

1 Phenotype-specific information improves prediction of
2 functional impact for noncoding variants

3 Corneliu A. Bodea^{1,2,4}, Adele A. Mitchell¹, Heiko Runz^{1,5}, and Shamil R.

4 Sunyaev^{2,3,4,5}

5 ¹Department of Genetics and Pharmacogenomics, MRL, Boston, Massachusetts,
6 USA.

7 ²Division of Genetics, Department of Medicine, Brigham and Women's Hospital and
8 Harvard Medical School, Boston, Massachusetts, USA.

9 ³Department of Biomedical Informatics, Harvard Medical School, Boston,
10 Massachusetts, USA.

11 ⁴The Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA.

12 ⁵These authors jointly supervised this work. Correspondence should be addressed to
13 S.R.S. (ssunyaev@rics.bwh.harvard.edu).

14 **Abstract**

15 A myriad of noncoding genetic association signals are now awaiting the identification of
16 causal alleles and their functional interpretation. We introduce the novel computational frame-
17 work PINES (Phenotype-Informed Noncoding Element Scoring), which evaluates the functional
18 impact of noncoding variants by integrating diverse epigenetic annotations. A unique feature of
19 PINES is that it directs the analysis towards genomic annotations most relevant to phenotypes
20 of interest. We show that PINES identifies functional noncoding variation more accurately than
21 methods that do not use phenotype-specific knowledge. We apply PINES to fine map noncoding
22 alleles at GWAS loci across a range of diseases, and predict new causal risk alleles for Parkinson's
23 disease and inflammatory bowel disease. We also use PINES to confirm several high-penetrance
24 variants implicated in Mendelian traits, as well as variants residing within known enhancer
25 regions. PINES consistently identifies functional variants in fine mapping analyses, dissecting

26 pathogenic loci while avoiding the resource-intensive traditional fine mapping studies. Due to
27 its flexibility and ease of use through a dedicated web portal, PINES provides a powerful *in*
28 *silico* method to prioritize and fine map functional noncoding variants.

29 **1 Introduction**

30 A growing body of evidence suggests that DNA variants outside of protein-coding regions of the
31 genome (here termed noncoding variants) impact human phenotypes, including the risk for common
32 diseases. Many signals identified by genome-wide association studies (GWAS) point to regulatory
33 regions as key determinants of complex traits. Likewise, some human Mendelian phenotypes such as
34 hair, skin, and eye pigmentation are under the tight control of individual highly penetrant noncoding
35 alleles [1, 2]. Thanks to whole genome sequencing (WGS), our ability to uncover novel noncoding
36 alleles has increased substantially. However, studies of both common and rare phenotypes almost
37 never resolve findings in noncoding regions to individual functional causal SNPs [3]. Functional
38 prioritization of noncoding variants thus holds significant promise to assist in fine-mapping efforts
39 and identify genetic lesions underlying Mendelian diseases.

40 To better understand the architecture of the noncoding genome, several large-scale efforts, such
41 as ENCODE [4] and the Roadmap Epigenomics Project [5], have aimed to characterize the diverse
42 landscape of histone modifications and DNA accessibility based on a wide range of assays across 127
43 cell types. Other efforts such as the FANTOM5 project [6] have identified enhancer elements across
44 the genome, and computational tools such as TargetFinder [7] have been developed to link such
45 enhancers to the relevant gene promoters. Databases of chromatin interactions such as 4DGenome
46 [8] are also helpful in identifying potential regulatory elements. Yet, while most genomic regions are
47 by now annotated with a plethora of epigenetic features, the challenge remains to draw meaningful
48 conclusions from these annotations, especially since the data are highly dimensional and many
49 epigenetic features are correlated.

50 Several approaches have recently been introduced to prioritize potentially functional noncoding
51 variants and address the complexity of the annotation data in a principled manner. These include
52 Eigen [9], GWAVA [10], and CADD [11], as well as population genetics and conservation-inspired
53 models such as PhastCons [12] and INSIGHT [13]. One drawback of models such as GWAVA is
54 that they require a training dataset of both functional and non-functional examples. However,
55 high-quality and particularly experimentally validated training data sets for noncoding variants are
56 still very scarce and incomplete, thus limiting the prediction accuracy. A second drawback is that

57 current scoring methods globally merge annotations across many different cell types, which ignores
58 the observation that regulatory elements often operate in a very cell type-specific manner [14].

59 We hypothesized that predicting the functional relevance of noncoding variants could benefit
60 from taking into account the vast number of variants that have not yet been annotated as having
61 any function (background variants), and use these as a baseline to search for variants that devi-
62 ate significantly from this background. Such an approach falls under the category of PU (Positive
63 and Unlabeled example) learning, or one-class classification, and is a well-studied machine learning
64 method in the field of outlier detection [15, 16]. We further hypothesized that for identifying noncod-
65 ing variants of potential relevance to a specific phenotype, integration of prior biological knowledge
66 around the phenotype should increase prediction reliability. This prior knowledge can consist of
67 relevant cell types, but also of relevant genes or pathways. To the best of our knowledge, no current
68 method has the ability to integrate general phenotype-relevant knowledge into the scoring procedure
69 in such a principled and flexible manner.

70 Here we introduce the Phenotype-Informed Noncoding Element Scoring (PINES) framework.
71 PINES uses an unsupervised approach to systematically assess the functional significance of non-
72 coding SNVs and indels, and allows to customize the search towards annotations considered as of
73 highest relevance to a phenotype of interest. We apply PINES to *in silico* fine map noncoding
74 variants at GWAS loci and predict novel causative SNPs at assumed promoter and enhancer regions
75 for several common traits and diseases, including Parkinson’s disease, inflammatory bowel disease
76 (IBD), multiple sclerosis, and blood lipid levels among others. We further show that PINES can be
77 applied to assess a potential functional relevance of noncoding variants on Mendelian phenotypes.
78 With its ease of use via a customizable web server we expect PINES to become a valuable resource
79 for the interpretation of the noncoding human genome.

80 **2 Results**

81 **2.1 PINES integrates epigenetic information to score non-coding variants** 82 **in customizable queries**

83 Under a default setting, PINES integrates data from each of the 127 cell types analyzed in the
84 Roadmap and ENCODE projects, specifically information on histone modifications indicative of
85 promoters or enhancers (H3K4me1, H3K4me3, H3K27ac, H3K9ac), DNase I hypersensitive sites,
86 sequence constraint scores (GERP, SiPhy), and ChromHMM chromatin state segmentations [17].

87 These data are applied to compare individual or a set of user-defined SNPs relative to the genomic
88 background. Users may either apply PINES in a default mode, or pre-define sets of reference
89 variants that are known to be linked to phenotypes of interest, or tissues of specific relevance to
90 their respective scientific question. This allows that, if available, prior knowledge can be taken into
91 account during the subsequent scoring process (see Figure 1 and Methods for details).

92 An important feature of PINES is its ability to incorporate prior biological information into the
93 scoring procedure. When no prior knowledge is provided or available, equal weights are assigned
94 to all features to perform an undirected scoring of variants. When prior information is provided
95 (for instance a list of lead SNPs identified as significantly associated with a trait of interest through
96 GWAS), PINES searches for annotations that are enriched in the provided dataset in order to learn
97 which annotations are most relevant for the phenotype under consideration. Alternatively, users can
98 manually specify the most phenotypically relevant tissue, and all annotations relevant to that tissue
99 will be up-weighted by the scoring procedure. Prior information can additionally be specified in the
100 form of separate annotations. By relying on an angle-based distance from the vector of the maximal
101 possible annotation load in a de-correlated annotation space, PINES addresses both the correlation
102 structure as well as the high dimensionality of epigenetic annotation data (see also Methods).

103 **2.2 PINES predicts causal noncoding variants with high confidence**

104 We illustrate the potential of PINES to fine map noncoding genomic regions implicated in Mendelian
105 traits, as well as likely causal noncoding variants from GWAS loci. In order to test whether PINES
106 correctly identifies and prioritizes functional noncoding SNPs, we applied the algorithm to *in silico*
107 fine-map 20kb regions around seven noncoding alleles whose regulatory impact on a nearby gene
108 had been confirmed experimentally: rs12350739 [18], rs356168 [19], rs6801957 [20], rs12821256 [21],
109 rs138557689 [22], rs2473307 [23], and rs227727 [24] (Figure 2 and Supplementary Table 1). Such well-
110 studied noncoding variants, although still rare in the current literature, provide an optimal testbed
111 for scoring methods. Weighted PINES was used when genome-wide significant SNPs were identified
112 by GWAS for the trait under consideration (see Methods). In all cases, PINES scores peak around
113 the experimentally confirmed functional variants, and PINES assigns these causal variants the lowest
114 p-values. Notably, comparative analyses using these regions demonstrate that PINES outcompetes
115 GWAVA, CADD or Eigen in correctly identifying causal noncoding variants (see Supplementary
116 Figures 1-7).

117 **2.3 PINES detects signal at fine mapped GWAS risk loci**

118 We next aimed to evaluate the ability of PINES to predict causal noncoding variants at GWAS loci
119 across a range of conditions. As a first example for this, we extracted 3,625 candidate causal non-
120 coding variants reported by a large statistical fine mapping study spanning 40 autoimmune diseases
121 [25]. This study used densely-mapped genotyping data and the observed pattern of association at
122 the LD locus to estimate each SNP's probability of being a causal variant. Since this collection
123 aggregates multiple immune-related phenotypes we applied weighted PINES scores, where all an-
124 notations corresponding to immune cells were equally up-weighted (see Methods). We compared
125 the collective PINES signal on these 3,625 variants with the results obtained on an additional set
126 of 30,000 background variants that we randomly selected across the genome. Q-Q plots detailing
127 the results highlight a well-calibrated PINES null distribution and clear signal on the fine mapped
128 variant set (Supplementary Figure 8). A comparison of the weighted PINES results with GWAVA,
129 Eigen, and CADD show that PINES delivers the best predictive performance (Figure 3 first panel).

130 Next, we assessed the performance of PINES to nominate causal alleles from loci associated
131 with individual common traits and diseases, including non-immune phenotypes. We extracted fine
132 mapped candidate causal noncoding variants from [25] related to multiple sclerosis, celiac disease,
133 inflammatory bowel disease, and blood lipid levels. We determined weighted and unweighted PINES
134 scores for each of these variants, and compared the outcomes of PINES to those of GWAVA, Eigen
135 and CADD. Phenotype-based annotation weights were automatically assigned based on GWAS data
136 ([26, 27, 28, 29, 30, 31]). PINES consistently delivered the highest AUROC values, with up to 12%
137 improvement over GWAVA, Eigen, or CADD when running PINES in the phenotype-weighted mode
138 (Figure 3).

139 **2.4 PINES detects signal at expression-modulating variants identified in** 140 **a multiplexed reporter assay**

141 To test the performance of PINES to detect functional noncoding variants en masse, we used 230
142 variants that were found to directly affect gene expression in a massively parallel reporter assay
143 (MPRA) [32]. For the selection of variants we used an FDR cutoff (Benjamini-Hochberg) of 1%.
144 MPRA is an extension of the traditional reporter gene setup, whereby the use of unique barcodes in
145 the 3' UTR of the reporter differentiate expression of individual oligos and thus allow for testing of
146 many different sequences simultaneously. Since most variants identified through this approach have
147 not yet been linked to individual phenotypes, we performed an unweighted PINES analysis. PINES

148 delivered the highest AUROC values, with up to 43% improvement over GWAVA, Eigen, or CADD
149 (Figure 4).

150 **2.5 PINES detects functional evidence for variants residing in FANTOM5** 151 **enhancers**

152 We next tested the power of PINES to correctly prioritize 9,000 variants residing within enhancers
153 that have been identified by the FANTOM5 project through cap analysis of gene expression [6].
154 As expected, PINES correctly assigned noncoding variants residing in these regions the highest
155 relevance scores relative to genomic background variation. Importantly, in doing so PINES outscored
156 GWAVA, Eigen, and CADD considerably, as demonstrated by a Wilcoxon signed rank test between
157 the weighted PINES results and those of all other methods, which delivered p-values strictly below
158 10^{-60} (Figure 5).

159 **2.6 PINES predicts novel causal variants for Parkinson’s disease and IBD** 160 **through fine mapping of GWAS loci**

161 We next tested whether PINES can be applied to predict novel functional noncoding SNPs from
162 GWAS loci. For this, we applied PINES to all GWAS loci associated in recent meta-analyses at
163 genome-wide significance with Parkinson’s disease [33] and IBD [29]. We concentrated on those loci
164 where all SNPs in LD of $R^2 \geq 0.4$ to the GWAS lead SNP were either intronic or intergenic. For
165 both Parkinson’s disease and IBD we used PINES to determine enrichment-based phenotype-specific
166 weights from the complete set of significantly disease-associated GWAS SNPs. We then ran weighted
167 PINES to fine map the most likely causal SNP across the 12 selected Parkinson’s disease loci and
168 19 selected IBD loci. With this approach, PINES distinguished 16 novel noncoding alleles from the
169 background that can be assumed with a high likelihood of being causal for conferring risk for Parkin-
170 son’s disease (rs10878226, rs3756063, rs2301134, rs36121867, rs1954874, rs9275373, rs117896735) and
171 IBD (rs35493230, rs2187892, rs4672507, rs4845604, rs2019262, rs10489630, rs12622128, rs55776317,
172 rs7685642)(Figure 6 and Supplementary Tables 2 and 3).

173 **3 Discussion**

174 The field of human genetics has accumulated thousands of linkage and association signals. The focus
175 is now rapidly shifting towards the identification of functional DNA variants underlying these signals,

176 biological interpretation of their roles, and generation of mechanistic hypothesis of disease etiology.
177 However, the notion of biological function of an allele is diffuse. Genetically mediated phenotypic
178 presentations are usually limited to a specific organ system, tissue or even cell type. Some of them are
179 pleiotropic and affect several systems, but very few represent truly systemic disorders or traits that
180 impact every cell in the body. This suggests that, from the genetic perspective, the notion of function
181 only makes sense in the specific phenotypic context defined by cell type, developmental stage, and
182 stimulus response. This is especially true for regulatory variants involved in transcriptional control.
183 A number of recent studies showed that genetic association signals are enriched in putative regulatory
184 elements, and that this enrichment is highly cell type-specific [25]. However, many experimental and
185 almost all computational approaches to probe the functional effects of allelic variants are agnostic
186 about the context of the phenotypic presentation.

187 Functional genomics is now actively embracing the multitude of contexts, starting from cell type
188 variability in epigenetic annotations [5]. PINES leverages this annotation richness and attempts to
189 predict the actual functional effect in the most relevant context rather than in the abstract framework
190 of ubiquitous functional relevance. We note that simply restricting the analysis to the most relevant
191 cell type is not the optimal approach. From purely statistical perspective, noisy correlated data
192 provide information and should not be completely neglected. More importantly, from the biological
193 perspective, many alleles are pleiotropic and many phenotypes are influenced by different biological
194 processes in different organs, tissues and cells. For example, risk of myocardial infarction is partly
195 influenced by blood lipids, but many genetic contributions are unrelated to blood lipid levels and
196 are likely mediated by the vascular effects. All autoimmune diseases are influenced by the adaptive
197 immune system, but individual conditions are limited to specific organs. PINES addresses this
198 complexity to some degree through its customizable weighting of annotations. Additionally, many
199 cell types that are relevant to a phenotype are currently not represented in the ENCODE and
200 Roadmap datasets. The ability of PINES to leverage information from related cell types and tissues
201 enables the analysis of noncoding variants even for such phenotypes. Finally, the noncoding genome
202 has been consistently linked to human phenotypes through our knowledge of conservation, GWAS
203 peak localization, and eQTLs, yet so far only few noncoding loci have been experimentally validated.
204 This lack of unambiguously-defined functional noncoding loci makes the unsupervised approach used
205 by PINES very versatile.

206 PINES can be easily queried through a web server in a similar manner to PolyPhen. This portal
207 allows for scoring of noncoding SNVs based on user-defined weighting schemes, making PINES
208 immediately applicable across a wide range of phenotypes. In addition, since the web server performs

209 all data processing, users can query PINES with minimal computational overhead. PINES allows for
210 the addition of epigenetic annotations as they become available without requiring significant changes
211 to the underlying statistical model or software implementation. Due to this ease of upgrading the
212 underlying annotation database, we aim PINES to become an always-up-to-date resource for the
213 scientific community.

214 In conclusion, PINES' ability to take advantage of a wide range of prior biological information
215 allows it to improve on the predictive power of other methods, and to provide an enhanced priori-
216 tization of phenotype-relevant variants. PINES avoids biases stemming from inaccurate labeling of
217 training datasets, and benefits from increased power when prior information is available to direct
218 analyses towards relevant annotations. There is a great need for such methods since identification
219 of regulatory activity specific to a subgroup of cell types or tissues can greatly increase our un-
220 derstanding of disease mechanisms. We have shown that PINES can assist in identifying functional
221 noncoding variants in fine mapping analyses, both for complex disease and Mendelian traits, without
222 requiring the significant resource expenditure involved in a traditional fine mapping study.

223 4 Acknowledgements

224 We would like to thank Dr. Ivan Adzhubey for help with setting up the PINES web server. We
225 also acknowledge Dr. Alex Bloemendal for the suggestion to reduce the noise in the annotation
226 covariance matrix by relying on a low-rank approximation of this matrix.

227 5 Online Methods

228 5.1 Annotation sources

229 PINES uses a wide range of annotations as part of the scoring algorithm. Open chromatin and
230 histone modifications for 127 cell types and tissues were obtained from ENCODE and Roadmap
231 Epigenomics ChIP-seq and DNase-seq peak sets. In order to capture combinatorial interactions
232 between different chromatin marks in their spatial context, we used ChromHMM chromatin state
233 segmentations from Roadmap Epigenomics computed via the standard 15-state HMM model. Chro-
234 matin interaction data from a variety of assays (3C, 4C-Seq, 5C, Hi-C, ChIA-PET, Capture-C) were
235 obtained from the 4DGenome database [8]. Additional DNaseI regions inferred via HMM from EN-
236 CODE and Roadmap Epigenomics data were obtained from the Reg2Map database. Conservation
237 was evaluated via GERP [34] and SiPhy [35].

238 Noncoding background variants were selected randomly across the genome, and we used the
239 ClinVar database [36] and GWAS Catalog [37] to ensure that no overlap exists with known functional
240 loci. For the analysis in Figure 3 we used GWAS studies of IBD [28, 29], celiac disease [27], blood
241 lipid levels [30, 31], and multiple sclerosis [26] to determine enrichment-based weights for weighted
242 PINES. We then used fine mapped variants on the corresponding phenotypes from [25] as our test
243 set. All immune-related fine mapped variants in [25] with posterior probability ≥ 0.2 were used to
244 generate the first panel in Figure 3. A list of FANTOM5 enhancers [6] was used to create Figure
245 5. The analysis of all noncoding Parkinson’s disease and IBD loci (Figure 6) was based on regions
246 identified in [33] and [29].

247 5.2 Working with a high-dimensional correlated annotation space

248 Individual variants are assigned a score of 0 or 1 for each of the annotations referenced above. In
249 particular, each variant is characterized by a vector of length 639 composed as follows:

- 250 • Presence or absence of H3K4me1, H3K4me3, H3K27ac, H3K9ac, and DNase annotations for
251 each of the 127 epigenomes (635 values).
- 252 • Presence or absence of a conserved region as predicted by GERP and SiPhy (2 values).
- 253 • Presence or absence of a DHS region as predicted by the ChromHMM 15 state model trained
254 on all epigenomes (1 value).
- 255 • Presence or absence of a region involved in chromatin interactions with other regions as re-
256 ported in the 4DGenome database (1 value).

257 In our annotation dataset, each variant is thus characterized by a vector of 635 cell type-specific
258 scores and 4 cell type-independent scores. The joint distribution of this vector is difficult to ascertain
259 explicitly due to its complex correlation structure; there are few outlier detection techniques that
260 are robust to correlated data. One such approach is the one-class support vector machine (SVM)
261 [38, 39, 40], which fits a hyperplane or hypersphere to the data in an attempt to isolate outlying
262 points. One-class SVMs however suffer from a few disadvantages, such as difficulty in choosing
263 tuning parameters, and the inability to add user-specified feature weights.

The alternative approach used in PINES is based on angular distances in a de-correlated annotation space. Let \mathbf{X} be the 639-dimensional matrix of annotations with covariance matrix Σ and mean vector μ , and let \mathbf{W} be a diagonal matrix of annotation weights (which in an unweighted analysis is the identity matrix). Since Σ is a noisy estimate of the true correlation structure of \mathbf{X} , we perform

the spectral decomposition $\Sigma = \sum_{i=1}^{639} \lambda_i \mathbf{u}_i \mathbf{u}_i^T$, and compute a low-rank approximation of the estimated covariance matrix based on the first 30 eigenvectors (chosen by visual inspection of the scree plot): $\hat{\Sigma} = \sum_{i=1}^{30} \lambda_i \mathbf{u}_i \mathbf{u}_i^T$. The matrix $\hat{\Sigma}^{-1}$ is obtained via the Moore-Penrose pseudo-inverse and is used to project annotation vectors corresponding to individual variants into a decorrelated annotation space via a Cholesky transformation: $\hat{\Sigma}^{-1} = \sum_{i=1}^{30} \frac{1}{\lambda_i} \mathbf{u}_i \mathbf{u}_i^T$. If \mathbf{x} is a vector of annotations, then the length of the vector projected into the decorrelated space is $\sqrt{(\mathbf{x} - \mu)^T \mathbf{W} \hat{\Sigma}^{-1} \mathbf{W} (\mathbf{x} - \mu)}$, which corresponds to the reweighed Mahalanobis distance to the mean vector μ [41]. The cosine of the angle between the projections of two vectors \mathbf{x} and \mathbf{y} into the decorrelated annotation space is given by

$$\frac{(\mathbf{x} - \mu)^T \mathbf{W} \hat{\Sigma}^{-1} \mathbf{W} (\mathbf{y} - \mu)}{\sqrt{(\mathbf{x} - \mu)^T \mathbf{W} \hat{\Sigma}^{-1} \mathbf{W} (\mathbf{x} - \mu)} \cdot \sqrt{(\mathbf{y} - \mu)^T \mathbf{W} \hat{\Sigma}^{-1} \mathbf{W} (\mathbf{y} - \mu)}}$$

which corresponds to the correlation between the projections of \mathbf{x} and \mathbf{y} . In PINES we are specifically interested in the angle between the projection of a variant of interest and the projection of the all-1 annotation vector $\mathbf{1}$. This vector is significant since it provides the direction of a point with maximal annotation load, and thus the greatest evidence for functionality. Following the approach presented in [42], we additionally scale this angle by the length of the projected \mathbf{x} vector, resulting in the following PINES score:

$$\text{PINES}(\mathbf{x}) = \frac{\arccos \left[\frac{(\mathbf{x} - \mu)^T \mathbf{W} \hat{\Sigma}^{-1} \mathbf{W} (\mathbf{1} - \mu)}{\sqrt{(\mathbf{x} - \mu)^T \mathbf{W} \hat{\Sigma}^{-1} \mathbf{W} (\mathbf{x} - \mu)} \cdot \sqrt{(\mathbf{1} - \mu)^T \mathbf{W} \hat{\Sigma}^{-1} \mathbf{W} (\mathbf{1} - \mu)}} \right]}{\sqrt{(\mathbf{x} - \mu)^T \mathbf{W} \hat{\Sigma}^{-1} \mathbf{W} (\mathbf{x} - \mu)} \cdot \sqrt{(\mathbf{1} - \mu)^T \mathbf{W} \hat{\Sigma}^{-1} \mathbf{W} (\mathbf{1} - \mu)}}$$

264 Smaller scores indicate greater evidence for functionality. Supplementary Figure 9 presents the re-
 265 lationship between weighted PINES scores and annotation load in the original highly correlated
 266 space on a simulated example, as well as the effect of introducing weights in the model. For this
 267 simulation we generated 20 correlated Bernoulli random variables (representing the epigenetic anno-
 268 tations). The diagonal weight matrix \mathbf{W} was constructed to assign weight 4 to one of the annotations
 269 and weight 1 to the remaining features. The plot shows the relationship between annotation count
 270 and PINES score, as well as shift in the score profile of observations that have the up-weighted
 271 annotation (red points). Finally, to determine the significance level of a given score we compute an
 272 empirical p-value based on a large set of background variants. These background variants represent
 273 150,000 common variant sites (since rare variant sites can potentially harbor more penetrant muta-
 274 tions) selected randomly across the genome and that are not represented in ClinVar or the GWAS
 275 Catalogue.

276 **5.3 Background variants and the null distribution**

277 PINES makes use of a set of background SNPs against which to compare the score for new variants.
278 We randomly selected 150,000 common variants across the genome to serve as background, none of
279 which have been previously tied to phenotypes. Based on this background, PINES reports a one-
280 sided p-value for each input variant, but the scores reported by CADD, GWAVA, and Eigen do not
281 have an absolute unit of meaning and are thus not directly comparable. To enable this comparison
282 we use the collection of background variants to determine a null distribution for each scoring method,
283 and transform the raw CADD, GWAVA, and Eigen scores into empirical one-sided p-values based on
284 their respective null distribution. This approach is similar to the one used to compute scaled CADD
285 scores by transforming raw into rank-based scores [11]. Having performed this normalization, we
286 can compare the results of PINES, CADD, GWAVA, and Eigen directly.

287 **5.4 Choice of feature weights**

288 Weighting of features can be performed manually, as was the case for the pigmentation variants
289 presented in Figure 2, the first panel in Figure 3, and the simulation in Supplementary Figure 9, by
290 setting the weight for biologically relevant annotations to a user-specified constant. Alternatively,
291 when GWAS peaks are available, the weights used by PINES to differentiate between the different cell
292 type-specific annotations are automatically computed based on the enrichment of each annotation
293 across the GWAS loci. Enrichment within GWAS peaks for each annotation is used to set weights
294 for the Parkinson’s disease, QRS prolongation, schizophrenia, and cleft lip and cleft palate variants
295 presented in Figure 2. Such enrichment-based approaches are frequently used in predicting cell
296 types contributing to specific phenotypes when GWAS or fine mapping data is available [25], with
297 highly enriched annotations indicating potentially relevant cell types and disease mechanisms. In
298 particular, we used the corresponding $-\log_{10}(\text{enrichment p-value})$ as weight for every annotation,
299 although different functional forms are possible. Regardless of whether a manual or enrichment-
300 based weighting is employed to construct the matrix \mathbf{W} , no annotation will be completely excluded
301 from the model. For example in a study of pigmentation, the objective is for variants that have
302 melanocyte-related annotations as well as exhibit evidence of functionality in other cell types to
303 receive more significant scores than variants that only have melanocyte-related annotations. Another
304 reason to rely on data from multiple cell types, even when the phenotypic effect of variants is limited
305 to a single, well characterized cell type, is to gain statistical power from accumulating noisy correlated
306 datasets.

307 5.5 Other Methods

308 A few approaches have been recently proposed to score noncoding regions and address the complexity
309 of the annotation data in a principled manner. The Genome-Wide Annotation of Variants method
310 (GWAVA) [10] aims to predict the impact of noncoding genetic variants based on a random forest
311 classifier, using variants reported in the Human GeneMutation Database (HGMD) as deleterious
312 training data, and common SNPs from the 1000 Genomes Project as benign examples. The CADD
313 approach [11] is based on the premise that harmful mutations are edged out of the gene pool over
314 time via natural selection and that variation that has not been selected against is thus less likely to
315 be deleterious. Notable for CADD is that it uses a dataset of simulated mutations for training, which
316 is then compared to observed variants. A score of deleteriousness is assigned to every possible SNP in
317 the human genome. One of the most recent methods, Eigen [9], is an unsupervised scoring framework
318 that uses the eigen-decomposition of the covariance matrix associated with a collection of functional
319 annotations to compute variant scores representing weighted sums of individual annotations.

References

- 320
- 321 [1] Mijke Visser, Manfred Kayser, and Robert-Jan Palstra. Herc2 rs12913832 modulates human
322 pigmentation by attenuating chromatin-loop formation between a long-range enhancer and the
323 oca2 promoter. *Genome research*, 22(3):446–455, 2012.
- 324 [2] Hans Eiberg, Jesper Troelsen, Mette Nielsen, Annemette Mikkelsen, Jonas Mengel-From,
325 Klaus W Kjaer, and Lars Hansen. Blue eye color in humans may be caused by a perfectly
326 associated founder mutation in a regulatory element located within the herc2 gene inhibiting
327 oca2 expression. *Human genetics*, 123(2):177–187, 2008.
- 328 [3] Dimitri J Stavropoulos, Daniele Merico, Rebekah Jobling, Sarah Bowdin, Nasim Monfared,
329 Bhooma Thiruvahindrapuram, Thomas Nalpathamkalam, Giovanna Pellecchia, Ryan KC Yuen,
330 Michael J Szego, et al. Whole-genome sequencing expands diagnostic utility and improves
331 clinical management in paediatric medicine. *npj Genomic Medicine*, 1:15012, 2016.
- 332 [4] ENCODE Project Consortium et al. The encode (encyclopedia of dna elements) project. *Sci-*
333 *ence*, 306(5696):636–640, 2004.
- 334 [5] Anshul Kundaje, Wouter Meuleman, Jason Ernst, Misha Bilenky, Angela Yen, Alireza Heravi-
335 Moussavi, Pouya Kheradpour, Zhizhuo Zhang, Jianrong Wang, Michael J Ziller, et al. Integra-
336 tive analysis of 111 reference human epigenomes. *Nature*, 518(7539):317–330, 2015.
- 337 [6] Erik Arner, Carsten O Daub, Kristoffer Vitting-Seerup, Robin Andersson, Berit Lilje, Finn
338 Drabløs, Andreas Lennartsson, Michelle Rønnerblad, Olga Hrydziusko, Morana Vitezic, et al.
339 Transcribed enhancers lead waves of coordinated transcription in transitioning mammalian cells.
340 *Science*, 347(6225):1010–1014, 2015.
- 341 [7] Sean Whalen, Rebecca M Truty, and Katherine S Pollard. Enhancer-promoter interactions are
342 encoded by complex genomic signatures on looping chromatin. *Nature genetics*, 48(5):488–496,
343 2016.
- 344 [8] Li Teng, Bing He, Jiahui Wang, and Kai Tan. 4dgenome: a comprehensive database of chromatin
345 interactions. *Bioinformatics*, 31(15):2560–2564, 2015.
- 346 [9] Iuliana Ionita-Laza, Kenneth McCallum, Bin Xu, and Joseph D Buxbaum. A spectral approach
347 integrating functional genomic annotations for coding and noncoding variants. *Nature genetics*,
348 2016.

- 349 [10] Graham RS Ritchie, Ian Dunham, Eleftheria Zeggini, and Paul Flicek. Functional annotation
350 of noncoding sequence variants. *Nature methods*, 11(3):294–296, 2014.
- 351 [11] Martin Kircher, Daniela M Witten, Preti Jain, Brian J O’Roak, Gregory M Cooper, and Jay
352 Shendure. A general framework for estimating the relative pathogenicity of human genetic
353 variants. *Nature genetics*, 46(3):310, 2014.
- 354 [12] Adam Siepel, Gill Bejerano, Jakob S Pedersen, Angie S Hinrichs, Minmei Hou, Kate Rosen-
355 bloom, Hiram Clawson, John Spieth, LaDeana W Hillier, Stephen Richards, et al. Evolution-
356 arily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome research*,
357 15(8):1034–1050, 2005.
- 358 [13] Ilan Gronau, Leonardo Arbiza, Jaaved Mohammed, and Adam Siepel. Inference of natural
359 selection from interspersed genomic elements based on polymorphism and divergence. *Molecular*
360 *biology and evolution*, page mst019, 2013.
- 361 [14] Gosia Trynka, Cynthia Sandor, Buhm Han, Han Xu, Barbara E Stranger, X Shirley Liu, and
362 Soumya Raychaudhuri. Chromatin marks identify critical cell types for fine mapping complex
363 trait variants. *Nature genetics*, 45(2):124–130, 2013.
- 364 [15] Bing Liu, Yang Dai, Xiaoli Li, Wee Sun Lee, and Philip S Yu. Building text classifiers using
365 positive and unlabeled examples. In *Data Mining, 2003. ICDM 2003. Third IEEE International*
366 *Conference on*, pages 179–186. IEEE, 2003.
- 367 [16] Peng Yang, Xiao-Li Li, Jian-Ping Mei, Chee-Keong Kwoh, and See-Kiong Ng. Positive-
368 unlabeled learning for disease gene identification. *Bioinformatics*, 28(20):2640–2647, 2012.
- 369 [17] Jason Ernst and Manolis Kellis. Chromhmm: automating chromatin-state discovery and char-
370 acterization. *Nature methods*, 9(3):215–216, 2012.
- 371 [18] Mijke Visser, Robert-Jan Palstra, and Manfred Kayser. Human skin color is influenced by an
372 intergenic dna polymorphism regulating transcription of the nearby *bnc2* pigmentation gene.
373 *Human molecular genetics*, page ddu289, 2014.
- 374 [19] Frank Soldner, Yonatan Stelzer, Chikdu S Shivalila, Brian J Abraham, Jeanne C Latourelle,
375 M Inmaculada Barrasa, Johanna Goldmann, Richard H Myers, Richard A Young, and Rudolf
376 Jaenisch. Parkinson-associated risk variant in distal enhancer of α -synuclein modulates target
377 gene expression. *Nature*, 533(7601):95–99, 2016.

- 378 [20] Malou van den Boogaard, Scott Smemo, Ozanna Burnicka-Turek, David E Arnolds, Harmen JG
379 van de Werken, Petra Klous, David McKean, Jochen D Muehlschlegel, Julia Moosmann, Okan
380 Toka, et al. A common genetic variant within *scn10a* modulates cardiac *scn5a* expression. *The*
381 *Journal of clinical investigation*, 124(4):1844–1852, 2014.
- 382 [21] Catherine A Guenther, Bosiljka Tasic, Liqun Luo, Mary A Bedell, and David M Kingsley. A
383 molecular basis for classic blond hair color in europeans. *Nature genetics*, 46(7):748–752, 2014.
- 384 [22] Nevena Cvjetkovic, Lorena Maili, Katelyn S Weymouth, S Shahrukh Hashmi, John B Mulliken,
385 Jacek Topczewski, Ariadne Letra, Qiuping Yuan, Susan H Blanton, Eric C Swindell, et al.
386 Regulatory variant in *fzd6* gene contributes to nonsyndromic cleft lip and palate in an african-
387 american family. *Molecular genetics & genomic medicine*, 3(5):440–451, 2015.
- 388 [23] William P Gilks, Matthew Hill, Michael Gill, Gary Donohoe, Aiden P Corvin, and Derek W
389 Morris. Functional investigation of a schizophrenia gwas signal at the *cdc42* gene. *The World*
390 *Journal of Biological Psychiatry*, 2012.
- 391 [24] Elizabeth J Leslie, Margaret A Taub, Huan Liu, Karyn Meltz Steinberg, Daniel C Koboldt,
392 Qunyuhan Zhang, Jenna C Carlson, Jacqueline B Hetmanski, Hang Wang, David E Larson, et al.
393 Identification of functional variants for cleft lip with or without cleft palate in or near *pax7*,
394 *fgr2*, and *nog* by targeted sequencing of gwas loci. *The American Journal of Human Genetics*,
395 96(3):397–411, 2015.
- 396 [25] Kyle Kai-How Farh, Alexander Marson, Jiang Zhu, Markus Kleinewietfeld, William J Hous-
397 ley, Samantha Beik, Noam Shores, Holly Whitton, Russell JH Ryan, Alexander A Shishkin,
398 et al. Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature*,
399 518(7539):337–343, 2015.
- 400 [26] International Multiple Sclerosis Genetics Consortium, Wellcome Trust Case Control Consortium
401 2, et al. Genetic risk and a primary role for cell-mediated immune mechanisms in multiple
402 sclerosis. *Nature*, 476(7359):214–219, 2011.
- 403 [27] Patrick CA Dubois, Gosia Trynka, Lude Franke, Karen A Hunt, Jihane Romanos, Alessandra
404 Curtotti, Alexandra Zhernakova, Graham AR Heap, Róza Ádány, Arpo Aromaa, et al. Mul-
405 tiple common variants for celiac disease influencing immune gene expression. *Nature genetics*,
406 42(4):295–302, 2010.

- 407 [28] Luke Jostins, Stephan Ripke, Rinse K Weersma, Richard H Duerr, Dermot P McGovern,
408 Ken Y Hui, James C Lee, L Philip Schumm, Yashoda Sharma, Carl A Anderson, et al. Host-
409 microbe interactions have shaped the genetic architecture of inflammatory bowel disease. *Nature*,
410 491(7422):119–124, 2012.
- 411 [29] Jimmy Z Liu, Suzanne van Sommeren, Hailiang Huang, Siew C Ng, Rudi Alberts, Atsushi
412 Takahashi, Stephan Ripke, James C Lee, Luke Jostins, Tejas Shah, et al. Association analyses
413 identify 38 susceptibility loci for inflammatory bowel disease and highlight shared genetic risk
414 across populations. *Nature genetics*, 47(9):979–986, 2015.
- 415 [30] Tanya M Teslovich, Kiran Musunuru, Albert V Smith, Andrew C Edmondson, Ioannis M
416 Stylianou, Masahiro Koseki, James P Pirruccello, Samuli Ripatti, Daniel I Chasman, Cristen J
417 Willer, et al. Biological, clinical and population relevance of 95 loci for blood lipids. *Nature*,
418 466(7307):707–713, 2010.
- 419 [31] Global Lipids Genetics Consortium et al. Discovery and refinement of loci associated with lipid
420 levels. *Nature genetics*, 45(11):1274–1283, 2013.
- 421 [32] Ryan Tewhey, Dylan Kotliar, Daniel S Park, Brandon Liu, Sarah Winnicki, Steven K Reilly,
422 Kristian G Andersen, Tarjei S Mikkelsen, Eric S Lander, Stephen F Schaffner, et al. Direct
423 identification of hundreds of expression-modulating variants using a multiplexed reporter assay.
424 *Cell*, 165(6):1519–1529, 2016.
- 425 [33] Mike A Nalls, Nathan Pankratz, Christina M Lill, Chuong B Do, Dena G Hernandez, Mohamad
426 Saad, Anita L DeStefano, Eleanna Kara, Jose Bras, Manu Sharma, et al. Large-scale meta-
427 analysis of genome-wide association data identifies six new risk loci for parkinson’s disease.
428 *Nature genetics*, 46(9):989–993, 2014.
- 429 [34] Eugene V Davydov, David L Goode, Marina Sirota, Gregory M Cooper, Arend Sidow, and
430 Serafim Batzoglou. Identifying a high fraction of the human genome to be under selective
431 constraint using *gerp++*. *PLoS Comput Biol*, 6(12):e1001025, 2010.
- 432 [35] Manuel Garber, Mitchell Guttman, Michele Clamp, Michael C Zody, Nir Friedman, and Xi-
433 aohui Xie. Identifying novel constrained elements by exploiting biased substitution patterns.
434 *Bioinformatics*, 25(12):i54–i62, 2009.
- 435 [36] Melissa J Landrum, Jennifer M Lee, Mark Benson, Garth Brown, Chen Chao, Shanmuga
436 Chitipiralla, Baoshan Gu, Jennifer Hart, Douglas Hoffman, Jeffrey Hoover, et al. Clinvar: public

437 archive of interpretations of clinically relevant variants. *Nucleic acids research*, 44(D1):D862–
438 D868, 2016.

439 [37] Danielle Welter, Jacqueline MacArthur, Joannella Morales, Tony Burdett, Peggy Hall, Heather
440 Junkins, Alan Klemm, Paul Flicek, Teri Manolio, Lucia Hindorff, et al. The nhgri gwas catalog,
441 a curated resource of snp-trait associations. *Nucleic acids research*, 42(D1):D1001–D1006, 2014.

442 [38] Bernhard Schölkopf, John C Platt, John Shawe-Taylor, Alex J Smola, and Robert C Williamson.
443 Estimating the support of a high-dimensional distribution. *Neural computation*, 13(7):1443–
444 1471, 2001.

445 [39] Larry M Manevitz and Malik Yousef. One-class svms for document classification. *the Journal*
446 *of machine Learning research*, 2:139–154, 2002.

447 [40] David MJ Tax and Robert PW Duin. Support vector data description. *Machine learning*,
448 54(1):45–66, 2004.

449 [41] Robert A Greevy, Carlos G Grijalva, Christianne L Roumie, Cole Beck, Adriana M Hung,
450 Harvey J Murff, Xulei Liu, and Marie R Griffin. Reweighted mahalanobis distance match-
451 ing for cluster-randomized trials with missing data. *Pharmacoepidemiology and drug safety*,
452 21(S2):148–154, 2012.

453 [42] Hans-Peter Kriegel, Arthur Zimek, et al. Angle-based outlier detection in high-dimensional data.
454 In *Proceedings of the 14th ACM SIGKDD international conference on Knowledge discovery and*
455 *data mining*, pages 444–452. ACM, 2008.

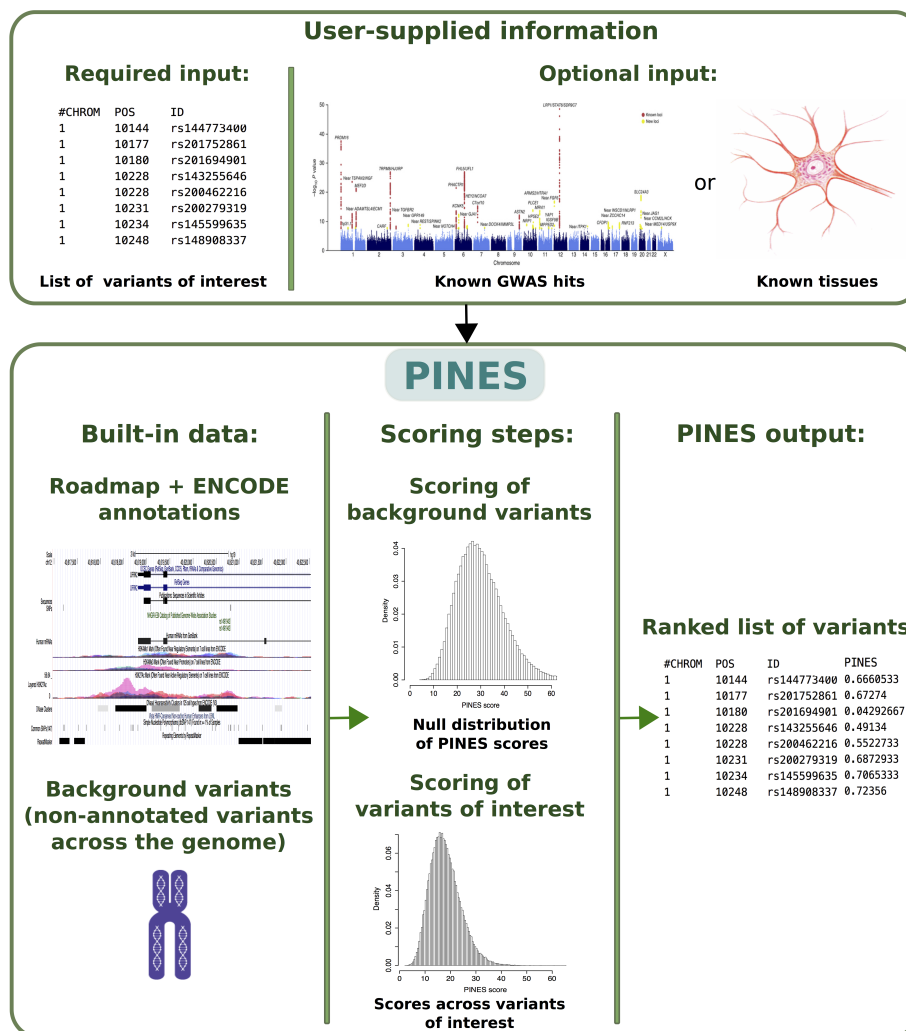


Figure 1: Overview of the PINES framework. PINES aims to systematically predict and rank the functional relevance of noncoding genomic variants. It can either work in a default (“unweighted”) mode and compare user-defined variants against the genomic background. Alternatively, users can customize searches towards annotations considered as of highest relevance to a phenotype of interest, for instance by providing a list of SNPs associated with a disease of interest through GWAS, or by highlighting disease-relevant tissues (“weighted” PINES mode). Scores of genomic background variants serve as an empirical null distribution against which significance levels for each variant of interest are computed and scored in an output file.

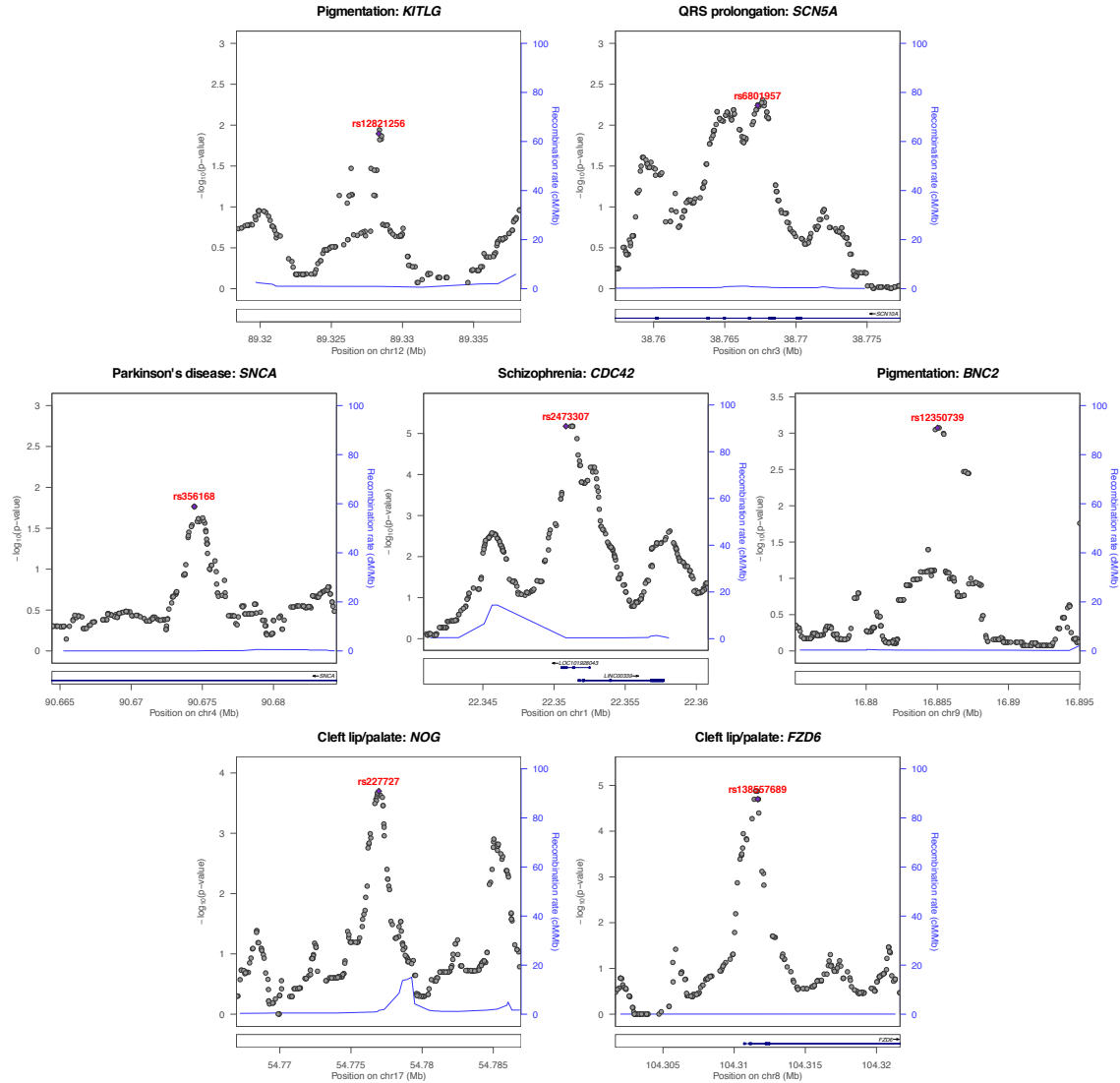


Figure 2: PINES prioritizes experimentally validated functional noncoding variants. We fine map 20kb regions surrounding functional noncoding variants (purple dots) and show that all of the variants validated experimentally as regulating expression of a nearby trait-associated gene are also assigned the highest PINES scores. Supplementary Figures 1-7 show that PINES outperforms existing methods on all loci.

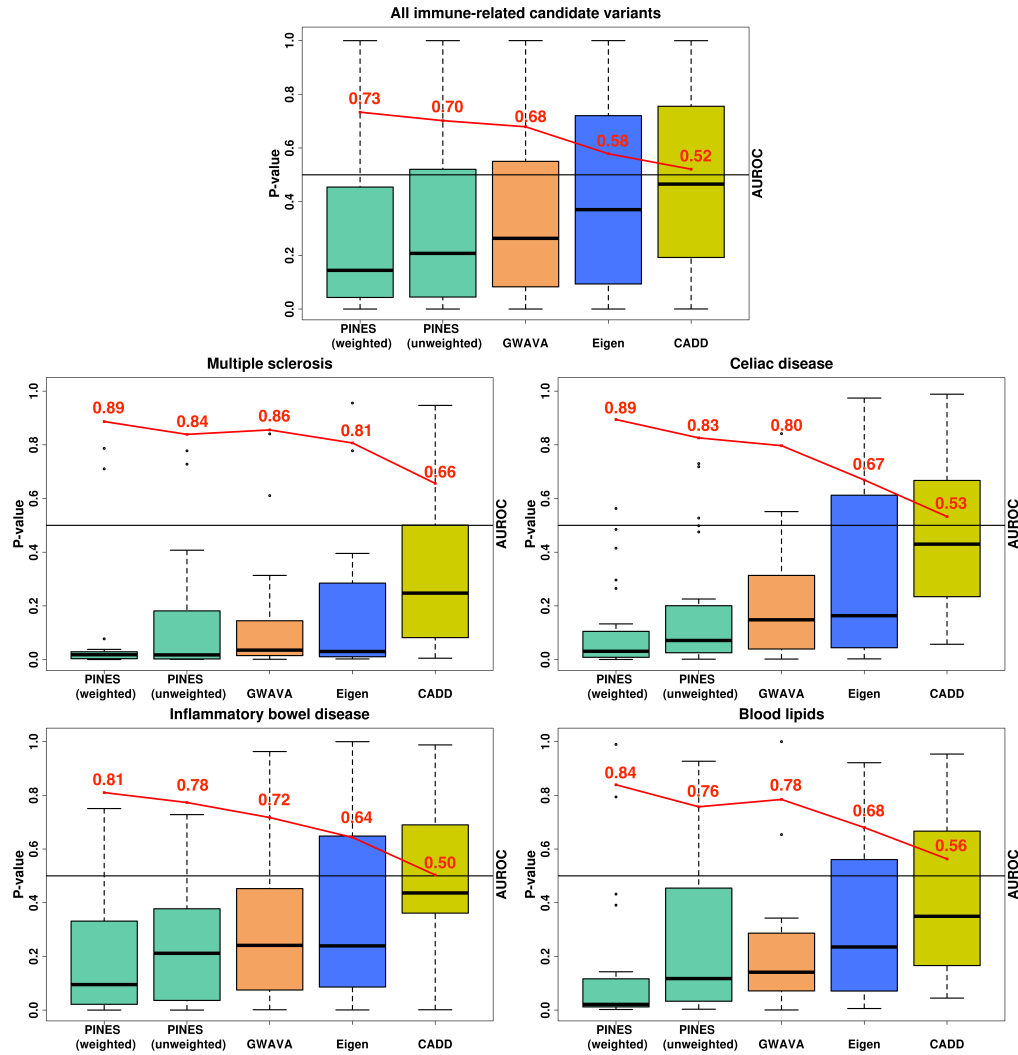


Figure 3: PINES improves statistical power to detect fine mapped variants across common neurologic, immune, and metabolic traits and diseases. AUROC values (red) were computed by selecting 20,000 background variants as negative examples, and the fine mapped variants relevant to each disease as positive examples. PINES achieves the best AUROC values, of to 12% higher than the other methods, based on its inclusion of weights encoding prior disease knowledge (in this case relevant cell types).

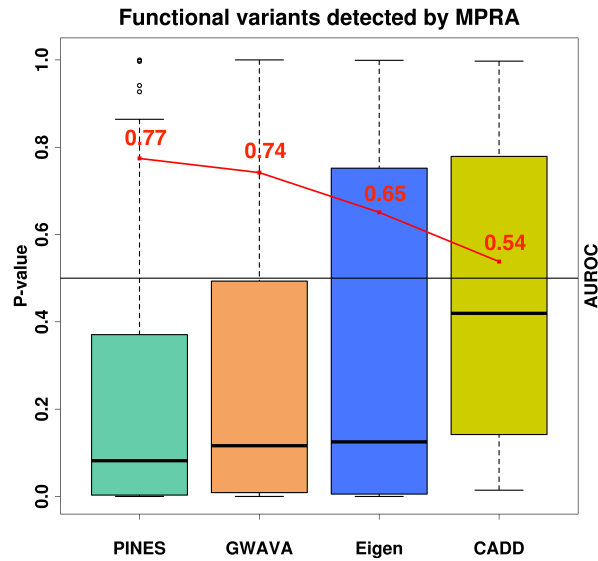


Figure 4: PINES delivers improved statistical power to identify functional noncoding variants detected by a massively parallel reporter assay. The AUROC values (red) were computed by selecting 20,000 background variants as negative examples, and the reported functional variants as positive examples. With AUROC values up to 43% higher than the other methods, PINES outperforms GWAVA, Eigen, and CADD in its ability to detect the functional variants.

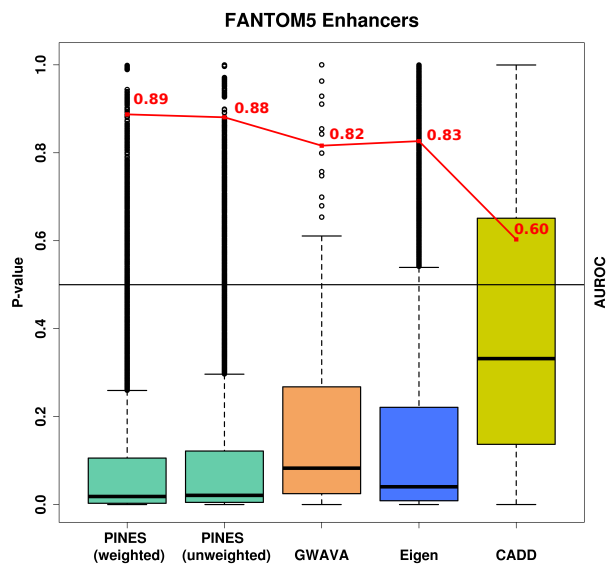


Figure 5: PINES improves the prioritization of variants residing in experimentally validated enhancer regions. The AUROC values (red) were computed by selecting 20,000 background variants as negative examples, and the variants residing in enhancer loci as positive examples. Based on AUROC values, the weighted PINES approach outperforms GWAVA, Eigen, and CADD in its ability to pinpoint enhancer variants. Additionally, testing whether the weighted PINES significance levels are smaller than those of other methods via a Wilcoxon signed rank test delivers p-values that are all below 10^{-60} .

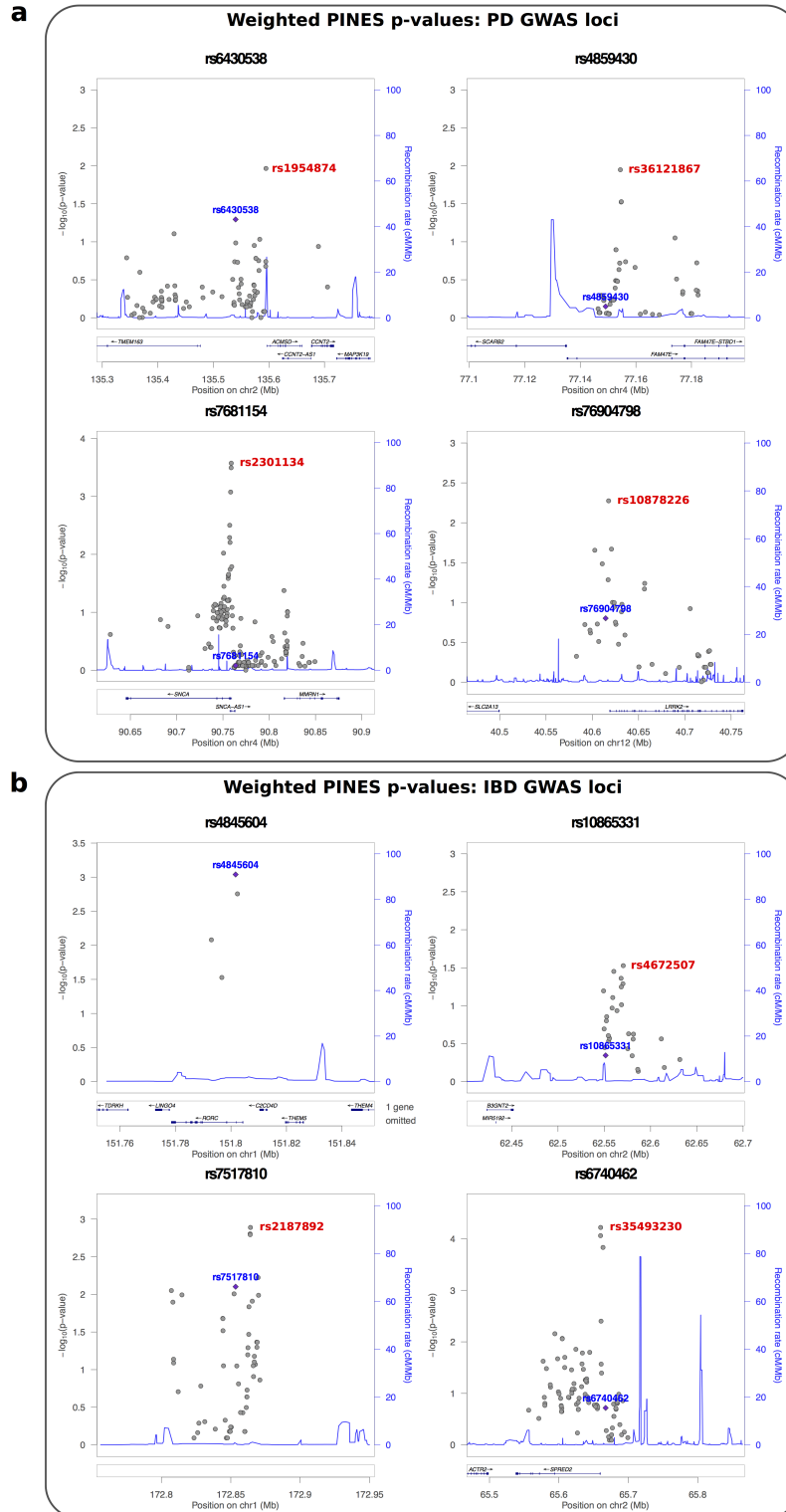


Figure 6: PINES predicts novel noncoding pathogenic variants through fine mapping of Parkinson’s disease and IBD GWAS loci. Loci were extracted from [33] and [29]. For each lead SNP, all variants with $LD \geq 0.4$ were selected, and loci were discarded if this list overlapped any coding regions or 3’ or 5’ UTRs of coding genes. All variants in LD to the lead SNP were scored via weighted PINES. The GWAS lead SNP is marked blue, and the variant predicted as likely causal through PINES fine mapping is marked red. For the rs4845604 locus the GWAS and PINES lead SNP overlaps.