0.7.17-r1188<sup>19</sup>). Supplementary hits were marked as secondary for Picard compatibility. Alignment results were sorted and indexed using samtools (version 1.9<sup>20</sup>). Potential PCR duplicates were marked using MarkDuplicates from GATK (version 4.1.1.0<sup>26</sup>).

Alignments with mapping quality scores less than 30 were discarded. In-house Python scripts were developed to mark overlapping alignments and to keep only one of them. Qualimap (version 2.2.1<sup>22</sup>), bedtools (version 2.25.0<sup>23</sup>), and IGV (version 2.4.19<sup>24</sup>) were used to summarize the alignment results and to calculate genome coverage and sequencing depth.

Mutect2 from GATK (version 4.1.1.0<sup>26–28</sup>) was used in tumor-only mode for variant calling. All samples were assigned the same sample name to combine them before variant calling. FilterMutectCalls was used to exclude variants with weak evidence, strand bias, or contamination. bcftools (version 1.9<sup>34</sup>) was used to filter, summarize, and compare VCF (Variant Call Format) files. VEP (version 98<sup>30</sup>) was used to predict the effects of variants and to annotate them against dbSNP (version 151<sup>29</sup>) and ClinVar (version 20191105<sup>32</sup>) databases. Principal component analysis was performed using PLINK (version 1.9<sup>31</sup>).

The whole analysis pipeline and parameter settings can be found in the attached Python scripts.

**Ethics approval and consent to participate.** The study was approved by the institutional ethics committee of the University of Medicine and Pharmacy, Ho Chi Minh city, Vietnam. The study has followed the guidelines set by the University of Medicine and Pharmacy, Ho Chi Minh city, Vietnam, in handling human genetic data of the participants. The participants who performed NIPT triSure at Medical Genetics Institute, Vietnam, have approved and given written informed consent to the anonymous re-use of their genomic data for this study.

**Consent for publication.** All authors have read and approved the manuscript for publication.

### Data availability

The database of genetic variants identified from NIPT data and the Python scripts for bioinformatics analysis pipeline is available on GitHub: [https://github.com/nh2tran/NIPT\\_WGS](https://github.com/nh2tran/NIPT_WGS).

Received: 24 April 2020; Accepted: 26 October 2020

Published online: 05 November 2020

### References

1. The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
2. Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
3. Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
4. The UK10K Consortium. The UK10K project identifies rare variants in health and disease. *Nature* **526**, 82–90 (2015).
5. Gudbjartsson, D. F. *et al.* Large-scale whole-genome sequencing of the Icelandic population. *Nat. Genet.* **47**, 435–444 (2015).
6. Maretty, L. *et al.* Sequencing and de novo assembly of 150 genomes from Denmark as a population reference. *Nature* **548**, 87–91 (2017).
7. The Genome of the Netherlands Consortium. Whole-genome sequence variation, population structure and demographic history of the Dutch population. *Nat. Genet.* **46**, 814–825 (2014).
8. Wu, D. *et al.* Large-scale whole-genome sequencing of three diverse Asian populations in Singapore. *Cell* **179**, 736–749 (2019).
9. Editorial. Diversity matters. *Nat. Rev. Genet.* **20**, 495 (2019).
10. Gurdasani, D., Barroso, I., Zeggini, E. & Sandhu, M. S. Genomics of disease risk in globally diverse populations. *Nat. Rev. Genet.* **20**, 520–535 (2019).
11. Le, V. S. *et al.* A Vietnamese human genetic variation database. *Hum. Mutat.* **40**, 1664–1675 (2019).
12. Liu, S. *et al.* Genomic analyses from non-invasive prenatal testing reveal genetic associations, patterns of viral infections, and Chinese population history. *Cell* **175**, 347–359 (2018).
13. Budis, J. *et al.* Non-invasive prenatal testing as a valuable source of population specific allelic frequencies. *J. Biotechnol.* **299**, 72–78 (2019).
14. Phan, M. D. *et al.* Establishing and validating noninvasive prenatal testing procedure for fetal aneuploidies in Vietnam. *J. Matern. Fetal Neonatal Med.* **32**, 4009–4015 (2019).
15. Li, Y., Sidore, C., Kang, H. M., Boehnke, M. & Abecasis, G. R. Low-coverage sequencing: Implications for design of complex trait association studies. *Genome Res.* **21**, 940–951 (2011).
16. Fumagalli, M. Assessing the effect of sequencing depth and sample size in population genetics inferences. *PLoS ONE* **8**, e79667 (2013).
17. FastQC. <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>. Accessed 24 Apr 2020.
18. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: A flexible trimmer for Illumina sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
19. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. <https://arxiv.org/1303.3997v2> [q-bio.GN].
20. Li, H. *et al.* The sequence alignment/map format and SAMtools. *Bioinformatics* **25**, 2078–2079 (2009).
21. Picard. <https://broadinstitute.github.io/picard/>. Accessed 24 Apr 2020.
22. Okonechnikov, K., Conesa, A. & Garcia-Alcalde, F. Qualimap 2: Advanced multi-sample quality control for high-throughput sequencing data. *Bioinformatics* **32**, 292–294 (2016).
23. Quinlan, A. R. & Hall, I. M. BEDTools: A flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
24. Robinson, J. T. *et al.* Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).
25. Karimzadeh, M., Ernst, C., Kundaje, A. & Hoffman, M. M. Umap and Bismap: Quantifying genome and methylome mappability. *Nucleic Acids Res.* **46**, e120 (2018).
26. Van der Auwera, G. A. *et al.* From FastQ data to high confidence variant calls: The genome analysis toolkit best practices pipeline. *Curr. Protoc. Bioinform.* **43**, 1–33 (2013).
27. Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nat. Biotechnol.* **31**, 213–219 (2013).

28. DePristo, M. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
29. Sherry, S. T. *et al.* dbSNP: The NCBI database of genetic variation. *Nucleic Acids Res.* **29**, 308–311 (2001).
30. McLaren, W. The Ensembl variant effect predictor. *Genome Biol.* **17**, 122 (2016).
31. Chang, C. C. *et al.* Second-generation PLINK: Rising to the challenge of larger and richer datasets. *GigaScience*. <https://doi.org/10.1186/s13742-015-0047-8> (2015).
32. Landrum, M. J. *et al.* ClinVar: Public archive of relationships among sequence variation and human phenotype. *Nucleic Acids Res.* **42**, D980–D985 (2014).
33. Gouya, L. *et al.* Contribution of a common single-nucleotide polymorphism to the genetic predisposition for erythropoietic protoporphyria. *Am. J. Hum. Genet.* **78**, 2–14 (2006).
34. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–2993 (2011).

## Acknowledgements

The authors thank Dr. Kwok Pui Choi for critical reading of the manuscript.

## Author contributions

T.B.V., H.A.T.P., T.H.T.D., N.M.N., Y.L.T.V., V.U.T., H.G.V., Q.T.N.B. and P.A.N.V. performed experiments. V.T.N., N.T.T., T.H.N.T., T.T.T.D. and N.V.L. recruited patients and performed clinical analysis. H.N.N., Q.T.T.N., P.C.T.N. and D.K.T. designed experiments and analyzed data. H.N.N. supervised the project. N.H.T., H.G. and M.D.P. analyzed the data and wrote the manuscript, designed experiments and analyzed sequencing data.

## Funding

This study was funded by Gene Solutions, Vietnam. The funder did not have any additional role in the study design, data collection and analysis, decision to publish, or preparation of the manuscript.

## Competing interests

NHT, TBV, HATP, THTD, NMN, YLTV, VUT, HGV, QTNB, PANV, HNN, HG and MDP are current employees of Gene Solutions, Vietnam. The other authors declare no competing interests.

## Additional information

**Supplementary information** is available for this paper at <https://doi.org/10.1038/s41598-020-76245-5>.

**Correspondence** and requests for materials should be addressed to H.-N.N., H.G. or M.-D.P.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2020