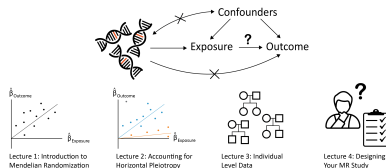# SSGG Short Course: A Introduction to Mendelian Randomization
## Lecture 1: Introduction to Mendelian Randomization[1]

Ting Ye

Department of Biostatistics, University of Washington
tingye1@uw.edu

February 12, 2024

Causation    Instrumental variable    What is MR?    IVW, dIVW, MR-raps    Diagnosis    Basic workflow and software    Discussion and Summary    References

●○○      ○○○      ○○○○      ○○○○      ○○      ○○      ○○○

# Hierarchy of evidence

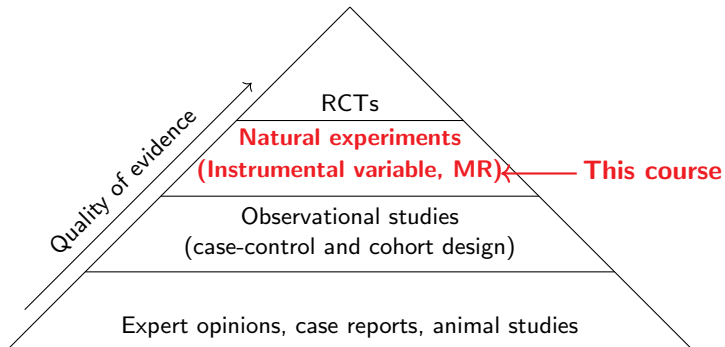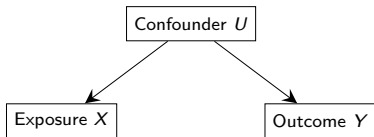When the goal is to infer causation...



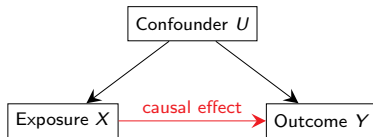Figure: (A rough) Hierarchy of evidence in medical studies.

# Fundamental challenges in observational studies

"Correlation does not imply causation"

▶ Correlation/association describes the **statistical relationship** in the data, indicating **difference in one variable is associated with difference in another**.

▶ Causation requires **mechanistic understanding**, indicating **intervention in one variable leads to change in another**.



(a) Correlated but not causal                               (b) Causal

# Fundamental challenges in observational studies

One idea: adjusting for possible sources of spurious correlation.

▶ Example: Possible confounders between low density lipoprotein cholesterol (LDL-C) and coronary heart disease (CHD): age, sex, BMI, ...

▶ **Fundamental challenge: We can never be sure this list is complete.**

▶ The promise of instrumental variables: estimating causal effect without enumerating confounders.

# What is an instrumental variable (IV)?



## Core IV assumptions

1. **Relevance**: $Z$ is associated with the exposure ($X$).
2. **Independence**: $Z$ is independent of unmeasured confounder ($U$).
3. **Exclusion restriction**: $Z$ cannot have any direct effect on the outcome ($Y$).

## Examples of IVs

▶ Encouragement, physician/hospital preference, distance to care provider, calendar time, genetic variants... (Baiocchi et al., 2014)

▶ In an RCT, patients are randomized to take Statin for lowering LDL-C.

Causation
000

**Instrumental variable**
00●

What is MR?
0000

IVW, dIVW, MR-raps
0000

Diagnosis
00

Basic workflow and software
00

Discussion and Summary
000

References

# The Wald ratio

### How it works?

- ▶ Suppose 1 unit ↑ in $Z$ ⇒ $\gamma$ unit ↑ in $X$
- ▶ Suppose 1 unit ↑ in $X$ ⇒ $\beta_0$ unit ↑ in $Y$ (Causal effect)
- ▶ Then, 1 unit ↑ in $Z$ ⇒ ?? unit ↑ in $Y$

# The Wald ratio

### How it works?

- Suppose 1 unit $\uparrow$ in $Z \Rightarrow \gamma$ unit $\uparrow$ in $X$
- Suppose 1 unit $\uparrow$ in $X \Rightarrow \beta_0$ unit $\uparrow$ in $Y$ (Causal effect)
- Then, 1 unit $\uparrow$ in $Z \Rightarrow$ ?? unit $\uparrow$ in $Y$

$$\textbf{Wald ratio:} \text{ Causal effect of } X \text{ on } Y = \frac{\text{Effect of } Z \text{ on } Y}{\text{Effect of } Z \text{ on } X}$$

# What is MR?

▶ MR uses **genetic variants (SNPs) as IVs** to infer causation

1. Relevance: many traits are influenced by genetics
2. Independence: SNPs are randomly inherited from parents (Mendel's laws of inheritance)
3. Exclusion restriction: SNPs do not have a direct effect on the outcome **(no horizontal pleiotropy)** → **Lecture 2**



**Figure 3** Mendelian randomization in parent–offspring design

Offspring should have an equal chance of receiving either of the alleles that the parents have at any particular locus

(Davey Smith and Ebrahim, 2003)

# MR in non-familial studies

*Of course populations share much common ancestry and the genetic make-up of individuals can be traced back through the random segregation of alleles during a sequence of matings, but associating genetic markers with disease risk or phenotype within such populations is not as well protected against potential distorting factors as are parent–offspring comparisons. Thus the Mendelian randomization in genetic association studies is approximate, rather than absolute.*

*(Davey Smith and Ebrahim; 2003)*

▶ Within-family MR → **Lecture 3**

Causation    Instrumental variable    What is MR?    IVW, dIVW, MR-raps    Diagnosis    Basic workflow and software    Discussion and Summary    References

000      000      0000      0000      00      00      000

# Example: MR of LDL-C on CHD

▶ **Key idea**: people who inherited certain alleles tends to have higher LDL



▶ MR emulates a hypothetical RCT

## Two-sample summary-data MR

Combine **publically available** summary data on **gene-exposure** and **gene-outcome** associations from **two separate samples**.

---

Zhao et al. (2020)

Causation    Instrumental variable    **What is MR?**    IVW, dIVW, MR-raps    Diagnosis    Basic workflow and software    Discussion and Summary    References

000     000     000●     0000     00     00     000

## Two-sample summary-data MR

Combine **publically available** summary data on **gene-exposure** and **gene-outcome** associations from **two separate samples**.

**Example**: estimate the effect of LDL-C on CHD using $p = 160$ independent SNPs.

1. **Exposure dataset**: A GWAS for LDL-C, $\text{lm}(X \sim Z_j) \Rightarrow \hat{\gamma}_j, \sigma_{Xj}, j = 1, \dots, p$.
2. **Outcome dataset**: A GWAS for CHD, $\text{lm}(Y \sim Z_j) \Rightarrow \hat{\Gamma}_j, \sigma_{Yj}, j = 1, \dots, p$.

**These two datasets are independent.**

## MR Assumptions

$\hat{\gamma}_j \sim N(\gamma_j, \sigma_{Xj}^2), \hat{\Gamma}_j \sim N(\Gamma_j, \sigma_{Yj}^2), j = 1, \dots, p$, are all independent, and $\Gamma_j / \gamma_j = \beta_0$ for all $j$.

- ▶ Reasonable when all SNPs are independent (from LD clumping), and no overlapping sample between exposure and outcome datasets.
- ▶ For continuous outcome: $\beta_0$ is average causal effect from one unit $\uparrow$ in exposure
- ▶ For binary outcome: $\beta_0$ is a conservative causal odds ratio from one unit $\uparrow$ in exposure

Zhao et al. (2020)

Causation    Instrumental variable    What is MR?    **IVW, dIVW, MR-raps**    Diagnosis    Basic workflow and software    Discussion and Summary    References

000         000         0000         ●000         00         00         000

# Inverse-variance weighted estimator (IVW)

Wald estimator: $\hat{\beta}_j = \hat{\Gamma}_j / \hat{\gamma}_j$, with $\mathrm{var}(\hat{\beta}_j) \approx \sigma_{Yj}^2 / \hat{\gamma}_j^2$.

$$\hat{\beta}_{\mathrm{IVW}} = \frac{\sum_{j=1}^{p} \hat{\gamma}_j^2 \sigma_{Yj}^{-2} \hat{\beta}_j}{\sum_{j=1}^{p} \hat{\gamma}_j^2 \sigma_{Yj}^{-2}}$$

---

Burgess et al. (2013)

# Inverse-variance weighted estimator (IVW)

Wald estimator: $\hat{\beta}_j = \hat{\Gamma}_j / \hat{\gamma}_j$, with $\mathrm{var}(\hat{\beta}_j) \approx \sigma_{Yj}^2 / \hat{\gamma}_j^2$.

$$\hat{\beta}_{\mathrm{IVW}} = \frac{\sum_{j=1}^{p} \hat{\gamma}_j^2 \sigma_{Yj}^{-2} \hat{\beta}_j}{\sum_{j=1}^{p} \hat{\gamma}_j^2 \sigma_{Yj}^{-2}}$$

▶ Alternatively, fit $\mathrm{lm}(\hat{\Gamma}_j \sim 0 + \hat{\gamma}_j, \text{weights}= \sigma_{Yj}^{-2})$

---

Burgess et al. (2013)

# Inverse-variance weighted estimator (IVW)

Wald estimator: $\hat{\beta}_j = \hat{\Gamma}_j / \hat{\gamma}_j$, with $\mathrm{var}(\hat{\beta}_j) \approx \sigma_{Yj}^2 / \hat{\gamma}_j^2$.

$$\hat{\beta}_{\mathrm{IVW}} = \frac{\sum_{j=1}^{p} \hat{\gamma}_j^2 \sigma_{Yj}^{-2} \hat{\beta}_j}{\sum_{j=1}^{p} \hat{\gamma}_j^2 \sigma_{Yj}^{-2}}$$

▶ Alternatively, fit $\mathrm{lm}(\hat{\Gamma}_j \sim 0 + \hat{\gamma}_j, \text{weights} = \sigma_{Yj}^{-2})$
▶ Well-known that IVW is **biased** with weak IVs.

---

Burgess et al. (2013)

Causation    Instrumental variable    What is MR?    **IVW, dIVW, MR-raps**    Diagnosis    Basic workflow and software    Discussion and Summary    References

○○○     ○○○     ○○○○     ●○○○     ○○     ○○     ○○○

# Inverse-variance weighted estimator (IVW)

Wald estimator: $\hat{\beta}_j = \hat{\Gamma}_j / \hat{\gamma}_j$, with $\mathrm{var}(\hat{\beta}_j) \approx \sigma_{Yj}^2 / \hat{\gamma}_j^2$.

$$\hat{\beta}_{\mathrm{IVW}} = \frac{\sum_{j=1}^p \hat{\gamma}_j^2 \sigma_{Yj}^{-2} \hat{\beta}_j}{\sum_{j=1}^p \hat{\gamma}_j^2 \sigma_{Yj}^{-2}}$$

▶ Alternatively, fit $\mathrm{lm}(\hat{\Gamma}_j \sim 0 + \hat{\gamma}_j, \text{weights} = \sigma_{Yj}^{-2})$

▶ Well-known that IVW is **biased** with weak IVs.

▶ To mitigate the bias, usually pre-screen for strong IVs (e.g. with genome-wide significance p-value $5 \times 10^{-8}$), but using weak IVs may increase power

---

Burgess et al. (2013)

# Inverse-variance weighted estimator (IVW)

Wald estimator: $\hat{\beta}_j = \hat{\Gamma}_j / \hat{\gamma}_j$, with $\mathrm{var}(\hat{\beta}_j) \approx \sigma_{Yj}^2 / \hat{\gamma}_j^2$.

$$\hat{\beta}_{\mathrm{IVW}} = \frac{\sum_{j=1}^p \hat{\gamma}_j^2 \sigma_{Yj}^{-2} \hat{\beta}_j}{\sum_{j=1}^p \hat{\gamma}_j^2 \sigma_{Yj}^{-2}}$$

▶ Alternatively, fit $\mathrm{lm}(\hat{\Gamma}_j \sim 0 + \hat{\gamma}_j,\ \mathrm{weights}= \sigma_{Yj}^{-2})$
▶ Well-known that IVW is **biased** with weak IVs.
▶ To mitigate the bias, usually pre-screen for strong IVs (e.g. with genome-wide significance p-value $5 \times 10^{-8}$), but using weak IVs may increase power
▶ **Debiased IVW** (Ye et al., 2021) uses a simple modification to effectively de-bias:

$$\hat{\beta}_{\mathrm{dIVW}} = \frac{\sum_{j=1}^p \hat{\Gamma}_j \hat{\gamma}_j \sigma_{Yj}^{-2}}{\sum_{j=1}^p (\hat{\gamma}_j^2 - \sigma_{Xj}^2) \sigma_{Yj}^{-2}}.$$

Standard error can be computed with a simple formula.

---

Burgess et al. (2013)

# Profile likelihood method (MR-raps)

▶ Profile log-likelihood:

$$\ell(\beta) = -\frac{1}{2} \sum_{j=1}^{p} \frac{(\hat{\Gamma}_j - \beta\hat{\gamma}_j)^2}{\sigma_{Yj}^2 + \beta^2 \sigma_{Xj}^2}$$

▶ $\hat{\beta}_{\text{raps}} = \arg\max_\beta \ell(\beta)$

---

Zhao et al. (2020)

# IV strength

▶ The average F-statistic is commonly used to measure IV strength:

$$\text{F-stat} = \frac{1}{p} \sum_{j=1}^{p} \frac{\hat{\gamma}_j^2}{\sigma_{Xj}^2} - 1$$

▶ IVW requires F-stat $> 10$, while dIVW requires F-stat $\cdot \sqrt{p} > 20$.

Causation
000

Instrumental variable
000

What is MR?
0000

IVW, dIVW, MR-raps
000●

Diagnosis
00

Basic workflow and software
00

Discussion and Summary
000

References

# Simulations: IVW, dIVW, MR-raps

To closely mirror real applications, we take the real BMI-CAD dataset (available in the *mr.divw*) package as our simulation parameters. We use $p = 1119$ independent SNPs (pre-selected).

1. **Exposure dataset**: A GWAS for BMI in the UK BioBank ($n = 336, 107$);
   $\Rightarrow \{\gamma_j, \sigma_{Xj}^2, j \in [p]\}$
2. **Outcome dataset**: A GWAS for CAD by the CARDIoGRAMplusC4D consortium
   ($n = 185, 000$). $\Rightarrow \{\sigma_{Yj}^2, j \in [p]\}$

We set $\beta_0 = 0.4, \Gamma_j = \beta_0 \gamma_j$ (i.e., no pleiotropy). We have F-stat $= 7.8$ and F-stat $\cdot \sqrt{p} = 260.2$.

| Method | mean | SD | SE | CP | |
|--------|------|-----|-----|------|---|
| **IVW** | 0.352 | 0.047 | 0.047 | 82.6 | ← **Biased and poor CP** |
| **dIVW** | 0.400 | 0.054 | 0.054 | 94.7 | ← **Unbiased and adequate CP** |
| **MR-raps** | 0.400 | 0.054 | 0.054 | 94.9 | |

Causation     Instrumental variable     What is MR?     IVW, dIVW, MR-raps     **Diagnosis**     Basic workflow and software     Discussion and Summary     References

000          000         0000         0000         ●0         00         000
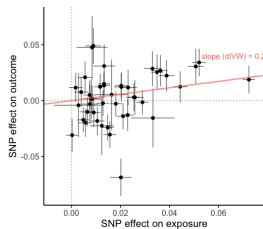
## Diagnosis: tests for MR assumptions

▶ **F-test for weak IVs**: estimator-specific (e.g., F-stat $> 10$ for IVW and F-stat $> 20/\sqrt{p}$ for dIVW)

▶ **Modified Cochran's Q test for heterogeneity** (Bowden et al., 2019):

$$\text{Q statistic}: \quad Q = \sum_{j=1}^{p} \frac{(\hat{\Gamma}_j - \hat{\beta}\hat{\gamma}_j)^2}{\sigma_{Yj}^2 + \hat{\beta}^2 \sigma_{Xj}^2}$$
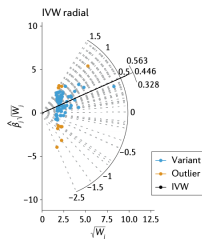
If $Q > \chi_{1-\alpha,p-1}^2$, then reject $H_0$ : same $\Gamma_j/\gamma_j$ across $j$.

▶ **MR-PRESSO global test, outlier test, and distortion test**:
   - Global test: $\text{RSS}_{\text{obs}} = \sum_{j=1}^{p}(\hat{\Gamma}_j - \hat{\beta}_{-j}\hat{\gamma}_j)^2$ compared against a simulated distribution under no heterogeneity
   - Outlier test: $\text{RSS}_{\text{obs},j} = (\hat{\Gamma}_j - \hat{\beta}_{-j}\hat{\gamma}_j)^2$ compared against a simulated distribution under no heterogeneity with Bonferroni correction
   - Distortion test: $D = 100 \times (\hat{\beta}_{\text{all}} - \hat{\beta}_{\text{sub}})/|\hat{\beta}_{\text{sub}}|$

▶ **Steiger filtering**: assess whether IVs primarily affect exposure or outcome (Hemani et al., 2017) (details $\rightarrow$ **Lecture 4**)
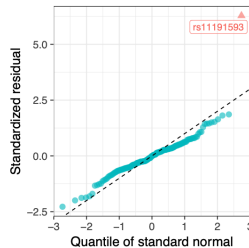
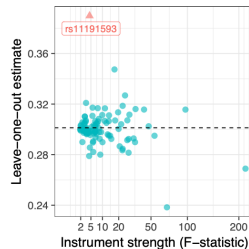# Diagnosis: visualization tools for MR assumptions



(a) Scatter plot

(b) Radial plot

(c) QQ plot

(d) Leave-one-out

Bowden et al. (2018); Zhao et al. (2020); Sanderson et al. (2022)

Causation    Instrumental variable    What is MR?    IVW, dIVW, MR-raps    Diagnosis    **Basic workflow and software**    Discussion and Summary    References

000       000        0000        0000        00       ●0        000

## Basic workflow and software: BMI on CHD

1. Installation and load package: `library(TwoSampleMR)`
2. Select IVs for the exposure:
   `exposure_dat <- extract_instruments("ieu-a-2", p1 = 5e-08,`
   `        clump = TRUE, r2 = 0.001, kb = 10000)`
3. Extract IVs for the outcome
   `outcome_dat <- extract_outcome_data(exposure_dat$SNP, "ieu-a-7")`
4. Harmonize the effect sizes
   `dat <- harmonise_data(exposure_dat, outcome_dat)`

|  | Exposure GWAS | | | | Outcome GWAS | | | |
|---|---|---|---|---|---|---|---|---|
| SNP | Effect | Effect allele | Other allele | Effect allele frequency | Effect | Effect allele | Other allele | Effect allele frequency |
| rs123456 | -0.485 | G | T | 0.41 | 0.056 | T | G | 0.61 |

Harmonize

|  | Exposure GWAS | | | | Outcome GWAS | | | |
|---|---|---|---|---|---|---|---|---|
| SNP | Effect | Effect allele | Other allele | Effect allele frequency | Effect | Effect allele | Other allele | Effect allele frequency |
| rs123456 | -0.485 | G | T | 0.41 | -0.056 | G | T | 0.39 |

5. MR analysis and diagnosis

More complete workflow → **Lecture 4**

## Practice in R ($\sim$20min)

Causation    Instrumental variable    What is MR?    IVW, dIVW, MR-raps    Diagnosis    Basic workflow and software    **Discussion and Summary**    References

000        000        0000        0000        00        00        ●○○

# Strengths and challenges of MR

**Strengths:**

▶ Less susceptible to conventional unmeasured confounding
- Mendel's laws of inheritance
▶ Less susceptible to reverse causation
- Genetics are fixed at conception
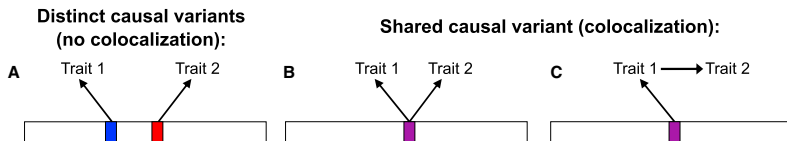▶ Has a summary-data and a two-sample option

**Challenges:**

▶ Weak IV bias
▶ Genetic-outcome confounding
▶ Widespread horizontal pleiotropy can cause bias
- Each variant has multiple biological functions
▶ Low power
▶ Assumes constant treatment effect
▶ Based on gene-environment equivalence
▶ Only applicable to heritable exposures

---

Song et al. (2023)

# Connections to related methods

▶ **Colocalization** (Wallace, 2021; Zuber et al., 2022):
   - whether two traits share common genetic causal variant(s) in a gene region



   - can be used in combination with MR (Case A → violation of MR assumptions; Case B → strengthened MR conclusion)

▶ **Polygenic risk scores (PRS)**:
   - predicted risk for disease based on genome-wide genotypes
   - common to use PRS as an IV in MR analysis

▶ **Transcriptome-wide association study (TWAS)** (Gusev et al., 2016; Gamazon et al., 2015):
   - study expression-trait associations through genetic imputed expression level
   - essentially an MR of gene expression on trait (Zhu and Zhou, 2021; Zhao et al., 2024)

# Summary

▶ MR leverages genetic variants as instruments to address causal questions
  - Emulates an RCT
  - Triangulation across multiple sources of evidence for causal inference
▶ MR assumptions
▶ Methods when all IVs are valid: IVW, dIVW, MR-raps
▶ Diagnosis: F-test for weak IVs, Q test and MR-PRESSO for heterogeneity, visualizations
▶ Basic workflow using the `TwoSampleMR` package
▶ Discussion: strengths and challenges, connection to other methods

Baiocchi, M., Cheng, J., and Small, D. S. (2014). Instrumental variable methods for causal inference. *Statistics in medicine*, 33(13):2297–2340.

Bowden, J., Del Greco M, F., Minelli, C., Zhao, Q., Lawlor, D. A., Sheehan, N. A., Thompson, J., and Davey Smith, G. (2019). Improving the accuracy of two-sample summary-data mendelian randomization: moving beyond the nome assumption. *International journal of epidemiology*, 48(3):728–742.

Bowden, J., Spiller, W., Del Greco M, F., Sheehan, N., Thompson, J., Minelli, C., and Davey Smith, G. (2018). Improving the visualization, interpretation and analysis of two-sample summary data mendelian randomization via the radial plot and radial regression. *International journal of epidemiology*, 47(4):1264–1278.

Burgess, S., Butterworth, A., and Thompson, S. G. (2013). Mendelian randomization analysis with multiple genetic variants using summarized data. *Genetic Epidemiology*, 37(7):658–665.

Davey Smith, G. and Ebrahim, S. (2003). 'Mendelian randomization': can genetic epidemiology contribute to understanding environmental determinants of disease? *International Journal of Epidemiology*, 32(1):1–22.

Gamazon, E. R., Wheeler, H. E., Shah, K. P., Mozaffari, S. V., Aquino-Michaels, K., Carroll, R. J., Eyler, A. E., Denny, J. C., Consortium, G., Nicolae, D. L., et al. (2015). A gene-based association method for mapping traits using reference transcriptome data. *Nature genetics*, 47(9):1091–1098.

Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B. W., Jansen, R., De Geus, E. J., Boomsma, D. I., Wright, F. A., et al. (2016). Integrative approaches for large-scale transcriptome-wide association studies. *Nature genetics*, 48(3):245–252.

Hemani, G., Tilling, K., and Davey Smith, G. (2017). Orienting the causal relationship between imprecisely measured traits using gwas summary data. *PLoS genetics*, 13(11):e1007081.

Sanderson, E., Glymour, M. M., Holmes, M. V., Kang, H., Morrison, J., Munafò, M. R., Palmer, T., Schooling, C. M., Wallace, C., Zhao, Q., et al. (2022). Mendelian randomization. *Nature Reviews Methods Primers*, 2(1):6.

Song, Y., Ye, T., Roberts, L. R., Larson, N. B., and Winham, S. J. (2023). Mendelian randomization in hepatology: A review of principles, opportunities, and challenges. *Hepatology*, pages 10–1097.

Wallace, C. (2021). A more accurate method for colocalisation analysis allowing for multiple causal variants. *PLoS genetics*, 17(9):e1009440.

Ye, T., Shao, J., and Kang, H. (2021). Debiased inverse-variance weighted estimator in two-sample summary-data Mendelian randomization. *The Annals of Statistics*, 49(4):2079–2100.

Zhao, Q., Wang, J., Hemani, G., Bowden, J., and Small, D. S. (2020). Statistical inference in two-sample summary-data Mendelian randomization using robust adjusted profile score. *Annals of Statistics*, 48(3):1742–1769.

Zhao, S., Crouse, W., Qian, S., Luo, K., Stephens, M., and He, X. (2024). Adjusting for genetic confounders in transcriptome-wide association studies improves discovery of risk genes of complex traits. *Nature Genetics*, pages 1–12.

Zhu, H. and Zhou, X. (2021). Transcriptome-wide association studies: A view from mendelian randomization. *Quantitative Biology*, 9(2):107–121.

Zuber, V., Grinberg, N. F., Gill, D., Manipur, I., Slob, E. A., Patel, A., Wallace, C., and Burgess, S. (2022). Combining evidence from mendelian randomization and colocalization: Review and comparison of approaches. *The American Journal of Human Genetics*.