

In the format provided by the authors and unedited.

Genetic analyses of diverse populations improves discovery for complex traits

Genevieve L. Wojcik^{1,35}, Mariaelisa Graff^{2,35}, Katherine K. Nishimura^{3,35}, Ran Tao^{4,5,35}, Jeffrey Haessler^{3,35}, Christopher R. Gignoux^{1,6,35}, Heather M. Highland^{2,35}, Yesha M. Patel^{7,35}, Elena P. Sorokin¹, Christy L. Avery², Gillian M. Belbin^{8,9}, Stephanie A. Bien³, Iona Cheng¹⁰, Sinead Cullina^{8,9}, Chani J. Hodonsky², Yao Hu³, Laura M. Huckins¹¹, Janina Jeff^{8,9}, Anne E. Justice², Jonathan M. Kocarnik³, Unhee Lim¹², Bridget M. Lin², Yingchang Lu⁹, Sarah C. Nelson¹³, Sung-Shim L. Park⁷, Hannah Poisner^{8,9}, Michael H. Preuss⁹, Melissa A. Richard¹⁴, Claudia Schurmann^{9,15,16}, Veronica W. Setiawan⁷, Alexandra Sockell¹, Karan Vahi¹⁷, Marie Verbanck⁹, Abhishek Vishnu⁹, Ryan W. Walker⁹, Kristin L. Young², Niha Zubair³, Victor Acuña-Alonso¹⁸, Jose Luis Ambite¹⁷, Kathleen C. Barnes⁶, Eric Boerwinkle¹⁹, Erwin P. Bottinger^{9,15,16}, Carlos D. Bustamante¹, Christian Caberto¹², Samuel Canizales-Quinteros²⁰, Matthew P. Conomos¹³, Ewa Deelman¹⁷, Ron Do^{9,11}, Kimberly Doheny²¹, Lindsay Fernández-Rhodes^{2,22}, Myriam Fornage¹⁴, Benyam Hailu²³, Gerardo Heiss², Brenna M. Henn²⁴, Lucia A. Hindorf²⁵, Rebecca D. Jackson²⁶, Cecelia A. Laurie¹³, Cathy C. Laurie¹³, Yuqing Li^{10,27}, Dan-Yu Lin², Andres Moreno-Estrada²⁸, Girish Nadkarni⁹, Paul J. Norman⁶, Loreall C. Pooler⁷, Alexander P. Reiner¹³, Jane Romm²¹, Chiara Sabatti¹, Karla Sandoval²⁸, Xin Sheng⁷, Eli A. Stahl¹¹, Daniel O. Stram⁷, Timothy A. Thornton¹³, Christina L. Wassel²⁹, Lynne R. Wilkens¹², Cheryl A. Winkler³⁰, Sachi Yoneyama², Steven Buyske^{31,36}, Christopher A. Haiman^{32,36}, Charles Kooperberg^{3,36}, Loïc Le Marchand^{12,36}, Ruth J. F. Loos^{9,11,36}, Tara C. Matise^{33,36}, Kari E. North^{2,36}, Ulrike Peters^{3,36}, Eimear E. Kenny^{8,9,11,34,36*} & Christopher S. Carlson^{3,36*}

¹Department of Biomedical Data Science, Stanford University, Stanford, CA, USA. ²Department of Epidemiology, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA. ³Division of Public Health Science, Fred Hutchinson Cancer Research Center, Seattle, WA, USA. ⁴Department of Biostatistics, Vanderbilt University Medical Center, Nashville, TN, USA. ⁵Vanderbilt Genetics Institute, Vanderbilt University Medical Center, Nashville, TN, USA. ⁶Colorado Center for Personalized Medicine, University of Colorado Anschutz Medical Campus, Aurora, CO, USA. ⁷Department of Preventive Medicine, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. ⁸The Center for Genomic Health, Icahn School of Medicine at Mount Sinai, New York, NY, USA. ⁹The Charles Bronfman Institute of Personalized Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA. ¹⁰Department of Epidemiology and Biostatistics, University of California San Francisco, San Francisco, CA, USA. ¹¹Pamela Sklar Division of Psychiatric Genomics, Icahn School of Medicine at Mount Sinai, New York, NY, USA. ¹²Epidemiology Program, University of Hawaii Cancer Center, Honolulu, HI, USA. ¹³Department of Biostatistics, University of Washington, Seattle, WA, USA. ¹⁴Brown Foundation Institute for Molecular Medicine, The University of Texas Health Science Center, Houston, TX, USA. ¹⁵Hasso-Plattner-Institute for Digital Engineering, Digital Health Center, Potsdam, Germany. ¹⁶Hasso-Plattner-Institute for Digital Health at Mount Sinai, Icahn School of Medicine at Mount Sinai, New York, NY, USA. ¹⁷Information Sciences Institute, University of Southern California, Marina del Rey, CA, USA. ¹⁸Escuela Nacional de Antropología e Historia, Mexico City, Mexico. ¹⁹Human Genetics Center, School of Public Health, The University of Texas Health Science Center, Houston, TX, USA. ²⁰Instituto Nacional de Medicina Genómica, Mexico City, Mexico. ²¹Center for Inherited Disease Research, Johns Hopkins University, Baltimore, MD, USA. ²²Department of Biobehavioral Health, The Pennsylvania State University, University Park, PA, USA. ²³NIH National Institute on Minority Health and Health Disparities, Bethesda, MD, USA. ²⁴Department of Anthropology, University of California Davis, Davis, CA, USA. ²⁵NIH National Human Genome Research Institute, Bethesda, MD, USA. ²⁶Center for Clinical and Translational Science, Ohio State Medical Center, Columbus, OH, USA. ²⁷Cancer Prevention Institute of California, Fremont, CA, USA. ²⁸National Laboratory of Genomics for Biodiversity (UGA-LANGEBIO), Irapuato, Mexico. ²⁹College of Medicine, University of Vermont, Burlington, VT, USA. ³⁰Basic Science Program, Frederick National Laboratory, Frederick, MD, USA. ³¹Department of Statistics, Rutgers University, New Brunswick, NJ, USA. ³²Center for Genetic Epidemiology, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA. ³³Department of Genetics, Rutgers University, New Brunswick, NJ, USA. ³⁴Department of Medicine, Icahn School of Medicine at Mount Sinai, New York, NY, USA. ³⁵These authors contributed equally: Genevieve L. Wojcik, Mariaelisa Graff, Katherine K. Nishimura, Ran Tao, Jeffrey Haessler, Christopher R. Gignoux, Heather M. Highland, Yesha M. Patel. ³⁶These authors jointly supervised this work: Steven Buyske, Christopher A. Haiman, Charles Kooperberg, Loïc Le Marchand, Ruth J. F. Loos, Tara C. Matise, Kari E. North, Ulrike Peters, Eimear E. Kenny, Christopher S. Carlson. *e-mail: eimear.kenny@mssm.edu; ccarlson@fredhutch.org

Supplementary Information

1		
2		
3	1. Detailed PAGE Study Descriptions	3
4	2. Phenotype Harmonization and Modeling	4
5	3. Genotyping and Imputation	8
6	Supplementary Figure 1: Average info from imputed data from MEGA to 1000Genomes Project	
7	within PAGE by minor allele frequency.	9
8	4. Population Substructure	9
9	Supplementary Figure 2: Principal Component Analysis of PAGE Populations.	10
10	5. Meta-analysis versus Mega-analysis in Multi-ethnic Studies	10
11	Supplementary Figure 3. Comparisons of the <i>p</i> -values between meta-analysis and four types of	
12	mega-analysis for MCHC (N=19,803).	11
13	Supplementary Figure 4. Comparisons of the <i>p</i> -values between meta-analysis and four types of	
14	mega-analysis for MCHC with variants whose ethnic-specific minor allele frequency (MAF)	
15	differences are greater than 0.4 (N=19,803).	12
16	Supplementary Figure 5. Quantile-quantile plots of <i>p</i> -values for MCHC when trait values are	
17	assumed to be heterogeneous versus homogeneous. <i>P</i> -values estimated from Wald test.	
18	(N=19,803)	13
19	Supplementary Figure 6. Quantile-quantile plots of <i>p</i> -values for eGFR when trait values are	
20	assumed to be heterogeneous versus homogeneous. <i>P</i> -values estimated from Wald test.	
21	(N=27,900)	13
22	Supplementary Figure 7. Quantile-quantile plots of <i>p</i> -values for FG when trait values are assumed	
23	to be heterogeneous versus homogeneous. <i>P</i> -values estimated from Wald test. (N=23,963)	14
24	Supplementary Figure 8. Quantile-quantile plots of <i>p</i> -values for HbA1c when trait values are	
25	assumed to be heterogeneous versus homogeneous. <i>P</i> -values estimated from Wald test.	
26	(N=11,178)	14
27	Supplementary Figure 9. Quantile-quantile plots of <i>p</i> -values for PR interval when trait values are	
28	assumed to be heterogeneous versus homogeneous. <i>P</i> -values estimated from Wald test.	
29	(N=17,428)	15
30	6. Selecting Principal Components of Ancestry for Use as Covariates	15
31	Supplementary Figure 10: Standardized principal components by population.	16
32	Supplementary Figure 11: Correlation between SNP genotype and PC, by chromosome.	17
33	7. Genome-wide Association Analysis	18
34	Supplementary Figure 12: Comparison of <i>P</i> -values from GWAS for SUGEN (Wald test) vs.	
35	GENESIS across all traits. (N _{max} =49,781; see Extended Data Table 1)	20
36	8. Secondary Signals versus Fine-mapping	21
37	Supplementary Figure 13: Residual signals can represent either refinement of signal or secondary	
38	alleles.	21
39	9. Meta-analysis and Finemapping with GIANT, UK Biobank	22

40	Supplementary Figure 14: BMI PVE.	22
41	Supplementary Figure 15: Finemapping for BMI.	23
42	10. Comparison of novel and secondary variant allele frequencies in European populations	24
43	Supplementary Figure 16: European allele frequencies of novel and secondary findings in PAGE.	24
44		
45	11. Clinically-relevant variants and their distribution in PAGE	25
46	Supplementary Figure 17: World map of <i>HCP5-G</i> frequencies within PAGE groups.	25
47	12. Additional Acknowledgements	26
48	Supplementary Information Bibliography	27
49		
50		

1. Detailed PAGE Study Descriptions

BioMe Biobank: The Charles Bronfman Institute for Personalized Medicine at Mount Sinai Medical Center (MSMC), BioMe™ BioBank (BioMe) is an EMR-linked bio-repository drawing from Mount Sinai Medical Center consented patients which were drawn from a population of over 70,000 inpatients and 800,000 outpatients annually. ¹ The MSMC serves diverse local communities of upper Manhattan, including Central Harlem (86% African American), East Harlem (88% Hispanic/Latino), and Upper East Side (88% Caucasian/White) with broad health disparities. BioMe™ enrolled over 26,500 participants from September 2007 through August 2013, with 25% African American, 36% Hispanic/Latino (primarily of Caribbean origin), 30% Caucasian, and 9% of Other ancestry. The BioMe™ population reflects community-level disease burdens and health disparities with broad public health impact. Biobank operations are fully integrated in clinical care processes, including direct recruitment from clinical sites waiting areas and phlebotomy stations by dedicate Biobank recruiters independent of clinical care providers, prior to or following a clinician standard of care visit. Recruitment currently occurs at a broad spectrum of over 30 clinical care sites. Study participants of self-reported European ancestry were not included in this analysis. (dbGaP study accession number: phs000925).

HCHS/SOL: HCHS/SOL: The Hispanic Community Health Study / Study of Latinos (HCHS/SOL) is a multi-center study of Hispanic/Latino populations with the goal of determining the role of acculturation in the prevalence and development of diseases, and to identify other traits that impact Hispanic/Latino health. ² The study is sponsored by the National Heart, Lung, and Blood Institute (NHLBI) and other institutes, centers, and offices of the National Institutes of Health (NIH). Recruitment began in 2006 with a target population of 16,000 persons of Cuban, Puerto Rican, Dominican, Mexican or Central/South American origin. Household sampling was employed as part of the study design. Participants were recruited through four sites affiliated with San Diego State University, Northwestern University in Chicago, Albert Einstein College of Medicine in Bronx, New York, and the University of Miami. Researchers from seven academic centers provided scientific and logistical support. Study participants who were self-identified Hispanic/Latino and aged 18-74 years underwent extensive psycho-social and clinical assessments during 2008-2011. A re-examination of the HCHS/SOL cohort is conducted during 2015-2017. Annual telephone follow-up interviews are ongoing since study inception to determine health outcomes of interest. (dbGaP study accession number: phs000555).

MEC: The Multiethnic Cohort (MEC) is a population-based prospective cohort study including approximately 215,000 men and women from Hawaii and California. All participants were 45-75 years of age at baseline, and primarily of 5 ancestries: Japanese Americans, African Americans, European Americans, Hispanic/Latinos, and Native Hawaiians. ^{3,4} MEC was funded by the National Cancer Institute in 1993 to examine lifestyle risk factors and genetic susceptibility to cancer. All eligible cohort members completed baseline and follow-up questionnaires. Within the PAGE II investigation, MEC proposes to study: 1) diseases for which we have DNA available for large numbers of cases and controls (breast, prostate, and colorectal cancer, diabetes, and obesity); 2) common traits that are risk factors for these diseases (e.g., body mass index / weight, waist-to-hip ratio, height), and 3) relevant disease-associated biomarkers (e.g., fasting insulin and lipids, steroid hormones). The specific aims are: 1) to determine the population-based epidemiologic profile (allele frequency, main effect, heterogeneity by disease characteristics) of putative causal variants in the five racial/ethnic groups in MEC; 2) for variants displaying effect heterogeneity across ethnic/racial groups, we will utilize differences in LD to identify a more complete spectrum of associated variants at these loci; 3) investigate gene x gene and gene x environment interactions to identify modifiers; 4) examine the associations of putative causal variants with already measured intermediate phenotypes (e.g., plasma insulin, lipids, steroid hormones); and 5) for variants that do not fall within known genes, start to investigate their relationships with gene expression and epigenetic patterns in small genomic studies. For this project, MEC contributed African American, Japanese American, and Native Hawaiian samples. (dbGaP study accession number: phs000220).

PAGE Global Reference Panel: The Global Reference Panel (GRP) was created by Stanford-contributed samples that can act as a population reference dataset across the globe. Therefore, this dataset includes reference individuals, without phenotypes, chosen to help infer ancestry that will aid in understanding the

106 diverse samples available in PAGE. The complete dataset comprises individuals of European, African,
107 Asian, Oceanian, and Native American descent, from a total of over 50 populations. A subset of these
108 individuals from Puno, Peru and Easter Island (Rapa Nui), Chile, are included in the PAGE samples that
109 were whole genome sequenced in 2015. The Global Reference Panel comprises 6 sample sets: (1) a
110 population sample of Andean individuals primarily of Quechuan/Aymaran ancestry from Puno, Peru; (2) a
111 population sample of Easter Island (Rapa Nui), Chile; (3) individuals of indigenous origin from Oaxaca,
112 Mexico; (4) individuals of indigenous origin from Honduras; (5) individuals of indigenous origin from
113 Colombia; (6) individuals of indigenous origin from the Nama and Khomani KhoeSan populations of the
114 Northern Cape, South Africa. PAGE also used samples from the Human Genome Diversity Project (HGDP)
115 ⁵, a subset of the Maasai from HapMap, as well as individuals sampled by the Bustamante Lab and their
116 collaborators. The dataset comprises individuals of European, African, Asian, Oceanian, and Native
117 American descent, from over 50 populations. Study participants were selected to reflect a family history of
118 living in the region. The data are currently available through dbGaP (dbGaP study accession number:
119 phs001033).

120
121 **WHI:** The Women’s Health Initiative (WHI) is a long-term, prospective, multi-center cohort study
122 investigating post-menopausal women’s health in the US. ⁶ WHI was funded by the National Institutes of
123 Health and the National Heart, Lung, and Blood Institute to study strategies to prevent heart disease, breast
124 cancer, colon cancer, and osteoporotic fractures in women 50-79 years of age. WHI involves 161,808
125 women recruited between 1993 and 1998 at 40 centers across the US. The study consists of two parts: the
126 WHI Clinical Trial which was a randomized clinical trial of hormone therapy, dietary modification, and
127 calcium/Vitamin D supplementation, and the WHI Observational Study, which focused on many of the
128 inequities in women’s health research and provided practical information about incidence, risk factors, and
129 interventions related to heart disease, cancer, and osteoporotic fractures. For this project, women who self-
130 identified as European were excluded from the study sample (dbGaP study accession number: phs000227).

131 2. Phenotype Harmonization and Modeling

132
133 The phenotypes included in this study were previously harmonized across the PAGE studies.

134
135 **Anthropometry:** The following anthropometric traits were analyzed: height, body mass index (BMI), and
136 waist-to-hip ratio (WHR). Weight in kilograms and height in centimeters were measured by trained clinic
137 staff in the SOL and WHI studies at the time of enrollment. Waist and hip were also measured in SOL and
138 WHI to the nearest centimeter. In MEC and BioMe weight and height were self-reported by questionnaire
139 and in MEC waist and hip were also self-reported. BMI was then calculated as the ratio of weight to height
140 squared. Individuals <18 years of age and women who were pregnant were also excluded. For GWAS
141 analysis, measurements outside of 6 standard deviations from the mean (based on sex and race) were
142 removed. Then we created sex-specific residuals for each trait adjusted for age (and BMI for waist-to-hip
143 ratio), then inverse normally transformed these residuals. These inverse normally transformed residuals
144 were used in the final analysis and further adjustment was made for self-identified ancestry, study, study
145 center (for MEC and SOL only), and 10 principal components.

146
147 **C-Reactive Protein (CRP):** Serum CRP was reported in mg/L. CRP outliers (+/- 4 standard deviations)
148 were dropped, and CRP was +1 and then natural log transformed. Those who were pregnant at blood draw
149 were excluded from the analysis. There were 28,537 individuals in the final sample. Models were adjusted
150 by age at CRP measurement, sex, BMI, current smoking status, self-identified race/ethnicity, study, study
151 center (for MEC and SOL only), and 10 principal components.

152
153 **Cigarettes per Day (CPD):** The number of cigarettes smoked per day (CPD) was estimated among ever
154 smokers (n=15,8672) based on self-report and electronic health record data. To normalize the distribution
155 of CPD, we added one to the reported CPD and then log transformed this variable. Models were adjusted
156 for age, sex, study, study center (for MEC and SOL only), self-reported race/ethnicity, and the first 10
157 principal components.

158

159 **Chronic Kidney Disease (CKD):** CKD was defined as an eGFR (estimated by the CKD Epi Equation)
160 ≤ 60 ml/min/1.73m² or ICD-9 codes 585.1-585.6, or 585.9, or ICD-10 codes N18.1-N18.5, or N18.9.
161 Participants with end stage-renal disease (ESRD) were excluded from the analysis. CKD was modeled as
162 a binary outcome, and models were adjusted for age, sex, race/ethnicity, study, study center (for MEC and
163 SOL only), and 10 principal components.

164
165 **Coffee Consumption:** The coffee analysis included 35,902 subjects with coffee consumption measured
166 by number of cups per day which was natural log transformed. Models were adjusted age, sex, study, study
167 center (for MEC and SOL only), and first 10 principal components.

168
169 **Diastolic Blood Pressure (DBP):** Diastolic blood pressure was measured as the average of resting
170 measurements in mmHg. Diastolic blood pressure was adjusted by 10 mmHg for the self-reported use of
171 any antihypertensive medication. We winsorized outliers by setting measurements +/- 6 standard deviations
172 from the overall mean to that value. Models for diastolic blood pressure adjusted for age, sex, BMI, self-
173 identified race/ethnicity, study, study center (for MEC and SOL only), and 10 principal components.

174
175 **Electrocardiogram – PR interval:** PR interval is a heritable electrocardiographic measure of atrial and
176 atrioventricular nodal conduction. Resting, supine, or semi-recumbent ECGs were digitally recorded in each
177 study at baseline by certified technicians using standard 12-lead ECGs using either Marquette MAC12 or
178 MAC PC machines (GE Healthcare, Milwaukee, WI, USA; Supplemental Table 1). Comparable procedures
179 were used for preparing participants, placing electrodes, recording, transmitting, processing, and controlling
180 the quality of the ECGs. The PR interval was measured electronically using the Marquette 12SL algorithm.
181 Exclusion criteria included pregnancy, poor ECG quality, non-sinus rhythm including atrial fibrillation and
182 atrial flutter on ECG, pacemaker implantation, second or third degree heart block, extreme PR values (PR
183 ≤ 80 ms or ≥ 320 ms), prevalent heart failure or myocardial infarction, and Wolff-Parkinson-White syndrome
184 on ECG. All models were adjusted for age, sex, study, study center (for MEC and SOL only), self-identified
185 race/ethnicity, systolic blood pressure, height, body mass index, the use of beta-adrenergic blocking agents,
186 and the first 10 principal components.

187 **Electrocardiogram – QRS interval:** QRS interval, from the beginning of the Q wave to the end of the S
188 wave on an electrocardiogram, reflects ventricular depolarization and conduction time. Resting, supine, or
189 semi-recumbent ECGs were digitally recorded in each study at baseline by certified technicians using
190 standard 12-lead ECGs using either Marquette MAC12 or MAC PC machines (GE Healthcare, Milwaukee,
191 WI, USA; Supplemental Table 1). Comparable procedures were used for preparing participants, placing
192 electrodes, recording, transmitting, processing, and controlling the quality of the ECGs. The QRS interval
193 was measured electronically using the Marquette 12SL algorithm. Exclusion criteria included pregnancy,
194 poor ECG quality, non-sinus rhythm including atrial fibrillation and atrial flutter on ECG, pacemaker
195 implantation, second- or third-degree heart block, QRS duration $>$ or equal 120 ms, use of antiarrhythmic
196 medications, prevalent heart failure or myocardial infarction, and Wolff-Parkinson-White syndrome on ECG.
197 All models were adjusted for age, sex, study, study center (for MEC and SOL only), self-identified
198 race/ethnicity, heart rate, systolic blood pressure, height, body mass index, and the first 10 principal
199 components.

200
201 **Electrocardiogram (ECG) measures – QT interval:** QT interval is a measurement of ventricular
202 depolarization and repolarization. Resting, supine, or semi-recumbent ECGs were digitally recorded in each
203 study at baseline by certified technicians using standard 12-lead ECGs using either Marquette MAC12 or
204 MAC PC machines (GE Healthcare, Milwaukee, WI, USA; Supplemental Table 1). Comparable procedures
205 were used for preparing participants, placing electrodes, recording, transmitting, processing, and controlling
206 the quality of the ECGs. The QT interval was measured electronically using the Marquette 12SL algorithm.
207 Exclusion criteria included pregnancy, poor ECG quality, non-sinus rhythm including atrial fibrillation and
208 atrial flutter on ECG, pacemaker implantation, QRS duration $>$ or equal 120ms, and prevalent heart failure.
209 All models were adjusted for age, sex, study, study center (for MEC and SOL only), self-identified
210 race/ethnicity, heart rate, and the first 10 principal components.

211
212 **End-Stage Renal Disease (ESRD):** ESRD was defined as an eGFR (by the CKD-Epi Equation) of ≤ 15

213 ml/min/1.73m² or ICD-9 code 585.6 or ICD-10 code N18.6. Participants with chronic kidney disease (CKD)
214 were excluded from the analysis. ESRD was modeled as a binary outcome, and models were adjusted for
215 age, sex, race/ethnicity, study, study center (for MEC and SOL only), and 10 principal components.
216

217 **Estimated Glomerular Filtration Rate (eGFR) eGFR by CKD Epi Equation:** Continuous eGFR in
218 ml/min/1.73m² was estimated by the serum creatinine-based CKD-Epi equation ⁷, which has been
219 validated for Hispanics ⁸. Non-Hispanic White equations were used to estimate GFR in Hispanic and Asian
220 participants. eGFR was not transformed, and models were adjusted for age, sex, race/ethnicity, study,
221 study center (for MEC and SOL only), and 10 principal components.
222

223 **Fasting Glucose (FG):** FG was reported in mmol/L. Individuals that were pregnant, had a fasting glucose
224 greater than 7 mmol/L, had Type 2 Diabetes, or were non-fasting at measurement were excluded from
225 analysis. Rank normalized residuals were calculated after adjusting for age, sex, age*sex, study, smoking
226 status, and BMI. Association models were adjusted for self-identified race/ethnicity, 10 principal
227 components, and study center (for MEC and SOL only).
228

229 **Fasting Insulin (FI):** FI was reported in pmol/L. Individuals that were pregnant, had a fasting glucose
230 greater than 7 mmol/L, had Type 2 Diabetes, or were non-fasting at measurement were excluded from
231 analysis. Insulin levels were log-transformed. Rank normalized residuals were calculated after adjusting for
232 age, sex, ageXsex, study, smoking status, and BMI. Association models were adjusted for self-identified
233 race/ethnicity, 10 principal components, and study center (for MEC and SOL only).
234

235 **Glycated Hemoglobin (HbA1c):** Glycated Hemoglobin was reported in mmol/mol. Individuals that were
236 pregnant, had a fasting glucose greater than 7 mmol/L, or had Type 2 Diabetes were excluded from the
237 analysis. Rank normalized residuals were calculated after adjusting for age, sex, age by sex, study,
238 smoking status, and BMI. Association models were adjusted for self-identified race/ethnicity, 10 principal
239 components, and study center (for MEC and SOL only).
240

241 **High-Density Lipoprotein (HDL):** HDL measurements were reported in mg/dL, were untransformed, and
242 were adjusted for each individual's medication use by adding a constant based on the type of medication
243 used. If multiple medications were reported, only the correction factor with the largest effect was applied.
244 The constant used for adjustment was based on effects observed in previous publications, and included
245 adjustments for statins ⁹, fibrates ⁹, bile acid sequestrants ¹⁰, niacin ⁹, and cholesterol absorption inhibitors
246 ^{11,12}. An individual's raw HDL measurement was adjusted by the following values if the participant was
247 taking one of these medications: statins: -2.3; fibrates: -5.9; bile acid sequestrants: -1.9; niacin: -9.9,
248 cholesterol absorption inhibitors: +0.0. Those who were pregnant at blood draw, or who had fasted less
249 than 8 hours prior to lipid blood draw were excluded from the study sample. There were 33,063 individuals
250 in the final study sample. Models were adjusted by age at lipid measurement, sex, study, study center (for
251 MEC and SOL only), self-identified race/ethnicity, and 10 principal components.
252

253 **Hypertension (HT):** Hypertension cases were defined based on any of the following criteria: 1) measured
254 systolic blood pressure ≥140 mmHg, 2) measured diastolic blood pressure ≥90 mmHg, 3) reported use of
255 any antihypertensive medication, or 4) ICD-9 codes 401.x-405.x or ICD-10 codes I10.x-I15.x. Individuals
256 not meeting any of these criteria were considered normotensive (controls). Models for hypertension were
257 adjusted for age, sex, BMI, study, study center (for MEC and SOL only), self-identified race/ethnicity, and
258 10 principal components.
259

260 **Low-Density Lipoprotein (LDL):** LDL measurements were reported in mg/dL, and were calculated using
261 the Friedewald Equation ¹³, which subtracts the HDL measurement and the Triglyceride measurement
262 (divided by 5) from the Total Cholesterol value. LDL was not calculated if the triglyceride value was greater
263 than 400 mg/dL. LDL values were then adjusted for each individual's medication use by adding a constant
264 based on the type of medication used. If multiple medications were reported, only the correction factor with
265 the largest effect was applied. The constant used for adjustment was based on effects observed in previous
266 publications, and included adjustments for statins ⁹, fibrates ⁹, bile acid sequestrants ¹⁰, niacin ⁹, and
267 cholesterol absorption inhibitors ^{11,12}. An individual's raw LDL measurement was adjusted by the following
268 values if the participant was taking one of these medications: statins: +49.9; fibrates: +40.1; bile acid

269 sequestrants: +40.5; niacin: +24.7; cholesterol absorption inhibitors: +40.5. Those who were pregnant at
270 blood draw, or who had fasted less than 8 hours prior to lipid blood draw were excluded from the study
271 sample. There were 32,221 individuals in the final study sample. Models were adjusted by age at lipid
272 measurement, sex, study, study center (for MEC and SOL only), self-identified race/ethnicity, and 10
273 principal components.

274
275 **Mean Corpuscular Hemoglobin Concentration (MCHC):** Mean corpuscular hemoglobin concentration
276 was reported in g/dL and was untransformed. MCHC is calculated using the formula:
277 $100 \times \text{hemoglobin} / \text{hematocrit}$. MCHC outliers (± 4 standard deviations) were dropped, along with
278 observations for HIV+ individuals, participants with a reported hereditary anemia, and women pregnant at
279 time of blood draw. The MCHC model was adjusted for age at blood draw, sex, current smoking status,
280 self-identified race/ethnicity, study, study center (for MEC and SOL only), and 10 principal components.

281
282 **Platelet Count (PLT):** Platelet count was reported as cells $\times 10^9/L$ and was untransformed. PLT outliers
283 (± 4 standard deviations) were dropped, along with observations for HIV+ individuals, and women
284 pregnant at time of blood draw. The PLT model was adjusted for age at blood draw, sex, current smoking
285 status, self-identified race/ethnicity, study, study center (for MEC and SOL only), and 10 principal
286 components.

287
288 **Systolic Blood Pressure (SBP):** Systolic blood pressure was measured as the average of resting
289 measurements in mmHg. Systolic blood pressure was adjusted by 15 mmHg for the self-reported use of
290 any antihypertensive medication. We winsorized outliers by setting measurements ± 6 standard deviations
291 from the overall mean to that value. Models for systolic blood pressure adjusted for age, sex, BMI, self-
292 identified race/ethnicity, study, study center (for MEC and SOL only), and 10 principal components.

293
294 **Total Cholesterol (TC):** Total Cholesterol measurements were reported in mg/dL, were untransformed,
295 and were adjusted for each individual's medication use by adding a constant based on the type of
296 medication used. If multiple medications were reported, only the correction factor with the largest effect was
297 applied. The constant used for adjustment was based on effects observed in previous publications, and
298 included adjustments for statins⁹, fibrates⁹, bile acid sequestrants¹⁰, niacin⁹, and cholesterol absorption
299 inhibitors^{11,12}. An individual's raw total cholesterol measurement was adjusted by the following values if the
300 participant was taking one of these medications: statins: +52.1; fibrates: +46.1; bile acid sequestrants: +0.0;
301 niacin: +34.6; cholesterol absorption inhibitors: +40.5. Those who were pregnant at blood draw, or who had
302 fasted less than 8 hours prior to lipid blood draw were excluded from the study sample. There were 33,185
303 individuals in the final study sample. Models were adjusted by age at lipid measurement, sex, study, study
304 center (for MEC and SOL only), self-identified race/ethnicity, and 10 principal components.

305
306 **Type II Diabetes (T2D):** The Type 2 Diabetes analysis included 14,046 cases and 31,695 controls with
307 complete covariate data after excluding individuals who were pregnant at blood draw and those who were
308 classified as cases for Type 1 Diabetes. Controls with glucose values greater than 7 mmol/L were excluded,
309 as well as any cases that were younger than 20 years of age. The models were adjusted by age at T2D
310 diagnosis, sex, study, study center (for MEC and SOL only), BMI, self-identified race/ethnicity, and first 10
311 principal components.

312
313 **Triglycerides (TG):** Triglyceride measurements were reported in mg/dL, were adjusted for each individual's
314 medication use by adding a constant based on the type of medication used, and then natural log
315 transformed. If multiple medications were reported, only the correction factor with the largest effect was
316 applied. The constant used for adjustment was based on effects observed in previous publications, and
317 included adjustments for statins⁹, fibrates⁹, bile acid sequestrants¹⁰, niacin⁹, and cholesterol absorption
318 inhibitors^{11,12}. An individual's raw Triglyceride measurement was adjusted by the following values if the
319 participant was taking one of these medications: statins: +18.4; fibrates: +57.1; bile acid sequestrants: +0.0;
320 niacin: +89.4; cholesterol absorption inhibitors: +0.0. Those who were pregnant at blood draw, or who had
321 fasted less than 8 hours prior to lipid blood draw were excluded from the study sample. Individuals with a
322 Triglyceride value greater than 3000 mg/dL were dropped from the analysis (n=1). There were 33,096
323 individuals in the final study sample. Models were adjusted by age at lipid measurement, sex, study, study
324 center (for MEC and SOL only), self-identified race/ethnicity, and 10 principal components.

325
326 **White Blood Cell Count (WBC):** Total WBC count was measured in 10^9 cells/L. WBC count outliers (± 4
327 standard deviations) were dropped, and total WBC was natural log transformed. Those who were pregnant
328 a blood draw were excluded from the analysis. There were 28,534 individuals in the final sample. Models
329 were adjusted by age at WBC measurement, sex, BMI, study, current smoking status, self-identified
330 race/ethnicity, study, study center (for MEC and SOL only), and 10 principal components.

331 **3. Genotyping and Imputation**

332
333 One major challenge in multi-ethnic studies is the limited availability of genotyping arrays that
334 comparably tag variation in multiple genetic ancestries, especially in those with African ancestry. To
335 address this, a collaboration among PAGE, Illumina Inc., the Consortium on Asthma among African-
336 ancestry Populations in the Americas (CAAPA)¹⁴, and other academic partners developed the Multi-Ethnic
337 Genotyping Array (MEGA), which includes a GWAS scaffold designed to tag both common and low
338 frequency variants in global populations.¹⁵ (**Extended Data Figure 2**) Additionally, it contains enhanced
339 tagging in exonic regions, hand-curated content to interrogate clinically relevant variants, and enriched
340 coverage to fine-map known GWAS loci.¹⁶

341 DNA was isolated from blood (SOL, BioMe, GRP), buffy coat (WHI, MEC), mouthwash/saliva
342 (MEC, GRP), or lymphoblastoid cell line (GRP). There were 548 HapMap genotyping control samples, and
343 1,001 blind study duplicate samples. Samples were genotyped on complete or partial 96-well plates, over
344 three batches. MEC, BioMe, SOL and WHI samples were distributed across the three batches in a
345 proportion that represented the size of each study. Sex, race/ethnicity, recruitment site, and DNA source
346 were not stratified, but samples were randomly selected within each stratification level. To better understand
347 the rich genetic diversity within PAGE, particularly in underrepresented ancestries from the Africa and the
348 Americas, we genotyped an additional 1,553 samples on MEGA from a Global Reference Panel (GRP),
349 which were drawn from the Human Genome Diversity Project⁵ and supplemented with previously sampled
350 populations from the Americas and Africa^{14,17-20}. Some GRP samples were included on the plates in
351 Batches 1 and 2, but the majority of these samples were included on 12 separate plates. Each plate
352 contained one or two duplicates. Duplicate samples were placed on a different plate than the original sample.
353 Each plate contained an average of one HapMap sample. There were also 110 investigator controls
354 previously genotyped that were included with at most one of these controls per plate.

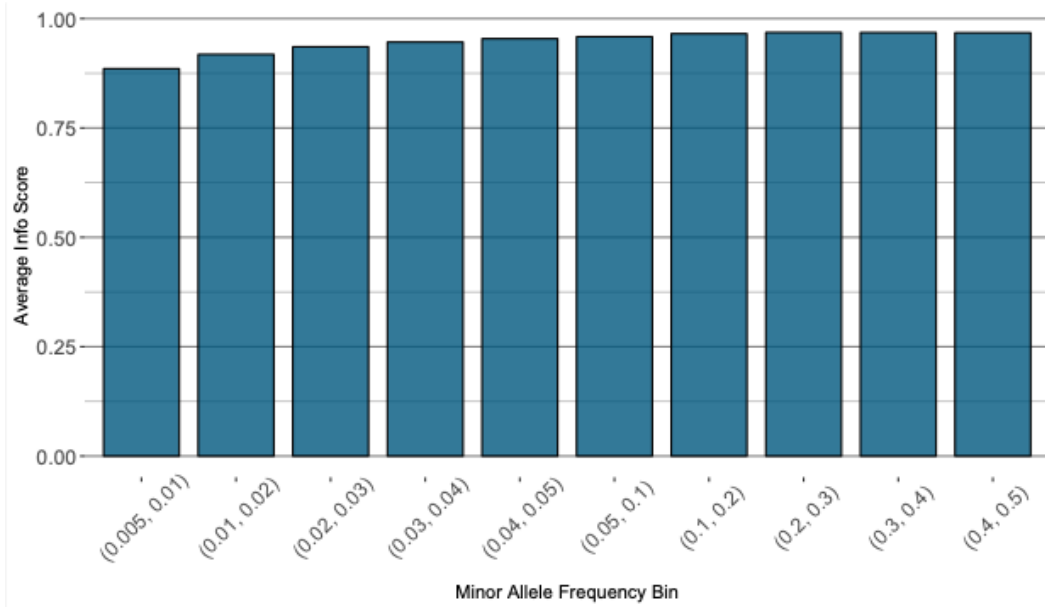
355 A total of 53,426 samples were genotyped at the Center for Inherited Disease Research (CIDR)
356 using the Multi-Ethnic Genotyping Array (MEGA), Consortium version. MEGA was designed through a
357 collaboration between PAGE, University of Michigan, CAAPA, and Illumina to provide broad coverage for
358 globally diverse populations, as well as provide enhanced exomic, functional, and clinically-relevant
359 content. Genotypes were called by the Center for Inherited Disease Research (CIDR) using the
360 GenomeStudio version 2001.1, Genotyping Module 1.9.4, and GenTrain version 1.0.

361 Genotyping data that passed initial quality control at CIDR, including sex discrepancies, Mendelian
362 inconsistencies, unexpected duplication, unexpected non-duplication, poor performance, or DNAmixture
363 observed were released to the Quality Assurance / Quality Control (QA/QC) analysis team at the University
364 of Washington Genetic Coordinating Center (UWGAC), the study investigator's team, and dbGaP. The
365 UWGAC QA/QC team used quality control methods previously described by Laurie et al.²¹ The UWGAC
366 QA/QC team further removed samples with identity issues, restricted consent, and duplicate scans to return
367 data for 51,520 subjects.

368 A total of 1,705,969 variants were genotyped on MEGA. Variant-level quality control (QC) was
369 completed by were filtered through various criteria, including the exclusion of (1) CIDR technical filters, (2)
370 variants with missing call rate $\geq 2\%$, (3) variants with more than 6 discordant calls in 988 study duplicates,
371 (4) SNPs with greater than 1 Mendelian errors in 282 trios and 1,439 duos, (5) variants with a Hardy-
372 Weinberg p-value less than 10^{-4} , (6) variants with sex difference in allele frequency ≥ 0.2 for
373 autosomes/XY, and (7) variants with sex difference in heterozygosity > 0.3 for autosomes/XY. After variant
374 QC, a total of 1,438,399 variants remained.

375 Imputation was conducted at the UWGAC. Sites were further restricted to variants with (1) known
376 chromosome and position; (2) located on chromosomes 1-22, X, or XY (pseudo-autosomal); (3) with
377 unique positions, which involved removing redundant and duplicate sites; and (4) sites with available

378 strand annotation. After these restrictions, a total of 1,402,653 sites remained. The study samples were
 379 phased with SHAPEIT2²² and imputed with IMPUTE2²³ to the 1000 Genomes Project Phase 3 data
 380 release²⁴. Reference panel variants were filtered to exclude all monomorphs and singletons (i.e.
 381 restricting to minor allele count (MAC) ≥ 2 across all 1000 Genomes Phase 2 samples). Imputed
 382 variants were excluded if the IMPUTE2 info score was less than 0.4.
 383



384

385 **Supplementary Figure 1: Average info from imputed data from MEGA to**
 386 **1000Genomes Project within PAGE by minor allele frequency.**

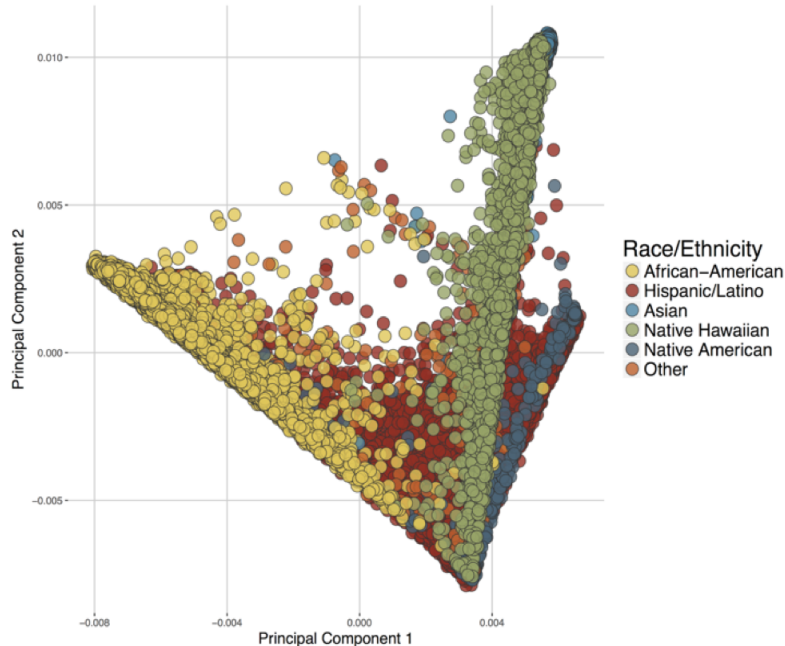
387 We show high imputation quality across all minor allele frequency bins from 0.5-50%.
 388

389 **4. Population Substructure**

390

391 Historically, analyses have been stratified by self-identified race/ethnicity to account for confounding by
 392 genetic ancestry. In PAGE, we conducted principal component analysis to evaluate population substructure
 393 and mapped self-identified racial/ethnic groups (Hispanic/Latino, African American, Asian, Native Hawaiian,
 394 Native American, and Other) onto the estimated principal components (PCs). The selection of unrelated
 395 individuals was essential for accurate estimation of the principal components within the global study
 396 population. Kinship coefficients were estimated using PC-Relate, as implemented in the R package
 397 GENESIS (Conomos et al. 2015; Conomos, Reiner, et al. 2016). The SNPRelate (Zheng et al. 2012)
 398 package in R was then used for principal components analysis using unrelated individuals, defined as
 399 pairwise kinship coefficients less than $2^{-(9/2)}$. Since principal components are required for unbiased
 400 kinship estimation in admixed populations, the two estimation procedures were iterated to ensure that the
 401 principal components were computed over unrelated individuals. Principal component scores were then
 402 estimated for all remaining individuals by projection.²⁵
 403

404 Most notably in Hispanics/Latinos, but evident to a lesser extent in all populations, genetic ancestry reveals
 405 greater demographic complexity compared with culturally assigned labels. Genetic ancestry appears as a
 406 continuum, demonstrating that it is not categorical in diverse populations that have varying degrees of
 407 admixture (**Supp Figure 1**). Stratifying by self-reported race/ethnicity would fail to separate groups with
 408 similar patterns of genetic ancestry and therefore would still require adjustment of PCs with reduced
 409 statistical power in a smaller sample size. For this reason, we pooled all samples in a single mega-analysis.



410

411 **Supplementary Figure 2: Principal Component Analysis of PAGE Populations.**
 412 *Scatterplot of PCs for PAGE racial/ethnic groups. Each point represents one individual, color-coded by self-*
 413 *identified race/ethnicity. Global variation for all PAGE participants for principal component (PC) 1 versus*
 414 *PC2. Genotyped individuals self-identified as Hispanic/Latino (N=22,216), African American (N=17,299),*
 415 *Asian (N=4,680), Native Hawaiian (N=3,940), Native American (N=652), or Other (N=1,052).*

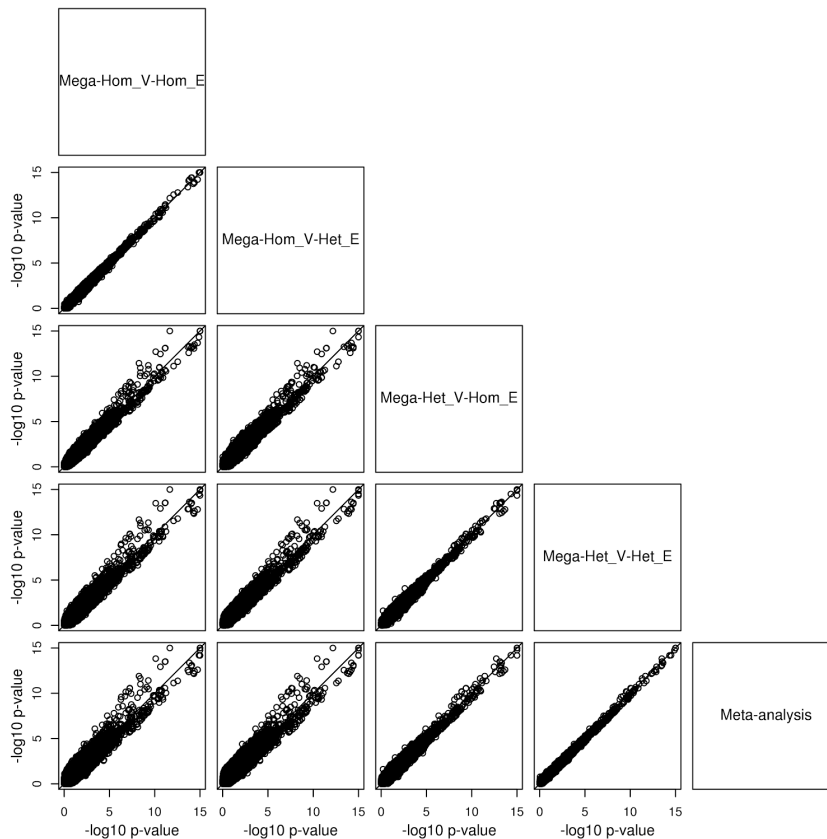
416 **5. Meta-analysis versus Mega-analysis in Multi-ethnic**
 417 **Studies**

418
 419 It has been shown, both theoretically and numerically, that meta-analysis and mega-analysis of
 420 independent studies are (asymptotically) equivalent, if mega-analysis allows nuisance parameters (i.e.,
 421 trait variances and covariate effects) to be different among studies^{26,27}. The comparison of meta and
 422 pooled analysis has some subtle aspects that are important. Considering a simple scenario with a variant
 423 with MAF=50% in population 1 and MAF=1% in population 2, and the same sample size and same effect
 424 size in each population, if allele frequency is the only difference between the two studies then it is actually
 425 slightly more powerful to perform a mega analysis (with no adjustment for study) rather than meta-
 426 analysis. In general, however, we expect that allele frequency will not constitute the only difference
 427 between studies, and we therefore always include study indicators in pooled analyses. Inclusion of study
 428 as a covariate ensures that the mega-analysis estimator for the genetic effect can be expressed as an
 429 inverse-variance weighted estimator, with weights that are asymptotically equivalent to the weights in the
 430 meta-analysis estimator; see Example 1 in Lin & Zeng (2010).²⁶ **Supplementary Figure 2** compares the
 431 results of meta-analysis to four types of mega-analysis (assuming heterogeneous vs homogeneous trait
 432 variances and heterogeneous vs homogeneous covariate effects) for MCHC. The mega-analysis that
 433 allows both trait variances and covariate effects to be different among studies (i.e., ethnicities in our case)
 434 fits the same model as meta-analysis does, so the *p*-values from the two methods are virtually identical;
 435 the *p*-values can be quite different when trait variances are allowed to be different across ethnicities vs
 436 when they are assumed to be the same across ethnicities; and the results are fairly similar when covariate
 437 effects are allowed to be different across ethnicities vs when they are assumed to be the same across
 438 ethnicities. **Supplementary Figure 3** shows that the same conclusions hold for variants whose ethnic-
 439 specific MAF differences are greater than 0.4. We chose mega-analysis over meta-analysis *because the*
 440 *former allows related individuals across studies (whereas the latter does not) and provides greater*

441 *flexibilities in modelling the effects of covariates.*

442

443 SUGEN and GENESIS allow trait variances to be different among studies (i.e., ethnicities in our case)
444 whereas other LMM methods do not. We chose to focus on methods adjusting for global ancestry as
445 these were more stable across the extremely heterogeneous mix of populations in PAGE, where local
446 ancestry estimation could be challenging to reconcile. Further, as we have previously shown, local
447 ancestry adjustment is expected to impair statistical power for discovery when compared to global
448 ancestry adjustment.²⁸ **Supplementary Figures 4-8** display the p -values for the five traits with the most
449 severe trait-variance heterogeneity, when the trait variances are assumed to be homogeneous vs
450 heterogeneous across ethnicities in the analysis. These results confirm the theoretical expectation that
451 allowing heterogeneous variances yields better control of type I error and higher power in association
452 tests.



453

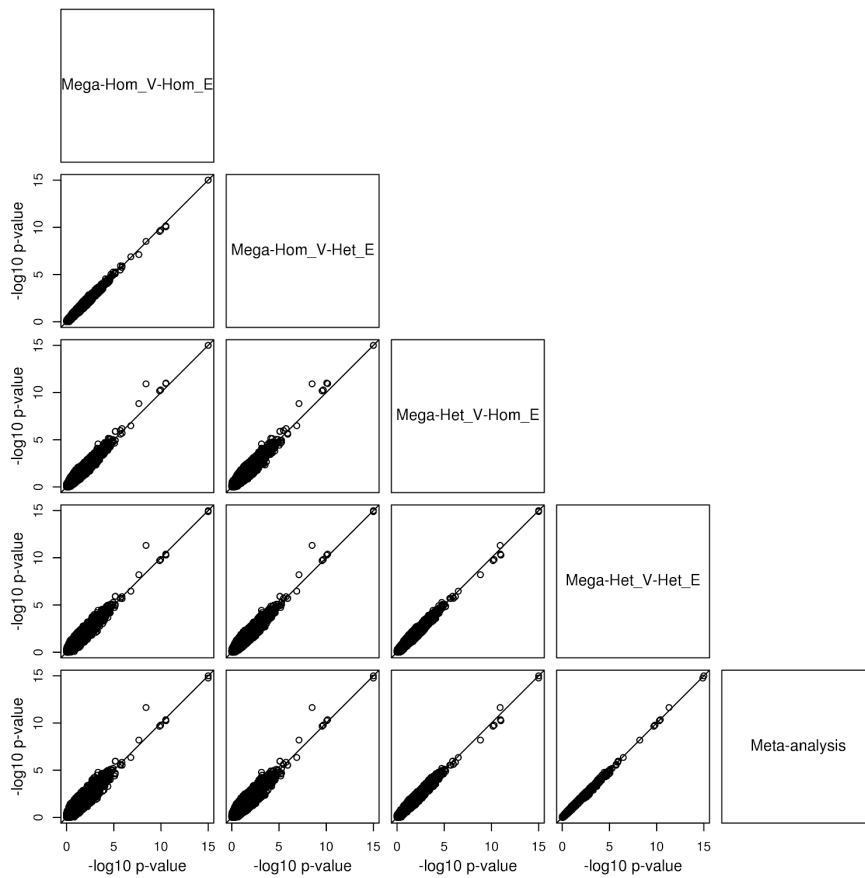
454

455 **Supplementary Figure 3. Comparisons of the p -values between meta-analysis and**
456 **four types of mega-analysis for MCHC (N=19,803).**

457 *Het_V and Hom_V denote, respectively, heterogeneous and homogeneous trait variances among ethnic*
458 *groups; Het_E and Hom_E denote, respectively, heterogeneous and homogeneous covariate effects*
459 *among ethnic groups. The p -values are estimated from Wald test, with values less than 110^{-15} winsorized*
460 *at 110^{-15} .*

461

462

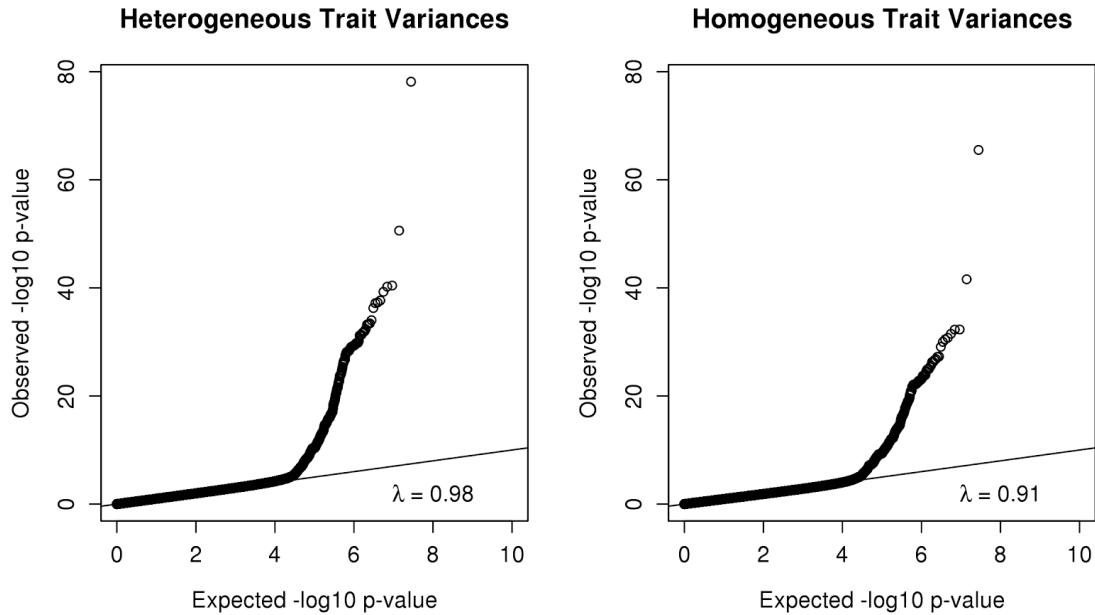


463
464

465 **Supplementary Figure 4. Comparisons of the p -values between meta-analysis and**
 466 **four types of mega-analysis for MCHC with variants whose ethnic-specific minor**
 467 **allele frequency (MAF) differences are greater than 0.4 (N=19,803).**

468 *Het_V* and *Hom_V* denote, respectively, heterogeneous and homogeneous trait variances among ethnic
 469 groups; *Het_E* and *Hom_E* denote, respectively, heterogeneous and homogeneous covariate effects
 470 among ethnic groups. The p -values are estimated from Wald test, with values less than 110^{-15} winsorized
 471 at 110^{-15} .

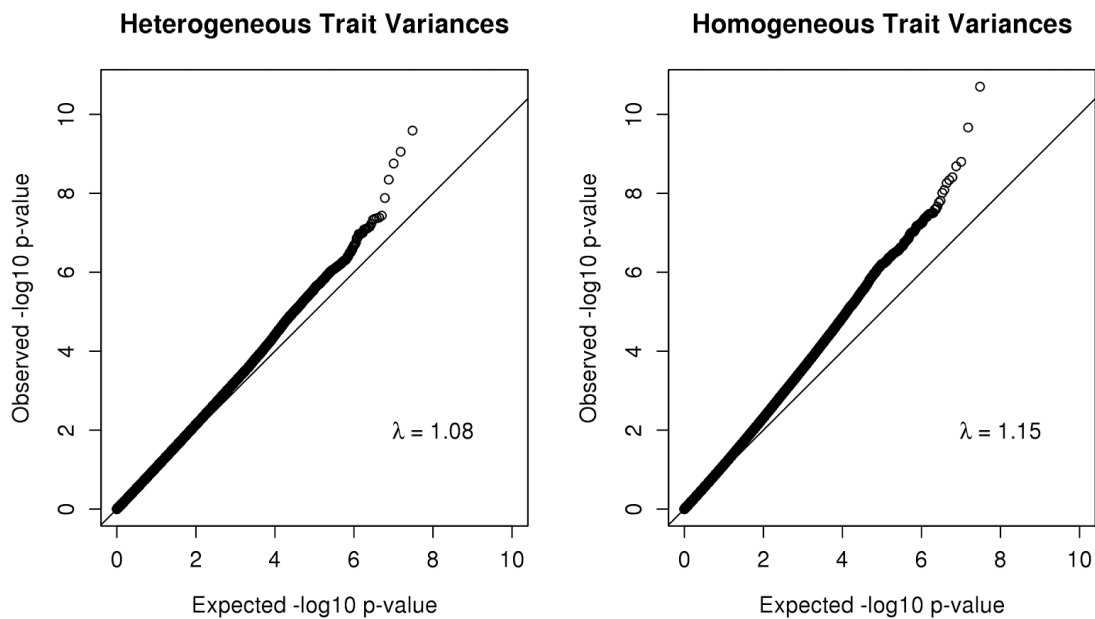
472
473



474

475 **Supplementary Figure 5. Quantile-quantile plots of p-values for MCHC when trait**
 476 **values are assumed to be heterogeneous versus homogeneous. P-values**
 477 **estimated from Wald test. (N=19,803)**

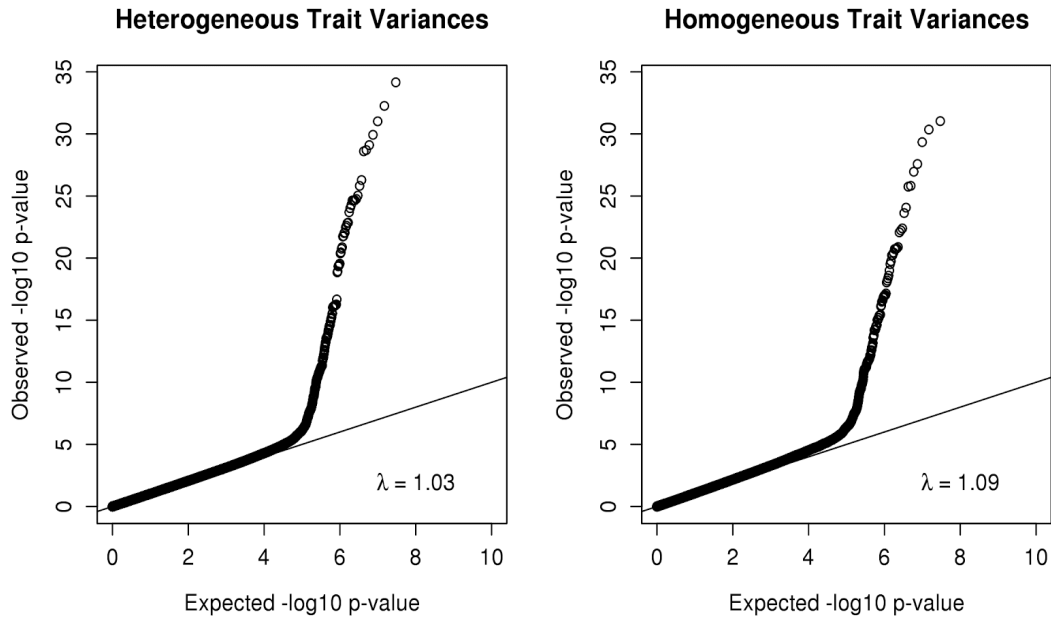
478



479

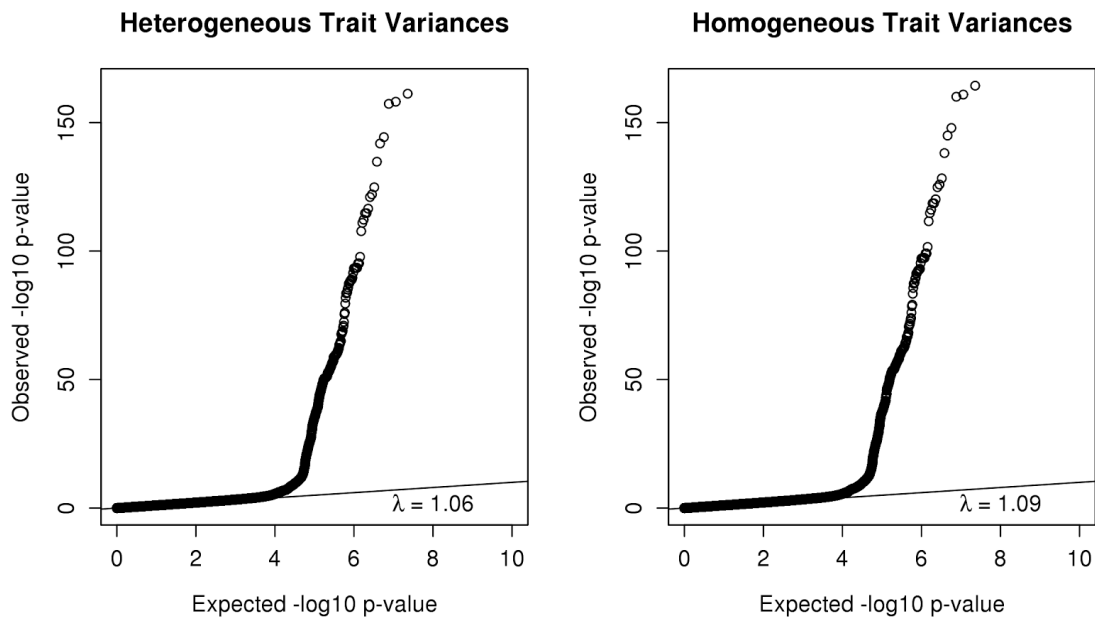
480 **Supplementary Figure 6. Quantile-quantile plots of p-values for eGFR when trait**
 481 **values are assumed to be heterogeneous versus homogeneous. P-values**
 482 **estimated from Wald test. (N=27,900)**

483



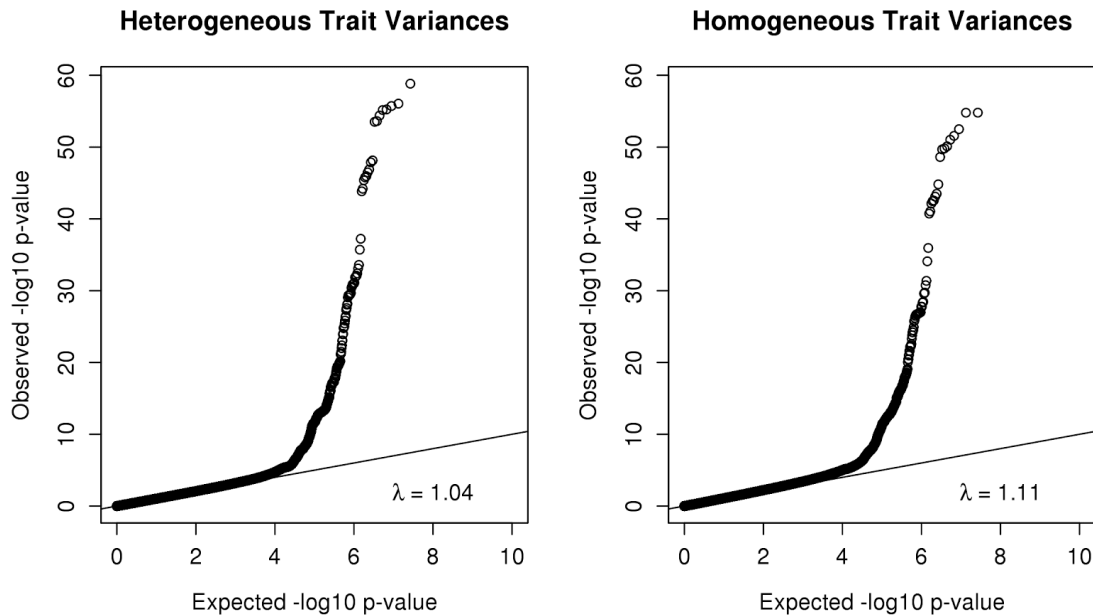
484

485 **Supplementary Figure 7. Quantile-quantile plots of p-values for FG when trait**
 486 **values are assumed to be heterogeneous versus homogeneous. P-values**
 487 **estimated from Wald test. (N=23,963)**
 488



489

490 **Supplementary Figure 8. Quantile-quantile plots of p-values for HbA1c when trait**
 491 **values are assumed to be heterogeneous versus homogeneous. P-values**
 492 **estimated from Wald test. (N=11,178)**
 493



494

495 **Supplementary Figure 9. Quantile-quantile plots of p-values for PR interval when**
 496 **trait values are assumed to be heterogeneous versus homogeneous. P-values**
 497 **estimated from Wald test. (N=17,428)**

498

499

500 **6. Selecting Principal Components of Ancestry for Use** 501 **as Covariates**

502

503 Previous GWAS analyses have generally used between 3 and 10 principal components (PCs) for
 504 adjustment for population structure, often simply determined by post-hoc interpretation of results. In our
 505 study, we used two additional criteria to calibrate the number of PCs to include in the analysis: reference
 506 population specificity (**Supp Fig 9**) and chromosome specificity (**Supp Fig 10**).

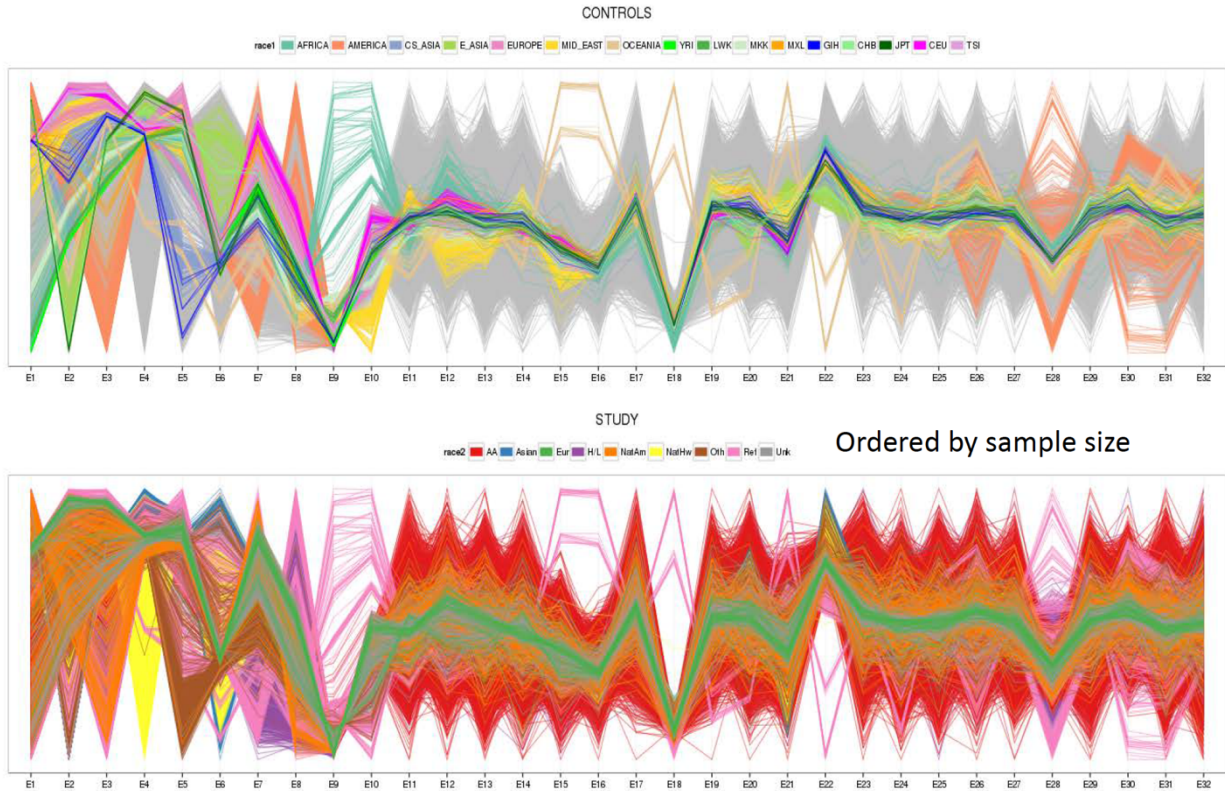
507 While PC1 through PC10 showed clear population specificity in the reference samples, most of the higher
 508 PCs showed much weaker population-specificity. PC9 and PC10 were clearly African in origin in the
 509 reference data (top panel), although these two PC did not vary tremendously within the PAGE study
 510 population. A few higher PC showed population specificity in the reference samples (PC15, 16, 18 and
 511 22), but we felt it was more appropriate to specify a single threshold (exclude all PC past PC10), rather
 512 than cherry picking a PC less-relevant to the samples with phenotypic data.

513 An alternative approach to assessing which principal components to include in an analysis uses
 514 chromosome-specificity. Principal components which load uniformly across the genome are likely to
 515 represent ancestral populations, while principal components that load heavily onto specific chromosomes
 516 are frequently artifacts in the data (e.g. large polymorphic inversions on chromosomes 8 and 17). A
 517 loading analysis is shown in the new **Supp Fig 10**, which shows the correlation between SNP genotype
 518 and PC1 through PC20. Again, PC1 through PC8 are clearly consistent across the genome, so their
 519 inclusion as covariables is justified. In contrast, PC11 through PC20 show significant chromosome

520 specificity, so we chose to exclude PC11 and higher due to the lack of discernible population specificity in
521 either the reference or study samples.

522 PC9 and PC10 posed a unique challenge. Population specificity (**Supp Fig 9**) suggests that these two
523 might be important in African populations, arguing for inclusion, but the chromosome loading (**Supp Fig**
524 **10**) suggests that these two could be artifactual. We conservatively chose to retain these two PC in our
525 analyses, but analyses with either eight or twelve PCs are almost entirely consistent with the ten PC
526 analysis, so PC9 and PC10 were not influential in the results that are presented.

527

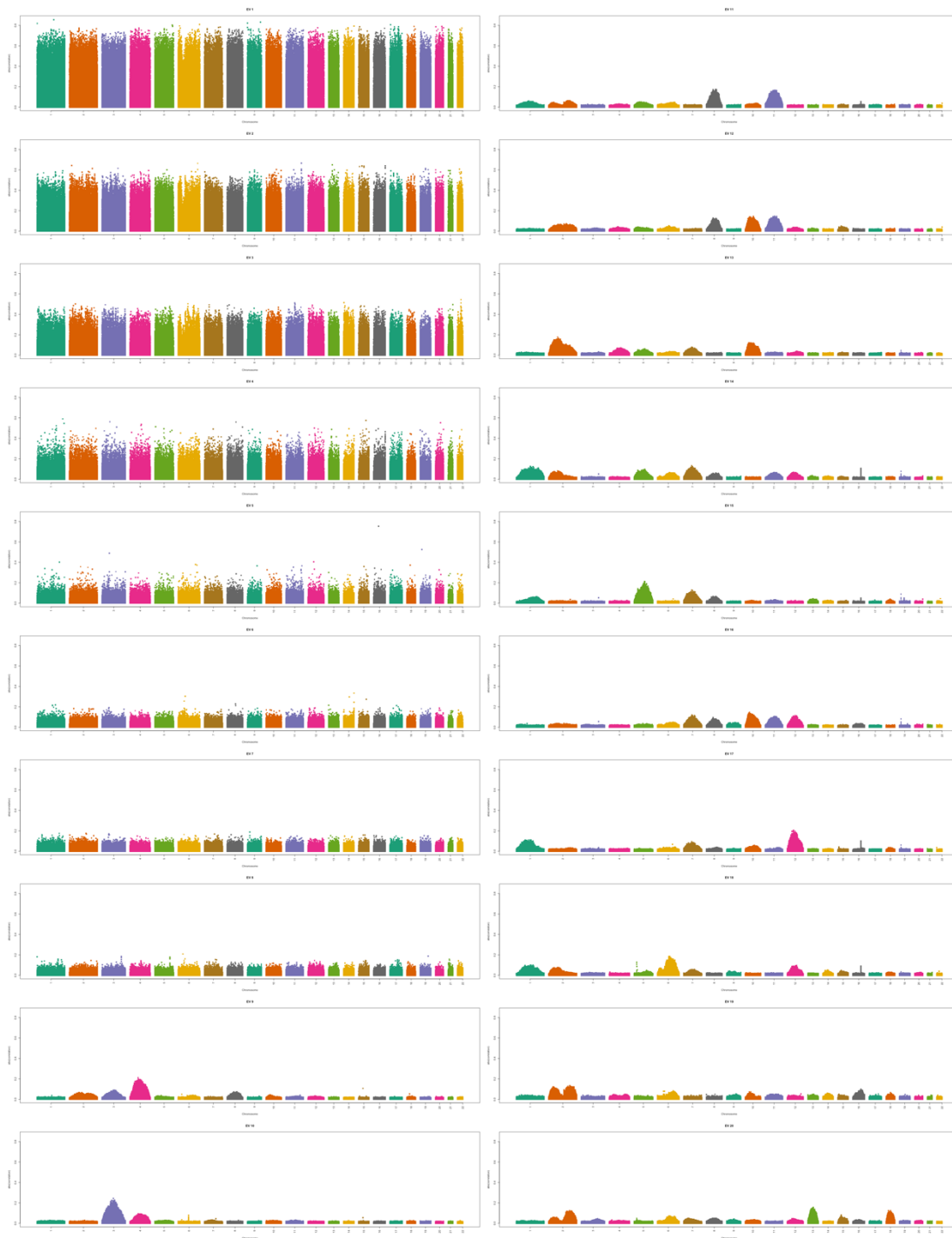


528

529 **Supplementary Figure 10: Standardized principal components by population.**
530 *After standardizing the ranges of principal component 1 (PC1) through PC32, we plotted the value for*
531 *each individual as a line (N=49,839). The top panel shows individuals within the reference population*
532 *color coded by population, with study samples in grey. The bottom panel shows PAGE participants*
533 *colored by their self-reported ancestry with the reference populations in pink. This allows us to see the*
534 *distribution of different race/ethnicity groups across the different principal components. For example, in*
535 *the top panel we see orange lines at the outer ranges of the distribution, indicating that principal*
536 *component 3 represents the spectrum of Native American ancestry, as orange denotes reference groups*
537 *from the Americas.*

538

539



540

541 **Supplementary Figure 11: Correlation between SNP genotype and PC, by**
 542 **chromosome.**

543 Genome position is shown on the X axis, and the correlation between genotype and PC is shown on the
 544 Y axis, range (0,1) in all panels. The first eight PCs are clearly consistent across the genome, while
 545 higher PCs tend to be more chromosome-specific. (N=49,839)

546

547

7. Genome-wide Association Analysis

Analysis:

We based our analysis on generalized linear models of form

$$M1: g(E(Y)) = X\alpha + G\beta$$

Where Y is the vector (of length n , the number of participants used in a given analysis) of observed outcomes which may be continuous or binary, $E()$ denotes expectation, g is a link function, X denotes a matrix of adjustment variables (of size $n \times p$, with p the number of adjustment variables) including age, sex, principal components (PCs), self-identified ethnicity as a proxy for cultural background, and other relevant variables, and G is a vector of length n of observed or imputed allele counts for a given variant of interest.

For binary traits we used the logit link function $g(x)=\log(x)-\log(1-x)$ so that it is the log odds that is linear in the G and X variables, while for continuous traits the identity function, $g(x)=x$, was used. Some continuous traits were inverse-normal transformed before model M1 (and other models) were applied to the sorted trait values. In particular, if $Y(1), \dots, Y(n)$ are the sorted trait values, then $Y(i)$ is replaced by $\Phi^{-1}(i/(n+1))$, where Φ^{-1} is the inverse cumulative normal distribution function.

We used a combination of methods to account for hidden population structure (e.g. admixture) and relatedness among study participants. For both continuous and binary traits, we included leading PCs of the genotype matrix as part of the adjustment variables in M1 to ensure that large scale population structure would not induce false positive associations. In all our analyses, we included ten PCs in M1. Based on our assessment, ten PCs was sufficient to account for all major ethnic variation, while not including too many PCs to negatively affect the power of the analyses. Limited experimentation (not shown) suggested that adding a few more PCs did not noticeably influence the results.

For continuous traits, we adopted a linear mixed models (LMM) and a generalized estimating equations (GEE) approach to correct for the effects of relatedness between individuals. For binary traits, we only used the GEE approach. Two programs, GENESIS and SUGEN, were developed by PAGE collaborators to implement these approaches for GWAS studies of populations with genetic admixture and known or cryptic relatedness.

GENESIS: The GENESIS program²⁹⁻³¹ is available as a Bioconductor package made available in R, and uses a LMM to test for SNP - phenotype associations. For continuous traits, the regression models were fit assuming a variance matrix model for the variance-covariance matrix of the outcomes Y of form $\sigma^2 I + \gamma^2 K$.

Here I is the identity matrix and K is the genetic relatedness matrix computed from the available SNP data, once for all of PAGE II. Score tests of $\beta=0$ were computed by replacing M1 with the null model

$$M0: E(Y) = X\alpha$$

and then in M1 performing a test for $\beta = 0$ with σ^2 and γ^2 held at their estimated values from the fit of M0. The elements of G (in turn) are simply the observed allele counts (0, 1, or 2), or for imputed data the estimated allele counts (taking values from 0-2) for each of the variants of interest. Using the variance model V1 corrects the score tests of $\beta = 0$ for both close and more distant relationships between individuals in the dataset. Estimation of α and the variance parameters σ^2, γ^2 only needs to be performed once which provides a great savings in computer time needed to use GENESIS.

Both GENESIS and SUGEN rely upon the estimated relationship matrix K . The GENESIS package includes the program PC-Relate, which uses a principal component analysis (PCA) based method to infer genetic relatedness in samples with unspecified and unknown population structure. By using individual-specific allele frequencies estimated with sample principal component eigenvectors, it provides estimates of kinship coefficients and identity by descent (IBD) sharing probabilities in samples with population structure, admixture, and HWE departures. It does not require additional reference population panels or prior specification of the number of ancestral subpopulations.

602 **SUGEN:** The SUGEN program ³² is a command-line software program developed for genetic association
603 analysis with complex survey sampling and relatedness patterns. It implements the generalized
604 estimating equation (GEE) approach, which empirically accounts for within-family correlations without
605 modeling the correlation structures of complex pedigrees.
606

607 Association analysis in SUGEN requires the study subjects to be grouped into “independent” families.
608 There is a complex pattern of relatedness in HCHS/SOL: individuals in the same household are related,
609 and there is endogamous mating within the Hispanic/Latino community, such that some households are
610 connected into large pedigrees. To address this challenge, we first used the genetic relationship matrix K
611 to identify pairs of individuals who are first-degree or second-degree relatives. We then formed extended
612 families by connecting the households who share first-degree relatives or either first- or second-degree
613 relatives. The trait values are assumed to be correlated within families but independent between families.
614 In our dataset, we found it sufficient to account for first-degree relatedness in association analysis.
615

616 The GEE approach uses a “sandwich” variance estimator to empirically estimate within-family
617 correlations. SUGEN adopts a modified version of the sandwich variance estimator, which replaces the
618 empirical covariance matrix of the score vectors by the Fisher information matrix for unrelated subjects.
619 This modified variance estimator is more accurate than the original sandwich variance estimator for low-
620 frequency variants.
621

622 SUGEN can perform Wald and score tests. We used the Wald test because it yields slightly better control
623 of the type I error than the score test. SUGEN can accommodate binary, continuous, and age-at-onset
624 traits. When analyzing continuous traits, we allowed the trait variances to be different among different
625 ethnic groups.
626

627 The GEE assumption of independence between families is more restrictive than the covariance model
628 assumed in the LMM. However, the GEE approach does not rely on the normality assumption and is
629 robust to model misspecification. While the original SUGEN version had general methodology, the
630 software has been extended to handle heterogeneous variances for the PAGE analyses
631 (<https://github.com/dragontaoran/SUGEN#>). In our dataset, SUGEN provides very accurate control of the
632 type I error, as judged by the QQ-plots and genomic control parameter. We used SUGEN as the primary
633 method for association testing and ran GENESIS in parallel for comparison purposes (**Supplemental**
634 **Table 2, Supplemental Figure 5**).
635

636 **Genetic Ancestry Interactions:** Besides simple tests of association, we were also interested in whether
637 allelic effects differ according to ethnicity or ancestral makeup. These presence of these effects can be
638 estimated by adding SNPxPC interaction effects into model M1 to form
639

$$640 \quad M2: g(E(Y)) = X\alpha + G\beta + (Gx)\theta$$

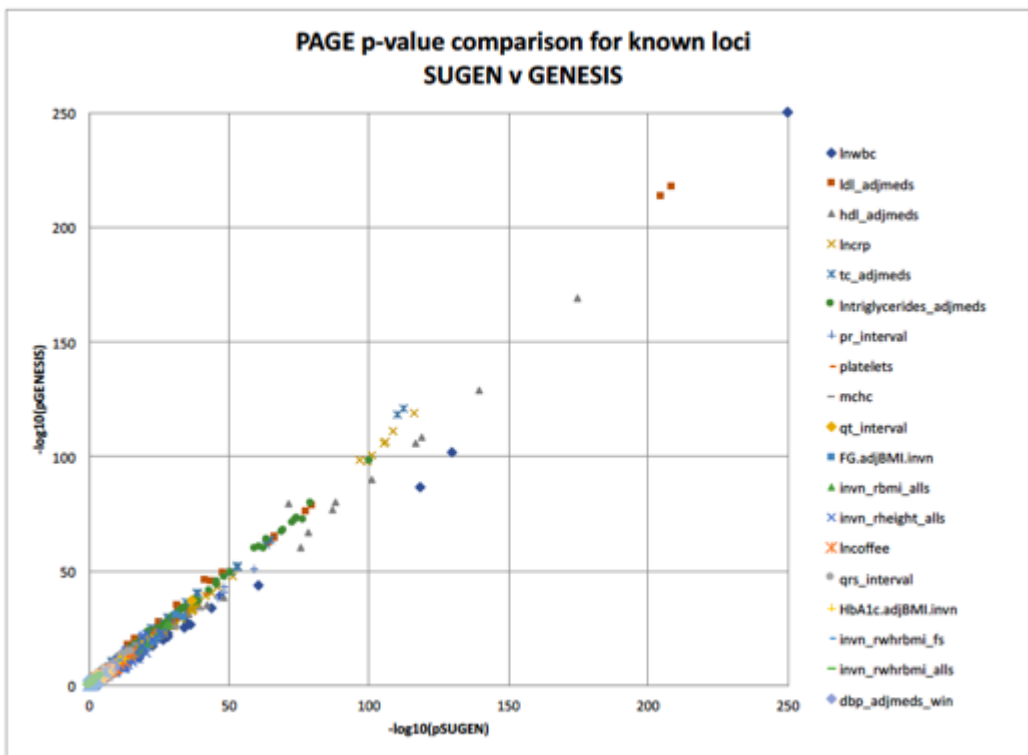
641 and then testing that $\theta = 0$, generally this will yield a multi degree of freedom test (F-test) for
642 heterogeneity of effects depending upon how many PCs are included in the hypothesis tests. We choose
643 to include interaction terms for all 10 PCs to account for the sub-continental differences that were
644 differentiated in the higher PCs. We explored using a smaller number of PCs in M2, but found that the P-
645 values for the F-test obtained with 5 PCs were extremely similar to those with 10 PCs for the vast majority
646 of SNPs (results not shown).
647
648

649 **Assessing Single Variant Results:** For each phenotype, QQ plots and genomic inflation factors
650 (λ) were used to assess inflation, using the full set of results, and results omitting previously known loci.
651 Inflation values ranged from 0.98 to 1.15 for all traits. Analyses were restricted to SNPs with an
652 imputation quality score greater than 0.4 and an effective sample size (effN) greater than 30 for
653 continuous traits, and greater than 50 for binary traits. The effN was calculated based on previous
654 publications²⁴:
655

$$656 \quad \text{effN} = 2 * \text{MAF} * (1 - \text{MAF}) * N * \text{info}$$

657

658 where MAF is the minor allele frequency among the set of individuals included in a phenotype-specific
659 model, N is the total sample size for a given phenotype, and info is the impute2 info quality score. The
660 SNP with the smallest p-value in a 1Mb region was considered the Lead SNP. A Lead SNP was
661 considered to be a Novel loci if it met the following criteria: 1) it was located greater than +/- 500 Kb away
662 from a previously known loci (per the phenotype-specific Known Loci list); 2) it had a SUGEN p-value less
663 than 5E-08; 3) it had a SUGEN conditional p-value less than 5E-08 after adjustment for all previously
664 known loci on the same chromosome; and 4) it had 2 or more neighboring SNPs (within +/- 500 Kb) with
665 a p-value less than 1E-05. A Lead SNP was considered to be a Residual signal in a previously known loci
666 if it met the following criteria: 1) it was located within +/- 500 Kb of a previously known loci; 2) it had a
667 SUGEN p-value less than 5E-08; and 3) it had a SUGEN conditional p-value less than 5E-08 after
668 adjustment for all previously known loci on the same chromosome. Full results for all Novel and Residual
669 findings are included in **Supplemental Tables 2-3**. Additionally, minor allele frequency-dependent
670 thresholds were used for genome-wide significance, as per guidelines in Fadista et al (2016).³³ For
671 common variants with MAF>5%, the standard $P < 5 \times 10^{-8}$ threshold was used to determine significance. For
672 low frequency and rare variants with MAF<5%, a more stringent $P < 3 \times 10^{-9}$ was utilized. This is reflected in
673 the 16 novel genome-wide significant trait-variant associations and the 11 low-frequency loci with
674 suggestive associations ($3 \times 10^{-9} > P > 5 \times 10^{-8}$).
675



676

677 **Supplementary Figure 12: Comparison of P-values from GWAS for SUGEN (Wald**
678 **test) vs. GENESIS across all traits. ($N_{max}=49,781$; see Extended Data Table 1)**
679

680

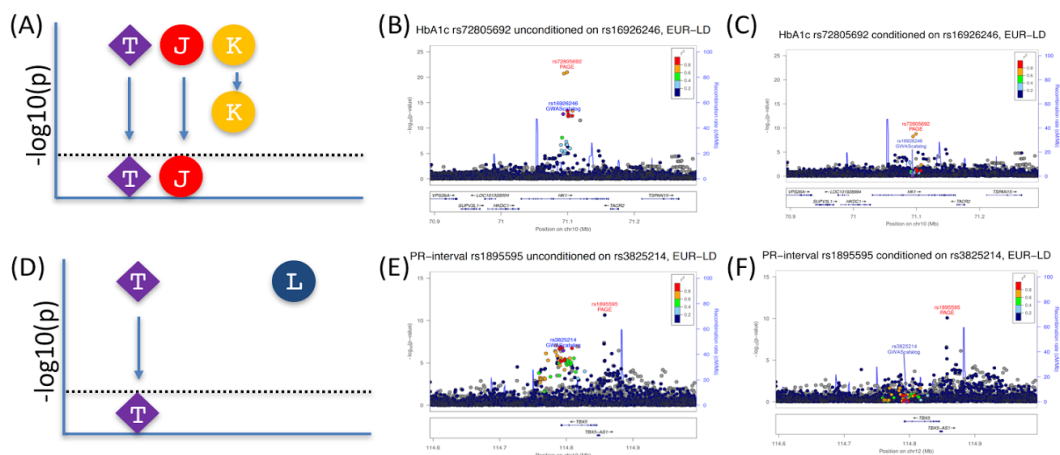
681

682

683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699

8. Secondary Signals versus Fine-mapping

To further illustrate the difference in mechanism between fine-mapping and secondary independent signals, we highlight two examples (**Supplementary Figure 12**). The first is a refinement of the association between hexokinase 1 (*HK1*) and HbA1c. The residual signal at rs72805692 ($P_{\text{unadj}}=9.22 \times 10^{-22}$, $N=11,178$, $\text{CAF}=0.061$) is in moderate LD in European ($r^2=0.61$) and Hispanic/Latino ($r^2=0.63$) populations with the previously implicated SNP (rs16926246) 5.7kb away. Therefore, after adjustment, the signal is greatly diminished but remains statistically significant ($P_{\text{cond}}=3.05 \times 10^{-9}$). This represents the refinement of a known locus (fine-mapping), as the high LD present in this area results in an attenuated, but still statistically significant, signal, and may represent only one underlying fSNP. In contrast, we found a residual signal for PR interval at rs1895595, upstream of *TBX5* ($P_{\text{unadj}}=2.16 \times 10^{-11}$, $N=17,428$, $\text{CAF}=0.17$). After adjustment for 5 known tagSNPs in this region (rs3825214, rs7312625, rs7135659, rs1895585, rs1896312), the signal remains largely unchanged ($P_{\text{cond}}=1.99 \times 10^{-11}$). This secondary signal at rs1895595 is independent of all 5 conditioned SNPs, with extremely low LD ($r^2 < 0.03$) across all global populations, and therefore likely represents an independent fSNP. Both fine-mapping of primary findings and knowledge of independent, secondary alleles are important to comprehensively characterize GWAS loci, particularly in diverse populations, thereby improving genetic risk prediction.



700

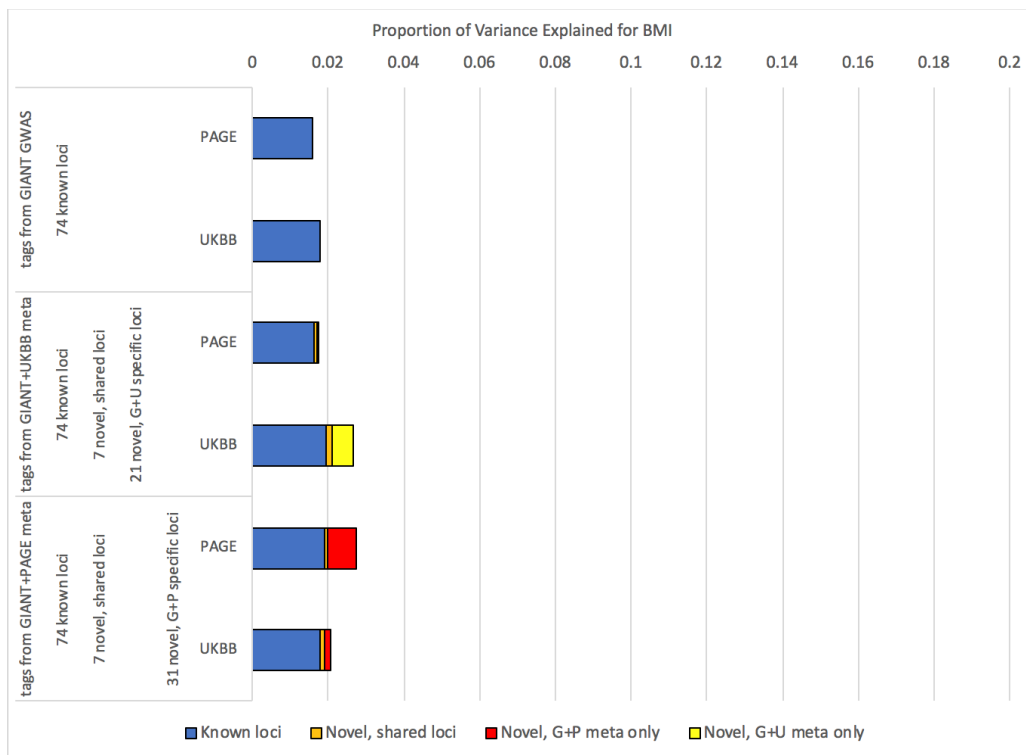
Supplementary Figure 13: Residual signals can represent either refinement of signal or secondary alleles.

(A) Fine-mapping: $-\log_{10} p$ values from SUGEN Wald test are plotted against position for a GWAS catalog tagSNP T, as well as two tagged SNPs: J is strongly tagged by T ($r^2=1$) in all populations, and K is variably tagged across populations. After adjustment, signal at T and J is no longer significant, but residual signal at K indicates that the original association has been fine-mapped. Unadjusted (B) and adjusted (C) results for trait HbA1c ($N=11,178$), showing weakened signal at residual SNP rs72805692 after adjusting for GWAS catalog tagSNP rs16926246, consistent with signal refinement. This tagSNP was first reported from a study of 46,368 Europeans³⁴, so LD with the tagSNP is shown from a European reference panel, illustrating how the set of strongly tagged SNPs (red/orange) is fine-mapped to the two strongest (residual) signals in the multi-ethnic population. (D) Secondary alleles are independent of known loci, so L is not in significant LD with T ($r^2 \sim 0$). After adjustment for T, signal at L is unchanged. Unadjusted (E) and adjusted (F) results for trait PR interval ($N=17,428$), showing no change in signal at residual SNP rs1895595 after adjusting for GWAS catalog tagSNP rs3825214, consistent with the residual signal being an independent secondary allele. Again, LD shown is from a European population, as the GWAS catalog report³⁵ was from 12,670 Europeans. P-values estimated from SUGEN Wald test.

717 **9. Meta-analysis and Finemapping with GIANT, UK**
 718 **Biobank**

719 For height³⁶, GIANT imputed their GWAS data to the HapMap in roughly 250,000 individuals, yielding 2.5M
 720 variants that overlapped with the PAGE dataset. All of these are common (MAF>5%) in at least one
 721 ancestry, so the traditional threshold of statistical significance ($P < 5 \times 10^{-8}$) is appropriate. In the GIANT BMI
 722 manuscript³⁷, the GWAS data were augmented with metabochip data (a focused platform targeting specific
 723 regions of the genome) from ~80,000 additional individuals. The previously published manuscripts used a
 724 more relaxed definition of “locus” than we have in this manuscript (associations less than 1Mbp apart were
 725 merged into a single locus, where we have used 500kbp), and also reported results from multiple analytic
 726 approaches (by subset, or conditioned on known loci). For clarity in our comparison, we use the same locus
 727 definition as we used earlier in this manuscript, and we limit the comparison to a single approach: the sex-
 728 combined joint analysis of all European individuals in GIANT, meta-analyzed with the sex-combined
 729 SUGEN results from either PAGE or the UK Biobank.

730
 731 To create a comparable sample size as PAGE, a total of 50,000 “White British” individuals were randomly
 732 subset from the larger UK Biobank (UKB50k). Height and BMI was adjusted for both sex and age identically
 733 to PAGE procedures. The linkage disequilibrium information for GIANT was generated from 9,700 ARIC
 734 individuals who were part of the larger GIANT consortium. For PAGE, correlation between sites were
 735 calculated separately for each race/ethnicity group. These matrices were then combined, weighted by the
 736 inverse sample size to create a combined weighted correlation matrix of sites. For the meta-analyses, the
 737 correlation matrices were again combined, weighted by the inverse of the sample size proportional to 250k
 738 GIANT, 50k PAGE, and 50k “White British” UK Biobank participants.
 739

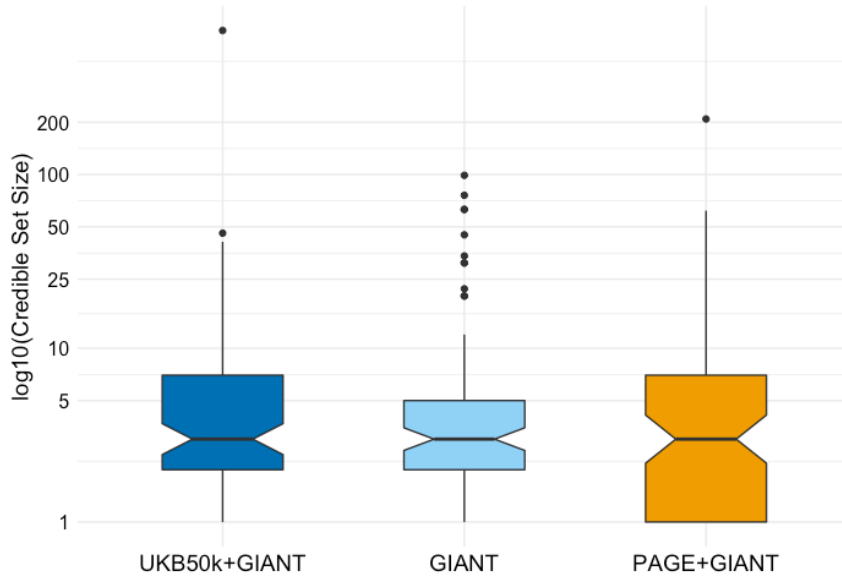


740

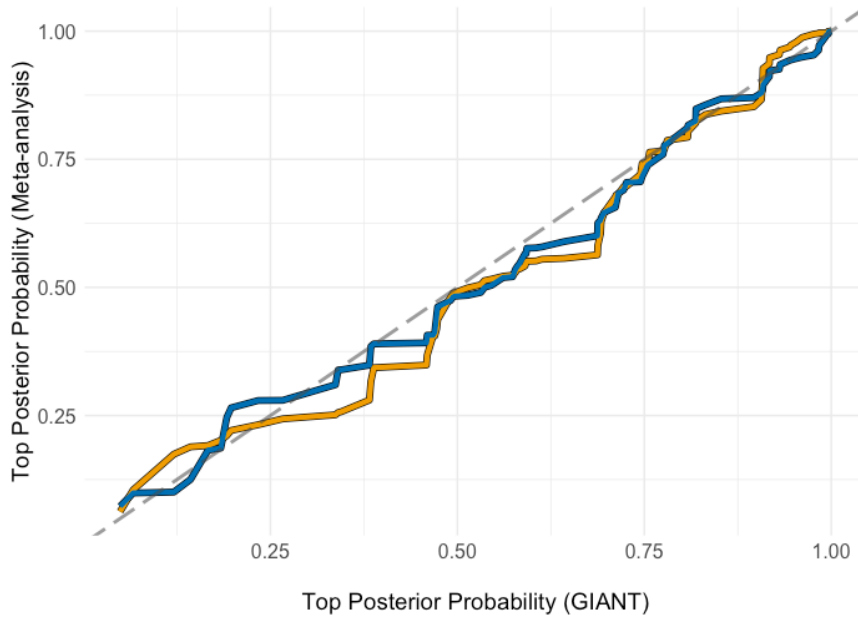
741 **Supplementary Figure 14: BMI PVE.**

742 Although less of the variance is explained for BMI than for height, results are broadly consistent: meta-
 743 analysis with more Europeans (GIANT+UKB) exacerbates existing disparities in PVE between Europeans
 744 and a multi-ethnic cohort (center pair of bars), while in this case, meta-analysis with the multi-ethnic
 745 cohort (GIANT+PAGE) actually yields improved PVE in the multi-ethnic cohort.

746 **A**



747 **B**
748



749

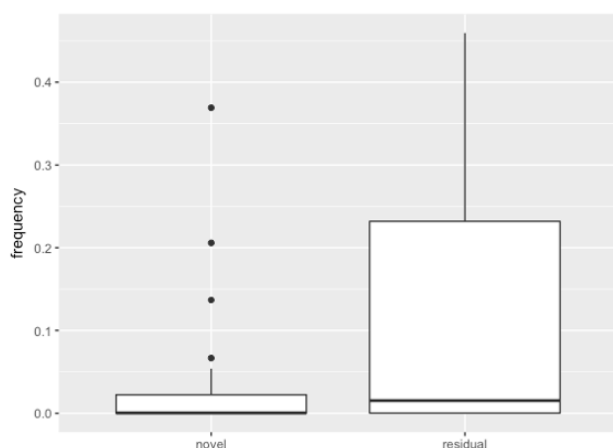
750 **Supplementary Figure 15: Finemapping for BMI.**

751 (A) Comparison of 95% credible sets for height, comparing GIANT alone (N=253,288) to UKB50k+GIANT
752 (N=303,288; paired sample t-test $P=0.60$) and PAGE+GIANT (N=303,069; paired sample t-test $P=0.50$).
753 Boxplots show the median at the notch, with the top and bottom of the box indicating the interquartile
754 range (IQR). Whiskers extend to either the minimum value or $1.5 \times \text{IQR}$. Notches indicate the 95%
755 confidence interval of the medians. (B) Top posterior probability from each 95% credible set for height,
756 comparing GIANT (N=253,288) to UKB50k+GIANT (N=303,288) and PAGE+GIANT (N=303,069).

757

758 10. Comparison of novel and secondary variant allele 759 frequencies in European populations 760

761 To interrogate the possibility of discovery of our novel and secondary findings in a European ancestry
762 sample, we downloaded the vcfs from the gnomAD browser (<https://gnomad.broadinstitute.org>) and
763 extracted the sites that matched our PAGE hits for the non-Finnish European group (NFE), as the largest
764 public repository of European-derived allele frequencies. Of these, we identified 24 novel and 35 residual
765 sites that were biallelic, did not contain repeat motifs, and within the callability mask. As can be seen in
766 the boxplots below, novel sites had significantly lower allele frequencies than secondary sites (median
767 minor allele frequency: 0.0050 vs 0.015, Wilcoxon p: 0.03). We also observed a weakly significant
768 increase in sites with a measured MAF of 0 in NFE: 7/24 vs 3/35 in novel and secondary, respectively:
769 logistic regression $P=0.05$, OR=4.4, 95% CI: 1.1-22, reflecting the small sample size of novel and
770 secondary sites. However as can be seen in the plot, even in the novel findings there are common
771 variants (maximum MAF: rs6730558, MAF~37%), indicating that lead variants still require fine mapping to
772 uncover the causal signals as described above.



773

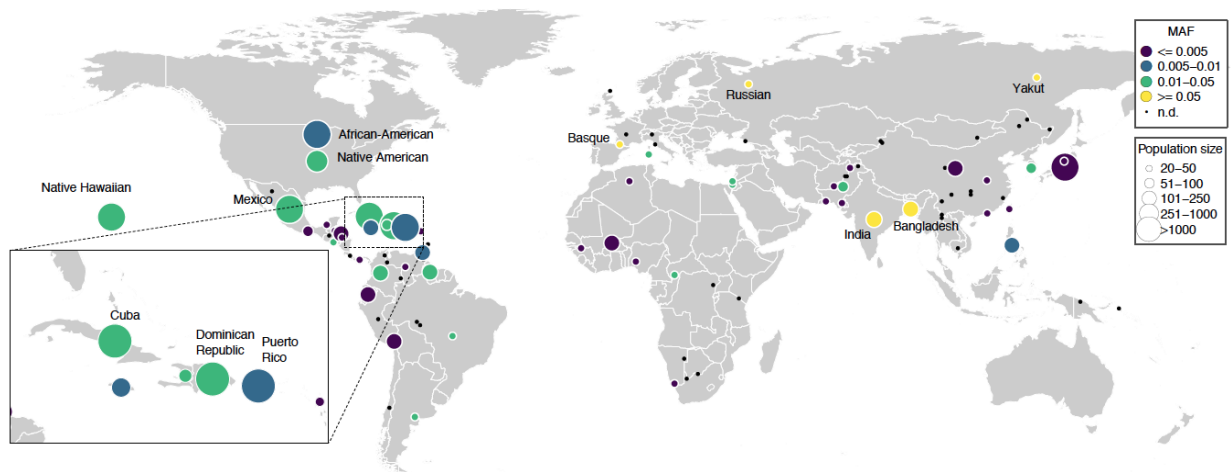
774 **Supplementary Figure 16: European allele frequencies of novel and secondary** 775 **findings in PAGE.**

776 Here we show the distribution of allele frequencies in the Non-Finnish European group in the gnomAD
777 browser (N=63,369) for our novel and secondary findings, demonstrating the preponderance of low
778 frequency variants in European populations which are now adequately powered in PAGE groups. The
779 median is denoted in bold with the top and bottom indicating the interquartile range (IQR). The whiskers
780 denote 1.5*IQR or the minimum/maximum value, with outliers displayed as dots.

781

782 **11. Clinically-relevant variants and their distribution in**
783 **PAGE**

784 We also investigated the HLA-B*57:01 allotype, which interacts with the HIV drug abacavir to trigger
785 a potentially life-threatening immune response³⁸⁻⁴¹ and therefore is recommended by the FDA for
786 screening prior to treatment initiation⁴². The rs2395029 (G) variant in *HCP5* is used to screen for abacavir
787 hypersensitivity⁴³, as it is a near perfect tag of *HLA-B*57:01* in Europeans and has utility ($r \sim 0.92$,⁴⁴)
788 across globally diverse populations in the 1000 Genomes Project. Using PAGE and Global Reference Panel
789 samples, we show that risk allele frequencies for rs2395029 rise above 5% in multiple large South Asian
790 populations, and above 1% within some admixed populations with Native American ancestry (**Figure 4**).
791 PAGE allele frequencies can therefore aid in expanding the reach of precision medicine to encompass
792 individuals of diverse ancestry, particularly when combined with other studies.^{17,45}
793
794



795
796 **Supplementary Figure 17: World map of *HCP5*-G frequencies within PAGE**
797 **groups.**

798 The histocompatibility protein variant HLA-B*57:01 interacts with the HIV drug abacavir to stimulate a
799 hypersensitivity response. A variant in a gene near HLA-B, *HCP5* rs2395029 (G allele), can be used to
800 genotype for the -B*57:01 allele as it is in high linkage disequilibrium (correlation ~ 0.92 in 1000 Genomes
801 Phase 1).^{43,44,46,47} This *HCP5* tag-SNP segregates within all continental populations of the PAGE study,
802 providing increased resolution of the global haplotype frequency, particularly within Latin America. Above,
803 minor allele (G) frequency is shown. Population size is indicated by the radius of the circle. Black dot
804 (MAF not displayed): population has less than twenty individuals or the variant is a singleton in that
805 population.

806
807

12. Additional Acknowledgements

BioMe Biobank: Samples and data of The Charles Bronfman Institute for Personalized Medicine (IPM) BioMe Biobank used in this study were provided by The Charles Bronfman Institute for Personalized Medicine at the Icahn School of Medicine at Mount Sinai (New York). Phenotype data collection was supported by The Andrea and Charles Bronfman Philanthropies.

HCHS/SOL: The following Institutes/Centers/Offices contribute to the HCHS/SOL through a transfer of funds to the NHLBI: National Institute on Minority Health and Health Disparities, National Institute on Deafness and Other Communication Disorders, National Institute of Dental and Craniofacial Research, National Institute of Diabetes and Digestive and Kidney Diseases, National Institute of Neurological Disorders and Stroke, NIH Institution-Office of Dietary Supplements.

PAGE Global Reference Panel: The specific datasets are:

Mexico: Samples of indigenous origin in Oaxaca were kindly provided by co-authors, and Samuel Canizales Quinteros and Victor Acuña Alonzo. Peru: Individuals from a primarily Quechuan and Aymaran-speaking community in Puno were kindly provided with funding support from the Burroughs Welcome Fund. Rapa Nui (Easter Island, Chile): Samples were kindly provided with funding from the Charles Rosenkranz Prize for Health Care Research in Developing Countries and the International Center for Genetic Engineering and Biotechnology (ICGEB) Grant CRP/MEX15-04_EC awarded to A.M.-E. South Africa: Samples of KhoeSan individuals from the ǀKhomani and Nama communities were kindly provided with funding from the Morrison Institute for Population and Resource Studies. Honduras and Colombia: Samples from communities in Honduras and Colombia were kindly provided by co-authors, Edwin Herrero-Paz (Universidad Católica de Honduras, San Pedro Sula, Honduras), Alvaro Mayorga (Universidad Católica de Honduras, San Pedro Sula, Honduras), Luis Caraballo (University of Cartagena), Javier Marrugo (university of Cartagena) Additional global samples: The following datasets are open access and available through the lab website of Carlos Bustamante (<https://bustamantelab.stanford.edu/>). The Human Genome Diversity Panel (HGDP-CEPH) is a group of cell lines maintained by the Centre d'Étude du Polymorphisme Humain, Fondation Jean Dausset (Paris, France) comprising 52 diverse populations across the world (Africa, Near East, Europe, South Asia, Central Asia, East Asia, Oceania and the Americas). Additional information on these datasets can be found on the CEPH website (http://www.cephb.fr/en/hgdp_panel.php), or originally at <http://www.ncbi.nlm.nih.gov/pubmed/11954565> and <http://www.ncbi.nlm.nih.gov/pubmed/12493913>, with numerous subsequent publications. Samples were filtered to include the H952 unrelated individuals as published here: <http://www.ncbi.nlm.nih.gov/pubmed/17044859>. Also available on the Bustamante Lab website is genotype data for the Maasai from Kinyawa, Kenya (MKK) samples maintained by the Coriell Institute for Medical Research (<https://catalog.coriell.org/1/NHGRI/Collections/HapMap-Collections/Maasai-in-Kinyawa-Kenya-MKK>) and genotyped as part of the International HapMap Project Phase 3 (<http://hapmap.ncbi.nlm.nih.gov/>, <http://www.sanger.ac.uk/resources/downloads/human/hapmap3.html>). We have genotyped a subset of unrelated individuals using the filters recommended in <http://www.ncbi.nlm.nih.gov/pubmed/20869033>.

WHI: The authors thank the WHI investigators and staff for their dedication, and the study participants for making the program possible. A listing of WHI investigators can be found at: <https://www.whi.org/researchers/Documents%20%20Write%20a%20Paper/WHI%20Investigator%20Short%20List.pdf>

Supplementary Information Bibliography

- 860
861
- 862 1. Gottesman, O. *et al.* The Electronic Medical Records and Genomics (eMERGE) Network: past,
863 present, and future. *Genet. Med.* **15**, 761–771 (2013).
 - 864 2. Sorlie, P. D. *et al.* Design and implementation of the Hispanic Community Health Study/Study of
865 Latinos. *Ann. Epidemiol.* **20**, 629–641 (2010).
 - 866 3. Lim, U. *et al.* Asian women have greater abdominal and visceral adiposity than Caucasian women
867 with similar body mass index. *Nutr. Diabetes* **1**, e6 (2011).
 - 868 4. Kolonel, L. N. *et al.* A multiethnic cohort in Hawaii and Los Angeles: baseline characteristics. *Am. J.*
869 *Epidemiol.* **151**, 346–357 (2000).
 - 870 5. Cann, H. M. *et al.* A human genome diversity cell line panel. *Science* **296**, 261–262 (2002).
 - 871 6. Design of the Women’s Health Initiative clinical trial and observational study. The Women’s Health
872 Initiative Study Group. *Control. Clin. Trials* **19**, 61–109 (1998).
 - 873 7. Levey, A. S. *et al.* A new equation to estimate glomerular filtration rate. *Ann. Intern. Med.* **150**, 604–
874 612 (2009).
 - 875 8. Stevens, L. A. *et al.* Evaluation of the Chronic Kidney Disease Epidemiology Collaboration equation
876 for estimating the glomerular filtration rate in multiple ethnicities. *Kidney Int.* **79**, 555–562 (2011).
 - 877 9. Kono, T. *et al.* Differential effects of cyclosporine A on ornithine decarboxylase activity induced by
878 ultraviolet-B and PUVA in mouse skin. *J. Invest. Dermatol.* **96**, 871–874 (1991).
 - 879 10. National Cholesterol Education Program (NCEP) Expert Panel on Detection, Evaluation, and
880 Treatment of High Blood Cholesterol in Adults (Adult Treatment Panel III). Third Report of the
881 National Cholesterol Education Program (NCEP) Expert Panel on Detection, Evaluation, and
882 Treatment of High Blood Cholesterol in Adults (Adult Treatment Panel III) final report. *Circulation* **106**,
883 3143–3421 (2002).
 - 884 11. Sweeney, M. E. & Johnson, R. R. Ezetimibe: an update on the mechanism of action,
885 pharmacokinetics and recent clinical trials. *Expert Opin. Drug Metab. Toxicol.* **3**, 441–450 (2007).
 - 886 12. Eaton, C. Selective cholesterol absorption inhibitors have effects on LDL-C very similar to bile
887 sequestrants. (2013).
 - 888 13. Friedewald, W. T., Levy, R. I. & Fredrickson, D. S. Estimation of the concentration of low-density
889 lipoprotein cholesterol in plasma, without use of the preparative ultracentrifuge. *Clin. Chem.* **18**, 499–
890 502 (1972).
 - 891 14. Galanter, J. M. *et al.* Genome-wide association study and admixture mapping identify different
892 asthma-associated loci in Latinos: the Genes-environments & Admixture in Latino Americans study.
893 *J. Allergy Clin. Immunol.* **134**, 295–305 (2014).
 - 894 15. Wojcik, G. L. *et al.* Imputation-Aware Tag SNP Selection To Improve Power for Large-Scale, Multi-
895 ethnic Association Studies. *G3 (Bethesda)* (2018). doi:10.1534/g3.118.200502
 - 896 16. Bien, S. A. *et al.* Strategies for enriching variant coverage in candidate disease loci on a multiethnic

- 897 genotyping array. *PLoS ONE* **11**, e0167758 (2016).
- 898 17. Henn, B. M. *et al.* Hunter-gatherer genomic diversity suggests a southern African origin for modern
899 humans. *Proc Natl Acad Sci USA* **108**, 5154–5162 (2011).
- 900 18. International HapMap 3 Consortium *et al.* Integrating common and rare genetic variation in diverse
901 human populations. *Nature* **467**, 52–58 (2010).
- 902 19. Drake, K. A. *et al.* A genome-wide association study of bronchodilator response in Latinos implicates
903 rare variants. *J. Allergy Clin. Immunol.* **133**, 370–378 (2014).
- 904 20. Uren, C. *et al.* Fine-Scale Human Population Structure in Southern Africa Reflects Ecogeographic
905 Boundaries. *Genetics* **204**, 303–314 (2016).
- 906 21. Laurie, C. C. *et al.* Quality control and quality assurance in genotypic data for genome-wide
907 association studies. *Genet. Epidemiol.* **34**, 591–602 (2010).
- 908 22. Delaneau, O., Marchini, J. & Zagury, J.-F. A linear complexity phasing method for thousands of
909 genomes. *Nat. Methods* **9**, 179–181 (2011).
- 910 23. Howie, B. N., Donnelly, P. & Marchini, J. A flexible and accurate genotype imputation method for the
911 next generation of genome-wide association studies. *PLoS Genet.* **5**, e1000529 (2009).
- 912 24. 1000 Genomes Project Consortium *et al.* A global reference for human genetic variation. *Nature* **526**,
913 68–74 (2015).
- 914 25. Zhu, X., Li, S., Cooper, R. S. & Elston, R. C. A unified association analysis approach for family and
915 unrelated samples correcting for stratification. *Am. J. Hum. Genet.* **82**, 352–365 (2008).
- 916 26. Lin, D. Y. & Zeng, D. Meta-analysis of genome-wide association studies: no efficiency gain in using
917 individual participant data. *Genet. Epidemiol.* **34**, 60–66 (2010).
- 918 27. Lin, D. Y. & Zeng, D. On the relative efficiency of using summary statistics versus individual-level
919 data in meta-analysis. *Biometrika* **97**, 321–332 (2010).
- 920 28. Zhang, J. & Stram, D. O. The role of local ancestry adjustment in association studies using admixed
921 populations. *Genet. Epidemiol.* **38**, 502–515 (2014).
- 922 29. Conomos, M. P., Miller, M. B. & Thornton, T. A. Robust inference of population structure for ancestry
923 prediction and correction of stratification in the presence of relatedness. *Genet. Epidemiol.* **39**, 276–
924 293 (2015).
- 925 30. Conomos, M. P. *et al.* Genetic diversity and association studies in US hispanic/latino populations:
926 applications in the hispanic community health study/study of latinos. *Am. J. Hum. Genet.* **98**, 165–
927 184 (2016).
- 928 31. Conomos, M. P., Reiner, A. P., Weir, B. S. & Thornton, T. A. Model-free Estimation of Recent
929 Genetic Relatedness. *Am. J. Hum. Genet.* **98**, 127–148 (2016).
- 930 32. Lin, D.-Y. *et al.* Genetic association analysis under complex survey sampling: the Hispanic
931 Community Health Study/Study of Latinos. *Am. J. Hum. Genet.* **95**, 675–688 (2014).
- 932 33. Fadista, J., Manning, A. K., Florez, J. C. & Groop, L. The (in)famous GWAS P-value threshold
933 revisited and updated for low-frequency variants. *Eur. J. Hum. Genet.* **24**, 1202–1205 (2016).

- 934 34. Soranzo, N. *et al.* Common variants at 10 genomic loci influence hemoglobin A₁(C) levels via
935 glycemic and nonglycemic pathways. *Diabetes* **59**, 3229–3239 (2010).
- 936 35. Holm, H. *et al.* Several common variants modulate heart rate, PR interval and QRS duration. *Nat.*
937 *Genet.* **42**, 117–122 (2010).
- 938 36. Wood, A. R. *et al.* Defining the role of common variation in the genomic and biological architecture of
939 adult human height. *Nat. Genet.* **46**, 1173–1186 (2014).
- 940 37. Locke, A. E. *et al.* Genetic studies of body mass index yield new insights for obesity biology. *Nature*
941 **518**, 197–206 (2015).
- 942 38. Mallal, S. *et al.* HLA-B*5701 screening for hypersensitivity to abacavir. *N. Engl. J. Med.* **358**, 568–
943 579 (2008).
- 944 39. Sousa-Pinto, B. *et al.* Pharmacogenetics of abacavir hypersensitivity: A systematic review and meta-
945 analysis of the association with HLA-B*57:01. *J. Allergy Clin. Immunol.* **136**, 1092–4.e3 (2015).
- 946 40. Hetherington, S. *et al.* Hypersensitivity reactions during therapy with the nucleoside reverse
947 transcriptase inhibitor abacavir. *Clin. Ther.* **23**, 1603–1614 (2001).
- 948 41. Illing, P. T. *et al.* Immune self-reactivity triggered by drug-modified HLA-peptide repertoire. *Nature*
949 **486**, 554–558 (2012).
- 950 42. Dean, L. in *Medical Genetics Summaries* (eds. Pratt, V., McLeod, H., Dean, L., Malheiro, A. &
951 Rubinstein, W.) (National Center for Biotechnology Information (US), 2012).
- 952 43. Martin, M. A. *et al.* Clinical Pharmacogenetics Implementation Consortium Guidelines for HLA-B
953 Genotype and Abacavir Dosing: 2014 update. *Clin. Pharmacol. Ther.* **95**, 499–500 (2014).
- 954 44. Pappas, D. J. *et al.* Significant variation between SNP-based HLA imputations in diverse populations:
955 the last mile is the hardest. *Pharmacogenomics J.* **18**, 367–376 (2018).
- 956 45. Baker, J. L., Shriner, D., Bentley, A. R. & Rotimi, C. N. Pharmacogenomic implications of the
957 evolutionary history of infectious diseases in Africa. *Pharmacogenomics J.* **17**, 112–120 (2017).
- 958 46. Colombo, S. *et al.* The HCP5 single-nucleotide polymorphism: a simple screening tool for prediction
959 of hypersensitivity reaction to abacavir. *J. Infect. Dis.* **198**, 864–867 (2008).
- 960 47. Sanchez-Giron, F. *et al.* Association of the genetic marker for abacavir hypersensitivity HLA-B*5701
961 with HCP5 rs2395029 in Mexican Mestizos. *Pharmacogenomics* **12**, 809–814 (2011).