



OPEN

Elucidating the genetic architecture of DNA methylation to identify promising molecular mechanisms of disease

Jiantao Ma^{1,11}✉, Roby Joehanes^{2,3,11}, Chunyu Liu^{4,5,6}, Amena Keshawarz^{2,3}, Shih-Jen Hwang^{2,3}, Helena Bui^{2,3}, Brandon Tejada^{2,3}, Meera Sooda^{2,3}, Peter J. Munson^{2,3}, Cumhur Y. Demirkale⁷, Paul Courchesne^{2,3}, Nancy L. Heard-Costa^{5,6,8}, Achilleas N. Pitsillides^{4,5,6}, Mike Feolo⁹, Nataliya Sharopova⁹, Ramachandran S. Vasan^{5,6,8}, Tianxiao Huan¹⁰ & Daniel Levy^{2,3}✉

DNA methylation commonly occurs at cytosine-phosphate-guanine sites (CpGs) that can serve as biomarkers for many diseases. We analyzed whole genome sequencing data to identify DNA methylation quantitative trait loci (mQTLs) in 4126 Framingham Heart Study participants. Our mQTL mapping identified 94,362,817 *cis*-mQTL variant-CpG pairs (for 210,156 unique autosomal CpGs) at $P < 1e-7$ and 33,572,145 *trans*-mQTL variant-CpG pairs (for 213,606 unique autosomal CpGs) at $P < 1e-14$. Using *cis*-mQTL variants for 1258 CpGs associated with seven cardiovascular disease (CVD) risk factors, we found 104 unique CpGs that colocalized with at least one CVD trait. For example, cg11554650 (*PPP1R18*) colocalized with type 2 diabetes, and was driven by a single nucleotide polymorphism (rs2516396). We performed Mendelian randomization (MR) analysis and demonstrated 58 putatively causal relations of CVD risk factor-associated CpGs to one or more risk factors (e.g., cg05337441 [*APOB*] with LDL; MR $P = 1.2e-99$, and 17 causal associations with coronary artery disease (e.g., cg08129017 [*SREBF1*] with coronary artery disease; MR $P = 5e-13$). We also showed that three CpGs, e.g., cg14893161 (*PM20D1*), are putatively causally associated with COVID-19 severity. To assist in future analyses of the role of DNA methylation in disease pathogenesis, we have posted a comprehensive summary data set in the National Heart, Lung, and Blood Institute's BioData Catalyst.

DNA methylation, the most frequently studied epigenetic modification, involves the transfer of a methyl group to the fifth carbon position of the cytosine DNA nucleotide to form 5-methylcytosine¹. DNA methylation is influenced both by genetic and environmental factors and may mediate gene-environment interactions; therefore, it may be used to determine the risk of many complex diseases through its critical role in gene expression regulation^{2,3}. Associations between DNA methylation and a wide range of phenotypes have been identified by epigenome-wide association studies (EWAS)⁴⁻⁶. DNA methylation therefore can serve both as a biomarker for disease and contribute to its pathogenesis.

Identification of genetic loci associated with the methylation of cytosine-phosphate-guanine sites (CpGs)—i.e., DNA methylation quantitative trait loci (mQTLs)—can facilitate the interpretation of the biological underpinnings of disease relations and causal inference regarding the role of DNA methylation in disease. Genome-wide association studies (GWAS) have successfully identified many disease-associated genetic variants⁷. Molecular

¹Division of Nutrition Epidemiology and Data Science, Friedman School of Nutrition Science and Policy, Tufts University, Boston, MA, USA. ²Present address: Population Sciences Branch, National Heart, Lung, and Blood Institute, National Institutes of Health, Bethesda, MD, USA. ³Framingham Heart Study, Framingham, MA, USA. ⁴Department of Biostatistics, School of Public Health, Boston University, Boston, MA, USA. ⁵Boston University, Boston, MA, USA. ⁶National Heart, Lung, and Blood Institute, Framingham Heart Study, Framingham, MA, USA. ⁷Critical Care Medicine Department, Clinical Center, National Institutes of Health, Bethesda, MD, USA. ⁸Boston University School of Medicine, Boston University, Boston, MA, USA. ⁹National Center for Biotechnology Information, Bethesda, MD, USA. ¹⁰Department of Ophthalmology and Visual Sciences, University of Massachusetts Medical School, Worcester, MA, USA. ¹¹These authors contributed equally: Jiantao Ma and Roby Joehanes. ✉email: jiantao.ma@tufts.edu; levyd@nhlbi.nih.gov

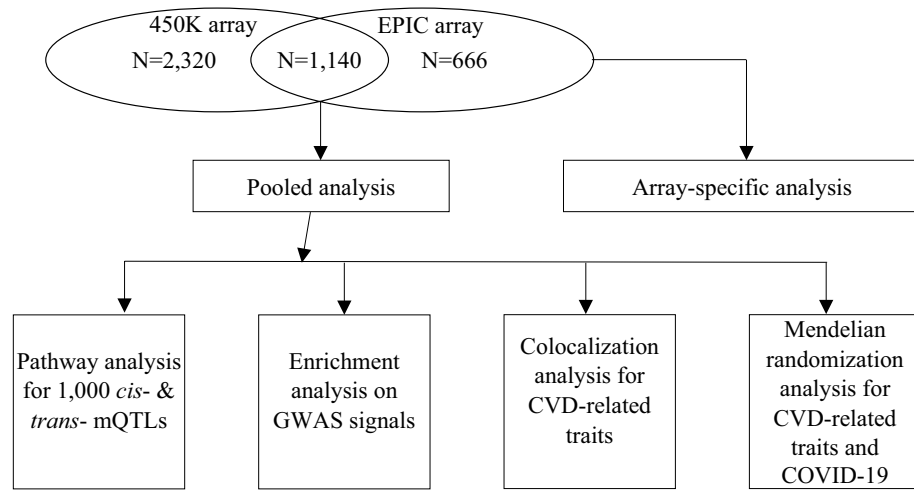


Figure 1. Study design.

mechanisms linking these variants to disease, however, are not fully understood. Exploring colocalization of disease-associated genetic variants from GWAS with mQTL variants may further reveal molecular mechanisms underlying the associations between genetic variants and diseases⁸. We hypothesize that by studying the overlap of mQTL variants with known disease-associated genetic variants from GWAS, we can further explore the joint contributions of genetic and environmental influences to diseases. Furthermore, by colocalizing mQTLs with genetic variants associated with gene expression (expression quantitative trait loci, eQTLs), we can better interpret the biological functions of disease-associated CpGs⁸. Utilizing effect sizes derived from GWAS for mQTL variants with different diseases, we can conduct causal inference testing to explore the putative causal roles of CpGs on a wide range of diseases^{9–14}.

In our earlier work, we performed GWAS of ~415,000 CpGs in whole blood derived DNA in Framingham Heart Study (FHS) participants with validation in the Atherosclerosis Risk in Communities (ARIC) study and the Grady Trauma Project (GTP)¹⁵. Genotyping was performed using commercial arrays with imputation across the genome. The present study greatly expands on our prior work by incorporating whole genome sequencing (WGS) data in FHS participants obtained as part of the National Heart, Lung, and Blood Institute's (NHLBI) Trans-Omics for Precision Medicine (TOPMed) Program (<https://www.nhlbiwgs.org/>). Use of WGS greatly reduces imputation uncertainty and vastly increases coverage of variation across the human genome. In this study, we utilized state-of-the-art WGS in conjunction with DNA methylation measured by commercial arrays to quantify single nucleotide polymorphism (SNP)-CpG associations in over 4000 FHS participants (Fig. 1). Our primary goal was to create a robust mQTL resource to better understand the genetic architecture of DNA methylation and facilitate the discovery of molecular mechanisms underlying a variety of diseases. We also provide examples of how mQTLs can be used in colocalization and Mendelian randomization (MR) analyses to infer the causal roles of DNA methylation in relation to disease phenotypes, with a focus on cardiovascular disease (CVD) risk factors and severity of coronavirus disease 2019 (COVID-19).

Results

Participant characteristics. As shown in Table 1, our pooled analysis included 4126 participants (2320 with DNA methylation data from the 450K array and 1806 with data from the EPIC array). In the FHS Offspring cohort, blood samples used for the 450K array measurements were collected ~6 years earlier than those for the EPIC array measurements, while blood samples for both arrays in the Third Generation cohort were obtained at the same visit. Therefore, the mean age for participants with EPIC array data was older than that for the 450K array. There were no substantial differences in sex, BMI, or other CVD risk factors.

mQTL mapping. Our primary pooled analysis examined association of 20,696,115 SNPs with 452,567 whole blood derived CpGs. In the pooled analysis, we identified 94,362,817 *cis*-mQTL variant-CpG pairs (details in "Methods") for 210,156 unique autosomal CpGs and at $P < 1e-7$ and 33,572,145 *trans*-mQTL-CpG pairs for 213,606 unique autosomal CpGs at $P < 1e-14$. The numbers of *cis*- and *trans*-mQTL variant-CpG pairs for each chromosome are presented in Supplemental Table 1. The *cis*-mQTL variants accounted for 0.7% to 79.9% (median 1.6%) of heritability of DNA methylation, and *trans*-mQTLs accounted for 1.4% to 78.7% (median 2.1%) of heritability. There were 1,080,716 *cis*-mQTL variants, associated with 31,422 unique CpGs (2,345,086 or 2.5% of the 94,362,817 *cis*-mQTL variant-CpG pairs), that accounted for $\geq 20\%$ of heritability of DNA methylation at the corresponding CpGs (Fig. 2). We also observed that 185,167 *trans*-mQTL variants accounted for $\geq 20\%$ of heritability of DNA methylation for 2711 unique CpGs (314,660 or 0.9% of the 33,572,145 *trans*-mQTL variant-CpG pairs; Fig. 2). The array-specific results are presented in Supplemental Table 2.

	Third generation cohort	Offspring cohort	Omni cohort
N	1945	2129	52
Women (%)	52.7	54.7	46.2
Age (years)	46 ± 8	67 ± 9	70 ± 9
BMI: kg/m ²	27.9 ± 5.7	28.2 ± 5.3	27.8 ± 4.9
Systolic blood pressure (mm Hg)	116 ± 14	129 ± 17	126 ± 16
Diastolic blood pressure (mm Hg)	74 ± 9	73 ± 10	68 ± 10
Hypertension (%)	33.9	53.2	42.3
Fasting glucose (mg/dL)	96.7 ± 20.8	106.3 ± 22.1	100.5 ± 12.6
Diabetes (%)	6.0	14.7	9.6
Triglyceride (mg/dL)	112 ± 79.3	119.8 ± 71.9	94.9 ± 43.9
High density lipoprotein (mg/dL)	59.3 ± 17.2	57.3 ± 18.1	65.5 ± 19.3
Low density lipoprotein (mg/dL)	104.5 ± 29.4	104.8 ± 31.1	87.9 ± 28.9

Table 1. Characteristics of the study population. Values are mean ± SD.

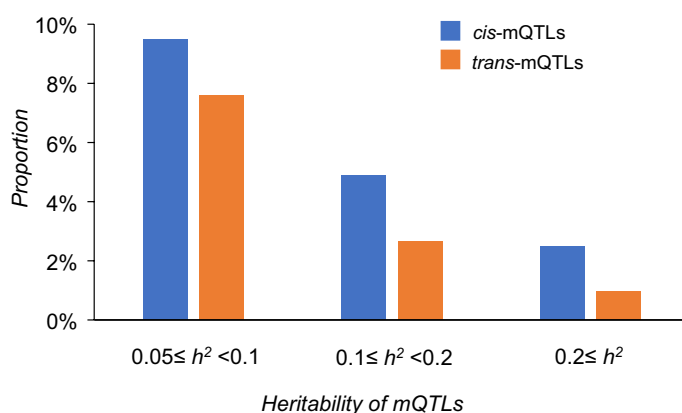


Figure 2. Heritability of DNA methylation explained by the *cis*- and *trans*-mQTLs identified in the pooled analysis. The total number of *cis*-mQTLs is 94,362,817 and the total number of *trans*-mQTLs is 32,434,987.

We examined whether whole blood derived mQTL variant-CpG pairs identified by other studies^{16,17} were significant in our dataset. For the top independent *cis*-mQTL variant-CpG pairs (168,675 pairs for 104,619 CpGs; $P < 1e-7$) identified by the pooled analysis in five Dutch biobanks¹⁶, 66.1% of the pairs (111,557 pairs for 79,099 CpGs) overlapped with our *cis*-mQTL variant-CpG pairs ($P < 1e-7$ with consistent effect direction). For the top independent *trans*-pairs (5865 pairs for 2066 CpGs with $P < 1e-14$) in the Dutch biobanks, 38.4% of the pairs (2250 pairs for 866 CpGs) overlapped with our *trans*-mQTL variant-CpG pairs ($P < 1e-14$ with consistent effect direction). Using blood samples collected from 3799 Europeans and 3195 South Asians, Hawe et al. identified 10,346,172 *cis*- and 819,387 *trans*-mQTL variant-CpG pairs at $P < 1e-14$ in a cross-ancestry analysis¹⁷. Compared to their study, at $P < 1e-14$, we identified 41,224,533 more *cis*-mQTL variant-CpG pairs and 32,752,758 more *trans*-mQTL variant-CpG pairs. Among the 10,346,172 *cis*- and 819,387 *trans*-mQTL variant-CpG pairs reported by Hawe et al.¹⁷, 78.3% ($n = 8,105,456$) of *cis*-mQTL variant-CpG pairs and 66% ($n = 540,851$) of *trans*-mQTL variant-CpG pairs were significant and had consistent effect direction in our mQTL database, respectively.

GO analysis for *cis*- and *trans*-mQTLs. To investigate the biological implications for diseases, we conducted Gene Ontology (GO) analysis to identify biological processes, cellular components, and molecular functions that are impacted by the detected mQTLs. Using the top 1000 unique *cis*-mQTL variants from the pooled analysis, we identified 19 significant GO pathways (16 for Biological Process and 3 for Cellular Component) at $FDR < 0.05$ (Supplemental Table 3); the top Biological Process term was dendrite development (GO:0016358; $P = 8.8e-7$; $FDR = 0.01$) and the top Cellular Component term was cell periphery (GO:0071944; $P = 3.4e-6$; $FDR = 0.01$). The top 1000 unique *trans*-mQTL variants from the pooled analysis were linked to nine significant GO pathways (six for Biological Process and three for Cellular Component) at $FDR < 0.05$ (Supplemental Table 4); the top Biological Process term was cellular component organization or biogenesis (GO:0071840; $P = 7.8e-7$; $FDR = 0.009$) and the top Cellular Component term was cytoplasm (GO: 0005737; $P = 3.4e-6$; $FDR = 0.009$).

Enrichment analysis of mQTL GWAS signals. To illustrate the potential health consequences associated with the mQTL variants, we examined the overlap between the detected mQTL variants with GWAS SNPs

CpG	Chr	Position	Annotated Gene of CpGs	Relation to CpG island*	Distance to TSS	CVD risk factors associated with CpGs in EWAS catalog	Colocalized traits	Colocalized <i>cis</i> -mQTL variants	PPFC	PPE %
cg14509967	17	48601772	<i>HOXB6</i>	S_shelf	6021	BMI	BMI	rs9299	1	100
cg01856529	12	54259306	<i>CBX5</i>		28,364	HDL	BMI	rs4759073	0.9869	74.42
cg27087650	19	44752538	<i>BCL3</i>	N_shore	4833	BMI	CAD	rs62117206	0.889	46.72
cg05337441	2	21043695	<i>APOB</i>	N_shore	42,266	LDL	CAD	rs668948	0.7996	41.29
cg03676485	7	2524181	<i>LFNG</i>	Island	11,652	HDL	DBP	rs4632959	0.9999	100
cg14099685	11	47524515	<i>CUGBP1</i>		58,578	SBP	DBP	rs34312154	0.993	57.7
cg21506299	6	136784086	<i>MAP3K5</i>		227,040	BMI	HDL	rs6924387	0.9978	99.74
cg27087650	19	44752538	<i>BCL3</i>	N_shore	4833	BMI	HDL	rs1531517	0.9961	59.67
cg27087650	19	44752538	<i>BCL3</i>	N_shore	4833	BMI	LDL	rs4803750	0.9989	100
cg03725309	1	109214962	<i>SARS1</i>	S_shore	1069	BMI:SBP:TG	LDL	rs4970829	0.9962	100
cg03676485	7	2524181	<i>LFNG</i>	Island	11,652	HDL	SBP	rs4632959	0.9999	99.99
cg20278790	20	59008418	<i>CTSZ</i>	S_shore	13,233	WC	SBP	rs151343	0.9979	100
cg11554650	6	30685413	<i>PPP1R18</i>	N_shore	9024	BMI	T2D	rs2516396	0.9835	100
cg00973118	16	324569	<i>AXINI</i>	N_shore	37,129	BMI	T2D	rs8049265	0.9505	99.33
cg27087650	19	44752538	<i>BCL3</i>	N_shore	4833	BMI	TG	rs4803750	0.9988	98.56
cg14099685	11	47524515	<i>CELFI</i>		58,578	SBP	TG	rs34312154	0.9927	54.52

Table 2. Top colocalization analysis results for CVD risk factor-associated CpGs. The top two colocalization results are presented for each outcome trait of interest (BMI, SBP, DBP, HDL, LDL, TG, T2D, & CAD). BMI, body mass index; SBP, systolic blood pressure; DBP, diastolic blood pressure; HDL, high density lipoprotein cholesterol; LDL, low density lipoprotein cholesterol; TG, triglyceride; T2D, type 2 diabetes; CAD, coronary artery disease; PPFC, posterior probability of colocalization; PPE, proportion of PPFC explained by the listed SNP. *Shore is region from 1 to 2 Kb away from a CpG island and shelf is located at 2–4 Kb from a CpG island. Prefixes N_ and S_ represent north (i.e., upstream) and south (i.e., downstream) regions relative to a CpG island. TSS, transcription start site.

included in the GWAS Catalog⁷. We examined 9,395,367 *cis*-mQTL variants and 6,039,960 *trans*-mQTL variants located in all autosomal chromosomes identified by the pooled analysis. The enrichment analysis showed that, at FDR < 0.05, the *cis*-mQTL variants were enriched with GWAS SNPs associated with 783 traits, representing 27.1% of the traits included in the GWAS Catalog⁷. For example, we found enrichment of SNPs associated with BMI (enrichment $P = 2e-305$ for BMI; Supplemental Table 5), systolic BP (enrichment $P = 2e-305$, triglyceride level (enrichment $P = 3.5e-231$), type 2 diabetes (enrichment $P = 2.4e-194$), and coronary artery disease (enrichment $P = 4.8e-118$). Compared to the *cis*-mQTL variants, the number of enriched GWAS traits for the *trans*-mQTL variants was lower with enrichment for nine GWAS traits (Supplemental Table 6).

Colocalization analysis. We tested 1258 CVD risk factor-associated CpGs for colocalization with five CVD-related traits to further explore the clinical implication of the detected mQTLs. We found that 104 unique CpGs colocalized with at least one CVD-related traits at PPFC threshold ≥ 0.7 (overall 155 colocalized pairs; Supplemental Tables 7 and 8). In Table 2, we present the top two CpGs that colocalized with each CVD-related trait. For example, cg11554650 (*PPP1R18*), a BMI-associated CpG on chromosome 6, colocalized with type 2 diabetes at SNP rs2516396 (PPFC = 0.98), which explained 100% of the observed PPFC; cg05337441 (*APOB*), an LDL-associated CpG at chromosome 2, colocalized with coronary artery disease at rs668948 (PPFC = 0.8), which explained 41% of the observed PPFC; and cg03676485 (*LFNG*), a HDL-associated CpG at chromosome 7, colocalized with systolic and diastolic BP at rs4632959 (PPFC = 0.99), which explained 100% of the observed PPFC.

Mendelian randomization analysis. To further demonstrate the clinical implications of the detected mQTLs in the development of complex diseases, we performed MR analysis to test the putative causal relationships between mQTL and CVD risk factors and COVID-19 severity. Using the *cis*-mQTL variants for the 1258 CVD risk factor-associated CpGs ($P < 1e-6$) reported in the EWAS catalog, we conducted MR analysis to test for putatively causal relations of CVD risk factor-associated CpGs with the corresponding CVD risk factors (e.g., HDL-associated-CpGs with HDL and fasting glucose associated-CpGs with type 2 diabetes). After Bonferroni correction for the number of tests in analysis for each trait (e.g., 0.05/566 or $8.8e-5$ in analysis for BMI), we identified 58 significant MR associations (Supplemental Tables 9). Information for the *cis*-mQTL variants and their corresponding CpGs is presented in Supplemental Table 8. The top three CpG-trait pairs reflected increased methylation levels at cg05337441 (*APOB*) with lower LDL (MR effect size: -2.94 ± 0.14 , $P = 1.2e-99$), increased methylation levels at cg26663590 (closest gene is *NFATC2IP* in UCSC genome browser) with lower BMI (MR effect size: -1.39 ± 0.13 , $P = 6.3e-26$), and increased methylation levels at cg14099685 (*CELFI*) with higher systolic BP (MR effect size: 138.64 ± 14.85 , $P = 9.9e-21$). We also demonstrated that 17 CVD risk factor-CpGs

CpG	Chr	Position	Gene	Relation to CpG island	Distance to TSS	N of IVs	Beta	SE	P
cg08129017	17	17825345	<i>SREBF1</i>	S_Shore	14,011	20	1.81	0.25	5.0e-13
cg00184953	6	31178444	<i>PSORS1C3</i>	N_Shelf	4709	7	-4.86	0.70	4.7e-12
cg02050917	1	2242131	<i>SKI</i>		13,812	13	2.65	0.39	1.4e-11
cg05337441	2	21043695	<i>APOB</i>	N_Shore	42,266	17	1.26	0.19	3.5e-11
cg21587837	6	31558116	<i>NFKBIL1</i>		11,265	33	-2.50	0.41	1.3e-9
cg03725309	1	109214962	<i>SARS1</i>	S_Shore	1069	1	12.59	2.25	2.2e-8
cg04545296	12	48351459	<i>ZNF641</i>	S_Shore	16,872	29	-0.89	0.16	2.6e-8
cg26562921	16	84726822	<i>USP10</i>		26,822	18	1.59	0.29	6.5e-8
cg20544516	17	17813868	<i>MIR33B;SREBF1</i>	S_Shore	32	1	7.50	1.43	1.5e-7
cg21242002	4	3263352	<i>MSANTD1</i>		19,079	4	-4.93	0.96	2.7e-7
cg08244301	19	17499941	<i>SLC27A1</i>	N_Shore	31,174	4	3.71	0.75	7.1e-7
cg21053741	6	31558083	<i>NFKBIL1</i>		11,232	7	-3.39	0.70	1.2e-6
cg27087650	19	44752538	<i>BCL3</i>	N_Shore	4833	1	7.36	1.63	6.1e-6
cg10101600	2	43251603	<i>THADA</i>		20,752	2	6.65	1.53	1.4e-5
cg19224164	4	2964656	<i>GRK4;NOP14</i>	S_Shore	1085	13	-1.67	0.40	2.3e-5
cg18933331	1	109643795		S_Shore		2	-5.92	1.42	3.0e-5
cg12467090	1	204490010	<i>PIK3C2B</i>		67,377	3	-4.65	1.12	3.3e-5

Table 3. Mendelian randomization analysis for CVD risk factor-associated CpGs (exposure) with coronary artery disease (outcome). IV: instrument variables, i.e., independent *cis*-mQTL variants with linkage disequilibrium $R^2 < 0.1$. Beta, SE, and *P* are derived from Mendelian randomization analysis using the inverse variance weighted (IVW) method.

CpG	CHR	BP	Gene	Relation to CpG island	Distance to TSS	COVID-19 Host Genetics Initiative GWAS (release 6)				GenOMICC study			
						N of IVs	Beta	SE	P	N of IVs	Beta	SE	P
cg07796016	1	152779584	<i>LCE1C</i>		2275	19	0.08	0.26	0.76	22	-0.35	0.28	0.20
cg14859874	1	154238265	<i>UBAP2L</i>		45,616	36	0.64	0.21	0.002	37	0.65	0.18	0.0002
cg17515347	1	159047163	<i>AIM2</i>		22,321	13	1.13	0.46	0.02	13	0.34	0.44	0.44
cg17178900	1	205818956	<i>PM20D1</i>	Island	21,802	31	-0.54	0.16	0.0007	31	-0.46	0.17	0.008
cg14893161	1	205819251	<i>PM20D1</i>	S_Shore	22,097	34	-0.69	0.17	6.1e-05	38	-0.54	0.17	0.002
cg08309069	6	31240651	<i>HLA-C</i>	S_Shore	4124	46	-0.40	0.19	0.04	45	0.31	0.28	0.28
cg05030953	6	31241000	<i>HLA-C</i>	S_Shore	4473	50	-0.24	0.15	0.10	43	0.47	0.18	0.01
cg02872426	6	110736772	<i>DDO</i>		27,247	28	0.77	0.27	0.004	27	0.39	0.28	0.16
cg12682382	8	74787918	<i>UBE2W</i>	N_Shelf	95,586	28	0.18	0.17	0.29	27	0.03	0.19	0.88
cg13571460	9	124989337	<i>LHX6</i>	Island	24,475	17	0.10	0.35	0.78	18	-0.21	0.33	0.53

Table 4. Mendelian randomization analysis of putatively causal relations of COVID-19 severity-associated CpGs to COVID-19 severity. IVs are independent *cis*-mQTL variants with linkage disequilibrium $R^2 < 0.1$. Inverse variance weighted (IVW) method was used to conduct Mendelian randomization analysis, using data from two COVID-19 GWAS (COVID-19 Host Genetics Initiative²⁰ and GenOMICC study¹⁹).

were associated with coronary artery disease (Table 3; corresponding $P < 3.9e-5$), e.g., at cg08129017 (*SREBF1*; reported as associated with BMI and triglyceride in the EWAS catalog) and cg02050917 (*SKI*; BMI-associated CpG), higher methylation levels were associated increased CVD risk, MR effect size: 1.81 ± 0.25 ; $P = 5e-13$ and 2.65 ± 0.39 , $P = 1.4e-11$, respectively.

A recent study conducted in 407 patients with COVID-19 showed that whole blood derived DNA methylation levels at 23 CpGs (annotated to 20 genes) were associated with COVID-19 severity¹⁸. We found that ten of the 23 COVID-19 severity-associated CpGs had at least one *cis*-mQTL variant in our database. We used independent *cis*-mQTL variants (linkage disequilibrium $R^2 < 0.1$), which overlapped with the SNPs tested by the two COVID-19 severity GWAS^{19,20}, to conduct MR analyses. As shown in Table 4, we observed that higher methylation levels at cg14893161 (*PM20D1*; $P = 6e-5$ and 0.002 for the two COVID GWAS, respectively), lower methylation levels at cg17178900 (*PM20D1*; $P = 7e-4$ and 0.008), and higher methylation levels at cg14859874 (*UBAP2L*; $P = 0.002$ and $2e-4$), were causally associated with COVID-19 severity after Bonferroni correction in analyses using both COVID GWAS databases.

Discussion

To create a cutting-edge genome wide resource of *cis*- and *trans*-mQTLs, we analyzed whole genome sequences in conjunction with array-based DNA methylation data from 4126 FHS participants. Our pooled analysis identified 94,362,817 *cis*-mQTL variant-CpG pairs (9,395,367 *cis*-mQTL variants; 210,156 unique autosomal CpGs; $P < 1e-7$) and 33,572,145 *trans*-mQTL variant-CpG pairs (6,039,960 *trans*-mQTL variants; 213,606 unique autosomal CpGs; $P < 1e-14$). This comprehensive database can bridge a GWAS knowledge gap regarding mechanisms of effects of disease-associated SNPs. For example, we demonstrated enrichment of mQTL variants for disease-associated SNPs from GWAS. Using *cis*-mQTL variants, our colocalization analyses support connections between CpGs with CVD traits. MR analyses further demonstrated that *cis*-mQTLs can be used to test causal relations of CpGs to multiple phenotypes. In particular, we showed that DNA methylation at several CpGs, e.g., cg14893161 (annotated to *PM20D1*), may play an important role in relation to COVID-19 severity. Taken together, our study created a robust mQTL repository to better understand the epigenetic mechanisms underlying a wide range of diseases. A comprehensive summary data set will be posted to the National Heart, Lung, and Blood Institute's BioData Catalyst site and will be freely accessible to the scientific community.

Consistent with our previous mQTL study¹⁵ and others^{21,22}, a majority of SNP-CpG pairs are *cis*. For example, the number of *cis*-mQTL-CpG pairs was 2.8 times of that of *trans*-mQTL-CpG pairs in our pooled analysis (1.5 times using $P < 1e-14$). To the best of our knowledge, our study is the largest mQTL mapping project using WGS, including ~ 20 million SNPs and INDELS and ~ 850 thousand CpGs. Our database expands the existing literature by adding ~ 40 million novel *cis*- and ~ 30 million *trans*-mQTL-CpG pairs based on WGS rather than imputed genotypes from array-based genotyping. In addition, our database included *cis*- and *trans*-mQTLs for 180,692 unique CpGs present on the EPIC array that are not on the 450K array. Compared to the older 450K array, the EPIC array increases CpG coverage of specific genomic regions such as enhancers and non-coding regions²³. Therefore, our data will facilitate future studies that examine the potential biological function and clinical impact of DNA methylation at these genomic regions.

To showcase the application of our mQTL database, we demonstrated the enrichment of mQTL variants for disease-associated SNPs from GWAS using the GWAS Catalog⁷. For example, analysis utilizing *cis*-mQTL variants showed enrichment for SNPs associated with CVD and multiple CVD risk factors including BMI, systolic BP, triglyceride, type 2 diabetes, and coronary artery disease. Our colocalization analysis using *cis*-mQTL variants for CpGs and GWAS summary statistics of these variants for CAD identified colocalization of an LDL-associated CpG, cg05337441 (*APOB*), with coronary artery disease. A intergenic SNP rs668948, mapped to *APOB* and *TDRD15*, explained 41% of the observed colocalization. The product encoded by *APOB* is the main apolipoprotein of LDL that serves as the ligand for the LDL receptor. The atherogenic potential of apolipoprotein B-100 has been demonstrated by many studies including MR analysis²⁴⁻²⁷. Our data are consistent with the notion that DNA methylation contributes to the atherogenicity of LDL and suggest that future studies are needed to examine the exact molecular underpinnings of these observations. Also, in line with these observations, our MR analysis showed that many CVD risk factor-associated CpGs are putatively causal for CVD and CVD risk factors (Supplemental Table 8). These findings provide epigenetic insights into associations reported in GWAS. For example, we observed that cg12816198 (*IRF5*) was associated with systolic BP (MR $P = 6.3e-8$). SNP rs4728142, an intergenic variant mapped to genes *IRF5* and *KCP*, has been reported to be associated with hypertension in previous GWAS²⁸. This SNP (rs4728142) is a strong *cis*-mQTL variant for cg12816198 (*IRF5*; $P = 7e-215$) and the leading instrumental variable in the MR analysis for systolic BP (single SNP MR analysis $P = 2.7e-9$), suggesting a causal pathway whereby rs4728142 modifies DNA methylation levels at cg12816198 with downstream effects on systolic BP. Interestingly, both colocalization analysis and MR analysis showed a connection between cg27087650 (*BCL3*) and coronary artery disease through *cis*-mQTL variant rs62117206 (intronic to *BCL3*; $P = 3.6e-15$; linkage disequilibrium $R^2 = 1$ with rs4803750, another *cis*-mQTL variant of cg27087650; $P = 1.8e-14$). CpG cg27087650 is located in the gene body of *BCL3*, which encodes a protein functioning as a transcriptional co-activator through its association with NF-kappa B homodimers. Expression of *BCL3* has been linked to CVD and cancer²⁹⁻³¹. These examples provide proof of principle that integrating *cis*-mQTLs with CpGs and traits can reveal biological pathways by linking DNA methylation to a variety of diseases.

Our mQTL database can also be used to screen candidate DNA methylation sites for further consideration in experimental and interventional studies. This is exemplified by our MR analysis that revealed a putatively causal effect of COVID-19 associated CpGs on disease severity. Our COVID analysis focused on ten CpGs that were identified in a case-control study of COVID-19 severity¹⁸. Because of the retrospective design of the study¹⁸, it could not infer causal relations between DNA methylation at these CpGs and the severity of COVID-19. Our analysis highlighted three COVID-related CpGs annotated to genes *PM20D1* and *UBAP2L* that were putatively causal for COVID-19 severity; more research is needed to understand if and how these CpGs might influence outcome in patients with the COVID-19.

In parallel with our mQTL project, our research team is examining eQTLs and expression quantitative trait methylation sites (eQTM) using WGS, RNA sequencing, and DNA methylation resources obtained in FHS participants. The eQTL and eQTM resources are also freely available online via the BioData Catalyst site. These molecular resources enable users to explore how DNA methylation affects transcriptional activities and pathways leading to a wide range of disease phenotypes. These molecular resources can be used in concert to reduce bias due to reverse causality and unmeasured confounding, particularly environmental confounders^{32,33}. Nonetheless, this study has several limitations that warrant discussion. Our analysis was conducted in a group of middle-aged and older, primarily white adults; therefore, the findings in this study may not be generalizable to other populations. Nonetheless, we demonstrated that mQTLs identified in other studies^{16,17}, including those identified in a cross-ancestry analysis¹⁷, were well replicated in our database. We captured whole blood-based DNA methylation profiles, which can serve as candidate biomarkers for diseases; however, they may not reflect tissue-specific

DNA methylation levels, which may be relevant to specific diseases. Utilizing a publicly available database, Lowe et al. compared DNA methylation profiles measured by the same commercial 450K array in multiple tissues and showed a higher number of tissue-specific differentially methylated positions in blood compared to other tissues³⁴. Their study provided evidence supporting a critical role of blood in crosstalk with other tissues. A recent study by Ng et al. further advanced this notion by showing that T cells are directly involved in the pathogenesis of cardiovascular comorbidities through increased interactions with endothelial cells in individuals with non-alcoholic fatty liver disease³⁵. MR analysis was used to showcase the potential application of our mQTL database; however, MR analysis is based on assumptions that may not be testable³⁶. Also, DNA methylation can be affected by both genetic and environmental factors. We did not attempt to test effect modification by environmental factors in this study. Future studies with larger sample sizes in diverse population are needed to replicate and expand our mQTL resource.

In conclusion, we have identified millions of *cis*- and *trans*-mQTL variant CpG pairs using state-of-the-art WGS data in conjunction with high-throughput DNA methylation data. We demonstrated the utility of this vast mQTL resource by conducting GWAS signal enrichment analyses, colocalization, and MR analyses. Our mQTL repository is freely available via the BioData Catalyst site for the scientific community to study the role of DNA methylation in health and disease.

Methods

Study population. The study sample included consenting participants from the FHS Offspring, Third Generation, and Omni cohorts. In 1971, the FHS recruited the offspring of participants in the Original FHS cohort as well as the spouses of offspring to form the FHS Offspring cohort³⁷. The children of the Offspring cohort participants were recruited to the Third Generation cohort beginning in 2002³⁸. Omni cohorts were established in parallel with the Offspring and the Third Generation cohorts. In the current investigation, the study sample included 4126 FHS participants with whole blood derived DNA methylation and WGS data; 2129 participants in the Offspring cohort (exam 8, N = 869; exam 9, N = 1260), 1945 participants in the Third Generation cohort (exam 2), and 52 participants in the Omni cohort. The FHS protocols and procedures were approved by the Institutional Review Board for Human Research at Boston University Medical Center, and all participants provided written informed consent. All research was performed in accordance with relevant guidelines/regulations.

Study design. A flow chart of the study design is presented in Fig. 1. The FHS had two sets of DNA methylation data, one set included 3460 participants assayed with the Illumina BeadChip 450K (450K array; 2009 Offspring exam 8 participants and 1451 Third Generation exam 2 participants) and the second set included 1806 participants assayed with the Illumina EPIC array (EPIC array; 1260 Offspring exam 9 participants, 494 Third Generation exam 2 participants, and 52 Omni cohort participants). To maximize the sample size, as our primary analysis we conducted a pooled analysis of the two data sets. Of note, 1140 Offspring participants were included in both sets, i.e., these participants had 450K array-based methylation data from exam 8 and EPIC array-based methylation data from exam 9. In the pooled analysis, we selected the EPIC array-based data for these 1140 participants to avoid any duplication. We also conducted array-specific analysis to explore if mQTLs were replicable and to examine mQTLs that are unique to the EPIC array. We then examined the top *cis*- and *trans*-mQTLs by conducting GO pathway analysis and enrichment analysis. We tested *cis*-mQTLs for colocalization and causal association using two-sample MR analysis with CVD traits and COVID-19 severity.

DNA methylation profiling. Preparation of whole blood samples for DNA methylation profiling was conducted as detailed previously¹⁵. Briefly, DNA was obtained from whole blood buffy coat samples and prepared using bisulfite conversion before whole-genome amplification, fragmentation, array hybridization, and single-base pair extension. DNA methylation was then measured in 4170 FHS participants using the Illumina Infinium Human Methylation-450 Beadchip (450K array) in three batches (Batch 1, N = 499; Batch 2, N = 2149; and Batch 3, N = 1522). Of these, 3460 participants also had WGS data. Additionally, the Illumina MethylationEPIC 850 K BeadChip (EPIC array) was used in 1806 FHS participants with WGS. All participants were with missing methylation levels of no more than 5% of CpGs (detection $P < 0.01$) and none of them were outliers in a multi-dimensional scaling plot. The CpGs have been prefiltered so that all CpGs had $< 5\%$ missing values (detection $P < 0.01$). We calculated DNA methylation beta values (range 0 to 1) as the ratio of mean methylated and sum of methylated and unmethylated probe signal intensities. We used the DASEN method^{39,40} to normalize the methylation beta values.

Whole genome sequencing. WGS of FHS participants was performed by the Broad Institute as part of the NHLBI's TOPMed program⁴¹. Genomic DNA from whole blood samples from 2194 FHS Offspring cohort and 1582 Third Generation cohort participants was sequenced at $> \times 30$ depth of coverage⁴¹. Genetic variations were identified in a joint calling of all samples by the TOPMed Informatics Resource Center at University of Michigan. Centralized read mapping, genotype calling, and quality control were also performed at the TOPMed Informatics Research Center. This analysis used genetic variants generated from TOPMed Freeze 10a. We analyzed 20,696,115 SNPs and insertion/deletion polymorphisms (INDELs) with minor allele count (MAC) ≥ 10 . WGS data acquisition is described on the Database of Genotype and Phenotype (dbGaP) website (https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000974.v4.p3).

mQTL mapping. The mQTL mapping was conducted separately for DNA methylation data generated using the 450K array and the EPIC array. In the analysis for the 450K array data, we calculated residuals for methylation beta values obtained within each of the three methylation batches after adjusting for relevant technical

covariates. Whereas, in the analysis for the EPIC array data, we derived residuals using all available samples, also adjusting for technical covariates. The residuals from separate datasets were then combined. We then used linear regression models to perform the association analyses between the SNPs and the CpGs, adjusting for sex, age, differential leukocyte counts (estimated using the Houseman method⁴²), along with the top 15 residual methylation principal components (PCs) and five genetic PCs. We chose to adjust for 15 methylation PCs and five genetic PCs because this strategy resulted in the highest replication rate between the 450K array data and the EPIC array data. Because of relatedness among FHS study participants, linear mixed models were used in mQTL mapping to account for family structure. The primary pooled analysis examined 452,567 CpGs that are common to both arrays. The 450K array-specific analysis analyzed the same 452,567 CpGs and the EPIC array-specific analysis examined 413,524 additional CpGs (i.e., CpGs not included in the pooled analysis). We defined SNPs residing within 1 million base pairs from a CpG site as *cis*-variants and those located ≥ 1 million base pairs away from the CpG site or on a different chromosome as *trans*-variants. We considered *cis*-variants as *cis*-mQTL variants if they were associated with DNA methylation levels at the corresponding CpG site with a two-sided $P < 1e-7$, whereas we considered variants as *trans*-mQTL variants when the variant-CpG associations had a two-sided $P < 1e-14$. The P value thresholds were selected for *cis*-mQTLs based on the Bonferroni correction for the number of CpGs tested (i.e., $n = 452,567$) and for *trans*-mQTLs based on an internal discovery-validation experiment that gave the highest *trans* replication rate. We counted the number of pairs of mQTL variants (*cis* or *trans*) with their corresponding CpGs at all autosomal chromosomes. R-squared values derived from a linear regression model were used to represent heritability (h_{SNP}^2) of each *cis*- or *trans*-mQTL variant.

mQTL replication. To explore consistency between our mQTLs with published databases, we examined whole blood derived mQTLs identified in two large studies, one conducted by Bonder et al. in 3841 individuals from five Dutch biobanks¹⁶ and the other conducted by Hawe et al. in 3799 European individuals and 3195 individuals from South Asia¹⁷. Both studies analyzed SNPs based on commercial arrays with imputation. Because the number of SNPs analyzed in the two studies (~ 5 and ~ 9 million, respectively) was smaller than that tested in the present study (~ 20 million), we examined whether mQTLs identified in the two studies were also significant in our database.

Gene Ontology analysis. We tested the over-representation of GO terms based on genes annotated to the top 1000 *cis*-mQTL variants (for 1000 CpGs) with Entrez IDs identified by the pooled analysis. The default setting in the *goana* function from the R *limma* (Linear Models for Microarray and RNA-seq Data) package was used to conduct the GO analysis⁴³. GO terms (Biological Process, Cellular Component, and Molecular Function) with false positive rate (FDR) < 0.05 were reported. We repeated the same analysis for the top 1000 *trans*-mQTL variants.

GWAS enrichment analysis. We analyzed all SNPs with association $P < 5e-8$ included in the NHGRI-EBI GWAS Catalog (<https://www.ebi.ac.uk/gwas/>)⁷. We identified 243,587 entries for 2960 GWAS traits. In this analysis, we examined all mQTL variants with unique RSIDs in *cis* or *trans* at $P < 1e-7$ or $P < 1e-14$, respectively. Fisher's exact test was used to perform the enrichment analysis for each trait, and traits with FDR < 0.05 were reported.

Colocalization analysis. We conducted colocalization analysis using the R *HyPrColoc* package, a highly efficient deterministic Bayesian algorithm based on GWAS summary statistics⁴⁴. We reported the posterior probability of full colocalization (PPFC). Default prior configuration parameters (*prior.l* = $1e-4$ and *prior.c* = 0.02) and threshold of 0.7 for PPFC were used. We extracted *cis*-mQTL variants ($P < 1e-7$) derived from the present pooled analysis for 1258 CpGs associated with CVD risk factors in the EWAS catalog ($P < 1e-6$) including BMI, waist circumference, fasting glucose, systolic blood pressure (systolic BP), diastolic blood pressure (diastolic BP), high-density lipoprotein cholesterol (HDL), low-density lipoprotein cholesterol (LDL), and triglyceride⁶. We examined the colocalization of these CVD risk factor-associated CpGs with CVD-related traits including BMI, BP, lipid concentrations, type 2 diabetes, and coronary artery disease. Summary statistics for associations between *cis*-mQTL variants and GWAS SNPs for CVD-related traits were obtained from published GWAS databases^{27,45-49}.

Mendelian randomization analysis. To showcase the potential use of the mQTL resource in causal inference analyses, we conducted MR analyses to infer causal associations of the CpGs with the abovementioned CVD-related traits and COVID-19 severity. In the MR analysis for CVD-related traits, we used the same *cis*-mQTL variants for the 1258 CVD risk factors. COVID-19-associated CpGs were obtained from a recently published EWAS of COVID-19 severity¹⁸. We performed MR analyses using a two-sample MR approach⁵⁰. We used independent *cis*-mQTL variants with pair-wise linkage disequilibrium (LD) $r^2 < 0.1$ as instrumental variables (IVs). Using the *TwoSampleMR* R package⁵¹, we performed the primary analysis using the inverse variance weighted (IVW) method and sensitivity analysis using the MR-Egger method. We tested for potential horizontal pleiotropy by examining the MR-Egger intercept P value. The effect sizes and standard errors for IV-CpG associations were obtained from the pooled mQTL analysis. The effect sizes and standard errors for associations between IVs and CVD-related traits were obtained from the published large GWAS described above^{27,45,46,48,49}. We obtained effect sizes and standard errors from two GWAS for COVID-19 severity conducted by the COVID-19 Host Genetics Initiative²⁰ and the Genetics of Mortality in Critical Care (GenOMICC) study¹⁹. The COVID-19 Host Genetics Initiative included 8779 cases (death or hospitalization requiring respiratory support due to COVID-19) and 1,001,875 population controls and the GenOMICC study included 7491 cases (confirmed

COVID-19 requiring continuous cardiorespiratory monitoring in intensive care units) and 48,400 population controls.

Data availability

The datasets analyzed in the present study are available at the dbGAP repository phs000007.v32.p13 (https://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000007.v30.p11).

Received: 4 May 2022; Accepted: 10 November 2022

Published online: 15 November 2022

References

1. Anway, M. D., Cupp, A. S., Uzumcu, M. & Skinner, M. K. Epigenetic transgenerational actions of endocrine disruptors and male fertility. *Science* **308**, 1466–1469 (2005).
2. Richards, E. J. Inherited epigenetic variation—revisiting soft inheritance. *Nat. Rev. Genet.* **7**, 395–401 (2006).
3. Jones, P. A. & Takai, D. The role of DNA methylation in mammalian epigenetics. *Science* **293**, 1068–1070 (2001).
4. Colicino, E. *et al.* Blood DNA methylation sites predict death risk in a longitudinal study of 12, 300 individuals. *Aging (Albany NY)* **12**, 14092–14124 (2020).
5. Agha, G. *et al.* Blood leukocyte DNA methylation predicts risk of future myocardial infarction and coronary heart disease. *Circulation* **140**, 645–657 (2019).
6. Battram, T. *et al.* The EWAS Catalog: A database of epigenome-wide association studies. *Wellcom Open Res.* **7**, 41–53 (2022).
7. 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).
8. Wen, X., Pique-Regi, R. & Luca, F. Integrating molecular QTL data into genome-wide genetic association analysis: Probabilistic assessment of enrichment and colocalization. *PLoS Genet.* **13**, e1006646 (2017).
9. Jamieson, E. *et al.* Smoking, DNA methylation, and lung function: a Mendelian randomization analysis to investigate causal pathways. *Am. J. Hum. Genet.* **106**, 315–326 (2020).
10. Yu, F. *et al.* Mendelian randomization identifies CpG methylation sites with mediation effects for genetic influences on BMD in peripheral blood monocytes. *Front. Genet.* **11**, 60 (2020).
11. Taylor, D. L. *et al.* Integrative analysis of gene expression, DNA methylation, physiological traits, and genetic variation in human skeletal muscle. *Proc. Natl. Acad. Sci. USA* **116**, 10883–10888 (2019).
12. Battram, T. *et al.* Appraising the causal relevance of DNA methylation for risk of lung cancer. *Int. J. Epidemiol.* **48**, 1493–1504 (2019).
13. Richardson, T. G. *et al.* Systematic Mendelian randomization framework elucidates hundreds of CpG sites which may mediate the influence of genetic variants on disease. *Hum. Mol. Genet.* **27**, 3293–3304 (2018).
14. Mendelson, M. M. *et al.* Association of body mass index with DNA methylation and gene expression in blood cells and relations to cardiometabolic disease: a mendelian randomization approach. *PLoS Med.* **14**, e1002215 (2017).
15. Huan, T. *et al.* Genome-wide identification of DNA methylation QTLs in whole blood highlights pathways for cardiovascular disease. *Nat. Commun.* **10**, 4267 (2019).
16. Bonder, M. J. *et al.* Disease variants alter transcription factor levels and methylation of their binding sites. *Nat. Genet.* **49**, 131–138 (2017).
17. Hawe, J. S. *et al.* Genetic variation influencing DNA methylation provides insights into molecular mechanisms regulating genomic function. *Nat. Genet.* **54**, 18–29 (2022).
18. Castro de Moura, M. *et al.* Epigenome-wide association study of COVID-19 severity with respiratory failure. *EBioMedicine* **66**, 103339 (2021).
19. Kousathanas, A. *et al.* Whole genome sequencing reveals host factors underlying critical Covid-19. *Nature* **607**, 97–103 (2022).
20. Initiative, C.-H.G. The COVID-19 Host Genetics Initiative, a global initiative to elucidate the role of host genetic factors in susceptibility and severity of the SARS-CoV-2 virus pandemic. *Eur. J. Hum. Genet.* **28**, 715–718 (2020).
21. McRae, A. F. *et al.* Identification of 55,000 replicated DNA methylation QTL. *Sci. Rep.* **8**, 17605 (2018).
22. Hannon, E. *et al.* Leveraging DNA-methylation quantitative-trait loci to characterize the relationship between methylomic variation, gene expression, and complex traits. *Am. J. Hum. Genet.* **103**, 654–665 (2018).
23. McEwen, L. M. *et al.* Systematic evaluation of DNA methylation age estimation with common preprocessing methods and the Infinium MethylationEPIC BeadChip array. *Clin. Epigenetics* **10**, 123 (2018).
24. Marston, N. A. *et al.* Association of apolipoprotein B-containing lipoproteins and risk of myocardial infarction in individuals with and without atherosclerosis: distinguishing between particle concentration, type, and content. *JAMA Cardiol.* **7**, 250–256 (2021).
25. Johannesen, C. D. L., Mortensen, M. B., Langsted, A. & Nordestgaard, B. G. Apolipoprotein B and Non-HDL cholesterol better reflect residual risk than LDL cholesterol in statin-treated patients. *J. Am. Coll. Cardiol.* **77**, 1439–1450 (2021).
26. Ference, B. A. *et al.* Association of triglyceride-lowering LPL variants and LDL-C-lowering LDLR variants with risk of coronary heart disease. *JAMA* **321**, 364–373 (2019).
27. Richardson, T. G. *et al.* Evaluating the relationship between circulating lipoprotein lipids and apolipoproteins with risk of coronary heart disease: A multivariable Mendelian randomisation analysis. *PLoS Med.* **17**, e1003062 (2020).
28. Surendran, P. *et al.* Trans-ancestry meta-analyses identify rare and common variants associated with blood pressure and hypertension. *Nat. Genet.* **48**, 1151–1161 (2016).
29. Prasongsukarn, K., Dechkhajorn, W., Benjathummarak, S. & Maneerat, Y. TRPM2, PDLIM5, BCL3, CD14, GBA genes as feasible markers for premature coronary heart disease risk. *Front. Genet.* **12**, 598296 (2021).
30. Miao, L. *et al.* BCL3-PVRL2-TOMM40 SNPs, gene-gene and gene-environment interactions on dyslipidemia. *Sci. Rep.* **8**, 6189 (2018).
31. Legge, D. N. *et al.* The role of B-cell lymphoma-3 (BCL-3) in enabling the hallmarks of cancer: implications for the treatment of colorectal carcinogenesis. *Carcinogenesis* **41**, 249–256 (2020).
32. Smith, G. D. & Ebrahim, S. “Mendelian randomization”: Can genetic epidemiology contribute to understanding environmental determinants of disease?. *Int. J. Epidemiol.* **32**, 1–22 (2003).
33. Richardson, T. G., Harrison, S., Hemani, G. & Davey Smith, G. An atlas of polygenic risk score associations to highlight putative causal relationships across the human phenome. *Elife* **8**, e43657 (2019).
34. Lowe, R., Slodkowitz, G., Goldman, N. & Rakyen, V. K. The human blood DNA methylome displays a highly distinctive profile compared with other somatic tissues. *Epigenetics* **10**, 274–281 (2015).
35. Ng, C. Y. *et al.* Endothelial-immune crosstalk contributes to vasculopathy in nonalcoholic fatty liver disease. *EMBO Rep.* **23**, e54271 (2022).
36. Davies, N. M., Holmes, M. V. & Davey Smith, G. Reading Mendelian randomisation studies: a guide, glossary, and checklist for clinicians. *BMJ* **362**, k601 (2018).

37. Kannel, W. B., Feinleib, M., McNamara, P. M., Garrison, R. J. & Castelli, W. P. An investigation of coronary heart disease in families. The Framingham offspring study. *Am J Epidemiol* **110**, 281–90 (1979).
38. Splansky, G. L. *et al.* The Third Generation Cohort of the National Heart, Lung, and Blood Institute's Framingham Heart Study: design, recruitment, and initial examination. *Am. J. Epidemiol.* **165**, 1328–1335 (2007).
39. Pidsley, R. *et al.* A data-driven approach to preprocessing Illumina 450K methylation array data. *BMC Genomics* **14**, 293 (2013).
40. Gorrie-Stone, T. J. *et al.* Bigmelon: tools for analysing large DNA methylation datasets. *Bioinformatics* **35**, 981–986 (2019).
41. Sarnowski, C. *et al.* Whole genome sequence analyses of brain imaging measures in the Framingham Study. *Neurology* **90**, e188–e196 (2018).
42. Houseman, E. A. *et al.* DNA methylation arrays as surrogate measures of cell mixture distribution. *BMC Bioinform.* **13**, 86 (2012).
43. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* **43**, e47 (2015).
44. Foley, C. N. *et al.* A fast and efficient colocalization algorithm for identifying shared genetic risk factors across multiple traits. *Nat. Commun.* **12**, 764 (2021).
45. Yengo, L. *et al.* Meta-analysis of genome-wide association studies for height and body mass index in approximately 700000 individuals of European ancestry. *Hum. Mol. Genet.* **27**, 3641–3649 (2018).
46. Evangelou, E. *et al.* Genetic analysis of over 1 million people identifies 535 new loci associated with blood pressure traits. *Nat. Genet.* **50**, 1412–1425 (2018).
47. Willer, C. J. *et al.* Discovery and refinement of loci associated with lipid levels. *Nat. Genet.* **45**, 1274–1283 (2013).
48. Xue, A. *et al.* Genome-wide association analyses identify 143 risk variants and putative regulatory mechanisms for type 2 diabetes. *Nat. Commun.* **9**, 2941 (2018).
49. van der Harst, P. & Verweij, N. Identification of 64 novel genetic loci provides an expanded view on the genetic architecture of coronary artery disease. *Circ. Res.* **122**, 433–443 (2018).
50. Hartwig, F. P., Davies, N. M., Hemani, G. & Davey Smith, G. Two-sample Mendelian randomization: avoiding the downsides of a powerful, widely applicable but potentially fallible technique. *Int. J. Epidemiol.* **45**, 1717–1726 (2016).
51. Hemani, G. *et al.* The MR-Base platform supports systematic causal inference across the human genome. *Elife* **7**, e34408 (2018).

Acknowledgements

The views and opinions 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.

Author contributions

D.L. and J.M. designed research, interpreted results, and had primary responsibility for final content; R.J. and J.M. conducted the analyses; J.M. and R.J. wrote the manuscript; H.B. contributed to the literature review used in the paper; C.L., A.K., H.S., H.B., B.T., M.S., P.J.M., C.Y.D., P.C., N.L.H., A.N.P., M.F., N.S., R.S.V., and T.H. critically reviewed the manuscript; and all authors read and approved the final manuscript.

Funding

The Framingham Heart Study was supported by NIH contracts N01-HC-25195, HHSN268201500001I, and 75N92019D00031. DNA methylation assays were supported in part by the Division of Intramural Research (D. Levy, Principal Investigator) and an NIH Director's Challenge Award (D. Levy, Principal Investigator). The analytical component of this project was funded by the NHLBI Division of Intramural Research (D. Levy, Principal Investigator). Whole genome sequencing for the TransOmics in Precision Medicine (TOPMed) program was supported by the NHLBI. Core support including centralized genomic read mapping and genotype calling, along with variant quality metrics and filtering were provided by the TOPMed Informatics Research Center (3R01HL-117626-02S1; contract HHSN268201800002I). Core support including phenotype harmonization, data management, sample identity QC, and general program coordination were provided by the TOPMed Data Coordinating Center (R01HL-120393; U01HL-120393; contract HHSN268201800001I). J. Ma is supported by NIH grants, K22HL135075 and R01AA028263.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary Information The online version contains supplementary material available at <https://doi.org/10.1038/s41598-022-24100-0>.

Correspondence and requests for materials should be addressed to J.M. or D.L.

Reprints and permissions information is available at www.nature.com/reprints.

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



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

© The Author(s) 2022