



Constantinescu, A-E., Mitchell, R. E., Zheng, J., Bull, C. J., Timpson, N. J., Amulic, B., Vincent, E. E., & Hughes, D. A. (2022). A framework for research into continental ancestry groups of the UK Biobank. *Human Genomics*, 16(1), [3]. <https://doi.org/10.1186/s40246-022-00380-5>

Publisher's PDF, also known as Version of record

License (if available):  
CC BY

Link to published version (if available):  
[10.1186/s40246-022-00380-5](https://doi.org/10.1186/s40246-022-00380-5)

[Link to publication record in Explore Bristol Research](#)  
PDF-document

This is the final published version of the article (version of record). It first appeared online via BMC at <https://doi.org/10.1186/s40246-022-00380-5>. Please refer to any applicable terms of use of the publisher.

## University of Bristol - Explore Bristol Research

### General rights

This document is made available in accordance with publisher policies. Please cite only the published version using the reference above. Full terms of use are available: <http://www.bristol.ac.uk/red/research-policy/pure/user-guides/ebr-terms/>

PRIMARY RESEARCH

Open Access



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

Andrei-Emil Constantinescu<sup>1,2,3</sup>, Ruth E. Mitchell<sup>1,2</sup>, Jie Zheng<sup>1,2</sup>, Caroline J. Bull<sup>1,2,3</sup>, Nicholas J. Timpson<sup>1,2</sup>, Borko Amulic<sup>4</sup>, Emma E. Vincent<sup>1,2,3</sup> and David A. Hughes<sup>1,2\*</sup> 

## Abstract

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

**Results:** Reference populations from the 1000 Genomes Project for Africa, Europe, East Asia, and South Asia were used to estimate ancestry for each individual. Those with at least 80% ancestry in one of these four continental ancestry groups were taken forward ( $N = 62,484$ ). Principal component and K-means clustering analyses were used to identify and characterize population structure within each ancestry group. Of the approximately 78,000 individuals in the UK Biobank that are of “non-white British” ancestry, 50,685, 6653, 2782, and 2364 individuals were associated to the European, African, South Asian, and East Asian continental ancestry groups, respectively. Each continental ancestry group exhibits prominent population structure that is consistent with self-reported country of birth data and geography.

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

**Keywords:** Ancestry, UK Biobank, Population structure

## Introduction

As the research community strives to understand the genetic architecture of disease [1], it has increasingly realized the necessity of inclusion and diversity—of ethnically, ancestrally, environmentally, and geographically diverse populations [2–5], not simply to enhance knowledge about health and disease, but to insure health equity. Epidemiological studies, including genome-wide associations studies (GWAS), have been overwhelmingly conducted in European populations [2]. However,

funding efforts and studies including the Human Heredity and Health in Africa (H3Africa) Initiative [6], the Population Architecture using Genomics and Epidemiology (PAGE) Consortium [7], Trans-Omics for Precision Medicine Consortium [8], Hispanic Community Health Study / Study of Latinos (SOL) [9], and the All of Us Research Program [10] are making concerted efforts to include and increase the number of under-represented populations in genomic epidemiology studies.

The UK Biobank project (UKBB) has phenotypic and genomic data from a prospective cohort of approximately 500,000 individuals from across the UK [11, 12]. It has become an outstanding resource for studies of health and disease, and genetic diversity within the UK. While it is

\*Correspondence: d.a.hughes@bristol.ac.uk

<sup>1</sup> MRC Integrative Epidemiology Unit at the University of Bristol, Bristol, UK  
Full list of author information is available at the end of the article



© The Author(s) 2022. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article’s Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article’s Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

made up of around 430,000 “white British ancestry” individuals, as defined by UKBB, it also contains a wealth of diversity from other self-described ethnicities (~78,000). This is a resource that should be utilized to help expand inclusion and diversity in epidemiological studies.

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

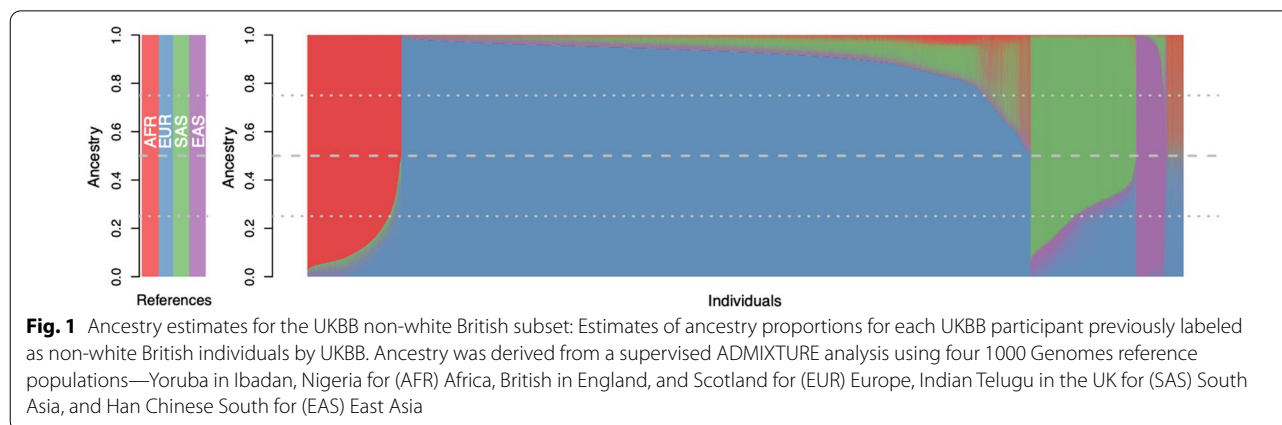
Here, we describe an approach to define continental ancestry groups and provide a description of the structure and population differentiation within them. We define “ancestry” here as genetic ancestry or the complex inheritance of one’s genetic material, but in practice we will be using methodologies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

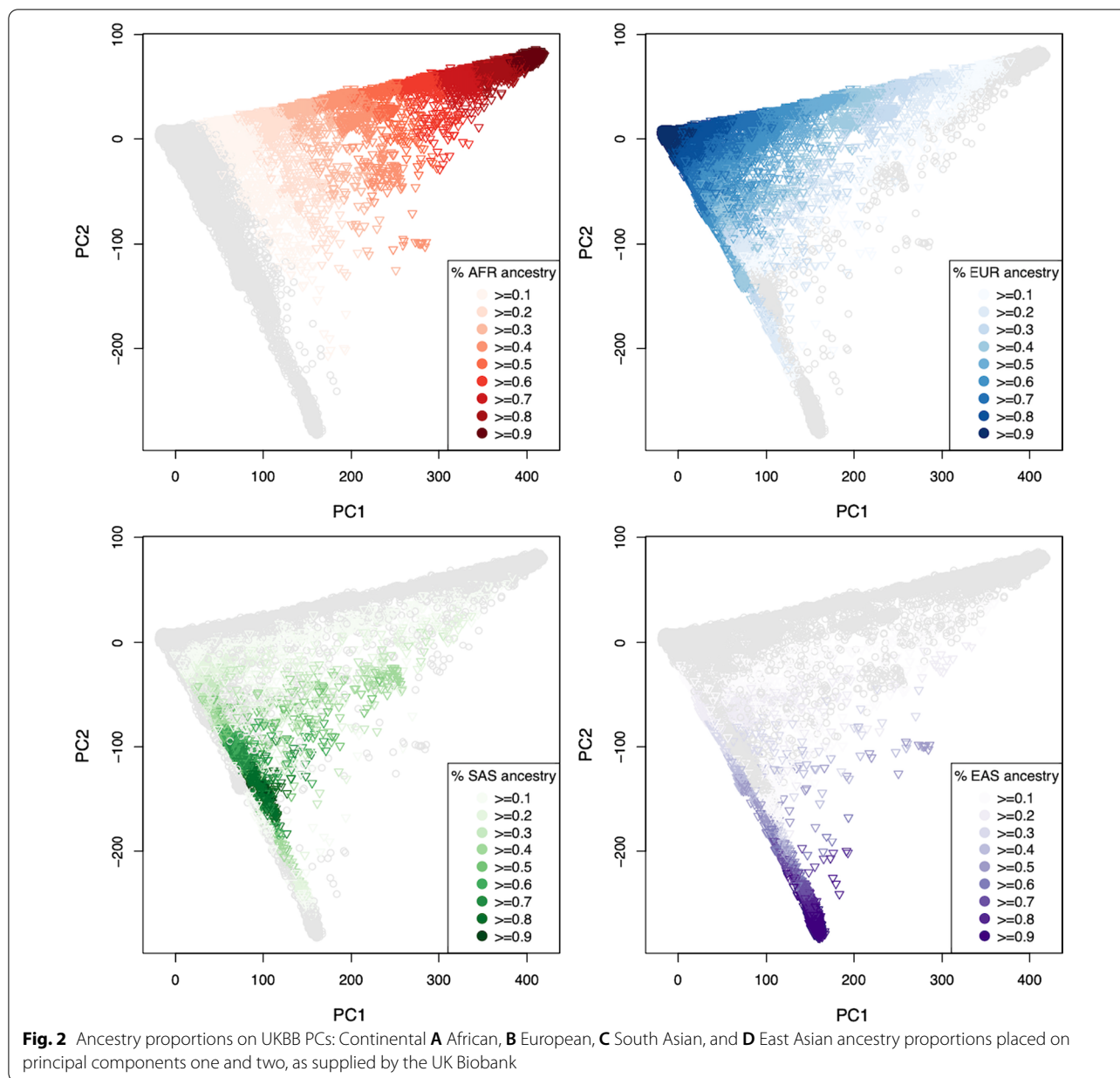
studies [14, 15]. We leverage public data from the 1000 Genomes Project (1KG) [16] to provide reference populations from four, therein described, super-populations or (sub)-continental ancestry groups (CAGs)—namely, Africa (AFR), Europe (EUR), South Asia (SAS), and East Asia (EAS). We note that we will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

## Results

### Estimations of continental ancestry

Each of the 78,296 UKBB “non-white British” was included in a supervised ADMIXTURE analysis to estimate a proportion of ancestry to each of African (AFR), European (EUR), South Asian (SAS), and East Asian (EAS) continental ancestry groups (Fig. 1). The proportion of continental ancestry is further illustrated, for each individual, within the context of UKBB population structure on principal components (PC) one and two as provided by the UKBB (Fig. 2). AFR ancestry (Fig. 2A) runs largely parallel with PC1, the major axis of variation. EUR ancestry runs at a roughly 135-degree angle (Fig. 2B) along PC1 and PC2, while SAS (Fig. 2C) and EAS (Fig. 2D) ancestry run, largely, along PC2. Of the approximately 78,000 UKBB samples included in the ADMIXTURE analysis 50,685, 6653, 2782, and 2364 individuals had 80% or more of their ancestry attributed to the EUR, AFR, SAS, and EAS continental super-populations, respectively. These individuals were carried forward into



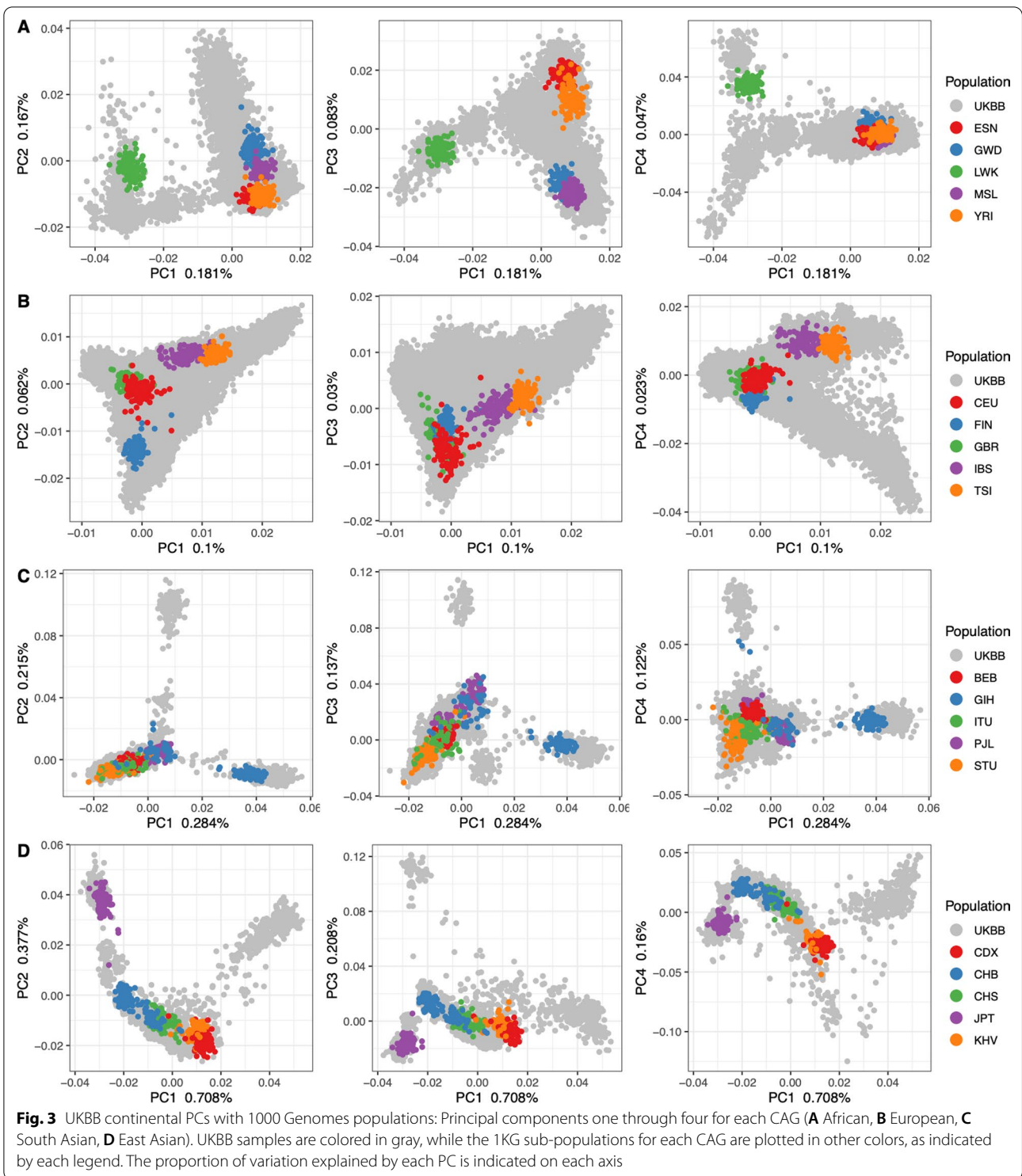


further analyses of population structure within these continental ancestry groups (CAGs). The 80% threshold was chosen to allow some error in the broader continental classification while also placing a limit on the complex structure and admixture evaluated in these subsets. A total of 15,812 “non-white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

**Population structure within continental regions**

To evaluate the level of population structure among the UKBB CAGs, we first re-estimated principal components

for each, while also projecting individuals from 1KG populations from each super-population, respectively, onto the newly derived PCs (Fig. 3, Additional file 1: Table S1). For each, there is considerable overlap between UKBB individuals and 1KG populations, providing some context for the diversity that is present within the UKBB. In the AFR continental ancestry group principal component one distinguishes West African from East African 1KG populations, while PC3 distinguishes among populations of West Africa (Fig. 3A). In the EUR continental ancestry group, the PCs and 1KG populations illustrate a strong North–South axis along PC2, with a similar but less



distinctive trend on PC1 (Fig. 3B). In the SAS continental ancestry group, there is a South-North trend along PC1, but no remarkable pattern can be attributed to the PCs (Fig. 3C). The 1KG sample populations in the EAS

ancestry group appear to indicate a North–South axis along PC1, and a West to East axis along PC2 (Fig. 3D).

### K-means clustering of PCs

Given that many population genetics and epidemiological analyses, such as genome-wide association studies, depend on limited population structure, a common desire is to have a relatively homogeneous population sample for these analyses. As such, we used an unsupervised algorithm to identify groups of individuals that approach Hardy–Weinberg population assumptions. To do so, we performed a K-means analysis on the top PCs (see Methods, Additional file 2: Fig. S1), from each CAG, to identify “K” subclusters or groups within each. An optimum number of K-clusters were determined by a silhouette analysis (see Methods, Additional file 2: Fig. S2). For each CAG, using only the UKBB participants, we identified seven, two, four, and three K-clusters of individuals for AFR, EUR, SAS, and EAS, respectively (Additional file 2: Fig. S3). However, for the EUR CAG we chose the second-best K-cluster ( $K=6$ ) for the remaining analyses to improve our ability to investigate the utility of this analytical method to discriminate population structure (Fig. 4).

### Country of birth

To evaluate the informativeness of these K-clusters, we mapped each individuals’ country of birth and United Nations (UN) geographic regions onto the PCs (Fig. 5 and Additional file 2: Figures S4–S5). These figures further illustrate the diversity and structure present in the sample. Each CAG presents an observable degree of population structure, and region of birth (ROB) data illustrate non-specific associations between CAGs and ROB (Fig. 5). For example, a large number of individuals have an East African ROB but are estimated to have more than 80% of their ancestry from South Asia (Fig. 5C and G). Nevertheless, ROB data illustrate structure across principal components for each CAG. Yet to ascertain if there is a correlation among the K-clusters identified above and the self-reported place of birth we performed a correspondence analysis for each CAG. The analyses indicate a correlation between K-means clusters and the UN regions for each continent: AFR (Dim1 53.29%, Dim2 41.88%), EUR (Dim1 58.25%, Dim2 28.67%), SAS (Dim1 80.00%, Dim2 18.2%), EAS (Dim1 92.11%, Dim2 7.89%) (Fig. 6A). When UN regions for a smaller geographical region were substituted, namely country of birth (COB; Additional file 2: Figs. S6–S9), an attenuated but correlated structure remained: AFR (Dim1 28.32%, Dim2 25.02%), EUR (Dim1 40.43%, Dim2 31.89%), SAS (Dim1 61.60%, Dim2 25.31%), EAS (Dim1 50.49%, Dim2 49.51%) (Fig. 6B, Additional file 2: Fig. S10).

### Population differentiation

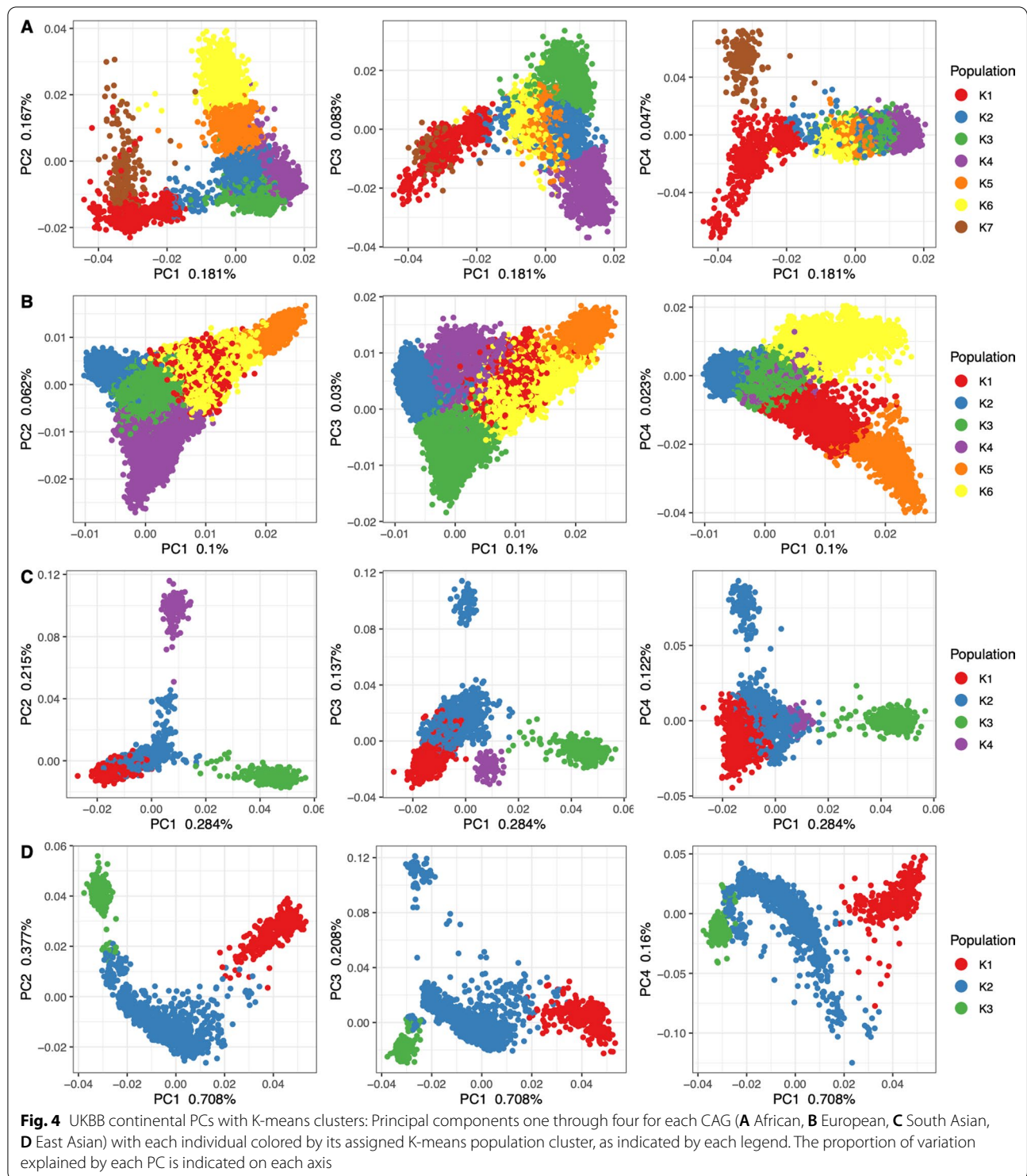
An evaluation of the degree of population differentiation within each CAG was performed by estimating  $F_{st}$ , or the fixation index between each pair of K-cluster groups and 1KG populations. All single-nucleotide polymorphisms (SNPs) that were included in each CAG’s principal component analysis were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

### Discussion

Here, we present an analytical pipeline to identify individual participants of the UKBB study with diverse and under-represented ancestries to be used in genomic epidemiology studies. While cohort studies centered in diverse geographic locations are essential for elucidating the effect of environment and genotype on disease, the diversity present in deeply phenotyped studies such as the UKBB should be utilized where possible. This study presents a description of some of the diversity present in the UKBB. Further, the methods presented here provide an approach to identify subsets of individuals to help broaden, inform, and improve the relevance of genetic epidemiological studies and their findings for those of, in this specific instance, a non-white British ancestry (Fig. 8).

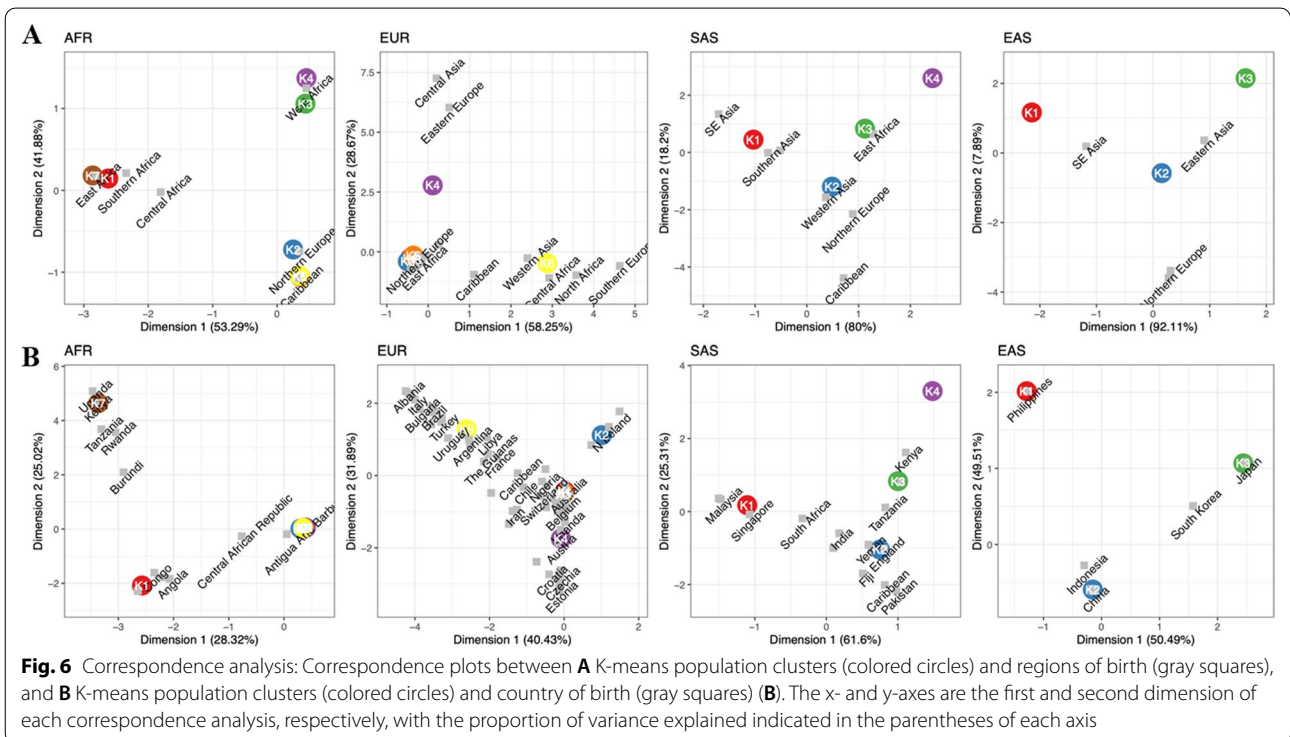
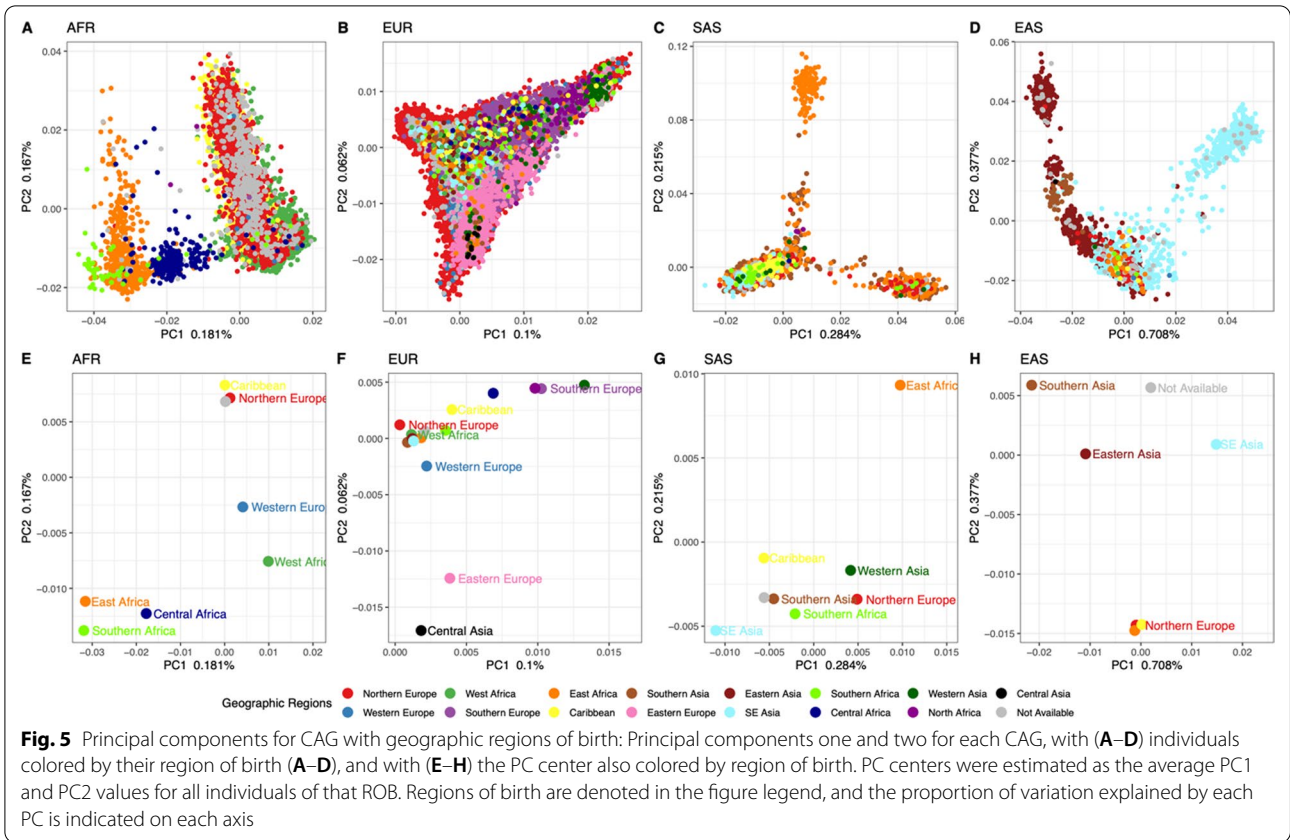
Throughout the paper, when we speak of ancestry, we are referring to “genetic ancestry,” or individuals who share a demographic history [13, 21, 22]. They should, at the population level, share a history of mutation, genetic drift, recombination, migration, natural selection, environment, and culture (niche construction [23]). As a product, they should have different genetic variants, allele frequencies, and patterns of linkage disequilibrium across their genomes [24–26].

The need to perform analyses like association studies, separately in unique ancestral populations, largely comes from the need to avoid correlations between

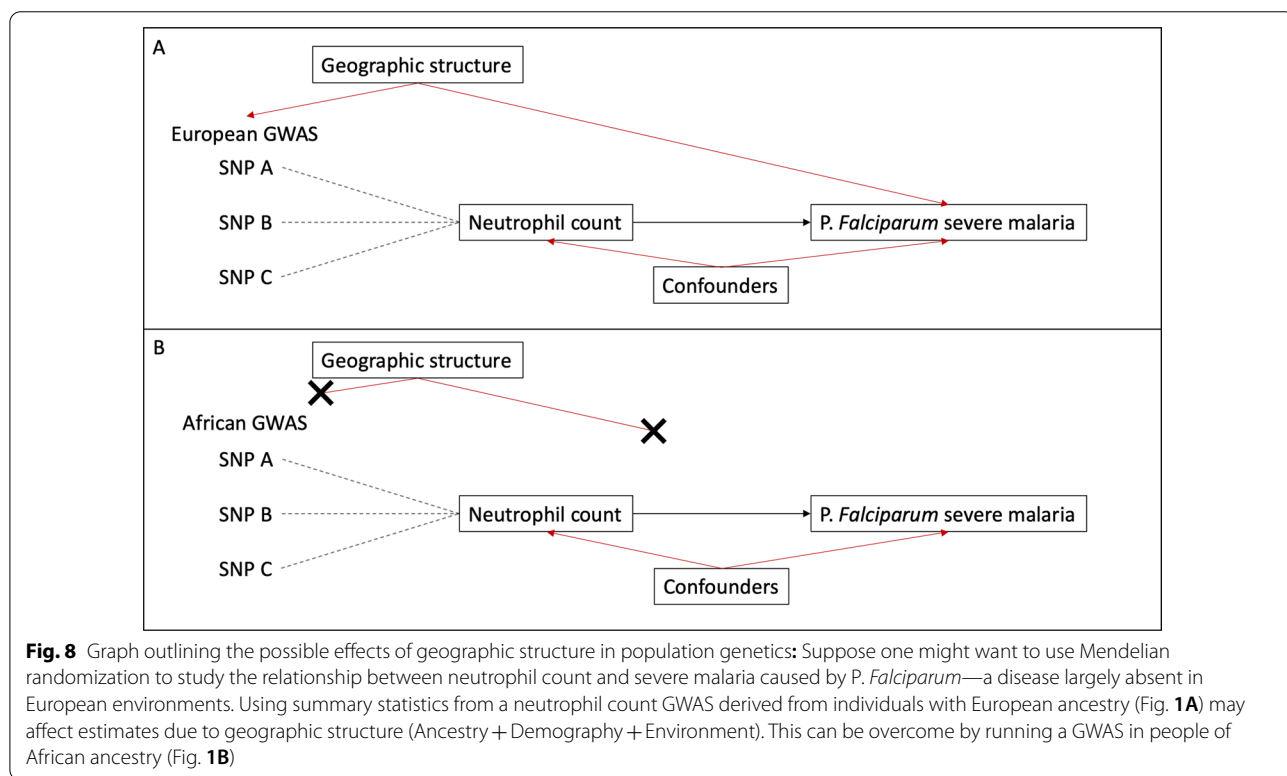
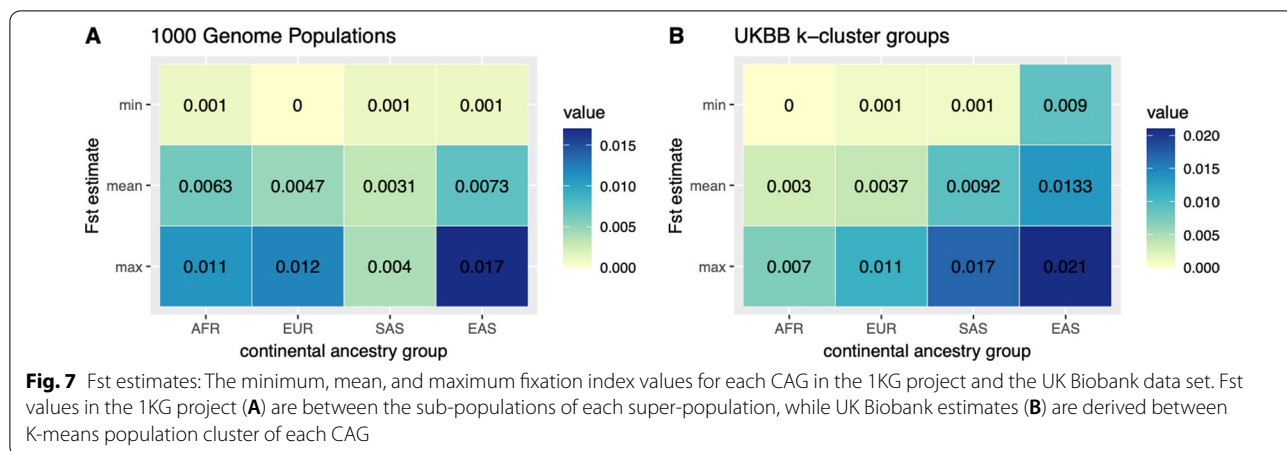


phenotype and genetic ancestry, or differences in allele frequencies among populations—i.e., population structure or population stratification [13, 27, 28]. For

example, if a disease (or environmentally influenced trait) is more frequent in ancestral population “A” than it is in “B” and if your association analysis pools these ancestral populations together you may erroneously







identify any allele that is more frequent in population “A” as a genetic variant associated with the disease. To avoid these confounding issues, analyses are commonly limited to relatively homogenous populations.

In genome-wide association studies, the aim is to derive accurate unbiased effect estimates for a genetic variant on a trait. However, the task becomes increasingly challenging, as variation in genetic ancestry comes with different allele frequencies, genetic backgrounds, and environments [29]. Methods such as the inclusion of

relatedness matrixes and principal components [30–33] are used to account for cryptic relatedness and undetected, fine-scale population stratification. In addition, they are also used to account for correlations between phenotype and genetic ancestry [34, 35]. However, is the inclusion of relatedness matrixes or principal components enough to control the structure present in the CAGs presented here? Or would smaller (K-means clusters) more homogenous populations be better suited to epidemiological analyses, like GWAS?

The problems introduced by population stratification persist even in populations like the “white British” subset of the UKBB, where individual genetic variants and polygenic scores for individual traits can retain correlations with geography, even after correcting for population structure [36, 37]. Moreover, when sampling populations across Europe—where genetic ancestry does mirror geography [38, 39]—and meta-analyzing independently run GWASs [40], effect estimates appear to retain a bias introduced by population structure [41, 42]. These fine-scale issues exemplify some of the reasons for performing separate epidemiological analysis, like GWAS, for populations with deeper population differentiations, i.e., unique ancestries, demographic histories, and environments. Other challenges and opportunities of population structure in biobank scale data are discussed further in Lawson et al. [43].

The complications of population stratification and opportunities for improving health outcomes for more people, even at the continental level, are precisely why a description of the structure within each continental ancestry group was provided here. Namely, the structure present within a CAG, as identified here, may also be too great to be properly accounted for with common methodologies and may thus need to be resolved into smaller more homogenous groups. At the very least, careful consideration is warranted when interpreting results where CAGs are used—because structure matters [44]. The unsupervised clustering performed within each CAG is not a perfect solution for identifying true “populations”—an exercise that may in fact be an impractical goal—but it is a method to identify groups of individuals with a more similar, homogeneous ancestry. Other techniques like uniform manifold approximation and projection [45] or more explicit leveraging of self-described ethnicity could help improve the identification of homogenous groups. Self-described ethnicity is not a synonym for genetic ancestry though, as it is a sociocultural construct. It would, however, help inform cultural, social, and other environmental influences—important aspects of a “population”—on phenotypes and disease [22].

In summary, we assigned individuals to continental ancestry groups (Figs. 1 and 2); illustrated the structure present among individuals within each CAG (Fig. 3); identified unsupervised clusters or groups of individuals within each (Fig. 4); and demonstrated that those clusters have an affinity to regions and countries of birth—i.e., the K-means clusters are consistent with geographic structure and isolation by distance models [46, 47] (Fig. 5). Notably, each CAG presents extensive structure, inconsistent with a randomly mating population, but rather with the sampling of unique, geographically distant populations. In particular, East Asian, South Asian, and

African CAGs have isolated, or discontinuous groups of individuals in the UKBB sample, exemplified in the K-means clustering analysis (Fig. 4) [19, 20]. For example, groups K1 and K3 in the EAS CAG (Fig. 4D) epitomize this discontinuous structure as they correspond to individuals born on the islands of Philippines and Japan, respectively (Fig. 5, Additional file 2: Fig. S8).

The methods employed here do have several limitations: First, a single 1KG population was used to represent each of four continental ancestry groups evaluated—Africa, Europe, South Asia, and East Asia. One population is a poor proxy for all of the variation present in any one (sub)-continent. However, as the 1KG project does not have optimal population coverage, including more or all the 1KG populations of a CAG would still poorly represent all the variation present in a (sub)-continent and would complicate the assignment of individuals to a single ancestry group. Second, our analysis was limited to four (sub)-continental ancestry groups, to the exclusion of the Americas (AMR, a 1KG superpopulation). Populations from the Americas often have a large and varying amount of recent admixture from various European and African populations [26, 48–52]. As such, including an AMR population in the ADMIXTURE analysis, as a reference population, could confound the genetic ancestries being estimated. However, while we limit this study to a few, broad, well-characterized ancestry groups, the approach presented here can be generalized to other, specific ancestries.

Third, the UKBB Axiom array used to genotype all UKBB participants was designed to optimize imputation of a European population while also including genetic variants previously associated with disease and other phenotypic traits derived from studies primarily conducted in European populations [11, 12]. As a product, the genomic data used here will have an ascertainment bias [53] that would influence imputation accuracy (although no imputation data were used here), allele frequency distributions, estimates of linkage disequilibrium, and diversity and divergence within and among populations. Each of these may influence estimations of population differentiation, principal component estimates, and the inferences made from them [54, 55]. Specific study designs [56, 57] have been made to remove ascertainment bias in genotype arrays so that unbiased inferences could be made for a wider range of genetic ancestries, but this was not available here.

Fourth, the principal components illustrated and used in the unsupervised K-means clustering analyses were derived from the UKBB participants only and resultantly represent the diversity (point three) and genetic ancestry found in that data set. The inclusion or use of other public data sets with more numerous sample populations,

that better represent regional, or continental diversity will provide alternative patterns of structure. Fifth, we are limited by the reference population used in the analyses. While the 1KG data set shall remain an essential reference panel for broad analyses like those conducted here, researchers with specific continental or geographically specific research questions could strengthen and refine the observations made here by including other geographically specific data sets. Finally, the unsupervised K-means clustering analysis is dependent upon the number of PCs included in it. Here, the number of PCs chosen did have an element of subjectivity (Additional file 2: Fig. S1). While analytical methods are available to select a number of informative PCs [58], we did not implement such methods here. Given that the K-means algorithm weights each PC equally, we sought to limit the PCs included to only those with the largest proportions of variance explained and not necessarily all that are analytically estimated to be informative.

## Conclusions

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

## Methods

### Description of working environment

All analyses were performed in a Linux environment supported by the University of Bristol's Advanced Computing Research Centre (ACRC) using the following publicly available software packages: PLINK v1.9 and v2.0 [59, 60], ADMIXTURE v1.3.0 [61, 62], and EIGENSOFT v8.0.0 [31, 32]. In addition, bespoke scripts, analyses, and figures were run and generated in the R environment using version 3.6.2 on the ACRC computer clusters and version 4.0.2 (Taking Off Again) on local computers [63].

### UK Biobank data

This research has been conducted using the UKBB Resource under Application Number 15825, from which directly genotyped SNP data ( $N=784,256$  SNPs) were made available. It includes data for a total of 78,296

individuals identified by UKBB as “non-white British” participants—our analyses were restricted to this subset. In addition to genotypic data, we also acquired several variables of interest (self-described ancestry, country of birth) data for this subset of individuals. 365 exclusions were made when filtering those with sex chromosome mismatch and/or aneuploidy, and outliers with high genetic heterozygosity and missing rates [64].

### 1000 Genomes data

Genetic data (v5a.20130502) from phase three of the 1KG, which includes data from 5 continental, or 1KG described super-populations [Europe (EUR), East Asia (EAS), South Asia (SAS), Africa (AFR), and the Americas (AMR)], were used to provide reference populations for admixture analyses and population structure inferences ([65] <http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/>). Our analyses did not include populations from the AMR super-population. This is to maintain a simplified analysis that avoided the complicating factors of the potentially recent admixture events that occurred in the Americas. Included in our analyses are five populations from 1KG super-population label: (AFR), also known as the continental Africa ancestry group (1) Yoruba in Ibadan, Nigeria (YRI); (2) Luhya in Webuye, Kenya (LWK); (3) Gambian in Western Division, The Gambia—Mandinka (GWD); (4) Mende in Sierra Leone (MSL); and (5) Esan in Nigeria (ESN). Five populations from the super-population label EUR or the continental Europe ancestry group: (1) Utah residents with Northern and Western European ancestry (CEU); (2) Toscani in Italia (TSI); (3) British in England and Scotland (GBR); (4) Finnish in Finland (FIN); and (5) Iberian populations in Spain (IBS). Five populations from the super-population label SAS or the continental South Asian ancestry group: (1) Gujarati Indian in Houston, Texas (GIH); (2) Punjabi in Lahore, Pakistan (PJI); (3) Bengali in Bangladesh (BEB); (4) Sri Lankan Tamil in the UK (STU); and (5) Indian Telugu in the UK (ITU). Finally, five populations from the super-population label EAS or the continental East Asian ancestry group: (1) Han Chinese in Beijing, China (CHB); (2) Japanese in Tokyo, Japan (JPT); (3) Han Chinese South (CHS); (4) Chinese Dai in Xishuangbanna, China (CDX); and (5) Kinh in Ho Chi Minh City, Vietnam (KHV).

### Merging UK Biobank and 1000 Genomes

The directly genotyped data from UKBB were used to identify SNPs with the same SNP identifier (RefSNP ID) present in the 1KG data set. A total of 718,711 SNPs were identified with the same ID and extracted from both data sets using PLINK v2.0. The two data sets were then merged using the `-bmerge` function in PLINK v2.0. After

removing problematic SNPs (e.g., multi-allelic, duplicate) in the merge step, a total of 718,487 SNPs remained.

### Linkage disequilibrium pruning

Prior to ancestry estimation, the merged data set was reduced to a set of independent SNPs based on linkage disequilibrium (LD) estimates using the PLINK v2.0 function and parameters “-indep-pairwise 50 10 0.025,” indicating an  $r^2$  threshold of 0.025, a window size of 50 kilobases and a window step size of 10 kilobases. In addition, 24 previously identified genomic regions with extensive linkage disequilibrium were also excluded [66, 67]. LD estimates in this analysis were limited to unrelated individuals from the 1KG YRI population sample. A total of 30,320 SNPs remained following LD pruning.

### Estimating African, European, South Asian, and East Asian ancestry

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

### Derivation of continental principal components

Unrelated individuals in each CAG including both 1KG and UKBB samples with  $\geq 80\%$  ancestry to that CAG were identified (using all 718,487 SNPs in the overlapping data set, and the PLINK (v1.9) function -rel-cutoff and a minor allele frequency (MAF) filter of 0.05 (-maf 0.05)). Then for each CAG and using all (1KG + UKBB) unrelated individuals assigned to the CAG, a list of approximately 40 thousand LD-independent SNPs were identified (using the PLINK (v2.0) function -indep-pairwise 50 10 0.025 (-indep-pairwise 50 10 0.02 for AFR and -indep-pairwise 50 10 0.05 for SAS) along with a MAF filter of 0.01, and the exclusion of the 24 previously identified genomic regions with extensive linkage

disequilibrium [66, 67]). New PLINK files including only the LD independent SNPs identified in step two were subsequently generated. smartrel from the EIGENSOFT (<https://github.com/DReichLab/EIG>) package was used to generate a new list of related individual pairs, along with our script “greedy\_unrelated\_selection.R” to identify a list of related individuals to exclude from principal component derivation [31, 32]. An exception this step was made for the European CAG as its sample size was prohibitively large to run smartrel; instead the list of unrelated individuals generated from step one was used. Finally, smartpca of the EIGENSOFT package was used to estimate principal components (PC), using only unrelated UKBB samples. Related and 1KG samples were subsequently projected upon these PCs by smartpca. Sample outliers were excluded from the PC analysis by smartpca with the following parameters: using 10 PCs to identify outliers (numoutlierevec), at six standard deviations from the mean (outliersigmathresh), and with 5 outlier removal iterations (numoutlieriter). Additional file 1: Table S1 provides numbers for each of these steps, for each CAG. The EUR CAG was treated uniquely due to its larger sample size. Smartpca was run twice as described above, once with “fastmode=NO” and then with “fastmode=YES.” The former provided estimates of the eigenvalues but not the eigenvectors, while the latter provided eigenvectors but not eigenvalues.

### K-means clustering of principal components

For each CAG, we estimated the variance explained by each principal component (PC) by dividing the eigenvalue of each PC by the sum of all eigenvalues. To identify the number of top PCs, we generated a scree plot, using the variance explained estimates, and identified the elbow or valley in each plot (Additional file 2: Fig. S1, Additional file 1: Table S2). The top PCs, and the top PCs only, were then used in an unsupervised K-means clustering analysis (k set from 2 to 20; using the function “kmeans()” from the R stats package) to identify clusters of UKBB individuals that maximize between cluster sums of squares and minimize within cluster sums of squares. An optimum number of clusters (k) were identified by silhouette analysis using the function “pamk()” from the fpc R package (Additional file 2: Fig. S2) [68]. These analyses are implemented in our function “DetermineK()” found in this study’s GitHub repository.

### Correspondence analysis

Each UKBB study participants’ country of birth information was placed into United Nations defined geographic regions (Additional file 1: Table S3). To determine whether the K-means population clusters have any relationship with an individual’s country of

birth or country of birth UN-region, we performed correspondence analyses (CAs) using the function “ca()” from the R package “ca,” for each continental ancestry group [38]. In addition, a Chi-square test was performed on the contingency table used in the correspondence analysis. Any UN-region or country of birth with fewer than 10 observations was excluded. Individuals for which country of birth information was not available were also excluded.

### Population differentiation among K-means population clusters

For each CAG, we took the best K-means population clusters, as defined by the silhouette analysis, and reran smartpca. However, on this run smartpca provides for us only an estimation of the average fixation index (Fst) for each pair of populations in the data set, including 1KG populations and UKBB K-means clusters. This was done with the inclusion of the parameters “fstonly” and “phylipoutname” [58], the latter of which provides a distance matrix of mean Fst values between populations. Estimations of Fst, which range from 0 to 1, provide a measure of population differentiation among populations. In brief, these describe the proportion of total variation at a SNP that is explained by variation between populations. For any SNP, a value of 0 would indicate that minimal variation is attributable to variation between populations. A value of 1 would indicate a fixed difference, i.e., the two populations are both invariable but for alternative alleles.

### Abbreviations

1KG: 1000 Genomes Project; ACRC: Advanced Computing Research Centre; AFR: African; AMR: Americas; BEB: Bengali in Bangladesh; CA: Correspondence analysis; CAG: Continental ancestry group; CDX: Chinese Dai in Xishuangbanna, China; CEU: Utah residents with Northern and Western European ancestry; CHB: Han Chinese in Beijing, China; CHS: Han Chinese South; COB: Country of birth; EAS: East Asian; ESN: Esan in Nigeria; EUR: European; FIN: Finnish in Finland; Fst: Fixation index; GBR: British in England and Scotland; GIH: Gujarati Indian in Houston, Texas; GWAS: Genome-wide association study; GWD: Gambian in Western Division, The Gambia—Mandinka; IBS: Iberian populations in Spain; ITU: Indian Telugu in the UK; JPT: Japanese in Tokyo, Japan; KHV: Kinh in Ho Chi Minh City, Vietnam; LD: Linkage disequilibrium; LWK: Luhya in Webuye, Kenya; MAF: Minor allele frequency; MSL: Mende in Sierra Leone; PC: Principal component; PJI: Punjabi in Lahore, Pakistan; ROB: Region of birth; SAS: South Asian; SNP: Single-nucleotide polymorphism; STU: Sri Lankan Tamil in the UK; TSI: Toscani in Italia; UKBB: UK Biobank; UN: United Nations; YRI: Yoruba in Ibadan, Nigeria.

### Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s40246-022-00380-5>.

**Additional file 1. Table S1.** Study genetic processing steps with relevant numbers. **Table S2.** Eigenvalue and eigenvector for each continental

ancestry group. **Table S3.** United Nations geoscheme for each country and continent present in the study.

**Additional file 2. Figure S1.** Continental ancestry PCA scree plots. **Figure S2.** K-means k selection with silhouette analysis. **Figure S3.** UK Biobank continental ancestry group PCs with K-means clusters. **Figure S4.** Population structure by UN defined geographic region. **Figure S5.** Population structure centers, as defined by UN geographic region. **Figure S6.** Population structure by country of birth in Africa by region. **Figure S7.** Population structure by country of birth in Europe by region. **Figure S8.** Population structure by country of birth in South Asia and East Asia. **Figure S9.** Population structure centers by country of birth. **Figure S10.** Population structure centers by country of birth.

### Acknowledgements

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

### Authors' contributions

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

### Funding

AC acknowledges funding from a Medical Research Council PhD studentship (MR/N013794/1). NJT and REM acknowledge funding from the Medical Research Council (MC\_UU\_00011/1). NJT is the PI of the Avon Longitudinal Study of Parents and Children (Medical Research Council & Wellcome Trust 217065/Z/19/Z) and is supported by the University of Bristol NIHR Biomedical Research Centre (BRC-1215-2001). EEV, CJB, NJT, and DH acknowledge funding from the Wellcome Trust (202802/Z/16/Z). EEV, CJB, and NJT also acknowledge funding by the CRUK Integrative Cancer Epidemiology Programme (C18281/A29019). EEV and CJB are supported by Diabetes UK (17/0005587) and the World Cancer Research Fund (WCRF UK), as part of the World Cancer Research Fund International grant program (IIG\_2019\_2009). JZ is supported by the Academy of Medical Sciences (AMS) Springboard Award, the Wellcome Trust, the Government Department of Business, Energy and Industrial Strategy (BEIS), the British Heart Foundation and Diabetes UK (SBF006/1117). JZ is funded by the Vice-Chancellor Fellowship from the University of Bristol and is supported by Shanghai Thousand Talents Program. BA acknowledges funding from the Medical Research Council (MR/R02149X/1). The funders of the study had no role in the study design, data collection, data analysis, data interpretation, or writing of the report.

### Availability of data and materials

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

### Declarations

#### Ethics approval and consent to participate

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

#### Consent for publication

All authors consented to the publication of this work.

#### Competing interests

The authors declare no competing interests.

#### Author details

<sup>1</sup>MRC Integrative Epidemiology Unit at the University of Bristol, Bristol, UK. <sup>2</sup>Bristol Medical School, Population Health Sciences, University of Bristol, Bristol, UK. <sup>3</sup>School of Translational Health Sciences, University of Bristol, Bristol, UK. <sup>4</sup>School of Cellular and Molecular Medicine, University of Bristol, Bristol, UK.

Received: 17 December 2021 Accepted: 18 January 2022  
Published online: 29 January 2022

## References

- Timpson NJ, Greenwood CMT, Soranzo N, Lawson DJ, Richards JB. Genetic architecture: the shape of the genetic contribution to human traits and disease. *Nat Rev Genet.* 2017;19(2):110–24. <https://doi.org/10.1038/nrg.2017.101>.
- Sirugo G, Williams SM, Tishkoff SA. The missing diversity in human genetic studies. *Cell.* 2019;177:26–31. <https://doi.org/10.1016/j.cell.2019.02.048>.
- Bentley AR, Callier SL, Rotimi CN. Evaluating the promise of inclusion of African ancestry populations in genomics. *Npj Genomic Med.* 2020;5(1):9. <https://doi.org/10.1038/s41525-019-0111-x>.
- Cooke Bailey JN, Bush WS, Crawford DC. Editorial: the importance of diversity in precision medicine research. *Front Genet.* 2020. <https://doi.org/10.3389/fgene.2020.00875>.
- Green ED, Gunter C, Biesecker LG, Di Francesco V, Easter CL, Feingold EA, et al. Strategic vision for improving human health at The Forefront of Genomics. *Nature.* 2020;586(7831):683–92. <https://doi.org/10.1038/s41586-020-2817-4>.
- Consortium TH. Enabling the genomic revolution in Africa: H3Africa is developing capacity for health-related genomics research in Africa. *Science.* 2014;344:1346. <https://doi.org/10.1126/SCIENCE.1251546>.
- Matisse TC, Study for the P, Ambite JL, Study for the P, Buyske S, Study for the P, et al. The Next PAGE in Understanding Complex Traits: Design for the Analysis of Population Architecture Using Genetics and Epidemiology (PAGE) Study. *Am J Epidemiol* 2011;174:849–59. <https://doi.org/10.1093/AJE/KWR160>.
- Taliun D, Harris DN, Kessler MD, Carlson J, Szpiech ZA, Torres R, et al. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nat.* 2021;590(7845):290–9. <https://doi.org/10.1038/s41586-021-03205-y>.
- Gallo LC, Penedo FJ, Carnethon M, Isasi C, Sotres-Alvarez D, Malcarne VL, et al. The Hispanic Community Health Study/Study of Latinos Sociocultural Ancillary Study: Sample, Design, and Procedures. *Ethn Dis.* 2014;24:77.
- Investigators TA of URP. The “All of Us” Research Program. 2019;381:668–76. <https://doi.org/10.1056/NEJMSR1809937>.
- Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, et al. UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS Med.* 2015;12:e1001779–e1001779. <https://doi.org/10.1371/journal.pmed.1001779>.
- Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, et al. The UK Biobank resource with deep phenotyping and genomic data. *Nature.* 2018;562:203–9. <https://doi.org/10.1038/s41586-018-0579-z>.
- Mathieson I, Scally A. What is ancestry? *PLOS Genet.* 2020;16:e1008624. <https://doi.org/10.1371/JOURNAL.PGEN.1008624>.
- Rodriguez S, Gaunt TR, Day INM. Hardy–Weinberg Equilibrium Testing of Biological Ascertainment for Mendelian Randomization Studies. *Am J Epidemiol.* 2009;169:505–14. <https://doi.org/10.1093/AJE/KWN359>.
- Graffelman J, Weir BS. On the testing of Hardy–Weinberg proportions and equality of allele frequencies in males and females at biallelic genetic markers. *Genet Epidemiol.* 2018;42:34–48. <https://doi.org/10.1002/GEPI.22079>.
- Altshuler DL, Durbin RM, Abecasis GR, Bentley DR, Chakravarti A, Clark AG, et al. A map of human genome variation from population-scale sequencing. *Nature.* 2010;467:1061–73. <https://doi.org/10.1038/nature09534>.
- Rosenberg NA, Pritchard JK, Weber JL, Cann HM, Kidd KK, Zhivotovskiy LA, et al. Genetic structure of human populations. *Science* (80-). 2002;298:2381–5. [https://doi.org/10.1126/SCIENCE.1078311/SUPPL\\_FILE/ROSENBERG.SOM.PDF.PDF](https://doi.org/10.1126/SCIENCE.1078311/SUPPL_FILE/ROSENBERG.SOM.PDF.PDF).
- Berezovskii ND, Giria VN. Estimation of combining ability of specialized types of the big white breed. *Tsitol Genet.* 1991;25:56–60.
- Serre D, Pääbo S. Evidence for gradients of human genetic diversity within and among continents. *Genome Res.* 2004;14:1679. <https://doi.org/10.1101/GR.2529604>.
- Rosenberg NA, Mahajan S, Ramachandran S, Zhao C, Pritchard JK, Feldman MW. Clines, clusters, and the effect of study design on the inference of human population structure. *PLoS Genet.* 2005;1: e70. <https://doi.org/10.1371/JOURNAL.PGEN.0010070>.
- Birney E, Inouye M, Raff J, Rutherford A, Scally A. The language of race, ethnicity, and ancestry in human genetic research n.d.
- Peterson RE, Kuchenbaecker K, Walters RK, Chen CY, Popejoy AB, Periyasamy S, et al. Genome-wide association studies in ancestrally diverse populations: opportunities, methods, pitfalls, and recommendations. *Cell.* 2019;179:589–603. <https://doi.org/10.1016/j.cell.2019.08.051>.
- Laland KN, Odling-Smee J, Myles S. How culture shaped the human genome: bringing genetics and the human sciences together. *Nat Rev Genet.* 2010;11(2):137–48. <https://doi.org/10.1038/nrg2734>.
- Przeworski M, Wall JD. Why is there so little intragenic linkage disequilibrium in humans? *Genet Res.* 2001;77:143–51. <https://doi.org/10.1017/S0016672301004967>.
- Ptak SE, Voelpel K, Przeworski M. Insights into recombination from patterns of linkage disequilibrium in humans. *Genetics.* 2004;167:387. <https://doi.org/10.1534/GENETICS.167.1.387>.
- Auton A, Abecasis GR, Altshuler DM, Durbin RM, Bentley DR, Chakravarti A, et al. A global reference for human genetic variation. *Nature.* 2015;526:68–74. <https://doi.org/10.1038/nature15393>.
- Lander ES, Schork NJ. Genetic dissection of complex traits. *Science* (80-). 1994;265:2037–48. <https://doi.org/10.1126/SCIENCE.8091226>.
- Hellwege JN, Keaton JM, Giri A, Gao X, Velez Edwards DR, Edwards TL. Population Stratification in Genetic Association Studies. *Curr Protoc Hum Genet.* 2017;95(1):22. <https://doi.org/10.1002/CPHG.48>.
- Vilhjálmsón BJ, Nordborg M. The nature of confounding in genome-wide association studies. *Nat Rev Genet.* 2012;14(1):1–2. <https://doi.org/10.1038/nrg3382>.
- Loh P-R, Kichaev G, Gazal S, Schoech AP, Price AL. Mixed-model association for biobank-scale datasets. *Nat Genet.* 2018;50(7):906–8. <https://doi.org/10.1038/s41588-018-0144-6>.
- Patterson N, Price AL, Reich D. Population structure and eigenanalysis. *PLoS Genet.* 2006;2:2074–93. <https://doi.org/10.1371/journal.pgen.0020190>.
- Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet.* 2006;38:904–9. <https://doi.org/10.1038/ng1847>.
- Loh PR, Tucker G, Bulik-Sullivan BK, Vilhjálmsson BJ, Finucane HK, Salem RM, et al. Efficient Bayesian mixed-model analysis increases association power in large cohorts. *Nat Genet.* 2015;47:284–90. <https://doi.org/10.1038/ng.3190>.
- Price AL, Zaitlen NA, Reich D, Patterson N. New approaches to population stratification in genome-wide association studies. *Nat Rev Genet.* 2010;11:459. <https://doi.org/10.1038/NRG2813>.
- Zaidi AA, Mathieson I. Demographic history mediates the effect of stratification on polygenic scores. *Elife.* 2020;9:1–30. <https://doi.org/10.7554/ELIFE.61548>.
- Haworth S, Mitchell R, Corbin L, Wade KH, Dudding T, Budu-Aggrey A, et al. Apparent latent structure within the UK Biobank sample has implications for epidemiological analysis. *Nat Commun.* 2019. <https://doi.org/10.1038/s41467-018-08219-1>.
- Abdellaoui A, Hugh-Jones D, Yengo L, Kemper KE, Nivard MG, Veul L, et al. Genetic correlates of social stratification in Great Britain. *Nat Hum Behav.* 2019;3(12):1332–42. <https://doi.org/10.1038/s41562-019-0757-5>.
- Lao O, Lu TT, Nothnagel M, Junge O, Freitag-Wolf S, Caliebe A, et al. Correlation between Genetic and Geographic Structure in Europe. *Curr Biol.* 2008;18:1241–8. <https://doi.org/10.1016/j.cub.2008.07.049>.
- Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A, et al. Genes mirror geography within Europe. *Nature.* 2008;456:98. <https://doi.org/10.1038/NATURE07331>.
- Wood AR, Esko T, Yang J, Vedantam S, Pers TH, Gustafsson S, et al. Defining the role of common variation in the genomic and biological architecture of adult human height. *Nat Genet.* 2014;46:1173. <https://doi.org/10.1038/NG.3097>.
- Berg JJ, Harpak A, Sinnott-Armstrong N, Joergensen AM, Mostafavi H, Field Y, et al. Reduced signal for polygenic adaptation of height in UK biobank. *Elife.* 2019. <https://doi.org/10.7554/eLife.39725>.

42. Sohail M, Maier RM, Ganna A, Bloemendal A, Martin AR, Turchin MC, et al. Polygenic adaptation on height is overestimated due to uncorrected stratification in genome-wide association studies. *Elife*. 2019. <https://doi.org/10.7554/eLife.39702>.
43. Lawson DJ, Davies NM, Haworth S, Ashraf B, Howe L, Crawford A, et al. Is population structure in the genetic biobank era irrelevant, a challenge, or an opportunity? *Hum Genet*. 2020;139:23–41. <https://doi.org/10.1007/s00439-019-02014-8>.
44. Barton N, Hermisson J, Nordborg M. Why structure matters. *Elife*. 2019. <https://doi.org/10.7554/ELIFE.45380>.
45. Diaz-Papkovich A, Anderson-Trocmé L, Ben-Eghan C, Gravel S. UMAP reveals cryptic population structure and phenotype heterogeneity in large genomic cohorts. *PLoS Genet*. 2019. <https://doi.org/10.1371/journal.pgen.1008432>.
46. Morton NE. Isolation by distance. *Genetics*. 1943;28:114. <https://doi.org/10.1016/B978-0-12-374984-0.00820-2>.
47. Slatkin M. Isolation by distance in equilibrium and non-equilibrium populations. *Evolution*. 1993;47:264–79. <https://doi.org/10.1111/J.1558-5646.1993.TB01215.X>.
48. Kidd JM, Gravel S, Byrnes J, Moreno-Estrada A, Musharoff S, Bryc K, et al. Population genetic inference from personal genome data: impact of ancestry and admixture on human genomic variation. *Am J Hum Genet*. 2012;91:660. <https://doi.org/10.1016/J.AJHG.2012.08.025>.
49. Homburger JR, Moreno-Estrada A, Gignoux CR, Nelson D, Sanchez E, Ortiz-Tello P, et al. Genomic insights into the Ancestry and Demographic History of South America. *PLoS Genet*. 2015. <https://doi.org/10.1371/JOURNAL.PGEN.1005602>.
50. Moreno-Estrada A, Gravel S, Zakharia F, McCauley JL, Byrnes JK, Gignoux CR, et al. Reconstructing the Population Genetic History of the Caribbean. *PLoS Genet*. 2013. <https://doi.org/10.1371/JOURNAL.PGEN.1003925>.
51. Ongaro L, Scliar MO, Flores R, Raveane A, Marnetto D, Sarno S, et al. The Genomic Impact of European Colonization of the Americas. *Curr Biol*. 2019;29:3974–3986.e4. <https://doi.org/10.1016/J.CUB.2019.09.076/ATTACHMENT/8D05D549-D774-4CBA-9BE7-94D3B60AD79D/MMC3.XLSX>.
52. Montinaro F, Busby GBJ, Pascali VL, Myers S, Hellenthal G, Capelli C. Unravelling the hidden ancestry of American admixed populations. *Nat Commun*. 2015. <https://doi.org/10.1038/NCOMMS7596>.
53. Geibel J, Reimer C, Weigend S, Weigend A, Pook T, Simianer H. How array design creates SNP ascertainment bias. *PLoS ONE*. 2021;16:e0245178–e0245178. <https://doi.org/10.1371/journal.pone.0245178>.
54. Lachance J, Tishkoff SA. SNP ascertainment bias in population genetic analyses: Why it is important, and how to correct it. *BioEssays*. 2013;35:780–6. <https://doi.org/10.1002/bies.201300014>.
55. Albrechtsen A, Nielsen FC, Nielsen R. Ascertainment biases in SNP chips affect measures of population divergence. *Mol Biol Evol*. 2010;27:2534–47. <https://doi.org/10.1093/molbev/msq148>.
56. Patterson N, Moorjani P, Luo Y, Mallick S, Rohland N, Zhan Y, et al. Ancient Admixture in Human History. *Genetics*. 2012;192:1065. <https://doi.org/10.1534/GENETICS.112.145037>.
57. Lu Y, Patterson N, Zhan Y, Mallick S, Reich D. Technical design document for a SNP array that is optimized for population genetics n.d.
58. Reich D, Thangaraj K, Patterson N, Price AL, Singh L. Reconstructing Indian population history. *Nature*. 2009;461:489–94. <https://doi.org/10.1038/nature08365>.
59. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: A tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet*. 2007;81:559–75. <https://doi.org/10.1086/519795>.
60. Chang CC, Chow CC, Tellier LCAM, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: Rising to the challenge of larger and richer datasets. *Gigascience*. 2015. <https://doi.org/10.1186/s13742-015-0047-8>.
61. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals. *Genome Res*. 2009;19:1655–64. <https://doi.org/10.1101/gr.094052.109>.
62. Alexander DH, Lange K. Enhancements to the ADMIXTURE algorithm for individual ancestry estimation. *BMC Bioinformatics*. 2011;12:246. <https://doi.org/10.1186/1471-2105-12-246>.
63. Core R Team. R: A Language and Environment for Statistical Computing. *R Found Stat Comput* 2019;2:<https://www.R-project.org>. <http://www.r-project.org> (accessed March 2, 2021).
64. Mitchell RE, Hemani G, Dudding T, Corbin L, Harrison S, Paternoster L. UK Biobank Genetic Data: MRC-IEU Quality Control, version 2, 18/01/2019 n.d.
65. Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et al. An integrated map of structural variation in 2,504 human genomes. *Nature*. 2015;526:75–81. <https://doi.org/10.1038/nature15394>.
66. Price AL, Weale ME, Patterson N, Myers SR, Need AC, Shianna KV, et al. Long-Range LD Can Confound Genome Scans in Admixed Populations. *Am J Hum Genet*. 2008;83:132. <https://doi.org/10.1016/J.AJHG.2008.06.005>.
67. Weale ME. Quality Control for Genome-Wide Association Studies. In: Barnes MR, Breen G, editors. *Genet. Var. Methods Protoc.*, Humana Press, New York, NY; 2010, p. 31.
68. Batool F, Hennig C. Clustering with the Average Silhouette Width. *Comput Stat Data Anal*. 2021;158: 107190. <https://doi.org/10.1016/J.JCSDA.2021.107190>.

## Publisher's Note

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

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more [biomedcentral.com/submissions](https://biomedcentral.com/submissions)

