

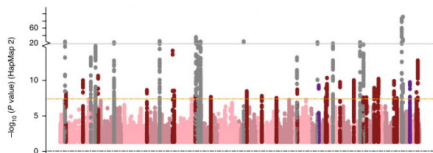
Large-scale genome-wide enrichment analyses of 31 human complex traits

Xiang Zhu
University of Chicago

ICSA 2017, June 28

Examining associations between variables is a useful tool in genetics.

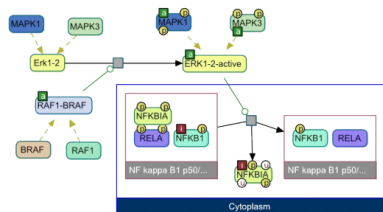
GWAS: SNP~Phenotype



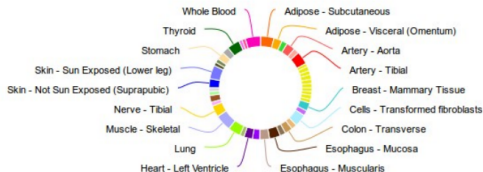
1000 Genomes: SNP~SNP



Biological Pathways: Gene~Gene



GTEX: Gene~Tissue



Enrichment analysis combines multiple sources of association.

- **SNP-trait:** GWAS summary statistics
- **SNP-SNP:** linkage disequilibrium (LD)
- **Gene-Gene:** biological pathways
- **Gene-Tissue:** RNA-seq across tissues

Let's keep this talk "jargon-free".

1. GWAS summary statistics
2. gene set enrichment analysis

What are GWAS summary statistics?

- **Data:** phenotype Y and genotype X
- **Size:** n (>10K) individuals and p (>1M) variants (SNPs)

$$Y := \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad X := \begin{bmatrix} x_{11} & \dots & x_{1j} & \dots & x_{1p} \\ x_{21} & \dots & x_{2j} & \dots & x_{2p} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ x_{n1} & \dots & x_{nj} & \dots & x_{np} \end{bmatrix}$$

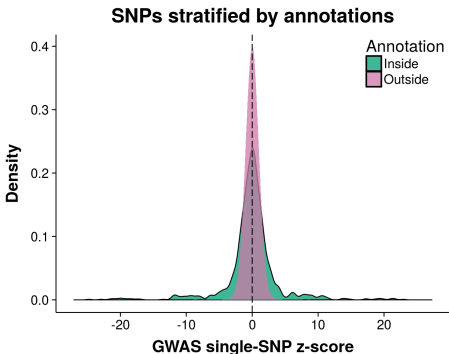
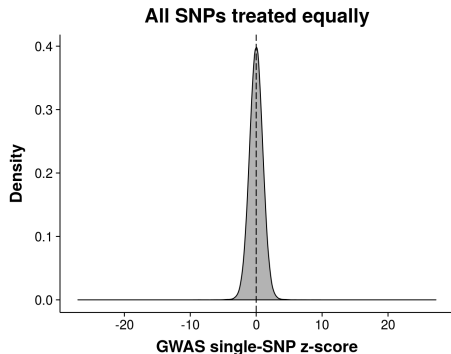
- **Model:** single-SNP association analysis

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \sim \begin{bmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{nj} \end{bmatrix} \rightsquigarrow \begin{cases} \hat{\beta}_j: & \text{marginal effect estimate} \\ \hat{\sigma}_j: & \text{standard error of } \hat{\beta}_j \end{cases}$$

- **Availability:** much more **accessible** than $\{Y, X\}$ (Nat. Genet., 2012)

What is enrichment analysis?

- **Phenotype:** low-density lipoprotein (Teslovich *et al.*, 2010)
- **Pathway:** *chylomicron-mediated lipid transport* (17 genes)
- **Annotation:** Is the SNP “near” a pathway gene? (yes or no)



Recent reviews: de Leeuw *et al.* (2016); Pers (2016); Mooney *et al.* (2014); Wang *et al.* (2010).

The “enrichment” idea is simple, but there are (at least) two technical issues.

1. If the gene set is truly enriched, we should relax significance threshold for “green” SNPs, but how much to relax?
 - **Data-driven** threshold ← Function (Pathway, Phenotype)
 - Maintained type 1 error + Improved power
2. The “inflated” green curve may be driven by correlation between SNPs (LD), rather than enrichment of signal.
 - SNP 1 has a large genetic effect on a trait.
 - SNPs 2-100 have zero effect, but all are in high LD with SNP 1.
 - Thus, SNPs 1-100 all show very large **single-SNP** z-scores.

We develop a method that systematically utilizes enrichment information.

Public Data

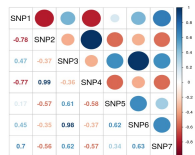
GWAS summary statistics



$$\hat{\beta} := (\hat{\beta}_1, \dots, \hat{\beta}_p)^T, \quad \hat{S} := \text{diag}\{(\hat{s}_1, \dots, \hat{s}_p)^T\}$$

$\hat{\beta}_j$: marginal effect estimate of SNP j
 \hat{s}_j : standard error of $\hat{\beta}_j$

External LD estimates



$$\tilde{R} := [\hat{r}_{ij}]_{p \times p}, \quad \hat{r}_{ij}: \text{LD between SNP } i \text{ and } j$$

Predefined gene sets

Gene	Chr	Start	End
BRAF	7	140419127	140624564
MAPK1	22	22108789	22221970
MAPK3	16	30125426	30134827
NFKB1	4	103422486	103538459

$$\alpha_j := 1 \text{ \{SNP } j \text{ is "near" a gene in the set\}}$$

Inference

Bayesian Model

Likelihood function

$$L(\beta) := \text{Normal}(\hat{\beta}; \hat{S} \hat{R} \hat{S}^{-1} \beta, \hat{S} \hat{R} \hat{S})$$

Prior distribution

$$\beta_j \sim \pi_j \cdot \text{Normal}(0, \sigma_\beta^2) + (1 - \pi_j) \delta_0$$

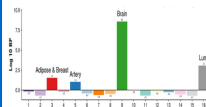
Baseline model (M_0)

$$\log_{10} \left(\frac{\pi_j}{1 - \pi_j} \right) = \theta_0$$

Enrichment model (M_1)

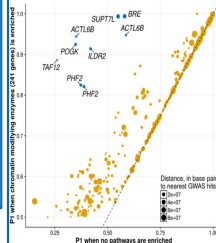
$$\log_{10} \left(\frac{\pi_j}{1 - \pi_j} \right) = \theta_0 + \alpha_j \theta$$

Gene set enrichment



$$\text{BF}(\text{gene set}) = \frac{\Pr(\text{Data} | M_1)}{\Pr(\text{Data} | M_0)}$$

Gene prioritization



$$\rho(\beta | \text{Data}, M_1) \text{ vs } \rho(\beta | \text{Data}, M_0)$$

Address Issue 1: Learning enrichment from data

Model-based approach:

- Assume that SNP j is trait-associated with probability π_j
- Represent π_j as a function of SNP-level annotation a_j

$$\log_{10} \left(\frac{\pi_j}{1-\pi_j} \right) := \theta_0 + a_j \theta$$

- Estimate enrichment parameter θ from data

Data-adaptive threshold:

- “Enriched” data $\rightsquigarrow \theta > 0 \rightsquigarrow$ larger π_j for $a_j = 1 \rightsquigarrow$ increased power
- “Null” data $\rightsquigarrow \theta \approx 0 \rightsquigarrow$ unchanged $\pi_j \rightsquigarrow$ maintained type 1 error

References: Veyrieras *et al.* (2008); Carbonetto & Stephens (2013); Pickrell (2014); Kichaev *et al.* (2014); Chen *et al.* (2016); Li & Kellis (2016); Wen *et al.* (2017).

Address Issue 2: Modeling linkage disequilibrium

Genome-wide multiple-SNP likelihood function:

$$L_{\text{rss}}(\beta) := \text{Normal}(\hat{\beta}; \hat{S}\hat{R}\hat{S}^{-1}\beta, \hat{S}\hat{R}\hat{S})$$

- multiple-SNP parameter: $\beta := (\beta_1, \dots, \beta_p)^\top$
- single-SNP summary data: $\hat{\beta} := (\hat{\beta}_1, \dots, \hat{\beta}_p)^\top$
- $\hat{S} := \text{diag}(\hat{s})$, $\hat{s} := (\hat{s}_1, \dots, \hat{s}_p)^\top$, $\hat{s}_j^2 := \hat{\sigma}_j^2 + n^{-1}\hat{\beta}_j^2 \approx \hat{\sigma}_j^2$
- \hat{R} : the shrinkage estimate of LD (Wen & Stephens, 2010; Li & Stephens, 2003)
- “Big Data” Genetics: **jointly** analyze **1.1** million common SNPs

Reference: Zhu and Stephens (2017 To appear). <http://dx.doi.org/10.1101/042457>.

We develop a method that systematically utilizes enrichment information.

Public Data

Inference

GWAS summary statistics

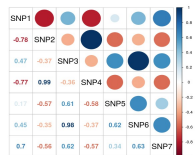


$$\hat{\beta} := (\hat{\beta}_1, \dots, \hat{\beta}_p)^T, \quad \hat{S} := \text{diag}\{(\hat{s}_1, \dots, \hat{s}_p)^T\}$$

$\hat{\beta}_j$: marginal effect estimate of SNP j

\hat{s}_j : standard error of $\hat{\beta}_j$

External LD estimates



$$\tilde{R} := [\hat{r}_{ij}]_{p \times p}, \quad \hat{r}_{ij}: \text{LD between SNP } i \text{ and } j$$

Predefined gene sets

Gene	Chr	Start	End
BRAF	7	140419127	140624564
MAPK1	22	22108789	22221970
MAPK3	16	30125426	30134827
NFKB1	4	103422486	103538459

$$\alpha_j := 1 \{ \text{SNP } j \text{ is "near" a gene in the set} \}$$

Bayesian Model

Likelihood function

$$L(\beta) := \text{Normal}(\hat{\beta}; \hat{S} \hat{R} \hat{S}^{-1} \beta, \hat{S} \hat{R} \hat{S})$$

Prior distribution

$$\beta_j \sim \pi_j \cdot \text{Normal}(0, \sigma_\beta^2) + (1 - \pi_j) \delta_0$$

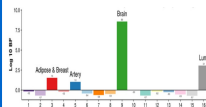
Baseline model (M_0)

$$\log_{10} \left(\frac{\pi_j}{1 - \pi_j} \right) = \theta_0$$

Enrichment model (M_1)

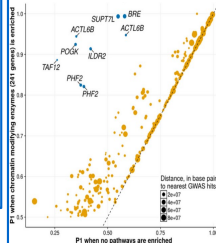
$$\log_{10} \left(\frac{\pi_j}{1 - \pi_j} \right) = \theta_0 + \alpha_j \theta$$

Gene set enrichment



$$\text{BF}(\text{gene set}) = \frac{\Pr(\text{Data} | M_1)}{\Pr(\text{Data} | M_0)}$$

Gene prioritization



$$\rho(\beta | \text{Data}, M_1) \text{ vs } \rho(\beta | \text{Data}, M_0)$$

We apply the method to analyze 31 complex traits and 4,026 gene sets.

This application is not small:

of Parameters = $31 \times (3,913+113) \times 1.1 \text{ Million} \approx 137 \text{ Billion}$

- **31** human phenotypes
- **3,913** biological pathways
- **113** tissue-based gene sets
- **1.1 million** common SNPs

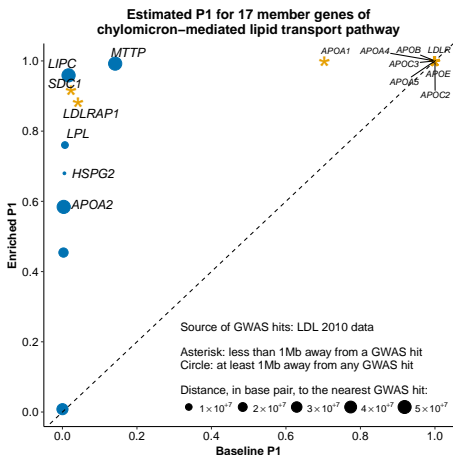
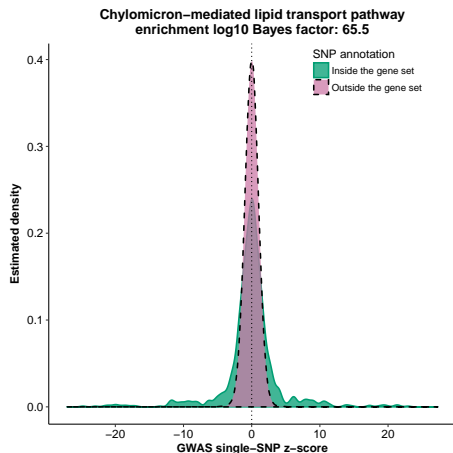
One student can get this done, aided by:

- Publicly available input data: GWAS + LD + gene sets
- Variational inference; Squared extrapolation method (SQUAREM)
- Parallel computing; Hierarchical data format (HDF5)

Our full results are publicly available.

- **Results**
<https://xiangzhu.github.io/rss-gsea/results>
- **Software**
<https://github.com/stephenslab/rss>
- **Demonstration**
<https://stephenslab.github.io/rss/Example-5>
- **R package (in progress)**
<https://github.com/stephenslab/rssr> (N. Knoblauch)

Example: Low-density lipoprotein & *MTTP*



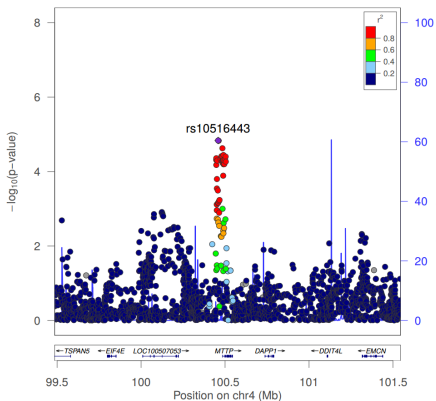
■ Gene detection: $P_1 := 1 - \Pr(\beta_j = 0, \forall j \in \text{locus} \mid \text{Data})$

■ *MTTP*: baseline $P_1 = \mathbf{0.14}$; enriched $P_1 = \mathbf{0.99}$

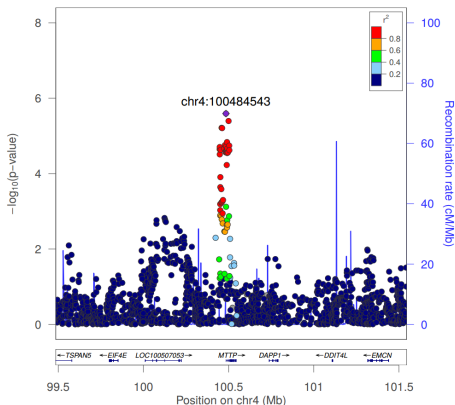
Example: Low-density lipoprotein & *MTTP*

- Total sample size: **95,454** (2010) → **173,082** (2013)
- GWAS *p*-value: **1.5×10^{-5}** (2010) → **2.6×10^{-6}** (2013) “<” 5×10^{-8}

LDL 2010 data, max sample size: 95,454

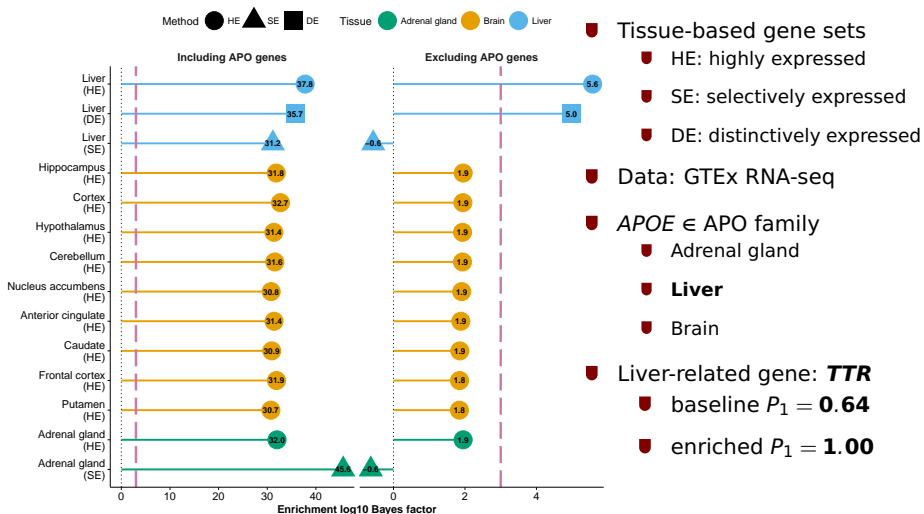


LDL 2013 data, max sample size: 173,082



References: Global Lipids Genetics Consortium (2013); Teslovich *et al.* (2010).

Example: Alzheimer's disease & Liver



References: Xi et al. (Submitted); Dey et al. (2017); The GTEx Consortium (2015); Lambert et al. (2013).

What's next?

We develop a new enrichment analysis method that:

- uses **publicly available** data as input
- **efficiently** assesses thousands of gene sets
- **automatically** identifies trait-associated genes

What do we plan to do next?

- one annotation at a time → jointly analyze **many** annotations
- gene-based annotations → **“finer-scale”** annotations
- spike-and-slab prior → more **flexible** effect size distributions

Acknowledgements

■ Joint work with **Matthew Stephens**

■ **GWAS summary statistics**

Teslovich *et al.* (2010); Manning *et al.* (2012); Morris *et al.* (2012); van der Harst *et al.* (2012); Köttgen *et al.* (2013); Lambert *et al.* (2013); Okada *et al.* (2014); Ripke *et al.* (2014); Wood *et al.* (2014); Day *et al.* (2015); Liu *et al.* (2015); Locke *et al.* (2015); Nikpay *et al.* (2015); Shungin *et al.* (2015); Okbay *et al.* (2016); van Rheenen *et al.* (2016)

■ **GTEx Consortium**

<https://gtexportal.org>

GORDON AND BETTY
MOORE
FOUNDATION



THE UNIVERSITY OF CHICAGO

**RESEARCH
COMPUTING
CENTER**

Thank you!

- **Preprint:** <https://doi.org/10.1101/160770>
- **Software:** <https://github.com/stephenslab/rss>
- **Contact:** xiangzhu [at] uchicago [dot] edu