# Supplemental Material

# Epigenomic signatures of tobacco smoking

Roby Joehanes[1,2*], Allan C. Just[3*], Riccardo E. Marioni[4,5,6*], Luke C. Pilling[7*], Lindsay M. Reynolds[8*], Pooja R. Mandaviya[9,10*], Weihua Guan[11*], Tao Xu[12*], Cathy E. Elks[13*], Stella Aslibekyan[14*], Hortensia Moreno-Macias[15,16*], Jennifer A. Smith[17*], Jennifer A Brody[18*], Radhika Dhingra[19*], Paul Yousefi[20], James S. Pankow[21], Sonja Kunze[12], Sonia Shah[6,22], Allan F. McRae[6,22], Kurt Lohman[23], Jin Sha[14], Devin M. Absher[24], Luigi Ferrucci[25], Wei Zhao[17], Ellen W. Demerath[20], Jan Bressler[26], Megan L. Grove[26], Tianxiao Huan[2], Chunyu Liu[2], Michael M. Mendelson[2,27], Chen Yao[2], Douglas P. Kiel[1], Annette Peters[12], Rui Wang-Sattler[12], Peter M. Visscher[4,6,22], Naomi R. Wray[6], John M. Starr[4,28], Jingzhong Ding[29], Carlos J. Rodriguez[8], Nicholas J. Wareham[13], Marguerite R. Irvin[14], Degui Zhi[30], Myrto Barrdahl[31], Paolo Vineis[32,33], Srikant Ambatipudi[16], André G. Uitterlinden[9], Albert Hofman[34], Joel Schwartz[35], Elena Colicino[35], Lifang Hou[36], Pantel S. Vokonas[37], Dena G. Hernandez[38], Andrew B. Singleton[38], Stefania Bandinelli[39], Stephen T. Turner[40], Erin B. Ware[17,41], Alicia K. Smith[42], Torsten Klengel[43,44], Elisabeth B. Binder[43,45], Bruce M. Psaty[18,47], Kent D. Taylor[47,48,49], Sina A. Gharib[50], Brenton R. Swenson[18], Liming Liang[51], Dawn L. DeMeo[52], George T. O'Connor[53], Zdenko Herceg[16], Kerry J. Ressler[42,44,54], Karen N. Conneely[55#], Nona Sotoodehnia[17#], Sharon L. R. Kardia[17#], David Melzer[7#], Andrea A. Baccarelli[35,56#], Joyce B. J. van Meurs[9#], Isabelle Romieu[16#], Donna K. Arnett[14#], Ken K. Ong[13#], Yongmei Liu[8#], Melanie Waldenberger[12#], Ian J. Deary[4,57#], Myriam Fornage[26,58#], Daniel Levy[2#], Stephanie J. London[59#]

*These authors contributed equally as first authors
#These authors contributed equally as senior authors
Correspondence is to be sent to Stephanie J. London (london2@niehs.nih.gov; T: 919-541-5772; F: 301-480-3290)

Affiliations
[1]Institute for Aging Research, Hebrew SeniorLife, Department of Medicine Beth Israel Deaconess Medical Center and Harvard Medical School, Boston, MA, USA
[2]National Heart, Lung, and Blood Institute, National Institutes of Health, Bethesda, MD, USA
[3]Department of Preventive Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA
[4]Centre for Cognitive Ageing and Cognitive Epidemiology, University of Edinburgh, UK
[5]Centre for Genomic and Experimental Medicine, Institute of Genetics and Molecular Medicine, University of Edinburgh, UK
[6]Queensland Brain Institute, University of Queensland, Australia
[7]Epidemiology and Public Health Group, Institute of Biomedical and Clinical Science, University of Exeter Medical School, Exeter, UK
[8]Department of Epidemiology & Prevention, Public Health Sciences, Wake Forest School of Medicine, Winston-Salem, NC, USA
[9]Department of Internal Medicine, Erasmus University Medical Center, Rotterdam, The Netherlands
[10]Department of Clinical Chemistry, Erasmus University Medical Center, Rotterdam, The Netherlands
[11]Division of Biostatistics, Schoold of Public Health, Univerisity of Minnesota, Minneapolis, MN, USA
[12]Research Unit of Molecular Epidemiology, Institute of Epidemiology II, Munich, Germany
[13]MRC Epidemiology Unit, Institute of Metabolic Science, University of Cambridge, Cambridge, UK
[14]Department of Epidemiology, University of Alabama at Birmingham, Birmingham, AL, USA
[15]Autonomous Metropolitan University-Iztapalapa, Mexico City, Mexico
[16]International Agency for Research on Cancer (IARC)
[17]Department of Epidemiology, School of Public Health, University of Michigan, Ann Arbor, MI, USA
[18]Cardiovascular Health Research Unit, Department of Medicine, Epidemiology, and Health Services, University of Washington, Seattle, WA, USA
[19]Department of Environmental Health, Rollins School of Public Health, Emory University, Atlanta, GA, USA
[20]School of Public Health, University of California, Berkeley, CA, USA
[21]Division of Epidemiology & Community Health, School of Public Health, University of Minnesota, Minneapolis, USA
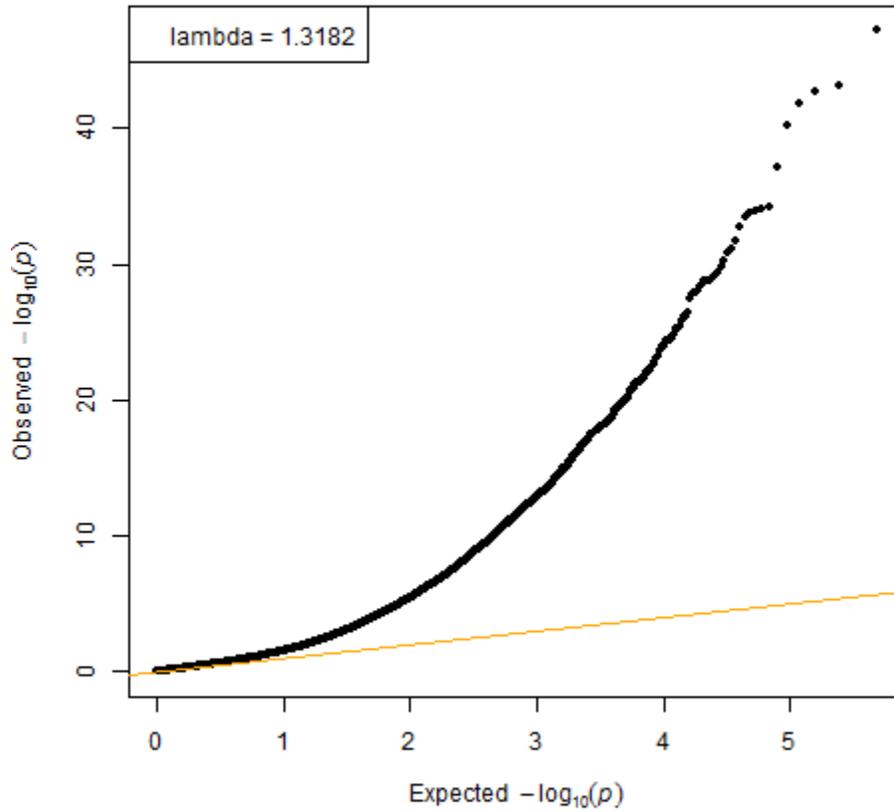[22]University of Queensland Diamantina Institute, Translational Research Institute, The University of Queensland, Brisbane 4072, QLD, Australia
[23]Department of Biostatistical Sciences, Division of Public Health Sciences, Wake Forest School of Medicine, Winston-Salem, NC, USA
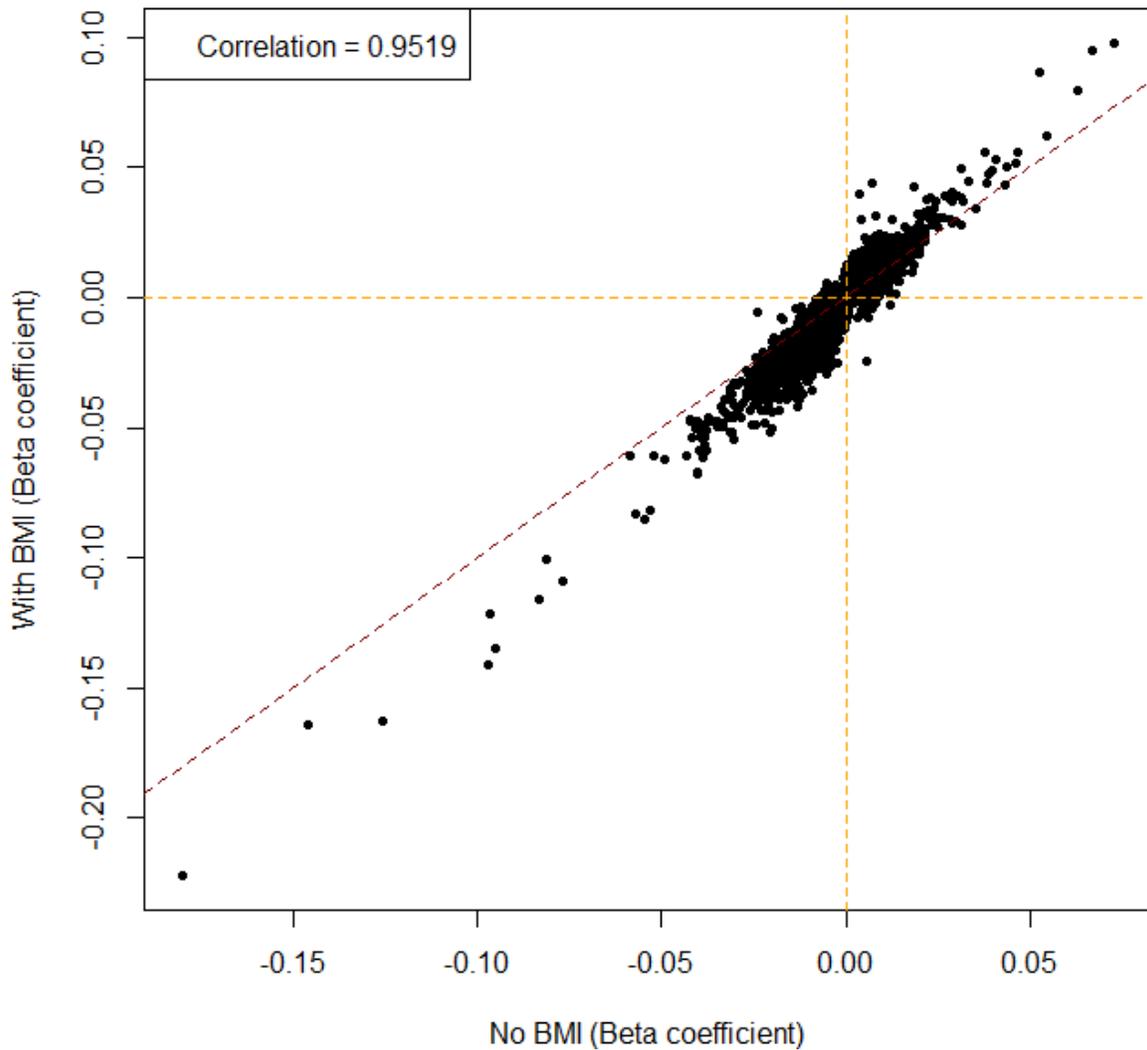
[24]HudsonAlpha Institute for Biotechnology, Huntsville, AL, USA

[25]Clinical Research Branch, National Institute on Aging, Baltimore, MD, USA

[26]Human Genetics Center, School of Public Health, The University of Texas Health Science Center at Houston, Houston, USA

[27]Children's Hospital, Boston, MA, USA

[28]Alzheimer Scotland Dementia Research Centre, University of Edinburgh, UK

[29]Department of Internal Medicine, Wake Forest School of Medicine, Winston-Salem, NC, USA

[30]Department of Biostatistics, University of Alabama at Birmingham, Birmingham, AL, USA

[31]Division of Cancer Epidemiology, German Cancer Research Center (DKFZ) Heidelberg, Im Neuenheimer Feld 581, D-69120, Heidelberg, Germany

[32]MRC/PHE Centre for Environment and Health, School of Public Health, Imperial College London, UK

[33]HuGeF Foundation, Torino, Italy

[34]Department of Epidemiology, Erasmus University Medical Center, Rotterdam, The Netherlands

[35]Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA, USA

[36]Department of Preventive Medicine and the Robert H. Lurie Comprehensive Cancer Center, Feinberg School of Medicine, Northwestern University, Chicago, IL, USA

[37]VA Normative Aging Study, VA Boston Healthcare System and Department of Medicine, Boston University School of Medicine, Boston, MA, USA

[38]Laboratory of Neurogenetics, National Institute on Aging, National Institutes of Health, Bethesda, MD, USA

[39]Geriatric Unit, Azienda Sanitaria di Firenze, Florence, Italy

[40]Division of Nephrology and Hypertension, Mayo Clinic, Rochester, MN, USA

[41]Research Center for Group Dynamics, Institute for Social Research, University of Michigan, Ann Arbor, MI

[42]Psychiatry and Behavioral Sciences, Emory University School of Medicine, Atlanta, GA, USA

[43]Department of Translational Research in Psychiatry, Max-Planck Institute of Psychiatry, Munich, Germany

[44]Division of Depression & Anxiety Disorders, McLean Hospital, Belmont, MA, USA

[45]Department of Psychiatry and Behavioral Sciences, Emory University School of Medicine, Atlanta, GA, USA

[47]Group Health Research Institute, Group Health Cooperative, Seattle, WA, USA

[47]Institute for Translational Genomics and Population Sciences, Los Angeles BioMedical Research Institute at Harbor-UCLA Medical Center, Torrance, CA, USA.

[48]Division of Genomic Outcomes, Department of Pediatrics, Harbor-UCLA Medical Center, Torrance, CA, USA.

[49]Departments of Pediatrics, Medicine, and Human Genetics, UCLA, Los Angeles, CA, USA

[50]Center for Lung Biology, Division of Pulmonary & Critical Care Medicine, Department of Medicine, University of Washington, Seattle, WA

[51]Harvard School of Public Health, Boston, MA, USA

[52]Channing Division of Network Medicine, Brigham and Women's Hospital, Harvard Medical School, Boston, MA, USA

[53]Boston University School of Medicine, Boston, MA, USA

[54]Department of Psychiatry, Harvard Medical School, Boston MA

[55]Department of Human Genetics, Emory University School of Medicine, Atlanta, GA, USA

[56]Department of Environmental Health, Harvard T.H. Chan School of Public Health, Boston, MA, USA

[57]Department of Psychology, University of Edinburgh, UK

[58]Institute of Molecular Medicine, The University of Texas Health Science Center McGovern Medical School, Houston, TX, USA

[59]Epidemiology Branch, National Institute of Environmental Health Sciences, National Institutes of Health, Department of Health and Human Services, Research Triangle Park, NC, USA

# Supplementary Materials

Supplementary Figure 1. Quantile-quantile (QQ) plot for CpG site association with respect to current versus never smoker
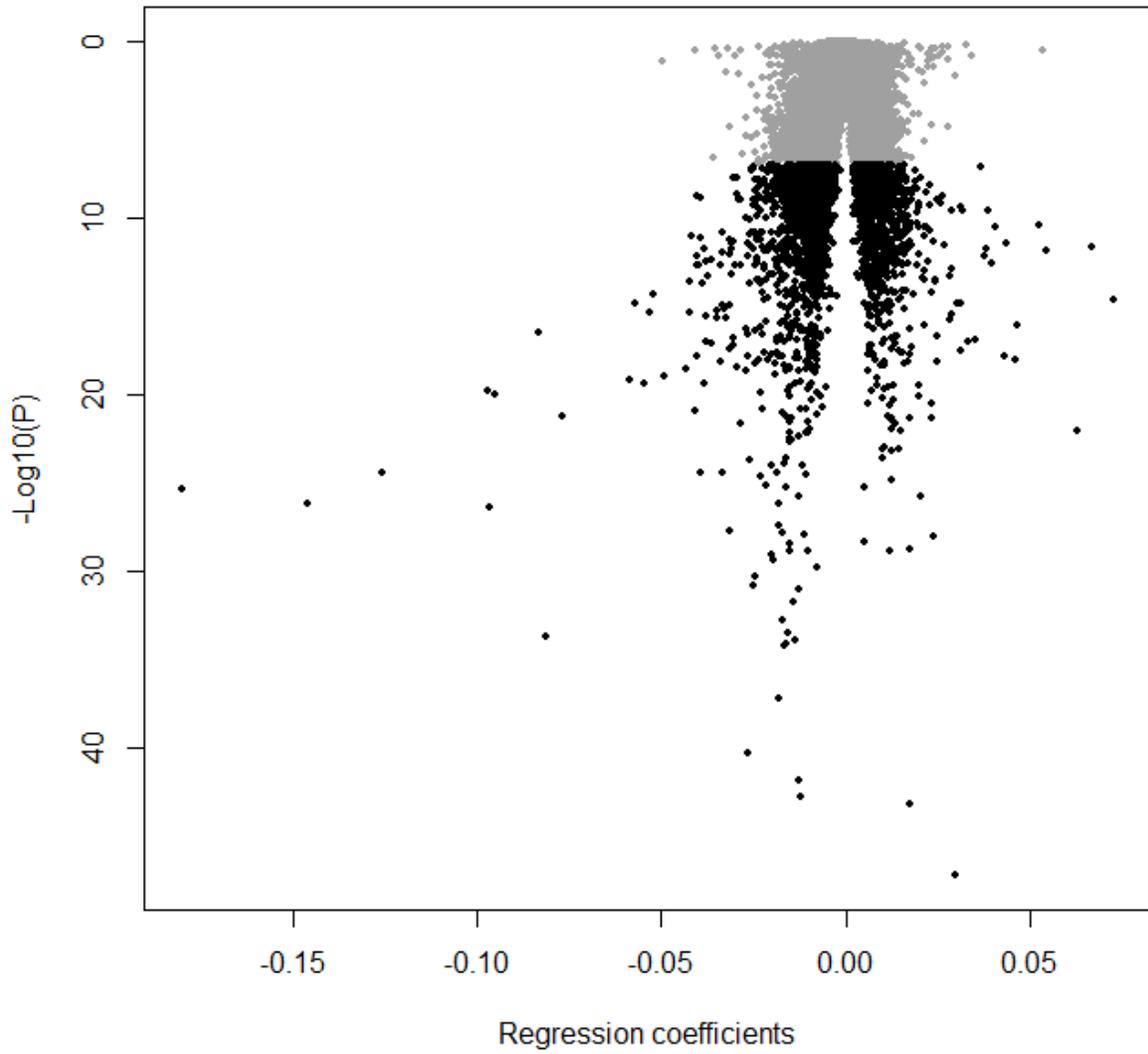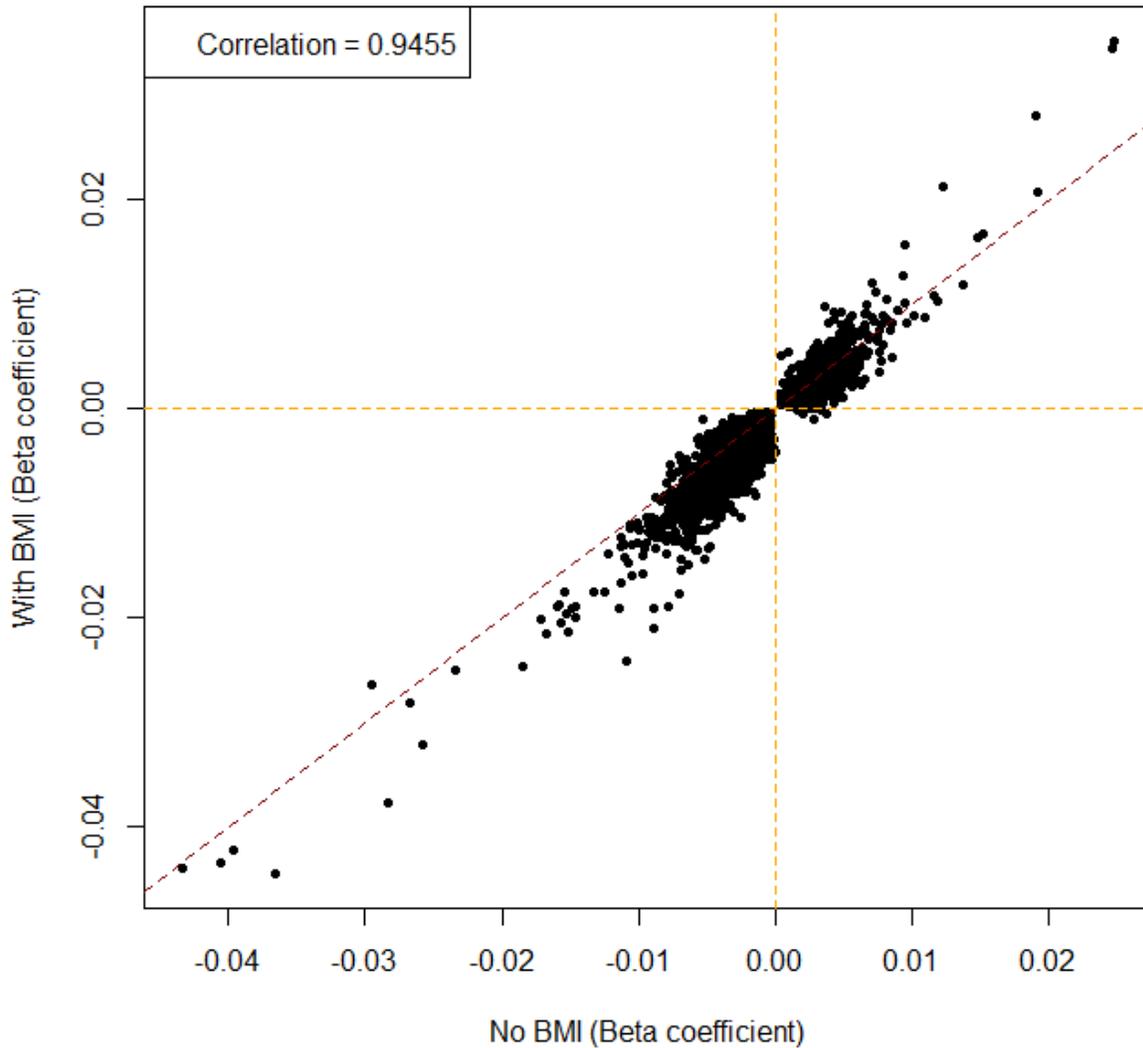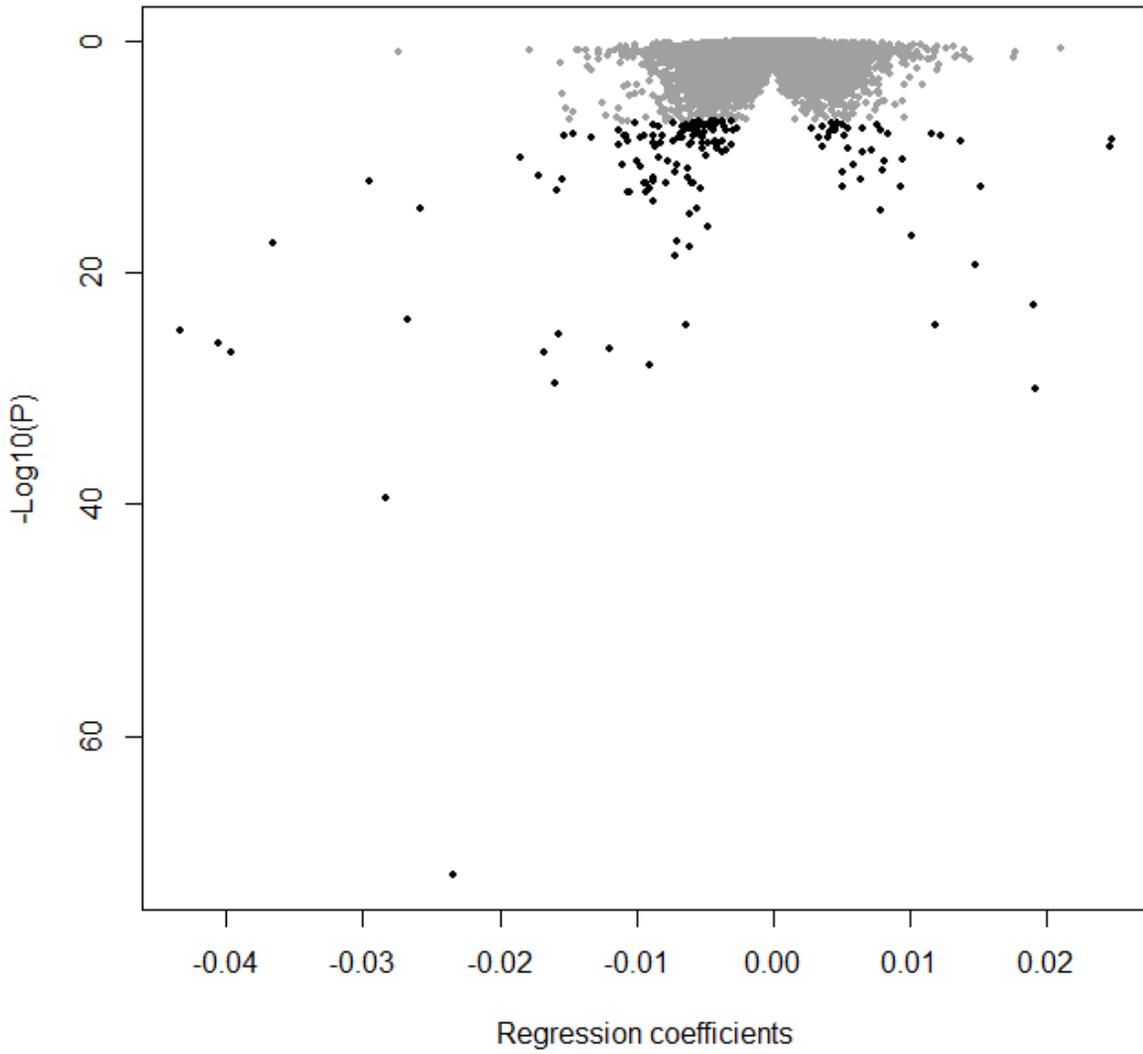
Supplementary Figure 2. Comparison of regression coefficients (beta) of significant 22,473 CpGs between two models (with and without BMI) in relation to current versus never smokers. The X axis indicates beta coefficients without body mass index (BMI). The y axis indicates beta coefficients with BMI added into the model. The CpGs are selected if they are significant (having false discovery rate < 0.05) in at least one of the models. The beta coefficients between the two models are correlated at 0.9519 level.
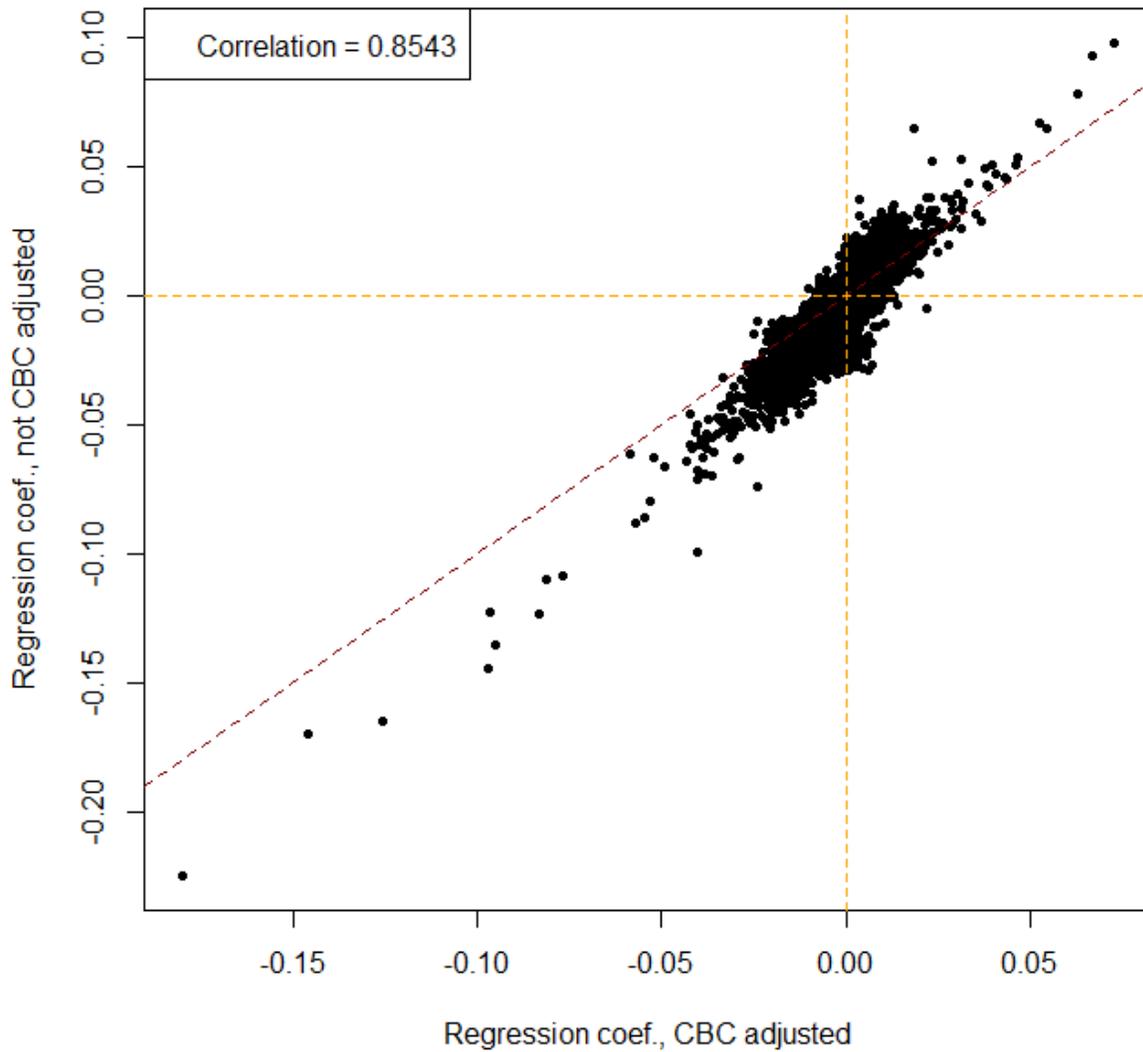
Supplementary Figure 3. Volcano plot for CpG site association with respect to current versus never smoker



**Current vs. Never Smokers**

Supplementary Figure 4. Quantile-quantile (QQ) plot for CpG site association with respect to former versus never smoker

Supplementary Figure 5. Comparison of regression coefficients (beta) of significant 2,998 CpGs between two models (with and without BMI) in relation to former versus never smokers. The X axis indicates beta coefficients without body mass index (BMI). The y axis indicates beta coefficients with BMI added into the model. The CpGs are selected if they are significant (having false discovery rate < 0.05) in at least one of the models. The beta coefficients between the two models are correlated at 0.9455 level.

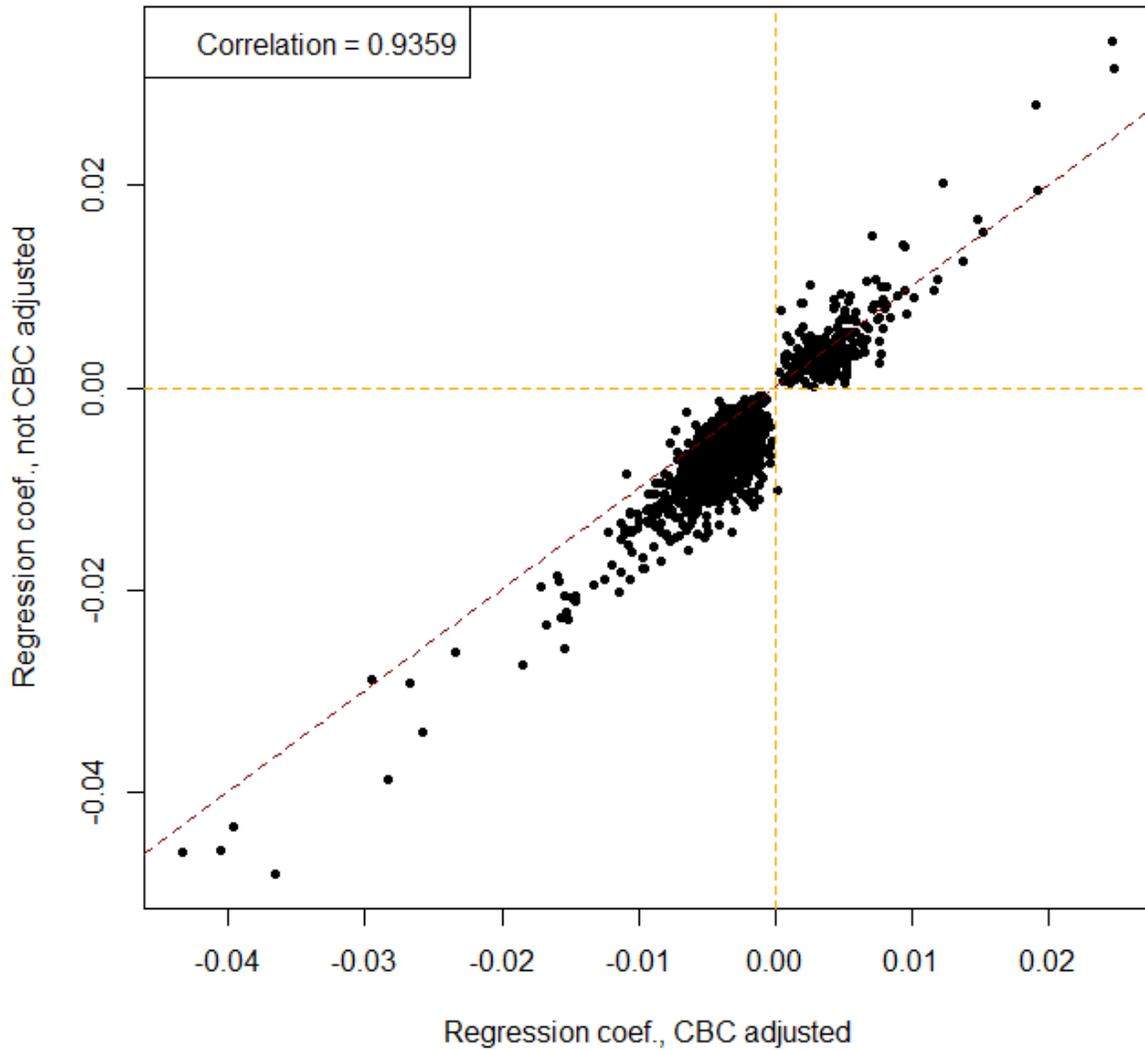Supplementary Figure 6. Volcano plot for CpG site association with respect to former versus never smoker
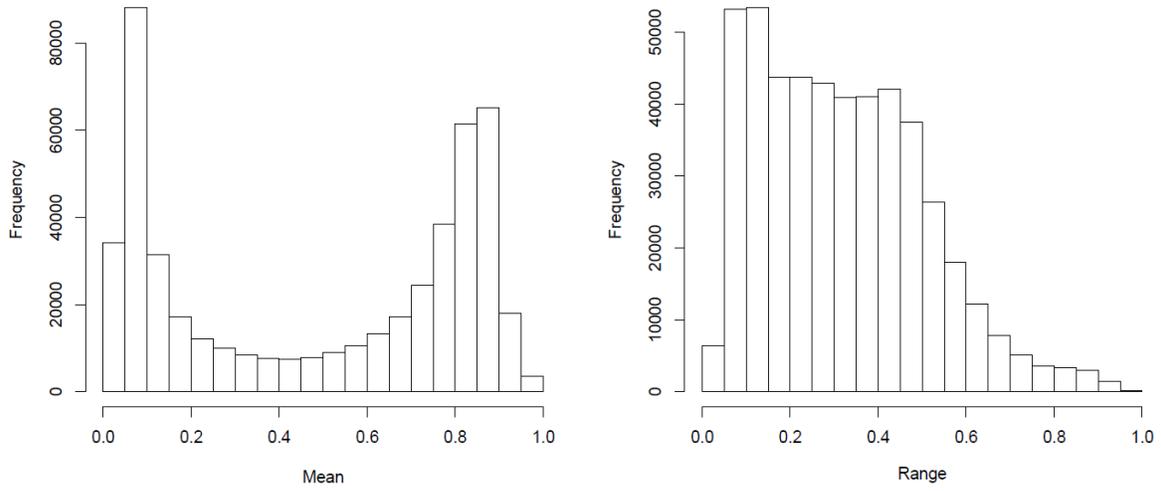
**Former vs. Never Smokers**

Supplementary Figure 7. Comparison of regression coefficients (beta) of significant 26,693 CpGs between two models (with and without blood cell type adjustment) in relation to current versus never smokers. The X axis indicates beta coefficients with complete blood count (CBC) adjustment. The y axis indicates beta coefficients with without CBC adjustment. The CpGs are selected if they are significant (having false discovery rate < 0.05) in at least one of the models. The beta coefficients between the two models are correlated at 0.8543 level.
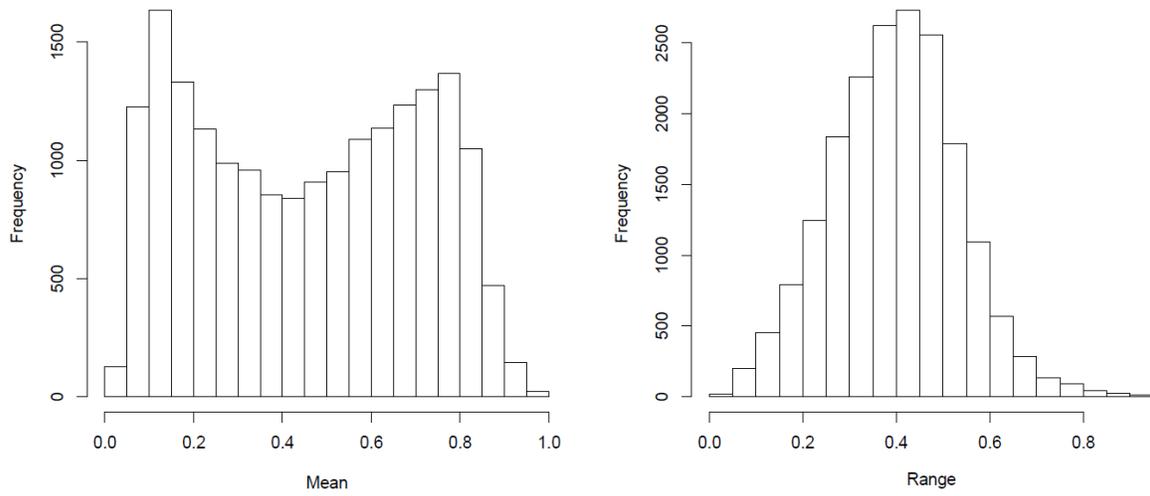
Supplementary Figure 8. Comparison of regression coefficients (beta) of significant 1,137 CpGs between two models (with and without blood cell type adjustment) in relation to former versus never smokers. The X axis indicates beta coefficients with complete blood count (CBC) adjustment. The y axis indicates beta coefficients with without CBC adjustment. The CpGs are selected if they are significant (having false discovery rate < 0.05) in at least one of the models. The beta coefficients between the two models are correlated at 0.9359 level.

Supplementary Figure 9. Histogram plot of mean and range of all 485,381 CpG sites in Framingham Heart Study (FHS) cohort, in methylation proportion (β) scale.

Supplementary Figure 10. Histogram plot of mean and range of 18,760 CpG sites significant in current *vs.* never smokers in Framingham Heart Study (FHS) cohort, in methylation proportion (β) scale.

## Supplementary Tables

*See separate Excel spreadsheet for all supplemental tables.*

Supplementary Table 1. Detailed participant characteristics by cohort.

Supplementary Table 2. Statistically significant CpGs in relation to current *vs.* never smoking status at false discovery rate (FDR)<0.05.

Supplementary Table 3. Statistically significant CpGs in relation to current *vs.* never smoking status that exhibit dose-response relationship (via pack years) at FDR<0.05.

Supplementary Table 4. Gene Ontology pathways of genes whose CpGs are statistically significant in relation to current *vs.* never smoking status.

Supplementary Table 5. Statistically significant CpGs in relation to former *vs.* never smoking status at false discovery rate (FDR)<0.05.

Supplementary Table 6. Statistically significant CpGs in relation to former *vs.* never smoking status that exhibit dose-response relationship (via pack years) at FDR<0.05.

Supplementary Table 7. Gene Ontology pathways of genes whose CpGs are statistically significant in relation to former *vs.* never smoking status.

Supplementary Table 8. List of genome-wide association study (GWAS) phenotypes or diseases for which statistically significant CpGs in relation to current *vs.* never smoking status are enriched.

Supplementary Table 9. List of genome-wide association study (GWAS) phenotypes or diseases for which statistically significant CpGs in relation to former *vs.* never smoking status are enriched.

Supplementary Table 10. List of genes that are GWAS-associated with CVD-related diseases or risk factors that are differentially methylated in relation to current *vs.* never smoking status.

Supplementary Table 11. Enrichment results for genomic features for which differentially methylated CpGs in relation to smoking status are enriched.

Supplementary Table 12. Differentially methylated CpGs in relation to current *vs.* never smoking status that exhibit transcriptomic control in *cis*.

Supplementary Table 13. Gene Ontology pathways of genes whose transcripts are associated in *cis* with the differentially methylated CpGs in relation to current *vs.* never smoking status.

Supplementary Table 14. Comparison of differentially methylated CpGs in relation to current *vs.* never smoking status between cohorts of African Ancestry (AA) and European Ancestry (EA).

Supplementary Table 15. Comparison of differentially methylated CpGs in relation to former *vs.* never smoking status between cohorts of African Ancestry (AA) and European Ancestry (EA).

Supplementary Table 16. Comparison of differentially methylated CpGs in relation to current *vs.* never smoking status between cohorts of whole blood and leukocyte samples.

Supplementary Table 17. Comparison of differentially methylated CpGs in relation to former *vs.* never smoking status between cohorts of whole blood and leukocyte samples.

Supplementary Table 18. Correlation among regression coefficients of CpGs showing significant associations on smoking status across different cell types. Numbers above the diagonal line are for current *vs.* never smoker status, while those below are for former *vs.* never smoker status.

## Supplementary Methods

### Cohort overview

This study of Cohorts for Heart and Aging Research in Genetic Epidemiology (CHARGE) comprises a total of 15,907 participants from 16 cohorts ARIC, FHS Offspring, KORA F4, GOLDN, LBC 1921, LBC 1936, NAS, Rotterdam, Inchianti, GTP, CHS European Ancestry (EA), CHS African Ancestry (AA), GENOA, EPIC Norfolk, EPIC, and MESA. The study was approved by institutional review committees for each cohort and all participants provided written informed consent for genetic research.

## Framingham Heart Study (FHS)

### Description
The <u>F</u>ramingham <u>H</u>eart <u>S</u>tudy (FHS) is a population-based study that began in 1948. The offspring cohort, consisting of 5,124 participants of European ancestry, was recruited in 1971[1]. Excluding control samples, DNA methylation was measured on 2,792 offspring cohort participants who attended the eighth examination cycle (2005-2008). Of these, 2,648 had both measurements on methylation and smoking status (274 current, 1,538 former, and 836 never smokers). All participants provided written informed consent for genetic research.

### DNA methylation sample, measurement, normalization, and quality control

Buffy coat fractions from peripheral whole blood samples were collected from 2,792 offspring cohort participants. Genomic DNA was extracted using the Puregene DNA extraction kit (Qiagen, Venlo, Netherlands) which subsequently bisulfite-converted using the EZ DNA Methylation kit (Zymo Research, Irvine, CA). The samples underwent whole genome amplification, fragmentation, array hybridization, single-base pair extension, and then assayed in two laboratories using the Infinium HumanMethylation 450 BeadChip, which contains 485,512 CpG sites in all. The first laboratory assayed 576 samples, while the second laboratory 2,270 samples.

Raw methylated and total probe intensities were extracted using the Illumina Genome Studio methylation module. Preprocessing of the methylated signal ($M$) and unmethylated signal ($U$) was conducted using DASEN of wateRmelon[2] version 3.0.2, an R package. The methylation beta ($\beta$) values were defined as $\beta = M/(M+U)$. We excluded low quality CpG sites (with detection p-value > 0.01). We excluded samples showing deviation from the first two principal components (PC1 and PC2), deviation from sex clusters (*i.e.*, male-labeled samples that cluster into female-sample cluster or vice versa), and deviation (>3*SD) from 5,997 SNPs showing the strongest cis methylation quantitative trait locus (mQTL).

### Smoking phenotype
Smoking phenotype is as described in the main paper.

### Analysis
In the first stage, we analyzed the data with two linear mixed effects models, as described in the main paper. For technical covariates, we included chip ID, row, and column effects as random effects, and PC1 and PC2 as fixed effects. The former factors are to account for technical artefact, the latter account for the inter-laboratory differences. As FHS is a cohort-based study, familial relationship was also included in the model. Thereby, we used pedigreemm package[3], instead of lme4. We also performed pack-year analysis, cessation analysis, and methylation by expression (MxE) analysis as described in the main paper.

## Genetics of Lipid Lowering Drugs and Diet Network (GOLDN)

### Description

The GOLDN family study, recruited ~1300 Caucasian men and women with at least two siblings and three generational pedigrees from the participants of the National Heart, Lung, and Blood Institute Family Heart Study in two genetically homogenous centers in Minneapolis, MN and Salt Lake City, UT. The trial aimed to identify genetic factors that mediated response to lipid-raising (*i.e.*, postprandial lipemia challenge) or lipid-lowering (fenofibrate therapy) among metabolically healthy individuals. Participants were asked to discontinue the use of lipid-lowering agents for at least 4 weeks, to fast for at least 8 hours, and to abstain from alcohol and smoking for at least 24 hours prior to study visits. The study protocol was approved by the Institutional Review Boards at the University of Minnesota, University of Utah, Tufts University/New England Medical Center and the University of Alabama at Birmingham, and written informed consent was obtained from all participants[4].

### DNA methylation sample, measurement, normalization, and quality control
*Epigenetic Phenotyping*
Details of the sample isolation are described in previous publications[5,6] and are as follows. CD4+ T-cells were isolated from frozen buffy coat samples using positive selection by antigen-specific magnetic beads (Invitrogen, Carlsbad, CA). DNA was isolated from the CD4+ T-cells using DNeasy kits (Qiagen, Venlo, Netherlands) (2). We used the Infinium Human Methylation 450 array (Illumina, San Diego, CA) to quantify genome-wide DNA methylation[5]. Prior to the standard manufacturer protocol steps of amplification, hybridization, and imaging steps, we treated 500ng of each DNA sample with sodium bisulfite (Zymo Research, Irvine, CA).  We used IlluminaGenomeStudio software to estimate β scores, defined as the proportion of total signal from the methylation-specific probe or color channel, and detection p-values, defined as the probability that the total intensity for a given probe falls within the background signal intensity. β scores with an associated detection p-value greater than 0.01 were removed, as were samples with more than 1.5% missing data points across ~470,000 autosomal CpGs.  Additionally, any CpG probes where more than 10% of samples failed to yield adequate intensity were removed[5]. Filtered β scores were normalized using the ComBat package for R software[7]. Normalization was performed on random subsets of 10,000 CpGs per run, where each array of 12 samples was used as a "batch."  Separate normalization of probes from the Infinium I and II chemistries was performed and subsequently the β scores

for Infinium II probes were adjusted using the equation derived from fitting a second order polynomial to the observed methylation values across all pairs of probes located <50bp apart (within-chemistry correlations > 0.99), where one probe was Infinium I and one was Infinium II.  Finally, any CpGs where the probe sequence mapped either to a location that did not match the annotation file, or to more than one locus were eliminated. Such markers were identified by re-aligning all probes (with unconverted Cs) to the human reference genome (2).  After quality control, we had data for 461,281 CpGs.  Principal components based on the beta scores of all autosomal CpGs passing QC were generated using the *prcomp* function in R (V 2.12.1).

*Genotyping*
A hybrid data set of 2,543,887 single nucleotide polymorphisms (SNPs), of which 484,029 were typed using the Affymetrix Genome-Wide Human 6.0 Array (Affymetrix, Santa Clara, CA) and the rest were imputed using MACH software (Version 1.0.16, Ann Arbor, MI) with Human Genome Build 36 as a reference. Prior to imputation, SNPs were excluded if they were monomorphic, had a call rate of less than 96%, exhibited Mendelian errors, had a minor allele frequency of <1%, or failed the Hardy-Weinberg equilibrium (HWE) test at the P-value threshold of less than $10^{-6}$.

## Smoking phenotype
Data for smoking variables in this study were collected based on self-reported information.  Smoking variables included current smoking status (smoke now, yes/no); number of pack years smoked (number of packs smoked per day x number of years smoked); ever smoked (current, past, or never).

## Analysis
Linear mixed effect models were used for analyses:
β = Smoking phenotype + Sex + Age + center + PCs (to account for T-cell purity) + family (random effect).

## Acknowledgements

## Rotterdam Study (RS)

### Description
The Rotterdam Study (RS) is a large prospective, population-based cohort study aimed at assessing the occurrence of and risk factors for chronic (cardiovascular, endocrine, hepatic, neurological, ophthalmic, psychiatric, dermatological, oncological, and respiratory) diseases in the elderly[8]. The study comprises 14,926 subjects in total, living in the well-defined Ommoord district in the city of Rotterdam in the Netherlands. In 1989, the first cohort, Rotterdam Study-I (RS-I) comprised of 7,983 subjects with age 55 years or above. In 2000, the second cohort, Rotterdam Study-II (RS-II) was included with 3,011 subjects who had reached an age of 45 years since 1989. In 2006, the third cohort, Rotterdam Study-III (RS-III) was further included with 3,932 subjects with age 45 years and above. Each participant gave an informed consent and the study was approved by the medical ethics committee of the Erasmus University Medical Center, Rotterdam, the Netherlands.


### DNA methylation sample, measurement, normalization, and quality control
At the Genetic Laboratory (Department of Internal Medicine, Erasmus University Medical Center, Rotterdam, the Netherlands), the DNA methylation dataset was generated for a subset of 747 individuals of RS-III at baseline. Genomic DNA was extracted from whole peripheral blood by standardized salting out methods. This was followed by a bisulfide conversion using the Zymo EZ-96 DNA-methylation kit (Zymo Research, Irvine, CA, USA). The genome for each sample was then amplified, fragmented and hybridized to the Infinium Illumina Human Methylation 450k arrays according to the manufacturer's protocol.

The quality control for samples was performed using the Methylation Module of the GenomeStudio software (http://www.illumina.com/applications/microarrays/microarray-software/genomestudio.html). Data was extracted into beta values from raw IDAT files. We excluded samples based on the detection p-value criteria >99% (n=7), poor bisulfite conversion based on control dashboard check (n=5) and failed chromosome X & Y clustering (n=4).

The data preprocessing was additionally performed using an R programming pipeline which is based on the pipeline developed by Tost & Toulemat[9], which includes additional parameters and options to preprocess and normalize methylation data directly from idat files. The beta values were extracted using the R package methylumi. We excluded probes which had a detection p-value >0.01 in >95% of samples. 11648 probes at X and Y chromosomes were excluded to avoid gender bias. This filtering criteria left 731 samples and 463,456 probes. The raw beta values were then background corrected and normalized using the DASEN option of the WateRmelon R-package[2].

## Smoking phenotype

Smoking phenotype is as described in the main paper.

## Analysis

In the first stage, data were analyzed with linear mixed effects models, as described in the main paper. Pack-year analysis was also performed.

## Acknowledgements

## InCHIANTI

### Description

The InCHIANTI population is a large population-based study based in the Chianti region of Tuscany, Italy[10]. The participants are aged between 30-104 years and underwent thorough examination every three years from 1998-2000. Whole blood samples were collected using the PAXgene system in 2007[11]. Ethical approval was granted by the Instituto Nazionale Riposo e Cura Anziani institutional review board.

### DNA methylation sample, measurement, normalization, and quality control

Genomic DNA was extracted from buffy coat samples using an AutoGen Flex and quantified on a Nanodrop1000 spectrophotometer prior to bisulfite conversion. Genomic DNA was bisulfite converted using Zymo EZ-96 DNA Methylation Kit per the manufacturer's protocol. CpG methylation status of 485,577 CpG sites was determined using Illumina Infinium HumanMethylation450 BeadChip per manufacturer's protocol and as previously described[12]. Initial data analysis was performed using GenomeStudio 2011.1 (Model M Version 1.9.0 Illumina, Inc. CA). Threshold call rate for inclusion of samples was 95%. Quality control of sample handling included comparison of clinically reported sex versus sex of the same samples determined by analysis of methylation levels of CpG sites on the X chromosome. Beta values were extracted for sites on the X chromosome. Subject mean methylation versus subject mean intensity levels were plotted in R V2.11.1. Based on methylation levels for chromosome X loci, these data split into two primary groups. Calls generated by this method were then compared with sample information reported by InChianti Study. Samples not matching between clinical reported sex and methylation data were excluded from analyses.

Quantile normalization of the methylation arrays was carried out using package "wateRmelon" for the R statistical computing language[2]. The DASEN method was applied, which performs the quantile normalization separately on M and U (methylated/un-methylated) values, and also separates the type 1 and type 2 Infinium probes. This minimises the technical variance between the arrays, whilst taking into consideration the different technologies present on the arrays. Methylation data was available for 506 InCHIANTI participants following quality control and data cleaning.

### Analysis
Linear mixed effects models were applied to each 450k array probe in turn with the following cofactors included: fixed effects: age, sex, total white blood cell counts, lymphocyte, monocyte, eosinophil and basophil proportions, platelet counts: included as random effects: sentrix ID, sentrix position, and array batch. Analyses were performed on current vs. never, former vs. never, pack-years smoked, and years since quitting (cessation).

## Cooperative health research in the Region of Augsburg (KORA)

### Description
Cooperative health research in the Region of Augsburg (KORA) is a population-based cohort study conducted in the region of Augsburg, Southern Germany[13,14]. The study has been conducted according to the principles expressed in the Declaration of Helsinki. Written informed consent has been given by each participant. The study was reviewed and approved by the local ethical committee (Bayerische Landesärztekammer). The baseline survey 4 (KORA S4) consists of 4,261 individuals (aged 25-74 years) examined between 1999 and 2001. During the years of 2006 to 2008, 3,080 participants took part in the follow-up survey 4 (KORA F4). Phenotypic data were retrieved from self-reports and medical records.

### DNA methylation sample, measurement, normalization, and quality control
In KORA F4, the analysis was performed using whole blood DNA of fasting participants (n=1776). Genomic DNA (1 μg) was bisulfite converted using the EZ-96 DNA Methylation Kit (Zymo Research, Orange, CA, USA) according to the manufacturer's procedure, with the alternative incubation conditions recommended when using the Illumina Infinium Methylation Assay. Genome-wide DNA methylation was assessed using the Illumina HumanMethylation450 BeadChip, following the Illumina Infinium HD Methylation protocol. This consisted of a whole genome amplification step using 4μl of each bisulfite converted sample, followed by enzymatic fragmentation and application of the samples to BeadChips (Illumina). The arrays were fluorescently stained and scanned with the Illumina HiScan SQ scanner. Raw methylation data were extracted with Illumina Genome Studio (version 2011.1) with methylation module (version 1.9.0). The percentage of methylation at a given cytosine is reported as a beta-value. Low-confidence probes, which has less than three functional beads or has a detection p-value larger than 0.01, were excluded. Sites representing or being located in a 50 bp proximity to SNPs with a minor allele frequency of at least 5% were also excluded from the data set. β-mixture quantile normalization[15] was applied to the DNA methylation data using the R package wateRmelon[2], version 1.0.3.8. KORA F4 samples were processed on 20/7 96-well plates in 9/4 batches, a plate effect representing 4.8% of the total variance of the methylation level was observed. Additionally, plate was included as a random effect in the analyses. Detailed quality control process was described the previous publication[16].

**Smoking phenotype**

**The smoking phenotype was defined based on self-reports, see main text for details.**
**Analysis**

The data was analyzed using two linear mixed effect models as described in the main text. In KORA F4, the blood count comprises the fractions of CD4+ T-cells, CD8+ T-cells, NK cells, and monocyte estimated using the Houseman et al. method. The chip number, the row and column number of the samples on the plates were included as the technical covariates with random effect in the model. No significant population stratification was found in the KORA F4 data, familial relationship was not adjusted in the model.

## Grady Trauma Project (GTP)

### Description

The Grady Trauma Project (GTP) is a population-based, prospective study of demographic characteristics, trauma exposure, and prevalence of post-traumatic stress disorder and major depressive disorder in an urban, predominantly African-American population[17]. Subjects were recruited prospectively from the waiting rooms of primary care and obstetrics-gynecology clinics of Grady Memorial Hospital in Atlanta, GA. Exclusion criteria included mental retardation, active psychosis, or the inability to give informed consent. Written and verbal informed consent was obtained for all participating subjects. All procedures in this study were approved by the Institutional Review Boards of Emory University School of Medicine and Grady Memorial Hospital. Since its inception in 2005, over 5000 subjects have been interviewed for the study.

### DNA methylation sample, measurement, normalization, and quality control

We extracted DNA from whole blood at the Max Planck Institute in Munich for 425 GTP participants using the Gentra Puregene Kit (Qiagen); for this study, we focus on 286 participants who are African American and have complete information for the smoking phenotype (described below). Genomic DNA was then bisulfite converted using the Zymo EZ-96 DNA Methylation Kit (Zymo Research). We assessed DNA methylation at >480,000 CpG sites using Illumina HumanMethylation450 BeadChip arrays, with hybridization and processing performed according to the instructions of the manufacturer. For each CpG site and individual, we collected two data points: M (the total methylated signal), and U (the total unmethylated signal). We set to missing data points with 1) a detection p-value greater than 0.001 or 2) a combined signal less than 25% of the total median signal and less than both the median unmethylated and median methylated signal. We removed individual samples from analysis if they were outliers in a hierarchical clustering analysis or had 1) a mean total signal less than half of the median of the overall mean signal or 2000 arbitrary units, or 2) a missingness rate above 5%. Similarly, we removed from analysis CpG sites with a missingness rate above 10%. We then computed β-values for each individual at each CpG site as the total methylated signal divided by the total signal: $\beta = M/(M+U)$. For quantile normalized data, the M and U signals were quantile normalized together prior to computation of β-values.

### Smoking phenotype

Smoking information was collected from GTP participants using an adapted KMSK questionnaire tool. This tool (originally described in Kellogg et al. 2003[18]) records, using a numerical scale, the current frequency of smoking, the duration of time that this frequency has been maintained and the amount of cigarettes smoked during this period. The adapted tool used in GTP recorded this information for both the 30 days prior and the time period where participant smoking was greatest for 425 individuals. Frequency (coded on a 0-5 point scale, where 5 = smoking at regular intervals most/all days; 4 = smoking at specific times of day most/all days; 3 = once a day most/all days; 2 = 20-100 times in lifetime; 1 = less than 20 times in lifetime; 0

= never smoked) for both time periods (hereafter referred to as '30-day' and 'maximum') was used to create a variable describing whether the individual is a current, former or never smoker (CFN).

The CFN scale was determined as follows:
- An individual was classified as a current smoker (N = 94), if their 30-day frequency was coded as a 3, 4 or 5 and their maximum frequency was coded as a 3, 4, 5 or missing.
- An individual was classified as a former smoker (N = 64), if 30-day frequency was coded as a 0, 1, or 2, and maximum frequency was coded as a 3, 4 or 5.
- An individual was classified as a never smoker (N = 128), if their maximum frequency was coded as a 0, 1, or 2 and their 30-day frequency was coded as a 0, 1, 2 or missing.
If an individual did not meet the above criteria (N = 45) or did not supply any smoking information (N = 61), their score on the CFN scale was recorded as missing.

### Analysis
Data were analyzed with two linear mixed effects models, as described in the main paper.

### Acknowledgements

## Lothian Birth Cohorts of 1921 and 1936 (LBC1921 and LBC1936)

### Description

The Lothian Birth Cohorts of 1921 and 1936 are two longitudinal studies of ageing[19–21]. They derive from the Scottish Mental Surveys of 1932 and 1947, respectively, when nearly all 11 year old children in Scotland completed a test of general cognitive ability[21]. Survivors living in the Lothian area of Scotland were recruited in late-life at mean age 79 for LBC1921 (n=550) and mean age 70 for LBC1936 (n=1,091). Follow-up has taken place at ages 70, 73, and 76 in LBC1936 and ages 79, 83, 87, and 90 in LBC1921. Collected data include genetic information, longitudinal epigenetic information, longitudinal brain imaging (LBC1936), and numerous blood biomarkers, anthropomorphic and lifestyle measures. Post QC, DNA methylation data were available for 920 LBC1936 participants at age 70, and for 446 LBC1921 participants at age 79.

### DNA methylation sample, measurement, normalization, and quality control

Detailed information about the collection and QC steps undertaken on the LBC methylation data have been reported previously[22]. Briefly, the Infinium HumanMethylation450 BeadChip (Illumina Inc, San Diego, CA) was used to measure DNA methylation in whole blood of consenting participants. Background correction was performed and QC was used to remove probes with a low detection rate (<95% at P < 0.01), low quality (manual inspection), low call rate (below 450,000 probes at P < 0.01), and samples with a poor match between genotypes and SNP control probes, and incorrect predicted sex. Background correction and internal normalisation were performed; the betas were modified such that the minimum was 0.001 and the maximum was 0.999.

### Smoking phenotype

Smoking was measured via self-response. Participants were asked if they were current smokers, never smokers, or former smokers.

### Analysis

Linear mixed effects models were used to analyze the data in both cohorts. Measured white blood cell counts (eosinophils, neutrophils, basophils, monocytes, and lymphocytes) were included as fixed effects along with age and sex; technical covariates (sample plate, BeadChip, position on BeadChip, and hybridisation date) were included as random effects.

### Acknowledgements

## The Multi Ethnic Study of Atherosclerosis (MESA)

### Description

The Multi-Ethnic Study of Atherosclerosis (MESA) was designed to investigate the prevalence, correlates, and progression of subclinical cardiovascular disease in a population cohort of 6,814 participants. Since its inception in 2000, five clinic visits collected extensive clinical, socio-demographic, lifestyle, behavior, laboratory, nutrition, and medication data[23]. DNA methylation and gene expression were measured in purified (CD14+) monocyte samples from the April 2010 – February 2012 examination (exam 5) of 1,264 randomly selected MESA participants from four MESA field centers (Baltimore, MD; Forsyth County, NC; New York, NY; and St. Paul, MN) as previously described[24]. The study protocol was approved by the Institutional Review Board at each site. All participants signed informed consent.

### DNA methylation sample, measurement, normalization, and quality control

As previously described[24], blood was initially collected in sodium heparin-containing Vacutainer CPT™ cell separation tubes (Becton Dickinson, Rutherford, NJ, USA) to separate peripheral blood mononuclear cells from other elements within 2 h from blood draw. Subsequently, monocytes were isolated with the anti-CD14-coated magnetic beads, using AutoMACs automated magnetic separation unit (Miltenyi Biotec, Bergisch Gladbach, Germany). Based on flow cytometry analysis of 18 specimens, monocyte samples were consistently >90% pure.   DNA and RNA were isolated from samples simultaneously using the AllPrep DNA/RNA Mini Kit (Qiagen, Inc., Hilden, Germany). DNA and RNA QC metrics included optical density measurements, using a NanoDrop spectrophotometer and evaluation of the integrity of 18s and 28s ribosomal RNA.

Illumina HumanMethylation450 BeadChips and HiScan reader were used to perform the epigenome-wide methylation analysis. Bead-level methylation data were summarized in GenomeStudio. Because a two-channel system and both Infinium I and II assays were used, normalization was performed in several steps using the lumi package. "Smooth quantile normalization" was used to adjust for color bias. Next, the data were background adjusted by subtracting the median intensity value of the negative control probes. Lastly, data were normalized across all samples by standard quantile normalization applied to the bead-type intensities and combined across Infinium I and II assays and both colors. QC measures included checks for sex and race/ethnicity mismatches, and outlier identification by multidimensional scaling plots. To estimate residual sample contamination for data analysis, we generated separate enrichment scores for neutrophils, B cells, T cells, monocytes, and natural killer cells. We implemented a Gene Set Enrichment Analysis[25] as previously described[24] to calculate the enrichment scores using the gene signature of each blood cell type from previously defined lists[26]. To remove technical error in methylation levels associated with batch effects across the multiple chips, positional effects of the sample on the chip, and residual sample contamination with non-monocyte cell types, we adjusted methylation values for chip, sample position on the chip, and estimated residual sample contamination with neutrophils, B cells, T cells, monocytes, and natural killer cells. The final methylation value for each methylation probe was computed as the beta-value, essentially the proportion of the methylated to the total intensity.

### Smoking phenotype

Smoking status was ascertained longitudinally (Exams 1-5). Current smokers reported to be current smokers at Exam 5, the time of the blood draw. Former smokers reported to be former smokers at any exam (1-5) or reported to ever smoke at least 100 cigarettes in their lifetime at exam 1. Never smokers reported never smoking at all exams.

### Analysis

Data were analyzed with two linear mixed effects models, as described in the main paper. A look-up for significant (FDR<0.001) methylation by expression (MxE) associations was performed using data previously reported in the same 1,264 CD14+ samples[24], including genes located within 1 MB of smoking-associated methylation identified in current vs. never and former vs. never analyses.

### Acknowledgements

## European Prospective Investigation into Cancer (EPIC)

### Description

The EPIC study is an on-going multi-center prospective cohort study designed to investigate the relation between nutrition and cancer occurrence. The cohort consists of 23 centers in 10 European countries (*i.e.*, Denmark, France, Germany, Greece, Italy, Netherlands, Norway, Spain, Sweden and United Kingdom). From 1992–2000, more than 500,000 individuals aged between 25 and 70 years were recruited. All participants gave written or oral informed consent. The study was approved by the International Agency for Research on Cancer (IARC) ethical review committee and by local ethical committees at the participating centers. The design of EPIC is described in detail elsewhere[27] DNA methylation was measured on 450 breast cancer cases and 450 matched controls among women using a nested case-control approach (2005-2008). Of these, 898 had both measurements on methylation and smoking status (196 current, 190 former, and 512 never smokers). All participants provided written informed consent for genetic research.

### DNA methylation sample, measurement, normalization, and quality control

DNA was isolated from white blood cells as per the standard DNA extraction procedure (Autopure LS, Qiagen). DNA methylome profiling was carried out using the Illumina Infinium HumanMethylation450 (HM450) BeadChip assay, which interrogates more than 480,000 methylation sites, essentially as described previously[28]. Briefly, 500 ng of extracted DNA was bisulfite-modified using the EZ DNA Methylation kit (Zymo Research, D5004), following the manufacturer's instructions for the HM450 BeadChip assay. The conversion was confirmed by performing PCR for *GAPDH* primers specific for modified/unmodified DNA samples. The hybridization and scanning of the arrays were performed as per the manufacturer's instructions.

Data pre-processing and analysis were performed using R (version 3.2.2) /Bioconductor packages. To avoid spurious associations, we excluded the cross-reactive probes and probes overlapping with a known single nucleotide polymorphism (SNP) with an allele frequency of ≥5% in the overall population (European ancestry[29]), leaving 423,066 probes. In any given sample, a probe with a detection *P*-value (a measure of an individual probe's performance) of greater than or equal to 0.05 was assigned missing status. If a probe was missing in more than 5% of samples, it was excluded from all samples. Thus, 1,625 probes were excluded on this basis. Finally, 421,441 probes were available for the analyses, which were corrected for probe colour bias, inter-sample quantile normalization followed by beta-mixture quantile normalization (BMIQ) to align Type I and Type II probe distributions[15]. The array annotations from FDb.InfiniumMethylation.hg19 were used to assign probes to their corresponding genes.

### Smoking phenotype

Each center participating in EPIC cohort had their own questionnaire, thus questions regarding smoking habits were slightly different. Responses were

harmonized in order to classified participants as never, former and current smoker according to the responses to their respective questionnaires with the rationale of clearly distinguish between groups. Some examples of questionnaire information are: "Do you currently smoke", "Have you ever smoked for over 3 month?"," Did you smoked 1 cigarette per day or more per day in the past?"," Have you ever smoked as much as one cigarette a day for as long as a year?"," Do you smoke cigarettes regularly? , etc. Women were classified as never, former and current smoker according their responses to their respective questionnaire with the rationale of clearly distinguish between groups.

## Analysis

First the proportion of methylation (beta values) for each CpG site were explored and then the range of beta values were checked. The vector of raw betas and the vector of normalized betas were used as the outcome in separate linear mixed models with center and pool ID as random effects. Models were adjusted for age, BMI, cancer status (case or control) and proportion of CD8 T lymphocytes, CD4 T lymphocytes, B cells, monocytes and natural killer cells. The genomic inflation factor (lambda) was calculated and the QQ-plot was generated for each model. The Benjamini & Hochberg[30] procedure was used for controlling False Discovery Rate because multiple testing. R software v.3.1.3 was used.

## The Atherosclerosis Risk in Communities (ARIC) Study

### Description

The Atherosclerosis Risk in Communities (ARIC) Study is a prospective cohort study of cardiovascular disease risk in four U.S. communities. Between 1987 and 1989, 7,082 men and 8,710 women aged 45–64 years were recruited from Forsyth County, North Carolina; Jackson, Mississippi (African Americans only); suburban Minneapolis, Minnesota; and Washington County, Maryland. The ARIC Study protocol was approved by the institutional review board of each participating institution. After written informed consent was obtained, including that for genetic studies, participants underwent a baseline clinical examination (Visit 1) and four subsequent follow-up clinical exams (Visits 2 – 5).

### DNA methylation sample, measurement, normalization, and quality control

At this time, DNA methylation data are available for African American members of the cohort from two centers (Forsyth County and Jackson). The present study comprises a cross-sectional analysis of smoking and methylation measured in samples collected at visit 2 and 3, with covariates obtained at the same visit.

Genomic DNA was extracted from peripheral blood leukocyte samples using the Gentra Puregene Blood Kit (Qiagen; Valencia, CA, USA) according to the manufacturer's instructions (www.qiagen.com). Bisulfite conversion of 1 ug genomic DNA was performed using the EZ-96 DNA Methylation Kit (Deep Well Format) (Zymo Research; Irvine, CA, USA) according to the manufacturer's instructions (www.zymoresearch.com). Bisulfite conversion efficiency was determined by PCR amplification of the converted DNA before proceeding with methylation analyses on the Illumina platform using Zymo Research's Universal Methylated Human DNA Standard and Control Primers.

Bisulfite-converted DNA was used for hybridization on the Illumina Infinium HumanMethylation450 (HM450) BeadChip, following the Illumina Infinium HD Methylation protocol (www.illumina.com). This consisted of a whole genome amplification step followed by enzymatic end-point fragmentation, precipitation and re-suspension. The re-suspended samples were hybridized to the complete set of bead-bound probes, followed by ligation and single-base extension during which a fluorescently-labeled nucleotide is incorporated, and scanned. The degree of methylation is determined for each CpG cytosine by measuring the amount of incorporated label for each probe. The intensities of the images were extracted using Illumina GenomeStudio 2011.1, Methylation module 1.9.0 software. The methylation score for each CpG was represented as a beta ($\beta$) value according to the fluorescent intensity ratio. Background subtraction was conducted with the GenomeStudio software using built-in negative control bead types on the array.

Positive and negative controls and sample replicates were included on each 96-well plate assayed. After exclusion of controls, replicates, and 22 samples that failed

bisulfite conversion, a total of 2,905 study participants had HM450 data available for further quality control analyses. We removed poor-quality samples with pass rate <99% (N=32). At the target level, we flagged poor-quality CpG sites with average detection p-value > 0.01, and calculated the percentage of samples having detection p-value > 0.01 for each autosomal and X chromosome CpG site. There were 5,174 autosomal and X chromosomal markers where >1% of samples showed detection p-value > 0.01, and these sites were excluded.

Methylation values were normalized using the Beta MIxture Quantile dilation (BMIQ) method[15].

### Smoking phenotype

Smoking phenotype is as described in the main paper.

### Analysis

Since white blood cell proportions were not directly measured in most participants in ARIC, they were imputed from the methylation data using the Houseman method. Specifically, the proportions of neutrophils, lymphocytes, monocytes, eosinophils, and basophils were estimated based on the measured differential cell counts available for a subset of ARIC participants at Visit 2 (n = 175). All association analyses were performed in R using linear mixed models with DNA methylation beta values as the outcome, as described in the main paper.

## Genetic Epidemiology Network of Arteriopathy (GENOA)

### Description

The Genetic Epidemiology Network of Arteriopathy (GENOA) study is a community-based study of hypertensive sibships that was designed to investigate the genetics of hypertension and target organ damage in African Americans from Jackson, Mississippi and non-Hispanic whites from Rochester, Minnesota[31]. In the initial phase of GENOA (Phase I: 1996-2001), all members of sibships containing ≥ 2 individuals with essential hypertension clinically diagnosed before age 60 were invited to participate, including both hypertensive and normotensive siblings. DNA methylation was measured on the peripheral blood leukocytes of 422 unrelated African American participants using stored blood samples collected during the Phase I examination. Participants were excluded if they were identified as an outlier in principal component plots generated during the methylation data cleaning process. A total of 420 African American GENOA participants were included in this analysis. All participants provided written informed consent for genetic research.

### DNA methylation sample, measurement, normalization, and quality control

Genomic DNA of 422 participants was extracted from stored peripheral blood leukocytes collected at the Phase I GENOA examination. The EZ DNA Methylation Gold Kit (Zymo Research, Irvine CA) was used for bisulfite conversion, and methylation was measured with the Illumina Infinium HumanMethylation450 BeadChip. The *minfi* R package was used to preprocess, normalize (SWAN), and calculate beta values. Principal components analysis was performed using the SWAN method to identify and exclude sample outliers (>6sd from the mean of the top 10 PCs). The proportion of each cell type was estimated using Houseman's method. Detection p-values were calculated for each sample at each CpG site, and values were set as missing when detection P-value was >0.01. CpG sites were excluded if >10% of samples had a detection P-value of >0.01. All samples had a call rate >90%.

### Smoking phenotype

Participants were categorized as being a current smoker (smoker within the past 1 year), former smoker (not having smoked in the past 1 year), or never smoker. A person was considered a never smoker if they answered "No" to the following question: "Have you ever smoked more than 100 cigarettes in your entire life?". A person was considered a former smoker if they answered "Yes" to "Have you ever smoked more than 100 cigarettes in your entire life?", answered "No" to "Do you now smoke cigarettes?", and there was greater than 1 year between their current age/date of exam and their answer to the question, "In what year or how old were you when you last quit smoking?" A person was considered a current smoker if they answered "Yes" to "Have you ever smoked more than 100 cigarettes in your entire life?" and answered "Yes" to "Do you now smoke cigarettes?". A person was also considered a current smoker if they answered "Yes" to "Have you ever smoked more than 100 cigarettes in your entire life?", answered "No" to "Do you now smoke

cigarettes?", and there was less than 1 year between their current age/date of exam and their answer to the question, "In what year or how old were you when you last quit smoking?".

## Analysis

GENOA data were analyzed with linear mixed effect models using the R software, as described in the main paper. DNA methylation beta values were used as the outcome variables.

## Acknowledgements

## Cardiovascular Health Study (CHS)

### Description

The CHS is a population-based cohort study of risk factors for coronary heart disease and stroke in adults ≥65 years conducted across four field centers[32]. The original predominantly European ancestry cohort of 5,201 persons was recruited in 1989-1990 from random samples of the Medicare eligibility lists; subsequently, an additional predominantly African-American cohort of 687 persons was enrolled for a total sample of 5,888.

DNA methylation was measured on 200 European ancestry and 200 African-American ancestry participants.  The samples were randomly selected among participants without presence of coronary heart disease, congestive heart failure, peripheral vascular disease, valvular heart disease, stroke or transient ischemic attack at study baseline or lack of available DNA at study year 5.

CHS was approved by institutional review committees at each field center and individuals in the present analysis had available DNA and gave informed consent including consent to use of genetic information for the study of cardiovascular disease.

### DNA methylation sample, measurement, normalization, and quality control

Methylation measurements were performed at the Institute for Translational Genomics and Population Sciences at the Harbor-UCLA Medical Center Institute for Translational Genomics and Population Sciences using the Infinium HumanMethylation450 BeadChip (Illumina Inc, San Diego, CA).

Quality control was performed in in the minfi R package[33–35] (version 1.12.0, http://www.bioconductor.org/packages/release/bioc/html/minfi.html). Samples with low median intensities of below 10.5 ($\log_2$) across the methylated and unmethylated channels, samples with a proportion of probes falling detection of greater than 0.5%, samples with QC probes falling greater than 3 standard deviation from the mean, sex-check mismatches, or failed concordance with prior genotyping were removed. In total, 11 samples were removed for sample QC resulting in a sample of 191 European-ancestry and 198 African-American samples.   Methylation values were normalized using the SWAN quantile normalization method[34].   Since white blood cell proportions were not directly measured in CHS they were estimated from the methylation data using the Houseman method[36].

### Smoking phenotype

Smoking phenotype is as described in the main paper.

### Analysis

All association analyses were performed in R using linear models with DNA methylation beta values as the outcome.  Analyses were stratified by race and all analyses were adjusted for age, gender, total white blood cell count, study clinic and

estimated white blood cell proportions, as well as chip number and position on the chip.

## European Prospective Investigation into Cancer and Nutrition-Norfolk (EPIC-Norfolk)

### Description

The European Prospective Investigation of Cancer (EPIC)-Norfolk study enrolled more than 25,000 community-based men and women at baseline (1993-1997), who were aged 40-79 years old and registered with a participating general practitioner in and around the city of Norwich (Norfolk, UK). The full details of the study design and follow up of participants has been reported previously[37]. Written informed consent was obtained from all participants. The study complies with the principles of the Declaration of Helsinki and ethical approval was given by the Norfolk Local Research Ethics Committee and the East Norfolk and Waveney NHS Research Governance Committee.

### DNA methylation sample, measurement, normalization, and quality control

DNA was isolated from white blood cells as per the standard DNA extraction procedure (Autopure LS, Qiagen). DNA methylome profiling was carried out using the Illumina Infinium HumanMethylation450 (HM450) BeadChip assay. 500 ng of extracted DNA was bisulfite-modified using the EZ DNA Methylation kit (Zymo Research, D5004) following the manufacturer's instructions. The minfi R package was used to preprocess, normalize (SWAN), and calculate beta values.

### Smoking phenotype

Personal medical history was assessed using the question in the Health and Lifestyle Questionnaire. Yes/no responses to the questions "Have you ever smoked as much as one cigarette a day for as long as a year?" and "Do you smoke cigarettes now?" were used to derive smoking history

### Analysis

All association analyses were performed in R using linear models with DNA methylation beta values as the outcome. Analyses were analyses were adjusted for age, gender, and estimated white blood cell proportions, as well as plate number and position.

### Acknowledgements

## Normative Aging Study (NAS)

### Description

The US Department of Veterans Affairs (VA) Normative Aging Study (NAS) is an ongoing longitudinal cohort established in 1963, which included men who were aged 21 to 80 years and free of known chronic medical conditions at entry[38]. Participants were subsequently invited to medical examinations every 3 to 5 years. At each visit, participants provided information on medical history, lifestyle, and demographic factors, and underwent a physical examination and laboratory tests. DNA samples were collected from 1999 to 2007 from the active participants and used for DNA methylation analysis.

DNA was extracted from buffy coat using the QIAamp DNA Blood Kit (QIAGEN, Valencia, CA, USA). A total of 500 ng of DNA was used to perform bisulfite conversion using the EZ-96 DNA Methylation Kit (Zymo Research, Orange, CA, USA). To limit chip and plate effects, a two-stage age-stratified algorithm was used to randomize samples and ensure similar age distributions across chips and plates; we randomized 12 samples - which were sampled across all the age quartiles - to each chip, then chips were randomized to plates (each housing eight chips). Quality control analysis was performed to remove samples where >1% of probes had a detection P value >0.05 and probes where >1% of passing samples had a detection P value >0.05. The passing samples were preprocessed using out-of-band background correction[39], dye bias adjustment, and probe type adjustment using the Beta MIxture Quantile dilation (BMIQ) method[15].

### Smoking phenotype

At each in-person examination visit, participants completed a questionnaire that included their smoking status that was classified as in the main paper.

### Analysis

Data were analyzed with linear mixed effects models, as in the main paper with main models adjusted with a fixed effect for age and an indicator for sentrix column (position on chip) and random effects for sentrix row (position on chip) and chip number. As the NAS does not include females there was no adjustment for sex.

## References

1. Feinleib M, Kannel WB, Garrison RJ, McNamara PM, Castelli WP. The Framingham Offspring Study. Design and preliminary data. *Prev Med*. 1975;4:518–525.

2. Pidsley R, Y Wong CC, Volta M, Lunnon K, Mill J, Schalkwyk LC. A data-driven approach to preprocessing Illumina 450K methylation array data. *BMC Genomics*. 2013;14:293.

3. Vazquez AI, Bates DM, Rosa GJM, Gianola D, Weigel KA. Technical note: an R package for fitting generalized linear mixed models in animal breeding. *J Anim Sci*. 2010;88:497–504.

4. Genetics of Lipid Lowering Drugs and Diet Network (GOLDN) study [Internet]. Available from: http://www.biostat.wustl.edu/goldn/

5. Absher DM, Li X, Waite LL, Gibson A, Roberts K, Edberg J, Chatham WW, Kimberly RP. Genome-wide DNA methylation analysis of systemic lupus erythematosus reveals persistent hypomethylation of interferon genes and compositional changes to CD4+ T-cell populations. *PLoS Genet*. 2013;9:e1003678.

6. Hidalgo B, Irvin MR, Sha J, Zhi D, Aslibekyan S, Absher D, Tiwari HK, Kabagambe EK, Ordovas JM, Arnett DK. Epigenome-Wide Association Study of Fasting Measures of Glucose, Insulin, and HOMA-IR in the Genetics of Lipid Lowering Drugs and Diet Network Study. *Diabetes*. 2014;63:801–807.

7. Johnson WE, Li C, Rabinovic A. Adjusting batch effects in microarray expression data using empirical Bayes methods. *Biostatistics*. 2006;8:118–127.

8. Hofman A, Brusselle GGO, Darwish Murad S, van Duijn CM, Franco OH, Goedegebure A, Ikram MA, Klaver CCW, Nijsten TEC, Peeters RP, Stricker BHC, Tiemeier HW, Uitterlinden AG, Vernooij MW. The Rotterdam Study: 2016 objectives and design update. *Eur J Epidemiol*. 2015;30:661–708.

9. Touleimat N, Tost J. Complete pipeline for Infinium(®) Human Methylation 450K BeadChip data processing using subset quantile normalization for accurate DNA methylation estimation. *Epigenomics*. 2012;4:325–341.

10. Ferrucci L, Bandinelli S, Benvenuti E, Di Iorio A, Macchi C, Harris TB, Guralnik JM. Subsystems contributing to the decline in ability to walk: bridging the gap between epidemiology and geriatric practice in the InCHIANTI study. *J Am Geriatr Soc*. 2000;48:1618–1625.

11. Debey-Pascher S, Eggle D, Schultze JL. RNA stabilization of peripheral blood and profiling by bead chip analysis. *Methods Mol Biol*. 2009;496:175–210.

12. Gibbs JR, van der Brug MP, Hernandez DG, Traynor BJ, Nalls MA, Lai S-L, Arepalli S, Dillman A, Rafferty IP, Troncoso J, Johnson R, Zielke HR, Ferrucci L, Longo DL, Cookson MR, Singleton AB. Abundant quantitative trait loci exist for DNA methylation and gene expression in human brain. *PLoS Genetics*. 2010;6:e1000952.

13. Holle R, Happich M, Löwel H, Wichmann HE, MONICA/KORA Study Group. KORA--a research platform for population based health research. *Gesundheitswesen*. 2005;67 Suppl 1:S19-25.

14. Wichmann H-E, Gieger C, Illig T, MONICA/KORA Study Group. KORA-gen--resource for population genetics, controls and a broad spectrum of disease phenotypes. *Gesundheitswesen*. 2005;67 Suppl 1:S26-30.

15. Teschendorff AE, Marabita F, Lechner M, Bartlett T, Tegner J, Gomez-Cabrero D, Beck S. A beta-mixture quantile normalization method for correcting probe design bias in Illumina Infinium 450 k DNA methylation data. *Bioinformatics*. 2013;29:189–196.

16. Pfeiffer L, Wahl S, Pilling LC, Reischl E, Sandling JK, Kunze S, Holdt LM, Kretschmer A, Schramm K, Adamski J, Klopp N, Illig T, Hedman ÅK, Roden M, Hernandez DG, Singleton AB, Thasler WE, Grallert H, Gieger C, Herder C, Teupser D, Meisinger C, Spector TD, Kronenberg F, Prokisch H, Melzer D, Peters A, Deloukas P, Ferrucci L, Waldenberger M. DNA methylation of lipid-related genes affects blood lipid levels. *Circ Cardiovasc Genet*. 2015;8:334–342.

17. Gillespie CF, Bradley B, Mercer K, Smith AK, Conneely K, Gapen M, Weiss T, Schwartz AC, Cubells JF, Ressler KJ. Trauma exposure and stress-related disorders in inner city primary care patients. *Gen Hosp Psychiatry*. 2009;31:505–514.

18. Kellogg SH, McHugh PF, Bell K, Schluger JH, Schluger RP, LaForge KS, Ho A, Kreek MJ. The Kreek-McHugh-Schluger-Kellogg scale: a new, rapid method for quantifying substance abuse and its possible applications. *Drug Alcohol Depend*. 2003;69:137–150.

19. Deary IJ, Gow AJ, Pattie A, Starr JM. Cohort profile: the Lothian Birth Cohorts of 1921 and 1936. *Int J Epidemiol*. 2012;41:1576–1584.

20. Deary IJ, Gow AJ, Taylor MD, Corley J, Brett C, Wilson V, Campbell H, Whalley LJ, Visscher PM, Porteous DJ, Starr JM. The Lothian Birth Cohort 1936: a study to examine influences on cognitive ageing from age 11 to age 70 and beyond. *BMC Geriatr*. 2007;7:28.

21. Deary IJ, Whiteman MC, Starr JM, Whalley LJ, Fox HC. The impact of childhood intelligence on later life: following up the Scottish mental surveys of 1932 and 1947. *J Pers Soc Psychol*. 2004;86:130–147.

22. Shah S, McRae AF, Marioni RE, Harris SE, Gibson J, Henders AK, Redmond P, Cox SR, Pattie A, Corley J, Murphy L, Martin NG, Montgomery GW, Starr JM, Wray NR, Deary IJ, Visscher PM. Genetic and environmental exposures constrain epigenetic drift over the human life course. *Genome Res.* 2014;24:1725–1733.

23. Bild DE, Bluemke DA, Burke GL, Detrano R, Diez Roux AV, Folsom AR, Greenland P, Jacob DR, Kronmal R, Liu K, Nelson JC, O'Leary D, Saad MF, Shea S, Szklo M, Tracy RP. Multi-Ethnic Study of Atherosclerosis: objectives and design. *Am J Epidemiol*. 2002;156:871–881.

24. Liu Y, Ding J, Reynolds LM, Lohman K, Register TC, De La Fuente A, Howard TD, Hawkins GA, Cui W, Morris J, Smith SG, Barr RG, Kaufman JD, Burke GL, Post W, Shea S, McCall CE, Siscovick D, Jacobs DR, Tracy RP, Herrington DM, Hoeschele I. Methylomics of gene expression in human monocytes. *Hum Mol Genet*. 2013;22:5065–5074.

25. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR, Lander ES, Mesirov JP. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci USA*. 2005;102:15545–15550.

26. Abbas AR, Baldwin D, Ma Y, Ouyang W, Gurney A, Martin F, Fong S, van Lookeren Campagne M, Godowski P, Williams PM, Chan AC, Clark HF. Immune response in silico (IRIS): immune-specific genes identified from a compendium of microarray expression data. *Genes Immun*. 2005;6:319–331.

27. Riboli E, Hunt KJ, Slimani N, Ferrari P, Norat T, Fahey M, Charrondière UR, Hémon B, Casagrande C, Vignat J, Overvad K, Tjønneland A, Clavel-Chapelon F, Thiébaut A, Wahrendorf J, Boeing H, Trichopoulos D, Trichopoulou A, Vineis P, Palli D, Bueno-De-Mesquita HB, Peeters PHM, Lund E, Engeset D, González CA, Barricarte A, Berglund G, Hallmans G, Day NE, Key TJ, Kaaks R, Saracci R. European Prospective Investigation into Cancer and Nutrition (EPIC): study populations and data collection. *Public Health Nutr*. 2002;5:1113–1124.

28. Hernandez-Vargas H, Castelino J, Silver MJ, Dominguez-Salas P, Cros M-P, Durand G, Calvez-Kelm FL, Prentice AM, Wild CP, Moore SE, Hennig BJ, Herceg Z, Gong YY, Routledge MN. Exposure to aflatoxin B1 in utero is associated with DNA methylation in white blood cells of infants in The Gambia. *Int J Epidemiol*. 2015;44:1238–1248.

29. Chen Y, Lemire M, Choufani S, Butcher DT, Grafodatskaya D, Zanke BW, Gallinger S, Hudson TJ, Weksberg R. Discovery of cross-reactive probes and polymorphic CpGs in the Illumina Infinium HumanMethylation450 microarray. *Epigenetics*. 2013;8:203–209.

30. Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *JSSRB*. 1995;57:289–300.

31. Daniels PR, Kardia SLR, Hanis CL, Brown CA, Hutchinson R, Boerwinkle E, Turner ST, Genetic Epidemiology Network of Arteriopathy study. Familial aggregation of hypertension treatment and control in the Genetic Epidemiology Network of Arteriopathy (GENOA) study. *Am J Med*. 2004;116:676–681.

32. Fried LP, Borhani NO, Enright P, Furberg CD, Gardin JM, Kronmal RA, Kuller LH, Manolio TA, Mittelmark MB, Newman A. The Cardiovascular Health Study: design and rationale. *Ann Epidemiol*. 1991;1:263–276.

33. Aryee MJ, Jaffe AE, Corrada-Bravo H, Ladd-Acosta C, Feinberg AP, Hansen KD, Irizarry RA. Minfi: a flexible and comprehensive Bioconductor package for the analysis of Infinium DNA methylation microarrays. *Bioinformatics*. 2014;30:1363–1369.

34. Maksimovic J, Gordon L, Oshlack A. SWAN: Subset-quantile within array normalization for illumina infinium HumanMethylation450 BeadChips. *Genome Biol*. 2012;13:R44.

35. Fortin J-P, Labbe A, Lemire M, Zanke BW, Hudson TJ, Fertig EJ, Greenwood CM, Hansen KD. Functional normalization of 450k methylation array data improves replication in large cancer studies. *Genome Biol*. 2014;15:503.

36. Houseman EA, Accomando WP, Koestler DC, Christensen BC, Marsit CJ, Nelson HH, Wiencke JK, Kelsey KT. DNA methylation arrays as surrogate measures of cell mixture distribution. *BMC Bioinformatics*. 2012;13:86.

37. Day N, Oakes S, Luben R, Khaw KT, Bingham S, Welch A, Wareham N. EPIC-Norfolk: study design and characteristics of the cohort. European Prospective Investigation of Cancer. *Br J Cancer*. 1999;80 Suppl 1:95–103.

38. Bell B, Rose CL, Damon A. The Veterans Administration longitudinal study of healthy aging. *Gerontologist*. 1966;6:179–184.

39. Triche TJ, Weisenberger DJ, Van Den Berg D, Laird PW, Siegmund KD. Low-level processing of Illumina Infinium DNA Methylation BeadArrays. *Nucleic Acids Res*. 2013;41:e90.