# Calibrated prediction intervals for polygenic scores across diverse contexts

3

Kangcheng Hou<sup>1</sup>, Ziqi Xu<sup>2</sup>, Yi Ding<sup>1</sup>, Arbel Harpak<sup>3,4</sup>, Bogdan Pasaniuc<sup>1,5,6,7</sup>

4 5

<sup>1</sup> Bioinformatics Interdepartmental Program, University of California, Los Angeles, Los Angeles, CA, USA

<sup>7</sup> <sup>2</sup> Department of Computer Science, University of California, Los Angeles, Los Angeles, CA, USA

<sup>3</sup> Department of Population Health, The University of Texas at Austin, Austin, TX, USA

<sup>9</sup> <sup>4</sup> Department of Integrative Biology, The University of Texas at Austin, Austin, TX, USA

<sup>5</sup> Department of Pathology and Laboratory Medicine, David Geffen School of Medicine, University of California, Los Angeles, Los Angeles, CA, USA

<sup>6</sup> Department of Computational Medicine, David Geffen School of Medicine, University of California, Los
 Angeles, Los Angeles, CA, USA

<sup>7</sup> Institute for Precision Health, University of California, Los Angeles, Los Angeles

15 To whom correspondence should be addressed: <u>houkc@ucla.edu</u>, <u>pasaniuc@ucla.edu</u>

# 16 Abstract

Polygenic scores (PGS) have emerged as the tool of choice for genomic prediction in a wide 17 range of fields from agriculture to personalized medicine. We analyze data from two large 18 biobanks in the US (All of Us) and the UK (UK Biobank) to find widespread variability in PGS 19 performance across contexts. Many contexts, including age, sex, and income, impact PGS 20 accuracies with similar magnitudes as genetic ancestry. PGSs trained in single versus multi-21 ancestry cohorts show similar context-specificity in their accuracies. We introduce trait prediction 22 intervals that are allowed to vary across contexts as a principled approach to account for context-23 specific PGS accuracy in genomic prediction. We model the impact of all contexts in a joint 24 framework to enable PGS-based trait predictions that are well-calibrated (contain the trait value 25 with 90% probability in all contexts), whereas methods that ignore context are mis-calibrated. We 26 show that prediction intervals need to be adjusted for all considered traits ranging from 10% for 27 diastolic blood pressure to 80% for waist circumference. Adjustment of prediction intervals 28 depends on the dataset; for example, prediction intervals for education years need to be adjusted 29 by 90% in All of Us versus 8% in UK Biobank. Our results provide a path forward towards 30 utilization of PGS as a prediction tool across all individuals regardless of their contexts while 31 highlighting the importance of comprehensive profile of context information in study design and 32 data collection. 33

## 34 Introduction

Accurate prediction of complex diseases/traits integrating genetic and non-genetic factors is essential for a wide range of fields from agriculture to personalized genomic medicine. The genetic contribution to common traits is typically predicted using polygenic scores (PGS) that summarize the joint contribution of many genetic factors<sup>1–4</sup>. A critical barrier in PGS use is their *context-specific accuracy* – their performance (and/or bias) varies across genetic ancestry<sup>5–9</sup>, age, sex, socioeconomic status and other factors<sup>10–12</sup>. This prevents equitable use of PGS across individuals of all contexts<sup>4,5,13</sup>.

PGS use data from large-scale genome-wide association studies (GWAS) to estimate linear 42 prediction models of traits based on genetic variants; these prediction models are then used for 43 new data that often has different context characteristics from the GWAS training data (e.g., 44 different distributions of genetic ancestry, social determinants of health, etc.)<sup>1,2,14</sup>. Even when 45 testing data is similar to training data, genetic effects themselves can vary by contexts (e.g., due 46 to genotype-environment interaction, across age<sup>15</sup>, sex<sup>16</sup>, genetic ancestry<sup>17–20</sup>) thus leading to 47 differential PGS performance (as traditional PGS do not model such interactions). Furthermore, 48 when genetic effects are unknown, allele frequency, linkage disequilibrium and differential tagging 49 of true latent genetic factors can also lead to context-specific accuracy of PGS-based 50 predictions<sup>10,15,21</sup>. 51

To account for PGS accuracy variability, we propose an approach to incorporate context-52 specificity into trait prediction intervals that are allowed to vary across contexts. Trait prediction 53 intervals denote the range containing true trait values with pre-specified confidence (e.g., 90%). 54 And they provide a natural approach to model variability in PGS accuracies – narrower prediction 55 intervals correspond to contexts where PGS attains higher accuracy - that can then be used in 56 applications of PGS-based trait predictions<sup>10,22,23</sup>. As an example, consider the case of two 57 individuals with the same PGS-based predictions for low-density lipoprotein cholesterol (LDL) of 58 120 mg/dL. If the two individuals have different contexts (e.g., sex) that are known to impact PGS 59 accuracy (e.g., R<sup>2</sup>=0.1 in men vs. 0.2 in women), their prediction intervals will also vary (e.g., 120 60 ± 40 mg/dL vs. 120 ± 10 mg/dL). In this example the second individual is more likely to meet a 61 decision criterion of LDL>100 mg/dL for clinical intervention. 62

To achieve calibration across all contexts, we propose a statistical model (*CalPred*) that jointly models the effects of all contexts on PGS accuracy leveraging calibration data. The key assumption is that new target individuals for whom PGS-based predictions will be employed have similar contexts as the calibration data. This is motivated by precision health efforts that created EHR-linked biobanks of patients from the same medical system in which the PGS-based prediction will be applied in the future<sup>24–27</sup>; in this context the assumption is that the biobank is representative of future patients entering the same medical system.

First, we analyze data across two large-scale biobanks (UK Biobank<sup>28</sup> and All of Us<sup>29</sup>) to find pervasive impact of context on PGS accuracy across a wide range of traits. All considered traits (N=72) have at least one context impacting their accuracies<sup>10,12</sup>. Socio-economic contexts have similar magnitudes of impact on PGS accuracies as genetic ancestry; for example, PGS accuracy varies by up to ~50% for individuals across the context of "education years" averaged across all considered traits in All of Us. Moreover, socio-economic contexts have greater impact on PGS accuracy in All of Us, a more diverse dataset, as compared to UK Biobank.

Second, we use simulations and real data analysis to find that CalPred provides well-calibrated prediction intervals across individuals of diverse contexts. For example, CalPred jointly models the impact of genetic ancestry, age and sex and other social determinants of health on LDL prediction to find that prediction intervals need adjustment by up to ~40% across contexts to achieve calibration. The context-specificity of PGS prediction varies across traits, with largest adjustments observed for traits including waist circumference and average mean spherical

equivalent (avMSE) where prediction intervals need adjustment by ~100% for individuals in certain contexts; meanwhile certain traits such as diastolic blood pressure only need a modest adjustment by ~20%. Notably, the context-specificity of the same trait also depends on the studied population; for example, prediction intervals for education years need adjustment by 90% in All of Us versus 8% in UK Biobank, reflecting the more diverse distribution of education years and other social determinants of health in All of Us. Overall, our approaches provide a path forward to modeling differential PGS accuracy by context in prediction of complex traits in humans.

## 90 **Results**

#### 91 Overview

We incorporate context-specific accuracy in PGS-based predictions using prediction intervals that 92 are allowed to vary across contexts to maintain calibration; the true phenotype is contained within 93 the prediction interval at a pre-specified probability (e.g., 90%; Figure 1a). Naturally, as accuracy 94 varies by context, the interval width needs to vary adaptively such that calibration is maintained 95 (Figure 1b). For illustrative purpose we distinguish among three types of prediction intervals 96 (Figure 2). First, standard errors of PGS weights can be used to estimate prediction intervals that 97 do not vary across contexts and/or individuals; these types of intervals are calibrated only when 98 target perfectly matches training which is hard to achieve in practice. Second, prediction intervals 99 can be estimated empirically using a calibration dataset across all data ignoring context<sup>1,30–34</sup>; 100 these types of intervals are robust to mismatches between training and testing, but are mis-101 calibrated in particular contexts due to the variability of PGS accuracy. Third, prediction intervals 102 that vary across contexts can be estimated using a calibration dataset by empirically quantifying 103 the impact of each context on prediction accuracy; context-specific prediction intervals are 104 adaptive and robust across contexts albeit at the expense of a more complex statistical model 105 and larger calibration data that spans all contexts. Motivated by clinical implementation of PGS-106 based predictions in medical systems where EHR-linked biobanks already exist, here we focus 107 on leveraging calibration data to estimate context-specific prediction intervals. In this scenario it 108 is natural to use existing EHR-linked biobanks as approximation for future patients within the 109 same medical system. For example, UCLA ATLAS biobank<sup>24</sup> contains data of ~150k patients 110 within the UCLA Health system that can be used to calibrate PGS-based predictors for future 111 visits of UCLA patients. 112

Mathematically, we model context-specific prediction accuracy via the error term  $\mathbb{E}[(y_i - \hat{y}_i)^2 | \mathbf{c}_i]$ 

for phenotype  $y_i$  and prediction mean (or point prediction)  $\hat{y}_i = \mathbb{E}[y_i | \mathbf{c}_i]$  as a function of context  $\mathbf{c}_i$ 

for each individual *i* in the calibration dataset. We parametrize the impact of all contexts on prediction intervals in a joint model as  $\mathbb{E}[(y_i - \hat{y}_i)^2 | \mathbf{c}_i] = \exp(\mathbf{c}_i^{\mathsf{T}} \mathbf{\beta}_{\sigma})$  where  $\mathbf{c}_i$  denotes contexts

- prediction intervals in a joint model as  $\mathbb{E}[(y_i \hat{y}_i)^2 | \mathbf{c}_i] = \exp(\mathbf{c}_i^{T} \boldsymbol{\beta}_{\sigma})$  where  $\mathbf{c}_i$  denotes contexts including age, sex, socioeconomic factors and top principal components (denoting major axes of
- genetic ancestry; Methods).  $\beta_{\sigma}$  quantifies the unique impact of each context on variation of the

prediction interval accounting for other contexts (Methods). This approach is a generalization of

the context-free approach. Denoting prediction standard deviation (SD) as  $\hat{\sigma}_i = \sqrt{\exp(\mathbf{c}_i^{\mathsf{T}} \hat{\boldsymbol{\beta}}_{\sigma})}$ , 90%

prediction intervals can be derived as  $(\hat{y}_i - 1.645 \times \hat{\sigma}_i, \hat{y}_i + 1.645 \times \hat{\sigma}_i)$ .

#### 122 Widespread context-specific PGS accuracy in diverse populations

Although PGS accuracy has been shown to vary across selected traits and contexts<sup>5,10–12</sup>, its 123 pervasiveness remains unclear. We analyzed two large-scale biobanks in the UK and US (UK 124 Biobank and All of Us) comprising >600K individuals spanning a wide range of contexts. We 125 trained PGS for 72 traits in individuals previous annotated as "White British"<sup>28</sup> (WB) from UK 126 Biobank and evaluated these PGSs in independent testing data from UK Biobank and All of Us. 127 We focused on 11 contexts that span genetic ancestry, sex, age, and socio-economic factors 128 such as educational attainment (Methods). We used *relative*  $\Delta R^2$  to quantify the impact of context 129 to PGS accuracy defined as  $\frac{R_{top quintile}^2 - R_{bottom quintile}^2}{R_{all}^2}$ , where  $R_{[subset]}^2$  denotes  $R^2$  between PGS 130 and residual phenotype computed in a given range of the context variable (top/bottom guintile for 131 continuous contexts; binary subgroups for binary contexts). We found widespread context-132 specific PGS accuracies across all traits and contexts studied (Figure 3, S1 and S2, Table S1 and 133 S2; Methods). 134

#### 135 Context-specific accuracy in UK Biobank

All 72 traits had at least one context impacting their accuracies in UK Biobank data; 264 (out of 136 792) PGS-context pairs had significant variable accuracy ( $p < 0.05 / (72 \times 11)$ ; Methods). Overall, 137 genetic ancestry had the most widespread impact on PGS accuracy: 70 of 72 traits had significant 138 differences in PGS accuracy, with an average relative  $\Delta R^2$  of -46% between top and bottom PC1 139 quintiles (Figure S3). Socioeconomic contexts also significantly impacted PGS accuracy; PGS 140 accuracy significantly differed for 62 traits, with an average relative  $\Delta R^2$  of -23% between top and 141 bottom deprivation index guintiles. The direction of context's impact depended on the trait being 142 studied. For example, age significantly impacted 19 traits; rather than consistently increasing or 143 decreasing accuracy, an older age led to increased accuracies for 13 traits (e.g., high-density 144 lipoprotein cholesterol and white blood cell count in Figure 3; HDL and WBC) and to decreased 145 accuracies for 6 traits (e.g., low-density lipoprotein cholesterol; LDL). 146

The widespread context-specificity retained even when testing data was matched to the training 147 data by genetic ancestry (Figure 3). 22 (out of 72) PGSs had at least one context significantly 148 impacting their prediction accuracies; 43 PGS-context pairs had significant variable accuracy (p 149  $< 0.05 / (72 \times 11)$ ). We replicated previously reported variable PGS accuracy in WB individuals 150 for diastolic blood pressure, body mass index, education years across contexts of sex, age and 151 deprivation index<sup>10</sup>. As an example, LDL was significantly impacted by six contexts in WB 152 individuals, with age having the strongest impact (relative  $\Delta R^2$  was more than 100% between top 153 and bottom age quintiles). 154

Next, we studied the unique impact of each context on variable PGS accuracy within CalPred model that jointly accounts for all contexts (Methods, Figure 3cd). Context contribution to variable accuracy conditional on all other contexts was quantified with  $\beta_{\sigma}$ , where larger absolute  $\beta_{\sigma}$ indicated more substantial variation in accuracy along a context variable (Methods). In general, the effects of contexts to traits were largely independent. For example, both PC1 and deprivation index significantly impacted PGS accuracy for a range of traits in the joint model, indicating both

had a unique contribution to variable PGS accuracy. We also found examples showing otherwise: 161 the impact of "wear glasses" context on LDL accuracy can be explained by its correlation with age 162 (Figure S4), while other contexts independently contributed to variable LDL accuracy. These 163 results indicated the importance of jointly considering all measured contexts to correctly assess 164 the unique contribution of each context. We found that contexts including sex, age, income, and 165 deprivation index had comparable impact on accuracy as genetic ancestry (Figure 3ef). The 166 distribution of estimated effects of  $\beta_{\sigma}$  suggested predominantly higher prediction accuracy for 167 individuals with higher income and lower deprivation indices; this can be partly explained by 168 different context distribution PGS training data: WB individuals had higher income and lower 169 deprivation indices compared to the rest of the UK Biobank<sup>35</sup> (Figure S5). 170

#### 171 Context-specific accuracy in All of Us

We next turned to All of Us, a diverse biobank across the US comprising more than 160K 172 participants (Figure S3 and S6). Due to challenges in phenotype matching across biobanks, we 173 restricted the analysis to 10 traits and 11 contexts matching the UK Biobank analyses (Methods). 174 All traits had at least one context that impacted their accuracies (Figure 4, Table S3 and S4). 81 175 176 PGS-context pairs were significant when considering all individuals, and 49 PGS-context pairs were significant when restricting to individuals with self-reported race/ethnicity (SIRE) as "White" 177 ("White SIRE") ( $p < 0.05 / (12 \times 11)$ ; Methods). Prediction of cholesterol and LDL were similarly 178 impacted by a broad range of contexts. Prediction of education years was impacted by contexts 179 including age, BMI, employment, income, both when considering all individuals and considering 180 "White SIRE" sample, consistent with evidence that socioeconomic contexts influence PGS of 181 socio-behavioral traits such as education<sup>10,36,37</sup>. 182

Interestingly, socioeconomic contexts had greater impact on context-specificity in All of Us as compared to UK Biobank. For example, years of education context significantly impacted 9 out of 11 traits with average relative  $\Delta R^2$ =50%, as compared to 2 out of 71 traits with average relative  $\Delta R^2$ =0.2% in UK Biobank (averaging across traits other than education years itself). This may be explained by larger variation of education years in the US and/or education being more correlated with latent social determinants of health in the US as compared to the UK.

For completeness we also evaluated PGSs for height<sup>38</sup> and LDL<sup>39</sup> derived from multi-ancestry 189 meta-analyses from PGS Catalog<sup>40</sup> (Figure 4). We found that multi-ancestry PGSs did not 190 alleviate widespread context-specific accuracy. Higher income, education years, better 191 employment, or lower BMI predominately led to higher prediction accuracy across traits (Figure 192 4ef). We formally compared and determined an overall consistency for fitted  $\beta_{\sigma}$  coefficients 193 across populations and biobanks (Figure S7). We determined that variable  $R^2$  across contexts 194 was not solely driven by differences of phenotype variance in context strata: context-specific  $R^2$ 195 can result from differences in either phenotypic variance or PGS predictiveness, and the extent 196 attributed to either component varied by each context-trait pair (Figure S8). 197

# 198 CalPred yields calibrated context-specific prediction in simulations

Having shown that context-specificity of PGS accuracy is pervasive across traits and biobanks, 199 we next turned to CalPred, an approach to estimate context-specific prediction intervals 200 accounting for context- and trait-specific variable accuracy (Methods). We first evaluated CalPred 201 in simulations where prediction accuracy varies across contexts similar to real data<sup>5,6,10</sup> (Figure 5; 202 Methods). We assessed calibration of prediction intervals at both the overall level and within each 203 context subgroup (Methods). First, we showed that generic prediction intervals without context-204 specific adjustment had severe over-/under-coverage when evaluated within each context 205 subgroup stratified by PC1, age, or sex. As expected, biases of coverage tracked closely with 206 accuracy across contexts (Figure 5). Second, we showed that CalPred context-specific prediction 207 intervals that were allowed to vary with each individual's context were calibrated across contexts 208 (Figure 5). This was due to the incorporation of context-specific prediction accuracy in the interval 209 estimation. CalPred performance depended on calibration sample size with  $N_{cal}$  >500 for accurate 210 model fitting (Figure S9). Next, we investigated the impact of unmeasured context and found that 211 CalPred was not calibrated across subgroups of individuals defined by the unmeasured context. 212 In simulations where we included excessive contexts that did not impact prediction accuracy, 213 coverages of prediction intervals were associated with larger standard errors, highlighting the 214 importance of selecting an appropriate set of contexts in calibration (Figure S9). We also 215 determined that parameter estimations of  $\beta_{\sigma}$  were accurate when the model was correctly 216 specified and remained robust in model mis-specification scenarios (Figure S10). Overall, 217 simulation results demonstrated that CalPred is able to produce well-calibrated and context-218 specific prediction intervals when contexts are measured and present in the data, and highlighted 219 the importance of comprehensive profiling of relevant context information. 220

#### 221 CalPred yields calibrated context-specific predictions in real data

Next we applied CalPred to produce context-specific prediction intervals for a wide range of traits 222 across UK Biobank and All of Us. We start by showcasing LDL, an important risk factor of 223 cardiovascular disease<sup>39</sup>. Calibration by context is particularly important because accuracy of 224 225 predicting LDL was impacted by many contexts, with largest impact from age (Figure 3 and 4). We modeled the prediction mean using PGS together with age, sex, and genetic ancestry, and 226 modeled context-specific prediction intervals using the set of contexts investigated in Figure 3 227 and 4 (Methods). Accuracy of LDL prediction decreased with age ( $R^2$ =17% in youngest quintile 228 vs.  $R^2$ =11% in oldest quintile; Figure 6a). Generic prediction intervals were mis-calibrated with 229 coverage of 93% and 86% for youngest and oldest quintiles instead of the nominal level of 90%. 230 In contrast, context-specific prediction intervals had the expected 90% coverage across all 231 considered contexts. This resulted from varying prediction interval length by context, with a wider 232 interval compensating for lower prediction accuracy. For example, as the model estimated a 233 positive impact of age to prediction uncertainty ( $\beta_{\sigma}$ =0.15; p<10<sup>-30</sup>), individuals in youngest/oldest 234 age quintiles had average prediction standard deviation (SD) of 27.9 vs. 34.5 mg/dL (24% 235 difference: Figure S11: Methods). These findings were replicated in All of Us and in other traits 236 (Figure S12 and S13), where  $R^2$  varied across contexts and context-specific prediction intervals 237 achieved well-calibration. 238

Next, we sought to examine the joint contribution of all considered contexts to variable prediction
 SD (instead of separately considering age, PC1 or sex; Figure 6b). Context-specific accuracy was

more pronounced by ranking individuals by prediction SD accounting for impact of all contexts (prediction SD ranged approximately from 20 mg/dL to 45mg/dL; Figure 6b): we detected a 39% difference comparing individuals in bottom and top deciles of prediction SD (26.0 mg/dL vs. 36.3 mg/dL; Figure 6c; Figure S14 and S15). This implied that individuals in top prediction SD decile (characterized by contexts of male, increased PC1 and age; see LDL column in Figure 4c) needs to have their prediction interval widths increased by 39% compared to those in bottom decile.

Extending analysis accounting for all contexts to all traits in UK Biobank and All of Us, we 247 determined a widespread large variation of context-specific prediction intervals across traits 248 (Figure 7). Average differences between top and bottom prediction SD deciles across traits were 249 31% and 43%, respectively for UK Biobank and All of Us. The trait with the highest prediction SD 250 difference was the average mean spherical equivalent (avMSE), a measure of refractive error, 251 that was impacted the most by "wear glasses" context. Individuals who wore glasses had a much 252 higher PGS-phenotype  $R^2$  (9.6%) than those who did not (2.2%), likely due to the reduced 253 variation in avMSE phenotypes among individuals who did not wear glasses. Comparing across 254 the two datasets, BMI, LDL, and cholesterol were more heavily influenced by context than 255 average, while diastolic blood pressure and HDL were less impacted, suggesting trait-specific 256 susceptibility to context-specific accuracy. Notably, there were also cases where context-257 specificity of the same trait was drastically different across datasets. For example, prediction SD 258 differences for predicting education years was 90% in All of Us versus 8% in UK Biobank. This 259 disparity likely reflected the more diverse distribution of education years and other social 260 determinants of health in the US population sampled in All of Us, in line with results in Figures 2 261 and 3. Such differences between datasets also highlight that context-specificity can be population-262 specific and the need to consider unique characteristics of different populations in calibration. 263 Taken together, our findings emphasize the importance of incorporating context information into 264 PGS-based models when applied in diverse populations. 265

# 266 **Discussion**

Our work adds to the literature of PGS-based prediction as follows. First, we show that context-267 specific accuracy of PGS is highly pervasive across traits and biobanks with socioeconomic 268 contexts often having larger impact than genetic ancestry<sup>5,10,12,23,41</sup>. Second, we introduce CalPred 269 to estimate context-specific prediction intervals that maintain calibration for all individuals across 270 contexts. Third, we show using real and simulated data how differential prediction intervals can 271 be used to incorporate uncertainty in predictions. Although we focused primarily on PGS-based 272 prediction, our approaches are general and can incorporate any other factors. Fourth, we focused 273 on trait prediction as the main output of our approach motivated by genomic medicine applications. 274 As PGSs are increasingly applied to diverse populations, we find it imperative to incorporate the 275 context-specific accuracy into PGS downstream analyses to avoid bias against certain contexts 276 due to differential prediction accuracy, especially for contexts that are correlated with 277 socioeconomic status. CalPred provides а principled framework to quantify 278 generalizability/portability of a given PGS and represent individualized context-specific accuracy 279 to be leveraged in downstream analyses. The prediction intervals can be interpreted as a 280 personalized reference range accounting for each individual's contexts (including age, sex, and 281

genetic variation via PGS). Such personalized reference range may prove useful in identifying individuals with outlier lab values in a personalized and equitable fashion to prevent under-/overdiagonsis<sup>42</sup>.

The observation that distribution of PGSs differs across genetic ancestry continuum<sup>41</sup> motivates 285 methods that regress out effects of variables representing genetic ancestry from PGS distribution 286 to facilitate comparison across individuals locating at different positions in genetic ancestry 287 continuum<sup>43,44</sup>. However, such approaches may unintentionally remove true biological differences 288 of PGS distribution across genetic ancestry continuum (e.g., African Americans have reduced 289 neutrophil count that can be explained by the large effect of a single Duffy-null SNP<sup>45</sup>) as they do 290 not consider phenotype value distribution in calibration procedure; in addition, these approaches 291 cannot represent different standard errors in PGS predictions of individuals across genetic 292 ancestry continuum. Our method leverages a set of calibration data to properly adjust point 293 predictions across contexts according to true phenotype distribution. Compared to other existing 294 calibration methods<sup>34</sup>, our approach provides a framework to incorporate context information. 295

We note several limitations and provide future directions of our work. First, we focused on 296 modeling and analyzing quantitative traits in this work. Context-specific accuracies can be further 297 incorporated in modeling case-control status and absolute risk of diseases, perhaps by modeling 298 the underlying disease liability using methods proposed in this study. Second, we made several 299 modeling assumptions, including the linear relationship between error terms and contexts, as well 300 as quantile normalization procedure to phenotype values to fit in normal assumption of CalPred 301 model. Future work may leverage models with fewer assumptions and calibration dataset with 302 larger sample size to enable more flexible modeling. Third, CalPred requires calibration data that 303 matches in distribution with the target data, including both the distribution of contexts and their 304 305 effects to phenotypes (in terms of both prediction mean and variance). Otherwise, there may be bias in target samples that are underrepresented in the calibration data. The magnitude of bias 306 due to mismatch between calibration and target data in realistic scenarios needs to be empirically 307 examined in future work. As shown in our simulation studies, missing contexts will also limit proper 308 calibration of PGS along such contexts; this observation advocates standardized and 309 comprehensive profiling of contexts across biobanks to better quantify the role of contexts to PGS 310 accuracy, especially for those related to social-economic status, to prevent further exacerbation 311 of health disparity. Relatedly, these results indicate that GWAS data collecting process not only 312 313 needs to prioritize diversity in genetic ancestry, but also promote diversity across social-economic contexts, because PGS may be estimated with different precision in different social-economic 314 contexts. Fourth, CalPred prediction intervals will benefit from improved modeling of the prediction 315 mean; this may be achieved by more fine-grained modeling of prediction factors to capture more 316 phenotype variation (Supplementary Note). For example, sex-specific SNP-level effects can be 317 estimated from individual-level GWAS data<sup>16</sup> and CalPred coupled with sex-specific PGS is likely 318 to produce more precise, and shorter, prediction intervals. 319

# 320 Figures



Figure 1: Calibrated and context-specific prediction intervals via CalPred. (a) Calibration of 323 prediction intervals. We consider a subset of individuals with the same point prediction (shaded area 324 in the left panel, dashed horizontal line in the right panel). Each dot denotes an individual's phenotype 325 326 value. Intervals with proper-coverage cover the true phenotype at pre-specified probability of 90%; intervals with over-coverage are incorrectly wide; intervals with under-coverage are incorrectly narrow. 327 (b) Context-specific calibration of prediction intervals. We consider two subpopulations in different 328 contexts (e.g., female and male). Context 1 (blue dots) has lower prediction accuracy and therefore 329 wider variation around the mean, while context 2 (red dots) has higher prediction accuracy and 330 therefore narrower variation around the mean. Context-specific intervals vary by context, providing 331 intervals with proper coverage in each context. 332

<sup>321</sup> 322



333 334

Figure 2: Different approaches for prediction intervals of PGS-based models. All approaches 335 start with a set of predefined PGS weights derived from existing GWAS. (a) prediction intervals can 336 be calculated using analytical formula without calibration data. However, these intervals are not 337 guaranteed to be well-calibrated. (b) Generic calibration methods do not consider context information; 338 they produce generic prediction intervals that are constant across individuals. (c) Context-specific 339 calibration leverages a set of calibration data to estimate the impact of each context to trait prediction 340 accuracy; the estimated impact can then be used to generate prediction intervals for any target 341 individuals matching in distribution with calibration data. 342



343 344

Figure 3: Widespread context-specific PGS prediction accuracy in UK Biobank. (a-b) Heatmaps 345 for context-specific PGS accuracy for all and WB individuals. Each row denotes a context and each 346 column denotes a trait; the squared correlation between PGS and residual phenotype ( $R^2$ ) is shown 347 in parentheses. Heatmap color denotes the PGS-phenotype relative  $\Delta R^2$  (defined as  $\frac{R_{\text{group1}}^2 - R_{\text{group2}}^2}{R_{\text{all}}^2}$ ), 348 where  $R_{lsubsetl}^2$  represents  $R^2$  computed in a given range of the context variable. For continuous 349 contexts, relative  $\Delta R^2$  denote differences of top quintile minus bottom quintile; for binary contexts 350 (including sex, smoking, wear glasses, alcohol), relative  $\Delta R^2$  denote differences of male minus female. 351 smoking minus not smoking, wearing glasses minus not wearing glasses, drinking alcohol minus not 352 drinking alcohol (these orders were arbitrarily chosen). Numerical values of relative  $R^2$  differences are 353 displayed for PGS-context pairs with statistically significant differences (multiple testing correction for 354 all 10×11 PGS-context pairs in this figure;  $p < 0.05 / (10 \times 11)$ ). <sup>(\*)</sup> are displayed for PGS-context pairs 355 with nominally significant differences (multiple testing correction for 11 contexts; p < 0.05 / 11). (c-d) 356 Heatmaps for effects to prediction accuracy in CalPred model (estimated  $\beta_{\sigma}$ ). Colormaps were 357 inversed to those of (**a-b**) to reflect that positive  $\beta_{\sigma}$  corresponds to lower prediction accuracy and vice 358 versa. (e) Distribution of estimated  $\beta_{\sigma}$  in the CalPred model for each context across traits. (f) Number 359 of significantly impacted traits by each context ( $p < 0.05 / (72 \times 11)$ ). 360



361 362

Figure 4: Widespread context-specific PGS prediction accuracy in All of Us. (a-b) Heatmaps for 363 context-specific PGS accuracy for all and white SIRE individuals. Each row denotes a context and 364 each column denotes a trait; overall  $R^2$  is shown in parentheses. Heatmap color denotes relative  $\Delta R^2$ : 365 differences of top quintile minus bottom quintile for continuous contexts and difference of male minus 366 female for binary context of sex. Numerical values of relative  $R^2$  differences are displayed for trait-367 context pairs with statistically significant differences (multiple testing correction for all 12×11 trait-368 context pairs in this figure;  $p < 0.05 / (12 \times 11)$ ). '\*' are displayed for trait-context pairs with nominally 369 significant differences (multiple testing correction for 11 contexts; p < 0.05 / 11). (c-d) Heatmaps for 370 estimated  $\beta_{\sigma}$  in CalPred model. (e) Distribution of estimated  $\beta_{\sigma}$  in CalPred model for each context 371 across traits. (f) Number of significantly impacted traits by each context ( $p < 0.05 / (12 \times 11)$ ). 372



373 374

Figure 5: Simulation studies of CalPred. Simulations were performed to reflect scenarios where 375 individuals have variable prediction accuracy by genetic PC1, age, and sex. For each simulation, we 376 first trained a calibration model using a random set of 5,000 training individuals and then evaluated 377 resulting prediction intervals on 5,000 target individuals (Methods). (a) Prediction  $R^2$  between y and  $\hat{y}$ 378 in simulated data both at the overall level, and in each context subgroup. (b) Coverage of generic vs. 379 context-specific 90% prediction intervals evaluated in each context subgroup. Generic intervals were 380 obtained by applying CalPred without context information; context-specific intervals were obtained by 381 applying CalPred together with context information. (c) Average length of generic vs. context-specific 382 prediction standard deviation (SD) in each context. Each box plot contains  $R^2$ /coverage/average length 383 evaluated across 100 simulations (100 points for each box plot), the center corresponds to the median; 384 the box represents the first and third quartiles of the points; the whiskers represent the minimum and 385 maximum points located within 1.5× interguartile range from the first and third guartiles, respectively. 386





Figure 6: CalPred PGS calibration of LDL in UK Biobank. (a) Top panel: Prediction  $R^2$  between 389 phenotype and point predictions (incorporating PGS and other covariates) both at the overall level, 390 and in each subgroup of individuals stratified by context. Middle panel: Coverage of generic vs. 391 context-specific 90% prediction intervals evaluated in each context subgroup. Generic intervals were 392 obtained by applying CalPred without context information; context-specific intervals were obtained by 393 applying CalPred together with context information. Bottom panel: Average length of generic vs. 394 context-specific 90% prediction intervals in each context. Each box plot contains  $R^2$ /coverage/average 395 length across 30 random samples with each sample of 5,000 training individuals and 5,000 target 396 individuals (30 points for each box plot) (b) Ordered LDL prediction SD in unit of mg/dL. Gray lines 397 denote prediction SD obtained with random sample of 5,000 training and applied to 5,000 target 398 individuals. Red line denote prediction SD obtained from all individuals. (c) Box plots of results in (b) 399 from individuals of LDL prediction SD quantile of 0-10%, 45-55%, 90-100%; the center corresponds 400 to the median; the box represents the first and third quartiles of the points; the whiskers represent the 401 minimum and maximum points located within 1.5× interguartile ranges from the first and third guartiles. 402 respectively. 403



<sup>404</sup> 405

Figure 7. Variation of prediction standard deviation (SD) accounting for all contexts. Relative difference of prediction SD between top and bottom prediction SD deciles (90-100% vs. 0-10%) for all traits in UK Biobank (a) and All of Us (b). Traits are ranked by prediction SD. The difference is calculated with the median prediction SD within decile of individuals with highest prediction SD s<sub>d1</sub> and decile of individuals with lowest prediction SD s<sub>d10</sub> using  $\left(\frac{s_{d1}-s_{d10}}{s_{d10}}-1\right) \times 100\%$ .

# 411 Methods

# 412 Constructing calibrated and context-specific prediction intervals

We first provide an overview of CalPred framework. CalPred takes as input from pre-trained PGS 413 weights, genotype, phenotype and contexts to train a calibration model to generate calibrated and 414 context-specific prediction intervals for target individuals. We consider a calibration dataset with 415  $N_{cal}$  individuals. For each individual i=1, ...,  $N_{cal}$ , we have measured genotype vector  $\mathbf{g}_i \in \{0,1,2\}^M$ 416 with M SNPs, and phenotype  $y_i$ . With pre-trained PGS weights for a given trait  $\hat{\beta}_g \in \mathbb{R}^M$ , we 417 calculate the PGS for everyone in the calibration data with  $\mathbf{g}_{l}^{\mathsf{T}} \boldsymbol{\beta}_{g}$ . Each individual's PGS, together 418 with other contexts, including age, sex, genetic ancestry and other socioeconomic factors, 419 compose each individual i's contexts  $c_i$  (all '1' intercepts are also included). Phenotypes are then 420 modeled as 421

422 423

$$y_i = \mathcal{N}(\mu(\mathbf{c}_i), \sigma^2(\mathbf{c}_i)), i = 1, ..., N_{cal}$$
  
$$\mu(\mathbf{c}_i) = \mathbf{c}_i^{\mathsf{T}} \boldsymbol{\beta}_{\mu}, \sigma^2(\mathbf{c}_i) = \exp(\mathbf{c}_i^{\mathsf{T}} \boldsymbol{\beta}_{\sigma}).$$

424 There are two main components in the model

•  $\mu(\mathbf{c}_i) = \mathbf{c}_i^{\mathsf{T}} \mathbf{\beta}_{\mu}$  models the baseline prediction mean. This term is commonly used to predict phenotypes using PGS together with other contexts.

- $\sigma^2(\mathbf{c}_i) = \exp(\mathbf{c}_i^{\mathsf{T}} \mathbf{\beta}_{\sigma})$  models the context-specific variance of *y* around prediction mean. Differential prediction accuracy across contexts can lead to variable variance around prediction mean across contexts. The use of  $\exp(\cdot)$  is to ensure that the variance term >= 0.
- 431

Model parameters  $\beta_{\mu}$ ,  $\beta_{\sigma}$  can be estimated leveraging a set of calibration data using restricted maximum likelihood for linear model with heteroskedasticity<sup>46</sup> implemented in statmod R package<sup>47</sup>. Then individual-level predictive distribution  $\mathcal{N}(\hat{\mu}(\mathbf{c}_i) = \mathbf{c}_i^T \hat{\boldsymbol{\beta}}_{\mu}, \widehat{\sigma^2}(\mathbf{c}_i) = \exp(\mathbf{c}_i^T \hat{\boldsymbol{\beta}}_{\sigma}))$  can be generated for any target individual  $\mathbf{c}_i$  using the fitted  $\hat{\boldsymbol{\beta}}_{\mu}, \hat{\boldsymbol{\beta}}_{\sigma}$ . The corresponding  $\alpha$ -level prediction interval (e.g.,  $\alpha$ =90% for 90% prediction interval) is

437  $\left[\hat{\mu}(\mathbf{c}_{i}) - \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)\hat{\sigma}(\mathbf{c}_{i}), \hat{\mu}(\mathbf{c}_{i}) + \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)\hat{\sigma}(\mathbf{c}_{i})\right]$ , where  $\Phi^{-1}$  is the inverse cumulative 438 distribution function of a standard normal distribution (e.g.,  $\Phi^{-1}\left(1 - \frac{\alpha}{2}\right) = 1.645$  for 90% 439 prediction interval). Since we fit a simple linear model, the extent of parameter overfitting is 440 minimal with moderate sample size for calibration data (e.g.,  $N_{cal} > 1,000$  as validated in our 441 simulation studies).

- 443 **Quantile normalization for non-normal phenotype distribution.** In the above, we have 444 assumed that prediction intervals can be properly modeled as a Gaussian distribution, which may 445 not be always valid for every phenotype. To reduce the impact of this assumption to real data 446 analysis, we apply a transformation function  $Q(\cdot)$  to y with ranked based inverse normal 447 transformation such that Q(y) follow a normal distribution; Q(y) can then be modeled using the 448 methods described above. Fitted prediction intervals can then be transformed back into the 449 original y space using  $Q^{-1}(y)$ .
- 450

# 451 **Quantifying context-specific** *R*<sup>2</sup> **of PGS**

We quantify context-specific prediction accuracy ( $R^2$ ) of PGS, that is, to what extent PGS have variable prediction accuracy across contexts (including age, sex, genetic ancestry, proxies for lifestyle, socioeconomic contexts that can influence traits<sup>48</sup>). Accurate quantification of contexts contributing to variable prediction accuracy is important in constructing calibration model. In detail, for each pair of context and trait in a population, we calculated the prediction accuracy  $R^2$  between

PGS  $\hat{y}_i$  and covariate-regressed phenotypes  $y_i$  (phenotypes for each trait were regressed out of 457 age, sex, age\*sex and top 10 PCs; this adjustment is to better separate the contribution of PGS) 458 across each subgroup of individuals defined by contexts. We summarized results using relative 459 differences of  $R^2$  across context groups to baseline  $R^2$  calculated across all evaluated individuals 460 (relative differences between two classes for binary contexts; differences between top and bottom 461 quintiles for continuous contexts). We calculated the Spearman's  $R^2$  between point predictions 462 and covariate-regressed phenotypes  $R^2(\hat{y}, y)$  within each context subgroup. We also calculated 463 the baseline Spearman's  $R^2$  denoted as  $R^2_{all}$  across all individuals regardless of contexts. We summarized the results for each pair of trait and context using the "relative  $\Delta R^2$ " defined as 464 465  $\frac{R_{group1}^2 - R_{group2}^2}{2}$ . We assessed statistical significance of  $\Delta R^2$  across context subgroups by testing 466  $R_{all}^2$ the null hypothesis  $H_0: \Delta R^2 = 0$  using 1,000 bootstrap samples of  $\Delta R^2$  (in each bootstrap sample, the whole dataset was resampled with replacement and  $\Delta R^2$  were then re-evaluated). Statistical 467 468 significance was assessed using two-sided p-values comparing the observed  $\Delta R^2$  to the bootstrap 469 samples of  $\Delta R^2$ . 470 471

472 **Relationship between CalPred model and**  $R^2$ . Population-level metrics such as  $R^2$  can be 473 derived from this model as a function of  $\beta_{\sigma}$  and distribution of  $c_i$ . Suppose  $y = \hat{y} + e, e \sim$ 474  $\mathcal{N}(0, \exp(\mathbf{c}^{\mathsf{T}} \boldsymbol{\beta}_{\sigma}))$ , where  $y, \hat{y}, e$  denote the phenotypes, point predictions and residual noises, 475 respectively. We have

476

$$R^{2}(y,\hat{y}) = R^{2}(\hat{y} + e, \hat{y}) = \frac{\operatorname{Var}[\hat{y}]}{\operatorname{Var}[\hat{y}] + \operatorname{Var}[e]}$$

Holding  $\operatorname{Var}[\hat{y}]$  as fixed,  $R^2(y, \hat{y})$  is a function of  $\operatorname{Var}[e]$ , which is determined by the distribution of c and values of  $\beta_{\sigma}$ . This indicates a correspondence between  $\beta_{\sigma}$  and  $R^2(y, \hat{y})$ . Therefore, estimated  $\beta_{\sigma}$  can also be used as a metric to quantify context-specific accuracy (as used in Figure 3-4). While relative  $\Delta R^2$  is easier to interpret, it assesses the marginal contribution of each context separately and require binning for continuous contexts. Meanwhile,  $\beta_{\sigma}$  in CalPred model jointly account for all contexts in parametric regression, and therefore can quantify the unique distribution of each context to variable accuracy.

484

On the other hand, even with constant prediction interval length (constant Var[*e*]), variable  $R^2$  can still result from variable Var[ $\hat{y}$ ] across context subgroups. While CalPred focus on modeling Var[*e*] as a function of contexts to represent variable  $R^2$ , Var[ $\hat{y}$ ] can also change as a function of contexts in certain scenarios. For example, Var[ $\hat{y}$ ] can vary with contexts if  $\hat{y} = PGS \times \beta_{slope}$  and the slope  $\beta$  varies as a function of context. For example, ref.<sup>16</sup> has reported  $\beta_{slope}$  can be different across contexts. Such variable slope term can be handled by modeling variable slope terms in prediction mean  $\hat{y}$  (Supplementary Note).

## 493 **Real data analysis**

We analyzed a diverse set of contexts and traits in UK Biobank and All of Us (1) to quantify the extent of context-specific prediction accuracy and (2) to evaluate context-specific prediction intervals via CalPred.

497

492

Training polygenic score weights. Polygenic scores were trained on 370K individuals in UK Biobank that were assigned to "white British" cluster and 1.1M HapMap3 SNPs. For each trait, we performed GWAS using plink2 --glm with age, sex and top 16 PCs as covariates. Then we estimated PGS weights using snp\_ldpred2\_auto in LDpred2<sup>49</sup> with input of GWAS summary statistics and in-sample LD matrix. These estimated PGS weights were then applied to target individuals in both UK Biobank and All of Us to obtain individual-level PGS. To train polygenic score weights to be used for individuals from All of Us, we overlapped 1.2M SNPs in
 All of Us quality-controlled microarray data to 12M SNPs in UK Biobank imputed data to obtain a
 set of 0.8M SNPs present in both datasets. Then we trained and applied polygenic scoring weights
 using these shared SNPs in UK Biobank to All of Us individuals. This procedure helps improve
 accuracy of the polygenic score in All of Us by ensuring all SNPs that have non-zero weights to
 present in the data.

510

529

UK Biobank dataset. We analyzed 490K genotyped individuals (including both training and 511 target individuals). We used 1.1M HapMap3<sup>50</sup> SNPs in all analyses. All UK Biobank individuals 512 are clustered into sub-continental ancestry clusters based on top 16 pre-computed PCs (data-513 field 22009 in ref.<sup>28</sup> as in ref.<sup>6</sup>). This procedure assigned 410K individuals into "white British" 514 cluster. A random subset of 370K "white British" individuals to perform GWAS and estimate PGS 515 weights (see above); we trained PGS weights starting with individual-level data to avoid overlap 516 of sample between training and target data. For evaluation, we used the rest of 120K individuals 517 with genotypes, phenotypes and contexts (including individuals from both ~40K "White British" 518 individuals and ~80K other individuals). We focused on analyzing 72 traits with  $R^2$ >0.05 in 40K 519 WB target individuals and/or biological importance). We followed https://github.com/privefl/UKBB-520 PGS/blob/main/code/prepare-pheno-fields.R and ref.<sup>6</sup> to perform basic preprocessing for trait 521 values (e.g., log-transformation and clipping of extreme values). For each trait, we quantile 522 normalized phenotype values; when performing calibration, phenotype quantiles were calculated 523 based on calibration data and were then used to normalize target data. We analyzed 11 contexts 524 representing a broad set of socioeconomic and genetic ancestry contexts, including binary 525 contexts (sex, ever smoked, wear glasses, drinking alcohol) and continuous contexts (top two 526 PCs, age, BMI, income, deprivation index, and education years). We note that income and 527 education years have been processed into 5 guintiles in the original data of UK Biobank. 528

All of Us dataset. We analyzed 165K genotyped individuals with diverse genetic ancestry 530 contexts (microarray data in release v6). We retained 1.2M SNPs from microarray data after basic 531 quality control using plink2 with plink2 -- geno 0.05 -- chr 1-22 -- max-alleles 2 -532 -rm-dup exclude-all --maf 0.001. We used microarray data because it contains more 533 individuals and can be analyzed with low computational cost. All individuals with microarray data 534 were used in the evaluation. We analyzed 10 heritable traits, including height, BMI, WHR, diastolic 535 blood pressure, systolic blood pressure, education years, LDL, cholesterol, HDL, triglyceride; they 536 are straightforward to phenotype and have large sample sizes. Physical measurement 537 phenotypes were extracted from Participant Provided Information. Lipid phenotypes (including 538 LDL, HDL, TC, TG) were extracted following https://github.com/all-of-us/ukb-cross-analysis-539 demo-project/tree/main/aou workbench siloed analyses, including procedures of extracting 540 most recent measurements per person, and correcting for statin usage. For each trait, we quantile 541 normalized phenotype values; when performing calibration, phenotype quantiles were calculated 542 based on calibration data and were then used to normalize target data. We included age, sex, 543 age\*sex, and top 10 in-sample principal components as covariates in the model. We also quantile 544 normalized each covariate and used the average of each covariate to impute missing values in 545 covariates. We analyzed 11 contexts, including binary contexts (sex) and continuous contexts 546 (top two PCs, age, BMI, smoking, alcohol, employment, education, income, number of years living 547 in current address). 548

549

**Population descriptor usage.** We explain our usage choices of population descriptor, including the use of top two PCs to capture genetic ancestry/similarity and the use of "white British" in analyses of UK Biobank and "white SIRE" in analyses of All of Us. We use the top two PCs computed across all individuals in UK Biobank or in All of Us, respectively, to capture the continuous genetic ancestry variation in each dataset. While these two PCs provide major axes

of genetic variation (Figure S3), we acknowledge that top two PCs alone are not sufficient to fully 555 capture all variation in the entire population. The discretized PC1 and PC2 subgroups used in 556 Figure 2-6 is to enable calculation of population-level statistics such as  $R^2$  while we acknowledge 557 that the underlying genetic variation is continuous. In UK Biobank, we intended to analyze a set 558 of individuals with relatively similar genetic ancestry to perform GWAS and derive PGS. We used 559 a set of individuals previously annotated with "white British" that were identified using a 560 combination of self-reported ethnic background and genetic information having very similar 561 ancestral backgrounds based on results of the PCA<sup>28</sup>. In All of Us, we selected a set of individuals 562 with self-reported race/ethnicity (SIRE) being "white", to study how PGS have different accuracy 563 across environmental contexts in such a sample defined by SIRE. Noting that SIRE is not 564 equivalent to genetic ancestry, the contrast of results from UK Biobank and All of Us helps 565 understand how the genetic, nongenetic factors impact PGS accuracy in a group of individuals 566 defined by SIRE or genetic ancestry. 567

568

Evaluating context-specific prediction intervals. Recall that the prediction mean and standard 569 deviation are  $\hat{\mu}(\mathbf{c}), \hat{\sigma}(\mathbf{c})$  for a target individual with contexts c. We evaluate the prediction intervals 570 with regard to phenotypes y using metrics of 571

Prediction accuracy:  $R^2(\hat{\mu}(\mathbf{c}), y)$ . • 572

573

•

coverage of prediction intervals: evaluating  $\Pr\left\{y \in \left[\hat{\mu}(\boldsymbol{c}_{i}) - \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)\hat{\sigma}(\boldsymbol{c}_{i}), \hat{\mu}(\boldsymbol{c}_{i}) + \Phi^{-1}\left(1 - \frac{\alpha}{2}\right)\hat{\sigma}(\boldsymbol{c}_{i})\right]\right\} \approx \alpha, \text{ i.e., whether prediction}$ 574 intervals cover the true phenotypes with pre-specified probability of  $\alpha$ . 575

Both metrics are evaluated both at the overall level for all individuals, and for each subgroup of 576 individuals defined by contexts. 577

578

We generated and evaluated context-specific intervals in both UK Biobank and All of Us. For both 579 datasets, we fit a model to simultaneously model the mean and variance where the mean term 580 includes PGS, age, sex, age\*sex, top 10 PCs so that this matches the baseline model that are 581 commonly fitted, and the variance term includes age, sex, top 2 PCs, and other contexts of interest 582 for each dataset (as shown in Figure 3-4 for UK Biobank and All of Us). For each trait, we 583 performed the evaluation by repeatedly randomly sampling 5,000 individuals as calibration data 584 to perform the calibration, and 5,000 individuals as target data to perform the evaluation (as 585 described in "Constructing calibrated and context-specific intervals"). 586 587

#### Simulations assessing coverage of context-specific prediction intervals 588

We simulated PGS point predictions  $\hat{y}$  and phenotype values y to simulate traits with variable 589 prediction accuracy across genetic ancestry continuum, age, and sex. We started with real 590 contexts from 76K UK Biobank individuals not used for PGS training (see section "Real data 591 analyses"). We used 3 contexts (PC1, age, and sex) in simulations. We quantile normalized each 592 context so they had mean 0 and variance 1. Such simulations preserved the correlation between 593 contexts. Given these processed contexts, we simulated point predictions  $\hat{y}$  using a normal 594 distribution  $\hat{y} \sim \mathcal{N}(0,1)$ , and we simulated phenotypes y with: 595

$$y \sim \mathcal{N}(\hat{y}, \exp\left(\beta_{\sigma,0} + \sum_{c} \beta_{\sigma,c} \times c\right),$$

where  $\beta_{\sigma,0}$  denoted the baseline variance of y, and  $\beta_{\sigma,c}$  was the effect of context c to the variance 597 of y. " $\Sigma_c$ " enumerated over PC1, age, sex. This procedure simulated different variance of y around 598  $\hat{y}$  for individuals with different contexts, as observed in real data. 599

600

596

In details, we first selected  $\beta_{\sigma,0}$  such that  $R^2(y, \hat{y}) = 30\%$  for individuals with average contexts 601 (such that  $\sum_{c} \beta_{\sigma,c} \times c = 0$ ). We simulated data with variable variances: we set  $\beta_{\sigma,age} =$ 602

603  $0.25, \beta_{\sigma,\text{sex}} = 0.2, \beta_{\sigma,\text{PC1}} = 0.15$ . These parameters were manually chosen to roughly reflect the 604 observed variable  $R^2$  in real data. In each simulation, we randomly sampled  $N_{\text{cal}}$ =100, 500, 2500, 605 5000 individuals used for estimating the calibration model and  $N_{\text{test}}$  = 5000 individuals for 606 evaluating the predictions from the set of 76K individuals. New point predictions and phenotypes 607  $\hat{y}, y$  were simulated in each simulation. Then we quantified the prediction accuracy and coverage 608 of prediction intervals in these simulations.

# **Data availability**

610 UK Biobank individual-level genotype and phenotype data are available through application at 611 <u>http://www.ukbiobank.ac.uk</u>. AoU individual-level genotype and phenotype are available through 612 application at <u>https://www.researchallofus.org</u>.

613

## 614 **Code availability**

615 Software implementing CalPred and code for replicating analyses: 616 <u>https://github.com/kangchenghou/CalPred</u>.

617

# 618 Acknowledgements

We thank Molly Przeworski for helpful suggestions. This research was funded in part by the 619 National Institutes of Health under awards R01HG009120, R01MH115676, and U01HG011715. 620 This research was conducted using the UK Biobank Resource under applications 33127. We 621 thank the participants of UK Biobank for making this work possible. The All of Us Research 622 Program is supported by the National Institutes of Health, Office of the Director: Regional Medical 623 Centers: 1 OT2 OD026549; 1 OT2 OD026554; 1 OT2 OD026557; 1 OT2 OD026556; 1 OT2 624 OD026550; 1 OT2 OD 026552; 1 OT2 OD026553; 1 OT2 OD026548; 1 OT2 OD026551; 1 OT2 625 OD026555; IAA #: AOD 16037; Federally Qualified Health Centers: HHSN 263201600085U; Data 626 and Research Center: 5 U2C OD023196; Biobank: 1 U24 OD023121; The Participant Center: 627 U24 OD023176; Participant Technology Systems Center: 1 U24 OD023163; Communications 628 and Engagement: 3 OT2 OD023205; 3 OT2 OD023206; and Community Partners: 1 OT2 629 OD025277: 3 OT2 OD025315: 1 OT2 OD025337: 1 OT2 OD025276. In addition, the All of Us 630 Research Program would not be possible without the partnership of its participants. 631

#### 632 **References**

633

- Chatterjee, N., Shi, J. & García-Closas, M. Developing and evaluating polygenic risk prediction
   models for stratified disease prevention. *Nat. Rev. Genet.* **17**, 392–406 (2016).
- Torkamani, A., Wineinger, N. E. & Topol, E. J. The personal and clinical utility of polygenic risk
   scores. *Nat. Rev. Genet.* **19**, 581–590 (2018).
- 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).
- 4. Kullo, I. J. et al. Polygenic scores in biomedical research. Nat. Rev. Genet. 23, 524–532 (2022).
- 5. Martin, A. R. et al. Clinical use of current polygenic risk scores may exacerbate health
- 642 disparities. *Nat. Genet.* **51**, 584–591 (2019).
- 643 6. 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**, 373 (2022).
- Weissbrod, O. *et al.* Leveraging fine-mapping and multipopulation training data to improve
   cross-population polygenic risk scores. *Nat. Genet.* 54, 450–458 (2022).
- 8. Ruan, Y. *et al.* Improving polygenic prediction in ancestrally diverse populations. *Nat. Genet.* 54,
  573–580 (2022).
- Bitarello, B. D. & Mathieson, I. Polygenic Scores for Height in Admixed Populations. *G3* 10,
  4027–4036 (2020).
- 10. Mostafavi, H. *et al.* Variable prediction accuracy of polygenic scores within an ancestry group. *Elife* 9, (2020).
- In Jiang, X., Holmes, C. & McVean, G. The impact of age on genetic risk for common diseases.
   *PLoS Genet.* 17, e1009723 (2021).
- Hui, D. *et al.* Quantifying factors that affect polygenic risk score performance across diverse
   ancestries and age groups for body mass index. *Pac. Symp. Biocomput.* 28, 437–448 (2023).
- 13. Ding, Y. *et al.* Large uncertainty in individual polygenic risk score estimation impacts PRS-based
- risk stratification. *Nat. Genet.* **54**, 30–39 (2022).

- 14. Wray, N. R. *et al.* Pitfalls of predicting complex traits from SNPs. *Nat. Rev. Genet.* 14, 507–515
  (2013).
- 15. Ge, T., Chen, C.-Y., Neale, B. M., Sabuncu, M. R. & Smoller, J. W. Phenome-wide heritability
  analysis of the UK Biobank. *PLoS Genet.* **13**, e1006711 (2017).
- 16. Zhu, C. *et al.* Amplification is the primary mode of gene-by-sex interaction in complex human
- 664 traits. *Cell Genom.* **3**, 100297 (2023).
- Brown, B. C., Ye, C. J., Price, A. L. & Zaitlen, N. Transethnic genetic-correlation estimates from
  summary statistics. *Am. J. Hum. Genet.* **99**, 76–88 (2016).
- 18. Shi, H. *et al.* Population-specific causal disease effect sizes in functionally important regions
   impacted by selection. *Nat. Commun.* **12**, 1098 (2021).
- Patel, R. A. *et al.* Genetic interactions drive heterogeneity in causal variant effect sizes for gene
  expression and complex traits. *Am. J. Hum. Genet.* **109**, 1286–1297 (2022).
- 20. Weine, E., Smith, S. P., Knowlton, R. K. & Harpak, A. Tradeoffs in modeling context
- dependency in complex trait genetics. *bioRxiv* 2023.06.21.545998 (2023)
- 673 doi:10.1101/2023.06.21.545998.
- Wang, Y. *et al.* Theoretical and empirical quantification of the accuracy of polygenic scores in
   ancestry divergent populations. *Nat. Commun.* **11**, 3865 (2020).
- 22. Lambert, S. A., Abraham, G. & Inouye, M. Towards clinical utility of polygenic risk scores. *Hum. Mol. Genet.* 28, R133–R142 (2019).
- Ding, Y. *et al.* Polygenic scoring accuracy varies across the genetic ancestry continuum in all
   human populations. *bioRxiv* 2022.09.28.509988 (2022) doi:10.1101/2022.09.28.509988.
- 24. 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**, 104 (2022).
- 25. Wiley, L. K. *et al.* Building a vertically-integrated genomic learning health system: The Colorado
- 683 Center for Personalized Medicine Biobank. *bioRxiv* (2022) doi:10.1101/2022.06.09.22276222.
- 26. Belbin, G. M. et al. Toward a fine-scale population health monitoring system. Cell 184, 2068-

685 2083.e11 (2021).

- Abul-Husn, N. S. & Kenny, E. E. Personalized medicine and the power of electronic health
   records. *Cell* **177**, 58–69 (2019).
- Bycroft, C. *et al.* The UK Biobank resource with deep phenotyping and genomic data. *Nature*562, 203–209 (2018).
- 29. Ramirez, A. H. *et al.* The All of Us Research Program: Data quality, utility, and diversity.
- 691 Patterns (N Y) **3**, 100570 (2022).
- Wand, H. *et al.* Improving reporting standards for polygenic scores in risk prediction studies.
   *Nature* 591, 211–219 (2021).
- 31. Wei, J. *et al.* Calibration of polygenic risk scores is required prior to clinical implementation:
- results of three common cancers in UKB. J. Med. Genet. 59, 243–247 (2022).
- wan Houwelingen, H. C. Validation, calibration, revision and combination of prognostic survival
   models. *Stat. Med.* **19**, 3401–3415 (2000).
- 33. Van Calster, B. *et al.* Calibration: the Achilles heel of predictive analytics. *BMC Med.* **17**, 230
  (2019).
- 34. Sun, J. *et al.* Translating polygenic risk scores for clinical use by estimating the confidence
   bounds of risk prediction. *Nat. Commun.* **12**, 5276 (2021).
- 35. Schoeler, T. *et al.* Participation bias in the UK Biobank distorts genetic associations and
   downstream analyses. *Nat. Hum. Behav.* (2023) doi:10.1038/s41562-023-01579-9.
- 36. Selzam, S. *et al.* Comparing within- and between-family polygenic score prediction. *Am. J. Hum. Genet.* **105**, 351–363 (2019).
- 37. Okbay, A. *et al.* Polygenic prediction of educational attainment within and between families from
   genome-wide association analyses in 3 million individuals. *Nat. Genet.* 54, 437–449 (2022).
- 38. Yengo, L. *et al.* A saturated map of common genetic variants associated with human height.
- 709 *Nature* **610**, 704–712 (2022).
- 39. Graham, S. E. *et al.* The power of genetic diversity in genome-wide association studies of lipids.
- 711 *Nature* **600**, 675–679 (2021).

- 40. Lambert, S. A. *et al.* The Polygenic Score Catalog as an open database for reproducibility and
  systematic evaluation. *Nat. Genet.* 53, 420–425 (2021).
- 41. Martin, A. R. *et al.* Human demographic history impacts genetic risk prediction across diverse
  populations. *Am. J. Hum. Genet.* **107**, 788–789 (2020).
- 42. Van Driest, S. L. *et al.* Association between a common, benign genotype and unnecessary bone
- marrow biopsies among African American patients. *JAMA Intern. Med.* **181**, 1100–1105 (2021).
- 43. Hao, L. *et al.* Development of a clinical polygenic risk score assay and reporting workflow. *Nat.*
- 719 *Med.* **28**, 1006–1013 (2022).
- 44. Khera, A. V. et al. Whole-genome sequencing to characterize monogenic and polygenic
- contributions in patients hospitalized with early-onset myocardial infarction. *Circulation* **139**,
- 722 1593–1602 (2019).
- 45. 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).
- 46. Smyth, G. K. An Efficient Algorithm for REML in Heteroscedastic Regression. *J. Comput. Graph. Stat.* 11, 836–847 (2002).
- 47. Giner, G. & Smyth, G. K. statmod: Probability Calculations for the Inverse Gaussian Distribution.
- 728 *arXiv* [stat.CO] (2016).
- 48. Yousefi, P. D. *et al.* DNA methylation-based predictors of health: applications and statistical
  considerations. *Nat. Rev. Genet.* 23, 369–383 (2022).
- 49. Privé, F., Arbel, J. & Vilhjálmsson, B. J. LDpred2: better, faster, stronger. *Bioinformatics* 36,
  5424–5431 (2020).
- 50. The International HapMap 3 Consortium. Integrating common and rare genetic variation in
  diverse human populations. *Nature* 467, 52–58 (2010).