

GigaScience

Korea4K: whole genome sequences of 4,157 Koreans with 107 phenotypes derived from extensive health check-ups --Manuscript Draft--

Manuscript Number:	GIGA-D-23-00109R1	
Full Title:	Korea4K: whole genome sequences of 4,157 Koreans with 107 phenotypes derived from extensive health check-ups	
Article Type:	Research	
Funding Information:	Ulsan National Institute of Science and Technology (1.200108.01)	Dr. Jong Hwa Bhak
	Ulsan National Institute of Science and Technology (1.200047.01)	Dr. Jong Hwa Bhak
	Small and Medium Business Administration (1425157253)	Dr. Jong Hwa Bhak
	Small and Medium Business Administration (1425157301)	Dr. Jong Hwa Bhak
	Small and Medium Business Administration (1425156792)	Dr. Jong Hwa Bhak
	Ministry of Trade, Industry and Energy (20016225)	Dr. Jong Hwa Bhak
Abstract:	<p>Background Phenome-wide association studies (PheWASs) have been conducted on Asian populations, including Koreans, but many were based on chip or exome genotyping data. Such studies have limitations regarding whole-genome-wide association analysis, making it crucial to have genome-to-phenome association information with the largest possible whole-genome and matched phenome data to conduct further population-genome studies and develop healthcare services based on population genomics.</p> <p>Results Here, we present 4,157 whole-genome sequences (Korea4K) coupled with 107 health check-up parameters as the largest genomic resource of the Korean Genome Project. It encompasses most of the variants with allele frequency > 0.001 in Koreans, indicating that it sufficiently covered most of the common and rare genetic variants with commonly measured phenotypes for Koreans. Korea4K provides 45,537,252 variants, and half of them were not present in Korea1K (1,094 samples). We also identified 1,356 new geno-phenotype associations which were not found by the Korea1K dataset. Phenomics analyses further revealed 24 significant genetic correlations, 14 pleiotropic associations, and 127 causal relationships based on Mendelian randomization among 37 traits. In addition, the Korea4K imputation reference panel, the largest Korean variants reference to date, showed a superior imputation performance to Korea1K across all allele frequency categories.</p> <p>Conclusions Collectively, Korea4K provides not only the largest Korean genome data but also corresponding health check-up parameters and novel genome-phenome associations. The large-scale pathological whole-genome-wide omics data will become a powerful set for genome-phenome level association studies to discover causal markers for the prediction and diagnosis of health conditions in future studies.</p>	
Corresponding Author:	Jong Hwa Bhak, Ph.D. Ulsan National Institute of Science and Technology Ulsan, Ulsan KOREA, REPUBLIC OF	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Ulsan National Institute of Science and Technology	

Corresponding Author's Secondary Institution:	
First Author:	Sungwon Jeon, Ph.D.
First Author Secondary Information:	
Order of Authors:	<p>Sungwon Jeon, Ph.D.</p> <p>Hansol Choi, Ph.D.</p> <p>Yeonsu Jeon, Ph.D.</p> <p>Whan-Hyuk Choi, Ph.D.</p> <p>Hyunjoo Choi, Master of Science</p> <p>Kyungwhan An, Bachelor of Science</p> <p>Hyojung Ryu, Ph.D.</p> <p>Jihun Bhak, Bachelor of Science</p> <p>Hyeonjae Lee, Bachelor of Science</p> <p>Yoonsung Kwon, Bachelor of Science</p> <p>Sukyeon Ha, Bachelor of Science</p> <p>Yeo Jin Kim, Ph.D.</p> <p>Asta Blazyte, Master of Science</p> <p>Changjae Kim, Ph.D.</p> <p>Yeonkyung Kim, Master of Science</p> <p>Younghui Kang, Bachelor of Science</p> <p>Yeong Ju Woo, Bachelor of Science</p> <p>Chanyoung Lee, Bachelor of Science</p> <p>Jeongwoo Seo, Bachelor of Science</p> <p>Changhan Yoon, Master of Science</p> <p>Dan Bolser, Ph.D.</p> <p>Orsolya Biro, Ph.D.</p> <p>Eun-Seok Shin, M.D., Ph.D.</p> <p>Byung Chul Kim, Ph.D.</p> <p>Seon-Young Kim, Ph.D.</p> <p>Ji-Hwan Park, Ph.D.</p> <p>Jongbum Jeon, Ph.D.</p> <p>Dooyoung Jung, Ph.D.</p> <p>Semin Lee, Ph.D.</p> <p>Jong Hwa Bhak, Ph.D.</p>
Order of Authors Secondary Information:	
Response to Reviewers:	<p>GIGA-D-23-00109</p> <p>Korea4K: whole genome sequences of 4,157 Koreans with 107 phenotypes derived from extensive health check-ups</p> <p>Sungwon Jeon; Hansol Choi; Yeonsu Jeon; Whan-Hyuk Choi; Hyunjoon Choi; Kyungwhan An; Hyojung Ryu; Jihun Bhak; Hyeonjae Lee; Yoonsung Kwon; Sukyeon Ha; Yeo Jin Kim; Asta Blazyte; Changjae Kim; Yeonkyung Kim; Younghui Kang; Yeong Ju Woo; Chanyoung Lee; Jeongwoo Seo; Changhan Yoon; Dan Bolser; Orsolya Biro; Eun-Seok Shin; Byung Chul Kim; Seon-Young Kim; Ji-Hwan Park; Jongbum Jeon;</p>

Dooyoung Jung; Semin Lee; Jong Hwa Bhak
GigaScience

Reviewer reports:

Reviewer #1: This manuscript describes the second phase of the Korean Genome Project (KGP) with 4,157 sets of whole-genome data (designated Korea4K). After error correction and sequencing data curation, the whole-genome sequencing (WGS) data from 3,614 unrelated were used in the analyses. They also analyzed 107 types of clinical traits from 2,685 healthy participants' health check-up reports over a 4-year period (2016-2019). They performed a range of analyses and claimed that this new data performed better than Korea1K, the first phase KGP dataset, in a number of ways. A larger Korean dataset adds to the global genome resource and provides further insights into the Korean population. However, the results are mostly descriptive and serve as a catalog without significant new insights. The results are as expected (Korea4K is a better imputation reference panel than Korea1K, new variants are identified in the population, new variants are found in association with various phenotypes, etc.) and this dataset is sufficiently large to capture all the common variants found in the homogeneous Korean population.

The authors should address several issues:

1. The use of whole genome sequencing data in GWAS. The Bonferroni correction the authors used in their analysis was that for SNP array studies. They must do a formal correction with the many more variants found in WGS data and use a statistically sound correction for their analysis. The severe penalty for multiple testing using WGS data for GWAS is why few such studies have been done. I suspect that many of the associations will not reach statistical significance after proper correction, as the dataset is quite small for most traits under study.

⇒ Thank you for your critical comments regarding the statistics on the GWAS results. As you pointed out, we agree that there should be a stricter correction. As one method, we have now employed the FDR correction (Benjamini-Hochberg) which can remove possible false positives. The FDR values for each variant are now included in Supplementary Table S6 (List of the GWAS variants which have association significance $P < 5E-8$). After the FDR correction, 2314 variants from 30 traits still maintained statistically significant associations ($FDR < 0.05$). We additionally noted the number of remaining variants and traits according to different FDR cutoffs in Table R1 below. Also, the results of the FDR correction were updated in the manuscript.

Page 11: "Among the significantly associated variants, 2,314 variants from 30 clinical traits still showed significance after false discovery rate (FDR) correction using the Benjamini-Hochberg approach ($FDR < 0.05$)."

Table R1. Number of significantly associated variants and traits according to FDR cutoffs

FDR cutoff	Number of variants	Number of traits
$FDR < 0.1$	2320	31
$FDR < 0.05$	2314	30
$FDR < 0.01$	2256	24
$FDR < 0.005$	2193	24
$FDR < 0.001$	1916	18

2. The authors should use the new genome references for their variant calling (T2T reference and the Human Pangenome Reference), as the GRCh38 is no longer the gold standard, and the results will be quite different with the most up-to-date references. Using the best human genome reference will make Korea4K more valuable.

⇒ We agree that we could potentially find more genetic markers that are related to the traits in our GWAS analysis, for example, when we use the T2T reference or the

Human Pangenome Reference. However, there are a few considerations that limit us from using these references:

The T2T reference lacks enough annotation data which is critical. Many major genomic databases, such as dbSNP, are based on GRCh38. Thus, even if we were to use the T2T reference, we would have limitations in interpreting or validating the variants/markers we could additionally discover.

The draft Human Pangenome reference genome contains genomic sequences of 47 genetically diverse individuals, which requires a totally different bioinformatics pipeline to analyze. The bioinformatics analysis using the Human Pangenome reference is not fully established currently, which means that the validation method of the genetic markers that could be discovered should also be investigated more. Also, the Human Pangenome reference was assembled based on long-read sequencing data such as PacBio and Oxford Nanopore Technologies (ONT). As the authors of the Human Pangenome reference paper mentioned, the 1-base level of the sequencing accuracy can be an issue, which makes it hard to know if additional discoveries using the Pangenome are true signals or artifacts.

Furthermore, mapping the whole-genome sequencing reads for 4K samples and jointly genotyping the variants technically requires more than a year of time to rebuild the dataset.

We understand the importance of more precise and complete genome references for the variant calling and we appreciate your suggestion. We will expand the variant call set using the T2T and Human pangenome references in our future studies. Unfortunately, as we mentioned above, due to several technical limitations, we were not able to apply the new genome references in our current study, although we revised the manuscript to add the importance of the usage of these references.

Page 22: "Moreover, utilizing recently introduced human genome references like the T2T reference [33] and Human Pangenome reference [34], which offer broader genomic coverage or have population-specific sequences compared to the existing GRCh38 reference, could help identify additional associations that might be overlooked. Nevertheless, these new references lack functional annotations and need to be connected to previous databases such as dbSNP and the GWAS catalog."

3. The authors should clarify how many of the participants who contributed clinical data are unrelated.

⇒ As described in the Methods, we filtered out a total of 540 individuals including 428 samples that have relatedness to other samples from 4,157 samples. Among the final unrelated 3,617 samples, 2,262 samples had clinical data. We updated the manuscript to provide a clear description of the participants who contributed to the clinical data.

Page 28: "Out of the final unrelated 3,617 samples, 2,374 samples had clinical data available and were included in the phenomics analyses."

Reviewer #2: The authors contribute 4,157 whole-genome sequences (Korea4K) coupled with 107 health check-up parameters as the largest genomic resource of the Korean Genome Project. It has likely characterized most of the common and very common genetic variants with commonly measured phenotypes for Koreans. It also discusses its applicability not only for the Korean population but also for other East Asian populations, and possibly to other national genome projects as well.

This work makes a significant contribution of data that can be used in future genome-wide association studies in the context of the Korean population. The manuscript appears to cover a lot of ground: from methodological issues to the real-world applications of the dataset in healthcare. The authors adopt innovative methods like GREML, which have been reported to have higher accuracy compared to older

methods.

The authors are transparent about the limitations of their study, such as sample size and lack of sufficient data for rare diseases. They also acknowledge that phenomics analyses were not powerful enough for novel discoveries, indicating areas for future research. However, given the increasing importance of genomic data in healthcare and personalized medicine, the paper appears to be highly relevant.

While the paper is well formulated, there are some issues that need to be addressed before is accepted for publication.

See below:

1. You referred to the UK Biobank data for some of your analyses. Were there any limitations or caveats in comparing your dataset to the UK Biobank? What about other national genomic projects that are out there? How transferable do you think the Korea4K dataset would be to studies focusing on other populations outside East Asia?

⇒ Yes, there can be many limitations/caveats. However, please take in consideration that we cannot report the limitations precisely since we have not fully utilized the raw genomic and extensive clinical data from the UK Biobank but only the GWAS summary data in the current study.

One clear limitation is the difference in allele frequencies of reported variants between two ethnic groups in comparison due to different population genomic structures. This has been demonstrated in our PCA results (please refer to Figure 2). Even if the two populations were to have the same sample size, disparities in the allele frequencies would lead to different summary statistics (i.e., beta- and p-values). Therefore, the downstream, namely the phenomic, analysis could suffer from different resultant statistics and a fair comparison between the two independent studies could be difficult for some trait pairs.

Another possible limitation is that the current GWAS summary data provided by the UK Biobank (and other national biobanks) is based on arrays and many of the variants are imputed genotypes. Here, we have utilized the whole-genomes. We found that the GWAS-significant variants in comparison are not often overlapping due to technical biases as reported in our prior study (Jeon YS et. al., 2023, Hum Genet.). A possible solution to this would be to perform the joint genotyping of the Korean whole-genomes in conjunction with the UK Biobank whole-genomes. However, the entire process is resource-intensive and time-consuming, such that only the institutes with sufficient computing power will be able to process the data. In addition, rare variants appearing in each ethnic group will be grossly undermined.

The content and amount of phenotypic and clinical information provided are another possible limitation. Our health check-up data have been collected from multiple sources/centers such that we had to process the heterogeneous physical and digital copies of the health records and standardize them. On the other hand, all the participants in the UK Biobank were sampled in a single-centered manner with a unified procedure. Hence, our phenotypic and clinical information is much more limited than the UK Biobank's, and the method of measurements may differ for a few specific categories although we did not deeply investigate.

Furthermore, Korea4K data is not as readily accessible as the UK Biobank data. Obtaining it requires navigating through IRB processes and administrative procedures, and the legislative framework in South Korea does not currently facilitate a straightforward download process. If improved in the future, we will be able to make them more accessible. As an initial step toward enhanced accessibility, we have deposited our dataset in the European Genome-Phenome Archive (EGA) under the study accession 'EGAS00001007580' and are actively working towards providing the WGS data openly and freely.

Moreover, we have included the accession number in the manuscript to facilitate easy reference.

Page 34: "The raw sequencing data that can be distributed were uploaded to the European Genome-Phenome Archive under the study accession

'EGAS00001007580'.

2. Could you expand on any ethical considerations that were taken into account, especially in terms of data privacy and informed consent?

⇒ Yes. In our project, we have taken ethical considerations, particularly concerning data privacy and informed consent, very seriously. The data we generated and used in our study comes from blood or saliva donations, and we have obtained explicit, full consent forms from the participants before getting the samples. These consent forms ensure that the participants are aware of how their data will be used and that they have willingly agreed to share this data for research purposes and IRB. As a result, we can make the data of 3,839 individuals publicly available while respecting the privacy and consent of the participants.

We have expanded the ethical considerations in "Ethics, consents and permissions":

Page 34: "The data employed in our study originates from voluntary blood or saliva donations, and we have diligently secured explicit, comprehensive consent forms from all participants prior to sample collection. These consent forms explicitly outline the intended use of their data for research purposes and underscore the voluntary nature of their participation. Furthermore, our study adheres to the ethical guidelines and regulations stipulated by the IRB. As a result, we can make the data of 3,839 individuals publicly available while respecting the privacy and consent of the participants."

3. How was the data cleaned and preprocessed, and were there any missing data points? If so, how were these handled? What number of reads(before and after QC), and other quality metrics do the sequenced reads have? What was the average coverage across the genome? What was the read length?

⇒ We cleaned and preprocessed by trimming possible adapter contamination on the sequencing reads and made the reads have at least 50bp of read length using the Cutadapt program (ver. 1.9.1). Then we confirmed the read counts, quality, and amount of bases using the FASTQC program. The average sequencing depth was $27.75 \times$ and initial read lengths were 151bp. However, it varied after the trimming. The number of reads, average coverage, average quality, and filtering percentage were visualized in Supplementary Figure S4 and we updated the preprocessing procedure in the method section of the manuscript.

Page 23: "All the sequencing data that we used in this study had 151bp as a read length. Average sequencing amount per sample was $20 \times$ (Supplementary Figure S4)."

Page 24: "Adapter contamination was trimmed using Cutadapt (RRID: SCR_011841, ver. 1.9.1) [35] with a forward adapter ('GATCGGAAGAGCACACGTCTGAACTCCAGTCAC') and reverse adapter ('GATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT') and with a minimum read length of 50 bp after trimming (Supplementary Figure S4)."

Page 24: "A total of 3,156 samples had a mapping depth of $\geq 20 \times$ (Supplementary Figure S4)."

4. How did you ensure the quality of the genomic data collected from different sources such as Korea1K and public data archives?

⇒ To ensure the quality, we rigorously utilized a standardized pipeline for processing the genomic data and performed batch effect removal.

We applied the same bioinformatics analysis and QC pipelines as in the Korea1K

(Jeon S et. al., 2020, Sci. Adv.). For example, we ensured the same version of the programs, and the same parameters through all the WGS data when we genotyped the samples.

We collected the WGS data from multiple sources or sequenced the whole genomes at different time points, which could suffer from the batch effect. After jointly genotyping the WGS data, we tried to reduce the batch effect from the different sequencing batches. As we noted in the manuscript, we applied the allele balance-based variant filtering.

The paper mentions mitigating batch effects through allele balance and manual checks. Could you provide more details on the methodology behind these checks and their efficiency?

⇒ As for the manual checks, we filtered out the undesired samples based on the following categories:

- (1) high missing genotype rate (>10%);
- (2) outlying heterozygous variants ratio (3 s.d.);
- (3) high relatedness;
- (4) non-Korean genetic background from PCA;
- (5) having a rare disease;
- (6) samples no longer available.

As for the allele balance, we first measured the average allele balance of the genotyped alleles (the read count of the allele divided by the total read count on a locus). Then, we excluded 12,713,580 variants that had an average allele balance of the loci out of the range of $\pm 1 \times$ standard deviation (SD) from a genome-wide average of allele balance to remove the sequencing batch effect. We confirmed that the batch effects were removed after the filtering and visualized it as PCA plots (Supplementary Figure S1).

We have now added the definition of the allele balance in our Methods:

Page 25: "To detect variants which were probably called because of a sequencing batch effect, we measured average allele balance of the genotyped alleles (the read count of the allele divided by the total read count on a locus)."

5. Could you provide more information about the control group? Was it matched for age, sex, or other variables?

⇒ In the Korea4K dataset, we do not have a specific control set. We conducted our GWAS analysis only for the quantitative clinical traits using a linear regression approach without separately categorizing the study participants into case and control groups. Age, Sex, and BMI were included as covariates in testing the significance of variants across the clinical variables of interest.

How was the sample size determined, and does it provide enough statistical power to support your conclusions?

The sample size employed for each trait was maximized by utilizing the available samples in the Korea4K. Regarding the statistical power of GWAS, we recognize the importance of ensuring an adequate sample size to obtain robust results. We calculated the statistical power and effect size based on a likelihood ratio test by the R package ("genpwr"). Out of the 90 WGwas traits analyzed, a majority of traits (77 traits) exhibited enough statistical powers exceeding 80% under the assumption of an effect size of 0.5 and a minor allele frequency (MAF) of 0.01 (added to Supplementary Table S14). Given the sufficient statistical power for the majority of traits examined. Our study's sample size of 4,157 individuals was appropriate for addressing the research objectives.

The detailed method is added to the manuscript:

Page 29: "Statistical powers of the 90 GWAS were calculated by the R package "genpwr" under the assumption of an effect size of 0.5 and a minor allele frequency of 0.01 (Supplementary Table S14)."

6. You mentioned that the statistical power of your study will increase with more participants. Would this have implications for other national genomes that are making similar projects?

⇒ Probably yes and no. Yes, for other populations which are as homogenous as Koreans. The current study is very specific in that the population is very genetically homogenous. Koreans are probably the most homogenous population in East Asia. The diversity is even less than the Japanese archipelago (and the mainland China) because the Korean peninsula has been geographically isolated compared to other islandic populations. No, if a population is extremely heterogeneous with mixed ethnicities, the statistical power of our analysis would be much lower with the same sample size. In general, the statistical power will increase as the sample size goes up as you already know. However, we are unsure if the behavior would be exactly as our claim, and further studies are warranted.

Please elaborate on how your sensitivity analysis could apply to other populations outside Korea.

⇒ Thank you for a very interesting question. We find it very hard to answer because as mentioned above the Korean population is extremely homogeneous. It is a special population in terms of genetic and environmental diversity. Therefore, it is very difficult to estimate what kind of insight we may apply to other populations outside Korea. It could be an advantage or a disadvantage depending on the application of the population.

7. The paper acknowledges the sample size as not sufficiently large for detecting weak associations and admits that the sample size was not large enough to detect weak association signals. Have you considered statistical methods that can boost power in small samples?

⇒ We appreciate your comment. We could think of two statistical methods to improve statistical power: Meta-analysis and Gene-based association test. 1) Meta-analysis is one of the methods to boost statistical power for small samples by combining GWAS summary statistics of multiple independent studies. However, our study focuses specifically on the whole-genome-wide associations in the Korean population. Although a meta-analysis could improve the sensitivity of the weak association signals, the analysis might be able to introduce potential biases due to the genetic heterogeneity across populations. Thus, we have focused on our analysis of the Korean population. The Korean population's genetic characteristics and their implications for clinical trait associations would provide an important aspect of our study, providing a unique and valuable. 2) Gene-based association test is another approach that increases the power to detect associations. This method aggregates the effects of multiple genetic variants within a gene to assess their collective impact. While the method is particularly advantageous for analyzing rare variants, we focused on identifying common variants associated with clinical traits in the Korean population. The application of gene-based association tests to our study could potentially yield additional insights, especially in the context of rare variant associations. We acknowledge the value of this approach and are considering its application in our subsequent study with a larger sample size.

8. Could you provide more details on the 107 clinical parameters used for the Korea4K phenome dataset? Were these parameters standardized across the different clinics and hospitals?

⇒ Yes, we standardized the parameters across the different clinics. For example,

discrepancies in unit measurements, such as micrograms (ug) and nanograms (ng), were unified for specific traits, posing a direct challenge to the analytical process if these variations were not duly reconciled. Also, some parameters such as e-GFR were calculated by different equations across the clinics. We re-calculated such parameters using a singular formula. More detailed methods were updated in the method section of the manuscript and Supplementary Table S13.

Page 28: "In the context of collecting data from over 200 diverse healthcare institutions, standardizing clinical information on 107 traits became imperative. We resolved discrepancies in unit measurements, such as micrograms and nanograms, for specific traits. Furthermore, certain clinical metrics, such as the estimated glomerular filtration rate (e-GFR), were found to exhibit variability contingent upon variables such as ethnicity, sex, and age. To maintain consistency and ensure methodological uniformity, we enforced the adoption of a singular clinical formula for the computation of e-GFR across all data samples. Such calculations were applied to 26 traits (Supplementary Table S13). Clinical traits that exhibited values characterized by inequalities likely due to the limit of detection (e.g., <5.0 and >99) were omitted from the analytical procedures, as such values have the potential to introduce disturbances to subsequent data analyses. Likewise, values that exhibited divergent formatting conventions across distinct healthcare institutions (e.g., 20 and a few or 999 and many) were harmonized to conform with prevailing standard criteria observed in most samples under investigation."

9. What criteria were used for initial sample filtering, particularly for excluding kinship? Could you clarify the steps taken to identify and filter the 64,301,272 SNVs and 8,776,608 Indels? How did you correct for batch effects arising from different Illumina NGS platforms and library preparations? Did you use specialized SNV calling software, or only GATK?

⇒ We have briefly mentioned this in our method section, we filtered out total 540 samples and the detailed filtering steps are noted in the method section as well. To exclude the samples who were in the kinship relationship, we first measured IBD values between the samples using Plink program (ver. 1.90b3n) and defined the family trees based on the IBD values which showed a PI_HAT value more than 0.05. Then, we filtered out samples in a family tree to have the maximum number of the remaining samples.

We also updated the manuscript to provide detailed procedures for defining the kinship relation between samples:

Page 25: "To explore kinship relations among the samples, we assessed Identical by Descent (IBD) using the Plink program (RRID:SCR_001757, ver. 1.90b3n) [30]. Samples with a PI_HAT value exceeding 0.05 were considered to be in a kinship relation."

About the variant calling, we did not use specialized SNV calling software. As we noted in the method section, we jointly genotyped the genotypes using only GATK 4.1.3 and identified the 64,301,272 SNVs and 8,776,608 Indels. With the jointly genotyped data, we measured allele balance of the loci. If the average allele balance of the loci was out of the range of $\pm 1 \times$ standard deviation (SD) from a genome-wide average of allele balance, the loci were treated as generated by possible batch effects due to different Illumina NGS platforms and library preparation and filtered out using an in-house script. A similar method was previously suggested by Muya F, et. al., 2019. Following the filtering method, we excluded 12,713,580 variants and confirmed that the batch effects were removed through the PCA plots in Figure S1.

10. How were allele frequencies calculated and what considerations were made to interpret their biological significance? You mention that more than half of the singleton and doubleton variants were newly discovered. Could you elaborate on the methodology used to confirm these as novel variants?

⇒ The allele frequencies were calculated by the number of alternative alleles divided by the number of called alleles in a position. Generally, the allele frequency distribution can reflect the genomic diversity of the population. The definition we used to assign the variant as a novel variant is whether the variant was reported in the dbSNP database or not. We updated the figure legend of Figure 1 of the manuscript to provide the definition.

Page 9-10: “dbSNP indicates the variants were reported in dbSNP database. Novel indicates the variants were not reported in dbSNP.”

11. The section on phenotypic correlations mentions 2,274 trait-trait relationships. How would you address the potential for population stratification affecting the results of your genetic and phenotypic correlations?

⇒ The problem of population stratification in GWAS especially arises when conducting GWAS in multi-ethnic countries or meta-analyses across multiple sources of data with mixed ancestries. Here, our study exclusively deals with samples of Korean ancestry, which means our dataset is genetically homogeneous compared to other studies. To make sure, we excluded any samples with non-Korean genetic backgrounds based on PCA analysis as we noted in the “Sample and variant filtering” section in the Methods. Also, we included PC1~10 as covariates in our regression models during GWAS to avoid any latent ancestral effects from differential ethnic subgroups as we noted in the “Whole genome-wide association study (WGWAS)” section in the Methods.

We argue that such factors may add to reducing spurious correlations introduced by population stratification. Our claim is supported by values for the genomic inflation factor (λ Median) (Supplementary Figure S4-19). As you may already know, the value of lambda below 1.1 is generally considered acceptable indicating minimal false positives caused by gross population structure (and systematic biases) (please refer to Yang et. al., 2011, Eur J Hum Genet.)

We have now included the calculation of genomic inflation factor to estimate the gross population structure (and systematic biases) in our data.

Page 29: “Calculating the genomic inflation factor (λ Median), we found that all of the traits in the test reside below 1.1 indicating there are minimal false positives caused by gross population structure or systematic biases (Supplementary Figure S4-19) [48].”

How did you account for multiple comparisons in determining significant genetic correlations, and what corrections were applied to maintain the FDR?

⇒ We see that we were insufficient in our details as to how we corrected for multiple comparisons while calculating genetic correlation. We put our best effort to avoid false discoveries by conducting Benjamini-Hochberg correction to maintain the FDR well and below 0.05. Consistently, we applied the same correction method for setting FDR for phenotypic correlation as well.

We have now added the Methods for multiple testing correction for phenotypic correlation and genetic correlation:

Page 29: “Benjamini-Hochberg method was used to adjust for multiple comparisons when documenting confident phenotypic correlations with FDR.”

Page 30: “The correction for multiple tests was done by Benjamini-Hochberg approach when reporting confident GCs that suffice the threshold of FDR below 0.05.”

What measures were taken to ensure that the traits considered in this section were not

subject to confounding and/or collider biases.

⇒ Thanks for pointing the important question. Confounding and collider biases are unavoidable when looking at multiple associations across numerous variables at once. First, we employed covariate adjustment to reduce confounding biases by traits that are highly correlated with other clinical variables. We focused on Age, Sex, and BMI which have previously been suggested by Shungin et al. (2015) (Shungin D et al., 2015, Nature). Second, we incorporated Mendelian Randomization (MR). MR is a statistical method to ascertain the direction of effect and imply possible causality devoid of confounders and colliders (Mitchell RE et. al., 2023, PLOS Genetics, and Ebrahim S and Davey Smith G, 2008, Human Genetics). To raise the confidence of our claim, we utilized three independent MR methods, namely IVW, MR-Egger, and MR-PRESSO, and documented the causal relationships if at least two of three were shown significant as noted in the Methods. However, we acknowledge that our methods of MR might still be prone to a collider bias, especially due to conditioning by covariates, which was to remove confounders and increase the power (Cai S et. al., 2022, Genetic Epidemiology). However, we would like to emphasize that such bias has minimal effect on the interpretation of phenotypic associations as previously reported by Pulit et al. (2019) (Pulit SL et. al., 2019, Hum Mol Genet).

We have added our responses to your comment into the Methods and Results, respectively:

Page 29: "Age and BMI were chosen especially due to their known shared associations with multiple traits as previously documented by Shungin and colleagues which could lead to confounding biases in the downstream interpretation of phenotypic relationships [47]."

Page 16: "In addition to the investigation on the general pleiotropic relationship, we employed Mendelian Randomization (MR) to detect vertical pleiotropy that can assert the direction of the phenotypic relationships [18]. This provides indirect evidence implying causality between the traits to discern spurious phenotypic associations, such as confounding and collider bias [19, 20]."

12. In your findings, Waist-Creatine showed opposite directions for genetic and phenotypic correlations. Could you elaborate on the potential implications or causes of this discrepancy?

⇒ Thanks for raising this point so that we could put attention here for a richer discussion. We suggest the discrepancy mainly comes from the shared environmental factors between Waist and Creatinine. This is well-reviewed by Sodini et al. (2018) (Sodini SM, et. al., 2018, Genetics). Most easy-to-understand case would be that from the dietary habits of individuals. It is well-known that the "meaty diet" readily elevates the serum creatinine level as well as the waist circumference (Khodayari S et. al., 2022, BMC Research Notes, and Pimenta E et. al., 2016, J Clin Epidemiol.). Hence, the higher the meat consumption, the higher the creatinine and wider the waist will be - a positive phenotypic correlation induced by the confounding effect of meat ingestion. We argue that the environmental effect would be exaggerated considering the relatively low heritability estimate of Creatinine.

We have now added the implication in our Results section:

Page 14: "Such discrepancies between the correlation estimates are possibly derived from the shared environmental factors between a pair of traits, such as dietary habits, that overwhelm the genotypic effects [12, 13]. This proves that the phenotypic correlation is not a mere proxy for the genetic correlation and consideration on the environmental effect is indispensable for the accurate interpretation of human phenomics [14]."

13. Were there any other surprising or unexpected correlations, and what are their

potential implications?

⇒ Yes, there were a few surprising correlations, and we would like to emphasize the following:

1) Utility of Secondary Body Measures: WHtR, WWtR, BMI

WHtR (Waist-to-height Ratio) and WWtR (Waist-to-weight Ratio) are secondary body measures that come from combining two bodily measures, similar to Body Mass Index (BMI). Our phenomics results also depict distinguishable patterns of association between these secondary body measures with other phenotypes. WHtR has a causal relationship with the C-reactive protein (CRP), Fat percentage, and HDL. On the other hand, WWtR showed associations with measures of lung capacity (FEV1 and FVC), liver function (AST), and inflammation (U_WBC). BMI positions as an intermediate phenotype, largely sharing its associations with WHtR and lightly with WWtR via left naked eyesight. These may reflect distinct biological mechanisms between the measurements warranting further studies. For example, WHtR is a well-known indicator of central adiposity, which serves as a better estimate obesity and related morbidities than BMI (Lee CMY et. al., 2008, J Clin Epidemiol.).

2) Complementary markers for cancer diagnosis: ALP and Amylase

CEA is a well-known biomarker for colon and lung cancer, while CA125 for ovarian. We found two independent causal relationships from these cancer biomarkers with other serum proteins that are seemingly irrelevant to cancerous phenotypes. Interestingly, our result suggests that alkaline phosphatase has influence on CEA. This implies that inflammation and dysfunction in the organs, such as liver or colon, precedes the alteration in the level of the cancer biomarker. It may be possible that ALP and CEA can both be used for detecting the presence of cancerous cells. Many early studies have reported their value in diagnosing cancer and monitoring metastasis in various types of cancer, including liver and colon, validating our finding (Aabo K et. al., 1986, Eur J Cancer Clin Oncol., Tartter PI et. al., 1981, Ann Surg., and Walach N et. al., 1989, J Surg Oncol.). Similarly, we could establish relationship between serum Amylase and CA125. Interestingly, we found cis-eQTLs underlying their pleiotropy are associated with the level of AMY2B expression in pancreas. Surprisingly, there have already been reports that patients with ovarian cancer manifest hyperamylasemia (Guo S et. al., 2018, Medicine, Shintani D, 2016, Eur J Gynaecol Oncol, and Zakrzewska I and Pietryńczak M, 1995). We argue that serum Amylase can be used as a complementary marker for ovarian cancer similar to ALP.

We elaborated our interesting correlations and their implications both in the Results and Discussion:

Page 18: "In our casual diagram (Figure 5, blue arrows), ALP and CEA showed potential causality, along with the shared genetic variants between them (pleiotropy near ABO gene). Numerous previous studies have consistently reported these markers together for diagnosing cancer and monitoring metastasis [21-23]. Similarly, CA125 and Amylase also displayed causality via shared genetic variants (pleiotropy near AMY2B gene). We propose that CA125 and Amylase might serve as complementary biomarkers for ovarian cancer, much like ALP and CEA. The biological relationships between these clinical blood measures remain unclear."

Page 18: "Our phenomics results also depicted distinguishable patterns of association between secondary body measures, such as WHtR (Waist-to-Height Ratio), WWtR (Waist-to-Weight Ratio), and BMI (Body Mass Index), with other phenotypes. WHtR exhibited a causal relationship with CRP (C-reactive protein), body fat percentage, and HDL. The result is concordant with previous reports that body fat percentage and CRP are correlated [24, 25]. Conversely, WWtR had casual associations with measures of lung capacity (FEV1 and FVC), liver function (AST), and inflammation (U_WBC). However, WWtR has yet to be proven its utility in clinical studies. BMI serves as an intermediate phenotype, sharing most of its associations with WHtR and, to a lesser extent, with WWtR via left naked eyesight. These findings suggest that the measurements reflect distinct biological mechanisms, warranting further studies. For instance, WHtR is a well-known indicator of central adiposity which provides a better estimate of obesity and related morbidities than BMI [26]."

	<p>Page 20: "Nevertheless, our findings bear important practical implications. We described the utility of secondary body measures, such as WHtR and WWtR, compared to BMI. We also elaborated on the diagnostic and prognostic value of other serum proteins, namely ALP and Amylase, in conjunction with the existing cancer biomarkers."</p> <p>14. You mentioned that phenomics analyses were not powerful enough for novel discoveries. Could you elaborate more on what would be needed to make them more effective? ⇒ Yes, we can elaborate on a few points of improvement for a more effective study in the future. The current dataset of the Korea4K includes 4,157 healthy samples with no apparent disease onset at the time of collection. Therefore, we could not see if our clinical variables, found from the phenomics analyses, have association on pathological conditions that is medically important. In future, we could use a wider variety of health-related categories to conduct a more powerful study, validating the current results with an enhanced scope which would bear invaluable medical and practical implications. Furthermore, collection of more samples for sequencing and health record data is also required for better chance of discovering new relationships. We have added the following sentence in our Discussion: Page 21: "However, we plan to collect more samples for sequencing and health record data with a wider variety of health-related categories to conduct a more powerful study in the future. This will allow us to not only validate our findings but also find correlations of medical importance that were missed in the present study."</p> <p>15. For the future implications, in terms of healthcare and personalized medicine, what do you see as the most immediate applications of the Korea4K dataset? ⇒ Thank you for asking these important points. As we mentioned in "Potential Implications", the Korea4K dataset contains both whole-genome scale genotypes and matched clinical information. Thus, the dataset can immediately be applied to discover novel genetic markers that are associated with several phenotypes, diseases, or drug responses for Korean and East Asians. As a reference panel, the expanded genotype dataset (1K to 4K) can support more accurate genotyping imputation which is essential for DNA chip-based genotyping that is still widely used for healthcare (such as genetic tests). Furthermore, the Korea4K dataset can be used as a control data set across many different studies if the proper control samples are not applicable.</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes

<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>Yes</p>

1 **Korea4K: whole genome sequences of 4,157 Koreans with 107** 2 **phenotypes derived from extensive health check-ups**

3

4 **Authors**

5 Sungwon Jeon^{1,2,†}, Hansol Choi^{1,3,†}, Yeonsu Jeon^{1,2}, Whan-Hyuk Choi^{1,3,4}, Hyunjoo Choi^{1,3},
6 Kyungwan An^{1,3}, Hyojung Ryu^{1,2}, Jihun Bhak^{1,3}, Hyeonjae Lee^{1,3}, Yoonsung Kwon^{1,3}, Sukyeon
7 Ha^{1,5}, Yeo Jin Kim², Asta Blazyte^{1,3}, Changjae Kim², Yeonkyung Kim², Younghui Kang^{1,2}, Yeong
8 Ju Woo², Chanyoung Lee^{1,3}, Jeongwoo Seo^{1,3}, Changhan Yoon^{1,3}, Dan Bolser⁶, Orsolya Biro⁷,
9 Eun-Seok Shin⁸, Byung Chul Kim², Seon-Young Kim⁹, Ji-Hwan Park⁹, Jongbum Jeon⁹, Dooyoung
10 Jung³, Semin Lee^{1,3,*}, and Jong Bhak^{1,2,3,10,*}

11

12 ¹Korean Genomics Center (KOGIC), Ulsan National Institute of Science and Technology (UNIST),
13 Ulsan 44919, Republic of Korea

14 ²Clinomics Inc., Ulsan 44919, Republic of Korea

15 ³Department of Biomedical Engineering, College of Information-Bio Convergence Engineering,
16 Ulsan National Institute of Science and Technology (UNIST), Ulsan 44919, Republic of Korea

17 ⁴Department of Mathematics, Kangwon National University, Chuncheon 24341, Republic of
18 Korea

19 ⁵Department of Computer Science & Engineering (CSE), College of Information-Bio
20 Convergence Engineering, Ulsan National Institute of Science and Technology (UNIST), Ulsan
21 44919, Republic of Korea

22 ⁶Geromics Ltd., 222 Mill Road, Cambridge, CB1 3NF, United Kingdom

23 ⁷Clinomics Europe Ltd., Budapest 1094, Hungary

1 ⁸ Department of Cardiology, Ulsan University Hospital, University of Ulsan College of
2 Medicine, Ulsan, 44033, Republic of Korea

3 ⁹Korea Bioinformation Center, Korea Research Institute of Bioscience and Biotechnology,
4 Daejeon 34141, Republic of Korea.

5 ¹⁰Personal Genomics Institute (PGI), Genome Research Foundation (GRF), Osong, 28160,
6 Republic of Korea

7 † These authors contributed equally to this work.

8 * Correspondence to: seminlee@unist.ac.kr (S.L.), and jongbhak@genomics.org (Jong B.
9 Sungwon Jeon [0000-0002-2729-9087]; Hansol Choi [0000-0002-3653-1474]; Yeonsu Jeon
10 [0000-0003-4560-4142]; Whan-Hyuk Choi [0000-0003-1587-5320]; Hyunjoo Choi [0000-0003-
11 2843-0419]; Kyungwhan An [0000-0001-9595-5130]; Hyojung Ryu [0000-0002-2276-850X];
12 Jihun Bhak [0009-0009-1965-8084]; Hyeonjae Lee [0000-0003-2683-3804]; Yoonsung Kwon
13 [0000-0002-7242-2734]; Sukyeon Ha [0009-0007-8192-3921]; Yeo Jin Kim [0000-0001-9122-
14 914X]; Asta Blazyte [0000-0001-7309-1482]; Changjae Kim [0000-0002-2018-1561];
15 Yeonkyung Kim [0000-0002-0649-916X]; Younghui Kang [0000-0002-0323-0851]; Yeong Ju
16 Woo [0000-0002-3049-7762]; Chanyoung Lee [0000-0002-8269-3191]; Jeongwoo Seo [0000-
17 0002-9163-4310]; Changan Yoon [0000-0003-0243-9853]; Dan Bolser [0000-0002-3991-
18 0859]; Orsolya Biro [0000-0002-4300-3602]; Eun-Seok Shin [0000-0002-9169-6968]; Byung
19 Chul Kim [0000-0002-4891-9679]; Seon-Young Kim [0000-0002-1030-7730]; Ji-Hwan Park
20 [0000-0002-6988-6239]; Jongbum Jeon [0000-0003-3533-1363]; Dooyoung Jung [0000-0002-
21 5381-4847]; Semin Lee [0000-0002-9015-6046]; Jong Hwa Bhak [0000-0002-4228-1299].

1 **Abstract**

2 **Background**

3 Phenome-wide association studies (PheWASs) have been conducted on Asian populations,
4 including Koreans, but many were based on chip or exome genotyping data. Such studies have
5 limitations regarding whole-genome-wide association analysis, making it crucial to have genome-
6 to-phenome association information with the largest possible whole-genome and matched
7 phenome data to conduct further population-genome studies and develop healthcare services based
8 on population genomics.

9 **Results**

10 Here, we present 4,157 whole-genome sequences (Korea4K) coupled with 107 health check-up
11 parameters as the largest genomic resource of the Korean Genome Project. It encompasses most
12 of the variants with allele frequency > 0.001 in Koreans, indicating that it sufficiently covered
13 most of the common and rare genetic variants with commonly measured phenotypes for Koreans.
14 Korea4K provides 45,537,252 variants, and half of them were not present in Korea1K (1,094
15 samples). We also identified 1,356 new geno-phenotype associations which were not found by the
16 Korea1K dataset. Phenomics analyses further revealed 24 significant genetic correlations, 14
17 pleiotropic associations, and 127 causal relationships based on Mendelian randomization among
18 37 traits. In addition, the Korea4K imputation reference panel, the largest Korean variants
19 reference to date, showed a superior imputation performance to Korea1K across all allele
20 frequency categories.

21 **Conclusions**

22 Collectively, Korea4K provides not only the largest Korean genome data but also corresponding
23 health check-up parameters and novel genome-phenome associations. The large-scale pathological

1 whole-genome-wide omics data will become a powerful set for genome-phenome level association
2 studies to discover causal markers for the prediction and diagnosis of health conditions in future
3 studies.

4

5 **Keywords**

6 Korean Genome Project, Genome, Phenome, Population genomics, Variome

7

1 **Background**

2 South Korea has perhaps one of the most extensive and convenient annual health check-up services.
3 Every year, almost all Koreans aged over 40 receive a standardized health check-up, yielding a
4 wealth of individual clinical data [1]. In 2020, we published 1,094 whole genomes with clinical
5 information (Korea1K) by providing all the participants with a free standard health check-up
6 showing the value of whole-genome data accompanied by clinical information mapping the
7 genome diversity with practical applications [2]. Here, we present the second phase of the Korean
8 Genome Project (KGP) with 4,157 sets of whole-genome data, Korea4K. It is accompanied by 107
9 types of clinical traits that have been donated by 2,685 healthy participants who acquired the health
10 check-up reports from the hospitals of their choice in the past years. We manually annotated
11 thousands of donated health reports that are matched with the whole-genome information.
12 Therefore, apart from the increased number of samples, the main difference between Korea1K and
13 Korea4K is that Korea4K's clinical information is from very heterogeneous but fairly standard
14 Korean health check-up centers, while Korea1K was from one very well-controlled university
15 hospital health check-up center. This was also a testbed to assess how difficult it would be to merge
16 data from the heterogeneous health check-up record system in a nation for a large-scale genome
17 to phenome association analysis.

18 Previously, there were a few phenome-wide association studies (PheWASs) on Asian populations,
19 but they were limited to chip or exome-based genotyping data. A Japanese PheWAS identified the
20 genetic links among clinical traits, complex diseases, and cell-type specific patterns [3]. Another
21 PheWAS using 10,000 Korean cohorts' health check-up data from multiple lab sources showed
22 network relationships between genes and phenotypes [4]. However, none of these studies covered
23 the entirety of genomic variation, and they have limitations on genome-wide data analyses [5, 6].

1 A scientific contribution of this version of KGP is that we provide extensive genome-to-phenome
2 association information with the largest genomic and clinical data from Korea to date to estimate
3 how many samples and clinical parameters cover the whole genomic and common phenotypic
4 diversity of Koreans. Korea4K contains 4,157 Korean genomes from East Asian ancestry, and
5 2,685 of them are accompanied by 107 types of clinical information such as height, waist
6 circumference, weight, albumin/globulin ratio, basophil, direct bilirubin, low-density lipoprotein,
7 high-density lipoprotein, mean corpuscular volume, and total cholesterol. The rest does not contain
8 such kind of data because the biobank does not have phenotype information, or we were not able
9 to collect it from the participants. Korea4K extends the efforts to completely map the totality of
10 Korean genomic diversity, which can be a useful scope reference for disease risk prediction,
11 diagnosis, and treatments in the future for personalized medicine.

12 As the second phase of the KGP, Korea4K not only extends the previously reported Korea1K [2]
13 but also includes new multi-phenotypic association analyses, that is, analyses on markers that are
14 associated with multiple phenotypes (pleiotropy), the genetic correlation between traits, and
15 estimated causality relationship among traits through Mendelian randomization (MR) and 3D
16 structure models for Korean specific missense variants. Combining these two omics data, we
17 provide the community with the most extensive geno-phenotype association of healthy Korean
18 participants. We have also applied the genomic variation data to the genotype imputation of low-
19 frequency variants in the Korean population.

20

1 **Data Description**

2 The goal of our project was to create a genome dataset for Korea4K, which included newly
3 sequenced genomic data from 2,848 participants as well as 1,309 whole-genome sequencing
4 (WGS) datasets from Korea1K and public data archives. Additionally, we established a phenome
5 dataset for Korea4K by gathering or computing 107 clinical parameters and genome data from
6 2,685 samples. We collected a total of 3,383 clinical datasets, including multiple time points per
7 sample, from regular health checkups conducted by various hospitals and clinics across Korea
8 between 2016 and 2019. The genome and phenome datasets were produced and curated by the
9 protocol in Material and Methods.

10

11 **Analyses**

12 **The largest Korean whole-genome variants data: Korea4K variome**

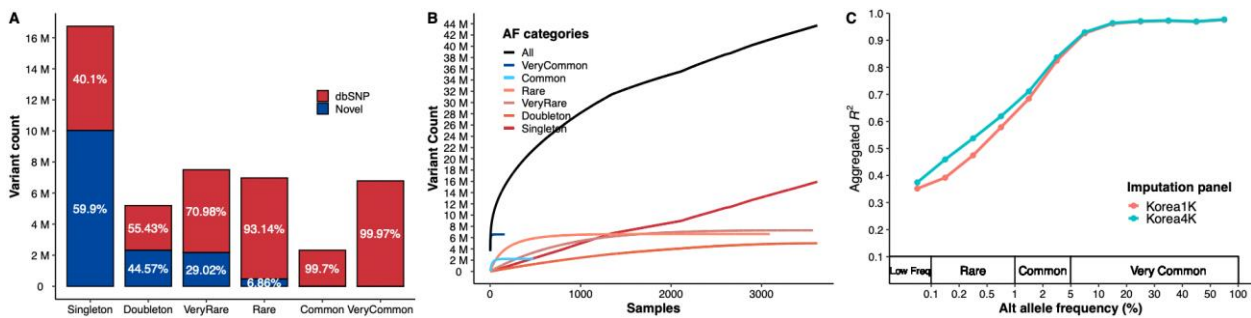
13 A total of 64,301,272 single nucleotide variants (SNVs) and 8,776,608 Indels were called against
14 the human genome reference (hg38) from the 4,157 Korean whole genomes, including 3,071
15 healthy controls (Supplementary Table S1 and S2). It contains 3,063 newly added whole genomes
16 sequenced by Illumina next-generation sequencing (NGS) platforms (HiSeq X10 and Novaseq
17 6000), in addition to the previous Korea1K dataset which was mostly generated by Illumina HiSeq
18 X10. Using the variant data, we selected 3,617 samples with no kinship after initial sample filtering
19 (see Methods). To exclude erroneous variants from sequencing batch effects from the
20 heterogeneous Illumina NGS platforms and library preparation, we applied an allele balance bias
21 measurement and finally acquired 12,713,580 erroneously called variant candidates
22 (Supplementary Fig. S1). After additional variant filtering (see Methods), we identified 45,537,252
23 variants including 42,124,137 SNVs, 36,029 double nucleotide variants (DNVs), 26,135 triple

1 nucleotide variants (TNVs), 3,261,682 indels, and 89,269 other types of small variants from the
2 3,617 unrelated samples. We named this filtered Korean dataset the Korea4K variome (Figure 1).
3 A total of 23,689,147 variants were not present in the previous Korea1K variome. This
4 unexpectedly large difference is likely derived from different batch effect filtering, and variant
5 calling and filtering procedures, as well as new variants from the larger sample size. Consistent
6 with the Korea1K study [2], most variants were located in intronic or intergenic regions and rarely
7 in splicing sites or coding regions (Supplementary Fig. S2), which is a sign of negative selection
8 pressure in the population. Half of the total variants (21,941,879; 48.2%) were singleton or
9 doubleton in the 3,617 unrelated samples, indicating that the Korean population's genetic diversity
10 is very low as the population diversity could be covered by fewer than 4,000 unrelated samples
11 (Figure 1A, Supplementary Table S3). Almost all the common (allele frequency of > 0.01 and
12 allele frequency of ≤ 0.05) and very common (allele frequency of > 0.05) variants were found to
13 be already reported in the dbSNP database (99.70% and 99.97%, respectively), while more than
14 half of the singleton and doubleton variants were newly discovered in this study (59.9% and
15 44.57%, respectively), indicating the new variant pool is well-exhausted in the Korean population
16 by the 3,617 samples resulting in a large portion of individual specific novel variants in the Korean
17 variome (Figure 1A, Supplementary Table S3). Only 3,092 and 3,569 unrelated individuals were
18 needed to discover all the rare (allele frequency of > 0.001 and allele frequency of ≤ 0.01) and
19 very rare (allele count of > 2 and allele frequency of ≤ 0.001) variants in the Korea4K variome,
20 respectively (Figure 1B) indicating that the Korea4K variome includes almost all the rare and very
21 rare variants of Korean people of East Asian ancestry. It is notable that in our previous Korea1K
22 data, the accumulated variant number curves did not reach a plateau [2]. Regarding common
23 variants, only 481 and 161 unrelated individuals were necessary for common and very common

1 variants, respectively, to cover the diversity which is close to the Korea1K statistics (440 and 132
 2 samples). Essentially, the Korea4K variome statistics indicate the saturation of population
 3 diversity detection among Koreans. However, as expected, in the case of singleton and doubleton
 4 variants, the Korea4K variant discovery curve did not reach a plateau. This is due to each
 5 individual's novel random variants and we will never reach a point of no novel variant discovery
 6 even with increased sample numbers.

7 As a practical application, we constructed a Korea4K imputation reference panel from the 3,614
 8 unrelated whole-genomes that showed a consistently better imputation performance than the
 9 Korea1K. The Korea4K panel was able to impute 198,805 more genotypes than the Korea1K panel
 10 (7,551,095 loci compared to 7,352,290) with the same dataset. Moreover, as expected, the
 11 Korea4K panel had better accuracy across all allele frequency categories than the Korea1K panel
 12 (Figure 1C). The difference in aggregated R^2 became larger for variants with allele frequency (AF)
 13 in Korea4K < 0.05 than for those in Korea1K, indicating higher accuracy in rare variants (Figure
 14 1C). In particular, the Korea4K imputation panel improved the imputation accuracy by 6% for the
 15 rare variants group compared to Korea1K on average.

16



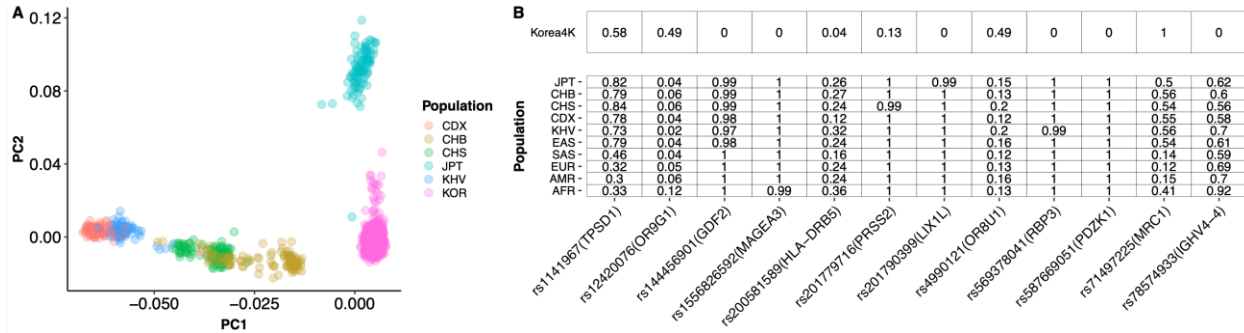
17

18 **Figure 1 Korean variome profile and imputation evaluation using Korea4K** (A) The number
 19 of variants in the Korea4K variome is categorized by allele frequencies (AFs) among unrelated
 20 Korea4K genomes. “dbSNP” indicates the variants were reported in dbSNP database. “Novel”

1 indicates the variants were not reported in dbSNP. Singleton, allele count = 1; doubleton, allele
2 count = 2; very rare, allele count of > 2 and allele frequency of ≤ 0.001 ; rare, allele frequency of
3 > 0.001 and allele frequency of ≤ 0.01 ; common, allele frequency of > 0.01 and allele frequency
4 of ≤ 0.05 ; very common, allele frequency of > 0.05. **(B)** The number of discovered variants as a
5 function of unrelated genomes. **(C)** Imputation performance evaluation using the Korea4K and
6 Korea1K panels. The X-axis indicates alternative (Alt) allele frequency in the Korea4K variome.
7 The Y-axis represents the aggregated R^2 values of variants. We used variants that were overlapped
8 by imputed results across two panels.

9
10 As in Korea1K, the Korean population is genetically distinct from the Chinese and Japanese
11 populations, confirmed by principal component analysis (PCA) with few outliers (Figure 2A). We
12 also found 62 missense variants out of 282,607 in Korea4K that had AFs significantly different
13 from ten populations in the 1000 genome project (1KGP) from European Bioinformatics Institute
14 (EBI), Cambridge, UK (Chi-squared test $P < 5 \times 10^{-5}$ against each of the ten populations, see
15 Methods; Supplementary Table S4). The genes containing such Korean-specific missense variants
16 included *LILRB3*, *HLA-DRB5*, *IGLV5-48*, and *IGHV4-4* that are known to be associated with
17 adaptive immunity, and *OR9G1* and *OR8U1* for olfactory receptors. Additionally, we found that
18 twelve Korean-specific missense variants were in protein functional domains (Figure 2B). Four of
19 them were predicted to facilitate increased structural stability calculated in the protein 3D models
20 built by AlphaFlod2 [7], while the other eight variants were predicted to cause decreased stability
21 (Supplementary Table S5).

1



2

3 **Figure 2 Comparison of Korea4K and 1KGP (A)** The results from principal component analysis
 4 of Korea4K and the 1KGP set of East Asian samples. **(B)** Allele frequency information of Korea4K
 5 and the populations in the 1KGP for the twelve Korean-specific missense variants located in
 6 protein functional domains. KOR: Korea4K; CDX: Dai Chinese; CHB: Han Chinese; CHS:
 7 Southern Han Chinese; JPT: Japanese; KHV: Kinh Vietnamese; EAS: East Asians; SAS: South
 8 Asians; EUR: European; AMR: American; AFR: African.

9

10

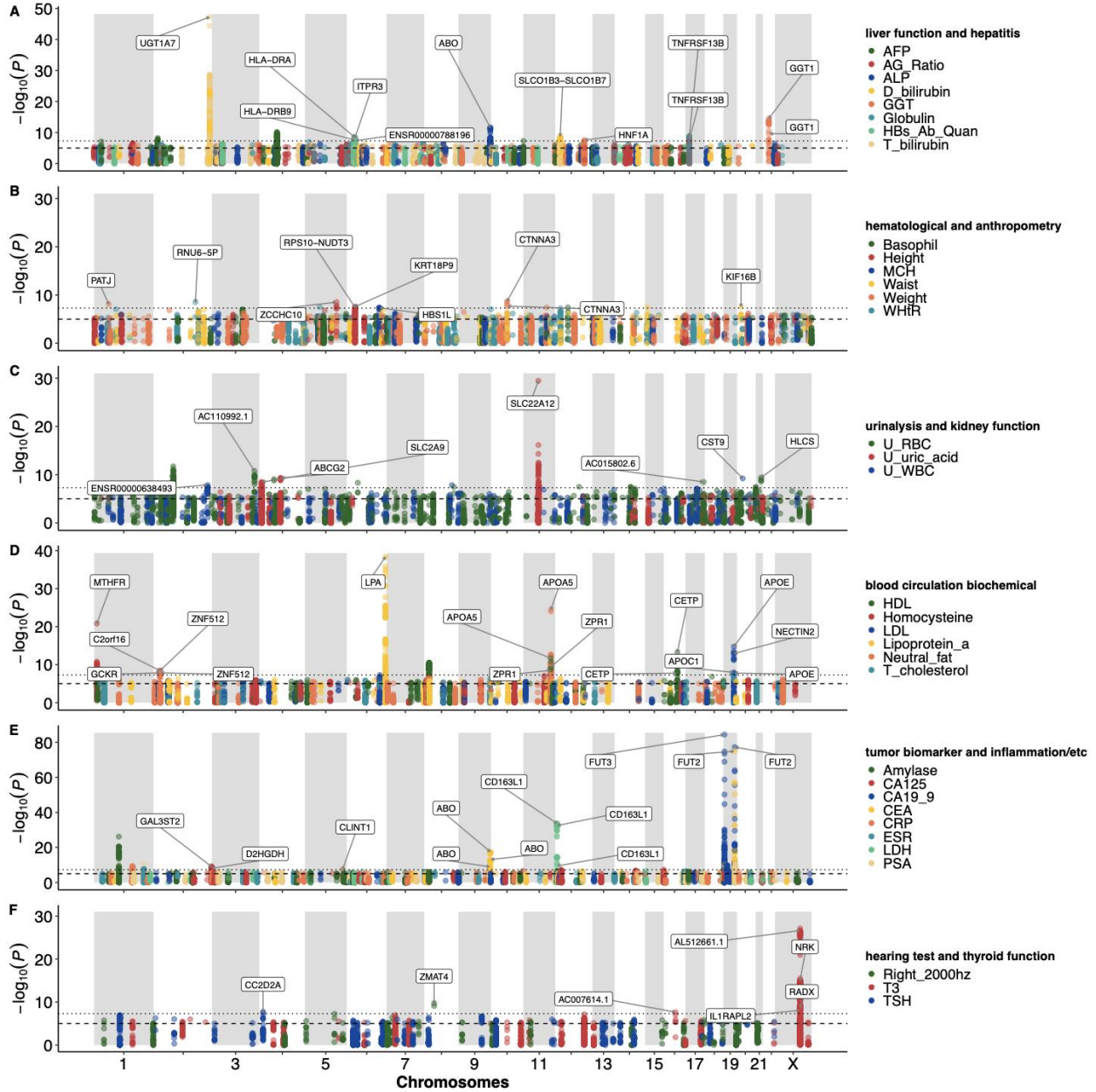
11 **Whole-genome-wide association study (WGWAS)**

12 Whole-genome-wide association studies (WGWASs) revealed that 2,324 variants from 157 unique
 13 loci had significant associations with 34 clinical traits from 37 WGWAS target traits ($P < 5 \times 10^{-8}$;
 14 Figure 3A-F, Supplementary Table S6). Among the significantly associated variants, 2,314
 15 variants from 30 clinical traits still showed significance after false discovery rate (FDR) correction
 16 using the Benjamini-Hochberg approach ($FDR < 0.05$). We used 90 clinical traits from the 107
 17 phenotypes after filtering 27 traits with a high missing rate and biased distribution for WGWASs
 18 (see Methods). Of the 90 traits, 54 were not confident in Quantile-Quantile plots and were excluded
 19 from further Mendelian randomization and pleiotropy analyses (see Methods). Among the 2,324

1 WGWAS significant variants, only 85 variants (31 loci) were reported in the GWAS catalog
2 database [8]. The trait with the largest number of significantly associated loci was carbohydrate
3 antigen 19-9 (CA19-9), a cancer antigen, with sixteen loci. Uric acid had the second highest
4 number of significant loci with fourteen loci.

5 Korea4K showed much stronger statistical power than the previous Korea1K study, identifying
6 1,356 new WGWAS variants (107 loci) from 28 common traits between Korea4K and Korea1K.
7 Also, Korea4K had much lower (i.e., more significant) P -values than Korea1K for all the
8 commonly found association variants between the two datasets (Supplementary Fig. S3). Among
9 the 107 loci containing the 1,356 new WGWAS variants, 798 Korea4K significant WGWAS
10 variants from 73 loci had not been significant in Korea1K (Supplementary Table S6). Furthermore,
11 twelve traits (albumin/globulin ratio, basophil, C-reactive protein, direct bilirubin, height, low-
12 density lipoprotein, mean corpuscular volume, right hearing at 2000hz, thyroid stimulating
13 hormone, total cholesterol, waist, weight) had 425 WGWAS variants that were significant
14 uniquely in Korea4K, meaning no significant WGWAS variants from the twelve traits in Korea1K
15 (Supplementary Table S6). For example, a missense variant, rs6431625 ($P = 1.41 \times 10^{-23}$, FDR =
16 5.23×10^{-18}), in *UGT1A3* was found to be associated with direct bilirubin in Korea4K. It was
17 previously reported to be associated with circulating bilirubin levels [9]. Another Korea4K-
18 specific missense variant is rs7412 ($P = 2.86 \times 10^{-14}$, FDR = 1.11×10^{-7}) in *APOE* which is
19 associated with low-density lipoprotein (LDL) levels. Its association with cholesterol levels has
20 been previously well-established [10]. Finding novel WGWAS variants in Korea4K was due to
21 the increased sample size and subsequently increased variant number compared to Korea1K.

1



2

3 **Figure 3 Whole-genome-wide association studies in Korea4K. (A-F)** Whole-genome-wide

4 association studies from 34 traits. Loci are presented only when index variants of the loci had

5 significant P -value ($P < 5 \times 10^{-8}$) from the WGWas. The dashed line indicates the suggestive

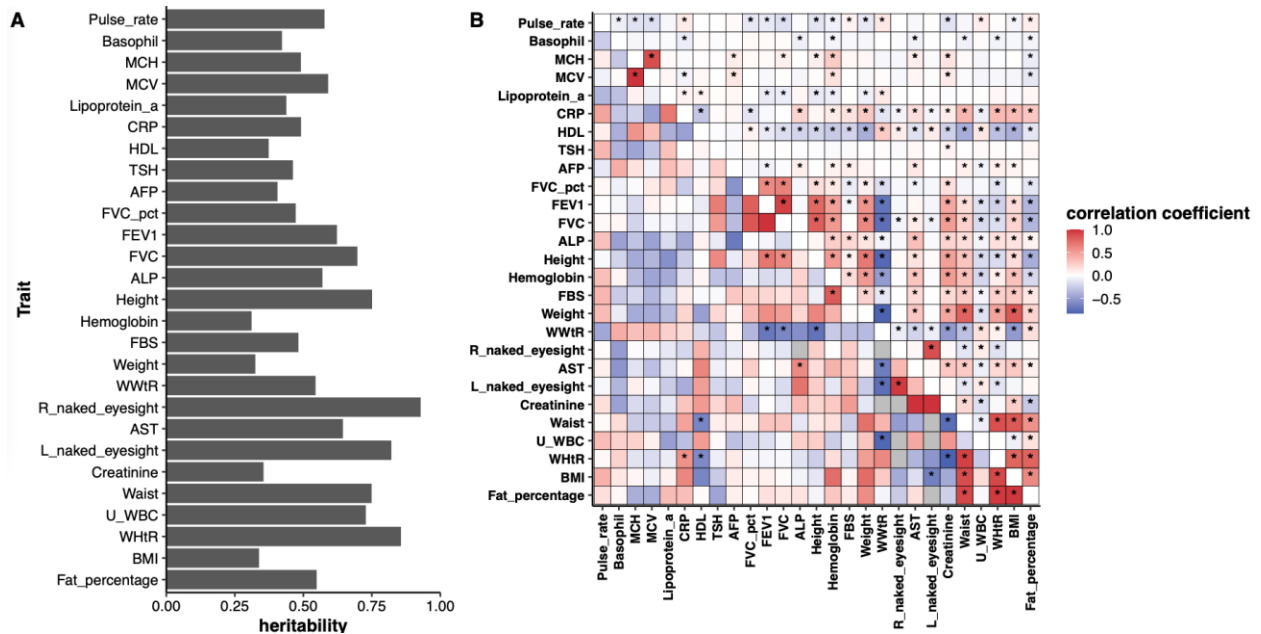
6 threshold ($P < 10^{-5}$). The dotted line indicates the significant threshold ($P < 5 \times 10^{-8}$).

7

1 **Genetic correlation (GC) and phenotypic correlation (PC)**

2 We found 27 traits with significant heritability among 89 quantitative traits (Figure 4A; the lower
3 boundary of genetic heritability > 0 with 95% confidence interval; Supplementary Table S7). A
4 total of 24 pairs of traits showed a significant genetic correlation ($FDR_{GC} < 0.05$), measured as r_G
5 value, among 351 trait pairs between the 27 traits that showed significant heritability (Figure 4,
6 Supplementary Table S8). We found consistent results of Weight-Waist and body mass index
7 (BMI)-Waist pairs, showing a significant genetic correlation in the UK Biobank data with the same
8 trend as our result ($r_G = 0.9$, $P = 10^{-308}$ in UK Biobank; $r_G = 0.9$, $P = 10^{-308}$ in UK Biobank,
9 respectively) [11]. We identified 2,274 trait-trait relationships that had significant phenotypic
10 correlation ($FDR_{PC} < 0.05$, its 95% CI does not include 0) from trait-trait associations between
11 3,916 pairs of 89 quantitative traits (Figure 4B, Supplementary Table S9). Most genetic and
12 phenotypic correlations showed the same direction of correlation. The only two exceptions were
13 waist/weight ratio (WWtR) – Urine white blood cell (U_WBC) and Waist-Creatinine which
14 showed opposite directions. This trend of Waist-Creatinine has also been reported in a correlation
15 database using UK-biobank data [12]. Such discrepancies between the correlation estimates are
16 possibly derived from the shared environmental factors between a pair of traits, such as dietary
17 habits, that overwhelm the genotypic effects [13, 14]. This proves that the phenotypic correlation
18 is not a mere proxy for the genetic correlation and consideration on the environmental effect is
19 indispensable for the accurate interpretation of human phenomics [15].

20



1
 2 **Figure 4 Genetic correlation and Phenotypic correlation in Korea4K.** (A) Genetic heritability
 3 of 27 traits that showed at least a marginal significance. (B) Genetic correlation and phenotypic
 4 correlation between the 27 traits. The upper triangle indicates phenotypic correlation coefficient
 5 (Pearson's) and lower triangle indicates genetic correlation coefficient (rG).

6
 7 **Pleiotropy and Mendelian randomization (MR)**

8 Out of the 37 GWAS target traits, we detected 1,131 variants from 21 traits having suggestive
 9 associations ($P_{GWAS} < 10^{-5}$) with at least two traits, indicating pleiotropic variants (Figure 5, red
 10 edges; Supplementary Table S10). We devised the Variant-Sharing Index (VSI) to measure the
 11 degree of intersection between two phenotypes (Table 1; see Methods). A VSI of zero signifies
 12 that two traits share no suggestively associated variants (SSVs), while 100 indicates they share all
 13 of them. The trait pairs with SSVs and the corresponding VSIs are listed in Table1. Notably, we
 14 had only one variant, rs77913154, that was shared among three traits: Globulin, AG_Ratio, and
 15 ESR (Supplementary Table S10). Interestingly, we found fifteen variants residing on *SOD2P1-*

1 *AC095032.2-AC095032.1* locus forming pleiotropy between the serum amylase level and the level
2 of CA125, a known ovarian cancer marker (Table 1, VSI=2.3). Fourteen variants of the fifteen
3 variants conform to the alteration of *AMY2B* expression level, as per cis-eQTL results from the
4 GTEx Portal (ver. 8), four of which were associated with expression in the pancreatic tissue (see
5 Methods). There have already been reports of hyperamylasemia in patients with ovarian cancer
6 [16-18]. In addition to the investigation on the general pleiotropic relationship, we employed
7 Mendelian Randomization (MR) to detect vertical pleiotropy that can assert the direction of the
8 phenotypic relationships [19]. This provides indirect evidence implying causality between the
9 traits to discern spurious phenotypic associations, such as confounding and collider bias [20, 21].
10 We found a total of 127 trait pairs among 1,332 pairs of the 37 GWAS traits were estimated to
11 have significant causal relationships (FDR < 0.05, Figure 5, Supplementary Table S11). These
12 findings were supported by at least two of three different Mendelian randomization (MR) analysis
13 methods (IVW: 166 pairs; MRPRESSO: 139; MR-Egger: 23). Among these, 59 trait pairs showed
14 unidirectional relationships while 68 exhibited bidirectional causal relationships (Supplementary
15 Table S11).

16
17
18
19
20
21
22
23

1 **Table 1 Pleiotropic associations and Variant-Sharing Index (VSI)**

Trait1	Trait2	Suggestive variants in trait1	Suggestive variants in trait2	Shared variants	Total variants	VSI
D_bilirubin	T_bilirubin	638	632	569	701	81.2
Globulin	AG_Ratio	294	230	147	377	39
HDL	Neutral_fat	348	398	191	555	34.4
CEA	CA19_9	221	264	74	411	18
T_cholesterol	LDL	74	238	38	274	13.9
WHtR	Waist	177	100	31	246	12.6
ALP	CEA	153	221	35	339	10.3
T3	GGT	542	125	23	644	3.6
CA125	Amylase	202	466	15	653	2.3
Weight	Waist	123	100	5	218	2.3
Height	Weight	173	123	2	294	0.7
ESR	AG_Ratio	163	230	1	392	0.3
Globulin	ESR	294	163	1	456	0.2
U_RBC	Globulin	627	294	1	920	0.1

2

3 **Summary results of the four phenomics analyses**

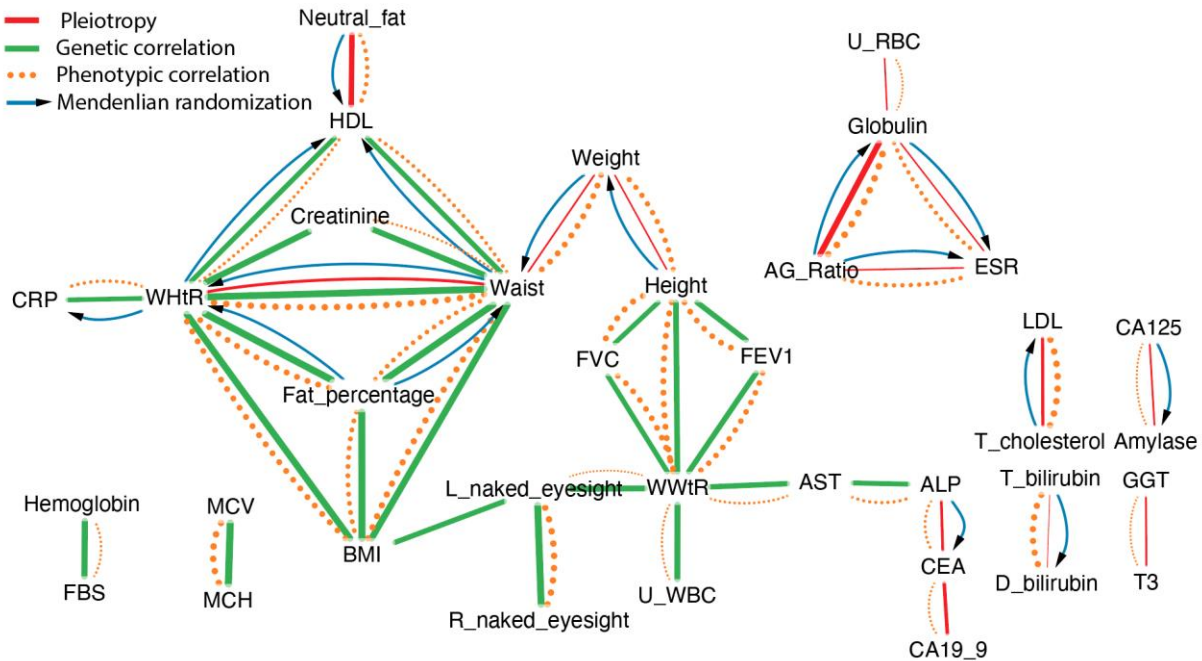
4 We summarized the results of four phenomics analyses (Genetic correlation, Phenotypic
5 correlation, Mendelian randomization, and pleiotropy) through visualizing them in network plots
6 (Figure5). In general, the identified trait-trait pairs of genetic correlation, Mendelian
7 randomization, and pleiotropy analyses did not often overlap. Genetic correlation and pleiotropy
8 are found to be exclusive of each other, even though both measures having shared genetic
9 components of two different traits. GC was primarily observed among body measures such as
10 waist circumference, weight, height, and left naked eyesight. On the other hand, pleiotropy was
11 more prevalent on the relationship between metabolites in blood such as LDL, bilirubin, or CEA.
12 The only overlap between these two was WHtR (Waist-to-Height Ratio)-Waist (circumference),
13 where one is derived from the other.

1 MR suggests a causal relationship between phenotypic correlations through mediation effect by a
2 genotype. In our causal diagram (Figure 5, blue arrows), ALP and CEA showed potential causality,
3 along with the shared genetic variants between them (pleiotropy near *ABO* gene). Numerous
4 previous studies have consistently reported these markers together for diagnosing cancer and
5 monitoring metastasis [22-24]. Similarly, CA125 and Amylase also displayed causality via shared
6 genetic variants (pleiotropy near *AMY2B* gene). We propose that CA125 and Amylase might serve
7 as complementary biomarkers for ovarian cancer, much like ALP and CEA. The biological
8 relationships between these clinical blood measures remain unclear.

9 Our phenomics results also depicted distinguishable patterns of association between secondary
10 body measures, such as WHtR (Waist-to-Height Ratio), WWtR (Waist-to-Weight Ratio), and BMI
11 (Body Mass Index), with other phenotypes. WHtR exhibited a causal relationship with CRP (C-
12 reactive protein), body fat percentage, and HDL. The result is concordant with previous reports
13 that body fat percentage and CRP are correlated [25, 26]. Conversely, WWtR had causal
14 associations with measures of lung capacity (FEV1 and FVC), liver function (AST), and
15 inflammation (U_WBC). However, WWtR has yet to be proven its utility in clinical studies. BMI
16 serves as an intermediate phenotype, sharing most of its associations with WHtR and, to a lesser
17 extent, with WWtR via left naked eyesight. These findings suggest that the measurements reflect
18 distinct biological mechanisms, warranting further studies. For instance, WHtR is a well-known
19 indicator of central adiposity which provides a better estimate of obesity and related morbidities
20 than BMI [27].

21

22



1
 2 **Figure 5 Graph visualization of genetic correlation, phenotypic correlation, pleiotropy, and**
 3 **Mendelian randomization.** Green line indicates significant genetic correlation (GC), and the edge
 4 thickness indicates the absolute value of the correlation coefficient. Red line indicates trait pairs
 5 that have pleiotropic variants. Dotted orange lines indicate phenotypic correlation (PC), and the
 6 edge thickness indicates the absolute value of Pearson’s correlation coefficient. Blue arrow line
 7 indicates a causal relationship from Mendelian randomization (MR). MR and PC were visualized
 8 only when at least one of GC or Pleiotropy relationships was observed between the traits.

9
 10 **Discussion**

11 Batch effect exacerbated by sequencing platform and library preparation bias is a critical problem
 12 in very large population genome association studies, especially with clinical data from
 13 heterogeneous health check-up centers. In the future, more and more diverse whole-genome data
 14 with extensive clinical data will be publicly available, and it is inevitable that they will be merged

1 for more precise whole genome-to-phenome association research. Korea4K is not an exception in
2 that regard, and in one homogeneous population WGWAS, it was necessary to consider and factor
3 in a great deal of sequencing and clinical data batch effects and errors. We attempted to minimize
4 the errors by using allele balance with optimal filtering criteria and time-consuming manual checks
5 on health reports that were donated by the participants (see Methods). The largest challenge of
6 Korea4K project was cleaning up heterogeneous clinical data from different health check-up
7 centers. Another major issue was that the health check-up data heterogeneity caused reduced
8 numbers of participants' common traits with which to compare. Some of the health data were from
9 past years' health check-ups from heterogeneous hospitals throughout South Korea. This
10 heterogeneity in location and time was not an intentional experimental design but was in order to
11 reduce the cost of performing expensive one-center health check-ups for the Korea4K participants.
12 Therefore, WGWAS along with standardized and unified national and public health check-up data
13 will greatly benefit future whole-genome-wide association studies.

14 Although 4,157 seems like a large number, we found the sample size in this study was still not
15 large enough to detect weak association signals. The Korea4K variome with matched phenotype
16 information has allowed us to estimate genomic correlation across various phenotypes using
17 GREML [28]. GREML has been reported to have higher accuracy compared to methods, such as
18 linkage disequilibrium score regression (LDSC), using only summary statistics from GWAS [29].
19 For example, the minimum heritability score was 0.34 (Degree of obesity) among the traits
20 detected as statistically significant. The statistical power of our maximum 2,685 subjects and FDR
21 < 0.05 is estimated to be 0.72 for detecting traits with heritability of 0.3 or higher (Calculated from
22 GCTA-GREML Power Calculator) [30]. This will increase to 0.97 with 4,000 subjects. In other

1 words, phenomics analyses were limited and not powerful enough to confidently discover novel
2 phenotypic associations with the current dataset.

3 Nevertheless, our findings bear important practical implications. We described the utility of
4 secondary body measures, such as WHtR and WWtR, compared to BMI. We also elaborated on
5 the diagnostic and prognostic value of other serum proteins, namely ALP and Amylase, in
6 conjunction with the existing cancer biomarkers. However, we plan to collect more samples for
7 sequencing and health record data with a wider variety of health-related categories to conduct a
8 more powerful study in the future. This will allow us to not only validate our findings but also find
9 correlations of medical importance that were missed in the present study. While chip-based GWAS
10 is a common approach, our study highlights the unique advantage of WGWAS (whole-genome-
11 wide association) in genotype-phenotype association studies. An illustrative advantage of
12 WGWAS is its whole-genome-wide, unbiased coverage of genetic variants, which allowed us to
13 assign specific variants accounting for pleiotropy. This was not achievable with conventional
14 methods. For example, we could identify the variants in the well-known pleiotropic relationships
15 such as ALP-CEA by *ABO* locus (35 variants), Neutral_Fat-HDL by *LPL* locus (181 variants) and
16 Total cholesterol-LDL by *TOMM40*, and *APOE* locus (4 and 2 variants, respectively)
17 (Supplementary Table S7). These loci and their corresponding trait pairs were previously reported
18 from chip-based GWAS summary results [31, 32]. Similarly, we anticipate the fine-mapping
19 analyses will also benefit from WGWAS, pinpointing novel genetic variants of phenotypic
20 importance, as demonstrated in our prior work [33]. Taken together, whole-genome sequencing
21 with its genomic completeness should be a well-considered choice for future genomic association
22 studies.

1 One of the main objectives of the Korea4K project was to build a genomic and phenomic reference
2 dataset to discover unknown whole-genome to phenome associations that can be detected from
3 samples of healthy people. This, however, is contradictory and it limited us in discovering clear
4 pathogenic associations because most of the participants examined in WGWAS were healthy
5 without any severe aberrant phenotypes or diseases that could bring us clues for interesting omics
6 analyses. Moreover, utilizing recently introduced human genome references like the T2T reference
7 [34] and Human Pangenome reference [35], which offer broader genomic coverage or have
8 population-specific sequences compared to the existing GRCh38 reference, could help identify
9 additional associations that might be overlooked. Nevertheless, these new references lack
10 functional annotations and need to be connected to previous databases such as dbSNP and the
11 GWAS catalog.

12 As for the future directions, there are several key limitations that have not been met in our current
13 study. The first is we failed to acquire long DNA sequencing reads from the healthy participants
14 for building a structural variation reference set for the Korean population. The second is the lack
15 of epigenomic data from the 4,157 samples. This was mostly due to high costs for generation and
16 computing long-read based assemblies and sequencing methylated DNA sites. The third one,
17 which is perhaps the most relevant for the purpose of performing association studies for healthcare
18 is that we failed to acquire more rare and severe disease data from patients, accompanied by precise
19 clinical and multiomics data. We have excluded a small number of rare disease cases, as those
20 required a large amount of sequencing data from genome, transcriptome, and methylome to
21 perform precise functional analyses. Large-scale pathological whole-genome-wide omics data will
22 become a powerful set for genome-phenome level association studies to detect causal markers for
23 the prediction and diagnosis of health conditions in future studies.

1

2 **Potential Implications**

3 The Korea4K dataset can be a valuable variome reference, as it contains matched phenome data
4 for personalized medicine, large-scale population genome studies, and the understanding of
5 anthropologic history in Korea. This large-scale Korean genome-phenome dataset can help
6 identify genetic basis for diseases and phenotypes, enabling personalized treatment plans for
7 individuals. Analyzing the genome-phenome association dataset can also be used to develop new
8 drugs that target specific genetic variations in the Korean population. The Korea4K dataset can
9 also be valuable for other populations, particularly East Asians, as it can be used to identify
10 population-specific genome-phenome patterns by comparing the population's genome-phenome
11 data to the Korea4K dataset. Furthermore, the Korea4K reference panel can be utilized for
12 genotype imputation of DNA chip genotyping data for the Korean population and other East
13 Asians.

14

15 **Materials and Methods**

16 **Sample collection and whole-genome sequencing**

17 We collected 2,848 blood samples or already processed DNA samples from Korean individuals.
18 A total of 1,094 whole-genome sequencing (WGS) datasets originating from our previous study
19 (Korea1K) and 215 WGS data from publicly available Clinical & Omics Data Archive (CODA)
20 were added to the aforementioned dataset [2]. The genomic DNA was extracted using the DNeasy
21 Blood & Tissue kit (Qiagen) from whole blood samples. We constructed the whole-genome
22 sequencing library from the DNA by using the TruSeq Nano DNA Sample Prep kit (Illumina) kit.
23 Whole-genome sequences of the 2,848 samples were generated by the Illumina Nova-seq 6000

1 platform. All the sequencing data that we used in this study had 151bp as a read length. Average
2 sequencing amount per sample was 27.75× (Supplementary Fig. S4).

3 **Joint genotype calling**

4 Adapter contamination was trimmed using Cutadapt (RRID:SCR_011841, ver. 1.9.1) [36] with a
5 forward adapter ('GATCGGAAGAGCACACGTCTGAACTCCAGTCAC') and reverse adapter
6 ('GATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT') and with a minimum read length of 50
7 bp after trimming (Supplementary Fig. S4). The quality of trimmed reads were screened by
8 FASTQC program (RRID:SCR_014583, ver. 0.11.5) We mapped the whole-genome sequencing
9 reads from 4,157 samples to the human reference genome (hg38) using BWA-mem
10 (RRID:SCR_010910, ver. 0.7.17) with the '-M' option and alt-aware mode [37]. The mapped
11 reads were sorted by genomic coordination using Picard (RRID:SCR_006525, ver. 2.20.3). We
12 marked the PCR-duplicates and recalibrated the base quality of the mapped reads using the
13 MarkDuplicates and BaseRecalibrator module in Picard (RRID:SCR_006525, ver. 2.20.3),
14 respectively. A total of 3,156 samples had a mapping depth of $\geq 20\times$ (Supplementary Fig. S4).
15 Individual genotypes were called in GVCF format by HaplotypeCaller in GATK
16 (RRID:SCR_001876, ver. 4.1.3) with '--genotyping-mode DISCOVERY -stand-call-conf 30 -
17 ERC GVCF' options [38]. We merged the individual genotypes to a single GVCF for each
18 chromosome using CombineGVCFs in GATK (RRID:SCR_001876, ver. 4.1.3) [38]. We jointly
19 genotyped the merged GVCF with the genotypeGVCF module in GATK (RRID:SCR_001876,
20 ver.4.1.3) [38]. Variant quality of the joint genotypes was recalibrated using the VQSR module in
21 GATK (RRID:SCR_001876, ver. 4.1.3) [38].

22 **Sample and variant filtering**

1 After joint genotyping, we filtered out a total of 540 participants on the criteria that are listed below
2 using SelectVariants in GATK (RRID:SCR_001876, ver. 4.1.3) with ‘--remove-unused-alternates’
3 option to remove unused variants [38]. To explore kinship relations among the samples, we
4 assessed Identical by Descent (IBD) using the Plink program (RRID:SCR_001757, ver. 1.90b3n)
5 [39]. Samples with a PI_HAT value exceeding 0.05 were considered to be in a kinship relation.

- 6 1. showing high missing genotype rate (>10%): nine samples
- 7 2. having too high or low heterozygous variants ratio compared to homozygous variants per
8 sample (3 s.d.): four samples
- 9 3. having relatedness to other samples: 428 samples
- 10 4. having non-Korean genetic background from PCA analysis with 1KGP set: seven samples
- 11 5. reported to have a rare disease: 40 samples
- 12 6. 52 samples who became not applicable for this study

13 Finally, the Korea4K variome data included 3,617 participants’ genomes. To detect variants which
14 were probably called because of a sequencing batch effect, we measured average allele balance of
15 the genotyped alleles (the read count of the allele divided by the total read count on a locus). Then,
16 we excluded 12,713,580 variants that had average allele balance of the loci out of the range of ± 1
17 \times standard deviation (SD) from a genome-wide average of allele balance to remove the sequencing
18 batch effect (Supplementary Fig.S1). We also excluded the variants which had a genotyping rate
19 of < 0.9 for downstream variant analysis. The variants in the final variome set were annotated
20 using Variant Effect Predictor (VEP) with Ensemble database (RRID:SCR_007931, ver. 101) [40].

21

22 **Principal Component Analysis (PCA) with the EBI’s 1KGP genome data**

1 The interpopulation genomic structure was evaluated by projecting the first two PCs determined
2 via PCA of SNVs from both Korea4K and East Asian populations from 1KGP. We merged variants
3 from the Korea4K and 1KGP sets and then filtered out variants with the following criteria: (i)
4 biallelic SNVs with a MAF < 1%; (ii) biallelic SNVs with an HWE $P < 10^{-6}$; (iii) biallelic SNVs
5 with a missing genotype rate of > 0.01. Extracted variants were LD pruned using “ --indep 200 4
6 0.1” option in PLINK (RRID:SCR_001757, ver. 1.90b3n) [39], yielding 330,350 sites. PCA was
7 carried out using PLINK (RRID:SCR_001757, ver. 1.90b3n) [39].

8

9 **Korean-specific missense variants**

10 We collected allele frequency data from ten populations (African (AFR), American (AMR),
11 European (EUR), South Asian (SAS), East Asian (EAS), Japanese in Tokyo (JPT), Kinh
12 Vietnamese (KHV), Han Chinese in Beijing (CHB), Han Chinese Southern (CHS), and Chinese
13 Dai in Xishuangbanna (CDX)) from EBI’s 1KGP database [41]. For each Korea4K variant, we
14 compared its allele frequency to the allele frequency of all of the ten populations using the Chi-
15 squared test. We selected variants that were specific to the Korean when the P -value of the Chi-
16 squared test to the ten populations was less than 5×10^{-5} .

17

18 **Protein structure modeling and thermodynamic stability measurement**

19 We constructed the mutant-type (MT) protein sequences of the Korean-specific missense variants
20 by substituting the reference protein sequences found in the Ensembl database
21 (RRID:SCR_002344, ver. 101) [42]. We modeled the structures of the wild-type (WT) and mutant-
22 type protein models using AlphaFold2 (ver. 2.0) with the ‘--max_template_data 2022-05-09 --

1 db_preset reduced_dbs' option with default databases downloaded by AlphaFold2 [7]. We used
2 the InterPro (RRID:SCR_006695) database [43] to determine whether a missense variant was
3 located in the domain region within the protein sequence. We extracted the domain region from
4 the WT and MT protein 3D models and excluded domains that had less than 50 amino acids.
5 Afterwards, we calculated ΔG_{WT} and ΔG_{MT} using the 'Stability' command of foldX
6 (RRID:SCR_008522) [44] to measure the protein thermodynamic stability. Finally, we measured
7 the change in protein thermodynamic stability between the two models by calculating the
8 difference between the WT and MT domain models ($\Delta\Delta G = \Delta G_{MT} - \Delta G_{WT}$).

9

10 **Imputation**

11 We constructed an imputation reference panel of Korea4K and Korea1K sets which includes 3,614,
12 and 873 Korean individuals, respectively. A total of 26,210,741 and 15,649,303 autosomal
13 biallelic variants with a missing genotype call rate of < 0.1 and minor allele count > 1 (not a
14 singleton) were extracted for the Korea4K and Korea1K panel, respectively. The extracted
15 variomes were phased into haplotype using SHAPEIT2 (ver. v2.r904) [45]. We used the same test
16 dataset as in the previous study [2]. The phased test data was imputed using the imputation
17 reference panel by Minimac3 (RRID:SCR_009292, ver. 2.0.1) [46]. We estimated imputation
18 accuracies using squared Pearson's correlation coefficients (R^2) between the true genotypes and
19 imputed genotype dosages.

20 **Clinical information**

21 We collected or calculated 107 clinical parameters (93 quantitative and 14 qualitative traits;
22 Supplementary Table S12) along with genome data from 2,685 samples among the Korea4K

1 samples. A total of 3,383 clinical datasets (including multiple time points per sample) from regular
2 health checkups carried out by various hospitals and clinics throughout Korea were collected from
3 2,685 participants between 2016 and 2019. When a single participant had multiple clinical datasets,
4 the most recent one was chosen for the following analysis. Out of the final unrelated 3,617 samples,
5 2,374 samples had clinical data available and were included in the phenomics analyses.

6 In the context of collecting data from over 200 diverse healthcare institutions, standardizing
7 clinical information on 107 traits became imperative. We resolved discrepancies in unit
8 measurements, such as micrograms and nanograms, for specific traits. Furthermore, certain clinical
9 metrics, such as the estimated glomerular filtration rate (e-GFR), were found to exhibit variability
10 contingent upon variables such as ethnicity, sex, and age. To maintain consistency and ensure
11 methodological uniformity, we enforced the adoption of a singular clinical formula for the
12 computation of e-GFR across all data samples. Such calculations were applied to 26 traits which
13 are shown in Supplementary Table S13. Clinical traits that exhibited values characterized by
14 inequalities likely due to the limit of detection (e.g., < 5.0 and > 99) were omitted from the
15 analytical procedures, as such values have the potential to introduce disturbances to subsequent
16 data analyses. Likewise, values that exhibited divergent formatting conventions across distinct
17 healthcare institutions (e.g., 20 and a few or 999 and many) were harmonized to conform with
18 prevailing standard criteria observed in most samples under investigation. Also, four quantitative
19 clinical traits and 12 qualitative traits were excluded from the further analysis, since the traits were
20 missing from more than 90% of participants due to health check-up reports heterogeneity, or the
21 traits that were qualitative and biased to one category (more than 1:4). Standard Weight was also
22 removed from the analysis, because the trait was not an inherently correct representation of the
23 sample's clinical data but rather a recommended value. Three traits (Hepatitis B virus antibody,

1 antigen, and hepatitis C antibody) contained both quantitative and qualitative values. Therefore,
2 both of the values were utilized for analysis, i.e, Hbs_Ab_Quan and Hbs_Ab_Binary. Phenotypic
3 correlations were calculated by Pearson's method. Benjamini-Hochberg method was used to adjust
4 for multiple comparisons when documenting confident phenotypic correlations with FDR.

6 **Whole genome-wide association study (GWAS)**

7 SNVs and indels with a MAF <1%, HWE $P < 10^{-6}$, and a missing genotype rate of > 0.01 were
8 excluded from the analysis using PLINK (ver. 1.90b3n) [39]. A total of 90 GWAS (88
9 quantitative and 2 qualitative traits) were performed with a total of 3,617 individuals and 7,782,381
10 variants. Each GWAS had a different number of individuals that included those who had the
11 target clinical traits. The GWAS was performed using linear and logistic regression under an
12 additive genetic model with PLINK (ver. 2.00 alpha) [47] for quantitative and qualitative traits,
13 respectively. Sex, age, age² (age squared), body mass index (BMI), and the top ten principal
14 components of SNV genotypes were included in the model as covariates. Age and BMI were
15 chosen especially due to their known shared associations with multiple traits as previously
16 documented by Shungin and colleagues which could lead to confounding biases in the downstream
17 interpretation of phenotypic relationships [48]. BMI was excluded from covariates in the GWAS
18 for BMI itself and degree of obesity. Calculating the genomic inflation factor (λ_{Median}), we found
19 that all of the traits in the test reside below 1.1 indicating there are minimal false positives caused
20 by gross population structure or systematic biases (Supplementary Fig. S4-19) [49]. We rejected
21 53 traits from further analysis based on QQ-plot analysis (Supplementary Fig. S5-S20). We used
22 5×10^{-8} for a whole-genome-wide significance threshold. The 7,782,381 variants were clumped
23 into 466,938 loci based on linkage disequilibrium (LD) information using PLINK (ver. 1.90b3n)

1 with ‘--clump-p1 1, --clump-p2 1, --clump-r2 0.1, --clump-kb 250, and --clump-index-first’
2 options [39]. Statistical powers of the 90 GWAS were calculated by the R package “genpwr”
3 under the assumption of an effect size of 0.5 and a minor allele frequency of 0.01 (Supplementary
4 Table S14).

5 **Measuring heritability and genetic correlation**

6 We calculated genetic relatedness among individuals from SNPs by genetic relationship matrix
7 (GRM) in genome-wide complex trait analysis (GCTA) (ver. 1.93.2) with ‘-autosome -maf 0.01
8 -make-grm’ options [28]. We estimated the genetic heritability of 87 quantitative traits using
9 GCTA (ver. 1.93.2) with ‘-reml -grm’ options [28]. We estimated the genetic correlations (GC)
10 using the bivariate genome-based restricted maximum likelihood (GREML) algorithm [50] in the
11 GCTA (ver. 1.93.2) with ‘-reml-bivar -grm -reml-bivar-lrt-rg’ options [28]. Two of the 253 trait
12 pairs were excluded since the log-likelihood did not converge. The correction for multiple tests
13 was done by Benjamini-Hochberg approach when reporting confident GCs that suffice the
14 threshold of FDR below 0.05.

15

16 **Calculation of Variant Sharing Index (VSI)**

17 The variant sharing index (VSI) is a Jaccard score to measure how many pleiotropic components
18 exist out of all significant variants from *i*-th and *j*-th traits, which is defined as

$$19 \text{VSI}(I,j)= |S_i \cap S_j| / |S_i \cup S_j|$$

20 where S_i and S_j denote sets of significant variants for the *i*-th and *j*-th traits, respectively. The VSI
21 increases as two traits have more pleiotropic variants among their significant variants.

22

1 **Pleiotropic variants with tissue-specific expression regulatory function**

2 We annotated the gene symbol of the pleiotropic variant by using Ensemble database (ver. 101)
3 [42]. In case of intergenic variants, we annotated the genes which were located the nearest in both
4 directions of the variant. The single tissue eQTL data (RRID: SCR_013042, ver. 8) from the GTEx
5 portal were used to investigate the eQTL of pleiotropic variants in Korea4K.

6 **Investigation of potential causal relationships between traits based on Mendelian**
7 **randomization (MR)**

8 We used the Mendelian randomization method to investigate potential causal relationships among
9 1,332 combinations of an exposure trait and an outcome trait among 37 clinical traits. MR is
10 computed from the linear regression analysis between the effects of SNPs on an exposure trait and
11 their effects on an outcome trait. We chose the SNPs with suggestive GWAS results (P -value $<$
12 10^{-5}) with exposure traits as the instrument variables. In case multiple SNPs existed in the LD
13 block, the one with the smallest P -value was chosen. We rejected 40 SNPs, which were detected
14 as outliers of linear regression from MR-PRESSO software (RRID: SCR_023697, ver. 1.0) [51]
15 with ‘NbDistribution=10000 and SignifThreshold=0.05’ options, from further analysis. MR
16 coefficients were computed using the chosen SNPs by three different methods: the Inverse-
17 variance weighted (IVW) and MR-Egger method of TwoSampleMR package
18 (RRID:SCR_019010, v.0.5.6) [52] and MR-PRESSO software (ver. 1.0) [51]. Finally, we selected
19 36 significant causal relationships that overlapped at least two of three methods (IVW, MR_Egger,
20 and MRPRESSO). All analyses were performed with default options.

1 **Acknowledgments**

2 We appreciate all participants and Ulsan citizens who supported this project. We also thank Ju
3 Yeon Park, Sangryoul Han, Jungae Shim, Nayoung Kim, Seung Gu Park, Byoung-Chul Kim,
4 Jungeun Kim, Neung-hwa Park, Suan Cho, and Yeshin Park for supporting this project. We thank
5 Young-woong Lee, Information Technology team in UNIST for supporting the data uploading to
6 EGA. This work was supported by Biodatafarm computing infrastructure funded by the Ulsan
7 metropolitan city government. We thank the Korea Institute of Science and Technology
8 Information (KISTI) for providing us with the Korea Research Environment Open NETwork
9 (KREONET). We thank our collaborators in NCSR, KRIS, and C.-G. Kim. The Genotype-
10 Tissue Expression (GTEx) Project was supported by the Common Fund of the Office of the
11 Director of the National Institutes of Health, and by NCI, NHGRI, NHLBI, NIDA, NIMH, and
12 NINDS. The data used for the analyses described in this manuscript were obtained from the GTEx
13 Portal on 11/25/2021. This work was supported by the U-K BRAND Research Fund (1.200108.01)
14 of UNIST (Ulsan National Institute of Science & Technology). This work was supported by the
15 Research Project Funded by Ulsan City Research Fund (1.200047.01) of UNIST. This work was
16 supported by the Promotion of Innovative Businesses for Regulation-Free Special Zones funded
17 by the Ministry of SMEs and Startups (MSS, Korea) (1425157301 and 1425156792). This work
18 was also supported by the Establishment of Demonstration Infrastructure for Regulation-Free
19 Special Zones funded by the Ministry of SMEs and Startups (MSS, Korea) (1425157253). This
20 research was also supported by the Technology Innovation Program (20016225, Development and
21 Dissemination on National Standard Reference Data) funded by the Ministry of Trade, Industry &
22 Energy (MOTIE, Korea). The biospecimens for this study were provided by Ulsan Medical Center
23 and the Biobanks of Gyeongsang National University Hospital, Chungbuk National University

1 Hospital (18-27, 20-04), and Kyung Hee University Hospital (2018-4, 2019-4, 2019-6), the
2 members of the National Biobank of Korea, which is supported by the Ministry of Health, Welfare
3 and Family Affairs. 215 whole-genome-seq data used in this study were provided by the Clinical
4 & Omics Data Archive (CODA), CODA accession number S000680 [53]. We thank Jaesu Bhak
5 and Maryana Bhak for editing the manuscript.

6
7 **Author contributions**

8 S.J., Hansol C., Y. J., W.C. and Jong B. wrote the manuscript. S.J., Hansol C., Y. J., Hyunjoo C.,
9 K.A., H. R., Jihun. B., H. L., Yoonsung. K., S. H., C.L., and J. S. conducted the data analysis. C.
10 K., Yeonkyung K., Younghui K., and Y. J. W. performed wet-lab experiments. S. J., Yeo Jin K.,
11 B. C. K., S.L., and Jong B. designed the study. S.J., Hansol C., Y.J. W.C., A.B., C.Y., D. B., O.
12 B., E. S., S. K., J. P., J. J., D. J., S. L., and Jong B. revised the manuscript. S. L. and Jong B. jointly
13 supervised the study.

14
15 **Ethics, consents, and permissions**

16 Sample collection and sequencing were approved by the Institutional Review Board (IRB) of the
17 Ulsan National Institute of Science and Technology (UNISTIRB-15-19-A and UNISTIRB-16-13-
18 C). The data employed in our study originates from voluntary blood or saliva donations, and we
19 have diligently secured explicit, comprehensive consent forms from all participants prior to sample
20 collection. These consent forms explicitly outline the intended use of their data for research
21 purposes and underscore the voluntary nature of their participation. Furthermore, our study adheres
22 to the ethical guidelines and regulations stipulated by the IRB. As a result, we can make the data
23 of 3,839 individuals publicly available while respecting the privacy and consent of the participants.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

Competing interests

S.J., Y. J., H. R., Y.J.K., C.K, Yeonkyung K., Younghui K., Y. J. W., and B. C. K. are employees and Jong B. is the CEO of Clinomics Inc. The authors declare no other competing interests.

Data Availability

Allele frequency information of variants is publicly available under the Korea4K webpage [54]. Raw sequencing data, individual genotype information, and clinical trait data will be as easily and freely available as possible upon request and after approval from the Korean Genomics Center’s review board in UNIST. The raw sequencing data that can be distributed were uploaded to the European Genome-Phenome Archive under the study accession ‘EGAS00001007580’. Information about the Korean Genome Project and other data sharing can be found at the Korea4K webpage All additional supporting data are available in the *GigaScience* repository, GigaDB [55].

Additional Files

Supplementary Fig. S1. Variants batch effect of DNA sequences.

Supplementary Fig. S2. Variants distribution based on variant location and allele frequency category in Korea4K.

Supplementary Fig. S3. Power comparison of whole-genome-wide association study between Korea4K and Korea1K.

Supplementary Fig. S4. Sequencing data quality metrics of Korea4K genomes.

Supplementary Fig. S5. QQplots for the whole-genome-wide association tests of the traits on the anthropometry category.

1 **Supplementary Fig. S6.** QQplots for the whole-genome-wide association tests of the traits on
2 blood circulation biochemical category.

3 **Supplementary Fig. S7.** QQplots for the whole-genome-wide association tests of the traits on
4 blood circulation physics category.

5 **Supplementary Fig. S8.** QQplots for the whole-genome-wide association tests of the traits on
6 diabetes category.

7 **Supplementary Fig. S9.** QQplots for the whole-genome-wide association tests of the traits on
8 electrolyte category.

9 **Supplementary Fig. S10.** QQplots for the whole-genome-wide association tests of the traits on
10 hearing test category.

11 **Supplementary Fig. S11.** QQplots for the whole-genome-wide association tests of the traits on
12 hematological category.

13 **Supplementary Fig. S12.** QQplots for the whole-genome-wide association tests of the traits on
14 hepatitis category.

15 **Supplementary Fig. S13.** QQplots for the whole-genome-wide association tests of the traits on
16 inflammation and etc category.

17 **Supplementary Fig. S14.** QQplots for the whole-genome-wide association tests of the traits on
18 kidney function category.

19 **Supplementary Fig. S15.** QQplots for the whole-genome-wide association tests of the traits on
20 liver function category.

21 **Supplementary Fig. S16.** QQplots for the whole-genome-wide association tests of the traits on
22 pulmonary function category.

1 **Supplementary Fig. S17.** QQplots for the whole-genome-wide association tests of the traits on
2 thyroid function category.

3 **Supplementary Fig. S18.** QQplots for the whole-genome-wide association tests of the traits on
4 tumor biomarker category.

5 **Supplementary Fig. S19.** QQplots for the whole-genome-wide association tests of the traits on
6 urinalysis category.

7 **Supplementary Fig. S20.** QQplots for the whole-genome-wide association tests of the traits on
8 vision category.

9 **Supplementary Table S1.** Sample count in Korea4K.

10 **Supplementary Table S2.** Variant count in Korea4K before sample and variant filtering.

11 **Supplementary Table S3.** Variant count based on variant categories and reported to dbSNP.

12 **Supplementary Table S4.** Allele frequency information of populations for 62 Korean-specific
13 missense variants

14 **Supplementary Table S5.** Prediction of changes in protein thermodynamic stability according to
15 missense variant

16 **Supplementary Table S6.** List of the GWAS variants which have association significance $P < 5E-$
17 8

18 **Supplementary Table S7.** Genetic heritability measurement

19 **Supplementary Table S8.** Genetic correlation measurement

20 **Supplementary Table S9.** Phenotypic correlation estimation

21 **Supplementary Table S10.** Pleiotropic variants

22 **Supplementary Table S11.** Mendelian randomization results

23 **Supplementary Table S12.** Statistics of clinical information

1 **Supplementary Table S13.** 26 traits with clinical calculations applied.

2 **Supplementary Table S14.** Statistical power measurement of 90 GWAS tests

3

4 **References**

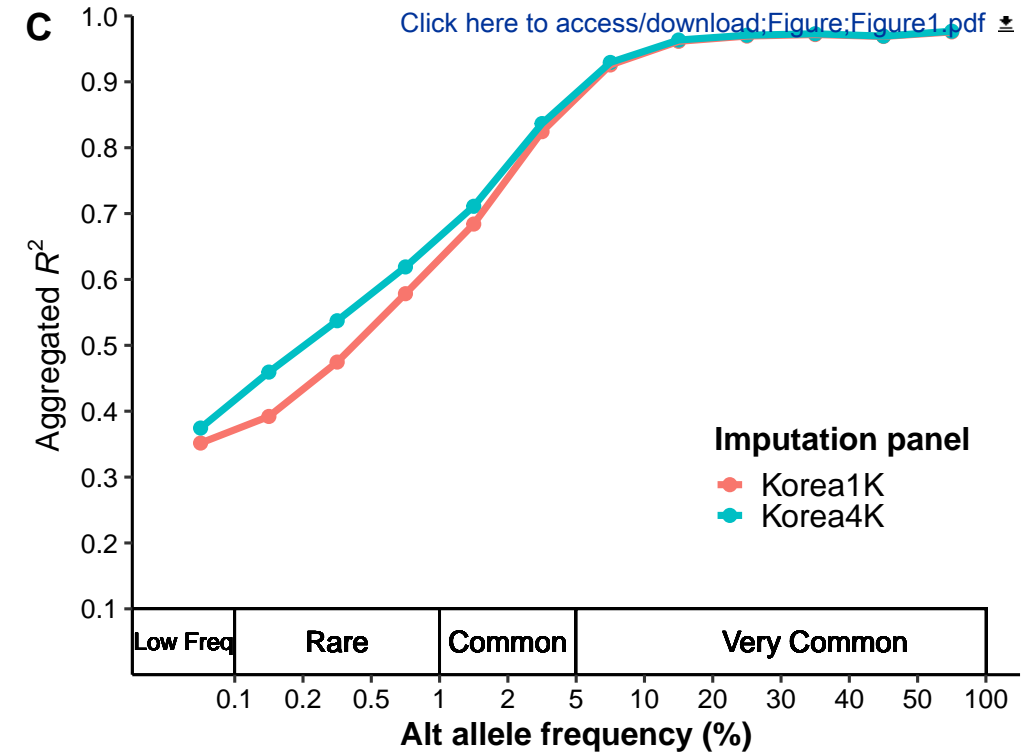
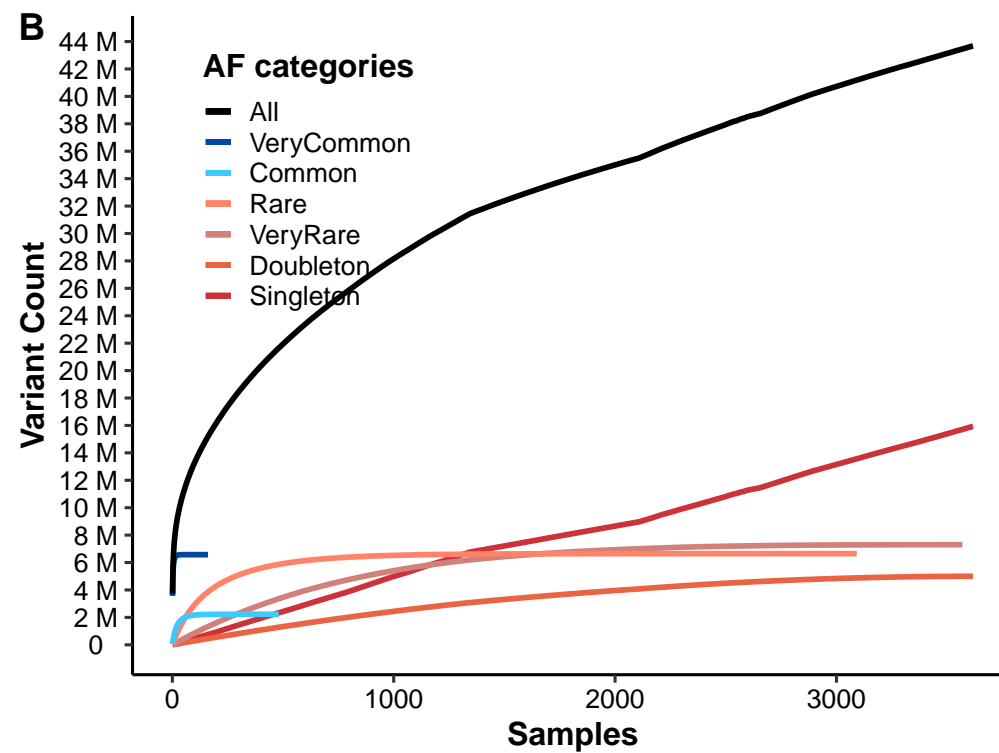
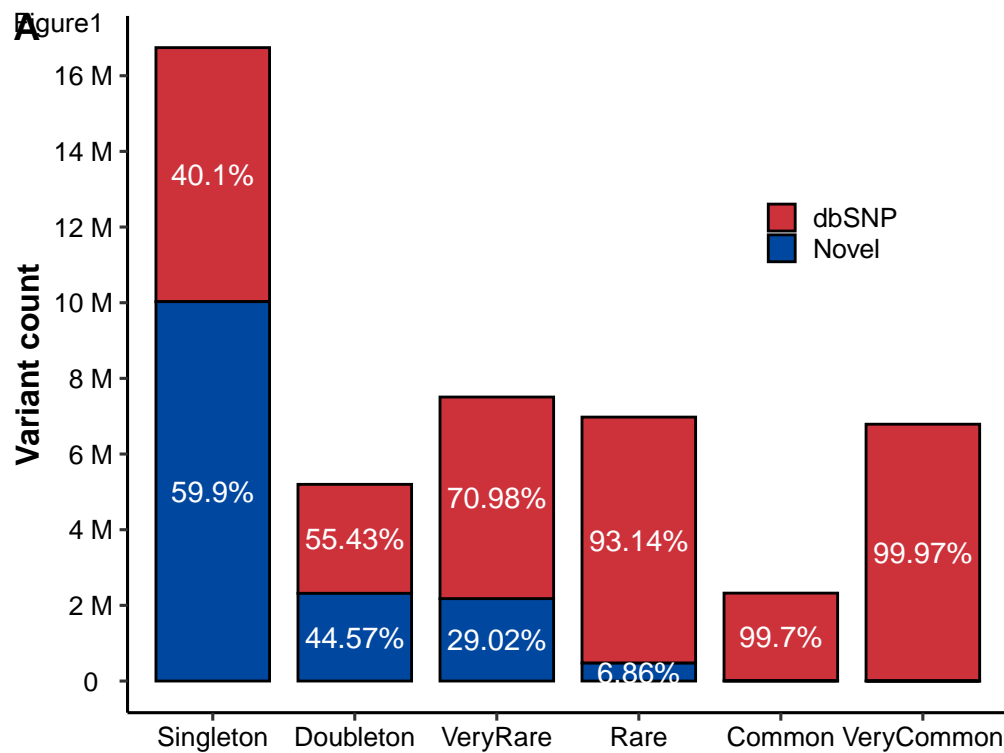
5

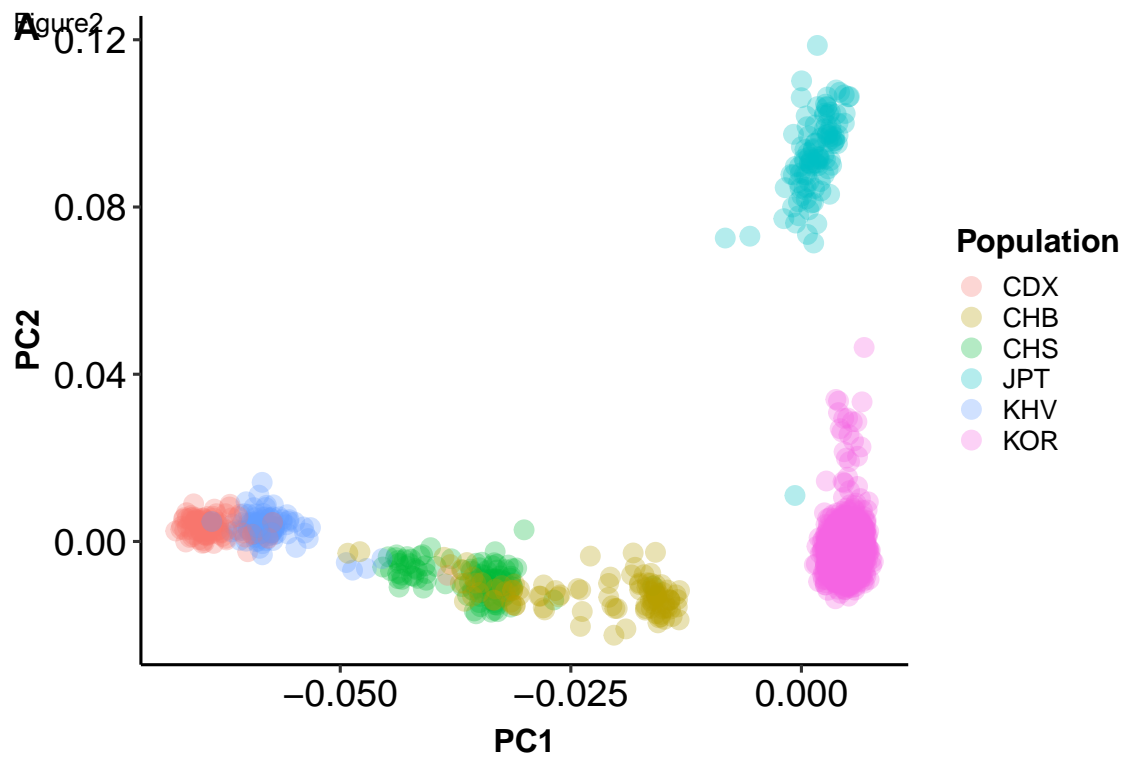
- 6 1. Song SO, Jung CH, Song YD, Park CY, Kwon HS, Cha BS, et al. Background and data
7 configuration process of a nationwide population-based study using the Korean national
8 health insurance system. *Diabetes Metab J.* 2014;38 5:395-403.
9 doi:10.4093/dmj.2014.38.5.395.
- 10 2. Jeon S, Bhak Y, Choi Y, Jeon Y, Kim S, Jang J, et al. Korean Genome Project: 1094
11 Korean personal genomes with clinical information. *Sci Adv.* 2020;6 22 doi:ARTN
12 eaaz783510.1126/sciadv.aaz7835.
- 13 3. Kanai M, Akiyama M, Takahashi A, Matoba N, Momozawa Y, Ikeda M, et al. Genetic
14 analysis of quantitative traits in the Japanese population links cell types to complex
15 human diseases. *Nat Genet.* 2018;50 3:390-+. doi:10.1038/s41588-018-0047-6.
- 16 4. Choe EK, Shivakumar M, Verma A, Verma SS, Choi SH, Kim JS, et al. Leveraging deep
17 phenotyping from health check-up cohort with 10,000 Korean individuals for phenome-
18 wide association study of 136 traits. *Sci Rep-Uk.* 2022;12 1 doi:ARTN
19 193010.1038/s41598-021-04580-2.
- 20 5. Van Hout CV, Tachmazidou I, Backman JD, Hoffman JD, Liu D, Pandey AK, et al.
21 Exome sequencing and characterization of 49,960 individuals in the UK Biobank. *Nature.*
22 2020;586 7831:749-+. doi:10.1038/s41586-020-2853-0.
- 23 6. Jiang LD, Zheng ZL, Qi T, Kemper KE, Wray NR, Visscher PM, et al. A resource-
24 efficient tool for mixed model association analysis of large-scale data. *Nat Genet.*
25 2019;51 12:1749-+. doi:10.1038/s41588-019-0530-8.
- 26 7. Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, et al. Highly
27 accurate protein structure prediction with AlphaFold. *Nature.* 2021;596 7873:583-+.
28 doi:10.1038/s41586-021-03819-2.
- 29 8. Buniello A, MacArthur JAL, Cerezo M, Harris LW, Hayhurst J, Malangone C, et al. The
30 NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted
31 arrays and summary statistics 2019. *Nucleic Acids Res.* 2019;47 D1:D1005-D12.
32 doi:10.1093/nar/gky1120.
- 33 9. Seyed Khoei N, Jenab M, Murphy N, Banbury BL, Carreras-Torres R, Viallon V, et al.
34 Circulating bilirubin levels and risk of colorectal cancer: serological and Mendelian
35 randomization analyses. *BMC Med.* 2020;18 1:229. doi:10.1186/s12916-020-01703-w.
- 36 10. Chang MH, Yesupriya A, Ned RM, Mueller PW and Dowling NF. Genetic variants
37 associated with fasting blood lipids in the U.S. population: Third National Health and
38 Nutrition Examination Survey. *BMC Med Genet.* 2010;11:62. doi:10.1186/1471-2350-
39 11-62.
- 40 11. UK Biobank. <http://www.nealelab.is/uk-biobank>. Accessed 3 April 2022.

- 1 12. Canela-Xandri O, Rawlik K and Tenesa A. An atlas of genetic associations in UK
2 Biobank. *Nat Genet.* 2018;50 11:1593-9. doi:10.1038/s41588-018-0248-z.
- 3 13. Khodayari S, Sadeghi O, Safabakhsh M and Mozaffari-Khosravi H. Meat consumption
4 and the risk of general and central obesity: the Shahedieh study. *BMC Res Notes.*
5 2022;15 1:339. doi:10.1186/s13104-022-06235-5.
- 6 14. Pimenta E, Jensen M, Jung D, Schaumann F, Boxnick S and Truebel H. Effect of Diet on
7 Serum Creatinine in Healthy Subjects During a Phase I Study. *J Clin Med Res.* 2016;8
8 11:836-9. doi:10.14740/jocmr2738w.
- 9 15. Sodini SM, Kemper KE, Wray NR and Trzaskowski M. Comparison of Genotypic and
10 Phenotypic Correlations: Cheverud's Conjecture in Humans. *Genetics.* 2018;209 3:941-8.
11 doi:10.1534/genetics.117.300630.
- 12 16. Guo S, Lv HT, Yan L and Rong FN. Hyperamylasemia may indicate the presence of
13 ovarian carcinoma A case report. *Medicine.* 2018;97 49 doi:ARTN e13520
14 10.1097/MD.00000000000013520.
- 15 17. Shintani D, Yoshida H, Imai Y and Fujiwara K. Acute pancreatitis induced by paclitaxel
16 and carboplatin therapy in an ovarian cancer patient. *Eur J Gynaecol Oncol.* 2016;37
17 2:286-7.
- 18 18. Zakrzewska I and Pietryńczak M. The activity of alpha-amylase and its salivary
19 isoenzymes in serum and urine of patients with neoplastic diseases of female
20 reproductive organs. *Roczniki Akademii Medycznej w Białymstoku (1995).* 1996;41
21 2:492-8.
- 22 19. Hemani G, Bowden J and Davey Smith G. Evaluating the potential role of pleiotropy in
23 Mendelian randomization studies. *Hum Mol Genet.* 2018;27 R2:R195-R208.
24 doi:10.1093/hmg/ddy163.
- 25 20. Pulit SL, Stoneman C, Morris AP, Wood AR, Glastonbury CA, Tyrrell J, et al. Meta-
26 analysis of genome-wide association studies for body fat distribution in 694 649
27 individuals of European ancestry. *Hum Mol Genet.* 2019;28 1:166-74.
28 doi:10.1093/hmg/ddy327.
- 29 21. Ebrahim S and Davey Smith G. Mendelian randomization: can genetic epidemiology help
30 redress the failures of observational epidemiology? *Hum Genet.* 2008;123 1:15-33.
31 doi:10.1007/s00439-007-0448-6.
- 32 22. Aabo K, Pedersen H and Kjaer M. Carcinoembryonic antigen (CEA) and alkaline
33 phosphatase in progressive colorectal cancer with special reference to patient survival.
34 *Eur J Cancer Clin Oncol.* 1986;22 2:211-7. doi:10.1016/0277-5379(86)90033-7.
- 35 23. Tartter PI, Slater G, Gelernt I and Aufses AH, Jr. Screening for liver metastases from
36 colorectal cancer with carcinoembryonic antigen and alkaline phosphatase. *Ann Surg.*
37 1981;193 3:357-60. doi:10.1097/0000658-198103000-00019.
- 38 24. Walach N, Guterman A, Zaidman JL, Kaufman S and Scharf S. Leukocyte alkaline
39 phosphatase and carcinoembryonic antigen in breast cancer patients: clinical correlation
40 with the markers. *J Surg Oncol.* 1989;40 2:85-7. doi:10.1002/jso.2930400205.
- 41 25. Forouhi NG, Sattar N and McKeigue PM. Relation of C-reactive protein to body fat
42 distribution and features of the metabolic syndrome in Europeans and South Asians. *Int J*
43 *Obes Relat Metab Disord.* 2001;25 9:1327-31. doi:10.1038/sj.ijo.0801723.
- 44 26. Lim S, Jang HC, Lee HK, Kimm KC, Park C and Cho NH. The relationship between
45 body fat and C-reactive protein in middle-aged Korean population. *Atherosclerosis.*
46 2006;184 1:171-7. doi:10.1016/j.atherosclerosis.2005.04.003.

- 1 27. Lee CM, Huxley RR, Wildman RP and Woodward M. Indices of abdominal obesity are
2 better discriminators of cardiovascular risk factors than BMI: a meta-analysis. *J Clin*
3 *Epidemiol.* 2008;61 7:646-53. doi:10.1016/j.jclinepi.2007.08.012.
- 4 28. Yang JA, Lee SH, Goddard ME and Visscher PM. GCTA: A Tool for Genome-wide
5 Complex Trait Analysis. *Am J Hum Genet.* 2011;88 1:76-82.
6 doi:10.1016/j.ajhg.2010.11.011.
- 7 29. Zhang YL, Cheng YS, Jiang W, Ye YX, Lu QS and Zhao HY. Comparison of methods
8 for estimating genetic correlation between complex traits using GWAS summary
9 statistics. *Brief Bioinform.* 2021;22 5 doi:ARTN bbaa44210.1093/bib/bbaa442.
- 10 30. Visscher PM, Hemani G, Vinkhuyzen AA, Chen GB, Lee SH, Wray NR, et al. Statistical
11 power to detect genetic (co)variance of complex traits using SNP data in unrelated
12 samples. *PLoS Genet.* 2014;10 4:e1004269. doi:10.1371/journal.pgen.1004269.
- 13 31. Li J, Gui L, Wu C, He Y, Zhou L, Guo H, et al. Genome-wide association study on serum
14 alkaline phosphatase levels in a Chinese population. *BMC Genomics.* 2013;14:684.
15 doi:10.1186/1471-2164-14-684.
- 16 32. Middelberg RP, Ferreira MA, Henders AK, Heath AC, Madden PA, Montgomery GW, et
17 al. Genetic variants in LPL, OASL and TOMM40/APOE-C1-C2-C4 genes are associated
18 with multiple cardiovascular-related traits. *BMC Med Genet.* 2011;12:123.
19 doi:10.1186/1471-2350-12-123.
- 20 33. Jeon Y, Jeon S, Choi WH, An K, Choi H, Kim BC, et al. Genome-wide analyses of early-
21 onset acute myocardial infarction identify 29 novel loci by whole genome sequencing.
22 *Hum Genet.* 2023;142 2:231-43. doi:10.1007/s00439-022-02495-0.
- 23 34. Nurk S, Koren S, Rhie A, Rautiainen M, Bizkadze AV, Mikheenko A, et al. The
24 complete sequence of a human genome. *Science.* 2022;376 6588:44-53.
25 doi:10.1126/science.abj6987.
- 26 35. Liao WW, Asri M, Ebler J, Doerr D, Haukness M, Hickey G, et al. A draft human
27 pangenome reference. *Nature.* 2023;617 7960:312-24. doi:10.1038/s41586-023-05896-x.
- 28 36. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
29 *EMBnet journal.* 2011;17 1:10-2.
- 30 37. Li H and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
31 transform. *Bioinformatics.* 2009;25 14:1754-60. doi:10.1093/bioinformatics/btp324.
- 32 38. Poplin R, Ruano-Rubio V, DePristo MA, Fennell TJ, Carneiro MO, Van der Auwera GA,
33 et al. Scaling accurate genetic variant discovery to tens of thousands of samples. *BioRxiv.*
34 2018:201178.
- 35 39. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK:
36 A tool set for whole-genome association and population-based linkage analyses. *Am J*
37 *Hum Genet.* 2007;81 3:559-75. doi:10.1086/519795.
- 38 40. McLaren W, Gil L, Hunt SE, Riat HS, Ritchie GR, Thormann A, et al. The ensembl
39 variant effect predictor. *Genome biology.* 2016;17 1:1-14.
- 40 41. Genomes Project C, Abecasis GR, Altshuler D, Auton A, Brooks LD, Durbin RM, et al.
41 A map of human genome variation from population-scale sequencing. *Nature.* 2010;467
42 7319:1061-73. doi:10.1038/nature09534.
- 43 42. Cunningham F, Allen JE, Allen J, Alvarez-Jarreta J, Amode MR, Armean IM, et al.
44 Ensembl 2022. *Nucleic Acids Res.* 2022;50 D1:D988-D95. doi:10.1093/nar/gkab1049.

- 1 43. Blum M, Chang HY, Chuguransky S, Grego T, Kandasamy S, Mitchell A, et al. The
2 InterPro protein families and domains database: 20 years on. *Nucleic Acids Res.* 2021;49
3 D1:D344-D54. doi:10.1093/nar/gkaa977.
- 4 44. Delgado J, Radusky LG, Cianferoni D and Serrano L. FoldX 5.0: working with RNA,
5 small molecules and a new graphical interface. *Bioinformatics.* 2019;35 20:4168-9.
6 doi:10.1093/bioinformatics/btz184.
- 7 45. McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, et al. A
8 reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet.* 2016;48
9 10:1279-83. doi:10.1038/ng.3643.
- 10 46. Das S, Forer L, Schonherr S, Sidore C, Locke AE, Kwong A, et al. Next-generation
11 genotype imputation service and methods. *Nat Genet.* 2016;48 10:1284-7.
12 doi:10.1038/ng.3656.
- 13 47. Chang CC, Chow CC, Tellier LCAM, Vattikuti S, Purcell SM and Lee JJ. Second-
14 generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience.*
15 2015;4 doi:ARTN 710.1186/s13742-015-0047-8.
- 16 48. Shungin D, Winkler TW, Croteau-Chonka DC, Ferreira T, Locke AE, Magi R, et al. New
17 genetic loci link adipose and insulin biology to body fat distribution. *Nature.* 2015;518
18 7538:187-96. doi:10.1038/nature14132.
- 19 49. Yang J, Weedon MN, Purcell S, Lettre G, Estrada K, Willer CJ, et al. Genomic inflation
20 factors under polygenic inheritance. *Eur J Hum Genet.* 2011;19 7:807-12.
21 doi:10.1038/ejhg.2011.39.
- 22 50. Lee S, Yang J, Goddard M, Visscher P and Wray N. Estimation of pleiotropy between
23 complex diseases using SNP-derived genomic relationships and restricted maximum
24 likelihood. *Bioinformatics.* 2012;28 19:2540-2.
- 25 51. Verbanck M, Chen CY, Neale B and Do R. Detection of widespread horizontal
26 pleiotropy in causal relationships inferred from Mendelian randomization between
27 complex traits and diseases. *Nat Genet.* 2018;50 5:693-+. doi:10.1038/s41588-018-0099-
28 7.
- 29 52. Hemani G, Zhengn J, Elsworth B, Wade KH, Haberland V, Baird D, et al. The MR-Base
30 platform supports systematic causal inference across the human phenome. *Elife.* 2018;7
31 doi:ARTN e3440810.7554/eLife.34408.
- 32 53. CODA. <http://coda.nih.go.kr>. Accessed October 2018.
- 33 54. Korea4K Genomes. http://koreangenome.org/Korea4K_Genomes. Accessed 13 March
34 2024.
- 35 55. Jeon S, Choi H, Jeon Y, Choi W, Choi H, An K, et al. Supporting data for "Korea4K:
36 whole genome sequences of 4,157 Koreans with 107 phenotypes derived from extensive
37 health check-ups". *GigaScience Database.* 2024. doi:dx.doi.org/10.5524/102507.
- 38

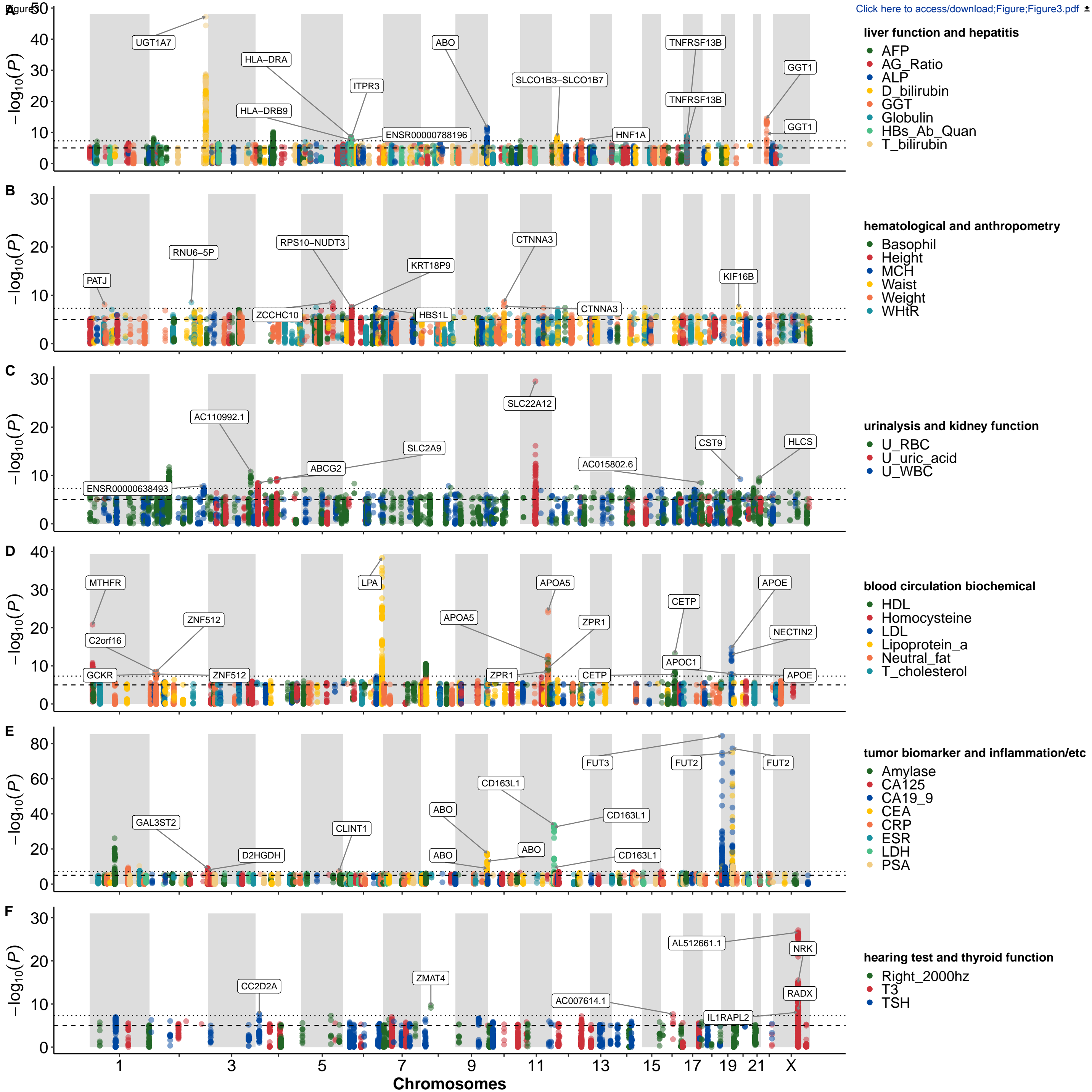


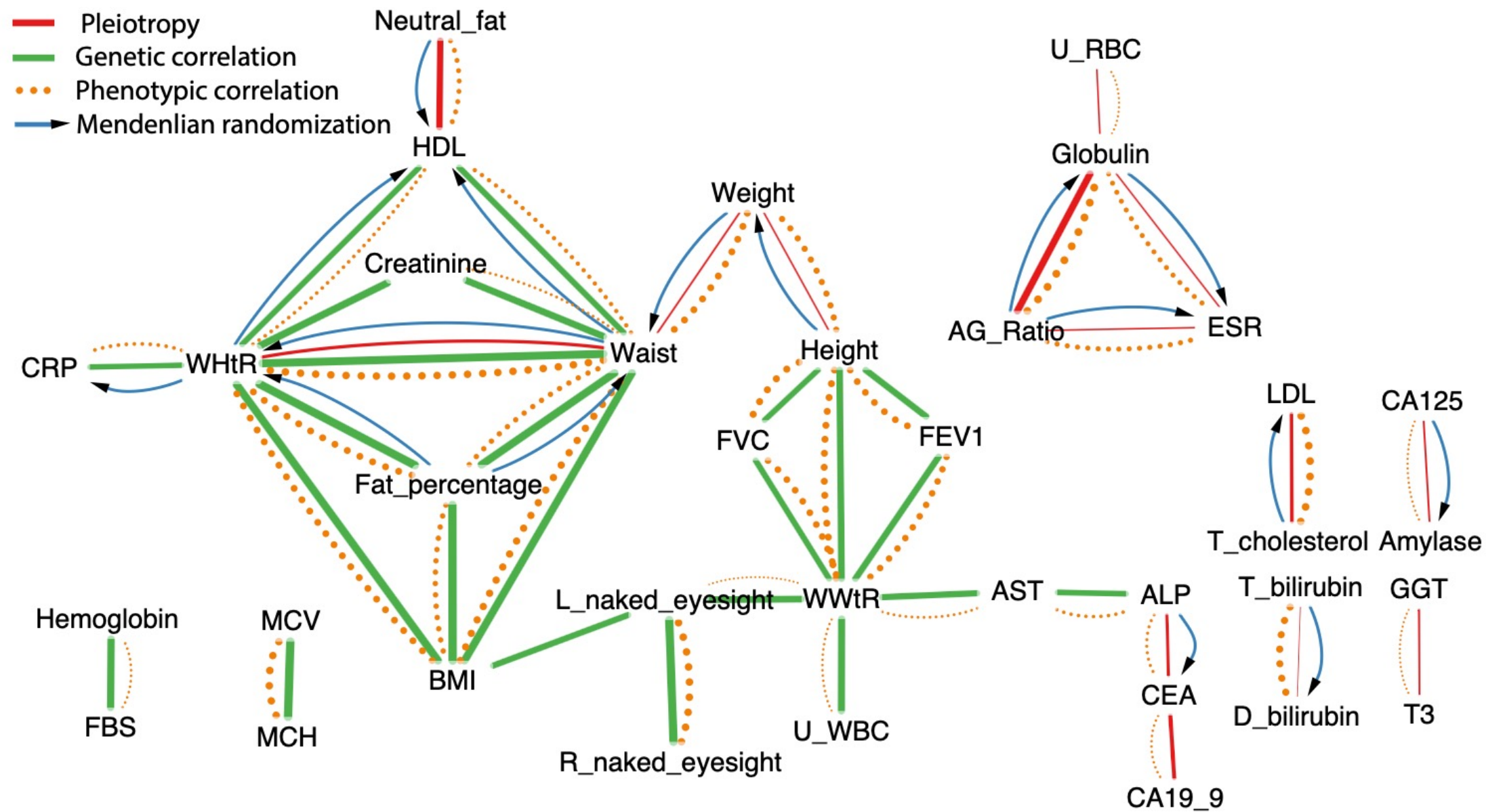


B

[Click here to access/download;Figure;Figure2.pdf](#)

Korea4K	0.58	0.49	0	0	0.04	0.13	0	0.49	0	0	1	0
JPT	0.82	0.04	0.99	1	0.26	1	0.99	0.15	1	1	0.5	0.62
CHB	0.79	0.06	0.99	1	0.27	1	1	0.13	1	1	0.56	0.6
CHS	0.84	0.06	0.99	1	0.24	0.99	1	0.2	1	1	0.54	0.56
CDX	0.78	0.04	0.98	1	0.12	1	1	0.12	1	1	0.55	0.58
KHV	0.73	0.02	0.97	1	0.32	1	1	0.2	0.99	1	0.56	0.7
EAS	0.79	0.04	0.98	1	0.24	1	1	0.16	1	1	0.54	0.61
SAS	0.46	0.04	1	1	0.16	1	1	0.12	1	1	0.14	0.59
EUR	0.32	0.05	1	1	0.24	1	1	0.13	1	1	0.12	0.69
AMR	0.3	0.06	1	1	0.24	1	1	0.16	1	1	0.15	0.7
AFR	0.33	0.12	1	0.99	0.36	1	1	0.13	1	1	0.41	0.92
	rs1141967(TPSD1)	rs12420076(OR9G1)	rs144456901(GDF2)	rs1556826592(MAGEA3)	rs200581589(HLA-DRB5)	rs201779716(PRSS2)	rs201790399(LIX1L)	rs4990121(OR8U1)	rs569378041(RBP3)	rs587669051(PDZK1)	rs71497225(MRC1)	rs78574933(IGHV4-4)







[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS1_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS2_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS3_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS4_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS5_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS6_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS7_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS8_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS9_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS10_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS11_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS12_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS13_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_TableS14_Supplementary_Material.xlsx](#)





[Click here to access/download](#)

Supplementary Material

[Korea4K_Figures_Supplementary_Material.docx](#)



Dear Editor,

We would like to submit our manuscript entitled “**Korea4K: whole genome sequences of 4,157 Koreans with 107 phenotypes derived from extensive health check-ups**” as a *Research Article* in *GigaScience*.

Since we reported Korea1K (1,094 Korean genomes with 79 clinical traits) in 2020 (Jeon et al., *Sci Adv.* 2020), we have pursued a more comprehensive study based on a larger cohort of Koreans (4,157 whole genomes with 107 clinical traits) as the second phase of the Korean Genome Project (KGP).

Here, we found that only around 4,000 whole genomes (Korea4K) were sufficient to cover the genomic diversity of the Korean population with East Asian ancestry by analyzing the statistics of common and rare SNP variants. We also present the Korea4K variome database as a part of the KGP, which could be a resource for a large-scale population genomics analysis of diverse ethnic groups in association with human evolution and diseases.

The major difference between Korea1K and Korea4K is not only in the sample size but also in the number of clinical traits derived from extensively curated reports covering the most common health check-up parameters. With the greater number of samples and clinical traits, we were able to identify 1,356 new associations between genotypes and phenotypes, which had not been detected in Korea1K. Furthermore, we performed genetic correlation, pleiotropy, and Mendelian randomization analyses to map the variome with the clinical traits from common health check-ups. We also confirmed that Korea4K, compared to Korea1K, could improve quality as a reference panel for genotype imputation.

As our study provides a possibly useful resource for exploring the relationship between the genome and the phenome, and the variome data will be publicly available as open as possible, we believe that this manuscript fits the scope of *GigaScience*.

All study participants provided informed consent, and the study design was approved by the appropriate ethics review board.

S.J., Y. J., H. R., Y.J.K., C.K, Yeonkyung K., Younghui K., Y. J. W., and B. C. K. are employees and Jong B. is the CEO of Clinomics Inc. The authors declare no other competing interests.

We confirm that all authors have approved the manuscript for submission and the content of the manuscript has not been published, or submitted for publication elsewhere.

We would like to suggest the following reviewers:

- Tim Hubbard, Ph.D., Professor of Bioinformatics and Head of Department, Department of Medical & Molecular Genetics at King’s College London, tim.hubbard@kcl.ac.uk
- Masao Nagasaki, Ph.D., Professor at the Center for Genomic Medicine, Kyoto University, nagasaki@genome.med.kyoto-u.ac.jp

Thank you for your consideration.

Sincerely,
Jong Bhak, Ph.D.

Korean Genomics Center
Ulsan National Institute of Science and Technology
Ulsan 44919, Republic of Korea
Email: jongbhak@genomics.org
Tel: +82 (0)10 4644 6754

Re: **“Korea4K: whole genome sequences of 4,157 Koreans with 107 phenotypes derived from extensive health check-ups”**

Dear *Giga Science* Editor,

We would like to express our appreciation for the constructive feedback provided by the reviewers and editorial team. Please find our revision note and a revised version of our manuscript **“Korea4K: whole genome sequences of 4,157 Koreans with 107 phenotypes derived from extensive health check-ups”**.

We have checked and tried to accommodate all the critical points and suggestions from the reviewers and you.

As a significant addition, we have clarified the methods of sequencing data and clinical data preprocessing. Also, we have investigated the results of the phenomics analyses and published literature and extended additional discussion points in the discussion sections. We have also been uploading the sequencing raw data of 4K samples to the European Genome-Phenome Archive (EGA) during the revision to make the scientific community access our data conveniently. Despite the very large size of the whole data (~120TB), we think we can finish the data upload on time.

Our manuscript has been modified in the abstract, results, discussion, and method to better represent the changes made according to all the reviewers' criticism and suggestions.

We think that these revisions address the concerns raised during the initial review and significantly enhance the scientific quality of the manuscript. Our study contributes to the understanding of the relationship between the genome and the phenome, providing a valuable resource for the scientific community.

Thank you very much for your consideration.

Sincerely,
Jong Bhak, Ph.D.

Korean Genomics Center
Ulsan National Institute of Science and Technology
Ulsan 44919, Republic of Korea
Email: jongbhak@genomics.org
Tel: +82 (0)10 4644 6754