

TVAR: Assessing Tissue-specific Functional Effects of Non-coding Variants with Deep Learning

Hai Yang

East China University of Science and Technology

Rui Chen

Vanderbilt University

Quan Wang

Vanderbilt University

Qiang Wei

Vanderbilt University

Ying Ji

Vanderbilt University <https://orcid.org/0000-0001-5691-1303>

Xue Zhong

Vanderbilt University

Bingshan Li (✉ bingshan.li@Vanderbilt.Edu)

Vanderbilt University <https://orcid.org/0000-0003-2129-168X>

Article

Keywords: whole genome-sequencing, TVAR, deep learning, non-coding variants

Posted Date: December 2nd, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-113771/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

16 **Abstract**

17 Analysis of whole genome-sequencing (WGS) for genetics of disease is still a challenge
18 due to lack of accurate functional annotation of noncoding variants, especially the rare
19 ones. As eQTLs have been extensively implicated in genetics of human diseases, we
20 hypothesize that noncoding rare variants discovered in WGS play a regulatory role in
21 predisposing disease risk. With thousands of tissue- and cell type-specific epigenomic
22 features, we propose TVAR, a multi-label learning based deep neural network that
23 predicts the functionality of noncoding variants in the genome based on eQTLs across
24 49 human tissues in GTEx. TVAR learns the relationships between high-dimensional
25 epigenomics and eQTLs across tissues, taking the correlation among tissues into
26 account to learn shared and tissue-specific eQTL effects. As a result, TVAR outputs
27 tissue-specific annotations, with an average of 0.77 across these tissues. We evaluate
28 TVAR's performance on four complex diseases (coronary artery disease, breast cancer,
29 Type 2 diabetes, and Schizophrenia), using TVAR's tissue-specific annotations, and
30 observe its superior performance in predicting functional variants for both common
31 and rare variants, compared to five existing state-of-the-art tools. We further evaluate
32 TVAR's G-score, a scoring scheme across all tissues, on ClinVar, fine-mapped GWAS loci,
33 Massive Parallel Reporter Assay (MPRA) validated variants, and observe consistently
34 better performance of TVAR compared to other competing tools.

35

36

37 **Introduction**

38 With the rapid development of Whole-genome sequencing (WGS), the interpretation
39 of non-coding regions accounting for 98% of the human genome has become a
40 significant focus of genetic studies¹. Studies have shown that the majority of non-
41 coding variants associated with complex diseases predispose disease risk by regulating
42 the expression of their target genes, often in a tissue-specific manner². Understanding
43 the regulatory role of the non-coding variants will help advance research in genetics,
44 evolution, and precision medicine³. Genome-wide association studies (GWAS) have
45 identified thousands of disease-associated variants, most of which are common in
46 population⁴. In contrast, analysis of rare noncoding variants suffers from low power
47 and severe multiple testing correction. Functional annotation of noncoding variants
48 becomes essential to prioritize potentially risk variants for complex diseases. Massive
49 genomics data, e.g. ENCODE⁵, Roadmap⁶ and FANTOM5⁷ have been used in various
50 computational methods (CADD⁸, GWAVA⁹, DANN¹⁰, LINSIGHT¹¹, EIGEN¹², DVAR¹³) for
51 such purposes. These methods are generally based on machine learning, supervised
52 or unsupervised. In general, supervised learning is powerful when equipped with well-
53 calibrated training set; however, for annotating noncoding variants, these methods are
54 often troubled by the lack of high-accuracy training set¹². On the other hand,
55 unsupervised methods e.g. Eigen¹², PINES¹⁴ and DVAR¹³, discover the internal
56 differences among the groups of variants in the training set, without the need of
57 labeling of true functional variants, and are therefore more robust in the situation
58 without curated functional variants that are representative of risk variants for complex

59 disease. It is well established that risk variants for complex diseases have weak effect
60 sizes, even for rare variants¹⁵. Known disease causing noncoding rare variants in Clinvar
61 for example are usually biased towards variants with large effect sizes, making the
62 machine learning approaches based on such training set unlikely to accurately
63 annotate functional variants for complex diseases. Strategies that are able to model
64 representative functional variants for complex diseases are needed.

65

66 Risk variants from GWAS exert disease risk largely through regulating risk genes.
67 Expression QTLs (eQTLs) have been shown to significantly co-localize with GWAS loci
68 and used to identify risk genes regulated by risk variants¹⁶. Thus, eQTLs provide a
69 robust means to functionally annotate risk variants for complex diseases. An added
70 layer of complexity is that both risk variants and gene expression exhibit strong tissue-
71 specificity¹⁷. It is therefore unlikely that a single annotation is equally effective for
72 different diseases. We argue that strategies taking tissue specificity of both the
73 functional variants and target genes are effective in prioritizing functional variants for
74 complex diseases. Large-scale tissue and cell type specific epigenomics data, e.g.
75 ENCODE⁵, Roadmap Epigenomics¹⁸, have been used in various machine learning
76 methods to infer functionality of noncoding variants. Transcriptome and matched WGS
77 data across a wide range of tissues in GTEx provide genome-wide tissue-specific eQTLs.
78 The reported eQTLs provide robust and unbiased association signals that provide
79 opportunities for accurately identifying tissue-specific risk variants.

80 For a long time, due to the limitation of the size of the existing annotation

81 database and the computer power, most computational models (such as the above
82 approaches) only describe the functional effects of variants at the organism level.
83 Until recently, a wide range of tissue /cell type-specific functional annotations (such as
84 ENCODE, Roadmap, and GTEx¹⁹) are published, and computational methods begin to
85 estimate the tissue-specific functionality of variants. Due to the paucity of known
86 noncoding variants associated with complex diseases, it is challenging to train machine
87 learning approaches using well-validated labels. For this reason, very few methods, e.g.
88 Delta-SVM²⁰ and GenoNet²¹, use supervised training for the prediction of tissue-
89 specific functional effects of noncoding variants. Other methods, e.g. FUN-LDA²²,
90 FitCons²⁴, and ExPecto²³, employed unsupervised learning approaches that do not rely
91 on accurately labeled data on various tissues/cell types. These methods, however,
92 train models individually for each tissue and do not consider the correlation between
93 tissues, resulting in loss of accuracy as well as potential incompatibility among tissues.
94 The challenge of estimating the functionality of variants in various tissues is far from
95 being considered resolved.

96 The following issues need to be addressed urgently to address the challenge of
97 predicting tissue-specific non-coding functional variants. As most eQTLs detectable
98 with sufficient power are common, how to leverage the eQTL resources to predict
99 functionality of rare noncoding variants is a key challenge and opportunity. An
100 additional challenge is how to leverage the genetic sharing and dissimilarity of eQTLs
101 across tissues to achieve tissue specific prediction. To address the first challenge, our
102 rationale is that common and rare regulatory variants share similar regulatory

103 machinery such that they exist similar functional genomics data across tissues and cell
104 types. Based on the rationale, we develop the prediction model to link genomics data
105 to eQTLs, and apply the model to predict functionality of rare variants based on the
106 shared genomics data. For the second challenge, we need to consider the connections
107 between multiple tissues simultaneously, and more specifically need to build a joint
108 model of each eQTL across multiple tissues, which requires a reasonable function that
109 is able to map multi-dimensional functional genomics data to multi-tissue eQTLs.

110 In this study, we propose TVAR, a deep neural network (DNN) that integrates
111 multi-label learning and multi-instance learning to solve the above problems. By using
112 the 1247-dimensional functional genomics features, TVAR accesses the tissue-specific
113 functional scores of each variant across the 49 GTEx tissues (from the GTEx dataset
114 V8). Based on the multi-label DNN, during the training, TVAR learns the differences
115 and connections between tissues, and jointly considers the functional utility of a
116 variant across 49 tissues simultaneously to leverage the sharing of eQTL among tissues.
117 Meanwhile, we also developed a multi-instance learning algorithm, G-score, to
118 provide an integrated functional score for each variant on the organism level. We
119 observed that TVAR achieves an average AUC of ROC of 0.77 on more than 70,000
120 variant datasets in 49 tissues of the GTEx dataset, and found that although the
121 functional input annotations are the same, the maximum and minimum performances
122 of TVAR across various tissues differ significantly, indicating that TVAR captures tissue-
123 specific regulatory machinery in high-dimensional genomics data to make tissue-
124 specific functional prediction of regulatory variants. We applied TVAR scoring on GWAS

125 data of 4 complex diseases, i.e. coronary artery disease, breast cancer, Type 2 diabetes,
126 and Schizophrenia, with tissue-specific annotations from heart, breast, pancreas, and
127 brains tissues in GTEx. We found that the top scoring rare variants (MAF<0.01) have
128 significantly smaller *P* values than the background variants for each of the diseases,
129 supporting that TVAR can be used to prioritize noncoding rare variants in a tissue-
130 specific manner for complex diseases. Finally, we compared the performance of the
131 TVAR G-score algorithm with other organism level score algorithms in four test
132 scenarios (Clinvar^{24,25} set, fine-mapped GWAS²⁶ set, GTEX eQTLs²⁷ set, and MPRA
133 validated variants²⁸ set), and found that TVAR has superior performance across all the
134 test scenarios. In particular, TVAR outperformed DeepSEA and DANN, two methods
135 that are also based on deep learning methods. Overall, we demonstrate that the
136 integration of multi-label learning and multi-instance learning with DNN implemented
137 in TVAR holds promise for scoring non-coding variants in a tissue-specific manner for
138 genetic analysis of complex diseases.

139
140

141 **RESULTS**

142 **Evaluation of TVAR's performance across 49 tissues in the GTEx dataset**

143 Through the TVAR framework, we can predict the functional variants corresponding to
144 the 49 different human tissues in the GTEx (V8 release) (See Fig. 1 for the list of the
145 tissues). Since TVAR is a supervised learning approach, to achieve the goal of
146 identifying tissue-specific functional variants, sufficient and high-quality labels of
147 training individuals are required. To obtain the training labels, we first collected the

148 variants that strongly associated with the eGenes in the GTEx dataset (with q-value
149 cutoff 0.01). Due to the effect of the linkage disequilibrium (LD), most detected eQTLs
150 are proxies to the functional regulatory variants ²⁹. We used a fine mapping strategy
151 to nominate the credible functional variants in an eQTL LD block. Specifically, we used
152 LINSIGHT, a state-of-art functional variant annotation tool, to filter these non-
153 functional variants (Methods). For each tissue, we retained the top 1500 variants with
154 the highest LINSIGHT scores, and merged the variant labels in 49 tissues into a matrix,
155 in which each row is a variant, and each column is a tissue. The label matrix is a 0-1
156 matrix, which is used for the multi-label DNN. Since the variants in the label matrix
157 show functionality in at least one tissue, we set these variants as the positive sample
158 set for model training. Next, we randomly selected variants with MAF > 0.05 in the
159 1000 Genomes Project ³⁰ to match the number of positive sample sets and removed
160 the variants that appeared in the positive set to form the negative sample set. The
161 labels of the negative sample set are all-zeros. Finally, we combined the matrix of the
162 positive and negative samples to obtain a highly reliable training data set (70232
163 variants) derived from the GTEx cohort. We used q-value cutoff ($q < 0.05$) in all tissues
164 to screen positive samples, and observed that out of the 35116 variants, 10421 (~30%)
165 appeared in at least two tissues and 3455 (~10%) appeared in at least four tissues. We
166 consider this strategy of accurately learning and successfully describing the
167 associations between variants and tissues will be the passkey to predicting tissue-
168 specific functional variants.

169 The TVAR framework is a deep feedforward network based on multi-label learning.
170 The network of TVAR describes the functionality of the variant-tissue pairs through the
171 fully connected layers, which can learn the differences and similarities among the 49
172 tissues. The output of TVAR is a 49-dimensional vector that represents the functional
173 scores of the variant-tissue pairs. To prove that TVAR can successfully learn the
174 functionality of the variants corresponding to each issue, we used five-fold cross-
175 validation to train and test the model. In each training process, we randomly selected
176 80% of the data for the model training and 20% of the data for the model testing (see
177 Supplementary Note 1). Receiver Operating Characteristic (ROC) curves were adopted
178 to distinguish the prediction power of different methods in all testing processes. As an
179 overall evaluation, we combined all variants across all tissues as a single evaluation set,
180 and obtained an average AUC = 0.770 across the GTEx 49 tissues, indicating that
181 through multi-label learning, the network can extract valid features from the input
182 1247-dimensional annotations predictive of eQTLs. Based on the average accuracy rate,
183 we divided the GTEx 49 tissues into the 'high accuracy group' and the 'low accuracy
184 group' (Fig. 1A).

185 In the high accuracy group (Fig. 1A), The AUCs of 15 tissues exceed 0.80, including
186 'Nerve_Tibial' (AUC=0.854), 'Cells_Transformed_fibroblasts' (AUC=0.851),
187 'Artery_Tibial' (AUC = 0.846) , 'Skin_Sun_Exposed_Lower_leg' (AUC = 0.843),
188 'Thyroid' (AUC = 0.841), 'Adipose_Subcutaneous' (AUC = 0.836), 'Muscle_Skeletal
189 '(AUC = 0.834), 'Whole_Blood' (0.826), 'Esophagus_Muscularis' (AUC = 0.826),
190 'Skin_Not_Sun_Exposed_Suprapubic' (AUC = 0.825), 'Esophagus_Mucosa' (AUC =

191 0.822), 'Lung' (AUC = 0.814), 'Testis' (AUC = 0.808), ' Artery_Aorta' (AUC = 0.802), and
192 'Adipose_VisceralO' (AUC = 0.802). On the remaining 10 tissues, the AUCs of TVAR
193 exceeded (or approached) the average accuracy rate (AUC = 0.77): 'Colon_Transverse'
194 (AUC = 0.795), 'Pancreas' (AUC = 0.792), 'Heart_Atrial_Appendage' (AUC = 0.790),
195 'Breast_Mammary_Ttissue' (AUC = 0.787), 'Esophagus_Gastroesophageal_Junction'
196 (AUC = 0.785), 'Heart_Left_Ventricle (AUC = 0.784)', 'Colon_Sigmoid' (AUC = 0.781),
197 'Stomach' (AUC = 0.779), 'Sple' (AUC = 0.777), 'Brain_Cerebellum' (AUC = 0.768). We
198 found that TVAR often achieved similar AUCs on similar tissues, such as
199 'Skin_Sun_Exposed_Lower_leg' and 'Skin_Not_Sun_Exposed_Suprapubic',
200 'Adipose_Subcutaneous' and 'Adipose_Visceral_Omentum', 'Heart_Left_Ventricle'
201 and 'Heart_Atrial_Appendage', ' Colon_Transverse ' and ' Colon_Sigmoid'.

202 In the low performance group (Fig. 1B), The AUCs of 80% tissues exceed 0.70,
203 including 'Pituitary' (AUC = 0.767), 'Adrenal_Gland' (AUC = 0.762), 'Pituitary'
204 (AUC=0.752), 'Brain_Cerebellar_Hemisphere' (AUC=0.751), 'Brain_Cortex'
205 (AUC=0.748), 'Brain_Nucleus_accumbens_basal_ganglia' (AUC=0.747),
206 'Brain_Caudate_basal_ganglia' (AUC=0.745), 'Brain_Frontal_Cortex_BA9' (AUC=0.742),
207 'Cells_EBV-transformed_lymphocytes' (AUC=0.742), 'Small_Intestine_Terminal_Ileum'
208 (AUC=0.733), 'Artery_Coronary' (AUC=0.733), 'Brain_Putamen_basal_ganglia'
209 (AUC=0.732), 'Liver' (AUC=0.724), 'Ovary' (AUC=0.724), 'Brain_Hypothalamus'
210 (AUC=0.722), 'Uterus' (AUC=0.719), 'Brain_Hippocampus' (AUC=0.719), 'Vagina'
211 (AUC=0.713), 'Brain_Anterior_cingulate_cortex_BA24' (AUC=0.712),
212 'Minor_Salivary_Gland' (AUC=0.709). TVAR's performance on the remaining 4 tissues

213 in the low accuracy group is as follows: 'Brain_Amygdala' (AUC = 0.699),
214 'Brain_Substantia_nigra' (AUC = 0.696), 'Kidney_Cortex' (AUC = 0.695),
215 'Brain_Spinal_cord_cervical_c-1' (AUC = 0.690). TVAR's predictions in related tissues
216 (such as brain tissues) are also similar (Table S2).

217

218 After the above analysis, we found that TVAR achieved the best performance on
219 'Nerve_Tibial' (AUC = 0.854), and the lowest performance on 'Brain_Substantia_nigra'
220 (AUC = 0.690). The vast differences in performance indicate that the input features of
221 the model have different interpretation capabilities for functional variants in different
222 tissues. Overall, TVAR achieved reasonable predict power (AUC > 0.7) on the majority
223 of the tissues (more than 91%).

224

225 **The relevance of TVAR scores across the 49 GTEx tissues**

226 After the training process with more than 70,000 variants on the GTEx dataset, TVAR
227 successfully achieved stable performance in predicting functional variants on the GTEx
228 49 tissues. We aim to investigate the predictive scores across the tissues. Overall, it is
229 clear that the more unrelated the tissues are, the higher dissimilarity of the TVAR
230 scores are. To explore this phenomenon further, we selected the most significant
231 eQTLs (q-value < 1e-30) separately in each of the 49 tissues of the GTEx dataset (we
232 removed the variants in the training set of the TVAR), and formed a testing dataset
233 (with 1240 variants); the variants in this dataset can be considered the most tissue-
234 specific functional variants. Next, we use the TVAR framework to score all variants in

235 this dataset. For each tissue, the TVAR outputs a 1240-dimensional scoring vector. We
236 use these scoring vectors to calculate the distance matrix of 49 tissues and use
237 agglomerative hierarchical clustering to analyze the patterns of TVAR scores across
238 different tissues. Similar to the results of cis-eQTLs across the GTEx tissues¹⁹, TVAR's
239 scoring pattern confers the specificity of tissues (Fig. 2). For example, most tissues
240 related to the brain are adjacent and located in the upper half of the figure, and two
241 heart-related tissues are adjacent (Fig. 2). Other tissues have similar clustering, e.g.
242 two colon related tissues, two adipose related tissues, and two esophagus related
243 tissues (Fig. 2). It should be emphasized that the tissue-specific patterns in these
244 results are solely derived from the knowledge learned by the TVAR model during the
245 previous training process. The clustering patterns suggest that the TVAR scoring results
246 are indeed tissue-specific, and the same variant will get similar scores in related tissues,
247 and vice versa.

248

249 **Identification of the tissue-specific functionality of GWAS Hits**

250 To further investigate the usefulness of TVAR in predicting disease risk variants, we
251 selected four complex diseases, i.e. coronary artery disease (CAD), breast cancer (BRC),
252 Type 2 diabetes (T2D), and Schizophrenia (SCZ), and mapped the diseases to relevant
253 tissues ('Heart_Left_Ventricle' for CAD, 'Breast_Mammary_Tissue' for BRC, 'Pancreas'
254 for T2D, and 'Brain_Hippocampus' for SCZ). We used the TVAR scoring algorithm to
255 predict GWAS hits on these diseases. We collected GWAS hits from the GWAS catalog
256 (All associations v1.0) with $P < 5 * 10^{-8}$, removing variants used in the TVAR training set.

257 In total, we obtained 325 GWAS hits for CAD, 345 GWAS hits for BRC, 200 GWAS hits
258 for T2D, and 285 GWAS hits for SCZ. We randomly selected an equal number of variants
259 from 1000 Genomes Project ³¹ (MAF> 0.05) as the negative samples for each disease.
260 To eliminate the influence of LD on the positive and negative sample sets, we set the
261 distance between all the variants in the negative sample set and the variants in the
262 positive sample set above 100kb. Previous studies have shown that the functional
263 scores of the positive sample set should be significantly higher than that of randomly
264 selected variants¹². We obtained the TVAR scores corresponding to these variants in
265 the four tissues related to CAD, BRC, T2D, and SCZ.

266 We compared TVAR's performance with five other state-of-art methods: CADD,
267 Eigen, DANN, DeepSEA, and FUN-LDA. CADD is a supervised learning method based on
268 SVM, which is a widely used variant score evaluation tool. DANN uses a deep neural
269 network to replace the SVM module in CADD, and performs better than CADD in
270 certain scenarios¹⁰; Eigen is a general variant evaluation method based on
271 unsupervised learning, and DeepSEA is a functional variant evaluation method based
272 on deep convolutional networks. FUN-LDA is a recent work that is based on
273 unsupervised learning and promotes tissue-specific functional variants prediction (it is
274 also the only method we investigated that supports scoring on a large number of
275 tissues and the corresponding scores on these tissues are available online). For CADD,
276 Eigen, DANN, and DeepSEA, we directly took the score corresponding to the variant
277 for the comparison. For FUN-LDA, it is based on the RoadMap 111 tissues, which are
278 different from the 49 tissues used in TVAR. Here we chose E095 ('Left ventricle') related

279 to CAD, E028 ('Breast vHMEC mammary epithelial') related to BRC, E098 ('Pancreas')
280 related to T2D, and E069 ('Brain cingulate gyrus') related to SCZ as the corresponding
281 tissues of FUN-LDA on the four diseases, and calculated the FUN-LDA scores of the
282 related variants.

283

284 We used one-sided Wilcoxon signed rank test to test whether each algorithm's scores
285 in the positive sample set (e.g. GWAS hits) are significantly higher than the scores in
286 the negative sample set. We observed that TVAR, along with four other methods
287 (CADD, Eigen, DANN, and DeepSEA), all achieved significance in distinguishing positive
288 and negative variants on all data sets (Fig. 3, Table 1). The exception is DANN, which
289 achieved significance on GWAS data of all diseases except SCZ (Fig. 3, Table 1). As a
290 direct comparison, TVAR outperformed all other methods in all of the diseases
291 investigated here. For example, for CAD, TVAR achieved the best performance ($P =$
292 $1.21e-24$), followed by FUN-LDA ($P = 3.25e-11$), DeepSEA ($P = 5.29e-10$), Eigen ($P =$
293 $1.93e-4$), DANN ($P = 1.14E-2$), CADD ($P = 2.30e-2$); for BRC, TVAR also obtained the best
294 performance ($P = 1.60e-21$), followed by DeepSEA ($P = 8.20e-8$), Eigen ($P = 2.02e-7$),
295 DANN ($P = 6.20e-6$), CADD ($P = 5.52e-5$), FUN-LDA ($P = 1.77E-4$); for T2D, TVAR achieved
296 the best performance ($P = 7.01e-20$), followed by DeepSEA ($P = 6.72e-8$), Eigen ($P =$
297 $5.37e-6$), CADD ($P = 7.54e-4$), DANN ($P = 1.80e-2$), FUN-LDA ($P = 1.94E-2$); for SCZ, the
298 performance of TVAR is the best ($P = 1.45e-4$), followed by Eigen ($P = 1.97e-4$), FUN-
299 LDA ($P = 2.59e-4$), DeepSEA ($P = 6.73e-4$), CADD ($P = 9.76e-4$), DANN (not significant).

300 Although the sample sizes of the four test scenarios are on the same scale, we found
301 that the performance advantage of TVAR is considerable higher on CAD, BRC, and T2D,
302 but comparable to other algorithms on SCZ. In fact, the low accuracy group of the TVAR
303 scores contains a large number of brain-related tissues, suggesting that the functional
304 prediction of variants on brain-related tissues is more challenging.

305

306 **Identification of the tissue-specific functionality of rare variants**

307 Since the training set of the TVAR were based significance of eQTLs, which are biased
308 towards common variants owing to their statistical power, we are interested in
309 evaluating TVAR's capability in predicting functionality of rare variants, based on the
310 rationale that common and rare regulatory variants share similar regulatory machinery
311 such that they exist similar functional genomics data across tissues and cell types. We
312 constructed four test sets based on the four previous GWAS of CAD, BRC, T2D, and
313 SCZ³¹⁻³⁴. We constructed the positive sets by including relatively significant rare
314 variants from the four GWAS results with $P < 1e-3$ and MAF < 0.01 . For CAD, BRC, T2D,
315 and SCZ, the numbers of positive sets are 191, 195, 53, 185, respectively. For each test
316 scenario, we randomly selected an equal number of variants from the 1000 Genomes
317 Project as the negative set (MAF < 0.01). Similar to the evaluation for common variants,
318 we used one-sided Wilcoxon signed-rank tests to explore whether each algorithm can
319 distinguish rare variants in the positive and the negative sets. From Fig. 3 and Table 2,
320 we observed that only TVAR could significantly distinguish the positive and negative
321 samples of rare variants for all of the four diseases investigated (P values are $3.76e-3$,

322 2.60e-6, 2.74e-2, 1.61e-2, for CAD, BRC, T2D, and SCZ, respectively). As a direct
323 contrast, all methods (CADD, Eigen, DeepSEA, DANN) that do not distinguish between
324 tissues, failed to separate the positive and negative samples on all data sets (Fig. 3,
325 Table 2). Even FUN-LDA, which has tissue-specific scores, has challenges to achieve a
326 significant discrimination, with T2D being close to significance with $P = 7.63e-2$.

327

328 **Organism-level Functional Predictions with TVAR**

329

330 Although TVAR is a tissue-specific functional variant prediction approach, the G-Score
331 scoring algorithm of TVAR can also be used to predict the variant function at the
332 organism level. To evaluate the performance of the TVAR-GScore algorithm, we
333 constructed four test scenarios: Clinvar (version: 20170130), fine-mapping GWAS
334 hits²⁶, GTEX eQTLs, and MPRA validated variants²⁸. The positive samples of Clinvar
335 include 2,713 clinically validated pathogenic noncoding variants. The positive samples
336 of the fine-mapping GWAS set includes 1,857 fine-mapped GWAS loci across 39
337 phenotypes. The positive samples of the GTEX eQTLs dataset include highly trusted
338 1,137 noncoding variants from the GTEx dataset¹³. The positive samples of the MPRA-
339 validated variants dataset include 234 validated noncoding variants reported based on
340 massively parallel reporter assays. The negative samples of the above four test sets
341 were all randomly selected in equal amounts from the variants reported by 1000
342 Genomes Project (filtered with $MAF > 0.05$). To ensure the unbiasedness of the test, on
343 the four data sets, we also removed the variants that have appeared in the training set
344 of TVAR. We chose four organism level methods as comparison algorithms: CADD,

345 DANN, DeepSEA, and Eigen. On the four data sets, we mainly use ROC curves to
346 evaluate the accuracy of various algorithms (Fig. 4). For comparison, we also give PR
347 curves of each method in the four test scenarios (see Supplementary Note 2).

348 On the ClinVar dataset, TVAR achieved the best performance (AUC = 0.981),
349 followed by DeepSEA (AUC = 0.975), CADD (AUC = 0.963), Eigen (AUC = 0.958), DANN
350 (AUC = 0.957). The AUCs of all algorithms on the ClinVar data set exceed 0.95, and the
351 gap is not large, indicating that the functional characteristics of the clinically validated
352 deleterious variants are clearly distinguishable so all methods have superior
353 performance. On the Refined GWAS Hits dataset, TVAR obtained the best performance
354 (AUC = 0.733), followed by Eigen (AUC = 0.616), DeepSEA (AUC = 0.615), CADD (AUC =
355 0.586), DANN (AUC = 0.545). TVAR has an absolute performance improvement of 11.7%
356 ~ 18.8% on this data set. The Refined GWAS Hits dataset contains many functional risk
357 variants of complex diseases, and the clear advantage of TVAR over other method
358 indicates the potential of TVAR in identifying causal variants for complex diseases. On
359 the GTEX eQTLs dataset, TVAR achieved the best performance (AUC = 0.809), followed
360 by DeepSEA (AUC = 0.622), Eigen (AUC = 0.583), CADD (AUC = 0.555), and DANN (AUC
361 = 0.521). Compared to other methods, TVAR has achieved an absolute performance
362 improvement of 18.7% ~ 28.8%. On the MPRA validated variants dataset, TVAR
363 obtained the best performance (AUC = 0.880), followed by DeepSEA (AUC = 0.765),
364 Eigen (AUC = 0.643), CADD (AUC = 0.627), DANN (AUC = 0.598). TVAR has achieved an
365 absolute performance improvement of 11.5% ~ 28.2% over the competing methods.
366 MPRA has validated the variants on this dataset, and the credibility is high. The AUC of

367 TVAR on this dataset also reaches 0.88, which shows that it can effectively predict
368 MPRA validated functional variants. TVAR's clear margins of improvement over other
369 methods in all scenarios investigated demonstrate its robustness in predicting the
370 functionality of noncoding variants.

371

372 **DISCUSSION**

373 We propose a deep neural network-based framework, TVAR, to evaluate the
374 functionality of variants in non-coding regions of the genome. TVAR supports the
375 prediction of tissue-specific functional effects of non-coding variants through multi-
376 label learning. Simultaneously, we combined the output of the TVAR model and used
377 multi-instance learning to perform TVAR's G-score scoring, which supports organismism
378 level functional prediction across tissues. We constructed a large-scale functional
379 variant dataset across 49 tissues (i.e. eQTLs) and the companion functional
380 annotations, which provides a training set suitable for multi-label learning and
381 supervised learning for generating tissue-specific functional prediction. In this set with
382 more than 70,000 variants as the training set, we used 5-fold cross-validation to train
383 and test the TVAR model. We found that in all GTEx tissues, combined with various
384 annotations as features, TVAR achieve an AUC of 0.77 of the functional variants across
385 all the tissues. In brain-related tissues, functional variants that TVAR can explain are
386 generally small (from 69% to 76%). In Esophagus-related tissues, the percentage of
387 variants that can be explained by TVAR is higher (more than 82%). Evaluating TVAR in
388 tissues relevant to 4 complex disease, i.e. CAD, BRC, T2D, and SCZ, we found that the

389 disease-associated variants have significantly higher TVAR scores than background
390 variants, supporting that the TVAR score can be effective in prioritizing disease-specific
391 functional variants. In the end, the G-Score algorithm based on multi-instance learning
392 achieved excellent performance on four well calibrated functional noncoding variants,
393 clarifying the reliability of the TVAR on the scoring of the functional variants at the
394 organism level.

395 For functional variants scoring, there are already methods that use deep learning
396 technology: DANN adopts the CADD features and training data and uses the DNN
397 model to replace CADD's SVM classification model; both DeepSEA and ExPecto use
398 deep convolutional networks to encode variants and DNA sequences into feature
399 vectors. However, none of them predict variants at the "multi-tissue" level (i.e.,
400 considering the differences and connections among all tissues). To the best of our
401 knowledge, TVAR is the first method that introduces DNN (with multi-label learning)
402 to the multi-tissue level functional variants prediction. We found that in the positive
403 training set of TVAR, more than 30% of the variants have the effect of regulating gene
404 expression in no less than two tissues, and more than 10% of the variants regulate the
405 expression of related genes in no less than four tissues. TVAR uses a large number of
406 fully connected layers in network design. The final output nodes of the model are
407 directly or indirectly related to all the nodes in all the previous layers, which
408 guarantees that TVAR uses non-linear functions to capture the differences and
409 connections among tissues.

410 Predicting the functional consequences of rare noncoding variants has been a

411 challenge. With the use of CRISPR-Cas9 and other technologies, some common
412 variants have been studied. However, the functional effects of rare variants and their
413 regulatory mechanisms in corresponding tissues are still elusive. As shown in previous
414 studies³⁵, randomly selected rare variants also tend to be more likely to be functional
415 than common variants, making it imperative to find real functional rare variants in
416 different tissues to facilitate genetics of common diseases. In this study, in the tissue-
417 specific rare variants test data set, we set the eQTL p -value threshold to $1e-3$ to
418 guarantee an adequate number of variants in the positive sample set. As a
419 consequence, this positive set is more difficult to separate from the negative set.
420 Among all the methods investigated, only TVAR can distinguish the positive and
421 negative sets significantly, demonstrating TVAR's advantages in learning the subtle
422 effects of rare noncoding variants. This benefits from the principle of TVAR's design,
423 i.e. rare regulatory variants share the same genomic profiles as common functional
424 variants, which enables the transfer of learned relationships of genomic features with
425 common eQTLs to effectively predict the function of rare regulatory variants.

426 TVAR's model training between tissues is performed at one time, so the number
427 of tissues supported by the TVAR can be expanded with a small computational cost. A
428 specific example is that the TVAR is updated with the GTEx (V7 to V8) and can also
429 support the corresponding new tissues. The earlier exploration of the TVAR was
430 carried out on GTEx V7 (see Supplementary Note 3). We found that on the 48 tissues
431 of the previous version of GTEx, the performance of the TVAR is almost the same as it
432 is now. The network structure used by TVAR is the same as previously. The

433 performance results illustrate the scalability and stability of the TVAR. In 49 tissues of
434 GTEx (V8), TVAR can achieve about an AUC of 0.77 to classify functional variants.
435 However, we found that there are still many functional variants in brain-related tissues
436 that cannot be explained by the functional annotation of the input features of TVAR.
437 This is not unexpected, however, as brain tissues are highly heterogeneous with
438 diverse cell types with different functions. We also uncovered that the use of multi-
439 instance learning for organism level functional scoring made significant
440 improvements over existing methods in specific scenarios, which revealed a new
441 direction in the development of organism level functional scoring algorithms. In the
442 current G-score algorithm, we have adopted the max function. In the future, we will
443 try to use unsupervised algorithms (such as the DPGMM³⁶ method) to better measure
444 the contribution of different tissues to the organism level score, to optimize the G-
445 score algorithm.
446

447 **METHODS**

448

449 **Overview of TVAR**

450

451 The input of TVAR is 1247-dimensional functional annotations of the variants across
452 multiple databases. The output of TVAR is the corresponding functional scores on the
453 49 tissues in the GTEx project that contain eQTL analysis results (Fig. 5). We preprocess
454 the input features to block the model from overfitting during training. The principal
455 component analysis (PCA) is used to retain the essential information in the input of
456 the TVAR. Next, we obtain the low-dimensional representation of or each variant and
457 apply a deep learning model to classify the preprocessed input. A deep neural network
458 (DNN) is also developed to accomplish the prediction task. The DNN model uses an
459 output layer containing 49 nodes (corresponding to 49 tissues) according to the
460 classification purposes. During the training process, the same variants in the training
461 data are used to represent the functionalities in different tissues, so that they can
462 simultaneously learn whether and in which tissues the variants are functional. TVAR
463 uses multi-label learning rather than multi-class learning methods to predict the
464 functional effects of noncoding variants. TVAR uses multiple fully connected layers and
465 Batch Normalization technique to regularize the mean and variance. The model also
466 uses dropout technique to randomly remove some network nodes and reduce
467 overfitting during network training. The total number of network layers is 23. A well-
468 trained network can learn the patterns of different tissues, and simultaneously output
469 the scores across GTEx 49 tissues (range: 0 ~ 1). To support the scoring in the case of
470 organicism level, we also developed the G-Score method, which connects different

471 scores of variants in different tissues, and finally outputs an overall functional score.

472

473 **The input and features of TVAR**

474

475 Many previous studies have shown that functional annotation across multiple
476 databases is essential for variant scoring systems. To assess the functionality of
477 variants as comprehensively as possible, we integrated three variant-related databases
478 of ENCODE, Roadmap, and FANTOM5. We attached four conservative scores of
479 PhastCons³⁷, PhyloP³⁸, GERP++³⁹ and SiPhy⁴⁰, and Transcription factor (TF) binding
480 sites, super-enhancer regions, transcription start site (TSS), CpG islands and other
481 information to provide functional evidence to support variants. Specifically, from the
482 ENCODE database, we selected the labels related to TF binding sites (132 values),
483 Histone modifications labels (18 values), DNase clusters labels (1 value); from the
484 Roadmap database, we selected labels such as DNase.macs2 (53 values), H2A.Z (23
485 values), H2AK5ac (7 values), H2AK9ac (1 values), H2BK120ac (7 values), H2BK12ac (6
486 values), H2BK15ac (5 values), H2BK20ac (3 values), H2BK5ac (7 values), H3K14ac (6
487 values), H3K18ac (7 values), H3K23ac (7 values), H3K23me2 (2 values), H3K27ac (98
488 values), H3K27me3 (127 values), H3K36me3 (127 values), H3K4ac (7 values),
489 H3K4me1 (127 values), H3K4me2 (24 values), H3K4me3 (126 values), H3K56ac (3
490 values), H3K79me1 (7 values), H3K79me2 (21 values), H3K9ac (62 values), H3K9me1
491 (4 values), H3K9me3 (127 values), H3T11ph (1 values), H4K12ac (1 values), H4K20me1
492 (19 values), H4K5ac (3 values), H4K8ac (7 values), H4K91ac (6 values). From the
493 FANTOM 5 database, we select CAGE peaks (1 value), permissive enhancers (1 value),

494 robust enhancers (1 value). From the 4DGenome⁴¹ database, we choose chromatin
495 interactions detected maker (1 value). From some other related studios, we collect
496 super-enhancers⁴² (1 value), CpG islands maker (1 value), the distance of the variant
497 to the nearest TSS (1 value). In addition, we also integrate PhastCons score (3 values),
498 PhyloP score (3 values), GERP ++ score (2 values), SiPhy score (1 values). Since variants
499 do not necessarily implicate their nearest genes^{43,44}, for each of them, we calculate
500 the average gene expression (49 values) of the two closest genes in 49 tissues.

501 For each variant, there are a total of 1247 input features. Over 1,000 functional
502 annotations enable the TVAR framework to comprehensively learn the functional
503 similarities and differences between variants in different tissues. The total number of
504 tissue-specific features in the TVAR framework is 1098, accounting for 88% of all
505 features. TVAR integrates these tissue-specific features and measures genome-wide
506 functional scores of variants across different tissues with deep learning technology. Of
507 all the inputs, 1188 values are bool-type data. For each variant, we use PCA to reduce
508 these 1188 values to 96 values and merge the other 59 values to obtain a 155-
509 dimensional vector in the training and test sets. Next, we use DNN to model the
510 various input annotations. To construct supervised learning labels on the training set,
511 for all variants in the training set, we determine the training labels of the variants in
512 the dataset based on the q-values (q-value < 0.01) of eQTLs in the GTEx tissues and the
513 LINSIGHT score (top 30%). A variant must meet these two thresholds concurrently to
514 be selected into the training set. Finally, the labels of the training set data are used for
515 the supervised multi-label learning DNN model.

516

517 **The DNN model of TVAR**

518

519 TVAR is designed to be a multi-label learning approach that measures the functional
520 effects of variants across tissues to learn the knowledge of the correlation among
521 tissues. As shown in Fig. 5, the input of the network is vector \mathbf{x} , corresponding to all
522 the functional annotations concatenated, for each variation. This information is passed
523 from the input layer to the output layer. The output is a 49-dimensional vector,
524 corresponding to 49 tissues in GTEx. Firstly, we use multiple fully connected layers to
525 transform the input \mathbf{x} :

$$526 \quad \begin{aligned} z_j &= \sum w_{ij}x_i + b_j \\ y_j &= f(z_j) \end{aligned} \quad (1)$$

527 Where i , j , and k represent the indexes of the nodes of each layer in the network
528 (including the input layer, hidden layer, and output layer). Let x_i be the input of node i ,
529 and y_j be the output of the corresponding hidden node j in the next layer, and z_j be the
530 value after the transformation. The final output layer is a 49-dimensional vector \mathbf{y} , and
531 $f()$ denotes a non-linear activation function. We use the ReLU function between the
532 input layer and the hidden layer (and also between the hidden layers):

$$533 \quad \text{ReLU} = \begin{cases} x & x \geq 0 \\ 0 & x < 0 \end{cases} \quad (2)$$

534 In addition, to prevent a large number of fully connected layers from overfitting the
535 model, we also added batch normalization layers and dropout layers between the fully
536 connected layers. The batch normalization layer performs the following operations on
537 each mini-batch $B = [x, \dots, x_m]$:

538

$$\begin{aligned}
\mu_B &= \frac{1}{m} \sum_{i=1}^m x_i \\
\sigma_B^2 &= \frac{1}{m} \sum_{i=1}^m (x_i - \mu_B)^2 \\
x_i^{norm} &= \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \varepsilon}} \\
y_i &= \gamma x_i^{norm} + \beta
\end{aligned}
\tag{3}$$

539 Where μ_B denotes the mean of the mini-batch, σ_B denotes the variance of the mini-
540 batch, \mathbf{x}_i^{norm} denotes the normalized input, γ and β denote the scale and shift of \mathbf{x}_i^{norm}
541 to get the output \mathbf{y}_i . Similarly, the dropout layer is denoted as:

542

$$\begin{aligned}
r_i &\sim \text{Bernoulli}(p) \\
y_i &= r_i x_i
\end{aligned}
\tag{4}$$

543 Where the Bernoulli function is used to generate a vector \mathbf{r} that is then multiplied to
544 the input \mathbf{x} . The resulting \mathbf{y}_i is set to 0 with p as the probability. Unlike the multi-class
545 classification task with neural networks, the loss function of the DNN model of TVAR
546 is designed to support the Sigmoid Cross-Entropy loss of multi-label learning:

547

$$L = - \frac{\sum_{i=1}^m \sum_{j=1}^n \left[l_{i,j} \log f(x_{i,j}) + (1 - l_{i,j}) \log (1 - f(x_{i,j})) \right]}{mn}
\tag{5}$$

548 Where $x_{i,j}$ denotes the i th output node (m nodes in total) of the j th variant (a total of
549 n variants) in the data set. Let l be the training label of (sample, issue), and L be the
550 total loss, $f()$ is the sigmoid function:

551

$$f(x_{i,j}) = \frac{1}{1 + e^{-x_{i,j}}}
\tag{6}$$

552 The specific network structure of the DNN model is shown in Table S1. The network
553 parameters, learning rate, and optimization methods are fitted during the training
554 process. After the training of the model, for each variant \mathbf{x} , a score vector \mathbf{y}

555 corresponding to the 49 GTEx tissues is generated.

556

557 **The G-score algorithm of TVAR**

558 Although the design of the TVAR is to score the variant-tissue pairs. However, with the

559 development of the G-score algorithm, we can make the TVAR support variant scoring

560 at the organism level. For a score that does not distinguish between tissues, only one

561 score s is needed for each variation. The G-score algorithm is designed as a multi-

562 instance learning approach: the score \mathbf{y} on 49 distinct tissues for each variant \mathbf{x} is

563 regarded as a bag. For bag \mathbf{y} , we use a function $f()$ to find its general score:

$$564 \quad s = f(\mathbf{y}) \quad (7)$$

565 In this study, we simply consider the maximum score among tissues as the G-score for

566 each variant, i.e. $f()$ takes the max function.

567

568 **Availability of Code and Functional Scores of TVAR**

569

570 The TVAR source code and its scores on the ClinVar catalog, fine mapped GWAS Loci,

571 high-confidence eQTLs from GTEx dataset, and MPRA validated functional variants are

572 available at <https://github.com/haiyang1986/TVAR>.

573

574 **Acknowledgements**

575 This work was partially supported by NIH grant U01HG009086 to RC, QWa, JY, QWe,

576 XZ and BL, as well as to HY during the period when HY was a Postdoctoral Fellow at

577 Vanderbilt University.

578

579 **Author contributions**

580 BL and HY designed the study. HY developed the TVAR and implemented it. RC, QWa,
581 JY, QWe and XZ collated functional annotations and contributed to the interpretation
582 of the results. HY and BL wrote the manuscript. All authors read and approved the
583 manuscript.

584

585 **Competing financial interests**

586 None of the authors declare competing financial interests.

587

588

589

590 **Reference**

591

592 1 Human Microbiome Project, C. Structure, function and diversity of the healthy human
593 microbiome. *Nature* **486**, 207-214, doi:10.1038/nature11234 (2012).

594 2 Zhang, F. & Lupski, J. R. Non-coding genetic variants in human disease. *Hum Mol Genet* **24**,
595 R102-110, doi:10.1093/hmg/ddv259 (2015).

596 3 Gloss, B. S. & Dinger, M. E. Realizing the significance of noncoding functionality in clinical
597 genomics. *Exp Mol Med* **50**, 97, doi:10.1038/s12276-018-0087-0 (2018).

598 4 Gulko, B. & Siepel, A. An evolutionary framework for measuring epigenomic information and
599 estimating cell-type-specific fitness consequences. *Nature genetics* **51**, 335-342,
600 doi:10.1038/s41588-018-0300-z (2019).

601 5 Skipper, M., Dhand, R. & Campbell, P. Presenting ENCODE. *Nature* **489**, 45,
602 doi:10.1038/489045a (2012).

603 6 Leung, D. *et al.* Integrative analysis of haplotype-resolved epigenomes across human tissues.
604 *Nature* **518**, 350-354, doi:10.1038/nature14217 (2015).

605 7 Lizio, M. *et al.* Gateways to the FANTOM5 promoter level mammalian expression atlas. *Genome*
606 *biology* **16**, 22, doi:10.1186/s13059-014-0560-6 (2015).

607 8 Kircher, M. *et al.* A general framework for estimating the relative pathogenicity of human
608 genetic variants. *Nature genetics* **46**, 310-315, doi:10.1038/ng.2892 (2014).

609 9 Ritchie, G. R., Dunham, I., Zeggini, E. & Flicek, P. Functional annotation of noncoding sequence
610 variants. *Nat Methods* **11**, 294-296, doi:10.1038/nmeth.2832 (2014).

611 10 Quang, D., Chen, Y. & Xie, X. DANN: a deep learning approach for annotating the pathogenicity
612 of genetic variants. *Bioinformatics* **31**, 761-763, doi:10.1093/bioinformatics/btu703 (2015).

613 11 Huang, Y. F., Gulko, B. & Siepel, A. Fast, scalable prediction of deleterious noncoding variants
614 from functional and population genomic data. *Nature genetics*, doi:10.1038/ng.3810 (2017).

615 12 Ionita-Laza, I., McCallum, K., Xu, B. & Buxbaum, J. D. A spectral approach integrating functional
616 genomic annotations for coding and noncoding variants. *Nature genetics* **48**, 214-220,
617 doi:10.1038/ng.3477 (2016).

618 13 Yang, H. *et al.* De novo pattern discovery enables robust assessment of functional
619 consequences of non-coding variants. *Bioinformatics* **35**, 1453-1460,
620 doi:10.1093/bioinformatics/bty826 (2019).

621 14 Bodea, C. A. *et al.* PINES: phenotype-informed tissue weighting improves prediction of
622 pathogenic noncoding variants. *Genome biology* **19**, 173, doi:10.1186/s13059-018-1546-6
623 (2018).

624 15 Janssens, A. C. *et al.* The impact of genotype frequencies on the clinical validity of genomic
625 profiling for predicting common chronic diseases. *Genet Med* **9**, 528-535,
626 doi:10.1097/gim.0b013e31812eece0 (2007).

627 16 Hannon, E. *et al.* Methylation QTLs in the developing brain and their enrichment in
628 schizophrenia risk loci. *Nature Neuroscience* **19**, 48-+, doi:10.1038/nn.4182 (2016).

629 17 Parker, S. C. J. *et al.* Chromatin stretch enhancer states drive cell-specific gene regulation and
630 harbor human disease risk variants. *Proceedings of the National Academy of Sciences of the*
631 *United States of America* **110**, 17921-17926, doi:10.1073/pnas.1317023110 (2013).

632 18 Roadmap Epigenomics, C. *et al.* Integrative analysis of 111 reference human epigenomes.
633 *Nature* **518**, 317-330, doi:10.1038/nature14248 (2015).

634 19 Consortium, G. T. *et al.* Genetic effects on gene expression across human tissues. *Nature* **550**,
635 204-213, doi:10.1038/nature24277 (2017).

636 20 Lee, D. *et al.* A method to predict the impact of regulatory variants from DNA sequence. *Nature*
637 *genetics* **47**, 955-961, doi:10.1038/ng.3331 (2015).

638 21 He, Z. H., Liu, L. X., Wang, K. & Ionita-Laza, I. A semi-supervised approach for predicting cell-
639 type specific functional consequences of non-coding variation using MPRA. *Nature*
640 *Communications* **9**, doi:10.1038/s41467-018-07349-w (2018).

641 22 Backenroth, D. *et al.* FUN-LDA: A Latent Dirichlet Allocation Model for Predicting Tissue-Specific
642 Functional Effects of Noncoding Variation: Methods and Applications. *Am J Hum Genet* **102**,
643 920-942, doi:10.1016/j.ajhg.2018.03.026 (2018).

644 23 Zhou, J. *et al.* Deep learning sequence-based ab initio prediction of variant effects on
645 expression and disease risk. *Nature genetics* **50**, 1171-1179, doi:10.1038/s41588-018-0160-6
646 (2018).

647 24 Landrum, M. J. *et al.* ClinVar: public archive of interpretations of clinically relevant variants.
648 *Nucleic acids research* **44**, D862-868, doi:10.1093/nar/gkv1222 (2016).

649 25 Landrum, M. J. *et al.* ClinVar: public archive of relationships among sequence variation and
650 human phenotype. *Nucleic acids research* **42**, D980-D985, doi:10.1093/nar/gkt1113 (2014).

651 26 Farh, K. K. *et al.* Genetic and epigenetic fine mapping of causal autoimmune disease variants.
652 *Nature* **518**, 337-343, doi:10.1038/nature13835 (2015).

653 27 Consortium, G. T. Human genomics. The Genotype-Tissue Expression (GTEx) pilot analysis:
654 multitissue gene regulation in humans. *Science* **348**, 648-660, doi:10.1126/science.1262110
655 (2015).

656 28 Tewhey, R. *et al.* Direct Identification of Hundreds of Expression-Modulating Variants using a
657 Multiplexed Reporter Assay. *Cell* **165**, 1519-1529, doi:10.1016/j.cell.2016.04.027 (2016).

658 29 Bulik-Sullivan, B. K. *et al.* LD Score regression distinguishes confounding from polygenicity in
659 genome-wide association studies. *Nature genetics* **47**, 291-295, doi:10.1038/ng.3211 (2015).

660 30 Genomes Project, C. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74,
661 doi:10.1038/nature15393 (2015).

662 31 Nikpay, M. *et al.* A comprehensive 1,000 Genomes-based genome-wide association meta-
663 analysis of coronary artery disease. *Nature genetics* **47**, 1121-1130, doi:10.1038/ng.3396
664 (2015).

665 32 Xue, A. *et al.* Genome-wide association analyses identify 143 risk variants and putative
666 regulatory mechanisms for type 2 diabetes. *Nat Commun* **9**, 2941, doi:10.1038/s41467-018-
667 04951-w (2018).

668 33 Pardini, A. F. *et al.* Common schizophrenia alleles are enriched in mutation-intolerant genes
669 and in regions under strong background selection. *Nature genetics* **50**, 381-389,
670 doi:10.1038/s41588-018-0059-2 (2018).

671 34 Michailidou, K. *et al.* Genome-wide association analysis of more than 120,000 individuals
672 identifies 15 new susceptibility loci for breast cancer. *Nature genetics* **47**, 373-380,
673 doi:10.1038/ng.3242 (2015).

674 35 Yang, H. *et al.* De Novo pattern discovery enables robust assessment of functional
675 consequences of noncoding variants. *Bioinformatics*, doi:10.1093/bioinformatics/bty826
676 (2018).

677 36 Yang, H., Wei, Q., Zhong, X., Yang, H. & Li, B. Cancer driver gene discovery through an

678 integrative genomics approach in a non-parametric Bayesian framework. *Bioinformatics* **33**,
679 483-490, doi:10.1093/bioinformatics/btw662 (2017).

680 37 Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast
681 genomes. *Genome research* **15**, 1034-1050, doi:10.1101/gr.3715005 (2005).

682 38 Cooper, G. M. *et al.* Distribution and intensity of constraint in mammalian genomic sequence.
683 *Genome research* **15**, 901-913, doi:10.1101/gr.3577405 (2005).

684 39 Davydov, E. V. *et al.* Identifying a high fraction of the human genome to be under selective
685 constraint using GERP++. *PLoS Comput Biol* **6**, e1001025, doi:10.1371/journal.pcbi.1001025
686 (2010).

687 40 Garber, M. *et al.* Identifying novel constrained elements by exploiting biased substitution
688 patterns. *Bioinformatics* **25**, i54-62, doi:10.1093/bioinformatics/btp190 (2009).

689 41 Teng, L., He, B., Wang, J. & Tan, K. 4DGenome: a comprehensive database of chromatin
690 interactions. *Bioinformatics* **31**, 2560-2564, doi:10.1093/bioinformatics/btv158 (2015).

691 42 Hnisz, D. *et al.* Super-enhancers in the control of cell identity and disease. *Cell* **155**, 934-947,
692 doi:10.1016/j.cell.2013.09.053 (2013).

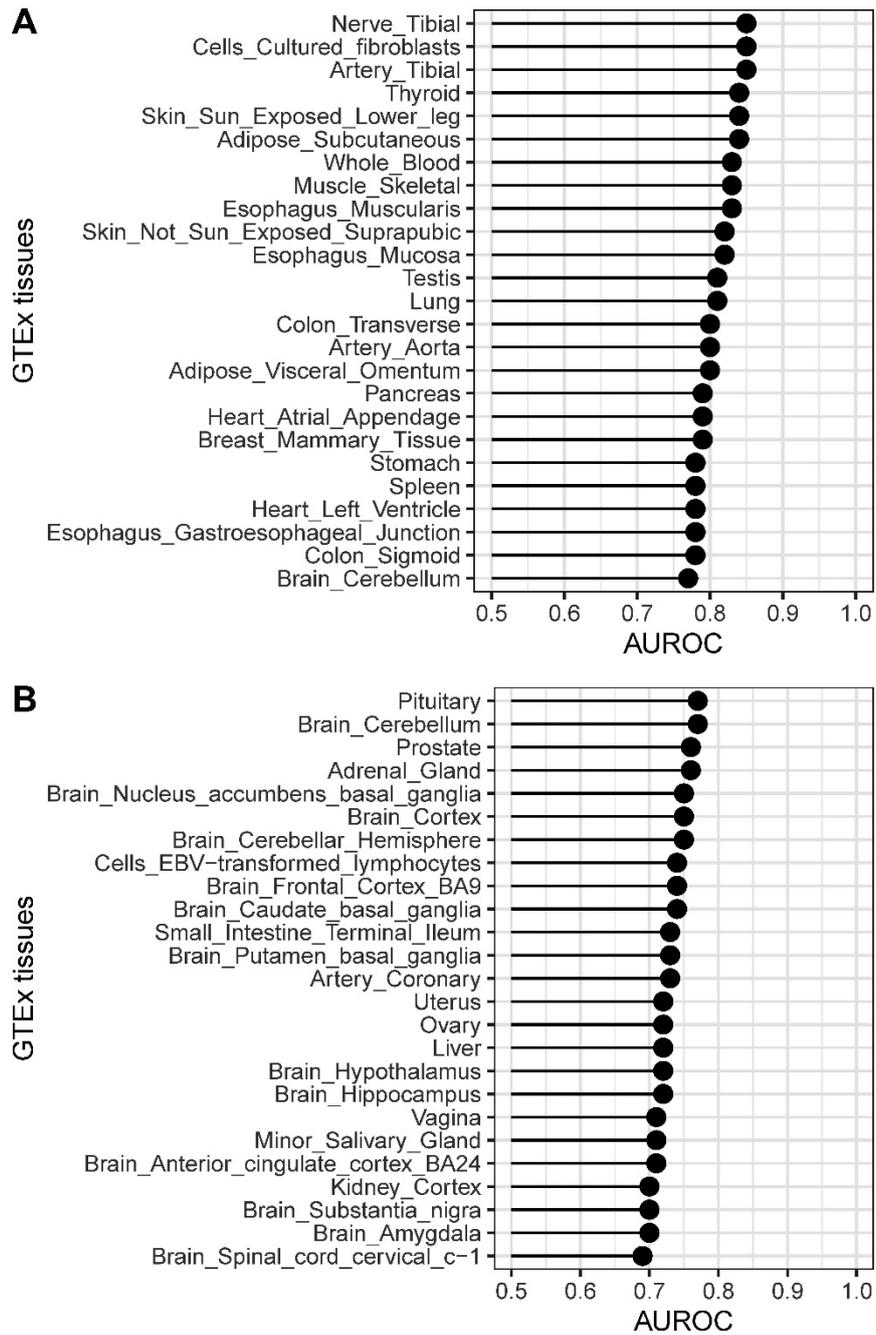
693 43 Wang, Q. *et al.* A Bayesian framework that integrates multi-omics data and gene networks
694 predicts risk genes from schizophrenia GWAS data. *Nat Neurosci* **22**, 691-699,
695 doi:10.1038/s41593-019-0382-7 (2019).

696 44 Smemo, S. *et al.* Obesity-associated variants within FTO form long-range functional
697 connections with IRX3. *Nature* **507**, 371-375, doi:10.1038/nature13138 (2014).

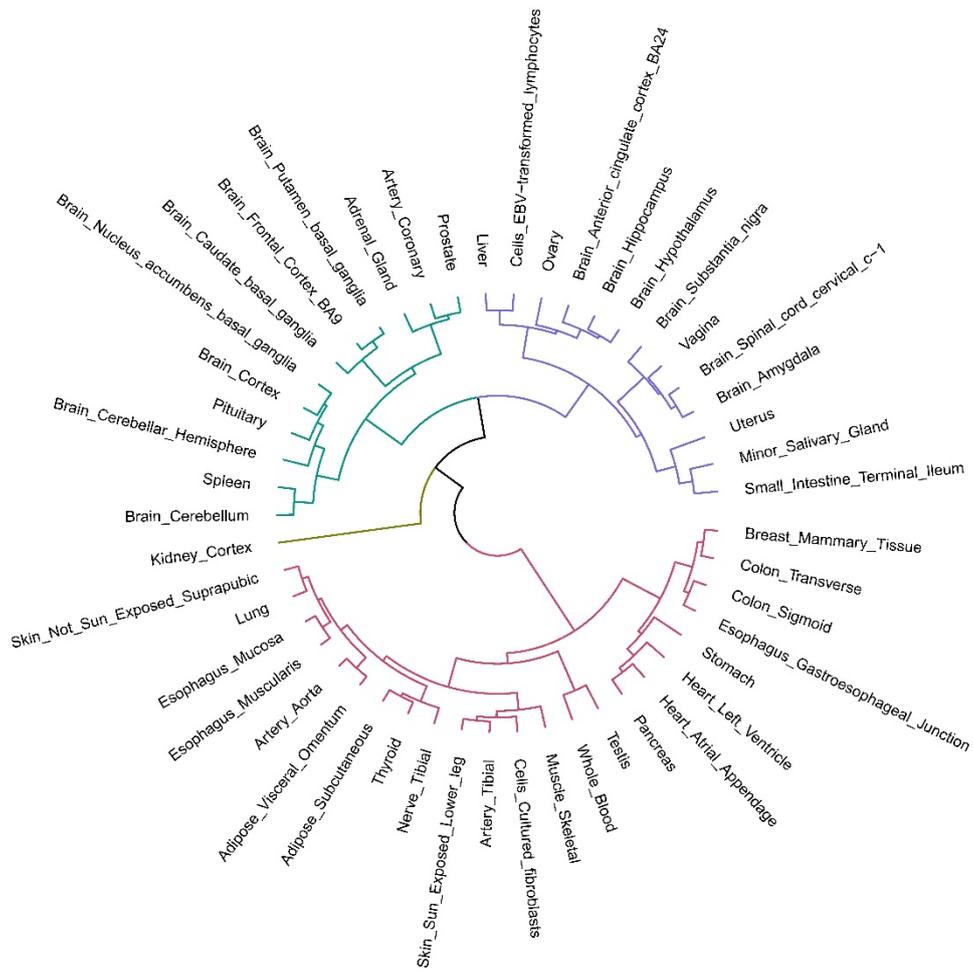
698

699

700 FIGURES
701

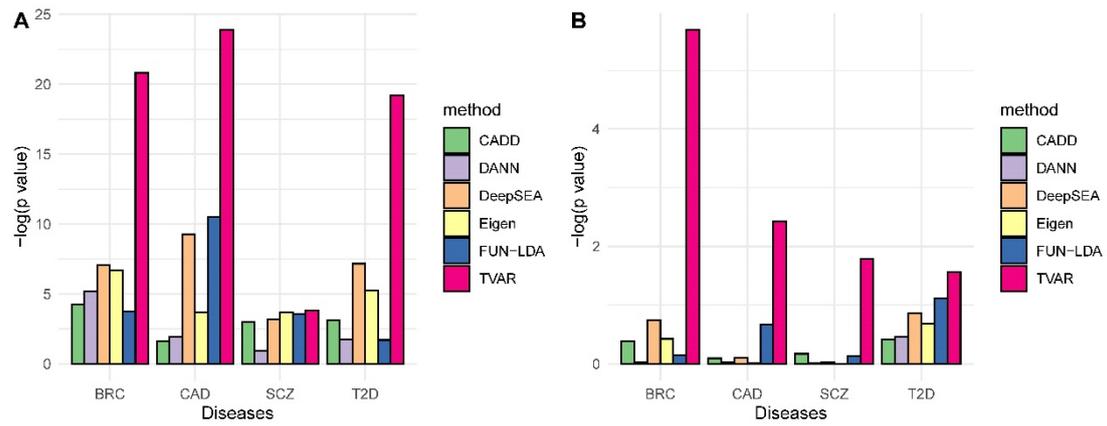


702
703 Fig. 1 Five-Fold cross-validation performance of TVAR on the GTEx 49 tissues: A). High
704 accuracy group. B) Low accuracy group.
705



706
 707
 708
 709

Fig. 2: The hierarchical clustering results of TVAR's functional scores across the GTEx 49 tissues (from the GTEx dataset V8).

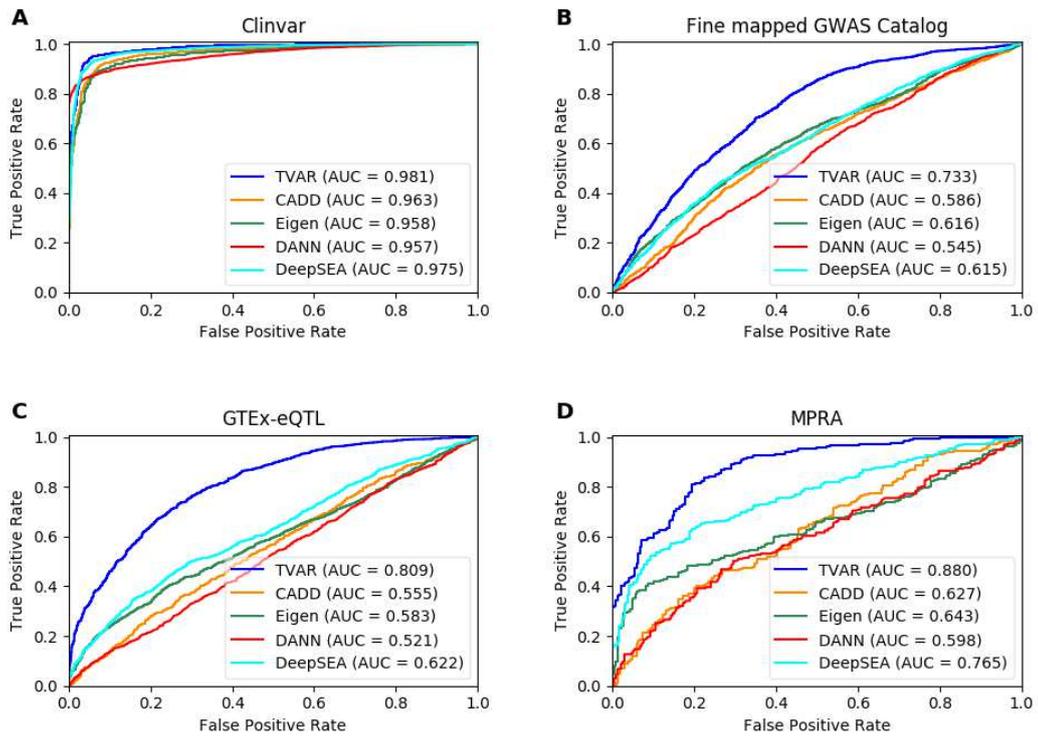


710

711 Fig. 3 Performance comparison of TVAR, CADD, DANN, DeepSEA, Eigen, FUN-LDA on
 712 datasets of four complex diseases: coronary artery disease (CAD), breast cancer (BRC),
 713 Type 2 diabetes (T2D), and Schizophrenia (SCZ): (A) GWAS hit data set; (B) rare variant
 714 data set. The prediction power was evaluated based on the $-\log_{10}(p\text{-value})$.

715

716

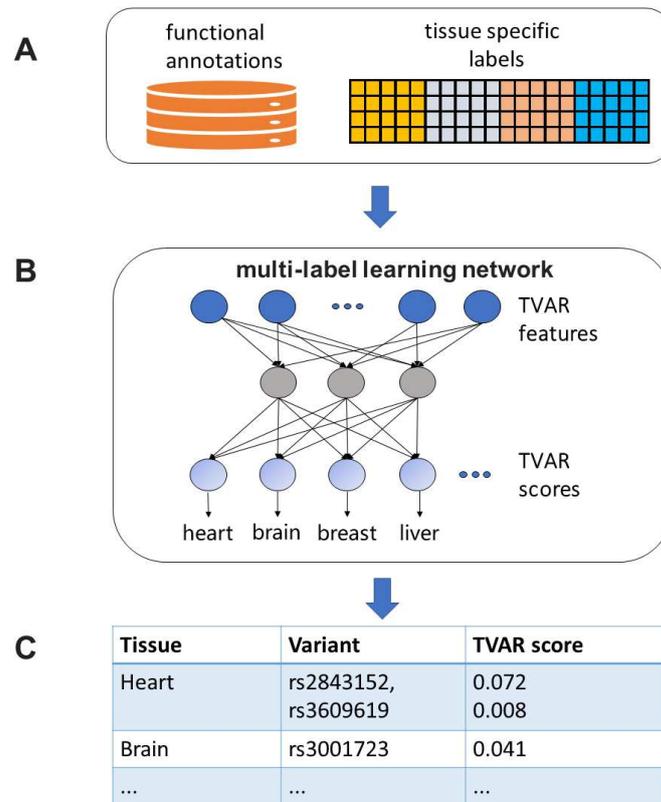


717

718 Fig. 4 Performance comparison of TVAR, CADD, DANN, DeepSEA, and Eigen on the
 719 extended dataset: A) clinically significant variants in Clinvar database, B) fine-mapped
 720 GWAS variants), C) GTEx-eQTLs, and D) MPRA validated variants. The identification
 721 accuracy is counted with the area under the ROC curves.

722

723



724

725

726

727

728

Fig. 5 Overview of the TVAR framework. (A) Take the functional annotations as the input of the deep learning network. (B) Multi-label learning with DNN. (C) Output the functional scores across multiple tissues.

729 Table 1 The p-values (Wilcoxon signed-rank test, one-sided with 'greater') for the
 730 analysis of common variants on four diseases (coronary artery disease, breast cancer,
 731 Type 2 diabetes, and Schizophrenia):
 732

Disease	N	Score	p-value
coronary artery disease (cardiovascular)	342	TVAR	1.21E-24
		CADD	2.30E-2
		Eigen	1.93E-4
		DANN	1.14E-2
		DeepSEA	5.29E-10
		FUN-LDA	3.25E-11
breast cancer (cancer)	354	TVAR	1.60E-21
		CADD	5.52E-5
		Eigen	2.02E-7
		DANN	6.20E-6
		DeepSEA	8.20E-08
		FUN-LDA	1.77E-4
Type 2 diabetes (metabolic disease)	210	TVAR	7.01E-20
		CADD	7.54E-4
		Eigen	5.37E-6
		DANN	1.80E-2
		DeepSEA	6.72E-8
		FUN-LDA	1.94E-2
Schizophrenia (mental)	286	TVAR	1.45E-4
		CADD	9.76E-4
		Eigen	1.97E-4
		DANN	1.13E-1
		DeepSEA	6.73E-4
		FUN-LDA	2.59E-4

733
 734

735 Table 2. The p-values (Wilcoxon signed-rank test, one-sided with 'greater') for the
 736 analysis of rare variants on four diseases (coronary artery disease, breast cancer, Type
 737 2 diabetes, and Schizophrenia):
 738

Disease	N	Score	p-value
coronary artery disease (cardiovascular)	191	TVAR	3.76E-3
		CADD	8.04E-1
		Eigen	9.71E-1
		DANN	9.26E-1
		DeepSEA	7.92E-1
		FUN-LDA	2.16E-1
breast cancer (cancer)	195	TVAR	2.06E-6
		CADD	4.14E-1
		Eigen	3.78E-1
		DANN	9.39E-1
		DeepSEA	1.79E-1
		FUN-LDA	7.19E-1
Type 2 diabetes (metabolic disease)	53	TVAR	2.74E-2
		CADD	3.86E-1
		Eigen	2.05E-1
		DANN	3.44E-1
		DeepSEA	1.36E-1
		FUN-LDA	7.63E-2
Schizophrenia (mental)	189	TVAR	1.61E-2
		CADD	6.74E-1
		Eigen	9.93E-1
		DANN	9.69E-1
		DeepSEA	9.33E-1
		FUN-LDA	7.42E-1

739

Figures

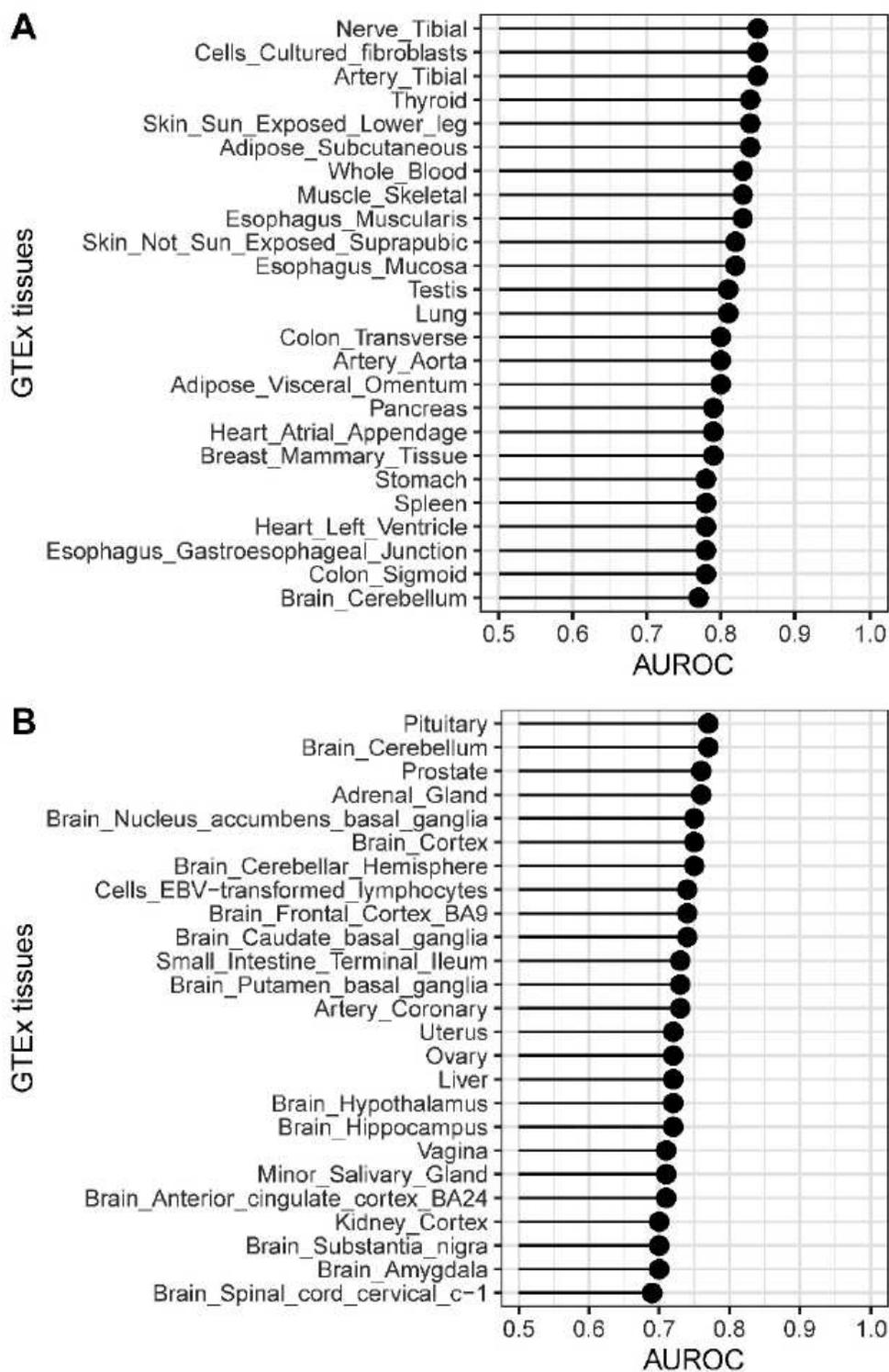


Figure 1

Five-Fold cross-validation performance of TVAR on the GTEx 49 tissues: A). High accuracy group. B) Low accuracy group.

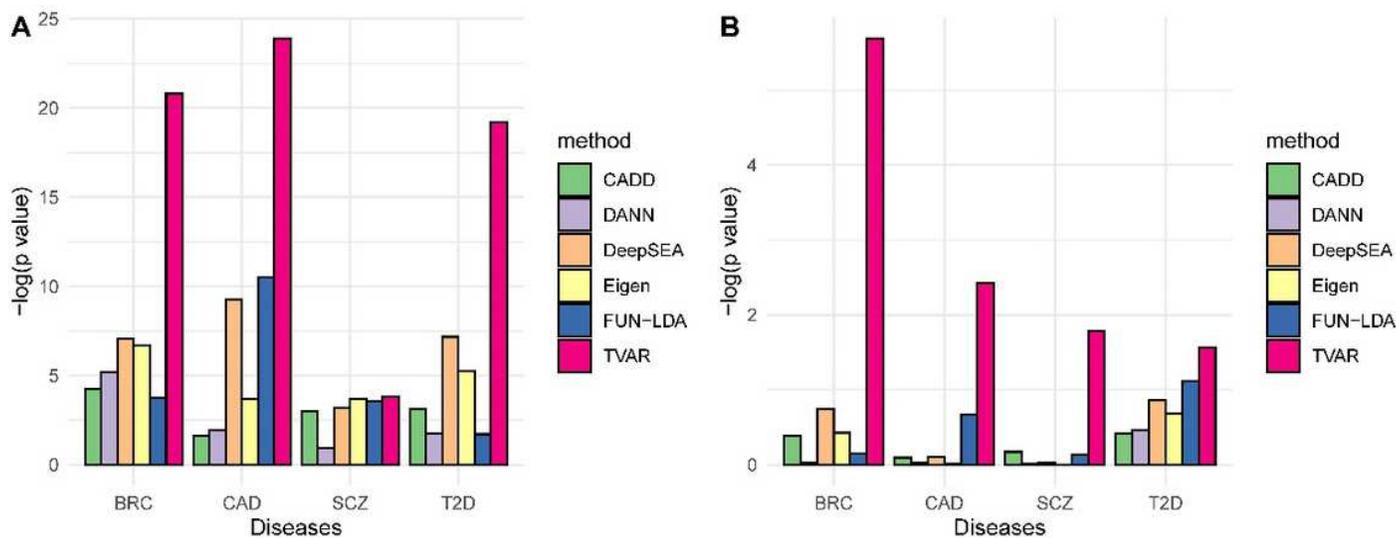


Figure 3

Performance comparison of TVAR, CADD, DANN, DeepSEA, Eigen, FUN-LDA on datasets of four complex diseases: coronary artery disease (CAD), breast cancer (BRC), Type 2 diabetes (T2D), and Schizophrenia (SCZ): (A) GWAS hit data set; (B) rare variant data set. The prediction power was evaluated based on the $-\log_{10}(\text{p value})$.

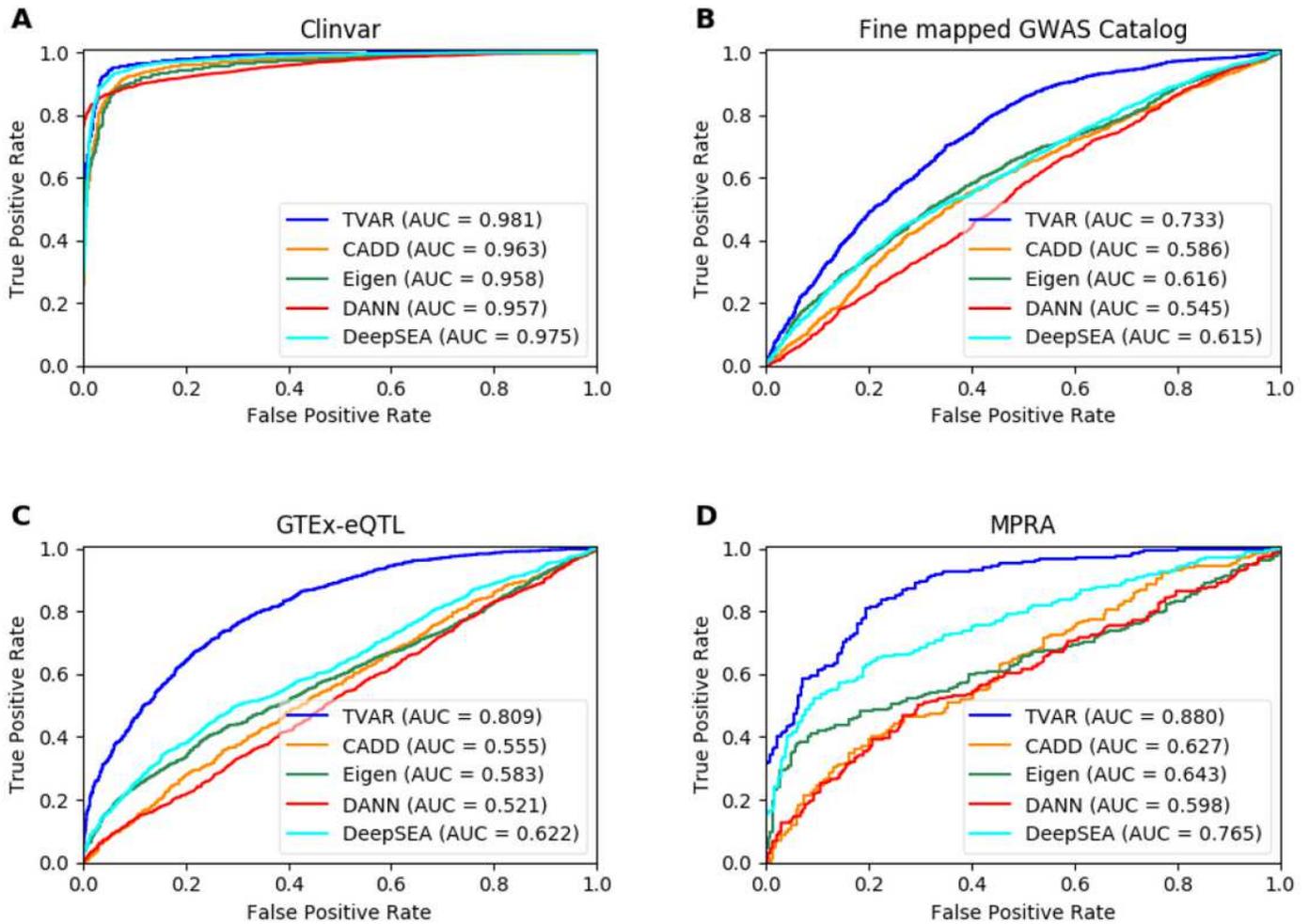


Figure 4

Performance comparison of TVAR, CADD, DANN, DeepSEA, and Eigen on the extended dataset: A) clinically significant variants in Clinvar database, B) fine-mapped GWAS variants), C) GTEx-eQTLs, and D) MPRA validated variants. The identification accuracy is counted with the area under the ROC curves.

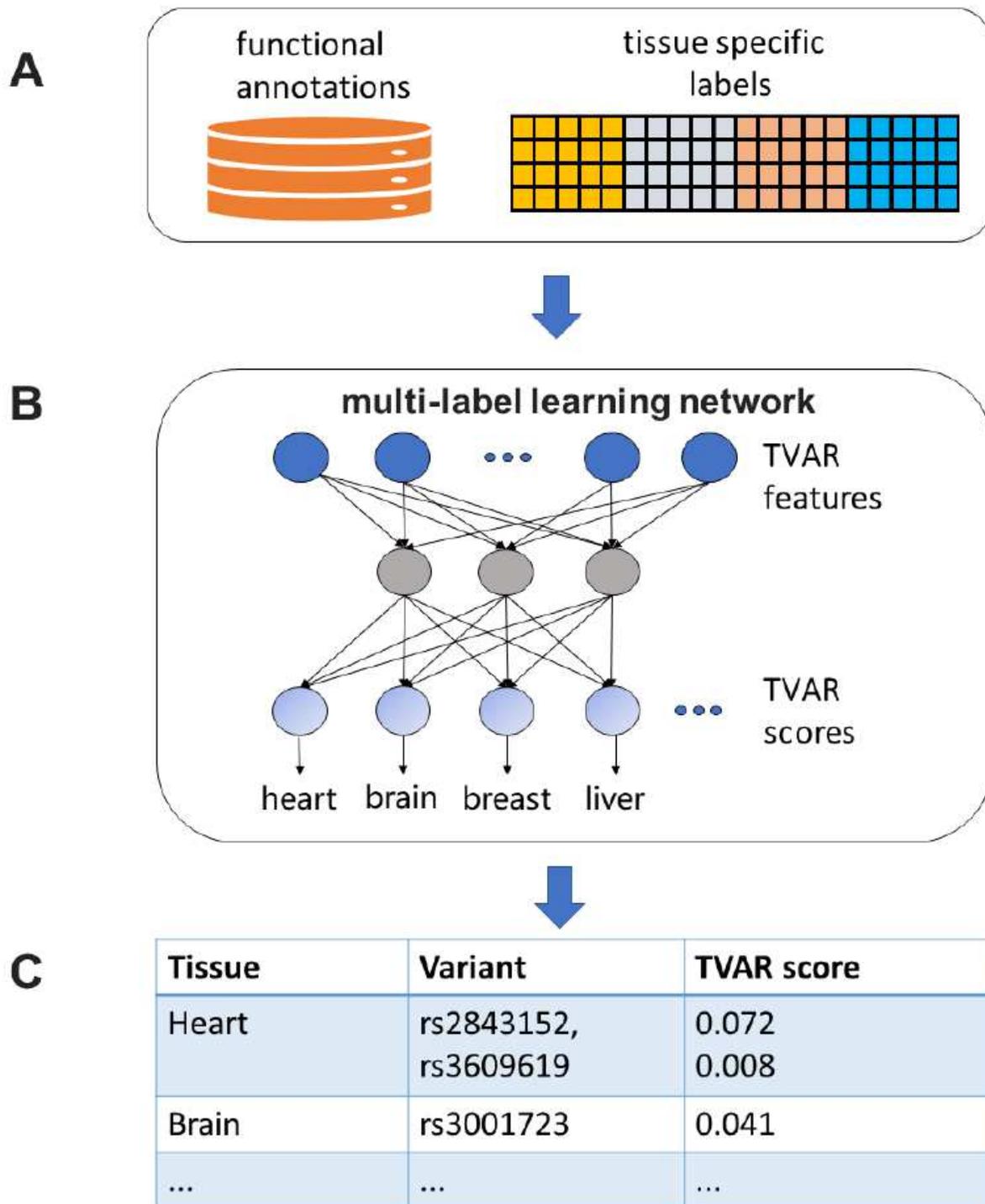


Figure 5

Overview of the TVAR framework. (A) Take the functional annotations as the input of the deep learning network. (B) Multi-label learning with DNN. (C) Output the functional scores across multiple tissues.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [02TVARsupplementarymaterials20201121.pdf](#)