

1 **Uncovering causal gene-tissue pairs and variants: A multivariable** 2 **TWAS method controlling for infinitesimal effects**

3 Yihe Yang¹, Noah Lorincz-Comi¹, Xiaofeng Zhu¹

4 ¹Department of Population and Quantitative Health Sciences, School of Medicine, Case Western
5 Reserve University, Cleveland, OH, 44106, USA

6 Corresponding author: X. Zhu. xxz10@case.edu

7 **Abstract**

8 Transcriptome-wide association studies (TWAS) are commonly used to prioritize causal genes
9 underlying associations found in genome-wide association studies (GWAS) and have been
10 extended to identify causal genes through multivariable TWAS methods. However, recent
11 studies have shown that widespread infinitesimal effects due to polygenicity can impair the
12 performance of these methods. In this report, we introduce a multivariable TWAS method named
13 Tissue-Gene pairs, direct causal Variants, and Infinitesimal effects selector (TGVIS) to identify
14 tissue-specific causal genes and direct causal variants while accounting for infinitesimal effects.
15 In simulations, TGVIS maintains an accurate prioritization of causal gene-tissue pairs and
16 variants and demonstrates comparable or superior power to existing approaches, regardless of the
17 presence of infinitesimal effects. In the real data analysis of GWAS summary data of 45
18 cardiometabolic traits and expression/splicing quantitative trait loci (eQTL/sQTL) from 31
19 tissues, TGVIS is able to improve causal gene prioritization and identifies novel genes that were
20 missed by conventional TWAS.

21 **Introduction**

22 Over the past two decades, genome-wide association studies (GWAS) have identified thousands
23 of genetic variants associated with complex traits¹⁻³. However, most GWAS signals are detected
24 in non-coding regions and have been shown to have complex regulatory landscapes across
25 different tissues and cell types⁴, making it challenging to pinpoint causal variants and genes
26 driving these GWAS signals. Joint GWAS and expression quantitative trait loci (eQTL) data
27 analysis methods, such as colocalization⁵, transcriptome-wide association studies (TWAS)⁶, and
28 *cis*-Mendelian randomization (*cis*-MR)⁷, have been developed to prioritize causal genes at
29 GWAS loci⁸. Colocalization simultaneously examines the expression of a gene and a trait to
30 determine whether they share common causal genetic variants at a locus⁵. Both TWAS and *cis*-
31 MR assume a causal diagram where eQTLs regulate tissue-specific gene expression that
32 subsequently affects a trait, and they identify these tissue-specific causal genes by testing the
33 significance of the causal effect estimates. Furthermore, these methods have been extended to a
34 broader range of molecular phenotypes, such as splicing events⁹ and protein abundance¹⁰, with
35 regulatory QTLs being splicing QTLs (sQTLs) and protein QTLs (pQTLs), which we call
36 xQTLs in general.

37 Nevertheless, colocalization, TWAS, and *cis*-MR are all univariable methods that statistically
38 measure the marginal correlations of genetic effect sizes between a trait and a tissue-specific
39 expression of a gene. Non-causal gene-tissue pairs may be falsely detected by these univariable
40 methods due to the *cis*-gene-tissue co-regulations with causal gene-tissue pairs^{8,11,12}. The
41 underlying mechanism may come in the following respects: the tissue-specific eQTLs of a causal
42 gene are in linkage disequilibrium (LD) with (1) the eQTLs of nearby non-causal genes¹³ and (2)
43 the eQTLs of causal genes expressed in non-causal tissues¹⁴. In addition, some variants can

44 influence a trait independently of causal gene-tissue pairs, which are frequently denoted as direct
45 causal variants¹³ and horizontal pleiotropy¹⁵. The non-causal gene-tissue pairs may be incorrectly
46 detected when their eQTLs are in LD with direct causal variants.

47 Multivariable TWAS methods, such as causal TWAS (cTWAS)¹³ and Tissue-Gene Fine-
48 Mapping (TGFM)¹⁴, have been proposed to address these issues. Specifically, cTWAS identifies
49 causal genes and direct causal variants among multiple candidates using the sum of single effects
50 (SuSiE)^{16,14} by examining tissues separately. TGFM extends cTWAS to allow multiple tissues to
51 be analyzed simultaneously and can identify the trait-relevant tissues beyond the causal variants
52 and genes. However, Cui et al.¹⁷ recently reported that current Bayesian fine-mapping methods,
53 including SuSiE^{16,18} and FINEMAP¹⁹, have a high replication failure rate (RFR) in practice. Cui
54 et al.¹⁷ discovered that the widespread infinitesimal effects, which may stem from the
55 polygenicity of complex traits, are the sources of the high RFR, and accounting for the
56 infinitesimal effects can reduce the RFR and improve statistical power. Similarly, the
57 polygenicity can also lead to inflating the test statistics in standard TWAS²⁰ and traditional
58 linkage studies²¹. Thus, due to the lack of modeling infinitesimal effects, it is expected that
59 cTWAS and TGFM can be vulnerable to spurious prioritization and reduced statistical power.

60 We present the Tissue-Gene pair, direct causal Variants, and Infinitesimal effect Selector
61 (TGVIS), a multivariable TWAS method to identify causal gene-tissue pairs and direct causal
62 variants while incorporating infinitesimal effects. TGVIS employs SuSiE^{16,14} for fine-mapping
63 causal gene-tissue pairs and direct causal variants, and uses restricted maximum likelihood
64 (REML)²² to estimate the infinitesimal effects. In addition, we introduce the Pratt index²³ to rank
65 the importance for improving the prioritization of causal genes and variants. We applied TGVIS

66 to identify causal *cis*-gene-tissue pairs and direct causal variants for 45 cardiometabolic traits
67 using GWAS datasets with the largest sample sizes to date^{3,24–32}, by incorporating the eQTL and
68 sQTL summary statistics from 28 tissues from genotype-tissue expression (GTEx)¹², and the
69 eQTL summary statistics of kidney tubulointerstitial³³, kidney glomerular³³, and pancreatic
70 islets³⁴ tissues. We summarized the causal gene-tissue pairs and direct causal variants,
71 highlighted the pleiotropic effects at the gene-tissue level, and demonstrated the different
72 functional activity³⁵ of eQTLs/sQTLs mediated through gene-tissues and the direct causal
73 variants. Moreover, we mapped the trait-relevant major tissues and demonstrated the
74 enrichments of genes identified by TGVIS in terms of colocalization⁵, on the silver standard of
75 lipid genes¹³, FDA-approved drug-target genes³⁶, and genes detected through pQTL summary
76 data¹⁰. Our study reveals a broader picture of gene and tissue co-regulations, which can provide
77 novel biological insights into complex traits.

78 **Results**

79 **Overview of method**

80 **Figure 1A** illustrates the causal diagram assumed in this report. Specifically, we hypothesize that
81 a set of xQTLs influence the products of genes (e.g., expressions and splicing events) at a locus.
82 Gene co-regulation^{8,12}, i.e., the correlation of xQTL effects among multiple gene products, can
83 emerge due to shared xQTLs or being in LD among them. Meanwhile, tissue co-regulation^{11,37,38},
84 defined as the correlation of gene expression across multiple tissues, can arise because of the
85 same mechanism. In the gene and tissue co-regulation network, certain gene-tissue pairs directly
86 influence a trait without mediation by other gene-tissue pairs, which are referred to as causal
87 gene-tissue pairs. In addition, some genetic variants may directly influence the trait, which we
88 consider as direct causal variants. Besides these direct causal variants which have relatively large

89 effects, we assume there are polygenic or infinitesimal effects that can be modeled through a
90 normal distribution with mean zero and small variance¹⁷ (**Method**).

91 The curse of dimensionality poses a substantial challenge in the multivariable TWAS model.

92 **Figure 1B** illustrates this challenge by an example of the association evidence with low-density

93 lipoprotein cholesterol (LDL-C)¹ at the *PCSK9* locus, where dozens of coding genes and long

94 non-coding ribonucleic acids (RNAs) are located, along with multiple potential direct causal

95 variants. Conventional statistical methods cannot precisely identify causal gene-tissue pairs and

96 variants because there are many correlated candidates which frequently range from hundreds to

97 thousands¹⁶. The proposed TGVIS enables to overcome the curse of dimensionality. **Figure 1C**

98 describes the workflow of TGVIS, where the inputs are the GWAS summary statistics of a

99 trait, xQTL summary statistics of gene-tissue pairs, and a reference LD matrix of the

100 variants at the locus. TGVIS utilizes a profile-likelihood approach to estimate the causal

101 effects of gene-tissue pairs and directly causal variant effects with SuSiE^{16,18} and model the

102 infinitesimal effects via REML²². This profile-likelihood iterates until all estimates are

103 converge. The details are described in the **Methods** and the **Supplementary Materials**.

104 In practice, another challenge arises when selecting a causal gene-tissue pair based solely on its

105 posterior inclusion probability (PIP) because many gene-tissue pairs share the same sets of

106 xQTLs at a locus, making them statistically indistinguishable. SuSiE groups these pairs into a

107 credible set during fine-mapping and introduces a single effect to describe the contributions of

108 the variables in the same credible set. Therefore, all inferences should be made based on the

109 single effects defined by SuSiE's credible sets. To quantify the contribution of each gene-tissue

110 pair and direct causal variant, we introduce the Pratt index²³ as a metric parallel to PIP. While

111 PIP measures the significance of variables from a Bayesian viewpoint, the Pratt index quantifies

112 their predictive importance. In the application, we calculated the cumulative Pratt index of
113 variables in a 95% credible set (CS-Pratt) and filtered out the credible sets with low CS-Pratt
114 values (**Methods** and **Supplementary Materials**). We observed that this procedure improved
115 the precision of causal gene and variant identification.

116 **Simulation**

117 We compared the TGVIS with 4 multivariable MR and TWAS methods: *cis*IVW³⁹, Grant2022⁴⁰,
118 cTWAS¹³, and TGF¹⁴. We applied the following criteria for determining the causality: the
119 95% credible set for TGVIS, TGF¹⁴, and cTWAS; $P < 0.05$ for *cis*IVW; and selection by lasso
120 in the Grant2022. We did not consider univariable methods because of their high type-I error
121 rates when the goal is to identify causal genes, given that xQTL effect sizes for multiple genes
122 are often correlated¹³. Detailed information on the settings and results can be found in the
123 **Methods** and **Supplementary Materials**.

124 We first assessed the accuracy of causal effect estimation for gene-tissue pairs. When
125 infinitesimal effects were absent, TGVIS showed a mean square error (MSE) for causal effect
126 estimates similar to that of cTWAS, and TGF¹⁴, while both *cis*IVW and Grant2022 exhibited
127 substantially larger MSE (as shown in the left two panels in **Figure 2A**). However, when
128 infinitesimal effects were present, TGVIS demonstrated a visibly lower MSE compared to the
129 other methods, with cTWAS and TGF¹⁴ showing approximately 32% higher MSE than TGVIS
130 (as shown in the right two panels in **Figure 2A**). These results indicate that TGVIS generally
131 outperforms its competitors by accounting for infinitesimal effects.

132 We then compared the true negative rate (TNR) and true positive rate (TPR) of these five
133 methods. A true negative is defined as a method that correctly identifies all 98 non-causal gene-

134 tissue pairs. Similarly, a true positive is defined as a method that correctly identifies the 2 causal
135 gene-tissue pairs. Across all the scenarios (**Figure 2B**), TGVIS achieved the highest TNR, with
136 an average of 0.614, followed by TGFm and cTWAS, with average TNRs of 0.513 and 0.499,
137 respectively. *CisIVW* and Grant2022 performed worst, with averages TNR of 0.064 and 0.013,
138 respectively, indicating that these two methods are prone to identifying a substantial number of
139 false positive gene-tissue pairs. On the other hand, TGVIS exhibited a similar TPR (average
140 TPR = 0.667) as TGFm, cTWAS, and *cisIVW* (average TPRs of 0.649, 0.667, and 0.661,
141 respectively), while Grant2022 had the highest TPR (averages TPR= 0.831) (**Figure 2C**), which
142 is not surprising given that Grant2022 also has lowest TNR.

143 We further assessed the performance in detecting direct causal variants. In scenarios where no
144 direct causal variants were present, the TGVIS identified fewer direct causal variants, with an
145 average number of 0.92, compared to 2.39 for TGFm and 2.38 for cTWAS (**Figure 2D**). When
146 there were two direct causal variants present, TGVIS identified an average of 2.86 direct causal
147 variants, compared to 3.58 for both cTWAS and TGFm. The averaged correlations between the
148 estimated and true direct causal effects across simulations were high for all three methods
149 (**Figure 2E**). However, predicting infinitesimal effects remains challenging, as evidenced by an
150 average correlation of 0.663 between the predicted and true infinitesimal effects in TGVIS
151 (**Figure 2F**).

152 **Searching potentially causal gene-tissue pairs and variants for 45 cardiometabolic traits**

153 We systematically analyzed 45 cardiometabolic traits and eQTL/sQTL summary statistics (**Table**
154 **S1-S2**) to identify potential causal gene-tissue pairs and direct causal variants. For the TGVIS,
155 we considered whether a gene-tissue pair or direct causal variant was causal if (1) it was within a
156 95% credible set and (2) had a CS-Pratt > 0.15. The criteria of CS-Pratt >0.15 was established

157 based on the empirical evidence (**Methods**). For TGFM, we followed the authors'
158 recommendation of considering individual PIP > 0.5 as indicative of causality. Additionally, we
159 did not compare with cTWAS because it analyzes tissues separately¹³.

160 TGVIS and TGFM identified a median of 119.5 and 227.5 causal gene-tissue pairs, and 42 and
161 183 causal variants per trait (**Figure 3A, Table S3-S6**), respectively. Additionally, TGVIS
162 detected a median of 0.313 causal gene-tissue pairs and 0.115 direct causal variants per locus,
163 while TGFM identified a median of 0.469 causal gene-tissue pairs and 0.466 direct causal
164 variants per locus (**Figure 3C**). Overall, TGVIS reduced the number of causal gene-tissue pairs
165 by a median of 55.7% and the number of direct causal variants by 24.5% per trait compared to
166 TGFM, representing the improved resolution of TGVIS over TGFM.

167 We expected that general causal gene-tissue pairs detected by TGVIS and TGFM would likely
168 be included among those identified by univariable TWAS methods such as S-PrediXcan⁴¹.

169 Surprisingly, among the causal pairs identified by TGVIS, a median of 34.3% were undetected
170 by S-PrediXcan, and this proportion was 60.1% for TGFM (**Figure 3A, Table S22**). For
171 example, TGVIS identified *SCN2A*-Nerve_Tibial as a novel causal gene-tissue pair for 17 traits
172 (**Figure S37**) but missed by TWAS. Both TGVIS and TWAS identified *SCN2A*-Nerve_Tibial for
173 type 2 diabetes. Our findings suggest *SCN2A* may regulate a wide range of metabolic traits.

174 These results indicate that TGVIS not only fine-maps causal genes detected but also uncovers
175 novel genes by modeling multiple tissue-gene pairs simultaneously.

176 We investigated how many traits can be influenced by a causal gene-tissue pair, reflecting the
177 pleiotropic effect at the gene-tissue level. Among the causal gene-tissue pairs falling in credible
178 sets of size less than 2, 22.4% identified by TGVIS and 16.7% by TGFM exhibit pleiotropic

179 effect (**Figure 3D, Table S7-S8**), indicating that many of these causal genes contribute to shared
180 biological mechanisms across multiple traits.

181 We further examined whether the direct causal variants and xQTLs mediated by causal gene-
182 tissue pairs differ in functionality using functional annotations³⁵ (**Methods**). Significant
183 differences were observed between these two types of variants identified by either TGVIS or
184 TGFM across multiple annotations (**Table S11**). As shown in **Figure 3E** and **3F**, the direct
185 causal variants generally have higher FathmmXF and h3k9me3 scores than the xQTLs mediated
186 by causal gene-tissue pairs (Wilcoxon signed-rank test, $P < 2.2E-16$), suggesting distinct
187 biological mechanisms for many of these variants.

188 We observed that multiple eGenes and sGenes often shared the same set of variants as their
189 xQTL, highlighting the importance of making inferences based on credible sets rather than
190 individual variables. Most credible sets consisted of 2 to 4 gene-tissue pairs (60.5%), although
191 some credible sets included more than 10 (11.5%) for TGVIS (**Figure 4A, Table S12**). In
192 comparison, TGFM resulted predominantly featured single gene-tissue pairs (56.0%) and 2 to 4
193 pairs (41.7%) per credible set (**Figure S18, Table S12**). On the other hand, most of the credible
194 sets only had one xQTL (66.6%), followed by two xQTLs (12.6%) for TGVIS (**Figure 4B,**
195 **Table S13**). As for TGFM, these percentages were 26.9% for one xQTL and 24.4% for two
196 xQTLs (**Figure S20, Table S13**). These differences arise because TGFM resampled all xQTLs
197 in the 95% credible sets, typically incorporating more variants, whereas TGVIS applied a stricter
198 criterion for selecting xQTLs (**Methods, Figure S18-S19**).

199 We investigated the proportions of identified causal eGenes and sGenes for the 45
200 cardiometabolic traits (**Methods**). TGVIS showed eGenes and sGenes proportions of 58.1% and
201 41.9%, respectively, while TGFM resulted in 63.5% for eGenes and 36.5% for sGenes (**Figure**

202 **4C, Figure S21**). These results align with the proportions observed in the GTEx Consortium
203 (63% *cis*-eQTL vs. 37% *cis*-sQTL)¹², with TGFm's proportions being slightly closer. A potential
204 explanation is that TGVIS' eGenes and sGenes were more likely enriched for causal genes
205 specific to cardiometabolic traits, leading to a slight difference, though this difference is not
206 substantial.

207 We calculated the Pratt index of gene-tissue pairs, direct causal variants, and infinitesimal effects
208 based on its additive property (**Figure 4D, Table S15**), which helps measure the contributions of
209 these three potentially correlated components (**Methods**). For TGVIS, the median of the Pratt
210 index was 0.161, 0.059, and 0.182 for gene-tissue pairs, direct causal variants, and infinitesimal
211 effects, respectively, with a median sum of the Pratt index of 0.403. In comparison, for TGFm,
212 the median of the Pratt index was 0.145 for gene-tissue pairs and 0.114 for direct causal variants,
213 with a median sum of the Pratt index of 0.262. These results support the existence of widespread
214 infinitesimal effects.

215 **Major relevant tissue map of cardiometabolic traits**

216 We searched the major relevant tissues by counting their numbers to the causal gene-tissue pairs
217 in credible sets identified by TGVIS and TGFm (**Methods**). We ranked the top relevant tissues
218 according to their contributions and clustered similar traits and tissues based on the similarity of
219 the identified causal gene-tissue pairs (**Figure 5A-5B, Figure S22-S23**). Overall, we observed
220 similar major relevant tissues and clustering patterns using both methods, although there were
221 some notable differences. TGVIS tended to cluster similar traits more closely together than
222 TGFm. For instance, TGVIS grouped all blood pressure traits into close clusters, placing them
223 near coronary artery disease (CAD), whereas TGFm positioned systolic and diastolic blood
224 pressures (SBP and DBP) farther from pulse pressure (PP) and CAD. Similarly, serum lipid traits

225 were clustered together by TGVIS, but not by TGFM. On the other hand, arterial tissues
226 consistently emerged as the major tissue for blood pressure traits and CAD, while heart tissues
227 were the major tissue for the QRS complex, atrial fibrillation, QT interval, and JT interval.
228 Fibroblasts were highlighted as an important tissue for many traits, aligning with recent findings
229 about their role in tissue integrity and chronic inflammation, alongside other tissues such as
230 adipose tissue and liver⁴².

231 We considered several lipids traits, including LDL-C, HDL-C, TC, triglyceride,
232 apolipoprotein A1 (APOA1), and apolipoprotein B (APOB), as examples to illustrate the
233 proportional counts of each tissue identified in the credible sets. For HDL-C and triglycerides,
234 the most relevant tissue was subcutaneous adipose (**Figure 5C**). In contrast, liver tissue was
235 consistently the most relevant tissue for LDL-C, APOB, and TC, despite the small sample size
236 for the liver tissue gene expression data¹². For APOA1, the two most relevant tissues were the
237 liver and subcutaneous adipose tissue. **Figures S24-S32** display the plots of major tissues for the
238 rest of the traits. Overall, the TGVIS and TGFM produced in general consistent results.

239 **Evaluation of the identified gene-tissue pairs**

240 To evaluate the accuracy of the prioritization of causal gene-tissue pairs, we first compared the
241 colocalization evidence of the causal credible sets identified by TGVIS and TGFM through
242 Coloc-SuSiE⁵. Since a credible set could include multiple tissue-gene pairs, we defined a
243 colocalization of a credible set in two criteria: (1) the credible set contained at least one gene-
244 tissue pair that is colocalized with the trait; (2) more than 50% of the gene-tissue pairs in the
245 credible set were colocalized with the trait (**Method**). TGVIS had much higher proportions of
246 colocalized credible sets (the median proportions across traits were 93.1% and 77.8% for the two
247 criteria, respectively) than TGFM (the median proportions across traits were both 40.9% for two

248 criteria) (**Figure 6A and Table S16-S17**), suggesting a substantial number of causal tissue-gene
249 pairs identified by TGFm do not have colocalization evidence.

250 We next followed the previous analysis strategy¹³ to assess the causal genes for LDL-C
251 identified by TGVIS and TGFm. Precision was evaluated using the 69 known lipid-related genes
252 as the silver standard positive gene set, and nearby genes within a 1MB-radius region as the
253 negative set, as studied by Zhao et al.¹³. We disregarded the tissue part of the identified causal
254 gene-tissue pairs and then calculated how many causal genes were within the lists of silver and
255 nearby genes. TGVIS demonstrated a precision of 60.0% (9 out of 15), outperforming TGFm,
256 which had a precision of 37.5% (10 out of 28) (**Table S18 and Figure S33**).

257 It is reasonable to assume that causal genes are more likely to be druggable targets. We utilized
258 the published list of 6,690 FDA/EMA-approved non-cancer drugs (**Table S1** provided by
259 Trajanoska et al.³⁶) to calculate the enrichment of the identified causal genes in the drug list
260 (**Figure 6B and Table S19-S20**). Although the number of causal genes identified by TGVIS in
261 the drug-targeted gene list was only 74.3% of that identified by TGFm, the enrichment identified
262 by TGVIS was 1.43 times more than that by TGFm ($P = 1.56E-3$).

263 We hypothesized that causal genes detected through eQTLs/sQTLs may be more likely to
264 demonstrate association evidence in protein data. To test this, we conducted univariable MR
265 analysis of protein abundances (pGenes) in blood tissue for genes identified by TGVIS and
266 TGFm, using both *trans*- and *cis*-pQTLs as instrument variables (**Figure 6C and Table S21**). On
267 average, 18.1% of pGenes identified by TGVIS showed significant causal evidence, compared to
268 13.7% of pGenes for TGFm ($P = 3.1E-3$). However, this proportion is lower than the estimated
269 true positive association rate of 27.8% between predicted *cis*-regulated gene expression and
270 plasma protein abundances⁴³. The discrepancy may arise from the fact that pGenes are

271 influenced by widespread *trans*-pQTLs¹⁰, whereas predicted gene expression is predominantly
272 contributed by *cis*-eQTLs, and its *trans*-regulated effects are much more difficult to detect. This
273 result suggests that eGenes/sGenes and pGenes may represent distinct biological processes
274 related to complex traits⁴³.

275 **Fine-mapping of causal gene-tissue pairs and variants in GWAS loci**

276 We exemplified four loci associated with LDL-C, CAD, and BMI. The first locus contains the
277 *PCSK9* gene for LDL-C (**Figure 7A**). TGVIS identified three 95% credible sets, including
278 *PCSK9*-Whole_Blood and two direct causal variants *rs11591147* and *rs11206517* (**Figure 7B**).
279 After applying the threshold of CS-Pratt > 0.15, *PCSK9*-Whole_Blood (CS-Pratt = 0.17) and
280 *rs11591147* (CS-Pratt = 0.492) remained. In contrast, TGFm identified nine gene-tissue pairs
281 and direct causal variants with PIPs > 0.5 (**Figure 7C**), including the *MROH7*-
282 Esophageal_Mucosa which has no clear connection to the biology of LDL-C. Applying the CS-
283 Pratt threshold, *PCSK9*-Whole_Blood (CS-Pratt = 0.204) and *rs11591147* (CS-Pratt = 0.524)
284 remained, consistent with the results yielded by TGVIS. This example demonstrates how TGVIS
285 reduces false positives by modeling infinitesimal effects and applies the Pratt Index as an
286 additional criterion.

287 The second locus contains the *HMGCR* gene causal⁴⁴ to LDL-C (**Figure 7D**). TGVIS identified
288 five 95% credible sets (**Figure 7E**). The first credible set (the darkest green) includes 9 gene-
289 tissue pairs, such as *HMGCR*-Muscle_Skeletal and five of its sGenes in esophagus mucosa,
290 nerve tibial, fibroblasts, and adipose visceral, all sharing the same xQTL *rs2112653*. When we
291 applied the threshold of individual PIP > 0.5, none of the pairs in this credible set were selected
292 although they were all in a 95% credible set. However, this set had the highest CS-Pratt of 0.322
293 among the five 95% credible sets. Conversely, TGFm identified *POLK*-Lung (CS-Pratt = 0.684)

294 but missed the crucial *HMGCR* gene (**Figure 7F**). This is likely a false discovery, as *HMGCR*
295 inhibitor is a key component of statins, which works by inhibiting HMG-CoA reductase and thus
296 reduces LDL-C in the blood⁴⁴.

297 In the third example, we focused on the *PHACTRI* locus related to CAD (**Figure 7G**). Both
298 TGVIS and TGM identified a major credible set at this locus, including *PHACTRI*-
299 Artery_Coronary and *PHACTRI*-Artery_Aorta, with CS-Pratt values of 0.632 and 0.612,
300 respectively (**Figure 7H-7I**). In TGVIS, the individual PIPs of them were both 0.5 and the
301 cumulative PIP for this credible set was 1. In contrast, TGM resampled both the eQTL effect
302 estimates and the PIPs (**Methods**), resulting in a higher individual PIP and Pratt index for
303 *PHACTRI*-Artery_Aorta (PIP = 0.597, Pratt = 0.472) than *PHACTRI*-Artery_Coronary (PIP =
304 0.222, Pratt = 0.053). However, as noted by Strober et al.¹⁴, this resampling process tends to
305 favor gene-tissue pairs with larger sample sizes, which may explain the exclusion of *PHACTRI*-
306 Artery_Coronary. TGVIS adheres to the original interpretation of SuSiE that the variables within
307 a credible set cannot be distinguished from the available data.

308 The final exemplary locus is the *FTO* locus associated with BMI (**Figure 7J**). TGVIS identified
309 only two direct causal variants *rs7206790* and *rs3751813* and did not find any gene-tissue pairs
310 at this locus (**Figure 6K**). In contrast, TGM identified four gene-tissue pairs:
311 *FTO*_Kidney_Glomerulus, *FTO*_Thyroid, *FTO*_Artery Tibial, and *IRX3*-
312 Adipose_Subcutaneous, and five direct causal variants (**Figure 7L**). However, the associations
313 between obesity and the expression of the *FTO* gene in the kidney glomerulus, thyroid, and tibial
314 artery are not well-established in the literature. After applying the Pratt index threshold, only two
315 direct causal variants *rs7206790* and *rs3751813*, remained, consistent with the result from the
316 TGVIS. When we reduced the locus radius from 1MB to 500KB and re-ran the analysis, both

317 TGVIS and TGM identified the sGene of *FTO*-Pancreas as causal, with CS-Pratt values of
318 0.345 and 0.407, respectively (**Figure S35**). The sQTLs of this *FTO* sGene are *rs7206790* and
319 *rs11642841*, which have been reported by Xu et al.⁴⁵. This example suggests that when applying
320 multivariable TWAS methods, the size of a *cis*-region can be sensitive and need to be calibrated.

321 **Discussion**

322 In this report, we developed TGVIS to identify causal gene-tissue pairs and direct causal variants
323 in loci identified through GWAS by integrating xQTL summary statistics. Compared to
324 cTWAS¹³ and TGM¹⁴, TGVIS not only analyzes multiple tissue-specific xQTL summary data
325 simultaneously to pinpoint causal gene-tissue pairs and direct causal variants, but also models the
326 widespread presence of infinitesimal effects underlying polygenic traits to reduce false discovery
327 rates in detecting causal molecular phenotypes¹⁷. In addition, TGVIS quantifies the importance
328 of a causal variable by the Pratt index, which has been well established in statistics and has
329 recently been applied to estimate the gene-by-environment contribution²³. Through simulations,
330 we demonstrated that under the presence of infinitesimal effects, TGVIS has lower MSE and
331 higher TPR and TNR compared to both cTWAS and TGM (**Figure 2**). In real data analysis,
332 TGVIS outperformed TGM in the following four aspects: (1) identifying more interpretable
333 major trait-relevant tissues (**Figure 5**); (2) resulting in a higher proportion of colocalized causal
334 credible sets (93.1% vs 40.9%, **Figure 6A**); (3) achieving notably higher precision in the 'silver
335 standard' sets of lipids (60.0% vs 37.5%, **Table S15**); and (4) demonstrating significantly greater
336 enrichment evidence based on druggable genes (1.43 times, **Figure 6B**) and causal proteins (1.31
337 times, **Figure 6C**).

338 Our analysis of 45 cardiometabolic traits provides several key insights. First, we identified a
339 median of 34.3% causal gene-tissue pairs that were missed in univariable TWAS analysis,

340 suggesting that TGVIS is able to identify novel genes besides fine-mapping the genes detected
341 by conventional TWAS (**Figure 3A**), representing a significant advance in TWAS. Second, we
342 observed that infinitesimal effects can make a substantial contribution to local genetic variation
343 of traits besides the gene-tissue pairs and direct variant (**Figure 4D**), which is consistent with
344 recent studies^{17,21,46}. Beyond underlying biological mechanisms such as the polygenicity of
345 human complex traits, the emergence of infinitesimal effects may also be attributed to non-
346 biological factors, particularly estimation errors in the LD matrix, xQTL effect sizes, and trait
347 GWAS imputation (**Method**). Both empirical observations and theoretical investigation
348 underscore the importance of including infinitesimal effects in future genetic researches and
349 methodological developments. Third, our study indicates that a significant proportion of causal
350 gene-tissue pairs (22.4%) exhibit pleiotropic effects at the gene-tissue level, suggesting shared
351 biological mechanisms across multiple traits (**Figure 3D, Table S7-S8**). Fourth, our findings
352 suggest that for most traits, only a limited number of relevant major tissues are involved (**Figure**
353 **4A**), implying that concentrating multi-omics data analyses on these relevant major tissues can
354 be more powerful and efficient, as well as it can make the findings more biologically
355 interpretable. For example, when the analysis is focused on the four major blood-pressure-
356 relevant tissues, i.e., adrenal gland, artery, heart, and kidney, it leads to the identifications of
357 more causal gene-tissue pairs, with an increased Pratt index for blood pressure traits⁴⁷. Fifth, our
358 results indicate that only 18.1% of causal genes from eQTL/sQTL analyses also show causal
359 evidence in univariable MR using pQTL summary data (**Figure 6C**), suggesting that gene
360 expressions and protein abundance represent distinct biological processes in complex traits⁴³.
361 Finally, we identified an average of 0.304 causal gene-tissue pairs per locus and failed to identify
362 any causal gene-tissue pairs in many GWAS loci (**Figure 3A**), which is consistent with the

363 recent study showing that the GWAS and eQTL studies are systematically biased toward
364 different types of variants⁴. Interestingly, the eQTLs/sQTLs of causal gene-tissue pairs and direct
365 causal variants have substantially different functional annotations (**Figure 3E-3F, Table S11**),
366 warranting further investigation.

367 Our study has some limitations. Due to the data and computational constraints, we only analyzed
368 genes using *cis*-eQTL/sQTL summary statistics, limiting our ability to distinguish between genes
369 that share *cis*-eQTLs/sQTLs, which may lead to false discoveries. This issue could potentially be
370 addressed by incorporating *trans*-eQTLs/sQTLs, although this would require much larger sample
371 sizes. In addition, we observed that a credible set often contains 2-4 gene-tissue pairs (**Figure**
372 **3C**), likely due to the small sample size in the GTEx data, which results in only 1 or 2
373 eQTL/sQTL for most gene-tissue pairs (**Figure 4B**). In other words, while TGVIS was able to
374 narrow down to a range of causal gene-tissue pairs, it could not always pinpoint the exact causal
375 pair(s) in some loci. Incorporating external information, such as colocalization evidence with
376 TGVIS, may aid in distinguishing these pairs⁴⁸. Furthermore, our eQTL/sQTL analysis relies on
377 bulk tissue expression data, which may limit our ability to identify cell-type-specific causal
378 genes⁴⁹. For example, recent studies increasingly suggest that *FTO* may not be the causal gene
379 for BMI; instead, experimental evidence indicates that *IRX3* and *IRX5* are the causal genes⁵⁰.
380 However, the causality of *IRX3* and *IRX5* was observed in experiments using preadipocytes,
381 rather than bulk subcutaneous adipose tissue, which may explain why TGVIS failed to identify
382 these genes (**Figure S36**). Moreover, we used the Pratt index²³ to rank the importance of
383 variables, but it has inherent statistical limitations²³. In simulations, the Pratt index slightly
384 underestimates the true contribution, although this underestimation becomes negligible as the
385 sample size increases (**Figure S1-S8**). In real data analysis, we used an empirical cutoff learned

386 by K-means (CS-Pratt = 0.15) to extract important causal variables, which gives us higher
387 precision but may have potentially hindered the discovery of causal gene-tissue pairs with small
388 to moderate causal effects. Last, as suggested in previous studies¹⁴, the inference of causality
389 based on statistical methods comes with a caveat, assuming no model misspecification and no
390 potential causal elements are missing from the model.

391 In summary, our developed TGVIS and accompany software pipeline provide a valuable tool in
392 fine-mapping and interpreting GWAS findings.

393 **Methods**

394 **Multivariable TWAS model**

395 The causal diagram shown in **Figure 1A** can be described by the following multivariable TWAS
396 model:

$$397 \quad y = \sum_{j=1}^J \sum_{t=1}^T X_{jt} \theta_{jt} + \mathbf{G}^T (\boldsymbol{\gamma} + \mathbf{v}) + \varepsilon, \quad (1)$$

398 where y is a trait, X_{jt} is the jt^{th} gene-tissue pair, $\mathbf{G} = (G_1, \dots, G_M)^T$ is an $(M \times 1)$ vector of
399 genetic variants in the *cis*-region, $\boldsymbol{\theta} = (\theta_{11}, \dots, \theta_{JT})^T$ is an $(JT \times 1)$ vector of causal effects with
400 θ_{jt} being the causal effect of the (j, t) th tissue-gene pair, $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_M)^T$ is an $(M \times 1)$ vector
401 of direct causal effects, $\mathbf{v} = (v_1, \dots, v_M)^T$ is an $(M \times 1)$ vector of infinitesimal effects, and ε is
402 the random error. Let $\boldsymbol{\beta}_{jt} = (\beta_{jt1}, \dots, \beta_{jtM})^T$ is an $(M \times 1)$ vector of the *cis*-eQTL effects of JT
403 tissue-gene pairs. Then we have

$$404 \quad X_{jt} = \boldsymbol{\beta}_{jt}^T \mathbf{G} + \epsilon_{jt}, \quad (2)$$

405 where ϵ_{jt} is the noise of the jt^{th} gene-tissue pair. The reduced form of (1) is then given by:

$$406 \quad y = \mathbf{G}^T \left(\sum_{j=1}^J \sum_{t=1}^T \boldsymbol{\beta}_{jt} \theta_{jt} + \boldsymbol{\gamma} + \mathbf{v} \right) + \varepsilon, \quad (3)$$

407 where mathematically $\epsilon = \varepsilon + \sum_{j=1}^J \sum_{t=1}^T \epsilon_{jt} \theta_{jt}$.

408 An alternative version of (1) based on summarized statistics⁵¹ is:

409
$$\hat{\mathbf{a}} \sim \mathcal{N} \left(\mathbf{R} \left(\sum_{j=1}^J \sum_{i=1}^T \boldsymbol{\beta}_{jt} \theta_{jt} + \boldsymbol{\gamma} + \mathbf{v} \right), \sigma_{\alpha}^2 \mathbf{R} \right), \quad (3)$$

410 where $\hat{\mathbf{a}} = (\hat{a}_1, \dots, \hat{a}_M)^{\top}$ represents the GWAS effects of the outcome, \mathbf{R} is an $(M \times M)$ LD
411 matrix of the M variants, and σ_{α}^2 is the variance of this model. The eQTL effect vector $\boldsymbol{\beta}_{jt}$
412 follows the model based on summarized statistics below:

413
$$\hat{\mathbf{b}}_{jt} \sim \mathcal{N} \left(\mathbf{R} \boldsymbol{\beta}_{jt}, \sigma_{\beta_{jt}}^2 \mathbf{R} \right), \quad (5)$$

414 where $\hat{\mathbf{b}}_{jt} = (\hat{b}_{jt1}, \dots, \hat{b}_{jtM})^{\top}$ represents the marginal *cis*-eQTL effect estimates for the jt^{th}
415 tissue-gene pair, and $\sigma_{\beta_{jt}}^2$ denotes the variance of this model.

416 To resolve this curse of dimensionality, we utilized the three sparsity conditions that are
417 commonly assumed in current fine-mapping methods^{16,19}: (SP1) one or small number of variants
418 causally contribute to tissue or cell-type specific gene expression¹²; (SP2) one or small number
419 of gene-tissue pairs causally contribute to the trait^{13,14}; (SP3) one or small number of direct
420 causal variants exist with relatively large effect sizes^{13,14}. In terms of statistical model: SP1
421 corresponds to $\boldsymbol{\beta}_{jt}$ being sparse for all j and t ; SP2 corresponds to $\boldsymbol{\theta}$ being sparse; SP3
422 corresponds to $\boldsymbol{\gamma}$ being sparse. In addition, we incorporated that variants can have infinitesimal
423 effects: \mathbf{v} is normally distributed with a mean 0 and a small, unknown variance¹⁷. To our best
424 knowledge, infinitesimal effects have not been modeled in current multivariable TWAS
425 methods.

426 **Estimation of *cis*-regulatory effect**

427 TGVIS first applies SuSiE¹⁸ to estimate the non-zero eQTL effect for each gene-tissue pair,
428 based on the fine-mapping model (Equation 5). Specifically, we set $L = 3$ for each pair and
429 determined the non-zero *cis*-regulatory effects based on two criteria: (1) if they are within any
430 95% credible set and their PIPs exceeds 0.25, and (2) if their individual PIPs are greater than 0.5.
431 The rationale behind this approach is that SuSiE's 95% credible set can sometimes include too
432 many weakly correlated variants (even after removing highly correlated ones using LD
433 clumping), leading to low PIPs for each variant. Therefore, we used a moderate threshold to filter
434 out credible sets with too many variants. Additionally, due to the low power of detection, the
435 maximum PIP of credible sets might fall below 0.95, so we retained variants with individual
436 PIPs greater than 0.5. Since a locus often contains over 10,000 gene-tissue pairs (mostly sGenes),
437 dynamically selecting using BIC would be computationally burdensome. Additionally, with
438 GTEx sample sizes under 200, only 1-2 gene-tissue pairs can be identified for most gene-tissue
439 pairs. Therefore, we choose to fix $L = 3$.

440 **Joint modelling of causal tissue-gene pairs, direct causal variants, and infinitesimal effects**
441 **using profile likelihood**

442 TGVIS estimates θ , γ , and \mathbf{u} using a profile likelihood approach. Given the estimate $\mathbf{u}^{(s)}$ from
443 the s th iteration, we considered the following fine-mapping model:

444
$$\hat{\mathbf{a}} - \mathbf{R}\mathbf{u}^{(s)} \sim \mathcal{N}(\mathbf{R}\boldsymbol{\gamma} + \mathbf{R}\hat{\mathbf{B}}\boldsymbol{\theta}, \sigma_{\alpha}^2 \mathbf{R}), \quad (6)$$

445 where $\hat{\mathbf{B}} = (\hat{\boldsymbol{\beta}}_{11}, \dots, \hat{\boldsymbol{\beta}}_{jt}, \dots, \hat{\boldsymbol{\beta}}_{JT})$ is an $M \times JT$ matrix consisting of estimated *cis*-regulatory
446 effects. To update $\boldsymbol{\gamma}$ and $\boldsymbol{\theta}$ simultaneously, we applied the same scheme as cTWAS and TGFM,
447 using the function `susie_rss(·)`. The input z-score vector is computed as:

448
$$z = \left(\frac{\hat{\boldsymbol{\beta}}_{11}^T \hat{\mathbf{a}}}{\sqrt{\hat{\boldsymbol{\beta}}_{11}^T \mathbf{R} \hat{\boldsymbol{\beta}}_{11}}}, \dots, \frac{\hat{\boldsymbol{\beta}}_{JT}^T \hat{\mathbf{a}}}{\sqrt{\hat{\boldsymbol{\beta}}_{JT}^T \mathbf{R} \hat{\boldsymbol{\beta}}_{JT}}}, \hat{a}_1, \dots, \hat{a}_M \right)^T, \quad (7)$$

449 and the other elements of input correlation matrix are computed as:

450
$$\text{cor}(\mathbf{G}_i \hat{\boldsymbol{\beta}}_{jt}, \mathbf{G}_i \hat{\boldsymbol{\beta}}_{j't'}) = \frac{\hat{\boldsymbol{\beta}}_{jt}^T \mathbf{R} \hat{\boldsymbol{\beta}}_{j't'}}{\sqrt{\hat{\boldsymbol{\beta}}_{jt}^T \mathbf{R} \hat{\boldsymbol{\beta}}_{jt}} \sqrt{\hat{\boldsymbol{\beta}}_{j't'}^T \mathbf{R} \hat{\boldsymbol{\beta}}_{j't'}}},$$

451
$$\text{cor}(\mathbf{G}_i \hat{\boldsymbol{\beta}}_{jt}, \mathbf{G}_i) = \frac{\mathbf{R} \hat{\boldsymbol{\beta}}_{jt}}{\sqrt{\hat{\boldsymbol{\beta}}_{jt}^T \mathbf{R} \hat{\boldsymbol{\beta}}_{jt}}}, \quad \text{cor}(\mathbf{G}_i) = \mathbf{R}, \quad (8)$$

452 The outputs are denoted as $\boldsymbol{\gamma}^{(s+1)}$ and $\boldsymbol{\theta}^{(s+1)}$.

453 Next, we consider the following model:

454
$$\hat{\mathbf{a}} - \boldsymbol{\eta}^{(s+1)} | \mathbf{v} \sim \mathcal{N}(\mathbf{R}\mathbf{v}, \sigma_\alpha^2 \mathbf{R}), \quad v \sim \mathcal{N}(0, \sigma_v^2 \mathbf{I}), \quad (9)$$

455 where $\boldsymbol{\eta}^{(s+1)} = \mathbf{R}(\hat{\mathbf{B}}\boldsymbol{\theta}^{(s+1)} + \boldsymbol{\gamma}^{(s+1)})$. The penalized quasi-likelihood (PQL) of \mathbf{v} is

456
$$\mathbf{v}^{(s+1)} = \underset{\mathbf{v}}{\text{argmin}} \left\{ (\hat{\mathbf{a}} - \boldsymbol{\eta}^{(s+1)} - \mathbf{R}\mathbf{v})^T \mathbf{R}^{-1} (\hat{\mathbf{a}} - \boldsymbol{\eta}^{(s+1)} - \mathbf{R}\mathbf{v}) + \frac{\sigma_\alpha^{(s)2}}{\sigma_v^{(s)2}} \|\mathbf{v}\|_2^2 \right\}, \quad (10)$$

457 which results in

458
$$\mathbf{v}^{(s+1)} = \left(\mathbf{R} + \frac{\sigma_\alpha^{(s)2}}{\sigma_v^{(s)2}} \mathbf{I} \right)^{-1} (\hat{\mathbf{a}} - \boldsymbol{\eta}^{(s+1)}), \quad (11)$$

459 where $\sigma_\alpha^{(s)2}$ is the current variance estimate. The variance $\sigma_v^{(s)2}$ is updated by REML:

460
$$\sigma_v^{(s+1)2} = \operatorname{argmin}_{\sigma_v^2} \left\{ \frac{\| \mathbf{v}^{(s+1)} \|_2^2}{\sigma_v^2} + M \log(\sigma_v^2) + \log \det \left(\frac{1}{\sigma_\alpha^{(s)2}} \mathbf{R} + \frac{1}{\sigma_v^{(s)2}} \mathbf{I} \right) \right\}, \quad (12)$$

461 which simplifies to

462
$$\sigma_v^{(s+1)2} = \frac{1}{M} \| \mathbf{v}^{(s+1)} \|_2^2 + \frac{1}{M} \operatorname{tr} \left(\left(\mathbf{R} + \frac{\sigma_\alpha^{(s)2}}{\sigma_v^{(s)2}} \mathbf{I} \right)^{-1} \right), \quad (13)$$

463 where M is the number of variants. We replace σ_v^2 in the last term by its current estimate $\sigma_v^{(s)2}$ to

464 obtain a close-form expression. Note that in Equation (13), $\frac{\sigma_\alpha^{(s)2}}{\sigma_v^{(s)2}}$ is usually replaced by $\frac{1}{\sigma_v^{(s)2}}$ to

465 avoid non-identifiability issues²².

466 When the profile likelihood converges, TGVIS estimates σ_α^2 as follows:

467
$$\sigma_\alpha^{(s+1)2}$$

468
$$= \frac{1}{M} \left(\hat{\mathbf{a}} - \mathbf{R}(\hat{\mathbf{B}}\boldsymbol{\theta}^{(s+1)} + \boldsymbol{\gamma}^{(s+1)} + \mathbf{v}^{(s+1)}) \right)^\top \mathbf{R}^{-1} \left(\hat{\mathbf{a}} - \mathbf{R}(\hat{\mathbf{B}}\boldsymbol{\theta}^{(s+1)} + \boldsymbol{\gamma}^{(s+1)} + \mathbf{v}^{(s+1)}) \right). \quad (14)$$

469 Bayesian Information Criterion for Summary Data

470 Based on Equation (3), we define the BIC for summary data:

471
$$\text{BIC} = \log(\sigma_\alpha^2) + \frac{\log M}{M} \text{df}, \quad (15)$$

472 where M is the number of IVs and df is the degree of freedom of the model⁵². In practice, σ_α^2 is

473 replaced by its empirical estimate $\hat{\sigma}_\alpha^2$, and df is the sum of non-zero causal effect estimates and

474 non-zero direct causal variant estimates. Our default setting assumes L can be 2,3,4,5,6,7, or 8

475 and uses BIC to select the optimal L among them. We found that when considering the

476 infinitesimal effect, it tends to capture variants with very small effects that SuSiE does not
477 identify, making it rare for L to exceed 8 in practice.

478 **Pratt index**

479 We use the Pratt index to assess the contribution of a gene-tissue pair. For a general linear
480 model: $y_i = \sum_{j=1}^p X_j \beta_j + \epsilon_i$, the Pratt index of x_{ij} is defined as $V_j = \beta_j \times b_j$, where $b_j =$
481 $\text{cov}(y, X_j)$. This definition assumes standardization where $E(y) = E(X_j) = 0$ and $\text{var}(y) =$
482 $\text{var}(X_j) = 1$, $1 \leq j \leq p$. The Pratt index measures the contribution of a variable in a linear
483 model because $R^2 = \sum_{j=1}^p V_j$ where $R^2 = \text{var}(\sum_{j=1}^p X_j \beta_j) / \text{var}(y)$. In practice, the Pratt index
484 can be estimated by $\hat{V}_j = \hat{\beta}_j \times \hat{b}_j$, where \hat{b}_j is the sample correlation between X_j and y .

485 The proportion of variance explained (PVE) is defined as $PVE_j = \beta_j^2$, assuming that all variables
486 are standardized. The Pratt index has two key advantages over PVE: (1) Pratt indices are additive
487 across variables, and (2) the sum of Pratt indices is the total trait variance explained by
488 covariates. In contrast, PVE lacks these advantages.

489 **Pratt Index in TGVIS**

490 We show how to yield the Pratt index V_{jt} in practice. We first estimate the marginal correlation:

491
$$\tilde{\delta}_{jt} = \widehat{\text{cor}}(\hat{\beta}_{jt}, \hat{\mathbf{a}}) = \frac{\hat{\beta}_{jt}^\top \hat{\mathbf{a}}}{\sqrt{\hat{\beta}_{jt}^\top \mathbf{R} \hat{\beta}_{jt} \sqrt{\hat{\mathbf{a}}^\top \mathbf{R}^{-1} \hat{\mathbf{a}}}}}. \quad (16)$$

492 As for the causal effect estimate $\hat{\theta}_{jt}$, we apply the transformation

493
$$\tilde{\theta}_{jt} = \hat{\theta}_{jt} \frac{\sqrt{\hat{\beta}_{jt}^\top \mathbf{R} \hat{\beta}_{jt}}}{\sqrt{\hat{\mathbf{a}}^\top \mathbf{R}^{-1} \hat{\mathbf{a}}}}, \quad (17)$$

494 since the Pratt index requires the covariates and trait are all standardized. Thus, the Pratt index of
495 the (j, t) th gene-tissue pair is

496
$$\hat{V}_{jt} = \tilde{\theta}_{jt} \times \tilde{\delta}_{jt} = \hat{\theta}_{jt} \frac{\hat{\boldsymbol{\beta}}_{jt}^T \hat{\mathbf{a}}}{\hat{\mathbf{a}}^T \mathbf{R}^{-1} \hat{\mathbf{a}}}. \quad (18)$$

497 Since Pratt indices are additive, the Pratt index of a credible set is simply calculated as

498
$$\hat{V}_{CS_l} = \sum_{j \in CS_l} \hat{V}_j. \quad (19)$$

499 Note that the Pratt index is only comparable within the same locus, as it represents the ratio of
500 the variance explained by the variable to the total variance of the trait.

501 It is worth comparing the gene-tissue pair, direct causal variant, and infinitesimal effect
502 contributions at a locus. To simplify the estimation, we consider the linear predictors of all gene-
503 tissue pairs and pleiotropy:

504
$$\tilde{\boldsymbol{\eta}}_{\theta} = \mathbf{R}^{\frac{1}{2}} \hat{\mathbf{B}} \hat{\boldsymbol{\theta}}, \quad \tilde{\boldsymbol{\eta}}_{\gamma} = \mathbf{R}^{\frac{1}{2}} \hat{\boldsymbol{\gamma}}, \quad \tilde{\boldsymbol{\eta}}_v = \mathbf{R}^{\frac{1}{2}} \hat{\mathbf{v}}. \quad (20)$$

505 and $\tilde{\mathbf{a}} = \mathbf{R}^{-\frac{1}{2}} \hat{\mathbf{a}}$, where $\mathbf{R}^{-\frac{1}{2}}$ is specified to remove the correlations of $\hat{\mathbf{B}}$ and $\hat{\mathbf{a}}$. Then, the Pratt
506 indices for the gene-tissue pairs, direct causal variants, and infinitesimal effects are

507
$$\hat{V}_{\theta} = \frac{\tilde{\boldsymbol{\eta}}_{\theta}^T \tilde{\mathbf{a}}}{\|\tilde{\mathbf{a}}\|_2^2}, \quad \hat{V}_{\gamma} = \frac{\tilde{\boldsymbol{\eta}}_{\gamma}^T \tilde{\mathbf{a}}}{\|\tilde{\mathbf{a}}\|_2^2}, \quad \hat{V}_v = \frac{\tilde{\boldsymbol{\eta}}_v^T \tilde{\mathbf{a}}}{\|\tilde{\mathbf{a}}\|_2^2}. \quad (21)$$

508 **Threshold of Pratt index**

509 We used the empirical data to determine the threshold for Pratt index to enhance the precision of
510 causal selection. Specifically, we employed K-means clustering with clusters to group the CS-
511 Pratt indices of all gene-tissue pairs and direct variants identified by TGVIS within the 95%

512 credible sets. We hypothesize that one cluster contains credible sets with smaller CS-Pratt
513 values, which are more likely to include falsely causal variables. Interestingly, regardless of
514 whether we focus on gene-tissue pairs, direct causal variants, or both, the minimum value in the
515 cluster with the larger centroid consistently remains at 0.15 (**Figure S34**). Consequently, we set
516 the threshold at CS-Pratt = 0.15 to prioritize the gene-tissue pairs and direct causal variants
517 identified by TGVIS, considering variables with CS-Pratt > 0.15 to have a higher likelihood of
518 being true causal.

519 **Score test of variance of infinitesimal effects**

520 In implementation, dynamically determining whether to consider the infinitesimal effect is a
521 clever empirical measure. Therefore, we apply the score test of the variance of the random effect
522 in the linear mixed model to test whether the variance of the infinitesimal effect is zero.
523 Specifically, we consider the following hypothesis testing problem:

$$524 \quad H_0: \sigma_v^2 = 0, \quad v. s. \quad H_1: \sigma_v^2 > 0. \quad (22)$$

525 The testing statistics of this hypothesis test is constructed according to Zhang and Lin⁵³. Let $\mathbf{A} =$
526 $(\mathbf{R}\widehat{\mathbf{B}}_{\mathcal{M}_\theta}, \mathbf{R}_{\mathcal{M}_\gamma})$ and $\boldsymbol{\vartheta} = (\boldsymbol{\theta}_{\mathcal{M}_\theta}^\top, \boldsymbol{\gamma}_{\mathcal{M}_\gamma}^\top)^\top$. When $\sigma_v^2 = 0$ and $\sigma_v^2 > 0$, the covariance matrix of $\widehat{\boldsymbol{\alpha}} -$
527 $\mathbf{A}\boldsymbol{\vartheta}$ are

$$528 \quad \text{cov}(\widehat{\boldsymbol{\alpha}} - \mathbf{A}\boldsymbol{\vartheta}) = \sigma_\alpha^2 \mathbf{R}, \quad \text{cov}(\widehat{\boldsymbol{\alpha}} - \mathbf{A}\boldsymbol{\vartheta}) = \sigma_\alpha^2 \mathbf{R} + \sigma_v^2 \mathbf{R}^2, \quad (23)$$

529 respectively. Similar to estimating σ_v , we replace σ_α^2 by 1 to avoid non-identifiability. The score
530 described in Zhang and Lin⁵³ defined the following three statistics:

$$531 \quad u = \frac{1}{2} \|\widehat{\boldsymbol{\alpha}} - \mathbf{A}\boldsymbol{\vartheta}\|_2^2, \quad e = \frac{1}{2} \text{tr}(\mathbf{P}\mathbf{R}^2), \quad h = \frac{1}{2} \text{tr}(\mathbf{P}\mathbf{R}^2\mathbf{P}\mathbf{R}^2), \quad (24)$$

532 where $\mathbf{P} = \mathbf{R}^{-1} - \mathbf{R}^{-1}\mathbf{A}(\mathbf{A}^T\mathbf{R}^{-1}\mathbf{A})^{-1}\mathbf{A}^T\mathbf{R}^{-1}$. Under the null hypothesis, $u \sim \kappa\chi_v^2$ where $\kappa =$
533 $h/(2e)$ and $v = 2e^2/h$. If the null hypothesis is accepted, we enforce $\mathbf{v}^{(s+1)} = 0$.

534 **Potential reasons leading to Infinitesimal effects**

535 Here we listed four possible reasons that can lead an infinitesimal effect. First, it has been
536 gradually understood that even within the same ethnic group, such as the European population,
537 different subgroups may have different genetic architectures, leading to different LD structures.
538 Therefore, it is natural to suspect that the LD structures of populations in the GTEx consortium
539 and those in traits GWAS differ, which results in

$$540 \quad E(\hat{\boldsymbol{\alpha}}) = \mathbf{R}_{\text{Meta}}(\mathbf{B}\boldsymbol{\theta} + \boldsymbol{\gamma}), \quad E(\hat{\boldsymbol{\beta}}_{jt}) = \mathbf{R}_{\text{GTEx}}\boldsymbol{\beta}_{jt}. \quad (25)$$

541 When we try to estimate $\boldsymbol{\beta}_{jt}$ using \mathbf{R}_{Meta} , then $\hat{\boldsymbol{\beta}}_{jt}$ is biased to $\boldsymbol{\beta}_{jt}$, which generates infinitesimal
542 effect $\mathbf{v} = \sum_{jt} (\boldsymbol{\beta}_{jt} - \hat{\boldsymbol{\beta}}_{jt}) \theta_{jt}$. It should be noted that the small sample size in the GTEx
543 consortium can also cause biased eQTL effect estimates, resulting in the appearance of
544 infinitesimal effects. There are other possible sources raising infinitesimal effects, such as (2) the
545 estimation errors of LD matrix, (3) the imputation errors of outcome GWAS effect sizes, and (4)
546 the absence of direct causal variants of outcome, which are shown in **Supplementary Materials**.

547 **cTWAS and TGFM programs**

548 For cTWAS, since its software is designed to be user-friendly to practical projects, it involves
549 complex settings that are not ideal for simulations, such as requiring a reference panel in BED
550 format and a .db file of eQTL fine-mapping data. Therefore, we directly utilize the principles of
551 cTWAS to develop an R function that employs SuSiE for the first-stage selection of *cis*-
552 regulatory effects and the second-stage selection of causal and horizontally pleiotropic effects.

553 At the time of writing this paper, the TGFM software has not yet been released. Therefore, we
554 did not consider the first step of cTWAS to estimate two universal prior parameters using the EM
555 algorithm across all loci in the genome. Instead, we restrict cTWAS simulations to a single locus.
556 In addition, we applied the following setting for cTWAS, TGFM and TGVIS: for the prior
557 weight π in SuSiE, we applied $\pi = p^{-1}$ for gene-tissue pairs and $\pi = M^{-1}$ for variants, where
558 p represents the number of gene-tissue pairs and M the number of variants.

559 To improve computational efficiency, we applied a slightly different resampling scheme
560 compared the original TGFM. Specifically, we first resampled the eQTL effects from the
561 posterior for 25 times, calculated their mean as $\hat{\beta}_{jt}^{t_i}$ and used these means to estimate $\hat{\theta}^{t_i}$
562 and $\hat{\gamma}^{t_i}$. This procedure was repeated 100 times, recording the estimates and PIP for each
563 iteration. We then compute the mean of $t_1 \times t_2$ resampled eQTL effects, $\hat{\beta}_{jt}^{TGFM}$, and estimate the
564 empirical $\hat{\theta}^{TGFM}$ and $\hat{\gamma}^{TGFM}$. The PIPs of $\hat{\theta}^{TGFM}$ and $\hat{\gamma}^{TGFM}$ were taken as the empirical PIPs
565 given by SuSiE in each resampling iteration. Finally, we recorded the credible sets of variables
566 from the final step and calculated the PIPs and Pratt indices of credible sets by summing the
567 individual PIPs and Pratt indices of variables within each credible set.

568 **Simulation Settings**

569 We simulated 20 genes across 5 tissues, resulting in $p = 100$ gene-tissue pairs. Correlations were
570 simulated both within and between genes across tissue. The first and last gene-tissue pairs were
571 designated as causal, with effect sizes of $\theta_1 = 1$ and $\theta_{100} = -1$, respectively. The total number
572 of variants was $M = 400$, with only 1,2,3 or 4 of them being eQTLs with non-zero effects for
573 each gene-tissue pair, while the remaining variants were associated with the trait due to LD. We
574 set 4 different sample sizes for the eQTL data ($n_{eQTL} = 100, 200, 400, 800$) and a fixed trait

575 GWAS sample size $n_{trait} = 0.5M$. Infinitesimal effect were generated from a normal
576 distribution, and gene-tissue pairs, direct causal variants, and infinitesimal effects together were
577 set to explain the trait heritability. For example, when only gene-tissue pairs and infinitesimal
578 effects are present, they each explain 50% of the local heritability for the outcome. When all
579 three are present, each explains 33% of the local heritability for the outcome. The detailed
580 settings, along with corresponding R codes, are provided in Section 2 of the **Supplementary**
581 **Materials**.

582 **GWAS summary data**

583 We conducted a meta-analysis on a subset of the 45 metabolic and cardiovascular traits. The
584 publicly available data for these traits are listed in **Table S1**, while the MVP GWAS summary
585 statistics can be accessed through dbGAP under accession number phs001672.v7.p1. For the
586 pleiotropy traits of SBP and DBP, we applied the approach developed in Zhu et al.³² using the
587 most recent GWAS summary statistics of SBP and DBP. To perform the meta-analysis, we used
588 METAL⁵⁴. We performed the meta-analysis on the Z-scores, weighting by the sample sizes of
589 the meta-analysis datasets. For binary trait, we always use the effective sample size n_{eff} . We
590 used CHR:BP (in GRCH37) as the identifier.

591 **EQTL summary data**

592 We utilized bulk eQTL and sQTL summary statistics from 28 tissues provided by GTEx¹² (with
593 sample size N ranging from 34.4 (Lymphocytes) to 179.5 (Muscle Skeletal)), as well as
594 additional eQTL summary statistics from tubulointerstitial³³ (N=311), kidney glomerular³³
595 (N=240), and islet³⁴ (N=420) tissues (**Table S2**).

596 **Linkage Disequilibrium Reference Panel**

597 Our study used variants from the UKBB project conducted by Neale's lab, which initially
598 includes 13 million SNPs. We selected approximately 9.3 million SNPs with a minor allele
599 frequency greater than 0.01 for our analysis. We also identified the top 9,620 unrelated
600 individuals from approximately 500,000 individuals in the UKBB (Field ID: 22828), consisting
601 of 5,205 females and 4,475 males. Data from these 9.3 million SNPs were extracted for these
602 individuals to construct our LD reference panel.

603 **Clumping and Thresholding**

604 We restricted the studied regions to those within 1MB of the genome-wide significant loci for
605 these traits. These loci were identified using the clumping and thresholding (C+T) method in
606 PLINK⁵⁵: --clump-kb 1000, --clump-p1 5E-8, --clump-p2 5E-8, and --clump-r2 0.01.

607 We recommend using C+T to filter out variants in high LD, which prevents the inclusion of
608 numerous highly correlated or redundant variants in the analysis, which can unnecessarily
609 complicate the model and result in multiple credible sets consisting of these variants. We
610 evaluated the minimum P-value of each variant across gene-tissue pairs in eQTL/sQTL data. In
611 PLINK, we applied the C+T with the following parameters: --clump-kb 1000, --clump-p1 1E-5, -
612 -clump-p2 1E-5, and --clump-r2 0.5. Given that the true causal variant for a trait might not be
613 included in the eQTL/sQTL variants, we combined these variants from outcome GWAS
614 satisfying $P < 5E-8$ and $r^2 < 0.5$.

615 **Removing gene-tissue pairs based on significance in S-Predixcan**

616 We used the minimum P-value from S-Predixcan and a modifier accounting for direct causal
617 variants (**Supplementary Materials**) to exclude eGenes/sGene with $P > 0.05$. These weak filters
618 will eliminate redundant gene-tissue pairs, thereby reducing the model's dimensionality. Since

619 our goal is fine-mapping the causal gene-tissue pairs and identifying direct causal variants on the
620 GWAS loci with significant signals, it will not induce a winner's curse.

621 **Searching causal gene-tissue pairs missed by univariable TWAS**

622 We compared the causal gene-tissue pairs identified by TGVIS and TGFM with the significant
623 gene-tissue pairs identified by S-PrediXcan. We considered genes with $P < 0.05/20,000$ as
624 significant gene-tissue pairs in tissue specific S-PrediXcan analysis. We did not adjust for
625 number of tissues. We then searched the gene-tissue pairs identified by TGVIS or TGFM but
626 missed by S-PrediXcan.

627 **Obtaining annotation scores from FAVOR and performing differential annotation tests**

628 We combined the direct causal variants and xQTLs of causal gene-tissue pairs identified by
629 TGVIS or TGFM across all 45 traits into two separate files and uploaded them to the FAVOR
630 online platform³⁵ to obtain annotation scores for these variants. We performed Wilcoxon signed-
631 rank test with both "less" and "greater" as alternative hypotheses for determining the direction of
632 shift location and calculated corresponding P-values. We used the R package FDREstimation to
633 convert the P-values to FDR Q-values using the Benjamini–Yekutieli (BY) procedure.
634 Annotations with Q-values less than 0.05 were considered to have significantly different scores.

635 **Mapping major trait relevant tissues**

636 For TGVIS, a 95% credible set often includes multiple gene-tissue pairs. In such cases, we
637 calculated the proportion of each tissue appearing among these pairs, allowing the number of
638 tissues in a causal credible set of gene-tissue pairs to be fractional. For TGFM, we first removed
639 the gene-tissue pairs with PIPs < 0.5 , and then applied the same procedure to map the dominant
640 tissues.

641 **Enrichment of identified causal genes in lipids silver gene list and druggable gene database**

642 We applied the following strategy to map silver genes. First, we checked each credible set to see
643 if any genes are part of the silver gene list; if so, we counted 1. If no silver gene was present, we
644 then checked if any genes in the credible set were among the nearby genes; if so, we also
645 counted 1. In other words, each credible set of gene-tissue pairs was counted only once. For
646 TGF_M, we first removed the gene-tissue pairs with PIPs < 0.5, and then applied the same
647 procedure as TGV_{IS} to count the silver and nearby genes. Similar to the mapping procedure for
648 silver genes, we examined each credible set identified as causal to see if it contained any
649 druggable genes. If a druggable gene is present, we count it once.

650 We used the following statistics to compare the enrichments of TGV_{IS} and TGF_M. For example,
651 regarding TGV_{IS} and a given trait, we consider three metrics: the number of causal genes
652 identified by TGV_{IS}, the overlap between genes identified by TGV_{IS} and those in the drug-
653 target list, and the ratio of these two metrics (referred to as Ratio hereafter). To compare whether
654 TGV_{IS} or TGF_M had a higher enrichment across traits, we performed a paired t-test using two
655 vectors of Ratio.

656 **Colocalization of credible sets**

657 We use colocalization to evaluate the causal credible sets identified by TGF_M and TGV_{IS}.
658 Within each region, we select variants from the outcome GWAS with P-values less than 5E-5
659 and $r^2 < 0.81$ for colocalization analysis. We perform fine-mapping on both the outcome and the
660 gene-tissue pairs within credible sets using SuSiE, then calculate the posterior probability for
661 hypothesis H_4 , i.e., both traits are associated and share the same single causal variant, between
662 each outcome and exposure pair using Coloc-SuSiE. We use a posterior probability of H_4 greater

663 than 0.5 as the threshold to determine colocalization evidence between gene-tissue pairs and the
664 outcome; notably, as long as at least one variant meets this criterion, it suffices.

665 **Mendelian Randomization using pQTL summary data**

666 We performed both univariable and multivariable MR using pQTLs of protein abundance as IVs
667 to evaluating the identified causal tissue-gene pairs. Because to the lack of tissue-specific protein
668 data, we focused on a subset of pGenes identified in blood tissues provided by Sun et al.¹⁰. We
669 selected independent, genome-wide significant pQTLs for each protein as IVs. The selection
670 method for independent IVs was C+T (--clump-kb 1000 --clump-p1 5e-8 --clump-p2 5e-8 --
671 clump-r2 0.01 using PLINK), with LD reference panels consisting of the 9,680 individuals and
672 9.3M variants from UKBB. We applied five univariable MR methods: MRMedian⁵⁶, IMRP⁵⁷,
673 MRCML⁵⁸, MRCUE⁵⁹, and MRBEE⁶⁰. Both MRCUE and MRBEE account for sample overlap,
674 with sample overlap correlations estimated using insignificant variants. We used the R package
675 FDREstimation to convert the P-values obtained by these methods to FDR Q-values, using “BY”
676 as the adjustment method. A pGene was considered significant if it was identified as such by at
677 least four methods. We also conducted an analysis comparing the enrichments of TGVIS and
678 TGFM, where the three corresponding metrics are: the overlap between causal genes identified
679 by TGVIS or TGFM and the pGenes reported in Sun et al.¹⁰, the number of significant pGenes
680 identified in univariable MR analysis, and the ratio of these two metrics.

681 **Author contributions**

682 Y.Y and X.Z conceived and designed the study. Y.Y. performed all analysis. Y.Y and X.Z.
683 drafted the manuscript. N.L. edited the manuscript. X.Z. supervised this project.

684 **Acknowledgments**

685 This work was supported by grant HG011052 and HG011052-03S1 (to X.Z.) from the National
686 Human Genome Research Institute (NHGRI), USA.

687 **Declarations of interests**

688 The authors declare no competing interests.

689 **Ethics approval**

690 The study was approved by the institutional review board (IRB number: STUDY20180592) at
691 Case Western Reserve University

692 **Data Availability**

693 The GWAS summary data, eQTL summary data, and pQTL summary data used in this study can
694 be downloaded from the “Data available” section of the literature listed in **Table S1-S2**. The
695 GTEx summary data can be obtained from <https://gtexportal.org/home/downloads/adult-gtex/ctl>.
696 The GWAS data in the million veteran program (MVP) are available through database of
697 genotypes and phenotypes (dbGAP) under accession number phs001672.v7.p1. The individual-
698 level data from the UKBB used for estimating the LD matrix was accessed through Application
699 ID: 81097.

700 **Code Availability**

701 The code for the analyses presented in this paper is available in the **Supplementary Materials**,
702 complete with step-by-step instructions. The TGVIS R package can be downloaded from
703 <https://github.com/harryiheyang/TGVIS>.

704 **References**

- 705 1. Graham, S. E. *et al.* The power of genetic diversity in genome-wide association studies of
706 lipids. *Nature* **600**, 675–679 (2021).
- 707 2. Yengo, L. *et al.* A saturated map of common genetic variants associated with human
708 height. *Nature* **610**, 704–712 (2022).
- 709 3. Suzuki, K. *et al.* Genetic drivers of heterogeneity in type 2 diabetes pathophysiology.
710 *Nature* **627**, 347–357 (2024).
- 711 4. Mostafavi, H., Spence, J. P., Naqvi, S. & Pritchard, J. K. Systematic differences in discovery
712 of genetic effects on gene expression and complex traits. *Nat. Genet.* **55**, 1866–1875
713 (2023).
- 714 5. Wallace, C. A more accurate method for colocalisation analysis allowing for multiple
715 causal variants. *PLOS Genet.* **17**, e1009440 (2021).
- 716 6. Gamazon, E. R. *et al.* A gene-based association method for mapping traits using
717 reference transcriptome data. *Nat. Genet.* **47**, 1091–1098 (2015).
- 718 7. Gkatzionis, A., Burgess, S. & Newcombe, P. J. Statistical methods for cis-Mendelian
719 randomization with two-sample summary-level data. *Genet. Epidemiol.* **47**, 3–25
720 (2023).
- 721 8. Wainberg, M. *et al.* Opportunities and challenges for transcriptome-wide association
722 studies. *Nat. Genet.* **51**, 592–599 (2019).
- 723 9. Li, Y. I. *et al.* RNA splicing is a primary link between genetic variation and disease.
724 *Science* **352**, 600–604 (2016).
- 725 10. Sun, B. B. *et al.* Plasma proteomic associations with genetics and health in the UK
726 Biobank. *Nature* **622**, 329–338 (2023).

- 727 11. Amariuta, T., Siewert-Rocks, K. & Price, A. L. Modeling tissue co-regulation estimates
728 tissue-specific contributions to disease. *Nat. Genet.* **55**, 1503–1511 (2023).
- 729 12. The GTEx Consortium *et al.* The GTEx Consortium atlas of genetic regulatory effects
730 across human tissues. *Science* **369**, 1318–1330 (2020).
- 731 13. Zhao, S. *et al.* Adjusting for genetic confounders in transcriptome-wide association
732 studies improves discovery of risk genes of complex traits. *Nat. Genet.* **56**, 336–347
733 (2024).
- 734 14. Strober, B. J., Zhang, M. J., Amariuta, T., Rossen, J. & Price, A. L. Fine-mapping causal
735 tissues and genes at disease-associated loci. Preprint at
736 <https://doi.org/10.1101/2023.11.01.23297909> (2023).
- 737 15. Burgess, S. & Thompson, S. G. *Mendelian Randomization*. (Chapman and Hall/CRC,
738 2021). doi:10.1201/9780429324352.
- 739 16. Wang, G., Sarkar, A., Carbonetto, P. & Stephens, M. A simple new approach to variable
740 selection in regression, with application to genetic fine mapping. *J. R. Stat. Soc. Ser. B*
741 *Stat. Methodol.* **82**, 1273–1300 (2020).
- 742 17. Cui, R. *et al.* Improving fine-mapping by modeling infinitesimal effects. *Nat. Genet.* **56**,
743 162–169 (2024).
- 744 18. Zou, Y., Carbonetto, P., Wang, G. & Stephens, M. Fine-mapping from summary data with
745 the “Sum of Single Effects” model. *PLOS Genet.* **18**, e1010299 (2022).
- 746 19. Benner, C. *et al.* FINEMAP: efficient variable selection using summary data from
747 genome-wide association studies. *Bioinformatics* **32**, 1493–1501 (2016).
- 748 20. Liang, Y., Nyasimi, F. & Im, H. K. On the problem of inflation in transcriptome-wide
749 association studies. Preprint at <https://doi.org/10.1101/2023.10.17.562831> (2023).

- 750 21. Sidorenko, J. *et al.* Genetic architecture reconciles linkage and association studies of
751 complex traits. *Nat. Genet.* (2024) doi:10.1038/s41588-024-01940-2.
- 752 22. Wood, S. N. Fast stable restricted maximum likelihood and marginal likelihood
753 estimation of semiparametric generalized linear models. *J. R. Stat. Soc. Ser. B Stat.*
754 *Methodol.* **73**, 3–36 (2011).
- 755 23. Aschard, H. A perspective on interaction effects in genetic association studies. *Genet.*
756 *Epidemiol.* **40**, 678–688 (2016).
- 757 24. Nielsen, J. B. *et al.* Biobank-driven genomic discovery yields new insight into atrial
758 fibrillation biology. *Nat. Genet.* **50**, 1234–1239 (2018).
- 759 25. Surendran, P. *et al.* Discovery of rare variants associated with blood pressure regulation
760 through meta-analysis of 1.3 million individuals. *Nat. Genet.* **52**, 1314–1332 (2020).
- 761 26. Sinnott-Armstrong, N. *et al.* Genetics of 35 blood and urine biomarkers in the UK
762 Biobank. *Nat. Genet.* **53**, 185–194 (2021).
- 763 27. Pazoki, R. *et al.* Genetic analysis in European ancestry individuals identifies 517 loci
764 associated with liver enzymes. *Nat. Commun.* **12**, 2579 (2021).
- 765 28. Young, W. J. *et al.* Genetic analyses of the electrocardiographic QT interval and its
766 components identify additional loci and pathways. *Nat. Commun.* **13**, 5144 (2022).
- 767 29. Ghouse, J. *et al.* Genome-wide meta-analysis identifies 93 risk loci and enables risk
768 prediction equivalent to monogenic forms of venous thromboembolism. *Nat. Genet.* **55**,
769 399–409 (2023).
- 770 30. Aragam, K. G. *et al.* Discovery and systematic characterization of risk variants and genes
771 for coronary artery disease in over a million participants. *Nat. Genet.* **54**, 1803–1815
772 (2022).

- 773 31. Roychowdhury, T. *et al.* Genome-wide association meta-analysis identifies risk loci for
774 abdominal aortic aneurysm and highlights PCSK9 as a therapeutic target. *Nat. Genet.*
775 **55**, 1831–1842 (2023).
- 776 32. Zhu, X., Zhu, L., Wang, H., Cooper, R. S. & Chakravarti, A. Genome-wide pleiotropy
777 analysis identifies novel blood pressure variants and improves its polygenic risk scores.
778 *Genet. Epidemiol.* **46**, 105–121 (2022).
- 779 33. Han, S. K. *et al.* Mapping genomic regulation of kidney disease and traits through high-
780 resolution and interpretable eQTLs. *Nat. Commun.* **14**, 2229 (2023).
- 781 34. Viñuela, A. *et al.* Genetic variant effects on gene expression in human pancreatic islets
782 and their implications for T2D. *Nat. Commun.* **11**, 4912 (2020).
- 783 35. Zhou, H. *et al.* FAVOR: functional annotation of variants online resource and annotator
784 for variation across the human genome. *Nucleic Acids Res.* **51**, D1300–D1311 (2023).
- 785 36. Trajanoska, K. *et al.* From target discovery to clinical drug development with human
786 genetics. *Nature* **620**, 737–745 (2023).
- 787 37. GTEx Consortium *et al.* Estimating the causal tissues for complex traits and diseases.
788 *Nat. Genet.* **49**, 1676–1683 (2017).
- 789 38. Arvanitis, M., Tayeb, K., Strober, B. J. & Battle, A. Redefining tissue specificity of genetic
790 regulation of gene expression in the presence of allelic heterogeneity. *Am. J. Hum. Genet.*
791 **109**, 223–239 (2022).
- 792 39. Burgess, S., Dudbridge, F. & Thompson, S. G. Combining information on multiple
793 instrumental variables in Mendelian randomization: comparison of allele score and
794 summarized data methods. *Stat. Med.* **35**, 1880–1906 (2016).

- 795 40. Grant, A. J. & Burgess, S. An efficient and robust approach to Mendelian randomization
796 with measured pleiotropic effects in a high-dimensional setting. *Biostatistics* **23**, 609–
797 625 (2022).
- 798 41. Barbeira, A. N. *et al.* Exploring the phenotypic consequences of tissue specific gene
799 expression variation inferred from GWAS summary statistics. *Nat. Commun.* **9**, 1825
800 (2018).
- 801 42. Kirk, T., Ahmed, A. & Rognoni, E. Fibroblast memory in development, homeostasis and
802 disease. *Cells* **10**, 2840 (2021).
- 803 43. Wittich, H. *et al.* Transcriptome-wide association study of the plasma proteome reveals
804 cis and trans regulatory mechanisms underlying complex traits. *Am. J. Hum. Genet.* **111**,
805 445–455 (2024).
- 806 44. Ross, S. D. *et al.* Clinical outcomes in statin treatment trials: A meta-analysis. *Arch.*
807 *Intern. Med.* **159**, 1793 (1999).
- 808 45. Xu, Y. *et al.* Rs7206790 and rs11644943 in FTO gene are associated with risk of obesity
809 in Chinese school-age population. *PLoS ONE* **9**, e108050 (2014).
- 810 46. Campos, A. I. *et al.* Boosting the power of genome-wide association studies within and
811 across ancestries by using polygenic scores. *Nat. Genet.* **55**, 1769–1776 (2023).
- 812 47. Yang, Y., Lee, D., Chakravarti, A. & Zhu, X. Partitioning tissue-specific heritability and
813 identifying tissue-specific causal genes of blood pressure traits using CRE annotations.
814 *Manuscript*.
- 815 48. Okamoto, J. *et al.* Probabilistic Fine-mapping of putative causal genes. Preprint at
816 <https://doi.org/10.1101/2024.10.27.620482> (2024).

- 817 49. Yazar, S. *et al.* Single-cell eQTL mapping identifies cell type-specific genetic control of
818 autoimmune disease. *Science* **376**, eabf3041 (2022).
- 819 50. Claussnitzer, M. *et al.* *FTO* obesity variant circuitry and adipocyte browning in humans.
820 *N. Engl. J. Med.* **373**, 895–907 (2015).
- 821 51. Zhu, X. & Stephens, M. Bayesian large-scale multiple regression with summary statistics
822 from genome-wide association studies. *Ann. Appl. Stat.* **11**, (2017).
- 823 52. Schwarz, G. Estimating the dimension of a model. *Ann. Stat.* **6**, (1978).
- 824 53. Zhang, D. & Lin, X. Hypothesis testing in semiparametric additive mixed models.
825 *Biostatistics* **4**, 57–74 (2003).
- 826 54. Willer, C. J., Li, Y. & Abecasis, G. R. METAL: fast and efficient meta-analysis of
827 genomewide association scans. *Bioinformatics* **26**, 2190–2191 (2010).
- 828 55. Purcell, S. *et al.* PLINK: A tool set for whole-genome association and population-based
829 linkage analyses. *Am. J. Hum. Genet.* **81**, 559–575 (2007).
- 830 56. Yavorska, O. O. & Burgess, S. MendelianRandomization: an R package for performing
831 Mendelian randomization analyses using summarized data. *Int. J. Epidemiol.* **46**, 1734–
832 1739 (2017).
- 833 57. Zhu, X., Li, X., Xu, R. & Wang, T. An iterative approach to detect pleiotropy and perform
834 Mendelian Randomization analysis using GWAS summary statistics. *Bioinformatics* **37**,
835 1390–1400 (2021).
- 836 58. Lin, Z., Xue, H. & Pan, W. Robust multivariable Mendelian randomization based on
837 constrained maximum likelihood. *Am. J. Hum. Genet.* **110**, 592–605 (2023).

- 838 59. Cheng, Q., Zhang, X., Chen, L. S. & Liu, J. Mendelian randomization accounting for
839 complex correlated horizontal pleiotropy while elucidating shared genetic etiology. *Nat.*
840 *Commun.* **13**, 6490 (2022).
- 841 60. Lorincz-Comi, N., Yang, Y., Li, G. & Zhu, X. MRBEE: A bias-corrected multivariable
842 Mendelian Randomization method. *Hum. Genet. Genomics Adv.* 100290 (2024)
843 doi:10.1016/j.xhgg.2024.100290.

It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

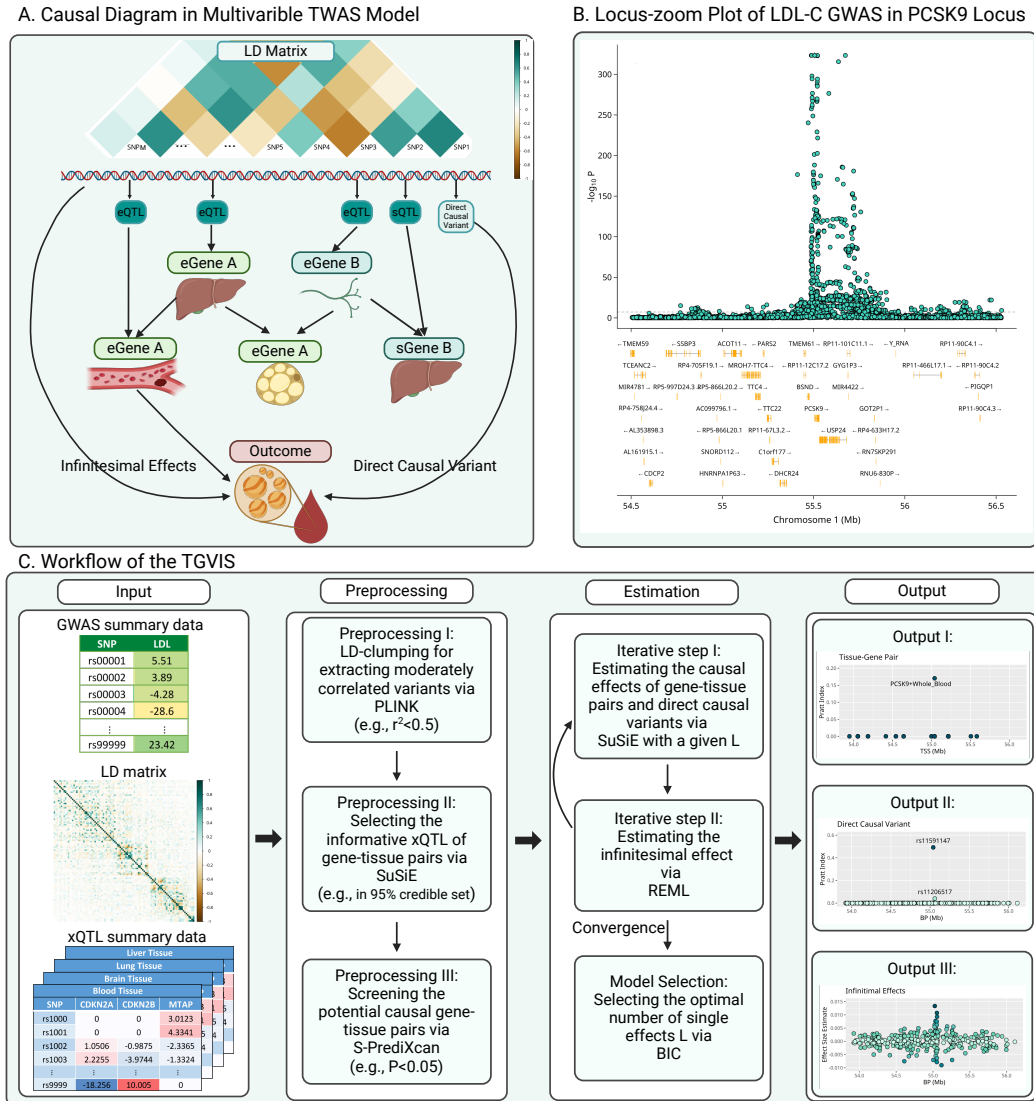


Figure 1: Overview of TGVIS. **A:** A hypothetical causal diagram illustrating the relationships between variants (including xQTLs, direct causal variants, and non-causal variants), tissue-specific gene expressions, and an outcome in a *cis*-region, where the arrows indicate the flow of causal effects in the causal diagram. Variants may be in LD, with only a subset having *cis*-regulatory effects. Gene expressions or splicing events are tissue-specific and form a complex co-regulation network. Only molecular phenotypes directly connected to the outcome are considered causal. **B:** Locus-zoom plot of the LDL-C GWAS in the *PCSK9* locus. The bottom panel displays the coding regions of genes located within this locus, including *PCSK9*, *UPS24*, *BSND*, etc. **C:** Workflow of TGVIS, consisting of three main steps. (I) Input, including GWAS summary data, eQTL summary data from multiple tissues, and LD matrix. (II) Preprocessing, including eQTL selection and pre-screening. We applied S-PrediXcan to pre-screen some noise pairs, aiming to reduce the dimension of the multivariable TWAS model to a reasonable scale. (III) Estimation, where TGVIS first selects the causal gene-tissue pairs and direct causal variants via SuSiE and then estimates the infinitesimal effect via REML. (IV) Output, including the causal effect estimate, direct causal effect estimate, and infinitesimal effect estimates. We output plots demonstrating the causal gene-tissue pairs, direct causal variants and predicted infinitesimal effects: (1) the Pratt indices and other statistics such as PIPs, estimates, SEs of causal gene-tissue pairs in the 95% credible sets, (2) the Pratt indices of the direct causal variants in the 95% credible sets, and (3) the best linear unbiased predictors of infinitesimal effects. The non-zero variance in output III in this figure suggests the non-zero contribution of infinitesimal effects.

It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

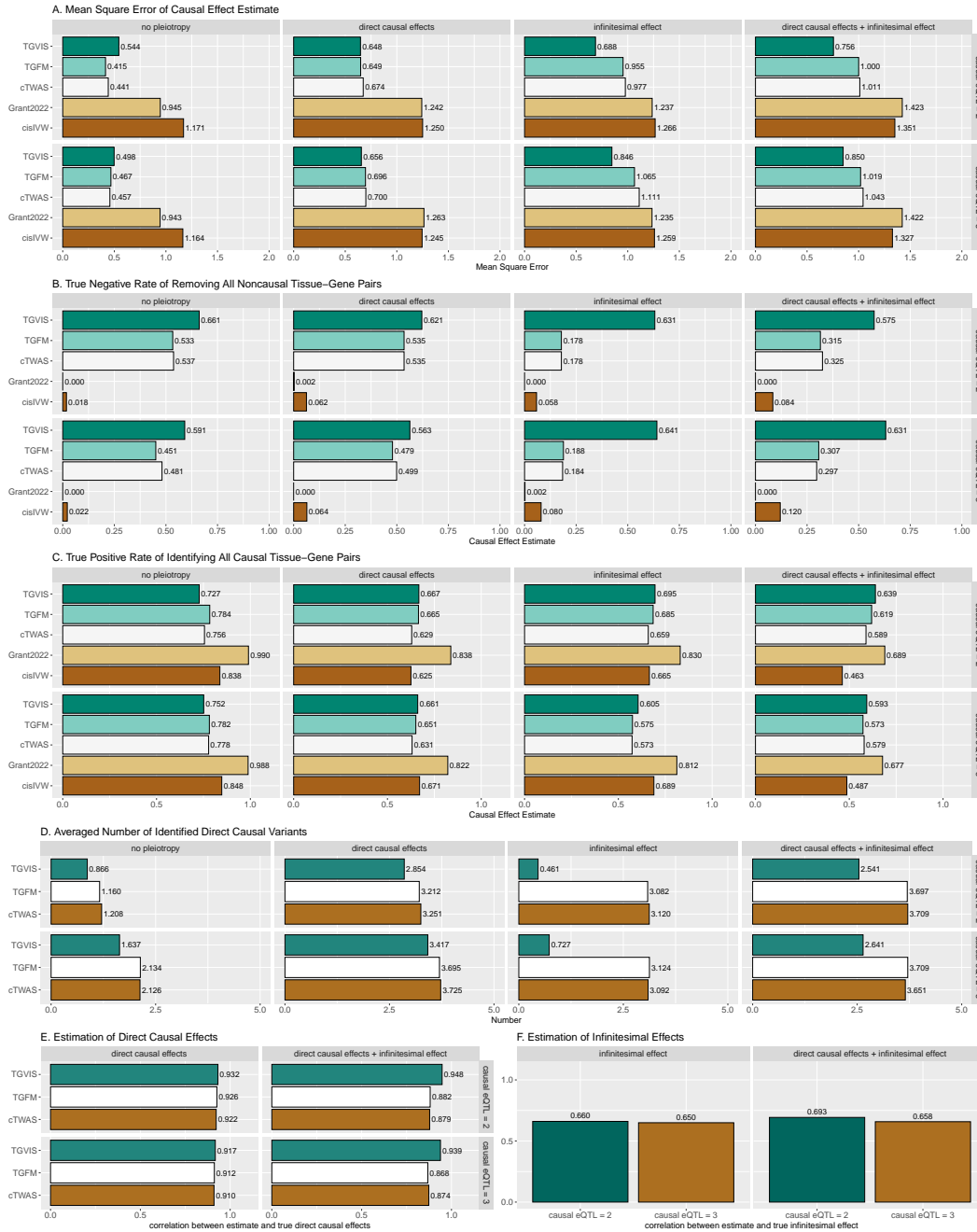


Figure 2: Simulation results comparing the performances of TGVIS, TGFm, cTWAS, Grant2022, and cisIVW with xQTL sample size = 200. **A:** The MSE of causal effect estimates under no pleiotropy, in the presence of direct causal variants, infinitesimal effects, and both. **B:** The true negative rate of identifying all 98 non-causal gene-tissue pairs under different scenarios i.e., no pleiotropy, in the presence of direct causal variants, infinitesimal effects, and both. This is equivalent to that if a method incorrectly identifies any non-causal pairs as causal, it will not be counted as a true negative event. **C:** Bar plots display the true positive rates of identifying all 2 causal gene-tissue pairs under different scenarios. **D:** The averaged number of identified direct causal variants by the different methods. The number of true causal variants were set to 0, 2, 0, and 2 for no-pleiotropy, direct-causal-variant, infinitesimal-effects, and direct-causal-variant and infinitesimal-effects, respectively. **E:** The averaged correlation of the true and estimated direct causal effects across simulations. **F:** The averaged correlation of the true and predicted infinitesimal effects across simulations.

It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

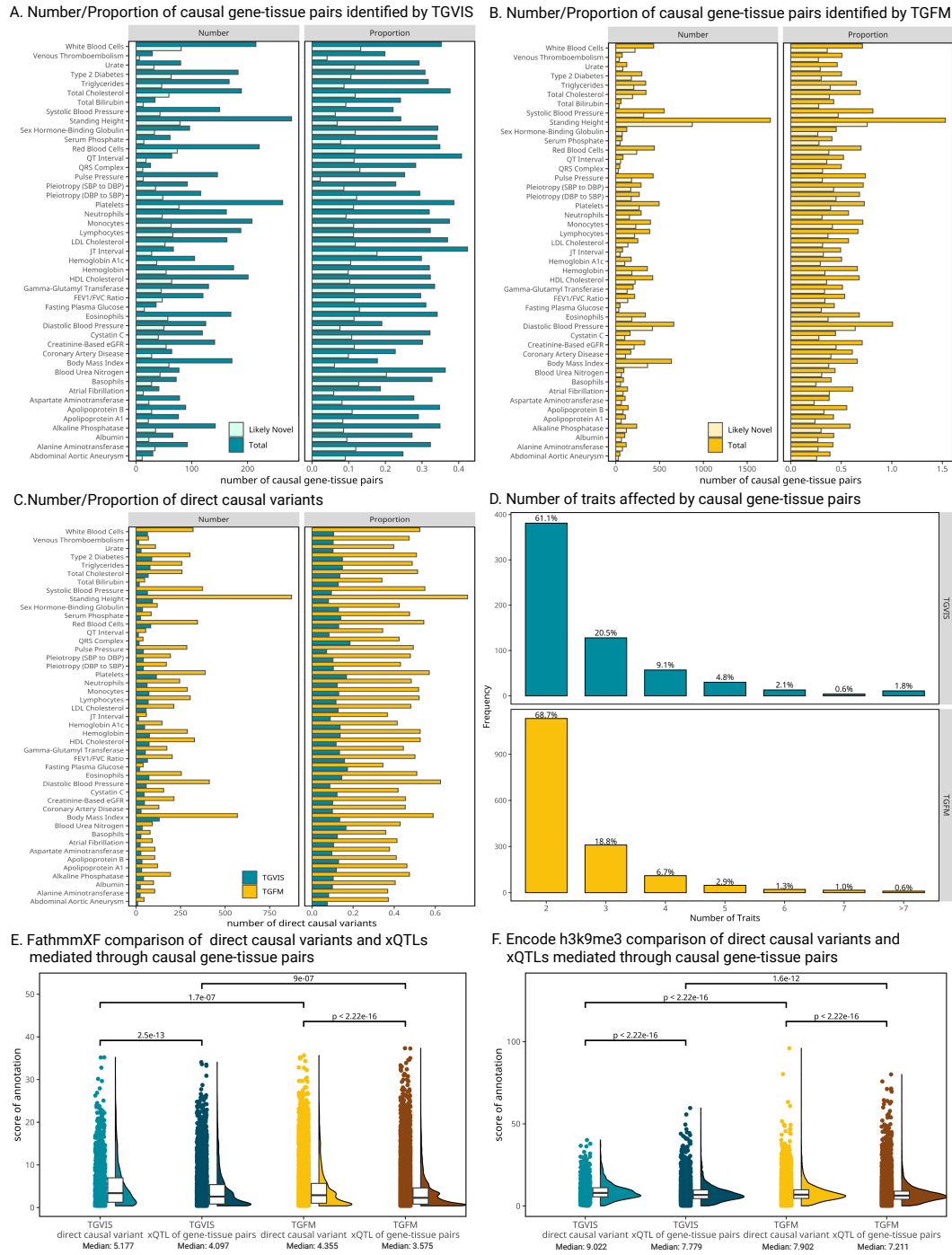


Figure 3: Summary of the identification of causal gene-tissue pairs and direct causal variants. **A-B**: The number and proportion of causal and likely novel causal gene-tissue pairs identified by TGVIS and TGF, respectively. Likely novel gene-tissue pairs are defined as those not present in the list of significant gene-tissue pairs identified by univariable S-PrediXcan ($P < 0.05/20000$). The proportion refers to the average number of causal and likely novel causal gene-tissue pairs per locus. **C**: The number and proportion of direct causal variants identified by TGVIS and TGF. **D**: The distribution of the number of traits affected by causal gene-tissue pairs. **E-F**: The distributions of scores for FathmmXF and Encode H3K9me3Sum annotations. Raincloud plots illustrate four classes: direct causal variants and xQTLs of causal gene-tissue pairs identified by TGVIS and TGF. Pairwise Wilcoxon signed-rank test P-values are displayed at the top, while medians of annotation scores are shown at the bottom.

It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

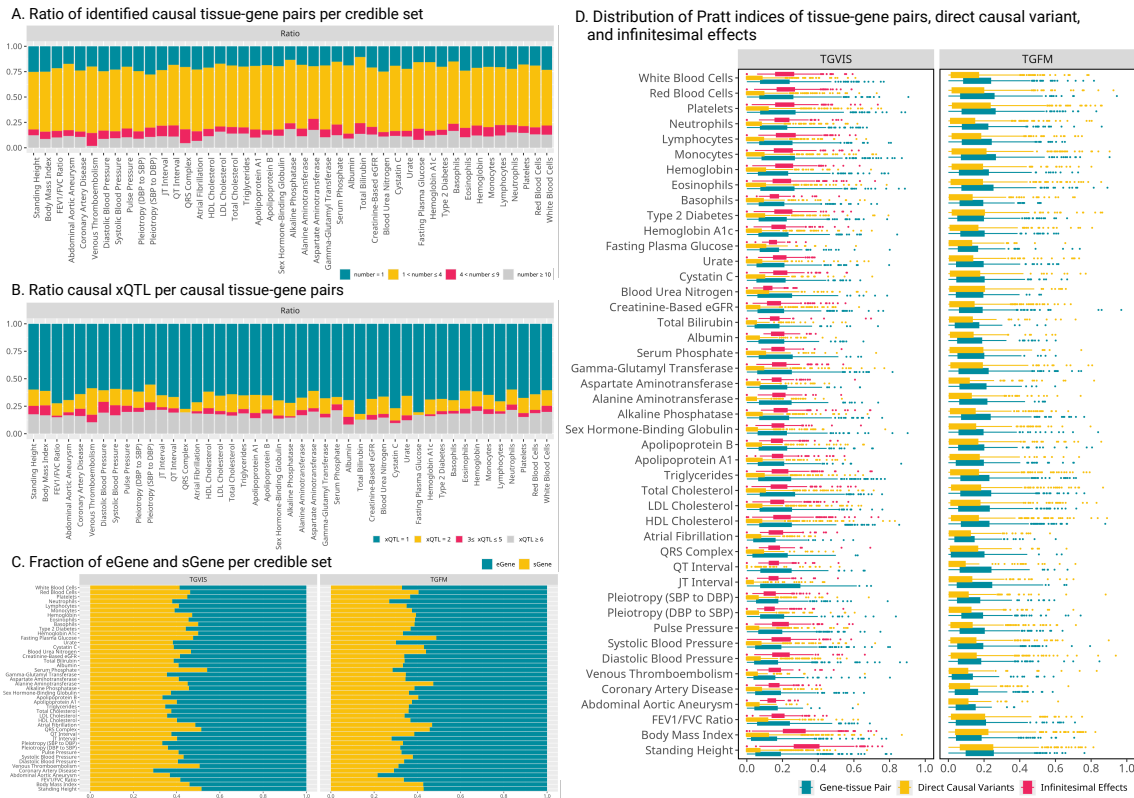


Figure 4: Genetic architecture inferred from the identification of causal gene-tissue pairs and direct causal variants. **A:** The ratio of identified causal gene-tissue pairs per credible set by TVGIS. Different gene-tissue pairs may share the same set of xQTLs, and end in the same credible set. **B:** The ratio of the number of causal eQTLs over the number of sQTLs per causal gene-tissue pair, indicating the distribution of eQTLs and sQTLs contributing to the gene-tissue pairs. **C:** The distribution of eGene and sGene in credible sets identified by TGVIS and TGMF. When a credible set contains multiple gene-tissue pairs, we calculate the proportion of eGenes and sGenes. **D:** The distribution of Pratt Index estimates for different traits, with a comparison between TGVIS and TGMF. In the boxplot, each point represents the Pratt Index of various molecular phenotypes within a single locus.

It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

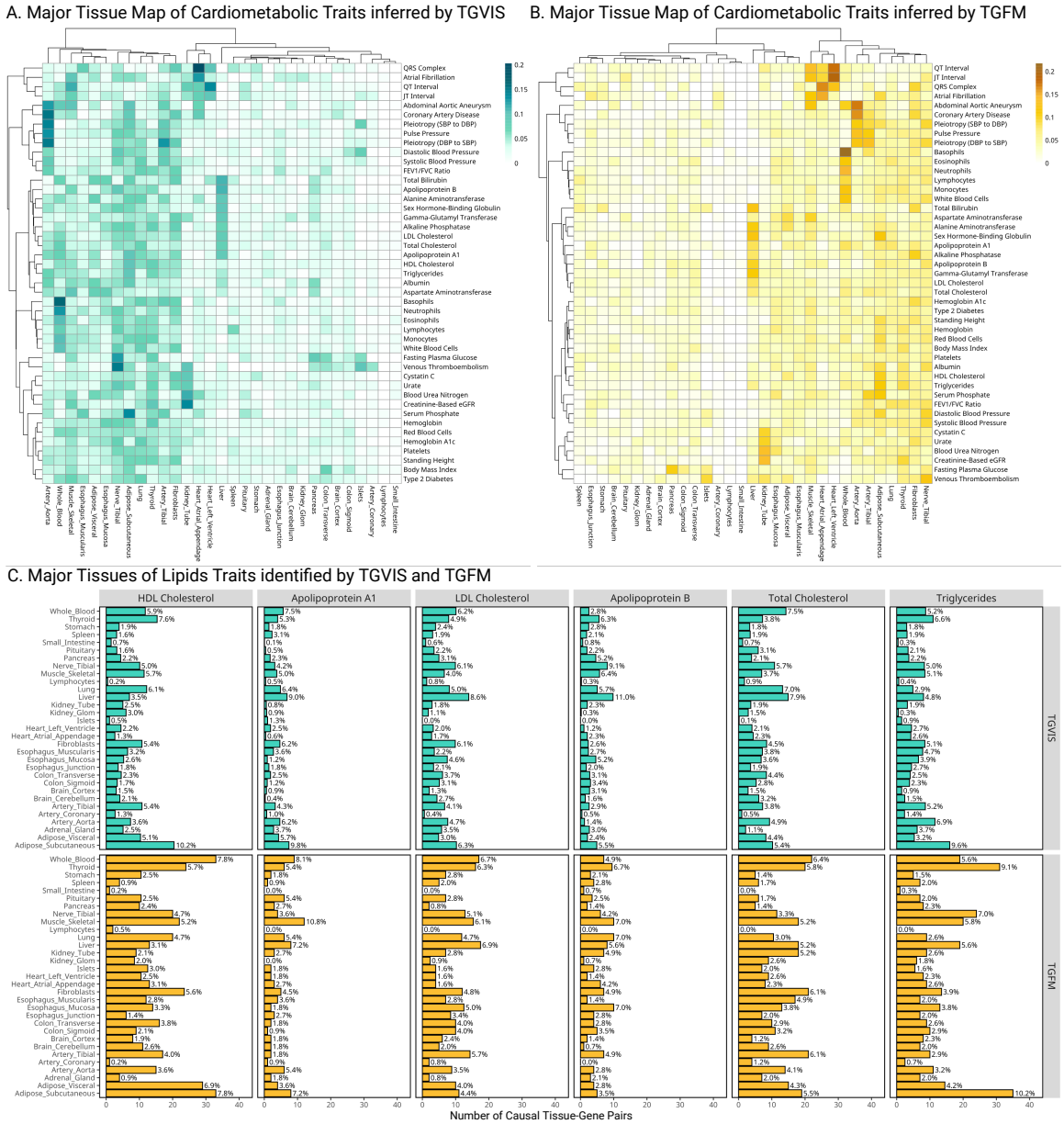
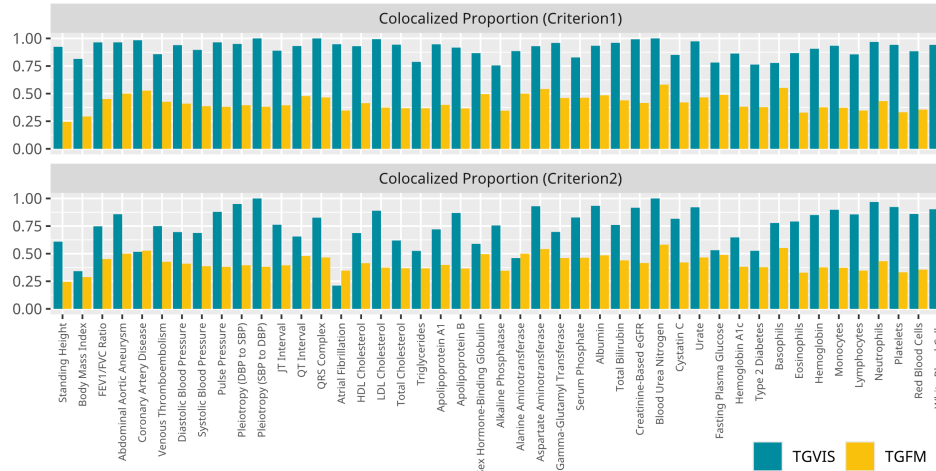


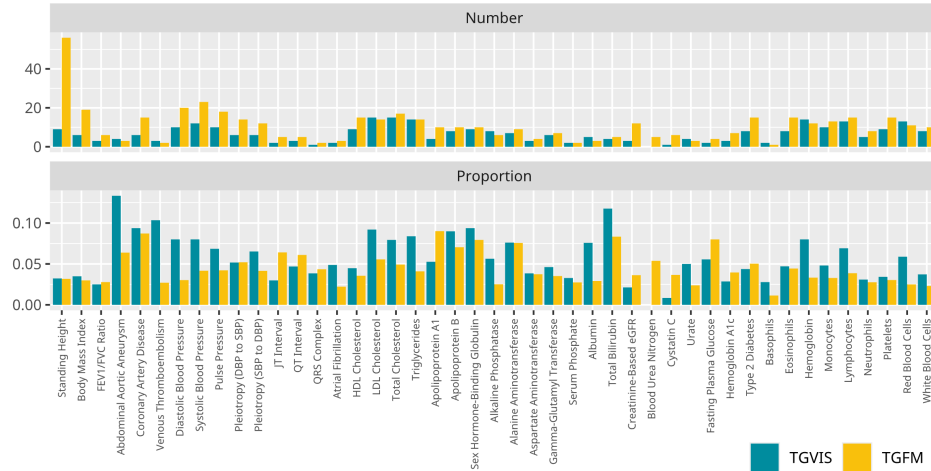
Figure 5: Distribution of major tissues for cardiometabolic traits. **A:** Heatmaps display the major tissues associated with each trait, identified by TGVIS. **B:** Heatmaps display the major tissues associated with each trait, identified by TGFm. The major gene-tissue pairs are cataloged based on stringent criteria ($CS\text{-}Pratt > 0.15$ for TGVIS and $PIP > 0.5$ for TGFm) and the proportions of major tissues derived from significant gene-tissue pairs for each trait are quantified. Hierarchical clustering is applied to arrange the heatmaps, utilizing the Ward2 method and Euclidean distance. **C:** Major tissues of lipid traits identified by TGVIS and TGFm. This panel shows bar plots detailing the number of causal gene-tissue pairs for various lipid traits, including HDL-C, LDL-C, TC, triglycerides, APOA1, and APOB, as identified by both TGVIS (top) and TGFm (bottom).

It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

A. Proportions of causal credible sets with colocalization evidence



B. Number/Proportion of causal genes in FDA-approved drug-target gene list



C. Number/Proportion of genes significant in univariable MR analysis

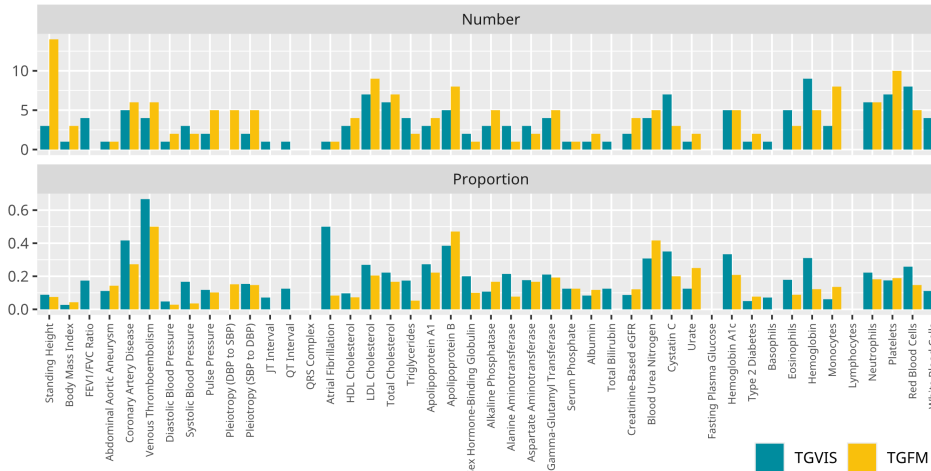


Figure 6: Evaluation of identified gene-tissue pairs. **A:** The colocalized proportions of causal credible sets (under two criteria) yielded by TGVIS and TGFM, respectively. **B:** The numbers and proportions of causal cis-genes in the list of FDA-approved drug-target genes provided by Trajanoska et al., identified by TGVIS (left) and TGFM (right), respectively. **C:** The number of significant pGenes in univariable MR analysis and the ratio of significant pGene in univariable MR analysis divided by significant eGenes/sGenes in eQTL/sQTL analysis.

It is made available under a [CC-BY-ND 4.0 International license](https://creativecommons.org/licenses/by-nd/4.0/).

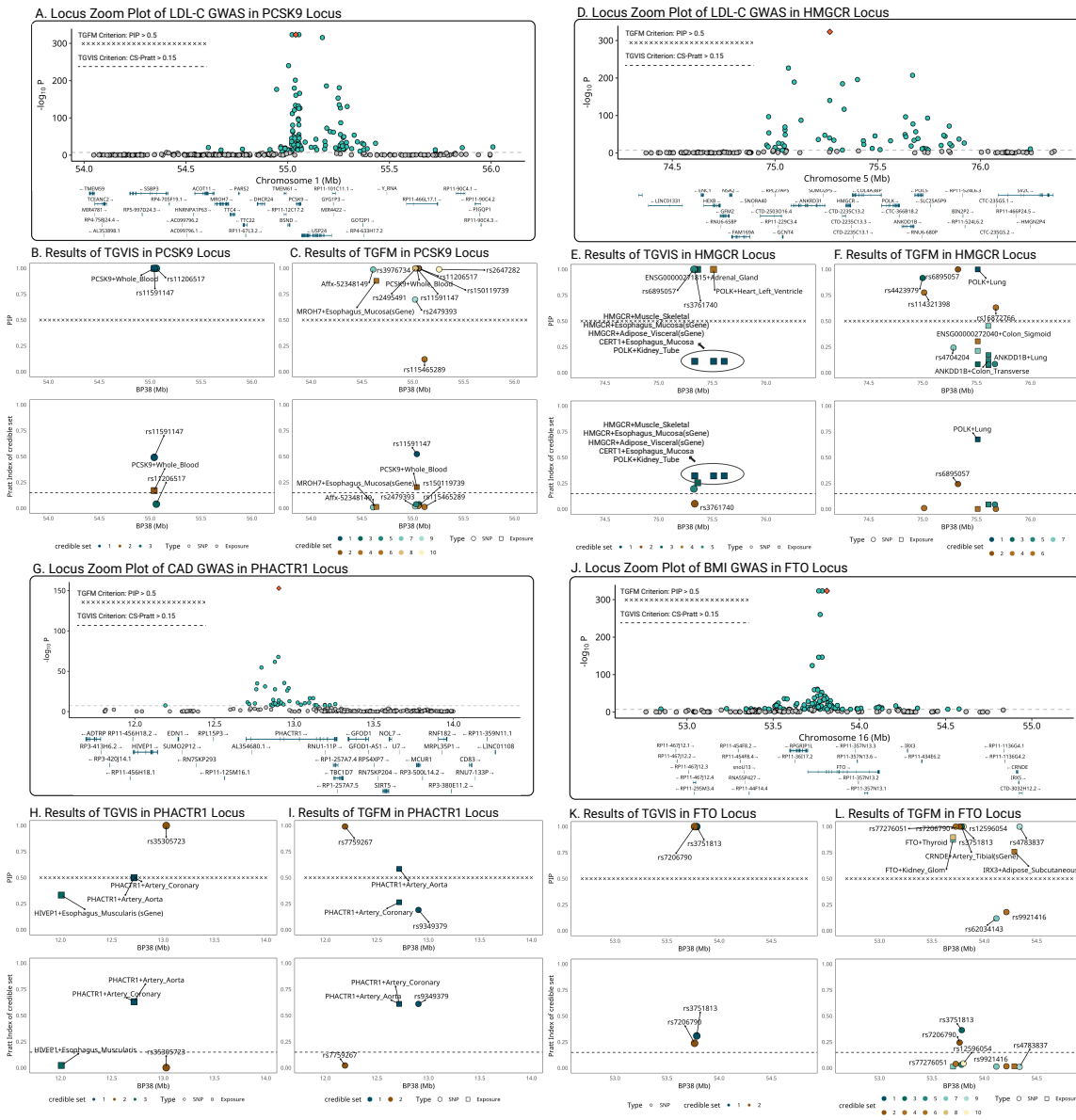


Figure 7: Locuszoom plots comparing the results of TGVIS and TGFm. **A-C**: LDL-C (*PCSK9* locus). **D-F**: LDL-C (*HMGCR* locus). **G-I**: CAD (*PHACTR1* locus), **K-L**: BMI (*FTO* locus). For each locus, we included three plots: (1) the GWAS of the trait, (2) the PIP of gene-tissue pairs and direct causal variants identified by the TGVIS and TGFm, and (3) the Pratt index of corresponding gene-tissue pairs and variants. For TGVIS, causality is determined by (1) the variables are in a 95% credible set and (2) the Pratt index of this credible set is larger 0.15. For TGFm, the causality is determined by (1) the individual PIP is larger than 0.5.