# Identification of cell-type-specific genetic regulation of gene expression for transcriptome-wide association studies

Duo Zhang*[1], Qiurui Ma*[2], Brandon Jew[3], Sriram Sankararaman[4]

[1]Shandong University; [2]Hong Kong University of Science and Technology; [3] Bioinformatics Interdepartmental Program, UCLA; [4]Department of Computer Science, Department of Human Genetics, UCLA

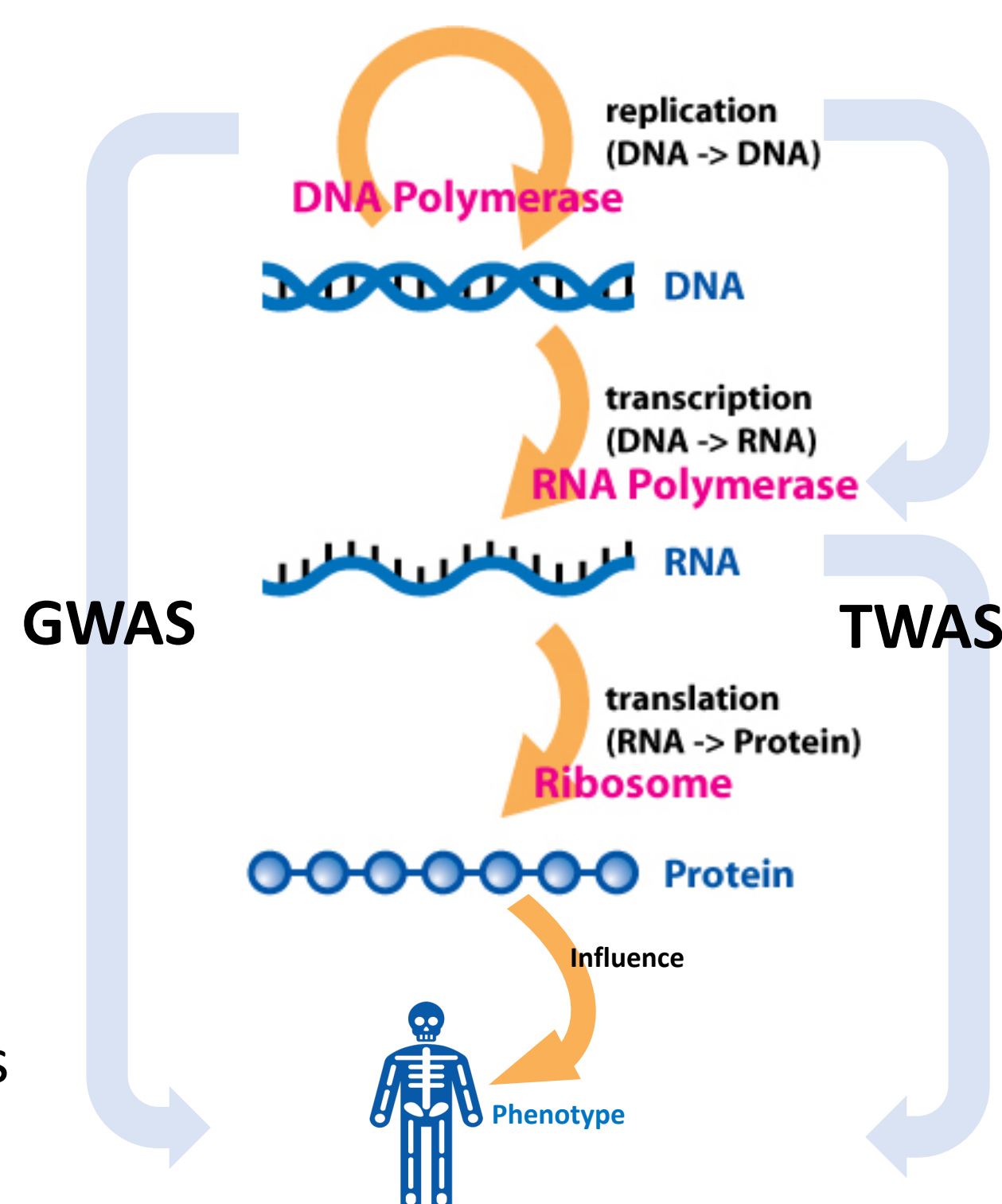* Equal contribution

## I. BACKGROUNDS

### i. The Central Dogma

- Single nucleotide polymorphisms (**SNPs**) are sites of variation in our DNA
- Gene expression (**GE (Z)**) is the level of mRNA in one cell type. **Bulk level GE (G)** is the combined GE of all cell types in a tissue.

**Encode DNA into SNP data**

**The Central Dogma**



### ii. Current Studies

- Genome-wide association studies (**GWAS**)[2] linearly associate SNPs with phenotypes
- Transcriptome-wide association studies (TWAS)[3] linearly characterize the association of **GE** regulated by SNPs and phenotypes

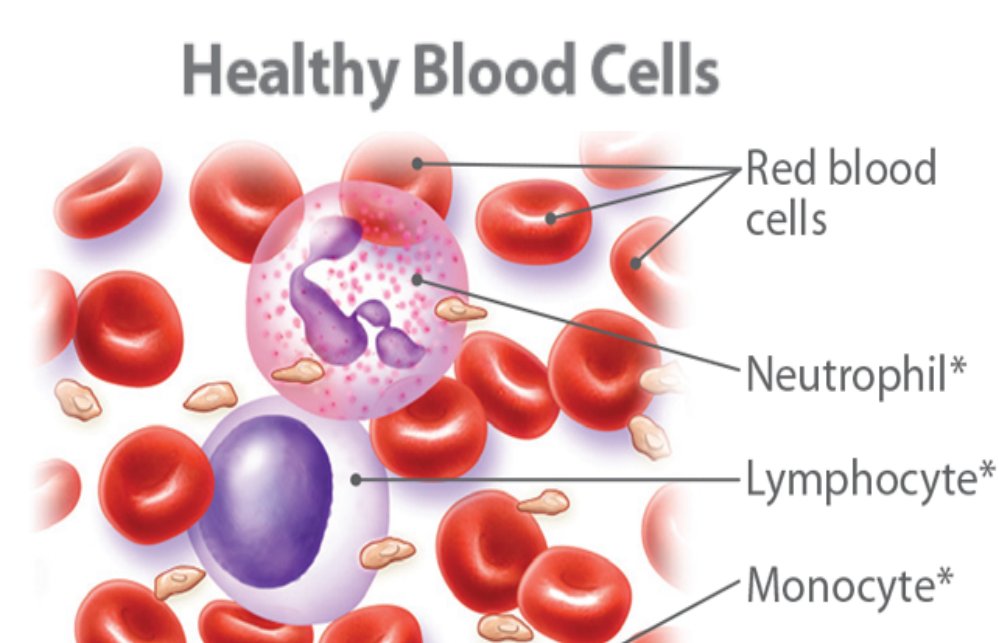### iii. Challenges & Goals

**Methodological**
Unclear how SNPs affect phenotypes
- Missing cell type information
- Identified associations do not indicate causality

**Our Goal**
Deconvolute bulk level GE into cell-specific GE with SNPs and cell-type weights. Associate cell-type specific GE with phenotypes
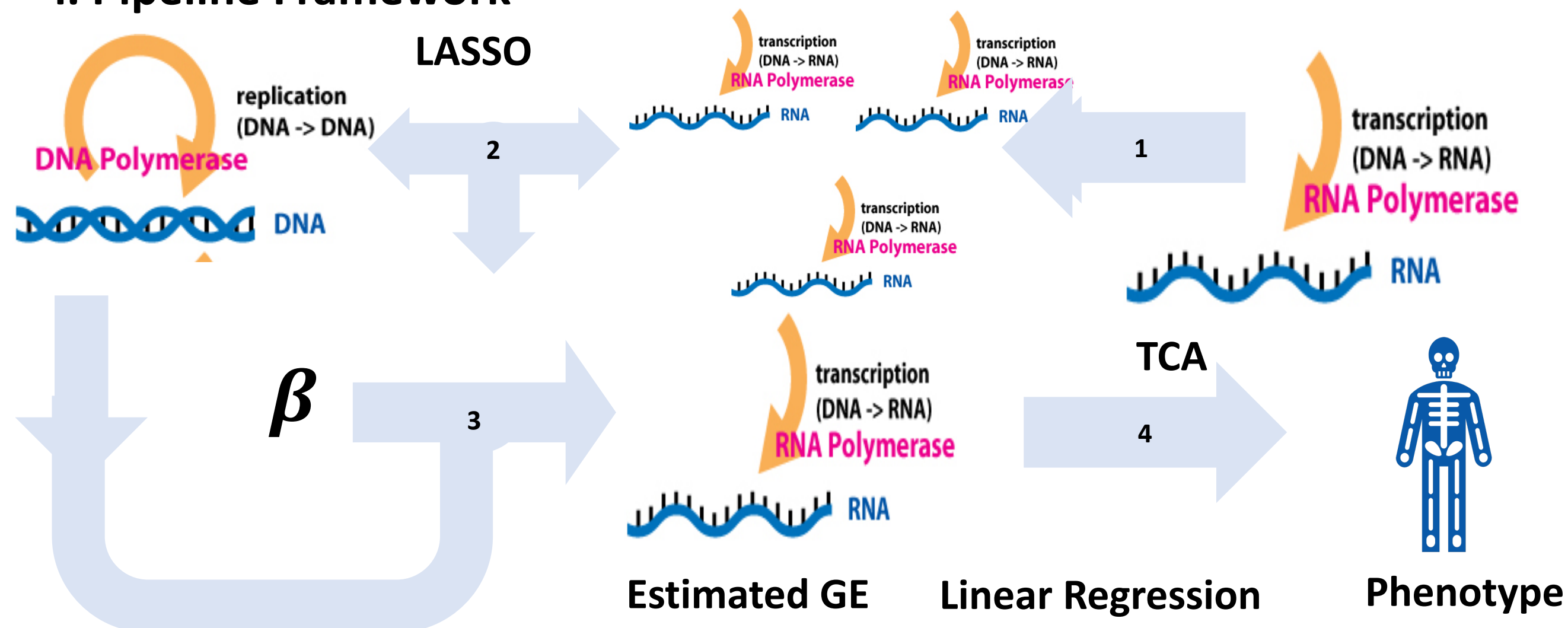
**Practical**
Cell-type-specific biological data is resource intensive and expensive to acquire.

**Healthy Blood Cells**



## II. METHODOLOGY

### i. Pipeline Framework



1. TCA deconvolutes bulk level GE into cell-type-specific ones
2. Effect size of SNPs on cell-type-specific GE imputed by LASSO
3. Cell-type-specific gene expression imputed from effect size for external cohorts
4. Estimated cell-type-specific GE is regressed into phenotype

### ii. TCA model

**Added SNPs effect: mdl1**

**TCA Model: mdl2**

$$z_h^i = \epsilon_z + \mu_h + \begin{bmatrix} 1 \\ 27 \\ 0 \end{bmatrix}^T \gamma_h + \begin{bmatrix} 0 \\ 2 \\ 1 \\ 0 \\ 1 \end{bmatrix}^T \beta_h \qquad \epsilon_z \in N(0, \sigma_z)$$

**Cell-Specific GE**   **Covariate:** (Gender, Age, Smoking)   **SNPs**

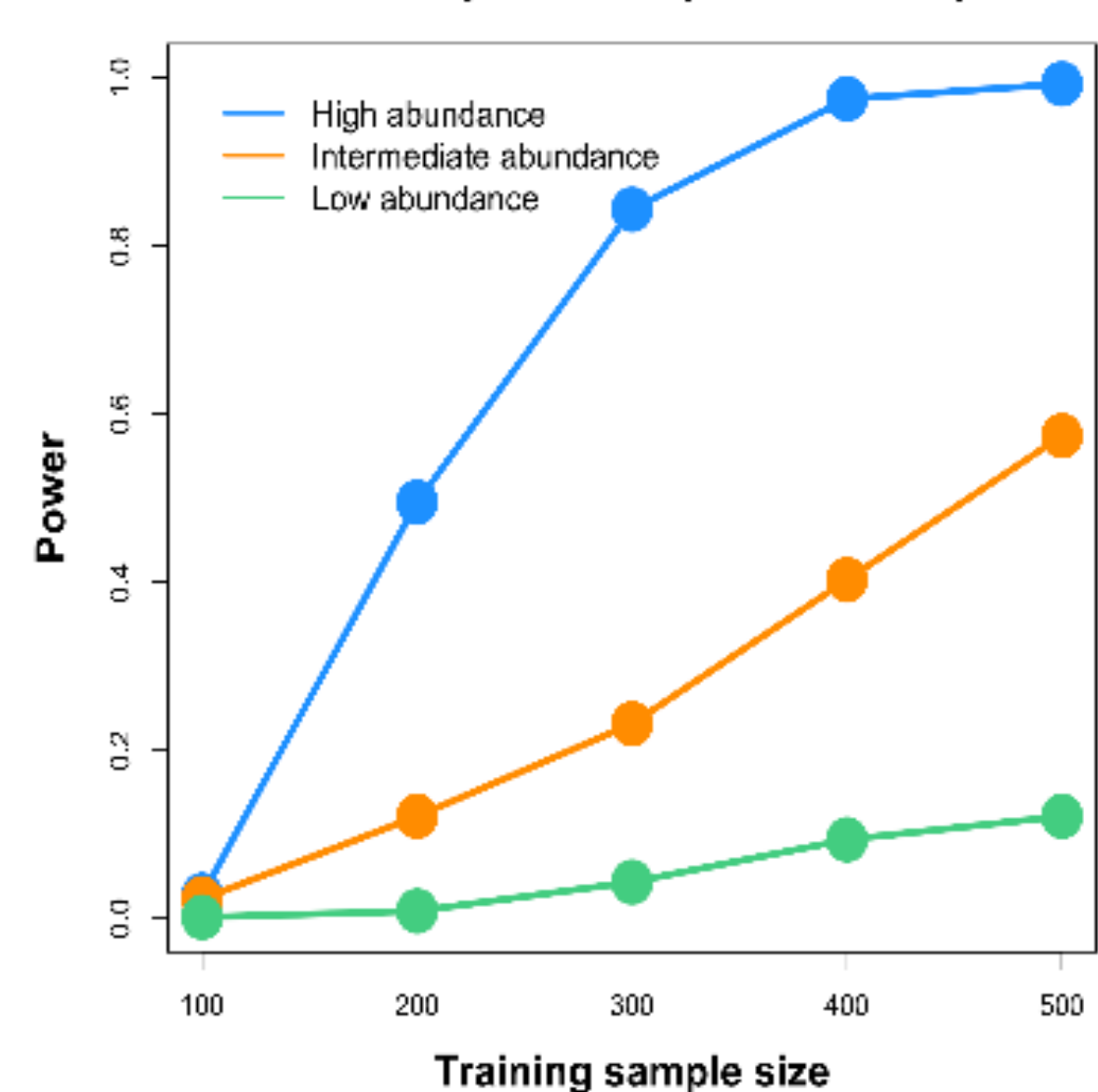$$G_i = c_i^2 \delta + \sum_{h=1}^{k} w_{hi} z_{hi} + \epsilon_g \qquad \epsilon_g \in N(0, \sigma_g)$$

TCA assumes bulk level **GE** is a linear combination of GEs
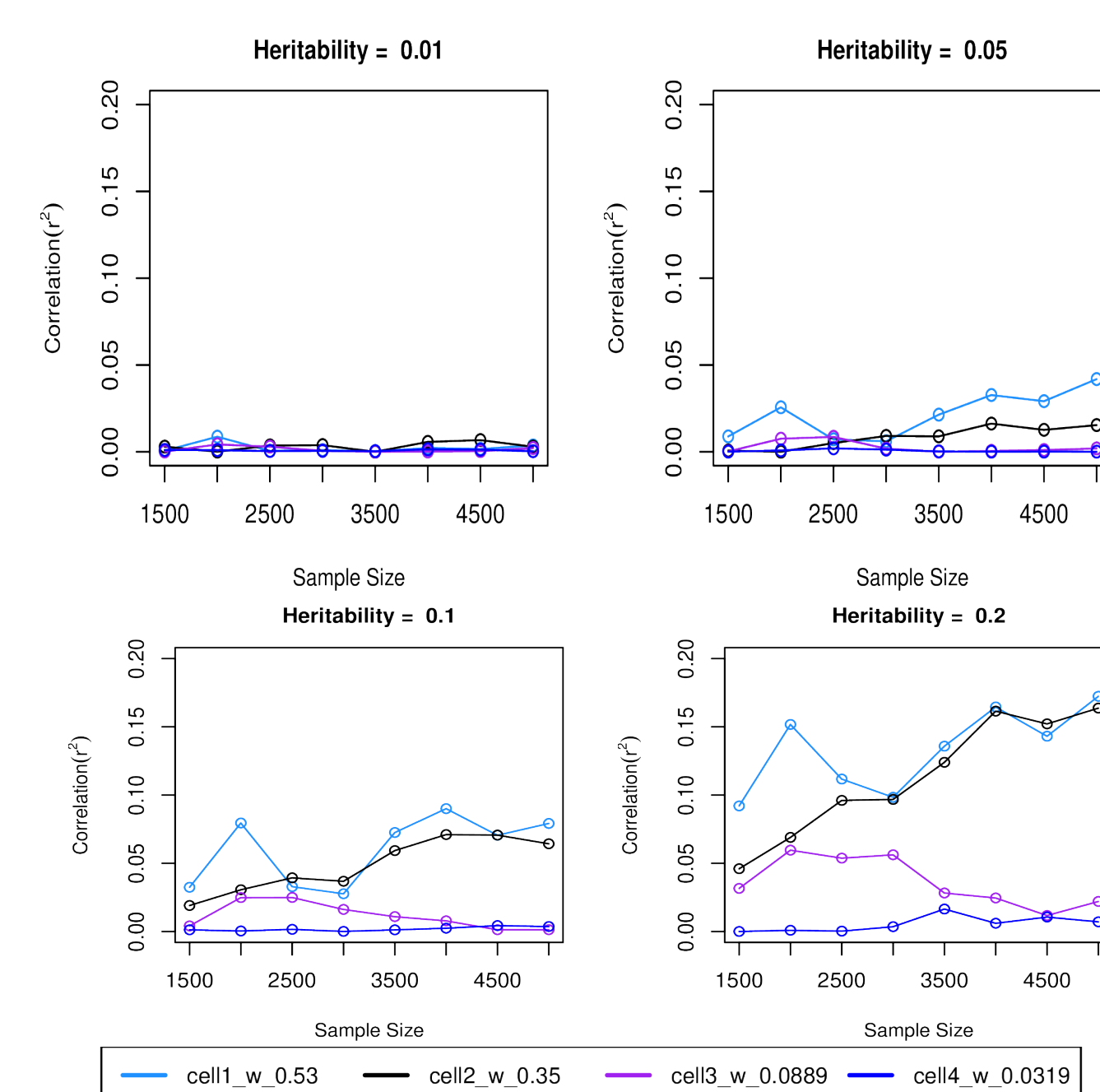
## III. RESULTS

### i. Simulated data

**The method possesses sufficient power to detect cell-specific expression-phenotype associations**



**The variance explained by the model is lower than the theoretical upper bound**
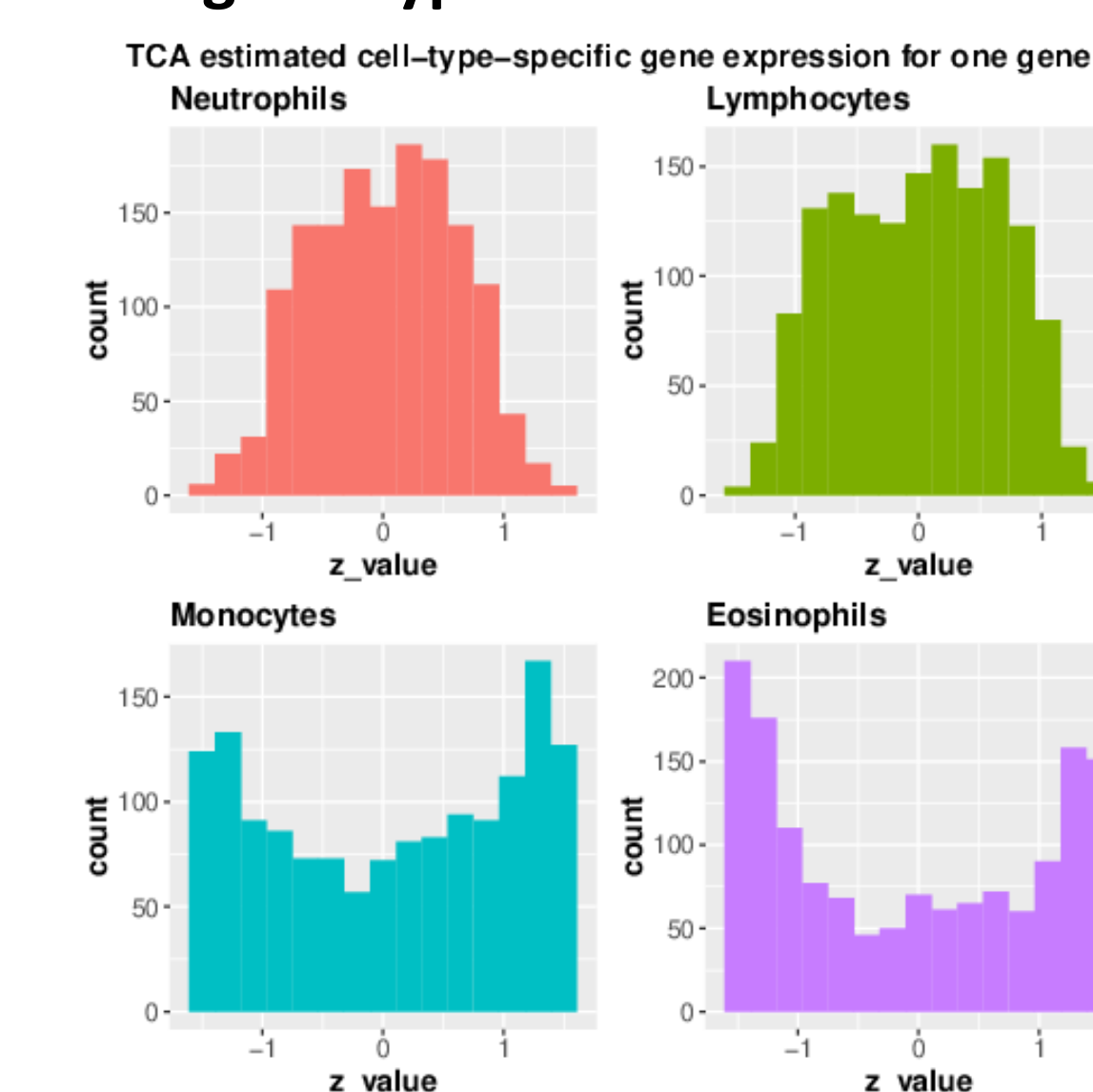


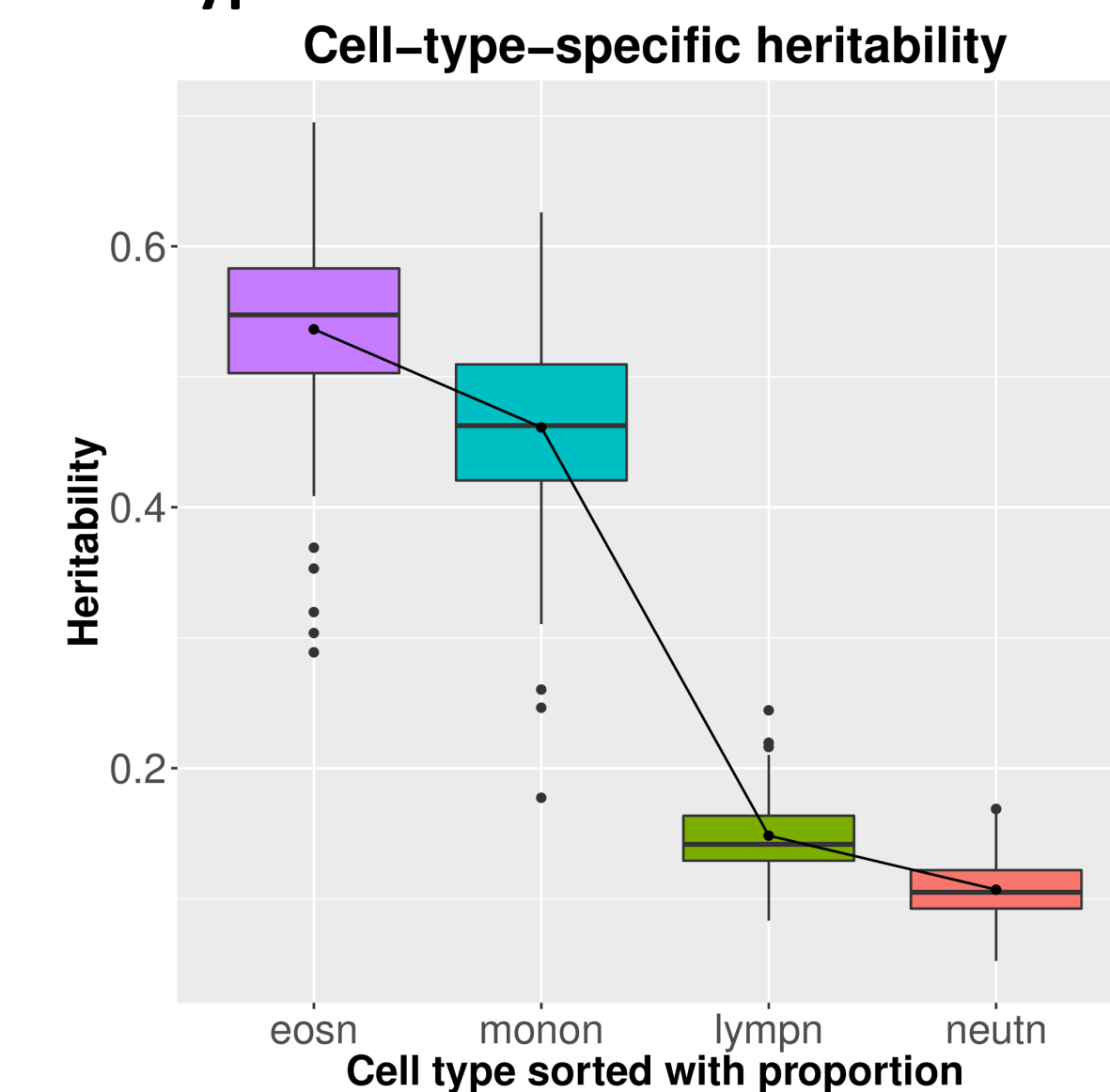**Our modified TCA performs better than the original TCA**



### ii. Real data

**The model's performance is inconsistent among cell types**



**The model overfits on less abundance cell types**



## IV. CONCLUSION

Cell-specific expression-phenotype associations in large datasets (UK BioBank)) could be learnt with its SNPs and readily available, abundant datasets with bulk level gene expressions.

**References**

[1] Rahmani, E., Schweiger, R., Rhead, B., Criswell, L. A., Barcellos, L. F., Eskin, E., ... & Halperin, E. (2019). Cell-type-specific resolution epigenetics without the need for cell sorting or single-cell biology. *BioRxiv*, 437368.

[2] Bush, W. S., & Moore, J. H. (2012). Chapter 11: Genome-wide association studies. *PLoS computational biology*, 8(12), e1002822. doi:10.1371/journal.pcbi.1002822

[3] Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B. W., ... & Sullivan, P. F. (2016). Integrative approaches for large-scale transcriptome-wide association studies. *Nature genetics*, 48(3), 245.