

1 Exploring the phenotypic consequences of tissue specific gene 2 expression variation inferred from GWAS summary statistics

3 Alvaro N. Barbeira¹, Scott P. Dickinson¹, Jason M. Torres², Jiamao Zheng¹, Eric S. Torstenson³,
4 Heather E. Wheeler⁴, Kaanan P. Shah¹, Rodrigo Bonazzola¹, Tzintzuni Garcia⁵, Todd Edwards³,
5 GTEx Consortium, Dan L. Nicolae¹, Nancy J. Cox³, Hae Kyung Im^{1,*}

6 1 Section of Genetic Medicine, The University of Chicago, Chicago, IL, USA

7 2 Committee on Molecular Metabolism and Nutrition, The University of Chicago,
8 Chicago, IL, USA

9 3 Vanderbilt Genetic Institute, Vanderbilt University Medical Center, Nashville, TN, USA

10 4 Departments of Biology and Computer Science, Loyola University Chicago, Chicago, IL,
11 USA

12 5 Center for Research Informatics, The University of Chicago, IL, USA

13 * E-mail: Corresponding haky@uchicago.edu

14 Abstract

15 Scalable, integrative methods to understand mechanisms that link genetic variants with phenotypes are
16 needed. Here we derive a mathematical expression to compute PrediXcan (a gene mapping approach)
17 results using summary data (S-PrediXcan) and show its accuracy and general robustness to misspecified
18 reference sets. We apply this framework to 44 GTEx tissues and 100+ phenotypes from GWAS and
19 meta-analysis studies, creating a growing public catalog of associations that seeks to capture the effects
20 of gene expression variation on human phenotypes. Replication in an independent cohort is shown. Most
21 of the associations were tissue specific, suggesting context specificity of the trait etiology. Colocalized
22 significant associations in unexpected tissues underscore the need for an agnostic scanning of multiple
23 contexts to improve our ability to detect causal regulatory mechanisms. Monogenic disease genes are
24 enriched among significant associations for related traits, suggesting that smaller alterations of these
25 genes may cause a spectrum of milder phenotypes.

26 Introduction

27 Over the last decade, GWAS have been successful in robustly associating genetic loci to human com-
28 plex traits. However, the mechanistic understanding of these discoveries is still limited, hampering the
29 translation of the associations into actionable targets. Studies of enrichment of expression quantitative
30 trait loci (eQTLs) among trait-associated variants [1–3] show the importance of gene expression regula-
31 tion. Functional class quantification showed that 80% of the common variant contribution to phenotype
32 variability in 12 diseases can be attributed to DNAase I hypersensitivity sites, further highlighting the
33 importance of transcript regulation in determining phenotypes [4].

34 Many transcriptome studies have been conducted where genotypes and expression levels are assayed
35 for a large number of individuals [5–8]. The most comprehensive transcriptome dataset, in terms of
36 examined tissues, is the Genotype-Tissue Expression Project (GTEx): a large-scale effort where DNA
37 and RNA were collected from multiple tissue samples from nearly 1000 individuals and sequenced to high
38 coverage [9,10]. This remarkable resource provides a comprehensive cross-tissue survey of the functional
39 consequences of genetic variation at the transcript level.

40 To integrate knowledge generated from these large-scale transcriptome studies and shed light on
41 disease biology, we developed PrediXcan [11], a gene-level association approach that tests the mediating
42 effects of gene expression levels on phenotypes. PrediXcan is implemented on GWAS or sequencing studies
43 (i.e. studies with genome-wide interrogation of DNA variation and phenotypes). It imputes transcriptome
44 levels with models trained in measured transcriptome datasets (e.g. GTEx). These predicted expression
45 levels are then correlated with the phenotype in a gene association test that addresses some of the key
46 limitations of GWAS [11].

47 Meta-analysis efforts that aggregate results from multiple GWAS have been able to identify an in-
48 creasing number of associations that were not detected with smaller sample sizes [12–14]. We will refer
49 to these results as GWAMA (Genome-wide association meta-analysis) results. In order to harness the
50 power of these increased sample sizes while keeping the computational burden manageable, methods that
51 use summary level data rather than individual level data are needed.

52 Methods similar to PrediXcan that estimate the association between intermediate gene expression
53 levels and phenotypes, but use summary statistics have been reported: TWAS (summary version) [15]
54 and SMR (Summary Mendelian Randomization) [16]. Another class of methods that integrate eQTL

55 information with GWAS results are based on colocalization of eQTL and GWAS signals. Colocalized
56 signals provide evidence of possible causal relationship between the target gene of an eQTL and the
57 complex trait. These include RTC [1], Sherlock [17], COLOC [18], and more recently eCAVIAR [19] and
58 ENLOC [20].

59 Here we derive a mathematical expression that allows us to compute the results of PrediXcan without
60 the need to use individual-level data, greatly expanding its applicability. We compare with existing
61 methods and outline a best practices framework to perform integrative gene mapping studies, which we
62 term MetaXcan.

63 We apply the MetaXcan framework by first training over 1 million elastic net prediction models of
64 gene expression traits, covering protein coding genes across 44 human tissues from GTEx, and then
65 performing gene-level association tests over 100 phenotypes from 40 large meta-analysis consortia and
66 dbGaP.

67 Results

68 Computing PrediXcan results using summary statistics

69 We have derived an analytic expression to compute the outcome of PrediXcan using only summary
70 statistics from genetic association studies. Details of the derivation are shown in the Methods section. In
71 Figure 1-A, we illustrate the mechanics of Summary-PrediXcan (S-PrediXcan) in relation to traditional
72 GWAS and the individual-level PrediXcan method [11].

73 We find high concordance between PrediXcan and S-PrediXcan results indicating that in most cases,
74 we can use the summary version without loss of power to detect associations. Figure 2 shows the
75 comparison of PrediXcan and S-PrediXcan Z-scores for a simulated phenotype (under the null hypothesis),
76 a cellular growth phenotype and two disease phenotypes: type 1 diabetes and bipolar disorder from the
77 WTCCC Consortium [21]. For the simulated phenotype, the study sets (in which GWAS is performed)
78 and the reference set (in which LD between SNPs is computed) were African, East Asian, and European
79 from 1000 Genomes. The training set (in which prediction models are trained) was European (DGN
80 Cohort [5]) in all cases. The high correlation between PrediXcan and S-PrediXcan demonstrates the
81 robustness of our method to mismatches between reference and study sets. Despite the generally good
82 concordance between the summary and individual level methods, there were a handful of false positive

83 results with S-PrediXcan much more significant than PrediXcan. This underscores the need to use closely
84 matched LD information whenever possible. Supplementary Figure 11 shows S-PrediXcan's performance
85 on a phenotype simulated under the alternative hypothesis.

86 Notice that we are not testing here whether PrediXcan itself is robust to population differences
87 between training and study sets. Robustness of the prediction across populations has been previously
88 reported [22]. We further corroborated this in Supplementary Figure 10.

89 Next we compare with other summary result-based methods such as S-TWAS, SMR, and COLOC.

90 Colocalization estimates complement PrediXcan results

91 One class of methods seeks to determine whether eQTL and GWAS signals are colocalized or are distinct
92 although linked by LD. This class includes COLOC [18], Sherlock [17], and RTC [1], and more recently
93 eCAVIAR [19], and ENLOC [20]. Thorough comparison between these methods can be found in [18,19].
94 HEIDI, the post filtering step in SMR that estimates heterogeneity of GWAS and eQTL signals, can
95 be included in this class. We focus here on COLOC, whose quantification of the probability of five
96 configurations complements well with S-PrediXcan results.

97 COLOC provides the probability of 5 hypotheses: H0 corresponds to no eQTL and no GWAS associ-
98 ation, H1 and H2 correspond to association with eQTL but no GWAS or vice-versa, H3 corresponds to
99 eQTL and GWAS association but independent signals, and finally H4 corresponds to shared eQTL and
100 GWAS association. P0, P1, P2, P3, and P4 are the corresponding probabilities for each configuration.
101 The sum of the five probabilities is 1. The authors [18] recommend to interpret H0, H1, and H2 as limited
102 power; we will aggregate these three hypothesis into one event with probability $1-P3-P4$ for convenience.

103 Figure 3 shows ternary plots [23] with P3, P4, and $1-P3-P4$ as vertices. The blue region, top sub-
104 triangle, corresponds to high probability of colocalized eQTL and GWAS signals (P4). The orange region,
105 bottom left, corresponds to high probability of distinct eQTL and GWAS signals (P3). The gray region,
106 center and bottom right, corresponds to low probability of both colocalization and independent signals.

107 Figure 3-B shows association results for all gene-tissue pairs with the height phenotype. We find that
108 most associations fall in the gray, 'undetermined', region. When we restrict the plot to S-PrediXcan
109 Bonferroni-significant genes (Figure 3-C), three distinct peaks emerge in the high P4 region ('colocalized
110 signals'), high P3 region ('independent signals' or 'non-colocalized signals'), and 'undetermined' region.
111 Moreover, when genes with low prediction performance are excluded (Supplementary Figure 6-D) the

112 ‘undetermined’ peak significantly diminishes.

113 These clusters provide a natural way to classify significant genes and complement S-PrediXcan results.
114 Depending on false positive/false negative trade-off choices, genes in the ‘independent signals’ or both
115 ‘independent signals’ and ‘undetermined’ can be filtered out. The proportion of colocated associations
116 ($P_4 > 0.5$) ranged from 5 to 100% depending on phenotype with a median of 27.6%. The proportion of
117 ‘non-colocalized’ associations ranged from 0 to 77% with a median of 27.0%. Supplementary Table 2
118 summarizes the percentages of significant associations that fall into the different colocalization regions.

119 This post-filtering idea was first implemented in the SMR approach using HEIDI. Comparison of
120 COLOC results with HEIDI is shown in Figure 6-E to -F.

121 **Comparison of S-PrediXcan to S-TWAS**

122 Gusev et al. have proposed Transcriptome-Wide Association Study based on summary statistics (S-
123 TWAS), which imputes the SNP level Z-scores into gene level Z-scores. This is not the same as computing
124 the results of individual level TWAS. We show (in Methods section) that the difference between the
125 individual level and summary level TWAS is given by the factor $\sqrt{\frac{1-R_l^2}{1-R_g^2}}$, where R_l is the proportion of
126 variance in the phenotype explained by a SNP’s allelic dosage, and R_g is the proportion explained by
127 gene expression (see Methods section). For most practical purposes we have found that this factor is very
128 close to 1 so that if the same prediction models were used, no substantial difference between S-TWAS
129 and S-PrediXcan should be expected.

130 Figure 4-A shows a diagram of S-PrediXcan and S-TWAS. Both use SNP to phenotype associations
131 results ($Z_{X,Y}$) and prediction weights (w_{X,T_g}) to infer the association between the gene expression level
132 (T_g) and phenotype (Y).

133 Figure 4-B compares S-TWAS significance (as reported in [24]) to S-PrediXcan significance. The
134 difference between the two approaches is mostly driven by the different prediction models: TWAS uses
135 BSLMM [25] whereas PrediXcan uses elastic net [26]. BSLMM allows two components: one sparse (small
136 set of large effect predictors) and one polygenic (all variants contribute some marginal effect to the
137 prediction). For PrediXcan we have chosen to use a sparse model (elastic net) based on the finding that
138 the genetic component of gene expression levels is mostly sparse [27].

139 Figure 4-C shows that the proportion of non-colocalized (independent) GWAS and eQTL signals is
140 larger among TWAS significant genes than among S-PrediXcan significant ones. This is likely due to

141 the polygenic component of BSLMM models, a wider set of SNPs increasing the chance of capturing
142 LD-contaminated (non-colocalized) association. Figure 4-D shows that, for most traits, the proportion
143 of colocalized signals is larger among S-PrediXcan significant genes than S-TWAS significant genes.

144 **Comparison of S-PrediXcan to SMR**

145 Zhu et al. have proposed Summary Mendelian Randomization (SMR) [16], a summary data based
146 Mendelian randomization that integrates eQTL results to determine target genes of complex trait-
147 associated GWAS loci. They derive an approximate χ_1^2 -statistic (Eq 5 in [16]) for the mediating effect of
148 the target gene expression on the phenotype.

149 Unfortunately, the derived statistic is mis-calibrated. A QQ plot comparing the SMR statistic (under
150 the null hypothesis of genome-wide significant eQTL signal and no GWAS association) shows deflation.
151 The sample mean of the statistic is ≈ 0.93 instead of 1, the expected value for the mean of a χ_1^2 random
152 variable. See Figure 5 (B and C) and Methods section for details. The χ^2 approximation is only valid
153 in two extreme cases: when the eQTL association is much stronger than the GWAS association or vice
154 versa, when the GWAS association is much stronger than the eQTL association. See Methods for details.

155 One limitation is that the significance of the SMR statistic is the lower of the top eQTL association
156 (genotype to expression) or the GWAS association (genotype to phenotype) as shown in Figure 5 (E and
157 F). Given the much larger sample sizes of GWAS studies, for most genes, the combined significance will be
158 determined by the eQTL association. The combined statistic forces us to apply multiple testing correction
159 for all genes, even those that are distant to GWAS associated loci, which is unnecessarily conservative.
160 Keep in mind that currently both SMR and PrediXcan only use cis associations. An example may clarify
161 this further. Let us suppose that for a given phenotype there is only one causal SNP and that the GWAS
162 yielded a highly significant p-value, say 10^{-50} . Let us also suppose that there is only one gene (gene A) in
163 the vicinity (we are only using cis predictors) associated with the causal SNP with $p=10^{-5}$. SMR would
164 compute the p-values of all genes and yield a p-value $\approx 10^{-5}$ for gene A (the less significant p-value).
165 However, after multiple correction this gene would not be significantly associated with the phenotype.
166 Here it is clear that we should not be adjusting for testing of all genes when we know a priori that only
167 one is likely to produce a gene level association. In contrast, the PrediXcan p-value would be $\approx 10^{-50}$
168 for gene A and would be distributed uniformly from 0 to 1 for the remaining genes. Most likely only
169 gene A (or perhaps a handful of genes, just by chance) would be significant after Bonferroni correction.

170 If we further correct for prediction uncertainty (here =eQTL association), a p-value of $\approx 10^{-5}$ would
171 remain significant since we only need to correct for the (at most) handful of genes that were Bonferroni
172 significant for the PrediXcan p-value.

173 Another potential disadvantage of this method is that only top-eQTLs are used for testing the gene
174 level association. This does not allow to aggregate the effect on the gene across multiple variants.

175 Figure 5-D compares S-PrediXcan (elastic net) and SMR association results. As expected, SMR p-
176 values tend to be less significant than S-PrediXcan's in large part due to the additional adjustment for
177 the uncertainty in the eQTL association. Figures 5-E and -F show that the SMR significance is bounded
178 by the eQTL and GWAS association strengths of the top eQTL.

179 SMR introduces a post filtering step via an approach called HEIDI, which is compared to COLOC in
180 Figure 3 and Supplementary Figure 6.

181 **MetaXcan framework**

182 Building on S-PrediXcan and existing approaches, we define a general framework (MetaXcan) to integrate
183 eQTL information with GWAS results and map disease-associated genes. This evolving framework can
184 incorporate models and methods to increase the power to detect causal genes and filter out false positives.
185 Existing methods fit within this general framework as instances or components (Figure 6).

186 The framework starts with the training of prediction models for gene expression traits followed by
187 a selection of high-performing models. Next, a mathematical operation is performed to compute the
188 association between each gene and the downstream complex trait. Additional adjustment for the uncer-
189 tainty in the prediction model can be added. To avoid capturing LD-contaminated associations, which
190 can occur when expression predictor SNPs and phenotype causal SNPs are different but in LD, we use
191 colocalization methods that estimate the probability of shared or independent signals.

192 PrediXcan implementations use elastic net models motivated by our observation that gene expression
193 variation is mostly driven by sparse components [27]. TWAS implementations have used Bayesian Sparse
194 Linear Mixed Models [25] (BSLMM). SMR fits into this scheme with prediction models consisting solely
195 of the top eQTL for each gene (weights are not necessary here since only one SNP is used at a time).

196 For the last step, we chose COLOC to estimate the probability of colocalization of GWAS and eQTL
197 signals. COLOC probabilities cluster more distinctly into different classes and thus, unlike other methods,
198 suggests a natural cut off threshold at $P=0.5$. Another advantage of COLOC is that for genes with low

199 probability of colocalization, it further distinguishes distinct GWAS and eQTL signals from low power.
200 This is a useful feature that future development of colocalization methods should also offer. SMR, on the
201 other hand, uses its own estimate of ‘heterogeneity’ of signals calculated by HEIDI.

202 **Suggested association analysis pipeline**

- 203 1. Perform PrediXcan or S-PrediXcan using all tissues. Use Bonferroni correction for all gene-tissue
204 pairs: keep $p < 0.05 / \text{number of gene-tissue pairs tested}$.
- 205 2. Keep associations with significant prediction performance adjusting for number of PrediXcan sig-
206 nificant gene-tissue pairs: keep prediction performance p-values $< 0.05 / (\text{number of significant as-}$
207 $\text{sociations from previous step})$.
- 208 3. Filter out LD-contaminated associations, i.e. gene-tissue pairs in the ‘independent signal’ (=‘non-
209 colocalized’) region of the ternary plot (See Figure 3-A): keep $\text{COLOC } P3 < 0.5$ (Blue and gray
210 regions in Figure 3-A).
- 211 4. If further reduction of number of genes to be taken to replication or validation is desired, keep only
212 hits with explicit evidence of colocalization: $P4 > 0.5$ (Blue region in Figure 3-A).

213 Any choice of thresholds has some level of arbitrariness. Depending on the false positive and negative
214 trade off, these number may be changed.

215 **Gene expression variation in humans is associated to diverse phenotypes**

216 We downloaded summary statistics from meta analyses of over 100 phenotypes from 40 consortia. The
217 full list of consortia and phenotypes is shown in Supplementary Table 4. We tested association between
218 these phenotypes and the predicted expression levels using elastic net models in 44 human tissues from
219 GTEx as described in the Methods section, and a whole blood model from the DGN cohort presented
220 in [11].

221 We used a Bonferroni threshold accounting for all the gene-tissue pairs that were tested ($0.05 / \text{total}$
222 $\text{number of gene-tissue pairs} \approx 2.5e-7$). This approach is conservative because the correlation between
223 tissues would make the total number of independent tests smaller than the total number of gene-tissue
224 pairs. Height had the largest number of significantly associated unique genes at 1,686 (based on a

225 GWAMA of 250K individuals). Other polygenic diseases with a large number of associations include
226 schizophrenia with 305 unique significant genes ($n = 150\text{K}$ individuals), low-density lipoprotein cholesterol
227 (LDL-C) levels with 296 unique significant genes ($n = 188\text{K}$), other lipid levels, glycemic traits, and
228 immune/inflammatory disorders such as rheumatoid arthritis and inflammatory bowel disease. For other
229 psychiatric phenotypes, a much smaller number of significant associations was found, with 8 significant
230 genes for bipolar disorder ($n = 16,731$) and one for major depressive disorder ($n = 18,759$), probably due
231 to smaller sample sizes, but also smaller effect sizes.

232 When step 2 from the suggested pipeline is applied, keeping only reliably predicted genes, we are left
233 with 739 genes for height, 150 for schizophrenia, 117 for LDL-C levels.

234 After step 3, which keeps genes that are without strong evidence of LD-contamination, these numbers
235 dropped to 264 for height, 58 for schizophrenia, and 60 for LDL-C levels. After step 4, which keeps only
236 genes with strong evidence of colocalization, we find 215 genes for height, 49 for schizophrenia, and 35
237 for LDL-C. The counts for the full set of phenotypes can be found in Supplementary Table 4.

238 Mostly, genome-wide significant genes tend to cluster around known SNP-level genome-wide signifi-
239 cant loci or sub-genome-wide significant loci. Regions with sub-genome-wide significant SNPs can yield
240 genome-wide significant results in S-PrediXcan, because of the reduction in multiple testing and the in-
241 crease in power arising from the combined effects of multiple variants. Supplementary Table 3 lists a few
242 examples where this occurs.

243 The proportion of colocalized associations ($P_4 > 0.5$) ranged from 5 to 100% depending on phenotype
244 with a median of 27.6%. The proportion of ‘non colocalized’ associations ranged from 0 to 77% with a
245 median of 27.0%.

246 See full set of results in our online catalog (gene2pheno.org). Significant gene-tissue pairs are included
247 in Supplementary Table 5. To facilitate comparison, the catalog contains all SMR results we generated
248 and the S-TWAS results reported by [24] for 30 GWAS traits and GTEx BSLMM models. Note that
249 SMR application to 28 phenotypes was reported by [28] using whole blood eQTL results from [29].

250 **Moderate changes in ClinVar gene expression is associated with milder phenotypes**

251 We reasoned that if complete knock out of monogenic disease genes cause severe forms of the disease, more
252 moderate alterations of gene expression levels (as effected by regulatory variation in the population) could
253 cause more moderate forms of the disease. Thus moderate alterations in expression levels of monogenic

254 disease genes (such as those driven by eQTLs) may have an effect on related complex traits, and this effect
255 could be captured by S-PrediXcan association statistics. To test this hypothesis, we obtained genes listed
256 in the ClinVar database [30] for obesity, rheumatoid arthritis, diabetes, Alzheimer's, Crohn's disease,
257 ulcerative colitis, age-related macular degeneration, schizophrenia, and autism. Figure 8 displays the QQ
258 plot for all associations and those in ClinVar database. As postulated, we found enrichment of significant
259 S-PrediXcan associations for ClinVar genes for all tested phenotypes except for autism and schizophrenia.
260 The lack of significance for autism is probably due to insufficient power: the distribution of p-values is
261 close to the null distribution. In contrast, for schizophrenia, many genes were found to be significant
262 in the S-PrediXcan analysis. There are several reasons that may explain this lack of enrichment: genes
263 identified with GWAS and subsequently with S-PrediXcan have rather small effect sizes, so that it would
264 not be surprising that they were missed until very large sample sizes were aggregated; ClinVar genes may
265 originate from rare mutations that are not well covered by our prediction models, which are based on
266 common variation (due to limited sample sizes of eQTL studies and the minor allele frequency -MAF-
267 filter used in GWAS studies); or the mechanism of action of the schizophrenia linked ClinVar genes may
268 be different than the alteration of expression levels. Also, the pathogenicity of some of the ClinVar entries
269 has been questioned [31]. The list of diseases in ClinVar used to generate the enrichment figures can be
270 found in Supplementary Table 1, along with the corresponding association results.

271 **Agnostic scanning of a broad set of tissues enabled by GTEx improves discovery**

272 Most genes were found to be significantly associated in a handful of tissues as illustrated in Figure 9-B. For
273 example, for LDL-C levels, liver was the most enriched tissue in significant associations as expected given
274 known biology of this trait. (See Supplementary Figure 5). This prominent role of liver was apparent
275 despite the smaller sample size available for building liver models ($n=97$), which was less than a third of
276 the numbers available for muscle ($n=361$) or lung ($n=278$).

277 However, in general, tissues expected to stand out as more enriched for diseases given currently
278 known biology did not consistently do so when we looked at the average across all (significant) genes,
279 using various measures of enrichment. For example, the enrichment in liver was less apparent for high-
280 density lipoprotein cholesterol (HDL-C) or triglyceride levels. We find for many significant associations
281 that the evidence is present across multiple tissues. This may be caused by a combination of context
282 specificity and sharing of regulatory mechanism across tissues.

283 Next, we illustrate the challenges of identifying disease relevant tissues based on eQTL information
284 using three genes with well established biology: *C4A* for schizophrenia [32] and *SORT1* [33] and *PCSK9*
285 both for LDL-C and cardiovascular disease. S-PrediXcan results for these genes and traits, and regulatory
286 activity by tissue (as measured by the proportion of expression explained by the genetic component), are
287 shown in Figure 9-A. Representative results are shown in Supplementary Tables 6, 7, and 8. Supple-
288 mentary Table 9 contains the full set MetaXcan results (i.e. association, colocalization, HEIDI) for these
289 genes.

290 *SORT1* is a gene with strong evidence for a causal role in LDL-C levels, and as a consequence, is
291 likely to affect risk for cardiovascular disease [33]. This gene is most actively regulated in liver (close to
292 50% of the expression level of this gene is determined by the genetic component) with the most significant
293 S-PrediXcan association in liver (p-value ≈ 0 , $Z = -28.8$), consistent with our prior knowledge of lipid
294 metabolism. In this example, tissue specific results suggest a causal role of *SORT1* in liver.

295 However, in the following example, association results across multiple tissues do not allow us to
296 discriminate the tissue of action. *C4A* is a gene with strong evidence of causal effect on schizophrenia
297 risk via excessive synaptic pruning in the brain during development [32]. Our results show that *C4A* is
298 associated with schizophrenia risk in all tissues (p < 2.5×10^{-7} in 36 tissue models and p < 0.05 for the
299 remaining 4 tissue models).

300 *PCSK9* is a target of several LDL-C lowering drugs currently under trial to reduce cardiovascular
301 events [34]. The STARNET study [35] profiled gene expression levels in cardiometabolic disease patients
302 and showed tag SNP rs12740374 to be a strong eQTL for *PCSK9* in visceral fat but not in liver. Consistent
303 with this, our S-PrediXcan results also show a highly significant association between *PCSK9* and LDL-C
304 (p $\approx 10^{-13}$) in visceral fat and not in liver (our training algorithm did not yield a prediction model for
305 *PCSK9*, i.e. there was no evidence of regulatory activity). In our results, however, the statistical evidence
306 is much stronger in tibial nerve (p $\approx 10^{-27}$). Accordingly, in our training set (GTEx), there is much
307 stronger evidence of regulation of this gene in tibial nerve compared to visceral fat.

308 Most associations highlighted here have high colocalization probabilities. See Supplementary tables
309 6, 7, and 8. However, visceral fat association shows evidence of non colocalization (probability of inde-
310 pendent signals P3=0.69 in LDL-C). It is possible that the relevant regulatory activity in visceral adipose
311 tissue was not detected in the GTEx samples for various reasons but it was detected in tibial nerve.
312 Thus by looking into all tissues' results we increase the window of opportunities where we can detect the

313 association.

314 *PCSK9* yields colocalized signals for LDL-C levels in Tibial Nerve, Lung and Whole blood. *SORT1*
315 shows colocalization with LDL-C in liver ($P4 \approx 1$) and pancreas ($P4 = 0.90$). *C4A* is colocalized with
316 schizophrenia risk for the majority of the tissues (29/40) with a median colocalization probability of 0.82.

317 These examples demonstrate the power of studying regulation in a broad set of tissues and contexts
318 and emphasize the challenges of determining causal tissues of complex traits based on in-silico analysis
319 alone. Based on these results, we recommend to scan all tissues' models to increase the chances to
320 detect the relevant regulatory mechanism that mediates the phenotypic association. False positives can
321 be controlled by Bonferroni correcting for the additional tests.

322 **Replication in an independent cohort, GERA**

323 We used data from the Resource for Genetic Epidemiology Research on Adult Health and Aging study
324 (GERA, phs000674.v1.p1) [36, 37]. This is a study led by the Kaiser Permanente Research Program on
325 Genes, Environment, and Health (RPGEH) and the UCSF Institute for Human Genetics with over 100,000
326 participants. We downloaded the data from dbGaP and performed GWAS followed by S-PrediXcan
327 analysis of 22 conditions available in the European subset of the cohort.

328 For replication, we chose Coronary Artery Disease (CAD), LDL cholesterol levels, Triglyceride levels,
329 and schizophrenia, which had closely related phenotypes in the GERA study and had a sufficiently large
330 number of Bonferroni significant associations in the discovery set. Analysis and replication of the type
331 2 diabetes phenotype can be found in [38]. Coronary artery disease hits were compared with 'Any
332 cardiac event', LDL cholesterol and triglyceride level signals were compared with 'Dyslipidemia', and
333 schizophrenia was compared to 'Any psychiatric event' in GERA.

334 High concordance between discovery and replication is shown in Figure 10 where dyslipidemia as-
335 sociation Z-scores are compared to LDL cholesterol Z-scores. The majority of gene-tissue pairs (92%,
336 among the ones with Z-score magnitude greater than 2 in both sets) have concordant direction of effects
337 in the discovery and replication sets. The high level of concordance is supportive of an omnigenic trait
338 architecture [39]

339 Following standard practice in meta-analysis, we consider a gene to be replicated when the following
340 three conditions are met: the p-value in the replication set is < 0.05 , the direction of discovery and
341 replication effects are the same, and the meta analyzed p-value is Bonferroni significant with the discovery

Discovery phenotype	Replication phenotype	# signif genes in disc set	# replicated genes	π_1 (all) in repl	π_1 (sig) in repl	% replicated genes	# replicated coloc or undeterm
Coronary artery disease	Any cardiac event	56	6	0.4%	49.1%	10.7%	6
LDL cholesterol	Dyslipidemia	282	219	5.8%	90.8%	78.5%	184
Triglycerides	Dyslipidemia	233	100	5.8%	73.1%	43.5%	69
Schizophrenia	Any psychiatric event	285	60	1.2%	47.6%	21.1%	51

Table 1. Replication of results in GERA. Significant genes/tissue pairs were replicated using a closely matched phenotype in an independent dataset from the GERA cohort [36]. The criteria consisted in significance threshold for replication at $p < 0.05$, concordant directions of effect, and meta analysis p-value less than the Bonferroni threshold in the discovery set. π_1 is an estimate of proportion of true positives in the replication set. π_1 (all) uses all gene-tissue pairs whereas π_1 (sig) is computed using only gene-tissue pairs that were significant in the discovery set. The column ‘# replicated genes coloc or undeterm’ is the number of replicated genes excluding the ones for which there was strong evidence of independent GWAS and eQTL signals.

342 threshold.

343 Among the 56 genes significantly associated with CAD in the discovery set, 6 (11%) were significantly
344 associated with ‘Any cardiac event’ in GERA. Using ‘Dyslipidemia’ as the closest matching phenotype,
345 78.5% and 43.5% of LDL and triglyceride genes replicated, respectively. Among the 285 genes associated
346 with schizophrenia in the discovery set, 51 (21%) replicated. The low replication rate for CAD and
347 Schizophrenia is likely due to the broad phenotype definitions in the replication.

348 We found no consistent replication pattern difference between colocalized and non-colocalized genes.
349 This is not unexpected if the LD pattern is similar between discovery and replication sets.

350 The full list of significant genes can be queried in gene2pheno.org.

351 Discussion

352 Here we derive a mathematical expression to compute PrediXcan results without using individual level
353 data, which greatly expands its applicability and is robust to study and reference set mismatches. This
354 has not been done before. TWAS, which for the individual level approach only differs from PrediXcan on
355 the prediction model used in the implementation, has been extended to use summary level data. When
356 Gaussian imputation is used, the relationship between individual level and summary versions of TWAS
357 is clear. This is not the case when extended to general weights (such as BSLMM). Our mathematical

358 derivation shows the analytic difference between them explicitly.

359 The larger proportion of non-colocalized signals from TWAS suggests that by using BSLMM (with
360 polygenic component), it could be more susceptible to LD-contamination than PrediXcan, which uses
361 elastic net (a sparser model). Improved colocalization methods without single causal variant assumption
362 may be needed to strengthen this argument. But the predominantly sparse genetic architecture of gene
363 expression traits [27] supports the benefit of elastic net over BSLMM predictors. We show that SMR
364 statistics needs to be calibrated and argue that by combining the eQTL and GWAS uncertainties into
365 one statistic, it forces the user to apply multiple correction that may be unnecessarily conservative.

366 We also add a post filtering step, to mitigate issues with LD-contamination. Based on consistency
367 with PrediXcan and interpretability of results, we have chosen to use COLOC for filtering. However,
368 colocalization estimation is an active area of research and improved versions or methods will be adopted
369 in the future.

370 Despite the generally good concordance between the summary and individual level methods, there
371 were a handful of false positive results with S-PrediXcan much more significant than PrediXcan. This
372 underscores the need to use closely matched LD information whenever possible.

373 We applied our framework to over 100 phenotypes using transcriptome prediction models trained in
374 44 tissue from the GTEx Consortium and generated a catalog of downstream phenotypic association
375 results of gene expression variation, a growing resource for the community.

376 The enrichment of monogenic disease genes among related phenotype associations suggests that mod-
377 erate alteration of expression levels as effected by common genetic variation may cause a continuum of
378 phenotypic changes. Alternatively, a more complex interplay between common and rare variation could
379 be taking place such as higher tolerance to loss of function mutations in lower expressing haplotypes
380 which could induce association with predicted expression.

381 We are finding that most trait associations are tissue specific; i.e. they are detected in a handful
382 of tissues. However, we also find that expected tissues given known biology do not necessarily rank
383 among the top enriched tissues. This suggests context specificity of the pathogenic mechanism; specific
384 developmental stage or environmental conditions may be necessary to detect the regulatory event. On
385 the other hand, we are detecting associations in unexpected tissues which suggests a sharing of regulation
386 across multiple tissues/contexts or perhaps novel biology that takes place in these tissues. In either case,
387 agnostic scanning of a broad set of tissues is necessary to discover these mechanisms.

388 Software and Resources

389 We make our software publicly available on a GitHub repository: <https://github.com/hakyimlab/>
390 **MetaXcan**. Prediction model weights and covariances for different tissues can be downloaded from Pre-
391 dictDB. A short working example can be found on the GitHub page; more extensive documentation can
392 be found on the project's wiki page. The results of MetaXcan applied to the 44 human tissues and a
393 broad set of phenotypes can be queried on gene2pheno.org.

394 Methods

395 Summary-PrediXcan formula

396 Figure 1-B shows the main analytic expression used by Summary-PrediXcan for the Z-score (Wald statis-
397 tic) of the association between predicted gene expression and a phenotype. The input variables are the
398 weights used to predict the expression of a given gene, the variance and covariances of the markers in-
399 cluded in the prediction, and the GWAS coefficient for each marker. The last factor in the formula can
400 be computed exactly in principle, but we would need additional information that is unavailable in typical
401 GWAS summary statistics output such as phenotype variance and sample size. Dropping this factor from
402 the formula does not affect the accuracy of the results as demonstrated in the close to perfect concordance
403 between PrediXcan and Summary-PrediXcan results on the diagonal of Figure 2-A.

404 The approximate formula we use is:

$$Zg \approx \sum_{l \in \text{Model}_g} w_{lg} \frac{\hat{\sigma}_l}{\hat{\sigma}_g} \frac{\hat{\beta}_l}{\text{se}(\hat{\beta}_l)} \quad (1)$$

405 where

- 406 • w_{lg} is the weight of SNP l in the prediction of the expression of gene g ,
- 407 • $\hat{\beta}_l$ is the GWAS regression coefficients for SNP l ,
- 408 • $\text{se}(\hat{\beta}_l)$ is standard error of $\hat{\beta}_l$,
- 409 • $\hat{\sigma}_l$ is the estimated variance of SNP l , and
- 410 • $\hat{\sigma}_g$ is the estimated variance of the predicted expression of gene g ,

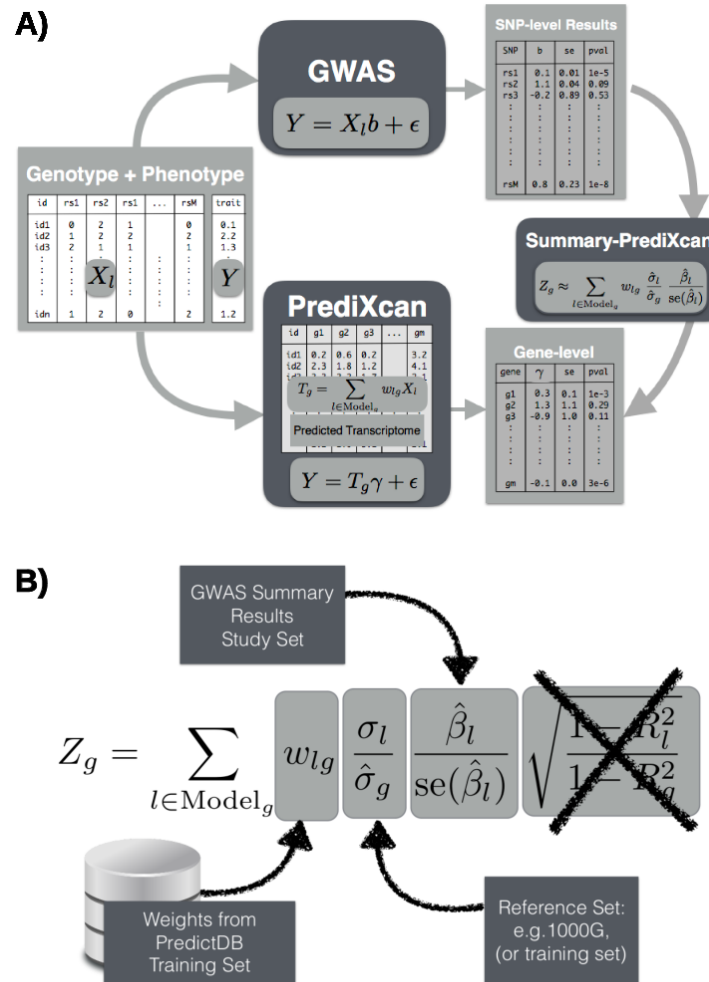


Figure 1. Panel A: Comparison between GWAS, PrediXcan, and Summary-PrediXcan. Both GWAS and PrediXcan take genotype and phenotype data as input. GWAS computes the regression coefficients of Y on X_l using the model $Y = a + X_l b + \epsilon$, where Y is the phenotype and X_l the individual SNP dosage. The output is a table of SNP-level results. PrediXcan, in contrast, starts first by predicting/imputing the transcriptome. Then it calculates the regression coefficients of the phenotype Y on each gene's predicted expression T_g . The output is a table of gene-level results. Summary-PrediXcan directly computes the gene-level association results using the output from GWAS. **Panel B: Components of the S-PrediXcan formula.** This figure shows the components of the formula to calculate PrediXcan gene-level association results using summary statistics. The different sets involved as input data are shown. The regression coefficient between the phenotype and the genotype is obtained from the study set. The training set is the reference transcriptome dataset where the prediction models of gene expression levels are trained. The reference set (1000G, or training set having some advantages) is used to compute the variances and covariances (LD structure) of the markers used in the predicted expression levels. Both the reference set and training set values are pre-computed and provided to the user so that only the study set results need to be provided to the software. The crossed out term was set to 1 as an approximation. We found this approximation to have negligible impact on the results.

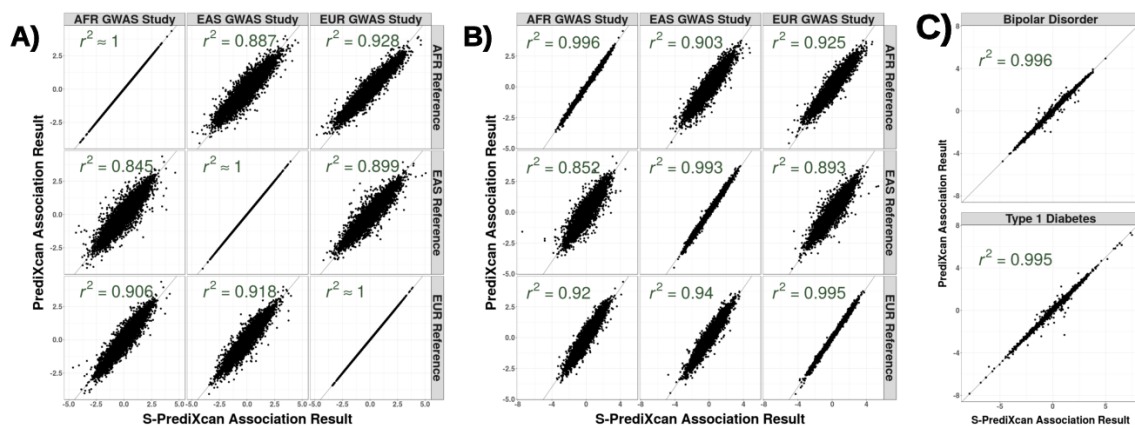


Figure 2. Comparison of PrediXcan vs. S-PrediXcan for A) a simulated phenotype under null hypothesis of no genetic component; B) a cellular phenotype (=intrinsic growth); and C) bipolar disorder and type 1 diabetes studies from Wellcome Trust Case Control Consortium (WTCCC). Gene expression prediction models were based on the DGN cohort presented in [11]. For the simulated phenotype, study sets (GWAS set) and reference sets (LD calculation set) consisted of African (661), East Asian (504) and European (503) individuals from the 1000 Genomes Project. When the same study set is used as reference set, we obtained a high correlation: $r^2 > 0.99999$. For the intrinsic growth phenotype, study sets were a subset of 140 individuals from each of the African, Asian and European groups from 1000 Genomes Project. The reference set was the same as for the simulated phenotype. For the disease phenotypes, the study set consisted of British individuals, and the LD calculation set was the European population subset of the 1000 Genomes Project.

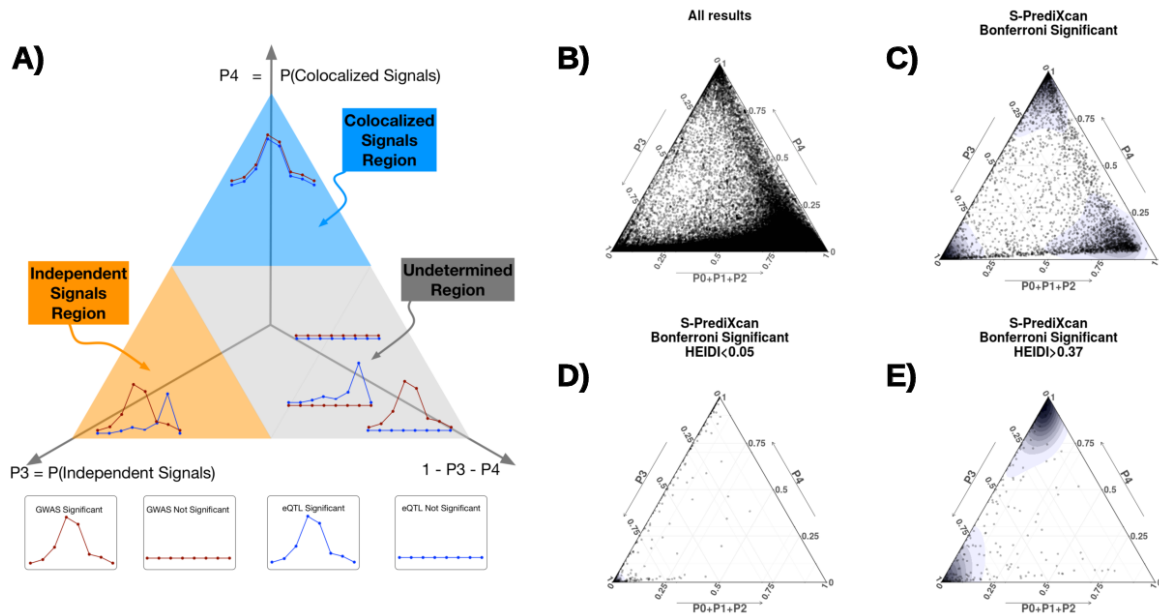


Figure 3. Colocalization status of S-PrediXcan results.

Panel A shows a ternary plot that represents the probabilities of various configurations from COLOC. This plot conveniently constrains the values such that the sum of the probabilities is 1. All points in a horizontal line have the same probability of ‘colocalized’ GWAS and eQTL signals (P_4), points on a line parallel to the right side of the triangle (NW to SE) have the same probability of ‘Independent signals’ (P_3), and lines parallel to the left side of the triangle (NE to SW) correspond to constant $P_1+P_2+P_3$. Top sub-triangle in blue corresponds to high probability of colocalization ($P_4 > 0.5$), lower left sub-triangle in orange corresponds to probability of independent signals ($P_3 > 0.5$), and lower right parallelogram corresponds to genes without enough power to determine or reject colocalization. The following panels present ternary plots of COLOC probabilities with a density overlay for S-PrediXcan results of the Height phenotype.

Panel B shows the colocalization probabilities for all gene-tissue pairs. Most results fall into the ‘undetermined’ region.

Panel C shows that if we keep only Bonferroni-significant S-PrediXcan results, associations tend to cluster into three distinct regions: ‘independent signals’, ‘colocalized’ and ‘undetermined’.

Panel D shows that HEIDI significant genes (to be interpreted as high heterogeneity between GWAS and eQTL signals, i.e. distinct signals) tightly cluster in the ‘independent signal’ region, in concordance with COLOC. A few genes fall in the ‘colocalized’ region, in disagreement with COLOC classification. Unlike COLOC results, HEIDI does not partition the genes into distinct clusters and an arbitrary cutoff p-value has to be chosen.

Panel E shows genes with large HEIDI p-value (no evidence of heterogeneity) which fall in large part in the ‘colocalized’ region. However a substantial number fall in ‘independent signal’ region, disagreeing with COLOC’s classification.

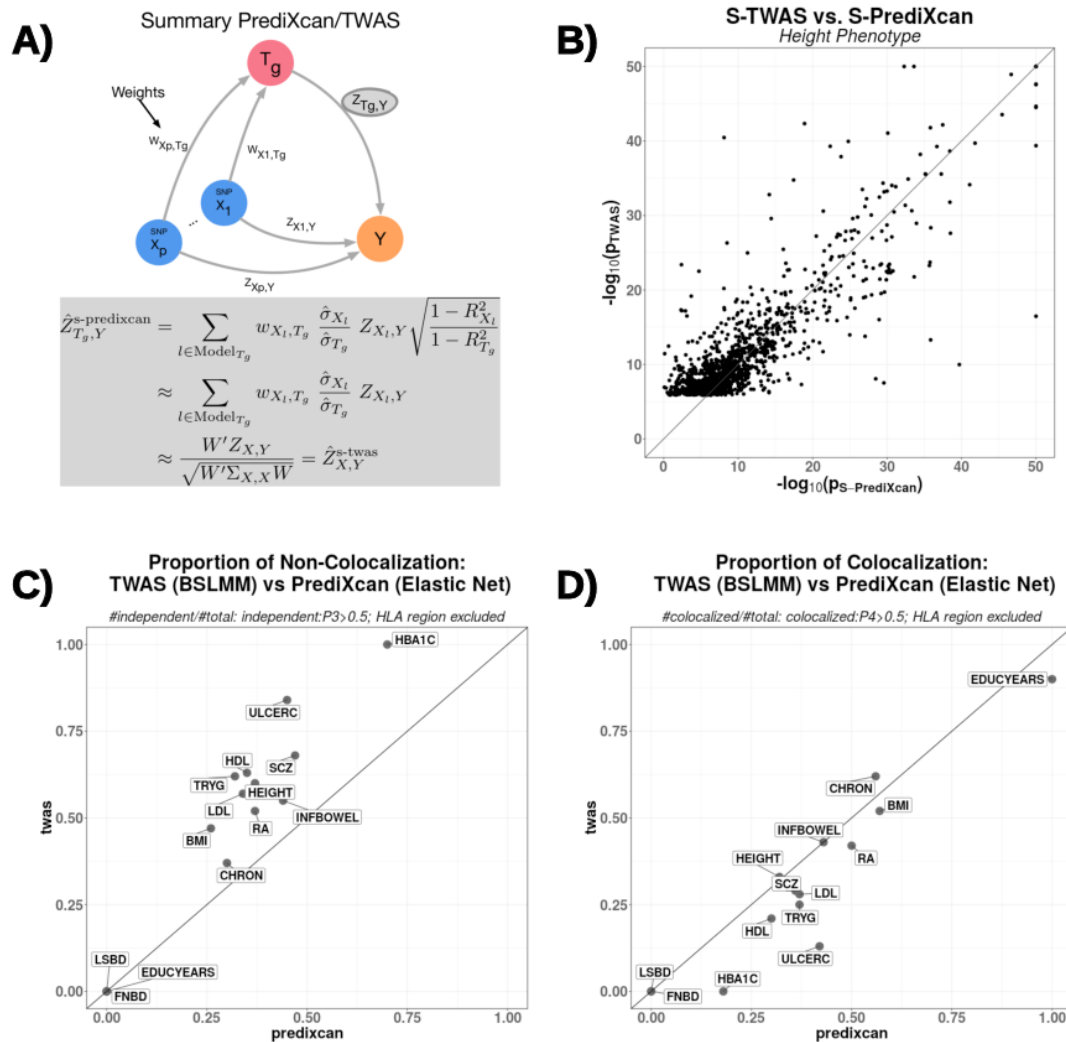


Figure 4. Comparison between S-PrediXcan and S-TWAS.

Panel A depicts how Summary- TWAS and PrediXcan test the mediating role of gene expression level T_g . Multiple SNPs are linked to the expression level of a gene via weights w_{X, T_g} .

Panel B shows the significance of Summary-TWAS (BSLMM) vs. Summary-PrediXcan (elastic net), for the height phenotype across 44 GTEx tissues. There is a small bias caused by using S-TWAS results available from [24], which only lists significant hits. S-PrediXcan tends to yield a larger number of significant associations (see Supplementary Figure 13). P-values were thresholded at 10^{-50} for visualization purposes.

Panel C shows the proportion of non-colocalized associations (distinct eQTL and GWAS signals) from S-TWAS significant vs S-PrediXcan significant results. For all phenotypes, S-TWAS has a higher proportion of LD-contaminated signals compared to S-PrediXcan, as estimated via COLOC.

Panel D shows the proportion of colocalized associations (shared eQTL and GWAS signals) from S-TWAS significant vs S-PrediXcan significant results. For most phenotypes, TWAS has lower proportion of colocalized signals compared to S-PrediXcan, as estimated via COLOC.

Phenotype Abbreviation: Femoral Neck Bone Density (FNBD), Lumbar Spine Bone Density (LSBD), Body Mass Index (BMI), Height (HEIGHT), Low-Density Lipoprotein Cholesterol (LDL), High-Density Lipoprotein Cholesterol (HDL), Tryglicerides (TRYG), Chron's Disease (CHRON), Inflammatory Bowel's Disease (INFBOWEL), Ulcerative Colitis (ULCERC), Hemoglobin Levels (HBA1C) HOMA Insulin Response (HOMA-IR) Schizophrenia (SCZ), Rheumatoid Arthritis (RA), College Completion (COLLEGE), Education Years (EDUCYEARS)

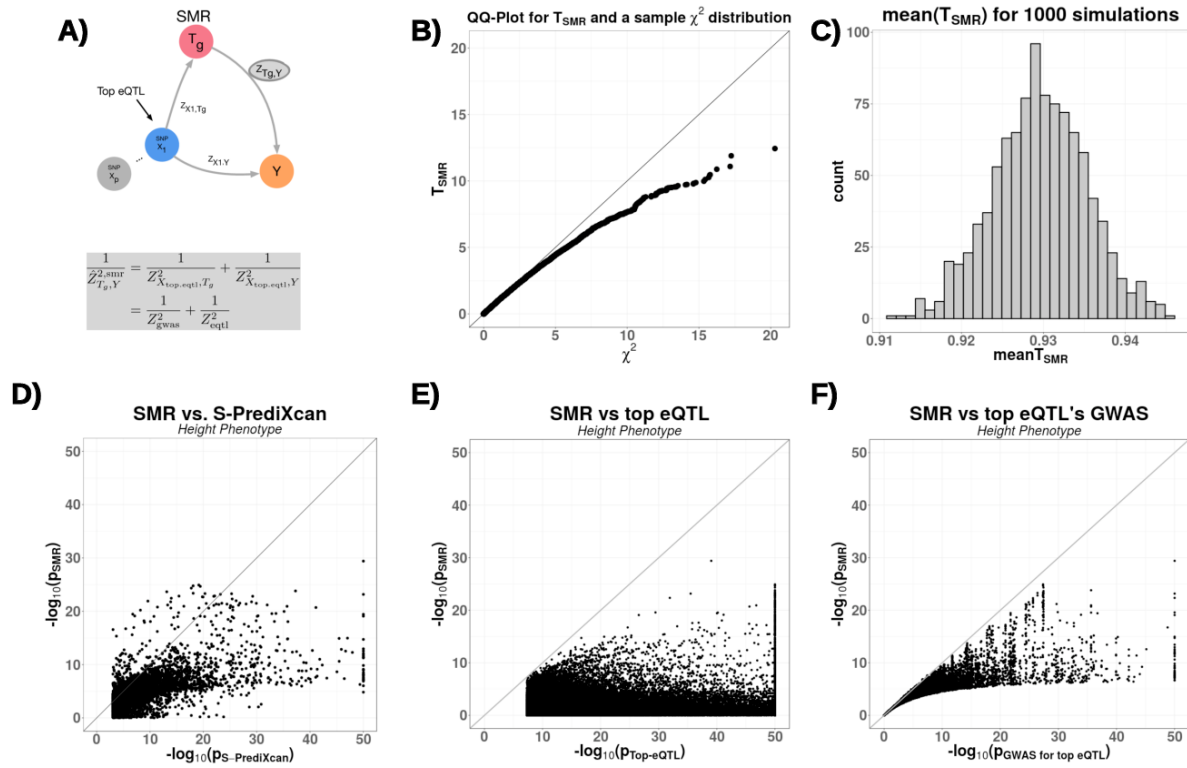


Figure 5. Comparison between Summary-PrediXcan and SMR.

Panel A depicts how SMR tests the mediating role of gene expression level T_g . The top eQTL is linked to the phenotype as an instrumental variable in a Mendelian Randomization approach.

Panel B shows a QQ plot for simulated values of T_{SMR} . Under the null hypothesis of significant eQTL signal and no GWAS association, we generated random values for Z_{GWAS}^2 and Z_{eQTL}^2 following the simulations from [16]. T_{SMR} statistic was calculated from these values, and compared to a χ_1^2 distribution to illustrate this statistics' deflation.

Panel C shows the sample mean of T_{SMR} from 1000 simulations, centered close to 0.93, instead of the expected value of 1 for a χ^2 -distributed variable.

Panel D shows the significance of SMR vs. the significance of Summary-PrediXcan. As expected, SMR associations tend to be smaller than S-PrediXcan ones.

Panels E and F show that the SMR statistics significance is bounded by GWAS and eQTL p-values. The p-values ($-\log_{10}$) of the SMR statistics are plotted against the GWAS p-value of the top eQTL SNP (**panel E**), and the gene's top eQTL p-value (**panel F**).

Some of the associations, GWAS and eQTL p-values were more significant than shown since they were thresholded at 10^{-50} to improve visualization.

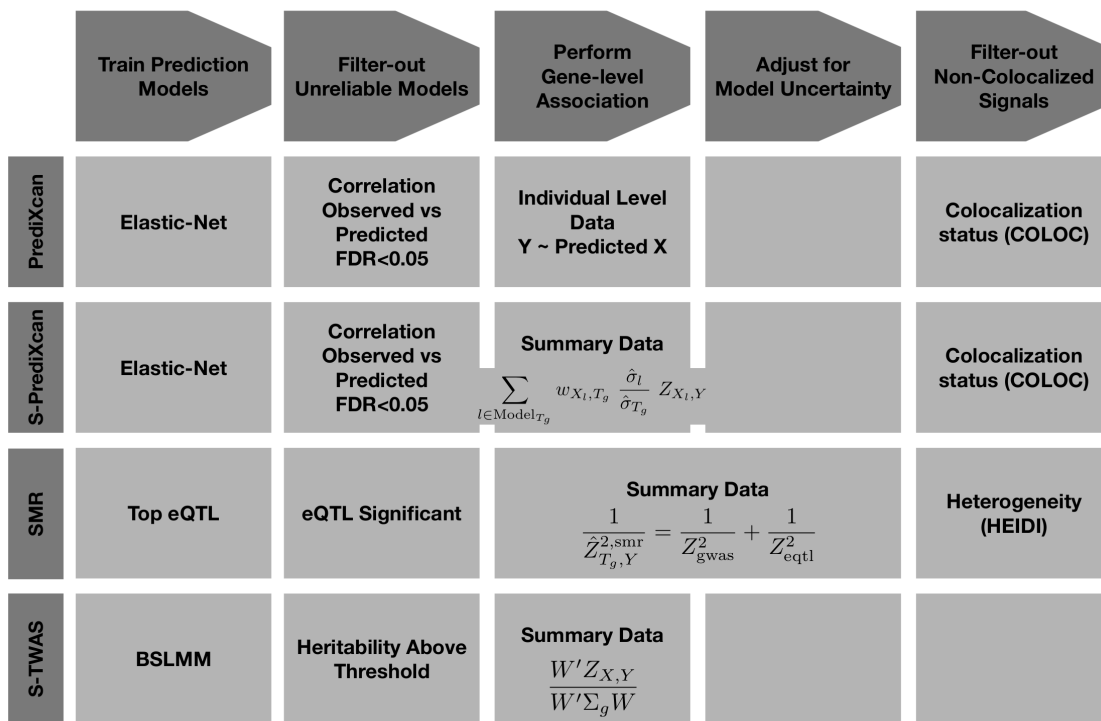


Figure 6. MetaXcan Framework The figure shows a general framework (MetaXcan) which encompasses methods such as PrediXcan, TWAS, SMR, COLOC among others.

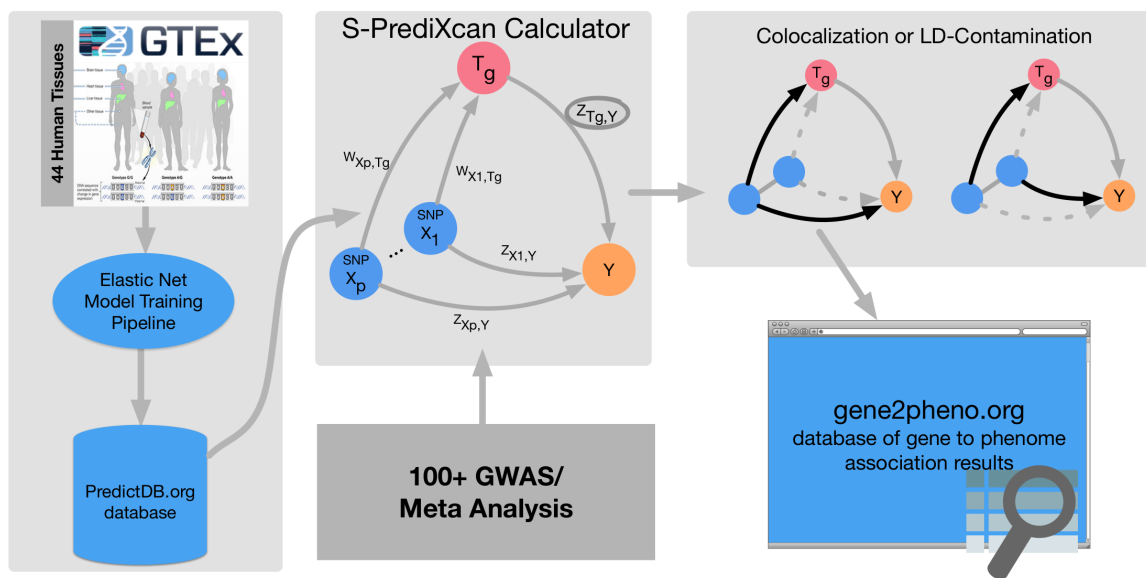


Figure 7. MetaXcan Framework application. This figure summarizes the application of the MetaXcan framework with S-PrediXcan using 44 GTEx tissue transcriptomes and over 100 GWAS and meta analysis results. We trained prediction models using elastic-net [26] and deposited the weights and SNP covariances in the publicly available resource (<http://predictdb.org/>). The weights, covariances and over 100 GWAS summary results were processed with S-PrediXcan. Colocalization status was computed and the full set of results was deposited in gene2pheno.org.

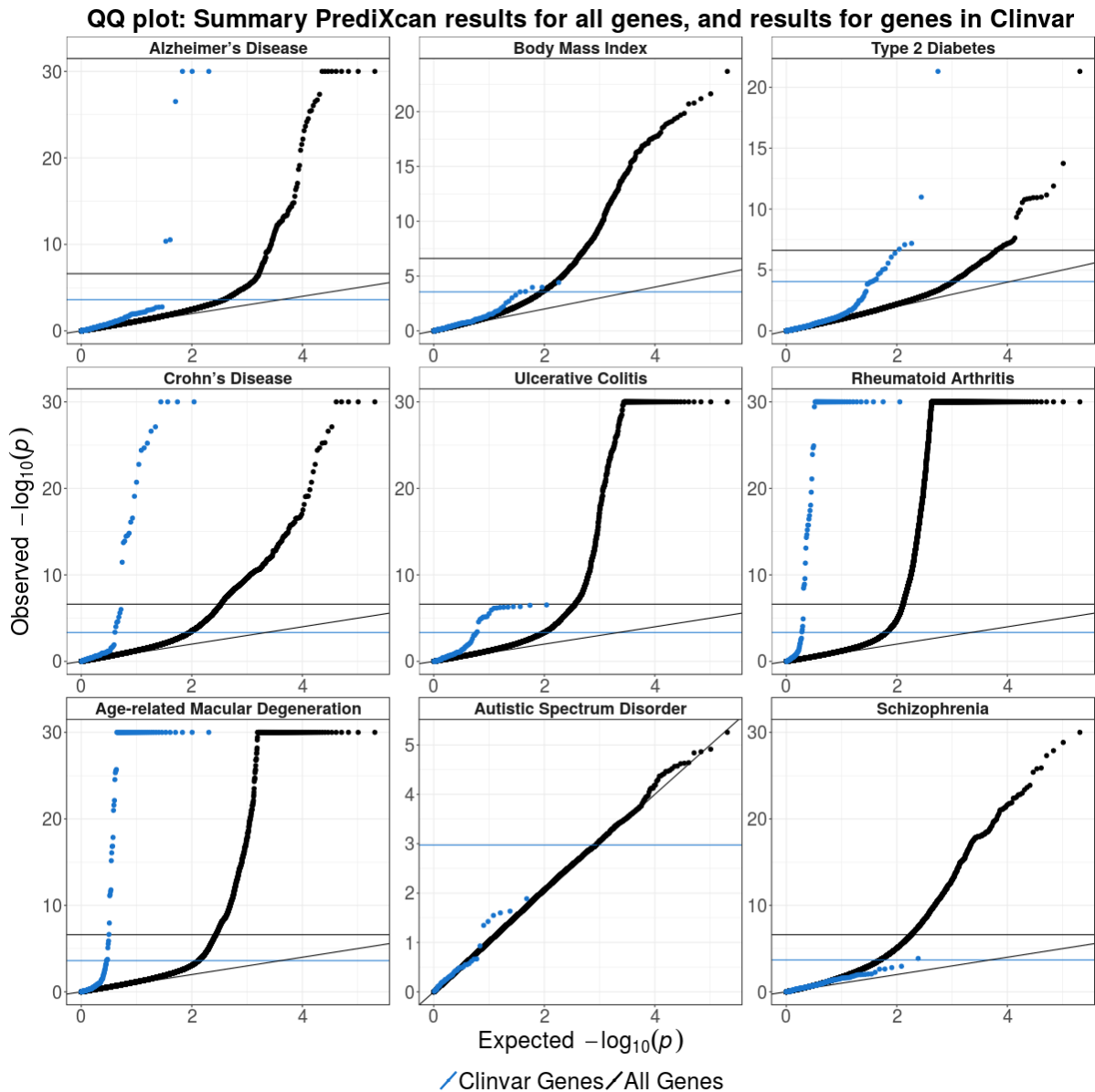


Figure 8. ClinVar genes show significant S-PrediXcan associations. Genes implicated in ClinVar tended to be more significant in S-PrediXcan for most diseases tested, except for schizophrenia and autism. This suggests that more moderate alteration of monogenic disease genes may contribute in a continuum of more moderate but related phenotypes. Alternatively, a more complex interplay between common and rare variation could be taking place such as higher tolerance to loss of function mutations in lower expressing haplotypes which could induce association with predicted expression. Blue circles correspond to the QQ plot of genes in ClinVar that were annotated with the phenotype and black circles correspond to all genes.

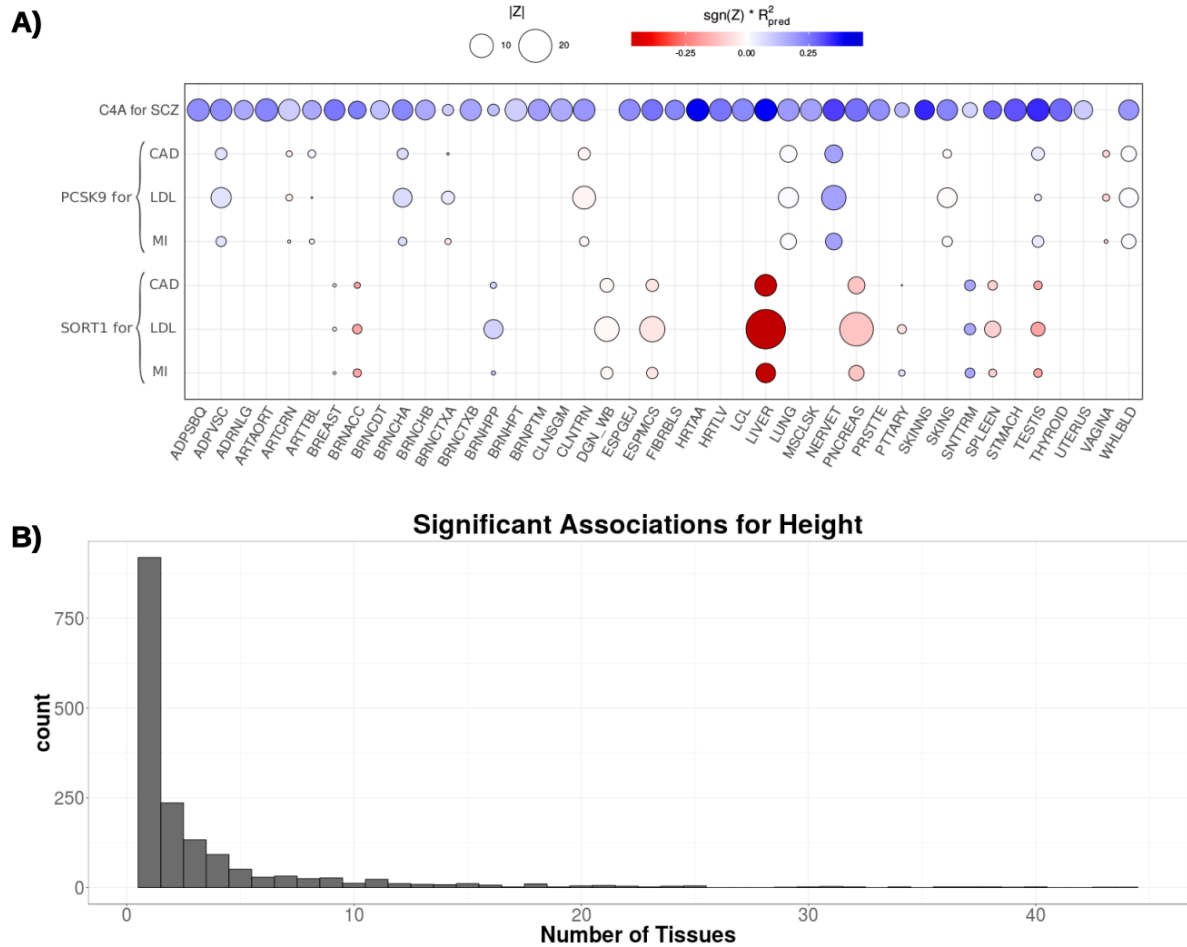


Figure 9. A) S-PrediXcan association for *PCSK9*, *SORT1*, and *C4A* by tissue. This figure shows the association strength between three well studied genes and corresponding phenotypes. *C4A* associations with schizophrenia (SCZ) are significant across most tissues. *SORT1* associations with LDL-C, coronary artery disease (CAD), and myocardial infarction (MI) are most significant in liver. *PCSK9* associations with LDL-C, coronary artery disease (CAD), and myocardial infarction (MI) are most significant in tibial nerve.

The size of the points represent the significance of the association between predicted expression and the traits indicated on the top labels. Red indicates negative correlation whereas blue indicates positive correlation. R_{pred}^2 is a performance measure computed as the correlation squared between observed and predicted expression, cross validated in the training set. Darker points indicate larger genetic component and consequently more active regulation in the tissue.

B) Tissue specificity of most trait associations. This figure shows a histogram of the number of tissues for which a gene is significantly associated with height (other phenotypes show similar pattern).

Tissue abbreviation: Adipose - Subcutaneous (ADPSBQ), Adipose - Visceral (Omentum) (ADPVSC), Adrenal Gland (ADRNLG), Artery - Aorta (ARTAORT), Artery - Coronary (ARTCRN), Artery - Tibial (ARTTBL), Bladder (BLDDER), Brain - Amygdala (BRNAMEY), Brain - Anterior cingulate cortex (BA24) (BRNACC), Brain - Caudate (basal ganglia) (BRNCDT), Brain - Cerebellar Hemisphere (BRNCHB), Brain - Cerebellum (BRNCHA), Brain - Cortex (BRNCTXA), Brain - Frontal Cortex (BA9) (BRNCTXB), Brain - Hippocampus (BRNHPP), Brain - Hypothalamus (BRNHPT), Brain - Nucleus accumbens (basal ganglia) (BRNNCC), Brain - Putamen (basal ganglia) (BRNPMT), Brain - Spinal cord (cervical c-1) (BRNSPC), Brain - Substantia nigra (BRNSNG), Breast - Mammary Tissue (BREAST), Cells - EBV-transformed lymphocytes (LCL), Cells - Transformed fibroblasts (FIBRBLs), Cervix - Ectocervix (CVXECT), Cervix - Endocervix (CVSEND), Colon - Sigmoid (CLNSGM), Colon - Transverse (CLNTRN), Esophagus - Gastroesophageal Junction (ESPEJ), Esophagus - Mucosa (ESPMCS), Esophagus - Muscularis (ESPMSL), Fallopian Tube (FLLPNT), Heart - Atrial Appendage (HRTAA), Heart - Left Ventricle (HRTL), Kidney - Cortex (KDNCTX), Liver (LIVER), Lung (LUNG), Minor Salivary Gland (SLVRYG), Muscle - Skeletal (MSCLSK), Nerve - Tibial (NERVET), Ovary (OVARY), Pancreas (PNCREAS), Pituitary (PTTARY), Prostate (PRSTTE), Skin - Not Sun Exposed (Suprapubic) (SKINNS), Skin - Sun Exposed (Lower leg) (SKINS), Small Intestine - Terminal Ileum (SNTTRM), Spleen (SPLEEN), Stomach (STMACH), Testis (TESTIS), Thyroid (THYROID), Uterus (UTERUS), Vagina (VAGINA), Whole Blood (WHLBLD).

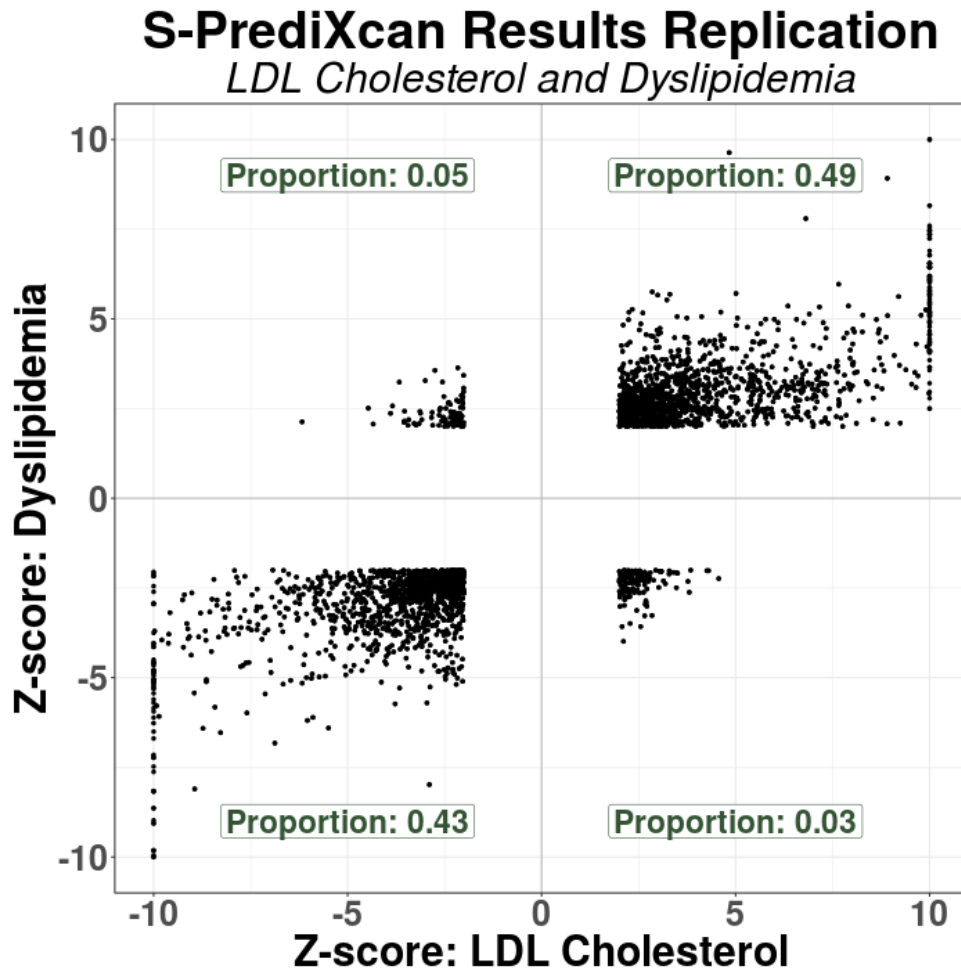


Figure 10. Discovery and replication Z-scores for lipid trait. This figure shows the Z-scores of the association between dyslipidemia (GERA) and predicted gene expression levels on the vertical axis and the Z-scores for LDL cholesterol on the horizontal axis. To facilitate visualization, very large Z-scores were thresholded to 10. Proportions in each quadrant were computed excluding Z-scores with magnitude smaller than 2 to filter out noise.

- 411 • dosage and alternate allele are assumed to be the same.

412 The inputs are based, in general, on data from three different sources:

- 413 • study set (e.g. GWAS study set),
- 414 • training set (e.g. GTEx, DGN),
- 415 • population reference set (e.g. the training set or 1000 Genomes).

416 The study set is the main dataset of interest from which the genotype and phenotypes of interest
417 are gathered. The regression coefficients and standard errors are computed based on individual-level
418 data from the study set or a SNP-level meta-analysis of multiple GWAS. Training sets are the reference
419 transcriptome datasets used for the training of the prediction models (GTEx, DGN, Framingham, etc.)
420 thus the weights w_{ig} are computed from this set. Training sets can also be used to generate variance and
421 covariances of genetic markers, which will usually be different from the study sets. When individual level
422 data are not available from the training set we use population reference sets such as 1000 Genomes data.

423 In the most common use scenario, users will need to provide only GWAS results using their study set.
424 The remaining parameters are pre-computed <https://github.com/hakyimlab/MetaXcan>.

425 Association enrichment

426 We display the enrichment for selected phenotypes in Supplementary Figure 5, measured as $\text{mean}(Z^2)$.
427 For visualization purposes, we selected 25 phenotypes from different categories such as anthropometric
428 traits, cardiometabolic traits, autoimmune diseases, and psychiatric conditions (please see figure caption
429 for the list of selected phenotypes). The simple mean of Z^2 for all gene-tissue pairs in a phenotype was
430 taken.

431 Derivation of Summary-PrediXcan Formula

432 The goal of Summary-PrediXcan is to infer the results of PrediXcan using only GWAS summary statistics.
433 Individual level data are not needed for this algorithm. We will introduce some notations for the derivation
434 of the analytic expressions of S-PrediXcan.

435 **Notation and Preliminaries**

436 Y is the n -dimensional vector of phenotype for individuals $i = 1, n$. X_l is the allelic dosage for SNP l .
 437 T_g is the predicted expression (or estimated GREx, genetically regulated expression). w_{lg} are weights to
 438 predict expression $T_g = \sum_{l \in \text{Model}_g} w_{lg} X_l$, derived from an independent training set.

We model the phenotype as linear functions of X_l and T_g

$$Y = \alpha_1 + X_l \beta_l + \eta$$

$$Y = \alpha_2 + T_g \gamma_g + \epsilon,$$

439 where α_1 and α_2 are intercepts, η and ϵ error terms independent of X_l and T_g , respectively. Let $\hat{\gamma}_g$ and $\hat{\beta}_l$
 440 be the estimated regression coefficients of Y regressed on T_g and X_l , respectively. $\hat{\gamma}_g$ is the result (effect
 441 size for gene g) we get from PrediXcan whereas $\hat{\beta}_l$ is the result from a GWAS for SNP l .

442 We will denote as $\widehat{\text{Var}}$ and $\widehat{\text{Cov}}$ the operators that compute the sample variance and covariance, i.e.
 443 $\widehat{\text{Var}}(Y) = \hat{\sigma}_Y^2 = \sum_{i=1, n} (Y_i - \bar{Y})^2 / (n - 1)$ with $\bar{Y} = \sum_{i=1, n} Y_i / n$. Let $\hat{\sigma}_l^2 = \widehat{\text{Var}}(X_l)$, $\hat{\sigma}_g^2 = \widehat{\text{Var}}(T_g)$ and
 444 $\Gamma_g = (\mathbf{X} - \bar{\mathbf{X}})'(\mathbf{X} - \bar{\mathbf{X}}) / n$, where \mathbf{X}' is the $p \times n$ matrix of SNP data and $\bar{\mathbf{X}}$ is a $n \times p$ matrix where
 445 column l has the column mean of \mathbf{X}_l (p being the number of SNPs in the model for gene g , typically
 446 $p \ll n$).

447 With this notation, our goal is to infer PrediXcan results ($\hat{\gamma}_g$ and its standard error) using only GWAS
 448 results ($\hat{\beta}_l$ and their standard error), estimated variances of SNPs ($\hat{\sigma}_l^2$), estimated covariances between
 449 SNPs in each gene model (Γ_g), and prediction model weights w_{lg} .

450

451 **Input:** $\hat{\beta}_l$, $\text{se}(\hat{\beta}_l)$, $\hat{\sigma}_l^2$, Γ_g , w_{lg} . **Output:** $\hat{\gamma}_g / \text{se}(\hat{\gamma}_g)$.

452

453 Next we list the properties and definitions used in the derivation:

$$\hat{\gamma}_g = \frac{\widehat{\text{Cov}}(T_g, Y)}{\widehat{\text{Var}}(T_g)} = \frac{\widehat{\text{Cov}}(T_g, Y)}{\hat{\sigma}_g^2} \quad (2)$$

454 and

$$\hat{\beta}_l = \frac{\widehat{\text{Cov}}(X_l, Y)}{\widehat{\text{Var}}(X_l)} = \frac{\widehat{\text{Cov}}(X_l, Y)}{\hat{\sigma}_l^2} \quad (3)$$

455 The proportion of variance explained by the covariate (T_g or X_l) can be expressed as

456

$$R_g^2 = \hat{\gamma}_g^2 \frac{\hat{\sigma}_g^2}{\hat{\sigma}_Y^2}$$

457

$$R_l^2 = \hat{\gamma}_l^2 \frac{\hat{\sigma}_l^2}{\hat{\sigma}_Y^2}$$

By definition

$$T_g = \sum_{l \in \text{Model}_g} w_{lg} X_l$$

$\widehat{\text{Var}}(T_g) = \hat{\sigma}_g^2$ can be computed as

$$\begin{aligned} \hat{\sigma}_g^2 &= \widehat{\text{Var}} \left(\sum_{l \in \text{Model}_g} w_{lg} X_l \right) \\ &= \widehat{\text{Var}}(\mathbf{W}_g \mathbf{X}_g) && \text{where } \mathbf{W}_g \text{ is the vector of } w_{lg} \text{ for SNPs in the model of } g \\ &= \mathbf{W}_g' \widehat{\text{Var}}(\mathbf{X}_g) \mathbf{W}_g && \text{where } \Gamma_g \text{ is the } \widehat{\text{Var}}(\mathbf{X}_g) = \text{sample covariance matrix of } \mathbf{X}_g \\ &= \mathbf{W}_g' \Gamma_g \mathbf{W}_g \end{aligned} \tag{4}$$

458 **Calculation of regression coefficient $\hat{\gamma}_g$**

$\hat{\gamma}_g$ can be expressed as

$$\begin{aligned} \hat{\gamma}_g &= \frac{\widehat{\text{Cov}}(T_g, Y)}{\hat{\sigma}_g^2} \\ &= \frac{\widehat{\text{Cov}}(\sum_{l \in \text{Model}_g} w_{lg} X_l, Y)}{\hat{\sigma}_g^2} \\ &= \sum_{l \in \text{Model}_g} \frac{w_{lg} \widehat{\text{Cov}}(X_l, Y)}{\hat{\sigma}_g^2} && \text{by linearity of } \widehat{\text{Cov}} \\ &= \sum_{l \in \text{Model}_g} \frac{w_{lg} \hat{\beta}_l \sigma_l^2}{\hat{\sigma}_g^2} && \text{using Eq 3} \end{aligned} \tag{5}$$

459 **Calculation of standard error of $\hat{\gamma}_g$**

460 Also from the properties of linear regression we know that

$$\text{se}^2(\hat{\gamma}_g) = \text{Var}(\hat{\gamma}_g) = \frac{\hat{\sigma}_\epsilon^2}{n \hat{\sigma}_g^2} = \frac{\hat{\sigma}_Y^2 (1 - R_g^2)}{n \hat{\sigma}_g^2} \tag{6}$$

461 In this equation, $\hat{\sigma}_Y^2/n$ is not necessarily known but can be estimated using the analogous equation (6)

462 for β_l :

$$\text{se}^2(\hat{\beta}_l) = \frac{\hat{\sigma}_Y^2(1 - R_l^2)}{n\hat{\sigma}_l^2} \quad (7)$$

463 Thus:

$$\frac{\hat{\sigma}_Y^2}{n} = \frac{\text{se}(\hat{\beta}_l)^2 \hat{\sigma}_l^2}{(1 - R_l^2)} \quad (8)$$

464 Notice that the right hand side of (8) is dependent on the SNP l while the left hand side is not. This
 465 equality will hold only approximately in our implementation since we will be using approximate values
 466 for $\hat{\sigma}_l^2$, i.e. from reference population, not the actual study population.

467 Calculation of Z-score

To assess the significance of the association, we need to compute the ratio of the estimated effect size $\hat{\gamma}_g$
 and standard error $\text{se}(\hat{\gamma}_g)$, or Z-score,

$$Z_g = \frac{\hat{\gamma}_g}{\text{se}(\hat{\gamma}_g)} \quad (9)$$

with which we can compute the p-value as $p = 2\Phi(-|Z_g|)$ where $\Phi(\cdot)$ is the normal CDF function.

$$\begin{aligned} Z_g &= \frac{\hat{\gamma}_g}{\text{se}(\hat{\gamma}_g)} \\ &= \sum_{l \in \text{Model}_g} \frac{w_{lg} \hat{\beta}_l \hat{\sigma}_l^2}{\hat{\sigma}_g^2} \sqrt{\frac{n}{\hat{\sigma}_Y^2} \frac{\hat{\sigma}_g^2}{(1 - R_g^2)}} && \text{using Eq. 5 and 6} \\ &= \sum_{l \in \text{Model}_g} \frac{w_{lg} \hat{\beta}_l \hat{\sigma}_l^2}{\hat{\sigma}_g} \sqrt{\frac{(1 - R_l^2)}{\text{se}(\hat{\beta}_l)^2 \hat{\sigma}_l^2}} \sqrt{\frac{1}{(1 - R_g^2)}} \\ &= \sum_{l \in \text{Model}_g} w_{lg} \frac{\hat{\sigma}_l}{\hat{\sigma}_g} \frac{\hat{\beta}_l}{\text{se}(\hat{\beta}_l)} \sqrt{\frac{1 - R_l^2}{1 - R_g^2}} \end{aligned} \quad (10)$$

$$\approx \sum_{l \in \text{Model}_g} w_{lg} \frac{\hat{\sigma}_l}{\hat{\sigma}_g} \frac{\hat{\beta}_l}{\text{se}(\hat{\beta}_l)} \quad (11)$$

468 Based on results with actual and simulated data for realistic effect size ranges, we have found that the
 469 last approximation does not affect our ability to identify the association. The approximation becomes
 470 inaccurate only when the effect sizes are very large. But in these cases, the small decrease in statistical

471 efficiency induced by the approximation is compensated by the large power to detect the larger effect
472 sizes.

473 **Calculation of σ_g in reference set**

474 The variance of predicted expression is computed using equation (4) which takes weights for each SNP
475 in the prediction model and the correlation (LD) between the SNPs. The correlation is computed in a
476 reference set such as 1000G or in the training set.

477 **Expression model training**

478 To train our prediction models, we obtained genotype data and normalized gene expression data collected
479 by the GTEx Project. We used 44 different tissues sampled by GTEx and thus generated 44 different
480 tissue-wide models (dbGaP Accession phs000424.v6.p1). Sample sizes for different tissues range from 70
481 (Uterus) to 361 (Muscle - Skeletal). The models referenced in this paper make use of the GTEx Project's
482 V6p data, a patch to the version 6 data and makes use of improved gene-level annotation. We removed
483 ambiguously stranded SNPs from genotype data, i.e. ref/alt pairs A/T, C/G, T/A, G/C. Genotype data
484 was filtered to include only SNPs with MAF > 0.01. For each tissue, normalized gene expression data
485 was adjusted for covariates such as gender, sequencing platform, the top 3 principal components from
486 genotype data and top PEER Factors. The number of PEER Factors used was determined by sample
487 size: 15 for $n < 150$, 30 for n between 150 and 250, and 35 for $n > 250$. Covariate data was provided by
488 GTEx. For our analysis, we used protein-coding genes only.

489 For each gene-tissue pair for which we had adjusted expression data, we fit an Elastic-Net model based
490 on the genotypes of the samples for the SNPs located within 1 Mb upstream of the gene's transcription
491 start site and 1 Mb downstream of the transcription end site. We used the R package glmnet with
492 mixing parameter alpha equal to 0.5, and the penalty parameter lambda was chosen through 10-fold
493 cross-validation.

494 Once we fit all models, we retained only the models q-value less than 0.05 [40]. For each tissue exam-
495 ined, we created a sqlite database to store the weights of the prediction models, as well as other statistics
496 regarding model training. These databases have been made available for download at PredictDB.org.

497 **Online Catalog and SMR, COLOC, TWAS**

498 Supplementary Table 4 shows the list of GWA/GWAMA studies we considered in this analysis. We
499 applied S-PrediXcan to these studies using the transcriptome models trained on GTEx studies for patched
500 version 6. For simplicity, S-PrediXcan only considers those SNPs that have a matching set of alleles in
501 the prediction model, and adjusts the dosages ($2 - \text{dosage}$) if the alleles are swapped.

502 To make the results of this study broadly accessible, we built a Postgre SQL relational database to
503 store S-PrediXcan results, and serve them via a web application <http://gene2pheno.org>.

504 We also applied SMR [16] to the same set of GWAMA studies, using the GTEx eQTL associations. We
505 downloaded version 0.66 of the software from the SMR website, and ran it using the default parameters.
506 We converted the GWAMA and GTEx eQTL studies to SMR input formats. In order to have SMR
507 compute the colocalization test, for those few GWAMA studies where allele frequency was not reported,
508 we filled in with frequencies from the 1000 Genomes Project [41] as an approximation. We also used the
509 1000 Genomes genotype data as reference panel for SMR.

510 Next we ran COLOC [18] over the same set of GWAMA and eQTL studies. We used the R package
511 available from CRAN. We used the Approximate Bayes Factor colocalization analysis, with effect sizes,
512 their standard errors, allele frequencies and sample sizes as arguments. When the frequency information
513 was missing from the GWAS, we filled in with data from the 1000 Genomes Project.

514 For comparison purposes, we have also included the results of the application of Summary-TWAS to
515 30 traits publicly shared by the authors [24].

516 **Comparison with TWAS**

Formal similarity with TWAS can be made more explicit by rewriting S-PrediXcan formula in matrix form. With the following notation and definitions

$$\begin{aligned}\tilde{\mathbf{W}}_g &= (\sigma_1 w_{1g}, \dots, \sigma_p w_{pg})' \\ \mathbf{Z}_{\text{SNPs}} &= (Z_1, \dots, Z_p)' \\ &= \left(\frac{\hat{\beta}_1}{\text{se}(\beta_1)}, \dots, \frac{\hat{\beta}_p}{\text{se}(\beta_p)} \right)'\end{aligned}$$

and correlation matrix of SNPs in the model for gene g

$$\Sigma_g = \text{diag}\left(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_p}\right) \cdot \Gamma_g \cdot \text{diag}\left(\frac{1}{\sigma_1}, \dots, \frac{1}{\sigma_p}\right)$$

517 it is quite straightforward to write the numerator in (1) and (11) as

$$\tilde{\mathbf{W}}_g \cdot \mathbf{Z}_{\text{SNPs}}$$

518 and in the denominator, the variance of the predicted expression level of gene g , as

$$\tilde{\mathbf{W}}_g' \cdot \Sigma_g \cdot \tilde{\mathbf{W}}_g$$

519 thus

$$Z_g = \frac{\tilde{\mathbf{W}}_g \cdot \mathbf{Z}_{\text{SNPs}}}{\sqrt{\tilde{\mathbf{W}}_g' \cdot \Sigma_g \cdot \tilde{\mathbf{W}}_g}}$$

520 This equation has the same form as the TWAS expression if we use the scaled weight vector $\tilde{\mathbf{W}}_g$ instead
 521 of \mathbf{W}_g . Summary-TWAS imputes the Z-score for the gene-level result assuming that under the null
 522 hypothesis, the Z-scores are normally distributed with the same correlation structure as the SNPs; whereas
 523 in S-PrediXcan we compute the results of PrediXcan using summary statistics. Thus, S-TWAS and S-
 524 PrediXcan yield equivalent mathematical expressions (after setting the factor $\sqrt{\frac{1-R_l^2}{1-R_g^2}} \approx 1$).

525 Summary-PrediXcan with only top eQTL as predictor

The S-PrediXcan formula when only the top eQTL is used to predict the expression level of a gene can be expressed as

$$\begin{aligned} Z_{\text{s-prediXcan}} &= \sum_{l \in \text{Model}_g} w_{lg} \frac{\hat{\sigma}_l}{\hat{\sigma}_g} \frac{\hat{\beta}_l}{\text{se}(\hat{\beta}_l)} && \text{using 11} \\ &= w_{1g} \frac{\hat{\sigma}_1}{\sqrt{w_{1g}^2 \sigma_1^2}} Z_1 && \text{only top eQTL is in the model} \\ &= Z_1 \end{aligned}$$

526 where Z_1 is the GWAS Z-score of the top eQTL in the model for gene. Thus

$$Z_{\text{top eQTL s-predixcan}}^2 = Z_{\text{GWAS}}^2 \quad (12)$$

527 Comparison with SMR

SMR quantifies the strength of the association between expression levels of a gene and complex traits with T_{SMR} using the following function of the eQTL and GWAS Z-score statistics:

$$T_{\text{SMR}} = \frac{Z_{\text{eQTL}}^2 Z_{\text{gwas}}^2}{Z_{\text{eQTL}}^2 + Z_{\text{gwas}}^2} \quad (13)$$

528 Here Z_{eQTL} is the Z-score (= effect size/standard error) of the association between SNP and gene
529 expression, and Z_{gwas} is the Z-score of the association between SNP and trait.

530 This SMR statistic (T_{SMR}) is not a χ_1^2 random variable as assumed in [16]. To prove this, we performed
531 simulations following those described in [16]. We generated 10^5 pairs of values for Z_{GWAS}^2 and Z_{eQTL}^2 .
532 Z_{GWAS}^2 was sampled from a χ_1^2 distribution. Z_{eQTL}^2 was sampled from a non-central χ_1^2 distribution with
533 parameter $\lambda = 29$ (a value chosen to mimic results from [29], see [16]). Only pairs with eQTLs satisfying
534 genome-wide significance ($p < 5 \times 10^{-8}$) were kept. We performed a QQ plot and observed deflation when
535 comparing to random values sampled from a χ_1^2 distribution (Supplementary Figure 5-B). This simulation
536 was repeated 1000 times, and we observed a mean of T_{SMR} close to 0.93 (Supplementary Figure 5-C).

537 Only in two extreme cases, the chi-square approximation holds: when $Z_{\text{eQTL}} \gg Z_{\text{gwas}}$ or $Z_{\text{eQTL}} \ll$
538 Z_{gwas} . In these extremes, we can apply Taylor expansions to find an interpretable form of the SMR
539 statistic.

540 If $Z_{\text{eQTL}} \gg Z_{\text{gwas}}$, i.e. if the eQTL association is much more significant than the GWAS association,

$$\begin{aligned}
 T_{\text{SMR}} &= \frac{Z_{\text{gwas}}^2}{1 + \frac{Z_{\text{gwas}}^2}{Z_{\text{eqtl}}^2}} \\
 &\approx Z_{\text{gwas}}^2 \left(1 - \frac{Z_{\text{gwas}}^2}{Z_{\text{eqtl}}^2} \right) \tag{14}
 \end{aligned}$$

so for large enough Z_{eqtl}^2 relative to Z_{gwas}^2

$$\approx Z_{\text{gwas}}^2 = Z_{\text{top eqtl s-predixcan}}^2 \tag{15}$$

541 with the last equality from 12. Thus, in this case, the SMR statistic is slightly smaller than the (top
 542 eQTL based) S-PrediXcan χ_1^2 -square. This reduced significance is accounting for the uncertainty in the
 543 eQTL association. As the evidence for eQTL association grows, the denominator Z_{eqtl}^2 increases and the
 544 difference tends to 0.

On the other extreme when the GWAS association is much stronger than the eQTL's, $Z_{\text{eqtl}} \ll Z_{\text{gwas}}$,

$$\begin{aligned}
 T_{\text{SMR}} &= \frac{Z_{\text{eqtl}}^2}{1 + \frac{Z_{\text{eqtl}}^2}{Z_{\text{gwas}}^2}} \\
 &\approx Z_{\text{eqtl}}^2 \left(1 - \frac{Z_{\text{eqtl}}^2}{Z_{\text{gwas}}^2} \right) \tag{16}
 \end{aligned}$$

so for large enough Z_{gwas}^2 relative to Z_{eqtl}^2

$$\approx Z_{\text{eqtl}}^2 \tag{17}$$

545 In both extremes, the SMR statistic significance is approximately equal to the less significant of the
 546 two statistics GWAS or eQTL, albeit strictly smaller.

547 In between the two extremes, the right distribution must be computed using numerical methods.
 548 When we look at the empirical distribution of the SMR statistic's p-value against the GWAS and eQTL
 549 (top eQTL for the gene) p-values, we find the ceiling of the SMR statistic is maintained as shown in
 550 Figure 5-E and -F. Supplementary Figure 12 shows a comparison of colocalization proportions between
 551 SMR and PrediXcan.

552 **GERA imputation**

553 Genotype files were obtained from dbGaP, and updated to release 35 of the probe annotations published
554 by Affymetrix via PLINK [42]. Probes were filtered out that had a minor allele frequency of <0.01 , were
555 missing in $>10\%$ of subjects, or did not fit Hardy-Weinberg equilibrium. Subjects were dropped that
556 had an unexpected level of heterozygosity ($F >0.05$). Finally the HRC-1000G-check-bim.pl script (from
557 <http://www.well.ox.ac.uk/~wrayner/tools/>) was used to perform some final filtering and split data
558 by chromosome. Phasing (via eagle v2.3 [43]) and imputation against the HRC r1.1 2016 panel [44] (via
559 minimac3) were carried out by the Michigan Imputation Server [45].

560 **GERA GWAS and MetaXcan Application**

561 European samples had been split into ten groups during imputation to ease the computational burden on
562 the Michigan server, so after obtaining the imputed .vcf files, we used the software PLINK [42] to convert
563 the genotype files into the PLINK binary file format and merge the ten groups of samples together,
564 while dropping any variants not found in all sample groups. For the association analysis, we performed
565 a logistic regression using PLINK, and following QC practices from [14] we filtered out individuals with
566 genotype missingness > 0.03 and filtered out variants with minor allele frequency < 0.01 , missingness $>$
567 0.05 , out of Hardy-Weinberg equilibrium significant at $1e-6$, or had imputation quality < 0.8 . We used
568 gender and the first ten genetic principal components as obtained from dbGaP as covariates. Following
569 all filtering, our analysis included 61,444 European samples with 7,120,064 variants. MetaXcan was then
570 applied to these GWAS results, using the 45 prediction models (GTEx and DGN).

571 **Acknowledgments**

572 **Grants**

573 We acknowledge the following US National Institutes of Health grants: R01MH107666 (H.K.I.), T32
574 MH020065 (K.P.S.), R01 MH101820 (GTEx), P30 DK20595 (Diabetes Research and Training Center),
575 F31 DK101202 (J.M.T.), P50 DA037844 (Rat Genomics), P50 MH094267 (Conte). H.E.W. was supported
576 in part by start-up funds from Loyola University Chicago.

577 The Genotype-Tissue Expression (GTEx) Project was supported by the Common Fund of the Office

578 of the Director of the National Institutes of Health. Additional funds were provided by the NCI, NHGRI,
579 NHLBI, NIDA, NIMH, and NINDS. Donors were enrolled at Biospecimen Source Sites funded by NCI
580 SAIC-Frederick, Inc. (SAIC-F) subcontracts to the National Disease Research Interchange (10XS170),
581 Roswell Park Cancer Institute (10XS171), and Science Care, Inc. (X10S172). The Laboratory, Data
582 Analysis, and Coordinating Center (LDACC) was funded through a contract (HHSN268201000029C)
583 to The Broad Institute, Inc. Biorepository operations were funded through an SAIC-F subcontract to
584 Van Andel Institute (10ST1035). Additional data repository and project management were provided
585 by SAIC-F (HHSN261200800001E). The Brain Bank was supported by a supplements to University
586 of Miami grants DA006227 & DA033684 and to contract N01MH000028. Statistical Methods devel-
587 opment grants were made to the University of Geneva (MH090941 & MH101814), the University of
588 Chicago (MH090951, MH090937, MH101820, MH101825), the University of North Carolina - Chapel Hill
589 (MH090936 & MH101819), Harvard University (MH090948), Stanford University (MH101782), Washing-
590 ton University St Louis (MH101810), and the University of Pennsylvania (MH101822). The data used for
591 the analyses described in this manuscript were obtained from dbGaP accession number phs000424.v6.p1
592 on 06/17/2016.

593 This work was completed in part with resources provided by the University of Chicago Research
594 Computing Center, Bionimbus [46], and the Center for Research Informatics. The Center for Research
595 Informatics is funded by the Biological Sciences Division at the University of Chicago with additional
596 funding provided by the Institute for Translational Medicine, CTSA grant number UL1 TR000430 from
597 the National Institutes of Health.

598 References

- 599 1. Nica AC, Montgomery SB, Dimas AS, Stranger BE, Beazley C, Barroso I, et al. Candidate causal
600 regulatory effects by integration of expression QTLs with complex trait genetic associations. *PLoS*
601 *Genetics*. 2010;6(4).
- 602 2. Nicolae DL, Gamazon E, Zhang W, Duan S, Eileen Dolan M, Cox NJ. Trait-associated SNPs are
603 more likely to be eQTLs: Annotation to enhance discovery from GWAS. *PLoS Genetics*. 2010;6(4).
- 604 3. Li YI, van de Geijn B, Raj A, Knowles DA, Petti AA, Golan D, et al. RNA splicing is a primary
605 link between genetic variation and disease. *Science*. 2016;352(6285):600–604. Available from: [http:](http://)

- 606 [//www.ncbi.nlm.nih.gov/pubmed/27126046](http://www.ncbi.nlm.nih.gov/pubmed/27126046).
- 607 4. Gusev A, Lee SH, Trynka G, Finucane H, Vilhjálmsson BJ, Xu H, et al. Partitioning heritability of
608 regulatory and cell-type-specific variants across 11 common diseases. *American Journal of Human*
609 *Genetics*. 2014;95(5):535–552.
- 610 5. Battle A, Mostafavi S, Zhu X, Potash JB, Weissman MM, McCormick C, et al. Characterizing
611 the genetic basis of transcriptome diversity through RNA-sequencing of 922 individuals. *Genome*
612 *Research*. 2014;24(1):14–24.
- 613 6. Lappalainen T, Sammeth M, Friedländer MR, 't Hoen PaC, Monlong J, Rivas Ma, et al.
614 Transcriptome and genome sequencing uncovers functional variation in humans. *Nature*.
615 2013;501(7468):506–11. Available from: [http://www.pubmedcentral.nih.gov/articlerender.](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3918453&tool=pmcentrez&rendertype=abstract)
616 [fcgi?artid=3918453&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3918453&tool=pmcentrez&rendertype=abstract).
- 617 7. Zhang X, Joehanes R, Chen BH, Huan T, Ying S, Munson PJ, et al. Identification of common
618 genetic variants controlling transcript isoform variation in human whole blood. *Nature Genetics*.
619 2015;47(4):345–352. Available from: <http://www.nature.com/doifinder/10.1038/ng.3220>.
- 620 8. Stranger BE, Montgomery SB, Dimas AS, Parts L, Stegle O, Ingle CE, et al. Patterns of Cis
621 regulatory variation in diverse human populations. *PLoS Genetics*. 2012;8(4).
- 622 9. The GTEx Consortium. The Genotype-Tissue Expression (GTEx) project. *Nature genetics*.
623 2013;45(6):580–5. Available from: [http://www.pubmedcentral.nih.gov/articlerender.fcgi?](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4010069&tool=pmcentrez&rendertype=abstract)
624 [artid=4010069&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4010069&tool=pmcentrez&rendertype=abstract).
- 625 10. Aguet F, Brown AA, Castel S, Davis JR, Mohammadi P, Segre AV, et al. Local genetic effects
626 on gene expression across 44 human tissues. *bioRxiv*. 2016; Available from: [http://biorxiv.org/](http://biorxiv.org/content/early/2016/09/09/074450)
627 [content/early/2016/09/09/074450](http://biorxiv.org/content/early/2016/09/09/074450).
- 628 11. Gamazon ER, Wheeler HE, Shah KP, Mozaffari SV, Aquino-Michaels K, Carroll RJ, et al. A gene-
629 based association method for mapping traits using reference transcriptome data. *Nature genetics*.
630 2015;47(9):1091–1098. Available from: <http://dx.doi.org/10.1038/ng.3367>.
- 631 12. Smoller JW, Craddock N, Kendler K, Lee PH, Neale BM, Nurnberger JI, et al. Identifica-
632 tion of risk loci with shared effects on five major psychiatric disorders: a genome-wide analy-

- 633 sis. *Lancet*. 2013;381(9875):1371–9. Available from: [http://discovery.ucl.ac.uk/1395494/\\$\](http://discovery.ucl.ac.uk/1395494/$\)
634 [delimiter"026E30F\\$nhttp://www.ncbi.nlm.nih.gov/pubmed/23453885](http://www.ncbi.nlm.nih.gov/pubmed/23453885).
- 635 13. Deloukas P, Kanoni S, Willenborg C, Farrall M, Assimes TL, Thompson JR, et al. Large-
636 scale association analysis identifies new risk loci for coronary artery disease. *Nature genetics*.
637 2013;45(1):25–33. Available from: <http://www.pubmedcentral.nih.gov/articlerender.fcgi?>
638 [artid=3679547&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3679547&tool=pmcentrez&rendertype=abstract).
- 639 14. Morris AP, Voight BF, Teslovich TM, Ferreira T, Segrè AV, Steinthorsdottir V, et al.
640 Large-scale association analysis provides insights into the genetic architecture and patho-
641 physiology of type 2 diabetes. *Nature Genetics*. 2012;44(9):981–990. Available
642 from: <http://www.ncbi.nlm.nih.gov/pubmed/22885922>
643 [\delimiter"026E30F\\$nhttp://www.nature.com/doifinder/10.1038/ng.2383](http://www.nature.com/doifinder/10.1038/ng.2383).
- 644 15. Gusev A, Lee SH, Trynka G, Finucane H, Vilhjálmsson BJ, Xu H, et al. Integrative approaches
645 for large-scale transcriptome-wide association studies. *Nature Genetics*. 2016;48:245–252.
- 646 16. Zhu Z, Zhang F, Hu H, Bakshi A, Robinson MR, Powell JE, et al. Integration of sum-
647 mary data from GWAS and eQTL studies predicts complex trait gene targets. *Nature genet-*
648 *ics*. 2016;48(5):481–7. Available from: <http://www.nature.com/doifinder/10.1038/ng.3538>
649 [\delimiter"026E30F\\$nhttp://www.ncbi.nlm.nih.gov/pubmed/27019110](http://www.ncbi.nlm.nih.gov/pubmed/27019110).
- 650 17. He X, Fuller CK, Song Y, Meng Q, Zhang B, Yang X, et al. Sherlock: Detecting Gene-Disease As-
651 sociations by Matching Patterns of Expression QTL and GWAS. *The American Journal of Human*
652 *Genetics*. 2013 May;92(5):667–680. Available from: [http://dx.doi.org/10.1016/j.ajhg.2013.](http://dx.doi.org/10.1016/j.ajhg.2013.03.022)
653 [03.022](http://dx.doi.org/10.1016/j.ajhg.2013.03.022).
- 654 18. Giambartolomei C, Vukcevic D, Schadt EE, Franke L, Hingorani AD, Wallace C, et al. Bayesian
655 test for colocalisation between pairs of genetic association studies using summary statistics. *PLoS*
656 *Genetics*. 2014 May;10(5):e1004383. Available from: [http://eutils.ncbi.nlm.nih.gov/entrez/](http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=24830394&retmode=ref&cmd=prlinks)
657 [eutils/elink.fcgi?dbfrom=pubmed&id=24830394&retmode=ref&cmd=prlinks](http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=24830394&retmode=ref&cmd=prlinks).
- 658 19. Hormozdiari F, van de Bunt M, Segrè AV, Li X, Joo JWJ, Bilow M, et al.; Los Angeles Los
659 Angeles CA 90095 USA. Department of Computer Science. Colocalization of GWAS and eQTL

- 660 Signals Detects Target Genes. *Am J Hum Genet.* 2016; Available from: <http://dx.doi.org/10.1016/j.ajhg.2016.10.003>.
- 661
- 662 20. Wen X, Pique-Regi R, Luca F. Integrating molecular QTL data into genome-wide genetic asso-
663 ciation analysis: Probabilistic assessment of enrichment and colocalization. *PLoS Genetics.* 2017
664 Mar;13(3):e1006646. Available from: <http://dx.plos.org/10.1371/journal.pgen.1006646>.
- 665 21. WTCCC. Genome-wide association study of 14,000 cases of seven common diseases and 3,000
666 shared controls. *Nature.* 2007;447(7145):661–78. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/17554300>.
- 667
- 668 22. Manor O, Segal E. Robust Prediction of Expression Differences among Human Individuals Using
669 Only Genotype Information. *PLoS Genetics.* 2013 Mar;9(3):e1003396. Available from: <http://dx.plos.org/10.1371/journal.pgen.1003396>.
- 670
- 671 23. Hamilton N. ggtern: An Extension to 'ggplot2', for the Creation of Ternary Diagrams; 2016. R
672 package version 2.2.0. Available from: <https://CRAN.R-project.org/package=ggtern>.
- 673 24. Mancuso N, Shi H, Goddard P, Kichaev G, Gusev A, Pasaniuc B. Integrating Gene Expression
674 with Summary Association Statistics to Identify Genes Associated with 30 Complex Traits. *The
675 American Journal of Human Genetics.* 2017 Mar;100(3):473–487. Available from: <http://dx.doi.org/10.1016/j.ajhg.2017.01.031>.
- 676
- 677 25. Zhou X, Carbonetto P, Stephens M. Polygenic Modeling with Bayesian Sparse Linear Mixed
678 Models. *PLoS Genetics.* 2013;9(2).
- 679 26. Zou H, Hastie T. Regularization and variable selection via the elastic-net. *Journal of the Royal
680 Statistical Society.* 2005;67:301–320.
- 681 27. Wheeler HE, Shah KP, Brenner J, Garcia T, Aquino-Michaels K, Cox NJ, et al. Survey of the
682 Heritability and Sparse Architecture of Gene Expression Traits across Human Tissues. *PLoS
683 Genetics.* 2016;12(11).
- 684 28. Pavlides JMW, Zhu Z, Gratten J, Mcrae AF, Wray NR, Yang J. Predicting gene targets from
685 integrative analyses of summary data from GWAS and eQTL studies for 28 human complex

- 686 traits. *Genome medicine*. 2016 Aug;8(1):1–6. Available from: <http://dx.doi.org/10.1186/s13073-016-0338-4>.
- 687
- 688 29. Westra HJ, Peters MJ, Esko T, Yaghootkar H, Schurmann C, Kettunen J, et al. Systematic
689 identification of trans eQTLs as putative drivers of known disease associations. *Nature Genetics*.
690 2013 Sep;45(10):1238–1243. Available from: [http://www.nature.com/doifinder/10.1038/ng.](http://www.nature.com/doifinder/10.1038/ng.2756)
691 2756.
- 692 30. Landrum MJ, Lee JM, Benson M, Brown G, Chao C, Chitipiralla S, et al. ClinVar:
693 public archive of interpretations of clinically relevant variants. *Nucleic acids research*.
694 2015;44(D1):D862–8. Available from: [http://www.pubmedcentral.nih.gov/articlerender.](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4702865&tool=pmcentrez&rendertype=abstract)
695 [fcgi?artid=4702865&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=4702865&tool=pmcentrez&rendertype=abstract).
- 696 31. Shah N, Hou YCC, Yu HC, Sainger R, Dec E, Perkins B, et al. Identification of misclassified
697 ClinVar variants using disease population prevalence. 2016 Sep;p. 1–23. Available from: [http:](http://biorxiv.org/lookup/doi/10.1101/075416)
698 [//biorxiv.org/lookup/doi/10.1101/075416](http://biorxiv.org/lookup/doi/10.1101/075416).
- 699 32. Sekar A, Bialas AR, de Rivera H, Davis A, Hammond TR, Kamitaki N, et al. Schizophre-
700 nia risk from complex variation of complement component 4. *Nature*. 2016;530(7589):177–
701 83. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/26814963>{%}5Cn[http://www.](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4752392)
702 [pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4752392](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC4752392).
- 703 33. Musunuru K, Strong A, Frank-Kamenetsky M, Lee NE, Ahfeldt T, Sachs KV, et al. From noncod-
704 ing variant to phenotype via SORT1 at the 1p13 cholesterol locus. *Nature*. 2010;466(7307):714–
705 9. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/20686566>{%}5Cn[http://www.](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3062476)
706 [pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3062476](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC3062476).
- 707 34. Dadu RT, Ballantyne CM. Lipid lowering with PCSK9 inhibitors. *Nature Publishing Group*. 2014
708 Jun;11(10):563–575. Available from: <http://dx.doi.org/10.1038/nrcardio.2014.84>.
- 709 35. Franzén O, Ermel R, Cohain A, Akers NK, Di Narzo A, Talukdar HA, et al. Cardiometabolic
710 risk loci share downstream cis- and trans-gene regulation across tissues and diseases. *Science*.
711 2016 Aug;353(6301):827–830. Available from: [http://www.sciencemag.org/cgi/doi/10.1126/](http://www.sciencemag.org/cgi/doi/10.1126/science.aad6970)
712 [science.aad6970](http://www.sciencemag.org/cgi/doi/10.1126/science.aad6970).

- 713 36. Hoffmann TJ, Ehret GB, Nandakumar P, Ranatunga D, Schaefer C, Kwok PY, et al. Genome-wide
714 association analyses using electronic health records identify new loci influencing blood pressure
715 variation. *Nature Genetics*. 2016 Nov;49(1):54–64. Available from: [http://www.nature.com/
716 doifinder/10.1038/ng.3715](http://www.nature.com/doifinder/10.1038/ng.3715).
- 717 37. Cook JP, Morris AP. Multi-ethnic genome-wide association study identifies novel locus for type 2
718 diabetes susceptibility. *European Journal of Human Genetics*. 2016 Aug;24(8):1175–1180. Available
719 from: <http://www.nature.com/doifinder/10.1038/ejhg.2016.17>.
- 720 38. Torres JM, Barbeira AN, Bonazzola R, Morris AP, Shah KP, Wheeler HE, et al. Integrative cross
721 tissue analysis of gene expression identifies novel type 2 diabetes genes. *bioRxiv*. 2017; Available
722 from: <http://biorxiv.org/content/early/2017/02/27/108134>.
- 723 39. Boyle EA, Li YI, Pritchard JK. An Expanded View of Complex Traits: From Polygenic to Om-
724 nigenic. *Cell*. 2017 Jun;169(7):1177–1186. Available from: [http://dx.doi.org/10.1016/j.cell.
725 2017.05.038](http://dx.doi.org/10.1016/j.cell.2017.05.038).
- 726 40. Storey JD. Statistical significance for genomewide studies. *Proceedings of the Na-
727 tional Academy of Sciences*. 2003 Jul;100(16):9440–9445. Available from: [http:
728 //scholar.google.com.proxy.uchicago.edu/scholar?hl=en&lr=&q=info:eSXwkHMI-nQJ:
729 scholar.google.com/&output=search](http://scholar.google.com.proxy.uchicago.edu/scholar?hl=en&lr=&q=info:eSXwkHMI-nQJ:scholar.google.com/&output=search).
- 730 41. Auton A, Altshuler DM, Durbin RMJA, Wang J, Yang H, Auton A, et al. A global
731 reference for human genetic variation. *Nature*. 2015;526(7571):68–74. Available from:
732 [http://www.nature.com/nature/journal/v526/n7571/fig_{_}tab/nature15393_{_}SF1.
733 html{_%}5Cnhttp://dx.doi.org/10.1038/nature15393{_%}5Cnhttp://www.ncbi.nlm.nih.gov/
734 pubmed/26432245{_%}5Cnhttp://www.nature.com/doifinder/10.1038/nature15393](http://www.nature.com/nature/journal/v526/n7571/fig_{_}tab/nature15393_{_}SF1.html).
- 735 42. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK:
736 rising to the challenge of larger and richer datasets. *GigaScience*. 2015 dec;4(1):7. Available from:
737 <http://gigascience.biomedcentral.com/articles/10.1186/s13742-015-0047-8>.
- 738 43. Loh PR, Danecek P, Palamara PF, Fuchsberger C, A Reshef Y, K Finucane H, et al. Reference-
739 based phasing using the Haplotype Reference Consortium panel. *Nature Genetics*. 2016

- 740 oct;48(11):1443–1448. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/27694958><http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC5096458><http://www.nature.com/doi/10.1038/ng.3679>.
- 741
- 742
- 743 44. McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, et al. A reference panel of
744 64,976 haplotypes for genotype imputation. *Nature Genetics*. 2016 aug;48(10):1279–1283. Available
745 from: <http://www.nature.com/doi/10.1038/ng.3643>.
- 746 45. Das S, Forer L, Schönherr S, Sidore C, Locke AE, Kwong A, et al. Next-generation genotype
747 imputation service and methods. *Nature Genetics*. 2016 aug;48(10):1284–1287. Available from:
748 <http://www.nature.com/doi/10.1038/ng.3656>.
- 749 46. Heath AP, Greenway M, Powell R, Spring J, Suarez R, Hanley D, et al. Bionimbus: a cloud
750 for managing, analyzing and sharing large genomics datasets. *Journal of the American Medical
751 Informatics Association : JAMIA*. 2014 Nov;21(6):969–975. Available from: [https://academic.
752 oup.com/jamia/article-lookup/doi/10.1136/amiajnl-2013-002155](https://academic.oup.com/jamia/article-lookup/doi/10.1136/amiajnl-2013-002155).
- 753 47. Im HK, Gamazon ER, Stark AL, Huang RS, Cox NJ, Dolan ME. Mixed effects modeling of
754 proliferation rates in cell-based models: Consequence for pharmacogenomics and Cancer. *PLoS
755 Genetics*. 2012;8(2).