# Functional interpretation of Genome-wide association studies

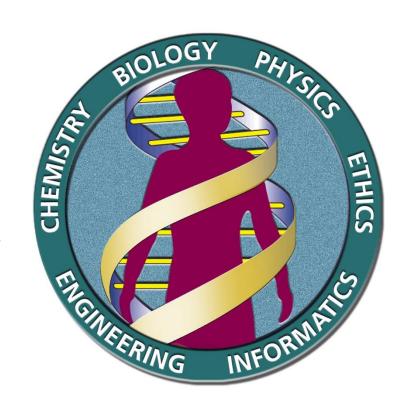
(GWAS2Function)

Kushal K. Dey

## Background to Human Genetic studies

### The Human Genome Project

- Human Genome Project (launched 1990, completed 2003)
- Generate the first sequence of the human genome
  - Reference genome: all base pairs in human genome
  - Map all genes observed ~22K protein-coding genes



Got the ball rolling in terms of genomic sequencing

HapMap Project: Cataloguing variations in the sequences of human DNA (2002-2010) (1,000 individuals)

DNA sequence of any two individuals is 99.5% similar, however the 0.5% difference drives differences in physiological traits and disease risk.

HapMap catalogued variation across ~1,000 individuals.

Sites in the DNA sequence where individuals differ at a single DNA base are called **single nucleotide polymorphisms (SNPs)**.

SNPs were identified at specific chromosomal		chrom.	physical position (bp)
positions (what nomenclature to use?)	rs10910034 rs1713712	1 1	2165898 2166021

## Genome wide Association studies

### Collecting genotype and phenotype data from many many individuals (order of 100, 000 individuals)

Large-scale retrospective studies (~100K-1M individuals)









### Collecting genotype and phenotype data from many many individuals (order of 100, 000 individuals)

Large-scale retrospective studies (~100K-1M individuals)









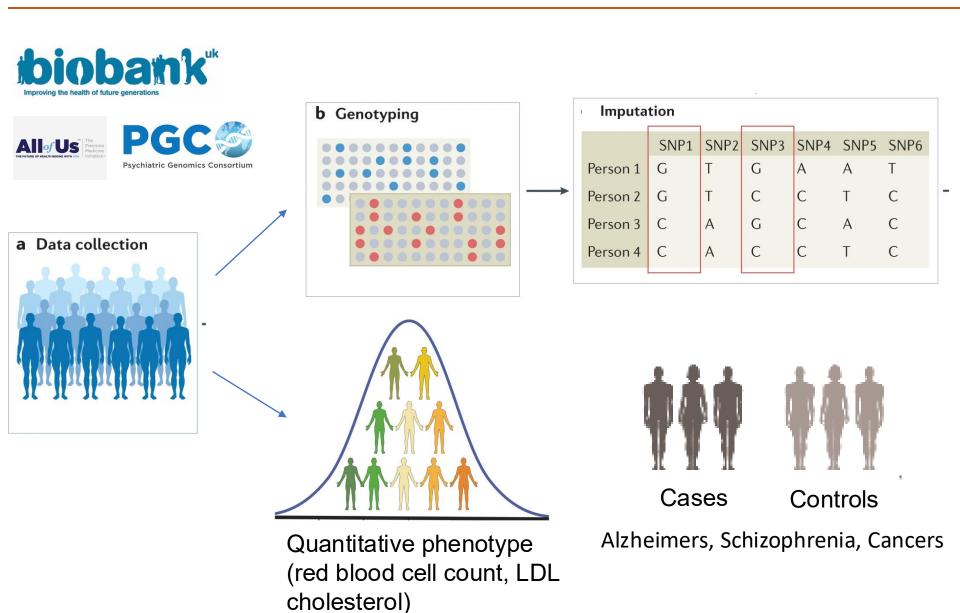
Disease-related prospective studies (~10K-100K)



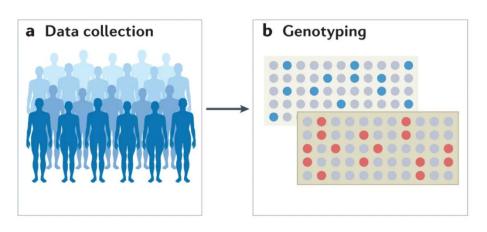




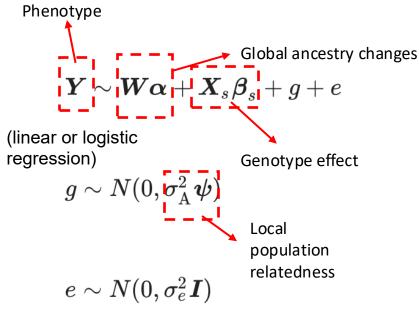
### Background: Genome-wide Association Studies



#### Mathematical model for Genome Wide association studies



Sequencing strategies: SNP array + imputation Whole exome sequencing and Whole Genome sequencing



- Y vector of phenotype values for all N individuals (for example: height or 1/0 for Type 2 diabetes status)
- X<sub>s</sub> vector of genotype values for all N individuals at SNP s (0/1/2 for unscaled: ore standardized)
- **W** matrix of covariates (age, sex, ancestry PCs)
- g represents polygenic effect of other SNPs
- e random effect of residual errors
- $\psi$  kinship or genetic relatedness matrix

### Calculating statistics from Genome Wide association studies

#### Estimates of the effect size

$$\hat{eta}_{ ext{snp}} = rac{\mathbf{x}_{ ext{snp}}^{\mathbf{T}} \mathbf{V}^{-1} \mathbf{y}}{\mathbf{x}_{ ext{snp}}^{\mathbf{T}} \mathbf{V}^{-1} \mathbf{x}_{ ext{snp}}} ext{with } ext{var} \left( \hat{eta}_{ ext{snp}} 
ight) = rac{1}{\mathbf{x}_{ ext{snp}}^{\mathbf{T}} \mathbf{V}^{-1} \mathbf{x}_{ ext{snp}}}$$

V =  $\sigma_g^2 \, \psi + \sigma_e^2 I$ 

Overall phenotypic variance-covariance matrix = genetic + error

Obtain z scores and p-values of the effect based on this.

**GCTA** 

a tool for Genome-wide Complex Trait Analysis

GCTA

SMR GSMR

OSCA

CTG forum

Yang Lab



**Download** 

FAQ

**Basic Options** 

**GREML** 

**GWAS Analysis** 

MLMA

fastGWA

#### fastGWA

#### fastGWA: A fast MLM-based Genome-Wide Association tool

fastGWA is an ultra-efficient tool for mixed linear model (MLM)-based GWAS analysis of biobank-scale data such as the UK Biobank (see Jiang et al. *Nature Genetics* 2019 for details of the method). Credits: Longda Jiang (method, simulation and analysis), Zhili Zheng (method, software and analysis) and Jian Yang (method and overseeing).

We have applied fastGWA to 2,173 traits on 456,422 array-genotyped and imputed individuals and 2,048 traits on 49,960 whole-examples of (WFS) individuals in the LIK Richark. All the summary statistics are available.

### Calculating statistics from Genome Wide association studies

#### Estimates of the effect size

$$\hat{eta}_{ ext{snp}} = rac{\mathbf{x}_{ ext{snp}}^{\mathbf{T}} \mathbf{V}^{-1} \mathbf{y}}{\mathbf{x}_{ ext{snp}}^{\mathbf{T}} \mathbf{V}^{-1} \mathbf{x}_{ ext{snp}}} ext{with } ext{var} \left( \hat{eta}_{ ext{snp}} 
ight) = rac{1}{\mathbf{x}_{ ext{snp}}^{\mathbf{T}} \mathbf{V}^{-1} \mathbf{x}_{ ext{snp}}}$$

V =  $\sigma_g^2 \, \psi + \sigma_e^2 I$ 

Overall phenotypic variance-covariance matrix = genetic + error

Obtain z scores and p-values of the effect based on this.

**GCTA** 

a tool for Genome-wide Complex Trait Analysis

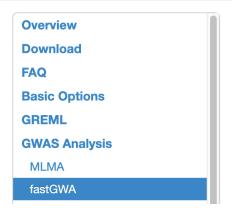
GCTA

SMR GSMR

OSCA

CTG forum

Yang Lab



fastGWA

fastGWA: A fast MLM-based Genome-Wide Association tool

fastGWA is an ultra-efficient tool for mixed linear model (MLM)-based GWAS analysis of biobank-scale data such as the UK Biobank (see Jiang et al. *Nature Genetics* 2019 for details of the method). Credits: Longda Jiang (method, simulation and analysis), Zhili Zheng (method, software and analysis) and Jian Yang (method and overseeing).

We have applied fastGWA to 2,173 traits on 456,422 array-genotyped and imputed individuals and 2,048 traits on 49,960 whole-expression and (WFS) individuals in the LIK Richark. All the summary statistics are available.

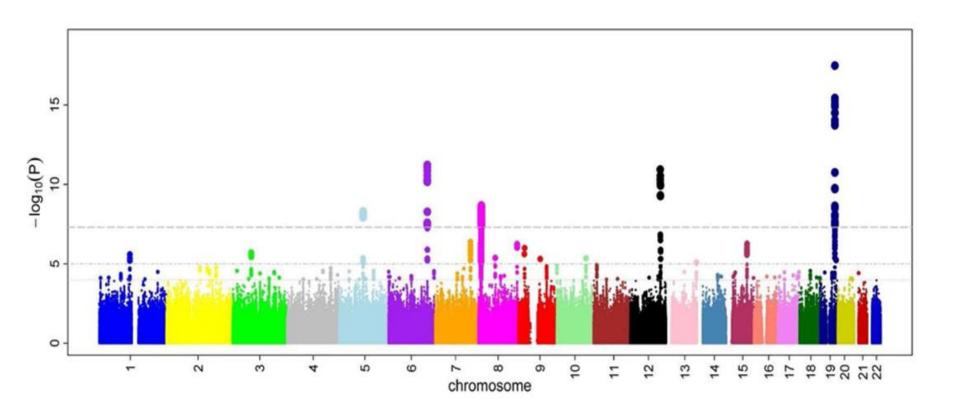
Check more recent approaches:

**SAIGE** (Zhou et al 2018 *Nat Genet* ), **REGENIE** (Mbatchou et al 2021 *Nat Genet*)

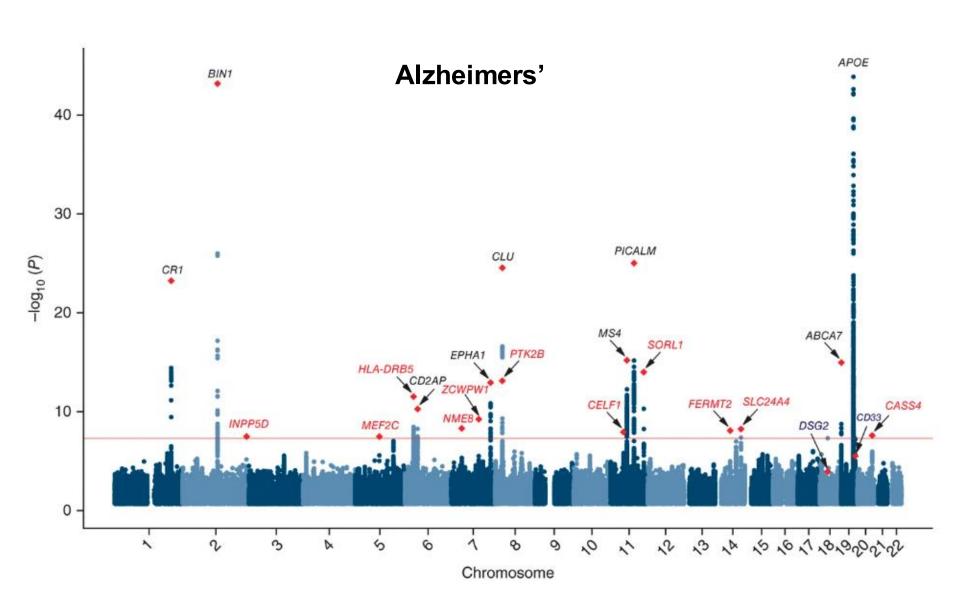
Jiang et al 2019 Nat Genet

### Standard visualization technique for GWAS results

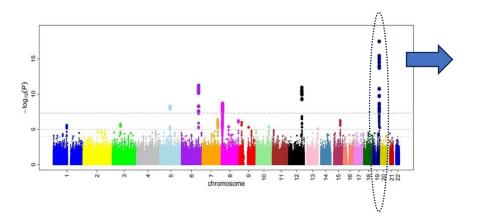
### Schizophrenia



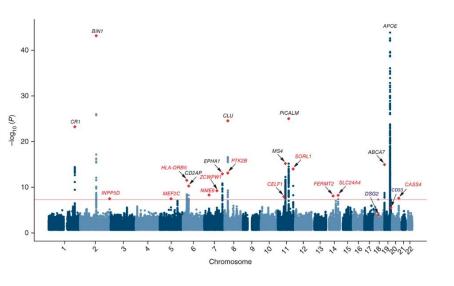
### Standard visualization technique for GWAS results



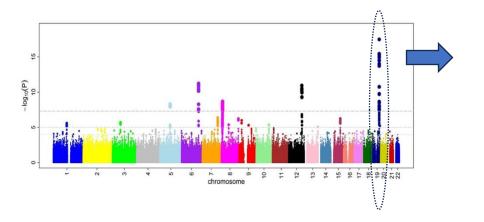
#### What is common between these GWAS-es?



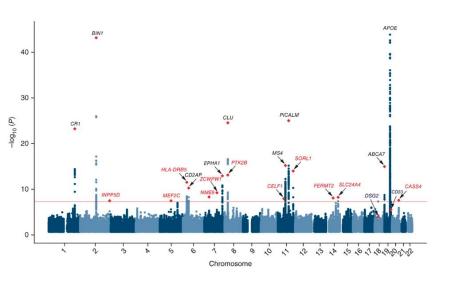
GWAS hits occur in clusters of variants all showing significant effects in same region – this is because of high linkage disequilibrium.



#### What is common between these GWAS-es?



GWAS hits occur in clusters of variants all showing significant effects in same region – this is because of high linkage disequilibrium.



GWAS signals are highly polygenic encompassing many genes.

We are likely missing out many weaker GWAS effect signals due to stringent p-value thresholds.

Can genotypes explain phenotypic variance across individuals?

**Heritability**: Proportion of phenotypic variance that can be attributed to genetic effects

Heritability of GWAS hits ( $h_{GWAS}^2$ ): Squared correlation between best fit linear model of all GWAS hits and the phenotype

$$\max_{w} [r^2 (\sum_{s \in GWAS \ hits} w_s X_{ns}, Y_n)]$$

**Heritability of all SNPs (h\_g^2):** Squared correlation between best fit linear model of all SNPs and the phenotype

$$\max_{w} [r^2(\sum_{s} w_s X_{ns}, Y_n)]$$

Can genotypes explain phenotypic variance across individuals?

**Heritability**: Proportion of phenotypic variance that can be attributed to genetic effects

Heritability of GWAS hits ( $h_{GWAS}^2$ ): Squared correlation between best fit linear model of all GWAS hits and the phenotype

$$\max_{w} [r^2 (\sum_{s \in GWAS \ hits} w_s X_{ns}, Y_n)]$$

**Heritability of all SNPs (h\_g^2):** Squared correlation between best fit linear model of all SNPs and the phenotype

$$\max_{w} [r^2(\sum_{s} w_s X_{ns}, Y_n)]$$

### How can we estimate SNP heritability $m{h}_g^2$ from full GWAS association summary statistics

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + u + \varepsilon$$

fixed effects random effects

 $Y = N \times 1$  vector of phenotypes

 $\mathbf{X} = N \times c$  matrix of c covariates

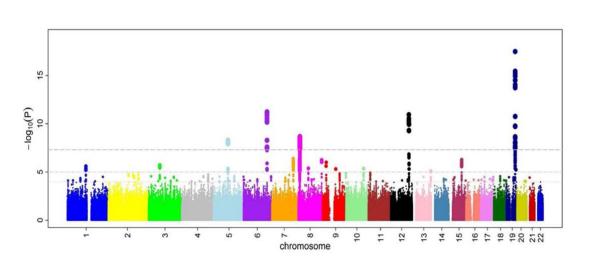
 $\mathbf{B} = c \times 1$  vector of effect sizes of c covariates

 $u \sim N(0, \sigma_g^2 \mathbf{A})$  is residual variance due to genetic effects  $\varepsilon \sim N(0, \sigma_e^2 \mathbf{I})$  is residual variance due to environmental effects

$$\mathbf{V} = \text{Var}(u + \varepsilon) = \sigma_g^2 \mathbf{A} + \sigma_e^2 \mathbf{I}$$

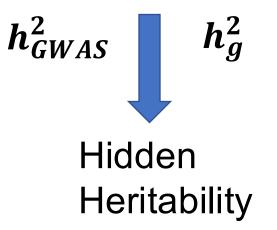
heritability 
$$h_g^2 = \sigma_g^2 / (\sigma_g^2 + \sigma_e^2)$$

### There is a big gap between only focusing on GWAS hits and looking at all of the GWAS association summary



### **Schizophrenia**

0.07 < 0.24



Lichtenstein et al 2009 *Lancet* Lee et al. 2012 *Nat Genet* Trubetskoy et al. 2022 *Nature* 

### This gap has been largely resolved for Adult height GWAS

### A saturated map of common genetic variants associated with human height

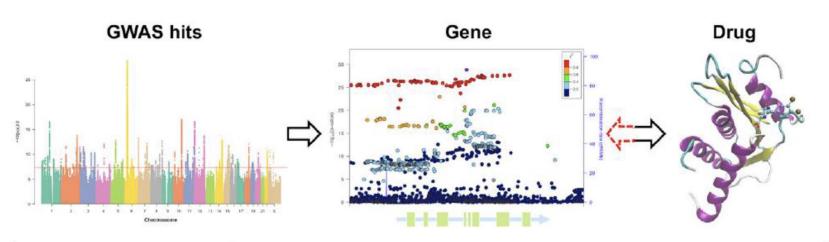
Loïc Yengo ☑, Sailaja Vedantam, Eirini Marouli, Julia Sidorenko, Eric Bartell, Saori Sakaue, Marielisa Graff, Anders U. Eliasen, Yunxuan Jiang, Sridharan Raghavan, Jenkai Miao, Joshua D. Arias, Sarah E. Graham, Ronen E. Mukamel, Cassandra N. Spracklen, Xianyong Yin, Shyh-Huei Chen, Teresa Ferreira, Heather H. Highland, Yingjie Ji, Tugce Karaderi, Kuang Lin, Kreete Lüll, Deborah E. Malden, 23andMe Research Team, VA Million Veteran Program, DiscovEHR (DiscovEHR and MyCode Community Health Initiative), eMERGE (Electronic Medical Records and Genomics Network), Lifelines Cohort Study, The PRACTICAL Consortium, Understanding Society Scientific Group, ... Joel N. Hirschhorn ☑ + Show authors

Nature (2022) Cite this article

"Here, using data from a genome-wide association study of <u>5.4 million individuals</u> of diverse ancestries, we show that <u>12,111 independent SNPs</u> that are significantly associated with height account for nearly all of the common SNP-based heritability."

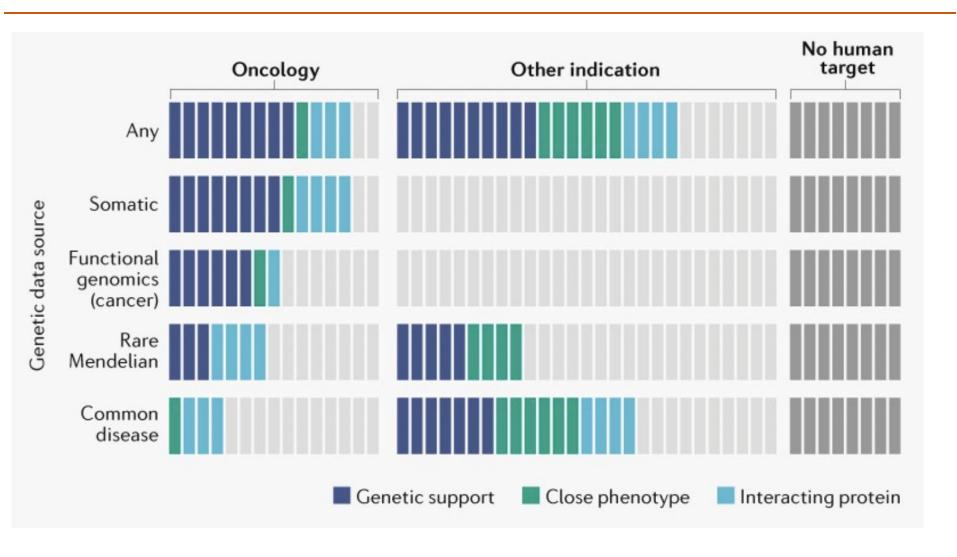
## Clinical and therapeutic implications of GWAS

### Is GWAS actually important? (GWAS hits to drugs)



Trait	Gene with GWAS hits	Known or candidate drug		
Type 2 Diabetes	SLC30A8/KCNJ11	ZnT-8 antagonists/Glyburide		
Rheumatoid Arthritis	PADI4/IL6R	BB-Cl-amidine/Tocilizumab		
Ankylosing Spondylitis(AS)	TNFR1/PTGER4/TYK2	TNF- inhibitors/NSAIDs/fostamatinib		
Psoriasis(Ps)	IL23A	Risankizumab		
Osteoporosis	RANKL/ESR1	Denosumab/Raloxifene and HR		
Schizophrenia	DRD2	Anti-psychotics		
LDL cholesterol	HMGCR	Pravastatin		
AS, Ps, Psoriatic Arthritis	IL12B	Ustekinumab		

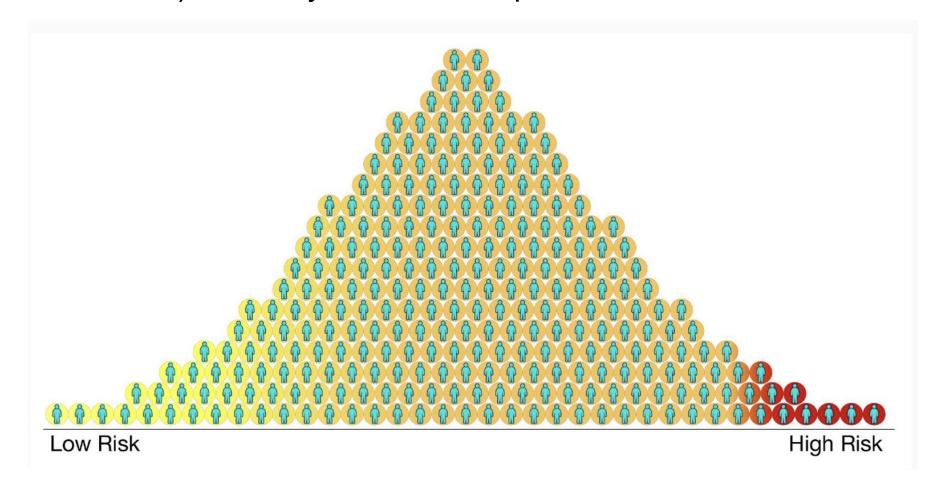
### Is GWAS actually important? (GWAS hits to drugs)



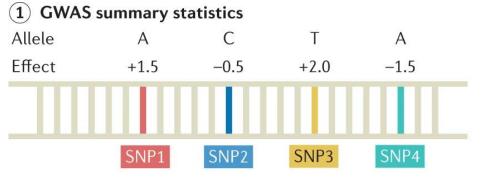
33 of 50 FDA approved drugs in 2021 have genetic support, with highest implicated from common disease GWAS.

### Is GWAS actually important? (Genetic risk score)

Identify the genetic risk for any individual for diseases and traits based on their genetic make-up (genotypes across all SNPs). Are they at risk for a specific disease?



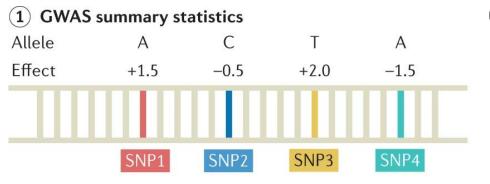
### How to calculate polygenic risk scores?



#### 2 Genotype data

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

### How to calculate polygenic risk scores?



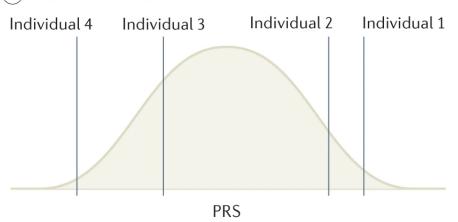
#### 3 Polygenic risk score

Individual 1	1.5	_	0.5	+	4.0	_	0.0	=	5.0
marviduat 1	1.5		0.5	·	7.0		0.0		5.0
Individual 2	1.5	_	0.0	+	2.0	_	1.5	=	2.0
Individual 3	0.0	-	1.0	+	2.0	-	1.5	=	-0.5
Individual 4	0.0	-	1.0	+	0.0	-	3.0	=	-4.0

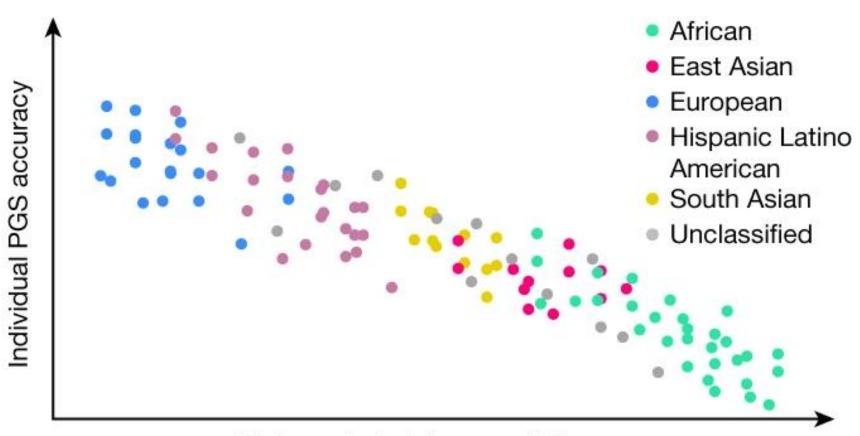
#### (2) Genotype data

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

#### 4 PRS distribution

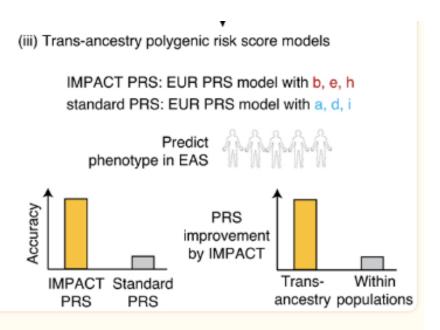


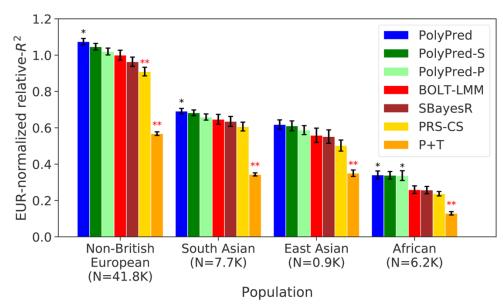
### A big challenge in polygenic risk scores (representation)



Distance to training population on genetic ancestry continuum

### Integrating functional annotations can resolve polygenic risk score transferability across populations





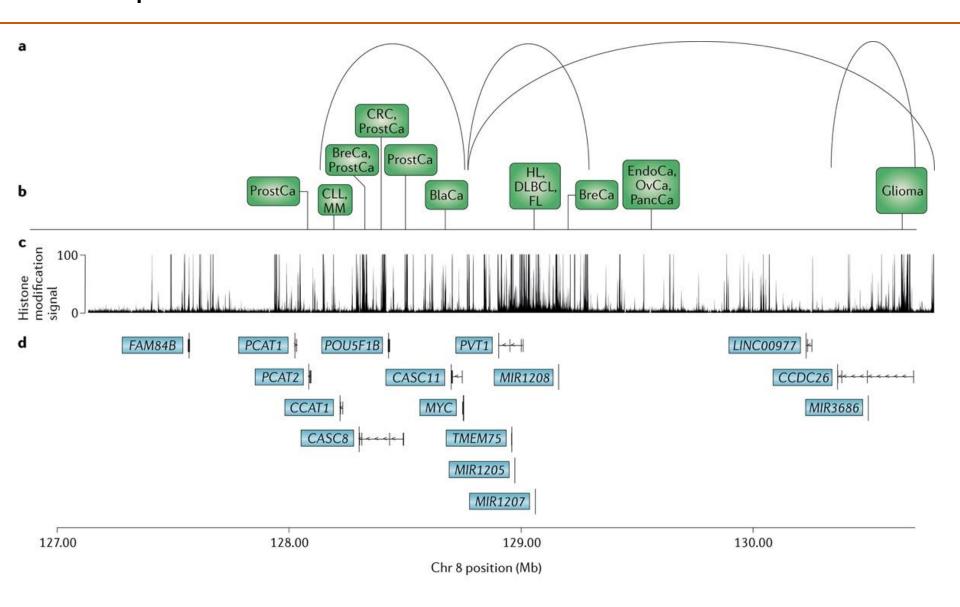
IMPACT: integrative combination of TF binding functional predictions

PolyPred: integrative combination of functional annotations to drive transancestry prediction risk accuracy

Amariuta..Dey et al 2020 Nat Genet

Weissbrod et al 2022 Nat Genet

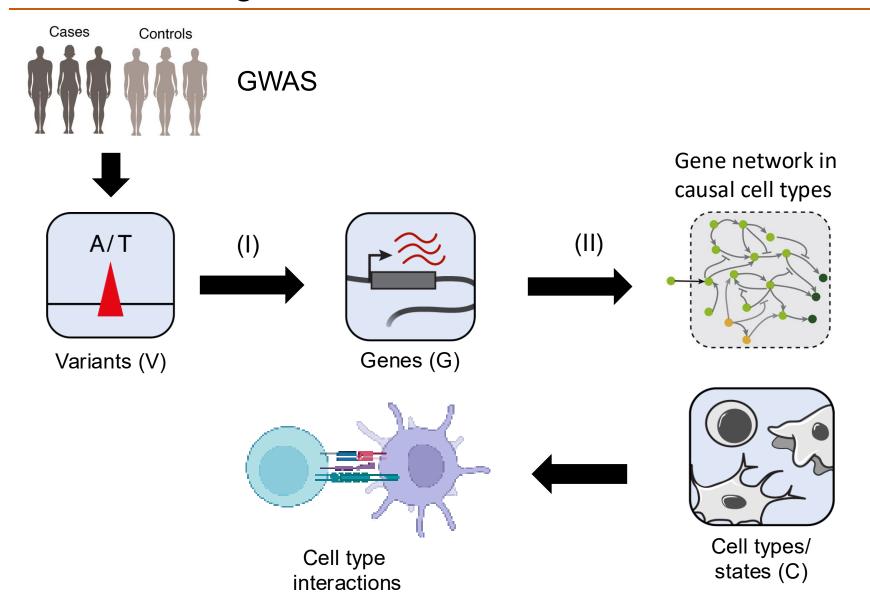
### Pleiotropic associations in Genome-wide associations



Cancer-related pleiotropy

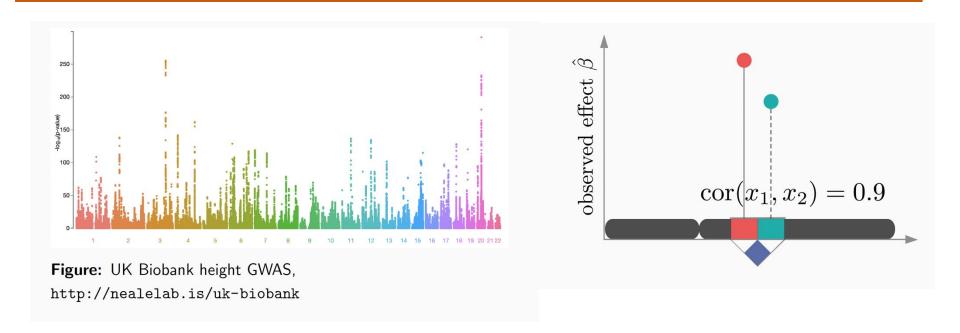
Learn shared associations across traits: **MTAG** (Turley et al 2018 *Nat Genet*)

### Understanding the functional basis of GWAS variants

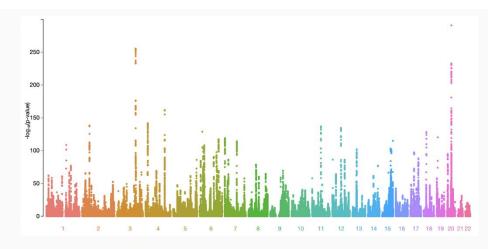


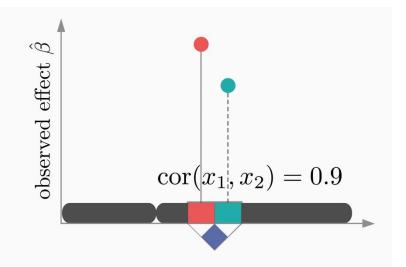
## GWAS-to-function (Overview)

### Linkage disequilibrium can hinder identification of causal variant for both GWAS and eQTL studies

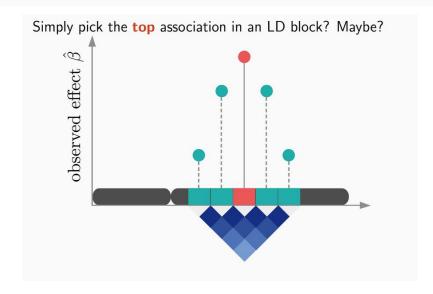


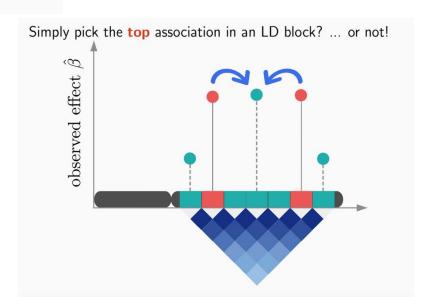
### Linkage disequilibrium can hinder identification of causal variant for both GWAS and eQTL studies



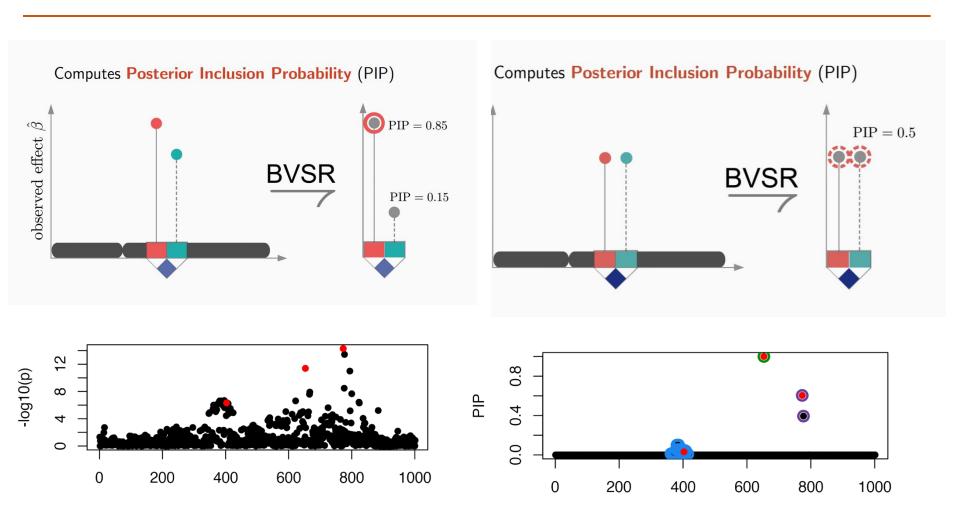


**Figure:** UK Biobank height GWAS, http://nealelab.is/uk-biobank





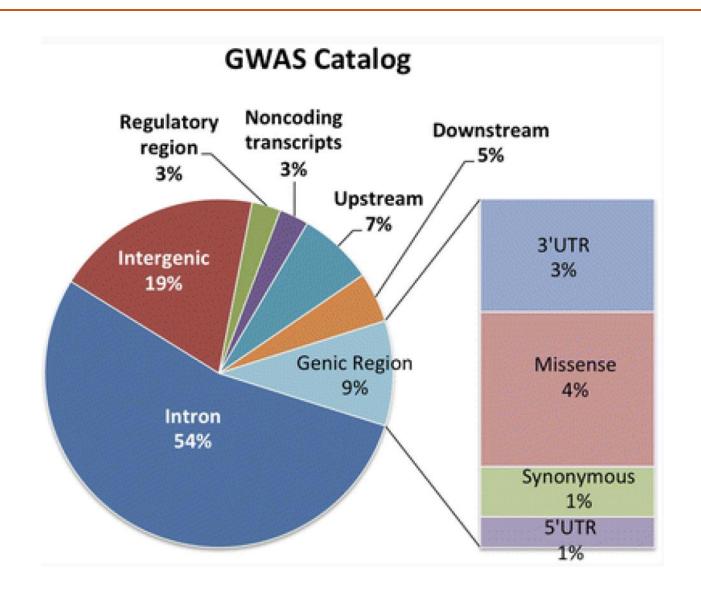
SuSIE: Method to perform Bayesian variable selection to identify independent causal GWAS variants or sets of variants when it is not sure



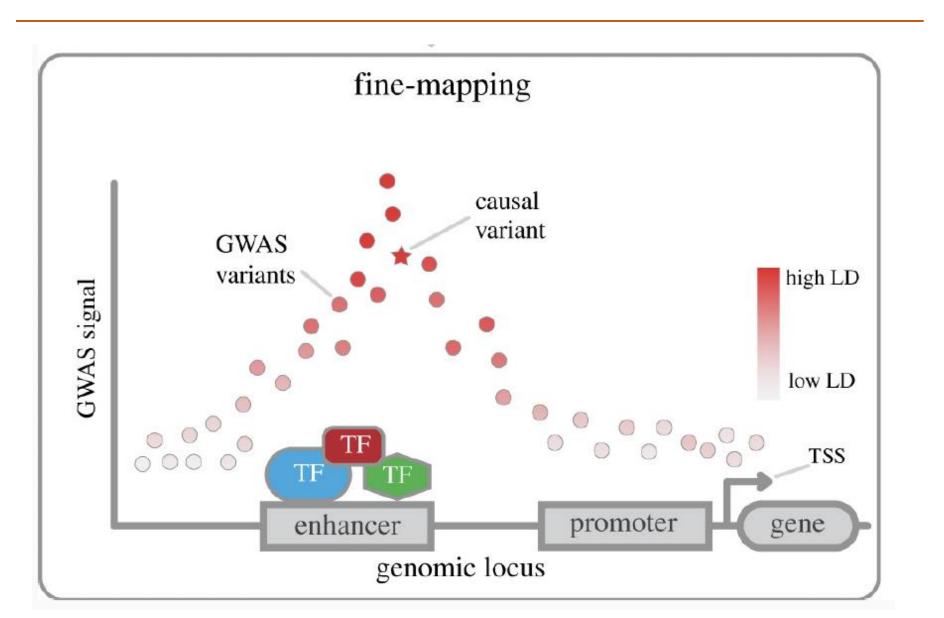
3 colors correspond to 95% <u>credible sets</u>: A credible set says a causal variant is within this set with 95% probability

SuSIE: Wang et al 2020 JRSS-B

### Making sense of the function of GWAS variants

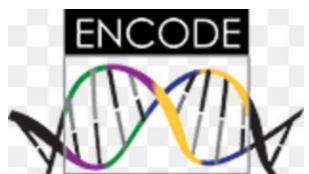


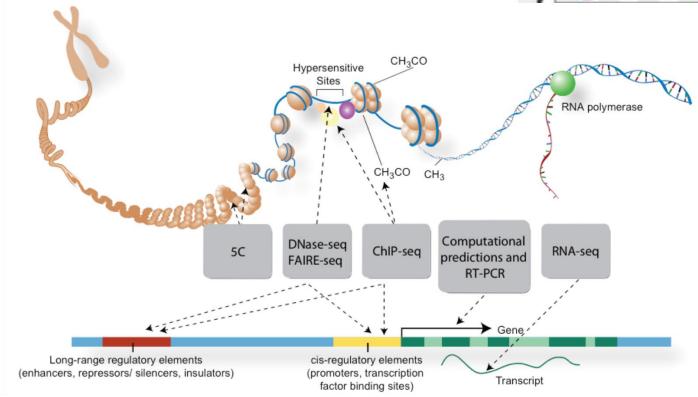
### GWAS signals can be confounded by LD. Can we use underlying function to find the causal variant?



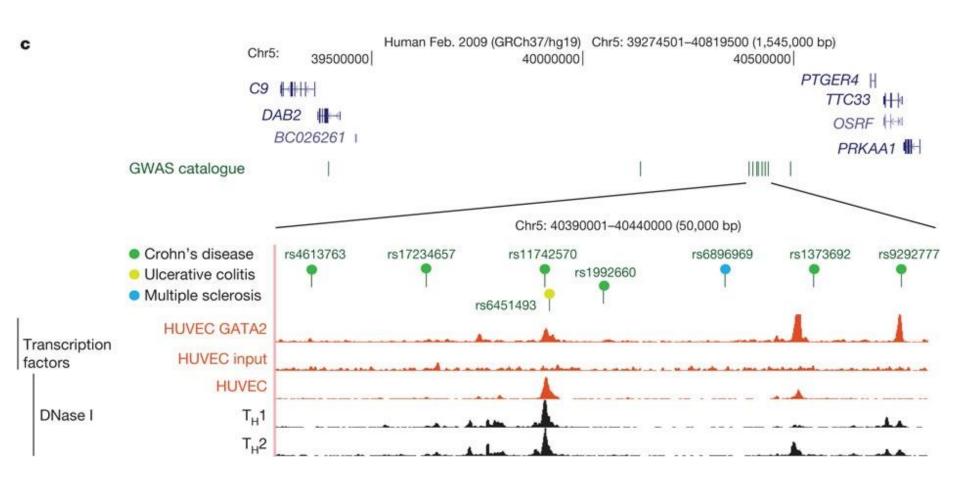
## Mapping the function of genomic elements in DNA (ENCODE: 2003-2023)

In 2003, the ENCODE consortium was launched by NIH to study the non-coding regions of the genome and identify functional elements





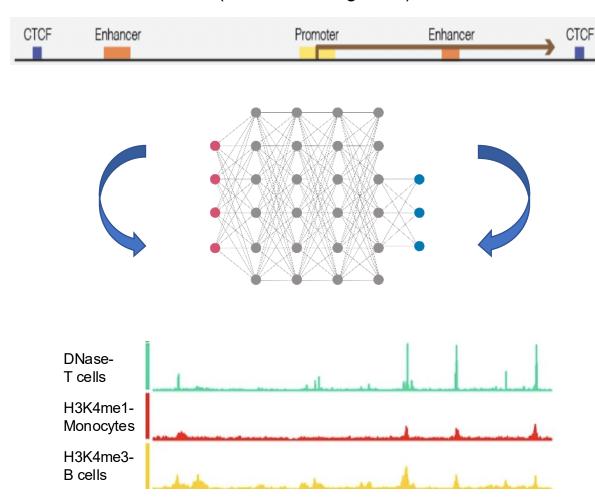
## Overlapping genome-wide functional annotation tracks against GWAS disease-associated variants



## Sequence-based deep learning models trained on epigenomic features

### **DNA** sequences

(1-hot encoding DNA)

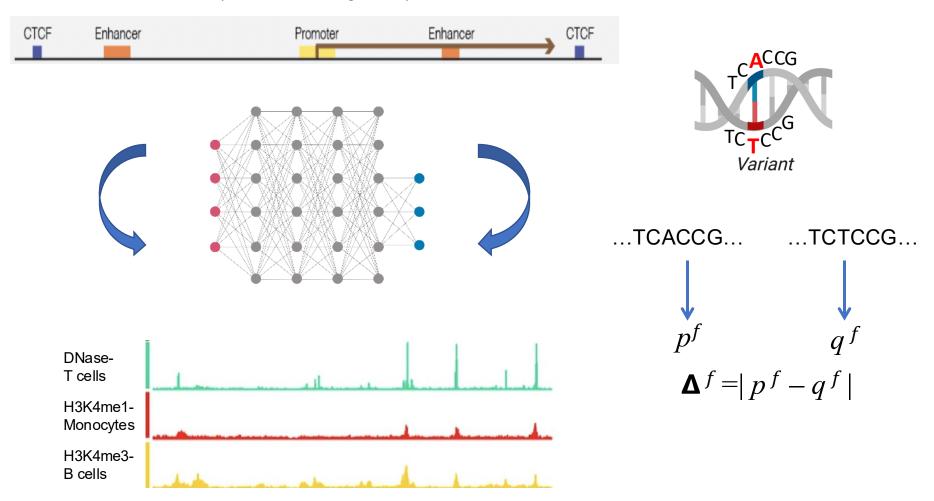


For each sequence, generates a prediction of affinity for each feature *f* at the site of the sequence.

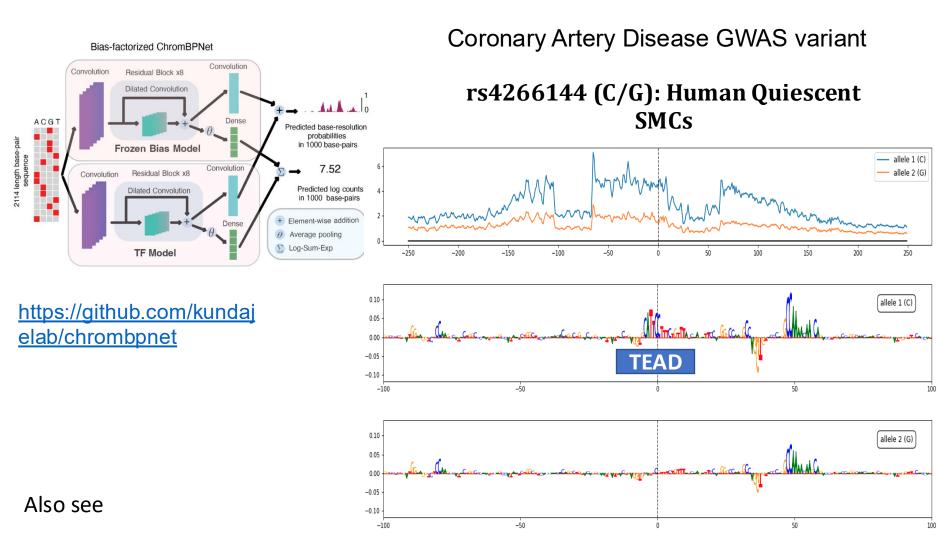
## Sequence-based deep learning models trained on epigenomic features

### **DNA** sequences

(1-hot encoding DNA)



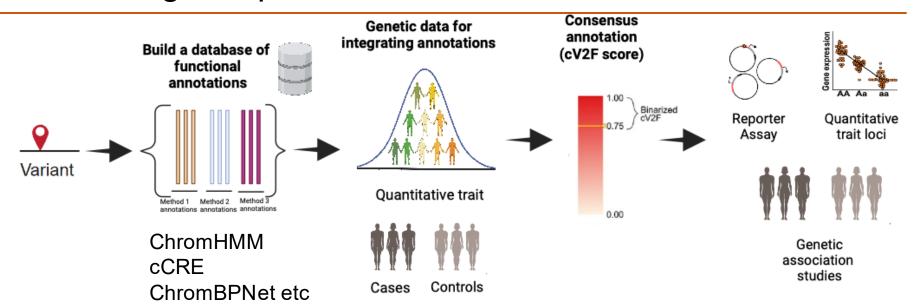
## ChromBPNet deep learning model captures sequence mediated function at GWAS variants



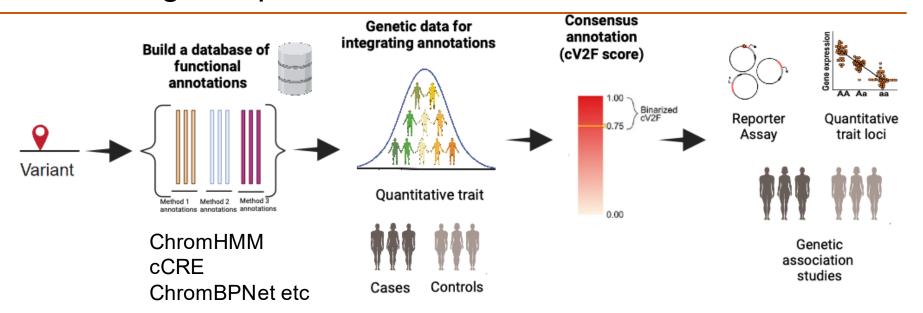
**Enformer**: Avsec et al 2021 Nat Methods **BPNet**: Avsec et al 2021 Nat Genet

Pampari et al 2024 bioRxiv Courtesy: Anshul Kundaje, Stanford

### Generating an optimal variant-to-disease function score



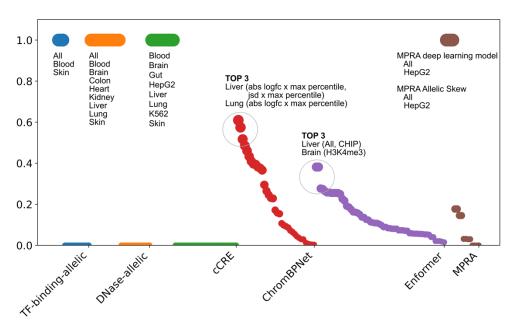
### Generating an optimal variant-to-disease function score



rs12740374

cV2F-score = 0.96 cV2F-Liver = 0.93

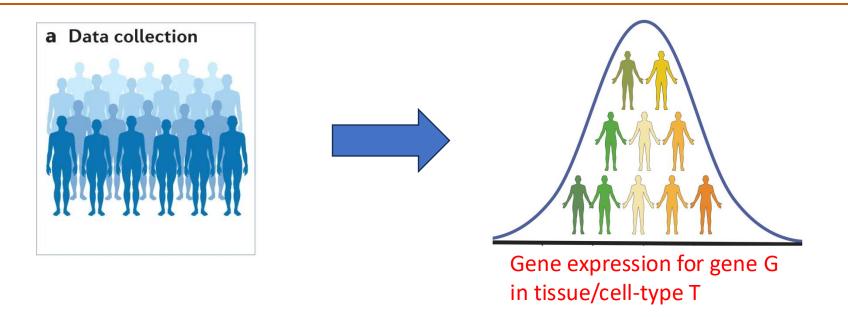
creates a C/EBP transcription factor binding site and alters the hepatic expression of the *SORT1* gene (Musunuru et al 2010 *Nature*)



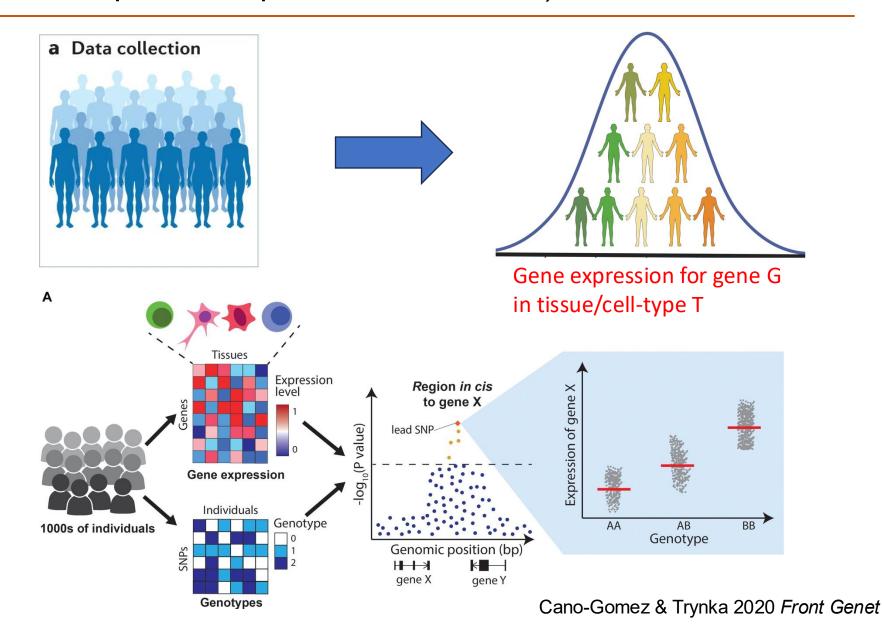
Fabiha et al 2024 bioRxiv, in rev Nat Genet

## Linking quantitative trait loci to disease GWAS

## Tracking genetic variation of gene expression phenotype (eQTL: expression quantitative trait loci)

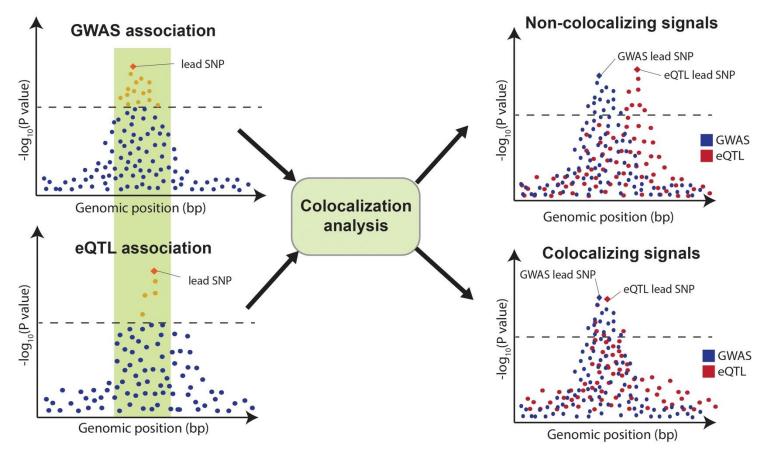


## Tracking genetic variation of gene expression phenotype (eQTL: expression quantitative trait loci)



## Statistical colocalization: Identifying shared causal variants between a disease trait and an eQTL

Typically performed for one gene and for one tissue separately against one focal disease GWAS.

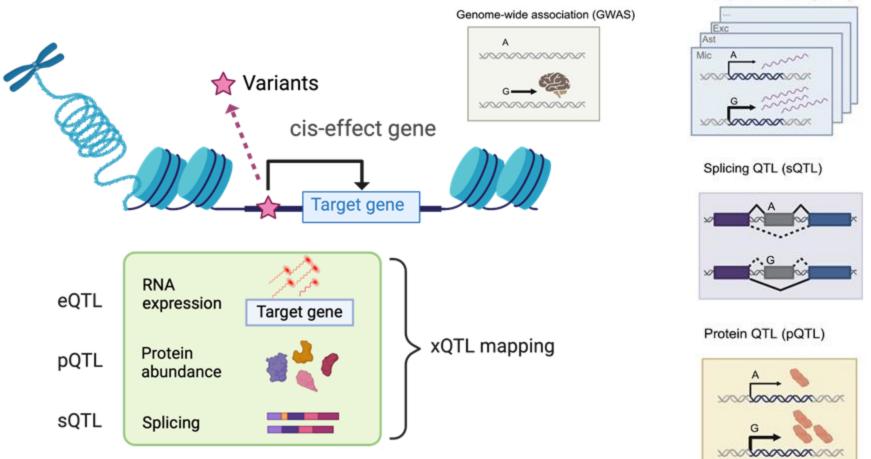


**Coloc**: standard method for colocalization. Does not scale well to more than 2 phenotypes.

## ColocBoost model to perform multimodal molecular phenotype QTL colocalization

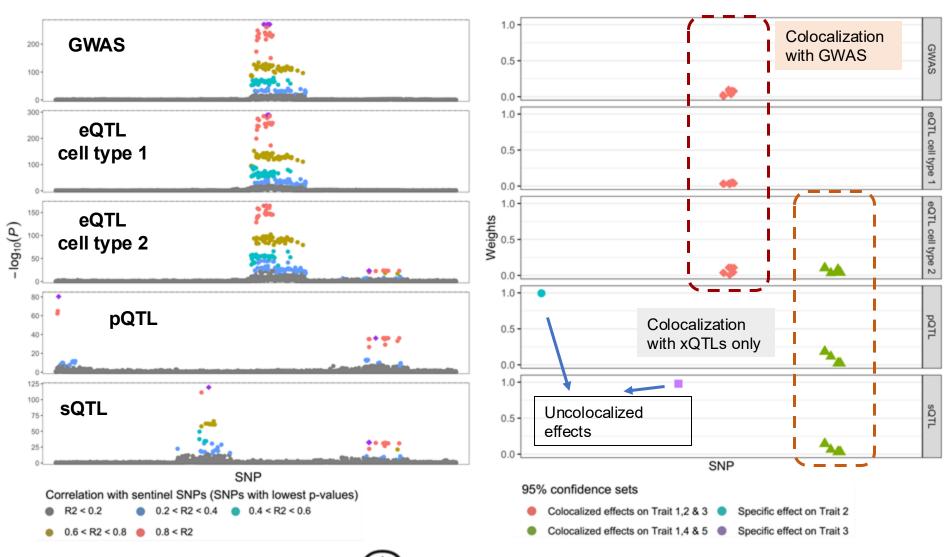
Recent advances in technology has made it easier to use other molecular phenotypes outside of gene expression, and also assess eQTL at cell type resolution for different cell types in a tissue

Expression QTL (eQTL)



Aguet et al. 2020. Nature Reviews Methods Primers

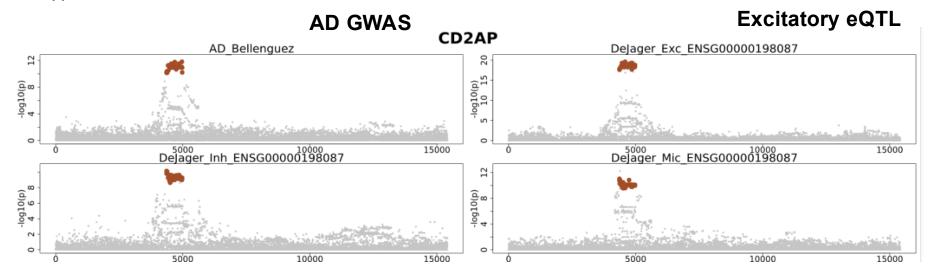
## Understanding colocalization: enhancing GWAS insights through shared genetic signals





## Shared genetic regulation across cell types observed for many disease risk variants are not indicative of cell-cell crosstalk

37.3% of AD causal risk variants show genetic regulation shared across multiple cell types in brain.



#### Inhibitory eQTL

Microglia eQTL

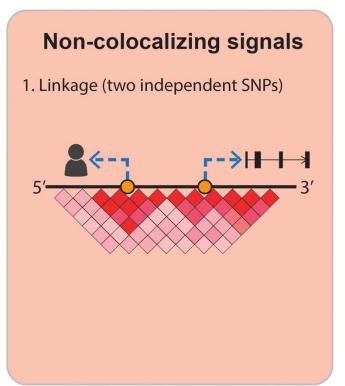
### Alzheimer's disease risk gene *CD2AP* is a dose-sensitive determinant of synaptic structure and plasticity

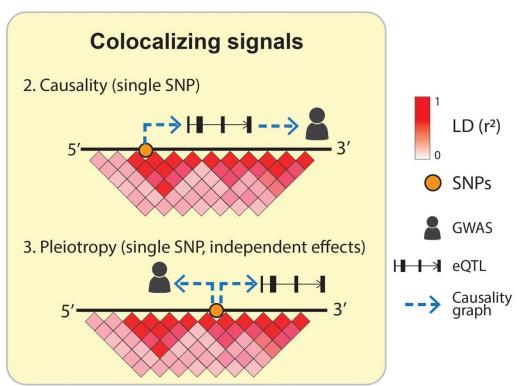
Matea Pavešković <sup>1,2,3</sup>, Ruth B De-Paula <sup>4,5,6</sup>, Shamsideen A Ojelade <sup>7,8</sup>, Evelyne K Tantry <sup>9,10</sup>, Mikhail Y Kochukov <sup>11,12</sup>, Suyang Bao <sup>13,14</sup>, Surabi Veeraragavan <sup>15,16</sup>, Alexandra R Garza <sup>17,18</sup>, Snigdha Srivastava <sup>19,20,21</sup>, Si-Yuan Song <sup>22,23</sup>, Masashi Fujita <sup>24</sup>, Duc M Duong <sup>25</sup>, David A Bennett <sup>26</sup>, Philip L De Jager <sup>27</sup>, Nicholas T Seyfried <sup>28</sup>, Mary E Dickinson <sup>29,30</sup>, Jason D Heaney <sup>31</sup>, Benjamin R Arenkiel <sup>32,33,34,#</sup>, Joshua M Shulman <sup>35,36,37,38,39,#,∞</sup>

### Microglial CD2AP deficiency exerts protection in an Alzheimer's disease model of amyloidosis

Lingliang Zhang <sup># 1</sup>, Lingling Huang <sup># 1</sup>, Yuhang Zhou <sup>1</sup>, Jian Meng <sup>1</sup>, Liang Zhang <sup>1</sup>, Yunqiang Zhou <sup>1</sup>, Naizhen Zheng <sup>1</sup>, Tiantian Guo <sup>1</sup>, Shanshan Zhao <sup>1</sup>, Zijie Wang <sup>1</sup>, Yuanhui Huo <sup>1</sup>, Yingjun Zhao <sup>1</sup>, Xiao-Fen Chen <sup>1</sup>, Honghua Zheng <sup>1</sup>, David M Holtzman <sup>2</sup>, Yun-Wu Zhang <sup>3</sup>

## Colocalization of GWAS and eQTL signals to identify shared causal variants between disease and expression





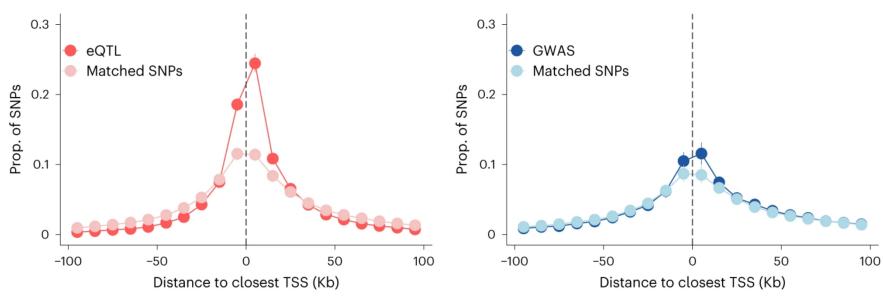
## Integrating eQTL and enhancer-gene maps to link variants to genes (OpenTargets)



Data type	Experiment type	Source	Weighting
<i>In silico</i> functional prediction	Transcript consequence	VEP	1.0
QTL	sQTL	GTEx v8	1.0
QTL	eQTL	many	0.66
QTL	pQTL	many	0.66
Interaction	PCHi-C	Javierre <i>et al.</i> (Cell, 2016)	0.33
Interaction	Enhancer-TSS correlation	Andersson <i>et al.</i> (Nature, 2014)	0.33
Interaction	DHS-promoter correlation	Thurman <i>et al.</i> (Nature, 2012)	0.33
Distance	Canonical TSS		0.33

### Systemic differences between eQTLs and GWAS

a Enrichment near transcription start sites (TSSs)



Nat Genet 2023 Nov:55(11):1866-1875 doi: 10.1038/s41588-023-01529-1 Epub 2023 Oct 19

### Systematic differences in discovery of genetic effects on gene expression and complex traits

Hakhamanesh Mostafavi $^{\,1}$ , Jeffrey P Spence $^{\,2}$ , Sahin Naqvi $^{\,2}\,^{\,3}$ , Jonathan K Pritchard  $^{\,4}\,^{\,5}$ 

Affiliations + expand

PMID: 37857933 DOI: 10.1038/s41588-023-01529-1

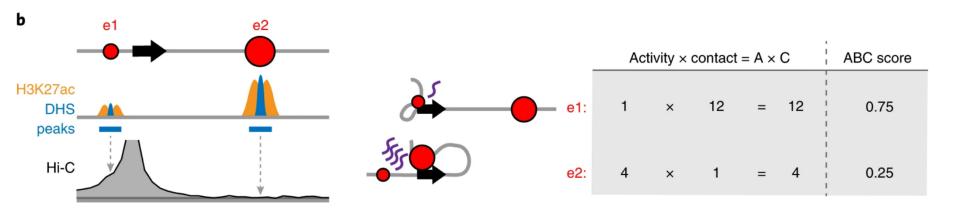
"GWAS and cis-eQTL hits are systematically different:
eQTLs cluster strongly near transcription start sites,
whereas GWAS hits do not. Genes near GWAS hits are
enriched in key functional annotations, are under strong
selective constraint and have complex regulatory
landscapes across different tissue/cell types, whereas
genes near eQTLs are depleted of most functional
annotations.

## Linking GWAS variants to target genes

### Broadening the scope of approaches to link variants to genes

S2G strategies	Description
5kb	SNPs in 5kb window around gene
100kb	SNPs in 100 kb window around gene
Promoter	SNPs in promoter region of the gene
TSS	SNPs in and around Transcription start sites
Coding	SNPs in coding regions of the gene

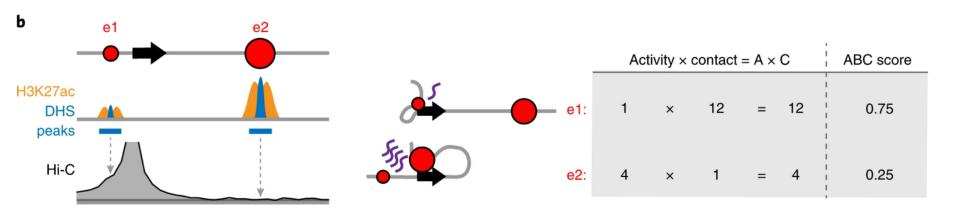
### The Activity-By-Contact element-gene linking method

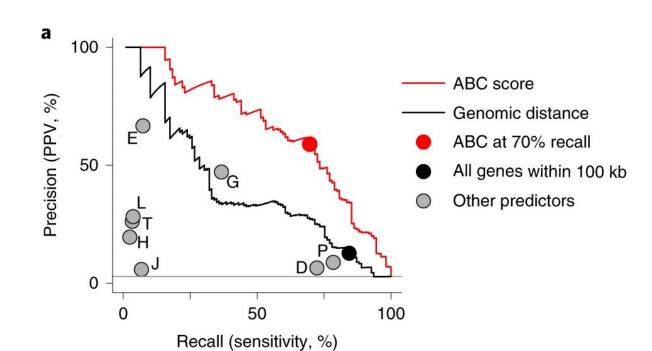


$$ext{ABC score}_{E,G} = rac{A_E imes C_{E,G}}{\sum\limits_{ ext{all elements e within 5 Mb of } G} A_e imes C_{e,G}}$$

Operationally, we estimated Activity (A) as the geometric mean of the read counts of DHS and H3K27ac chromatin immunoprecipitation sequencing (ChIP–seq) at element E, and Contact (C) as the KR-normalized Hi-C contact frequency between E and the promoter of gene G at 5-kb resolution (see Supplementary Note E and Supplementary Figs. E and E and E and E0.

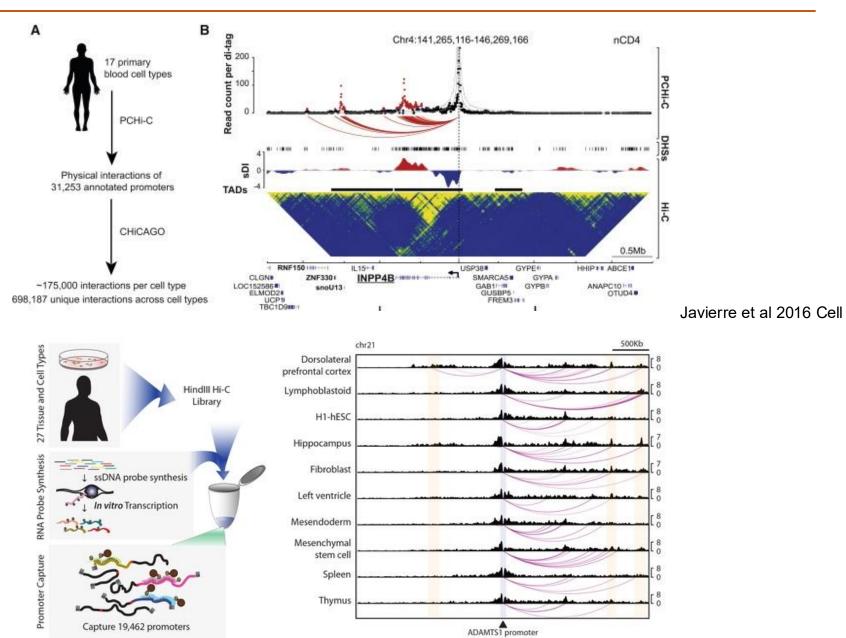
### The Activity-By-Contact element-gene linking method





Fulco et al 2019 Nat Genet Nasser et al 2021 Nature

### Promoter-capture Hi-C to link elements to genes



Jung et al 2020 Nat Genet

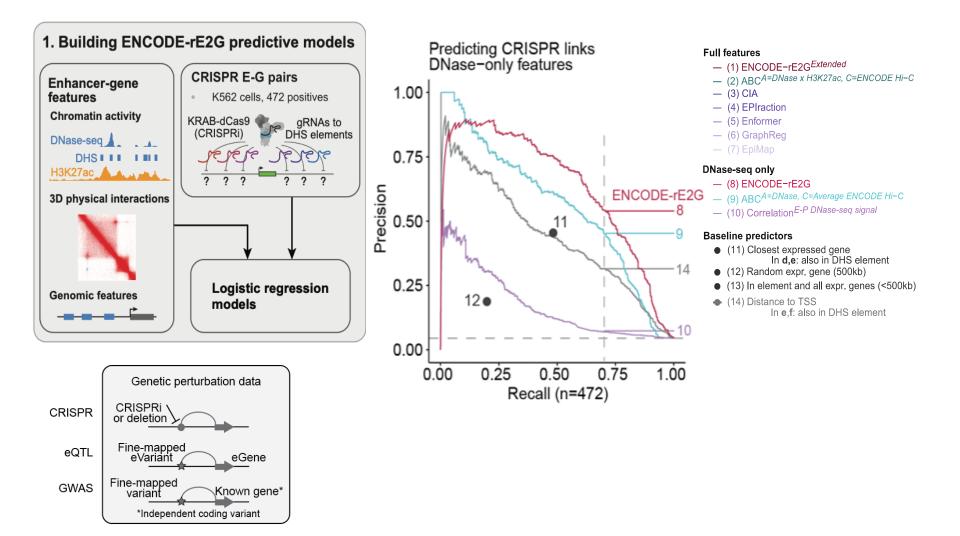
### Broadening the scope of approaches to link variants to genes

Naive S2G	Expression S2G	Hi-C S2G
	S2G strategies	Description
	5kb	SNPs in 5kb window around gene
	100kb	SNPs in 100 kb window around gene
	Promoter	SNPs in promoter region of the gene
	TSS	SNPs in and around Transcription start sites
	Coding	SNPs in coding regions of the gene
	eQTL	Max. post. causal probability in GTEx blood <sup>1,2</sup>
	ATAC	Correlated ATAC-seq peaks and gene expression in blood <sup>3</sup>
	Roadmap	Correlated enhancers and gene expression in blood <sup>4,5,6</sup>
	PC-HiC	Promoter Capture Hi-C <sup>7</sup>
	ABC	DHS ∩ H3K27ac ∩ Hi-C in blood <sup>8</sup>

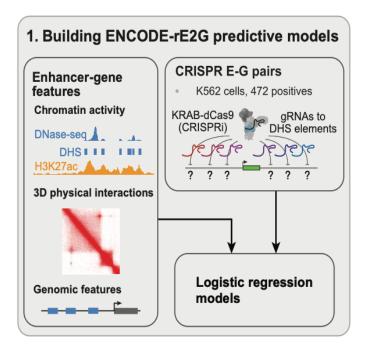
<sup>1</sup>Hormozdiari et al 2018 NG, <sup>2</sup>Aguet et al 2019 bioRxiv, <sup>3</sup>Yoshida et al 2019 Cell, <sup>4</sup>Liu et al 2017 Gen. Biol. , <sup>5</sup>Ernst et al 2011 Nat. meth., <sup>6</sup>Kundaje et al 2015 Nature , <sup>7</sup>Javierre et al 2016 Cell,, <sup>8</sup>Fulco et al 2019 Nat.Genet Dey et al 2022, *Cell Genomics*,

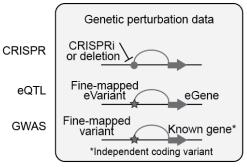
Gazal..Dey et al 2022 Nat Genet

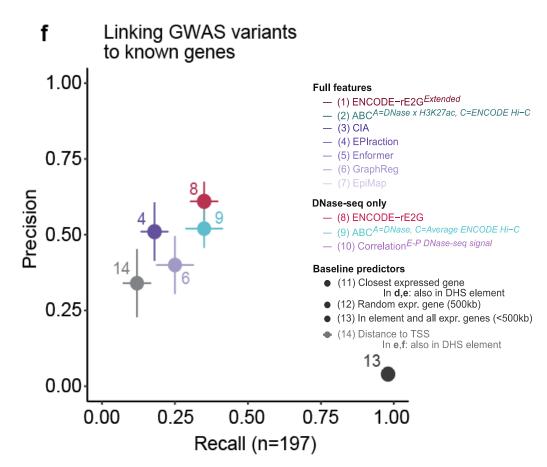
### Benchmarking different element-gene linking approaches



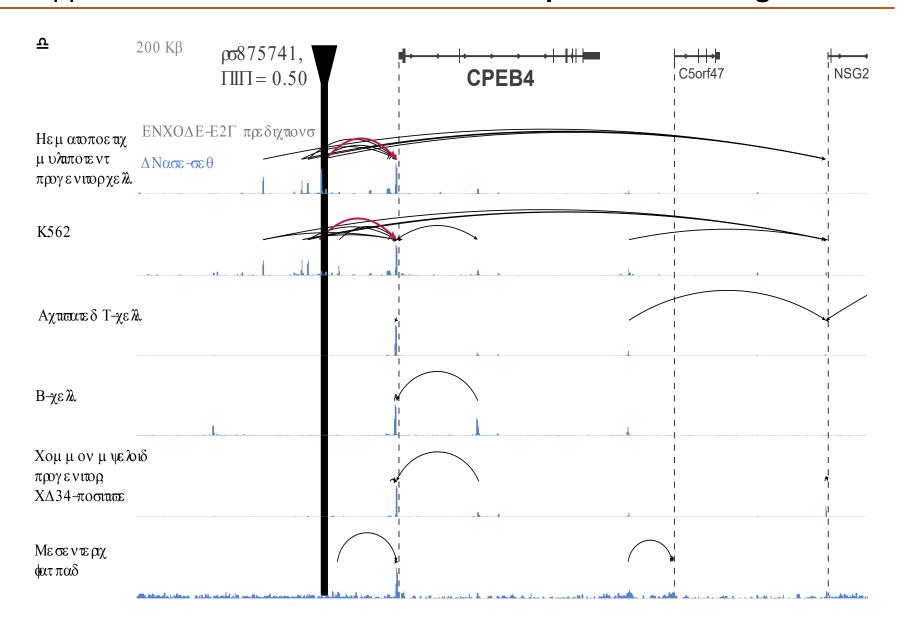
### Benchmarking different element-gene linking approaches



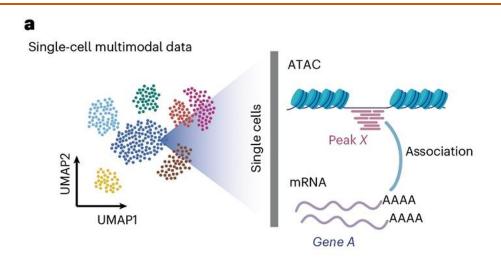


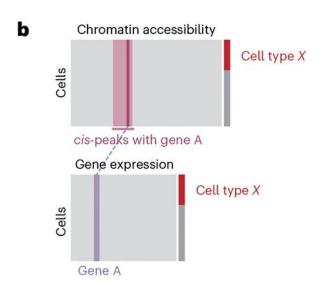


### Visualizing the element-gene links underlying **rs875741**: fine-mapped variant **PIP = 0.50 for mean corpuscular hemoglobin**.

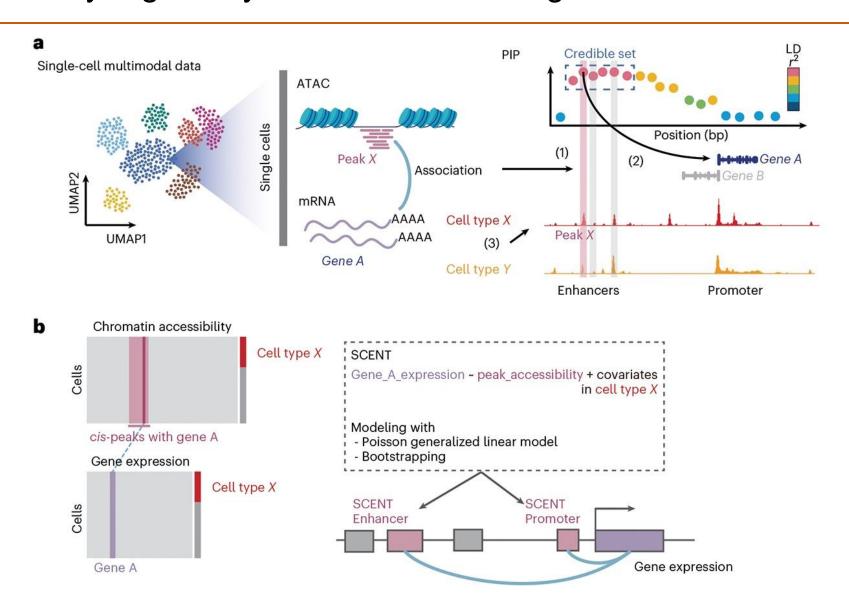


## Using single-cell RNA and ATAC data in the same cell to identify regulatory elements linked to genes

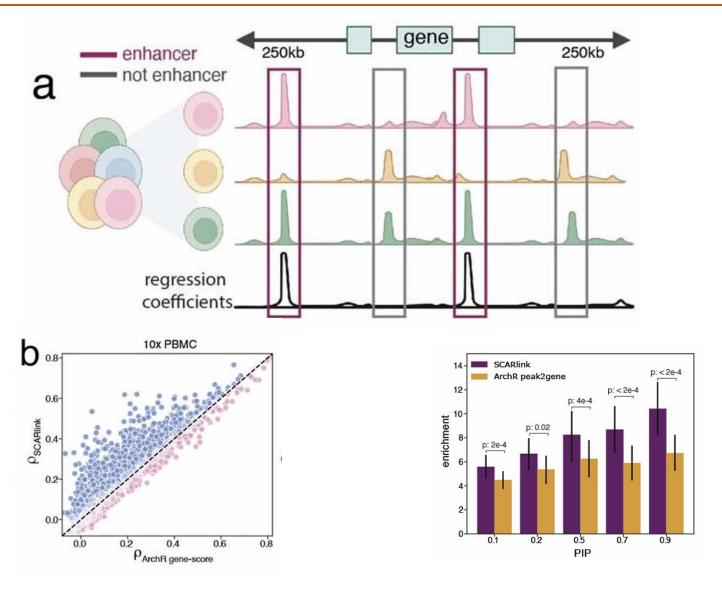




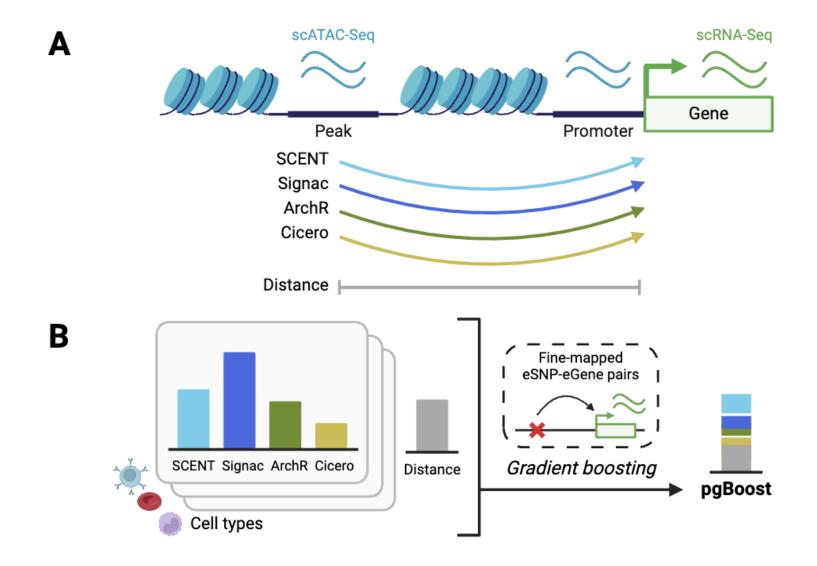
## Using single-cell RNA and ATAC data in the same cell to identify regulatory elements linked to genes



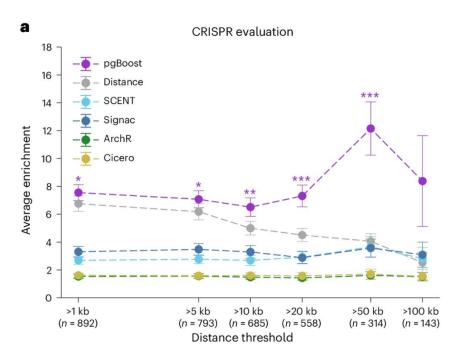
## Using single-cell RNA and ATAC data in the same cell to identify regulatory elements linked to genes

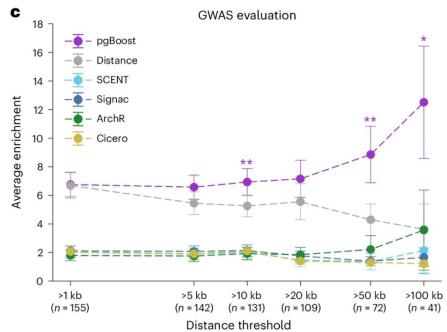


## Integrative model of multiome peak-gene links and genomic distance is informative for diseases



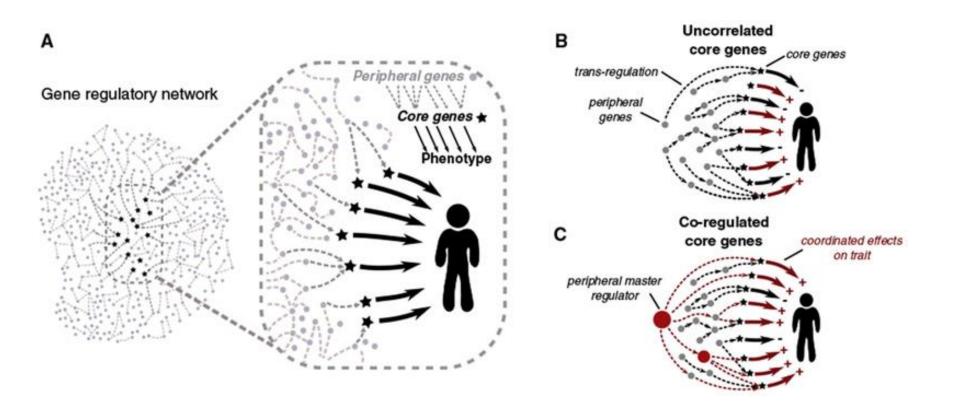
## Integrative model of multiome peak-gene links and genomic distance is informative for diseases





# Linking genes and gene programs to GWAS disease risk

### Beyond immediate target genes for GWAS variants, what are the gene programs or pathways affected by GWAS variants



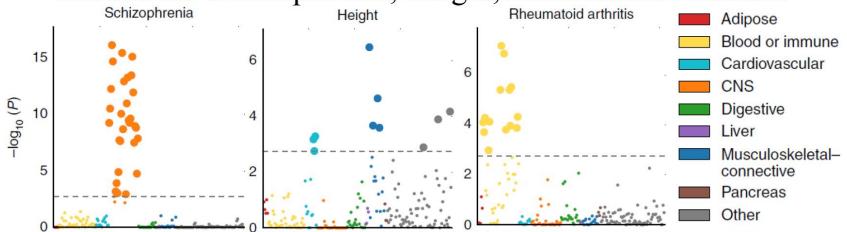
Liu, Li, Pritchard 2019 Cell

## Specifically expressed genes in a tissue show tissue-specific heritability enrichments

SEG genes: Genes in top 10% of most enriched expression in a focal tissue compared to other GTEx tissues.

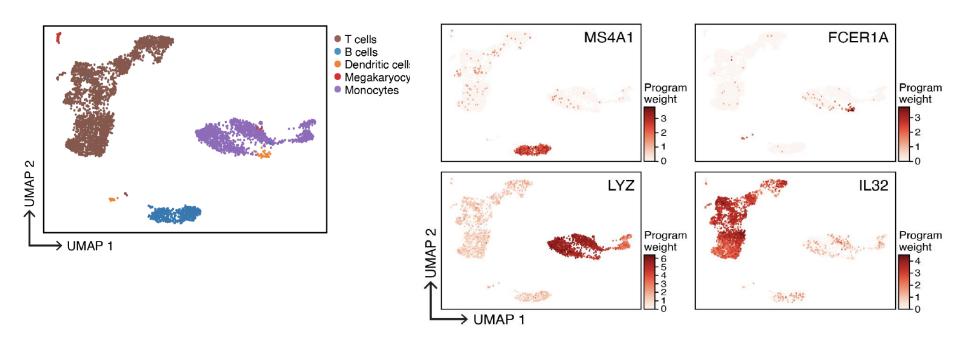
Annotate variants if they lie in a 100KB window around the SEG\_genes



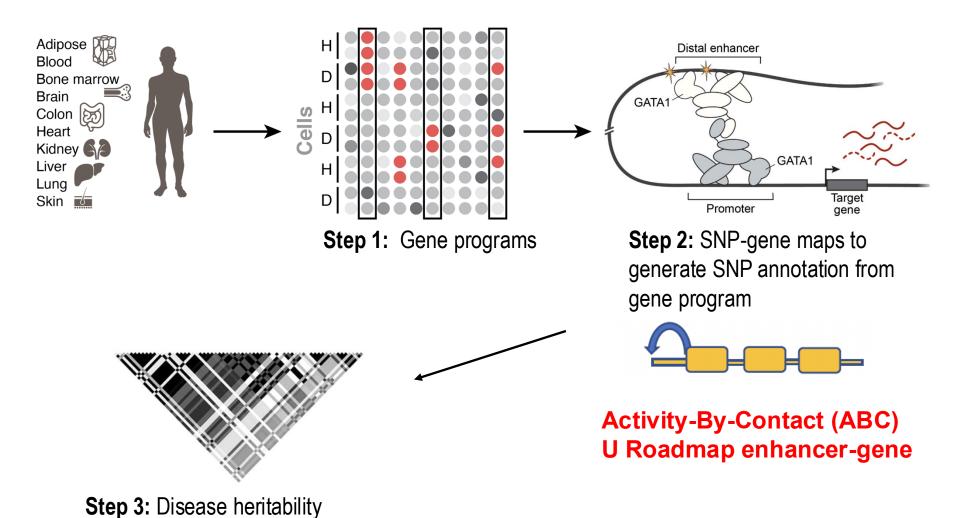


### Constructing gene programs to characterize cell type

**Cell type program**: Genes specifically expressed in an annotated cell type compared to other cell types in the tissue [nonparametric DE: Wilcoxon rank-sum test].



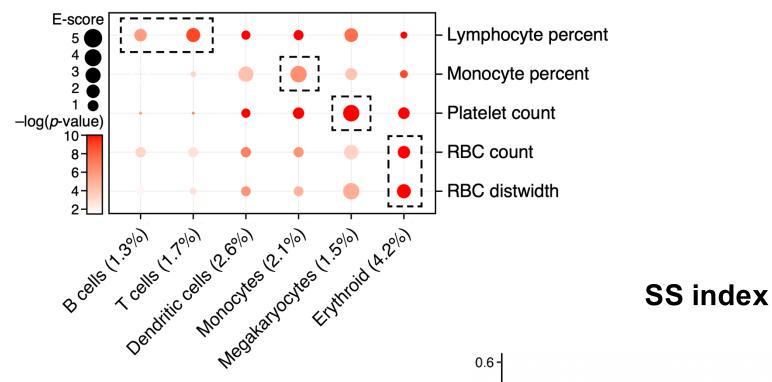
## sc-linker: Using single-cell RNA-seq to assess cell-type-specific heritability enrichments of specifically expressed genes



enrichment

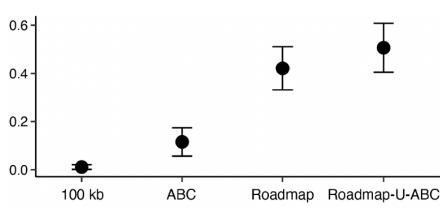
Jagadeesh\*, Dey\* et al 2022 Nat Genet; Delorey\* ... Dey\* et al 2021, Nature

### sc-linker: Using single-cell RNA-seq to assess cell-type-specific heritability enrichments of specifically expressed genes

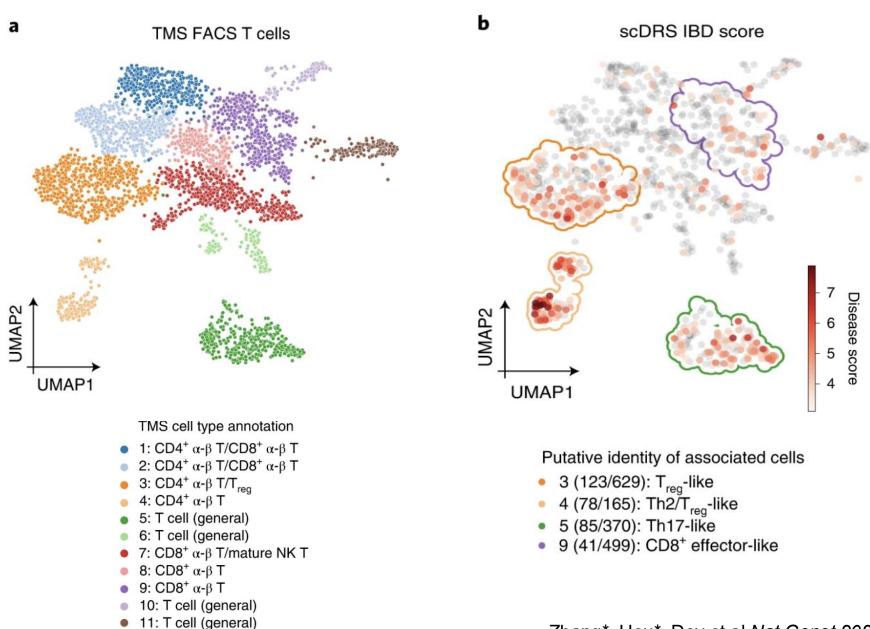


### Sensitivity-specificity (SS) index:

Average difference in disease signal between boxed pairs of cell types and traits versus all other pairs.

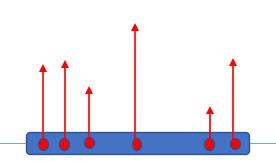


### Scoring a single cell for disease association using scDRS



Zhang\*, Hou\*, Dey et al Nat Genet 2022

### Prioritizing genes for a complex disease (MAGMA and PoPS)

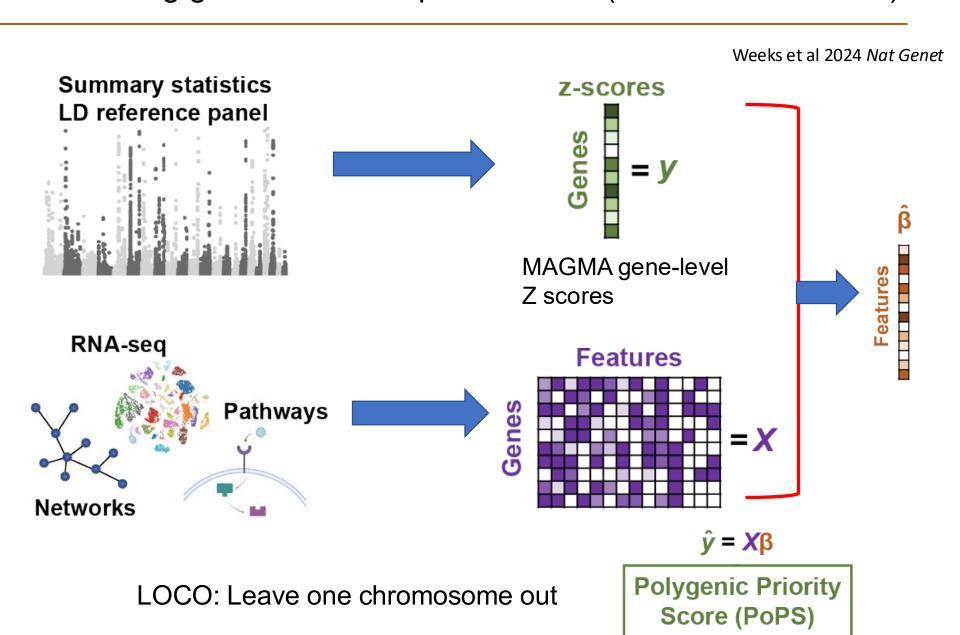


Two types of gene test statistics have been implemented in MAGMA:

- (a) The mean of the  $\chi^2$  statistic for the SNPs in a gene,
- (b) The top  $\chi^2$  statistic among the SNPs in a gene.

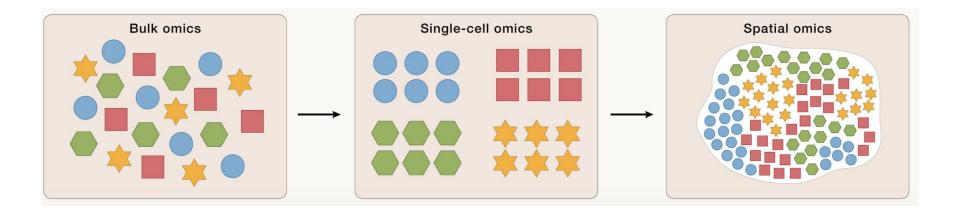
For the mean  $\chi^2$  statistic, a gene p-value is then obtained by using a known approximation of the sampling distribution

### Prioritizing genes for a complex disease (MAGMA and PoPS)

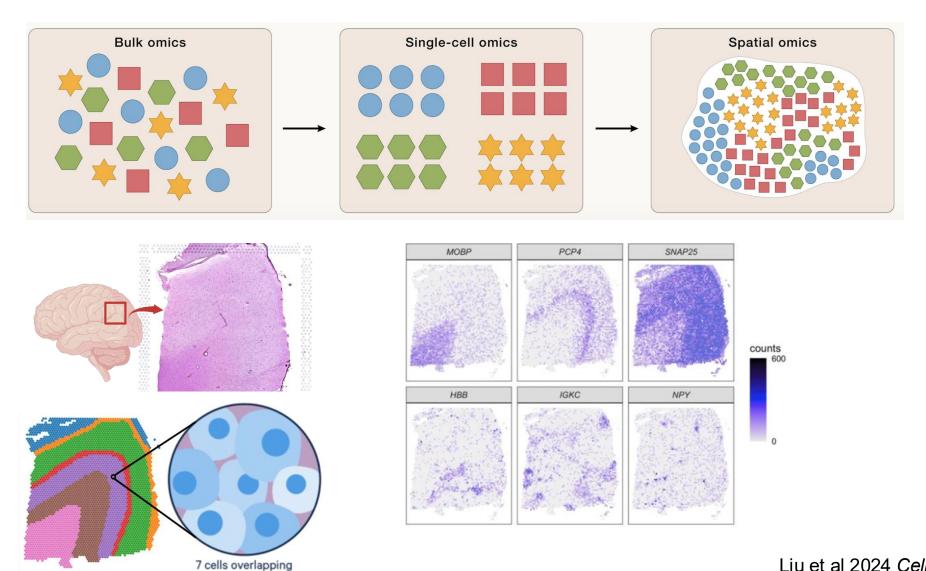


# Linking spatial transcriptomics to GWAS disease risk

# Spatial transcriptomics assays demonstrate the structure of cellular organization in a tissue



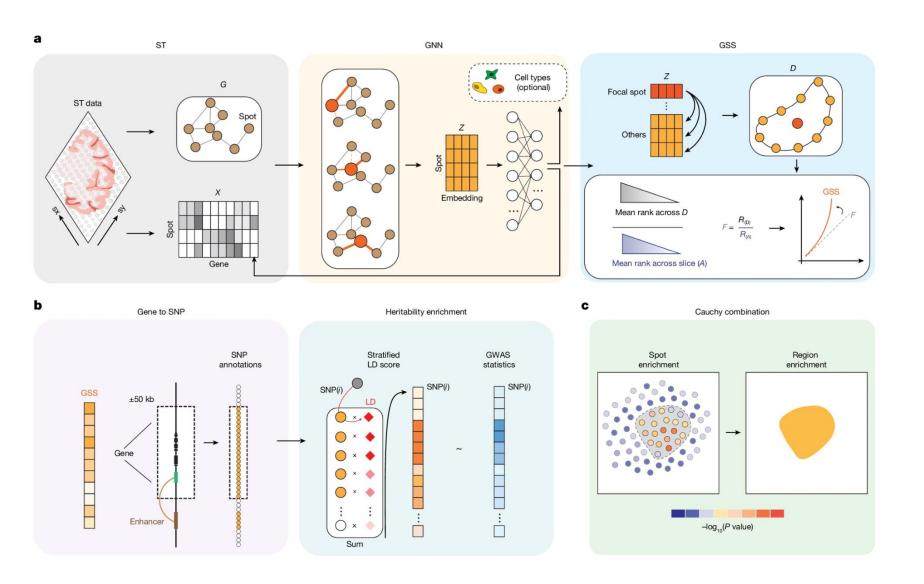
# Spatial transcriptomics assays demonstrate the structure of cellular organization in a tissue



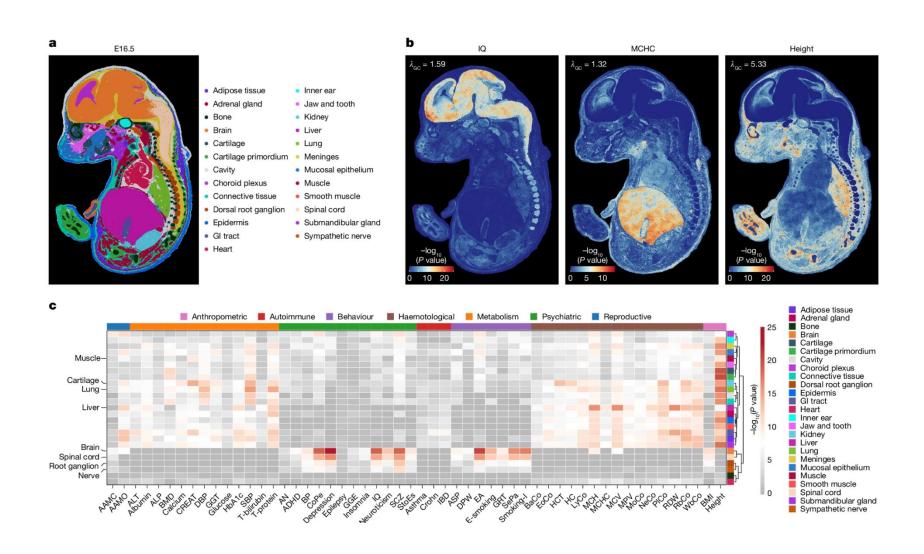
a single spot

Liu et al 2024 *Cell* Weber et al 2023 *Nat Commun* 

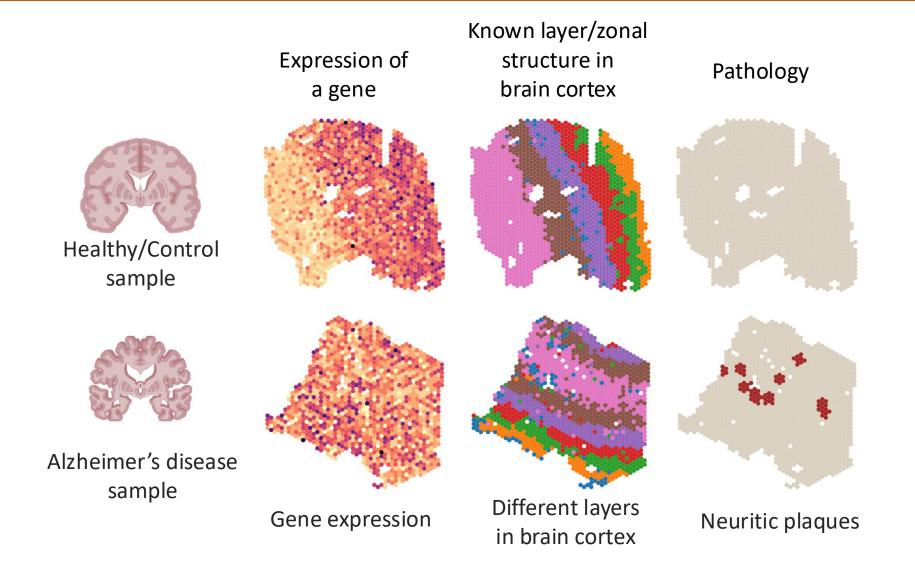
### gsMap: Scoring spatial cell-resolution disease maps



### gsMap: Scoring spatial cell-resolution disease maps



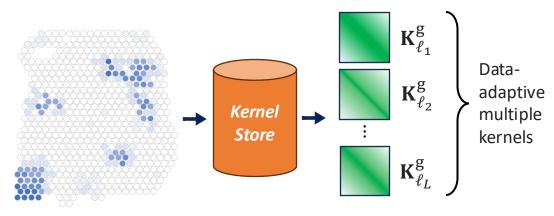
# Spatial variability patterns of a gene changes between healthy and disease states of a tissue



**Spacelink**: Lee...Dey 2025 bioRxiv, in rev Nat Commun

# Multi-scale spatial modeling of gene expression patterns can be modeled using data-driven kernels

#### Modeling of gene expression



$$\mathbf{K}_{l}^{\mathrm{g}}$$
: Exponential kernel, i.e.,  $\left[\mathbf{K}_{l}^{\mathrm{g}}\right]_{\mathrm{a,b}} = \exp\left(-\frac{\mathrm{distance\ between\ spot\ a\ and\ b}}{\mathrm{length\ scale\ }l}\right)$ 

$$\begin{array}{lll} \text{Model: } \mathbf{y}_{g} = \underbrace{\boldsymbol{\mu}_{g}} + \underbrace{\boldsymbol{U}_{1}^{g} + \cdots + \boldsymbol{U}_{L}^{g}}_{1} + \underbrace{\boldsymbol{\epsilon}_{g};}_{1} & \underbrace{\boldsymbol{U}_{l}^{g} \sim \textit{N}(\mathbf{0}, \sigma_{l}^{g} \mathbf{K}_{l}^{g})}_{\mathbf{c}_{g}} \\ \text{Fixed Spatial Nugget} & \text{mean random effect} \\ & \text{effects} & \\ \end{array}$$

Estimate  $\sigma_l^g$  using the method of moments estimator of covariance.

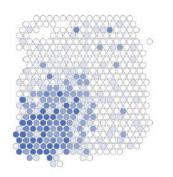
Identify SVG by testing  $H_0$ :  $\sigma_l^g = 0$  for all l  $H_a$ :  $\sigma_l^g > 0$  for some l

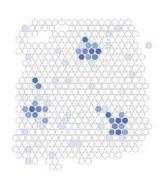
### Defining a score per gene that characterizes the effective spatial variability of a gene in a tissue

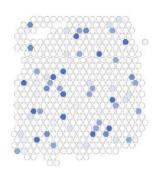
$$\begin{aligned} \mathbf{y}_{g} &= \mathbf{\mu}_{g} + \mathbf{U}_{1}^{g} + \cdots + \mathbf{U}_{L}^{g} + \mathbf{\epsilon}_{g}; & \mathbf{U}_{l}^{g} \sim \textit{N}(\mathbf{0}, \sigma_{l}^{g} \mathbf{K}_{l}^{g}) \\ & \text{Fixed} & \text{Spatial} & \text{Nugget} \\ & \text{mean} & \text{random} & \text{effect} \\ & & \text{effects} \end{aligned}$$

$$\rightarrow \text{ESV} = \frac{\sum_{l=1}^{L} w_l \sigma_l^g}{\sum_{l=1}^{L} \sigma_l^g + \tau^g}$$

$$\mathbf{w}_{l} = \frac{\left\|\mathbf{K}_{l}^{g} - \mathbf{I}\right\|_{F}}{\left\|\mathbf{K}_{l}^{g}\right\|_{F}} \in [0,1]$$





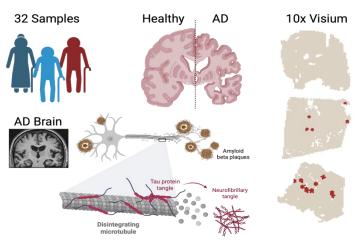


Global SVG with high ESV

Global SVG with low ESV

Not global SVG

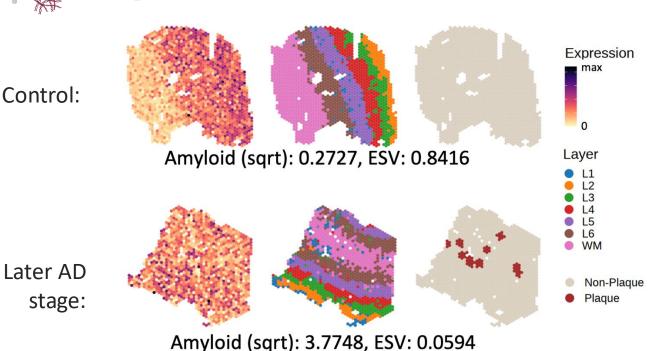
# Evaluating the performance of ESV on Alzheimer's spatial data along different pathology stages



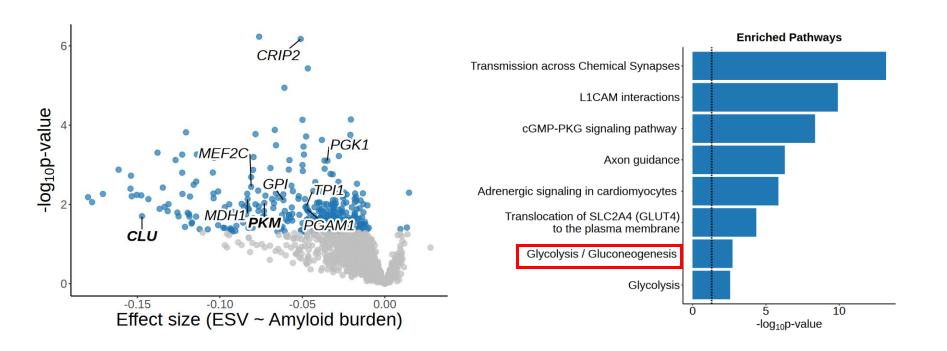
**PKM** 

#### Pathological measures:

- Total amyloid plaque burden
- Neurofibrillary tangles accumulation
- Number of neuritic plaques

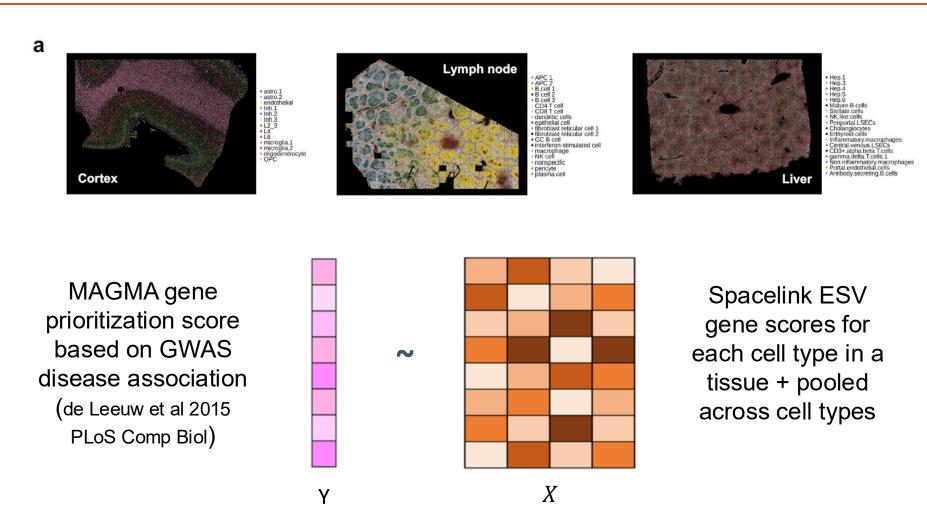


# We see higher numbers of genes showing trends of ESV change against pathology compared to other metrics



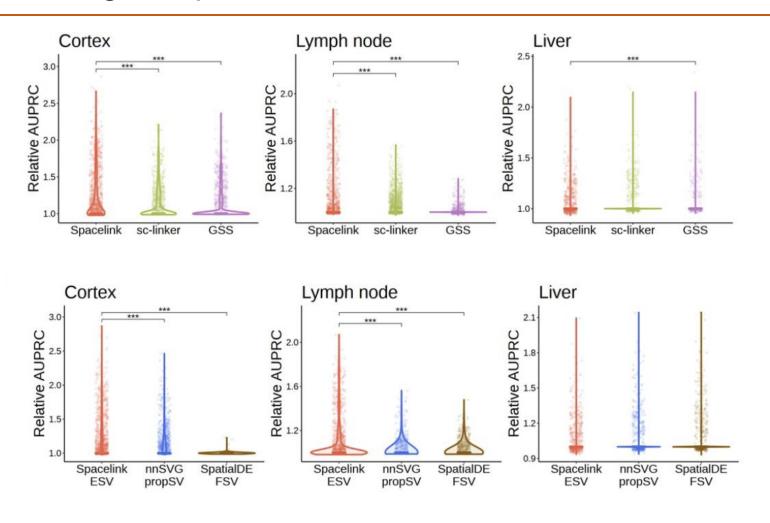
- **334** genes showed nominally significant ESV trends against AD disease pathology,
- **1.2-3.1x** higher than other spatial gene prioritization scores.

## Assessing predictive power for tissue-associated disease-related genes using Spacelink ESV features



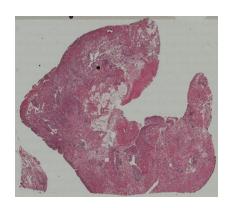
We perform this regression by leave-one-chromosome-out and predicting on the held-out chromosome to get a Spacelink integrated gene score fine-tuned for a disease.

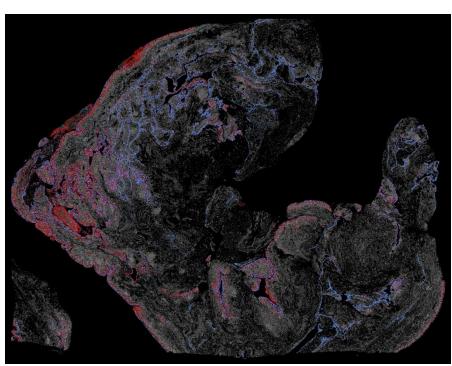
# Spacelink ESV scores across different cell types outperform other spatial and non-spatial scores in disease gene prediction task



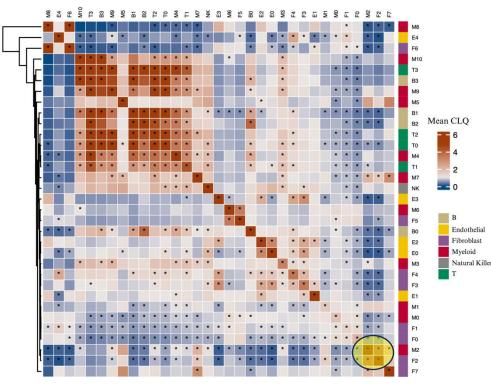
GSS: Song et al 2025 *Nature* sc-linker: Jagadeesh\*, Dey\* et al 2022 *Nat Genet* 

# Spatial co-localization of specific cell types is observed under certain disease pathological states

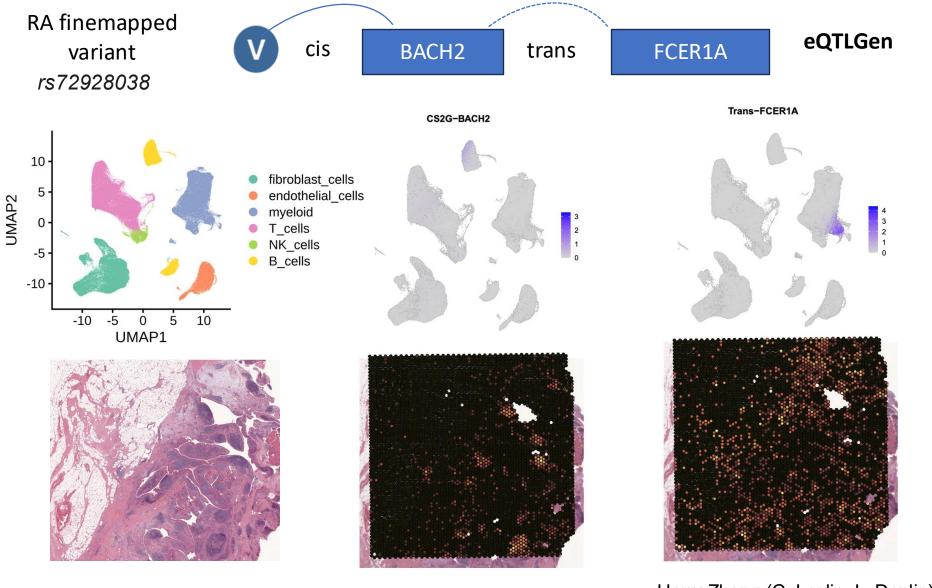




Example case of spatial colocalization of two cell types in disease pathology: SPP1-macrophages and lining fibroblasts



### Spatial patterns in gene expression can elucidate cross celltype cis-trans mechanisms downstream of genetic variation



Harry Zhang (C. Leslie, L. Donlin)

### Tip of the iceberg to the full iceberg

