

1 **Prioritization of disease genes from GWAS using ensemble based positive-**
2 **unlabeled learning**

3

4

5

6 Nikita Kolosov^{1,3}, Mark J. Daly^{2,3,4,*}, Mykyta Artomov^{1,2,3,4,*}

7

8 ¹ – ITMO University, St. Petersburg, Russia

9 ² – Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston,

10 USA

11 ³ – Broad Institute, Cambridge, USA

12 ⁴ – Institute for Molecular Medicine Finland (FIMM), Helsinki, Finland

13 * - correspondence: artomov@broadinstitute.org, mjdaly@atgu.mgh.harvard.edu

14

15 **Authors declare no conflict of interests.**

16

17 **Abstract**

18 Major complication in understanding disease biology from GWAS arises from
19 inability to identify a complete set of causal genes. Integration of multiple omics data
20 sources could provide an important functional link between associated variants and
21 candidate genes. Machine-learning could take advantage of this variety of data and
22 provide a solution for prioritization of disease genes. Yet, classical positive-negative
23 classifiers impose strong limitations on the gene prioritization procedure, such as lack of
24 reliable non-causal genes for training.

25 Here, we developed a novel gene prioritization tool - Gene Prioritizer (GPrior). It is
26 an ensemble of five positive-unlabeled bagging classifiers, that treat all genes of unknown
27 relevance as an unlabeled set. GPrior selects an optimal combination of algorithms to
28 tune the model for each specific phenotype.

29 Altogether, GPrior fills an important niche of methods for GWAS data post-
30 processing, significantly improving the ability to pinpoint disease genes compared to
31 existing solutions.

32 **Introduction**

33
34 Despite the tens of thousands of genetic associations identified using GWAS to
35 date, the ultimate goal - informing and guiding therapeutic development - has been
36 achieved for only a few phenotypes. A major complication in understanding disease
37 biology from GWAS often arises from inability to directly identify disease genes ¹.
38 Therefore, additional post-GWAS analysis is needed to first, identify a variant that drives
39 the signal within the locus, and then to connect this variant to a gene.

40 Fine-mapping, based on a Bayesian framework, sets out to prioritize variants
41 within the locus and, ultimately, identify the disease-causing variant ^{2,3,4}. Fine-mapping
42 algorithms - FINEMAP ⁵, PAINTOR ⁶, fGWAS ⁷, SUSIE ⁸ etc. had significant impact on
43 the field and had successfully identified causal variants for multiple traits. Importantly,
44 fine-mapping is done independently for each locus and in its current configuration does
45 not take advantage of biological relatedness (e.g., same pathway membership) of genes
46 involved in a phenotype ⁹.

47 At the same time, identification of the disease gene linked to a disease-causing
48 variant presents a major and yet unresolved challenge. Most GWAS associations
49 implicate a set of correlated genetic variants, none of which alter the protein-coding
50 sequence of a gene and which often physically span or are near to multiple genes. Since
51 our knowledge of regulatory sequence patterns of the genome, the relevant cells, tissues
52 and developmental time points relevant to disease are all incomplete, it is currently the
53 case that the vast majority of GWAS 'hits' do not have a certain link to a gene. Though

54 data sets with which to infer functional annotation and gene expression are growing
55 rapidly in their utility.

56 Post-processing of the GWAS results with inclusion of functional information is a
57 promising direction on the road to identify disease genes. For example, Post-GWAS
58 Analysis Platform (POSTGAP¹⁰) uses GWAS summary statistics along with LD-structure
59 and external functional databases (GTEx¹¹, FANTOM5¹², RegulomeDB¹³ etc.) to
60 prioritize SNPs within the locus and narrow down the list of potential gene candidates.
61 Yet, the gene prioritization utility of POSTGAP is still in early development.

62 Altogether, fine-mapping, functional annotations and known biologic relatedness
63 across putative disease genes become valuable data sources for gene prioritization,
64 defined as evaluation of the likelihood of gene involvement in generating a disease
65 phenotype¹⁴. Machine-learning (ML) based prioritization could take an advantage of
66 these data sources and provide a solution for novel disease gene identification.

67 Typically, existing ML solutions use Positive-Negative (PN) classification strategy.
68 In PN classification per-gene probabilities are obtained by using known disease genes as
69 a positive (P) training set and unknown genes as a negative (N) training set^{15, 16, 17, 18}.
70 Such an approach suffers from contamination of a negative set by hidden positives (HP),
71 represented by yet undiscovered disease genes. Additionally, it is challenging to find
72 reliable negative examples (i.e., genes that with certainty do not contribute to the
73 development of a phenotype). Most biological databases do not store negative evidence
74 (e.g., absence of gene interaction), rather they provide only observed positive evidence.
75 As a result, PN-classifiers could suffer from high false-negative prediction rates and
76 biased quality metrics.

77 It is feasible to design a model, where a limited number of reliable positive
78 examples (likely causal genes) will be used along with the rest of genes without treating
79 the latter as reliably negatives. Positive-Unlabeled (PU) learning has been developed to
80 overcome limitations of PN-learning. PU-learning treats unknown examples as a mixture
81 of P and N, called unlabeled (U) set.

82 PU-learning was first proposed by Denis et al.²¹, and several algorithms have been
83 published since then^{19,20,21}. A particular class of PU methods - PU-bagging, showed the
84 best stability of the learning algorithm²². Specifically, Mordlet et al.²² proposed the
85 “bagging SVM” approach that took advantage of a limited number of positive examples
86 and significantly improved performance and stability of classification using a new
87 aggregation technique.

88 Nevertheless, a single ML algorithm cannot fit all complex phenotypes and highly
89 heterogeneous biological data. To overcome this, Yang et al.²³ introduced a concept of
90 integration for several PU learning classifiers into one workflow using ensemble
91 technique. This technique was only tested with a specific family of PU algorithms - two-
92 step methods, heuristic in nature and sensitive to the initial choice of negative examples
93²⁴, significantly limiting applicability to GWAS data. Two-step PU algorithms first

94 attempted to identify negative examples in the unlabeled set, and then train a model from
95 the positive, unlabeled and likely negative examples. However, directly learning to
96 discriminate P from U with estimation of optimal misclassification costs leads to better
97 results ²⁵.

98 Therefore, combination of different ML algorithms (kernel-based and tree-based)
99 along with PU bagging is a promising strategy for building a gene-prioritization model
100 suitable for a large number of complex phenotypes and high variety of data sources, that
101 is still lacking in the field.

102 Here, we propose a novel PU-learning based gene prioritization tool - Gene
103 Prioritizer (GPrior), intended for post fine-mapping usage of GWAS results. In GPrior we
104 implemented the ensemble of 5 different ML classifiers for PU-bagging with further
105 selection of the optimal combination of predictions. Our approach returns probability
106 scores for the whole provided set of genes based on similarity level with positive examples
107 and is complementary to any other gene prioritization tools and fine-mapping techniques.

108 We illustrate the utility of PU-learning and GPrior with common public ML quality
109 evaluation dataset with known ground truth and a series of case studies. Comparison with
110 popular methods (TOPPGENE ²⁶; Bagging SVM ²²; MAGMA ²⁷) confirmed significantly
111 higher quality of predictions returned by GPrior.

112 **Methods**

113 GPrior was designed for prioritizing disease-relevant genes given a matrix of gene-
114 level features and a set of reliably causal genes. We integrated multiple ML techniques
115 in a single tool and a data-driven framework to select the most appropriate algorithm (or
116 combination of algorithms) on a case-by-case basis.

117 The prioritization scheme includes two independent steps. First, each ML algorithm
118 is used for positive-unlabeled bagging and generation of predictions for each gene.
119 Second, the best-performing combination of predictions is generated. To ensure the
120 independence of steps, a set of true genes that is used for training is separated into two
121 parts – set of genes for training individual ML algorithms and an algorithm evaluation set
122 (**Figure 1, Sup. Figure S1**). The latter is used to evaluate the quality of predictions from
123 ML algorithms and select predictions that will contribute to the optimal combination.
124 Altogether, such an approach allows to combine multiple learning algorithms and achieve
125 previously unattainable for an individual algorithm performance.

126

127 Input and Features

128

129 In addition to the described above true gene sets needed for training, GPrior
130 requires a data matrix with rows representing genes and columns representing features.

131 GWAS summary statistics contains only variant information which needs to be
132 converted into gene level data. Initially, we filtered out likely non-associated variants with

133 the threshold determined on case-by-case basis, ensuring inclusion of the majority of
134 potentially causal genes (even if no significant association was observed in GWAS) into
135 the prioritization analysis. Otherwise, GWAS p -values were not a part of the prioritization
136 model.

137 Next, we used POSTGAP¹⁰, which takes advantage of LD structure and variant
138 functional annotations to assign potential gene candidates for each variant. Such
139 preprocessing of GWAS summary statistics yields a variant-based data matrix with
140 mappings to a list of candidate genes.

141 A major challenge in transforming variant-based data into gene-based data matrix
142 for GPrior is preservation of valuable information about variant annotations. We used a
143 transformation of variant-level features (e.g. functional annotations, GERP scores, etc)
144 into gene-level features using a method proposed by Lehne et al.²⁸ to obtain a gene-
145 based data matrix.

146 In addition, we used gene expression and gene interaction data that proved their
147 utility for the gene identification problem in previous works^{15, 29, 30}. Specifically, we used
148 the GTEx database to obtain median gene expression levels for 53 tissues and Reactome
149³¹ and UCSC (GeneSorter)³² databases were used for gene interaction data (**Sup. Table**
150 **S1**).

151 Additional functional features and predictions of other prioritization algorithms
152 could be included in the data matrix to be used for GPrior model to boost the performance
153 quality³³. GPrior could take as an input either the raw output of POSTGAP (variant-based
154 data matrix) or any gene-based data matrix provided by a user.

155 We kept the same set of features for the case studies to preserve the fairness of
156 performance comparisons for different phenotypes. Although, for each phenotype
157 features could be selected in concordance with phenotype-specific needs, for example,
158 relevant cell type expression data. Overall, GPrior is not bound to a pre-specified set of
159 features and could be used with any user-defined set of features to boost the trait-specific
160 performance.

161

162 Feature preprocessing

163

164 As for any ML approach, prior to algorithm execution, features should undergo
165 preprocessing procedure to equalize scales and eliminate potential performance biases.

166 In gene prioritization, raw data sets can potentially reach hundreds or thousands
167 of features. Along with a limited number of positive examples, this enormity can lead to
168 the “curse of dimensionality”³⁴. Hypothetically, the more features are available in the data
169 the more accurate result should be expected. However, greater number of features leads
170 to the exponential growth of the training examples amount, required to cover the sparse
171 feature space and achieve acceptable prediction quality. In real-life applications, the
172 number of positive examples is limited, therefore, conventionally this problem is solved

173 by clustering raw features. GPrior uses agglomerative feature clustering as a
174 dimensionality reduction technique to extract appropriate number of features and achieve
175 the highest performance (**Suppl.Methods, Input and Features**).

176

177 GPrior algorithm

178

179 GPrior consists of five PU Bagging ensembles, each of them uses a different
180 classification algorithm: Logistic Regression (LR), Support-Vector Machine (SVM),
181 Decision Tree (DT), Random Forest (RF), Adaptive boosting (AB) (**Figure 1, Sup.**
182 **Methods, GPrior Algorithm**).

183 Each positive-unlabeled bagging procedure starts with a creation of a training set
184 with all positive (P) instances, treated as Positives, and a random subsample of
185 unlabeled, size of P, treated as negatives. Resulting in the size of a bootstrap sample
186 being equal to P. This way, on each iteration only a small portion of unlabeled instances
187 is treated as negatives, minimizing false negative error rate (**Sup.Figure S2**). Each
188 learning method is then fine-tuned by finding an optimal set of parameters (**Sup. Table**
189 **S2**).

190 After training and tuning, individual classifiers generate a probability score for each
191 gene to belong to the positive class. All the steps are repeated T times. Per-gene
192 probabilities are obtained by dividing the sum of all predictions by the number of times
193 each gene was sampled from the unlabeled set. All the predictions are averaged and
194 stored as a final PU Bagging result. All the steps are repeated for each classification
195 algorithm.

196 Next, GPrior selects the combination of predictions that shows the best
197 performance in prioritizing true genes from an independent algorithm evaluation set.
198 Since “true negative” data points that falsely were classified as positives could not be
199 identified in PU-data, any metric depending on false positives could not be applied for
200 quality evaluation. We used *PU-score* as a formal quality metric suitable for positive-
201 unlabeled data classes^{25, 35} (**Sup. Methods**).

202 All combinations from individual predictions are evaluated using *PU-score*
203 calculated for algorithm evaluation set and the best performing composition of methods
204 is then selected as the best fitting for a given phenotype. Selected combination is used to
205 return a vector of probabilities corresponding to the genes in the input matrix.

206

207 **Results**

208 Performance evaluation

209

210 The number of known true positive and negative data points is a critical information
211 parameter for gene prioritization. It is challenging to estimate both the number of genes

212 involved in a complex trait and the number of genes confidently irrelevant for the disease.
213 Height, as a classic example of a highly polygenic trait, shows effect size for a median
214 GWAS SNP in the genome about 10% of that for genome-wide significant hits. This
215 suggests that in current GWAS for height, significant associations are observed only for
216 a small proportion of true positive data points, while many others are yet to be confirmed.
217 However, additional alleles at known genes are a likely source of much of what is missing
218 so this does not easily translate into an estimate of how many relevant genes are
219 implicated from a GWAS. As shown in recent works, the genetic architecture of height is
220 broadly similar to that of a wide variety of other quantitative traits and diseases ranging
221 from diabetes or autoimmune diseases to BMI or cholesterol levels ³⁶ – for all of which
222 the evidence suggests many more positive genes exist in the ‘currently not associated’
223 gene set. Additionally, only some of the genome-wide significant loci were mapped to a
224 single gene, even further reducing the number of known true genes suitable for training
225 a model. Therefore, it is reasonable to assume that gene prioritization algorithms should
226 expect to be trained using only a small fraction of all disease-relevant genes.

227 Due to inability to obtain a GWAS dataset with confidently identified complete sets
228 of the disease genes and disease irrelevant genes, for initial GPrior performance
229 evaluation we used a public benchmark dataset. Breast Cancer Wisconsin (Diagnostic)
230 Data Set - a popular public dataset for machine learning tools benchmarking was used to
231 compare performance of GPrior to conventional PN-learning (biased SVM ³⁷) and single
232 PU-learning algorithm - bagging-SVM ²² (**Suppl.Methods, GPrior Performance**
233 **Evaluation**). The data with known true positives and true negatives enables calculation
234 of fair prediction error rates which is impossible for real life data with yet undiscovered
235 true data points. While this dataset is not related to the gene prioritization problem, it
236 clearly illustrates the benefits of ensemble PU-learning in case of only a small fraction of
237 known true instances to be used for training.

238 The dataset contained 569 samples and 30 features (**Figure 2A**). Initially, we
239 assumed that only 4% of true positive class samples (malignant tumors in the Breast
240 Cancer Wisconsin data), are known and available for training. Since GPrior requires two
241 independent true sets for training, the original known positive class was broken down into
242 two parts – 5 samples were used for training and 3 for algorithm evaluation set. For all
243 other methods the whole known positive class was used as a single batch.

244 We performed PCA to highlight how data classes are recognized in PN and PU-
245 learning (**Figure 2B**). PN-learning treats all instances as true negative, except those used
246 for training (known positives). This way, PN-learning attempts to identify samples falsely
247 classified to be true negative. In opposite, PU-learning treats all instances not used for
248 training as unlabeled and is, therefore, free from assumption that true negative class
249 exists. GPrior generated predictions using each algorithm and algorithm evaluation set
250 was used to estimate *PU*-scores (**Suppl. Figure S3A**). All combinations of 5 algorithms
251 available in GPrior were tested and the combination with the best *PU*-score estimated

252 from the algorithm evaluation set was selected to return the results (**Suppl. Figure**
253 **S3B,C**).

254 Precision-recall analysis was used to determine optimal decision threshold for
255 each method and ensure that best possible performance was extracted from each
256 algorithm (**Figure 2C,D**). As a result, GPrior carefully identified corresponding sample
257 classes ($F1 = 0.916$, **Suppl.Figure S4A, B**), while PN-learning fails to achieve similar
258 performance ($F1=0.551$, **Suppl.Figure S4A, C**).

259 Testing all 3 methods - biased SVM, bagging SVM and GPrior for different fractions
260 of known positives (100 simulations) results in superior performance of the latter (**Figure**
261 **2E**). Notably, PN-learning starts to behave equally well compared to PU-learning only if
262 more than 30% of the true positive data points are already known and used for training.

263 For the benchmark dataset, with pre-specified true positives it is possible to
264 compare $F1$ -score and PU -score metrics. Both metrics show similar results, justifying
265 further usage of PU -score for GPrior performance evaluation in case studies, where $F1$ -
266 score could not be computed (**Figure S4D, E**).

267 Additionally, we tried to fix contamination of unlabeled-class with hidden positives
268 and change only the number of KP without changing the percentage of contamination to
269 make the setup even more fair for PN and PU learning. Despite this, we obtained better
270 results from GPrior (**Figure S5**).

271 To illustrate utility and performance quality of GPrior we performed several case
272 studies using GWAS results for several phenotypes: inflammatory bowel disease (IBD),
273 educational attainment (EA), coronary artery disease (CAD) and schizophrenia (SCZ).

274 Case study 1: Inflammatory bowel disease (IBD)

275
276 We used GPrior and summary statistics from Huang et al. ³⁸, to construct gene
277 prioritization for IBD. Summary statistics was preprocessed to obtain a data matrix with
278 2,166 gene candidates found in loci with original p -value $< 10^{-8}$. A list of 31 genes with
279 known evidence to be likely causal for IBD was used as a positive training set. Algorithm
280 evaluation set consisted of 14 genes reported in monogenic loci with p -value $< 10^{-10}$ found
281 in GWAS catalog. Independent validation set used only for performance evaluation
282 included 51 genes found within monogenic loci with p -values falling in range $10^{-10} - 10^{-8}$
283 (**Figure 3A, Suppl. Table S3**).

284 We generated gene priorities using GPrior (**Sup. Table S4**) and a set of methods
285 for comparison - a single PU-learning (bagging-SVM ²²), PN-learning (weighted LR ²⁵)
286 and TOPPGENE ²⁶. While GPrior implies two training steps and usage of two training
287 gene sets – true gene set and algorithm evaluation set, for other methods we used a
288 union of the two gene sets for training.

289 Next, we compared the performance quality of the methods. PU -score is a formal
290 quality metric for a ML-based classifier, rather than for a prioritization itself, and it depends

291 on the decision threshold used to assign classes to the instances. Gene prioritization
292 implies only a construction of the ranked list of genes, but not the classification of the
293 genes into “disease” and “non-disease”. Yet, we evaluated the maximal possible
294 performance of the methods in the classification problem. We used an independent
295 validation gene set to estimate *PU*-scores for all possible decision thresholds (fraction of
296 positive predictions made by the classifier) and GPrior has significantly greater maximal
297 *PU*-score compared to others (**Figure 3B**).

298 To evaluate prioritization quality, we estimated cumulative gains. Gain chart shows
299 enrichment of the genes from the validation set at the top of the ranked list of predictions,
300 that is the sharper is the growth of gain in the beginning of the chart – the more enriched
301 are correct predictions at the top of the predictions list (**Figure 3C**).

302 Since original GWAS summary statistics was preprocessed to include only variants
303 with p -value $< 10^{-8}$, all 2,166 genes in the data matrix are found in or in a proximity of
304 significantly associated loci. GPrior does not use association strength or DNA location
305 information for gene prioritization. Yet, genes from the validation set are significantly
306 prioritized over the non-relevant neighbors (Mann-Whitney, one-sided, p -value = 2.5×10^{-13} ,
307 **Figure 3D**).

308 We evaluated non-randomness of the predictions, by estimating enrichment of the
309 validation set genes at the top of the ranked list produced by GPrior (permutation p -value
310 $< 1 \times 10^{-6}$; **Figure 3E**).

311 Treatment of all genes from the validation set as a finite set of disease genes,
312 implies that all genes that are not in the validation set are true negatives. In case of
313 polygenic traits, this is most likely a false assumption which would lead to an
314 underestimation of the true value of the area under the *ROC*-curve (*ROC AUC*). Thus
315 *AUC* values will illustrate only approximate quality measurement. In such settings, GPrior
316 demonstrated the most efficient predictive power out of all tools ($AUC = 0.8$, **Suppl. Table**
317 **S5**).

318

319 Case study 2: Educational attainment

320

321 We performed a control experiment to demonstrate that GPrior predictions are
322 disease specific and are driven by underlying biological similarities for disease related
323 genes. We considered two phenotypes with likely very modest overlap in underlying
324 biological causes – IBD and educational attainment (EA). We hypothesized that usage of
325 the training gene set fitted for IBD should fail to predict genes for EA.

326 GWAS summary statistics from Lee et al³⁹ for EA was preprocessed to obtain the
327 data matrix of candidate genes ($N=10,638$) and features (**Methods**). To eliminate
328 potential bias in the size of training sets for the two phenotypes, we used for GPrior
329 training only 18 genes (12 for ML training and 6 for algorithm selection) from the IBD
330 training gene set that were also found in EA GWAS loci with p -value $< 10^{-6}$. As a validation

331 set for IBD we used the original IBD validation genes (N=51), for EA we used 381 genes
332 found in monogenic loci from GWAS catalog EA results (**Sup. Table S6, Sup. Methods**).

333 Usage of appropriate training set for IBD resulted in significant enrichment
334 (permutation $P < 10^{-6}$) of validation set genes in the top predicted genes (**Sup. Figure**
335 **S6A**). Predictions based on the same list of training genes were constructed for EA and
336 demonstrated no enrichment of the EA-specific validation gene set (permutation p -value
337 p -value = 0.12, **Sup. Figure S6B**). Yet, usage of the EA-specific training gene set of the
338 same size (N=18, **Sup. Methods**) led to successful prioritization of EA-specific genes
339 (permutation $P < 10^{-6}$, **Sup. Figure S6C**).

340 We expanded the training set for EA by inclusion of all genes found in monogenic
341 loci in GWAS Catalog with (N=119) and repeated prioritization analysis. As a result, we
342 obtained even more significant enrichment for the validation set and confirmed superior
343 performance quality of GPrior in comparison with other methods (**Sup. Figure S7, Sup.**
344 **Tables S7, S8**).

345 Finally, we estimated how strongly initial GWAS summary statistics preprocessing
346 contributed to the overall success of the prioritization. While POSTGAP, that was used
347 for mapping variants to a list of candidate genes, is not designed specifically for gene
348 prioritization, it reports a variant-to-gene mapping score based on the sum of values for
349 7 features (**Sup. Figure S8**). We used the maximal variant-to-gene score for each gene
350 to construct a ranked list of genes. First, we estimated the largest possible PU -score for
351 the model using only POSTGAP-based gene ranking for educational attainment data (PU -
352 score = 3.82). We limited feature space to exactly the same 7 features and constructed
353 gene prioritization using GPrior, resulting in ~10% increase in PU -score (PU -score = 4.1).
354 GPrior uses all available features, while POSTGAP score is limited by only initial 7
355 features, yet expanding feature space for GPrior yields significant increase leading to the
356 maximal PU -score of 4.84 (**Sup. Figure S7**), showing ~27% increase in quality.

357 Case study 3: Coronary artery disease (CAD)

358
359 We used the summary statistics of coronary artery disease GWAS of 34,541 CAD
360 cases and 261,984 controls from UK Biobank followed by replication in 88,192 cases and
361 162,544 controls⁴⁰. After preprocessing we obtained a gene-based data matrix with 2,794
362 gene candidates found in loci with original p -value $< 10^{-8}$.

363 Recent review by Khera and Kathiresan⁴¹ was used to compile gene sets for
364 GPrior (**Figure 4A**). All genes with identified biological roles in any of the known disease
365 pathways were used for the training set (TS=18, AES=8). All other genes, implicated in
366 CAD, but with yet undiscovered molecular pathway membership became a validation set
367 (VS=37) (**Suppl. Table S9**).

368 Prioritization list obtained with GPrior (**Sup. Table S10**) has shown the best
369 accuracy with all quality metrics in comparison with other methods (**Figure 4B-E, Sup.**
370 **Table 11**).

371 Conclusively, using risk genes with known molecular pathway membership GPrior
372 successfully prioritizes genes with yet unknown biological contribution but confidently
373 implicated in the disease. Importantly, by further analyzing feature importance in the
374 prediction model it is possible to build testable biological hypotheses for novel genes
375 discovered in predictions.

376 Case study 4: Schizophrenia

377
378 We used GWAS Summary statistics from Pardiñas et al ⁴². This study used
379 genotypes of 105 318 individuals, 40,675 schizophrenia cases and 64,643 controls.

380 After preprocessing we obtained a gene-based data matrix with 3,831 gene
381 candidates found in loci with original p -value $< 10^{-6}$.

382 Training set was prepared using reported genes found in monogenic loci from
383 GWAS meta-analysis results ⁴². Training gene set for individual ML algorithms included
384 20 genes with p -values falling in range 10^{-44} - 10^{-14} , and algorithm evaluation set included
385 24 genes with p -values within 10^{-13} - 10^{-8} range. Validation set (VS) included 28 genes
386 and was obtained from the same study and included all genes from significant polygenic
387 loci (**Figure 5A, Sup. Table S12, S13**).

388 GPrior demonstrated superior results in comparison with other methods using all
389 quality metrics. GPrior achieved the highest PU score (9.64) and AUC (0.92) values. On
390 all the top intervals of the predictions list (1%,5%,15%,25%) GPrior showed the highest
391 enrichment of the validation set genes (**Figure 5B-E, Sup. Table S14**).

392 Conclusively, even for complex phenotypes with limited biological mechanism
393 knowledge, like schizophrenia, GPrior is well powered to detect the relevant signal of
394 biologic relatedness and prioritize likely disease genes.

395 **MAGMA comparison**

396
397 We compared GPrior with a commonly used method that takes GWAS summary
398 statistics as an input and attempts to pinpoint likely disease genes - MAGMA ²⁷. It
399 computes gene-based p -value (mean association of SNPs in the gene, corrected for LD).
400 We ran MAGMA with default parameters and compared performance quality with GPrior.
401 As an output MAGMA returns a list of genes and corresponding p -values, which we used
402 to sort the list for prioritization purposes.

403 One of the challenges for a non-biased comparison was the relatively small
404 number of gene candidates in output from MAGMA. Therefore, we took the same as

405 reported in MAGMA number of top genes from GPrior results to compare equal number
406 of gene predictions.

407 GPrior demonstrated enrichment of top ranked predictions for validation sets for
408 all phenotypes – EA (p -value= 9×10^{-3}), schizophrenia (p -value= 7×10^{-4}), CAD (p -
409 value= 3×10^{-3}) and IBD (p -value=0.019). MAGMA produced significantly enriched
410 predictions only for CAD (p -value= 9×10^{-3}) (**Suppl. Figure S9**).

411 Conclusively, GPrior demonstrated best performance out of all evaluated
412 approaches for gene prioritization in multiple settings and for various phenotypes.

413 **4. Discussion**

414
415 A large number of GWAS studies performed to date provide an invaluable source
416 of information for generating biological hypotheses for disease causes. The majority of
417 these studies greatly benefited from fine-mapping that implicated a limited number of
418 gene candidates. However, for highly polygenic phenotypes like schizophrenia, known
419 genes represent only a tiny segment of the disease biology.

420 The challenge of mapping “variants to function” could greatly benefit from machine
421 learning approaches. Especially, for those phenotypes for which a strictly genetic fine-
422 mapping approach has had limited success in conclusively identifying risk genes. As we
423 illustrate, conventional positive-negative machine learning approaches require a
424 substantial fraction of already known disease genes to achieve sufficient prioritization
425 quality for novel candidates. Additionally, it is nearly impossible at this point to confidently
426 state that a gene is not involved in a disease, therefore, directly assuming “negative”
427 examples for training is fated to include false negatives in a training set, further reducing
428 prediction quality.

429 Instead, we provide a tool that uses positive-unlabeled learning and requires
430 confidence in selection of only positive instances for training. Such genes are relatively
431 easy to identify based on association significance, previously reported functional studies,
432 etc. Importantly, PU-learning performs well even when the training set is quite small.

433 Additional challenge for a single-method-based solution is presented by phenotype
434 complexity. Phenotypes may present significantly different genetic architectures or
435 impose certain limitations on the set of available data sources; therefore, it is unlikely that
436 a single technique will be suitable for gene prioritization in all of them. We provide a
437 software package for gene prioritization – GPrior that takes advantage of the ensemble
438 of PU-learning techniques. Such approach overcomes unresolved challenges of PN-
439 learning and issues arising from phenotype complexity. In GPrior, two key steps of the
440 model training: PU-classifiers training and selection of optimal classifiers combination are
441 performed using two independent gene sets. The two-step strategy ensures independent
442 quality assessment for all classifiers and unbiased selection of the optimal prioritization
443 method, as well as delivering optimal prioritization results for the specific phenotype.

444 GPrior can be utilized with many sources of functional data. Data types used in our
445 case studies – tissue expression levels, Reactome pathway data and others represent
446 only a small part of possibilities. Each phenotype study would significantly benefit from
447 inclusion of additional features, such as – single cell expression data, specific protein-
448 protein interactions, gene conservation metrics (pLI, LOEUF) and others. In our study we
449 have not selected features to be specific to each phenotype, therefore, users can expect
450 to see even higher performance in case of thorough feature selection. Additionally, GPrior
451 can be straightforwardly integrated with conventional fine-mapping tools. One of the
452 limiting steps in our GWAS processing scheme was naïve selection of gene candidates
453 from each locus. More sophisticated preprocessing of the raw GWAS summary statistics
454 with methods such as SuSie or FINEMAP to improve variant-to-gene mapping could
455 significantly aid variant-level to gene-level features transformation. Finally, we used a
456 relatively conservative set of features for gene annotations, which could be significantly
457 expanded with phenotype specific annotations.

458 Altogether, GPrior fills an important and currently underdeveloped niche of
459 methods for GWAS data post-processing, significantly improving the ability to pinpoint
460 disease genes compared to existing solutions.

461

462 **Code availability**

463

464 <https://github.com/faramer86/GPrior>

465

466 **Acknowledgements**

467 Authors would like to thank Dr. Alexey Sergushichev (ITMO University) and Dr. Maxim
468 Artyomov (Washington University in St. Louis) for helpful discussions.

469

470 **References**

471

- 472 1. Ding, K. & Kullo, I. J. Methods for the selection of tagging SNPs: a comparison of tagging
473 efficiency and performance. *Eur. J. Hum. Genet.* **15**, 228–236 (2007).
- 474 2. Foulkes, A. S. *Applied Statistical Genetics with R: For Population-based Association*
475 *Studies*. (Springer Science & Business Media, 2009).
- 476 3. Spain, S. L. & Barrett, J. C. Strategies for fine-mapping complex traits. *Hum. Mol. Genet.*
477 **24**, R111–9 (2015).
- 478 4. Stephens, M. & Balding, D. J. Bayesian statistical methods for genetic association studies.
479 *Nature Reviews Genetics* vol. 10 681–690 (2009).
- 480 5. Benner, C. *et al.* FINEMAP: efficient variable selection using summary data from genome-
481 wide association studies. *Bioinformatics* **32**, 1493–1501 (2016).
- 482 6. Kichaev, G. *et al.* Integrating Functional Data to Prioritize Causal Variants in Statistical
483 Fine-Mapping Studies. *PLoS Genetics* vol. 10 e1004722 (2014).
- 484 7. Pickrell, J. K. Joint analysis of functional genomic data and genome-wide association
485 studies of 18 human traits. *Am. J. Hum. Genet.* **94**, 559–573 (2014).
- 486 8. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable

- 487 selection in regression, with application to genetic fine-mapping. doi:10.1101/501114.
488 9. Rossin, E. J. *et al.* Proteins encoded in genomic regions associated with immune-mediated
489 disease physically interact and suggest underlying biology. *PLoS Genet.* **7**, e1001273
490 (2011).
491 10. Peat, G. *et al.* The Open Targets Post-GWAS analysis pipeline. *Bioinformatics* (2020)
492 doi:10.1093/bioinformatics/btaa020.
493 11. Consortium, G. & GTEx Consortium. Erratum: Genetic effects on gene expression across
494 human tissues. *Nature* vol. 553 530–530 (2018).
495 12. Andersson, R. *et al.* An atlas of active enhancers across human cell types and tissues.
496 *Nature* **507**, 455–461 (2014).
497 13. Boyle, A. P. *et al.* Annotation of functional variation in personal genomes using
498 RegulomeDB. *Genome Res.* **22**, 1790–1797 (2012).
499 14. Bromberg, Y. Chapter 15: disease gene prioritization. *PLoS Comput. Biol.* **9**, e1002902
500 (2013).
501 15. Adie, E. A., Adams, R. R., Evans, K. L., Porteous, D. J. & Pickard, B. S. Speeding disease
502 gene discovery by sequence based candidate prioritization. *BMC Bioinformatics* **6**, 55
503 (2005).
504 16. Xu, J. & Li, Y. Discovering disease-genes by topological features in human protein–protein
505 interaction network. *Bioinformatics* vol. 22 2800–2805 (2006).
506 17. Smalter, A., Lei, S. F. & Chen, X.-W. Human Disease–Gene Classification with Integrative
507 Sequence-Based and Topological Features of Protein-Protein Interaction Networks. *2007*
508 *IEEE International Conference on Bioinformatics and Biomedicine (BIBM 2007)* (2007)
509 doi:10.1109/bibm.2007.47.
510 18. Isakov, O., Dotan, I. & Ben-Shachar, S. Machine Learning–Based Gene Prioritization
511 Identifies Novel Candidate Risk Genes for Inflammatory Bowel Disease. *Inflammatory*
512 *Bowel Diseases* vol. 23 1516–1523 (2017).
513 19. Denis, F. PAC Learning from Positive Statistical Queries. *Lecture Notes in Computer*
514 *Science* 112–126 (1998) doi:10.1007/3-540-49730-7_9.
515 20. Sriphaew, K., Takamura, H. & Okumura, M. Cool Blog Classification from Positive and
516 Unlabeled Examples. *Advances in Knowledge Discovery and Data Mining* 62–73 (2009)
517 doi:10.1007/978-3-642-01307-2_9.
518 21. Zhang, B. & Zuo, W. Learning from Positive and Unlabeled Examples: A Survey. *2008*
519 *International Symposiums on Information Processing* (2008) doi:10.1109/isip.2008.79.
520 22. Mordelet, F. & Vert, J. A bagging SVM to learn from positive and unlabeled examples.
521 *Pattern Recognition Letters* vol. 37 201–209 (2014).
522 23. Yang, P., Li, X., Chua, H.-N., Kwok, C.-K. & Ng, S.-K. Ensemble positive unlabeled learning
523 for disease gene identification. *PLoS One* **9**, e97079 (2014).
524 24. Scott, C., Blanchard, G. Novelty detection: unlabeled data definitely help. *AISTATS '09*
525 *Proceedings of the Twelfth International Conference on Artificial Intelligence and Statistics,*
526 *JMLR: W&CP* **5**, 464–471 (2009).
527 25. Liu, B., Lee, W.S. Learning with positive and unlabeled examples using weighted logistic
528 regression. *Proceedings of the 20th international conference on machine learning* **20**, 448–
529 455 (2003).
530 26. Chen, J., Bardes, E. E., Aronow, B. J. & Jegga, A. G. ToppGene Suite for gene list
531 enrichment analysis and candidate gene prioritization. *Nucleic Acids Res.* **37**, W305–11
532 (2009).
533 27. de Leeuw, C. A., Mooij, J. M., Heskes, T. & Posthuma, D. MAGMA: generalized gene-set
534 analysis of GWAS data. *PLoS Comput. Biol.* **11**, e1004219 (2015).
535 28. Lehne, B., Lewis, C. M. & Schlitt, T. From SNPs to Genes: Disease Association at the
536 Gene Level. *PLoS ONE* vol. 6 e20133 (2011).
537 29. Stranger, B. E. *et al.* Population genomics of human gene expression. *Nat. Genet.* **39**,

- 538 1217–1224 (2007).
- 539 30. Ala, U. *et al.* Prediction of human disease genes by human-mouse conserved coexpression
540 analysis. *PLoS Comput. Biol.* **4**, e1000043 (2008).
- 541 31. Jassal, B. *et al.* The reactome pathway knowledgebase. *Nucleic Acids Res.* **48**, D498–
542 D503 (2020).
- 543 32. Kent, W. J. *et al.* Exploring relationships and mining data with the UCSC Gene Sorter.
544 *Genome Res.* **15**, 737–741 (2005).
- 545 33. Fine, R. S., Pers, T. H., Amariuta, T., Raychaudhuri, S. & Hirschhorn, J. N. Benchmark:
546 An Unbiased, Association-Data-Driven Strategy to Evaluate Gene Prioritization Algorithms.
547 *Am. J. Hum. Genet.* **104**, 1025–1039 (2019).
- 548 34. Altman, N. & Krzywinski, M. The curse(s) of dimensionality. *Nat. Methods* **15**, 399–400
549 (2018).
- 550 35. Claesen, M., De Smet, F., Suykens, J. A. K. & De Moor, B. A robust ensemble approach to
551 learn from positive and unlabeled data using SVM base models. *Neurocomputing* vol. 160
552 73–84 (2015).
- 553 36. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An Expanded View of Complex Traits: From
554 Polygenic to Omnigenic. *Cell* **169**, 1177–1186 (2017).
- 555 37. Liu, B., Dai, Y., Li, X., Lee, W. S. & Yu, P. S. Building text classifiers using positive and
556 unlabeled examples. *Third IEEE International Conference on Data Mining*
557 doi:10.1109/icdm.2003.1250918.
- 558 38. Huang, H. *et al.* Fine-mapping inflammatory bowel disease loci to single-variant resolution.
559 *Nature* **547**, 173–178 (2017).
- 560 39. Lee, J. J. *et al.* Gene discovery and polygenic prediction from a genome-wide association
561 study of educational attainment in 1.1 million individuals. *Nat. Genet.* **50**, 1112–1121
562 (2018).
- 563 40. van der Harst, P. & Verweij, N. Identification of 64 Novel Genetic Loci Provides an
564 Expanded View on the Genetic Architecture of Coronary Artery Disease. *Circ. Res.* **122**,
565 433–443 (2018).
- 566 41. Khera, A. V. & Kathiresan, S. Genetics of coronary artery disease: discovery, biology and
567 clinical translation. *Nat. Rev. Genet.* **18**, 331–344 (2017).
- 568 42. Pardiñas, A. F. *et al.* Common schizophrenia alleles are enriched in mutation-intolerant
569 genes and in regions under strong background selection. *Nat. Genet.* **50**, 381–389 (2018).
- 570

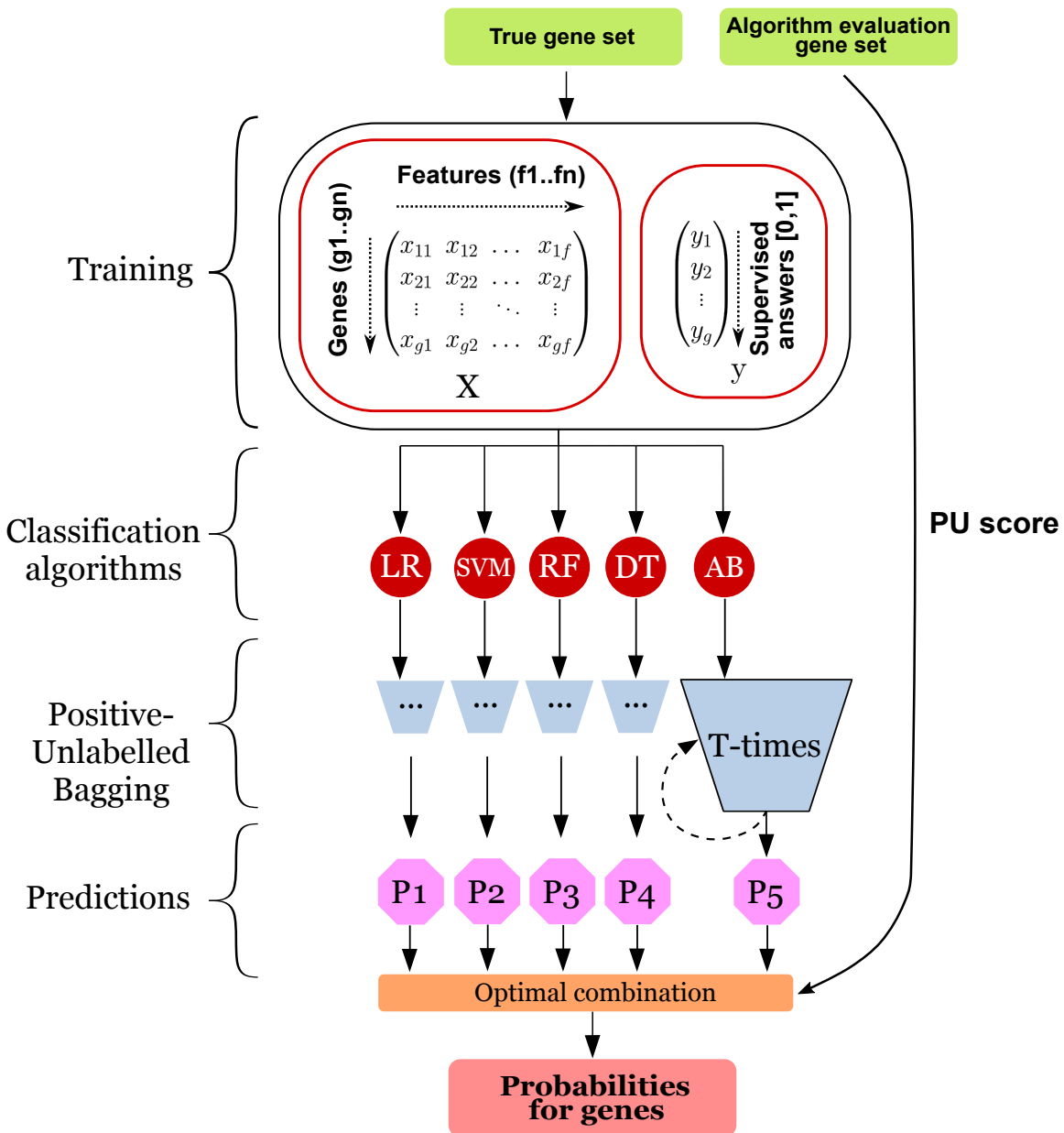


Figure 1. GPrior ensemble positive-unlabeled learning framework

Matrix of gene features along with vector of supervised answers is used to train 5 models using PU-bagging approach. Two independent gene sets are used for training – true set of genes for individual classification algorithms training and algorithm evaluation set of true genes for selecting the optimal combination of the predictions. Predictions are generated using positive-unlabeled bagging and further an optimal combination returning the largest *PU*-score is returned.

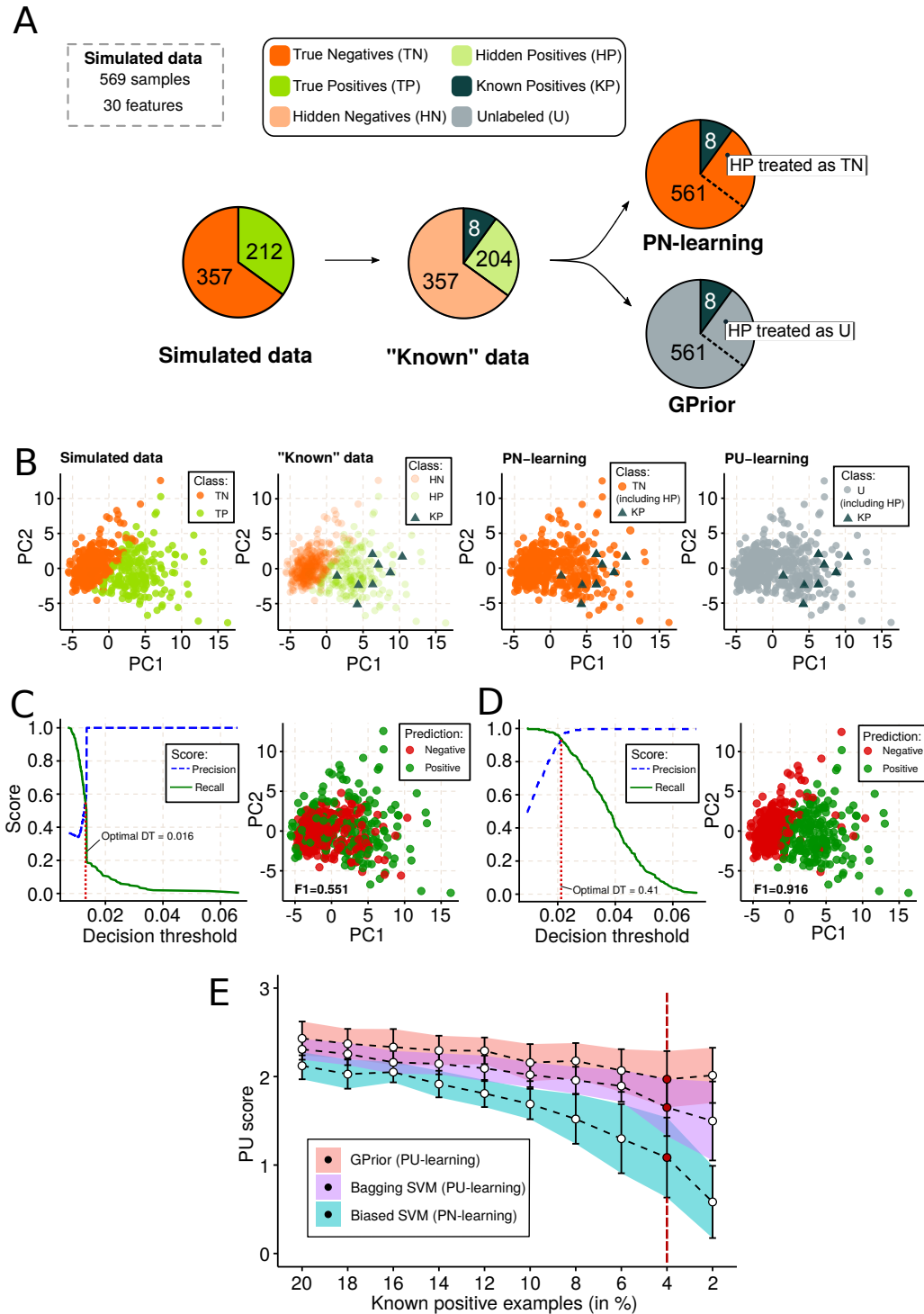


Figure 2. Benchmark dataset shows advantage of PU learning over PN learning in a scenario when few positive instances used for training.

(A) Dataset breakdown. Known data represents scenario when only a part of the true positive instances were discovered to date; (B) PCA of the simulated data with highlighted true instance classes; (C) Prediction results for PN-learning (biased SVM) approach; (D) Prediction results for GPrior; (E) Performance of PU and PN learning approaches with respect to a fraction of known positive data points.

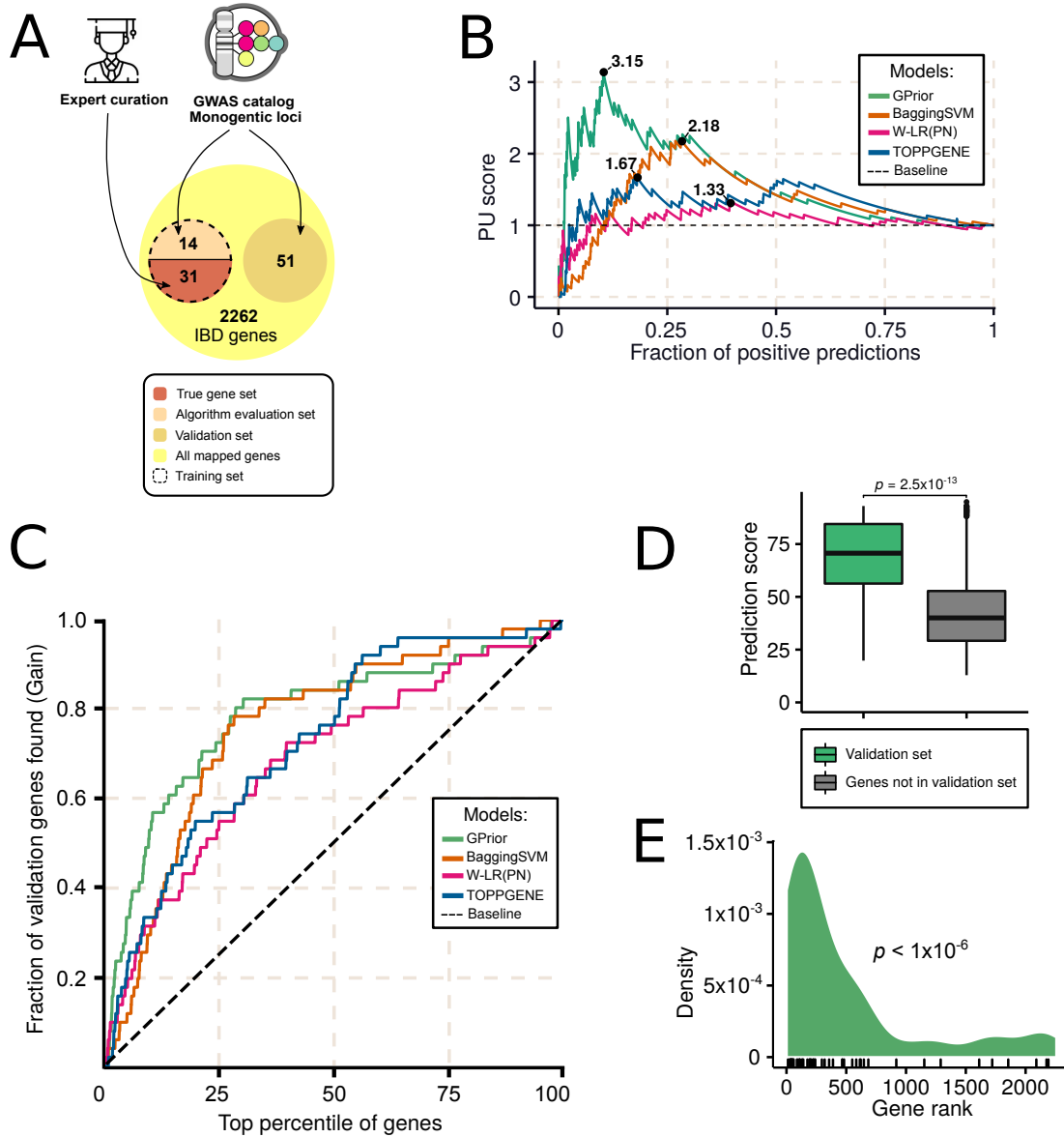


Figure 3. Gene prioritization for inflammatory bowel disease GWAS.

(A) Scheme for selection of training, algorithm evaluation and validation gene sets; **(B)** Classification quality comparison for GPrior, Bagging SVM and conventional PN-learning with weighted linear regression; **(C)** Cumulative gains curve shows better prioritization of true genes at the top of the candidate list using GPrior in comparison with other methods; **(D)** True genes from the independent validation gene set receive significantly higher scores than genes found within the same loci but not implicated in the disease; **(E)** Enrichment of true genes from independent validation gene set among top predictions from GPrior.

573
574

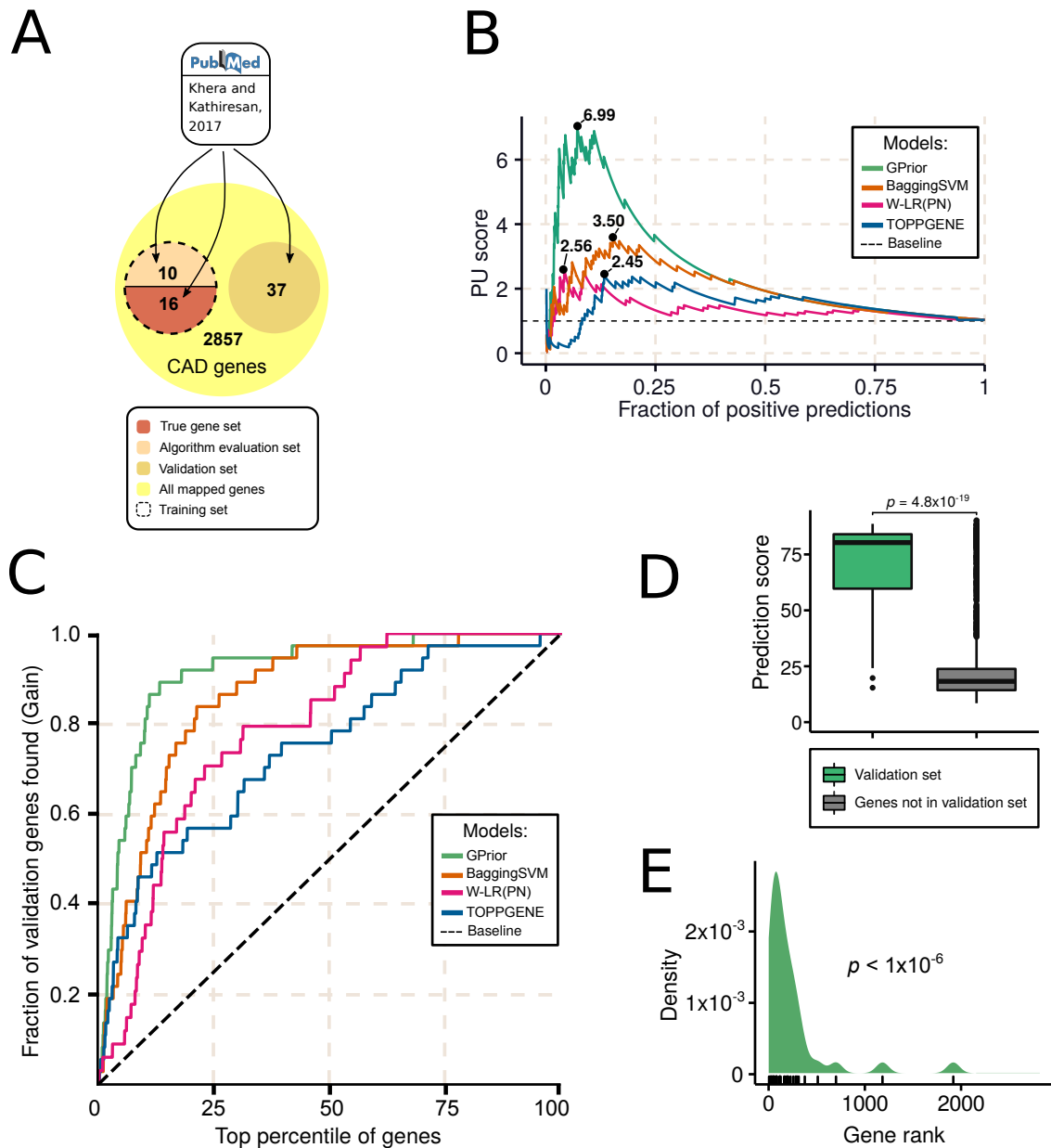


Figure 4. Gene prioritization for coronary artery disease GWAS.

(A) Scheme for selection of training, algorithm evaluation and validation gene sets; (B) Classification quality comparison for GPrior, Bagging SVM and conventional PN-learning with weighted linear regression; (C) Cumulative gains curve shows better prioritization of true genes at the top of the candidate list using GPrior in comparison with other methods; (D) True genes from the independent validation gene set receive significantly higher scores than genes found within the same loci but not implicated in the disease; (E) Enrichment of true genes from independent validation gene set among top predictions from GPrior.

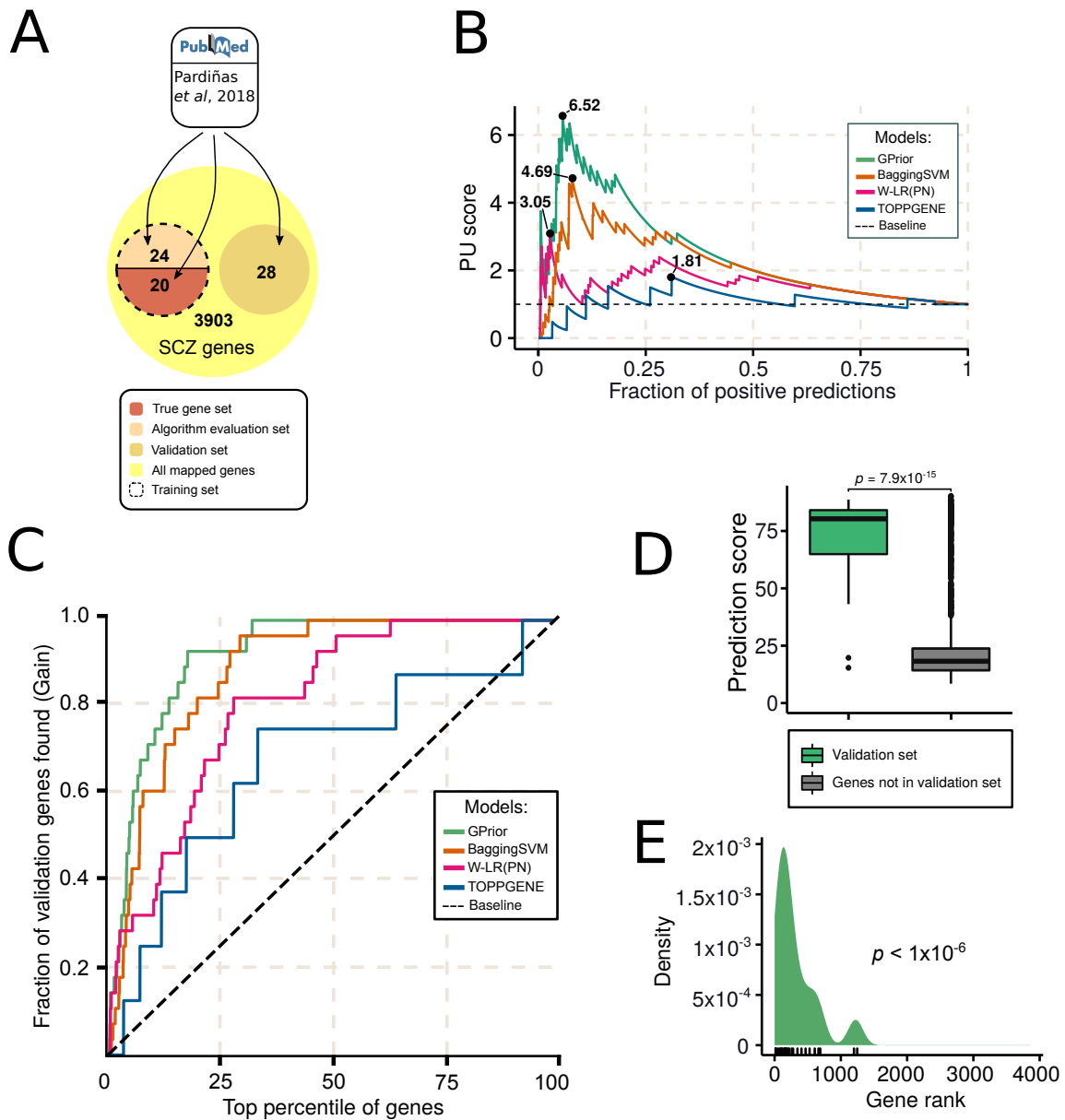


Figure 5. Gene prioritization for schizophrenia GWAS.

(A) Scheme for selection of training, algorithm evaluation and validation gene sets; (B) Classification quality comparison for GPrior, Bagging SVM and conventional PN-learning with weighted linear regression; (C) Cumulative gains curve shows better prioritization of true genes at the top of the candidate list using GPrior in comparison with other methods; (D) True genes from the independent validation gene set receive significantly higher scores than genes found within the same loci but not implicated in the disease; (E) Enrichment of true genes from independent validation gene set among top predictions from GPrior.