# TransferTWAS: A transfer learning framework for cross-tissue transcriptome-wide association study

### Authors

Daoyuan Lai, Han Wang, Tian Gu, Siqi Wu, Dajiang J. Liu, Pak Chung Sham, Yan Dora Zhang

### Correspondence doraz@hku.hk

TransferTWAS is a TWAS method that adaptively borrows information from multiple external tissues to boost geneexpression prediction in tissues with small sample sizes. It outperforms competing approaches both in gene-expression imputation accuracy and in the number of detected gene-trait associations.

> Lai et al., 2025, The American Journal of Human Genetics 112, 1–12 August 7, 2025 © 2025 American Society of Human Genetics. Published by Elsevier Inc. All rights are reserved, including those for text and data mining, Al training, and similar technologies. https://doi.org/10.1016/j.ajhg.2025.06.006



### ARTICLE

### TransferTWAS: A transfer learning framework for cross-tissue transcriptome-wide association study

Daoyuan Lai,<sup>1</sup> Han Wang,<sup>2</sup> Tian Gu,<sup>3</sup> Siqi Wu,<sup>1</sup> Dajiang J. Liu,<sup>4</sup> Pak Chung Sham,<sup>5,6</sup> and Yan Dora Zhang<sup>1,\*</sup>

#### Summary

Transcriptome-wide association studies (TWASs) utilize gene-expression data to explore the genetic basis of complex traits. A key challenge in TWASs is developing robust imputation models for tissues with limited sample sizes. This paper introduces transfer learningassisted TWAS (TransferTWAS), a framework that adaptively transfers information from multiple tissues to improve gene-expression prediction in the target tissue. TransferTWAS employs a data-driven strategy that assigns higher weights to genetically similar external tissues. It outperforms other multi-tissue TWAS methods, such as the Unified Test for Molecular Signatures (UTMOST), which neglects tissue similarity, and Joint-Tissue Imputation (JTI), which relies on functional annotations to represent tissue similarity. Simulation studies demonstrate that TransferTWAS achieves the highest imputation accuracy, and analyses using the ROS/MAP and GEUVADIS datasets show a substantial power gain while maintaining control over type-I errors. Furthermore, analysis of the low-density lipoprotein cholesterol GWAS dataset and other complex traits demonstrates that TransferTWAS effectively identifies more associations compared with existing methods.

#### Introduction

Genome-wide association studies (GWASs) have identified numerous single-nucleotide polymorphisms (SNPs) associated with complex diseases; however, they face significant challenges in pinpointing causal genes, particularly for SNPs in non-coding regions.<sup>1,2</sup> To address these limitations, transcriptome-wide association studies (TWASs) have emerged as a powerful approach, focusing on the association between predicted levels of genetically regulated gene expression (GReX) and phenotypes of interest.<sup>3,4</sup> TWAS leverages gene-expression reference panels, such as the Genotype-Tissue Expression (GTEx) project, 5-8 to explore the relationship between genotype and phenotype.<sup>9</sup> Central to TWAS is the hypothesis that SNPs influence complex traits through expression quantitative trait loci (eQTLs). A TWAS involves two key steps: imputing tissue-specific GReX using transcriptomic and genetic data from reference panels and conducting association analyses between **GReX** and the phenotypes.

TWAS methods can be broadly categorized into singletissue and multi-tissue approaches. Singe-tissue methods, which focus on gene expression in biologically relevant tissues, face several limitations. These include an inability to fully utilize the multi-tissue nature of gene-expression reference panels such as GTEx, disregard for cross-tissue transcriptional regulatory similarities, and poor performance in tissues with limited sample sizes.<sup>10–12</sup> In response, multi-tissue methods have been developed to leverage information across multiple tissues, aiming to improve performance in tissues with small effective sample sizes by incorporating data from larger, external tissues.<sup>10–13</sup>

The multi-tissue method, MultiXcan, regresses the complex phenotype of interest onto the principal components of the **GReX** from all available tissues.<sup>14</sup> However, MultiXcan does not enhance the quality of GReX in individual tissues, and the interpretability of its principal components is limited. Another approach, the Unified Test for Molecular Signatures (UTMOST), jointly models genotype and cross-tissue gene-expression data but fails to account for cross-tissue similarity.<sup>10,11</sup> To address this limitation, Zhou et al.<sup>11</sup> proposed the joint-tissue imputation (JTI) method, which leverages shared genetic regulation of gene expression across tissues and incorporates external annotations, such as tissue-level expression correlations and gene-level DNase I hypersensitive site similarity. While JTI improves prediction performance, its effectiveness depends on the quality and availability of the functional annotations.

To overcome these challenges, we propose TransferTWAS, a multi-tissue TWAS method that employs transfer learning to enhance *GReX* imputation. Unlike traditional methods, TransferTWAS does not depend on functional annotations; instead, it automatically prioritizes tissues with expression patterns similar to the target tissue, minimizing the impact of dissimilar tissues. Our extensive simulations show that TransferTWAS significantly enhances gene-expression

\*Correspondence: doraz@hku.hk

https://doi.org/10.1016/j.ajhg.2025.06.006.

<sup>&</sup>lt;sup>1</sup>Department of Statistics and Actuarial Science, School of Computing and Data Science, The University of Hong Kong, Hong Kong SAR, China; <sup>2</sup>College of Science, China Agricultural University, Beijing, China; <sup>3</sup>Department of Biostatistics, Columbia University Mailman School of Public Health, New York, NY, USA; <sup>4</sup>Institute for Personalized Medicine, College of Medicine, Pennsylvania State University, University Park, PA, USA; <sup>5</sup>Department of Psychiatry, Li Ka Shing Faculty of Medicine, The University of Hong Kong, Hong Kong SAR, China; <sup>6</sup>Centre for PanorOmic Sciences, Li Ka Shing Faculty of Medicine, The University of Hong Kong SAR, China

<sup>© 2025</sup> American Society of Human Genetics. Published by Elsevier Inc. All rights are reserved, including those for text and data mining, AI training, and similar technologies.



## Figure 1. The schematic workflow of TransferTWAS

Suppose the target tissue is tissue k. The first step is to take standard ridge regression to train the expression predictive models in other available tissues and obtain the imputation weights  $\hat{\boldsymbol{\beta}}^{(j)}$  for  $j \neq k$ . Next, we take the algorithm 1 to aggregate  $\hat{\boldsymbol{\beta}}^{(j)}$ 's information into  $\hat{\boldsymbol{\beta}}^{(-k)}$ . The third step uses  $\hat{\boldsymbol{\beta}}^{(-k)}$ , tissue k's gene expression  $\boldsymbol{E}^{(k)}$ , and genotype data  $\boldsymbol{G}^{(k)}$  as inputs to calculate the final estimator.

imputation accuracy compared to existing multi-tissue methods, resulting in greater TWAS power while maintaining control over type-I error rates. When applied to a quantile-transformed low-density lipoprotein cholesterol (LDL-C; MIM: 605028) GWAS dataset, TransferTWAS effectively identifies well-known causal genes, including the *SORT1* (MIM: 602458)-*PSRC1* (MIM: 613126)-*CELSR2* (MIM: 604265) cluster, *KPNB1* (MIM: 602738), and *LIPC* (MIM: 151670). Additionally, when tested on 30 other complex traits, TransferTWAS uncovers the highest number of significant associations, highlighting its wide-ranging applicability and effectiveness in genetic research.

#### Material and methods

#### **TWAS framework**

The first step of TWAS involves estimating *cis*-eQTL effect sizes using a gene-expression reference panel, which includes both gene-expression and genotype data. For tissue k ( $k \in \{1,...,K\}$ ), the relationship between gene expression and genotype is modeled as a multiple linear regression:

$$\boldsymbol{E}^{(k)} = \boldsymbol{G}^{(k)}\boldsymbol{\beta}^{(k)} + \boldsymbol{\epsilon}^{(k)}, \qquad (\text{Equation 1})$$

where  $\mathbf{E}^{(k)} \in \mathbb{R}^{n_k}$  is a vector of gene expression for  $n_k$  individuals in tissue  $k, \mathbf{G}^{(k)} \in \mathbb{R}^{n_k \times M}$  is the genotype matrix of the M cis-SNPs (within 1 MB of the gene's flanking regions),  $\boldsymbol{\beta}^{(k)} = \left(\boldsymbol{\beta}_1^{(k)}, \dots, \boldsymbol{\beta}_M^{(k)}\right)^{\top}$  is the M vector of the eQTL effect sizes, and  $\boldsymbol{\epsilon}^{(k)} \in \mathbb{R}^{n_k}$  denotes the residual error term. Gene expression  $\mathbf{E}^{(k)}$  is adjusted for non-genetic covariates and centered such that  $\mathbb{E}(\mathbf{E}^{(k)}) = 0$ . The genotype matrix  $\mathbf{G}^{(k)}$  is centered but not standardized. After estimating  $\hat{\boldsymbol{\beta}}^{(k)}$  for  $k \in \{1, \dots, K\}$ , the **GReX** for a GWAS dataset with genotype matrix  $\overline{\mathbf{G}}$  is imputed as

$$G\widehat{Re}X = \overline{G}\,\widehat{\beta}^{(k)}$$

#### **Overview of TransferTWAS**

We provide a visualization of the workflow of TransferTWAS in Figure 1. For simplicity, we first assume that we are working on

a gene that only has expression in two tissues. To calculate the cis-eQTL effect size in tissue k, TransferTWAS optimizes the following loss function to enhance *GReX* imputation:

$$\widehat{\boldsymbol{\beta}}^{(k)} = \operatorname{argmin}_{\boldsymbol{\beta}^{(k)}} \frac{1}{n_k} \| \boldsymbol{E}^{(k)} - \boldsymbol{G}^{(k)} \boldsymbol{\beta}^{(k)} \|_2^2 + \underbrace{\lambda \| \boldsymbol{\beta}^{(k)} \|_2^2}_{\text{Ridge penalty}}.$$
 (Equation 2)  
$$-\underbrace{2\eta \left( \widehat{\boldsymbol{\beta}}^{(-k)} \right)^{\mathsf{T}} \boldsymbol{\beta}^{(k)}}_{\text{Angle-based penalty}}, k = 1, ..., K$$

Here,  $\lambda, \eta \in \mathbb{R}$  are tuning parameters controlling the ridge and the angle-based penalties, respectively. Since we only have two tissues (one target and one external),  $\hat{\boldsymbol{\beta}}^{(-k)} \in \mathbb{R}^M$  is the estimated eQTL effect size from that external tissue. The angle-based penalty, motivated by Gu et al.,<sup>15</sup> encourages alignment between  $\boldsymbol{\beta}^{(k)}$  and informative external effect directions.

Now we turn to the scenario that multiple (larger than one) external tissues are available. In this case,  $\hat{\beta}^{(-k)}$  is adaptively aggregated across external tissues using algorithm 1, which assigns higher weights to tissues with stronger predictive utility. This ensures that only relevant tissues contribute to the model, reducing noise from uninformative sources. Readers may refer to Note S1 for a detailed explanation of the details and the logic of algorithm 1.

#### Solving the loss function

Equation 2 has a closed-form solution given by

$$\widehat{\boldsymbol{\beta}}^{(k)} = \left(\boldsymbol{G}^{(k)\top}\boldsymbol{G}^{(k)} + n_k \lambda \boldsymbol{I}\right)^{-1} \left(\boldsymbol{G}^{(k)\top}\boldsymbol{E}^{(k)} + n_k \eta \widehat{\boldsymbol{\beta}}^{(-k)}\right). \quad \text{(Equation 3)}$$

However, calculating the inverse of  $\mathbf{G}^{(k)\top}\mathbf{G}^{(k)}$  becomes computationally expensive for large *M*. To address this, we derive an alternative formulation of Equation 3,

$$\widehat{\boldsymbol{\beta}}^{(k)} = \boldsymbol{V}_t \boldsymbol{\Sigma}_1 \boldsymbol{U}_t^{\top} \boldsymbol{E}^{(k)} + n_k \eta \left( \boldsymbol{V}_t \boldsymbol{\Sigma}_2 \boldsymbol{V}_t^{\top} + \frac{(\boldsymbol{I} - \boldsymbol{V}_t \boldsymbol{V}_t^{\top})}{n_k \lambda} \right) \widehat{\boldsymbol{\beta}}^{(-k)}.$$
(Equation 4)

In this formulation, the matrices  $U_t$  and  $V_t$  consist of the first t columns of the matrices U and V, which are obtained through the singular value decomposition (SVD) of  $G^{(k)}$ . We define

$$\mathbf{\Sigma}_1$$
: = diag $\left(\frac{d_1}{d_1^2 + n_k \lambda}, \dots, \frac{d_t}{d_t^2 + n_k \lambda}\right)$ ,

$$\mathbf{\Sigma}_2$$
: = diag $\left(\frac{1}{d_1^2 + n_k \lambda}, \dots, \frac{1}{d_t^2 + n_k \lambda}\right)$ ,

where  $(d_1, ..., d_t)$  are the first t singular values of  $\mathbf{G}^{(k)}$ , and  $\mathbf{I}$  is an identity matrix of size  $M \times M$ . Throughout this paper, we set  $t = \min(n_k, M)$ .

Equation 4 offers two key advantages: it avoids the computationally expensive matrix inversion, making it scalable for large M, and it significantly reduces computational time by using the economic SVD instead of the full SVD. A detailed derivation of Equation 4 can be found in Note S2.

#### Choice of tuning parameters

The selection of the optimal tuning parameters  $(\lambda, \eta)$  is a critical step in the implementation of TransferTWAS. The tuning process starts with  $\lambda$ . We first perform standard ridge regression using the R function cv.glmnet() with alpha = 0 and nfold = 5. This function generates a sequence of  $\lambda$  values, from which we select five equally spaced candidates within the range  $(\lambda_{\min}, \lambda_{\max})$ .

Next, candidate values for  $\eta$  are generated based on the selected  $\lambda$  candidates. As suggested by Gu et al.,<sup>15</sup> the theoretical optimal  $\eta_{\text{opt}}$  is related to  $\lambda_{\text{opt}}$  through the following equation:

$$\eta_{\text{opt}} = \lambda_{\text{opt}} \rho \frac{\alpha_k}{\alpha_{-k}}.$$
 (Equation 5)

Here,  $\rho$  is the Pearson correlation between  $\widehat{\boldsymbol{\beta}}^{(k)}$  and  $\widehat{\boldsymbol{\beta}}^{(-k)}$ . The two terms  $\alpha_k := \mathbb{E}(\|\widehat{\boldsymbol{\beta}}^{(k)}\|_2^2)$  and  $\alpha_{-k} := \mathbb{E}(\|\widehat{\boldsymbol{\beta}}^{(-k)}\|_2^2)$  measure the signal strength of the target and external datasets, respectively. Using Equation 5, five  $\eta$  candidates are generated from the  $\lambda$  candidates. Finally, the best combination of  $\lambda$  and  $\eta$  is determined through cross-validation (CV), ensuring optimal performance of the model.

#### GTEx, GEUVADIS, and ROS/MAP data preprocessing

The preprocessing of the GTEx data follows Zhou et al.,<sup>11</sup> Wang et al.,<sup>16</sup> and Hu et al..<sup>10</sup> Gene-expression imputation models were trained using genotype and normalized gene-expression data from 48 GTEx tissues. SNPs with ambiguous alleles, minor allele frequency (MAF) less than 0.05, or Hardy-Weinberg equilibrium (HWE) *p* values less than 0.05 were excluded. The gene-expression data were adjusted for potential confounding effects, including sex, sequencing platform, the top three principal components of genotype data, and the top probabilistic estimation of expression residuals (PEER) factors. The number of PEER factors included in the adjustment was determined by tissue sample size: 15 ( < 150 samples), 30 (150–250 samples), and 35 (> 250 samples). Covariates were sourced from the GTEx portal, and biallelic SNPs within a 1-MB region of the target gene were selected as features.

TransferTWAS predictive performance was evaluated using two distinct gene-expression reference panels: Religious Orders Study and Rush Memory Aging Project (ROS/MAP) with brain (prefrontal cortex) tissue,<sup>17</sup> and Genetic European variation in disease (GEUVADIS) with lymphoblastoid cell lines.<sup>18</sup> These datasets enable a comprehensive assessment of TransferTWAS. Their preprocessing follows protocols from Wang et al.<sup>19</sup> and Keys et al.,<sup>20</sup> with detailed methods provided in Note S3.

# Simulation study I: Gene-expression imputation accuracy

Simulation study I evaluates the accuracy of gene-expression imputation for UTMOST, JTI, and TransferTWAS. The design followed Khunsriraksakul et al.,<sup>21</sup> Feng et al.,<sup>22</sup> and Nagpal et al.<sup>23</sup> We measured imputation accuracy using the squared Pearson correlation (i.e.,  $r^2$ ) between the observed and predicted gene-expression levels. We focused on the gene *CPTP* (MIM: 615467), which has expression across all GTEx tissues, and examined three types of tissues: causal tissue, where genetic variants directly affect expression levels; genetically correlated tissues, which influence expression in the target tissue to a lesser extent; and genetically uncorrelated tissues with no genetic relationship to the causal tissue. In the simulations, the brain (prefrontal cortex) was designated as the causal tissue.

#### Simulation of gene expression in causal tissue

The gene expression in the causal tissue,  $\boldsymbol{E}$ , was simulated using the formula

$$\boldsymbol{E} = \boldsymbol{G}\boldsymbol{\beta} + \boldsymbol{\epsilon}. \qquad (\text{Equation 6})$$

Here,  $\mathbf{G} \in \mathbb{R}^{n \times M}$  is the normalized genotype matrix for n = 130GTEx individuals with *CPTP* expression data in the target tissue, where the matrix has a mean of 0 and a variance of 1. The vector  $\boldsymbol{\beta} = (\beta_1, ..., \beta_M)^\top$  contains the eQTL effect sizes with M = 2,212. We randomly selected  $p_{\text{causal}} = 5\%$  SNPs as causal SNPs. The set of non-zero coefficients (i.e., the causal SNPs) in  $\boldsymbol{\beta}$  is denoted by  $S = \{j : w_j \neq 0\}$ . For the SNPs in *S*, we generated effect sizes  $\beta_j$ (where  $j \in S$ ) from a standard normal distribution N(0, 1), while the effect sizes for the remaining non-causal SNPs were set to 0. We then rescaled the effect sizes  $\boldsymbol{\beta}$  to ensure that the gene-expression heritability (i.e., the proportion of gene-expression variance explained by SNPs) is  $h_e^2$ . The residual error  $\epsilon$  follows a normal distribution  $N(\mathbf{0}, (1 - h_e^2)\mathbf{I})$ .

#### Simulation of gene expression in correlated tissues

The gene expressions in  $N_{\rm corr}$  randomly selected correlated tissues were simulated next. We assumed a uniform genetic correlation  $\rho$ between the causal tissue and each correlated tissue. The effect sizes for the *k*-th correlated tissue were denoted as  $\boldsymbol{\beta}^{(k)}$  and simulated as follows:

$$eta_j^{(k)} \sim N\Big(
hoeta_j, ig(1-
ho^2ig) imes h_e^2\Big), k=1,...,N_{ ext{corr}}, j \in S_k.$$

Here,  $\beta_j^{(k)}$  represents the *j*-th coordinate of  $\beta^{(k)}$ , and  $S_k$  is the active set of  $\beta^{(k)}$ . We assumed that  $S_k$  is a random subset of *S*, with  $|S_k| = q_k |S|$ . The percentage  $q_k$  represents proportion of the shared causal SNPs between the causal tissue and tissue *k* in GTEx, with specific values provided in Khunsriraksakul et al..<sup>21</sup> The residual gene-expression values in correlated tissues were simulated as

$$\epsilon^{(k)} \sim N\left(\mathbf{0}, \operatorname{diag}\left(\sqrt{1-h_e^2}\right) \times \mathbf{\Sigma} \times \operatorname{diag}\left(\sqrt{1-h_e^2}\right)\right)$$

where  $\Sigma$  represents the residual correlation among the geneexpression levels across tissues. Following Khunsriraksakul et al.,<sup>21</sup> we set  $\Sigma = I$ . The gene-expression levels in the *k*-th correlated tissues,  $E^{(k)}$ , were then simulated as

$$\mathbf{E}^{(k)} = \mathbf{G}^{(k)} \mathbf{\beta}^{(k)} + \boldsymbol{\epsilon}^{(k)}, \qquad (\text{Equation 7})$$

where  $\mathbf{G}^{(k)} \in \mathbb{R}^{n_k \times M}$  is the genotype matrix of the  $n_k$  GTEx individuals that have gene expression in the *k*-th tissue.

#### Simulation of gene expression in uncorrelated tissues

Gene-expression data  $(\mathbf{E}^{(l)}, \mathbf{G}^{(l)})$  in the *l*-th uncorrelated tissues were simulated using a similar model as in Equation 6, but with  $\rho = 0$ .

Simulations were repeated 1,000 times for each of the six pairs of  $(\rho, N_{\rm corr}, p_{\rm causal})$ , i.e., (0.3, 0, 5%), (0.3, 24, 5%), (0.3, 47, 5%), (0.8, 0, 5%), (0.8, 24, 5%), and (0.8, 47, 5%). Results were presented as averages across 1,000 replicates for each  $(\rho, N_{\rm corr}, p_{\rm causal})$  combination. The proportion of causal SNPs,  $p_{\rm causal}$ , also varied among values in vector  $p_{\rm causal} = (0.1\%, 2\%, 5\%, 10\%)$ .

#### Simulation study II: TWAS power analysis

Simulation study II compares the performance of UTMOST, JTI, and TransferTWAS in terms of TWAS power, following the design outlined by Zhou et al..<sup>11</sup> We assumed *P* causal genes for the brain (prefrontal cortex) based on the ROS/MAP panel, and the true expression  $\mathbf{E}_{g,\text{true}}$  for each causal gene g = 1, 2, ..., P was simulated from a standard normal distribution  $N(\mathbf{0}, \mathbf{I})$ . The phenotype  $\mathbf{Y}$  was genetically determined by these *P* causal genes and generated using the equation

$$\boldsymbol{Y} = \sum_{g=1}^{p} \alpha_{g} \boldsymbol{E}_{g,\text{true}} + \boldsymbol{\epsilon}.$$
 (Equation 8)

In this equation, the coefficients  $\alpha_g$  were drawn from  $N(0, h_p^2/P)$ , while the residual  $\epsilon$  followed  $N(0, (1 - h_p^2)\mathbf{I})$ . The term  $h_p^2$  represents phenotypic heritability, which is the proportion of phenotypic variance explained by gene-expression levels. Equation 8 indicates the overall phenotypic variance explained by gene expression is  $h_p^2$ , and each causal gene contributes, on average,  $\mathbb{E}(\alpha_i) = h_p^2/P$  to this variance.

We constructed predicted gene expressions, denoted as  $\widehat{GReX}_{UTMOST}$ ,  $\widehat{GReX}_{JTI}$ , and  $\widehat{GReX}_{TransferTWAS}$ . To achieve this, we utilized eQTL effect sizes estimated from UTMOST, JTI, and TransferTWAS, which were derived from GTEx data. These effect sizes were applied to predict gene expression in the ROS/MAP dataset. Next, we calculated the empirical correlations for the three methods:  $r_{UTMOST}$ ,  $r_{JTI}$ , and  $r_{TransferTWAS}$ . These correlations represent the Pearson correlation between observed and predicted gene expressions based on each method's performance in ROS/MAP. An empirical correlation matrix  $\Phi$  was constructed as follows:

$$\mathbf{\Phi} = \begin{pmatrix} 1 & r_{\text{UTMOST}} & r_{\text{JTI}} & r_{\text{TransferTWAS}} \\ r_{\text{UTMOST}} & 1 & 0 & 0 \\ r_{\text{JTI}} & 0 & 1 & 0 \\ r_{\text{TransferTWAS}} & 0 & 0 & 1 \end{pmatrix}.$$

The predicted gene expressions  $G\widehat{ReX}_{UTMOST}, G\widehat{ReX}_{JTI}, G\widehat{ReX}_{TransferTWAS}$  were determined using the following formula:

$$\left( \boldsymbol{E}_{g,\text{true}}, \boldsymbol{G} \widehat{\boldsymbol{R}} \boldsymbol{e} \boldsymbol{X}_{\text{UTMOST}}, \boldsymbol{G} \widehat{\boldsymbol{R}} \boldsymbol{e} \boldsymbol{X}_{\text{JTI}}, \boldsymbol{G} \widehat{\boldsymbol{R}} \boldsymbol{e} \boldsymbol{X}_{\text{TransferTWAS}} \right) = \\ \left( \boldsymbol{E}_{g,\text{true}}, \boldsymbol{Z}_1, \boldsymbol{Z}_2, \boldsymbol{Z}_3 \right) \times \text{cholesky}(\boldsymbol{\Phi}) \times \text{SD} \left( \boldsymbol{E}_{g,\text{true}} \right) \\ + \text{mean} \left( \boldsymbol{E}_{g,\text{true}} \right),$$

where  $Z_1$ ,  $Z_2$ , and  $Z_3$  were independently simulated from a standard normal distribution  $N(\mathbf{0}, \mathbf{I})$ . The goal is to achieve specified correlations  $r_{\text{UTMOST}}$ ,  $r_{\text{JTI}}$ , and  $r_{\text{TransferTWAS}}$  between  $\mathbf{E}_{g,\text{true}}$  and the respective predicted gene expressions, mimicking the behavior of the prediction models under study. The predicted gene expressions were generated independently while ensuring they correlate with the true expression  $\mathbf{E}_{g,\text{true}}$ .

Subsequently, 1,000 simulations were performed to test the association between the predicted gene expression  $G\widehat{ReX}_{UTMOST}$ ,  $G\widehat{ReX}_{TTI}, G\widehat{ReX}_{TransferTWAS}$ , and the phenotype Y. The TWAS power was estimated as the proportion of simulations achieving sta-

tistical significance ( $p_{\text{Bonferroni}} < 0.05$ ). To explore different scenarios, we varied the expression heritability  $h_p^2$  across the values (0.05, 0.1, 0.2, 0.3) and the number of causal genes *P* among (40, 50, 60, 70). Additionally, we varied the causal tissues by using the GEUVADIS dataset as the reference panel and selected Epstein-Barr virus (EBV) transformed lymphocytes as the causal tissue to construct  $\Phi$  and simulate  $GReX_{\text{UTMOST}}$ ,  $GReX_{\text{TransferTWAS}}$ , and *Y*.

#### Simulation study III: Type-I error analysis

To assess the type-I error rates of the three methods, we assume no association between the phenotype  $\mathbf{Y}$  and the true gene expression  $\mathbf{E}_{g,\text{true}}$ . Therefore,  $\mathbf{Y}$  was directly simulated from a standard normal distribution  $N(\mathbf{0}, \mathbf{I})$ . For each gene in the designated tissue, we simulated  $\mathbf{Y}$  and regressed it on the predicted gene expression  $\mathbf{GReX}$  derived directly from GTEx. This process was replicated 1,000 times for each gene, and a significance threshold of 0.05 was used for the type-I error analysis. We considered two different GTEx tissues as the causal tissue: the brain (prefrontal cortex) and EBV-transformed lymphocytes.

We also evaluated TWAS power and type-I error using alternative simulation designs from Nagpal et al.,<sup>23</sup> Feng et al.,<sup>22</sup> and Khunsriraksakul et al.<sup>21</sup> These are designated as simulation studies IV (power) and V (type-I error), with detailed information provided in Note S4.

#### Model assessment in GTEx

The gene-expression imputation performances of UTMOST, JTI, and TransferTWAS were compared based on their  $r^2$ . The dataset was randomly divided into five equal-sized groups, with 3/5 as training set, 1/5 as validation set, and 1/5 as test set. A 5-fold CV was conducted on the training set to select the best tuning parameter that minimizes the prediction mean squared error on the validation set. The models were then trained on the training and validation sets using the selected tuning parameter, and prediction performance was evaluated on the test set through the Pearson correlation. The final correlation was calculated based on the average of the five Pearson correlation estimates, with  $r^2$ set to 0 if the training model is null. This assessment procedure follows Khunsriraksakul et al.<sup>21</sup>

#### TWAS with summary-level GWAS

When working with summary-level data in a GWAS, S-PrediXcan<sup>24</sup> is applied to calculate the TWAS statistic using the following formula:

$$Z_g = \sum_{m=1}^{M} \widehat{\beta}_m \frac{\widehat{\sigma}_m}{\widehat{\sigma}_g} \frac{\widehat{\gamma}_m}{\operatorname{se}(\widehat{\gamma}_m)}, \qquad (\text{Equation 9})$$

where  $\hat{\beta}_m$  is the prediction weight of gene *g*'s SNP *m* obtained in the first step of TWAS,  $\hat{\sigma}_m$  is SNP *m*'s variance,  $\hat{\sigma}_g$  is an estimate of gene *g*'s predicted expression's variance, and  $\hat{\gamma}_m$  and se( $\hat{\gamma}_m$ ) are the GWAS regression coefficient for SNP *m* and corresponding standard error.

#### Results

#### Simulation studies

Simulation showed TransferTWAS achieved the highest TWAS power, controlled type-I error, and improved gene-expression imputation accuracy.



Figure 2. Comparison of gene-expression imputation accuracy ( $r^2$ ) among TransferTWAS, UTMOST, and JTI in simulation study I

The average Pearson correlations ( $r^2$ ) between observed and predicted gene expression in the test dataset by TransferTWAS, UTMOST, and JTI, with various combinations of genetic correlation between causal and correlated tissues  $\rho = (0.3, 0.8)$  and number of correlated tissues  $N_{\text{corr}} = (0.24, 47)$ . We assumed the proportion of causal SNPs is 5% in this simulation.

Simulation study I examined TransferTWAS's geneexpression imputation performance under various scenarios. Figure 2 shows that across different  $(\rho, N_{\text{corr}})$  combinations, TransferTWAS consistently outperformed UTMOST and JTI in terms of  $r^2$  if 5% of the SNPs were causal. Notably, we considered some extreme cases-for example, there is only one causal tissue and no correlated tissue  $(N_{\text{corr}} = 0)$ . In this situation, TransferTWAS still achieved higher imputation accuracy compared to the other methods. This highlights the robustness of TransferTWAS, as it effectively avoids negative transfer, where a transfer learning method performs worse than a target-only method.<sup>25</sup> We also considered the scenario where  $\rho = 0.8$ , indicating a non-zero correlation between the correlated tissues. Under this scenario, TransferTWAS still achieved improved performance.

TransferTWAS demonstrated robust performance across varying causal proportions. When  $p_{causal} > 2\%$ , it consistently achieved higher imputation  $r^2$  on the test dataset compared with UTMOST and JTI across all levels of expression heritability  $h_e^2$  and number of causal genes P (Figures 2 and S4–S6). At  $p_{causal} = 2\%$ , TransferTWAS outperformed both methods for  $h_e^2 = (0.05, 0.1, 0.15, 0.2)$ , while maintaining an advantage over JTI at  $h_e^2 = 0.25$  despite UTMOST's slightly better performance in this specific scenario (Figure S4). However, under a sparse model with  $p_{causal} = 0.1\%$ , UTMOST and JTI yield higher  $r^2$  in the test dataset compared with TransferTWAS (Figure S5). These patterns suggest that TransferTWAS achieves optimal performance when  $p_{causal} \ge 2\%$  and

 $h_e^2 \leq 0.2$ , with performance comparable to that of UTMOST at  $h_e^2 = 0.25$ . Since UTMOST does not incorporate tissue-similarity information, these results indicate that leveraging external tissue data provides greatest benefit when  $p_{\text{causal}} \geq 2\%$  and  $h_e^2 \leq 0.2$ . The observed performance differences reflect each method's underlying assumptions. TransferTWAS assumes a non-sparse architecture, while UTMOST and JTI assume a sparse one. We will expand on these implications in the discussion.

Simulation studies II and IV indicated that TransferTWAS exhibits significantly higher statistical power than other methods across sample sizes from 5,000 to 500,000 (Figures 3, S1, and S2) when analyzing the brain (prefrontal cortex) using ROS/MAP data. This advantage remained consistent across varying levels of  $h_e^2$  and P. A similar trend was observed with EBV-transformed lymphocytes, as shown in Figure S2. Additionally, TransferTWAS maintained robust performance under varying conditions, including the number of tissues correlated with the causal tissue ( $N_{corr}$ ), expression heritability ( $h_e^2$ ), and correlation strength ( $\rho$ ), as detailed in Table S1.

In addition to its enhanced power, TransferTWAS controlled type-I error rates in Simulation studies III and V. Evaluations in the brain (prefrontal cortex) and EBV-transformed lymphocytes revealed that TransferTWAS maintains well-controlled type-I error rates, as shown in Figures 4 and S3. While UTMOST and JTI exhibited comparable type-I error rates in the brain (prefrontal cortex), UTMOST showed inflation in EBV-transformed



Figure 3. Power comparison of UTMOST, JTI, and TransferTWAS based on simulation study II using ROS/MAP data

We simulated 40 causal genes (P = 40) explaining  $h_p^2 = 2\%$  of the total phenotypic variance. True gene-expression levels and their effects on the trait were simulated, with each gene contributing  $h_p^2/P$  variance. Predicted expression levels were generated using the actual prediction performance ( $r^2$ ) from ROS/MAP for each method. Power was calculated as the proportion of simulations with Bonferroni-corrected significance  $p_{\text{Bonferroni}} < 0.05$ . More scenarios were evaluated in Figure S1.

lymphocytes. Further analysis (simulation study V) confirms that TransferTWAS's type-I error remains well controlled across varying levels of  $h_e^2$ , $\rho$ , and  $N_{corr}$ , as summarized in Table S2.

#### Real application to GTEx

We first compared the transcriptome-wide 5-fold CV  $r^2$  of TransferTWAS, UTMOST, and JTI in GTEx. Figures 5A and S7 illustrate that TransferTWAS and UTMOST improve  $r^2$ over JTI on the test dataset, with TransferTWAS showing increased imputation accuracy as GTEx tissue sample sizes decrease. Notably, TransferTWAS achieved a mean  $\Delta r_{\text{TransferTWAS}}^2$  ( $r^2$  difference between TransferTWAS and JTI) of 0.017, surpassing UTMOST's mean  $\Delta r_{\text{UTMOST}}^2$  ( $r^2$  difference between UTMOST and JTI) of 0.006 (Table S3). While TransferTWAS's  $r^2$  was lower than JTI's in four tissues with large sample size, it showed improvement over JTI in the remaining 44 GTEx tissues. The enhancement over UTMOST can be attributed to the inclusion of tissue similarity information, which UTMOST does not consider. Additionally, TransferTWAS's improvement over JTI showed that its data-driven approach to aggregating external tissue information appears to be more effective in most GTEx tissues.

Second, we compared the number of imputable genes (iGenes, defined as  $r^2 > 0.01$ , as suggested by multiple

studies<sup>3,23,26</sup>). TransferTWAS showed an average of 10,744 iGenes, exceeding UTMOST's 7,927 and JTI's 6,668. TransferTWAS consistently outperformed JTI across all GTEx tissues, while UTMOST failed to do so in larger-sample-size tissues (Figure 5B and Table S3). Although TransferTWAS may not exceed JTI in  $\Delta r^2$  for some larger tissues, it effectively captured more iGenes, indicating strong imputation capability.

Figure 5C and Table S3 analyze the proportion of iGenes captured. TransferTWAS captured an average of 81.68% of JTI's iGenes compared to UTMOST's 75.46%. Thus, TransferTWAS not only identified a substantial number of iGenes of JTI but also those previously unaccounted for.

Focusing on tissues with sample sizes smaller than 300, TransferTWAS's superiority became more evident, achieving a mean  $\Delta r_{\text{TransferTWAS}}^2 = 0.028$ , versus UTMOST's 0.013 (Table S3). It identified an average of 11,390 iGenes, exceeding UTMOST's 8,419 and JTI's 6,361, and captured 89.42% of JTI's iGenes compared to UTMOST's 87.50%. This confirms the effectiveness of TransferTWAS's transfer learning approach for tissues with limited sample sizes.

In challenging contexts involving large sample sizes, TransferTWAS consistently outperformed or matched UTMOST and JTI. For muscle (skeletal) tissue (n =706), TransferTWAS has  $\Delta r_{\text{TransferTWAS}}^2 = -0.001$  and identifies 8,380 iGenes, outperforming UTMOST  $(\Delta r_{\rm UTMOST}^2 = -0.005 \text{ with } 6,077 \text{ iGenes}) \text{ and JTI } (6,454)$ iGenes). In the testis tissue, TransferTWAS demonstrated an even more impressive  $\Delta r_{\text{TransferTWAS}}^2 = 0.004$ , surpassing UTMOST's  $\Delta r_{\rm UTMOST}^2 = -$  0.004. Additionally, TransferTWAS identified a substantially larger number of significant iGenes, totaling 11,908, compared to UTMOST's 8,932 and JTI's 8,276. Overall, while TransferTWAS may show a slight disadvantage in imputation  $r^2$  for tissues with relatively large sample size, it consistently allows for a greater number of genes to be classified as imputable.

Additionally, we conducted a replication study using weights trained on GTEx samples to predict expression levels in 373 European individuals from the GEUVADIS dataset. TransferTWAS achieved higher prediction  $r^2$  and identified more iGenes compared to UTMOST and JTI (Table S4).

Overall, TransferTWAS enhanced imputation accuracy in GTEx tissues, which was consistent with our simulation result.

# Real application to quantile-transformed LDL-C GWAS dataset

We applied the gene-expression imputation models from GTEx data to identify potential risk genes of quantile-transformed LDL-C (N = 343,621) using the UK Biobank GWAS dataset. The SNP-SNP covariance matrices for Equation 9 were estimated using the GTEx v.8 samples, and identified associations are validated against existing literature.



Figure 4. Comparison of type-I error rates for UTMOST, JTI, and TransferTWAS in simulation study III using ROS/MAP prefrontal cortex data

Quantile-quantile plot of TWAS *p* values from TransferTWAS, UTMOST, and JTI are generated to visualize the type-I error rates of these models in brain (prefrontal cortex) compared to the expected values, with the blue dashed lines representing the 95% confidence intervals of the expected  $-\log(p)$  values.

As shown in Figure 6A, TransferTWAS identified the largest number of significant associations (1,385) in liver tissue, outperforming JTI (375) and UTMOST (483). The significance threshold was set at a false discovery rate (FDR)-corrected *p* value of less than 0.05 ( $p_{FDR} < 0.05$ ). Among associations identified by JTI, 54.67% (205) were also nominally significant (*p* < 0.05) under UTMOST, while TransferTWAS increased this proportion to 67.73% (254).

We examined TransferTWAS's ability to replicate wellknown LDL-C-associated genes. Among the 59 LDL-Crelated genes reported by Zhou et al.,<sup>11</sup> TransferTWAS identified 22, while UTMOST and JTI captured 11 (Figure 6B and Table S5). UTMOST, JTI, and TransferTWAS consistently captured many well-known LDL-C-related genes. For example, all three methods show similar strong association signals for the potential LDL-C-related genes, including *PCSK9* (MIM: 607786), *SORT1-PSRC1-CELSR2* cluster, *KPNB1*, and *LIPC* (Table S6). For other LDL-C genes, TransferTWAS showed a boosted performance. TransferTWAS uniquely identified *ANGPTL3* (MIM: 603874) as imputable ( $r^2 = 1.85\%$ ), leading to a significant TWAS association (p = 0) and supporting findings that inhibiting *ANGPTL3* lowers LDL-C levels.<sup>27</sup>

We identified 898 additional associations through the TransferTWAS method (Table S7). Based on the suggestion of Zhou et al.,<sup>11</sup> we defined the additional association usings the following criteria: TransferTWAS  $p_{FDR} < 0.05$ ; UTMOST p > 0.05 or not imputable; and JTI p > 0.05 or not imputable. Among these, several associations merit discussion. An improved signal was detected for *APOA1* (MIM: 107680) (TransferTWAS:  $r^2 = 2.95\%$ ), whereas the other two methods reported this gene as not imputable. Such improvement on gene-expression

imputation may contribute to the significant associations from TransferTWAS ( $p = 5.79 \times 10^{-7}$ ). Similarly, TransferTWAS showed an improved imputation quality in the APOB (MIM: 107730) gene (TransferTWAS:  $r^2 =$ 17%; UTMOST: not imputable; JTI: not imputable), leading to a significant association (TransferTWAS: p = 0; UTMOST: not imputable; JTI: not imputable). This finding replicated the observations of Peloso et al.,<sup>28</sup> who reports that mutations in APOB may be associated with lower LDL-C. An enhanced imputation quality for ABCA6 (MIM: 612504) was suggested by TransferTWAS  $(r^2 = 3.87\%)$ , and the corresponding TWAS p value is  $5.10 \times 10^{-4}$ . This finding is in line with those of Francis et al.,<sup>29</sup> who associated a variant of this gene with LDL-C. UTMOST and JTI failed to impute this gene.

#### **Enriched pathways in LDL-C**

To assess the biological relevance between LDL-C and the significant genes identified by TransferTWAS, we performed functional enrichment analysis. As shown in Figure S8, the most significant pathways are directly tied to lipid metabolism and cardiovascular mechanisms. Specifically, the top seven significant pathways are all closely related to LDL-C, which include total cholesterol ( $p = 2.24 \times$  $10^{-25}$ ), LDL-C ( $p = 6.23 \times 10^{-21}$ ), triglycerides (p = $1.89 \times 10^{-15}$ ), metabolite levels ( $p = 1.08 \times 10^{-13}$ ), lipid metabolism phenotypes ( $p = 1.93 \times 10^{-12}$ ), cholesterol metabolism ( $p = 1.82 \times 10^{-12}$ ), and high-density lipoprotein cholesterol (HDL-C) ( $p = 7.41 \times 10^{-11}$ ). Other pathways such as cholesterol metabolic process ( $p = 5.81 \times$  $10^{-6}$ ) and cholesterol homeostasis ( $p = 1.46 \times 10^{-6}$ ) are also among the top 50 significant pathways, aligning with LDL-C's central role in lipid regulation. Additionally,



Figure 5. Gene-expression imputation accuracy improvement over JTI in GTEx tissues (A) The average  $r^2$  increment of UTMOST and TransferTWAS compared to JTI. The average  $r^2$  values are calculated over all expressed genes in each tissue. (B) The average iGene ( $r^2 > 0.01$ ) number increment of UTMOST and TransferTWAS compared with JTI.

(b) The average form (7 > 0.07) number due utility (0 ST and Transfer WAS

(C) The proportion of JTI iGenes captured by UTMOST and TransferTWAS.

pathways linked to cardiovascular disease risk are also strongly enriched, including coronary heart disease ( $p = 5.14 \times 10^{-10}$ ), coronary artery disease ( $p = 3.08 \times 10^{-7}$ ), and cardiovascular disease risk factors ( $p = 1.72 \times 10^{-7}$ ). This is consistent with the clinical implications of elevated LDL-C. Interestingly, the enrichment of MHC class II antigen presentation and immune-related pathways (e.g., graft-versus-host disease) may reflect emerging links between lipid metabolism and inflammation. The strong functional overlap with established LDL-C-relevant pathways and disease mechanisms further validates the ability of TransferTWAS in capturing trait-relevant genes.

#### Real application to other complex traits

We tested the TWAS performance of UTMOST, JTI, and TransferTWAS in 30 other complex traits, including depressive symptoms, schizophrenia, and Alzheimer disease ( $N_{total}\approx 2.5$  million without adjusting for crossstudy sample overlap). These GWAS datasets were previously employed in Hu et al..<sup>10</sup> To identify the biologically most related tissues for each analyzed trait, Hu et al.<sup>10</sup> employed linkage-disequilibrium-score regression<sup>30</sup> and tissue-specific functional genome predicted by GenoSkyline-Plus annotations.<sup>31</sup> The results are listed in their Supplementary Table 24, and we used them to define the causal tissue of each trait.

TransferTWAS identified the greatest number of significant associations within biologically relevant tissues across 30 complex traits. As illustrated in Figure 6C, TransferTWAS outperformed competing methods, detecting substantially more associations in the most biologically relevant tissue for each trait. Specifically, TransferTWAS exhibited a 192.17% increase in associations compared to UTMOST and a 213.14% increase compared to JTI. Applying paired one-sided Wilcoxon tests on the number of associations identified by each method confirmed these improvements: TransferTWAS significantly found more associations compared with UTMOST ( $p = 2.64 \times 10^{-2}$ ), and JTI ( $p = 2.39 \times 10^{-2}$ ). In contrast, while UTMOST identified 23.3% more associations than JTI, this difference was not statistically significant (p = 0.2618). We list the number of associations identified in each trait in Table S8.

#### Discussion

The proposed TransferTWAS method aims to enhance gene-expression imputation accuracy by leveraging tissue-tissue similarity information. This approach borrows information from tissues with substantial sample sizes to improve predictions in tissues with limited samples. The performance of TransferTWAS was evaluated through extensive simulations and real data analysis using GTEx, GEUVADIS, ROS/MAP, and multiple GWAS datasets. We found that TransferTWAS can enhance the power of TWAS, and no evidence of inflated type-I error was observed. An enrichment analysis was conducted to



#### Figure 6. TWAS results of studying low-density lipoprotein cholesterol

(A) The number of genes that were significant under TransferTWAS, UTMOST, or JTI. Here, significance was defined as false discovery rate (FDR)-corrected p value of less than 0.05 ( $p_{\text{FDR}} < 0.05$ ).

(B) The number of predefined known LDL-C-related genes detected by the three methods.

(C) The number of genes identified in the biologically relevant tissue for each of the 30 complex traits. Each box includes two horizontal borders that represent the upper and lower quartiles and a solid line that represents the median. The highest and lowest points indicate the maximum and minimum values.

clarify the biological relevance between LDL-C and the iGenes identified by TransferTWAS.

Transfer learning has been applied in various areas of statistical genetics, such as enhancing prediction accuracy by leveraging pretrained polygenic risk score models.<sup>32,33</sup> TransferTWAS demonstrated improved TWAS power compared to other methods by leveraging eQTL effect-size information from multiple external tissues with similar genetic regulation profiles. The method's ability to effectively utilize external tissue information across various scenarios reinforces its potential as a powerful tool for enhancing TWAS imputation performance.

In simulation I on the GTEx dataset, TransferTWAS outperformed competing methods across various scenarios in terms of imputation accuracy. Its shrinkage-based approach, which avoids SNP selection during optimization, aligns with its strength under an infinitesimal model (many SNPs with weak effects), whereas regularization-based methods like UTMOST excel under sparse architectures (a few causal SNPs with strong effects). TransferTWAS outperformed UTMOST when more than 2% of SNPs were causal (Figures 2 and S4–S6), highlighting the limitations of regularization-based methods (such as UTMOST and JTI) as default choices. This is also supported by TIGAR (transcriptome-integrated genetic association resource),<sup>23,34</sup> which shows improved prediction accuracy over PrediXcan—a method relying on an elastic net model.<sup>35</sup>

In simulation studies II-V, TransferTWAS was evaluated for power and type-I error using ROS/MAP and GEUVADIS panels. It demonstrated superior TWAS power compared to UTMOST and JTI by leveraging eQTL effectsize information from tissues with similar genetic regulation profiles. While UTMOST lacks tissue similarity modeling and JTI relies on functional annotation, TransferTWAS's data-driven approach proved more effective, achieving the highest TWAS power across scenarios, even in tissue like EBV-transformed lymphocytes (Figure S1) with highly specific gene expression.<sup>11</sup> That is, its regulation is less influenced by cross-tissue expression information. Simulation studies III and V confirmed no inflated type-I error, and pathway enrichment analysis of iGenes revealed significant overlap with LDL-C-related pathways, indicating minimal false positives.

While TransferTWAS achieved lower imputation  $r^2$  than JTI in four tissues (Figure 5A), it increased the number of iGenes across all tissues (Figure 5B). This enabled more genes to enter the second step of TWAS, thereby enhancing the likelihood of identifying significant associations. In LDL-C analysis, TransferTWAS identified 898 associations missed by other methods. Given its primary goal of improving imputation in tissues with limited sample sizes, the lower  $r^2$  in specific tissues is less critical.

Several future directions for TransferTWAS warrant consideration. First, the method does not currently account for uncertainty in weight estimation during

gene-expression prediction. Recent studies have successfully incorporated such uncertainty by identifying ciseQTLs, performing fine-mapping to pinpoint key variants, and using the multivariate adaptive shrinkage (MASH) method to jointly estimate eQTL effects across tissues while incorporating tissue-specific uncertainty and correlations.<sup>36–42</sup> While TransferTWAS currently relies on standard ridge regression for tissue-specific effect estimation, integrating MASH could be a promising extension, although penalized regression methods (e.g., ridge, LASSO, and elastic net) face challenges in providing valid uncertainty due to biased estimates.43,44 Second, TransferTWAS could be enhanced by incorporating external eQTL summary-level data on tissue expression similarity. For instance, Zhang et al.<sup>45</sup> proposed a TWAS method that leverages eQTL summary-level data to improve gene-expression prediction accuracy. Since TransferTWAS only requires tissue-specific point estimates as input, it appears well suited for integrating such data. Third, addressing potential false-positive inflation in TWAS, as highlighted by recent studies,<sup>46,47</sup> could further improve TransferTWAS's reliability and accuracy.

In summary, we introduced TransferTWAS, a transfer learning algorithm that leverages GTEx data for geneexpression imputation. By improving imputation accuracy and TWAS power, TransferTWAS has the potential to advance our understanding of the genetic underpinnings of complex traits.

#### Data and code availability

- Project name: TransferTWAS
- Project homepage: https://github.com/daoyuan-lai/Transfer TWAS
- License: MIT license

#### Acknowledgments

This work was supported by the Hong Kong Research Grants Council General Research Fund (17307324). The Genotype-Tissue Expression (GTEx) project was supported by the Common Fund of the Office of the Director of the National Institutes of Health and by NCI, NHGRI, NHLBI, NIDA, NIMH, and NINDS. All protected data of the GTEx project are available through the database of Genotypes and Phenotypes (dbGaP) (accession number phs000424.v8.p2). We also thank Xiang Li for many insightful discussions.

#### Author contributions

Y.D.Z. and D.Y.L. conceived the project. D.Y.L. and H.W. implemented the method and performed the analyses. D.Y.L., H.W., T. G., S.W., D.J.L., P.C.S., and Y.D.Z. interpreted the results. D.Y.L. and Y.D.Z. drafted the original manuscript. All authors read and approved the final manuscript.

#### **Declaration of interests**

The authors declare no competing interests.

# Declaration of generative AI and AI-assisted technologies in the writing process

During the preparation of this work, the authors used ChatGPT in order to improve readability and language only. After using this tool, the authors reviewed and edited the content as needed and take full responsibility for the content of the publication.

#### Supplemental information

Supplemental information can be found online at https://doi.org/10.1016/j.ajhg.2025.06.006.

#### Web resources

- 1000 Genomes Project on GRCh38, https://www.international genome.org/
- bigsnpr, https://privefl.github.io/bigsnpr/
- dbSNP (build 151) with GRCh37.p13 as reference assembly, https://ftp-ncbi-nih-gov.eproxy.lib.hku.hk/snp/organisms/ human\_9606\_b151\_GRCh37p13/
- GEUVADIS, https://www.ebi.ac.uk/arrayexpress/experiments/E-GEUV-1
- GTEx v.8, https://gtexportal.org/home/
- JTI, https://github.com/gamazonlab/MR-JTI/tree/master/model\_ training/JTI

LDL-C, http://biobank.ctsu.ox.ac.uk/crystal/field.cgi?id=30790

Online Mendelian Inheritance in Man, https://www.omim.org/

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

UTMOST, https://github.com/yiminghu/CTIMP/blob/master/main.R

Received: December 26, 2024 Accepted: June 10, 2025

#### References

- 1. Gallagher, M.D., and Chen-Plotkin, A.S. (2018). The post-GWAS era: from association to function. Am. J. Hum. Genet. *102*, 717–730.
- 2. Maurano, M.T., Humbert, R., Rynes, E., Thurman, R.E., Haugen, E., Wang, H., Reynolds, A.P., Sandstrom, R., Qu, H., Brody, J., et al. (2012). Systematic localization of common disease-associated variation in regulatory DNA. Science *337*, 1190–1195.
- **3.** 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., et al.; GTEx Consortium (2015). A gene-based association method for mapping traits using reference transcriptome data. Nat. Genet. *47*, 1091–1098.
- Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B. W.J.H., Jansen, R., de Geus, E.J.C., Boomsma, D.I., Wright, F. A., et al. (2016). Integrative approaches for large-scale transcriptome-wide association studies. Nat. Genet. 48, 245–252.
- Lonsdale, J., Thomas, J., Salvatore, M., Phillips, R., Lo, E., Shad, S., Hasz, R., Walters, G., Garcia, F., Young, N., et al. (2013). The genotype-tissue expression (GTEx) project. Nat. Genet. 45, 580–585.

- 6. GTEx Consortium, Ardlie, K.G., Deluca, D.S., Segrè, A.V., Sullivan, T.J., Young, T.R., Gelfand, E.T., Trowbridge, C.A., Maller, J.B., and Tukiainen, T. (2015). The Genotype-Tissue Expression (GTEx) pilot analysis: multitissue gene regulation in humans. Science *348*, 648–660.
- Aguet, F., Brown, A.A., Castel, S.E., Davis, J.R., He, Y., Jo, B., Mohammadi, P., Park, Y., Parsana, P., Segrè, A.V., et al. (2017). Genetic effects on gene expression across human tissues. Nature 550, 204–213.
- GTEx Consortium (2020). The GTEx Consortium atlas of genetic regulatory effects across human tissues. Science 369, 1318–1330.
- **9.** Wainberg, M., Sinnott-Armstrong, N., Mancuso, N., Barbeira, A.N., Knowles, D.A., Golan, D., Ermel, R., Ruusalepp, A., Quertermous, T., Hao, K., et al. (2019). Opportunities and challenges for transcriptome-wide association studies. Nat. Genet. *51*, 592–599.
- Hu, Y., Li, M., Lu, Q., Weng, H., Wang, J., Zekavat, S.M., Yu, Z., Li, B., Gu, J., Muchnik, S., et al. (2019). A statistical framework for cross-tissue transcriptome-wide association analysis. Nat. Genet. *51*, 568–576.
- 11. Zhou, D., Jiang, Y., Zhong, X., Cox, N.J., Liu, C., and Gamazon, E.R. (2020). A unified framework for joint-tissue transcriptome-wide association and Mendelian randomization analysis. Nat. Genet. *52*, 1239–1246.
- 12. Li, B., Veturi, Y., Verma, A., Bradford, Y., Daar, E.S., Gulick, R. M., Riddler, S.A., Robbins, G.K., Lennox, J.L., Haas, D.W., and Ritchie, M.D. (2021). Tissue specificity-aware TWAS (TSA-TWAS) framework identifies novel associations with metabolic, immunologic, and virologic traits in HIV-positive adults. PLoS Genet. 17, e1009464.
- **13.** Mai, J., Lu, M., Gao, Q., Zeng, J., and Xiao, J. (2023). Transcriptome-wide association studies: recent advances in methods, applications and available databases. Commun. Biol. *6*, 899.
- Barbeira, A.N., Pividori, M., Zheng, J., Wheeler, H.E., Nicolae, D.L., and Im, H.K. (2019). Integrating predicted transcriptome from multiple tissues improves association detection. PLoS Genet. 15, e1007889.
- Gu, T., Han, Y., and Duan, R. (2024). Robust angle-based transfer learning in high dimensions. J. Roy. Stat. Soc. B Stat. Methodol. qkae111. https://doi.org/10.1093/jrsssb/ qkae111.
- **16.** Wang, A., Tian, P., and Zhang, Y.D. (2024). TWAS-GKF: a novel method for causal gene identification in transcriptome-wide association studies with knockoff inference. Bioinformatics *40*, btae502.
- 17. 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*, S161–S189.
- 18. Lappalainen, T., Sammeth, M., Friedländer, M.R., 't Hoen, P. A.C., Monlong, J., Rivas, M.A., Gonzàlez-Porta, M., Kurbatova, N., Griebel, T., Ferreira, P.G., et al. (2013). Transcriptome and genome sequencing uncovers functional variation in humans. Nature *501*, 506–511.
- **19.** Wang, H., Li, X., Li, T., Li, Z., Sham, P.C., and Zhang, Y.D. (2025). MAAT: a new nonparametric Bayesian framework for incorporating multiple functional annotations in transcriptome-wide association studies. Genome Biol. *26*, 21.
- Keys, K.L., Mak, A.C.Y., White, M.J., Eckalbar, W.L., Dahl, A. W., Mefford, J., Mikhaylova, A.V., Contreras, M.G., Elha-

wary, J.R., Eng, C., et al. (2020). On the cross-population generalizability of gene expression prediction models. PLoS Genet. *16*, e1008927.

- **21.** Khunsriraksakul, C., McGuire, D., Sauteraud, R., Chen, F., Yang, L., Wang, L., Hughey, J., Eckert, S., Dylan Weissenkampen, J., Shenoy, G., et al. (2022). Integrating 3D genomic and epigenomic data to enhance target gene discovery and drug repurposing in transcriptome-wide association studies. Nat. Commun. *13*, 3258.
- **22.** Feng, H., Mancuso, N., Gusev, A., Majumdar, A., Major, M., Pasaniuc, B., and Kraft, P. (2021). Leveraging expression from multiple tissues using sparse canonical correlation analysis and aggregate tests improves the power of transcriptome-wide association studies. PLoS Genet. *17*, e1008973.
- 23. Nagpal, S., Meng, X., Epstein, M.P., Tsoi, L.C., Patrick, M., Gibson, G., De Jager, P.L., Bennett, D.A., Wingo, A.P., Wingo, T.S., and Yang, J. (2019). TIGAR: an improved Bayesian tool for transcriptomic data imputation enhances gene mapping of complex traits. Am. J. Hum. Genet. *105*, 258–266.
- 24. Barbeira, A.N., Dickinson, S.P., Bonazzola, R., Zheng, J., Wheeler, H.E., Torres, J.M., Torstenson, E.S., Shah, K.P., Garcia, T., Edwards, T.L., et al. (2018). Exploring the phenotypic consequences of tissue specific gene expression variation inferred from GWAS summary statistics. Nat. Commun. 9, 1825.
- **25.** Weiss, K., Khoshgoftaar, T.M., and Wang, D. (2016). A survey of transfer learning. J. Big Data *3*, 9–40.
- **26.** Wu, L., Shi, W., Long, J., Guo, X., Michailidou, K., Beesley, J., Bolla, M.K., Shu, X.-O., Lu, Y., Cai, Q., et al. (2018). A transcriptome-wide association study of 229,000 women identifies new candidate susceptibility genes for breast cancer. Nat. Genet. *50*, 968–978.
- 27. Gaudet, D., Gipe, D.A., Pordy, R., Ahmad, Z., Cuchel, M., Shah, P.K., Chyu, K.-Y., Sasiela, W.J., Chan, K.-C., Brisson, D., et al. (2017). ANGPTL3 inhibition in homozygous familial hypercholesterolemia. N. Engl. J. Med. *377*, 296–297.
- **28.** Peloso, G.M., Nomura, A., Khera, A.V., Chaffin, M., Won, H.-H., Ardissino, D., Danesh, J., Schunkert, H., Wilson, J.G., Samani, N., et al. (2019). Rare protein-truncating variants in APOB, lower low-density lipoprotein cholesterol, and protection against coronary heart disease. Circ. Genom. Precis. Med. *12*, e002376.
- **29.** Francis, M., Li, C., Sun, Y., Zhou, J., Li, X., Brenna, J.T., and Ye, K. (2021). Genome-wide association study of fish oil supplementation on lipid traits in 81,246 individuals reveals new gene-diet interaction loci. PLoS Genet. *17*, e1009431.
- **30.** Finucane, H.K., Bulik-Sullivan, B., Gusev, A., Trynka, G., Reshef, Y., Loh, P.-R., Anttila, V., Xu, H., Zang, C., Farh, K., et al. (2015). Partitioning heritability by functional annotation using genome-wide association summary statistics. Nat. Genet. *47*, 1228–1235.
- **31.** Lu, Q., Powles, R.L., Abdallah, S., Ou, D., Wang, Q., Hu, Y., Lu, Y., Liu, W., Li, B., Mukherjee, S., et al. (2017). Systematic tissue-specific functional annotation of the human genome highlights immune-related DNA elements for late-onset Alzheimer's disease. PLoS Genet. *13*, e1006933.
- Zhao, Z., Fritsche, L.G., Smith, J.A., Mukherjee, B., and Lee, S. (2022). The construction of cross-population polygenic risk scores using transfer learning. Am. J. Hum. Genet. *109*, 1998–2008.
- 33. Tian, P., Chan, T.H., Wang, Y.-F., Yang, W., Yin, G., and Zhang, Y.D. (2022). Multiethnic polygenic risk prediction

in diverse populations through transfer learning. Front. Genet. 13, 906965.

- 34. Parrish, R.L., Buchman, A.S., Tasaki, S., Wang, Y., Avey, D., Xu, J., De Jager, P.L., Bennett, D.A., Epstein, M.P., and Yang, J. (2024). SR-TWAS: leveraging multiple reference panels to improve transcriptome-wide association study power by ensemble machine learning. Nat. Commun. 15, 6646.
- Zou, H., and Hastie, T. (2005). Regularization and variable selection via the elastic net. J. Roy. Stat. Soc. B Stat. Methodol. 67, 301–320.
- **36.** Gao, G., Fiorica, P.N., McClellan, J., Barbeira, A.N., Li, J.L., Olopade, O.I., Im, H.K., and Huo, D. (2023). A joint transcriptome-wide association study across multiple tissues identifies candidate breast cancer susceptibility genes. Am. J. Hum. Genet. *110*, 950–962.
- 37. Gao, G., McClellan, J., Barbeira, A.N., Fiorica, P.N., Li, J.L., Mu, Z., Olopade, O.I., Huo, D., and Im, H.K. (2024). A multi-tissue, splicing-based joint transcriptome-wide association study identifies susceptibility genes for breast cancer. Am. J. Hum. Genet. *111*, 1100–1113.
- 38. Li, J.L., McClellan, J.C., Zhang, H., Gao, G., and Huo, D. (2024). Multi-tissue transcriptome-wide association studies identified 235 genes for intrinsic subtypes of breast cancer. J. Natl. Cancer Inst. *116*, 1105–1115.
- **39.** McClellan, J.C., Li, J.L., Gao, G., and Huo, D. (2024). Expression-and splicing-based multi-tissue transcriptome-wide association studies identified multiple genes for breast cancer by estrogen-receptor status. Breast Cancer Res. *26*, 51.
- 40. Araujo, D.S., Nguyen, C., Hu, X., Mikhaylova, A.V., Gignoux, C., Ardlie, K., Taylor, K.D., Durda, P., Liu, Y., Papanicolaou, G., et al. (2023). Multivariate adaptive shrinkage improves

cross-population transcriptome prediction and association studies in underrepresented populations. HGG Adv. *4*, 100216.

- **41.** Chen, D.M., Dong, R., Kachuri, L., Hoffmann, T.J., Jiang, Y., Berndt, S.I., Shelley, J.P., Schaffer, K.R., Machiela, M.J., Freedman, N.D., et al. (2024). Transcriptome-wide association analysis identifies candidate susceptibility genes for prostate-specific antigen levels in men without prostate cancer. HGG Adv. *5*, 100315.
- **42.** Urbut, S.M., Wang, G., Carbonetto, P., and Stephens, M. (2019). Flexible statistical methods for estimating and testing effects in genomic studies with multiple conditions. Nat. Genet. *51*, 187–195.
- **43.** Zhang, C.-H., and Zhang, S.S. (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. J. Roy. Stat. Soc. B Stat. Methodol. *76*, 217–242.
- Javanmard, A., and Montanari, A. (2014). Confidence intervals and hypothesis testing for high-dimensional regression. J. Mach. Learn. Res. *15*, 2869–2909.
- **45.** Zhang, Z., Bae, Y.E., Bradley, J.R., Wu, L., and Wu, C. (2022). SUMMIT: An integrative approach for better transcriptomic data imputation improves causal gene identification. Nat. Commun. *13*, 6336.
- **46.** de Leeuw, C., Werme, J., Savage, J.E., Peyrot, W.J., and Posthuma, D. (2023). On the interpretation of transcriptomewide association studies. PLoS Genet. *19*, e1010921.
- 47. Liang, Y., Nyasimi, F., and Im, H.K. (2024). Pervasive polygenicity of complex traits inflates false positive rates in transcriptome-wide association studies. Preprint at bioRxiv. https://doi.org/10.1101/2023.2010.2017.562831.