1 **PGSbuilder: An end-to-end platform for human genome association analysis**

2 **and polygenic risk score predictions**

3

4 Ko-Han Lee[1,†], Yi-Lun Lee[1,†], Tsung-Ting Hsieh[1,†], Yu-Chuan Chang[1,†], Su-Shia Wang[1], Geng-

5 Zhi Fann[1], Wei-Che Lin[1], Hung-Ching Chang[1], Ting-Fu Chen[1], Peng-Husan Li[1], Ya-Ling Kuo[1],

6 Pei-Lung Chen[2,3,4,5], Hsueh-Fen Juan[1,6,7], Huai-Kuang Tsai[1,8], Chien-Yu Chen[1,7,9,*], Jia-Hsin

7 Huang[1,*]

8

9 [1]Taiwan AI Labs & Foundation, Taipei 10351, Taiwan

10 [2]Graduate Institute of Medical Genomics and Proteomics, National Taiwan University College of

11 Medicine, Taipei 10617, Taiwan

12 [3]Department of Medical Genetics, National Taiwan University Hospital, Taipei 10617, Taiwan

13 [4] Genome and Systems Biology Degree Program, National Taiwan University and Academia

14 Sinica, Taipei 11529, Taiwan

15 [5]Graduate Institute of Clinical Medicine, National Taiwan University College of Medicine,

16 Taipei 10051, Taiwan

17 [6]Department of Life Science, National Taiwan University, Taipei 10617, Taiwan

18 [7]Center for Computational and Systems Biology, National Taiwan University, Taipei 10617,

19 Taiwan

20 [8]Institute of Information Science, Academia Sinica, Taipei, 11529, Taiwan

21 [9]Department of Biomechatronics Engineering, National Taiwan University, Taipei 10617,

22 Taiwan

23

24 [†]Ko-Han Lee, Tsung-Ting Hsieh, Yi-Lun Lee, and Yu-Chuan Chang contributed equally to this

25 work.

26

27 [*]Correspondence: Chien-Yu Chen (chienyuchen@ntu.edu.tw); Jia-Hsin Huang

28 (jiahsin.huang@ailabs.tw)

29

30

## Abstract

Understanding the genetic basis of human complex diseases is increasingly important in the development of precision medicine. Over the last decade, genome-wide association studies (GWAS) have become a key technique for detecting associations between common diseases and single nucleotide polymorphisms (SNPs) present in a cohort of individuals. Alternatively, the polygenic risk score (PRS), which often applies results from GWAS summary statistics, is calculated for the estimation of genetic propensity to a trait at the individual level. Despite many GWAS and PRS tools being available to analyze a large volume of genotype data, most clinicians and medical researchers are often not familiar with the bioinformatics tools and lack access to a high-performance computing cluster resource. To fill this gap, we provide a publicly available web server, PGSbuilder, for the GWAS and PRS analysis of human genomes with variant annotations. The user-friendly and intuitive PGSbuilder web server is developed to facilitate the discovery of the genetic variants associated with complex traits and diseases for medical professionals with limited computational skills. For GWAS analysis, PGSbuilder provides the most renowned analysis tool PLINK 2.0 package. For PRS, PGSbuilder provides six different PRS methods including Clumping and Thresholding, Lassosum, LDPred2, GenEpi, PRS-CS, and PRSice2. Furthermore, PGSbuilder provides an intuitive user interface to examine the annotated functional effects of variants from known biomedical databases and relevant literature using advanced natural language processing approaches. In conclusion, PGSbuilder offers a reliable platform to aid researchers in advancing the public perception of genomic risk and precision medicine for human disease genetics. PGSbuilder is freely accessible at http://pgsb.tw23.org.

## Keywords

61

## Introduction

An ultimate goal of human genetics is to understand the genetic basis of human diseases, diagnosis, and management. Results from a large amount of genome-wide association studies (GWAS) have vastly demonstrated that many single nucleotide polymorphisms (SNP) genetic variants are associated with various complex traits[1]. In early 2023, more than 6,300 studies have conducted to map over 496,000 associations between human SNPs and diseases/traits in the GWAS catalog[2]. In the past two decades, the successes of GWAS not only drive the discovery of deleterious mutations linked to certain disease phenotypes but also imply a general pattern of polygenicity of common diseases[3,4]. Many common diseases that conform to polygenic inheritance are underpinned by multiple genetic variants with small or moderate effects[5]. After the realization of a large proportion of the variance in genetic liability to common diseases, utilization of causative risk alleles based on the GWAS discoveries for disease risk prediction has become the potential to stratify patients for precision prevention[6,7].

Polygenic risk score (also known as polygenic scores; PRS) is an important methodology to leverage the genetic contribution of an individual's genotype to measure the genetic liability to complex traits or diseases[8,9]. Clumping and thresholding (C+T)[10] is the primary PRS method based on the summary statistics from GWAS by pruning SNPs through a process of Linkage Disequilibrium (LD) clumping and selecting a *P*-value threshold. Still, it has limitations in the predictive performance without considering other genetic factors. Currently, several PRS methods based on the summary statistics apply a different selection of the prior distribution on the effect sizes of the SNPs under the Bayesian framework. For example, LDpred[11] and LDpred2[12] improve the prediction performance by enhancing LD modeling based on the normality assumption. PRS-CS[13] introduces a different concept to provide a continuous shrinkage (CS) prior to accommodate diverse underlying genetic architectures. Alternatively, SBayesR[14] and SDPR[15] assume a different mixture of normal distributions on the individual-level data as input for adaptive modeling of SNP effect size. Lassosum[16] implements a penalized regression approach with a Lasso-type penalty. Empirical evidence from benchmark experiments shows that not a single method clearly outperforms all other methods in the prediction accuracy for all the simulated data and disease traits[12–14,17]. Nevertheless, each different PRS method can

92    potentially improve the development of PRS construction with specific optimization procedures.

93    Recent studies have demonstrated that the comparison of many PRS methods could facilitate the

94    future implementation of PRS in clinical settings[18,19]. Although a few practical guidelines have

95    introduced how best to perform PRS analyses[20–22], a steep learning curve of implementing those

96    PRS packages and the computing resources required by some tools are impractical for doctors

97    and clinical professionals.

98

99    As the popularity of PRS increases, over 400 publications report more than 3,200 polygenic

100   scores in the Polygenic Score Catalog (https://www.PGSCatalog.org)[23]. However, those PRS

101   studies were predominantly conducted on individuals of European descent[24]. Due to the poor

102   transferability of PRS across populations[25,26], one critical step toward effectiveness in PRS

103   accuracy is to conduct PRS development for the diversity of participants from different

104   ancestries. Along with the cost of a single genetic test per individual plummeting to less than

105   US$50, it becomes feasible to acquire a sufficient cohort size for PRS from the population with

106   underrepresented ancestries by the medical institutes in different countries. In addition, the

107   current consensus about the refinement of PRS should include other informative clinical factors

108   based on their healthy records. To facilitate genetic analysis and PRS development, a

109   sophisticated analysis platform could enable the construction of PRS in clinical research

110   efficiently. For example, impute.me is a recently developed web tool to provide basic PRS

111   estimation using a single method of LDpred to predict individual polygenic risks[27]. To increase

112   the clinical practice of PRS, a comprehensive comparison of different PRS methods could

113   leverage the extent of predictive values into a better understanding of the genetic liability for

114   disease traits.

115

116   In this study, we present PGSbuilder which is an integrated cloud-based platform to analyze

117   human genotype data. PGSbuilder provides a one-stop service to conduct both GWAS and PRS

118   analyses and interactively visualize the analysis results. In PGSbuilder, users can run six

119   different PRS methods as well as the PRS models with clinical factors to compare their

120   performances concurrently. To the best of our knowledge, no other existing web server offers the

121   possibility to compare multiple PRS models. Further, the interpretation of PRS is needed to

122   apply the scores into biological explanations and clinical use. Notably, PGSbuilder also

123 integrates the variant annotation automatically for the candidate SNPs from GWAS and PRS

124 analyses using Ensembl Variant Effect Predictor (VEP)[28] and biomedical literature mining from

125 pubmedKB[29]. In addition, our web interface allows easy access to link all genetic analysis results

126 and candidate SNP information with interactive displays. Finally, users can download all the

127 analysis output files for further exploration.

128

129

130 **Materials and Methods**

131 Data privacy and security

132 Because genetic data will be uploaded to our server, a wide array of security measures are in

133 force to ensure data privacy and security. Our local server has ISO 27001 certification for

134 implementing an information security management system (ISMS). In addition, our server is

135 designed based on the express MVC (Model-View-Controller) framework that encapsulates our

136 features surrounded by powerful security layers. All interactions with the server are protected

137 and secured with HTTPS. Any input data is deleted from our server once the analysis is

138 completed. With the encryption by a firm one-time password, all analyzed results can only be

139 accessed by the data uploader via an encrypted connection, within a 14 days timeframe.

140

141 GWAS

142 To conduct quality control (QC) procedures and following genome-wide association studies

143 (GWAS), we utilize PLINK 2.0, a comprehensive genome association analysis tool for

144 population genetics[30]. There are three major steps for QC and two for GWAS. QC consists of

145 variant filtering, individual filtering, and population stratification while GWAS analysis consists

146 of principal component analysis (PCA) and association test.

147

148 First, unqualified SNPs are filtered out according to the minor allele frequency, Hardy-Weinberg

149 equilibrium, and missingness. Secondly, individuals with the high missing rate of SNPs, large

150 deviation of heterozygosity rate, and high kinship coefficient[31] are also removed. Finally, to

151 exclude individuals with different populations, population stratification is conducted against the

152 population in HapMap 3[32]. Most of the QC criteria and recommended thresholds are referred to

153 Marees et al[33].

154

155    For the GWAS analysis, the top 10 principal components extracted from PCA are used to correct

156    the genetic difference between in-group individuals[34]. Of note, the population stratification

157    during the QC analysis is also conducted via PCA to remove outliers at the level of population,

158    such as Asians, Africans, or Europeans. Next, the principal components and other provided

159    covariates are included to correct the genetic effect during association tests. Only the effect size

160    of autosomal SNPs is calculated using the "glm" function of PLINK 2.0[30,35].

161

162    ## PRS methods

163    In PGSbuilder, the input dataset is separated into the base, target, and test sets, respectively.

164    First, QC is applied on both base and target sets, and then GWAS is only performed on the base

165    set to get the summary statistics. Combining the summary statistics with the target set which is

166    used for the calculation of linkage disequilibrium (LD) and the selection of hyperparameters,

167    PGSbuilder performs PRS analysis to build models based on different methods. This pipeline of

168    PRS analysis is referred to Choi et al[21]. There are six PRS methods provided in PGSbuilder,

169    including clumping and thresholding, PRSice2, LDpred2, Lassosum, PRS-CS, and GenEpi. Five

170    methods, except GenEpi, are selected to produce PRS prediction from the external summary

171    statistics without individual genetic data. On the other hand, GenEpi method is included due to

172    its consideration of gene-based epistasis, which is a distinct machine learning-based algorithm to

173    estimate PRS, for comparison.

174

175    *Clumping and Thresholding*: Clumping and thresholding (C+T) is the classical algorithm that

176    adjusts the LD using clumping and selects SNPs with *P*-value less than a specified threshold to

177    calculate the PRS for each individual[10]. In PGSbuilder, SNPs within 250 kb away from the index

178    SNP and have the R-squared over 0.1 with it are assigned to the clump of the index SNP. Nine

179    thresholds, including $10^{-8}$, $10^{-7}$, $10^{-6}$, $10^{-5}$, $10^{-4}$, $10^{-3}$, $10^{-2}$, $10^{-1}$, and 1, are applied to the clumped

180    SNPs to build PRS models. Beta scores derived from the summary statistics are set as the effect

181    size estimates directly. The model with the best performance on the target set is selected as the

182    final PRS model.

183

184     *PRSice2*: PRSice2 is also a clumping and thresholding-based PRS algorithm with a higher
185     resolution of thresholds[17]. SNPs with a minor allele frequency lower than 0.01 are filtered out.
186     Like C+T, beta scores are set as the effect size estimates directly.

187

188     *Lassosum*: Lassosum uses penalized regression to adjust the effect size of SNPs for a PRS
189     model[16]. The summary statistics provide the SNP-wise correlation with the phenotype and the
190     initial effect size of SNPs. LD blocks are defined from the subpopulation of the 1000 Genome
191     database, and the LD matrix is calculated from the target set. Additionally, the target set is used
192     for the selection of hyperparameters to get the best PRS model.

193

194     *LDpred2*: LDpred2 is a Bayesian PRS predictor by adjusting the effect size of SNPs from the
195     summary statistics[12]. The target set provides the correlation between SNPs for LD estimation
196     within 3 centimorgan. In PGSbuilder, for summary statistics having more than 10 SNPs with *P*-
197     value$<10^{-8}$, we implement the "LDpred2-grid" mode to select the best hyperparameters,
198     including the proportion of causal variants and the heritability. On the other hand, for those with
199     less significant SNPs, we implement the "LDpred2-inf" mode, an infinitesimal model.

200

201     *PRS-CS*: PRS-CS is a Bayesian polygenic prediction method that infers the posterior effect size
202     of SNPs from the summary statistics using continuous shrinkage priors[13]. In PGSbuilder, we use
203     the 1000 Genome dataset as the reference panel for LD estimation. The global shrinkage
204     parameter is fixed at 0.2 and other parameters are left as defaults.

205

206     *GenEpi*: GenEpi, a machine learning approach, takes both additive effect and SNP-SNP
207     interactions into consideration to build a PRS model from the raw genomic data[36]. GenEpi uses
208     two-stage feature selection to select a single SNP, intragenic interaction, and intergenic
209     interaction and then applies a regression model to fit the selected features. In PGSbuilder, we
210     only train the GenEpi model on the base set.

211

212     Covariates

213     In GWAS analysis, covariates are used to adjust the genetic effect on the target phenotype.
214     PGSbuilder performs PCA before GWAS, and the top ten principal components (PCs) are served

215  as covariates. In addition, users can provide a covariate file, and covariates with the variance

216  inflation factor (VIF) less than 50 or a missing rate over 20% are removed. Finally, the effect

217  size of each SNP is corrected with PCs and provided covariates during the association test.

218

219  On the other hand, to provide a comprehensive risk assessment for individuals, features other

220  than genetic factors should be taken into consideration. After building a PRS model, PGSbuilder

221  combines the PRS score as a genetic factor and user-provided covariates as clinical factors to

222  build a regression model trained on the target set. Then, PGSbuilder predicts each individual

223  using this regression model to stratify the risk of the target phenotype.

224

225  Variant annotation tools

226  The annotation of significant SNP from GWAS or other genomic analysis is of great importance.

227  Annotation of variants is vital for the translation of genomic results to the functional level for

228  further analysis. The Ensembl Variant Effect Predictor (VEP) is an open-source, powerful, and

229  versatile toolset for the annotation and prioritization of genomic variants for a transcript or even

230  non-coding region[28]. We select VEP (version 106) because of its broad collection of databases,

231  scalability, and free open license. In order to display the important variant information to show

232  first on the web page, PGSbuilder sorts the VEP results by several criteria, including transcript

233  consensus, mutation consequence, mutation severity, and feature biotype. The complete VEP result

234  is provided in the downloaded file. In addition, allele frequencies from Taiwan Biobank[37] and

235  1000 Genome Project are provided in the VEP annotation.

236

237  Moreover, we integrate our literature mining engines, variant2literature[38] and pubmedKB[29], by

238  retrieving entity mentions and odds ratio statistics to create a report of textual evidence for each

239  variant-phenotype pair. The literature report contains an overall summary and single paper

240  snippets. For the overall summary, we first collect sentences and clinical case sentences where

241  the target variant and phenotype are both mentioned. We then present the most important

242  sentences and clinical cases identified by page rank[39]. For single paper snippets, we present the

243  paragraph describing odds ratio statistics of the target variant and phenotype.

244

245     Example Data

246     *Taiwan Biobank*: Taiwan Biobank (TWB) is a prospective cohort study with genomic data and a

247     variety of phenotypes collected from Taiwanese population[37]. The TWB cohort contains 27,500

248     individuals genotyped for 653,288 SNPs on the TWB v1.0 array as well as 68,978 individuals

249     genotyped for 748,344 SNPs on the TWB v2.0 array.

250

251     *NIA ADC Cohort*: The NIA ADC Cohort consists of individuals evaluated clinically from

252     National Institute on Aging (NIA)-funded Alzheimer Disease Centers (ADC)[40]. Inclusion criteria

253     of late-onset Alzheimer's disease are autopsied subjects with age >60 or cases diagnosed with

254     DSM-IV or Clinical Dementia Rating >1[40]. All the seven ADC datasets downloaded from

255     NIAGADS (https://www.niagads.org/datasets) were merged directly as a joint analysis. In total,

256     there are 10,256 samples, including 5,334 cases, 3,973 controls, and 949 unknowns, genotyped

257     for 914,402 SNPs.

258

259

260     **Results**

261     PGSbuilder analysis workflow

262     PGSbuilder is a web-based server to provide end-to-end analysis for genetic cohort data

263     including GWAS, PRS, and variant annotation. The GWAS analysis aims to figure out the

264     significant SNPs associated with a specific phenotype while the PRS analysis aims to build a

265     model for the estimation of the individual risk. After the analysis, SNP-level annotation and

266     literature exploration using pubmedKB[29] are performed to provide useful insights into causal

267     variants.

268

269     The GWAS pipeline (Figure 1), which is applied to the whole input dataset, consists of quality

270     control (QC) and association tests. As for the PRS pipeline (Figure 1), the input dataset is firstly

271     separated into training and test sets. The training set is undergone QC steps and split into base

272     and target subsets. The base subset is used to obtain the summary statistics of GWAS, while the

273     target subset is used to build the PRS models. On the other hand, the option of an external

274     summary statistics file is available in PGSbuilder. When the external summary statistics file is

275    provided, it replaces the base subset to provide GWAS results and the entire training set serves
276    as the target subset alternatively. To build a PRS model, most PRS methods combine the
277    summary statistics providing the initial SNP effect sizes with the linkage disequilibrium (LD)
278    estimation derived from the target subset. Of note, GenEpi is unavailable for building a PRS
279    model from the external summary statistics. Finally, to validate the model performance, the
280    estimated risks of individuals in the test set are independently calculated by the adjusted effect
281    size.

282

283    System implementation

284    We used Kubernetes and docker technology to group our applications including web interfaces,
285    data processing, GWAS and PRS pipelines, and variant annotation into a service platform. For
286    the web interface, we adopted React architecture and Node.js for the frontend and backend
287    respectively. For the analysis, after users upload genotype data, PGSbuilder will create pods for
288    GWAS and PRS pipelines dynamically and instantly. Significant variants derived from GWAS
289    and PRS pipelines will be annotated through the VEP and pubmedKB to determine the effect of
290    variants in the public database and academic literature.

291

292    For the security of private genomic data, users have to sign up via email activation. After login,
293    two studies including a binary trait (classification model) and a quantitative trait (regression
294    model) are demonstrated on the analysis page. To create a new study, users have to upload
295    genotype data in PLINK format and fill in relevant information such as population, genome
296    build, and prediction method (classification or regression). PGSbuilder provides flexibility for
297    users to modify some quality control parameters and select multiple PRS methods (Figure 2A). If
298    the data is successfully uploaded to the PGSbuilder server, the job is added to the analysis queue
299    and will be processed as soon as possible. Users will receive an email notice to check the state of
300    jobs on the running page. Once the job is completed, users can download a comprehensive report
301    for GWAS and PRS results. PGSbuilder also provides an interactive interface to view the result
302    in detail.

303

304    On the GWAS result page, PCA plot, quantile-quantile plot (Q-Q plot), Manhattan plot, and the
305    variant table are demonstrated (Figure 2B). PCA is used for the correction of population

10

306    stratification, and the top 10 principal components (PCs) are selected as covariates for GWAS.

307    The paired distributions of the top 3 PCs are shown interactively, and users can arbitrarily switch

308    between three figures through arrow buttons. In addition, each dot represents a sample whose ID

309    will be displayed via a mouseover event, which can help users discriminate outliers. The Q-Q

310    plot is provided to evaluate the deviation of observed $P$-values from expected $P$-values under a

311    uniform distribution. For the Manhattan plot and variant table, we set a suggestive $P$-value

312    threshold of $1 \times 10^{-5}$ and a strict $P$-value threshold of $5 \times 10^{-8}$. SNPs with a $P$-value smaller than

313    the threshold are colored in orange and listed in the variant table. The SNPs in the Manhattan

314    plot and the variant table are interactive. Clicking on an orange point on the Manhattan plot

315    navigates the variant table to the corresponding SNP with its information, and vice versa.

316    Besides, users can search for a specific SNP through the search bar. More detailed information of

317    all SNPs including their $P$-values and annotated information are compressed as a zip file to be

318    downloaded.

319

320    On the PRS result page, we compare the performance of selected PRS methods. The quantile

321    plot shows the risk stratification (Figure 2C). For each method, samples in the test set are divided

322    into 10 quantiles of increasing PRS. Then, in each quantile, the odds ratio is calculated for binary

323    phenotypes while the mean of values is calculated for quantitative phenotypes. A great difference

324    between the first and the last group represents a good risk stratification. Of note, all individuals

325    in the test set serve as the baseline for odds ratio calculation for binary tracts. In the classification

326    analysis for a binary tract, the receiver operating characteristic (ROC) curve and distribution plot

327    for each method are demonstrated (Figure 2C). The area under the ROC curve illustrates the

328    performance and the distribution plots illustrate the prediction distribution for cases against

329    controls. In the regression analysis for a quantitative tract, Spearman correlations and scatter

330    plots are shown (Figure 2C). The Spearman correlation is performed to evaluate the performance

331    and the scatter plot with a regression line illustrates the relationship between phenotypes and

332    prediction rankings for each method. The tabs of method lists allow users to switch results

333    between different methods. Users can click one of them to view the corresponding performance

334    and variant table.

335

336  Furthermore, analysis beyond genetic factors is also available in PGSbuilder. If the covariate file

337  is provided, covariates will be used to correct the effect size of SNPs during GWAS, and then

338  serve as clinical factors combined with PRSs to build a regression model for risk prediction. The

339  performance with or without clinical factors is also demonstrated in the figures for comparison.

340  The weight of each clinical factor is shown in a table for users to figure out important factors.

341

342  ## Variant annotation panel

343  In order to help interpret GWAS and PRS results, PGSbuilder provides a comprehensive variant

344  annotation panel for users to explore biological significance. There are often a large number of

345  SNPs associated with a phenotype. PGSbuilder will automatically sort the important SNPs at the

346  top of the panel according to several annotation information including transcript consensus,

347  mutation consequence, mutation severity, and feature biotype. Figure 3 displays an example of the

348  significant SNP information from the GWAS results. Accordingly, three key features are present

349  including variant effect prediction information, external links about the variant, and the related

350  literature. PGSbuilder uses ClinVar[41] and VEP[28] for variant interpretation (Fig. 3B). Several

351  external links are provided to easily navigate the further variant information (Fig. 3C). Lastly,

352  PGSbuilder integrates the literature mining results from the pubmedKB[29] to assist researchers

353  and clinical professionals in obtaining the related literature.

354

355  ## System performance

356  For benchmarking, we recorded execution time, average memory, and CPU usage for QC,

357  GWAS, and PRS methods with 680k SNPs given 20k, 50k, and 110k samples (Table 1). The

358  resource for each execution was limited to 20 GB and 10 CPUs. Obviously, more resources were

359  needed as the sample size increased. Table 1 shows the comparison between six PRS methods.

360  PRSice2, PRS-CS, and GenEpi took much more execution time than the others, but PRSice2 and

361  GenEpi used the least CPU and memory respectively. In conclusion, it takes about three days to

362  complete a comprehensive PRS analysis for a dataset with 110k samples and 680k SNP.

363

364

**Table 1.** The system performance, including execution time, average CPU, and memory of PGSbuilder. We performed QC, GWAS, and six PRS methods (classification for a binary trait) on a dataset with the same number of SNPs but different sample sizes.

| Sample | STATS | QC | GWAS | C+T | Lassosum | LDpred2 | PRSice2 | PRS-CS | GenEpi | Total |
|---|---|---|---|---|---|---|---|---|---|---|
| 20k | Time (min) | 8.0 | 8.0 | 15.5 | 18.8 | 27.3 | 79.2 | 171.8 | 450.9 | 779.5 |
| 20k | Avg. CPU | 2.9 | 6.9 | 5.6 | 4.7 | 6.3 | 3.3 | 7.5 | 7.7 | |
| 20k | Memory (GB) | 5.1 | 2.0 | 10.9 | 10.5 | 10.9 | 10.7 | 8.4 | 3.8 | |
| 50k | Time (min) | 43.5 | 20.0 | 33.4 | 41.6 | 73.9 | 180.8 | 268.6 | 715.0 | 1376.8 |
| 50k | Avg. CPU | 4.7 | 6.8 | 7.9 | 7.6 | 7.6 | 5.6 | 8.0 | 7.4 | |
| 50k | Memory (GB) | 16.0 | 1.7 | 16.2 | 16.1 | 15.8 | 14.8 | 12.9 | 3.8 | |
| 110k | Time (min) | 139.5 | 77.0 | 200.3 | 220.3 | 251.3 | 510.6 | 967.3 | 2086.9 | 4453.1 |
| 110k | Avg. CPU | 4.8 | 8.6 | 7.4 | 7.3 | 7.6 | 5.5 | 8.0 | 7.2 | |
| 110k | Memory (GB) | 19.4 | 13.8 | 17.0 | 17.0 | 16.7 | 13.9 | 12.0 | 10.6 | |

## Case Study

To demonstrate the capability of PGSbuilder, we performed two case studies using the cohorts with a large number of individuals and corresponding phenotypes. Firstly, in the Taiwan Biobank (TWB)[37], a Taiwanese cohort composed of healthy adults, we previously defined nine quantitative traits and five binary traits related to some common chronic diseases, such as type 2 diabetes or dyslipidemia, according to their phenotypic measures (see https://github.com/chienyuchen/TWB-PRS for more information). The presented GWAS and PRS models across fourteen traits in the TWB were built by using PGSbuilder. Among them, low-density lipoprotein (LDL), a quantitative trait, was selected here to demonstrate the usage of

13

378 adding covariates and the leverage of external summary statistics to run PGSbuilder. Secondly,

379 for the cohort with a specific disease, we performed GWAS and PRS analysis on the National

380 Institute on Aging (NIA)-funded Alzheimer Disease Centers (ADC) Cohort[40] to demonstrate the

381 result of a binary trait.

382

383 *Low-density lipoprotein*: Low-density lipoprotein (LDL), which is a kind of lipoprotein to

384 transport fat molecules around the body, acts as the primary driver of atherogenesis resulting in

385 cardiovascular diseases[42]. Several genes, such as LDLR, PCSK9, and APOB, affecting the

386 quantity of LDL in circulation have been reported[43]. Recognizing people with a genetic tendency

387 for high LDL could help them by providing early intervention to avoid the progression of severe

388 cardiovascular diseases. Therefore, in this study, we applied GWAS and PRS analysis using

389 PGSbuilder on the TWB data. The covariates, including age, sex, and body mass index (BMI),

390 were added to correct GWAS for genetic factors and then serve as clinical factors to build

391 regression models for risk prediction.

392

393 With the default QC settings of PGSbuilder, 55,412 samples and 276,068 SNPs were passed the

394 quality control (Table S1-2). To control the population stratification, PGSbuilder always

395 performs PCA analysis and applies the top ten principal components (PCs) as covariates during

396 GWAS. Figure 4A demonstrates the distribution of PC1 and PC2 to confirm SNPs without

397 unusual differentiation between quantiles in the TWB data. The interactive Manhattan plot is

398 shown in Figure 4B and the significant SNPs with a *P*-value $< 10^{-5}$ are highlighted in orange for

399 clicking to navigate variant information. Notably, in comparison with the previous study using

400 the same TWB data[44], highly similar results were observed in PGSbuilder as shown that more

401 than 80% (89/111) of significant SNPs in the TWB arrays were identically found to associate

402 with the LDL trait. That is, the pipeline in PGSbuilder is indeed reproducible.

403

404 In addition, PGSbuilder allows users to provide external summary statistics to build PRS models.

405 Herein, the external summary statistics from the BioBank Japan[45] to identify significant variants

406 and stratify people by the risk of high LDL were applied to estimate PRS in the TWB data.

407 Figure 4C shows the performance on the test set of each PRS method with and without clinical

408 factors. Overall, PRS combined with clinical factors performs better than PRS-only and clinical

409    factors-only models. These results indicate that the genetic factor combined with clinical factors

410    provide a better prediction effect. Figure 4D depicts the risk stratification of models using

411    clinical factors. "PRS + clinical factors" models stratified the test set better than the "clinical

412    factors-only" model. In the "PRS + clinical factors" models, the difference in average LDL

413    between the first and last groups is up to forty. Furthermore, the weight of each feature in the

414    "PRS-clinical factors" model is listed in Table 2, where PRS has the largest contribution in all

415    the models.

416

417    **Table 2.** The weight of PRS and clinical factors for "PRS + clinical factors" models of LDL.

|         | C+T  | PRSice2 | Lassosum | LDpred2 | PRS-CS |
|---------|------|---------|----------|---------|--------|
| **PRS** | 7.96 | 7.93    | 8.70     | 5.29    | 5.87   |
| **Sex** | 2.57 | 2.57    | 2.64     | 2.55    | 2.53   |
| **Age** | 4.29 | 4.29    | 4.26     | 4.31    | 4.29   |
| **BMI** | 4.41 | 4.41    | 4.45     | 4.33    | 4.34   |

418

419    *Alzheimer's disease*: Alzheimer's disease (AD), the major cause of dementia, is a complex

420    disorder associated with genetic factors and environmental factors[46]. Several genetic loci, such as

421    APOE, have been identified at the level of association study[47,48]. Combining the effects of these

422    genetic loci to build a PRS model could provide individuals with the disease risk for further

423    preventive strategies[49]. In this study, to build PRS models based on different methods and

424    compare the performance of them, we analyzed the National Institute on Aging (NIA)-funded

425    Alzheimer Disease Centers (ADC) cohort using PGSbuilder.

426

427    Figure 5 shows the performance of PRS analysis from PGSbuilder. There are two obvious

428    groups with different performances. C+T, PRSice2, Lassosum, and GenEpi have better auROC

429    than LDpred2 and PRS-CS (Figure 5A). Figure 5B depicts the prediction distribution of cases

430    and controls; the more distance between the distributions the better performance of the model.

431    For further comparison of different methods, an UpSet plot depicts the intersection of top-100

432    valuable SNPs from each method (Figure 5C). Notably, LDpred2 and PRS-CS have some

15

433    distinct SNPs than others, which might cause noise for the PRS prediction and decrease the
434    model performance.

435

436    To investigate the information of SNPs, PGSbuilder annotates SNPs using VEP[28] and
437    pubmedKB[29]. For example, Figure 5D shows the annotation of rs157580, which is an intron
438    variant of gene TOMM40 with average allele frequency across different populations. A previous
439    study (PMID: 21867541) also reported that rs157580 was significantly associated with AD[50].
440    The literature mining of PubMed abstracts by pubmedKB facilitates users to interpret the
441    variants more readily.

442

443

## Discussion

445    PGSbuilder is a cloud-based platform that offers comprehensive genotyping analyses, including
446    GWAS and PRS, all in one place. Our goal for GWAS is to help identify significant SNPs
447    associated with the target phenotype, while for PRS, we aim to assist evaluation of the prediction
448    performance of polygenic models. Customized settings are available for users to adjust the
449    analytic process, such as quality control, population stratification, and the selection of PRS
450    methods. With PGSbuilder's interactive interfaces, users can easily interpret their results. For
451    instance, users can select specific SNPs on the Manhattan plot and view the corresponding
452    annotations in the table. Additionally, PGSbuilder integrates pubmedKB for variant
453    interpretation by providing literature support. With these features, PGSbuilder is a
454    comprehensive and user-friendly platform for GWAS and PRS.

455

456    In addition to the analytic pipeline, PGSbuilder offers various visualization plots to compare the
457    performance of different PRS methods. To evaluate risk stratification, the quantile plot is a key
458    interpretation tool. The UpSet plot enables users to observe the intersection of important SNPs
459    selected from each method. Additionally, PGSbuilder incorporates our original GenEpi
460    software[36], which provides a unique method to uncover the genetic epistasis associated with
461    phenotypes, as demonstrated in other recent studies[51,52]. Finally, as clinical factors are provided,
462    PGSbuilder will rank the weights of them and PRS to highlight the most predictive feature,
463    which helps users investigate the risk factor precisely.

16

464

465     While PGSbuilder provides a range of useful features, there are some limitations to its

466     functionality. First, it is important to consider the limitations of hardware resources when dealing

467     with large datasets. For example, some imputed files containing 10 million SNPs and 50K

468     samples may not be immediately accessible due to these restrictions. However, computationally

469     efficient methods such as C+T, Lassosum, and PRSice2 can eb effectively applied to such

470     datasets, based on our internal experiments. It is worth noting that building a predictive model

471     using some PRS methods may require a significant amount of time. On the other hand, GenEpi,

472     which discovers the gene-based epistasis, is not practical for imputed data due to its

473     computational complexity. Secondly, some known PRS methods, such as those based on a

474     mixture model for SNP effective size (e.g. SBayesR[14], DPR[53], DBSLMM[54]), are currently not

475     included in PGSbuilder. Lastly, PRS models can only be downloaded from PGSbuilder output

476     directly. Going forward, we are planning to implement a prediction module that allows users to

477     upload other datasets and then automatically obtain predictions of available PRS models .

478

479     The field of PRS development is growing rapidly, with mounting evidence using the wealth of

480     data collected in biobanks[55–58]. As the proof of concept is solidly demonstrated, an effective and

481     comprehensive platform is necessary to perform GWAS and PRS analysis for diseases that are

482     not covered by biobanks. PGSbuilder provides researchers with the ability to identify significant

483     loci with annotations and investigate the polygenicity of a target phenotype across a specific

484     population effectively. By leveraging genotypes, a PRS model has the clinical potential to offer

485     risk evaluations to individuals. This, in turn, can facilitate early surveillance for severe diseases.

486

487

488     **Conclusion**

489     PGSbuilder is an end-to-end platform that seamlessly integrates QC of genotype data, GWAS,

490     PRS, SNP annotation, and visualizations. This platform is versatile, allowing the incorporation of

491     external GWAS summary statistics to run PRS using various methods, thereby enabling the

492     estimation of genetic risk in smaller cohort samples. In addition, PGSbuilder's user-friendly

493     interface is designed to be accessible to users without programming experiences. In the future,

494 we plan to further augment and broaden PGSbuilder by introducing a prediction module that

495 allows users to directly run their PRS models for specific disease phenotypes.

496

497

505

506

## Data and software availability

508 All genetic and phenotype data in TWB described in this paper are publicly available via the

509 Taiwan Biobank data access protocol. Fourteen PRS models using TWB data, including five

510 binary phenotypes and nine quantitative traits, are freely available on the GitHub project

511 repository (https://github.com/chienyuchen/TWB-PRS). The AD data is publicly available to

512 registered researchers by request from the National Institute on Aging Genetics of Alzheimer's

513 Disease Data Storage Site (NIAGADS). The source codes for GWAS and PRS analyses were

514 deposited to Github and is available at https://github.com/ailabstw/PGSbuilder.

515

## Ethics approval and consent to participate

517 The application number of TWB data is TWBR10411-03. This application of NIA ADC Cohort

518 dataset has been filed with the IRB ( 202106049RINA) in order to get approval from NIAGADS.

519

520

## Competing interests

522 The authors declare that they have no competing interests.

523

524

## Authors' contributions

KHL, YLL, TTH, YCC, and HCC conceived and implemented the pipeline development. YCC inspired team members to unite as a product manager, and designed all the frameworks of this web service, including wireframe, prototype, and database schema. SSW, WCL, and GZF implemented the web design and interface. TFC and PHL implemented the literature mining. YLK served as liaisons to user communities. YCC and JHH helped project development and management. PLC led the application of TWB data. HFJ, HKT, CYC, and JHH supervised the project. KHL and JHH led the writing of the manuscript. All authors discussed the results and implications and commented on the manuscript. All authors read and approved the final manuscript.

## References

1.  Mills, M. C. & Rahal, C. A scientometric review of genome-wide association studies. *Commun Biol* **2**, 9 (2019).

2.  Buniello, A. *et al.* The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Research* vol. 47 D1005–D1012 Preprint at https://doi.org/10.1093/nar/gky1120 (2019).

3.  Visscher, P. M. *et al.* 10 Years of GWAS Discovery: Biology, Function, and Translation. *Am. J. Hum. Genet.* **101**, 5–22 (2017).

4.  Stranger, B. E., Stahl, E. A. & Raj, T. Progress and promise of genome-wide association studies for human complex trait genetics. *Genetics* **187**, 367–383 (2011).

5.  Watanabe, K. *et al.* A global overview of pleiotropy and genetic architecture in complex traits. *Nat. Genet.* **51**, 1339–1348 (2019).

6.  Chatterjee, N., Shi, J. & García-Closas, M. Developing and evaluating polygenic risk prediction models for stratified disease prevention. *Nat. Rev. Genet.* **17**, 392–406 (2016).

7.  Wray, N. R. *et al.* From Basic Science to Clinical Application of Polygenic Risk Scores: A Primer.

552     *JAMA Psychiatry* **78**, 101–109 (2021).

553     8.    Ma, Y. & Zhou, X. Genetic prediction of complex traits with polygenic scores: a statistical review.

554     *Trends Genet.* **37**, 995–1011 (2021).

555     9.    Wray, N. R., Goddard, M. E. & Visscher, P. M. Prediction of individual genetic risk to disease from

556     genome-wide association studies. *Genome Res.* **17**, 1520–1528 (2007).

557     10.   International Schizophrenia Consortium *et al.* Common polygenic variation contributes to risk of

558     schizophrenia and bipolar disorder. *Nature* **460**, 748–752 (2009).

559     11.   Vilhjálmsson, B. J. *et al.* Modeling Linkage Disequilibrium Increases Accuracy of Polygenic Risk

560     Scores. *Am. J. Hum. Genet.* **97**, 576–592 (2015).

561     12.   Privé, F., Arbel, J. & Vilhjálmsson, B. J. LDpred2: better, faster, stronger. *Bioinformatics* (2020)

562     doi:10.1093/bioinformatics/btaa1029.

563     13.   Ge, T., Chen, C.-Y., Ni, Y., Feng, Y.-C. A. & Smoller, J. W. Polygenic prediction via Bayesian

564     regression and continuous shrinkage priors. *Nat. Commun.* **10**, 1776 (2019).

565     14.   Lloyd-Jones, L. R. *et al.* Improved polygenic prediction by Bayesian multiple regression on

566     summary statistics. *Nat. Commun.* **10**, 5086 (2019).

567     15.   Zhou, G. & Zhao, H. A fast and robust Bayesian nonparametric method for prediction of complex

568     traits using summary statistics. *PLoS Genet.* **17**, e1009697 (2021).

569     16.   Mak, T. S. H., Porsch, R. M., Choi, S. W., Zhou, X. & Sham, P. C. Polygenic scores via penalized

570     regression on summary statistics. *Genet. Epidemiol.* **41**, 469–480 (2017).

571     17.   Choi, S. W. & O'Reilly, P. F. PRSice-2: Polygenic Risk Score software for biobank-scale data.

572     *Gigascience* **8**, (2019).

573     18.   Ni, G. *et al.* A Comparison of Ten Polygenic Score Methods for Psychiatric Disorders Applied

574     Across Multiple Cohorts. *Biol. Psychiatry* **90**, 611–620 (2021).

575     19.   Pain, O. *et al.* Evaluation of polygenic prediction methodology within a reference-standardized

576     framework. *PLoS Genet.* **17**, e1009021 (2021).

577     20.   Collister, J. A., Liu, X. & Clifton, L. Calculating Polygenic Risk Scores (PRS) in UK Biobank: A

578    Practical Guide for Epidemiologists. *Front. Genet.* **13**, 818574 (2022).

579    21.  Choi, S. W., Mak, T. S.-H. & O'Reilly, P. F. Tutorial: a guide to performing polygenic risk score

580         analyses. *Nat. Protoc.* **15**, 2759–2772 (2020).

581    22.  Wray, N. R. *et al.* Research review: Polygenic methods and their application to psychiatric traits. *J.*

582         *Child Psychol. Psychiatry* **55**, 1068–1087 (2014).

583    23.  Lambert, S. A. *et al.* The Polygenic Score Catalog as an open database for reproducibility and

584         systematic evaluation. *Nat. Genet.* **53**, 420–425 (2021).

585    24.  Martin, A. R. *et al.* Clinical use of current polygenic risk scores may exacerbate health disparities.

586         *Nat. Genet.* **51**, 584–591 (2019).

587    25.  Scutari, M., Mackay, I. & Balding, D. Using Genetic Distance to Infer the Accuracy of Genomic

588         Prediction. *PLoS Genet.* **12**, e1006288 (2016).

589    26.  Wang, Y. *et al.* Theoretical and empirical quantification of the accuracy of polygenic scores in

590         ancestry divergent populations. *Nat. Commun.* **11**, 3865 (2020).

591    27.  Folkersen, L. *et al.* Impute.me: An Open-Source, Non-profit Tool for Using Data From Direct-to-

592         Consumer Genetic Testing to Calculate and Interpret Polygenic Risk Scores. *Front. Genet.* **11**, 578

593         (2020).

594    28.  McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122 (2016).

595    29.  Li, P.-H. *et al.* pubmedKB: an interactive web server for exploring biomedical entity relations in the

596         biomedical literature. *Nucleic Acids Res.* (2022) doi:10.1093/nar/gkac310.

597    30.  Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets.

598         *Gigascience* **4**, 7 (2015).

599    31.  Manichaikul, A. *et al.* Robust relationship inference in genome-wide association studies.

600         *Bioinformatics* **26**, 2867–2873 (2010).

601    32.  Consortium, T. I. H. 3. & The International HapMap 3 Consortium. Integrating common and rare

602         genetic variation in diverse human populations. *Nature* vol. 467 52–58 Preprint at

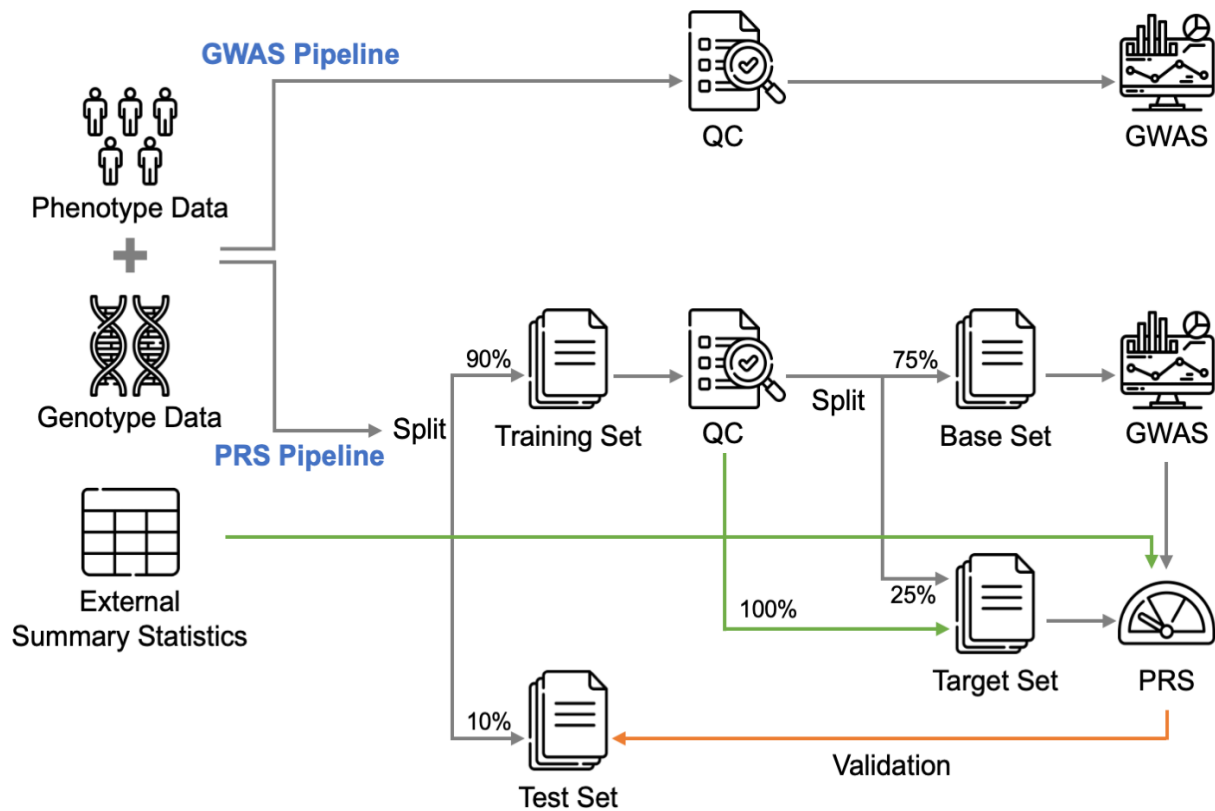603         https://doi.org/10.1038/nature09298 (2010).

604    33. Marees, A. T. *et al.* A tutorial on conducting genome-wide association studies: Quality control and

605        statistical analysis. *International Journal of Methods in Psychiatric Research* **27**, e1608 (2018).

606    34. Price, A. L. *et al.* Principal components analysis corrects for stratification in genome-wide

607        association studies. *Nat. Genet.* **38**, 904–909 (2006).

608    35. Purcell, S. *et al.* PLINK: a tool set for whole-genome association and population-based linkage

609        analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).

610    36. Chang, Y.-C. *et al.* GenEpi: gene-based epistasis discovery using machine learning. *BMC*

611        *Bioinformatics* **21**, 68 (2020).

612    37. Feng, Y.-C. A. *et al.* Taiwan Biobank: a rich biomedical research database of the Taiwanese

613        population. Preprint at https://doi.org/10.1101/2021.12.21.21268159.

614    38. Lin, Y.-H. *et al.* variant2literature: full text literature search for genetic variants. *bioRxiv* 583450

615        (2019) doi:10.1101/583450.

616    39. Page, L., Brin, S., Motwani, R. & Winograd, T. The PageRank Citation Ranking: Bringing Order to

617        the Web. (1999).

618    40. Naj, A. C. *et al.* Common variants at MS4A4/MS4A6E, CD2AP, CD33 and EPHA1 are associated

619        with late-onset Alzheimer's disease. *Nat. Genet.* **43**, 436–441 (2011).

620    41. Landrum, M. J. *et al.* ClinVar: improving access to variant interpretations and supporting evidence.

621        *Nucleic Acids Res.* **46**, D1062–D1067 (2018).

622    42. Borén, J. *et al.* Low-density lipoproteins cause atherosclerotic cardiovascular disease:

623        pathophysiological, genetic, and therapeutic insights: a consensus statement from the European

624        Atherosclerosis Society Consensus Panel. *Eur. Heart J.* **41**, 2313–2330 (2020).

625    43. Borén, J. *et al.* Low-density lipoproteins cause atherosclerotic cardiovascular disease:

626        pathophysiological, genetic, and therapeutic insights: a consensus statement from the European

627        Atherosclerosis Society Consensus Panel. *Eur. Heart J.* **41**, 2313–2330 (2020).

628    44. Chen, C.-Y. *et al.* Analysis across Taiwan Biobank, Biobank Japan and UK Biobank identifies

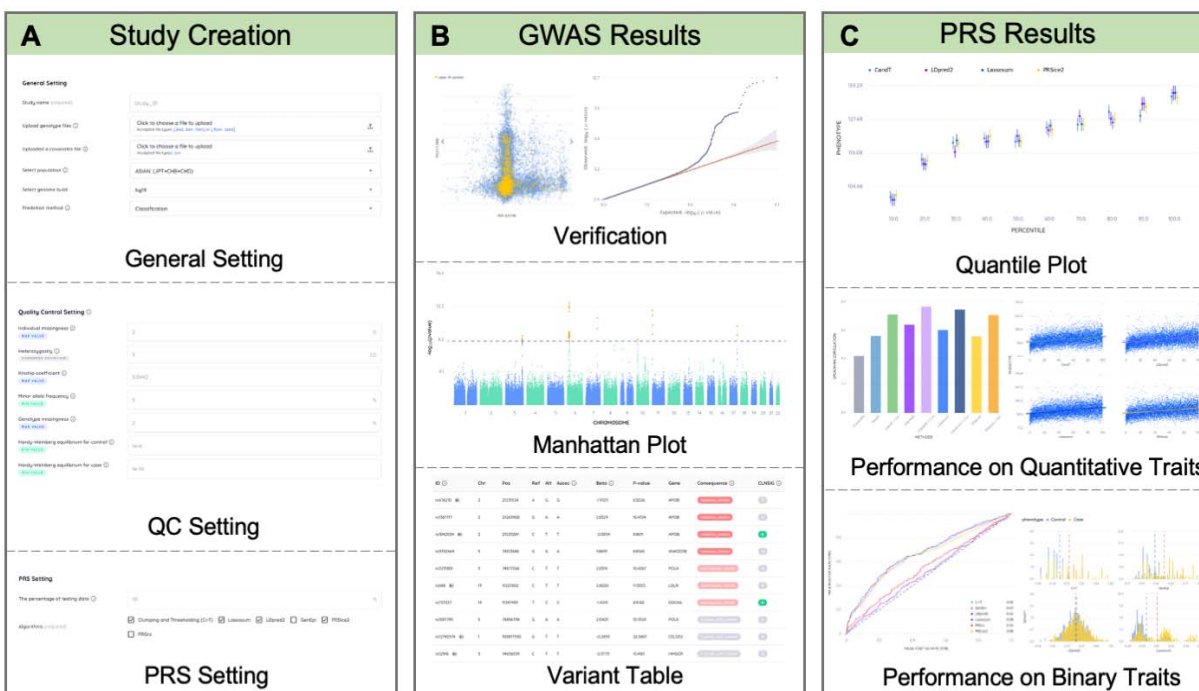629        hundreds of novel loci for 36 quantitative traits. Preprint at

630      https://doi.org/10.1101/2021.04.12.21255236.

631  45.  Sakaue, S. *et al.* A cross-population atlas of genetic associations for 220 human phenotypes. *Nat.*

632      *Genet.* **53**, 1415–1424 (2021).

633  46.  Breijyeh, Z. & Karaman, R. Comprehensive Review on Alzheimer's Disease: Causes and Treatment.

634      *Molecules* **25**, (2020).

635  47.  Wightman, D. P. *et al.* A genome-wide association study with 1,126,563 individuals identifies new

636      risk loci for Alzheimer's disease. *Nat. Genet.* **53**, 1276–1282 (2021).

637  48.  Bellenguez, C. *et al.* New insights into the genetic etiology of Alzheimer's disease and related

638      dementias. *Nat. Genet.* **54**, 412–436 (2022).

639  49.  de Rojas, I. *et al.* Common variants in Alzheimer's disease and risk stratification by polygenic risk

640      scores. *Nat. Commun.* **12**, 3417 (2021).

641  50.  Simmons, C. R., Zou, F., Younkin, S. G. & Estus, S. Evaluation of the global association between

642      cholesterol-associated polymorphisms and Alzheimer's disease suggests a role for rs3846662 and

643      HMGCR splicing in disease risk. *Mol. Neurodegener.* **6**, 62 (2011).

644  51.  Yashin, A. I. *et al.* Roles of interacting stress-related genes in lifespan regulation: insights for

645      translating experimental findings to humans. *J Transl Genet Genom* **5**, 357–379 (2021).

646  52.  Rodrigo, L. M. & Nyholt, D. R. Imputation and Reanalysis of ExomeChip Data Identifies Novel,

647      Conditional and Joint Genetic Effects on Parkinson's Disease Risk. *Genes* **12**, (2021).

648  53.  Zeng, P. & Zhou, X. Non-parametric genetic prediction of complex traits with latent Dirichlet

649      process regression models. *Nat. Commun.* **8**, 456 (2017).

650  54.  Yang, S. & Zhou, X. Accurate and Scalable Construction of Polygenic Scores in Large Biobank

651      Data Sets. *Am. J. Hum. Genet.* **106**, 679–693 (2020).

652  55.  Zhang, R. *et al.* Novel disease associations with schizophrenia genetic risk revealed in ~400,000 UK

653      Biobank participants. *Mol. Psychiatry* **27**, 1448–1454 (2022).

654  56.  Richardson, T. G., Harrison, S., Hemani, G. & Davey Smith, G. An atlas of polygenic risk score

655      associations to highlight putative causal relationships across the human phenome. *Elife* **8**, (2019).

656    57.  Sakaue, S. *et al.* Trans-biobank analysis with 676,000 individuals elucidates the association of

657         polygenic risk scores of complex traits with human lifespan. *Nat. Med.* **26**, 542–548 (2020).

658    58.  Shen, X. *et al.* A phenome-wide association and Mendelian Randomisation study of polygenic risk

659         for depression in UK Biobank. *Nat. Commun.* **11**, 2301 (2020).

660    59.  Sherry, S. T. dbSNP: the NCBI database of genetic variation. *Nucleic Acids Research* vol. 29 308–

661         311 Preprint at https://doi.org/10.1093/nar/29.1.308 (2001).

662    60.  Karczewski, K. J. *et al.* The mutational constraint spectrum quantified from variation in 141,456

663         humans. *Nature* **581**, 434–443 (2020).

664    61.  Safran, M. *et al.* The GeneCards Suite. in *Practical Guide to Life Science Databases* (eds.

665         Abugessaisa, I. & Kasukawa, T.) 27–56 (Springer Nature Singapore, 2021).

666    62.  Lex, A., Gehlenborg, N., Strobelt, H., Vuillemot, R. & Pfister, H. UpSet: Visualization of

667         Intersecting Sets. *IEEE Trans. Vis. Comput. Graph.* **20**, 1983–1992 (2014).

24

## Figures and Figure legends



**Figure 1. Analysis pipelines of PGSbuilder.** PGSbuilder performs GWAS and PRS analysis respectively on the input dataset. For the GWAS pipeline, PGSbuilder applies QC followed by GWAS on the whole input dataset. For the PRS pipeline, PGSbuilder splits the input dataset into training and test sets with the default ratio of 9:1 and applies QC on the training set. The training set is later split into base and target subsets with a ratio of 3:1, and the GWAS result is obtained from the base set. Combining the target set with the summary statistics derived from the base set, PGSbuilder builds PRS models based on different PRS methods. Alternatively, users could provide external summary statistics and the entire training set will be used to build the PRS model. Finally, the independent test set is used to evaluate the performance of the PRS model.

**Figure 2. PGSbuilder interface and visualizations.** (A) First of all, users can create a new study with customization, including general, QC, and PRS settings. (B) After analysis, GWAS results are composed of verification, including the PCA and Q-Q plots, and significant SNPs, including the Manhattan plot and variant table. (C) On the other hand, PRS results show a quantile plot for risk stratification and performance comparison for quantitative or binary traits.

## (A) SNP Information

| ID | Chromosome | Position | Reference | Alternative |
|---|---|---|---|---|
| rs440446 | 19 | 44905910 | C | G |

## (B) VEP Annotation

| Gene | | APOE |
|---|---|---|
| **Transcript** | | ENST00000434152 |
| **Ensembl** | **Biotype** | protein coding |
| | **Consequence** | missense variant |
| **ClinVar** | **CLNSIG** | benign |
| | **CLNREVSTAT** | ★★★★ no assertion |
| **Allele Frequency** | **Taiwan Biobank** | 39% |
| | **1000 Genome** | 64.4% All |

## (C) External Websites

- dbSNP
- gnomAD
- GWAS Catalog
- NCBI Gene
- GeneCards

## (D) pubmedKB

**Summary**

PubmedKB Summary

There is only one paper mentions rs440446 and leprosy. The paper title is Phenotypic severity in a family with MEND syndrome is directly associated with the accumulation of potentially functional variants of cholesterol homeostasis genes..

It indicates that "The SNPs rs440446 and rs429358 were associated with leprosy when we compared the patients with leprosy with the healthy con- trols, but the significance did not survive Bonferroni correc- tion (Table 3)."

**Literatures**

1 .A pleiotropic effect of the APOE gene: association of APOE polymorphisms with multibacillary leprosy in Han Chinese from Southwest China.

PMID: 28977675

Wang D, Zhang DF, Li GD, Bi R, Fan Y, Wu Y, Yu XF, Long H, Li YY, Yao YG • Br. J. Dermatol. • 2017 • Impact Factor: 6.71399999
Results excerpt: The SNPs rs440446 and rs429358 were associated with leprosy when we compared the patients with leprosy with the healthy con- trols, but the significance did not survive Bonferroni correc- tion (Table 3). As the leprosy-associated SNPs (rs405509 and rs439401) identified in the Yuxi sample were not covered by target sequencing, we checked the LD pattern of the five SNPs (rs405509,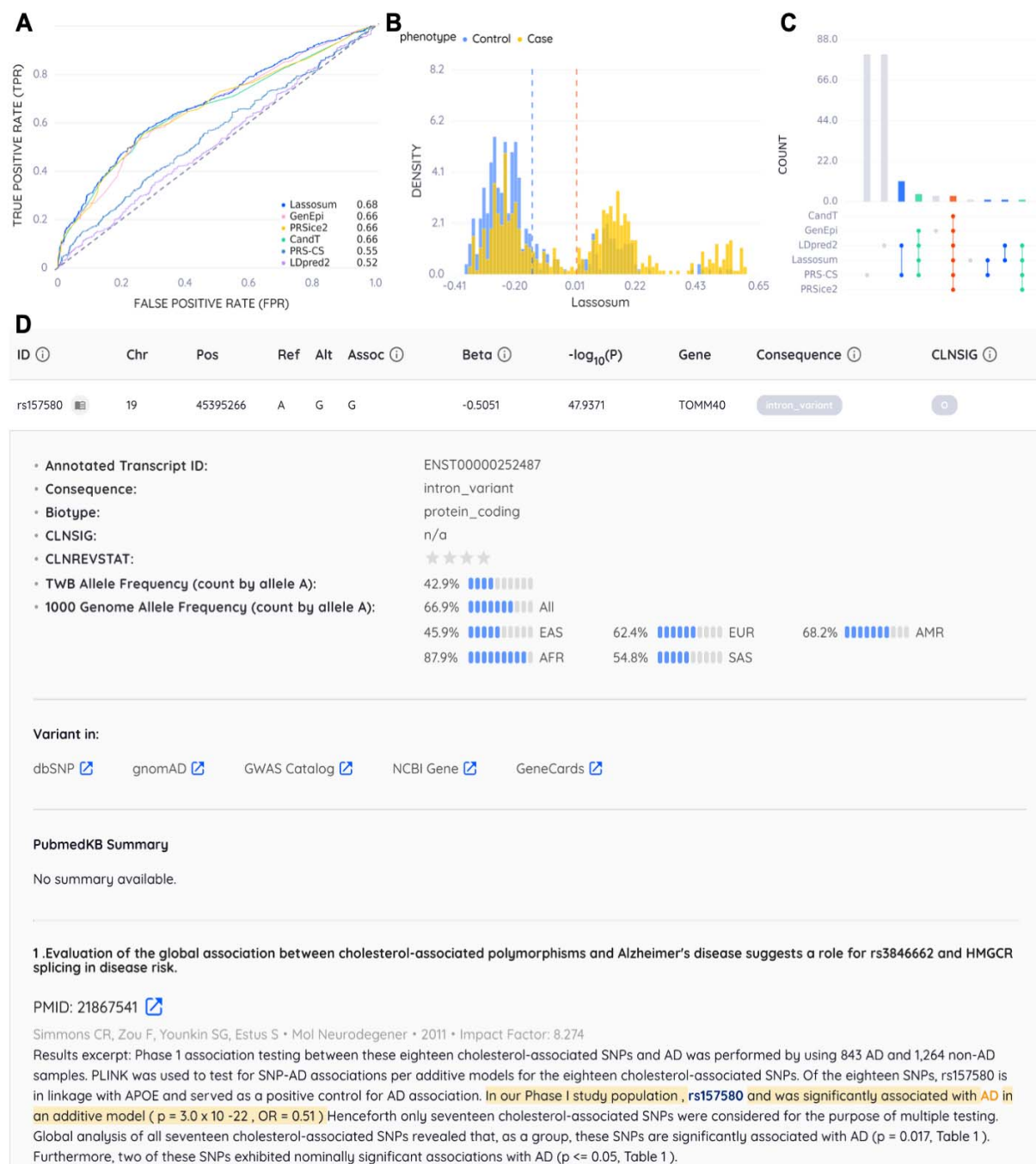 rs439401, rs440446, rs429358 and rs7412) in the CHB population from the 1000 Genomes dataset. 43 We found that rs440446 was linked with rs439401 (r 2 = 0 86). The SNP rs440446 in the Wenshan sample had an OR in the same direc- tion ( OR 1 193 , 95 % confidence interval 1 045 - 1 363 , P = 0 010 ) as that of rs439401 in the Yuxi sample , providing further evidence for the association of APOE SNPs with leprosy The leprosy-risk single-nucleotide polymorphisms affected APOE expression in human tissue Next, we tested the eQTL effect of the five SNPs (rs405509, rs769450, rs429358, rs7412 and rs439401) genotyped in the of Southwest China a (c) 2017 British Association of Dermatologists British Journal of Dermatology (2018) 178, pp931 - 939 APOE gene and leprosy susceptibility, D. Wang et al. 935 Yuxi sample and the three common SNPs (rs440446, rs429358 and rs7412) identified in the Wenshan sample in human blood and skin tissue using the dataset from the GTEx project. 35 We found that the two leprosy-risk SNPs in the Yuxi sample were significant cis eQTL (rs405509, P = 3 80 9 10 6 , Fig. 1a; rs439401, P = 3 40 9 10 12 , Fig. 1b) in skin tissue.

688

**Figure 3**. **Example annotation result of SNP "rs440446" on PGSbuilder.** （A) There is the basic information and statistics (e.g. GWAS $P$-value) of the variant. (B) We apply different colors on consequence (the red one) and ClinVar significance (the green one) according to tables

27

692     provided by Ensembl and ClinVar, respectively,for a better presentation of SNP importance

693     level. The following block is the transcript ID, ClinVar significance, and allele frequency from

694     VEP. (C) We also provided links to external websites with more variant or gene information,

695     such as dbSNP[59], gnomAD[60], GWAS Catalog[2], and GeneCards[61]. (D) The block at the bottom is

696     the results from pubmedKB. The summary presents the sentence where the SNP and the

697     phenotype co-occur, and we show the paper snippet of odds ratio statistics.

698

**Figure 4. Results of LDL GWAS and PRS analyses.** (A) PCA plot of the first and second PCs. To view any deviation of PCs among the samples, values of the quantitative phenotype are separated into four quantiles. (B) Manhattan plot of −log10(P-value). GWAS is performed on autosomal SNPs, and SNPs with P-value <10-5 are colored in orange. (Source data in Table S3) (C) Bar plot of Spearman's correlation of each PRS model. Models derived from different methods with or without covariates (Cov) are demonstrated simultaneously. (Source data in Table S4) (D) Quantile plot for risk stratification. The "covariate-only (Cov)" model and "PRS + covariate" models are plotted to compare the usage of genetic factors. (Source data in Table S5).

**Figure 5.** R**esults of AD across different PRS methods.** (A) ROC curve of each PRS model on the test set. (Source data in Table S6) (B) Prediction distributions of the Lassosum PRS model for cases (yellow) and controls (blue). The dashed line represents the mean of each group. (C) UpSet plot[62] for the intersection of important SNPs derived from different PRS methods. The intersection, or the combination, of methods are presented as the matrix layout while the variant counts of each intersection are shown as the histogram. Different colors represent the number of

717    PRS methods. (corresponding output data in Table S7) (D) Annotations for SNP "rs157580". On

718    the top is the basic information and statistics of the variant. The following block is the transcript

719    ID, ClinVar significance and allele frequency from VEP[28]. In addition, we also provided links to

720    external websites with more variant information, such as dbSNP[59] and gnomAD[60]. The block in

721    the bottom is the results from pubmedKB[29] which highlights the odds ratio of AD in the presence

722    of this variant.