

1 Polygenic scoring accuracy varies across the genetic ancestry 2 continuum in all human populations

3 Yi Ding¹, Kangcheng Hou¹, Ziqi Xu², Aditya Pimplaskar¹, Ella Petter², Kristin Boulier¹, Florian
4 Privé³, Bjarni J. Vilhjálmsson^{3,4,5}, Loes Olde Loohuis^{6,7}, Bogdan Pasaniuc^{1,8,9,10}

- 5
- 6 1. Bioinformatics Interdepartmental Program, UCLA, Los Angeles, CA, USA
- 7 2. Department of Computer Science, UCLA, Los Angeles, CA, USA
- 8 3. National Center for Register-Based Research, Aarhus University, Aarhus, Denmark
- 9 4. Bioinformatics Research Centre, Aarhus University, Aarhus, Denmark
- 10 5. Novo Nordisk Foundation Center for Genomic Mechanisms of Disease, Broad Institute, Cambridge, MA, USA
- 11 6. Department of Psychiatry, David Geffen School of Medicine at UCLA, Los Angeles, CA, USA
- 12 7. Program in Neurobehavioral Genetics, David Geffen School of Medicine at UCLA, Los Angeles, CA, USA
- 13 8. Department of Computational Medicine, David Geffen School of Medicine at UCLA, Los Angeles, CA, USA
- 14 9. Department of Human Genetics, David Geffen School of Medicine at UCLA, Los Angeles, CA, USA
- 15 10. Department of Pathology and Laboratory Medicine, David Geffen School of Medicine at UCLA, Los Angeles,
16 CA, USA

17
18 Correspondence: Y.D. (yiding920@ucla.edu); B.P. (pasaniuc@ucla.edu).

19 Abstract

20 Polygenic scores (PGS) have limited portability across different groupings of individuals (e.g., by genetic
21 ancestries and/or social determinants of health), preventing their equitable use. PGS portability has typically
22 been assessed using a single aggregate population-level statistic (e.g., R^2), ignoring inter-individual
23 variation within the population. Here we evaluate PGS accuracy at individual-level resolution, independent
24 of its annotated genetic ancestries. We show that PGS accuracy varies between individuals across the
25 genetic ancestry continuum in all ancestries, even within traditionally “homogeneous” genetic ancestry
26 clusters. Using a large and diverse Los Angeles biobank (ATLAS, $N=36,778$) along with the UK Biobank
27 (UKBB, $N=487,409$), we show that PGS accuracy decreases along a continuum of genetic ancestries in all
28 considered populations and the trend is well-captured by a continuous measure of genetic distance (GD)
29 from the PGS training data; Pearson correlation of -0.95 between GD and PGS accuracy averaged across
30 84 traits. When applying PGS models trained in UKBB “white British” individuals to European-ancestry
31 individuals of ATLAS, individuals in the highest GD decile have 14% lower accuracy relative to the lowest
32 decile; notably the lowest GD decile of Hispanic/Latino American ancestry individuals showed similar PGS
33 performance as the highest GD decile of European ancestry ATLAS individuals. GD is significantly
34 correlated with PGS estimates themselves for 82 out of 84 traits, further emphasizing the importance of
35 incorporating the continuum of genetic ancestry in PGS interpretation. Our results highlight the need for
36 moving away from discrete genetic ancestry clusters towards the continuum of genetic ancestries when
37 considering PGS and their applications.

38 Introduction

39 Polygenic scores (PGS)—estimates of an individual’s genetic predisposition for complex traits/diseases (i.e.
40 genetic value)—are a promising application of large-scale genome-wide association studies (GWAS) to
41 personalized genomic medicine¹⁻⁴, disease risk prediction and prevention⁵⁻⁸. The portability of PGS across

42 different ancestry and socio-demographic groups is limited due to Euro-centric sampling of GWAS data
43 coupled with differences in linkage disequilibrium (LD), minor allele frequency (MAF) and/or disease
44 genetic architecture^{3,9-13}, which poses a critical equity barrier that has prevented widespread adoption of
45 PGS for personalized medicine. For example, PGS are significantly more accurate for individuals of
46 European ancestries as compared to other genetic ancestries^{10,14}; furthermore, PGS accuracy varies across
47 socio-genomic features (e.g., sex, age and social economic status)¹¹, thus complicating interpretability of
48 PGS across groups with different environmental exposures.

49 PGS accuracy is traditionally assayed using population-level metrics of accuracy (e.g., R^2), thus assuming
50 some level of homogeneity across individuals within the considered population^{2,11,15}. However,
51 homogeneous populations are an idealized concept that only roughly approximate human populations;
52 human diversity exists along a genetic ancestry continuum without clearly defined clusters and with various
53 correlations between genetic and socio-environmental factors¹⁵⁻²⁰. Grouping individuals into genetic
54 ancestry clusters obscures the impact of individual variation on PGS accuracy. This is evident for
55 individuals with recently admixed genomes, where genetic ancestries vary individual-to-individual and
56 locus-to-locus in the genome. For example, a single population-level PGS accuracy estimated across all
57 African Americans greatly overestimates PGS accuracy for African Americans with large proportions of
58 African genetic ancestries²¹; likewise, coronary artery disease PGS performs poorly in Hispanic individuals
59 with high proportions of African ancestry²². The genetic ancestry continuum impacts PGS accuracy even
60 in traditionally-labeled “homogeneous/non-admixed” populations; for example, PGS accuracy decays
61 across a gradient of subcontinental ancestries within Europe as the target cohorts are more genetically
62 dissimilar from the data used to train the PGS^{19,23}. Assessing PGS accuracy using population-level metrics
63 is further complicated by technical issues in assigning individuals to discrete clusters of genetic ancestries.
64 Different algorithms and/or reference panels may assign the same individual to different clusters^{15,23,24} and
65 thus to different PGS accuracy classes. Moreover, many individuals are not assigned to a cluster due to
66 limited reference panels used for genetic ancestry inference^{23,25}, leaving such individuals outside PGS
67 accuracy characterization; this poses equity concerns as it limits PGS applications only to individuals within
68 well-defined clusters of genetic ancestries.

69 Here we leverage methods that characterize PGS performance at the level of a single target individual²⁶ to
70 evaluate the impact of the genetic ancestry continuum on PGS accuracy. We use simulation and real data
71 analysis to show that PGS accuracy decays continuously individual-to-individual across the genetic
72 continuum of ancestry as function of genetic distance (GD) from the PGS training data; GD is defined as a
73 principal component analysis (PCA) projection of the target individual on the training data used to estimate
74 the PGS weights. We leverage a large and diverse Los Angeles biobank at UCLA (ATLAS, N= 36,778)
75 joint with UK Biobank (UKBB, N= 487,409) to investigate the interplay between genetic ancestries and
76 PGS for 84 complex traits and diseases. The accuracy of PGS models trained in UKBB “white British”
77 individuals (N= 371,018) is negatively correlated with GD for all considered traits (average Pearson $R=-$
78 0.95 across 84 traits), demonstrating pervasive individual variation in PGS accuracy. The negative
79 correlation remained significant even when restricted to traditionally defined “homogenous” clusters of
80 genetic ancestries (ranging from $R=-0.43$ for East Asian cluster to $R=-0.85$ for the African American cluster
81 in ATLAS). On average across the 84 traits, when rank-ordering individuals according to distance from
82 training data, PGS accuracy decreased by 14% in the closest vs furthest decile in the European genetic
83 cluster; notably the furthest decile of European ancestry individuals showed similar accuracy to the closest
84 decile of Hispanic Latino individuals. Characterizing PGS accuracy across continuum of GD allows for

85 inclusion of individuals unassigned to a given genetic ancestry clusters (6% of all ATLAS), thus allowing
86 for more individuals to be included in PGS applications. Finally, we explore the relationship between GD
87 and PGS estimates themselves; 82 (out of 84) PGS show significant correlation between GD and PGS with
88 30 showing opposite correlation (GD, trait) vs (GD, PGS); we exemplify the importance of incorporating
89 GD in interpretation of PGS using height and neutrophils in the ATLAS data. Overall, our results
90 demonstrate the need to incorporate the genetic ancestry continuum on PGS performance and/or bias.

91 **Results**

92 **Overview of the study**

93 PGS accuracy has traditionally been assessed at the level of discrete genetic ancestry clusters using
94 population-level metrics of accuracy (e.g., R^2). Individuals from diverse genetic backgrounds are routinely
95 grouped into discrete genetic ancestry clusters using computational inference methods such as PCA²⁷ and/or
96 admixture analysis²⁸ (Figure 1a). Population-level metrics of PGS accuracy are then estimated for each
97 genetic ancestry cluster and generalized to each individual in the cluster (Figure 1b). This approach has
98 three major limitations: (1) the inter-individual variability within each cluster is ignored; (2) the genetic
99 ancestry cluster boundary is sensitive to algorithms and reference panels used for clustering; and (3) a
100 significant proportion of individuals may not be assigned to any cluster due to a lack of reference panels
101 for genetic ancestry inference (e.g., individuals of uncommon or admixed ancestries).

102 In this work, we evaluate PGS accuracy across the genetic ancestry continuum at level of a single target
103 individual. We model the phenotype of individual i as $y_i = x_i^T \beta + \epsilon_i$, where x_i is a $M \times 1$ genotype vector
104 indicating allele counts, β is a $M \times 1$ allelic causal effects vector and ϵ_i is random noise. Under a random
105 effects model $gv_i = x_i^T \beta$ and $\widehat{PGS}_i = E(x_i^T \beta | D)$ are random variables where the randomness comes from
106 β and training data D ($D = (X_{train}, y_{train})$). We define the individual PGS accuracy as the correlation of
107 an individual's genetic value and PGS estimates as:

$$108 \quad r_i^2(gv_i, \widehat{PGS}_i) = \frac{cov_{\beta, D}(gv_i, \widehat{PGS}_i)}{var_{\beta}(gv_i)var_{\beta, D}(\widehat{PGS}_i)} = 1 - \frac{E_D \left(var_{\beta | D}(x_i^T \beta) \right)}{var_{\beta}(x_i^T \beta)} \quad (\text{equation 1})$$

109 We use Ldpred2 to estimate $E_D \left(var_{\beta | D}(x_i^T \beta) \right)$ ^{26,29} and approximate $var_{\beta}(x_i^T \beta)$ as the heritability of the
110 phenotype (Methods)³⁰; equation 1 can be further simplified assuming all variants are causal drawn from a
111 normal distribution (infinitesimal model, see Methods). As continuous genetic distance (GD) we use $d_i =$
112 $\sqrt{\sum_{j=1}^J (x_i^T v_j)^2}$ where v_j is the j_{th} eigenvector of training genotype data (Figure 1c). Individuals that are
113 clustered into the same genetic ancestry clusters may have different genetic distance from training data and
114 different individual PGS accuracy (Figure 1d). We use theory and empirical data analyses to show that PGS
115 accuracy decay is well-approximated by the continuous metric of genetic distance.

116 We organize the manuscript as follows. First, we show the relation between genetic distance and PGS
117 accuracy in simulations using real genotype data from UK biobank. Next, we show that existing PGS have
118 accuracy that decreases individual-to-individual as function of genetic distance in a diverse biobank from
119 UCLA. Finally we showcase the impact of genetic distance on interpretability of PGS using height and
120 neutrophil count as example traits.

121 Individual PGS performance is calibrated across the genetic ancestry continuum in 122 simulations

123 First, we evaluated calibration of $E_D \left(\text{var}_{\beta|D} (x_i^T \beta) \right)$ estimated by LDpred2 for individuals at various
124 genetic distances from the UKBB “white British” individuals used to train PGS by checking the calibration
125 the of 90% credible intervals (Figure 2a). We simulated 100 phenotypes at heritability $h_g^2 = 0.25$ and
126 proportion of causal variants $p_{causal} = 1\%$ for all individuals in UK Biobank, assuming shared causal
127 variants and homogenous allelic effect sizes for individuals from various genetic backgrounds (see
128 Methods). Overall, the 90% credible intervals are approximately well-calibrated, i.e. the credible interval
129 overlaps with the true genetic value across 90 out of 100 replicates, for all individuals, regardless of their
130 distance from the training population or genetic ancestry labels (Figure 2a). For example, when individuals
131 are binned into 10 deciles based on their GD from the training population, the average empirical coverage
132 of the 90% credible intervals is 89.7% (s.d. 2.6%) for individuals from the lowest decile (composed of 96.9%
133 individuals labeled as “white British”, 3.1% labeled as “Poland” under discrete view of ancestries)
134 compared with the average empirical coverage of 82.4 % (s.d. 4.6%) for individuals from the highest decile
135 (composed of 19.9% individuals labeled as “Caribbean” and 80.1% labeled as “Nigeria”).

136 Next, we investigated the impact of GD on individual-level PGS accuracy. As expected, the credible
137 interval width increases linearly with GD reflecting reduced predictive accuracy for the PGS (Figure 2b).
138 The average width of 90% credible interval is 1.83 in the highest decile of GD, a 1.8-fold increase over the
139 average width in the lowest decile of GD. In contrast to the credible interval width, the individual-level
140 PGS accuracy \widehat{r}_i^2 decreases with genetic distance from training data (Figure 2c); the average estimated
141 accuracy of individuals in the lowest decile GD is 4-fold higher than that of individuals in the highest decile.
142 Even among the most homogenous grouping of individuals traditionally labeled as white British, we
143 observe a 5% relative decrease in accuracy for individuals at the highest decile of GD as compared to those
144 in the lowest decile. Similar results are observed when using population-level PGS metric of accuracy,
145 albeit at expense of binning individuals according to GD; we find a high degree of concordance between
146 the average \widehat{r}_i^2 within the bin and the population-level R^2 estimated within the bin (Figure 2d,
147 Supplementary Figure 1a). Similarly, we observe a high consistency between average \widehat{r}_i^2 and squared
148 correlation between PGS and simulated phenotype, ($R = 0.87$, $P < 2.2e-16$, Supplementary Figure 1b).
149 Taken together, our results show that 90% credible interval remains calibrated for individuals that are
150 genetically distant from the training population at the expense of a wider credible intervals while
151 \widehat{r}_i^2 captures the PGS accuracy decay across genetic distance.

152 Individual PGS accuracy varies across the genetic ancestry continuum in all ancestries

153 After having validated our approach in simulations, we next turn to empirical data. For illustration purposes
154 we use height as example focusing on ATLAS biobank as target population with PGS trained in the 371,018
155 “white British” individuals from UKBB (Methods); other traits show similar trends and are presented in
156 the next sections. PGS accuracy at the individual level varies with GD across the entire biobank as well as
157 within all genetically inferred ancestry clusters (Figure 3, Supplementary Fig 2). For example, GD strongly
158 correlates with PGS accuracy of individuals in the genetic ancestry cluster labeled as Hispanic/Latino
159 American (HL, $R = -0.83$) and African American (AA, $R = -0.88$) in ATLAS. Notably, GD correlates with
160 PGS accuracy even in non-admixed genetic clusters of ancestry with correlations as -0.66, -0.66 and -0.35,
161 for European Americans (EA), South Asian Americans (SAA) or East Asian Americans (EAA),

162 respectively. Similar qualitative results are also observed when applying PGS in a test data from UKBB
163 with significant negative correlations between GD and individual PGS accuracy in all the sub-continental
164 genetic clusters in UKBB (Supplementary Figure 2) ranging from $R = -0.031$ for the “white British” cluster
165 to $R = -0.62$ for the Caribbean cluster.

166 Next, we focused on the impact of GD on PGS accuracy across all ATLAS individuals regardless of genetic
167 ancestry clustering ($R = -0.96$, $P < 2.2 \times 10^{-16}$, Figure 3b). Notably, we find a strong overlap of PGS accuracies
168 across individuals from different genetical ancestry clusters demonstrating the limitation of using a single
169 cluster-specific metric of accuracy. For example, when rank-ordering by GD, we find the individuals from
170 the closest GD decile in HL cluster have similar estimated accuracy as the individuals from the farthest GD
171 decile in EA cluster (average \widehat{r}_i^2 of 0.71 vs 0.71). This shows that GD enables identification of HL
172 individuals with similar PGS performance as the EA cluster thus partly alleviating inequities due to lack of
173 access to accurate PGS. Most notably, GD can be used to evaluate PGS performance for individuals that
174 cannot be easily clustered by current genetic inference methods (6% of all individuals in ATLAS, Figure
175 3b) partly due to limitations of reference panels and algorithms for assigning ancestries. Among this
176 traditionally overlooked group of individuals, we find the GD ranging from 0.02 to 0.64 and their
177 corresponding estimate PGS accuracy \widehat{r}_i^2 ranging from 0.63 to 0.21. In addition to evaluating PGS accuracy
178 with respect to the genetic value, we also evaluated accuracy with respect to the residual height after
179 regressing out sex, age, PC1-10 on the ATLAS from the actual measured trait. Using equally spaced bins
180 across the GD continuum, we find that correlation between PGS and the measured height tracks
181 significantly with GD ($R = -0.9$, $P\text{-value} = 5.9 \times 10^{-8}$, Figure 3c).

182 **The continuous decay of PGS accuracy across genetic distance is pervasive across all traits**

183 Having established the coupling of GD with PGS accuracy in simulations and for height, we next turn to
184 the question of whether such relationship is pervasive across complex traits using PGS for a broad set of
185 84 traits (Supplementary Table 1). We find consistent and pervasive correlations of GD with PGS accuracy
186 across all considered traits in both ATLAS and UK Biobank (Figure 4). For example, the correlations
187 between GD and individual PGS accuracy range from -0.71 to -0.97 with an average of -0.95 across the 84
188 PGS in ATLAS with similar results in UKBB. Traits with sparser genetic architectures and fewer non-zero
189 weights in the PGS yield to a lower correlation between GD and PGS accuracy; we hypothesize this is
190 because GD represents genome-wide genetic variation patterns that may not reflect a limited number of
191 causal SNPs well. For example, PGS for Lipoprotein A (`log_lipoA`) has the lowest polygenicity estimate
192 (0.02%) among the 84 traits and has the lowest correlation in ATLAS (-0.71) and UKBB (-0.85). In contrast,
193 we observe a high correlation between GD and PGS accuracy (>0.9) for all traits with an estimated
194 polygenicity $> 0.1\%$. Next, we show that the fine-scale population structures accountable for the individual
195 PGS accuracy variation is also prevalent within the traditionally defined genetic ancestry group. For
196 example, in ATLAS we find 501 out of 504 (84 traits across 6 genetic ancestry clusters) trait-ancestry pairs
197 have a significant association between GD and individual PGS accuracy after Bonferroni correction. In
198 UKBB, we find 572 out of the 756 (84 traits across 9 subcontinental genetic ancestry clusters) trait-ancestry
199 pairs have significant association between genetic distance and PGS accuracy after Bonferroni correction.
200 We also find that a more stringent definition of homogenous genetic clusters results in a lower correlation
201 magnitude (Supplementary Figure 3).

202 Genetic distance correlates with PGS estimates across most traits

203 We focused so far on investigating the relationship between GD and PGS accuracy. Next, we turn to
204 evaluating the impact of GD on PGS estimates themselves. We find that GD is significantly correlated with
205 PGS estimates for 82 out of 84 traits in UKBB ranging from $R=-0.52$ to $R=0.74$ (Supplementary Figure 4);
206 this broad range of correlations is in stark contrast with the highly consistent negative correlation of GD
207 and PGS r_i^2 . To gain insights into whether PGS coupling with GD is due to stratification or true signal, we
208 next contrasted the correlation of GD to PGS estimates ($cor(d_i, \widehat{PGS}_i)$) with correlation of GD to the
209 measured phenotype values ($cor(d_i, y_i)$). We find a wide-range of couplings reflecting trait-specific signals;
210 30 traits GD correlate in opposite directions with PGS vs. phenotype; 40 trait GD correlates in the same
211 direction with PGS vs. phenotype but differ in correlation magnitude (Supplementary Figure 4). For
212 example, GD shows opposite and significantly different correlations PGS vs trait for years of education
213 (years_of_edu, $cor(y_i, d_i) = 0.03$, $cor(\widehat{PGS}_i, d_i) = -0.18$). Other traits such as hair color show highly
214 consistent impact of GD on PGS vs trait (darker_hair, $cor(y_i, d_i) = 0.59$, $cor(\widehat{PGS}_i, d_i) = 0.74$); while for
215 monocyte percentage GD shows different magnitudes albeit with the same directions (monocyte_perc,
216 $cor(y_i, d_i) = -0.03$, $cor(\widehat{PGS}_i, d_i) = -0.52$). The correlation between GD and phenotype/PGS is also
217 observed in ATLAS. For example, both height phenotype and PGS for height decrease along GD within in
218 ATLAS (Figure 5a); this holds true even if restricted to the European American genetic ancestry cluster
219 (Figure 5b). This is consistent with genetic value driving difference in phenotypes but could also be
220 explained by residual stratification. For neutrophil counts, phenotype and PGS varies in opposite direction
221 along GD across the ATLAS (Figure 5c), although the trend is similar for phenotype and PGS in European
222 American cluster (Figure 5d). This could be explained by genetic value driving signal in Europeans with
223 stratification for other groups. Neutrophil counts have been reported to vary greatly across ancestry groups
224 with reduced counts in individuals of African ancestries³¹. In ATLAS, we observe a negative correlation (-
225 0.04) between GD and neutrophil counts in agreement with the previous reports, while GD is positively
226 correlated (0.08) with PGS estimates with genetically distant individuals traditionally labeled as African
227 American having higher PGS than average. The opposite directions in phenotype/PGS-distance correlations
228 are partly attributed to Duffy-null SNP rs2814778 on chromosome 1q23.2. This variant has a large
229 association with neutrophil counts among individuals traditionally identified as African ancestry, but it is
230 rare and excluded in our training data. This exemplifies the potential bias in PGS due to non-shared causal
231 variants and urges ancestral diversity in genetic studies.

232 Since PGS can vary across GD either as reflection of true signal (i.e. genetic value varying with ancestry)
233 or due to biases in PGS estimation ranging from unaccounted residual population stratification to
234 incomplete data (e.g., partial ancestry-specific tagging of causal effects), our results emphasize the need to
235 consider GD in PGS interpretation beyond adjusting for PGS r_i^2 .

236 Discussion

237 In this work, we showed that continuous genetic ancestry impacts PGS accuracy and its interpretability
238 across a continuum of genetic ancestries. We proposed individual PGS accuracy as an approach to
239 individualize PGS performance to each target individual. We use a PCA-based genetic distance²³ from the
240 center of training data to describe an individual's unique location on the genetic ancestry continuum. In
241 simulations and real data analyses, we showed that individual PGS accuracy tracks well with genetic
242 distance. We demonstrate the pervasive continuous decay of PGS performance as the target individual is
243 further away from the training population. Our demonstration of the continuous PGS accuracy decay

244 directly leads to two conclusions: first, PGS accuracy decay already happens within traditionally defined
245 genetic ancestry groups; second, PGS accuracy can be similar for genetically adjacent individuals that are
246 separated into distinct genetic ancestry groups (usually with different population PGS accuracy). Individual
247 PGS accuracy also enables the evaluation of PGS performance for individuals who cannot be clustered into
248 reference populations, obviating the necessity for genetic ancestry clustering and PGS accuracy evaluation
249 under a discrete view of genetic ancestry.

250 Our results have several implications for applying PGS to populations with diverse genetic ancestries. First,
251 we highlight the variability in PGS performance along the continuum of genetic ancestry, even within
252 traditionally defined homogenous populations. With the increasing recognition that genetic ancestries are
253 not discrete but rather continuous¹⁵⁻¹⁹, the individual-level accuracy introduced here provides a powerful
254 tool to study PGS performance along the genetic ancestry continuum. Given the pervasive variable PGS
255 accuracy across individuals, incorporating individual-level metrics of PGS performance can improve the
256 utility of PGS. For example, by using individual-level PGS accuracy, we can identify individuals from
257 Hispanic/Latino genetic ancestry cluster who have similar PGS accuracy with European individuals thus
258 partly alleviating inequities due to lack of access to accurate PGS.

259 Second, our simulation and real data analysis show that the individual PGS accuracy is highly correlated
260 with genetic distance from training data. The increased genetic distance corresponds to a lower relatedness
261 of the testing individual with the training population¹³. This provides a finer resolution compared with
262 previous theoretical studies that investigate population level PGS portability based on LD/MAF difference,
263 F_{st} and mean kinship between training and targeting population^{12,13}. Along with previous studies, our results
264 emphasize the importance of powerful PGS training in non-European cohorts³² to improve the PGS
265 performance for individuals from diverse genetic background. To narrow the prediction gap and ameliorate
266 consequent health disparities between European and non-European ancestries individuals, concerted global
267 effort and equitable collaborations are needed to increase the sample size of underrepresented
268 individuals^{32,33}. Equally important is the development of multi-ancestry PGS methods that can effectively
269 leverage ancestrally diverse populations to train PGS models³⁴. Some examples of such methods recently
270 developed include PRS-CSx³⁵, vilma³⁶ and CT-SLEB³⁷.

271 Third, our results highlight the pervasive correlation between PGS estimates and genetic distance from the
272 training data, which usually displays opposite direction or different magnitude compared with the
273 correlation between phenotype and genetic distance. This observation provides a finer resolution of the
274 previously reported mean shift of PGS estimates across genetic ancestry groups⁹. We note that the
275 correlation between genetic distance and phenotype can stem from both potential bias in PGS estimates and
276 true biological difference such as continuous genetic variation. We provide neutrophil counts PGS as an
277 example of potential bias due to low allele frequency of Duffy-null SNP rs2814778³¹ in the training data,
278 however we cannot rule out the impact of true biological differences for most traits. More effort is needed
279 to investigate the PGS bias especially in the context of continuous genetic ancestry.

280 We note several limitations and future directions of our work. First, individual PGS accuracy is derived
281 from individual PGS uncertainty with approximations under strong assumptions that the causal variants and
282 effects are the same across all genetic ancestries. In reality, despite the abundance of shared causal variants³⁸
283 and the strong transethnic genetic effect sizes correlation³⁹, population-specific causal variants and effects
284 still exist and limit the transferability of PGS. Future work could investigate the impact of the population-
285 specific components of genetic architecture on the calibration of PGS accuracy. Second, we approximate

286 the variance of genetic value with heritability and set the value fixed for all individuals. Further work can
287 be done to quantify the genetic value variance for individuals at different genetic distance and assess its
288 impact on accuracy. Third, individual PGS accuracy evaluates how well PGS estimates the genetic value
289 instead of how accurate the PGS predicts the phenotype. Quantifying the individual accuracy of PGS with
290 respect to phenotype can be achieved by modeling environments to calibrate over phenotypes. Fourth, there
291 can be misspecification of model assumptions for the individual-level PGS uncertainty. Future work can be
292 done to investigate the impact of the genetic architectures on the calibration of PGS uncertainty/accuracy.
293 Fifth, limited by the sample size, we train PGS on white British individuals in UKBB, and inevitably define
294 genetic distance relative to European individuals. This work should be replicated while training PGS in
295 non-European individuals in future works. Alternative definitions of genetic distance such as genetic
296 relatedness¹³ and other multi-dimensional descriptions of genetic ancestry continuum¹⁶ can also be explored
297 in the future.

298

299 Methods

300 **Model.** We model the phenotype of an individual with a standard linear model $y_i = x_i^T \beta + \epsilon_i$ where x_i is
 301 an $M \times 1$ genotype vector indicating allele counts β is an $M \times 1$ vector of allelic genetic effects and ϵ_i is
 302 random noise. Under a random effects model, β is a vector of random variable sampled from a prior
 303 distribution $p(\beta)$ which differs under different genetic architecture assumptions⁴⁰ and PGS methods^{29,41-43}.
 304 The PGS weights $\hat{\beta} = E_{\beta|D}(\beta)$ are estimated to be the posterior mean given the observed data D ($D =$
 305 (X_{train}, y_{train})) with access to individual-level genotypes X_{train} and phenotype y_{train} or $D = (\hat{\beta}_{GWAS}, \hat{R})$
 306 with access to marginal association statistics $\hat{\beta}_{GWAS}$ and \hat{R} . An individual i 's genetic value ($gv_i = x_i^T \beta$)
 307 is estimated to be $\widehat{PGS}_i = E_{\beta|D}(x_i^T \beta)$, the uncertainty of which is estimated as the posterior variance of
 308 genetic value $var(\widehat{PGS}_i) = var_{\beta|D}(x_i^T \beta)$ ²⁶.

309 **Individual PGS accuracy.** Under a random effects model both the genetic value and PGS estimate for
 310 individual i are random variables. The randomness of $gv_i = x_i^T \beta$ comes from the randomness in β and the
 311 randomness of $\widehat{PGS}_i = x_i^T \hat{\beta}$ comes from the randomness of both β and the training data D . Individual PGS
 312 accuracy measures the correlation between gv_i and \widehat{PGS}_i , which can be computed with the following
 313 equation:

$$314 \quad r_i^2 = 1 - \frac{E_D \left(var_{\beta|D}(x_i^T \beta) \right)}{var_{\beta}(x_i^T \beta)} \quad (\text{equation 1})$$

315 where $var_{\beta|D}(x_i^T \beta)$ is the posterior variance of genetic value given the training data and $var_{\beta}(x_i^T \beta)$ is the
 316 genetic variance. The equation is derived as follows:

317 First, we show that under the random effects model, $cov_{\beta,D}(x_i^T \hat{\beta}, x_i^T \beta) = var_D(x_i^T \hat{\beta})$ (where $\hat{\beta} =$
 318 $E_{\beta|D}(\beta)$) following equation 5.149 in ref⁴⁴:

$$\begin{aligned} 319 \quad cov_{\beta,D}(\hat{\beta}, \beta^T) &= E_{\beta,D}(\hat{\beta} \beta^T) - E_{\beta,D}(\hat{\beta}) E_{\beta,D}(\beta^T) \\ 320 \quad &= E_D \left(E_{\beta|D}(\hat{\beta} \beta^T) \right) - E_{D,\beta}(\hat{\beta}) E_D \left(E_{\beta|D}(\beta^T) \right) \\ 321 \quad &= E_D \left(E_{\beta|D}(E_{\beta|D}(\beta) \beta^T) \right) - E_D \left(E_{\beta|D}(\beta) \right) E_D \left(E_{\beta|D}(\beta^T) \right) \\ 322 \quad &= E_D \left(E_{\beta|D}(\beta) E_{\beta|D}(\beta^T) \right) - E_D \left(E_{\beta|D}(\beta) \right) E_D \left(E_{\beta|D}(\beta^T) \right) \\ 323 \quad &= var_D \left(E_{\beta|D}(\beta) \right) \\ 324 \quad &= var_D(\hat{\beta}) \end{aligned}$$

325 Multiply x_i on both sides of equation, we obtain:

$$\begin{aligned} 326 \quad x_i^T cov_{\beta,D}(\hat{\beta}, \beta) x_i &= x_i^T var_D(\hat{\beta}) x_i \\ 327 \quad cov_{\beta,D}(x_i^T \hat{\beta}, x_i^T \beta) &= var_D(x_i^T \hat{\beta}) \quad (\text{equation 2}) \end{aligned}$$

328 Next, by applying the law of total variance, we show that:

$$329 \quad \text{var}_{\beta,D}(gv_i) = \text{var}_{\beta,D}(x_i^T \beta) = E_D \left(\text{var}_{\beta|D}(x_i^T \beta) \right) + \text{var}_D \left(E_{\beta|D}(x_i^T \beta) \right)$$

$$330 \quad \text{var}_D(x_i^T \hat{\beta}) = \text{var}_{\beta,D}(gv_i) - E_D \left(\text{var}_{\beta|D}(x_i^T \beta) \right) \text{ (equation 3)}$$

331 Third, we derive the correlation between gv_i and \widehat{PRS}_i as:

$$332 \quad r_i^2 = \frac{\text{cov}_{\beta,D}(gv_i, \widehat{gv}_i)}{\text{var}_{\beta,D}(gv_i)\text{var}_{\beta,D}(\widehat{gv}_i)}$$

$$333 \quad = \frac{\text{var}_D(x_i^T \hat{\beta})^2}{\text{var}_D(x_i^T \beta)\text{var}_D(x_i^T \hat{\beta})} \text{ by applying equation 2}$$

$$334 \quad = \frac{\text{var}_D(x_i^T \hat{\beta})}{\text{var}_D(x_i^T \beta)}$$

$$335 \quad = \frac{\text{var}_{\beta,D}(x_i^T \beta) - E_D \left(\text{var}_{\beta|D}(x_i^T \beta) \right)}{\text{var}_D(x_i^T \beta)} \text{ by applying equation 3}$$

$$336 \quad = 1 - \frac{E_D \left(\text{var}_{\beta|D}(x_i^T \beta) \right)}{\text{var}_D(x_i^T \beta)}$$

$$337 \quad$$

$$338 \quad = 1 - \frac{E_D \left(\text{var}_{\beta|D}(x_i^T \beta) \right)}{\text{var}_{\beta}(x_i^T \beta)}$$

339 Equation 1 is widely used in animal breeding theory to compute the reliability of estimated breeding value
 340 for each individual³⁰. In this work, we use individual PGS uncertainty $\text{var}(\widehat{PRS}_i) = \text{var}_{\beta|D}(x_i \beta)$ as an
 341 unbiased estimator of $E_D(\text{var}_{\beta|D}(x_i^T \beta))$. We also use estimated heritability to approximate $\text{var}_{\beta}(x_i^T \beta)$ in
 342 simulation where the phenotype has unit variance. In real data analysis, since the phenotypes does not
 343 necessarily have unit variance, we approximate $\text{var}_{\beta}(x_i^T \beta)$ by scaling the estimated heritability multiplied
 344 by the residual phenotypic variance in the training population after regressing GWAS covariates including
 345 sex, age and precomputed UKBB PC1-16 (Data-Field 22009).

346 **Analytical form of individual PGS accuracy under infinitesimal assumption.** Without loss of generality,
 347 we assume a prior distribution of genetic effects as follows:

$$348 \quad p(\beta | \sigma_g^2) = \text{MVN}(0, \sigma_g^2 I_M)$$

349 With access to individual genotype data X_{train} and phenotypes y_{train} , the likelihood of the data is

$$350 \quad p(y_{train} | X_{train}, \beta, \sigma_e^2) = \text{MVN}(X_{train}\beta, \sigma_e^2 I_N)$$

351 The posterior distribution of genetic effects given the data is proportional to the product of the prior and the
 352 likelihood:

$$353 \quad p(\beta | X_{train}, y_{train}, \sigma_g^2, \sigma_e^2) \propto p(\beta | \sigma_g^2) p(y_{train} | X_{train}, \beta, \sigma_e^2)$$

$$354 \quad \propto \text{MVN}(0, \sigma_g^2 I_M) \text{MVN}(X_{train}\beta, \sigma_e^2 I_N)$$

355 $\propto MVN(\mu_\beta, \Sigma_\beta)$

356 Where $\mu_\beta = \left(\frac{\sigma_e^2}{\sigma_g^2} I_M + X_{train}^T X_{train}\right)^{-1} X_{train}^T y_{train}$ and $\Sigma_\beta = \sigma_e^2 \left(\frac{\sigma_e^2}{\sigma_g^2} I_M + X_{train}^T X_{train}\right)^{-1}$

357 For a new target individual, the posterior variance of the genetic value is:

358 $var(x_i \beta | x_i, X_{train}, y_{train}, \sigma_g^2, \sigma_e^2) = x_i^T \Sigma_\beta x_i = \sigma_e^2 x_i^T \left(\frac{\sigma_e^2}{\sigma_g^2} I_M + X_{train}^T X_{train}\right)^{-1} x_i$

359 After performing eigendecomposition on $X_{train}^T X_{train} = \sum_{j=1}^J \lambda_j v_j v_j^T$, we can rewrite

360 $\left(\frac{\sigma_e^2}{\sigma_g^2} I_M + X_{train}^T X_{train}\right)^{-1} = \left(\frac{\sigma_e^2}{\sigma_g^2} + \sum_{j=1}^J \lambda_j v_j v_j^T\right)^{-1} = \sum_{j=1}^J \left(\frac{\sigma_e^2}{\sigma_g^2} + \lambda_j\right)^{-1} v_j v_j^T$

361 where v_j and λ_j corresponds to the j_{th} eigenvalue and unit-length eigenvector of training genotype X_{train}

362 Thus, we can rewrite the posterior variance of genetic value as

363 $var(x_i \beta | x_i, X_{train}, y_{train}, \sigma_g^2, \sigma_e^2) = \sigma_e^2 \sum_{j=1}^J \left(\frac{\sigma_e^2}{\sigma_g^2} + \lambda_j\right)^{-1} x_i^T v_j v_j^T x_i$

364 Replacing $E_D\left(var_{\beta|D}(x_i^T \beta)\right)$ in equation 1 with analytical form of $var(x_i \beta | x_i, X_{train}, y_{train}, \sigma_g^2, \sigma_e^2)$, we
365 get

366 $r_i^2 = 1 - \frac{var(x_i^T \beta | x_i, X_{train}, y_{train}, \sigma_g^2, \sigma_e^2)}{var(x_i^T \beta)} = 1 - \frac{\sigma_e^2 \sum_{j=1}^J \left(\frac{\sigma_e^2}{\sigma_g^2} + \lambda_j\right)^{-1} x_i^T v_j v_j^T x_i}{\sigma_g^2 x_i^T x_i}$

367

368 As the eigenvalue of $X_{train}^T X_{train}$ increases linearly with training sample size N^{45} . At the UKBB level
369 sample size (e.g. $N = 371,018$ for our UKBB white British training data), the eigenvalue for the top PCs are
370 usually larger than the ratio of environmental noise variance and genetic variance $\frac{\sigma_e^2}{\sigma_g^2}$. Thus, we can further
371 approximate the analytical form with:

372 $r_i^2 = 1 - \frac{\sigma_e^2 \sum_{j=1}^J \frac{1}{\lambda_j} x_i^T v_j v_j^T x_i}{\sigma_g^2 x_i^T x_i} = 1 - \frac{\sigma_e^2}{\sigma_g^2} \frac{\sum_{j=1}^J \frac{1}{\lambda_j} x_i^T v_j v_j^T x_i}{x_i^T x_i}$

373 The term $\sum_{j=1}^J \frac{1}{\lambda_j} x_i^T v_j v_j^T x_i$ is the Mahalanobis distance of the testing individual i from the center of the
374 training genotype data on its PC space and $x_i^T x_i$ is the sum of squared allelic counts across all variants.

375 Empirically, the ratio between the two is highly correlated with the Euclidean distance of the individual
376 from the training data on that PC space ($R=1$, $P\text{-value} < 2.2e-16$ in UKBB).

377 **Genetic Distance.** The genetic distance is defined as the Euclidean distance between a target individual
378 and the center of training data on the PC space of training data.

379
$$d_i = \sqrt{\sum_{j=1}^J (x_i^T v_j)^2}$$

380 where d_i is the genetic distance of a testing individual i from the training data, x_i is an $M \times 1$ genotype
381 vector for testing individual i , v_j is the j_{th} eigenvector for the genotype matrix of training individuals and
382 J is set to 20.

383 **Ancestry ascertainment in UKBB.** The UKBB individuals are clustered into nine sub-continental ancestry
384 clusters White British, Poland, Iran, Italy, Ashkenazi, India, China, Caribbean and Nigeria based on the top
385 16 precomputed PCs (Data-Field 22009) as described in ref²³. The center of each ancestry group on PC
386 space is obtained from ref²³. Each individual is assigned to one of the nine ancestral groups based on their
387 Euclidean distance to the centers on the PC space. The genetic ancestry of an individual is labeled as
388 unknown if its distance to any genetic ancestry center is larger than one eighth of maximum distance
389 between any pairs of sub-continental ancestry clusters. We are able to cluster 91% of the UKBB participants
390 into 411,018 British white, 4127 Polish, 1169 Iran, 6499 Italy, 2352 Ashkenazi, 1798 China, 2472
391 Caribbean and 3894 Nigeria.

392 **Ancestry ascertainment in ATLAS.** The ATLAS individuals are clustered into five genetic ancestry
393 clusters - European Americans (EA), Hispanic and Latino Americans (HL), South Asian Americans (SAA)
394 and East Asian Americans (ESA) and African Americans (AA) as described in ref²⁵ based on their
395 proximity with 1000 Genome super populations on the PC space. First, we filter the ATLAS typed
396 genotypes with plink2 by Mendel error rate ('plink --me 1 1 --set-me-missing'), founders ('--filter-
397 founders'), minor allele frequency ('-maf 0.15'), genotype missing call rate ('--geno 0.05'), and Hardy-
398 Weinberg equilibrium test p-value ('-hwe 0.001'). Next, ATLAS genotypes were merged with the 1000
399 Genomes phase 3 dataset. Then, LD pruning was performed on the merged dataset ('--indep 200 5 1.15 --
400 indep-pairwise 100 5 0.1'). The top10 PCs were computed with the flashpca²⁴⁶ software with all default
401 parameters. Next, we use the super population label and PCs of the 1000 Genome individuals to train the
402 K-nearest neighbors (KNN) model to assign genetic ancestry labels to each ATLAS individual. For each
403 ancestry cluster, we run KNN on the pair of PCs that capture the most variation for each genetic ancestry
404 group: European, East Asian, and African ancestry groups utilize PCs 1 and 2, the Admixed American
405 group use PCs 2 and 3, and the South Asian group use PCs 4 and 5. In each analysis, we use 10-fold cross-
406 validation to select the 'k' hyper-parameter from k=5, 10, 15, 20. If an individual is assigned to multiple
407 ancestries with probability larger than 0.5 or is not assigned to any clusters, it's labeled as unknown. We
408 relabel the five 1000 genome super population as EA for EUR, HL for AMR, SAA for SAS, AA for AFR
409 and ESA for EAS. We can cluster 95% of the ATLAS participants into 22,380 EA, 6973 HL, 625 SAA,
410 3331 EAA, 1995 AA and 2332 individuals are labeled as unknown.

411 **Genotype data.** In simulations, we use 1,054,151 UKBB HapMap3 SNPs for simulating phenotypes,
412 training PGS models and calculating PGS for testing individuals in UKBB. For real data analysis, we use

413 an intersection of UKBB HapMap 3 SNPs and ATLAS imputed SNPs for the training of PGS in UKBB
414 and calculating PGS for remaining UKBB individuals and ATLAS individuals. We start from 1,054,151
415 UKBB HapMap3 SNPs and 8,048,268 ATLAS imputed SNPs. Since UKBB is on genome build hg37 and
416 ATLAS is on hg38, we first lift all ATLAS SNPs from hg38 to hg37 with *snp_modifyBuild* function in
417 *bigsnpr* R package. Next, we match UKBB SNPs and ATLAS SNPs by chromosome and position with
418 *snp_match* function in *bigsnpr*. Then, we recode ATLAS SNPs using UKBB reference alleles with *plink2*
419 `--recode` flag. In the end, 979,457 SNPs remain for training the LDpred2 models in real data analysis.

420 **Simulated phenotypes.** We use simulations on all UKBB individuals to investigate the impact of genetic
421 distance from training data on the various metrics of PGS. We fix the proportion of causal SNPs $p_{causal} =$
422 0.01 and heritability as $h_g^2 = 0.25$. The simulated genetic effects and phenotype are generated as follows:
423 First, we randomly sample

$$424 \quad \beta_m \sim \begin{cases} N\left(0, \frac{h_g^2}{\text{var}(x_m) M p_{causal}}\right) & c_j = 1, \text{ with probability } p_{causal} \\ 0 & c_j = 0, \text{ with probability } 1 - p_{causal} \end{cases}$$

425 where $\text{var}(x_m)$ is the variance of allele counts for SNP m among all UKBB individuals. Second, we
426 compute the genetic value for each individual as $g v_i = \sum_{m=1}^M x_{im} \beta_m$ and randomly sample environmental
427 noise $\epsilon_i \sim N(0, 1 - h_g^2)$. Third, generate phenotype as $y_i = g v_i + \epsilon_i$. We repeat the process 100 times to
428 generate 100 sets of genetic values and phenotypes.

429 **Genetic distance from PGS training data.** To compute the genetic distance of testing individuals from
430 the training population, we perform PCA on the 371,018 UKBB white British training individuals and
431 project the 48,586 UKBB testing individuals and 36,778 ATLAS training individuals on the PC space. We
432 start from the 979,457 SNPs that are overlapped in UKBB and ATLAS. First, we perform LD pruning with
433 *plink2* (`--indep-pairwise 1000 50 0.05`) and exclude the long-range LD regions. Next, we perform PCA
434 analysis with *flashpca*⁴⁶ on the 371,018 UKBB white British training individuals to obtain the top 20 PCs.
435 Then, we project the remaining 48,586 UKBB individuals that are not included in the training data and
436 36,778 ATLAS individuals onto the PC space of training data by using SNP loadings (`--outload loadings.txt`)
437 and their means and standard deviations (`--outmeansd meansd.txt`) output from *flashpca2*. In the end, we
438 compute the genetic distance for each individual as the norm of its projection on the PC space.

439 **LDpred2 PGS model training.** The PGS models were trained on 371,018 UKBB individuals labeled as
440 white British with the LDpred2²⁹ method for both simulation and real data analysis. For simulation analysis,
441 we use 1,054,151 UKBB HapMap3 variants. For real data analysis, we use 979,457 SNPs that are
442 overlapped in UKBB HapMap3 variants and ATLAS imputed genotypes.

443 First, we obtain GWAS summary statistics by performing GWAS on the training individuals with *plink2*
444 using sex, age and precomputed PC-1-16 as covariates. Second, we calculate the in-sample LD matrix with
445 the function *snp_cor* from R package *bigsnpr*⁴⁷. Next, we use the GWAS summary statistics and LD matrix
446 as input for *snp_ldpred2_auto* function in *bigsnpr* to sample from the posterior distribution of genetic effect
447 sizes. Instead of using a held-out validation dataset to select hyperparameters p (proportion of causal
448 variants) and h_2 (heritability), *snp_ldpred2_auto* estimates the two parameters from data with MCMC
449 directly. We run 10 chains with different initial sparsity p from 10^{-4} to 1 equally spaced in log space. For
450 all chains, we set the initial heritability as the LD score regression heritability⁴⁸ estimated by the built-in

451 function *snp_ldsc*. We perform quality control of the 10 chains by filtering out chains with estimated
452 heritability that are smaller than 0.7 times of the median heritability of the 10 chains or with estimated
453 sparsity that are smaller than 0.5 times of the median sparsity or 2 times of the median sparsity. For each
454 chain that passes filtering, we remove the first 100 MCMC iterations as burn-in and thin the next 500
455 iterations by selecting every 5th iteration to reduce autocorrelation between MCMC samples. In the end,
456 we obtain a $M \times B$ matrix $[\tilde{\beta}^{(1)}, \tilde{\beta}^{(2)}, \dots, \tilde{\beta}^{(B)}]$, where each column of the matrix $\tilde{\beta}^{(b)}$ is a sample of
457 posterior causal effects of the M SNPs. Due to the quality control of MCMC chains, the total number of
458 posterior samples B ranges from 500 to 1000.

459 **Calculate PGS and accuracy.** We use the score function in *plink2* to compute the PGS for 48,586 and
460 36,778 testing individuals in UKBB and ATLAS, respectively. For each $\tilde{\beta}^{(b)}$, we compute the PGS for
461 each individual i as $x_i^T \tilde{\beta}^{(b)}$ with *plink2* (--score). For each individual with genotype x_i , we compute
462 $x_i^T \tilde{\beta}^{(1)}, x_i^T \tilde{\beta}^{(2)}, \dots, x_i^T \tilde{\beta}^{(B)}$ to approximate its posterior distribution of genetic value. The genotype x_i^T
463 is centered to the average allele count (--read-freq) in training data to reduce the uncertainty from the
464 unmodeled intercept. We estimate the PGS with the posterior mean of the genetic value as $\widehat{PGS}_i =$
465 $E_{\beta|D}(x_i^T \beta) = \frac{1}{B} \sum_{b=1}^B x_i^T \tilde{\beta}^{(b)}$. We estimate the individual level PGS uncertainty as $var(\widehat{PGS}_i) =$
466 $var_{\beta|D}(x_i^T \beta) = \frac{1}{B} \sum_{b=1}^B (x_i^T \tilde{\beta}^{(b)} - \widehat{PGS}_i)^2$. The individual level PGS accuracy is calculated as $\widehat{r}_i^2 = 1 -$
467 $\frac{var(\widehat{PGS}_i)}{h_g^2}$ for simulation (h_g^2 is the heritability estimated by the LDpred2 model) and as $\widehat{r}_i^2 = 1 -$
468 $\frac{var(\widehat{PGS}_i)}{h_g^2 var(y_{train} - \hat{y}_{train})}$ for real data analysis, where $var(y_{train} - \hat{y}_{train})$ refers to the variance of residual
469 phenotype in training data after regressing out GWAS covariates).

470 **Calibration of credible interval in simulation.** We run the LDpred2 model on 371,018 white British
471 training individuals for the 100 simulation replicates. In each simulation r , for individual with genotype
472 x_i , we compute $x_i^T \tilde{\beta}_r^{(1)}, x_i^T \tilde{\beta}_r^{(2)}, \dots, x_i^T \tilde{\beta}_r^{(B)}$ to approximate its posterior distribution of genetic value,
473 generate 90% credible interval $CI - GV_{ir}$ with 5% and 95% quantile of the distribution and check if its
474 genetic value is contained in the credible interval $I(gv_{ir} \in CI - GV_{ir})$. We compute the empirical coverage
475 for each individual as the mean across the 100 simulation replicates $coverage_i = \frac{1}{100} \sum_{r=1}^{100} I(gv_{ir} \in CI -$
476 $GV_{ir})$.

477

478 **Ethics declarations**

479 All research performed in this study conformed with the principles of the Helsinki Declaration. All
480 individuals provided written informed consent to participate in the study. Patient Recruitment and Sample
481 Collection for Precision Health Activities at UCLA is an approved study by the UCLA Institutional
482 Review Board (UCLA IRB). IRB#17-001013.

483 **Data availability**

484 The individual-level genotype and phenotype data are available by application from the UKBB
485 <http://www.ukbiobank.ac.uk/>. Summary statistics for UCLA ATLAS data are available at [https://atlas-](https://atlas-phewas.mednet.ucla.edu/)
486 [phewas.mednet.ucla.edu/](https://atlas-phewas.mednet.ucla.edu/).

487

488 **URLs**

489 LDPred2 software implementing individual PRS uncertainty:

490 <https://privefl.github.io/bigsnpr/articles/LDpred2.html>

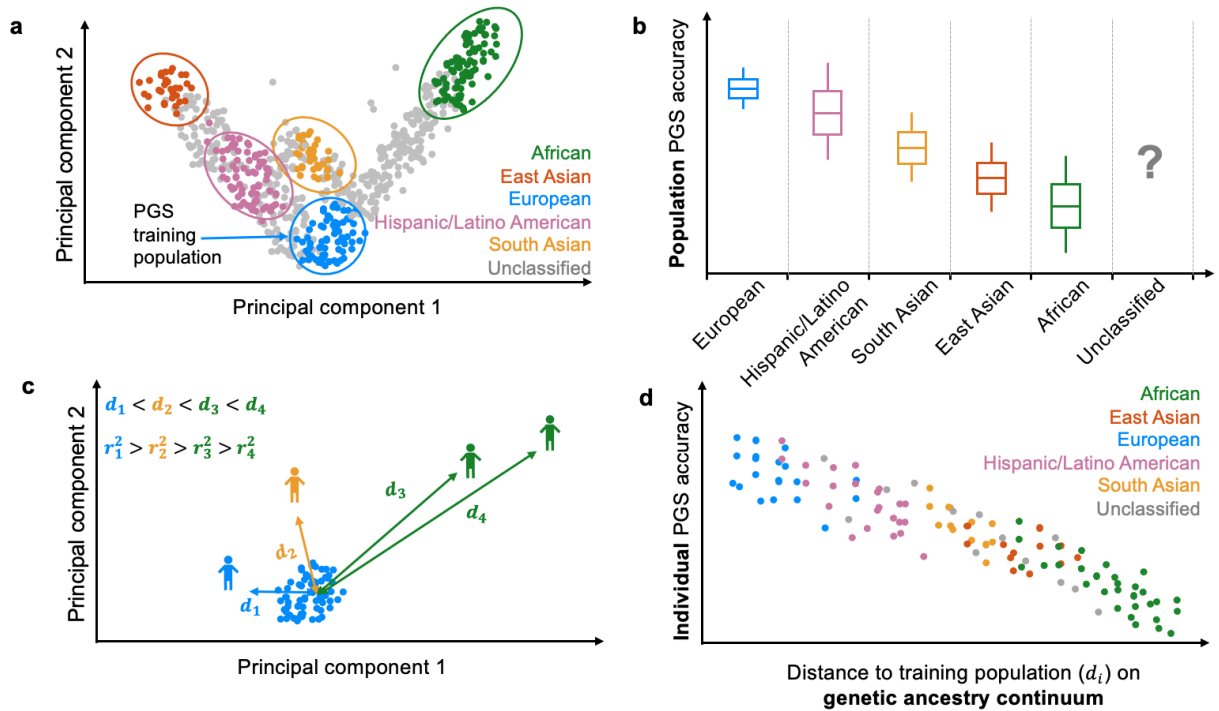
491

492 **Acknowledgments**

493 This research was conducted using the UK Biobank Resource under application 33297. We thank the
494 participants of UK Biobank for making this work possible. This work was funded in part by NIH awards
495 U01HG011715, R01HG009120, R01MH115676. The content is solely the responsibility of the authors and
496 does not necessarily represent the official views of the NIH.

497

498 **Figures**

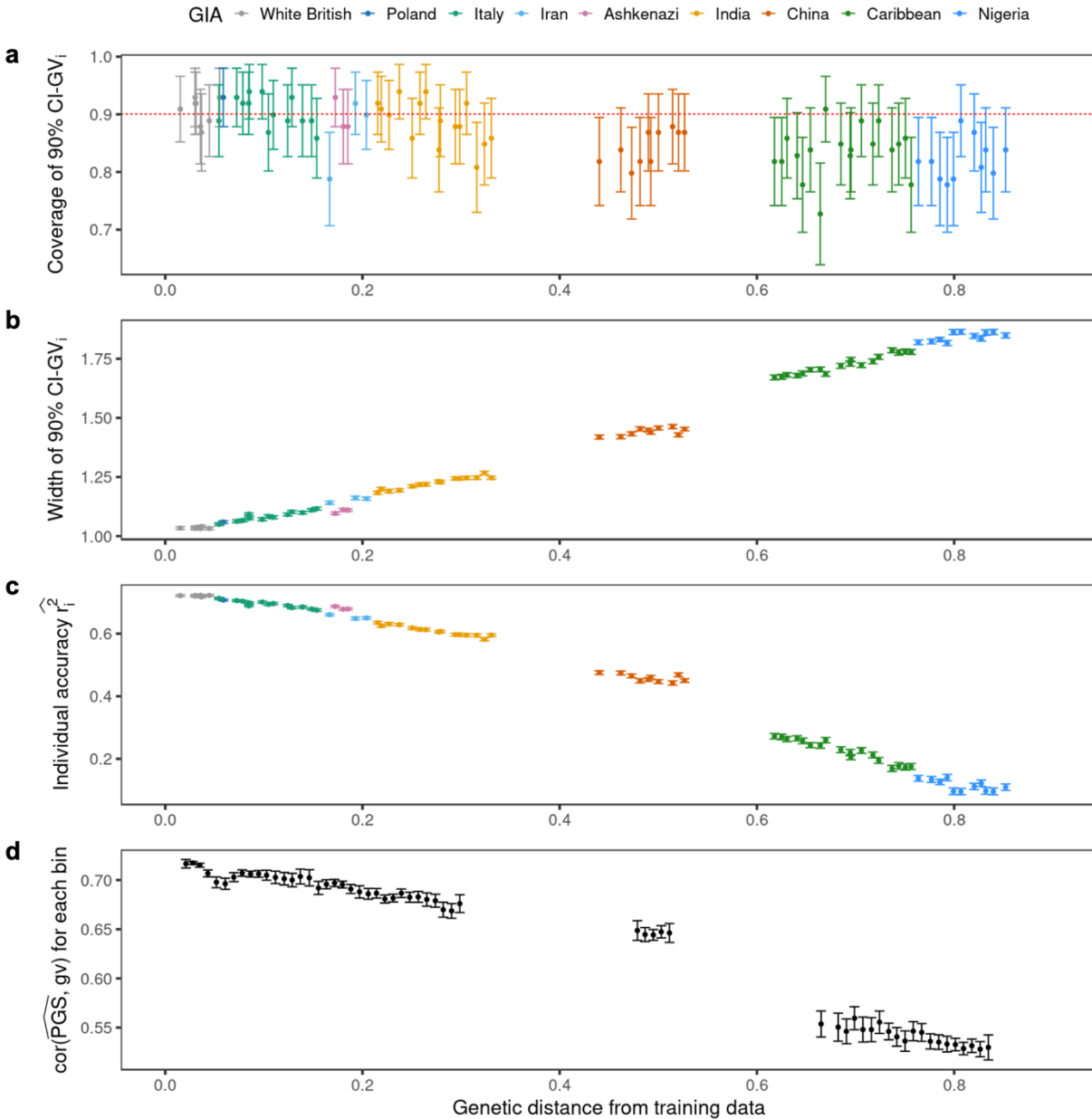


499
 500 **Figure 1. Population-level vs individual-level PGS accuracy.** (a) Discrete labeling of genetic ancestry
 501 with PCA-based clustering. Each dot represents an individual. The circles represent arbitrary boundaries
 502 imposed on the genetic ancestry continuum to divide individuals into different genetic ancestry clusters.
 503 The color represents the assigned genetic ancestry cluster label. The gray dots are individuals who are left
 504 unclassified. (b) Population-level PGS accuracy varies across clusters. The box plot represents the PGS
 505 accuracy (e.g., R^2) measured at population level. The question mark emphasizes that the PGS accuracy for
 506 unclassified individuals is unknown due to the lack of a reference group. Gray dashed lines emphasize the
 507 categorical nature of genetic ancestry clustering. (c) Continuous labeling of each individual's unique
 508 position on genetic ancestry continuum with a PCA-based genetic distance. The genetic distance is defined
 509 as the Euclidean distance of an individual's genotype from the center of the training data when projected
 510 on the PC space of training genotype data. Each individual has its own unique genetic distance d_i and
 511 individual PGS accuracy r_i^2 . (d) Individual-level PGS accuracy decays along the genetic ancestry
 512 continuum. Each dot represents an individual and its color represents the genetic ancestry label assigned in
 513 panel a. Individuals labeled with the same ancestry spread out on the genetic ancestry continuum and there
 514 are no clear boundaries between genetic ancestry groups.

515

516

517



518

519 **Figure 2. PGS performance is calibrated across genetic distance in simulations using UKBB data.**

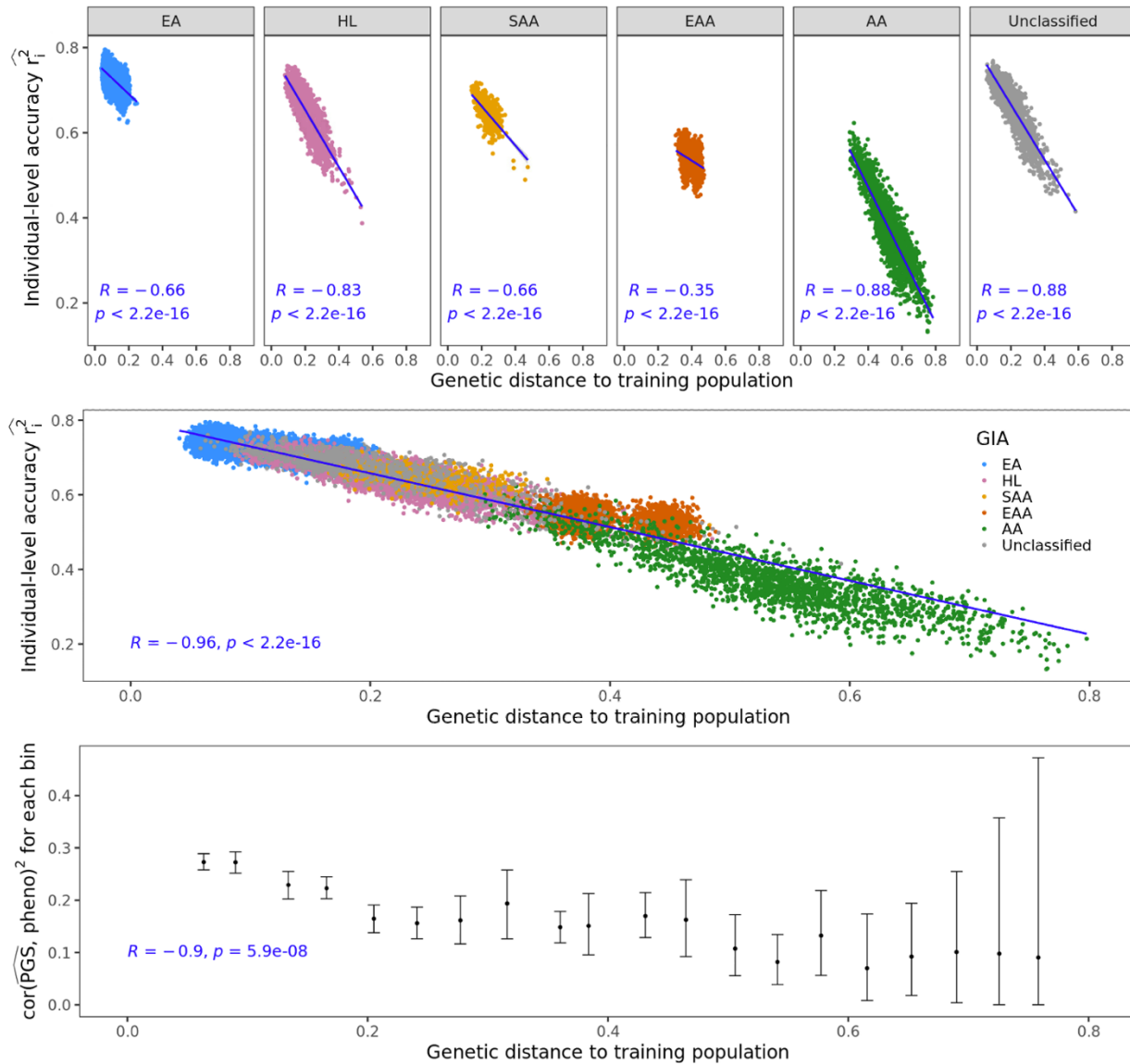
520 (a) 90% credible intervals are well calibrated for testing individuals at all genetic distances. The red
 521 dotted line represents the expected coverage of 90% credible interval. Each dot represents a randomly
 522 selected UKBB testing individual. For each dot, the x-axis is its genetic distance from training data, the y-
 523 axis is the empirical coverage of 90% credible interval calculated as the proportion of simulation
 524 replicates where the 90% credible intervals contain the individual's true genetic value, and the error bars
 525 represent mean ± 1.96 standard error of the mean (s.e.m) of the empirical coverage calculated from 100
 526 simulations. (b) The width of 90% credible interval increases with genetic distance. For each dot, the y-
 527 axis is the width of 90% credible interval across 100 simulation replicates, and the error bars represent
 528 ± 1.96 s.e.m. (c) Individual PGS accuracy decreases with genetic distance. For each dot, the y-axis is the
 529 average individual level PGS accuracy across 100 simulation replicates, and the error bars represent ± 1.96

530 s.e.m. (d) Population-level metrics of PGS accuracy recapitulates the decay in PGS accuracy across
531 genetic continuum. All UKBB testing individuals are divided into 100 equal-interval bins based on their
532 genetic distance. The x-axis is the average genetic distance for the bin and the y-axis is the squared
533 correlation between genetic value and PGS estimates for the individuals within the bin. The dot and error
534 bars represent the mean and ± 1.96 s.e.m from 100 simulations.

535

536

537



538

539

540

541

542

543

544

545

546

547

548

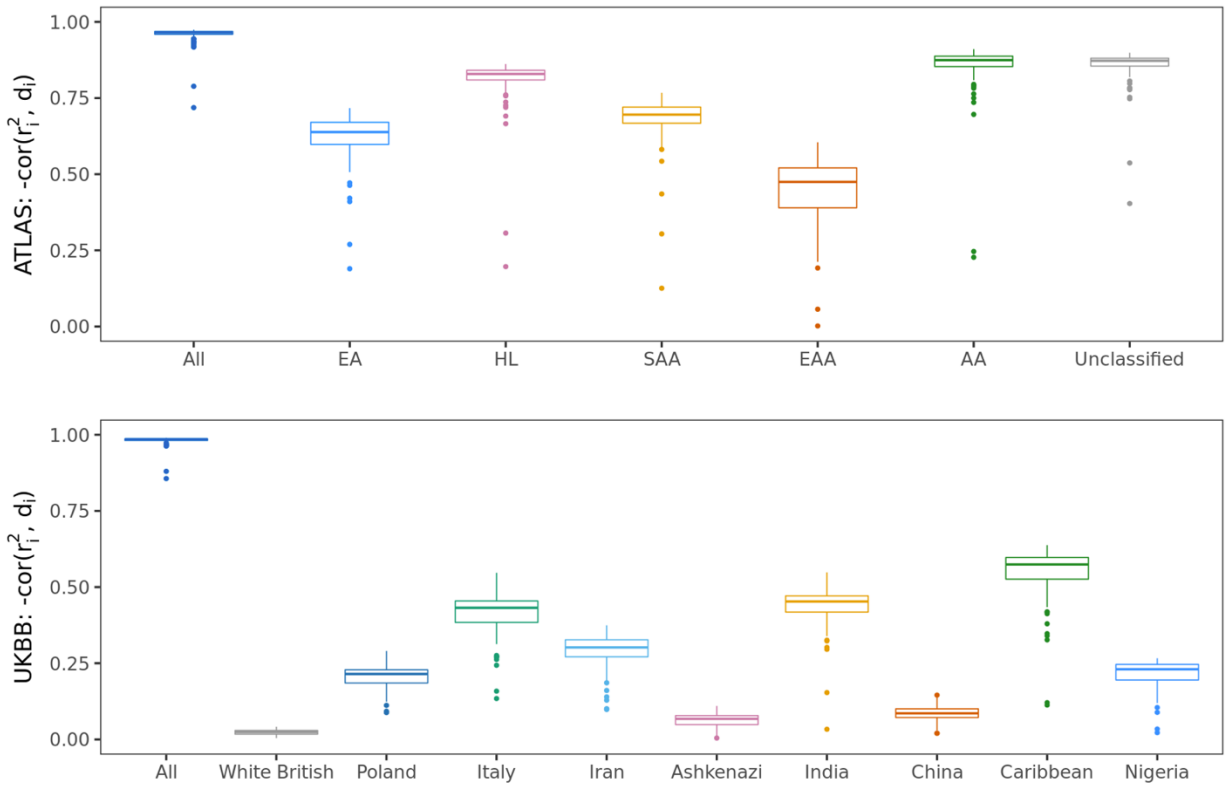
549

550

551

Figure 3. The individual-level accuracy for height PGS decreases across the genetic ancestry continuum in ATLAS. (a) Individual PGS accuracy decreases within both homogenous and admixed genetic ancestry clusters. Each dot represents a testing individual from ATLAS. For each dot, the x-axis represents its distance from the training population on the genetic continuum; the y-axis represents its PGS accuracy. The color represents the inferred genetic ancestry cluster. R and p refer to the correlation between genetic distance and individual-level PGS accuracy and its significance from two-sided t-tests. (b) Individual PGS accuracy decreases across the entire ATLAS. (c) Population-level PGS accuracy decreases with the average genetic distance in each genetic distance bin. All ATLAS individuals are divided into 20 equal-interval genetic distance bins. The x-axis is the average genetic distance within the bin, the y-axis is the squared correlation between PGS and phenotype for individuals in the bin; The dot and error bar show mean and 95% confidence interval from 1000 bootstrap samples. (EA, European American; HL, Hispanic/Latino American; SAA, South Asian American; EAA, East Asian American; AA, African American.)

552

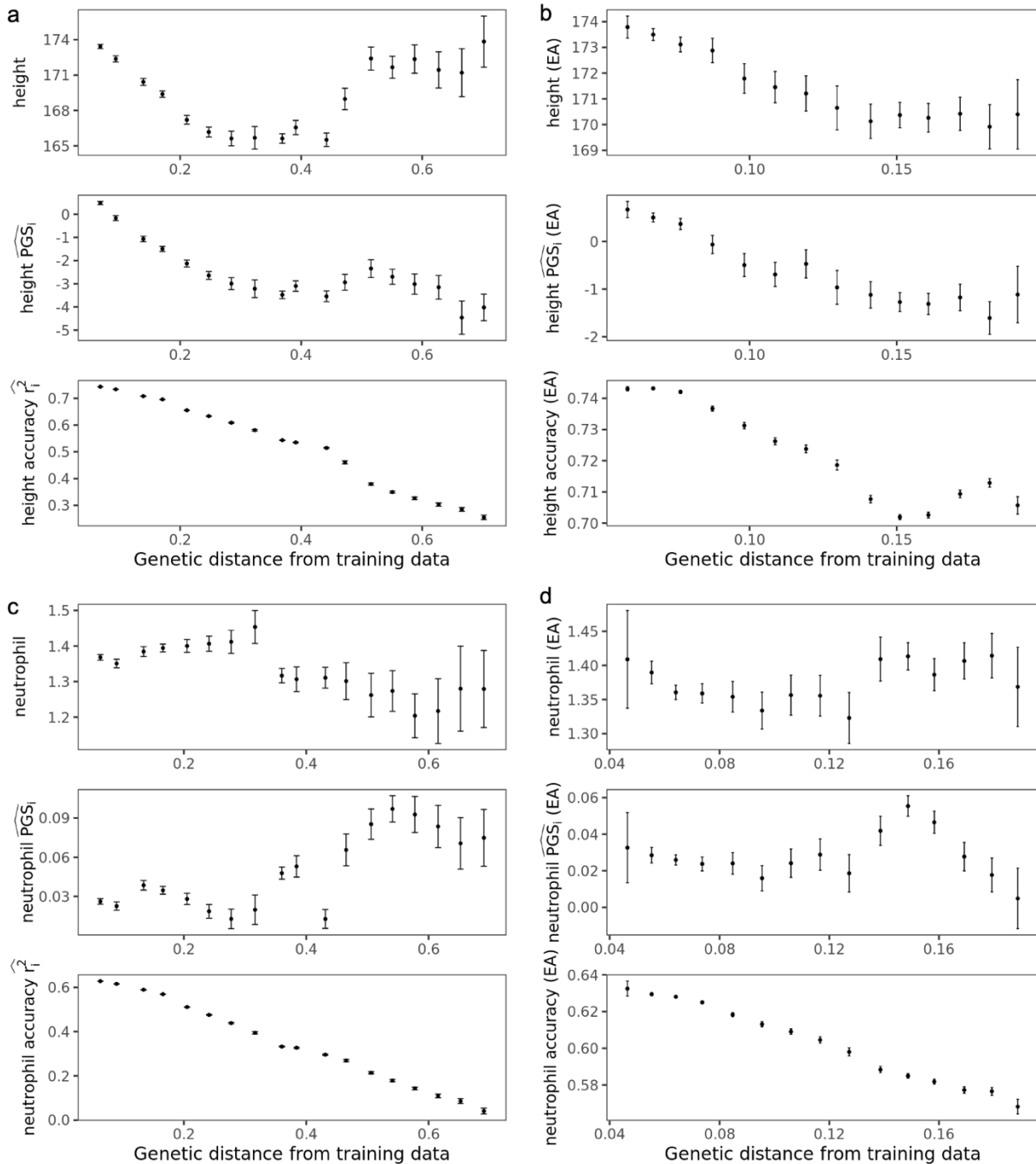


553

554 **Figure 4. The correlation between individual PGS accuracy and genetic distance is pervasive across**
555 **84 traits across ATLAS and UKBB.** (a) The distribution of correlation between PGS accuracy and genetic
556 distance for 84 traits in ATLAS. (b) The distribution of correlation between PGS accuracy and genetic
557 distance for 84 traits in UKBB. Each boxplot contains 84 points corresponding to the correlation between
558 PGS accuracy and genetic distance within the group specified by x-axis for each of the 84 traits. The box
559 shows the first, second and third quartile of the 84 correlations, and whiskers extend to the minimum and
560 maximum estimates located within $1.5 \times$ IQR from the first and third quartiles, respectively. Numerical
561 results are reported in Supplementary Table 2 and 3.

562

563



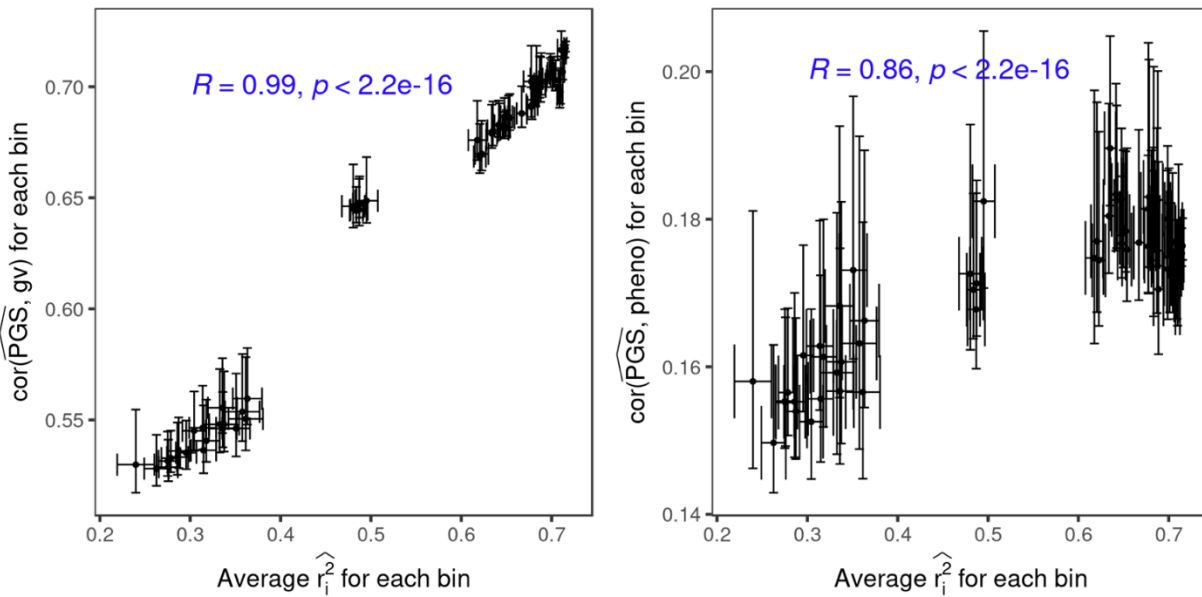
564

565 **Figure 5. Measured phenotype, PGS estimates, and accuracy varies across the ATLAS and within**
 566 **European American genetic ancestry clusters.** (a) Variation of height phenotype, PGS estimates and
 567 accuracy across different genetic distance bins in ATLAS. The 36,778 ATALS individuals are divided into
 568 20 equal-interval genetic distance bins. The x-axis is the average genetic distance within the bin, the y-axis
 569 is the average phenotype (top), PGS (middle) and individual PGS accuracy (bottom). The error bars
 570 represent +/- 1.96 standard error of the mean. Bins with fewer than 50 individuals are not shown due to
 571 large standard error of the mean. (b) Variation of height phenotype, PGS estimates and accuracy across
 572 different genetic distance bins within European American (EA) genetic ancestry clusters in ATLAS. The

573 22,380 EA individuals are divided into 20 equal-interval genetic distance bins. The x-axis is the average
574 genetic distance within the bin, the y-axis is the average phenotype (top), PGS (middle) and individual PGS
575 accuracy (bottom). The error bars represent +/- 1.96 standard error of the mean. Bins with fewer than 50
576 individuals are not shown due to large standard error of the mean. (c) Variation of log neutrophil counts,
577 PGS estimates and accuracy across different genetic distance bins across ATLAS. (d) Variation of log
578 neutrophil counts, PGS estimates and accuracy across different genetic distance bins within European
579 American (EA) genetic ancestry clusters in ATLAS.
580

581

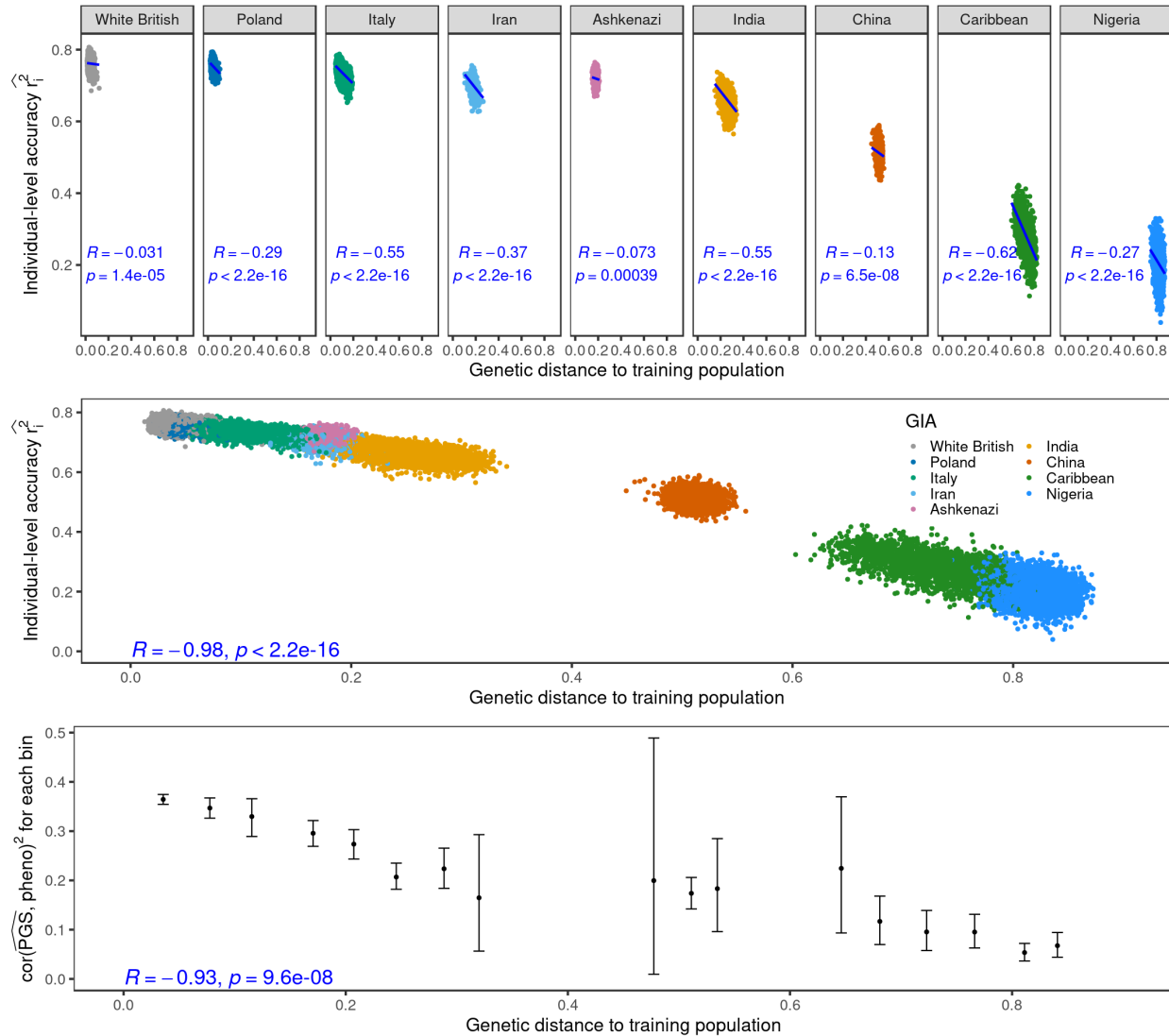
582 **Supplementary Figures**



583

584 **Supplementary Figure 1. The individual level accuracy is highly correlated with population level**
585 **accuracy.** All UKBB testing individuals are divided into 100 bins based on their genetic distance. The x-
586 axis is the average individual-level PGS accuracy for the individuals within the bin and the y-axis is (a)
587 the squared correlation between simulated genetic value and PGS estimates for the individuals within the
588 bin (b) the squared correlation between simulated phenotype and PGS estimates. The dot and error bars
589 represent the mean and ± 1.96 s.e.m from 100 simulations.

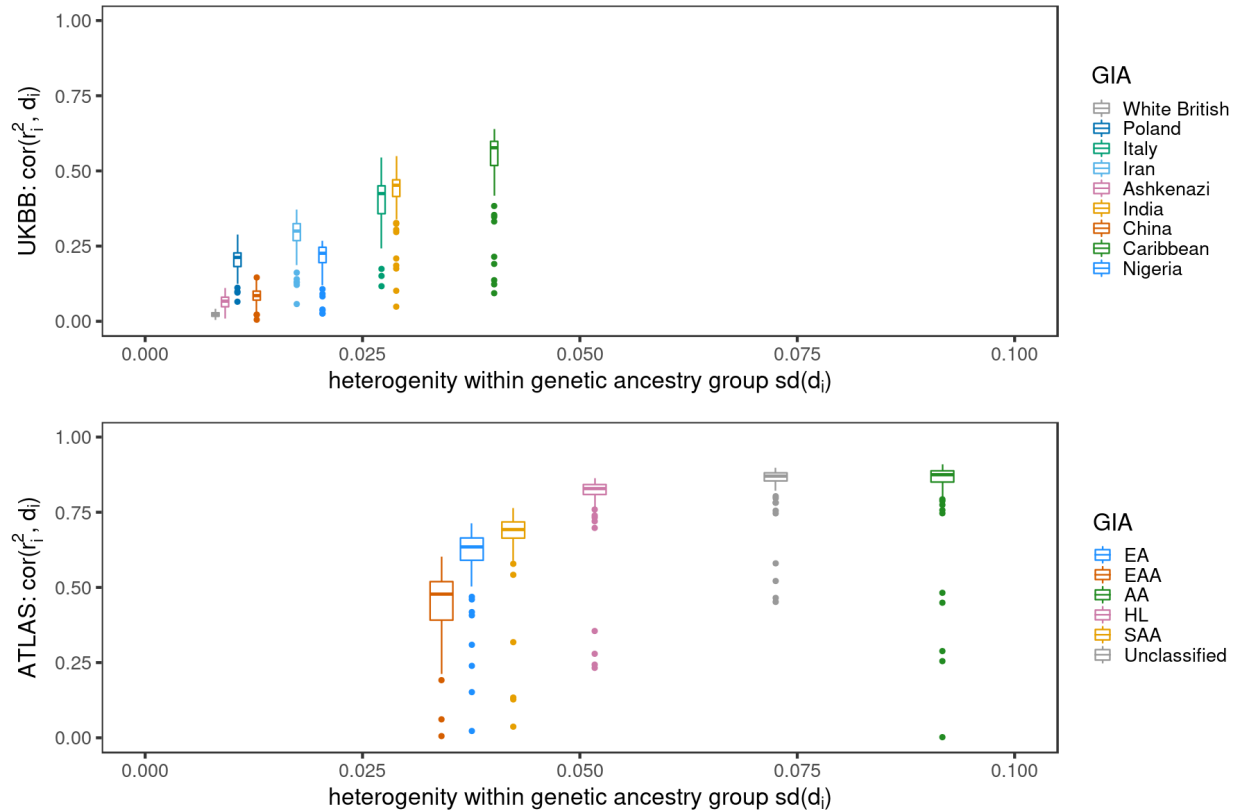
590



591

592 **Supplementary Figure 2. The individual-level accuracy for height PGS decreases across the genetic**
 593 **ancestry continuum in UKBB.** (a) Individual PGS accuracy decreases within subcontinental admixed
 594 genetic ancestry clusters. Each dot represents a testing individual from UKBB. For each dot, the x-axis
 595 represents its distance from the training population on the genetic continuum; the y-axis represents its PGS
 596 accuracy. The color represents the inferred genetic ancestry cluster. R and p refer to the correlation between
 597 genetic distance and individual-level PGS accuracy and its significance from two-sided t-tests. (b)
 598 Individual PGS accuracy decreases across the entire UKBB. (c) The population PGS accuracy decreases
 599 with the average genetic distance in each genetic distance bin. All UKBB individuals are divided into 20
 600 equal-interval genetic distance bins. The x-axis is the average genetic distance within the bin; the y-axis is
 601 the squared correlation between PGS and phenotype for individuals in the bin. The dot and error bar show
 602 mean and 95% confidence interval from 1000 bootstrap samples.

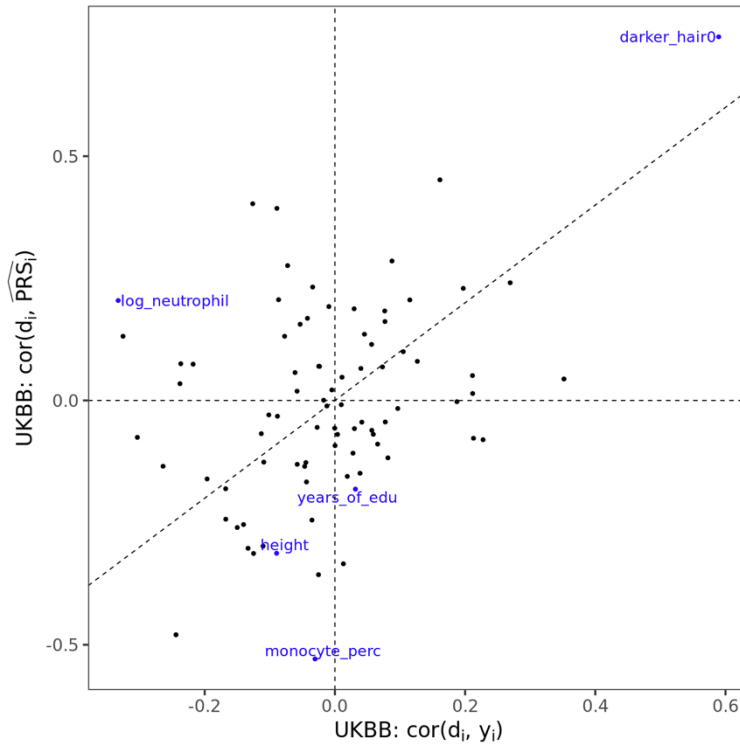
603



604

605 **Supplementary Figure 3. Lower heterogeneity within the genetic ancestry group corresponds to a**
606 **lower correlation between genetic distance and individual PGS accuracy** (a) The distribution of
607 correlations between PGS accuracy and genetic distance for 84 traits in ATLAS. (b) The distribution of
608 correlations between PGS accuracy and genetic distance for 84 traits in UKBB. The x-axis is the
609 homogeneity of the genetic ancestry clusters measured as standard deviation of genetic distance within a
610 genetic ancestry cluster; a larger $sd(d_i)$ indicates a larger variation of genetic background. Each boxplot
611 contains 84 points corresponding to the correlation between PGS accuracy and genetic distance within the
612 group specified by x-axis for each of the 84 traits. The box shows the first, second and third quartile of the
613 84 correlations, and whiskers extend to the minimum and maximum estimates located within $1.5 \times$ IQR
614 from the first and third quartiles, respectively.

615



616

617 **Supplementary Figure 4. Discordant directions of phenotype/PGS-distance correlations in UKBB.**

618 The x axis is the correlation between phenotype and genetic distance and the y axis is the correlation
619 between PGS estimates and genetic distance for all 48,586 testing individuals in UKBB. Numerical results
620 are reported in Supplementary Table 4.

621

622

623 **Supplementary Table 1. The training sample size, proportion of causal variants and heritability of**
624 **the 84 traits.**

625

626 **Supplementary Table 2. The correlation between individual PGS accuracy and genetic distance**
627 **from training data across ATLAS and within each genetic ancestry clusters**

628

629 **Supplementary Table 3. The correlation between individual PGS accuracy and genetic distance**
630 **from training data across UKBB and within each genetic ancestry clusters**

631

632 **Supplementary Table 4. The correlation between measured phenotype/PGS and genetic distance**
633 **from training data across UKBB**

Reference

1. Torkamani, A., Wineinger, N. E. & Topol, E. J. The personal and clinical utility of polygenic risk scores. *Nat. Rev. Genet.* **19**, 581–590 (2018).
2. Lambert, S. A., Abraham, G. & Inouye, M. Towards clinical utility of polygenic risk scores. *Hum. Mol. Genet.* **28**, R133–R142 (2019).
3. Kullo, I. J. *et al.* Polygenic scores in biomedical research. *Nat. Rev. Genet.* (2022)
doi:10.1038/s41576-022-00470-z.
4. Li, R., Chen, Y., Ritchie, M. D. & Moore, J. H. Electronic health records and polygenic risk scores for predicting disease risk. *Nat. Rev. Genet.* **21**, 493–502 (2020).
5. Natarajan, P. *et al.* Polygenic Risk Score Identifies Subgroup With Higher Burden of Atherosclerosis and Greater Relative Benefit From Statin Therapy in the Primary Prevention Setting. *Circulation* **135**, 2091–2101 (2017).
6. Lee, A. *et al.* BOADICEA: a comprehensive breast cancer risk prediction model incorporating genetic and nongenetic risk factors. *Genet. Med.* **21**, 1708–1718 (2019).
7. Khera, A. V. *et al.* Polygenic Prediction of Weight and Obesity Trajectories from Birth to Adulthood. *Cell* **177**, 587-596.e9 (2019).
8. Perkins, D. O. *et al.* Polygenic Risk Score Contribution to Psychosis Prediction in a Target Population of Persons at Clinical High Risk. *Am. J. Psychiatry* **177**, 155–163 (2020).
9. Martin, A. R. *et al.* Human Demographic History Impacts Genetic Risk Prediction across Diverse Populations. *Am. J. Hum. Genet.* **100**, 635–649 (2017).
10. Martin, A. R. *et al.* Clinical use of current polygenic risk scores may exacerbate health disparities.

- Nat. Genet.* **51**, 584–591 (2019).
11. Mostafavi, H. *et al.* Variable prediction accuracy of polygenic scores within an ancestry group. *Elife* **9**, e48376 (2020).
 12. Wang, Y. *et al.* Theoretical and empirical quantification of the accuracy of polygenic scores in ancestry divergent populations. *Nat. Commun.* **11**, 3865 (2020).
 13. Scutari, M., Mackay, I. & Balding, D. Using Genetic Distance to Infer the Accuracy of Genomic Prediction. *PLoS Genet.* **12**, e1006288 (2016).
 14. Duncan, L. *et al.* Analysis of polygenic risk score usage and performance in diverse human populations. *Nat. Commun.* **10**, 1–9 (2019).
 15. Coop, G. Genetic similarity and genetic ancestry groups. *arXiv [q-bio.PE]* (2022).
 16. Lewis, A. C. F. *et al.* Getting genetic ancestry right for science and society. *Science* **376**, 250–252 (2022).
 17. Mathieson, I. & Scally, A. What is ancestry? *PLoS Genet.* **16**, e1008624 (2020).
 18. Krainc, T. & Fuentes, A. Genetic ancestry in precision medicine is reshaping the race debate. *Proceedings of the National Academy of Sciences* **119**, e2203033119 (2022).
 19. Belbin, G. M. *et al.* Toward a fine-scale population health monitoring system. *Cell* **184**, 2068–2083.e11 (2021).
 20. Wojcik, G. L. *et al.* Genetic analyses of diverse populations improves discovery for complex traits. *Nature* **570**, 514–518 (2019).
 21. Bitarello, B. D. & Mathieson, I. Polygenic Scores for Height in Admixed Populations. *G3* **10**, 4027–4036 (2020).

22. Clarke, S. L. *et al.* Race and Ethnicity Stratification for Polygenic Risk Score Analyses May Mask Disparities in Hispanics. *Circulation* **146**, 265–267 (2022).
23. Privé, F. *et al.* Portability of 245 polygenic scores when derived from the UK Biobank and applied to 9 ancestry groups from the same cohort. *Am. J. Hum. Genet.* **109**, 12–23 (2022).
24. Johnson, R. D. *et al.* The UCLA ATLAS Community Health Initiative: promoting precision health research in a diverse biobank. *medRxiv* (2022).
25. Johnson, R. *et al.* Leveraging genomic diversity for discovery in an electronic health record linked biobank: the UCLA ATLAS Community Health Initiative. *Genome Med.* **14**, 1–23 (2022).
26. Ding, Y. *et al.* Large uncertainty in individual polygenic risk score estimation impacts PRS-based risk stratification. *Nat. Genet.* **54**, 30–39 (2021).
27. Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* **38**, 904–909 (2006).
28. Alexander, D. H. & Lange, K. Enhancements to the ADMIXTURE algorithm for individual ancestry estimation. *BMC Bioinformatics* **12**, 246 (2011).
29. Privé, F., Arbel, J. & Vilhjálmsson, B. J. LDpred2: better, faster, stronger. *Cold Spring Harbor Laboratory* 2020.04.28.066720 (2020) doi:10.1101/2020.04.28.066720.
30. Walsh, B. & Lynch, M. Evolution and Selection of Quantitative Traits. in *Evolution and Selection of Quantitative Traits* (Oxford University Press, 2018).
31. Reich, D. *et al.* Reduced neutrophil count in people of African descent is due to a regulatory variant in the Duffy antigen receptor for chemokines gene. *PLoS Genet.* **5**, e1000360 (2009).
32. Fatumo, S. *et al.* A roadmap to increase diversity in genomic studies. *Nat. Med.* **28**, 243–250 (2022).

33. Martin, A. R. *et al.* Increasing diversity in genomics requires investment in equitable partnerships and capacity building. *Nat. Genet.* **54**, 740–745 (2022).
34. Wang, Y., Tsuo, K., Kanai, M., Neale, B. M. & Martin, A. R. Challenges and Opportunities for Developing More Generalizable Polygenic Risk Scores. *Annu. Rev. Biomed. Data Sci.* **5**, 293–320 (2022).
35. Ruan, Y. *et al.* Improving polygenic prediction in ancestrally diverse populations. *Nat. Genet.* **54**, 573–580 (2022).
36. Spence, J. P., Sinnott-Armstrong, N., Assimes, T. L. & Pritchard, J. K. A flexible modeling and inference framework for estimating variant effect sizes from GWAS summary statistics. *bioRxiv* 2022.04.18.488696 (2022) doi:10.1101/2022.04.18.488696.
37. Zhang, H. *et al.* Novel Methods for Multi-ancestry Polygenic Prediction and their Evaluations in 3.7 Million Individuals of Diverse Ancestry. *bioRxiv* 2022.03.24.485519 (2022) doi:10.1101/2022.03.24.485519.
38. Shi, H. *et al.* Localizing Components of Shared Transethnic Genetic Architecture of Complex Traits from GWAS Summary Data. *Am. J. Hum. Genet.* **106**, 805–817 (2020).
39. Shi, H. *et al.* Population-specific causal disease effect sizes in functionally important regions impacted by selection. *Nat. Commun.* **12**, 1098 (2021).
40. Hou, K. *et al.* Accurate estimation of SNP-heritability from biobank-scale data irrespective of genetic architecture. *Nat. Genet.* **51**, 1244–1251 (2019).
41. Ge, T., Chen, C.-Y., Ni, Y., Feng, Y.-C. A. & Smoller, J. W. Polygenic prediction via Bayesian regression and continuous shrinkage priors. *Nat. Commun.* **10**, 1776 (2019).
42. Zhang, Q., Privé, F., Vilhjálmsson, B. & Speed, D. Improved genetic prediction of complex traits

from individual-level data or summary statistics. *Cold Spring Harbor Laboratory*

2020.08.24.265280 (2020) doi:10.1101/2020.08.24.265280.

43. Lloyd-Jones, L. R. *et al.* Improved polygenic prediction by Bayesian multiple regression on summary statistics. *Nat. Commun.* **10**, 1–11 (2019).
44. Sorensen, D. & Gianola, D. *Likelihood, Bayesian, and MCMC Methods in Quantitative Genetics*. (Springer New York).
45. Abraham, G. & Inouye, M. Fast principal component analysis of large-scale genome-wide data. *PLoS One* **9**, e93766 (2014).
46. Abraham, G., Qiu, Y. & Inouye, M. FlashPCA2: principal component analysis of Biobank-scale genotype datasets. *Bioinformatics* **33**, 2776–2778 (2017).
47. Privé, F., Aschard, H., Ziyatdinov, A. & Blum, M. G. B. Efficient analysis of large-scale genome-wide data with two R packages: bigstatsr and bigsnpr. *Bioinformatics* **34**, 2781–2787 (2018).
48. Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nat. Genet.* **47**, 291–295 (2015).