# Assessing the genetic effect mediated through gene expression from summary eQTL and GWAS data

Richard Barfield[1], Helian Feng[2], Alexander Gusev[3,4], Lang Wu[5], Wei Zheng[5], Bogdan Pasaniuc[6,7,8], Peter Kraft[2,9,10]

Author Affiliations

1) Public Health Sciences Division, Fred Hutchinson Cancer Research Center, Seattle, WA, USA.
2) Department of Biostatistics, Harvard T.H. Chan School of Public Health, Boston, MA, USA
3) Division of Population Sciences, Dana-Farber Cancer Institute
4) Division of Genetics, Brigham & Women's Hospital
5) Division of Epidemiology, Department of Medicine, Vanderbilt Epidemiology Center, Vanderbilt-Ingram Cancer Center, Vanderbilt University School of Medicine, Nashville, TN, USA.
6) Department of Pathology & Laboratory Medicine, David Geffen School of Medicine, University of California, Los Angeles, Los Angeles, CA, USA
7) Bioinformatics Interdepartmental Program, University of California, Los Angeles, Los Angeles, CA, USA
8) Department of Human Genetics, David Geffen School of Medicine, University of California, Los Angeles, Los Angeles, CA, USA
9) Program in Genetic Epidemiology and Statistical Genetics, Harvard T.H. Chan School of Public Health, Boston, MA, USA
10) Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, MA, USA

**Corresponding Author: Peter Kraft. (617-432-4271)**

**655 Huntington Avenue**

**Building II Room 249A**

**Boston, Massachusetts 02115**

# Abstract

Integrating genome-wide association (GWAS) and expression quantitative trait locus (eQTL) data can boost power to detect novel disease loci or pinpoint the susceptibility gene at a known disease locus. However, it is often the case that multiple eQTL genes co-localize at disease loci (an effect of linkage disequilibrium, LD), making the identification of the true susceptibility gene challenging. To distinguish between true susceptibility genes (i.e. when the genetic effect on phenotype is mediated through expression) and spurious co-localizations, we developed an approach to quantify the genetic effect mediated through expression. Our approach can be viewed as an extension of the standard Mendelian randomization Egger technique to incorporate LD among variants while only requiring summary association data (both GWAS and eQTL) along with LD from reference panels. Through simulations we show that when eQTLs have pleiotropic or LD-confounded effects on disease, our approach provides adequate control of Type I error, more power, and less bias than previously-proposed methods. When there is no effect of gene expression on disease, our method has the desired Type I Error, while LD-aware Mendelian randomization, which assumes no pleiotropy, can have inflated Type I Error. In the presence of direct effect of genetic variants on traits, our approach attained up to 3x greater power than the standard approaches while properly controlling Type I error. To illustrate our method, we analyzed recent large scale breast cancer GWAS with gene expression in breast tissue from GTEx.

Key Words: Genome-wide association study; Gene Expression, Mendelian Randomization, transciptome-wide association study

# Introduction

Integrating data from genome-wide association studies (GWAS) and expression quantitative trait loci (eQTL) studies can help in detecting novel disease loci and pinpointing genes of interest. This is done by aggregating association signals across multiple SNPs associated with transcript levels, under the assumption that the association between SNPs and disease may be mediated through the expression of particular genes. However, there is often insufficient data on the same set of individuals for the disease of interest with SNPs and gene expression in the relevant tissue . In contrast, there are large amounts of publicly available summary statistics data from separate GWAS and eQTL studies and there has been a rise in the number of statistical methods available to utilize and jointly analyze these summary data [Barbeira, et al. 2016; Gamazon, et al. 2015; Giambartolomei, et al. 2014; Gusev, et al. 2016; Zhu, et al. 2016].  While these methods can identify association between genetically-predicted expression levels and disease, they cannot distinguish between mediation (the SNPs affect disease risk through their effects on a particular gene's expression) and co-localization (e.g. the eQTLs and causal disease SNPs are distinct but in linkage disequilibrium causing spurious associations). As studies are now identifying hundreds of loci with putative expression-disease associations [Mancuso, et al. 2017], this misspecification may lead to many significant associations even though the tested gene expression levels do not affect the outcome. In contrast, more recent work has proposed methods to estimate the genetic effect mediated by gene expression on complex traits but only do so in aggregate across the genome [O'Connor, et al. 2017].

Previous approaches to estimate the association between a mediator and an outcome using summary statistics on SNP-mediator and SNP-outcome associations are not appropriate when the mediator is gene expression; some assume the SNPs are not in linkage disequilibrium [Bowden, et al. 2015] while others account for linkage disequilibrium but assume the SNPs have no direct effect on outcome [Burgess, et al. 2016]. To distinguish between true susceptibility genes (i.e. when the genetic effect on phenotype is mediated through expression) and spurious co-localizations, we developed LD-aware Mendelian Randomization Egger regression (LDA MR-Egger), an extension of MR-Egger regression [Bowden, et al. 2015] to multiple SNPs in

linkage disequilibrium. We first introduce the model relating the outcome and gene expression to the SNPs. We next discuss four existing approaches for testing and/or estimating the magnitude of? the association between the outcome and the gene of interest, and introduce our new LDA MR-Egger regression. The statistical properties of these estimates and their performance are then examined in presence or absence of co-localization. We perform an empirical study to assess the type I error, power, and bias of the estimates, showing that when there is a constant direct effect the LDA MR-Egger has correct type I error while existing approaches do not. Finally, we apply the various approaches to summary statistics from a GWAS on breast cancer [Michailidou, et al. 2017] with eQTL data from a breast tissue panel in GTEx [Consortium 2013]. Our results demonstrate that the LDA MR-Egger has the proper type I error without loss of power compared to the LDA MR under more simulation scenarios, and proper type I error when there is a constant direct effect of SNPs on outcome.

## Methods

### Main Models

Let $Y$ denote our outcome ($nx1$), $M$ the mediator ($nx1$), and $G$ the SNP matrix ($nxJ$) of interest. We assume that the columns of $G$ have been standardized to have mean zero and variance one. If $G$ has not been standardized and the GWAS effect estimates are on the minor allele counts, we can transform the effects to what would have been observed if the SNPs had been standardized (Appendix A). We denote the LD structure of the SNPs $G$ as $\Sigma$, a $JxJ$ symmetric positive definite matrix. This matrix can be after a regularization such as ridge regression to make the matrix symmetric positive definite [Pasaniuc and Price 2017]. For a link function $g$, we have the following models relating $M, G,$ and $Y$.

$$g\big(E(Y|M,G)\big) = \gamma_0 + M\gamma + G\theta \tag{i.a.}$$

$$g\big(E(Y|G)\big) = \gamma_0^* + G\beta_G \tag{i.b.}$$

$$M = \beta_0 + G\beta_E + \epsilon_M; \epsilon_M \sim N(0, I_n\sigma^2) \tag{i.c.}$$

In the above model, $\theta$ is a $J$-column vector of $G$ effects on $Y$ conditional on $M$, $\gamma$ is the effect of $M$ on $Y$ conditional on $G$, $\beta_E$ is the $J$-column vector of SNP effects on the mediator, and $\beta_G$ is the $J$-vector of the $G$ effects marginal over $M$. $\beta_G$ and $\beta_E$ represent mutually-conditioned SNP

effects on outcome and gene expression respectively. $\epsilon_M$ represents the residual variance in $M$ and $I_n$ represents the $n$ x $n$ identity matrix. We are interested in the situation where we cannot directly estimate the parameters in model $(i.a)$, as we do not have complete data on $Y$, $M$ and $G$ from (sufficiently many) individuals. We want to derive inference on $\gamma$, because if $\gamma \neq 0$, the gene affects the trait. To relate the parameters in the marginal model of $Y$ and the marginal model of $M$ we assume one of the following for the remainder of the paper:

1) $g$ is either the log or identity link function.
2) $Y$ is a sufficiently rare binary trait and $g$ is the logit link.

If either of the two conditions above hold, we will have the following:

$$\boldsymbol{\beta}_G \approx \boldsymbol{\beta}_E \gamma + \boldsymbol{\theta}. \textbf{ (ii.)}$$

If $g$ is the log or linear link, the approximation will be exact. Equation (ii) suggests that the effects of $G$ marginal over $M$ are a function of the eQTL statistics ($\boldsymbol{\beta}_E$), the effect of the gene expression on the outcome conditional on $G$ ($\gamma$), and the effect of $G$ on the outcome conditional on gene expression ($\boldsymbol{\theta}$). We will call $\boldsymbol{\theta}$ the direct effect and $\boldsymbol{\beta}_E \gamma$ the mediated effect. If the necessary causal assumptions are met, these parameters could be interpreted in the causal/counterfactual framework [Valeri and VanderWeele 2013], but for this article it is not necessary.

In practice, we do not have an overlap of data on $M$, $Y$, and $G$. Instead, we have a sample of size $N$ that the GWAS is run on to estimate $\boldsymbol{\beta}_G$ and an independent sample of size $N_E$ that was used to estimate $\boldsymbol{\beta}_E$. We therefore cannot estimate $\boldsymbol{\theta}$ directly. Moreover, we typically only have estimates of individual SNP effects marginal over the other SNPs, $\widehat{\boldsymbol{\beta}}_E^*$ and $\widehat{\boldsymbol{\beta}}_G^*$. For our purposes, $\hat{\beta}_{G,j}^*$ and $\hat{\beta}_{E,j}^*$ were estimated with the same reference allele for SNP $j$. If they were not, the sign of the effect can be changed so as to refer to the same reference allele.

The formulas given above relating the mean of $Y$ and $M$ to $G$ were given on the conditional level. We therefore transform our marginal estimates ($\widehat{\boldsymbol{\beta}}_E^*$ and $\widehat{\boldsymbol{\beta}}_G^*$) to the conditional scale. Given an estimate of the LD matrix ($\boldsymbol{\Sigma}$) we estimate the conditional eQTL and GWAS effects as $\widehat{\boldsymbol{\beta}}_E = \boldsymbol{\Sigma}^{-1}\widehat{\boldsymbol{\beta}}_E^*$ and $\widehat{\boldsymbol{\beta}}_G = \boldsymbol{\Sigma}^{-1}\widehat{\boldsymbol{\beta}}_G^*$ [Shi, et al. 2016]. If the marginal effect estimates were not calculated on the standardized genotypes they can be transformed using the MAF (Appendix

A). We assume that $\widehat{\boldsymbol{\beta}}_E$ and $\widehat{\boldsymbol{\beta}}_G$ are unbiased for $\boldsymbol{\beta}_E$ and $\boldsymbol{\beta}_G$ respectively.

We do not know the form of $\boldsymbol{\theta}$; it could be a constant or vary by SNP. All we know is that it is a vector of length $J$. Our estimated GWAS effects, (given our assumptions above) are a function of $\boldsymbol{\theta}, \boldsymbol{\beta}_E, \gamma$ and the sampling error:

$$\widehat{\boldsymbol{\beta}}_G \approx \boldsymbol{\theta} + \boldsymbol{\beta}_E \gamma + \boldsymbol{\epsilon}_G, \boldsymbol{\epsilon}_G \sim N(\mathbf{0}, \boldsymbol{\Sigma}_G).$$

If the SNPs are not in LD, then the marginal and the conditional will be equal ($\widehat{\boldsymbol{\beta}}_G \approx \widehat{\boldsymbol{\beta}}_G^*$).

We next derive the covariance of $\widehat{\boldsymbol{\beta}}_G$ ($\boldsymbol{\Sigma}_G$). Let $\boldsymbol{\sigma}_G^*$ denote a $Jx1$ vector of the marginal standard errors of $\widehat{\boldsymbol{\beta}}_G^*$. Let "·" denote element wise multiplication between two matrices. As $\boldsymbol{G}$ has been standardized, $cov(\widehat{\boldsymbol{\beta}}_G^*) = \boldsymbol{\Sigma} \cdot \boldsymbol{\sigma}_G^* \boldsymbol{\sigma}_G^{*T}$. This gives the covariance of our conditional GWAS estimates:

$$\boldsymbol{\Sigma}_G = cov(\widehat{\boldsymbol{\beta}}_G) = cov(\boldsymbol{\Sigma}^{-1}\widehat{\boldsymbol{\beta}}_G^*) = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\Sigma} \cdot \boldsymbol{\sigma}_G^* \boldsymbol{\sigma}^{*T}_G)\boldsymbol{\Sigma}^{-1}.$$

If $\boldsymbol{\sigma}_G^* = v\mathbf{1}_J$, where $\mathbf{1}_J$ is a column vector of ones, then $\boldsymbol{\Sigma}_G = \boldsymbol{\Sigma}^{-1}v^2$. Our goal is to derive a valid test and also estimate $\gamma$ in the presence of direct effects of the SNPs onto the outcome. We next go over approaches for testing for an association between gene expression and outcome.

**Transcriptome Wide Association Studies (TWAS)**

The Transcriptome Wide Association Study (TWAS) statistic uses summary statistics to test for an association between genetically predicted gene expression and a phenotype of interest [Gusev, et al. 2016]. The TWAS does not necessarily estimate the $\gamma$ above but does provide a valid test for the association between the gene of interest and the outcome. The TWAS test statistic is:

$$Z_{TWAS} = \frac{\widehat{\boldsymbol{\beta}}_E^{*T} \boldsymbol{\Sigma}^{-1} \boldsymbol{Z}_G^*}{\sqrt{\widehat{\boldsymbol{\beta}}_E^{*T} \boldsymbol{\Sigma}^{-1} \widehat{\boldsymbol{\beta}}_E^*}},$$

where $\boldsymbol{Z}_G^* = \widehat{\boldsymbol{\beta}}_G^* \cdot (\boldsymbol{\sigma}_G^*)^{-1}$, is a column vector of the marginal test statistics. The test statistic is then compared to a standard normal to assess significance. The TWAS can use either the marginal or conditional eQTL estimates as weights. If using the conditional, the equation above becomes:

$$Z_{TWAS} = \frac{\hat{\boldsymbol{\beta}}_E^T \boldsymbol{Z}_G^*}{\sqrt{\hat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma} \hat{\boldsymbol{\beta}}_E}}.$$

## The Summary Mendelian Randomization Estimator

The summary Mendelian Randomization estimator (also known as the "Toby Johnson" estimator [Johnson 2011]) is:

$$\hat{\gamma}_{MR} = \frac{\hat{\boldsymbol{\beta}}_E^{*T} \boldsymbol{V}^{-1} \hat{\boldsymbol{\beta}}_G^*}{\hat{\boldsymbol{\beta}}_E^{*T} \boldsymbol{V}^{-1} \hat{\boldsymbol{\beta}}_E^*},$$

where $\boldsymbol{V}$ is a diagonal matrix with $v_{jj} = \sigma_{G,j}^{2*} = var(\hat{\beta}_{G,j}^*)$. If there are no direct effects, i.e. $\boldsymbol{\theta} = 0$, and the SNPs are not in LD with each other ($\boldsymbol{\Sigma} = \boldsymbol{I}_J$), then this will be an unbiased estimate of $\gamma$. This estimate (from here referred to as the MR estimate) can be rewritten as:

$$\hat{\gamma}_{MR} = \frac{\sum_{j=1}^{J} \hat{\beta}_{E,j}^* \hat{\beta}_{G,j}^* \sigma_{G,j}^{-2*}}{\sum_{j=1}^{J} \hat{\beta}_{E,j}^{2*} \sigma_{G,j}^{-2*}}.$$

We estimate its variance as:

$$var(\hat{\gamma}_{MR}) = \frac{\hat{\sigma}_{MR}^2}{\sum_{j=1}^{J} \hat{\beta}_{E,j}^{2*} \sigma_{G,j}^{-2*}},$$

$$\hat{\sigma}_{MR}^2 = \frac{1}{J-1} \sum_{j=1}^{J} \sigma_{G,j}^{-2*} \left( \hat{\beta}_{G,j}^* - \hat{\beta}_{E,j}^* \hat{\gamma}_{MR} \right)^2.$$

The MR test statistic is then:

$$Z_{MR} = \frac{\hat{\gamma}_{MR}}{\sqrt{var(\hat{\gamma}_{MR})}} = \frac{\sum_{j=1}^{J} \hat{\beta}_{E,j}^* \hat{\beta}_{G,j}^* \sigma_{G,j}^{-2*}}{\hat{\sigma}_{MR} \sqrt{\sum_{j=1}^{J} \hat{\beta}_{E,j}^{2*} \sigma_{G,j}^{-2*}}}.$$

For testing, we compare $Z_{MR}$ to the quantiles of a $t$-distribution with $J$-1 degrees of freedom. If the SNPs are in LD or if there are direct effects, $\hat{\gamma}_{MR}$ can lead to incorrect inference. The MR estimate can be framed as a weighted linear regression without an intercept of the marginal GWAS estimates on the marginal eQTL estimates with weights equal to $\sigma_{G,j}^{-2*}$.

We note here that Zhu et al refer to their single SNP MR test as a summary Mendelian randomization, that approach is not the same as the Toby Johnson but just a difference in naming [Zhu, et al. 2016]. The Zhu et al approach provides an estimate for pleiotropy (HEIDI) based on the summary statistics and conditioning on the lead SNP.

**MR-Egger Estimate**

If the MR estimate can be thought of as a weighted linear regression without an intercept, the MR-Egger extends the MR by including an intercept ($\alpha$) to the weighted linear regression. It assumes that $E\left(\widehat{\boldsymbol{\beta}}_G^* \middle| \widehat{\boldsymbol{\beta}}_E^*\right) = \alpha \mathbf{1}_J + \widehat{\boldsymbol{\beta}}_E^* \gamma$. The estimates are (same $\mathbf{V}$ as the MR estimate):

$$\begin{bmatrix} \hat{\alpha} \\ \hat{\gamma}_{MRE} \end{bmatrix} = \begin{bmatrix} \mathbf{1}_J^T V^{-1} \mathbf{1}_J & \mathbf{1}_J^T V^{-1} \widehat{\boldsymbol{\beta}}_E^* \\ \widehat{\boldsymbol{\beta}}_E^{*T} V^{-1} \mathbf{1}_J & \widehat{\boldsymbol{\beta}}_E^{*T} V^{-1} \widehat{\boldsymbol{\beta}}_E^* \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{1}_J^T V^{-1} \\ \widehat{\boldsymbol{\beta}}_E^{*T} V^{-1} \end{bmatrix} \widehat{\boldsymbol{\beta}}_G^*,$$

The estimate of $\gamma$ is:

$$\hat{\gamma}_{MRE} = \frac{\left(\sum_{j=1}^J \sigma_{G,j}^{-2*}\right)\left(\sum_{j=1}^J \hat{\beta}_{E,j}^* \hat{\beta}_{G,j}^* \sigma_{G,j}^{-2*}\right) - \left(\sum_{j=1}^J \hat{\beta}_{E,j} \sigma_{G,j}^{-2*}\right)\left(\sum_{j=1}^J \hat{\beta}_{G,j}^* \sigma_{G,j}^{-2*}\right)}{\left(\sum_{j=1}^J \sigma_{G,j}^{-2*}\right)\left(\sum_{j=1}^J \hat{\beta}_{E,j}^{*2} \sigma_{G,j}^{-2*}\right) - \left(\sum_{j=1}^J \hat{\beta}_{E,j}^* \sigma_{G,j}^{-2*}\right)^2}.$$

The variance of the estimate is:

$$var(\hat{\gamma}_{MRE}) = \hat{\sigma}_{MRE}^2 \frac{\sum_{j=1}^J \sigma_{G,j}^{-2*}}{\left(\sum_{j=1}^J \sigma_{G,j}^{-2*}\right)\left(\sum_{j=1}^J \hat{\beta}_{E,j}^{*2} \sigma_{G,j}^{-2*}\right) - \left(\sum_{j=1}^J \hat{\beta}_{E,j}^* \sigma_{G,j}^{-2*}\right)^2},$$

$$\hat{\sigma}_{MRE}^2 = \frac{1}{J-2} \sum_{j=1}^J \sigma_{G,j}^{-2*} \left(\hat{\beta}_{G,j}^* - \hat{\alpha} - \hat{\beta}_{E,j}^* \hat{\gamma}_{MRE}\right)^2.$$

The MR-Egger test statistic is then:

$$Z_{MRE} = \frac{\left(\sum_{j=1}^J \sigma_{G,j}^{-2*}\right)\left(\sum_{j=1}^J \hat{\beta}_{E,j}^* \hat{\beta}_{G,j}^* \sigma_{G,j}^{-2*}\right) - \left(\sum_{j=1}^J \hat{\beta}_{E,j} \sigma_{G,j}^{-2*}\right)\left(\sum_{j=1}^J \hat{\beta}_{G,j}^* \sigma_{G,j}^{-2*}\right)}{\hat{\sigma}_{MRE} \sqrt{\sum_{j=1}^J \sigma_{G,j}^{-2*}} \sqrt{\left(\sum_{j=1}^J \sigma_{G,j}^{-2*}\right)\left(\sum_{j=1}^J \hat{\beta}_{E,j}^{*2} \sigma_{G,j}^{-2*}\right) - \left(\sum_{j=1}^J \hat{\beta}_{E,j}^* \sigma_{G,j}^{-2*}\right)^2}}.$$

To test, we compare to the quantiles of a *t*-distribution with *J*-2 degrees of freedom. If the SNPs are in LD, the test for $\hat{\gamma}_{MRE}$ may lead to incorrect inference due to the variance of $\hat{\gamma}_{MRE}$ being misspecified.

## LD-Aware MR (LDA MR) Estimate

The LDA MR estimator of $\gamma$ extends the MR estimator by relaxing the assumption of the SNPs being independent. It still requires that there are no direct effects. Recall that $\boldsymbol{\Sigma}_G = cov(\widehat{\boldsymbol{\beta}}_G)$. The LDA MR estimator is then:

$$\hat{\gamma}_{LDMR} = \frac{\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \widehat{\boldsymbol{\beta}}_G}{\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \widehat{\boldsymbol{\beta}}_E},$$

And we estimate the variance as:

$$var(\hat{\gamma}_{LDMR}) = \frac{\hat{\sigma}^2_{LDMR}}{\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \widehat{\boldsymbol{\beta}}_E},$$

$$\hat{\sigma}^2_{LDMR} = \frac{1}{J-1} \left(\widehat{\boldsymbol{\beta}}_G - \widehat{\boldsymbol{\beta}}_E \hat{\gamma}_{LDMR}\right)^T \boldsymbol{\Sigma}_G^{-1} \left(\widehat{\boldsymbol{\beta}}_G - \widehat{\boldsymbol{\beta}}_E \hat{\gamma}_{LDMR}\right).$$

The test statistic is then:

$$Z_{LDMR} = \frac{\hat{\gamma}_{LDA-MR}}{\sqrt{var(\hat{\gamma}_{LDMR})}} = \frac{\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \widehat{\boldsymbol{\beta}}_G}{\hat{\sigma}_{LDMR}\sqrt{\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \widehat{\boldsymbol{\beta}}_E}}.$$

Similar to the MR, we compare to a $t$ distribution with $J$-1 df. If there are direct effects, $\hat{\gamma}_{LDMR}$ can lead to incorrect inference. Just as the MR estimate is a weighted linear regression without an intercept, the LDA MR estimate is a weighted linear regression without an intercept using the weight matrix $\boldsymbol{\Sigma}_G^{-1}$. If $\sigma^*_{G,j} = v$ for all SNPs, the LDA MR and the TWAS test-statistic run on the marginal eQTL will be proportional to each other by $\hat{\sigma}^{-1}_{LDMR}$. The proof is provided in Appendix B.

## LDA MR-Egger Estimate

The LDA MR-Egger is the same extension of the MR-Egger as the LDA MR was to the MR, by incorporating the LD structure of the SNPs. We include an intercept, in the aims of accounting for the direct effect of the SNPS. The estimate is:

$$\begin{bmatrix} \hat{\alpha}_{LD} \\ \hat{\gamma}_{LDMRE} \end{bmatrix} = \left( \begin{bmatrix} \mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J & \mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \widehat{\boldsymbol{\beta}}_E \\ \widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J & \widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \widehat{\boldsymbol{\beta}}_E \end{bmatrix} \right)^{-1} \begin{bmatrix} \mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \\ \widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \end{bmatrix} \widehat{\boldsymbol{\beta}}_G,$$

and our estimate of γ:

$$\hat{\gamma}_{LDMRE} = \frac{\left(\mathbf{1}_J \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J\right)\left(\hat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_G\right) - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_E\right)\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_G\right)}{\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J\right)\left(\hat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_E\right) - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_E\right)^2}.$$

The variance is estimated as:

$$var(\hat{\gamma}_{LDMRE}) = \hat{\sigma}_{LDMRE}^2 \frac{\mathbf{1}_J \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J}{\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J\right)\left(\hat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_E\right) - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_E\right)^2},$$

$$\hat{\sigma}_{LDMRE}^2 = \frac{1}{J-2}\left(\hat{\boldsymbol{\beta}}_G - \hat{\alpha}_{LD} \mathbf{1}_J - \hat{\boldsymbol{\beta}}_E \hat{\gamma}_{LDMRE}\right)^T \boldsymbol{\Sigma}_G^{-1}\left(\hat{\boldsymbol{\beta}}_G - \hat{\alpha}_{LD} \mathbf{1}_J - \hat{\boldsymbol{\beta}}_E \hat{\gamma}_{LDMRE}\right).$$

The test statistic is then:

$$Z_{LDMRE} = \frac{\hat{\gamma}_{LDMRE}}{\sqrt{var(\hat{\gamma}_{LDMRE})}} = \frac{\left(\mathbf{1}_J \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J\right)\left(\hat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_G\right) - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_E\right)\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_G\right)}{\hat{\sigma}_{LDMRE} \sqrt{\mathbf{1}_J \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J} \sqrt{\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J\right)\left(\hat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_E\right) - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \hat{\boldsymbol{\beta}}_E\right)^2}}.$$

For testing, we compare to a *t*-distribution with *J*-2 degrees of freedom. This estimate is a combination of the approaches from LDA MR and the MR-Egger. The LDA MR-Egger will provide valid inference in the same scenarios as the MR-Egger, but in addition when the SNPs are in LD. Details of all tests are provided in Table I.

**Biases of Estimates**

We first focus on the estimates that do not incorporate an intercept, the MR and the LDA MR. We now examine the bias of the MR estimate:

$$\hat{\gamma}_{MR} = \frac{\hat{\boldsymbol{\beta}}_E^{*T} \boldsymbol{V}^{-1} \hat{\boldsymbol{\beta}}_G^*}{\hat{\boldsymbol{\beta}}_E^{*T} \boldsymbol{V}^{-1} \hat{\boldsymbol{\beta}}_E^*},$$

Note that $\hat{\boldsymbol{\beta}}_G^* = \boldsymbol{\Sigma}\hat{\boldsymbol{\beta}}_G$ and $\hat{\boldsymbol{\beta}}_E^* = \boldsymbol{\Sigma}\hat{\boldsymbol{\beta}}_E$. Thus $E(\hat{\boldsymbol{\beta}}_G^*)$=E$(\boldsymbol{\Sigma}\hat{\boldsymbol{\beta}}_G)$=$\boldsymbol{\Sigma}\boldsymbol{\beta}_E \gamma + \boldsymbol{\Sigma}\boldsymbol{\theta}$. We assume that $\hat{\boldsymbol{\beta}}_E$ and $\hat{\boldsymbol{\beta}}_G$ were estimated from different samples and thus are independent, E$(\hat{\boldsymbol{\beta}}_G | \hat{\boldsymbol{\beta}}_E)$=E$(\hat{\boldsymbol{\beta}}_G)$. Using this gives us that:

$$E(\hat{\gamma}_{MR} | \hat{\boldsymbol{\beta}}_E^*) = \gamma \frac{\hat{\boldsymbol{\beta}}_E^{*T} \boldsymbol{V}^{-1} \boldsymbol{\Sigma} \boldsymbol{\beta}_E}{\hat{\boldsymbol{\beta}}_E^{*T} \boldsymbol{V}^{-1} \hat{\boldsymbol{\beta}}_E^*} + \frac{\hat{\boldsymbol{\beta}}_E^{*T} \boldsymbol{V}^{-1} \boldsymbol{\Sigma} \boldsymbol{\theta}}{\hat{\boldsymbol{\beta}}_E^{*T} \boldsymbol{V}^{-1} \hat{\boldsymbol{\beta}}_E^*}.$$

Then use that $\widehat{\boldsymbol{\beta}}_E^* = \boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E$:

$$E\left(\hat{\gamma}_{MR}\big|\widehat{\boldsymbol{\beta}}_E^*\right) = \gamma\frac{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}\boldsymbol{V}^{-1}\boldsymbol{\Sigma}\boldsymbol{\beta}_E}{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}\boldsymbol{V}^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E} + \frac{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}\boldsymbol{V}^{-1}\boldsymbol{\Sigma}\boldsymbol{\theta}}{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}\boldsymbol{V}^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E}$$

The term $\widehat{\boldsymbol{\beta}}_E$ was estimated from a sample size of $N_E$. As $N_E \to \infty$, we again have that $\widehat{\boldsymbol{\beta}}_E \to \boldsymbol{\beta}_E$ and the first term goes $\gamma$. Assuming $N_E$ is sufficiently large:

$$E\left(\hat{\gamma}_{MR}\big|\widehat{\boldsymbol{\beta}}_E^*\right) \approx \gamma + \frac{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}\boldsymbol{V}^{-1}\boldsymbol{\Sigma}\boldsymbol{\theta}}{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}\boldsymbol{V}^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E}.$$

Now unless $\boldsymbol{\theta} = 0$ or the transformation $\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}\boldsymbol{V}^{-1}\boldsymbol{\Sigma}$ is orthogonal to $\boldsymbol{\theta}$, the MR will be biased. Next assessing the LDA MR:

$$\hat{\gamma}_{LDMR} = \frac{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_G}{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E},$$

$$E\left(\hat{\gamma}_{LDMR}\big|\widehat{\boldsymbol{\beta}}_E\right) = \gamma\frac{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}_G^{-1}\boldsymbol{\beta}_E}{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E} + \frac{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}_G^{-1}\boldsymbol{\theta}}{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E}.$$

As $N_E \to \infty$, $\widehat{\boldsymbol{\beta}}_E \to \boldsymbol{\beta}_E$ and the first term will go to $\gamma$. Assuming that $N_E$ is sufficiently large enough for this to occur, we have:

$$E\left(\hat{\gamma}_{LDMR}\big|\hat{\beta}_E\right) \approx \gamma + \frac{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}_G^{-1}\boldsymbol{\theta}}{\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E}.$$

Again, if $\boldsymbol{\theta} \neq 0$ or if the transformation $\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}_G^{-1}$ is not orthogonal to $\boldsymbol{\theta}$, then the estimate will be biased. As $J$ increases, the $N_E$ needed for $\widehat{\boldsymbol{\beta}}_E \to \boldsymbol{\beta}_E$ will also increase.

We next look at the MR-Egger and again use that $\widehat{\boldsymbol{\beta}}_G^* = \boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_G$ and $\widehat{\boldsymbol{\beta}}_E^* = \boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E$.

$$\hat{\gamma}_{MRE} = \frac{\left(\mathbf{1}_J^T\boldsymbol{V}^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^{T*}\boldsymbol{V}^{-1}\widehat{\boldsymbol{\beta}}_G^*\right) - \left(\mathbf{1}_J^T\boldsymbol{V}^{-1}\widehat{\boldsymbol{\beta}}_E^*\right)\left(\mathbf{1}_J^T\boldsymbol{V}^{-1}\widehat{\boldsymbol{\beta}}_G^*\right)}{\left(\mathbf{1}_J^T\boldsymbol{V}^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^{T*}\boldsymbol{V}^{-1}\widehat{\boldsymbol{\beta}}_E^*\right) - \left(\mathbf{1}_J^T\boldsymbol{V}^{-1}\widehat{\boldsymbol{\beta}}_E^*\right)^2}$$

$$= \frac{\left(\mathbf{1}_J^T\boldsymbol{V}^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}\boldsymbol{V}^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_G\right) - \left(\mathbf{1}_J^T\boldsymbol{V}^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right)\left(\mathbf{1}_J^T\boldsymbol{V}^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_G\right)}{\left(\mathbf{1}_J^T\boldsymbol{V}^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T\boldsymbol{\Sigma}\boldsymbol{V}^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right) - \left(\mathbf{1}_J^T\boldsymbol{V}^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right)^2}.$$

$$E\left(\hat{\gamma}_{MRE}\middle|\widehat{\boldsymbol{\beta}}_E^*\right) = \gamma\frac{\left(\mathbf{1}_J^T V^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma} V^{-1}\boldsymbol{\Sigma}\boldsymbol{\beta}_E\right) - \left(\mathbf{1}_J^T V^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right)\left(\mathbf{1}_J^T V^{-1}\boldsymbol{\Sigma}\boldsymbol{\beta}_E\right)}{\left(\mathbf{1}_J^T V^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma} V^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right) - \left(\mathbf{1}_J^T V^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right)^2} +$$

$$\frac{\left(\mathbf{1}_J^T V^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma} V^{-1}\boldsymbol{\Sigma}\boldsymbol{\theta}\right) - \left(\mathbf{1}_J^T V^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right)\left(\mathbf{1}_J^T V^{-1}\boldsymbol{\Sigma}\boldsymbol{\theta}\right)}{\left(\mathbf{1}_J^T V^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma} V^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right) - \left(\mathbf{1}_J^T V^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right)^2}$$

As we have done previously, we assume $N_E$ is sufficiently large so that $\widehat{\boldsymbol{\beta}}_E \approx \boldsymbol{\beta}_E$, making the first term $\gamma$.

$$E\left(\hat{\gamma}_{MRE}\middle|\widehat{\boldsymbol{\beta}}_E^*\right) \approx \gamma + \frac{\left(\mathbf{1}_J^T V^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma} V^{-1}\boldsymbol{\Sigma}\boldsymbol{\theta}\right) - \left(\mathbf{1}_J^T V^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right)\left(\mathbf{1}_J^T V^{-1}\boldsymbol{\Sigma}\boldsymbol{\theta}\right)}{\left(\mathbf{1}_J^T V^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma} V^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right) - \left(\mathbf{1}_J^T V^{-1}\boldsymbol{\Sigma}\widehat{\boldsymbol{\beta}}_E\right)^2}$$

Thus, the MR-Egger estimate will be unbiased if there is no direct effect ($\boldsymbol{\theta} = 0$), or if $\boldsymbol{\theta}$ is a constant, or if the numerator in the second term is equal to 0. Note that if $\boldsymbol{\Sigma}=\mathbf{I}$ and $V \propto \frac{1}{N}\mathbf{I}$, as is the case in the setting where the MR-Egger estimate was originally proposed, then the numerator is equal to the empirical covariance between $\widehat{\boldsymbol{\beta}}_E$ and $\boldsymbol{\theta}$. This condition is referred to as the Instrument Strength Independent of Direct Effect (InSIDE) condition [Bowden, et al. 2015]. If the number of SNPs in the instrument for the mediator (*J*) is large, and there is no systematic relationship between $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$, then this empirical covariance will be small. However, if *J* is small, as may be the case when gene expression is the mediator and the genetic instrument is limited to cis SNPs, then the numerator in the second term may be non-zero even when there is no systematic relationship between $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$. Even if the INSIDE condition holds but $\boldsymbol{\Sigma} \neq \mathbf{I}$, then the MR-Egger estimate will be biased, as it does not account for linkage disequilibrium. Finally, the inference will be incorrect due to misspecification of the variance of $\widehat{\boldsymbol{\beta}}_G^*$.

We now look at the LDA MR-Egger estimate as:

$$\hat{\gamma}_{LDMRE} = \frac{\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_G\right) - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right)\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_G\right)}{\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right) - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right)^2},$$

Again, using that $\widehat{\boldsymbol{\beta}}_E$ and $\widehat{\boldsymbol{\beta}}_G$ were estimated from different studies, we have that

$$E\left(\hat{\gamma}_{LDMRE}\middle|\widehat{\boldsymbol{\beta}}_E\right) = \gamma\frac{\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\mathbf{1}_J\right)\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1}\boldsymbol{\beta}_E - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right)\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\boldsymbol{\beta}_E}{\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right) - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right)^2}$$

$$+\frac{\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J\right)\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1}\boldsymbol{\theta} - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right)\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\boldsymbol{\theta}}{\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right) - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right)^2}.$$

As with the other estimates: $N_E \to \infty$, $\widehat{\boldsymbol{\beta}}_E \to \boldsymbol{\beta}_E$. If we assume that $N_E$ is sufficiently large:

$$E\left(\hat{\gamma}_{LDMRE}\big|\widehat{\boldsymbol{\beta}}_E\right) \approx \gamma + \frac{\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J\right)\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1}\boldsymbol{\theta} - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right)\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\boldsymbol{\theta}}{\left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1} \mathbf{1}_J\right)\left(\widehat{\boldsymbol{\beta}}_E^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right) - \left(\mathbf{1}_J^T \boldsymbol{\Sigma}_G^{-1}\widehat{\boldsymbol{\beta}}_E\right)^2}.$$

The numerator of the second term is a function of the sample univariate covariance between $\widehat{\boldsymbol{\beta}}_E$ and $\boldsymbol{\theta}$ weighted by $\boldsymbol{\Sigma}_G^{-1}$. The LDA MR-Egger estimate will be unbiased in the same situations as the MR-Egger estimate, but it does not require the SNPs be independent.

**Simulation Study**

We examine the empirical performance of four methods: the MR, MR-Egger, LDA MR, and the LDA MR-Egger. Under our simulation scenario, the TWAS and the LDA MR are approximately the same ($\sigma_{G,j}^*$ is a constant). We generated the data such that all SNPs have a small causal effect. We note, that this procedure is valid for common variants, and likely would not work for rare variants. For each simulation, we generated summary eQTL and GWAS statistics from a multivariate normal distribution as opposed to individual level data [Han, et al. 2009]. We fixed the sample of the eQTL study to 1000 ($N_E$) and the sample size of the GWAS to 5000 (N). We varied the number of SNPs at the loci ($J$), the proportion of variation in $Y$ explained by $G$ ($h_{G\to Y}^2$) and $M$ ($h_{E\to Y}^2$), proportion of variation in $M$ explained by $G$ ($h_e^2$), the variance of the direct effect ($\tau$), and the LD matrix ($\boldsymbol{\Sigma}$, AR($J$) structure, $\Sigma_{i,j} = \rho^{|i-j|}$). This is similar to decreasing LD with distance. For all simulations, the expression effects $\boldsymbol{\beta}_E$ and direct effects $\boldsymbol{\theta}$ were sampled independently. More details along with values taken are given in Table 2. For a variable direct effect, $\tau$ was set to 0. For a constant direct effect, $\tau = 10$. The process and order for the generation of the simulation data is provided in Table 3. We performed 50K simulation for each combination of parameters (486 different combinations). In each simulation, we generated new true $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$, which are functions of $h_e^2, h_{G\to Y}^2, J$ and $\tau$. Generating this new "true" parameter value better represents the different eQTL and GWAS patterns across the genome, and therefore more resembles a standard GWAS or eQTL, as the distribution of eQTL and GWAS parameter values differ across the genome. The procedure detailed in Table 3 is thus

repeated 50K times for all combinations of the parameters in Table 2. Type I error and power were evaluated at 0.05.

As $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$ were drawn independently, the average bias for the LDA MR-Egger estimates across many independent simulations of $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$ should be 0. (This will also be true for the MR-Egger estimates in the absence of LD.) This is analogous to saying the average bias across all the genes in the genome will be zero. However, in practice, we will usually be interested in the test statistic applied to a particular gene. For particular gene with modest $J$, the bias for the LDA MR-Egger need not be 0. To account for this, we next performed a set of 10K simulations with a fixed truth to highlight potential bias. We generate one true $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$ for each value of $J$ that is then held constant for all simulations while we vary the other parameters in Table 2. Therefore, Steps 1 and 5 of Table 3 are not performed in this set of simulations, or more accurately only performed once for J=50 or 300.

**Application to Breast Cancer GWAS Summary Data**

We next compared standard and LDA-MR Egger analyses applied to a breast cancer GWAS. We previously conducted a separate breast cancer TWAS using a different approach to build expression weights and including both validation of predicted expression and functional follow up of significant genes [Wu, et al. 2017]. Here our focus is to compare different analysis approaches where each analysis uses the same set of simple expression weights and the same GWAS summary statistics. These expression weights are different from that of the Wu et al paper [Wu, et al. 2017]. The marginal GWAS summary statistics were from a recent GWAS on breast cancer within women of European descent [Michailidou, et al. 2017]. SNP data was meta-analyzed from 13 platforms: OncoArray [Amos, et al. 2017], iCOGS, and eleven other GWAS. More details on the QC procedure can be found in Michailidou et al [Michailidou, et al. 2017]. After QC, the study consisted of 11.8 million SNPs, with 105,974 controls and 122,977 cases. The GWAS estimates were calculated on the non-standardized minor allele counts, and therefore were transformed using the minor allele frequency as detailed in Appendix A.

Expression weights were calculated from GTEx along with LD information in breast tissue in an overall sample of 183 individuals [Consortium 2013]. We restricted our analysis to the set of transcripts that were deemed heritable using GCTA [Yang, et al. 2011] and examined SNPs within 500kb of the gene boundary. The expression weights were calculated on standardize

minor allele counts of SNPs (mean zero, variance one) and were conditionally estimated using the BSLMM approach [Zhou, et al. 2013]. A gene was deemed heritable if the GCTA p-value for each tissue from GTEx was less than the Bonferroni threshold of 0.05 (after adjusting for 27,945 tests). We were left with 683 potential gene transcripts to analyze.

We analyze the Breast Cancer data for these genes using the TWAS, LDA MR, and LDA MR-Egger. We did not examine the MR and MR-Egger as we were testing for cis-signals as opposed to genome wide and therefore the SNPs were in LD. We took the overlap of the $\widehat{\boldsymbol{\beta}}_G$ SNPs from the GWAS with the available SNP correlations from GTEx. If the effects were estimated with respect to a different reference allele, we reversed the sign for that eQTL effect estimate. In total, the 683 genes corresponded to 191,583 unique SNPs.

## Results

### Simulation Study

First, we examined the type I error rate in simulations and observed that when there is little LD and no direct effect of the SNPs, all of the approaches have the correct type I error (Figure 1, Sup Figure 1). When there is little LD and the direct effect is variable, the MR and MR-Egger approaches have modestly inflated Type I Error rates. When there is low LD and a constant direct effect, only the LDA MR-Egger has correct type I error. If the SNPs are in high LD (bottom row of Figure 1) and there is no direct effect, the MR and the MR-Egger have inflated type I error due to misspecification of the variance. When there is a variable direct effect, all four approaches have inflated type I error. Finally, when there is strong linkage and a constant direct effect, only the LDA MR-Egger has the correct type I error.

We next examined the power when there is little to no LD between the SNPs and no direct effect (Figure 2). Under this situation, all four approaches had correct type I error (Figure 1 top left plot, Sup Figure 1 top left plot) and are valid tests. Regardless of the magnitude of the effect of **M** on **Y**, we see similar power amongst all four methods. There is a slightly smaller power for the LDA MR-Egger compared to the LDA MR or MR, when the SNPs only explain 20% of the variation in the gene expression, but once the SNPs explaining 50% of the variation in *M*, all approaches have approximately equal power. In each individual plot, we see that as the SNPs explain more of the variation (thus better instruments), we have an increase in power.

The decrease in power from J=50 to J=300 is due to an increase in the signal to noise ratio when predicting expression levels using SNPs. For J=50, a larger proportion of the variation explained by SNPs is shared by each individual SNP leading to more precise estimates. When J=300, a smaller proportion of that same amount of variation is explained by each SNP in a larger set of SNPs. Assuming the sample size in the reference panel used to estimate $\widehat{\boldsymbol{\beta}}_E$ is the same, the sampling error in the SNP-specific estimates $\hat{\beta}_{E,j}$ is the same for J=50 and J=300. We have held the proportion of variance explained by the SNPs constant, while increasing the noise due to sampling error (as we are using estimated expression effects from 300 SNPs rather than 50).

When there is low LD and a constant direct effect of the SNPs on the outcome (Table 4), then the MR-Egger has less power than the LDA MR-Egger to detect an association when J=50. This may be due to the LDA MR-Egger gaining some information by accounting for the weak LD structure. If J=300, there is a decrease in power compared to when J=50 regardless of the presence of a direct effect in these two methods. At J=300, the LDA MR-Egger has slightly higher power than the MR-Egger. When the SNPs explain 50% of the variation in $\boldsymbol{M}$, both methods have power greater than 80% regardless of the effect of $\boldsymbol{M}$ on $\boldsymbol{Y}$ (Table 4). We did not report the power of the MR or LDA MR, as they did not have proper type I error when there is a constant direct effect (Figure 1 and Supplement Figure 1).

Finally, we examine the power when there is strong LD amongst the SNPs. We do not assess the MR or MR-Egger for this scenario as they do not have correct type I error for correlated SNPs. When J=50, and there is no direct effect, the LDA MR has more power than the LDA MR-Egger, though the difference in power is less pronounced when $h_E^2$=0.5 and $\boldsymbol{M}$ has a strong effect on $\boldsymbol{Y}$. When J=300, the two methods have comparable power (Table 4), with the LDA MR-Egger having slightly less power than the LDA MR. When there is a direct effect, the LDA MR-Egger has approximately equal power to when no direct effect.

Comparing the power when there is strong LD vs small LD, and J=50, when there is a small mediated effect and $h_E^2 = 0.5$, we have less power compared to when there is little LD (LDA MR-Egger 0.914 to 0.784). When J=300, there is a smaller drop-in power for small LD to strong LD, with the LDA MR-Egger power going from 0.869 to 0.840. The LDA MR has equal power regardless of whether the SNPs are strongly correlated vs weakly correlated.

We next examine the bias of our estimates (Figure 3). We here show the results when there is strong LD (Figure 3). The results for low are in LD Supplement Figure 2. In Figure 3, the first two rows show the results when J=50, and the next two rows when J=300. The first column shows when there is no direct effect, the second column when there is a variable direct effect, and the third column when there is a constant direct effect. The first and third row are when there is no effect of the mediator on the outcome ($\gamma = 0$) and the second and fourth row show when there is a large effect ($\gamma \neq 0$). Here we have fixed $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$ effects for all simulations.

As expected based on the section on bias of estimates sections, when $\gamma = 0$ and $\boldsymbol{\theta} = \boldsymbol{0}$, all of the estimates are unbiased. When there is an effect on the outcome and no direct effect, we see that the non LDA aware approaches converge faster to the truth than the LD aware. While the non LDA converge faster, recall that they misspecify the variance and lead to improper inference. We also see attenuation bias when $\gamma \neq 0$, with estimates improving as the SNPs become better instruments. The attenuation bias is larger for J=300 relative to J=50, for the same reason that we saw a decrease in power from J=50 to J=300: an increase in signal to noise. When there is a variable direct effect, all of the approaches are biased (second column). Our empirical bias result depends on the particular values of $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$. As described in the methods ("Bias"), all of the estimates for the mediated effect of gene expression include a weighted covariance between $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$. Even if the average across all genes for this covariance is zero, for any particular gene it is non-zero, and can be large. If the $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$ for each gene can be thought of as independent draws from their respective distributions, then the average absolute magnitude of the bias term is a decreasing function of $J$. This is why we see smaller bias when $J$=300. Finally, when there is a constant direct effect, only the LDA MR-Egger is unbiased when $\gamma = 0$ or J=50. When J=300 and $\gamma \neq 0$, we see the attenuation bias in the LDA-MR Egger, with it biased downward from the null.

**Application to Breast Cancer GWAS Summary Data**

Of the 683 genes tested, 74 were called significant (p< $7.32*10^{-5}$) by at least one approach (TWAS, LDA MR or LDA MR-Egger, Supplementary Tables). Comparing the TWAS vs the LDA MR-Egger (Figure 4A), there were 18 genes that were significant by the TWAS but not by LDA MR-Egger, 22 genes that were by LDA MR-Egger and not TWAS, and 14 that were called by both. With the LDA MR and the LDA MR-Egger, there was much more agreement due

to the same weight matrix being used, but still the LDA MR called 23 genes as significant that the LDA-MR Egger did not (Figure 4B). Twenty-seven gene transcripts were called significant by both the LDA-MR and the LDA-MR Egger (Supplement Tables 1 and 3). There were 12 genes called significant by all three approaches. A detailed list of which gene transcripts were found significant by which method is provided in Supplementary Tables 1 through 7. Examining the spearman correlation between the p-values, LDA MR and LDA MR –Egger had an $r^2$ of 0.49, LDA MR and TWAS of .27, and LDA MR-Egger and TWAS of 0.208. The kappa statistic for calling a gene transcript significant between LDA MR and LDA MR-Egger was 0.6, between LDA MR and TWAS was 0.33, and between LDA MR-Egger and TWAS was 0.38. We did not see much agreement between the TWAS and LDA MR as the TWAS was calculated using the conditional expression weights, and the variances of all SNPs were not the same at every loci (comparing test statistics Supplement Figure 3).

A gene that we will highlight that was called significant by the TWAS was SET Domain Containing 9 (*SETD9*), with a p-value of 2.47e-23 at cytoband 5q11.2. The p-value for this gene transcript for the LDA MR and the LDA MR-Egger was 0.012 and 0.278 respectively. This is potentially a TWAS false-positive as the LDA MR-Egger was not associated with gene. A fine mapping analysis of this locus found four functional candidate SNPs in a sample of approximately 100K women of European descent [Glubb, et al. 2015]. These four candidate functional SNPs were associated with an increase in activity of *MAP3K1* (Mitogen-Activated Protein Kinase 1), another gene at this locus. *SETD9* was ruled as not the gene of interest as it had no association with these four candidate SNPs. *MAP3K1* is located 94 kbp from *SETD9*. (*MAP3K1* did not pass our cis-heritable tissue threshold and was not included in the analysis.)

Another association of interest was at the 2q33.1 locus, where Caspase 8 (*CASP8*) had a p-value using the TWAS method of 4.63e-07, but was not significant after accounting for the number of genes tested by the LDA MR-Egger (p-value 0.01). These results are consistent with published fine mapping analyses: after mutual conditioning, SNPs near *ALS2CR12* (1 kbp from *CASP8*) were more significantly associated with breast cancer risk than those near *CASP8* [Lin, et al. 2015]. (*ALS2CR12* did not pass our cis-heritable threshold in GTEx data and was not included in our analyses.) Both of these examples highlight situations where the significant

TWAS and LDA MR results may be due to co-localization, while LDA MR-Egger was not strongly affected by this confounding.

## Discussion

In this work, we propose a new LD aware MR-Egger estimate of the effect of gene expression on an outcome using just summary statistics. Our method properly accounts for both LD and when there is a constant direct effect. The LDA MR and the summary MR and MR-Egger are not proper estimates under these situations as they either do not account for the direct effect (LDA MR, MR) or the correlation between SNPs (MR and MR-Egger). For scenarios where the LDA MR is a valid test, the LDA MR-Egger has comparable or slightly lower power. However, the researchers can gain confidence that they are correctly adjusting for pleiotropy when the pleiotropic effect is constant. We also provided a 1-1 relationship between the TWAS and the traditional LDA-MR when the standard errors of the marginal GWAS effects are all equal. In our real data application, the three approaches (TWAS, LDA MR and LDA MR-Egger) all picked up different results, with the LDA MR-Egger potentially correctly calling some results as null that the TWAS did not.

As we were finalizing this paper, we were made aware of a recent publication done independently of this group by Burgess and Thompson [Burgess and Thompson 2017]. In the discussion of their work, they describe the extension of the MR-Egger to the case of correlated instruments and give a proof in their appendix of its convergence. They did not provide a real data application or examine the empirical performance. Here, we provide simulations detailing the empirical performance of this approach in comparison to the existing methods. We motivated the LDA MR-Egger in its relationship to the TWAS and how it can help with potential problems in the wide scale use of these summary data approaches. We finally gave a real data example of the application of this method to a large scale Breast Cancer GWAS.

While we focused on the case of gene expression data, the LDA MR-Egger can easily be extended to an arbitrary mediator when the instruments are correlated. The approaches we consider are a function of three things: the correlation of the instruments, the strength of the instruments, and the association between the instruments and the outcome not taking into account the intermediate variable. It is thus easily extendable and easily executable.

The advantages of the LDA MR-Egger estimator over the LDA MR and MR-Egger approaches are that it relaxes the strong assumption regarding no pleiotropy and accounts for linkage disequilibrium among SNP instruments. Just as with MR-Egger regression, LDA MR-Egger regression does not provide robustness against general forms of pleiotropy: it is unbiased when the direct SNP effects are identical or when the empirical weighted covariance between the instrument strengths ($\boldsymbol{\beta}_E$) and direct effects ($\boldsymbol{\theta}$) is 0. When the number of SNPs in the instrument is small, this correlation may be non-zero even if there is no systematic, mechanistic relationship between $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$.

Our simulations assumed an "infinitesimal" genetic architecture for both the gene expression and outcome phenotypes. Departures from this model—for example, if gene expression was causally influenced by only one or a small proportion of nearby SNPs—could affect the performance of the proposed tests. Future work will consider the impact of local genetic architecture on these tests. We also focused on the case where there may be pleiotropic direct effects on the outcome, but these are not systematically related to the SNPs effects on the mediator. This is reflected in our simulations when we draw $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$ independently. If there is a systematic relationship between $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$, then none of the methods we have discussed here will provide unbiased estimates of $\gamma$ or valid tests of $\gamma = 0$. There could be a systematic relationship between $\boldsymbol{\beta}_E$ and $\boldsymbol{\theta}$ even when $\gamma = 0$ if, for example, the SNPs influence a third (unobserved) trait, which in turn influences both the mediator and outcome. This is arguably unlikely to be true when the mediator is expression of a gene and the genetic instrument consists of cis SNPs, except in the case of co-localization. Evaluating the causal effect of the mediator on outcome when the instrument and direct effects are systematically related is not possible without additional information on the mediator-outcome relationship: large samples with data on genetic factors, the mediator, outcome and possible confounders will likely be needed.

In practice, we caution interpretation of the $\gamma$ terms due to QC and pre-processing of the data make it difficult to infer meaningful biological interpretation of $\gamma$. We are also faced with the issue of attenuation bias, as we have estimates of the true eQTL effects. Despite this, the LDA MR-Egger is still a valid test for the effect of the gene on outcome in most scenarios and can inform the direction of the effect.

In summary, we have extended the use of LDA MR and MR-Egger into LDA MR-Egger, a useful tool for when the genetic instruments are correlated and disease and eQTL SNPs have colocalized. This method of incorporating summary statistics from different sources can help in discovering novel loci as well as narrowing in on the susceptible gene in the region. We provided equations for their bias as well as evaluated their performance empirically. Further work can be done to account for the variation in the eQTL estimates and in parsing out the direct effects.

## Acknowledgments

## References

Amos CI, Dennis J, Wang Z, Byun J, Schumacher FR, Gayther SA, Casey G, Hunter DJ, Sellers TA, Gruber SB and others. 2017. The OncoArray Consortium: A Network for Understanding the Genetic Architecture of Common Cancers. Cancer Epidemiol Biomarkers Prev 26(1):126-135.

Barbeira A, Shah KP, Torres JM, Wheeler HE, Torstenson ES, Edwards T, Garcia T, Bell GI, Nicolae D, Cox NJ and others. 2016. MetaXcan: Summary Statistics Based Gene-Level Association Method Infers Accurate PrediXcan Results. bioRxiv.

Bowden J, Davey Smith G, Burgess S. 2015. Mendelian randomization with invalid instruments: effect estimation and bias detection through Egger regression. Int J Epidemiol 44(2):512-25.

Burgess S, Dudbridge F, Thompson SG. 2016. Combining information on multiple instrumental variables in Mendelian randomization: comparison of allele score and summarized data methods. Stat Med 35(11):1880-906.

Burgess S, Thompson SG. 2017. Interpreting findings from Mendelian randomization using the MR-Egger method. Eur J Epidemiol 32(5):377-389.

Consortium GT. 2013. The Genotype-Tissue Expression (GTEx) project. Nat Genet 45(6):580-5.

Gamazon ER, Wheeler HE, Shah KP, Mozaffari SV, Aquino-Michaels K, Carroll RJ, Eyler AE, Denny JC, Consortium GT, Nicolae DL and others. 2015. A gene-based association method for mapping traits using reference transcriptome data. Nat Genet 47(9):1091-8.

Giambartolomei C, Vukcevic D, Schadt EE, Franke L, Hingorani AD, Wallace C, Plagnol V. 2014. Bayesian test for colocalisation between pairs of genetic association studies using summary statistics. PLoS Genet 10(5):e1004383.

Glubb DM, Maranian MJ, Michailidou K, Pooley KA, Meyer KB, Kar S, Carlebur S, O'Reilly M, Betts JA, Hillman KM and others. 2015. Fine-scale mapping of the 5q11.2 breast cancer locus reveals at least three independent risk variants regulating MAP3K1. Am J Hum Genet 96(1):5-20.

Gusev A, Ko A, Shi H, Bhatia G, Chung W, Penninx BW, Jansen R, de Geus EJ, Boomsma DI, Wright FA and others. 2016. Integrative approaches for large-scale transcriptome-wide association studies. Nat Genet 48(3):245-52.

Han B, Kang HM, Eskin E. 2009. Rapid and accurate multiple testing correction and power estimation for millions of correlated markers. PLoS Genet 5(4):e1000456.

Johnson T. 2011. Efficient calculation for multi-SNP genetic risk scores. Technical Report, Queen Mary University of London.

Lin WY, Camp NJ, Ghoussaini M, Beesley J, Michailidou K, Hopper JL, Apicella C, Southey MC, Stone J, Schmidt MK and others. 2015. Identification and characterization of novel associations in the CASP8/ALS2CR12 region on chromosome 2 with breast cancer risk. Hum Mol Genet 24(1):285-98.

Mancuso N, Shi H, Goddard P, Kichaev G, Gusev A, Pasaniuc B. 2017. Integrating Gene Expression with Summary Association Statistics to Identify Genes Associated with 30 Complex Traits. Am J Hum Genet 100(3):473-487.

Michailidou K, Lindstrom S, Dennis J, Beesley J, Hui S, Kar S, Lemacon A, Soucy P, Glubb D, Rostamianfar A and others. 2017. Association analysis identifies 65 new breast cancer risk loci. Nature.

O'Connor LJ, Gusev A, Liu X, Loh P-R, Finucane HK, Price AL. 2017. Estimating the proportion of disease heritability mediated by gene expression levels. bioRxiv.

Pasaniuc B, Price AL. 2017. Dissecting the genetics of complex traits using summary association statistics. Nat Rev Genet 18(2):117-127.

Shi H, Kichaev G, Pasaniuc B. 2016. Contrasting the Genetic Architecture of 30 Complex Traits from Summary Association Data. Am J Hum Genet 99(1):139-53.

Valeri L, VanderWeele TJ. 2013. Mediation analysis allowing for exposure-mediator interactions and causal interpretation: theoretical assumptions and implementation with SAS and SPSS macros. Psychol Methods 18(2):137-50.

Wu L, Shi W, Long J, others a. 2017. Identification of novel susceptibility loci and genes for breast cancer risk: A large transcriptome-wide association study in nearly 230,000 women of European descent. Submitted.

Yang J, Lee SH, Goddard ME, Visscher PM. 2011. GCTA: a tool for genome-wide complex trait analysis. Am J Hum Genet 88(1):76-82.

Zhou X, Carbonetto P, Stephens M. 2013. Polygenic modeling with bayesian sparse linear mixed models. PLoS Genet 9(2):e1003264.

Zhu Z, Zhang F, Hu H, Bakshi A, Robinson MR, Powell JE, Montgomery GW, Goddard ME, Wray NR, Visscher PM and others. 2016. Integration of summary data from GWAS and eQTL studies predicts complex trait gene targets. Nat Genet 48(5):481-7.

# Figure Legends

**Figure 1: Type I error when J=50.** Each bar represents results over $5 \times 10^4$ simulations. Evaluated at $\alpha = 0:05$. First panel represent when low LD (plots with A). Second panel represents when strong LD (plots with B). From left to right correspond to: no direct effect, variable direct effect, and a constant direct effect.
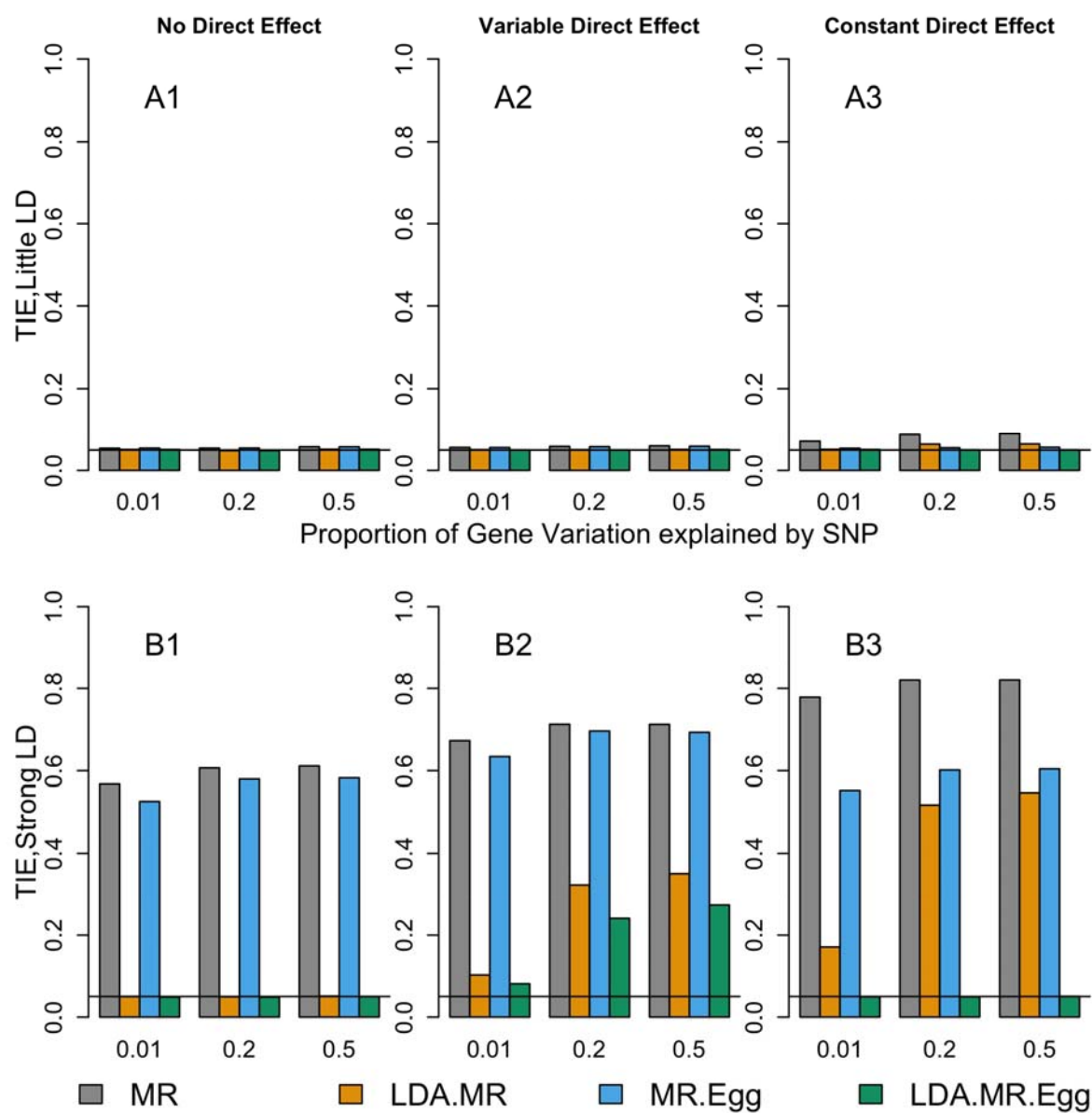
**Figure 2: Power when little to no LD and J=50, 300.** Power results when there is little to no LD and no direct effect. Each bar represents results over $5 \times 10^4$ simulations. Evaluated at $\alpha = 0{:}05$. First row represents J=50 and second row when J=300. From left to right: when $\gamma^2 = 0.005$ and $\gamma^2 = 0.01$.

**Figure 3: Bias when strong LD for J=50, 300.** Bias plots for when there is strong LD in the SNP set. First row corresponds to J=50, $= 0$ (plots with A). Second panel (plots with B) when J $= 50$ and $\gamma^2 = 0.01$. Third panel (plots with C) when J $= 300$ and $= 0$. Final panel (plots with D) J $= 300$ and $\gamma^2 = 0.01$. From left to right: no direct effect, variable direct effect, constant direct effect.

**Figure 4: Comparing –log 10 p-values.** Shows for 683 genes between LDA MR Egger and TWAS (A), LDA MR Egger and LDA MR (B), and LDA MR and TWAS (C). Red line is the Bonferroni cutoff of -$\log_{10}(.05/683)$.

## Table

**Table 1: Commonly used approaches for testing using summary statistics.** Details four common Summary Mendelian Randomization approaches and the TWAS for testing for an association between gene expression and outcome through GWAS (*J* is the number of SNPs in the loci).

| Method | Adjusts for Direct Effect | Accounts for LD | Distribution |
|---|---|---|---|
| TWAS | No | Yes | $Z_{TWAS} \sim N(0,1)$ |
| MR | No | No | $Z_{MR} \sim t(df = J - 1)$ |
| MR-Egger | Yes | No | $Z_{MRE} \sim t(df = J - 2)$ |
| LDA MR | No | Yes | $Z_{LDMR} \sim t(df = J - 1)$ |
| LDA MR-Egger | Yes | Yes | $Z_{LDMRE} \sim t(df = J - 2)$ |

**Table 2: Simulation parameters that were modified.** Parameters that were varied for all simulations performed. A description of the parameter and potential values could be taken are given. All combinations of parameters were examined.

| Parameter | Definition | Values taken |
|---|---|---|
| $h_e^2$ | Proportion of variability in E explained by G | 0.01, 0.20, 0.50 |
| $h_{E \to Y}^2$ | Proportion of Variability in Y explained by E | 0, 0.005, 0.01 |
| $h_{G \to Y}^2$ | Proportion of Variability in Y explained by G | 0, 0.005, 0.01 |
| $N$ | Sample Size of GWAS | 5000 |
| $N_E$ | Sample Size of EQTL | 1000 |
| $\rho$ | Correlation between SNPs. AR structure | 0.125, 0.9 |
| $J$ | Number of SNPs in the Loci | 50,300 |
| $\tau$ | Strength of Pleiotropic Effects | 0,10 |
| **L** | Cholesky Decomposition of $\Sigma$ | Function of $\rho$ |

**Table 3: Procedure for generating simulation study.** Gives the steps in order for generating the type I error and power simulations detailed in the paper and how the parameters were generated.

| Step | Procedure | Mathematically |
|------|-----------|----------------|
| 1 | Generate true eQTL | $\boldsymbol{\beta}_E \sim N_J(0, \boldsymbol{I}_J)$ |
| 2 | Set the proportion of variability in expression due to SNPs | $\sigma_E = \sqrt{\dfrac{h_e^2}{\boldsymbol{\beta}_E^T \boldsymbol{\Sigma} \boldsymbol{\beta}_E}}$ |
| 3 | Rescale the eQTL effects | $\boldsymbol{\beta}_E = \sigma_E \boldsymbol{\beta}_E$ |
| 4 | Set Expression to outcome Effect | $\gamma = \sqrt{h_{E \to Y}^2}$ |
| 5 | Generate potential Direct Effect | $\boldsymbol{\theta} = \boldsymbol{e} + \tau; \; \boldsymbol{e} \sim N_J(0, \boldsymbol{I}_J)$ |
| 6 | Set the proportion of variability in outcome due directly to SNPs | $\sigma = \sqrt{\dfrac{h_{G \to Y}^2}{\boldsymbol{\theta}^T \boldsymbol{\Sigma} \boldsymbol{\theta}}}$ |
| 7 | Rescale the Direct Effects | $\boldsymbol{\theta} = \sigma * \boldsymbol{\theta}$ |
| 8 | Generate the GWAS effects | $\boldsymbol{\beta}_G = \boldsymbol{\theta} + \gamma \boldsymbol{\beta}_E$ |
| 9 | Generate Observed eQTL Values | $\widehat{\boldsymbol{\beta}}_E^* \sim \boldsymbol{\Sigma} \boldsymbol{\beta}_E + \boldsymbol{L}^T \boldsymbol{\epsilon}_E; \; \boldsymbol{\epsilon}_E \sim N_J\left(0, \dfrac{1 - h_E^2}{N_E} \boldsymbol{I}_J\right)$ |
| 10 | Generate Observed GWAS values | $\widehat{\boldsymbol{\beta}}_G^* \sim \boldsymbol{\Sigma} \boldsymbol{\beta}_G$ $+ \boldsymbol{L}^T \boldsymbol{\epsilon}_G; \; \boldsymbol{\epsilon}_G \sim N_J\left(0, \boldsymbol{I}_J \dfrac{1 - h_E^2 * h_{E \to Y}^2 - h_{G \to Y}^2}{N}\right)$ |

**Table 4: Power results from simulation study.** Details the power results at an $\alpha$ level of 0.05 under varying levels of LD, number of SNPs in the loci, and presence of a direct effect. TWAS not shown as is equivalent to the LDA MR under are simulation procedure.

| LD | Number of SNPS | % Variation E explained by SNPs Strength of Mediated Effect | Method | No Direct effect | | Constant Direct Effect | |
|----|----|----|----|----|----|----|----|
| | | | | Small Effect | Large Effect | Small Effect | Large Effect |
| | | | MR | 0.058 | 0.063 | ---- | ---- |
| | | | LDA MR | 0.054 | 0.058 | ---- | ---- |
| | | 0.01 | MR-Egger | 0.057 | 0.062 | 0.055 | 0.057 |
| | | | LDA MR - Egger | 0.054 | 0.059 | 0.055 | 0.059 |
| | | | MR | 0.524 | 0.804 | ---- | ---- |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | J=50 | 0.2 | LDA MR | 0.510 | 0.797 | ---- | ---- |
| | | | MR-Egger | 0.514 | 0.795 | 0.302 | 0.508 |
| | | | LDA MR - Egger | 0.499 | 0.787 | 0.498 | 0.788 |
| | | 0.5 | MR | 0.922 | 0.997 | ---- | ---- |
| | | | LDA MR | 0.921 | 0.997 | ---- | ---- |
| | | | MR-Egger | 0.915 | 0.997 | 0.658 | 0.911 |
| Small LD | | | LDA MR - Egger | 0.914 | 0.997 | 0.915 | 0.997 |
| | | 0.01 | MR | 0.054 | 0.057 | ---- | ---- |
| | | | LDA MR | 0.051 | 0.053 | ---- | ---- |
| | | | MR-Egger | 0.054 | 0.057 | 0.052 | 0.051 |
| | | | LDA MR - Egger | 0.050 | 0.052 | 0.052 | 0.051 |
| | J=300 | 0.2 | MR | 0.339 | 0.576 | ---- | ---- |
| | | | LDA MR | 0.321 | 0.560 | ---- | ---- |
| | | | MR-Egger | 0.338 | 0.573 | 0.288 | 0.501 |
| | | | LDA MR - Egger | 0.321 | 0.557 | 0.325 | 0.562 |

| | | | | | | |
|---|---|---|---|---|---|---|
| | | MR | 0.877 | 0.992 | ---- | ---- |
| | | LDA MR | 0.870 | 0.992 | ---- | ---- |
| | 0.5 | MR-Egger | 0.875 | 0.992 | 0.812 | 0.980 |
| | | LDA MR - Egger | 0.869 | 0.991 | 0.869 | 0.992 |
| | 0.01 | LDA MR | 0.054 | 0.059 | ---- | ---- |
| | | LDA MR - Egger | 0.052 | 0.055 | 0.053 | 0.055 |
| J=50 | 0.2 | LDA MR | 0.513 | 0.798 | ---- | ---- |
| | | LDA MR - Egger | 0.389 | 0.629 | 0.389 | 0.637 |
| | 0.5 | LDA MR | 0.919 | 0.998 | ---- | ---- |
| Large LD | | LDA MR - Egger | 0.784 | 0.943 | 0.788 | 0.941 |
| | 0.01 | LDA MR | 0.049 | 0.049 | ---- | ---- |
| | | LDA MR - Egger | 0.050 | 0.049 | 0.051 | 0.051 |
| J=300 | 0.2 | LDA MR | 0.318 | 0.560 | ---- | ---- |
| | | LDA MR - | 0.294 | 0.523 | 0.300 | 0.524 |

|      |             | Egger |       |       |       |
| ---- | ----------- | ----- | ----- | ----- | ----- |
| 0.5  | LDA MR      | 0.869 | 0.992 | ----  | ----  |
|      | LDA MR - Egger | 0.840 | 0.985 | 0.842 | 0.986 |