

1 Three Open Questions in Polygenic Score Portability

2 Joyce Y. Wang¹, Neeka Lin¹, Michael Zietz², Jason Mares³, Vagheesh M. Narasimhan^{1,4},
3 Paul J. Rathouz^{4,5} and Arbel Harpak^{1,5,+}

4 ¹ Department of Integrative Biology, The University of Texas at Austin, Austin, TX

5 ² Department of Biomedical Informatics, Columbia University, New York, NY

6 ³ Department of Neurology, Columbia University, New York, NY

7 ⁴ Department of Statistics and Data Science, The University of Texas at Austin, Austin, TX

8 ⁵ Department of Population Health, The University of Texas at Austin, Austin, TX

9 ⁺ Correspondence should be addressed to A.H. (arbelharpak@utexas.edu)

10 **Abstract**

11 A major obstacle hindering the broad adoption of polygenic scores (PGS) is their lack of
12 “portability” to people that differ—in genetic ancestry or other characteristics—from the
13 GWAS samples in which genetic effects were estimated. Here, we use the UK Biobank to
14 measure the change in PGS prediction accuracy as a continuous function of individuals’
15 genome-wide genetic dissimilarity to the GWAS sample (“genetic distance”). Our results
16 highlight three gaps in our understanding of PGS portability. First, prediction accuracy
17 is extremely noisy at the individual level and not well predicted by genetic distance. In
18 fact, variance in prediction accuracy is explained comparably well by socioeconomic mea-
19 sures. Second, trends of portability vary across traits. For several immunity-related traits,
20 prediction accuracy drops near zero quickly even at intermediate levels of genetic distance.
21 This quick drop may reflect GWAS associations being more ancestry-specific in immunity-
22 related traits than in other traits. Third, we show that even qualitative trends of portability
23 can depend on the measure of prediction accuracy used. For instance, for white blood cell
24 count, a measure of prediction accuracy at the individual level (reduction in mean squared
25 error) increases with genetic distance. Together, our results show that portability cannot
26 be understood through global ancestry groupings alone. There are other, understudied fac-
27 tors influencing portability, such as the specifics of the evolution of the trait and its genetic
28 architecture, social context, and the construction of the polygenic score. Addressing these
29 gaps can aid in the development and application of PGS and inform more equitable genomic
30 research.

31 Introduction

32 Polygenic scores (PGS), genetic predictors of complex traits based on genome-wide associa-
33 tion studies (GWAS), are gaining traction among researchers and practitioners^{13,15,25}. Yet a
34 major problem hindering their broad application is their highly variable performance across
35 prediction samples^{19,6,10,14}. Often, prediction accuracy appears to decline in groups unlike the
36 GWAS sample—in genetic ancestry, social context or environmental exposures^{17,19,10,32,42,30},
37 restricting the contexts in which PGS can be used reliably.

38 This so-called “portability” problem is a subject of intense study. Typically, portability
39 is evaluated through variation in the within-group phenotypic variance explained by a PGS
40 (i.e., the coefficient of determination, R^2) among genetic ancestry groups. Indeed, population
41 genetics theory gives clear predictions for the relationship between genetic dissimilarity to
42 the GWAS sample and PGS prediction accuracy under some models (neutral evolution^{27,44,3},
43 directional²³, or stabilizing selection^{44,23}), all else being equal (including, e.g., assumptions
44 about environmental effects).

45 However, inference based on empirical variation in R^2 can be misleading for various
46 reasons. For one, it can be arbitrarily low even when the model fitted to the data is correct.
47 It also cannot be compared across transformations of the data. R^2 is not comparable across
48 datasets, because, for instance, it depends on the extent of variation in the independent
49 variable^{39,16,34}. In the context of inference about the causes of PGS portability, these issues
50 can manifest in different ways. For example, heterogeneity in within-group genetic variance
51 and environmental variance can each greatly affect group differences in R^2 .

52 A related issue is that the impacts of environmental and social factors on portability are
53 not well understood, despite evidence illustrating these impacts can be substantial^{19,10,43,20}.
54 To complicate matters, such factors may be confounded with genetic ancestry, limiting our
55 ability to make inferences based on the typical decay of R^2 between PGS and trait value in
56 ancestries less represented in GWAS samples^{19,25,10,43}.

57 With these limitations of R^2 , and the possible confounding with environmental and social

58 factors, it remains unclear how well genetic ancestry would predict the applicability of PGS
59 for individuals. Recent work implied that individual-level prediction accuracy should be
60 largely explained by genome-wide genetic dissimilarity to the GWAS sample (see figure 3
61 in [7] and figure 5 in [38]). However, we note that this work focused on the relationship
62 between genetic distance and the prediction interval, i.e. expected uncertainty in prediction
63 under an assumed model, rather than the relationship with the realized prediction accuracy.
64 Understanding the drivers of variation in prediction accuracy is especially pertinent for
65 personalized clinical risk predictions and decisions regarding their reporting to patients^{15,12}.

66 This motivated us to empirically study PGS prediction accuracy at the individual level.
67 In what follows, we highlight three puzzling observations that also point to three gaps in our
68 understanding of the portability problem: (1) Genetic dissimilarity to the GWAS sample
69 poorly predicts portability at the individual level, (2) portability trends (with respect to
70 genetic distance) can be trait-specific; and (3) portability trends depend on the measure of
71 prediction accuracy. Informed by our results, we suggest avenues of future research that can
72 help bridge these gaps.

73 Results

74 **Portability and individual-level genetic distance from the GWAS sample.** We ex-
75 amined PGS portability as a function of genetic distance from the GWAS sample in the UK
76 Biobank (UKB). For each of 15 continuous physiological traits, we performed a GWAS in
77 a sample of 350,000 individuals. For 129,279 individuals not included in the GWAS sample
78 (henceforth referred to as “prediction sample”), we predicted the trait value using the PGS
79 and covariates. Using a Principal Component Analysis (PCA) of the genotype matrix of
80 the entire sample, we quantify each individual’s genetic distance from the GWAS sample as
81 distance from the centroid of GWAS individuals’ coordinates in PCA space (**Fig. 1A**). This
82 measure is quicker to compute, yet highly correlated with F_{st} between the GWAS sample and
83 single individuals in the prediction sample ($r > 0.98$), albeit noticeably less reflective of F_{st}
84 at intermediate genetic distances (**Fig. 1B**). The imperfect correlation may be a result of our

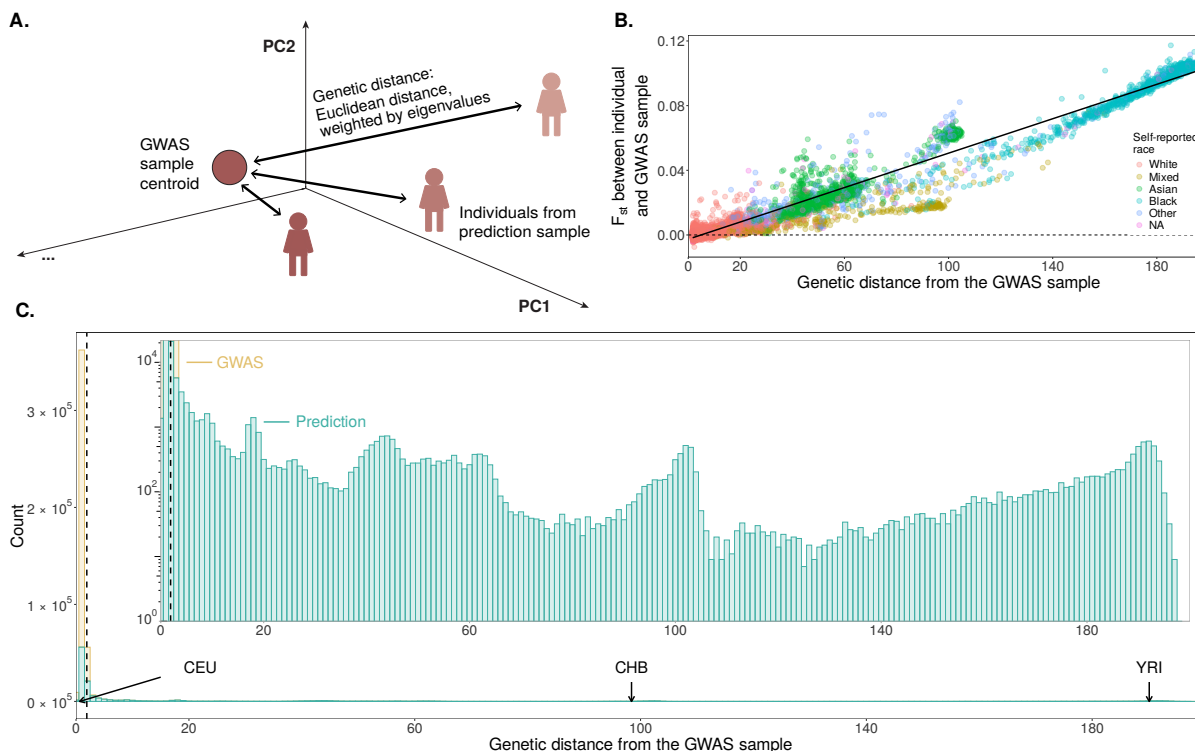


Figure 1: Measuring “genetic distance” from the GWAS sample. **A.** Across 350,000 individuals in the GWAS sample and 129,279 individuals in the prediction set, we measure “genetic distance” from the GWAS sample as the weighted Euclidean distance from the centroid of GWAS individuals in PCA space, with each PC weighted by its respective eigenvalue. **B.** Across 10,000 individuals from the prediction set, genetic distance to the GWAS sample (calculated with 40 PCs) is highly correlated with F_{st} between the GWAS sample and the individual (**Fig. S1**). Under a theoretical model where portability is driven by genetic ancestry alone and the trait evolves neutrally, F_{st} should perfectly predict variation in prediction accuracy. We note that genetic distance is less reflective of F_{st} for intermediate genetic distances. **C.** The distribution of genetic distance. For reference, we show the mean genetic distances for subsets of the 1000 Genomes dataset²: CEU, Utah residents of primarily Northern and Western European descent; CHB, Han Chinese in Beijing, China; YRI, Yoruba in Ibadan, Nigeria. The dashed line represents the 95th percentile of genetic distance from among GWAS sample individuals. In what follows, our reports are based on individuals with genetic distances larger than this value. The inset is a zoomed-in view of a smaller range and on a log-scale, to better visualize the distribution within the prediction sample.

85 use of only the top 40 PCs^{24,27}. Under some theoretical conditions (such as neutral evolu-
 86 tion, additive contribution of genotype and environment, fixed environmental variance)— F_{st}
 87 should perfectly predict variation in prediction accuracy due to genetic ancestry^{26,27,3}. We
 88 standardized genetic distance such that its mean is 1 across GWAS sample individuals.

89 In the prediction sample, we observed a continuum of genetic distance from the GWAS
 90 sample with several clear modes, the main one at short distances: 96,457 individuals have

91 a genetic distance of up to 10 and the remaining 32,822 individuals at distances between
92 10-197.6 (**Fig. 1B, C**). To ground our expectations, we estimated the mean genetic distance
93 for three 1000 Genomes² subsamples: Utah residents of primarily Northern and Western
94 European descent (CEU) average at 0.6, Han Chinese in Beijing, China (CHB) average at
95 98.4, and Yoruba in Ibadan, Nigeria (YRI) average at 190.0 (**Fig. 1C**).

96 For each of the 15 continuous physiological traits, we measure the prediction accuracy
97 at the group and individual level with slightly different prediction models (**Methods**). In
98 both cases, we fit a prediction model regressing the trait to the polygenic score and other
99 covariates. To evaluate group-level accuracy, we split individuals into 500 bins of genetic
100 distance comprising of 258-259 individuals each. Within each bin we measure the partial R^2
101 of the polygenic score and the trait value. To evaluate individual-level accuracy, we measure
102 the squared difference between the PGS-predicted value and the trait, after residualizing the
103 trait for covariates.

104 **Prediction accuracy is weakly predicted by genetic distance.** For some traits,
105 such as height, group-level prediction accuracy decayed monotonically with genetic distance
106 from the GWAS sample, as expected and reported previously (**Fig. 2A**)^{40,27,7}. A major
107 factor driving this decay appears to be an associated decay in heterozygosity in the PGS
108 marker SNPs (**Figs. 4B,S22**; see [23, 40]). Lower heterozygosity in PGS markers impacts the
109 genetic variance a polygenic score can capture because it makes for a less variable predictor.
110 The impact of genetic distance on LD with causal variation is less straightforward⁴⁰.

111 Previous work implied that variation in individual-level prediction accuracy should be
112 largely explained by genetic distance^{7,38}. However, that was not the case in our analysis.
113 While individual-level accuracy generally decayed with distance for most traits, this corre-
114 lation was weak (**Figs. 2B,S2**). Even a flexible cubic spline fit of genetic distance explains
115 little of the variance in prediction accuracy ($R^2 = 0.31\%$).

116 In fact, individual-level prediction accuracy is explained comparably well by socioeco-
117 nomic measures (**Fig. S18-S21**). For example, we observed a steady mean increase in
118 squared prediction error across quantiles of Townsend Deprivation Index³⁷ for 9/15 of the

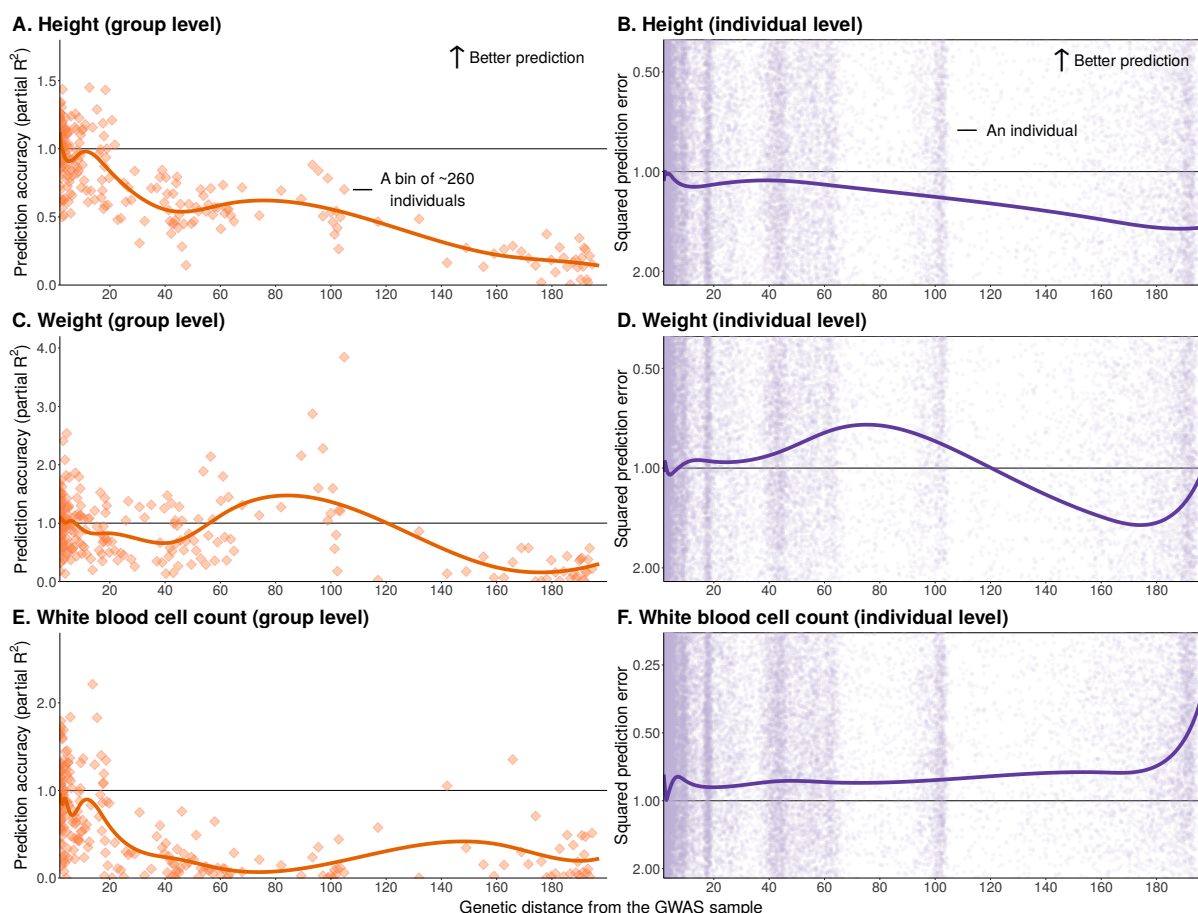
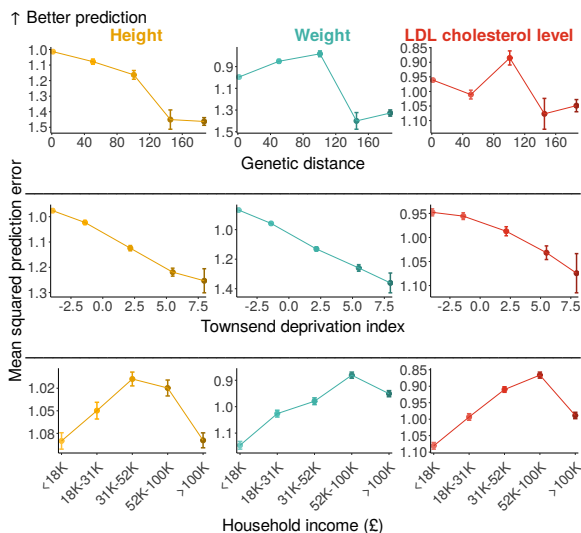


Figure 2: Trends of portability vary across traits and measures. At the group level (left panels), we measured prediction accuracy with the squared partial correlation between the PGS and the trait value in 500 bins of 258-259 individuals each. At the individual level (right panels), we measured the squared prediction error. Curves show cubic spline fits, with 8 knots placed based on the density of data points. **A**, **B**. For height, prediction accuracy decays nearly monotonically with genetic distance at both the group (A) and individual (B) levels. **C**, **D**. For weight, prediction accuracy does not monotonically decay with genetic distance. **E**, **F**. For white blood cell count, at the group level, prediction accuracy drops near zero at a short genetic distance from the GWAS sample (E); yet at the individual level, it increases (F). See **Fig. S2-S5** for other traits and **Fig. S6-S9** for plots showing the full ranges of individual-level prediction accuracy.

119 traits examined, suggesting poorer prediction in individuals of lower socioeconomic status
 120 (**Figs. 3A,S14,S17,S16**; the four exceptions being white blood cell-related traits, **Fig. S15**;
 121 see also similar reports in [19, 10]). Like genetic distance, the Townsend Deprivation Index
 122 only explains between 0.02% and 0.53% of the variance in squared prediction error across
 123 traits with a cubic spline. Notably, however, for the majority of traits, more variance is
 124 explained by this measure of socioeconomic status than by genetic distance (**Fig. 3B**).

A. Mean trends in individual-level prediction accuracy



B. Deprivation index and genetic distance explain portability comparably well

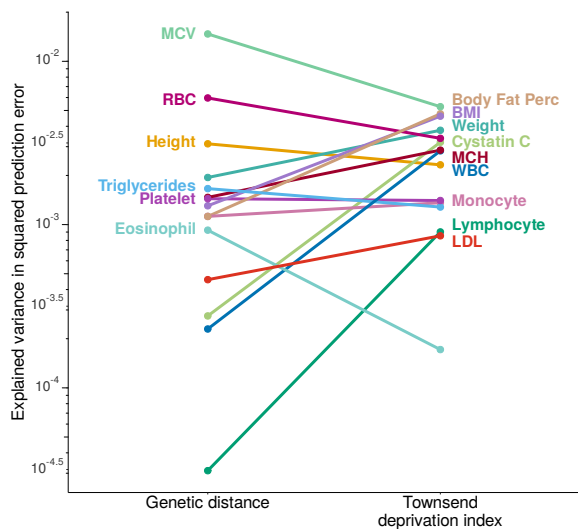


Figure 3: Genetic distance and socioeconomic factors explain individual-level prediction accuracy comparably well. **A.** Data points confer to mean (\pm SE) squared prediction errors of individuals in the prediction sample (divided by a constant, the mean squared prediction error in a reference group), binned into 5 equidistant strata. The x-axis shows the median measure value for each stratum. “Household income” refers to average yearly total household income before tax. See **Fig. S14-S17** for other traits. **B.** We compared the variance in squared prediction error explained by a cubic spline fit to genetic distance to the variance explained by a cubic spline fit to the Townsend deprivation index. MCV: mean corpuscular volume. MCH: mean corpuscular hemoglobin. RBC: red blood cell count. Body fat perc: body fat percentage. WBC: white blood cell count. LDL: LDL cholesterol level. See **Fig. S18-S21** for the variance explained by other genetic and socioeconomic measures.

125 **Trends of portability vary across traits.** Previous reports suggested that the re-
 126 lationship between genetic distance and prediction accuracy is similar across traits^{17,27,7}.
 127 However, we observed variation in this relationship among traits. Unlike the case of height,
 128 the prediction accuracy for many other traits did not decay monotonically with genetic dis-
 129 tance. Weight, mean corpuscular volume (MCV), mean corpuscular hemoglobin (MCH) and
 130 body fat percentage peaked in accuracy at intermediate genetic distances (**Fig. 2D, Fig. S2,**
 131 **S4**).

132 In other traits we examined, in particular white blood cell-related traits, group-level pre-
 133 diction accuracy dropped near zero even at a short genetic distance (**Fig. 2E,S3**). There
 134 are multiple possible drivers of trait-specific portability trends. We considered, in partic-
 135 ular, variable selective pressures on the immune system across time and geography. We

136 hypothesized that these would lead to less portable genetic associations (across ancestry)
137 compared to other traits. To test this prediction, we re-estimated the effects of index SNPs
138 (SNPs included in the PGS, ascertained in the original GWAS sample) in two subsets of the
139 prediction sample, one closer and another farther (in terms of genetic distance) from the
140 GWAS sample. The prediction sample based allelic effect estimates were least consistent
141 with the original GWAS for lymphocyte count, compared, e.g., to triglyceride levels, a trait
142 of similar SNP heritability (**Fig. 4A**). To further illustrate this point, 30.8% of index SNPs
143 for lymphocyte count had a different sign when estimated in the original GWAS and in the
144 “closer” GWAS, compared to 3.1% for triglyceride levels.

145 The rapid turnover of allelic effects may also interact with statistical biases. Consider, for
146 example, “winners curse”, whereby effect estimates are inflated due to the ascertainment of
147 index SNPs and the estimation of their effects in the same sample¹⁸. Winners curse would be
148 most severe in large effect PGS index SNPs: These SNPs are typically at lower frequencies
149 in the GWAS sample than small effect index SNPs, because GWAS power scales with the
150 product of squared allelic effect and heterozygosity^{35,22,23}. If causal effects on lymphocyte
151 count change rapidly, then large effect index SNPs may be under weaker selective constraint
152 in the prediction sample than in the GWAS sample, and segregate at high allele frequencies.
153 Indeed, for lymphocyte count, the heterozygosity of large effect variants increases with ge-
154 netic distance from the GWAS sample (**Fig. 4B**; see **Fig. S22** for other traits). As a result
155 of the trends of heterozygosity, the variance in the polygenic score (a sum over index SNP
156 heterozygosity multiplied by their squared effect estimates) quickly increases with genetic
157 distance for white blood cell count, lymphocyte count, and monocyte count, despite decreas-
158 ing for the remaining 12 traits we have examined (**Fig. 4C**). And so, taken together, the
159 PGS variance increases quickly and allelic effect estimates become non-predictive even close
160 to the GWAS sample (**Fig. S24**). Together, this may drive the immediate drop in prediction
161 accuracy of white blood-cell related traits.

162 **The measure of predictive performance can alter our view of portability.** Fi-
163 nally, the qualitative trends of portability can even depend on the measure of prediction

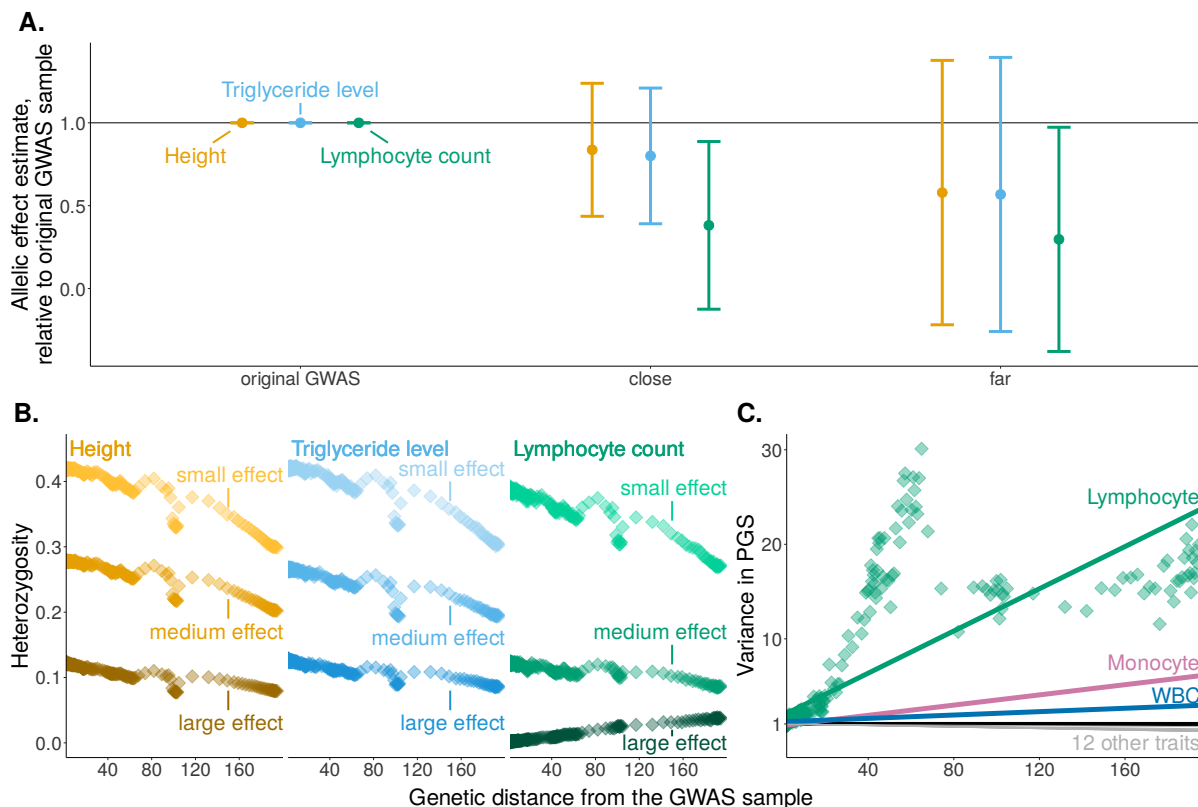


Figure 4: Lymphocyte count as an example of trait-dependent factors influencing portability. **A.** We re-estimated the allelic effects of PGS index SNPs in subsamples of the prediction set: “close” (genetic distance ≤ 10 , with 96,457 individuals), and “far” (genetic distance > 10 , with 32,822 individuals). For each index SNP of each PGS, we computed the allelic effect estimate relative to the effect estimate in the original GWAS sample. Shown are means \pm standard deviations across PGS index SNPs for three traits, highlighting the poorer agreement between allelic effect estimates for lymphocyte count. **B.** We compared the mean heterozygosity of index SNPs for height, triglycerides, and lymphocyte count. For each trait, SNPs are stratified into three equally-sized bins of squared allelic effect estimate (**Fig. S23**). Each data point is the mean heterozygosity of a stratum in a bin of genetic distance. Unlike other traits, the heterozygosity of large effect variants for lymphocyte count increases with genetic distance from the GWAS sample. See **Fig. S22** for other traits. **C.** We compared the variance of PGS, in each bin, relative to the variance of PGS in the reference group, across traits. Among the 15 traits we have examined, only for lymphocyte count, monocyte count, and white blood cell count (WBC) the PGS variance increased with genetic distance. Green points show the PGS variance for lymphocyte count in genetic ancestry bins. Lines show the ordinary least squares linear fit to the respective bin-level data for each trait.

164 accuracy. For triglyceride levels, lymphocyte count, and white blood cell count, group-level
 165 prediction accuracy is near zero far from the GWAS sample (**Figs. 2E,S3**) whereas at the
 166 individual level, prediction accuracy increases (**Figs. 2F,S3,S11,S13**).

167 Discussion

168 Through an examination of empirical trends of portability at the individual level, we high-
169 lighted three gaps in our current understanding of the portability problem. Below, we discuss
170 possible avenues towards filling these gaps.

171 The driver of portability that has been extensively discussed in the literature is ancestral
172 similarity to the GWAS sample^{17,27,40,3,15,7}. Yet our results show that, at the individual
173 level, prediction accuracy is poorly predicted by genome-wide genetic ancestry. We note
174 that our measure of genetic distance (also similar to that used in other studies^{27,7,10,9}) is
175 plausibly sub-optimal, as suggested, for example, by the noisiness of its relationship with F_{st}
176 at intermediate genetic distances (**Fig. 2B**). Therefore, one path forward is to ask whether
177 refined measures of genetic distance from the GWAS sample, in particular ones that capture
178 local ancestry^{11,31} (e.g., in the genomic regions containing the PGS index SNPs), better
179 explain portability. Another direction is in quantifying how environmental and social context,
180 such as access to healthcare, affect portability (See [10] for a recent method in this vein).
181 The relative importance of these factors will also inform the efforts to diversify participation
182 in GWAS.

183 Second, we observed some trait-specific trends in portability, and hypothesize that they
184 reflect the specifics of natural selection and evolutionary history of genetic variants affecting
185 the trait. While previous work considered the impact of directional^{3,8,5,23} and stabilizing se-
186 lection^{44,40,23} on portability, the trait-specific (and PGS-specific) impact—notably for disease
187 prediction—is yet to be studied empirically. Evolutionary perspectives on genetic architec-
188 tures and other facets of GWAS data have been transformative³³. This may also prove to
189 be the case for understanding PGS portability.

190 Third, we show that individual-level measures, which are arguably the most relevant to
191 eventual applications of PGS, can yield different results to group level measures that are
192 widely used. PGS research has been focused on coefficient of determination (R^2) analyzed
193 at the group level^{17,19,40,7,27,44,3}. More generally, different applications and questions call

194 for different measures of prediction accuracy, for instance when considering the utility of
195 a public health intervention applied to communities, as opposed to asking about the cost-
196 effectiveness of an expensive drug for an individual patient (see [1] for related discussion).
197 Therefore, future research of predictive performance could benefit from more focus on the
198 metrics most relevant to the intended application.

199 Addressing these gaps in our understanding of PGS portability will be key for evaluating
200 the utility of a PGS, and for its equitable application in the clinic and beyond.

201 Acknowledgements

202 We thank Ipsita Agarwal, Maryn Carlson, Yi Ding, Doc Edge, Kangcheng Hou, Hakhamanesh
203 Mostafavi, Bogdan Pasaniuc, Molly Przeworski, Sam Smith, Jeff Spence and members of
204 the Harpak Lab for helpful feedback. This work was funded by NIH grant R35GM151108,
205 a fellowship from the Simons Foundation's Society of Fellows (#633313) and a Pew Schol-
206 arship to A.H. This study was conducted using the UK Biobank resource under application
207 61666, as approved by the University of Texas at Austin institutional review board (protocol
208 2019-02-0125). We acknowledge the Texas Advanced Computing Center (TACC) at The
209 University of Texas at Austin for providing computational resources that have contributed
210 to the research results reported within this paper.

211 Methods

212 Data

213 **Data overview.** All analyses were conducted with data from the UK Biobank, a large-scale
214 biomedical database with a sample size of 502,490 individuals³⁶. In this study, we considered
215 479,406 individuals who passed quality control (QC) checks, which included the removal of
216 651 samples identified by the UK Biobank as having sex chromosome aneuploidy (data field
217 22019), and an additional 14,433 individuals whose self-reported biological sex (data field 31)

218 differed from sex determined from that implied by their sex chromosome karyotype (data
219 field 22001). We removed 963 individuals who are outliers in heterozygosity or genotype
220 missingness (data field 22027) and 6,854 individuals with genotype missingness greater than
221 2% (data field 22005). To prevent biased estimations of the effect sizes of SNPs, we excluded
222 183 individuals with 10 or more 3rd-degree relatives (data field 22021).

223 **Genotype data.** We started with 765,067 biallelic variants out of a total of 784,256
224 genotyped variants on the autosomes. We first removed 10,543 SNPs within the major
225 histocompatibility complex (MHC) and extended region in strong LD with it (chromosome
226 6, positions 26,477,797-35,448,354 in the GRCh37 genome build). We excluded variants with
227 a Hardy-Weinberg equilibrium p-value (`--hwe`) lower than 1×10^{-10} among White British
228 (WB) individuals (see in Section **GWAS** below), removing another 46,854 variants. We
229 also removed an additional 39,939 variants by setting the minor allele frequency threshold
230 (`--maf`) among WB to $> 0.01\%$. After filtering, we had 667,731 variants which we analyzed
231 going forward.

232 **Phenotype data.** We analyzed 15 highly heritable traits, as determined based on
233 Neale Lab SNP heritability estimates²¹ (**Table S1**). These included both physiological
234 measurements to biomarkers: standing height (data field 50), cystatin C level (data field
235 30720), platelet count (data field 30080), mean corpuscular volume (MCV, data field 30040),
236 weight (data field 21002), mean corpuscular hemoglobin (MCH, data field 30050), body
237 mass index (BMI, data field 21001), red blood cell count (RBC, data field 30010), body fat
238 percentage (data field 23099), monocyte count (UKB data field 30130), triglyceride level
239 (data field 30870), lymphocyte count (data field 30120), white blood cell count (WBC, data
240 field 30000), eosinophil count (data field 30150), and LDL cholesterol level (data field 30780)
241 (**Table S1**). For all analyses, we removed individuals with missing trait data.

242 Genetic distance calculations

243 The fixation index (F_{st}) is a natural metric, a single number, to measure the divergence be-
244 tween two sets of chromosomes and we considered using it to measure the distance between

245 the pair of chromosomes of an individual and chromosomes in the GWAS sample. However,
246 calculating F_{st} was computationally costly. Since previous work²⁷ showed it is tightly corre-
247 lated with Euclidean distance in the PC space in the UKB, we used Euclidean distance as a
248 single number proxying genetic distance from the GWAS sample. We used the pre-computed
249 PCA provided by the UK Biobank (data field 22009). To calculate individual-level scores on
250 each PC, we used the genotype matrix of the full post-filtering sample of individuals (data
251 field 22009).

252 The genetic distance is the weighted PC distance between an individual coordinates
253 vector in PCA subspace of the first K PCs, x , and the centroid of M individuals $\{x^m\}_{m=1}^M$
254 in the GWAS sample, $C = \frac{\sum_{m=1}^M x^m}{M}$, is

$$\sqrt{\sum_{k=1}^K w_k (x_k - c_k)^2}$$

with weights

$$w_k = \frac{\lambda_k}{\sum_{n=1}^{40} \lambda_n},$$

255 where λ_k is the k 'th eigenvalue.

256 To identify K , the number of PCs we used and to confirm the approximation is reasonable
257 for our data, we examined the correlation of genetic distance with F_{st} as a function of K on
258 a small subset of the prediction sample.

259 We randomly selected 10,000 prediction sample individuals with a weighted PC distance
260 greater than the weighted PC distance of 95% of the GWAS set (based on weighted PC
261 distance calculated from the $K = 10$). For those individuals, we estimated their F_{st} and
262 weighted PC distance to the GWAS centroid for $K \in 1, \dots, 40$. We estimated F_{st} in this
263 subsample with the Weir and Cockerham method⁴¹ using the `--fst` flag in *PLINK 1.9*^{28,4}

264 Since the PC distance calculated from using $K = 40$ correlated most strongly with F_{st}
265 ($r = 0.98$) (**Fig. S1**), we used this number of PC to estimate the genetic distance for
266 all test individuals (**Fig. 1B-C**). We note that genetic distance is less reflective of F_{st} for
267 intermediate genetic distances (**Fig. 1B**).

268 We divided the raw genetic distances by the (raw) mean genetic distance among GWAS
269 sample individuals. To gain intuition about these standardized units of genetic distance, we
270 wished to estimate where on this scale we would find individuals from three subsamples from
271 the 1000 Genomes Phase 3 dataset²: CEU, Utah residents (CEPH) from primarily Northern
272 and Western European descent; CHB, Han Chinese in Beijing, China; and YRI, Yoruba in
273 Ibadan, Nigeria. To this end, we ran a PCA with a dataset that includes both the UKB
274 individuals and the CEU, CHB, and YRI individuals. We identified the UKB individuals
275 with the shortest weighted Euclidean distance to the centroid of each of the three 1000
276 Genomes populations, and used the genetic distance of those three UKB individuals in our
277 PCA of only UKB individuals as a proxy of where the three 1000 Genomes subsamples fall
278 on the scale of our genetic distance measurement (**Fig. 1C**).

279 The distribution of genetic distance is heavily right-skewed, with most individuals falling
280 close to the GWAS centroid. Since we wanted to focus on the individuals far away from the
281 GWAS set, we only analyzed data for individuals with a genetic distance greater than the
282 95th percentile of genetic distance from among GWAS sample individuals (**Fig. 1C**), with
283 the exception of the analyses behind **Fig. 3** and **Fig. S14-S21**.

284 For group level analyses, we binned the prediction samples by genetic distance using 500
285 equally-sized bins, with 258-259 individuals per bin.

286 **PGS and evaluating PGS prediction accuracy**

287 **GWAS.** In the selection of the GWAS sample, we used the WB classification as provided
288 by the UKB. This classification includes two criteria: an individual must self-identified as
289 White British (data field 21000) and Caucasian (data field 22006). All other individuals
290 are “Non-White British” (NWB). We randomly selected 350,000 WB as the GWAS sample.
291 We considered the remaining 52,281 WB all (77,125, after filtering) NWB as the prediction
292 sample. Next, for each trait, we used the `--glm` flag from *PLINK 2.0*^{29,4} to run GWAS on the
293 GWAS set. We used the following covariates: the first 20 PCs from UKB (data field 22009,
294 age (data field 21022), age², sex (data field 31), age*sex, and age²*sex, where the asterisk

295 (*) denotes the product of two variables, referring to an interaction term. We clumped the
296 SNPs with the `--clump` flag from *PLINK 1.9*^{28,4}, setting the association p-value threshold
297 for clumping to 0.01, LD r^2 threshold to 0.2, and window size to 250 kb.

298 **PGS construction.** After clumping and thresholding the SNPs with marginal asso-
299 ciation $p < 1 \times 10^{-5}$, we calculated PGS for each individual for every phenotype. The
300 calculations were carried out with the `--score` flag in *PLINK 2.0*^{29,4}.

301 **PGS prediction accuracy at the group level.** To evaluate prediction accuracy at
302 the group level, we linearly regressed the phenotype on the covariates (array type (data
303 field 22000), age, age², sex, age*sex, and age²*sex) and PGS within each genetic distance
304 bin (phenotype \sim covariates + PGS), which is the full model. We then performed another
305 linear regression of the phenotype on the covariates, excluding the PGS, within each bin
306 (phenotype \sim covariates), which is the reduced model. Using these two squared correlations,
307 we calculated partial R^2 for the PGS with the sum of squared errors (SSE) of these two
308 models as

$$R_{\text{partial}}^2 = \frac{SSE(\text{reduced_model}) - SSE(\text{full_model})}{SSE(\text{reduced_model})},$$

309 which represents the prediction accuracy of PGS for each bin.

310 As a baseline prediction accuracy, we identified the 50 bins (of 269 individuals each)
311 with the median genetic distance most similar to the mean genetic distance for GWAS
312 individuals; This reference group represents individuals from the prediction set that are
313 most similar to “typical” GWAS individuals in terms of genetic distance. The mean PGS
314 prediction accuracy across these 50 bins served as the baseline value. Throughout the paper,
315 we report the prediction accuracy at the group level as a bin’s squared partial correlation
316 between the PGS and the trait divided by this baseline value.

317 **Prediction error at the individual level.** For the individual-level prediction error, we
318 first derived phenotypic values adjusted for covariates Z in two steps, involving residualizing
319 some covariates in each genetic ancestry bin independently and some covariates globally.

320 First, we regress raw phenotype values Y independently in each bin on covariates,

$$Y \sim \text{array type} + \text{age}^2 + \text{sex} + \text{age} * \text{sex} + \text{age}^2 * \text{sex}.$$

321 In bins in which only a single individual was genotyped with a particular array type, we did
322 not include array type as a covariate. We then regress the residual X (where $X = Y - \hat{Y}$
323 and \hat{Y} is the fitted value from the first step) globally on covariates,

$$X \sim \text{genetic_distance_polynomial} + \text{sex} + \text{sex} * \text{genetic_distance},$$

324 where genetic_distance_polynomial is a 20-degree polynomial in genetic distance. Finally,
325 we regress the residual of this second regression, $Z = X - \hat{X}$ onto the PGS in a simple
326 (univariate) linear regression. We refer to the squared residual of this regression,

$$\left(Z - \hat{Z} \right)^2,$$

327 as the unstandardized squared prediction error. Similar to the group-level analysis, we
328 computed the the mean unstandardized squared prediction error in the 50 reference bins
329 as a baseline values (**Table S1** details the baseline values across traits). The squared
330 prediction error, the measure of individual-level prediction accuracy we refer to throughout,
331 is the unstandardized squared prediction error divided by the baseline value.

332 **Spline fits.** For both the individual-level and group-level analysis, **Fig. 2, Fig. S2-S9**
333 show cubic spline fits. We fitted these splines using 8 knots. The knot positions were chosen
334 based on the density of the individual genetic distances, such that there is an equal number
335 of samples between any two knots. This resulted in knots at genetic distances of 1.91, 2.25,
336 3.12, 5.02, 9.39, 18.74, 43.11, 61.96, and 160.82.

337 **Mean trends in individual-level prediction accuracy.** In the **Results** section
338 of the main text, we discuss various individual predictors of squared prediction error. In
339 addition to genetic distance, we considered 2 measures of individual-level prediction error:
340 The Townsend Deprivation Index (data field 189) and average yearly total household income

341 before tax (data field 738). For genetic distance and the Townsend Deprivation Index, we
342 considered five uniformly-spaced bins, and computed the mean squared prediction error and
343 the standard error of this mean (**Fig. 3A, Fig. S14-S17**). For household income, which
344 the UK Biobank provides as categorical data conferring to ranges in British Pounds, we
345 converted the categories into an ordinal variable coded as 1,2,3,4 and 5, and computed the
346 mean squared prediction error, and standard error of the mean, in each. This also allowed us
347 to use the income categories directly as measures in the regression models used for comparison
348 of variance in prediction error explained that we discuss below.

349 **Comparison of variance in prediction accuracy explained across measures.**
350 For this analysis, we used all the individuals in the prediction set and did not filter for
351 the individuals with a genetic distance greater than 95th percentile of genetic distance from
352 among GWAS sample individuals. We compared the variance in squared prediction error
353 explained for 8 raw measures: genetic distance, Townsend Deprivation Index (data field 189),
354 average yearly total household income before tax (data field 738), educational attainment
355 (data field 6138), which we converted into years of education, minor allele counts for SNPs
356 with different with different magnitudes of effects (three equally-sized bins of small, medium,
357 and large squared effect sizes, see **Fig. S23**), and minor allele counts of all SNPs. “Minor”
358 here is with respect to the GWAS sample, and the count is the total sum of minor alleles
359 across index SNPs of the magnitude category. Namely, for each measure, we independently
360 fit three different models:

- 361 – A linear predictor, fit using Ordinary Least Squares (OLS).
- 362 – A discretized predictor, using one predicted value per each of the 5 bins where all five
363 bins had identical widths.
- 364 – A cubic spline. 16 knots were placed based on the density of data points, such that
365 there was an equal number of data points between each pair of consecutive knots.

366 After fitting the models, we calculated the R^2 values to determine the variance explained by
367 each measure-method combination. We then computed 95% central confidence intervals for

368 these R^2 values to assess the reliability of the estimates (**Fig. 3B, Fig. S18-S21**).

369 **Additional analyses on lymphocyte count**

370 **Comparing allelic estimates across three GWASs.** To test whether allelic effect es-
371 timates are similar across genetic ancestry, we performed two additional GWASs for each
372 trait in two subsets of the prediction sample: “close” group (genetic distance ≤ 10 , with
373 96,457 individuals) and “far” (genetic distance > 10 , with 32,822 individuals). For both
374 groups, we adjusted for 20 PCs of the genotype matrix of the respective set of individuals,
375 using the `--pca approx 20` flag in *PLINK 2.0*^{29,4}. After running GWAS independently
376 in the two groups, for each index SNP of the original PGS, we divided allelic effect esti-
377 mates in the original GWAS / close / far set by the allelic effect estimate in the original
378 GWAS. **Fig. 4A** shows the mean \pm standard deviation across PGS index SNPs for each of
379 three traits, highlighting the poorer agreement of the allelic effect estimates for lymphocyte
380 count.

381 **Heterozygosity at index SNPs as a function of genetic distance.** For each PGS,
382 we calculated the heterozygosity of each index SNP in each bin from allele counts using the
383 `--freq` flag from *PLINK 1.9*^{28,4}. We stratified index SNPs into three equally-sized bins
384 based on their squared effect sizes (**Fig. S23**). **Figs. 4B, S22** show the mean heterozygosity
385 (across stratum SNPs) for each stratum of in each genetic distance bin.

386 **Variance of PGS as a function of genetic distance.** For each phenotype, we
387 calculated the variance of PGS in each bin relative to the mean of the variance of PGS in
388 the 50 bins close to the GWAS set. In **Figs. 4C**, we plotted the values in each bin as well as
389 a linear fit for lymphocyte count. For other traits, we only plotted the linear fit.

390 **Heritability associated with each index SNP.** We estimated the heritability ex-
391 plained by each index SNP as

$$\hat{h}_{index}^2 = 2p(1-p)\hat{\beta}^2,$$

392 where $\hat{\beta}$ is the estimated allelic effect and p is the allele frequency. In **Fig. S24**, we compared
393 the distribution of index SNP heritability across traits and with allelic effect estimates and

394 heterozygosities calculated both in the original GWAS sample, the “close” prediction sample
395 and the “far” prediction sample. For each trait, the SNPs used are also stratified into three
396 equal-sized strata (small, medium, and large) based on their squared effect sizes, as discussed
397 above.

398 **Code availability**

399 The scripts for the analyses and figures are available at [https://github.com/harpak-lab/](https://github.com/harpak-lab/Portability_Questions)
400 `Portability_Questions`.

401 References

- 402 [1] Abramowitz, S. A., Boulier, K., Keat, K., Cardone, K. M., Shivakumar, M., *et al.*,
403 2024. Population Performance and Individual Agreement of Coronary Artery Disease
404 Polygenic Risk Scores. *medRxiv*, pages 2024–07.
- 405 [2] Auton, A., Abecasis, G. R., Altshuler, D. M., Durbin, R. M., Abecasis, G. R., *et al.*, 10
406 2015. A global reference for human genetic variation. *Nature*, 526:68–74.
- 407 [3] Carlson, M. O., Rice, D. P., Berg, J. J., and Steinrücken, M., 5 2022. Polygenic score
408 accuracy in ancient samples: Quantifying the effects of allelic turnover. *PLOS Genetics*,
409 18:e1010170.
- 410 [4] Chang, C. C., Chow, C. C., Tellier, L. C., Vattikuti, S., Purcell, S. M., *et al.*, 02
411 2015. Second-generation PLINK: rising to the challenge of larger and richer datasets.
412 *GigaScience*, 4(1):s13742–015–0047–8.
- 413 [5] Cox, S. L., Moots, H. M., Stock, J. T., Shbat, A., Bitarello, B. D., *et al.*, 1 2022. Predict-
414 ing skeletal stature using ancient DNA. *American Journal of Biological Anthropology*,
415 177:162–174.
- 416 [6] Ding, Y., Hou, K., Burch, K. S., Lapinska, S., Privé, F., *et al.*, 1 2022. Large uncertainty
417 in individual polygenic risk score estimation impacts PRS-based risk stratification. *Nature*
418 *Genetics*, 54:30–39.
- 419 [7] Ding, Y., Hou, K., Xu, Z., Pimplaskar, A., Petter, E., *et al.*, 6 2023. Polygenic scoring
420 accuracy varies across the genetic ancestry continuum. *Nature*, 618:774–781.
- 421 [8] Durvasula, A. and Lohmueller, K. E., 4 2021. Negative selection on complex traits limits
422 phenotype prediction accuracy between populations. *The American Journal of Human*
423 *Genetics*, 108:620–631.

- 424 [9] Habtewold, T. D., Wijesiriwardhana, P., Biedrzycki, R. J., and Tekola-Ayele, F., 7 2024.
425 Genetic distance and ancestry proportion modify the association between maternal ge-
426 netic risk score of type 2 diabetes and fetal growth. *Human Genomics*, 18:81.
- 427 [10] Hou, K., Xu, Z., Ding, Y., Mandla, R., Shi, Z., *et al.*, 7 2024. Calibrated prediction
428 intervals for polygenic scores across diverse contexts. *Nature Genetics*, 56:1386–1396.
- 429 [11] Hu, S., Ferreira, L. A. F., Shi, S., Hellenthal, G., Marchini, J., *et al.*, 2023. Leveraging
430 fine-scale population structure reveals conservation in genetic effect sizes between human
431 populations across a range of human phenotypes. *bioRxiv*.
- 432 [12] Kamiza, A. B., Toure, S. M., Vujkovic, M., Machipisa, T., Soremekun, O. S., *et al.*, 6
433 2022. Transferability of genetic risk scores in African populations. *Nature Medicine*,
434 28:1163–1166.
- 435 [13] Khera, A. V., Chaffin, M., Aragam, K. G., Haas, M. E., Roselli, C., *et al.*, 9 2018.
436 Genome-wide polygenic scores for common diseases identify individuals with risk equiv-
437 alent to monogenic mutations. *Nature Genetics*, 50:1219–1224.
- 438 [14] Kullo, I. J., 8 2024. Promoting equity in polygenic risk assessment through global
439 collaboration. *Nature Genetics*.
- 440 [15] Lewis, A. C. F., Perez, E. F., Prince, A. E. R., Flaxman, H. R., Gomez, L., *et al.*,
441 10 2022. Patient and provider perspectives on polygenic risk scores: implications for
442 clinical reporting and utilization. *Genome Medicine*, 14:114.
- 443 [16] Lewontin, R. C., 5 1974. Annotation: the analysis of variance and the analysis of causes.
444 *American journal of human genetics*, 26:400–11.
- 445 [17] Martin, A. R., Kanai, M., Kamatani, Y., Okada, Y., Neale, B. M., *et al.*, 4 2019. Clinical
446 use of current polygenic risk scores may exacerbate health disparities. *Nature Genetics*,
447 51:584–591.

- 448 [18] Martin, G. and Lenormand, T., 12 2006. The fitness effect of mutations across environ-
449 ments: a survey in light of fitness landscape models. *Evolution; international journal*
450 *of organic evolution*, 60:2413–27.
- 451 [19] Mostafavi, H., Harpak, A., Agarwal, I., Conley, D., Pritchard, J. K., *et al.*, 1 2020.
452 Variable prediction accuracy of polygenic scores within an ancestry group. *eLife*, 9.
- 453 [20] Nagpal, S. and Gibson, G., 2024. Dual exposure-by-polygenic score interactions high-
454 light disparities across social groups in the proportion needed to benefit. *medRxiv*, pages
455 2024–07.
- 456 [21] Neale Lab, 10. UK Biobank. URL <http://www.nealelab.is/uk-biobank>.
- 457 [22] O’Connor, L. J., Schoech, A. P., Hormozdiari, F., Gazal, S., Patterson, N., *et al.*, 9
458 2019. Extreme Polygenicity of Complex Traits Is Explained by Negative Selection. *The*
459 *American Journal of Human Genetics*, 105:456–476.
- 460 [23] Patel, R. A., Weiß, C. L., Zhu, H., Mostafavi, H., Simons, Y. B., *et al.*, 2024. Conditional
461 frequency spectra as a tool for studying selection on complex traits in biobanks. *bioRxiv*.
- 462 [24] Peter, B. M., 6 2022. A geometric relationship of F_2 , F_3 and F_4 -statistics with prin-
463 cipal component analysis. *Philosophical Transactions of the Royal Society B: Biological*
464 *Sciences*, 377.
- 465 [25] Polygenic Risk Score Task Force of the International Common Disease Alliance,
466 Adeyemo, A., Balaconis, M. K., Darnes, D. R., Fatumo, S., *et al.*, 11 2021. Responsible
467 use of polygenic risk scores in the clinic: potential benefits, risks and gaps. *Nature*
468 *Medicine*, 27:1876–1884.
- 469 [26] Pritchard, J. K. and Przeworski, M., 7 2001. Linkage Disequilibrium in Humans: Models
470 and Data. *The American Journal of Human Genetics*, 69:1–14.

- 471 [27] Privé, F., Aschard, H., Carmi, S., Folkersen, L., Hoggart, C., *et al.*, 1 2022. Portability
472 of 245 polygenic scores when derived from the UK Biobank and applied to 9 ancestry
473 groups from the same cohort. *The American Journal of Human Genetics*, 109:12–23.
- 474 [28] Purcell, S. and Chang, C. PLINK 1.9. URL www.cog-genomics.org/plink/1.9.
- 475 [29] Purcell, S. and Chang, C. PLINK 2.0. URL www.cog-genomics.org/plink/2.0.
- 476 [30] Ragsdale, A. P., Nelson, D., Gravel, S., and Kelleher, J., 10 2020. Lessons Learned
477 from Bugs in Models of Human History. *The American Journal of Human Genetics*,
478 107:583–588.
- 479 [31] Ried, T., 9 1998. Chromosome painting: a useful art. *Human Molecular Genetics*,
480 7:1619–1626.
- 481 [32] Saitou, M., Dahl, A., Wang, Q., and Liu, X., 2022. Allele frequency differences of causal
482 variants have a major impact on low cross-ancestry portability of PRS. *medRxiv*.
- 483 [33] Sella, G. and Barton, N. H., 8 2019. Thinking About the Evolution of Complex Traits in
484 the Era of Genome-Wide Association Studies. *Annual Review of Genomics and Human*
485 *Genetics*, 20:461–493.
- 486 [34] Shalizi, C. R., 2024. Advanced Data Analysis from an Elementary Point of View. URL
487 www.stat.cmu.edu/~cshalizi/ADAfaEPoV.
- 488 [35] Simons, Y. B., Bullaughey, K., Hudson, R. R., and Sella, G., 3 2018. A population
489 genetic interpretation of GWAS findings for human quantitative traits. *PLOS Biology*,
490 16:e2002985.
- 491 [36] Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., *et al.*, 3 2015. UK Biobank:
492 An Open Access Resource for Identifying the Causes of a Wide Range of Complex
493 Diseases of Middle and Old Age. *PLOS Medicine*, 12:e1001779.
- 494 [37] Townsend, P., 4 1987. Deprivation. *Journal of Social Policy*, 16:125–146.

- 495 [38] Tsuo, K., Shi, Z., Ge, T., Mandla, R., Hou, K., *et al.*, 2024. All of Us diversity
496 and scale improve polygenic prediction contextually with greatest improvements for
497 underrepresented populations. *bioRxiv*.
- 498 [39] Tukey, J. W., 2 1969. Analyzing data: Sanctification or detective work? *American*
499 *Psychologist*, 24:83–91.
- 500 [40] Wang, Y., Guo, J., Ni, G., Yang, J., Visscher, P. M., *et al.*, 7 2020. Theoretical
501 and empirical quantification of the accuracy of polygenic scores in ancestry divergent
502 populations. *Nature Communications*, 11:3865.
- 503 [41] Weir, B. S. and Cockerham, C. C., 11 1984. Estimating F-Statistics for the Analysis of
504 Population Structure. *Evolution*, 38:1358.
- 505 [42] Weissbrod, O., Kanai, M., Shi, H., Gazal, S., Peyrot, W. J., *et al.*, 4 2022. Leveraging
506 fine-mapping and multipopulation training data to improve cross-population polygenic
507 risk scores. *Nature Genetics*, 54:450–458.
- 508 [43] Westerman, K. E. and Sofer, T., 4 2024. Many roads to a gene-environment interaction.
509 *The American Journal of Human Genetics*, 111:626–635.
- 510 [44] Yair, S. and Coop, G., 6 2022. Population differentiation of polygenic score predictions
511 under stabilizing selection. *Philosophical Transactions of the Royal Society B: Biological*
512 *Sciences*, 377.

513 Supplementary Materials for: Three Open Questions in Polygenic
514 Score Portability

515 Joyce Y. Wang¹, Neeka Lin¹, Michael Zietz², Jason Mares³, Vagheesh M. Narasimhan^{1,4},
516 Paul J. Rathouz^{4,5} and Arbel Harpak^{1,5,+}

517 ¹ Department of Integrative Biology, The University of Texas at Austin, Austin, TX

518 ² Department of Biomedical Informatics, Columbia University, New York, NY

519 ³ Department of Neurology, Columbia University, New York, NY

520 ⁴ Department of Statistics and Data Science, The University of Texas at Austin, Austin, TX

521 ⁵ Department of Population Health, The University of Texas at Austin, Austin, TX

522 + Correspondence should be addressed to A.H. (arbelharpak@utexas.edu)

523 **Contents**

524 **List of Tables**

525 S1 Characteristics of traits and PGS analyzed. 29

526 **List of Figures**

527 S1 Correlation between PC distance and F_{st} 30

528 S2 Trends of portability vary across traits and measures, for anthropometric mea-
529 surements. 31

530 S3 Trends of portability vary across traits and measures, for white blood-cell
531 related traits. 32

532 S4 Trends of portability vary across traits and measures, for red blood cell-related
533 traits. 33

534 S5 Trends of portability vary across traits and measures, for other biomarkers. . 34

535 S6 Full range of individual level prediction error for anthropometric measurements. 35

536 S7 Full range of individual level prediction error for white blood cell-related traits. 36

537 S8 Full range of individual level prediction error for red blood cell-related traits. 37

538 S9 Full range of individual level prediction error for other biomarkers. 38

539 S10 Divergence of group- and individual-level prediction accuracy for anthropo-
540 metric measurements. 39

541 S11 Divergence of group- and individual-level prediction accuracy for white blood
542 cell-related traits. 40

543 S12 Divergence of group- and individual-level prediction accuracy for red blood
544 cell-related traits. 41

545	S13	Divergence of group- and individual-level prediction accuracy for other biomark-	
546		ers.	42
547	S14	Mean trends in individual-level prediction accuracy by different measures for	
548		anthropometric measurements.	43
549	S15	Mean trends in individual-level prediction accuracy by different measures for	
550		white blood cell-related traits.	44
551	S16	Mean trends in individual-level prediction accuracy by different measures for	
552		red blood cell-related traits.	45
553	S17	Mean trends in individual-level prediction accuracy by different measures for	
554		other biomarkers.	46
555	S18	Individual-level prediction error explained by different measures for anthro-	
556		pometric measurements.	47
557	S19	Individual-level prediction error explained by different measures for white	
558		blood cell-related traits.	48
559	S20	Individual-level prediction error explained by different measures for red blood	
560		cell-related traits.	49
561	S21	Individual-level prediction error explained by different measures for other	
562		biomarkers.	50
563	S22	Mean heterozygosity of SNPs, stratified by effect size.	51
564	S23	Squared allelic effect estimate of SNPs.	52
565	S24	Heritability explained by index SNPs.	53

Phenotype	SNP h^2	Mean squared prediction error (MSE)	Variance of residualized phenotype ($Var[Z]$)	$1 - \frac{MSE}{Var[Z]}$	Prediction accuracy (partial R^2)
Height	0.4852	30.2588	39.5079	0.2341	0.2417
Cystatin C level	0.3214	0.0241	0.0255	0.0533	0.0686
Platelet count	0.3079	2828.9020	3290.2100	0.1402	0.1527
Mean corpuscular volume	0.2667	15.8315	18.2524	0.1326	0.1380
Weight	0.2654	182.9963	195.4927	0.0639	0.0654
Mean corpuscular hemoglobin	0.2530	2.5789	2.9190	0.1165	0.1269
BMI	0.2482	21.1426	22.1774	0.0467	0.0482
Red blood cell count	0.2337	0.1066	0.1177	0.0941	0.0959
Body fat percentage	0.0472	37.6601	39.4749	0.0460	0.0585
Monocyte count	0.2305	0.0521	0.0539	0.0335	0.0928
Triglyceride level	0.2182	0.8744	0.9382	0.0680	0.0724
Lymphocyte count	0.2103	1.4708	1.4775	0.0045	0.0264
White blood cell count	0.1910	3.9834	4.1347	0.0366	0.0563
Eosinophil count	0.1840	0.0163	0.0172	0.0517	0.0561
LDL cholesterol level	0.0825	0.6659	0.7201	0.0753	0.0781

Table S1: Characteristics of traits and PGS analyzed. SNP heritabilities (SNP h^2) are taken from the Neale Lab's UKB analysis²¹. For the 50 bins with a genetic distance most similar to the mean genetic distance of the GWAS group, we calculated the group-level prediction accuracy (partial genetic correlation of the PGS and the trait value), mean individual prediction error (squared prediction error), the variance of residualized phenotype, and the ratio between the two subtracted from one, as another measure of phenotype variance explained by the PGS close to the GWAS sample.

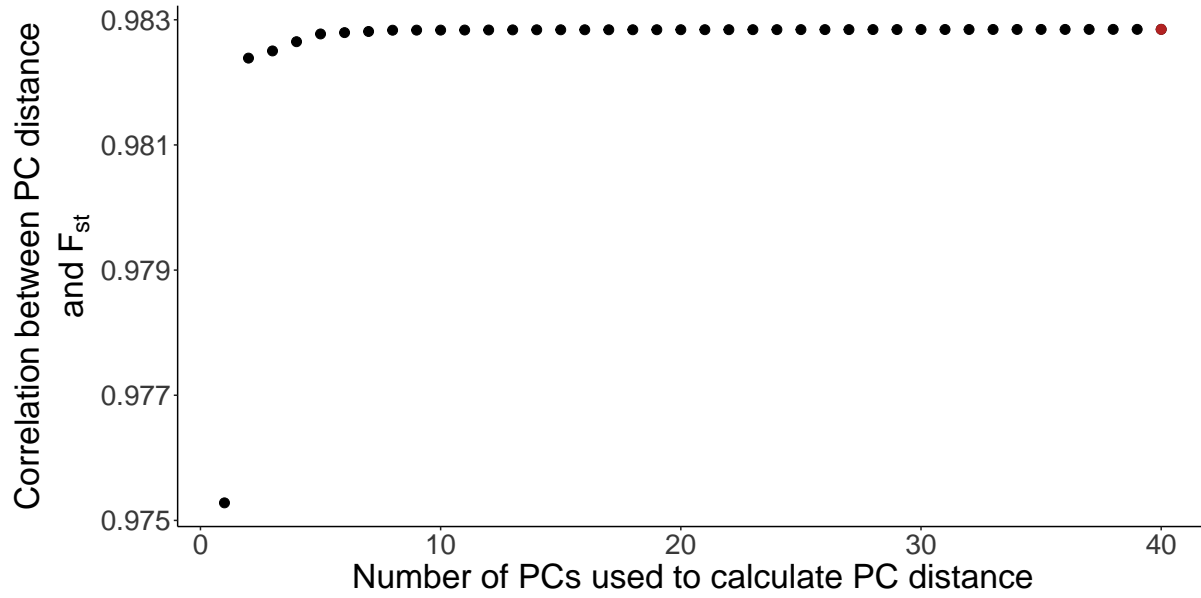


Figure S1: Correlation between PC distance and F_{st} . This figure presents the correlation between PC distance and F_{st} , calculated using different numbers of UKB PCs, from 1 to 40. Using 40 PCs produces the highest correlation between PC distance and F_{st} (red dot, $r = 0.98$).

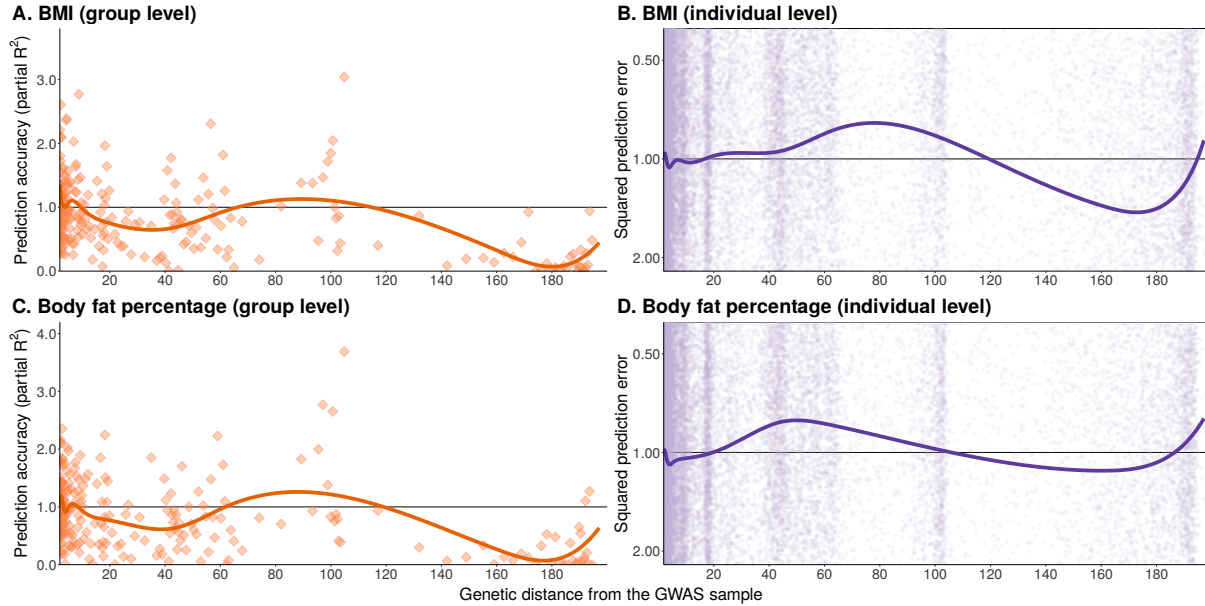


Figure S2: Trends of portability vary across traits and measures, for anthropometric measurements. This figure presents the same analysis as **Fig. 2** in the main text, but for other anthropometric traits. At the group level (left panels), we measured prediction accuracy with the squared partial correlation between the PGS and the trait value in 500 bins of 258-259 individuals each. At the individual level (right panels), we measured the squared prediction error. Curves show cubic spline fits, with 8 knots placed based on the density of data points.

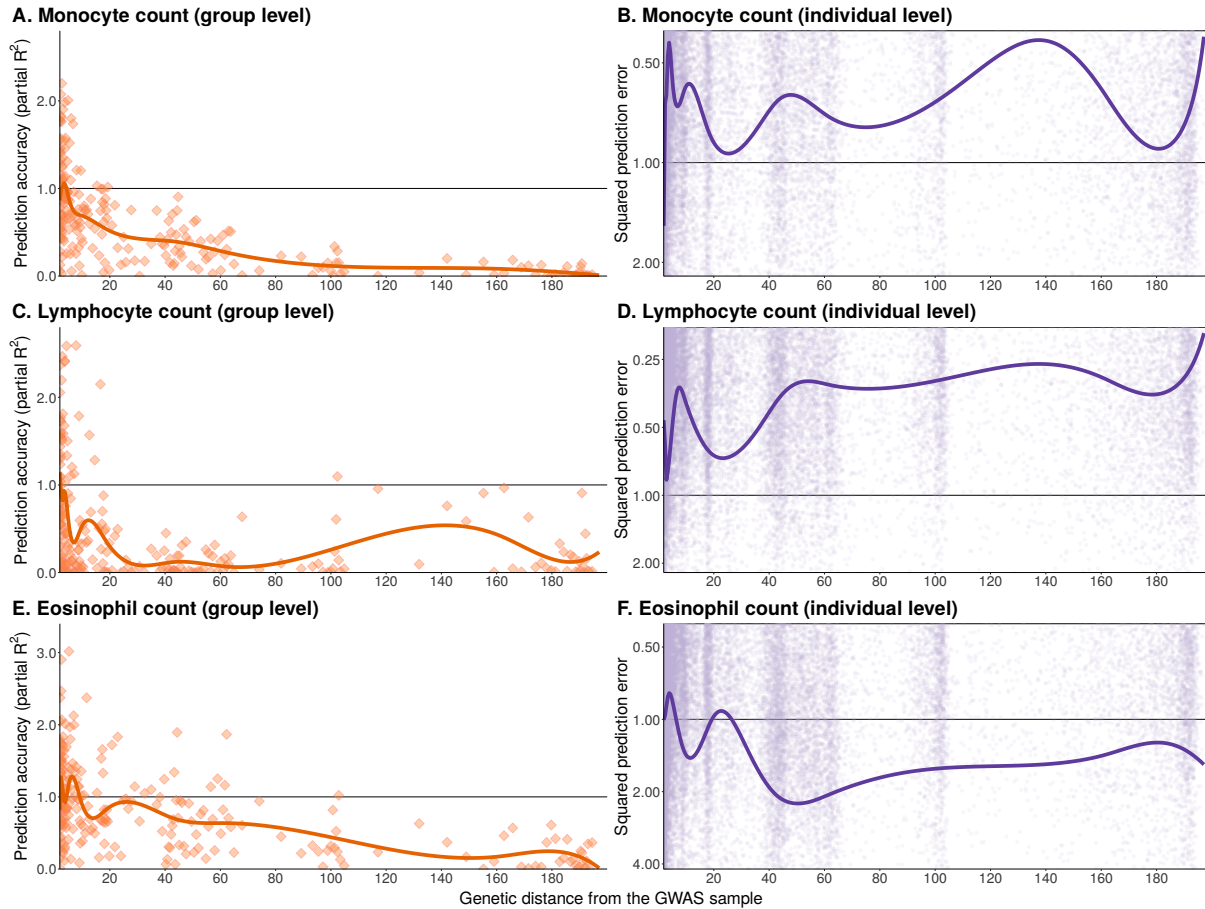


Figure S3: Trends of portability vary across traits and measures, for white blood-cell related traits. This figure presents the same analysis as **Fig. 2** in the main text, but shows other white blood cell-related traits. At the group level (left panels), we measured prediction accuracy with the squared partial correlation between the PGS and the trait value in 500 bins of 258-259 individuals each. At the individual level (right panels), we measured the squared prediction error. Curves show cubic spline fits, with 8 knots placed based on the density of data points.

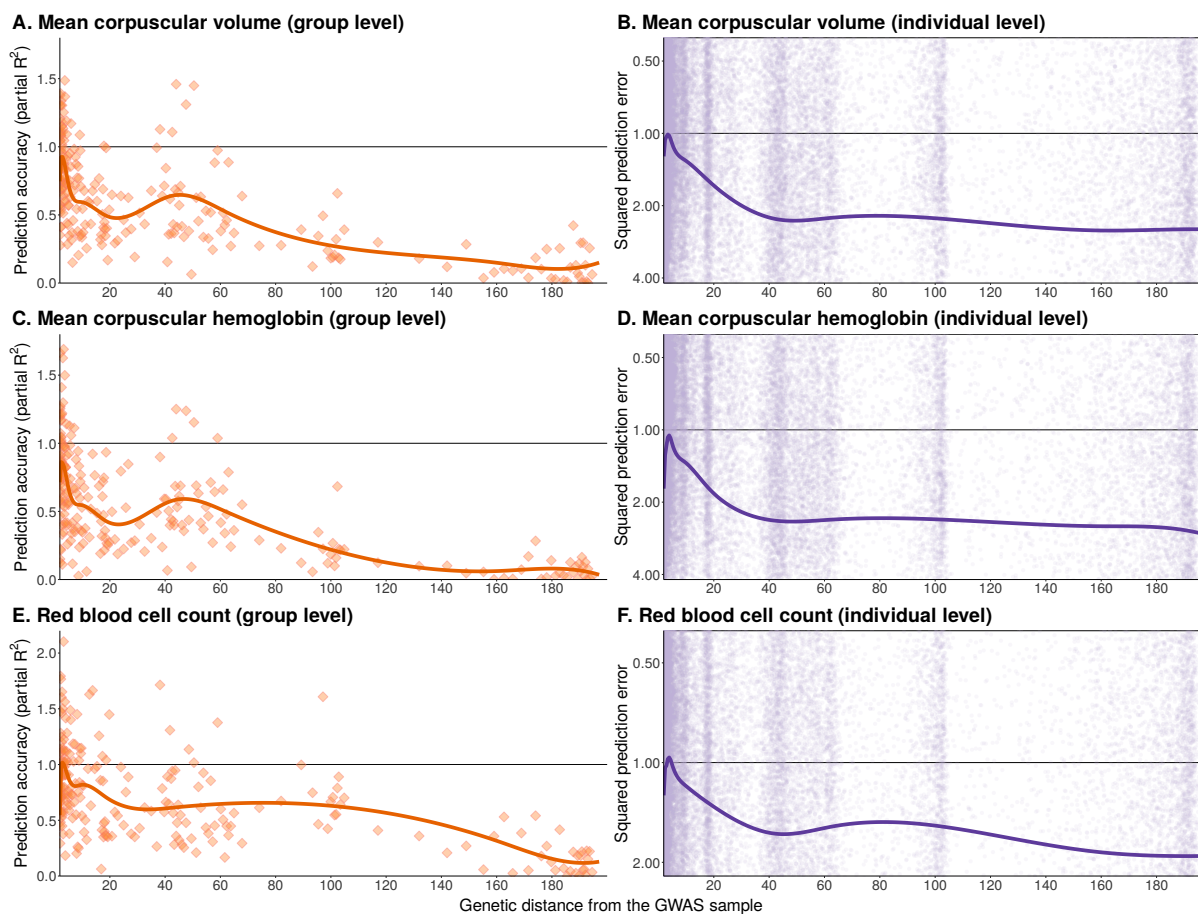


Figure S4: Trends of portability vary across traits and measures, for red blood-cell related traits. This figure presents the same analysis as **Fig. 2** in the main text, but shows red blood cell-related traits. At the group level (left panels), we measured prediction accuracy with the squared partial correlation between the PGS and the trait value in 500 bins of 258-259 individuals each. At the individual level (right panels), we measured the squared prediction error. Curves show cubic spline fits, with 8 knots placed based on the density of data points.

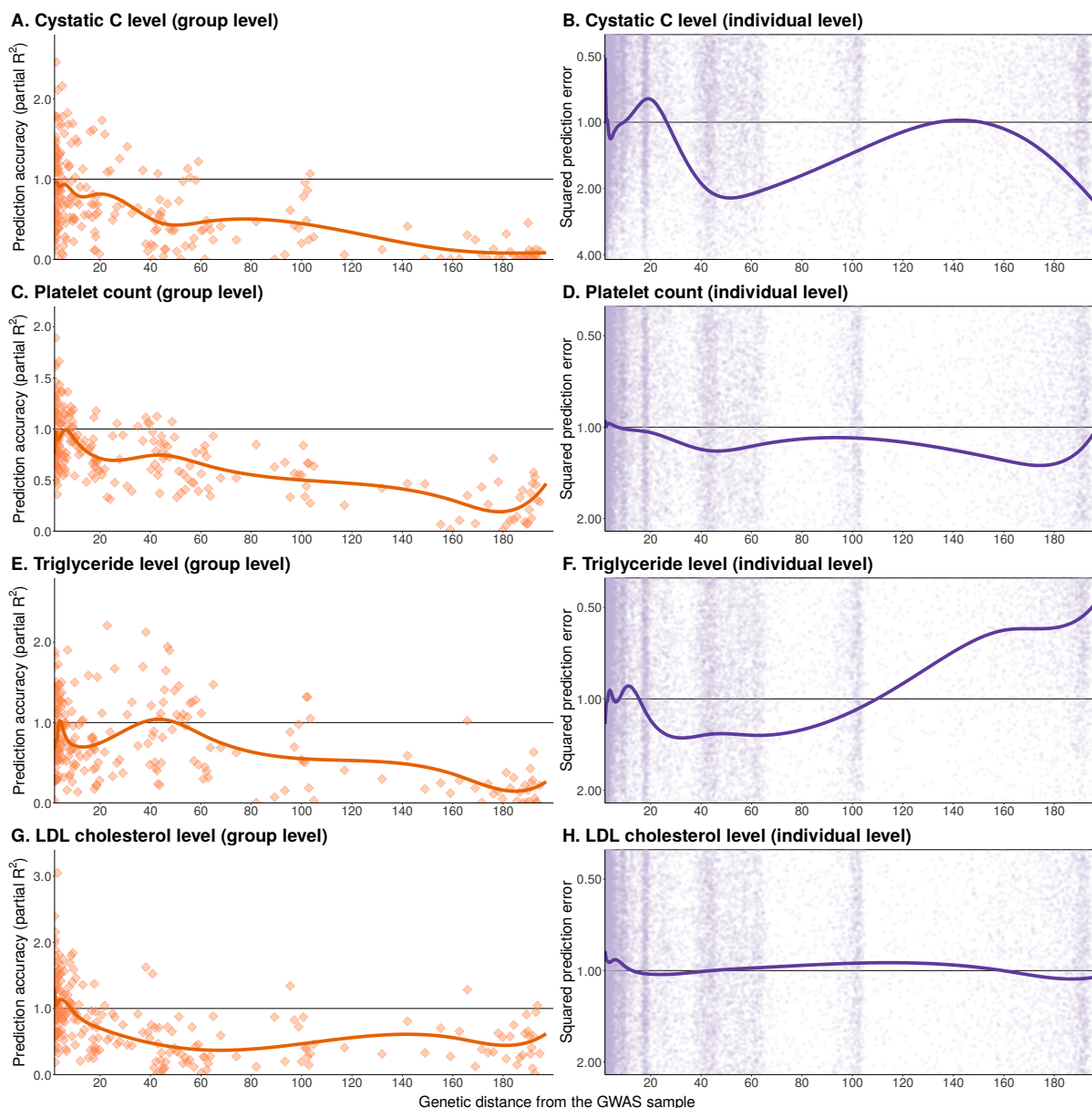


Figure S5: Trends of portability vary across traits and measures, for other biomarkers. This figure presents the same analysis as **Fig. 2** in the main text, but shows other biomarker-related traits. At the group level (left panels), we measured prediction accuracy with the squared partial correlation between the PGS and the trait value in 500 bins of 258-259 individuals each. At the individual level (right panels), we measured the squared prediction error. Curves show cubic spline fits, with 8 knots placed based on the density of data points.

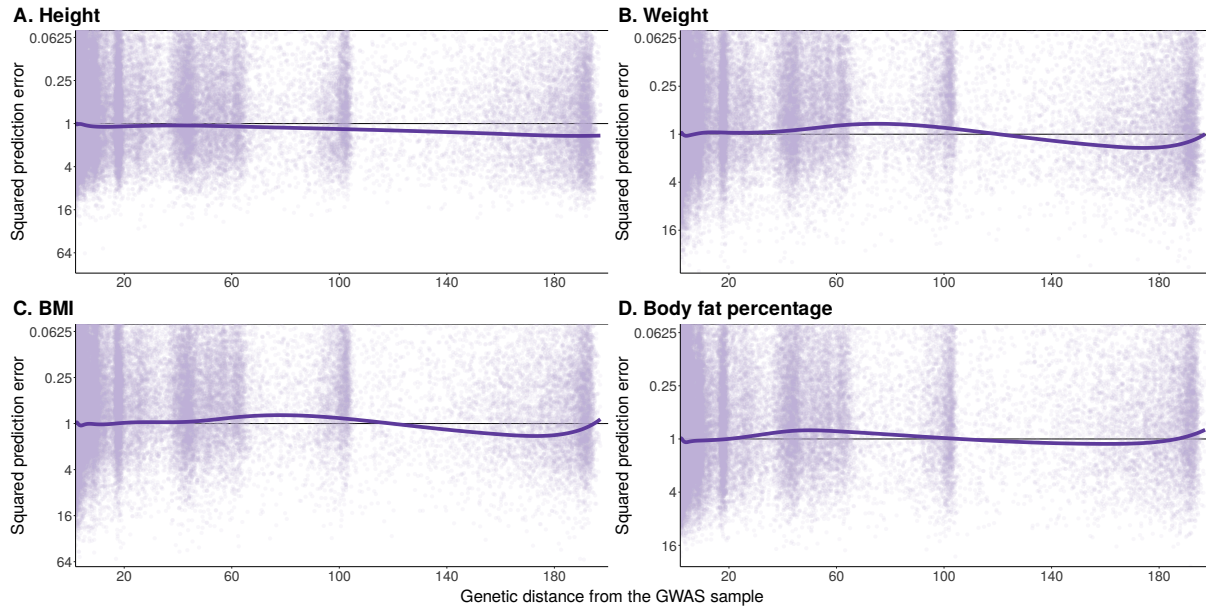


Figure S6: Full range of individual level prediction error for anthropometric measurements. This figure presents the same analysis as **Fig. 2** in the main text, but shows the full range of prediction error for other anthropometric traits. At the individual level, we measured the squared prediction error. Curves show cubic spline fits, with 8 knots placed based on the density of data points.

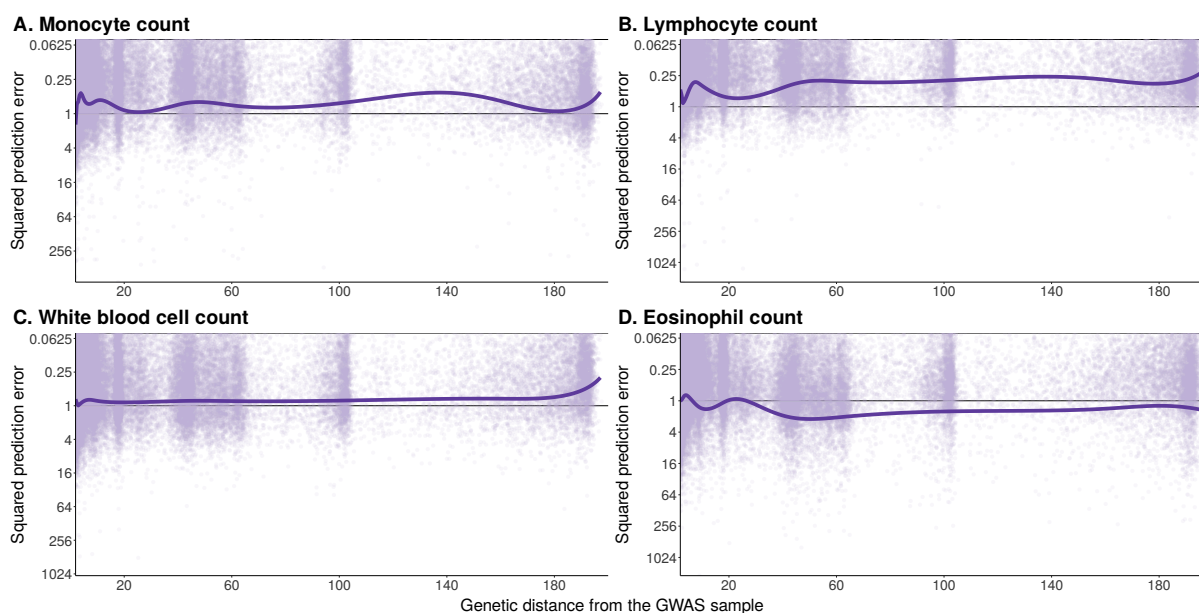


Figure S7: Full range of individual level prediction error for white blood cell-related traits. This figure presents the same analysis as **Fig. 2** in the main text, but shows the full range of prediction error for white blood cell-related traits. At the individual level, we measured the squared prediction error. Curves show cubic spline fits, with 8 knots placed based on the density of data points.

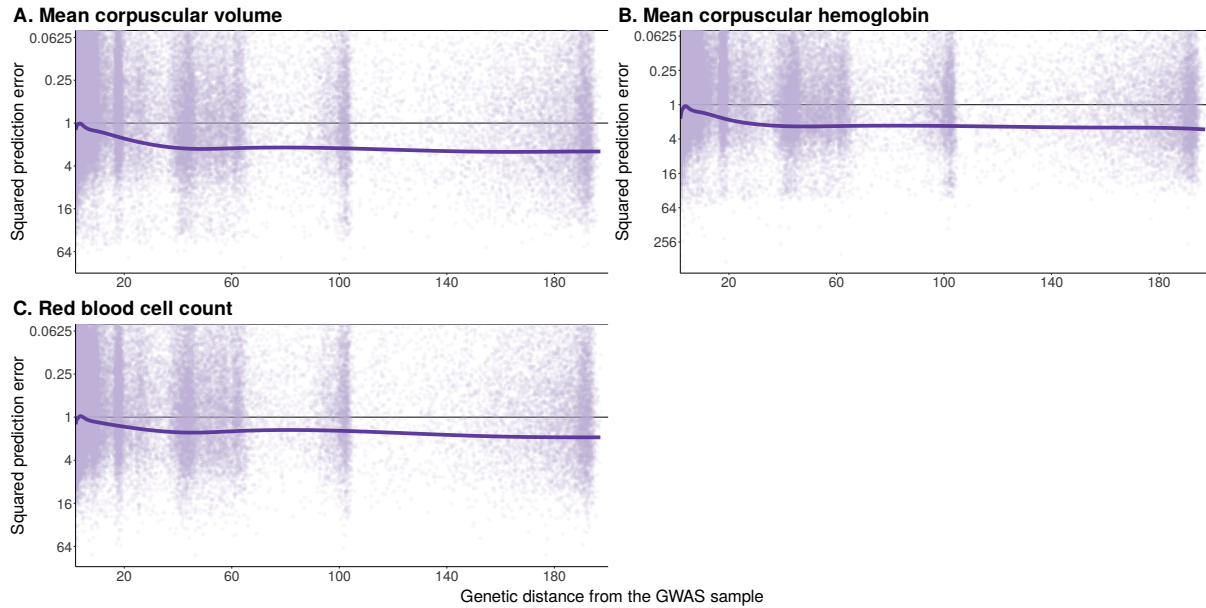


Figure S8: Full range of individual level prediction error for red blood cell-related traits. This figure presents the same analysis as **Fig. 2** in the main text, but shows the full range of prediction error for red blood cell-related traits. At the individual level, we measured the squared prediction error. Curves show cubic spline fits, with 8 knots placed based on the density of data points.

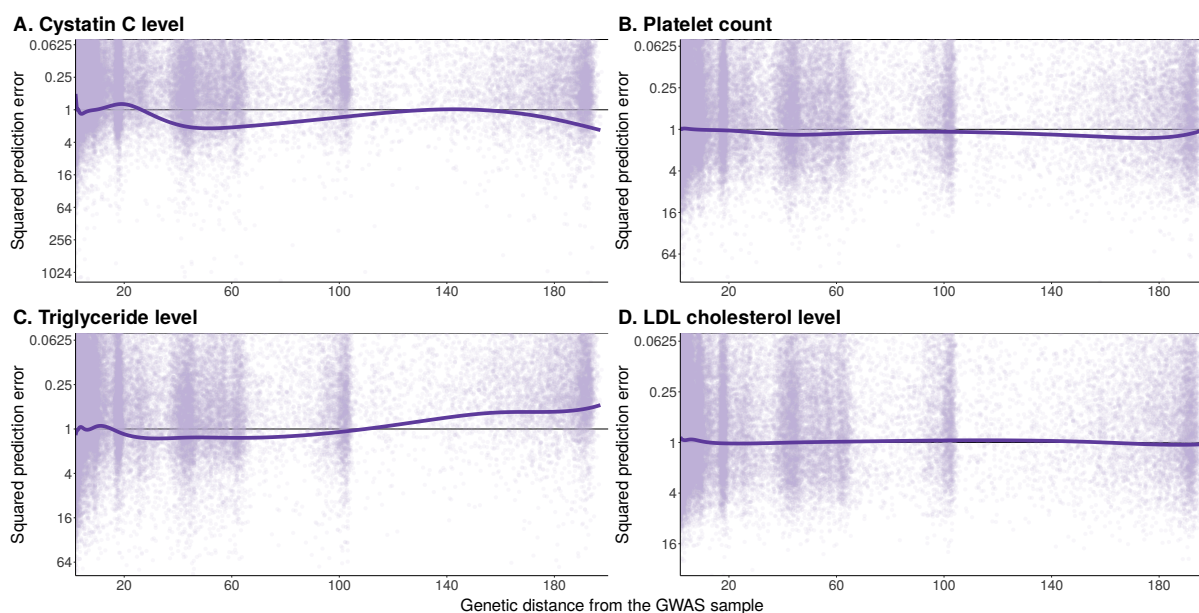


Figure S9: Full range of individual level prediction error for other biomarkers. This figure presents the same analysis as **Fig. 2** in the main text, but shows the full range of prediction error for other biomarkers. At the individual level, we measured the squared prediction error. Curves show cubic spline fits, with 8 knots placed based on the density of data points.

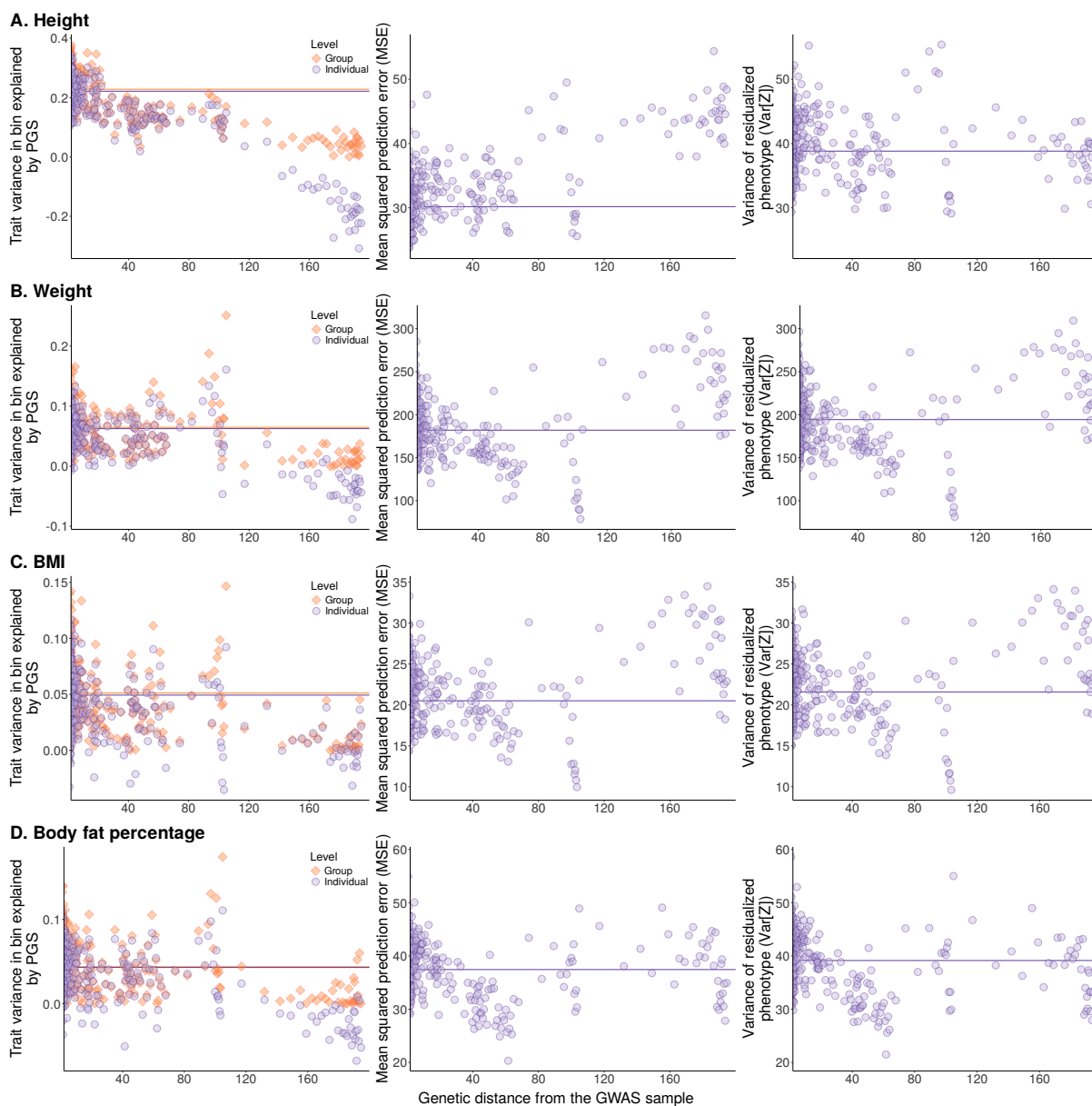


Figure S10: Divergence of group- and individual-level prediction accuracy for anthropometric measurements. The squared prediction error is with respect to the PGS as a predictor of a phenotypic value residualized for covariates (Z). Its mean (MSE) and the variation of residualized phenotype ($Var[Z]$) in each bin are shown in the middle and right column, respectively. In the left column, we show measures of the variance explained in bin of about 260 individuals, binned by genetic distance. “Group” level refers to the unstandardized partial R^2 between PGS and phenotype. “Individual” level refers to $1 - \frac{MSE}{Var[Z]}$. Horizontal lines show mean values across the 50 bins with a genetic distance closest to the mean genetic distance of the GWAS sample individuals (reference bins). The values of the 50 reference bins in all panels can be found in Table S1.

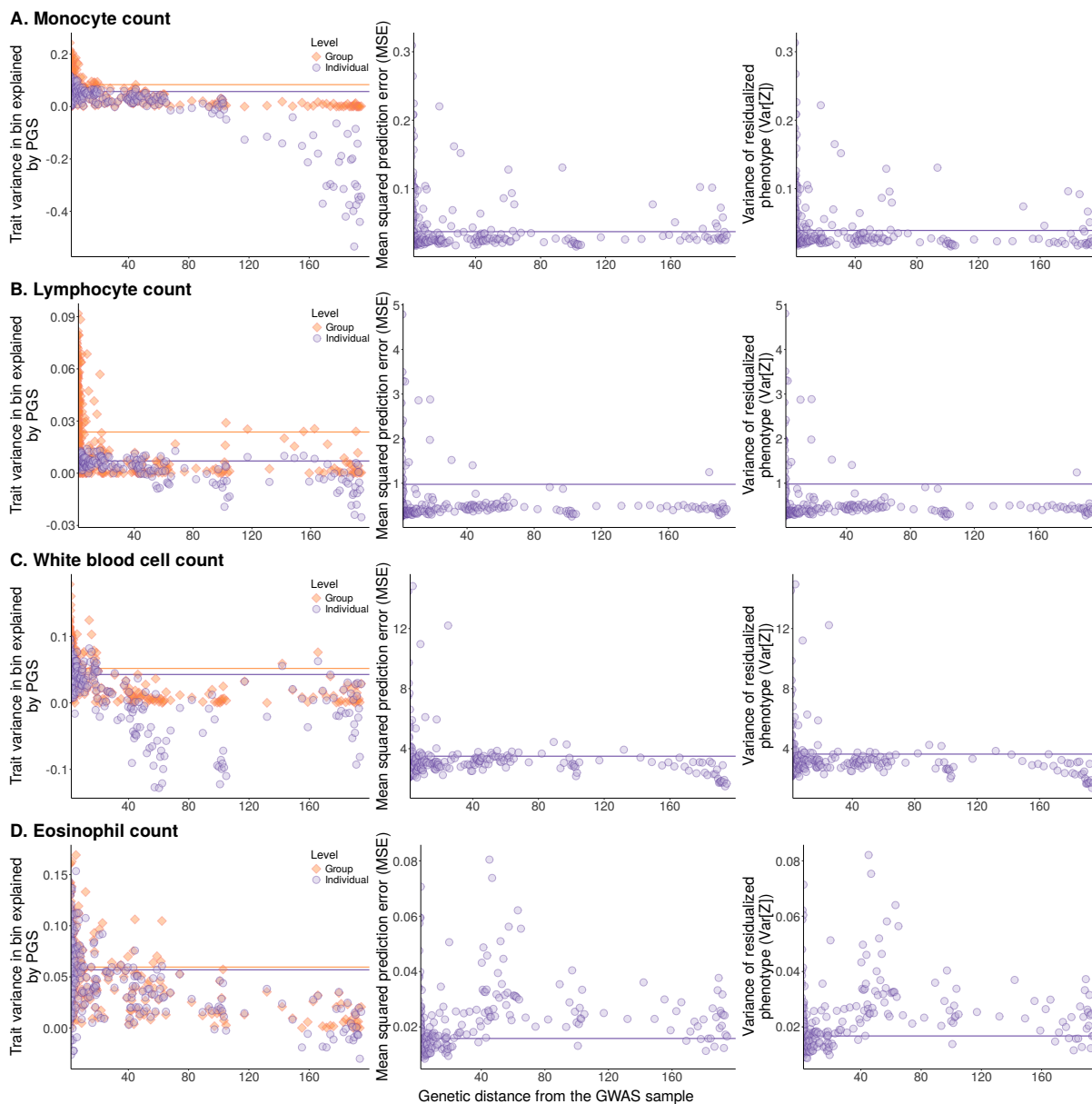


Figure S11: Divergence of group- and individual-level prediction accuracy for white blood cell-related traits. The squared prediction error is with respect to the PGS as a predictor of a phenotypic value residualized for covariates (Z). Its mean (MSE) and the variation of residualized phenotype ($Var[Z]$) in each bin are shown in the middle and right column, respectively. In the left column, we show measures of the variance explained in bin of about 260 individuals, binned by genetic distance. “Group” level refers to the unstandardized partial R^2 between PGS and phenotype. “Individual” level refers to $1 - \frac{MSE}{Var[Z]}$. Horizontal lines show mean values across the 50 bins with a genetic distance closest to the mean genetic distance of the GWAS sample individuals (reference bins). The values of the 50 reference bins in all panels can be found in **Table S1**.

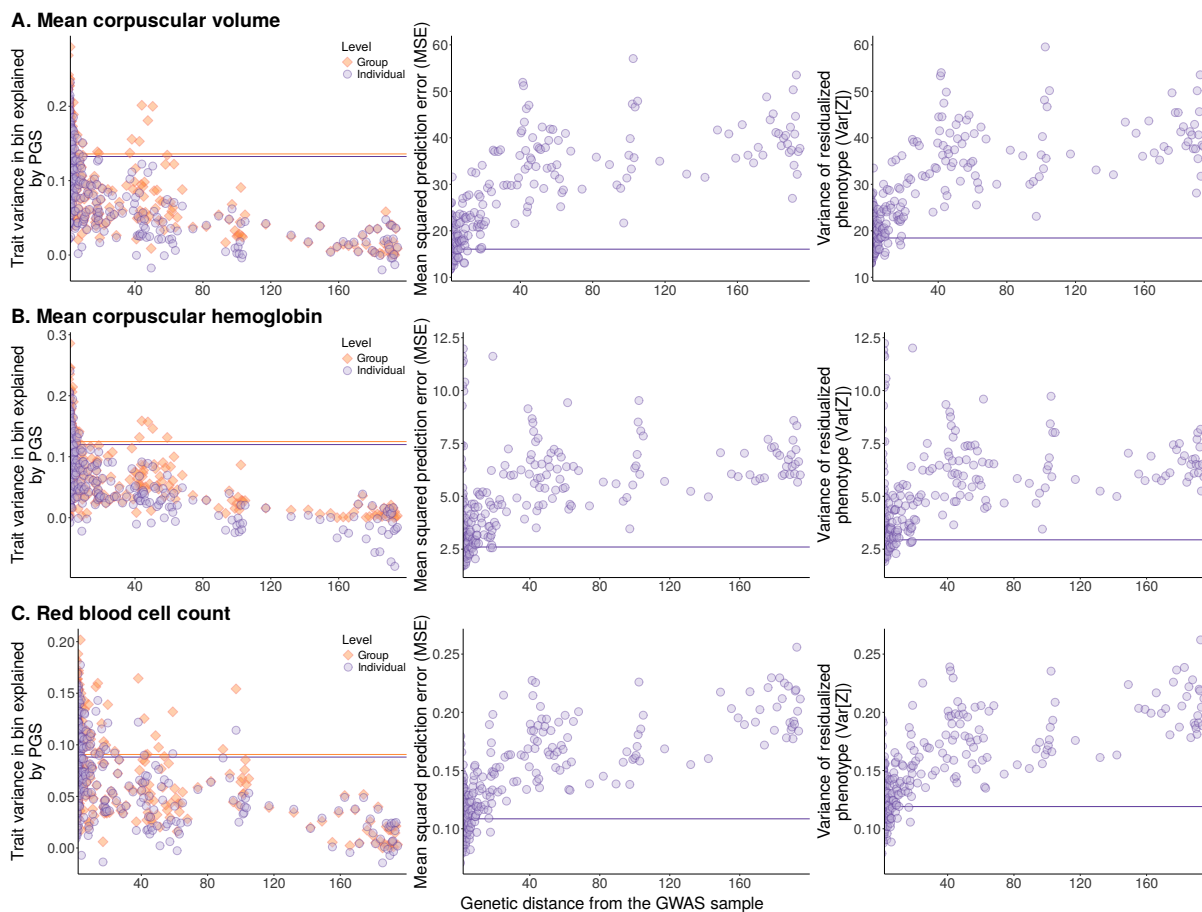


Figure S12: Divergence of group- and individual-level prediction accuracy for red blood cell-related traits. The squared prediction error is with respect to the PGS as a predictor of a phenotypic value residualized for covariates (Z). Its mean (MSE) and the variation of residualized phenotype ($Var[Z]$) in each bin are shown in the middle and right column, respectively. In the left column, we show measures of the variance explained in bin of about 260 individuals, binned by genetic distance. “Group” level refers to the unstandardized partial R^2 between PGS and phenotype. “Individual” level refers to $1 - \frac{MSE}{Var[Z]}$. Horizontal lines show mean values across the 50 bins with a genetic distance closest to the mean genetic distance of the GWAS sample individuals (reference bins). The values of the 50 reference bins in all panels can be found in **Table S1**.

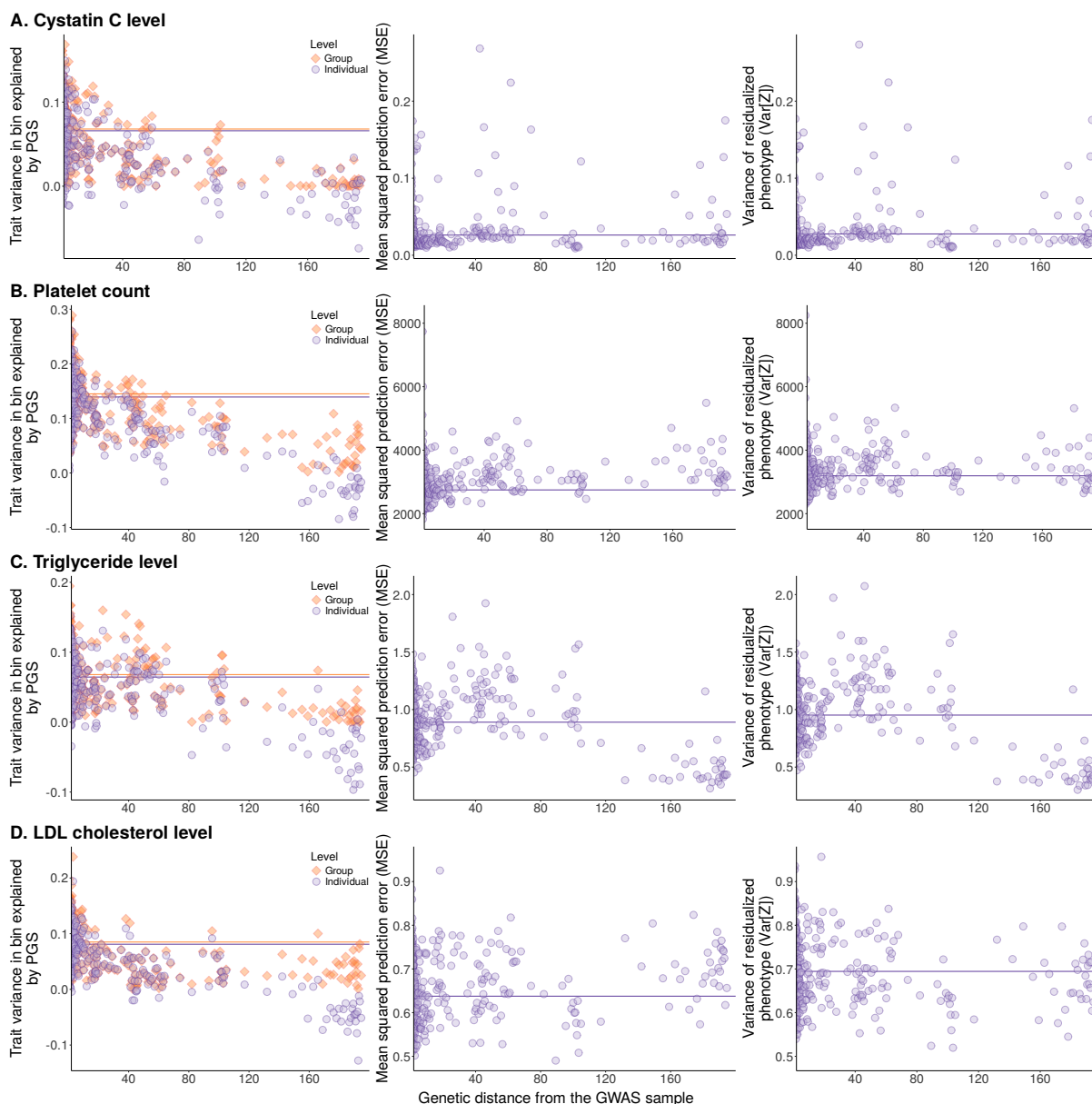


Figure S13: Divergence of group- and individual-level prediction accuracy for other biomarkers. The squared prediction error is with respect to the PGS as a predictor of a phenotypic value residualized for covariates (Z). Its mean (MSE) and the variation of residualized phenotype ($Var[Z]$) in each bin are shown in the middle and right column, respectively. In the left column, we show measures of the variance explained in bin of about 260 individuals, binned by genetic distance. “Group” level refers to the unstandardized partial R^2 between PGS and phenotype. “Individual” level refers to $1 - \frac{MSE}{Var[Z]}$. Horizontal lines show mean values across the 50 bins with a genetic distance closest to the mean genetic distance of the GWAS sample individuals (reference bins). The values of the 50 reference bins in all panels can be found in **Table S1**.

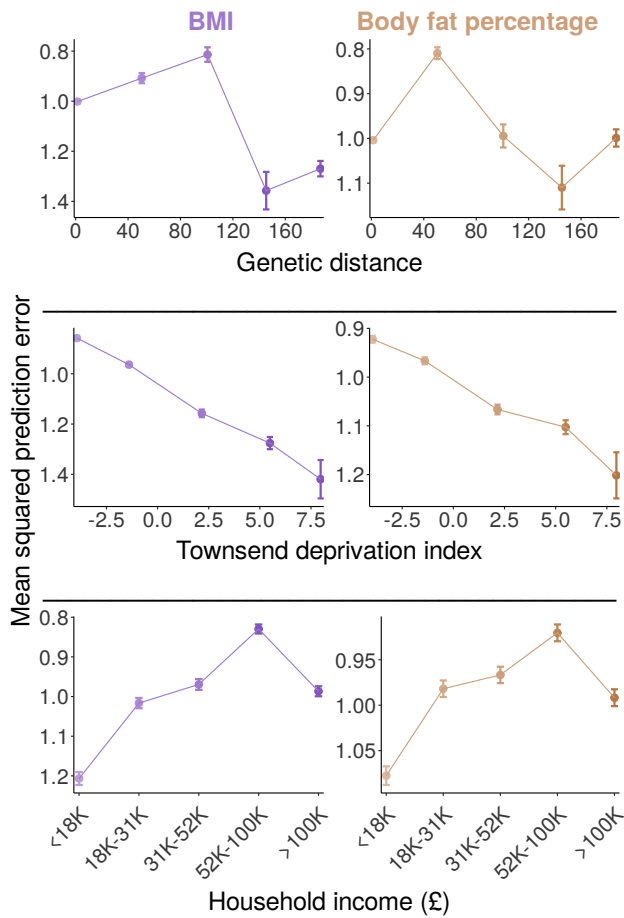


Figure S14: Mean trends in individual-level prediction accuracy by different measures for anthropometric measurements. This figure presents the same analysis as **Fig. 3A** in the main text, but shows other anthropometric traits: body mass index and body fat percentage. Data points confer to mean (\pm SE) squared prediction errors of individuals in the prediction sample, binned into 5 equidistant strata. The x-axis shows the median measure value for each stratum. “Household income” refers to average yearly total household income before tax.

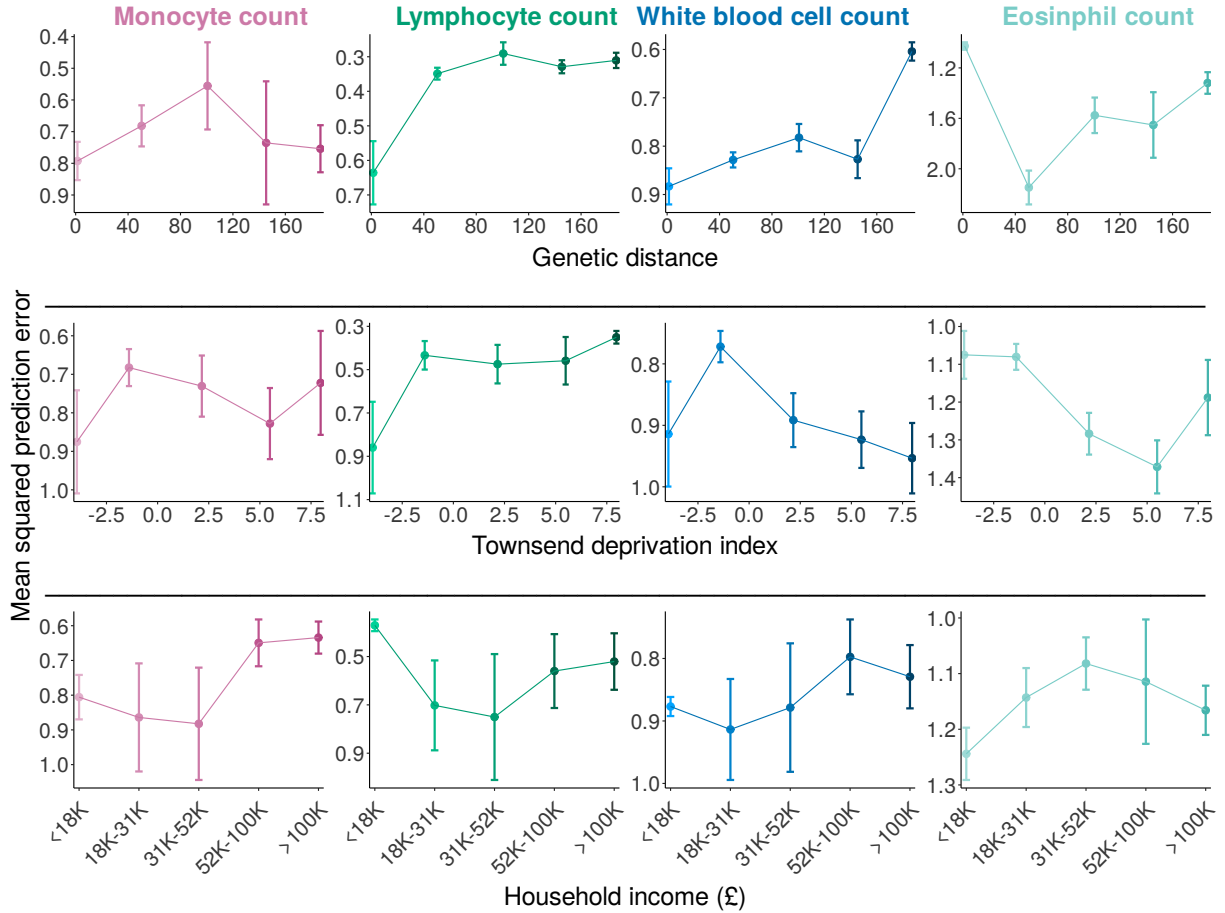


Figure S15: Mean trends in individual-level prediction accuracy by different measures for white blood cell-related traits. This figure presents the same analysis as **Fig. 3A** in the main text, but shows the mean trends for white blood cell-related traits. Data points confer to mean (\pm SE) squared prediction errors of individuals in the prediction sample, binned into 5 equidistant strata. The x-axis shows the median measure value for each stratum. “Household income” refers to average yearly total household income before tax.

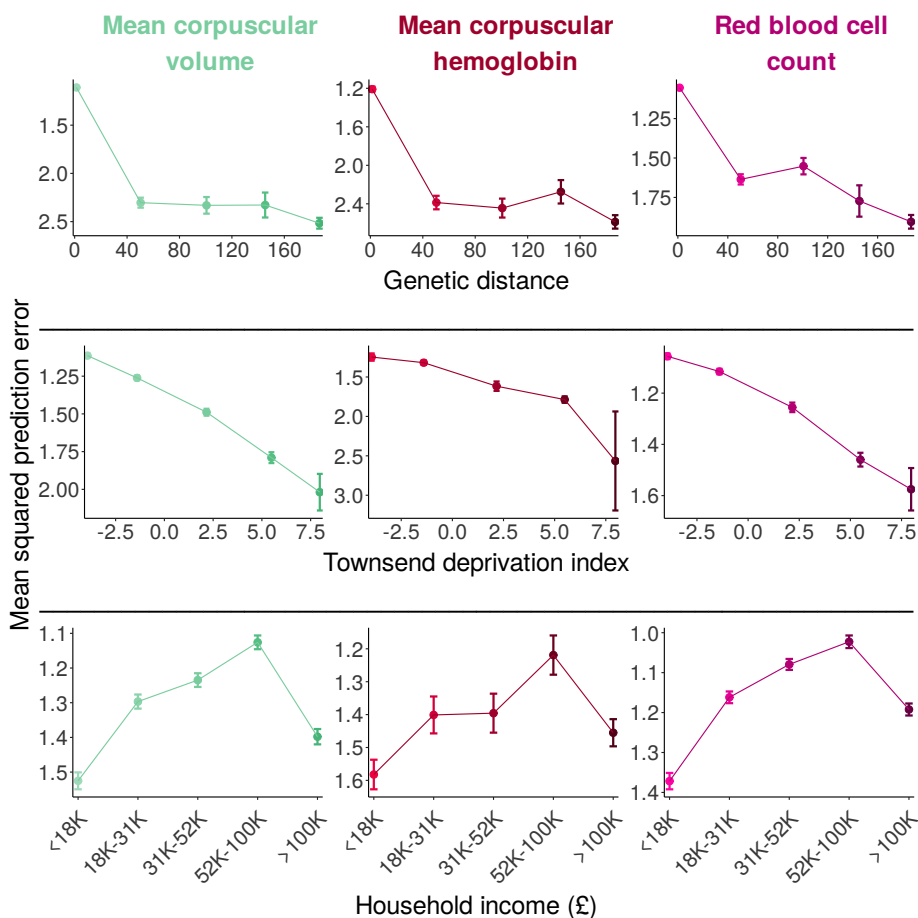


Figure S16: Mean trends in individual-level prediction accuracy by different measures for red blood cell-related traits. This figure presents the same analysis as **Fig. 3A** in the main text, but shows different traits: mean corpuscular volume, mean corpuscular hemoglobin, and red blood cell count. Data points confer to mean (\pm SE) squared prediction errors of individuals in the prediction sample, binned into 5 equidistant strata. The x-axis shows the median measure value for each stratum. “Household income” refers to average yearly total household income before tax.

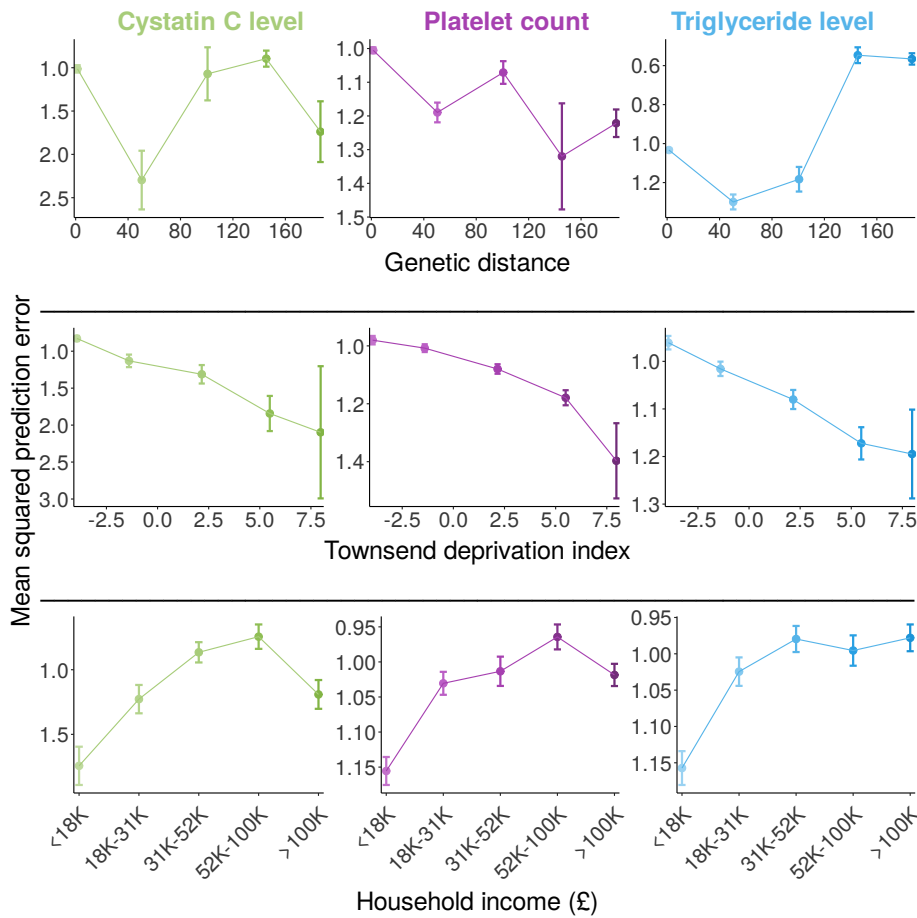


Figure S17: Mean trends in individual-level prediction accuracy by different measures for other biomarkers. This figure presents the same analysis as **Fig. 3A** in the main text, but shows different traits: cystatin C level, platelet count, and triglyceride level. Data points confer to mean (\pm SE) squared prediction errors of individuals in the prediction sample, binned into 5 equidistant strata. The x-axis shows the median measure value for each stratum. “Household income” refers to average yearly total household income before tax.

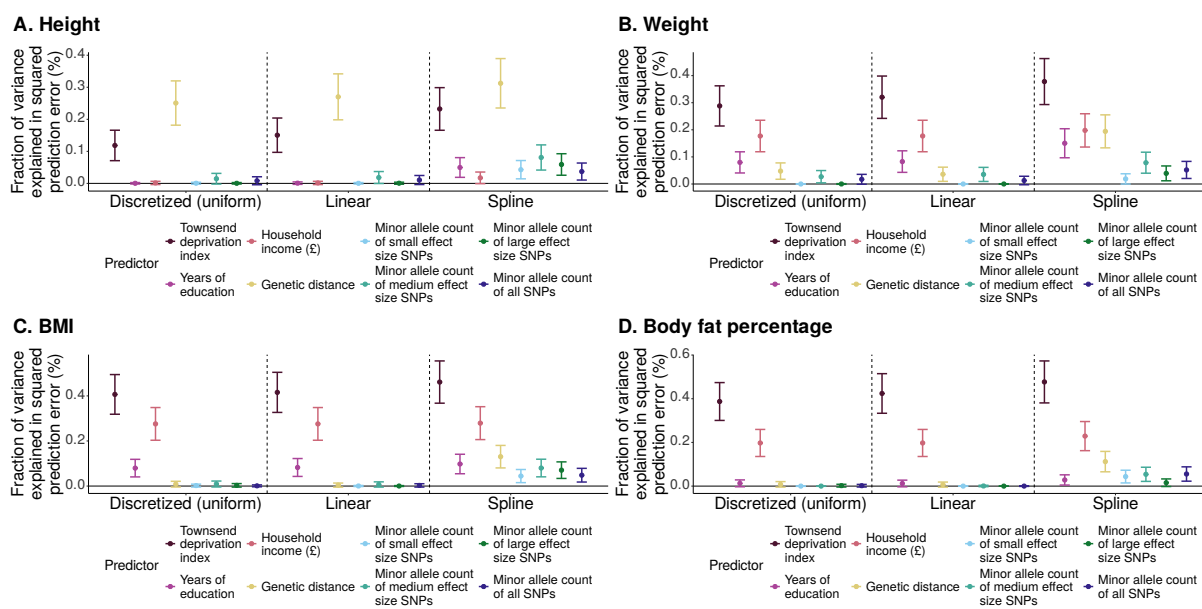


Figure S18: Individual-level prediction error explained by different measures for anthropometric measurements. This figure presents the same analysis as **Fig. 3B** in the main text, but shows more measures and methods for fitting the measures, and focusing on anthropometric traits. “Minor allele” refers to the minor allele with respect to the GWAS sample. For each measure, we independently fit three different models. “Discretized (uniform)” refers to a discretized predictor, using one predicted value per each of the 5 bins where all 5 bins had identical widths. “Linear” refers to a linear predictor, fit using Ordinary Least Squares (OLS). “Spline” refers to a cubic spline, for which 16 knots were placed based on the density of data points, such that there was an equal number of data points between each pair of consecutive knots.

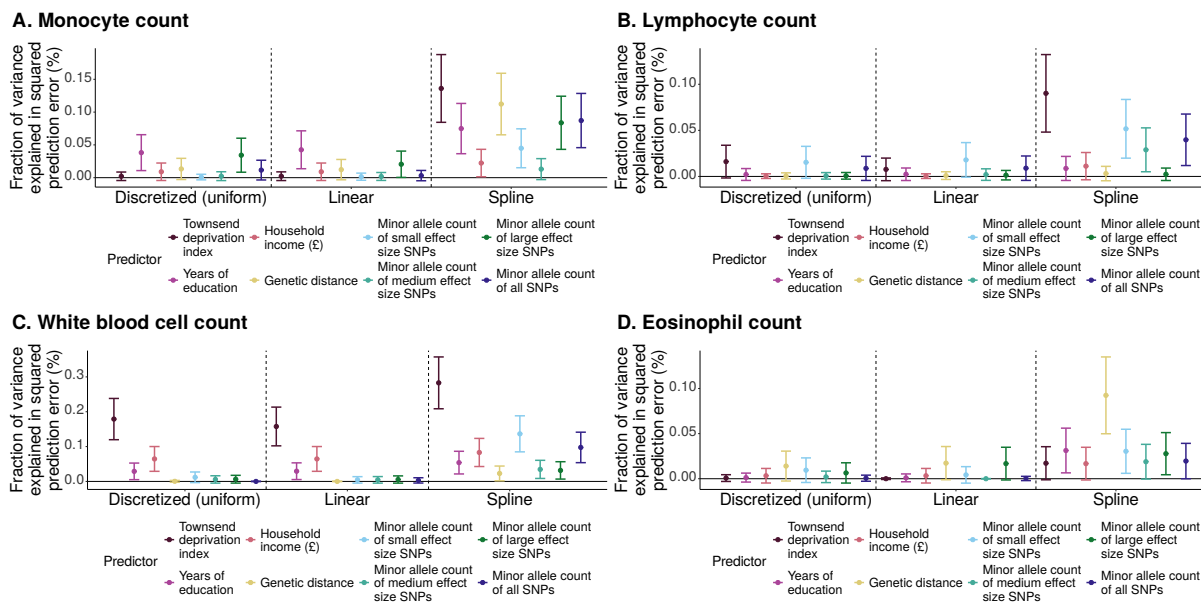


Figure S19: Individual-level prediction error explained by different measures for white blood cell-related traits. This figure presents the same analysis as **Fig. 3B** in the main text, but shows more measures and methods for fitting the measures. “Minor allele” refers to the minor allele with respect to the GWAS sample. For each measure, we independently fit three different models. “Discretized (uniform)” refers to a discretized predictor, using one predicted value per each of the 5 bins where all 5 bins had identical widths. “Linear” refers to a linear predictor, fit using Ordinary Least Squares (OLS). “Spline” refers to a cubic spline, for which 16 knots were placed based on the density of data points, such that there was an equal number of data points between each pair of consecutive knots.

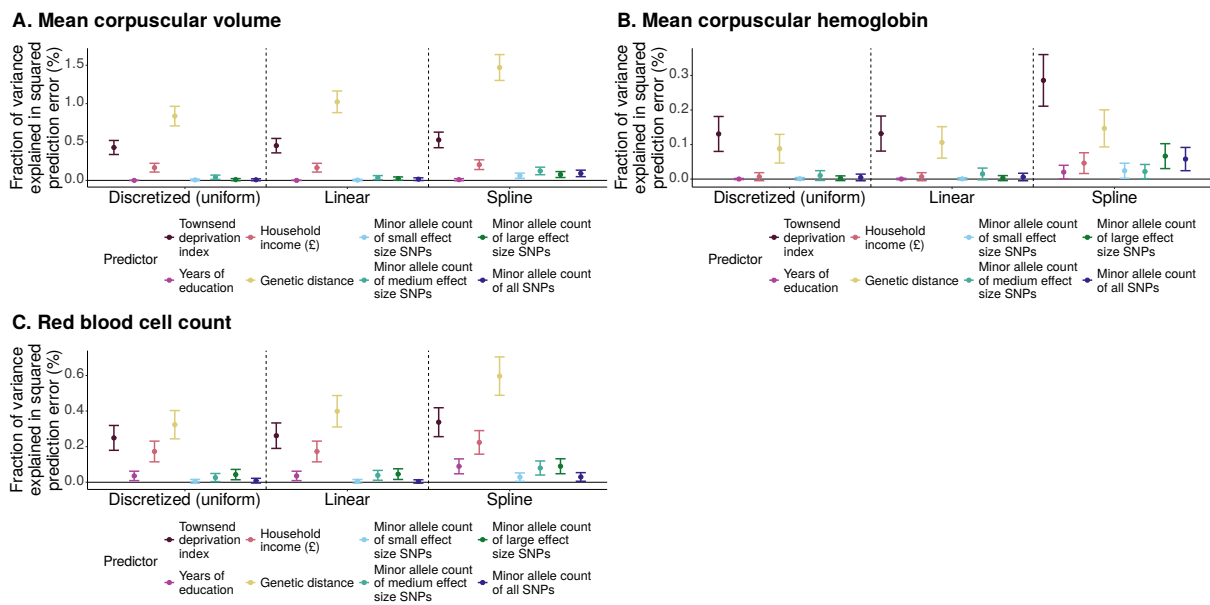


Figure S20: Individual-level prediction error explained by different measures for red blood cell-related traits. This figure presents the same analysis as **Fig. 3B** in the main text, but shows more measures and methods for fitting the measures. “Minor allele” refers to the minor allele with respect to the GWAS sample. For each measure, we independently fit three different models. “Discretized (uniform)” refers to a discretized predictor, using one predicted value per each of the 5 bins where all 5 bins had identical widths. “Linear” refers to a linear predictor, fit using Ordinary Least Squares (OLS). “Spline” refers to a cubic spline, for which 16 knots were placed based on the density of data points, such that there was an equal number of data points between each pair of consecutive knots.

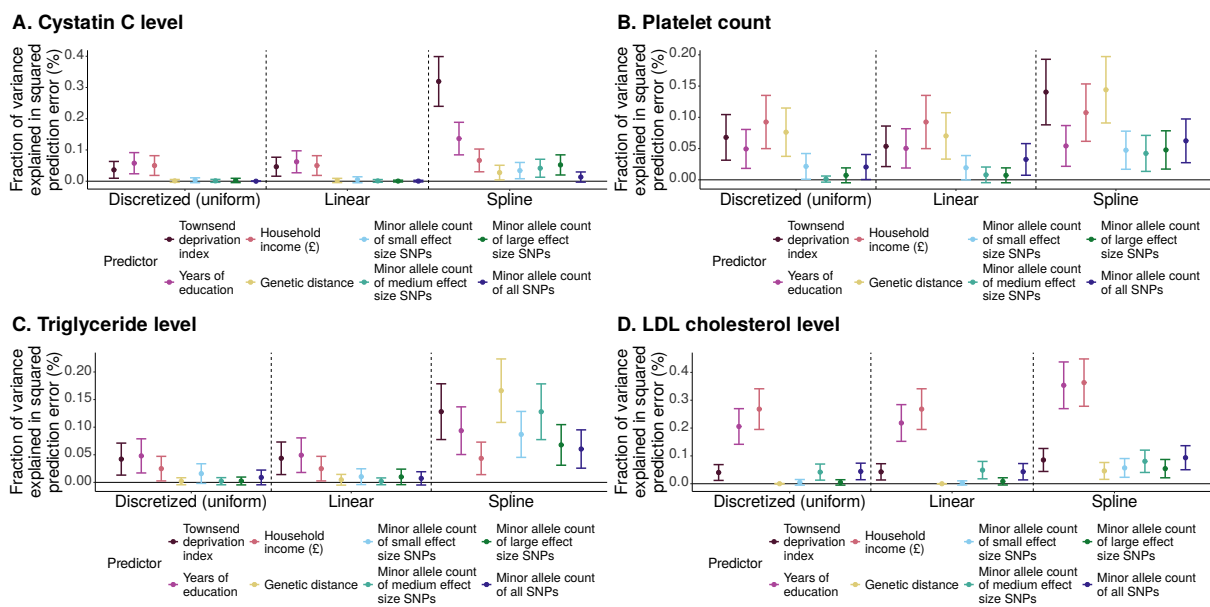


Figure S21: Individual-level prediction error explained by different measures for other biomarkers. This figure presents the same analysis as **Fig. 3B** in the main text, but shows more measures and methods for fitting the measures. “Minor allele” refers to the minor allele with respect to the GWAS sample. For each measure, we independently fit three different models. “Discretized (uniform)” refers to a discretized predictor, using one predicted value per each of the 5 bins where all 5 bins had identical widths. “Linear” refers to a linear predictor, fit using Ordinary Least Squares (OLS). “Spline” refers to a cubic spline, for which 16 knots were placed based on the density of data points, such that there was an equal number of data points between each pair of consecutive knots.

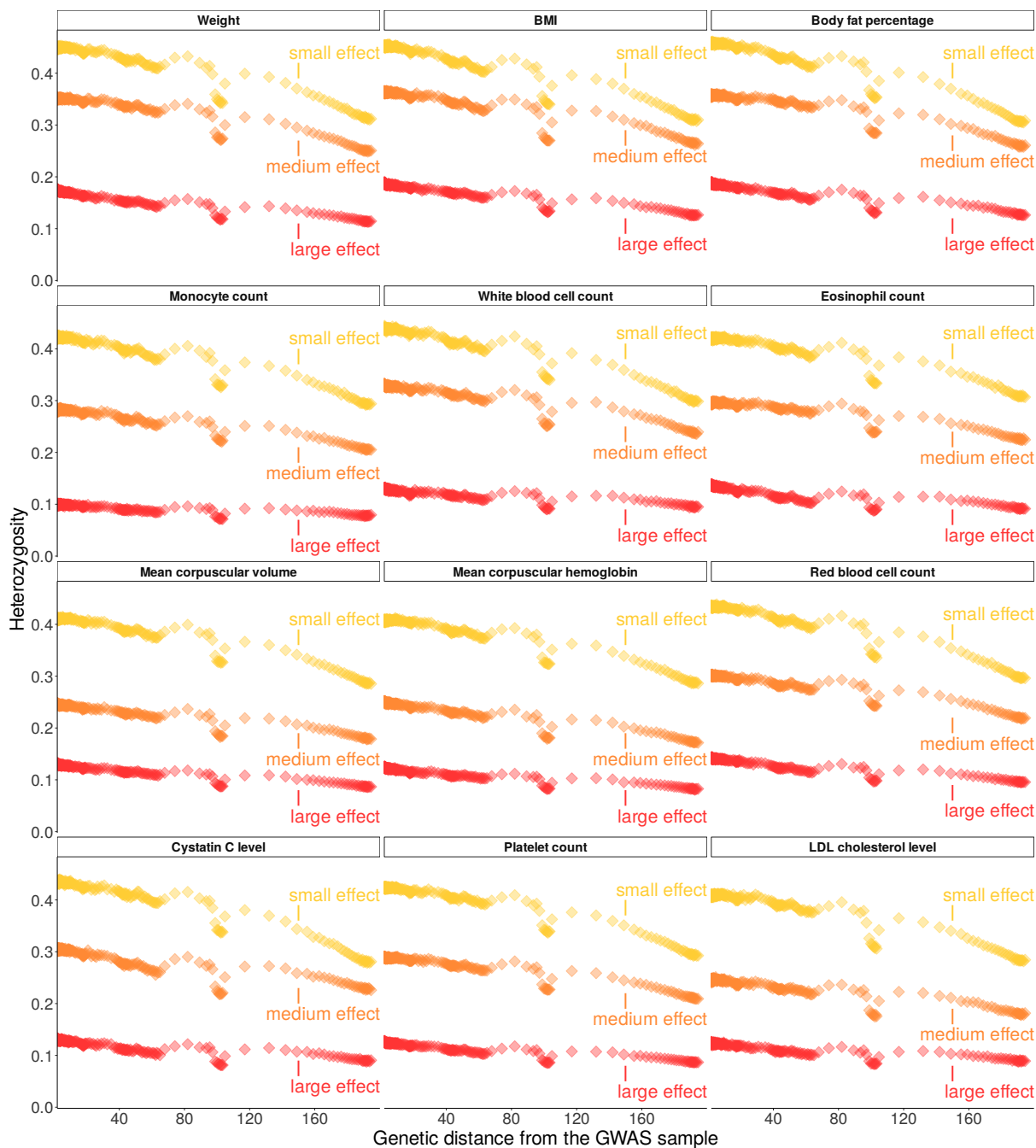


Figure S22: Mean heterozygosity of SNPs, stratified by effect size. This figure presents the same analysis as **Fig. 4B** in the main text, but for the 12 phenotypes not included there. For each trait, SNPs are stratified into three equal-sized strata (small, medium, and large) based on squared effect sizes (**Fig. S23**). Each data point is the mean heterozygosity of a stratum in a bin of genetic distance.

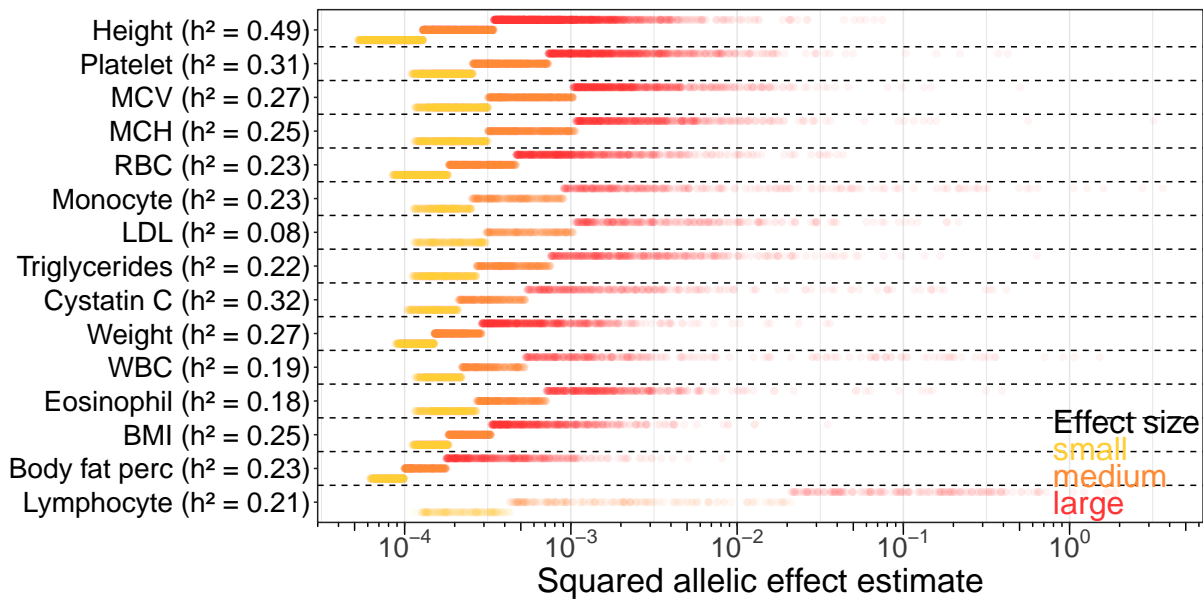


Figure S23: Squared allelic effect estimate. For each trait, the index SNPs of the respective PGS are stratified into three equal-sized strata (small, medium, and large) based on squared effect sizes. The x-axis represents the squared effect sizes in units of trait variance in the GWAS set. SNP heritabilities (h^2) are taken from the Neale Lab's UKB analysis²¹. Each data point represents a SNP. MCV: mean corpuscular volume. MCH: mean corpuscular hemoglobin. RBC: red blood cell count. Body fat perc: body fat percentage. WBC: white blood cell count. LDL: LDL cholesterol level.

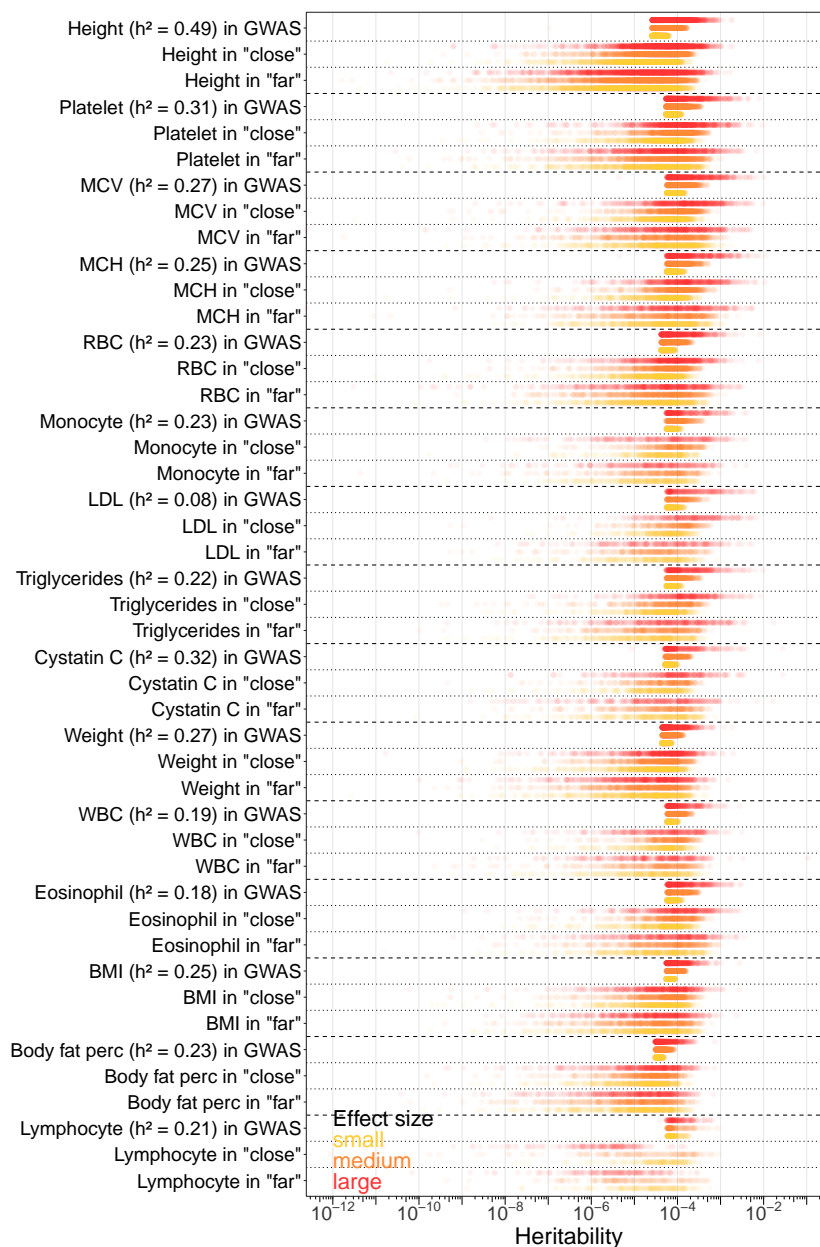


Figure S24: Heritability explained by index SNPs. For each trait, index SNPs of the respective PGS are stratified by their squared effect sizes. Data points confer to index SNPs. In each of the GWAS sample, the “close” subset of the prediction sample (genetic distance ≤ 10 in the prediction set, with 96,457 individuals) and the “far” subsample of the prediction sample (genetic distance > 10 , with 32,822 individuals), we estimate the allelic effect of each SNP (in units of trait standard deviations) and its heterozygosity. The product of the two is the estimated heritability. SNP heritability estimates (h^2) on the y-axis are taken from the Neale Lab’s UKB analysis²¹.