

# Portability of 245 polygenic scores when derived from the UK Biobank and applied to 9 ancestry groups from the same cohort

Florian Privé,<sup>1,\*</sup> Hugues Aschard,<sup>2,3</sup> Shai Carmi,<sup>4</sup> Lasse Folkersen,<sup>5</sup> Clive Hoggart,<sup>6</sup> Paul F. O'Reilly,<sup>6</sup> and Bjarni J. Vilhjálmsson<sup>1,7</sup>

## Summary

The low portability of polygenic scores (PGSs) across global populations is a major concern that must be addressed before PGSs can be used for everyone in the clinic. Indeed, prediction accuracy has been shown to decay as a function of the genetic distance between the training and test cohorts. However, such cohorts differ not only in their genetic distance but also in their geographical distance and their data collection and assaying, conflating multiple factors. In this study, we examine the extent to which PGSs are transferable between ancestries by deriving polygenic scores for 245 curated traits from the UK Biobank data and applying them in nine ancestry groups from the same cohort. By restricting both training and testing to the UK Biobank data, we reduce the risk of environmental and genotyping confounding from using different cohorts. We define the nine ancestry groups at a sub-continental level, based on a simple, robust, and effective method that we introduce here. We then apply two different predictive methods to derive polygenic scores for all 245 phenotypes and show a systematic and dramatic reduction in portability of PGSs trained using Northwestern European individuals and applied to nine ancestry groups. These analyses demonstrate that prediction already drops off within European ancestries and reduces globally in proportion to genetic distance. Altogether, our study provides unique and robust insights into the PGS portability problem.

## Introduction

Ever larger genetic datasets are becoming more readily available. This enables researchers to derive polygenic scores (PGSs), which summarize an individual's genetic component for a particular trait or disease by aggregating information from many genetic variants into a single score. In human genetics, polygenic scores are usually derived from summary statistics from a large meta-analysis of multiple genome-wide association studies (GWASs) and an ancestry-matched linkage disequilibrium (LD) reference panel.<sup>1</sup> Polygenic scores can also be derived directly from individual-level data when available, i.e., from the genetic and phenotypic information of many individuals.<sup>2</sup> When using a single individual-level dataset with only moderate sample size, deriving polygenic scores usually results in poor prediction for most phenotypes, e.g., for autoimmune diseases with moderately large effects.<sup>3,4</sup> Fortunately, biobank datasets such as the UK Biobank now link genetic data for half a million individuals with phenotypic data for hundreds of traits and diseases.<sup>5</sup> Thanks to the availability of these large datasets and to efficient methods recently developed to handle such data,<sup>4,6,7</sup> individual-level data may be used to derive competitive PGSs for hundreds of phenotypes.

A major concern about PGSs is that they usually transfer poorly to other ancestries, e.g., a PGS derived from individ-

uals of European ancestry is not likely to predict as well in individuals of African ancestry. Prediction in another ancestry has been shown to decay with genetic distance to the training population<sup>8,9</sup> and with increasing proportion of admixture with a distant ancestry.<sup>10,11</sup> This portability issue is suspected to be primarily due to differences in LD and allele frequencies between populations, and not so much about differences in effects and positions of causal variants.<sup>9,11</sup> Individual-level data from the UK Biobank offers an opportunity to further investigate this problem of PGS portability in a more controlled setting.<sup>9,12</sup> Indeed, while the UK Biobank data contain genetic information for more than 450K British or European individuals, it also contains the same data for tens of thousands of individuals of non-British ancestry.<sup>5</sup> Of particular interest, those individuals of diverse ancestries all live in the UK and had their genetic and phenotypic information derived in the same way as people of UK ancestry. Our study design circumvents potential confounding bias that might arise in comparative analyses from independent studies and makes the UK Biobank data very well suited for comparing and evaluating predictive performance of derived PGSs in diverse ancestries and across multiple phenotypes. Indeed, the UK Biobank has been shown to offer a much more controlled setting (compared to published GWAS meta-analyses) in the case of studying (for example) polygenic adaptation.<sup>13,14</sup> Note that these analyses are not completely free

<sup>1</sup>National Centre for Register-Based Research, Aarhus University, Aarhus 8210, Denmark; <sup>2</sup>Department of Computational Biology, Institut Pasteur, Paris 75015, France; <sup>3</sup>Program in Genetic Epidemiology and Statistical Genetics, Harvard T.H. Chan School of Public Health, Boston, MA 02115, USA; <sup>4</sup>Braun School of Public Health and Community Medicine, The Hebrew University of Jerusalem, Jerusalem 9112102, Israel; <sup>5</sup>Danish National Genome Center, Copenhagen 2300, Denmark; <sup>6</sup>Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY 10029, USA; <sup>7</sup>Bioinformatics Research Centre, Aarhus University, Aarhus 8000, Denmark

\*Correspondence: [florian.prive.21@gmail.com](mailto:florian.prive.21@gmail.com)

<https://doi.org/10.1016/j.ajhg.2021.11.008>

© 2021 American Society of Human Genetics.



of bias since, on average, genotyped variants are more common and imputed variants are more accurately imputed in European ancestries. We also acknowledge that some residual structure may remain when deriving PGSs.<sup>15</sup>

To investigate portability of PGSs to other ancestries, we must first define groups of different ancestries from the data. Principal component analysis (PCA) has been widely used to correct for population structure in association studies and has been shown to mirror geography in Europe.<sup>16,17</sup> Due to its popularity, many methods have been developed for efficiently performing PCA<sup>18–20</sup> as well as appropriately projecting samples onto a reference PCA space,<sup>20,21</sup> making it possible to perform these analyses for ever increasing datasets. Naturally, PCA has also been used for ancestry inference.<sup>21–23</sup> However, among the studies where we have seen PCA used for ancestry inference, there does not seem to be a consensus on what is the most appropriate method for inferring ancestry using PCA. For example, there are divergences on which distance metric to use and the number of PCs to use to compute these distances. The ancestry of an individual can also be inferred based on other approaches, including the ADMIXTURE model, its various extensions, and haplotype-based methods.<sup>24–31</sup> However, we focus on PCA here because it is very fast and effective.

In this study, we examine the extent to which PGSs are transferable between ancestries by deriving 245 polygenic scores from the UK Biobank data and applying them in nine ancestry groups from the same cohort. We first propose simple, robust, and effective methods for global ancestry inference and grouping from PCA of genetic data, and we use them to define nine ancestry groups in the UK Biobank data. We then apply a computationally efficient implementation of penalized regression<sup>4</sup> to derive PGSs for 245 traits using the UK Biobank genetic and phenotypic data only. As an alternative method, we also run LDpred2-auto,<sup>32</sup> for which we directly derive the summary statistics from the individual-level data available. We show a dramatically low portability of PGSs from UK ancestry to other ancestries. For example, on average, the phenotypic variance explained by the PGSs is only 64.7% in South Asia (the “India” ancestry group defined here), 48.6% in East Asia (“China”), and 18% in West Africa (“Nigeria”) compared to in individuals of Northwestern European ancestry (“United Kingdom”). These results are presented at a finer scale than the usual continental level, which allows us to show that prediction already drops within Europe, e.g., for Northeast and South Europe (the “Poland” and “Italy” ancestry groups) compared to Northwest Europe. We find that this decay in variance explained by the PGSs is roughly linear in the PC distance to the training population and is remarkably consistent across most phenotypes and for both prediction methods applied. The few exceptions include traits such as hair color, tanning, and some blood measurements. We also explore using more than HapMap3 variants when fitting PGSs, it proves useful when large effects are poorly tagged by HapMap3 var-

iants, e.g., for lipoprotein(a), but not in the general case. We also explore the performance of PGS trained using a mixture of European and non-European ancestry samples, but do not observe any significant gain in prediction here.

## Material and methods

### Data

We derive polygenic scores for 245 phenotypes using the UK Biobank (UKBB) data only.<sup>5</sup> We read dosages data from UKBB BGEN files using function `snp_readBGEN()` of R package `bigsnpr`.<sup>19</sup> We divide the UKBB data in eight ancestry groups (Note A) and restrict to 437,669 individuals without second-degree relatives (KING kinship  $< 2^{-3.5}$ ). We also define a ninth ancestry group composed of 1,709 unrelated Ashkenazi (see below). For the variants, we use 1,040,096 HapMap3 variants used in the LD reference provided in Privé et al.<sup>32</sup> and that were also present in the iPSYCH2015 data<sup>33</sup> with imputation INFO score larger than 0.6. Even though the iPSYCH data is not used in this study, we plan to use the PGSs derived here for iPSYCH in the future.

To define phenotypes, we first map ICD10 and ICD9 codes (UKBB fields 40001, 40002, 40006, 40013, 41202, 41270, and 41271) to phecodes using R package `PheWAS`.<sup>34,35</sup> We filter down to 142 phecodes of interest that showed potential genetic signals in the PheWeb results from the SAIGE UKBB GWAS.<sup>36,37</sup> We further filter down to 106 phecodes with sufficient power for penalized regression to include at least a few variants in the predictive models. We then look closely at all 2,408 UKBB fields that we have access to and filter down to defining 111 continuous and 28 binary phenotypes based on manual curation.

### Additional data: Genotyped data

For the genotyped data used in some follow-up analyses, we restrict to variants that have been genotyped on both chips used by the UK Biobank, that pass quality control (QC) for all batches and QC for possible mismappings,<sup>38</sup> with a minor allele frequency (MAF) larger than 0.01 and imputation INFO score of 1. There are 586,534 such high-quality variants, which we read from the BGEN imputed data so that there is no missing value.

### Additional data: 8M+ variants

We also design a larger set of imputed variants to compare against using only HapMap3 variants for prediction. We first restrict to UKBB variants with MAF  $> 0.01$  and INFO  $> 0.6$ . We then compile frequencies and imputation INFO scores from other datasets, iPSYCH, and summary statistics for breast cancer, prostate cancer, coronary artery disease, and type 1 diabetes.<sup>33,39–42</sup> We restrict to variants with a mean INFO  $> 0.5$  in these other datasets and also compute the median frequency. To exclude potential mismappings in the genotyped data<sup>38</sup> that might have propagated to the imputed data, we compare median frequencies in the external data to the ones in UKBB (Figure S20). As we expect these potential errors to be localized around errors in the genotype data (confirmed in Figure S21), we apply a moving-average smoothing on the frequency differences to increase power to detect these errors and also reduce false positives. We define the threshold on these smoothed differences based on visual inspection of their histogram. This is the same method we have previously applied to PC loadings to detect long-range LD regions when computing PCA.<sup>19,20</sup> This results in a set of 8,238,692 variants.

## Ashkenazi Jewish ancestry group

First, we refer the reader to [Note A](#) on ancestry grouping for the details on how we define the other eight ancestry groups, and also to better understand how we infer the “Ashkenazi Jewish” ancestry group. Briefly, we project the UKBB data onto the PCA space of a reference dataset composed of many Jewish and non-Jewish individuals.<sup>43</sup> We then compute the robust center (geometric median) of the Ashkenazi Jewish reference individuals and compute the PC distance to this center for all projected UKBB individuals. Based on visual inspection of the histogram of these distances and on the fact that the closest non-Ashkenazi Jewish reference individual, an Italian Jew ([Figure S22](#)), is at distance 12.7, we use a threshold of 12.5 under which to assign to the “Ashkenazi Jewish” ancestry group. 1,709 unrelated UKBB individuals are then assigned to this group. Note that, within the already defined eight ancestry groups, the closest individual to this new group belongs to the Italian group, and is at distance 17.3, so this new Ashkenazi group is not overlapping with any of the other groups defined previously.

## Penalized regression

To derive polygenic scores based on individual-level data from the UKBB, we use the fast implementation of penalized linear and logistic regressions from R package `bigstatsr`.<sup>4</sup> We have also considered the recently developed R package `snpnet` for fitting penalized regressions on large genetic data; however, we provide theoretical and empirical evidence that `bigstatsr` is much faster than `snpnet` ([Note B](#)). Our implementation allows for lasso and elastic-net penalizations; yet, for the sake of simplicity and because the UKBB data is very large, we have decided to only use the lasso penalty.<sup>4</sup> We recall that fitting a penalized linear regression with lasso penalty corresponds to finding the vector of effects  $\beta$  (also  $\mu$  and  $\gamma$ ) that minimizes

$$L(\lambda) = \underbrace{\|y - (\mu + G\beta + X\gamma)\|_2^2}_{\text{Loss function}} + \underbrace{\lambda \|\beta\|_1}_{\text{Penalisation}},$$

where  $\mu$  is an intercept,  $G$  is the genotype matrix,  $X$  is the matrix of covariates,  $y$  is the (quantitative) phenotype of interest, and  $\lambda$  is a hyper-parameter that controls the strength of the regularization and needs to be chosen. We use sex (field 22001), age (field 21022), birth date (fields 34 and 52), Townsend deprivation index (field 189), and the first 16 genetic principal components (field 22009),<sup>20</sup> as unpenalized covariates when fitting the lasso models.

We have extended our implementation in two ways by allowing

for using different penalties for the variants (i.e., having  $\sum_j \lambda_j |\beta_j|$  instead of  $\lambda \|\beta\|_1$ ). First, this enables us to use a different scaling for genotypes. By default, variants in  $G$  are implicitly scaled. By using  $\lambda_j \propto (\text{SD}_j)^{(\xi-1)}$ , this effectively scales variant  $j$  by dividing it by  $(\text{SD}_j)^\xi$  in our implementation. The default uses  $\xi = 1$  but we also test  $\xi = 0$  (no scaling) and  $\xi = 0.5$  (Pareto scaling). We introduce a new parameter `power_scale` for which the user can provide a vector of values to test; the best value is chosen within the Cross-Model Selection and Averaging (CMSA) procedure.<sup>4</sup> We also introduce a second parameter, `power_adaptive`, which can be used to put less penalization on variants with the largest marginal effects,<sup>44</sup> we try three values here (0 the default, 0.5, and 1.5) and the best one is also chosen within the CMSA procedure.

## LDpred2-auto

Using the individual-level data from the training set in the UK Biobank, we run a linear regression GWAS using function `big_univLin`

Reg of R package `bigstatsr`,<sup>19</sup> accounting for the same covariates as in the penalized regression above. As LD reference, we use the one provided in [Privé et al.](#)<sup>32</sup> based on UKBB data for European ancestry. We use these summary statistics and this LD reference as input for LDpred2-auto. LDpred2 assumes a point-normal mixture distribution for effect sizes, where only a proportion of causal variants  $p$  contributes to the SNP heritability  $h^2$ . In LDpred2-auto, these two parameters are directly estimated from the data.<sup>32</sup> We use the sparse option in LDpred2-auto to also obtain a vector of effects that is potentially sparse, i.e., effects of some variants are exactly 0. Also note that, as we use linear regression for all phenotypes, we use the total sample size instead of the effective sample size ( $4/(1/n_{\text{case}} + 1/n_{\text{control}})$ ) for binary phenotypes as input to LDpred2. This means that heritability estimates from both LD score regression and LDpred2-auto must be transformed to the liability scale using both the prevalence in the GWAS and in the population; this can be performed using function `coef_to_liab` from R package `bigsnpr`. For simplicity, we assume here that the prevalence in the population is the same as the prevalence in the training set.

## New formula used in LDpred2

We also slightly modify the formula used in [Privé et al.](#),<sup>32</sup> we have previously used

$$\text{se}(\hat{\gamma}_j)^2 = \frac{\left(\check{y} - \hat{\gamma}_j \check{G}_j\right)^T \left(\check{y} - \hat{\gamma}_j \check{G}_j\right)}{\frac{(n-K-1) \check{G}_j^T \check{G}_j \approx \check{y}^T \check{y}}{n \check{G}_j^T \check{G}_j \approx \frac{\text{var}(\check{y})}{\text{var}(\check{G}_j)}}},$$

where  $\hat{\gamma}_j$  is the marginal effect of variant  $j$ , and where  $\check{y}$  and  $\check{G}_j$  are the vectors of phenotypes and genotypes for variant  $j$  residualized from  $K$  covariates, e.g., centering them. The first approximation expects  $\hat{\gamma}_j$  to be small, while the second approximation assumes the effects from covariates are small. However, we have found here that some variants can have very large effects, e.g., one variant explains about 30% of the variance in bilirubin log-concentration. Then, instead we compute

$$\begin{aligned} \left(\check{y} - \hat{\gamma}_j \check{G}_j\right)^T \left(\check{y} - \hat{\gamma}_j \check{G}_j\right) &= \check{y}^T \check{y} - 2\hat{\gamma}_j \check{G}_j^T \check{y} + \hat{\gamma}_j^2 \check{G}_j^T \check{G}_j \\ &= \check{y}^T \check{y} - \hat{\gamma}_j^2 \check{G}_j^T \check{G}_j, \end{aligned}$$

which now gives

$$(n-K-1) \text{se}(\hat{\gamma}_j)^2 = \frac{\check{y}^T \check{y} - \hat{\gamma}_j^2 \check{G}_j^T \check{G}_j}{\frac{\check{G}_j^T \check{G}_j \approx \check{y}^T \check{y}}{\check{G}_j^T \check{G}_j - \hat{\gamma}_j^2 \approx \frac{\text{var}(\check{y})}{\text{var}(\check{G}_j)} - \hat{\gamma}_j^2}},$$

finally giving (note the added term  $\hat{\gamma}_j^2$ )

$$\text{sd}(\check{G}_j) \approx \frac{\text{sd}(\check{y})}{\sqrt{n \text{se}(\hat{\gamma}_j)^2 + \hat{\gamma}_j^2}}. \quad (\text{Equation 1})$$

[Figure S23](#) shows that the updated formula [Equation 1](#) is better; we now use it in the code of LDpred2, and also recommend using it for the QC procedure proposed in [Privé et al.](#)<sup>32</sup>

## Using more than HapMap3 variants in LDpred2

Here we also run LDpred2 using more than HapMap3 variants, based on a set of 8M+ variants (see above). However, LDpred2 cannot be run on 8M variants because the implementation is quadratic with the number of variants in terms of time and

**Table 1. Overview of sets of individuals used in this study**

Set	UK1	UK2	UK3	Poland	Italy	Iran	India	China	Caribbean	Nigeria	Ashkenazi Jewish
Training 1	367,063	24,061	–	–	–	–	–	–	–	–	–
Test 1	–	–	20,000	4,136	6,660	1,200	6,331	1,810	2,484	3,924	1,709
Training 2	367,063	–	–	4,136	6,660	1,200	6,331	1,810	–	3,924	–
Test 2	–	–	20,000	–	–	–	–	–	2,484	–	–

In total, 439,378 unrelated individuals are used here. Most analyses in this paper use UK1 + UK2 (391,124 individuals) as training set and the other groups as test sets. Secondary analyses in section “Training with a mixture of ancestries” involve multiple ancestry training and keep only the UK3 and Caribbean groups as test sets; UK2 is removed from the training so that sample size from training 2 is the same as training 1 (391,124 individuals). Note that the names of the first eight ancestry groups we define here refer to the country names from the UK Biobank (field 20115) that we use to define the centers of each ancestry group; therefore, these groups also include individuals from nearby countries. For example, the “United Kingdom” ancestry group also includes many individuals who self-identify as Irish, and the “India” ancestry group also includes many individuals who self-identify as Pakistani (Note A).

memory requirements. Thus, we employ another strategy consisting in keeping only the 1M most significant variants. To correct for winner’s curse, we employ the maximum likelihood estimator used in Zhong and Prentice<sup>45</sup> and Shi et al.:<sup>46</sup>

$$Z = Z^* + \frac{\phi(Z^* - Z_{thr}) - \phi(-Z^* - Z_{thr})}{\Phi(Z^* - Z_{thr}) + \Phi(-Z^* - Z_{thr})},$$

where  $\phi$  is the standard normal density function,  $\Phi$  is the standard normal cumulative density function,  $Z$  is the Z-score obtained from the GWAS,  $Z_{thr}$  is the threshold used on (absolute) Z-scores for filtering, and  $Z^*$  is the corrected Z-score that we estimate and use. As input for LDpred2, instead of using  $\beta$  (along with  $SE(\beta)$  and  $N$ ), we use  $\beta^* = \beta \cdot Z^* / Z$  where  $Z = \beta / SE(\beta)$ . This is now implemented in function `snp_thr_correct` of package `bigsnpr`.

### Performance metric

Here we use the partial correlation as the performance metric, which is the correlation between the PGS and the phenotype after they have been both residualized using the covariates used in this paper, i.e., sex, age, birth date, deprivation index, and 16 PCs. To derive 95% confidence intervals for these correlations, we use Fisher’s Z-transformation. We implement this in function `pcor` of R package `bigstatsr` and use it here.

## Results

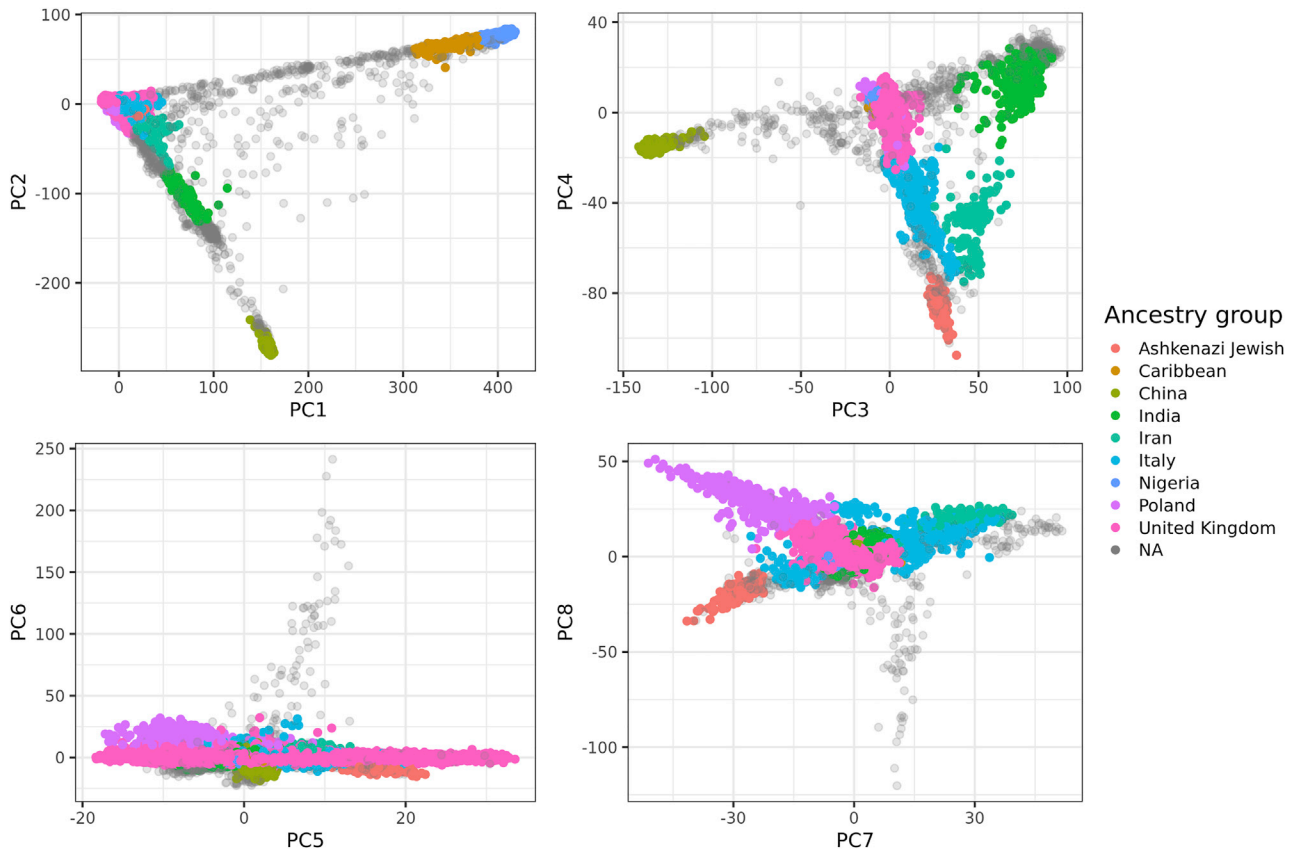
### Overview of study

Here, we use the UK Biobank (UKBB) data only.<sup>5</sup> We first infer nine ancestry groups in the UKBB. Then we use 391,124 individuals of Northwestern European ancestry to train polygenic scores (PGSs) for 245 phenotypes (about half being diseases; see categories in Figure S1) based on UKBB individual-level genotypes and phenotypes, and we assess portability of these PGSs in the remaining individuals of diverse ancestries (Table 1). As additional analyses, we also investigate using more variants than the HapMap3 variants used in the main analyses, and we train models using a mixture of multiple ancestries. To derive PGSs in this study, we use two different methods, penalized regression and LDpred2-auto, and finally compare them.

### Ancestry grouping

We investigate various approaches to classify individuals in ancestry groups based on principal component analysis (PCA) of genome-wide genotype data. Detailed results can be found in the corresponding Note A; we recall main results here. First, we show that (squared) Euclidean distances in the PCA space of genetic data are approximately proportional to  $F_{ST}$  between populations, and we therefore recommend using this simple distance. We also provide evidence that using only two PCs, or even four PCs, is not enough to distinguish between some less-distant populations, and we recommend using all PCs visually capturing some population structure. Then, we use this PCA-based distance to infer ancestry in the UK Biobank and the POPRES datasets. We propose two solutions to do so, either relying on projection of PCs to reference populations such as the 1000 Genomes Project, or by directly using internal data only. We show that these solutions are simple, robust, and effective methods for inferring global ancestry and for grouping genetically homogeneous individuals.

Here, we first use the second solution presented in Note A, relying on PCs computed within the UK Biobank and individual information on the countries of birth, for inferring the first eight ancestry groups presented in Table 1. These groups were chosen on the basis of being distant enough from the other groups, and including enough individuals (e.g., >1,000) to draw meaningful conclusions. Note that the names of the ancestry groups we define here refer to the country names from the UK Biobank (field 20115) that we use to define the centers of each ancestry group; therefore, these groups also include individuals from nearby countries. For example, the “United Kingdom” ancestry group also includes many individuals who self-identify as Irish, and the “India” ancestry group also includes many individuals who self-identify as Pakistani (Note A). Then, for inferring the “Ashkenazi Jewish” ancestry group, we use the first solution, projecting UKBB individuals onto the PCA space of a reference dataset composed of many Jewish and non-Jewish individuals.<sup>43</sup> We identify a ninth group of 1,709 unrelated individuals, which is entirely non-overlapping with the other eight groups previously defined (Material and methods). This



**Figure 1. The first eight PC scores of the UK Biobank (field 22009) colored by the homogeneous ancestry group we infer for these individuals**

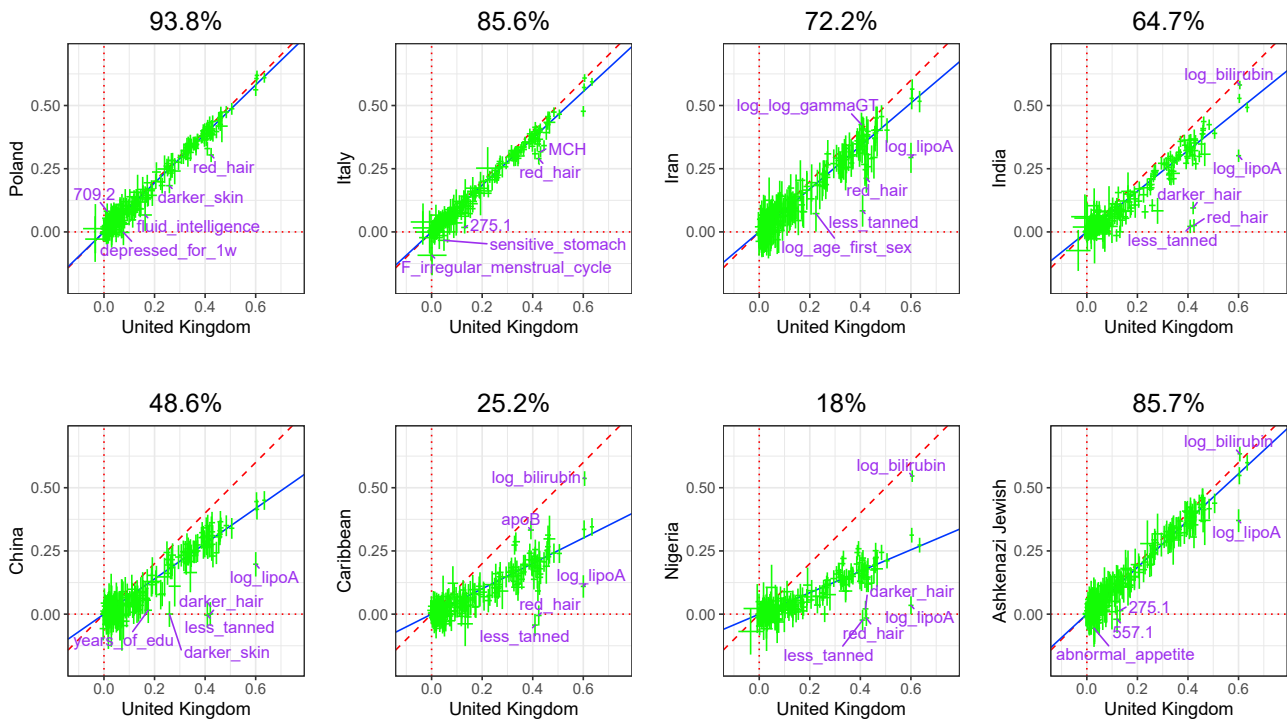
Only 50,000 individuals are represented at random. “NA” means that the corresponding individual is not categorized in any of the nine ancestry groups.

group is largely overlapping with the 1,719 presumably British Jews identified from IBD segments in Naseri et al.<sup>47</sup> (personal correspondence with the authors). Finally, we run ADMIXTURE (with  $k = 8$  and  $k = 5$ ) on 200 individuals from each of the nine ancestry groups defined here.<sup>24</sup> The results are consistent with the PCA analysis (Figure 1), e.g., showing that the Caribbean group we define is mostly composed of admixed individuals with mostly African ancestry and some small percentage of European ancestry (Figure S2). Moreover, the other groups we define have distinct ADMIXTURE profiles (consistently with being distinct on PCA), except for the “United Kingdom” and “Poland” ancestry groups, which cannot be distinguished based on this analysis.

#### Portability of polygenic scores to other ancestries

Figure 2 presents the results when fitting penalized regression using a training set composed of Northwestern European individuals from the UK Biobank (“United Kingdom,” hereinafter also referred to as “the UK individuals” or “the UK” for simplicity purposes) and testing in nine different ancestry groups from the same cohort (Table 1). Averaged over 245 phenotypes, compared to prediction performance in individuals of Northwestern European ancestry, relative predictive ability in terms of partial- $r^2$  (Material and

methods) is 93.8% in the “Poland” ancestry group (North-east Europe), 85.6% in “Italy” (South Europe), 72.2% in “Iran” (Middle East), 64.7% in “India” (South Asia), 48.6% in “China” (East Asia), 25.2% in the “Caribbean,” 18% in “Nigeria” (West Africa), and 85.7% for the Ashkenazi Jewish group. As a follow-up analysis to ensure that this drop in performance in other ancestries is not due to differences in imputation quality across ancestries, we perform the same analysis for 83 of the continuous phenotypes using high-quality genotyped variants only (Material and methods) instead of the (mostly imputed) HapMap3 variants; results are highly consistent (Figure S3). We also run the previous follow-up analysis while removing third-degree relatives, which leaves us with 349,991 individuals for training (instead of 391,124) and 43,631 for testing (instead of 46,545); results are practically unchanged (Figure S4). These results are also very similar when using LDpred2-auto instead of penalized regression for training predictive models for all phenotypes (Figure S5). A few phenotypes deviate from this global trend, e.g., prediction of bilirubin concentration ranges between 0.537 and 0.619 (partial- $r$ ) for all ancestries except for “China,” for which it is 0.415 (95% CI: 0.374–0.453, see Material and methods). In contrast, for example for hair and skin color, partial correlations decrease quickly and are not significantly different



**Figure 2. Partial correlation and 95% CI in the UK test set versus in a test set from another ancestry group**

Each point represents a phenotype and training has been performed with penalized regression on UK individuals (training 1 in Table 1) and HapMap3 variants. The slope (in blue) is computed using Deming regression accounting for standard errors in both x and y, fixing the intercept at 0. The square of this slope is provided above each plot, which we report as the relative predictive performance compared to testing in the “United Kingdom” ancestry group.

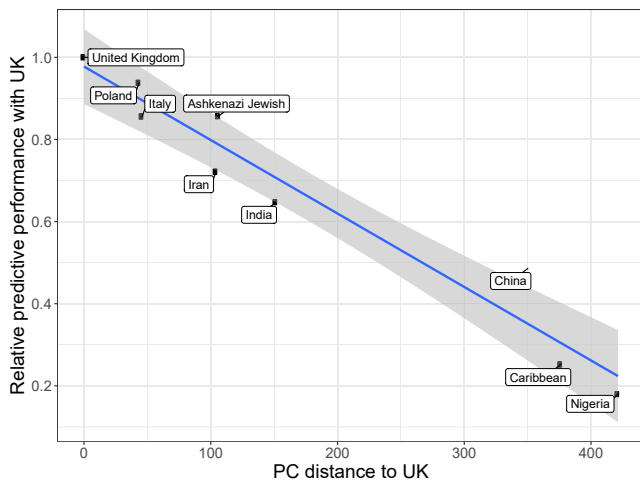
from 0 for both “China” and “Nigeria,” while of 0.420 (95% CI: 0.409–0.432) for “darker hair” in the “United Kingdom” ancestry group (Figure 2). Overall, relative predictive performance decreases approximately linearly with PC distance to the training set (Figure 3). A similar pattern is observed when computing PCA based on more balanced ancestry groups, as recommended in Privé et al.<sup>20</sup> (Figure S6).

### Using more than HapMap3 variants?

We investigate some of the outlier phenotypes in Figure 2, especially the ones from blood biochemistry which have some variants with large effects. We hypothesize that using a denser set of variants could improve tagging of the causal variants with large effect sizes, resulting in an improved prediction in all ancestries. We focus on “total bilirubin,” “lipoprotein(a)” (lipoA), and “apolipoprotein B” (apoB). We perform a localized GWAS which includes all variants around the most significant variant (hereinafter denoted as “top hit”) from the GWAS in the training set 1 (UK individuals and HapMap3 variants only) in each of the first eight ancestry groups defined here. More precisely, we include all variants with an imputation INFO score larger than 0.3 and within a window of 500 kb from the HapMap3 top hit in the UK; there are approximately 30K such variants for all three phenotypes. For bilirubin, the overall top hit is a HapMap3 variant and explains around 30% of the phenotypic variance (Figure S8). Effects from the three top hits are fairly consistent within all ancestry

groups (Figure S9) explaining why genetic prediction is highly consistent in all ancestries, except for “China” (Figure 2), for which these variants are rarer. For lipoA, results are very different across ancestries; HapMap3 variants are far from being the top hits for the UK individuals, where the top HapMap3 variant explains 5% of phenotypic variance compared to 29% for the (non-HapMap3) top hit (Figure 4). Note that this top hit is more than 200 kb away from the HapMap3 top hit from the UK group. Moreover, the three top hits for lipoA do not have very consistent effect sizes across ancestries (Figure S10). Finally, for apoB, effects from the three top hits, which are not part of HapMap3 variants, are fairly consistent across ancestries and explain up to 8.5% of the phenotypic variance (Figures S11 and S12).

We then investigate whether the use of a larger set of variants than the HapMap3 set is beneficial; we use more than 8M common variants (Material and methods) and apply LDpred2-auto after restricting to the 1M most significant variants and applying winner’s curse correction (Material and methods). Except for lipoA for which we get a large improvement in predictive accuracy compared to using HapMap3 variants only, it is not beneficial for the other seven phenotypes analyzed here (Figure 5). Remarkably, while the partial correlation for lipoA is about 75% in the UK test set when using this prioritized set of variants, it is still not different from 0 when applied to the “Nigeria” group. For height and BMI, estimated SNP heritability is



**Figure 3. Relative variance explained compared to the UK versus PC distance from the UK**

PC distances are computed using Euclidean distance between geometric medians of the first 16 reported PC scores (field 22009) of each ancestry group. Relative performance values are the ones reported in Figure 2. The slope and standard errors are computed internally by function `geom_smooth(method = "lm")` of R package `ggplot2`.

reduced when using this set of most significant variants only, and all these variants are estimated to be causal, i.e., the estimate of the proportion of causal variants  $p$  is 1 (Table S1). As height and BMI are very polygenic traits ( $p$  is estimated to be  $\sim 2\%$  and  $\sim 4\%$ , respectively, when using HapMap3 variants), contribution from less significant causal variants is missed due to this thresholding selection. For the three binary phenotypes of breast cancer (phecode: 174.1), prostate cancer (185), and coronary artery disease (411.4), although heritability estimates are larger when using this set of prioritized variants (Table S1), predictive accuracy does not improve compared to when using HapMap3 variants (Figure 5).

### Training with a mixture of ancestries

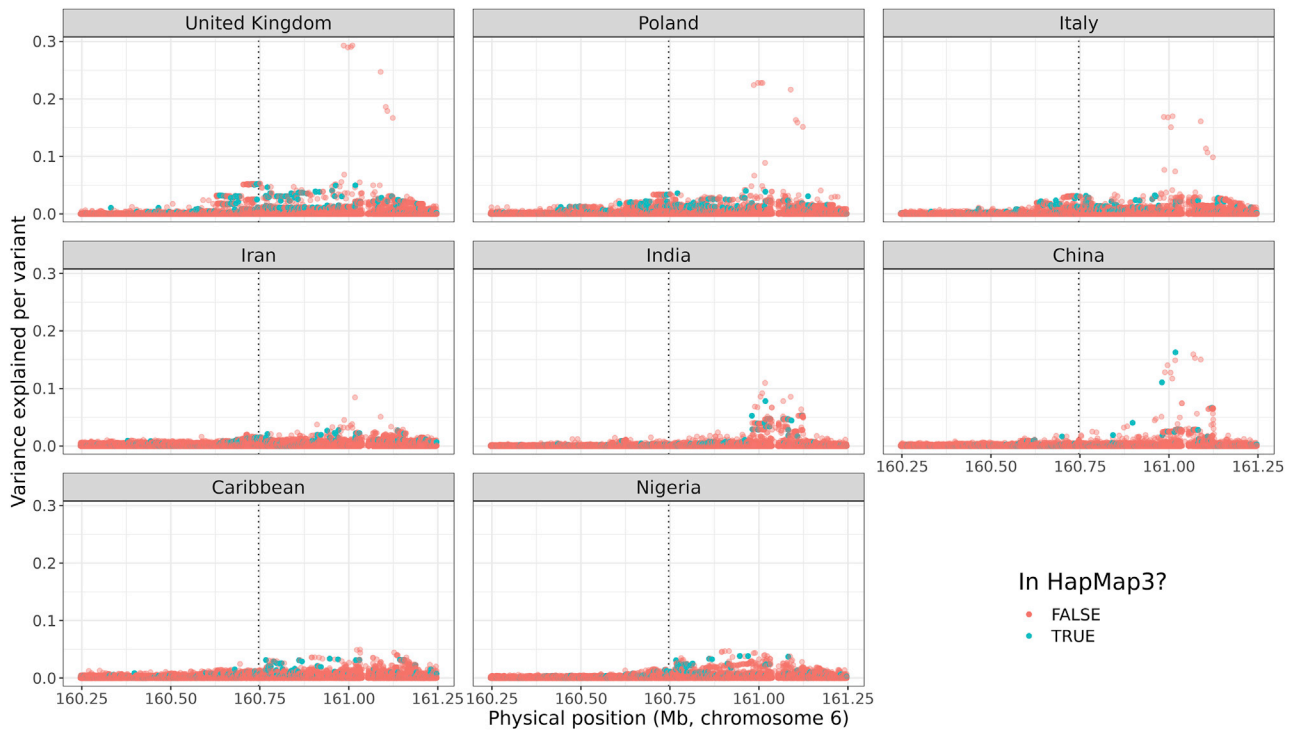
We hypothesize that using individuals from diverse ancestries could improve tagging of the causal variants, resulting in an improved prediction in all ancestries. Indeed, power improvements for both association and prediction have been reported when using even a small set of individuals from different ancestries.<sup>11,48,49</sup> Here we use all ancestry groups except for the Caribbean and Ashkenazi for training penalized regressions; we remove the same number of UK individuals to keep the same training sample size as before (training 2 in Table 1). We recall that Caribbean individuals are mostly admixed between African, European, and Native American ancestries,<sup>50</sup> which are almost all represented here in the training set 2. In Figure S13, we investigate nine phenotypes of interest, either because they are highly studied diseases or are outliers in Figure 2: breast cancer (phecode: 174.1), prostate cancer (186), type 2 diabetes (250.2), hypertension (401),

coronary artery disease (411.4), skin tone, total bilirubin concentration, lipoprotein(a) concentration, and years of education. We predict in the test sets from the UK and the Caribbean (test set 2); overall, the predictive performance is highly similar when using this multi-ancestry training compared to when using only UK individuals, in both the UK and the Caribbean target samples. Prediction is only improved for lipoprotein(a) concentration when the mixed ancestry training data is used in application to the Caribbean target data (Figure S13). Discrepancies between our results and results from Márquez-Luna et al.<sup>51</sup> and Cavazos and Witte<sup>11</sup> may be explained by the fact that we use the exact same sample size when training with multiple ancestries (by removing some UK individuals; see Table 1), whereas these studies use extra (non-European) individuals, making it hard to know if the improved predictions come from using non-European individuals, or just from using more individuals. We also run the newly developed PRS-CSx method<sup>49</sup> using individuals from training 2, deriving the GWAS summary statistics from the UK Biobank individual-level data (as for LDpred2-auto). PRS-CSx provides lower predictive performance than using the penalized regression on training 2 for both the UK and Caribbean test sets, except when predicting years of education for both sets as well as “darker skin” and coronary artery disease (phecode 411.4) in the Caribbean test set (Figure S13). Predictive performance of PRS-CSx is particularly lower for traits with large effects (bilirubin and lipoprotein(a) concentrations) and moderate effects (breast and prostate cancers; phecodes 174.1 and 185).

### Comparison of predictive models

Penalized regression and LDpred2-auto provide approximately similar predictive performance across all traits and ancestries considered here (Figure S14); there are only four pairs of phenotype-ancestry (out of nearly 2,000 pairs) for which 95% CIs for partial- $r$  from penalized regression and LDpred2 are not overlapping: “615: endometriosis” in the “China” ancestry group with 0.065 (0.0074 to 0.122) versus  $-0.051$  ( $-0.108$  to 0.0068); “hard falling asleep” in UK with  $-0.0349$  ( $-0.742$  to 0.0045) versus 0.071 (0.031 to 0.110); height in UK with 0.634 (0.626 to 0.643) versus 0.613 (0.605 to 0.622); and log-bilirubin in “Nigeria” with 0.546 (0.523 to 0.569) versus 0.475 (0.449 to 0.500). For prediction in UK ancestry, penalized regression tends to provide better predictive performance than LDpred2 for phenotypes for which partial- $r > 0.3$ , and LDpred2 tends to outperform penalized regression for phenotypes harder to predict (Figure S14).

Both methods allow for fitting sparse effects, i.e., some resulting effects are exactly 0. Sparse models may be beneficial because they may be more easily implemented. The sparse option in LDpred2-auto provides similar performance as LDpred2-auto without this option (Figure S15). Sparsity of resulting effects follows a very different pattern for



**Figure 4. Zoomed Manhattan plot for lipoprotein(a) concentration**

The phenotypic variance explained per variant is computed as  $r^2 = t^2/(n+t^2)$ , where  $t$  is the t-score from GWAS and  $n$  is the degrees of freedom (the sample size minus the number of variables in the model, i.e., the covariates used in the GWAS, the intercept, and the variant). The GWAS includes all variants with an imputation INFO score larger than 0.3 and within a 500 kb radius around the top hit from the GWAS performed in the UK training set and on the HapMap3 variants, represented by a vertical dotted line.

penalized regression compared to LDpred2-auto-sparse. Indeed, penalized regression tends not to include variants if it is uncertain that they have a non-zero effect, i.e., when effects are very small and prediction is difficult (Figure S16). In contrast, LDpred2-auto-sparse tends not to discard variants, only when  $h^2$  is large enough it sets lots of effects to 0 if  $p$  is small (Figure S17). Finally, running each penalized regression model takes between a few minutes and a few days depending on the number of non-zero effects in the resulting model (Figure S18). In contrast, LDpred2-auto should take the same computation time for all phenotypes; it completed under 7 h for most phenotypes (Figure S19).

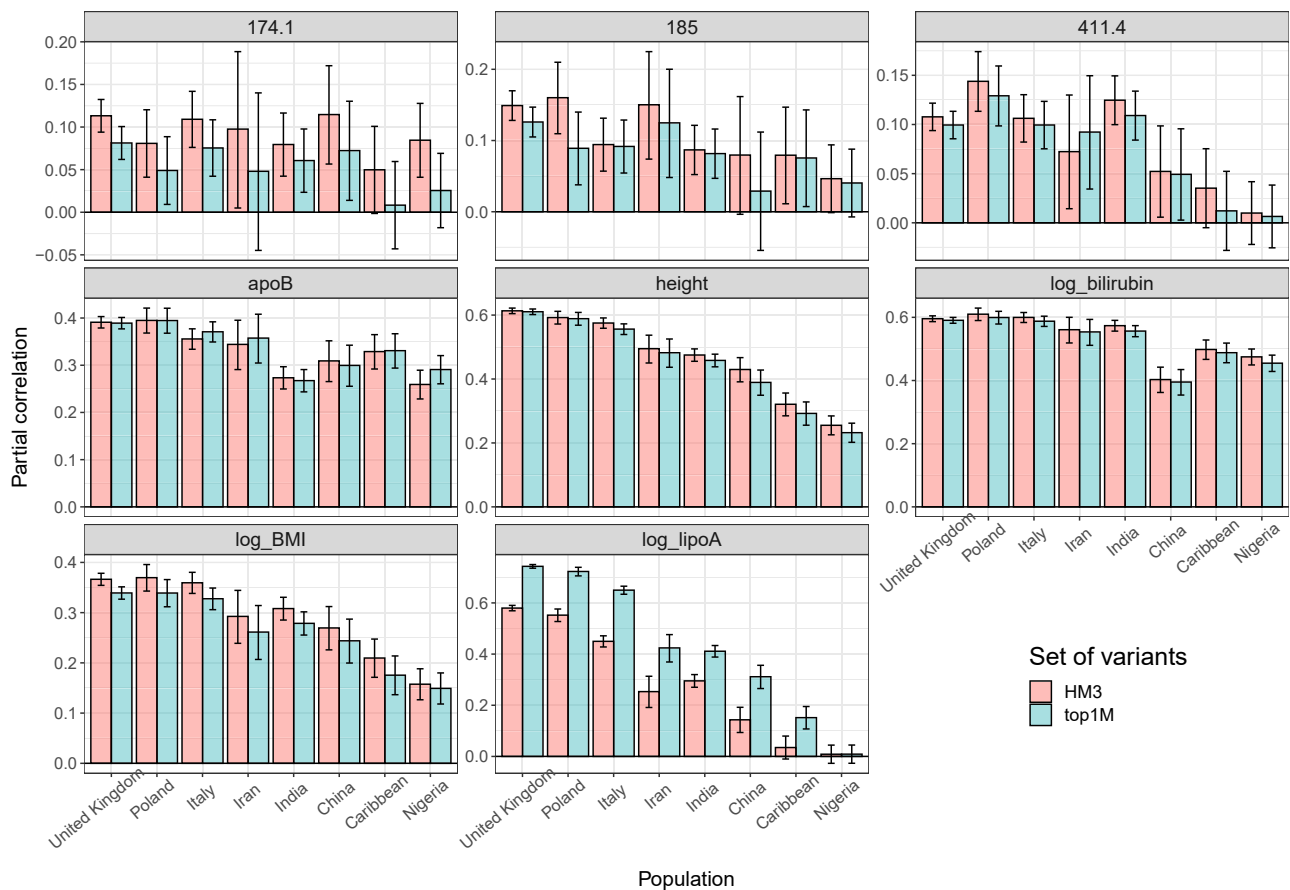
## Discussion

In this paper, we have conducted an extensive assessment of PGS portability across ancestries using hundreds of phenotypes. Our analysis demonstrates a canonical relation between genetic distance and predictive performance for most phenotypes. The reported poor portability is in agreement with three previous studies;<sup>9,52,53</sup> we show a relative predictive performance compared to Europeans of ~18% for Africans (versus 22%, 42%, and 24%), ~49% for East Asians (versus 50%, 95%, and 64%), and ~65% for South Asians (versus 60%, 62.5%, and 72%). However, our results also provide a significant addition to the current literature

in many ways. First, we show that the portability issue remains strong even when PGSs are derived and applied in the same cohort. Second, the presented results are averaged over 245 phenotypes, which is much more than what has been typically used, and should capture a broad range of the phenotypic spectrum. Portability results are highly consistent across most phenotypes (with a few exceptions) and could therefore be used to predict the expected loss of accuracy for other phenotypes. Third, we provide this result at a finer scale than the usual continental level by proposing a simple, robust, and effective method for grouping UKBB individuals in nine ancestry groups. This allows us to show, for example, that predictive performance already decreases within Europe with only ~94% for Northeast Europe and ~86% for South Europe of the performance reached within Northwest Europe.

We showcase two methods for deriving polygenic scores when large individual-level datasets are available. Although LDpred2-auto is a method based on summary statistics, it provides good predictive performance compared to penalized regression, when applied to individual-level data. Moreover, portability results shown here are similar when using either the individual-level penalized regression or the summary statistics based LDpred2 method. Fitting of penalized models is relatively fast when using 1M HapMap3 variants. We have also tried fitting penalized regression using 8M variants (>3 TB of data); this was possible but took several days for the phenotypes we tried,





**Figure 5. Predictive performance with LDpred2-auto for eight phenotypes, when using either HapMap3 variants or the 1M most significant variants**

One phenotype shown in each panel. Bars represent the 95% confidence intervals. Phencode 174.1: breast cancer; 185: prostate cancer; 411.4: coronary artery disease. HM3, HapMap3; top1M, the 1M most significant variants out of more than 8M common variants (see [Material and methods](#)).

so we have not investigated this further. To the best of our knowledge, we use the most efficient penalized regression implementation currently available. Recently, Qian et al.<sup>7</sup> proposed *snnpnet*, a new R package for fitting penalized regressions on large individual-level genetic datasets, but we have found it to be much less efficient than R package *bigstatsr* on UKBB data ([Note B](#)). As for LDpred2, it currently cannot be run using 8M variants, but we show how to use a subset of 1M prioritized variants out of these 8M. Using this new set of variants provides a large improvement in predicting lipoprotein(a) concentration (lipoA), but not for the other seven phenotypes studied in this analysis. This improvement for lipoA is not surprising given that the top HapMap3 variant explains 5% of phenotypic variance compared to 29% for the (non-HapMap3) top hit ([Figure 4](#)).

Here we use only the UK Biobank data to fit polygenic scores. We do not use external information such as functional annotations; those could be used to improve the heritability model assumed by predictive methods in order to improve predictive performance.<sup>54</sup> Moreover, we do not use external summary statistics, which means that polygenic scores derived from large GWAS meta-analyses would probably outperform the ones we derived here. Neverthe-

less, Albiñana et al.<sup>55</sup> have shown that an efficient strategy to improve predictive ability of polygenic scores consists in combining two different polygenic scores, one derived using external summary statistics and another one derived using internal individual-level data. Therefore, the polygenic scores we derived here could be combined with polygenic scores derived using external summary statistics; we will release these PGSs publicly and share them in databases such as the PGS Catalog and the Cancer-PRSweb.<sup>56,57</sup>

#### Data and code availability

The UK Biobank data are available through a procedure described at <https://www.ukbiobank.ac.uk/using-the-resource/>. All code used for this paper is available at <https://github.com/privefl/UKBB-PGS/tree/main/code>. Links to the code used for the [Notes A and B](#) are provided there. Code to reproduce our nine ancestry groups is available at <https://github.com/privefl/UKBB-PGS#code-to-reproduce-ancestry-groups>.

We have extensively used R packages *bigstatsr* and *bigsnpr*<sup>19</sup> for analyzing large genetic data, packages from the future framework<sup>58</sup> for easy scheduling and parallelization of analyses on the HPC cluster, and packages from the *tidyverse* suite<sup>59</sup> for shaping

and visualizing results. We have also used R package deming for fitting Deming regressions.

## Supplemental information

Supplemental information can be found online at <https://doi.org/10.1016/j.ajhg.2021.11.008>.

## Acknowledgments

Authors thank the reviewers for their comments and suggestions. Authors thank Abdel Abdellaoui for his help with defining the “years of education” phenotype and Alex Diaz-Papkovich and others for their useful feedback on the ancestry inference. Authors thank GenomeDK and Aarhus University for providing computational resources and support that contributed to these research results. This research has been conducted using the UK Biobank Resource under Application Number 58024.

E.P. and B.J.V. are supported by the Danish National Research Foundation (Niels Bohr Professorship to Prof. John McGrath) and also acknowledge the Lundbeck Foundation Initiative for Integrative Psychiatric Research, iPSYCH (R248-2017-2003). B.J.V. is also supported by a Lundbeck Foundation Fellowship (R335-2019-2339).

## Declaration of interests

S.C. is a paid consultant to MyHeritage. The other authors declare no competing interests.

Received: May 26, 2021

Accepted: November 4, 2021

Published: January 6, 2022

## Web resources

bigsnpr, tutorial on LDpred2, <https://privefl.github.io/bigsnpr/articles/LDpred2.html>

bigstatsr, tutorial on penalized regressions, <https://privefl.github.io/bigstatsr/articles/penalized-regressions.html>

PGS Catalog, effect sizes of PGSs derived here, <https://www.pgscatalog.org/publication/PGP000263/>

UK Biobank, quality control information on genetic variants, [https://biobank.ctsu.ox.ac.uk/crystal/crystal/auxdata/ukb\\_snp\\_qc.txt](https://biobank.ctsu.ox.ac.uk/crystal/crystal/auxdata/ukb_snp_qc.txt)

UKBB-PGS, description of the 245 phenotypes used in this study, <https://github.com/privefl/UKBB-PGS/blob/main/phenotype-description.xlsx>

UKBB-PGS, other information on the phenotypes (e.g., sample sizes), <https://github.com/privefl/UKBB-PGS/blob/main/phenotype-info.xlsx>

## References

- Choi, S.W., Mak, T.S.-H., and O'Reilly, P.F. (2020). Tutorial: a guide to performing polygenic risk score analyses. *Nat. Protoc.* *15*, 2759–2772.
- de los Campos, G., Gianola, D., and Allison, D.B. (2010). Predicting genetic predisposition in humans: the promise of whole-genome markers. *Nat. Rev. Genet.* *11*, 880–886.
- Abraham, G., Tye-Din, J.A., Bhalala, O.G., Kowalczyk, A., Zobel, J., and Inouye, M. (2014). Accurate and robust genomic prediction of celiac disease using statistical learning. *PLoS Genet.* *10*, e1004137.
- Privé, F., Aschard, H., and Blum, M.G.B. (2019). Efficient implementation of penalized regression for genetic risk prediction. *Genetics* *212*, 65–74.
- Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L.T., Sharp, K., Motyer, A., Vukcevic, D., Delaneau, O., O'Connell, J., et al. (2018). The UK Biobank resource with deep phenotyping and genomic data. *Nature* *562*, 203–209.
- Loh, P.-R., Kichaev, G., Gazal, S., Schoech, A.P., and Price, A.L. (2018). Mixed-model association for biobank-scale datasets. *Nat. Genet.* *50*, 906–908.
- Qian, J., Tanigawa, Y., Du, W., Aguirre, M., Chang, C., Tibshirani, R., Rivas, M.A., and Hastie, T. (2020). A fast and scalable framework for large-scale and ultrahigh-dimensional sparse regression with application to the UK Biobank. *PLoS Genet.* *16*, e1009141.
- Scutari, M., Mackay, I., and Balding, D. (2016). Using genetic distance to infer the accuracy of genomic prediction. *PLoS Genet.* *12*, e1006288.
- Wang, Y., Guo, J., Ni, G., Yang, J., Visscher, P.M., and Yengo, L. (2020). Theoretical and empirical quantification of the accuracy of polygenic scores in ancestry divergent populations. *Nat. Commun.* *11*, 3865.
- Bitarello, B.D., and Mathieson, I. (2020). Polygenic scores for height in admixed populations. *G3 (Bethesda)* *10*, 4027–4036.
- Cavazos, T.B., and Witte, J.S. (2021). Inclusion of variants discovered from diverse populations improves polygenic risk score transferability. *HGG Adv* *2*, 100017.
- Sinnott-Armstrong, N., Tanigawa, Y., Amar, D., Mars, N., Benner, C., Aguirre, M., Venkataraman, G.R., Wainberg, M., Ollila, H.M., Kiiskinen, T., et al. (2021). Genetics of 35 blood and urine biomarkers in the UK Biobank. *Nat. Genet.* *53*, 185–194.
- Berg, J.J., Harpak, A., Sinnott-Armstrong, N., Joergensen, A.M., Mostafavi, H., Field, Y., Boyle, E.A., Zhang, X., Racimo, F., Pritchard, J.K., and Coop, G. (2019). Reduced signal for polygenic adaptation of height in UK Biobank. *eLife* *8*, e39725.
- Sohail, M., Maier, R.M., Ganna, A., Bloemendal, A., Martin, A.R., Turchin, M.C., Chiang, C.W., Hirschhorn, J., Daly, M.J., Patterson, N., et al. (2019). Polygenic adaptation on height is overestimated due to uncorrected stratification in genome-wide association studies. *eLife* *8*, e39702.
- Haworth, S., Mitchell, R., Corbin, L., Wade, K.H., Dudding, T., Budu-Aggrey, A., Carslake, D., Hemani, G., Paternoster, L., Smith, G.D., et al. (2019). Apparent latent structure within the UK Biobank sample has implications for epidemiological analysis. *Nat. Commun.* *10*, 333.
- Price, A.L., Patterson, N.J., Plenge, R.M., Weinblatt, M.E., Shadick, N.A., and Reich, D. (2006). Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* *38*, 904–909.
- Novembre, J., Johnson, T., Bryc, K., Kutalik, Z., Boyko, A.R., Auton, A., Indap, A., King, K.S., Bergmann, S., Nelson, M.R., et al. (2008). Genes mirror geography within Europe. *Nature* *456*, 98–101.
- Abraham, G., Qiu, Y., and Inouye, M. (2017). FlashPCA2: principal component analysis of Biobank-scale genotype datasets. *Bioinformatics* *33*, 2776–2778.

19. Privé, F., Aschard, H., Ziyatdinov, A., and Blum, M.G.B. (2018). Efficient analysis of large-scale genome-wide data with two R packages: bigstatsr and bigsnpr. *Bioinformatics* 34, 2781–2787.
20. Privé, F., Luu, K., Blum, M.G.B., McGrath, J.J., and Vilhjálms-son, B.J. (2020). Efficient toolkit implementing best practices for principal component analysis of population genetic data. *Bioinformatics* 36, 4449–4457.
21. Zhang, D., Dey, R., and Lee, S. (2020). Fast and robust ancestry prediction using principal component analysis. *Bioinformatics* 36, 3439–3446.
22. Chen, C.-Y., Pollack, S., Hunter, D.J., Hirschhorn, J.N., Kraft, P., and Price, A.L. (2013). Improved ancestry inference using weights from external reference panels. *Bioinformatics* 29, 1399–1406.
23. Byun, J., Han, Y., Gorlov, I.P., Busam, J.A., Seldin, M.F., and Amos, C.I. (2017). Ancestry inference using principal component analysis and spatial analysis: a distance-based analysis to account for population substructure. *BMC Genomics* 18, 789.
24. Alexander, D.H., Novembre, J., and Lange, K. (2009). Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.* 19, 1655–1664.
25. Lawson, D.J., Hellenthal, G., Myers, S., and Falush, D. (2012). Inference of population structure using dense haplotype data. *PLoS Genet.* 8, e1002453.
26. Raj, A., Stephens, M., and Pritchard, J.K. (2014). fastSTRUCTURE: variational inference of population structure in large SNP data sets. *Genetics* 197, 573–589.
27. Frichot, E., Mathieu, E., Trouillon, T., Bouchard, G., and François, O. (2014). Fast and efficient estimation of individual ancestry coefficients. *Genetics* 196, 973–983.
28. Haller, T., Leitsalu, L., Fischer, K., Nuotio, M.-L., Esko, T., Boomsma, D.I., Kyvik, K.O., Spector, T.D., Perola, M., and Metspalu, A. (2017). MixFit: Methodology for computing ancestry-related genetic scores at the individual level and its application to the Estonian and Finnish population studies. *PLoS ONE* 12, e0170325.
29. Cheng, J.Y., Mailund, T., and Nielsen, R. (2017). Fast admixture analysis and population tree estimation for SNP and NGS data. *Bioinformatics* 33, 2148–2155.
30. Jin, Y., Schaffer, A.A., Feolo, M., Holmes, J.B., and Kattman, B.L. (2019). GRAF-pop: a fast distance-based method to infer subject ancestry from multiple genotype datasets without principal components analysis. *G3 (Bethesda)* 9, 2447–2461.
31. Cabrerós, I., and Storey, J.D. (2019). A likelihood-free estimator of population structure bridging admixture models and principal components analysis. *Genetics* 212, 1009–1029.
32. Privé, F., Arbel, J., and Vilhjálms-son, B.J. (2020). LDpred2: better, faster, stronger. *Bioinformatics* 36, 5424–5431.
33. Bybjerg-Grauholm, J., Pedersen, C.B., Baekvad-Hansen, M., Pedersen, M.G., Adamsen, D., Hansen, C.S., Agerbo, E., Grove, J., Als, T.D., Schork, A.J., et al. (2020). The iPSYCH2015 case-cohort sample: updated directions for unravelling genetic and environmental architectures of severe mental disorders. *medRxiv*. <https://doi.org/10.1101/2020.11.30.20237768>.
34. Carroll, R.J., Bastarache, L., and Denny, J.C. (2014). R PheWAS: data analysis and plotting tools for phenome-wide association studies in the R environment. *Bioinformatics* 30, 2375–2376.
35. Wu, P., Gifford, A., Meng, X., Li, X., Campbell, H., Varley, T., Zhao, J., Carroll, R., Bastarache, L., Denny, J.C., et al. (2019). Mapping ICD-10 and ICD-10-CM codes to phecodes: work-flow development and initial evaluation. *JMIR Med. Inform.* 7, e14325.
36. Zhou, W., Nielsen, J.B., Fritsche, L.G., Dey, R., Gabrielsen, M.E., Wolford, B.N., LeFaive, J., VandeHaar, P., Gagliano, S.A., Gifford, A., et al. (2018). Efficiently controlling for case-control imbalance and sample relatedness in large-scale genetic association studies. *Nat. Genet.* 50, 1335–1341.
37. Gagliano Taliun, S.A., VandeHaar, P., Boughton, A.P., Welch, R.P., Taliun, D., Schmidt, E.M., Zhou, W., Nielsen, J.B., Willer, C.J., Lee, S., et al. (2020). Exploring and visualizing large-scale genetic associations by using PheWeb. *Nat. Genet.* 52, 550–552.
38. Kunert-Graf, J.M., Sakhanenko, N.M., and Galas, D.J. (2020). Allele frequency mismatches and apparent mismappings in UK Biobank SNP data. *bioRxiv*. <https://doi.org/10.1101/2020.08.03.235150>.
39. Michailidou, K., Lindström, S., Dennis, J., Beesley, J., Hui, S., Kar, S., Lemaçon, A., Soucy, P., Glubb, D., Rostamianfar, A., et al.; NBCS Collaborators; ABCTB Investigators; and ConFab/AOCS Investigators (2017). Association analysis identifies 65 new breast cancer risk loci. *Nature* 551, 92–94.
40. Schumacher, F.R., Al Olama, A.A., Berndt, S.I., Benlloch, S., Ahmed, M., Saunders, E.J., Dadaev, T., Leongamornlert, D., Anokian, E., Cieza-Borrella, C., et al.; Profile Study; Australian Prostate Cancer BioResource (APCB); IMPACT Study; Canary PASS Investigators; Breast and Prostate Cancer Cohort Consortium (BPC3); PRACTICAL (Prostate Cancer Association Group to Investigate Cancer-Associated Alterations in the Genome) Consortium; Cancer of the Prostate in Sweden (CAPS); Prostate Cancer Genome-wide Association Study of Uncommon Susceptibility Loci (PEGASUS); and Genetic Associations and Mechanisms in Oncology (GAME-ON)/Elucidating Loci Involved in Prostate Cancer Susceptibility (ELLIPSE) Consortium (2018). Association analyses of more than 140,000 men identify 63 new prostate cancer susceptibility loci. *Nat. Genet.* 50, 928–936.
41. Nikpay, M., Goel, A., Won, H.-H., Hall, L.M., Willenborg, C., Kanoni, S., Saleheen, D., Kyriakou, T., Nelson, C.P., Hopewell, J.C., et al. (2015). A comprehensive 1,000 Genomes-based genome-wide association meta-analysis of coronary artery disease. *Nat. Genet.* 47, 1121–1130.
42. Censin, J.C., Nowak, C., Cooper, N., Bergsten, P., Todd, J.A., and Fall, T. (2017). Childhood adiposity and risk of type 1 diabetes: A Mendelian randomization study. *PLoS Med.* 14, e1002362.
43. Behar, D.M., Metspalu, M., Baran, Y., Kopelman, N.M., Yunusbayev, B., Gladstein, A., Tzur, S., Sahakyan, H., Bahmanimehr, A., Yepiskoposyan, L., et al. (2013). No evidence from genome-wide data of a Khazar origin for the Ashkenazi Jews. *Hum. Biol.* 85, 859–900.
44. Zou, H. (2006). The adaptive lasso and its oracle properties. *J. Am. Stat. Assoc.* 101, 1418–1429.
45. Zhong, H., and Prentice, R.L. (2008). Bias-reduced estimators and confidence intervals for odds ratios in genome-wide association studies. *Biostatistics* 9, 621–634.
46. Shi, J., Park, J.-H., Duan, J., Berndt, S.T., Moy, W., Yu, K., Song, L., Wheeler, W., Hua, X., Silverman, D., et al.; MGS (Molecular Genetics of Schizophrenia) GWAS Consortium; GECCO (The Genetics and Epidemiology of Colorectal Cancer Consortium); GAME-ON/TRICL (Transdisciplinary Research in Cancer of the Lung) GWAS Consortium; PRACTICAL (Prostate cancer Association group To Investigate Cancer Associated aAlterations) Consortium; PanScan Consortium; and

- GAME-ON/ELLIPSE Consortium (2016). Winner's curse correction and variable thresholding improve performance of polygenic risk modeling based on genome-wide association study summary-level data. *PLoS Genet.* *12*, e1006493.
47. Naseri, A., Tang, K., Geng, X., Shi, J., Zhang, J., Shakya, P., Liu, X., Zhang, S., and Zhi, D. (2021). Personalized genealogical history of UK individuals inferred from biobank-scale IBD segments. *BMC Biol.* *19*, 32.
  48. Wojcik, G.L., Graff, M., Nishimura, K.K., Tao, R., Haessler, J., Gignoux, C.R., Highland, H.M., Patel, Y.M., Sorokin, E.P., Avery, C.L., et al. (2019). Genetic analyses of diverse populations improves discovery for complex traits. *Nature* *570*, 514–518.
  49. Ruan, Y., Lin, Y.-F., Feng, Y.-C.A., Chen, C.-Y., Lam, M., Guo, Z., Stanley Global Asia Initiatives, He, L., Sawa, A., Martin, A.R., Qin, S., Huang, H., and Ge, T. (2021). Improving polygenic prediction in ancestrally diverse populations. medRxiv. <https://doi.org/10.1101/2020.12.27.20248738>.
  50. Moreno-Estrada, A., Gravel, S., Zakharia, F., McCauley, J.L., Byrnes, J.K., Gignoux, C.R., Ortiz-Tello, P.A., Martínez, R.J., Hedges, D.J., Morris, R.W., et al. (2013). Reconstructing the population genetic history of the Caribbean. *PLoS Genet.* *9*, e1003925.
  51. Márquez-Luna, C., Loh, P.-R., Price, A.L.; South Asian Type 2 Diabetes (SAT2D) Consortium; and SIGMA Type 2 Diabetes Consortium (2017). Multiethnic polygenic risk scores improve risk prediction in diverse populations. *Genet. Epidemiol.* *41*, 811–823.
  52. Martin, A.R., Kanai, M., Kamatani, Y., Okada, Y., Neale, B.M., and Daly, M.J. (2019). Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat. Genet.* *51*, 584–591.
  53. Duncan, L., Shen, H., Gelaye, B., Meijssen, J., Ressler, K., Feldman, M., Peterson, R., and Domingue, B. (2019). Analysis of polygenic risk score usage and performance in diverse human populations. *Nat. Commun.* *10*, 3328.
  54. Zhang, Q., Privé, F., Vilhjálmsson, B., and Speed, D. (2021). Improved genetic prediction of complex traits from individual-level data or summary statistics. *Nat. Commun.* *12*, 4192.
  55. Albiñana, C., Grove, J., McGrath, J.J., Agerbo, E., Wray, N.R., Bulik, C.M., Nordentoft, M., Hougaard, D.M., Werge, T., Børglum, A.D., et al. (2021). Leveraging both individual-level genetic data and GWAS summary statistics increases polygenic prediction. *Am. J. Hum. Genet.* *108*, 1001–1011.
  56. Fritsche, L.G., Patil, S., Beesley, L.J., VandeHaar, P., Salvatore, M., Ma, Y., Peng, R.B., Taliun, D., Zhou, X., and Mukherjee, B. (2020). Cancer PRSweb: An online repository with polygenic risk scores for major cancer traits and their evaluation in two independent biobanks. *Am. J. Hum. Genet.* *107*, 815–836.
  57. Lambert, S.A., Gil, L., Jupp, S., Ritchie, S.C., Xu, Y., Buniello, A., McMahon, A., Abraham, G., Chapman, M., Parkinson, H., et al. (2021). The Polygenic Score Catalog as an open database for reproducibility and systematic evaluation. *Nat. Genet.* *53*, 420–425.
  58. Bengtsson, H. (2021). A Unifying Framework for Parallel and Distributed Processing in R using Futures. arXiv. arXiv:2008.00553.
  59. Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L.D., François, R., Golemund, G., Hayes, A., Henry, L., Hester, J., et al. (2019). Welcome to the tidyverse. *J. Open Source Software* *4*, 1686.