

# Leveraging GTEx to Unveil the Biology of Complex Traits

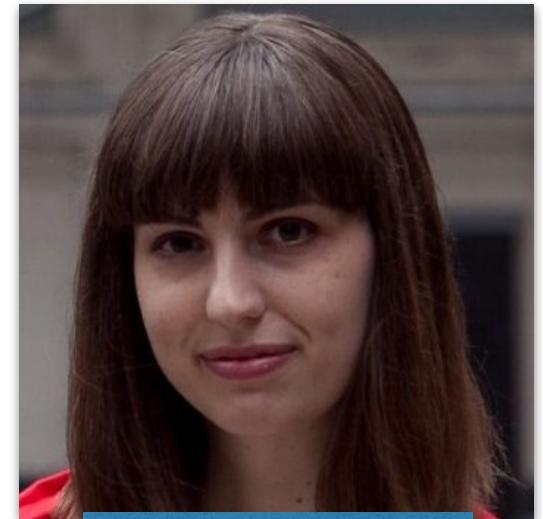
Hae Kyung Im, PhD



THE UNIVERSITY OF  
**CHICAGO**

October 5, 2015

# Heritability of Gene Expression Traits



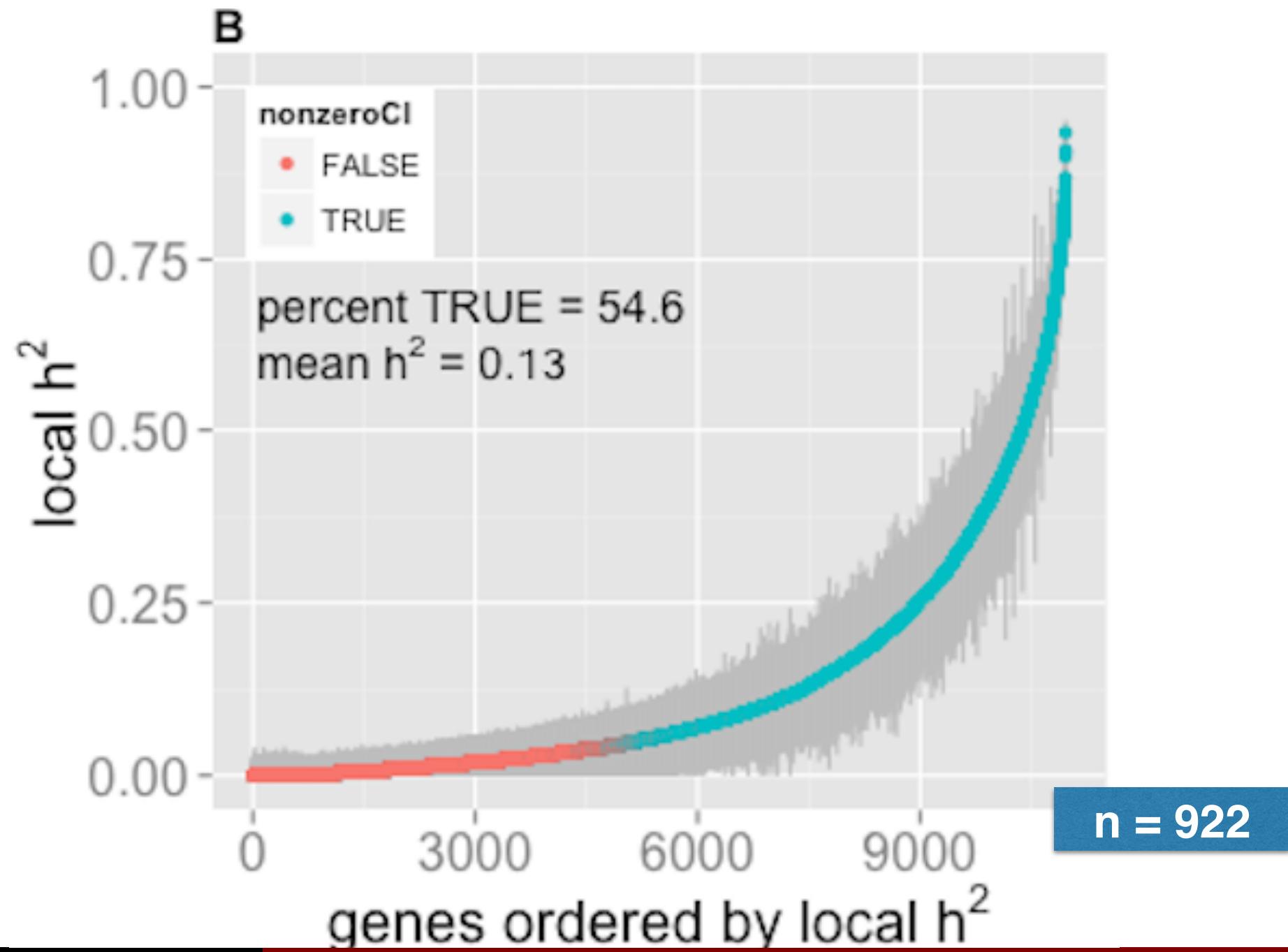
Heather Wheeler

# Local + Distant Heritability Estimation

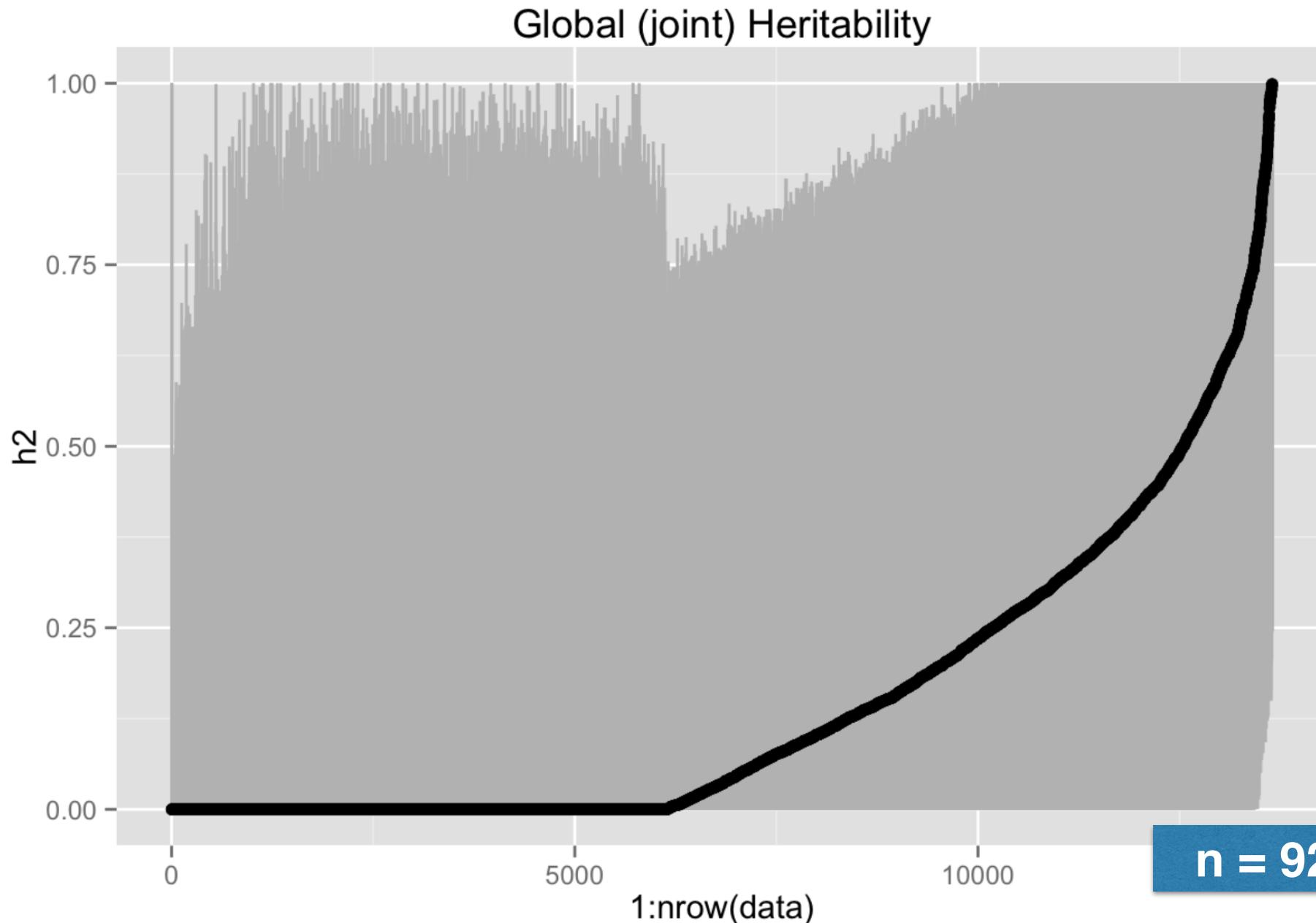
$$Y = \sum_{\text{local}} \beta_k^{\text{local}} X_k + \sum_{\text{distant}} \beta_k^{\text{distant}} X_k + \epsilon$$

**Total Heritability = Local H<sub>2</sub> + Distant H<sub>2</sub>**

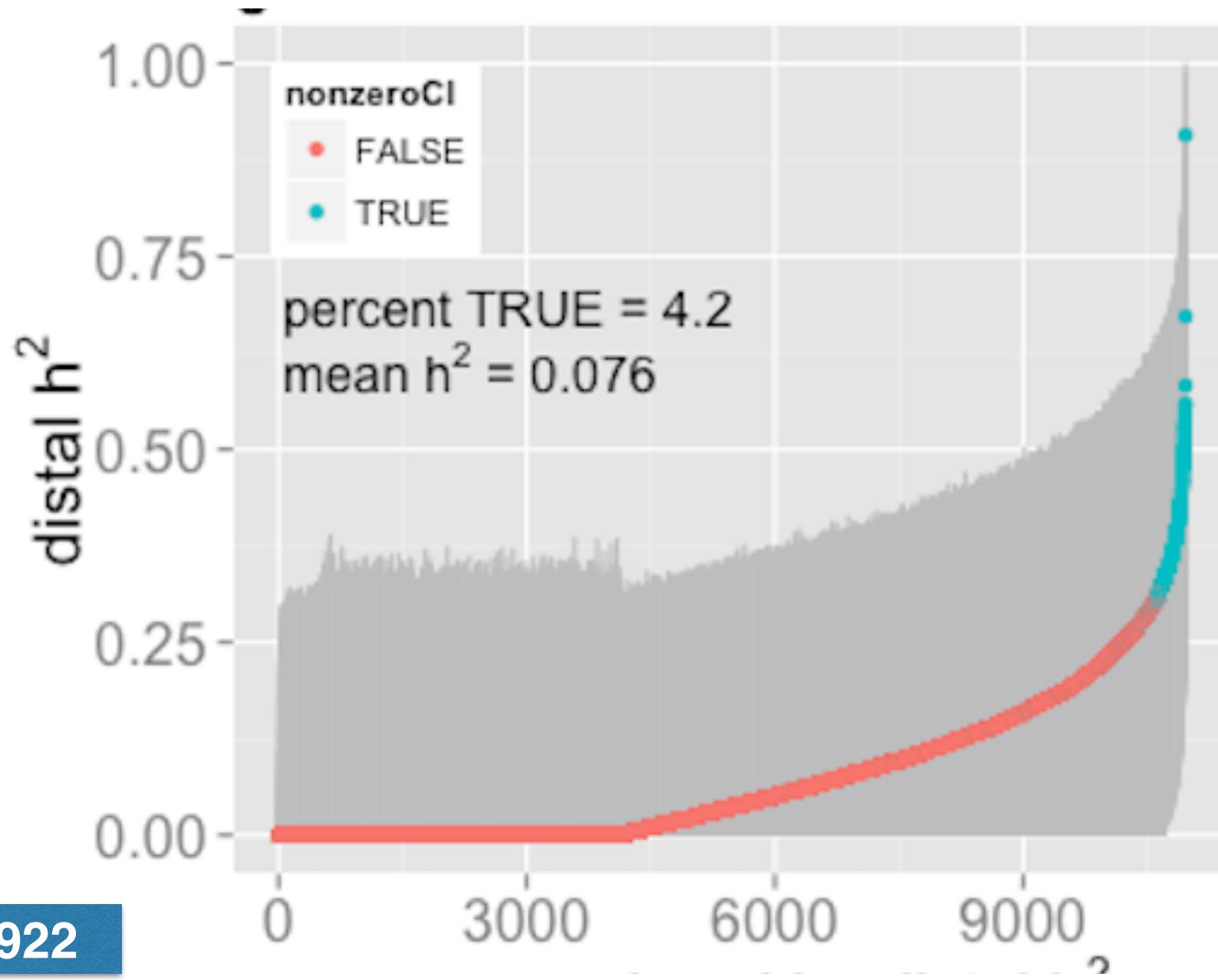
# Local Heritability Well Estimated



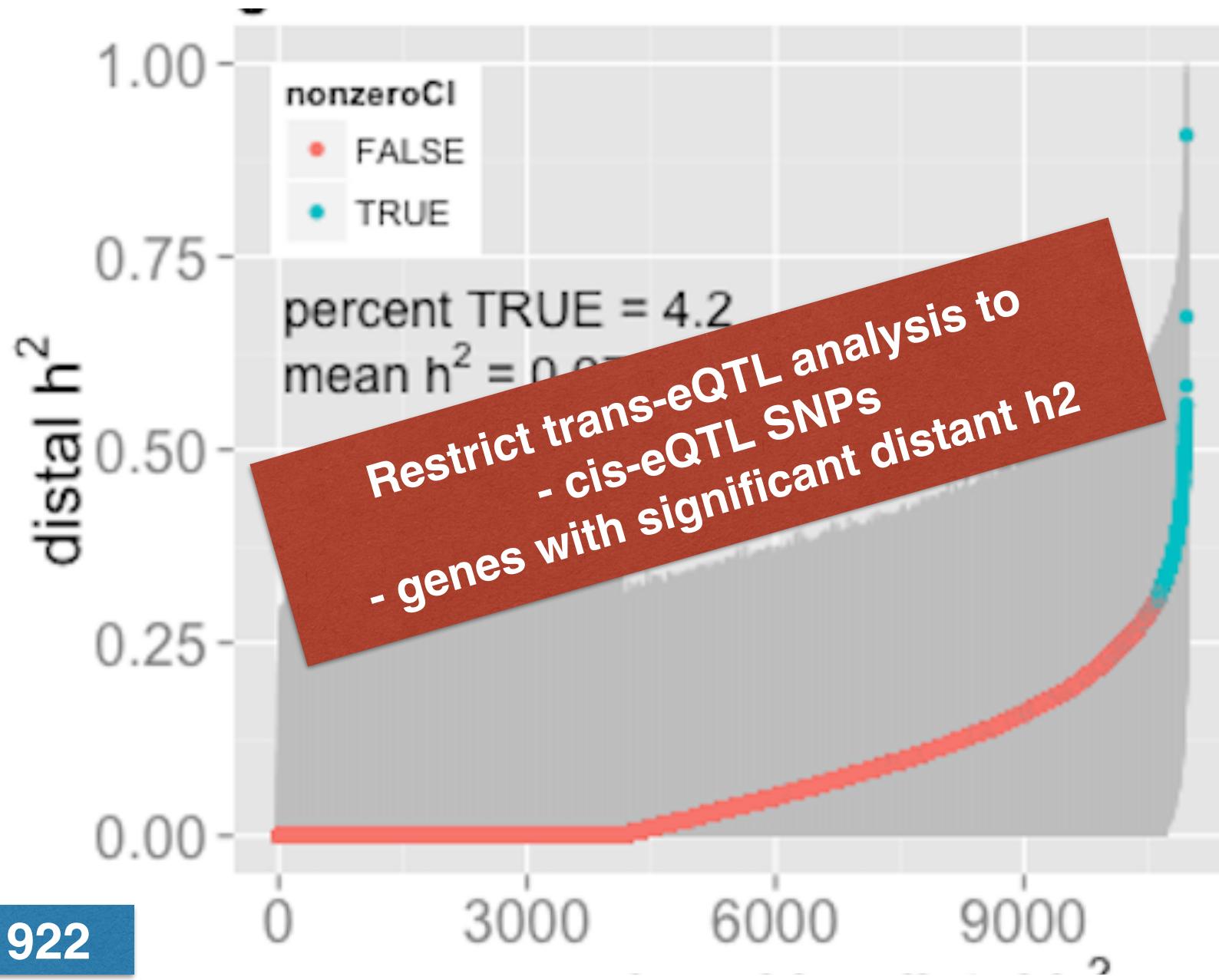
# Distant h<sup>2</sup> Not Well Estimated (n=922)



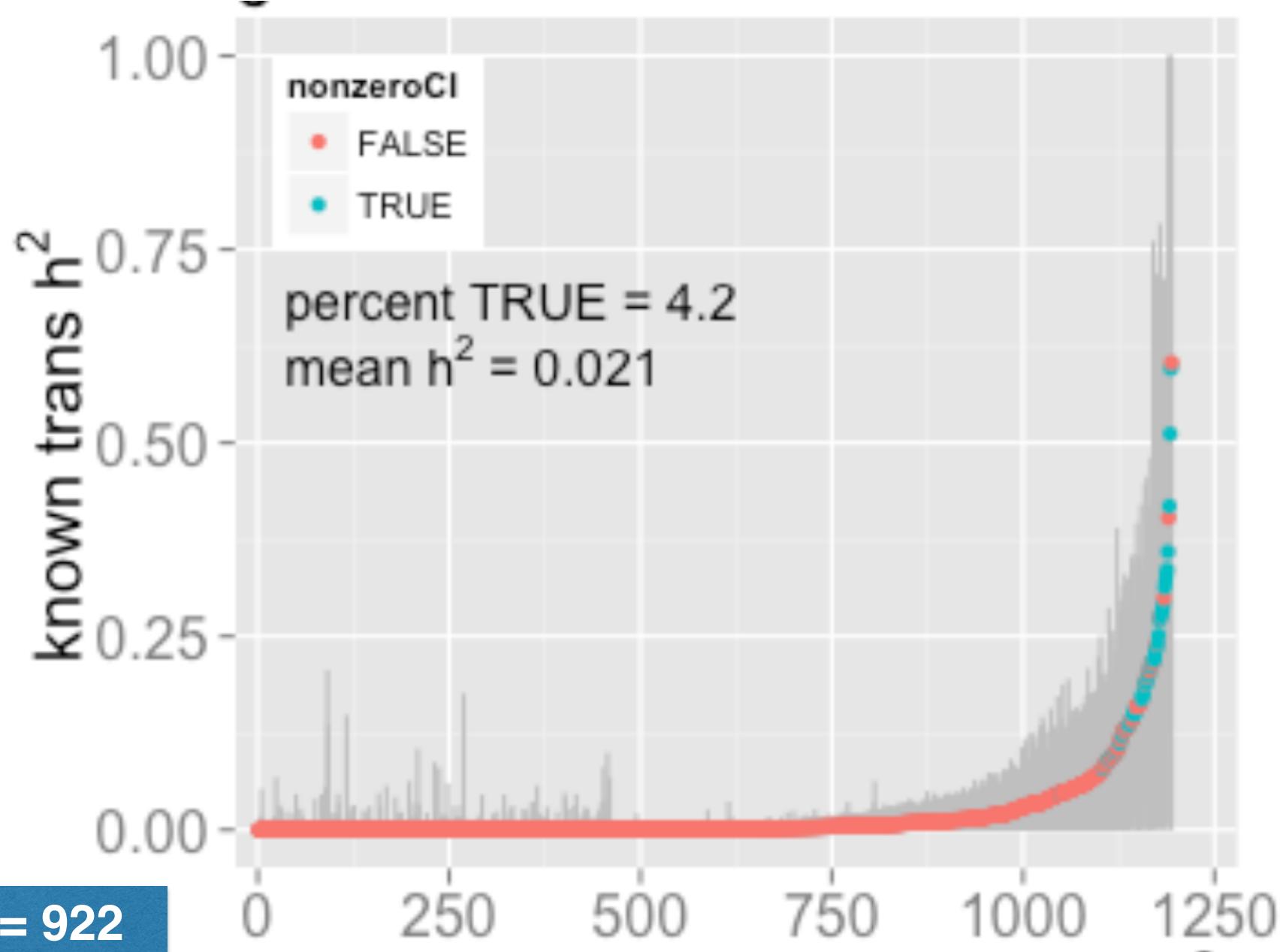
# cis-eQTL Priors (FHS) Improve Distant h<sup>2</sup> Estimates



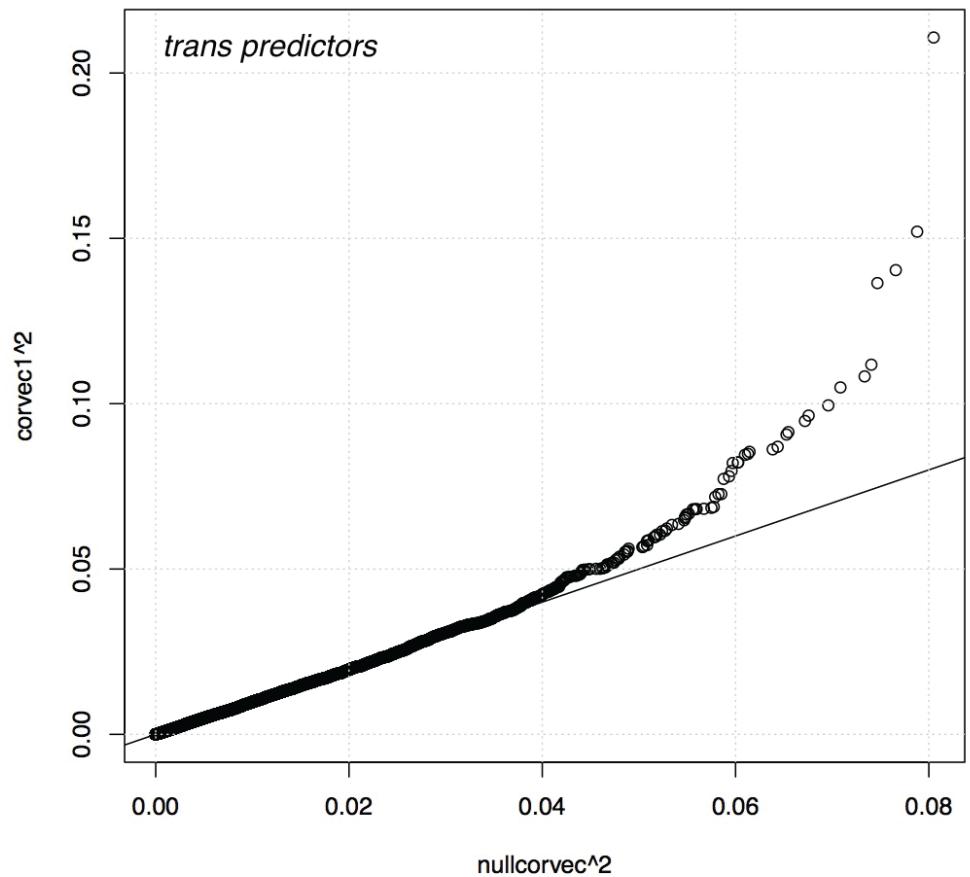
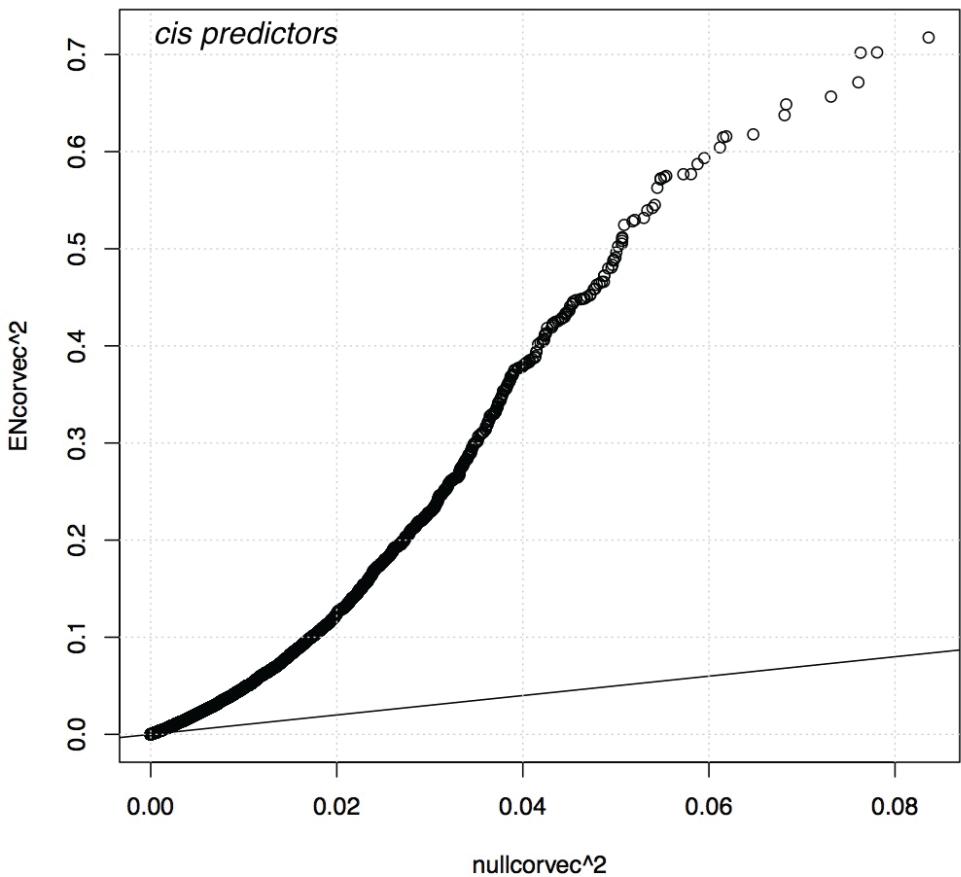
# cis-eQTL Priors (FHS) Improve Distant $h^2$ Estimates



# Distant eQTL Prior for Specific Genes Improves $h^2$



# Performance of Local vs. Distant Variants



# Sparsity of Gene Expression Traits

# Estimating Sparse & Polygenic Components with BSLMM

OPEN  ACCESS Freely available online



## Polygenic Modeling with Bayesian Sparse Linear Mixed Models

Xiang Zhou<sup>1\*</sup>, Peter Carbonetto<sup>1</sup>, Matthew Stephens<sup>1,2\*</sup>

**1** Department of Human Genetics, University of Chicago, Chicago, Illinois, United States of America, **2** Department of Statistics, University of Chicago, Chicago, Illinois, United States of America

### Abstract

Both linear mixed models (LMMs) and sparse regression models are widely used in genome-wide association studies. These two approaches make different assumptions about the underlying genetic architecture.

$$\mathbf{y} = \mathbf{1}_n \boldsymbol{\mu} + \mathbf{X} \tilde{\boldsymbol{\beta}} + \mathbf{u} + \boldsymbol{\varepsilon},$$

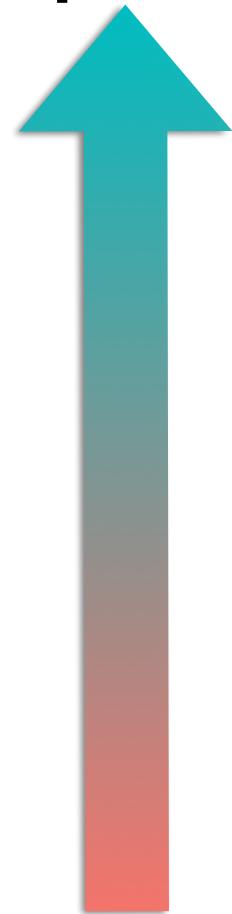
$$\mathbf{u} \sim \text{MVN}_n(0, \sigma_b^2 \tau^{-1} \mathbf{K}),$$

$$\boldsymbol{\varepsilon} \sim \text{MVN}_n(0, \tau^{-1} \mathbf{I}_n),$$

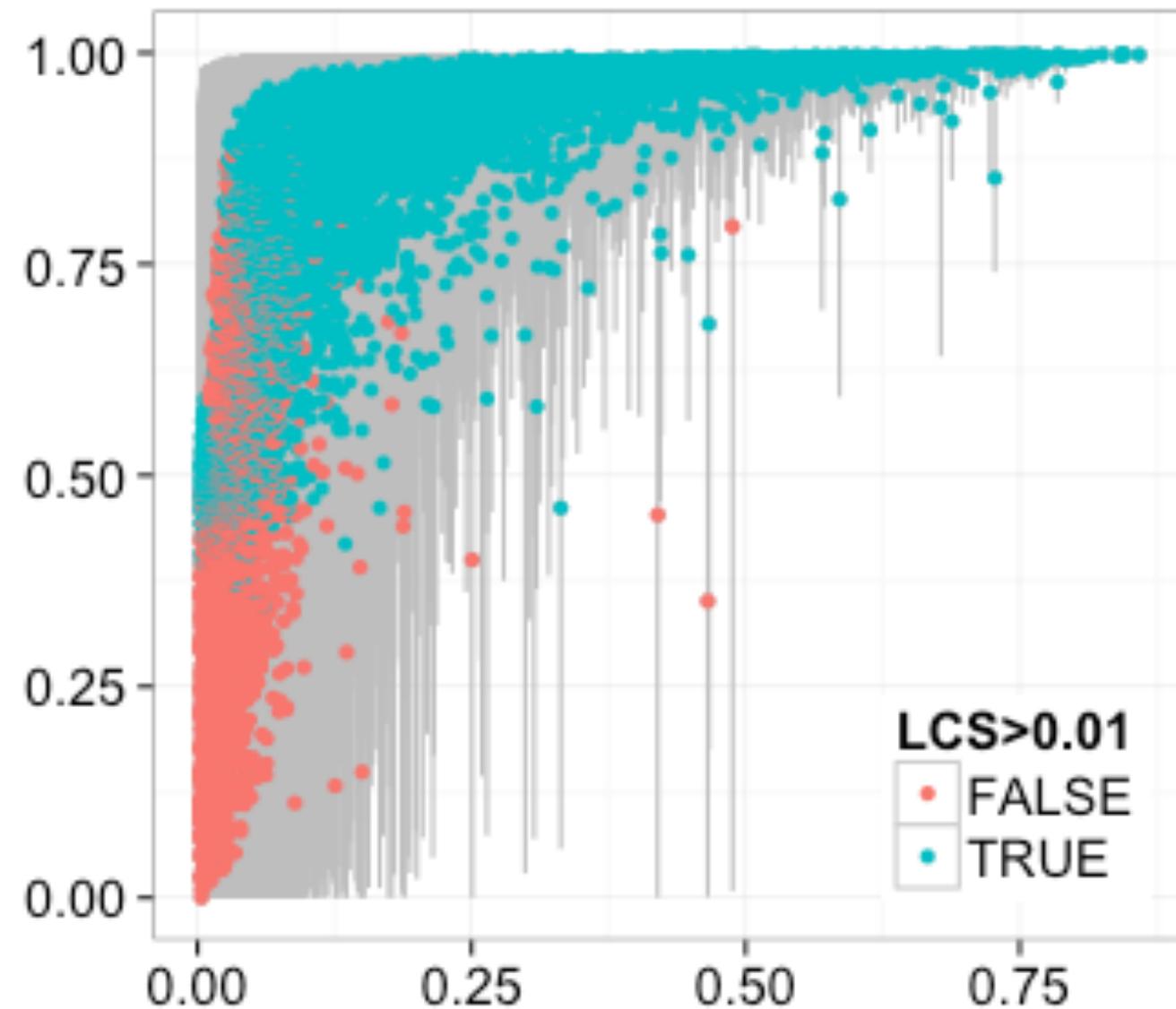
$$\tilde{\boldsymbol{\beta}}_i \sim \pi \mathcal{N}(0, \sigma_a^2 \tau^{-1}) + (1 - \pi) \delta_0,$$

# Highly Heritable Genes are Sparse

Sparse

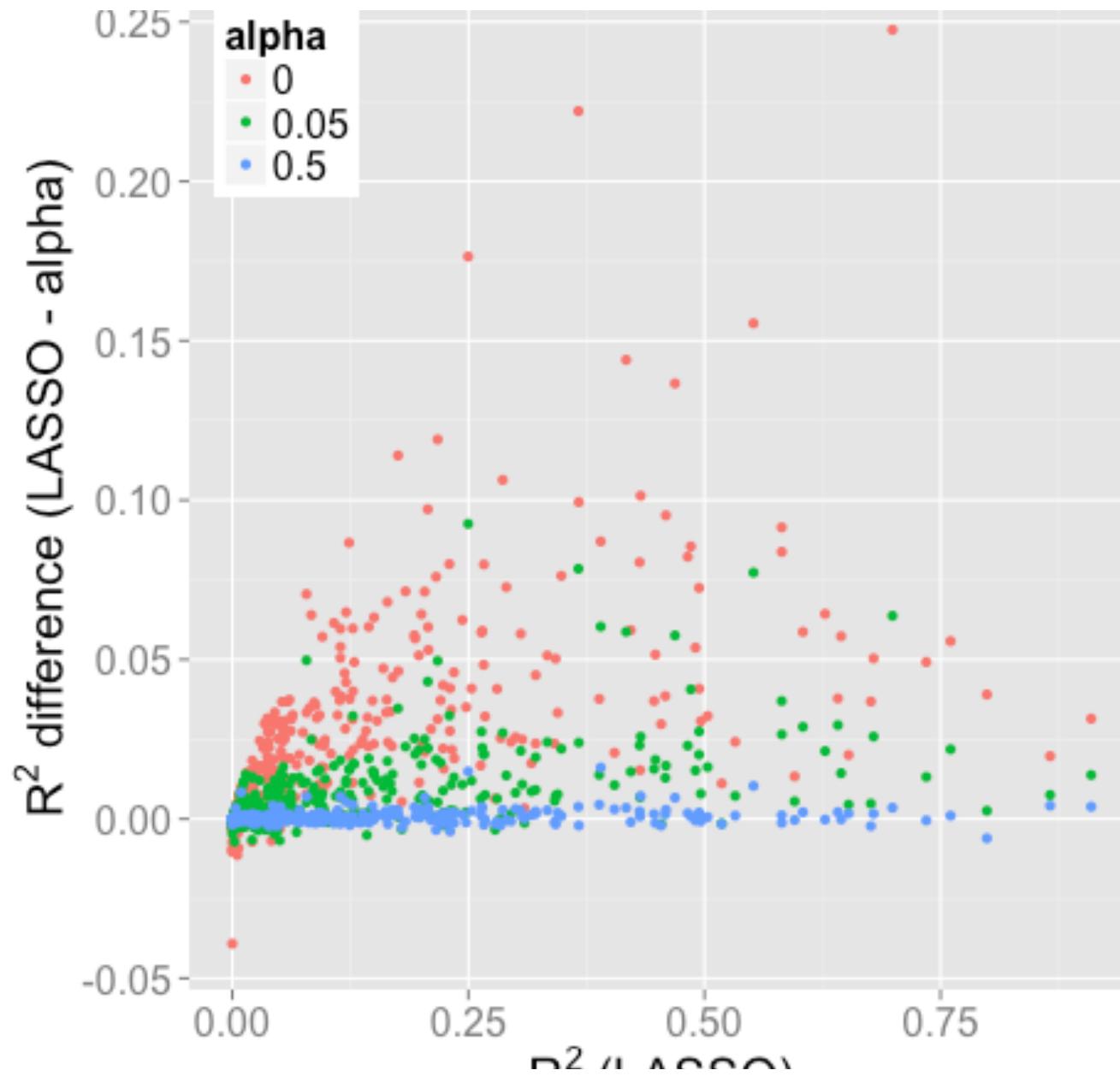


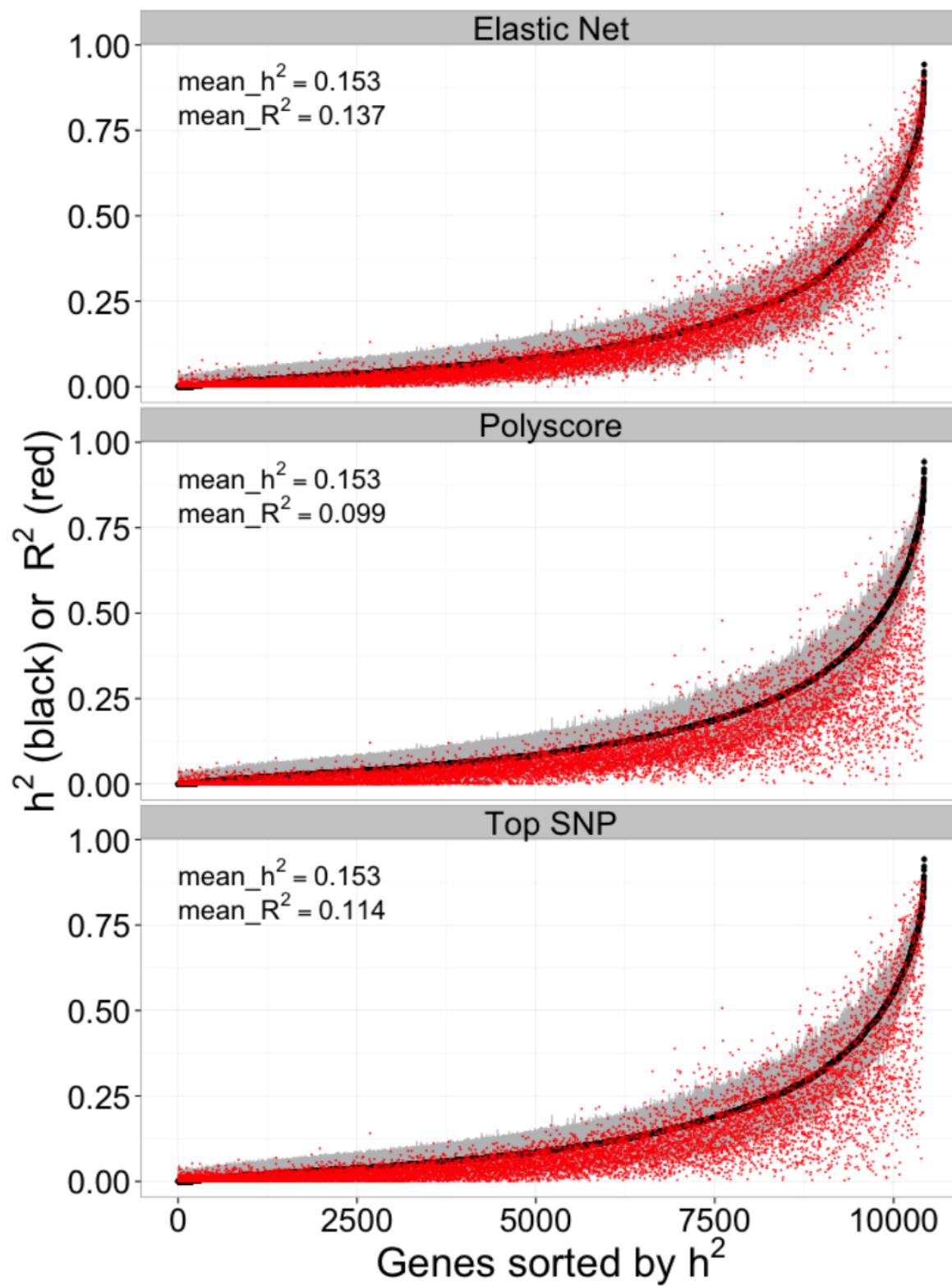
Polygenic



PVE (Local Heritability)

# Sparse Models Outperform Polygenic Models in Prediction



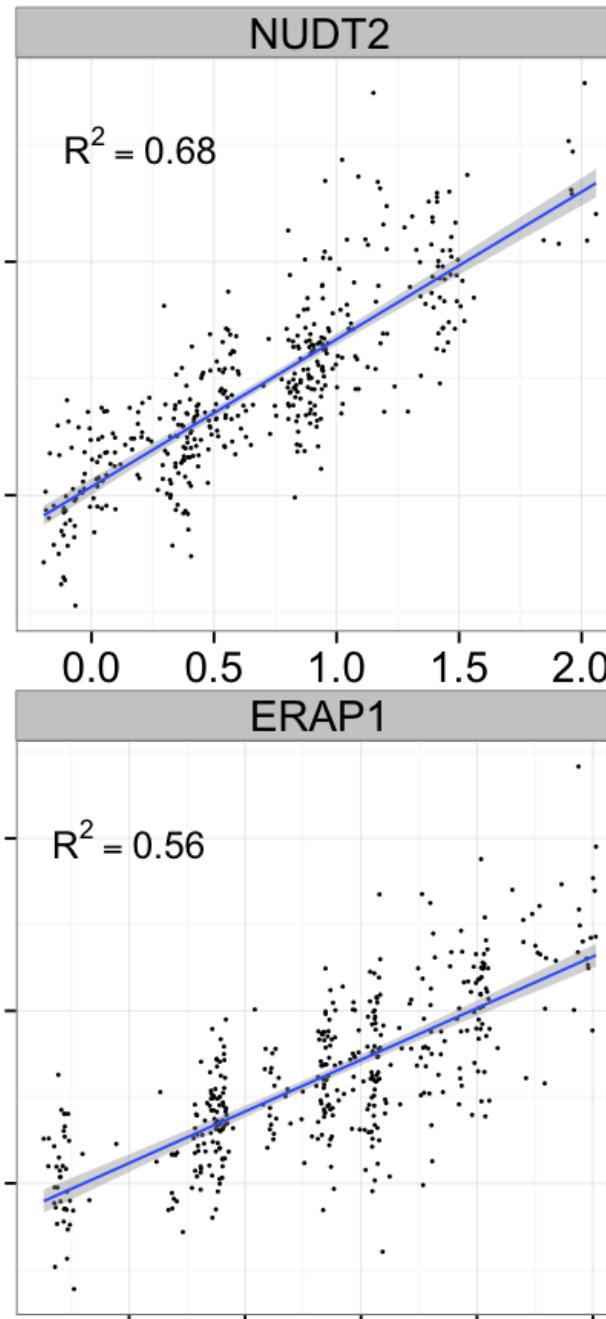
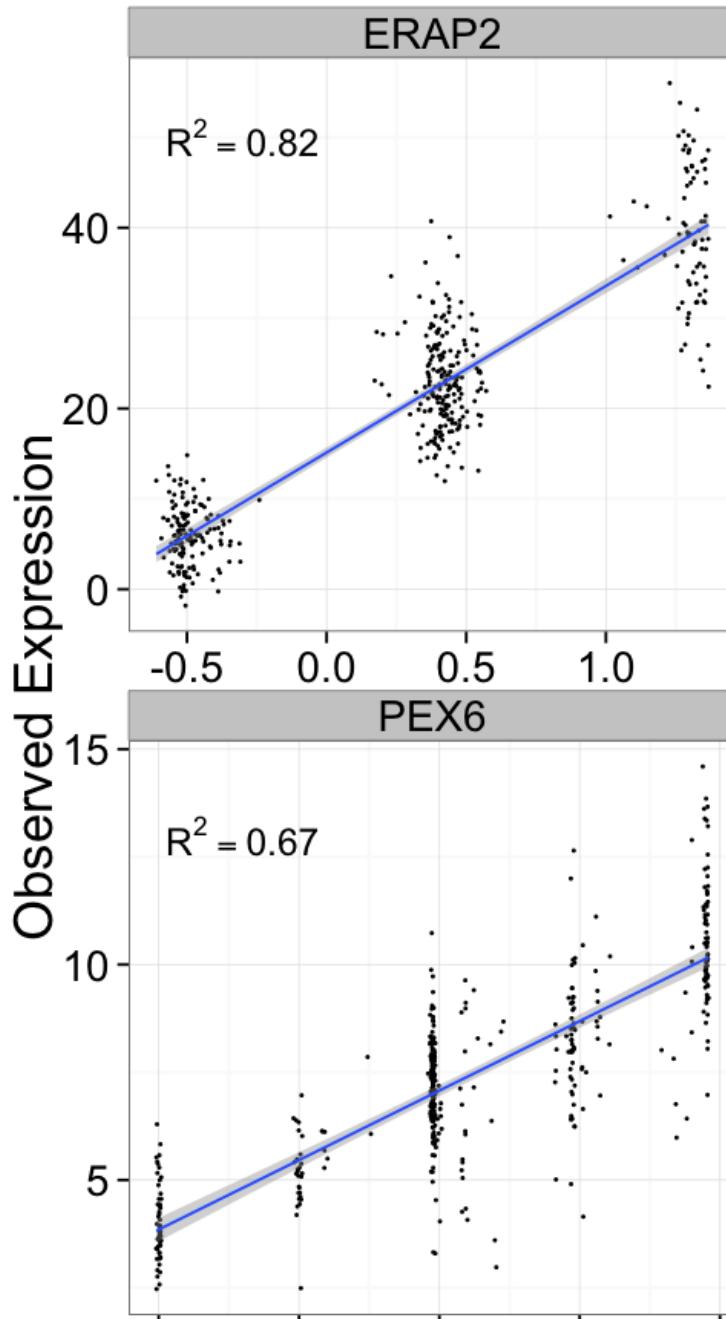


**Sparse Models Outperform Polygenic Models**

**Polygenic architecture not supported**

**Using top eQTL suboptimal**

# Examples of well predicted genes

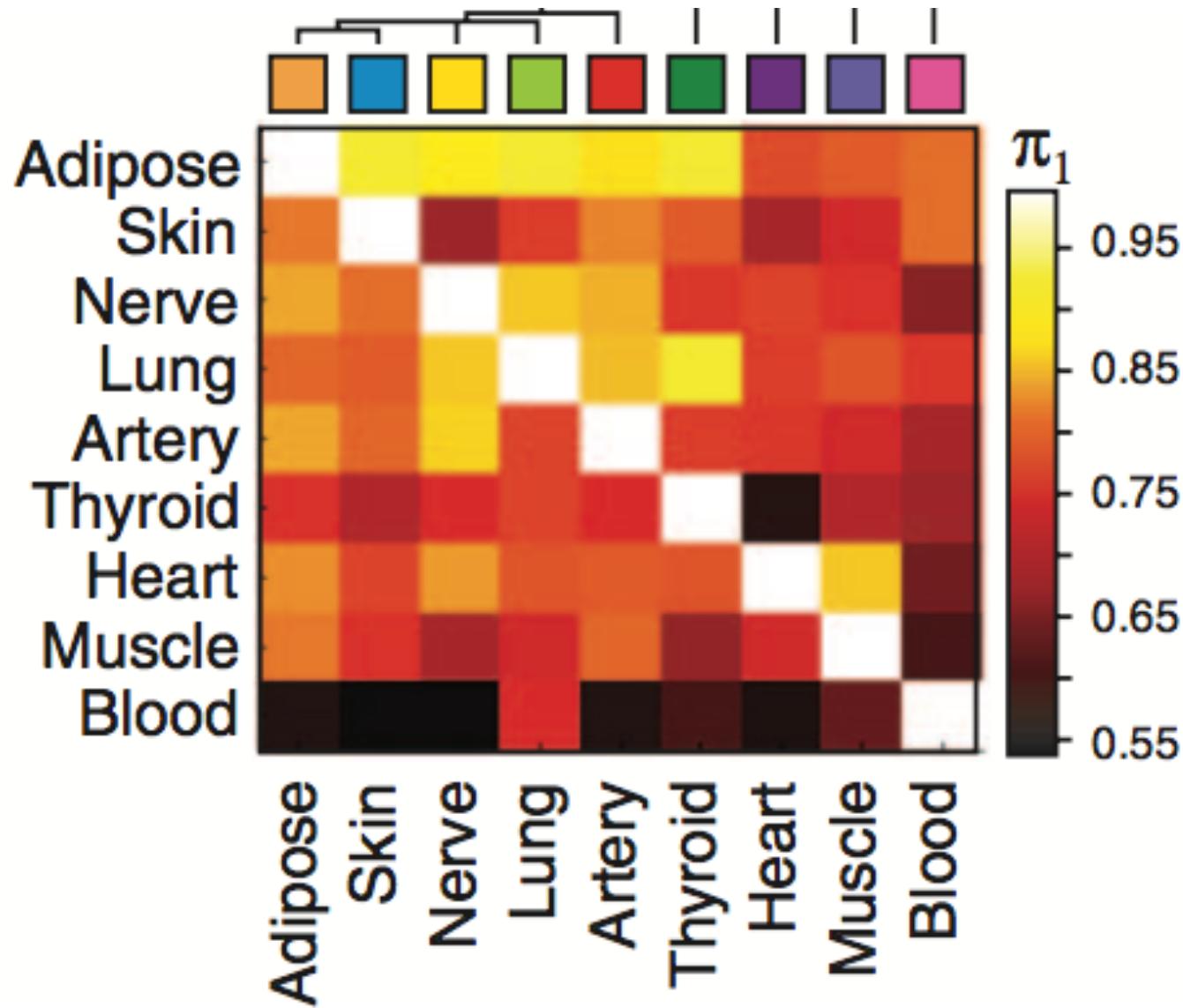


ERAP2: one causal variant

PEX6: two causal variants

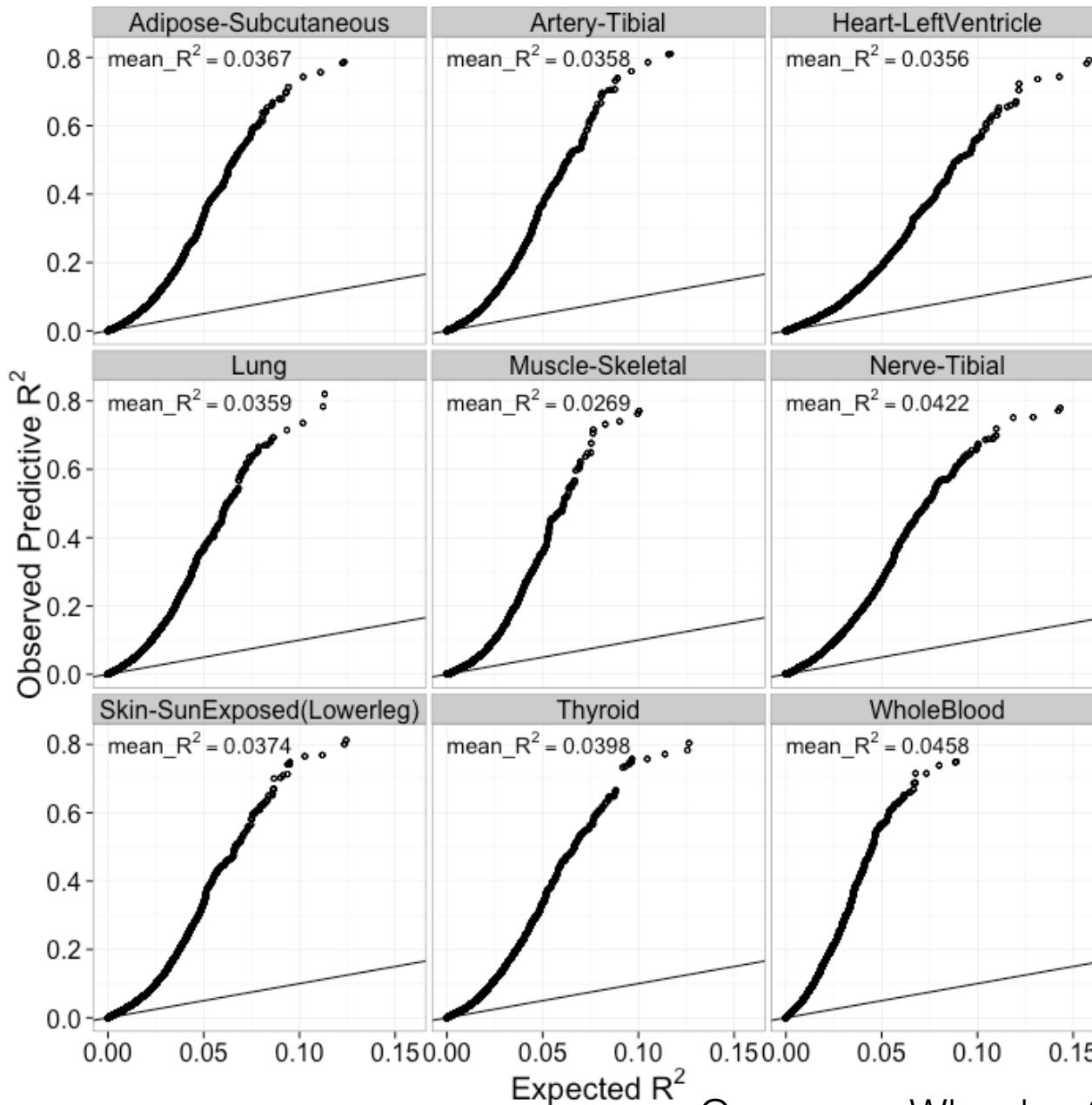
# Cross Tissue and Tissue Specific Architecture

# eQTLs Are Highly Shared Across Tissues



GTEX Consortium, Science 2015

# Predicted Expression Performs Well Across Tissues



Gamazon, Wheeler, Shah et al NG 2015

# Splitting Tissue Specific and Cross Tissue Components

Adipose						
id	g1	g2	g3	...	g20K	
id1	0.1	0.1	0.2		3.2	
id2	1.1	3.1	1.2		4.3	
id3	1.2	2.0	2.1		2.1	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
idn	1.2	2.2	3.1		2.1	

Brain						
id	g1	g2	g3	...	g20K	
id1	0.1	0.1	0.2		3.2	
id2	2.2	1.7	1.2		4.1	
id3	1.3	2.0	1.7		2.1	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
idn	1.2	2.2	3.1		2.1	

WB						
id	g1	g2	g3	...	g20K	
id1	0.1	0.1	0.2		3.2	
id2	2.2	1.7	1.2		4.1	
id3	1.3	2.0	1.7		2.1	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
idn	1.2	2.2	3.1		2.1	

=

Cross Tissue						
id	g1	g2	g3	...	g20K	
id1	0.1	0.1	0.2		3.2	
id2	2.2	1.7	1.2		4.1	
id3	1.3	2.0	1.7		2.1	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
idn	1.2	2.2	3.1		2.1	

+

Adipose						
id	g1	g2	g3	...	g20K	
id1	0.1	0.1	0.2		3.2	
id2	1.1	3.1	1.2		4.3	
id3	1.2	2.0	2.1		2.1	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
idn	1.2	2.2	3.1		2.1	

Brain						
id	g1	g2	g3	...	g20K	
id1	0.1	1.7	1.2		4.1	
id2	2.2	1.7	1.2		4.1	
id3	1.3	2.0	1.7		2.1	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
idn	1.2	2.2	3.1		2.1	

WB						
id	g1	g2	g3	...	g20K	
id1	0.1	0.1	0.2		3.2	
id2	2.2	1.7	1.2		4.1	
id3	1.3	2.0	1.7		2.1	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
:	:	:	:	:	:	
idn	1.2	2.2	3.1		2.1	

$$Y_{it} = \mu + c_i + r_{it}$$

$$c_i \sim N(0, \Sigma_{n \times n})$$

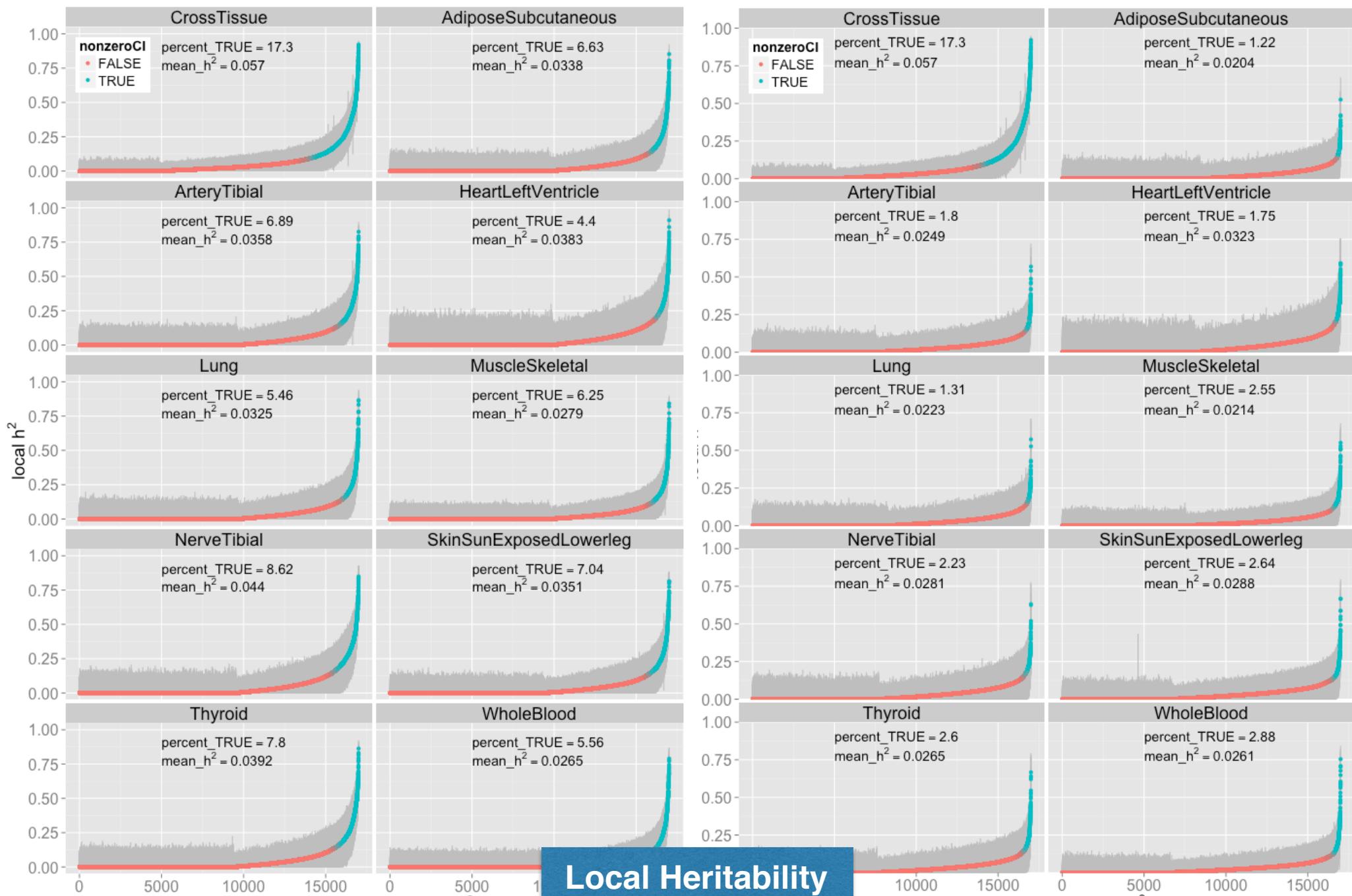
$$r_{it} \sim N(0, \Sigma_{nt \times nt})$$

# Cross-Tissue Less Noisy & Larger Sample Size

Tissue	n
<b>Cross-Tissue</b>	<b>450</b>
Muscle-Skeletal	361
WholeBlood	339
Skin-SunExposed(Lowerleg)	303
Adipose-Subcutaneous	298
Artery-Tibial	285
Lung	279
Thyroid	279
Nerve-Tibial	256
Heart-LeftVentricle	190

## Tissue Whole

## Tissue Specific



**Tissue Whole**

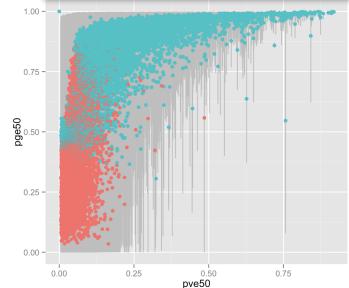
**Tissue Specific**

**Distant Heritability with FHS cis-eQTL subset**

**Tissue Whole**

**Tissue Specific**

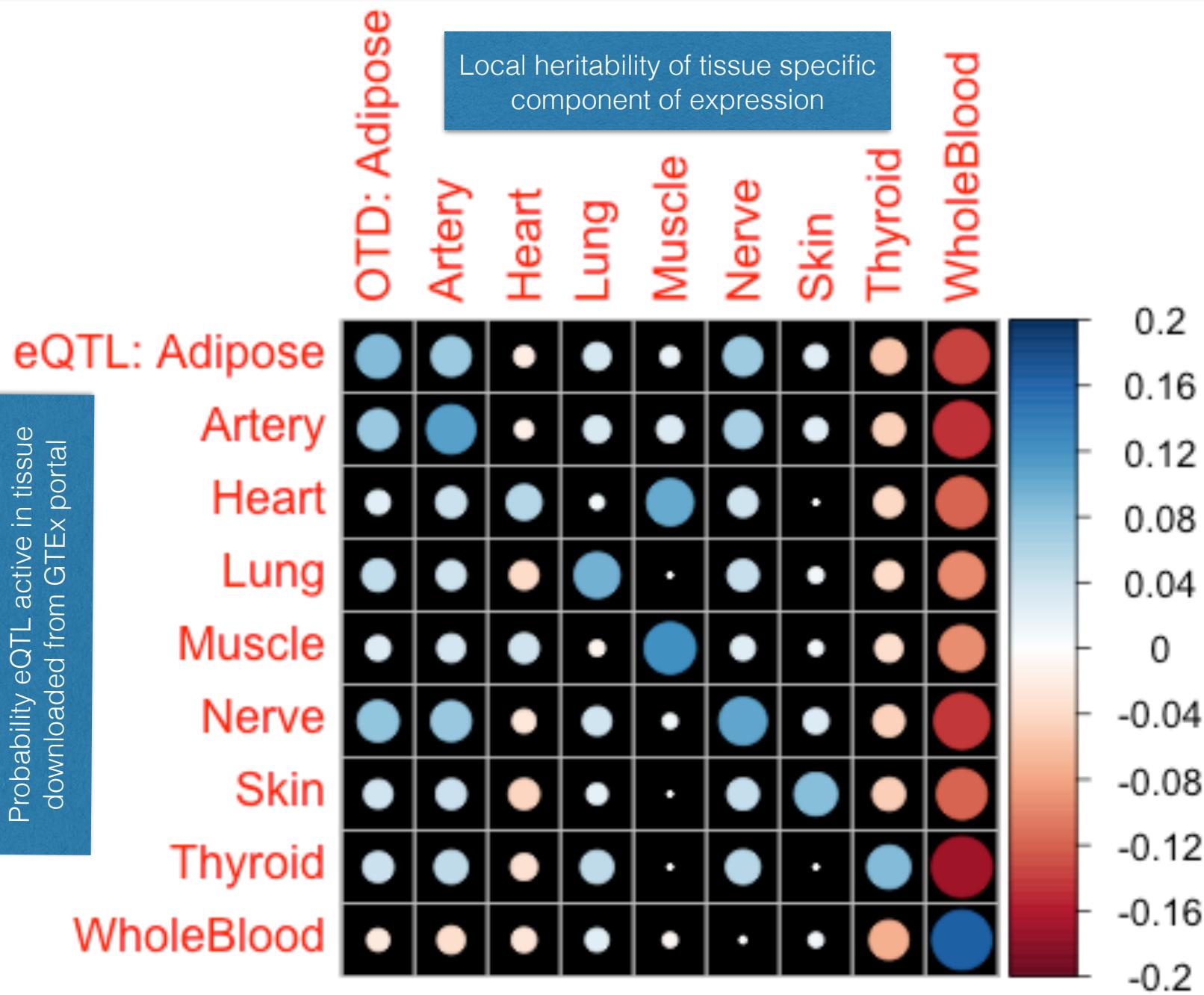
**Cross Tissue**



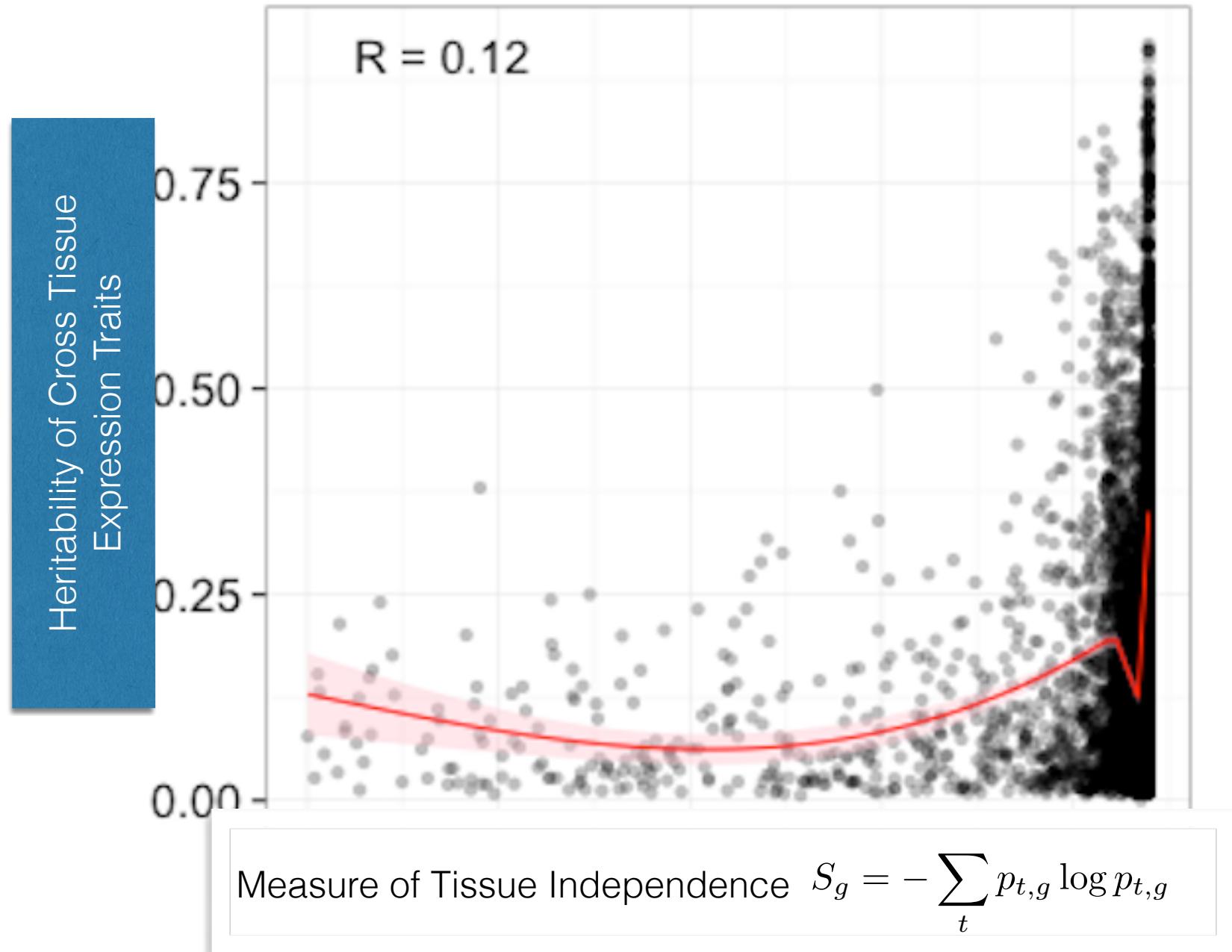
**Sparsity**

# Correlation bw Tissue Specific h<sup>2</sup> vs Prob eGene in Tissue

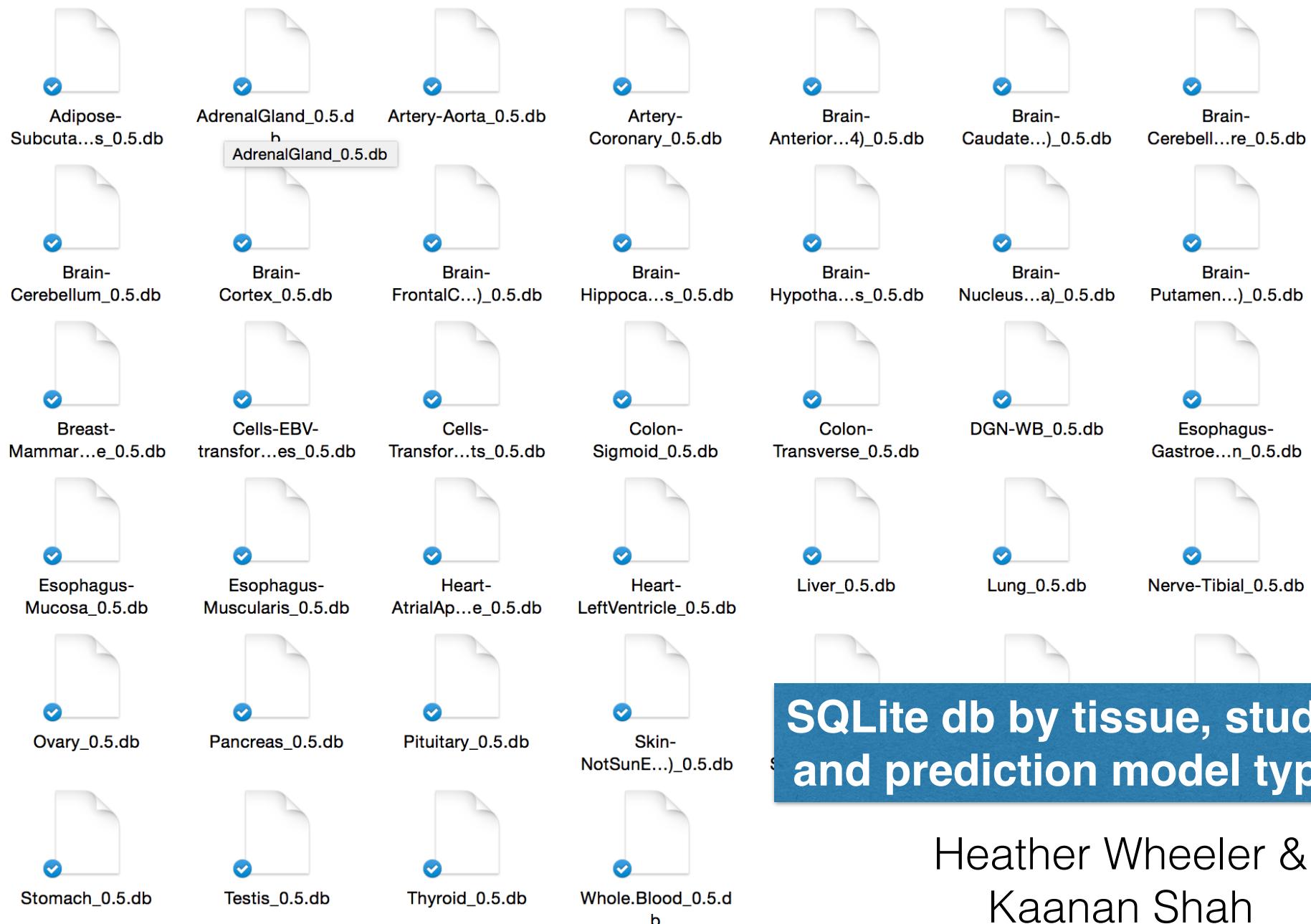
Probability eQTL active in tissue  
downloaded from GTEx portal



# Tissue Independence Correlates with Cross Tissue h<sup>2</sup>



# Whole Blood DGN (n=922) + 38 GTEx Tissue Models



Heather Wheeler &  
Kaanan Shah

# Summary

- Quantified local and distant heritability of gene expression traits
- Local heritability is well estimated
- Distant heritability can only be estimated using strong functional priors to reduce number of genetic markers
- For heritable genes architecture is sparse rather than polygenic
- Orthogonal Tissue Decomposition
  - Defined new phenotypes cross tissue & orthogonal tissue specific expressions
  - Found consistency with multi-tissue

ASHG - Invited Session #13  
Wednesday 11-1pm Room 316  
“Secure, Efficient, and Scalable Computational  
Genetics via Summary Statistics”

# MetaXcan

## a summary statistics based gene-level association test



Alvaro Barbeira

# PrediXcan Imputes Transcriptome & Tests Assoc.

PrediXcan on GWAS Data

## Genetic Variation

M SNPs

id	rs1	rs2	rs1	...	rsM
id1	0	2	1		0
id2	1	2	2		2
id3	2	1	1		1
:	:	:	:		:
:	:	:	:		:
:	:	:	:		:
:	:	:	:		:
idn'	1	2	0		2

n<sup>'</sup> individuals

## "Imputed" Transcriptome

m genes

id	g1	g2	g3	...	gm	trait
id1	0.2	0.6	0.2		3.2	0.1
id2	2.3	1.8	1.2		4.1	2.2
id3	3.3	2.2	1.7		2.1	1.3
:	:	:	:		:	:
:	:	:	:		:	:
:	:	:	:		:	:
:	:	:	:		:	:
idn'	2.2	2.0	3.1		2.1	1.2

Association  
Test

# MetaXcan Formula

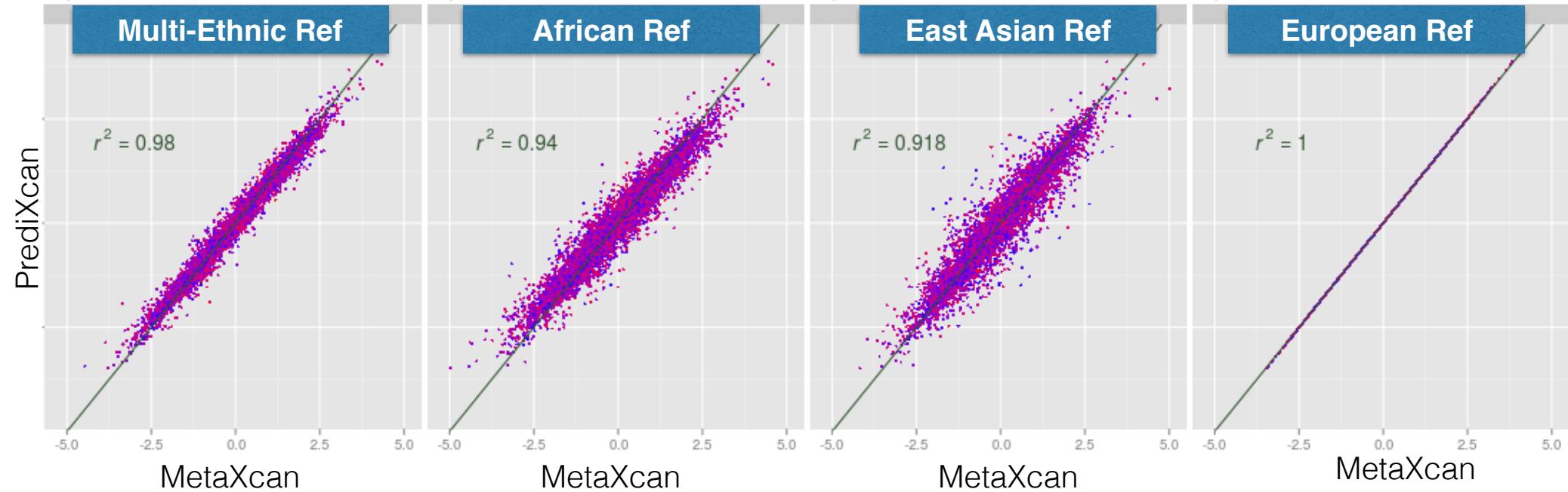
$$\begin{aligned}\hat{\gamma}_g &= \frac{\text{Cov}(T_g, Y)}{\hat{\sigma}_g^2} \\ &= \frac{\text{Cov}(\sum_{l \in \text{Model}_g} w_{lg} X_l, Y)}{\hat{\sigma}_g^2} \\ &= \sum_{l \in \text{Model}_g} \frac{w_{lg} \text{Cov}(X_l, Y)}{\hat{\sigma}_g^2} \\ &= \sum_{l \in \text{Model}_g} \frac{w_{lg} \hat{\beta}_l \sigma_l^2}{\hat{\sigma}_g^2}\end{aligned}$$

# MetaXcan Formula

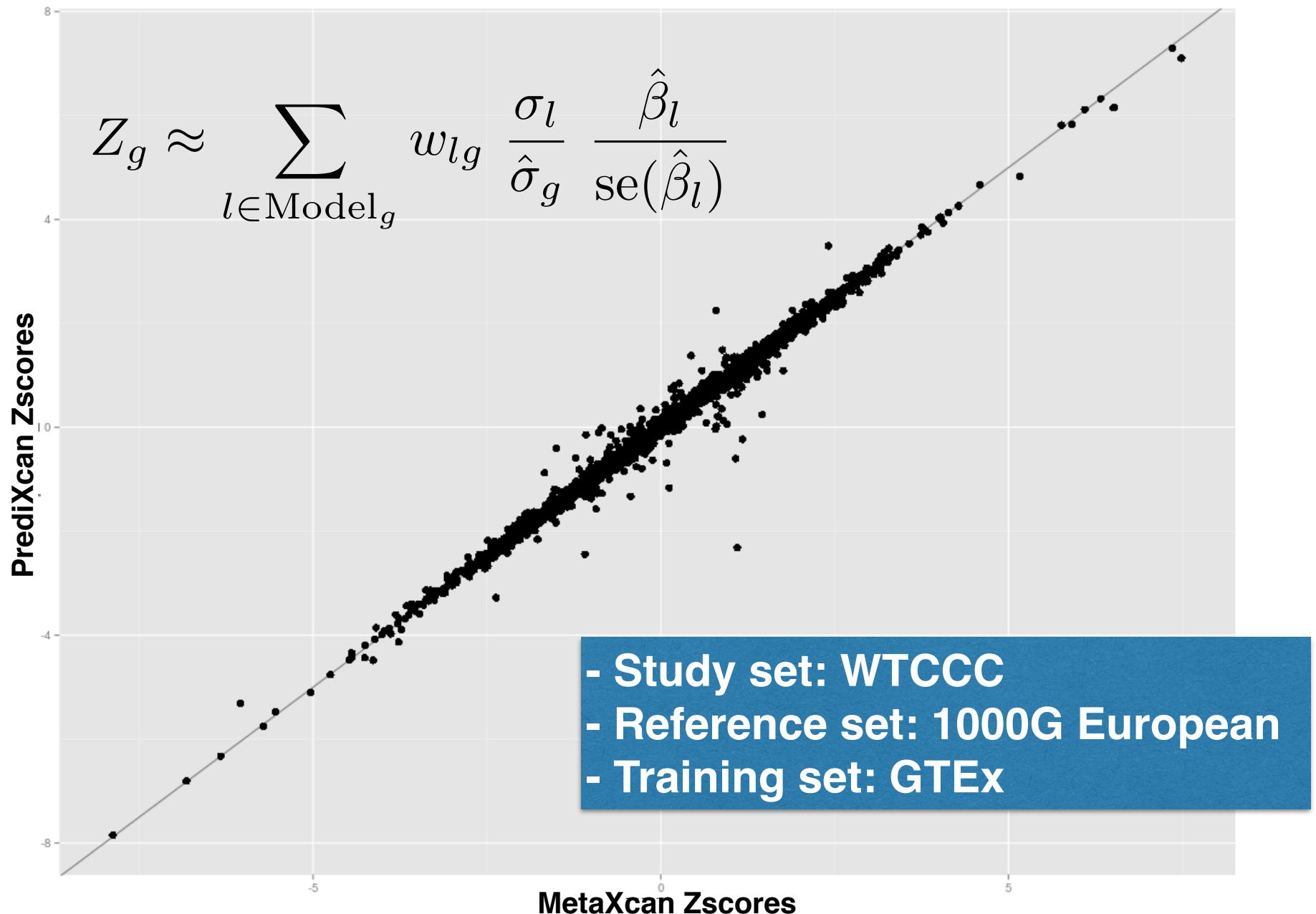
$$Z_g = \sum_{l \in \text{Model}_g} w_{lg} \frac{\sigma_l}{\hat{\sigma}_g} \frac{\hat{\beta}_l}{\text{se}(\hat{\beta}_l)} \sqrt{\frac{1 - R_l^2}{1 - R_g^2}}$$

# Robustness to Reference Population Differences

**PrediXcan vs. MetaXcan**  
**Study set: European**  
**Training set: GTEx**

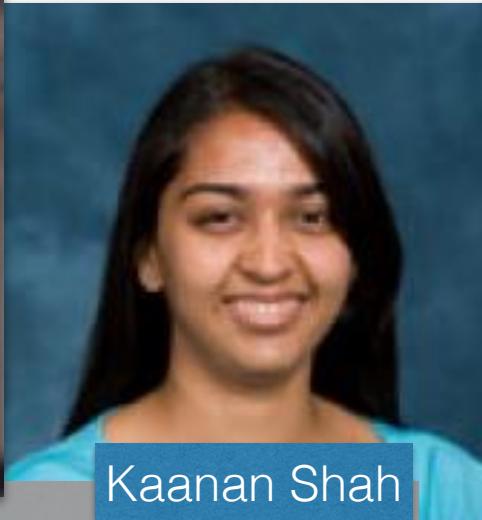


# WTCCC T1D PrediXcan vs MetaXcan





Heather Wheeler



Kaanan Shah



Sahar Mozaffari



Nancy Cox



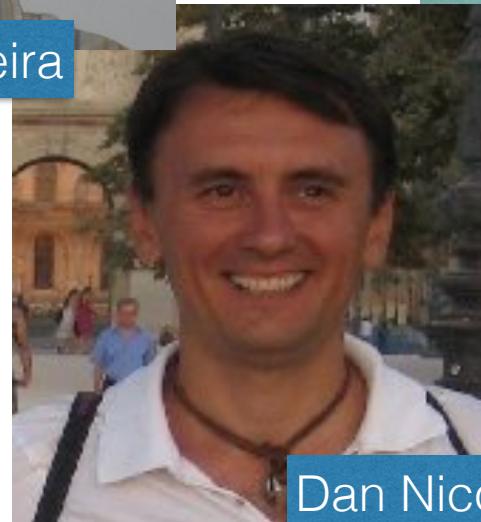
Alvaro Barbeira



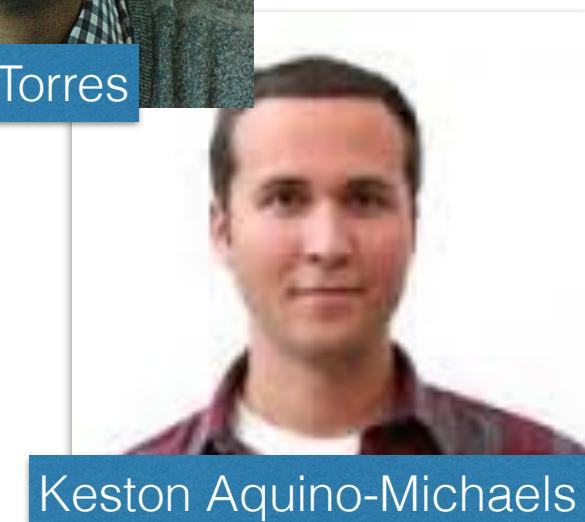
Jason Torres



Eric Gamazon



Dan Nicolae



Keston Aquino-Michaels

# Thank You

## Contributors

- Heather Wheeler
- Alvaro Barbeira
- Nancy Cox
- Kaanan P. Shah
- Eric Gamazon
- Dan Nicolae
- Keston Aquino Michaels
- Sahar Mozaffari
- Nicholas Knoblauch
- GTEx Consortium
- Josh Denny, Robert Carroll, Anne Eyler

## Data sources

- WTCCC
- NIMH: Depression Genes & Networks

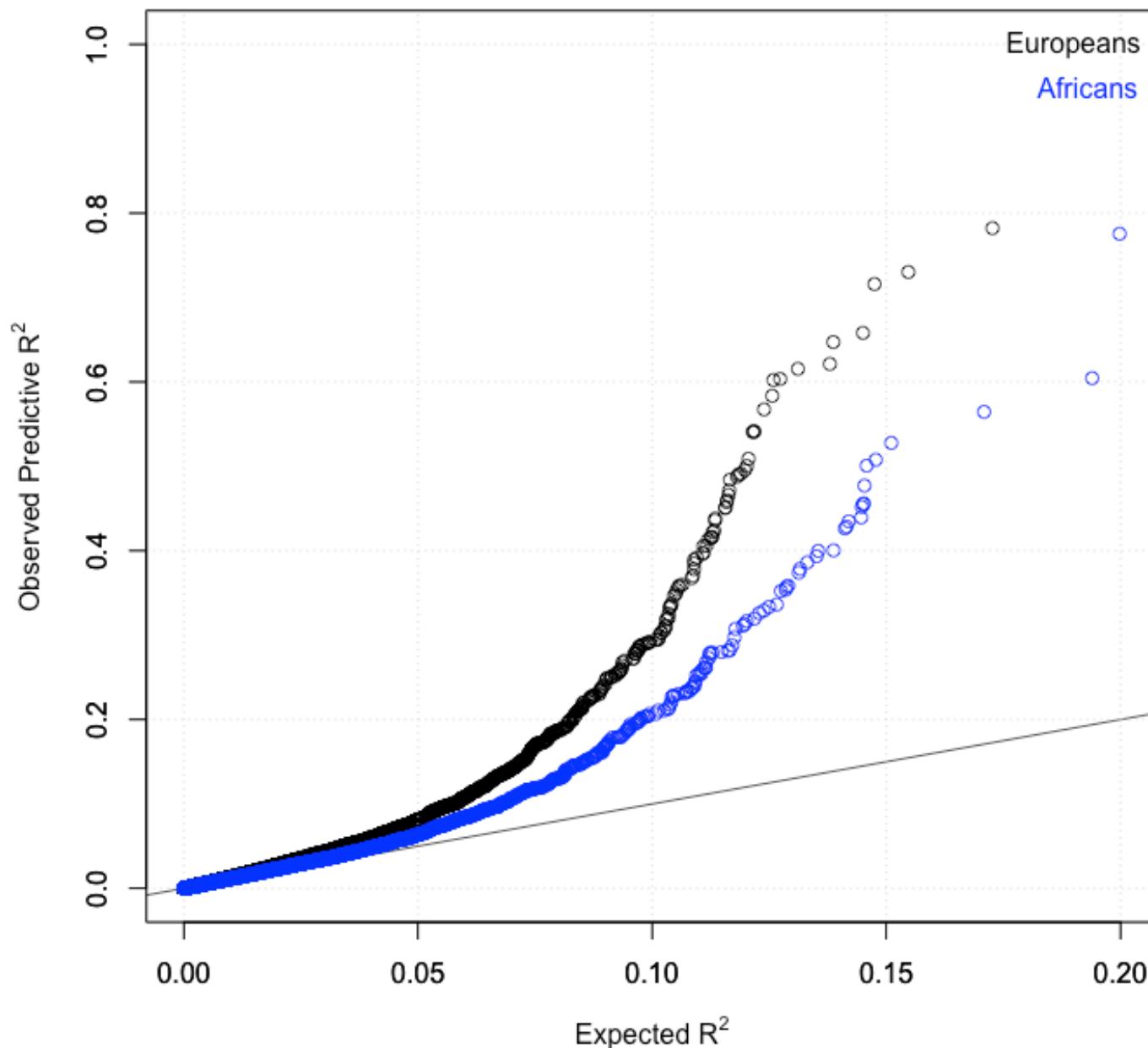
## Funding

- HKI was funded in part by UChicago CTSA NCI K12CA139160
- University of Chicago Diabetes Research and Training Center: P60 DK20595, P30 DK020595
- Genotype of Tissue Expression GTEx R01 MH090937 and R01 MH101820
- P50DA037844 - Integrated GWAS of complex behavioral and gene expression traits in outbred rats
- Pharmacogenomics of Anticancer Agents PAAR UO1GM61393
- Pharmacogenomics Research Network (PGRN) Statistical Analysis Resource (P-STAR) U19 HL065962
- Conte Center grant P50MH094267

# PrediXcan vs. MetaXcan - T1D WTCCC

$$Z_g = \sum_{l \in \text{Model}_g} w_{lg} \frac{\sigma_l}{\hat{\sigma}_g} \frac{\hat{\beta}_l}{\text{se}(\hat{\beta}_l)} \sqrt{\frac{1 - R_l^2}{1 - R_g^2}}$$

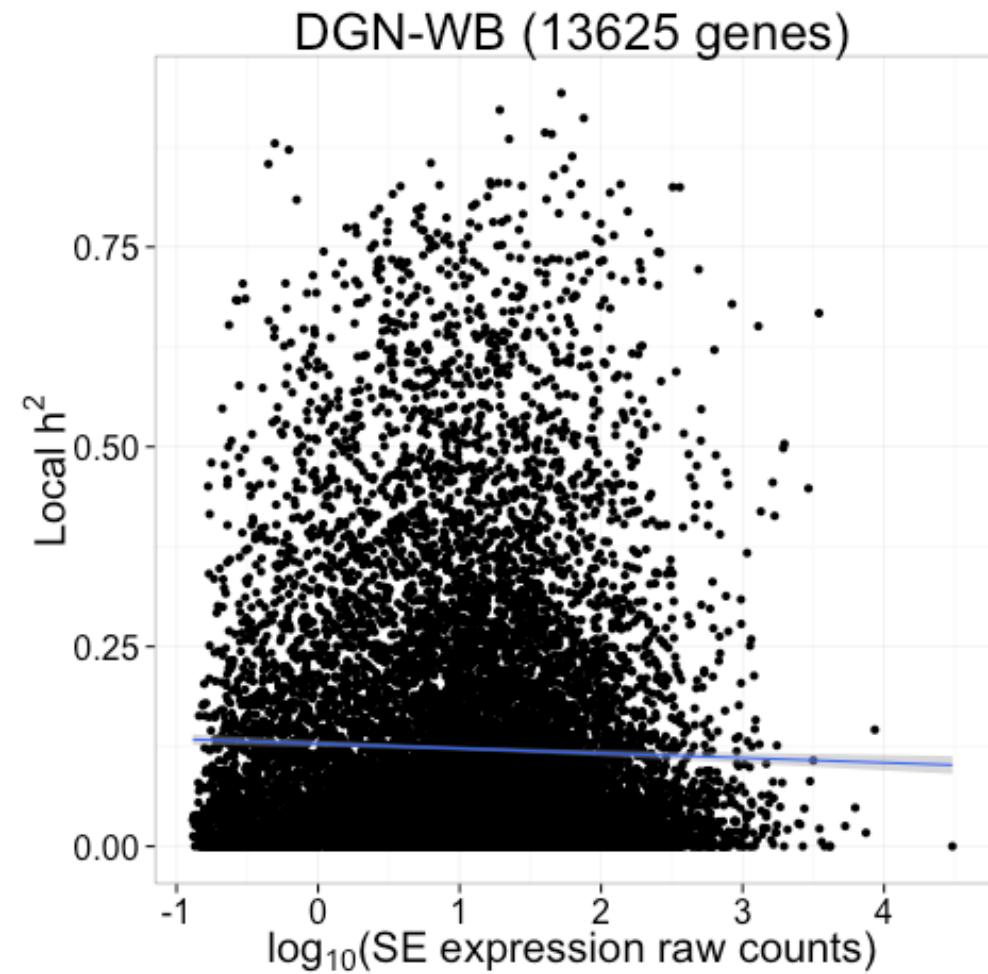
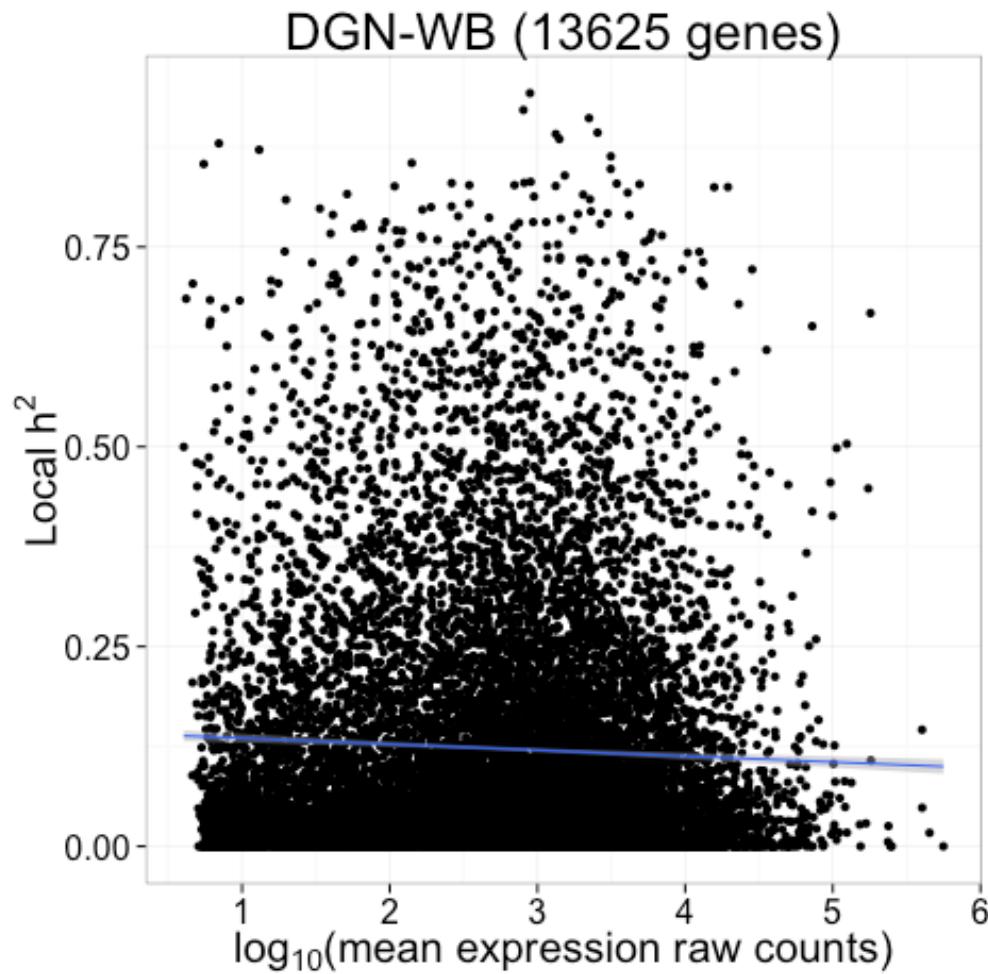
# Cross Population Prediction Performance



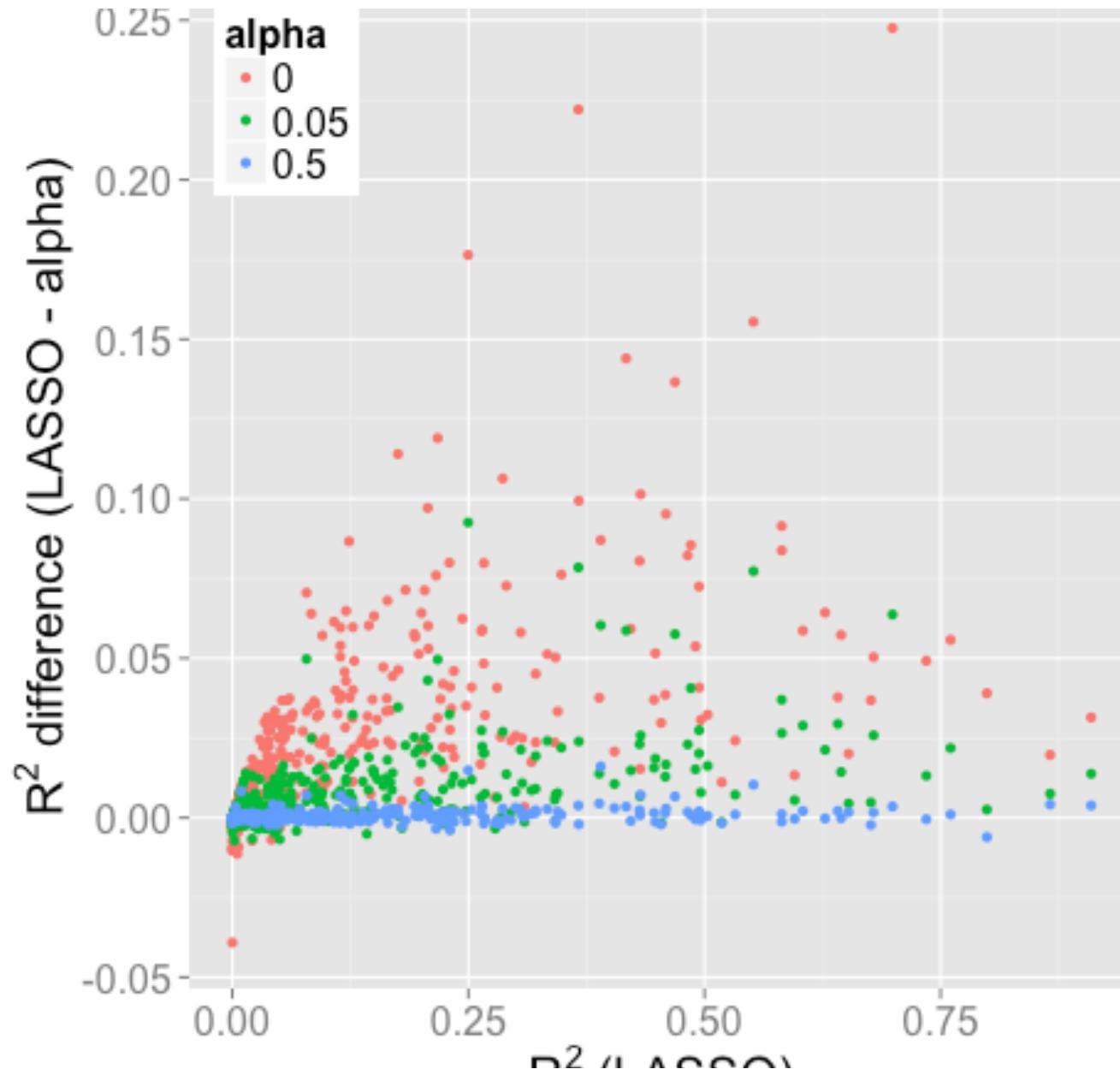
# MetaXcan Formula

$$Z_g = \sum_{l \in \text{Model}_g} w_{lg} \frac{\sigma_l}{\sqrt{\mathbf{W}'_g \Gamma_g \mathbf{W}_g}} \frac{\hat{\beta}_l}{\text{se}(\hat{\beta}_l)} \sqrt{\frac{1 - R_l^2}{1 - R_g^2}}$$

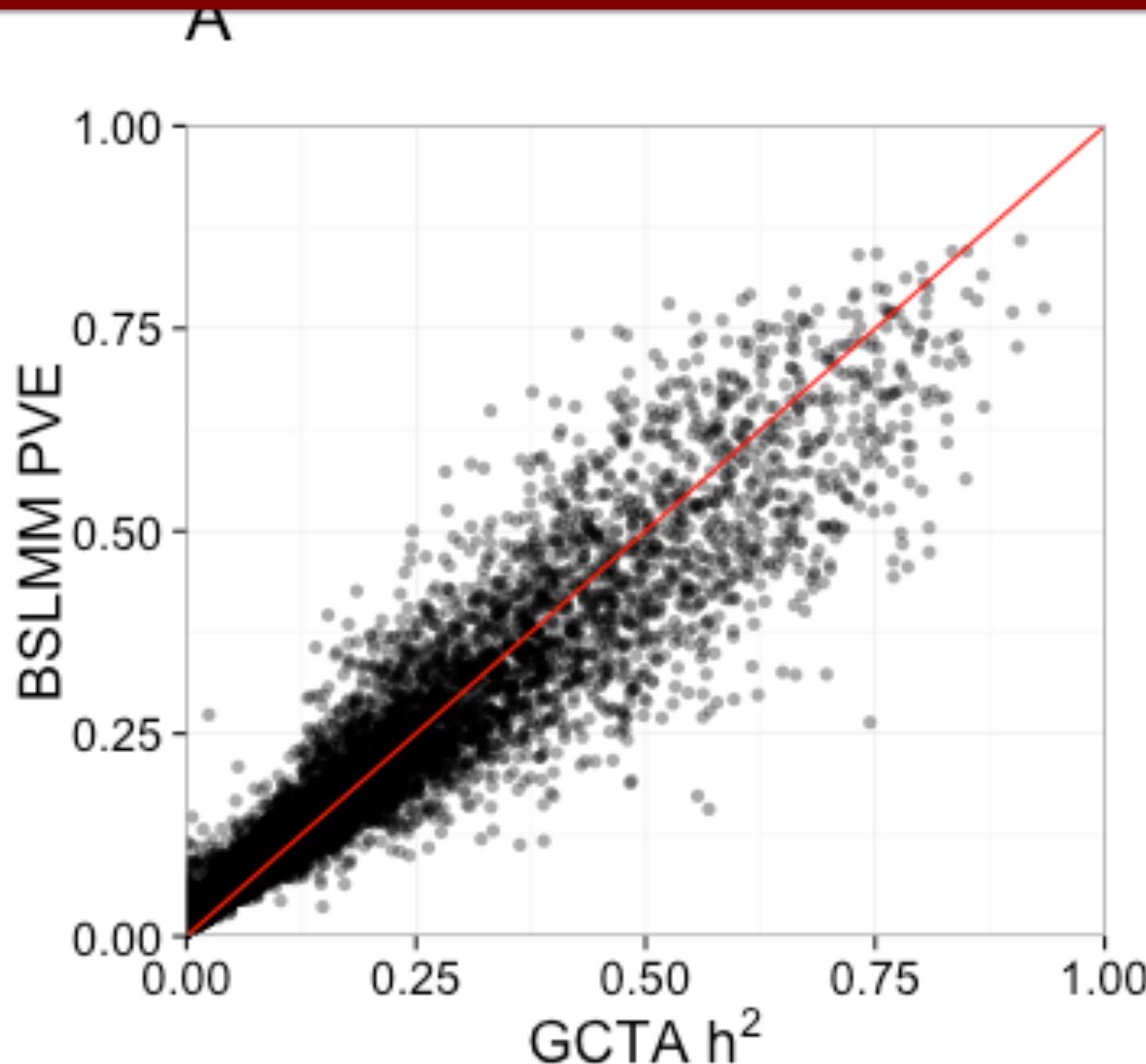
# Heritability vs. Expression Level



# Sparse Models Outperform Polygenic Models



# Similar Estimates of $h^2$ with BSMM and GCTA



# PrediXcan Results WTCCC - Gene2Pheno.org

gene2pheno.org

## G2P: Gene to Phenotype Query

Search Tab2

Example: Try searching for AKTIP in gname. Currently available diseases are BD, T2D, T1D, RA, CAD, CD, HT

Gene name:

Phenotype:

Tissue:

Submit

Your Input:

Gname: GATA6

7 results returned

gname	pval	tval	phenotype	source	tissue	cohort	model	R2	n.snps
GATA6	0.00021965	-3.695274229	CAD	DGN	WB	WTCCC	E-N0.5	0.0160770084023387	14
GATA6	0.004841634	-2.817385701	T2D	DGN	WB	WTCCC	E-N0.5	0.0160770084023387	14
GATA6	0.005981928	-2.748770339	T1D	DGN	WB	WTCCC	E-N0.5	0.0160770084023387	14
GATA6	0.030829947	-2.159260705	BD	DGN	WB	WTCCC	E-N0.5	0.0160770084023387	14
GATA6	0.138561263	-1.481170057	HT	DGN	WB	WTCCC	E-N0.5	0.0160770084023387	14
GATA6	0.523548294	0.637885382	RA	DGN	WB	WTCCC	E-N0.5	0.0160770084023387	14
GATA6	0.763356148	-0.301076516	CD	DGN	WB	WTCCC	E-N0.5	0.0160770084023387	14

