# GIFT: Conditional TWAS for fine-mapping candidate causal genes

**Speaker: Yuekai Li**

**Major: Biostatistics**

# CONTENTS

1. Background

2. Challenges

3. Methods

4. Results

# 01
# Background

# ■ Outline

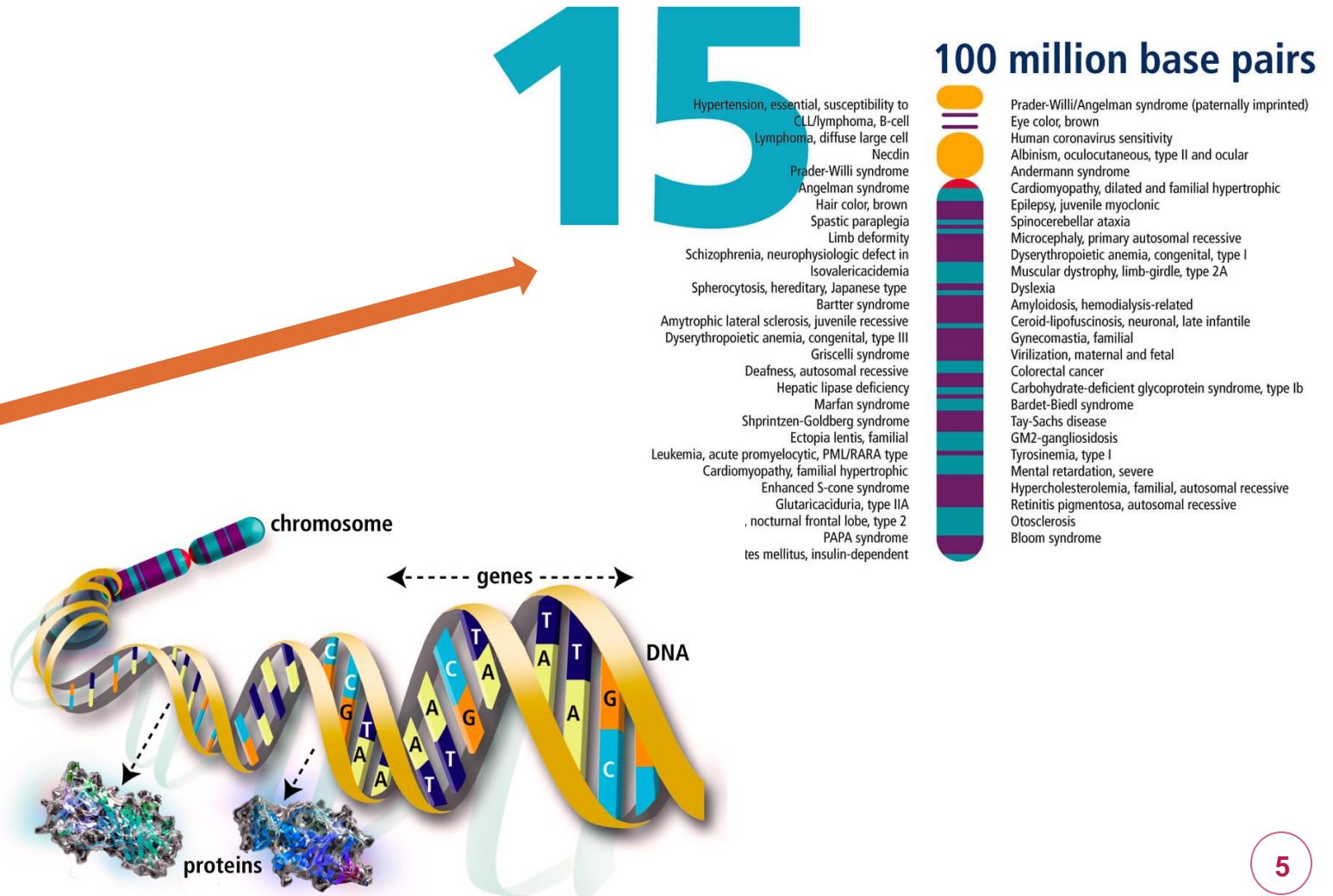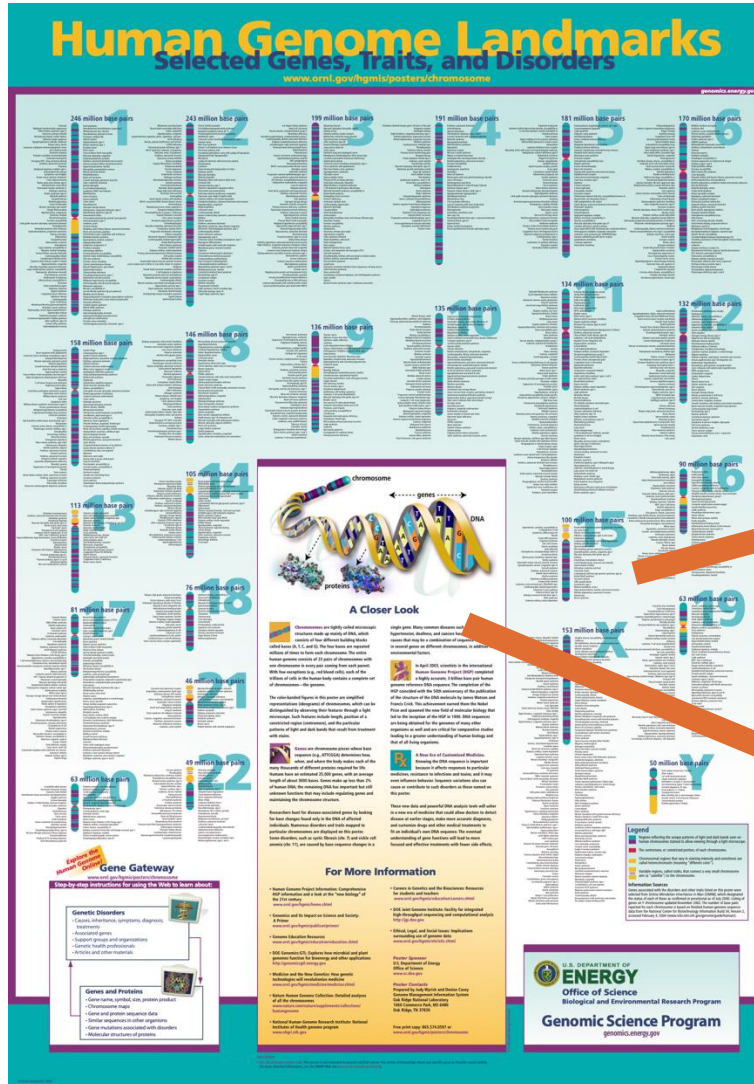**01**     The relationship between traits, genes, SNPs

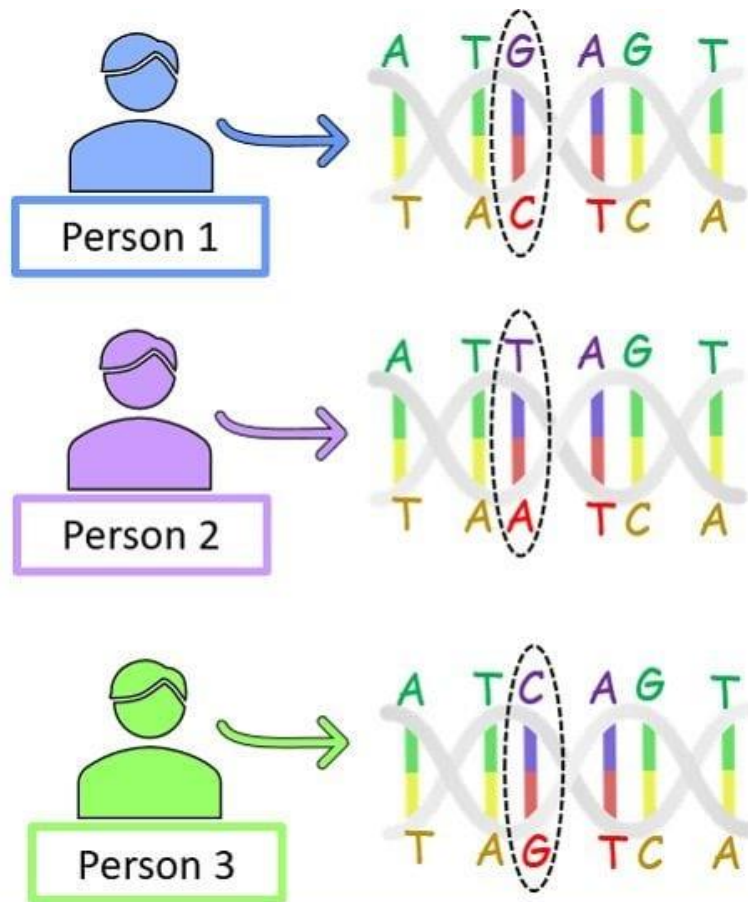**02**     What is GWAS

**03**     The motivation of TWAS

## ■ The relationship between traits, genes, SNPs

# ■ The relationship between traits, genes, SNPs



Single Nucleotide Polymorphism

- SNP is the replacement of a single base pair in the DNA sequence.

- SNP is the most common type of genetic variation.

- More than 600 million SNPs have been identified across the human genome in the world's population.

6

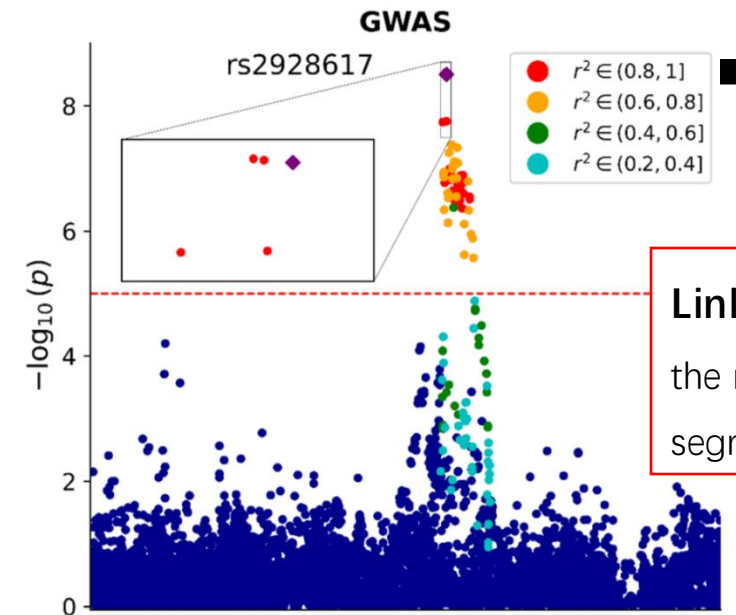# What is GWAS



**Linkage disequilibrium (LD)** : the non-independent segregation of genetic variants.

## ■ What is GWAS



**Genotype data:**



| | SNP1 | SNP2 | SNP3 | SNP4 |
|---|---|---|---|---|
| Individual 1 | AT | CG | TT | CC |
| Individual 2 | TA | GG | GT | CA |
| Individual 3 | TT | CC | GT | CA |
| Individual 4 | TT | CC | GG | AA |

Major=2     Heterozygous=1     Minor=0

For example, if we assume A is the major allele, then A:A=2, A:C/C:A=1, CC=0

8

# ■ What is GWAS

Linear regression models for GWAS can be

written as follows:

Fixed effect

$$Y \sim W\alpha + X_s \boxed{\beta_s} + g + e \tag{1}$$

Its p-value measures the strength of the

$$g \sim N(0, \sigma_A^2 \psi) \tag{2}$$

association between SNPs and trait.

$$e \sim N(0, \sigma_e^2 I) \tag{3}$$

$Y$ : the phenotype value

$W$ : the vector of covariates including an intercept term

$\alpha$ : the corresponding vector of effect

$X_s$ : the genotype value for the genetic variant $s$

$\beta_s$ : the corresponding fixed effect

$g$ : the random effect that captures the polygenic effect of other SNPs

$e$ : the random effect of residual errors

$\sigma_A^2$ : the additive genetic variation of the phenotype

$\psi$ : the standard genetic relationship matrix

$\sigma_e^2$ : residual variance

# ■ The motivation of TWAS



- SNPs in different tissues have different regulatory effects on gene expression.



- Most SNPs (about 98.5%) are located in the non-coding regions of DNA.

# ■ The motivation of TWAS

- Most SNPs (about 98.5%) are located in the non-coding regions of DNA.

# ■ The motivation of TWAS

**Genotype data**

| ID | $SNP_1$ | $SNP_2$ | ... | $SNP_p$ |
|---|---|---|---|---|
| $id_1$ | 2 | 0 | ... | 0 |
| $id_2$ | 0 | 1 | ... | 2 |
| $id_3$ | 2 | 1 | ... | 0 |
| ... | ... | ... | ... | ... |
| $id_{n_2}$ | 1 | 0 | ... | 0 |

**GWAS data** ($\tilde{Z}$)

**Expression data**

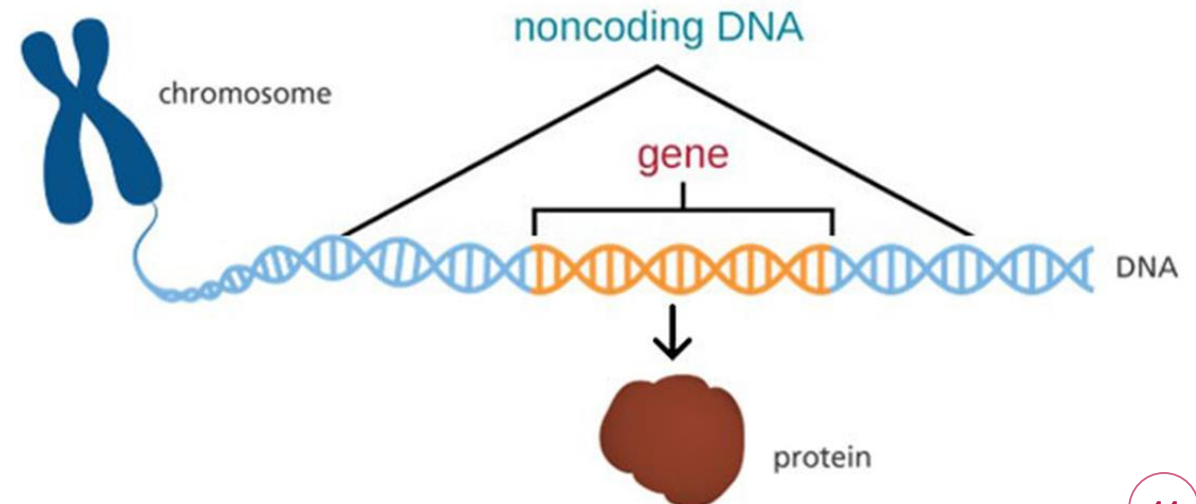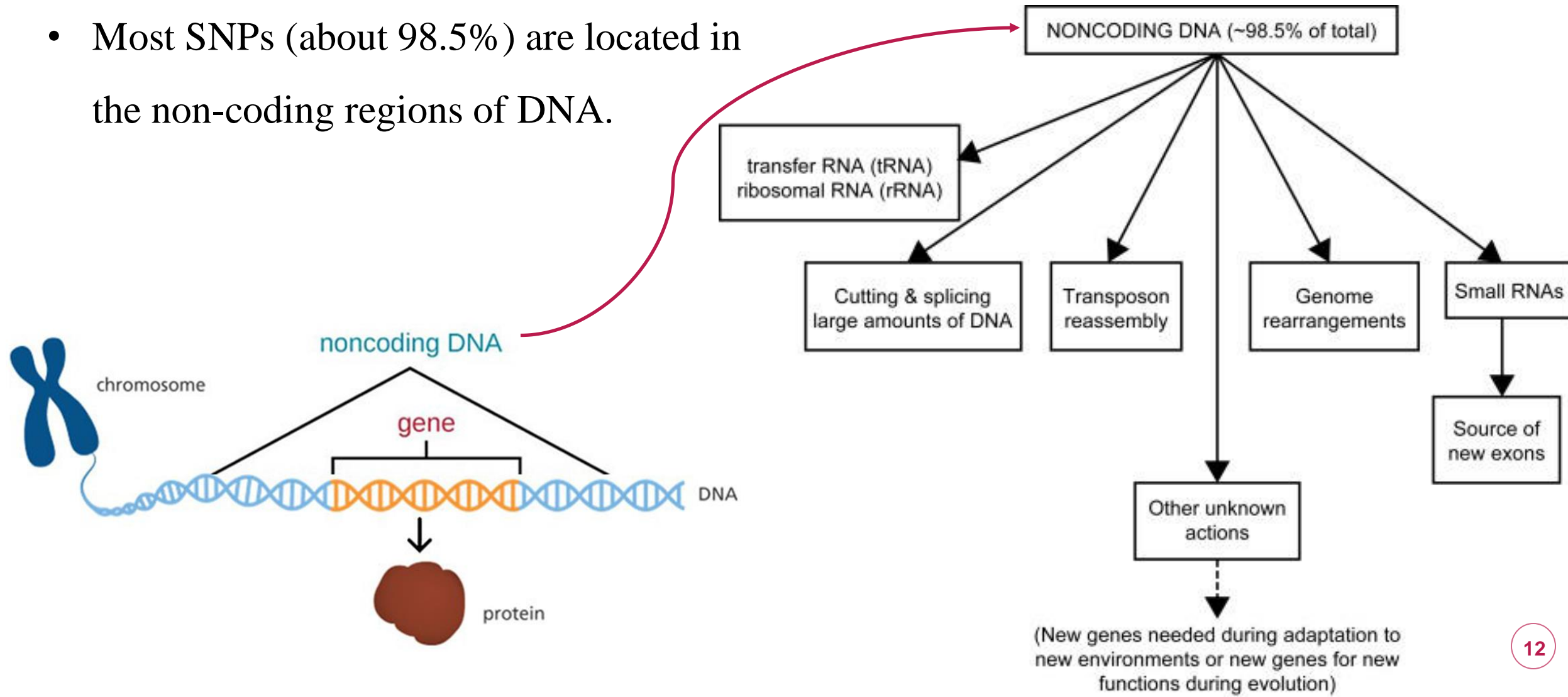| ID | $\hat{x}_1$ | $\hat{x}_2$ | ... | $\hat{x}_m$ |
|---|---|---|---|---|
| $id_1$ | | | ... | |
| $id_2$ | | | ... | |
| $id_3$ | | | ... | |
| ... | ... | ... | ... | ... |
| $id_{n_2}$ | | | ... | |

$\hat{x}_i = \tilde{Z}\beta_i$

② **Impute**

**?**

**Phenotype data**

| ID | Trait |
|---|---|
| $id_1$ | 1.23 |
| $id_2$ | 4.56 |
| $id_3$ | 7.89 |
| ... | ... |
| $id_{n_2}$ | 2.33 |

③ **Associate**

$y = \alpha_i \hat{x}_i + \tilde{e}$

$$\beta = [\beta_1, \beta_2, ..., \beta_m]$$

① **Train**

$x_i = Z\beta_i + e_i$

**Reference panel** ($Z$)

| ID | $SNP_1$ | $SNP_2$ | ... | $SNP_p$ |
|---|---|---|---|---|
| $id_1$ | 0 | 0 | ... | 0 |
| $id_2$ | 1 | 2 | ... | 1 |
| ... | ... | ... | ... | ... |
| $id_{n_1}$ | 1 | 0 | ... | 1 |

| ID | $x_1$ | $x_2$ | ... | $x_m$ |
|---|---|---|---|---|
| $id_1$ | 0.1 | 0.5 | ... | 1.3 |
| $id_2$ | 1.2 | 2.2 | ... | 0.1 |
| ... | ... | ... | ... | ... |
| $id_{n_1}$ | 0.2 | 0.1 | ... | 1.0 |

# ■ The motivation of TWAS



**1. Training stage**: Estimate regulatory effect sizes of multiple SNPs on the gene expression level from a small reference panel with genotype and expression data.

**2. Imputation stage**: Obtain the predicted gene expression of GWAS individuals.

**3. Association stage**: Implement hypothesis tests between predicted gene expression and the target trait

# 02

# Challenges

# ■ Outline

☑ **LD and expression correlation lead to confounding**

☑ **Two-Step Inference Procedure lead to power loss**

# ■ LD and expression correlation lead to confounding



**Linkage disequilibrium (LD)** :

the non-independent

segregation of genetic variants.

Boxplot displays the maximum of the absolute value of **expression correlation**

estimates in each region across the 22 chromosomes

# ■ LD and expression correlation lead to confounding

There are no unobserved exposures.

Such situations rarely occur.

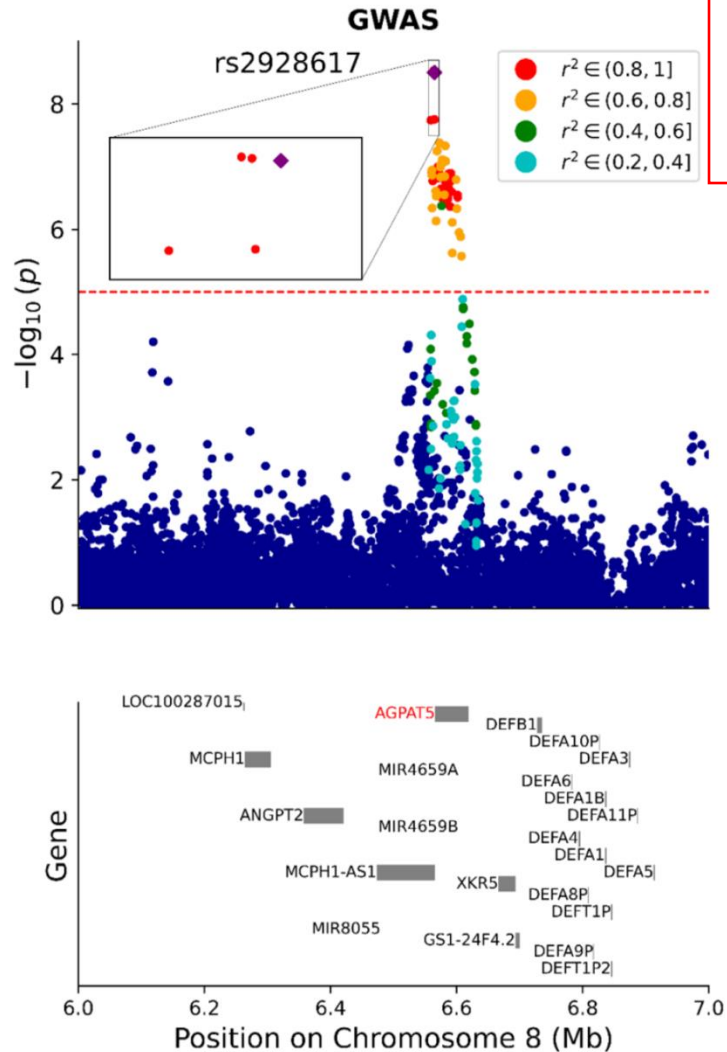**Confounding**

**Pleiotropy:**

A genetic variant affects the outcome through a pathway that does not involve the risk factor of interest.



**a** Causality with SNPs in LD

Two causal eQTL SNPs in LD → Exposure (gene expression)

SNP A   SNP B

Two correlated SNPs affect the exposure and outcome, no pleiotropy

SNP A → Exposure → Outcome
SNP B        $b_E$

**b** Causality with pleiotropy through LD

Unobserved exposure     Observed exposure

SNP B   SNP A

Outcome is affected through two pathways from two correlated SNPs

SNP A → Observed exposure → $b_E$ Outcome
SNP B → Unobserved exposure → $b_U$

**c** Causality with pleiotropy through overlap

Unobserved exposure     Observed exposure

SNP

Outcome is affected through two pathways from the same SNP

SNP → Observed exposure → $b_E$ Outcome
SNP → Unobserved exposure → $b_U$

18

# ■ Two-Step Inference Procedure lead to power loss



Most of the existing marginal TWAS methods consist of two separate analytical steps.

The characteristic of these methods is that they estimate the weights $\boldsymbol{\beta}$ from the reference panel in advance. TWAS (or TWAS fine-mapping) is performed given weights $\boldsymbol{\beta}$.

■ **Two-Step Inference Procedure lead to power loss**

The point estimation of $\boldsymbol{\beta}_i$ has more uncertainty.

$$\mathbf{x}_i = \mathbf{Z}\boldsymbol{\beta}_i + e_i$$

Plug in

$$\boldsymbol{y} = \alpha_i \hat{\mathbf{x}}_i + \tilde{e}$$

**Two-Step Inference**

Most of the existing marginal TWAS methods consist of two separate analytical steps.

The characteristic of these methods is that they estimate the weights $\boldsymbol{\beta}$ from the reference panel in advance. TWAS (or TWAS fine-mapping) is performed given weights $\boldsymbol{\beta}$.

20

**03**

**Methods**

# ■ Outline

**01**　　TWAS method timeline

**02**　　The model of GIFT

**03**　　The advantages of GIFT

# ■ TWAS method timeline

### PredictDB & FUSION

PredictDB and FUSION are commonly-used models for calculating eQTL weights.

### MA-FOCUS

MA-FOCUS extends the FOCUS model to the case of **multi-ancestry**.

**2019**

**2024**

**2016**

**2022**

FOCUS is the first to perform fine-mapping of **multiple genes** in a genomic region based on the probabilistic model.

GIFT jointly models the entire process of TWAS and performs **conditional TWAS** on all genes in a genomic region.

### FOCUS

### GIFT

## ■ The model of GIFT (individual-level)

GIFT (Gene-based Integrative Fine-mapping through conditional TWAS),

**jointly models** all $k$ genes residing in the focal region and carries out TWAS conditional analysis:

$$\begin{cases} \mathbf{x}_i = Z_i \boldsymbol{\beta}_i + \mathbf{e}_i, i = 1, \cdots, k & (1) \\ \mathbf{y} = \sum_{i=1}^{k} \alpha_i (\tilde{Z}_i \boldsymbol{\beta}_i) + \tilde{\mathbf{e}} & (2) \end{cases}$$

$\mathbf{x}_i \in \mathbb{R}^{n_1 \times 1}$: expression vector for the $i$-th gene

$\mathbf{y} \in \mathbb{R}^{n_2 \times 1}$ : phenotype vector

$\mathbf{Z}_i \in \mathbb{R}^{n_1 \times p_i}$ : genotype matrix in the reference panel for the $i$-th gene

$\widetilde{\mathbf{Z}}_i \in \mathbb{R}^{n_2 \times p_i}$: genotype matrix in the GWAS data for the $i$-th gene

$\boldsymbol{\beta}_i \in \mathbb{R}^{p_i \times 1}$ : eQTL random effects on the $i$-th gene expression

$\alpha_i \in \mathbb{R}$ : effects of predicted expression for the $i$-th gene

Assumed that $\mathbf{y}$, $\mathbf{x}_i$ and each column of $\mathbf{Z}_i$ and $\widetilde{\mathbf{Z}}_i$ have all been standardized to have a **mean of zero** and standard **deviation of 1**.

24

## ■ The model of GIFT (individual-level)

GIFT (Gene-based Integrative Fine-mapping through conditional TWAS),

**jointly models** all $k$ genes residing in the focal region and carries out TWAS conditional analysis:

$$\begin{cases} \mathbf{x}_i = Z_i \boldsymbol{\beta}_i + \mathbf{e}_i, i = 1, \cdots, k & (1) \\ \mathbf{y} = \sum_{i=1}^{k} \alpha_i (\tilde{Z}_i \boldsymbol{\beta}_i) + \tilde{\mathbf{e}} & (2) \end{cases}$$

Due to $p_i > n_1$

$\mathbf{x}_i \in \mathbb{R}^{n_1 \times 1}$: expression vector for the $i$-th gene

$\mathbf{y} \in \mathbb{R}^{n_2 \times 1}$ : phenotype vector

$\boldsymbol{\beta}_i \sim N(\mathbf{0}, \sigma_{\beta_i}^2 \cdot \boldsymbol{I}_{p_i})$

$\mathbf{Z}_i \in \mathbb{R}^{n_1 \times p_i}$ : genotype matrix in the reference panel for the $i$-th gene

Using the posterior distribution of $\boldsymbol{\beta}_i$ for eQTL effect

$\tilde{\mathbf{Z}}_i \in \mathbb{R}^{n_2 \times p_i}$: genotype matrix in the GWAS data for the $i$-th gene

instead of the point estimate $\hat{\boldsymbol{\beta}}_i$

$\boldsymbol{\beta}_i \in \mathbb{R}^{p_i \times 1}$ : eQTL random effects on the $i$-th gene expression

$\alpha_i \in \mathbb{R}$ : effects of predicted expression for the $i$-th gene

26

# ■ The model of GIFT (individual-level)

GIFT (Gene-based Integrative Fine-mapping through conditional TWAS),

**jointly models** all *k* genes residing in the focal region and carries out TWAS conditional analysis:

$$\begin{cases} \mathbf{x}_i = Z_i\boldsymbol{\beta}_i + \mathbf{e}_i, i = 1, \cdots, k \qquad (1) \\ \mathbf{y} = \sum_{i=1}^{k} \alpha_i(\tilde{Z}_i\boldsymbol{\beta}_i) + \tilde{\mathbf{e}} \qquad\qquad (2) \end{cases}$$

$\mathbf{e}_i \in \mathbb{R}^{n_1 \times 1}$: residual errors for the *i*-th gene,

where $\left(\mathbf{e}_{l,1}, \mathbf{e}_{l,2}, \ldots, \mathbf{e}_{l,k}\right)^T \sim N_k(0, \Omega)$ for the same individual *l*

$\tilde{\mathbf{e}} \in \mathbb{R}^{n_2 \times 1}$ : residual error with each element *i.i.d.* from the same normal distribution $N(0, \sigma_y^2)$

GIFT takes the correlation of gene expressions into account.

27

## ■ The model of GIFT (summary-level)

GIFT can also be extended to perform inference using summary statistics only.

The corresponding model for summary statistics are:

$$\begin{cases} \widehat{\boldsymbol{\beta}}^*_{\boldsymbol{x}_i} = \boldsymbol{\Sigma}_{1i}\boldsymbol{\beta}_i + \boldsymbol{e}_{\boldsymbol{x}_i}, i = 1, \cdots, k \\ \widehat{\boldsymbol{\beta}}^*_{\boldsymbol{y}} = \boldsymbol{\Sigma}_2\big(\alpha_1\boldsymbol{\beta}_1^T, \cdots, \alpha_k\boldsymbol{\beta}_k^T\big)^T + \boldsymbol{e}_{\boldsymbol{y}} \end{cases}$$

$\widehat{\boldsymbol{\beta}}^*_{\boldsymbol{x}_i} \in \mathbb{R}^{1 \times p_i}$: the estimates for the marginal SNP effects on the $i$-th gene expression

$\widehat{\boldsymbol{\beta}}^*_{\boldsymbol{y}} \in \mathbb{R}^{1 \times p}$: the estimates for the marginal SNP effects on the trait

$\boldsymbol{\Sigma}_{1i} \in \mathbb{R}^{p_i \times p_i}$: correlation matrix of all cis-SNPs for the $i$-th gene in the <mark>reference panel</mark>

$\boldsymbol{\Sigma}_2 \in \mathbb{R}^{p \times p}$: correlation matrix of all cis-SNPs for all the genes in the focal region in the <mark>GWAS data</mark>.

$\boldsymbol{\beta}_i \in \mathbb{R}^{p_i \times 1}$ : eQTL effects on the $i$-th gene experssion

$\alpha_i \in \mathbb{R}$ : effects of predicted expression for the $i$-th gene

# ■ The advantages of GIFT

## 01

### Conditional TWAS analysis

GIFT performs TWAS fine-mapping conditional on the effects of the other genes to avoid confounding.

## 02

### Joint likelihood inference framework

The joint inference framework accounts for the uncertainty in the SNP effect-size estimates on gene expression and the uncertainty in the predicted expression.

## 03

### PX-EM algorithm

GIFT introduces the auxiliary parameter $\lambda$ through the parameter expansion method to significantly improve the convergence speed.

**04**

**Results**

■ **Outline**

☑ **Data input**

☑ **GIFT produces calibrated P values under the null simulations**

☑ **GIFT is powerful under a range of alternative simulations**

☑ **Real-data applications**

# ■ Data input

# ▪ GIFT produces calibrated P values under the null simulations



**a** All genes in the region are null

**b** One gene in the region is casual

**c**

**a, b**: Quantile–quantile plots of −log10(P values) from the three frequentist methods,

which are both displayed for the **non-causal genes**.

**c**: Boxplot from FOCUS displays the PIPs from causal genes and non-causal genes.

*All of the simulation results are from 1,000 random region.

# ■ GIFT is powerful under a range of alternative simulations

When there is one causal gene in the region and it explains 1% of phenotypic variance, the FDR and power under the recommended thresholds from different methods as follows:

The 3 frequentist methods (GIFT, FOGS and MV-IWAS) are based on **Bonferroni's adjusted P-value threshold (0.05/m)** and the Bayesian method (FOCUS) is based on 90% credible sets.

**Bonferroni correction:**

The most common way to control the familywise error rate. You will find the critical value (alpha) for an individual test by dividing the familywise error rate (usually 0.05) by the number of tests.

| Methods | FDR | power |
|---------|-----|-------|
| GIFT | ↓ 0% | 46.8% |
| FOCUS | 42.1% | ↑ 70.6% |
| FOGS | 39.6% | 49.6% |
| MV-IWAS | 0.5% | 56.2% |

*All of the simulation results are from 1,000 random region.

## ■ GIFT is powerful under a range of alternative simulations

As the threshold for GIFT corresponds to a much lower FDR than the other three methods, such a threshold naturally leads to a lower power for GIFT.

To allow for fair, we further computed power **based on a true FDR of 0.05**:



One gene in the region is casual　　　Two genes in the region are casual

**a, b**: Power comparisons for different methods based on a **true FDR** of 0.05.

34

*All of the simulation results are from 1,000 random region.

# ■ GIFT is powerful under a range of alternative simulations

As the true FDR is known only in simulations but unknown for any real dataset, we also used P-value to compared power **based on the estimated FDR of 0.05**:

**Benjamini–Hochberg method:**

Order the m hypothesis by ascending p-values, where $P_i$ is the p-value at the $i$-th position with the associated hypothesis $H_i$. Let $k$ be the largest $i$ for which:

$$P_i = \frac{i}{m} q$$

Reject hypotheses $i = 1, 2, 3,..., k$. The Benjamini–Hochberg method has been proven to control the FDR for all tests at a level of q

One gene in the region is casual          Two genes in the region are casual

**a, b**: Number of genes identified by different methods based on a **estimated FDR** of q=0.05.

Colors represent the number of detected causal genes.

*All of the simulation results are from 1,000 random region.

35

# Real-data applications

# ■ Discussion



Individual-level data

eQTL data

| Sample | rs$_1$ | ... | rs$_m$ | ... | rs$_p$ |
|--------|--------|-----|--------|-----|--------|
| $n_1$ | 0 | ... | 1 | ... | 1 |
| $n_2$ | 2 | ... | 0 | ... | 1 |
| ... | ... | ... | ... | ... | ... |

| Sample | $x_1$ | $x_2$ | ... | $x_k$ |
|--------|-------|-------|-----|-------|
| $n_1$ | −0.27 | −0.16 | ... | 0.09 |
| $n_2$ | 0.93 | −0.23 | ... | 1.86 |
| ... | ... | ... | ... | ... |

Expression correlation matrix $\Omega$
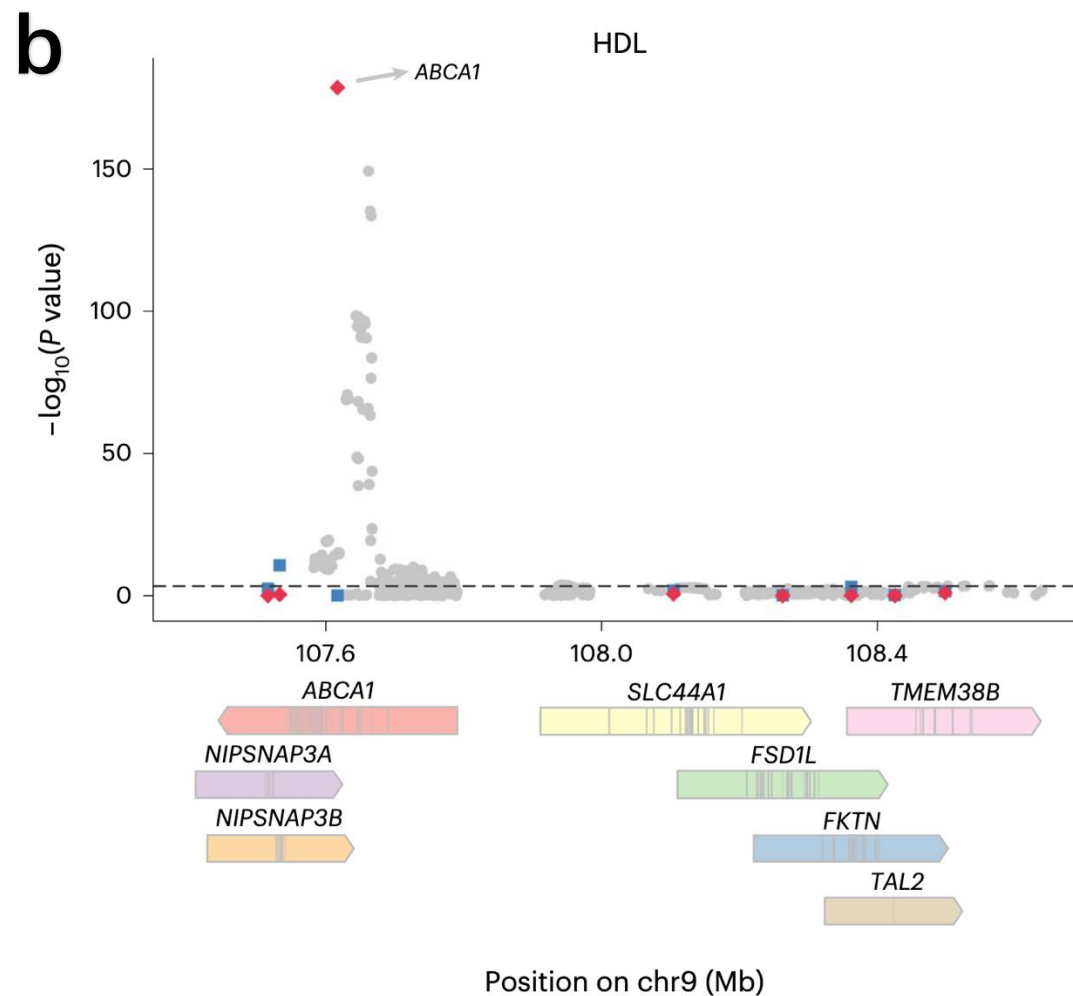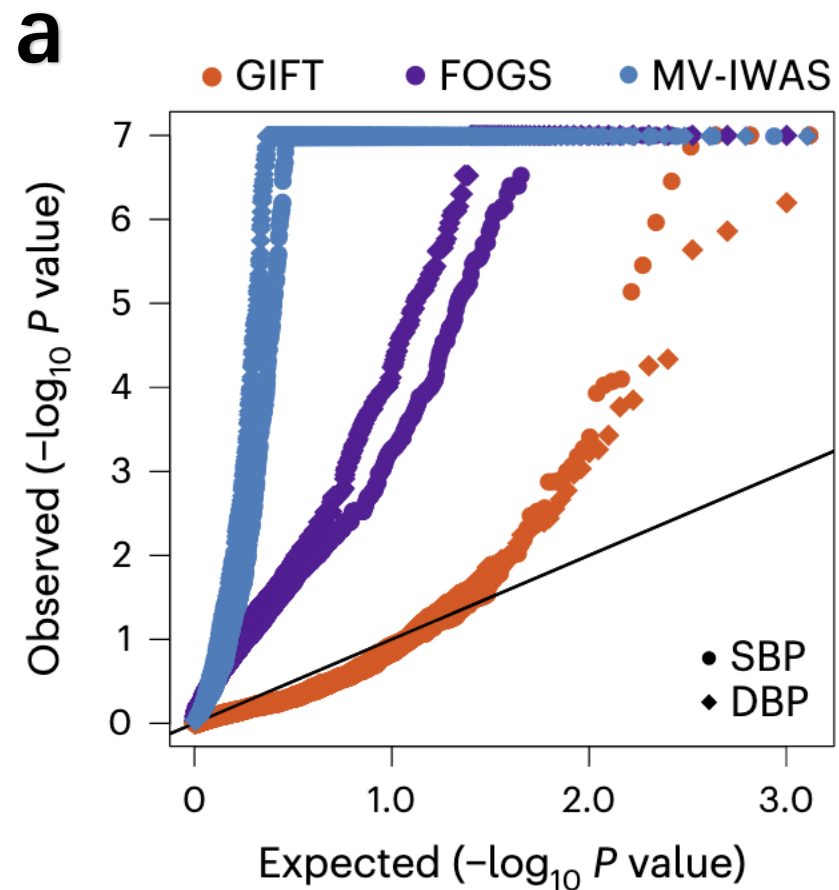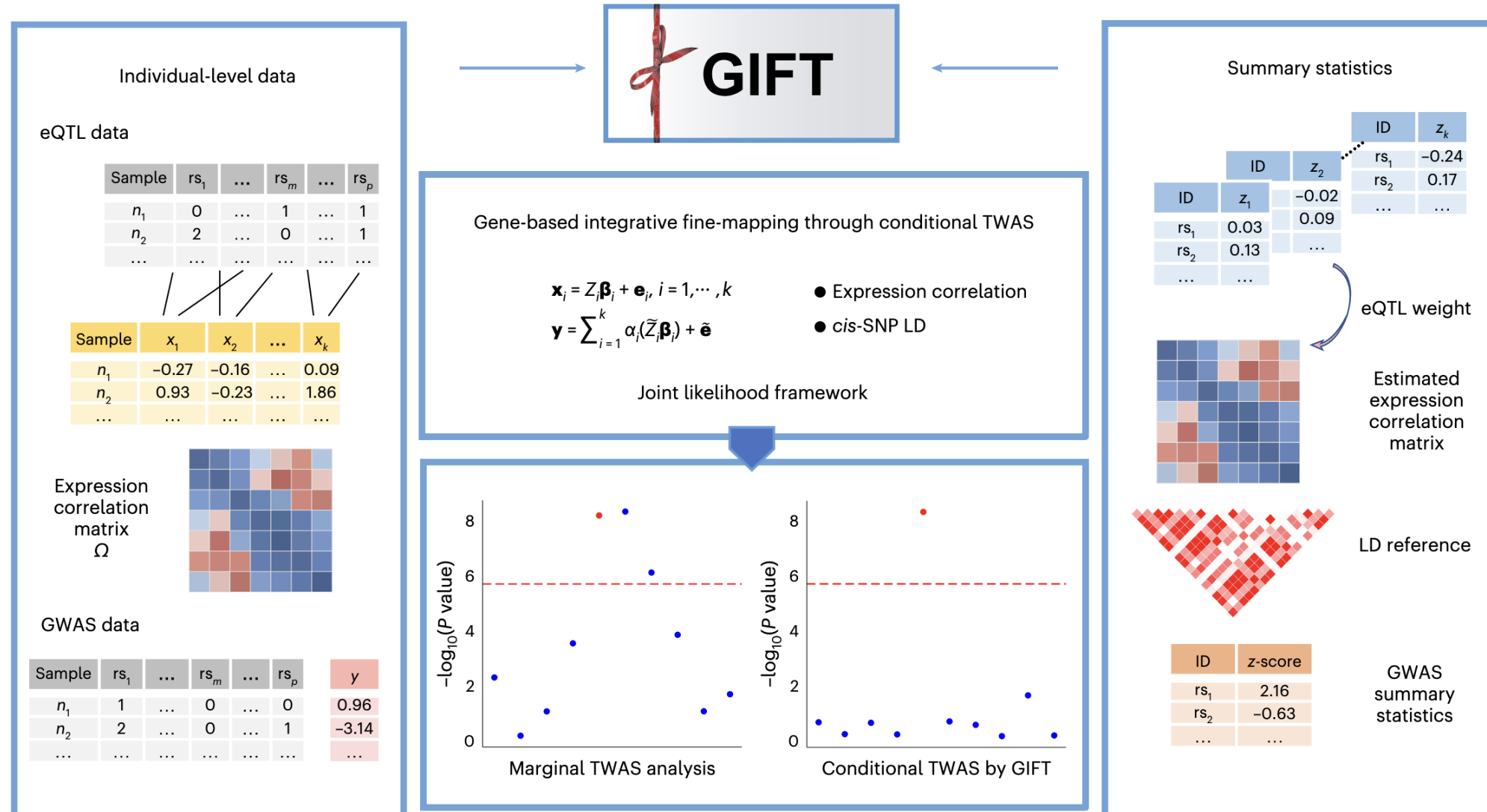
GWAS data

| Sample | rs$_1$ | ... | rs$_m$ | ... | rs$_p$ | $y$ |
|--------|--------|-----|--------|-----|--------|-----|
| $n_1$ | 1 | ... | 0 | ... | 0 | 0.96 |
| $n_2$ | 2 | ... | 0 | ... | 1 | −3.14 |
| ... | ... | ... | ... | ... | ... | ... |

**GIFT**

Gene-based integrative fine-mapping through conditional TWAS

$$\mathbf{x}_i = Z_i\boldsymbol{\beta}_i + \mathbf{e}_i, \; i = 1,\cdots,k$$

$$\mathbf{y} = \sum_{i=1}^{k} \alpha_i(\tilde{Z}_i\boldsymbol{\beta}_i) + \tilde{\mathbf{e}}$$

● Expression correlation

● *cis*-SNP LD

Joint likelihood framework

Marginal TWAS analysis

Conditional TWAS by GIFT

Summary statistics

| ID | $z_k$ |
|----|-------|
| rs$_1$ | −0.24 |
| rs$_2$ | 0.17 |
| ... | ... |

| ID | $z_2$ |
|----|-------|
| | −0.02 |
| | 0.09 |
| ... | |

| ID | $z_1$ |
|----|-------|
| rs$_1$ | 0.03 |
| rs$_2$ | 0.13 |
| ... | ... |

eQTL weight

Estimated expression correlation matrix

LD reference

| ID | z-score |
|----|---------|
| rs$_1$ | 2.16 |
| rs$_2$ | −0.63 |
| ... | ... |

GWAS summary statistics

37

# THANK YOU

**Speaker: Yuekai Li**

**Major: Biostatistics**