# Genome, Phenome, and What Happens in Between

Hae Kyung Im, PhD

THE UNIVERSITY OF CHICAGO

GSK, Philadelphia
June 8, 2017

To develop statistical and computational methods to sift through large amounts of genomic and other high dimensional data to make discoveries that can be translated to improve human health.

To catalog the phenotypic consequences of gene expression variation in humans

# Model Organism Knock Out



@hakyim      Natural Experiments to Function     4

# Model Organism Knockouts

# Natural Human Knockouts

- People with loss of function mutations in both copies of the gene

- Natural experiments

- We can measure phenotypes to learn function of the gene



Approximately 20,000 genes → Identification of complete gene knockouts → Disease pathology / Biochemical testing / Analyse gene's role

**R. M. Plenge, "Human genes lost and their functions found," Nature, 2017.**

## Human knockouts and phenotypic analysis in a cohort with a high rate of consanguinity

Danish Saleheen, Pradeep
Khetarpal, Hong-Hee Won
Samocha, Benjamin Weisl
Shahid Abbas, Faisal Maje
Mucksavage, Nadeem Qar

Systematic effort to understand the consequence of complete disruption of every human gene

**Knockout**

0

Systematic effort to understand the consequence of complete disruption of every human gene

**"Knockdown"**

**Knockout**

0

Systematic effort to understand the consequence
of ~~complete~~ disruption of every human gene
**partial**

# Human Knockout vs "Knockdown"

## Knockout

- Large effect sizes

- Small sample size
- Need to sequence large number of individuals

## Knockdown

- Small effect sizes
- LD-contamination
- Pleiotropy

- Large sample size
- Cheaper genotyping may be enough
- Sequence data can be used

- Increase sample size

- to address burden of larger sample sizes
  - use Summary-PrediXcan

- Compute colocalization and discard if GWAS and eQTL signals are independent
    - COLOC (Giambartolomei et al, PLoS Genetics 2013)
    - RTC (Nica, …, Dermitzakis et al 2010)
    - eCAVIAR (Hormozdiari … Eskin et al, AJGH 2017)
    - ENLOC (Wen et al, PLoS Genetics 2017)
    - HEIDI (Zhu, …, Visscher, Yang, Nature Gen. 2016)

# Use Causal rather than Associated SNPs



ERAP2

$R^2 = 0.82$

**Observed Expression**

**Predicted Expression**

ERAP2 is predicted by 80 SNPs but has one underlying causal variant

- Post filtering step with COLOC or other methods

- Current prediction models are purely statistical
    - Use causal predictors to reduce chance of LD
    - S-PrediXcan needs to efficiently impute GWAS results for causal variants
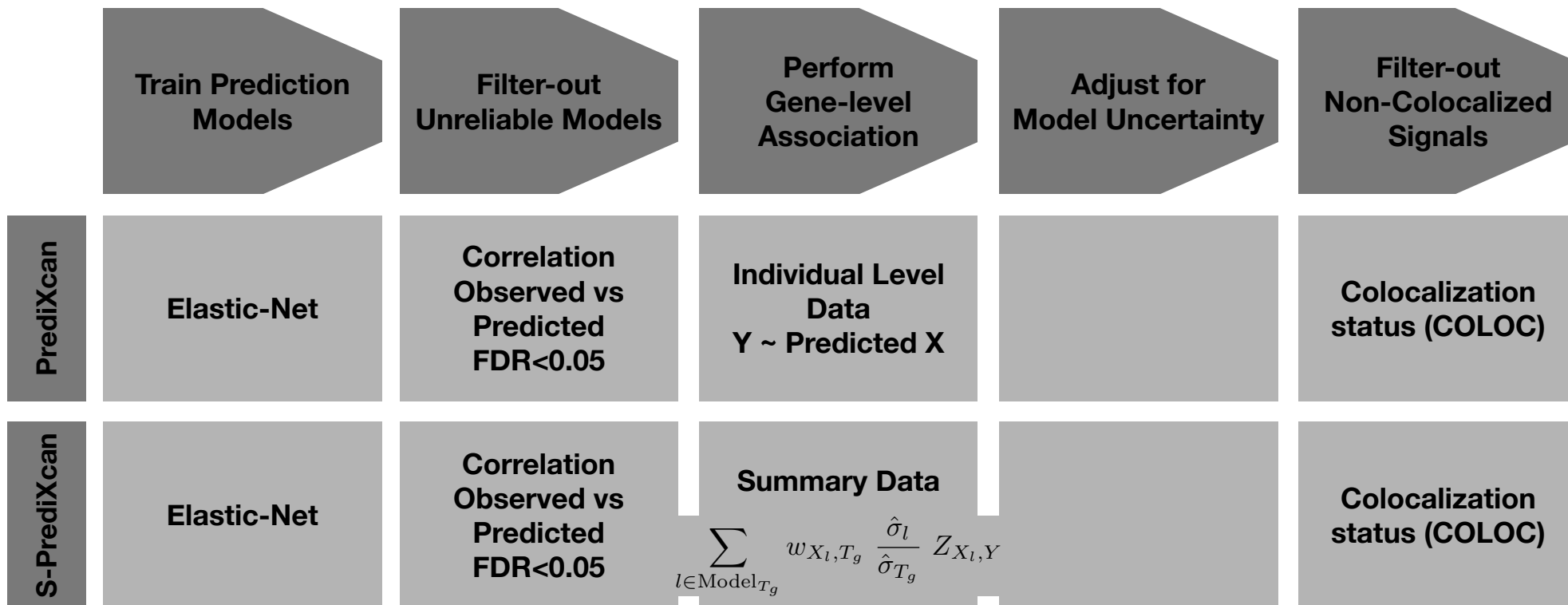
- Experiments in model systems

# Best Practices Framework: MetaXcan

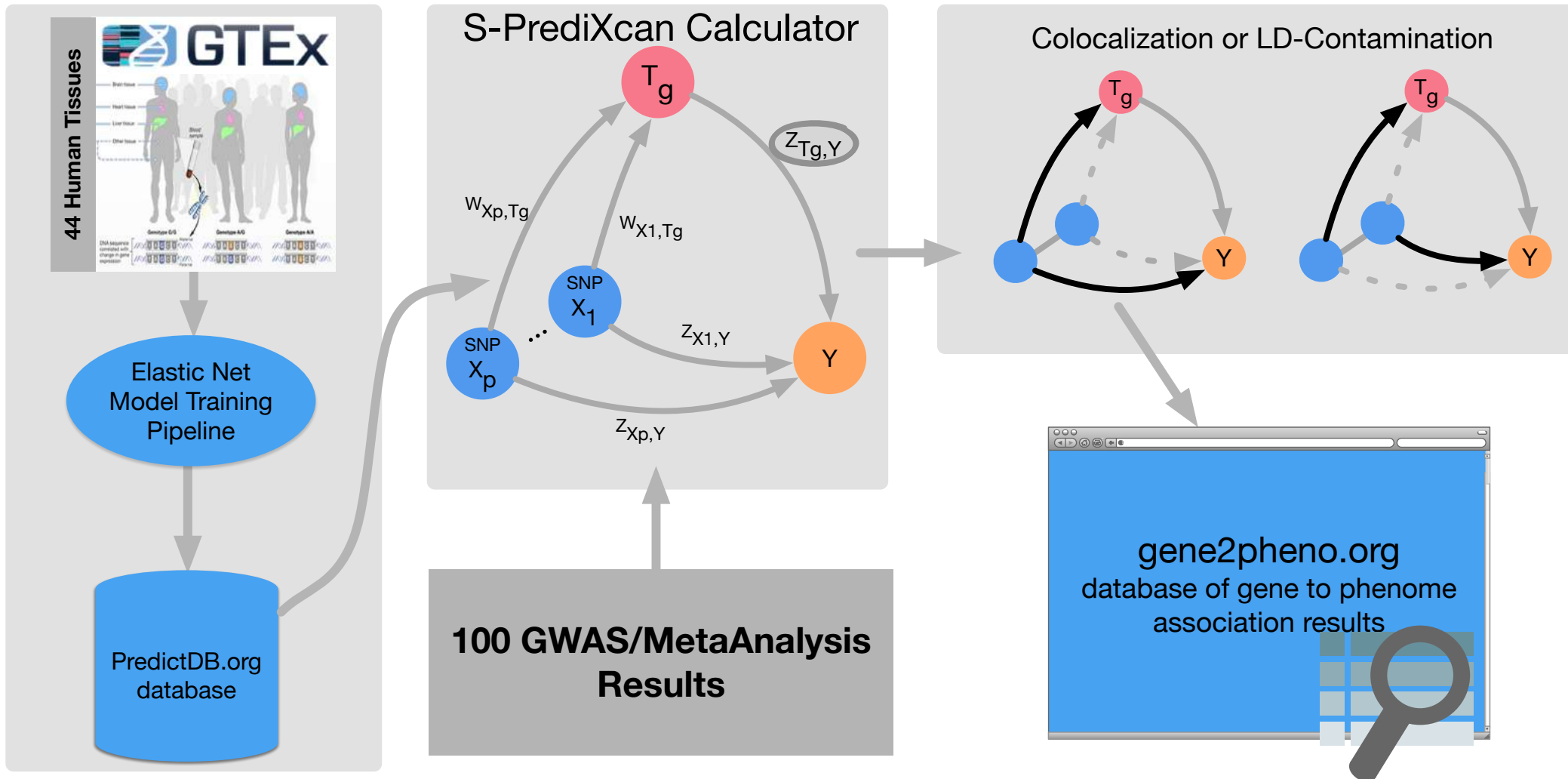| | Train Prediction Models | Filter-out Unreliable Models | Perform Gene-level Association | Adjust for Model Uncertainty | Filter-out Non-Colocalized Signals |
|---|---|---|---|---|---|
| PrediXcan | Elastic-Net | Correlation Observed vs Predicted FDR<0.05 | Individual Level Data Y ~ Predicted X | | Colocalization status (COLOC) |
| S-PrediXcan | Elastic-Net | Correlation Observed vs Predicted FDR<0.05 | Summary Data $\sum_{l \in \mathrm{Model}_{T_g}} w_{X_l,T_g} \frac{\hat{\sigma}_l}{\hat{\sigma}_{T_g}} Z_{X_l,Y}$ | | Colocalization status (COLOC) |

**Causal models**

**Incorporate Model Uncertainty**

**Imputation of GWAS results**

**Allow for multiple causal variants in COLOC**

# Computing Phenotypic Consequences with Summary Stat



44 Human Tissues

**GTEx**

Elastic Net Model Training Pipeline

PredictDB.org database

## S-PrediXcan Calculator

$T_g$

$W_{Xp,Tg}$

$W_{X1,Tg}$

$Z_{Tg,Y}$

SNP $X_1$

SNP $X_p$

...

$Z_{X1,Y}$

$Z_{Xp,Y}$

Y

**100 GWAS/MetaAnalysis Results**

## Colocalization or LD-Contamination

$T_g$ Y $T_g$ Y

**gene2pheno.org**
database of gene to phenome association results

https://github.com/hakyimlab/PredictDBPipeline

https://github.com/hakyimlab/MetaXcan

## gene2pheno.org

### Metaxcan Association

Data Release: September 7, 2016.

Prediction models and covariances built with GTEx V6P and DGN on HapMap SNPs.

### Results:

**What to show:**

Results ▼

| Gene Name: | ☑ Ordered | Phenotype: | Tissue: | R2 threshold: | Pvalue threshold: | Record limit: |
|---|---|---|---|---|---|---|
| | | All ▼ | All ▼ | 0.01 | 0.05 | 100 |
| | | Patterns: | Patterns: | | | |

Show 20 ⬍ entries                                                                                   Search: [          ]

| | gene_name ⬍ | zscore ⬍ | effect_size ⬍ | pval ⬍ | phenotype ⬍ | tissue | pred_perf_r2 ⬍ | pred_perf_pval ⬍ | pre |
|---|---|---|---|---|---|---|---|---|---|
| 1 | HLA-DQA2 | 38.23 | 0.46 | 0 | RA_OKADA_TRANS_ETHNIC | TW_Artery_Aorta_Elastic_Net_0.5 | 0.47 | 2.3e-28 | |
| 2 | HLA-DQA2 | 38.62 | 0.5 | 0 | RA_OKADA_TRANS_ETHNIC | TW_Colon_Sigmoid_Elastic_Net_0.5 | 0.48 | 6.7e-19 | |
| 3 | CFH | -37.04 | | 2.8e-300 | AdvancedAMD_2015 | DGN_WB_Elastic_Net_0.5 | 0.01 | 0.015 | |
| 4 | HLA-DQA2 | 36.67 | 0.49 | 1.9e-294 | RA_OKADA_TRANS_ETHNIC | DGN_WB_Elastic_Net_0.5 | 0.76 | 5.4e-286 | |
| 5 | HLA-DRB1 | -36.58 | -0.54 | 6.5e-293 | RA_OKADA_TRANS_ETHNIC | DGN_WB_Elastic_Net_0.5 | 0.75 | 3.5e-277 | |
| 6 | CFHR3 | 36.47 | | 3.8e-291 | AdvancedAMD_2015 | TW_Adrenal_Gland_Elastic_Net_0.5 | 0.4 | 2.7e-15 | |
| 7 | HLA-DQA2 | 36.41 | 0.49 | 3.1e-290 | RA_OKADA_TRANS_ETHNIC | TW_Prostate_Elastic_Net_0.5 | 0.38 | 1.9e-10 | |
| 8 | CFHR3 | 35.95 | | 4.6e-283 | AdvancedAMD_2015 | TW_Breast_Mammary_Tissue_Elastic_Net_0.5 | 0.08 | 0.00017 | |
| 9 | HLA-DQA2 | 35.92 | 0.4 | 1.5e-282 | RA_OKADA_TRANS_ETHNIC | TW_Pancreas_Elastic_Net_0.5 | 0.37 | 1.2e-16 | |
| 10 | CFHR1 | 35.73 | | 1.5e-279 | AdvancedAMD_2015 | TW_Brain_Putamen_basal_ganglia_Elastic_Net_0.5 | 0.08 | 0.0085 | |
| 11 | CFHR1 | 35.58 | | 3.1e-277 | AdvancedAMD_2015 | TW_Adipose_Visceral_Omentum_Elastic_Net_0.5 | 0.05 | 0.0028 | |

https://github.com/hakyimlab/MetaXcan

# Tissue Specificity

Target Genes Found across Multiple Tissues

Mean Z2

Histogram for tissue–specificity
Phenotype: GIANT_HEIGHT
cutoff1 = 2.5e–07
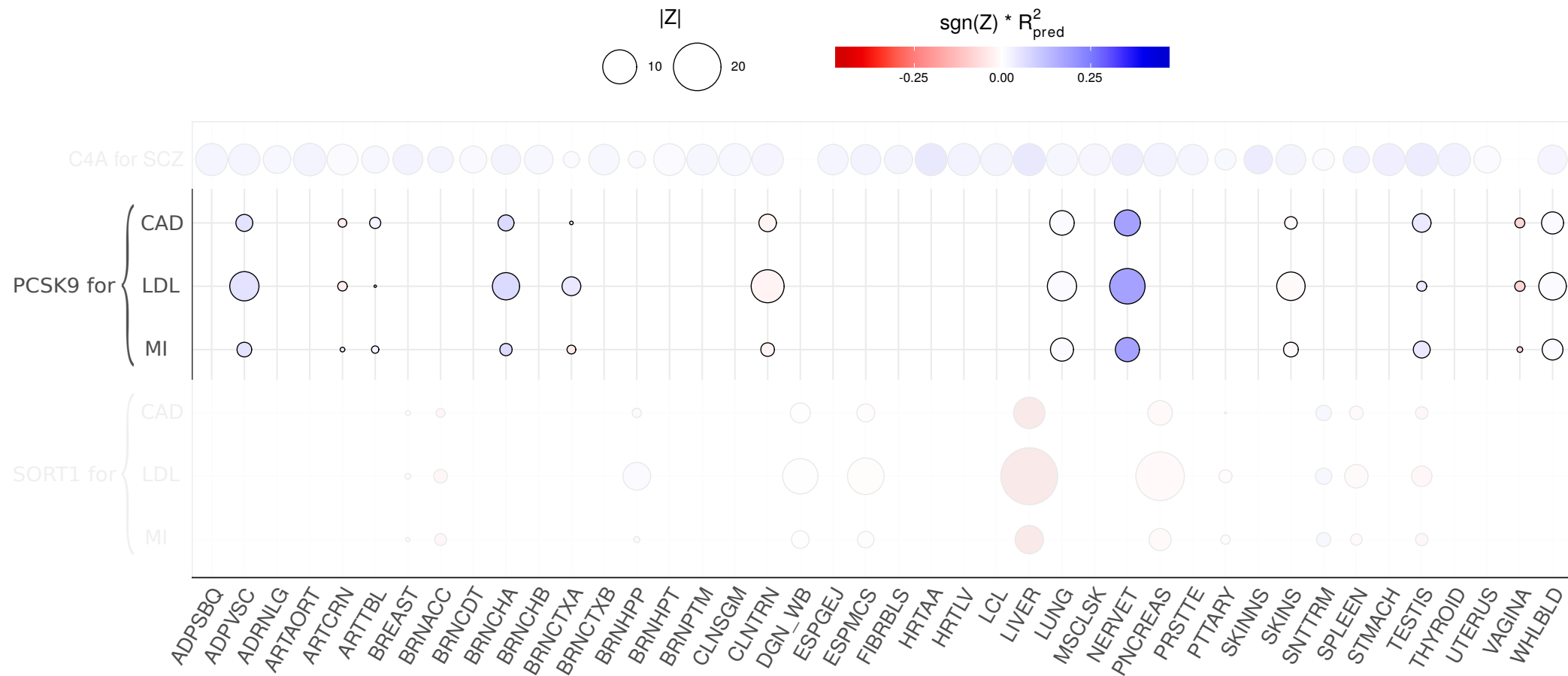cutoff2 = 0.05 / N_models, 1690 genes with more than 0 models.

# From noncoding variant to phenotype via *SORT1* at the 1p13 cholesterol locus

Kiran Musunuru[1,2,3]*, Alanna Strong[4]*, Maria Frank-Kamenetsky[5], Noemi E. Lee[1], Tim Ahfeldt[1,6], Katherine V. Sachs[4], Xiaoyu Li[4], Hui Li[4], Nicolas Kuperwasser[1], Vera M. Ruda[1], James P. Pirruccello[1,2], Brian Muchmore[7], Ludmila Prokunina-Olsson[7], Jennifer L. Hall[2,8], Eric E. Schadt[9], Carlos R. Morales[10], Sissel Lund-Katz[11], Michael C. Phillips[11], Jamie Wong[5], William Cantley[5], Timothy Racie[5], Kenechi G. Ejebe[1,2]

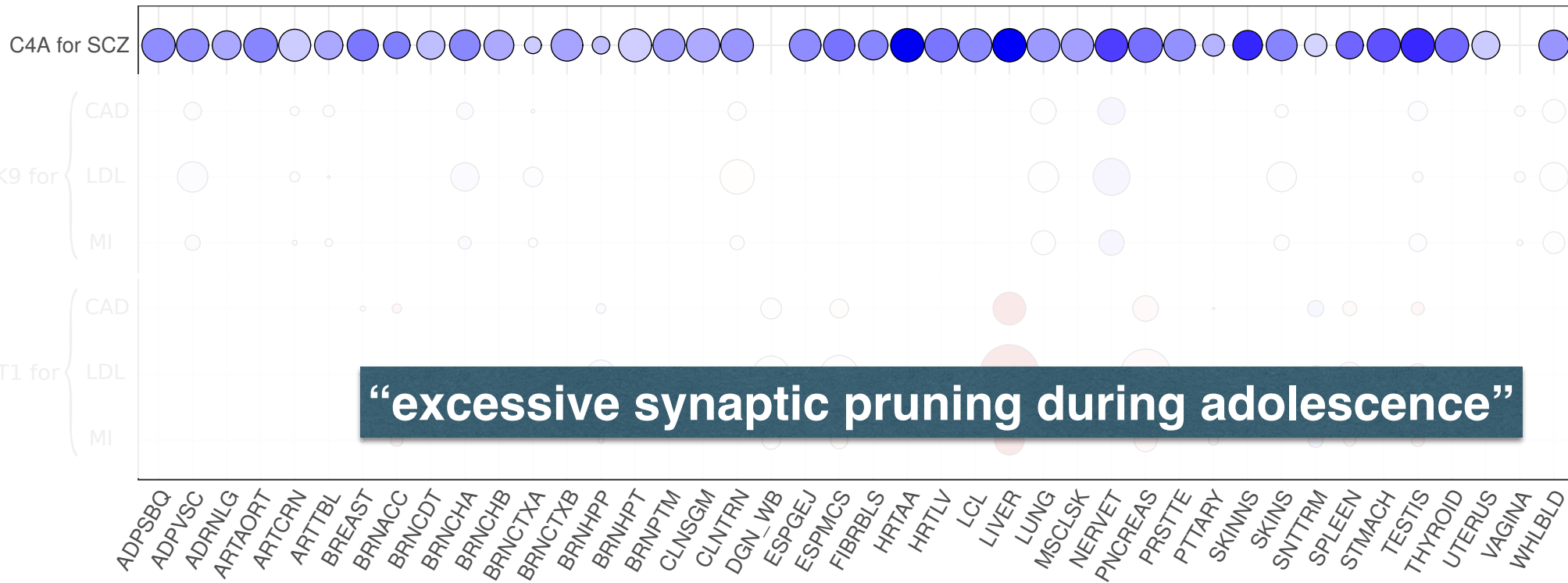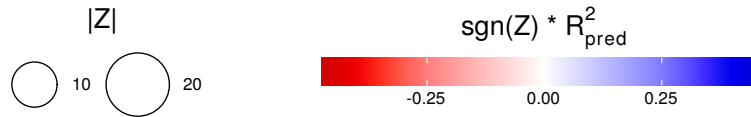**HUMAN GENETICS**

# Cardiometabolic risk loci share downstream cis- and trans-gene regulation across tissues and diseases

Oscar Franzén,[1,2*] Raili Ermel,[3,4*] Ariella Cohain,[1*] Nicholas K. Akers,[1]

"excessive synaptic pruning during adolescence"
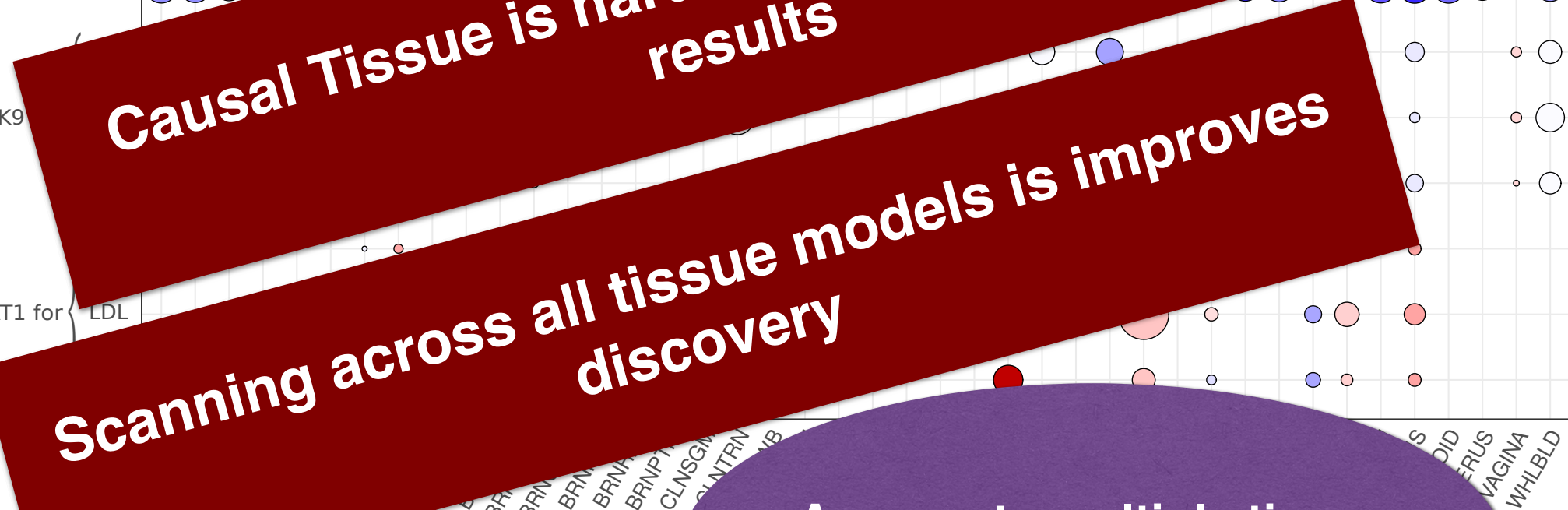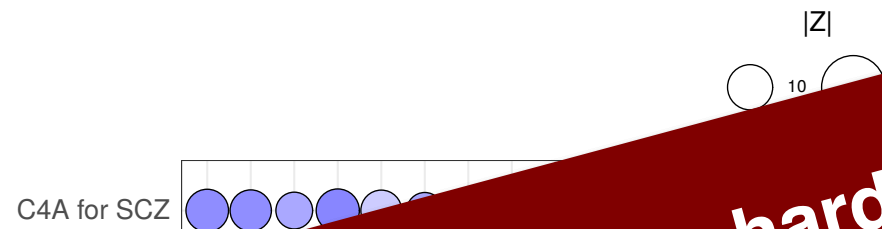
# Schizophrenia risk from complex variation of complement component 4

Aswin Sekar[1,2,3], Allison R. Bialas[4,5], Heather de Rivera[1,2], Avery Davis[1,2], Timothy R. Hammond[4], Nolan Kamitaki[1,2], Katherine Tooley[1,2], Jessy Presumey[5], Matthew Baum[1,2,3,4], Vanessa Van Doren[1], Giulio Genovese[1,2], Samuel A. Rose[2], Robert E. Handsaker[1,2], Schizophrenia Working Group of the Psychiatric Genomics Consortium*, Mark J. Daly[2,6], Michael C. Carroll[5], Beth Stevens[2,4] & Steven A. McCarroll[1,2]

**Causal Tissue is hard to establish from these results**

**Scanning across all tissue models is improves discovery**

**Aggregate multiple tissue results into one gene level one**

# Multi-Tissue PrediXcan

$$Y = a + b_1 X_g^{\text{tissue}_1} + b_2 X_g^{\text{tissue}_2} + \cdots + b_k X_g^{\text{tissue}_k} + \epsilon$$

$$Y = a + b_1 X_g^{\text{tissue}_1} + b_2 X_g^{\text{tissue}_2} + \cdots + b_k X_g^{\text{tissue}_k} + \epsilon$$

What if we only have univariate regression coefficients?

$$Y = a + b_1 X_g^{\text{tissue}_1} + b_2 X_g^{\text{tissue}_2} + \cdots + b_k X_g^{\text{tissue}_k} + \epsilon$$

What if we only have univariate regression coefficients?

$$Y = a + \beta_1 X_g^{\text{tissue}_1} + \epsilon'$$

$$Y = a + \beta_2 X_g^{\text{tissue}_2} + \epsilon''$$

$$\cdots$$

$$Y = a + \beta_k X_g^{\text{tissue}_k} + \epsilon'''$$

$$Y = a + b_1 X_g^{\text{tissue}_1} + b_2 X_g^{\text{tissue}_2} + \cdots + b_k X_g^{\text{tissue}_k} + \epsilon$$

What if we only have univariate regression coefficients?

$$Y = a + \beta_1 X_g^{\text{tissue}_1} + \epsilon'$$
$$Y = a + \beta_2 X_g^{\text{tissue}_2} + \epsilon''$$
$$\ldots$$
$$Y = a + \beta_k X_g^{\text{tissue}_k} + \epsilon'''$$

$$\hat{\boldsymbol{b}} = (\boldsymbol{X}'\boldsymbol{X})^{-1}\boldsymbol{D}\widehat{\boldsymbol{\beta}}$$
$$\text{var}(\widehat{\boldsymbol{b}}) = \sigma_j(\boldsymbol{X}'\boldsymbol{X})^{-1}$$

$$\chi_k^2 = \hat{\boldsymbol{b}}'(\boldsymbol{X}'\boldsymbol{X})^{-1}\hat{\boldsymbol{b}}$$

$$D_t = \sum_i X_{it}^2$$

$$\chi_t^2 = \hat{\boldsymbol{b}}'(\boldsymbol{X}'\boldsymbol{X})^{-1}\hat{\boldsymbol{b}}$$
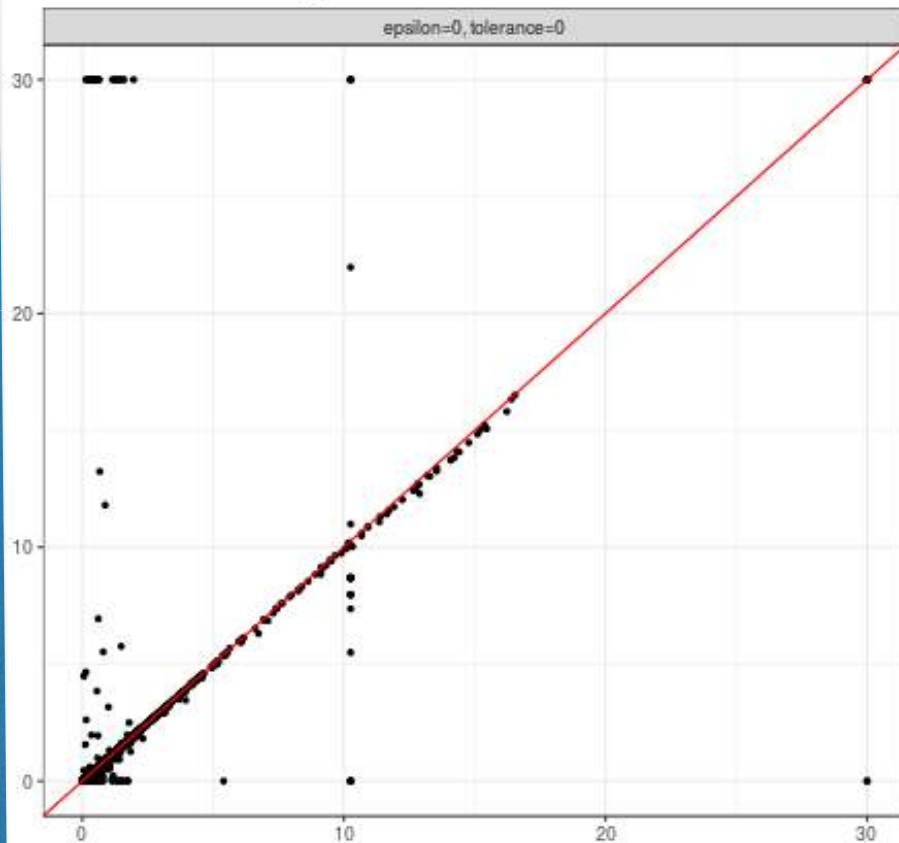
Predicted
expression of a gene
across tissues

We can predict
expression in reference
population
(1000G, GTEx, etc)

Inverse may not be
computable because
of correlation between
tissues

WTCCC T1D Phenotype: PrediXcan MultiTissue vs Combined Univariate PrediXcan

**Combined Univeraite PrediXcan -log10 p**

**Multivariate PrediXcan -log10 p**

$$\hat{b} = (X'X)^{-1} D \hat{\beta}$$

$$\text{var}(\hat{b}) = \sigma_j (X'X)^{-1}$$

$$\chi_k^2 = \hat{b}'(X'X)^{-1}\hat{b}$$

WTCCC T1D Phenotype: PrediXcan MultiTissue vs Combined Univariate PrediXcan
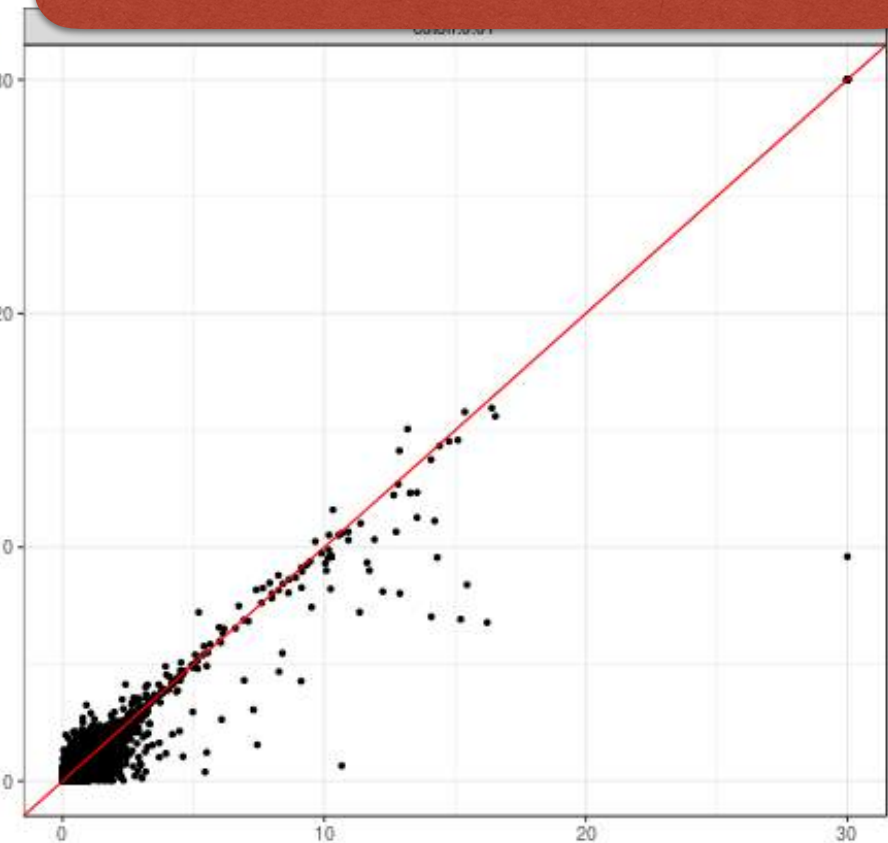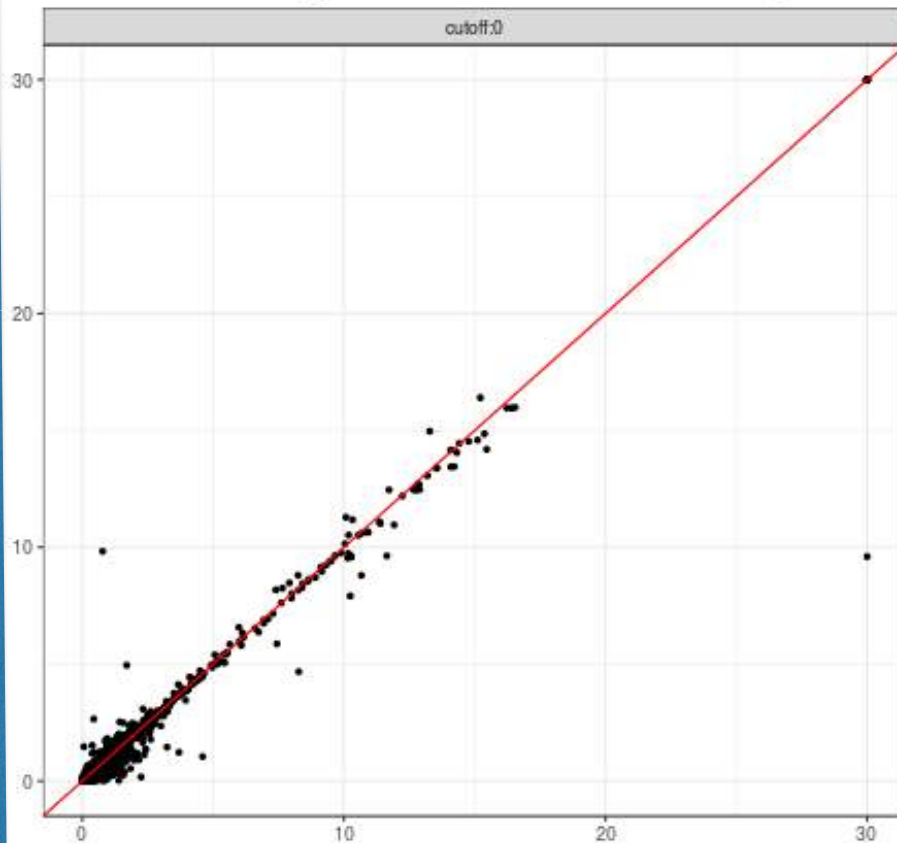
**Combined Univariate -log10 p**

**Multivariate**

**Avoid singularity of covariance eliminating axis of variations with small eigenvalues**

# Combined Summary PrediXcan vs Multivariate



**Combined Summary PrediXcan -log10 p**

WTCCC T1D Phenotype: PrediXcan MultiTissue vs Summary-Multitissue-P

cutoff:0

**Predicted expression is estimated in study sample**

**Multivariate PrediXcan -log10 p**

$$\hat{b} = (X'X)^{-1}D\widehat{\beta}$$
$$\mathrm{var}(\widehat{b}) = \sigma_j(X'X)^{-1}$$

$$\chi_k^2 = \hat{b}'(X'X)^{-1}\hat{b}$$

WTCCC T1D Phenotype: PrediXcan MultiTissue vs Summary-Multitissue-P

cutoff:0

**Predicted expression is estimated in different samples**

**Combined Univeraite PrediXcan -log10 p**

**Multivariate PrediXcan -log10 p**

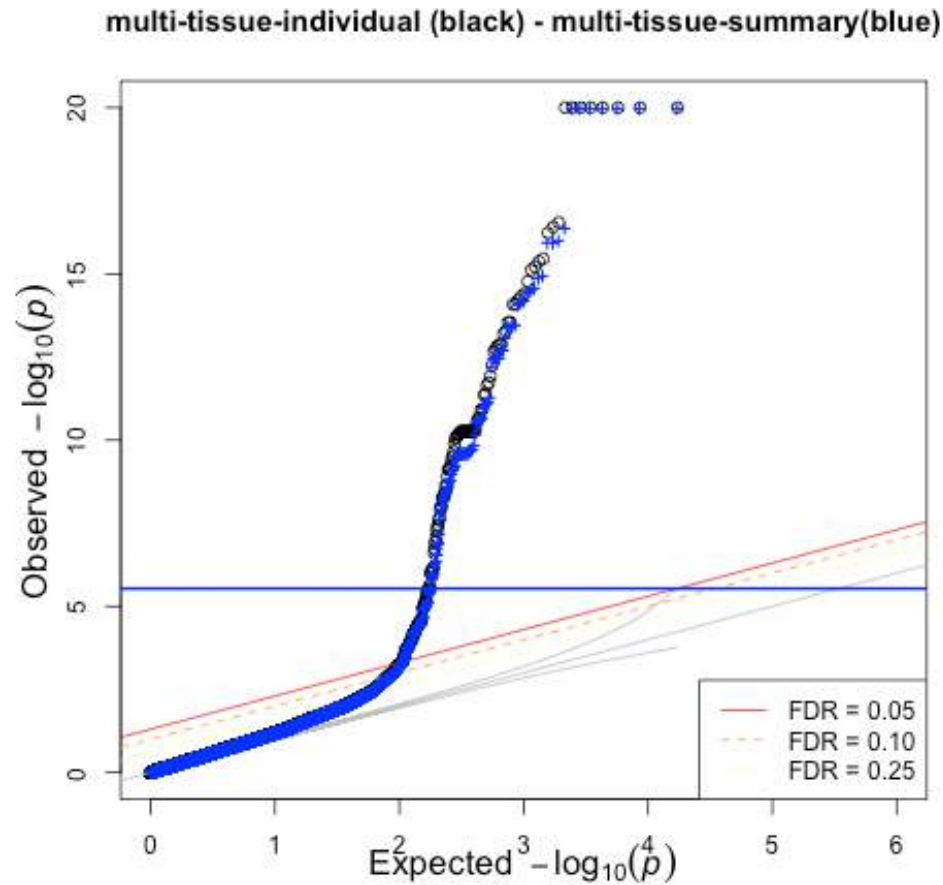WTCCC T1D Phenotype: PrediXcan MultiTissue vs Summary-Multitissue-PrediXcan
(GTEx SNP covariance, Expression covariance from GTEX SNP intersection to GWAS)
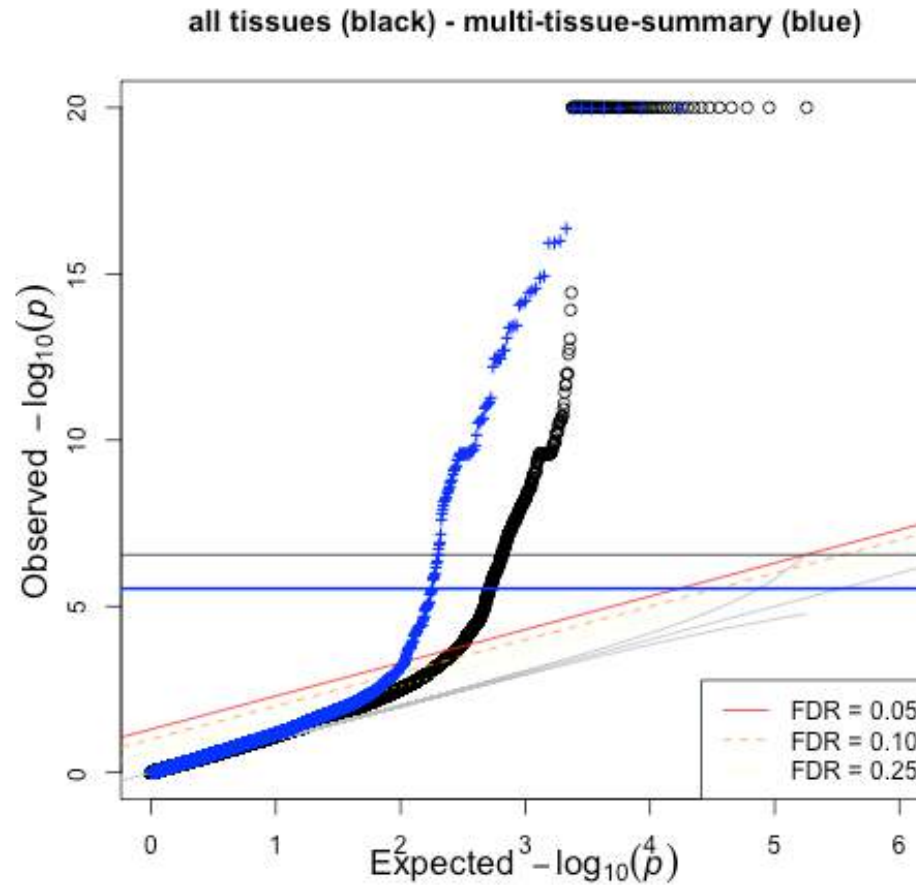
**Combined Univeraite PrediXcan -log10 p**

**Covariance estimated directly from SNPs rather than predicting**

**Multivariate PrediXcan -log10 p**

multi-tissue-individual (black) - multi-tissue-summary(blue)

all tissues (black) - multi-tissue-summary (blue)

- Human knockouts are invaluable experiments of nature that provides information on function of genes

- Human "knockdown" gene2pheno.org, related and complementary

- Need to develop new methods to address challenges

- Summary Multi Tissue PrediXcan

**Haky Im Lab**

- Alvaro Barbeira
- Jiamao Zheng
- Scott Dickinson
- Rodrigo Bonazzola
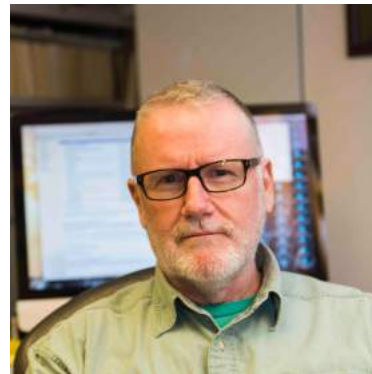- Milton Pividori

**Collaborators**

- Jason Torres
- Kaanan Shah
- Heather Wheeler
- Eric Tortesson
- Tzintzuni Garcia
- Nancy Cox
- Dan Nicolae
- Graeme Bell

- Funding
  - R01MH107666 (HKI)
  - R01MH101820 (GTEx)
  - P30DK020595 (DRTC)

No conflicts of interests