

# UCLA

## UCLA Previously Published Works

### Title

Whole genome sequence analysis of blood lipid levels in >66,000 individuals

### Permalink

<https://escholarship.org/uc/item/39k4g98h>

### Journal

Nature Communications, 13(1)

### ISSN

2041-1723

### Authors

Selvaraj, Margaret Sunitha

Li, Xihao

Li, Zilin

et al.

### Publication Date

2022

### DOI

10.1038/s41467-022-33510-7

### Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>

Peer reviewed

# Whole genome sequence analysis of blood lipid levels in >66,000 individuals

Received: 27 September 2021

Accepted: 21 September 2022

Published online: 11 October 2022

 Check for updates

A list of authors and their affiliations appears at the end of the paper

Blood lipids are heritable modifiable causal factors for coronary artery disease. Despite well-described monogenic and polygenic bases of dyslipidemia, limitations remain in discovery of lipid-associated alleles using whole genome sequencing (WGS), partly due to limited sample sizes, ancestral diversity, and interpretation of clinical significance. Among 66,329 ancestrally diverse (56% non-European) participants, we associate 428M variants from deep-coverage WGS with lipid levels; ~400M variants were not assessed in prior lipids genetic analyses. We find multiple lipid-related genes strongly associated with blood lipids through analysis of common and rare coding variants. We discover several associated rare non-coding variants, largely at Mendelian lipid genes. Notably, we observe rare *LDLR* intronic variants associated with markedly increased LDL-C, similar to rare *LDLR* exonic variants. In conclusion, we conducted a systematic whole genome scan for blood lipids expanding the alleles linked to lipids for multiple ancestries and characterize a clinically-relevant rare non-coding variant model for lipids.

The discovery of rare alleles linked to plasma lipids (i.e., low-density lipoprotein cholesterol [LDL-C], high-density lipoprotein cholesterol [HDL-C], total cholesterol [TC], and triglycerides [TG]) continue to yield important translational insights toward coronary artery disease (CAD), including *PCSK9* and *ANGPTL3* inhibitors now available in clinical practice<sup>1–5</sup>. The monogenic and polygenic bases of plasma lipids are well-suited to population-based discovery analyses and confer broader insights for genetic analyses of complex traits. We now evaluate numerous newly catalogued, largely rare, alleles never previously systematically analyzed with lipids.

Analyses of imputed array-derived genome-wide genotypes and whole exome sequences in hundreds of thousands of increasingly diverse individuals continue to uncover low-frequency protein-coding variants linked to lipids. Due to purifying selection, causal variants conferring large effects tend to occur relatively more recently, and are thus rare and often specific to families or communities<sup>6</sup>. Most discovery analyses for large-effect rare alleles have focused on the analysis of disruptive protein-coding variants given (1) well-recognized constraint in coding regions, (2) incomplete genotyping of rare non-coding sequence given relative sparsity of deep-coverage (i.e., >30X) whole genome sequencing (WGS), and (3) better prediction of coding versus non-coding sequence

variation consequence<sup>1,7–12</sup>. We recently described a statistical framework incorporating multi-dimensional reference datasets paired with genomic data to improve rare coding and non-coding variant analyses for WGS analysis of lipids and other complex traits<sup>13,14</sup>. Furthermore, including individuals of non-European ancestry facilitates the discovery of both novel alleles at established loci as well as novel loci<sup>14–16</sup>.

Here, we examine the full allelic spectrum with plasma lipids using whole genome sequences and harmonized lipids from the National Heart, Lung, and Blood Institute (NHLBI) Trans-Omics for Precision Medicine (TOPMed) program<sup>17,18</sup>. We studied 66,329 participants and 428 million variants across multiple ancestry groups—44.48% European, 25.60% Black, 21.02% Hispanic, 7.11% Asian, and 1.78% Samoan. We identified robust allelic heterogeneity at known loci with several novel variants at these loci; we additionally identified novel loci and pursued replication in independent cohorts. We then explored the association of genome-wide rare variants with lipids, with detailed explorations of rare coding and non-coding variant models at known Mendelian dyslipidemia genes. Our systemic effort yields new insights for plasma lipids and provides a framework for population-based WGS analysis of complex traits.

✉ e-mail: [gpeloso@bu.edu](mailto:gpeloso@bu.edu); [pnatarajan@mgh.harvard.edu](mailto:pnatarajan@mgh.harvard.edu)

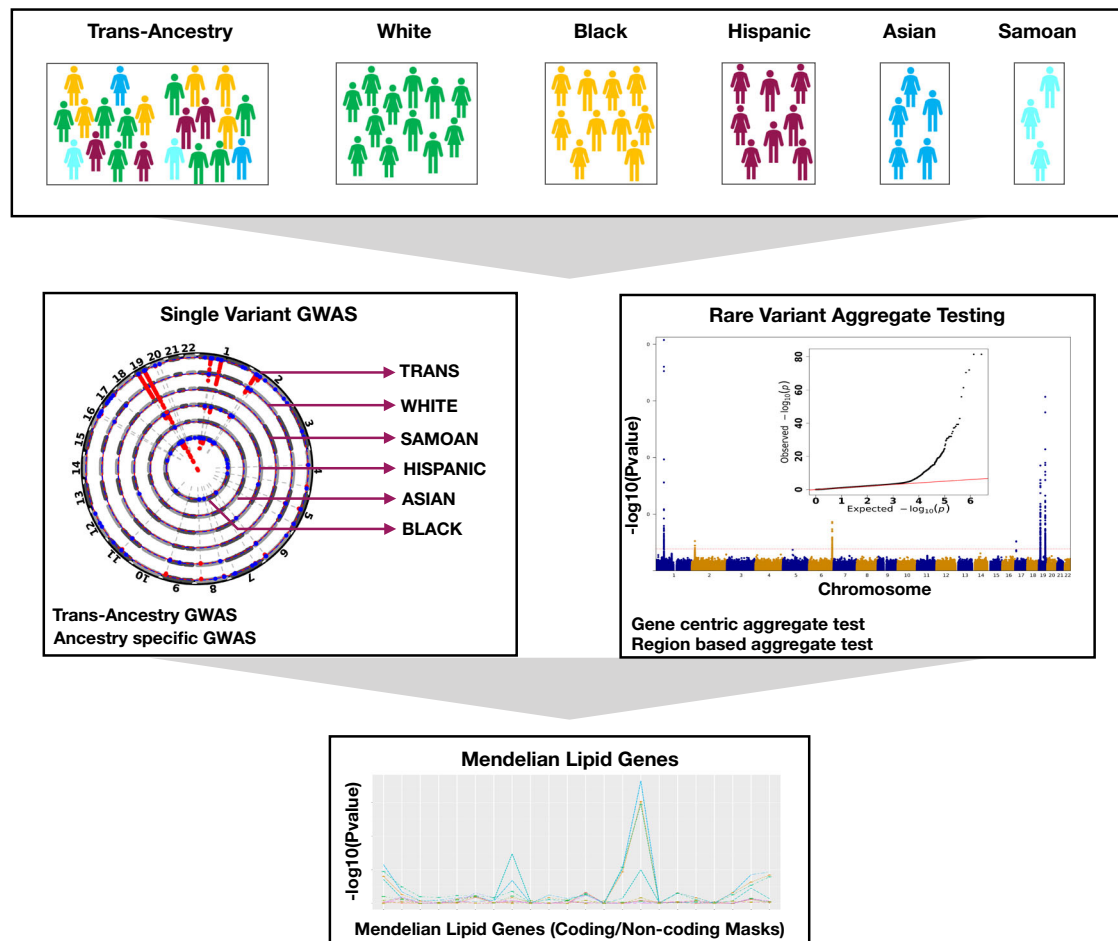
## Results

### Overview

We studied the TOPMed Freeze8 dataset of 66,329 samples from 21 studies and performed genome-wide association studies (GWAS) separately for the four plasma lipid phenotypes (i.e., LDL-C, HDL-C, TC, and TG) using 28 M individual autosomal variants (minor allele count [MAC] >20) and aggregated rare autosomal variant (minor allele frequency [MAF] <1%) association testing for 417 M variants (Fig. 1, Supplementary Fig. 1). Secondly, we associated individual variants with minor allele frequencies (MAF) >0.01% within each ancestry group to detect ancestry-specific lipid-associated alleles. We intersected our results with currently published array-based GWAS results<sup>15</sup> to identify novel associations with lipids. We performed replication analyses for the putative novel associations identified, in up to ~45,000 independent samples with array-based genotyping imputed to TOPMed and 400 K samples from UK Biobank (UKB) imputed genotypes. Finally, we conducted rare variant association studies as multiple aggregate tests across the genome to identify gene-specific functional categories and non-coding genomic regions influencing plasma lipid concentrations. We replicated the significant rare variant aggregates in ~130 K whole genomes from UKB.

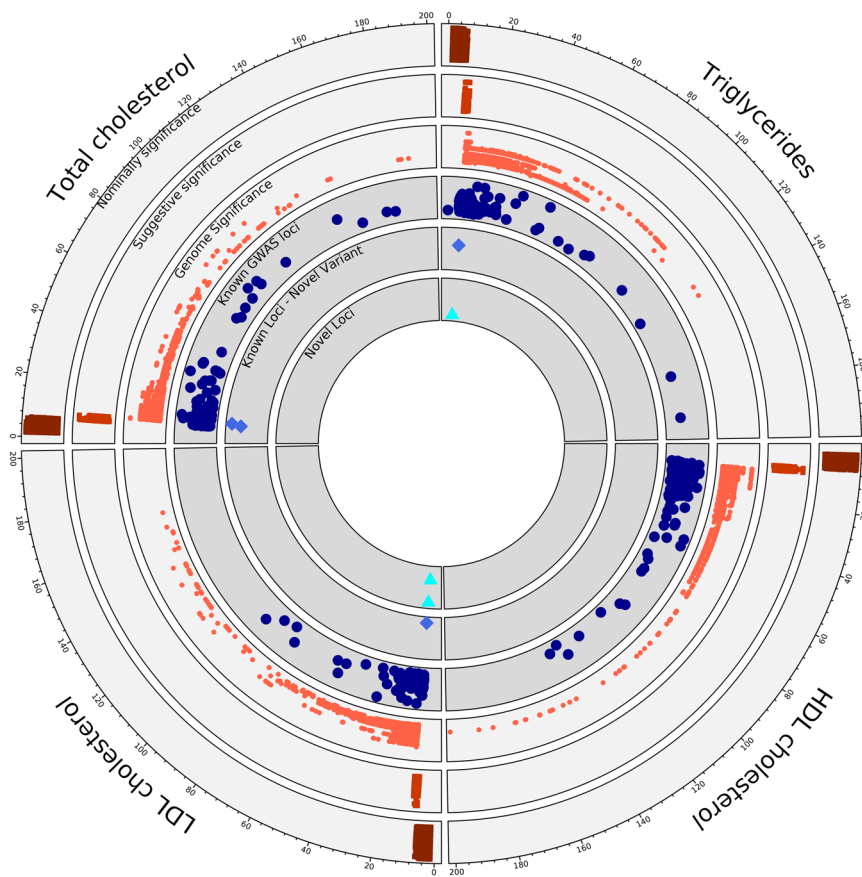
### TOPMed baseline characteristics

The TOPMed Informatics Research Center (IRC) and TOPMed Data Coordinating Center (DCC) performed quality control, variant calling, and calculated the relatedness of population structures of Freeze 8 data<sup>17</sup>. We studied 66,329 samples across 21 cohorts, and 41,182 (62%) were female. The ancestry distribution was 29,502 (44.46%) White, 16,983 (25.60%) Black, 13,943 (21.02%) Hispanic, 4719 (7.11%) Asian, and 1182 (1.78%) Samoan (Supplementary Data 1). The mean (standard deviation [SD]) age of the full cohort was 53 (15.00) years which varied by cohort from 25 (3.56) years for Coronary Artery Risk Development in Young Adults (CARDIA) to 73 (5.38) years for Cardiovascular Health Study (CHS). The Amish cohort had a higher-than-average concentration of LDL-C (140 [SD 43] mg/dL) and HDL-C (56 [SD 16] mg/dL) as well as lower TG (median 63 [IQR 50] mg/dL) consistent with the known founder mutations in *APOB* and *APOC3*<sup>7,8,14</sup>. In the Women's Health Initiative (WHI) cohort, the TC (230 [SD 41] mg/dL) and TG (median 129 [IQR 87] mg/dL) concentrations were higher than for other cohorts as previously described<sup>12</sup>. We accounted for lipid-lowering medications and fasting status and inverse rank normalized the phenotypes as before<sup>12,14</sup> which are further detailed in the Methods. The adjusted normalized lipid concentrations for the four lipids were similar across the cohorts.



**Fig. 1 | Overall study schematic.** The analyses were conducted using the multi-ancestral TOPMed freeze8 data to associate whole genome sequence variation with lipid phenotypes (i.e., LDL-C, HDL-C, TC, and TG). A total of 66,329 samples with lipids quantified data from five ancestry groups were analyzed. Single variant GWAS were carried out using SAIGE on the Encore platform using SNPs with MAC >20. Both trans-ancestry and ancestry-specific GWAS were conducted. Genome-wide rare variant (MAF <1%) gene-centric and region-based aggregate tests were

grouped and analyzed using STAARpipeline. Finally, single variant and rare variant associations at Mendelian dyslipidemia genes were investigated in further detail. TOPMed Trans-Omics for Precision Medicine, HDL-C high-density lipoprotein cholesterol, LDL-C low-density lipoprotein cholesterol, TC total cholesterol, TG triglycerides, GWAS genome wide association study, SAIGE Scalable and Accurate Implementation of GEneralized mixed model, MAC minor allele count, MAF minor allele frequency, SNPs single nucleotide polymorphisms.



**Fig. 2 | Summary of single variant genome-wide association.** Representation of the single variant GWAS results from TOPMed Freeze 8 whole genome sequenced data of 66,329 samples. Each quarter represents a different lipid phenotype, and dots extending in clock-wise fashion represent variants with increasing evidence of association as noted by  $-\log_{10}(p\text{-value})$ , which was truncated at 200. The outer three circles show the GWAS data from TOPMed freeze8 where variants binned to nominally significant ( $p\text{-value } 0.05\text{--}5 \times 10^{-7}$ ), suggestive significant ( $p\text{-value } 5 \times 10^{-7}\text{--}5 \times 10^{-9}$ ) and genome wide significant ( $p\text{-value } < 5 \times 10^{-9}$ ). The inner three

circles compare our TOPMed results with known significantly associated lipid loci and variants from the MVP summary statistics and GWAS catalog to the identified novel variants and loci that are genome-wide significant from the current study, respectively. The figure represents the outputs from two-sided genetic association testing performed using SAIGE-QT model, where the model was adjusted for all the covariates; see Methods. TOPMed Trans-Omics for Precision Medicine, GWAS genome wide association study, MVP million veteran program.

A total of 428 M variants passed the quality criteria with an average depth  $>30X$  in 22 autosomes. 202 M variants were singletons, 417 M were rare variants (MAF  $<1\%$ ), and 11 M were common or low-frequency variants (MAF  $>1\%$ ) with differences by cohort (Supplementary Data 2).

### Individual variant associations with lipids

We performed single variant analysis of  $\sim 28$  M variants with a MAC  $> 20$  for four lipid phenotypes. We identified significant genomic risk loci for each lipid level (Supplementary Data 3) and considered a  $p\text{-value } < 5 \times 10^{-9}$  to claim significance as previously recommended for whole genome sequencing common variant association studies<sup>14,19</sup>. The total numbers of variants that met our significance threshold were 2214, 2314, 2697, and 2442 for LDL-C, HDL-C, TC and TG, respectively, and after clumping<sup>20</sup> the numbers of variants were 357, 338, 324, and 289, respectively. Of these variants, 99% were previously demonstrated to be associated with plasma lipids either at the variant- or locus-level<sup>15</sup> (Supplementary Data 4, Supplementary Fig. 2).

To identify putative novel variant associations, we compared our results to a recent multi-ethnic lipid GWAS among 312,571 participants of the Million Veteran Program (MVP)<sup>15</sup> as well as the GWAS Catalog (All associations (v1.0) file dated 06/04/2020) (Fig. 2). We clumped (window 250 kb,  $r^2$  0.5) significant variants using Plink<sup>20</sup> and queried these in the GWAS Catalog and MVP. Among genome-wide significant

variants, we tabulated 'known-position' (variant previously associated), 'known-loci' (variants not previously significantly associated with the corresponding lipid phenotype but within 500 kb of a known locus, thereby representing additional allelic heterogeneity), and 'novel' variants (variants not in a known lipid locus) (Supplementary Data 4).

The novel variants, tabulated in Table 1, are divided into two subsets—'novel variants' or variants at established lipid loci for another lipid phenotype, and 'novel loci,' representing new loci associations for any lipid phenotype. For example, the *CETP* locus is well-known for its link to HDL-C, but we now found that rs183130 (16:56957451:C:T, MAF 28.3%) at the locus is associated with LDL-C. Similarly, the variants rs7140110 (13:113841051:T:C, MAF 27.8%) *GAS6* and rs73729083 (7:137875053:T:C, MAF 4.5%) *CREB3L2* are newly associated with TC, while previous studies showed that rs73729083 associates with LDL-C<sup>21</sup> and rs7140110 associates with LDL-C<sup>22</sup> and TG<sup>23</sup>. Index variants at novel loci were typically low-frequency variants often observed in non-European ancestries, so we also conducted ancestry-specific association analyses for these alleles (Supplementary Data 5). For example, 12q23.1 (12:97352354:T:C, MAF 0.3%) and 4q34.2 (4:176382171:C:T, MAF 0.2%) associations with LDL-C are specific to Hispanic (MAF 1.3%) and Black (MAF 0.6%) populations, respectively and among Asians (MAF 1.5%) alone, 11q13.3 (11:69219641:C:T, MAF 0.2%) was associated with TG. One variant initially passing the novel locus filter for HDL-C (*RNF111*

**Table 1 | Putative novel variants identified in TOPMed and evidence for replication**

Associated lipid phenotype	Novel variant class	Variants (Gene)	Discovery Cohort TOPMed Freeze8 (N = 66,329)			Replication Cohort Meta Analysis (METASOFT) MGB Biobank (N = 25,137); Penn Medicine Biobank (N = 20,079); UK Biobank (N = 424,955)		
			Effect estimate	p-value	MAF	Beta	p-value	Std.Err
LDL-C	Novel locus	12:97352354:T:C	-12.439	$4.88 \times 10^{-09}$	0.003	3.316	$3.62 \times 10^{-01}$	3.634
LDL-C	Novel variant	16:56957451:C:T ( <i>CETP</i> )	-1.568	$2.88 \times 10^{-09}$	0.283	-1.459	$8.74 \times 10^{-84}$	0.075
LDL-C	Novel locus	4:176382171:C:T	-16.086	$2.82 \times 10^{-09}$	0.002	-0.980	$7.80 \times 10^{-01}$	3.514
TC	Novel variant	13:113841051:T:C ( <i>GAS6</i> )	1.731	$1.12 \times 10^{-09}$	0.278	1.262	$1.29 \times 10^{-38}$	0.097
TC	Novel variant	7:137875053:T:C ( <i>CREB3L2</i> )	-4.106	$7.54 \times 10^{-11}$	0.045	-3.538	$7.70 \times 10^{-07}$	0.716
TG	Novel locus	11:69219641:C:T	0.232	$1.98 \times 10^{-09}$	0.002	-0.030	$6.04 \times 10^{-01}$	0.059
TG	Novel variant	13:107551611:C:T ( <i>FAM155A</i> )	0.052	$6.78 \times 10^{-10}$	0.045	0.015	$2.20 \times 10^{-02}$	0.006

Variants identified as novel after comparing with the GWAS catalog and MVP summary statistics for associations with lipid phenotypes, including LDL-C, TC, and TG. All effect estimates are in mg/dL units, except for TG which was log-transformed in analysis thereby representing fractional change. Variants are categorized as novel loci or novel variant (i.e., known locus associated with another lipid phenotype) and the genes assigned to the variants per TOPMed whole genome sequence annotations (WGSa) are listed. Data is provided for the discovery (TOPMed freeze8) and replication cohorts (Imputed datasets from MGB Biobank, Penn Medicine Biobank and UK Biobank). Meta-analysis with the replication cohorts was carried out and the corresponding beta, p-values and standard-errors are provided. All the effect-estimates and p-values are reported from two-sided association testing with all independent samples from each cohort (Discovery-TOPMed: 66,329; Replication-MGB Biobank: 25,137; UK Biobank: 424,955; Penn Biobank: 20,079).

GWAS genome wide association study, MVP million veteran program, LDL-C low-density lipoprotein cholesterol, TC total cholesterol, TG triglycerides, TOPMed trans-omics for precision medicine, WGSa whole genome sequence annotations.

- rs12147665, beta = 8.664, p-value =  $6.51 \times 10^{-10}$ , was in LD ( $r = 0.7$ ) with LIPC p.Thr405Met (rs113298164) which is known to be associated with HDL-C. The lead variant from MVP was 604 kb away from the *RNF111* variant but the rare LIPC missense variant p.Thr405Met was 421 kb away. Conditional analysis accounting for LIPC p.Thr405Met rendered the non-coding variant near *RNF111* variant non-significant (beta = 4.351, p-value =  $2.47 \times 10^{-02}$ ), therefore we reclassified *RNF111* variant as a known-position variant. Ancestry-specific GWAS did not yield additional novel loci beyond our larger trans-ancestry GWAS. The majority of genome-significant single variants were captured by previous lipid GWAS<sup>15</sup>, but ancestry-specific novel-hits are unique to WGS TOPMed data.

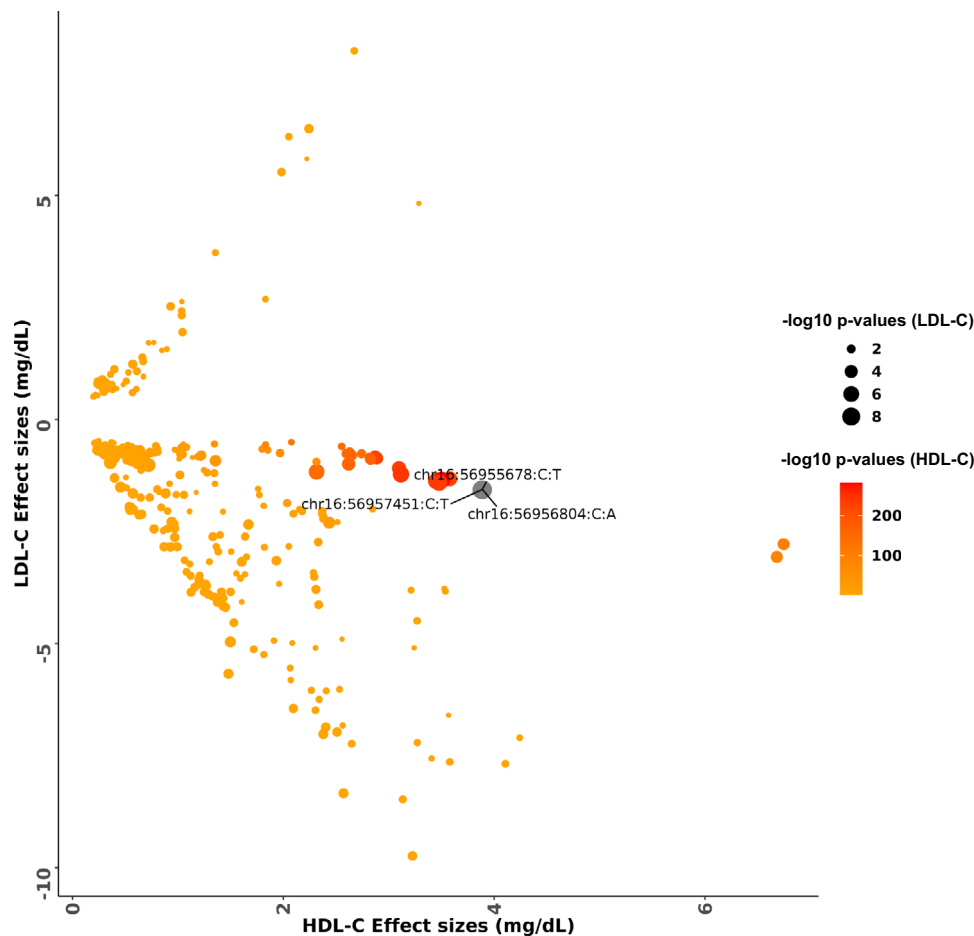
For the single variant GWAS, we pursued replication with two genome-wide array-based genotyped datasets imputed to TOPMed WGS<sup>17,24</sup>: Mass General Brigham (MGB) Biobank (N = 25,137) and Penn Medicine Biobank (N = 20,079)<sup>25,26</sup>, these replication cohorts had diverse ancestry distribution, where non-European samples accounted for 15.77% in MGB Biobank and 51.20% in Penn Medicine Biobank. We also conducted replication using UKB imputed data which accounted for 16.10% of non-European samples (Supplementary Data 6). We brought seven putative novel variants with Type="Italic">p-values <  $5 \times 10^{-9}$  forward for replication. The three common variants, rs183130 (*CETP*), rs7140110 (*GAS6*), and rs73729083 (*CREB3L2*), that were associated with both LDL-C and TC in TOPMed replicated in MGB and UKB along with rs77687061 for TG and two of these (rs183130, rs73729083) replicated in Penn Biobank at an alpha level of 0.05 and consistent direction of effect (Supplementary Data 5). The two variants that were associated in all three replication studies were most significantly associated among African Americans in TOPMed (rs183130: beta = -2.762 mg/dL, p-value =  $5.71 \times 10^{-07}$ ; rs73729083: beta = -3.725 mg/dL, p-value =  $5.25 \times 10^{-07}$ ). We meta-analyzed the single variant replication from the three cohorts and identified three common variants with suggestive p-value ( $5 \times 10^{-5}$ ) (Table 1). Low-frequency variants from specific ancestry groups associated with lipids in TOPMed were not replicated but we cannot rule out the possibility of reduced power due to the general underrepresentation of non-white ancestry groups in the replication data. In exploratory analyses, we extended the same approach for variants discovered to have  $5 \times 10^{-9} < p\text{-value} < 5 \times 10^{-7}$  but did not observe replication (Supplementary Data 7).

### In-silico analysis to gain mechanistic insights from single variant GWAS results

**Prioritization and functional enrichment analysis.** We first mapped the variants to genes and to functional regions using ANNOVAR. Second, we determined gene tissue specificity, relating tissue-specific gene expression with disease-gene associations, using MAGMA. Significantly associated variants were enriched in intronic and intergenic regions (Supplementary Fig. 3). Using GTEx, tissue-specific gene expression was enriched among liver, stomach, and pancreatic tissues (Supplementary Fig. 4) with top tissue-gene sets tabulated in Supplementary Data 8. Using the STRING protein-protein interaction database examining liver-specific genes, we highlight that the HDL-C protein network uniquely harbored metal-ions related genes (*MTIA*, *MTIB*, *MFIF*, *MTIG*, *MTIH*) and anticipated *LCAT-CETP* interactions (Supplementary Fig. 5). Enriched pathways from Reactome, GeneOntology and other curated and canonical pathways (Supplementary Data 9) with a p-value <  $2.5 \times 10^{-06}$  were observed including response to metal ions, lipoprotein assembly, and chylomicron remodeling.

The enrichment analysis was carried out with the full single variant summary statistics, where we identified that most of the prioritized loci/genes were previously documented for lipid associations. Next, we specifically investigated the novel variants that we identified from this study. Out of the seven variants documented in Table 1, four were low frequency variants, 12:97352354:T:C (rs189010847) closest to *NEDD1*, 4:176382171:C:T (rs115489644) closest to *SPCS3*, 11:69219641:C:T (rs74791751) near to *SMIM38*, are all intergenic variants and 13:107551611:C:T (rs77687061) is an intronic variant in *FAM155A*. We did not find any information for these variants in the Open Target Genetics database<sup>27</sup>. Finally, two of the common novel-loci variants (rs183130 and rs7140110) were present in eQTL and sQTL databases<sup>28</sup>, therefore, we performed analysis to determine the correlation among effects and the importance of these variants more in detail.

***CETP* locus, HDL-C, and LDL-C.** *CETP* is a well-recognized Mendelian HDL-C gene and the locus was previously known to be significantly associated with HDL-C, TC, and TG at genome-wide significance<sup>15</sup>. Pharmacologic *CETP* inhibitors have shown strong associations with increased HDL-C but mixed effects for LDL-C reduction in clinical trials<sup>29-32</sup>. We found that the *CETP* locus variant rs183130 (chr16:56957451:C:T, MAF 28.3%, intergenic variant) was associated



**Fig. 3 | Comparison of effects estimates for HDL-C and LDL-C among variants in the *CETP* locus.** The color scale of the data points was based on  $-\log_{10}$  p-values from HDL-C association and the size of each data point was based on  $-\log_{10}$  p-values of LDL-C association. Variants which are genome wide significant with LDL-C are represented as chromosome:position:reference allele:alternate allele.

The effect estimates and p-values were calculated from two-sided genetic association testing performed using SAIGE-QT model, where the model was adjusted for all the covariates; see Methods. HDL-C high-density lipoprotein cholesterol, LDL-C low-density lipoprotein cholesterol.

with reduced LDL-C concentration (beta =  $-1.568$  mg/dL, SE = 0.264,  $p$ -value =  $2.88 \times 10^{-09}$ ). The lead HDL-C-associated variant at the locus, rs3764261 (chr16:56959412:C:A, MAF 30.3%), was associated with 3.5 mg/dL increased HDL-C ( $p$ -value =  $8.03 \times 10^{-283}$ ), and rs183130 was associated with 3.9 mg/dL increased HDL-C ( $p$ -value  $< 1 \times 10^{-284}$ ) as well. Among the ancestry groups analyzed, rs183130 was most significantly associated with LDL-C among those of African ancestry (beta =  $-2.762$  mg/dL,  $p$ -value =  $5.71 \times 10^{-07}$ ) (Supplementary Data 10). We next investigated variants by their HDL-C and LDL-C effects within this locus ( $\pm 500$  kb of rs183130 and rs3764261) (Fig. 3). We identified five variants showing at least suggestive ( $p$ -value  $< 5 \times 10^{-07}$ ) association with both HDL-C and LDL-C. Though variants with strong LD (linkage disequilibrium) existed, ancestry-specific analyses showed that the stronger LDL-C effects were among those of African ancestry.

To better understand the mechanisms for HDL-C and LDL-C effects at the *CETP* locus, we pursued colocalization with eQTLs from three tissues (Liver, Adipose Subcutaneous and Adipose Visceral [Omentum]) from GTEx<sup>28</sup>. We analyzed 5 LDL-C and 441 HDL-C associated ( $p$ -values  $< 5 \times 10^{-07}$ ) variants. We correlated eQTL effect estimates for genes at the locus with lipid outcome effect estimates. Indeed, *CETP* gene expression effects were strongly negatively correlated with HDL-C effects (Liver:  $\rho$   $-0.933$ ,  $p$ -value  $4.01 \times 10^{-17}$ ; Adipose Subcutaneous:  $\rho$   $-0.762$ ,  $p$ -value  $8.87 \times 10^{-12}$ ; Adipose Visceral:  $\rho$   $-0.739$ ,  $p$ -value  $5.52 \times 10^{-10}$ ) (Supplementary Fig. 6). However, *CETP* expression effects were not significantly correlated with LDL-C (Liver:  $\rho$  0.007,  $p$ -value 0.99; Adipose

Subcutaneous:  $\rho$  0.344,  $p$ -value 0.57; Adipose Visceral:  $\rho$   $-0.59$ ,  $p$ -value 0.29). Given the possibility that the observed lack of correlation for LDL-C could be due to reduced power from a limited number of variants attaining a suggestive  $p$ -value ( $< 5 \times 10^{-07}$ ), we repeated the analysis with a subset of 122 nominally significant ( $p$ -value  $< 0.05$ ) LDL-C associated variants in this locus. Indeed, *CETP* gene expression effects were strongly positively correlated with LDL-C effects (Liver:  $\rho$  0.957,  $p$ -value  $2.28 \times 10^{-08}$ ; Adipose Subcutaneous:  $\rho$  0.922,  $p$ -value  $1.34 \times 10^{-15}$ ; Adipose Visceral:  $\rho$  0.868,  $p$ -value  $6.09 \times 10^{-11}$ ).

**GAS6 locus, LDL-C/TG, and TC.** Variants at *GAS6* were previously associated with LDL-C and TG<sup>22,23</sup>, but in our analysis, rs7140110 was now significantly associated with TC. We performed colocalization analysis of the variants  $\pm 500$  Kb from rs7140110 in liver and adipose tissues from GTEx. Across the three lipid-related tissues (liver, adipose subcutaneous, and adipose visceral), strong colocalization was observed in liver for all three lipid phenotypes (TG 46.6%; LDL-C 33.3%; TC 28%). The TG and LDL-C-associated variants were eQTLs for the *GAS6* gene only. However, the TC-associated eQTLs at this locus influenced the *cis* expression of multiple genes, including *GAS6*, anti-sense genes of *GAS6* (AS1, AS2) as well as other genes (i.e., *TFDP1*, *CHAMP1*, *LINCO0565*, *ADPRHL1*, *RASA3*, *UPF3A*, *GRTPI*, *AL442125.1*, *CI3orf46*, *DCUNID2*, *CDC16*, *TMEM255B*, *GRTPI-AS1*, *ATP4B*, *TMCO3*). In addition to *GAS6*, the TC-associated rs7140110 is an sQTL for *TMEM255B* in adipose subcutaneous tissue ( $p$ -value  $5.6 \times 10^{-08}$ ), with

further support from TC colocalization analysis and was not significant for other lipid levels.

**Phenome-wide association with complex traits.** We conducted a phenome-wide association (PheWAS) of 1572 binary complex traits using UK Biobank for the three replicated common variants (16:56957451:C:T (*CETP*); 13:113841051:T:C (*GAS6*); 7:137875053:T:C (*CREB3L2*)) adjusting for PC1–10, age, age<sup>2</sup>, sex, and race. We claimed significance at FDR of 0.05 and identified various complex traits significant, including ischemic heart disease for the *CETP* variant and heart failure/atherosclerosis, hypercholesterolemia traits for *GAS6* variant. The summary statistics from PheWAS analysis for the significant complex traits are tabulated in Supplementary Data 11.

### Rare variant aggregates associated with lipids

**Gene-Centric associations.** We next evaluated the association of aggregated rare (MAF < 1%) variants, linked to protein-coding genes ('gene-centric'). We employed a Bonferroni-corrected significance threshold of  $0.05/20,000 = 2.5 \times 10^{-6}$  for coding and non-coding gene-centric rare variant analyses (Supplementary Fig. 7). We identified 102 coding and 160 non-coding gene-centric rare variant aggregates significantly associated with at least one of the four plasma lipid phenotypes in nonconditional analysis (Supplementary Data 12, 13). We secondarily conditioned our significant aggregate sets on variants individually associated with lipid levels from the GWAS catalog, MVP summary statistics and the TOPMed data. We identified 74 coding and 25 non-coding rare variants aggregates associated with at least one lipid level after conditional analyses (Supplementary Data 14, 15).

Most of the coding gene-centric sets remained significant after secondary conditioning, while a minority of non-coding gene-centric sets remained significant after conditioning. Significant genes identified from coding rare variant analyses included multiple known Mendelian lipid genes including *LCAT*, *LDLR*, and *APOB* (Supplementary Data 13). *RFC2* putative loss-of-function mutations (combined allele frequency < 0.002%) were significantly associated with triglycerides ( $p$ -value  $2 \times 10^{-6}$ ) representing a putative novel association for triglycerides. The *RFC2* aggregate set (plof) was associated with reduced TG (beta = -0.89 for log[*TG*]). The persistently significant regions identified from non-coding rare variant analyses linked to genes included the UTR (untranslated region) for *CETP* and promoter-CAGE (CAGE–Cap Analysis of Gene Expression sites) around *APOA1* for HDL-C, and *APOE* promoter-CAGE, *APOE* enhancer-DHS (DHS–DNase hypersensitivity sites), and *EHD3* promoter-DHS for total cholesterol (Supplementary Data 15). Most of the coding aggregates had larger effects compared to non-coding aggregates, and among the non-coding aggregates *SPC24* non-coding aggregate (enhancer-CAGE) at the *LDLR* locus had the strongest effect for LDL-C (beta = 2.320 mg/dL;  $p$ -value =  $1.75 \times 10^{-05}$ ).

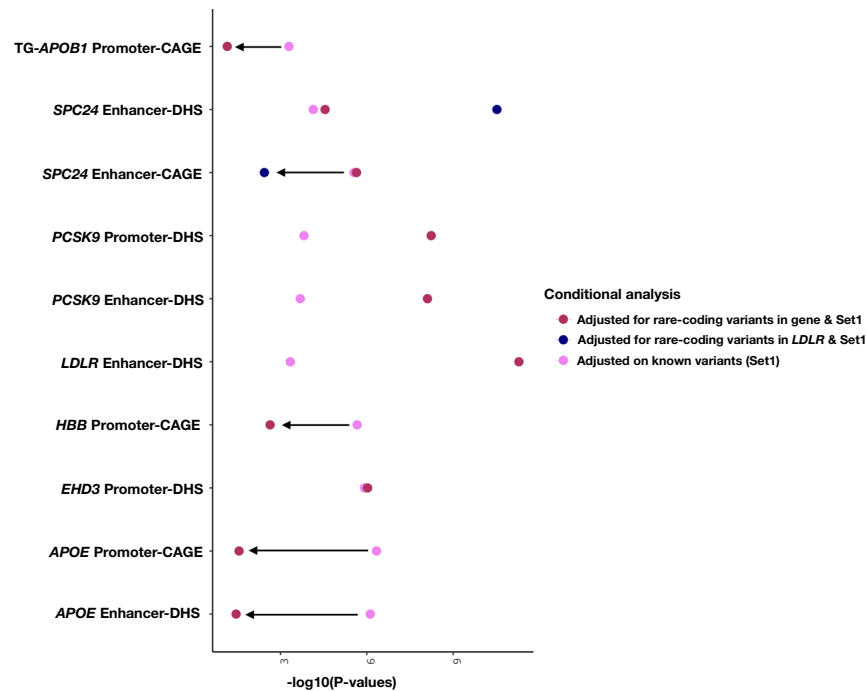
We analyzed the UK Biobank whole genome sequences among ~130 K participants to provide evidence of replication for the significant coding and non-coding aggregate sets. We used a Bonferroni-corrected significance threshold based on the number of genes tested in each type of aggregate-based test. For gene centric-coding aggregates, we conducted replication of 21 genes ( $p$ -value <  $0.05/21 = 2.38 \times 10^{-03}$ ) and for non-coding aggregates we replicated the findings from 13 genes ( $p$ -value <  $0.05/13 = 3.85 \times 10^{-03}$ ). At Bonferroni significance, 71% and 62% of genes replicated for at least one coding and non-coding aggregate set, respectively (Supplementary Data 14, 15). We observed that most of the Mendelian lipid genes replicated for coding aggregates including *ABCA1*, *ABCG5*, *LCAT*, *APOB*, *LDLR*, *PCSK9*, and *LPL*. For the non-coding aggregate set, the most significant replications were observed for the *APOB*, *LDLR* (*SPC24*), and *PCSK9* loci, further corroborating the observation that both coding and noncoding rare variant signals contribute to variation in lipid levels at these loci.

**Region-based associations.** We also performed unbiased region-based rare variant association analyses tiled across the genome with both static and dynamic window sizes. We first evaluated 2.6 M regions statically at 2 kb size and 1 kb window overlap by the sliding window approach. Statistical significance was assigned at  $0.05/(2.6 \times 10^6) = 1.88 \times 10^{-08}$ . We identified 28 significantly associated windows with at least one lipid phenotype. After conditioning on variants individually associated with the corresponding lipid phenotype, we identified two regions at *LDLR* still significantly associated with both total cholesterol and LDL-C, although these regions included both intronic and exonic variants (Supplementary Data 16). *LDLR* intron 1, which encodes *LDLR-ASI* (*LDLR* antisense RNA 1) on the minus strand, had suggestive evidence for association with TC ( $p$ -value  $3.17 \times 10^{-6}$ ) with -2.76 mg/dL reduction in TC. A prior study identified that a common variant (rs6511720, MAF 0.11) in *LDLR* intron 1 is associated with increased *LDLR* expression in a luciferase assay and reduction in LDL-C<sup>33</sup>. When adjusting for rs6511720, the significance improved ( $p$ -value  $1.43 \times 10^{-8}$ ) with -3.35 mg/dL reduction in TC.

For dynamic window scanning of the genome, we implemented the SCANG method<sup>34</sup>. The SCANG procedure accounts for multiple testing by controlling the genome-wide error rate (GWER) at 0.1<sup>34</sup>. In the dynamic window-based workflow, STAAR-O detected 51 regions significantly associated with at least one lipid phenotype after conditioning on known variants (Supplementary Data 17). Most of the regions mapped to known Mendelian lipid genes, including *LCAT* ( $8.7 \times 10^{-13}$ ) for HDL-C, and *LDLR* ( $2.4 \times 10^{-28}$ ,  $7.3 \times 10^{-26}$ ) and *PCSK9* ( $2.9 \times 10^{-12}$ ,  $5.5 \times 10^{-12}$ ) for LDL-C and TC, respectively. Exon 4 aggregates of *LDLR* were specifically associated with 20 mg/dL increase in LDL-C. *PCSK9* Exon2-Intron2 region spanning chr1:55043782–55045960 had significantly reduced LDL-C by 6 mg/dL ( $p$ -value =  $3 \times 10^{-13}$ ), and the effect persisted even with only Intron 2 rare variants of *PCSK9* (-5 mg/dL,  $p$ -value =  $2 \times 10^{-8}$ ). Strikingly, in secondary analyses, we found evidence for very large effects for rare variants in *LDLR* Introns 2 and 3 (+21 mg/dL,  $p$ -value =  $7 \times 10^{-4}$ ) and *LDLR* Introns 16 and 17 (+17 mg/dL,  $p$ -value = 0.02), similar to rare coding *LDLR* mutations. While 32 of the significant dynamic windows also included exonic regions, there were also several dynamic windows significantly independently associated with lipids not containing exonic regions. For example, four non-coding windows (two overlapping) at 2p24.1, which harbors the Mendelian *APOB* gene, were significantly associated with LDL-C. Intronic non-coding regions were associated with both LDL-C and TC -associated windows at *LPAL2-LPA-SLC22A3*; for example, *LPAL2* Intron 3 was associated with a 3.7 mg/dL increase in TC. Non-coding TC-associated significant dynamic windows were near *TOMM40/APOE*. One rare variant signal observed was at *TOMM40* Intron 6, where the 'poly-T' variant in this region is on the *APOE4* haplotype and influences expressivity for Alzheimer's disease age-of-onset<sup>35,36</sup>. For HDL-C, we identified significant non-coding windows at an intergenic region near *LPL* and *CD36* Intron 4. In the generation of the spontaneously hypertensive rat model, the deletion of intron 4 in *CD36* with resultant *CD36* deficiency has been mapped to defective fatty acid metabolism in this model<sup>37</sup>. Several regions significant in SCANG were not even nominally significant in burden association analyses indicating the likelihood of causal variants with bidirectional effects.

We replicated 28 sliding and 51 dynamic window aggregate sets using UKB whole genomes, at a Bonferroni-corrected alpha threshold of 0.05/no. of regions for each approach separately. At Bonferroni significance, 61% of the regions from each of the sliding window ( $p$ -value <  $0.05/28 = 1.79 \times 10^{-03}$ ) and dynamic window ( $p$ -value <  $0.05/51 = 9.80 \times 10^{-04}$ ) approaches significantly replicated (Supplementary Data 16, 17). Multiple regions linked to *LDLR*, *PCSK9*, *CETP*, *APOC3*, and *ABCA1* were highly significant.

Several gene-centric non-coding aggregates associated with lipids near known monogenic lipid genes but mapped to another gene at the



**Fig. 4 | Conditional analysis of coding rare-variants from the same gene and a near-by gene.** Non-coding rare variant sets significantly associated with TC and TG after the conditional analysis on known variants are shown with additional adjustment on rare-coding variants. The additional adjustment for rare-coding variants were carried out for the same gene of the aggregate set and for certain gene aggregates (SPC24) the conditional analysis was carried out with a nearby Mendelian gene. After adjusting for rare-coding variants and known variants, *EHD3* signal drops minimally, whereas signal from *PCSK9* (promoter-DHS, enhancer-DHS), *LDLR*-loci (enhancer-DHS, *SPC24* enhancer-DHS) enhances significantly.

*APOB1*, *SPC24* (enhancer-CAGE), *HBB* and *APOE* signal drops after the conditional analysis on rare-coding variants. The different colored dots on the plot represents the conditional STAAR-O p-values when adjusting for known variants (Set1) and rare-coding variants of the same or near-by gene. The p-values were calculated from two-sided aggregate testing preformed using STAAR gene-centric model, where the model was adjusted for all the covariates; see Methods. STAAR variant-Set Test for Association using annotation information, TC total cholesterol, TG triglycerides, CAGE cap analysis of gene expression, DHS DNase hypersensitivity.

locus via annotations. Therefore, we performed downstream conditional analyses adjusting the gene-centric non-coding results for rare coding variants (MAF < 1%) within known lipid monogenic genes (Supplementary Data 18). When accounting for both common and rare coding variants at the nearby familial hypercholesterolemia *LDLR* gene, *SPC24*-enhancer DHS was significantly associated with total cholesterol ( $p$ -value =  $3.01 \times 10^{-11}$ ) and with suggestive evidence for LDL-C ( $p$ -value =  $1.57 \times 10^{-06}$ ). In a similarly adjusted model, *LDLR*-enhancer-DHS showed a strong association with TC ( $p$ -value  $5.18 \times 10^{-12}$ ). When adjusting for known common variants as well as rare coding variants in *PCSK9*, both *PCSK9*-enhancer DHS and *PCSK9*-promoter DHS were significantly associated with total cholesterol (Fig. 4, Supplementary Fig. 8). Through this procedure, *CETP* UTR retained significance for its independent association with HDL-C as well as the putatively novel gene *EHD3*-promoter DHS association with TC. However, the non-coding gene-centric *APOC3* and *APOE* associations were rendered non-significant for HDL-C and TC, respectively.

Since we cannot rule out the possibility of reduced power for genome-wide rare variant analyses, we leveraged current knowledge of 22 Mendelian lipid genes for more focused exploratory analyses<sup>14</sup>. We validated most genes in rare variant coding analyses. The genes with the strongest coding signals typically had at least nominal evidence of gene-centric non-coding rare variant associations (Supplementary Data 19, Supplementary Fig. 9). When rare coding variants were introduced into the model, the evidence for non-coding rare variant associations were largely unchanged. Our findings expanding the currently described genetic basis for hypercholesterolemia to include rare non-coding variation at *LDLR* and *PCSK9* (Fig. 5).

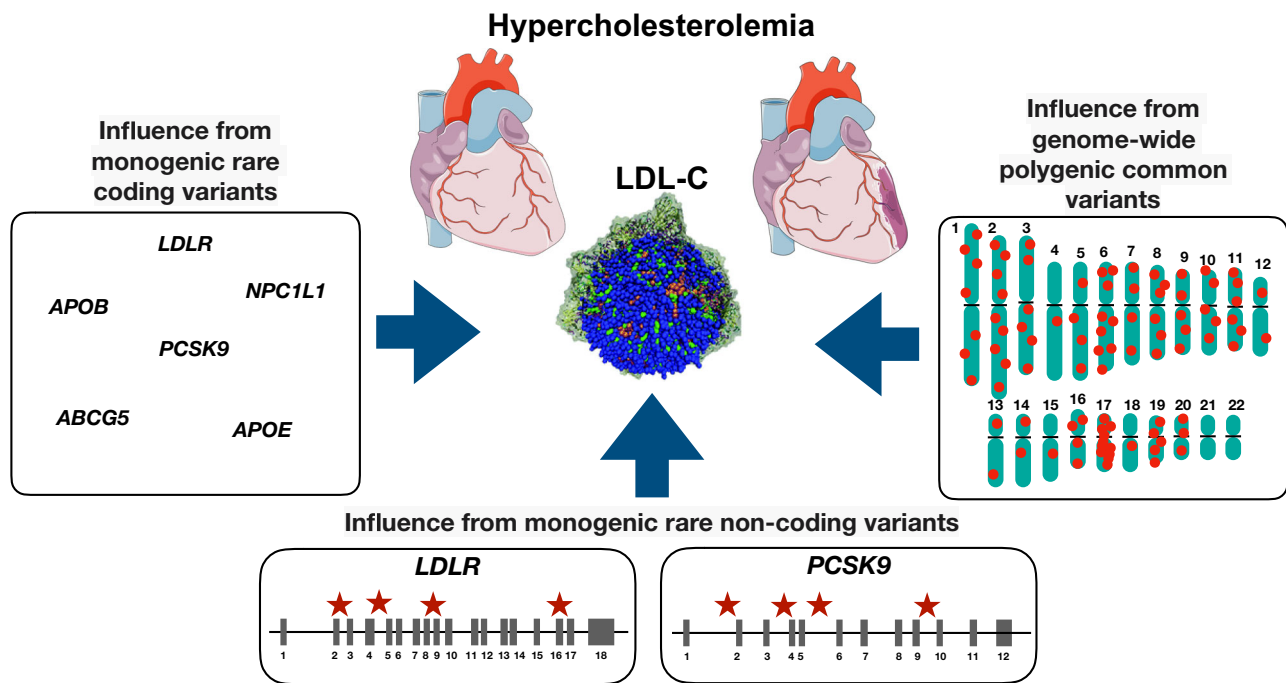
### Heritability contributions from rare variants

To understand the contribution of rare variants towards lipid trait heritability, we examined heritability of lipids by variant allele frequency across three ancestral samples (White, Black, and Hispanic) in TOPMed. We calculated trait heritability using Greml-LDMS<sup>38</sup> following the steps as implemented by Wainschein et al.<sup>39</sup>. Using the TOPMed WGS, we grouped the variants into 4 MAF bins for the three ancestral samples. In each MAF bin, we grouped variants based on the LD scores into four quartiles and calculated variance contributed by the SNPs ( $h^2$ ) for each of the lipids using unrelated individuals from each ancestral group (Supplementary Fig. 10) and set negative estimate to zero. We observed that rare variants from the lower MAF bins contributed to trait heritability but have large standard errors (Supplementary Data 20). We observed an increase in  $h^2$  values including WGS variants relative to estimates obtained from array-genotypes as reported by Cadby et al.<sup>40</sup> for the European samples. We also compared the  $h^2$  estimates from all the variants from WGS TOPMed cohort against array-genotypes captured in MGB Biobank to understand the differences contributed by these two sequencing methods. As expected, the  $h^2$  estimates from array-genotypes were reduced corresponding to missing heritability from the lower MAF bins captured by WGS. The heritability estimates from array-genotypes were markedly higher for European samples relative to African and Hispanic sample sets indicating that WGS better captured heritability for the latter groups.

### Discussion

Conducting one of the largest population-based WGS association analyses, we now simultaneously interrogate and establish a common, rare coding, and rare non-coding variant model for a complex trait. Utilizing 66,329 diverse individuals with deep-coverage WGS, we





**Fig. 5 | Influence of common and rare variants with hypercholesterolemia.** In addition to monogenic contributions from rare variants in Mendelian hypercholesterolemia genes, multiple genome-wide significant LDL-C-associated common variants also yield a polygenic basis for hypercholesterolemia. In the present work, we now identify rare non-coding variants in proximity of Mendelian

hypercholesterolemia genes, specifically *LDLR* and *PCSK9*, that also contribute to the genetic basis of hypercholesterolemia. Parts of the figure were generated using pictures from Servier Medical Art. Servier Medical Art by Servier is licensed under a Creative Commons Attribution 3.0 Unported License (<https://creativecommons.org/licenses/by/3.0/>). LDL-C low-density lipoprotein cholesterol.

interrogated 428 M variants with plasma lipids expanding the allelic series to rare non-coding variants, often within introns, of Mendelian lipid genes with prior robust rare coding variant support. Our observations have important implications for plasma lipids as well as the genetic basis of complex traits more broadly.

WGS of diverse ancestries enables both allelic and locus heterogeneity for complex traits. Population genetic analyses have largely been enriched for individuals of European descent<sup>41</sup>. Genetic association of plasma lipids using arrays or whole exome sequencing among Europeans have yielded several important insights regarding plasma lipids and the causal determinants of CAD<sup>4,5,42–44</sup>. Similar increasingly larger studies among non-Europeans have often yielded new genetic loci and sometimes new genes, such as *PCSK9*<sup>1,15,16,45,46</sup>. Such differences have also led to concerns about the use of polygenic risk scores gleaned from much larger European GWAS of complex traits for non-Europeans<sup>47</sup>. Aided by the availability of WGS data, we identify new putative loci associated with lipids in non-Europeans. Furthermore, our study enabled the discovery of several novel alleles at known loci, with richly distinct allelic heterogeneity across ancestry groups. For example, HDL-C-raising *CETP* locus variants linked to *CETP* gene expression were only associated with LDL-C reduction among those of African ancestry. While all pharmacologic *CETP* inhibitors increase HDL-C, only those that decrease LDL-C also reduce cardiovascular disease risk<sup>29–32</sup>. Given the contribution of genetic differences, clinical trials with more diverse samples would show insights.

Our study now provides increasingly robust evidence for a rare non-coding variant model for complex traits. Our rare non-coding variant associations in both gene-centric and sliding window models were largely restricted to the introns of Mendelian lipid genes with prior robust rare coding variant support consistent with biologic plausibility<sup>48</sup>. Rare intronic variants, often impacting splicing, have been previously implicated in afflicted Mendelian families or small exceptional case series, often through candidate gene approaches<sup>49–52</sup>. We discovered one example of a rare non-coding signal without prior

rare coding support—i.e., *EHD3* which also nominally replicated in the independent UKB WGS cohort. We obtained estimates of phenotypic effect using burden tests. For most regions, even nominal significance was not detected using burden testing indicating the likelihood of variants with bidirectional effects further complicating clinical interpretation. When burden signals were detected, observed effects were typically larger than common non-coding variants and less than rare coding variants, with the exception of *LDLR*, consistent with whole genome mutational constraint models<sup>53–55</sup>.

The detection of independent rare non-coding variant signals has remained elusive largely due to limited sample sizes with requisite WGS and limitations in the interpretation of rare non-coding variation functional consequence. Previously, we used annotated functional non-coding sequence in 16,324 TOPMed participants, and found that rare non-coding gene regions associated with lipid levels, but they were not independent of individually associated single variants<sup>14</sup>. Using STAAR, we observed putative rare non-coding variant associations for lipids independent of individual variants associated with lipids in TOPMed.

WGS can improve diagnostic yield beyond the current standard of next-generation gene panel sequencing for dyslipidemias. A very small fraction with severe hypercholesterolemia and features consistent with strong genetic predisposition have a familial hypercholesterolemia variant in *LDLR*, *APOB*, or *PCSK9*<sup>56,57</sup>. The presence of familial hypercholesterolemia variants is independently prognostic for CAD, beyond lipids, and merits the consideration of more costly lipid-lowering medications<sup>56–59</sup>. We now observe that rare *LDLR* variants in Introns 2, 3, 16, and 17 lead to -0.5 standard deviation increase in LDL-C, approximating effects observed with clinically reported exonic familial hypercholesterolemia variants in *LDLR*<sup>59</sup>. Small studies have indicated the possibility of rare intronic *LDLR* variants causing familial hypercholesterolemia due to altered splicing, which we now observe in our unbiased population-based WGS study<sup>60,61</sup>. A WGS approach to lipid disorders, particularly for familial

hypercholesterolemia, will markedly improve the diagnostic yield beyond existing limited approaches.

Our dynamic window approach may also improve the clinical curation of exonic variants. Among the data used to curate exonic variants is the use of *in silico* functional prediction tools<sup>62</sup>. Although evolutionary constraint measures are typically employed, such tools are largely agnostic to functional domain. As it relates to lipids, disruptive *APOB* and *PCSK9* exonic variants can lead to strikingly opposing directions with large effects for LDL-C depending on locations<sup>1,8,63,64</sup>. Using SCANG<sup>34</sup>, we detect a significant association with large effect for *LDLR* Exon 4 itself. This observation supports the pathogenicity of *LDLR* Exon 4 disruptive variants among patients with severe hypercholesterolemia. The majority of familial hypercholesterolemia variants worldwide occur in Exon 4 of *LDLR*<sup>65–68</sup>. Conventional rare coding variant analyses aggregate all exonic variants for a transcript. Here, we demonstrate an opportunity for exon-level rare variant association testing.

Our discovery analyses with replication as well as heritability assessment are consistent with the notion that both rare coding and non-coding alleles, not well-captured by genome-wide arrays. Furthermore, we observe that heritability gains relative to genome-wide genotyping arrays are more significant for individuals of European-ancestry likely indicative of Eurocentric array designs. A tradeoff for WGS, however, is the greater cost. However, as costs continue to decrease as well as cheaper WGS implementations via reduced coverage, cost may no longer be a downside.

Our study has important limitations. First, while our study is large for a WGS study by contemporary standards, it is dwarfed by existing GWAS datasets limiting power for novel discovery. Nevertheless, by using WGS in diverse ancestries, we can study hundreds of millions new variants. Second, prediction of rare non-coding variation consequence to prioritize causal variants remains a challenge thereby limiting power<sup>69</sup>. The striking difference for most STAAR and burden results also highlights bidirectional effects for rare non-coding variants within the same region and further challenges for clinical utility. Third, given the paucity of multi-ancestral WGS datasets with lipids, our analyses are largely restricted to TOPMed and replication to European rich UK Biobank WGS data. For single variant associations, we pursued TOPMed-imputed GWAS datasets but were limited by the lack of ancestral diversity. As TOPMed is a consortium of multiple different cohorts, we demonstrate consistencies by cohort. Furthermore, rare variant non-coding signals were largely restricted to regions with rare variant coding signals supporting biological plausibility.

In conclusion, using WGS and lipids among 66,329 ancestrally diverse individuals we expand the catalog of alleles associated with lipids, including allelic heterogeneity at known loci and locus heterogeneity by ancestry. We characterize the common, rare coding, and rare non-coding variant model for lipids and replicated the results. Lastly, we now demonstrate a monogenic-equivalent model for rare *LDLR* intronic variants predisposing to marked alterations in LDL-C, currently not recognized in current population or clinical models for LDL-C.

## Methods

### Dataset

**Contributing studies.** The discovery cohort includes the whole genome sequenced (WGS) data of 66,329 samples from 21 studies of the Trans-Omics for Precision Medicine (TOPMed) program with blood lipids available<sup>17</sup>. The overall goal of TOPMed is to generate and use trans-omics, including whole genome sequencing, of large numbers of individuals from diverse ancestral backgrounds with rich phenotypic data to gain novel insights into heart, lung, blood, and sleep disorders. The Freeze 8 data includes 140,306 samples out of which 66,329 samples qualified with lipid phenotype. Freeze 8 dataset passed the central quality control protocol implemented by the TOPMed

Informatics Research Core (described below) and was deposited in the dbGaP TOPMed Exchange Area.

The studies included in the current dataset, along with their abbreviations and sample sizes, contains the Old Order Amish (Amish,  $n = 1083$ ), Atherosclerosis Risk in Communities study (ARIC,  $n = 8016$ ), Mt Sinai BioMe Biobank (BioMe,  $n = 9848$ ), Coronary Artery Risk Development in Young Adults (CARDIA,  $n = 3,056$ ), Cleveland Family Study (CFS,  $n = 579$ ), Cardiovascular Health Study (CHS,  $n = 3,456$ ), Diabetes Heart Study (DHS,  $n = 365$ ), Framingham Heart Study (FHS,  $n = 3992$ ), Genetic Studies of Atherosclerosis Risk (GeneSTAR,  $n = 1757$ ), Genetic Epidemiology Network of Arterio-pathy (GENOA,  $n = 1046$ ), Genetic Epidemiology Network of Salt Sensitivity (GenSalt,  $n = 1772$ ), Genetics of Lipid-Lowering Drugs and Diet Network (GOLDN,  $n = 926$ ), Hispanic Community Health Study - Study of Latinos (HCHS\_SOL,  $n = 7714$ ), Hypertension Genetic Epidemiology Network and Genetic Epidemiology Network of Arterio-pathy (HyperGEN,  $n = 1853$ ), Jackson Heart Study (JHS,  $n = 2847$ ), Multi-Ethnic Study of Atherosclerosis (MESA,  $n = 5290$ ), Massachusetts General Hospital Atrial Fibrillation Study (MGH\_AF,  $n = 683$ ), San Antonio Family Study (SAFS,  $n = 619$ ), Samoan Adiposity Study (SAS,  $n = 1182$ ), Taiwan Study of Hypertension using Rare Variants (THRV,  $n = 1982$ ) and Women's Health Initiative (WHI,  $n = 8263$ ) (Please see Supplementary Note 1 for additional details). The multi-ancestral data set included individuals from White (44%), Black (26%), Hispanic (21%), Asian (7%), and Samoan (2%) ancestries. Study participants granted consent per each study's Institutional Review Board (IRB) approved protocol. Secondly, these data were analyzed through a protocol approved by the Massachusetts General Hospital IRB. Supplementary Data 1 details the number of samples across different studies and ancestral group.

The replication cohorts for single variant GWAS include TOPMed-imputed genome-wide array data from the Mass General Brigham (MGB), Penn Medicine Biobanks and UK Biobank (UKB) imputed data which consist of 25,137, 20,079, and 424,955 samples respectively<sup>25,26,70</sup>. The replication cohort for rare variant aggregates test include UKB whole genome sequenced data which consists of a subset of 133,360 UKB participants, where we removed unconsented and related individuals. We curated the MGB Biobank and Penn Medicine Biobank phenotype data from the corresponding electronic health record databases in accordance with corresponding institutional IRB approvals. The UKB data included volunteer residents of the UK aged 40–69 and were recruited between 2006 and 2010. Consent was previously obtained from each participant regarding storage of biological specimens, genetic sequencing, access to all available electronic health record (EHR) data, and permission to recontact for future studies. All UKB participants gave written informed consent per UKB primary protocol. The MGB Biobank consists of 54%, Penn Medicine Biobank consist of 52% and UK Biobank imputed data consist of 54% of female samples and average ages were 55.89, 58.35 and 56.55 years, respectively (Supplementary Data 6).

**Phenotypes.** The primary outcomes in this study included LDL cholesterol (LDL-C), HDL cholesterol (HDL-C), total cholesterol (TC), and triglycerides (TG) phenotypes. LDL-C was either directly measured or calculated by the Friedewald equation when triglycerides were  $<400$  mg/dL. Given the average effect of lipid lowering-medicines, when lipid-lowering medicines were present, we adjusted the total cholesterol by dividing by 0.8 and LDL-C by dividing by 0.7, as previously done<sup>14</sup>. Triglycerides remained natural log transformed for analysis. Fasting status was accounted for with an indicator variable.

We harmonized the phenotypes across each cohort<sup>18</sup> and inverse rank normalization of the residuals of each race within each cohort scaled by the standard deviation of the trait and adjusted for covariates<sup>12</sup>. We included covariates such as age, age<sup>2</sup>, sex, PCI-11, study-groups as well as Mendelian founder lipid variants *APOB*

p.R3527Q and *APOC3* p.R19X for the Amish cohort<sup>7,8,71</sup>. Supplementary Data 1 provides the distributions of each of the four lipid phenotypes by cohort, ancestral groups, and gender. For the UK Biobank, we curated the first instance of the four lipids (data field numbers: HDL-C-30760; LDL-C-30780; TC-30690; TG-30870). The lipid measurements from mmol/L were converted to mg/dL by multiplying TG measurements by 88.57 and for other lipids by multiplying by 38.67. We executed similar steps of phenotype harmonization and normalization for the replication cohorts. In addition, we adjusted the MGB Biobank for study-center and array-type, and Penn Medicine Biobank for ancestry and BMI in addition to the other common covariates.

**Genotypes.** Whole genome sequencing of goal >30X coverage was performed at seven centers (Broad Institute of MIT and Harvard, Northwest Genomics Center, New York Genome Center, Illumina Genomic Services, PSOMAGEN [formerly Macrogen], Baylor College of Medicine Human Genome Sequencing Center, and McDonnell Genome Institute [MGI] at Washington University). In most cases, all samples for a given study within a given Phase were sequenced at the same center (Supplementary Note 1). The reads were aligned to human genome build GRCh38 using a common pipeline across all centers (BWA-MEM).

The TOPMed Informatics Research Core at the University of Michigan performed joint genotype calling on all samples in Freeze 8. The variant calling “GotCloud” pipeline ([https://github.com/statgen/topmed\\_variant\\_calling](https://github.com/statgen/topmed_variant_calling)) is under continuous development and details on each step can be accessed through TOPMed website for Freeze<sup>8,17</sup>. The resulting BCF files were split by study and consent group for distribution to approved dbGaP users. Quality control was performed centrally by the TOPMed IRC and the TOPMed Data Coordinating Center (DCC) as previously described<sup>17</sup>. Briefly, the two sequence quality criteria used in freeze 8 are: estimated DNA sample contamination below 10%, and 95% or more of the genome covered to 10× or greater. The variant filtering in TOPMed Freeze 8 is performed by (1) first calculating Mendelian consistency scores using known familial relatedness and duplicates, and (2) training a Support Vector Machine (SVM) classifier between known variant sites (positive labels) and Mendelian inconsistent variants. A small number of sex mismatches were detected as annotated females with low X and high Y chromosome depth or annotated males with high X and low Y chromosome depth. These samples were either excluded from the sample set to be released on dbGaP or their sample identities were resolved using information from prior array genotype comparisons and/or pedigree checks. Details regarding WGS data acquisition, processing and quality control vary among the TOPMed data freezes. Freeze-specific methods are described on the TOPMed website (<https://www.nhlbiwgs.org/data-sets>) and in documents included in each TOPMed accession released on dbGaP. The VCF/BCF files were converted to GDS (Genomic Data Structure) format by the DCC and were deposited into the dbGaP TOPMed Exchange Area.

The genetic relationship matrix (GRM) is an N\*N matrix of relatedness information of the samples included in the study and was computed centrally using ‘PC-relate’ R package (version: 1.24.0)<sup>72</sup>. Using the ‘Genesis’ R package (version:2.20.1)<sup>73</sup> we generated subsetted GRM for the samples with plasma lipid profiles. The GDS files with the variants were annotated internally by curating data from multiple database sources using Functional Annotation of Variant–Online Resource (FAVOR (<http://favor.genohub.org>))<sup>13</sup>. This study used the resultant aGDS (annotation GDS) files.

The MGB Biobank replication cohort was genotyped using three different arrays (Multiethnic Exome Global (Meg), Human multi-ethnic array (Mega), and Expanded multi-ethnic genotyping array (Megex)), and we separately imputed the data using TOPMed imputation server with default parameters<sup>74,75</sup>. This study applied the Version-r2 of the imputation panel, it includes 97,256 reference samples and ~300 M

genetic variants. The Illumina Global Screening array was used to genotype the Penn Medicine Biobank. Penn Medicine Biobank TOPMed imputation was performed using EAGLE<sup>75</sup> and Minimac<sup>76</sup> software. For this study, we downloaded variants that passed a min R<sup>2</sup> threshold of 0.3. The TOPMed imputation panel is robust, built from 97,256 deeply sequenced human genomes and contains 308,107,085 genetic variants from multi-ethnic samples. Imputation was performed in independent non-overlapping samples agnostic to phenotypes. The UKB imputed data was derived using merged UK10K<sup>77</sup>, 1000 Genomes phase2 reference panels and was combined to the Haplotype reference Consortium<sup>78</sup> (HRC) using IMPUTE 4 program (<https://jmarchini.org/software/>). The UKB WGS data consist of whole genomes of 150,119 UKB participants with an average coverage of 32.5X. We used joint called VCFs from GraphTyper, which consist of 710,913,648 variants<sup>79</sup>. We used VCFs provided on the UK Biobank and conducted all the analysis in UKB Research Analysis Platform (UKB RAP).

### Single variant association

We performed genome-wide single variant association analyses for autosomal variants with minor allele frequency (MAF) >0.1% across the dataset with each of the four lipid phenotypes. We implemented the SAIGE-QT<sup>80</sup> method, which employs fast linear mixed models with kinship adjustment, in Encore (<https://encore.sph.umich.edu/>) for single variant association analyses. We additionally adjusted the model for covariates (PCI-PC11, age, sex, age<sup>2</sup>, and study-groups [cohort-race subgrouping]).

We conducted single variant association replications for putative novel variants. After comparing the results with published lipid GWAS summary statistics, we filtered putative novel GWAS variants based on a stringent whole genome-wide significant threshold ( $\alpha = 5 \times 10^{-9}$ )<sup>81</sup>. Replication was performed in the MGB, Penn Medicine Biobanks and UK Biobank where linear regression models were fitted and adjusted for covariates as indicated above. In addition, we adjusted the MGB Biobank for study recruitment center and array and Penn Medicine Biobank for ancestry and BMI. In the MGB Biobank, we selected lipid concentrations closest to the sample acquisition time point and adjusted for statins if prescribed within one year prior to sample acquisition. In the Penn Biobank, we utilized each participant’s median lipid concentration for replication; statins prescribed prior to lipid concentration used were adjusted in the models. In addition, we carried out meta-analysis using fixed effects model based on inverse-variance-weighted effect size for the two replication cohorts using METASOFT<sup>82</sup>.

### Rare variant association test

We performed rare variant association (RVA) using the Variant-Set Test for Association using Annotation infoRmation (STAAR) pipeline<sup>13,83</sup>. STAARpipeline is a regression-based framework that permits adjustment of covariates, population structure, and relatedness by fitting linear and logistic mixed models for quantitative and dichotomous traits<sup>83–85</sup>. We chose STAAR to leverage the annotation information and associated scores that were available for TOPMed Freeze 8 data to incorporate the analysis of rare non-coding variants from whole genome sequencing. The method implements genome-wide scanning of rare variants (MAF <0.01) in gene-centric and region-based workflows. For each variant set, STAARpipeline calculates a set-based *p*-value using the STAAR method, which increases the analysis power by incorporating multiple in silico variant functional annotation scores capturing diverse genomic features and biochemical readouts<sup>13</sup>. We aggregated rare variants into multiple groups for coding and non-coding analyses. For the coding region, we defined five different aggregate masks of rare variants 1) plof (putative loss-of-function), plof-Ds (putative loss-of-function or disruptive missense), missense, disruptive-missense, and synonymous. For the non-coding regions, we used seven rare variant masks: (1) promoter-CAGE (promoter variants

within Cap Analysis of Gene Expression [CAGE] sites<sup>86</sup>), (2) promoter-DHS (promoter variants within DNase hypersensitivity [DHS] sites<sup>87</sup>), (3) enhancer-CAGE (enhancer within CAGE sites<sup>88,89</sup>), (4) enhancer-DHS (enhancer variants within DHS sites<sup>87,89</sup>), (5) UTR (rare variants in 3' untranslated region [UTR] and 5' UTR untranslated region), (6) upstream, and (7) downstream. Detailed explanations of the regions defined based on these masks is discussed within STAARpipeline<sup>13,83</sup>.

In the gene-centric workflows, for both coding (within exonic boundaries) and non-coding (promoter: +/-3 kb window of transcription starting site (TSS), enhancer: GeneHancer predicted regions, UTR (both 5' and 3' UTR regions)/upstream/downstream: GENCODE Variant Effect Predictor (VEP) categories) regions, we considered only genes with at least two rare variants (i.e., 18,445 genes in all 22 autosomes). In the region-based workflows, we implemented two protocols: (1) a 'sliding window' approach, where we aggregated rare variants within 2-kb sliding windows and with 1-kb overlap length, and (2) a 'dynamic window' approach, where we executed SCANG<sup>34</sup> method and aggregated dynamically variant-sets between 40–300 variants per set, where the method systematically scans the whole genome with overlapping windows of varying sizes. The STAARpipeline R-package implements multiple rare-variants aggregate tests including SKAT<sup>90</sup>, Burden<sup>91</sup> and ACAT<sup>92</sup> and integrates them as STAAR-O<sup>13,83</sup>. We performed gene-centric and region-based rare variant tests using annotated GDS files of TOPMed.

We completed aggregate tests as three-step process. In the first step, we fitted a null model using glmkin() function. The null model was fitted for each of the four lipid phenotypes adjusted for all covariates and relatedness except the genotype of interest. In the second step, we ran genome-wide gene-centric and region-based rare-variant aggregate tests. The third step directed conditional analyses, where the results were adjusted for previously known significantly lipid-associated (i.e.,  $p < 5 \times 10^{-8}$  in external datasets) individual variants from GWAS Catalog<sup>93</sup> and Million Veterans Program (MVP)<sup>15</sup> GWAS summary statistics. To obtain effect estimates of significant aggregate sets, we associated the cumulative genotypes (binary scores) based on the variants forming the aggregates and used Glim.Wald test from GMMAT R package<sup>83</sup> (version 1.3.1). For significantly associated window-based rare variant aggregations, we trimmed the exonic variants and estimated the effects with only non-coding variants.

For the rare variant replication in UKB WGS data, we curated the rare variant aggregate sets in UKB RAP for the gene-centric coding/non-coding and region-based significant sets and applied STAAR workflow as demonstrated by the STAARpipeline (<https://github.com/xihaoli/STAARpipeline>) and describe above.

### Computational mining of single variant GWAS

**Gene-set enrichment using FUMA.** We performed enrichment analysis with single variant GWAS summary stats from the four lipids using FUMA<sup>94</sup> (version 1.3.7) with default parameters and significance at  $5 \times 10^{-9}$ . FUMA is an integrated platform which efficiently facilitates functional mapping and enrichment of GWAS-associated genes using multiple useful resources. The method uses 18 different biological data repositories and tools to process GWAS data. We additionally used MAGMA<sup>95</sup> (version 1.08) gene-based analysis enrichment workflow within FUMA with the complete GWAS summary data for eQTL based tissue enrichment. The functionally prioritized genes were visualized based on their protein-protein interaction networks using the STRING database<sup>96</sup>.

**CETP and GAS6 gene expression and lipid trait colocalization.** We studied the correlation of LDL-C and HDL-C effects with eQTL effects at chromosome 16q13, which includes *CETP* and correlation of LDL-C and TC with eQTLs at rs7140110 of *GAS6*. We downloaded GTEx eQTL build 38 (version8) data for liver, adipose subcutaneous, and adipose visceral (omentum) tissues from GTEx on 16/APR/2020<sup>97</sup>.

For the *CETP* variant analysis, we selected eQTLs with nominal significance ( $p$ -value  $< 0.05$ ) and utilized the eQTL-gene pairs with the most significant  $p$ -values. Genes with at least 5 eQTLs were selected for the colocalization analysis. We selected variants with a suggestive significance ( $p$ -value  $< 5 \times 10^{-7}$ ) for LDL-C or HDL-C effects within 500 kb of the lead locus variant. For the *GAS6* variant analysis, we curated all the GWAS variants within 500 kb of the lead variant with nominal significance ( $p$ -value  $< 0.05$ ) and matched them to eQTL data where the transcription starting site of the corresponding gene is within +/-500 kb. We conducted colocalization analysis using the coloc.abf() function<sup>98</sup> and identified nominally significant (PP.H4  $> 1 \times 10^{-3}$ ) genes-eQTL pairs. The coloc methodology implements an efficient statistical framework to identify shared variants from two association signals through posteriors probabilities. Finally, we used the colocalized signals and compared the significant genes using STRING<sup>96</sup>, a protein-protein interaction database. All the correlation tests were conducted in R, where we calculated Pearson correlations between the lipid effect estimates and gene expression effects (slope) from GTEx.

**Phenome wide association analysis.** The complex trait information was curated from UK Biobank resource, where we curated multiple disease phenotypes for UKB samples into International Classification of Diseases (ICD)-based phecodes based on phecode map (<https://phewascatalog.org>) using the PheWAS R package (version PheWAS\_0.99.5-4). We conducted a phenome-wide association analysis (PheWAS) using a logistic regression model glm() in R. We adjusted the models for PCI-10, age, age<sup>2</sup>, sex, and race.

### Calculation of heritability estimates from TOPMed WGS data

We calculated heritabilities estimated for the four lipids using TOPMed WGS data using Greml-LDMS approach<sup>39</sup>, where we binned the variants into four MAF bins based on minor allele frequency and grouped the variants to four LD quartiles based on LD score calculated by GCTA method<sup>99</sup>. The four MAF bins used in this study includes  $\geq 0.05$ ,  $\geq 0.01$  to  $< 0.05$ ,  $\geq 0.001$  to  $< 0.01$  and  $\geq 0.0001$  to  $< 0.001$ . We excluded any variant with MAF  $< 0.0001$  from this analysis. The hereditary estimation was calculated for three ancestral groups (African, European, Hispanic) where only unrelated samples (kinship score  $< 0.025$ ) were included in the analysis. We excluded the other two ancestral groups (i.e., Asian and Samoan) from this analysis due to insufficient sample sizes. In total we included 9640, 21568 and 10631 in African, European and Hispanic ancestries respectively. For each MAF bin, we implemented certain quality control (QC) measures using PLINK software<sup>20</sup>, which includes; genotype missingness (--geno 0.05), sample missingness (--mind 0.05), Hardy-Weinberg equilibrium (--hwe  $10^{-6}$ ) and LD pruned variants (--indep-pairwise 50 5 0.1) as implemented by Wainschein et al.<sup>39</sup>. Next, we implemented Greml-LDMS with LD score region as 200 and GRM cut-off as 0.05 for the four lipid phenotypes. We calculated 20 principal components from the QC passed variants in each MAF bin and implemented GCTA workflow with --reml-no-constrain, --reml-no-lrt and --reml-maxit 10,000 parameters to avoid the no-convergence issues and negative  $h^2$  estimates. For comparing the  $h^2$  estimates between variants from WGS data and array-genotypes, first, we used QC passed WGS variants as mentioned above, second, we curated the variants from MGB Biobank array data and intersected them with WGS variants from TOPMed. Next, we calculated heritability estimates for array-genotype variants and compared with  $h^2$  estimates from WGS variants for the three ancestral groups.

### Reporting summary

Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Data availability

Individual whole-genome sequence data for TOPMed and harmonized lipids at individual sample level are available through restricted access via the TOPMed dbGaP Exchange area. Summary level genotype data from TOPMed are available through the BRAVO browser (<https://bravo.sph.umich.edu/>). The UK Biobank (UKB) whole-genome sequence data can be accessed through UKB Research Analysis Platform (RAP), through the UKB approval system (<https://www.ukbiobank.ac.uk>). The Mass General Brigham Biobank (MGBB) individual-level data are available from <https://personalizedmedicine.partners.org/Biobank/Default.aspx>, where the data is available through institutional review board (IRB) approval, therefore not publicly available. Individual-level data from Penn Medicine BioBank is not publicly available due to research participants privacy concerns. The summary data captured using whole exome sequencing can be accessed through PMBB Genome Browser (<https://pmbb.med.upenn.edu/allele-frequency/>). The dbGaP accessions for TOPMed cohorts are as follows: Old Order Amish (Amish) *phs000956* and *phs00039*; Atherosclerosis Risk in Communities study (ARIC) *phs001211* and *phs000280*; Mt Sinai BioMe Biobank (BioMe) *phs001644* and *phs000925*; Coronary Artery Risk Development in Young Adults (CARDIA) *phs001612* and *phs000285*; Cleveland Family Study (CFS) *phs000954* and *phs000284*; Cardiovascular Health Study (CHS) *phs001368* and *phs000287*; Diabetes Heart Study (DHS) *phs001412* and *phs001012*; Framingham Heart Study (FHS) *phs000974* and *phs000007*; Genetic Studies of Atherosclerosis Risk (GeneSTAR) *phs001218* and *phs000375*; Genetic Epidemiology Network of Arteriopathy (GENOA) *phs001345* and *phs001238*; Genetic Epidemiology Network of Salt Sensitivity (GenSalt) *phs001217* and *phs000784*; Genetics of Lipid-Lowering Drugs and Diet Network (GOLDN) *phs001359* and *phs000741*; Hispanic Community Health Study - Study of Latinos (HCHS\_SOL) *phs001395* and *phs000810*; Hypertension Genetic Epidemiology Network and Genetic Epidemiology Network of Arteriopathy (HyperGEN) *phs001293* and *phs001293*; Jackson Heart Study (JHS) *phs000964* and *phs000286*; Multi-Ethnic Study of Atherosclerosis (MESA) *phs001416* and *phs000209*; Massachusetts General Hospital Atrial Fibrillation Study (MGH\_AF) *phs001062* and *phs001001*; San Antonio Family Study (SAFS) *phs001215* and *phs000462*; Samoan Adiposity Study (SAS) *phs000972* and *phs000914*; Taiwan Study of Hypertension using Rare Variants (THRV) *phs001387* and *phs001387*; Women's Health Initiative (WHI) *phs001237* and *phs000200*.

## Code availability

Codes used to implement STAAR workflows are available at <https://github.com/xihaoli/STAAR> and [https://github.com/xihaoli/STAAR\\_pipeline](https://github.com/xihaoli/STAAR_pipeline). Workflow implemented for whole genome heritability calculations are available at [https://github.com/CNSGenomics/Heritability\\_WGS](https://github.com/CNSGenomics/Heritability_WGS).

## References

- Cohen, J. C., Boerwinkle, E., Mosley, T. H. & Hobbs, H. H. Sequence variations in PCSK9, low LDL, and protection against coronary heart disease. *N. Engl. J. Med.* **354**, 1264–1272 (2006).
- Cohen, J. et al. Low LDL cholesterol in individuals of African descent resulting from frequent nonsense mutations in PCSK9. *Nat. Genet.* **37**, 161–165 (2005).
- Musunuru, K. et al. Exome sequencing, *ANGPTL3* mutations, and familial combined hypolipidemia. *N. Engl. J. Med.* **363**, 2220–2227 (2010).
- Stitzel, N. O. et al. *ANGPTL3* deficiency and protection against coronary artery disease. *J. Am. Coll. Cardiol.* **69**, 2054–2063 (2017).
- Dewey, F. E. et al. Genetic and pharmacologic inactivation of *ANGPTL3* and cardiovascular disease. *N. Engl. J. Med.* **377**, 211–221 (2017).
- Manolio, T. A. et al. Finding the missing heritability of complex diseases. *Nature* **461**, 747–753 (2009).
- Pollin, T. I. et al. A null mutation in human APOC3 confers a favorable plasma lipid profile and apparent cardioprotection. *Science* **322**, 1702–1705 (2008).
- Shen, H. et al. Familial defective apolipoprotein B-100 and increased low-density lipoprotein cholesterol and coronary artery calcification in the old order amish. *Arch. Intern. Med.* **170**, 1850–1855 (2010).
- Saleheen, D. et al. Human knockouts and phenotypic analysis in a cohort with a high rate of consanguinity. *Nature* **544**, 235–239 (2017).
- Exome Aggregation Consortium. et al. Analysis of protein-coding genetic variation in 60,706 humans. *Nature* **536**, 285–291 (2016).
- Samocha, K. E. et al. A framework for the interpretation of de novo mutation in human disease. *Nat. Genet.* **46**, 944–950 (2014).
- Natarajan, P. et al. Chromosome Xq23 is associated with lower atherogenic lipid concentrations and favorable cardiometabolic indices. *Nat. Commun.* **12**, 2182 (2021).
- Li, X. et al. Dynamic incorporation of multiple in silico functional annotations empowers rare variant association analysis of large whole-genome sequencing studies at scale. *Nat. Genet.* **52**, 969–983 (2020).
- Natarajan, P. et al. Deep-coverage whole genome sequences and blood lipids among 16,324 individuals. *Nat. Commun.* **9**, 3391 (2018).
- Klarin, D. et al. Genetics of blood lipids among ~300,000 multi-ethnic participants of the Million Veteran Program. *Nat. Genet.* **50**, 1514–1523 (2018).
- Hu, Y. et al. Minority-centric meta-analyses of blood lipid levels identify novel loci in the Population Architecture using Genomics and Epidemiology (PAGE) study. *PLoS Genet.* **16**, e1008684 (2020).
- NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium. et al. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature* **590**, 290–299 (2021).
- Stilp, A. M. et al. A System for Phenotype Harmonization in the NHLBI Trans-Omics for Precision Medicine (TOPMed) Program. *Am. J. Epidemiol.* <https://doi.org/10.1093/aje/kwab115> (2021).
- Fadista, J., Manning, A. K., Florez, J. C. & Groop, L. The (in)famous GWAS P-value threshold revisited and updated for low-frequency variants. *Eur. J. Hum. Genet.* **24**, 1202–1205 (2016).
- Purcell, S. et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- Bentley, A. R. et al. Multi-ancestry genome-wide gene-smoking interaction study of 387,272 individuals identifies new loci associated with serum lipids. *Nat. Genet.* **51**, 636–648 (2019).
- Ripatti, P. et al. Polygenic hyperlipidemias and coronary artery disease risk. *Circ. Genom. Precis. Med.* **13**, e002725 (2020).
- van Leeuwen, E. M. et al. Meta-analysis of 49 549 individuals imputed with the 1000 Genomes Project reveals an exonic damaging variant in *ANGPTL4* determining fasting TG levels. *J. Med. Genet.* **53**, 441–449 (2016).
- Nielsen, J. B. et al. Loss-of-function genomic variants highlight potential therapeutic targets for cardiovascular disease. *Nat. Commun.* **11**, 6417 (2020).
- Aragam, K. G. et al. Limitations of contemporary guidelines for managing patients at high genetic risk of coronary artery disease. *J. Am. Coll. Cardiol.* **75**, 2769–2780 (2020).
- Park, J. et al. Exome-wide evaluation of rare coding variants using electronic health records identifies new gene-phenotype associations. *Nat. Med.* **27**, 66–72 (2021).
- Mountjoy, E. et al. An open approach to systematically prioritize causal variants and genes at all published human GWAS trait-associated loci. *Nat. Genet.* **53**, 1527–1533 (2021).

28. Lonsdale, J. et al. The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.* **45**, 580–585 (2013).
29. Barter, P. J. et al. Effects of torcetrapib in patients at high risk for coronary events. *N. Engl. J. Med.* **357**, 2109–2122 (2007).
30. Schwartz, G. G. et al. Effects of dalcetrapib in patients with a recent acute coronary syndrome. *N. Engl. J. Med.* **367**, 2089–2099 (2012).
31. The HPS3/TIMI55–REVEAL Collaborative Group. Effects of anacetrapib in patients with atherosclerotic vascular disease. *N. Engl. J. Med.* **377**, 1217–1227 (2017).
32. Lincoff, A. M. et al. Evacetrapib and cardiovascular outcomes in high-risk vascular disease. *N. Engl. J. Med.* **376**, 1933–1942 (2017).
33. Fairroozy, R. H., White, J., Palmén, J., Kalea, A. Z. & Humphries, S. E. Identification of the functional variant(s) that explain the low-density lipoprotein receptor (LDLR) GWAS SNP rs6511720 association with lower LDL-C and risk of CHD. *PLoS ONE* **11**, e0167676 (2016).
34. Li, Z. et al. Dynamic scan procedure for detecting rare-variant association regions in whole-genome sequencing studies. *Am. J. Hum. Genet.* **104**, 802–814 (2019).
35. Roses, A. D. et al. A TOMM40 variable-length polymorphism predicts the age of late-onset Alzheimer’s disease. *Pharmacogenomics J.* **10**, 375–384 (2010).
36. Li, G. et al. TOMM40 intron 6 poly-T length, age at onset, and neuropathology of AD in individuals with APOE  $\epsilon$ 3/ $\epsilon$ 3. *Alzheimers Dement. J. Alzheimers Assoc.* **9**, 554–561 (2013).
37. Glazier, A. M., Scott, J. & Aitman, T. J. Molecular basis of the Cd36 chromosomal deletion underlying SHR defects in insulin action and fatty acid metabolism. *Mamm. Genome. J. Int. Mamm. Genome Soc.* **13**, 108–113 (2002).
38. The LifeLines Cohort Study. et al. Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. *Nat. Genet.* **47**, 1114–1120 (2015).
39. Wainschein, P. et al. Assessing the contribution of rare variants to complex trait heritability from whole-genome sequence data. *Nat. Genet.* **54**, 263–273 (2022).
40. Cadby, G. et al. Heritability of 596 lipid species and genetic correlation with cardiovascular traits in the Busselton Family Heart Study. *J. Lipid Res.* **61**, 537–545 (2020).
41. Global Lipids Genetics Consortium. Discovery and refinement of loci associated with lipid levels. *Nat. Genet.* **45**, 1274–1283 (2013).
42. Willer, C. J. et al. Discovery and refinement of loci associated with lipid levels. *Nat. Genet.* **45**, 1274–1283 (2013).
43. ENGAGE Consortium. et al. The impact of low-frequency and rare variants on lipid levels. *Nat. Genet.* **47**, 589–597 (2015).
44. The Myocardial Infarction Genetics Consortium Investigators. Inactivating mutations in *NPC1L1* and protection from coronary heart disease. *N. Engl. J. Med.* **371**, 2072–2082 (2014).
45. GLGC Consortium. et al. Exome chip meta-analysis identifies novel loci and East Asian-specific coding variants that contribute to lipid levels and coronary artery disease. *Nat. Genet.* **49**, 1722–1730 (2017).
46. Hoffmann, T. J. et al. A large electronic-health-record-based genome-wide study of serum lipids. *Nat. Genet.* **50**, 401–413 (2018).
47. Martin, A. R. et al. Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat. Genet.* **51**, 584–591 (2019).
48. Peloso, G. M. & Natarajan, P. Insights from population-based analyses of plasma lipids across the allele frequency spectrum. *Curr. Opin. Genet. Dev.* **50**, 1–6 (2018).
49. Kremer, L. S. et al. Genetic diagnosis of Mendelian disorders via RNA sequencing. *Nat. Commun.* **8**, 15824 (2017).
50. Cummings, B. B. et al. Improving genetic diagnosis in Mendelian disease with transcriptome sequencing. *Sci. Transl. Med.* **9**, eaal5209 (2017).
51. Genome Aggregation Database Production Team. et al. Transcript expression-aware annotation improves rare variant interpretation. *Nature* **581**, 452–458 (2020).
52. Mendes de Almeida, R. et al. Whole gene sequencing identifies deep-intronic variants with potential functional impact in patients with hypertrophic cardiomyopathy. *PLoS ONE* **12**, e0182946 (2017).
53. Vitsios, D., Dhindsa, R. S., Middleton, L., Gussow, A. B. & Petrovski, S. Prioritizing non-coding regions based on human genomic constraint and sequence context with deep learning. *Nat. Commun.* **12**, 1504 (2021).
54. di Iulio, J. et al. The human noncoding genome defined by genetic diversity. *Nat. Genet.* **50**, 333–337 (2018).
55. Genome Aggregation Database Consortium. et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
56. Khera, A. V. et al. Diagnostic yield and clinical utility of sequencing familial hypercholesterolemia genes in patients with severe hypercholesterolemia. *J. Am. Coll. Cardiol.* **67**, 2578–2589 (2016).
57. Benn, M., Watts, G. F., Tybjaerg-Hansen, A. & Nordestgaard, B. G. Mutations causative of familial hypercholesterolaemia: screening of 98 098 individuals from the Copenhagen General Population Study estimated a prevalence of 1 in 217. *Eur. Heart J.* **37**, 1384–1394 (2016).
58. Grundy, S. M. et al. 2018 AHA/ACC/AACVPR/AAPA/ABC/ACPM/ADA/AGS/APhA/ASPC/NLA/PCNA Guideline on the Management of Blood Cholesterol: Executive Summary. *J. Am. Coll. Cardiol.* **73**, 3168–3209 (2019).
59. Sturm, A. C. et al. Clinical genetic testing for familial hypercholesterolemia. *J. Am. Coll. Cardiol.* **72**, 662–680 (2018).
60. Reeskamp, L. F. et al. A Deep intronic variant in *LDLR* in familial hypercholesterolemia: time to widen the scope? *Circ. Genomic Precis. Med.* **11**, e002385 (2018).
61. Calandra, S., Tarugi, P. & Bertolini, S. Altered mRNA splicing in lipoprotein disorders. *Curr. Opin. Lipidol.* **22**, 93–99 (2011).
62. on behalf of the ACMG Laboratory Quality Assurance Committee. et al. Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet. Med.* **17**, 405–423 (2015).
63. Peloso, G. M. et al. Rare protein-truncating variants in APOB, lower low-density lipoprotein cholesterol, and protection against coronary heart disease. *Circ. Genom. Precis. Med.* **12**, e002376 (2019).
64. Abifadel, M. et al. Mutations in PCSK9 cause autosomal dominant hypercholesterolemia. *Nat. Genet.* **34**, 154–156 (2003).
65. Jiang, L. et al. The distribution and characteristics of LDL receptor mutations in China: a systematic review. *Sci. Rep.* **5**, 17272 (2015).
66. Arráiz, N. et al. Novel mutations identification in exon 4 of LDLR gene in patients with moderate hypercholesterolemia in a Venezuelan population. *Am. J. Ther.* **17**, 325–329 (2010).
67. Gudnason, V. et al. Identification of recurrent and novel mutations in exon 4 of the LDL receptor gene in patients with familial hypercholesterolemia in the United Kingdom. *Arterioscler. Thromb. J. Vasc. Biol.* **13**, 56–63 (1993).
68. Goldmann, R. et al. Genomic characterization of large rearrangements of the LDLR gene in Czech patients with familial hypercholesterolemia. *BMC Med. Genet.* **11**, 115 (2010).
69. Zuk, O. et al. Searching for missing heritability: Designing rare variant association studies. *Proc. Natl Acad. Sci. USA* **111**, E455–E464 (2014).
70. Bycroft, C. et al. The UK Biobank resource with deep phenotyping and genomic data. *Nature* **562**, 203–209 (2018).
71. Soria, L. F. et al. Association between a specific apolipoprotein B mutation and familial defective apolipoprotein B-100. *Proc. Natl Acad. Sci. USA* **86**, 587–591 (1989).
72. Conomos, M. P., Reiner, A. P., Weir, B. S. & Thornton, T. A. Model-free estimation of recent genetic relatedness. *Am. J. Hum. Genet.* **98**, 127–148 (2016).

73. Gogarten, S. M. et al. Genetic association testing using the GENESIS R/Bioconductor package. *Bioinforma. Oxf. Engl.* **35**, 5346–5348 (2019).
74. Das, S. et al. Next-generation genotype imputation service and methods. *Nat. Genet.* **48**, 1284–1287 (2016).
75. Loh, P.-R. et al. Reference-based phasing using the Haplotype Reference Consortium panel. *Nat. Genet.* **48**, 1443–1448 (2016).
76. Fuchsberger, C., Abecasis, G. R. & Hinds, D. A. minimac2: faster genotype imputation. *Bioinformatics* **31**, 782–784 (2015).
77. UK10K Consortium. et al. Improved imputation of low-frequency and rare variants using the UK10K haplotype reference panel. *Nat. Commun.* **6**, 8111 (2015).
78. the Haplotype Reference Consortium. A reference panel of 64,976 haplotypes for genotype imputation. *Nat. Genet.* **48**, 1279–1283 (2016).
79. Halldorsson, B. V. et al. The sequences of 150,119 genomes in the UK biobank. *Nature* **607**, 732–740 (2022).
80. Zhou, W. et al. Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nat. Genet.* **50**, 1335–1341 (2018).
81. Pulit, S. L., de With, S. A. J. & de Bakker, P. I. W. Resetting the bar: statistical significance in whole-genome sequencing-based association studies of global populations. *Genet. Epidemiol.* **41**, 145–151 (2017).
82. Han, B. & Eskin, E. Random-effects model aimed at discovering associations in meta-analysis of genome-wide association studies. *Am. J. Hum. Genet.* **88**, 586–598 (2011).
83. Li, Z. et al. A framework for detecting noncoding rare variant associations of large-scale whole-genome sequencing studies. <https://doi.org/10.1101/2021.11.05.467531> (2021).
84. Chen, H. et al. Control for population structure and relatedness for binary traits in genetic association studies via logistic mixed models. *Am. J. Hum. Genet.* **98**, 653–666 (2016).
85. Chen, H. et al. Efficient variant set mixed model association tests for continuous and binary traits in large-scale whole-genome sequencing studies. *Am. J. Hum. Genet.* **104**, 260–274 (2019).
86. The FANTOM Consortium and the RIKEN PMI and CLST (DGT). A promoter-level mammalian expression atlas. *Nature* **507**, 462–470 (2014).
87. The ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
88. The FANTOM Consortium. et al. An atlas of active enhancers across human cell types and tissues. *Nature* **507**, 455–461 (2014).
89. Fishilevich, S. et al. GeneHancer: genome-wide integration of enhancers and target genes in GeneCards. *Database J. Biol. Databases Curation* **2017**, (2017).
90. Wu, M. C. et al. Rare-variant association testing for sequencing data with the sequence kernel association test. *Am. J. Hum. Genet.* **89**, 82–93 (2011).
91. Madsen, B. E. & Browning, S. R. A groupwise association test for rare mutations using a weighted sum statistic. *PLoS Genet.* **5**, e1000384 (2009).
92. Liu, Y. et al. ACAT: a fast and powerful p value combination method for rare-variant analysis in sequencing studies. *Am. J. Hum. Genet.* **104**, 410–421 (2019).
93. Buniello, A. et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–D1012 (2019).
94. Watanabe, K., Taskesen, E., van Bochoven, A. & Posthuma, D. Functional mapping and annotation of genetic associations with FUMA. *Nat. Commun.* **8**, 1826 (2017).
95. de Leeuw, C. A., Mooij, J. M., Heskes, T. & Posthuma, D. MAGMA: Generalized Gene-Set Analysis of GWAS Data. *PLoS Comput. Biol.* **11**, e1004219 (2015).
96. Szklarczyk, D. et al. STRING v11: protein–protein association networks with increased coverage, supporting functional discovery in genome-wide experimental datasets. *Nucleic Acids Res.* **47**, D607–D613 (2019).
97. GTEx Consortium. Genetic effects on gene expression across human tissues. *Nature* **550**, 204–213 (2017).
98. Giambartolomei, C. et al. Bayesian test for colocalisation between pairs of genetic association studies using summary statistics. *PLoS Genet.* **10**, e1004383 (2014).
99. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).

## Acknowledgements

Whole genome sequencing (WGS) for the Trans-Omics in Precision Medicine (TOPMed) program was supported by the National Heart, Lung and Blood Institute (NHLBI). P.N. is supported by grants from the National Heart, Lung, and Blood Institute (R01HL142711, R01HL148050, R01HL151283, R01HL148565, R01HL135242, R01HL151152), Fondation Leducq (TNE-18CVDO4), and Massachusetts General Hospital (Paul and Phyllis Fireman Endowed Chair in Vascular Medicine). G.M.P. is supported by NIH grants R01HL142711 and R01HL127564. X.Lin is supported by grants R35-CA197449, U19-CA203654, R01-HL113338, and U01-HG009088. Prior to his employment at Novartis and during this work S.A.L. was supported by NIH grants R01HL139731, R01HL157635, and American Heart Association 18SFRN34250007. We like to acknowledge all the grants that supported this study, R01 HL121007, U01 HL072515, R01 AG18728, X01HL134588, HL 046389, HL113338, and 1R35HL135818, K01 HL135405, R03 HL154284, U01HL072507, R01HL087263, R01HL090682, P01HL045522, R01MH078143, R01MH078111, R01MH083824, U01DK085524, R01HL113323, R01HL093093, R01HL140570, R01HL142711, R01HL127564, R01HL148050, R01HL148565, HL105756, and Leducq TNE-18CVDO4. The views expressed in this manuscript are those of the authors and do not necessarily represent the views of the National Heart, Lung, and Blood Institute; the National Institutes of Health; or the U.S. Department of Health and Human Services. Detailed acknowledgements provided in Supplementary Note 2.

## Author contributions

M.S.S., G.M.P., and P.N. designed the study. M.S.S. carried out all the primary analysis with critical inputs from G.M.P. and P.N. M.S.S., Xih.L, Z.L., A.P., D.Y.Z., J.P., S.A., J.C.B., J.A.B., B.E.C., L.M.C., R.H.C., J.E.C., L.F., P.S.V., R.D., B.I.F., M.G., X.G., N.H.C., B.H., C.M.H., M.R.I., T.N.K., B.G.K., L.L., Xia.L, M.L., S.A.L., A.W.M., P.M., M.E.M., A.C.M., T.N., J.R.O.C., N.D.P., P.A.P., M.S.R., J.A.S., X.S., K.D.T., R.P.T., M.Y.T., Z.W., Y.W., B.W., J.T.W., L.R.Y., W.Z., D.K.A., J. Blanger, E.B., D.W.B., Y.I.C., A.C., L.A.C., S.K.D., P.T.E., M.F., S. Gabriel, S. Germer, R.G., J.H., R.C.K., S.L.R.K., R. Kim, C.K., R.J.F.L., K.M., R.A.M., S.T.M., B.D.M., D.N., K.E.N., B.M.P., S. Redline, A.P.R., R.S.V., S.S.R., C.W., J.I.R., D.J.R., X.Lin., G.M.P., and P.N. acquired, analyzed or interpreted data. M.S.S., G.M.P. and P.N. wrote the first draft of the manuscript and all others provided intellectual revisions. G.M.P. and P.N. and NHLBI TOPMed Lipids Working Group provided administrative, technical, or material support.

## Competing interests

P.N. reports investigator-initiated grant support from Amgen, Apple, AstraZeneca, and Boston Scientific, personal fees from Apple, AstraZeneca, Blackstone Life Sciences, Foresite Labs, Genentech, TenSixteen Bio, and Novartis, scientific advisory board membership of geneXwell and TenSixteen Bio, and spousal employment at Vertex, all unrelated to the present work. B.P. serves on the Steering Committee of the Yale Open Data Access Project funded by Johnson & Johnson.

M.E.M. receives funding from Regeneron Pharmaceutical Inc. unrelated to this work. S.A. has employment and equity in 23andMe, Inc. The spouse of C.J.W. works at Regeneron. S.A.L. is a full-time employee of Novartis as of July 18, 2022. S.A.L. has received sponsored research support from Bristol Myers Squibb, Pfizer, Boehringer Ingelheim, Fitbit, Medtronic, Premier, and IBM, and has consulted for Bristol Myers Squibb, Pfizer, Blackstone Life Sciences, and Invitae. X. Lin is a consultant of AbbVie Pharmaceuticals and Verily Life Sciences. The remaining authors declare no competing interests.

## Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41467-022-33510-7>.

**Correspondence** and requests for materials should be addressed to Gina M. Peloso or Pradeep Natarajan.

**Peer review information** *Nature Communications* thanks David Meyre and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. Peer reviewer reports are available.

**Reprints and permission information** is available at <http://www.nature.com/reprints>

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

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

© The Author(s) 2022

Margaret Sunitha Selvaraj<sup>1,2,3</sup>, Xihao Li<sup>4</sup>, Zilin Li<sup>4</sup>, Akhil Pampana<sup>2</sup>, David Y. Zhang<sup>5,6</sup>, Joseph Park<sup>5,6</sup>, Stella Aslibekyan<sup>7</sup>, Joshua C. Bis<sup>8</sup>, Jennifer A. Brody<sup>8</sup>, Brian E. Cade<sup>9</sup>, Lee-Ming Chuang<sup>10</sup>, Ren-Hua Chung<sup>11</sup>, Joanne E. Curran<sup>12</sup>, Lisa de las Fuentes<sup>13,14</sup>, Paul S. de Vries<sup>15</sup>, Ravindranath Duggirala<sup>12</sup>, Barry I. Freedman<sup>16</sup>, Mariaelisa Graff<sup>17</sup>, Xiuqing Guo<sup>18</sup>, Nancy Heard-Costa<sup>19</sup>, Bertha Hidalgo<sup>7</sup>, Chii-Min Hwu<sup>20</sup>, Marguerite R. Irvin<sup>7</sup>, Tanika N. Kelly<sup>21,22</sup>, Brian G. Kral<sup>23</sup>, Leslie Lange<sup>24</sup>, Xiaohui Li<sup>18</sup>, Martin Lisa<sup>25</sup>, Steven A. Lubitz<sup>1,26</sup>, Ani W. Manichaikul<sup>27</sup>, Preuss Michael<sup>28</sup>, May E. Montasser<sup>29</sup>, Alanna C. Morrison<sup>15</sup>, Take Naseri<sup>30</sup>, Jeffrey R. O'Connell<sup>29</sup>, Nicholette D. Palmer<sup>31</sup>, Patricia A. Peyser<sup>32</sup>, Muagututia S. Reupena<sup>33</sup>, Jennifer A. Smith<sup>32</sup>, Xiao Sun<sup>21</sup>, Kent D. Taylor<sup>18</sup>, Russell P. Tracy<sup>34</sup>, Michael Y. Tsai<sup>35</sup>, Zhe Wang<sup>28</sup>, Yuxuan Wang<sup>36</sup>, Wei Bao<sup>37</sup>, John T. Wilkins<sup>38</sup>, Lisa R. Yanek<sup>23</sup>, Wei Zhao<sup>32</sup>, Donna K. Arnett<sup>39</sup>, John Blangero<sup>12</sup>, Eric Boerwinkle<sup>15</sup>, Donald W. Bowden<sup>31</sup>, Yii-Der Ida Chen<sup>40</sup>, Adolfo Correa<sup>41</sup>, L. Adrienne Cupples<sup>36</sup>, Susan K. Dutcher<sup>42</sup>, Patrick T. Ellinor<sup>1,26</sup>, Myriam Fornage<sup>43</sup>, Stacey Gabriel<sup>44</sup>, Soren Germer<sup>45</sup>, Richard Gibbs<sup>46</sup>, Jiang He<sup>21,22</sup>, Robert C. Kaplan<sup>47,48</sup>, Sharon L. R. Kardia<sup>32</sup>, Ryan Kim<sup>49</sup>, Charles Kooperberg<sup>48</sup>, Ruth J. F. Loos<sup>28,50</sup>, Karine A. Viaud-Martinez<sup>51</sup>, Rasika A. Mathias<sup>23</sup>, Stephen T. McGarvey<sup>52</sup>, Braxton D. Mitchell<sup>29,53</sup>, Deborah Nickerson<sup>54</sup>, Kari E. North<sup>17</sup>, Bruce M. Psaty<sup>8,55,56</sup>, Susan Redline<sup>9</sup>, Alexander P. Reiner<sup>55,48</sup>, Ramachandran S. Vasan<sup>57,58,59</sup>, Stephen S. Rich<sup>27</sup>, Cristen Willer<sup>60</sup>, Jerome I. Rotter<sup>18</sup>, Daniel J. Rader<sup>5,6,61</sup>, Xihong Lin<sup>2,4,62</sup>, NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium\*, Gina M. Peloso<sup>36,220</sup> ✉ & Pradeep Natarajan<sup>1,2,3,220</sup> ✉

<sup>1</sup>Cardiovascular Research Center, Massachusetts General Hospital, Boston, MA 02114, USA. <sup>2</sup>Program in Medical and Population Genetics, Broad Institute of Harvard and MIT, Cambridge, MA 02142, USA. <sup>3</sup>Department of Medicine, Harvard Medical School, Boston, MA 02115, USA. <sup>4</sup>Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA 02115, USA. <sup>5</sup>Department of Genetics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA 19104, USA. <sup>6</sup>Department of Medicine, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA 19104, USA. <sup>7</sup>Department of Epidemiology, University of Alabama at Birmingham School of Public Health, Birmingham, AL, USA. <sup>8</sup>Cardiovascular Health Research Unit, Department of Medicine, University of Washington, Seattle, WA, USA. <sup>9</sup>Department of Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA. <sup>10</sup>Department of Internal Medicine, National Taiwan University Hospital, Taipei, Taiwan. <sup>11</sup>Institute of Population Health Sciences, National Health Research Institutes, Zhunan 350, Taiwan. <sup>12</sup>Department of Human Genetics and South Texas Diabetes and Obesity Institute, University of Texas Rio Grande Valley School of Medicine, Brownsville, TX 78520, USA. <sup>13</sup>Department of Medicine, Cardiovascular Division, Washington University School of Medicine, St. Louis, MO, USA. <sup>14</sup>Division of Biostatistics, Washington University School of Medicine, St. Louis, MO, USA. <sup>15</sup>Human Genetics Center, Department of Epidemiology, Human Genetics, and Environmental Sciences, School of Public Health, The University of Texas Health Science Center at Houston, Houston, TX, USA. <sup>16</sup>Department of Internal Medicine, Section on Nephrology, Wake Forest School of Medicine, Winston-Salem, NC 27157, USA. <sup>17</sup>Department of Epidemiology, UNC Chapel Hill, Chapel Hill, NC, USA. <sup>18</sup>The Institute for Translational Genomics and Population Sciences, Department of Pediatrics, The Lundquist Institute for Biomedical Innovation at Harbor-UCLA Medical Center, Torrance, CA, USA. <sup>19</sup>Department of Neurology, Boston University School of Medicine, Boston, MA, USA. <sup>20</sup>Section of Endocrinology and Metabolism, Department of Medicine, Taipei Veterans General Hospital, Taipei, Taiwan. <sup>21</sup>Department of Epidemiology, Tulane University School of Public Health and Tropical Medicine, New Orleans, LA 70112, USA. <sup>22</sup>Tulane University Translational Science Institute, New Orleans, LA 70112, USA. <sup>23</sup>Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, MD 21205, USA. <sup>24</sup>Division of Biomedical



Informatics and Personalized Medicine, Department of Medicine, University of Colorado Anschutz Medical Campus, Aurora, CO, USA. <sup>25</sup>Department of Medicine, George Washington University, Washington, DC, USA. <sup>26</sup>Cardiovascular Disease Initiative, The Broad Institute of MIT and Harvard, Cambridge, MA 02124, USA. <sup>27</sup>Department of Public Health Sciences, Center for Public Health Genomics, University of Virginia, Charlottesville, VA, USA. <sup>28</sup>The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA. <sup>29</sup>Department of Medicine, University of Maryland School of Medicine, Baltimore, MD, USA. <sup>30</sup>Ministry of Health, Government of Samoa, Samoa, USA. <sup>31</sup>Department of Biochemistry, Wake Forest School of Medicine, Winston-Salem, NC 27157, USA. <sup>32</sup>Department of Epidemiology, University of Michigan, Ann Arbor, MI 48109, USA. <sup>33</sup>Lutia i Puava ae Mapu i Fagalele, Apia, Samoa. <sup>34</sup>Departments of Pathology & Laboratory Medicine and Biochemistry, Larner College of Medicine at the University of Vermont, Colchester, VT, USA. <sup>35</sup>Department of Laboratory Medicine and Pathology, University of Minnesota, Minneapolis, MN, USA. <sup>36</sup>Department of Biostatistics, Boston University School of Public Health, Boston, MA 02118, USA. <sup>37</sup>Institute of Public Health, Division of Life Sciences and Medicine, University of Science and Technology of China, Hefei, Anhui 230026, China. <sup>38</sup>Department of Medicine (Cardiology) and Department of Preventive Medicine, Northwestern University Feinberg School of Medicine, Chicago, IL, USA. <sup>39</sup>Dean's Office, University of Kentucky College of Public Health, Lexington, KY, USA. <sup>40</sup>Lundquist Institute for Biomedical Innovation at Harbor-UCLA Medical Center, Torrance, CA, USA. <sup>41</sup>Department of Population Health Science, University of Mississippi Medical Center, Jackson, MS, USA. <sup>42</sup>The McDonnell Genome Institute, Washington University School of Medicine, St. Louis, MO 63108, USA. <sup>43</sup>Brown Foundation Institute of Molecular Medicine, McGovern Medical School, The University of Texas Health Science Center at Houston, Houston, TX 7722, USA. <sup>44</sup>Broad Institute, Cambridge, MA 02142, USA. <sup>45</sup>New York Genome Center, New York, NY 10013, USA. <sup>46</sup>Baylor College of Medicine Human Genome Sequencing Center, Houston, TX 77030, USA. <sup>47</sup>Department of Epidemiology and Population Health, Albert Einstein College of Medicine, Bronx, NY 10461, USA. <sup>48</sup>Division of Public Health Sciences, Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA. <sup>49</sup>Psomagen, Inc. (formerly MacroGen USA), Rockville, MD, USA. <sup>50</sup>NNF Center for Basic Metabolic Research, University of Copenhagen, Copenhagen, Denmark. <sup>51</sup>Illumina Laboratory Services, Illumina, Inc, San Diego 92122, USA. <sup>52</sup>Department of Epidemiology, International Health Institute, Brown University, Providence, RI, USA. <sup>53</sup>Geriatrics Research and Education Clinical Center, Baltimore Veterans Administration Medical Center, Baltimore, MD, USA. <sup>54</sup>University of Washington, Department of Genome Sciences, Seattle, WA 98195, USA. <sup>55</sup>Department of Epidemiology, University of Washington, Seattle, WA, USA. <sup>56</sup>Department of Health Systems and Population Health, University of Washington, Seattle, WA, USA. <sup>57</sup>Sections of Preventive medicine and Epidemiology, Cardiovascular medicine, Department of Medicine, Boston University School of Medicine, Boston, MA, USA. <sup>58</sup>Department of Epidemiology, Boston University School of Public Health, Boston, MA, USA. <sup>59</sup>Framingham Heart Study, Framingham, MA, USA. <sup>60</sup>University of Michigan, Internal Medicine, Ann Arbor, MI 48109, USA. <sup>61</sup>Institute for Translational Medicine and Therapeutics, Perelman School of Medicine, University of Pennsylvania, Philadelphia, PA 19104, USA. <sup>62</sup>Department of Statistics, Harvard University, Cambridge, MA 02138, USA. <sup>220</sup>These authors jointly supervised this work: Gina M. Peloso, Pradeep Natarajan. \*A list of authors and their affiliations appears at the end of the paper. ✉ e-mail: [gpeloso@bu.edu](mailto:gpeloso@bu.edu); [pnatarajan@mgh.harvard.edu](mailto:pnatarajan@mgh.harvard.edu)

## NHLBI Trans-Omics for Precision Medicine (TOPMed) Consortium

**Namiko Abe<sup>63</sup>, Gonçalo Abecasis<sup>64</sup>, Francois Aguet<sup>65</sup>, Christine Albert<sup>66</sup>, Laura Almasy<sup>67</sup>, Alvaro Alonso<sup>68</sup>, Seth Ament<sup>69</sup>, Peter Anderson<sup>70</sup>, Pramod Anugu<sup>71</sup>, Deborah Applebaum-Bowden<sup>72</sup>, Kristin Ardlie<sup>65</sup>, Dan Arking<sup>73</sup>, Allison Ashley-Koch<sup>74</sup>, Tim Assimes<sup>75</sup>, Paul Auer<sup>76</sup>, Dimitrios Avramopoulos<sup>73</sup>, Najib Ayas<sup>77</sup>, Adithya Balasubramanian<sup>78</sup>, John Barnard<sup>79</sup>, Kathleen Barnes<sup>80</sup>, R. Graham Barr<sup>81</sup>, Emily Barron-Casella<sup>73</sup>, Lucas Barwick<sup>82</sup>, Terri Beaty<sup>73</sup>, Gerald Beck<sup>83</sup>, Diane Becker<sup>84</sup>, Lewis Becker<sup>73</sup>, Rebecca Beer<sup>85</sup>, Amber Beitelshes<sup>69</sup>, Emelia Benjamin<sup>86</sup>, Takis Benos<sup>87</sup>, Marcos Bezerra<sup>88</sup>, Larry Bielak<sup>64</sup>, Thomas Blackwell<sup>64</sup>, Russell Bowler<sup>89</sup>, Ulrich Broeckel<sup>90</sup>, Jai Broome<sup>70</sup>, Deborah Brown<sup>91</sup>, Karen Bunting<sup>63</sup>, Esteban Burchard<sup>92</sup>, Carlos Bustamante<sup>93</sup>, Erin Buth<sup>94</sup>, Jonathan Cardwell<sup>95</sup>, Vincent Carey<sup>96</sup>, Julie Carrier<sup>97</sup>, Cara Carty<sup>98</sup>, Richard Casaburi<sup>99</sup>, Juan P. Casas Romero<sup>100</sup>, James Casella<sup>73</sup>, Peter Castaldi<sup>101</sup>, Mark Chaffin<sup>65</sup>, Christy Chang<sup>69</sup>, Yi-Cheng Chang<sup>102</sup>, Daniel Chasman<sup>103</sup>, Sameer Chavan<sup>95</sup>, Bo-Juen Chen<sup>63</sup>, Wei-Min Chen<sup>104</sup>, Yii-Der Ida Chen<sup>105</sup>, Michael Cho<sup>96</sup>, Seung Hoan Choi<sup>65</sup>, Mina Chung<sup>106</sup>, Clary Clish<sup>107</sup>, Suzy Comhair<sup>108</sup>, Matthew Conomos<sup>94</sup>, Elaine Cornell<sup>109</sup>, Carolyn Crandall<sup>99</sup>, James Crapo<sup>110</sup>, L. Adrienne Cupples<sup>111</sup>, Jeffrey Curtis<sup>64</sup>, Brian Custer<sup>112</sup>, Coleen Damcott<sup>69</sup>, Dawood Darbar<sup>113</sup>, Sean David<sup>114</sup>, Colleen Davis<sup>70</sup>, Michelle Daya<sup>95</sup>, Mariza de Andrade<sup>115</sup>, Michael DeBaun<sup>116</sup>, Ranjan Deka<sup>117</sup>, Dawn DeMeo<sup>96</sup>, Scott Devine<sup>69</sup>, Huyen Dinh<sup>78</sup>, Harsha Doddapaneni<sup>78</sup>, Qing Duan<sup>118</sup>, Shannon Dugan-Perez<sup>78</sup>, Ravi Duggirala<sup>119</sup>, Jon Peter Durda<sup>109</sup>, Charles Eaton<sup>120</sup>, Lynette Ekunwe<sup>71</sup>, Adel El Boueiz<sup>121</sup>, Leslie Emery<sup>70</sup>, Serpil Erzurum<sup>79</sup>, Charles Farber<sup>104</sup>, Jesse Farek<sup>78</sup>, Tasha Fingerlin<sup>122</sup>, Matthew Flickinger<sup>64</sup>, Nora Franceschini<sup>123</sup>, Chris Frazar<sup>70</sup>, Mao Fu<sup>69</sup>, Stephanie M. Fullerton<sup>70</sup>, Lucinda Fulton<sup>124</sup>, Weiniu Gan<sup>85</sup>, Shanshan Gao<sup>95</sup>, Yan Gao<sup>71</sup>, Margery Gass<sup>125</sup>, Heather Geiger<sup>126</sup>, Bruce Gelb<sup>127</sup>, Mark Geraci<sup>128</sup>, Robert Gerszten<sup>129</sup>, Auyon Ghosh<sup>96</sup>, Chris Gignoux<sup>75</sup>, Mark Gladwin<sup>87</sup>, David Glahn<sup>130</sup>, Stephanie Gogarten<sup>70</sup>, Da-Wei Gong<sup>69</sup>, Harald Goring<sup>131</sup>, Sharon Graw<sup>80</sup>, Kathryn J. Gray<sup>132</sup>, Daniel Grine<sup>95</sup>, Colin Gross<sup>64</sup>, C. Charles Gu<sup>124</sup>, Yue Guan<sup>69</sup>, Namrata Gupta<sup>65</sup>, David M. Haas<sup>133</sup>, Jeff Haessler<sup>125</sup>, Michael Hall<sup>134</sup>, Yi Han<sup>78</sup>, Patrick Hanly<sup>135</sup>, Daniel Harris<sup>136</sup>, Nicola L. Hawley<sup>137</sup>, Ben Heavner<sup>94</sup>, Susan Heckbert<sup>138</sup>, Ryan Hernandez<sup>92</sup>, David Herrington<sup>139</sup>, Craig Hersh<sup>140</sup>, Bertha Hidalgo<sup>141</sup>, James Hixson<sup>142</sup>, Brian Hobbs<sup>96</sup>, John Hokanson<sup>95</sup>, Elliott Hong<sup>69</sup>, Karin Hoth<sup>143</sup>, Chao Agnes Hsiung<sup>144</sup>, Jianhong Hu<sup>78</sup>, Yi-Jen Hung<sup>145</sup>, Haley Huston<sup>146</sup>, Chii Min Hwu<sup>147</sup>, Rebecca Jackson<sup>148</sup>, Deepti Jain<sup>70</sup>, Cashell Jaquish<sup>85</sup>, Jill Johnsen<sup>149</sup>, Andrew Johnson<sup>85</sup>, Craig Johnson<sup>70</sup>, Rich Johnston<sup>68</sup>, Kimberly Jones<sup>73</sup>, Hyun Min Kang<sup>150</sup>, Shannon Kelly<sup>151</sup>, Eimear Kenny<sup>127</sup>, Michael Kessler<sup>69</sup>, Alyna Khan<sup>70</sup>, Ziad Khan<sup>78</sup>, Wonji Kim<sup>152</sup>, John Kimoff<sup>153</sup>, Greg Kinney<sup>154</sup>, Barbara Konkle<sup>146</sup>, Holly Kramer<sup>155</sup>, Christoph Lange<sup>156</sup>, Ethan Lange<sup>95</sup>, Cathy Laurie<sup>70</sup>, Cecelia Laurie<sup>70</sup>, Meryl LeBoff<sup>96</sup>, Jiwon Lee<sup>96</sup>, Sandra Lee<sup>78</sup>, Wen-Jane Lee<sup>147</sup>, Jonathon LeFaive<sup>64</sup>, David Levine<sup>70</sup>, Dan Levy<sup>85</sup>, Joshua Lewis<sup>69</sup>, Yun Li<sup>118</sup>, Henry Lin<sup>105</sup>, Honghuang Lin<sup>157</sup>, Simin Liu<sup>158</sup>, Yongmei Liu<sup>159</sup>, Yu Liu<sup>160</sup>,**

Kathryn Lunetta<sup>157</sup>, James Luo<sup>85</sup>, Ulysses Magalang<sup>161</sup>, Michael Mahaney<sup>162</sup>, Barry Make<sup>73</sup>, Alisa Manning<sup>163</sup>, JoAnn Manson<sup>96</sup>, Lisa Martin<sup>164</sup>, Melissa Marton<sup>126</sup>, Susan Mathai<sup>95</sup>, Susanne May<sup>94</sup>, Patrick McArdle<sup>69</sup>, Merry-Lynn McDonald<sup>141</sup>, Sean McFarland<sup>152</sup>, Daniel McGoldrick<sup>165</sup>, Caitlin McHugh<sup>94</sup>, Becky McNeil<sup>166</sup>, Hao Mei<sup>71</sup>, James Meigs<sup>167</sup>, Vipin Menon<sup>78</sup>, Luisa Mestroni<sup>80</sup>, Ginger Metcalf<sup>78</sup>, Deborah A. Meyers<sup>168</sup>, Emmanuel Mignot<sup>169</sup>, Julie Mikulla<sup>85</sup>, Nancy Min<sup>71</sup>, Mollie Minear<sup>170</sup>, Ryan L. Minster<sup>87</sup>, Matt Moll<sup>101</sup>, Zeineen Momin<sup>78</sup>, Courtney Montgomery<sup>171</sup>, Donna Muzny<sup>78</sup>, Josyf C. Mychaleckyj<sup>104</sup>, Girish Nadkarni<sup>127</sup>, Rakhi Naik<sup>73</sup>, Sergei Nekhai<sup>172</sup>, Sarah C. Nelson<sup>94</sup>, Bonnie Neltner<sup>95</sup>, Caitlin Nessner<sup>78</sup>, Osuji Nkechinyere<sup>78</sup>, Jeff O'Connell<sup>173</sup>, Tim O'Connor<sup>69</sup>, Heather Ochs-Balcom<sup>174</sup>, Geoffrey Okwuonu<sup>78</sup>, Allan Pack<sup>175</sup>, David T. Paik<sup>176</sup>, James Pankow<sup>177</sup>, George Papanicolaou<sup>85</sup>, Cora Parker<sup>178</sup>, Juan Manuel Peralta<sup>119</sup>, Marco Perez<sup>75</sup>, James Perry<sup>69</sup>, Ulrike Peters<sup>179</sup>, Lawrence S. Phillips<sup>68</sup>, Jacob Pleiness<sup>64</sup>, Toni Pollin<sup>69</sup>, Wendy Post<sup>180</sup>, Julia Powers Becker<sup>181</sup>, Meher Preethi Boorgula<sup>95</sup>, Michael Preuss<sup>127</sup>, Pankaj Qasba<sup>85</sup>, Dandi Qiao<sup>96</sup>, Zhaohui Qin<sup>68</sup>, Nicholas Rafaels<sup>182</sup>, Laura Raffield<sup>183</sup>, Mahitha Rajendran<sup>78</sup>, Ramachandran S. Vasan<sup>157</sup>, D. C. Rao<sup>124</sup>, Laura Rasmussen-Torvik<sup>184</sup>, Aakrosh Ratan<sup>104</sup>, Robert Reed<sup>69</sup>, Catherine Reeves<sup>185</sup>, Elizabeth Regan<sup>110</sup>, Alex Reiner<sup>186</sup>, Muagututia S. Reupena<sup>33</sup>, Ken Rice<sup>70</sup>, Rebecca Robillard<sup>187</sup>, Nicolas Robine<sup>126</sup>, Dan Roden<sup>188</sup>, Carolina Roselli<sup>65</sup>, Ingo Ruczinski<sup>73</sup>, Alexi Runnels<sup>126</sup>, Pamela Russell<sup>95</sup>, Sarah Ruuska<sup>146</sup>, Kathleen Ryan<sup>69</sup>, Ester Cerdeira Sabino<sup>189</sup>, Danish Saleheen<sup>190</sup>, Shabnam Salimi<sup>69</sup>, Sejal Salvi<sup>78</sup>, Steven Salzberg<sup>73</sup>, Kevin Sandow<sup>191</sup>, Vijay G. Sankaran<sup>192</sup>, Jireh Santibanez<sup>78</sup>, Karen Schwander<sup>124</sup>, David Schwartz<sup>95</sup>, Frank Sciurba<sup>87</sup>, Christine Seidman<sup>193</sup>, Jonathan Seidman<sup>194</sup>, Frédéric Sériès<sup>195</sup>, Vivien Sheehan<sup>196</sup>, Stephanie L. Sherman<sup>197</sup>, Amol Shetty<sup>69</sup>, Aniket Shetty<sup>95</sup>, Wayne Hui-Heng Sheu<sup>147</sup>, M. Benjamin Shoemaker<sup>198</sup>, Brian Silver<sup>199</sup>, Edwin Silverman<sup>96</sup>, Robert Skomro<sup>200</sup>, Albert Vernon Smith<sup>201</sup>, Josh Smith<sup>70</sup>, Nicholas Smith<sup>138</sup>, Tanja Smith<sup>63</sup>, Sylvia Smoller<sup>202</sup>, Beverly Snively<sup>203</sup>, Michael Snyder<sup>75</sup>, Tamar Sofer<sup>96</sup>, Nona Sotoodehnia<sup>70</sup>, Adrienne M. Stilp<sup>70</sup>, Garrett Storm<sup>204</sup>, Elizabeth Streeten<sup>69</sup>, Jessica Lasky Su<sup>96</sup>, Yun Ju Sung<sup>124</sup>, Jody Sylvia<sup>96</sup>, Adam Szpiro<sup>70</sup>, Daniel Taliun<sup>64</sup>, Hua Tang<sup>205</sup>, Margaret Taub<sup>73</sup>, Matthew Taylor<sup>80</sup>, Simeon Taylor<sup>69</sup>, Marilyn Telen<sup>74</sup>, Timothy A. Thornton<sup>70</sup>, Machiko Threlkeld<sup>206</sup>, Lesley Tinker<sup>125</sup>, David Tirschwell<sup>70</sup>, Sarah Tishkoff<sup>207</sup>, Hemant Tiwari<sup>208</sup>, Catherine Tong<sup>209</sup>, Dhananjay Vaidya<sup>73</sup>, David Van Den Berg<sup>210</sup>, Peter VandeHaar<sup>64</sup>, Scott Vrieze<sup>177</sup>, Tarik Walker<sup>95</sup>, Robert Wallace<sup>143</sup>, Avram Walts<sup>95</sup>, Fei Fei Wang<sup>70</sup>, Heming Wang<sup>211</sup>, Jiongming Wang<sup>201</sup>, Karol Watson<sup>99</sup>, Jennifer Watt<sup>78</sup>, Daniel E. Weeks<sup>87</sup>, Joshua Weinstock<sup>150</sup>, Bruce Weir<sup>70</sup>, Scott T. Weiss<sup>212</sup>, Lu-Chen Weng<sup>213</sup>, Jennifer Wessel<sup>214</sup>, Kayleen Williams<sup>94</sup>, L. Keoki Williams<sup>215</sup>, Carla Wilson<sup>96</sup>, James Wilson<sup>216</sup>, Lara Winterkorn<sup>126</sup>, Quenna Wong<sup>70</sup>, Joseph Wu<sup>176</sup>, Huichun Xu<sup>69</sup>, Ivana Yang<sup>95</sup>, Ketian Yu<sup>64</sup>, Seyedeh Maryam Zekavat<sup>65</sup>, Yingze Zhang<sup>217</sup>, Snow Xueyan Zhao<sup>110</sup>, Wei Zhao<sup>218</sup>, Xiaofeng Zhu<sup>219</sup>, Michael Zody<sup>63</sup> & Sebastian Zoellner<sup>64</sup>

<sup>63</sup>New York Genome Center, New York, NY 10013, USA. <sup>64</sup>University of Michigan, Ann Arbor, MI 48109, USA. <sup>65</sup>Broad Institute, Cambridge, MA 02142, USA. <sup>66</sup>Cedars Sinai, Boston, MA 02114, USA. <sup>67</sup>Children's Hospital of Philadelphia, University of Pennsylvania, Philadelphia, PA 19104, USA. <sup>68</sup>Emory University, Atlanta, GA 30322, USA. <sup>69</sup>University of Maryland, Baltimore, MD 21201, USA. <sup>70</sup>University of Washington, Seattle, WA 98195, USA. <sup>71</sup>University of Mississippi, Jackson, MS 38677, USA. <sup>72</sup>National Institutes of Health, Bethesda, MD 20892, USA. <sup>73</sup>Johns Hopkins University, Baltimore, MD 21218, USA. <sup>74</sup>Duke University, Durham, NC 27708, USA. <sup>75</sup>Stanford University, Stanford, CA 94305, USA. <sup>76</sup>University of Wisconsin Milwaukee, Milwaukee, WI 53211, USA. <sup>77</sup>Providence Health Care, Medicine, Vancouver, USA. <sup>78</sup>Baylor College of Medicine Human Genome Sequencing Center, Houston, TX 77030, USA. <sup>79</sup>Cleveland Clinic, Cleveland, OH 44195, USA. <sup>80</sup>University of Colorado Anschutz Medical Campus, Aurora, CO 80045, USA. <sup>81</sup>Columbia University, New York, NY 10032, USA. <sup>82</sup>The Emmes Corporation, LTRC, Rockville, MD 20850, USA. <sup>83</sup>Cleveland Clinic, Quantitative Health Sciences, Cleveland, OH 44195, USA. <sup>84</sup>Johns Hopkins University, Medicine, Baltimore, MD 21218, USA. <sup>85</sup>National Heart, Lung, and Blood Institute, National Institutes of Health, Bethesda, MD 20892, USA. <sup>86</sup>Boston University, Massachusetts General Hospital, Boston University School of Medicine, Boston, MA 02114, USA. <sup>87</sup>University of Pittsburgh, Pittsburgh, PA 15260, USA. <sup>88</sup>Fundação de Hematologia e Hemoterapia de Pernambuco - Hemope, Recife 52011-000, Brazil. <sup>89</sup>National Jewish Health, National Jewish Health, Denver, CO 80206, USA. <sup>90</sup>Medical College of Wisconsin, Milwaukee, WI 53226, USA. <sup>91</sup>University of Texas Health at Houston, Pediatrics, Houston, TX 77030, USA. <sup>92</sup>University of California, San Francisco, San Francisco, CA 94143, USA. <sup>93</sup>Stanford University, Biomedical Data Science, Stanford, CA 94305, USA. <sup>94</sup>University of Washington, Biostatistics, Seattle, WA 98195, USA. <sup>95</sup>University of Colorado at Denver, Denver, CO 80204, USA. <sup>96</sup>Brigham & Women's Hospital, Boston, MA 02115, USA. <sup>97</sup>University of Montreal, Quebec, Canada. <sup>98</sup>Washington State University, Pullman, WA 99164, USA. <sup>99</sup>University of California, Los Angeles, Los Angeles, CA 90095, USA. <sup>100</sup>Brigham & Women's Hospital, Boston, USA. <sup>101</sup>Brigham & Women's Hospital, Medicine, Boston, MA 02115, USA. <sup>102</sup>National Taiwan University, Taipei 10617, Taiwan. <sup>103</sup>Brigham & Women's Hospital, Division of Preventive Medicine, Boston, MA 02215, USA. <sup>104</sup>University of Virginia, Charlottesville, VA 22903, USA. <sup>105</sup>Lundquist Institute, Torrance, CA 90502, USA. <sup>106</sup>Cleveland Clinic, Cleveland Clinic, Cleveland, OH 44195, USA. <sup>107</sup>Broad Institute, Metabolomics Platform, Cambridge, MA 02142, USA. <sup>108</sup>Cleveland Clinic, Immunity and Immunology, Cleveland, OH 44195, USA. <sup>109</sup>University of Vermont, Burlington, VT 05405, USA. <sup>110</sup>National Jewish Health, Denver, CO 80206, USA. <sup>111</sup>Boston University, Biostatistics, Boston, MA 02115, USA. <sup>112</sup>Vitalant Research Institute, San Francisco, CA 94118, USA. <sup>113</sup>University of Illinois at Chicago, Chicago, IL 60607, USA. <sup>114</sup>University of Chicago, Chicago, IL 60637, USA. <sup>115</sup>Mayo Clinic, Health Quantitative Sciences Research, Rochester, MN 55905, USA. <sup>116</sup>Vanderbilt University, Nashville, TN 37235, USA. <sup>117</sup>University of Cincinnati, Cincinnati, Ohio 45220, USA. <sup>118</sup>University of North Carolina, Chapel Hill, NC 27599, USA. <sup>119</sup>University of Texas Rio Grande Valley School of Medicine, Edinburg, TX 78539, USA. <sup>120</sup>Brown University, Providence, RI 02912, USA. <sup>121</sup>Harvard University, Channing Division of Network Medicine, Cambridge, MA 02138, USA. <sup>122</sup>National Jewish Health, Center for Genes, Environment and Health, Denver, CO 80206, USA. <sup>123</sup>University of North Carolina, Epidemiology, Chapel Hill, NC 27599, USA. <sup>124</sup>Washington University in St Louis, St Louis, MO 63130, USA. <sup>125</sup>Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA. <sup>126</sup>New York Genome Center, New York City, NY 10013, USA. <sup>127</sup>Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA. <sup>128</sup>University of Pittsburgh, Pittsburgh, PA, USA. <sup>129</sup>Beth Israel Deaconess Medical Center, Boston, MA 02215, USA. <sup>130</sup>Boston Children's Hospital, Harvard Medical School, Department of Psychiatry, Boston, MA 02115, USA. <sup>131</sup>University of Texas Rio Grande Valley School of Medicine, San Antonio,

TX 78229, USA. <sup>132</sup>Mass General Brigham, Obstetrics and Gynecology, Boston, MA 02115, USA. <sup>133</sup>Indiana University, OB/GYN, Indianapolis, Indiana 46202, USA. <sup>134</sup>University of Mississippi, Cardiology, Jackson, MS 39216, USA. <sup>135</sup>University of Calgary, Medicine, Calgary, Canada. <sup>136</sup>University of Maryland, Genetics, Philadelphia, PA 19104, USA. <sup>137</sup>Yale University, Department of Chronic Disease Epidemiology, Connecticut 06520, USA. <sup>138</sup>University of Washington, Epidemiology, Seattle, WA 98195, USA. <sup>139</sup>Wake Forest Baptist Health, Winston-Salem, NC 27157, USA. <sup>140</sup>Brigham & Women's Hospital, Channing Division of Network Medicine, Boston, MA 02115, USA. <sup>141</sup>University of Alabama, Birmingham, AL 35487, USA. <sup>142</sup>University of Texas Health at Houston, Houston, TX 77225, USA. <sup>143</sup>University of Iowa, Iowa City, IA 52242, USA. <sup>144</sup>National Health Research Institute Taiwan, Institute of Population Health Sciences, NHRI, Miaoli County 350, Taiwan. <sup>145</sup>Tri-Service General Hospital National Defense Medical Center, Taipei, Taiwan. <sup>146</sup>Blood Works Northwest, Seattle, WA 98104, USA. <sup>147</sup>Taichung Veterans General Hospital Taiwan, Taichung City 407, Taiwan. <sup>148</sup>Oklahoma State University Medical Center, Internal Medicine, Division of Endocrinology, Diabetes and Metabolism, Columbus, OH 43210, USA. <sup>149</sup>Blood Works Northwest, Research Institute, Seattle, WA 98104, USA. <sup>150</sup>University of Michigan, Biostatistics, Ann Arbor, MI 48109, USA. <sup>151</sup>University of California, San Francisco, San Francisco, CA 94118, USA. <sup>152</sup>Harvard University, Cambridge, MA 02138, USA. <sup>153</sup>McGill University, Montréal, QC H3A 0G4, Canada. <sup>154</sup>University of Colorado at Denver, Epidemiology, Aurora, CO 80045, USA. <sup>155</sup>Loyola University, Public Health Sciences, Maywood, IL 60153, USA. <sup>156</sup>Harvard School of Public Health, Biostats, Boston, MA 02115, USA. <sup>157</sup>Boston University, Boston, MA 02215, USA. <sup>158</sup>Brown University, Epidemiology and Medicine, Providence, RI 02912, USA. <sup>159</sup>Duke University, Cardiology, Durham, NC 27708, USA. <sup>160</sup>Stanford University, Cardiovascular Institute, Stanford, CA 94305, USA. <sup>161</sup>Ohio State University, Division of Pulmonary, Critical Care and Sleep Medicine, Columbus, OH 43210, USA. <sup>162</sup>University of Texas Rio Grande Valley School of Medicine, Brownsville, TX 78520, USA. <sup>163</sup>Broad Institute, Harvard University, Massachusetts General Hospital, Cambridge, USA. <sup>164</sup>George Washington University, cardiology, Washington, DC 20037, USA. <sup>165</sup>University of Washington, Genome Sciences, Seattle, WA 98195, USA. <sup>166</sup>RTI International, North Carolina, USA. <sup>167</sup>Massachusetts General Hospital, Medicine, Boston, MA 02114, USA. <sup>168</sup>University of Arizona, Tucson, AZ 85721, USA. <sup>169</sup>Stanford University, Center For Sleep Sciences and Medicine, Palo Alto, CA 94304, USA. <sup>170</sup>National Institute of Child Health and Human Development, National Institutes of Health, Bethesda, MD 20892, USA. <sup>171</sup>Oklahoma Medical Research Foundation, Genes and Human Disease, Oklahoma City, OK 73104, USA. <sup>172</sup>Howard University, Washington, DC 20059, USA. <sup>173</sup>University of Maryland, Baltimore, MD 21201, USA. <sup>174</sup>University at Buffalo, Buffalo, NY 14260, USA. <sup>175</sup>University of Pennsylvania, Division of Sleep Medicine/Department of Medicine, Philadelphia, PA 19104-3403, USA. <sup>176</sup>Stanford University, Stanford Cardiovascular Institute, Stanford, CA 94305, USA. <sup>177</sup>University of Minnesota, Minneapolis, MN 55455, USA. <sup>178</sup>RTI International, Biostatistics and Epidemiology Division, Research Triangle Park, North Carolina 27709-2194, USA. <sup>179</sup>Fred Hutchinson Cancer Research Center, Fred Hutch and UW, Seattle, WA 98109, USA. <sup>180</sup>Johns Hopkins University, Cardiology/Medicine, Baltimore, MD 21218, USA. <sup>181</sup>University of Colorado at Denver, Medicine, Denver, CO 80204, USA. <sup>182</sup>University of Colorado at Denver, Denver, CO 80045, USA. <sup>183</sup>University of North Carolina, Genetics, Chapel Hill, NC 27599, USA. <sup>184</sup>Northwestern University, Chicago, IL 60208, USA. <sup>185</sup>New York Genome Center, New York Genome Center, New York City, NY 10013, USA. <sup>186</sup>Fred Hutchinson Cancer Research Center, University of Washington, Seattle, WA 98109, USA. <sup>187</sup>University of Ottawa, Sleep Research Unit, University of Ottawa Institute for Mental Health Research, Ottawa, ON K1Z 7K4, Canada. <sup>188</sup>Vanderbilt University, Medicine, Pharmacology, Biomedical Informatics, Nashville, TN 37235, USA. <sup>189</sup>Universidade de Sao Paulo, Faculdade de Medicina, Sao Paulo 01310000, Brazil. <sup>190</sup>Columbia University, New York, NY 10027, USA. <sup>191</sup>Lundquist Institute, TGPS, Torrance, CA 90502, USA. <sup>192</sup>Harvard University, Division of Hematology/Oncology, Boston, MA 02115, USA. <sup>193</sup>Harvard Medical School, Genetics, Boston, MA 02115, USA. <sup>194</sup>Harvard Medical School, Boston, MA 02115, USA. <sup>195</sup>Université Laval, Quebec City G1V 0A6, Canada. <sup>196</sup>Emory University, Pediatrics, Atlanta, GA 30307, USA. <sup>197</sup>Emory University, Human Genetics, Atlanta, GA 30322, USA. <sup>198</sup>Vanderbilt University, Medicine/Cardiology, Nashville, TN 37235, USA. <sup>199</sup>UMass Memorial Medical Center, Worcester, MA 01655, USA. <sup>200</sup>University of Saskatchewan, Saskatoon, SK S7N 5C9, USA. <sup>201</sup>University of Michigan, Ann Arbor, USA. <sup>202</sup>Albert Einstein College of Medicine, New York, NY 10461, USA. <sup>203</sup>Wake Forest Baptist Health, Biostatistical Sciences, Winston-Salem, NC 27157, USA. <sup>204</sup>University of Colorado at Denver, Genomic Cardiology, Aurora, CO 80045, USA. <sup>205</sup>Stanford University, Genetics, Stanford, CA 94305, USA. <sup>206</sup>University of Washington, University of Washington, Department of Genome Sciences, Seattle, WA 98195, USA. <sup>207</sup>University of Pennsylvania, Genetics, Philadelphia, PA 19104, USA. <sup>208</sup>University of Alabama, Biostatistics, Birmingham, AL 35487, USA. <sup>209</sup>University of Washington, Department of Biostatistics, Seattle, WA 98195, USA. <sup>210</sup>University of Southern California, USC Methylation Characterization Center, University of Southern California, California 90033, USA. <sup>211</sup>Brigham & Women's Hospital, Mass General Brigham, Boston, MA 02115, USA. <sup>212</sup>Brigham & Women's Hospital, Channing Division of Network Medicine, Department of Medicine, Boston, MA 02115, USA. <sup>213</sup>Massachusetts General Hospital, Boston, MA 02114, USA. <sup>214</sup>Indiana University, Epidemiology, Indianapolis, Indiana 46202, USA. <sup>215</sup>Henry Ford Health System, Detroit, MI 48202, USA. <sup>216</sup>Beth Israel Deaconess Medical Center, Cardiology, Cambridge, MA 02139, USA. <sup>217</sup>University of Pittsburgh, Medicine, Pittsburgh, PA 15260, USA. <sup>218</sup>University of Michigan, Department of Epidemiology, Ann Arbor, MI 48109, USA. <sup>219</sup>Case Western Reserve University, Department of Population and Quantitative Health Sciences, Cleveland, OH 44106, USA.