

# Evaluation of polygenic scoring methods in five biobanks shows larger variation between biobanks than methods and finds benefits of ensemble learning

Remo Monti<sup>1,2\*</sup>, Lisa Eick<sup>3\*</sup>, Georgi Hudjashov<sup>4</sup>, Kristi Läll<sup>4</sup>, Stavroula Kanoni<sup>5</sup>, Brooke N. Wolford<sup>6</sup>, Benjamin Wingfield<sup>7</sup>, Oliver Pain<sup>8</sup>, Sophie Wharrie<sup>9</sup>, Bradley Jermy<sup>3</sup>, Aoife McMahon<sup>7</sup>, Tuomo Hartonen<sup>3</sup>, Henrike Heyne<sup>1</sup>, Nina Mars<sup>3,10,11</sup>, Samuel Lambert<sup>12,14,15,17</sup>, Genes & Health Research Team, Kristian Hveem<sup>6,12</sup>, Michael Inouye<sup>13,14,15,16,17,18</sup>, David A. van Heel<sup>19</sup>, Reedik Mägi<sup>4</sup>, Pekka Marttinen<sup>9</sup>, Samuli Ripatti<sup>3,20,20</sup>, Andrea Ganna<sup>3,21</sup>, Christoph Lippert<sup>1,22,23,24\*</sup>  
\*Equal contribution

\*Corresponding author: christoph.lippert@hpi.de

## 10 Abstract

Methods to estimate polygenic scores (PGS) from genome-wide association studies are increasingly utilized. However, independent method evaluation is lacking, and method comparisons are often limited. Here, we evaluate polygenic scores derived using seven methods in five biobank studies (totaling about 1.2 million participants) across 16 diseases and quantitative traits, building on a reference-standardized framework. We conducted meta-analyses to quantify the effects of method choice, hyperparameter tuning, method ensembling and target biobank on PGS performance. We found that no single method consistently outperformed all others. PGS effect sizes were more variable between biobanks than between methods within biobanks when methods were well-tuned. Differences between methods were largest for the two investigated autoimmune diseases, seropositive rheumatoid arthritis and type 1 diabetes. For most methods, cross-validation was more reliable for tuning hyperparameters than automatic tuning (without the use of target data). For a given target phenotype, elastic net models combining PGS across methods (ensemble PGS) tuned in the UK Biobank provided consistent, high, and cross-biobank transferable performance, increasing PGS effect sizes ( $\beta$ -coefficients) by a median of 5.0% relative to LDpred2 and MegaPRS (the two best performing single

<sup>1</sup> Hasso Plattner Institute, University of Potsdam, Digital Engineering Faculty, Potsdam, Germany

<sup>2</sup> Max-Delbrück-Center for Molecular Medicine in the Helmholtz Association, Berlin Institute for Medical Systems Biology, Berlin, Germany

<sup>3</sup> Institute for Molecular Medicine Finland, FIMM, HiLIFE, University of Helsinki, Helsinki, Finland

<sup>4</sup> Estonian Genome Centre, Institute of Genomics, University of Tartu, Tartu, Estonia

<sup>5</sup> William Harvey Research Institute, Barts and the London School of Medicine and Dentistry, Queen Mary University of London, London, UK

<sup>6</sup> K. G. Jebsen Center for Genetic Epidemiology, Department of Public Health and Nursing, Faculty of Medicine and Health, Norwegian University of Science and Technology, Trondheim, Norway

<sup>7</sup> European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK

<sup>8</sup> Maurice Wohl Clinical Neuroscience Institute, Department of Basic and Clinical Neuroscience, Institute of Psychiatry, Psychology and Neuroscience, King's College London, London, United Kingdom

<sup>9</sup> Aalto University, Department of Computer Science, Espoo, Finland

<sup>10</sup> Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA, USA

<sup>11</sup> Stanley Center for Psychiatric Research and Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>12</sup> Levanger Hospital, Nord-Trøndelag Hospital Trust, Levanger, Norway

<sup>13</sup> Cambridge Baker Systems Genomics Initiative, Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK

<sup>14</sup> Cambridge Baker Systems Genomics Initiative, Baker Heart and Diabetes Institute, Melbourne, Victoria, Australia

<sup>15</sup> British Heart Foundation Cardiovascular Epidemiology Unit, Department of Public Health and Primary Care, University of Cambridge, Cambridge, UK

<sup>16</sup> Victor Phillip Dahdaleh Heart and Lung Research Institute, University of Cambridge, Cambridge, UK

<sup>17</sup> British Heart Foundation Cambridge Centre of Research Excellence, School of Clinical Medicine, University of Cambridge, Cambridge, UK

<sup>18</sup> Health Data Research UK Cambridge, Wellcome Genome Campus and University of Cambridge, Cambridge, UK

<sup>19</sup> Blizard Institute, Queen Mary University of London, UK

<sup>20</sup> Department of Public Health, University of Helsinki, Helsinki, Finland

<sup>21</sup> Massachusetts General Hospital and Broad Institute of MIT and Harvard, Cambridge, MA, USA

<sup>22</sup> Windreich Dept. of Artificial Intelligence & Human Health, Icahn School of Medicine at Mount Sinai, New York, NY, USA

<sup>23</sup> Hasso Plattner Institute for Digital Health at Mount Sinai, Icahn School of Medicine at Mount Sinai, New York, NY, USA

<sup>24</sup> Department of Diagnostic, Molecular, and Interventional Radiology, Icahn School of Medicine at Mount Sinai, New York, NY, USA

25 methods when tuned with cross-validation). Our interactively browsable online-results  
(<https://methodscomparison.intervene-geneticscores.org/>) and open-source workflow prspipe  
(<https://github.com/intervene-EU-H2020/prspipe>) provide a rich resource and reference for the  
analysis of polygenic scoring methods across biobanks.

## Introduction

30 Polygenic scores (PGS), also referred to as polygenic risk scores (PRS), have become a major  
application of genome-wide association studies (GWAS). PGS are constructed by scoring  
individuals based on their genotype, adding up effects of many genetic variants genome-  
wide. They can improve existing disease risk models that rely on family history and  
established biomarkers<sup>1-4</sup>, and individuals in the upper tail of the PGS distribution have an  
35 elevated disease risk similar to that caused by rare damaging monogenic mutations for some  
diseases<sup>5</sup>. PGS have received attention in areas ranging from disease prevention to clinical  
trials, owing to their wide applicability to personalized medicine<sup>6-9</sup>.

Various methods to derive PGS weights from GWAS summary statistics (effect sizes and  
their correlation structure) have been developed. These methods are of particular interest as  
40 they do not rely on access to individual-level data, which is typically restricted. Furthermore,  
the largest GWAS are meta-analyses, for which direct access to all individual-level source  
data is not feasible.

The construction of a PGS from summary statistics can be divided into two main stages: A  
public stage, that relies only on publicly available data and tools, and a private stage that  
45 requires access to individual-level target data, i.e., genotypes and phenotypes. The public  
stage uses variant correlation (linkage disequilibrium; LD) from reference panels that are  
matched in ancestry to the GWAS sample to adjust the marginal effect size estimates of  
genetic variants and derive the per-variant PGS weights. These adjustments include  
frequentist shrinkage<sup>10</sup>, Bayesian approaches<sup>11-15</sup>, or other strategies like thresholding, which  
50 depend on one or more hyperparameters (e.g., p-value thresholds, heritability estimates, or  
shrinkage parameters).

Many methods allow automatically setting suitable parameters without the use of phenotype  
data (we refer to this generally as automatic tuning). Alternatively, target data can be used to  
empirically determine hyperparameters based on, for example, cross-validation (CV). The  
55 adjusted variant effect sizes (PGS weights) are used in the private stage to score individuals  
based on their genotypes using a linear additive model, i.e., to calculate their PGS.

PGS method authors usually claim superior performance to other methods. However,  
comparisons are often limited to a small number of methods, traits, or target datasets.  
Furthermore, the input summary statistics used in those comparisons may not reflect the

60 properties of (messy) real-world data, especially those from meta-analyses. In practice, other factors also affect performance, e.g., ease of use and documentation. Few studies have compared a large number of PGS methods<sup>16-18</sup>. Yet, evaluation either only covered few traits in specialized cohorts<sup>17</sup> or was largely limited to within-biobank comparisons<sup>16,18</sup>.

The INTERVENE consortium<sup>19</sup> seeks to develop risk scoring methods that integrate PGS  
65 with other health-related information. For this reason, we compared summary-statistics-based PGS methods. Building on an updated version of the GenoPred suite originally introduced by Pain et al. that implements different PGS methods in a reference standardized framework<sup>16</sup>, we developed prspipe, a snakemake<sup>20</sup> workflow that runs seven polygenic scoring methods. A full evaluation including hyperparameter tuning with cross-validation was performed in the  
70 UK Biobank<sup>25</sup> (UKBB) and replicated in FinnGen<sup>21</sup>, Estonian Biobank<sup>22</sup> (EBB), the Trøndelag Health Study<sup>23</sup> (HUNT) and Genes & Health<sup>24</sup> (GNH). In total, we meta-analyzed performances for ten harmonized binary disease traits and six quantitative traits in two replicated ancestry groups European (EUR) and South Asian (SAS). Replication in multiple biobanks allowed us to estimate how much PGS effect sizes vary within biobanks (between  
75 methods) and how this compares to the variation between biobanks.

We publish our workflow, summary data, and PGS weights, allowing others to replicate analyses e.g., for methods comparisons or developing new polygenic scores from summary statistics. The results of this analysis are made available in a browsable online resource at <https://methodscomparison.intervenegeneticscores.org/>.

## 80 **Methods**

### **Participating Studies**

Data from five biobanks were considered: The UK Biobank<sup>25</sup>, FinnGen<sup>21</sup>, Estonian Biobank<sup>22</sup>, Trøndelag Health Study (HUNT)<sup>23</sup> and Genes & Health<sup>24</sup>. All biobanks independently performed genotyping, imputation, and variant quality control (Supplemental Methods).

### 85 **GWAS summary statistics selection and processing**

We selected summary statistics from the GWAS catalog for eight binary traits and for five continuous traits. Table 1 shows GWAS catalog study identifiers and traits. Where available, we directly used the pre-harmonized summary statistics provided by the GWAS catalog. For GWAS catalog studies GCST90013445<sup>26</sup> type 1 diabetes (T1D), GCST008972<sup>27</sup> (urate),  
90 GCST007954<sup>28</sup> glycated haemoglobin (HbA1c) and GCST004773<sup>29</sup> type 2 diabetes (T2D) we used the MungeSumstats R package<sup>30</sup> (version 1.0.1) to retrieve missing fields (e.g., variant positions). GWAS variants were matched to the HapMap3-1KG variants based on positions

and allele codes and renamed accordingly. Other quality control steps are flipping of variants to match the HapMap3-1KG reference, variant frequency filtering (>1%), removal of variants  
95 with invalid p-values (>1 or <0), ambiguous variants, variants with missing data, duplicate variants, or variants with sample size more than three standard deviations away from the median per-variant sample size (if available), as previously described<sup>16</sup>.

We selected GWAS studies with predominantly European ancestry discovery samples, because the evaluated biobanks primarily contain individuals of European ancestry. Because  
100 we use subsets of the UKBB for evaluation and tuning, we selected for studies with large sample sizes that preferably did not include the UKBB in the discovery sample. Yet, the selected summary statistics for Alzheimer's disease (AD) and height came from GWAS which included the UKBB-EUR sample. Therefore, we did not use the UKBB-EUR sample for tuning or evaluation in these phenotypes.

## 105 Reference genotype harmonization

We constructed our own definition of the HapMap3-variants<sup>31</sup> to avoid favoring one of the definitions used by the PGS methods. We retrieved HapMap3 variant rsIDs, and downloaded genotypes for the 1000 Genomes reference from  
[ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20140708\\_previous\\_phase3/v5\\_vcfs/](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20140708_previous_phase3/v5_vcfs/) (v5). We retrieved updated rsIDs for all 1000 Genomes variants using the  
110 Bioconductor SNPlocs.Hsapiens.dbSNP144.GRCh37 R-package<sup>32</sup> (the latest version for the GRCh37 genome-build at the time) based on GRCh37 variant positions and allele codes, and intersected them with the HapMap3 variants based on rsIDs. We then mapped these variants from GRCh37 to GRCh38 using liftOver<sup>33</sup> and retrieved rsIDs in that genome build too,  
115 based on location and allele codes, using the SNPlocs.Hsapiens.dbSNP151.GRCh38 R-package<sup>32</sup> (the latest version for the GRCh38 genome-build at the time). We retained variants with an allele frequency of at least 1% in any of the 1000 Genomes superpopulations. These variants (HapMap3-1KG, N = 1,330,821) form the basis for subsequent analyses.

The list of variants including GRCh37 (hg19) and GRCh38 (hg38) coordinates, rsIDs, and  
120 allele frequencies in the 1000 Genomes superpopulations is available on [https://github.com/intervene-EU-H2020/prspipe/blob/main/resources/1kg/1KGPhase3\\_hm3\\_hg19\\_hg38\\_mapping\\_cached.tsv.gz](https://github.com/intervene-EU-H2020/prspipe/blob/main/resources/1kg/1KGPhase3_hm3_hg19_hg38_mapping_cached.tsv.gz). Scripts to reproduce these steps are available as part of the prspipe workflow

([https://github.com/intervene-EU-H2020/prspipe/blob/main/workflow/rules/1kg\\_hm3\\_processing.smk](https://github.com/intervene-EU-H2020/prspipe/blob/main/workflow/rules/1kg_hm3_processing.smk)). The filtered and intersected 1000 Genomes genotypes are  
125

provided as a separate resource. Variants are further filtered when constructing polygenic scores, as described below.

### Target genotype harmonization

130 Target genotype data were intersected with the HapMap3-1KG variants based on positions and allele codes, renamed, and converted to PLINK1 format. The harmonized data served as input to all subsequent analyses involving target genetic data, i.e., ancestry reference matching and polygenic scoring. Target harmonization is part of the prspipe workflow and corresponding steps are defined in

[https://github.com/intervene-EU-H2020/prspipe/blob/main/workflow/rules/  
135 genotype\\_harmonization.smk](https://github.com/intervene-EU-H2020/prspipe/blob/main/workflow/rules/genotype_harmonization.smk).

### Binary disease phenotype harmonization

We used expert curated ICD-code based definitions<sup>21,34</sup> developed at FinnGen to define binary disease traits (referred to as endpoints). Individuals were counted as cases for a specific endpoint if they matched ICD-9- or ICD-10-code-based inclusion/exclusion criteria  
140 (13). The remaining (non-matching) individuals for that endpoint were counted as controls. All data used to define binary disease endpoints were registry based. breast cancer was only evaluated when the reported sex was female, and prostate cancer only when the reported sex was male.

For the UKBB, we considered both main (data-fields 41202 and 41203) and secondary (data-  
145 fields 41204 and 41205) ICD-9 and ICD-10 diagnosis codes derived from hospital inpatient admissions.

### Continuous trait definitions

For the UKBB, we used the following data-fields to define continuous traits: 50 for height, 21001 for body mass index (BMI), 30700 for creatinine, 30750 for HbA1c and 30880 for  
150 urate.

For GNH, we considered all instances where a continuous trait was measured per individual through their primary and secondary health records. We removed outliers based on a 6SD deviation per trait and calculated the mean value per trait per individual to use in the analysis. For HUNT, the latest value was chosen when continuous traits were measured at more than  
155 one baseline enrollment or sub-study screening over three recruitment waves since 1984. Standard quality assessment measures were taken across variables and are described at the HUNT Databank (<https://hunt-db.medisin.ntnu.no/hunt-db/variablelist>). BMI was defined by height and weight measured at screening. High-density lipoprotein (HDL) and creatinine

were measured from serum in non-fasting individuals. HbA1C was measured in mmol/mol  
160 according to The International Federation of Clinical Chemistry and Laboratory Medicine (IFCC)  
standard.

For Estonian Biobank, the earliest value available were chosen for BMI and height, as some  
individuals are repeatedly measured. Height values larger than 260cm or smaller than 100cm  
were omitted. Similarly, BMI values less than 10kg/m<sup>2</sup> or larger than 200 kg/m<sup>2</sup> were  
165 discarded. Metabolic profiles for HDL and creatinine were obtained with NMR for a random  
subset of Estonian biobank (n=10681).

For each biobank in which creatinine measurements were available, we calculated the  
estimated glomerular filtration rate (eGFR) based according to the diet in renal disease study  
equation<sup>35</sup>, as follows:

170 
$$eGFR = \alpha \times S_{cr}^{-1.154} \times age^{-0.203} \times \sigma$$

Where  $\alpha$  is 30849 if creatinine was measured in  $\mu\text{mol/l}$ , or 175 if measured in  $\text{mg/dl}$ ,  $S_{cr}$  is the  
serum creatinine measurement and  $\sigma$  is 0.742 if the reported sex is female, or 1 if sex is male.  
We did not include the multiplier for “ethnicity” as we only perform comparisons within  
ancestry-matched populations. During evaluation, all continuous traits are standardized to  
175 mean 0 and unit standard deviation.

### Polygenic score weight derivation

We derived polygenic scoring weights with pT+clump, lassosum, PRSCs, SBayesR (robust  
parameterization), LDpred2 and DBSLMM using the settings described previously<sup>16</sup>. For  
MegaPRS, we used the author-recommended BLD-LDAK heritability model and specified  
180 “--model mega” to fit many different scores with different tools (lasso, bolt, ridge, bayesr)  
and included the HLA region, as recommended. Software versions and sources for each tool  
are listed in 14.

Besides letting methods determine suitable hyperparameters based on the summary statistics  
alone (automatic tuning), we generated scores over grids of hyperparameters for target-data-  
185 based tuning with 10-fold cross-validation (see below).

We used European ancestry LD reference panels for all analyses, as the selected GWAS were  
performed in majority European ancestry samples. In contrast to Pain et al., we use PGS  
method author-provided LD-references for DBSLMM, lassosum, LDpred2, SBayesR and  
PRSCs. DBSLMM and lassosum LD-references are based on the 1000 Genomes data. For  
190 LDpred2, SBayesR, and PRSCs they are based on UKBB data. We use the 1000 Genomes  
EUR-subset to calculate LD when running pT+clump and MegaPRS. Scripts to download

PGS method software and data are part of the prspipe workflow. The workflow uses GenoPred scripts to generate PLINK2-compatible scoring files.

### Ancestry matching and genetic outlier removal

195 Rather than directly inferring genetic ancestry, we score individuals according to their similarity with groups defined in the 1000 Genomes reference<sup>36</sup>. We use GenoPred Ancestry\_identifier.R to project target genetic data into the 1000 Genomes genetic principal component space, and match individuals to one of the five 1000 Genomes superpopulations (AFR, AMR, EAS, EUR, SAS). Both the target and 1000 Genomes genotype data are filtered  
200 to variants available in both samples (subset of HapMap3-1KG) with allele frequencies above 5%, missingness below 2%, and variants that do not violate Hardy Weinberg equilibrium ( $p < 1e-6$ ). Regions of long LD are excluded<sup>37</sup> and variants are LD-pruned based on the 1000 Genomes reference (using PLINK “--indep-pairwise 1000 5 0.2”). Genetic PCs are derived in the 1000 Genomes based on the filtered variants, and an elastic net classifier is fit with 5-fold  
205 cross validation to place individuals into one of the five groups based on 100 PCs. Target genotype data are projected into the same PC space using PLINK, and the classifier is used to predict the most matching superpopulation for all individuals.

Additionally, we used GenoPred Population\_outlier.R to remove extreme outliers within the assigned ancestry-matched groups in the target data based on the first eight genetic principal  
210 components constructed within those groups. We used the same variant filters described above (except that LD-pruning was now performed within the target data), calculate genetic PCs within the assigned groups, and define up to ten centroids in the PC space using R NbClust<sup>38</sup> (distance = ‘euclidian’, method = ‘kmeans’). For each centroid, the Euclidian distance of individuals to the center is calculated and those with distances that are larger than  
215 the 75<sup>th</sup> percentile + 30 IQR are removed (i.e., extreme outliers). The UKBB was the only biobank that had more than one ancestry group well-represented. Our analyses focus on the replicated groups EUR (UKBB, EBB, HUNT, FinnGen) and SAS (UKBB, GNH).

### Polygenic scoring

We performed polygenic scoring for scores derived by single methods with PLINK2 using  
220 GenoPred Scaled\_polygenic\_scorer\_plink2.R. Polygenic scoring is part of the prspipe workflow. For the evaluation of the ensemble PGS, we performed scoring with PLINK2. In both cases, missing genotypes are imputed using the 1000 Genomes matched superpopulation allele frequencies as previously described<sup>16</sup>.

## Hyperparameter tuning and ensemble PGS

225 For the methods that generate scores over a range of hyperparameters (pT+clump, lassosum, PRSs, LDpred2) we used 10-fold cross validation to select the score with the largest correlation with the trait, as described previously<sup>16</sup>. Where available, we included scores produced by methods' automatic settings in the selection process. We perform cross-validation using 80% of the UKBB EUR data and retain 20% for evaluation (we used  
230 different subsamples for each trait in order to perform stratified sampling).

For pT+clump, we define the score given by the p-value threshold of 1e-8 as the automatically tuned score. SBayesR and DBSLMM only use automatic settings, i.e., they produce just a single set of weights and are not tuned with CV.

To fit the ensemble PGS, we include all scores from all methods across hyperparameters, and  
235 use 10-fold CV to determine suitable shrinkage parameters for an elastic net model combining the different scores with the caret R-package<sup>39</sup> (which relies on glmnet<sup>40</sup>). For tuning of the ensemble PGS, we used non-nested scores for pT+clump, i.e., scores with disjoint variant sets corresponding to 10 p-value bins. These steps were performed with GenoPred Model\_builder\_V2.R and are part of the prspipe workflow.

240 To score other biobanks with the UKBB ensemble PGS, we generated PLINK2-compatible scoring files by multiplying the PGS weights of every variant with their corresponding weights in the ensemble PGS model and adding them (yielding a single weight for each variant).

## Performance evaluation within biobanks

245 All PGS were standardized to mean zero and unit standard deviation within biobanks and ancestries for performance evaluation. We calculated the following metrics for binary disease traits:  $\beta$  coefficients, i.e., the change in log-odds ratios per PGS standard deviation, the change in odds ratio per PGS standard deviation ( $OR = \exp(\beta)$ ), fraction of variance explained on the observed scale ( $r^2_{obs}$ ) and the area under the receiver operating characteristic  
250 curve (AUROC). The variance explained on the liability scale ( $r^2_{liab}$ ) was calculated from  $r^2_{obs}$ <sup>41</sup> using the median prevalence within ancestries as the population prevalence estimate (1). We retrieved DeLong 95% confidence intervals<sup>42</sup> for the AUROC using the ci.auc-function in the pROC R-package. Confidence intervals for  $r^2_{obs}$  were derived from 1000 bootstrap samples of  $r_{obs}$  (the Pearson correlation on the observed scale) for binary traits.

255 For continuous traits, we calculated  $\beta$  coefficients, i.e., the change in standard deviations of the trait per standard deviation of the PGS and the fraction of variance explained ( $r^2_{obs}$ ).



When comparing two effect sizes of scores  $\beta_{a,i}$  and  $\beta_{b,i}$  within biobank “i”, we use the two-sided z-test and adjust for the correlation between scores, with test statistic:

$$z = \frac{\beta_{a,i} - \beta_{b,i}}{\sigma}$$

260 ,where

$$\sigma = \sqrt{\sigma_{a,i}^2 + \sigma_{b,i}^2 - 2\rho_{(a,b),i}\sigma_{a,i}\sigma_{b,i}}$$

With  $\sigma_{a,i}$  and  $\sigma_{b,i}$  denoting the standard deviations of  $\beta_{a,i}$  and  $\beta_{b,i}$ , respectively and  $\rho_{(a,b),i}$  denoting the correlation between scores “a” and “b” measured in biobank “i”.

### Data Exclusions before meta-analyses

265 None of the scores evaluated in GNH-SAS for T1D reached nominal significance ( $p < 0.05$ , two-sided z-test) for association with the endpoint (1), and all effect sizes were close to zero. We removed these data from further analysis. We found strongly reduced effect sizes of scores for HbA1c in GNH-SAS compared to the UKBB-SAS (2, 1) and decided not to include these data for meta-analyses (it was unclear if reduced performance was due to a phenotyping issue). We found low effect sizes compared to other biobanks for T1D in HUNT (1), determined it was likely due to a phenotyping issue, and excluded those data from meta-analyses.

### Meta-analysis for methods comparisons

275 All Meta-analyses were performed in R (version 4.1.1) with the metafor package (rma.mv function, version 3.8-1), using the V-argument to account for the dependence of effect-sizes within biobanks (see below), and models were fit with REML.

We meta-analyzed the  $\beta$  coefficients of scores across biobanks within ancestries and traits using meta-analytic mixed effects models. The observed  $\beta$  coefficients are modelled as follows:

280 
$$\beta_{s,b} = \alpha \mathbf{w}_s + \zeta_b + \epsilon_{s,b}$$

Where  $\beta_{s,b}$  is the observed coefficient for PGS “s” in biobank “b” for a specific trait.  $\beta_{s,b}$  is modelled as a combination of fixed effects (moderators) with realizations  $\mathbf{w}_s$  and parameters  $\alpha$  (bold characters indicate vectors) and two error terms: the sampling error  $\epsilon_{b,s}$  and a biobank-specific random intercept  $\zeta_b$  (shared by all observed coefficients in that biobank).  $\tau_{\text{biobank}}^2 = \text{var}(\zeta)$  is the random effect, where  $\text{var}(\zeta)$  denotes the variance of the biobank-specific random intercepts  $\zeta$ .

285 For every trait, we meta-analyzed up to 13 PGS in the same model. PGS-choice is modelled with the fixed effects, i.e.,  $\mathbf{w}_s$  only contains a single non-zero entry of 1 indicating which PGS

“s” produced  $\beta_{s,b}$ . With this parameterization, parameters in  $\alpha$  directly correspond to the meta-analyzed effect sizes for the different PGS after inverse variance weighting, and the formula above can equivalently be written as:

$$\beta_{s,b} = \beta_{s^*i} + \zeta_b + \epsilon_{s,b,i}$$

Where  $\beta_{s^*}$  is the average effect size for score “s” across all biobanks “\*”. To test whether two meta-analyzed effect sizes are significantly different, we compare parameters  $\alpha_a$  and  $\alpha_b$  with  $H_0: \alpha_a - \alpha_b = 0$  using the z-test. We also report results for the t-test (produced by the anova function applied to metafor rma.mv objects) in 5-6.

We retrieved 95% likelihood-based confidence intervals for  $\tau^2_{\text{biobank}}$  using the confint function (all values are reported in 4). We further report meta-analyzed AUROC,  $r^2_{\text{obs}}$ , and  $r^2_{\text{liab}}$  values produced by weighting studies by their effective sample size<sup>43</sup> in 5-6.

### 300 Method ranking

To rank methods across traits, we considered just the traits for which CV-tuning and ensemble PGS were available (i.e., all except height and AD), and ranked scores based on their meta-analyzed effect sizes  $\beta_{s^*}$  (see definition above). To avoid counting scores produced by the same summary statistics twice for eGFR/CKD and urate/gout (e.g., for the ranks shown in Figure 3A), we applied the following rule: If the continuous phenotype was available in the same number of biobanks as its binary counterpart, we used the continuous phenotype (higher power), otherwise we used the binary phenotype (larger target diversity). This led to consideration of eGFR for SAS, CKD in EUR, urate for SAS, and gout in EUR. We applied the same reasoning when calculating mean and median values of method performances across all traits.

### 3-level meta-analytic random effects models

For the 3-level meta-analysis, the observed effect-sizes are modelled as follows:

$$\beta_{b,m} = \mu_\beta + \zeta_{(2)b,m} + \zeta_{(3)b} + \epsilon_{b,m}$$

Where  $\beta_{b,m}$  is the observed  $\beta$ -coefficient for method “m” in biobank “b”,  $\mu_\beta$  is the mean of the distribution of true effect sizes across biobanks and methods,  $\zeta_{(2)b,m}$  is the within-biobank random intercept due to the choice of method (level 2) and  $\zeta_{(3)b}$  is the random intercept due to target biobank (level 3, shared by all observations in biobank). The estimated parameter  $\tau^2_{\text{biobank}} = \tau^2_{(3)} = \text{var}(\zeta_{(3)})$  quantifies the heterogeneity of effect sizes due to target biobank, and  $\tau^2_{\text{method}} = \tau^2_{(2)} = \text{var}(\zeta_{(2)})$  quantifies the heterogeneity of effect sizes due to choice of method within biobank<sup>44</sup>. In contrast to the model introduced in the previous section, method effects are considered nested within biobanks (independent between biobanks).

All models were fit using restricted maximum-likelihood with the metafor package in R (rma.mv function), using the V-argument to account for the dependence of effect-sizes measured within the same biobanks (see below). We retrieved 95% likelihood-based confidence intervals for  $\tau^2_{\text{biobank}}$  and  $\tau^2_{\text{method}}$  using the confint function. To calculate  $I^2_{\text{biobank}} = I^2_{(3)}$  and  $I^2_{\text{method}} = I^2_{(2)}$ , we used the implementation provided in dmetar<sup>44</sup> (var.comp/mlm.variance.distribution function, <https://github.com/MathiasHarrer/dmetar/blob/master/R/mlm.variance.distribution.R>, commit 21bde652cbae5677b56b0ff848eb96c9bea877d8) based on the three-level extension of the  $I^2$  metric<sup>45</sup>.  $I^2_{\text{biobank}}$  captures the fraction of the overall variance in effect sizes (including sampling error) attributable to biobank (level 3) and  $I^2_{\text{method}}$  captures the fraction of the overall variance in effect sizes attributable to methods within biobank (level 2).

### Accounting for dependent effect sizes in meta-analytic models

Within each biobank, ancestry, and trait, we calculated pairwise correlations between polygenic scores based on up to 50,000 randomly sampled individuals. We use the resulting score-score correlation matrix  $R_b$  (where “b” indicates the biobank) to estimate  $V_b$ , the variance-covariance matrix capturing the dependency of errors of effect size estimates for biobank “b”:

$$V_b = S_b R_b S_b$$

Where  $S_b$  is a diagonal matrix containing the standard errors of the estimated effect sizes corresponding to the rows/columns of  $R_b$ . The effect sizes measured in different biobanks are considered independent, therefore, the full matrix  $V$  supplied to the rma.mv function is a block diagonal matrix containing all  $V_b$  for the different biobanks  $b$  from 1 to  $n$  on the diagonal:

$$V = \begin{matrix} & V_1 & 0 & 0 \\ 0 & & \ddots & 0 \\ & 0 & 0 & V_n \end{matrix}$$

Where 0 denotes a matrix of zeros with the same shape the different  $V_b$ .

### Calculating PGS variance in the HLA-region

For the phenotypes T1D and rheumatoid arthritis (RA), we scored individuals in the 1000 Genomes EUR subset using either all PGS variants, or only variants contained in the HLA region (defined as the interval 28,000,000-34,000,000 on chromosome 6). The fraction of variance was then computed by dividing the variance of HLA-only PGS by that of the full PGS.

## Results

### 355 Prspipe workflow and experimental setup

We created a snakemake workflow prspipe to run different polygenic risk scoring methods based on GWAS summary statistics. Prspipe makes it possible to automate the within-biobank analyses from Pain et. al.<sup>16</sup> based on the GenoPred suite of scripts (<https://github.com/intervene-EU-H2020/GenoPred>). Notable differences include an updated  
360 set of methods (p-value thresholding and clumping (pT+clump), lassosum<sup>10</sup>, PRScs<sup>11</sup>, LDpred2<sup>13</sup>, DBSLMM<sup>14</sup>, SBayesR<sup>12</sup> (robust parameterization) and MegaPRS<sup>15</sup>), the use of LD reference panels provided by the methods' authors, and software managed partially with containers. We used prspipe to derive PGS weights using both methods' automatic settings (auto), and grids of hyperparameters (MegaPRS, LDpred2, PRScs, lassosum, and  
365 pT+clump)<sup>16</sup>. For the baseline method pT+clump, we considered the score with the most stringent p-value threshold ( $p < 1e-8$ , i.e., keeping only highly significant variants) as the automatically tuned score.

The workflow defines steps to set up PGS methods, download and process summary statistics from the NHGRI-EBI GWAS Catalog<sup>46</sup>, run PGS methods (i.e., the derivation of PGS  
370 weights), target genotype harmonization, ancestry matching based on the 1000 Genomes superpopulations<sup>36</sup>, and target polygenic scoring with PLINK2<sup>47</sup>. As PGS performance depends on the genetic similarity of the target and GWAS samples<sup>48</sup>, performance evaluation is stratified according to the matched superpopulation. Using CV, elastic net models combining scores from different methods are fit (ensemble PGS), and the best single PGS  
375 weights are selected for each method (hyperparameter tuning)<sup>16</sup>.

We applied this workflow to 14 sets of summary statistics from the GWAS Catalog to derive PGS and predict six continuous traits and 10 binary disease traits derived from harmonized ICD-code-based definitions<sup>21</sup> (Methods, Table 1). Our main analyses focus on the two replicated ancestry-reference-matched superpopulations: EUR and SAS. The number of cases  
380 used for performance evaluation across biobanks ranged from 5,384 (T1D) to 81,487 (type 2 diabetes; T2D) for EUR, and 60 (RA, available in GNH only) to 8,696 (T2D) for SAS ancestry matched target data. The total sample size for performance evaluation for continuous traits ranged from 85,973 (urate, available in UKBB only) to 524,056 (height) for EUR target data, and 13,572 (urate) to 43,197 (height) for SAS target data (Table 2).

385 Using 80% of the UKBB EUR target data (training set), we selected the best performing weights for each method and fit ensemble PGS (full workflow). PGS weights were shared with other biobanks, in which we still performed target data harmonization, ancestry matching, and polygenic scoring steps needed for performance evaluation (Figure 1).

### Browsable results, meta-analysis and ranking

390 As outlined in Figure 2 for T2D, we calculated PGS effect sizes (Figure 2A) for continuous and binary traits across all target biobanks and ancestries (Table S1-3, Figures S1-5) and performed mixed model meta-analyses within ancestry groups to determine the best performing PGS (the one with the largest effect size) for each trait across biobanks (Figure 2B-C, 4-6). Additionally to scores produced by single methods, we evaluated the UKBB-tuned ensemble PGS in other biobanks after projecting them back to the variant-level  
395 (Methods). For each trait, we meta-analyzed  $\beta$ -coefficients (i.e., the change in the trait per PGS standard deviation, on the log-odds scale for binary traits, in standard deviations for continuous traits) of up to 13 PGS corresponding to different tuning types (auto or CV) for seven methods and the UKBB-EUR-tuned ensemble PGS (6).

400 Other than the ensemble PGS, we found that CV-tuned PGS from LDpred2 and MegaPRS ranked highly across traits (Figure 3A-B). The median relative increase in PGS effect size over CV-tuned pT+clump was 29.2% for CV-tuned LDpred2 (mean 30.9%±10.3sd, N=12, EUR) and 29.9% for CV-tuned MegaPRS (mean 31.2%±12.6 sd, N=12, EUR), showing overall comparable performance (the median relative difference between the two was 0.1% in  
405 favor of MegaPRS).

Scores produced by automatic tuning appeared overall less reliable, especially for LDpred2 (see Discussion) and SBayesR (as previously described<sup>16</sup>). Although automatic tuning typically outperformed the baseline method pT+clump (even when the latter was tuned with CV), we observed seemingly non-systematic cases of reduced relative performance (e.g.,  
410 SBayesR for urate/gout, DBSLMM for Alzheimer's disease, or LDpred2 for HbA1c or RA) (11). MegaPRS was the best automatically tuned method (median 23.3% relative increase over CV-tuned pT+clump, mean 27.4%±15.9 sd, EUR), yet PGS effect sizes were comparatively low for some continuous traits (e.g., BMI, HDL, or height, 11).

Effect size differences between the top PGS by single methods were mostly not significant  
415 (7, FWER ≤ 0.05, two-sided z-test). We also provide these data on the level of individual biobanks (9-10), revealing that the best single method for a given trait was not necessarily consistent between biobanks.

## UKBB-tuned ensemble PGS outperforms other methods

The ensemble PGS ranked favorably for all traits in EUR- and SAS-matched target data (Figure 3A-B) except for T1D in EUR (driven by lower performance in FinnGen) and stroke in SAS (the trait with the overall lowest performance). For EUR target data, effect sizes were significantly greater than those of all other PGS for 6/9 binary and 5/5 continuous traits (FWER  $\leq 0.05$ , two-sided z-test) and the largest overall in 13/14 traits for which we fit ensemble PGS. These results stood out compared to those for single methods, which did not produce a consistent best method and for which differences were mostly not significant. Compared to the best single methods, the median relative increase in effect size was 3.7% over CV-tuned LDpred2 and 4.5% over CV-tuned MegaPRS for binary disease traits (N=9). Median relative increases for continuous traits were larger (5.2% and 7.9%, respectively, N=5). When measured in terms of variance explained, relative differences were larger. We observed median relative increases of 7.4% and 9% for binary traits (liability scale) and 10.7% and 16.1% for continuous traits, respectively (8, 7). Similar trends were observed for SAS target data, with the ensemble PGS having the largest effect size in 12/13 traits, albeit its effect size was only significantly larger than all others for continuous traits urate, eGFR and HDL (FWER $\leq 0.05$ , two-sided z-test). We report relative effect sizes of all methods relative to the ensemble PGS in Table 3.

## CV-tuning increases PGS robustness

Hyperparameter tuning with cross-validation using the UKBB EUR data was often beneficial and rarely harmful when evaluated on EUR target data (Figure 3C). CV-hyperparameter tuning strongly increased effect sizes in a subset of traits for specific methods, rather than providing large benefits across traits (12-13). pT+clump benefited most from CV-tuning when evaluated on EUR target data, i.e., selecting p-value thresholds larger than the baseline  $1e-8$  was always beneficial (median 12.8% increase in effect size), followed by lassosum (median 6.2% increase) and LDpred2 (median 4.1% increase). MegaPRS and PRScs benefited the least (median 1.2% and 0.2% increase, respectively). For SAS target data, the median benefits were smaller, except for PRScs (8) and overall less consistent (Figure 3C). Mean increases were larger for all methods except for PRScs in EUR target data, often dominated by few instances in which automatic tuning had comparatively low performance. The performance increases seen by CV-tuning were by and large significant when evaluated in EUR target data (Figure 3C, FWER  $\leq 0.05$ , two-sided z-test), except for PRScs which only saw an improvement for two phenotypes. Significant negative effects of CV-tuning for

EUR data were only observed for RA (PRSCs and MegaPRS, driven by FinnGen) and CKD (PRSCs, 12). For SAS target data, we observed fewer significant differences, and PRSCs was the only method for which we observed a significant reduction in effect size (BMI, 13). A more detailed description of these comparisons is provided in the Supplemental Results.

#### 455 Tuned PGS performance varies more between biobanks than between methods within biobanks

We estimated PGS effect size heterogeneity between biobanks and how it compares to the heterogeneity between methods within biobanks using 3-level meta-analytic random effects models in EUR target data (Methods, Figure 4). These nested models have two random effect  
460 parameters:  $\tau^2_{\text{biobank}}$  and  $\tau^2_{\text{method}}$ .  $\tau^2_{\text{biobank}}$  captures effect-size heterogeneity due to differences between biobanks, and  $\tau^2_{\text{method}}$  captures heterogeneity due to differences between methods within biobanks (we report their square roots  $\tau_{\text{biobank}}$  and  $\tau_{\text{method}}$ , as they are on the same scale as the PGS effect sizes). Additionally, we estimated  $I^2_{\text{biobank}}$  and  $I^2_{\text{method}}$ , which quantify the overall fraction of variance (between 0 and 1) in effect sizes attributable to biobank or choice  
465 of method within biobank, respectively.

We focused on scores selected via cross-validation in the UKBB-EUR sample (if available) and excluded scores from SBayesR that performed poorly in the UKBB-EUR 80% training data (RA, T1D, BMI, urate/gout). We did not consider the ensemble PGS or baseline method pT+clump, meaning that up to 6 scores were considered per trait. This setting was chosen to  
470 mimic the case in which multiple validated PGS from standard methods are available.

We found significant heterogeneity of PGS effect sizes in all 13 traits replicated in at least two biobanks (FWER $\leq$ 0.05, Cochran's  $Q$ -test, accounting for 13 tests). The target biobank had a larger influence on the PGS effect size than the choice of method within biobank across all traits (i.e.,  $\tau_{\text{method}} < \tau_{\text{biobank}}$ , Figure 4, Table 4). When adjusting for covariates, sex, age and  
475 genetic PCs 1-10 this effect was slightly reduced, but  $\tau_{\text{method}} < \tau_{\text{biobank}}$  remained true for the majority of traits (10 out of 13; with T1D, Stroke and T2D having  $\tau_{\text{method}} > \tau_{\text{biobank}}$ ) (Supplemental Results, Table S16). However, likelihood-based 95% confidence intervals for  $\tau_{\text{biobank}}$  were large and sometimes included the estimate for  $\tau_{\text{method}}$  (RA, stroke, T2D) and 0 (T1D, breast cancer). The variation in PGS effect sizes could to a large degree be explained  
480 by heterogeneity between biobanks (average  $I^2_{\text{biobank}} = 82.9\% \pm 14.3$  sd,  $N = 13$ ) and, to a lesser degree by heterogeneity between methods (average  $I^2_{\text{method}} = 11.97\% \pm 12.4$  sd,  $N=13$ ). Effect sizes for inflammatory bowel disease and RA varied most between biobanks ( $\tau_{\text{biobank}}$ ), also when adjusting for the average effect size ( $\tau_{\text{biobank}}/\mu_{\beta}$ ). Effect sizes for BMI varied the

least between biobanks, both absolutely ( $\tau_{\text{biobank}}$ ) and relative to the average effect size  
485 ( $\tau_{\text{biobank}}/\mu_{\beta}$ ). For binary traits, effect sizes for breast cancer varied the least across biobanks,  
both absolutely and relative to the average effect size.

Across traits,  $\tau_{\text{method}}$  was correlated with the average effect size (Pearson correlation 0.54,  $p =$   
0.0558,  $t$ -statistic = 2.1377, 11 degrees of freedom), especially when removing T1D and RA  
(Pearson correlation 0.85,  $p = 0.00094$ ,  $t$ -statistic = 4.83, 9 degrees of freedom), i.e., we  
490 found a linear relationship between the differences between methods and overall effect size,  
especially in the set of non-autoimmune traits.

$\tau_{\text{biobank}}$  was less correlated with the meta-analyzed average effect size (Pearson correlation  
0.367,  $p = 0.219$ ,  $t$ -statistic = 1.30, 11 degrees of freedom), i.e., large PGS effect sizes  
weren't necessarily associated with higher variability between biobanks.

495 For SAS ancestry target data, we did not find significant heterogeneity of PGS effect sizes in  
CKD, stroke, prostate cancer, or breast cancer ( $\text{FWER} \leq 0.05$ , Cochran's  $Q$ -test, accounting  
for 11 tests) and  $\tau_{\text{biobank}}$  could never reliably be estimated (14). Likelihood-based 95%  
confidence intervals for  $\tau_{\text{biobank}}$  included 0 for 9/11 replicated traits (all but T2D and eGFR).

### High variability between PGS methods for autoimmune diseases

500 Effect sizes were most variable between methods for autoimmune diseases T1D and RA  
( $\tau_{\text{method}}$ ) (15), even when accounting for the average effect size in those traits ( $\tau_{\text{method}}/\mu_{\beta}$ ), or  
relative to the total variation of effect sizes ( $I^2_{\text{method}}$ ). The scores for T1D and RA also had the  
largest fraction of PGS variance originating in the HLA region (mean  $0.7 \pm 0.12$  sd and  $0.54$   
 $\pm 0.21$  sd, respectively) (16). For T1D, the method with the largest effect size appeared to be  
505 biobank-specific, with FinnGen favoring PRSCs, while the UKBB and EBB had significantly  
larger effect sizes for LDpred2 and MegaPRS (forest plots for all 3-level meta-analytic  
models are available in the Supplemental Data). In contrast, effect sizes for BMI and HDL  
varied the least between methods.

Regarding SAS-matched target data, RA was only available in GNH with a limited number  
510 of cases (60) but displayed the highest heterogeneity of effect sizes due to method ( $\tau_{\text{method}} =$   
0.137, 95% CI: 0.046-0.379), consistent with the findings in EUR ancestry. T1D scores were  
not predictive in GNH (1), and not evaluated in the UKBB due to small sample size,  
therefore, we couldn't replicate the related findings from the EUR subset.

## Discussion

515 With this study, we have provided a comprehensive systematic PGS method comparison,  
with over one million individuals across multiple biobanks.



By publishing our workflow, we aim to increase access to PGS methods and facilitate future research. We believe that PGS method software could be greatly improved by support for standard formats (e.g., those maintained by the GWAS Catalog and PGS Catalog<sup>49</sup>) alongside software containerization (containers were supported in all the research environments that contributed to this study).

Our analysis was based on a previously published framework<sup>16</sup> which we automated, and expanded application and evaluation to multiple biobanks. Recent methods explicitly tailored for diverse target populations or source GWAS<sup>50-52</sup> were missing in this framework, and diverse ancestries were not well represented in our target data, which provides a limitation of this study. PGS tuning was performed in one biobank (UKBB-EUR) relying largely on author-provided LD reference panels. This approach more closely resembles real-world PGS application and allowed us to harness the full sample sizes in other biobanks/ancestries to maximize statistical power, and test transferability.

Importantly, we were unable to identify a single method that consistently outperformed all others (not counting the ensemble PGS), and the two highest performing methods (CV-tuned MegaPRS and LDpred2) were virtually tied. The best automatically tuned method was MegaPRS, albeit like other automatic methods it suffered sporadic cases of comparatively lower performance. Which method performs best may vary based on the specifics of the GWAS summary statistics, trait, and target sample. Given that the best methods performed so similarly, other modelling choices not investigated here (such as the set of included variants and their availability in the target sample) may well tip the balance in favor of one or the other when starting from the same GWAS summary statistics. Based on our results, we recommend tuning with cross-validation (with sufficiently large ranges of hyperparameters) instead of using methods' automatic settings, primarily to prevent cases of comparatively lower performance, rather than providing large improvements across traits. These findings are in line with previous comparisons showing moderate gains when tuning and evaluation are performed within biobanks<sup>16,18</sup>.

One reason for the lower performance of automatic tuning could be model misspecification, e.g., mismatched LD-references, or misreported fields in the input summary statistics. These inconsistencies may not be considered when tools are developed. The variable performance of LDpred2-auto stood out particularly against the high performance of CV-tuned PGS from same method. We note that LDpred2-auto has been updated at the time of writing including an optional new parameterization, which could affect its performance<sup>53</sup>. Limiting ourselves to the implementation of methods provided by GenoPred (which implements default method

parameters) meant that we did not evaluate CV-tuning for DBSLMM, which has since been recommended by the authors (with performance gains over the default automatic version of about 1.13%<sup>18</sup>).

555 These cases highlight a challenge faced by any method comparison: The frequent emergence of new tools, methods and related recommendations means that comparisons risk becoming outdated shortly after execution. Method evaluation across multiple biobanks can hardly match the pace of new developments. We therefore caution against using the results of this study to make definitive claims about relative method performance of actively developed methods. A more sustainable approach to method comparisons would be decentralized, with  
560 researchers individually submitting performance estimates for published scores (starting from the same summary statistics and variants) to a central repository and receiving credit by having such submissions be referenceable.

Using meta-analytic mixed models, we found that the performances of well-tuned PGS varied more between biobanks than within biobanks. This trend held true for most traits even when  
565 including covariates age, sex and genetic PCs 1-10. This likely reflects heterogeneity in phenotyping (e.g., disease diagnosis practices) rather than differences in population structure or genotyping. Effect sizes for BMI, which presumably is consistently measured, varied the least between biobanks, supporting this hypothesis. Yet, we cannot exclude a genetic contribution to the heterogeneity between biobanks, as PGS performance has been shown to  
570 vary with the distance to the GWAS sample even within genetically similar groups matched to reference populations<sup>54</sup>. The variability between biobanks for some traits implies that scores need to be re-evaluated when switching between different target data even when comparing ancestry-matched populations.

We note that the parameters by which we quantified variability are sensitive to which  
575 biobanks and PGS are included. The setting we chose mimics the case in which multiple UKBB-EUR-validated PGS are available. The variability between methods could increase if poorly performing (non-validated) scores are included in the analysis. On the other hand, the variability between biobanks could decrease if, e.g., phenotype definitions were further refined.

580 We found particularly large differences between methods for autoimmune diseases T1D and RA. This could be driven by the way methods handle the HLA-region, as well as genotyping differences in the target biobanks. Our analyses highlight modelling of the HLA-region as an area in which methods could potentially be improved.

One of the most useful insights from this study is that ensemble PGS tuned in the UKBB-  
585 EUR sample provided consistently strong performance, albeit at the cost of higher  
computational demand during training. This shows that benefits seen within a target  
sample<sup>16,18</sup> can be transferred to other samples without re-tuning ensemble weights. We see  
this method as complimentary to cross-trait prediction strategies (MultiPGS)<sup>55-57</sup> that use PGS  
constructed from multiple sets of GWAS summary statistics (from different traits).  
590 Considering the small differences in performance we observe for well-tuned scores from  
single methods, we see ensemble PGS and MultiPGS as promising avenues to further  
improve PGS performances beyond what is currently possible with single methods. Future  
research needs to assess how well EUR-trained ensemble PGS transfer to other genetic  
ancestries. It is possible that training needs to be performed in a population similar to the  
595 target population to ensure optimal performance and avoid exacerbating already existing  
issues with current PGS<sup>7</sup>.

In Summary, while no single method outperformed all others, method ensembles provided  
consistently strong performance (with few exceptions). PGS effect size heterogeneity  
between biobanks was larger than between methods within biobanks, likely pointing to  
600 challenges with phenotyping. Large heterogeneity between methods was observed for  
autoimmune diseases, indicating that special care should be taken for PGS which rely heavily  
on the HLA region. Our open-source workflow, analyses framework and online results  
provide a rich ground for future method benchmarking and development.

### Data and code availability

605 The prspipe workflow used to generate polygenic score weights, perform  
polygenic scoring and ancestry matching is available on GitHub  
(<https://github.com/intervene-EU-H2020/prspipe>).

Non-sensitive experimental data exported from the biobanks are permissively licensed and  
deposited in an open data repository (<https://zenodo.org/doi/10.5281/zenodo.10012995>).

610 Processed summary statistics are permissively licensed and hosted on GitHub and accessible  
through in an R data package (<https://github.com/intervene-EU-H2020/pgsCompaR>). A  
website containing an interactive results browser is permissively licensed and available on  
GitHub (<https://github.com/intervene-EU-H2020/pgs-method-compare>), hosted at  
<https://methodscomparison.intervenegeneticscores.org/>.

615 15

## Consortia

Current Genes & Health Research Team (in alphabetical order by surname): Shaheen Akhtar,  
620 Mohammad Anwar, Omar Asgar, Samina Ashraf, Saeed Bidi, Gerome Breen, James Broster,  
Raymond Chung, David Collier, Charles J Curtis, Shabana Chaudhary, Grainne Colligan,  
Panos Deloukas, Ceri Durham, Faiza Durrani, Fabiola Eto, Sarah Finer, Joseph Gafton, Ana  
Angel, Chris Griffiths, Joanne Harvey, Teng Heng, Sam Hodgson, Qin Qin Huang, Matt  
Hurles, Karen A Hunt, Shapna Hussain, Kamrul Islam, Vivek Iyer, Benjamin M Jacobs,  
625 Georgios Kalantzis, Ahsan Khan, Claudia Langenberg, Cath Lavery, Sang Hyuck Lee, Daniel  
MacArthur, Sidra Malik, Daniel Malawsky, Hilary Martin, Dan Mason, Rohini Mathur,  
Mohammed Bodrul Mazid, John McDermott, Caroline Morton, Bill Newman, Elizabeth  
Owor, Asma Qureshi, Shwetha Ramachandrappa, Mehru Raza, Jessry Russell, Nishat Safa,  
Miriam Samuel, Moneeza Siddiqui, Michael Simpson, John Solly, Marie Spreckley. Daniel  
630 Stow, Michael Taylor, Richard C Trembath, Karen Tricker, David A van Heel, Klaudia  
Walter, Caroline Winckley, Suzanne Wood, John Wright, Ishevanhu Zengeya, Julia Zöllner.  
The current members of HUNT All-In Research Team (in alphabetical order by surname):  
Bjørn Olav Åsvold, Ben Brumpton, Maiken Elvestad Gabrielsen, Kristian Hveem, Ida  
Surakka, Laurent Thomas, Wei Zhou

635 Estonian Biobank research team members are Andres Metspalu, Lili Milani, Tõnu Esko,  
Reedik Mägi, Mari Nelis and Georgi Hudjashov

The current members of FinnGen can be found in the attached table "FinnGen-  
banner\_Authors Jan2024.xlsx"

## Acknowledgements

640 See Supplements

## Author Disclosures and Declaration of Interests

M.I. is a trustee of the Public Health Genomics (PHG) Foundation, a member of the  
Scientific Advisory Board of Open Targets, and has a research collaboration with  
AstraZeneca PLC which is unrelated to this study. M.I. is supported by core funding from the  
645 British Heart Foundation (RG/18/13/33946) and NIHR Cambridge Biomedical Research  
Centre (BRC-1215-20014; NIHR203312). The views expressed are those of the authors and  
not necessarily those of the NIHR or the Department of Health and Social Care.

K.L. has participated as an analyst in a collaboration research project at the Institute of  
Genomics, University of Tartu, which was funded by Geneto OÜ.

650 O.P. provides consultancy services for UCB pharma company.

#### Author Contributions

C.L. and A.G. conceptualized the study. R. Monti, S.W. and O.P. wrote the prspipe workflow. R. Monti, L.E., G.H., K.L., S.K. and B. Wolford performed analyses in biobanks. R. Monti and L.E. performed statistical meta-analyses. B. Wingfield implemented the  
655 companion website. R. Monti wrote the manuscript with assistance from all co-authors. L.E. lead revisions. L.E., R. Monti, G.H., K.L., S.K., and B. Wolford performed revisions. All authors contributed to regular discussions and provided critical feedback regarding the study design/results and contributed to review/editing of the manuscript.

#### Web Resources

660 Polygenic score weights for scores that were at least nominally significantly associated with the phenotype ( $p < 0.05$ ) in all EUR target data samples are made publicly available through the GWAS catalog (<https://www.ebi.ac.uk/gwas/>) with publication ID PGP000517. All evaluated scores except the one produced by LDpred2-auto for RA met this threshold. A list of PGS catalog score IDs is provided in 15.

665

study	GWAS trait	$N_{cas}$	$N_{con}$	$N_{variants}$	target traits
<b>GCST005838</b> <sup>66</sup>	Stroke	67,162	454,450	1,121,867	Stroke
<b>GCST90012877</b> <sup>67</sup>	AD or family history of AD	53,042	355,900	1,136,233	AD
<b>GCST90013534</b> <sup>68</sup>	RA	22,628	288,664	778,275	RA
<b>GCST004773</b> <sup>29</sup>	T2D	26,676	132,532	1,071,786	T2D
<b>GCST004988</b> <sup>69</sup>	Breast cancer	76,192	63,082	1,137,481	Breast cancer
<b>GCST006085</b> <sup>70</sup>	Prostate cancer	79,148	61,106	1,139,693	Prostate cancer
<b>GCST90013445</b> <sup>26</sup>	T1D	22,153	37,374	63,204	T1D
<b>GCST004131</b> <sup>71</sup>	IBD	25,042	34,915	1,103,333	IBD
<b>GCST008059</b> <sup>72</sup>	eGFR	567,460	-	1,141,659	CKD, eGFR
<b>GCST90018959</b> <sup>73</sup>	Height	525,444	-	1,119,889	Height
<b>GCST008972</b> <sup>27</sup>	Urate levels	457,690	-	1,005,478	Gout, Urate
<b>GCST002783</b> <sup>74</sup>	BMI	236,781	-	1,039,042	BMI
<b>GCST007140</b> <sup>75</sup>	HDL	94,674	-	1,138,452	HDL
<b>GCST007954</b> <sup>28</sup>	HbA1c	88,355	-	1,009,664	HbA1c

670 **Table 1: GWAS summary statistics used to derive PGS weights**

Entries are ordered by the total sample size and type of trait (binary, continuous). From left to right: GWAS catalog study identifiers (study), the respective reported GWAS traits, number of cases ( $N_{cas}$ ) and controls ( $N_{con}$ ), the number of variants after intersection with HapMap3-1KG and quality control ( $N_{variants}$ ), and the evaluated target traits. Scores constructed from  
675 urate and eGFR summary statistics were also evaluated for gout and CKD, respectively. The

GWAS for T1D considered only a small panel of variants of which 84% remained after intersection and QC.

	EUR total	SAS total	EUR EBB	EUR FinnGen	EUR HUNT	EUR UKBB (test)	EUR UKBB (train)	SAS GNH	SAS UKBB
AD	15,940	-	555	13,823	1,562	-	-	-	-
RA	13,060	60	2,384	9,332	1,139	205	820	60	-
Breast cancer	23,610	393	2,685	16,076	1,729	3,120	12,483	197	196
CKD	19,714	1,609	4,224	9,314	2,802	3,374	13,496	1,131	478
Gout	22,399	488	10,646	8,759	1,318	1,676	6,704	282	206
IBD	13,016	634	2,097	7,815	1,769	1,335	5,340	466	168
Prostate cancer	20,492	205	2,227	13,606	2,242	2,417	9,671	95	110
Stroke	37,920	635	4,515	26,166	5,204	2,035	8,142	424	211
T1D	5,384	443	501	4,286	396	201	804	443	-
T2D	81,487	8,696	12,344	59,345	3,861	5,937	23,748	6,630	2,066
BMI	346,290	42,243	189,651	-	66,663	89,976	359,913	33,146	9,097
HDL	139,248	37,693	10,642	-	49,824	78,782	315,135	29,628	8,065
HbA1c	120,242	21,696	-	-	34,192	86,050	344,209	12,948	8,748
Height	524,056	43,197	190,013	267,343	66,700	-	-	34,089	9,108
Urate	85,973	13,572	-	-	-	85,973	343,904	4,730	8,842
eGFR	152,793	38,916	-	-	66,759	86,034	344,140	3,061	8,855

680 **Table 2: Target sample sizes across traits.**

For each trait and replicated ancestry group (EUR, SAS) the number of cases (binary disease traits) or sample size are shown, either combined (“total”, excluding UKBB training data) or separated by biobank. For the UKBB-EUR, data were split into train (80%, used to tune hyperparameters and ensemble PGS) and test sets (20%, used for evaluation and meta-  
685 analyses). UKBB EUR data were excluded for Alzheimer’s disease and height due to sample



overlap and could therefore not be used for tuning (leaving 14 traits for a full evaluation).  
Dashes “-” indicate the phenotype was unavailable.

method	tuning type	trait	N (EUR)	N (SAS)	median (EUR)	median (SAS)	mean (EUR)	sd (EUR)	mean (SAS)	sd (SAS)
ldpred2	CV	binary	9	8	0.965	0.963	0.943	0.045	0.972	0.127
megaprs	CV	binary	9	8	0.957	0.934	0.947	0.041	0.958	0.124
lassosum	CV	binary	9	8	0.921	0.914	0.913	0.061	0.920	0.111
prscs	CV	binary	9	8	0.903	0.900	0.896	0.100	0.909	0.259
pt.clump	CV	binary	9	8	0.735	0.734	0.721	0.077	0.748	0.186
megaprs	auto	binary	9	8	0.948	0.950	0.933	0.050	0.964	0.104
ldpred2	auto	binary	9	8	0.927	0.955	0.838	0.265	0.904	0.314
prscs	auto	binary	9	8	0.925	0.876	0.915	0.052	0.877	0.264
sbayesr	auto	binary	9	8	0.907	0.895	0.873	0.083	0.841	0.181
dbslmm	auto	binary	9	8	0.904	0.865	0.890	0.092	0.815	0.199
lassosum	auto	binary	9	8	0.891	0.870	0.861	0.103	0.753	0.262
pt.clump	auto	binary	9	8	0.629	0.627	0.607	0.098	0.527	0.302
ldpred2	CV	continuous	5	5	0.950	0.936	0.948	0.016	0.925	0.049
megaprs	CV	continuous	5	5	0.927	0.931	0.940	0.024	0.946	0.034
prscs	CV	continuous	5	5	0.923	0.926	0.909	0.031	0.875	0.090
lassosum	CV	continuous	5	5	0.906	0.920	0.914	0.026	0.916	0.021
pt.clump	CV	continuous	5	5	0.735	0.729	0.743	0.037	0.718	0.048
prscs	auto	continuous	5	5	0.923	0.928	0.907	0.028	0.903	0.076
dbslmm	auto	continuous	5	5	0.901	0.904	0.891	0.039	0.897	0.066
sbayesr	auto	continuous	5	5	0.887	0.862	0.868	0.075	0.806	0.184
megaprs	auto	continuous	5	5	0.883	0.907	0.885	0.067	0.922	0.055
lassosum	auto	continuous	5	5	0.873	0.890	0.877	0.021	0.901	0.054
ldpred2	auto	continuous	5	5	0.851	0.778	0.823	0.121	0.781	0.114
pt.clump	auto	continuous	5	5	0.643	0.614	0.606	0.088	0.635	0.112

690 **Table 3: PGS meta-analyzed  $\beta$  coefficients relative to the ensemble PGS ( $\beta_{s^*}/\beta_{\text{EnsPGS}^*}$ )**

For the 14 traits for which we tuned hyperparameters with CV, relative PGS effect sizes relative the ensemble PGS are shown ( $\beta_{s^*}/\beta_{\text{EnsPGS}^*}$ ) stratified by PGS method, tuning type

(CV/auto), ancestry (EUR, SAS) and type of trait (binary/continuous). The number of traits (N), medians, means and standard deviations (sd) are shown. Methods are ordered by the median EUR relative effect size within traits and tuning types.

trait	$N_{\text{biobank}}$	$N_{\text{method}}$	$\mu_{\beta} \pm \text{sd}$	$\tau_{\text{biobank}}$ (95% CI)	$\tau_{\text{method}}$ (95% CI)	$I^2_{\text{biobank}}$ (%)	$I^2_{\text{method}}$ (%)
T1D	3	5*	0.815 $\pm 0.05$	0.072 (0-0.383)	0.069(0.046- 0.112)	48.8	45
Prostate cancer	4	6	0.664 $\pm 0.029$	0.054 (0.024- 0.167)	0.02(0.015- 0.029)	82.2	11
Breast cancer	4	6	0.537 $\pm 0.017$	0.029 (0-0.1)	0.014(0.012- 0.02)	67.1	16.5
Gout	4	5*	0.522 $\pm 0.05$	0.098 (0.049- 0.294)	0.018(0.014- 0.028)	94.7	3.3
IBD	4	6	0.513 $\pm 0.089$	0.177 (0.092- 0.525)	0.019(0.015- 0.028)	97.8	1.1
RA	4	5*	0.458 $\pm 0.087$	0.165 (0.067- 0.522)	0.095(0.068- 0.144)	74.4	24.7
T2D	4	6	0.428 $\pm 0.017$	0.031 (0.013- 0.099)	0.015(0.012- 0.02)	77.8	17.2
CKD	4	6	0.213 $\pm 0.031$	0.059 (0.028- 0.18)	0.006(0.006- 0.01)	93.4	0.9
Stroke	4	6	0.133 $\pm 0.016$	0.029 (0.01- 0.099)	0.01(0.01-0.016)	78.5	10.4
HDL	3	6	0.303 $\pm 0.02$	0.035 (0.016- 0.148)	0.005(0.005- 0.009)	96.2	2.3
BMI	3	5*	0.282 $\pm 0.006$	0.01 (0.01-0.046)	0.004(0.004- 0.004)	80.2	13.8
eGFR	2	6	0.267 $\pm 0.046$	0.065 (0.025- 0.733)	0.009(0.009- 0.015)	97.9	1.8
HbA1c	2	6	0.172 $\pm 0.013$	0.018 (0.011- 0.211)	0.005(0.005- 0.009)	88.3	7.5

**Table 4: 3-level meta-analytical random effects model results (EUR)**

700 Table corresponding to Figure 4. From left to right, the target trait, the number of biobanks  
with the trait ( $N_{\text{biobank}}$ ), the number of methods/scores considered ( $N_{\text{method}}$ ), the meta-analyzed  
average PGS effect size across methods/scores and biobanks ( $\mu_{\beta}$ ) with standard deviation  
(sd), the standard deviation of the random intercepts specific to biobanks ( $\tau_{\text{biobank}}$ ) including  
95% likelihood-based confidence intervals (95% CI), the standard deviation of the random  
705 intercepts specific to methods within biobanks ( $\tau_{\text{method}}$ ) including 95% CI, the fraction of total  
effect size variance due to heterogeneity between biobanks ( $I^2_{\text{biobank}}$ ) in %, and the fraction of  
the total effect size variance due to heterogeneity between methods ( $I^2_{\text{method}}$ ) in %. Endpoints  
are ordered by type (binary/continuous) and  $\mu_{\beta}$ . SBayesR was excluded for RA, T1D, gout,

and BMI (\*). Full results for EUR and SAS are given in 9-12

710

**Figure 1: prspipe workflow and application.**

Prspipe is a snakemake workflow that automates within-biobank method comparisons introduced by Pain et al.<sup>16</sup> The public stage uses only public data (e.g., summary statistics, 715 ancestry reference, PGS software) to derive PGS weights using seven methods from GWAS summary statistics. The private stage requires access to target genotype and phenotype data and includes data harmonization, polygenic scoring, and PGS tuning using cross-validation (CV). We used prspipe to generate PGS weights and tune hyperparameters in the UKBB EUR data (full run). PGS weights were shared with other biobanks for evaluation/replication 720 (skipping the public stage). Other biobanks were not used for hyperparameter tuning. Downstream analyses were conducted to determine PGS performance using a meta-analytical framework (not part of the workflow), and results were published as an online resource at <https://methodscomparison.intervenegeneticscores.org>.

**Figure 2: Meta-analysis workflow for methods comparison, example: type 2 diabetes.**

a) PGS effect sizes  $\beta_{s,b}$  (i.e., the change in log odds-ratio per PGS standard deviation measured for scores “s” across biobanks “b”, see Methods) with 95% confidence intervals for all PGS methods (x-axis) stratified by biobank, replicated ancestries (EUR, SAS) and tuning types (auto, CV) serve as the inputs for the meta-analysis (shown for example trait type 2 diabetes). We evaluated scores for seven methods shown on the right, as well as the UKBB-EUR-tuned ensemble PGS (EnsPGS). The largest effect size for each ancestry and biobank is marked with a triangle (given by the ensemble in all cases).  $\beta_{s,b}$  for all target data and traits are displayed in 1 and browsable online. b) PGS effect-sizes are meta-analyzed within ancestries across biobanks (yielding a single  $\beta_{s^*}$  for each score “s”). Effect-size differences relative to the largest meta-analyzed effect size ( $\beta_{top^*}$ , given by the ensemble) and 95% confidence intervals are shown. All pairwise differences are available in 5-6, and browsable online. c) Meta-analyzed effect sizes  $\beta_{s^*}$  are compared, and significance testing is performed. Heatmaps show both the effect-size relative to the largest ( $\beta_{s^*} / \beta_{top^*}$ , left) as well as corresponding two-sided z-test significance levels at which  $H_0: \beta_{s^*} - \beta_{top^*} = 0$  can be rejected (right). Significant differences at a FWER  $\leq 0.05$  are marked with an asterisk (\*), accounting for all 351 tests performed across traits and ancestries. The score against which comparisons are performed with effect size  $\beta_{top^*}$  is marked with a “1” and black border. Arrows indicate two example comparisons: against PRSs-CV (significant difference in SAS and EUR) and LDpred2-auto (significant only in EUR). Data for all PGS and traits are provided in 6

750 **Figure 3: Relative meta-analyzed PGS effect sizes across 14 traits**

a) For the 14 traits for which we tuned hyperparameters using CV (x-axis), we show heatmaps of meta-analyzed  $\beta$ -coefficients relative to the highest within traits ( $\beta_s^*/\beta_{top}^*$ , left) as well as significance levels for the two-sided z-test ( $H_0: \beta_s^* - \beta_{top}^* = 0$ , right) stratified by ancestry (EUR, SAS). The top score with the largest effect size for each trait ( $\beta_{top}^*$ ) is marked with a “1” and black box. Differences significant at FWER  $\leq 0.05$  are marked with asterisks (\*), accounting for all 351 pairwise tests performed across traits, replicated ancestries and tuning types (auto, CV) (full data for all traits and scores are displayed in 6). b) Bar plot counting PGS ranks across traits (1 is the highest), stratified by ancestry, method, and tuning-type (auto, CV). Methods are ordered by the average rank-sum across tuning types (highest to lowest). c) For each method (y-axis), dot plots showing the relative meta-analyzed effect-size of the score derived using methods’ automatic settings against the CV-tuned scores ( $\beta_{auto}^*/\beta_{CV}^*$ ). Colors denote sign and significance of the two-sided z-test ( $H_0: \beta_{CV}^* - \beta_{auto}^* = 0$ ) at FWER  $\leq 0.05$  after accounting for 114 tests across traits and ancestries (12-13). Methods are ordered by the median difference.

765



#### Figure 4: 3-level meta-analysis of PGS effect sizes in EUR target data

For all 13 traits replicated in at least 2 biobanks in EUR ancestry target data and CV-tuned in UKBB, from left to right: 1) PGS effect sizes ( $\beta$ -coefficients,  $\beta_{m,b}$ ) with 95% confidence intervals for three example traits within biobanks (T1D: high variability between methods, IBD: high variability between biobanks, T2D: intermediate to low variability between methods and biobanks), 2) the meta-analyzed average effect sizes across biobanks and methods ( $\mu_\beta$ ) with bars denoting the square roots of the variance components ( $\tau$ ), i.e., the standard deviations of the random intercepts for biobanks or methods, 3)  $\tau$ -values with likelihood-based 95% confidence intervals and 4)  $I^2$  estimates, i.e., the fraction of variance of effect sizes explained by heterogeneity between biobanks or methods within biobanks.  $\tau$  and  $I^2$  are colored according to the levels of the meta-analytic 3-level random effects model (Methods).

1. Lee, A., Mavaddat, N., Wilcox, A.N., Cunningham, A.P., Carver, T., Hartley, S., Babb de Villiers, C., Izquierdo, A., Simard, J., Schmidt, M.K., et al. (2019). BOADICEA: a comprehensive breast cancer risk prediction model incorporating genetic and nongenetic risk factors. *Genetics in Medicine* *21*, 1708–1718. 10.1038/s41436-018-0406-9.
- 785 2. Weale, M.E., Riveros-Mckay, F., Selzam, S., Seth, P., Moore, R., Tarran, W.A., Gradovich, E., Giner-Delgado, C., Palmer, D., Wells, D., et al. (2021). Validation of an Integrated Risk Tool, Including Polygenic Risk Score, for Atherosclerotic Cardiovascular Disease in Multiple Ethnicities and Ancestries. *The American Journal of Cardiology* *148*, 157–164. 10.1016/j.amjcard.2021.02.032.
- 790 3. Mars, N., Koskela, J.T., Ripatti, P., Kiiskinen, T.T.J., Havulinna, A.S., Lindbohm, J.V., Ahola-Olli, A., Kurki, M., Karjalainen, J., Palta, P., et al. (2020). Polygenic and clinical risk scores and their impact on age at onset and prediction of cardiometabolic diseases and common cancers. *Nat Med* *26*, 549–557. 10.1038/s41591-020-0800-0.
- 795 4. Mars, N., Lindbohm, J.V., Parolo, P. della B., Widén, E., Kaprio, J., Palotie, A., and Ripatti, S. (2022). Systematic comparison of family history and polygenic risk across 24 common diseases. *The American Journal of Human Genetics* *109*, 2152–2162. 10.1016/j.ajhg.2022.10.009.
- 800 5. Khera, A.V., Chaffin, M., Aragam, K.G., Haas, M.E., Roselli, C., Choi, S.H., Natarajan, P., Lander, E.S., Lubitz, S.A., Ellinor, P.T., et al. (2018). Genome-wide polygenic scores for common diseases identify individuals with risk equivalent to monogenic mutations. *Nature genetics* *50*, 1219–1224.
6. Lewis, C.M., and Vassos, E. (2020). Polygenic risk scores: From research tools to clinical instruments. *Genome Medicine* *12*, 1–11. 10.1186/s13073-020-00742-5.
- 805 7. Adeyemo, A., Balaconis, M.K., Darnes, D.R., Fatumo, S., Granados Moreno, P., Hodonsky, C.J., Inouye, M., Kanai, M., Kato, K., Knoppers, B.M., et al. (2021). Responsible use of polygenic risk scores in the clinic: potential benefits, risks and gaps. *Nat Med* *27*, 1876–1884. 10.1038/s41591-021-01549-6.
- 810 8. Marston, N.A., Kamanu, F.K., Nordio, F., Gurmu, Y., Roselli, C., Sever, P.S., Pedersen, T.R., Keech, A.C., Wang, H., Lira Pineda, A., et al. (2020). Predicting Benefit From Evolocumab Therapy in Patients With Atherosclerotic Disease Using a Genetic Risk Score. *Circulation* *141*, 616–623. 10.1161/CIRCULATIONAHA.119.043805.
- 815 9. Damask, A., Steg, P.G., Schwartz, G.G., Szarek, M., Hagström, E., Badimon, L., Chapman, M.J., Boileau, C., Tsimikas, S., Ginsberg, H.N., et al. (2020). Patients With High Genome-Wide Polygenic Risk Scores for Coronary Artery Disease May Receive Greater Clinical Benefit From Alirocumab Treatment in the ODYSSEY OUTCOMES Trial. *Circulation* *141*, 624–636. 10.1161/CIRCULATIONAHA.119.044434.
10. Mak, T.S.H., Porsch, R.M., Choi, S.W., Zhou, X., and Sham, P.C. (2017). Polygenic scores via penalized regression on summary statistics. *Genetic Epidemiology* *41*, 469–480. 10.1002/gepi.22050.

- 820 11. Ge, T., Chen, C.-Y., Ni, Y., Feng, Y.-C.A., and Smoller, J.W. (2019). Polygenic prediction via Bayesian regression and continuous shrinkage priors. *Nat Commun* *10*, 1776. 10.1038/s41467-019-09718-5.
12. Lloyd-Jones, L.R., Zeng, J., Sidorenko, J., Yengo, L., Moser, G., Kemper, K.E., Wang, H., Zheng, Z., Magi, R., Esko, T., et al. (2019). Improved polygenic prediction by  
825 Bayesian multiple regression on summary statistics. *Nat Commun* *10*, 5086. 10.1038/s41467-019-12653-0.
13. Privé, F., Arbel, J., and Vilhjálmsson, B.J. (2020). LDpred2: Better, faster, stronger. *Bioinformatics* *36*, 5424–5431. 10.1093/bioinformatics/btaa1029.
14. Yang, S., and Zhou, X. (2020). Accurate and Scalable Construction of Polygenic Scores  
830 in Large Biobank Data Sets. *The American Journal of Human Genetics* *106*, 679–693. 10.1016/j.ajhg.2020.03.013.
15. Zhang, Q., Privé, F., Vilhjálmsson, B., and Speed, D. (2021). Improved genetic prediction of complex traits from individual-level data or summary statistics. *Nature Communications* *12*, 1–9. 10.1038/s41467-021-24485-y.
- 835 16. Pain, O., Glanville, K.P., Hagenaars, S.P., Selzam, S., Furtjes, A.E., Gaspar, H.A., Coleman, J.R.I., Rinfeld, K., Breen, G., Plomin, R., et al. (2021). Evaluation of polygenic prediction methodology within a reference-standardized framework. *PLoS Genetics* *17*, 1–22. 10.1371/journal.pgen.1009021.
17. Ni, G., Zeng, J., Revez, J.A., Wang, Y., Zheng, Z., Ge, T., Restuadi, R., Kiewa, J.,  
840 Nyholt, D.R., Coleman, J.R.I., et al. (2021). A Comparison of Ten Polygenic Score Methods for Psychiatric Disorders Applied Across Multiple Cohorts. *Biological Psychiatry* *90*, 611–620. 10.1016/j.biopsych.2021.04.018.
18. Yang, S., and Zhou, X. (2022). PGS-server: accuracy, robustness and transferability of polygenic score methods for biobank scale studies. *Briefings in Bioinformatics* *23*,  
845 bbac039. 10.1093/bib/bbac039.
19. Jermy, B., Läll, K., Wolford, B., Wang, Y., Zguro, K., Cheng, Y., Kanai, M., Kanoni, S., Yang, Z., Hartonen, T., et al. (2023). A unified framework for estimating country-specific cumulative incidence for 18 diseases stratified by polygenic risk. Preprint at medRxiv, 10.1101/2023.06.12.23291186 10.1101/2023.06.12.23291186.
- 850 20. Köster, J., Mölder, F., Jablonski, K.P., Letcher, B., Hall, M.B., Tomkins-Tinch, C.H., Sochat, V., Forster, J., Lee, S., Twardziok, S.O., et al. (2021). Sustainable data analysis with Snakemake. *F1000Research* *10*. 10.12688/f1000research.29032.2.
21. Kurki, M.I., Karjalainen, J., Palta, P., Sipilä, T.P., Kristiansson, K., Donner, K.M., Reeve, M.P., Laivuori, H., Aavikko, M., Kaunisto, M.A., et al. (2023). FinnGen provides genetic  
855 insights from a well-phenotyped isolated population. *Nature* *613*, 508–518. 10.1038/s41586-022-05473-8.
22. Leitsalu, L., Haller, T., Esko, T., Tammesoo, M.-L., Alavere, H., Snieder, H., Perola, M., Ng, P.C., Mägi, R., Milani, L., et al. (2015). Cohort Profile: Estonian Biobank of the Estonian Genome Center, University of Tartu. *Int J Epidemiol* *44*, 1137–1147.  
860 10.1093/ije/dyt268.

23. Åsvold, B.O., Langhammer, A., Rehn, T.A., Kjelvik, G., Grøntvedt, T.V., Sørgerd, E.P., Fenstad, J.S., Heggland, J., Holmen, O., Stuijbergen, M.C., et al. (2023). Cohort Profile Update: The HUNT Study, Norway. *International Journal of Epidemiology* 52, e80–e91. 10.1093/ije/dyac095.
- 865 24. Finer, S., Martin, H.C., Khan, A., Hunt, K.A., MacLaughlin, B., Ahmed, Z., Ashcroft, R., Durham, C., MacArthur, D.G., McCarthy, M.I., et al. (2020). Cohort Profile: East London Genes & Health (ELGH), a community-based population genomics and health study in British Bangladeshi and British Pakistani people. *International Journal of Epidemiology* 49, 20–21i. 10.1093/ije/dyz174.
- 870 25. Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L.T., Sharp, K., Motyer, A., Vukcevic, D., Delaneau, O., O’Connell, J., et al. (2018). The UK Biobank resource with deep phenotyping and genomic data. *Nature* 562, 203–209. 10.1038/s41586-018-0579-z.
- 875 26. Robertson, C.C., Inshaw, J.R.J., Onengut-Gumuscu, S., Chen, W.-M., Santa Cruz, D.F., Yang, H., Cutler, A.J., Crouch, D.J.M., Farber, E., Bridges, S.L., et al. (2021). Fine-mapping, trans-ancestral and genomic analyses identify causal variants, cells, genes and drug targets for type 1 diabetes. *Nat Genet* 53, 962–971. 10.1038/s41588-021-00880-5.
- 880 27. Tin, A., Marten, J., Halperin Kuhns, V.L., Li, Y., Wuttke, M., Kirsten, H., Sieber, K.B., Qiu, C., Gorski, M., Yu, Z., et al. (2019). Target genes, variants, tissues and transcriptional pathways influencing human serum urate levels. *Nat Genet* 51, 1459–1474. 10.1038/s41588-019-0504-x.
- 885 28. Wheeler, E., Leong, A., Liu, C.-T., Hivert, M.-F., Strawbridge, R.J., Podmore, C., Li, M., Yao, J., Sim, X., Hong, J., et al. (2017). Impact of common genetic determinants of Hemoglobin A1c on type 2 diabetes risk and diagnosis in ancestrally diverse populations: A transethnic genome-wide meta-analysis. *PLOS Medicine* 14, e1002383. 10.1371/journal.pmed.1002383.
29. Scott, R.A., Scott, L.J., Mägi, R., Marullo, L., Gaulton, K.J., Kaakinen, M., Pervjakova, N., Pers, T.H., Johnson, A.D., Eicher, J.D., et al. (2017). An Expanded Genome-Wide Association Study of Type 2 Diabetes in Europeans. *Diabetes* 66, 2888–2902. 10.2337/db16-1253.
- 890 30. Skene, N.G., Bryois, J., Bakken, T.E., Breen, G., Crowley, J.J., Gaspar, H.A., Giusti-Rodriguez, P., Hodge, R.D., Miller, J.A., Muñoz-Manchado, A.B., et al. (2018). Genetic identification of brain cell types underlying schizophrenia. *Nat Genet* 50, 825–833. 10.1038/s41588-018-0129-5.
- 895 31. Altshuler, D.M., Gibbs, R.A., Peltonen, L., Schaffner, S.F., Yu, F., Dermitzakis, E., Bonnen, P.E., De Bakker, P.I.W., Deloukas, P., Gabriel, S.B., et al. (2010). Integrating common and rare genetic variation in diverse human populations. *Nature* 467, 52–58. 10.1038/nature09298.
- 900 32. Gentleman, R.C., Carey, V.J., Bates, D.M., Bolstad, B., Dettling, M., Dudoit, S., Ellis, B., Gautier, L., Ge, Y., Gentry, J., et al. (2004). Bioconductor: open software development for computational biology and bioinformatics. *Genome Biology* 5, R80. 10.1186/gb-2004-5-10-r80.

- 905 33. Hinrichs, A.S., Karolchik, D., Baertsch, R., Barber, G.P., Bejerano, G., Clawson, H., Diekhans, M., Furey, T.S., Harte, R.A., Hsu, F., et al. (2006). The UCSC Genome Browser Database: update 2006. *Nucleic Acids Research* 34, D590–D598. 10.1093/nar/gkj144.
34. World Health Organization (2004). ICD-10: international statistical classification of diseases and related health problems: tenth revision (World Health Organization).
- 910 35. Levey, A.S., Coresh, J., Greene, T., Marsh, J., Stevens, L.A., Kusek, J.W., and Van Lente, F. (2007). Expressing the modification of diet in renal disease study equation for estimating glomerular filtration rate with standardized serum creatinine values. *Clinical Chemistry* 53, 766–772. 10.1373/clinchem.2006.077180.
36. Auton, A., Abecasis, G.R., Altshuler, D.M., Durbin, R.M., Bentley, D.R., Chakravarti, A., Clark, A.G., Donnelly, P., Eichler, E.E., Flicek, P., et al. (2015). A global reference for human genetic variation. *Nature* 526, 68–74. 10.1038/nature15393.
- 915 37. Price, A.L., Weale, M.E., Patterson, N., Myers, S.R., Need, A.C., Shianna, K.V., Ge, D., Rotter, J.I., Torres, E., Taylor, K.D., et al. (2008). Long-Range LD Can Confound Genome Scans in Admixed Populations. *The American Journal of Human Genetics* 83, 132–135. 10.1016/j.ajhg.2008.06.005.
- 920 38. Charrad, M., Ghazzali, N., Boiteau, V., and Niknafs, A. (2014). NbClust: An R Package for Determining the Relevant Number of Clusters in a Data Set. *Journal of Statistical Software* 61, 1–36. 10.18637/jss.v061.i06.
39. Kuhn, M. (2008). Building Predictive Models in R Using the caret Package. *Journal of Statistical Software* 28, 1–26. 10.18637/jss.v028.i05.
- 925 40. Friedman, J.H., Hastie, T., and Tibshirani, R. (2010). Regularization Paths for Generalized Linear Models via Coordinate Descent. *Journal of Statistical Software* 33, 1–22. 10.18637/jss.v033.i01.
41. Lee, S.H., Goddard, M.E., Wray, N.R., and Visscher, P.M. (2012). A better coefficient of determination for genetic profile analysis. *Genetic epidemiology* 36, 214–224. 10.1002/gepi.21614.
- 930 42. DeLong, E.R., DeLong, D.M., and Clarke-Pearson, D.L. (1988). Comparing the Areas under Two or More Correlated Receiver Operating Characteristic Curves: A Nonparametric Approach. *Biometrics* 44, 837–845. 10.2307/2531595.
- 935 43. Wang, Y., Namba, S., Lopera, E., Kerminen, S., Tsuo, K., Läll, K., Kanai, M., Zhou, W., Wu, K.-H., Favé, M.-J., et al. (2023). Global Biobank analyses provide lessons for developing polygenic risk scores across diverse cohorts. *Cell Genomics* 3, 100241. 10.1016/j.xgen.2022.100241.
44. Ebert, M.H., Pim Cuijpers, Toshi Furukawa, David (2021). *Doing Meta-Analysis with R: A Hands-On Guide* (Chapman and Hall/CRC) 10.1201/9781003107347.
- 940 45. Cheung, M.W.-L. (2014). Modeling dependent effect sizes with three-level meta-analyses: a structural equation modeling approach. *Psychol Methods* 19, 211–229. 10.1037/a0032968.

- 945 46. MacArthur, J., Bowler, E., Cerezo, M., Gil, L., Hall, P., Hastings, E., Junkins, H., McMahon, A., Milano, A., Morales, J., et al. (2017). The new NHGRI-EBI Catalog of published genome-wide association studies (GWAS Catalog). *Nucleic Acids Research*. 10.1093/nar/gkw1133.
47. Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., and Lee, J.J. (2015). Second-generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience* 4, s13742-015.
- 950 48. Martin, A.R., Kanai, M., Kamatani, Y., Okada, Y., Neale, B.M., and Daly, M.J. (2019). Clinical use of current polygenic risk scores may exacerbate health disparities. *Nat Genet* 51, 584–591. 10.1038/s41588-019-0379-x.
- 955 49. Lambert, S.A., Gil, L., Jupp, S., Ritchie, S.C., Xu, Y., Buniello, A., McMahon, A., Abraham, G., Chapman, M., Parkinson, H., et al. (2021). The Polygenic Score Catalog as an open database for reproducibility and systematic evaluation. *Nature Genetics* 53, 420–425.
50. Ruan, Y., Lin, Y.-F., Feng, Y.-C.A., Chen, C.-Y., Lam, M., Guo, Z., He, L., Sawa, A., Martin, A.R., Qin, S., et al. (2022). Improving polygenic prediction in ancestrally diverse populations. *Nat Genet* 54, 573–580. 10.1038/s41588-022-01054-7.
- 960 51. Cai, M., Xiao, J., Zhang, S., Wan, X., Zhao, H., Chen, G., and Yang, C. (2021). A unified framework for cross-population trait prediction by leveraging the genetic correlation of polygenic traits. *The American Journal of Human Genetics* 108, 632–655. 10.1016/j.ajhg.2021.03.002.
- 965 52. Hoggart, C.J., Choi, S.W., García-González, J., Souaiaia, T., Preuss, M., and O’Reilly, P.F. (2024). BridgePRS leverages shared genetic effects across ancestries to increase polygenic risk score portability. *Nat Genet* 56, 180–186. 10.1038/s41588-023-01583-9.
53. Privé, F., Albiñana, C., Arbel, J., Pasaniuc, B., and Vilhjálmsson, B.J. (2023). Inferring disease architecture and predictive ability with LDpred2-auto. *The American Journal of Human Genetics*. 10.1016/j.ajhg.2023.10.010.
- 970 54. Ding, Y., Hou, K., Xu, Z., Pimplaskar, A., Petter, E., Boulier, K., Privé, F., Vilhjálmsson, B.J., Olde Loohuis, L.M., and Pasaniuc, B. (2023). Polygenic scoring accuracy varies across the genetic ancestry continuum. *Nature* 618, 774–781. 10.1038/s41586-023-06079-4.
- 975 55. Norland, K., Schaid, D.J., and Kullo, I.J. (2023). A linear weighted combination of polygenic scores for a broad range of traits improves prediction of coronary heart disease. *Eur J Hum Genet*, 1–6. 10.1038/s41431-023-01463-0.
56. Krapohl, E., Patel, H., Newhouse, S., Curtis, C.J., von Stumm, S., Dale, P.S., Zabaneh, D., Breen, G., O’Reilly, P.F., and Plomin, R. (2018). Multi-polygenic score approach to trait prediction. *Mol Psychiatry* 23, 1368–1374. 10.1038/mp.2017.163.
- 980 57. Albiñana, C., Zhu, Z., Schork, A.J., Ingason, A., Aschard, H., Brikell, I., Bulik, C.M., Petersen, L.V., Agerbo, E., Grove, J., et al. (2023). Multi-PGS enhances polygenic prediction by combining 937 polygenic scores. *Nat Commun* 14, 4702. 10.1038/s41467-023-40330-w.

58. Howie, B., Fuchsberger, C., Stephens, M., Marchini, J., and Abecasis, G.R. (2012). Fast and accurate genotype imputation in genome-wide association studies through pre-phasing. *Nat Genet* *44*, 955–959. 10.1038/ng.2354.
59. McCarthy, S., Das, S., Kretzschmar, W., Delaneau, O., Wood, A.R., Teumer, A., Kang, H.M., Fuchsberger, C., Danecek, P., Sharp, K., et al. (2016). A reference panel of 64,976 haplotypes for genotype imputation. *Nat Genet* *48*, 1279–1283. 10.1038/ng.3643.
60. Huang, J., Howie, B., McCarthy, S., Memari, Y., Walter, K., Min, J.L., Danecek, P., Malerba, G., Trabetti, E., Zheng, H.-F., et al. (2015). Improved imputation of low-frequency and rare variants using the UK10K haplotype reference panel. *Nat Commun* *6*, 8111. 10.1038/ncomms9111.
61. Loh, P.-R., Danecek, P., Palamara, P.F., Fuchsberger, C., A Reshef, Y., K Finucane, H., Schoenherr, S., Forer, L., McCarthy, S., Abecasis, G.R., et al. (2016). Reference-based phasing using the Haplotype Reference Consortium panel. *Nat Genet* *48*, 1443–1448. 10.1038/ng.3679.
62. Guo, Y., He, J., Zhao, S., Wu, H., Zhong, X., Sheng, Q., Samuels, D.C., Shyr, Y., and Long, J. (2014). Illumina human exome genotyping array clustering and quality control. *Nat Protoc* *9*, 2643–2662. 10.1038/nprot.2014.174.
63. Das, S., Forer, L., Schön herr, S., Sidore, C., Locke, A.E., Kwong, A., Vrieze, S.I., Chew, E.Y., Levy, S., McGue, M., et al. (2016). Next-generation genotype imputation service and methods. *Nature genetics* *48*, 1284–1287.
64. Browning, B.L., Zhou, Y., and Browning, S.R. (2018). A One-Penny Imputed Genome from Next-Generation Reference Panels. *The American Journal of Human Genetics* *103*, 338–348. 10.1016/j.ajhg.2018.07.015.
65. Taliun, D., Harris, D.N., Kessler, M.D., Carlson, J., Szpiech, Z.A., Torres, R., Taliun, S.A.G., Corvelo, A., Gogarten, S.M., Kang, H.M., et al. (2021). Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature* *590*, 290–299.
66. Malik, R., Chauhan, G., Traylor, M., Sargurupremraj, M., Okada, Y., Mishra, A., Rutten-Jacobs, L., Giese, A.-K., van der Laan, S.W., Gretarsdottir, S., et al. (2018). Multiancestry genome-wide association study of 520,000 subjects identifies 32 loci associated with stroke and stroke subtypes. *Nat Genet* *50*, 524–537. 10.1038/s41588-018-0058-3.
67. Schwartzenuber, J., Cooper, S., Liu, J.Z., Barrio-Hernandez, I., Bello, E., Kumasaka, N., Young, A.M.H., Franklin, R.J.M., Johnson, T., Estrada, K., et al. (2021). Genome-wide meta-analysis, fine-mapping and integrative prioritization implicate new Alzheimer’s disease risk genes. *Nat Genet* *53*, 392–402. 10.1038/s41588-020-00776-w.
68. Ha, E., Bae, S.-C., and Kim, K. (2021). Large-scale meta-analysis across East Asian and European populations updated genetic architecture and variant-driven biology of rheumatoid arthritis, identifying 11 novel susceptibility loci. *Annals of the Rheumatic Diseases* *80*, 558–565. 10.1136/annrheumdis-2020-219065.

69. Michailidou, K., Lindström, S., Dennis, J., Beesley, J., Hui, S., Kar, S., Lemaçon, A., Soucy, P., Glubb, D., Rostamianfar, A., et al. (2017). Association analysis identifies 65 new breast cancer risk loci. *Nature* *551*, 92–94. 10.1038/nature24284.
- 1025 70. Schumacher, F.R., Al Olama, A.A., Berndt, S.I., Benlloch, S., Ahmed, M., Saunders, E.J., Dadaev, T., Leongamornlert, D., Anokian, E., Cieza-Borrella, C., et al. (2018). Association analyses of more than 140,000 men identify 63 new prostate cancer susceptibility loci. *Nat Genet* *50*, 928–936. 10.1038/s41588-018-0142-8.
- 1030 71. de Lange, K.M., Moutsianas, L., Lee, J.C., Lamb, C.A., Luo, Y., Kennedy, N.A., Jostins, L., Rice, D.L., Gutierrez-Achury, J., Ji, S.-G., et al. (2017). Genome-wide association study implicates immune activation of multiple integrin genes in inflammatory bowel disease. *Nat Genet* *49*, 256–261. 10.1038/ng.3760.
- 1035 72. Wuttke, M., Li, Y., Li, M., Sieber, K.B., Feitosa, M.F., Gorski, M., Tin, A., Wang, L., Chu, A.Y., Hoppmann, A., et al. (2019). A catalog of genetic loci associated with kidney function from analyses of a million individuals. *Nat Genet* *51*, 957–972. 10.1038/s41588-019-0407-x.
- 1040 73. Sakaue, S., Kanai, M., Tanigawa, Y., Karjalainen, J., Kurki, M., Koshihara, S., Narita, A., Konuma, T., Yamamoto, K., Akiyama, M., et al. (2021). A cross-population atlas of genetic associations for 220 human phenotypes. *Nat Genet* *53*, 1415–1424. 10.1038/s41588-021-00931-x.
74. Locke, A.E., Kahali, B., Berndt, S.I., Justice, A.E., Pers, T.H., Day, F.R., Powell, C., Vedantam, S., Buchkovich, M.L., Yang, J., et al. (2015). Genetic studies of body mass index yield new insights for obesity biology. *Nature* *518*, 197–206. 10.1038/nature14177.
- 1045 75. Hoffmann, T.J., Theusch, E., Haldar, T., Ranatunga, D.K., Jorgenson, E., Medina, M.W., Kvale, M.N., Kwok, P.-Y., Schaefer, C., Krauss, R.M., et al. (2018). A large electronic-health-record-based genome-wide study of serum lipids. *Nat Genet* *50*, 401–413. 10.1038/s41588-018-0064-5.



1050 123456789101112131415112131415265758Table 39310311512131415161615