

# UC Davis

## UC Davis Previously Published Works

### Title

TIGAR: An Improved Bayesian Tool for Transcriptomic Data Imputation Enhances Gene Mapping of Complex Traits.

### Permalink

<https://escholarship.org/uc/item/5dq47483>

### Journal

American Journal of Human Genetics, 105(2)

### Authors

Nagpal, Sini

Meng, Xiaoran

Epstein, Michael

et al.

### Publication Date

2019-08-01

### DOI

10.1016/j.ajhg.2019.05.018

Peer reviewed

# TIGAR: An Improved Bayesian Tool for Transcriptomic Data Imputation Enhances Gene Mapping of Complex Traits

Sini Nagpal,<sup>1,11</sup> Xiaoran Meng,<sup>2,3,11</sup> Michael P. Epstein,<sup>2,3</sup> Lam C. Tsoi,<sup>4</sup> Matthew Patrick,<sup>5</sup> Greg Gibson,<sup>1</sup> Philip L. De Jager,<sup>6</sup> David A. Bennett,<sup>7</sup> Aliza P. Wingo,<sup>8,9</sup> Thomas S. Wingo,<sup>3,10</sup> and Jingjing Yang<sup>3,\*</sup>

The transcriptome-wide association studies (TWASs) that test for association between the study trait and the imputed gene expression levels from *cis*-acting expression quantitative trait loci (*cis*-eQTL) genotypes have successfully enhanced the discovery of genetic risk loci for complex traits. By using the gene expression imputation models fitted from reference datasets that have both genetic and transcriptomic data, TWASs facilitate gene-based tests with GWAS data while accounting for the reference transcriptomic data. The existing TWAS tools like PrediXcan and FUSION use parametric imputation models that have limitations for modeling the complex genetic architecture of transcriptomic data. Therefore, to improve on this, we employ a nonparametric Bayesian method that was originally proposed for genetic prediction of complex traits, which assumes a data-driven nonparametric prior for *cis*-eQTL effect sizes. The nonparametric Bayesian method is flexible and general because it includes both of the parametric imputation models used by PrediXcan and FUSION as special cases. Our simulation studies showed that the nonparametric Bayesian model improved both imputation  $R^2$  for transcriptomic data and the TWAS power over PrediXcan when  $\geq 1\%$  *cis*-SNPs co-regulate gene expression and gene expression heritability  $\leq 0.2$ . In real applications, the nonparametric Bayesian method fitted transcriptomic imputation models for 57.8% more genes over PrediXcan, thus improving the power of follow-up TWASs. We implement both parametric PrediXcan and nonparametric Bayesian methods in a convenient software tool “TIGAR” (Transcriptome-Integrated Genetic Association Resource), which imputes transcriptomic data and performs subsequent TWASs using individual-level or summary-level GWAS data.

## Introduction

Genome-wide association studies (GWASs) have successfully identified thousands of genetic risk loci for complex traits. However, the majority of these loci are located within noncoding regions whose molecular mechanisms remain unknown.<sup>1–3</sup> Recent studies have shown that these associated regions were enriched for regulatory elements such as enhancers (H3K27ac marks)<sup>4,5</sup> and expression of quantitative trait loci (eQTL),<sup>6,7</sup> suggesting that the genetically regulated gene expression might play a key role in explaining the etiology of complex traits. Multiple studies have recently generated rich transcriptomic datasets for diverse tissues of the human body (besides genotype data), e.g., the Genotype-Tissue Expression (GTEx) project for >44 human tissues,<sup>6</sup> Genetic European Variation in Health and Disease (GEUVADIS) for lymphoblastoid cell lines,<sup>8</sup> Depression Genes and Networks (DGN) for whole-blood samples,<sup>9</sup> and the North American Brain Expression Consortium (NABEC) for cortex tissues.<sup>10</sup> Previous studies<sup>11–16</sup> have also shown that integrating transcriptomic data in GWASs can help identify functional loci.

The majority of GWAS projects do not profile transcriptomic data and thus cannot enable direct integrative analysis. However, existing studies<sup>11,12</sup> have shown that one can impute the genetically regulated gene expression (GReX) within such GWAS projects by using reference datasets like GTEx<sup>6</sup> and GEUVADIS<sup>8</sup> to train gene expression imputation models, and then test for the association between imputed GReX for GWAS samples and the trait of interest—referred to as transcriptome-wide association studies (TWASs).<sup>11,12</sup> Specifically, the gene expression imputation models are fitted by regressing assayed gene expression levels on *cis*-eQTL genotypes with reference dataset. For examples, the PrediXcan<sup>11</sup> method uses an Elastic-Net<sup>17</sup> variable selection model and the FUSION<sup>12</sup> tool implements a Bayesian sparse linear mixed model (BSLMM)<sup>18</sup> to estimate the *cis*-eQTL effect sizes with reference dataset. The estimated *cis*-eQTL effect sizes are then used to impute the GReX for GWAS samples.

In short, the Elastic-Net<sup>17</sup> model used by PrediXcan<sup>11</sup> assumes a combination of LASSO<sup>19</sup> ( $L_1$ ) and Ridge<sup>20</sup> ( $L_2$ ) penalties on the *cis*-eQTL effect sizes, which is equivalent to a Bayesian model with a mixture Gaussian and Laplace

<sup>1</sup>School of Biology, Georgia Institute of Technology, Atlanta, GA 30322, USA; <sup>2</sup>Department of Biostatistics and Bioinformatics, Emory University School of Public Health, Atlanta, GA 30322, USA; <sup>3</sup>Center for Computational and Quantitative Genetics, Department of Human Genetics, Emory University School of Medicine, Atlanta, GA 30322, USA; <sup>4</sup>Department of Dermatology; Department of Computational Medicine & Bioinformatics; Department of Biostatistics, University of Michigan, Ann Arbor, MI 48109, USA; <sup>5</sup>Department of Dermatology, University of Michigan Medical School, Ann Arbor, MI 48109, USA; <sup>6</sup>Medical Center Neurological Institute, Columbia University, New York, NY 10032, USA; <sup>7</sup>Rush Alzheimer's Disease Center, Rush University Medical Center, Chicago, IL 60612, USA; <sup>8</sup>Division of Mental Health, Atlanta VA Medical Center, Decatur, GA, USA; <sup>9</sup>Department of Psychiatry and Behavioral Sciences, Emory University School of Medicine, Atlanta, GA 30322, USA; <sup>10</sup>Department of Neurology, Emory University School of Medicine, Atlanta, GA 30322, USA

<sup>11</sup>These authors contributed equally to this work

\*Correspondence: [jingjing.yang@emory.edu](mailto:jingjing.yang@emory.edu)

<https://doi.org/10.1016/j.ajhg.2019.05.018>

© 2019 American Society of Human Genetics.



prior.<sup>21</sup> In contrast, the BSLMM<sup>18</sup> used by FUSION<sup>12</sup> is a combination of Bayesian variable selection model (BVSR)<sup>22</sup> and linear mixed model (LMM)<sup>23</sup> by assuming a normal mixture prior. Since a parametric prior is assumed for the *cis*-eQTL effect sizes by both Elastic-Net and BSLMM, it restricts the capability of PrediXcan and FUSION for handling the underlying complex genetic architecture of transcriptomes. Existing studies<sup>11,12</sup> have also shown that both PrediXcan<sup>11</sup> and FUSION<sup>12</sup> estimated the average regression  $R^2$  (i.e., the percentage of gene expression variation that can be explained by *cis*-genotypes) as ~5% for human whole-blood transcriptome, while the average genome-wide heritability of gene expression in human whole-blood transcriptome is estimated to be more than double that quantity.<sup>24,25</sup>

Therefore, to flexibly model *cis*-eQTL distributions, we use a nonparametric Bayesian method that was originally proposed for genetic prediction of complex traits,<sup>26</sup> where the prior for effect sizes is nonparametric and can be estimated from the data by assuming a Dirichlet process prior on effect-size variance. This Bayesian model is also known as latent Dirichlet process regression (DPR) model,<sup>26</sup> which can flexibly model the underlying complex genetic architecture of transcriptomes. Thus, DPR is a more generalized model that includes Elastic-Net (implemented in PrediXcan<sup>11</sup>) and BSLMM (implemented in FUSION<sup>12</sup>) as special cases. Consequently, DPR can robustly estimate *cis*-eQTLs and then improve imputation  $R^2$  (the squared Pearson correlation between the observed and imputed values on test samples). Moreover, a variational Bayesian algorithm<sup>26–28</sup> can be employed as an alternative of Monte Carlo Markov Chain (MCMC)<sup>29</sup> to efficiently fit the Bayesian model.

Similar to PrediXcan<sup>11</sup> and FUSION<sup>12</sup> methods, we employ DPR to estimate *cis*-eQTLs effect sizes from a reference dataset, which can then be used for downstream TWASs using either individual-level or summary-level GWAS data. In subsequent sections, we first describe the DPR<sup>26</sup> approach for estimating *cis*-eQTL effect sizes from a reference dataset and how we can then use these effect sizes for a downstream TWAS. We then compare the performance of DPR with PrediXcan using both simulated data and real GWAS and transcriptomic data from the Religious Orders Study and Rush Memory Aging Project (ROS/MAP)<sup>30–33</sup> for studying Alzheimer disease (AD).

Our in-depth simulation studies demonstrated that the DPR method obtained higher imputation  $R^2$  on test samples, when  $\geq 1\%$  *cis*-SNPs are true causal and the true expression heritability is  $\leq 0.2$ . Consequently, better imputation  $R^2$  resulted in improved power for follow-up association studies. Meanwhile, application of DPR to the ROS/MAP study imputed GReX for 57.8% more genes than PrediXcan. Using DPR, we also found a potentially associated gene *TRAPPC6A* for AD pathology indices, which was missed by PrediXcan. Further, by using the transcriptomic imputation models fitted from ROS/MAP data and summary-level GWAS data generated from the Inter-

national Genomics of Alzheimer's Project (IGAP),<sup>34</sup> we identified three known AD loci<sup>34–38</sup> that potentially affect the late-onset AD risk through transcript abundance. We conclude with a discussion of future topics and further describe our software tool TIGAR (Transcriptome-Integrated Genetic Association Resource) implementing both parametric Elastic-Net and nonparametric Bayesian DPR methods for public use.

## Material and Methods

Here, we briefly describe the underlying statistical model of gene-expression imputation. Consider the following linear regression model for estimating the *cis*-eQTL effect sizes from a reference study that has both genetic and transcriptomic data available,

$$\mathbf{E}_g = \mathbf{X}\mathbf{w} + \boldsymbol{\varepsilon}, \boldsymbol{\varepsilon} \sim \mathbf{N}(0, \sigma_\varepsilon^2 \mathbf{I}) \quad (\text{Equation 1})$$

where  $\mathbf{E}_g$  denotes the gene expression levels (after corrections for confounding covariates such as age, sex, and principal components) for gene  $g$ ,  $\mathbf{X}$  denotes the genotype matrix for all *cis*-genotypes (encoded as the number of minor alleles or genotype dosages),  $\mathbf{w}$  denotes the corresponding *cis*-eQTL effect-size vector, and  $\boldsymbol{\varepsilon}$  denotes the error term. The intercept term is dropped in Equation 1 for assuming both  $\mathbf{E}_g$  and  $\mathbf{X}$  are centered at 0. Generally, SNPs within 1 Mb of the flanking 5' and 3' ends (*cis*-SNPs) are included in this regression model and non-zero  $\hat{\mathbf{w}}$  will be used for follow-up analysis. The GReX will be imputed by

$$\widehat{\mathbf{GReX}} = \mathbf{X}_{new} \hat{\mathbf{w}},$$

with *cis*-SNP data  $\mathbf{X}_{new}$  for GWAS samples.

### Nonparametric Bayesian Method

Following the nonparametric Bayesian DPR model proposed in previous studies for genetic prediction of complex traits,<sup>26</sup> a normal prior  $N(0, \sigma_w^2)$  is assumed for the *cis*-eQTL effect sizes ( $w_i, i = 1, \dots, p$ ) and a Dirichlet process (DP) prior<sup>39</sup> is assumed for the effect-size variance  $\sigma_w^2$  (as in Equation 1):

$$\mathbf{w}_i \sim \mathbf{N}(0, \sigma_w^2), \sigma_w^2 \sim \mathbf{D}, \mathbf{D} \sim \mathbf{DP}(\mathbf{IG}(\mathbf{a}, \mathbf{b}), \xi). \quad (\text{Equation 2})$$

The prior distribution  $D$  deviates from the DP with base distribution as an inverse gamma (IG) distribution and concentration parameter  $\xi$ . Note that  $\sigma_w^2$  can be viewed as a latent variable and integrating out  $\sigma_w^2$  will induce a nonparametric prior distribution for  $w_i$ , which is equivalent to a DP normal mixture model,<sup>26–28</sup>

$$\begin{aligned} \mathbf{w}_i &\sim \sum_{k=0}^{+\infty} \pi_k \mathbf{N}(0, \sigma_k^2), \sigma_k^2 \sim \mathbf{IG}(\mathbf{a}_k, \mathbf{b}_k), \pi_k = \nu_k \prod_{l=0}^{k-1} (1 - \nu_l), \nu_k \\ &\sim \mathbf{Beta}(1, \xi). \end{aligned} \quad (\text{Equation 3})$$

Here, the nonparametric prior distribution on  $w_i$  is equivalently represented by a mixture normal prior that is a weighted sum of an infinitely number of normal distributions ( $\mathbf{N}(0, \sigma_k^2)$ ,  $k = 0, \dots, +\infty$ ), corresponding weight  $\pi_k$  is determined by  $(\nu_l, l = 0, \dots, k)$  with a Beta prior, and  $\xi$  in the Beta prior (the same concentration parameter as in Equation 2) determines the number of components with non-zero weights in the mixture normal prior. Conjugate hyper priors  $\xi \sim \text{Gamma}(a_\xi, b_\xi)$  and  $\sigma_\varepsilon^2 \sim \text{IG}(a_\varepsilon, b_\varepsilon)$  are assumed.

Generally, the hyper parameters  $a_k, b_k, a_\varepsilon, b_\varepsilon$  in the inverse gamma distributions can be set as 0.1 and  $(a_\xi, b_\xi)$  in the gamma distribution can be set as (1, 0.1) to induce non-informative priors for  $(\sigma_k^2, \sigma_\varepsilon^2, \xi)$ . That is, the parameters  $(\sigma_k^2, \sigma_\varepsilon^2, \xi)$  will be adaptively estimated from the data and the nonparametric prior on  $w_i$  will be data driven. The posterior estimates for  $\mathbf{w}$  can be obtained by the MCMC<sup>29</sup> or variational Bayesian algorithm,<sup>28,40</sup> from the following joint conditional posterior distribution

$$P(\mathbf{w}, \boldsymbol{\pi}, \nu, \xi, \sigma_\varepsilon^2 | \mathbf{E}_g, \mathbf{X}) \propto$$

$$P(\mathbf{E}_g | \mathbf{w}, \mathbf{X}, \sigma_\varepsilon^2) P(\mathbf{w} | \boldsymbol{\pi}, \sigma_1^2, \dots, \sigma_k^2, \dots) \left( \prod_{k=0}^{+\infty} P(\sigma_k^2 | a_k, b_k) \right) \times \\ P(\boldsymbol{\pi} | \nu) P(\nu | \xi) P(\xi | a_\xi, b_\xi) P(\sigma_\varepsilon^2 | a_\varepsilon, b_\varepsilon).$$

Particularly, the variational Bayesian algorithm<sup>28,40</sup> is an approximation for the MCMC<sup>29</sup> with greatly improved computational efficiency, which is also used in our tool. Please refer to the [Supplemental Material and Methods](#) for technical details of both MCMC sampling and variational inference algorithms for obtaining the Bayesian posterior estimates for the *cis*-eQTL effect sizes.

### Elastic-Net and BSLMM Methods

The Elastic-Net model<sup>17</sup> (used by PrediXcan<sup>11</sup>) estimates the *cis*-eQTL effect sizes  $\hat{\mathbf{w}}$  in Equation 1 with a combination of  $L_1$  (LASSO)<sup>19</sup> and  $L_2$  (Ridge)<sup>20</sup> penalties by

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} \left( \|\mathbf{E}_g - \mathbf{X}\mathbf{w}\|_2^2 + \lambda \left( \alpha \|\mathbf{w}\|_1 + \frac{1}{2}(1 - \alpha) \|\mathbf{w}\|_2^2 \right) \right),$$

where  $\|\cdot\|_2$  denotes  $L_2$  norm,  $\|\cdot\|_1$  denotes  $L_1$  norm,  $\alpha \in [0, 1]$  denotes the proportion of  $L_1$  penalty, and  $\lambda$  denotes the penalty parameter. Particularly, PrediXcan<sup>11</sup> takes  $\alpha = 0.5$  and tunes the penalty parameter  $\lambda$  by a 5-fold cross validation.

As pointed out by previous studies,<sup>17,21</sup> the Elastic-Net model is equivalent to a Bayesian model with a mixture Gaussian and Laplace (mixture normal) prior for  $\mathbf{w}$ , that is,  $p(\mathbf{w}) \propto \exp\left(-\lambda\left(\alpha\|\mathbf{w}\|_1 + \frac{1}{2}(1-\alpha)\|\mathbf{w}\|_2^2\right)\right)$ . In contrast, the BSLMM<sup>18</sup> assumes a mixture of two normal as the prior for *cis*-eQTL effect sizes,  $w_i \sim \pi N(0, (\sigma_1^2 + \sigma_2^2)) + (1 - \pi)N(0, \sigma_2^2)$ . That is, the BSLMM<sup>18</sup> assumes all *cis*-SNPs have at least a small effect, which are normally distributed with variance  $\sigma_2^2$ , and some proportion ( $\pi$ ) of *cis*-SNPs have an additional effect, normally distributed with variance  $\sigma_1^2$ . Particularly, with  $\sigma_2^2 = 0$ , BSLMM becomes BVSR,<sup>22</sup> and with  $\pi = 0$ , the BSLMM becomes the LMM.<sup>23</sup> Therefore, the DP normal mixture<sup>26–28</sup> as assumed by the DPR method includes the parametric (mixture normal) priors used by Bayesian Elastic-Net<sup>21</sup> and BSLMM<sup>18</sup> as special cases, which is the main reason why DPR is a more generalized model including Elastic-Net and BSLMM as special cases. This is also why the DPR method can robustly model complex genetic architecture and improve the imputation  $R^2$ .

### Association Study with Univariate Phenotype

Given individual-level GWAS data (genotype data  $\mathbf{X}_{new}$ , phenotype  $\mathbf{Y}$ , covariant matrix  $\mathbf{C}$ ) and *cis*-eQTL effect size estimates  $\hat{\mathbf{w}}$ , the follow-up TWAS (using a burden type gene-based test<sup>41</sup>) is to test the association between  $\mathbf{GReX} = \mathbf{X}_{new}\hat{\mathbf{w}}$  and  $\mathbf{Y}$  based on the following generalized linear regression model

$$\mathbf{f}(\mathbf{E}[\mathbf{Y} | \mathbf{X}, \mathbf{C}]) = \boldsymbol{\eta}\mathbf{C} + \boldsymbol{\beta}\mathbf{GReX}. \quad (\text{Equation 4})$$

Here,  $f(\cdot)$  is a pre-specified link function, which can be set as identity function for quantitative phenotype or set as logit function for dichotomous phenotype. The gene-based association test is equivalent to test  $H_0 : \boldsymbol{\beta} = 0$  in Equation 4.

If only summary-level GWAS data are available, we can take the same approach as implemented by the FUSION<sup>12</sup> method. Let  $\mathbf{Z}$  denote the vector of Z-scores generated by single variant tests (Wald, likelihood ratio, score tests, etc.) for all *cis*-SNPs. The burden Z-score for gene-based association test is defined as

$$\tilde{\mathbf{Z}} = \frac{\mathbf{Z}\hat{\mathbf{w}}}{\sqrt{\mathbf{Z}\mathbf{w}}} = \frac{\mathbf{Z}\hat{\mathbf{w}}}{\sqrt{\hat{\mathbf{V}}\mathbf{V}\mathbf{w}}}, \quad (\text{Equation 5})$$

where  $\mathbf{V}$  denotes the covariance matrix of analyzed SNPs that can be estimated from training data or reference panels such as 1000 Genomes Project<sup>42</sup> (of the same ethnicity).

### Association Study with Multivariate Phenotype

To test the association between multivariate phenotypes and imputed GReX of the focal gene, we take a similar approach as the MultiPhen method.<sup>43</sup> For example, consider two phenotypes  $(\mathbf{Y}_1, \mathbf{Y}_2)$  and a covariate matrix  $\mathbf{C}$ , we first adjust for the covariates by taking the residuals  $(\tilde{\mathbf{Y}}_1, \tilde{\mathbf{Y}}_2)$  respectively from the linear regression models  $\mathbf{Y}_j = \boldsymbol{\eta}\mathbf{C} + \boldsymbol{\varepsilon}$ ,  $j = 1, 2$ . Then we test whether the regression  $R^2$  is significantly greater than zero ( $H_0 : R^2 = 0$ ) for the following regression model

$$\mathbf{GReX}_g = \boldsymbol{\beta}_1 \tilde{\mathbf{Y}}_1 + \boldsymbol{\beta}_2 \tilde{\mathbf{Y}}_2 + \boldsymbol{\varepsilon}. \quad (\text{Equation 6})$$

That is, we test whether the multivariate phenotypes can jointly explain a non-zero percentage of variance in the imputed GReX. The p value can be calculated by using the F-statistic for the regression  $R^2$  in Equation 6.

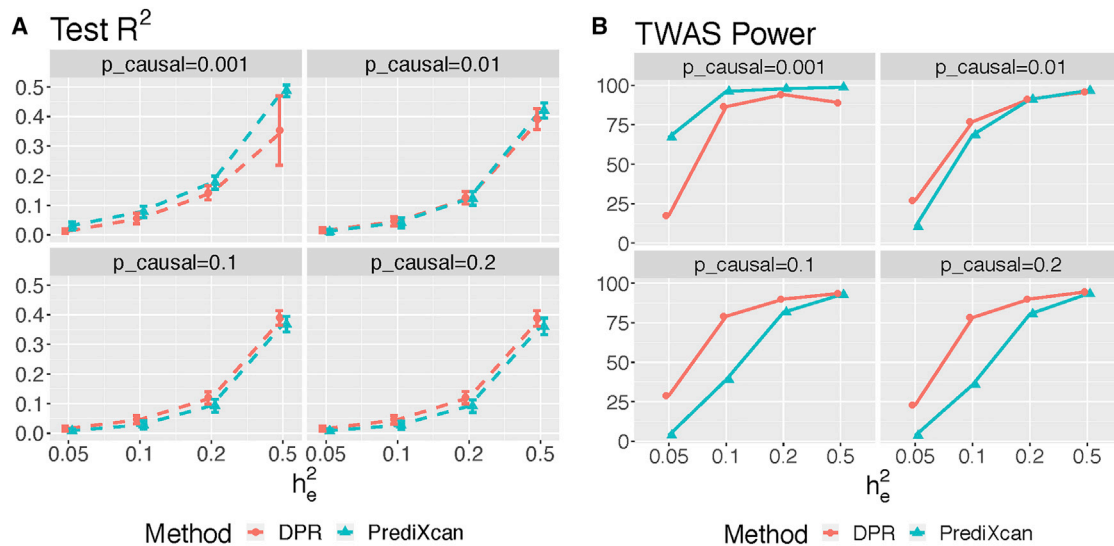
Even when only summary-level GWAS data are available, we can first obtain a burden Z-score per phenotype from Equation 5, i.e.,  $\tilde{\mathbf{Z}} = (\tilde{\mathbf{Z}}_1, \tilde{\mathbf{Z}}_2)$  with two phenotypes. Then, a similar burden approach can be used to obtain a joint Z-score for multi-phenotype test,

$$\tilde{\mathbf{Z}}_{\text{joint}} = \frac{\tilde{\mathbf{Z}}\mathbf{J}}{\sqrt{\tilde{\mathbf{Z}}\mathbf{J}}} = \frac{\tilde{\mathbf{Z}}\mathbf{J}}{\sqrt{\mathbf{J}'\mathbf{V}\mathbf{J}}}, \mathbf{J} = (1, \dots, 1)',$$

where  $\mathbf{V}_Y$  is the covariance matrix among multiple traits.

### Simulation Study Design

We conducted in-depth simulation studies to compare the performance of both PrediXcan and DPR methods with respect to imputation  $R^2$  in the test data and the power of TWASs. Specifically, we used data from 499 ROS/MAP participants<sup>44</sup> which contains both RNA-sequencing and genotype data as training data, and genotype data from an additional 1,200 ROS/MAP participants<sup>44</sup> as test data. The test sample size (1,200) was chosen arbitrarily (randomly selected from the ROS/MAP study) to be comparable with the sample size (1,164) in the real association study of AD pathology indices. The genotyped and imputed genetic data for 2,799 *cis*-SNPs (with minor allele frequency (MAF) > 5% and Hardy-Weinberg p value >  $10^{-5}$ ) of the arbitrarily chosen gene *ABCA7* (see Figure S1 for the LD block structure) were used to simulate gene expression levels.



**Figure 1. Performance Comparison of DPR versus PrediXcan**

Plots of average imputation  $R^2$  (A) and TWAS power (B) in test samples by DPR and PrediXcan, with various proportions of true causal SNPs  $p_{\text{causal}} = (0.001, 0.01, 0.1, 0.2)$  and true expression heritability  $h_e^2 = (0.05, 0.1, 0.2, 0.5)$ . TWAS power was evaluated with paired expression and phenotype heritability  $(h_e^2, h_p^2) = ((0.05, 0.8), (0.1, 0.5), (0.2, 0.25), (0.5, 0.1))$ .

We performed comprehensive scenarios that varied the proportion of causal SNPs (out of 2,799 SNPs, influenced gene expression) among values in the vector  $p_{\text{causal}} = (0.001, 0.01, 0.1, 0.2)$ . We varied the proportion of gene expression variance explained by causal SNPs (i.e., expression heritability), along with the proportion of phenotypic variance explained by simulated gene expression levels (i.e., phenotypic heritability), among values in the vector  $(h_e^2, h_p^2) = ((0.05, 0.8), (0.1, 0.5), (0.2, 0.25), (0.5, 0.1))$ . The phenotypic heritability was selected arbitrarily with respect to expression heritability such that the follow-up association study power fell within the range of (25%, 85%). We also considered various training sample sizes (100, 300, 499) for simulation scenario with  $p_{\text{causal}} = 0.2$  and  $(h_e^2, h_p^2) = (0.2, 0.25)$ .

With genotype matrix  $\mathbf{X}_g$  of the randomly selected causal SNPs (according to  $p_{\text{causal}}$ ), we generated effect sizes  $w_i$  from  $N(0,1)$  and then re-scaled the effect sizes to ensure the targeted  $h_e^2$ . Gene expression levels were generated by  $\mathbf{E}_g = \mathbf{X}_g \mathbf{w} + \epsilon$ , with  $\epsilon \sim N(0, (1 - h_e^2))$ . Then the phenotype values were generated by  $\mathbf{Y} = \beta \mathbf{E}_g + \epsilon$ , where  $\beta$  was selected with respect to  $h_p^2$  and  $\epsilon \sim N(0, (1 - h_p^2))$ .

For each scenario, we repeated simulations for 1,000 times, where we applied both PrediXcan<sup>11</sup> and DPR methods to obtain imputation models with training samples, impute the GRex for test samples, and then conduct follow-up association studies using the imputed GRex. We did not compare with FUSION<sup>12</sup> using BSLMM because of the computational burden of estimating *cis*-eQTL effect sizes by MCMC (~2 h per gene). The association study power was calculated as the proportion of 1,000 repeated simulations with p value  $< 2.5 \times 10^{-6}$  (genome-wide significance threshold adjusting for testing 20K independent genes).

### ROS/MAP Data

Samples in the ROS/MAP data were collected from participants of the Religious Orders Study (ROS) and the Rush Memory and Aging Project (MAP), which are prospective cohort studies of studying aging and dementia.<sup>30,31,33</sup> The ROS/MAP study recruited senior adults without known dementia at enrollment who underwent

annual clinical evaluation. Brain autopsy was done at the time of death for each participant. All participants signed an informed consent and Anatomic Gift Act, and the studies were approved by the Institutional Review Board of Rush University Medical Center, Chicago, IL. Specifically, microarray genotype data generated for 2,093 European-descent participants<sup>44</sup> were further imputed to the 1000 Genomes Project Phase 3<sup>42</sup> in our analysis. The post-mortem brain samples (gray matter of the dorsolateral prefrontal cortex) from ~30% these participants were profiled for transcriptomic data by next-generation RNA sequencing.<sup>45</sup> In this paper, we conducted TWASs for two important indices of AD pathology that were quantified with  $\beta$ -antibody specific immunostains:<sup>30,31,33</sup> neurofibrillary tangle density (tangles) with stereology and  $\beta$ -amyloid load (amyloid) with image analysis. The neurofibrillary tangle density quantifies the average Tau tangle density within two or more 20  $\mu\text{m}$  sections from eight brain regions—hippocampus, entorhinal cortex, midfrontal cortex, inferior temporal, angular gyrus, calcarine cortex, anterior cingulate cortex, and superior frontal cortex. The  $\beta$ -amyloid load quantifies the average percent area of cortex occupied by  $\beta$ -amyloid protein in adjacent sections from the same eight brain regions.

## Results

### Simulation Studies

In the simulation studies, we observed that the DPR method performed robustly with respect to different causal proportions and gene expression heritability. Specifically, when  $p_{\text{causal}} > 0.01$  DPR outperformed PrediXcan across all expression heritability values, giving higher imputation  $R^2$  in test data (Figure 1A). For example, when  $p_{\text{causal}} = 0.2$ , the average imputation  $R^2$  of 1,000 simulations was estimated as 4.55% by using DPR versus 2.64% by using PrediXcan with  $h_e^2 = 0.1$ , while the average imputation  $R^2$  was estimated as 12.02% by using DPR versus 9.13% by



**Table 1. Simulation Prediction  $R^2$  Comparison**

$h_e^2$	Causal Proportion 0.01		Causal Proportion 0.2	
	DPR	PrediXcan	DPR	PrediXcan
0.05	1.60%*	1.12%	1.54%*	0.76%
0.1	4.54%*	4.13%	4.55%*	2.64%
0.2	12.54%*	12.29%	12.02%*	9.13%
0.5	39.31%	42.05%*	38.78%*	36.04%

Various simulation scenarios were considered, with the proportion of true causal SNPs  $p_{\text{causal}} = (0.01, 0.2)$  and expression heritability  $h_e^2 = (0.05, 0.1, 0.2, 0.5)$ . The best prediction  $R^2$  per scenario is indicated with asterisk (\*).

using PrediXcan with  $h_e^2 = 0.2$  (Table 1). When  $p_{\text{causal}} = 0.01$ , DPR performed slightly outperformed PrediXcan with  $h_e^2 = (0.05, 0.1, 0.2)$  and PrediXcan outperformed DPR with  $h_e^2 = 0.5$  (Table 1, Figure 1). On the other hand, under a sparse *cis*-eQTL causality model with  $p_{\text{causal}} = 0.001$  (i.e., with 3 true causal *cis*-eQTL), the Elastic-Net method resulted in higher imputation  $R^2$  and TWAS power on test data (Figure 1).

Consequently, when  $p_{\text{causal}} \geq 0.01$  and  $h_e^2 \leq 0.2$ , the power of association studies was higher by using DPR than using PrediXcan imputation models (Figure 1B). When  $h_e^2 = 0.5$ , using both imputation models led to comparable power for association studies (Figure 1B). Even though both methods had similar over-estimated training  $R^2$  (Figure S2), the DPR method resulted in higher imputation  $R^2$  for test data (Table 1; Figures 1A) and higher power for association studies under *cis*-eQTL causality models with  $p_{\text{causal}} \geq 0.01$  and  $h_e^2 \leq 0.2$  (Figure 1B). In addition, from the simulation studies with various training sample sizes (100, 300, 499),  $p_{\text{causal}} = 0.2$ , and  $(h_e^2, h_p^2) = (0.2, 0.25)$ , the imputation  $R^2$  and TWAS power increases as sample size increases while the DPR method consistently outperforms PrediXcan (Figure 2). Overall, these results demonstrated the advantages of the DPR method for modeling the complex genetic architecture of transcriptomes, especially when the causal proportions  $\geq 0.01$  and the expression heritability  $\leq 0.2$ .

### Real Applications to ROS/MAP Data

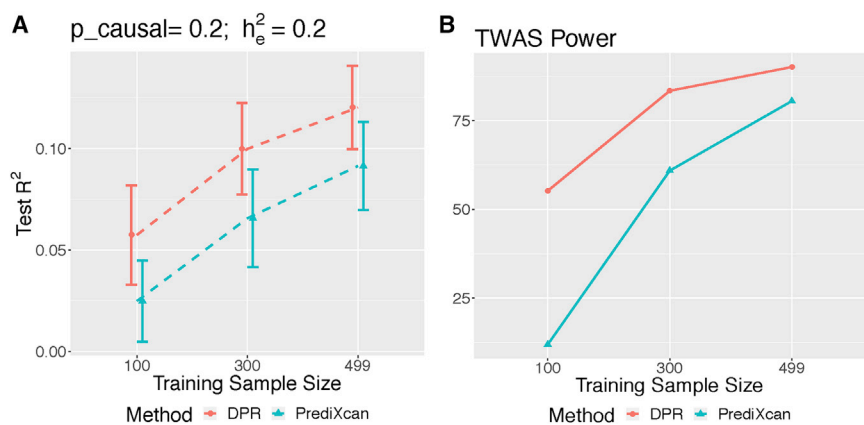
To illustrate the performance of the DPR method in real studies, we applied both DPR and PrediXcan on the ROS/MAP data (see Material and Methods). We trained the gene expression imputation models using 499 samples that have both transcriptomic data for prefrontal cortex tissues and genotype data (imputed to 1000 Genomes Phase 3, with MAF > 5%, Hardy-Weinberg p value >  $10^{-5}$ , and genotype imputation  $R^2 > 0.3$ ). A total of 15,583 genes had gene expression levels after standard RNA-sequencing quality control. The gene expression levels were first adjusted for age at death, sex, postmortem interval, study (ROS or MAP), batch effects, RNA integrity number scores, and cell type proportions (with respect to oligodendrocytes, astrocytes, microglia, neurons) by linear

regression models. For each gene, *cis*-SNPs within the 1 Mb of the flanking 5' and 3' ends were used in the imputation models as predictors.

First, we compared transcriptome-wide 5-fold cross validation (CV) regression  $R^2$  estimated by using both DPR and PrediXcan methods. Specifically, we randomly split 499 training samples into 5 folds, where the imputation  $R^2$  of each fold was calculated using the model trained with the other 4-fold samples. If the training model is null, we take the imputation  $R^2$  as 0 and take the average imputation  $R^2$  across all 5-fold test samples as 5-fold CV  $R^2$ . The transcriptome-wide median of 5-fold CV  $R^2$  is 0.013 by DPR versus 0.005 by PrediXcan. The 5-fold CV  $R^2$  was used as the criterion for selecting significant imputation models ( $R^2 > 0.01$  as used by previous studies<sup>11,46</sup>). From Figure 3A, we can see that the DPR method obtained more imputation models and higher imputation  $R^2$  when 5-fold CV  $R^2$  is in the range of (0.01, 0.05), which is also consistent with our simulation studies. Overall, the DPR method obtained significant imputation models for 8,752 genes versus 5,547 genes by PrediXcan (with 57.8% increases). Thus, the DPR method featuring data-driven nonparametric prior for the *cis*-eQTL is preferred in real studies for identifying more genes with imputable expression levels.

Second, to investigate how both DPR and PrediXcan methods perform in real studies with independent prediction cohort, we used the ROS cohort (256 samples) to train gene expression imputation models and then used the MAP cohort (243 samples) as a test dataset. Specifically, we compared the median prediction  $R^2$  by both DPR and PrediXcan with MAP test cohort. As shown in Table 2, the DPR method obtained higher median prediction  $R^2$  than PrediXcan among 8,752 genes that have 5-fold CV  $R^2 > 0.01$  by DPR (0.011 versus 0.003), performed similarly as PrediXcan among 5,547 genes that have 5-fold CV  $R^2 > 0.01$  by PrediXcan (0.026 versus 0.026), obtained slightly lower median prediction  $R^2$  among 4,819 genes that have 5-fold CV  $R^2 > 0.01$  by both DPR and PrediXcan (0.033 versus 0.036). These results are also consistent with our simulation results and 5-fold cross validation results with ROS/MAP data. That is, PrediXcan method is preferred for genes with sparse causal eQTL that have relatively large effect sizes, whereas DPR is preferred for genes with less sparse causal eQTL that have minor effect sizes due to low expression heritability.

Third, we used all 499 training samples to fit imputation models for genes with respective 5-fold CV  $R^2 > 0.01$  by both DPR and PrediXcan, and then used these models to impute the GReX for all GWAS samples. We conducted univariate phenotype association studies (Material and Methods) using all GWAS samples ( $n = 1,164$ ) that have the AD pathology indices (neurofibrillary tangle density and  $\beta$ -amyloid load, with Pearson correlation 0.48) quantified. Possible confounding covariates including age at death, sex, study (ROS or MAP), smoking, education, and first three genotype principle components were adjusted



**Figure 2. Performance of DPR and PrediXcan with Respect to Various Training Sample Sizes**

Test  $R^2$  (A) and TWAS power (B) from simulation studies with causal proportion  $p_{\text{causal}} = 0.2$ , expression heritability and phenotype heritability  $(h_e^2, h_p^2) = (0.2, 0.25)$ , and various training sample sizes (100, 300, 499).

in the association studies. Interestingly, the association studies for both AD pathology indices using the DPR imputation models identified the same top significant gene *TRAPPC6A* (within the 2 Mb region from the major risk gene *APOE*, encoding apolipoprotein E, but independent of *APOE*) with  $p$  values  $1.64 \times 10^{-5}$  and  $5.35 \times 10^{-5}$  (Figures S3A and S4A). Moreover, the multivariate phenotype association studies (Material and Methods) for both AD pathology indices identified *TRAPPC6A* as the most significant gene with  $p$  value  $5.81 \times 10^{-6}$  and FDR 0.08 (Figure 3C). On the other hand, the PrediXcan failed to obtain a transcriptomic imputation model for *TRAPPC6A* (Figures S3B, S4B, and S6). Quantile-quantile plots for these TWAS  $p$  values were presented in Figure S5.

In addition, for 14 known common and rare loci of late-onset AD<sup>34–38</sup> with significant imputation models, we conducted association studies using transcriptomic imputation models (DPR and PrediXcan) fitted from ROS/MAP data and summary-level GWAS data from IGAP.<sup>34</sup> Using the imputation models fit by DPR, we identified three significant loci with FDR < 0.05 (Figure 3B)—*ADAM10*, *CD2AP*, and *TREM2*—that potentially affect late-onset AD risk through transcriptomic changes. Here, *TREM2* was also identified by using the PrediXcan imputation model (Figure 3B). Particularly, the PrediXcan method imputed GReX for only 5 out of these 14 loci. In summary, these results show that the DPR method has superior power for follow-up TWASs.

## Discussion

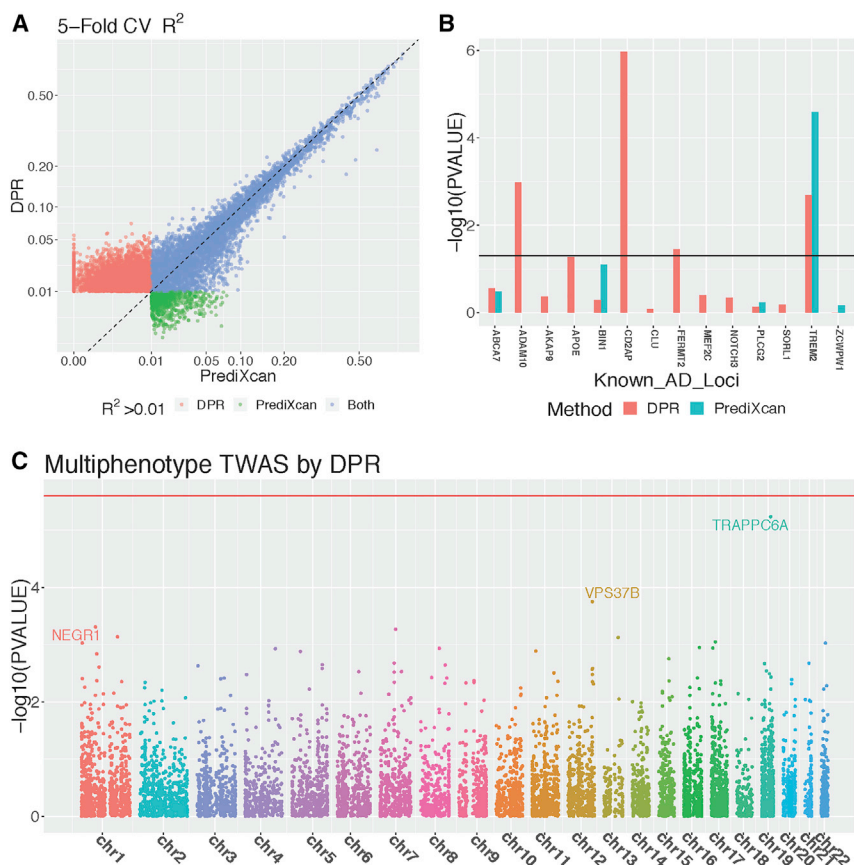
In this paper, by both in-depth simulations and real applications using individual-level ROS/MAP<sup>30–33</sup> and summary-level IGAP<sup>34</sup> GWAS data, we demonstrated that the nonparametric Bayesian DPR method is preferred for imputing gene expression when the proportion of causal *cis*-eQTL  $\geq 0.01$  and the true gene expression heritability  $\leq 0.2$ . The advantage of DPR model is due to the flexible nonparametric modeling of *cis*-eQTL effect sizes that results in improved imputation  $R^2$  for gene expression levels and higher power for TWASs. Here, we provide an

integrated tool (freely available on GITHUB), referred as Transcriptome-Integrated Genetic Association Resource (TIGAR), which integrates both parametric Elastic-Net and non-parametric Bayesian DPR models as two options for transcriptomic data imputation, along with TWAS options using individual-level and summary-level GWAS data for univariate and multi-variate phenotypes. TIGAR also conducts 5-fold cross validation by default and output significant imputation models with  $CV R^2 > 0.01$ .

With respect to user-friendly interface and computational efficiency, TIGAR can (1) take standard input files such as genotype files in VCF and dosage formats, phenotype files in PED format, and a combined text file for gene annotations and expression levels; (2) load input data per gene by TABIX for memory efficiency; (3) filter SNPs based on input thresholds of MAF and Hardy-Weinberg  $p$  value; (4) provide options of training both Elastic-Net (use Python3 scripts) and DPR (generate input files and call the executable tool developed with C++<sup>26</sup>) imputation models with unified output format; and (5) implement multi-threaded computation to take full advantage of multi-core clusters. These features make TIGAR a preferred tool for saving tedious data preparation and computation time for users. For example, TIGAR can complete training imputation models for  $\sim 20K$  genes and  $\sim 1K$  samples within  $\sim 20$  h and TWAS within  $\sim 1$  h with a 2.4 GHz 16-core CPU.

It is important to notice that imputing GReX with *cis*-eQTL effect sizes estimated from a training dataset is analogous to the idea of estimating polygenic risk scores (PRSs).<sup>47</sup> Even though studies of population heterogeneity are lacked for imputing GReX, the same philosophy of estimating PRSs still applies because of the same underlying statistical models. That is, given both genetic and transcriptomic heterogeneities across different populations, one needs to be cautious not using training dataset of a different ethnicity for a TWAS.<sup>47</sup>

As observed in the real ROS/MAP studies, there remains a large gap between the 5-fold CV  $R^2$  using *cis*-eQTL predictors ( $\sim 5\%$ ) and the average genome-wide heritability of gene expression levels (21.8% estimated by GCTA<sup>48</sup> based on a LMM). This is likely due to the large *trans*-acting contribution to transcript abundance documented for most genes. Thus, we hypothesize that it is promising to further improve the imputation  $R^2$  by fitting



**Figure 3. TWAS Results of Studying Alzheimer's Disease**

Transcriptome-wide 5-fold cross validation  $R^2$  (A) by PrediXcan and DPR with 499 ROS/MAP training samples, with different colors denoting whether the imputation  $R^2 > 0.01$  by DPR, PrediXcan, or both methods (genes with  $R^2 > 0.01$  by both DPR and PrediXcan were excluded from the plot). TWAS results (B) at known AD loci using GWAS summary-level statistics from IGAP and imputation models fitted from ROS/MAP data, where missing values are due to NULL imputation models by PrediXcan. Manhattan plot (C) for the multiphenotype TWAS (with neurofibrillary tangle density and  $\beta$ -amyloid load), using individual-level ROS/MAP data.

rating environmental contributions. The imputed transcript abundance levels can then be used for gene network analysis, differential gene expression analysis, and transcriptome mediation analysis with GWAS data. Validation of transcriptomic prediction accuracy in independent datasets will be critical in this regard, but unfortunately multiple large and similar datasets are not yet generally available for tissues other than peripheral blood.

transcriptomic imputation models with genome-wide variants as predictors. Scalable Bayesian inference techniques such as the Expectation Maximization MCMC (EM-MCMC) algorithm<sup>49</sup> are required for incorporating genome-wide variants.

Another limitation of existing TWAS methods is that the uncertainty of *cis*-eQTL effect-size estimates has not been taken into account in the association studies. A Bayesian framework can also be derived by taking the standard errors of these *cis*-eQTL effect-size estimates as prior standard deviations, which is part of our continuing research.

Besides the follow-up gene-based association studies (i.e., TWASs) described in this paper, the transcriptomic imputation models can be further extended by incorpo-

rating environmental contributions. In conclusion, we expect our work will provide a convenient and improved tool for transcriptomic imputation using the currently available rich reference datasets, as well as enhanced gene mapping for better understanding the genetic etiology of complex traits.

### Supplemental Data

Supplemental Data can be found online at <https://doi.org/10.1016/j.ajhg.2019.05.018>.

### Acknowledgments

J.Y. was supported by the startup funding from Department of Human Genetics at Emory University School of Medicine. A.P.W. and T.S.W. were supported by National Institutes of Health (NIH) R01AG056533. M.P.E. was supported by NIH R01GM11796. L.C.T. was supported by the Dermatology Foundation, the Arthritis National Research Foundation, the National Psoriasis Foundation, and NIH K01AR072129. ROS/MAP study data were provided by the Rush Alzheimer's Disease Center, Rush University Medical Center, Chicago, IL. Data collection was supported through funding by NIA grants P30AG10161, R01AG15819, R01AG17917, R01AG30146, R01AG36836, U01AG32984, and U01AG46152, the Illinois Department of Public Health, and the Translational Genomics Research Institute. In addition, we thank Thanneer Perumal and Benjamin Logsdon for performing quality control of the ROS/MAP RNA-sequencing data and for creating the brain cell type proportions.

**Table 2. Real Study Prediction  $R^2$  Comparison**

Number of Genes	DPR	PrediXcan
8,752 <sup>a</sup>	0.011	0.003
5,547 <sup>b</sup>	0.026	0.026
4,819 <sup>c</sup>	0.033	0.036

Median prediction  $R^2$  in MAP test cohort by using imputation models trained with ROS cohort with both DPR and PrediXcan methods.

<sup>a</sup>Genes that have 5-fold CV  $R^2 > 0.01$  by DPR.

<sup>b</sup>Genes that have 5-fold CV  $R^2 > 0.01$  by PrediXcan.

<sup>c</sup>Genes that have 5-fold CV  $R^2 > 0.01$  by both DPR and PrediXcan.



## Declaration of Interests

The authors declare no competing interests.

Received: December 27, 2018

Accepted: May 23, 2019

Published: June 20, 2019

## Web Resources

FUSION, <http://gusevlab.org/projects/fusion/>

IGAP data, [http://web.pasteur-lille.fr/en/recherche/u744/igap/igap\\_download.php](http://web.pasteur-lille.fr/en/recherche/u744/igap/igap_download.php)

PrediXcan, <https://github.com/hakyim/PrediXcan>

RADC Research Resource Sharing Hub, <http://www.radc.rush.edu/>

ROS/MAP data, <https://www.synapse.org/#!Synapse:syn3219045>

TIGAR, <https://github.com/yanglab-emory/TIGAR>

## References

1. Visscher, P.M., Brown, M.A., McCarthy, M.I., and Yang, J. (2012). Five years of GWAS discovery. *Am. J. Hum. Genet.* *90*, 7–24.
2. McCarthy, M.I., Abecasis, G.R., Cardon, L.R., Goldstein, D.B., Little, J., Ioannidis, J.P., and Hirschhorn, J.N. (2008). Genome-wide association studies for complex traits: consensus, uncertainty and challenges. *Nat. Rev. Genet.* *9*, 356–369.
3. Huang, Q. (2015). Genetic study of complex diseases in the post-GWAS era. *J. Genet. Genomics* *42*, 87–98.
4. Farh, K.K., Marson, A., Zhu, J., Kleinewietfeld, M., Housley, W.J., Beik, S., Shores, N., Whitton, H., Ryan, R.J., Shishkin, A.A., et al. (2015). Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature* *518*, 337–343.
5. Tsoi, L.C., Stuart, P.E., Tian, C., Gudjonsson, J.E., Das, S., Zawistowski, M., Ellinghaus, E., Barker, J.N., Chandran, V., Dand, N., et al. (2017). Large scale meta-analysis characterizes genetic architecture for common psoriasis associated variants. *Nat. Commun.* *8*, 15382.
6. Battle, A., Brown, C.D., Engelhardt, B.E., Montgomery, S.B.; GTEx Consortium; Laboratory, Data Analysis & Coordinating Center (LDACC)—Analysis Working Group; Statistical Methods groups—Analysis Working Group; Enhancing GTEx (eGTEx) groups; NIH Common Fund; NIH/NCI; NIH/NHGRI; NIH/NIMH; NIH/NIDA; Biospecimen Collection Source Site—NDRI; Biospecimen Collection Source Site—RPCI; Biospecimen Core Resource—VARI; Brain Bank Repository—University of Miami Brain Endowment Bank; Leidos Biomedical—Project Management; ELSI Study; Genome Browser Data Integration & Visualization—EBI; Genome Browser Data Integration & Visualization—UCSC Genomics Institute, University of California Santa Cruz; Lead analysts; Laboratory, Data Analysis & Coordinating Center (LDACC); NIH program management; Biospecimen collection; Pathology; and eQTL manuscript working group (2017). Genetic effects on gene expression across human tissues. *Nature* *550*, 204–213.
7. Nicolae, D.L., Gamazon, E., Zhang, W., Duan, S., Dolan, M.E., and Cox, N.J. (2010). Trait-associated SNPs are more likely to be eQTLs: annotation to enhance discovery from GWAS. *PLoS Genet.* *6*, e1000888.
8. Lappalainen, T., Sammeth, M., Friedländer, M.R., 't Hoen, P.A., Monlong, J., Rivas, M.A., González-Porta, M., Kurbatova, N., Griebel, T., Ferreira, P.G., et al.; Geuvadis Consortium (2013). Transcriptome and genome sequencing uncovers functional variation in humans. *Nature* *501*, 506–511.
9. Battle, A., Mostafavi, S., Zhu, X., Potash, J.B., Weissman, M.M., McCormick, C., Haudenschild, C.D., Beckman, K.B., Shi, J., Mei, R., et al. (2014). Characterizing the genetic basis of transcriptome diversity through RNA-sequencing of 922 individuals. *Genome Res.* *24*, 14–24.
10. Gibbs, J.R., van der Brug, M.P., Hernandez, D.G., Traynor, B.J., Nalls, M.A., Lai, S.L., Arepalli, S., Dillman, A., Rafferty, I.P., Troncoso, J., et al. (2010). Abundant quantitative trait loci exist for DNA methylation and gene expression in human brain. *PLoS Genet.* *6*, e1000952.
11. Gamazon, E.R., Wheeler, H.E., Shah, K.P., Mozaffari, S.V., Aquino-Michaels, K., Carroll, R.J., Eyler, A.E., Denny, J.C., Nicolae, D.L., Cox, N.J., Im, H.K.; and GTEx Consortium (2015). A gene-based association method for mapping traits using reference transcriptome data. *Nat. Genet.* *47*, 1091–1098.
12. Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B.W., Jansen, R., de Geus, E.J., Boomsma, D.I., Wright, F.A., et al. (2016). Integrative approaches for large-scale transcriptome-wide association studies. *Nat. Genet.* *48*, 245–252.
13. Zhu, Z., Zhang, F., Hu, H., Bakshi, A., Robinson, M.R., Powell, J.E., Montgomery, G.W., Goddard, M.E., Wray, N.R., Visscher, P.M., and Yang, J. (2016). Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. *Nat. Genet.* *48*, 481–487.
14. Mancuso, N., Shi, H., Goddard, P., Kichaev, G., Gusev, A., and Pasaniuc, B. (2017). Integrating Gene Expression with Summary Association Statistics to Identify Genes Associated with 30 Complex Traits. *Am. J. Hum. Genet.* *100*, 473–487.
15. Su, Y.R., Di, C., Bien, S., Huang, L., Dong, X., Abecasis, G., Berndt, S., Bezieau, S., Brenner, H., Caan, B., et al. (2018). A Mixed-Effects Model for Powerful Association Tests in Integrative Functional Genomics. *Am. J. Hum. Genet.* *102*, 904–919.
16. Hu, Y., Li, M., Lu, Q., Weng, H., Wang, J., Zekavat, S.M., Yu, Z., Li, B., Gu, J., Muchnik, S., et al.; Alzheimer's Disease Genetics Consortium (2019). A statistical framework for cross-tissue transcriptome-wide association analysis. *Nat. Genet.* *51*, 568–576.
17. Zou, H., and Hastie, T. (2005). Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Series B Stat. Methodol.* *67*, 301–320.
18. Zhou, X., Carbonetto, P., and Stephens, M. (2013). Polygenic modeling with bayesian sparse linear mixed models. *PLoS Genet.* *9*, e1003264.
19. Tibshirani, R. (1996). Regression Shrinkage and Selection via the Lasso. *J. R. Stat. Soc. B* *58*, 267–288.
20. Hoerl, A.E., and Kennard, R.W. (2000). Ridge Regression: Biased Estimation for Nonorthogonal Problems. *Technometrics* *42*, 80–86.
21. Li, Q., and Lin, N. (2010). The Bayesian elastic net. *Bayesian Anal.* *5*, 151–170.
22. Guan, Y.T., and Stephens, M. (2011). Bayesian Variable Selection Regression for Genome-Wide Association Studies and Other Large-Scale Problems. *Ann. Appl. Stat.* *5*, 1780–1815.
23. Yu, J., Pressoir, G., Briggs, W.H., Vroh Bi, I., Yamasaki, M., Doebley, J.F., McMullen, M.D., Gaut, B.S., Nielsen, D.M., Holland, J.B., et al. (2006). A unified mixed-model method for association mapping that accounts for multiple levels of relatedness. *Nat. Genet.* *38*, 203–208.
24. Huan, T., Liu, C., Joehanes, R., Zhang, X., Chen, B.H., Johnson, A.D., Yao, C., Courchesne, P., O'Donnell, C.J., Munson, P.J., and

- Levy, D. (2015). A systematic heritability analysis of the human whole blood transcriptome. *Hum. Genet.* 134, 343–358.
25. Lloyd-Jones, L.R., Holloway, A., McRae, A., Yang, J., Small, K., Zhao, J., Zeng, B., Bakshi, A., Metspalu, A., Dermitzakis, M., et al. (2017). The Genetic Architecture of Gene Expression in Peripheral Blood. *Am. J. Hum. Genet.* 100, 371.
  26. Zeng, P., and Zhou, X. (2017). Non-parametric genetic prediction of complex traits with latent Dirichlet process regression models. *Nat. Commun.* 8, 456.
  27. Blei, D.M., and Jordan, M.I. (2006). Variational inference for Dirichlet process mixtures. *Bayesian Anal.* 1, 121–143.
  28. Blei, D.M., Kucukelbir, A., and McAuliffe, J.D. (2017). Variational Inference: A Review for Statisticians. *J. Am. Stat. Assoc.* 112, 859–877.
  29. Casella, G. (2001). Empirical Bayes Gibbs sampling. *Biostatistics* 2, 485–500.
  30. Bennett, D.A., Schneider, J.A., Arvanitakis, Z., and Wilson, R.S. (2012). Overview and findings from the religious orders study. *Curr. Alzheimer Res.* 9, 628–645.
  31. Bennett, D.A., Schneider, J.A., Buchman, A.S., Barnes, L.L., Boyle, P.A., and Wilson, R.S. (2012). Overview and findings from the rush Memory and Aging Project. *Curr. Alzheimer Res.* 9, 646–663.
  32. Ng, B., White, C.C., Klein, H.U., Sieberts, S.K., McCabe, C., Patrick, E., Xu, J., Yu, L., Gaiteri, C., Bennett, D.A., et al. (2017). An xQTL map integrates the genetic architecture of the human brain's transcriptome and epigenome. *Nat. Neurosci.* 20, 1418–1426.
  33. Bennett, D.A., Buchman, A.S., Boyle, P.A., Barnes, L.L., Wilson, R.S., and Schneider, J.A. (2018). Religious Orders Study and Rush Memory and Aging Project. *J. Alzheimers Dis.* 64 (s1), S161–S189.
  34. Lambert, J.C., Ibrahim-Verbaas, C.A., Harold, D., Naj, A.C., Sims, R., Bellenguez, C., DeStafano, A.L., Bis, J.C., Beecham, G.W., Grenier-Boley, B., et al.; European Alzheimer's Disease Initiative (EADI); Genetic and Environmental Risk in Alzheimer's Disease; Alzheimer's Disease Genetic Consortium; and Cohorts for Heart and Aging Research in Genomic Epidemiology (2013). Meta-analysis of 74,046 individuals identifies 11 new susceptibility loci for Alzheimer's disease. *Nat. Genet.* 45, 1452–1458.
  35. Reitz, C. (2014). Genetic loci associated with Alzheimer's disease. *Future Neurol.* 9, 119–122.
  36. Reitz, C. (2015). Novel susceptibility loci for Alzheimer's disease. *Future Neurol.* 10, 547–558.
  37. Sims, R., van der Lee, S.J., Naj, A.C., Bellenguez, C., Badarinarayan, N., Jakobsdottir, J., Kunkle, B.W., Boland, A., Raybould, R., Bis, J.C., et al.; ARUK Consortium; and GERAD/PERADES, CHARGE, ADGC, EADI (2017). Rare coding variants in PLCG2, ABI3, and TREM2 implicate microglial-mediated innate immunity in Alzheimer's disease. *Nat. Genet.* 49, 1373–1384.
  38. Yuan, X.Z., Sun, S., Tan, C.C., Yu, J.T., and Tan, L. (2017). The Role of ADAM10 in Alzheimer's Disease. *J. Alzheimers Dis.* 58, 303–322.
  39. Müller, P., and Mitra, R. (2013). Bayesian Nonparametric Inference - Why and How. *Bayesian Anal.* 8, 8.
  40. Carbonetto, P., and Stephens, M. (2012). Scalable Variational Inference for Bayesian Variable Selection in Regression, and Its Accuracy in Genetic Association Studies. *Bayesian Anal.* 7, 73–107.
  41. Li, B., and Leal, S.M. (2008). Methods for detecting associations with rare variants for common diseases: application to analysis of sequence data. *Am. J. Hum. Genet.* 83, 311–321.
  42. Abecasis, G.R., Auton, A., Brooks, L.D., DePristo, M.A., Durbin, R.M., Handsaker, R.E., Kang, H.M., Marth, G.T., McVean, G.A.; and 1000 Genomes Project Consortium (2012). An integrated map of genetic variation from 1,092 human genomes. *Nature* 491, 56–65.
  43. O'Reilly, P.F., Hoggart, C.J., Pomyen, Y., Calboli, F.C., Elliott, P., Jarvelin, M.R., and Coin, L.J. (2012). MultiPhen: joint model of multiple phenotypes can increase discovery in GWAS. *PLoS ONE* 7, e34861.
  44. De Jager, P.L., Shulman, J.M., Chibnik, L.B., Keenan, B.T., Raj, T., Wilson, R.S., Yu, L., Leurgans, S.E., Tran, D., Aubin, C., et al.; Alzheimer's Disease Neuroimaging Initiative (2012). A genome-wide scan for common variants affecting the rate of age-related cognitive decline. *Neurobiol. Aging* 33, 1017.e1–1017.e15.
  45. De Jager, P.L., Srivastava, G., Lunnon, K., Burgess, J., Schalkwyk, L.C., Yu, L., Eaton, M.L., Keenan, B.T., Ernst, J., McCabe, C., et al. (2014). Alzheimer's disease: early alterations in brain DNA methylation at ANK1, BIN1, RHBDF2 and other loci. *Nat. Neurosci.* 17, 1156–1163.
  46. Wu, L., Shi, W., Long, J., Guo, X., Michailidou, K., Beesley, J., Bolla, M.K., Shu, X.O., Lu, Y., Cai, Q., et al.; NBCS Collaborators; and kConFab/AOCS Investigators (2018). A transcriptome-wide association study of 229,000 women identifies new candidate susceptibility genes for breast cancer. *Nat. Genet.* 50, 968–978.
  47. Martin, A.R., Gignoux, C.R., Walters, R.K., Wojcik, G.L., Neale, B.M., Gravel, S., Daly, M.J., Bustamante, C.D., and Kenny, E.E. (2017). Human Demographic History Impacts Genetic Risk Prediction across Diverse Populations. *Am. J. Hum. Genet.* 100, 635–649.
  48. Yang, J., Lee, S.H., Goddard, M.E., and Visscher, P.M. (2011). GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* 88, 76–82.
  49. Yang, J., Fritsche, L.G., Zhou, X., Abecasis, G.; and International Age-Related Macular Degeneration Genomics Consortium (2017). A Scalable Bayesian Method for Integrating Functional Information in Genome-wide Association Studies. *Am. J. Hum. Genet.* 101, 404–416.