

2017-04-24

# Methods for Estimating Hidden Structure and Network Transitions in Genomics

Daniel Schlauch, PhD Candidate

Department of Biostatistics  
Harvard School of Public Health

April 24, 2017

## Table of Contents

2017-04-24

### Table of Contents

Table of Contents

- 1 Batch Effect on Covariance Structure
- 2 Identification of genetic outliers
- 3 State Transitions Using Gene Regulatory Network Models



# Batch effect on covariance structure confounds gene coexpression

Daniel Schlauch<sup>1,2</sup>, Joseph N. Paulson<sup>2</sup>, Kimberly Glass<sup>2,3</sup>, and John Quackenbush<sup>1,3</sup>

<sup>1</sup>Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute and Department of Biostatistics, Harvard TH Chan School of Public Health, Boston, MA

<sup>2</sup>Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA

<sup>3</sup>Department of Medicine, Harvard Medical School, Boston, MA.

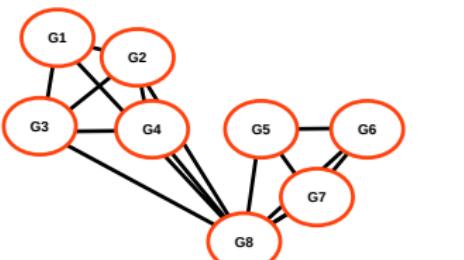
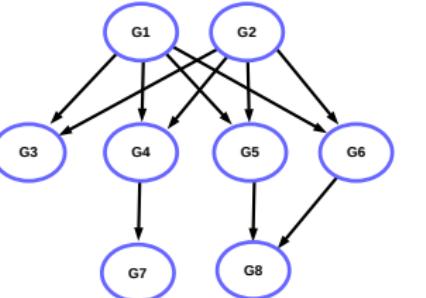
Department of Medicine, Harvard Medical School, Boston, MA



## Background: Differential Network Inference

How do we model functional interactions?

- Gene Regulatory/Coexpression Networks (GRN/GCN)
- Directed/undirected graph
- May imply a sort of physical interaction
- Guilt by association

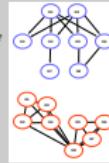


Methods for Estimating Hidden Structure and Network Transitions in Genomics

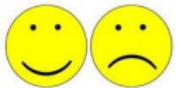
- Batch Effect on Covariance Structure
  - Gene Networks
  - Background: Differential Network Inference

2017-04-24

Background: Differential Network Inference

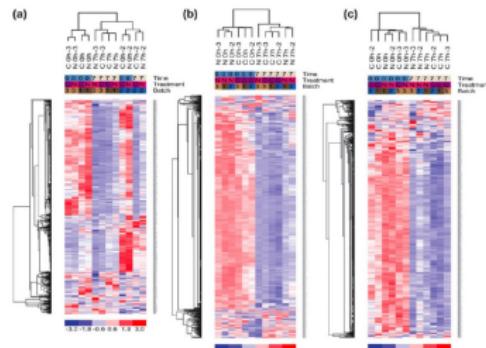


## Background: Batch Effect



Batches

Experimental Conditions



Johnson et al. (Biostatistics 2007)

Unwanted variation from:

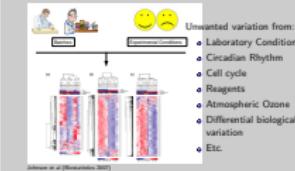
- Laboratory Conditions
- Circadian Rhythm
- Cell cycle
- Reagents
- Atmospheric Ozone
- Differential biological variation
- Etc.

Methods for Estimating Hidden Structure and Network Transitions in Genomics

- └ Batch Effect on Covariance Structure
- └ Controlling for Batch Effect
- └ Background: Batch Effect

2017-04-24

Background: Batch Effect



# Methods for Controlling Batch Effect

Location scale model:

$$Y_{ijg} = \alpha_g + X\beta_g + \gamma_{ig} + \delta_{ig}\epsilon_{ijg}$$

Limitations:

- Gene-specific location/scale assumptions
- Independent effects
- *Differential coexpression*

Batch effect removal methods typically return a corrected gene expression matrix (e.g. ComBat) or a correction vector (e.g. SVA).

2017-04-24

Methods for Estimating Hidden Structure and Network Transitions in Genomics  
└ Batch Effect on Covariance Structure  
  └ Controlling for Batch Effect  
    └ Methods for Controlling Batch Effect

Methods for Controlling Batch Effect

Location scale model:

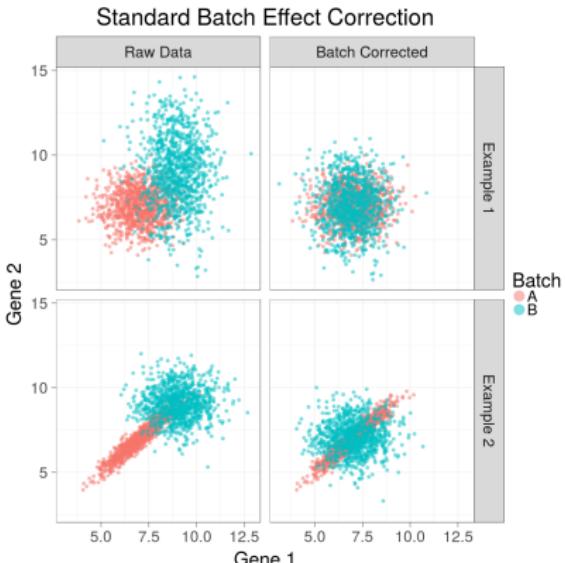
$$Y_{ijg} = \alpha_g + X\beta_g + \gamma_{ig} + \delta_{ig}\epsilon_{ijg}$$

Limitations:

- Gene-specific location/scale assumptions
- Independent effects
- Differential coexpression

Batch effect removal methods typically return a corrected gene expression matrix (e.g. ComBat) or a correction vector (e.g. SVA).

# Limitations to common batch effect correction methods



## Standard corrections:

$$f[Gene1|BatchA] = f[Gene1|BatchB]$$

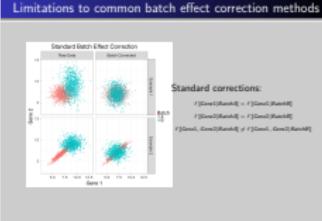
$$f[Gene2|BatchA] = f[Gene2|BatchB]$$

$$f[Gene1, Gene2|BatchA] \neq f[Gene1, Gene2|BatchB]$$

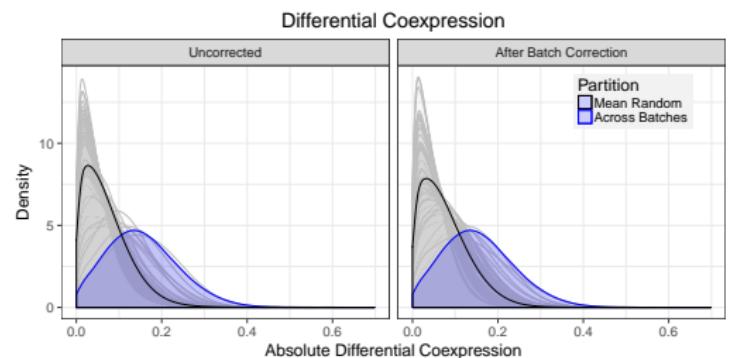
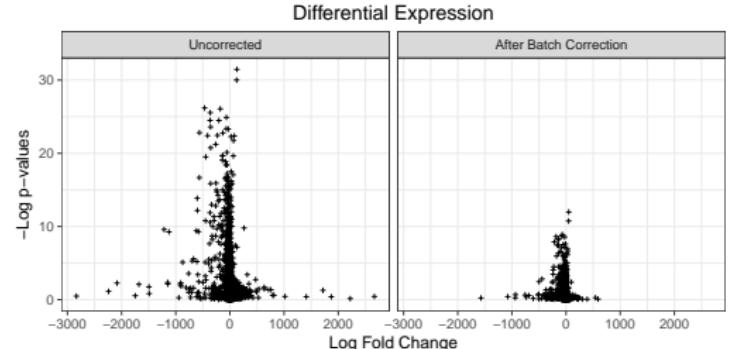
## Batch Effect on Covariance Structure

### Controlling for Batch Effect

#### Limitations to common batch effect correction methods



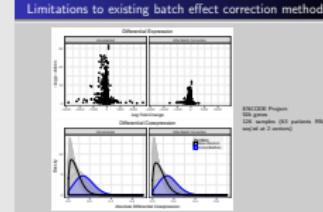
# Limitations to existing batch effect correction methods



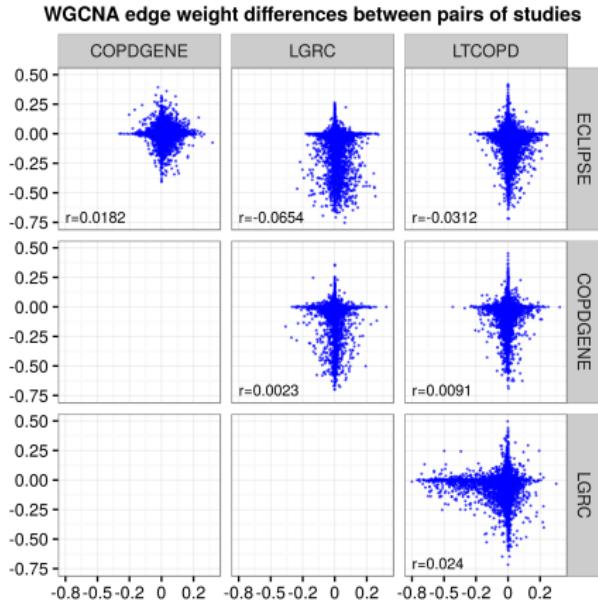
ENCODE Project:  
50k genes  
126 samples (63 patients RNA-seq'ed at 2 centers)

## Methods for Estimating Hidden Structure and Network Transitions in Genomics

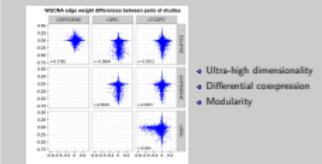
- Batch Effect on Covariance Structure
- Controlling for Batch Effect
- Limitations to existing batch effect correction methods



# Challenges with batch effect on differential coexpression



- Ultra-high dimensionality
- Differential coexpression
- Modularity



# Estimating the conditional coexpression matrix

2017-04-24

Motivating concepts:

- Provide a regression framework for the coexpression matrix
- Estimate a reduced number of parameters
- Exploit modular nature of gene expression patterns

Our proposal:

- Define our parameters as functions of components of variation.
- Estimate the eigenvalue contribution of each eigenvector.



## Model

Consider a set of  $N$  samples with  $q$  covariates measuring gene expression across  $p$  genes. Let  $\mathbf{x}_i = (x_{i1}, \dots, x_{iq})$  denote the covariates for sample  $i$  and let  $\mathbf{g}_i = (g_{i1}, \dots, g_{ip})^T$  denote the gene expression values for sample  $i$  for the  $p$  genes.

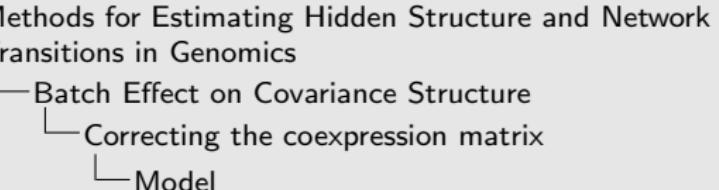
We can express a model for the gene expression as

$$\mathbf{g}_i = \beta^T \mathbf{x}_i + \epsilon_i \text{ for } i = 1, \dots, N$$

where  $\epsilon_i \sim MVN_p(\mathbf{0}, \Sigma_i)$ . Notably, the covariance of  $\epsilon_i$  differ according to  $i$ .

$$\Sigma_i = \mathbf{Q} \mathbf{D}_i \mathbf{Q}^T$$

where  $\mathbf{D}_i$  is a diagonal matrix with diagonal defined as  $\mathbf{X}_i \Psi_{q \times p}$



2017-04-24

Model

Consider a set of  $N$  samples with  $q$  covariates measuring gene expression across  $p$  genes. Let  $\mathbf{x}_i = (x_{i1}, \dots, x_{iq})$  denote the covariates for sample  $i$  and let  $\mathbf{g}_i = (g_{i1}, \dots, g_{ip})^T$  denote the gene expression values for sample  $i$  for the  $p$  genes. We can express a model for the gene expression as

$$\mathbf{g}_i = \beta^T \mathbf{x}_i + \epsilon_i \text{ for } i = 1, \dots, N$$

where  $\epsilon_i \sim MVN_p(\mathbf{0}, \Sigma_i)$ . Notably, the covariance of  $\epsilon_i$  differ according to  $i$ .

$$\Sigma_i = \mathbf{Q} \mathbf{D}_i \mathbf{Q}^T$$

where  $\mathbf{D}_i$  is a diagonal matrix with diagonal defined as  $\mathbf{X}_i \Psi_{q \times p}$

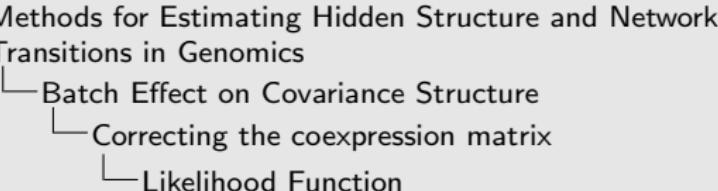
# Likelihood Function

$$\mathcal{L}(\mu, \Sigma) = \prod_{i=1}^N \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma_i|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{G}_i - \mu)^T \Sigma_i^{-1} (\mathbf{G}_i - \mu)}$$

Where we define  $\Sigma_i$ ,

$$\Sigma_i = \mathbf{Q} \mathbf{D}_i \mathbf{Q}^T$$

Where  $\mathbf{Q}$  is a matrix with columns defined as the eigenvectors of the estimated coexpression matrix,  $\mathbf{G}^* \mathbf{G}^{*T} / N$ .



2017-04-24

- point 1
- point 2

Likelihood Function

$$\mathcal{L}(\mu, \Sigma) = \prod_{i=1}^N \frac{1}{(2\pi)^{\frac{p}{2}} |\Sigma_i|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{G}_i - \mu)^T \Sigma_i^{-1} (\mathbf{G}_i - \mu)}$$

Where we define  $\Sigma_i$

$$\Sigma_i = \mathbf{Q} \mathbf{D}_i \mathbf{Q}^T$$

Where  $\mathbf{Q}$  is a matrix with columns defined as the eigenvectors of the estimated coexpression matrix,  $\mathbf{G}^* \mathbf{G}^{*T} / N$ .



# Least Squares Estimator

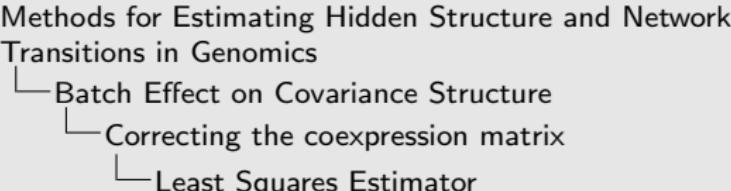
To calculate least squares solution,  $\hat{\Psi}$ , we solve separately for each column of  $\mathbf{Q}$ .

and note that the residual matrices should be orthogonal to the hyperplane spanned by  $X^T$ .

Recall  $0 = \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\hat{\beta})$

$$\mathbf{0}_q = \sum_{i=1}^N \mathbf{X}_i^T \left[ \mathbf{Q}_h^T \left[ \mathbf{G}_i^* \mathbf{G}_i^{*T} - \mathbf{Q}_h \mathbf{X}_i \hat{\Psi}_h \mathbf{Q}_h^T \right] \mathbf{Q}_h \right] \quad (1)$$

$$\hat{\Psi}_h = (\mathbf{X}^T \mathbf{X})^{-1} \sum_{i=1}^N [\mathbf{X}_i^T \mathbf{Q}_h^T \mathbf{G}_i^* \mathbf{G}_i^{*T} \mathbf{Q}_h] \quad (2)$$



2017-04-24

## Least Squares Estimator

To calculate least squares solution,  $\hat{\Psi}$ , we solve separately for each column of  $\mathbf{Q}$  and note that the residual matrices should be orthogonal to the hyperplane spanned by  $X^T$ .  
Recall  $0 = \mathbf{X}^T (\mathbf{Y} - \mathbf{X}\hat{\beta})$

$$\mathbf{0}_q = \sum_{i=1}^N \mathbf{X}_i^T \left[ \mathbf{Q}_h^T \left[ \mathbf{G}_i^* \mathbf{G}_i^{*T} - \mathbf{Q}_h \mathbf{X}_i \hat{\Psi}_h \mathbf{Q}_h^T \right] \mathbf{Q}_h \right] \quad (1)$$

$$\hat{\Psi}_h = (\mathbf{X}^T \mathbf{X})^{-1} \sum_{i=1}^N [\mathbf{X}_i^T \mathbf{Q}_h^T \mathbf{G}_i^* \mathbf{G}_i^{*T} \mathbf{Q}_h] \quad (2)$$

- point 1
- point 2



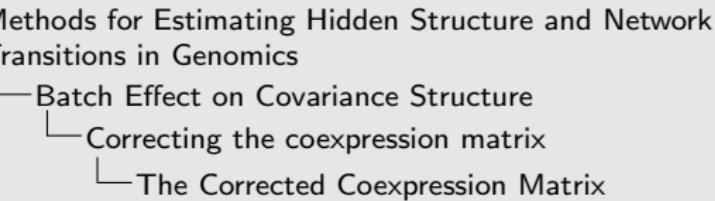
# The Corrected Coexpression Matrix

With the estimates obtained with our method, it is straightforward to see how fitted values for the coexpression matrix for each sample or experimental condition can be obtained. Given an estimate for  $\Psi$ ,  $\hat{\Psi}$ , we can now estimate the batch-independent coexpression structure as

$$\hat{\mathbf{S}} = \mathbf{Q} \text{diag}(\bar{\mathbf{X}}\hat{\Psi}) \mathbf{Q}^T \text{ or } \hat{\mathbf{S}} = \sum_{i=1}^p \bar{\mathbf{X}}\hat{\Psi}_i \mathbf{Q}_i \mathbf{Q}_i^T$$

The differential coexpression matrix between two conditions, defined in binary as column 2 of  $\mathbf{X}$ , is computed

$$\hat{\mathbf{W}} = \mathbf{Q} \text{diag}(\hat{\Psi}_{2,.}) \mathbf{Q}^T$$



2017-04-24

The Corrected Coexpression Matrix

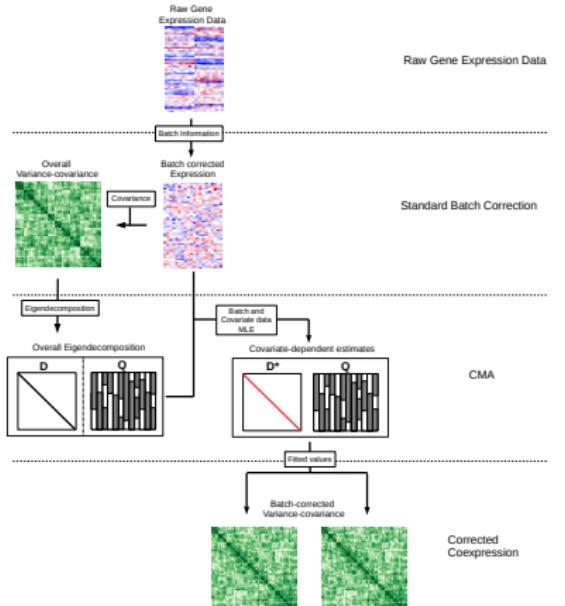
With the estimates obtained with our method, it is straightforward to see how fitted values for the coexpression matrix for each sample or experimental condition can be obtained. Given an estimate for  $\Psi$ ,  $\hat{\Psi}$ , we can now estimate the batch-independent coexpression structure as

$$\hat{\mathbf{S}} = \mathbf{Q} \text{diag}(\bar{\mathbf{X}}\hat{\Psi}) \mathbf{Q}^T \text{ or } \hat{\mathbf{S}} = \sum_{i=1}^p \bar{\mathbf{X}}\hat{\Psi}_i \mathbf{Q}_i \mathbf{Q}_i^T$$

The differential coexpression matrix between two conditions, defined in binary as column 2 of  $\mathbf{X}$ , is computed

$$\hat{\mathbf{W}} = \mathbf{Q} \text{diag}(\hat{\Psi}_{2,.}) \mathbf{Q}^T$$

# Example Workflow

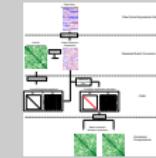


Methods for Estimating Hidden Structure and Network Transitions in Genomics

- Batch Effect on Covariance Structure
- Correcting the coexpression matrix
- Example Workflow

2017-04-24

Example Workflow



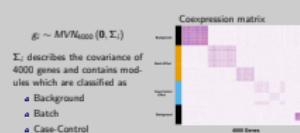
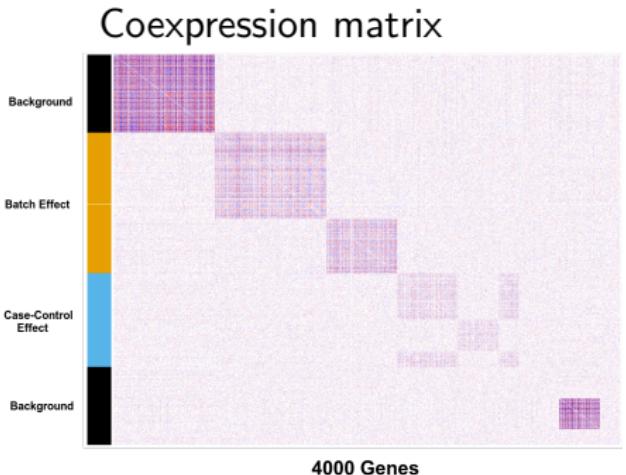
## Simulations

2017-04-24

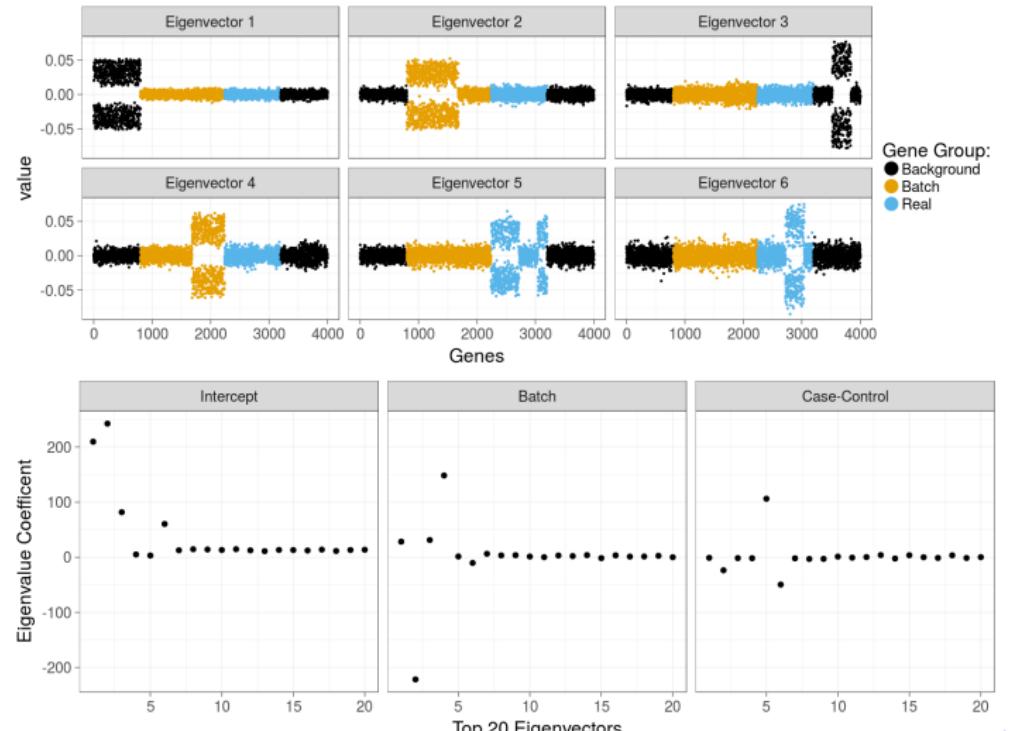
$$g_i \sim MVN_{4000}(\mathbf{0}, \Sigma_i)$$

$\Sigma_i$  describes the covariance of 4000 genes and contains modules which are classified as

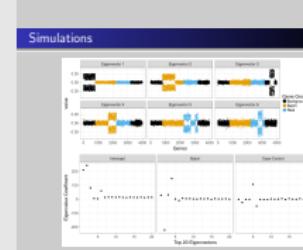
- Background
- Batch
- Case-Control



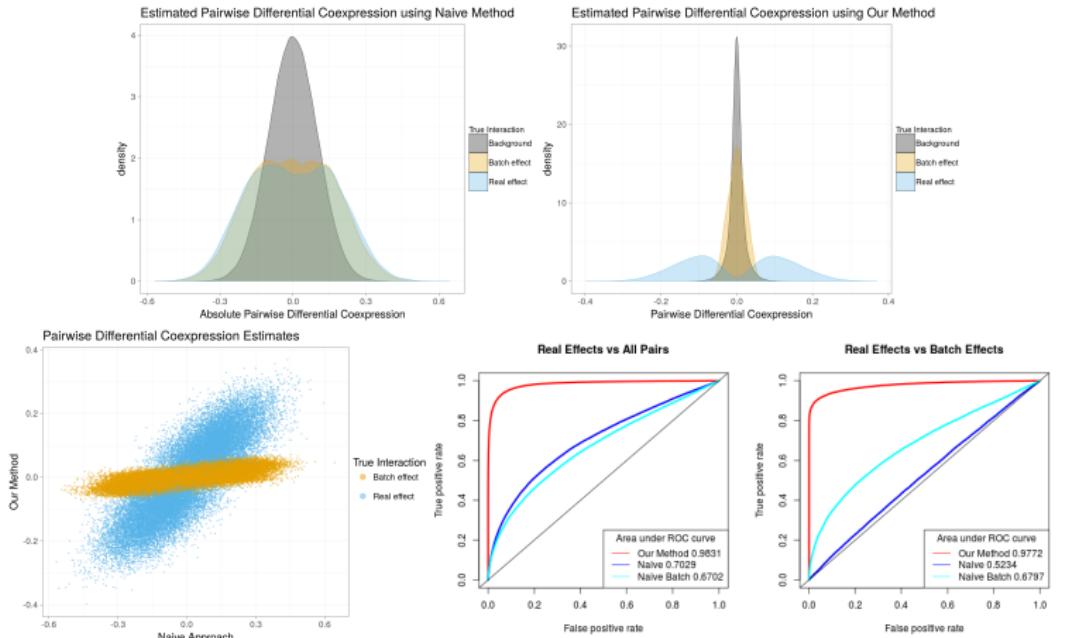
# Simulations



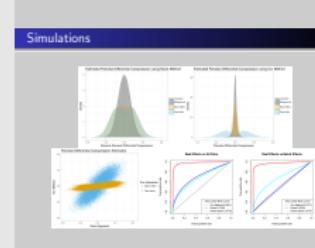
2017-04-24



# Simulations



2017-04-24



# Application to data from COPDGene Study

| GO Term  | Count | %    | Enrichment | FDR      |
|--|-------|------|------------|----------|
| anatomical structure development               | 309   | 0.26 | 1.29       | 2.58E-05 |
| single-organism developmental process          | 309   | 0.26 | 1.29       | 2.73E-05 |
| anatomical structure morphogenesis             | 168   | 0.14 | 1.46       | 1.60E-04 |
| single-multicellular organism process          | 324   | 0.27 | 1.25       | 4.01E-04 |
| system process                                 | 132   | 0.11 | 1.50       | 1.46E-03 |
| regulation of cellular process                 | 514   | 0.43 | 1.12       | 5.86E-03 |
| single organism signaling                      | 328   | 0.28 | 1.21       | 7.89E-03 |
| regulation of localization                     | 151   | 0.13 | 1.40       | 1.15E-02 |
| regulation of multicellular organismal process | 156   | 0.13 | 1.35       | 6.56E-02 |

Table : GO categories for differential coexpression in COPDGene identified with CPBA found with FDR<0.1.

2017-04-24

| GO Term  | Count | %    | Enrichment | FDR      |
|--|-------|------|------------|----------|
| anatomical structure development               | 309   | 0.26 | 1.29       | 2.58E-05 |
| single-organism developmental process          | 309   | 0.26 | 1.29       | 2.73E-05 |
| anatomical structure morphogenesis             | 168   | 0.14 | 1.46       | 1.60E-04 |
| single-multicellular organism process          | 324   | 0.27 | 1.25       | 4.01E-04 |
| system process                                 | 132   | 0.11 | 1.50       | 1.46E-03 |
| regulation of cellular process                 | 514   | 0.43 | 1.12       | 5.86E-03 |
| single organism signaling                      | 328   | 0.28 | 1.21       | 7.89E-03 |
| regulation of localization                     | 151   | 0.13 | 1.40       | 1.15E-02 |
| regulation of multicellular organismal process | 156   | 0.13 | 1.35       | 6.56E-02 |

Table : GO categories for differential coexpression in COPDGene identified with CPBA found with FDR<0.1.

# Identifying Genetic Outliers

## Identification of genetic outliers due to sub-structure and cryptic relationships

Daniel Schlauch<sup>1,2,4</sup>, Heide Fier<sup>1,3</sup> and Christoph Lange<sup>1,4</sup>

<sup>1</sup>Department of Biostatistics, Harvard TH Chan School of Public Health, Boston, MA 02115

<sup>2</sup>Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute, Boston, MA 02115

<sup>3</sup>Institute of Genomic Mathematics, University of Bonn, Bonn, Germany

<sup>4</sup>Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA 02115

## Methods for Estimating Hidden Structure and Network Transitions in Genomics

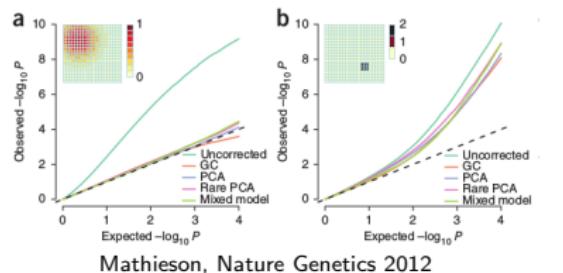
- └ Identification of genetic outliers

- └ Identifying Genetic Outliers



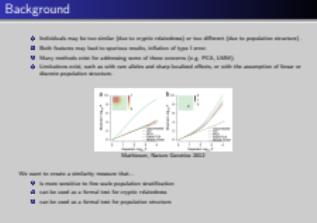
# Background

- Individuals may be too similar (due to cryptic relatedness) or too different (due to population structure).
- Both features may lead to spurious results, inflation of type I error.
- Many methods exist for addressing some of these concerns (e.g. PCA, LMM).
- Limitations exist, such as with rare alleles and sharp localized effects, or with the assumption of linear or discrete population structure.



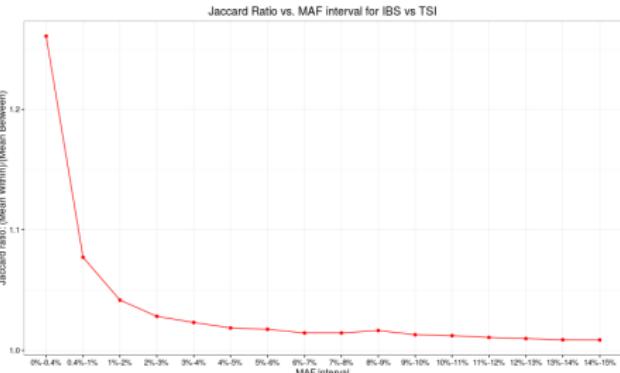
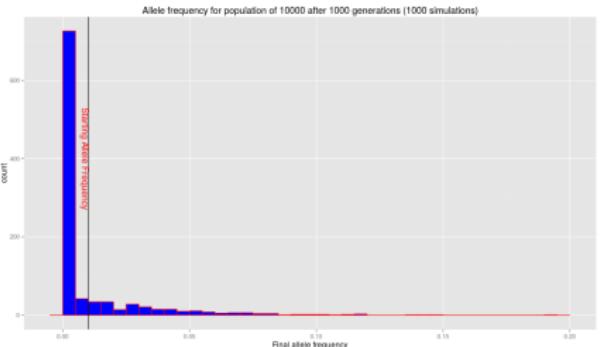
We want to create a similarity measure that...

- is more sensitive to fine scale population stratification
- can be used as a formal test for cryptic relatedness
- can be used as a formal test for population structure



## Basis for measure

- Rare variants are recent variants.
- In the absence of selection, rare variants become fixed at 0% with high probability over a relatively short timeframe.

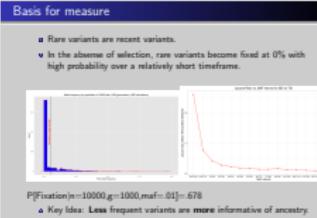
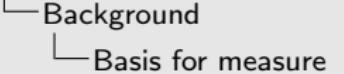


$$P[\text{Fixation} | n=10000, g=1000, \text{maf}=.01] = .678$$

- Key Idea: **Less frequent variants are more informative of ancestry.**



### Identification of genetic outliers



# Test Statistic

$$s_{i,j} = \frac{\sum_{k=1}^N w_k \mathbf{G}_{i,k} \mathbf{G}_{j,k}}{\sum_{k=1}^N I \left[ \sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1 \right]}$$

where

$$w_k = \begin{cases} \frac{\binom{2n}{2}}{\binom{\sum_{l=1}^{2n} \mathbf{G}_{l,k}}{2}} & \sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1 \\ 0 & \sum_{l=1}^{2n} \mathbf{G}_{l,k} \leq 1 \end{cases}$$

$$E [s_{i,j}] = 1$$

## Methods for Estimating Hidden Structure and Network Transitions in Genomics

- Identification of genetic outliers

- Similarity measure
- Test Statistic

### Test Statistic

$$s_{i,j} = \frac{\sum_{k=1}^N w_k \mathbf{G}_{i,k} \mathbf{G}_{j,k}}{\sum_{k=1}^N I \left[ \sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1 \right]}$$

where

$$w_k = \begin{cases} \frac{\binom{2n}{2}}{\binom{\sum_{l=1}^{2n} \mathbf{G}_{l,k}}{2}} & \sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1 \\ 0 & \sum_{l=1}^{2n} \mathbf{G}_{l,k} \leq 1 \end{cases}$$

$$E [s_{i,j}] = 1$$



## Test Statistic

In the absence of population structure, cryptic relatedness and dependence between loci the distribution of the similarity index,  $s_{i,j}$

$$s_{i,j} \sim N(1, \sigma_{i,j}^2)$$

Where the variance of  $s_{ij}$  can be estimated by

$$\hat{\sigma}_{i,j}^2 = \hat{Var}(s_{i,j}) = \frac{\sum_{k=1}^N (w_k - 1)}{\left(\sum_{k=1}^N I\left[\sum_{l=1}^{2n} \mathbf{G}_{l,k} > 1\right]\right)^2}$$

$$s_{i,j}^{(diploid)} = \frac{\sum_{k=1}^N [w_k \mathbf{H}_{i,k} \mathbf{H}_{j,k}] / 4}{\sum_{k=1}^N I[(\sum_{l=1}^n \mathbf{H}_{l,k}) > 1]}$$

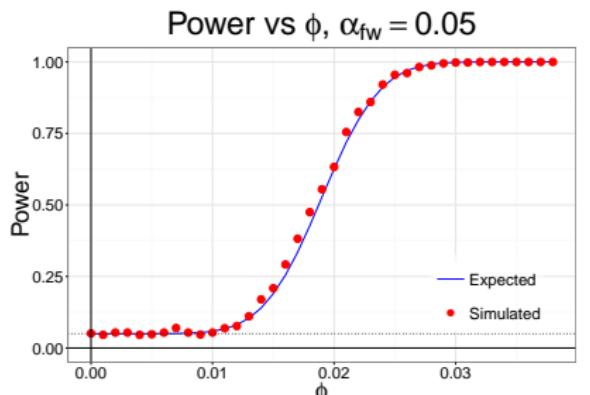


# Tests of Heterogeneity

$$\hat{\phi}_{i,j} = \frac{s_{i,j} - 1}{\left[ \frac{\sum_{k=1}^N \hat{\rho}_k w_k}{\sum_{k=1}^N I[\sum_{l=1}^{2n} G_{l,k} > 1]} - 1 \right]}$$

$$R : \max(s_{i,j}) > 1 - \text{probit} \left( \frac{\alpha}{\binom{n}{2}} \right)$$

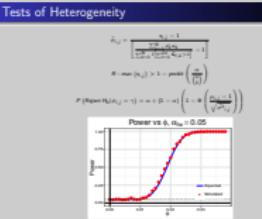
$$P(\text{Reject } H_0 | \phi_{i,j} = \gamma) = \alpha + (1 - \alpha) \left( 1 - \Phi \left( \frac{\mu_{i,j} - 1}{\sqrt{\sigma^2_{i,j}}} \right) \right)$$



Methods for Estimating Hidden Structure and Network  
Transitions in Genomics

- Identification of genetic outliers
- Similarity measure
- Tests of Heterogeneity

2017-04-24



# Tests of Heterogeneity

$$H_0 : \mu_{i,j} = 1 \forall i, j \in 1 \dots n$$

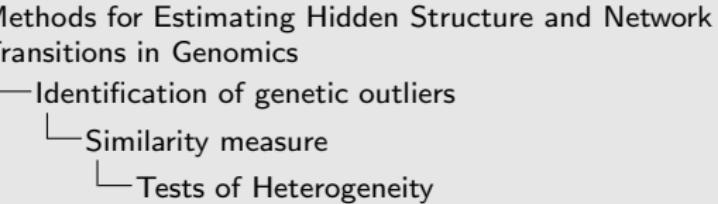
$$H_A : \exists i, j \in 1 \dots n | \mu_{i,j} \neq 1$$

Test for population structure:

$$K = \sup_x |F_s(x) - \Phi(x)|$$

Test for cryptic relatedness:

$$R : \max(s_{i,j}) > 1 - \text{probit} \left( \frac{\alpha}{\binom{n}{2}} \right)$$



2017-04-24

Tests of Heterogeneity

$$\begin{aligned} H_0 : \mu_{i,j} &= 1 \forall i, j \in 1 \dots n \\ H_A : \exists i, j \in 1 \dots n | \mu_{i,j} &\neq 1 \end{aligned}$$

Test for population structure:

$$K = \sup_x |F_s(x) - \Phi(x)|$$

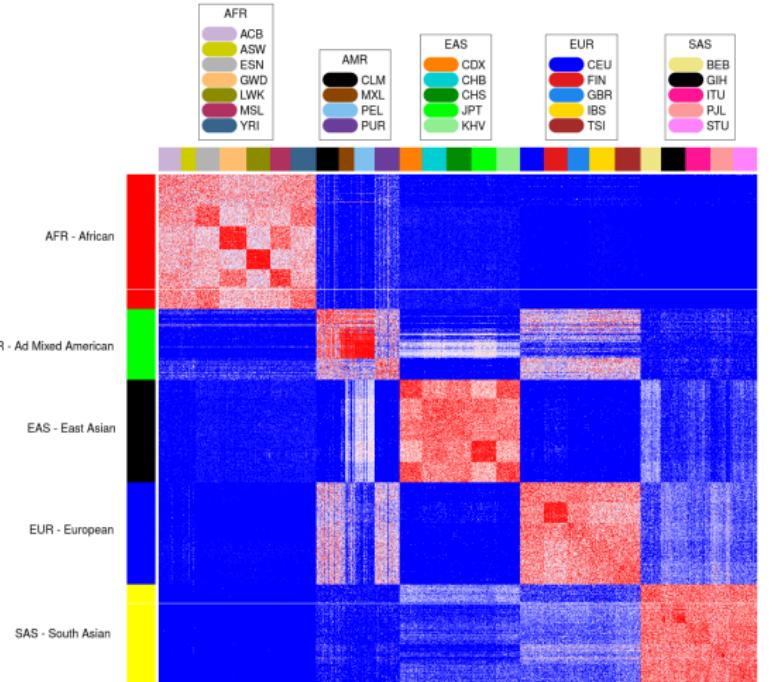
Test for cryptic relatedness:

$$R : \max(s_{i,j}) > 1 - \text{probit} \left( \frac{\alpha}{\binom{n}{2}} \right)$$



# Application to 1000 Genomes Project

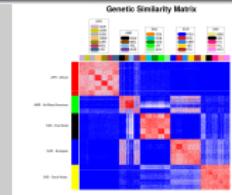
## Genetic Similarity Matrix



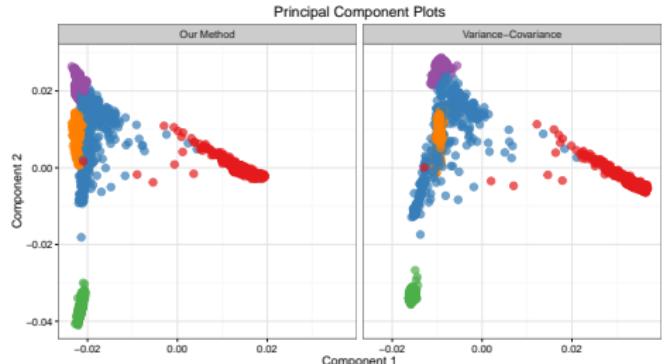
Methods for Estimating Hidden Structure and Network  
Transitions in Genomics  
└ Identification of genetic outliers  
└ Results  
└ Application to 1000 Genomes Project

2017-04-24

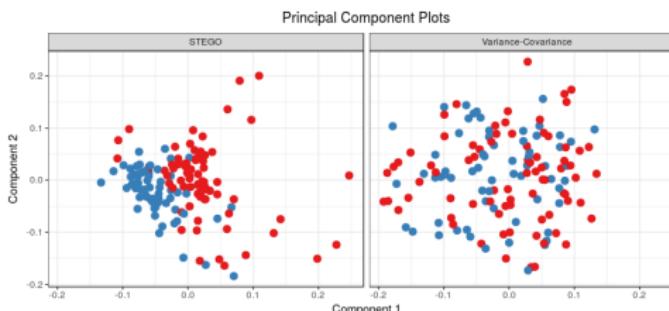
Application to 1000 Genomes Project



# Application to 1000 Genomes Project



STEGO is comparable  
to PCA when applied  
on a global scale.

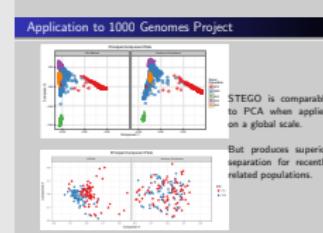


But produces superior  
separation for recently  
related populations.

Methods for Estimating Hidden Structure and Network  
Transitions in Genomics

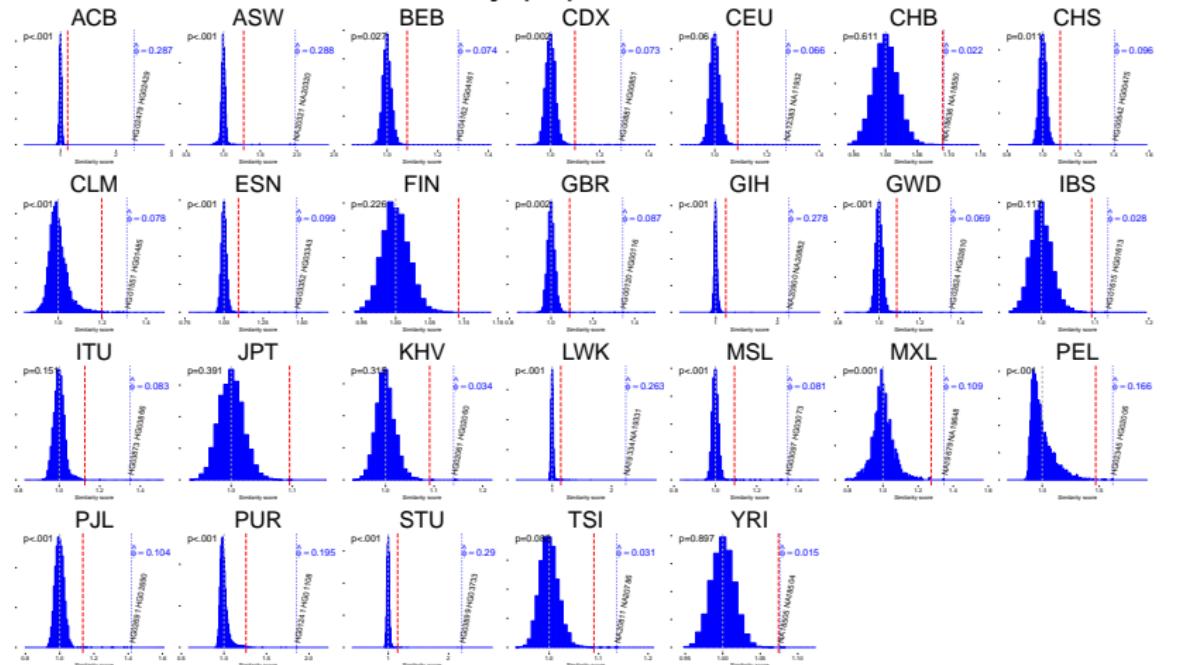
- Identification of genetic outliers
- Results
  - Application to 1000 Genomes Project

2017-04-24



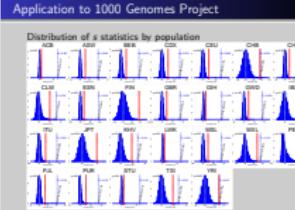
# Application to 1000 Genomes Project

## Distribution of $s$ statistics by population



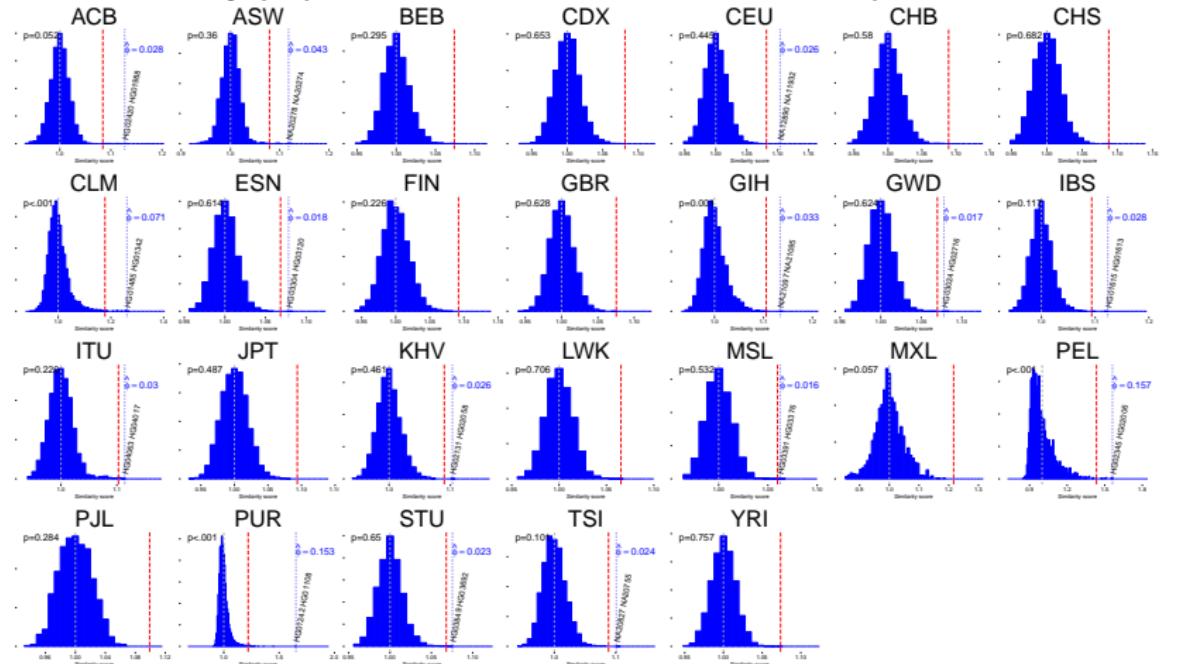
Methods for Estimating Hidden Structure and Network  
Transitions in Genomics  
└ Identification of genetic outliers  
└ Results  
└ Application to 1000 Genomes Project

2017-04-24



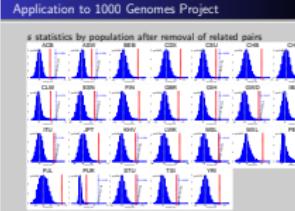
# Application to 1000 Genomes Project

s statistics by population after removal of related pairs



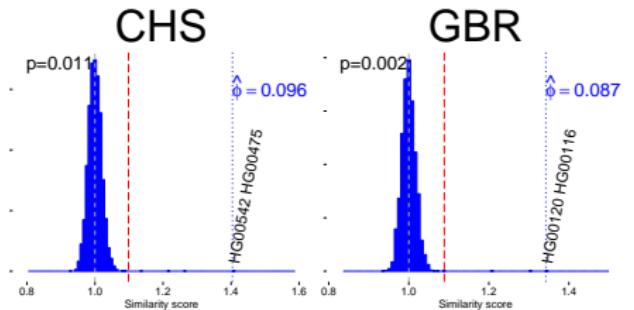
Methods for Estimating Hidden Structure and Network  
Transitions in Genomics  
└ Identification of genetic outliers  
└ Results  
└ Application to 1000 Genomes Project

2017-04-24

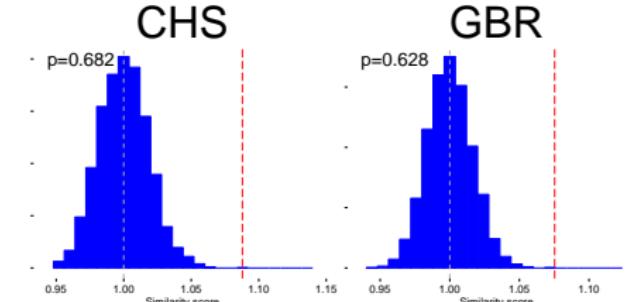


# Application to 1000 Genomes Project

Original Data

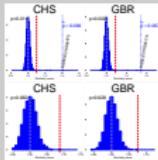


After removal of  
related pairs

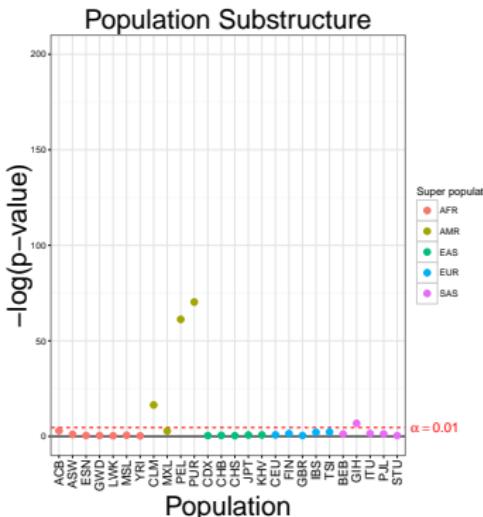
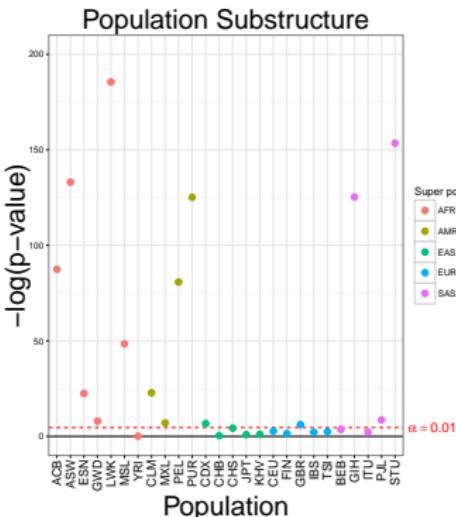


Original Data

After removal of  
related pairs

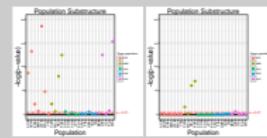


# Application to 1000 Genomes Project



Methods for Estimating Hidden Structure and Network  
Transitions in Genomics

- Identification of genetic outliers
- Results
  - Application to 1000 Genomes Project



# State Transitions Using Gene Regulatory Network Models

Methods for Estimating Hidden Structure and Network  
Transitions in Genomics  
└ State Transitions Using Gene Regulatory Network Models  
└ State Transitions Using Gene Regulatory Network  
Models

2017-04-24

State Transitions Using Gene Regulatory Network Models

Estimating Drivers of Cell State Transitions  
Using Gene Regulatory Network Models

Daniel Schlauch<sup>1,2</sup>, Kimberly Glass<sup>2,3</sup>, Craig P. Hersh<sup>2,3,4</sup>, Edwin  
K. Silverman<sup>2,3,4</sup> and John Quackenbush<sup>1,2,3</sup>

## Estimating Drivers of Cell State Transitions Using Gene Regulatory Network Models

Daniel Schlauch<sup>1,2</sup>, Kimberly Glass<sup>2,3</sup>, Craig P. Hersh<sup>2,3,4</sup>, Edwin  
K. Silverman<sup>2,3,4</sup> and John Quackenbush<sup>1,2,3</sup>

<sup>1</sup>Department of Biostatistics and Computational Biology, Dana-Farber Cancer Institute and Department of  
Biostatistics, Harvard TH Chan School of Public Health, Boston, MA

<sup>2</sup>Channing Division of Network Medicine, Brigham and Women's Hospital, Boston, MA

<sup>3</sup>Department of Medicine, Harvard Medical School, Boston, MA

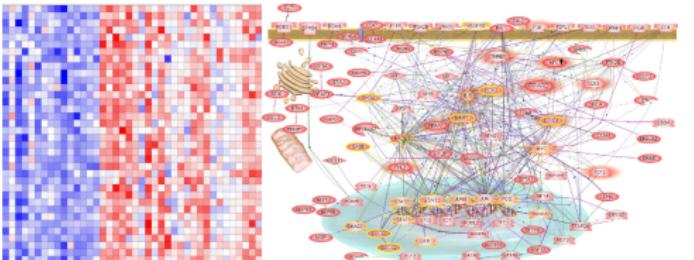
<sup>4</sup>Pulmonary and Critical Care Division, Brigham and Women's Hospital and Harvard Medical School, Boston, MA



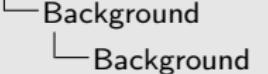
## Background

### Why Study Gene Regulatory Networks?

- Genes are not independent objects.
- Regulation of higher level pathways and processes.

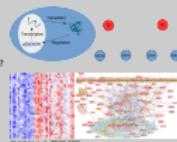


### State Transitions Using Gene Regulatory Network Models



2017-04-24

- point 1
- point 2



## Background

### Biological Challenges

- Measurements of gene expression are at the mRNA level.
- Measurements only consist of mRNA abundance.
- Experimental data is collected as static snapshots.
- Biological variability can be difficult to induce

### Statistical Challenges

- Gene expression measurements are noisy.
- Model complexity may require the estimate of too many model parameters.
- May be computationally intractable.
- May be statistically undetermined. “The curse of dimensionality”

2017-04-24



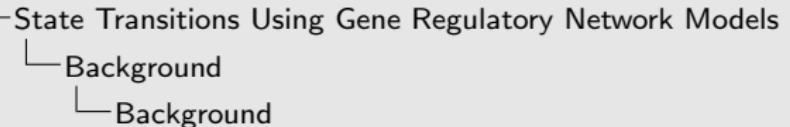
## Background

### Biological Challenges

- Measurements of gene expression are at the mRNA level.
- Measurements only consist of mRNA abundance.
- Experimental data is collected as static snapshots.
- Biological variability can be difficult to induce

### Statistical Challenges

- Gene expression measurements are noisy.
- Model complexity may require the estimate of too many model parameters.
- May be computationally intractable.
- May be statistically undetermined. “The curse of dimensionality”



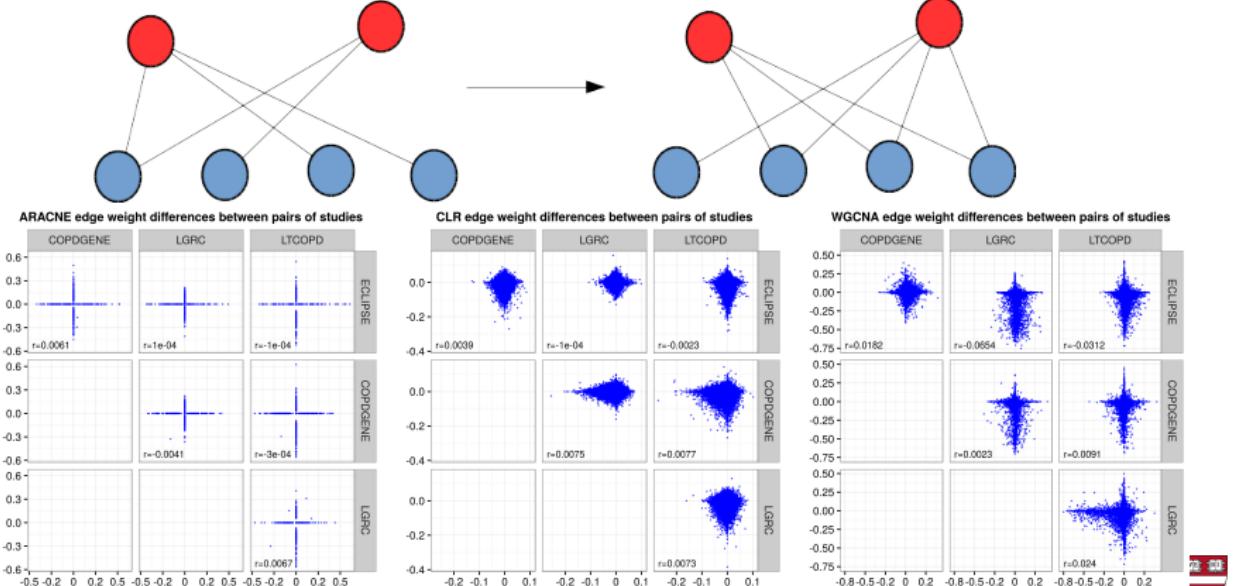
2017-04-24

- Biological Challenges
- Measurements of gene expression are at the mRNA level.
  - Measurements only consist of mRNA abundance.
  - Experimental data is collected as static snapshots.
  - Biological variability can be difficult to induce
- Statistical Challenges
- Gene expression measurements are noisy.
  - Model complexity may require the estimate of too many model parameters.
  - May be computationally intractable.
  - May be statistically undetermined. “The curse of dimensionality”

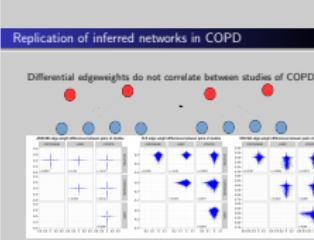


# Replication of inferred networks in COPD

Differential edgeweights do not correlate between studies of COPD



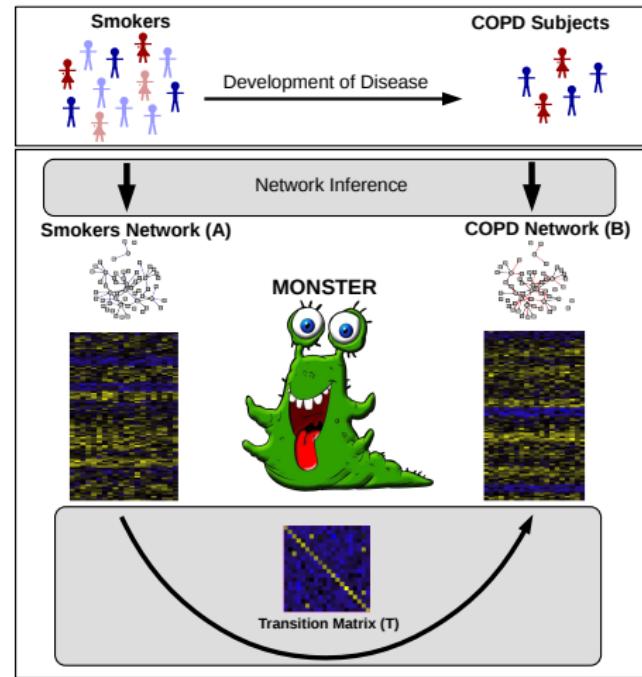
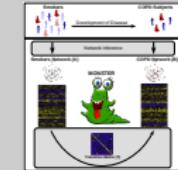
Methods for Estimating Hidden Structure and Network Transitions in Genomics  
State Transitions Using Gene Regulatory Network Models  
Background  
Replication of inferred networks in COPD



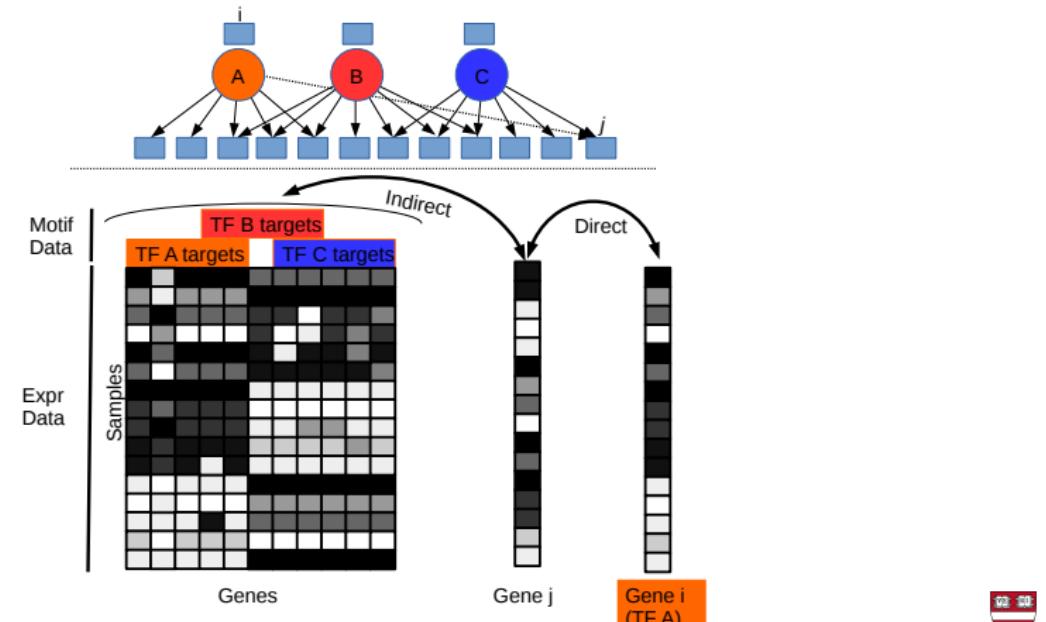
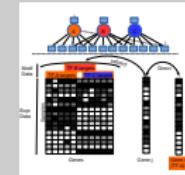
# Algorithm Overview

Methods for Estimating Hidden Structure and Network  
Transitions in Genomics  
└ State Transitions Using Gene Regulatory Network Models  
  └ Network Inference  
    └ Algorithm Overview

Algorithm Overview



# Network Inference



## Network Inference

Direct Evidence:

$$d_{i,j} = \text{cor}(g_i, g_j | \{g_{k,-i} : k \neq i, k \in \text{TF}\})^2$$

Indirect Evidence:

$$\text{logit}(E[M_i]) = \beta_0 + \beta_1 g_{(1)} + \cdots + \beta_N g_{(n)}$$

$$e_{i,j} = \frac{1}{1 + e^{\beta_0 + \beta_1 g_{j,(1)} + \cdots + \beta_k g_{j,(k)}}}$$

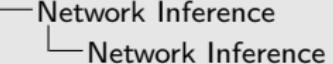
Edgeweight:

$$w_{i,j} = (1 - \alpha) [d_{i,j}] + \alpha [e_{i,j}]$$

## Methods for Estimating Hidden Structure and Network

### Transitions in Genomics

#### State Transitions Using Gene Regulatory Network Models



2017-04-24

## Network Inference

### Direct Evidence:

$$d_{i,j} = \text{cor}(g_i, g_j | \{g_{k,-i} : k \neq i, k \in \text{TF}\})^2$$

### Indirect Evidence:

$$\text{logit}(E[M_i]) = \beta_0 + \beta_1 g_{(1)} + \cdots + \beta_N g_{(n)}$$

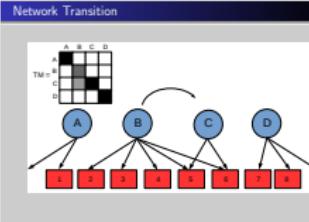
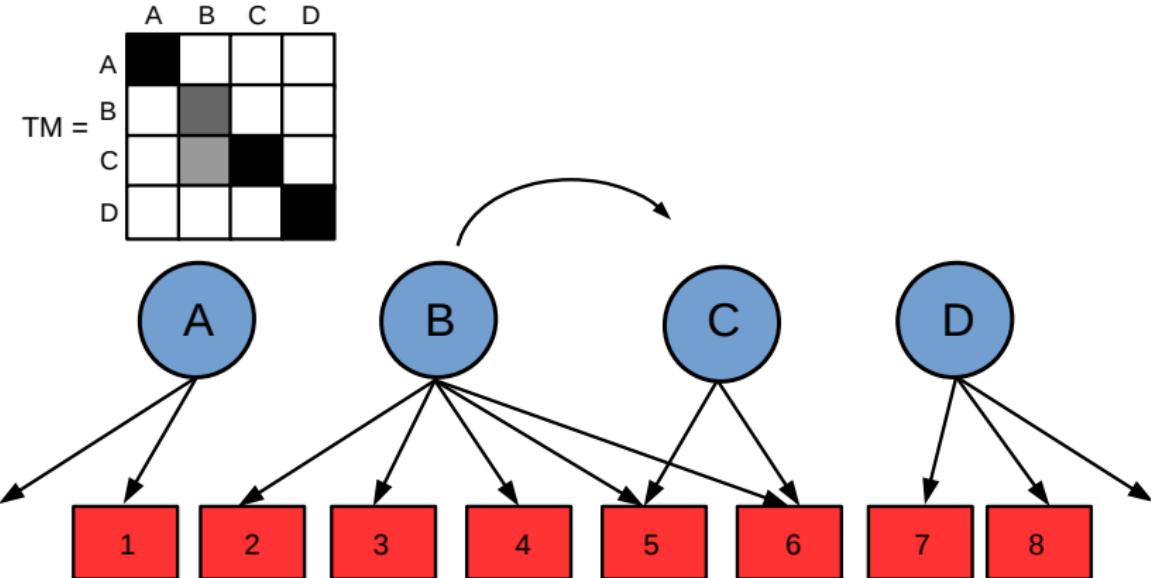
### Edgeweight:

$$w_{i,j} = \frac{1}{1 + e^{\beta_0 + \beta_1 g_{j,(1)} + \cdots + \beta_k g_{j,(k)}}}$$

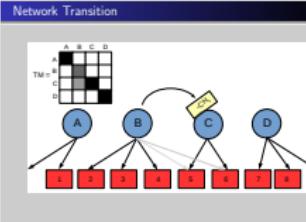
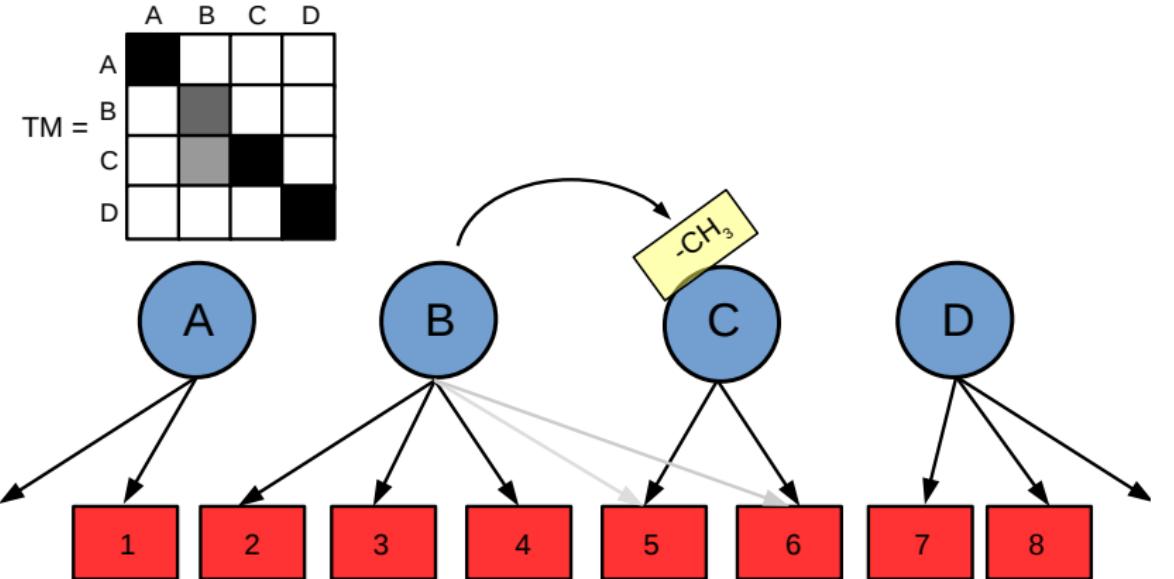
$$w_{i,j} = (1 - \alpha) [d_{i,j}] + \alpha [e_{i,j}]$$



# Network Transition



# Network Transition



# Network Transition

$$E[b_i - a_i] = \tau_{1,i}a_1 + \cdots + \tau_{m,i}a_m$$

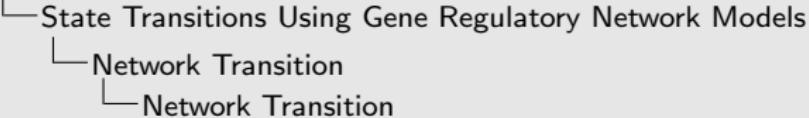
where  $b_i$  and  $a_i$  are column-vectors in  $\mathbf{B}$  and  $\mathbf{A}$  that describe the regulatory targeting of transcription factor  $i$  in the final and initial networks, respectively.

In the simplest case, this can be solved with normal equations,

$$\hat{\tau}_i = (A^T A)^{-1} A^T (b_i - a_i)$$

to generate each of the columns of the transition matrix  $\mathbf{T}$  such that

$$\hat{\mathbf{T}} = [\hat{\tau}_1, \hat{\tau}_2, \dots, \hat{\tau}_m]$$



2017-04-24

$E[b_i - a_i] = \tau_{1,i}a_1 + \cdots + \tau_{m,i}a_m$   
where  $b_i$  and  $a_i$  are column-vectors in  $\mathbf{B}$  and  $\mathbf{A}$  that describe the regulatory targeting of transcription factor  $i$  in the final and initial networks, respectively.  
In the simplest case, this can be solved with normal equations,  
 $\hat{\tau}_i = (A^T A)^{-1} A^T (b_i - a_i)$   
to generate each of the columns of the transition matrix  $\mathbf{T}$  such that  
 $\hat{\mathbf{T}} = [\hat{\tau}_1, \hat{\tau}_2, \dots, \hat{\tau}_m]$

# Network Transition

Regularization:

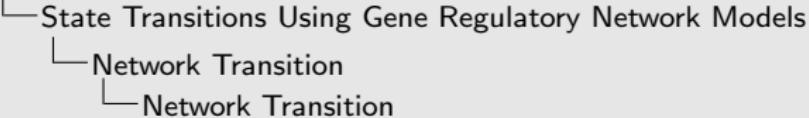
$$\mathbf{Q}_{i,j} = \begin{cases} 1 & \text{for } i = j \neq k \\ 0 & \text{elsewhere} \end{cases},$$

which results in the minimization of the penalized residual sum of squares

$$PRSS(\mathbf{T}_{\cdot,k}) = \sum_{i=1}^p \left( \mathbf{B}_{i,k} - \sum_{j=1}^m A_{i,j} \mathbf{T}_{j,k} \right)^2 + \lambda \sqrt{\mathbf{T}'_{\cdot,k} \mathbf{Q} \mathbf{T}_{\cdot,k}}$$

An implementation of this extension is available in the R package MONSTER.

## Methods for Estimating Hidden Structure and Network Transitions in Genomics

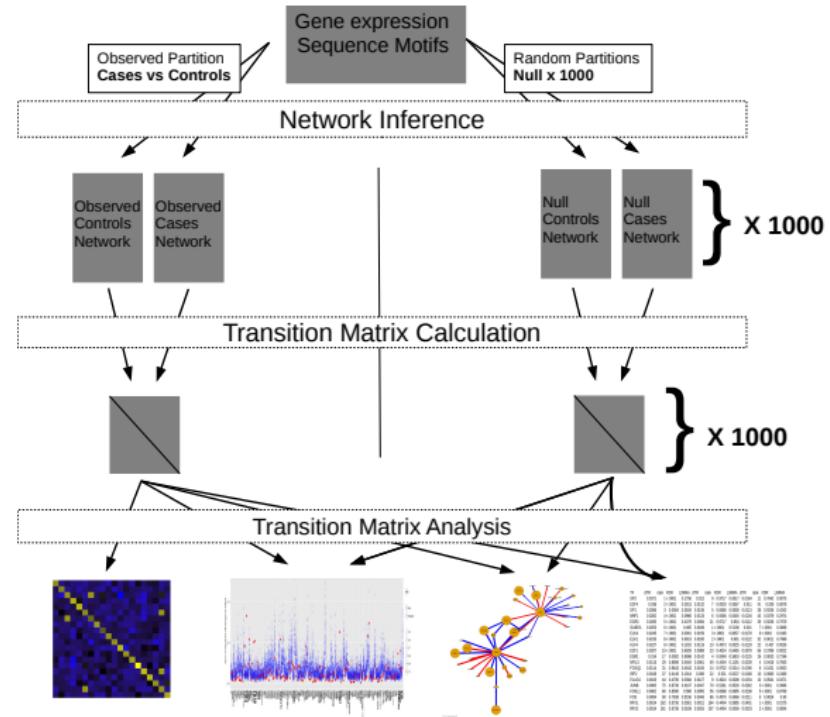


2017-04-24

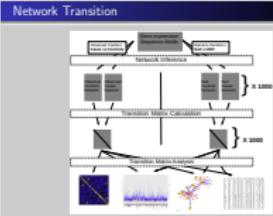
### Network Transition

Regularization:  
 $\mathbf{Q}_{i,j} = \begin{cases} 1 & \text{for } i = j \neq k \\ 0 & \text{elsewhere} \end{cases}$ ,  
which results in the minimization of the penalized residual sum of squares  
$$PRSS(\mathbf{T}_{\cdot,k}) = \sum_{i=1}^p \left( \mathbf{B}_{i,k} - \sum_{j=1}^m A_{i,j} \mathbf{T}_{j,k} \right)^2 + \lambda \sqrt{\mathbf{T}'_{\cdot,k} \mathbf{Q} \mathbf{T}_{\cdot,k}}$$
  
An implementation of this extension is available in the R package MONSTER.

# Network Transition

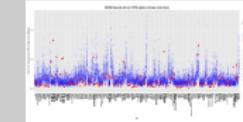


Methods for Estimating Hidden Structure and Network  
Transitions in Genomics  
└ State Transitions Using Gene Regulatory Network Models  
  └ Network Transition  
    └ Network Transition



2017-04-24

# Evaluating Transition Matrix



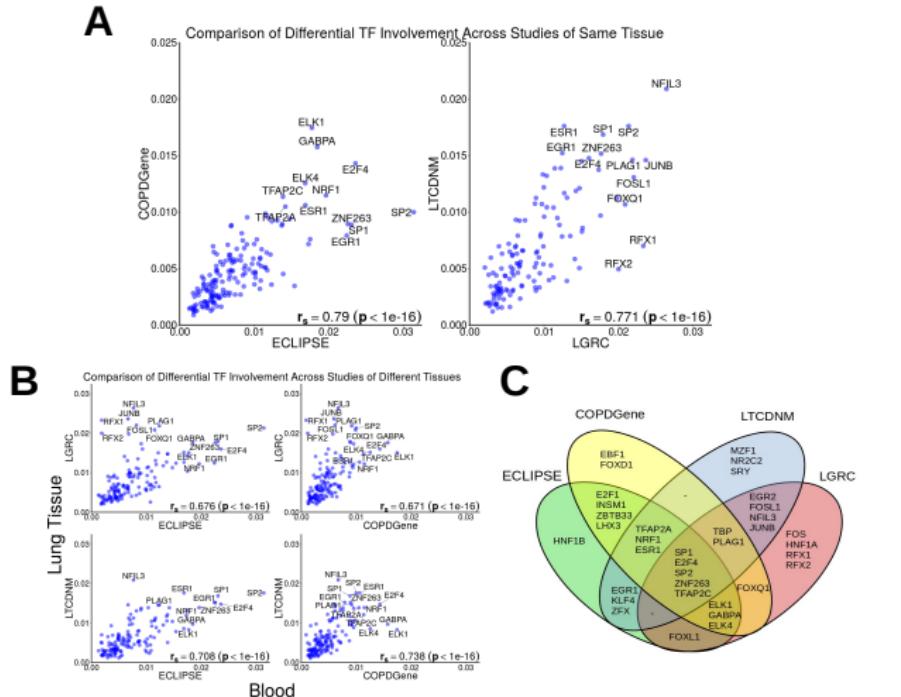
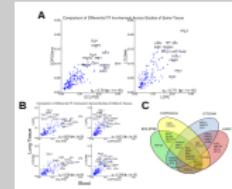
$$d\hat{TFI}_j = \frac{\sum_{i=1}^m I(i \neq j) \hat{\tau}_{ij}^2}{\sum_{i=1}^m \hat{\tau}_{ij}^2}$$

$$d\hat{TFI}_j = \frac{\sum_{i=1}^m I(i \neq j) \hat{\tau}_{ij}^2}{\sum_{i=1}^m \hat{\tau}_{ij}^2}$$



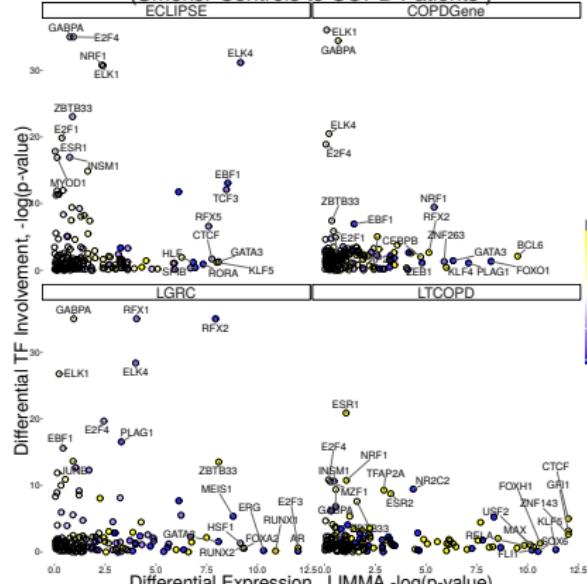
# Reproducibility and novel results

- └ State Transitions Using Gene Regulatory Network Models
- └ Results
- └ Reproducibility and novel results

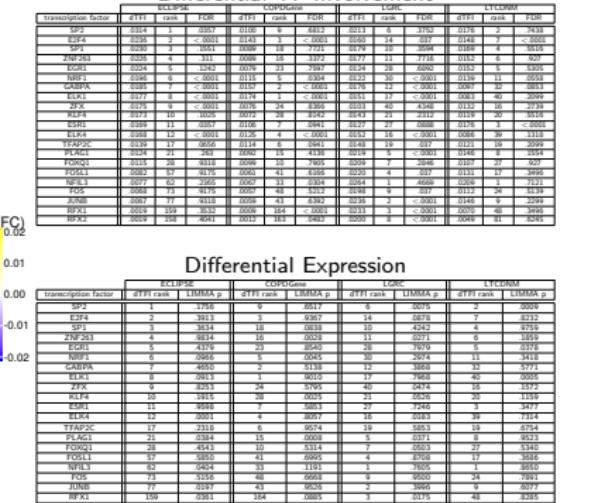


# Reproducibility and novel results

Differential Involvement vs Differential Expression  
(Smoker Controls to COPD Patients)

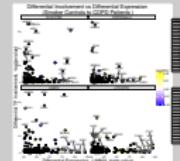


Differential TF Involvement



## Methods for Estimating Hidden Structure and Network Transitions in Genomics

- └ State Transitions Using Gene Regulatory Network Models
- └ Results
- └ Reproducibility and novel results



## Acknowledgements

### Dissertation Committee

- John Quackenbush
  - Christoph Lange
  - Kimberly Glass

Channing Division of Network Medicine

- Ed Silverman
  - Craig Hersh

JQ Lab

- John Quackenbush
  - Joe Barry
  - Joey Chen
  - Maude Fagny
  - Marieke Kuijjer
  - Camila Lopes-Ramos
  - Megha Padi
  - Joe Paulson
  - John Platig
  - Heather Selby
  - Nicole Trotman



| Elimination Committee | JZ Lab                |
|-----------------------|-----------------------|
| John Quackenbush      | John Quackenbush      |
| Christophe Lengyel    | Joe Harry             |
| Karenna Giese         | Mike Hause            |
|                       | Marcus Kujala         |
|                       | Caroline Losos-Raines |
|                       | Magnus Pall           |
|                       | Matthew Pfeifer       |
|                       | John Plueck           |
|                       | Wadeer Selly          |
|                       | Nicole Trutman        |