

Integrative Methods

DCEG Statistical Genetics Workshop

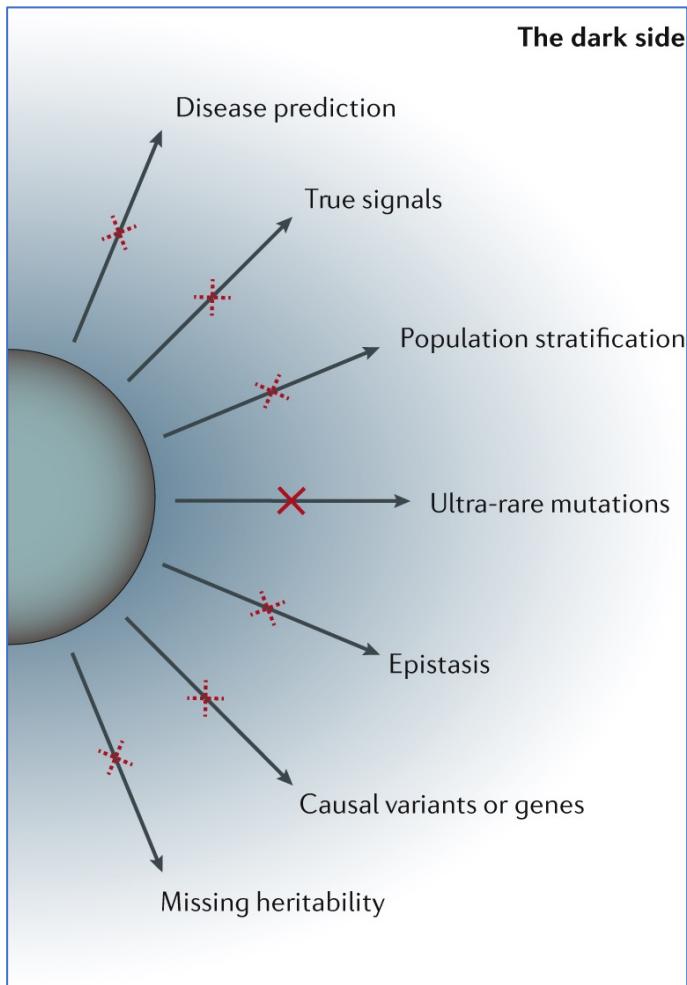
Lecture 6 (Part 1)

15 November 2023

Outline

- Transcriptome-Wide Association Studies (TWAS)
 - “General Steps”
 - Software
 - Why Perform a TWAS?
 - Limitations / Caveats
- Resources
- Other “–WAS” Approaches
 - pWAS
 - Metabo-WAS
 - Epigenome-WAS
 - Regulome-WAS
- Example Applications in Practice – Breast Cancer Risk

Background: Moving Beyond GWAS...

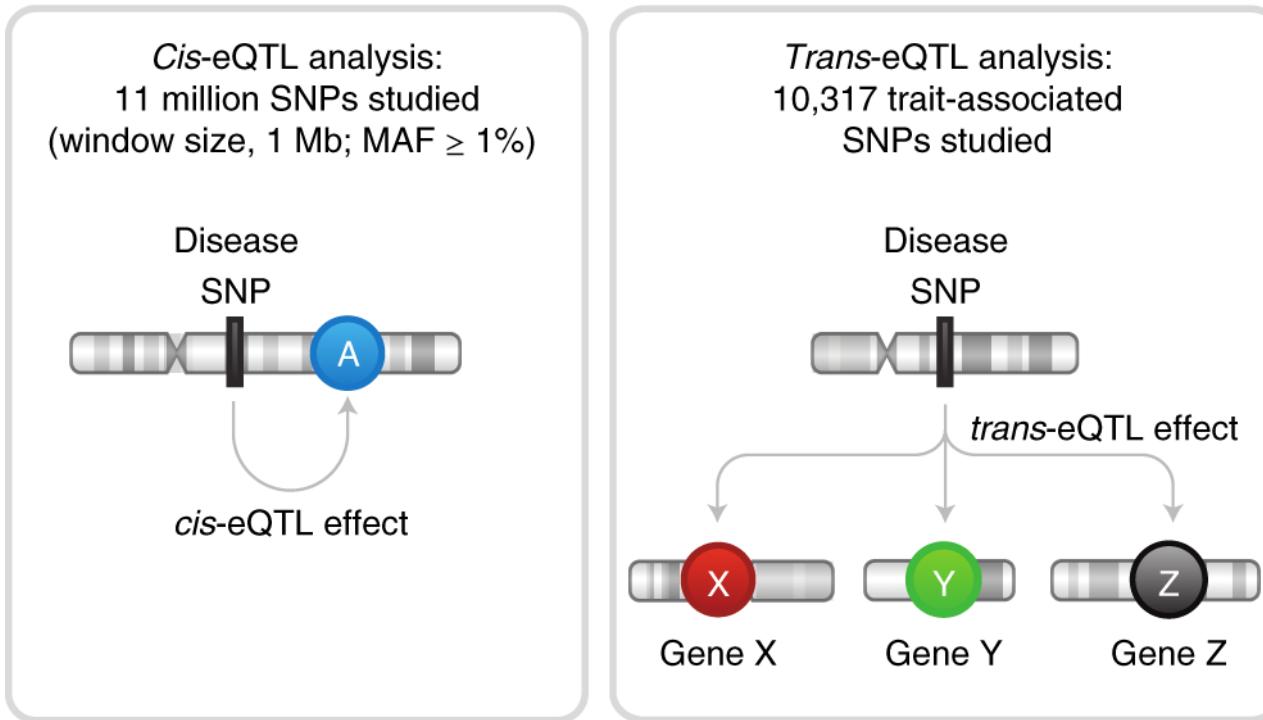


“The Dark Side” of GWAS:

=> **Most associated variants not present in coding regions of the genome...**

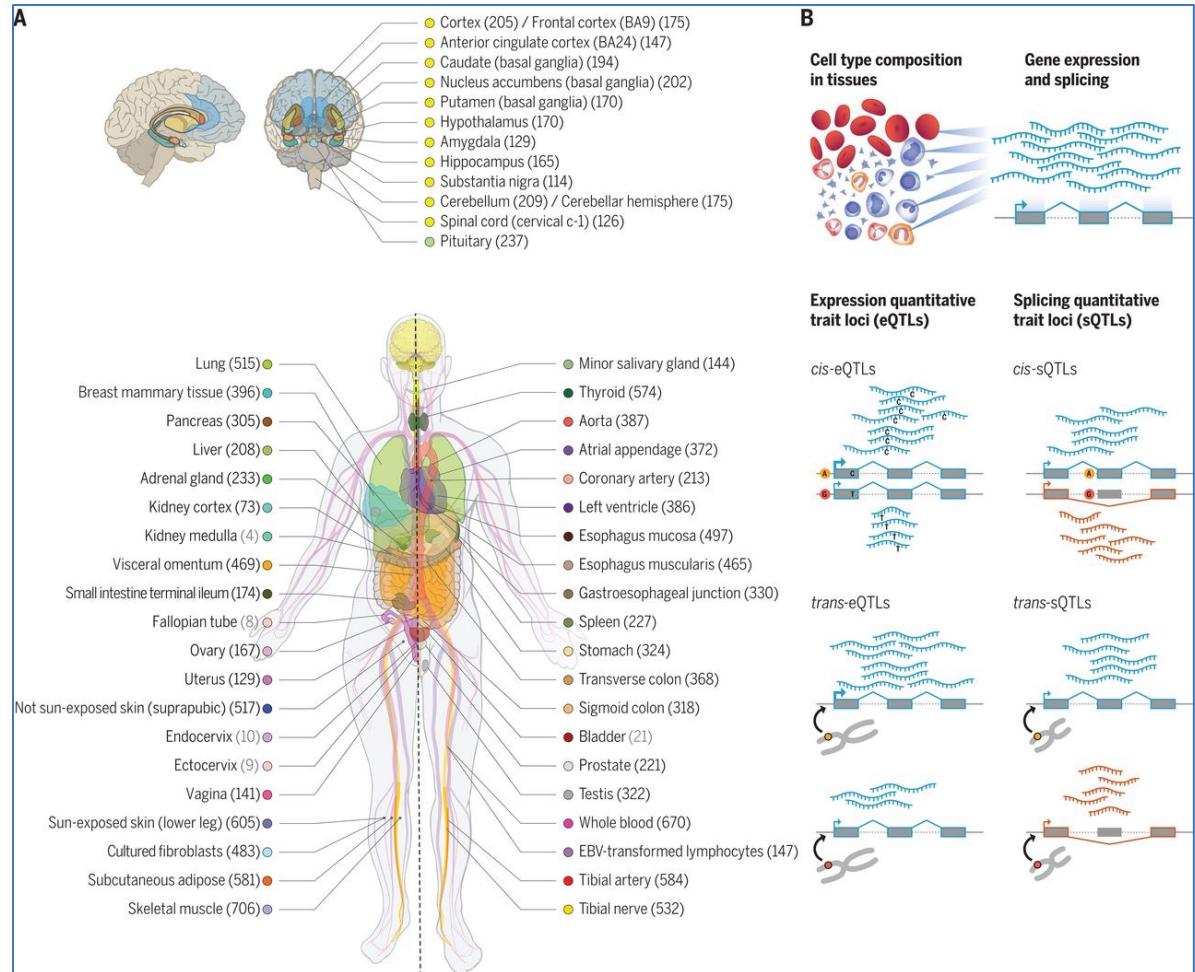
Background: Expression Quantitative Trait Loci (eQTLs)

Locus (typically non-coding) that explains variation in mRNA / gene expression



Background: Genotype-Tissue Expression Database (GTEx v8)

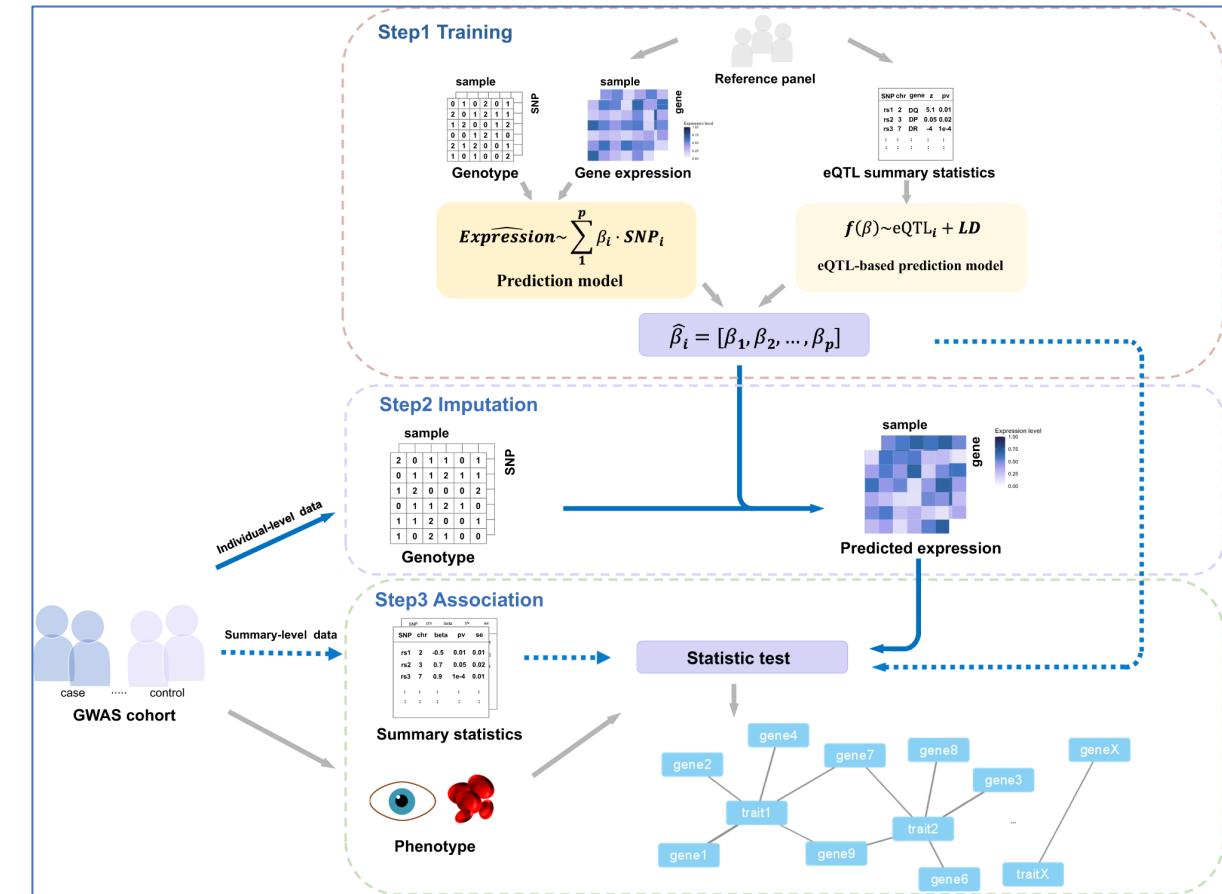
- Estimation of “normal” effects of genotype -> expression change per tissue
- Tissue-specific gene expression from rapid autopsy donors
- **54 tissues, 948 donors**



“General Steps” of a Transcriptome-Wide Association Study (TWAS)

- Define **genotype** and **phenotype**
- Step 1: Use **reference transcriptome / prediction models with regulatory weights** to train genotype -> expression
- Step 2: Impute genotype change or summary statistics -> predicted expression change
- Step 3: Associate imputed expression changes with phenotype of interest

RESULT: Gene-phenotype associated mediated through the transcriptome



Prediction Model: Basic Design

$$E_g = \mu + X_{np} * \beta + \varepsilon$$

Expression level

X: n*p genotype matrix
p-vector SNPs
n individuals

β : p-vector of
SNP weights based on eQTL
effect sizes

Association Testing: Basic Design

- Gene-trait (per tissue) associations
 - Linear regression
 - Logistic regression
 - Cox regression
 - Spearman model

“1st Generation” TWAS Software

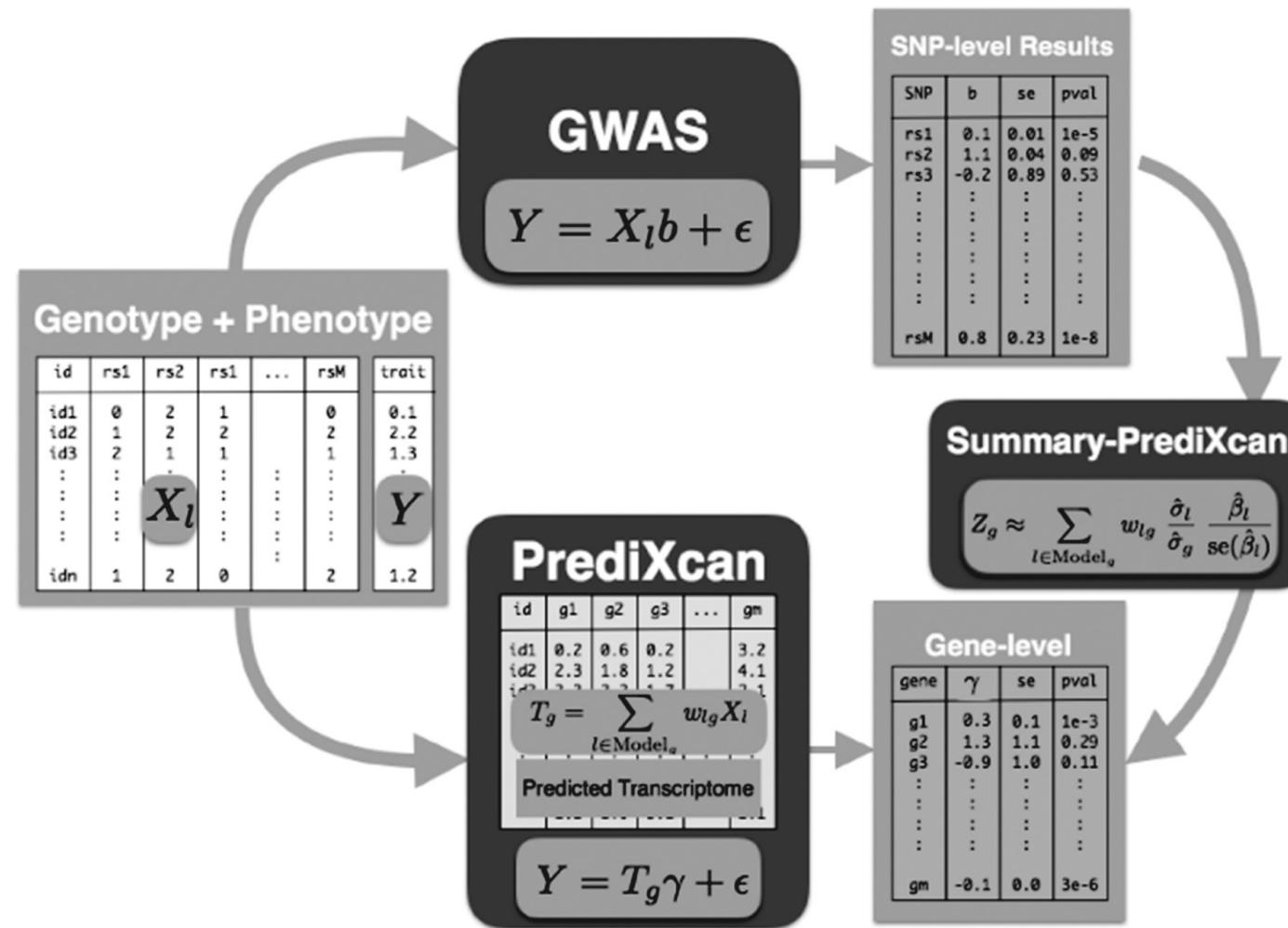
PrediXcan -> S-PrediXcan

- PredictDB for prediction models (leverages lasso/elastic net)
- *Cis*-eQTLs
- Individual tissues
- Imputation / Summary statistics

FUSION

- Bayesian sparse linear mixed model (BSLMM) for prediction models
- *Cis*-eQTLs
- Individual tissues
- Summary statistics

PrediXcan/S-PrediXcan Framework



“2nd Generation” TWAS Software

MultiXcan / S-MultiXcan

- PredictDB
- *Cis*-eQTLs
- Across tissues with an F test
- Summary statistics

UTMOST

- Across tissue prediction model with Lasso, Berk-Jones test
- *Cis*-eQTLs
- Summary statistics

TIGAR

- Non-parametric Dirichlet process regression for prediction models
- *Cis*-eQTLs
- Individual tissues
- Imputation

“3rd Generation” TWAS Software

MOSTWAS: Multi-Omics

- Chromatin status, transcription factor binding sites, microRNA effects, and CpG methylation also included
- *Cis-* and *trans*-eQTLs
- Individual tissues
- Summary statistics

Multi-Ancestry

- MATS
- TESLA

Why Perform a TWAS?

- Identifies genes regulated by disease-associated variants
- Gene-level association to phenotype of interest =>
 - Statistical power (improves multiple testing burden of variants)
 - Easier experimental validation (gene changes vs. specific variant changes)
- Lower computing complexity vs. GWAS
- Tissue effect specificity

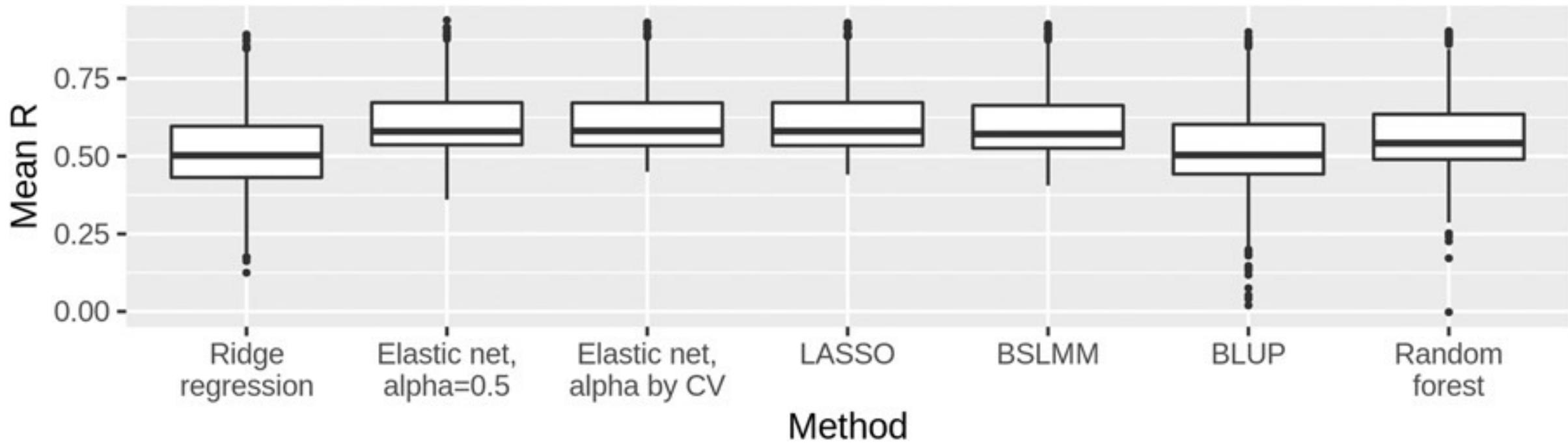
Limitations / Caveats

- Accuracy of TWAS Depends On:
 - Prediction Models!!
 - Underlying GWAS assumptions
- Correlated prediction: Contamination from linkage disequilibrium
 - Fine mapping and colocalization => **Return back to lecture 3!**
- Correlated expression: True co-regulation

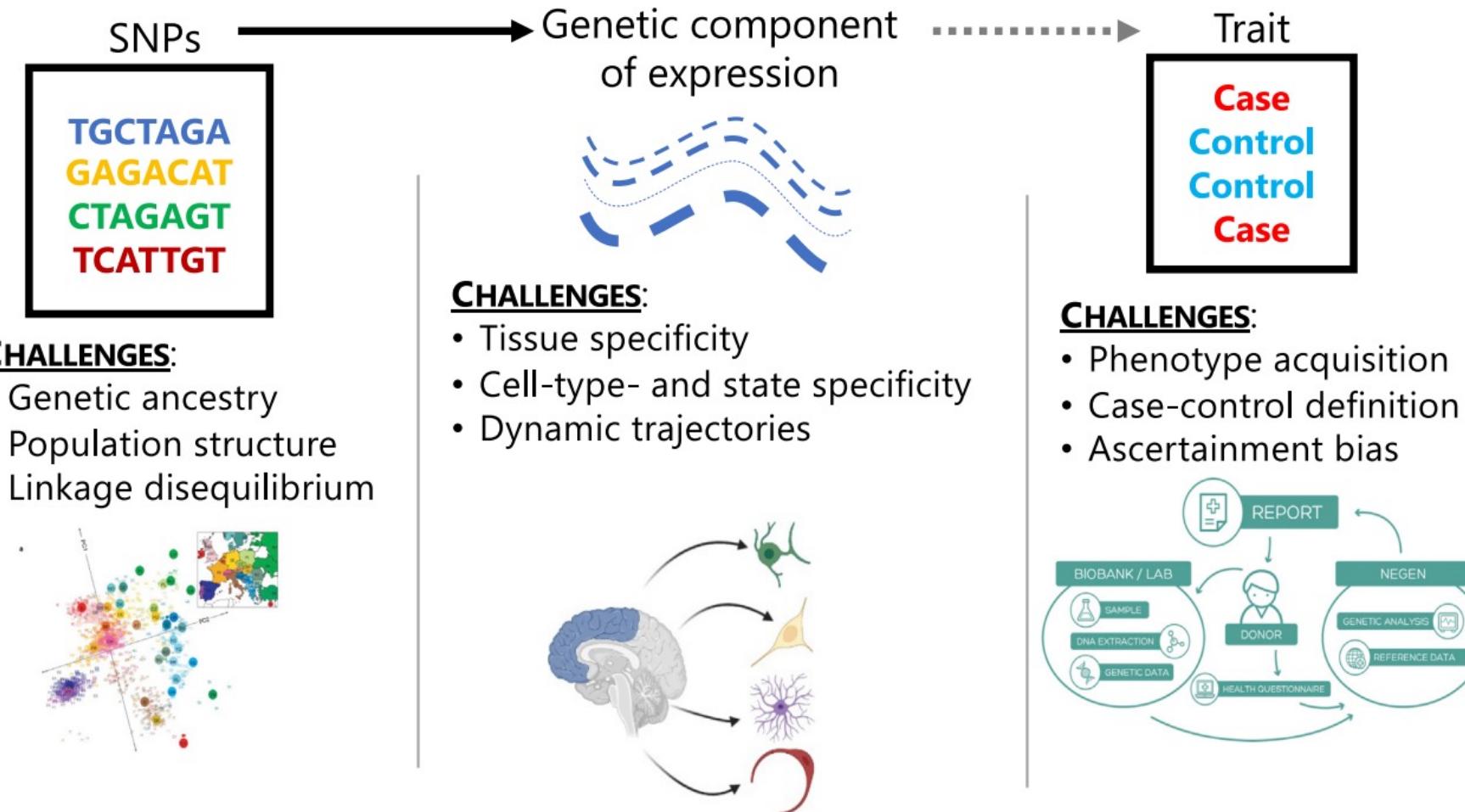
Accuracy of TWAS Approaches: Prediction Models

- Prediction model ancestry needs to correspond to TWAS population ancestry or be multi-ancestry
- Focus on *cis*-eQTLs biases against *trans*-eQTLs
- SUMMIT: Estimates *cis*-eQTL effect size (weights for prediction models) with eQTL summary statistics, LD information w/ shrinkage estimator
- Multi-ancestry modeling

Accuracy of TWAS Approaches: Prediction Models – Statistical Design

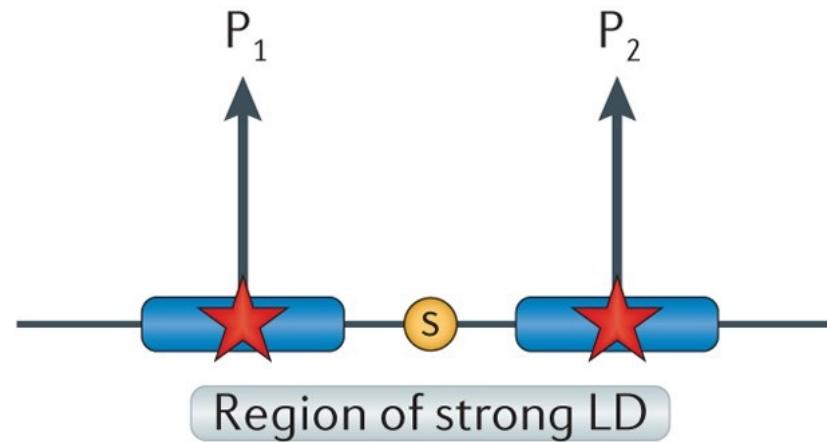


Accuracy of TWAS Approaches: Relation to GWAS



Limitations: LD Contamination

Spurious pleiotropy:
causal variants in different genes



MUST USE: Fine mapping and colocalization to resolve

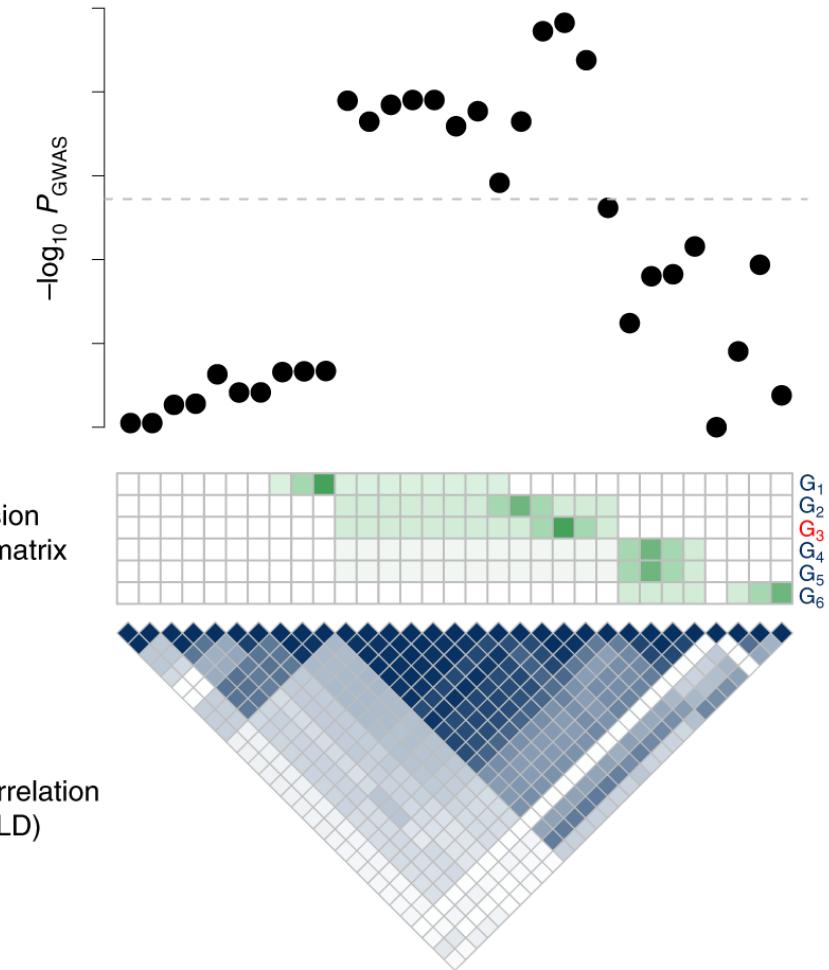
Revisiting: Fine Mapping and Colocalization (in TWAS)

Fine Mapping

- Goal: To refine more likely causal genes within a given LD block
- Prior fine mapping tools could be applied (e.g. eCAVIAR, DAP-G, SUSIE-R)

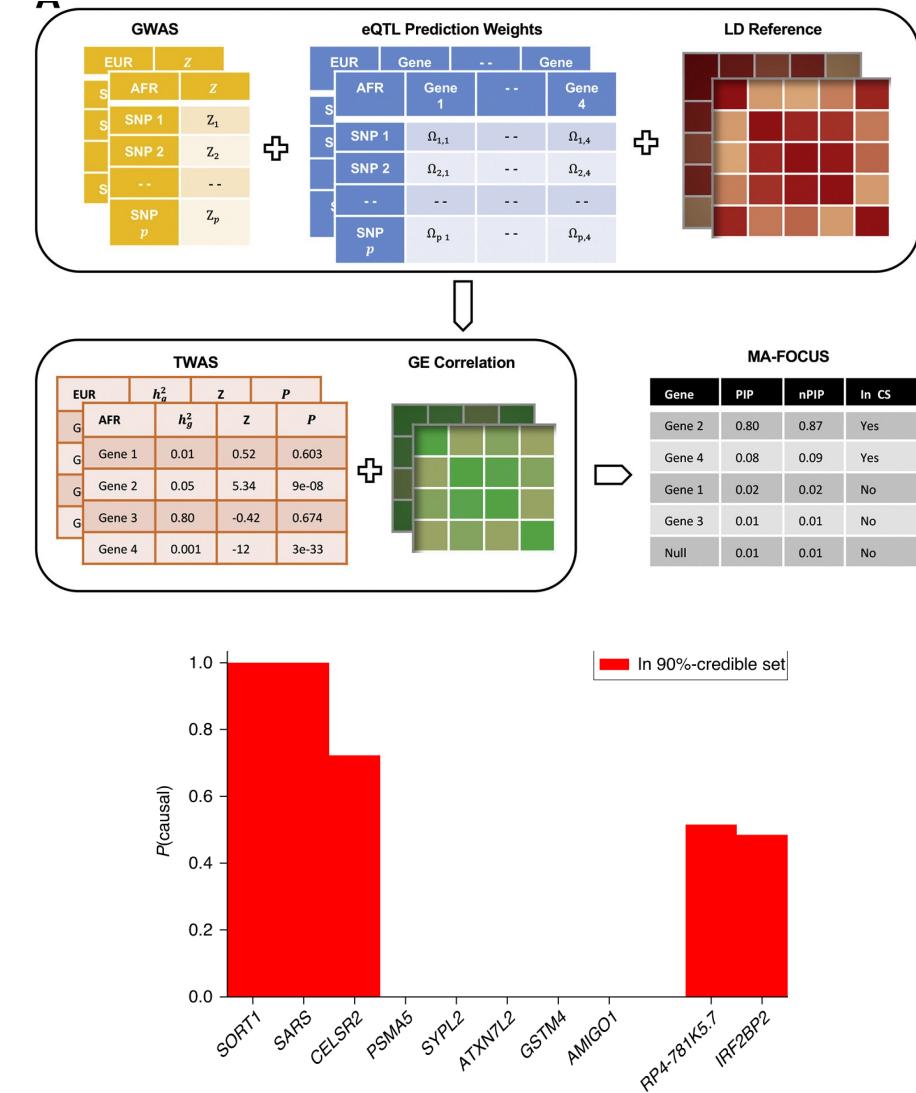
Colocalization

- Goal: To limit effect of LD contamination by prioritizing most likely signal using posterior probability estimates based on regional structure
- Prior colocalization tools can be applied (coloc, enloc, eCAVIAR)



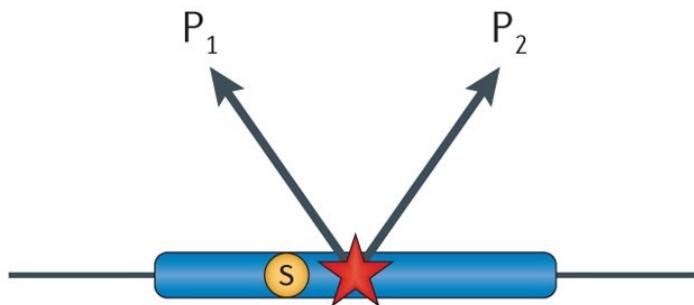
FOCUS/MA-FOCUS: TWAS-Specific Estimates of Causal Probabilities

- Models predicted expression correlations and uses them to assign genes posterior probabilities of causality (MA = Multi-Ancestry)
- **Input:** GWAS summary data, expression prediction weights, LD among SNPs in region
- **Output:** Probability for each gene in region to explain TWAS signal (PIP and CS from fine-mapping)

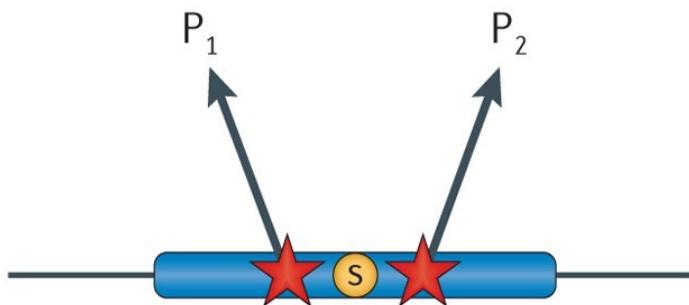


Limitations: True Co-Regulation

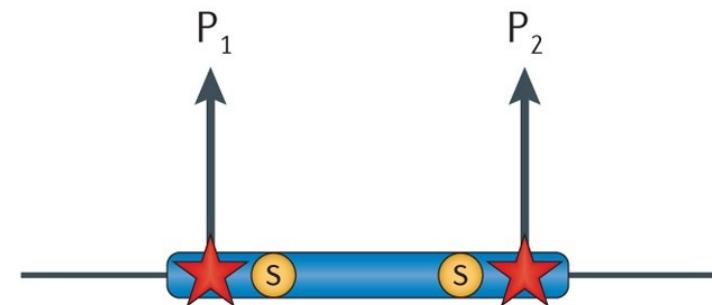
Biological pleiotropy:
single causal variant



Biological pleiotropy: different causal
variants colocalizing in same gene and
tagged by the same genetic variant



Biological pleiotropy: different causal
variants colocalizing the same gene



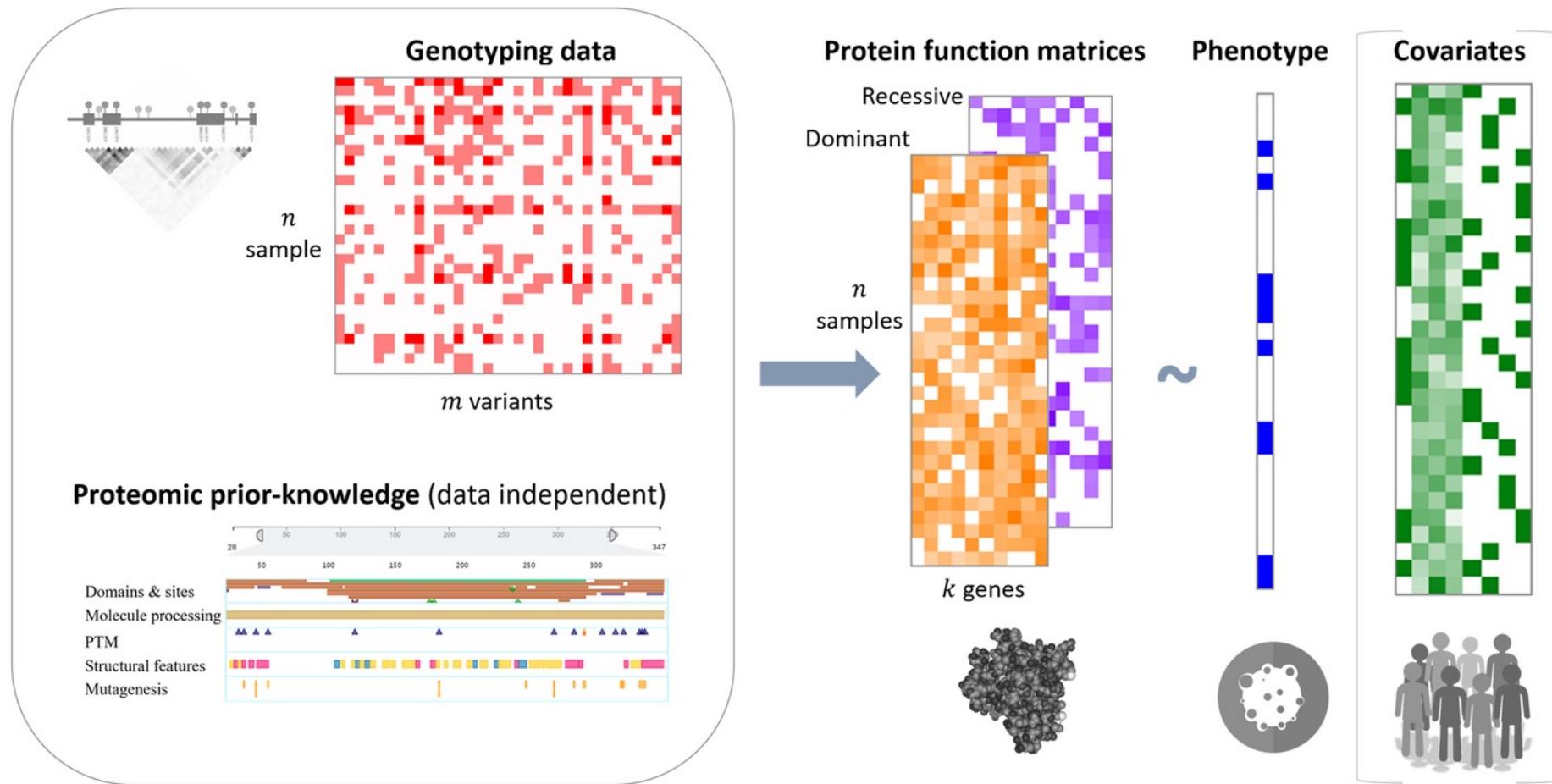
Resources: Reference Prediction Models / Datasets

- GTEx: <https://gtexportal.org/home/>
- eQTLGen Consortium: <https://www.eqtlgen.org/>
- Predict DB: <https://predictdb.org/>
- FUSION Prediction Models:
<http://gusevlab.org/projects/fusion/#gtex-v8-multi-tissue-expression>

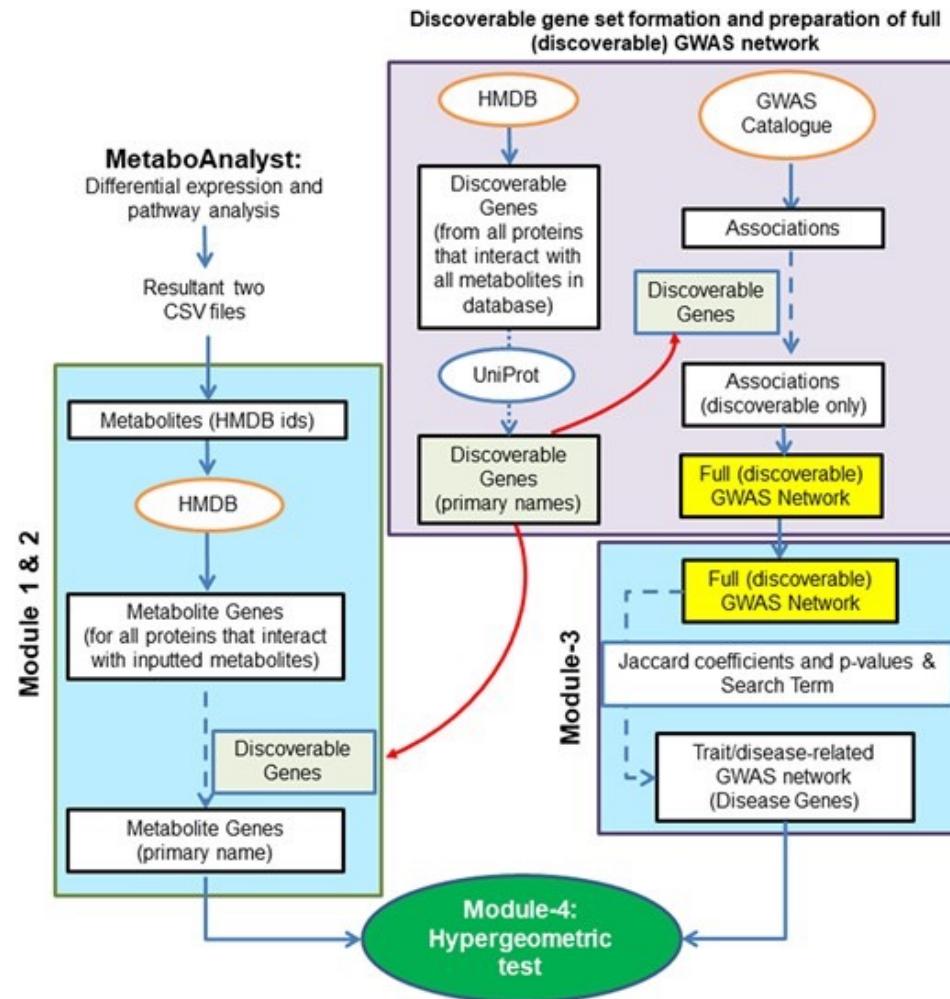
Resources: TWAS Databases

- <http://twas-hub.org/>
- <http://www.webtwas.net/>
- <https://ngdc.cncb.ac.cn/twas/>

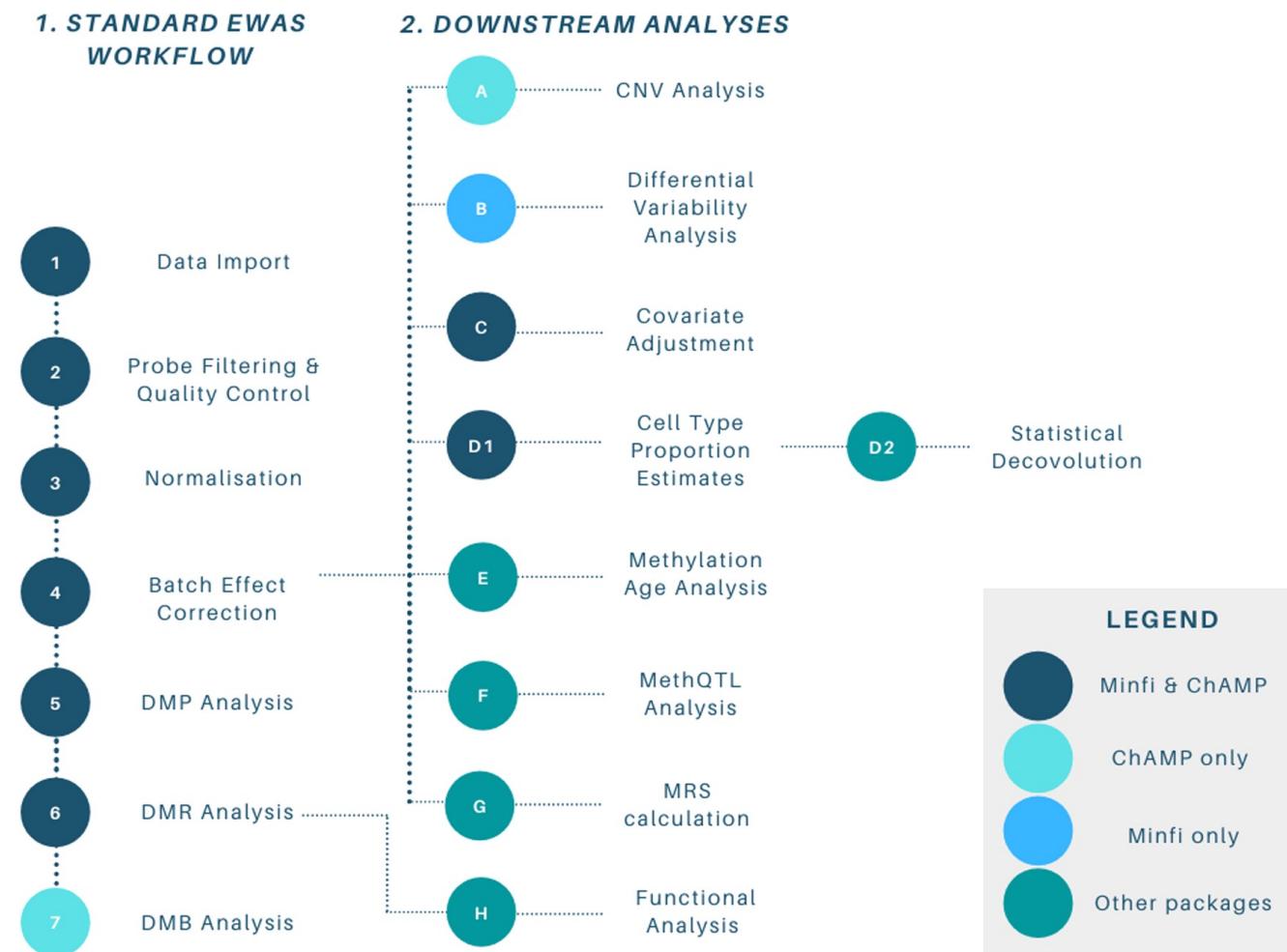
Other “-WAS” Approaches: Proteome-WAS



Other “-WAS” Approaches: Metabo-WAS

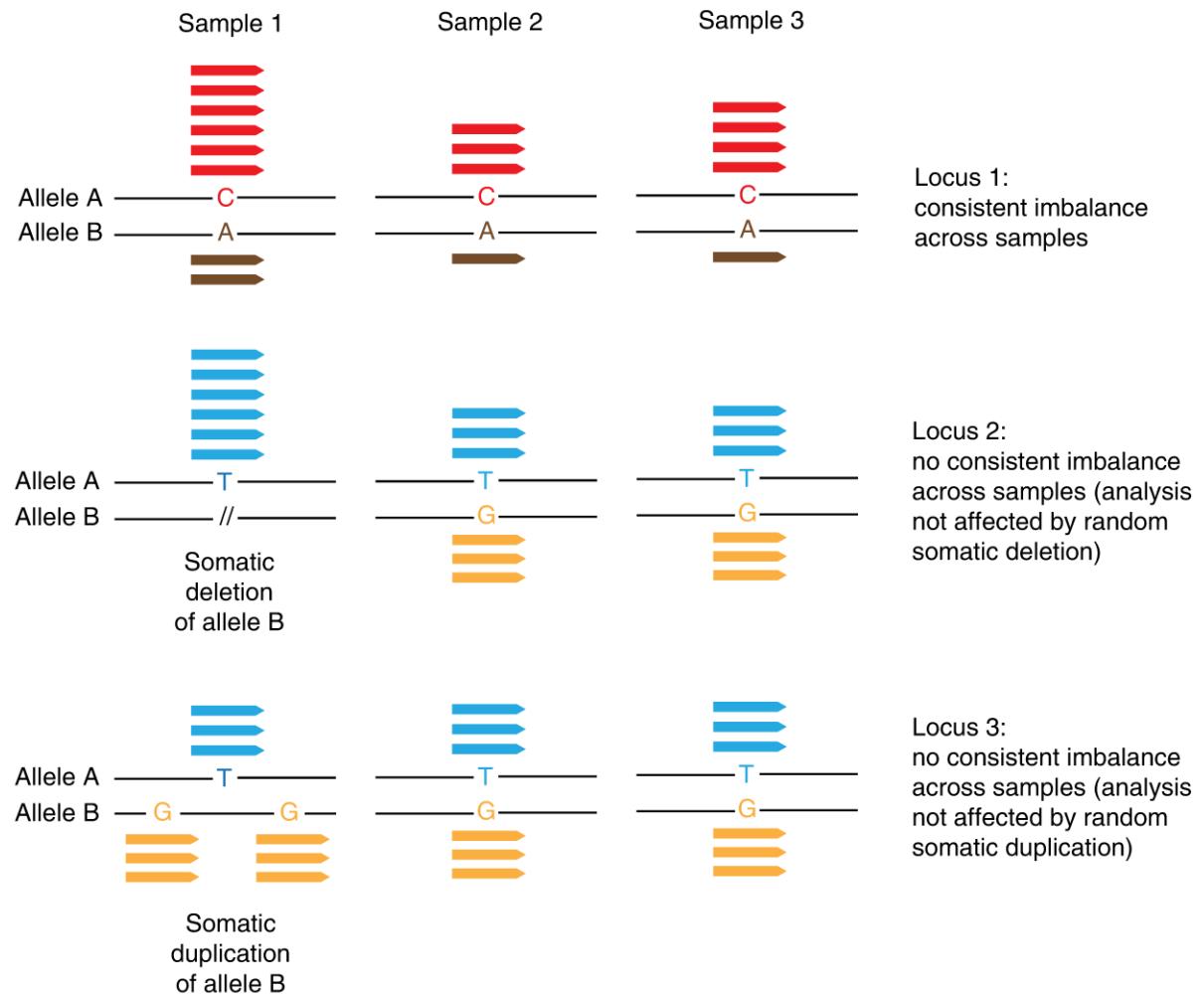


Other “-WAS” Approaches: Epigenome-WAS



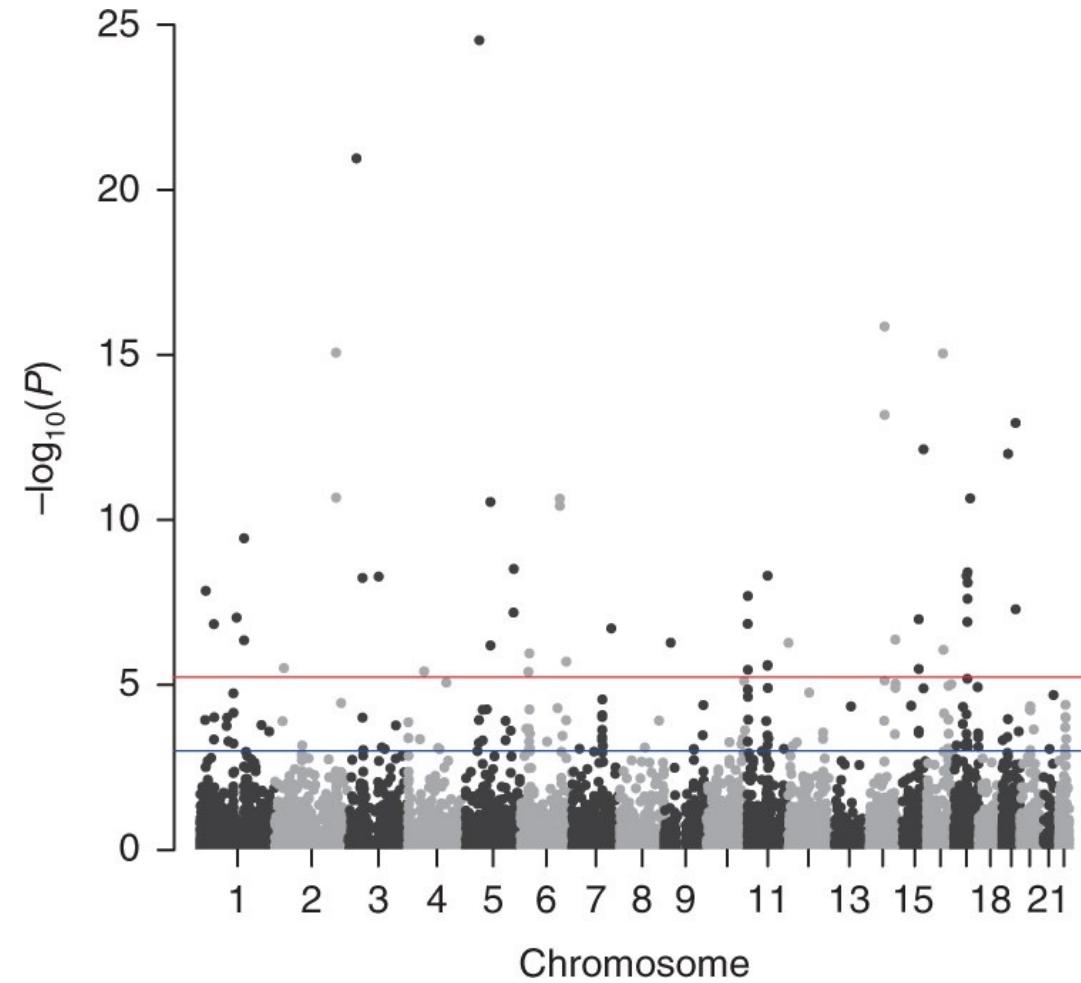
Other “-WAS” Approaches: Regulome-WAS

- Leverages ATAC-Seq samples for allele-specific accessible eQTLs



Applications in Practice: Breast Cancer Risk

| Region | Gene |
|----------|----------------|
| 1p34.1 | ZSWIM5 |
| 3p24.1 | LRRC3B |
| 4q12 | SPATA18 |
| 6p22.1 | UBD |
| 7p32.2 | KLHDC10 |
| 9p21.3 | MIR31HG |
| 11p15.5 | RIC8A |
| 11q13.2 | B3GNT1 |
| 11q13.2 | RP11-867G23.10 |
| 12q13.33 | RP11-218M22.1 |
| 14q24.1 | GALNT16 |
| 14q24.1 | PLEKHD1 |
| 15q24.2 | MAN2C1 |
| 15q24.2 | CTD-2323K18.1 |

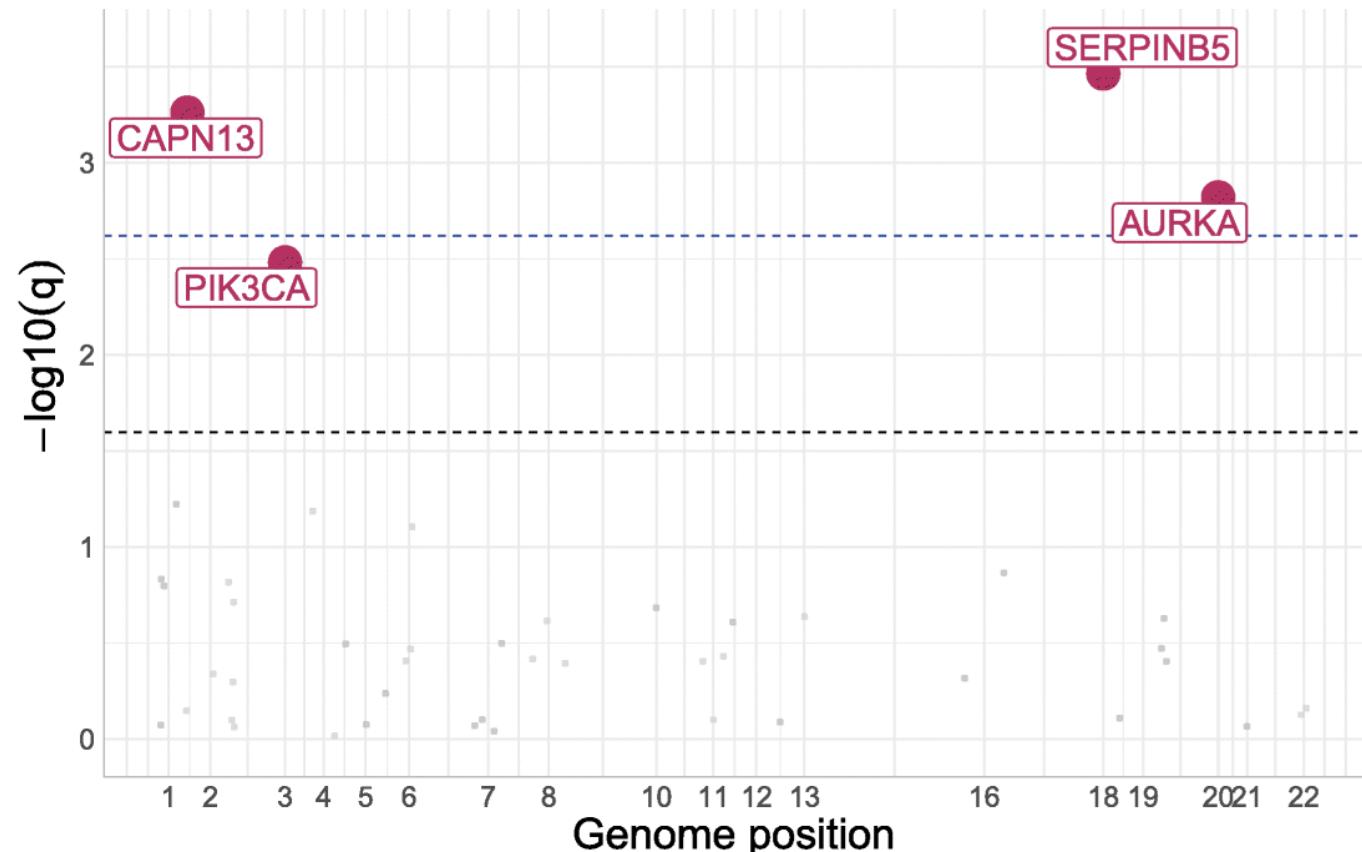


Applications in Practice: Breast Cancer Risk - Splicing-Associated Variants

- Conditional joint TWAS across multiple tissues
- At least 1MB from prior GWAS sites and noncoding -> potential splicing related variants = Greater likelihood of association with splicing

Applications in Practice: Breast Cancer Risk - Ancestry-Specific

- Race-stratified predictive models of germline genotype \rightarrow tumor expression in 406 genes related to breast cancer
- Models not applicable across race / subtype
- **Identified associations in Black women near AURKA, CAPN13, PIK3CA, and SERPINB5 via TWAS that are underpowered in GWAS**



Overview

- Transcriptome-Wide Association Studies (TWAS)
 - “General Steps”
 - Software
 - Why Perform a TWAS?
 - Limitations / Caveats
- Resources
- Other “–WAS” Approaches
 - pWAS
 - Metabo-WAS
 - Epigenome-WAS
 - Regulome-WAS
- Applications in Practice – Breast Cancer Risk