

 Open access • Posted Content • DOI:10.1101/2020.07.19.211151

Power analysis of transcriptome-wide association study — Source link

Bowei Ding, Chen Cao, Qing Li, Jingjing Wu ...+1 more authors

Institutions: University of Calgary, Alberta Children's Hospital

Published on: 26 Jul 2020 - bioRxiv (Cold Spring Harbor Laboratory)

Related papers:

- [Power analysis of transcriptome-wide association study: implications for practical protocol choice](#)
- [Model Checking via Testing for Direct Effects in Mendelian Randomization and Transcriptome-wide Association Studies](#)
- [Controlling bias and inflation in epigenome- and transcriptome-wide association studies using the empirical null distribution](#)
- [Testing and controlling for horizontal pleiotropy with the probabilistic Mendelian randomization in transcriptome-wide association studies](#)
- [The Future of and Beyond GWAS](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/power-analysis-of-transcriptome-wide-association-study-vebyqlbvgu>

1 **Power analysis of transcriptome-wide association study: implications for** 2 **practical protocol choice**

3 Chen Cao^{1,¶}, Bowei Ding^{2,¶}, Qing Li¹, Devin Kwok², Jingjing Wu^{2,*}, Quan Long^{1,2,3,4,*}

4

5 ¹Department of Biochemistry & Molecular Biology, Alberta Children's Hospital Research
6 Institute, University of Calgary, Calgary, Canada.

7 ²Department of Mathematics & Statistics, University of Calgary, Calgary, Canada.

8 ³Department of Medical Genetics, University of Calgary, Calgary, Canada.

9 ⁴Hotchkiss Brain Institute, O'Brien Institute for Public Health, University of Calgary, Calgary,
10 Canada.

11

12 ¶ = Joint first authors: C.C. and B.D.

13 * = Joint corresponding authors: J.W. (jinwu@ucalgary.ca) and Q.L. (quan.long@ucalgary.ca)

14

15 **Abstract**

16

17 The transcriptome-wide association study (TWAS) has emerged as one of several promising
18 techniques for integrating multi-scale 'omics' data into traditional genome-wide association
19 studies (GWAS). Unlike GWAS, which associates phenotypic variance directly with genetic
20 variants, TWAS uses a reference dataset to train a predictive model for gene expressions, which
21 allows it to associate phenotype with variants through the mediating effect of expressions.

22 Although effective, this core innovation of TWAS is poorly understood, since the predictive
23 accuracy of the genotype-expression model is generally low and further bounded by expression
24 heritability. This raises the question: to what degree does the accuracy of the expression model
25 affect the power of TWAS? Furthermore, would replacing predictions with actual,

26 experimentally determined expressions improve power? To answer these questions, we

27 compared the power of GWAS, TWAS, and a hypothetical protocol utilizing real expression

28 data. We derived non-centrality parameters (NCPs) for linear mixed models (LMMs) to enable
29 closed-form calculations of statistical power that do not rely on specific protocol

30 implementations. We examined two representative scenarios: causality (genotype contributes to

31 phenotype through expression) and pleiotropy (genotype contributes directly to both phenotype
32 and expression), and also tested the effects of various properties including expression
33 heritability. Our analysis reveals two main outcomes: (1) Under pleiotropy, the use of predicted
34 expressions in TWAS is superior to actual expressions. This explains why TWAS can function
35 with weak expression models, and shows that TWAS remains relevant even when real
36 expressions are available. (2) GWAS outperforms TWAS when expression heritability is below a
37 threshold of 0.04 under causality, or 0.06 under pleiotropy. Analysis of existing publications
38 suggests that TWAS has been misapplied in place of GWAS, in situations where expression
39 heritability is low.

40

41 **Keywords:** Power analysis, GWAS, TWAS, Non-centrality parameter, Expression heritability

42 **Author Summary**

43 We compared the effectiveness of three methods for finding genetic effects on disease in
44 order to quantify their strengths and help researchers choose the best protocol for their data. The
45 genome-wide association study (GWAS) is the standard method for identifying how the genetic
46 differences between individuals relate to disease. Recently, the transcriptome-wide association
47 study (TWAS) has improved GWAS by also estimating the effect of each genetic variant on the
48 activity level (or expression) of genes related to disease. The effectiveness of TWAS is
49 surprising because its estimates of gene expressions are very inaccurate, so we ask if a method
50 using real expression data instead of estimates would perform better. Unlike past studies, which
51 only use simulation to compare these methods, we incorporate novel statistical calculations to
52 make our comparisons more accurate and universally applicable. We discover that depending on
53 the type of relationship between genetics, gene expression, and disease, the estimates used by
54 TWAS could be actually more relevant than real gene expressions. We also find that TWAS is
55 not always better than GWAS when the relationship between genetics and expression is weak
56 and identify specific turning points where past studies have incorrectly used TWAS instead of
57 GWAS.

58

59 **Introduction**

60 High-throughput sequencing instruments have enabled the rapid profiling of
61 transcriptomes (RNA expression of genes) [1-4], proteomes (proteins) [5-7] and other ‘omics’
62 data [8-10]. These ‘omics’ provide insight into the intermediary effects of genotypes on
63 endophenotypes, and can improve the ability of genome-wide association studies (GWAS) to
64 find associations between genetic variants and disease phenotypes. [11-13]. The integration of
65 diverse ‘omics’ data sources remains a challenging and active field of research [14-17].

66 One approach to integrating ‘omics’ and GWAS is the transcriptome-wide association
67 study (TWAS), which quantitatively aggregates multiple genetic variants into a single test using
68 transcriptome data. Pioneered by Gamazon *et al* [18], the TWAS protocol typically has two
69 steps. First, a model is trained to predict gene expressions from local genetic variants near the
70 focal genes, using a reference dataset containing both genotype and expression data. Second, the
71 pretrained model is used to predict expressions from genotypes in the association mapping
72 dataset under study, which contains genotypes and phenotypes (but not expression). The

73 predicted expressions are then associated to the phenotype of interest. TWAS can also be
74 conducted with summary statistics from GWAS datasets (i.e. meta-analysis) as first
75 demonstrated by Gusev *et al.* [19] [20]. TWAS has since achieved significant popularity and
76 success in identifying the genetic basis of complex traits [21-27], inspiring similar protocols for
77 other endophenotypes such as IWAS for images [28] and PWAS for proteins [29].

78 Despite its demonstrated effectiveness, important questions remain regarding the
79 theoretical conditions under which TWAS is superior to GWAS. **First:** TWAS mapping relies
80 entirely on predicted expressions, but as shown by many methodological papers, the mean R^2
81 between predicted and actual expressions is very low (around 0.02 ~ 0.05). This is in part due to
82 low expression heritability [18], which bounds the maximum predictive accuracy attainable by
83 the genotype-expression model. Naturally, one can ask: given sufficiently low expression
84 heritability, is there is a point at which TWAS performs worse than GWAS? Indeed in real data,
85 genes discovered with significant TWAS p -values tend to have a higher R^2 , and thus expression
86 heritability, than on average [18, 19, 30-32]. We therefore investigate the effect of expression
87 heritability on the power of TWAS, as well as its interactions with trait heritability, phenotypic
88 variance from expressions, number of causal genes, and genetic architecture. **Second:** as
89 described by Gamazon *et al.* [18], the key insight of TWAS is that it aggregates sensible genetic
90 variants to estimate “genetically regulated gene expression”, or GReX [18], for use in
91 downstream GWAS. Given this hypothesis, one may ask if actual expression data would further
92 improve the power of downstream GWAS over predicted expressions. This is not a trivial
93 question, as although actual expressions do not suffer from prediction errors, they also include
94 experimental or environmental noise which masks the genetic component of expression. To test
95 this problem, we invent a hypothetical protocol associating real expressions to phenotype, which
96 we call “expression mediated GWAS” or emGWAS. While emGWAS is not in practical use due
97 to the difficulties of accessing relevant tissues (e.g., in the studies of brain diseases), it can
98 potentially be applied to future analyses of diseases where tissues are routinely available (e.g.,
99 blood or cancerous tissues). More importantly, emGWAS serves as a useful benchmark for
100 evaluating the theoretical properties of TWAS-predicted expressions against ground truth
101 expression data. By analyzing the power of TWAS, GWAS, and emGWAS, we develop practical
102 guidelines for choosing each protocol given different expression heritability and genetic
103 architectures.

104 While there has been an existing study comparing the power of GWAS, TWAS, and a
105 protocol which integrates eQTLs with GWAS [33], the existing study is purely simulation-based,
106 whereas we determine power directly using traditional closed-form analysis. We derive non-
107 centrality parameters (NCPs) for the relevant statistical tests and the linear mixed model (LMM)
108 in particular (**Methods**). Our derivation uses a novel method to convert an LMM into a linear
109 regression by decorrelating the covariance structure of the LMM response variable (**Methods**).
110 To our best knowledge, this is the first closed-form derivation of the NCP for LMMs in current
111 literature, with potential for broad applications as LMMs are the dominant models used in
112 GWAS and portions of the TWAS pipeline.

113 Unlike pure simulations, which stochastically resample the alternative hypothesis to
114 estimate statistical power, our closed-form derivation directly calculates power from a particular
115 configuration of association mapping data. As a result, our method saves computational
116 resources, yields more accurate power estimations, and adapts easily to similar protocols such as
117 IWAS [28] and PWAS [29, 34]. Moreover, as the closed-form derivation avoids conducting the
118 actual regression, our power calculations do not depend on specific implementations of GWAS
119 and TWAS, which could otherwise cause our results to vary due to differences in filtering inputs
120 or parameter optimizations. Our work therefore characterizes the theoretical power of the
121 protocols across all LMM-based implementations and datasets, although we are unable to
122 account for power losses due to practical implementation issues.

123 In the following section we describe our novel derivation of NCPs for LMMs and our
124 power analyses of GWAS, TWAS, and emGWAS. We present guidelines on the applicability of
125 each protocol under different input conditions and discuss potential limitations of our approach
126 as well as areas for future research.

127

128 **Materials & Methods**

129 **Mathematical definitions of GWAS, TWAS, and emGWAS protocols**

130 While there are many variations of GWAS and TWAS [18, 19, 35-39], in this work we
131 assume that multiple genes contribute to phenotypic variation, and for each causal gene, multiple
132 single nucleotide polymorphisms (SNPs) contribute to both gene expression and phenotype. This
133 setting is motivated by the fact that most complex traits are known to have multiple contributing
134 loci, and TWAS fundamentally assumes that genes have multiple local causal variants. To ensure

135 consistency, we apply the same assumptions in the design of the hypothetical protocol
136 emGWAS. Specifically, we define the following models:

137 **GWAS.** For GWAS, we adopted a standard LMM similar to EMMAX [35]

$$138 \quad Y = \beta_{j_0} \mathbf{1} + X_j \beta_{j_1} + u + \varepsilon, \quad j = 1, 2, \dots, n_x, \quad (1)$$

139 where n is the number of individuals, n_x is the total number of genetic variants, Y is an $n \times 1$
140 vector of phenotypes, $\mathbf{1}$ is an $n \times 1$ vector of ones, X_j is an $n \times 1$ genotype vector with $X_{ij} \in$
141 $\{0, 1, 2\}$ representing the number of minor allele copies for the i^{th} individual and j^{th} genetic
142 variant, β_{j_0} and β_{j_1} are the intercept and effect size of the genetic variant, u is an $n \times 1$ vector of
143 random effects following the multivariate normal distribution, i.e. $u \sim N(0, \sigma_g^2 K_x)$, and ε is an
144 $n \times 1$ vector of errors with $\varepsilon \sim N(0, \sigma_e^2 I)$. In the distributions of u and ε , σ_g^2 and σ_e^2 are their
145 respective variance components, I is an $n \times n$ identity matrix, and K_x is the genomic relationship
146 matrix (GRM), which is a known $n \times n$ real symmetric matrix. Following Patterson *et al* [40],
147 K_x is calculated by

$$148 \quad K_x = \frac{1}{n_x} \tilde{X} \tilde{X}^T, \quad (2)$$

149 where n_x is the total number of genetic variants and \tilde{X} is a standardized $n \times n_x$ matrix. For
150 example, an element \tilde{X}_{ij} in the j^{th} genetic variant column is calculated as

$$151 \quad \tilde{X}_{ij} = \frac{X_{ij} - \bar{X}_{.j}}{S_{Xj}}, \quad (3)$$

152 where $\bar{X}_{.j} = \frac{1}{n} \sum_{i=1}^n X_{ij}$ and $S_{Xj}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_{ij} - \bar{X}_{.j})^2$ are the sample mean and sample variance
153 of the j^{th} variant, respectively.

154 **emGWAS.** For emGWAS, we first regress the phenotype on the actual (not predicted)
155 expressions, and then regress the expressions on individual local genetic variants in a similar
156 manner as a cis-eQTL analysis. We chose the LMM to associate phenotype with expression,
157 since under the assumption that multiple genes contribute to phenotype, we expect that the
158 random term of the LMM can capture the effects of non-focal genes. We calculate the GRM
159 from DNA instead of expressions because they provide better estimates of pairwise relationships
160 between study participants than correlations based on predicted expression data. ~~using the~~
161 ~~assumption that the ultimate goal is to identify genetic variants underlying expressions.~~ We
162 chose to use linear regression (LM) to model the association between expression and local

163 genetic variants (which correspond to cis-eQTLs), as it is the most common model used in cis-
164 eQTL analyses.

165 Specifically, the phenotype-expression model is

$$166 \quad Y = \beta_{l0}\mathbf{1} + \beta_{l1}Z_l + u + \varepsilon, \quad l = 1, 2, \dots, n_z, \quad (4)$$

167 where n , Y , $\mathbf{1}$, u and ε have identical interpretations as in the GWAS model from (1), n_z is the
168 total number of genes, Z_l is an $n \times 1$ gene expression vector for the l^{th} gene, and β_{l0} and β_{l1} are
169 the intercept and effect size of the gene.

170 The linear regression associating gene expression with local genetic variants is

$$171 \quad Z_l = \beta_{lk0}\mathbf{1} + \beta_{lk1}X_{lk} + \varepsilon_{el}, \quad l = 1, 2, \dots, n_z, k = 1, 2, \dots, n_{el}, \quad (5)$$

172 where X_{lk} is an $n_{el} \times 1$ vector of the k^{th} local genetic variants for the l^{th} gene, $\varepsilon_{el} \sim N(0, \sigma_{el}^2 I)$
173 is a $n \times 1$ vector of errors with variance component σ_{el}^2 , n_{el} is the total number of local genetic
174 variants in the l^{th} gene, and β_{lk0} and β_{lk1} are the intercept and effect size of the variant.

175 **TWAS.** For TWAS, we apply an analysis similar to emGWAS, except that gene
176 expressions are predicted using a pretrained elastic-net model. Specifically,

$$177 \quad Y = \beta_{Pl0}\mathbf{1} + \beta_{Pl1}\hat{Z}_l + u + \varepsilon, \quad l = 1, 2, \dots, n_z, \quad (6)$$

178 where \hat{Z}_l is the altered notation representing an $n \times 1$ vector of predicted gene expressions for
179 the l^{th} gene, and β_{Pl0} and β_{Pl1} are the intercept and effect size of the predicted gene expression.

180 There are several methods to estimate gene expression including least absolute shrinkage
181 and selection operator (LASSO) and elastic-net. Gamazon *et al.* has shown that elastic-net has
182 good performance and is more robust to minor changes in the input variants [18]. We therefore
183 use the “glmnet” package in R to train a predictive model using elastic-net. The objective
184 function in “glmnet” is

$$185 \quad L_{enet}(\beta) = \frac{1}{2n} \|Z - X\beta\|^2 + \lambda \left(\frac{1 - \alpha}{2} \|\beta\|^2 + \alpha \|\beta\|_1 \right) \quad (7)$$

186 where λ and α are tuning parameters. The penalty term is a convex (linear) combination of
187 LASSO and ridge penalties, where $\alpha = 1$ is equivalent to the LASSO objective function, and
188 $\alpha = 0$ is equivalent to ridge regression. Optimal values of λ and α were chosen by minimizing
189 the cross-validated squared-error. Readers are referred to **S1 Appendix** for details.

190 In practice, the specific regression model varies depending on the tool in use. For
191 example, the leading TWAS tool PrediXcan [18] does not include the random effects of a mixed
192 model, and many TWAS tools can also analyze summary statistics instead of subject-level

193 genotypes [19]. The motivation of this work is to reveal the key issues of using gene expressions
194 as mediations, therefore has to adapt comparable framework. In other works, we do not intend to
195 compare LMM against linear regression, which will mislead the comparison between GWAS
196 and TWAS. Since LMMs are dominant in GWAS, we chose LMMs as the underlying model for
197 all of the protocols we analyze, which allows us to compare them under an equivalent statistical
198 framework. We believe that LMMs are a sensible approach for TWAS, since the random term
199 can capture the genetic contributions of non-focal genes.

200

201 **Closed-form derivation of NCP and power calculation**

202 The non-centrality parameter (NCP) measures the distance between a non-central
203 distribution and a central distribution under a specific alternative hypothesis. The NCP enables
204 calculation of the probability of rejecting the null hypothesis, assuming the central distribution,
205 when the alternative hypothesis is correct. As such, the NCP naturally allows the power of a
206 statistical test to be determined in a closed form. We have developed the following method to
207 derive the NCP for LMMs, which we believe is new to the literature.

208 For a standard simple linear regression, the NCP of a t -test of the coefficient of the
209 predictor variable can be derived similarly to a one-sample t -test statistic as follows: if
210 $X_1, \dots, X_n \sim N(\mu, \sigma)$ is a simple random sample, then the one-sample t -test statistic for evaluating
211 the null hypothesis $H_0: \mu = \mu_0$ is

$$212 \quad T = \frac{\bar{X} - \mu_0}{S/\sqrt{n}} = \frac{\sqrt{n}(\bar{X} - \mu_0)/\sigma}{\sqrt{\frac{(n-1)S^2/\sigma^2}{n-1}}} \sim t_{n-1}, \quad (8)$$

213 where \bar{X} and S are the sample mean and (unbiased) sample standard deviation respectively.

214 Under H_0 , $\sqrt{n}(\bar{X} - \mu_0)/\sigma \sim N(0,1)$ and $(n-1)S^2/\sigma^2 \sim \chi_{n-1}^2$, and thus $T \sim t_{n-1}$. Under the
215 alternative hypothesis $H_a: \mu = \mu_a$, the test statistic $T = \frac{\sqrt{n}[(\bar{X} - \mu_a) + (\mu_a - \mu_0)]/\sigma}{\sqrt{\frac{(n-1)S^2/\sigma^2}{n-1}}}$ follows a non-

216 central t distribution with NCP given by

$$217 \quad v = \frac{\mu_a - \mu_0}{\sigma/\sqrt{n}} \quad (9)$$

218 To derive a closed-form NCP for LMMs, we convert the LMM to a linear regression
219 without intercept by decorrelating the response variable and the predictors, a technique that has
220 previously been applied to mixed models [41, 42]. The procedure is as follows: we first fit the

221 null model $Y_c = u + \varepsilon$ with no genetic variants, following an existing innovation for reducing the
 222 computational cost of repeatedly factorizing the GRM when analyzing many variants [35, 42].

223 We then estimate σ_g^2 using the Newton-Raphson method detailed in **S2 Appendix**. Denoting the
 224 eigen decomposition of the GRM as $K_x = U_x \Lambda_x U_x^{-1}$, we construct the de-correlation matrix as

$$225 \quad D_x = (\sigma_g^2 \Lambda_x + \sigma_e^2 I)^{-\frac{1}{2}} U_x^T. \quad (10)$$

226 By left multiplying both X and Y by D_x , and denoting $X^* = D_x X = (X_1^*, X_2^*, \dots, X_n^*)^T$ and
 227 $Y^* = D_x Y = (Y_1^*, Y_2^*, \dots, Y_n^*)^T$, the covariance structure in Y^* is thus removed and a linear
 228 regression of Y^* on X^* is equivalent to the original LMM model. A proof of the validity of this
 229 decorrelation structure is presented in **S3 Appendix**.

230 Based on the closed-form NCP for linear regression, we derive the estimated NCP of the
 231 LMM from (1), which is given by

$$232 \quad \hat{v}_{Gj} = \frac{\sum_{i=1}^n \hat{X}_{ij}^* \hat{Y}_i^* \sum_{i=1}^n \hat{D}_{xi}^2 - \sum_{i=1}^n \hat{Y}_i^* \hat{D}_{xi} \cdot \sum_{i=1}^n \hat{X}_{ij}^* \hat{D}_{xi}}{\sqrt{\sum_{i=1}^n (\hat{X}_{ij}^*)^2 (\sum_{i=1}^n \hat{D}_{xi})^2 - (\sum_{i=1}^n \hat{D}_{xi} \cdot \hat{X}_{ij}^*)^2 \sum_{i=1}^n \hat{D}_{xi}^2}}, \quad (11)$$

233 where $\hat{X}_j^* = \hat{D}_x X_j = (\hat{X}_{1j}^*, \hat{X}_{2j}^*, \dots, \hat{X}_{nj}^*)^T$, $\hat{Y}^* = \hat{D}_x Y = (\hat{Y}_1^*, \hat{Y}_2^*, \dots, \hat{Y}_n^*)^T$, and $\hat{D}_{xi} = \sum_{j=1}^n \hat{D}_{xij}$. A
 234 proof of this expression of the NCP for LMMs is in **S4 Appendix**.

235 The above result allows us to derive the statistical power of the GWAS, emGWAS, and
 236 TWAS protocols. For GWAS, we use the Bonferroni-corrected significance level $\alpha_x = \frac{0.05}{n_x}$ to
 237 account for multiple testing [43], where n_x is the total number of SNPs. Throughout this paper,
 238 we use $f(t; v)$ to denote the probability density function of the non-central t distribution with $n-2$
 239 degrees of freedom and NCP v . The statistical power of the j^{th} SNP can then be estimated by
 240 $P_{Gj} = \int_{F_0^{-1}(1-\alpha_x)}^{+\infty} f(t; \hat{v}_{Gj}) dt$ using the estimated NCP \hat{v}_{Gj} , where $F_0(t)$ is the cumulative
 241 distribution function of the central t distribution with $n-2$ degrees of freedom, and $F_0^{-1}(1 - \alpha_x)$
 242 gives the critical value for the central distribution. We directly implement this power
 243 computation in R via the function “pt”, which takes the critical value, NCP, and degrees of
 244 freedom as parameters.

245 For emGWAS, we assume that the powers of the expression-phenotype and genotype-
 246 expression regression models (4) and (5) are independent of each other. For the model $Y =$
 247 $\beta_{l0} \mathbf{1} + Z_l \beta_{l1} + u + \varepsilon$ from (4), we left multiply the estimated \hat{D}_x to both sides of the equation so
 248 that the estimated NCP for the l^{th} gene expression is given by

$$\hat{v}_{eZl} = \frac{\sum_{i=1}^n \hat{Z}_{il}^* \hat{Y}_i^* \sum_{i=1}^n \hat{D}_{xi}^2 - \sum_{i=1}^n \hat{Y}_i^* \hat{D}_{xi} \cdot \sum_{i=1}^n \hat{Z}_{il}^* \hat{D}_{xi}}{\sqrt{\sum_{i=1}^n (\hat{Z}_{il}^*)^2 (\sum_{i=1}^n \hat{D}_{xi}^2)^2 - (\sum_{i=1}^n \hat{D}_{xi} \cdot \hat{Z}_{il}^*)^2 \sum_{i=1}^n \hat{D}_{xi}^2}}, \quad (12)$$

where $\hat{Z}_l^* = \hat{D}_x Z_l = (\hat{Z}_{1l}^*, \hat{Z}_{2l}^*, \dots, \hat{Z}_{nl}^*)^T$. We use the significance level $\alpha_z = \frac{0.05}{n_z}$ for each individual test, where n_z is the total number of genes. The statistical power of detecting the l^{th} gene expression is then estimated by $P_{eZl} = \int_{F_0^{-1}(1-\alpha_z)}^{+\infty} f(t; \hat{v}_{eZl}) dt$. For the model from (5), we simply calculate the estimated NCP of the standard linear regression, which is

$$\hat{v}_{eXlk} = \frac{\sum_{i=1}^n (X_{ilk} - \bar{X}_{\cdot lk}) Z_{il}}{\sqrt{\sum_{i=1}^n (X_{ilk} - \bar{X}_{\cdot lk})^2 \hat{\sigma}_{el}}}, \quad (13)$$

where

$$\hat{\sigma}_{el} = \frac{1}{n-2} \sum_{i=1}^n \left(Z_{il} - \bar{Z}_{\cdot l} + \hat{\beta}_{lk} (X_{ilk} - \bar{X}_{\cdot lk}) \right)^2. \quad (14)$$

Again, we use the significance level $\alpha_{el} = \frac{0.05}{n_{el}}$, where n_{el} is the total number of local genetic variants in the l^{th} gene, so that the power of detecting X_{lk} is estimated by $P_{eXlk} = \int_{F_0^{-1}(1-\alpha_{el})}^{+\infty} f(t; \hat{v}_{eXlk}) dt$. Since we assume the power of (4) and (5) are independent, the power of detecting the l^{th} gene and the k^{th} variants in the l^{th} gene simultaneously is given by $P_{eZl} P_{eXlk}$. If the independence assumption is violated, i.e., the powers of these two steps are positively correlated, then the estimated power for emGWAS will be conservative.

For TWAS, the NCP is estimated in a similar manner as the first step of emGWAS, i.e.

$$\hat{v}_{Tl} = \frac{\sum_{i=1}^n \hat{Z}_{il}^* \hat{Y}_i^* \sum_{i=1}^n \hat{D}_{xi}^2 - \sum_{i=1}^n \hat{Y}_i^* \hat{D}_{xi} \cdot \sum_{i=1}^n \hat{Z}_{il}^* \hat{D}_{xi}}{\sqrt{\sum_{i=1}^n (\hat{Z}_{il}^*)^2 (\sum_{i=1}^n \hat{D}_{xi}^2)^2 - (\sum_{i=1}^n \hat{D}_{xi} \cdot \hat{Z}_{il}^*)^2 \sum_{i=1}^n \hat{D}_{xi}^2}}, \quad (15)$$

where the only difference between (12) and (15) is that $\hat{Z}_{il}^* = \hat{D}_x Z_{il}$ in (15) is replaced by $\hat{Z}_{il}^* = \hat{D}_x \hat{Z}_{il}$ in (15). The significance level is again $\alpha_z = \frac{0.05}{n_z}$ and the power is estimated by $P_{Tl} =$

$$\int_{F_0^{-1}(1-\alpha_z)}^{+\infty} f(t; \hat{v}_{Tl}) dt.$$

268

269 Simulation of phenotype and expression

270 As the statistical power of each protocol depends on the magnitude of the genetic effect,
271 we simulated input data at various effect sizes. While effect size depends on a combination of

272 many factors, we chose to focus on the following three aspects. 1) We considered two genetic
273 architectures: causality and pleiotropy (**Fig 1**). In the causality scenario, the contribution of
274 genotype to phenotype is mediated through expression (**Fig 1a**), whereas in the pleiotropy
275 scenario, genotype contributes to both expression and phenotype directly (**Fig 1b**). We did not
276 consider the scenario where phenotype is causal to expression. 2) We considered the strength of
277 three different variant components: trait heritability (the variance component of phenotype
278 explained by genotype, denoted $h_{x \Rightarrow y}^2$), expression heritability (the variance component of
279 expression explained by genotype, denoted $h_{x \Rightarrow z}^2$), and the phenotypic variance component
280 explained by expression, denoted $h_{z \Rightarrow y}^2$ and abbreviated as *PVX*. 3) We also considered the
281 number of genes contributing to phenotype and the number of local genetic variants contributing
282 to expression.

283 In all our simulations, we use real genotypes from the 1000 Genomes Project ($N = 2504$).
284 Although there are multiple existing datasets containing both expressions and genotype, we
285 chose to use simulated expressions instead as it is difficult to match real data exactly to desired
286 properties such as expression heritability or the number of contributing genetic variants. By
287 simulating expressions, we can perform a consistent power analysis across a comprehensive
288 range of prespecified input conditions.

289 In the causality scenario, phenotypes were simulated with the following procedure. First,
290 several genes ($n_{z-sig} = 4, 9, \text{ or } 13$) were selected as causal genes. For each gene (indexed by $l =$
291 $1, 2, \dots, n_{z-sig}$), several common and independent genetic variants were selected as causal
292 variants ($n_{z(l)-sig} = 4 \sim 9$, $MAF > 0.05$, and $R^2 < 0.01$). A linear combination of local variants
293 in the l^{th} gene is generated to produce the expression values $Z_{(l)}$, and a linear combination of
294 these gene expressions \mathbf{Z} is generated as the genomic contribution to phenotype. Note that at
295 each step, we ensure the simulated linear combinations of variants and expressions match our
296 desired values for expression heritability $h_{x \Rightarrow z}^2$ and *PVX* $h_{z \Rightarrow y}^2$ (**S5 Appendix**).

297 In the pleiotropy scenario, we followed a similar procedure except that the phenotype Y
298 was directly generated from a linear combination of genotypes, instead of expressions (**S6**
299 **Appendix**). Note that although the expressions Z and phenotype Y are unrelated by genuine
300 biological causality, they are generated from the same genetic variants and are therefore
301 statistically correlated. Therefore, if the trait heritability and expression heritability are

302 sufficiently large, TWAS can still identify causal genes using the statistical correlation between
 303 genetic variants and expression.

304 We simulated both scenarios with expression heritability $h_{x \Rightarrow z}^2$ from the values (2.5%,
 305 3%, 4%, 6%, 8%, 10%, 30%), and with trait heritability $h_{x \Rightarrow y}^2$ in the pleiotropy scenario or PVX
 306 $h_{z \Rightarrow y}^2$ in the causality scenario from the values (0.5%, 1%, 2.5%, 5%, 10%). Although we
 307 initially tested more extreme values, our **Results** show that the turning points where TWAS
 308 outperforms GWAS are well within the range of values presented here, and the relative
 309 performance of the protocols remains consistent under more extreme conditions. We therefore
 310 chose to restrict our discussion to the most relevant values for protocol selection, noting that the
 311 expression heritability values we examine are at the high-end of real observed values [18], while
 312 the trait heritability values are lower than typically found in GWAS.

313 Finally, as each simulation involves multiple variants and genes, the overall power of
 314 each protocol is defined as follows: the power of GWAS is the probability of detecting at least
 315 one causal variant in any causal gene, the power of emGWAS is the probability of detecting at
 316 least one gene and one local SNP of that gene simultaneously, and the power of TWAS is the
 317 probability that at least one predicted gene expression is significant. Specifically,

$$318 \quad P_{GWAS} = 1 - \prod_{j=1}^{n_{x-sig}} (1 - P_{G(j)}), \quad (16)$$

$$319 \quad P_{emGWAS} = 1 - \prod_{l=1}^{n_{z-sig}} (1 - P_{eZ(l)} P_{eX(l)}), \text{ where } P_{eX(l)} = 1 - \prod_{k=1}^{n_{z(l)-sig}} (1 - P_{eX(l)(k)}), \quad (17)$$

$$320 \quad P_{TWAS} = 1 - \prod_{l=1}^{n_{z-sig}} (1 - P_{T(l)}), \quad (18)$$

321 where n_{x-sig} , n_{z-sig} and $n_{z(l)-sig}$ denote the numbers of significant SNPs, genes, and SNPs in
 322 the l^{th} significant gene respectively, $G(j)$ denotes the j^{th} significant SNP identified by GWAS,
 323 $Z(l)$ and $X(l)(k)$ denote the l^{th} significant gene and the k^{th} significant SNP of the l^{th}
 324 significant gene identified by emGWAS, and $T(l)$ denotes the l^{th} significant gene identified by
 325 TWAS.

326

327

328 Results

329 As a quality control measure, we compared the actual expression heritability and the
330 mean R^2 of the predicted expressions (**Table 1**). As expected, the mean R^2 grows closer to the
331 actual heritability value as expression heritability increases.

332
333 **Table 1: Comparisons of R^2 of imputed gene expression under different levels of expression**
334 **heritability and number of genetic variants.**

335

	Mean of R^2	Sample Standard Deviation of R^2
$h_1^2 = 0.025$	0.007847616	0.007415877
$h_1^2 = 0.03$	0.01259302	0.008410582
$h_1^2 = 0.04$	0.02319834	0.009481371
$h_1^2 = 0.06$	0.04415579	0.01083593
$h_1^2 = 0.08$	0.06465895	0.01175991
$h_1^2 = 0.1$	0.08518152	0.01264175
$h_1^2 = 0.3$	0.2886779	0.01514781

336

337 Causality scenario

338 We first analyzed cases where expression heritability is high ($h_{x \rightarrow z}^2 = 0.1$ or 0.3) but the
339 PVX is low (**Fig 2**). Overall, emGWAS clearly outperforms both GWAS and TWAS by a large
340 margin, and TWAS also generally outperforms GWAS. Note that although the PVX is low and
341 favors GWAS, TWAS is still more powerful due to the high expression heritability, which shows
342 that expression heritability affects the performance of TWAS more than the PVX. Consistent
343 with intuition, we observed that GWAS and TWAS have higher power as expression heritability
344 increases, whereas this increase is much smaller for emGWAS. The power of GWAS and
345 emGWAS reduces as the number of causal genes grows, whereas TWAS is largely unaffected by
346 the number of causal genes. This is also consistent with intuition since TWAS uses GR \hat{X} (\hat{Z}) to
347 aggregate genetic effects, avoiding the burden of multiple-testing correction.

348 We then analyzed cases where the PVX is high, but expression heritability is relatively
349 low ($h_{x \Rightarrow z}^2 = 0.025, 0.03, 0.04$ or 0.08). Evidently, emGWAS performs best with powers
350 consistently at 1.0. The comparison between TWAS and GWAS is more nuanced, as TWAS is
351 suboptimal to GWAS when the expression heritability is 0.025 or 0.03 (**Figs 3a and 3b**), begins
352 to outperform GWAS when the expression heritability is 0.04 (**Fig 3c**), and clearly outperforms
353 GWAS when the expression heritability is 0.08 (**Fig 3d**). This quantifies an important turning
354 point in that GWAS is superior to TWAS when expression heritability is less than 0.04, even if
355 PVX is high (favoring TWAS).

356

357 **Pleiotropy scenario**

358 Again, we first analyze cases where expression heritability is high and trait heritability is
359 low (**Fig 4**). Unlike in the causality scenario, the power of emGWAS is very low compared to
360 TWAS and GWAS. A potential explanation is that when the effect of genetic variants on
361 phenotype is not mediated through expressions, the non-genetic effects within the actual
362 expressions add noise to emGWAS predictions. In contrast, the elastic-net model in TWAS
363 captures only the genetic component of expressions, meaning the predicted expressions are a
364 more accurate model of the direct genetic effect on phenotype. While errors are unavoidable in
365 the elastic-net training process (as revealed in **Table 1**), our results show that the loss of power
366 due to non-genetic effects is overwhelmingly greater than the loss due to training errors. As in
367 the casualty scenario, TWAS generally outperforms GWAS except in the case where trait
368 heritability is extremely low and the number of contributing genes is large, which is rare in
369 practice. We therefore conclude that in both scenarios, TWAS has better power than GWAS
370 when expression heritability is high.

371 We finally analyze cases where expression heritability is low but trait heritability is high.
372 Here, emGWAS continues to be the least powerful of the three protocols. As in the causality
373 scenario, we again observe a turning point where TWAS outperforms GWAS: TWAS has lower
374 power than GWAS when the expression heritability is 0.025 or 0.04 (**Figs 5a and 5b**), TWAS
375 has comparable power when the expression heritability is 0.06 (**Fig 5c**), and TWAS outperforms
376 GWAS when the expression heritability is 0.08 (**Fig 5d**).

377 Our results can be summarized in two observations (**Fig 6**). First, emGWAS outperforms
378 TWAS and GWAS in the casualty scenario, but is less powerful in the pleiotropy scenario

379 regardless of the accuracy of the predicted expressions (**Table 1**). This demonstrates that when
380 non-genetic components in expression do not contribute to phenotype (i.e. pleiotropy scenario),
381 predicted expressions capture genetic contributions better than actual expressions (which include
382 non-genetic components). Second, the turning point at which traditional GWAS outperforms
383 TWAS is an expression heritability of less than 0.04 in the causality scenario, or 0.06 in the
384 pleiotropy scenario.

385 These turning points are immediately relevant to the practical conduct of association
386 mapping studies, as shown by the following analysis of expression heritability in existing TWAS
387 publications. As few publications disclose their estimated expression heritability, we use
388 published R^2 values of the correlation between predicted and actual expressions to approximate
389 the underlying expression heritability. We use the difference between expression heritability and
390 R^2 as calculated from our simulations (**Table 1**) to map these R^2 values to an estimated
391 expression heritability (i.e. R^2 of 0.023 and 0.044 give expression heritability values 0.04 and
392 0.06, respectively), although in practice the true difference may vary depending on the predictive
393 model used in each study. Table 1 of the PrediXcan publication lists significant results from their
394 paper, in which 14 out of 41 discovered genes have R^2 values less than 0.044, with 2 values less
395 than 0.023. Additionally, our review of recent TWAS publications shows that most of the genes
396 presented have mean R^2 values less than 0.044 or 0.023 (**Table 2**). As our power analysis
397 indicated, GWAS may have better power than TWAS given these low expression heritability
398 conditions. Although we are unable to determine if the genes discovered by these publications
399 follow the causality or pleiotropy scenario, other advanced statistical models [44] may be used to
400 determine appropriate thresholds to distinguish between pleiotropy and causality.

401 In summary, we suggest the following modifications to the TWAS protocol. First, one
402 may estimate expression heritability in the reference panel and filter out genes with expression
403 heritability less than 0.04. Second, after conducting TWAS association mapping, determine the
404 underlying causality scenario (causality or pleiotropy) in order to choose an appropriate
405 expression heritability threshold (0.04 or 0.06). Finally, conduct GWAS for each gene with an
406 expression heritability below the given threshold.

407

408

409 **Table 2: Mean R^2 in published TWAS projects.**

Title of the publication	Description of prediction accuracy
Large-scale transcriptome-wide association study identifies new prostate cancer risk regions [22]	The mean $R^2 = 0.07$ for measured and predicted gene expression for TCGA normal prostate samples using models fitted in GTEx normal prostate.
A framework for transcriptome-wide association studies in breast cancer in diverse study populations [45]	The median CV R^2 for the 153 genes is 0.011 in both African American and white women.
Evaluation of PrediXcan for prioritizing GWAS associations and predicting gene Expression [46]	The average of prediction accuracy (R^2) is 0.023 for the DGN model and 0.02 for the GTEx model, with both using whole blood model.
A gene-based association method for mapping traits using reference transcriptome data [18]	The average prediction R^2 value is 0.0197 for GEUVADIS LCLs. For GTEx tissues, the prediction R^2 values are 0.0367 (adipose), 0.0358 (tibial artery), 0.0356 (left-ventricular heart), 0.0359 (lung), 0.0269 (muscle), 0.0422 (tibial nerve), 0.0374 (sun-exposed skin), 0.0398 (thyroid) and 0.0458 (whole blood).

410

411 **Application to the power estimation of EpiXcan**

412 Our NCP-based framework can be applied to estimate the power of other protocols. To
413 demonstrate this point, we estimated the power of EpiXcan [27], a novel TWAS-like protocol
414 integrating epigenetic functional annotations to improve the accuracy of predicted expressions
415 and therefore overall TWAS power. The original EpiXcan paper demonstrated that (1) the
416 predictive accuracy of expressions is significantly increased, and (2) EpiXcan enabled the
417 discovery of novel genes [27]. We present here the first rigorous power analysis of EpiXcan. We
418 first conduct simulations where a subset of SNPs are assigned increased effects, which reflects
419 the main insight of the EpiXcan paper that epigenetic-relevant functional SNPs have higher
420 impact on variation in gene expression. In particular, we assume the real effect size follows a

421 standard normal distribution $N(0,1)$, and sample effect sizes from this distribution. Assuming
422 these functional SNPs are known (based on various techniques of annotating SNP functions), we
423 relieve their penalty in training the predictive model. Using the predicted expressions, we
424 calculate power using our derived NCP, and compare the resulting analysis with the standard
425 TWAS protocol. **Supplementary Fig. S1-S4** depict this quantitative evaluation of the
426 improvement in power due to the contribution of epigenetic-relevant functional SNPs. Evidently,
427 under the causality model EpiXcan indeed increases power by improving expression predictions
428 (**Supplementary Fig. S1, S2**). However, under the pleiotropy model, EpiXcan only shows a very
429 small increase in power over TWAS (**Supplementary Fig. S3, S4**). This observation suggests
430 that when DNA mutations contribute to phenotype directly, the benefit of more accurate
431 predictions for expressions may not be substantial.

432

433 **Discussion**

434 In this work, we produced a novel derivation of the NCP for LMMs based on the
435 decorrelation procedure, allowing us to calculate closed-form estimates of statistical power for
436 three protocols: GWAS, emGWAS, and TWAS. Our power analysis revealed two practical
437 insights. **First**, in the pleiotropy scenario, the use of predicted expressions in TWAS is
438 overwhelmingly more powerful than the use of actual expressions in emGWAS, regardless of the
439 accuracy of the predicted expressions *per se* (**Table 1**). This suggests that even if real
440 expressions can be experimentally determined, TWAS is still superior for the analysis of some
441 genes. While this appears counterintuitive, in statistical terms it is a direct result of the lack of a
442 causal relationship between expression and phenotype under pleiotropy. This result reinforces
443 the key insight, as presented by some publications [18], that TWAS uses expression as an
444 objective function to select a linear combination of genetic variants, rather than attempting to
445 accurately predict expressions. We note that this is equivalent to denoising in the field of
446 machine learning [47]. **Second**, expression heritability determines the relative power of TWAS
447 and GWAS. When the expression heritability is lower than 0.04 (in the casualty scenario) or 0.06
448 (in the pleiotropy scenario), GWAS outperforms TWAS despite not utilizing gene expression
449 information. This suggests that in practice, TWAS may often be suboptimal when expression
450 heritability is low (**Table 2 & Table 1 in [18]**), which can be mitigated by choosing the optimal
451 association mapping protocol according to this work's quantitative guidelines.

452 A recent publication has also compared the statistical powers of GWAS and TWAS using
453 pure simulations [33]. However, since we calculate power from a closed-form NCP derivation,
454 our work establishes theoretical benchmarks for the performance of each protocol, independent
455 of their implementations. Our work also has a different focus: rather than comparing techniques
456 for training the genotype-expression predictive model and the impact of the actual number of
457 causal genetic variants, we rank the effectiveness of GWAS, TWAS and emGWAS to better
458 guide the practical application of TWAS. We analyze the theoretical effectiveness of real
459 expressions as utilized by emGWAS, but exclude the protocol eGWAS as analyzed in [33],
460 which uses eQTLs to assist association mapping. Our conclusions also differ slightly, as while
461 the previous publication highlighted the importance of expression heritability, they concluded
462 that expression heritability affects power only under the causality scenario, and not pleiotropy. In
463 contrast, we concluded that expression heritability affects both scenarios.

464 Finally, our closed-form derivation is readily adaptable to other methods utilizing middle
465 ‘omics’ (endophenotypes) such as IWAS [28] and PWAS [29, 34]. In fact, the variable Z in
466 formula (15) can already represent such data as images or proteins, and thus no further
467 modifications of the NCPs are necessary to adapt this work.

468 The present NCP framework only focuses on statistical power for detecting associations,
469 and is not able to determine causality in the framework of Mendelian randomization such as in
470 SMR and its extensions [48, 49]. As a future work, we may attempt to derive closed-form power
471 analyses for the MR framework.

472 There are several limitations in the present study. Although our closed-form derivation is
473 easily adaptable and works independently of specific implementations, it is unable to capture
474 power loss due to implementation limitations or bias in specific datasets. Additionally, closed-
475 form derivations are more sensitive to model assumptions than simulation-based methods. Our
476 calculation of the NCP also requires the variance component σ_g^2 to be estimated from data, in
477 order to form the decorrelation matrix D_x . Although this approximation introduces extra
478 variability and may therefore cause a decrease in power, we have omitted this variability from
479 our analyses as the estimation of σ_g^2 is generally well-established, and has high accuracy in
480 practice when given thousands of samples. Finally, we only compared linear models for GWAS
481 and TWAS. As a future work, we may explore kernel-based nonparametric and semiparametric
482 methods for conducting both GWAS [50, 51] and TWAS [52].

483

484

485 **Acknowledgements**

486 J.W. is supported by an NSERC Discovery Grant (RGPIN-2018-04328). Q.L. is supported by an
487 NSERC Discovery Grant (RGPIN-2017-04860), a Canada Foundation for Innovation JELF grant
488 (36605), a New Frontiers in Research Fund (NFRFE-2018-00748), a Clinical Research Fund and
489 a Startup grant supported by Alberta Children's Hospital Research Institute (ACHRI). C.C. is
490 supported by ACHRI postdoctoral scholarship.

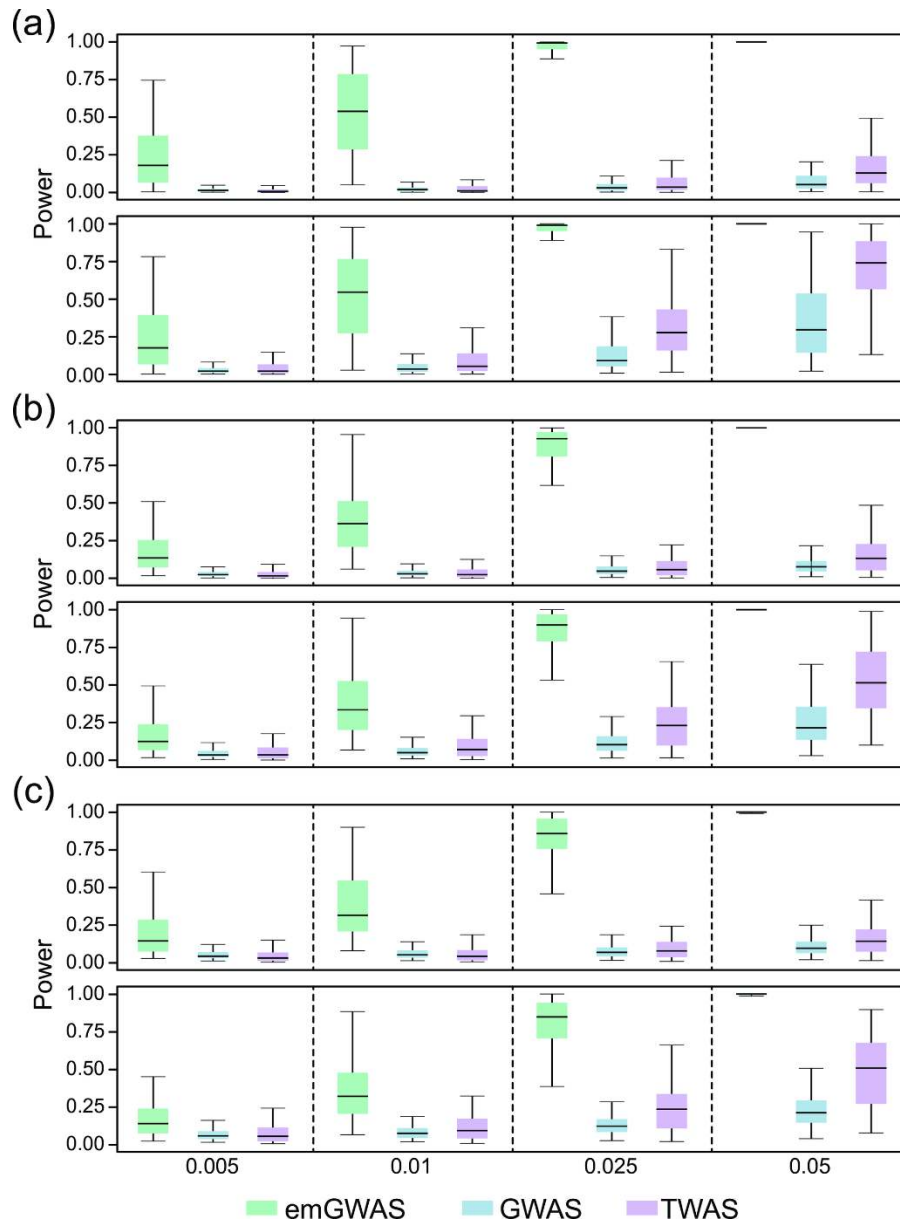
491 **References**

- 492 1. Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al. Transcript
493 assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching
494 during cell differentiation. *Nat Biotechnol.* 2010;28(5):511-5. Epub 2010/05/04. doi:
495 10.1038/nbt.1621. PubMed PMID: 20436464; PubMed Central PMCID: PMC3146043.
- 496 2. Wang Z, Gerstein M, Snyder M. RNA-Seq: a revolutionary tool for transcriptomics. *Nat*
497 *Rev Genet.* 2009;10(1):57-63. Epub 2008/11/19. doi: 10.1038/nrg2484. PubMed PMID:
498 19015660; PubMed Central PMCID: PMC3146043.
- 499 3. McGettigan PA. Transcriptomics in the RNA-seq era. *Curr Opin Chem Biol.*
500 2013;17(1):4-11. Epub 2013/01/08. doi: 10.1016/j.cbpa.2012.12.008. PubMed PMID: 23290152.
- 501 4. Ozsolak F, Milos PM. RNA sequencing: advances, challenges and opportunities. *Nat*
502 *Rev Genet.* 2011;12(2):87-98. Epub 2010/12/31. doi: 10.1038/nrg2934. PubMed PMID:
503 21191423; PubMed Central PMCID: PMC3146043.
- 504 5. Selevsek N, Chang CY, Gillet LC, Navarro P, Bernhardt OM, Reiter L, et al.
505 Reproducible and consistent quantification of the *Saccharomyces cerevisiae* proteome by
506 SWATH-mass spectrometry. *Mol Cell Proteomics.* 2015;14(3):739-49. Epub 2015/01/07. doi:
507 10.1074/mcp.M113.035550. PubMed PMID: 25561506; PubMed Central PMCID:
508 PMC3146043.
- 509 6. Pible O, Armengaud J. Improving the quality of genome, protein sequence, and
510 taxonomy databases: a prerequisite for microbiome meta-omics 2.0. *Proteomics.*
511 2015;15(20):3418-23. Epub 2015/06/04. doi: 10.1002/pmic.201500104. PubMed PMID:
512 26038180.
- 513 7. Bell AW, Deutsch EW, Au CE, Kearney RE, Beavis R, Sechi S, et al. A HUPO test
514 sample study reveals common problems in mass spectrometry-based proteomics. *Nat Methods.*
515 2009;6(6):423-30. Epub 2009/05/19. doi: 10.1038/nmeth.1333. PubMed PMID: 19448641;
516 PubMed Central PMCID: PMC3146043.
- 517 8. Zhang A, Sun H, Wang P, Han Y, Wang X. Modern analytical techniques in
518 metabolomics analysis. *Analyst.* 2012;137(2):293-300. Epub 2011/11/22. doi:
519 10.1039/c1an15605e. PubMed PMID: 22102985.
- 520 9. Coats VC, Rumpho ME. The rhizosphere microbiota of plant invaders: an overview of
521 recent advances in the microbiomics of invasive plants. *Front Microbiol.* 2014;5:368. Epub
522 2014/08/08. doi: 10.3389/fmicb.2014.00368. PubMed PMID: 25101069; PubMed Central
523 PMCID: PMC3146043.
- 524 10. Teperino R, Lempradl A, Pospisilik JA. Bridging epigenomics and complex disease: the
525 basics. *Cell Mol Life Sci.* 2013;70(9):1609-21. Epub 2013/03/07. doi: 10.1007/s00018-013-
526 1299-z. PubMed PMID: 23463237.
- 527 11. Klein RJ, Zeiss C, Chew EY, Tsai JY, Sackler RS, Haynes C, et al. Complement factor H
528 polymorphism in age-related macular degeneration. *Science.* 2005;308(5720):385-9. doi:
529 10.1126/science.1109557. PubMed PMID: 15761122; PubMed Central PMCID: PMC1512523.
- 530 12. Ozaki K, Ohnishi Y, Iida A, Sekine A, Yamada R, Tsunoda T, et al. Functional SNPs in
531 the lymphotoxin-alpha gene that are associated with susceptibility to myocardial infarction. *Nat*
532 *Genet.* 2002;32(4):650-4. Epub 2002/11/12. doi: 10.1038/ng1047. PubMed PMID: 12426569.
- 533 13. Mills MC, Rahal C. A scientometric review of genome-wide association studies.
534 *Commun Biol.* 2019;2:9. Epub 2019/01/10. doi: 10.1038/s42003-018-0261-x. PubMed PMID:
535 30623105; PubMed Central PMCID: PMC3146043.
- 536 14. Eddy S, Mariani LH, Kretzler M. Integrated multi-omics approaches to improve
537 classification of chronic kidney disease. *Nat Rev Nephrol.* 2020. Epub 2020/05/20. doi:
538 10.1038/s41581-020-0286-5. PubMed PMID: 32424281.

- 539 15. Hasin Y, Seldin M, Lusia A. Multi-omics approaches to disease. *Genome Biol.*
540 2017;18(1):83. Epub 2017/05/10. doi: 10.1186/s13059-017-1215-1. PubMed PMID: 28476144;
541 PubMed Central PMCID: PMC5418815.
- 542 16. Yan J, Risacher SL, Shen L, Saykin AJ. Network approaches to systems biology
543 analysis of complex disease: integrative methods for multi-omics data. *Brief Bioinform.*
544 2018;19(6):1370-81. Epub 2017/07/07. doi: 10.1093/bib/bbx066. PubMed PMID: 28679163;
545 PubMed Central PMCID: PMC6454489.
- 546 17. Fukushima A, Kusano M, Redestig H, Arita M, Saito K. Integrated omics approaches in
547 plant systems biology. *Curr Opin Chem Biol.* 2009;13(5-6):532-8. Epub 2009/10/20. doi:
548 10.1016/j.cbpa.2009.09.022. PubMed PMID: 19837627.
- 549 18. Gamazon ER, Wheeler HE, Shah KP, Mozaffari SV, Aquino-Michaels K, Carroll RJ, et
550 al. A gene-based association method for mapping traits using reference transcriptome data.
551 *Nature genetics.* 2015;47(9):1091-8. doi: 10.1038/ng.3367. PubMed PMID: 26258848; PubMed
552 Central PMCID: PMC4552594.
- 553 19. Gusev A, Ko A, Shi H, Bhatia G, Chung W, Penninx BW, et al. Integrative approaches
554 for large-scale transcriptome-wide association studies. *Nature genetics.* 2016;48(3):245-52. doi:
555 10.1038/ng.3506. PubMed PMID: 26854917; PubMed Central PMCID: PMC4767558.
- 556 20. Barbeira AN, Dickinson SP, Bonazzola R, Zheng J, Wheeler HE, Torres JM, et al.
557 Exploring the phenotypic consequences of tissue specific gene expression variation inferred
558 from GWAS summary statistics. *Nat Commun.* 2018;9(1):1825. Epub 2018/05/10. doi:
559 10.1038/s41467-018-03621-1. PubMed PMID: 29739930; PubMed Central PMCID:
560 PMC5940825.
- 561 21. Gusev A, Mancuso N, Won H, Kousi M, Finucane HK, Reshef Y, et al. Transcriptome-
562 wide association study of schizophrenia and chromatin activity yields mechanistic disease
563 insights. *Nat Genet.* 2018;50(4):538-48. Epub 2018/04/11. doi: 10.1038/s41588-018-0092-1.
564 PubMed PMID: 29632383; PubMed Central PMCID: PMC5942893.
- 565 22. Mancuso N, Gayther S, Gusev A, Zheng W, Penney KL, Kote-Jarai Z, et al. Large-scale
566 transcriptome-wide association study identifies new prostate cancer risk regions. *Nat Commun.*
567 2018;9(1):4079. Epub 2018/10/06. doi: 10.1038/s41467-018-06302-1. PubMed PMID:
568 30287866; PubMed Central PMCID: PMC6172280.
- 569 23. Theriault S, Gaudreault N, Lamontagne M, Rosa M, Boulanger MC, Messika-Zeitoun D,
570 et al. A transcriptome-wide association study identifies PALMD as a susceptibility gene for
571 calcific aortic valve stenosis. *Nat Commun.* 2018;9(1):988. Epub 2018/03/08. doi:
572 10.1038/s41467-018-03260-6. PubMed PMID: 29511167; PubMed Central PMCID:
573 PMC5840407.
- 574 24. Gong L, Zhang D, Lei Y, Qian Y, Tan X, Han S. Transcriptome-wide association study
575 identifies multiple genes and pathways associated with pancreatic cancer. *Cancer Med.*
576 2018;7(11):5727-32. Epub 2018/10/20. doi: 10.1002/cam4.1836. PubMed PMID: 30334361;
577 PubMed Central PMCID: PMC6247024.
- 578 25. Ratnapriya R, Sosina OA, Starostik MR, Kwicklis M, Kappahn RJ, Fritsche LG, et al.
579 Retinal transcriptome and eQTL analyses identify genes associated with age-related macular
580 degeneration. *Nat Genet.* 2019;51(4):606-10. Epub 2019/02/12. doi: 10.1038/s41588-019-0351-
581 9. PubMed PMID: 30742112; PubMed Central PMCID: PMC6441365.
- 582 26. Atkins I, Kinnersley B, Ostrom QT, Labreche K, Ilyasova D, Armstrong GN, et al.
583 Transcriptome-Wide Association Study Identifies New Candidate Susceptibility Genes for
584 Glioma. *Cancer Res.* 2019;79(8):2065-71. Epub 2019/02/03. doi: 10.1158/0008-5472.CAN-18-
585 2888. PubMed PMID: 30709929; PubMed Central PMCID: PMC6522343.
- 586 27. Zhang W, Voloudakis G, Rajagopal VM, Readhead B, Dudley JT, Schadt EE, et al.
587 Integrative transcriptome imputation reveals tissue-specific and shared biological mechanisms
588 mediating susceptibility to complex traits. *Nat Commun.* 2019;10(1):3834. Epub 2019/08/25. doi:

- 589 10.1038/s41467-019-11874-7. PubMed PMID: 31444360; PubMed Central PMCID:
590 PMCPMC6707297.
- 591 28. Xu Z, Wu C, Pan W, Alzheimer's Disease Neuroimaging I. Imaging-wide association
592 study: Integrating imaging endophenotypes in GWAS. *Neuroimage*. 2017;159:159-69. Epub
593 2017/07/25. doi: 10.1016/j.neuroimage.2017.07.036. PubMed PMID: 28736311; PubMed
594 Central PMCID: PMCPMC5671364.
- 595 29. Brandes N, Linal N, Linal M, editors. PWAS: Proteome-Wide Association Study2020;
596 Cham: Springer International Publishing.
- 597 30. Mancuso N, Shi H, Goddard P, Kichaev G, Gusev A, Pasaniuc B. Integrating Gene
598 Expression with Summary Association Statistics to Identify Genes Associated with 30 Complex
599 Traits. *Am J Hum Genet*. 2017;100(3):473-87. Epub 2017/02/28. doi:
600 10.1016/j.ajhg.2017.01.031. PubMed PMID: 28238358; PubMed Central PMCID:
601 PMCPMC5339290.
- 602 31. Wainberg M, Sinnott-Armstrong N, Mancuso N, Barbeira AN, Knowles DA, Golan D, et
603 al. Opportunities and challenges for transcriptome-wide association studies. *Nat Genet*.
604 2019;51(4):592-9. Epub 2019/03/31. doi: 10.1038/s41588-019-0385-z. PubMed PMID:
605 30926968; PubMed Central PMCID: PMCPMC6777347.
- 606 32. Mogil LS, Andaleon A, Badalamenti A, Dickinson SP, Guo X, Rotter JI, et al. Genetic
607 architecture of gene expression traits across diverse populations. *PLoS Genet*.
608 2018;14(8):e1007586. Epub 2018/08/11. doi: 10.1371/journal.pgen.1007586. PubMed PMID:
609 30096133; PubMed Central PMCID: PMCPMC6105030.
- 610 33. Veturi Y, Ritchie MD. How powerful are summary-based methods for identifying
611 expression-trait associations under different genetic architectures? *Pac Symp Biocomput*.
612 2018;23:228-39. Epub 2017/12/09. PubMed PMID: 29218884; PubMed Central PMCID:
613 PMCPMC5785784.
- 614 34. Okada H, Ebhardt HA, Vonesch SC, Aebersold R, Hafen E. Proteome-wide association
615 studies identify biochemical modules associated with a wing-size phenotype in *Drosophila*
616 *melanogaster*. *Nat Commun*. 2016;7:12649. Epub 2016/09/02. doi: 10.1038/ncomms12649.
617 PubMed PMID: 27582081; PubMed Central PMCID: PMCPMC5025782.
- 618 35. Kang HM, Sul JH, Service SK, Zaitlen NA, Kong SY, Freimer NB, et al. Variance
619 component model to account for sample structure in genome-wide association studies. *Nat*
620 *Genet*. 2010;42(4):348-54. Epub 2010/03/09. doi: 10.1038/ng.548. PubMed PMID: 20208533;
621 PubMed Central PMCID: PMCPMC3092069.
- 622 36. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al. PLINK: a
623 tool set for whole-genome association and population-based linkage analyses. *American journal*
624 *of human genetics*. 2007;81(3):559-75. doi: 10.1086/519795. PubMed PMID: 17701901;
625 PubMed Central PMCID: PMC1950838.
- 626 37. Aulchenko YS, Ripke S, Isaacs A, van Duijn CM. GenABEL: an R library for genome-
627 wide association analysis. *Bioinformatics*. 2007;23(10):1294-6. Epub 2007/03/27. doi:
628 10.1093/bioinformatics/btm108. PubMed PMID: 17384015.
- 629 38. Gogarten SM, Bhangale T, Conomos MP, Laurie CA, McHugh CP, Painter I, et al.
630 GWASTools: an R/Bioconductor package for quality control and analysis of genome-wide
631 association studies. *Bioinformatics*. 2012;28(24):3329-31. Epub 2012/10/12. doi:
632 10.1093/bioinformatics/bts610. PubMed PMID: 23052040; PubMed Central PMCID:
633 PMCPMC3519456.
- 634 39. Zhou X, Stephens M. Genome-wide efficient mixed-model analysis for association
635 studies. *Nat Genet*. 2012;44(7):821-4. Epub 2012/06/19. doi: 10.1038/ng.2310. PubMed PMID:
636 22706312; PubMed Central PMCID: PMCPMC3386377.
- 637 40. Patterson N, Price AL, Reich D. Population structure and eigenanalysis. *PLoS Genet*.
638 2006;2(12):e190. Epub 2006/12/30. doi: 10.1371/journal.pgen.0020190. PubMed PMID:
639 17194218; PubMed Central PMCID: PMCPMC1713260.

- 640 41. Rakitsch B, Lippert C, Stegle O, Borgwardt K. A Lasso multi-marker mixed model for
641 association mapping with population structure correction. *Bioinformatics*. 2013;29(2):206-14.
642 Epub 2012/11/24. doi: 10.1093/bioinformatics/bts669. PubMed PMID: 23175758.
- 643 42. Long Q, Zhang Q, Vilhjalmsón BJ, Forai P, Seren U, Nordborg M. JAWAMix5: an out-
644 of-core HDF5-based java implementation of whole-genome association studies using mixed
645 models. *Bioinformatics*. 2013;29(9):1220-2. Epub 2013/03/13. doi:
646 10.1093/bioinformatics/btt122. PubMed PMID: 23479353.
- 647 43. Shaffer JP. Multiple hypothesis testing. *Annual review of psychology*. 1995;46(1):561-84.
- 648 44. Schadt EE, Lamb J, Yang X, Zhu J, Edwards S, Guhathakurta D, et al. An integrative
649 genomics approach to infer causal associations between gene expression and disease. *Nature*
650 *genetics*. 2005;37(7):710-7. doi: 10.1038/ng1589. PubMed PMID: 15965475; PubMed Central
651 PMCID: PMC2841396.
- 652 45. Bhattacharya A, Garcia-Closas M, Olshan AF, Perou CM, Troester MA, Love MI. A
653 framework for transcriptome-wide association studies in breast cancer in diverse study
654 populations. *Genome Biol*. 2020;21(1):42. Epub 2020/02/23. doi: 10.1186/s13059-020-1942-6.
655 PubMed PMID: 32079541; PubMed Central PMCID: PMC7033948.
- 656 46. Li B, Verma SS, Veturi YC, Verma A, Bradford Y, Haas DW, et al. Evaluation of
657 PrediXcan for prioritizing GWAS associations and predicting gene expression. *Pac Symp*
658 *Biocomput*. 2018;23:448-59. Epub 2017/12/09. PubMed PMID: 29218904; PubMed Central
659 PMCID: PMC5749400.
- 660 47. Tian C, Fei L, Zheng W, Xu Y, Zuo W, Lin C-W. Deep learning on image denoising: An
661 overview. *arXiv preprint arXiv:1912.13171*. 2019.
- 662 48. Zhu Z, Zhang F, Hu H, Bakshi A, Robinson MR, Powell JE, et al. Integration of summary
663 data from GWAS and eQTL studies predicts complex trait gene targets. *Nat Genet*.
664 2016;48(5):481-7. Epub 2016/03/29. doi: 10.1038/ng.3538. PubMed PMID: 27019110.
- 665 49. Hauberg ME, Zhang W, Giambartolomei C, Franzen O, Morris DL, Vyse TJ, et al. Large-
666 Scale Identification of Common Trait and Disease Variants Affecting Gene Expression. *Am J*
667 *Hum Genet*. 2017;101(1):157. Epub 2017/07/08. doi: 10.1016/j.ajhg.2017.06.003. PubMed
668 PMID: 28686855; PubMed Central PMCID: PMC5501865.
- 669 50. Wu MC, Lee S, Cai T, Li Y, Boehnke M, Lin X. Rare-variant association testing for
670 sequencing data with the sequence kernel association test. *Am J Hum Genet*. 2011;89(1):82-
671 93. Epub 2011/07/09. doi: 10.1016/j.ajhg.2011.05.029. PubMed PMID: 21737059; PubMed
672 Central PMCID: PMC3135811.
- 673 51. Wu MC, Kraft P, Epstein MP, Taylor DM, Chanock SJ, Hunter DJ, et al. Powerful SNP-
674 set analysis for case-control genome-wide association studies. *American journal of human*
675 *genetics*. 2010;86(6):929-42. doi: 10.1016/j.ajhg.2010.05.002. PubMed PMID: 20560208;
676 PubMed Central PMCID: PMC3032061.
- 677 52. Cao C, Kwok D, Edie S, Li Q, Ding B, Kossinna P, et al. kTWAS: integrating kernel-
678 machine with transcriptome-wide association studies improves statistical power and reveals
679 novel genes. *bioRxiv*. 2020.
- 680



684

685 **Fig 2: Causality scenario when expression heritability is high and PVX is low.**

686 The PVX is 0.005, 0.01, 0.025, and 0.05 in the four columns as indicated by the X-axis labels.

687 The number of genes contributing to phenotype for (a), (b) and (c) are 4, 9, and 13 respectively.

688 The expression heritability for the top and bottom rows of (a), (b) and (c) are 0.1 and 0.3

689 respectively. The number of causal variants per gene is randomly sampled from the interval

690 [4,9].

691

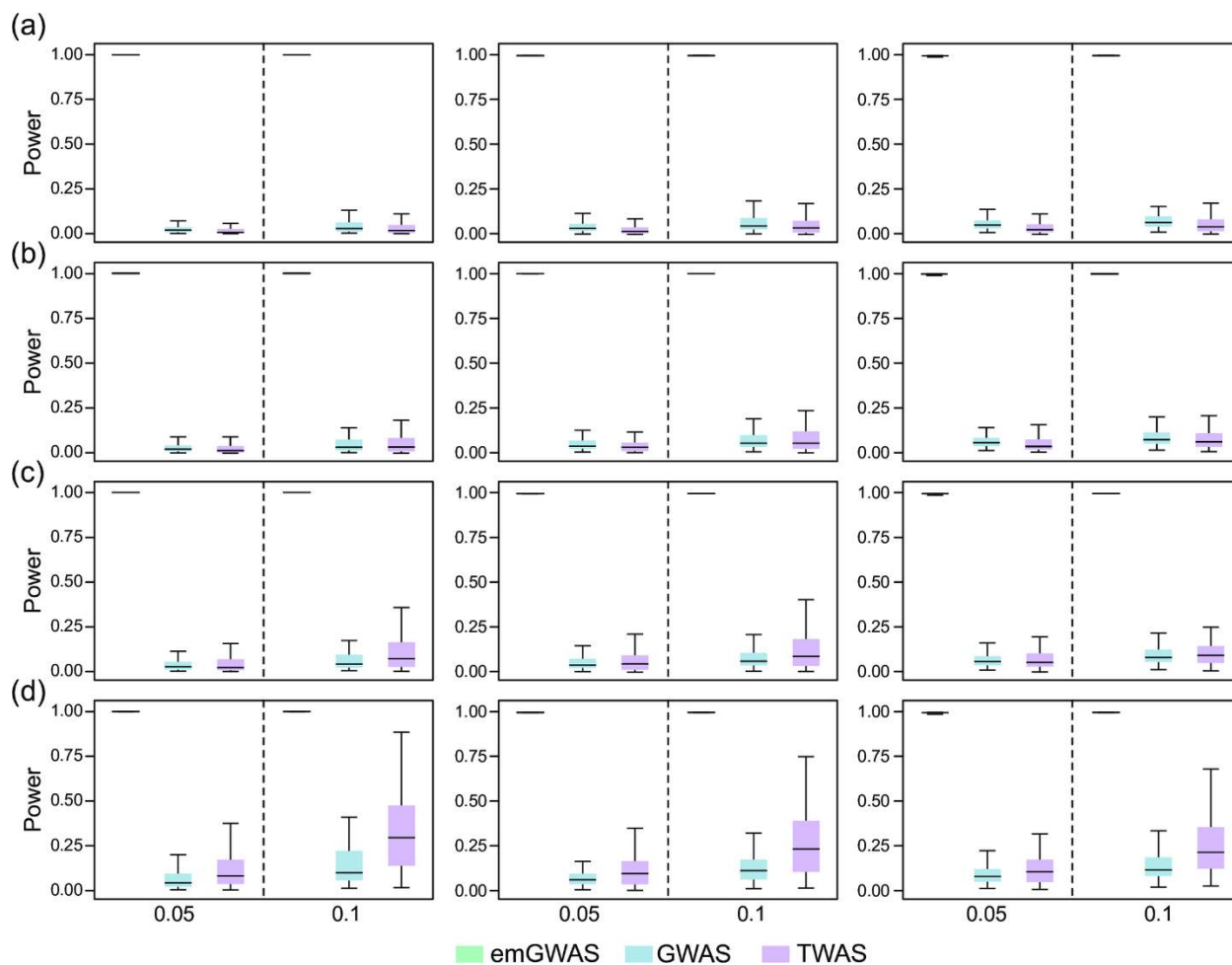
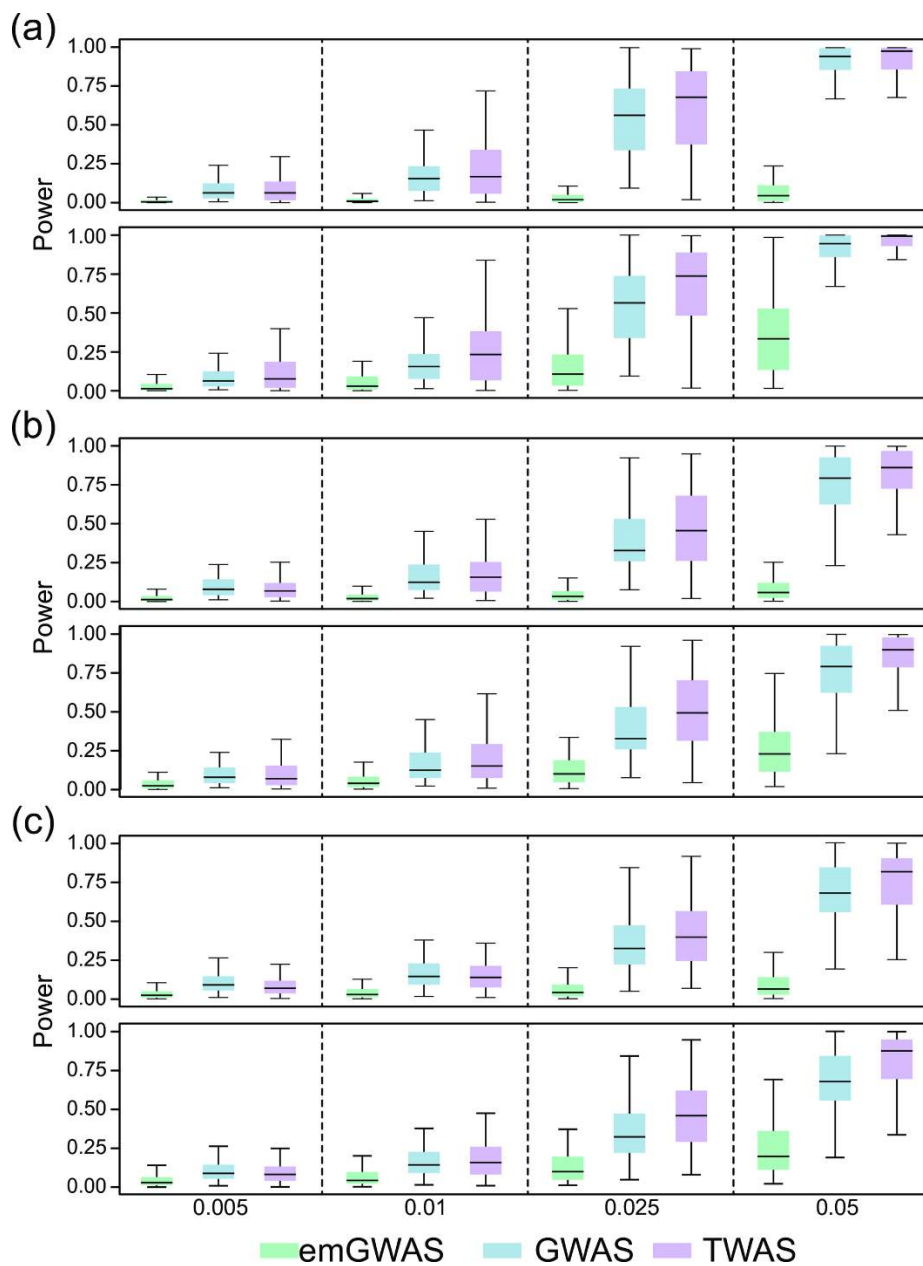


Fig 3: Causality scenario when expression heritability is low and PVX is high.

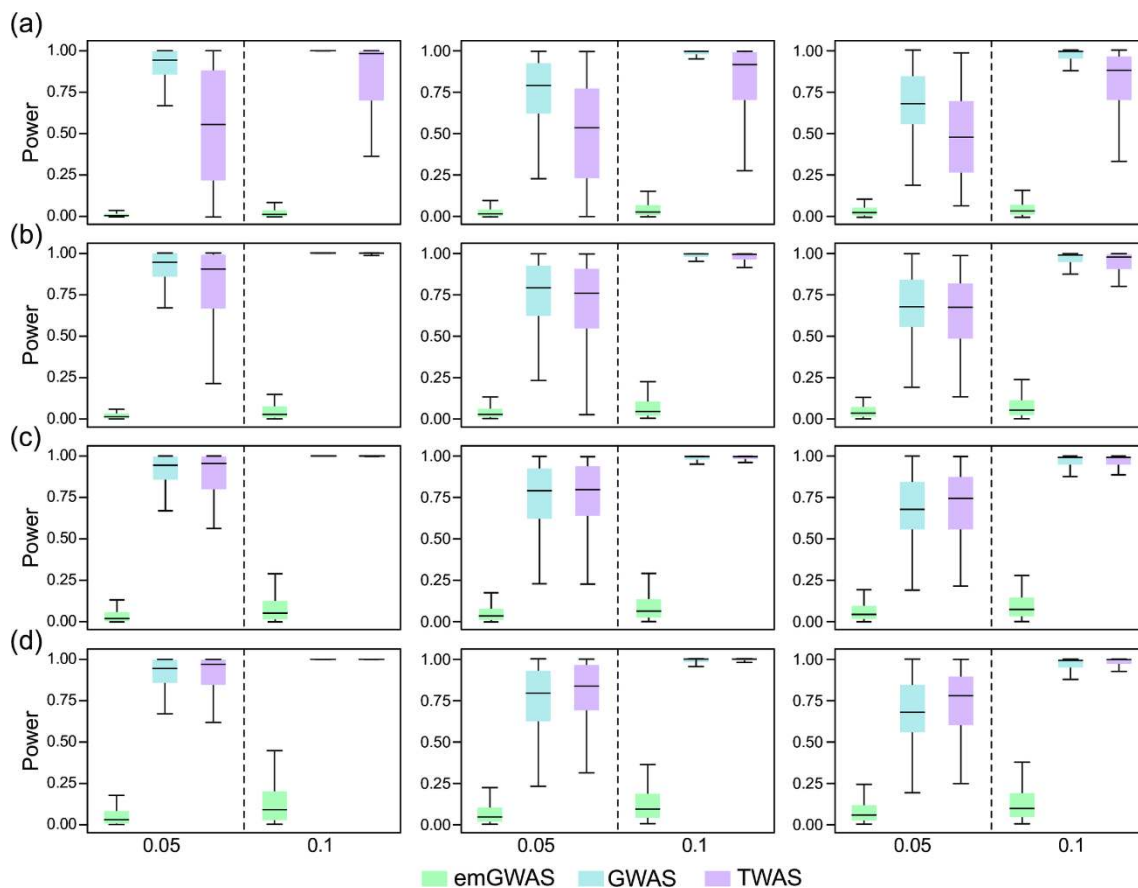
The PVX is 0.05 and 0.1 in the two columns as indicated by the X-axis labels. The numbers of genes contributing to phenotype in the left, middle and right panels are 4, 9, and 13 respectively. The expression heritability levels in (a), (b), (c) and (d) are 0.025, 0.03, 0.04, and 0.08 respectively. The number of causal variants per gene is randomly sampled from the interval [4,9].



699

700 **Fig 4: Pleiotropy scenario when expression heritability is high and trait heritability is low.**

701 The trait heritability is 0.005, 0.01, 0.025, and 0.05 in the four columns as indicated by the X-
702 axis labels. The numbers of genes contributing to phenotype for (a), (b) and (c) are 4, 9, and 13
703 respectively. The expression heritability for the top and bottom rows of (a), (b) and (c) are 0.1
704 and 0.3 respectively. The number of causal variants per gene is randomly sampled from the
705 interval [4,9].



706

707 **Fig 5: Pleiotropy scenario when expression heritability is low and trait heritability is high.**

708 The PVX is 0.05 and 0.1 in the two columns as indicated by the X-axis labels. The numbers of
709 genes contributing to phenotype for the left, middle and right panels are 4, 9, and 13 respectively.

710 The expression heritability levels in (a), (b), (c) and (d) are 0.025, 0.04, 0.06, and 0.08

711 respectively. The number of causal variants per gene is randomly sampled from the interval

712 [4,9].