1   FULL TITLE: Reference trait analysis reveals correlations between gene expression

2   and quantitative traits in disjoint samples

3

4   SHORT TITLE: Reference traits for systems genetics in disjoint samples

5

6   AUTHORS: Daniel A. Skelly[1], Narayanan Raghupathy[1], Raymond F. Robledo[1], Joel H.

7   Graber[1,2], Elissa J. Chesler[1,3]

8

9   AFFILIATIONS:

10  [1]The Jackson Laboratory, Bar Harbor, ME 04609, USA

11  [2]MDI Biological Laboratory, Bar Harbor, ME 04609, USA

12  [3]Corresponding author:

13  Elissa J. Chesler

14  The Jackson Laboratory

15  600 Main St.

16  Bar Harbor ME 04609

17  207-288-6453

18  elissa.chesler@jax.org

19

20 ABSTRACT

21

22 Systems genetic analysis of complex traits involves the integrated analysis of

23 genetic, genomic, and disease related measures. However, these data are often

24 collected separately across multiple study populations, rendering direct correlation

25 of molecular features to complex traits impossible. Recent transcriptome-wide

26 association studies (TWAS) have harnessed gene expression quantitative trait loci

27 (eQTL) to associate unmeasured gene expression with a complex trait in genotyped

28 individuals, but this approach relies primarily on strong eQTLs. We propose a

29 simple and powerful alternative strategy for correlating independently obtained

30 sets of complex traits and molecular features. In contrast to TWAS, our approach

31 gains precision by correlating complex traits through a common set of continuous

32 phenotypes instead of genetic predictors, and can identify transcript-trait

33 correlations for which the regulation is not genetic. In our approach, a set of

34 multiple quantitative "reference" traits is measured across all individuals, while

35 measures of the complex trait of interest and transcriptional profiles are obtained in

36 disjoint sub-samples. A conventional multivariate statistical method, canonical

37 correlation analysis, is used to relate the reference traits and traits of interest in

38 order to identify gene expression correlates. We evaluate power and sample size

39 requirements of this methodology, as well as performance relative to other

40 methods, via extensive simulation and analysis of a behavioral genetics experiment

41 in 258 Diversity Outbred mice involving two independent sets of anxiety-related

42 behaviors and hippocampal gene expression. After splitting the dataset and hiding

43    one set of anxiety-related traits in half the samples, we identified transcripts

44    correlated with the hidden traits using the other set of anxiety-related traits and

45    exploiting the highest canonical correlation ($R = 0.69$) between the trait datasets.

46    We demonstrate that this approach outperforms TWAS in identifying associated

47    transcripts. Together, these results demonstrate the validity, reliability, and power

48    of the reference trait method for identifying relations between complex traits and

49    their molecular substrates.

50

51      AUTHOR SUMMARY

52

53      Systems genetics exploits natural genetic variation and high-throughput

54      measurements of molecular intermediates to dissect genetic contributions to

55      complex traits. An important goal of this strategy is to correlate molecular features,

56      such as transcript or protein abundance, with complex traits. For practical,

57      technical, or financial reasons, it may be impossible to measure complex traits and

58      molecular intermediates on the same individuals. Instead, in some cases these two

59      sets of traits may be measured on independent cohorts. We outline a method,

60      reference trait analysis, for identifying molecular correlates of complex traits in this

61      scenario. We show that our method powerfully identifies complex trait correlates

62      across a wide range of parameters that are biologically plausible and experimentally

63      practical. Furthermore, we show that reference trait analysis can identify

64      transcripts correlated to a complex trait more accurately than approaches such as

65      TWAS that use genetic variation to predict gene expression. Reference trait analysis

66      will contribute to furthering our understanding of variation in complex traits by

67      identifying molecular correlates of complex traits that are measured in different

68      individuals.

69

70    INTRODUCTION

71

72    A major goal of complex trait analysis is to discover pathways and mechanisms

73    associated with disease. By definition, these traits exhibit hallmarks of genetic

74    complexity including pleiotropy, epistasis, and gene-environment interaction.

75    Genetic mapping is a powerful approach for detecting quantitative trait loci that

76    influence complex trait variation, but it has limited power for detecting small effect

77    loci and can suffer from poor mapping resolution, hindering the identification of

78    causal genes. Moreover, these causal genetic variants do not always reside in

79    relevant therapeutic targets. Therefore, many systems genetic strategies have

80    emerged to correlate complex traits directly with molecular phenotypic variation,

81    with the goal of constructing molecular networks that are correlated with trait

82    variation from a trait-relevant tissue or cell population.

83

84    Ideally, trait correlation networks are constructed using direct phenotypic

85    measurements for each member of a population. However, there are wide-ranging

86    questions for which this approach is infeasible or impossible because it is physically,

87    technically, or financially impossible to obtain all of the measures in the same

88    individuals. To refer to phenotypes whose measurement on the same individual is

89    infeasible or impossible, we will use the term incompatible phenotypes.

90    Incompatible phenotypes arise in common experimental designs such as studies of

91    susceptibility to exposure effects where the exposure affects physiology (e.g.

92    predisposition to psychostimulant addiction) or studies of disease that relate early

93    stage changes to late stage outcomes (e.g. early molecular correlates predictive of

94    Alzheimer's disease risk). Moreover, incompatible phenotypes arise when the

95    original study population no longer is available but there is a desire to extend the

96    study to a new set of traits, a situation that is common in human genetic analyses.

97    Finally, phenotypes could be incompatible for strictly financial or logistical reasons,

98    for example due to prohibitively high costs of genomic assays in large cohorts,

99    leading to fractional collection of data on some samples and more thorough

100    characterization of others.

101

102    One emerging approach for relating gene expression and complex traits measured

103    in different cohorts of genetically diverse individuals is to exploit genetic variants

104    that affect gene expression (eQTL) to impute transcript abundance from genotypes

105    alone (Gamazon *et al.* 2015; Gusev *et al.* 2016a; b; Mancuso *et al.* 2017; Barbeira *et*

106    *al.* 2017). This enables estimation of the association between imputed gene

107    expression and complex traits, an approach that has been called a transcriptome-

108    wide association study (TWAS; Gusev *et al.* 2016a). However, the TWAS approach

109    suffers from several limitations, most notably a reliance on strong local (presumably

110    *cis*-acting) eQTL and consequent inability to impute transcript abundance for genes

111    without detected eQTL. In contrast to using sparse, discrete *genotypes* to impute

112    per-individual gene expression and infer correlation to complex traits, our approach

113    uses shared variation across a rich set of quantitative, multidimensional *phenotypes*

114    to infer gene expression correlates of phenotypic variability.

115

116   Rather than impute gene expression from genetic data, another strategy is to impute

117   phenotypic data from other phenotypes. Hormozdiari et al. (2016a) used this

118   approach to impute unmeasured phenotypes in the context of genome-wide

119   association studies (GWAS; Hormozdiari *et al.* 2016a). Specifically, the method of

120   Hormozdiari et al. (2016a) uses the correlation structure in one set of traits to

121   predict a single unmeasured target trait in a second cohort using only phenotypic

122   data. In the present study, we extend this strategy to multivariate phenotyping and

123   apply it to transcriptomics, providing a precise transcript-to-trait correlation

124   approach that can be compared to the TWAS method.

125

126   We outline a simple method, reference trait analysis, to study relations between a

127   set of complex traits of interest (*target traits*) and a set of high-dimensional

128   molecular traits obtained in disjoint subsets of individuals. Reference trait analysis

129   relates these two incompatible, multidimensional sets of phenotypes indirectly

130   through the use of a shared set of *reference traits* measured in all individuals. Since

131   target and molecular traits are not measured in the same individuals, direct

132   comparisons are impossible. Instead, we relate these traits through reference traits.

133   Reference traits are best chosen with *a priori* knowledge that they share biological

134   underpinnings with target traits. This relationship between reference and target

135   traits is exploited to compute scores from reference traits that capture variation in

136   unmeasured target traits and can be directly related to transcriptional profiles. By

137   design, our method is robust to the detection of transcript-trait associations for

138   which the regulation is not genetic or is characterized by multiple weak, indirect

139     genetic effects. Therefore, it captures biological variability associated with both

140     genetic and environmental sources of vulnerability, and has the potential to identify

141     molecular networks of complex trait variation even when there is insufficient power

142     to detect a quantitative trait locus or genome-wide significant SNP association.

143

144     In this study we develop and evaluate the reference trait analysis method using data

145     from a previously published behavioral study of Diversity Outbred mice (Logan *et*

146     *al.* 2013). Diversity Outbred mice are genetically unique; consequently, per subject

147     terminal traits such as brain gene expression can only be obtained in a single

148     exposure condition. However, the approach we propose can be useful in any

149     heterogeneous population for which a common reference set of traits is assessed.

150     Our assessment data set consists of multiple measures of anxiety-related traits in a

151     sample of Diversity Outbred mice, all of whom have been subjected to brain

152     transcriptional profiling as well as measurements of two sets of related behaviors.

153     We present an overview of our method, use these data to assess sample size

154     requirements, and quantify the method's reliability across a range of target-

155     reference trait correlations. Finally, we test whether the reference trait method

156     more faithfully recovers trait-gene expression correlations than the TWAS

157     approach.

158

159

160

161    RESULTS AND DISCUSSION

162

163    *Outline of Approach*

164

165    The reference trait analysis procedure is straightforward, and relies on well-

166    characterized canonical correlation analysis. Beginning with a population of

167    individuals, reference traits (labeled using the variable $U$) are measured on all

168    individuals, target traits (labeled with $V$) on the *training* cohort, and high

169    dimensional molecular traits (labeled with $M$) on the *testing* cohort (Figure 1).

170    Although target traits and their molecular correlates are of primary interest, the

171    choice of reference traits is an important aspect of the method. First, as we will

172    show, the strength of the multivariate relationship between target and reference

173    traits is a key parameter determining the power to detect trait-transcript

174    correlations. Second, because our method leverages shared variation between target

175    and reference traits, it identifies trait-transcript correlations driven by the portion

176    of target trait variation that is shared with reference traits. For example, studying

177    addiction-related traits using novelty behaviors as reference traits would be

178    expected to uncover transcripts associated with addiction behaviors through

179    biological pathways that also contribute to the etiology of novelty-seeking

180    behaviors.

181

182    To conduct reference trait analysis, we employ canonical correlation (Hotelling

183    1936), which can be thought of as a parent analysis of the more familiar multiple

184    regression. A multiple regression of $Y$ on $X$ models the relationships between

185    multiple $X$ measures $X_1, X_2, \dots, X_p$ and univariate $Y$. In contrast, canonical correlation

186    reveals the magnitude and nature of relationships between multivariate $U$ and $V$, e.g.

187    $U_1, U_2, \dots, U_p$ and $V_1, V_2, \dots, V_q$. Specifically, canonical correlation identifies linear

188    combinations of two multivariate measures $U$ and $V$ such that the (univariate) linear

189    combinations of each measure $\vec{u}$ and $\vec{v}$, known as canonical variables, are maximally

190    correlated. In this study we use canonical correlation to build linear combinations of

191    reference traits (transforming $U$ to $\vec{u}$) that maximize shared variance with target

192    traits ($V$) in the set of training individuals. The possible number of canonical

193    variables is limited to the size of the smaller of $U$ and $V$, and each successive

194    covariate captures a diminishing proportion of the shared variance between the

195    traits. In this study we focus on the first canonical variable, $\vec{u}_1$ or $\vec{v}_1$, which explains

196    the largest fraction of shared variance between $U$ and $V$. This quantity can be

197    thought of as a summary of each set of traits analogous to their first principal

198    component, but rather than being aligned with the axis of maximal variation *among*

199    a single set of variables, it is aligned in the direction of maximal shared variation

200    *between* the two sets of traits $U$ and $V$. For datasets with a very large number of

201    reference and/or target traits (i.e. $p \gg n$), sparse canonical correlation analysis

202    (Witten and Tibshirani 2009; Wilms and Croux 2016) may reduce over-fitting, but

203    this situation is not common when relating two sets of traits $U$ and $V$ that contain

204    organism-level phenotypes as opposed to molecular features.

205

206    The analysis of training data defines canonical coefficients that can be used to

207    compute first canonical variables from individual-level trait data (i.e. transform $U$ to

208    $\vec{u}_1$ or $V$ to $\vec{v}_1$). We use these coefficients learned from the training data (Figure 1, top)

209    to transform reference trait data from the testing cohort $U'$, which projects these

210    data in the direction of maximal shared variation with target traits. Thus, these

211    "projected" traits $\vec{u}_1'$ optimally capture the portion of variation shared between

212    reference and target traits due to their underlying genetic and environmental

213    covariation. Projected traits are then compared to high-dimensional genomic

214    measurements to extract molecular phenotypes in one sample set that co-vary with

215    target traits from another group (Figure 1, bottom right).

216

217    *Transitive reliability captures global patterns of covariation between incompatible*

218    *traits*

219

220    Reference trait analysis reveals covariation between molecular phenotypes and

221    target trait variation. There are many possible applications of this strategy. For

222    example, in addiction research, many studies evaluate transcriptional response to

223    drug exposure but are unable to evaluate the predisposing characteristics of a drug

224    naïve brain that associate with addiction-related behaviors. Using a reference trait

225    strategy, one can evaluate the transcriptomes of drug naïve brains and relate them

226    to the response to drug self-administration through a set of reference traits that do

227    not involve drug exposure. We have previously estimated the association of novelty

228    seeking and drug self-administration in mice, revealing a canonical correlation of

229    0.61 among these sets of traits (Dickson *et al.* 2015).

230

231    To evaluate whether the reference traits strategy could be applied to find

232    transcriptional correlates of drug self-administration, we used a dataset where

233    reference, target, and molecular trait profiling were performed on the same

234    individuals to allow for assessment of the accuracy and robustness of the method. In

235    this data set, transcriptional profiles, target traits, and reference traits are available

236    for all individuals. This allows evaluation of the properties of the reference trait

237    strategy, including robustness and sample size requirements. Specifically, we

238    studied relationships between two distinct sets of anxiety-related traits and

239    hippocampal gene expression, where all traits were measured in each of $N = 258$

240    Diversity Outbred mice (Logan *et al.* 2013). The anxiety-related traits consisted of

241    eleven measurements of open-field arena exploration behaviors and five

242    measurements of light-dark box behaviors (Supplementary Table 1). A canonical

243    correlation analysis of these two sets of traits yielded a statistically significant

244    model ($F_{55,1123.75} = 4.48$, $p < 2 \times 10^{-16}$, Wilk's $\lambda = 0.400$) that had a first canonical

245    correlation coefficient of magnitude 0.69. This was higher than all univariate

246    correlations between open-field and light-dark box traits (median 0.11, maximum

247    0.65), and similar in magnitude to the shared variation revealed by the first

248    canonical variable in the motivating analysis of novelty-related behaviors and

249    cocaine self-administration (Dickson *et al.* 2015). We arbitrarily designated the

250    open-field traits as target traits and light-dark box traits as reference traits. For

251    reference trait analysis, we hid gene expression data for some mice (training set)

252    and open-field data for the remaining mice (testing set).

253

254    In this evaluation of reference trait analysis, we know the true values of all hidden

255    data and can directly evaluate the power of the method to reveal gene expression

256    patterns associated with target trait variation. Specifically, we estimate canonical

257    coefficients (weights to calculate canonical variables) from the training set and use

258    them to calculate projected traits $\vec{u}_1'$ in the testing set. To quantify the performance

259    of reference trait analysis when the true answer is known, we computed

260    correlations in testing set animals between gene expression $E$ and either (1) the first

261    projected trait $\vec{u}_1'$, cor $(E, \vec{u}_1')$ or (2) the first canonical variable computed using

262    hidden target traits $\vec{v}_1'$, cor$(E, \vec{v}_1')$. The latter quantity, the "truth", is unavailable in a

263    real application of reference trait analysis. A set of reference traits that perfectly

264    captures all variation in target traits would result in a vector of gene expression-

265    trait correlations that is identical whether the target traits were known or projected

266    from reference traits (i.e. the reference traits serve as a perfect surrogate for target

267    traits). We define *transitive reliability* as the correlation between these vectors i.e.

268    cor$[\text{cor}(E, \vec{u}_1'), \text{cor}(E, \vec{v}_1')]$. High transitive reliability would indicate that strong

269    correlations between gene expression and target traits are likely to be identified

270    using projected traits.

271

272    Transitive reliability, estimated using real gene expression data and simulated

273    canonical variables with known correlation, scales linearly with the magnitude of

274    the canonical correlation coefficient (Figure 2A), confirming our intuition that

275    greater sharing of variation between target and reference traits increases the utility

276    of leveraging reference traits to understand target trait variation. We divided the

277    anxiety dataset into equally sized subsets (partially overlapping for larger sample

278    sizes) to examine the dependence of transitive reliability on sample size. The

279    canonical correlation was upwardly biased for small sample sizes ($N < 90$; data not

280    shown), as has previously been recognized (e.g. Thompson 1990). When we used

281    Wherry's correction as suggested by Thompson (1990), canonical correlations no

282    longer depended on sample size (linear model; $p > 0.8$). Overall, transitive reliability

283    asymptotically approached the magnitude of the canonical correlation coefficient

284    calculated from the full dataset (Figure 2B, black line), demonstrating that global

285    patterns of trait-gene expression correlation can be recovered with relatively

286    modest sample sizes using the reference trait approach. In contrast, weights from

287    the smallest (fifth) canonical covariate, which captures little shared variation

288    between datasets, produced low transitive reliabilities (median 0.11).

289

290    *Reference trait analysis successfully identifies known trait correlations*

291

292    Ultimately, the primary goal of reference trait analysis is to identify molecular

293    correlates of unmeasured phenotypes. To discover these correlates, individual gene

294    expression levels are correlated to projected traits. To test this strategy, we first

295    employed reference trait analysis on the anxiety-related phenotype data described

296    above. After randomly splitting the dataset and withholding open-field data

297    (arbitrarily designated as target traits) in half the individuals, we identified gene

298    expression levels correlated to projected reference traits. We found high overlap

299    between the genes most strongly correlated with hidden target trait canonical

300    variable 1 and those most strongly correlated with projected traits (23% overlap

301    among genes with top 5% of correlations to each trait, compare to 2.5% expected

302    overlap; $p < 1 \times 10^{-15}$, Fisher's Exact Test). Across all genes, including those with

303    weaker correlations, we found that the vector of trait-gene expression correlations

304    computed using reference trait analysis showed significant similarity to the true

305    correlations ($p < 0.001$, permutation test using generalized Jaccard similarity

306    statistic). Moreover, in contrast to the alternative methods for identifying trait-gene

307    expression correlations discussed above, some correlations detected using

308    reference trait analysis involved genes with no significant eQTL (e.g. 42% of top 50

309    correlations). These genes, which are demonstrably associated with trait variation,

310    would not be detectable using TWAS type approaches.

311

312    To examine the power and robustness of reference trait analysis across a wide

313    range of biologically plausible parameter values, we conducted extensive

314    simulations. We simulated data across a range of sample sizes (100, 200, 300, …,

315    1000, 1200, 1400, …, 2000) and enforced a similar covariance structure to the

316    observed data. Specifically, data were simulated using observed covariances within

317    each set of anxiety traits, but we perturbed covariances between the two sets of

318    traits in order to generate datasets with varying canonical correlations. We then

319    simulated gene expression levels with known correlation to the first target trait

320  canonical variable, $\vec{v}_1$ ($\rho$ = 0.2, 0.225, 0.25, ..., 0.9 with 20 genes each). We simulated

321  trait data and gene expression data at random for each of 1,000 simulations for each

322  sample size.

323

324  For each simulation, after hiding target traits in half the individuals and gene

325  expression data in the other half, we conducted reference trait analysis. We

326  computed projected reference traits, correlated to gene expression, and quantified

327  performance as the fraction of true trait-gene expression correlations that were

328  detected using a 10% false discovery rate (FDR) threshold. For high trait-gene

329  correlations ($\rho$ > 0.6) and strong target-reference trait canonical correlations (R =

330  0.7 or 0.9), the correlation of interest was essentially always detected (Figure 3). For

331  lower target-reference trait canonical correlations (R = 0.5), even relatively modest

332  true trait-gene expression correlations (e.g. $\rho$ = 0.3) were often detected with

333  sample sizes above ~300 individuals (Figure 3). Thus, reference trait analysis was a

334  highly effective means for identifying trait-gene expression correlations across a

335  diverse range of practical sample sizes, typical values for trait-to-gene expression

336  correlation, and canonical correlation parameters.

337

338  *Comparison of reference trait analysis to related approaches*

339

340  An alternative approach to identifying genes associated with complex traits is to

341  make use of known genetic variation that regulates gene expression (gene

342  expression QTL or eQTL). There has been considerable recent interest in methods

343     that integrate complex trait associations and gene expression genetics in order to

344     identify genes whose expression is associated with trait variation (Nica *et al.* 2010;

345     Wallace *et al.* 2012; He *et al.* 2013; Gamazon *et al.* 2015; Gusev *et al.* 2016a; Zhu *et*

346     *al.* 2016; Hormozdiari *et al.* 2016b; Wen *et al.* 2017; Hauberg *et al.* 2017). Several

347     methods perform tests of the hypothesis that genome-wide association (GWA)

348     signals and eQTLs are truly colocalized versus independent but appearing

349     colocalized due to linkage disequilibrium  (Nica *et al.* 2010; Wallace *et al.* 2012;

350     Giambartolomei *et al.* 2014; Fortune *et al.* 2015; Zhu *et al.* 2016; Hormozdiari *et al.*

351     2016b; Wen *et al.* 2017; Hauberg *et al.* 2017). Another approach that is more

352     directly applicable to the experimental designs studied herein is to harness strong

353     genetic predictors of gene expression variation (eQTL) to impute transcriptomes in

354     genotyped and phenotyped cohorts, which allows detection of trait-expression

355     correlations (the TWAS approach; Gamazon *et al.* 2015; Gusev *et al.* 2016a; b;

356     Mancuso *et al.* 2017; Barbeira *et al.* 2017). TWAS is an approach that is

357     complementary to reference trait analysis, and has been a particularly powerful

358     method for discovery of candidate genes driving GWA signals detected in very large

359     human cohorts (tens or hundreds of thousands of individuals). Supplementary

360     Figure 1 provides a comparison of genotype, phenotype, and gene expression data

361     in the reference traits and TWAS strategies. One weakness of the TWAS approach is

362     that it hinges on the presence of detectable eQTL (typically local, presumably cis-

363     acting eQTL; but see He *et al.* 2013; Vervier and Michaelson 2016). In humans, even

364     panels of 1,000 individuals with gene expression measurements only result in a

365     modest number of genes (500-4,000) with significant *cis*-heritability that can be

366    imputed in the cohort lacking gene expression data (Gusev *et al.* 2016a). In contrast,

367    reference trait analysis has no requirement for detection of eQTLs, and therefore it

368    is amenable to detect of correlation of transcripts with complex expression

369    regulatory mechanisms to traits of similarly complex regulation, and retains

370    performance across lower sample sizes, as we demonstrate below.

371

372    Although TWAS and reference trait analysis utilize different data types, both are

373    tools inferring relations between complex traits and transcript abundance, so we

374    sought to compare their performance on the same dataset. For TWAS, we used

375    methods implemented in the software suite PrediXcan (Gamazon *et al.* 2015). We

376    randomly divided our anxiety dataset in half and considered open-field

377    measurements as target traits. We withheld gene expression measurements in half

378    the animals; therefore, only genotype and reference trait data were visible for all

379    animals. We built predictive models of gene expression from the training cohort of

380    mice, applied these models to impute gene expression in the testing cohort, and

381    calculated correlations between imputed gene expression and a summary measure

382    of the target traits (first canonical variable). We conducted 1,000 permutations with

383    random 50:50 divisions of the anxiety dataset to account for stochastic sampling

384    effects. For each replicate, we compared global trait-gene expression correlations

385    for PredictDB-imputed gene expression versus those computed using projected

386    traits obtained with our new method. In the former case, trait data is available and

387    gene expression data is imputed, while in the latter case gene expression data is

388    available and trait data is imputed.

389

390    For direct comparisons between reference trait analysis and TWAS, we ran

391    reference trait analysis using only genes that were significantly predicted by the

392    PredictDB module of PrediXcan (FDR < 5%; see Methods). Across the 1,000

393    permutations, we imputed gene expression for a mean 12,250 genes (range 11,640-

394    12,750; mean represents ~70% of total 17,539 genes measured), indicating that a

395    substantial fraction of genes has insufficient local genetic signal for accurate

396    imputation. An advantage of reference trait analysis is that it is not limited by the

397    presence of strong eQTL and all genes can be tested for association with projected

398    reference traits. For each of the 1,000 permutations, we computed the transitive

399    reliability of TWAS and of reference trait analysis. Reference trait analysis more

400    accurately captured global patterns of trait-transcript correlation than TWAS

401    (Figure 4). Specifically, transitive reliability for target trait first canonical covariate-

402    gene expression correlations was higher using the reference trait approach

403    (measured gene expression and projected reference traits) compared to the TWAS

404    approach (imputed gene expression and measured traits) for 92.7% of simulations

405    (Figure 4; Supplementary Figure 2 shows an example of results from one

406    permutation). Thus, we show empirically that reference trait analysis outperforms

407    TWAS in the mouse anxiety dataset.

408

409    In addition to the quantitative comparison of the methods, we sought to determine

410    which approach provided the best retrieval of known anxiety related genes. To

411    perform this analysis we made use of GeneWeaver's database of gene sets curated

412    from multiple sources (Baker *et al.* 2016). The top four hundred genes identified

413    using each analysis method were entered as three gene lists, and each gene list was

414    compared to every gene set in the GeneWeaver database via Jaccard similarity. For

415    each, the top 249 similar gene sets were exported, and a rater with expertise in

416    behavioral neuroscience who was blind to the analysis methods scored a combined

417    list of all similar gene sets obtained in these three analyses. Gene sets were

418    categorized discretely based on relevance to anxiety, with categories including

419    irrelevant, generally relevant to brain or behavior, and specifically relevant to

420    anxiety. We found that true open-field first canonical variable—gene expression

421    correlations had highest relevance to anxiety. The top truly correlated genes were

422    similar to gene sets more relevant to anxiety than those genes identified using

423    reference traits or those using TWAS ($p = 0.0065$ and $p = 1.5 \times 10^{-14}$, respectively;

424    two-sided Fisher's Exact Test). Nevertheless, reference trait analysis performed

425    significantly better than TWAS at identifying genes with similarity to anxiety-

426    relevant gene sets ($p = 7.3 \times 10^{-6}$).

427

428    Finally, another alternative to relating traits and transcripts between population

429    cohorts is to make use of polygenic risk predictors trained using genome-wide

430    genotypes and phenotypes, and applied to individuals with genotypes but missing

431    phenotypes (in this case, samples with only transcriptional profiles available)

432    (Makowsky *et al.* 2011; Dudbridge 2013; Wray *et al.* 2013). However, theoretical

433    considerations and empirical results suggest that this approach generally requires

434    sample sizes much larger than 1,000 individuals to obtain accurate predictions

435    (Dudbridge 2013). In the context of reference trait analysis, relating complex

436    reference and target traits that share high canonical correlation implicitly leverages

437    the common polygenic or omnigenic (Boyle *et al.* 2017) basis of these traits by

438    making use of all of the information contained in continuous quantitative variation.

439

440

441    *Conclusions*

442

443    We have described a general method for exploring trait covariation among

444    incompatible and independently collected phenotypes studied in disjoint samples of

445    genetically diverse individuals to extract molecular networks associated with

446    disease. Our method utilizes canonical correlation analysis, a standard multivariate

447    statistical method, to relate incompatible phenotypes using a set of reference traits

448    measured on all individuals. Our analyses demonstrate that this approach performs

449    well over a range of parameters typically encountered in the study of trait

450    correlations, and under sample size requirements that are practical to obtain. This

451    approach can be useful both for capturing global patterns of covariation between

452    target traits and high-dimensional molecular phenotypes, as well as for identifying

453    specific molecular correlates to target traits. Our method identifies trait-gene

454    expression associations and we do not assert that these associations are necessarily

455    causal, as has been recognized by studies relating GWAS results and eQTL (Gamazon

456    *et al.* 2015; Gusev *et al.* 2016a; Hauberg *et al.* 2017).

457

458    When will reference trait analysis be a useful tool? Intuitively, and as demonstrated

459    in Figure 3, large sample sizes, precise trait measurements, and high shared

460    variance between reference and target traits would allow for the most accurate

461    estimation of canonical correlation coefficients and high power to detect

462    correlations to molecular phenotypes. Although our method could be applied in a

463    wide variety of scenarios, it is likely to be particularly useful for studies of highly

464     complex, polygenic, multidimensional traits (e.g. behavior, physiology, and

465     morphology) in cohorts of modest size. As with any method that applies information

466     learned from one cohort to biological measures from another cohort, reference trait

467     analysis requires the absence of systematic differences (i.e. heterogeneity in

468     population characteristics) between the training and testing cohorts. For very large

469     cohorts of individuals where obtaining suitable reference traits may be difficult,

470     polygenic scores based on either genetic predictors alone or on a combination of

471     genetic and environmental risk factors (Dudbridge *et al.* 2017) may be a valuable

472     approach for predicting phenotypic variation in a test cohort that can then be

473     correlated with molecular networks.

474

475     Although our application of reference trait analysis involves correlations to high

476     dimensional molecular phenotypes, the method could, in principle, be applied to any

477     sets of phenotypes that are multivariate in nature. Moreover, the high relative

478     performance of our method underscores the importance of extensive phenotyping

479     using quantitative traits rather than relying on binary indicators of disease and

480     disease-related phenotypes that may mask complex underlying etiologies. We

481     anticipate that the framework outlined in this study will be increasingly useful as

482     studies of diverse, genetically unique populations become more widespread. A

483     useful future extension to this approach would incorporate statistical techniques

484     such as sparse canonical correlation analysis (Witten and Tibshirani 2009; Wilms

485     and Croux 2016), which could permit inference in phenome-level studies where the

486     target or reference traits are high dimensional. Overall, our approach is likely to be

487     particularly important in functional genomics studies, those utilizing post-mortem

488     subjects, and large population studies in which individuals are unavailable for

489     further characterization.

490

491

492

493    **Materials and Methods**

494

495    *Mouse rearing and phenotyping*

496

497    Diversity Outbred mice (J:DO, The Jackson Laboratory) are a heterogeneous stock

498    derived from the same eight founder strains as the Collaborative Cross (Svenson *et*

499    *al.* 2012; Churchill *et al.* 2012; Gatti *et al.* 2014; Chesler *et al.* 2016). In this study we

500    used a subset ($N$ = 258) of the 283 Diversity Outbred mice studied by Logan et al.

501    (2013) with hippocampal gene expression profiled by RNA-Seq (see below). Mice in

502    this study were from generations 4 to 5 (G4-G5) of the DO population. Briefly, each

503    mouse was acclimated to the housing area, and subject to a brief testing battery

504    which included a 20 minute novel open-field test and a 10 minute light-dark test,

505    among other common behavioral tasks. The open-field and light-dark tests are used

506    to measure exploratory activity and approach-avoidance behavior. Many complex

507    trait measures can be extracted from these tasks. For this analysis, we chose two

508    sets of informative measures (Supplementary Table 1). Complete details of animal

509    rearing, husbandry and phenotyping are presented in Logan et al. (2013). Mice were

510    sacrificed using decapitation which was necessary to preserve fresh brain tissue in

511    the absence of drug or asphyxiation. All procedures and protocols were approved by

512    The Jackson Laboratory Animal Care and Use Committee, and were conducted in

513    compliance with the National Institutes of Health Guidelines for the Care and Use of

514    Laboratory Animals.

515

516 *Genotyping*

517

518 DNA was prepared from tail biopsies and samples were genotyped using the Mouse

519 Univeral Genotyping Array (MUGA) (Morgan *et al.* 2016). We obtained genotypes at

520 7,802 markers from arrays processed by GeneSeek (Lincoln, NE). We used

521 intensities from each array to infer the haplotype blocks in each individual DO

522 genome using a hidden Markov model (Gatti *et al.* 2014).

523

524 *Gene expression profiling*

525

526 Total hippocampal RNA was isolated using the TRIzol® Plus RNA purification kit

527 (Life Technologies Corp., Carlsbad, CA) with on-column DNase digestion. Samples

528 for RNA-Seq analysis were prepared using the TruSeq kit (Illumina Inc., San Diego,

529 CA) according to the manufacturer's protocols and subjected to paired-end 100 base

530 pair sequencing on the HiSeq 2000 (Illumina) per manufacturer's

531 recommendations. RNA sequencing was performed in nine sequencing runs with

532 two technical replicates for each sample, resulting in an averaging sequencing depth

533 of approximately 24 million reads per sample after pooling technical replicates. To

534 obtain estimates of gene expression, we aligned reads to individualized diploid

535 genomes using the bowtie aligner (Langmead *et al.* 2009) and quantified transcript

536 abundance by allocating multi-mapping reads using the EM algorithm with RSEM (Li

537 and Dewey 2011) as described in Munger et al. (2014). Raw counts in each sample

538 were normalized to the upper quartile value and transformed to normal scores.

539    *Reference trait analysis*

540

541    We conducted reference trait analysis using R version 3.3.2 (R Core Team 2016).

542    Canonical correlation analysis was carried out using the cancor function in base R.

543    We regressed out the effect of sex on each phenotype because it is not of primary

544    interest in this study. An example walk-through of a reference trait analysis and

545    code to carry out the analyses described in this paper are available at

546    https://daskelly.github.io/reference_traits/reference_trait_analysis_walkthrough.ht

547    ml.

548

549    To examine the power and robustness of reference trait analysis, we simulated data

550    with varying sample sizes and canonical correlation coefficients. We based our

551    simulations on the anxiety phenotype data, consisting of open-field exploration and

552    light-dark box behavioral measures. Specifically, for each of 1,000 simulations we

553    started with the covariance matrix computed from five open-field and five light-dark

554    box traits and randomly increased or decreased each of the 5×5=25 inter-dataset

555    covariances by 20%. We then simulated multivariate normal phenotype data with

556    the specified covariance matrix. This procedure resulted in two multivariate

557    datasets (simulated open-field and light-dark box traits), where the covariance

558    structure *within* each dataset was similar to that in the real data but with different

559    covariances *between* datasets. When a canonical correlation analysis was carried out

560    on each pair of simulated datasets, the magnitude of the first canonical correlation

561    coefficient varied between $R = 0.35$ and $R = 0.98$, due to the variation in inter-

562    dataset covariances.

563

564    We simulated gene expression traits with exact correlation to the first target trait

565    canonical variable $\vec{v}_1$ in the simulated dataset. In order to simulate a random vector

566    of observations with defined correlation to an existing vector, we took advantage of

567    the geometric property that the cosine between two mean-centered vectors equals

568    their correlation. Therefore, a random vector with defined correlation to an existing

569    vector can be computed by starting with random draws from a normal distribution,

570    mean-centering, and applying standard linear algebra operations.

571

572    After hiding target traits in half the individuals and gene expression data in the

573    other half, we conducted reference trait analysis and quantified performance as the

574    fraction of the time true trait-gene expression correlations were detected using a

575    10% FDR threshold. $P$-values for trait-gene expression correlations were calculated

576    using a two-sided $T$ statistic and correlations deemed significant at a 10% FDR were

577    identified using $q$-values (Storey and Tibshirani 2003).

578

579    *Imputing gene expression using TWAS*

580

581    We divided the anxiety dataset in half and considered open-field measurements as

582    target traits, hiding gene expression measurements for the animals where we did

583    not hide open-field traits. For the TWAS strategy, our training cohort consisted of

584 animals with genotypes and gene expression data, and our testing cohort consisted

585 of animals with genotypes and open-field traits (i.e. training/testing labels are

586 reversed from reference trait analysis, see Supplementary Figure 1). Diversity

587 Outbred mice are an outbred population with genomic ancestry derived from eight

588 inbred founder strains. We used methods implemented in R/qtl2 software

589 (http://kbroman.org/qtl2/) to impute single nucleotide polymorphism (SNP)

590 variation in each mouse from array-based genotypes obtained at coarser resolution

591 (see above) using known SNP genotypes present in founder haplotypes. This

592 resulted in genotypes for ~30 million SNPs. Given the limited number of

593 generations of outbreeding, haplotype blocks in Diversity Outbred mice typically

594 stretch for megabases (Svenson *et al.* 2012), leading to strong local linkage

595 disequilibrium (LD). As such, we used PLINK version 1.9 (Purcell *et al.* 2007) to

596 prune variants in very strong LD in the eight founder strains, using the parameters -

597 -indep-pairwise 200kb 40kb 0.95. This procedure reduced the number of SNPs to

598 235,335 with minimal loss of information.

599

600 To impute gene expression, we used the PredictDB module of PrediXcan (Gamazon

601 *et al.* 2015) to build predictive models of gene expression from local genotypes

602 within 10Mb of each gene, with sex included as a covariate. We conducted 1,000

603 permutations with random 50:50 divisions of the anxiety dataset to account for

604 stochastic sampling effects. For each replicate we obtained predictive models of

605 gene expression by running PredictDB on the training cohort and applied them to

606 the testing cohort in order to impute gene expression. Following Gamazon et al.

607    (2015; https://github.com/hakyimlab/PrediXcan), we considered only genes with

608    models that were significantly predictive of gene expression (FDR ≤ 5%). Finally, we

609    calculated correlations between imputed gene expression and a summary measure

610    of the target traits (first canonical variable) in the testing cohort. Results were

611    nearly identical whether we correlated to the first canonical variable or first

612    principal component of the target traits (median transitive reliability 45% vs. 44%),

613    but correlations to first canonical variable allow for direct comparison with results

614    from reference trait analysis.

615

616    *Scoring gene sets to assess retrieval of known anxiety-related genes*

617

618    To score gene sets for relevance to anxiety, a rater with expertise in behavioral

619    neuroscience who was blind to the analysis methods scored a combined list of all

620    gene sets obtained herein. We assigned a score of zero to irrelevant data sets, a

621    score of two to gene sets with general brain or behavior relevance, and a score of

622    four to anxiety relevant data sets in which either the gene set was generated in an

623    anxiety relevant experiment, the gene set consisted of genes interacting with a

624    compound known to be anxiolytic or anxiogenic, or the gene set was a Gene

625    Ontology annotation set with direct biological relevance to anxiety. For compounds,

626    a single MEDLINE query of the compound name and 'anxiety' was performed and

627    the results of the query were examined for overall conceptual relevance.

628

629    *Data availability*

630

631    Raw RNA-Seq gene expression data from the hippocampus of 258 Diversity Outbred

632    mice are available from ArrayExpress (accession number XXX). A processed and

633    normalized gene expression matrix is available as Supplementary Dataset 1.

634    Phenotype data acquired via the open-field and light-dark box paradigms are

635    available as Supplementary Datasets 2 and 3.

636

645

646    **Author contributions:**

647    DAS – Conceptualization, Data Curation, Formal Analysis, Methodology, Software,

648    Visualization, Writing

649    NR – Data Curation, Formal Analysis, Software

650    RFR – Investigation, Methodology

651    JHG – Formal Analysis

652     EJC – Conceptualization, Funding Acquisition, Methodology, Project Administration,

653     Supervision, Writing

654

## References

Baker E., Bubier J. A., Reynolds T., Langston M. A., Chesler E. J., 2016 GeneWeaver: data driven alignment of cross-species genomics in biology and disease. Nucleic Acids Res. 44: D555-559.

Barbeira A. N., Dickinson S. P., Torres J. M., Bonazzola R., Zheng J., *et al.*, 2017 Exploring the phenotypic consequences of tissue specific gene expression variation inferred from GWAS summary statistics. bioRxiv: 045260.

Boyle E. A., Li Y. I., Pritchard J. K., 2017 An Expanded View of Complex Traits: From Polygenic to Omnigenic. Cell 169: 1177–1186.

Chesler E. J., Gatti D. M., Morgan A. P., Strobel M., Trepanier L., *et al.*, 2016 Diversity Outbred Mice at 21: Maintaining Allelic Variation in the Face of Selection. G3 Bethesda Md 6: 3893–3902.

Churchill G. A., Gatti D. M., Munger S. C., Svenson K. L., 2012 The Diversity Outbred mouse population. Mamm. Genome Off. J. Int. Mamm. Genome Soc. 23: 713–718.

Dickson P. E., Ndukum J., Wilcox T., Clark J., Roy B., *et al.*, 2015 Association of novelty-related behaviors and intravenous cocaine self-administration in Diversity Outbred mice. Psychopharmacology (Berl.) 232: 1011–1024.

Dudbridge F., 2013 Power and Predictive Accuracy of Polygenic Risk Scores. PLOS Genet. 9: e1003348.
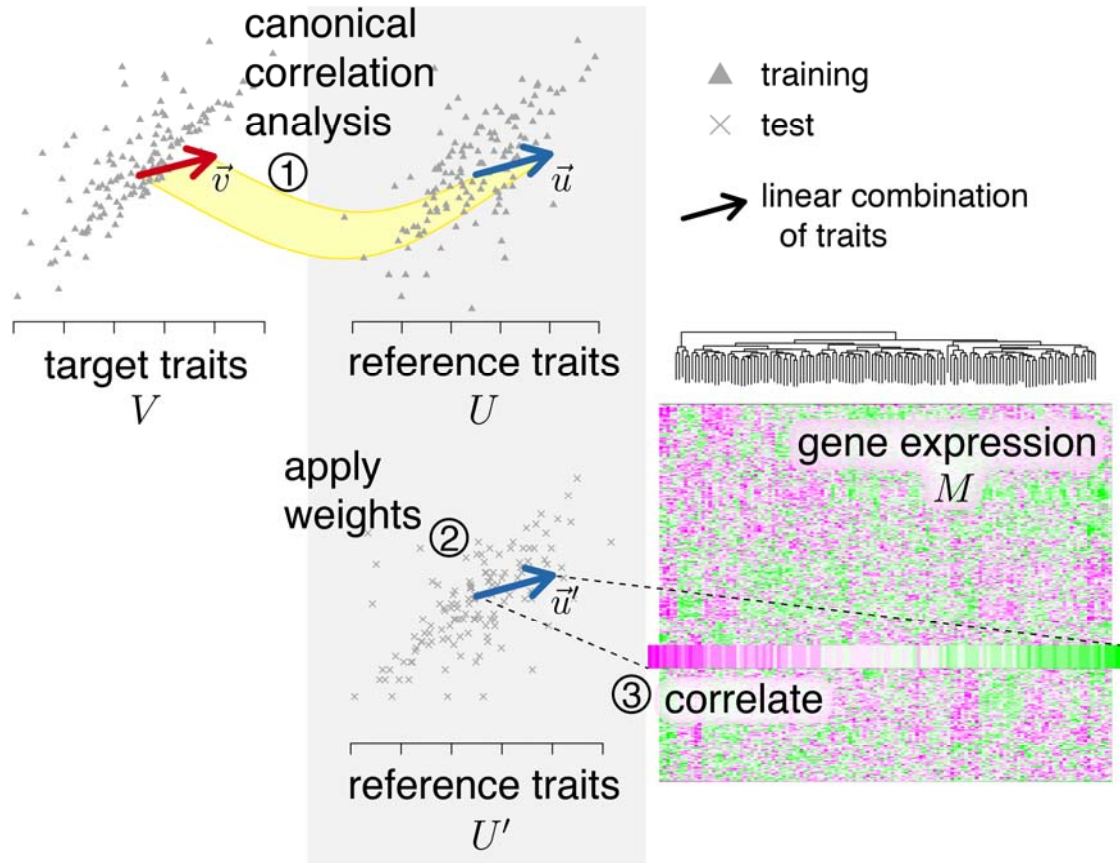
676    Dudbridge F., Pashayan N., Yang J., 2017 Predictive accuracy of combined genetic

677         and environmental risk scores. Genet. Epidemiol.: 1–16.

678    Fortune M. D., Guo H., Burren O., Schofield E., Walker N. M., *et al.*, 2015 Statistical

679         colocalization of genetic risk variants for related autoimmune diseases in the

680         context of common controls. Nat. Genet. 47: 839–846.

681    Gamazon E. R., Wheeler H. E., Shah K. P., Mozaffari S. V., Aquino-Michaels K., *et al.*,

682         2015 A gene-based association method for mapping traits using reference

683         transcriptome data. Nat. Genet. 47: 1091–1098.

684    Gatti D. M., Svenson K. L., Shabalin A., Wu L.-Y., Valdar W., *et al.*, 2014 Quantitative

685         trait locus mapping methods for diversity outbred mice. G3 Bethesda Md 4:

686         1623–1633.

687    Giambartolomei C., Vukcevic D., Schadt E. E., Franke L., Hingorani A. D., *et al.*, 2014

688         Bayesian Test for Colocalisation between Pairs of Genetic Association Studies

689         Using Summary Statistics. PLOS Genet. 10: e1004383.

690    Gusev A., Ko A., Shi H., Bhatia G., Chung W., *et al.*, 2016a Integrative approaches for

691         large-scale transcriptome-wide association studies. Nat. Genet. 48: 245–252.

692    Gusev A., Mancuso N., Finucane H. K., Reshef Y., Song L., *et al.*, 2016b Transcriptome-

693         wide association study of schizophrenia and chromatin activity yields

694         mechanistic disease insights. bioRxiv: 067355.

695    Hauberg M. E., Zhang W., Giambartolomei C., Franzén O., Morris D. L., *et al.*, 2017

696        Large-Scale Identification of Common Trait and Disease Variants Affecting

697        Gene Expression. Am. J. Hum. Genet. 100: 885–894.

698    He X., Fuller C. K., Song Y., Meng Q., Zhang B., *et al.*, 2013 Sherlock: Detecting Gene-

699        Disease Associations by Matching Patterns of Expression QTL and GWAS. Am.

700        J. Hum. Genet. 92: 667–680.

701    Hormozdiari F., Kang E. Y., Bilow M., Ben-David E., Vulpe C., *et al.*, 2016a Imputing

702        Phenotypes for Genome-wide Association Studies. Am. J. Hum. Genet. 99: 89–

703        103.

704    Hormozdiari F., van de Bunt M., Segrè A. V., Li X., Joo J. W. J., *et al.*, 2016b

705        Colocalization of GWAS and eQTL Signals Detects Target Genes. Am. J. Hum.

706        Genet. 99: 1245–1260.

707    Hotelling H., 1936 Relations Between Two Sets of Variates. Biometrika 28: 321–377.

708    Langmead B., Trapnell C., Pop M., Salzberg S. L., 2009 Ultrafast and memory-efficient

709        alignment of short DNA sequences to the human genome. Genome Biol. 10:

710        R25.

711    Li B., Dewey C. N., 2011 RSEM: accurate transcript quantification from RNA-Seq data

712        with or without a reference genome. BMC Bioinformatics 12: 323.

713 Logan R. W., Robledo R. F., Recla J. M., Philip V. M., Bubier J. A., *et al.*, 2013 High-

714        precision genetic mapping of behavioral traits in the diversity outbred mouse

715        population. Genes Brain Behav. 12: 424–437.

716 Makowsky R., Pajewski N. M., Klimentidis Y. C., Vazquez A. I., Duarte C. W., *et al.*,

717        2011 Beyond Missing Heritability: Prediction of Complex Traits. PLOS Genet.

718        7: e1002051.

719 Mancuso N., Shi H., Goddard P., Kichaev G., Gusev A., *et al.*, 2017 Integrating Gene

720        Expression with Summary Association Statistics to Identify Genes Associated

721        with 30 Complex Traits. Am. J. Hum. Genet. 100: 473–487.

722 Morgan A. P., Fu C.-P., Kao C.-Y., Welsh C. E., Didion J. P., *et al.*, 2016 The Mouse

723        Universal Genotyping Array: From Substrains to Subspecies. G3 Genes

724        Genomes Genet. 6: 263–279.

725 Munger S. C., Raghupathy N., Choi K., Simons A. K., Gatti D. M., *et al.*, 2014 RNA-Seq

726        Alignment to Individualized Genomes Improves Transcript Abundance

727        Estimates in Multiparent Populations. Genetics 198: 59–73.

728 Nica A. C., Montgomery S. B., Dimas A. S., Stranger B. E., Beazley C., *et al.*, 2010

729        Candidate Causal Regulatory Effects by Integration of Expression QTLs with

730        Complex Trait Genetic Associations. PLOS Genet. 6: e1000895.

731    Purcell S., Neale B., Todd-Brown K., Thomas L., Ferreira M. A. R., *et al.*, 2007 PLINK: a

732            tool set for whole-genome association and population-based linkage

733            analyses. Am. J. Hum. Genet. 81: 559–575.

734    R Core Team, 2016 *R: A Language and Environment for Statistical Computing*. R

735            Foundation for Statistical Computing, Vienna, Austria.

736    Storey J. D., Tibshirani R., 2003 Statistical significance for genomewide studies. Proc.

737            Natl. Acad. Sci. U. S. A. 100: 9440–9445.

738    Svenson K. L., Gatti D. M., Valdar W., Welsh C. E., Cheng R., *et al.*, 2012 High-

739            resolution genetic mapping using the Mouse Diversity outbred population.

740            Genetics 190: 437–447.

741    Thompson B., 1990 Finding a Correction for the Sampling Error in Multivariate

742            Measures of Relationship: A Monte Carlo Study. Educ. Psychol. Meas. 50: 15–

743            31.

744    Vervier K., Michaelson J. J., 2016 SLINGER: large-scale learning for predicting gene

745            expression. Sci. Rep. 6: 39360.

746    Wallace C., Rotival M., Cooper J. D., Rice C. M., Yang J. H. M., *et al.*, 2012 Statistical

747            colocalization of monocyte gene expression and genetic risk variants for type

748            1 diabetes. Hum. Mol. Genet. 21: 2815–2824.

749    Wen X., Pique-Regi R., Luca F., 2017 Integrating molecular QTL data into genome-

750        wide genetic association analysis: Probabilistic assessment of enrichment

751        and colocalization. PLOS Genet. 13: e1006646.

752    Wilms I., Croux C., 2016 Robust sparse canonical correlation analysis. BMC Syst. Biol.

753        10: 72.

754    Witten D. M., Tibshirani R. J., 2009 Extensions of Sparse Canonical Correlation

755        Analysis with Applications to Genomic Data. Stat. Appl. Genet. Mol. Biol. 8: 1–

756        27.

757    Wray N. R., Yang J., Hayes B. J., Price A. L., Goddard M. E., *et al.*, 2013 Pitfalls of

758        predicting complex traits from SNPs. Nat. Rev. Genet. 14: 507–515.

759    Zhu Z., Zhang F., Hu H., Bakshi A., Robinson M. R., *et al.*, 2016 Integration of summary

760        data from GWAS and eQTL studies predicts complex trait gene targets. Nat.

761        Genet. 48: 481–487.
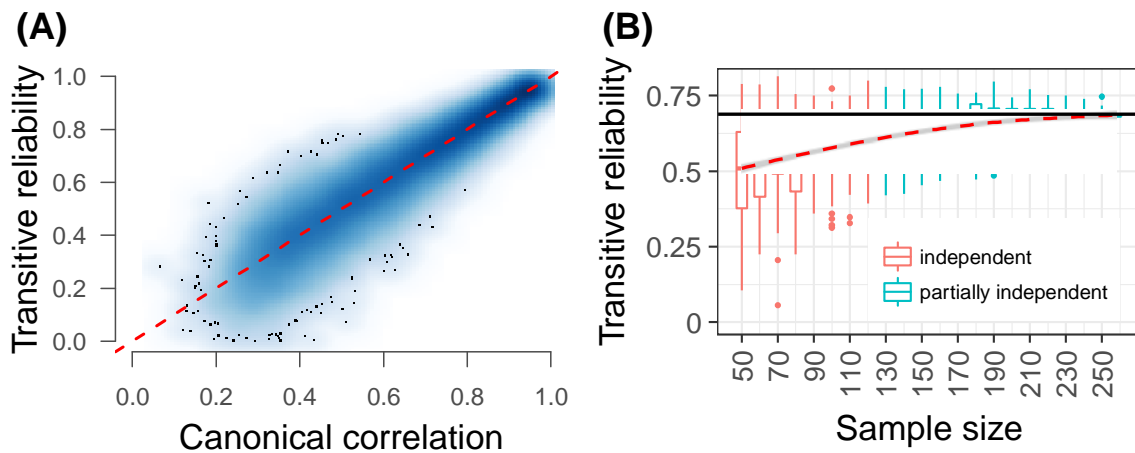
762

763

764
765
766 **Figure 1:** Overview of reference trait analysis. Target and reference traits are
767 measured in a set of training individuals (top plots; grey triangles), while reference
768 traits and gene expression are measured in test individuals (bottom plots and X
769 symbols). (1) Canonical correlation is used to identify a linear combination of
770 reference traits (top blue arrow) that best captures variation in the traits of interest
771 (red arrow; yellow curve connecting arrows represents canonical correlation
772 analysis). (2) The weights derived from canonical correlation analysis are applied to
773 reference traits in the testing population to derive reference trait scores for each
774 individual (projected reference traits; bottom blue arrow). (3) Projected reference
775 traits are correlated with molecular phenotypes.
776
777
778

779



780
781

**Figure 2:** Reference trait analysis reveals overall patterns of covariation between incompatible traits. (A) Relationship between canonical correlation and transitive reliability. To evaluate the mathematical relationship between these quantities, we simulated two vectors with known correlation to represent the canonical covariates, and calculated transitive reliability with real gene expression data. Canonical correlation shown is absolute value, and transitive reliability is sign-matched. (B) Sample size increases lead to higher and more precise transitive reliability. Plot shows transitive reliability estimated using anxiety data with animals subsampled as described in the main text. Sample size on *x*-axis indicates the number of individuals used in each of the training and testing groups (the number of individuals phenotyped for target traits and the number with high-dimensional molecular phenotypes, respectively). Black line indicates magnitude of first canonical correlation calculated from full dataset. Color indicates whether training and testing groups were fully or partially independent.
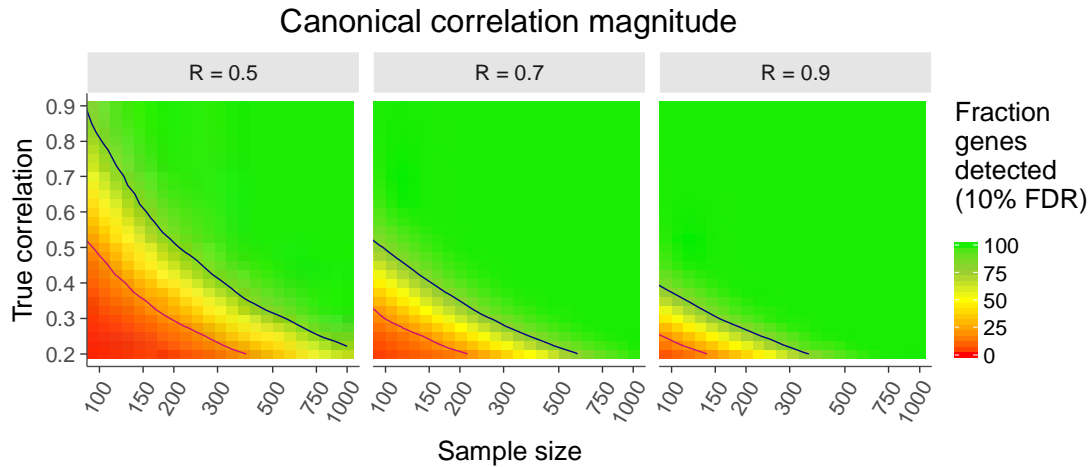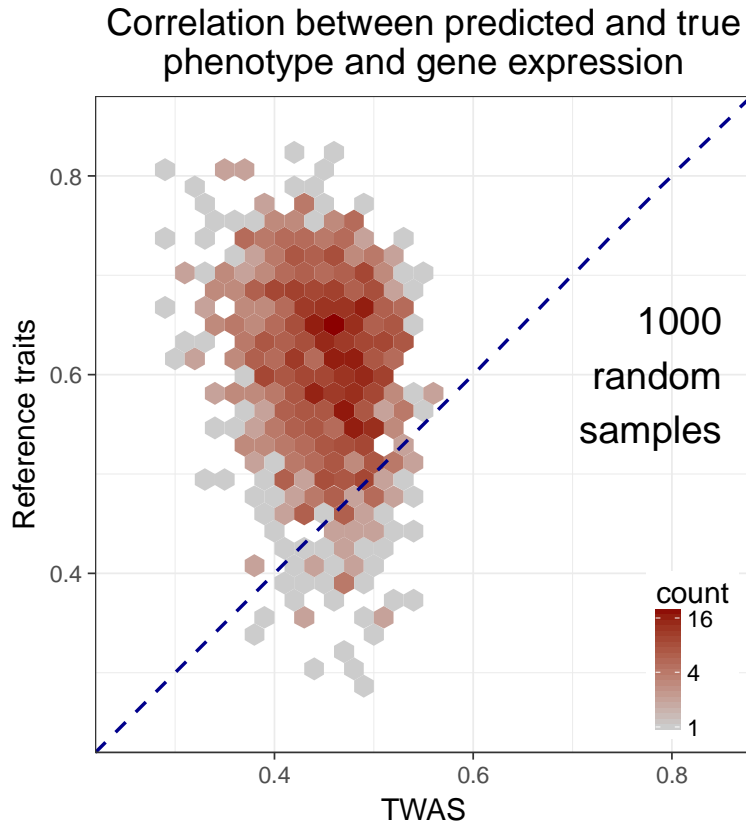
796
797
798

799



**Figure 3:** Reference trait analysis identifies simulated trait-gene expression correlations across a wide variety of parameter values. Sample size plotted along *x*-axis is the number of individuals used in each of the training and testing groups (equal sample size for the two groups, where the training group consists of individuals phenotyped for target traits and the testing group those with high-dimensional molecular phenotypes). True correlation (*y*-axis) indicates correlation between first target trait canonical variable ($\vec{v}_1$) and simulated gene expression. Facets indicate magnitude of canonical correlation coefficient between reference and target traits (*R* listed along grey strips, ±0.02). Navy and magenta contour lines depict regions above/below which trait-gene expression correlations are detected >80% and <20% of the time, respectively.

**Figure 4:** Reference trait analysis recovers true trait-gene expression correlations more accurately than TWAS. Binned hexagon plot shows the results of 1,000 random samples where the anxiety dataset was split into two halves randomly designated the training and testing groups. Reference trait analysis and TWAS were used to recover trait-gene expression correlations. The true values of both the trait and gene expression are known in this dataset, but were hidden when running reference trait analysis or TWAS. For each method, the correlation across all genes between predicted and true values was computed.

825

826    **Supplementary Table 1:** Anxiety-related traits measured on 258 Diversity Outbred

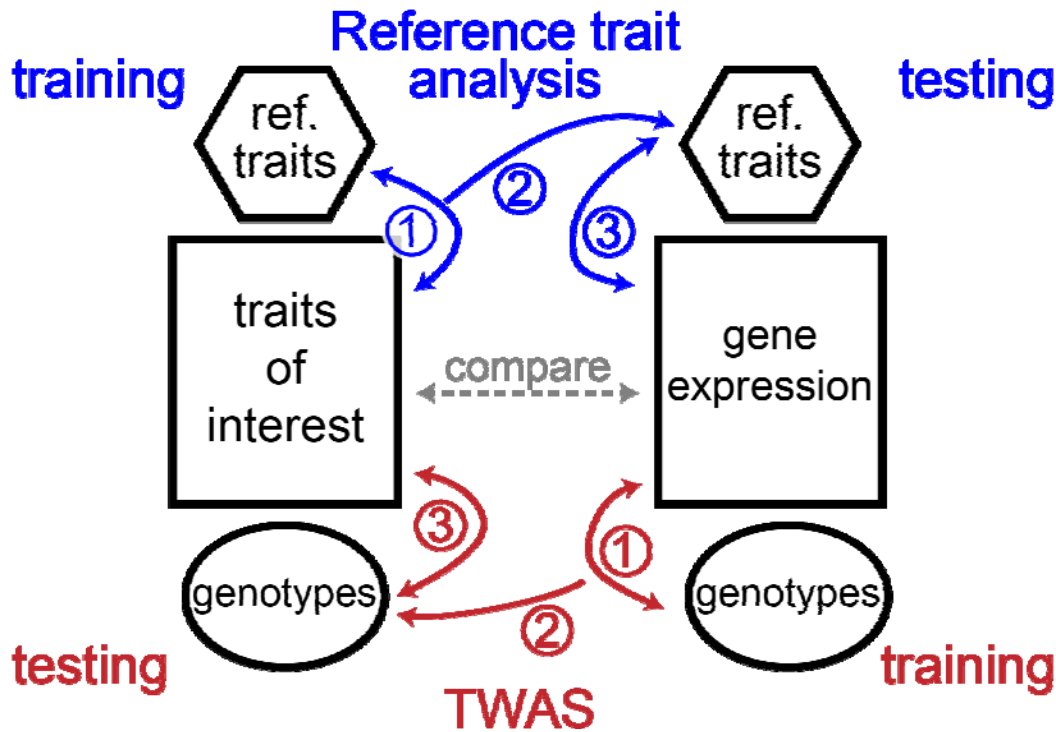827    mice used in case study of reference trait analysis.

828

| Group | Phenotype |
|---|---|
| Light-dark box | Distance traveled |
| Light-dark box | Light-dark transitions |
| Light-dark box | Percent time in light (first four minutes) |
| Light-dark box | Percent time in light (total) |
| Light-dark box | Percent time in light, slope |
| Open-field | Distance traveled (first four minutes) |
| Open-field | Distance traveled (total) |
| Open-field | Distance change (first – last) |
| Open-field | Percent time in corner |
| Open-field | Percent time in corner, slope |
| Open-field | Percent time in periphery |
| Open-field | Percent time in periphery, slope |
| Open-field | Percent time in center (square-root transformed) |
| Open-field | Percent time in center, slope |
| Open-field | Percent time mobile |
| Open-field | Fecal boli count |

829

830

831



832
833 **Supplementary Figure 1:** Schematic comparing overall strategies of reference trait
834 analysis and TWAS. For reference trait analysis, canonical correlation analysis is
835 used to relate traits of interest to reference traits (blue, 1) and coefficients derived
836 from this model are applied to reference traits in the cohort without measurements
837 of traits of interest (blue, 2). Finally, these projected reference traits are compared
838 to gene expression to identify trait-gene expression correlations (blue, 3). In the
839 TWAS approach, genotypes are used to build models that predict gene expression
840 through eQTL (red, 1). These models are applied to genotypes in the cohort without
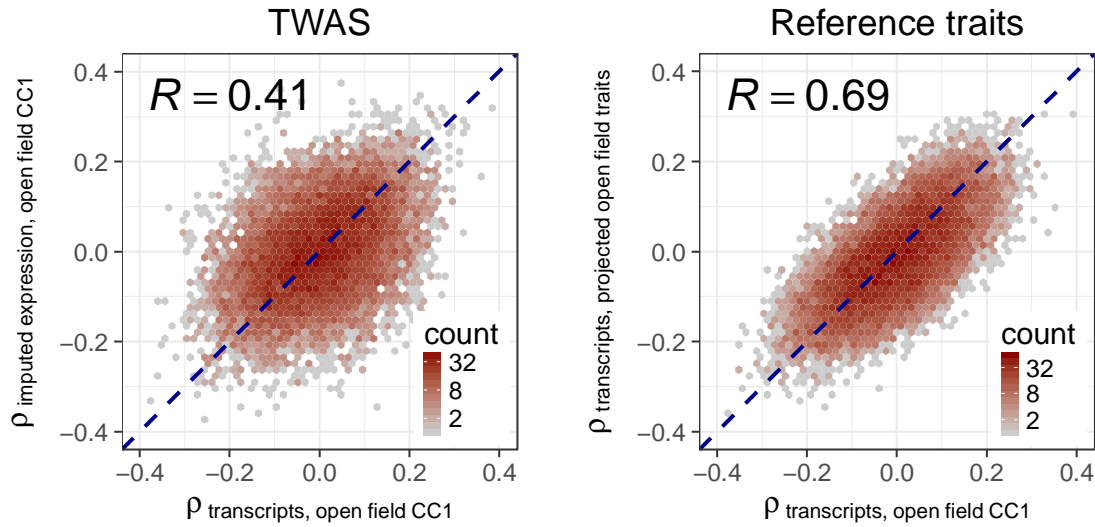841 gene expression measurements (red, 2) and imputed gene expression is compared
842 with traits of interest to identify trait-gene expression correlations (red, 3). Note
843 that training and testing cohort labels are switched for the two methods but that the
844 end result of each is to compare traits of interest with gene expression (grey dashed
845 line, middle).
846
847

848
849
850 **Supplementary Figure 2:** Comparison of TWAS and reference trait analysis using a
851 single random division of the mouse anxiety dataset. For both panels we take the
852 true trait of interest to be the first canonical covariate of open-field traits (open-field
853 CC1). For TWAS we used genotypes to impute gene expression. Left panel shows
854 correlation of individual transcripts to the trait of interest, where the $x$-axis plots
855 correlations based on true transcript abundance and the $y$-axis plots correlations
856 based on imputed transcript abundance. Right panel shows the analogous result but
857 using reference trait analysis, where gene expression is fixed and predictors of
858 open-field behavior are represented by projected traits.
859

860    **Supplementary Datasets**

861

862    **Supplementary Dataset 1:** Normalized hippocampal gene expression matrix. RNA-

863    Seq data were processed as described (Methods). To obtain normalized gene

864    expression matrix, raw counts in each sample were normalized to the upper quartile

865    value and transformed to normal scores.

866

867    **Supplementary Dataset 2:** Traits derived from open-field arena exploration assay

868    and used in case study of reference trait analysis. Supplementary Table 1 provides

869    basic information on phenotypes, while complete details of animal rearing,

870    husbandry and phenotyping are presented in Logan et al. (2013).

871

872    **Supplementary Dataset 3:** Traits derived from light-dark box behavior assay and

873    used in case study of reference trait analysis. Supplementary Table 1 provides basic

874    information on phenotypes, while complete details of animal rearing, husbandry

875    and phenotyping are presented in Logan et al. (2013).

876