

6.874, 6.802, 20.390, 20.490, HST.506

Computational Systems Biology

Deep Learning in the Life Sciences

# Lecture 17 - Systems genetics: Deep Learning, eQTLs, Polygenicity, Heritability, LMMs, LDSC, PRS, Networks

Prof. Manolis Kellis



Massachusetts  
Institute of  
Technology

<http://mit6874.github.io>

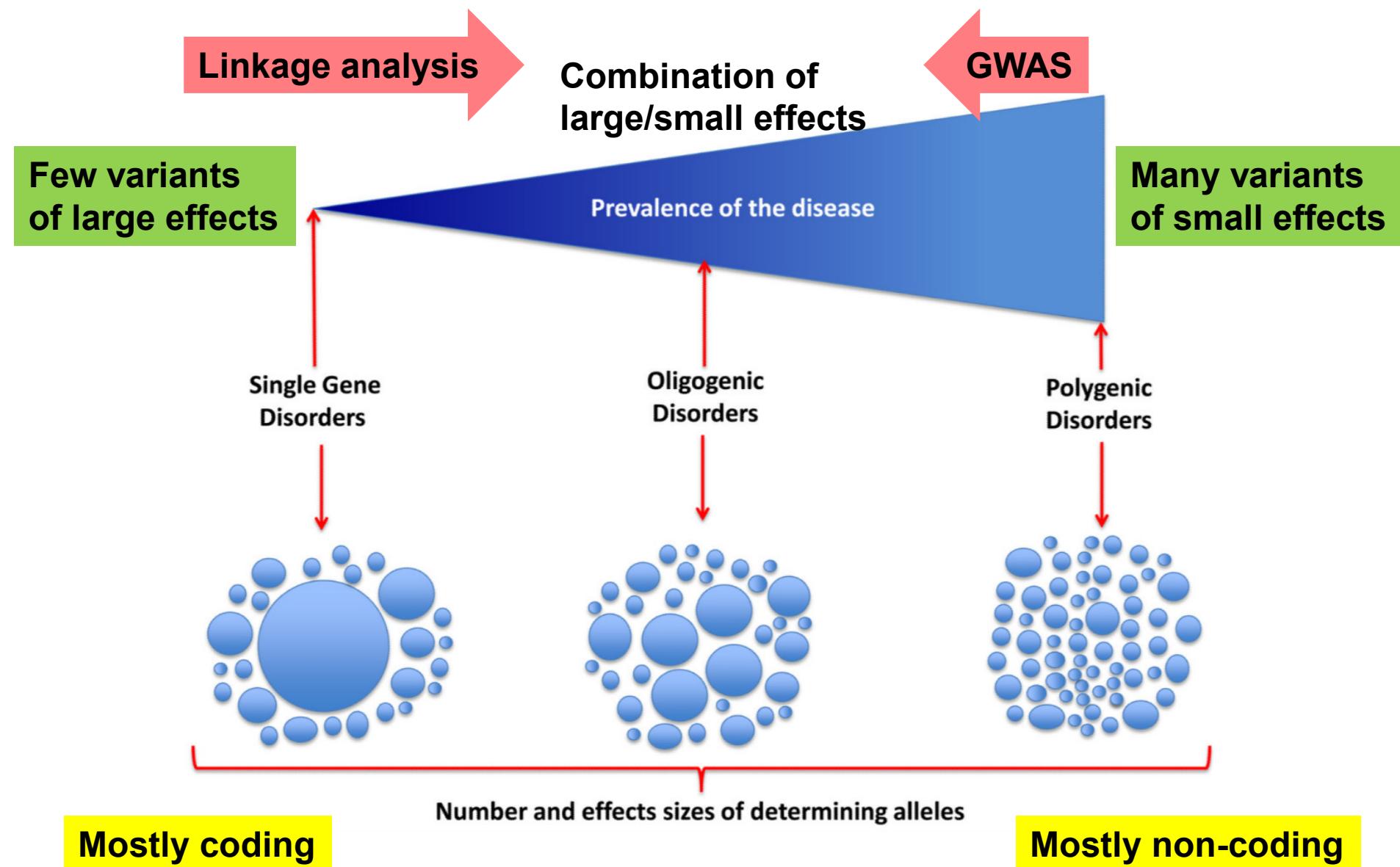
**Slides credit:** Yongjin Park, Abhishek Sarkar,  
Mark Daly, David Gifford, et al

# Today: Deep Learning for Human Genetics and Disease

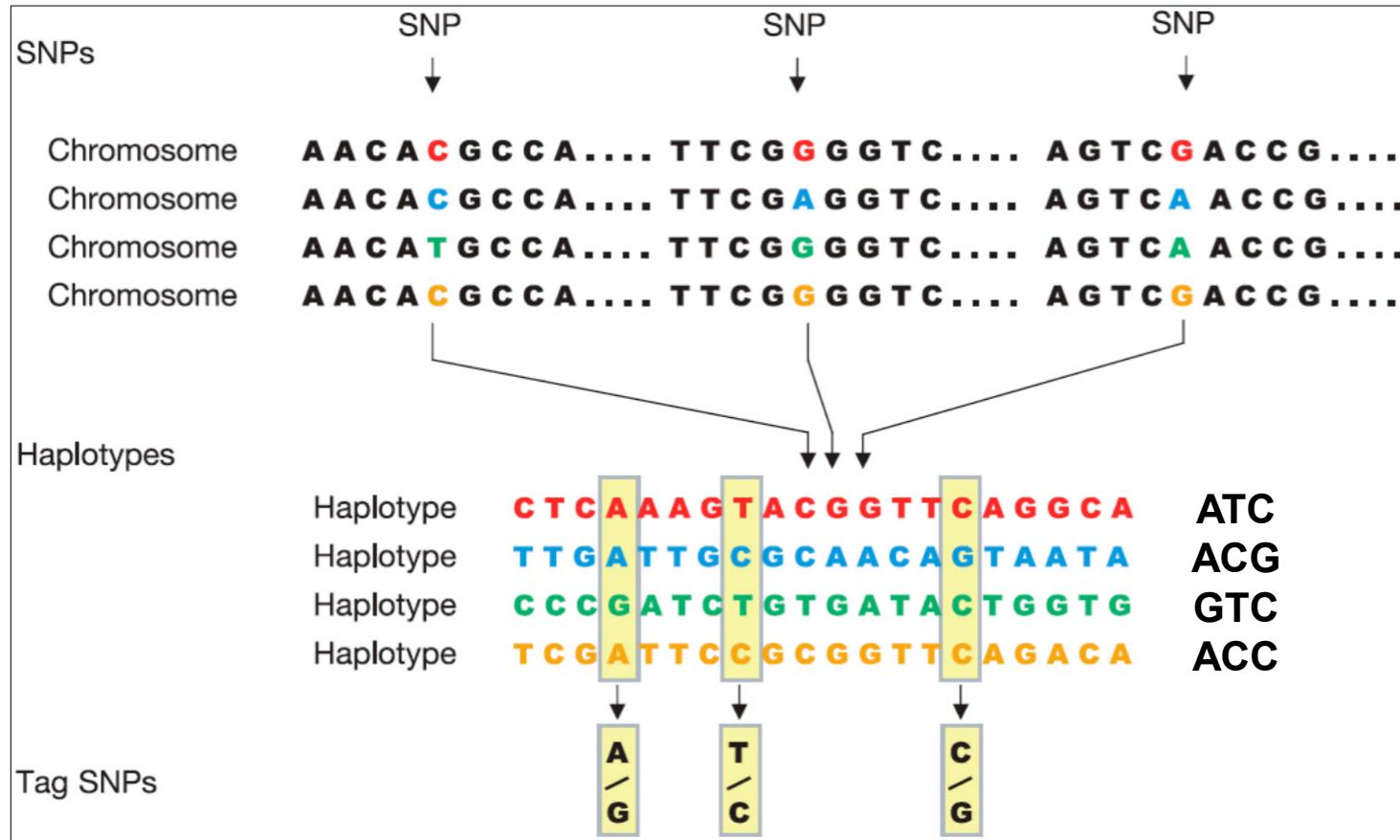
1. Review: GWAS, fine-mapping, Bayesian variant prioritization
2. Deep Learning for GWAS: calling SNPs, prioritize function
3. eQTLs/Mediation: intermediate molecular phenotypes
4. Linear Mixed Models (LMMs) for GWAS and for eQTL calling
5. Polygenic Risk Scores (PRS): summing over many variants
6. Heritability: definition(s), missing heritability, partitioning
7. LD SCore regression (LDSC) for fast heritability partitioning
8. Polygenic/Omnigenic disease models: core vs. periphery
9. Disease gene networks from GWAS evidence boosting

# **1. Review: GWAS, fine-mapping, Bayesian methods for variant prioritization**

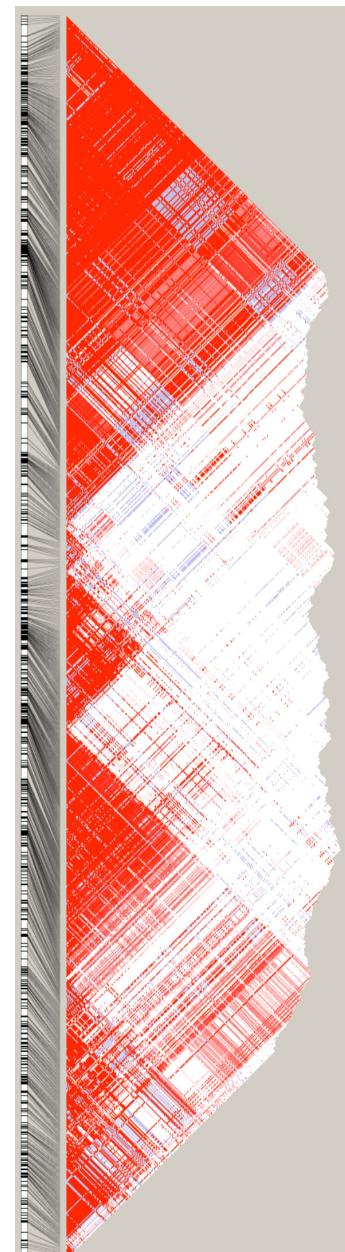
# Monogenic vs. oligogenic vs. polygenic disorders



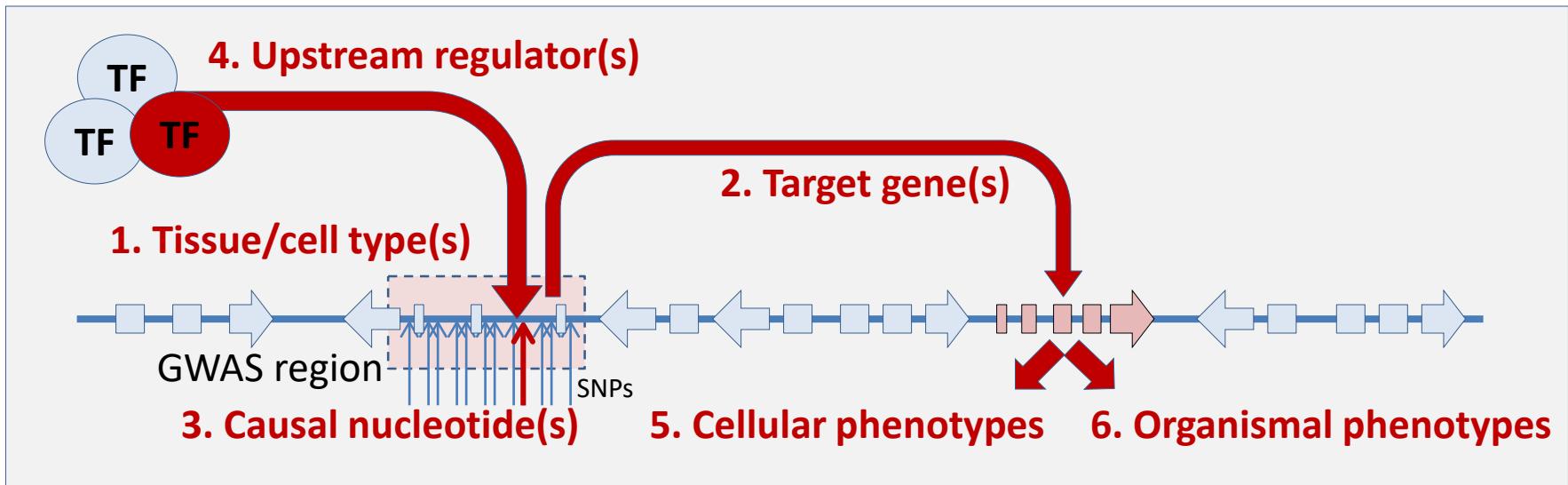
# Common variants (SNPs) live in Haplotypes



- Common SNPs only once every 1000 nucleotides or so
- These are co-inherited, so only need to profile a subset
- Markers selected for haplotype profiling are “tag” SNPs



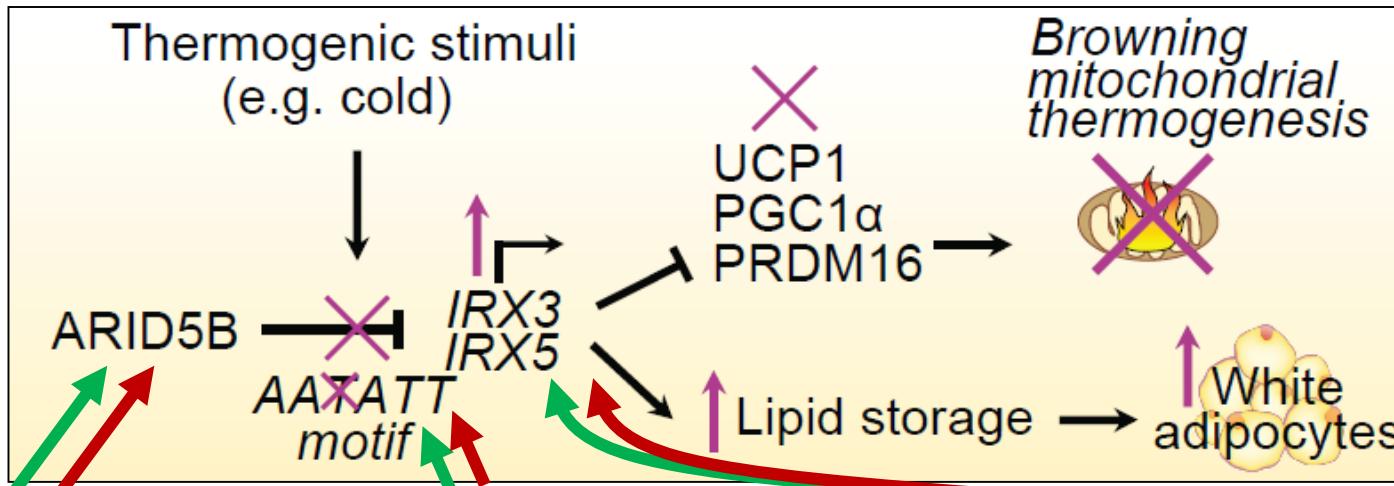
# Dissecting non-coding genetic associations



1. Establish relevant **tissue/cell type**
2. Establish downstream **target gene(s)**
3. Establishing **causal** nucleotide variant
4. Establish upstream **regulator** causality
5. Establish **cellular** phenotypic consequences
6. Establish **organismal** phenotypic consequences

Goal:  
Apply these to  
the FTO locus  
in obesity

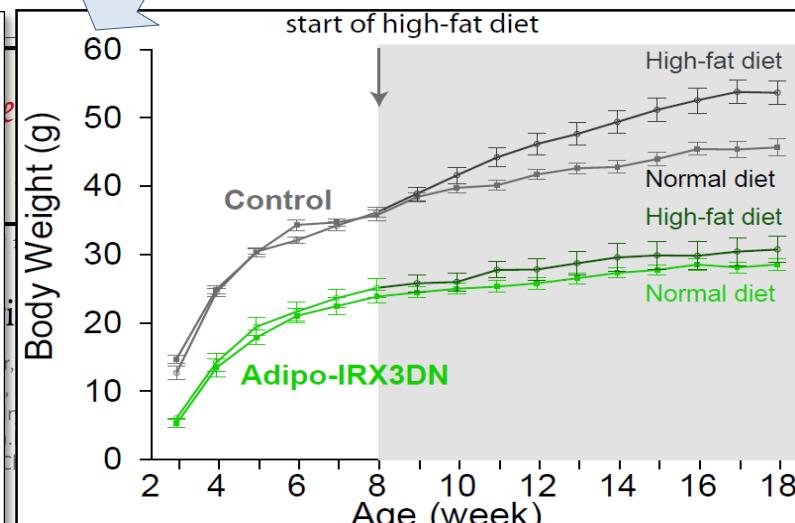
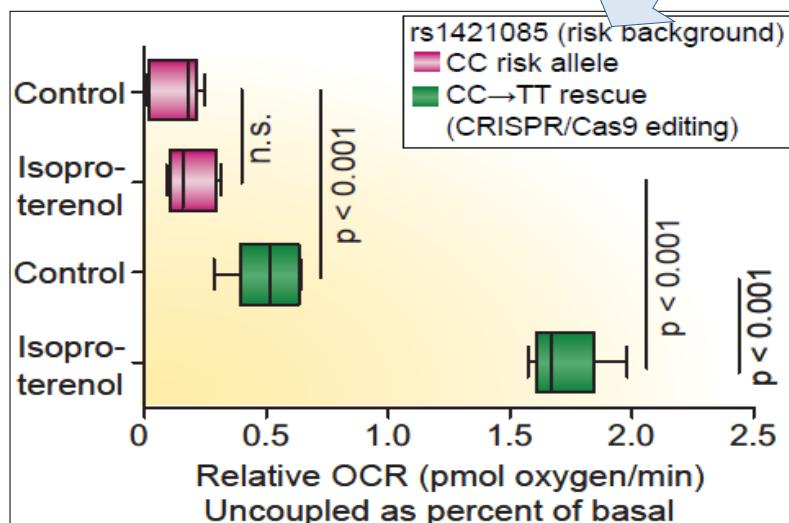
# Manipulate circuitry → reverse disease phenotypes



Incr. ARID5B → Lean  
Decr ARID5B → Obese

C-to-T → Lean  
T-to-C → Obese

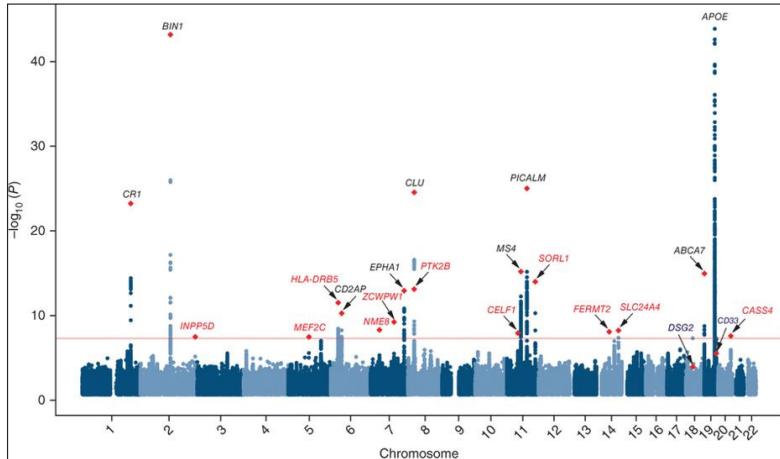
Decrease IRX3, IRX5 → Lean  
Increase IRX3, IRX5 → Obese



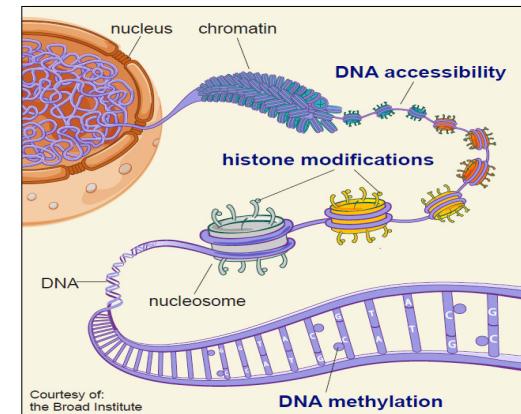
CRISPR-edit human fat cells  
→ able to burn calories again

IRX3 KD → Burn calories in their sleep  
→ 54% weight loss. Can't gain weight

# Dissect mechanisms of disease-associated regions

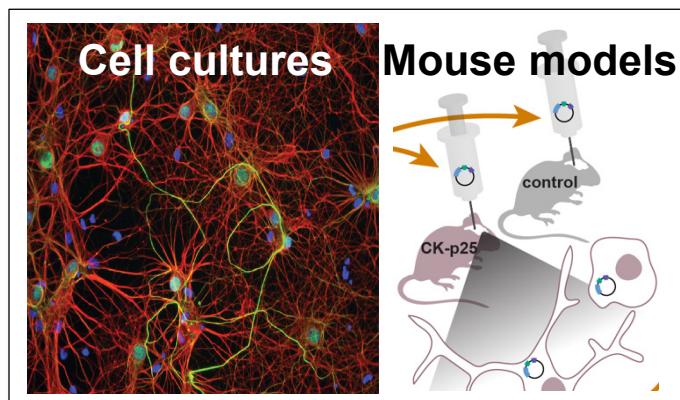
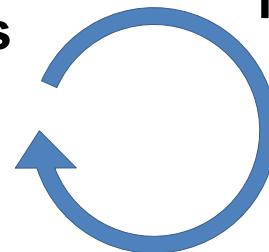


1. Disease genetics reveals common + rare variants/regions

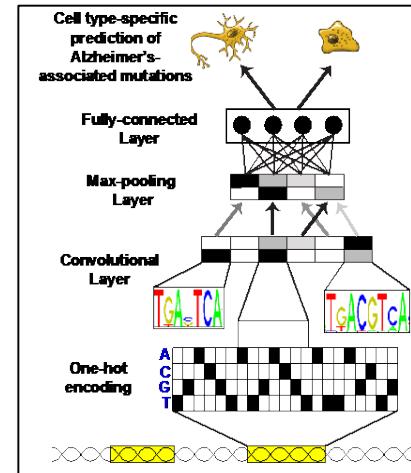
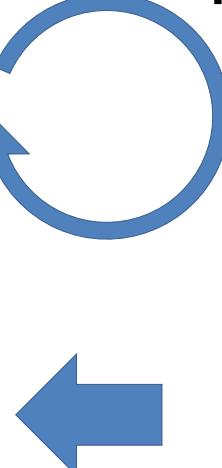


2. Profile RNA + Epigenome in healthy + disease samples

5. Disseminate results

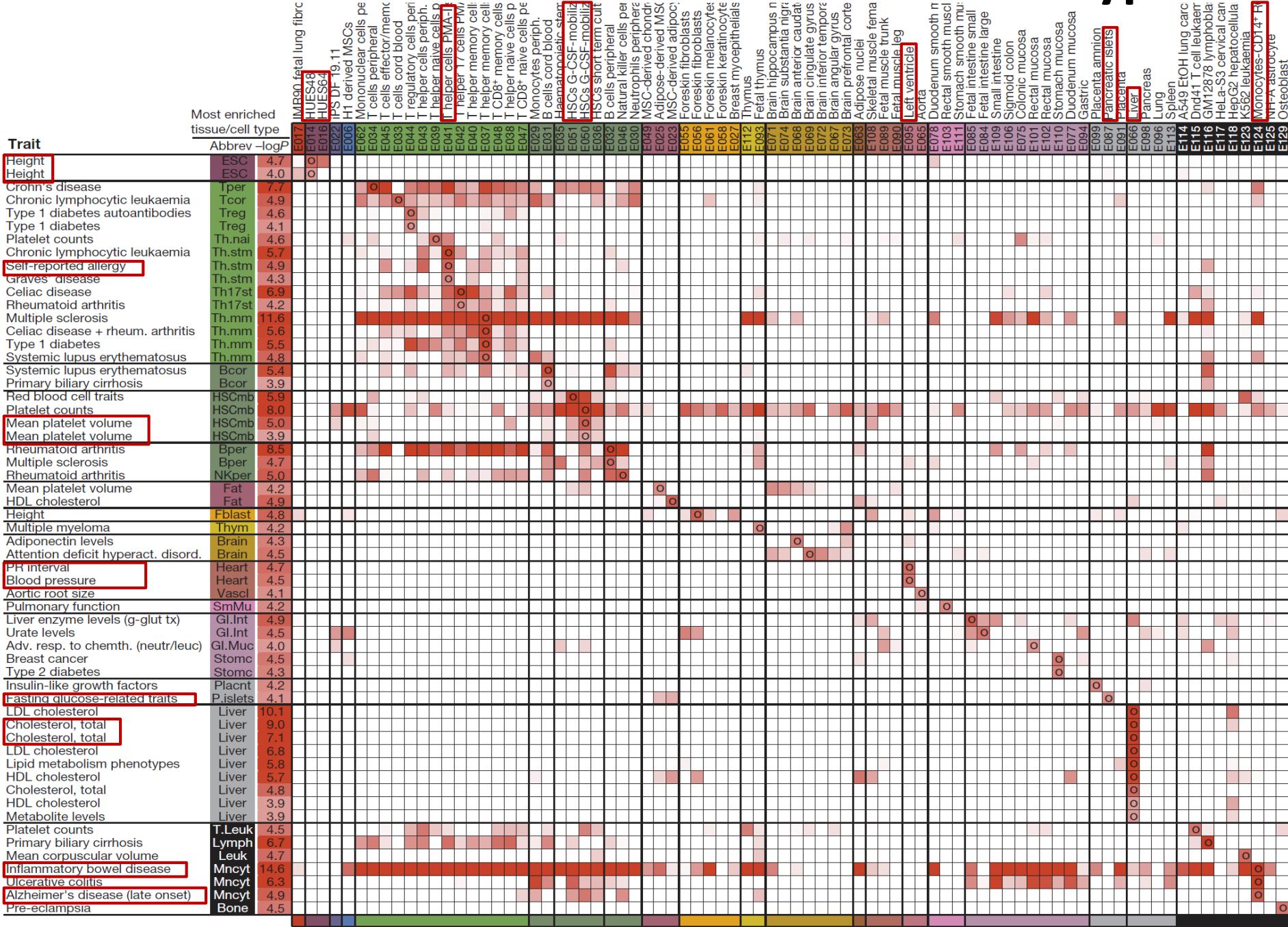


4. Validate predictions in human cells + mouse models

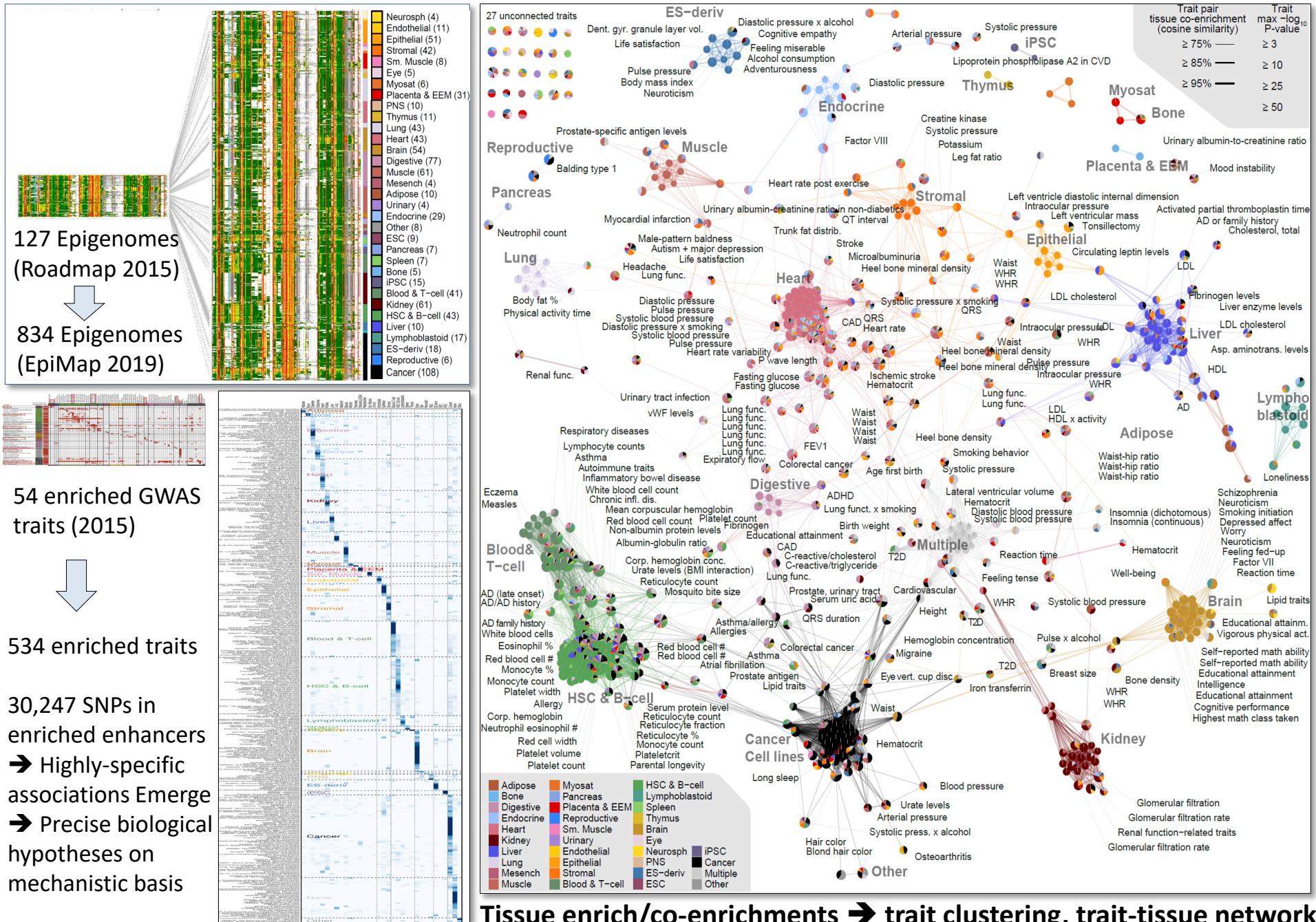


3. Integrate data to predict driver genes, regions, cell types<sup>8</sup>

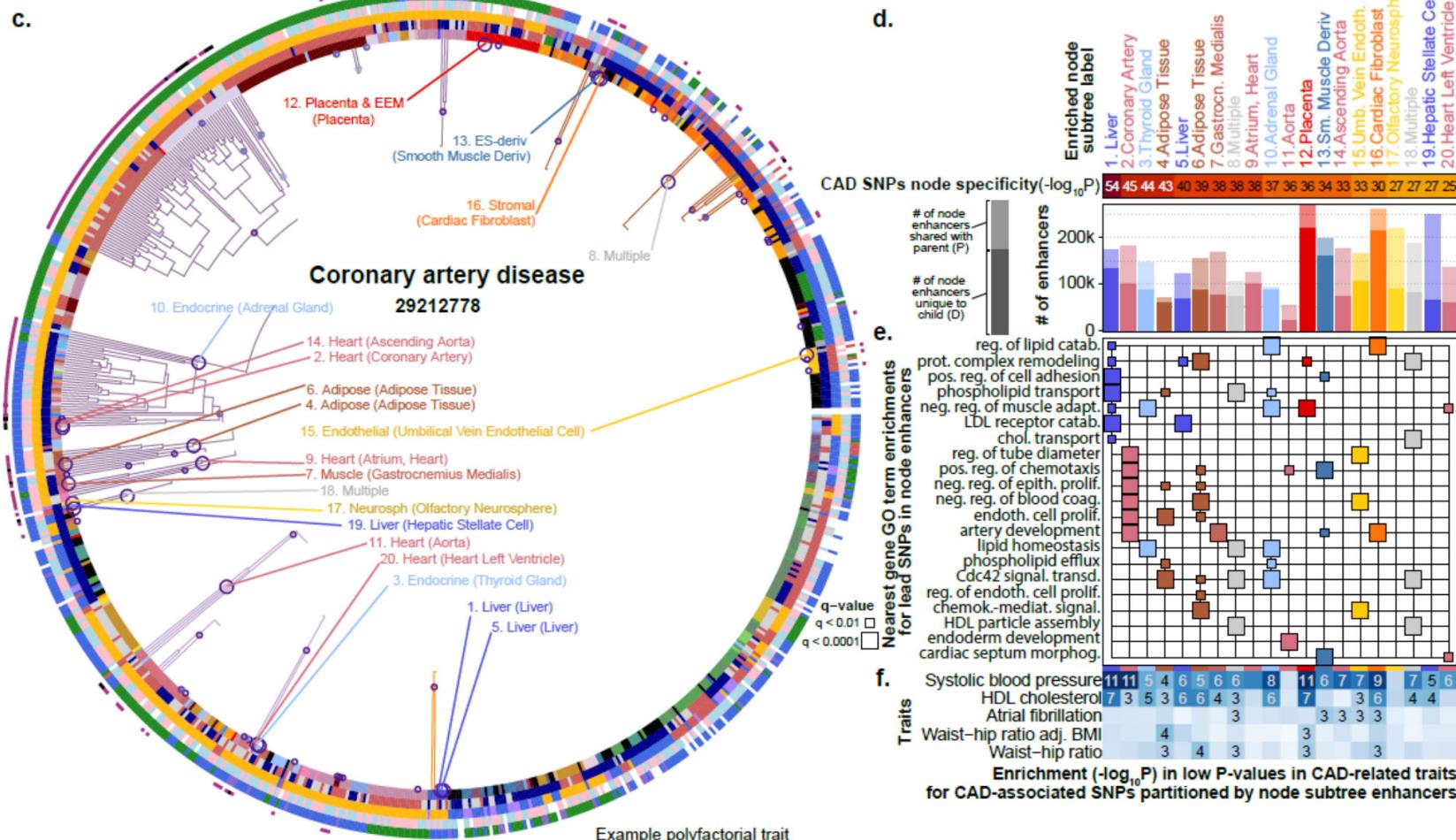
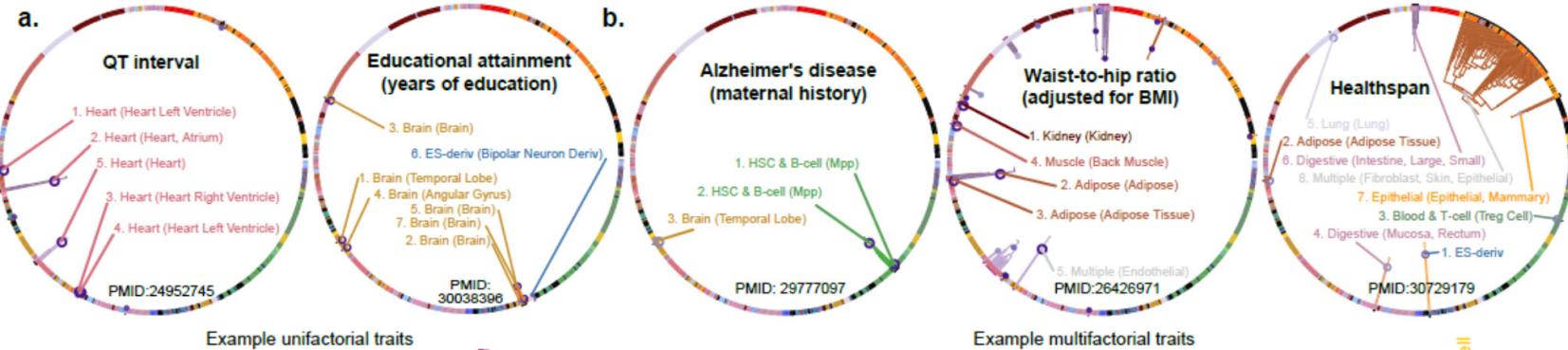
# Disease hits in enhancers of relevant cell types



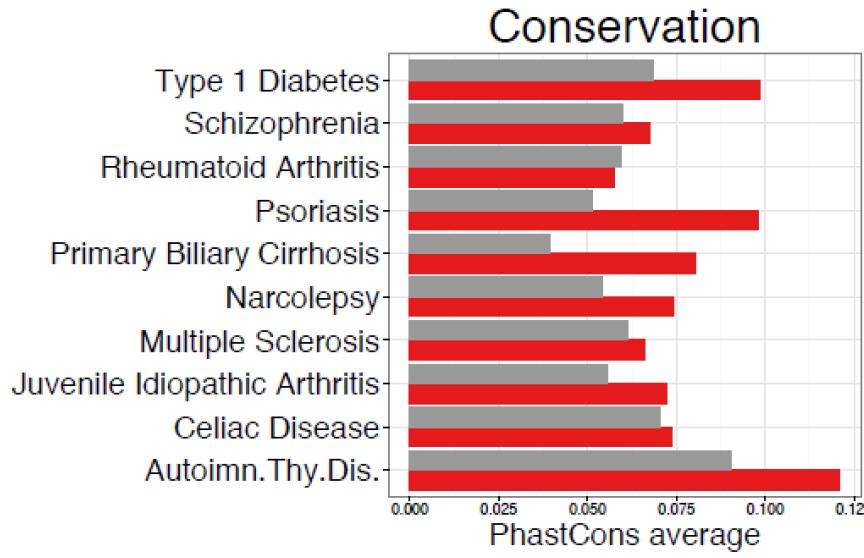
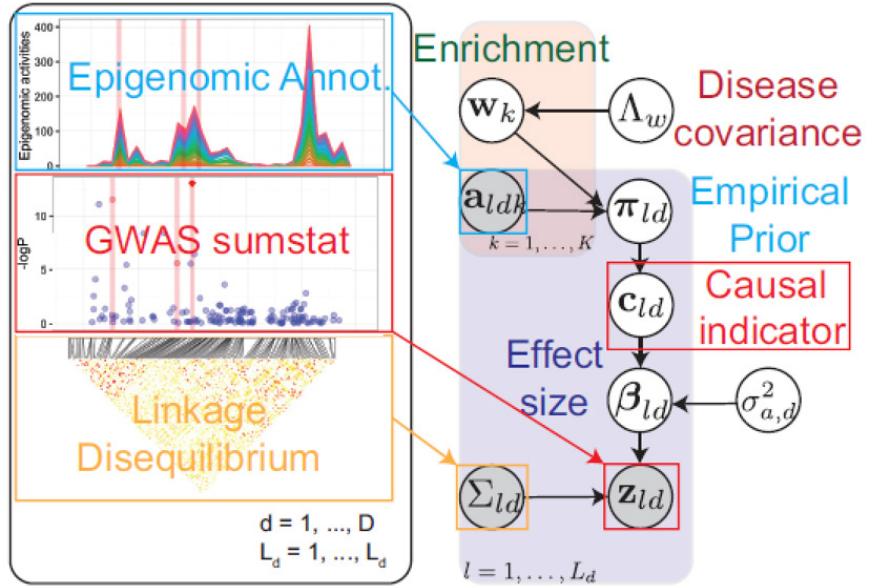
**Scale up: 834 tissue/cell types → 30k GWAS SNPs in 534 traits**



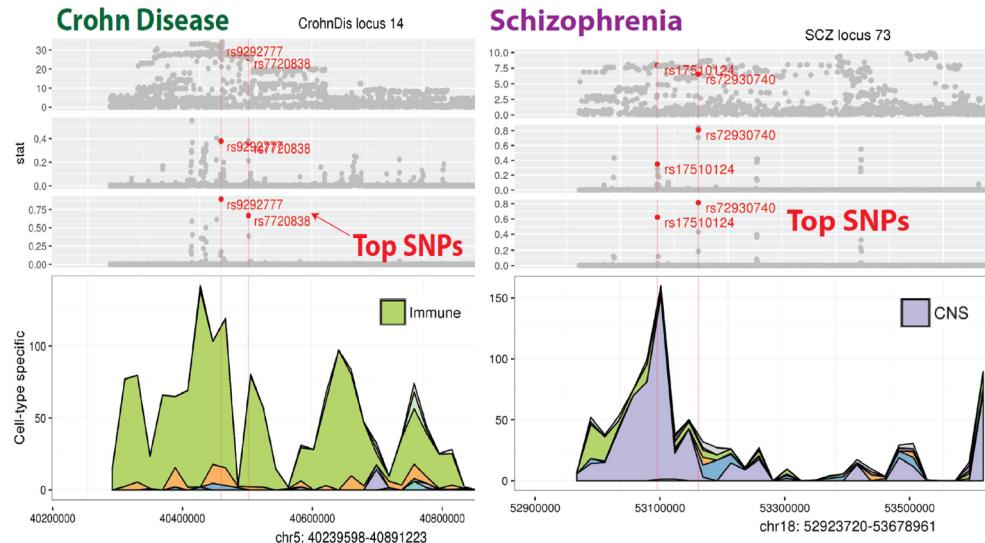
# Partitioning complex traits across multiple tissues of action



# Bayesian fine-mapping: Predict causal variant and cell type

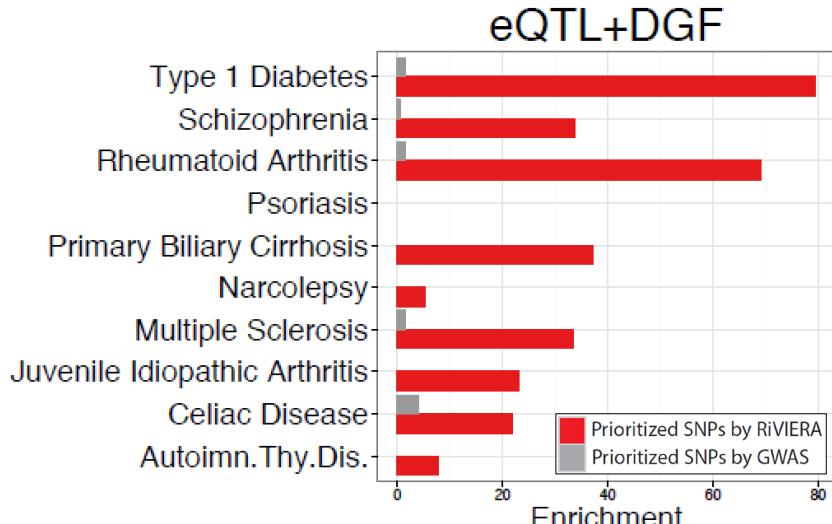


## RiVIERA: multi-trait GWAS integration



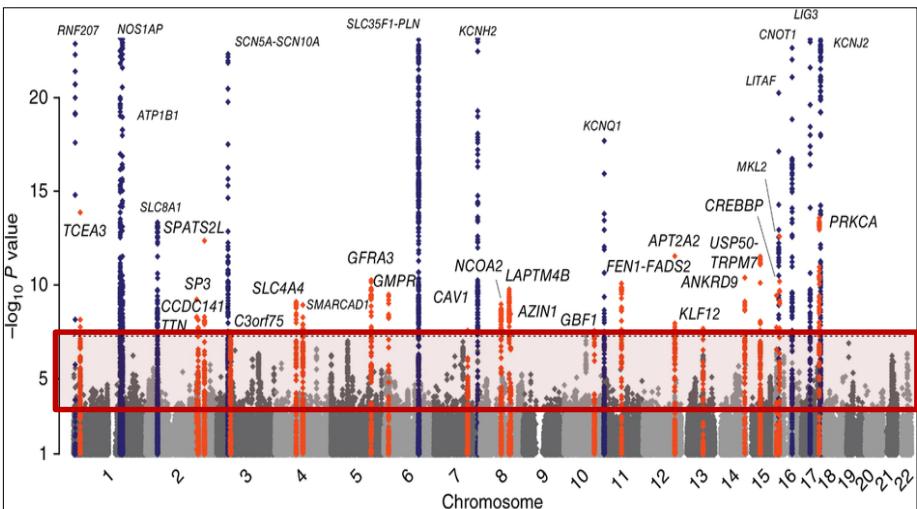
Predict causal variants and cell types

## Capture conserved elements



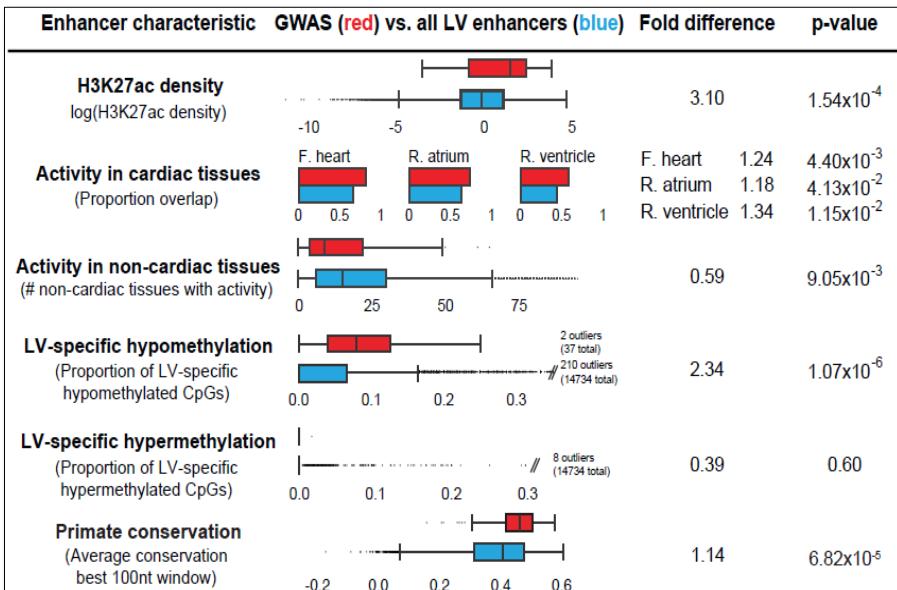
Capture eQTLs from GTEx

# Combine GWAS+Epig to find new target genes/SNPs



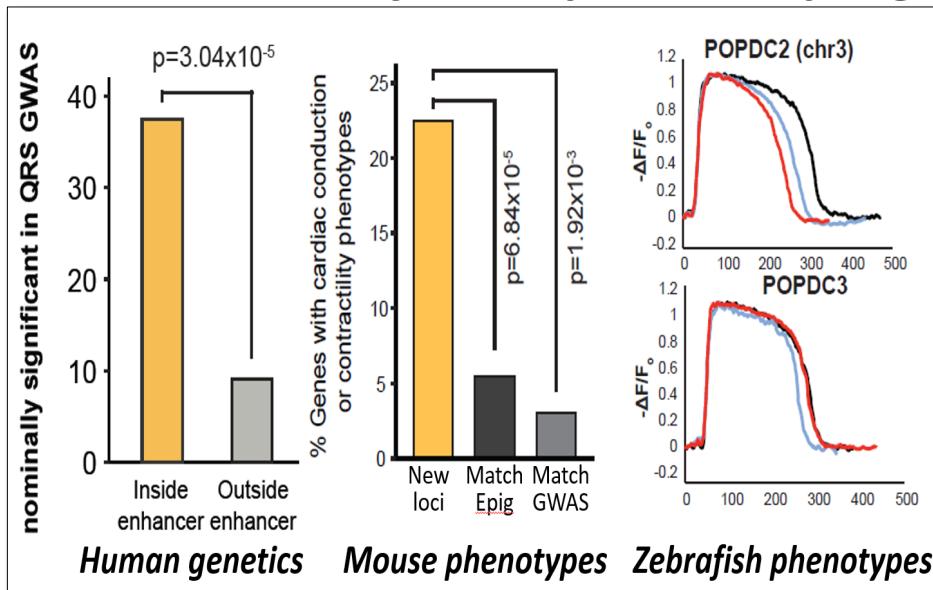
Lead SNP	p-value	Enhancer	1. Luciferase reporter	2. 4C-seq interactions
rs1886512	$4.30 \times 10^{-8}$	chr13:74,520,000-74,520,400	0.015	No interactions
rs1044503	$5.13 \times 10^{-7}$	chr14:102,965,400-102,972,000	$4.70 \times 10^{-9}$	CINP, RCOR1
rs10030238	$6.21 \times 10^{-7}$	chr4:141,807,800-141,809,600 chr4:141,900,800-141,908,000	$1.35 \times 10^{-14}$ -	RNF150 RNF150
rs6565060	$1.52 \times 10^{-5}$	chr16:82,746,400-82,750,800	$5.00 \times 10^{-3}$	No interactions
rs3772570	$1.73 \times 10^{-5}$	chr3:148,733,200-148,738,600	0.67	-
rs3734637	$2.23 \times 10^{-5}$	chr6:126,081,200-126,081,800	$1.06 \times 10^{-4}$	HDDC2
rs1743292	$6.48 \times 10^{-5}$	chr6:105,706,600-105,710,200 chr6:105,720,200-105,723,000	$3.20 \times 10^{-4}$ -	BVES, POPDC3 BVES, POPDC3
rs11263841	$6.87 \times 10^{-5}$	chr1:35,307,600-35,312,200	0.22	GJA4, DLGAP3
rs11119843	$7.14 \times 10^{-5}$	chr1:212,247,600-212,248,600	0.031	-
rs6750499	$7.37 \times 10^{-5}$	chr2:11,559,600-11,563,000 (split into two 2kb fragments)	0.54 $3.26 \times 10^{-7}$	ROCK2
rs17779853	$7.73 \times 10^{-5}$	chr17:30,063,800-30,066,800	$4.33 \times 10^{-3}$	No interactions

Prioritize sub-threshold loci ( $<10^{-4}$ )



Machine learning predictive features

Validate new enhancers:  
allelic activity, enh-prom looping



Validate new genes in hum/mou/zb

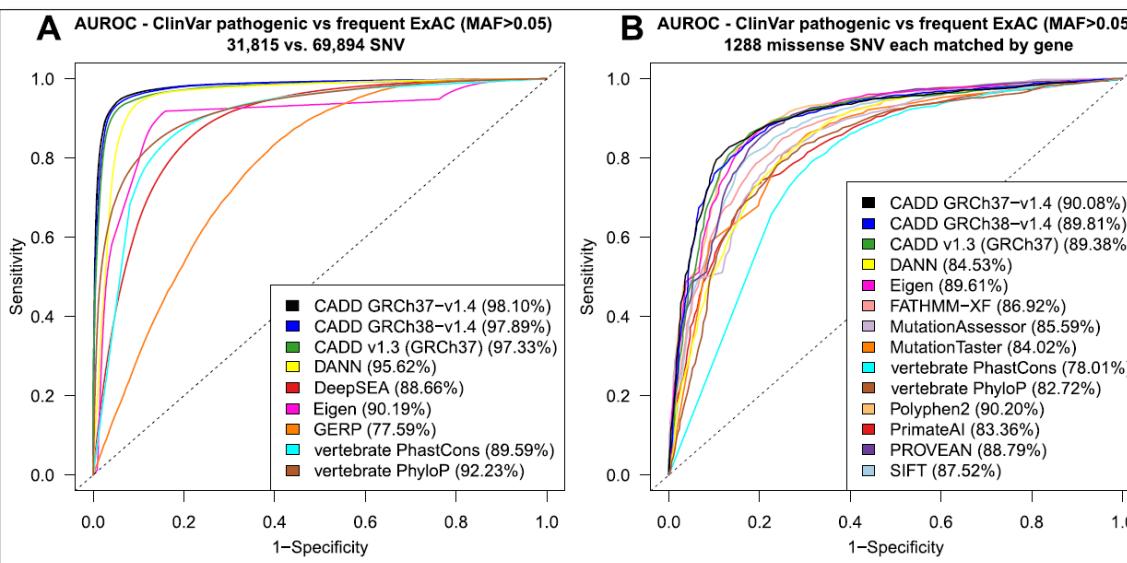
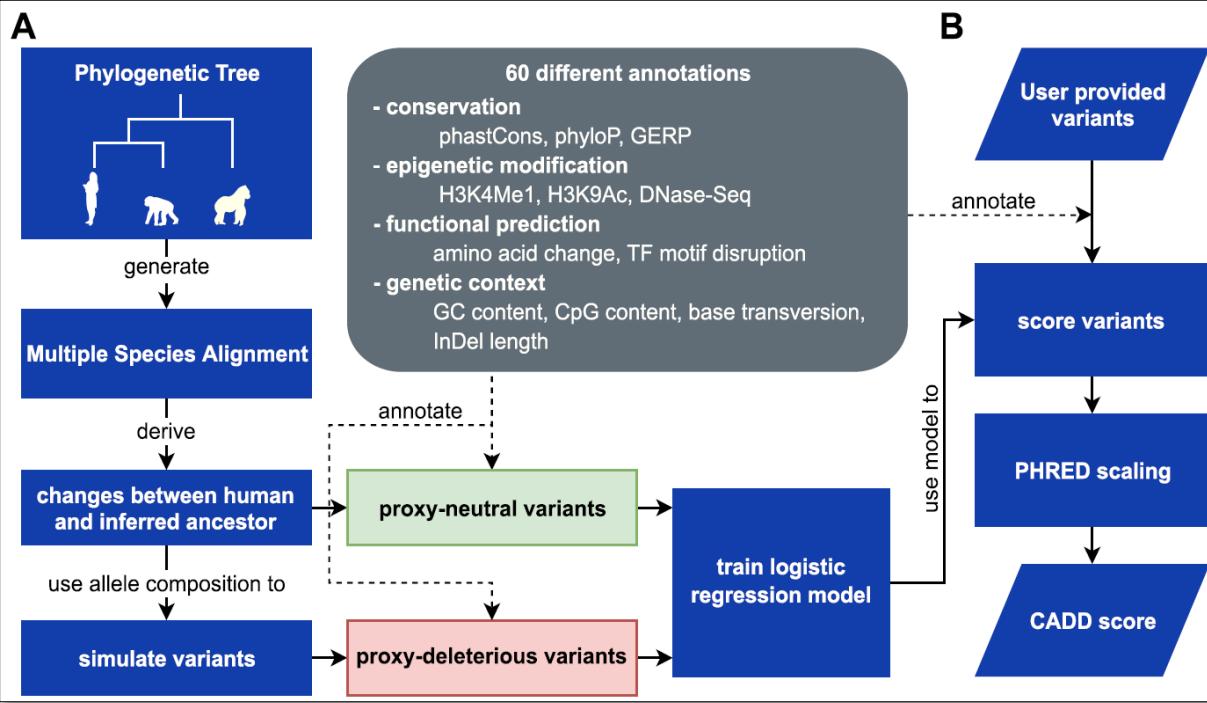
# Today: Deep Learning for Human Genetics and Disease

1. Review: GWAS, fine-mapping, Bayesian variant prioritization
2. Deep Learning for GWAS: calling SNPs, prioritize function
3. eQTLs/Mediation: intermediate molecular phenotypes
4. Linear Mixed Models (LMMs) for GWAS and for eQTL calling
5. Polygenic Risk Scores (PRS): summing over many variants
6. Heritability: definition(s), missing heritability, partitioning
7. LD SCore regression (LDSC) for fast heritability partitioning
8. Polygenic/Omnigenic disease models: core vs. periphery
9. Disease gene networks from GWAS evidence boosting

## **2. Deep Learning methods for GWAS**

### Calling variants, prioritizing functional SNPs

# CADD: combine evidence to predict variant function



Nucleic Acids Research, 2018 1  
doi: 10.1093/nar/gky1016

**CADD: predicting the deleteriousness of variants throughout the human genome**

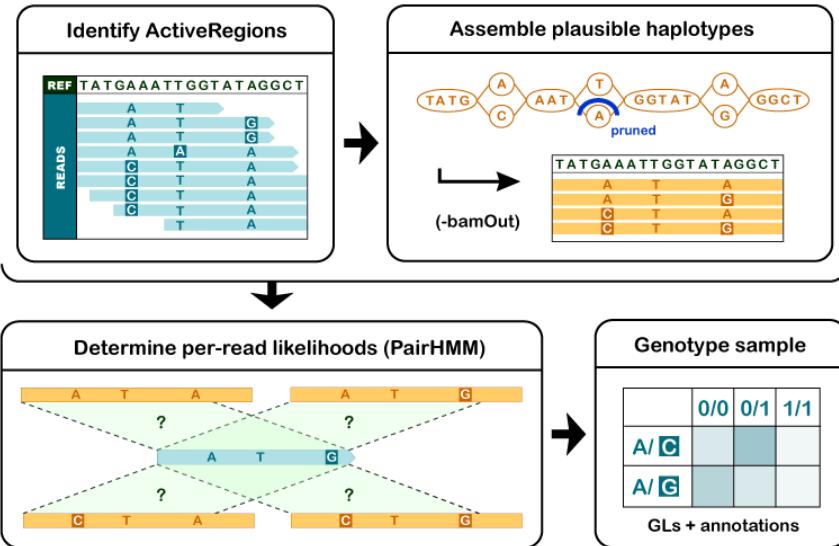
Philipp Rentzsch <sup>1,2</sup>, Daniela Witten<sup>3</sup>, Gregory M. Cooper <sup>2,4</sup>, Jay Shendure <sup>2,5,6,\*</sup> and Martin Kircher <sup>1,2,5,\*</sup>

# Large number of methods for variant prioritization

Score	Data sources	Approach	Ref.
Eigen	<ul style="list-style-type: none"> <li>Uses data from the ENCODE and Roadmap Epigenomics projects</li> </ul>	<ul style="list-style-type: none"> <li>Weighted linear combination of individual annotations</li> <li>Unsupervised learning method</li> <li>Weighted scoring system</li> </ul>	(14)
FunSeq2	<ul style="list-style-type: none"> <li>Inter- and Intra-species conservation</li> <li>Loss- and gain-of-function events for transcription factor binding</li> <li>Enhancer-gene linkage</li> </ul>		(15)
LINSIGHT	<ul style="list-style-type: none"> <li>Conservation scores (phastCons, phyloP), predicted binding sites (TFBS, RNA), regional annotations (ChIP-seq, RNA-seq)</li> </ul>	<ul style="list-style-type: none"> <li>Graphical model</li> <li>Selection parameter fitting using generalized linear model based on 48 genomic features</li> </ul>	(16)
CADD	<ul style="list-style-type: none"> <li>Ensembl variant effect predictor</li> <li>Protein-level scores: Grantham, SIFT, PolyPhen</li> <li>DNase hypersensitivity, TFBS, transcript information</li> <li>GC content, CpG content, histone methylation</li> <li>46-way sequence conservation</li> <li>ChIP-seq, TFBS, DNase-seq</li> <li>FAIRE, footprints, GC content</li> </ul>	<ul style="list-style-type: none"> <li>Support vector machine</li> </ul>	(11)
FATHMM		<ul style="list-style-type: none"> <li>Hidden Markov models</li> </ul>	(17)
ReMM	<ul style="list-style-type: none"> <li>Predict potential of non-coding variant to cause a Mendelian disease if mutated</li> <li>26 features: PhastCons, PhyloP, CpG, GC, regulation annotations</li> </ul>	<ul style="list-style-type: none"> <li>Random forest classifier</li> </ul>	(18)
Orion	<ul style="list-style-type: none"> <li>Predict potential of non-coding variant to cause a Mendelian disease if mutated</li> </ul>	<ul style="list-style-type: none"> <li>Expected and observed site-frequency spectrum of a given stretch of sequence</li> </ul>	(19)
CDTS	<ul style="list-style-type: none"> <li>Independent from annotation and features</li> <li>Identify constrained non-coding regions in the human genome and deleteriousness of variants</li> <li>Independent from annotation and features. Uses k-mers</li> </ul>	<ul style="list-style-type: none"> <li>Expected and observed site-frequency spectrum of a given heptamer</li> </ul>	(8)

# Whole genome variant calling: GATK HaplotypeCaller

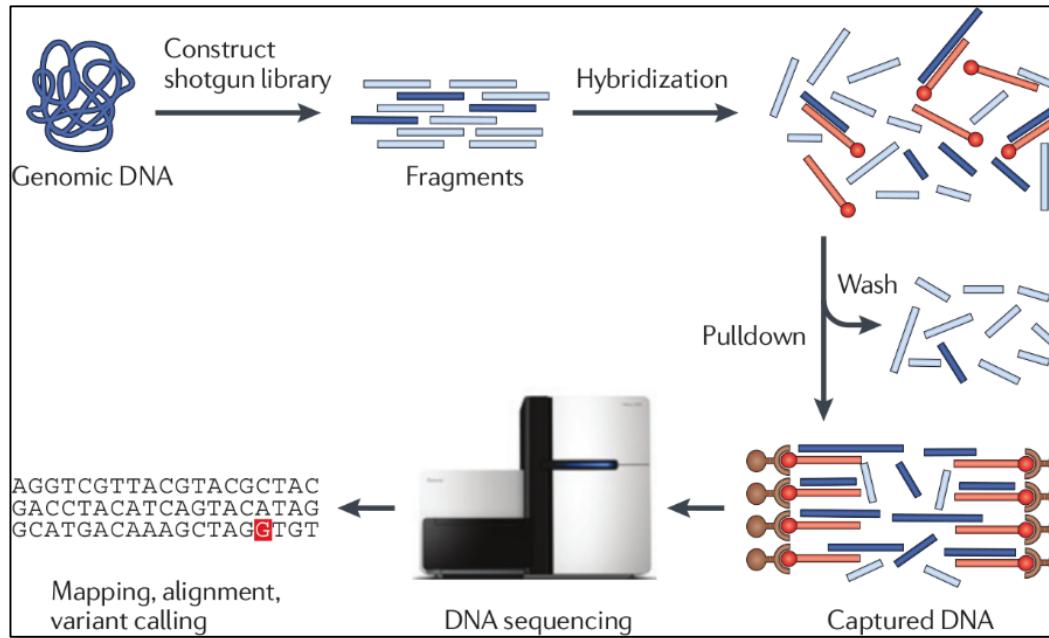
1. Use heuristic to find mismatches not explained by noise
2. Use assembly graph to identify possible haplotypes
3. For each haplotype, estimate:  
**P(read | haplotype)**  
using *probabilistic sequence alignment*
  - Hidden Markov Model
  - States: insertion, deletion, substitution
  - Emissions: pairs of aligned nucleotides/gaps
  - Transitions: equivalent to insertion/deletion/gap penalties from Smith-Waterman algorithm (DP alignment)
  - Get **P(read | haplotype)** using forward-backward algorithm
4. Use Bayes rule to get **P(haplotype | read)**
5. Assign genotypes to each sample based on the max a posteriori haplotypes



Tour de Force, combining many methods:

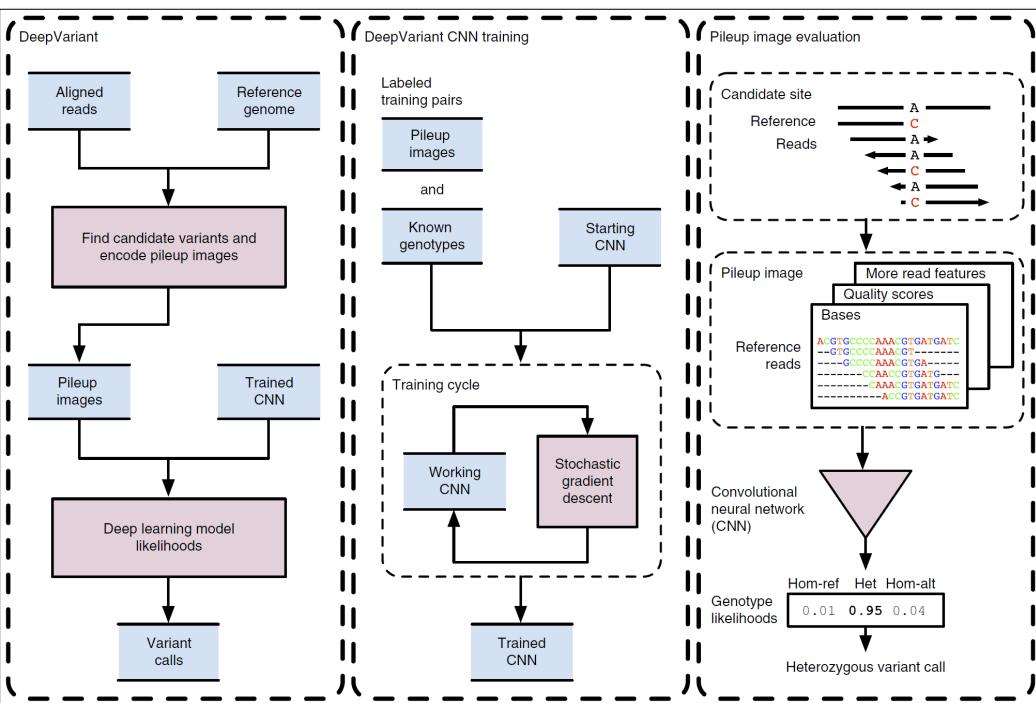
- **Logistic regression** to model base errors
- **Hidden Markov models** to compute read likelihoods
- **Naive Bayes** classification to identify variants
- **Gaussian mixture model** with hand-crafted features to filter likely false positive variants, capturing common error modes

# Exome variant calling: atlas2



- Motivation: the exome has different sequence properties than the rest of the genome (e.g., substitution rates, GC content).
- Train **logistic regression classifier** to predict which mismatches are errors and which are variants
  - Training data: 1KG Exome project sequencing reads where >2 reads align with a mismatch
  - True positives: Reads where mismatch is also discovered in 1KG Exon pilot project
  - True negatives: Remaining reads
  - Features: mismatch quality score, flanking quality score, whether neighboring nucleotides were swapped, normalized distance to 3' end of the read
- Much faster than full Bayesian model (e.g. HaplotypeCaller), lower false positive rate in validation data

# DeepVariant: Combine evidence to call variants



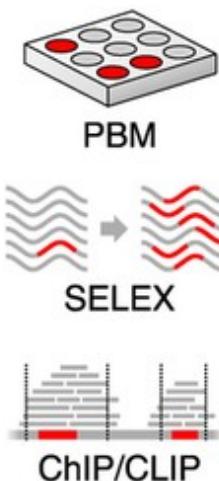
A universal SNP and small-indel variant caller using deep neural networks

Ryan Poplin<sup>1,2</sup>, Pi-Chuan Chang<sup>2</sup>, David Alexander<sup>2</sup>, Scott Schwartz<sup>2</sup>, Thomas Colthurst<sup>2</sup>, Alexander Ku<sup>2</sup>, Dan Newburger<sup>1</sup>, Jojo Dijamco<sup>1</sup>, Nam Nguyen<sup>1</sup>, Pegah T Afshar<sup>1</sup>, Sam S Gross<sup>1</sup>, Lizzie Dorfman<sup>1,2</sup>, Cory Y McLean<sup>1,2</sup> & Mark A DePristo<sup>1,2</sup>

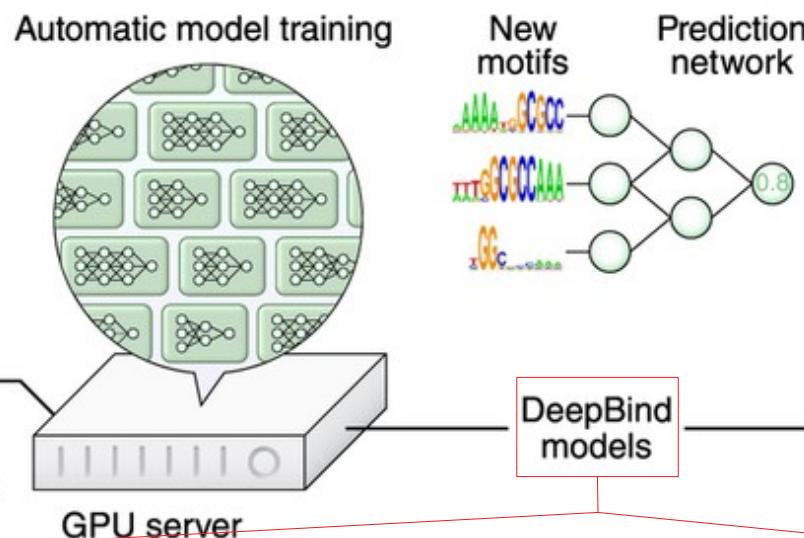
Method	Type	F1	Recall	Precision	TP	FN	FP	FP.gt	FP.al	Version
DeepVariant (live GitHub)	Indel	0.99507	0.99347	0.99666	357,641	2350	1,198	217	840	Latest GitHub v0.4.1-b4e8d37d
GATK (raw)	Indel	0.99366	0.99219	0.99512	357,181	2810	1,752	377	995	3.8-0-ge9d806836
Strelka	Indel	0.99227	0.98829	0.99628	355,777	4214	1,329	221	855	2.8.4-3-gbe58942
DeepVariant (pFDA)	Indel	0.99112	0.98776	0.99450	355,586	4405	1,968	846	1,027	pFDA submission May 2016
GATK (VQSR)	Indel	0.99010	0.98454	0.99573	354,425	5566	1,522	343	909	3.8-0-ge9d806836
GATK (flt)	Indel	0.98229	0.96881	0.99615	348,764	11227	1,349	370	916	3.8-0-ge9d806836
FreeBayes	Indel	0.94091	0.91917	0.96372	330,891	29,100	12,569	9,149	3,347	v1.1.0-54-g49413aa
16GT	Indel	0.92732	0.91102	0.94422	327,960	32,031	19,364	10,700	7,745	v1.0-34e8f934
SAMtools	Indel	0.87951	0.83369	0.93066	300,120	59,871	22,682	2,302	20,282	1.6
DeepVariant (live GitHub)	SNP	0.99982	0.99975	0.99989	3,054,552	754	350	157	38	Latest GitHub v0.4.1-b4e8d37d
DeepVariant (pFDA)	SNP	0.99958	0.99944	0.99973	3,053,579	1,727	837	409	78	pFDA submission May 2016
Strelka	SNP	0.99935	0.99893	0.99976	3,052,050	3,256	732	87	136	2.8.4-3-gbe58942
GATK (raw)	SNP	0.99914	0.99973	0.99854	3,054,494	812	4,469	176	257	3.8-0-ge9d806836
16GT	SNP	0.99583	0.99850	0.99318	3,050,725	4,581	20,947	3,476	3,899	v1.0-34e8f934
GATK (VQSR)	SNP	0.99436	0.98940	0.99937	3,022,917	32,389	1,920	80	170	3.8-0-ge9d806836
FreeBayes	SNP	0.99124	0.98342	0.99919	3,004,641	50,665	2,434	351	1,232	v1.1.0-54-g49413aa
SAMtools	SNP	0.99021	0.98114	0.99945	2,997,677	57,629	1,651	1,040	200	1.6
GATK (flt)	SNP	0.98958	0.97953	0.99983	2,992,764	62,542	509	168	26	3.8-0-ge9d806836

# DeepBind

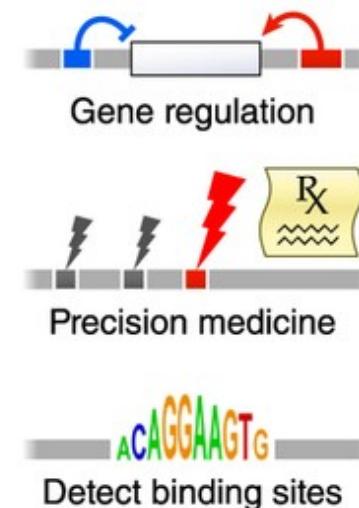
## 1. High-throughput experiments



## 2. Massively parallel deep learning

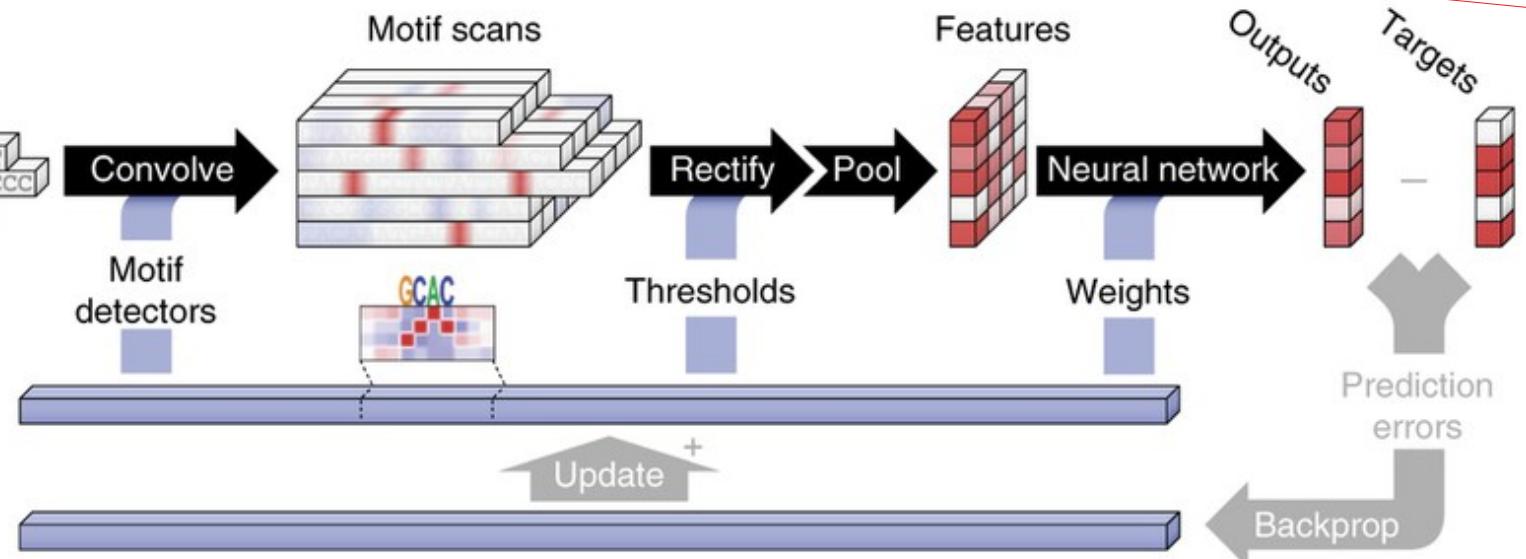


## 3. Community needs

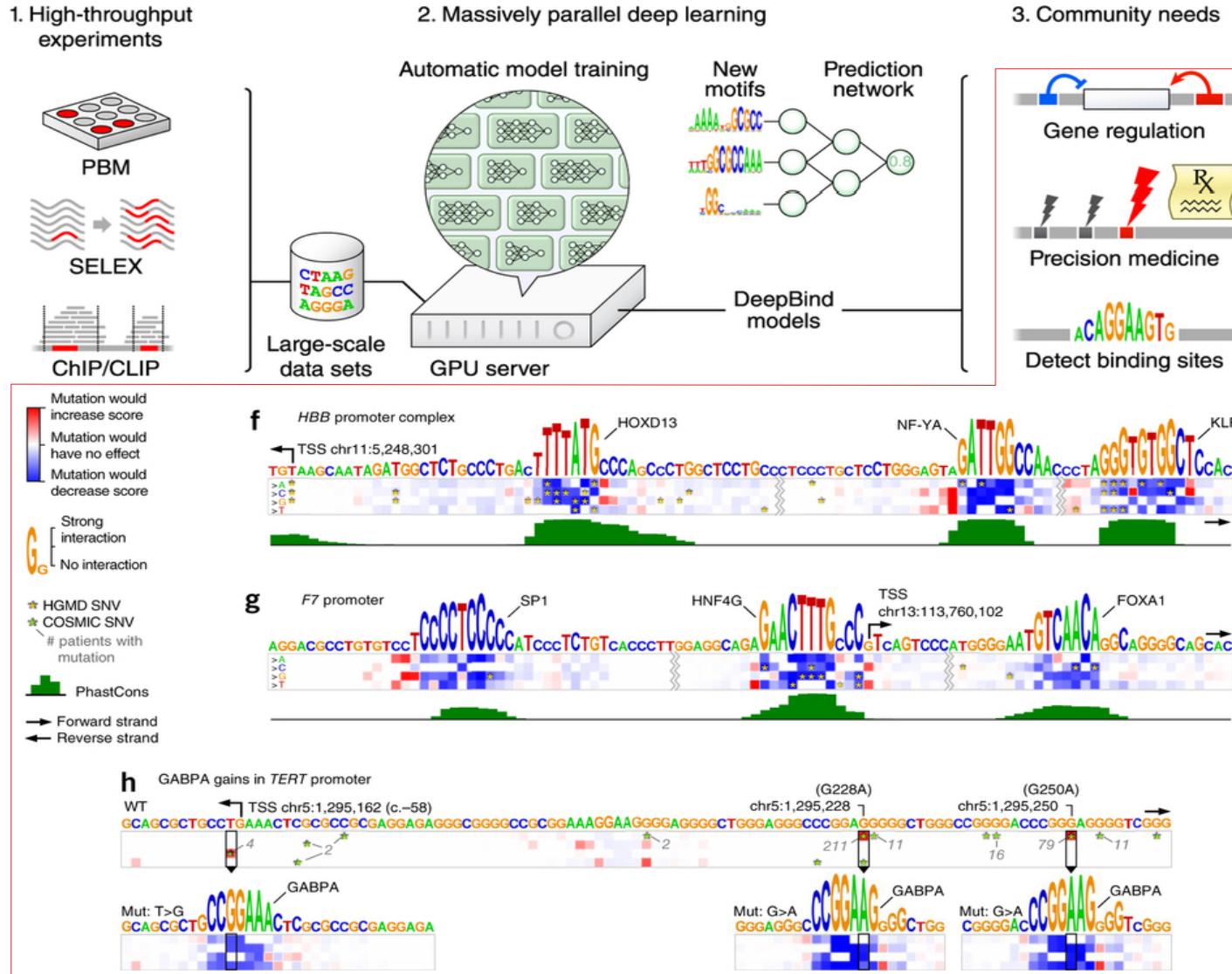


### Current batch of inputs

CTAACGCACCGTCT  
TTAGGGGCACCACTACT  
TAGCACCTCTATTGACACC  
CTCGGGGCCCTGCAAT  
TACAAATGAGCACAA



# Predicting disease mutations



[Alipanahi et al., 2015]

# DeepBind summary

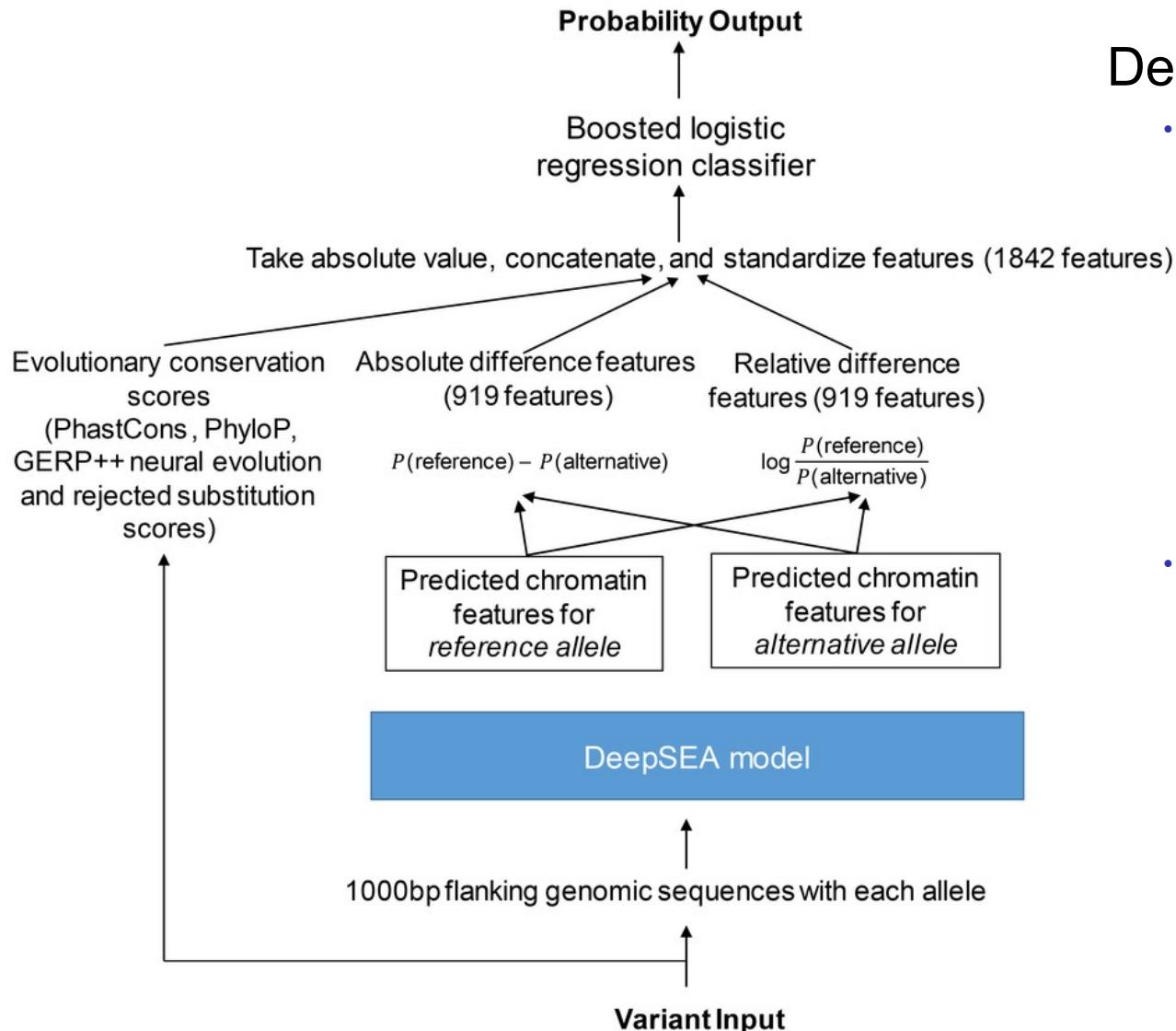
The key deep learning techniques:

- Convolutional learning
- Representational learning
- Back-propagation and stochastic gradient
- Regularization and dropout
- Parallel GPU computing especially useful for hyperparameter search

Limitations in DeepBind:

- Require defining negative training examples, which is often arbitrary
- Using observed mutation data only as post-hoc evaluation
- Modeling each regulatory dataset separately

# DeepSea



## DeepSea:

- Similar as DeepBind but trained a separate CNN on each of the ENCODE/Roadmap Epigenomic chromatin profiles 919 chromatin features (125 DNase features, 690 TF features, 104 histone features).
- It uses the  $\Delta s$  mutation score as input to train a linear logistic regression to predict GWAS and eQTL SNPs defined from the GRASP database with a P-value cutoff of 1E-10 and GWAS SNPs from the NHGRI GWAS Catalog

## CNNs for DNA-binding prediction from sequence

DanQ: a hybrid convolutional and recurrent deep neural network for quantifying the function of DNA sequences. Uses convolution layers to capture regulatory motifs, and a recurrent layer to discover a 'grammar' for how these single motifs work together. Based on Keras/Theano.

Basset—learning the regulatory code of the accessible genome with deep convolutional neural networks. CNN to discover regulatory sequence motifs to predict the accessibility of chromatin. Accounts for cell-type specificity using multi-task learning.

DeepBind and DeeperBind—predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning. Based on ChIP-seq, ChIP-chip, RIP-seq, protein-binding microarrays and others. Deeperbind adds a recurrent sequence learning module (LSTM) after the convolutional layer(s).

DeepMotif—visualizing genomic sequence classifications. Predicting binding specificities of proteins to DNA motifs. Makes use of a convolutional layers with more layers than the DeepBind network.

Convolutional neural network architectures for predicting DNA–protein binding. Systematic exploration of CNN architectures for predicting DNA sequence binding using a large compendium of transcription factor data sets.

## Predicting enhancers, 3d interactions and cis-regulatory regions

PEDLA: predicting enhancers with a deep-learning-based algorithmic framework. Predicting enhancers based on heterogeneous features from (e.g.) the ENCODE project using a deep learning, HMM hybrid model.

DEEP: a general computational framework for predicting enhancers. Predicting enhancers based on data from the ENCODE project.

Genome-wide prediction of cis-regulatory regions using supervised deep-learning methods. toolkit based on the Theano) for applying different deep-learning architectures to cis-regulatory elements.

FIDDLE: an integrative deep-learning framework for functional genomic data inference. Prediction of transcription start site and regulatory regions. FIDDLE stands for Flexible Integration of Data with Deep Learning that models several genomic signals using convolutional networks (DNase-seq, ATAC-seq, ChIP-seq, TSS-seq, RNA-seq signals).

## DNA methylation

DeepCpG—predicting DNA methylation in single cells. Neural network for predicting DNA methylation in multiple cells.

Predicting DNA methylation state of CpG dinucleotide using genome topological features and deep networks. Uses a stacked autoencoder with a supervised layer on top of it to predict whether CpG islands are methylated.

## Variant callers, pathogenicity scores and identification of genomic elements

DeepVariant—a variant caller in germline genomes. Uses a deep neural network architecture (Inception-v3) to identify SNP and small indel variants from next-generation DNA sequencing data.

DeepLNC, a long non-coding RNA prediction tool using deep neural network. Identification of lncRNA-based on k-mer profiles.

evoNet—deep learning for population genetic inference [code][paper]. Jointly inferring natural selection and demographic history

DANN. Uses the same feature set and training data as CADD to train a deep neural network

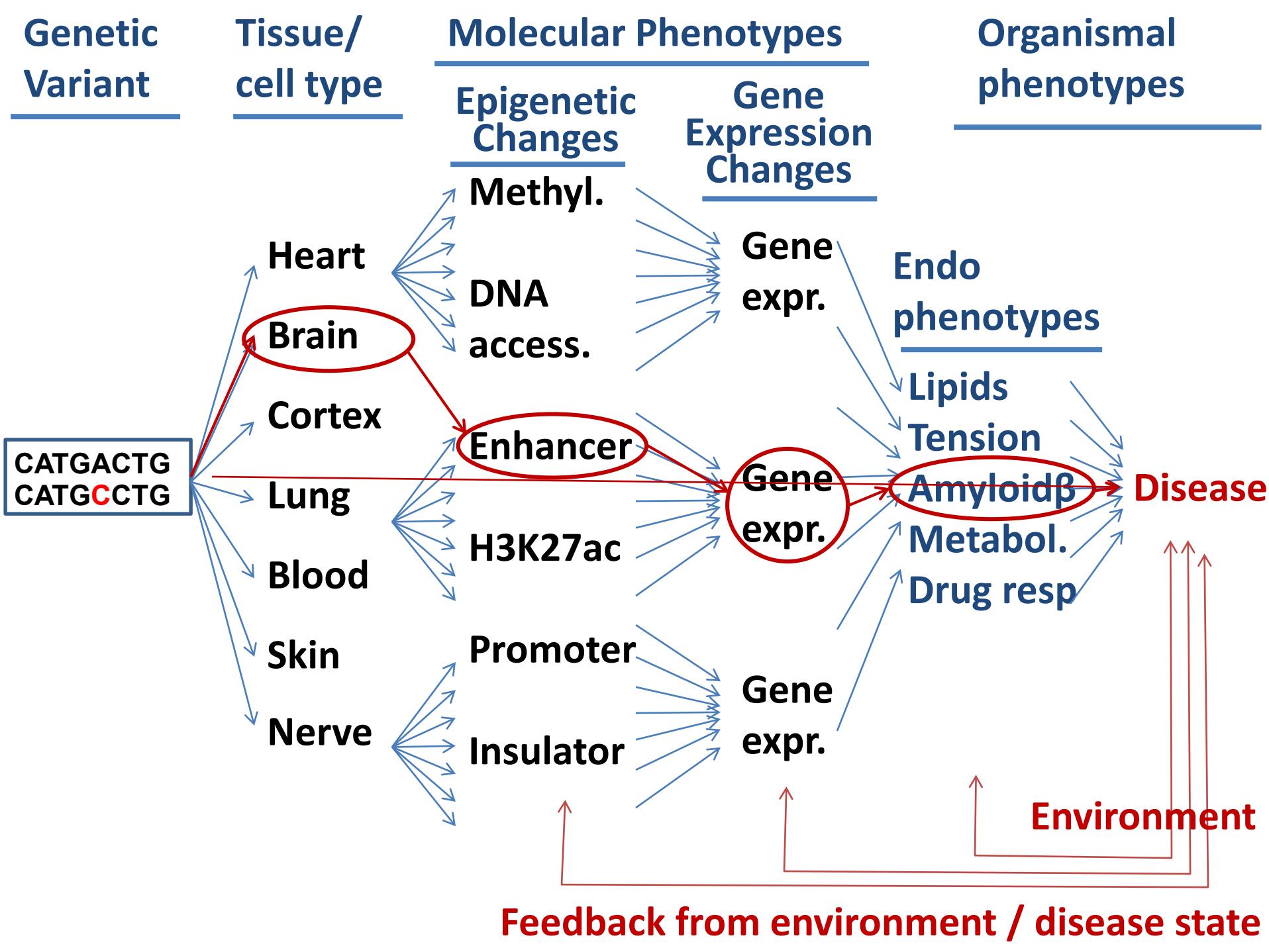
DeepSEA—predicting effects of non-coding variants with deep-learning-based sequence model. Models chromatin accessibility as well as the binding of transcription factors, and histone marks associated with changes in accessibility.

# Today: Deep Learning for Human Genetics and Disease

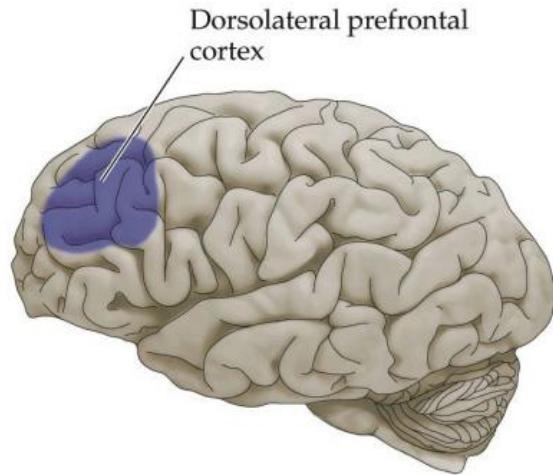
1. Review: GWAS, fine-mapping, Bayesian variant prioritization
2. Deep Learning for GWAS: calling SNPs, prioritize function
3. eQTLs/Mediation: intermediate molecular phenotypes
4. Linear Mixed Models (LMMs) for GWAS and for eQTL calling
5. Polygenic Risk Scores (PRS): summing over many variants
6. Heritability: definition(s), missing heritability, partitioning
7. LD SCore regression (LDSC) for fast heritability partitioning
8. Polygenic/Omnigenic disease models: core vs. periphery
9. Disease gene networks from GWAS evidence boosting

### **3. eQTLs and Mediation Analysis**

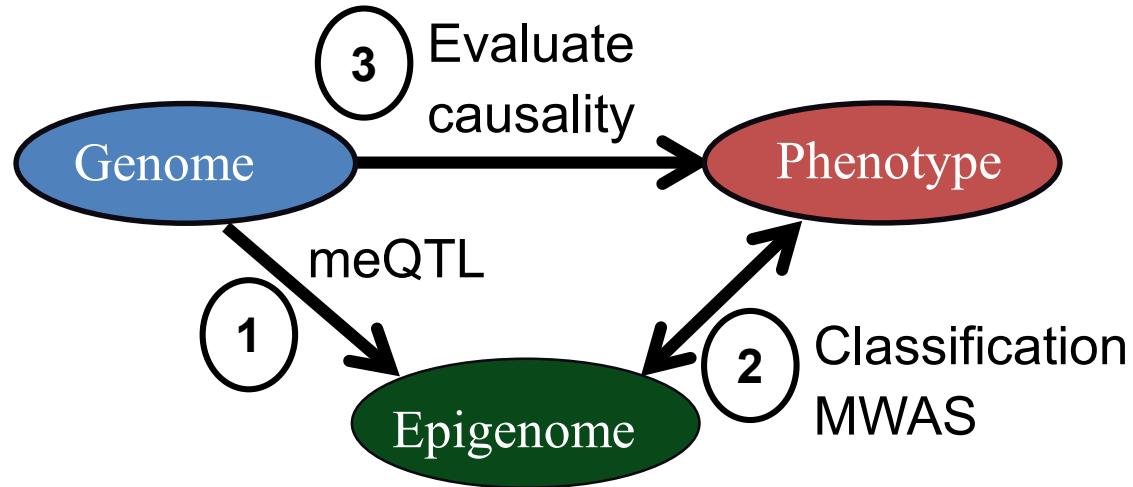
Intermediate molecular phenotypes to disease



# Methylation in 750 Alzheimer patients/controls



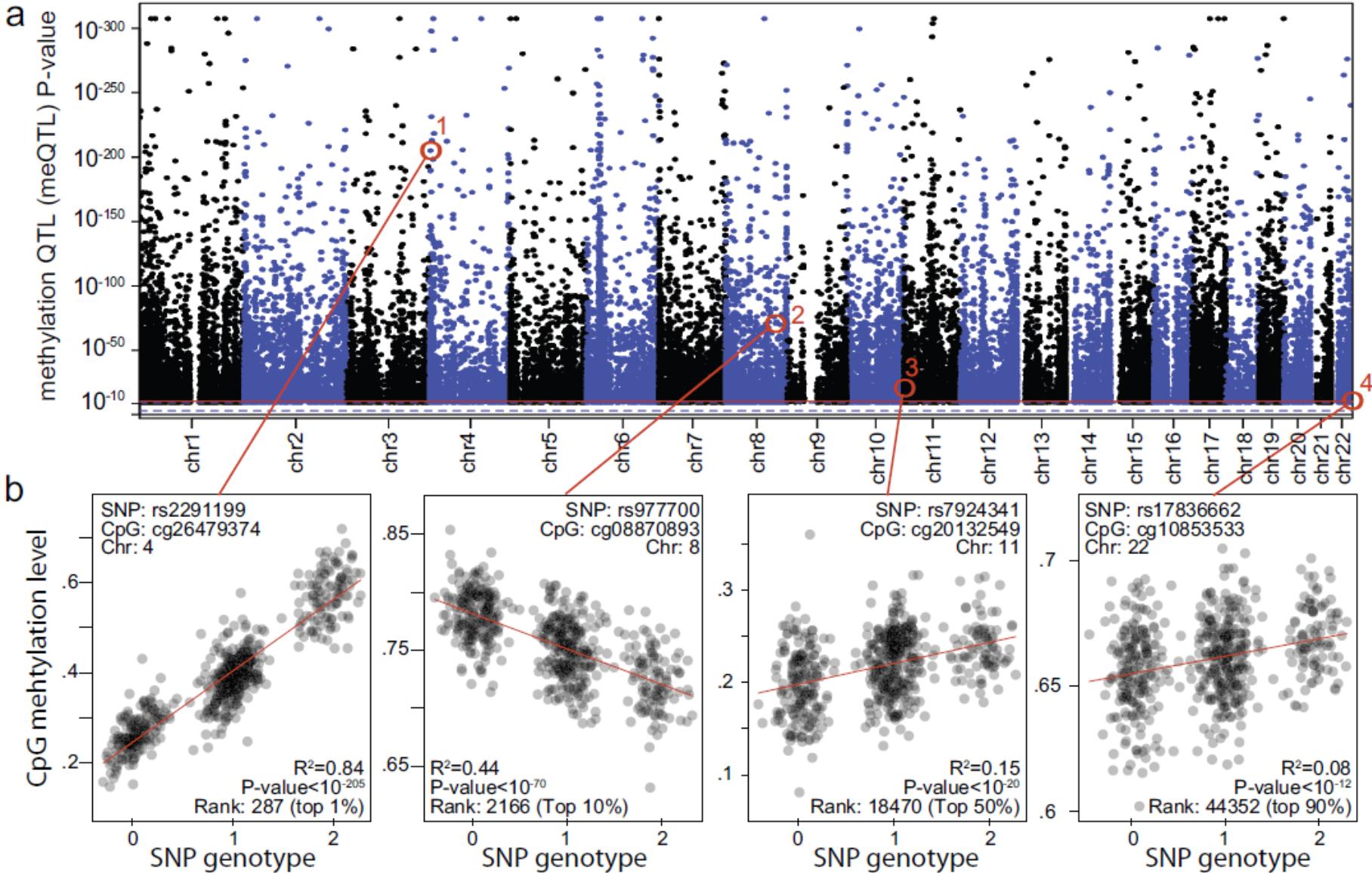
Methylation variation  
in 723 individuals



Relate to genotype and AD variation

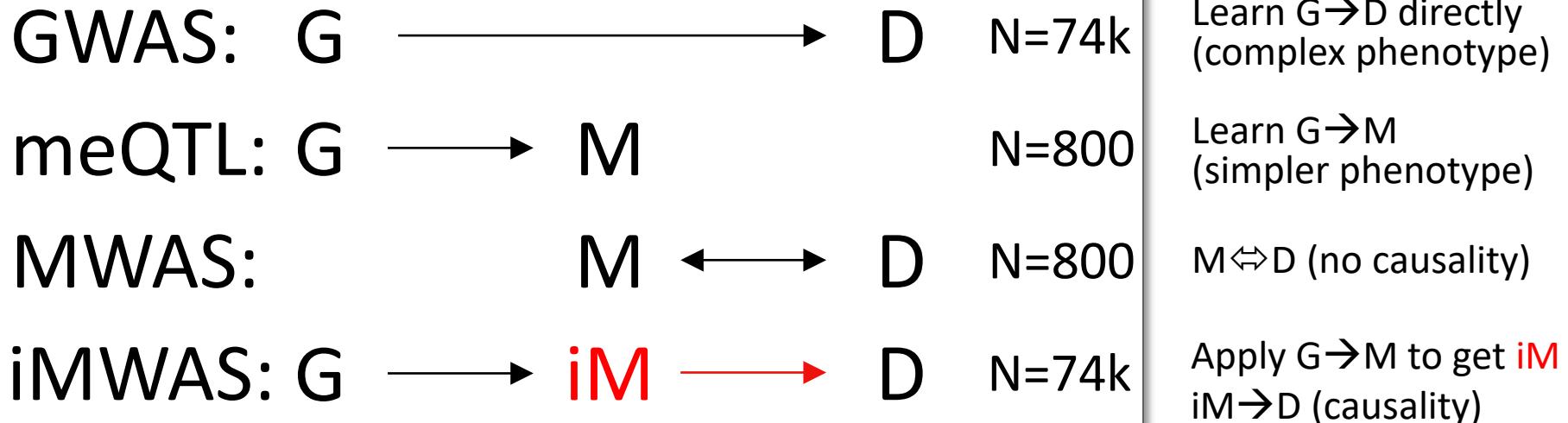
- *ROS-MAP cohort (RUSH: David Bennett, HMS: Phil De Jager)*
  - *Patients followed for 10+ years with cognitive evaluations*
  - *Brain samples donated post-mortem methylation/genotype*
- *Seek predictive features: SNPs, QTLs, mQTLs, regulation*

# 50,000 significant meQTLs after Bonferroni



- Strong effects across entire range of discovery values

# Imputed MWAS: increased power, genetic component



## Key Idea:

- Learn G→M model (ROSMAP n=800) Fewer indiv. Simpler phenotype
- Impute methylation iM for GWAS cohort (n=74k)
- iMWAS between genotype-driven M and AD phenotype (n=47k)

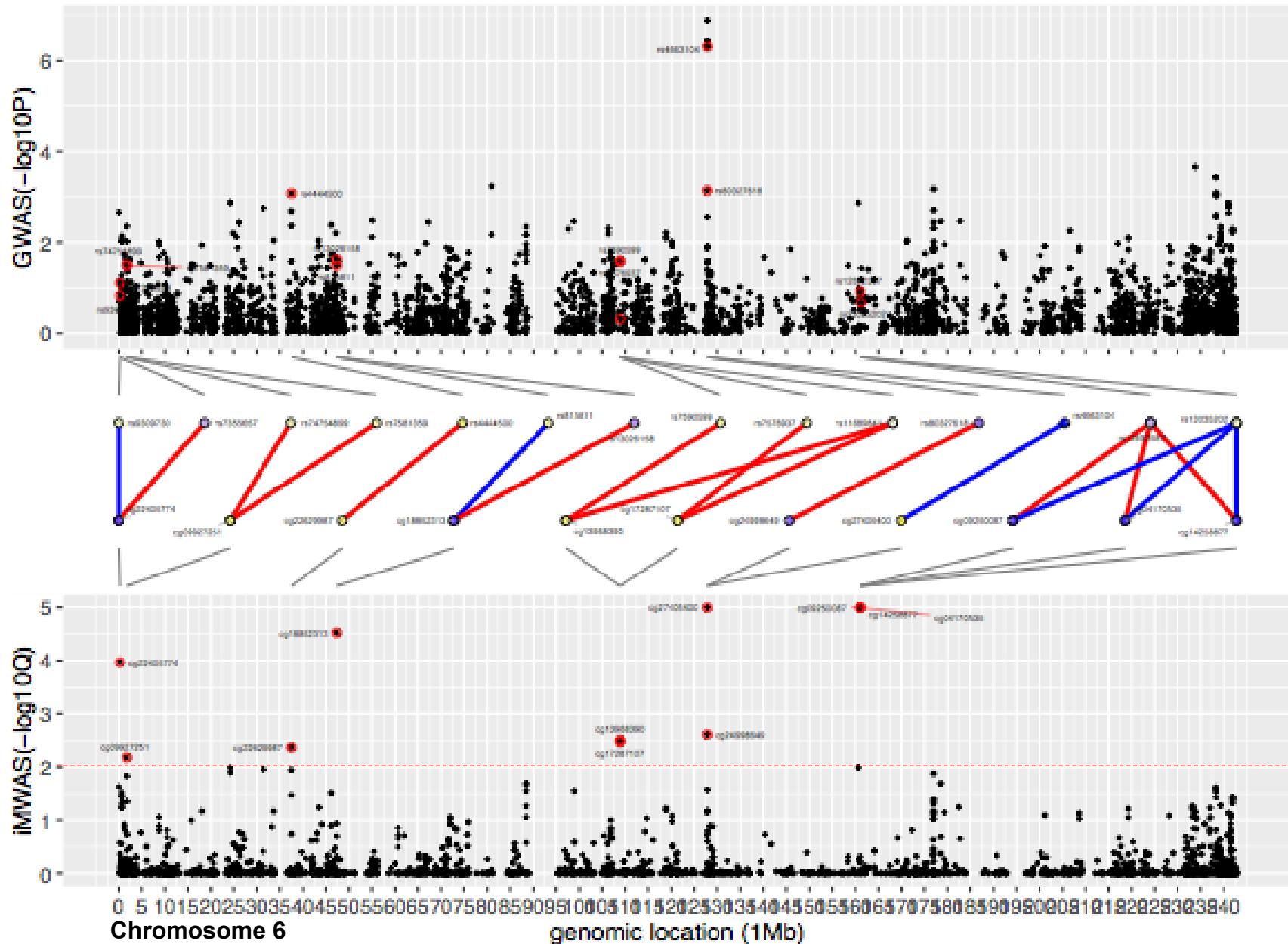
## Advantage:

- Much larger GWAS cohorts (>>MWAS): increased power
- Genetic component of methyl. variation

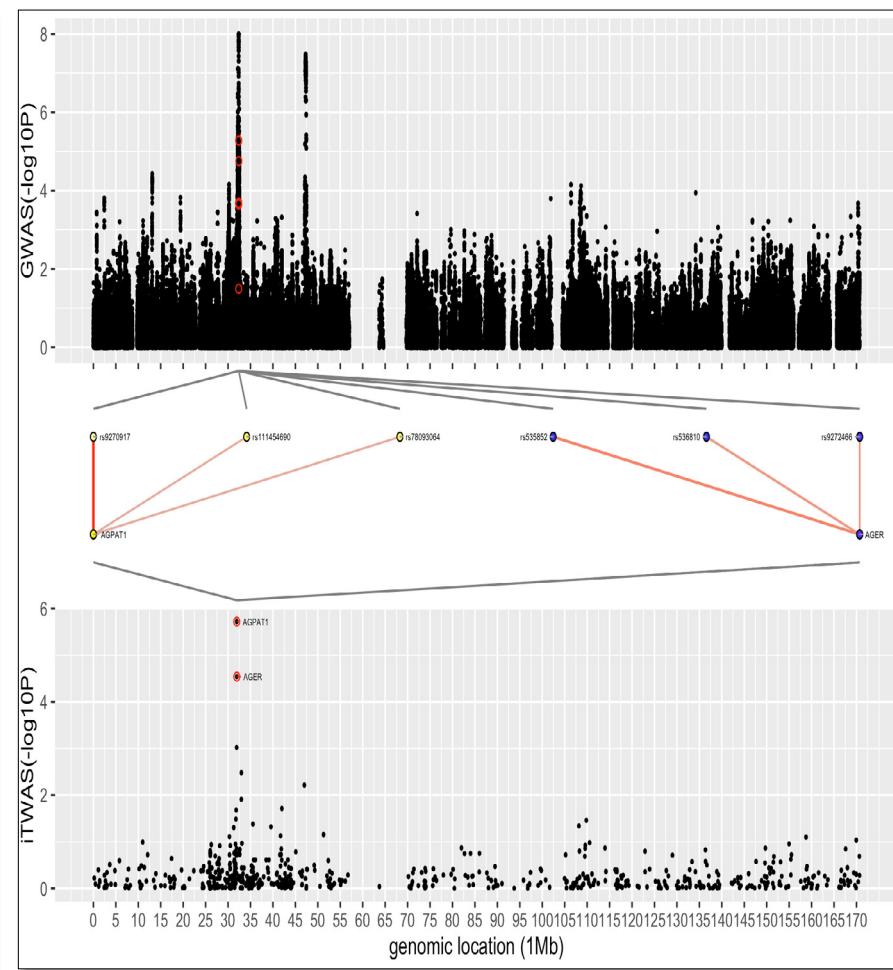
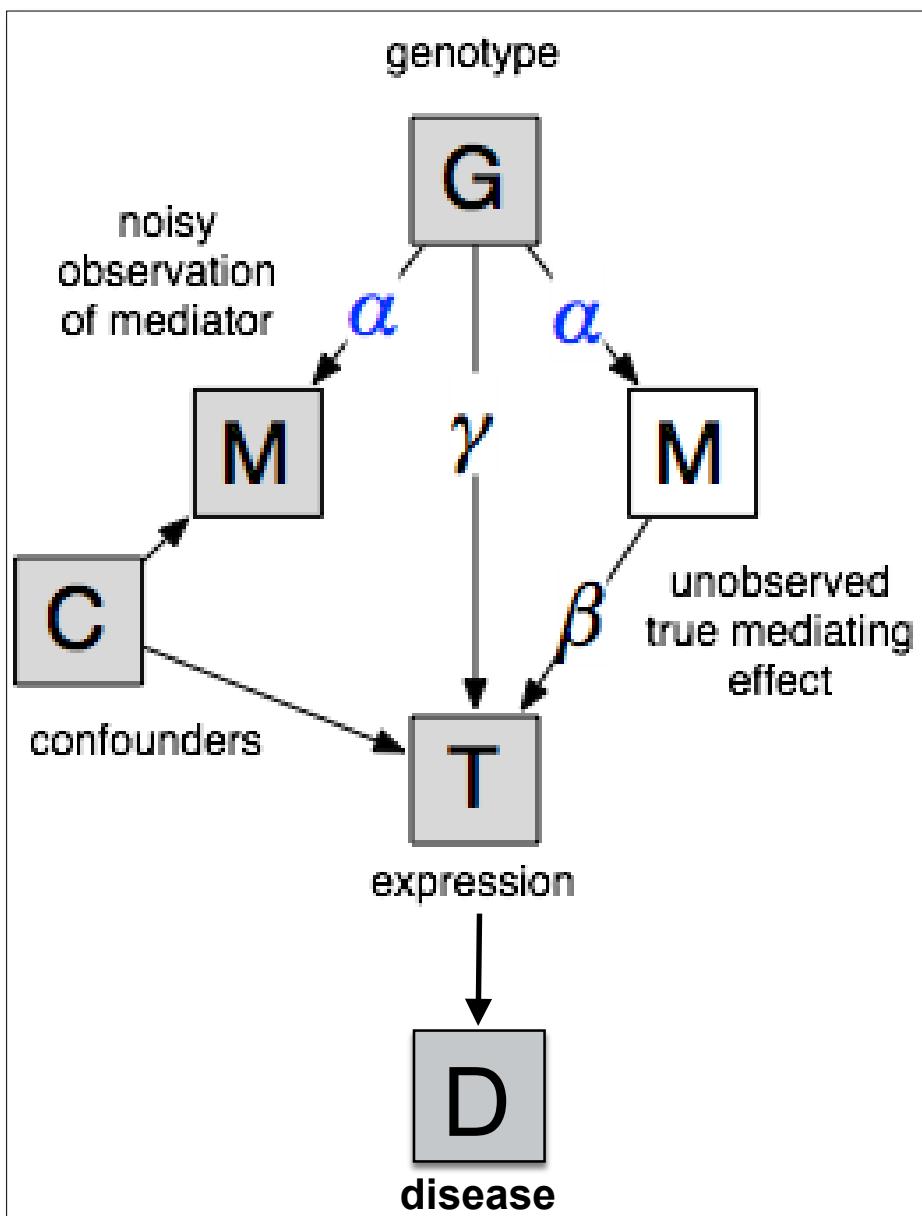
## Logistical challenge:

- Summary stats, not full genotypes → Linear model, impute stats direct

# iMWAS results: new loci, multiple contributing SNPs



# iMTWAS: Imputation across multiple intermediate variables



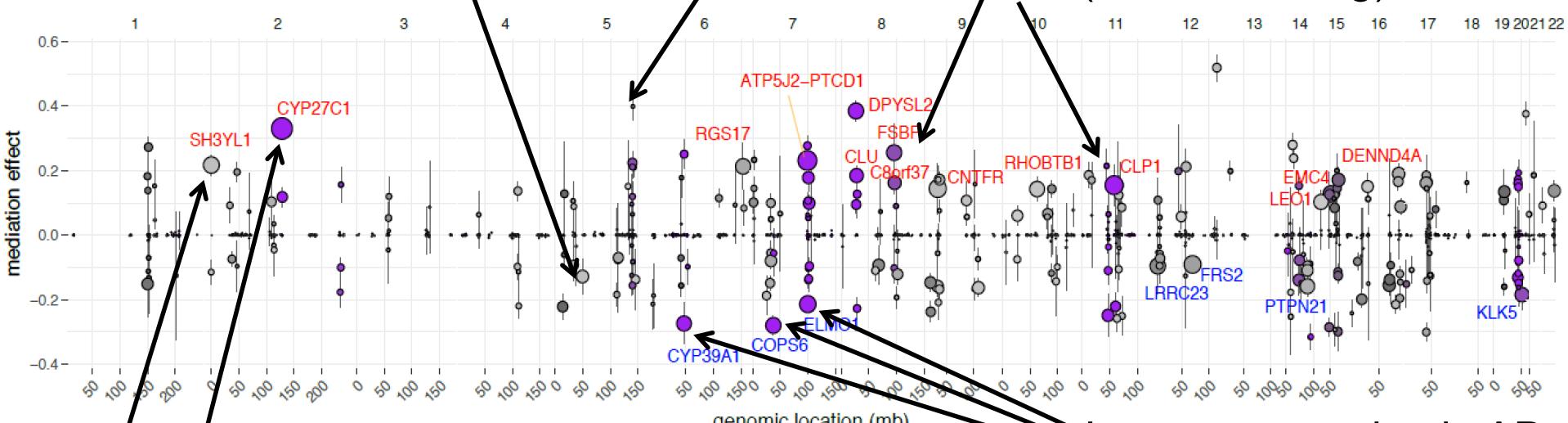
Model multiple mediator variables  
SNP → Methylation → Expression → Disease  
Predict new loci, increased power  
Predict regulatory regions & target genes

# CaMMEL: 206 significant mediating genes in AD

Small expression change (short),  
large variance explained (big circle)

Large expression change (tall),  
little variance explained (small circle)

Higher-expression in AD  
(risk increasing)

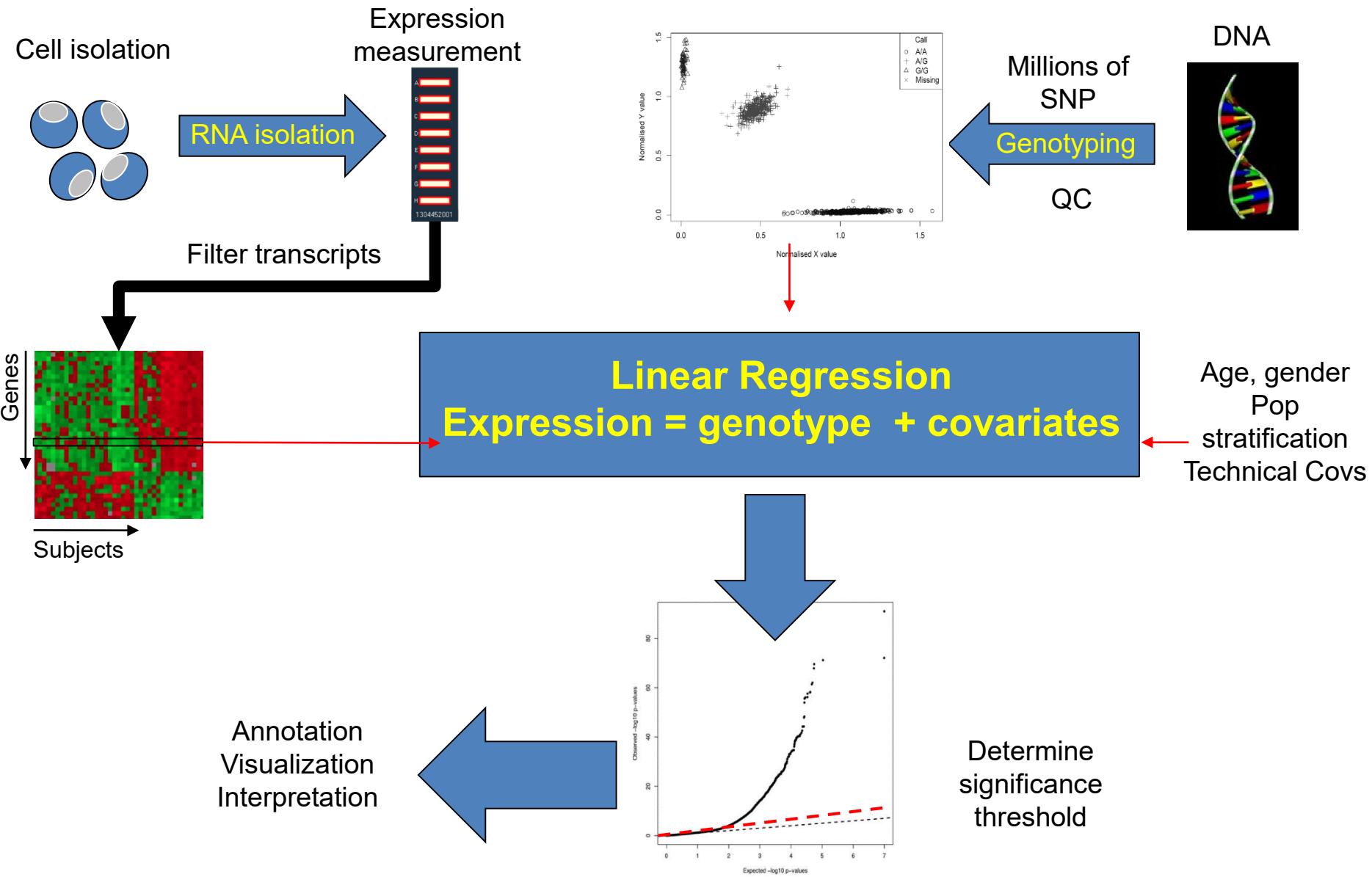


Genome-wide significant locus (purple)

Lower-expression in AD  
(protective)

Sub-threshold locus (grey)

# The nuts and bolts of an eQTL study



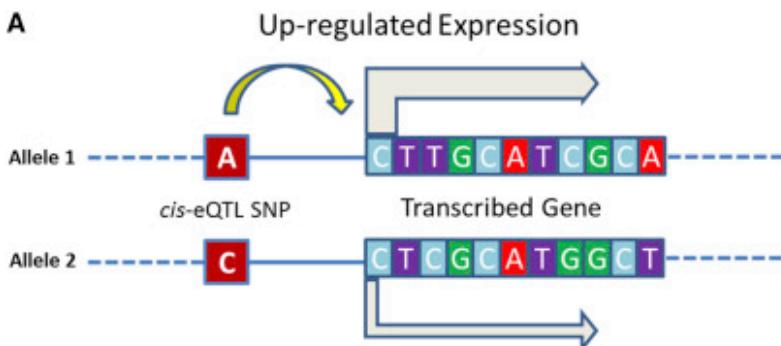
# Expanded eQTL models

$$Y_{ij} = \alpha + \beta_{ijs} \text{genotype} + \varepsilon$$

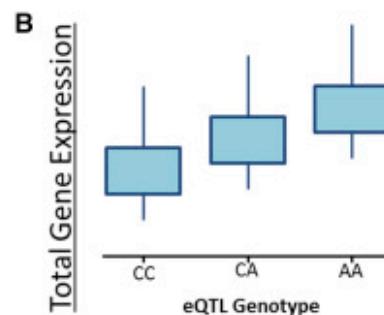
$$Y_{ij} = \alpha + \beta_{1ijs} \text{genotype} + \beta_{2i} \text{gender} + \beta_{3i} \text{age} +$$
$$\beta_{4i} g\text{PC1} + \beta_{5i} g\text{PC2} + \beta_{6i} g\text{PC3} + \beta_{7i} g\text{PC4} + \left. \right] \text{Genotype PCs}$$
$$\beta_{8i} e\text{PC1} + \beta_{9i} e\text{PC2} + \beta_{10i} e\text{PC3} + \beta_{11i} e\text{PC4} + \left. \right] \text{Expression PCs}$$
$$\beta_{12i} e\text{PC5} + \beta_{13i} e\text{PC6} + \beta_{14i} e\text{PC7} \left. \right]$$
$$+ \varepsilon$$

# Allelic analysis complements eQTLs

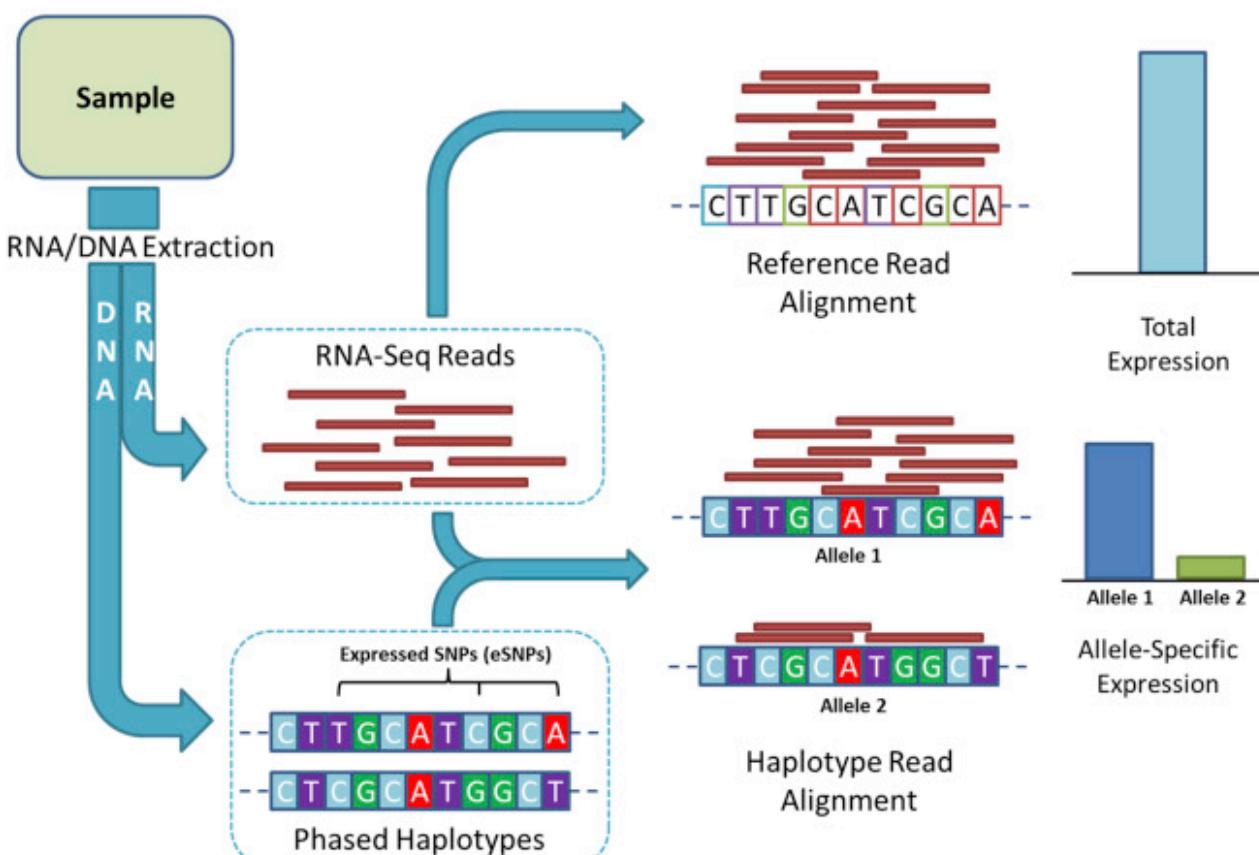
A



B

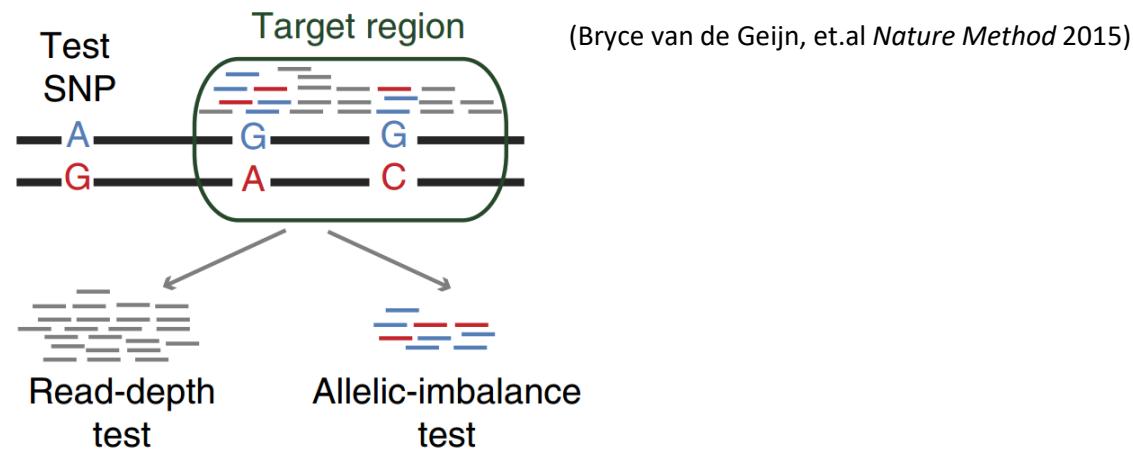


C



Distinguish reads  
within the same  
heterozygous individual

# Combined Haplotype Test

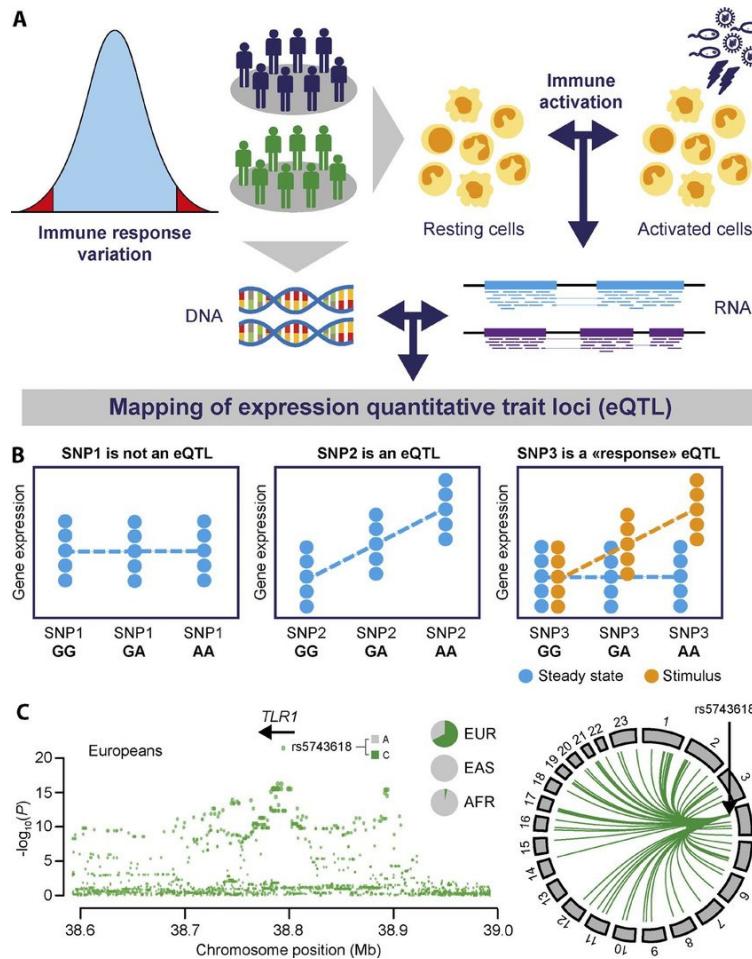


Maximize likelihood of two observed components:

$$\mathbf{L}(\alpha_h, \beta_h, \phi_j | D) = \prod_i \left[ \Pr_{\text{BNB}}(X = x_{ij} | \lambda_{hi}, \Omega_i, \phi_j) \prod_k \Pr_{\text{BB-mix}}(Y = y_{ik} | p_h, n_{ik}, \Upsilon_i) \right]$$

Beta-Negative-Binomial                            Beta-Binomial

# “Response eQTLs”: Trait-conditional eQTLs



# Today: Deep Learning for Human Genetics and Disease

1. Review: GWAS, fine-mapping, Bayesian variant prioritization
2. Deep Learning for GWAS: calling SNPs, prioritize function
3. eQTLs/Mediation: intermediate molecular phenotypes
4. Linear Mixed Models (LMMs) for GWAS and for eQTL calling
5. Polygenic Risk Scores (PRS): summing over many variants
6. Heritability: definition(s), missing heritability, partitioning
7. LD SCore regression (LDSC) for fast heritability partitioning
8. Polygenic/Omnigenic disease models: core vs. periphery
9. Disease gene networks from GWAS evidence boosting

## **4. Linear Mixed Models (LMMs)**

### for GWAS and for eQTL calling

# Formal definition of a linear model

n individuals

p SNPs

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}, X = \begin{pmatrix} X_{11} & \cdots & X_{1p} \\ X_{21} & \cdots & X_{2p} \\ \vdots & \cdots & \vdots \\ X_{n1} & \cdots & X_{np} \end{pmatrix}, \boldsymbol{\theta} = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{pmatrix}$$

In matrix notation, phenotype  $y$  as a factor of genetic information  $x$

$$\mathbf{y} = X\boldsymbol{\theta} + \epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 I).$$

$\theta$  = effect size (can be itself sampled from a normal prior)

# What are we missing in the previous multivariate model?

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boldsymbol{\epsilon}, \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}).$$

Assume IID individuals.  
This may not be true.

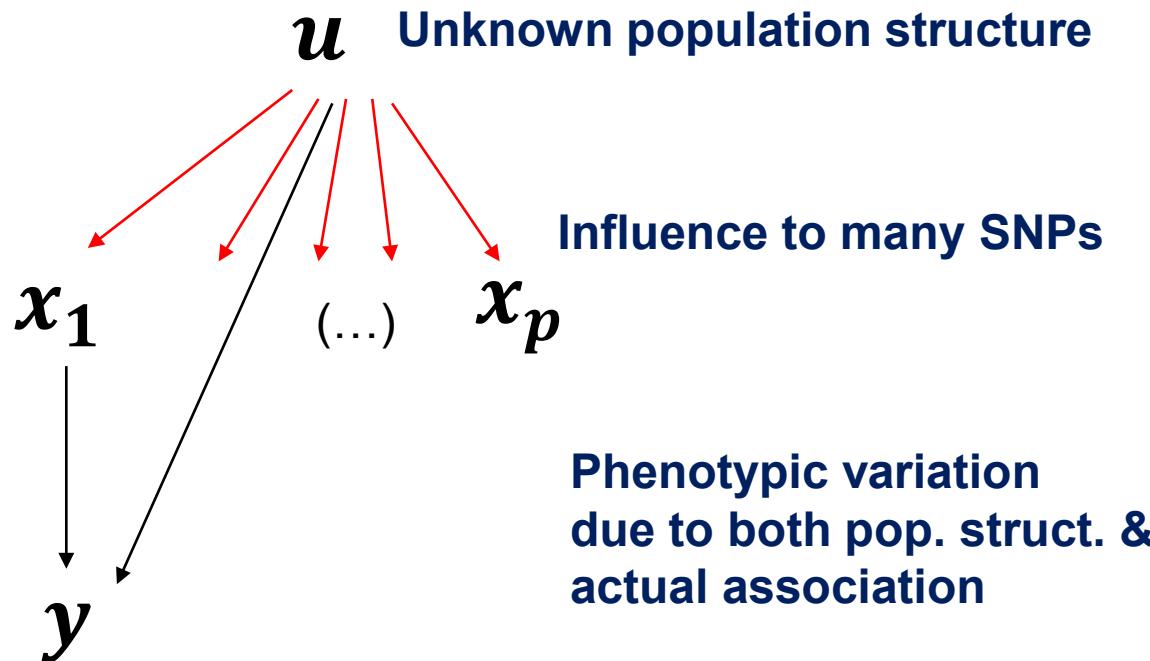
$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boxed{\boldsymbol{u}} + \boldsymbol{\epsilon}. \quad \text{Add random effects to account for the unknown}$$

$$\boxed{\boldsymbol{u}} \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$$

We assume this random effect can be captured by Kinship covariance.

In GWAS problems, the most influential/spurious random effect stems from population structure.

# Why do we need a random effect?



# A Bayesian approach to account for the random effect $\underline{u}$

Likelihood model:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \boxed{\mathbf{u}} + \epsilon.$$

(Empirical) prior knowledge:

$$\boxed{\mathbf{u}} \sim \mathcal{N}(\mathbf{0}, \mathbf{K})$$

A Bayesian method ≈ Address/remove uncertainty by averaging out

$$p(\mathbf{y}|X\boldsymbol{\theta}) = \int p(\mathbf{y}|X\boldsymbol{\theta}, \mathbf{u})p(\mathbf{u})d\mathbf{u}$$

A Linear mixed effect model:

**two components  
in covariance matrix**

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \tilde{\epsilon}$$

with

$$\tilde{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I} + \boxed{\tau^2 \mathbf{K}})$$

IID error      Kinship  
components

# Linear mixed models

$$\begin{aligned} p &\sim N(0, h^2 G + (1 - h^2) I) \\ G &= X X' / p \end{aligned}$$

- Joint model of all SNPs explains more heritability (Yang 2010)
- Idea: under suitable assumptions,  $V[a] = \sum \beta_j^2$
- Under the infinitesimal assumption  $\beta_j \sim N(0, h^2/p)$ , we can estimate  $V[a]$  without estimating individual  $\beta_j$  using residual maximum likelihood (REML)
- REML avoids using ML fit of parameters, instead uses transformed data so that nuisance parameters have no effect.
- In variance components analysis (random effects model), transformation focuses on differences, sum of variances
- **This works despite not knowing the causal variants**
- Example (height): ;  $h^2_{\text{GWAS}} = 0.16$ ,  $h^2 = 0.73$ ,  $h^2_g = 0.5$

# Linear mixed models

$$p \sim N(0, h^2 G - (1 - h^2) I)$$

$$G = XX' / p$$

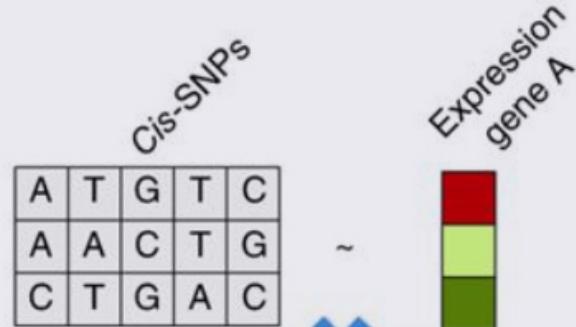
$$E[p_i p_j] = h^2 G_{ij}$$

- We can generalize Haseman-Elston regression to estimate heritability for unrelated individuals using LMM
- Intuition: genetic relationship matrix  $G$  captures identity by state in unrelated individuals
- This is again the probability of sharing the same allele at the causal variants
- This is called **PCGC regression** (Golan 2015)  
(phenotype correlation – genotype correlation regression)

# Imputation-based association

**1 = learn eQTLs in reference panel**

Reference panel



Individual TWAS

Cis-SNPs

A	T	G	T	C
A	A	C	T	G
C	T	G	A	C
C	T	G	A	C
A	A	C	A	C
C	A	G	T	G

Predicted expression  
gene A

Trait

A

SNP-trait  
standardized  
effects

$z_1$	$z_2$	$z_3$
...		

Predicted  
[gene A]-trait  
effect

$$w_1 z_1 + w_2 z_2 + w_3 z_3 + \dots$$

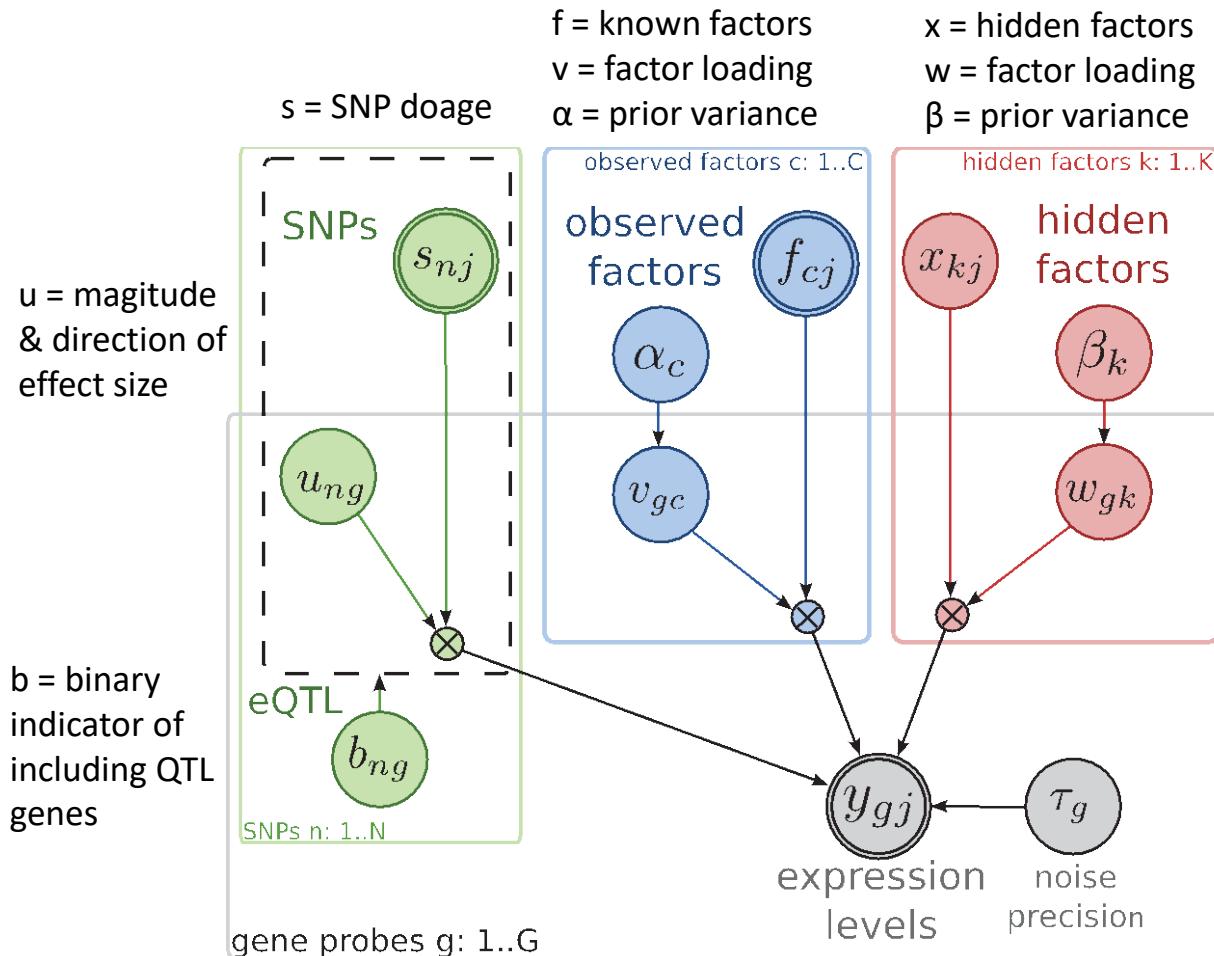


SNP LD  
reference

**2 = impute expression for each person in a genotyped cohort**

**3 = use summary statistics to get to associations directly**

# Bayesian linear regression for eQTL modeling



# Bayesian extension to ordinary regression models

1. Spike-slab prior to select relevant variables
2. Random effect models
3. Bayesian sparse linear mixed effect model
4. Fine mapping causal variants in LD correlation

# Extension 1: spike-slab prior on $\theta$

$$p(\theta | z=1) \sim N(0, 1/\tau)$$

Fat Gaussian for true effects  
(slab; magnitude and direction)

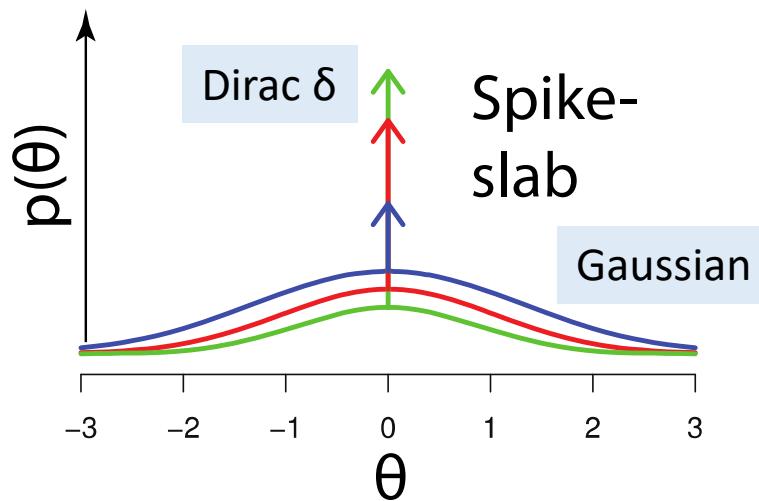
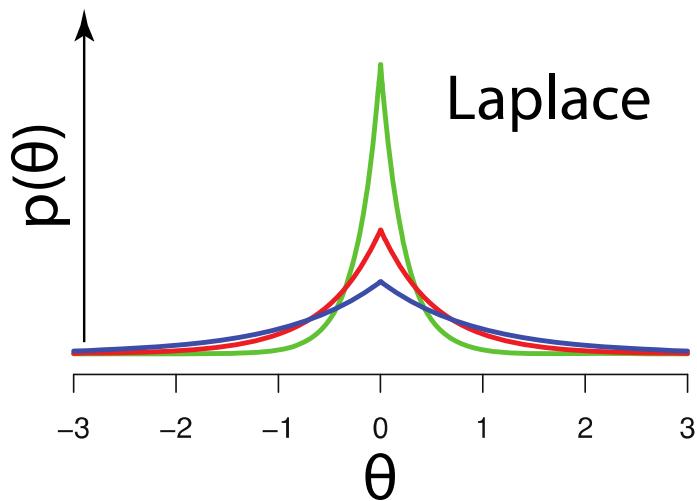
$$p(\theta | z=0) = \delta(\theta)$$

Completely set to zero  
if not selected

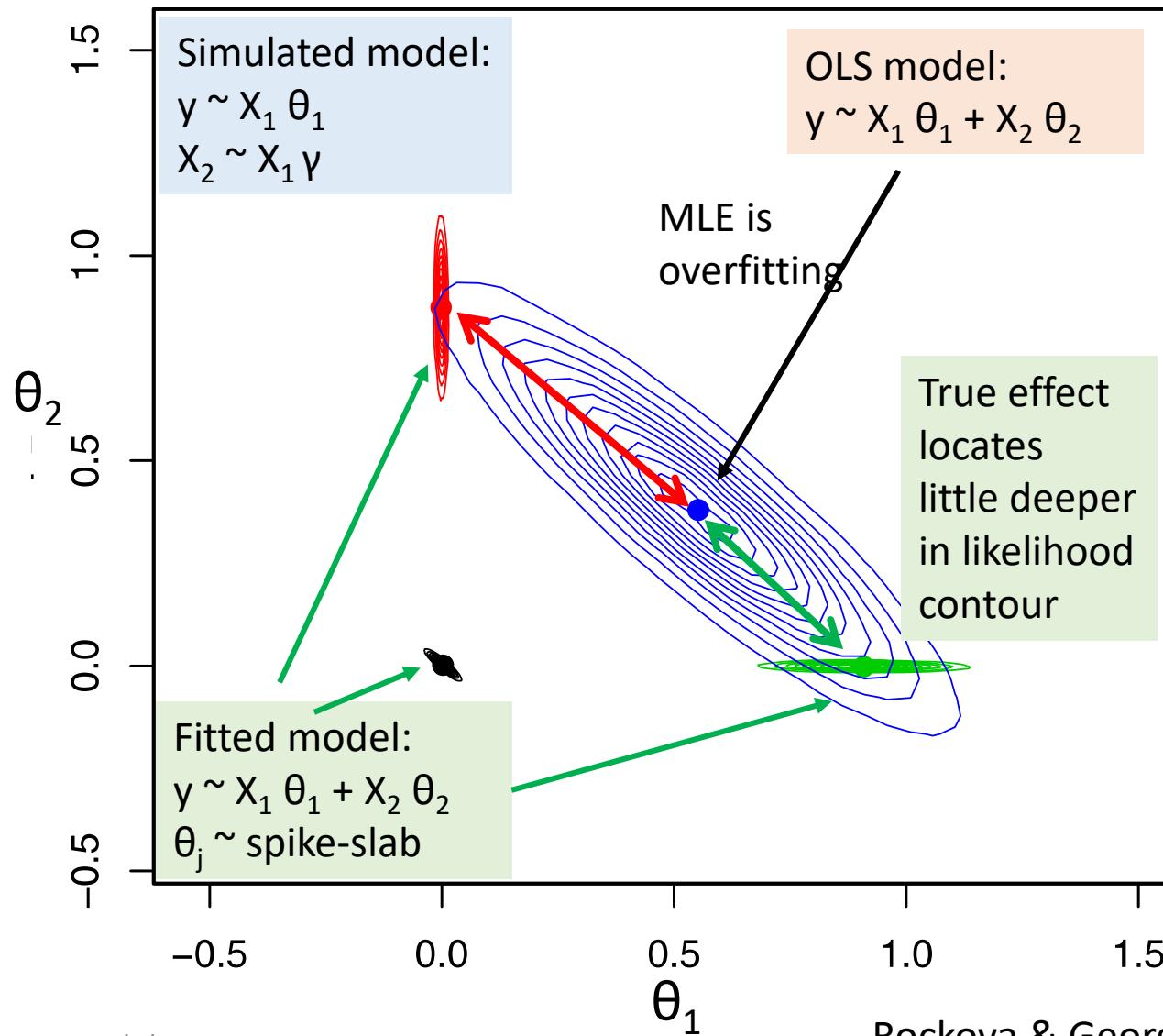
$$z = 1 \sim \text{Bernoulli}(\pi)$$

$\pi$  determines prior prob.  
of including variables  
(usually  $< .1$ ; spike;  
prescribed or optimized)

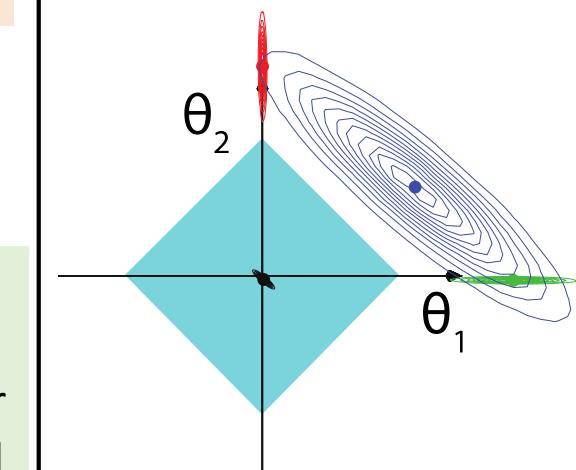
$$p(\theta) \sim \exp(-\lambda|\theta|)$$



# Spike-slab prior model effectively avoid colinearity



Can L1-regularized one handle this?



If correlation between  $X_1 \sim X_2$  is strong, probably not ...  
(best solution within the box is still non-zero for both vars).

# Ext 2: random-effect for pop. stratification

Additive effect of random vector  $u$  ( $n \times 1$ ):

$$\mathbf{y} = X\boldsymbol{\theta} + \boxed{\mathbf{u}} + \boldsymbol{\epsilon}$$

The random effect captures population structure  $K$  (kinship matrix):

$$\boxed{\mathbf{u}} \sim \mathcal{N}(0, \tau^2 \boxed{K})$$

$n \times n$   
covar.  
( $\sim$ PCs)

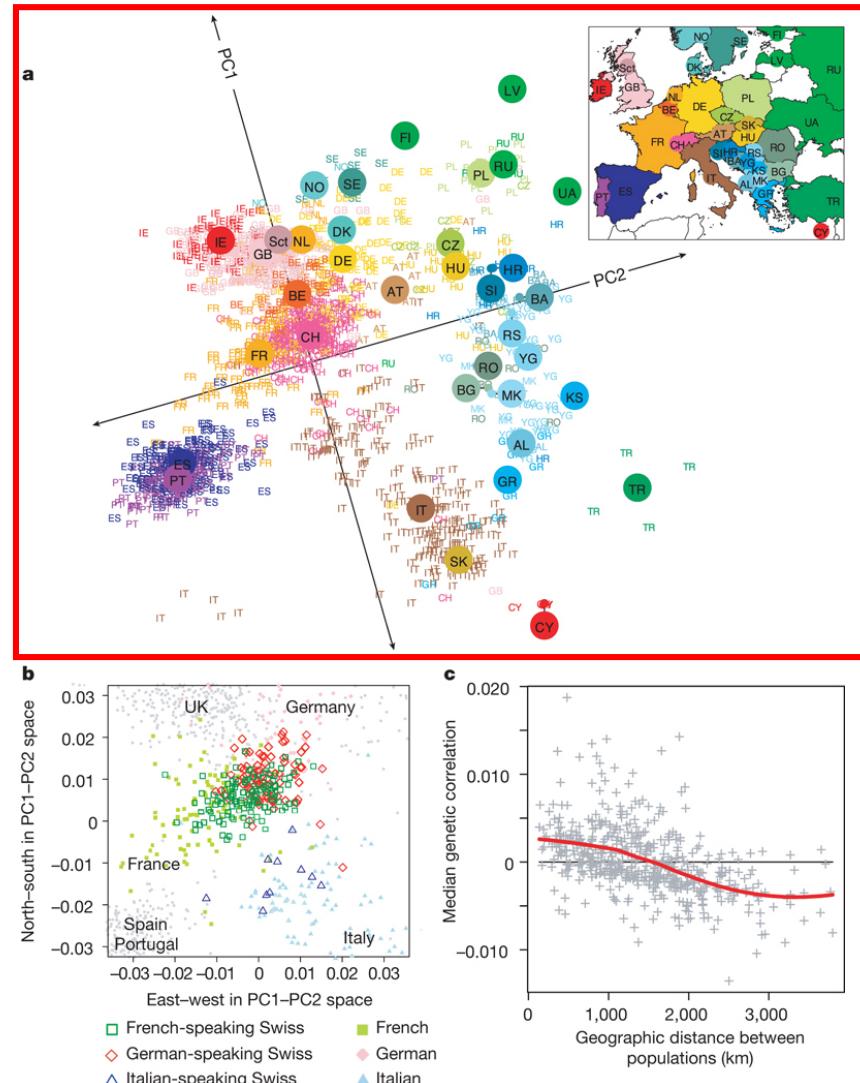
Integrate out uncertain random effect  $u$ :

$$\int p(\mathbf{y}|X, \boldsymbol{\theta}, \mathbf{u})p(\mathbf{u}|\boldsymbol{\tau}, K)d\mathbf{u} \\ = \mathcal{N}(\mathbf{y}|X\boldsymbol{\theta}, \tau^2 K + \sigma^2 I)$$

population  
structure

random noise

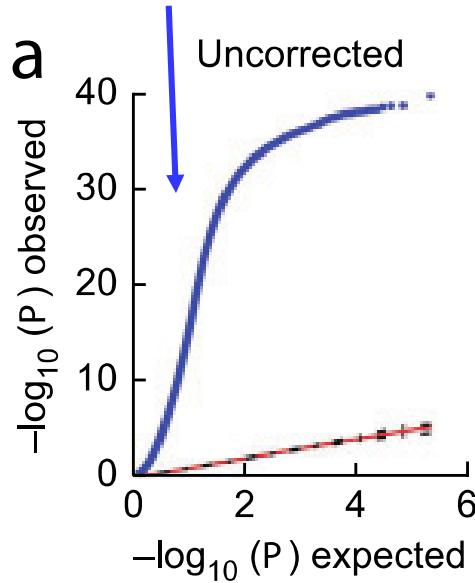
Linear Gaussian model with two variance components.



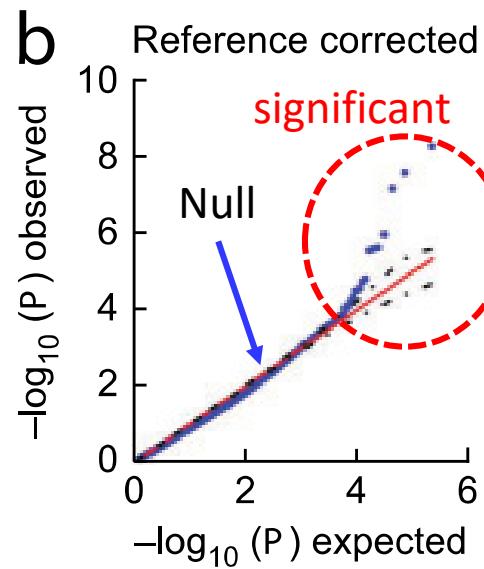
J Novembre *et al.* *Nature* 000, 1-4 (2008)

# Extension 2: random effect model

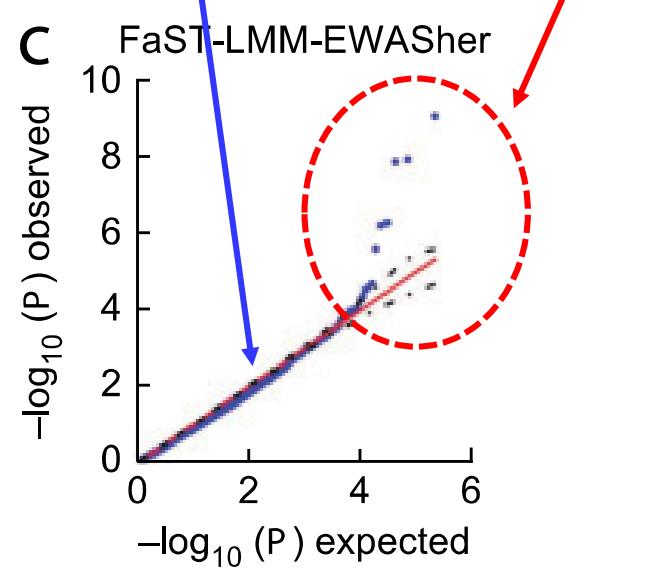
Inflated statistics  
due to unknown  
population structure  
(almost all loci are  
significant)



Adjusted GWAS  
qq-plot with  
correct  
structure



Linear mixed-  
effect  
calibrated the  
null distrib.



LMM can  
correctly  
capture  
significant  
ones.

# Extension 3: Bayesian sparse linear mixed effect model

Random effect

$$\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{u} + \boldsymbol{\epsilon},$$

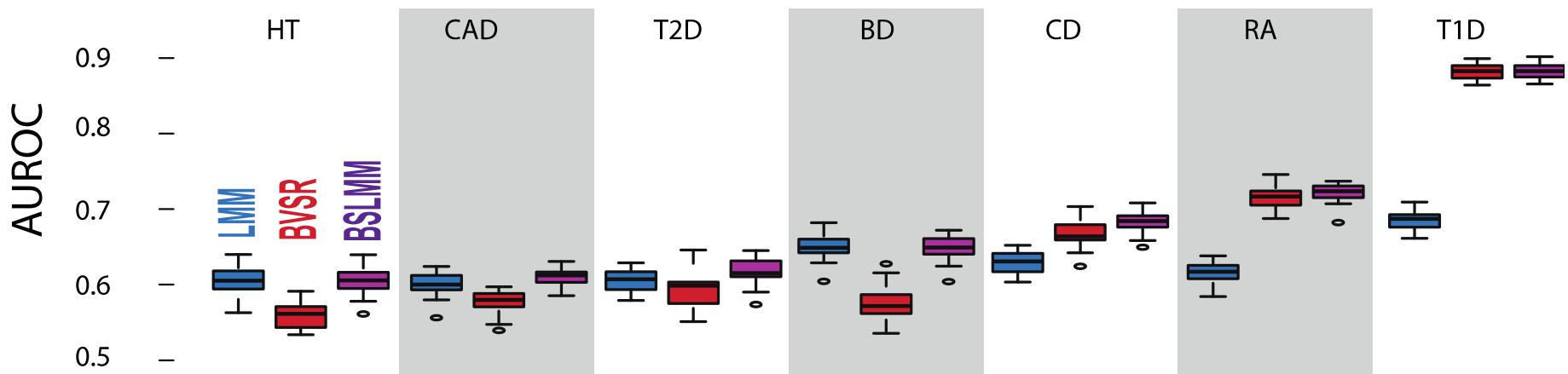
$$\mathbf{u} \sim \mathcal{N}(0, K),$$

A sort of spike-slab (two mixture model)

$$\theta_j \sim \pi\mathcal{N}(0, \tau_1^2) + (1 - \pi)\mathcal{N}(0, \tau_2^2)$$

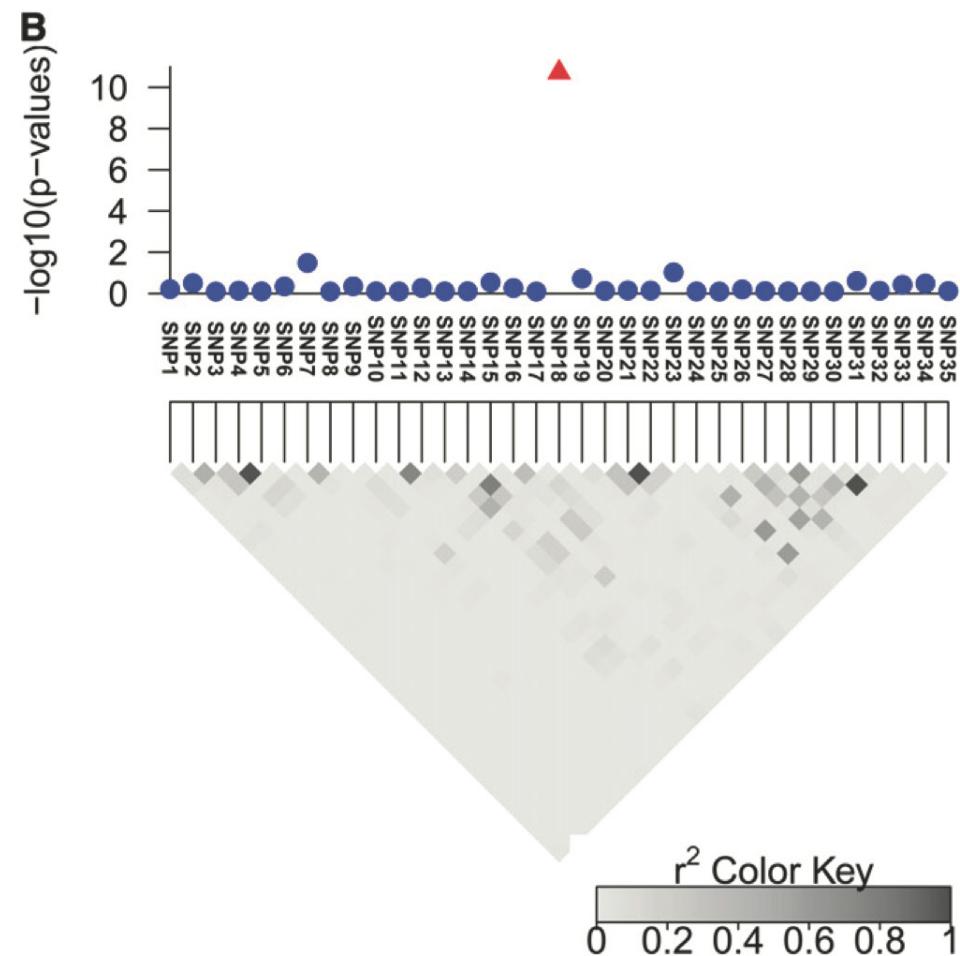
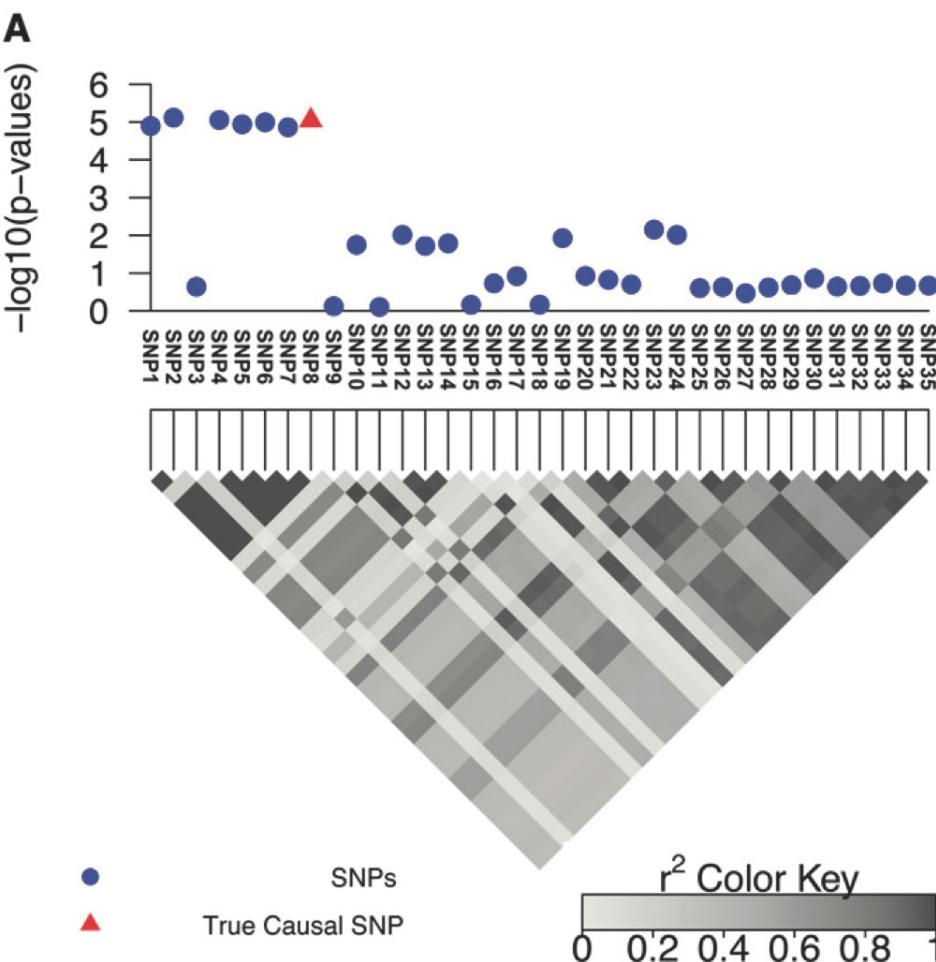
causal effect

infinitesimal  
background effect



Zhou, Carbonetto, Stephens, *PLoS Gen.* (2013)

# Extension 4: Fine-mapping causal variants



Hormozdiari *et al.* (2014)

# Extension 4: Fine-mapping under the hood

summary z-score obs.

unknown genotype

unkonwn phenotype y vector

$$\mathbf{z} \approx \mathbf{X}^T \mathbf{y} / \sqrt{n} \sigma$$

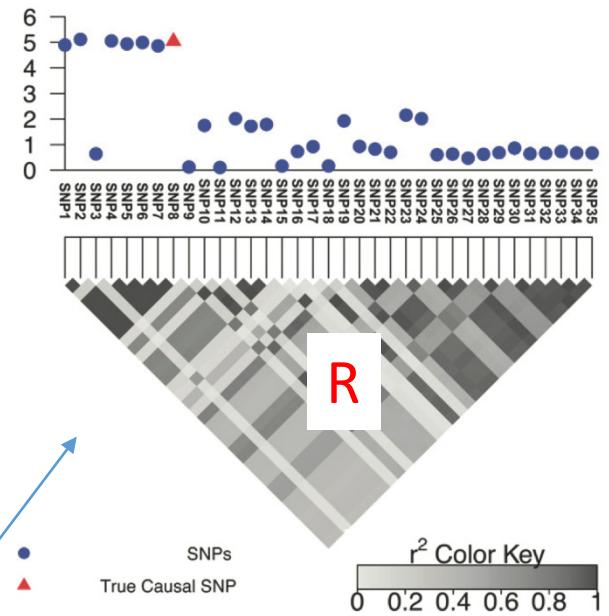
We assume phenotype vector were generated by

$$\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\theta}, \sigma^2 I).$$

Therefore  $p \times 1$  vector follows

$$\mathbf{z} \sim \mathcal{N}\left(\frac{\mathbf{X}^T \mathbf{X} \boldsymbol{\theta}}{\sqrt{n} \sigma}, \frac{\mathbf{X}^T \mathbf{X}}{n}\right) \approx \mathcal{N}(\lambda \mathbf{R} \boldsymbol{\theta}, \mathbf{R}).$$

where LD matrix  $R = n^{-1} \mathbf{X}^T \mathbf{X}$  and  $\lambda = (n\sigma^2)^{-1/2}$  absorbs all scaling factors.



- Considering potential colinearity embedded in the R matrix,  $\boldsymbol{\theta}$  desperately needs spike-slab prior.
- For computational efficiency, previously developed algorithms restrict number of causal variants (e.g., at most 3).

Hormozdiari *et al.* (2014)

# Bayesian inference algorithms

	Exact inference	Markov Chain Monte Carlo	Variational Bayes
Accuracy	correct	approximate, stochastic	approximate, deterministic
Convergence	sure	Global optima at equilibrium	Local optima in finite time
Flexibility	very limited	high	high
Examples	HMM's forward-backward, Dynamic programming	Importance sampling, Metropolis-Hastings, Gibbs, Hamiltonian MC, Elliptical slice sampling	Laplace, Mean-field approx., Belief propagation, Expectation propagation

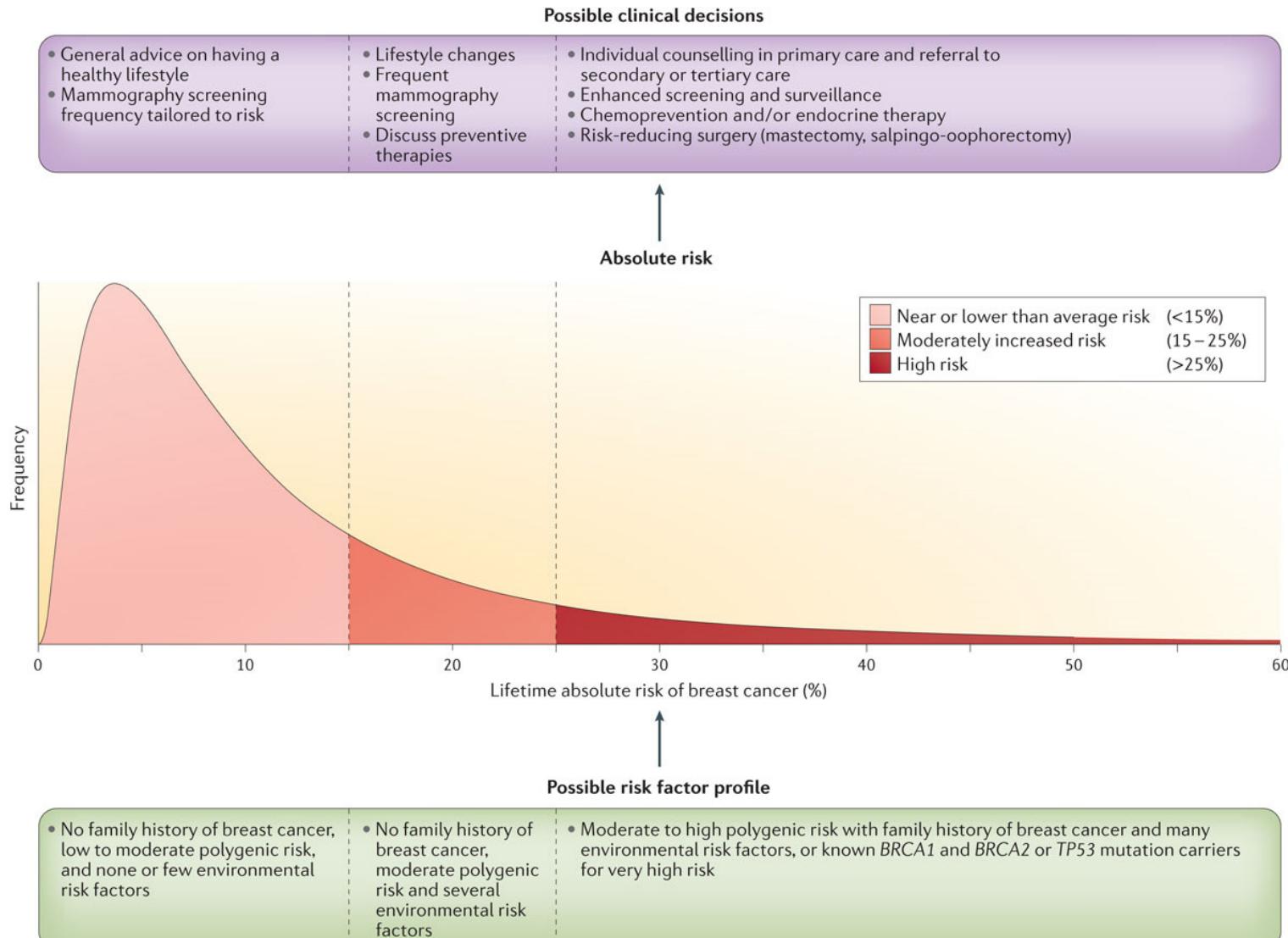
# Today: Deep Learning for Human Genetics and Disease

1. Review: GWAS, fine-mapping, Bayesian variant prioritization
2. Deep Learning for GWAS: calling SNPs, prioritize function
3. eQTLs/Mediation: intermediate molecular phenotypes
4. Linear Mixed Models (LMMs) for GWAS and for eQTL calling
5. Polygenic Risk Scores (PRS): summing over many variants
6. Heritability: definition(s), missing heritability, partitioning
7. LD SCore regression (LDSC) for fast heritability partitioning
8. Polygenic/Omnigenic disease models: core vs. periphery
9. Disease gene networks from GWAS evidence boosting

## **5. Polygenic Risk Scores (PRS):**

Summing over all variants (and more)

# Estimate absolute risk combining genetic and environmental risk factors



# How do we estimate polygenic risk score?

Univariate GWAS statistics teach us:

$$\beta_j = \log(\text{odds ratio of SNP } j)$$

$$g_j = \text{genotype (dosage)}$$

Predict overall risk by combining many, many variants!

$$\text{PRS} = \sum_{j \in \{\text{SNPs}\}} \beta_j g_j$$

**Can we just combine all the SNPs? Why not?**

- Is correlation between  $g_1$  and  $g_2$  zero?
- Can we trust the estimate  $\beta$  of all the SNPs?
- Can we just select GWAS significant SNPs?

# A common practice of PRS estimation

Univariate GWAS statistics:

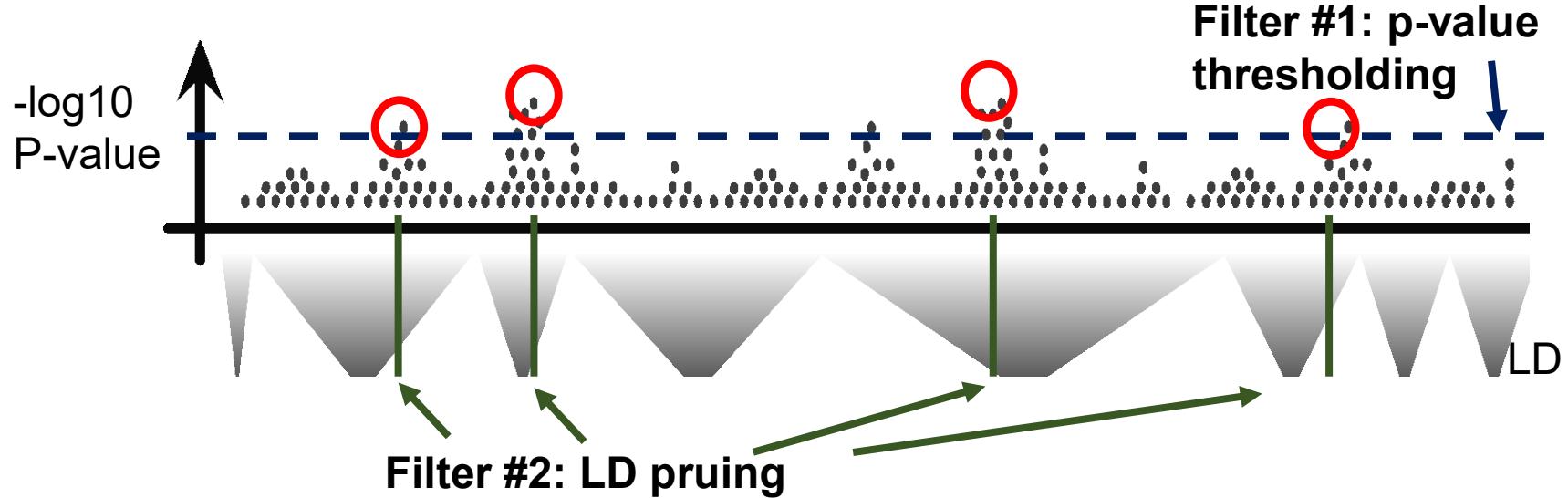
$$\beta_j = \log(\text{OR of SNP } j)$$

$g_j$  = genotype (dosage)

PRS model:

$$\text{PRS}[i] = \sum_{j \in \{\text{SNPs}\}} \beta_j g_j[i]$$

Goal: Tuning this parameter



# A common practice of PRS estimation: Cross-validation with observed phenotype

Univariate GWAS statistics:

$$\beta_j = \log(\text{OR of SNP } j)$$

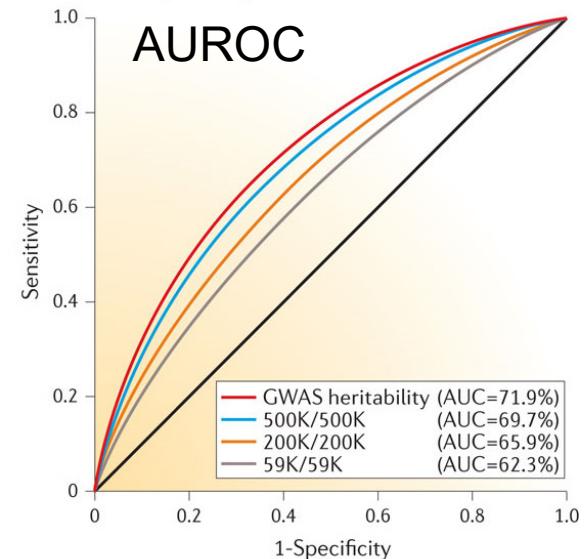
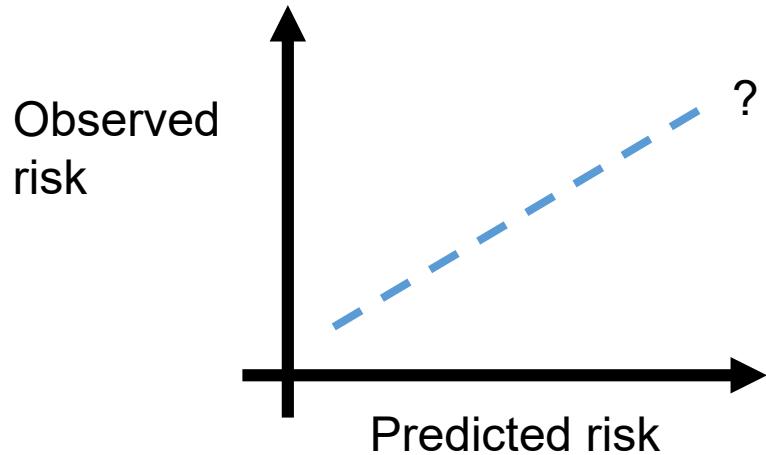
$g_j$  = genotype (dosage)

PRS model:

$$\text{PRS}[i] = \sum_{j \in \{\text{SNPs}\}} \beta_j g_j[i]$$

Goal: Tuning this parameter

How do we know the selected SNPs are good?



# An alternative method for estimating PRS (and a simpler and more powerful way)

Univariate GWAS statistics:

$$\beta_j = \log(\text{OR of SNP } j)$$

$g_j$  = genotype (dosage)



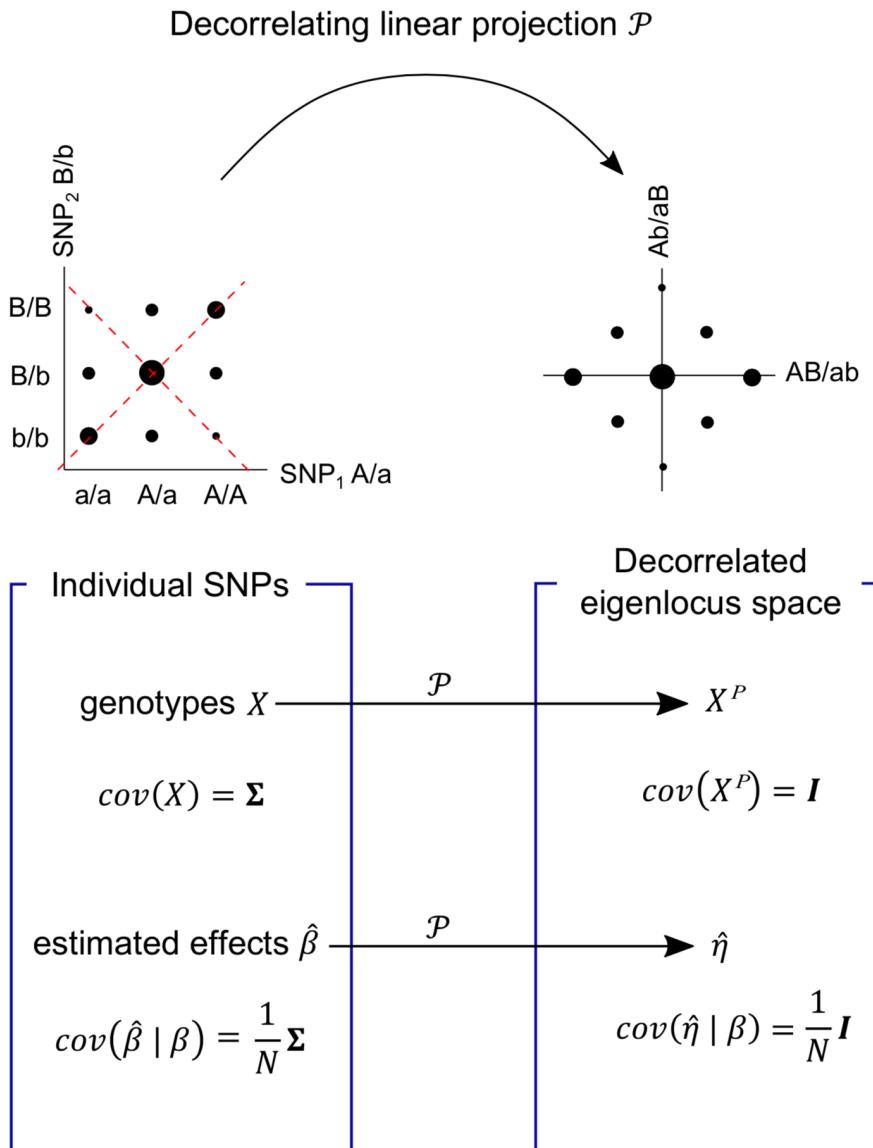
PRS model:

$$\text{PRS}[i] = \sum_{j \in \{\text{SNPs}\}} \beta_j g_j[i]$$



What's wrong with using all  
the SNPs? LD between them.  
Adjust spurious weak effects.

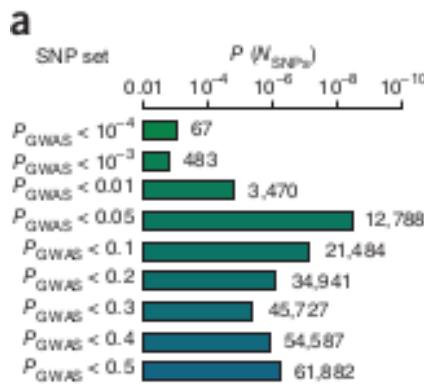
# Idea: Decorrelate LD structure



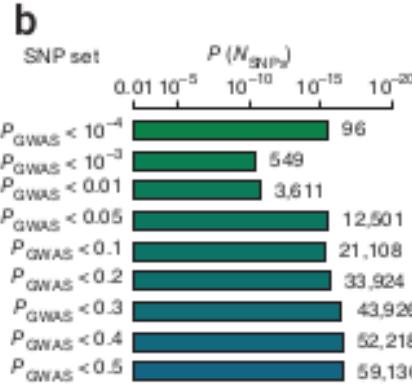
- Transform SNP space to multi-SNP space (SVD)
- Select independent & orthogonal factors.
- Or regularize eigenvalues to smooth out spurious associations.
- We don't need much tuning with regularization.

Chun .. Sunyeav, BioRxiv (2019)  
Baker *et al.*, Genetic Epidemiology (2017)

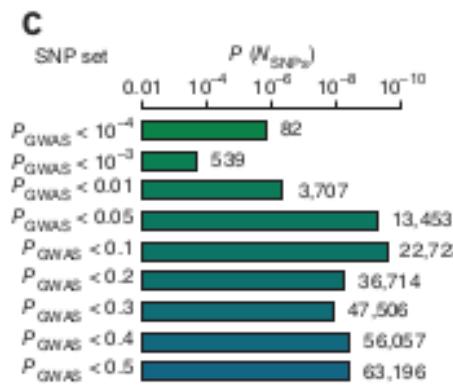
# Polygenic risk scores



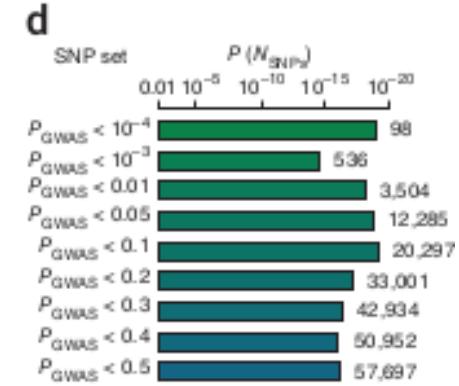
Rheumatoid  
Arthritis



Celiac



Myocardial  
infarction



Coronary  
artery disease

- Aggregate burden of sub-threshold SNPs to improve prediction performance (Stahl 2012)
- As we include more SNPs in the risk score, the association with RA, celiac disease, MI, CAD gets stronger
- In practice, requires tuning of p-value threshold, LD pruning threshold

# Phasing diploid genomes is hard

- Humans are **diploid** organisms
- Each individual carries two **homologous** copies of each chromosome
- Therefore, they carry two copies of each variant (called the **maternal/paternal allele**)
- Variants co-occur in **haplotypes** which are inherited as a unit
- Experimentally possible, but currently infeasible, to directly measure haplotypes over the whole genome
- Cheaper and more efficient to measure **genotypes** (counts of minor allele)
- Genotyping loses information, which we need algorithms and statistical models to recover (**phasing, imputation**)

## Haplotypes

0 0 1 0 1 1 0 (maternal)

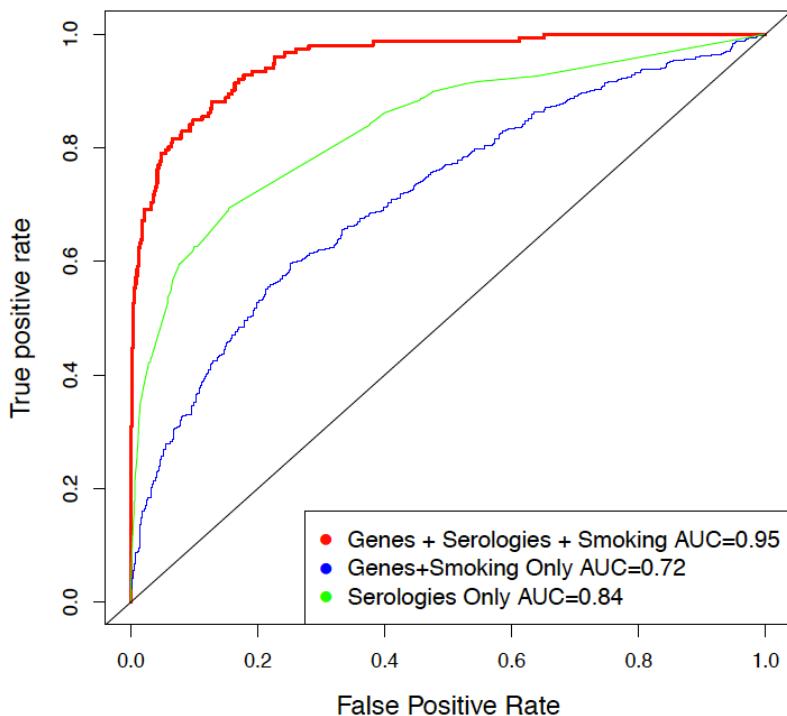
0 1 1 0 0 1 0 (paternal)

## Genotypes

0 1 2 0 1 2 0

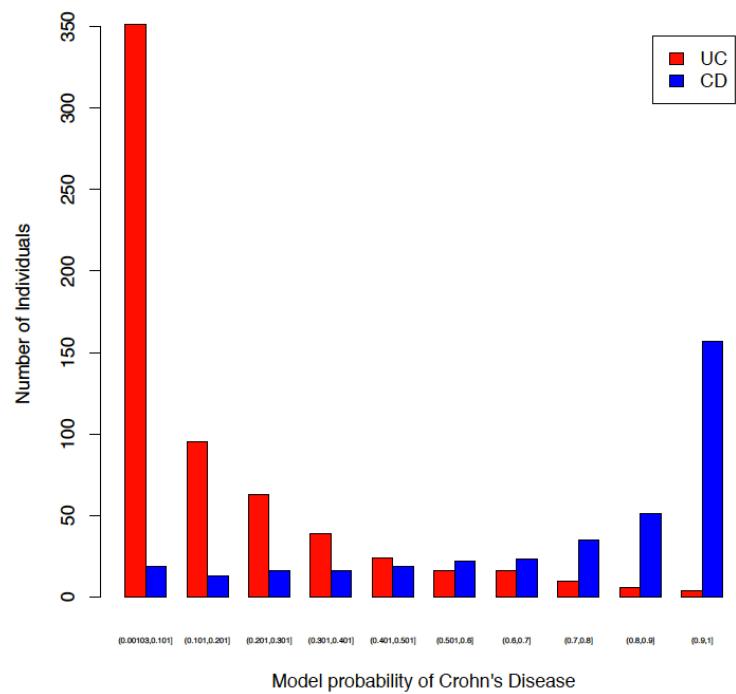
# Molecular diagnostics in IBD

ROC Curves For A Model That Discriminates CD from UC Patients



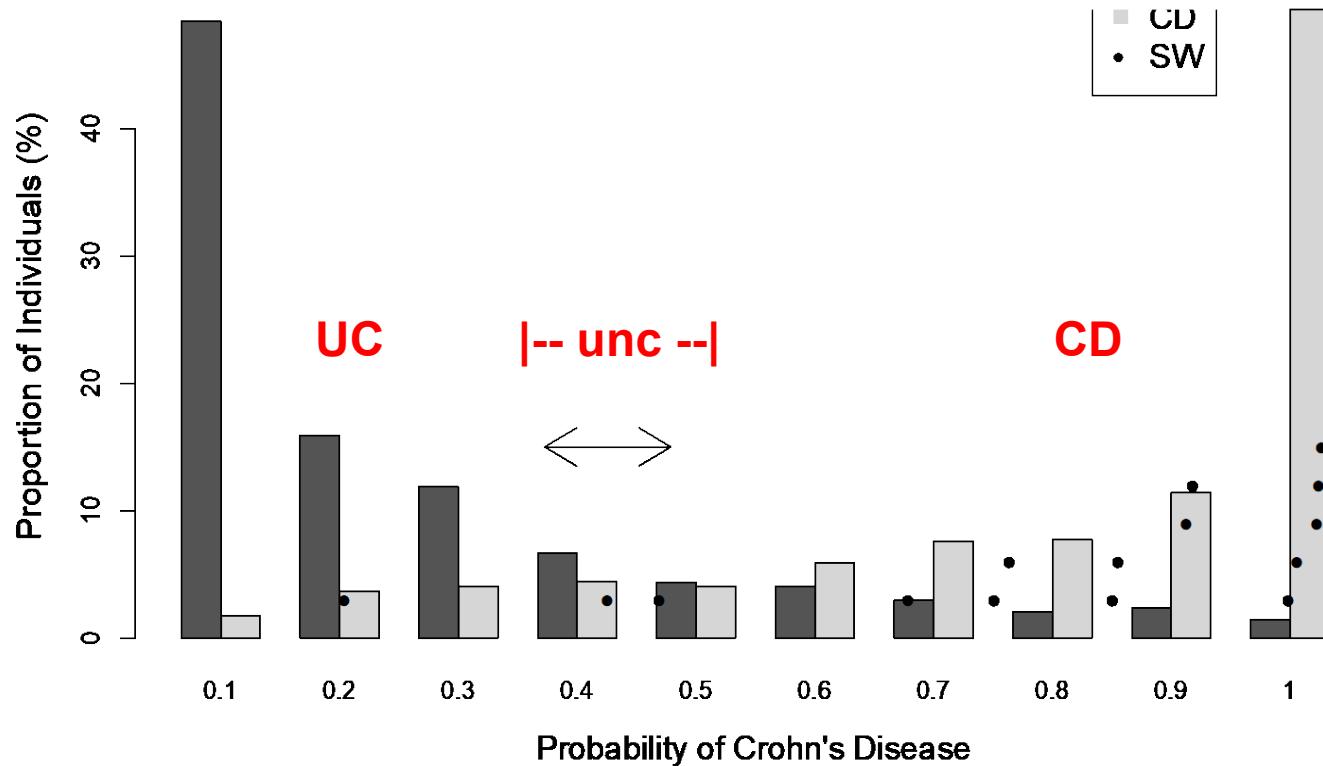
'Molecular' diagnosis (based on GWAS SNPs & serologic biomarkers)  
concordant with GI dx: CD & UC  
patients can be distinguished accurately

Model Calibration



>90% of patients correctly classified  
with >90% reliability

# Molecular diagnostics flag patients with worst outcome



Black dots represent patients diagnosed with UC who later underwent colectomy and then developed full-blown Crohn's disease

# Today: Deep Learning for Human Genetics and Disease

1. Review: GWAS, fine-mapping, Bayesian variant prioritization
2. Deep Learning for GWAS: calling SNPs, prioritize function
3. eQTLs/Mediation: intermediate molecular phenotypes
4. Linear Mixed Models (LMMs) for GWAS and for eQTL calling
5. Polygenic Risk Scores (PRS): summing over many variants
6. Heritability: definition(s), missing heritability, partitioning
7. LD SCore regression (LDSC) for fast heritability partitioning
8. Polygenic/Omnigenic disease models: core vs. periphery
9. Disease gene networks from GWAS evidence boosting

# **6. Heritability:**

## Definition, Missing Heritability, Partitioning

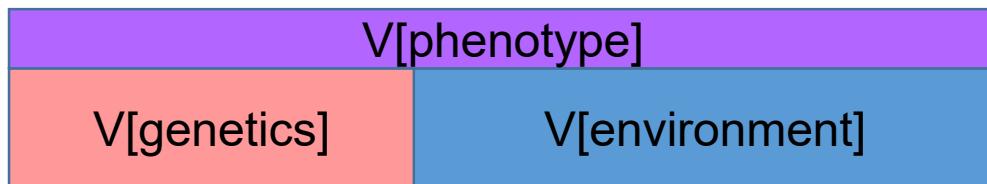
# Lessons of GWAS



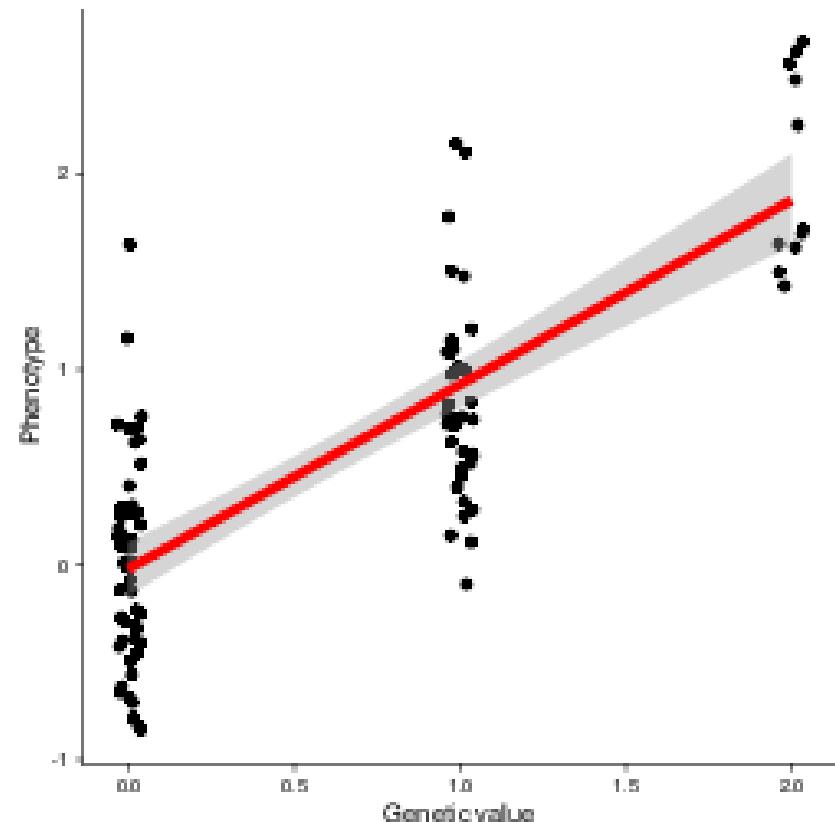
1. **We haven't found all causal loci:** known loci explain little phenotypic variance
2. **Most loci affect transcriptional regulation:** they don't tag coding variation

# Components of phenotypic variance

- Assume  $p$  (phenotype) =  $g$  (genetic) +  $e$  (environment)
- Then,  $V[p] = V[g] + V[e] + 2\text{Cov}(G,E)$   
(assume no gene-environment interactions)

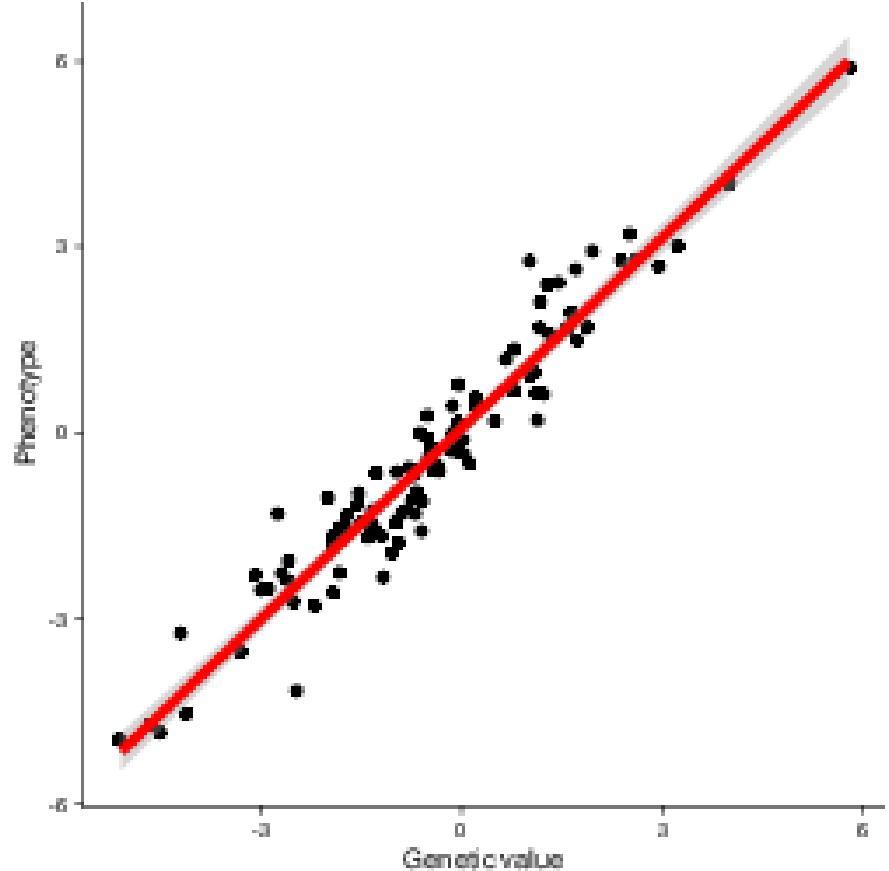


- Example: one causal variant
- Three possible **genetic values** in the population
- Intuition:  $V[g]$  is the variance of mean phenotype across different genetic values
- $V[e]$  is the variance of phenotype for the same genetic value



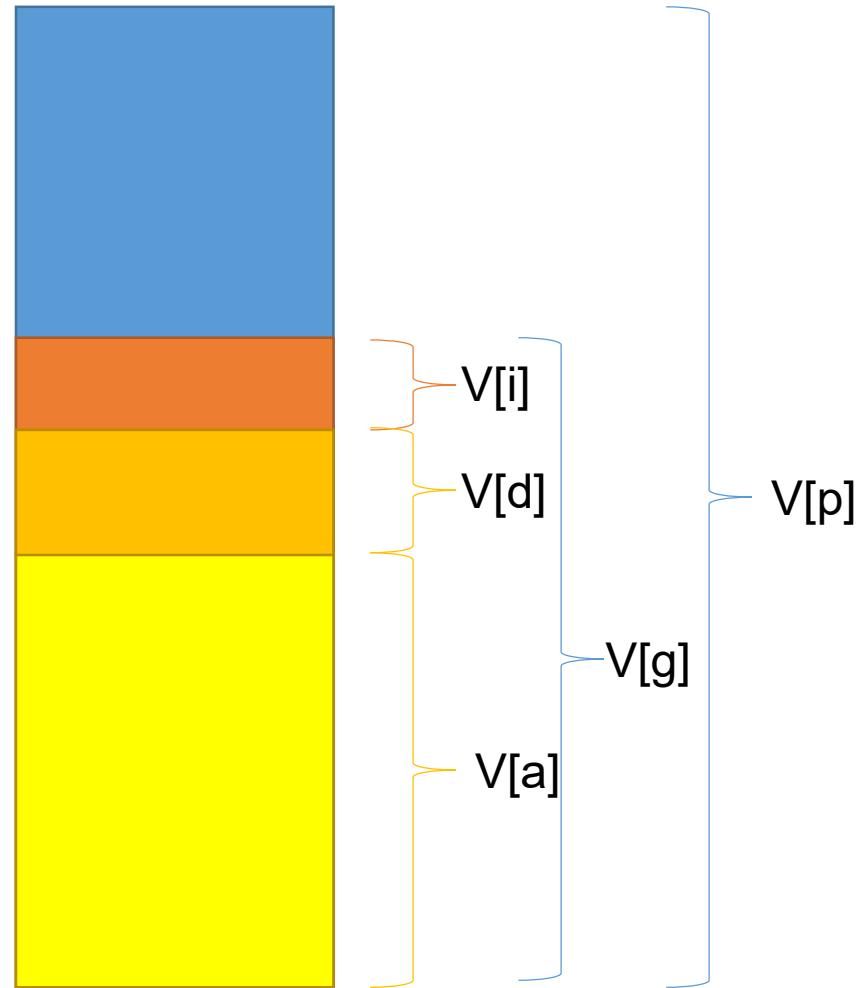
# Components of genetic variance

- Assume  $V[g] = V[a]$  (additive)  
+  $V[d]$  (dominance) +  $V[i]$   
(interactions)
- The additive component corresponds to a linear model
- As we add more causal variants, phenotypes become closer to Gaussian
- We could further decompose interactions
- We could include variance due to *de novo* mutations



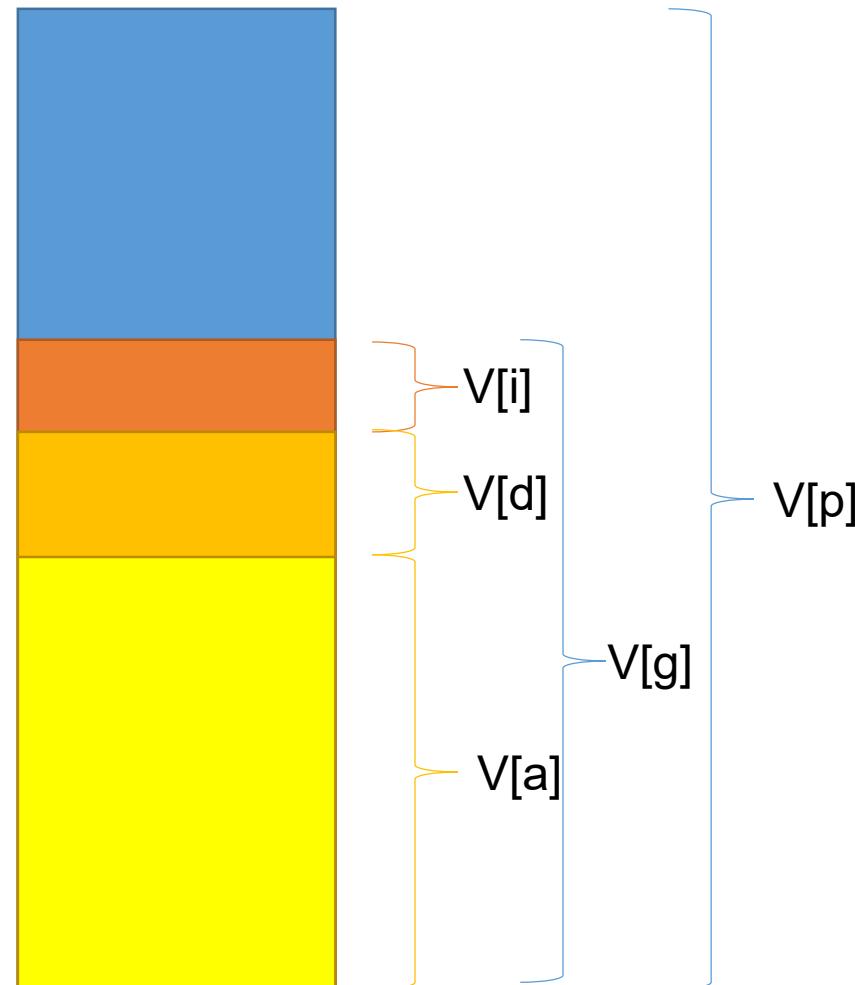
# Heritability is a ratio of variances

- $V[p] = V[g] + V[e]$
- $V[g] = V[a] + V[d] + V[i]$
- **Broad sense heritability**  
 $H^2 = V[g] / V[p]$
- Broad sense captures all genetic factors
- **Narrow sense heritability**  
 $h^2 = V[a] / V[p]$
- Narrow sense captures only additive effects
- Ongoing debate about the relative importance of additive vs. other effects in disease, selection, etc.



# Why study heritability?

- Quantify the importance of genetics vs. environment in traits of interest
- Learn about *genetic architecture*: how many causal variants, effect sizes, allele frequencies
- Narrow sense heritability is the fundamental parameter needed for phenotype prediction (and is the theoretical best possible prediction performance with a linear model)



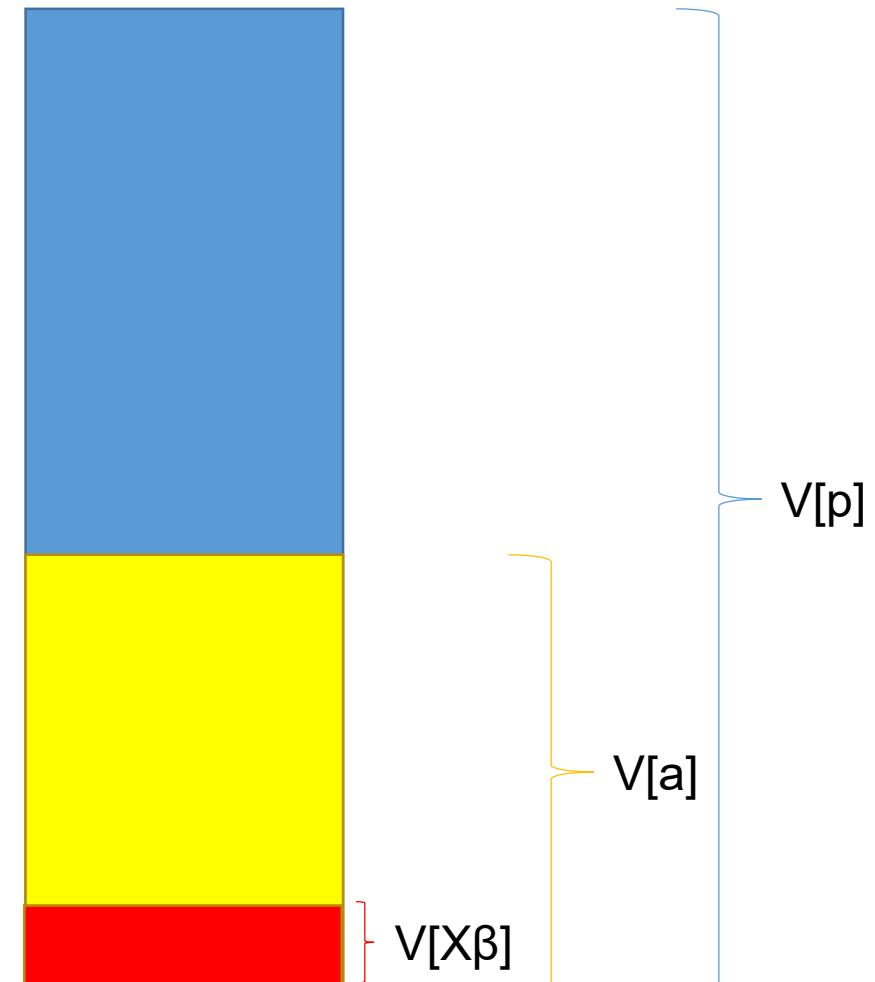
# Estimating heritability in relatives

$$p = g + e$$
$$E[p_i p_j] = h^2 E[g_i g_j]$$

- Intuition: heritability relates phenotypic correlations to genotypic correlations
- If two individuals have the same allele at each of the causal variants, they will have the same phenotype
- **Haseman-Elston regression:** fit linear regression of phenotypic correlations against genotypic correlations
- Derive genotypic correlation from family relationships: monozygotic twins share 100% of genome, siblings share 50%, etc.
- Example (height):  $h^2 = 0.73$

# Estimating heritability from GWAS

- Linear model  $g = X\beta$
- We can estimate SNP effect sizes  $\beta$  from GWAS
- The variance explained by each SNP depends on effect size and MAF
- $V[X_j \beta_j] = 2 f_j (1 - f_j) \beta_j^2$
- If we do this with genome-wide significant SNPs, we usually  $h^2_{GWAS} < h^2$
- Example (height): 253,288 samples; 697 genome-wide significant loci;  $h^2_{GWAS}=0.16$ ,  $h^2=0.73$
- Known as the **missing heritability problem**

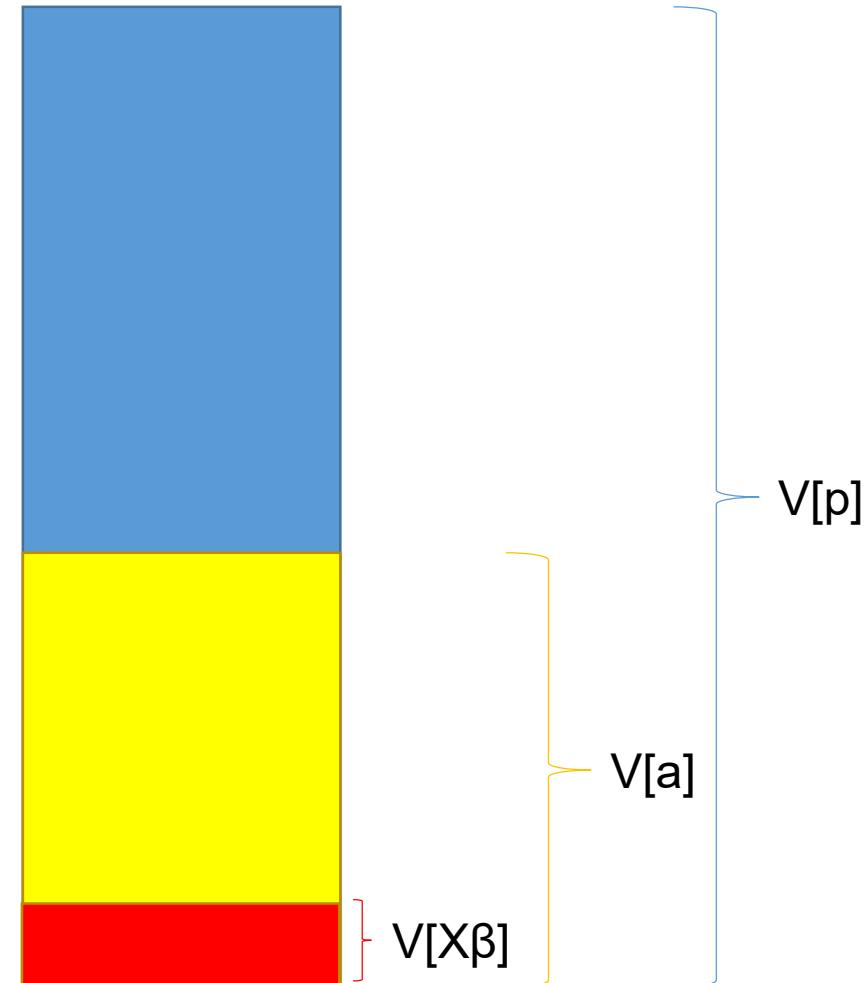


# Sources of missing heritability

Ongoing debate about several possible explanations for the missing heritability problem.

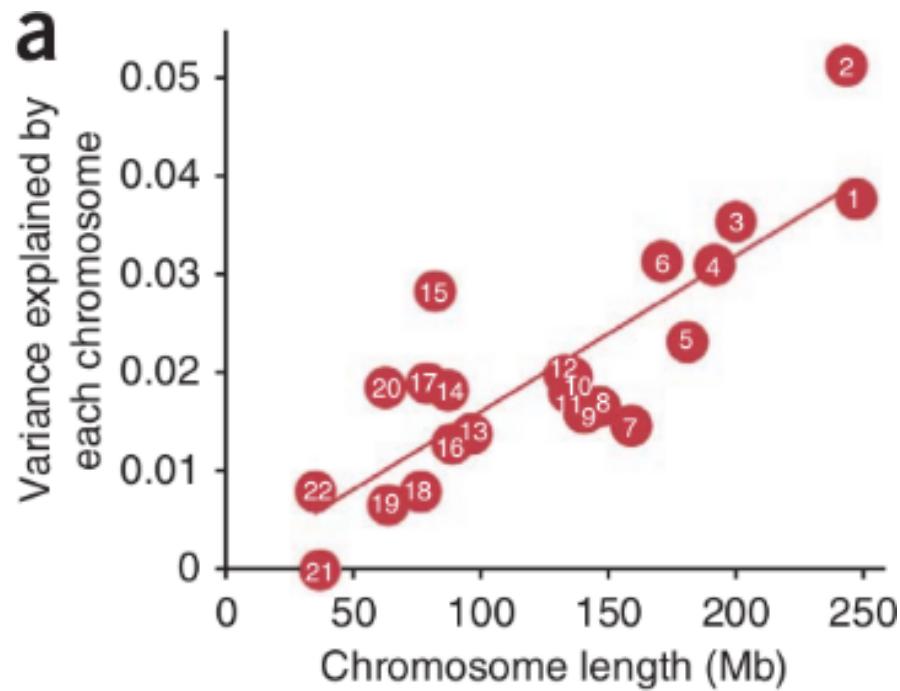
1. Many common variants, small effects
2. Unobserved rare variants, large effects
3. Wrong model assumptions

Each has very different implications for the future of human genetics studies.

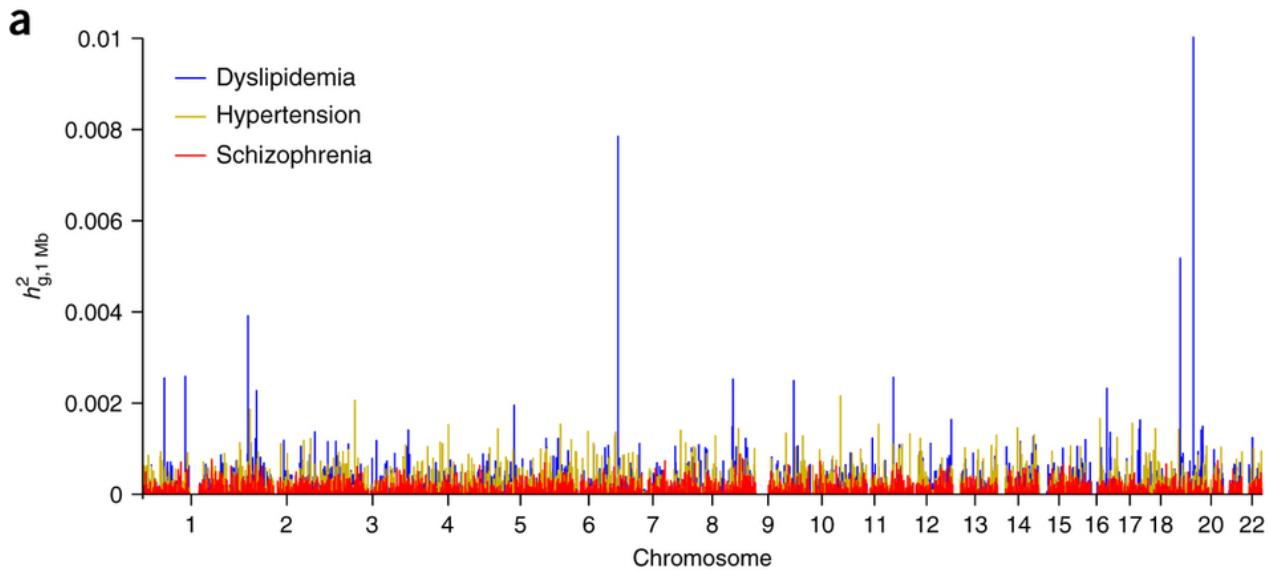


# Partitioning heritability

- Extend the model so chromosomes can explain different proportions of variance
- Intuition: add more variance parameters for each partition of SNPs
- Each partition induces a different genetic relationship matrix
- Longer chromosomes explain more heritability
- Suggests causal variants are spread uniformly through the genome

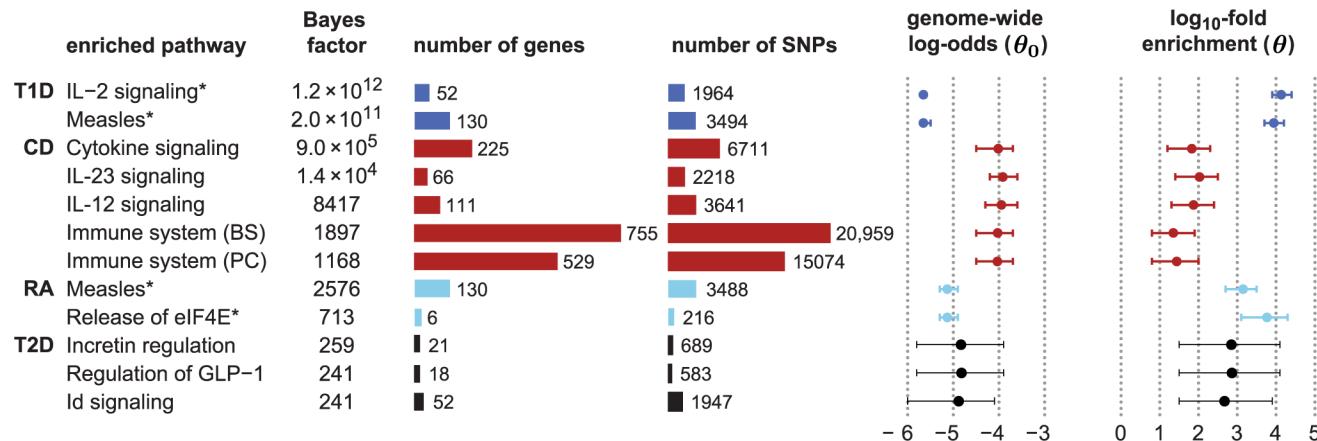


# Partitioning heritability



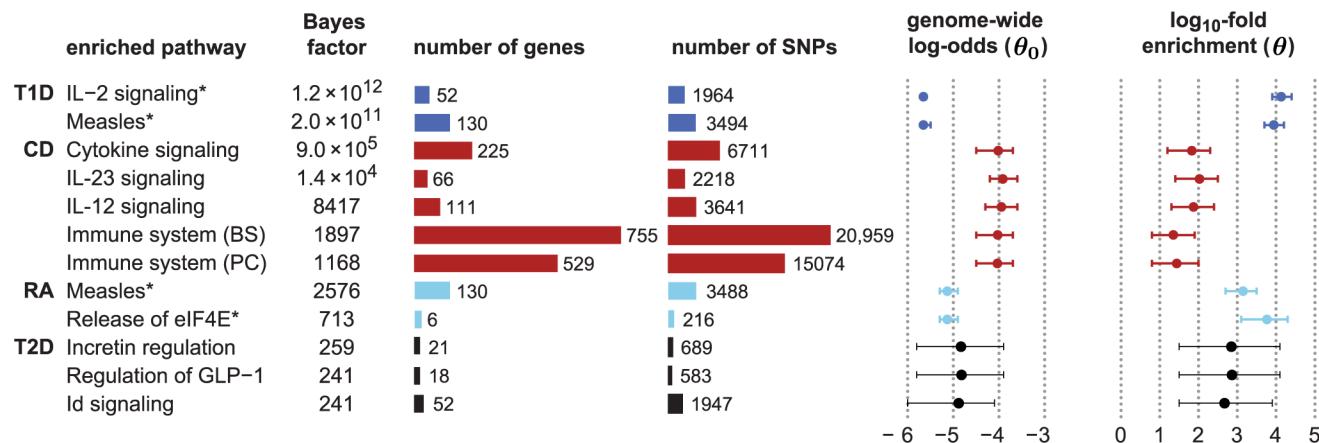
- Fit a model with one component per 1MB window (Loh 2015)
- Bound cumulative heritability explained to estimate number of regions
- Most of the genome explains non-zero heritability

# Bayesian variable selection



- Directly fitting the underlying linear model is ill-posed: we have  $n < p$  so there are infinitely many solutions
- Idea: use **spike and slab** prior to force many effects to be exactly 0 and regularize the problem (one solution)
- Inference goal: estimate the effect sizes and the level of sparsity (Carbonetto 2013)

# Pathways-informed prior from enrichments

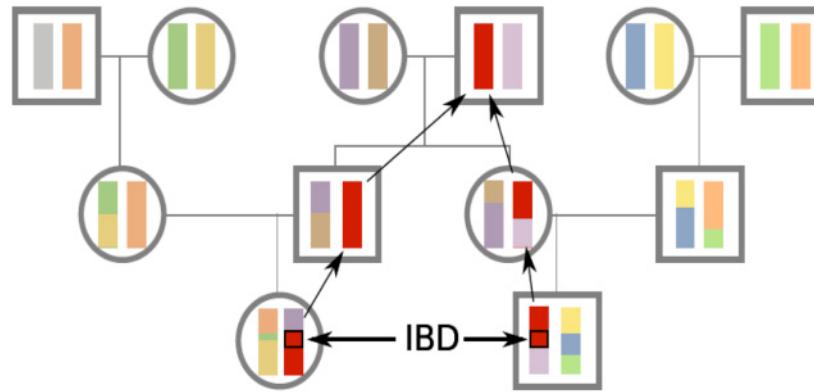


- Extension: some pathways contain more causal variants than the rest of the genome
- Incorporate into the prior
- Identifies relevant immune signaling pathways which are not found using existing methods
- Identifies tens of thousands of SNPs which could be affecting those pathways

# Evidence for other explanations

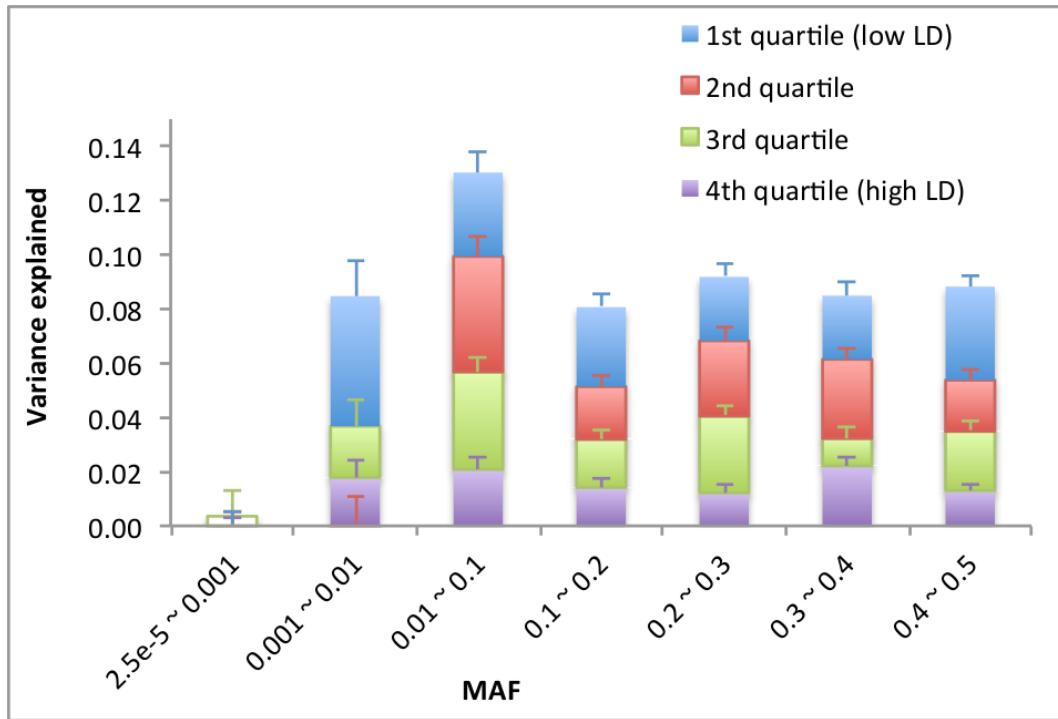
- Incorporating Identity by Descent (IBD) in unrelated individuals
- Partitioning SNPs by MAF, LD
- Assumptions do not hold in real data

# Estimating heritability: shared haplotypes



- Shared haplotypes explain more heritability than tag SNPs
- There is still a discrepancy between  $h^2_g$  and  $h^2$
- If two individuals share a chromosomal segment, unobserved variants should also be shared (Bhatia 2015)
- Idea: Identify IBD segments by quickly scanning SNPs and finding stretches of identical alleles
- Inferring shared segments captures rarer variants more effectively than LD

# Partitioning SNPs by MAF/LD



- Low frequency/low LD variants are poorly tagged by observed/imputed variants, so estimate variance for them separately (Yang 2015)
- Partitioning appears to explain all of the heritability of height using only common/low frequency variants!

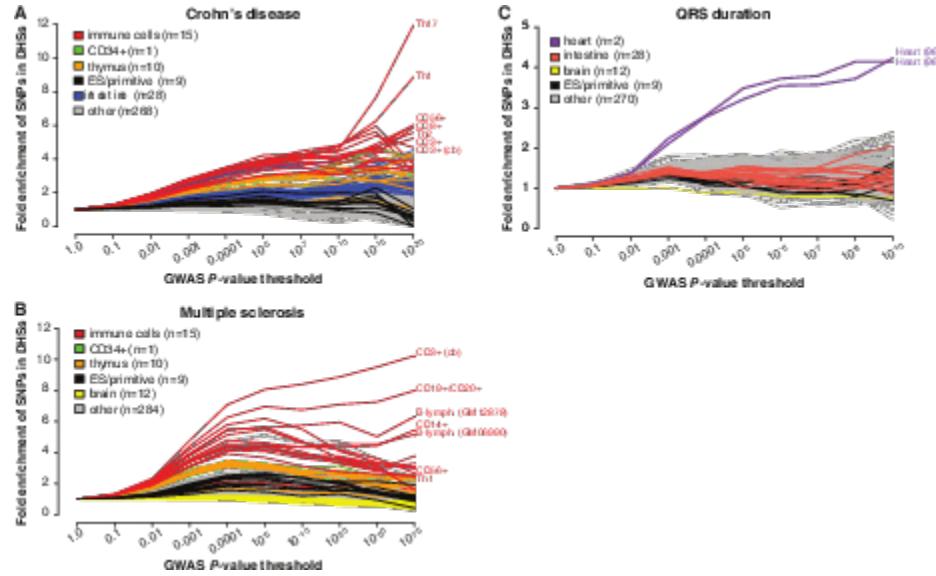
# Examining model assumptions

- Phenotypes might not be Gaussian
- GWAS samples are not independent and identically distributed
- SNPs are not independent
- Not all SNPs have an effect
- Not all causal SNPs have equal effects
- There are gene-environment interactions
- There are gene-gene interactions

# Limitations of heritability

- Explaining all of the heritability of complex traits is not enough
- As sample size goes to infinity, will the entire genome be associated with all traits? (Goldstein 2009)
- **Goal:** Find biological pathways recurrently disrupted by non-coding variation

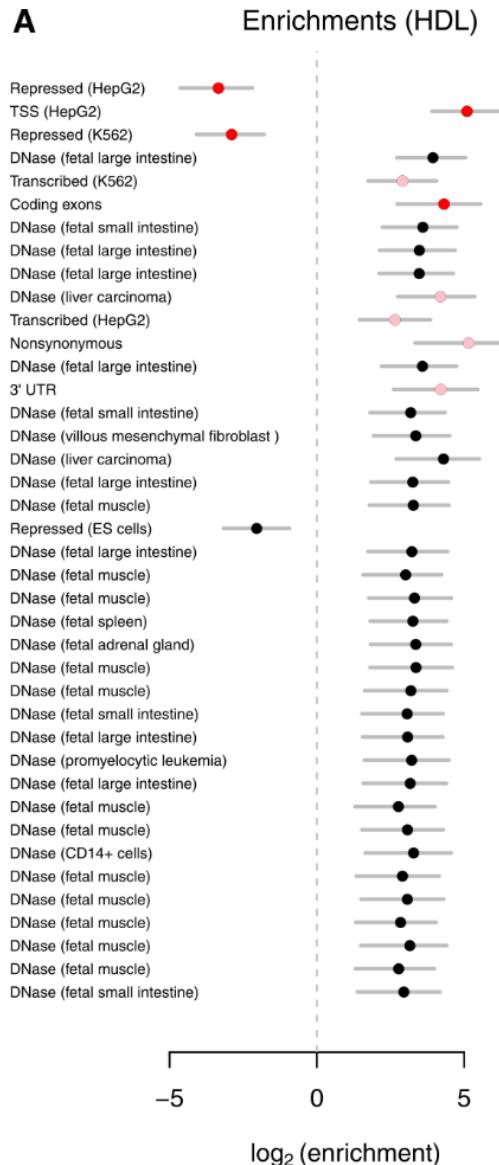
# Regulatory enrichments



- Weakly associated variants overlap accessible chromatin more often than expected by chance (Maurano 2012)
- Same trend observed in other predicted regulatory elements: histone peaks, ChromHMM segments, super enhancer clusters

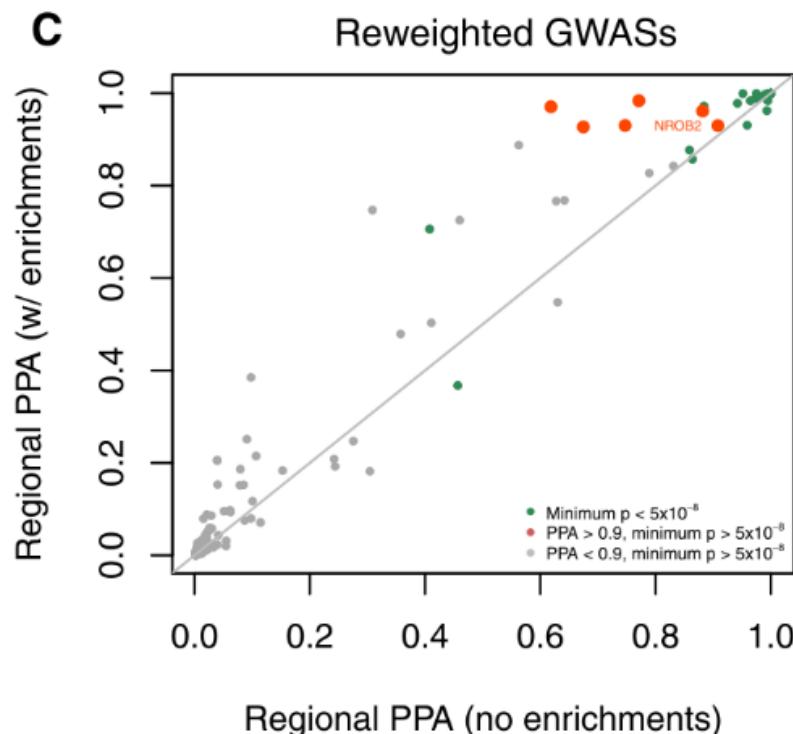
# Joint model of SNPs and annotations

- Use **penalized stepwise regression** to pick relevant annotations (Pickrell 2014)
- Use approximate Bayes factors to compute posterior probability of association
- Forward steps: add annotations to the model until they don't explain enough variance
- Backward steps: remove annotations from the fitted model until variance explained drops too much



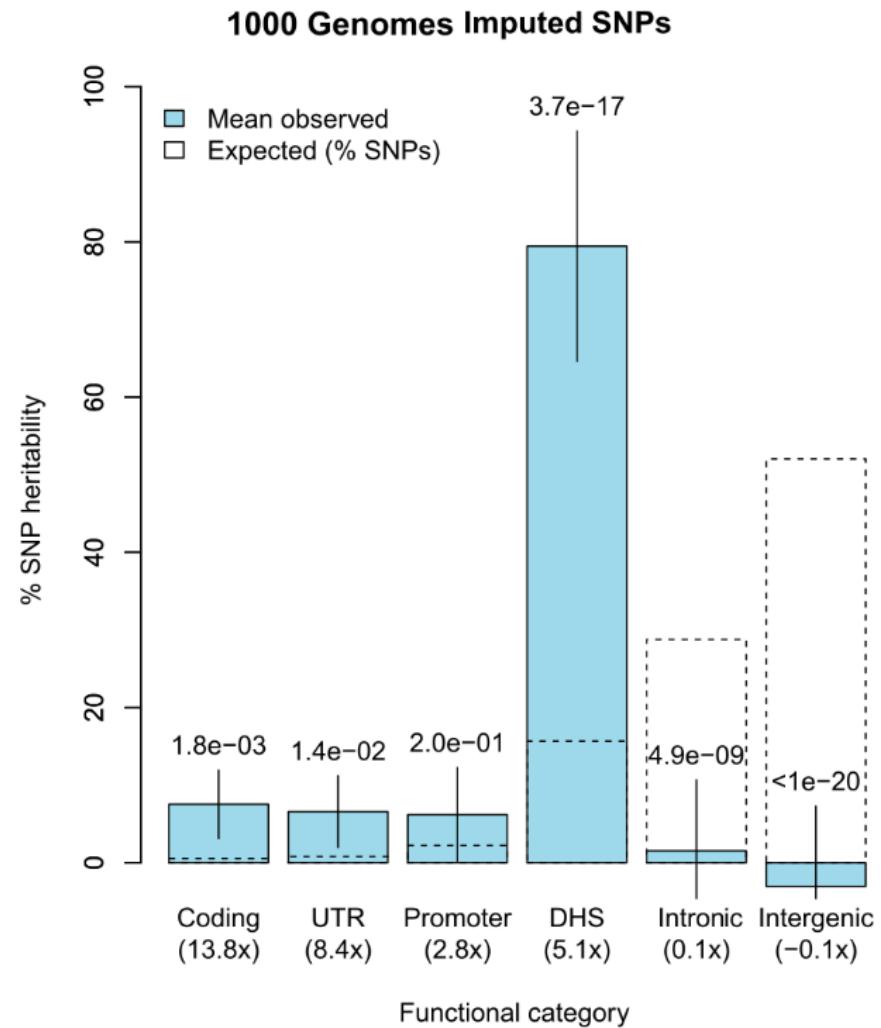
# Joint model of SNPs and annotations

- Use approximate Bayes factors to compute posterior probability of association
- Posterior probability of association re-prioritizes new GWAS loci



# Partitioning heritability by annotation

- Accessible chromatin explains more heritability
- Combine DHS in >100 cell types: 70% of genome is accessible in some cell type, but only 16% is accessible in multiple cell types
- Implies non-coding SNPs explain more variance per SNP than coding SNPs



# Today: Deep Learning for Human Genetics and Disease

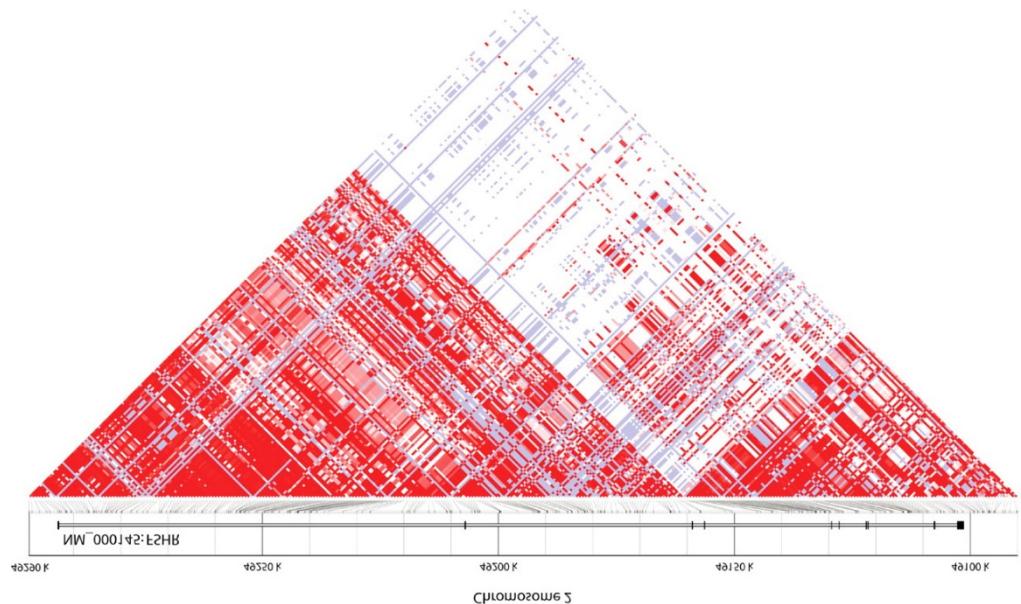
1. Review: GWAS, fine-mapping, Bayesian variant prioritization
2. Deep Learning for GWAS: calling SNPs, prioritize function
3. eQTLs/Mediation: intermediate molecular phenotypes
4. Linear Mixed Models (LMMs) for GWAS and for eQTL calling
5. Polygenic Risk Scores (PRS): summing over many variants
6. Heritability: definition(s), missing heritability, partitioning
7. LD SCore regression (LDSC) for fast heritability partitioning
8. Polygenic/Omnigenic disease models: core vs. periphery
9. Disease gene networks from GWAS evidence boosting

## **7. LD SCore regression (LDSC):**

Computing and partitioning\* heritability quickly  
(\* with stratified LD SCore regression)

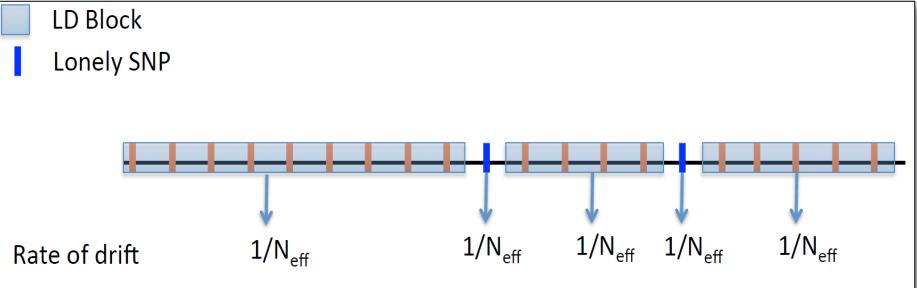
# LD SCore regression (LDSC)

$$E[z_j^2] = N l_j h^2 / M$$

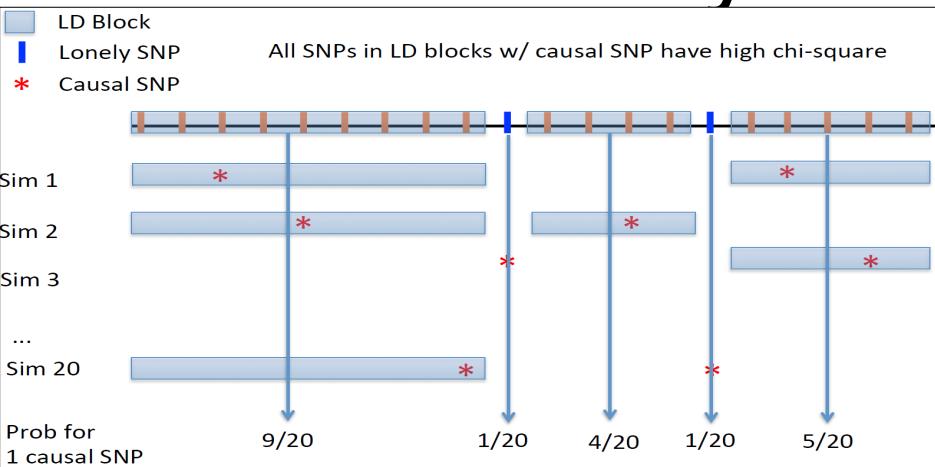


- Intuition: Causal variants drawn uniformly at random from the genome are more likely to come from larger LD blocks (Bulik-Sullivan 2014)
- Linear regression of summary statistics against LD score gives  $h^2$  without access to individual-level genotype matrix

# Intuition: LD score $\leftrightarrow$ heritability



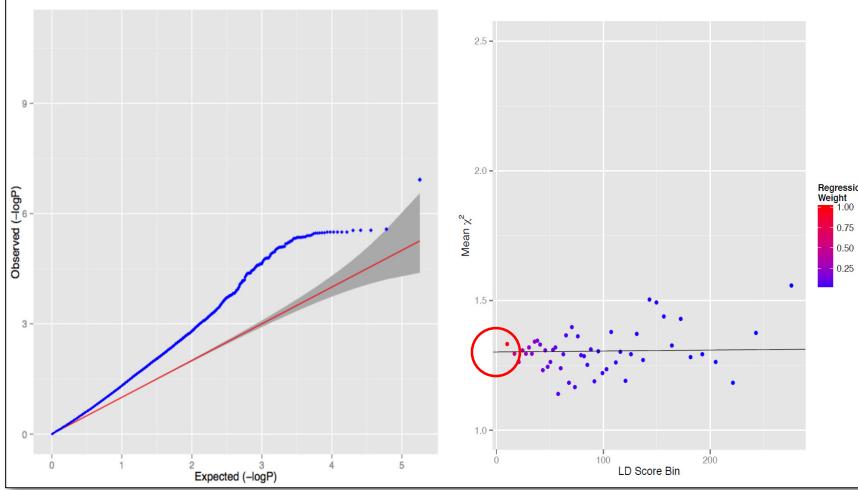
Under pure drift, LD is uncorrelated to magnitude of allele frequency differences between populations



Assuming *i.i.d.* (standardized) effect sizes, more LD yields higher chi-square (on average)  
 More tags  $\rightarrow$  more causal SNPs.  
 More shots  $\rightarrow$  more shots on goal

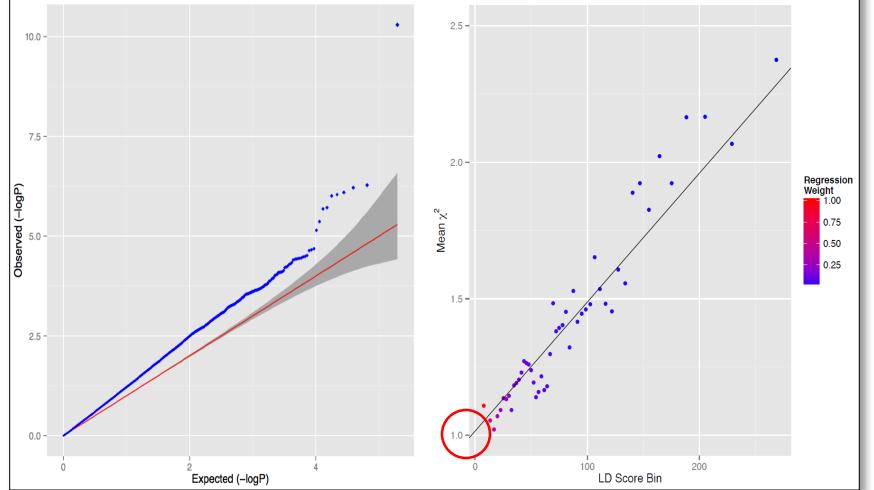
## Simulation under stratification

- $\lambda_{GC} = 1.30$ ; LD Score Regression intercept = 1.32



## Simulation under association

- $\lambda_{GC} = 1.30$ ; LD Score Regression intercept = 1.02



# Linkage disequilibrium: D and D'

- Genetic variants do not segregate independently
- $D = \text{coeff. of linkage disequilibrium between alleles A and B at loci L1 and L2}$ 
  - $D_{AB} = P_{11}P_{00} - P_{10}P_{01} = 0.07$
  - Property of the specific **alleles**. Different alleles at these loci will have diff  $D_{AB}$
- If independent, then  $D_{AB}=0$  ( $P_{11}P_{00}=P_{10}P_{01}$ )
- Linkage disequilibrium measures the degree of departure from Mendel's laws of independent assortment

## How to interpret actual values?

- Relative to  $D_{AB\max}$ , which depends on frequencies of individual alleles at A, B
- $D_{AB\max} = P_0^*P_{*1} - P_1^*P_{*0} = 0.138$
- $D' = D/D_{\max} = 0.51$
- ➔ 51% of max possible disequilibrium

Haplotype	Marginal allele frequency
AB	
0*	0.54
1*	0.46
*0	0.30
*1	0.60

Haplotype	Expected	Observed
00	0.162	0.24**
01	0.324	0.31
10	0.138	0.07**
11	0.276	0.39**

# Linkage disequilibrium: $r^2$

- Define
- $r^2 = \frac{D^2}{P(A=0)P(B=0)P(A=1)P(B=1)} = 0.37$
- This really is the squared Pearson correlation of the two SNPs
- In practice, Pearson correlation is efficiently computed for all SNPs in windows as  $X'X/n$
- This is a fundamental quantity for modeling GWAS z-scores

Haplotype	Marginal allele frequency	
AB		
0*	0.54	
1*	0.46	
*0	0.30	
*1	0.60	
Haplotype	Expected	Observed
00	0.162	0.24
01	0.324	0.31
10	0.138	0.07
11	0.276	0.39

Key property:  $r^2$  correlation for individual SNPs is exactly the  $r^2$  of the GWAS association summary statistics of these SNPs

# LD score regression estimates heritability from summary data

A multivariate model for phenotype variation

**phenotype  
indiv.  $i$**   $y_i = \sum_j X_{ij} \beta_j + \varepsilon_i$  **non-genetic  
for indiv.  $i$**   
**multivar.  
effect on SNP  $j$**

Assuming  $E[X_j] = 0$  and  $V[X_j] = 1$ ,  
**heritability** =  $V[\mathbf{X}\boldsymbol{\beta}] \approx \Sigma \mathbf{X}^2 \boldsymbol{\beta}^2 \approx \Sigma \boldsymbol{\beta}^2$

$$h^2 = \sum_j \beta_j^2$$

Heritability by partitioning  
(restricting on a set  $C$ ):

$$h^2(C) = \sum_{j \in C} \beta_j^2$$

# LD score regression estimates heritability from summary data

A multivariate model

$$y_i = \sum_j X_{ij} \beta_j + \varepsilon_i$$



Summary statistics data

$$\chi_j^2$$
  
$$r_{jk}^2$$

(1) X-square tests statistic for all SNP  $j$   
and (2) LD matrix (or correlation between SNP  $j$  and  $k$ )

Assuming  $E[X_j] = 0$  and  $V[X_j] = 1$ ,  
**heritability** =  $V[\mathbf{X}\boldsymbol{\beta}] \approx \Sigma \mathbf{X}^2 \boldsymbol{\beta}^2 \approx \Sigma \boldsymbol{\beta}^2$

$$h^2 = \sum_j \beta_j^2$$

Heritability by partitioning  
(restricting on a set  $C$ ):

$$h^2(C) = \sum_{j \in C} \beta_j^2$$

# Idea: Reverse-engineer summary data to find multivar. parameters

A univariate effect (GWAS)

$$\begin{aligned}\hat{\beta}_j &= \frac{1}{N} X_j^T (X\beta + \epsilon) \\ &= \sum_k \boxed{\hat{r}_{jk}} \beta_k + \epsilon'_j\end{aligned}$$

LD between  
SNP  $j$  and  $k$

A univariate chi-square (GWAS)

$$\begin{aligned}\chi_j^2 &= N \hat{\beta}_j^2 \\ \text{E}[\chi_j^2] &= N \text{E} \left( \sum_k \hat{r}_{jk} \beta_k + \epsilon'_j \right)^2\end{aligned}$$

# Idea: Reverse-engineer summary data to find multivar. parameters

A univariate effect (GWAS)

$$\begin{aligned}\hat{\beta}_j &= \frac{1}{N} X_j^T (X\beta + \epsilon) \\ &= \sum_k \hat{r}_{jk} \beta_k + \epsilon'_j\end{aligned}$$

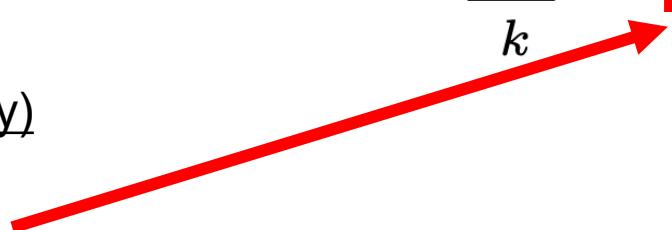
LD between  
SNP  $j$  and  $k$

A univariate chi-square (GWAS)

$$\begin{aligned}\chi_j^2 &= N \hat{\beta}_j^2 \\ \mathbb{E}[\chi_j^2] &= N \mathbb{E} \left( \sum_k \hat{r}_{jk} \beta_k + \epsilon'_j \right)^2 \\ &= N \sum_k \hat{r}_{jk}^2 \mathbb{E}[\beta_k^2] + N \mathbb{E}[\epsilon'_j]^2\end{aligned}$$

Per SNP variance (heritability)

$$\begin{aligned}\text{Var}(\beta_j) &= \sum_{c:j \in \mathcal{C}_c} \tau_c \\ &= \mathbb{E}[\beta_j^2] \text{ (assuming } \mathbb{E}[\beta_j] \approx 0)\end{aligned}$$



# Idea: Reverse-engineer summary data to find multivar. parameters

A univariate effect (GWAS)

$$\begin{aligned}\hat{\beta}_j &= \frac{1}{N} X_j^T (X\beta + \epsilon) \\ &= \sum_k \hat{r}_{jk} \beta_k + \epsilon'_j\end{aligned}$$

**LD between  
SNP  $j$  and  $k$**

A univariate chi-square (GWAS)

$$\begin{aligned}\chi_j^2 &= N \hat{\beta}_j^2 \\ \mathbb{E}[\chi_j^2] &= N \mathbb{E} \left( \sum_k \hat{r}_{jk} \beta_k + \epsilon'_j \right)^2 \\ &= N \sum_k \hat{r}_{jk}^2 \mathbb{E}[\beta_k^2] + N \mathbb{E}[\epsilon'_j]^2\end{aligned}$$

Per SNP variance (heritability)

$$\begin{aligned}\text{Var}(\beta_j) &= \sum_{c:j \in \mathcal{C}_c} \tau_c \\ &= \mathbb{E}[\beta_j^2] \text{ (assuming } \mathbb{E}[\beta_j] \approx 0)\end{aligned}$$

$$\mathbb{E}[\chi_j^2] = N \sum_c \tau_c \sum_{k \in \mathcal{C}_c} \hat{r}_{jk}^2 + \sigma_e^2$$

# Regression of chi-square statistics on LD scores

$$E[\chi_j^2] = N \sum_c \tau_c \sum_{k \in C_c} \hat{r}_{jk}^2 + \sigma_e^2$$

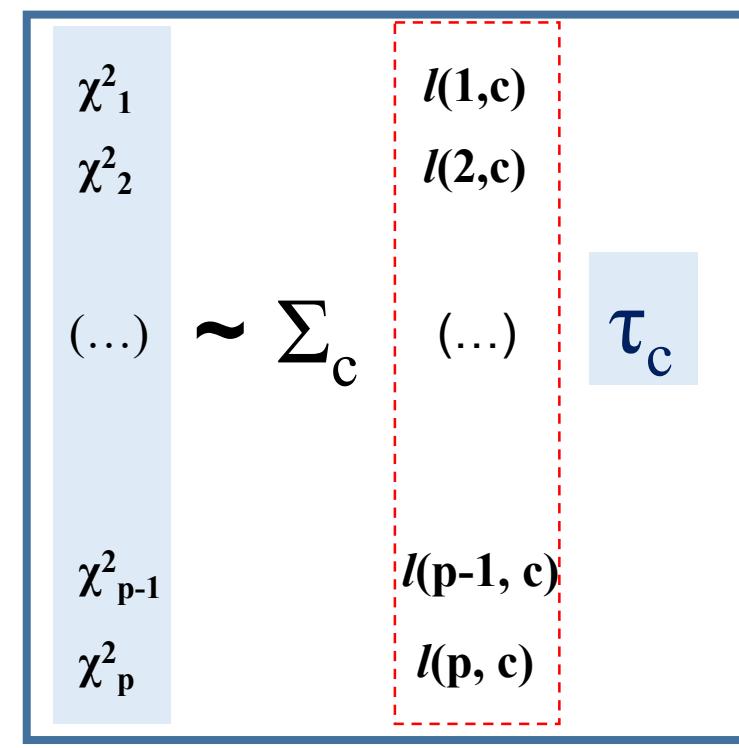
$$E[\chi_j^2] = N \sum_c \tau_c \ell(j, c) + 1$$

$$\ell(j, c) := \sum_{k \in C_c} r_{jk}^2$$

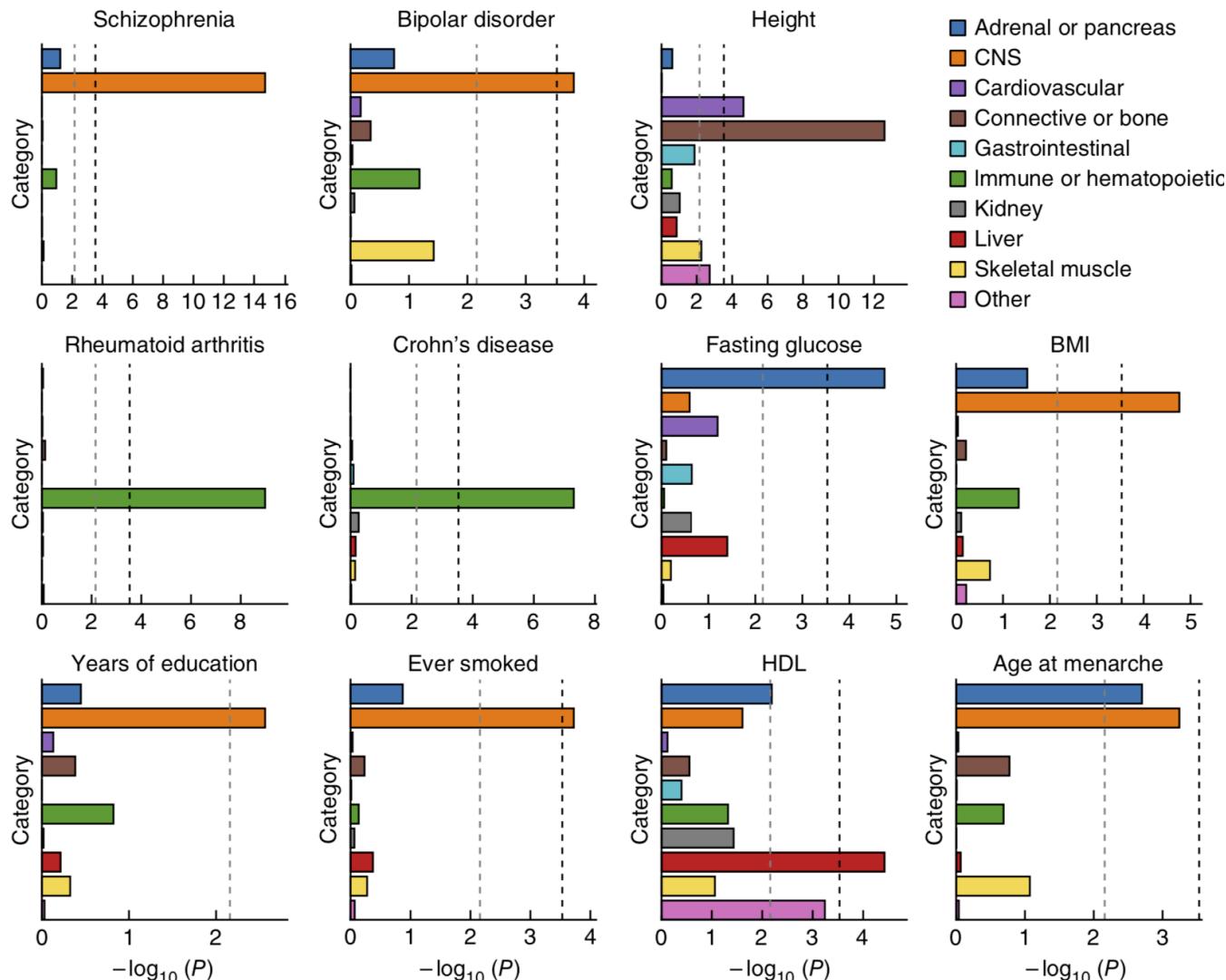
LD-scores between SNP  $j$  and other SNP  $k$  in annotation  $c$

*Intuition: Remove unwanted “double-counting” of annotation enrichment due to LD*

Regression to estimate  $\tau_c$ :



# Stratified LDSC partitions heritability of complex trait GWAS summary



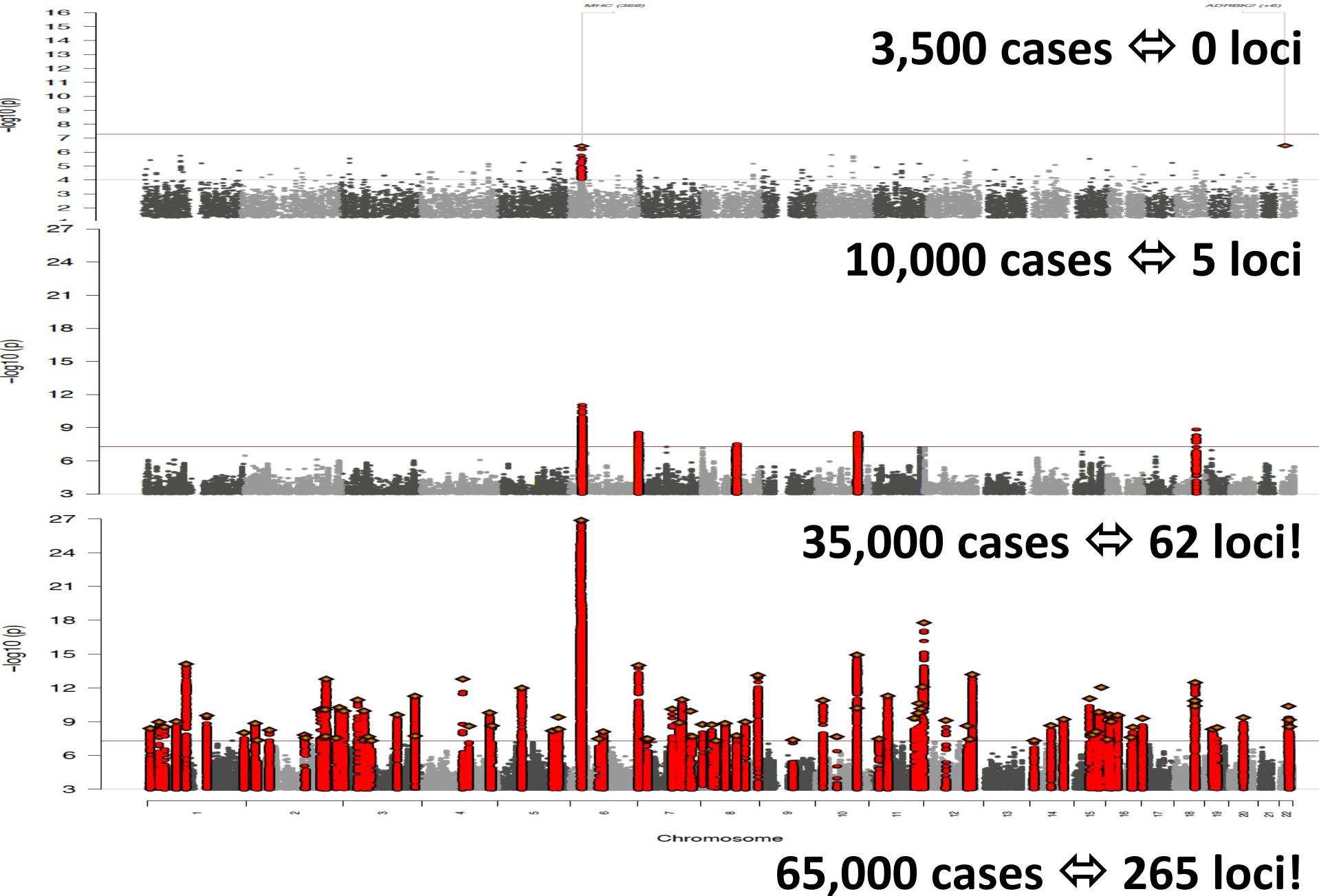
# Today: Deep Learning for Human Genetics and Disease

1. Review: GWAS, fine-mapping, Bayesian variant prioritization
2. Deep Learning for GWAS: calling SNPs, prioritize function
3. eQTLs/Mediation: intermediate molecular phenotypes
4. Linear Mixed Models (LMMs) for GWAS and for eQTL calling
5. Polygenic Risk Scores (PRS): summing over many variants
6. Heritability: definition(s), missing heritability, partitioning
7. LD SCore regression (LDSC) for fast heritability partitioning
8. Polygenic/Omnigenic disease models: core vs. periphery
9. Disease gene networks from GWAS evidence boosting

## **8. Polygenic → Omnigenic models of disease**

Recognizing “core” vs. “periphery” pathways

# Schizophrenia GWAS: Number of significant loci

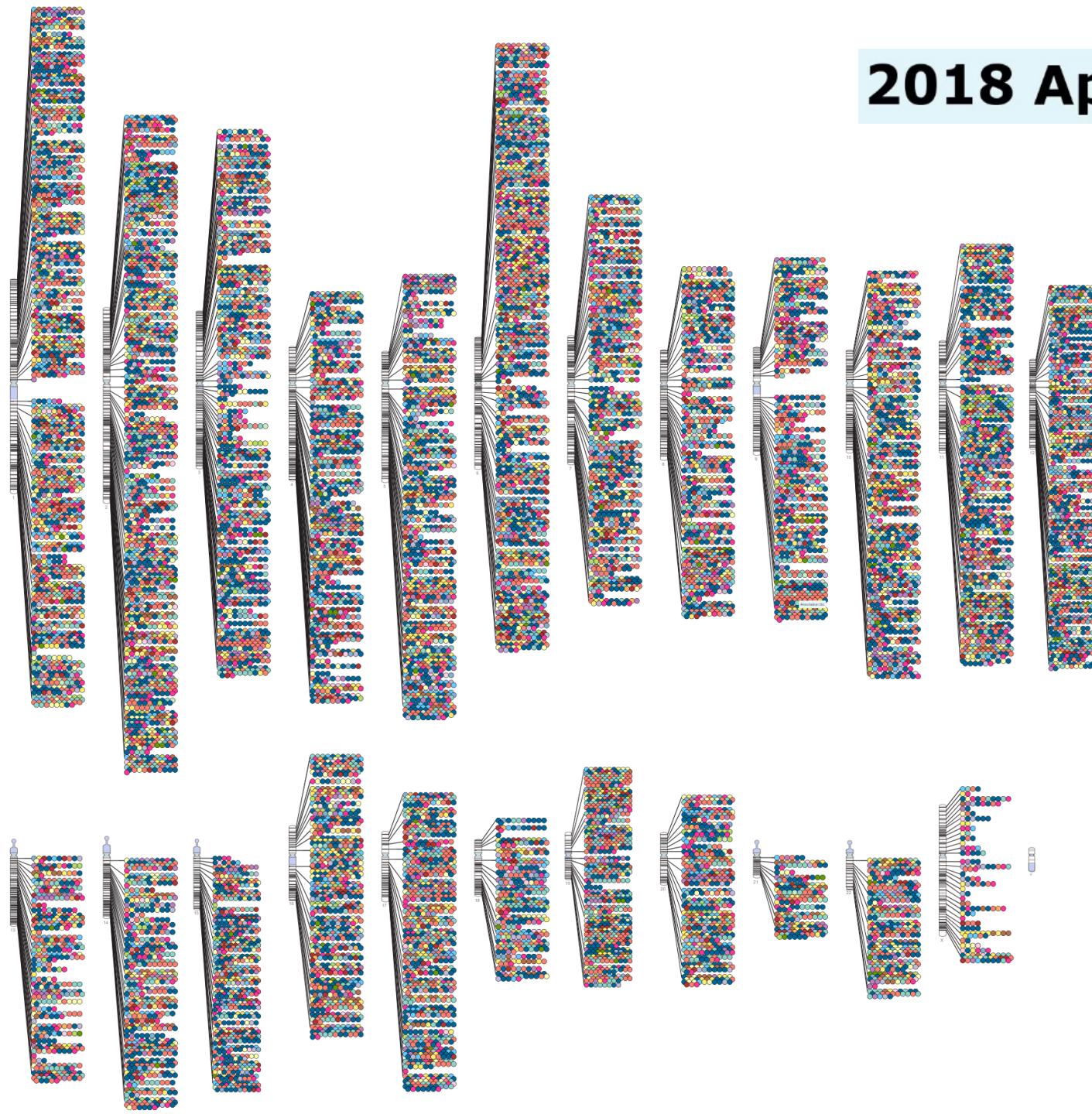


Associations: 69,885

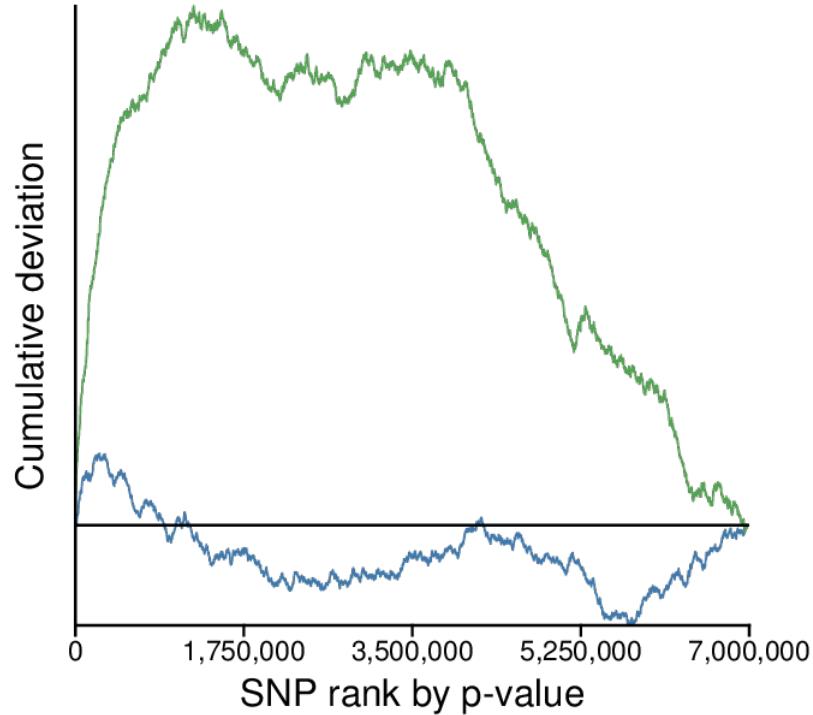
2018 Apr

Studies: 5,152

Papers: 3,378

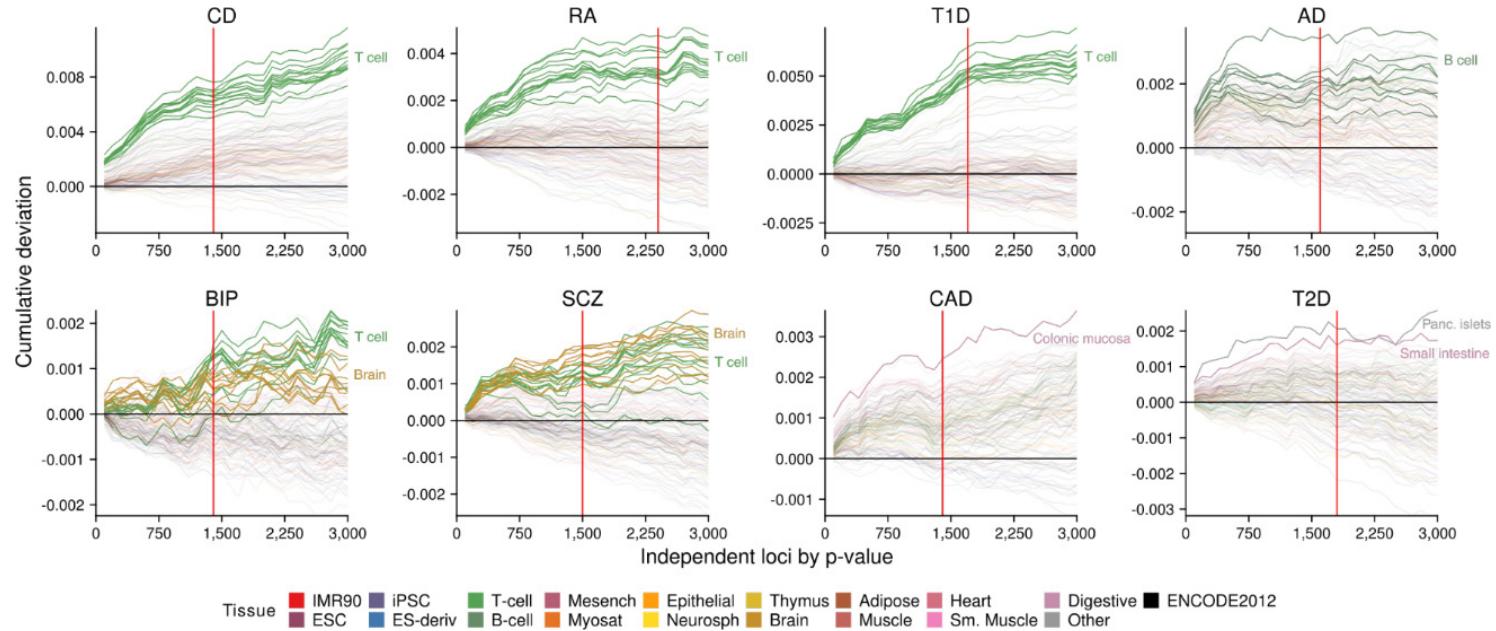


# How far down the SNP list does enrichment go?



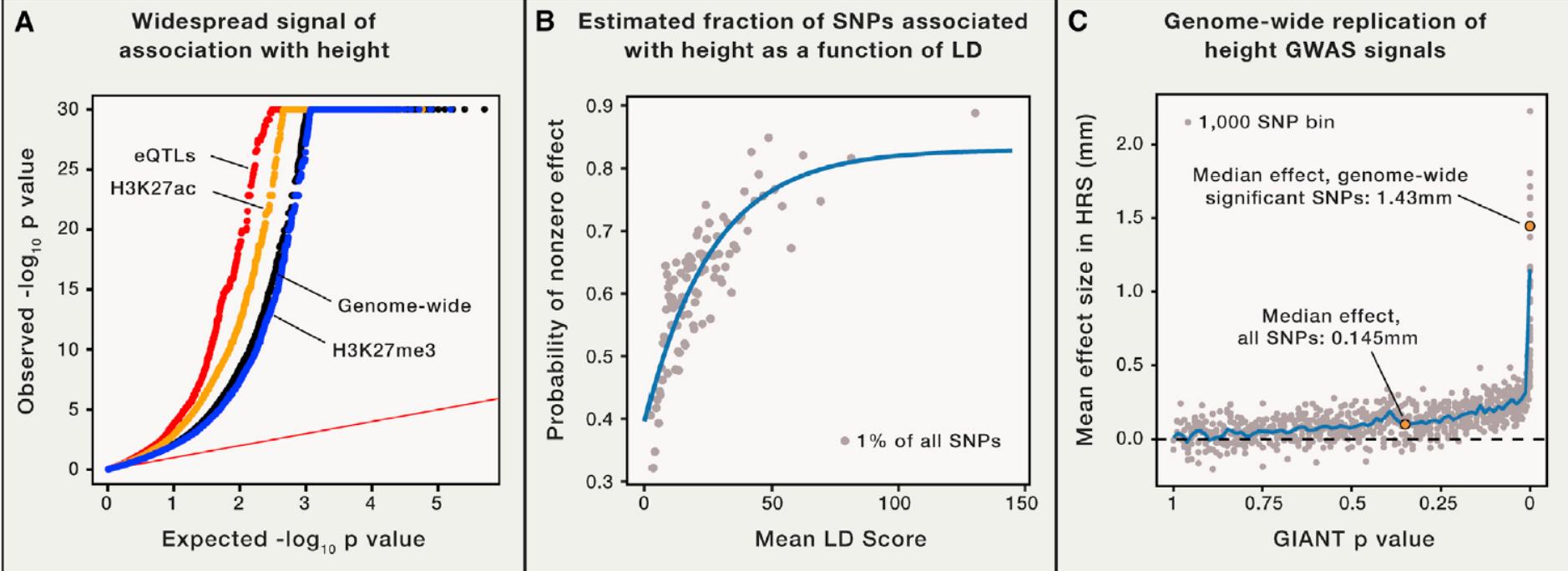
- Use functional enrichment to gain insight into genetic architecture (Sarkar 2016)
- Idea: as we consider more SNPs beyond genome-wide significance, relevant regulatory regions will be disrupted more often than irrelevant regions

# Long tails of enrichment for 8 diseases



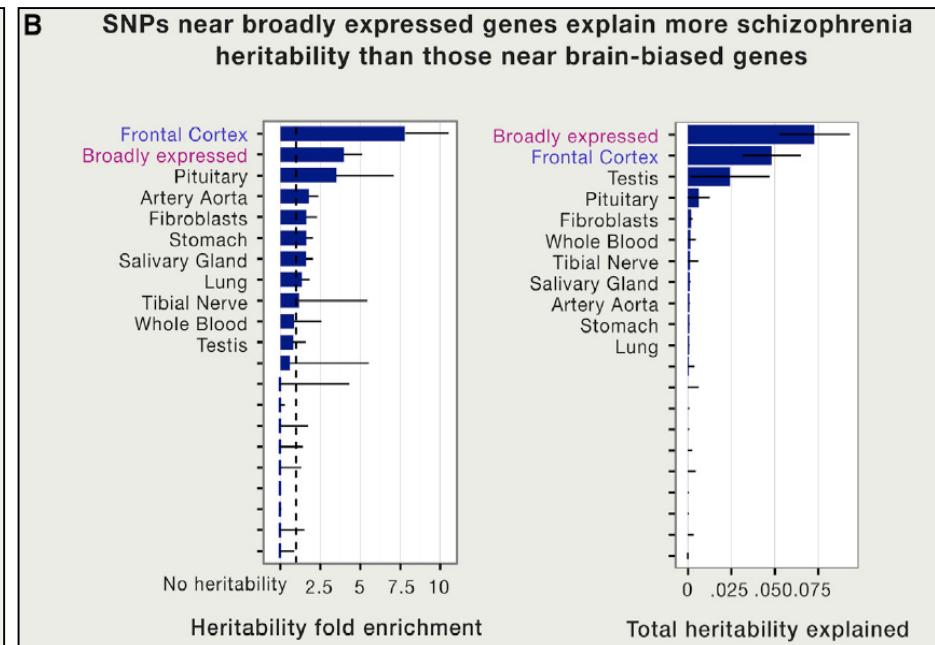
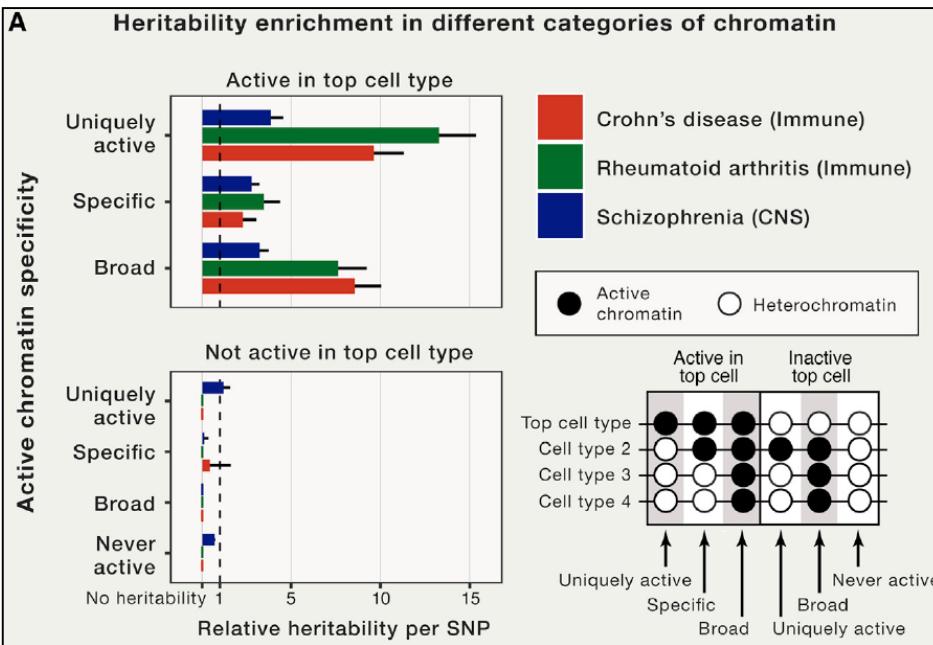
- Use functional enrichment to gain insight into genetic architecture (Sarkar 2016)
- Idea: as we consider more SNPs beyond genome-wide significance, relevant regulatory regions will be disrupted more often than irrelevant regions

# Omnigenic model of heritability



- (A) Genome-wide inflation of small p values from the GWAS for height, with particular enrichment among expression quantitative trait loci and single-nucleotide polymorphisms (SNPs) in active chromatin (H3K27ac).
- (B) Estimated fraction of SNPs associated with non-zero effects on height (Stephens, 2017) as a function of linkage disequilibrium score (i.e., the effective number of SNPs tagged by each SNP; Bulik-Sullivan et al., 2015b). Each dot represents a bin of 1% of all SNPs, sorted by LD score. Overall, we estimate that 62% of all SNPs are associated with a non-zero effect on height. The best-fit line estimates that 3.8% of SNPs have causal effects.
- (C) Estimated mean effect size for SNPs, sorted by GIANT p value with the direction (sign) of effect ascertained by GIANT. Replication effect sizes were estimated using data from the Health and Retirement Study (HRS). The points show averages of 1,000 consecutive SNPs in the p-value-sorted list. The effect size on the median SNP in the genome is about 10% of that for genome-wide significant hits.

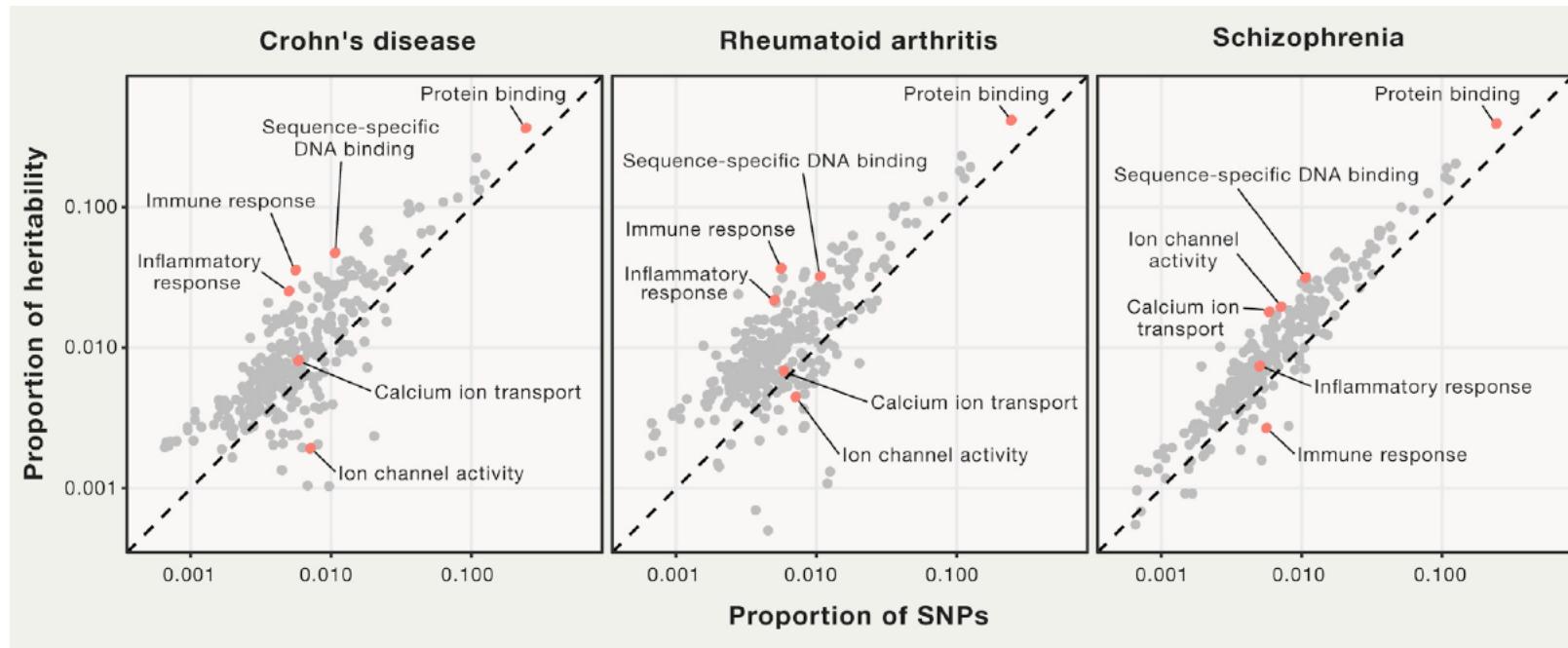
# More heritability in broad classes



- Contributions to heritability (relative to random SNPs) as a function of chromatin context. There is enrichment for signal among SNPs that are in chromatin active in the relevant tissue, regardless of the overall tissue breadth of activity

- Genes with brain-specific expression show the strongest enrichment of schizophrenia signal (left), but broadly expressed genes contribute more to total heritability due to their greater number (right)

# Most GO categories are enriched



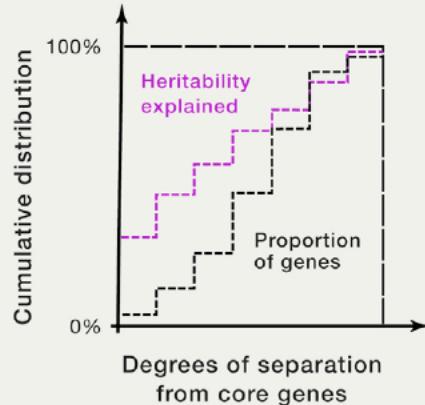
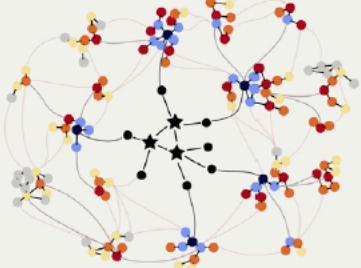
- Gene Ontology Enrichments for Three Diseases, with Categories of Particular Interest Labeled. The x axis indicates the fraction of SNPs in each category; the y axis shows the fraction of heritability assigned to each category as a fraction of the heritability assigned to all SNPs. Note that the diagonal indicates the genome-wide average across all SNPs; most GO categories lie above the line due to the general enrichment of signal in and around genes. Analysis by stratified LD score regression

# Core genes vs. periphery

**A**

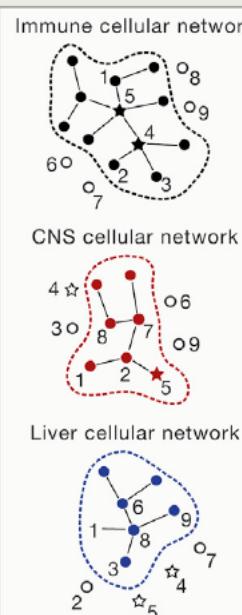
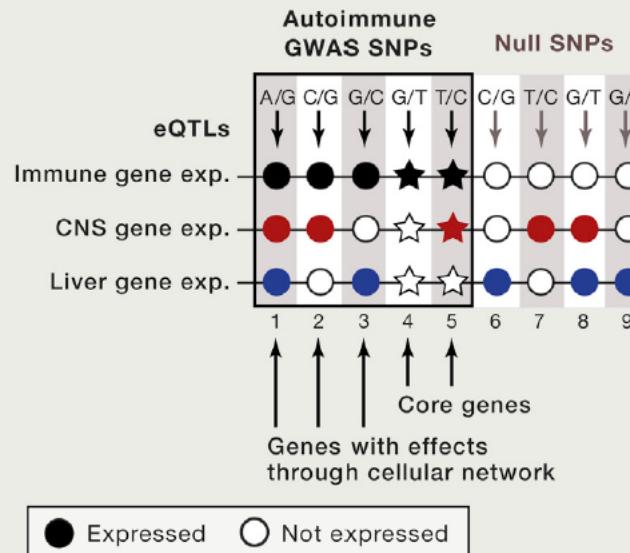
**Model: Most genes affect disease risk through highly connected cellular networks**

Degrees of separation from core genes  
Low [Black] 1 2 3 4 5 6 >7 High



**B**

**Autoimmune GWAS hits affect shared and tissue-specific regulation of immune cells**



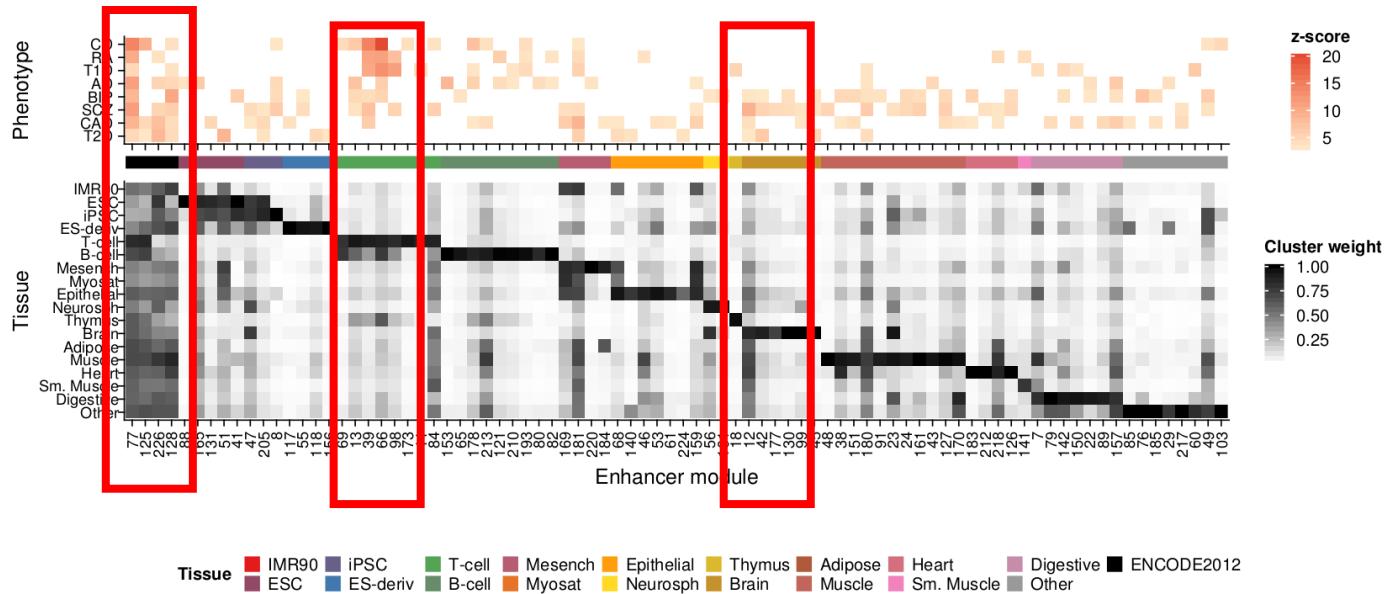
- **Omnigenic Model of Complex Traits**
- (A) For any given disease phenotype, a limited number of genes have direct effects on disease risk. However, by the small world property of networks, most expressed genes are only a few steps from the nearest core gene and thus may have non-zero effects on disease. Since core genes only constitute a tiny fraction of all genes, most heritability comes from genes with indirect effects.
- (B) Diseases are generally associated with dysfunction of specific tissues; genetic variants are only relevant if they perturb gene expression (and hence network state) in those tissues. For traits that are mediated through multiple cell types or tissues, the overall effect size of any given SNP would be a weighted average of its effects in each cell type.

# Today: Deep Learning for Human Genetics and Disease

1. Review: GWAS, fine-mapping, Bayesian variant prioritization
2. Deep Learning for GWAS: calling SNPs, prioritize function
3. eQTLs/Mediation: intermediate molecular phenotypes
4. Linear Mixed Models (LMMs) for GWAS and for eQTL calling
5. Polygenic Risk Scores (PRS): summing over many variants
6. Heritability: definition(s), missing heritability, partitioning
7. LD SCore regression (LDSC) for fast heritability partitioning
8. Polygenic/Omnigenic disease models: core vs. periphery
9. Disease gene networks from GWAS evidence boosting

## **9. GWAS networks for evidence boosting**

# Enhancer modules: constitutive, cell type specific



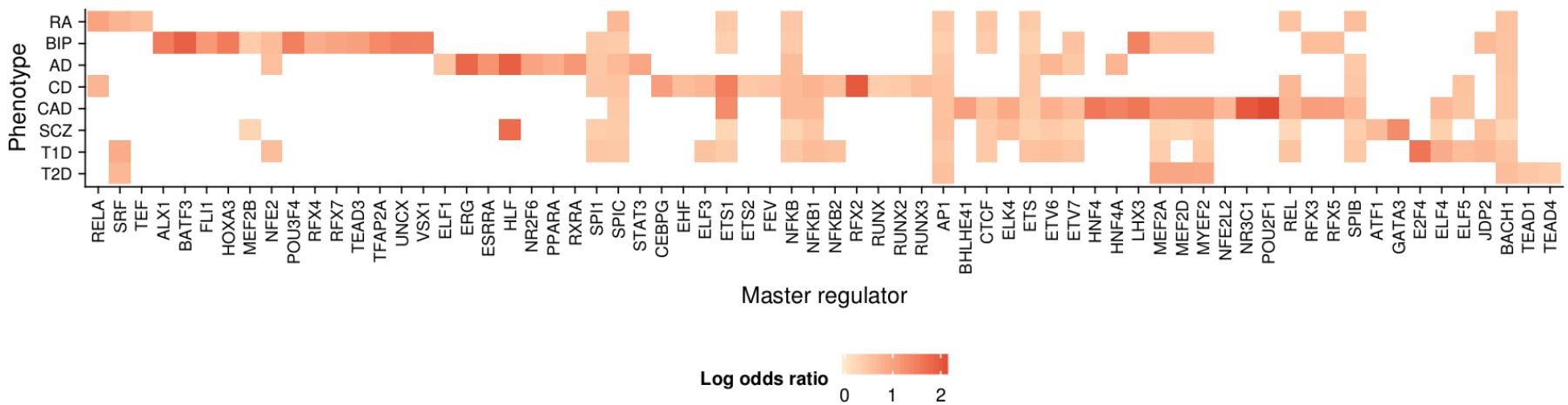
- Challenge: annotations learned one cell type at a time can't account for sharing of elements across cell types
- Use k-means clustering to define modules of enhancer activity
- Functional enrichments highlight importance of both constitutive and lineage-specific enhancers

# From enhancers to genes to pathways

Trait	Known pathways	Total genes	Total pathways
AD	Cyclic GMP signaling, immune response	220	216
BIP	Glucocorticoid signaling	217	230
CAD	Cholesterol/triglyceride metabolism, IgA	248	215
CD	CD8 T cell proliferation, IgE, IL4	224	359
RA	NFKB, actin nucleation	196	146
SCZ	Dendritic spine development	271	183
T1D	MHC I/II, JAK-STAT, IFNG	266	245
T2D	Pancreatic beta cell apoptosis	281	177

- Link enhancers to their downstream target genes
- Target genes enriched in known disease pathways, but through previously unknown mechanisms
- Reveals broad similarities at pathway level between classes of diseases (e.g. signaling in autoimmune traits), but also specific pathways important to each disease
- Potentially implicate novel genes in enriched pathways

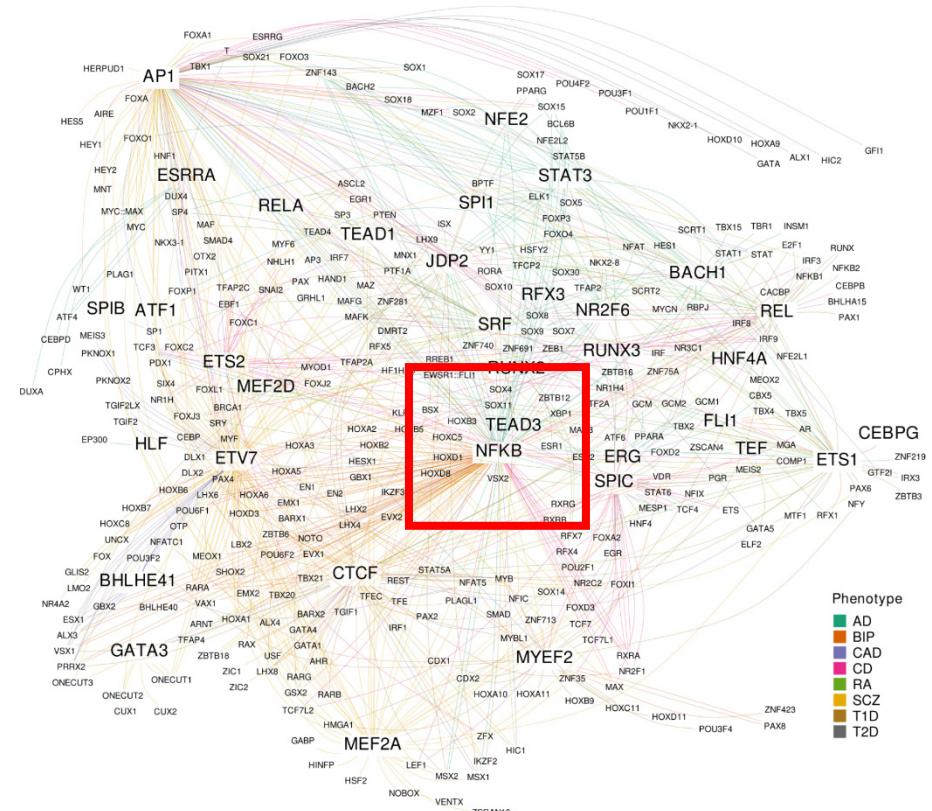
# From genes/pathways to upstream regulators



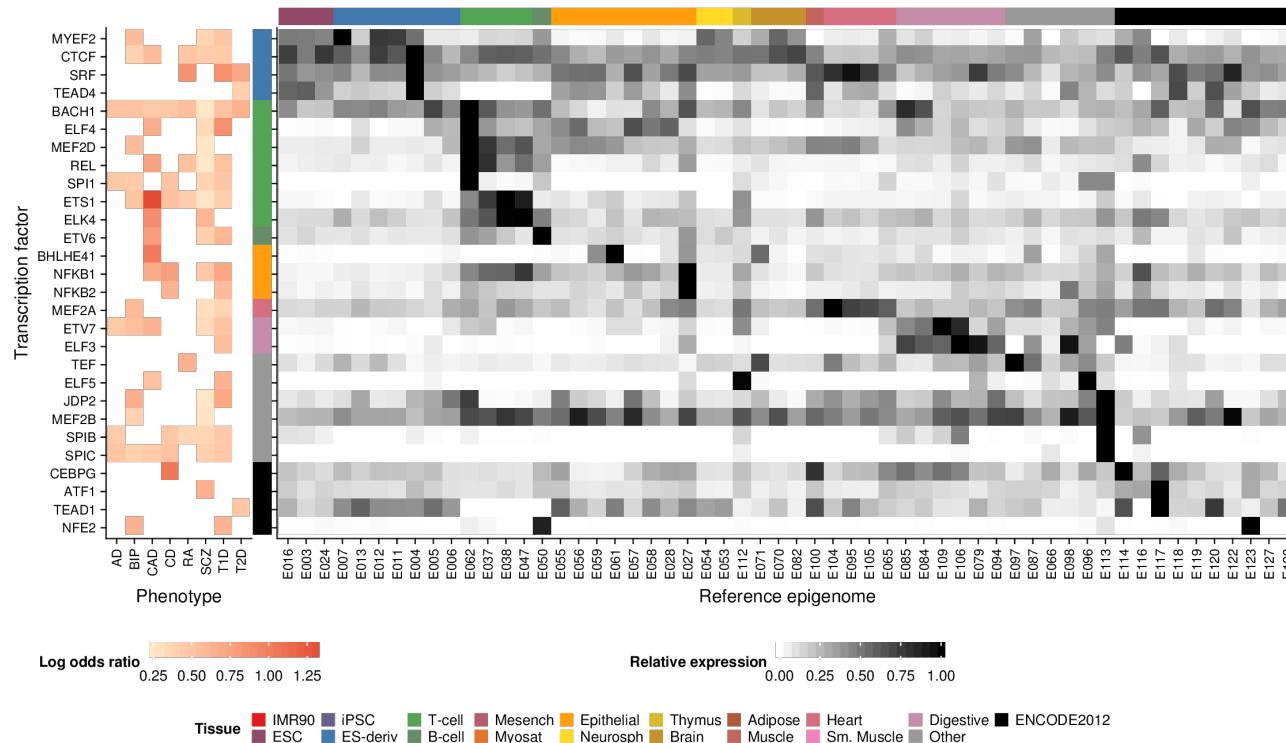
- Challenge: heritability-based methods can't identify specific enhancer regions
- Our method can implicate specific enhancers, so we can dissect their mechanism
- Predict the upstream regulator using sequence-based enrichment (Kheradpour 2013) without considering GWAS
- Find master regulators recurrently disrupted by sub-threshold SNPs
- Many disease-specific regulators, but interesting shared regulators

# Regulator → gene networks across diseases

- GWAS associated SNP often does not directly disrupt the predicted master regulator
- Instead, falls in a different motif instance for a putative co-factor
- Explains how master regulators can be shared across very different phenotypes (NFKB in schizophrenia, T1D)

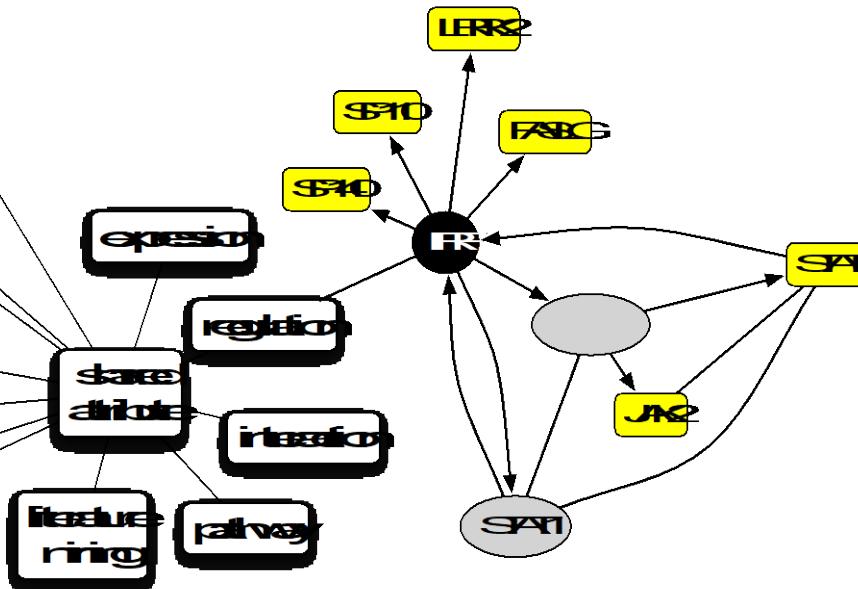
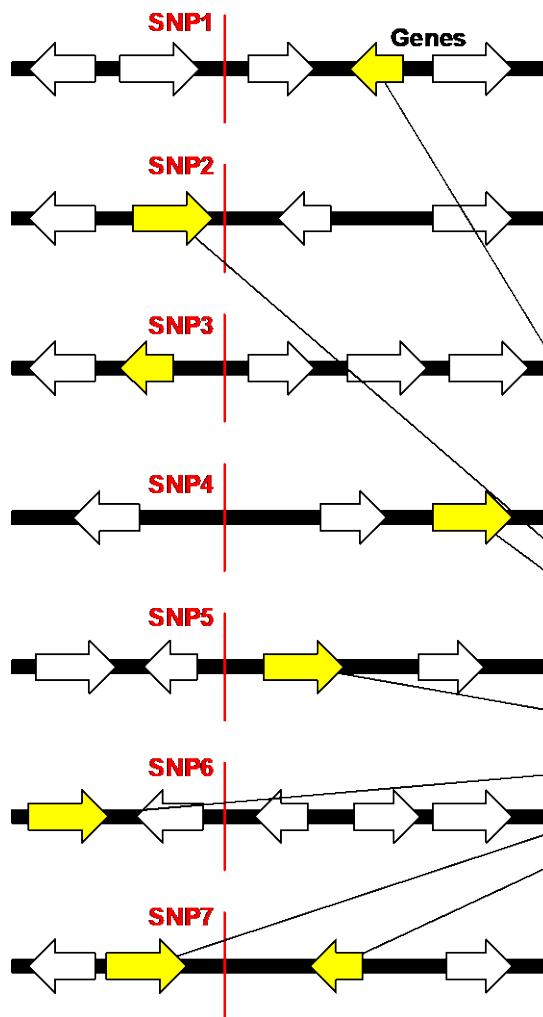


# Upstream regulators add cell-type-specificity



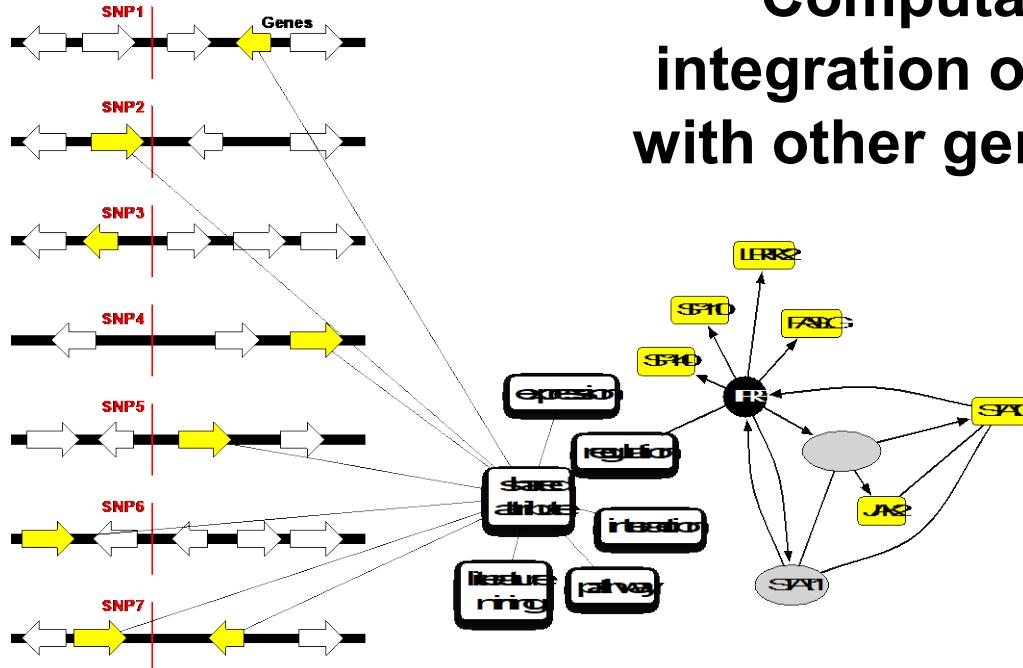
- Many predicted master regulators found in predicted constitutive enhancers rather than cell type-specific regulators
- Although enhancers might be constitutively marked, expression of the upstream regulator is cell type-specific
- Additional insight into transcriptional regulation needed to interpret non-coding disease associations

# Hypothesis: Many associated genes implicate limited number of pathways



**Proof: Statistically significant excess connectivity of genes in GWAS regions**

# Computational tools enable the integration of ‘human genetic screens’ with other genome-scale screening data

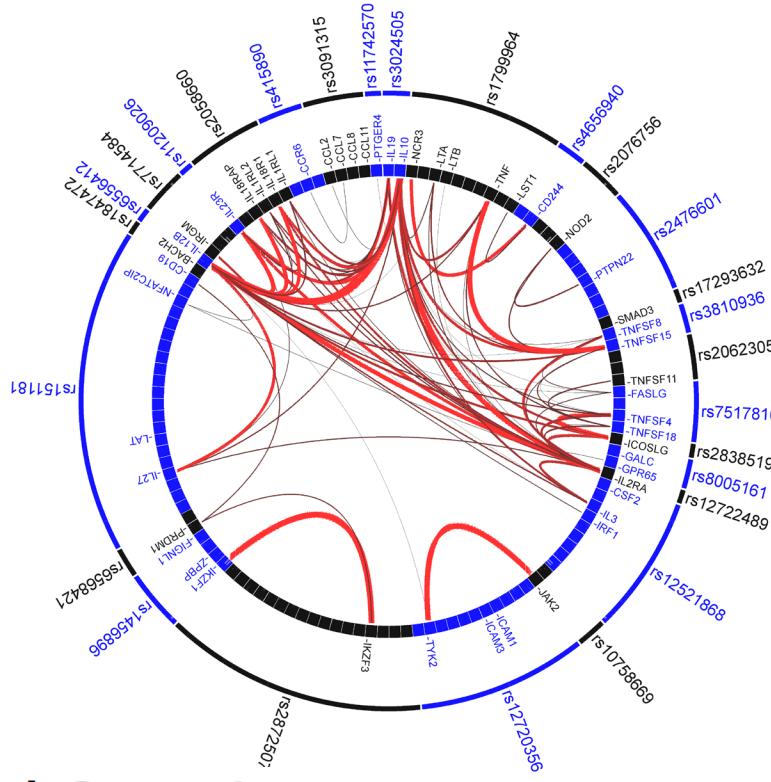


# Proteins Encoded in Genomic Regions Associated with Immune-Mediated Disease Physically Interact and Suggest Underlying Biology DAPPLE

# DAPPLE

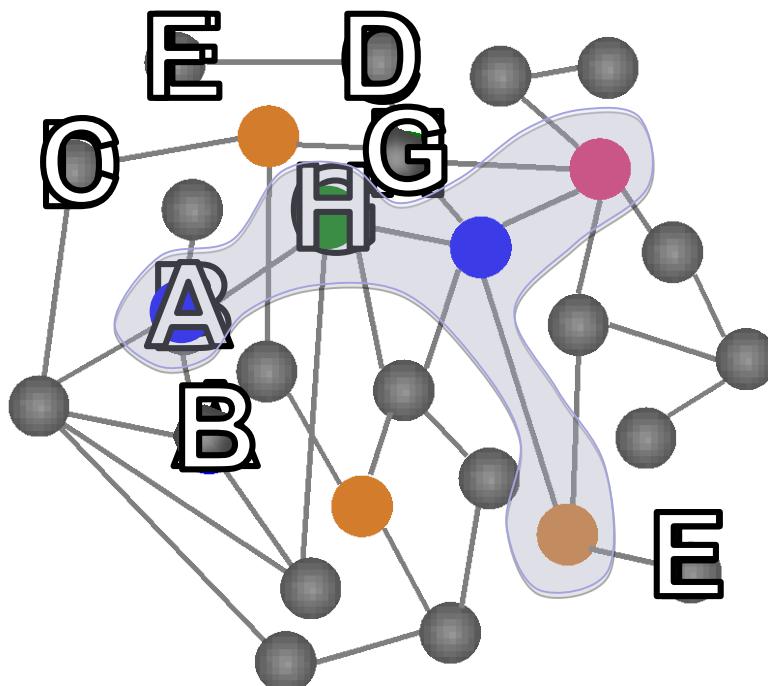
Elizabeth J. Rossin<sup>1,2,3,4,5</sup>, Kasper Lage<sup>2,3,6,7</sup>, Soumya Raychaudhuri<sup>1,2,8</sup>, Ramnik J. Xavier<sup>1,2,3</sup>, Diana Tatar<sup>6</sup>, Yair Benita<sup>1</sup>, International Inflammatory Bowel Disease Genetics Consortium<sup>1</sup>, Chris Cotsapas<sup>1,2,9</sup>, Mark J. Daly<sup>1,2,3,4,5,\*</sup>

# Common Inherited Variation in Mitochondrial Genes Is Not Enriched for Associations with Type 2 Diabetes or Related Glycemic Traits MAGENTA

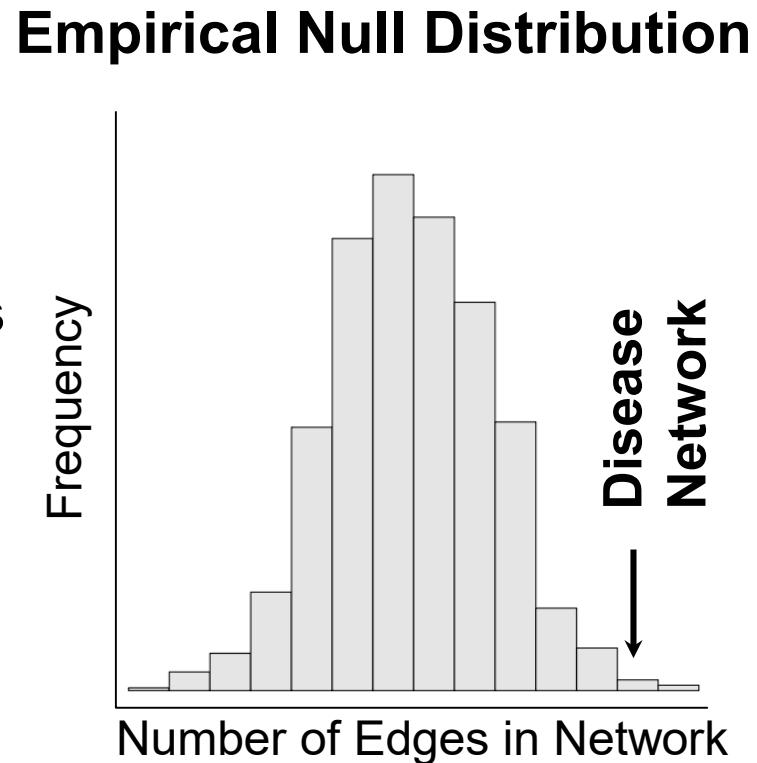


## **GRAIL** plot from Franke et al 2010

# Evaluating Significance



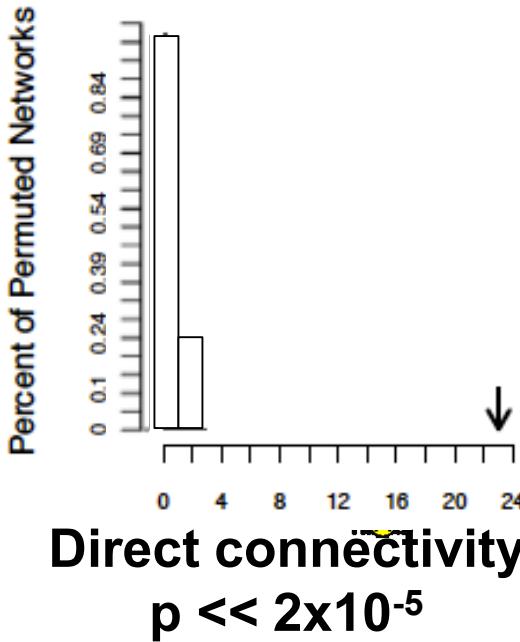
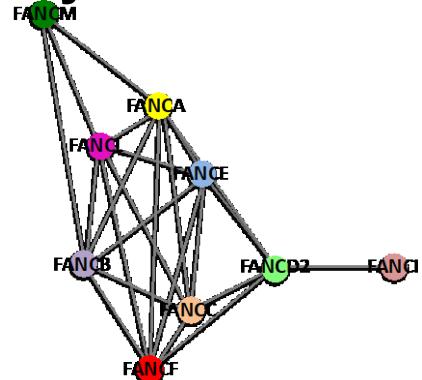
Repeat full  
permutation  
50,000 times



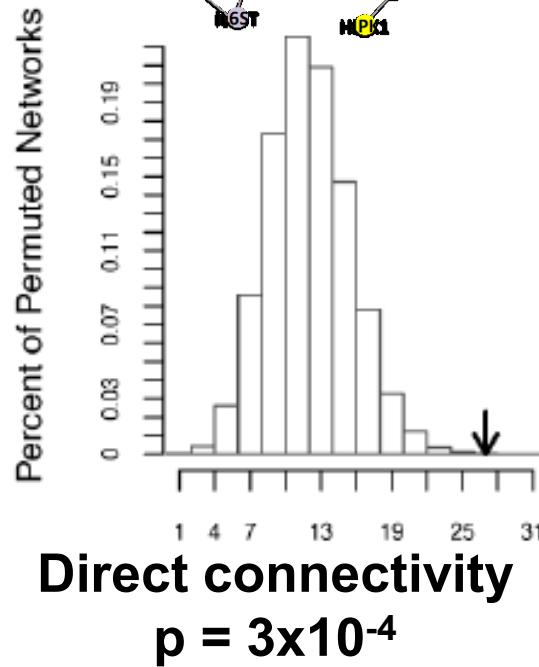
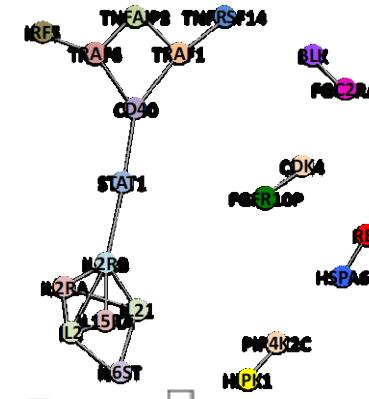
...keep moving labels  
until the network has  
been fully permuted

# PPI Networks identify specific genes and pathways

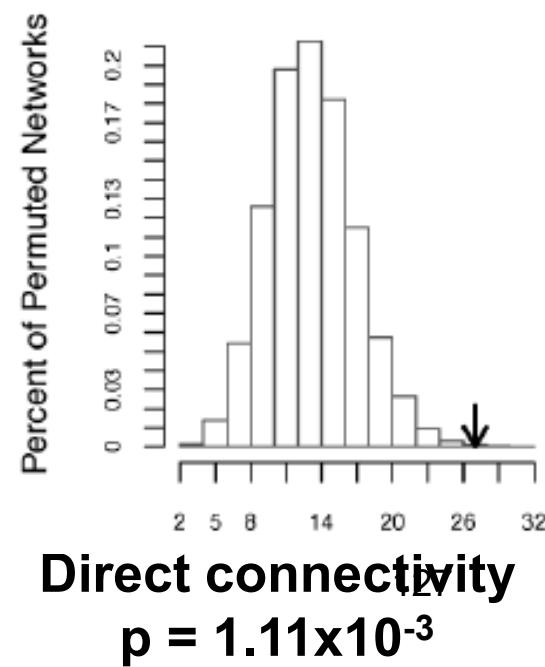
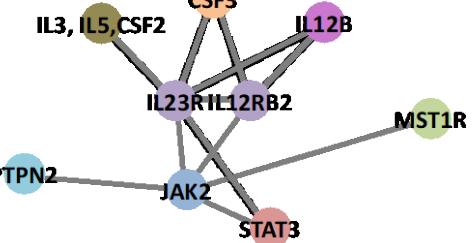
Fanconi anemia  
9 synthetic loci



Rheumatoid arthritis  
27 loci

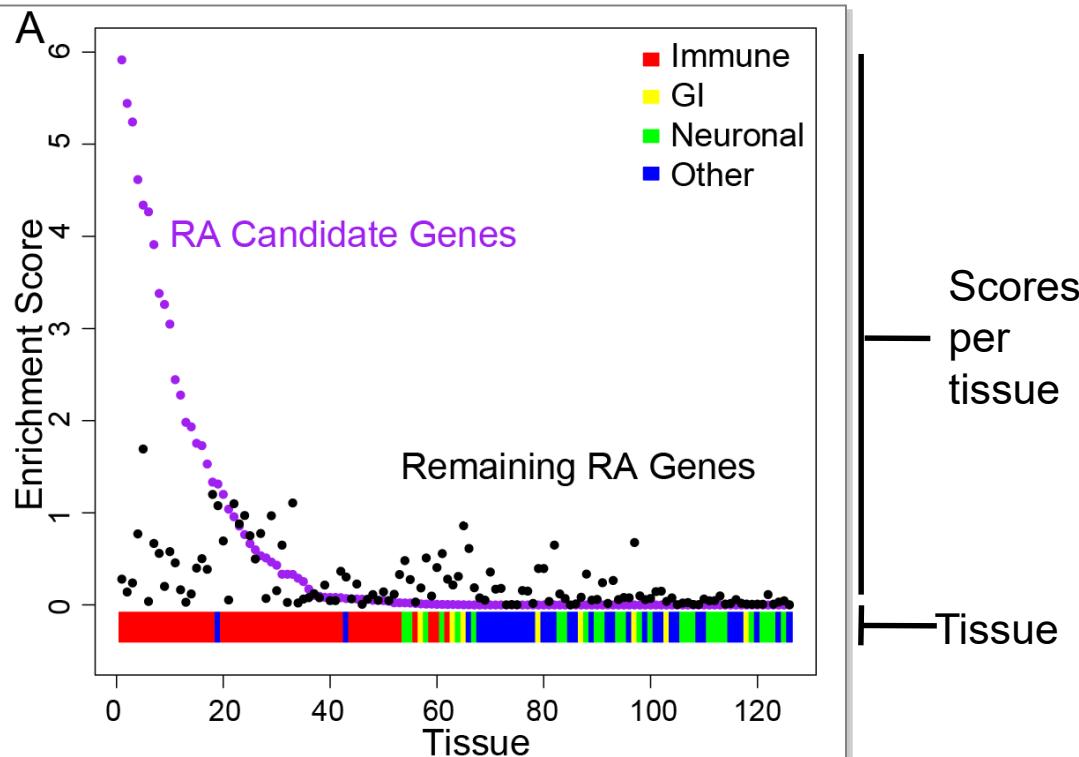


Crohn's disease  
25 loci



# Validation of PPI networks

Further experimental support that the non-random networks are truly implicating the underlying genes



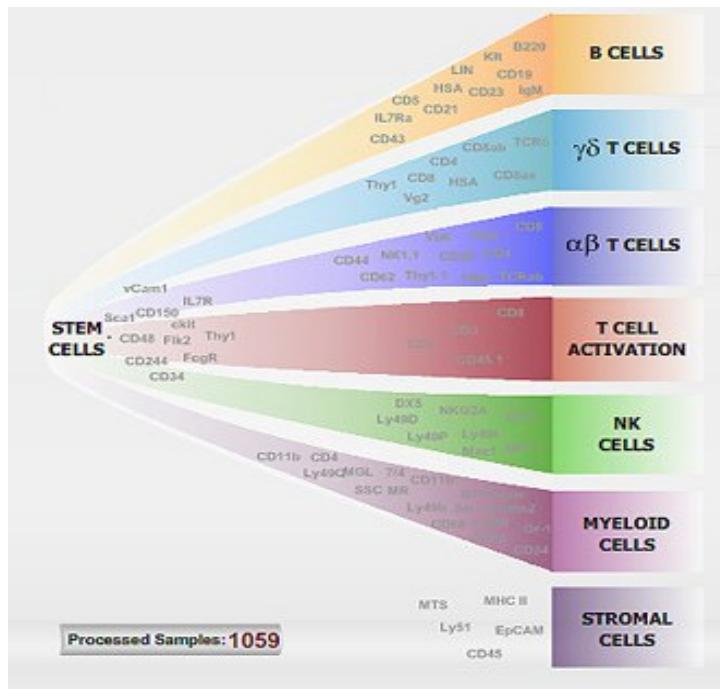
**Network genes are co-expressed**

**Connected proteins are enriched for newly confirmed associated genes ( $p=6.5 \times 10^{-4}$ )**

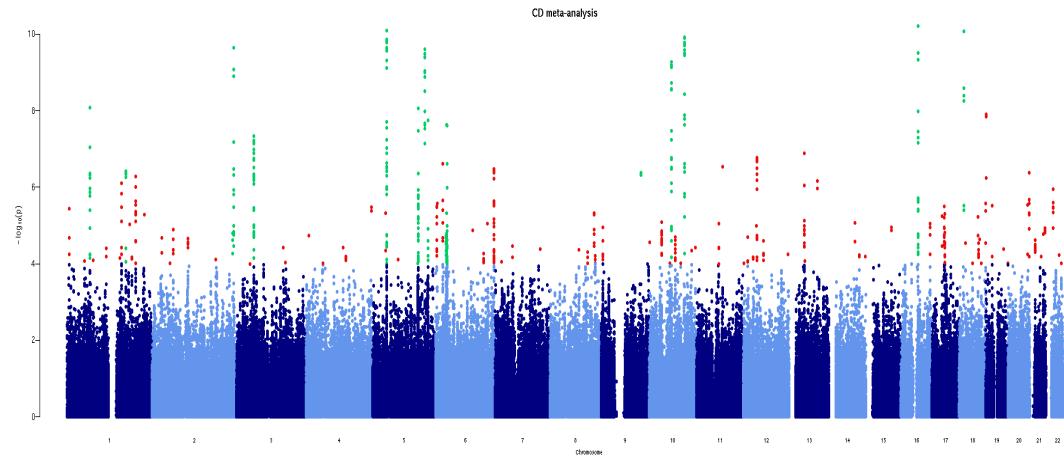
# Integrating Autoimmune Risk Loci with Gene-Expression Data Identifies Specific Pathogenic Immune Cell Subsets

Xinli Hu,<sup>1,2,3,4</sup> Hyun Kim,<sup>1,2</sup> Eli Stahl,<sup>1,2,3</sup> Robert Plenge,<sup>1,2,3</sup> Mark Daly,<sup>3,5</sup> and Soumya Raychaudhuri<sup>1,2,3,6,\*</sup>

The American Journal of Human Genetics 89, 481–482, October 7, 2011

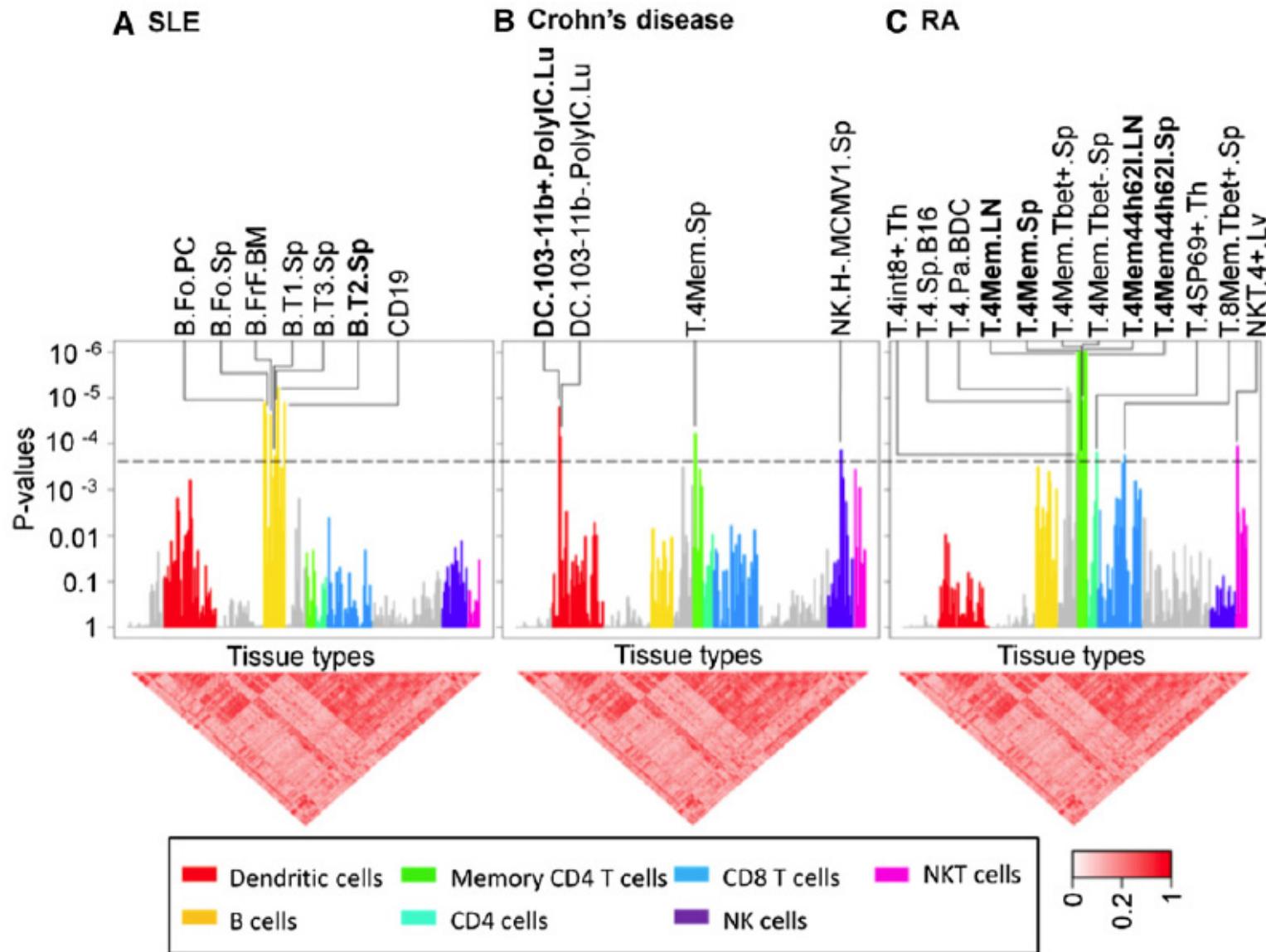


ImmGen data set:  
223 murine immune cell subsets  
Expression measured on 15,149  
human homologs

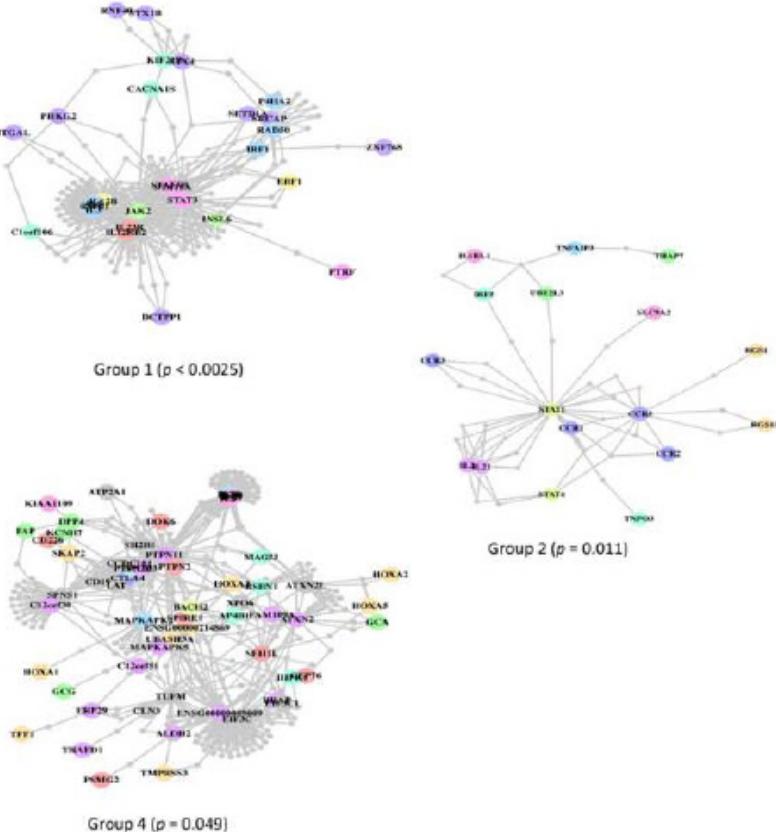
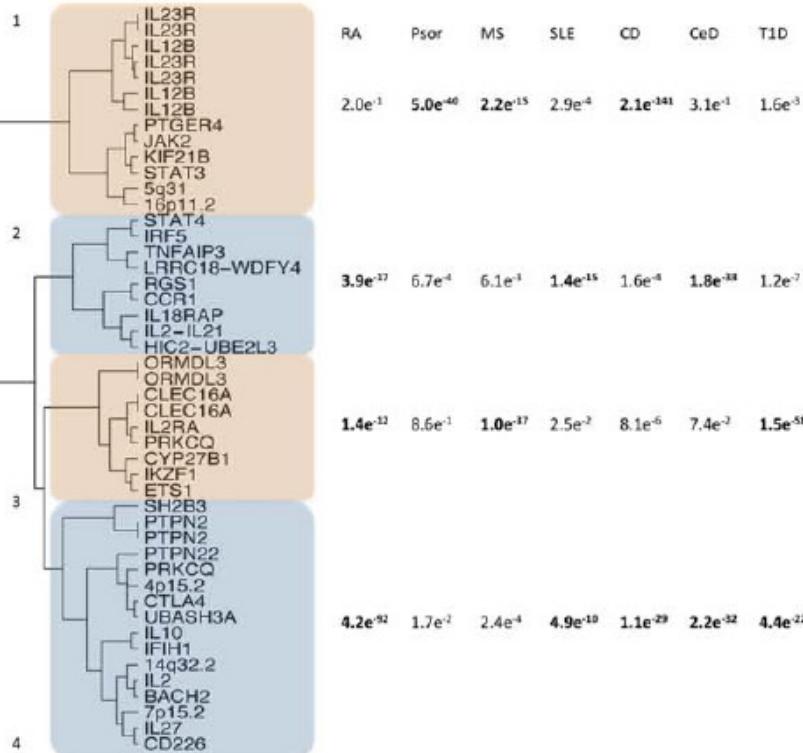


**Are human GWAS hits harboring loci significantly co-expressed in specific immune cell subsets?**

# **GWAS hits significantly co-expressed in specific immune cell subsets**



# Other opportunities: Cross-disease information



Genes coordinately associated to multiple disease are tightly functionally linked

Cotsapas et al, August 2011 *PLoS Genetics*