



# TESTING FOR GENE-ENVIRONMENT (GxE) INTERACTION USING P-VALUE AGGREGATION IDENTIFIES MANY GxE LOCI

---

Saurabh Mishra

Advisor: Dr. Arunabha Majumdar

Department of Mathematics, IIT Hyderabad

- **Background**
  - Genetics and inheritance models
  - Model misspecification
  - Gene–environment interactions
- **Our approach**
  - Objective
  - P-value aggregation framework
- **Results**
  - Simulation studies
  - UK Biobank applications
- **Conclusions and Future Directions**

## BASICS OF GENETICS

---

# PHENOTYPE = GENETICS + ENVIRONMENT + ...



=



+



**Variation in a Phenotype**

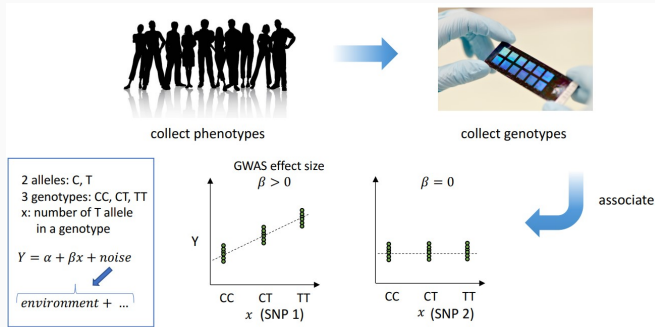
**Genetic Variation**

**Variation Due to  
Environmental Factors**

e.g., cholesterol level = genetics + diet + exercise + ...

# GENOME-WIDE ASSOCIATION STUDY (GWAS)

A statistical approach used to identify genomic variants that are statistically associated with a risk for a disease or a particular trait.



- T: risk allele, C: reference allele

- For **each SNP**, test:

$$H_0 : \beta = 0 \quad \text{vs} \quad H_1 : \beta \neq 0$$

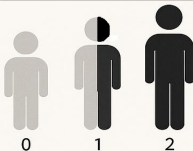
- Apply multiple-testing correction



# Genetic inheritance Models

## Additive Model

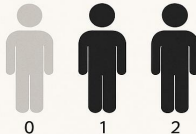
Each risk allele contributes additively to trait



Height increases incrementally with each risk allele

## Dominant Model

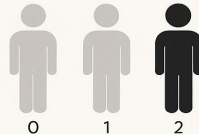
One copy of the risk allele is enough to affect the trait



Huntington's disease

## Recessive Model

Both alleles must be risk variants to influence the trait



Cystic fibrosis

Model	CC	CT	TT
Additive ( $G_A$ )	0	1	2
Dominant ( $G_D$ )	0	1	1
Recessive ( $G_R$ )	0	0	1
Genotypic ( $G_{\text{Het}}, G_{\text{Hom}}$ )	(0,0)	(1,0)	(0,1)

Model	Model specification	Null hypothesis
Additive	$g(\mathbb{E}[Y_i]) = \beta_0 + \beta_1 G_{A,i} + \boldsymbol{\gamma}^\top \mathbf{C}_i$	$H_0 : \beta_1 = 0$
Dominant	$g(\mathbb{E}[Y_i]) = \beta_0 + \beta_1 G_{D,i} + \boldsymbol{\gamma}^\top \mathbf{C}_i$	$H_0 : \beta_1 = 0$
Recessive	$g(\mathbb{E}[Y_i]) = \beta_0 + \beta_1 G_{R,i} + \boldsymbol{\gamma}^\top \mathbf{C}_i$	$H_0 : \beta_1 = 0$
Genotypic (2df)*	$g(\mathbb{E}[Y_i]) = \beta_0 + \beta_1 G_{\text{Het},i} + \beta_2 G_{\text{Hom},i} + \boldsymbol{\gamma}^\top \mathbf{C}_i$	$H_0 : \beta_1 = \beta_2 = 0$

$Y_i$ : phenotype;  $g(\cdot)$  identity (continuous) or logit (binary);  $\mathbf{C}_i$ : covariates.

*\*A model-free test that evaluates SNP–phenotype association without assuming a specific inheritance pattern.*

The assumed genetic model does not match the true mode of inheritance.

## Why does it occur?

- True inheritance patterns are unknown *a priori*.
- It is common to assume a fixed model (e.g., additive) for all SNPs.
- A dominant or recessive SNP modeled additively may lead to poor fit.

## Why does it matter?

- Loss of statistical power.
- Biased effect estimates.

---

Gaye, Amadou, and Sharon K. Davis. *Genetic model misspecification in genetic association studies*. BMC Research Notes 10.1 (2017): 569.



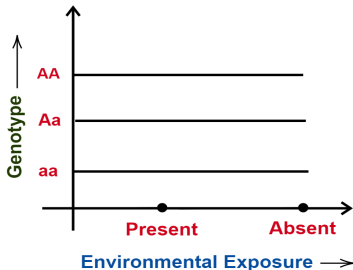
## GENE-ENVIRONMENT INTERACTIONS

---

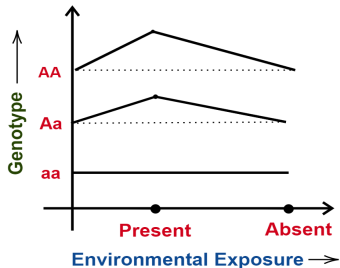
## Gene-Environment Interaction (Ottman,1996)<sup>1</sup>

"A different effect of an environmental exposure on disease risk in persons with different genotypes, and vice versa.

**Effects when GxE  
Interaction is Absent**

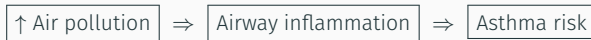


**Effects when GxE  
Interaction is Present**

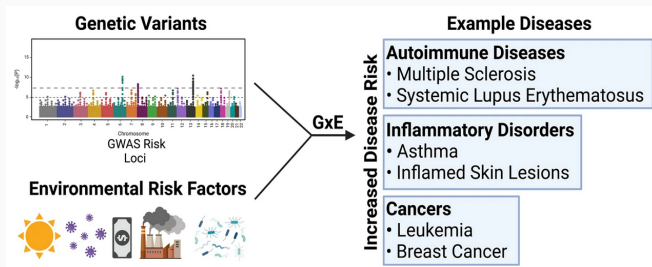


<sup>1</sup>Ottman, R. (1996). Gene-environment interaction: definitions and study design. Preventive medicine.

- Air pollutants (e.g.,  $\text{PM}_{2.5}$ ,  $\text{NO}_2$ , ozone) are known risk factors for asthma.



- The magnitude of pollution-induced asthma risk **differs by genotype**.



G×E model (Additive / Dominant / Recessive):

$$g(\mathbb{E}[Y_i]) = \beta_0 + \beta_1 G_i + \beta_2 E_i + \beta_3 (G_i \cdot E_i) + \boldsymbol{\gamma}^T \mathbf{C}_i$$

Where  $G_i \in \{G_{\text{add}}, G_{\text{dom}}, G_{\text{rec}}\}$

**Hypothesis:** G×E interaction (1 df)

$$H_0 : \beta_3 = 0 \quad \text{vs.} \quad H_A : \beta_3 \neq 0$$

- Interaction effect sizes are typically small.
- Statistical power is substantially lower than for main effects.
- Beyond these challenges, misspecifying the genetic model further penalizes  $G \times E$  tests, leading to additional loss of statistical power.

$$g(\mathbb{E}[Y_i]) = \beta_0 + \beta_1 G_{\text{Het},i} + \beta_2 G_{\text{Hom},i} + \beta_E E_i \\ + \beta_3 (G_{\text{Het},i} \cdot E_i) + \beta_4 (G_{\text{Hom},i} \cdot E_i) + \gamma^\top \mathbf{C}_i$$

**Hypothesis:** G×E interaction (2 df)

$$H_0 : \beta_3 = \beta_4 = 0 \quad \text{vs.} \quad H_A : \text{At least one of } \beta_3, \beta_4 \neq 0$$

---

<sup>2</sup>Moore, Camille M., Sean A. Jacobson, and Tasha E. Fingerlin. "Power and sample size calculations for genetic association studies in the presence of genetic model misspecification." *Human heredity* 84.6 (2020): 256-271.

$$g(\mathbb{E}[Y_i]) = \beta_0 + \beta_1 G_{\text{Het},i} + \beta_2 G_{\text{Hom},i} + \beta_E E_i \\ + \beta_3 (G_{\text{Het},i} \cdot E_i) + \beta_4 (G_{\text{Hom},i} \cdot E_i) + \gamma^\top \mathbf{C}_i$$

**Hypothesis:** G×E interaction (2 df)

$$H_0 : \beta_3 = \beta_4 = 0 \quad \text{vs.} \quad H_A : \text{At least one of } \beta_3, \beta_4 \neq 0$$

**Strengths:**

- Robust to model misspecification.
- Higher power when the true model is recessive, overdominant.
- Allows genotype-specific environmental effects.

**Drawbacks:**

- Lower power under true additive/ dominant inheritance.
- Requires larger sample sizes.

---

<sup>2</sup>Moore, Camille M., Sean A. Jacobson, and Tasha E. Fingerlin. "Power and sample size calculations for genetic association studies in the presence of genetic model misspecification." *Human heredity* 84.6 (2020): 256-271.

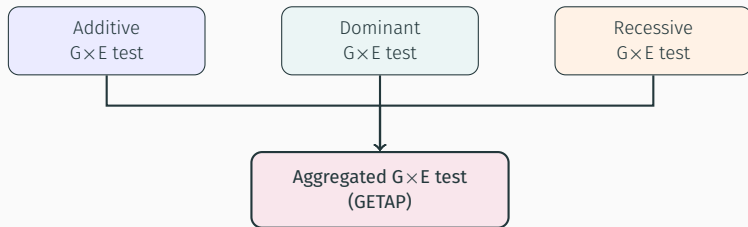
- G×E analyses typically assume a **single genetic model**, and **misspecification** leads to substantial power loss.
- The 2df test is **robust**, but **inefficient** under true additive or dominant models.
- Testing multiple single models separately introduces **multiple-testing burden**.



## OUR SOLUTION: AGGREGATE EVIDENCE ACROSS GENETIC MODELS

- Instead of selecting a single genetic model, we test  $G \times E$  under multiple models
- Evidence across models is combined into a single omnibus test

*This avoids committing to a possibly misspecified genetic model.*



Aggregation performed using ACAT or HMP

## P-VALUE AGGREGATION

---

Combining multiple p-values into a single test of a **global null hypothesis**.

Why is aggregation useful?

- **Power:** Detect weak but consistent signals across tests.
- **Robustness:** Protect against model misspecification.
- **Multiplicity:** Avoid repeated multiple-testing corrections.

Why is this important in genomics?

- Tests are often **dependent** (e.g., LD, correlated models).
- Standard combination methods may fail under dependence.

# THE AGGREGATED CAUCHY ASSOCIATION TEST (ACAT)<sup>3</sup>

Consider  $k$  hypothesis tests with p-values  $p_1, p_2, \dots, p_k$ , where  $p_i$  is from the  $i$ th test.

## Step 1: Transform individual P-values

$$C_i = \tan [(0.5 - p_i)\pi]$$

## Step 2: ACAT test statistic

$$T_{\text{ACAT}} = \sum_{i=1}^k w_i C_i = \sum_{i=1}^k w_i \tan [(0.5 - p_i)\pi]$$

Here,  $w_i \geq 0$  are user-specified weights (uniform weights  $w_i = 1/k$  used in our analysis).

## Step 3: Combined P-value (Cauchy tail approximation)

$$p_{\text{ACAT}} \approx 1 - \frac{1}{\pi} \arctan \left( \frac{T_{\text{ACAT}}}{\bar{W}} \right) \quad \text{where} \quad \bar{W} = \sum_{i=1}^k w_i$$

- *Very fast and analytically tractable, while remaining robust to **arbitrary dependency**.*

---

<sup>3</sup>Liu, Yaowu, et al. "ACAT: a fast and powerful p value combination method for rare-variant analysis in sequencing studies." The American Journal of Human Genetics 104.3 (2019): 410-421.

## Step 1: HMP Statistic

$$\overset{\circ}{p} = \frac{\sum_{i=1}^k w_i}{\sum_{i=1}^k \frac{w_i}{p_i}}, \quad \left( \sum_{i=1}^k w_i = 1 \right)$$

For equal weights:  $\overset{\circ}{p} = \frac{k}{\sum_{i=1}^k \frac{1}{p_i}}$

## Step 2: Final P-value (Two Methods)

- **Method A:** Compare  $\overset{\circ}{p}$  to critical value  $\alpha_k$  (approximate).
- **Method B:** Asymptotically exact p-value:

$$p_{\overset{\circ}{p}} = \int_{1/\overset{\circ}{p}}^{\infty} f_{\text{Landau}}(x \mid \mu, \sigma) dx$$

where  $\mu \approx \log(k) + 0.874$ ,  $\sigma = \pi/2$

---

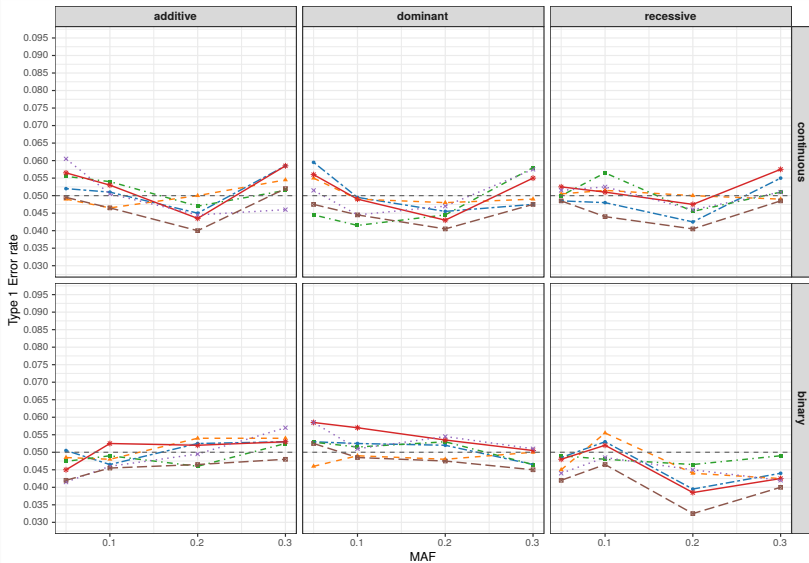
<sup>4</sup>Wilson, Daniel J. "The harmonic mean p-value for combining dependent tests." Proceedings of the National Academy of Sciences 116.4 (2019): 1195-1200.

## SIMULATION RESULTS

---

- **True genetic model:** Additive, dominant, recessive
- **Trait type:** Continuous and binary
- **Environmental exposure:** Continuous and binary
- **Sample size:**  $n = 10,000$
- **Replicates:** 2,000 simulations
- **Minor allele frequency (MAF):** 0.05, 0.10, 0.20, 0.30
- **Significance level:**  $\alpha = 0.05$

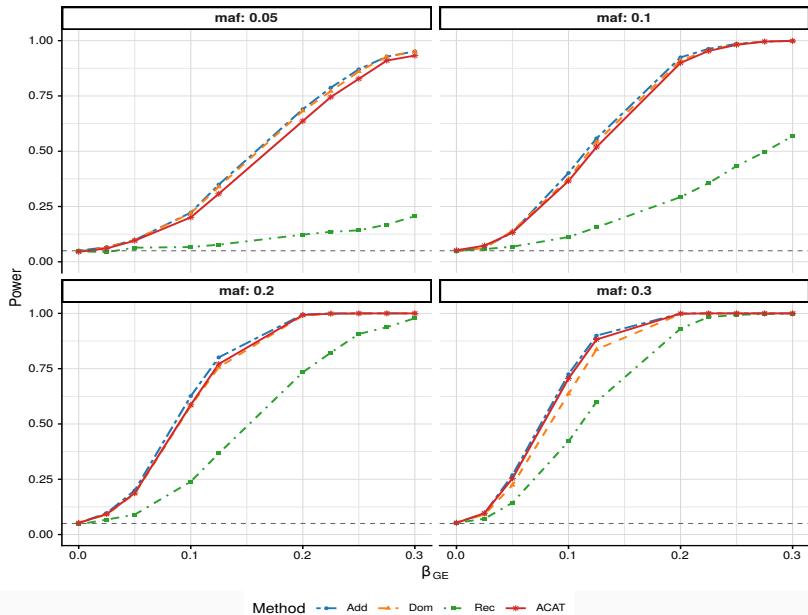
Estimated T1ER for continuous phenotype (n = 10000)

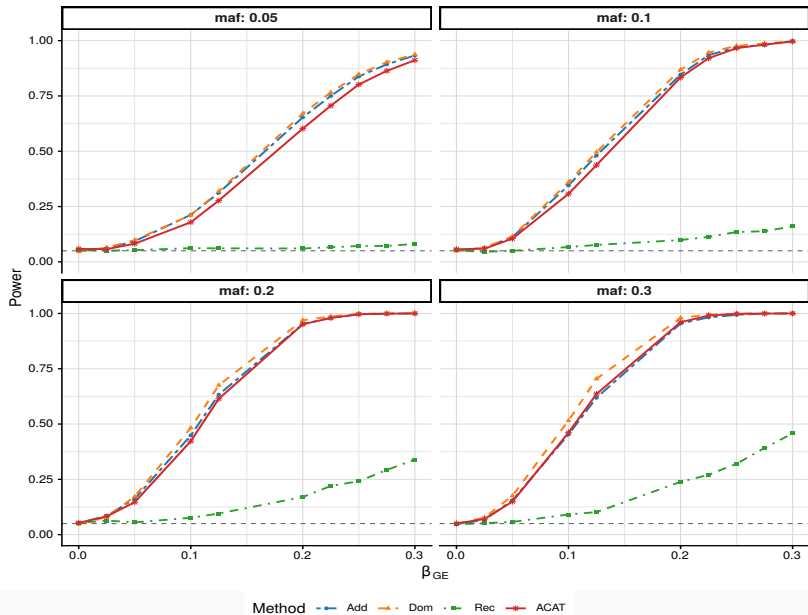


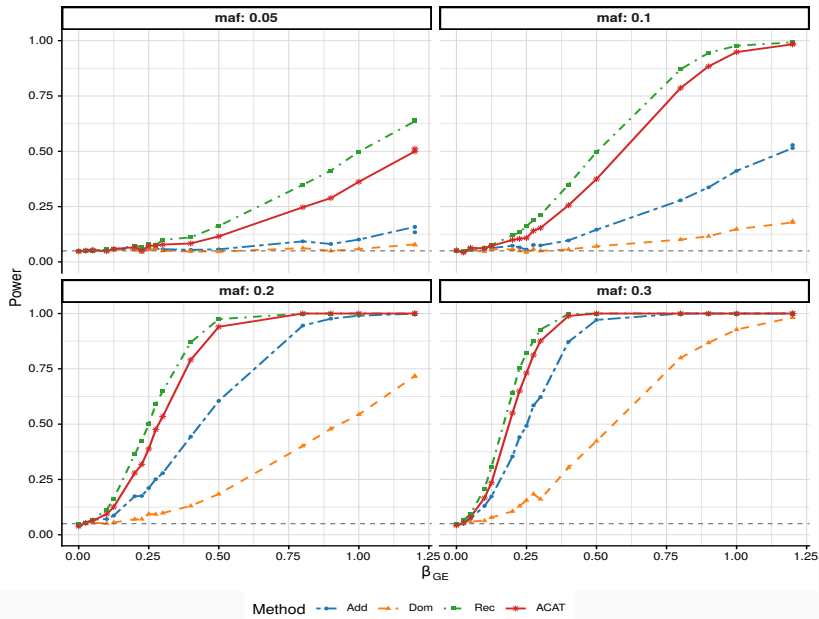
Method

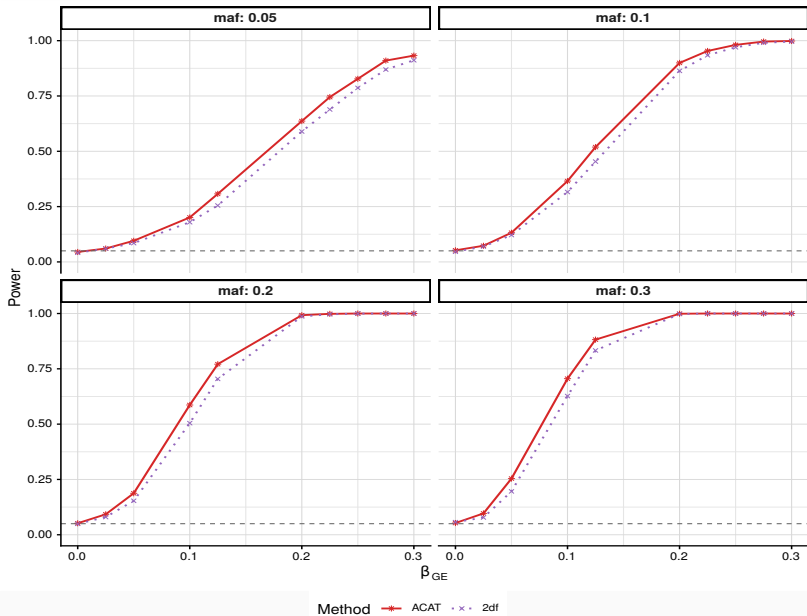
- Add
- Dom
- Rec
- ACAT
- 2df
- HMP

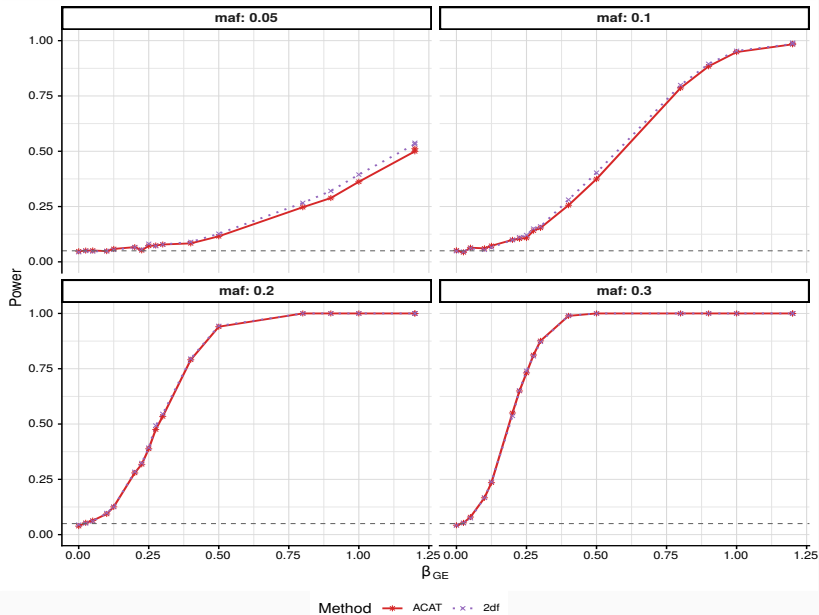


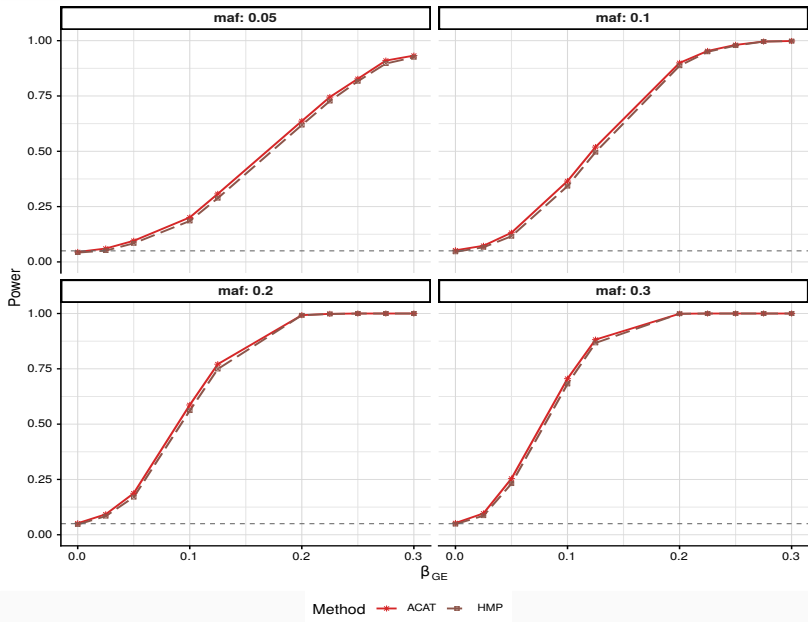












## REAL DATA APPLICATION

---

## UK Biobank

- ~500,000 participants
- Unrelated White British individuals
- Autosomal SNPs after QC (~600K)

## Phenotypes and environments

- Continuous and binary traits
- Lifestyle and behavioral exposures

Phenotype + Genotype + Environment  
data



Additive / Dominant / Recessive  
G×E tests



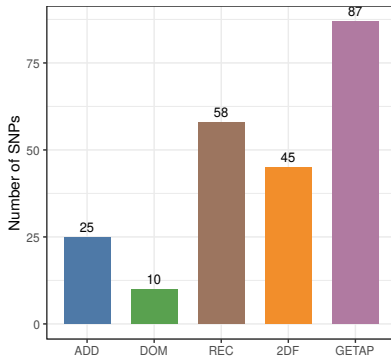
GETAP (ACAT)  
P-value aggregation



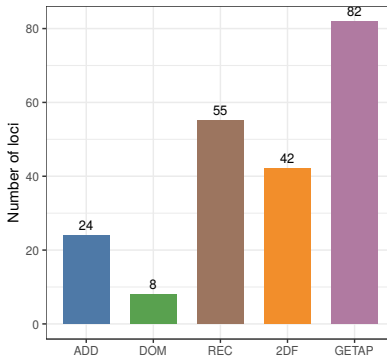
Significant G×E loci



**G×E Discoveries (FDR < 0.05)**

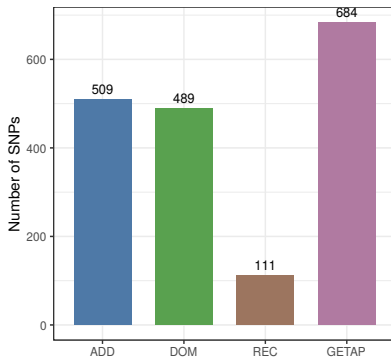


**LD-independent loci**

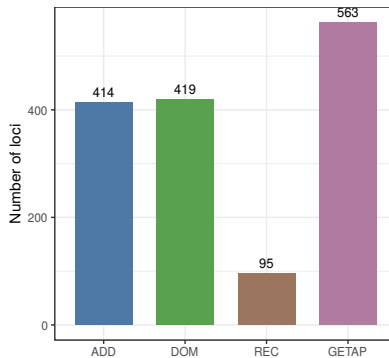


# PHENOTYPE: TYPE 2 DIABETES, ENV.: SLEEP DURATION

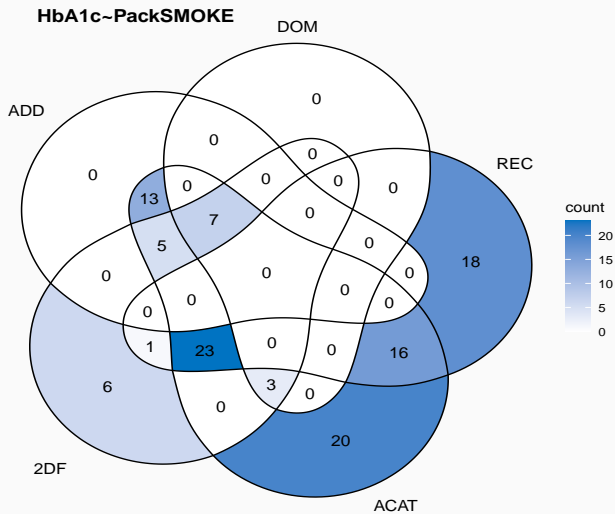
**G×E Discoveries (FDR < 0.05)**



**LD-independent loci**



# OVERLAP STRUCTURE

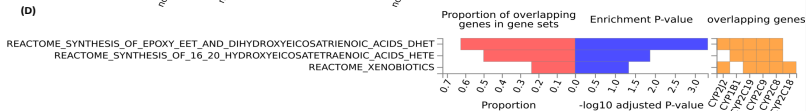
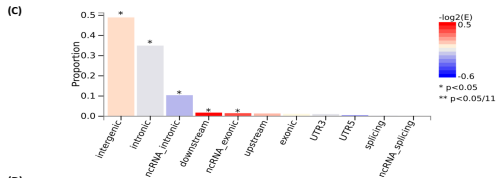
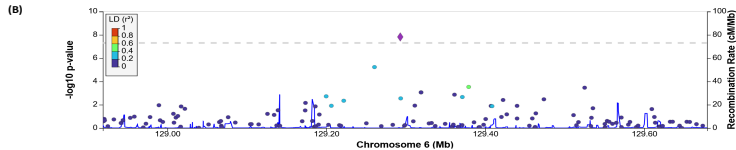
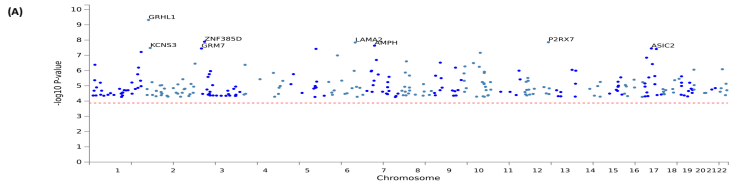


**Table 1:** Selected SNPs showing genome-wide significant G×E signals.

Phenotype	Environment	SNP	Chr	$P_{\text{ADD}}$	$P_{\text{DOM}}$	$P_{\text{REC}}$	$P_{\text{GETAP}}$	$P_{2\text{DF}}$
HbA1c	Pack-years (Smoking)	rs407423	8	$2.7 \times 10^{-2}$	$1.7 \times 10^{-1}$	$5.1 \times 10^{-10}$	$1.5 \times 10^{-9}$	$3.7 \times 10^{-9}$
FEV <sub>1</sub> /FVC	Pack-years (Smoking)	rs13180	15	$1.0 \times 10^{-8}$	$2.9 \times 10^{-6}$	$2.8 \times 10^{-6}$	$3.0 \times 10^{-8}$	$5.8 \times 10^{-8}$
T2D	Sleep duration	rs2801198	1	$1.4 \times 10^{-3}$	$6.7 \times 10^{-2}$	$3.0 \times 10^{-8}$	$8.9 \times 10^{-8}$	–

## BIOLOGICAL SIGNIFICANCE: T2D

---



- Across scenarios, the aggregated test behaves close to the best-performing model without knowing that model in advance.
- GETAP provides robust and scalable inference under genetic model uncertainty.
- P-value aggregation recovers  $G \times E$  loci missed by single-model tests.

## Previous Work

- Mishra, S. and Majumdar, A., 2025. A Multi-Phenotype Approach to Joint Testing of Main Genetic and Gene-Environment Interaction Effects. *Statistics in Medicine*, 44(20-22), p.e70253.

- **R package:** *MvGGE* (implements the above method).

Available open source on GitHub: <https://github.com/SauMStats/MvGGE>



## WAY FORWARD



- **Ancestry-specific  $G \times E$  analysis:**







Apply our multivariate  $G \times E$  framework (MvGGE) to UK Biobank South Asian ancestry data to study population-specific interaction effects.

- **Multiple-environment  $G \times E$  methods:**

Extend current models to jointly incorporate multiple environmental exposures, moving beyond single-environment interaction analyses.

- **Software dissemination:**

Publish the existing *MvGGE* R package on Bioconductor and expand it into a scalable, user-friendly tool for large-scale multivariate  $G \times E$  analysis.

-  Zeng, Ping, et al. "Aggregating multiple expression prediction models improves the power of transcriptome-wide association studies." *Human Molecular Genetics* 30.10 (2021): 939-951.
-  Gaye, Amadou, and Sharon K. Davis. "Genetic model misspecification in genetic association studies." *BMC research notes* 10.1 (2017): 569.
-  Moore, Camille M., Sean A. Jacobson, and Tasha E. Fingerlin. "Power and sample size calculations for genetic association studies in the presence of genetic model misspecification." *Human heredity* 84.6 (2020): 256-271.
-  Haas, Cameron B., et al. "Interactions between folate intake and genetic predictors of gene expression levels associated with colorectal cancer risk." *Scientific Reports* 12.1 (2022): 18852.
-  Evans, Luke M., et al. "Transcriptome-wide gene-gene interaction associations elucidate pathways and functional enrichment of complex traits." *PLoS Genetics* 19.5 (2023): e1010693.
-  Gao, Guimin, et al. "A joint transcriptome-wide association study across multiple tissues identifies candidate breast cancer susceptibility genes." *The American Journal of Human Genetics* 110.6 (2023): 950-962.

*THANK YOU*

## BACKUP SLIDES

---

# WALD AND LIKELIHOOD RATIO TESTS (2 DF)

## 1. Wald Test:

$$W = \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix}^\top \left[ \widehat{\text{Var}} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{pmatrix} \right]^{-1} \begin{bmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \end{bmatrix} \sim \chi^2_2$$

- Uses estimated coefficients and their variance-covariance matrix.

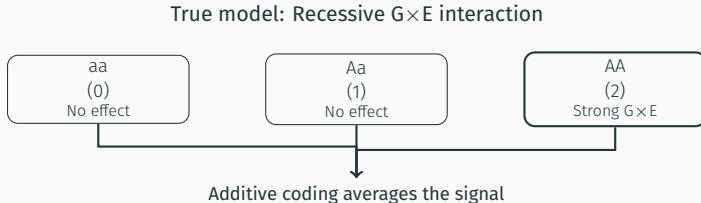
## 2. Likelihood Ratio Test (LRT):

$$\Lambda = -2 \left[ \ell(\hat{\theta}_0) - \ell(\hat{\theta}_1) \right] \sim \chi^2_2$$

- $\ell(\hat{\theta}_0)$ : log-likelihood under null model (no genotype effect)
- $\ell(\hat{\theta}_1)$ : log-likelihood under full model (with  $\beta_1, \beta_2$ )
- Tests improvement in model fit when including genotype indicators

# GENETIC MODEL MISSPECIFICATION CAUSES POWER LOSS

- Genetic effects may be additive, dominant, or recessive
- Incorrect genotype coding dilutes interaction signals
- This problem is amplified for rare variants and  $G \times E$  effects



## WHY BH FDR IN GxE ANALYSIS?

- Balances power and error control in high-dimensional genomic data
- Controls expected proportion of false positives ( $q < 0.05$ ), unlike conservative FWER methods
- Maintains robust control under moderate LD dependence typical in post-QC SNPs
- Standard practice in large biobank GxE studies (e.g., UK Biobank)
- Avoids Bonferroni's type II error inflation and BY's unnecessary power loss



## COMPARISON OF MULTIPLE TESTING CORRECTIONS

Multiple testing is an analysis in which multiple independent hypotheses are tested. The overall combined probability of making a type I error will increase.

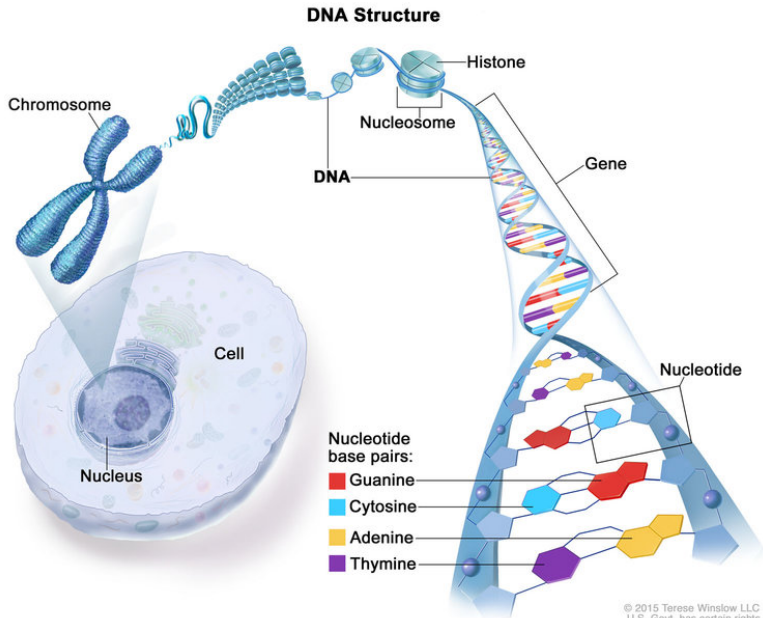
Method	Controls	Suitable for	Stringency	Ideal Usage Scenario
Bonferroni	FWER	Small/independent tests	Very conservative	High-stakes studies where false positives are critical
BH	FDR	Large/high-dimensional data	Less conservative	Exploratory studies with large test numbers (e.g., genomics)
BY	FDR under dependence	Correlated/multiple related tests	Moderate	Studies with known dependencies among tests (e.g., pathway analyses)

The Benjamini-Hochberg (BH) procedure ranks the  $p$ -values  $p_{(1)}, p_{(2)}, \dots, p_{(m)}$  in ascending order and finds the largest  $k$  such that:

$$p_{(k)} \leq \frac{k}{m} \cdot \alpha$$

where  $m$  is the total number of tests and  $\alpha$  is the desired false discovery rate. All  $p$ -values  $p_{(1)}, \dots, p_{(k)}$  are considered significant.

# CELL → NUCLEUS → CHROMOSOME → DNA



- **SNP:** Single nucleotide polymorphism (SNP, pronounced "snip") is a genomic variant at a single base position in the DNA.
- **Allele:** One of two versions of DNA sequence at a given genomic location. Example: C, T.
- **Genotype:** The overall genetic makeup of an individual.
  - The 3 genotypes for alleles C and T will be CC, CT, and TT.
  - Genotype values vary from one person to another.
- **Phenotype:** A characteristic of an individual which can be observed/ measured.

