TCA-TWAS

# Identification of cell-type-specific genetic regulation of gene expression for transcriptome-wide association studies

Qiurui Ma; Duo Zhang; Brandon Jew; Sriram Sankararaman

TCA-TWAS

Background

Model

Experiment

Conclusion



File:Central Dogma of Molecular Biochemistry with Enzymes.jpg. (2019, September 4). Wikimedia Commons, the free media repository. Retrieved 00:50, September 5, 2019 from https://commons.wikimedia.org/w/index.php?title=File:Central_Dogma_of_Molecular_Biochemistry_with_Enzymes.jpg&oldid=364559532.

# The Central Dogma

TCA-TWAS

Background

Model

Experiment

Conclusion

## SNPs

| Ind 1 | Ind 2 | Ind 3 |
|-------|-------|-------|
| AA    | AG    | AG    |
| CT    | CC    | CC    |

| Geno1 | Geno2 | Geno3 |
|-------|-------|-------|
| 0     | 1     | 1     |
| 1     | 1     | 2     |

replication
(DNA -> DNA)

**DNA Polymerase**

**DNA**

transcription
(DNA -> RNA)

**RNA Polymerase**

**RNA**

translation
(RNA -> Protein)

**Ribosome**

**Protein**

**Influence**

**Phenotype**

# The Central Dogma

TCA-TWAS

Background

Model

Experiment

Conclusion

**SNPs**

|  | Ind 1 | Ind 2 | Ind 3 |
|---|---|---|---|
|  | AA | AG | AG |
|  | CT | CC | CC |

|  | Geno1 | Geno2 | Geno3 |
|---|---|---|---|
|  | 0 | 1 | 1 |
|  | 1 | 1 | 2 |

replication
(DNA -> DNA)

**DNA Polymerase**

DNA

transcription
(DNA -> RNA)

**RNA Polymerase**

RNA

translation
(RNA -> Protein)

**Ribosome**

Protein

**Influence**

**Phenotype**

**Gene Expression (GE):**
cell 1, cell 2, cell 3, cell 4

**Phenotype:**
Height, Skin Color

TCA-TWAS

Background

Model

Experiment
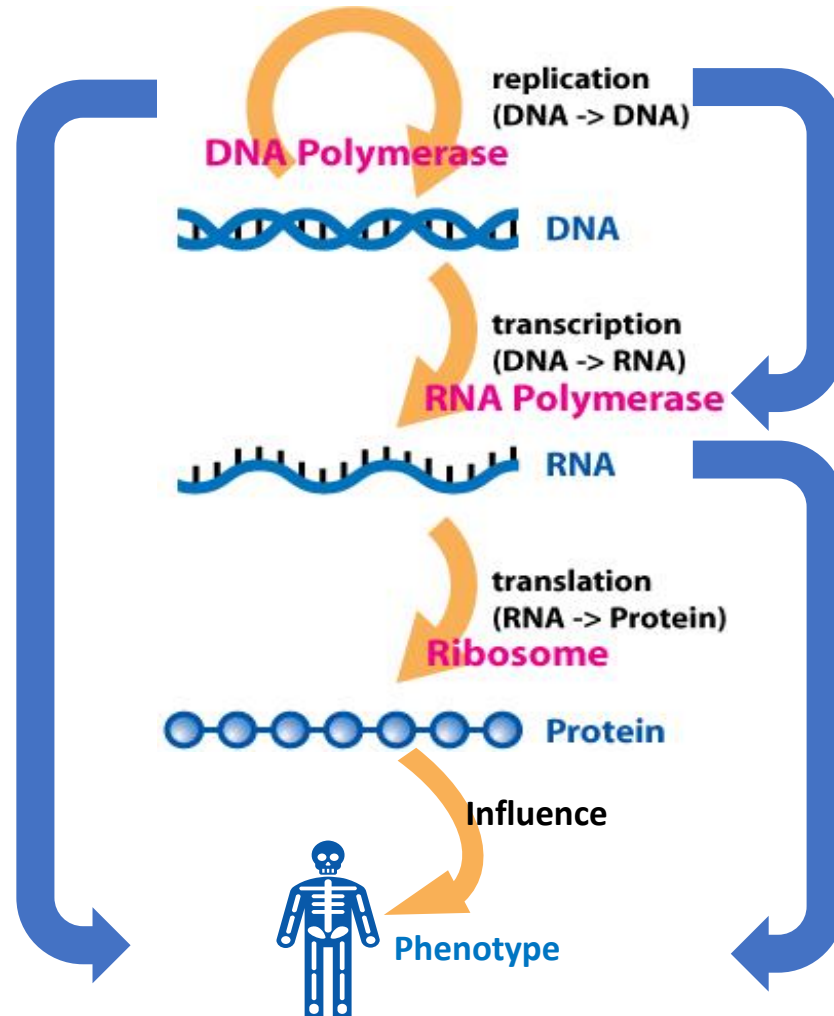
Conclusion

**Linear Regression of SNPs and phenotype**



replication
(DNA -> DNA)

**DNA Polymerase**

DNA

transcription
(DNA -> RNA)

**RNA Polymerase**

RNA

translation
(RNA -> Protein)

**Ribosome**

Protein

**Influence**

**Phenotype**

Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B. W., ... & Sullivan, P. F. (2016). Integrative approaches for large-scale transcriptome-wide association studies. *Nature genetics*, *48*(3), 245

TCA-TWAS

Background

Model

Experiment

Conclusion



Linear Regression of SNPs and GE

Linear Regression of SNPs and phenotype

Linear Regression of GE and Phenotype

Gusev, A., Ko, A., Shi, H., Bhatia, G., Chung, W., Penninx, B. W., ... & Sullivan, P. F. (2016). Integrative approaches for large-scale transcriptome-wide association studies. *Nature genetics*, *48*(3), 245

TCA-TWAS

**Background**

Model

Experiment

Conclusion



Manhattan Plot of UADT(Upper Aero-Digestive Tract) Cancer GWAS Discovery Phase

Mckay et, al. (2011). A Genome-Wide Association Study of Upper Aerodigestive Tract Cancers Conducted within the INHANCE Consortium. PLoS genetics. 7. e1001333. 10.1371/journal.pgen.1001333.

TCA-TWAS

Background

Model

Experiment

Conclusion

**Healthy Blood Cells**

Red blood cells

Neutrophil*

Lymphocyte*

Monocyte*

Basophil*

Platelet

*White blood cells

replication (DNA -> DNA)

**DNA Polymerase**

DNA

transcription (DNA -> RNA)

**RNA Polymerase**

RNA

**Bulk Level Gene Expression**

**Weighted Cell-Type-Specific Gene Expression**

transcription (DNA -> RNA)

**RNA Polymerase**

RNA

transcription (DNA -> RNA)

**RNA Pol**

RNA

transcription (DNA -> RNA)

**RNA Polymerase**

RNA

Healthy Blood Cells

Red blood cells
Neutrophil*
Lymphocyte*
Monocyte*
Basophil*
Platelet
*White blood cells

**Methodological:**
Unclear how SNPs affect phenotypes
- Missing cell type information
- Fail to tell causality from correlation

**Practical:**
Cell-Specific Biological Data being resource intensive, expensive to acquire

TCA-TWAS

Background

Model

Experiment

Conclusion

**Healthy Blood Cells**



Red blood cells

Neutrophil*

Lymphocyte*

Monocyte*

Basophil*

Platelet

*White blood cells

**Methodological:**

Unclear how SNPs affect phenotypes
- Missing cell type information
- Fail to tell causality from correlation

**Practical:**

Cell-Specific Biological Data being resource intensive, expensive to acquire

**Goal:**

Impute cell-type specific gene expressions from SNPs and bulk level gene expressions to perform downstream TWAS

TCA-TWAS

Background

Model

Experiment

Conclusion



**Train Dataset**    **Target Dataset**

**Bulk Level GE**

TCA-TWAS

Background
Model
Experiment
Conclusion

**Train Dataset**

**TCA Deconvolution**



1

**Cell-Type-Specific Gene Expression**

**Bulk Level Gene Expression**

# Pipeline Framework
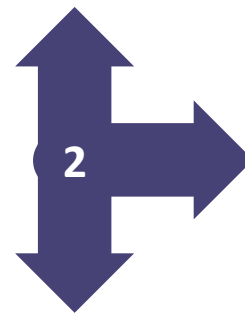
**Train Dataset**

replication
(DNA -> DNA)

**DNA Polymerase**

**DNA**

**TCA**

transcription
(DNA -> RNA)
RNA Polymerase
RNA

transcription
(DNA -> RNA)
RNA Polymerase
RNA

transcription
(DNA -> RNA)
RNA Polymerase
RNA

**2**

**1**

transcription
(DNA -> RNA)
**RNA Polymerase**
**RNA**

**LASSO**

$\beta$

**Estimated effect size of SNPs on GE**

TCA-TWAS

Background

Model

Experiment

Conclusion

**Train Dataset**

**TCA**

**Linear Regression**

2

1

**LASSO**

$\beta$

3

4

**Estimated GE**

**Phenotype**

14

TCA-TWAS

Background

Model

Experiment

Conclusion

**Target Dataset**

replication
(DNA -> DNA)

**DNA Polymerase**

**DNA**

**Linear Regression**

**LASSO**

$\boldsymbol{\beta}$    **3**

transcription
(DNA -> RNA)

**RNA Polymerase**

**4**

**RNA**

**Estimated GE**

**Phenotype**

15

TCA-TWAS

Background

Model

Experiment

Conclusion

**added SNPs effect: mdl1**

**TCA Model: mdl2**

$$Z_h^i = \epsilon_z + \mu_h + \begin{bmatrix} 1 \\ 27 \\ 0 \end{bmatrix}^T \gamma_h + \begin{bmatrix} 0 \\ 2 \\ 1 \\ \vdots \\ 0 \\ 1 \end{bmatrix}^T \beta_h$$

**Cell-Specific GE**

**Covariate:**
**(Gender, Age, Smoking)**

**SNPs**

Rahmani, E., Schweiger, R., Rhead, B., Criswell, L. A., Barcellos, L. F., Eskin, E., ... & Halperin, E. (2019). Cell-type-specific resolution epigenetics without the need for cell sorting or single-cell biology. *BioRxiv*, 437368.

TCA-TWAS

Background

**Model**

Experiment

Conclusion

**added SNPs effect: mdl1**

**TCA Model: mdl2**

$$Z_h^i = \epsilon_z + \mu_h + \begin{bmatrix} 1 \\ 27 \\ 0 \end{bmatrix}^T \gamma_h + \begin{bmatrix} 0 \\ 2 \\ 1 \\ \vdots \\ 0 \\ 1 \end{bmatrix}^T \beta_h$$

**Cell-Specific GE**

**Covariate:**
**(Gender, Age, Smoking)**

**SNPs**

$$G_i = c_i^2 \delta + \sum_{h=1}^{k} w_{hi} z_{hi} + \epsilon_g$$

**Cell Type Weight Extra Information**

Rahmani, E., Schweiger, R., Rhead, B., Criswell, L. A., Barcellos, L. F., Eskin, E., ... & Halperin, E. (2019). Cell-type-specific resolution epigenetics without the need for cell sorting or single-cell biology. *BioRxiv*, 437368.

TCA-TWAS

Background

Model

Experiment

Conclusion

**added SNPs effect: mdl1**

**TCA Model: mdl2**

$$Z_h^i = \epsilon_z + \mu_h + \begin{bmatrix} 1 \\ 27 \\ 0 \end{bmatrix}^T \gamma_h + \begin{bmatrix} 0 \\ 2 \\ 1 \\ \vdots \\ 0 \\ 1 \end{bmatrix}^T \beta_h$$

**Cell-Specific GE**

**Covariate:**
**(Gender, Age, Smoking)**

**SNPs**

$$G_i = c_i^2 \delta + \sum_{h=1}^{k} w_{hi} z_{hi} + \epsilon_g$$

**Cell Type Weight Extra Information**

$$Pr(Z_h^i | G_i, w_i, \mu_h, \sigma_z, \sigma_g, \sigma_\delta)$$

**EM to infer parameters**

Rahmani, E., Schweiger, R., Rhead, B., Criswell, L. A., Barcellos, L. F., Eskin, E., ... & Halperin, E. (2019). Cell-type-specific resolution epigenetics without the need for cell sorting or single-cell biology. *BioRxiv*, 437368.

TCA-TWAS

Background

Model

Experiment

Conclusion

**Is the prediction result of proposed model significant?**



transcription
(DNA -> RNA)
**RNA Polymerase**
RNA

**Estimated GE**

**4**

**Phenotype**

**Power of cell-specific expression imputation**



Legend:
— High abundance
— Intermediate abundance
— Low abundance

X-axis: Training sample size (100, 200, 300, 400, 500)
Y-axis: Power (0.0 to 1.0)

**Power calculated as percentage of p value less than threshold**
**Heritability 0.3, 1000 genes, 100000 sample size**

TCA-TWAS

Background

Model

Experiment

Conclusion

## Is the model underfitting or overfitting?



LASSO $\beta$

**3**

**Estimated GE**

**Heritability being theoretical upper bound for r^2**

## Could TCA recover the ground truth data distribution?

TCA-TWAS

Background

Model

Experiment

Conclusion

**Step** 1



Simulated cell–type–specific gene expression for one gene / TCA estimated cell–type specific gene expression for one gene. Cell_type1 W: 0.53, Cell_type2 W: 0.349, Cell_type3 W: 0.0883, Cell_type4 W: 0.0321

**TCA recover abundant cell type distribution, but not the lesser ones**

**TCA-TWAS**

Background

Model

Experiment

Conclusion

## Compare our modified TCA and the original TCA, which is better?



Correlation of estimated SNPs effect size and the ground truth

LASSO $\beta$
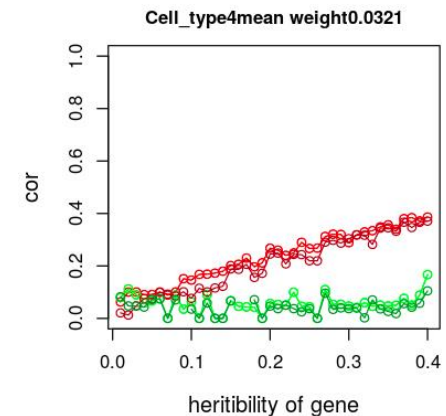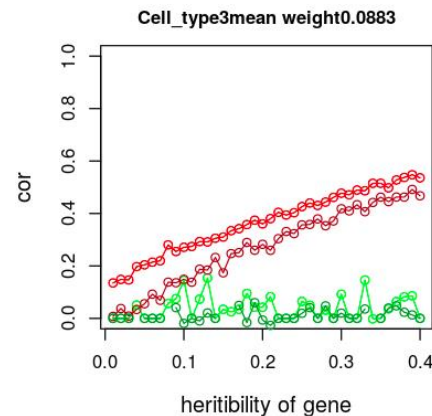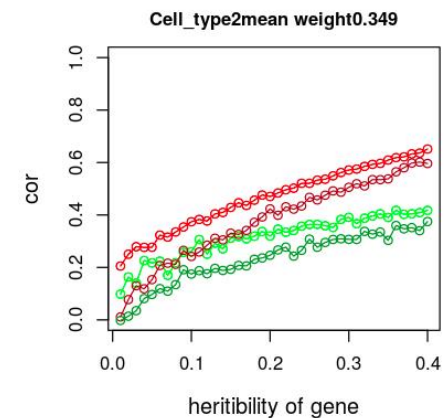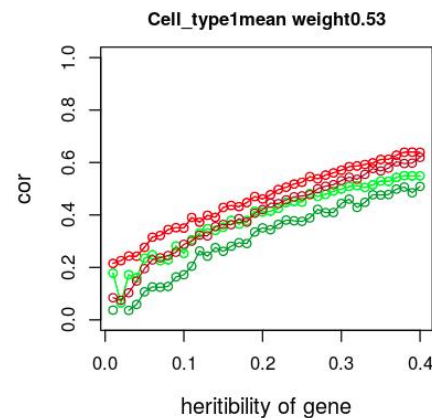
## Compare our modified TCA and the original TCA, which is better?



**LASSO** $\beta$

**Estimated GE**

**Can the model leverage other structural information inside data?**

TCA-TWAS

Background

Model

Experiment

Conclusion

**Is the model's correctly identifying causal effects?**



Precision and recall distribution of three model

TCA-TWAS

Background

Model

Experiment

Conclusion

## Train Dataset: Dutch Dataset with GE

- 5000 individuals, blood, 4 cell types
- 10201 Genes, 801501 SNPs, 10201 GE

**Ancestry Pruning** → **LD Pruning** → **Missing Value Imputation** →
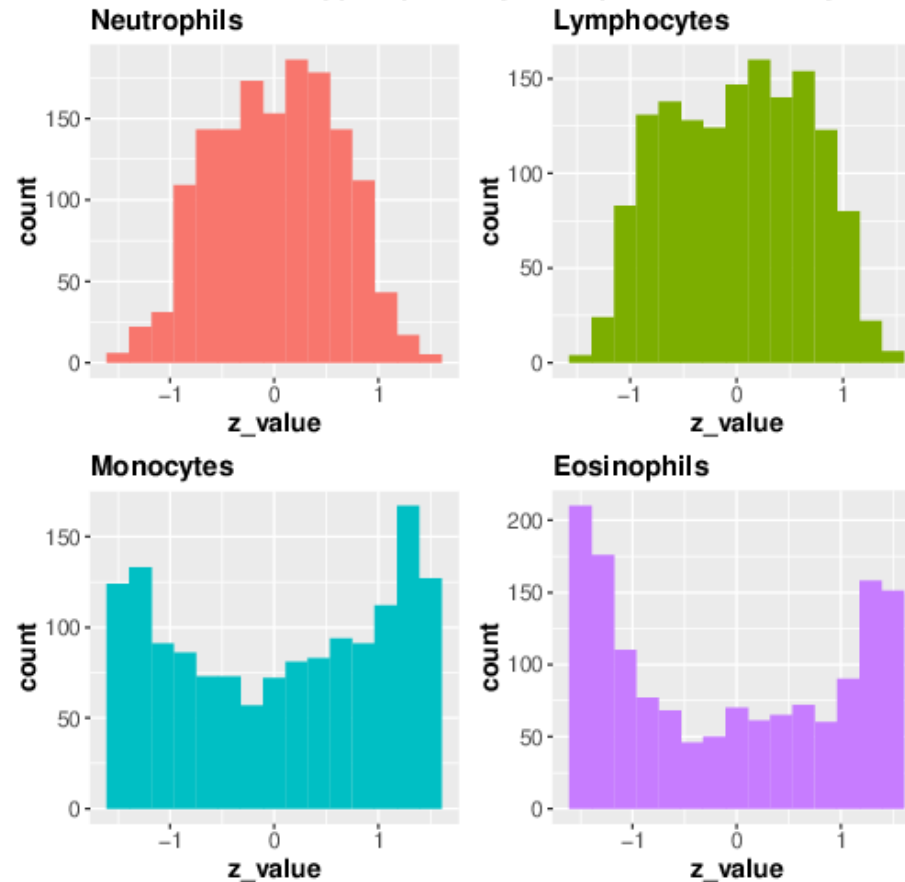
**CIS-SNPs locating** → **Pipeline**

## Target Dataset: UKBiobank Dataset without GE

- 500,000 individuals
- SNPs, Phenotypes

TCA-TWAS

Background

Model

Experiment

Conclusion

**Is the model's performance consistent between cell types?**
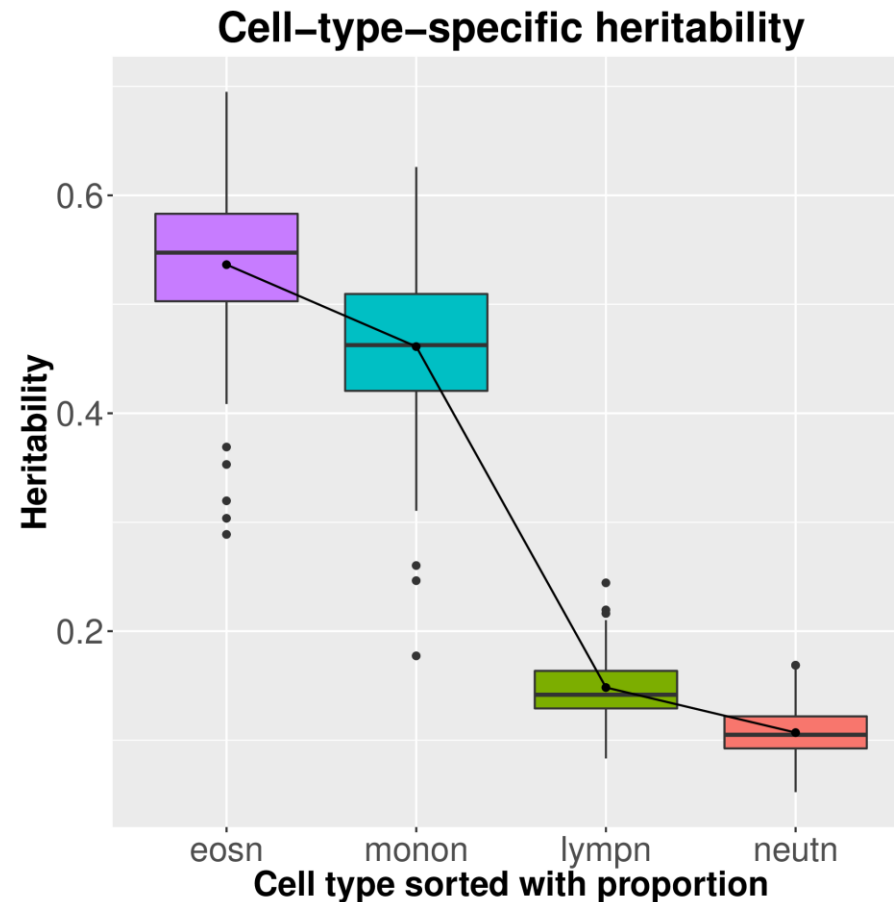


TCA estimated cell-type-specific gene expression for one gene

**TCA's estimation result fall short of normal distribution, which is the original assumption**

**TCA-TWAS**

Background

Model

Experiment

Conclusion

**Is the model's performance consistent between cell types?**



Cell–type–specific heritability

**Heritability calculated from LASSO's prediction for gene expression. Overfitting occurs for low cell types**

TCA-TWAS

Background

Model

Experiment

Conclusion

## Summary

- Cell-type-specific GE can be imputed from SNPs in a relative low cost to conduct the downstream phenotype association study

## Key Contribution

- Incorporated genetic effects into TCA to deconvolute bulk level GE into cell-type specific ones
- Produced effect size estimators on Dutch Dataset to impute cell-type specific gene expressions on UKBiobank Dataset
- Functional R package on Cran for standard usage

## Future Work

- Add sparsity constraints to TCA parameter estimates
- Utilize extra correlation structure to enhance TCA performance
- Consider batch effects when transferring effect size from train dataset to target dataset

TCA-TWAS

Background

Model

Experiment

Conclusion

**Q&A**