# 1 GENETICS OF GENE EXPRESSION

# Chromosome map of disease-associated regions



2019 July

www.ebi.ac.uk/gwas

# "GWAS have so far identified only a small fraction of the heritability of common diseases, so the ability to make meaningful predictions is still quite limited"

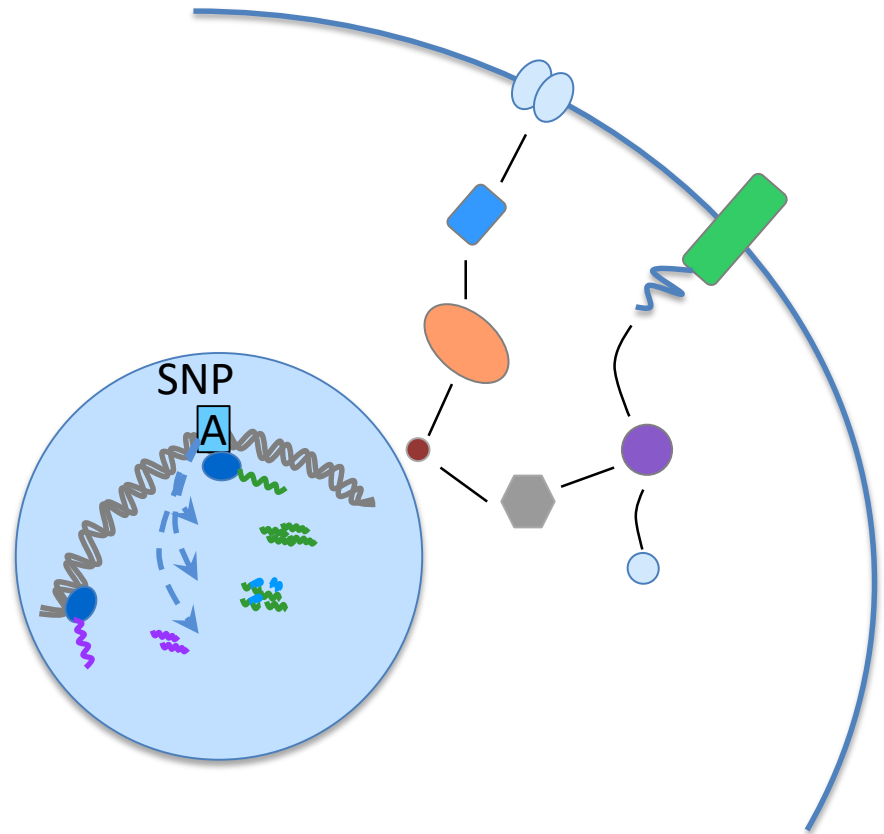Francis Collins, Director of the NIH, *Nature*, April 2010

| Trait | Heritability (Family base) | Individuals studied | Heritability explained |
|---|---|---|---|
| Coronary artery disease | 40% | 86995 | 10% |
| Type 2 Diabetes | 40% | 47117 | 10% |
| BMI | 50% | 249796 | 3% |
| Blood pressure | 50% | 34433 | 1% |
| Circulating lipids | 50% | 100000 | 25% |
| Height | 80% | 183727 | 12.5% |

# Motivation

How can we use gene expression and epigenetics to help us understand complex trait genetics?

Majority of trait-associated variation is non-coding.

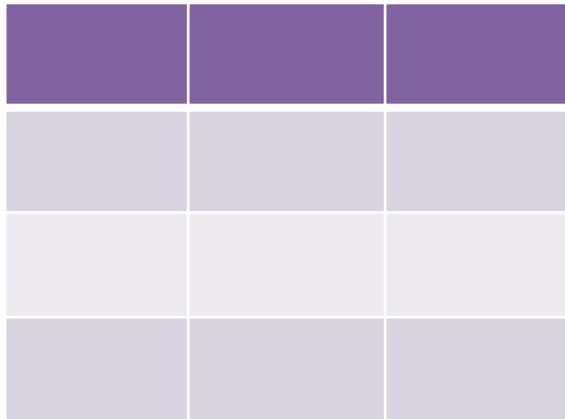Common hypothesis is that most of these function by altering gene expression.
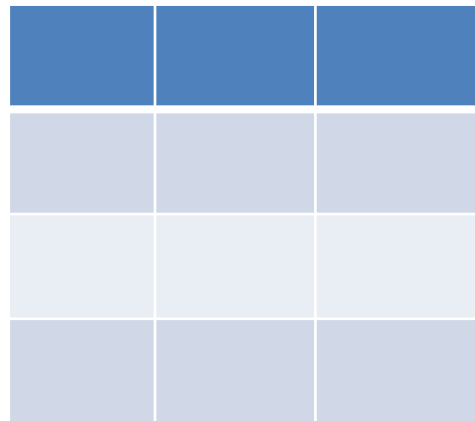
# eQTL analysis Statistics

- Regression: find the coefficients for the effect of expression on genotype when conditioned on the covariates in a linear model and test if they are significant diffetent than 0

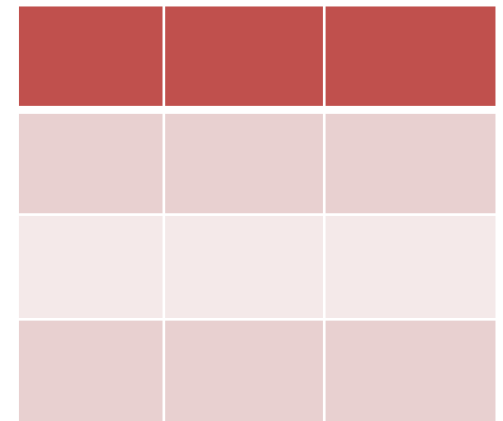$$gene\ expression = \beta_0 + \beta_1 genotype + \beta_2 covaraites$$
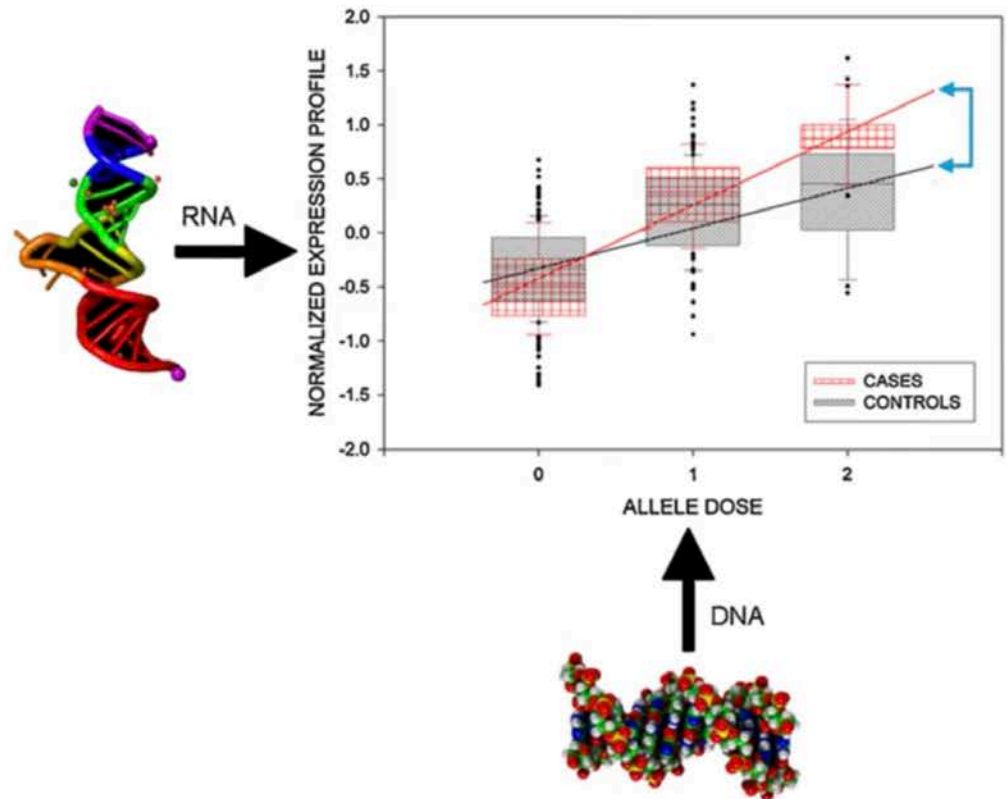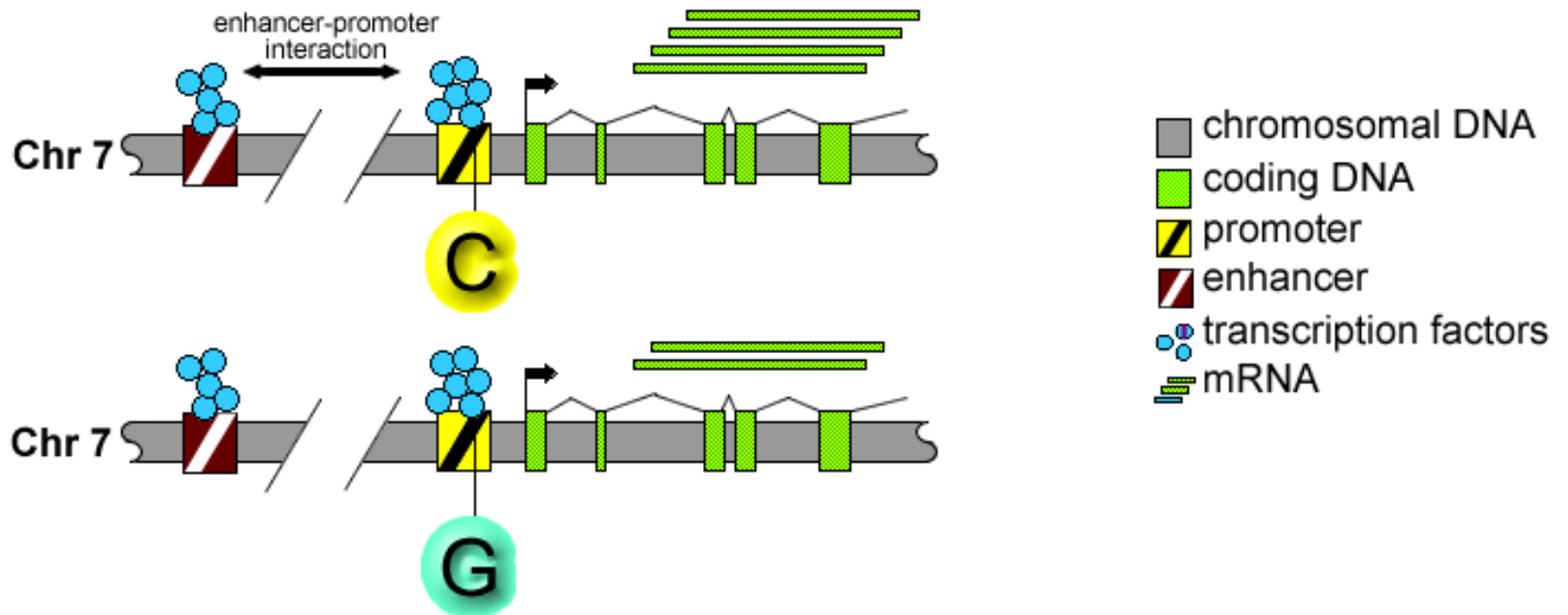
Expression

Genotype

Covaraite

# Mapping expression (e)QTL

- RNA expression levels can be treated like any other quantitative trait in QTL mapping.

- 30,000 genes by 10,000 SNPs = 300,000,000 comparisons!

- eQTL studies are sometimes called genetical genomics



Myers, A.J. The age of the "ome": Genome, transcriptome and proteome data set collection and analysis. Brain Research Bulletin Volume 88, Issue 4 2012 294 - 301
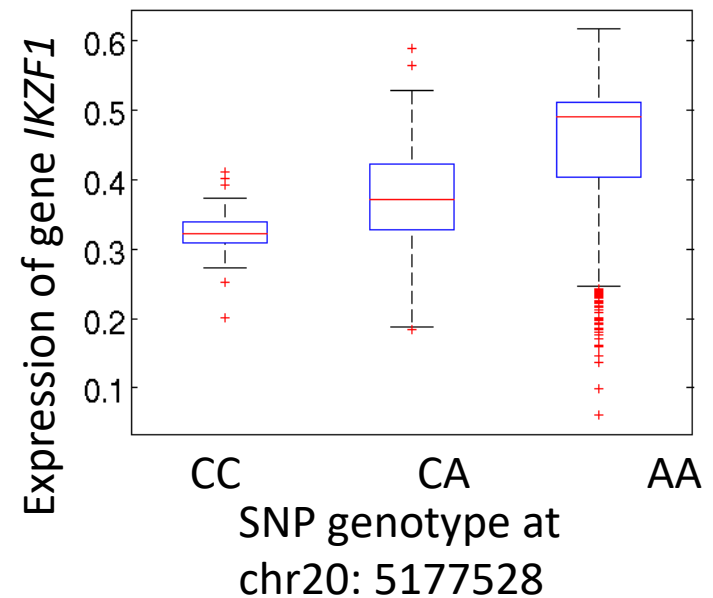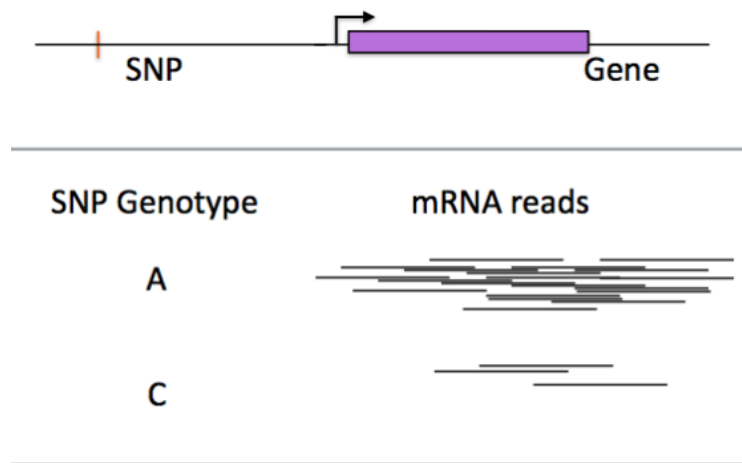
*cis*- effect
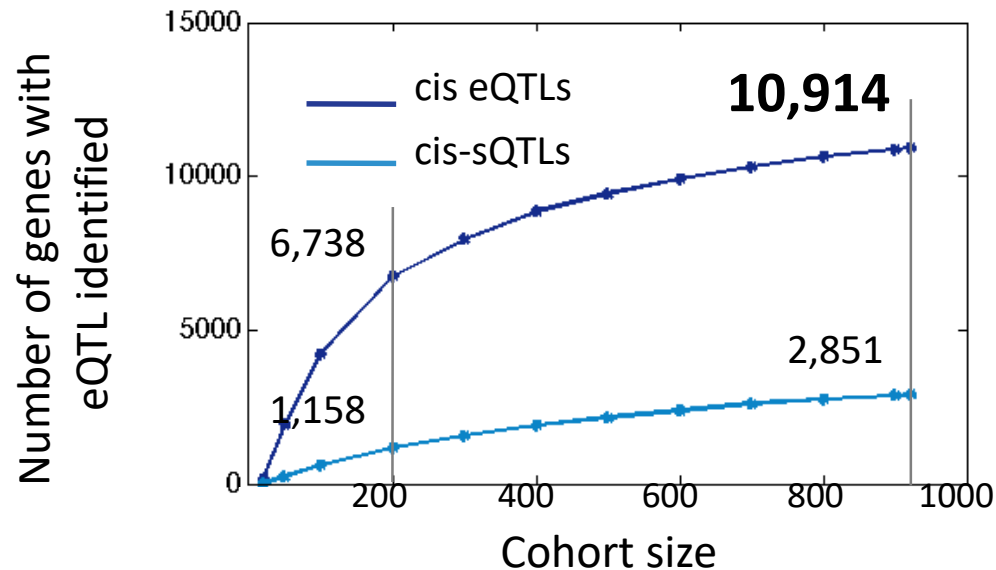
Canonical model

# Genetic variants affect gene expression

**eQTL (expression Quantitative Trait Locus) analysis:**
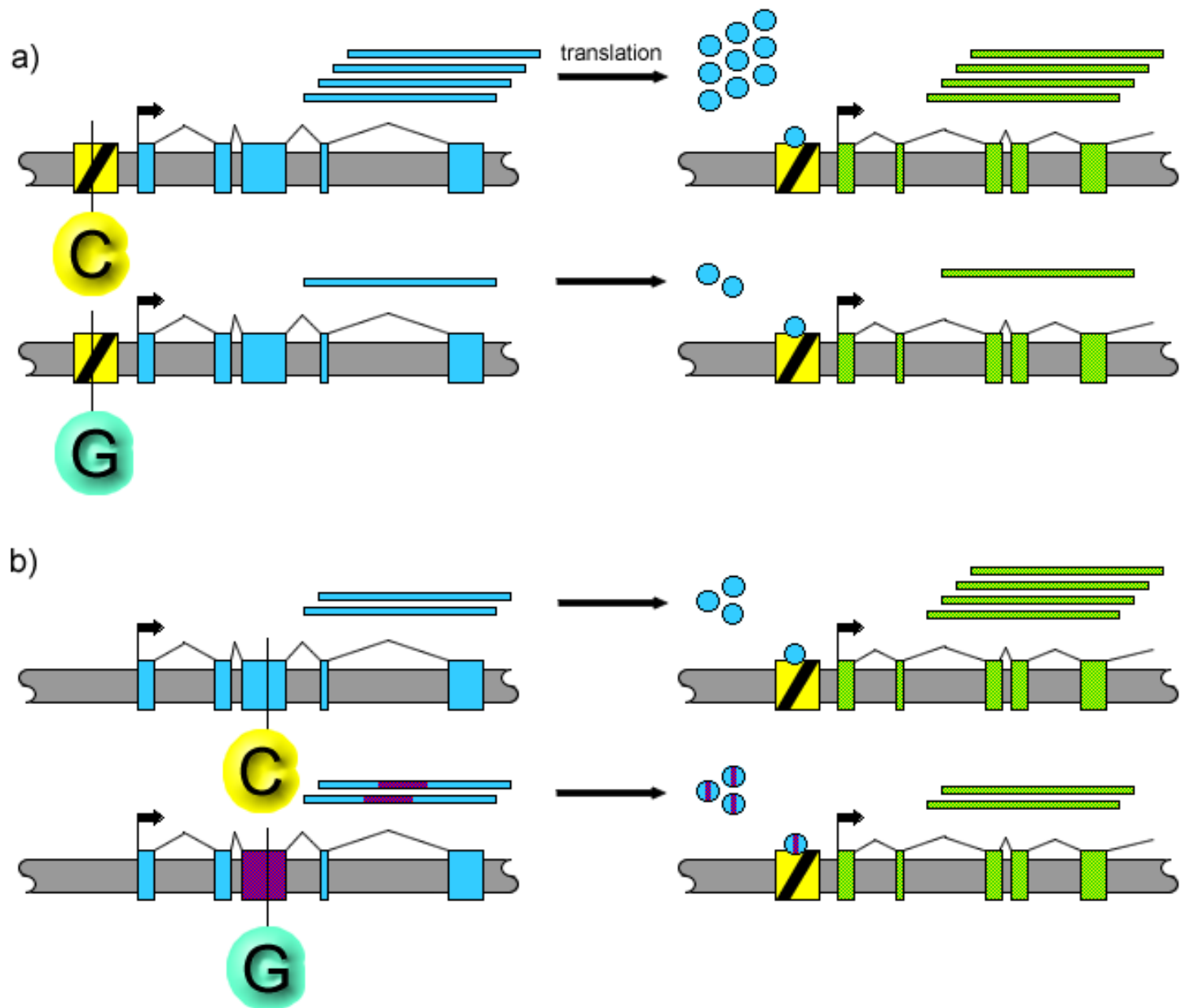Association between genotype and RNA expression levels
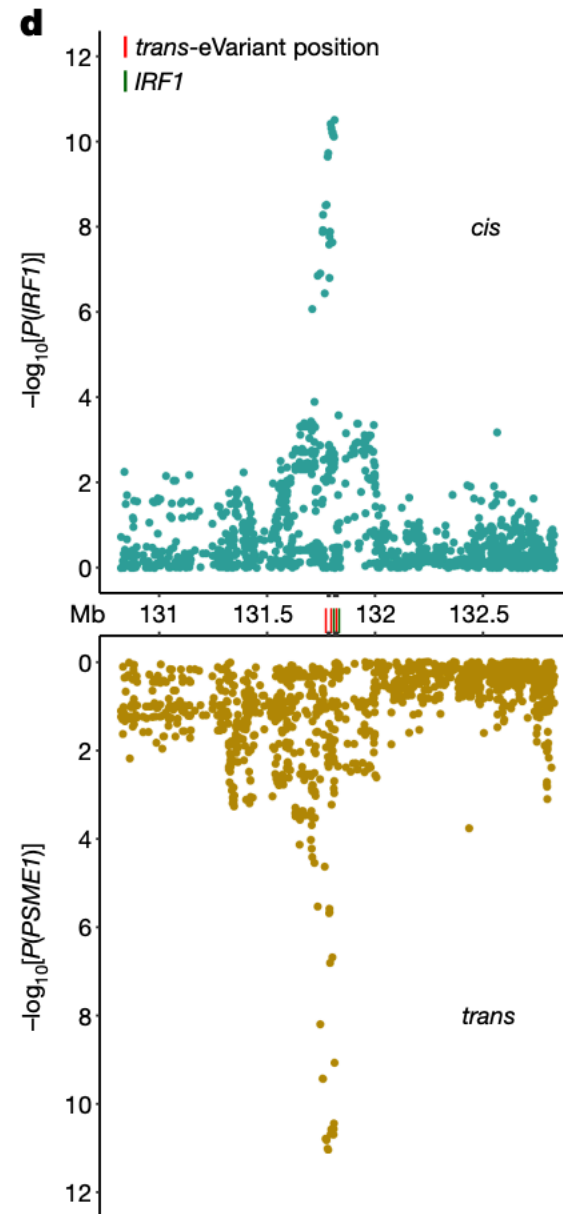
# Prevalence of eQTLs

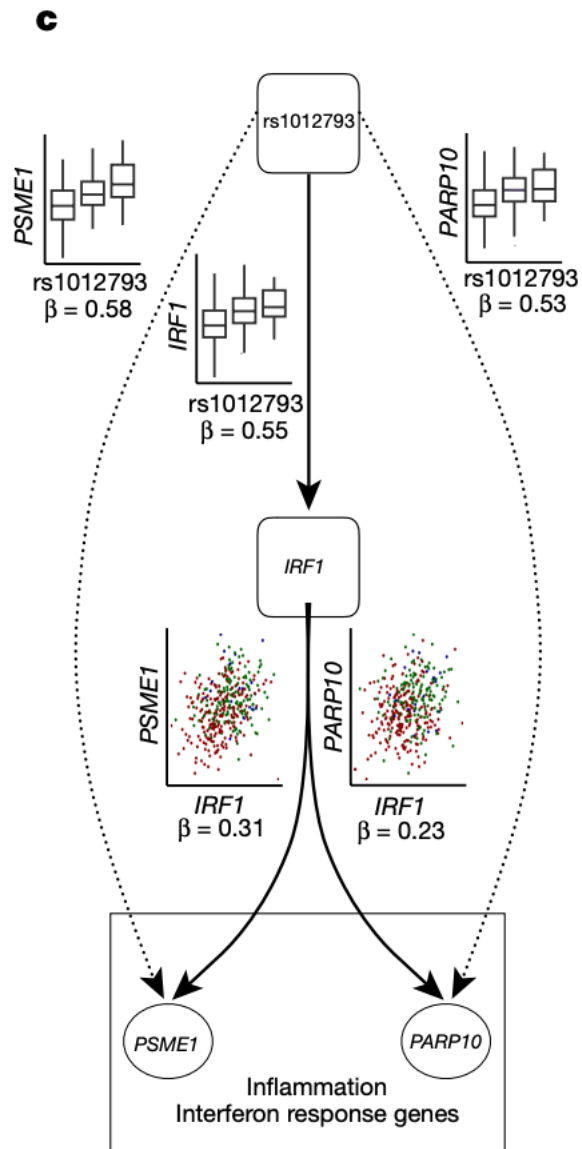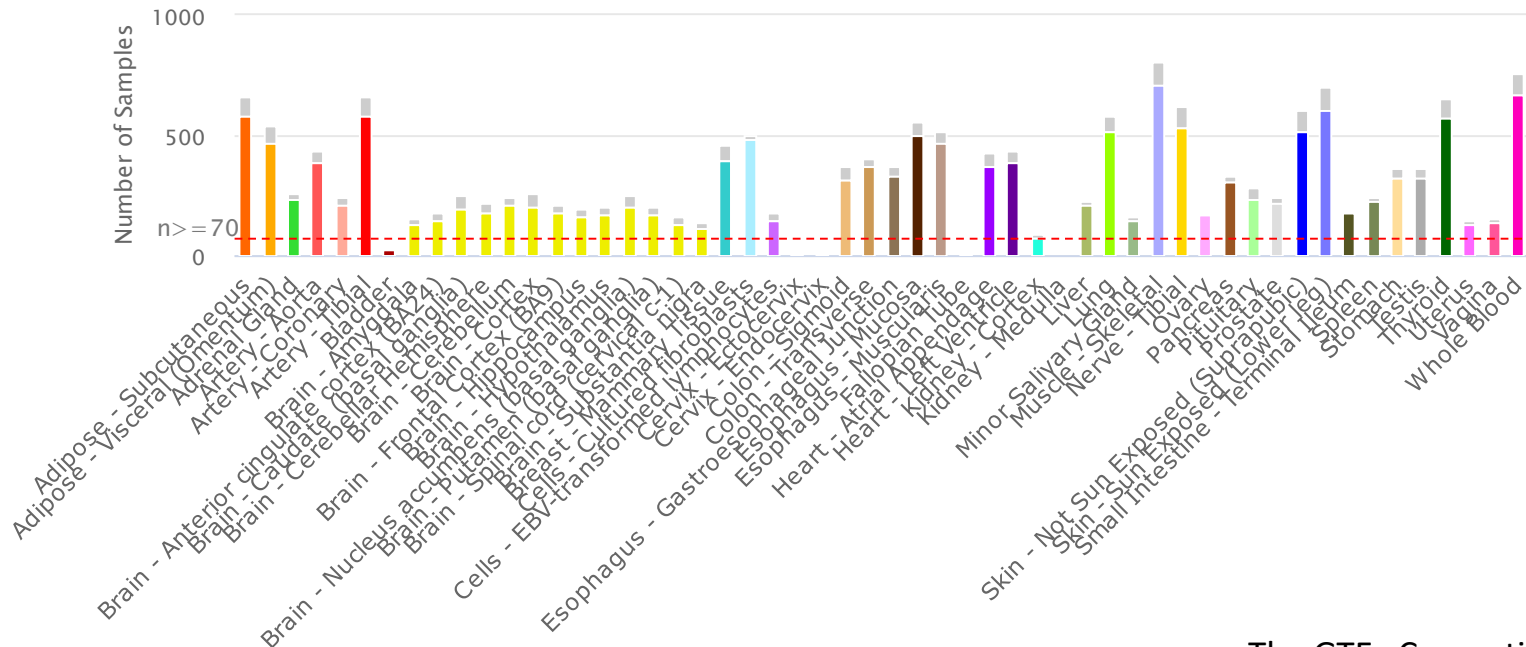Cis-eQTLs have now been identified for nearly every human gene, with numerous large studies available



Battle, Genome Research, 2014

## *trans-* effect

Trans-eQTL example

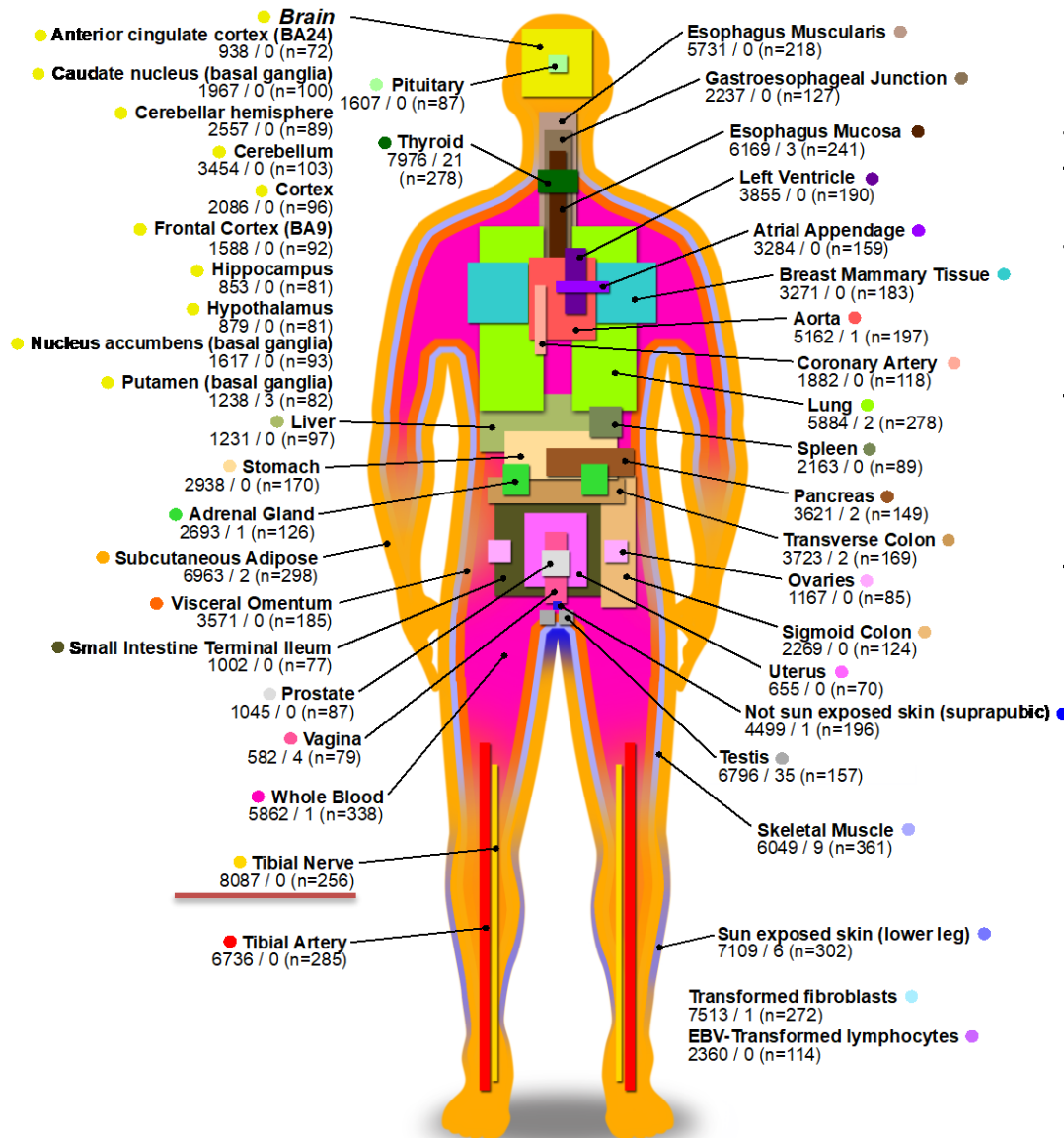The GTEx Consortium, 2017

# GTEx Project

## GTEx Consortium v8 data

- 838 genotyped donors

- 17832 gene expression samples



The GTEx Consortium

# Genetic effects across human tissues



Total unique eQTL genes:
Cis: 19,725 (FDR 5%)
Trans: 93 (FDR 10%)
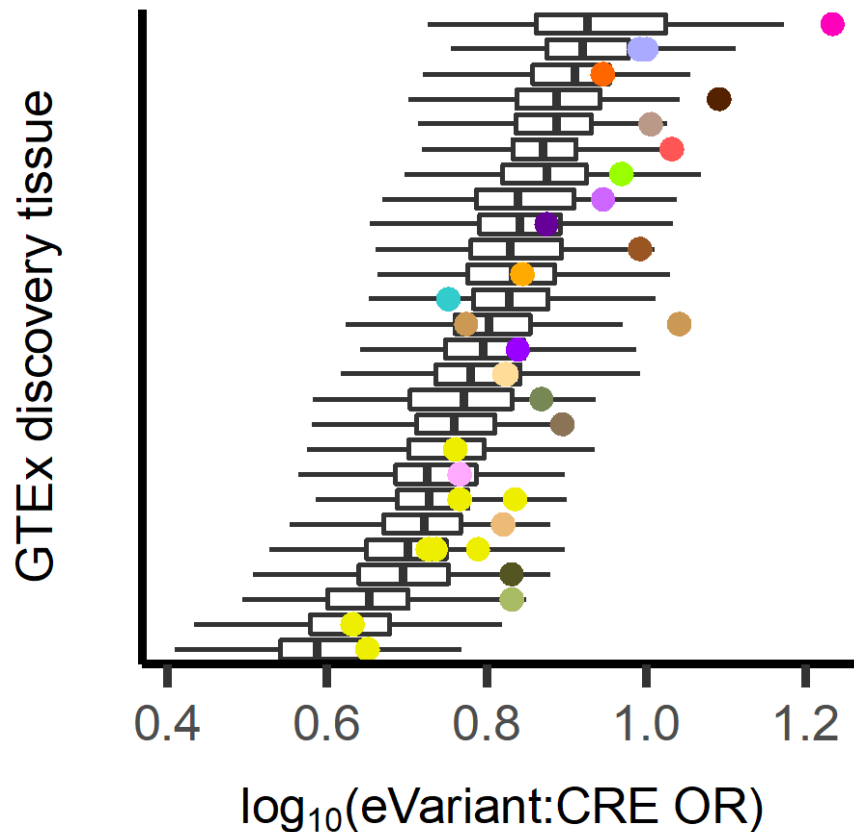
Most cis per tissue:
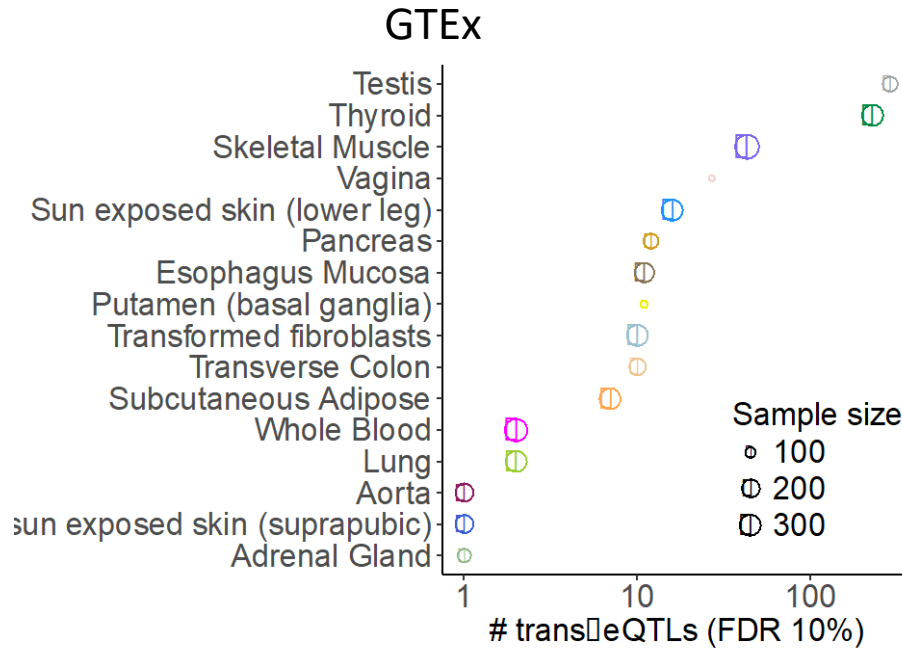8,087 Tibial nerve (N=256)
（脛神經）
Most trans per tissue:
35 Testis (N=157)

The GTEx Consortium, Nature 2017

# Characterizing eQTLs across tissues

- Cis-eQTL variants fall in tissue-specific regulatory elements (from Roadmap Epigenomics)

# Trans-eQTLs



Large studies:

Westra et al (N=5,311, using GWAS variants only)
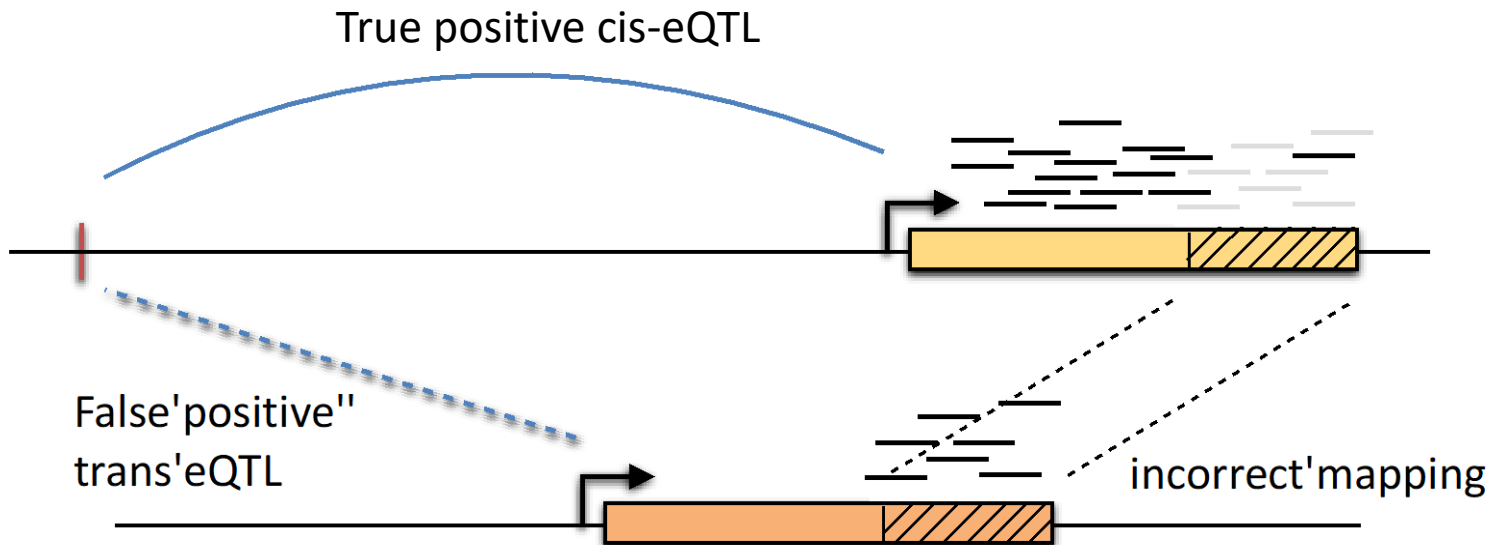ALSPAC (N=869)
MUTHER (N=850)
DGN (N=922)
Framingham (N=5257)

Studies report wildly different # hits (10s–10000s)
Replication and validation remains poor

# Challenges for trans-eQTL detection

- Power

- False positives from many sources e.g. over and under correcting confounders (Dahl et al, 2017)

- Mapping error (similar to probe cross-hybrid.)

True positive cis-eQTL

False'positive''
trans'eQTL

incorrect'mapping

# Introduction of GTEx protal

- https://gtexportal.org/home/

# 2 MANY TYPES OF QTLS

# Next generation sequencing has increased our ability to survey the transcriptome.



RNA-Seq    Montgomery, Nature 2010
Pickrell, Nature 2010

ChIP-Seq

McDaniell, Science 2010

# RNA-seq provides resolution of more QTLs

RNA-sequencing in 60 Europeans (HapMap genotypes; LCLs)

**Found 2x more expression Quantitative Trait Loci (eQTLs) and...**



**Exon-eQTLs**

**UTR Length-QTLs**

**Splicing eQTLs**

**Rare eQTLs with allele specific expression-based approaches**

# Splicing eQTL

Can investigate relative transcript ratios or reads across junctions.



Katz et al, Nature Methods, 2010

- Splicing also affected for many genes



Battle et al, Genome Research, 2014

# Epigenetic data



- ENCODE, Roadmap Epigenomics
- Regulatory elements: promoters, enhancers
- Transcription factor binding sites
- CpG sites
- ChromHMM

ENCODE Project Consortium. Plos Biology 2011.

# Epigenetic data informs heritability

LD score regression, related approaches partition $h^2$



Large scale epigenetic data (Roadmap, ENCODE) enable analysis, indicate contribution of gene regulation

Figure from Finucane, NG, 2015

# Ommigenic model

- Most/all expressed genes in disease-relevant cell types affect trait



- Highlights potential role of eQTLs, trans effects

Boyle et al., Cell, 2017

# Advantages of ASE

- Test within an individual allelic imbalance, given one has sufficient reads.

# Using ASE to detect GWAS signals driven by <u>multiple</u> causal variants

**GWAS variant genotype**



LACK OF ASE FOR HOMS

ASE ⟹ ABUNDANT ASE FOR HETS

LACK OF ASE FOR HOMS

**Tests functional differences between alleles in population**

Lucia Conde et al, *AJHG*, 2013

# Coloc: A Bayesian test for colocalization of pairs of association signals

H1 is the hypothesis that there is only an eQTL signal at a locus

H2 is the hypothesis that there is only a GWAS signal at a locus.

H3 is the hypothesis that there are two independent eQTL and GWAS signals in linkage.

H4 is the strong hypothesis that the same SNP (not just the locus) is responsible for both the GWAS and eQTL.

# GWAS eQTL colocalization

- **Coloc**
- **eQTpLot**
- hypercoloc

# Examples of H3 and H4



On the left, the profile of association at the FRK locus with LDL (top) is very different from that with *FRK* expression.

H3 is the supported hypothesis.

On the right, even though there are two different peak SNPs, they are in the same strong LD region and the profiles are almost the same for Total Cholesterol and *Soc1* expression.

H4 is the supported hypothesis.

Bayesian analysis evaluate each H relative to the other four and generates a confidence level for the most likely one.

# Coloc results
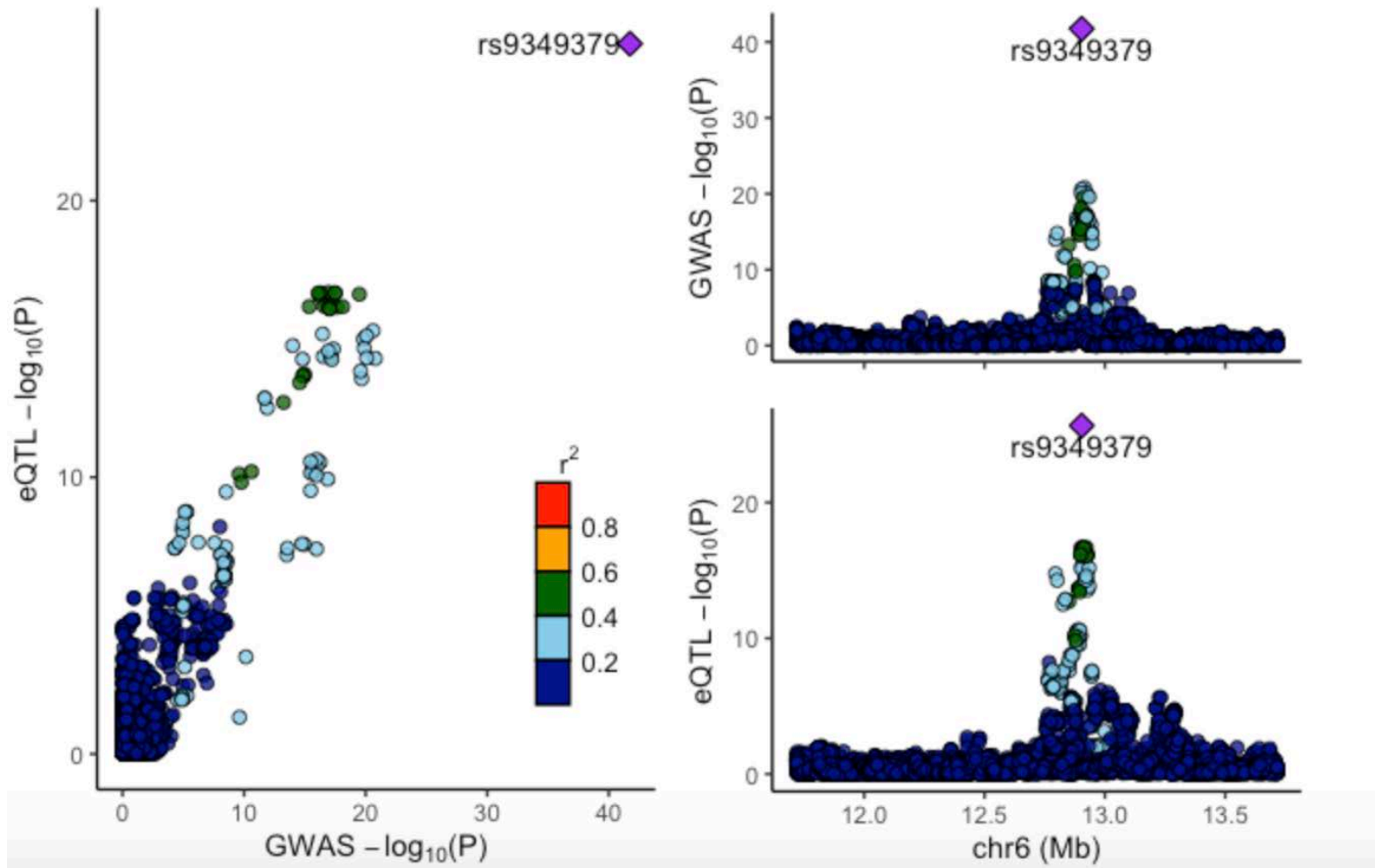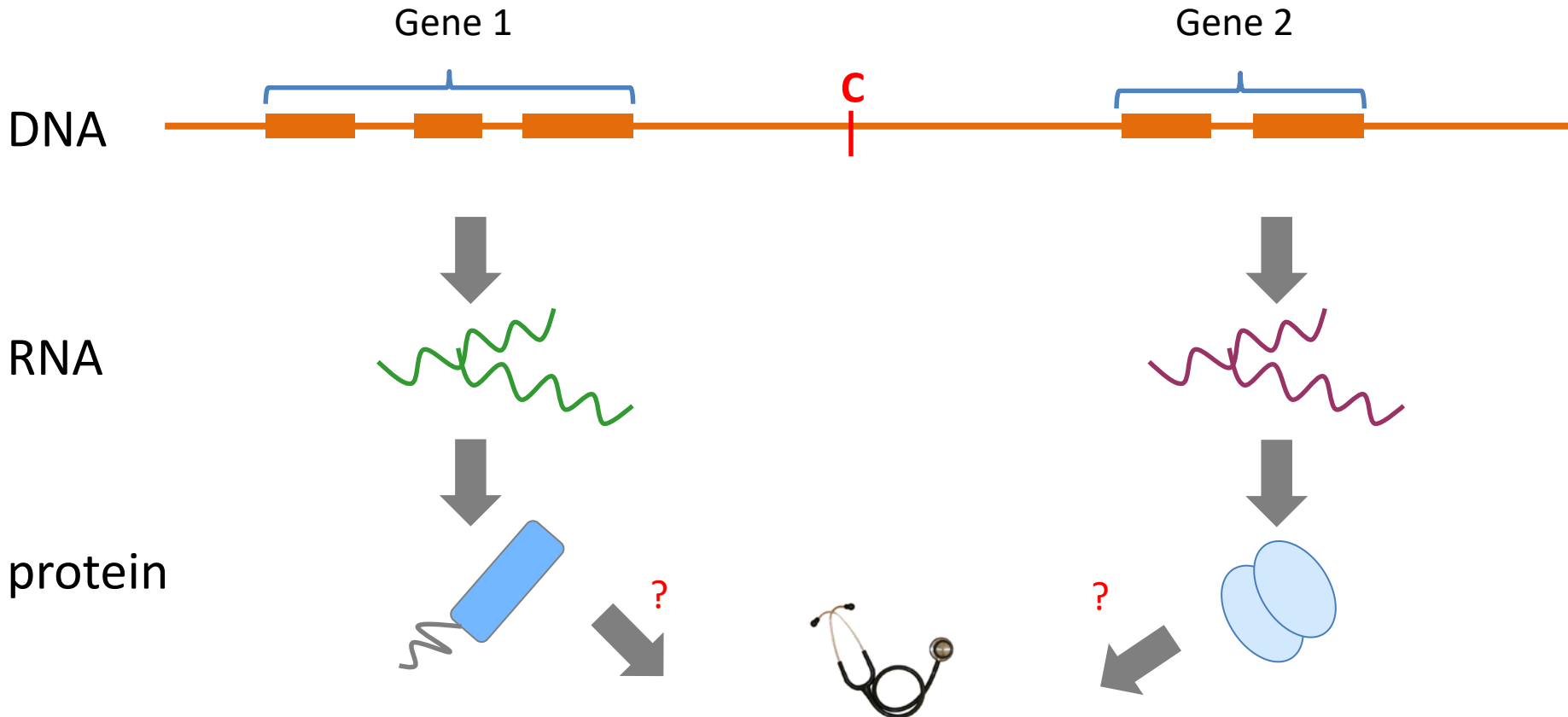
# eQTLs and complex disease genetics

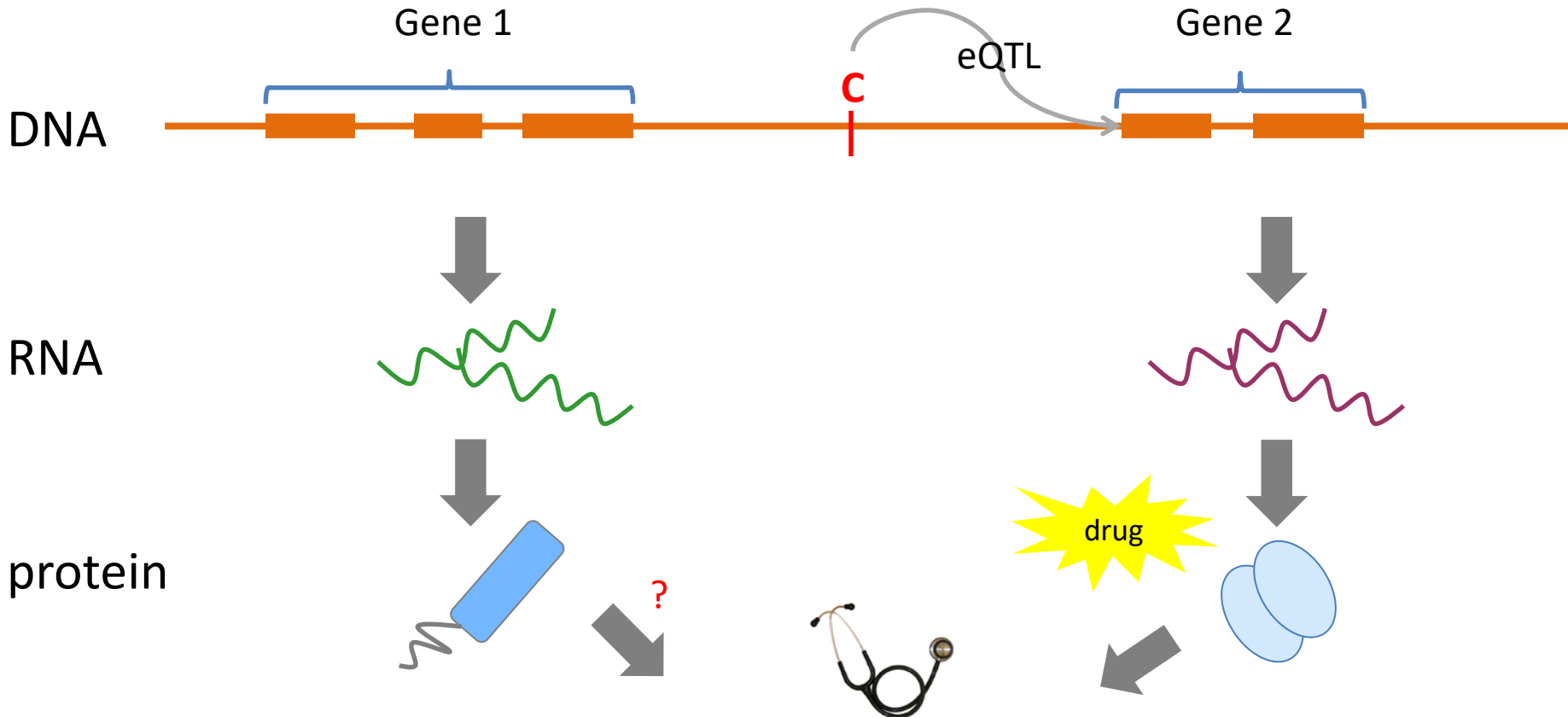Help interpret GWAS variants (especially non-coding):

- understand mechanism
- guide interventions

# eQTLs and complex disease genetics
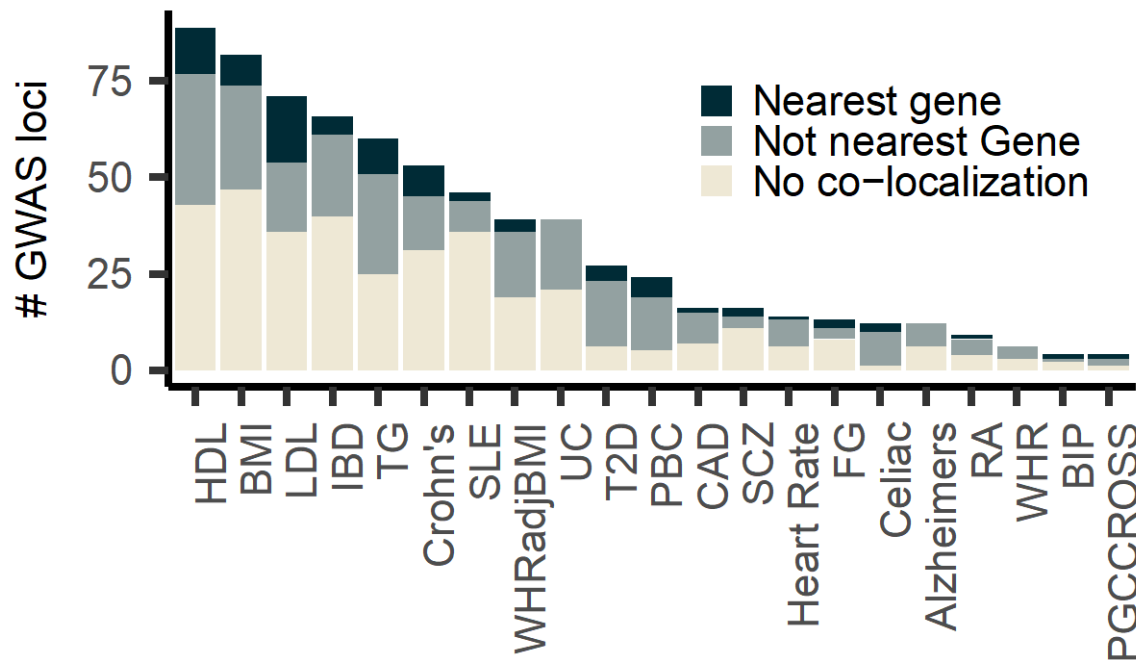
Help interpret GWAS variants (especially non-coding):
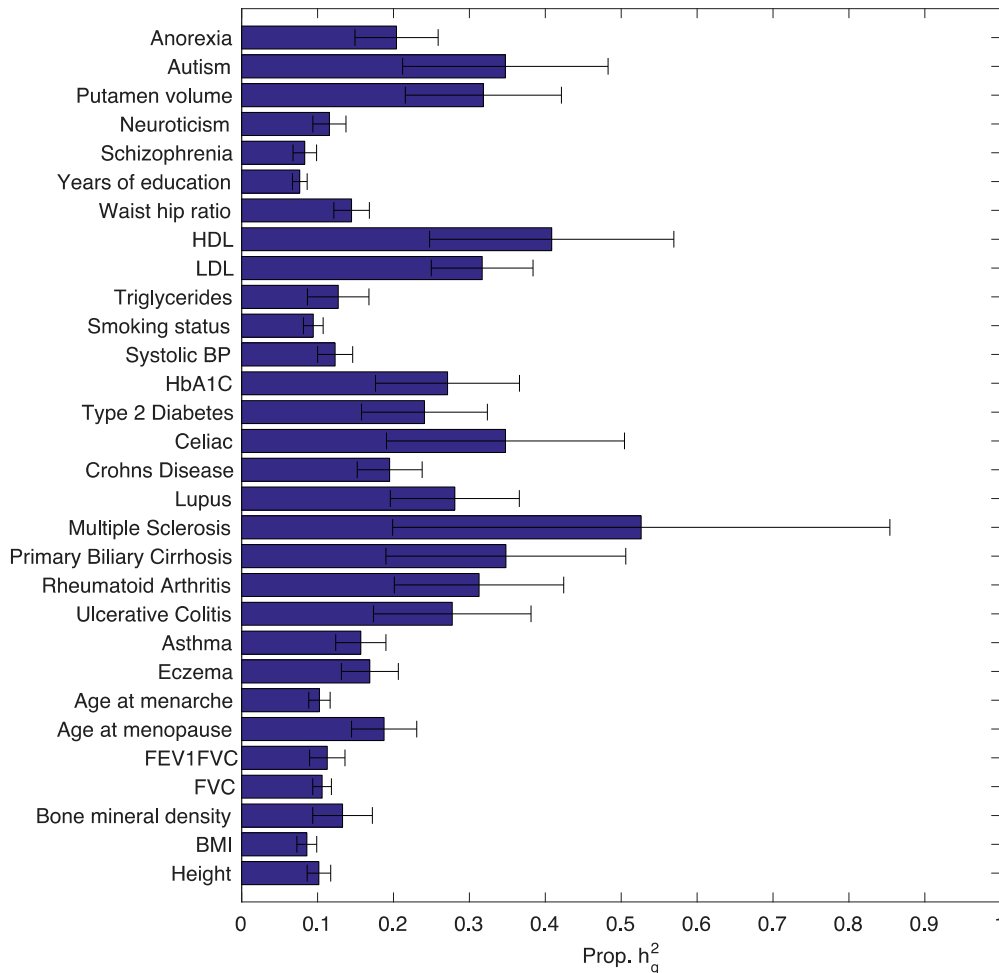
- understand mechanism
- guide interventions

# eQTLs and complex disease genetics

52% of genetic variants associated with human disease co-localize with an eQTL

# eQTL data informs heritability

GE co-score regression indicates cis-eQTLs explain mean 21% of $h^2$ across a set of complex traits
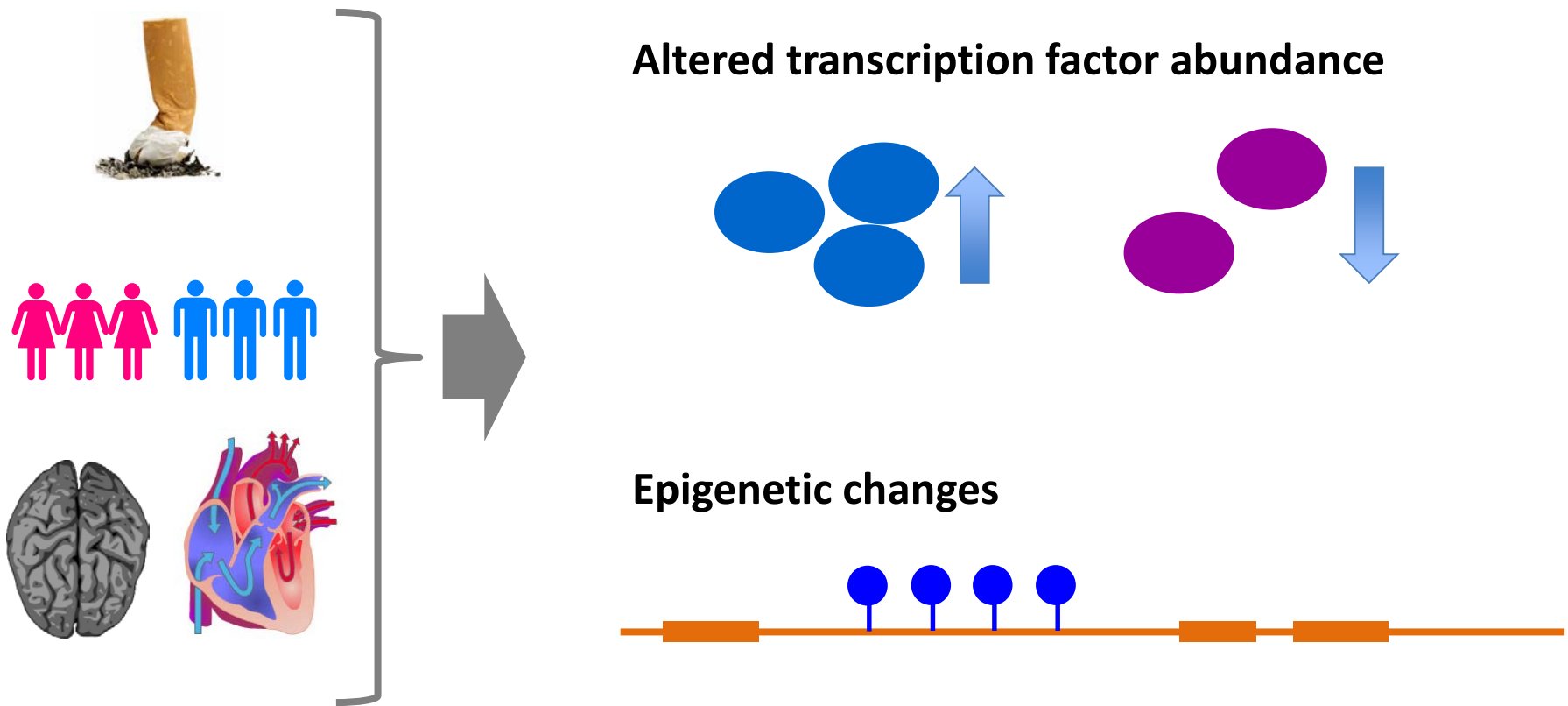


O'Connor et al. bioRxiv, 2017

# 3. Complex effects of genetic variation on gene expression

# What are we missing?

- Most studies are done on steady-state total expression measurements at a single adult or post-mortem time point

- Disease-relevant states include different developmental stages, environmental exposures, cell types

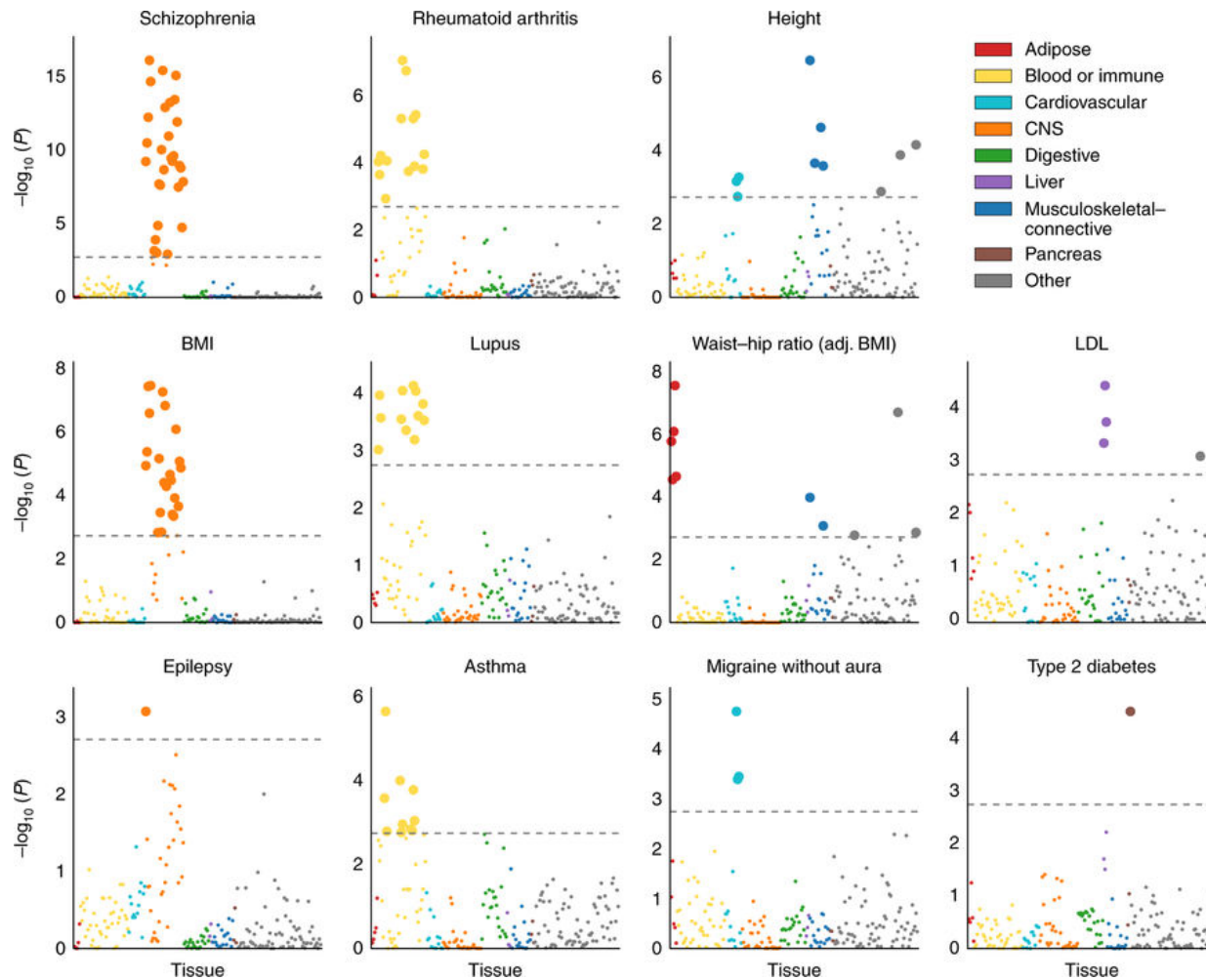- Other variant classes and regulatory effects

# GTEx tissue-specificity of cis and trans



Median spearman correlation=0.547 (cis)

Median spearman correlation=0.138 (trans)

trans ρ

0.4
0.3
0.2
0.1

cis ρ

0.8
0.7
0.6
0.5
0.4

Brain

Muscle/heart

Skin

Artery

Adipose/breast

Trans eQTLs appear more highly tissue-specific than cis-eQTLs

# Tissue specificity and heritability



From Finucane et al, NG, 2018

# Rare variants

## Recent work emphasizes importance of rare variation in driving extreme expression levels



Li et al, Nature, 2017

# Rare variants

Preprint (Hernandez et al 2017) suggests rare variants explain a large fraction of heritability of gene expression

# 4. Conclusions

# Why delve deeper into expression?

- Help determine when and how much to invest in WGS, expression, epigenetic data

- To continue understanding implicated
  - Genes
  - Tissue and cell types
  - Epigenetic and other regulatory mechanisms

- Challenges and caveats
  - Ambiguity: many variants affect multiple genes
  - Interpretability: missing relevant cell types
  - Power: trans-eQTLs also require large sample sizes

# 5. eQTL tool demo

matrixEQTL

# Statistical Clinics



1. It is a special service for free for researchers in Academia Sinica.

2. The service is offered at 14:00-16:00, Thursday, in Institute of Statistical Science Building, Room 401.

3. You are welcome to apply for the service.
http://disc.stat.sinica.edu.tw/statistical-clinic-service-appointment/

# Thank you