

1 Title

2 A framework for research into continental ancestry groups of the UK Biobank

3 Authors

4 Andrei-Emil Constantinescu^{1,2,3}, Ruth E. Mitchell^{1,2}, Jie Zheng^{1,2}, Caroline J. Bull^{1,2,3},

5 Nicholas J. Timpson^{1,2}, Borko Amulic⁴, Emma E. Vincent^{1,2,3}, David A. Hughes^{1,2*}

6

7 1 MRC Integrative Epidemiology Unit at the University of Bristol, Bristol, United Kingdom.

8 2 Population Health Sciences, Bristol Medical School, University of Bristol, Bristol, United

9 Kingdom.

10 3 School of Translational Health Sciences, University of Bristol, Bristol, United Kingdom.

11 4 School of Cellular and Molecular Medicine, University of Bristol, Bristol, United

12 Kingdom.

13

14 *Corresponding author: Dr. David A. Hughes (d.a.hughes@bristol.ac.uk)

15

16 Keywords

17 Ancestry, UK Biobank, Population structure

18

19 **Declarations**

20 *Ethics approval and consent to participate*

21 UK Biobank received ethical approval from the NHS National Research Ethics Service North
22 West (11/NW/0382; 16/NW/0274) and was conducted in accordance with the Declaration of
23 Helsinki. All participants provided written informed consent before enrolment in the study.

24

25 *Consent for publication*

26 All authors consented to the publication of this work.

27

28 *Availability of data and material*

29 Genetic data from UK Biobank were made available as part of project code 15825. Analytical
30 code is available on GitHub at <https://github.com/andrewcon/popgen-biobank>.

31

32 *Competing interests*

33 None to declare.

34

35 *Funding*

36 AC acknowledges funding from a Medical Research Council PhD studentship
37 (MR/N013794/1). NJT and REM acknowledge funding from the Medical Research Council
38 (MC_UU_00011/1). NJT is the PI of the Avon Longitudinal Study of Parents and Children
39 (Medical Research Council & Wellcome Trust 217065/Z/19/Z) and is supported by the
40 University of Bristol NIHR Biomedical Research Centre (BRC-1215-2001). EEV, CJB, NJT
41 and DH acknowledge funding from the Wellcome Trust (202802/Z/16/Z). EEV, CJB and

42 NJT also acknowledge funding by the CRUK Integrative Cancer Epidemiology Programme
43 (C18281/A29019). EEV and CJB are supported by Diabetes UK (17/0005587) and the World
44 Cancer Research Fund (WCRF UK), as part of the World Cancer Research Fund
45 International grant program (IIG_2019_2009). JZ is supported by the Academy of Medical
46 Sciences (AMS) Springboard Award, the Wellcome Trust, the Government Department of
47 Business, Energy and Industrial Strategy (BEIS), the British Heart Foundation and Diabetes
48 UK (SBF006\1117). JZ is funded by the Vice-Chancellor Fellowship from the University of
49 Bristol and is supported by Shanghai Thousand Talents Program. BA acknowledges funding
50 from the Medical Research Council (MR/R02149x/1). The funders of the study had no role in
51 the study design, data collection, data analysis, data interpretation or writing of the report.

52

53 *Authors' contributions*

54 AC, DH, and REM conceived the idea for the paper. AC and DH conducted the analysis. All
55 authors contributed to the interpretation of the findings. AC and DH wrote the manuscript.
56 All authors critically revised the paper for intellectual content and approved the final version
57 of the manuscript.

58

59 *Acknowledgements*

60 We are grateful to the UK Biobank study and its participants. This research has been
61 conducted using the UK Biobank resource under Application 15825.

62

63 *Authors' information (optional)*

64 MRC Integrative Epidemiology Unit at the University of Bristol, Bristol, United Kingdom.

65 AC, REM, JZ, CJB, NJT, EEV and DH

66 Population Health Sciences, Bristol Medical School, University of Bristol, Bristol, United
67 Kingdom.

68 AC, REM, JZ, CJB, NJT, EEV and DH

69 School of Translational Health Sciences, University of Bristol, Bristol, United Kingdom.

70 AC, CJB and EEV

71 School of Cellular and Molecular Medicine, University of Bristol, Bristol, United Kingdom.

72 BA

73

74 [Abstract](#)

75 **Background**

76 The UK Biobank is a large prospective cohort, based in the United Kingdom, that has deep
77 phenotypic and genomic data on roughly a half a million individuals. Included in this
78 resource are data on approximately 78,000 individuals with “non-white British ancestry”.
79 Whilst most epidemiology studies have focused predominantly on populations of European
80 ancestry, there is an opportunity to contribute to the study of health and disease for a broader
81 segment of the population by making use of the UK Biobank’s “non-white British ancestry”
82 samples. Here we present an empirical description of the continental ancestry and population
83 structure among the individuals in this UK Biobank subset.

84 **Results**

85 Reference populations from the 1000 Genomes Project for Africa, Europe, East Asia, and
86 South Asia were used to estimate ancestry for each individual. Those with at least 80%
87 ancestry in one of these four continental ancestry groups were taken forward (N=62,484).
88 Principal component and K-means clustering analyses were used to identify and characterize
89 population structure within each ancestry group. Of the approximately 78,000 individuals in
90 the UK Biobank that are of “non-white British” ancestry, 50,685, 6,653, 2,782, and 2,364

91 individuals were associated to the European, African, South Asian, and East Asian
92 continental ancestry groups, respectively. Each continental ancestry group exhibits prominent
93 population structure that is consistent with self-reported country of birth data and geography.

94 **Conclusions**

95 Methods outlined here provide an avenue to leverage UK Biobank’s deeply phenotyped data
96 allowing researchers to maximise its potential in the study of health and disease in
97 individuals of non-white British ancestry.

98

99 **Introduction**

100 As the research community strives to understand the genetic architecture of disease
101 [1], it has increasingly realized the necessity of inclusion and diversity – of ethnically,
102 ancestrally, environmentally, and geographically diverse populations [2–5]. Not simply to
103 enhance knowledge about health and disease, but to insure health equity. Epidemiological
104 studies, including genome-wide associations studies (GWAS), have been overwhelmingly
105 conducted in European populations [2]. However, funding efforts and studies including the
106 Human Heredity and Health in Africa (H3Africa) Initiative [6], the Population Architecture
107 using Genomics and Epidemiology (PAGE) Consortium [7], Trans-Omics for Precision
108 Medicine Consortium [8], Hispanic Community Health Study / Study of Latinos (SOL) [9],
109 and the All of Us Research Program [10] are making concerted efforts to include and increase
110 the number of under-represented populations in genomic epidemiology studies.

111

112 The UK Biobank project (UKBB) has phenotypic and genomic data from a
113 prospective cohort of approximately 500,000 individuals from across the United Kingdom
114 [11,12]. It has become an outstanding resource for studies of health and disease, and genetic
115 diversity within the United Kingdom. Whilst it is made up of around 430,000 “white British

116 ancestry” individuals, as defined by UKBB, it also contains a wealth of diversity from other
117 self-described ethnicities (~78,000). This is a resource that should be utilized to help expand
118 inclusion and diversity in epidemiological studies.

119

120 The Pan-UK Biobank, or the Pan-ancestry genetic analysis of the UKBB, has
121 leveraged the diversity present in UKBB and is freely providing GWAS summary statistics
122 for over seven thousand phenotypes in six continental ancestry groups
123 (<https://pan.ukbb.broadinstitute.org>). The genetic “ancestry” groups identified by Pan-UK
124 Biobank and within our study refer to groups of individuals with a shared genetic ancestry
125 and demographic history. Studies and public resources like Pan-UK Biobank are vital to the
126 goal of increasing under-represented populations and the larger goal of describing and
127 understanding the genetic architecture of phenotypic traits and disease. However, the limited
128 information on intra-population structure and non-specific use of covariates in Pan-UK
129 Biobank GWAS models may influence association effect estimates. A description of the
130 continental diversity and population structure present in the UKBB will aid future study
131 design, methodological choice(s) and ultimately improve our understanding of how genotype
132 influences phenotype.

133

134 Here, we describe an approach to define continental ancestry groups and provide a
135 description of the structure and population differentiation within them. We define “ancestry”
136 here as genetic ancestry or the complex inheritance of one’s genetic material, but in practice
137 we will be using methodologies that use genetic similarity to identify groups of individuals
138 with high (genetic) affinity or likeness [13]. The aim is to identify relatively homogenous
139 groups of individuals that approach populations consistent with a Hardy-Weinberg model and
140 are resultantly more appropriate for many of the assumptions built into many of the methods

141 used in genomic epidemiology studies [14,15]. We leverage public data from the 1000
142 Genomes Project (1KG) [16] to provide reference populations from four, therein described,
143 superpopulations or (sub)-continental ancestry groups (CAGs) – namely, Africa (AFR),
144 Europe (EUR), South Asia (SAS), and East Asia (EAS). We note that we will refer to the
145 groupings or clusters of individuals derived by this work, not as populations, but as groups or
146 clusters of individuals. Further, the groups and clusters identified here are used as discrete
147 units, but ancestry does not have decisive boundaries and is a continuum [17–20]. The use of
148 discrete units is an analytical simplification. Finally, the overarching purpose of our study is
149 to provide a description of the population structure present in the UKBB as an aid to future
150 research investigating the health of individuals from diverse ancestries.

151

152 Results

153 *Estimations of continental ancestry*

154 Each of the 78,296 UKBB “non-white British” were included in a supervised
155 ADMIXTURE analysis to estimate a proportion of ancestry to each of African (AFR),
156 European (EUR), South Asian (SAS), and East Asian (EAS) continental ancestry groups
157 (**Figure 1**). The proportion of continental ancestry is further illustrated, for each individual,
158 within the context of UKBB population structure on principal components (PC) one and two
159 as provided by the UKBB (**Figure 2**). AFR ancestry (**Figure 2A**) runs largely parallel with
160 PC1, the major axis of variation. EUR ancestry runs at a roughly 135-degree angle (**Figure**
161 **2B**) along PC1 and PC2, while SAS (**Figure 2C**) and EAS (**Figure 2D**) ancestry run, largely,
162 along PC2. Of the approximately 78,000 UKBB samples included in the ADMIXTURE
163 analysis 50,685, 6,653, 2,782, and 2,364 individuals had 80% or more of their ancestry
164 attributed to the EUR, AFR, SAS, and EAS continental super-populations, respectively.
165 These individuals were carried forward into further analyses of population structure within

166 these continental ancestry groups (CAGs). The 80% threshold was chosen to allow some
167 error in the broader continental classification whilst also placing a limit on the complex
168 structure and admixture evaluated in these subsets. A total of 15,812 “non-white British”
169 UKBB study participants were not included in any of the four CAGs, given the methods and
170 cut-offs used here.

171

172 *Population structure within continental regions.*

173 To evaluate the level of population structure among the UKBB CAGs, we first re-
174 estimated principal components for each, while also projecting individuals from 1KG
175 populations from each super-population respectively, onto the newly derived PCs (**Figure 3,**
176 **Supplementary Table 1**). For each there is considerable overlap between UKBB individuals
177 and 1KG populations, providing some context for the diversity that is present within the
178 UKBB. In the AFR continental ancestry group principal component one distinguishes West
179 African from East African 1KG populations, while PC3 distinguishes among populations of
180 West Africa (**Figure 3A**). In the EUR continental ancestry group, the PCs and 1KG
181 populations illustrate a strong North-South axis along PC2, with a similar but less distinctive
182 trend on PC1 (**Figure 3B**). In the SAS continental ancestry group, there is a South-North
183 trend along PC1, but no remarkable pattern can be attributed to the PCs (**Figure 3C**). The
184 1KG sample populations in the EAS ancestry group appear to indicate a North-South axis
185 along PC1, and a West to East axis along PC2 (**Figure 3D**).

186

187 *K-means clustering of PCs*

188 Given that many population genetics and epidemiological analyses, such as genome-
189 wide association studies, depend on limited population structure, a common desire is to have
190 a relatively homogeneous population sample for these analyses. As such, we used an

191 unsupervised algorithm to identify groups of individuals that approach Hardy-Weinberg
192 population assumptions. To do so we performed a K-means analysis on the top PCs (see
193 Methods, **Supplementary Figure 1**), from each CAG, to identify ‘K’ subclusters or groups
194 within each. An optimum number of K-clusters was determined by a silhouette analysis (see
195 Methods, **Supplementary Figure 2**). For each CAG, using only the UKBB participants, we
196 identified seven, two, four, and three K-clusters of individuals for AFR, EUR, SAS, and
197 EAS, respectively (**Supplementary Figure 3**). However, for the EUR CAG we choose the
198 second-best K-cluster (K=6) for the remaining analyses to improve our ability to investigate
199 the utility of this analytical method to discriminate population structure (**Figure 4**).

200

201 *Country of birth*

202 To evaluate the informativeness of these K-clusters we mapped each individuals’
203 country of birth and United Nations (UN) geographic regions onto the PCs (**Figure 5** and
204 **Supplementary Figure 4-5**). These figures further illustrate the diversity and structure
205 present in the sample. Each CAG presents an observable degree of population structure, and
206 region of birth (ROB) data illustrates non-specific associations between CAGs and ROB
207 (Figure 5). For example, a large number of individuals have an East African ROB but are
208 estimated to have more than 80% of their ancestry from South Asia (Figure 5 C and G).
209 Nevertheless, ROB data illustrates structure across principal components for each CAG. Yet
210 to ascertain if there is a correlation among the K-clusters identified above and the self-
211 reported place of birth we performed a correspondence analysis for each CAG. The analyses
212 indicate a correlation between K-means clusters and the UN regions for each continent: AFR
213 (Dim1 53.29%, Dim2 41.88%), EUR (Dim1 58.25%, Dim2 28.67%), SAS (Dim1 80.00%,
214 Dim2 18.2%), EAS (Dim1 92.11%, Dim2 7.89%) (**Figure 6A**). When UN regions for a
215 smaller geographical region were substituted, namely country of birth (COB; Supplementary

216 Figures 6-9), an attenuated but correlated structure remained: AFR (Dim1 28.32%, Dim2
217 25.02%), EUR (Dim1 40.43%, Dim2 31.89%), SAS (Dim1 61.60%, Dim2 25.31%), EAS
218 (Dim1 50.49%, Dim2 49.51%) (**Figure 6B, Supplementary Figure 10**).

219

220 *Population differentiation*

221 An evaluation of the degree of population differentiation within each CAG was
222 performed by estimating F_{st} , or the fixation index between each pair of K-cluster groups and
223 1KG populations. All single-nucleotide polymorphisms (SNPs) that were included in each
224 CAG's principal component analysis were used here. An average, minimum, and maximum
225 estimate was used to summarize the distribution of estimates between pairs (**Figure 7**).
226 Relative to the population differentiation observed in the 1KG sample populations we
227 observed, on average, a small degree a population differentiation among AFR and EUR K-
228 means clusters, and larger average estimates among SAS and EAS groups. Among the UKBB
229 samples average F_{st} estimates indicate that the EAS CAG has the largest amount of
230 population differentiation with an average F_{st} of 0.0133. This is followed by SAS with an
231 average estimate of 0.0092, EUR with 0.0037, and finally AFR with the smallest average
232 estimate of 0.003. However, we note that these estimates were derived from SNPs with a
233 European ascertainment bias and as such they may not coincide with analyses using an
234 unbiased set of genetic variants.

235

236 Discussion

237 Here we present an analytical pipeline to identify individual participants of the UKBB
238 study with diverse and under-represented ancestries to be used in genomic epidemiology
239 studies. Whilst cohort studies centred in diverse geographic locations are essential for
240 elucidating the effect of environment and genotype on disease, the diversity present in deeply

241 phenotyped studies such as the UKBB should be utilized where possible. This study presents
242 a description of some of the diversity present in the UKBB. Further, the methods presented
243 here provide an approach to identify subsets of individuals to help broaden, inform, and
244 improve the relevance of genetic epidemiological studies and their findings for those of, in
245 this specific instance, a non-white British ancestry (**Figure 8**).

246 Throughout the paper, when we speak of ancestry, we are referring to “genetic
247 ancestry”, or individuals who share a demographic history [13,21,22]. They would, at the
248 population level, share a history of mutation, genetic drift, recombination, migration, natural
249 selection, environment, and culture (niche construction). As a product, they would have
250 different genetic variants, allele frequencies, and patterns of linkage disequilibrium across
251 their genomes [23–25].

252 The need to perform analyses, like association studies, separately in unique ancestral
253 populations is largely born from the need to avoid correlations between phenotype and
254 genetic ancestry, or differences in allele frequencies among populations [13,26]. For
255 example, if a disease (or environmentally influenced trait) is more frequent in ancestral
256 population ‘A’ than it is in ‘B’ and your association analysis pools these ancestral
257 populations together you may erroneously identify any allele that is more frequent in
258 population ‘A’ as a genetic variant associated with the disease. To avoid these confounding
259 issues, analyses are commonly limited to relatively homogenous populations.

260 In genome-wide association studies, the aim is to derive accurate unbiased effect
261 estimates for a genetic variant on a trait. However, the task becomes increasingly
262 challenging, as variation in genetic ancestry comes with different allele frequencies, genetic
263 backgrounds and environments [27]. Methods such as the inclusion of relatedness matrixes
264 and principal components [28–31] are used to account for cryptic relatedness and undetected,
265 fine-scale population stratification. In addition, they are also used to account for correlations

266 between phenotype and genetic ancestry [32,33]. However, are the inclusion of relatedness
267 matrixes or principal components enough to control the structure present in the CAGs
268 presented here? Or would smaller (K-means clusters) more homogenous populations be
269 better suited to epidemiological analyses, like GWAS?

270 The problems introduced by population stratification persist even in populations like
271 the “white British” subset of the UKBB, where individual genetic variants and polygenic
272 scores for individual traits can retain correlations with geography, even after correcting for
273 population structure [34,35]. Moreover, when sampling populations across Europe – where
274 genetic ancestry does mirror geography [36,37] – and meta-analysing independently run
275 GWASs [38], effect estimates appear to retain a bias introduced by population structure
276 [39,40]. These fine scale issues exemplify some of the reasons for performing separate
277 epidemiological analysis, like GWAS, for populations with deeper population
278 differentiations, i.e. unique ancestries, demographic histories, and environments.

279 The complications of population stratification and opportunities for improving health
280 outcomes for more people, even at the continental level, are precisely why a description of
281 the structure within each continental ancestry group was provided here. Namely, the structure
282 present within a CAG, as identified here, may also be too great to be properly accounted for
283 with common methodologies and may thus need to be resolved into smaller more
284 homogenous groups. At the very least, careful consideration is warranted when interpreting
285 results where CAGs are used - because structure matters [41]. The unsupervised clustering
286 performed within each CAG is not a perfect solution for identifying true “populations” – an
287 exercise that may in fact be an impractical goal – but it is a method to identify groups of
288 individuals with a more similar, homogeneous ancestry. Other techniques like uniform
289 manifold approximation and projection [42] or more explicit leveraging of self-described
290 ethnicity could help improve the identification of homogenous groups. Self-described

291 ethnicity is not a synonym for genetic ancestry though, as it is a sociocultural construct. It
292 would however help inform cultural, social, and other environmental influences – important
293 aspects of a “population” - on phenotypes and disease [22].

294 In summary, we assigned individuals to continental ancestry groups (**Figure 1 and 2**);
295 illustrated the structure present among individuals within each CAG (**Figure 3**), identified
296 unsupervised clusters or groups of individuals within each (**Figure 4**) and demonstrated that
297 those clusters have an affinity to regions and countries of birth – i.e. the K-means clusters are
298 consistent with geographic structure and isolation by distance models [43,44] (**Figure 5**).
299 Notably, each CAG presents extensive structure, inconsistent with a randomly mating
300 population, but rather with the sampling of unique, geographically distant populations. In
301 particular, East Asian, South Asian, and African CAGs have isolated, or discontinuous
302 groups of individuals in the UKBB sample, exemplified in the K-means clustering analysis
303 (Figure 4) [19,20]. For example, groups K1 and K3 in the EAS CAG (**Figure 4D**) epitomizes
304 this discontinuous structure as they correspond to individuals born on the islands of
305 Philippines and Japan, respectively (**Figure 5, Supplementary Figure 8**).

306 The methods employed here do have several limitations: First, a single 1KG
307 population was used to represent each of four continental ancestry groups evaluated – Africa,
308 Europe, South Asia, and East Asia. One population is a poor proxy for all of the variation
309 present in any one (sub)-continent. However, as the 1KG project does not have optimal
310 population coverage, including more or all the 1KG populations of a CAG would still poorly
311 represent all the variation present in a (sub)-continent and would complicate the assignment
312 of individuals to a single ancestry group. Second, our analysis was limited to four (sub-
313)continental ancestry groups, to the exclusion of the Americas (AMR, a 1KG
314 superpopulation). Populations from the Americas often have a large and varying amount of
315 recent admixture from various European and African populations [25,45–49]. As such,

316 including an AMR population in the ADMIXTURE analysis, as a reference population, could
317 confound the genetic ancestries being estimated. However, whilst we limit this study to a
318 few, broad, well characterized ancestry groups the approach presented here can be generalised
319 to other, specific ancestries.

320 Third, the UKBB Axiom array used to genotype all UKBB participants was designed
321 to optimize imputation of a European population while also including genetic variants
322 previously associated with disease and other phenotypic traits derived from studies primarily
323 conducted in European populations [11,12]. As a product, the genomic data used here will
324 have an ascertainment bias [50] that would influence imputation accuracy (although no
325 imputation data was used here), allele frequency distributions, estimates of linkage
326 disequilibrium and diversity and divergence within and among populations. Each of these
327 may influence estimations of population differentiation, principal component estimates and
328 the inferences made from them [51,52]. Specific study designs [53,54] have been made to
329 remove ascertainment bias in genotype arrays so that unbiased inferences could be made for a
330 wider range of genetic ancestries, but this was not available here.

331 Fourth, the principal components illustrated and used in the unsupervised K-means
332 clustering analyses were derived from the UKBB participants only and resultantly represents
333 the diversity (point three) and genetic ancestry found in that data set. The inclusion or use of
334 other public data sets with more numerous sample populations, that better represent regional,
335 or continental diversity will provide alternative patterns of structure. Fifth, we are limited by
336 the reference population used in the analyses. Whilst the 1KG data set shall remain an
337 essential reference panel for broad analyses like those conducted here, researchers with
338 specific continental or geographically specific research questions could strengthen and refine
339 the observations made here by including other geographically specific data sets. Finally, the
340 unsupervised K-means clustering analysis is dependent upon the number of PCs included in

341 it. Here the number of PCs chosen did have an element of subjectivity (Supplementary Figure
342 1). Whilst analytical methods are available to select a number of informative PCs [55], we
343 did not implement such methods here. Given that the K-means algorithm weights each PC
344 equally, we sought to limit the PCs included to only those with the largest proportions of
345 variance explained and not necessarily all that are analytically estimated to be informative.
346

347 Conclusions

348 The approach presented here demonstrates a method to leverage the deeply
349 phenotyped and widely used UKBB data set to help improve the inclusion and equity of
350 epidemiological studies for under-represented populations. Careful considerations must be
351 given to the diversity present within continental ancestry groups. However, given the
352 thousands of individuals present in the genetic ancestry groups identified here, the UKBB
353 data set shall prove insightful for studies of health and disease in populations beyond the
354 British Isles. While the methods presented here do not describe a perfect solution to identify
355 populations, we hope that they provide an avenue to leverage the diverse data available in
356 UKBB and a methodological platform to improve and build upon.

357

358 Methods

359 *Description of working environment*

360 All analyses were performed in a Linux environment supported by the University of
361 Bristol's Advanced Computing Research Centre (ACRC) using the following publicly
362 available software packages: Plink v1.9 and v2.0 [56,57], ADMIXTURE v1.3.0 [58,59], and
363 EIGENSOFT v8.0.0 [29,30]. In addition, bespoke scripts, analyses, and figures were run and

364 generated in the R environment using version 3.6.2 on the ACRC computer clusters and
365 version 4.0.2 (Taking Off Again) on local computers [60].

366

367 *UK Biobank data*

368 This research has been conducted using the UKBB Resource under Application
369 Number 15825, from which directly genotyped SNP data (N=784,256 SNPs) were made
370 available. It includes data for a total of 78,296 individuals identified by UKBB as “non-white
371 British” participants – our analyses were restricted to this subset. In addition to genotypic
372 data, we also acquired several variables of interest (self-described ancestry, country of birth)
373 data for this subset of individuals. 365 exclusions were made when filtering those with sex
374 chromosome mismatch and/or aneuploidy, and outliers with high genetic heterozygosity and
375 missing rates [61].

376

377 *1000 Genomes data*

378 Genetic data (v5a.20130502) from phase three of the 1KG, which includes data from
379 5 continental, or 1KG described super-populations [Europe (EUR), East Asia (EAS), South
380 Asia (SAS), Africa (AFR), and the Americas (AMR)], were used to provide reference
381 populations for admixture analyses and population structure inferences ([62]
382 <http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/>). Our analyses did not include populations from
383 the AMR superpopulation. This is to maintain a simplified analysis that avoided the
384 complicating factors of the potentially recent admixture events that occurred in the Americas.
385 Included in our analyses are five populations from 1KG super-population label: (AFR), also
386 known as the continental Africa ancestry group (1) Yoruba in Ibadan, Nigeria (YRI); (2)
387 Luhya in Webuye, Kenya (LWK); (3) Gambian in Western Division, The Gambia -
388 Mandinka (GWD); (4) Mende in Sierra Leone (MSL) and (5) Esan in Nigeria (ESN). Five

389 populations from the super-population label EUR or the continental Europe ancestry group:
390 (1) Utah residents with Northern and Western European ancestry (CEU); (2) Toscani in Italia
391 (TSI); (3) British in England and Scotland (GBR); (4) Finnish in Finland (FIN) and (5)
392 Iberian populations in Spain (IBS). Five populations from the super-population label SAS or
393 the continental South Asian ancestry group: (1) Gujarati Indian in Houston, Texas (GIH); (2)
394 Punjabi in Lahore, Pakistan (PJI); (3) Bengali in Bangladesh (BEB); (4) Sri Lankan Tamil in
395 the UK (STU) and (5) Indian Telugu in the UK (ITU). Finally, five populations from the
396 super-population label EAS or the continental East Asian ancestry group: (1) Han Chinese in
397 Beijing, China (CHB); (2) Japanese in Tokyo, Japan (JPT); (3) Han Chinese South (CHS);
398 (4) Chinese Dai in Xishuangbanna, China (CDX) and (5) Kinh in Ho Chi Minh City,
399 Vietnam (KHV).

400

401 *Merging UK Biobank and 1000 Genomes*

402 The directly genotyped data from UKBB was used to identify SNPs with the same
403 SNP identifier (RefSNP ID) present in the 1KG data set. A total of 718,711 SNPs were
404 identified with the same ID and extracted from both data sets using PLINK v2.0. The two
405 datasets were then merged using the -bmerge function in PLINK v2.0. After removing
406 problematic SNPs (e.g. multi-allelic, duplicate) in the merge step, a total of 718,487 SNPs
407 remained.

408

409 *Linkage disequilibrium pruning*

410 Prior to ancestry estimation the merged dataset was reduced to a set of independent
411 SNPs based on linkage disequilibrium (LD) estimates using the PLINK v2.0 function and
412 parameters "--indep-pairwise 50 10 0.025", indicating an r^2 threshold of 0.025, a window size
413 of 50 kilobases and a window step size of 10 kilobases. In addition, 24 previously identified

414 genomic regions with extensive linkage disequilibrium were also excluded [63,64]. LD
415 estimates in this analysis were limited to unrelated individuals from the 1KG YRI population
416 sample. A total of 30,320 SNPs remained following LD pruning.

417

418 *Estimating African, European, South Asian, and East Asian ancestry*

419 Four 1KG populations were included as reference populations in a supervised
420 Admixture (v1.3.0) analysis. They were (1) British in England and Scotland (GBR), of the
421 European ancestry (EUR) superpopulation, (2) Yoruba in Ibadan, Nigeria (YRI), of the
422 African ancestry (AFR) superpopulation, (3) Indian Telugu in the UK (ITU), of the South
423 Asian ancestry (SAS) superpopulation, and (4) Han Chinese South (CHS), of the East Asian
424 ancestry (EAS) superpopulation. These singular population samples were chosen to broadly
425 represent each of their four respective continental (superpopulation) ancestry groups, with an
426 average population differentiation (F_{st} , or fixation index) value of 0.1055 amongst them, as
427 estimated by ADMIXTURE. The supervised ADMIXTURE analysis provides, for each
428 UKBB sample, a proportion of ancestry for each of the four reference populations. Those
429 individuals with at least 80% of their ancestry attributed to one continental ancestry group, or
430 1KG defined superpopulation, were carried forward into further analyses.

431

432 *Derivation of continental principal components*

433 Unrelated individuals in each CAG, including both 1KG and UKBB samples with
434 $\geq 80\%$ ancestry to that CAG were identified (using all 718,487 SNPs in the overlapping data
435 set, and the plink (v1.9) function `--rel-cutoff` and a minor allele frequency (MAF) filter of
436 0.05 (`--maf 0.05`)). Then for each CAG and using all (1KG + UKBB) unrelated individuals
437 assigned to the CAG, a list of approximately 40 thousand LD independent SNPs were
438 identified (using the plink (v2.0) function `--indep-pairwise 50 10 0.025` (`--indep-pairwise 50`

439 10 0.02 for AFR and --indep-pairwise 50 10 0.05 for SAS) along with a MAF filter of 0.01,
440 and the exclusion of the 24 previously identified genomic regions with extensive linkage
441 disequilibrium [63,64]). New plink files including only the LD independent SNPs identified
442 in step two were subsequently generated. smartrel from the EIGENSOFT
443 (<https://github.com/DReichLab/EIG>) package was used to generate a new list of related
444 individual pairs, along with our script “greedy_unrelated_selection.R” to identify a list of
445 related individuals to exclude from principal component derivation [29,30]. An exception this
446 step was made for the European CAG as its sample-size was prohibitively large to run
447 smartrel, instead the list of unrelated individuals generated from step one was used. Finally,
448 smartpca of the EIGENSOFT package was used to estimate principal components (PC), using
449 only unrelated UKBB samples. Related and 1KG samples were subsequently projected upon
450 these PCs by smartpca. Sample outliers were excluded from the PC analysis by smartpca with
451 the following parameters: using 10 PCs to identify outliers (numoutlierevec), at six standard
452 deviations from the mean (outliersigmathresh), and with 5 outlier removal iterations
453 (numoutlieriter). **Supplementary Table 1** provides numbers for each of these steps, for each
454 CAG. The EUR CAG was treated uniquely due to its larger sample-size. Smartpca was run
455 twice as described above, once with “fastmode=NO” and then with “fastmode=YES”. The
456 former provided estimates of the eigenvalues but not the eigenvectors, while the latter
457 provided eigenvectors but not eigenvalues.

458

459 *K-means clustering of principal components*

460 For each CAG, we estimated the variance explained by each principal component
461 (PC) by dividing the eigenvalue of each PC by the sum of all eigenvalues. To identify the
462 number of top PCs we generated a scree plot, using the variance explained estimates, and
463 identified the elbow or valley in each plot (**Supplementary Figure 1, Supplementary Table**

464 2). The top PCs, and the top PCs only, were then used in an unsupervised K-means clustering
465 analysis (k set from 2 to 20; using the function “kmeans()” from the R stats package) to
466 identify clusters of UKBB individuals that maximize between cluster sums of squares and
467 minimize within cluster sums of squares. An optimum number of clusters (k) was identified
468 by silhouette analysis using the function “pamk()” from the fpc R package (**Supplementary**
469 **Figure 2**) [65]. These analyses are implemented in our function “DetermineK()” found in this
470 study’s GitHub repository.

471

472 *Correspondence analysis*

473 Each UKBB study participants’ country of birth information was placed into United
474 Nations defined geographic regions (**Supplementary Table 3**). To determine if the K-means
475 population clusters have any relationship with an individual’s country of birth or country of
476 birth UN-region we performed correspondence analyses (CAs) using the function “ca()” from
477 the R package “ca”, for each continental ancestry group [36]. In addition, a chi-square test
478 was performed on the contingency table used in the correspondence analysis. Any UN-region
479 or country of birth with fewer than 10 observations was excluded. Individuals for which
480 country of birth information was not available were also excluded.

481

482 *Population differentiation among K-means population clusters*

483 For each CAG, we took the best K-means population clusters, as defined by the
484 silhouette analysis, and re-ran smartpca. However, on this run we had smartpca provide for us
485 only an estimation of the average fixation index (Fst) for each pair of populations in the data
486 set, including 1KG populations and UKBB K-means clusters. This was done with the
487 inclusion of the parameters “fstonly” and “phylipoutname” [55], the latter of which provides
488 a distance matrix of mean Fst values between populations. Estimations of Fst, which range

489 from 0 to 1, provide a measure of population differentiation among populations. In brief,
490 these describe the proportion of total variation at a SNP that is explained by variation
491 between populations. For any SNP a value of 0 would indicate that minimal variation is
492 attributable to variation between populations. A value of 1 would indicate a fixed difference
493 i.e., the two populations are both invariable but for alternative alleles.

494 [List of abbreviations](#)

495 1KG = 1000 Genomes Project
496 ACRC = Advanced Computing Research Centre
497 AFR = African
498 AMR = Americas
499 BEB = Bengali in Bangladesh
500 CA = correspondence analysis
501 CAG = continental ancestry group
502 CDX = Chinese Dai in Xishuangbanna, China
503 CEU = Utah residents with Northern and Western European ancestry
504 CHB = Han Chinese in Beijing, China
505 CHS = Han Chinese South
506 COB = country of birth
507 EAS = East Asian
508 ESN = Esan in Nigeria
509 EUR = European
510 FIN = Finnish in Finland
511 Fst = fixation index
512 GBR = British in England and Scotland
513 GIH = Gujarati Indian in Houston, Texas

- 514 GWAS = Genome-wide association study
- 515 GWD = Gambian in Western Division, The Gambia - Mandinka
- 516 IBS = Iberian populations in Spain
- 517 ITU = Indian Telugu in the UK
- 518 JPT = Japanese in Tokyo, Japan
- 519 KHV = Kinh in Ho Chi Minh City, Vietnam
- 520 LD = linkage disequilibrium
- 521 LWK = Luhya in Webuye, Kenya
- 522 MAF = minor allele frequency
- 523 MSL = Mende in Sierra Leone
- 524 PC = principal component
- 525 PJL = Punjabi in Lahore, Pakistan
- 526 ROB = region of birth
- 527 SAS = South Asian
- 528 SNP = single-nucleotide polymorphism
- 529 STU = Sri Lankan Tamil in the UK
- 530 TSI = Toscani in Italia
- 531 UKBB = UK Biobank
- 532 UN = United Nations
- 533 YRI = Yoruba in Ibadan, Nigeria

534

535 [References](#)

- 536 [1] Timpson NJ, Greenwood CMT, Soranzo N, Lawson DJ, Richards JB. Genetic
537 architecture: the shape of the genetic contribution to human traits and disease. *Nat Rev*
538 *Genet* 2017 19:110–24. <https://doi.org/10.1038/nrg.2017.101>.

- 539 [2] Sirugo G, Williams SM, Tishkoff SA. The Missing Diversity in Human Genetic
540 Studies. *Cell* 2019;177:26–31. <https://doi.org/10.1016/J.CELL.2019.02.048>.
- 541 [3] Bentley AR, Callier SL, Rotimi CN. Evaluating the promise of inclusion of African
542 ancestry populations in genomics. *Npj Genomic Med* 2020 51 2020;5:1–9.
543 <https://doi.org/10.1038/s41525-019-0111-x>.
- 544 [4] Cooke Bailey JN, Bush WS, Crawford DC. Editorial: The Importance of Diversity in
545 Precision Medicine Research. *Front Genet* 2020;0:875.
546 <https://doi.org/10.3389/FGENE.2020.00875>.
- 547 [5] Green ED, Gunter C, Biesecker LG, Di Francesco V, Easter CL, Feingold EA, et al.
548 Strategic vision for improving human health at The Forefront of Genomics. *Nat* 2020
549 5867831 2020;586:683–92. <https://doi.org/10.1038/s41586-020-2817-4>.
- 550 [6] Consortium TH. Enabling the genomic revolution in Africa: H3Africa is developing
551 capacity for health-related genomics research in Africa. *Science* 2014;344:1346.
552 <https://doi.org/10.1126/SCIENCE.1251546>.
- 553 [7] Matisse TC, Study for the P, Ambite JL, Study for the P, Buyske S, Study for the P, et
554 al. The Next PAGE in Understanding Complex Traits: Design for the Analysis of
555 Population Architecture Using Genetics and Epidemiology (PAGE) Study. *Am J*
556 *Epidemiol* 2011;174:849–59. <https://doi.org/10.1093/AJE/KWR160>.
- 557 [8] Taliun D, Harris DN, Kessler MD, Carlson J, Szpiech ZA, Torres R, et al. Sequencing
558 of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nat* 2021 5907845
559 2021;590:290–9. <https://doi.org/10.1038/s41586-021-03205-y>.
- 560 [9] Gallo LC, Penedo FJ, Carnethon M, Isasi C, Sotres-Alvarez D, Malcarne VL, et al.
561 The Hispanic Community Health Study/Study of Latinos Sociocultural Ancillary
562 Study: Sample, Design, and Procedures. *Ethn Dis* 2014;24:77.
- 563 [10] Investigators TA of URP. The “All of Us” Research Program.

- 564 <https://doi.org/10.1056/NEJMSr1809937> 2019;381:668–76.
- 565 <https://doi.org/10.1056/NEJMSR1809937>.
- 566 [11] Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, et al. UK biobank: an
567 open access resource for identifying the causes of a wide range of complex diseases of
568 middle and old age. *PLoS Med* 2015;12:e1001779–e1001779.
569 <https://doi.org/10.1371/journal.pmed.1001779>.
- 570 [12] Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, et al. The UK
571 Biobank resource with deep phenotyping and genomic data. *Nature* 2018;562:203–9.
572 <https://doi.org/10.1038/s41586-018-0579-z>.
- 573 [13] Mathieson I, Scally A. What is ancestry? *PLOS Genet* 2020;16:e1008624.
574 <https://doi.org/10.1371/JOURNAL.PGEN.1008624>.
- 575 [14] Rodriguez S, Gaunt TR, Day INM. Hardy-Weinberg Equilibrium Testing of Biological
576 Ascertainment for Mendelian Randomization Studies. *Am J Epidemiol* 2009;169:505–
577 14. <https://doi.org/10.1093/AJE/KWN359>.
- 578 [15] Graffelman J, Weir BS. On the testing of Hardy-Weinberg proportions and equality of
579 allele frequencies in males and females at biallelic genetic markers. *Genet Epidemiol*
580 2018;42:34–48. <https://doi.org/10.1002/GEPI.22079>.
- 581 [16] Altshuler DL, Durbin RM, Abecasis GR, Bentley DR, Chakravarti A, Clark AG, et al.
582 A map of human genome variation from population-scale sequencing. *Nature*
583 2010;467:1061–73. <https://doi.org/10.1038/nature09534>.
- 584 [17] Rosenberg NA, Pritchard JK, Weber JL, Cann HM, Kidd KK, Zhivotovsky LA, et al.
585 Genetic structure of human populations. *Science* (80-) 2002;298:2381–5.
586 https://doi.org/10.1126/SCIENCE.1078311/SUPPL_FILE/ROSENBERG.SOM.PDF.P
587 DF.
- 588 [18] Berezovskiĭ ND, Giria VN. [Estimation of combining ability of specialized types of

- 589 the big white breed]. *Tsitol Genet* 1991;25:56–60.
- 590 [19] Serre D, Pääbo S. Evidence for Gradients of Human Genetic Diversity Within and
591 Among Continents. *Genome Res* 2004;14:1679. <https://doi.org/10.1101/GR.2529604>.
- 592 [20] Rosenberg NA, Mahajan S, Ramachandran S, Zhao C, Pritchard JK, Feldman MW.
593 Clines, Clusters, and the Effect of Study Design on the Inference of Human Population
594 Structure. *PLoS Genet* 2005;1:e70.
595 <https://doi.org/10.1371/JOURNAL.PGEN.0010070>.
- 596 [21] Birney E, Inouye M, Raff J, Rutherford A, Scally A. The language of race, ethnicity,
597 and ancestry in human genetic research n.d.
- 598 [22] Peterson RE, Kuchenbaecker K, Walters RK, Chen CY, Popejoy AB, Periyasamy S, et
599 al. Genome-wide Association Studies in Ancestrally Diverse Populations:
600 Opportunities, Methods, Pitfalls, and Recommendations. *Cell* 2019;179:589–603.
601 <https://doi.org/10.1016/J.CELL.2019.08.051>.
- 602 [23] Przeworski M, Wall JD. Why is there so little intragenic linkage disequilibrium in
603 humans? *Genet Res* 2001;77:143–51. <https://doi.org/10.1017/S0016672301004967>.
- 604 [24] Ptak SE, Voelpel K, Przeworski M. Insights into recombination from patterns of
605 linkage disequilibrium in humans. *Genetics* 2004;167:387.
606 <https://doi.org/10.1534/GENETICS.167.1.387>.
- 607 [25] Auton A, Abecasis GR, Altshuler DM, Durbin RM, Bentley DR, Chakravarti A, et al.
608 A global reference for human genetic variation. *Nature* 2015;526:68–74.
609 <https://doi.org/10.1038/nature15393>.
- 610 [26] Lander ES, Schork NJ. Genetic dissection of complex traits. *Science* (80-)
611 1994;265:2037–48. <https://doi.org/10.1126/SCIENCE.8091226>.
- 612 [27] Vilhjálmsson BJ, Nordborg M. The nature of confounding in genome-wide association
613 studies. *Nat Rev Genet* 2012 141 2012;14:1–2. <https://doi.org/10.1038/nrg3382>.

- 614 [28] Loh P-R, Kichaev G, Gazal S, Schoech AP, Price AL. Mixed-model association for
615 biobank-scale datasets. *Nat Genet* 2018 507 2018;50:906–8.
616 <https://doi.org/10.1038/s41588-018-0144-6>.
- 617 [29] Patterson N, Price AL, Reich D. Population structure and eigenanalysis. *PLoS Genet*
618 2006;2:2074–93. <https://doi.org/10.1371/journal.pgen.0020190>.
- 619 [30] Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal
620 components analysis corrects for stratification in genome-wide association studies. *Nat*
621 *Genet* 2006;38:904–9. <https://doi.org/10.1038/ng1847>.
- 622 [31] Loh PR, Tucker G, Bulik-Sullivan BK, Vilhjálmsson BJ, Finucane HK, Salem RM, et
623 al. Efficient Bayesian mixed-model analysis increases association power in large
624 cohorts. *Nat Genet* 2015;47:284–90. <https://doi.org/10.1038/ng.3190>.
- 625 [32] Price AL, Zaitlen NA, Reich D, Patterson N. New approaches to population
626 stratification in genome-wide association studies. *Nat Rev Genet* 2010;11:459.
627 <https://doi.org/10.1038/NRG2813>.
- 628 [33] Zaidi AA, Mathieson I. Demographic history mediates the effect of stratification on
629 polygenic scores. *Elife* 2020;9:1–30. <https://doi.org/10.7554/ELIFE.61548>.
- 630 [34] Haworth S, Mitchell R, Corbin L, Wade KH, Dudding T, Budu-Aggrey A, et al.
631 Apparent latent structure within the UK Biobank sample has implications for
632 epidemiological analysis. *Nat Commun* 2019;10. [https://doi.org/10.1038/S41467-018-](https://doi.org/10.1038/S41467-018-08219-1)
633 [08219-1](https://doi.org/10.1038/S41467-018-08219-1).
- 634 [35] Abdellaoui A, Hugh-Jones D, Yengo L, Kemper KE, Nivard MG, Veul L, et al.
635 Genetic correlates of social stratification in Great Britain. *Nat Hum Behav* 2019 312
636 2019;3:1332–42. <https://doi.org/10.1038/s41562-019-0757-5>.
- 637 [36] Lao O, Lu TT, Nothnagel M, Junge O, Freitag-Wolf S, Caliebe A, et al. Correlation
638 between Genetic and Geographic Structure in Europe. *Curr Biol* 2008;18:1241–8.

- 639 <https://doi.org/10.1016/J.CUB.2008.07.049>.
- 640 [37] Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A, et al. Genes mirror
641 geography within Europe. *Nature* 2008;456:98.
642 <https://doi.org/10.1038/NATURE07331>.
- 643 [38] Wood AR, Esko T, Yang J, Vedantam S, Pers TH, Gustafsson S, et al. Defining the
644 role of common variation in the genomic and biological architecture of adult human
645 height. *Nat Genet* 2014;46:1173. <https://doi.org/10.1038/NG.3097>.
- 646 [39] Berg JJ, Harpak A, Sinnott-Armstrong N, Joergensen AM, Mostafavi H, Field Y, et al.
647 Reduced signal for polygenic adaptation of height in UK biobank. *Elife* 2019;8.
648 <https://doi.org/10.7554/eLife.39725>.
- 649 [40] Sohail M, Maier RM, Ganna A, Bloemendal A, Martin AR, Turchin MC, et al.
650 Polygenic adaptation on height is overestimated due to uncorrected stratification in
651 genome-wide association studies. *Elife* 2019;8. <https://doi.org/10.7554/eLife.39702>.
- 652 [41] Barton N, Hermisson J, Nordborg M. Why structure matters. *Elife* 2019;8.
653 <https://doi.org/10.7554/ELIFE.45380>.
- 654 [42] Diaz-Papkovich A, Anderson-Trocme L, Ben-Eghan C, Gravel S. UMAP reveals
655 cryptic population structure and phenotype heterogeneity in large genomic cohorts.
656 *PLoS Genet* 2019;15. <https://doi.org/10.1371/journal.pgen.1008432>.
- 657 [43] Morton NE. Isolation by Distance. *Genetics* 1943;28:114.
658 <https://doi.org/10.1016/B978-0-12-374984-0.00820-2>.
- 659 [44] Slatkin M. ISOLATION BY DISTANCE IN EQUILIBRIUM AND NON-
660 EQUILIBRIUM POPULATIONS. *Evolution* 1993;47:264–79.
661 <https://doi.org/10.1111/J.1558-5646.1993.TB01215.X>.
- 662 [45] Kidd JM, Gravel S, Byrnes J, Moreno-Estrada A, Musharoff S, Bryc K, et al.
663 Population Genetic Inference from Personal Genome Data: Impact of Ancestry and

- 664 Admixture on Human Genomic Variation. *Am J Hum Genet* 2012;91:660.
665 <https://doi.org/10.1016/J.AJHG.2012.08.025>.
- 666 [46] Homburger JR, Moreno-Estrada A, Gignoux CR, Nelson D, Sanchez E, Ortiz-Tello P,
667 et al. Genomic Insights into the Ancestry and Demographic History of South America.
668 *PLoS Genet* 2015;11. <https://doi.org/10.1371/JOURNAL.PGEN.1005602>.
- 669 [47] Moreno-Estrada A, Gravel S, Zakharia F, McCauley JL, Byrnes JK, Gignoux CR, et
670 al. Reconstructing the Population Genetic History of the Caribbean. *PLoS Genet*
671 2013;9. <https://doi.org/10.1371/JOURNAL.PGEN.1003925>.
- 672 [48] Ongaro L, Scliar MO, Flores R, Raveane A, Marnetto D, Sarno S, et al. The Genomic
673 Impact of European Colonization of the Americas. *Curr Biol* 2019;29:3974-3986.e4.
674 [https://doi.org/10.1016/J.CUB.2019.09.076/ATTACHMENT/8D05D549-D774-
675 4CBA-9BE7-94D3B60AD79D/MMC3.XLSX](https://doi.org/10.1016/J.CUB.2019.09.076/ATTACHMENT/8D05D549-D774-4CBA-9BE7-94D3B60AD79D/MMC3.XLSX).
- 676 [49] Montinaro F, Busby GBJ, Pascali VL, Myers S, Hellenthal G, Capelli C. Unravelling
677 the hidden ancestry of American admixed populations. *Nat Commun* 2015;6.
678 <https://doi.org/10.1038/NCOMMS7596>.
- 679 [50] Geibel J, Reimer C, Weigend S, Weigend A, Pook T, Simianer H. How array design
680 creates SNP ascertainment bias. *PLoS One* 2021;16:e0245178–e0245178.
681 <https://doi.org/10.1371/journal.pone.0245178>.
- 682 [51] Lachance J, Tishkoff SA. SNP ascertainment bias in population genetic analyses: Why
683 it is important, and how to correct it. *BioEssays* 2013;35:780–6.
684 <https://doi.org/10.1002/bies.201300014>.
- 685 [52] Albrechtsen A, Nielsen FC, Nielsen R. Ascertainment biases in SNP chips affect
686 measures of population divergence. *Mol Biol Evol* 2010;27:2534–47.
687 <https://doi.org/10.1093/molbev/msq148>.
- 688 [53] Patterson N, Moorjani P, Luo Y, Mallick S, Rohland N, Zhan Y, et al. Ancient

- 689 Admixture in Human History. *Genetics* 2012;192:1065.
690 <https://doi.org/10.1534/GENETICS.112.145037>.
- 691 [54] Lu Y, Patterson N, Zhan Y, Mallick S, Reich D. Technical design document for a SNP
692 array that is optimized for population genetics n.d.
- 693 [55] Reich D, Thangaraj K, Patterson N, Price AL, Singh L. Reconstructing Indian
694 population history. *Nature* 2009;461:489–94. <https://doi.org/10.1038/nature08365>.
- 695 [56] Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al.
696 PLINK: A tool set for whole-genome association and population-based linkage
697 analyses. *Am J Hum Genet* 2007;81:559–75. <https://doi.org/10.1086/519795>.
- 698 [57] Chang CC, Chow CC, Tellier LCAM, Vattikuti S, Purcell SM, Lee JJ. Second-
699 generation PLINK: Rising to the challenge of larger and richer datasets. *Gigascience*
700 2015;4. <https://doi.org/10.1186/s13742-015-0047-8>.
- 701 [58] Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in
702 unrelated individuals. *Genome Res* 2009;19:1655–64.
703 <https://doi.org/10.1101/gr.094052.109>.
- 704 [59] Alexander DH, Lange K. Enhancements to the ADMIXTURE algorithm for individual
705 ancestry estimation. *BMC Bioinformatics* 2011;12:246. [https://doi.org/10.1186/1471-](https://doi.org/10.1186/1471-2105-12-246)
706 [2105-12-246](https://doi.org/10.1186/1471-2105-12-246).
- 707 [60] Core R Team. R: A Language and Environment for Statistical Computing. *R Found*
708 *Stat Comput* 2019;2:[https://www.R--project.org](https://www.R-project.org). <http://www.r-project.org> (accessed
709 March 2, 2021).
- 710 [61] Mitchell RE, Hemani G, Dudding T, Corbin L, Harrison S, Paternoster L. UK Biobank
711 Genetic Data: MRC-IEU Quality Control, version 2, 18/01/2019 n.d.
- 712 [62] Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et al.
713 An integrated map of structural variation in 2,504 human genomes. *Nature*

- 714 2015;526:75–81. <https://doi.org/10.1038/nature15394>.
- 715 [63] Price AL, Weale ME, Patterson N, Myers SR, Need AC, Shianna K V., et al. Long-
716 Range LD Can Confound Genome Scans in Admixed Populations. *Am J Hum Genet*
717 2008;83:132. <https://doi.org/10.1016/J.AJHG.2008.06.005>.
- 718 [64] Weale ME. Quality Control for Genome-Wide Association Studies. In: Barnes MR,
719 Breen G, editors. *Genet. Var. Methods Protoc.*, Humana Press, New York, NY; 2010,
720 p. 31.
- 721 [65] Batool F, Hennig C. Clustering with the Average Silhouette Width. *Comput Stat Data*
722 *Anal* 2021;158:107190. <https://doi.org/10.1016/J.CSDA.2021.107190>.

723

724

725 **Figure Legends**

726

727 **Figure 1 Ancestry estimates for the UKBB non-white British subset:** Estimates of
728 ancestry proportions for each UKBB participant previously labeled as non-white British
729 individuals by UKBB. Ancestry was derived from a supervised ADMIXTURE analysis using
730 four 1000 Genomes reference populations - Yoruba in Ibadan, Nigeria for (AFR) Africa,
731 British in England, and Scotland for (EUR) Europe, Indian Telugu in the UK for (SAS) South
732 Asia, and Han Chinese South for (EAS) East Asia.

733

734 **Figure 2 Ancestry proportions on UKBB PCs:** Continental (A) African, (B) European, (C)
735 South Asian, and (D) East Asian ancestry proportions placed on principal components one
736 and two, as supplied by the UK Biobank.

737

738 **Figure 3 UKBB continental PCs with 1000 Genomes populations:** Principal components
739 one through four for each CAG. UKBB samples are colored in gray, while the 1KG sub-
740 populations for each CAG are plotted in other colors, as indicated by each legend. The
741 proportion of variation explained by each PC is indicated on each axis.

742

743 **Figure 4 UKBB continental PCs with K-means clusters:** Principal components one
744 through four for each CAG with each individual colored by its assigned K-means population
745 cluster, as indicated by each legend. The proportion of variation explained by each PC is
746 indicated on each axis.

747

748 **Figure 5 Principal components for CAG with geographic regions of birth:** Principal
749 components one and two for each CAG, with (A-D) individuals colored by their region of
750 birth (A-D), and with (E-H) the PC center also colored by region of birth. PC centers were
751 estimated as the average PC1 and PC2 values for all individuals of that ROB. Regions of
752 birth are denoted in the figure legend, and the proportion of variation explained by each PC is
753 indicated on each axis.

754

755 **Figure 6 Correspondence analysis:** Correspondence plots between (A) K-means population
756 clusters (colored circles) and regions of birth (grey squares), and (B) K-means population
757 clusters (colored circles) and country of birth (grey squares) (B). The x and y axes are the
758 first and second dimension of each correspondence analysis, respectively, with the proportion
759 of variance explained indicated in the parentheses of each axis.

760

761 **Figure 7 Fst estimates:** The minimum, mean and maximum fixation index values for each
762 CAG in the 1KG project and the UK Biobank dataset. Fst values in the 1KG project are

763 between the sub-populations of each super-population, while UK Biobank estimates are
764 derived between K-means population cluster of each CAG.

765

766 **Figure 8 Graph outlining the possible effects of geographic structure in population**
767 **genetics:** Suppose one might want to use Mendelian randomization to study the relationship
768 between neutrophil count and severe malaria caused by *P. Falciparum* – a disease largely
769 absent in European environments. Using summary statistics from a neutrophil count GWAS
770 derived from individuals with European ancestry (Box 1A) may affect estimates due to
771 geographic structure (Ancestry + Demography + Environment). This can be overcome by
772 running a GWAS in people of African ancestry (Box 1B).

773

774 **Supplementary Figure 1 Continental ancestry PCA Scree plots:** Legend: Scree plots
775 illustrating the proportion of variation explained by each of the top 20 PCs, in each UKBB
776 continental ancestry principal component analyses. The Scree plots were used to identify the
777 number of top PCs to carry forward into the K-means clustering analysis. The continental
778 ancestry supergroups are Africa (AFR), Europe (EUR), South Asia (SAS), and East Asia
779 (EAS). The number of PCs selected as top PCs are AFR = 4, SAS = 5, EAS = 4, EUR = 5.
780 The horizontal line in each plot denotes where 10% variance explained is in each plot to aid
781 in inter-CAG comparisons.

782

783 **Supplementary Figure 2 K-means k selection with silhouette analysis:** Selection of an
784 optimum number of k clusters in the K-means analysis of the top PCs, by silhouette analysis.
785 A silhouette plot for each UKBB continental supergroup (AFR) African, (EUR) European,
786 (SAS) South Asian, and (EAS) East Asian is provided. The x-axis indicates the number of k
787 clusters evaluated, and the y-axis provides an estimate of the average silhouette width

788 (ASW). ASW is an estimation of cluster quality, or intra- and inter- cluster distances derived
789 from a partitioning around medoids (PAM). The optimum number of k clusters in each
790 UKBB continental supergroup were identified as AFR = 7, EUR = 2, SAS = 4, and EAS = 3.

791

792 **Supplementary Figure 3 UKBB continental ancestry group PCs with K-means clusters:**

793 UK Biobank continental ancestry group PCs with K-means population clusters color coded:

794 AFR (A), EUR (B), SAS (C), EAS (D).

795

796 **Supplementary Figure 4 Population structure by UN defined geographic region: UK**

797 Biobank continental ancestry group PCs colored by the region of birth color coded: AFR (A-

798 C), EUR (D-F), SAS (G-I), EAS (J-L).

799

800 **Supplementary Figure 5 Population structure centers, as defined by UN geographic**

801 **region:** UK Biobank continental ancestry group PCs 1-4 with region of birth centers

802 (averaged across all individuals from each ROB) colour coded: AFR (A-C), EUR (D-F), SAS

803 (G-I), EAS (J-L).

804

805 **Supplementary Figure 6 Population structure by country of birth in AFR by region: UK**

806 Biobank continental ancestry group PCs 1-4 for the AFR CAG divided by UN regions of

807 birth: Eastern (A), Central (B), Western (C), Northern/Southern (D). Samples are color coded

808 by their country of birth.

809

810 **Supplementary Figure 7 Population structure by country of birth in EUR by region:**

811 UK Biobank continental ancestry group PCs 1-4 for the EUR CAG divided by UN regions of

812 birth: Northern (A), Eastern (B), Southern (C), Western (D). Samples are color coded by their
813 country of birth.

814

815 **Supplementary Figure 8 Population structure by country of birth in SAS and EAS: UK**

816 Biobank continental ancestry group PCs 1-4 for the SAS and EAS CAGs divided by UN

817 regions of birth: Southern Asia (A), East Asia and South-eastern Asia (B). Samples are color

818 coded by their country of birth.

819

820 **Supplementary Figure 9 Population structure centers by country of birth: UK Biobank**

821 continental ancestry group centers colored by the country of birth: AFR (A-C), EUR (D-F),

822 SAS (G-I), EAS (J-L).

823

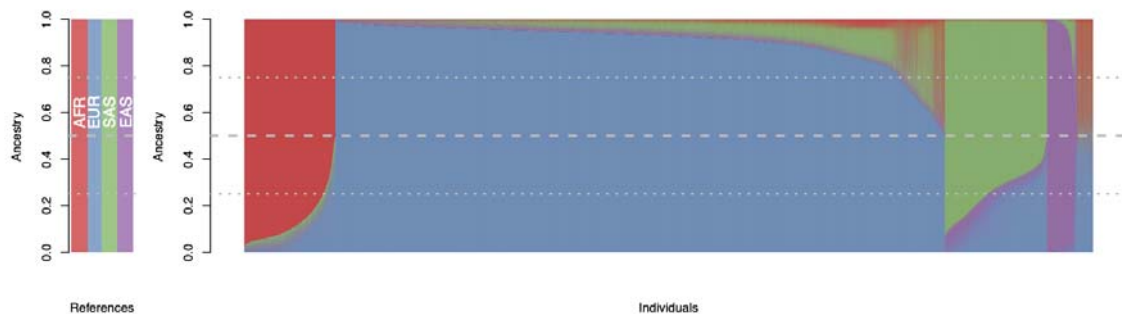
824 **Supplementary Figure 10 Population structure centers by country of birth: UK Biobank**

825 continental ancestry group centers in the correspondence analysis, colored by K-means

826 population clusters, overlapping with country of birth data in grey: AFR (A-C), EUR (D-F),

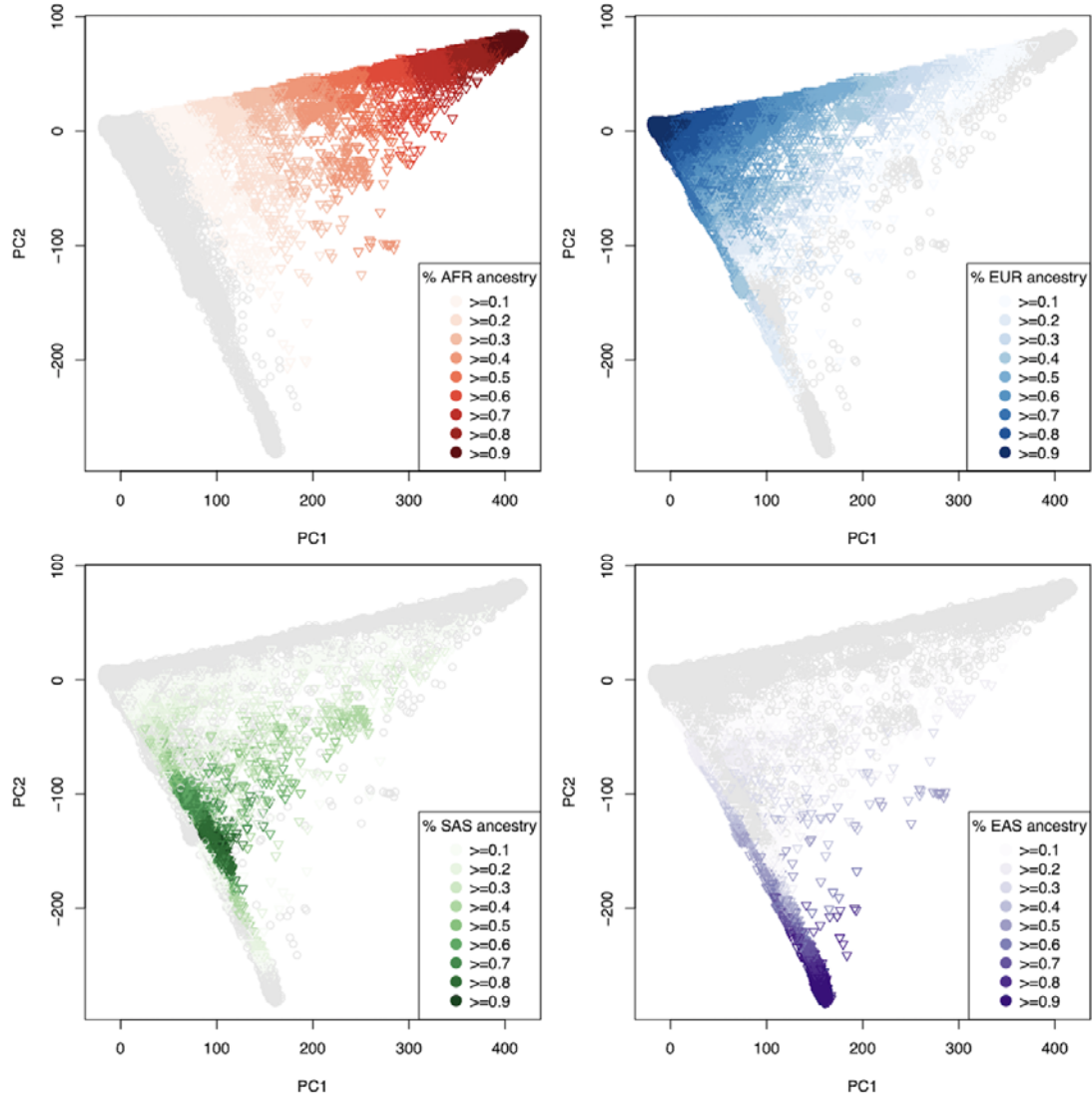
827 SAS (G-I), EAS (J-L).

828 **Figure 1 Ancestry estimates for the UKBB non-white British subset**



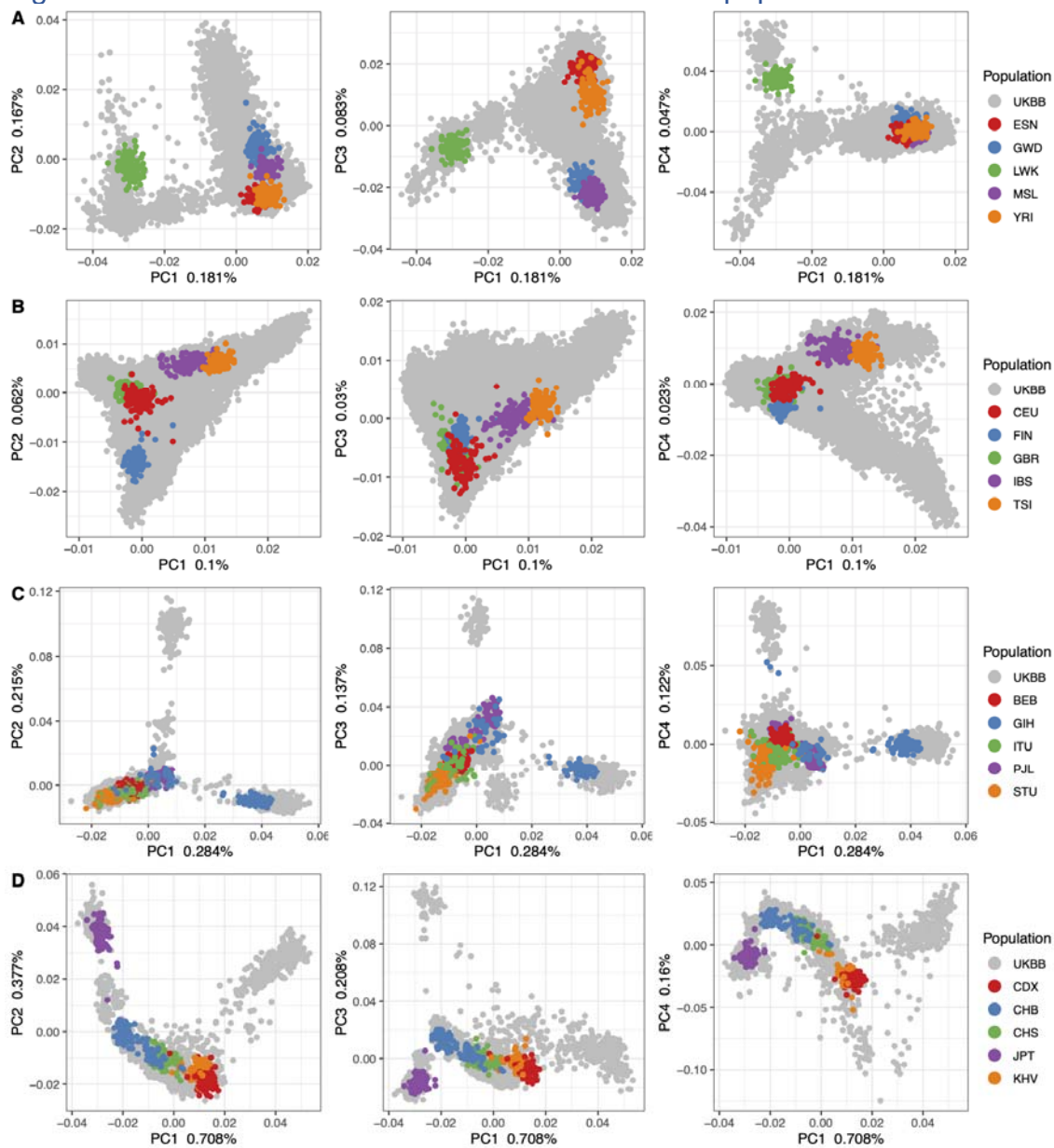
829
830

831 Figure 2: Ancestry proportions on UKBB PCs



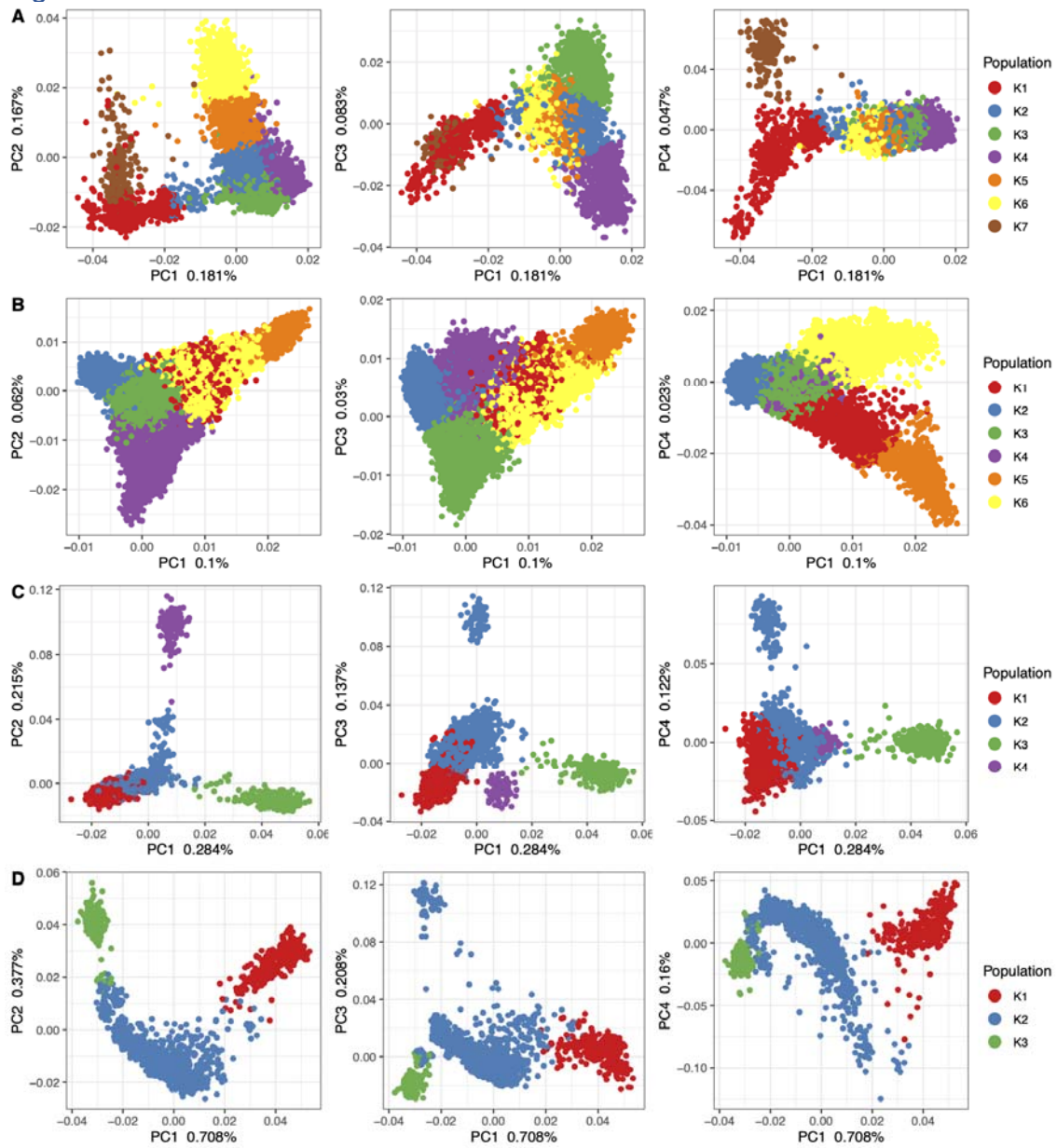
832
833

834 Figure 3: UKBB continental PCs with 1000 Genomes populations



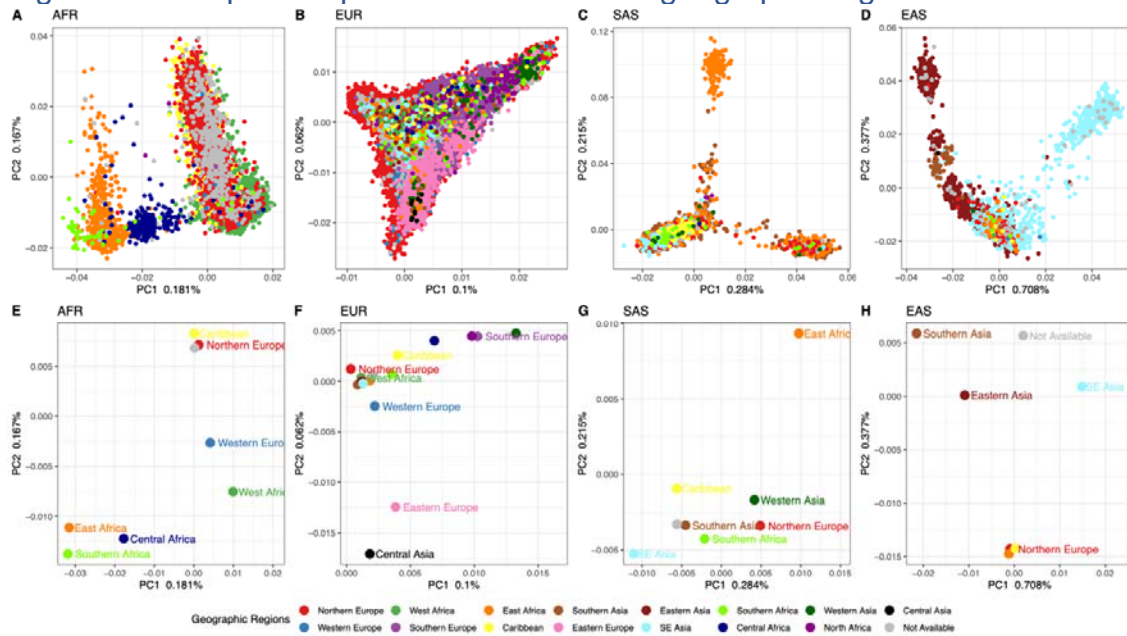
835

836 Figure 4: UKBB continental PCs with K-means clusters



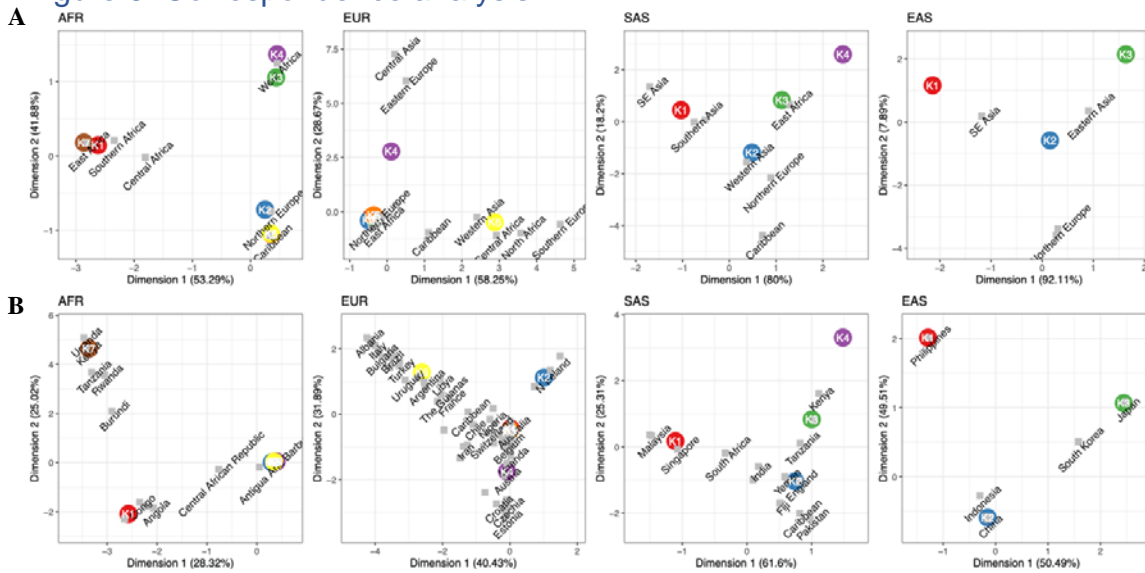
837
838

839 Figure 5: Principal components for CAG with geographic regions of birth.



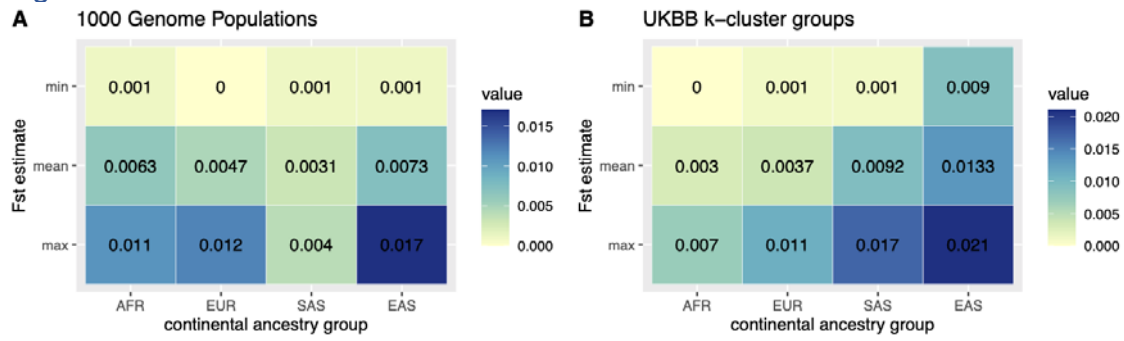
840
841
842
843

Figure 6: Correspondence analysis



844
845

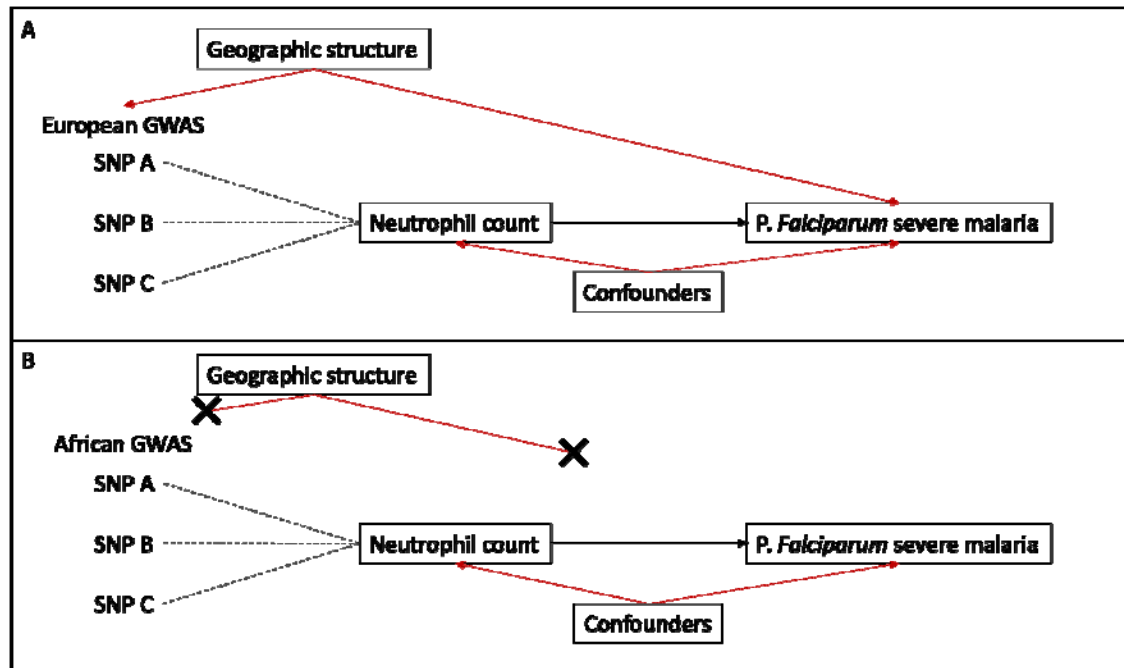
846 **Figure 7: Fst estimates**



847
848
849

850 **Figure 8: Graph outlining the possible effects of geographic structure in**
851 **population genetics**

852



853
854