

1 **PennPRS: a centralized cloud computing platform for efficient**  
2 **polygenic risk score training in precision medicine**

3

4 **Running title: PennPRS Platform**

5

6 Jin Jin<sup>1,2,\*</sup>, Bingxuan Li<sup>3</sup>, Xiyao Wang<sup>4</sup>, Xiaochen Yang<sup>5</sup>, Yujue Li<sup>5</sup>, Ruofan Wang<sup>1</sup>,  
7 Chenglong Ye<sup>6</sup>, Juan Shu<sup>7</sup>, Zirui Fan<sup>7</sup>, Fei Xue<sup>5</sup>, Tian Ge<sup>8</sup>, Marylyn D. Ritchie<sup>9,10</sup>,  
8 Bogdan Pasaniuc<sup>9</sup>, Genevieve Wojcik<sup>11</sup>, and Bingxin Zhao<sup>2,5,7,10,\*</sup>

9

10 <sup>1</sup>Department of Biostatistics, Epidemiology and Informatics, Perelman School of Medicine,  
11 University of Pennsylvania, Philadelphia, PA 19104, USA.

12 <sup>2</sup>Penn Center for Eye-Brain Health, Perelman School of Medicine, University of  
13 Pennsylvania, Philadelphia, PA 19104, USA.

14 <sup>3</sup>UCLA Samueli School of Engineering, Los Angeles, CA 90095, USA.

15 <sup>4</sup>Department of Computer Science, Columbia University, New York, NY 10027, USA.

16 <sup>5</sup>Department of Statistics, Purdue University, West Lafayette, IN 47907, USA.

17 <sup>6</sup>Department of Statistics, University of Kentucky, Lexington, KY 40536, USA.

18 <sup>7</sup>Department of Statistics and Data Science, University of Pennsylvania, Philadelphia, PA  
19 19104, USA.

20 <sup>8</sup>Psychiatric and Neurodevelopmental Genetics Unit, Center for Genomic Medicine,  
21 Massachusetts General Hospital, Boston, MA 02114, USA.

22 <sup>9</sup>Department of Genetics, Perelman School of Medicine, University of Pennsylvania,  
23 Philadelphia, PA 19104, USA.

24 <sup>10</sup>Penn Institute for Biomedical Informatics, Perelman School of Medicine, University of  
25 Pennsylvania, Philadelphia, PA 19104, USA.

26 <sup>11</sup>Department of Epidemiology, Johns Hopkins Bloomberg School of Public Health,  
27 Baltimore, MD 21205, USA.

28

29 *\*Corresponding authors:*

30 Jin Jin ([jin.jin@penmedicine.upenn.edu](mailto:jin.jin@penmedicine.upenn.edu)) and

31 Bingxin Zhao ([bxzhao@wharton.upenn.edu](mailto:bxzhao@wharton.upenn.edu))

1 **Abstract**

2 Polygenic risk scores (PRS) are becoming increasingly vital for risk prediction and  
3 stratification in precision medicine. However, PRS model training presents  
4 significant challenges for broader adoption of PRS, including limited access to  
5 computational resources, difficulties in implementing advanced PRS methods, and  
6 availability and privacy concerns over individual-level genetic data. Cloud  
7 computing provides a promising solution with centralized computing and data  
8 resources. Here we introduce PennPRS (<https://pennprs.org>), a scalable cloud  
9 computing platform for online PRS model training in precision medicine. We  
10 developed novel pseudo-training algorithms for multiple PRS methods and  
11 ensemble approaches, enabling model training without requiring individual-level  
12 data. These methods were rigorously validated through extensive simulations and  
13 large-scale real data analyses involving over 6,000 phenotypes across various  
14 data sources. PennPRS supports online single- and multi-ancestry PRS training  
15 with seven methods, allowing users to upload their own data or query from more  
16 than 27,000 datasets in the GWAS Catalog, submit jobs, and download trained  
17 PRS models. Additionally, we applied our pseudo-training pipeline to train PRS  
18 models for over 8,000 phenotypes and made their PRS weights publicly  
19 accessible. In summary, PennPRS provides a novel cloud computing solution to  
20 improve the accessibility of PRS applications and reduce disparities in  
21 computational resources for the global PRS research community.

22

23 **Keywords:** Cloud computing; GWAS Catalog; Polygenic risk scores; Precision  
24 medicine; Resampling-based pseudo-training.

## 1 Introduction

2 The last two decades have seen remarkable growth in genome-wide association  
3 studies (GWAS), yielding extensive data resources valuable for genetic risk  
4 prediction<sup>1,2</sup>. Polygenic risk scores (PRS), calculated as the sum of the number of  
5 alleles of genetic variants weighted by their effect sizes, encapsulate cumulative  
6 genome-wide risks for complex traits and diseases<sup>3-5</sup>. Numerous studies have  
7 highlighted the utility of PRS in precision medicine to help disease risk stratification  
8 and inform clinical intervention decisions<sup>6-9</sup>. To improve the accuracy and  
9 robustness of PRS, a wide range of methods, software, standards, and web  
10 resources have been developed<sup>3,10-14</sup>. Recent initiatives aim to further extend PRS  
11 applications to more diverse and admixed global populations<sup>15-18</sup>. Such efforts  
12 have been reflected by the establishment of a series of NHGRI-funded consortia,  
13 including the PRIMED<sup>19</sup>, which aims to develop and evaluate methods to improve  
14 the use of PRS for predicting disease risks in diverse ancestry populations, and  
15 the eMERGE Network<sup>20,21</sup>, which supports genomic medicine translation by  
16 returning PRS results to individuals along with healthcare recommendations in  
17 diverse clinical settings. The combination of methodological advancements,  
18 increasingly rich discovery GWAS data, and decreasing costs in biotechnology are  
19 anticipated to persistently and substantially improve both the capabilities and  
20 accessibility of PRS-based disease risk prediction and stratification.

21

22 The accessibility and scalability of PRS applications, however, are often hindered  
23 by significant challenges in the PRS model training process, particularly for users  
24 of advanced PRS algorithms (**Fig. 1a**). For example, access to high-performance  
25 computational resources required to run these algorithms and store large-scale  
26 GWAS summary data is often dependent on existing institutional infrastructure,  
27 which may not be readily available to all PRS researchers across diverse  
28 organizations and scientific fields. Additionally, managing and testing various PRS  
29 methods within local pipelines can involve a steep learning curve and make it  
30 difficult to keep up with the frequent updates to new methods. A further  
31 complication arises from the need for an independent individual-level dataset  
32 during the training process, which is typically used as tuning data for optimizing

1 model parameters and training ensemble models. This dataset must be sufficiently  
2 large and independent from the one used to generate GWAS summary statistics.  
3 Due to privacy concerns surrounding the sharing of individual-level genetic data,  
4 obtaining such a dataset can present logistical challenges in PRS applications.  
5 Furthermore, for certain traits, even when individual-level datasets are available,  
6 their sample sizes may be insufficient to produce reliable and stable parameter  
7 tuning results.

8

9 The use of cloud computing is gaining increasing momentum in biomedical  
10 research<sup>22-28</sup>, especially with centralized research resources, given the energy-  
11 efficient nature of cloud computing for hosting large-scale computing and data  
12 sources with high scalability and security<sup>29</sup>. Several biobanks have recently  
13 developed study-centric cloud computing platforms, such as the UK Biobank  
14 Research Analysis Platform (<https://ukbiobank.dnanexus.com/>) and the All of Us  
15 Researcher Workbench (<https://www.researchallofus.org/data-tools/workbench/>), to  
16 increase their data accessibility across diverse research communities. Cloud  
17 computing provides a promising next-generation solution to address the  
18 challenges in the widespread expansion of PRS applications. By leveraging robust  
19 online resources (such as Amazon Web Services [AWS]), cloud computing can  
20 provide a well-organized platform for diverse PRS users, facilitating efficient data  
21 analysis through centralized data storage and unified pipelines. However, a key  
22 barrier to implementing PRS model training on online servers is the reliance of  
23 many PRS methods on individual-level genetic data, which raises concerns about  
24 availability and data privacy. Recent advances have introduced the pseudo-  
25 training approach for PRS model development<sup>30-33</sup>. This approach allows for the  
26 sampling of pseudo-training and validation summary statistics from the underlying  
27 probability distribution of GWAS summary statistics. These sampled statistics  
28 closely mimic what would be obtained if there were access to two subsets of the  
29 GWAS samples, enabling parameter tuning and the derivation of PRS models.  
30 This “self-training” approach makes it possible to generate PRS weights without  
31 the need for individual-level genetic data.

32

1 Building on these advancements, this paper aims to integrate cloud computing with  
2 pseudo-training approaches to enable online training of PRS models, providing a  
3 secure, efficient, and scalable solution for the PRS research community. We first  
4 developed pseudo-training versions for multiple single- and multi-ancestry PRS  
5 methods and rigorously showed their robust performance across thousands of  
6 phenotypes from various data sources. Based on their reliable numerical  
7 performance, we introduced PennPRS (<https://pennprs.org/>), a scalable cloud  
8 computing platform for online training of PRS models using summary statistics only  
9 (**Fig. 1b**). PennPRS provides a wide range of user options and supports both  
10 single- and multi-ancestry analyses across the five super populations<sup>34</sup>, including  
11 European (EUR), African and African American (AFR), Admixed American (AMR),  
12 East Asian (EAS), and South Asian (SAS). Users can input GWAS summary  
13 statistics, submit a job with selected methods and customized settings, and  
14 download the trained PRS models upon completion. As a centralized PRS online  
15 training platform, PennPRS provides cloud computing functionalities, extensive  
16 data resources, and offline pipelines, offering an efficient solution to PRS model  
17 development in precision medicine (**Fig. 2a**). It is designed to accommodate the  
18 training needs of various research groups with diverse requirements and  
19 computational resources.

20

## 21 **Results**

### 22 **Summary data-based PRS model tuning and ensemble learning**

23 We developed single-ancestry pseudo-training pipelines with summary data-  
24 based parameter tuning for three PRS methods: Clumping and Thresholding  
25 (C+T)<sup>4,5</sup>, Lassosum2<sup>35,36</sup>, and LDpred2<sup>37,38</sup>, which we denote by C+T-pseudo,  
26 Lassosum2-pseudo, and LDpred2-pseudo, respectively (**Fig. 2b**). The PUMAS<sup>31</sup>  
27 workflow was used to derive pseudo training and validation subsamples from the  
28 input GWAS summary statistics, enabling the selection of optimal tuning  
29 parameters. While the general framework of PUMAS pseudo-training has been  
30 established, applying it to specific PRS methods is challenging due to the  
31 complexities involved in implementing different PRS methods and non-universality  
32 of the original PUMAS software to the general GWAS summary data. Therefore,

1 we made a series of important modifications to both the methodology and the  
2 software to ensure proper alignment between the PUMAS algorithm and each of  
3 the implemented PRS methods. For example, the original Lassosum2 and  
4 LDpred2 algorithms may generate non-convergent PRS weights under some  
5 tuning parameter settings, which can result in problematic parameter tuning. To  
6 address this potential issue, we developed a data-driven approach to improve  
7 robustness of the summary data-based parameter optimization in Lassosum2-  
8 pseudo and LDpred2-pseudo (see **Methods**).

9

10 For single-ancestry analysis, we further developed two ensemble approaches  
11 within our pseudo-training framework: Ensemble-pseudo and Ensemble-ARM-  
12 pseudo. These approaches combine PRS models trained by different methods  
13 using a linear combination strategy<sup>39</sup> (Ensemble-pseudo) or an adaptive  
14 regression by mixing approach<sup>40</sup> (Ensemble-ARM-pseudo) (**Fig. 2b**). The two  
15 ensemble learning methods were originally designed for use with the need for  
16 individual-level tuning datasets. Here we redeveloped them for pseudo-training  
17 within the PUMA-CUBS<sup>31</sup> framework, incorporating a series of modifications to  
18 ensure robustness. Details are provided in the **Methods** section.

19

20 Notably, our pipeline supports multi-ancestry PRS training based on ancestry-  
21 stratified GWAS summary data from multiple ancestral populations, a process that  
22 is typically computationally intensive and requires more learning resources (**Fig.**  
23 **2c**). We developed PROSPER-pseudo, a pseudo-training pipeline for PROSPER<sup>41</sup>  
24 which is an ensemble learning-assisted multi-ancestry PRS method. PROSPER-  
25 pseudo will generate two complementary population-specific models: PROSPER-  
26 Single-pseudo and PROSPER-pseudo, where the former provides the best single  
27 PRS generated before the final ensemble step in PROSPER, and the latter  
28 provides the final PRS that combines multiple PRS across different ancestries and  
29 tuning parameter settings (see **Methods**). In summary, we have developed three  
30 single-ancestry methods, two ensemble approaches, and one multi-ancestry  
31 method as the primary methods for implementation on our cloud computing  
32 platform, all based on pseudo-training that eliminates the need for individual-level

1 data. A summary of our pseudo-training algorithms for implementing these single-  
2 and multi-ancestry PRS methods is provided in **Supplementary Fig. 1**.  
3 Additionally, we have developed several complementary methods for the offline  
4 pipeline, which will be introduced in later sections.

5

## 6 **Large-scale evaluation of the PRS pseudo-training approach**

7 We evaluated the performance of our PRS pseudo-training methods through  
8 extensive simulations and real data analyses. Simulation results of the single-  
9 ancestry methods under various settings of genetic architecture of the phenotype  
10 (such as heritability and causal variant proportion) and GWAS sample size<sup>42</sup>  
11 demonstrate that our PRS pseudo-training methods perform comparably to those  
12 original methods that tuned model parameters with a sufficiently large hold-out  
13 individual-level dataset (e.g.,  $N_{tuning} = 1000$ , **Fig. 3a**, **Supplementary Fig. 2**, and  
14 **Supplementary Table 1**). As training GWAS sample size increases, our pseudo-  
15 training methods tend to better approximate the PRS under the optimal tuning  
16 parameter setting. Pseudo-training versions of the ensemble PRS tend to have  
17 slightly lower prediction R-squared ( $R^2$ ) than individual-level data-based ensemble.  
18 It is important to note that the above comparisons assume access to sufficiently  
19 large individual-level tuning data. However, when the number of individual-level  
20 tuning samples is insufficient (e.g.,  $N_{tuning} < 1,000$ ), pseudo PRS training notably  
21 outperforms traditional PRS training methods that rely on individual-level tuning  
22 data (**Figs. 3b** and **3c**). For example, compared to the traditional PRS training  
23 pipelines using an individual-level tuning dataset of size  $N_{tuning} = 400$  or 100, our  
24 PRS pseudo-training pipeline for the same PRS methods achieved a 6.5% or 44.7%  
25 higher  $R^2$ , respectively. These findings highlight the utility of the PRS pseudo-  
26 training methods we developed, particularly in scenarios where individual-level  
27 data is limited or unavailable for parameter tuning.

28

29 We examined the performance of our pseudo-training pipelines for single-ancestry  
30 PRS model training across different phenotypes and data sources (**Fig. 2d**). First,  
31 we used 2,106 multi-organ multi-modal imaging phenotypes with GWAS summary  
32 data available from the UK Biobank (UKB)<sup>43</sup> study, including 1,432 brain structural

1 magnetic resonance imaging (sMRI)<sup>44</sup>, 674 diffusion MRI (dMRI), 82 resting-state  
2 functional MRI (rfMRI), 41 abdominal MRI<sup>45</sup>, 82 cardiovascular MRI<sup>46</sup>, and 46 eye  
3 optical coherence tomography images (OCT)<sup>47</sup> (average  $N = 32,859$ ). These  
4 imaging phenotypes are well-established clinical biomarkers with widespread  
5 practical applications in precision medicine<sup>48</sup>. We found that all three single-  
6 ancestry pseudo-training methods, as well as the two pseudo-training ensemble  
7 approaches, demonstrated strong performance across these diverse imaging  
8 modalities (**Figs. 4a and 5a, Supplementary Tables 2-5**; mean  $R^2 = 0.0562$  vs.  
9  $0.0567$ ,  $R^2$  correlation =  $0.955$ ). Consistent with our simulation studies, we  
10 observed that pseudo-training methods outperform individual-level data-based  
11 tuning as the individual-level tuning sample size decreases (**Fig. 4b and**  
12 **Supplementary Fig. 3**). For example, with a tuning sample size  $N_{tuning} = 1,000$ ,  
13 300, and 100, our pseudo-training methods produced PRS with  $R^2$  values that were  
14 0.4%, 6.9%, and 18.5% higher, respectively, compared to methods using  
15 individual-level tuning data for eye OCT phenotypes (**Fig. 4b**).

16  
17 Furthermore, we trained PRS for 29 binary disease phenotypes based on GWAS  
18 summary statistics from the FinnGen<sup>49</sup> study and evaluated their performance on  
19 matched clinical outcomes using UKB testing individuals<sup>50</sup> (average  $N_{case} = 23,048$ ,  
20 **Supplementary Table 6**). Again, all three single-ancestry pseudo-training  
21 methods demonstrated performance comparable to the traditional methods using  
22 individual-level tuning data (area under the ROC curve [AUC] correlations =  $0.880$ ).  
23 Notably, the pseudo-training ensembles outperformed the individual-level data  
24 ensembles, even though the individual-level tuning datasets are large (**Fig. 5b and**  
25 **Supplementary Table 7**; mean AUC =  $0.564$  vs.  $0.555$ , one-sided  $P = 1.51 \times 10^{-7}$ ).  
26 In addition, we evaluated the PRS performance on 2,734 Olink plasma proteins  
27 with GWAS data from the UKB-PPP project<sup>51</sup> (average  $N = 40,852$ , **Fig. 5c and**  
28 **Supplementary Tables 8-9**). Plasma proteins, which are crucial for disease  
29 diagnosis and treatment<sup>52,53</sup>, exhibit a unique special architecture, generally  
30 showing higher heritability and with *cis*-loci accounting for a large proportion of  
31 genetic variation<sup>51</sup>. For such genetic architecture, our analysis revealed that C+T-  
32 pseudo and LDpred2-pseudo showed highly consistent performance with training



1 based on individual-level tuning data (mean  $R^2 = 0.0562$  vs.  $0.0567$ ,  $R^2$  correlation  
2 =  $0.998$ ), whereas Lassosum2-pseudo consistently delivered sub-optimal  
3 performance for proteins with high prediction  $R^2$  (e.g.,  $> 0.40$ ) (mean  $R^2 = 0.475$   
4 vs.  $0.557$ ). These findings suggest that C+T-pseudo and LDpred2-pseudo may be  
5 more suitable for deriving genetic scores<sup>53</sup> for these proteins and other molecular  
6 traits with similar genetic architecture.

7

8 For multi-ancestry PRS training, our simulation studies suggested that PROSPER-  
9 Single-pseudo, the pseudo-trained best single PROSPER PRS without  
10 implementing the final ensemble step, approximates its individual-level data-based  
11 version (PROSPER-Single) well, while the PROSPER-pseudo PRS (with the final  
12 ensemble step) tends to perform slightly worse than PROSPER PRS, its individual-  
13 level data-based version, if large hold-out individual-level dataset exists (**Fig. 6a**  
14 and **Supplementary Table 10**). We further evaluated their performance in multi-  
15 ancestry real data applications using ancestry-stratified GWAS summary statistics  
16 (**Supplementary Tables 11-14**). We first analyzed four blood lipids, including high-  
17 density lipoprotein (HDL), low-density lipoprotein (LDL), log-transformed  
18 triglycerides (logTG), and total cholesterol (TC). We used GWAS summary data  
19 for EUR, AFR, AMR, EAS, and SAS from the Global Lipids Genetics Consortium<sup>54</sup>  
20 (GLGC) ( $N = 33,658-930,671$ ). The performance was evaluated on UKB validation  
21 individuals of EUR, AFR, AMR, EAS, and SAS ancestries<sup>54,55</sup>. Our results showed  
22 that pseudo-training had consistent performance across all ancestries (**Fig. 6b**,  
23 mean  $R^2$ :  $0.084$  [PROSPER-Single-pseudo] and  $0.088$  [PROSPER-pseudo] vs.  
24  $0.089$  [PROSPER],  $R^2$  correlation =  $0.93$  and  $0.95$ , respectively). We also  
25 evaluated the performance of multi-ancestry pseudo-training using GWAS  
26 summary statistics for 1,413 brain dMRI and sMRI phenotypes from the Chinese  
27 Imaging Genetics (CHIMGEN) study<sup>56</sup>, jointly with matched imaging phenotypes  
28 in the UKB study. Specifically, the inputs were the CHIMGEN summary statistics  
29 (average  $N = 7,058$ ) and UKB European summary statistics (average  $N = 32,620$ ),  
30 with performance evaluated in independent testing data from hold-out UKB  
31 samples (average  $N = 2,510$  for EUR ancestry and  $N = 222$  for EAS ancestry). We  
32 found that pseudo-training outperformed individual-level data training for both EUR

1 ( $R^2 = 0.027$  vs.  $0.023$ ), which has sufficient tuning samples, and EAS ( $R^2 = 0.010$   
2 vs.  $0.008$ ), which has limited tuning samples, although the results for EAS had  
3 larger uncertainty due to the much smaller testing sample sizes (**Fig. 6b** and  
4 **Supplementary Fig. 4**). As expected, analyses of GLGC and CHIMGEN data also  
5 showed improved prediction accuracy when GWAS training data from both  
6 ancestries were integrated (**Fig. 6c**) and the consistent pattern of relative  
7 performance of PROSPER-Single-pseudo and PROSPER-pseudo across  
8 different data resources (**Supplementary Fig. 5**).

9

10 Overall, our large-scale numerical results demonstrate that, without access to an  
11 independent individual-level tuning dataset, the developed summary data-only  
12 pseudo-training methods can produce PRS weights comparable to those  
13 generated using a large individual-level tuning dataset. Furthermore, pseudo-  
14 training may even outperform traditional methods, particularly when the tuning  
15 dataset is limited in size. These findings lay the methodological groundwork for the  
16 development of PennPRS as a centralized cloud computing solution for online  
17 PRS model training.

18

### 19 **PennPRS: a cloud computing platform for the global PRS community**

20 We established a centralized cloud computing platform hosted on AWS to  
21 implement the developed pseudo-training methods, enabling users to freely train  
22 PRS weights online using GWAS summary statistics (**Fig. 1b**). Upon completing  
23 registration, users can upload GWAS summary statistics or query over 27,000  
24 harmonized summary statistics from the GWAS Catalog<sup>57</sup>, select PRS methods  
25 and model parameters, and submit jobs. These jobs are managed by the queuing  
26 system and processed on our server, and users can download the generated PRS  
27 weights and log files once the job is completed. In addition to the set of newly  
28 developed pseudo-training PRS methods (C+T-pseudo, Lassosum2-pseudo,  
29 LDpred2-pseudo, and PROSPER-pseudo) and pseudo-training ensemble  
30 methods (Ensemble-pseudo and Ensemble-ARM-pseudo), PennPRS supports  
31 three existing tuning-parameter-free methods (PRS-CS-auto<sup>58</sup>, LDpred2-auto<sup>37</sup>,  
32 and DBSLMM<sup>59</sup>). Our platform presents a robust frontend-to-backend web

1 infrastructure with detailed tutorials and a comprehensive data harmonization  
2 pipeline to ensure regularized PRS training and an efficient user experience (see  
3 **Methods**).

4  
5 Similar to many biomedical cloud computing platforms in other fields<sup>22-27</sup>, the  
6 PennPRS team will cover data analysis expenses for all users. This approach  
7 aligns with our mission to make PRS accessible to more researchers and,  
8 ultimately, to study participants in precision medicine, while reducing disparities in  
9 computational resources within the global PRS research community and the  
10 broader fields of genetic and medical research. To optimize the efficiency of our  
11 computational infrastructure, we conducted various tests to determine the optimal  
12 configurations for our platform, such as RAM and CPU deployment for different  
13 PRS methods implemented. We also conducted extensive tests to validate the  
14 platform's stability and computational performance. For example, using the three  
15 single-ancestry pseudo-training methods and their two ensemble approaches as a  
16 case study, we analyzed the runtime for each step in the algorithm. We found that  
17 using 2 CPUs (with 30 GB RAM) allowed a typical job to complete in approximately  
18 two and a half hours, while increasing to 4 CPUs reduced the runtime to around  
19 two hours (**Supplementary Table 15**). Based on these empirical observations, we  
20 have optimized our configuration to make efficient use of AWS resources, currently  
21 supporting up to eight concurrent user jobs. The AWS cloud computing service,  
22 additionally supported by our local computing IT teams, provides a flexible  
23 management system for CPU and RAM, enabling us to easily maintain the server  
24 and adjust resource allocation for scaling up or down as needed.

### 25 26 **Public sources: GWAS Catalog data querying and working examples**

27 The GWAS Catalog<sup>57</sup> has become an invaluable database of public GWAS  
28 summary statistics, with a fast-growing collection of data curated and harmonized  
29 for post-GWAS applications. We have developed a feature that links PennPRS  
30 directly to the GWAS Catalog database to enable efficient PRS model training.  
31 This allows users to query data from the GWAS Catalog for PRS pseudo training  
32 directly without the need to download, preprocess, and upload the data to

1 PennPRS. To enable this functionality with high quality, we focused on harmonized  
2 datasets from the GWAS Catalog and ensured that the provided data meet the  
3 basic requirements for implementing the various PRS methods supported, such as  
4 having the necessary GWAS summary-level data information and excluding data  
5 from exome studies. Currently, we provide access to over 27,000 datasets for  
6 users to query directly through PennPRS.

7

8 We provide two examples of querying GWAS Catalog datasets for efficient PRS  
9 pseudo-training on PennPRS. The first example demonstrates the use of the  
10 PennPRS single-ancestry data analysis pipeline to train a PRS model for height in  
11 Hispanics. In this example, we first navigated to the PennPRS GWAS Queryable  
12 Database (<https://pennprs.org/data>) and searched for “height”. We identified the  
13 dataset from the study “GCST90095033” as a suitable input for PRS training,  
14 which provided GWAS summary statistics from 59,771 Hispanic or Latin American  
15 individuals<sup>60</sup>. We then created a single-ancestry data analysis job on PennPRS,  
16 entering the relevant dataset information (e.g., study accession ID) to enable direct  
17 querying from the GWAS Catalog. Next, we selected three pseudo-training  
18 methods (C+T-pseudo, Lassosum2-pseudo, and LDpred2-pseudo) along with the  
19 ensemble option, which would utilize two ensemble methods, Ensemble-pseudo  
20 and Ensemble-ARM-pseudo, to train two ensemble PRS models combining PRS  
21 trained by the three selected methods and submitted the job. PennPRS completed  
22 the job in approximately two and a half hours, returning five PRS models along  
23 with a detailed log file. Similarly, the second example demonstrates the use of the  
24 PennPRS multi-ancestry data analysis pipeline to train PRS models for height  
25 across four ancestries (EUR, AFR, AMR, and EAS). We queried four  
26 corresponding ancestry-specific GWAS datasets from the GWAS Catalog  
27 (“GCST90029008”, “GCST90013468”, “GCST90095033”, and “GCST90018739”)  
28 and used the multi-ancestry method, PROSPER-pseudo, for online training.  
29 PennPRS completed this job in approximately ten and a half hours, generating  
30 eight PRS models for the four ancestries (PROSPER-Single-pseudo PRS and  
31 PROSPER-pseudo PRS for each ancestry). Detailed steps and illustrations for

1 these examples are available in the tutorial (<https://pennprs.gitbook.io/pennprs>),  
2 serving as quick-start guides for new users.

3

#### 4 **PennPRS offline pipeline and pretrained PRS models**

5 In addition to establishing the online PRS training server, we have developed a  
6 comprehensive pipeline for offline implementation of the supported PRS pseudo-  
7 training and tuning-parameter-free methods. The cloud computing server is  
8 designed as a convenient and eco-friendly<sup>29</sup> online tool for PRS users, particularly  
9 those who face challenges in setting up local computational clusters, while the  
10 offline pipeline is powerful for large-scale analyses if researchers have access to  
11 high-performance computing clusters. In our offline pipeline, we have additionally  
12 developed novel pseudo-training versions of three single- and multi-ancestry PRS  
13 methods, including PRS-CS-grid<sup>58</sup>-pseudo, PRS-CSx<sup>61</sup>-pseudo, and MUSSEL<sup>55</sup>-  
14 pseudo. Due to the nature of these methods, they have much higher memory and  
15 computational demands than other methods and are therefore included exclusively  
16 in our offline pipeline. By offering both online and offline options, we aim to  
17 accommodate the diverse needs of research groups and help reduce disparities in  
18 computational resources for the PRS application community.

19

20 To demonstrate the power and efficiency of our offline pipeline, we applied it to a  
21 wide range of phenotypes, including those mentioned in our model evaluations (for  
22 which we had access to individual-level testing data for performance assessment),  
23 as well as many more GWAS summary datasets from the GWAS Catalog<sup>57</sup>,  
24 Biobank Japan (BBJ)<sup>62</sup>, the Million Veteran Program (MVP) study<sup>63</sup>, and the Global  
25 Biobank Meta-analysis Initiative (GBMI) consortium<sup>64</sup>. Specifically, we have  
26 conducted single-ancestry analysis with all three single-ancestry pseudo-training  
27 methods (C+T-pseudo, Lassosum2-pseudo, and LDpred2-pseudo) and their  
28 ensembles (Ensemble-pseudo and Ensemble-ARM-pseudo) using default tuning  
29 parameter settings on over 8,000 harmonized GWAS Catalog datasets and 169  
30 phenotypes from BBJ. We have also conducted multi-ancestry analysis with  
31 PROSPER-pseudo on 181 ancestry-stratified GWAS summary datasets from the  
32 MVP and nine ancestry-stratified GWAS summary datasets from the GBMI. We

1 have made these pretrained PRS models publicly available in the PennPRS public  
2 resource hub (<https://pennprs.org/result>), allowing users to freely download and  
3 utilize them in their research. As the GWAS Catalog and other databases continue  
4 to expand, harmonizing and making more curated GWAS summary statistics  
5 publicly available, we will leverage the established PennPRS pipeline to analyze  
6 these summary datasets and share the trained PRS models with the PRS research  
7 community. These resources will accelerate the application of PRS across various  
8 fields.

9

## 10 **Discussion**

11 PRS training methods and their cluster-based implementations have been  
12 traditionally handled by local servers, typically established by individual research  
13 groups. However, the fast-paced evolution of PRS methodologies, along with the  
14 growing volume of GWAS resources, presents logistical, computational, and  
15 environmental challenges for hosting these PRS pipelines locally. This is  
16 particularly true for smaller research groups that may lack sufficient computational  
17 resources or are new to PRS. In this paper, we developed a series of pseudo-  
18 training algorithms, data resources, and cloud computing functionalities to enable  
19 online single- and multi-ancestry PRS pseudo-training using summary data only,  
20 eliminating the need for local setups. Our platform aims to lower the barriers to  
21 PRS applications across various phenotypes and ancestry populations, while also  
22 reducing disparities in computational resources within the global PRS research  
23 community.

24

25 The development of cloud computing platforms and centralized resources have  
26 provided significant environmental benefits<sup>29,65</sup> and opened new opportunities  
27 across the broad fields of biomedical data science<sup>22-28</sup>. The novel pseudo-training  
28 methods we developed provide several advantages for cloud-based PRS model  
29 training solutions. First, pseudo-training mitigates the privacy risks and concerns  
30 associated with uploading or sharing individual-level genetic datasets online.  
31 Second, individual-level validation data is not always available for PRS model  
32 development, and our pseudo-training pipeline provides a more accessible

1 solution for PRS training across a wider range of disease and health outcomes.  
2 Third, pseudo-training could deliver PRS with better prediction performance,  
3 especially for those disease outcomes with limited individual-level tuning data  
4 available (e.g.,  $N_{tuning} < 1,000$ ). The pseudo-training approach we developed would  
5 lend power for these understudied outcomes. Fourth, the pseudo-training  
6 approach facilitates seamless integration with online GWAS data resources, such  
7 as the GWAS Catalog, providing a centralized data resource for PRS model  
8 training. In the future, we plan to extend the functionality of pseudo-training and  
9 PennPRS in several directions. For example, the current training procedure relies  
10 on the five super ancestral population labels<sup>34</sup> (e.g., EUR and AFR). We aim to  
11 expand our framework to include additional population labels and further integrate  
12 flexible genetic ancestry continuum information<sup>66</sup> as the field increasingly  
13 incorporates diverse and admixed ancestry information<sup>18</sup>. We will also provide  
14 unified PRS models for the general population instead of ancestry-specific PRS  
15 models that require categorizing individuals into discrete ancestry groups<sup>15</sup>.  
16 Furthermore, beyond generating PRS model training, we will additionally develop  
17 pseudo-training methods to produce additional accuracy metrics and uncertainty  
18 measures for the generated PRS models, such as confidence intervals<sup>67</sup>, which  
19 are increasingly critical for downstream clinical applications of PRS<sup>68-71</sup>.

20

21 Notably, the applicability of the FAIR data principles<sup>72</sup> to our cloud computing  
22 platform highlights the broad impact of PennPRS on future translational research  
23 involving PRS. By providing standardized computing pipelines, curated data  
24 resources, detailed documentation, and accessible PRS weight files, PennPRS  
25 facilitates transparency and ensures the **F**indability, **A**ccessibility, **I**nteroperability,  
26 and **R**eusability of PRS resources. This is particularly important as the adoption  
27 and application of PRS continue to expand in precision medicine, a process that  
28 typically involves multiple steps, from PRS model development and assessment  
29 to implementation and translation in clinical settings. Efficient data and information  
30 sharing between these stages will be critical for successful translation. In addition,  
31 the computing methods and data resources provided by our PennPRS platform  
32 enable efficient between study comparisons of PRS within the same ancestry

1 background. This is particularly important given the diversity of biobanks in the US  
2 and globally. By providing a centralized platform, PennPRS anchors comparisons  
3 to a consistent linkage disequilibrium and allele frequency background, reducing  
4 bias or noise that might hinder cross-study comparisons. This allows researchers  
5 to focus on other factors affecting performance<sup>18,73</sup>, improving the reliability and  
6 generalizability of PRS analyses across studies.

7

8 There are also limitations to cloud computing platforms. For example, if users have  
9 already established powerful local computational clusters—typically supported by  
10 their research institutions' infrastructure—they might consider setting up the  
11 pipeline locally, although this approach may be more time consuming and less  
12 environmentally friendly<sup>65</sup>. To meet this complementary need, we have developed  
13 a ready-to-use offline version of PennPRS pipeline. Another challenge of cloud  
14 computing platform is the maintenance of the server and pipelines. To ensure  
15 sustainable development of our platform, we have formed a dedicated  
16 interdisciplinary team of researchers centered around the PRS research  
17 community at the University of Pennsylvania, in collaboration with researchers  
18 from other institutions, to support regular updates to PennPRS. These updates will  
19 include incorporating additional PRS methods, generating new data resources,  
20 and more efficient method implementations. The long-term maintenance of the  
21 platform is supported by the robust AWS server, with additional technical  
22 assistance from our local IT teams. We welcome user feedback and suggestions  
23 to improve the PennPRS platform and better meet the diverse needs of the global  
24 PRS research community.

25

## 26 **Acknowledgements**

27 Research reported in this publication was supported by the National Human  
28 Genome Research Institute under Award Number 5R00HG012223 (J.J.), National  
29 Institute of Mental Health under Award Number R01MH136055 (B.Z.), and  
30 National Institute on Aging under Award Number RF1AG082938 (B.Z.). The  
31 content is solely the responsibility of the authors and does not necessarily  
32 represent the official views of the National Institutes of Health. The study has also



1 been partially supported by funding from the Purdue University Statistics  
2 Department, Department of Statistics and Data Science at the University of  
3 Pennsylvania, Wharton Dean's Research Fund, Analytics at Wharton, Wharton AI  
4 & Analytics Initiative, Perelman School of Medicine CCEB Innovation Center  
5 Grant, and the University Research Foundation at the University of Pennsylvania.  
6 The individual-level genotype and phenotype data for UK Biobank validation  
7 samples used in this study were obtained under application 76139 subject to a  
8 data transfer agreement. We would like to thank the individuals who represented  
9 themselves in the UK Biobank for their participation and the research teams for  
10 their efforts in collecting, processing, and disseminating these datasets. We would  
11 like to thank the research computing and IT groups at the Wharton School of the  
12 University of Pennsylvania and the Rosen Center for Advanced Computing at the  
13 Purdue University for providing computational resources and support that have  
14 contributed to these research results.

15

#### 16 **Author Contributions Statement**

17 J.J. and B.Z. conceived the project. J.J. developed the pseudo-training methods  
18 and algorithms. B.L. and X.W. set up the cloud computing server with the help from  
19 local IT teams. J.J., B.L., X.W., and B.Z. carried out data analyses, interpreted the  
20 results, designed the cloud computing platform, and developed the offline pipeline.  
21 X.Y., Y.L., J.S., R.W., and Z.F. processed the GWAS summary statistics,  
22 developed the curated datasets, and contributed to the development and testing  
23 of the computing platform and offline pipeline. C.Y., F.X., T.G., M.D.R., G.W., B.P.,  
24 and R.W. contributed to the design of the methods, cloud computing platform, and  
25 offline pipeline. J.J. and B.Z. drafted the manuscript with feedback from all authors.

26

#### 27 **Competing Interests Statement**

28 The authors declare no competing interests.

29

#### 30 **References**

- 31 1. Uffelmann, E. *et al.* Genome-wide association studies. *Nature Reviews Methods*  
32 *Primers* **1**, 1-21 (2021).

- 1 2. Abdellaoui, A., Yengo, L., Verweij, K.J. & Visscher, P.M. 15 years of GWAS  
2 discovery: Realizing the promise. *The American Journal of Human Genetics*  
3 (2023).
- 4 3. Choi, S.W., Mak, T.S.-H. & O'Reilly, P.F. Tutorial: a guide to performing polygenic  
5 risk score analyses. *Nature protocols* **15**, 2759-2772 (2020).
- 6 4. Wray, N.R., Goddard, M.E. & Visscher, P.M. Prediction of individual genetic risk  
7 to disease from genome-wide association studies. *Genome research* **17**, 1520-  
8 1528 (2007).
- 9 5. Consortium, I.S. Common polygenic variation contributes to risk of schizophrenia  
10 and bipolar disorder. *Nature* **460**, 748-752 (2009).
- 11 6. Torkamani, A., Wineinger, N.E. & Topol, E.J. The personal and clinical utility of  
12 polygenic risk scores. *Nature Reviews Genetics* **19**, 581-590 (2018).
- 13 7. Kullo, I.J. *et al.* Polygenic scores in biomedical research. *Nature Reviews Genetics*  
14 **23**, 524-532 (2022).
- 15 8. Chatterjee, N., Shi, J. & García-Closas, M. Developing and evaluating polygenic  
16 risk prediction models for stratified disease prevention. *Nature Reviews Genetics*  
17 **17**, 392 (2016).
- 18 9. Lewis, C.M. & Vassos, E. Polygenic risk scores: from research tools to clinical  
19 instruments. *Genome medicine* **12**, 1-11 (2020).
- 20 10. Wand, H. *et al.* Improving reporting standards for polygenic scores in risk  
21 prediction studies. *Nature* **591**, 211-219 (2021).
- 22 11. Ma, Y. & Zhou, X. Genetic prediction of complex traits with polygenic scores: a  
23 statistical review. *Trends in Genetics* **37**, 995-1011 (2021).
- 24 12. Yang, S. & Zhou, X. PGS-server: accuracy, robustness and transferability of  
25 polygenic score methods for biobank scale studies. *Briefings in Bioinformatics* **23**,  
26 bbac039 (2022).
- 27 13. Lambert, S.A. *et al.* The Polygenic Score Catalog as an open database for  
28 reproducibility and systematic evaluation. *Nature Genetics* **53**, 420-425 (2021).
- 29 14. Page, M.L. *et al.* The Polygenic Risk Score Knowledge Base offers a centralized  
30 online repository for calculating and contextualizing polygenic risk scores.  
31 *Commun Biol* **5**, 899 (2022).

- 1 15. Kachuri, L. *et al.* Principles and methods for transferring polygenic risk scores  
2 across global populations. *Nature Reviews Genetics*, 1-18 (2023).
- 3 16. Duncan, L. *et al.* Analysis of polygenic risk score usage and performance in  
4 diverse human populations. *Nature communications* **10**, 1-9 (2019).
- 5 17. Hou, K. *et al.* Admix-kit: an integrated toolkit and pipeline for genetic analyses of  
6 admixed populations. *Bioinformatics* **40**, btae148 (2024).
- 7 18. Ding, Y. *et al.* Polygenic scoring accuracy varies across the genetic ancestry  
8 continuum. *Nature* **618**, 774-781 (2023).
- 9 19. Kullo, I.J. *et al.* The PRIMED Consortium: Reducing disparities in polygenic risk  
10 assessment. *The American Journal of Human Genetics* **111**, 2594-2606 (2024).
- 11 20. Linder, J.E. *et al.* Returning integrated genomic risk and clinical  
12 recommendations: The eMERGE study. *Genetics in Medicine* **25**, 100006 (2023).
- 13 21. Lennon, N.J. *et al.* Selection, optimization and validation of ten chronic disease  
14 polygenic risk scores for clinical implementation in diverse US populations.  
15 *Nature Medicine* **30**, 480-487 (2024).
- 16 22. Watanabe, K., Taskesen, E., Bochoven, A. & Posthuma, D. Functional mapping  
17 and annotation of genetic associations with FUMA. *Nature Communications* **8**,  
18 1826 (2017).
- 19 23. Das, S. *et al.* Next-generation genotype imputation service and methods. *Nature*  
20 *genetics* **48**, 1284-1287 (2016).
- 21 24. Li, Y. *et al.* Analyzing bivariate cross-trait genetic architecture in GWAS summary  
22 statistics with the BIGA cloud computing platform. *bioRxiv*, 2023.04. 28.538585  
23 (2023).
- 24 25. Dai, C. *et al.* quantms: a cloud-based pipeline for quantitative proteomics  
25 enables the reanalysis of public proteomics data. *Nature Methods*, 1-5 (2024).
- 26 26. Hayashi, S. *et al.* brainlife. io: a decentralized and open-source cloud platform to  
27 support neuroscience research. *Nature methods* **21**, 809-813 (2024).
- 28 27. Artomov, M., Loboda, A.A., Artyomov, M.N. & Daly, M.J. Public platform with  
29 39,472 exome control samples enables association studies without genotype  
30 sharing. *Nature Genetics* **56**, 327-335 (2024).
- 31 28. Langmead, B. & Nellore, A. Cloud computing for genomic data analysis and  
32 collaboration. *Nature Reviews Genetics* **19**, 208-219 (2018).

- 1 29. Lannelongue, L., Grealey, J. & Inouye, M. Green Algorithms: Quantifying the  
2 Carbon Footprint of Computation. *Adv Sci (Weinh)* **8**, 2100707 (2021).
- 3 30. Chen, T., Zhang, H., Mazumder, R. & Lin, X. Fast and scalable ensemble learning  
4 method for versatile polygenic risk prediction. *Proceedings of the National  
5 Academy of Sciences* **121**, e2403210121 (2024).
- 6 31. Zhao, Z. *et al.* PUMAS: fine-tuning polygenic risk scores with GWAS summary  
7 statistics. *Genome biology* **22**, 1-19 (2021).
- 8 32. Zhang, Q., Privé, F., Vilhjálmsson, B. & Speed, D. Improved genetic prediction of  
9 complex traits from individual-level data or summary statistics. *Nature  
10 communications* **12**, 4192 (2021).
- 11 33. Zhao, Z. *et al.* Optimizing and benchmarking polygenic risk scores with GWAS  
12 summary statistics. *Genome Biol* **25**, 260 (2024).
- 13 34. Genomes Project, C. *et al.* A global reference for human genetic variation. *Nature*  
14 **526**, 68-74 (2015).
- 15 35. Mak, T.S.H., Porsch, R.M., Choi, S.W., Zhou, X. & Sham, P.C. Polygenic scores via  
16 penalized regression on summary statistics. *Genet Epidemiol* **41**, 469-480 (2017).
- 17 36. Privé, F., Arbel, J., Aschard, H. & Vilhjálmsson, B.J. Identifying and correcting for  
18 misspecifications in GWAS summary statistics and polygenic scores. *Human  
19 Genetics and Genomics Advances* **3**(2022).
- 20 37. Privé, F., Arbel, J. & Vilhjálmsson, B.J. LDpred2: better, faster, stronger.  
21 *Bioinformatics* **36**, 5424-5431 (2020).
- 22 38. Vilhjalmsjon, B.J. *et al.* Modeling Linkage Disequilibrium Increases Accuracy of  
23 Polygenic Risk Scores. *Am J Hum Genet* **97**, 576-92 (2015).
- 24 39. Prive, F., Vilhjalmsjon, B.J., Aschard, H. & Blum, M.G.B. Making the Most of  
25 Clumping and Thresholding for Polygenic Scores. *Am J Hum Genet* **105**, 1213-  
26 1221 (2019).
- 27 40. Yang, Y. Adaptive Regression by Mixing. *Journal of the American Statistical  
28 Association* **96**, 574-588 (2001).
- 29 41. Zhang, J. *et al.* An ensemble penalized regression method for multi-ancestry  
30 polygenic risk prediction. *Nature Communications* **15**, 3238 (2024).
- 31 42. Zhang, H. *et al.* A new method for multiancestry polygenic prediction improves  
32 performance across diverse populations. *Nat Genet* **55**, 1757-1768 (2023).

- 1 43. Littlejohns, T.J. *et al.* The UK Biobank imaging enhancement of 100,000  
2 participants: rationale, data collection, management and future directions.  
3 *Nature communications* **11**, 1-12 (2020).
- 4 44. Alfaro-Almagro, F. *et al.* Image processing and Quality Control for the first 10,000  
5 brain imaging datasets from UK Biobank. *NeuroImage* **166**, 400-424 (2018).
- 6 45. Fan, Z. *et al.* The role of sleep in the human brain and body: insights from multi-  
7 organ imaging genetics. *medRxiv*, 2022.09. 08.22279719 (2022).
- 8 46. Bai, W. *et al.* A population-based phenome-wide association study of cardiac and  
9 aortic structure and function. *Nature Medicine* **26**, 1654-1662 (2020).
- 10 47. Zhao, B. *et al.* Eye-brain connections revealed by multimodal retinal and brain  
11 imaging genetics. *Nature Communications* **15**, 6064 (2024).
- 12 48. Yang, X. *et al.* Multi-organ imaging-derived polygenic indexes for brain and body  
13 health. *medRxiv*, 2023.04. 18.23288769 (2023).
- 14 49. Kurki, M.I. *et al.* FinnGen provides genetic insights from a well-phenotyped  
15 isolated population. *Nature* **613**, 508-518 (2023).
- 16 50. Sun, B.B. *et al.* Genetic associations of protein-coding variants in human disease.  
17 *Nature* **603**, 95-102 (2022).
- 18 51. Sun, B.B. *et al.* Plasma proteomic associations with genetics and health in the UK  
19 Biobank. *Nature* **622**, 329-338 (2023).
- 20 52. Deng, Y.-T. *et al.* Atlas of the plasma proteome in health and disease in 53,026  
21 adults. *Cell* **188**, 253-271. e7 (2025).
- 22 53. Xu, Y. *et al.* An atlas of genetic scores to predict multi-omic traits. *Nature* **616**,  
23 123-131 (2023).
- 24 54. Graham, S.E. *et al.* The power of genetic diversity in genome-wide association  
25 studies of lipids. *Nature* **600**, 675-679 (2021).
- 26 55. Jin, J. *et al.* MUSSEL: Enhanced Bayesian Polygenic Risk Prediction Leveraging  
27 Information across Multiple Ancestry Groups. *bioRxiv* (2023).
- 28 56. Fu, J. *et al.* Cross-ancestry genome-wide association studies of brain imaging  
29 phenotypes. *Nature Genetics*, 1-11 (2024).
- 30 57. Sollis, E. *et al.* The NHGRI-EBI GWAS Catalog: knowledgebase and deposition  
31 resource. *Nucleic Acids Res* **51**, D977-D985 (2023).

- 1 58. Ge, T., Chen, C.-Y., Ni, Y., Feng, Y.-C.A. & Smoller, J.W. Polygenic prediction via  
2 Bayesian regression and continuous shrinkage priors. *Nature Communications*  
3 **10**, 1776 (2019).
- 4 59. Yang, S. & Zhou, X. Accurate and scalable construction of polygenic scores in  
5 large biobank data sets. *The American Journal of Human Genetics* **106**, 679-693  
6 (2020).
- 7 60. Fernández-Rhodes, L. *et al.* Ancestral diversity improves discovery and fine-  
8 mapping of genetic loci for anthropometric traits—The Hispanic/Latino  
9 Anthropometry Consortium. *Human Genetics and Genomics Advances* **3**(2022).
- 10 61. Ruan, Y. *et al.* Improving polygenic prediction in ancestrally diverse populations.  
11 *Nat Genet* **54**, 573-580 (2022).
- 12 62. Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human  
13 phenotypes. *Nature genetics* **53**, 1415-1424 (2021).
- 14 63. Verma, A. *et al.* Diversity and scale: Genetic architecture of 2068 traits in the VA  
15 Million Veteran Program. *Science* **385**, eadj1182 (2024).
- 16 64. Zhou, W. *et al.* Global Biobank Meta-analysis Initiative: Powering genetic  
17 discovery across human disease. *Cell Genomics* **2**, 100192 (2022).
- 18 65. Lannelongue, L. *et al.* GREENER principles for environmentally sustainable  
19 computational science. *Nature Computational Science* **3**, 514-521 (2023).
- 20 66. Ruan, Y. *et al.* Leveraging genetic ancestry continuum information to interpolate  
21 PRS for admixed populations. *medRxiv*, 2024.11.09.24316996 (2024).
- 22 67. Ding, Y. *et al.* Large uncertainty in individual polygenic risk score estimation  
23 impacts PRS-based risk stratification. *Nature genetics* **54**, 30-39 (2022).
- 24 68. Abramowitz, S.A. *et al.* Evaluating Performance and Agreement of Coronary  
25 Heart Disease Polygenic Risk Scores. *JAMA* (2024).
- 26 69. Momin, M.M., Lee, S., Wray, N.R. & Lee, S.H. Significance tests for R<sup>2</sup> of out-of-  
27 sample prediction using polygenic scores. *The American Journal of Human*  
28 *Genetics* **110**, 349-358 (2023).
- 29 70. Wang, X. *et al.* Impact of individual level uncertainty of lung cancer polygenic risk  
30 score (PRS) on risk stratification. *Genome medicine* **16**, 22 (2024).
- 31 71. Fu, H., Huang, J., Fan, Z. & Zhao, B. Uncertainty of high-dimensional genetic data  
32 prediction with polygenic risk scores. *arXiv preprint arXiv:2412.20611* (2024).

- 1 72. Wilkinson, M.D. *et al.* The FAIR Guiding Principles for scientific data management  
2 and stewardship. *Scientific data* **3**, 1-9 (2016).
- 3 73. Wang, Y. *et al.* Theoretical and empirical quantification of the accuracy of  
4 polygenic scores in ancestry divergent populations. *Nature communications* **11**,  
5 1-9 (2020).
- 6 74. van der Laan, M.J., Polley, E.C. & Hubbard, A.E. Super learner. *Stat Appl Genet*  
7 *Mol Biol* **6**, Article25 (2007).
- 8 75. Zhao, Z. *et al.* One score to rule them all: regularized ensemble polygenic risk  
9 prediction with GWAS summary statistics. *bioRxiv*, 2024.11. 27.625748 (2024).
- 10 76. International HapMap, C. *et al.* Integrating common and rare genetic variation in  
11 diverse human populations. *Nature* **467**, 52-8 (2010).
- 12 77. Chang, C.C. *et al.* Second-generation PLINK: rising to the challenge of larger and  
13 richer datasets. *Gigascience* **4**, 7 (2015).
- 14 78. Bulik-Sullivan, B.K. *et al.* LD Score regression distinguishes confounding from  
15 polygenicity in genome-wide association studies. *Nat Genet* **47**, 291-5 (2015).
- 16 79. Su, Z., Marchini, J. & Donnelly, P. HAPGEN2: simulation of multiple disease SNPs.  
17 *Bioinformatics* **27**, 2304-5 (2011).
- 18 80. Smith, S.M. *et al.* An expanded set of genome-wide association studies of brain  
19 imaging phenotypes in UK Biobank. *Nature neuroscience* **24**, 737-745 (2021).
- 20 81. Elliott, L.T. *et al.* Genome-wide association studies of brain imaging phenotypes  
21 in UK Biobank. *Nature* **562**, 210-216 (2018).
- 22 82. Beckmann, C.F. & Smith, S.M. Probabilistic independent component analysis for  
23 functional magnetic resonance imaging. *IEEE transactions on medical imaging*  
24 **23**, 137-152 (2004).
- 25 83. Hyvarinen, A. Fast and robust fixed-point algorithms for independent component  
26 analysis. *IEEE transactions on Neural Networks* **10**, 626-634 (1999).
- 27 84. Zhao, B. *et al.* Heart-brain connections: Phenotypic and genetic insights from  
28 magnetic resonance images. *Science* **380**, abn6598 (2023).
- 29 85. Liu, Y. *et al.* Genetic architecture of 11 organ traits derived from abdominal MRI  
30 using deep learning. *Elife* **10**(2021).

- 1 86. Sorokin, E.P. *et al.* Analysis of MRI-derived spleen iron in the UK Biobank  
2 identifies genetic variation linked to iron homeostasis and hemolysis. *Am J Hum*  
3 *Genet* **109**, 1092-1104 (2022).
- 4 87. Wilman, H.R. *et al.* Characterisation of liver fat in the UK Biobank cohort. *PLoS*  
5 *One* **12**, e0172921 (2017).
- 6 88. Mojtahed, A. *et al.* Reference range of liver corrected T1 values in a population  
7 at low risk for fatty liver disease—a UK Biobank sub-study, with an appendix of  
8 interesting cases. *Abdom Radiol (NY)* **44**, 72-84 (2019).
- 9 89. Karlsson, A. *et al.* Automatic and quantitative assessment of regional muscle  
10 volume by multi-atlas segmentation using whole-body water-fat MRI. *J Magn*  
11 *Reson Imaging* **41**, 1558-69 (2015).
- 12 90. Borga, M. *et al.* Validation of a fast method for quantification of intra-abdominal  
13 and subcutaneous adipose tissue for large-scale human studies. *NMR Biomed* **28**,  
14 1747-53 (2015).
- 15 91. West, J. *et al.* Feasibility of MR-Based Body Composition Analysis in Large Scale  
16 Population Studies. *PLoS One* **11**, e0163332 (2016).
- 17 92. Linge, J. *et al.* Body Composition Profiling in the UK Biobank Imaging Study.  
18 *Obesity (Silver Spring)* **26**, 1785-1795 (2018).
- 19 93. Borga, M. *et al.* Reproducibility and repeatability of MRI-based body composition  
20 analysis. *Magn Reson Med* **84**, 3146-3156 (2020).
- 21 94. Langner, T. *et al.* Kidney segmentation in neck-to-knee body MRI of 40,000 UK  
22 Biobank participants. *Sci Rep* **10**, 20963 (2020).
- 23 95. Wu, C., Zhang, Z., Yang, X. & Zhao, B. Large-scale imputation models for multi-  
24 ancestry proteome-wide association analysis. *bioRxiv*, 2023.10. 05.561120  
25 (2023).

26

## 27 **Methods**

### 28 **Single-ancestry PRS pseudo-training**

29 Single-ancestry PRS training aims to develop PRS models for a target genetic  
30 ancestral population based on a single GWAS summary dataset generated from  
31 training samples of the same population. We developed a general summary data-  
32 based parameter optimization approach for multiple single-ancestry PRS methods



1 that avoids the need for individual-level tuning data (**Supplementary Fig. 1a**). We  
2 have implemented the approach to develop pseudo-training versions of three  
3 single-ancestry PRS methods: C+T-pseudo, LDpred2-pseudo, and Lassosum2-  
4 pseudo, which are included in our PennPRS cloud computing platform. We have  
5 also developed the pseudo-training version of an additional method, PRS-CS-grid-  
6 pseudo, which has a much higher computational demand and is included in our  
7 offline pipeline. Our PRS pseudo-training pipeline follows the general framework  
8 of PUMAS<sup>31,33</sup>. Specifically, in Step 1, we use the subsampling approach in  
9 PUMAS to sample marginal association statistics for two “pseudo” subsets of  
10 training and validation individuals from the full GWAS summary data<sup>31</sup>. This  
11 approach enables us to generate GWAS summary statistics for pseudo training  
12 and validation sets for PRS training and parameter tuning, respectively, without  
13 the need to collect an independent individual-level dataset for parameter tuning. In  
14 Step 2, we apply each selected PRS method to train PRS models on the pseudo  
15 summary-level training dataset. In Step 3, we conduct parameter tuning on the  
16 pseudo summary-level validation dataset. This summary data-based parameter  
17 tuning is conducted using the method in PUMAS that allows estimation of the  
18 prediction  $R^2$  of PRS using summary statistics only. This step selects the best  
19 tuning parameter setting for each method based on performance on the pseudo  
20 validation summary dataset. If multiple PRS methods are implemented in Step 2,  
21 we will proceed to Step 4, which offers the option to train ensemble PRS models.  
22 These models combine the PRS models generated by various methods using the  
23 pseudo-validation dataset and two ensemble approaches: Ensemble-pseudo and  
24 Ensemble-ARM-pseudo, which will be introduced in the next section. Finally, in  
25 Step 5, we train the final PRS models on the full GWAS summary dataset with  
26 selected tuning parameter settings obtained from Step 3 and trained ensemble  
27 weights for different methods in the ensemble PRS models from Step 4. To  
28 increase stability of the parameter tuning results, we repeat the training-validation  
29 data splitting procedure in Step 2  $k=2$  times and conduct Steps 2 to 4 with  $k$ -fold  
30 cross-validation. Specifically, for parameter tuning, we select the parameter setting  
31 that correspond to the highest estimated prediction  $R^2$  on the pseudo validation

1 data averaged across the k folds; and for ensemble PRS training, we obtain the  
2 weights in the ensemble model as the average across the k folds.

3

4 We have identified several potential issues of the original PUMAS algorithm when  
5 incorporating it with different PRS methods and have made substantial  
6 modifications accordingly to ensure the applicability and increase the robustness  
7 of our pipeline. For example, the original Lassosum2 and LDpred2 algorithms may  
8 generate non-convergent or problematic PRS weights (e.g., overly large  $|\hat{\beta}_j|$ )  
9 under some tuning parameter settings, which can lead to inflated  $R^2$  estimate for  
10 these settings, resulting in problematic parameter tuning by PUMAS. We resolved  
11 this issue by discarding the tuning parameter settings in which there exist genetic  
12 effect estimates  $|\hat{\beta}_j| > 1$ . Furthermore, it is likely that the final PRS model trained  
13 on the full GWAS summary data has non-convergent variant weights, even though  
14 the selected optimal tuning parameter setting gives a converged model when  
15 trained on the pseudo training dataset. This issue is due to the unstable  
16 performance of the PRS methods, not the PUMAS pseudo-training algorithm itself.  
17 To avoid this inconsistency in the PRS models trained based on the pseudo-  
18 training dataset and the original summary dataset, we select optimal tuning  
19 parameters only from the settings that lead to converged variant weights on the  
20 original summary data. We also noticed that the selected tuning parameter setting  
21 may be far from the optimal setting if its adjacent tuning parameter settings led to  
22 nonconvergent results. We thus only consider parameter settings for which the  
23 adjacent settings also lead to converged results. If no such candidate setting  
24 exists, then we will skip this step and just select the setting that gives the highest  
25  $R^2$  on the pseudo validation set. Finally, for traits that have minimal heritability or  
26 have a small GWAS sample size, the PRS model trained by some of the methods  
27 may have limited power, reflected by negative prediction  $R^2$  estimates on the  
28 pseudo validation data. In this case, we still output the trained PRS models but will  
29 also print a warning message to let users know about this issue. We will also  
30 exclude the corresponding PRS models from the pseudo ensemble learning step.

31

32 **Pseudo ensemble learning combining multiple single-ancestry PRS models**

1 As mentioned in the previous section, if multiple PRS methods are implemented in  
2 single-ancestry analysis, we will provide an option to conduct pseudo-training of  
3 ensemble PRS models combining PRS trained by the various methods based on  
4 the pseudo validation dataset. We propose two approaches for the pseudo  
5 ensemble PRS training. The first approach trains a linear combination<sup>39</sup> of the PRS  
6 models obtained from the various methods (“Ensemble-pseudo”). This approach  
7 was proposed in the PUMA-CUBS framework<sup>33</sup>. We notice that this approach  
8 sometimes generates a PRS that has a lower power than the best single PRS  
9 model, possibly due to the suboptimal performance of some of the single PRS  
10 models. Therefore, we propose an alternative approach adopting a model  
11 combination method named adaptive regression by mixing (ARM)<sup>40</sup>, which, under  
12 certain conditions, can approximate the optimal performance among a set of single  
13 models (“Ensemble-ARM-pseudo”). We observe from our simulation studies and  
14 data applications that either one of the two approaches outperforms the other on  
15 different phenotypes and with different training GWAS datasets. We thus provide  
16 both ensemble PRS models to users to further increase the power of the “best”  
17 PRS model provided by PennPRS.

18

### 19 **Multi-ancestry PRS pseudo-training**

20 Multi-ancestry PRS training jointly analyzes ancestry-stratified GWAS summary  
21 statistics from  $K$  ancestral populations (a subset of [EUR, AFR, AMR, EAS, and  
22 SAS]) and generates ancestry-specific PRS models for the  $K$  populations. We  
23 developed a GWAS summary data-based parameter tuning approach for multi-  
24 ancestry PRS training that avoids the need for individual-level tuning data  
25 (**Supplementary Fig. 1b**). We have implemented this approach to develop the  
26 pseudo-training version of PROSPER on our PennPRS cloud computing platform  
27 and two other methods, PRS-CSx-pseudo and MUSSEL-pseudo, which require  
28 much larger memory and/or are computationally more intensive and are thus only  
29 included in our offline pipeline.

30

31 Our general multi-ancestry PRS pseudo-training framework follows that of  
32 PUMAS<sup>31,33</sup>. Specifically, in Step 1, we use the subsampling approach in PUMAS

1 to generate summary statistics for pseudo training and validation sets for PRS  
2 training and parameter tuning, respectively, for each of the  $K$  ancestry populations.  
3 In Step 2, we apply each selected method to train PRS models on the pseudo  
4 training dataset. For PROSPER-pseudo and MUSSEL-pseudo which prerequire  
5 implementation of the single-ancestry Lassosum2 and LDpred2 algorithms,  
6 respectively, we use a procedure similar to the one in single-ancestry pseudo-  
7 training for selecting optimal parameters of Lassosum2-pseudo and LDpred2-  
8 pseudo. In Step 3, we conduct parameter optimization of the multi-ancestry joint  
9 modeling step in PROSPER or MUSSEL on the pseudo summary-level validation  
10 dataset. This step selects a best PRS model for each ancestry based on its  
11 performance (the estimated prediction  $R^2$ ) on the pseudo validation dataset. All  
12 three methods have a final ensemble learning step (Step 4, **Supplementary Fig.**  
13 **1**), where PRS-CSx trains a linear combination of the  $K$  best ancestry-specific PRS  
14 models from Step 3, while PROSPER and MUSSEL use the super learner  
15 algorithm<sup>74</sup> with base learners including linear regression, elastic-net regression,  
16 and ridge regression to train an “optimal” linear combination of all PRS models  
17 across all tuning parameter settings and ancestries. For PRS-CSx, we use the  
18 ensemble approach in PUMAS to train the final PRS model for each ancestry  
19 based on the pseudo validation dataset of that ancestry. For PROSPER and  
20 MUSSEL, a summary data version of the super learner can be implemented for  
21 the final ensemble step utilizing a recently developed approach<sup>75</sup>. For now, we  
22 consider an alternative strategy, where we train a linear combination of a subset  
23 of  $L$  best performing single PRS models without regularization. Specifically, for  
24 each ancestry, we first select the top  $L$  of the  $\sum_{k=1}^K M_k$  PRS models that have the  
25 highest prediction  $R^2$  on the pseudo validation dataset of that ancestry, where  $M_k$   
26 denotes the number of tuning parameter settings, i.e., number of candidate PRS  
27 models, generated for ancestry  $k$ ,  $k = 1, 2, \dots, K$ . We then train a linear combination  
28 of these  $L$  top PRS models on the pseudo validation dataset. We set  $L$  to the  
29 minimum between five and the number of converged PRS models among the  
30  $\sum_{k=1}^K M_k$  models. Finally, in Step 5, we train the final PRS models on the full original  
31 GWAS summary data based on the selected optimal parameter settings from Step  
32 3 and trained ensemble weights for different single PRS models from Step 4. We

1 notice that the single best PRS model may have a higher prediction power than  
2 the final ensemble PRS. We thus provide both the best single PRS model  
3 (“PROSPER-Single-pseudo PRS”, **Figs. 1c and 5, Supplementary Figs. 4-5**) and  
4 the final PRS model with the ensemble step (“PROSPER-pseudo PRS”, **Figs. 1c**  
5 **and 5, Supplementary Figs. 4-5**) to the user. Again, we repeat the training-  
6 validation data splitting procedure in Step 2  $k=2$  times and conduct Steps 2-4 with  
7  $k$ -fold cross-validation to increase stability of the results. Specifically, for parameter  
8 tuning, we select the parameter setting that correspond to the highest estimated  
9 prediction  $R^2$  on the pseudo validation data averaged across the  $k$  folds; and for  
10 the ensemble step, we obtain the weights in the ensemble model as the average  
11 across the  $k$  folds. In our multi-ancestry analysis pipeline, we also consider the  
12 various modifications to the original PUMAS algorithm implemented in our single-  
13 ancestry analysis pipeline to improve robustness of the pseudo-training.

14

### 15 **Configuration of the online PennPRS pipeline**

16 Our PennPRS cloud computing pipeline currently supports seven PRS methods,  
17 including pseudo-training versions of three single-ancestry methods, C+T-pseudo,  
18 lassosum2-pseudo, and LDpred2-pseudo; three tuning-parameter-free single-  
19 ancestry methods, PRS-CS-auto, LDpred2-auto, and DBSLMM; and pseudo-  
20 training version of one multi-ancestry method, PROSPER-pseudo. PennPRS  
21 supports PRS model development based on genetic variants in the HapMap 3<sup>76</sup>.  
22 For implementation of methods that have tuning parameters, we set up default  
23 tuning parameter settings based on the ones in the original algorithms of the  
24 methods but with slight modifications to balance between prediction performance  
25 and computational efficiency of our online PRS training. We use genotype data of  
26 unrelated individuals from the Phase 3 1000 Genomes project as the linkage  
27 disequilibrium (LD) reference data. We now introduce the parameter settings and  
28 other relevant details of the newly developed pseudo-training methods supported  
29 by PennPRS.

30

31 C+T-pseudo. C+T-pseudo first conducts an LD clumping step to select relatively  
32 independent genetic variants with an absolute pairwise correlation lower than  $r^2 =$

1 0.1 within a genetic distance 500kb calculated based on the reference genotype  
2 dataset for the same ancestral population from the Phase 3 1000 Genomes  
3 Project<sup>34</sup>. It then selects the remaining genetic variants that reach varying p-value  
4 cutoffs (tuning parameter:  $pt$ )<sup>4,5</sup>. Our default candidate values for  $pt$  are  $5 \times 10^{-8}$ ,  
5  $5 \times 10^{-7}$ ,  $5 \times 10^{-6}$ ,  $5 \times 10^{-5}$ ,  $5 \times 10^{-4}$ ,  $5 \times 10^{-3}$ ,  $5 \times 10^{-2}$ , and  $5 \times 10^{-1}$ . C+T-pseudo then  
6 selects the score with the “optimal” p-value threshold based on the performance  
7 on the pseudo validation dataset with respect to the prediction  $R^2$ . PennPRS runs  
8 C+T-pseudo using PLINK 1.90<sup>77</sup>.

9

10 Lassosum2-pseudo. Lassosum2-pseudo is a penalized regression-based  
11 approach that estimates joint genetic effect sizes based on GWAS summary  
12 statistics. Tuning parameters include (i)  $\lambda$ : shrinkage parameter in the  $L_2$   
13 regularization (default candidate values: 0.001, 0.01, 0.1, and 1); (ii) number of  
14 candidate values for  $\lambda$ , the shrinkage parameter in the  $L_1$  regularization (default:  
15 30); and (iii) ratio between the lowest and highest candidate values of  $\lambda$  (default:  
16 0.01). The current version of PennPRS implements Lassosum2-pseudo with R  
17 package “bigsnpr” (version 1.6.1, last updated Jun 8, 2023).

18

19 LDpred2-pseudo. LDpred2-pseudo is a Bayesian approach that jointly analyzes  
20 correlated genetic variants across the genome and accounts for LD<sup>37,38</sup>. It uses a  
21 spike-and-slab prior on genetic effect sizes, assuming a proportion ( $p$ ) of the  
22 genetic variants have non-zero effect on the phenotype. Tuning parameters  
23 include (1) the causal variant proportion  $p$  (default candidate values:  $10^{-5}$ ,  $3.2 \times 10^{-5}$ ,  
24  $10^{-4}$ ,  $3.2 \times 10^{-4}$ ,  $10^{-3}$ ,  $3.2 \times 10^{-3}$ ,  $10^{-2}$ ,  $3.2 \times 10^{-2}$ ,  $10^{-1}$ ,  $3.2 \times 10^{-1}$ , and 1), (2)  
25 heritability-related parameter,  $\alpha$ : the total heritability is set to  $H_2 = \alpha H_2^0$ , where  $H_2^0$   
26 is the heritability estimated by LD score regression<sup>78</sup> (default candidate values:  
27  $\alpha = 0.7$ , 1.0, and 1.4), and an additional sparse option (default: FALSE) to shrink  
28 the posterior genetic effects that exceed  $p$  to zero. The current version of  
29 PennPRS implements LDpred2-pseudo with R package “bigsnpr” (version 1.6.1,  
30 last updated Jun 8, 2023).

31

1 PROSPER-pseudo. PROSPER-pseudo is a penalized regression-based multi-  
2 ancestry PRS method that utilizes an  $L_1$  penalty to induce sparsity of genetic  
3 variants with non-zero effects and an  $L_2$  penalty to induce correlation in genetic  
4 effects between ancestries<sup>41</sup>. Tuning parameters include (i) the number of  
5 candidate values of the shrinkage parameter in the  $L_1$  penalty (default: 5) and (ii)  
6 the number of candidate values of the shrinkage parameter in the  $L_2$  penalty  
7 (default: 5).

8

### 9 **Simulation studies**

10 We evaluated the performance of our proposed pseudo-training approach for both  
11 single- and multi-ancestry PRS development in comparison to the traditional,  
12 individual-level tuning data-based PRS training in various data settings based on  
13 a large-scale synthetic GWAS data previously generated<sup>42</sup>. The synthetic  
14 genotype data were generated for all five super populations (EUR, AFR, AMR,  
15 EAS, and SAS) using HAPGEN2 (version 2.1.2)<sup>79</sup> to closely mimic the reference  
16 genotype data from the Phase 3 1000 Genomes Project. Phenotype data were  
17 generated assuming a causal variant proportion of 1%, 0.1%, of 0.05% and GWAS  
18 sample size of 15,000, 45,000, or 80,000 across the five ancestral populations.  
19 Details of the simulation procedure were previously described<sup>42</sup>.

20

### 21 **Real data analyses for evaluation of pseudo PRS training**

22 UKB imaging data analysis. We conducted a large-scale evaluation of single-  
23 ancestry pseudo-training methods using multi-organ multi-modality imaging  
24 data<sup>48,80,81</sup> from the UK Biobank (UKB) study, covering brain, heart, eye, and  
25 abdominal organs (**Supplementary Tables 2-5**). For the brain, we used imaging-  
26 derived phenotypes from three major modalities: structural MRI (sMRI), diffusion  
27 MRI (dMRI), and resting-state functional MRI (rfMRI). For example, brain sMRI  
28 included 1,432 phenotypes generated from the FIRST, FAST, and FreeSurfer  
29 pipelines<sup>44</sup>. Brain dMRI data included 674 phenotypes processed with TBSS and  
30 ProbtrackX, while brain rfMRI encompassed 82 phenotypes from whole-brain  
31 spatial independent component analysis<sup>44,81-83</sup>, covering regional amplitude and  
32 global functional connectivity. For cardiac MRI, we used 82 phenotypes related to

1 the heart and aorta<sup>46,84</sup>. Additionally, 41 abdominal MRI phenotypes<sup>85-94</sup> were  
2 included, covering kidney, liver, and abdominal organ or tissues. In addition, we  
3 analyzed 46 phenotypes derived from eye optical coherence tomography  
4 images<sup>47</sup>. The GWAS summary data for these imaging phenotypes were obtained  
5 from subjects of self-reported British European ancestry, with average sample  
6 sizes of 32,634 for brain, 30,506 for heart, 29,849 for abdomen, and 50,465 for  
7 eye. Age, sex, the top 40 genetic principal components (PCs), and imaging-  
8 specific covariates were adjusted for, as detailed in a previous study<sup>45,48</sup>. PRS  
9 performance was assessed on 2,227 to 8,172 European non-British subjects,  
10 using the same set of covariates as those in the corresponding GWAS.

11

12 FinnGen disease data analysis. We evaluated the performance of single-ancestry  
13 pseudo-training methods on binary phenotypes using GWAS summary statistics  
14 from the FinnGen study (R9)<sup>49</sup>. Following the FinnGen-phencode mapping  
15 approach used in previous studies<sup>48,50</sup>, we mapped 29 disease pairs, with an  
16 average of 333,355 cases and controls per phenotype (**Supplementary Tables 6-**  
17 **7**). PRS performance was assessed on 1,225 to 155,170 European cases from the  
18 UKB, with adjustments for effects of age and sex.

19

20 UKB-PPP Olink plasma protein data analysis. We evaluated the performance of  
21 single-ancestry pseudo PRS training methods on 2,734 Olink plasma proteins from  
22 the UKB-PPP<sup>51</sup> project (**Supplementary Tables 8-9**). The GWAS summary  
23 statistics were obtained from a previous study<sup>95</sup>, which included 40,852 subjects  
24 of British European ancestry with adjustments for age, sex, and the top 40 genetic  
25 PCs. PRS performance was assessed on 2,517 to 2,923 European non-British  
26 subjects, using the same set of covariates as in the GWAS.

27

28 GLGC blood lipids data analysis. We trained ancestry-specific PRS models by both  
29 pseudo-training and individual-level tuning versions of the various methods for four  
30 blood lipids, including high-density lipoprotein (HDL), low-density lipoprotein (LDL),  
31 log-transformed triglycerides (logTG), and total cholesterol (TC)<sup>54</sup>. The ancestry-  
32 stratified training GWAS summary data were obtained from the Global Lipids



1 Genetics Consortium<sup>54</sup> (GLGC) on five ancestry groups including EUR ( $N =$   
2 840,018-927,975), AFR or admixed AFR ( $N = 87,759-92,554$ ), Hispanic/Latino  
3 ( $N = 33,989-48,056$ ), EAS ( $N = 80,676-145,512$ ), and SAS ( $N = 33,658-34,135$ )<sup>54</sup>.  
4 We validated method performance on a random set of 20,000 UKB individuals of  
5 EUR ancestry and all UKB individuals of AFR ( $N = 9,169$ ), AMR ( $N = 750$ ), EAS  
6 ( $N = 2,019$ ), and SAS ( $N = 10,853$ ) ancestry. We inferred the ancestry of the UKB  
7 individuals by a genetic component analysis<sup>41</sup>. We used 50% of these UKB  
8 samples to conduct individual-level data-based parameter tuning, ensemble PRS  
9 training, and conducting the ensemble step in PROSPER, and used the remaining  
10 50% (testing set) to evaluate PRS performance of the various methods. GWAS  
11 sample sizes, validation sample sizes, and the number of genetic variants  
12 analyzed are reported in **Supplementary Table 11**. Detailed data quality control  
13 procedures were previously described<sup>42,55</sup>. Age, gender, and the top 10 genetic  
14 PCs were adjusted for as covariates when calculating prediction  $R^2$  of the PRS  
15 models.

16  
17 UKB/CHIMGEN brain imaging data analysis. We evaluated the performance of  
18 multi-ancestry pseudo-training methods on brain imaging phenotypes using  
19 GWAS summary statistics from both the Chinese Imaging Genetics (CHIMGEN)  
20 study<sup>56</sup> for East Asians (average  $N = 7,058$ ) and the UKB study for British  
21 European ancestry (average  $N = 34,286$ ). Similar to our single-ancestry analysis  
22 on UKB, we included 968 phenotypes from sMRI and 445 from dMRI. PRS  
23 performance was assessed on 443 Asian subjects from the UKB study (half were  
24 used as tuning samples for parameter optimization for traditional PRS training  
25 methods, half were used as testing data to evaluate performance of all methods)  
26 with adjustments for the same set of covariates as in the single-ancestry analysis  
27 on the same phenotypes. GWAS sample sizes, validation sample sizes, and the  
28 number of genetic variants analyzed are reported in **Supplementary Table 13**.

29

### 30 **Other GWAS datasets on which we have generated PRS models**

31 We applied our offline pipeline to GWAS summary statistics from the Biobank  
32 Japan (BBJ)<sup>62</sup>, the Million Veteran Program (MVP) study<sup>63</sup>, the Global Biobank

1 Meta-analysis Initiative (GBMI) consortium<sup>64</sup>, and the GWAS Catalog<sup>57</sup>. For the  
2 BBJ, we conducted single-ancestry analysis on GWAS summary statistics for 169  
3 phenotypes available at <https://pheweb.jp/downloads>. For the GWAS Catalog, we  
4 analyzed nearly 8,000 harmonized datasets in single-ancestry analysis. For the  
5 GBMI, we performed multi-ancestry analysis on nine phenotypes with ancestry-  
6 stratified GWAS summary statistics from the five super populations<sup>34</sup>. For the MVP  
7 study, we carried out multi-ancestry analysis on 181 phenotypes with ancestry-  
8 stratified GWAS summary statistics from AFR, AMR, EAS, and EUR populations.  
9 Using default parameter settings, our pipelines were successfully applied to these  
10 data resources, and the generated PRS models have been shared in the PennPRS  
11 public resource hub.

12

### 13 **Cloud computing platform development**

14 PennPRS is a cloud-based platform hosted on AWS that consists of two main  
15 components: the frontend and the backend. For the frontend, we used Next.js  
16 (<https://nextjs.org/>) and MUI (<https://mui.com/>) to create a clean, intuitive interface.  
17 Users can easily input their data (through file uploads or data queries), choose the  
18 type of analysis, PRS methods, and parameter setting they want, and view job  
19 status and results. We used Next.js to ensure that the platform loads quickly and  
20 that all interactions (such as submitting data or viewing outputs) are smooth and  
21 responsive. The backend provides the infrastructure for PRS model development,  
22 including data harmonization and QC pipelines, GWAS Catalog data querying, and  
23 various PRS methods and training mechanisms. We developed the backend with  
24 FastAPI (<https://fastapi.tiangolo.com/>), which allows us to process multiple tasks  
25 efficiently, supporting both simple requests and more complex data processing.  
26 For example, when a user uploads data for PRS analysis, FastAPI sends this data  
27 to the job queue, ensuring that requests are processed in a fair and timely manner.

28

29 In addition, Redis (<https://redis.io/>) is used for job management and queue,  
30 keeping track of all incoming requests and organizing them so that the system can  
31 handle multiple tasks simultaneously. Redis also helps prevent delays and keeps  
32 the platform running smoothly even during busy times. Moreover, since different

1 types of analyses have different resource requirements, we organized the  
2 computing infrastructure into distinct subgroups to optimize resources. Each  
3 subgroup is tailored to handle specific types of jobs, ensuring that the right  
4 resources (such as memory and CPU power) are available for the task at hand,  
5 which optimizes resource allocation and improves overall efficiency. Once the  
6 analysis is completed, the results are sent back to the frontend so users can  
7 access and download them. To ensure reliability and scalability, the platform  
8 incorporates monitoring tools for system performance checks, automated testing,  
9 and continuous integration pipelines. This setup enables quick future updates and  
10 secure data handling, ensuring a smooth user experience as demand grows.

11

### 12 **Code availability**

13 The developed PRS pseudo-training methods and PennPRS pipelines can be  
14 freely accessed at <https://pennprs.org/> and <https://github.com/PennPRS/pipeline>.

15

### 16 **Data availability**

17 The simulated genotype and phenotype data used in our simulations are available  
18 at: <https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/C>  
19 [OXHAP](https://dataverse.harvard.edu/dataset.xhtml?persistentId=doi:10.7910/DVN/C). GWAS summary statistics used in our PRS training and evaluation can  
20 be obtained from their respective data sources, subject to data sharing policies  
21 and approvals. Specifically, the harmonized GWAS summary statistics from the  
22 GWAS Catalog are available at [https://www.ebi.ac.uk/gwas/downloads/summary-](https://www.ebi.ac.uk/gwas/downloads/summary-statistics)  
23 [statistics](https://www.ebi.ac.uk/gwas/downloads/summary-statistics). The EUR GWAS summary statistics for the UKB imaging phenotypes  
24 across different organs are available from previous study<sup>45,48</sup>. The EUR protein  
25 GWAS summary statistics from the UKB-PPP project are available from previous  
26 study<sup>95</sup>. The EUR GWAS summary statistics from the FinnGen study are available  
27 at [https://www.finnngen.fi/en/access\\_results](https://www.finnngen.fi/en/access_results). The EAS GWAS summary statistics  
28 from BBJ are available at <https://pheweb.jp/>. The EAS GWAS summary statistics  
29 for brain imaging phenotypes from the CHIMGEN study are available from  
30 previous study<sup>56</sup>. Ancestry-stratified GWAS summary statistics from the GBMI are  
31 available at <https://www.globalbiobankmeta.org/resources>. Ancestry-stratified  
32 GWAS summary statistics for blood lipids across five super populations from

1 GLGC are available at [http://csg.sph.umich.edu/willer/public/glgc-](http://csg.sph.umich.edu/willer/public/glgc-lipids2021/results/ancestry_specific)  
2 [lipids2021/results/ancestry\\_specific](http://csg.sph.umich.edu/willer/public/glgc-lipids2021/results/ancestry_specific). Ancestry-stratified GWAS summary statistics  
3 from the MVP study are available from previous study<sup>63</sup>. The individual-level UK  
4 Biobank data used in this study can be requested from  
5 <https://www.ukbiobank.ac.uk/>. The PRS model weights generated by the  
6 PennPRS pipeline have been made publicly available through the PennPRS public  
7 resource hub at <https://pennprs.org/result>.

1 **Figure legends**

2

3 **Fig. 1: The Challenges of Traditional PRS Model Training and the Promise of**  
4 **PennPRS Cloud Computing Platform.**

5 **a.** Left: A figurative representation of the key challenges in performing PRS model  
6 training with local computing servers and pipelines. Right: Our proposed cloud  
7 computing approach for online PRS model training, which leverages centralized  
8 computing and data resources alongside novel pseudo-training algorithms and  
9 pipelines to overcome these challenges. **b.** An overview of the cloud computing  
10 platform of PennPRS and its major impacts on PRS applications in precision  
11 medicine.

12

13 **Fig. 2: Development and Distribution of the PennPRS Cloud Computing**  
14 **Platform and Accompanying Data and Computational Resources.**

15 **a.** A summary of the main contributions of our study, including the distribution and  
16 large-scale validation of PRS pseudo-training pipelines, establishment of the  
17 PennPRS cloud computing platform, and distribution of queryable GWAS  
18 summary data sources, pretrained PRS models, and offline pipeline. **b.** Workflow  
19 of the single-ancestry PRS training supported by PennPRS. **c.** Workflow of the  
20 multi-ancestry PRS training supported by PennPRS. **d.** Highlighted features of  
21 PennPRS: (i) new PRS pseudo-training pipelines supporting three single-ancestry  
22 methods, two ensemble approaches combining different single-ancestry methods,  
23 and one multi-ancestry method; and (ii) large-scale application and validation of  
24 the PRS pseudo-training pipeline across nine data resources and over 6,000  
25 phenotypes.

26

27 **Fig. 3: Comparison of Single-ancestry PRS Pseudo-training and Traditional**  
28 **PRS methods with Individual-level Tuning Data of Various Sample Sizes**  
29 **under various settings of causal SNP proportion and heritability.**

30 We compared the prediction  $R^2$  of the PRS models trained by C+T-pseudo,  
31 Lassosum2-pseudo, LDpred2-pseudo, Ensemble-pseudo, and Ensemble-ARM-  
32 pseudo ( $R^2_{sum}$ ) with those of PRS models trained based on individual-level tuning

1 dataset ( $R^2_{ind}$ ) that has a sample size **a.**  $N_{tuning}=1,000$ , **b.**  $N_{tuning}=400$ , or **c.**  
2  $N_{tuning}=100$ . Results were summarized across 10 training GWAS summary  
3 datasets of  $N_{GWAS}=15,000$  and averaged across 100 random splits, with each split  
4 having  $N_{tuning}$  tuning samples for individual data-based parameter tuning and  
5  $N_{val}=2,500$  validation samples for calculating prediction  $R^2$  for all models. Detailed  
6 results are reported in Supplementary Table 1.

7

8 **Fig. 4: Evaluation of Single-ancestry PRS Pseudo-training on Body Imaging**  
9 **Phenotypes Using GWAS Summary Data and Validation Data from the UK**  
10 **Biobank (UKB) study.**

11 **a.** We compared our PRS pseudo training approaches, C+T-pseudo, Lassosum2-  
12 pseudo, LDpred2-pseudo, Ensemble-pseudo, and Ensemble-ARM-pseudo ( $R^2_{sum}$ )  
13 with the original methods that use individual-level tuning dataset ( $R^2_{ind}$ ) on 41  
14 abdominal MRI (average  $N_{GWAS}=29,849$ ), 82 cardiac MRI (average  $N_{GWAS}=30,506$ ),  
15 and 46 eye OCT (average  $N_{GWAS}=50,465$ ) phenotypes and evaluated their  
16 performance on hold-out independent UKB samples of EUR origin ( $N_{val}=5,760$ ). **b.**  
17 We assessed the relative performance of the pseudo-training methods to their  
18 original versions utilizing individual-level tuning datasets of different sizes  $N_{tuning}=$   
19 1,000, 300, or 100, on the abdominal MRI, cardiac MRI, and eye OCT phenotypes.  
20 Results were averaged across 100 random splits, with each split having  $N_{tuning}$   
21 tuning samples for individual data-based parameter tuning and the remaining  
22 samples for calculating prediction  $R^2$  for all models. Detailed data information and  
23 results are summarized in Supplementary Tables 2-3.

24

25 **Fig. 5: Additional Evaluation of Single-ancestry PRS Pseudo-training across**  
26 **Various Phenotypes and Data Sources.**

27 We compared our PRS pseudo-training approaches, C+T-pseudo, Lassosum2-  
28 pseudo, LDpred2-pseudo, Ensemble-pseudo, and Ensemble-ARM-pseudo ( $R^2_{sum}$ )  
29 with the original methods that use individual-level tuning dataset ( $R^2_{ind}$ ) on **a.** 2,363  
30 brain multi-modal imaging phenotypes based on GWAS summary statistics of EUR  
31 ancestry from the UK Biobank (UKB) study ( $N_{GWAS}=32,620$ ) and evaluated their  
32 performance on hold-out independent UKB samples of EUR ancestry ( $N_{val}=5,020$ );

1 **b.** 29 binary disease phenotypes based on GWAS summary statistics of EUR  
2 ancestry from the FinnGen study (a total of 333,355 cases and controls on average)  
3 and evaluated their performance on UKB samples of EUR ancestry (23,048 cases  
4 on average); and **c.** 2,734 Olink plasma proteins based on GWAS summary  
5 statistics of EUR ancestry from the UKB-PPP project ( $N_{GWAS}=40,852$ ) and  
6 evaluated their performance on hold-out independent UKB samples of EUR  
7 ancestry ( $N_{val}=2,731$ ). We used half of the UKB validation samples for individual-  
8 level parameter tuning and the remaining half to report AUC (for binary disease  
9 phenotypes) and prediction  $R^2$  (for continuous phenotypes) for both our pseudo-  
10 training approach and the individual-level tuning data-based training approach.  
11 Here rfMRI stands for resting-state functional MRI, dMRI stands for diffusion MRI,  
12 and sMRI stands for structural MRI. Detailed data information and results are  
13 summarized in Supplementary Tables 4-9.

14

15 **Fig. 6: Evaluation of Multi-ancestry PRS Pseudo-training by Simulation**  
16 **Studies and Applications on Various Phenotypes and Data Sources.**

17 **a.** Results show the comparison of the PRS trained by pseudo-training methods  
18 ( $R^2_{sum}$ , PROSPER-Single-pseudo and PROSPER-pseudo) with PROSPER-Single  
19 PRS and PROSPER PRS trained with individual-level tuning datasets ( $R^2_{ind}$ ) on **a.**  
20 simulated datasets under different settings of heritability, negative selection  
21 patterns, and causal genetic variant proportions assuming a 100,000 GWAS  
22 sample size for EUR and varying GWAS sample sizes for each non-EUR  
23 population (15,000, 45,000, or 80,000) with 2,500 tuning samples for individual  
24 data-based parameter tuning and 2,500 validation samples for calculating  
25 prediction  $R^2$  for the various models; **b.** Four blood lipid phenotypes based on  
26 GWAS summary statistics of EUR, AFR, AMR, EAS, and SAS ancestries from the  
27 GLGC study ( $N_{GWAS}=33,658-930,671$ ) and evaluated their performance on  
28 independent UK Biobank (UKB) samples ( $N_{val}=1,752-19,030$ ); and 382 brain  
29 diffusion MRI (dMRI) phenotypes based on GWAS summary data of EUR ancestry  
30 from the UKB study ( $N_{GWAS}=28,626-32,744$ ) and GWAS summary data of EAS  
31 ancestry from the CHIMGEN study ( $N_{GWAS}=7,058$ ) and evaluated their  
32 performance on hold-out independent UKB samples ( $N_{val}=4,955$  for EUR and

- 1  $N_{val}=413-444$  for EAS, half for parameter tuning for the original PROSPER method
- 2 and the remaining half for calculating  $R^2$  for all models). **c.** Comparison of the
- 3 performance of multi-ancestry method, PROSPER-pseudo, and its single-ancestry
- 4 analogue, Lassosum2-pseudo, on the 382 dMRI phenotypes. Detailed data
- 5 information and results are summarized in Supplementary Tables 11-14.



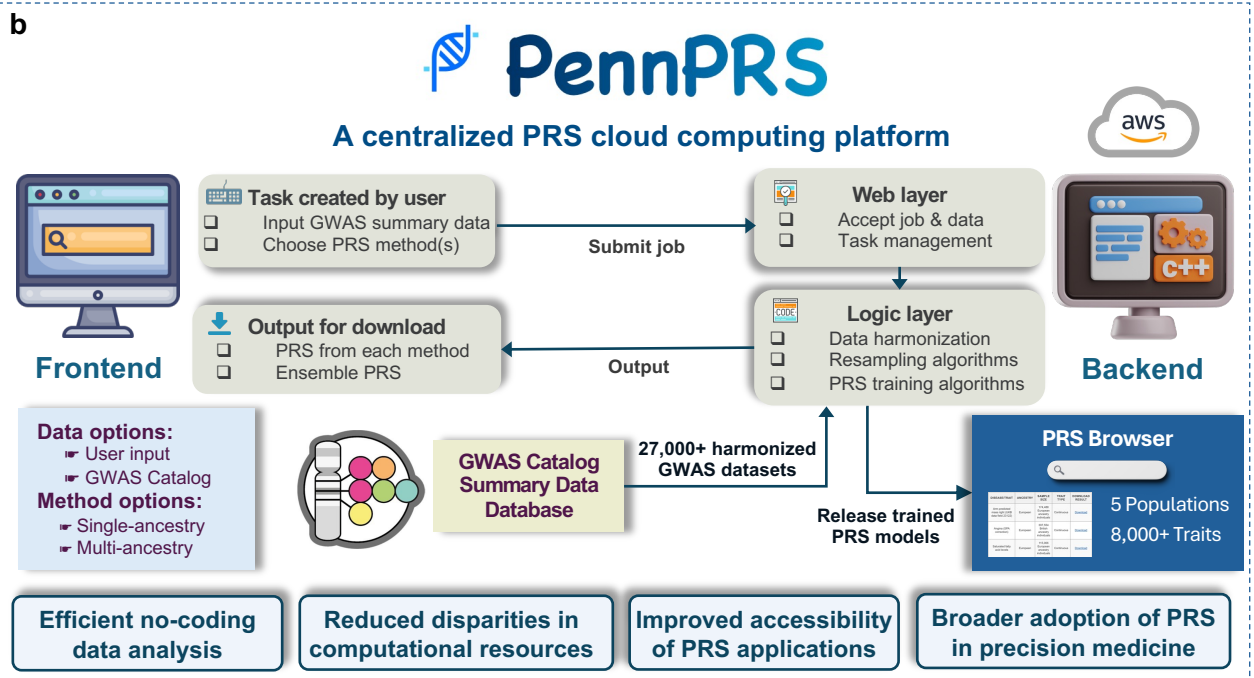
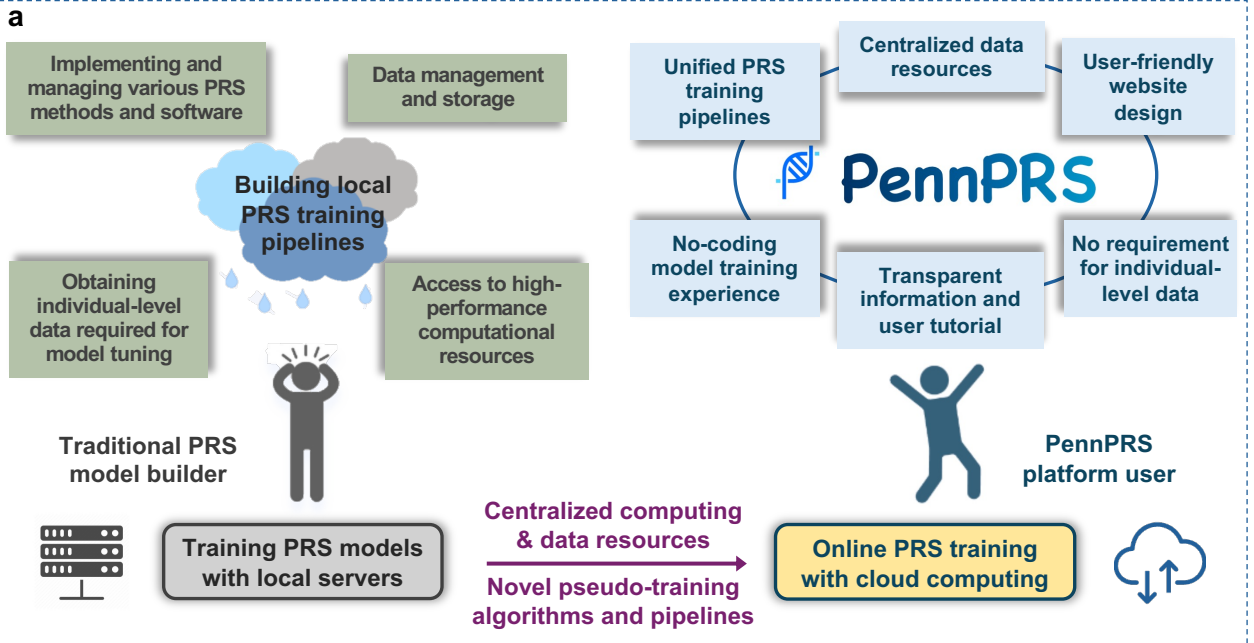


Fig. 1

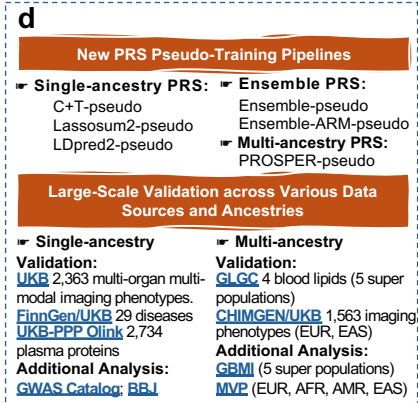
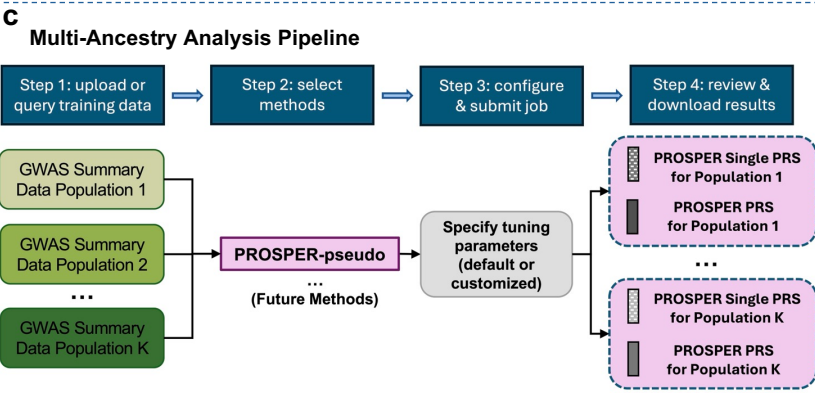
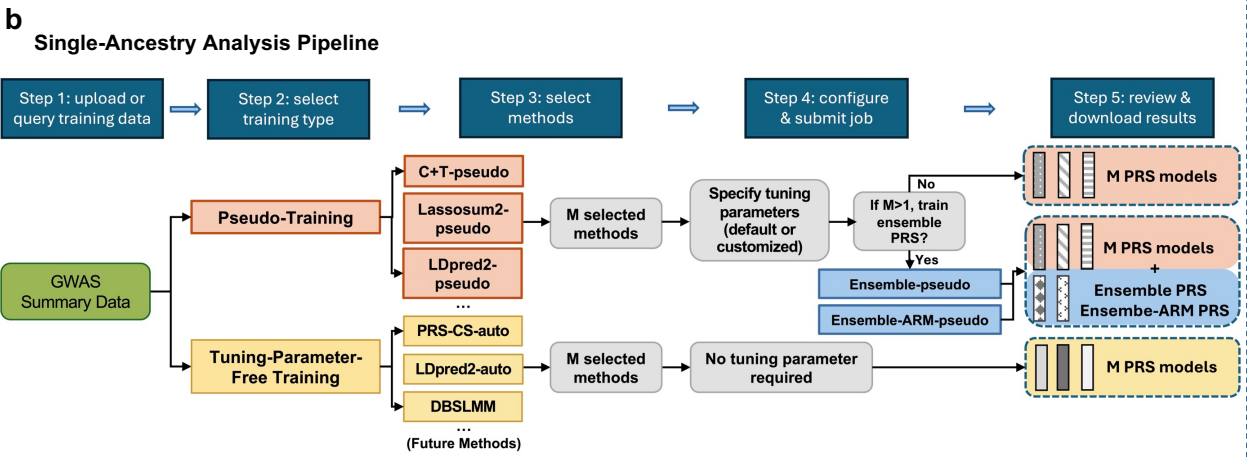
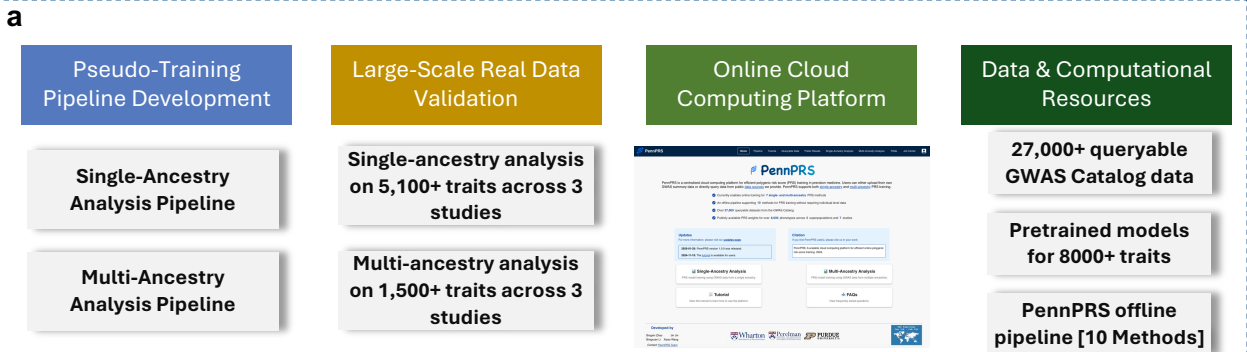
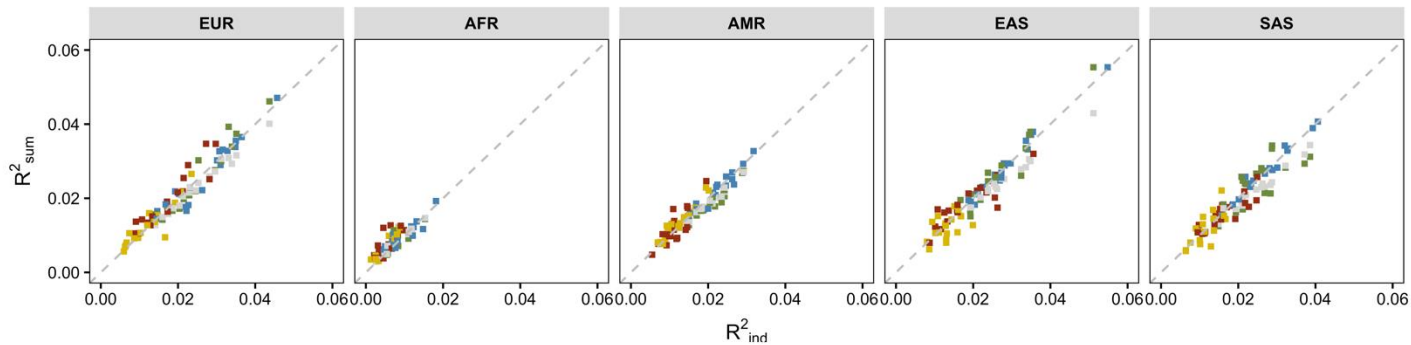
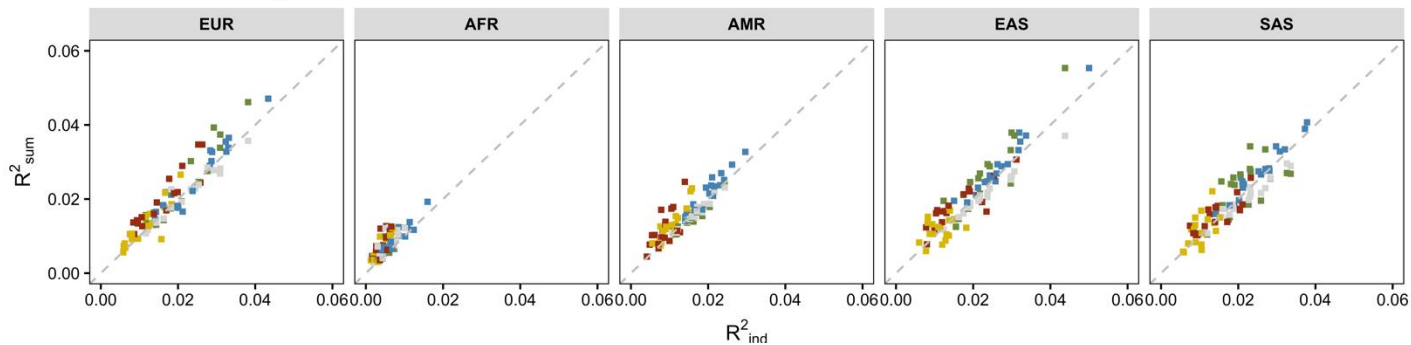
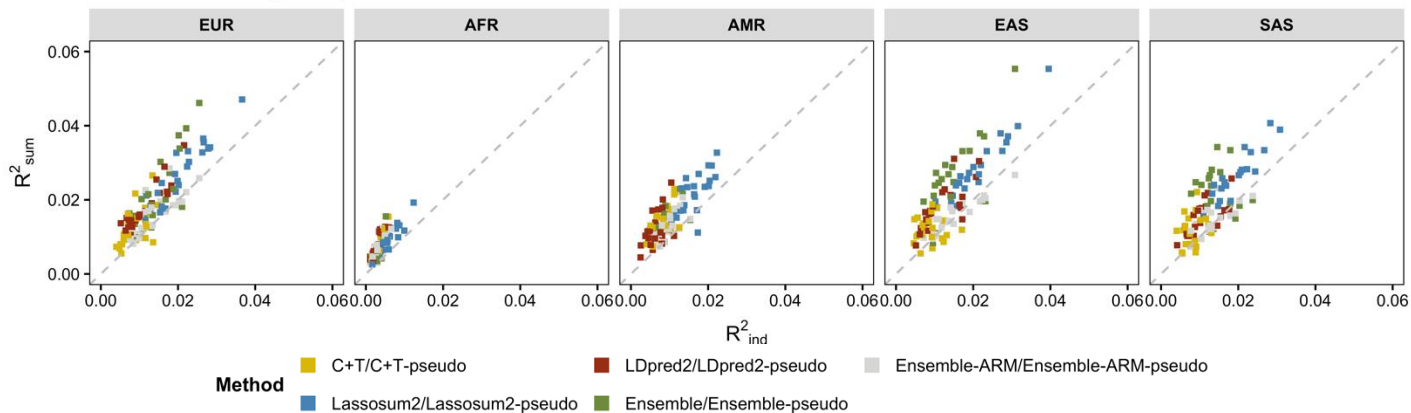
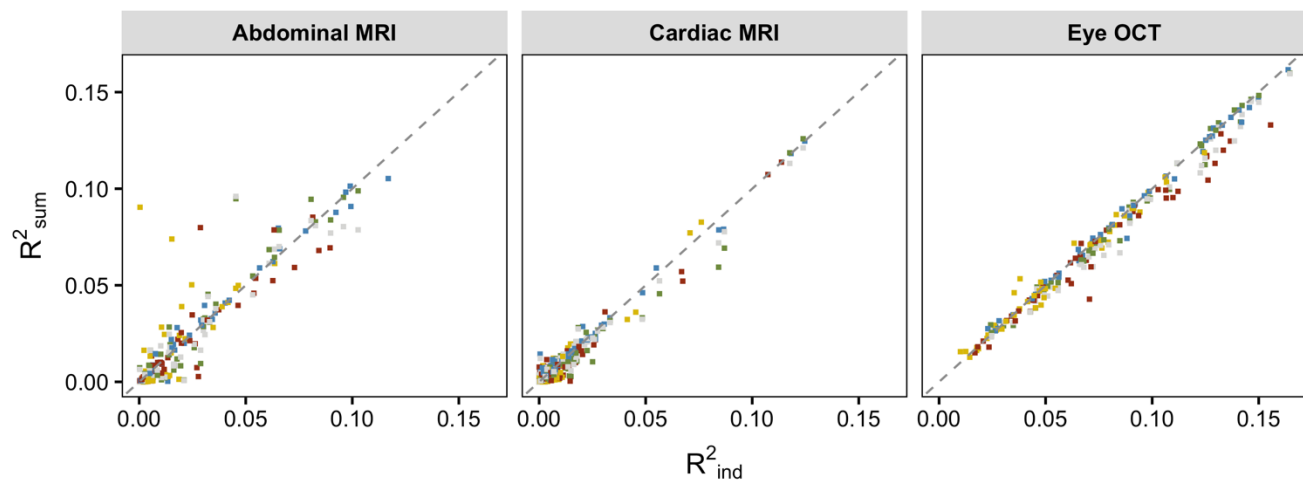
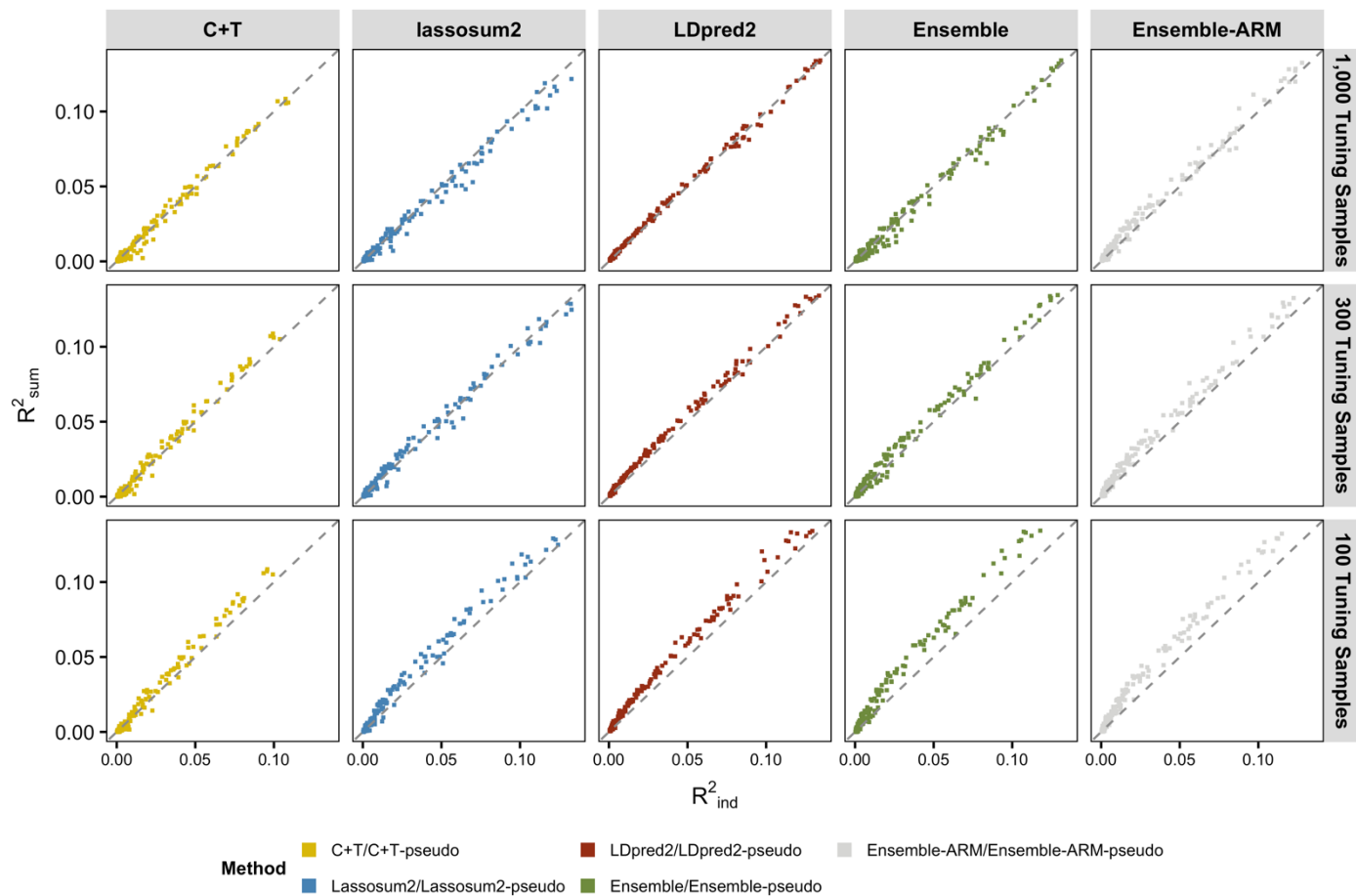
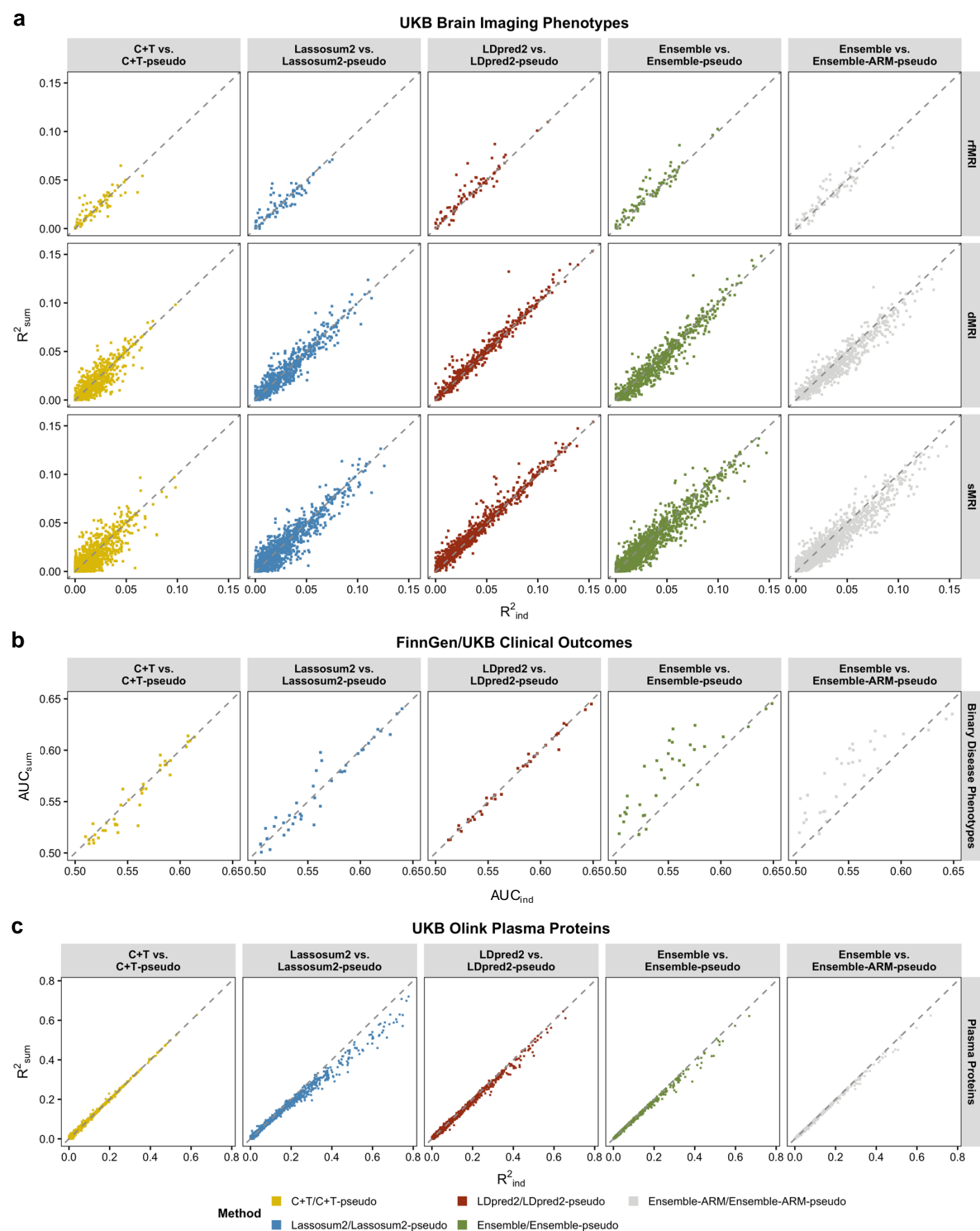


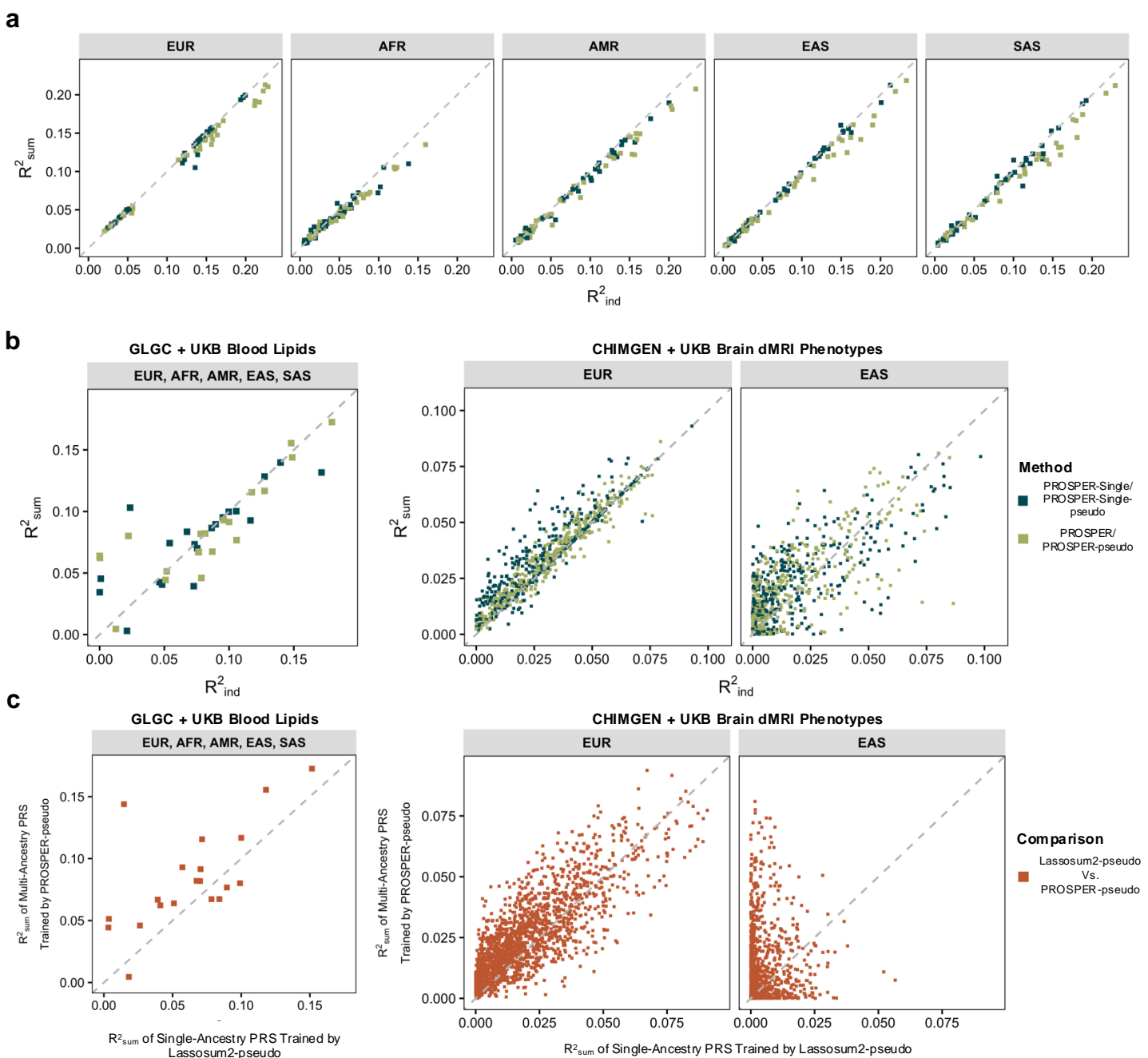
Fig. 2

**a. Individual-Level Tuning Sample Size = 1000****b. Individual-Level Tuning Sample Size = 400****c. Individual-Level Tuning Sample Size = 100****Fig. 3**

**a****UKB Body Imaging Phenotypes****b****Fig. 4**



**Fig. 5**



**Fig. 6**