

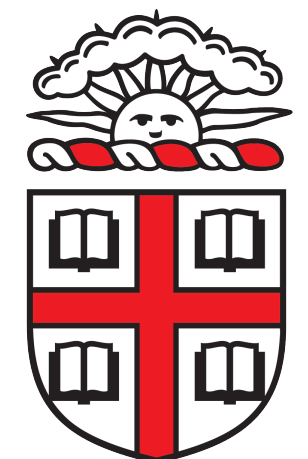
Multivariate MAPIT

Leveraging the genetic correlation between traits improves the detection of epistasis in genome-wide association studies

Julian Stamp
Center for Computational Molecular Biology
Brown University



CCMB



BROWN



Outline



Introduction



Methods



Simulated Data



Real Data



Summary

Genome Wide Association Studies

- Genotype large cohorts
- Map traits through statistical tests

20 years later & many open problems:

- Rare variants
- Inaccurate predictions
- Missing heritability



Phenotypic Variance

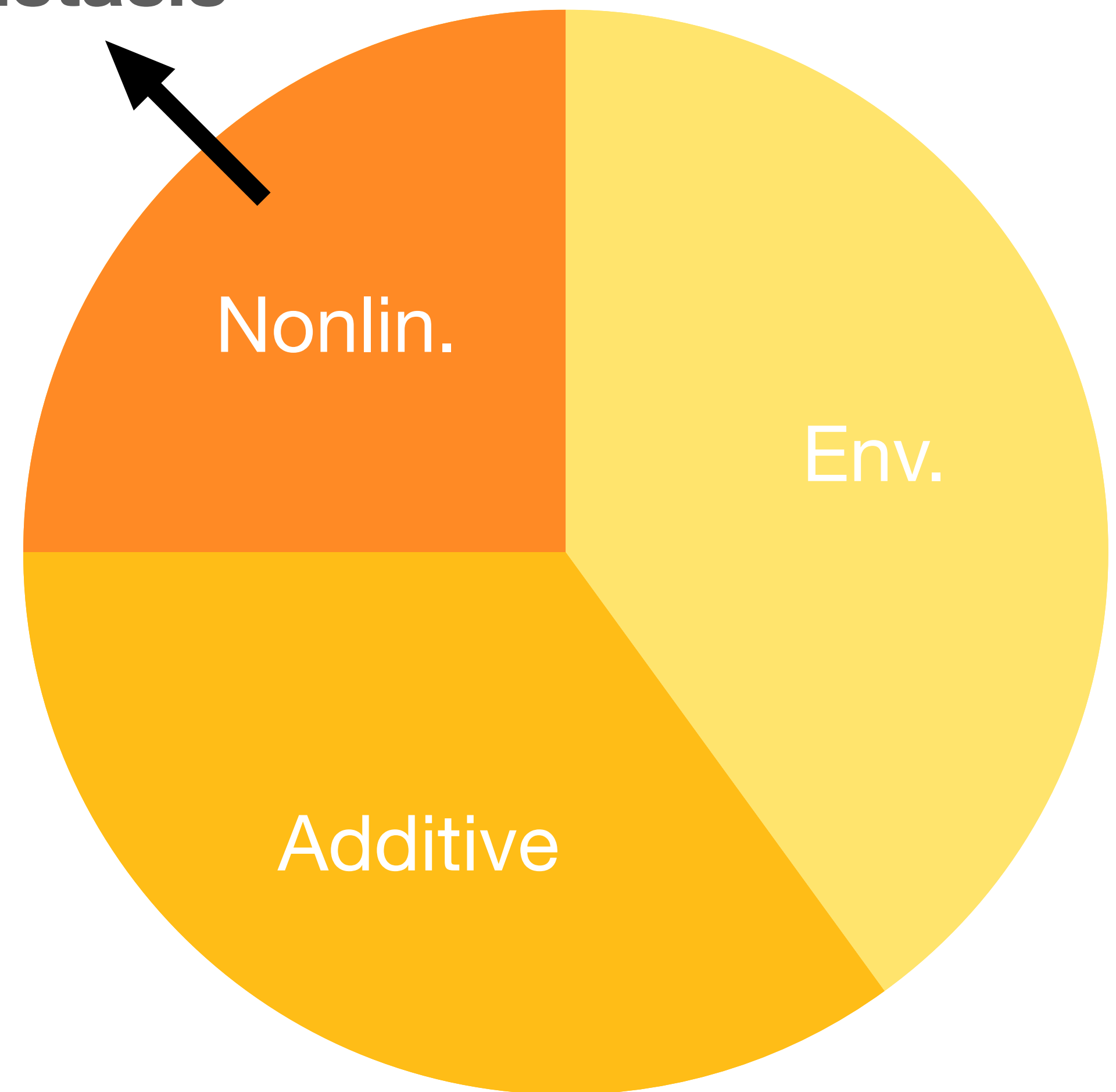
Genetic & Environmental Factors

$$P = G + E$$

Broad sense Heritability

$$H^2 = \frac{\text{Var}[G]}{\text{Var}[P]}$$

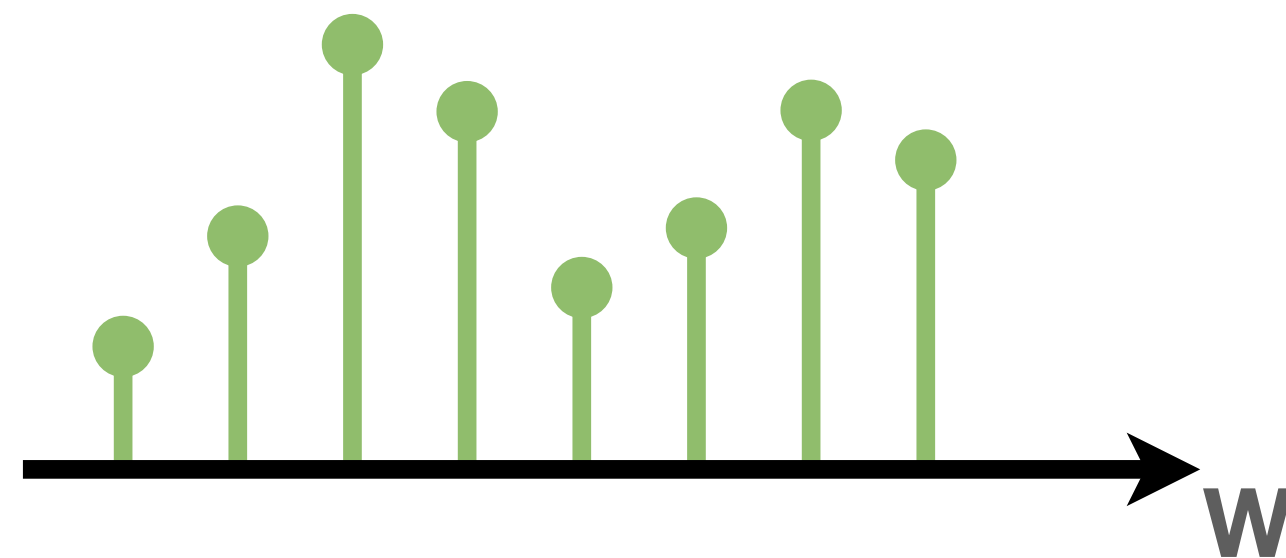
Epistasis



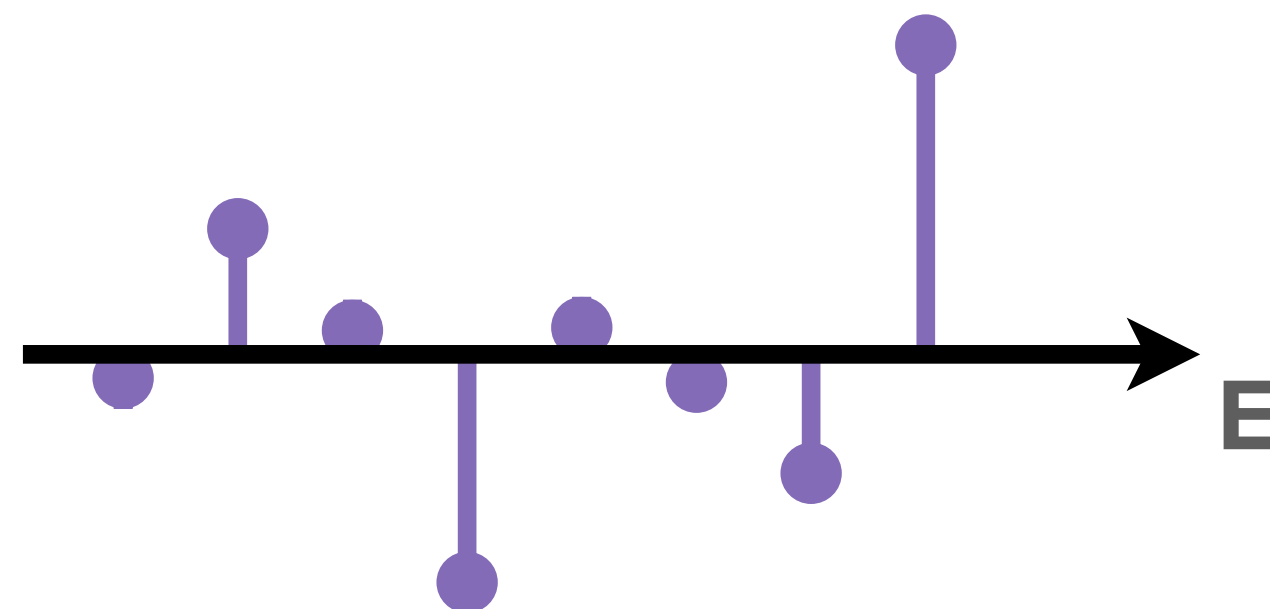
Hadamard-Walsh Transform

- For combinatorially complete data transform linearly from trait to epistatic effects
- Rows are mapped to genotypes

Trait domain



Effect size domain



$$\frac{1}{2^L} \Psi \vec{W} = \vec{E}_W$$

$$\frac{1}{2^L} \cdot \begin{pmatrix} +1 & +1 & +1 & \dots & +1 \\ +1 & -1 & +1 & \dots & -1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ +1 & -1 & -1 & \dots & -1 \end{pmatrix} \cdot \begin{pmatrix} 1.07 \\ 2.54 \\ \vdots \\ 0.41 \end{pmatrix} = \begin{pmatrix} 2.28 \\ 0.0857 \\ \vdots \\ -0.1051 \end{pmatrix}$$

Hadamard matrix Trait values Effect size of interactions

Biobank Scale Data

- $\sim 10^5$ to 10^6 variants
- $\sim 10^4$ to 10^5 samples
- $\sim 3^{1000000}$ genotype combinations*

\implies underdetermined & combinatorially incomplete

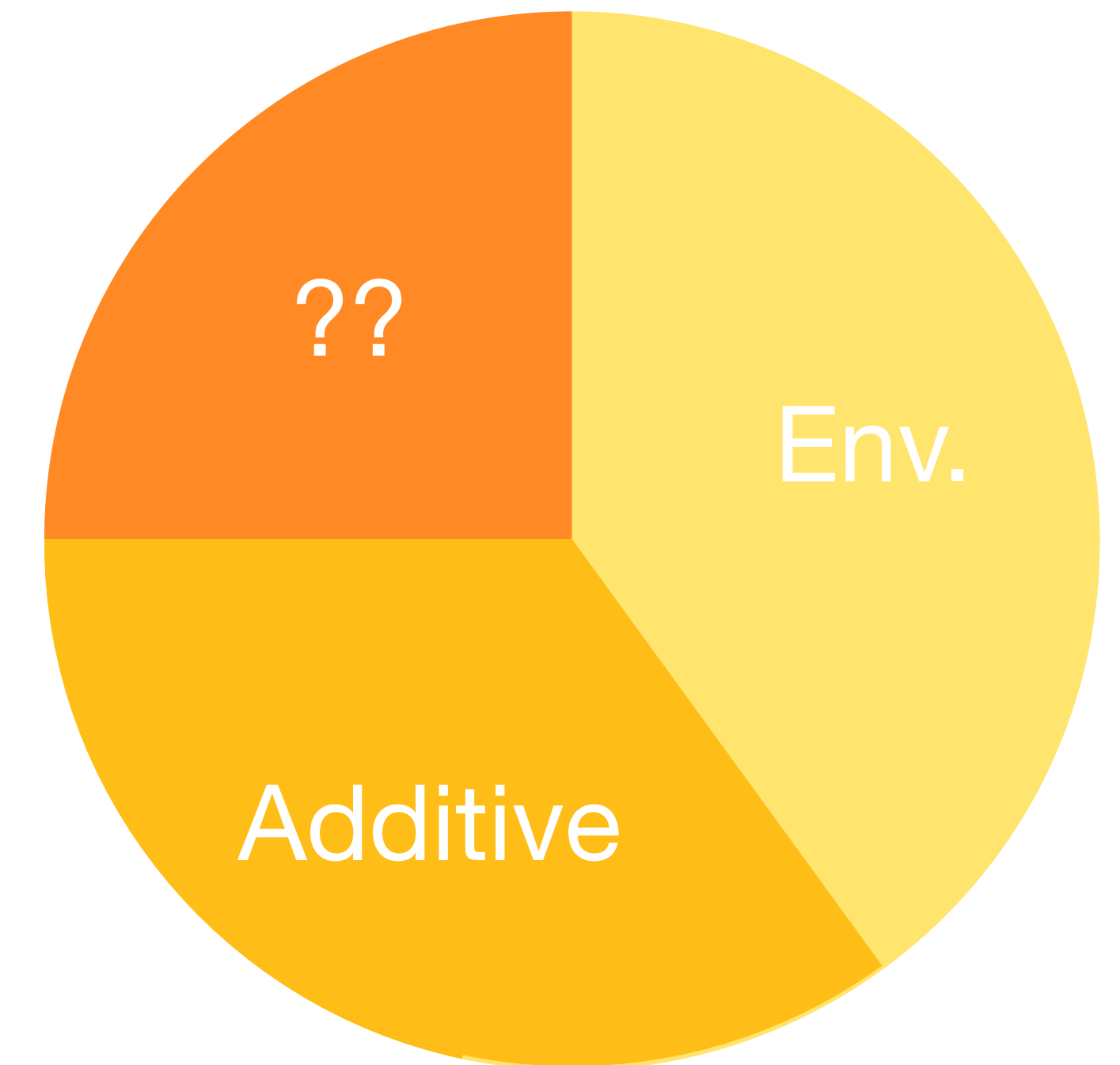


*Number of atoms in the universe: $\sim 10^{80}$

Motivation

Multivariate approach to studying non-linear contributions in complex traits

- More than 11 million SNPs in human genome¹, ~400k trait associations²
- Majority of the heritability of complex traits “missing”³
- Epistasis could explain missing heritability
- Computational methods to detect epistasis are underpowered or computationally resource intensive⁴



1 Madsen et al. (2007), *Genome Research*

2 Sollis et al. (2022), *Nucleic Acids Research*

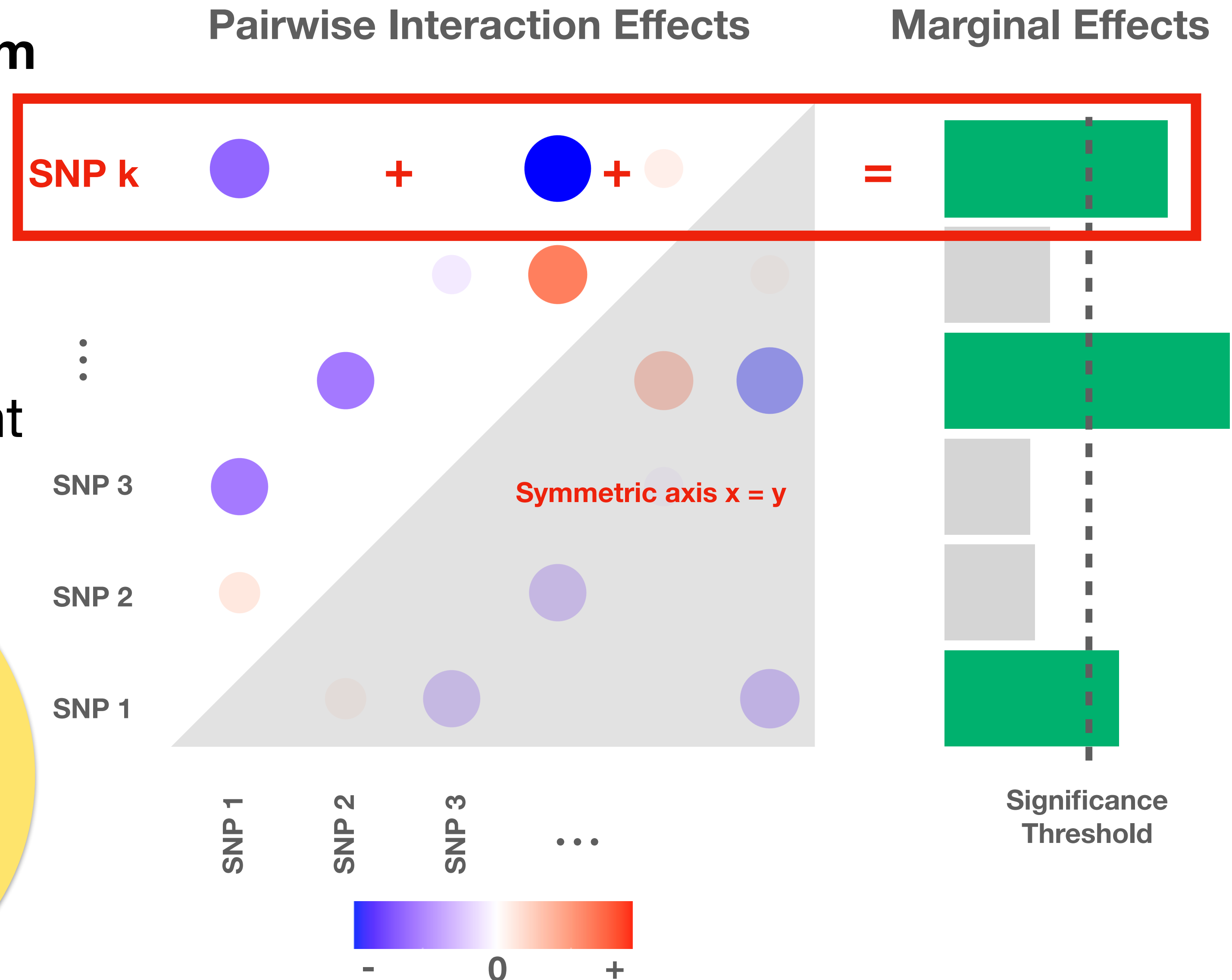
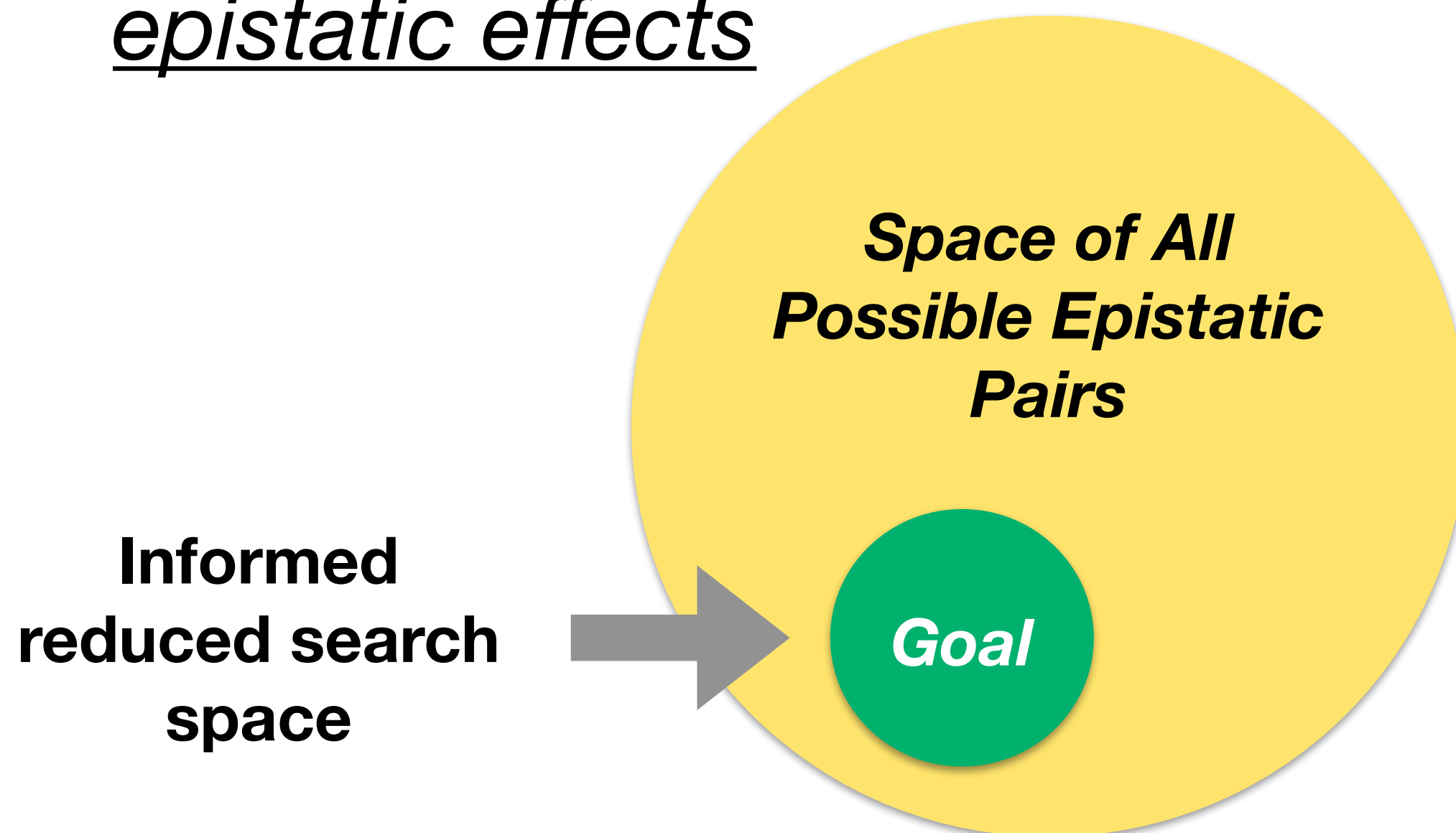
3 Young (2019), *PLOS Genetics*

4 Crawford et al. (2017), *PLOS Genetics*

Explicit search space

Epistasis as combinatorial problem

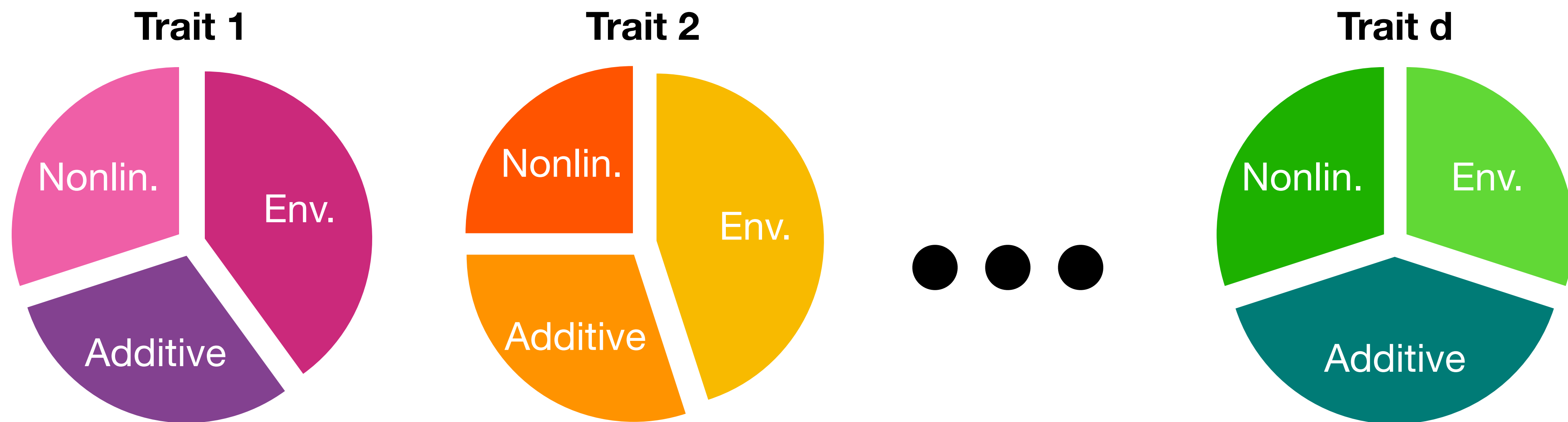
- There are $p(p - 1)/2$ possible interacting pairs for p SNPs
- **Idea:** Prioritize search for variant interactions using marginal epistatic effects



Multivariate LMM

- Genetic correlations between traits maintained by pleiotropy¹
- Multivariate modelling improves GWAS²

⇒ Can we leverage **genetic correlations** to improve detection of epistasis?

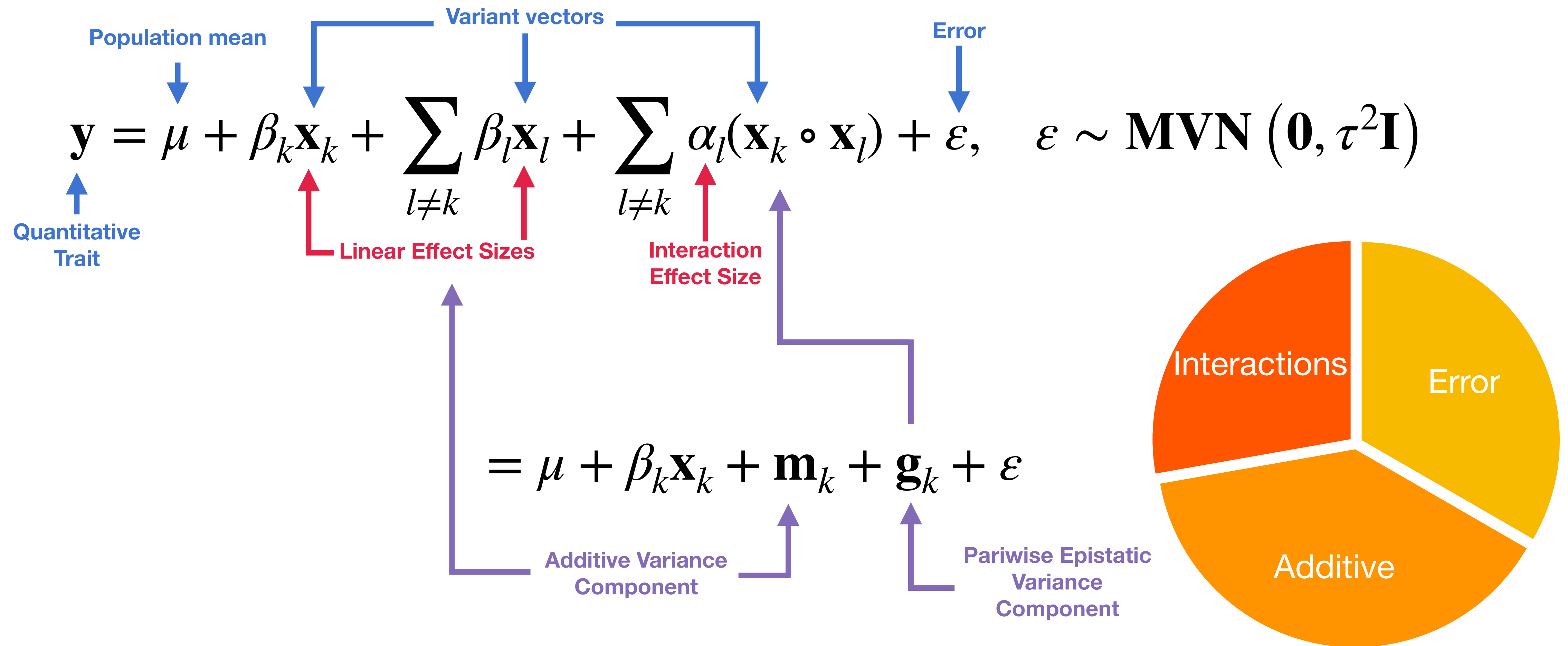


¹ Chebib and Guillaume (2021), *Genetics*

² Zhou and Stephens (2014), *Nature*

Approach

Starting point: The Marginal Epistasis Test (MAPIT)



Approach

Normal assumption for effect size trick for underdetermined data

- Genetic Relatedness Matrix

$$\mathbf{K} = \mathbf{X}_{-k} \mathbf{X}_{-k}^T$$

- Covariance of the interaction of SNP k with it's background

$$\mathbf{G} = \mathbf{D}_k \mathbf{K} \mathbf{D}_k \text{ with}$$

$$\mathbf{D}_k = \text{diag}(\mathbf{x}_k)$$

- Estimate variance parameters jointly using MQS

$$\mathbf{y} = \mu + \beta_k \mathbf{x}_k + \mathbf{m}_k + \mathbf{g}_k + \varepsilon$$

$$\mathbf{m}_k \sim \text{MVN}(\mathbf{0}, \omega^2 \mathbf{K})$$

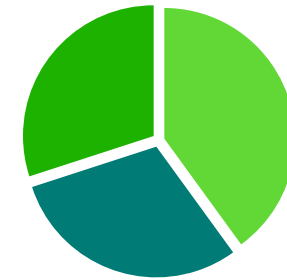
$$\mathbf{g}_k \sim \text{MVN}(\mathbf{0}, \sigma^2 \mathbf{G})$$

$$\varepsilon \sim \text{MVN}(\mathbf{0}, \tau^2 \mathbf{I})$$

Approach

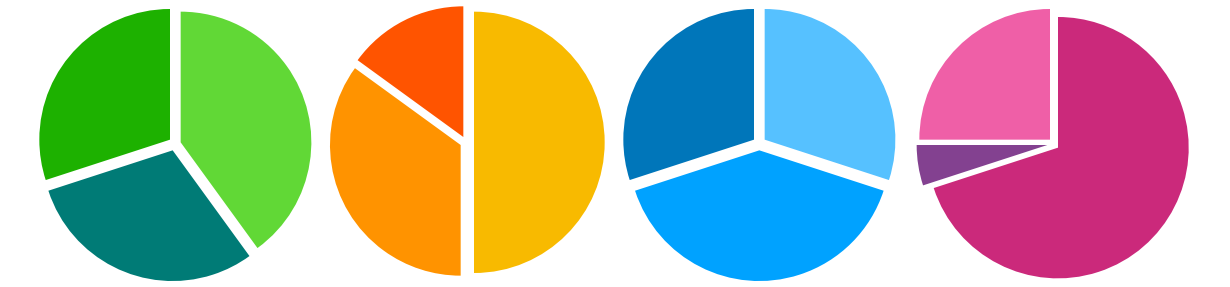
Multivariate extension of MAPIT (mvMAPIT)

MAPIT



- One trait $\mathbf{y} = (y_1, \dots, y_n)^\top$
- Only covariance between samples
 $\mathbf{g}_k \sim \text{MVN}(\mathbf{0}, \sigma^2 \mathbf{G})$
- Estimate variance components
 $\hat{\sigma}^2 = \mathbf{y}^\top \mathbf{A}_k \mathbf{y}$

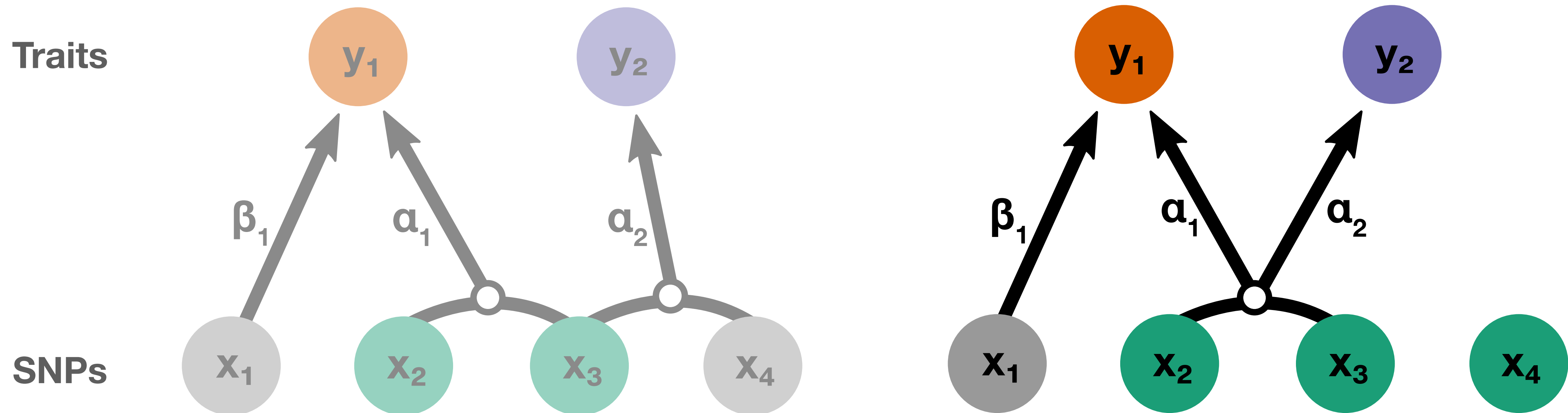
mvMAPIT



- Many traits $\mathbf{Y} = \begin{pmatrix} y_{11} & \cdots & y_{1d} \\ \vdots & \ddots & \vdots \\ y_{n1} & \cdots & y_{nd} \end{pmatrix}$
- Covariance between samples and variance components
 $\mathbf{g}_k \sim \text{MN}_{n \times d}(\mathbf{0}, \mathbf{V}_G, \sigma^2 \mathbf{G})$
- Estimate d choose 2 variance and covariance components $\hat{\sigma}_{12}^2 = \mathbf{y}_1^\top \mathbf{A}_k \mathbf{y}_2$

mvMAPIT

Modelling cross-trait genetic correlations of interaction effects



MAPIT

Simulations of complex traits

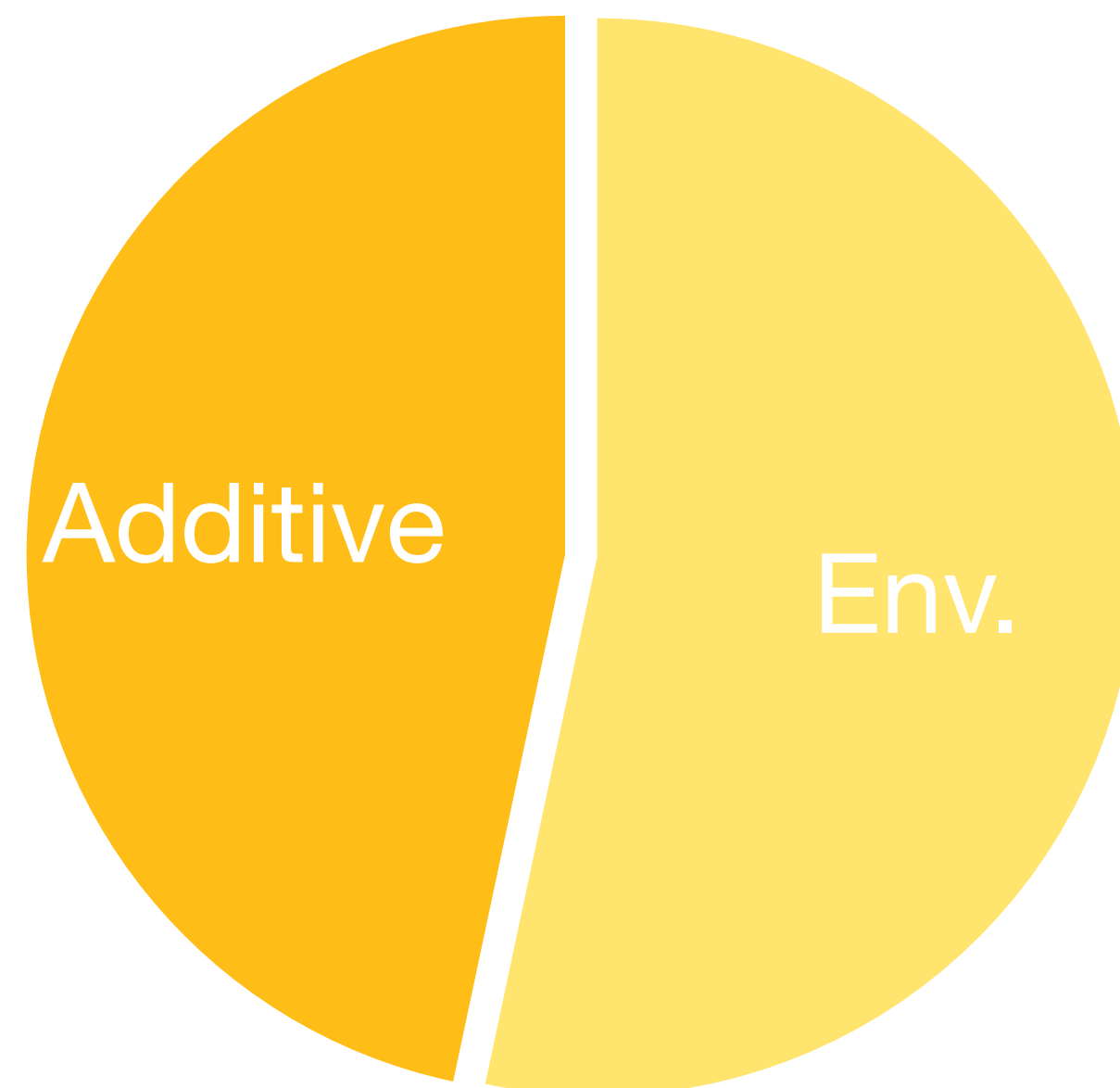


Scenarios

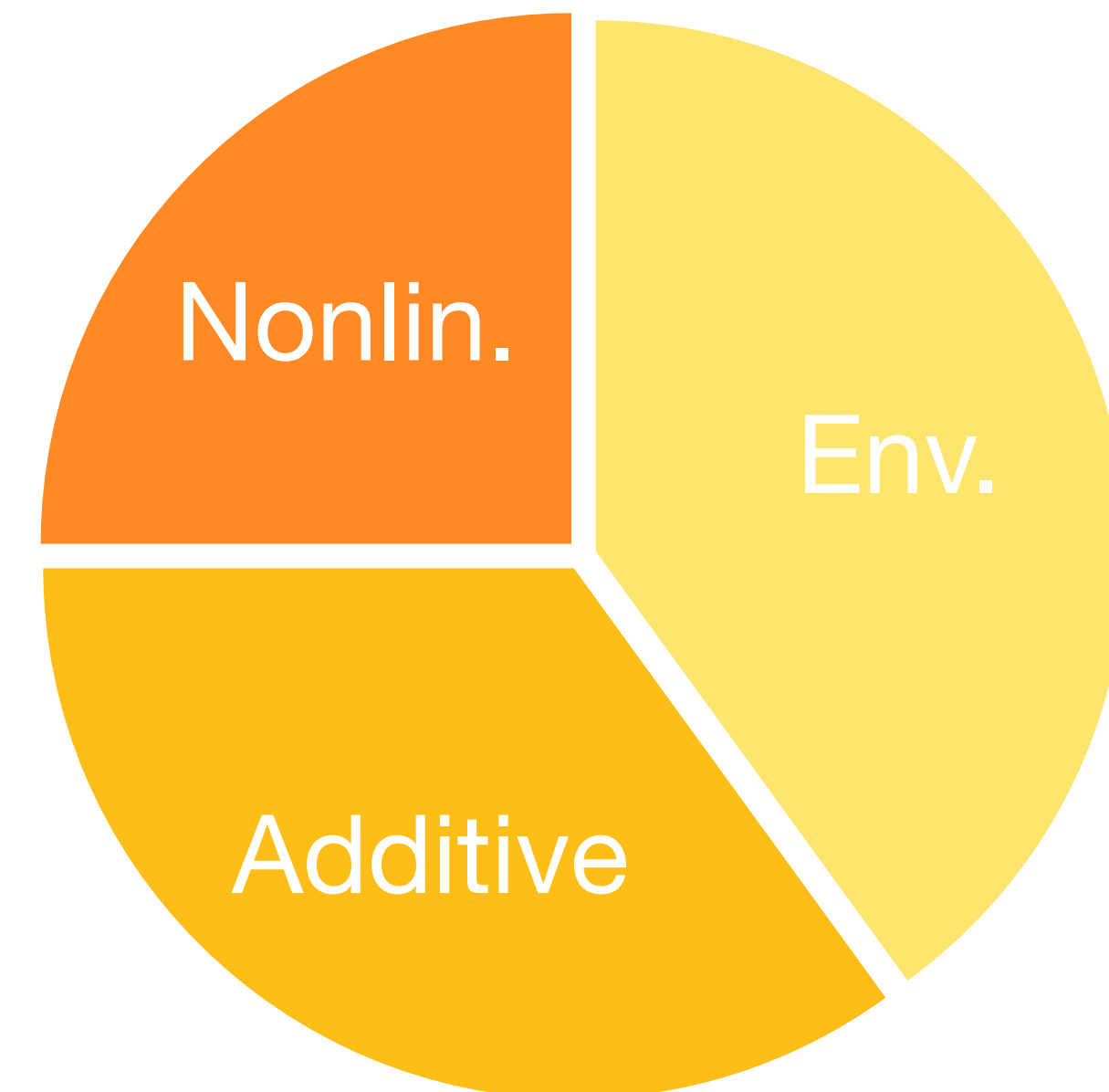
- Null Hypothesis true: no epistasis
- Epistasis with varying parameters

Parameters

- Broad sense heritability H^2
- Proportion of heritable variance due to epistasis $H^2(1 - \rho)$



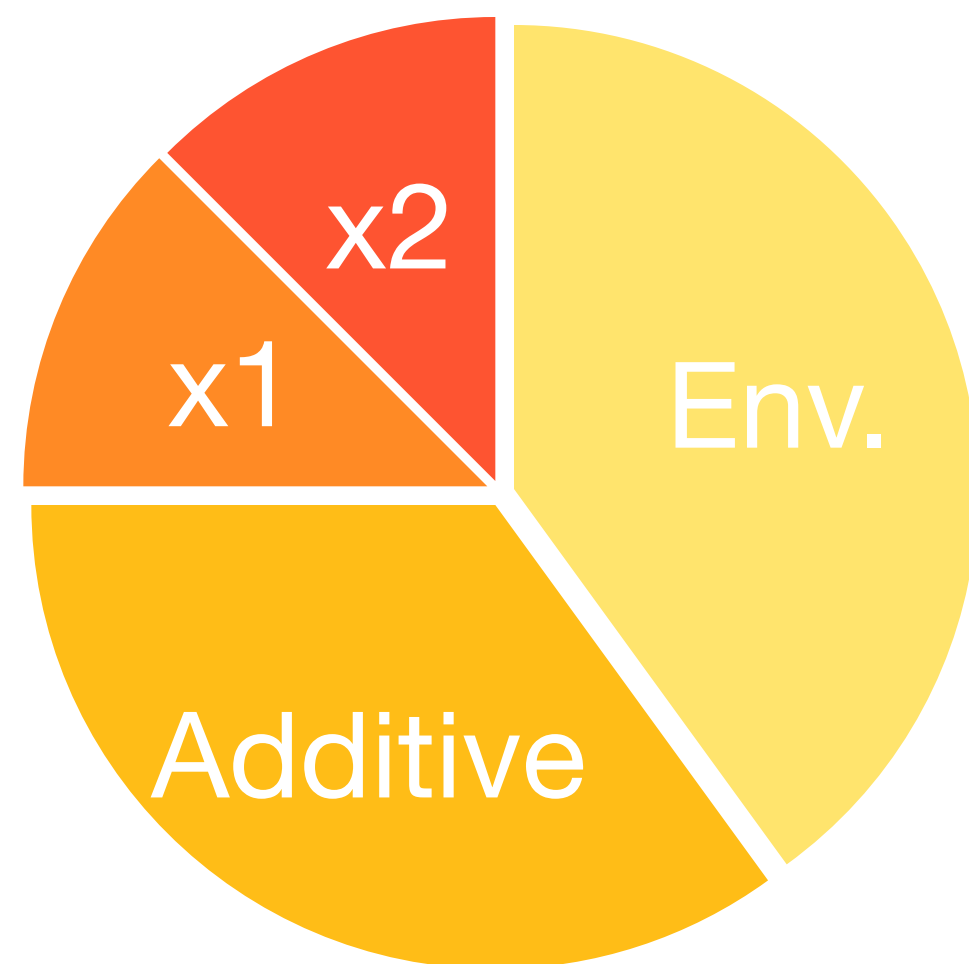
$$H_0 : \sigma^2 = 0$$



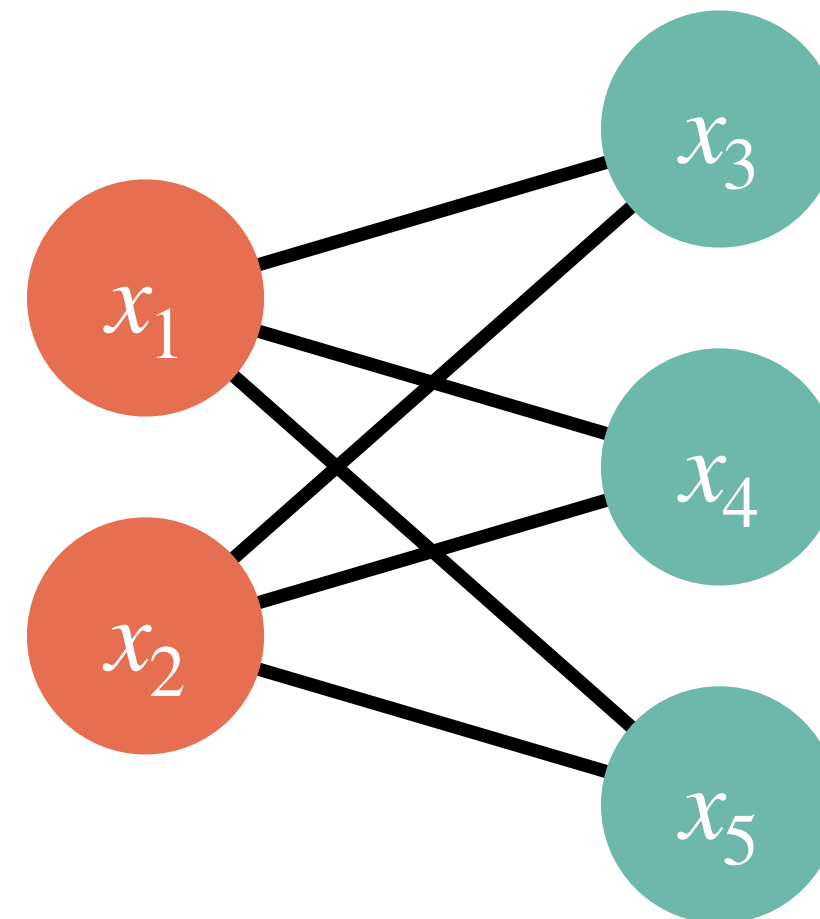
$$H_1 : \sigma^2 \neq 0$$

MAPIT

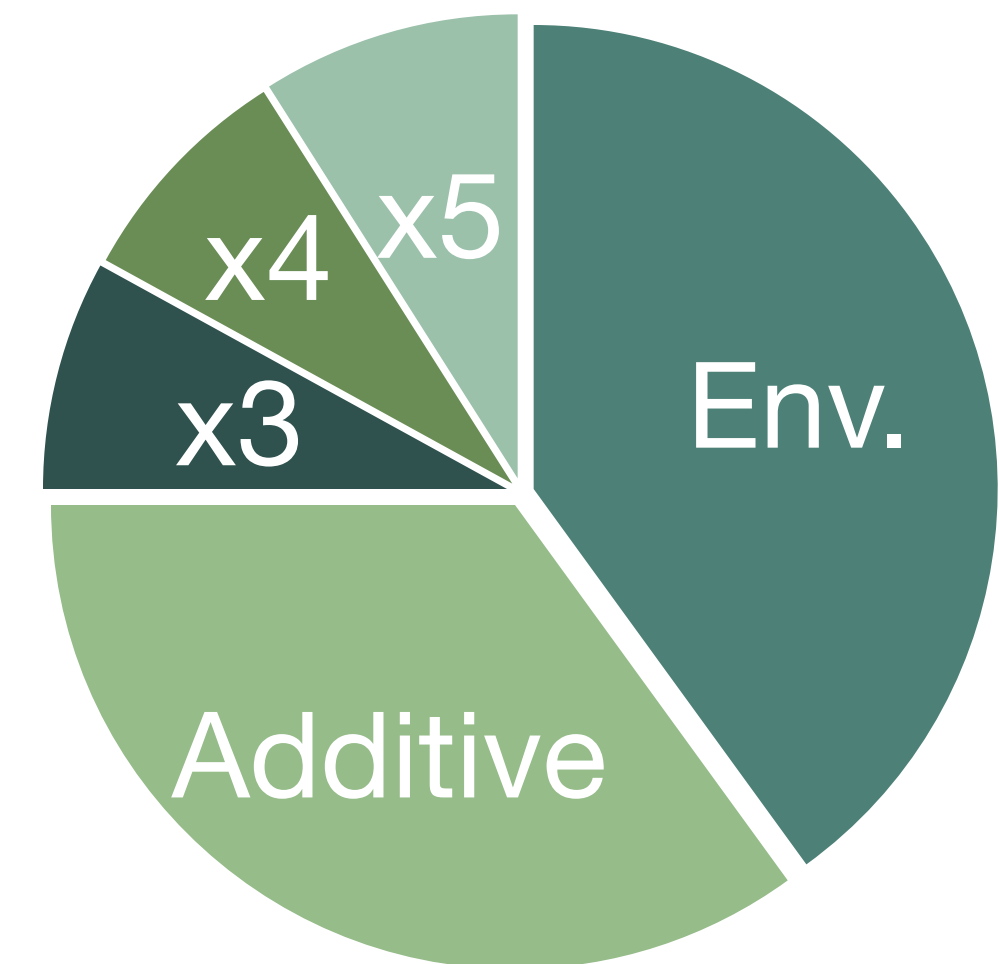
Simulations of complex traits



Group 1



Group 2



- Additive SNPs
- Epistatic Group 1
- Epistatic Group 2

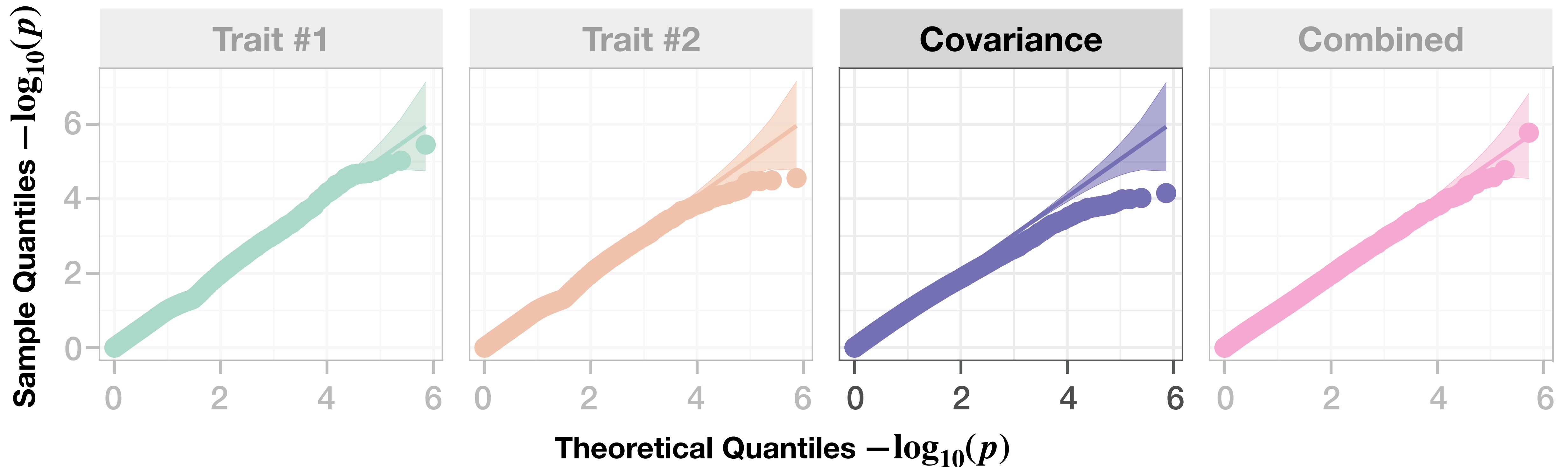
Marginal epistasis e.g.

$$\bullet \mathbf{g}_{x_1} = (\mathbf{x}_1 \circ \mathbf{x}_3) \cdot \alpha_{13} + (\mathbf{x}_1 \circ \mathbf{x}_4) \cdot \alpha_{14} + (\mathbf{x}_1 \circ \mathbf{x}_5) \cdot \alpha_{15}$$

$$\bullet \mathbf{g}_{x_3} = (\mathbf{x}_1 \circ \mathbf{x}_3) \cdot \alpha_{13} + (\mathbf{x}_2 \circ \mathbf{x}_3) \cdot \alpha_{23}$$

QQ-Plots*

mvMAPIT is well calibrated

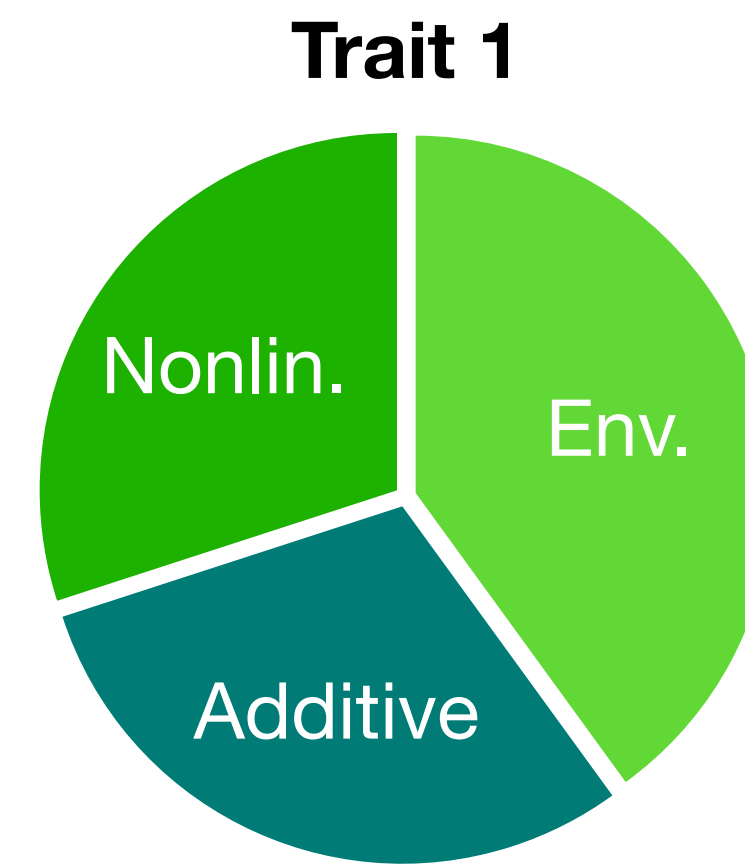


* Simulated data with null hypothesis true

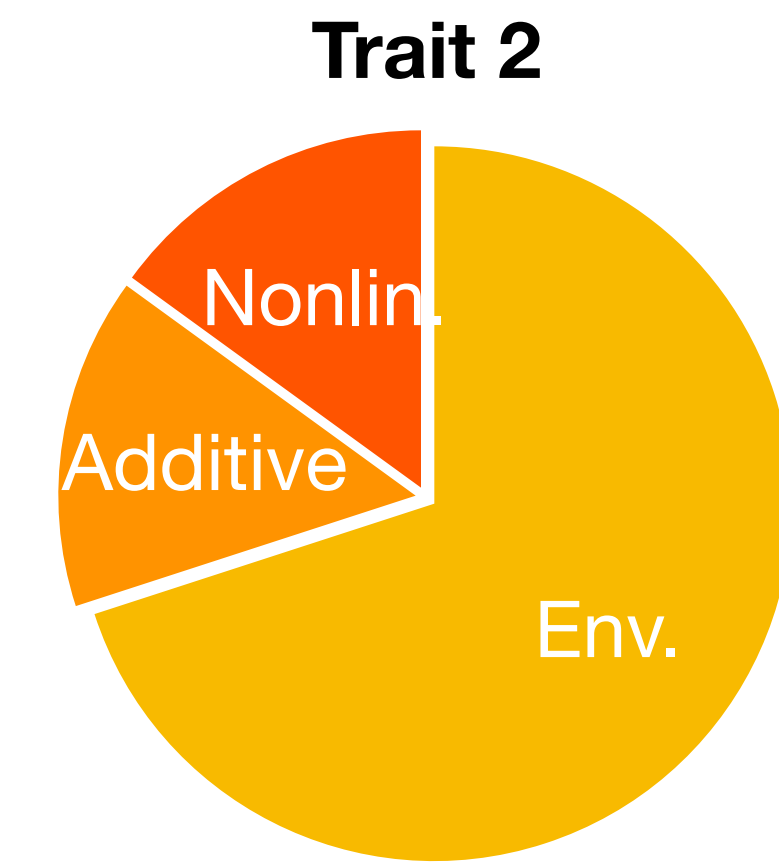
Empirical Power

Genetic correlations improve power of mvMAPIT

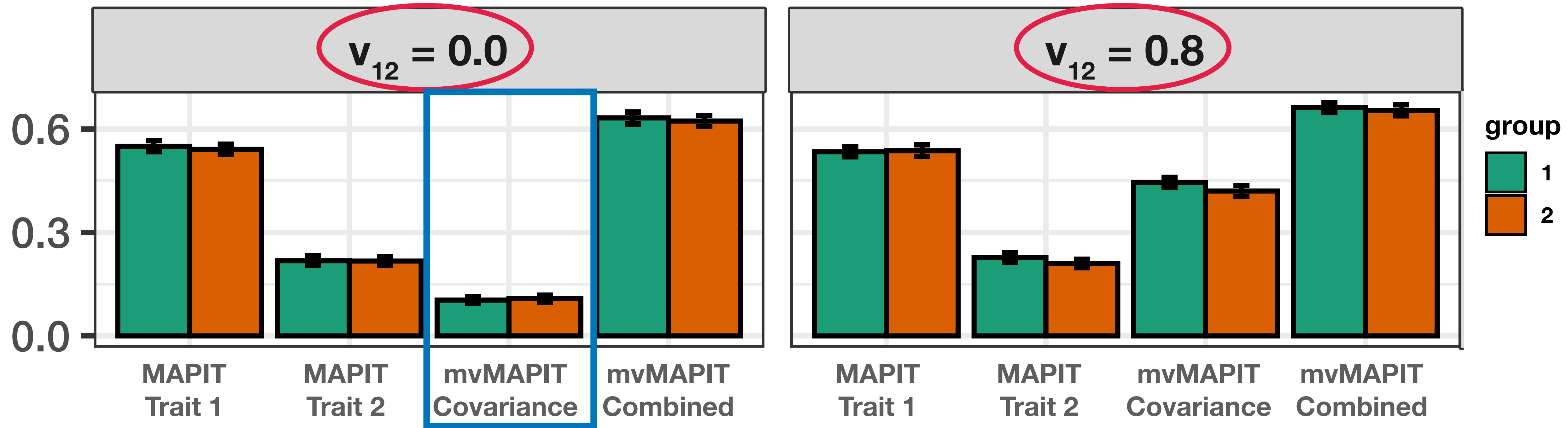
Correlation between
epistatic effect sizes V_{12}



$$H^2 = 0.6$$



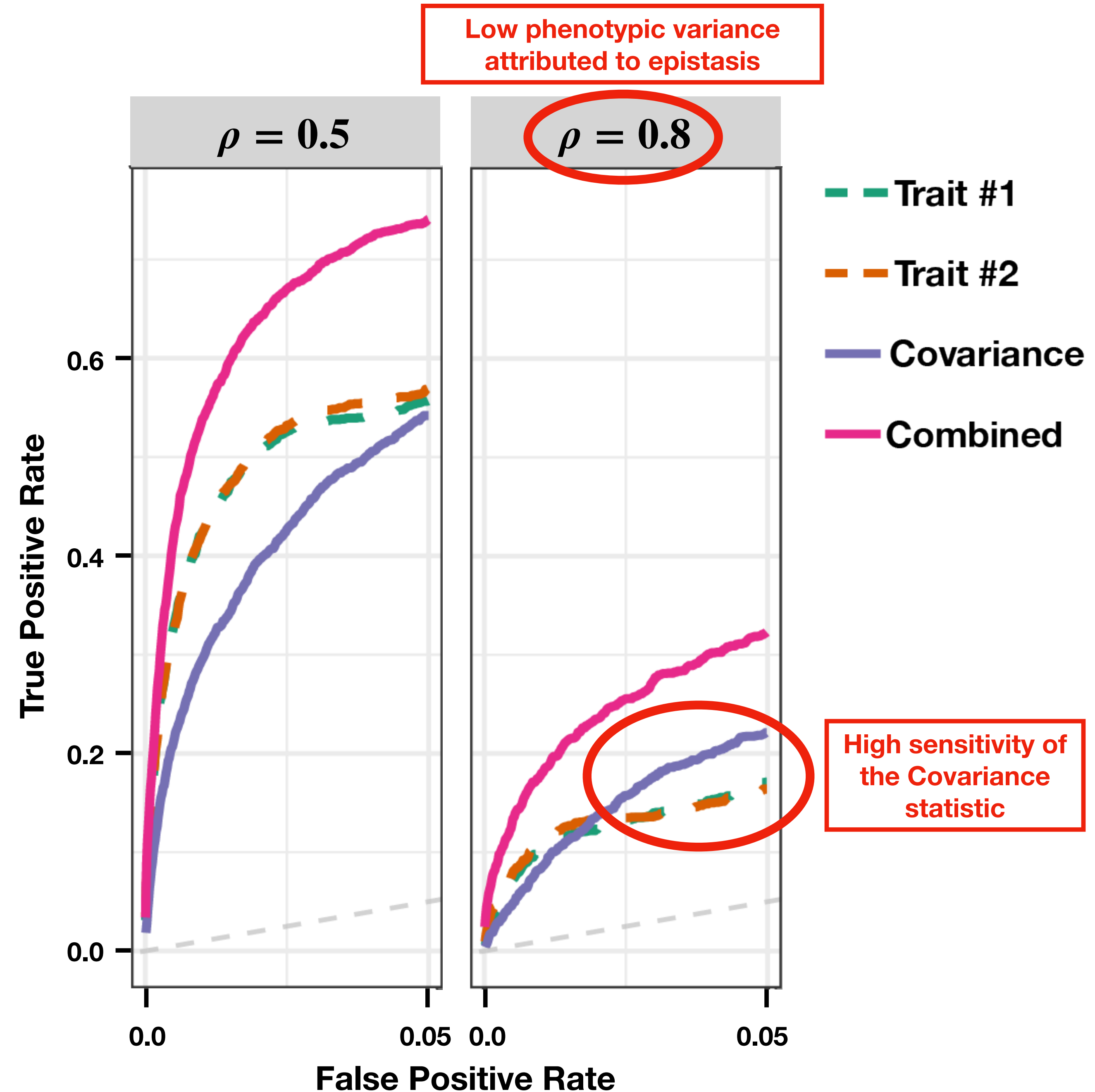
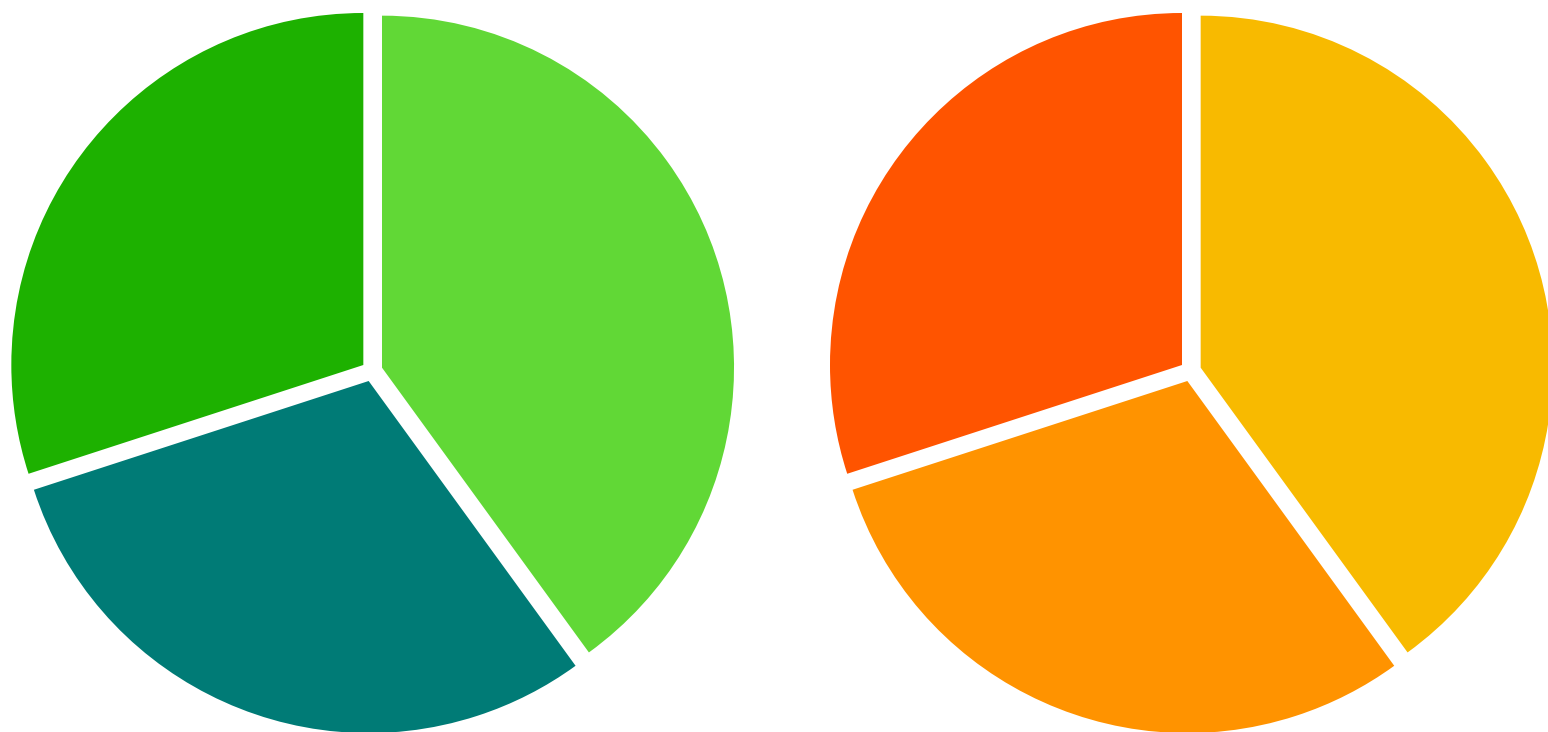
$$H^2 = 0.3$$



ROC Curves

Genetic correlations increase sensitivity of mvMAPIT

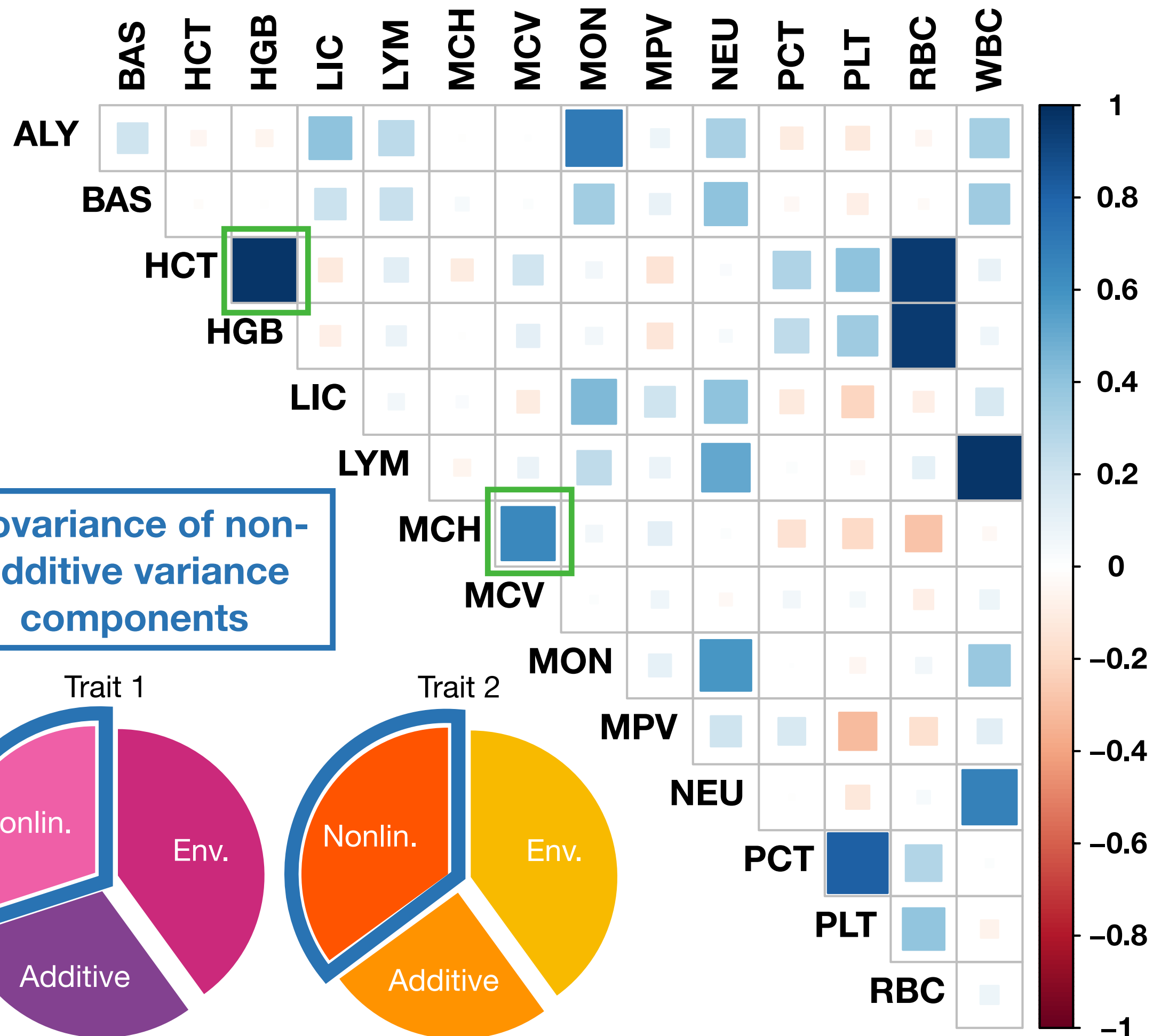
- Heritability $H^2 = 0.6$
- High correlation $\nu_{12} = 0.8$
- Steeper curves indicate higher sensitivity



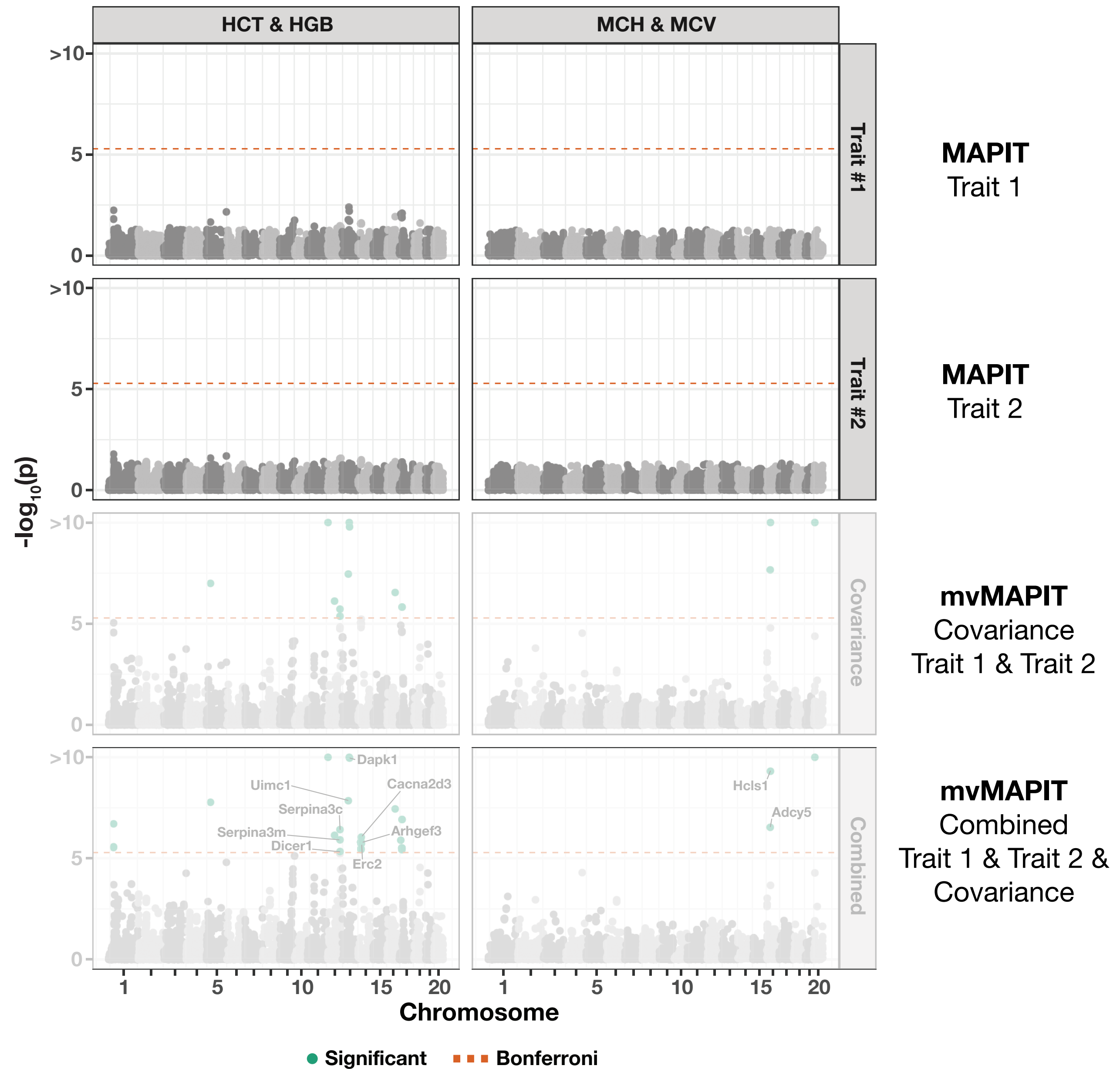
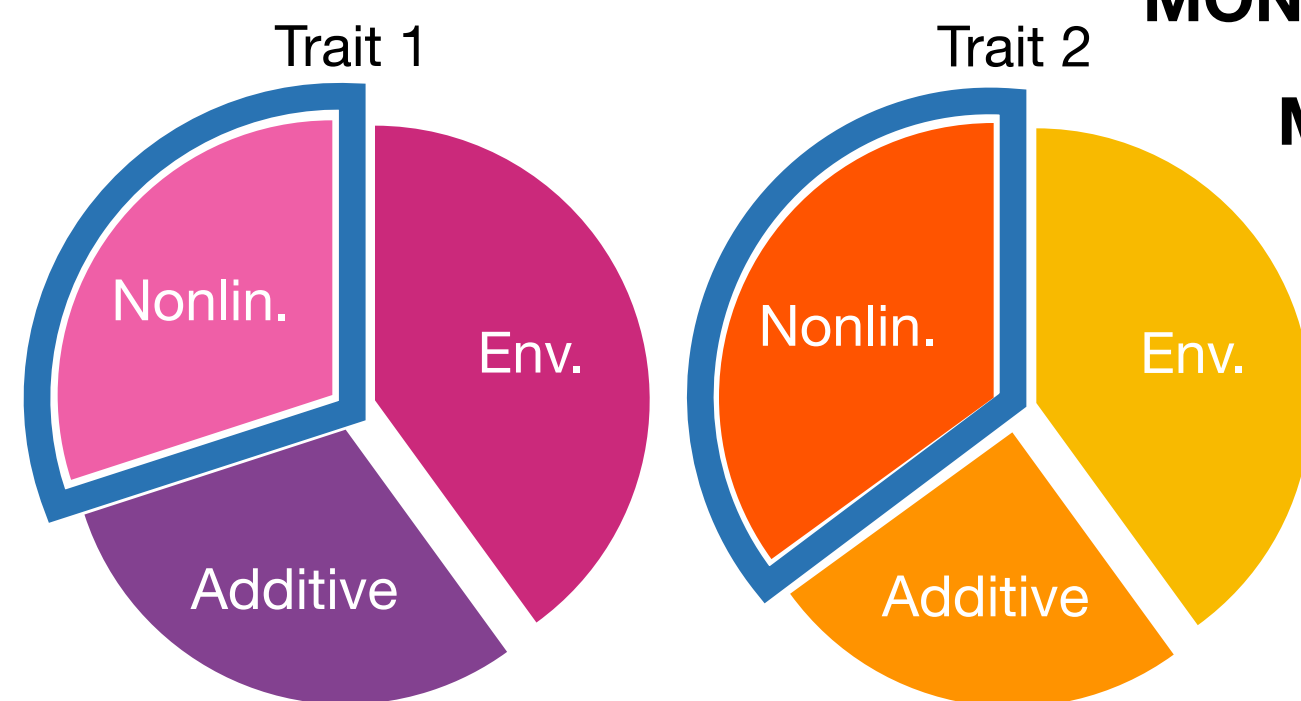
*Trait 1 & Trait 2 correspond to baseline (MAPIT)

Real Data*

Genetic correlations reveal strong signal of epistasis



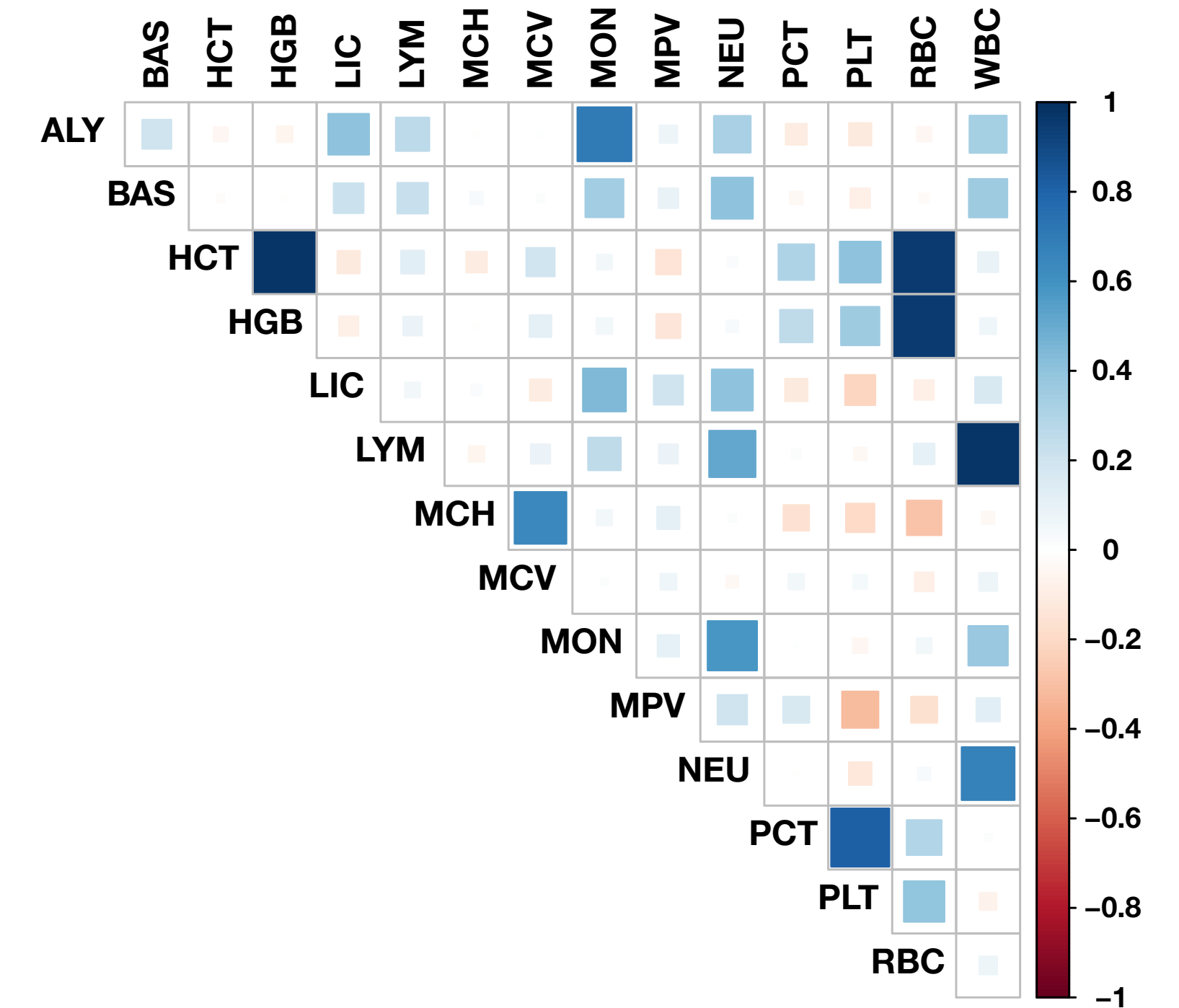
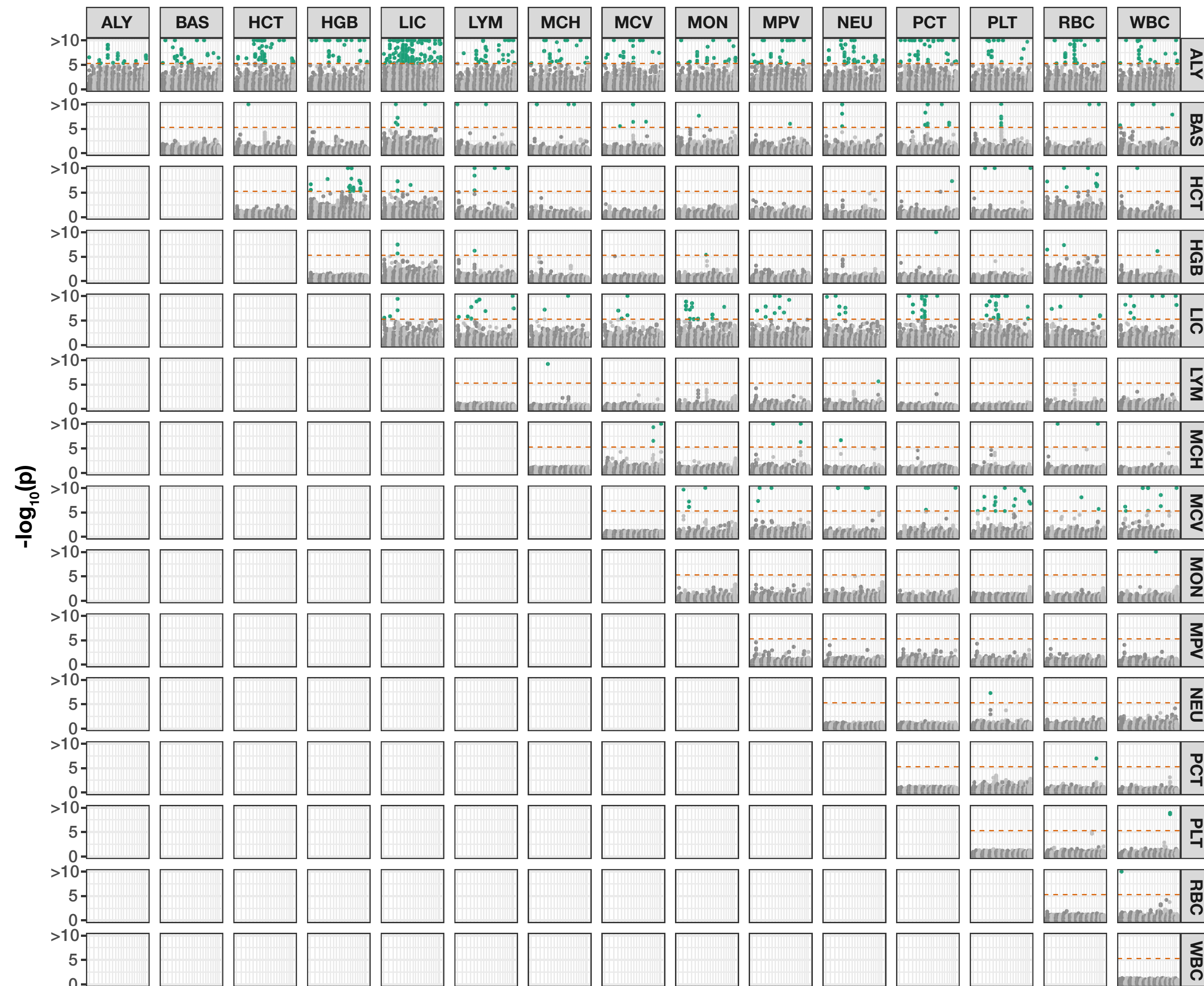
Covariance of non-additive variance components



* Hematology traits of WTCCC Mice

Real Data*

New significant associations across many trait pairs

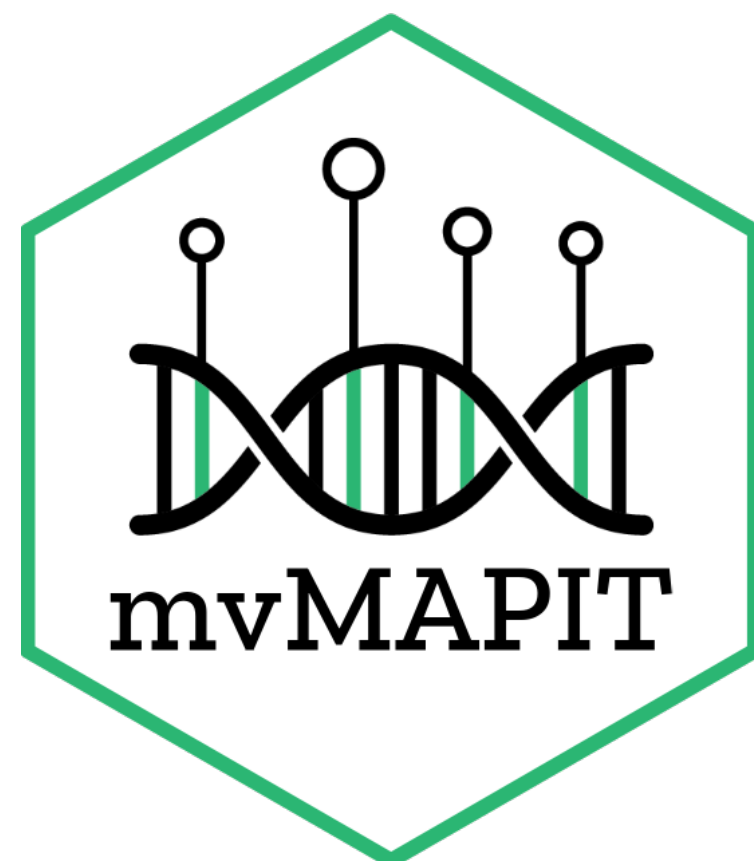


* Hematology traits of WTCCC Mice

● Significant ■ Bonferroni

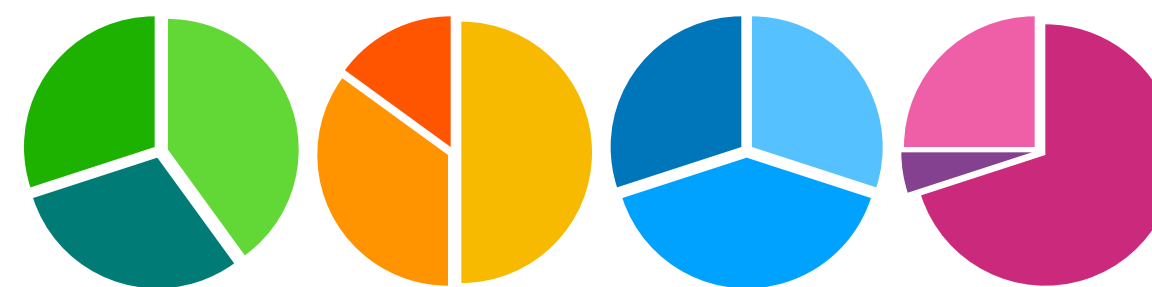
Weaknesses

- Time complexity scaling with sample size
- Unknown interaction partner
- Meta analysis p-value interpretability



Strengths

- Analysis of shared genetic architecture of traits
- Correlation between effects improves sensitivity
- Marginalisation accumulates weak effects to strong signal
- Marginalisation reduces search space



Acknowledgements

Advisors

Lorin Crawford
Dan Weinreich

Crawford Lab and Dave

Chibuikem Nwizu
Dave Peede
Ria Vinod
Alex Wong
Emily Winn
Ashley Conard

Dana Edwin
Wai Shing Tang
Whitney Sloneker
Yu Zhong
Collin Small
Ryan Huang



CCMB



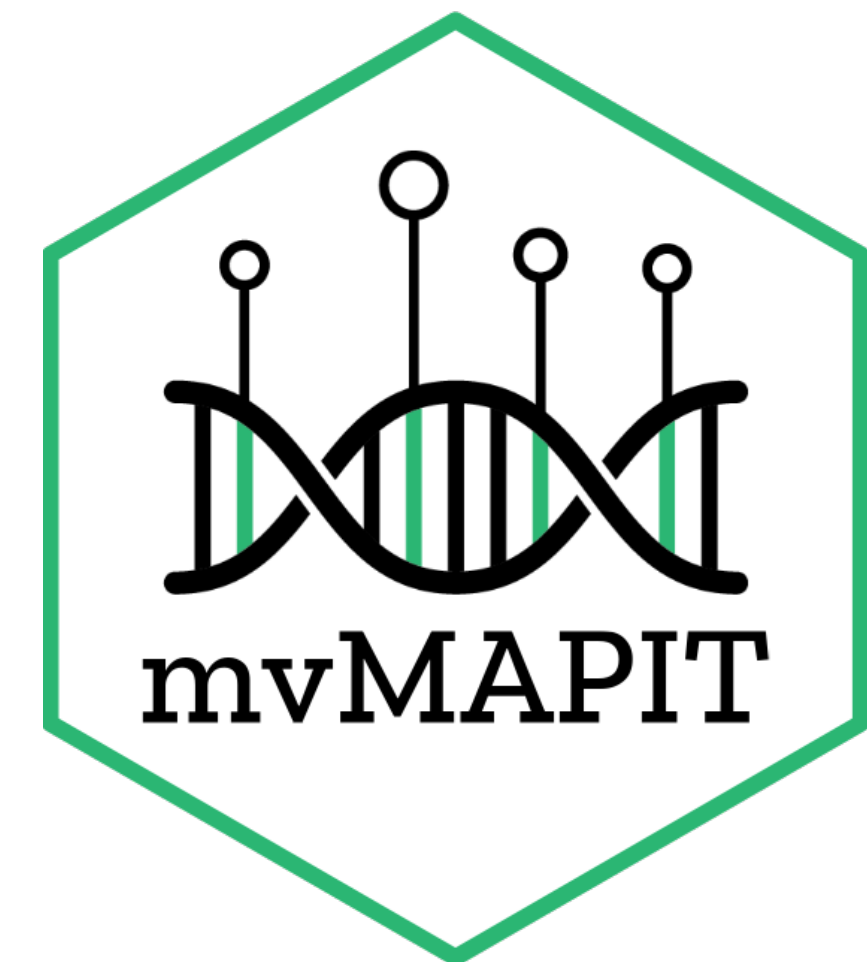
BROWN



mvMAPIT

- Code and documentation on GitHub: <https://lcrawlab.github.io/mvMAPIT/>
- R package published on CRAN: <https://cran.r-project.org/package=mvMAPIT>

```
install.packages ( 'mvMAPIT' )
```



Relevant References

Variance Component Estimation

- X. Zhou. "A unified framework for variance component estimation with summary statistics in genome-wide association studies." *Ann. Appl. Stat.* 11 (4) 2027 - 2051, December 2017. <https://doi.org/10.1214/17-AOAS1052>

Marginal Epistasis Detection

- L. Crawford, P. Zeng, S. Mukherjee, & X. Zhou, (2017). Detecting epistasis with the marginal epistasis test in genetic mapping studies of quantitative traits. *PLOS Genetics*, 13(7), e1006869. <https://doi.org/10.1371/journal.pgen.1006869>
- **J. Stamp**, A. DenAdel, D. Weinreich, & L. Crawford, (2023). Leveraging the Genetic Correlation between Traits Improves the Detection of Epistasis in Genome-wide Association Studies. *G3 Genes|Genomes|Genetics*, jkad118. <https://doi.org/10.1093/g3journal/jkad118>

Related Software/Source Code:

- mvMAPIT: <https://lcrawlab.github.io/mvMAPIT/>

