

1 Integrating Predicted Transcriptome From Multiple Tissues 2 Improves Association Detection

3 Alvaro N. Barbeira¹, Milton D. Pividori¹, Jiamao Zheng¹, Heather E. Wheeler^{2,3}, Dan L. Nicolae^{1,4,5},
4 Hae Kyung Im^{1,5,*}

5 **1 Section of Genetic Medicine, The University of Chicago, Chicago, IL, USA**

6 **2 Department of Biology, Loyola University Chicago, Chicago, IL, USA**

7 **3 Department of Computer Science, Loyola University Chicago, Chicago, IL, USA**

8 **4 Department of Statistics, The University of Chicago, Chicago, IL, USA**

9 **5 Department of Human Genetics, The University of Chicago, Chicago, IL, USA**

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

11 Abstract

12 Integration of genome-wide association studies (GWAS) and expression quantitative trait loci (eQTL)
13 studies is needed to improve our understanding of the biological mechanisms underlying GWAS hits, and
14 our ability to identify therapeutic targets. Gene-level association test methods such as PrediXcan can
15 prioritize candidate targets. However, limited eQTL sample sizes and absence of relevant developmental
16 and disease context restricts our ability to detect associations. Here we propose an efficient statistical
17 method that leverages the substantial sharing of eQTLs across tissues and contexts to improve our ability
18 to identify potential target genes: MulTiXcan. MulTiXcan integrates evidence across multiple panels
19 while taking into account their correlation. We apply our method to a broad set of complex traits available
20 from the UK Biobank and show that we can detect a larger set of significantly associated genes than
21 using each panel separately. To improve applicability, we developed an extension to work on summary
22 statistics: S-MulTiXcan, which we show yields highly concordant results with the individual level version.
23 Results from our analysis as well as software and necessary resources to apply our method are publicly
24 available.

25 Introduction

26 Recent technological advances allow interrogation of the genome to a high level of coverage and precision,
27 enabling experimental studies that query the effect of genotype on both complex and molecular traits.
28 Among these, GWAS have successfully associated genetic loci to human complex traits. GWAS meta-
29 analyses with ever increasing sample sizes allow the detection of associated variants with smaller effect
30 sizes [1–3]. However, understanding the mechanism underlying these associations remains a challenging
31 problem, requiring follow-up studies and a wide array of techniques such as prioritization [4] and pathway
32 analysis [5].

33 Another approach is the study of quantitative trait loci (eQTLs), measuring association between
34 genotype and gene expression. These studies provide a wealth of biological information but tend to have
35 smaller sample sizes. A similar observation applies to QTL studies of other traits such methylation,
36 metabolites, or protein levels.

37 The importance of gene expression regulation in complex traits [6–9] has motivated the integration
38 of eQTL studies and GWAS. To examine these mechanisms we developed PrediXcan [10], a method that
39 tests the mediating role of gene expression variation in complex traits. We also developed an extension
40 that accurately infers its gene-level association results using summary statistics data: S-PrediXcan [11].
41 This allows the rapid exploration of information available in publicly available GWAS summary statistics
42 results, at a dramatically reduced computational burden.

43 Due to sharing of eQTLs across multiple tissues, we have shown the benefits of an agnostic scanning
44 across all available tissues [11]. Despite the increased multiple testing burden (for Bonferroni correction,
45 the total number of gene-tissue pairs must be used when determining the threshold) we gain considerably
46 in number of significant genes. However, given the substantial correlation between different tissues [12],
47 Bonferroni correction can be too stringent increasing the false negative rate.

48 In order to aggregate evidence more efficiently, here we present a method termed MultiXcan that
49 tests the joint effects of gene expression variation from different tissues. Furthermore, we develop
50 and implement a method that only needs summary statistics from a GWAS: Summary-MulTiXcan (S-
51 MulTiXcan for short). We make our implementation publicly available to the research community in
52 <https://github.com/hakyimlab/MetaXcan>. We apply this method to traits from the UK Biobank
53 study and over a hundred public GWAS, and publish the results in <http://gene2pheno.org>.

54 Results

55 Combining Information Across Tissues Through Multivariate Regression

56 To combine information across tissues, we regress the phenotype of interest on the predicted expression
57 of the gene in multiple tissues as follows:

$$Y = \boldsymbol{\mu} + \mathbf{t}_1g_1 + \mathbf{t}_2g_2 + \cdots + \mathbf{t}_pg_p + \mathbf{e} \quad (1)$$

58 where Y is the phenotype, $\boldsymbol{\mu}$ is an intercept term, \mathbf{t}_i is predicted expression for a model trained on tissue
59 study i , g_i is its effect size, and \mathbf{e} an error term.

60 Expression predictions from many of these tissues are highly correlated. To avoid numerical issues
61 caused by collinearity, we compute the principal components of the predicted expression matrix and
62 discard the axes of smallest variation. Additional covariates can be added to the regression seamlessly.
63 Fig. 1-a displays an overview of the method; see further details in the Methods section. We illustrate
64 prediction correlation across models in Supp. Fig. 1.

65 MulTiXcan Detects More Associations Than Single-Tissue PrediXcan

66 We applied our method to 222 traits from UK Biobank [13]. We used prediction models for 44 tissues
67 trained with Genotype-Tissue Expression (GTEx) samples [12]. We found that it can detect many more
68 associations than PrediXcan using a specific tissue or even when aggregating results from all tissues
69 (Bonferroni-corrected for all gene-tissue pairs tested).

70 More specifically, we compared three approaches for assessing the significance of a gene jointly across
71 all tissues: 1) running PrediXcan using the most relevant tissue; 2) running PrediXcan using all tissues,
72 one tissue at a time; 3) running MulTiXcan. Fig. 1-b illustrates a schematic representation of the results
73 from each approach. Overall, MulTiXcan detects more associations than PrediXcan, as shown in Fig.
74 2-c. Supplementary Data 1 contains a summary of associations per trait. See Supplementary Data 2 and
75 3 for the list of significant MulTiXcan and PrediXcan results respectively.

76 We examined the High-Cholesterol trait results in closer detail. We used 50,497 cases and 100,994
77 controls. After Bonferroni correction, MulTiXcan was able to detect a larger number of significantly
78 associated genes (251) than PrediXcan using all tissues (196) or only a single tissue (whole blood, 33). 172

79 genes were detected by both PrediXcan and MulTiXcan. Fig. 2-a compares the number of significantly
80 associated genes in MulTiXcan, and PrediXcan both for a single tissue (whole blood) and all tissues.
81 Fig. 2-b shows the QQ-plot for associations in these three approaches. There are 79 genes associated to
82 high cholesterol via MulTiXcan and not PrediXcan. Among them, we find many genes related to lipid
83 metabolism (*APOM* [14], *PAFAH1B2* [15]), glucose transport(*SLC5A6* [16]), and vascular processes
84 (*NOTCH4* [17]).

85 **MulTiXcan Results Can Be Inferred From GWAS Summary Results**

86 To expand the applicability of our method to massive sample sizes and to studies where individual
87 level data are not available, we have extended our method to use GWAS summary results rather than
88 individual-level data.

89 We call this extension Summary-MulTiXcan (S-MulTiXcan). Here we derive an analytic expression
90 that relates the joint regression estimates (g_j) to the marginal regression estimates (γ_j as obtained from
91 S-PrediXcan), assuming a known LD structure from a reference panel. We display a conceptual overview
92 of the method in Fig. 4-a. See details in the Methods Section.

93 Figure 3 displays a few examples of the general agreement between the individual-level MulTiXcan and
94 S-MulTiXcan. The summary-based version's results tend to be slightly less significant than MulTiXcan,
95 as illustrated in Supplementary Figure 2.

96 To reduce false positives due to LD misspecification, we discard any significant association result
97 for a gene when the best single tissue result has p-value greater than 10^{-4} . This is rather conservative
98 since it is possible that evidence with modest significance from weakly correlated tissues can lead to very
99 significant combined association. We have, in fact, found several such instances using individual level
100 data.

101 **Application to a broad set of complex traits**

102 We applied S-MulTiXcan to over a 100 traits on publicly available GWAS. As with the individual level
103 method, we observed that S-MulTiXcan detects more associations than S-PrediXcan in most cases, as
104 shown in figure 4-b, after discarding suspicious associations. We also show the QQ-plots and total number
105 of detected association for a sample trait (Schizophrenia) on Figure 4-c and 4-d.

106 These results have been incorporated into the publicly available catalog at <http://gene2pheno.org>.

107 The list of analyzed traits can be found in Supplementary Data 4. Supplementary Data 5 contains a
108 summary of significant associations for each trait. Supplementary Data 6 lists the significant S-MulTiXcan
109 results for each trait.

110 **Highlight Of Associations Identified By Summary-MulTiXcan**

111 We examined the biological relevance of some of the genes detected by our new method that was missed
112 by looking at one tissue at a time (S-PrediXcan).

113 For example, in the Early Growth Genetics (EGG) Consortium’s Body-Mass Index (BMI) study,
114 S-MulTiXcan detects three genes not significant in S-PrediXcan: *POMC* (p-value= 1.4×10^{-6} , tied to
115 childhood obesity [18]); *RACGAP1* (p-value= 1.2×10^{-10} ; embryogenesis [19], cell growth and differentia-
116 tion, [20]); and *TUBA1B* (p-value= 1.23×10^{-09} , circadian cycle processes and psychological disorders [21],
117 suggesting a behavioral pathway).

118 In the CARDIoGRAM+C4D Coronary Artery Disease (CAD) study, S-MulTiXcan detected 12 as-
119 sociations not significant in S-PrediXcan. The top result was *AS3MT* (p-value= 4.3×10^{-9}), related
120 to arsenic metabolism; interestingly, environmental and toxicological studies link arsenic exposure and
121 *AS3MT* polymorphisms with cardiovascular disease [22, 23]. Associations previously linked to CAD in-
122 cluded *CDKN2B* (p-value $< 1.0 \times 10^{-6}$, [24]) *HECTD4* (p-value $< 2.3 \times 10^{-6}$, [25]). Other interesting
123 S-MulTiXcan findings were *CLCC1* (pvalue= 1.2×10^{-7} , a gene for chloride channel activity); *IREB2*
124 (p-value= 2.1×10^{-7} , recently linked to pulmonary conditions, [26]), and *ADAM15* (p-value= 2.5×10^{-07} ,
125 from the disintegrin and metalloproteinase family, linked to atherosclerosis [27], atrial fibrillation [28],
126 and other vascular processes [29, 30]).

127 The list of significant S-MulTiXcan and S-PrediXcan results for all traits can be found in Supplemen-
128 tary Data 6 and 7.

129 **Discussion**

130 Motivated by the widespread sharing of regulatory processes across tissues [12], we propose here a method
131 (MulTiXcan) that aggregates information across multiple tissues and improves the identification of genes
132 significantly associated with complex traits. To expand the applicability of our approach we have extended
133 the method to accommodate GWAS studies where only summary results are available (S-MulTiXcan).

134 We show through applications to hundreds of traits the performance of both individual and summary
135 based methods. We also show that the summary based method provides a reasonably good approximation
136 to the individual level results.

137 As any method relying on a reference panel, S-MulTiXcan might be inaccurate when the study
138 population has a different Linkage Disequilibrium (LD) structure than the reference panel. For example,
139 should two models for a gene yield predicted expressions that are lowly correlated in the reference panel
140 but highly correlated in the study population, then this method underestimates their correlation. To avoid
141 this misspecification, a reference panel matching the study population should be used when available (i.e.
142 using East Asian population from 1000 Genomes if the study set is composed of East Asian individuals).

143 A limitation of PrediXcan and S-PrediXcan is LD contamination, i.e. when causal loci for the trait
144 and expression are different but in LD. We have addressed this in S-PrediXcan through an additional
145 colocalization filtering step. For MulTiXcan, this could be avoided by restricting the analysis to gene-
146 tissue pairs with high colocalization probability.

147 Here we showed the advantages of our joint estimation method through application to multiple traits
148 with publicly available GWAS results as well as the ones available in the UK Biobank. These results
149 include many novel associations of interest, which we make publicly available to the research community
150 in <http://gene2pheno.org>.

151 **Software And Resources**

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

157 **Methods**

158 **Definitions, Notation And Preliminaries**

159 Let us consider a GWAS study of n samples, and assume availability of prediction models in p different
160 tissues. Each model j is a collection of prediction weights w_i^j .

161 Let:

- 162 • \mathbf{y} be an n -vector of phenotypes, assumed to be centered for convenience.
- 163 • \mathbf{X} the genotype matrix, where each column X_l is the n -vector genotype for SNP l . We assume it
164 coded in the range $[0,2]$ but it can be defined in another range, or normalized.
- 165 • Let $\tilde{\mathbf{t}}_j = \sum_{i \in \text{model}_j} w_i^j X_i$ be the predicted expression for model j . Let \mathbf{t}_j be the standardization of
166 $\tilde{\mathbf{t}}_j$.

167 In our application, we will consider $p = 44$ models for a given gene's expression, trained on GTEx
168 data. This method is easily extensible to support incorporation of other covariates, or correction by them.

169 MulTiXcan

170 MulTiXcan consists of fitting a linear regression of the phenotype on predicted expression from multiple
171 tissue models jointly:

$$\begin{aligned} \mathbf{y} &= \sum_{j=1}^p \mathbf{t}_j g_j + \mathbf{e} \\ &= \mathbf{T}\mathbf{g} + \mathbf{e}, \end{aligned} \tag{2}$$

172 where \mathbf{y} is a vector of phenotypes for n individuals, \mathbf{t}_j is an n -vector of normalized predicted gene
173 expression for model j , g_j is the effect size for the predicted gene expression j , \mathbf{e} is an error term, and p
174 is the number of tissues; thus \mathbf{T} is a data matrix where each column j contains the values from \mathbf{t}_j , and
175 \mathbf{g} is the p -vector of effect sizes g_j . One of this columns is a constant intercept term.

176 The high degree of eQTL sharing between different tissues induces a high correlation between pre-
177 dicted expression levels. In order to avoid collinearity issues and numerical instability, we decompose the
178 predicted expression matrix into principal components and keep only the eigenvectors of non negligible
179 variance. To select the number of components, we used a condition number threshold of $\frac{\lambda_{\max}}{\lambda_i} < 30$, where
180 λ_i is an eigenvalue of the matrix $\mathbf{T}^t \mathbf{T}$. A range of values between 10 and 100 yielded similar results. We
181 use an F-test to quantify the significance of the joint fit.

182 We use Bonferroni correction to determine the significance threshold. For MulTiXcan, we use the total
183 number of genes with a prediction model in at least one tissue, which yields a threshold approximately at
184 $0.05/17500 \sim 2.9 \times 10^{-6}$. For PrediXcan across all tissues, we use the total number of gene-tissue pairs,
185 which yields a threshold approximately at $0.05/200,000 \sim 2.5 \times 10^{-7}$.

186 Application To UK Biobank Data

187 We used the same covariates reported in [31], which include the first ten genotype principal components,
188 sex, age, genotyping array, and depending on the trait others such as body mass index (BMI), weight
189 or height. We used 44 models trained on GTEx tissues from release version v6p. For diseases, we used
190 twice as many healthy individuals as controls, selected at random.

191 Summary-MulTiXcan

192 We have demonstrated that S-PrediXcan can accurately infer PrediXcan results from GWAS Summary
193 Statistics and LD information from a reference panel [11], with the added benefits of reduced computa-
194 tional and regulatory burden. Here we extend MulTiXcan in a similar fashion.

195 Summary-MulTiXcan (S-MulTiXcan) infers the individual-level MulTiXcan results, using univariate
196 S-PrediXcan results and LD information from a reference panel. It consists of the following steps:

- 197 • Computation of single tissue association results with S-PrediXcan.
- 198 • Estimation of the correlation matrix of predicted gene expression for the models using the Linkage
199 Disequilibrium (LD) information from a reference panel (typically GTEx or 1000 Genomes [32])
- 200 • Discarding components of smallest variation from this correlation matrix to avert collinearity and
201 numerical problems (Singular Value Decomposition, analogue to PC analysis in individual-level
202 data).
- 203 • Estimation of joint effects from the univariate (single-tissue) results and expression correlation.
- 204 • Discarding suspicious results, suspect to be false positives arising from LD-structure mismatch.

205 Joint Analysis Estimation From Marginal Effects

206 To derive the multivariate regression (2) effect sizes and variances using the marginal regression (3)
207 estimates, we employ a technique presented in [33]. We use the phenotype variance as a conservative
208 approximation to the residual variance σ_e .

209 More specifically, we want to obtain the multivariate regression coefficient estimates for g_j (2) using
210 the estimates from the marginal regression:

$$\mathbf{y} = \mathbf{t}_j \gamma_j + \epsilon. \quad (3)$$

211 where we assume \mathbf{y} centered for convenience.

First, notice that the solution to the multivariate regression in eq. (2) is

$$\hat{\mathbf{g}} = (\mathbf{T}^t \mathbf{T})^{-1} \mathbf{T}^t \mathbf{y} \quad (4)$$

$$\text{var}(\hat{\mathbf{g}}) = \sigma_e^2 (\mathbf{T}^t \mathbf{T})^{-1} \quad (5)$$

, whereas the solution to the marginal regression in eq. (3) is:

$$\hat{\gamma} = \mathbf{D}^{-1} \mathbf{T}^t \mathbf{y} \quad (6)$$

$$\text{var}(\hat{\gamma}) = \sigma_e^2 \mathbf{D}^{-1} \quad \text{with } \mathbf{D} = \text{diag}(\mathbf{T}^t \mathbf{T}) \quad (7)$$

212 From (6) we get $\mathbf{T}^t \mathbf{y} = \mathbf{D} \hat{\gamma}$, which we replace in (4) and obtain the relationship between marginal
213 and joint estimates:

$$\hat{\mathbf{g}} = (\mathbf{T}^t \mathbf{T})^{-1} \mathbf{D} \hat{\gamma} \quad (8)$$

214 To compute the variance of the estimated effect sizes (5) we use the variance of the phenotype as a
215 conservative estimate of σ_e^2 and LD information from reference samples as described next.

216 Estimating Expression Correlation From A Reference Panel

217 As the genotypes from most GWAS are typically unavailable, we must use a reference panel to compute
218 $\mathbf{T}^t \mathbf{T}$, using only those SNPS available in the GWAS results. To do so, notice that:

$$\begin{aligned}
 (\mathbf{T}^t \mathbf{T})_{ij} &= \text{Cor}(\mathbf{t}_i, \mathbf{t}_j) \\
 &= \text{Cov}(\mathbf{t}_i, \mathbf{t}_j) \\
 &= \text{Cor}(\tilde{\mathbf{t}}_i, \tilde{\mathbf{t}}_j) \\
 &= \frac{\text{Cov}(\tilde{\mathbf{t}}_i, \tilde{\mathbf{t}}_j)}{\sqrt{\widehat{\text{var}}(\tilde{\mathbf{t}}_i) \widehat{\text{var}}(\tilde{\mathbf{t}}_j)}} \\
 &= \frac{\text{Cov}\left(\sum_{a \in \text{model}_i} w_a^i X_a, \sum_{b \in \text{model}_j} w_b^j X_b\right)}{\sqrt{\widehat{\text{var}}(\tilde{\mathbf{t}}_i) \widehat{\text{var}}(\tilde{\mathbf{t}}_j)}} \\
 &= \frac{\sum_{\substack{a \in \text{model}_i \\ b \in \text{model}_j}} w_a^i w_b^j \text{Cov}(X_a, X_b)}{\sqrt{\widehat{\text{var}}(\tilde{\mathbf{t}}_i) \widehat{\text{var}}(\tilde{\mathbf{t}}_j)}} \\
 &= \frac{\sum_{\substack{a \in \text{model}_i \\ b \in \text{model}_j}} w_a^i w_b^j \Gamma_{ab}}{\sqrt{\widehat{\text{var}}(\tilde{\mathbf{t}}_i) \widehat{\text{var}}(\tilde{\mathbf{t}}_j)}}, \tag{9}
 \end{aligned}$$

where Γ_{ij} are the elements of the covariance matrix $\mathbf{\Gamma} = \widehat{\text{var}}(\mathbf{X}) = (\mathbf{X} - \bar{\mathbf{X}})^t (\mathbf{X} - \bar{\mathbf{X}}) / n$. We compute the variances as in the S-PrediXcan analysis:

$$\begin{aligned}
 \widehat{\text{var}}(\tilde{\mathbf{t}}_j) &= \hat{\sigma}_j^2 \\
 &= (\mathbf{W}^j)^t \mathbf{\Gamma}^j \mathbf{W}^j \\
 &= \sum_{\substack{a \in \text{model}_j \\ b \in \text{model}_j}} w_a^j w_b^j \Gamma_{ab} \tag{10}
 \end{aligned}$$

219 Addressing Singularity Of The Correlation Matrix

220 Given the high degree of correlation among many of the prediction models, $\mathbf{T}^t \mathbf{T}$ is close to singular
 221 and its inverse cannot be reliably calculated for many genes. To address this problem, we compute the
 222 pseudo-inverse via Singular Value Decomposition, decomposing the covariance matrix into its principal
 223 components and removing those with small eigenvalues. In other terms, we will restrict the analysis to
 224 axes or largest variation of the expression data. This is analogous to the principal components-based
 225 approach used with individual level data. We denote with Σ^+ the pseudo-inverse for any matrix S . We
 226 use the same condition number from individual-level MultiXcan ($\frac{\lambda_{\max}}{\lambda_i} < 30$) as threshold.

227 Estimating Significance

To quantify significance, we use the fact that the regression coefficient estimates follow a (approximate) multivariate normal distribution: $\hat{\mathbf{g}} \sim \mathcal{N}(\hat{\mathbf{g}}, \sigma^2 (\mathbf{T}^t \mathbf{T})^{-1})$. Under the null hypothesis of no association, it follows that $\hat{\mathbf{g}}^t \frac{\mathbf{T}^t \mathbf{T}}{\sigma^2} \hat{\mathbf{g}} \sim \chi_p^2$. We can then replace $\hat{\mathbf{g}}$ with its estimate from the marginal regression:

$$\begin{aligned} \frac{\hat{\mathbf{g}}^t (\mathbf{T}^t \mathbf{T}) \hat{\mathbf{g}}}{\sigma_e^2} &= \frac{\hat{\boldsymbol{\gamma}}^t \mathbf{D} (\mathbf{T}^t \mathbf{T})^{-1} \mathbf{T}^t \mathbf{T} (\mathbf{T}^t \mathbf{T})^{-1} \mathbf{D} \hat{\boldsymbol{\gamma}}}{\sigma_e^2} \\ &= \frac{\hat{\boldsymbol{\gamma}}^t \mathbf{D}}{\sigma_e} (\mathbf{T}^t \mathbf{T})^{-1} \frac{\mathbf{D} \boldsymbol{\gamma}^t}{\sigma_e} \\ &\approx \hat{\mathbf{z}}^t (\mathbf{T}^t \mathbf{T})^{-1} \hat{\mathbf{z}}, \end{aligned}$$

228 where \mathbf{z} is the p -vector of marginal analysis z-scores, $\gamma_j/se(\gamma_j)$. We have used $\sigma_Y^2 \approx \sigma_e^2 \approx \sigma_{\epsilon_j}^2$ in the
229 last step as an approximation. This simplification is conservative, and based on our comparison to the
230 individual multivariate results we consider the loss of efficiency acceptable.

231 Implementation And Computation

232 Prediction Models were obtained from PredictDB.org resource. These models were trained using Elastic
233 Net as implemented in R's package *glmnet* [34], with a mixing parameter $\alpha = 0.5$, over 44 tissue studies
234 from GTEx' release version 6p. The underlying GTEx study data was obtained from dbGaP with accession
235 number phs000424.v6.p1. Please see [11] for details. We implemented MulTiXcan and S-MulTiXcan
236 working up from existing software in the MetaXcan package.

237 UK Biobank genotype data for 487,409 individuals was downloaded and processed in the Bionimbus
238 Protected Data Cloud (PDC), a secure biomedical cloud operated at FISMA moderate as IaaS with an
239 NIH Trusted Partner status for analyzing and sharing protected datasets. We computed GWAS results
240 using BGENIE, a program for efficient GWAS for multiple continuous traits [35]. We selected 222 traits
241 available for these individuals, covering continuous phenotypes such as height and self reported diseases
242 such as asthma. We used different covariate groups for these phenotypes as in [31]. Age, sex and the
243 top ten principal components were used in all cases. For diseases, we randomly sampled twice as many
244 healthy controls as there were cases. Gene expression prediction was computed on the genotype data
245 using the 44 GTEx models.

246 When running MulTiXcan, we used the same covariates and data as in the GWAS. On most continuous

247 phenotypes, there were between 300,000 and 400,000 individuals with available data determined by the
248 intersection of covariates and traits. For the case of self reported diseases, we found a number of cases
249 ranging from a few hundreds (i.e. Acne) to 50,000 (i.e. High Cholesterol). We also ran S-PrediXcan on
250 105 public GWAS traits (the same analyzed in [11], see Supplementary Data 4 for details).

251 Acknowledgments

252 Grants

253 We acknowledge the following US National Institutes of Health grants: R01MH107666 (H.K.I.), R01
254 MH101820 (GTEx)

255 The Genotype-Tissue Expression (GTEx) Project was supported by the Common Fund of the Office
256 of the Director of the National Institutes of Health. Additional funds were provided by the NCI, NHGRI,
257 NHLBI, NIDA, NIMH, and NINDS. Donors were enrolled at Biospecimen Source Sites funded by NCI
258 SAIC-Frederick, Inc. (SAIC-F) subcontracts to the National Disease Research Interchange (10XS170),
259 Roswell Park Cancer Institute (10XS171), and Science Care, Inc. (X10S172). The Laboratory, Data
260 Analysis, and Coordinating Center (LDACC) was funded through a contract (HHSN268201000029C)
261 to The Broad Institute, Inc. Biorepository operations were funded through an SAIC-F subcontract to
262 Van Andel Institute (10ST1035). Additional data repository and project management were provided
263 by SAIC-F (HHSN261200800001E). The Brain Bank was supported by a supplements to University
264 of Miami grants DA006227 & DA033684 and to contract N01MH000028. Statistical Methods devel-
265 opment grants were made to the University of Geneva (MH090941 & MH101814), the University of
266 Chicago (MH090951, MH090937, MH101820, MH101825), the University of North Carolina - Chapel Hill
267 (MH090936 & MH101819), Harvard University (MH090948), Stanford University (MH101782), Washing-
268 ton University St Louis (MH101810), and the University of Pennsylvania (MH101822). The data used for
269 the analyses described in this manuscript were obtained from dbGaP accession number phs000424.v6.p1
270 on 06/17/2016.

271 This work was completed in part with resources provided by Bionimbus [36], and the Center for
272 Research Informatics. The Center for Research Informatics is funded by the Biological Sciences Division
273 at the University of Chicago with additional funding provided by the Institute for Translational Medicine,
274 CTSA grant number UL1 TR000430 from the National Institutes of Health.

275 **Author Contributions**

276 A.N.B. contributed to method development, automated its execution on UK Biobank and public GWAS
277 traits, performed analysis, contributed to the text and figures. M.D.P. processed UK Biobank data,
278 generated predicted expression and contributed to the text. J.Z. processed UK Biobank data and reviewed
279 the manuscript. D.L.N. contributed to the main text and analysis. H.E.W. reviewed the manuscript and
280 contributed to the text and figures. H.K.I. contributed to method development, supervised the project,
281 performed analysis, contributed to the text and figures.

282 **Figures**

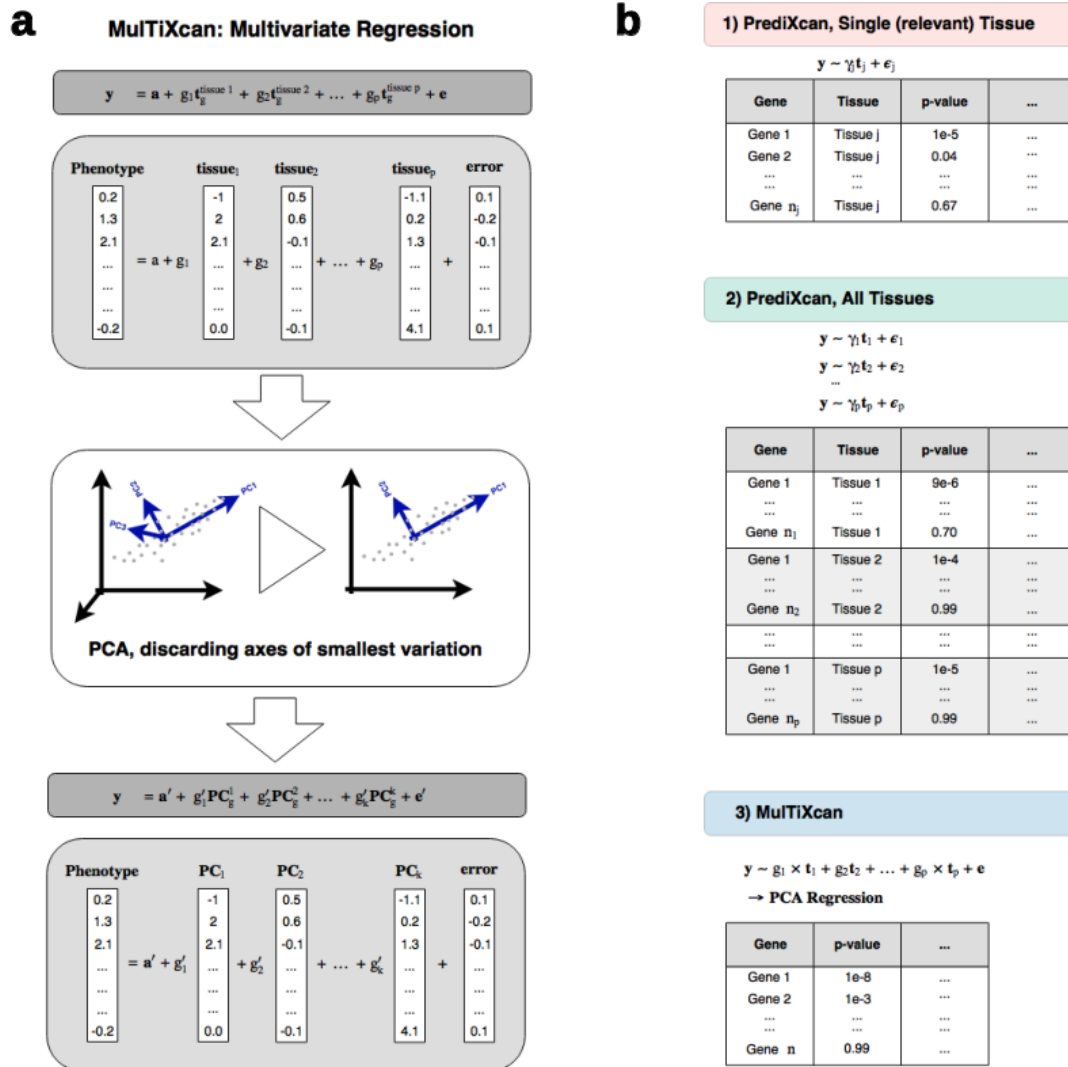


Figure 1. MuTiXcan method.

Panel a illustrates MuTiXcan method. Predicted expression from all available tissue models are used as explanatory variables. To avoid multicollinearity, we use the first k Principal Components of the predicted expression. y is a vector of phenotypes for n individuals, $t_g^{\text{tissue } j}$ is the normalized predicted gene expression for tissue j , g_j is its effect size, a is an intercept and e is an error term.

Panel b shows a schematic representation of MuTiXcan results compared to classical PrediXcan, both for a single relevant tissue and all available tissues in agnostic scanning. y is a vector of phenotypes for n individuals, t_j is the standardized predicted gene expression for model j , g_j is its effect size in the joint regression, γ_j is its effect size in the marginal regression using only prediction j , e and ϵ_j are error terms.

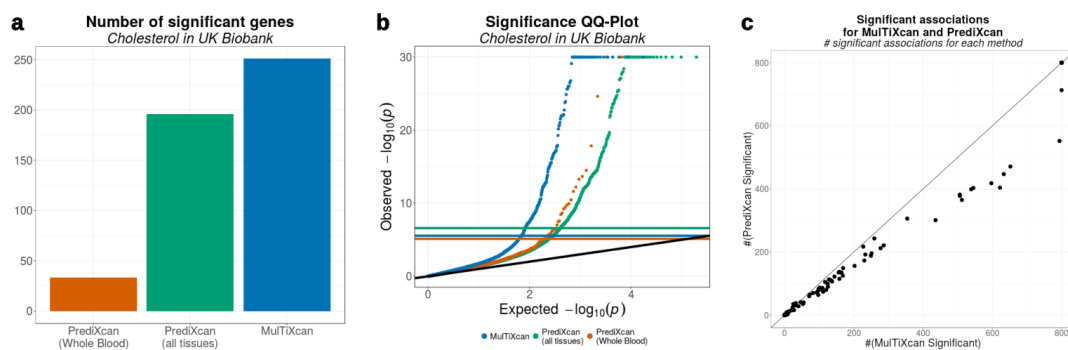


Figure 2. Joint testing across all tissues increases number of significant genes.

Panel a shows the number of discoveries in each method for Cholesterol trait. MuTiXcan is able to detect more findings (251 significant associations) than either of PrediXcan approaches (33 using only Whole Blood and 196 using all 44 GTEx tissues).

Panel b compares the distribution of MuTiXcan's p-values to PrediXcan's p-values for the Cholesterol trait in the UK Biobank cohort. Both PrediXcan with a single tissue model (GTEx Whole Blood) and 44 models (GTEx v6p models) are shown. Notice that Bonferroni-significance levels are different for each case, since 6588 genes were tested in PrediXcan for Whole Blood, 195532 gene-tissue pairs for all GTEx tissues, and 17434 genes in MuTiXcan. P-values were truncated at 10^{-30} for visualization convenience.

Panel c compares the number of significant associations discovered by MuTiXcan and PrediXcan for 222 traits from UK Biobank. These numbers were thresholded at 800 for visualization purposes.

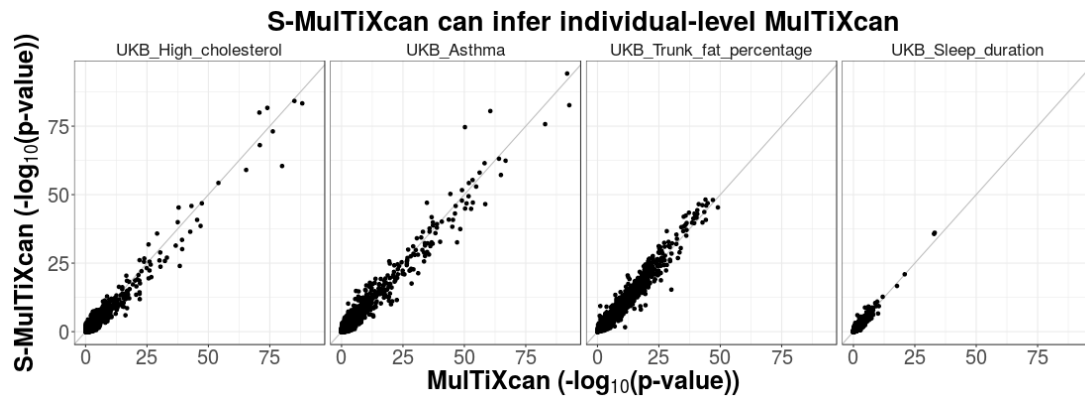


Figure 3. MulTiXcan results can be inferred from GWAS summary statistics and a reference panel. This figure compares S-MulTiXcan to MulTiXcan in three UK Biobank phenotypes. GTEx individuals were used as a reference panel for estimating expression correlation in the study population. The summary data-based method shows a good level of agreement with the individual-based method. In cases where the LD-structure between reference and study cohorts is mismatched, the summary-based method becomes less accurate. For example in Asthma, two genes are significantly overestimated; however it tends to be conservative for most genes.

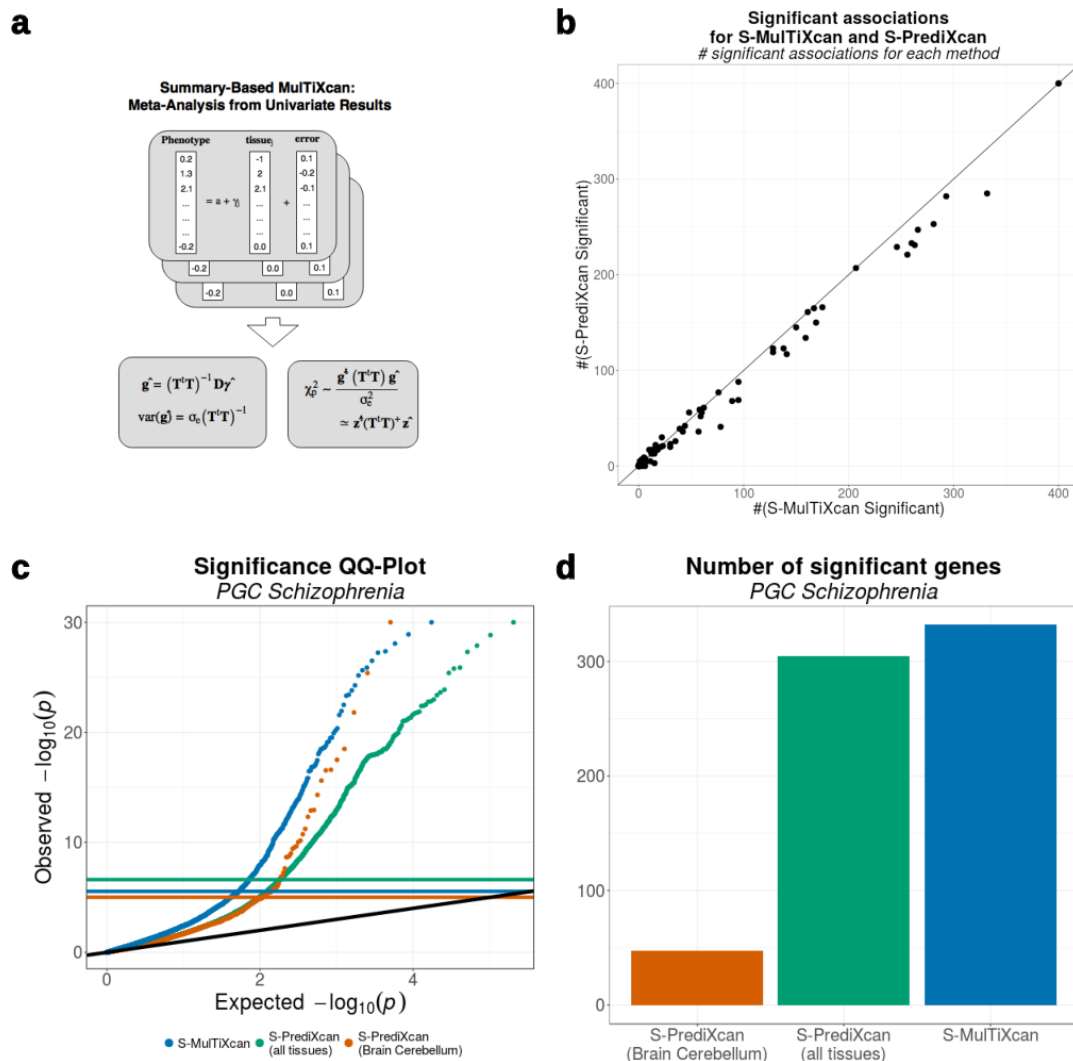


Figure 4. Comparison between S-PrediXcan and S-MultiXcan.

Panel a illustrates the S-MultiXcan method: the marginal univariate S-PrediXcan effect sizes are computed, then the joint effect sizes are estimated from them. The significance is quantified through an omnibus test.

Panel b compares the number of associations significant via S-MultiXcan versus those significant via S-PrediXcan, for the same GWAS Studies. In most cases, S-MultiXcan detects a larger number of exclusive significant associations. The number of discoveries was thresholded at 200 for visualization purposes.

Panel c displays QQ-Plots for the association p-values from S-MultiXcan and S-PrediXcan in Schizophrenia, using a model trained on brain's cerebellum, and S-PrediXcan associations for all 44 GTEx tissues.

Panel d shows the number of significant associations in Schizophrenia for each method as a bar plot.

283 **Tables**

284 **References**

- 285 1. Smoller JW, Craddock N, Kendler K, Lee PH, Neale BM, Nurnberger JI, et al. Identifica-
286 tion of risk loci with shared effects on five major psychiatric disorders: a genome-wide analy-
287 sis. *Lancet*. 2013;381(9875):1371–9. Available from: [http://discovery.ucl.ac.uk/1395494/\\$\delimitter"026E30F\\$nhttp://www.ncbi.nlm.nih.gov/pubmed/23453885](http://discovery.ucl.ac.uk/1395494/$\delimitter).
288
- 289 2. Deloukas P, Kanoni S, Willenborg C, Farrall M, Assimes TL, Thompson JR, et al. Large-
290 scale association analysis identifies new risk loci for coronary artery disease. *Nature genetics*.
291 2013;45(1):25–33. Available from: [http://www.pubmedcentral.nih.gov/articlerender.fcgi?
292 artid=3679547{%&}tool=pmcentrez{%&}rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=3679547{%&}tool=pmcentrez{%&}rendertype=abstract).
- 293 3. Morris AP, Voight BF, Teslovich TM, Ferreira T, Segrè AV, Steinthorsdottir V, et al.
294 Large-scale association analysis provides insights into the genetic architecture and patho-
295 physiology of type 2 diabetes. *Nature Genetics*. 2012;44(9):981–990. Available
296 from: [http://www.ncbi.nlm.nih.gov/pubmed/22885922\\$\delimitter"026E30F\\$nhttp://www.
297 nature.com/doifinder/10.1038/ng.2383](http://www.ncbi.nlm.nih.gov/pubmed/22885922$\delimitter).
- 298 4. Moreau Y, Tranchevent LC. Computational tools for prioritizing candidate genes: boosting disease
299 gene discovery. *Nature Reviews; Genetics*. 2012;13(8):523–536.
- 300 5. Cantor RM, Lange K, Sinsheimer JS. Prioritizing GWAS Results: A Review of Statistical Methods
301 and Recommendations for Their Application. *American Journal of Human Genetics*. 2010;86(1):6–
302 22.
- 303 6. Nica AC, Montgomery SB, Dimas AS, Stranger BE, Beazley C, Barroso I, et al. Candidate causal
304 regulatory effects by integration of expression QTLs with complex trait genetic associations. *PLoS
305 Genetics*. 2010;6(4).
- 306 7. Nicolae DL, Gamazon E, Zhang W, Duan S, Eileen Dolan M, Cox NJ. Trait-associated SNPs are
307 more likely to be eQTLs: Annotation to enhance discovery from GWAS. *PLoS Genetics*. 2010;6(4).

- 308 8. Li YI, van de Geijn B, Raj A, Knowles DA, Petti AA, Golan D, et al. RNA splicing is a primary
309 link between genetic variation and disease. *Science*. 2016;352(6285):600–604. Available from: <http://www.ncbi.nlm.nih.gov/pubmed/27126046>.
310
- 311 9. Gusev A, Lee SH, Trynka G, Finucane H, Vilhjálmsson BJ, Xu H, et al. Partitioning heritability of
312 regulatory and cell-type-specific variants across 11 common diseases. *American Journal of Human
313 Genetics*. 2014;95(5):535–552.
- 314 10. Gamazon ER, Wheeler HE, Shah KP, Mozaffari SV, Aquino-Michaels K, Carroll RJ, et al. A gene-
315 based association method for mapping traits using reference transcriptome data. *Nature genetics*.
316 2015;47(9):1091–1098. Available from: <http://dx.doi.org/10.1038/ng.3367>.
- 317 11. Barbeira A, Dickinson SP, Torres JM, Bonazzola R, Zheng J, Torstenson ES, et al. Integrating
318 tissue specific mechanisms into GWAS summary results. *bioRxiv*. 2017; Available from: <http://www.biorxiv.org/content/early/2017/05/21/045260>.
319
- 320 12. Aguet F, Brown AA, Castel S, Davis JR, Mohammadi P, Segre AV, et al. Local genetic effects
321 on gene expression across 44 human tissues. *bioRxiv*. 2016; Available from: [http://biorxiv.org/
322 content/early/2016/09/09/074450](http://biorxiv.org/content/early/2016/09/09/074450).
- 323 13. Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, et al. UK Biobank: An Open Access
324 Resource for Identifying the Causes of a Wide Range of Complex Diseases of Middle and Old Age.
325 *PLoS Medicine*. 2015;12(3).
- 326 14. Xu N, Dahlbäck B. A novel human apolipoprotein (apoM). *The Journal of biological chem-*
327 *istry*. 1999;274(44):31286–90. Available from: [http://www.jbc.org.ezproxy.lib.ucalgary.ca/
328 content/274/44/31286.full.pdf](http://www.jbc.org.ezproxy.lib.ucalgary.ca/content/274/44/31286.full.pdf)<http://www.ncbi.nlm.nih.gov/pubmed/10531326>.
- 329 15. Peloso GM, Auer PL, Bis JC, Voorman A, Morrison AC, Stitzel NO, et al. Association of low-
330 frequency and rare coding-sequence variants with blood lipids and coronary heart disease in 56,000
331 whites and blacks. *American Journal of Human Genetics*. 2014;94(2):223–232.
- 332 16. Wright EM, Turk E. The sodium/glucose cotransport family SLC5; 2004.
- 333 17. Gridley T. Notch signaling in vascular development and physiology. *Development (Cambridge,
334 England)*. 2007;134(15):2709–2718.

- 335 18. Kuehnen P, Mischke M, Wiegand S, Sers C, Horsthemke B, Lau S, et al. An alu element-associated
336 hypermethylation variant of the POMC gene is associated with childhood obesity. *PLoS Genetics*.
337 2012;8(3).
- 338 19. Grewal S, Carver JG, Ridley AJ, Mardon HJ. Implantation of the human embryo requires Rac1-
339 dependent endometrial stromal cell migration. *Proceedings of the National Academy of Sciences*
340 of the United States of America. 2008;105(42):16189–16194. Available from: [http://eutils.
341 ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=18838676&retmode=
342 ref&cmd=prlinks%5Cnpapers2://publication/doi/10.1073/pnas.0806219105](http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=18838676&retmode=ref&cmd=prlinks%5Cnpapers2://publication/doi/10.1073/pnas.0806219105).
- 343 20. Hallstrom TC, Mori S, Nevins JR. An E2F1-Dependent Gene Expression Program that Determines
344 the Balance between Proliferation and Cell Death. *Cancer Cell*. 2008;13(1):11–22.
- 345 21. Byrne EM, Heath AC, Madden PAF, Pergadia ML, Hickie IB, Montgomery GW, et al. Testing
346 the role of circadian genes in conferring risk for psychiatric disorders. *American Journal of Medical*
347 *Genetics, Part B: Neuropsychiatric Genetics*. 2014;165(3):254–260.
- 348 22. Gong G, O’Byrant SE. Low-level arsenic exposure, AS3MT gene polymorphism and cardiovascular
349 diseases in rural Texas counties. *Environmental Research*. 2012;113:52–57.
- 350 23. Moon K, Guallar E, Navas-Acien A. Arsenic exposure and cardiovascular disease: An updated
351 systematic review; 2012.
- 352 24. Samani NJ, Erdmann J, Hall AS, Hengstenberg C, Mangino M, Mayer B, et al. Genomewide Associ-
353 ation Analysis of Coronary Artery Disease. *New England Journal of Medicine*. 2007;357(5):443–453.
354 Available from: <http://www.nejm.org/doi/abs/10.1056/NEJMoa072366>.
- 355 25. Lu X, Wang L, Chen S, He L, Yang X, Shi Y, et al. Genome-wide association study in Han Chinese
356 identifies four new susceptibility loci for coronary artery disease. *Nature Genetics*. 2012;44(8):890–
357 894.
- 358 26. DeMeo DL, Mariani T, Bhattacharya S, Srisuma S, Lange C, Litonjua A, et al. Integration of
359 Genomic and Genetic Approaches Implicates IREB2 as a COPD Susceptibility Gene. *American*
360 *Journal of Human Genetics*. 2009;85(4):493–502.

- 361 27. Oksala N, Levula M, Airla N, Pelto-Huikko M, Ortiz RM, JÄdrvinen O, et al. ADAM-9, ADAM-15,
362 and ADAM-17 are upregulated in macrophages in advanced human atherosclerotic plaques in aorta
363 and carotid and femoral arteriesÄTTampere vascular study. *Annals of Medicine*. 2009;41(4):279–
364 290. Available from: <http://dx.doi.org/10.1080/07853890802649738>.
- 365 28. Arndt M, Lendeckel U, Röcken C, Nepple K, Wolke C, Spiess A, et al. Altered expression of ADAMs
366 (A Disintegrin And Metalloproteinase) in fibrillating human atria. *Circulation*. 2002;105(6):720–
367 725.
- 368 29. Xie B, Shen J, Dong A, Swaim M, Hackett SF, Wyder L, et al. An Adam15 amplification
369 loop promotes vascular endothelial growth factor-induced ocular neovascularization. *FASEB*
370 *journal : official publication of the Federation of American Societies for Experimental Biol-*
371 *ogy*. 2008;22(8):2775–83. Available from: [http://www.pubmedcentral.nih.gov/articlerender.](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2493454&tool=pmcentrez&rendertype=abstract)
372 [fcgi?artid=2493454&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=2493454&tool=pmcentrez&rendertype=abstract).
- 373 30. Komiya K, Enomoto H, Inoki I, Okazaki S, Fujita Y, Ikeda E, et al. Expression of
374 ADAM15 in rheumatoid synovium: up-regulation by vascular endothelial growth factor
375 and possible implications for angiogenesis. *Arthritis research & therapy*. 2005;7(6):R1158–
376 R1173. Available from: [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1297561&tool=pmcentrez&rendertype=abstract)
377 [1297561&tool=pmcentrez&rendertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=1297561&tool=pmcentrez&rendertype=abstract).
- 378 31. Ge T, Chen CY, Neale BM, Sabuncu MR, Smoller JW. Phenome-wide heritability analysis of the
379 UK Biobank. *PLoS Genetics*. 2017;13(4).
- 380 32. 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM,
381 et al. A global reference for human genetic variation. *Nature*. 2015;526(7571):68–74. Available from:
382 <http://www.nature.com/doifinder/10.1038/nature15393>
383 <http://www.ncbi.nlm.nih.gov/pubmed/26432245>.
- 384 33. Yang J, Ferreira T, Morris AP, Medland SE, Madden PAF, Heath AC, et al. Conditional and
385 joint multiple-SNP analysis of GWAS summary statistics identifies additional variants influencing
386 complex traits. *Nature Genetics*. 2012;44(4):369–375. Available from: [http://www.nature.com/](http://www.nature.com/doifinder/10.1038/ng.2213)
387 [doifinder/10.1038/ng.2213](http://www.nature.com/doifinder/10.1038/ng.2213).

- 388 34. Friedman J, Hastie T, Tibshirani R. Regularization Paths for Generalized Linear Models via
389 Coordinate Descent. *Journal of Statistical Software*. 2010;33(1):1–22. Available from: [http://](http://www.jstatsoft.org/v33/i01/)
390 www.jstatsoft.org/v33/i01/.
- 391 35. Bycroft C, Freeman C, Petkova D, Band G, Elliott LT, Sharp K, et al. Genome-wide genetic
392 data on ~500,000 UK Biobank participants. *bioRxiv*. 2017;p. 166298. Available from: [https:](https://www.biorxiv.org/content/early/2017/07/20/166298)
393 [//www.biorxiv.org/content/early/2017/07/20/166298](https://www.biorxiv.org/content/early/2017/07/20/166298).
- 394 36. Heath AP, Greenway M, Powell R, Spring J, Suarez R, Hanley D, et al. Bionimbus: a cloud
395 for managing, analyzing and sharing large genomics datasets. *Journal of the American Medical*
396 *Informatics Association : JAMIA*. 2014 Nov;21(6):969–975. Available from: [https://academic.](https://academic.oup.com/jamia/article-lookup/doi/10.1136/amiajnl-2013-002155)
397 [oup.com/jamia/article-lookup/doi/10.1136/amiajnl-2013-002155](https://academic.oup.com/jamia/article-lookup/doi/10.1136/amiajnl-2013-002155).

398 **Supplementary Material**

399 **Supplementary Data**

400 **Supplementary Data 1. Summary statistics for UK Biobank traits used in the MulTiXcan**
401 **analysis.** MulTiXcan was run for 222 traits on UK Biobank. Summary statistics for significant results
402 included in **supp-data-ukb-multixcan-stats.txt**. Columns are: **tag**: trait, gene2pheno.org display
403 name; **n_predixcan_significant**: Number of Bonferroni-significant PrediXcan results; **n_MulTiXcan_significant**
404 number of Bonferroni-significant results for MulTiXcan; **n_predixcan_only** number of results only sig-
405 nificant in PrediXcan; **n_MulTiXcan_only** number of results only significant in MulTiXcan.

406 **Supplementary Data 2. Significant associations for MulTiXcan on UK Biobank.** Signifi-
407 cant results included in **supp-data-ukb-multixcan-significant.txt**. Columns are: **phenotype**: trait,
408 gene2pheno.org display name; **gene**: Ensembl id; **gene_name**: HUGO name; **pvalue**: p-value of the
409 S-MulTiXcan association; **n_models** number of prediction models available for the gene; **n_used** num-
410 ber of independent components surviving PCA selection; **n_samples**: number of individuals available.

411
412 **Supplementary Data 3. Significant associations for PrediXcan on UK Biobank.** Significant
413 results included in **supp-data-ukb-p-significant.txt**. Columns are: **Phenotype**: trait, gene2pheno.org
414 display name; **model**: GTEx tissue where the model was trained; **gene**: Ensembl Id; **gene_name**:
415 HUGO name; **model** GTEx tissue where model was trained; **zscore** PrediXcan association Z-score,
416 **pvalue** PrediXcan association p-value; **n_samples**: number of individuals available.

417 **Supplementary Data 4. List of Genome-wide Association Meta Analysis (GWAMA) Con-**
418 **sortia and phenotypes.** Data included in **supp-data-gwas-traits.txt**. Columns are consortium
419 name, study name, gene2pheno.org display name, study sample size, study population, URL of portal
420 where data was downloaded from, link to pubmed entry if available.

421 **Supplementary Data 5. Summary statistics for traits used in the MulTiXcan analysis.**
422 MulTiXcan was run for 105 public GWAS. Summary statistics for significant results included in **supp-**
423 **data-gwas-smultixcan-stats.txt**. Columns are: **tag**: gene2pheno.org display name; **consortium**:
424 Consortium Name; **name**: study name; **n_spredixcan_significant**: Number of Bonferroni-significant
425 S-PrediXcan results; **n_sMulTiXcan_significant** number of Bonferroni-significant results for MulTi-

426 Xcan; **n_spredixcan_only** number of results only significant in S-PrediXcan; **n_sMulTiXcan_only**
427 number of results only significant in S-MulTiXcan.

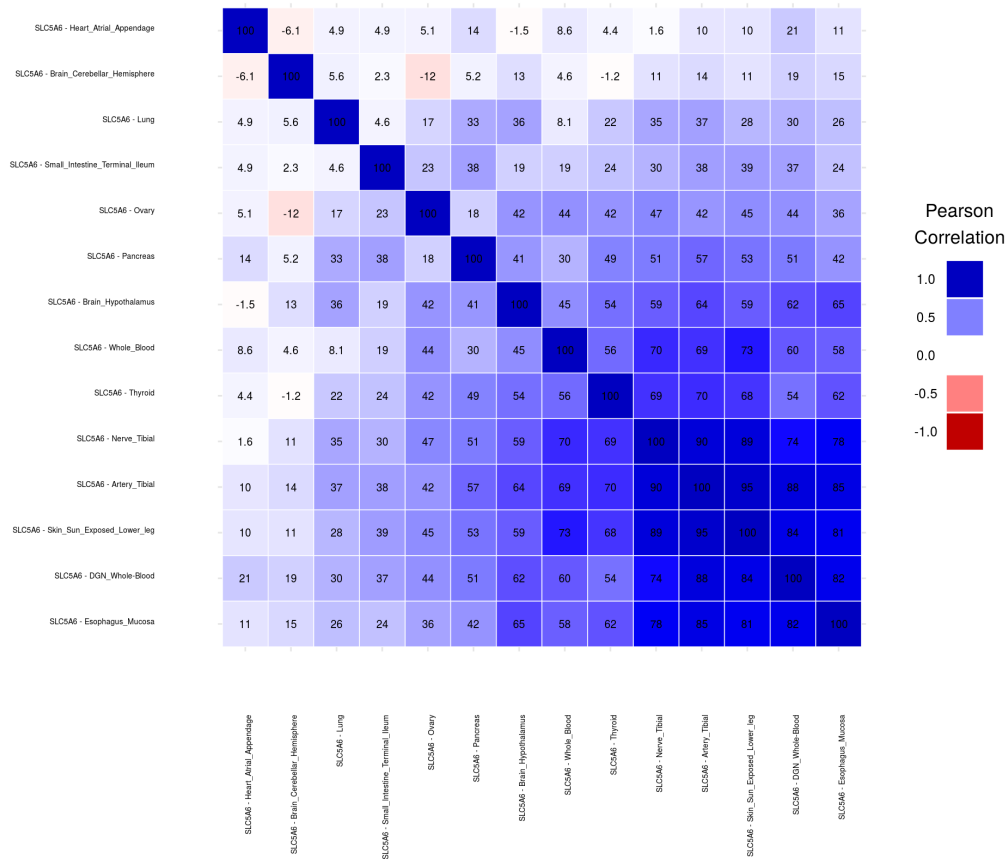
428 **Supplementary Data 6. Significant associations for Summary-MulTiXcan on public GWAS.**

429 Significant results included in **supp-data-gwas-smultixcan-significant.txt**. Columns are: **tag**: gene2pheno.org
430 display name; **consortium**: Consortium Name; **name**: study name; **gene**: Ensembl id; **gene_name**:
431 HUGO name; **pvalue**: p-value of the S-MulTiXcan association; **n** number of S-PrediXcan results avail-
432 able for the gene; **n_indep** number of independent components surviving SVD; **p_i_best** best p-value
433 of S-PrediXcan; **t_i_best** tissue that presented best S-PrediXcan result; **p_i_worst** worst p-value of S-
434 PrediXcan; **t_i_worst** tissue that presented worst S-PrediXcan result; **suspicious**: whether the result
435 was discarded as a potential false positive.

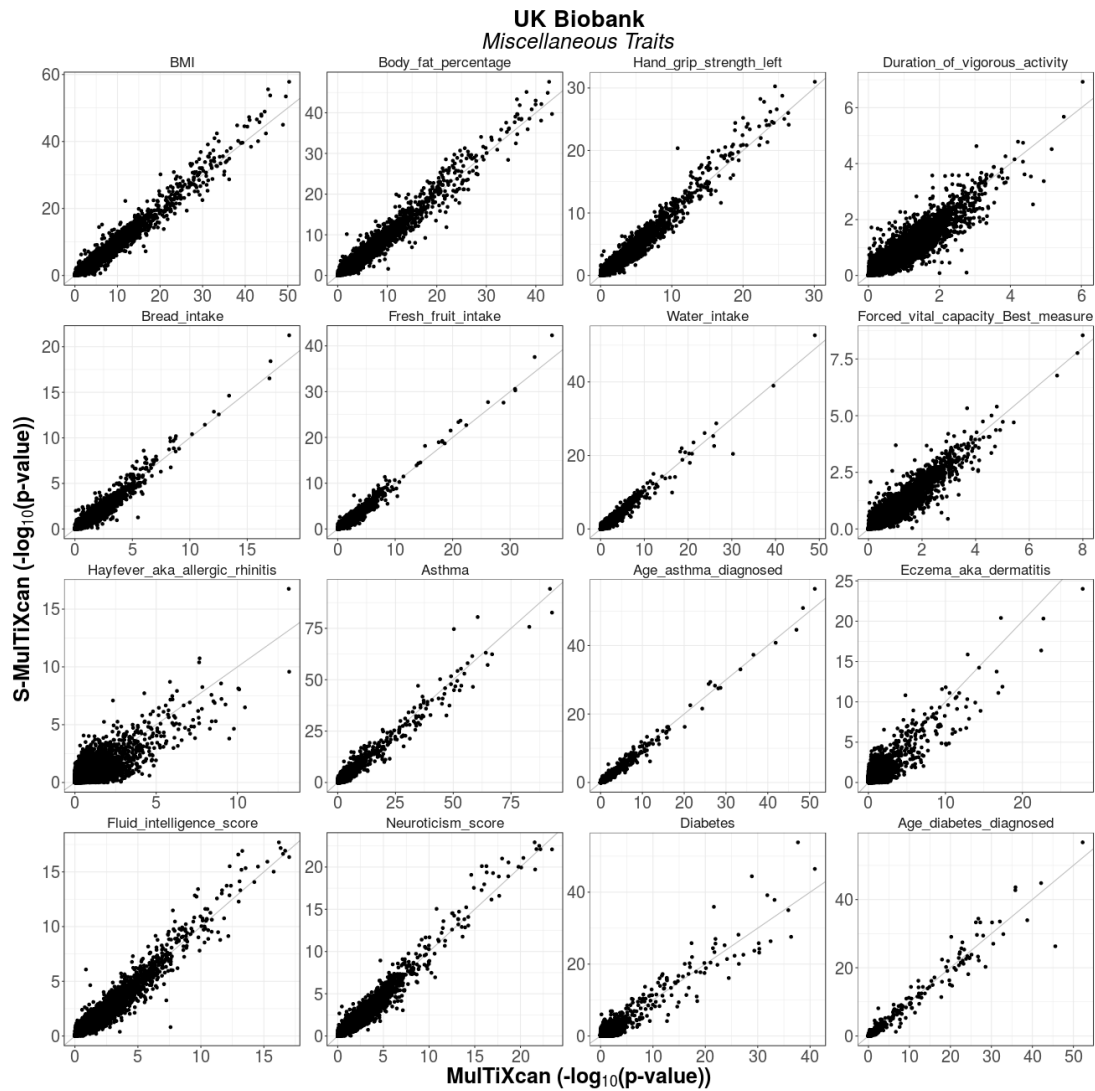
436 **Supplementary Data 7. Significant associations for Summary-PrediXcan on public GWAS.**

437 Significant results included in **supp-data-gwas-sp-significant.txt**. Columns are: **consortium**: Con-
438 sortium Name; **name**: study name; **tag**: gene2pheno.org display name; **gene**: Ensembl Id; **gene_name**:
439 HUGO name; **model** GTEx tissue where model was trained; **zscore** S-PrediXcan association Z-score,
440 **pvalue** S-PrediXcan association p-value.

⁴⁴¹ Supplementary Figures



Supplementary Figure 1. Predicted expression correlation for gene *SLC5A6*. We observe a high degree of predicted expression correlation, in agreement with recent publications on the high degree of mechanism sharing across tissues [12]. This behavior is exhibited in most genes.



Supplementary Figure 2. Summary-MulTiXcan vs MulTiXcan for Miscellaneous Traits.

There is a satisfactory agreement between the individual-level and the summary-level versions of MulTiXcan in UK Biobank traits.