



# Using information across tissues and genes to predict gene expression in Transcriptome-wide Association Studies (TWAS)

Fabio Morgante<sup>1</sup>, Gao Wang<sup>2</sup>, Matthew Stephens<sup>2,3</sup> and Yang I Li<sup>1,2</sup>

<sup>1</sup>Section of Genetic Medicine - Department of Medicine, <sup>2</sup>Department of Human Genetics, <sup>3</sup>Department of Statistics  
The University of Chicago, Chicago, IL

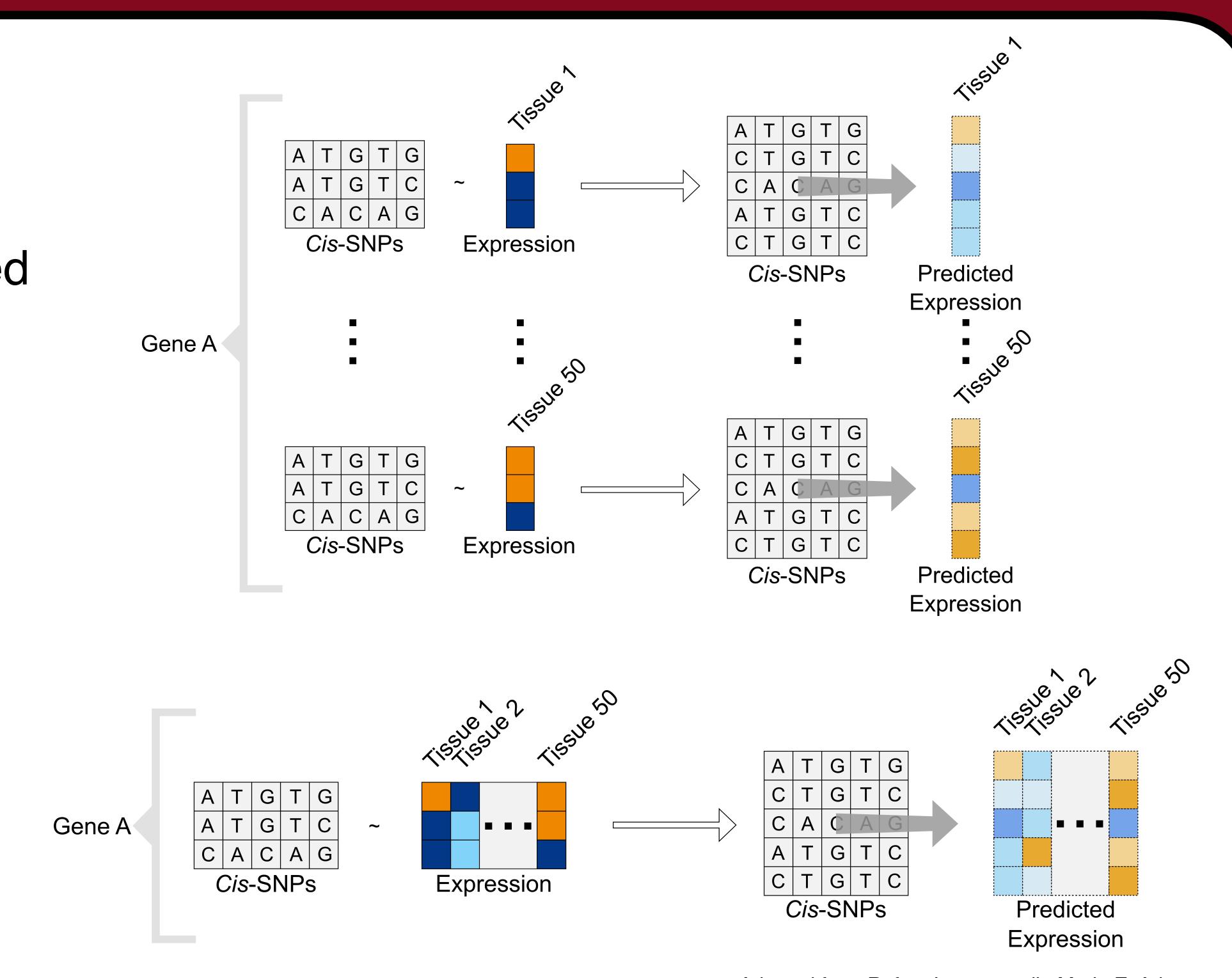
## Prediction methods in TWAS

### Single-tissue

- TWAS imputes gene expression based on genotypes in samples of genotyped and phenotyped individuals by leveraging genotyped expression reference panels (e.g., GTEx). Predicted gene expression is then tested for association with traits of interest.
- Methods initially developed for TWAS (e.g., PrediXcan, FUSION) perform gene expression imputation in one tissue at a time.
- However, these methods do not exploit patterns of expression Quantitative Trait Loci (eQTL) sharing among tissues.

### Multi-tissue

- Multivariate methods that perform gene expression imputation in all tissues simultaneously (e.g., fql, UTMOST) have been shown to provide greater accuracy than single-tissue methods.
- However, these methods analyze every gene separately: they do not exploit the information that multiple genes can provide.



## Our proposed multi-tissue prediction method (M&M)

- We have developed a multivariate Bayesian variable selection method for fine-mapping. Here, we evaluate its prediction performance.

$$\mathbf{Y} = \mathbf{XB} + \mathbf{E}$$

$$\mathbf{B} = \sum_{l=1}^L \mathbf{B}_l$$

$$\mathbf{B}_l = \gamma_l \mathbf{b}_l^T$$

$$\gamma \sim \text{Mult}(1, \pi)$$

$$\mathbf{b} \sim \sum_{k=1}^K \omega_k N_R(\mathbf{0}, \mathbf{U}_k)$$

$$\mathbf{E} \sim MN_{N \times R}(\mathbf{0}, \mathbf{I}, \mathbf{V})$$

- Main idea: Build a mixture with several covariance matrices,  $\mathbf{U}_k$ , that captures patterns of sharing and specificity among tissues. Then, let the data inform us which components are important by estimating their weights,  $\omega_k$ .
- Covariance matrices can be either data-driven (i.e., estimated from the data considering information from multiple genes) or canonical (e.g., identity).

For more details about the method, visit poster 2940W.

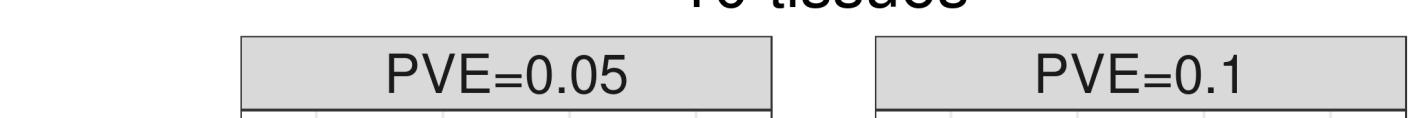
## Comparison of methods via simulation

- We simulated gene expression from actual GTEx genotypes:
  - 500 randomly sampled genes with ~600 individuals.
  - 1000 randomly sampled cis-SNPs (MAF > 0.01).
  - Either 5 or 10 tissues.
  - ~50% of genes = 1 causal variant (CV), ~30% of genes = 2 CVs and ~20% of genes = 3 CVs.
  - Effects were sampled from a multivariate normal distribution with different covariance structure across tissues:
    - Independent effects.
    - Shared effects.
    - Effects in only one tissue.
  - Per tissue proportion of variance explained (PVE) by CV(s) = 0.05 or 0.1.
  - Independent residuals among tissues, for simplicity.
- We compared our method to other methods:
  - Factored QTL (fql): multivariate regression method developed for TWAS<sup>1</sup>.
  - Sum of Single Effects (SuSiE): univariate regression method developed for fine-mapping<sup>2</sup>.
  - Elastic Net (enet): univariate regression method used in PrediXcan<sup>3</sup> and FUSION<sup>4</sup>.
- Prediction setting details:
  - Model trained on 80% of the data, prediction performance evaluated in the remaining 20%.
  - Accuracy evaluated as  $r^2(\text{true}, \text{predicted})$ .

## Results of the simulations

### Independent effects among all tissues

1	0	0
0	1	0
0	0	1



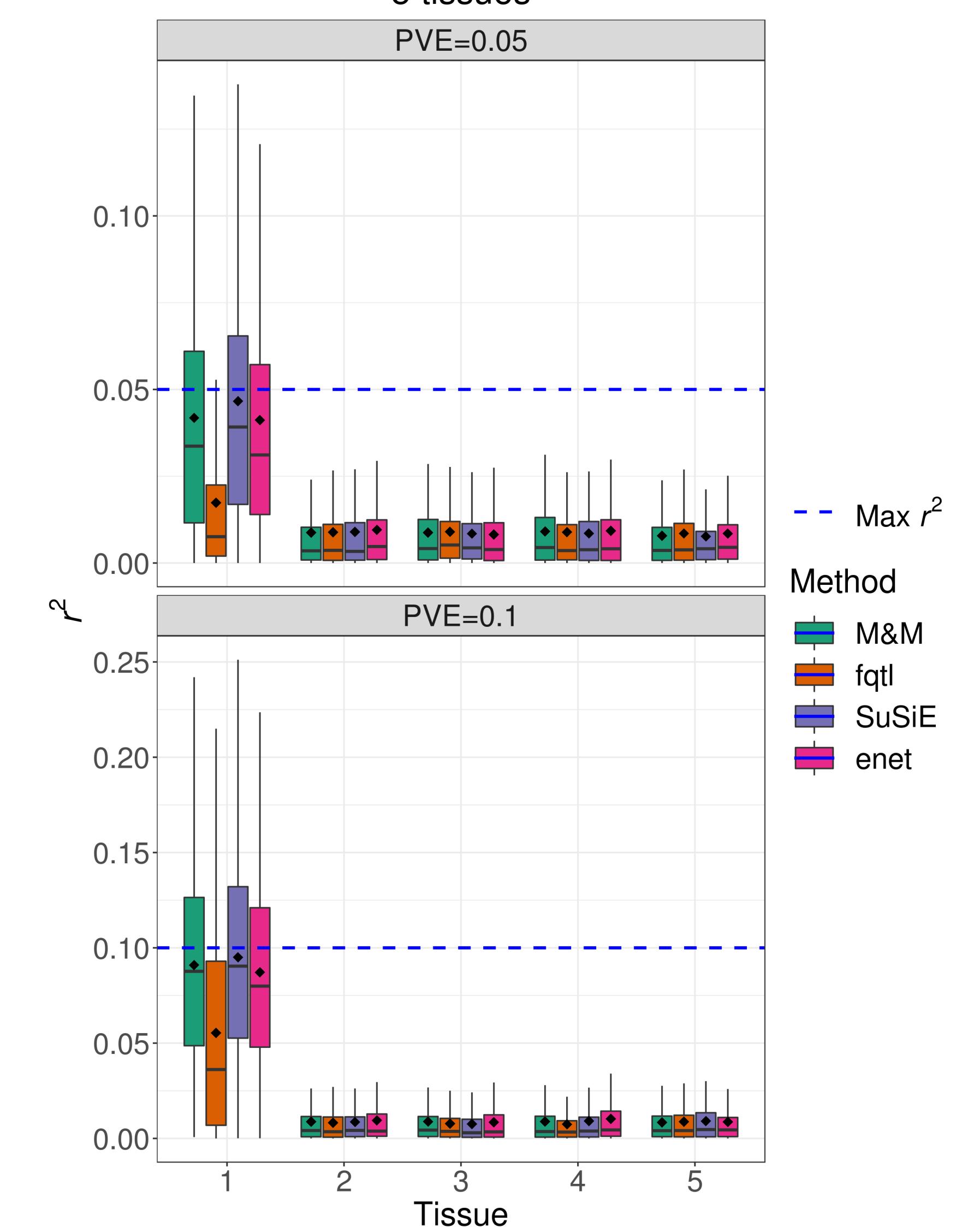
### Shared effects among all tissues

1	1	1
1	1	1
1	1	1



### Effects in only one tissue

1	0	0
0	0	0
0	0	0



All the methods performed similarly, although M&M performed slightly better with lower PVE

Multi-tissue methods performed better than single-tissue methods, especially with lower PVE

M&M retained performance in the presence of noise from tissues with no genetic effects

## Main conclusions

- Multi-tissue methods achieve higher accuracy than single-tissue methods, especially with a larger number of tissues with genetic effects and lower PVE.
- Our method performed well across all the simulated scenarios, even when the other multi-tissue method did not.
- Our method gives interpretable results due to concurrent fine-mapping (visit poster 2940W for details).

## References

- Gamazon *et al.* (2015). A gene-based association method for mapping traits using reference transcriptome data. *Nature genetics*, 47(9), 1091.
- Gusev *et al.* (2016). Integrative approaches for large-scale transcriptome-wide association studies. *Nature genetics*, 48(3), 245.
- Park *et al.* (2017). Multi-tissue polygenic models for transcriptome-wide association studies. *bioRxiv*, 107623.
- Wang *et al.* (2018). A simple new approach to variable selection in regression, with application to genetic fine-mapping. *bioRxiv*, 501114.