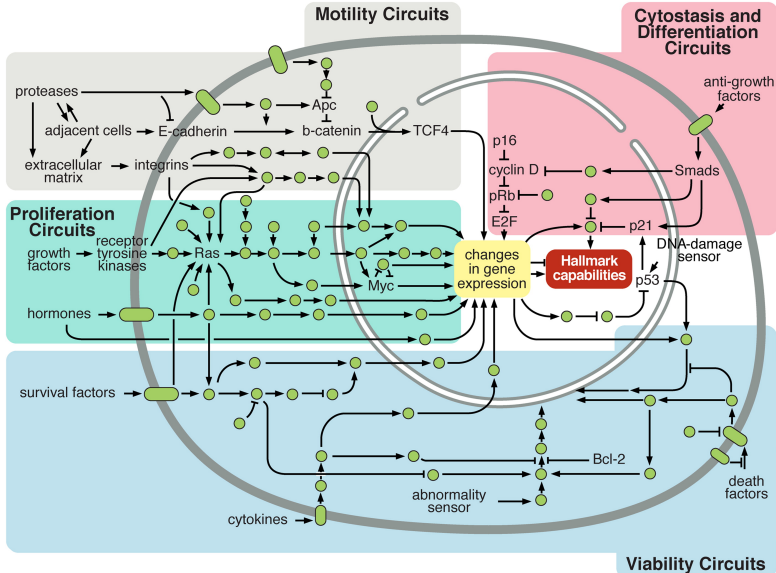


# CS 5854: Supervised Inference of Gene Regulatory Networks from Single-Cell Gene Expression Data

T. M. Murali

February 16, 18, 2021

# Signaling Pathways and Gene Expression



# Gene Expression is a Dynamic Process

**B**

**if (F = 1 or E = 1 or CD = 1) and (Z = 1)**      Repression functions of modules F, E, and DC mediated by Z site

$$\alpha = 1$$

**else**       $\alpha = 0$

**if (P = 1 and CG<sub>1</sub> = 1)**

Both P and CG<sub>1</sub>, needed for synergistic link with module B

$$\beta = 2$$

**else**       $\beta = 0$

**if (CG<sub>2</sub> = 1 and CG<sub>3</sub> = 1 and CG<sub>4</sub> = 1)**

Final step up of system output

$$\gamma = 2$$

**else**       $\gamma = 1$

$$\delta(t) = B(t) + G(t)$$

Positive input from modules B and G

$$\varepsilon(t) = \beta * \delta(t)$$

Synergistic amplification of module B output by CG<sub>1</sub>-P subsystem

**if ( $\varepsilon(t) = 0$ )**

$$\xi(t) = Otx(t)$$

Switch determining whether Otx site in module A, or upstream modules (i.e., mainly module B), will control level of activity

**else**       $\xi(t) = \varepsilon(t)$

**if ( $\alpha = 1$ )**

Repression function inoperative in endoderm but blocks activity elsewhere

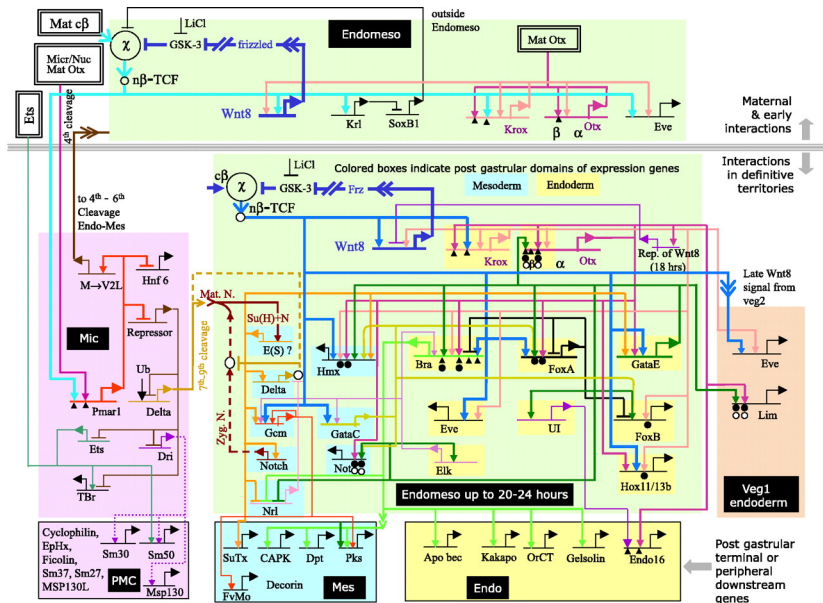
$$\eta(t) = 0$$

**else**       $\eta(t) = \xi(t)$

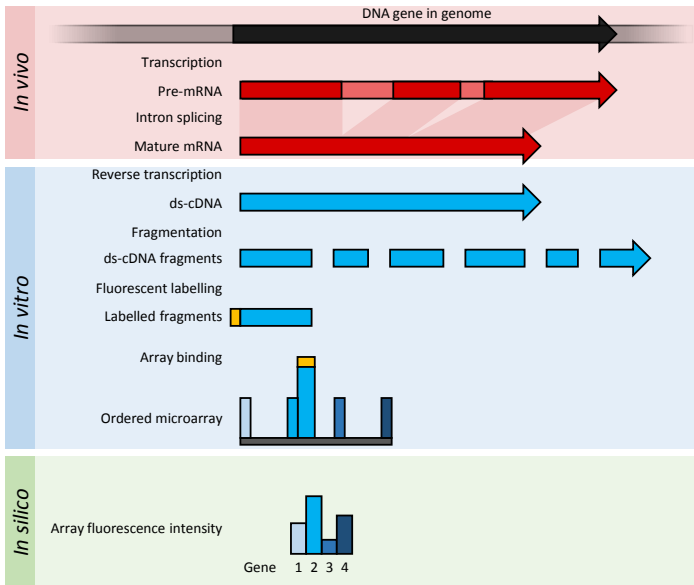
$$\Theta(t) = \gamma * \eta(t)$$

Final output communicated to BTA

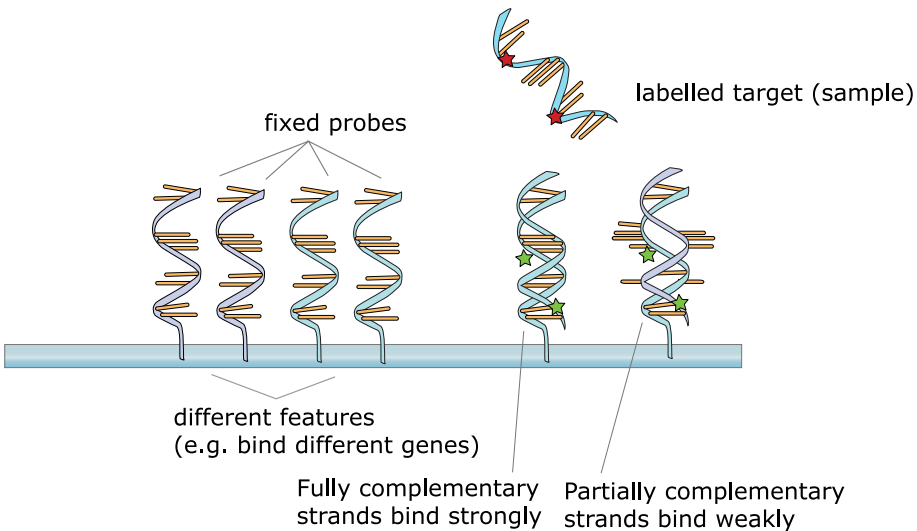
# Gene Expression is a Dynamic Process



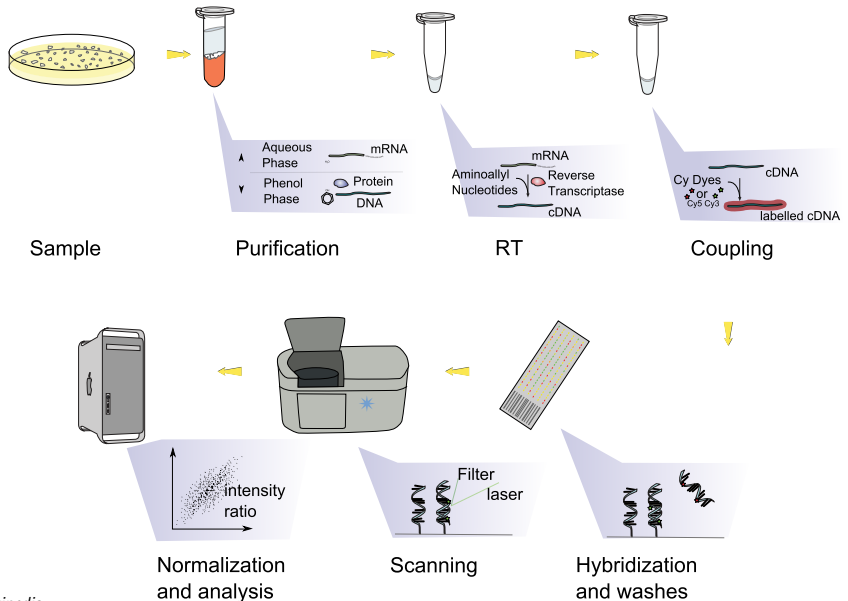
# Measuring Genomewide Gene Expression: DNA Microarrays



## Measuring Genomewide Gene Expression: DNA Microarrays



## Measuring Genomewide Gene Expression: DNA Microarrays



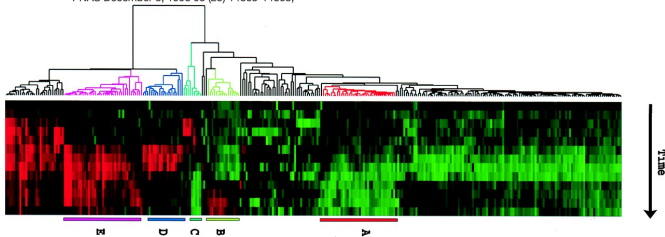
# Applications of DNA Microarray Data

RESEARCH ARTICLE

## Cluster analysis and display of genome-wide expression patterns

Michael B. Eisen, Paul T. Spellman, Patrick O. Brown, and David Botstein

PNAS December 8, 1998 95 (25) 14863-14868;





# Applications of DNA Microarray Data

## RESEARCH ARTICLE

### Cluster analysis and display of genome-wide expression patterns

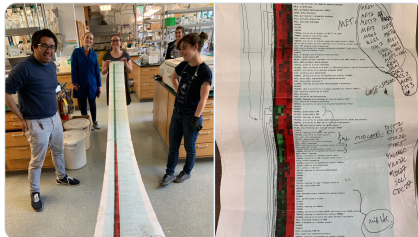
Michael B. Eisen, Paul T. Spellman, Patrick O. Brown, and David Botstein

PNAS De



Michael Eisen   
@mbeisen

Inspired by @UCSDCooperLab's question about origins of the red/green color scheme in microarray clustering, I present THE FIRST dna microarray cluster analysis made by me in 1997 for [ncbi.nlm.nih.gov/m/pubmed/97841...](https://ncbi.nlm.nih.gov/m/pubmed/97841...) w/handwritten notes from Pat Brown and the late Ira Herskowitz.



6:27 PM · Jun 4, 2019 · Twitter for iPhone

# Applications of DNA Microarray Data

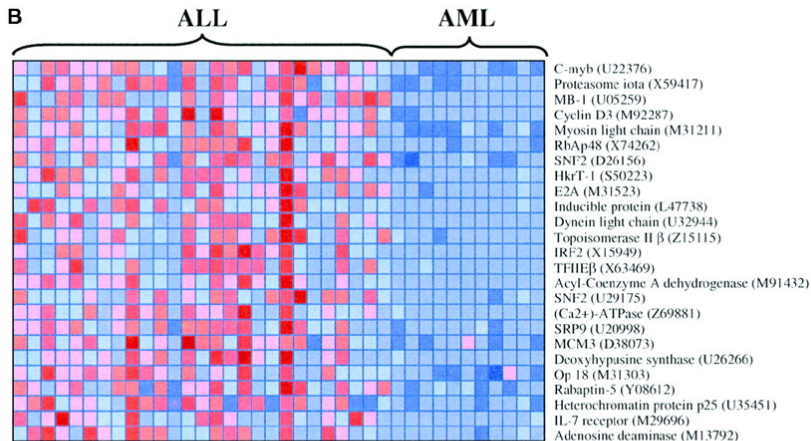
REPORT

## Molecular Classification of Cancer: Class Discovery and Class Prediction by Gene Expression Monitoring

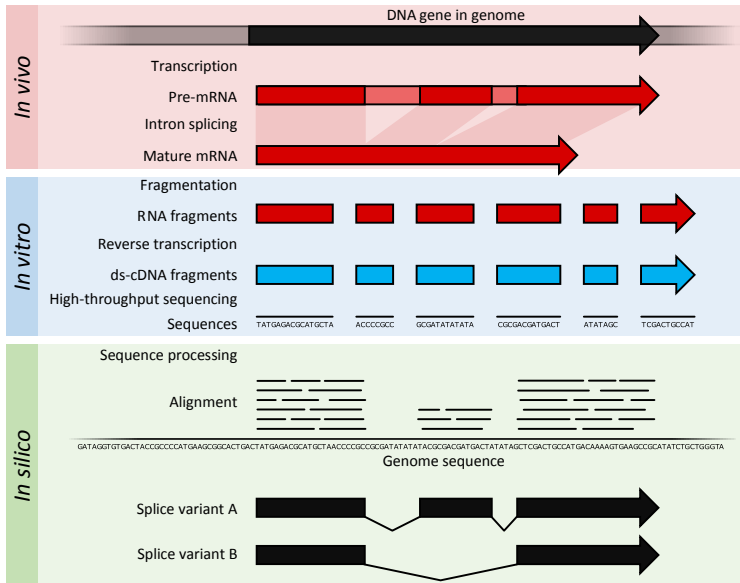
T. R. Golub<sup>1,2,3,4</sup>, D. K. Slonim<sup>1,4</sup>, P. Tamayo<sup>1</sup>, C. Huard<sup>1</sup>, M. Gaasenbeek<sup>1</sup>, J. P. Mesirov<sup>1</sup>, H. Coller<sup>1</sup>, M. L. Loh<sup>2</sup>, J. R. Downing<sup>1</sup>

• See all authors and affiliations

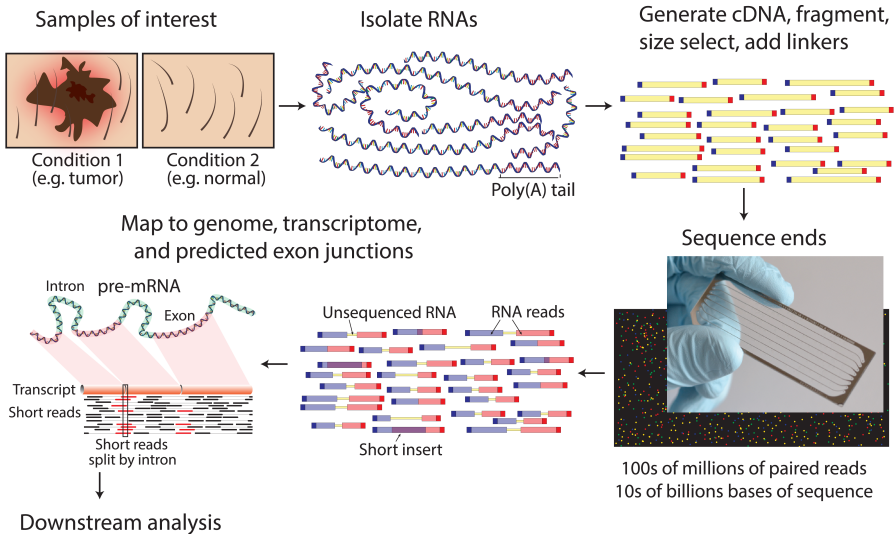
Science 15 Oct 1999  
 Vol. 286, Issue 5439, pp. 531-537  
 DOI: 10.1126/science.286.5439.531



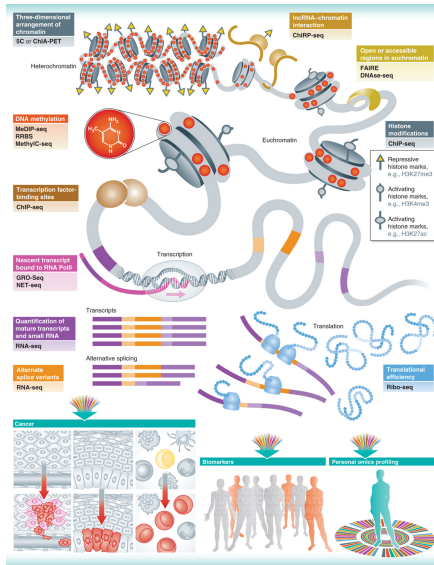
## Measuring Genomewide Gene Expression: RNA-seq



## Measuring Genomewide Gene Expression: RNA-seq



# \*-Seq Techniques



Soon, Hariharan, and Snyder. High-throughput sequencing for biology and medicine. *Mol. Sys. Bio.* 2013.

# Single-Cell RNA-Seq

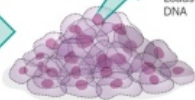
## ONE GENOME FROM MANY

Sequencing the genomes of single cells is similar to sequencing those from multiple cells — but errors are more likely.

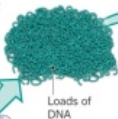
### ► Standard genome sequencing



A sample containing thousands to millions of cells is isolated.



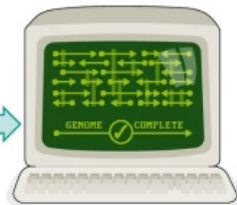
DNA is extracted from all the nuclei.



Loads of DNA



DNA is broken into fragments and then sequenced.

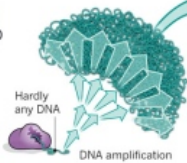


The sequences are assembled to give a common, 'consensus' sequence.

### ► Single-cell sequencing



A single cell is difficult to isolate, but it can be done mechanically or with an automated cell sorter.

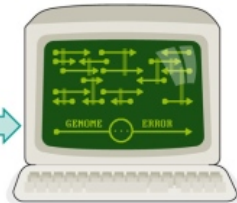


Hardly any DNA  
DNA amplification

The DNA is extracted and amplified, during which errors can creep in.



Amplified DNA is sequenced.

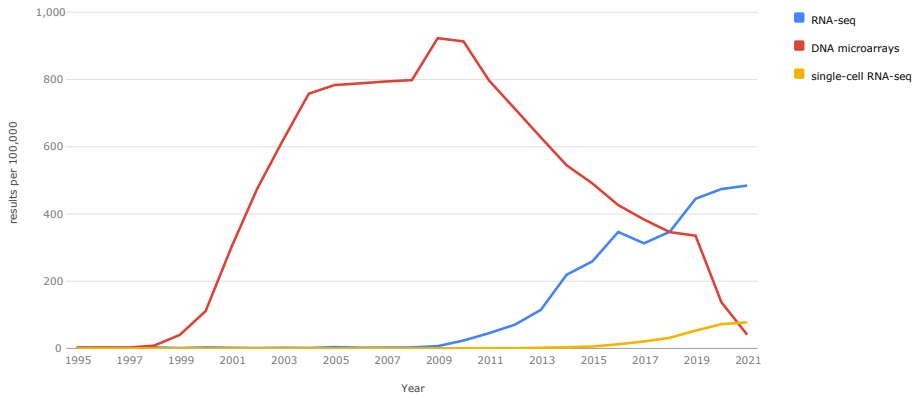


Errors introduced in earlier steps make sequence assembly difficult; the final sequence can have gaps.

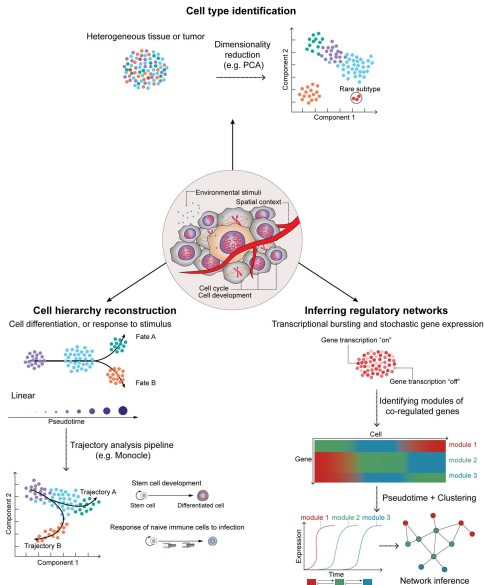
Owens, *Genomics: The single life*, Nature, 2012.

# Technology Trends

Results per 100,000 citations in PubMed  
proportion for each search by year, 1995 to 2021



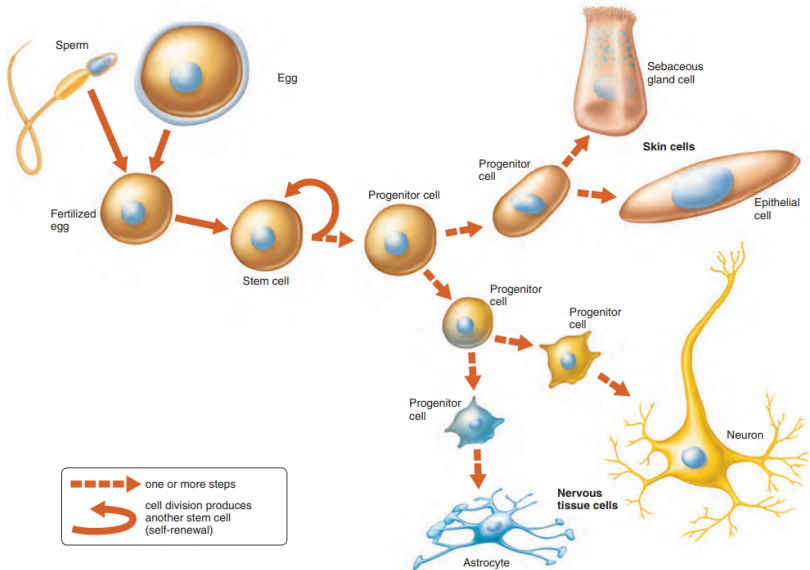
# Applications of scRNA-seq Data



Hwang, Lee, and Bang, *Single-cell RNA sequencing technologies and bioinformatics pipelines*, *Exp. Mol. Med.*, 2018

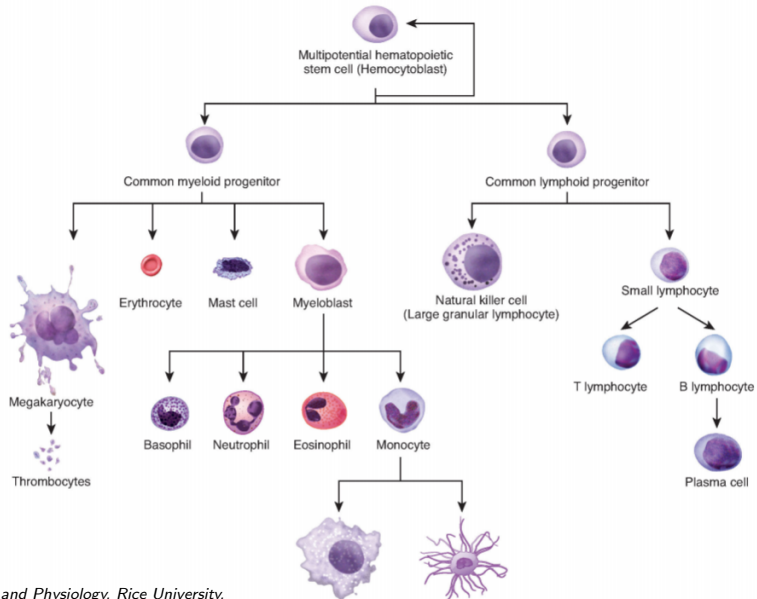


# Cellular Differentiation

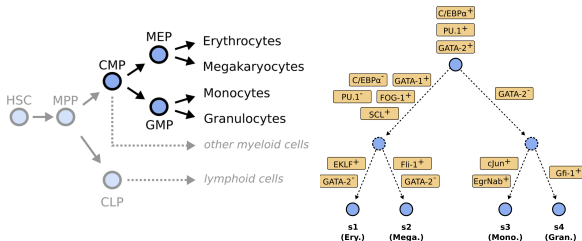


Shier et al., (2015) "Hole's Essentials of Human Anatomy and Physiology". McGraw-Hill

# Cellular Differentiation



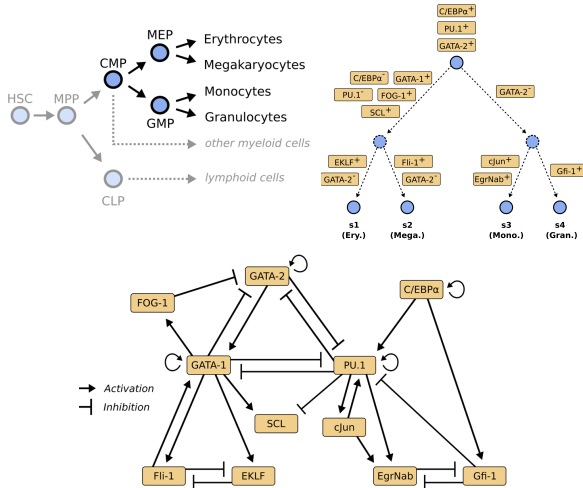
# Cellular Differentiation



- Cells in different states express different sets of genes.
- Cells move from one “state” to another.

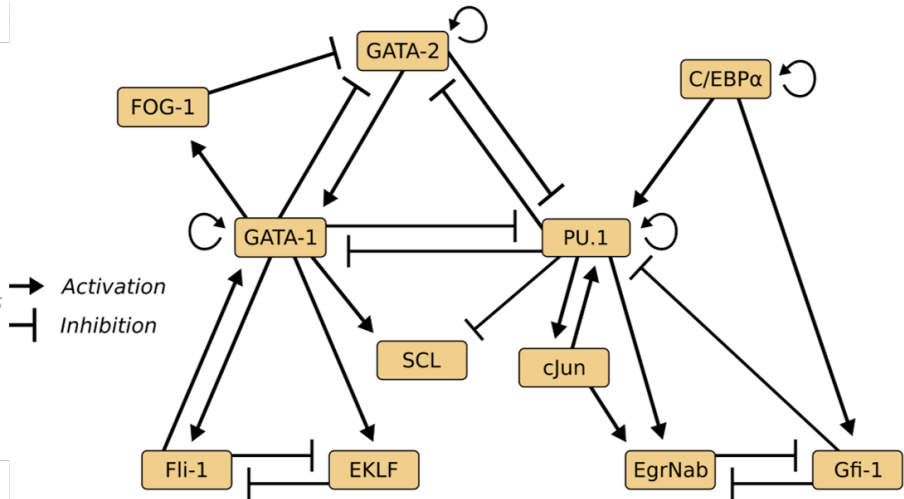
Krumsiek et al. (2010). “Hierarchical Differentiation of Myeloid Progenitors...” PLoS ONE

# Cellular Differentiation



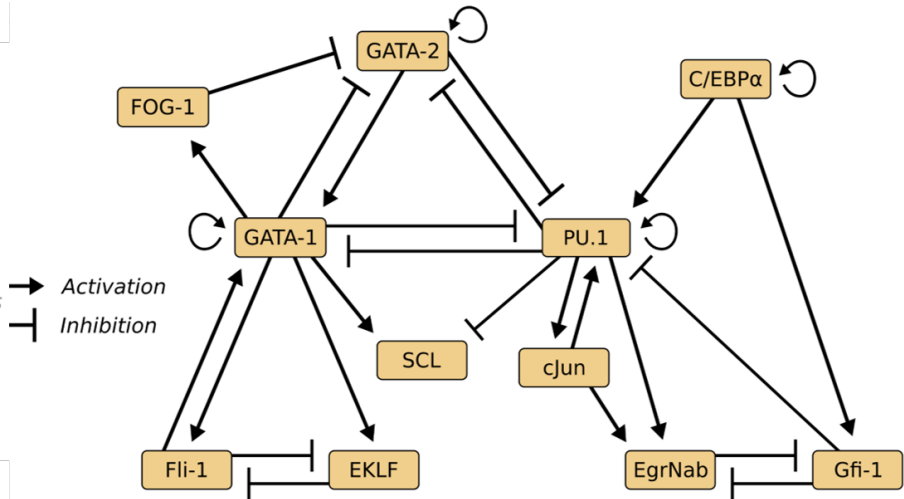
- Transcription factors activate/inhibit genes to effect cell transition from one state to another.

# Gene Regulatory Network (GRN)



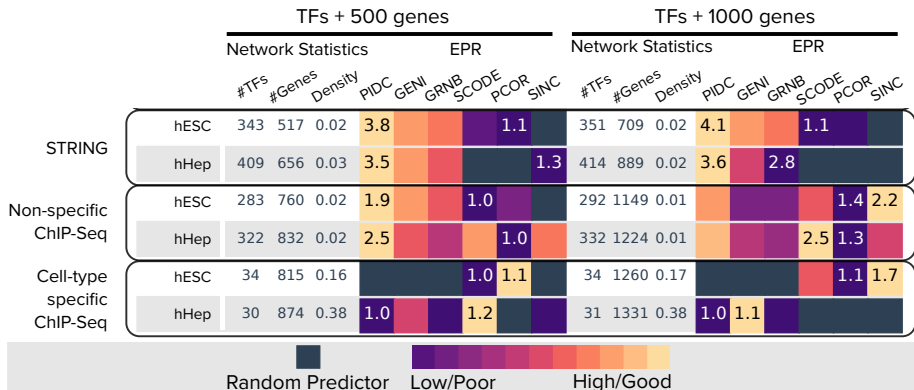
Krumsiek et al. (2010). "Hierarchical Differentiation of Myeloid Progenitors..." PLoS ONE

# Gene Regulatory Network (GRN)

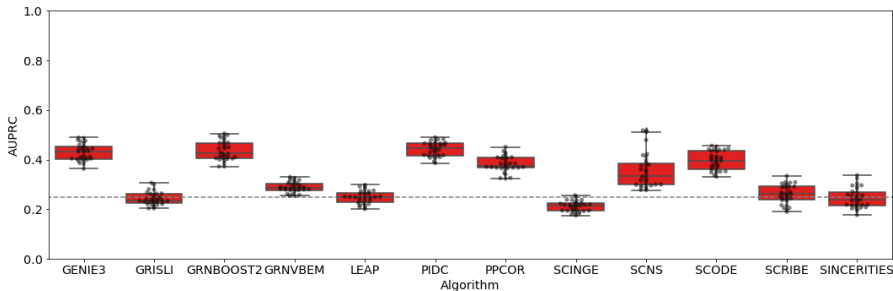
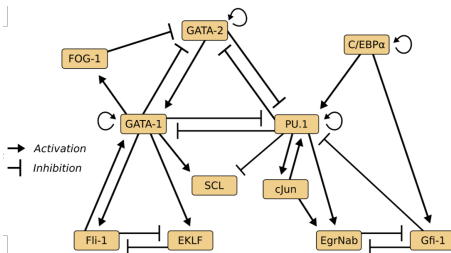


How do we build GRNs using computational techniques?

# BEELINE Results for Human Datasets



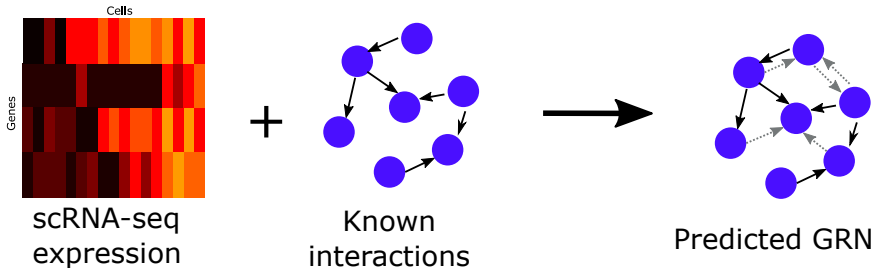
# Poor AUPRC performance





# Supervised GRN inference

Can supervised learning methods take advantage of known regulatory interactions for GRN inference from scRNA-seq data?



# Existing approaches for supervised GRN

- Generate TF-gene features and build a classifier (e.g., SVM)
  - ▶ Concatenate expression vectors<sup>1</sup>
  - ▶ Outer product<sup>2</sup>
  - ▶ Kernels<sup>3</sup>

---

<sup>1</sup> Cerulo et al. (2010) "Learning gene regulatory ..." *BMC Bioinfo.*, 11(1):228

<sup>2</sup> Maetschke et al. (2014) "Supervised, semi- ..." *Brief. Bioinfo.*, 15(2):195–211

<sup>3</sup> Cuong et al. (2008) "Supervised inference ..." *BMC Bioinfo.*, 9(1):2

# Drawbacks

Drawbacks of fixed TF-gene feature representation:

- 1 Dropouts + noise in the input expression data
  - ▶ Dropout: where a gene is observed in one cell but is not detected in another cell of the same cell type
  - ▶ Unclear how fixed feature representation can overcome these problems

# Drawbacks

Drawbacks of fixed TF-gene feature representation:

- 1 Dropouts + noise in the input expression data
  - ▶ Dropout: where a gene is observed in one cell but is not detected in another cell of the same cell type
  - ▶ Unclear how fixed feature representation can overcome these problems
- 2 Do not scale well for datasets with large number of cells

# Proposed solutions

- ① Dropouts + noise in the input expression data → Denoise and impute data using network propagation<sup>1 2 3</sup>

---

<sup>1</sup>Ronen *et al.* (2018) "netSmooth ..." *F1000 Res.* 7.

<sup>2</sup>Ye *et al.* (2019) "scNPF " *BMC Genomics* 20, 347.

<sup>3</sup>Elyanow *et al.* (2020) "netNMF-sc ..." *Gen Res* 30.2: 195-204.

# Proposed solutions

- 1 Dropouts + noise in the input expression data → Denoise and impute data using network propagation<sup>1 2 3</sup>
- 2 Do not scale well for datasets with large number of cells → Dimensionality reduction

---

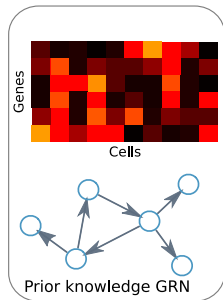
<sup>1</sup>Ronen *et al.* (2018) "netSmooth ..." *F1000 Res.* 7.

<sup>2</sup>Ye *et al.* (2019) "scNPF " *BMC Genomics* 20, 347.

<sup>3</sup>Elyanow *et al.* (2020) "netNMF-sc ..." *Gen Res* 30.2: 195-204.

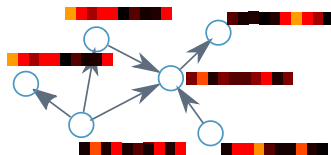
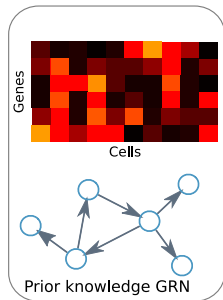
# Denosing the input data

## Inputs



# Denoising the input data

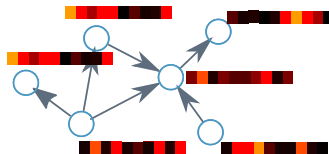
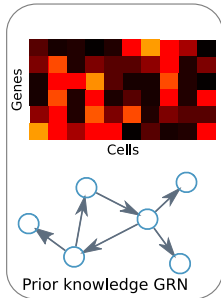
## Inputs



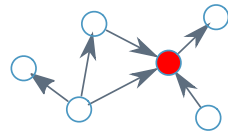


# Denoising the input data

## Inputs

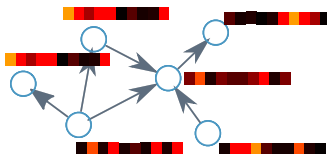
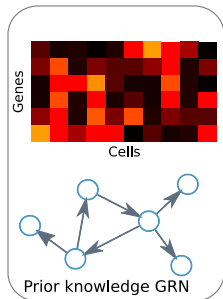


Update expression vector for gene in red

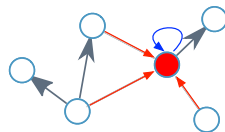


# Denoising the input data

## Inputs

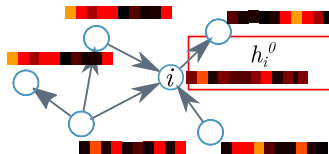
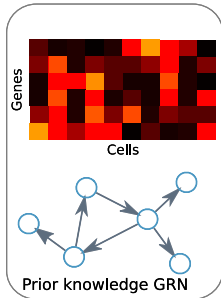


Update expression vector for gene in red

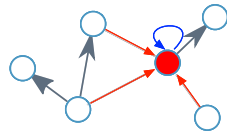


# Denoising the input data

## Inputs

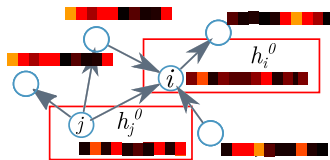
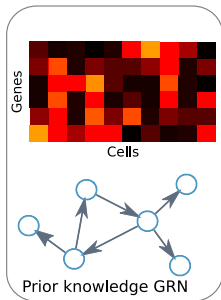


Update expression vector for gene in red

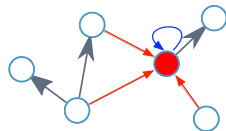


# Denoising the input data

## Inputs

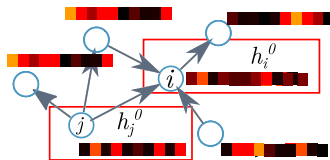
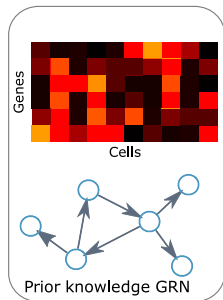


Update expression vector for gene in red

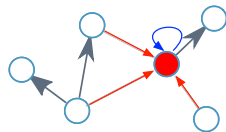


# Denoising the input data

## Inputs



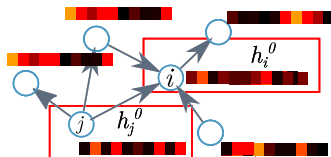
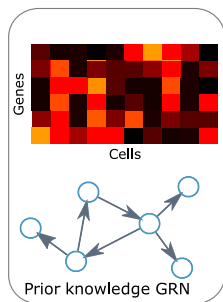
Update expression vector for gene in red



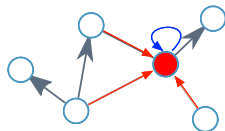
$$\text{Update: } h_i = h_i^0 w_0 + \sum_{j \in N_i} h_j^0 w_1$$

# Denoising the input data

## Inputs



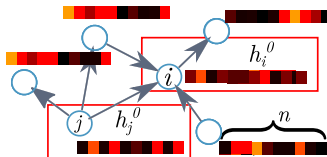
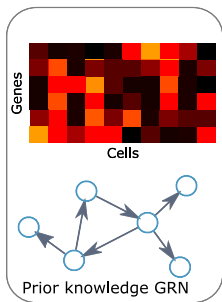
Update expression vector for gene in red



$$\text{Update: } h_i^1 = h_i^0 W_0 + \sum_{j \in N_i} h_j^0 W_1$$

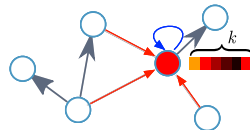
# Reducing dimensions of the input data

## Inputs



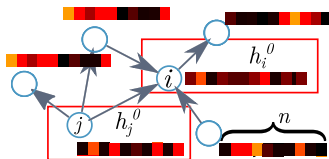
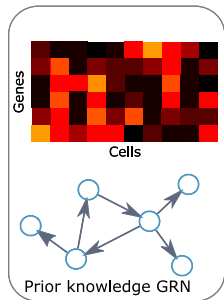
$$\text{Update: } h_i^1 = h_i^0 W_0 + \sum_{j \in N_i} h_j^0 W_1$$

Update expression vector for gene in red



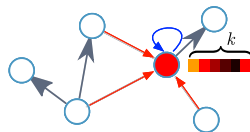
$$|W| = n \times k, k < n$$

## Inputs



$$\text{Update: } h_i^1 = h_i^0 W_0 + \sum_{j \in N_i} h_j^0 W_1$$

Update expression vector for gene in red

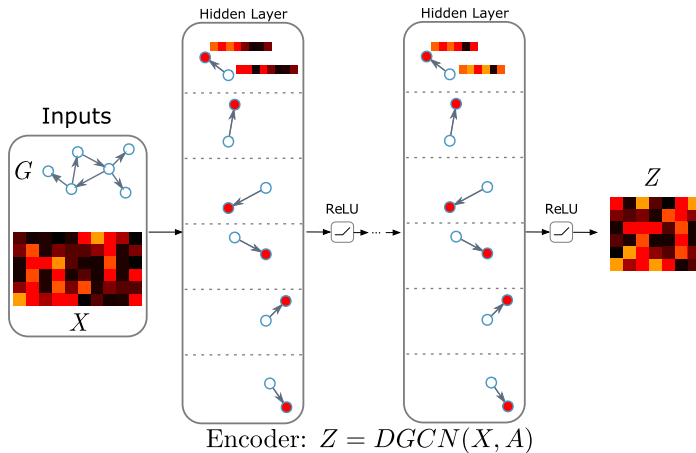


$$|W| = n \times k, k < n$$

## Graph convolutional networks (GCNs)

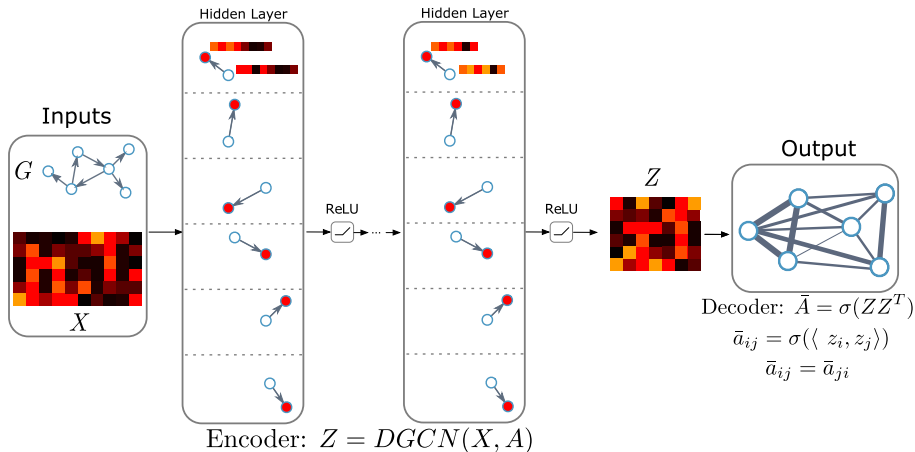


# GCN-based Autoencoders



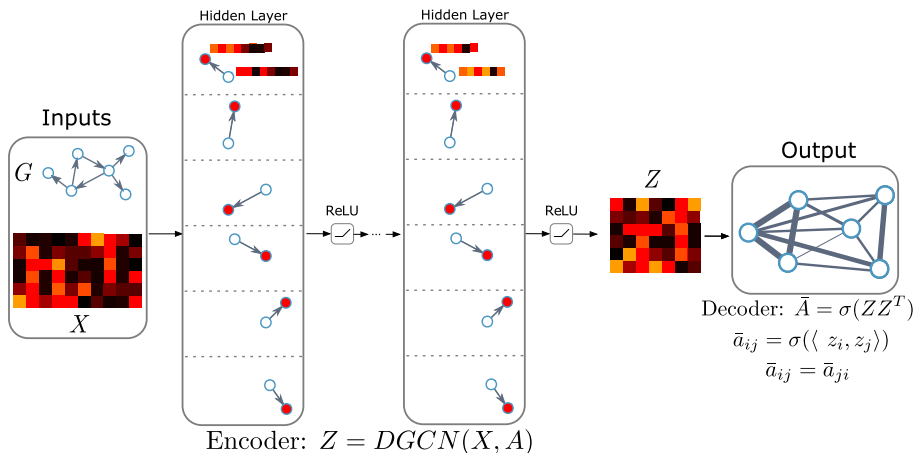
Kipf et al. (2016) "Semi-supervised classification with graph ...", CoRR, 1609.02907.

# GCN-based Autoencoders



Kipf et al. (2016) "Semi-supervised classification with graph ...", CoRR, 1609.02907.

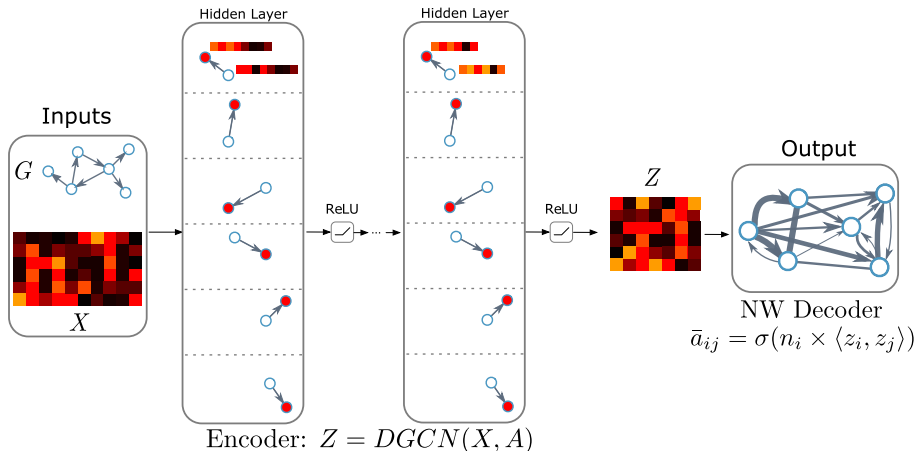
# GCN-based Autoencoders



**Learning objective:** minimize the cross-entropy loss between  $A$  and  $\bar{A}$

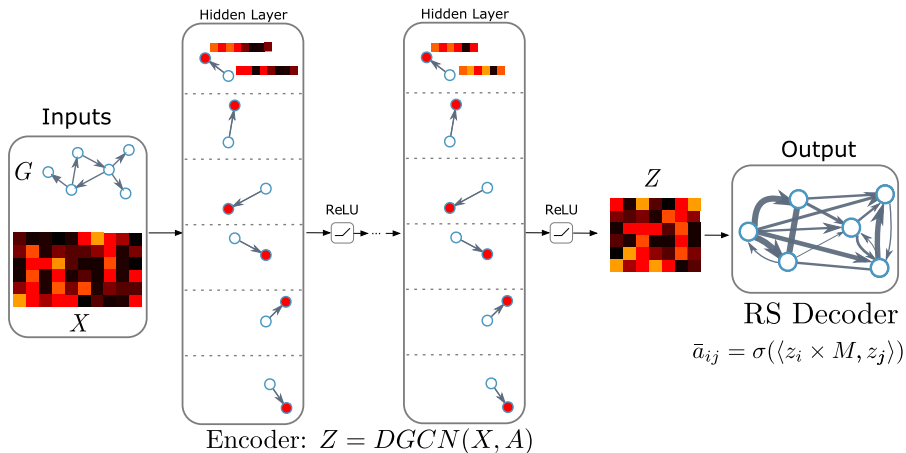
$$\mathcal{L} = \frac{1}{|E|} \left( - \sum_{(i,j) \in E} \log \bar{a}_{ij} - \sum_{(p,q) \in \bar{E}} \log(1 - \bar{a}_{pq}) \right)$$

# Directed GCN + Node Weight Decoder (NW)



$n_i$ : Learned node weight for node  $i$

# Directed GCN + RESCAL Decoder (RS)<sup>1</sup>



$M$ : Learned weight matrix

<sup>1</sup>Nickel *et al.* (2012) "Factorizing YAGO...", *In Proc. WWW*, pp. 271–280

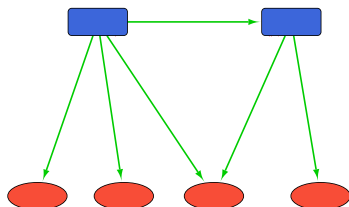
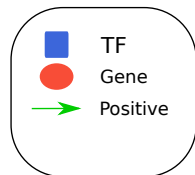
# Variants of GCN-based Autoencoders

Six encoder-decoder combinations:

- Encoders
  - ▶ Undirected GCN-based encoder (**GCN**)
  - ▶ Directed GCN-based encoder (**DGCN**)
- Decoders
  - ▶ Inner Product decoder (**IP**)
  - ▶ Node Weight decoder (**NW**)
  - ▶ RESCAL Decoder (**RS**)

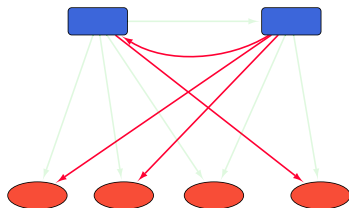
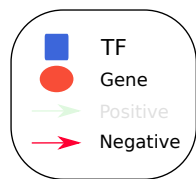
# Evaluation

- $k$ -fold cross validation: Edges, TFs
- Positive edges: TF-gene edges present in the input GRN  $G = (V, E)$



# Evaluation

- $k$ -fold cross validation: Edges, TFs
- Positive edges: TF-gene edges present in the input GRN  $G = (V, E)$
- Negative edges: TF-gene edges not in  $G$





# 10-Fold Edge Holdout Cross-Validation

- Randomly partition positive edges into 10 sets
- Holdout one set of edges as testing positives
- Use the remaining edges as training positives
- Sample uniformly at random as many training (testing) negatives as there are training (testing) positives

# 10-Fold TF Holdout Cross-Validation

- Randomly partition TF nodes in the GRN into 10 sets
- Holdout one set of TFs and all the edges adjacent to them in the GRN as testing positives
- Use the remaining edges in the GRN as training positives
- How do we sample negatives?

# 10-Fold TF Holdout Cross-Validation

- Randomly partition TF nodes in the GRN into 10 sets
- Holdout one set of TFs and all the edges adjacent to them in the GRN as testing positives
- Use the remaining edges in the GRN as training positives
- How do we sample negatives? For each TF, we randomly sample as many negatives as there are positives adjacent to that TF, once for the set training and once for the testing set.

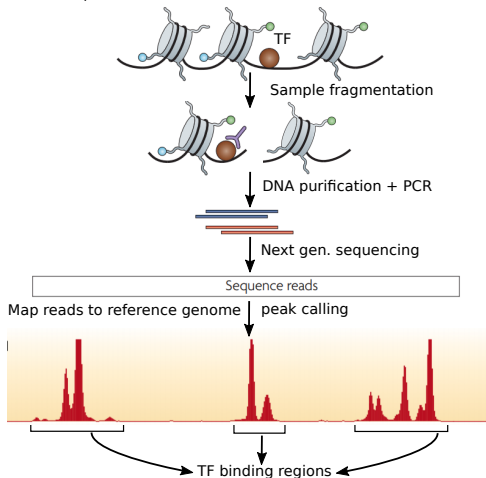
# Ground-Truth Networks

- Cell-type specific ChIP-seq network
- Non-specific ChIP-seq network

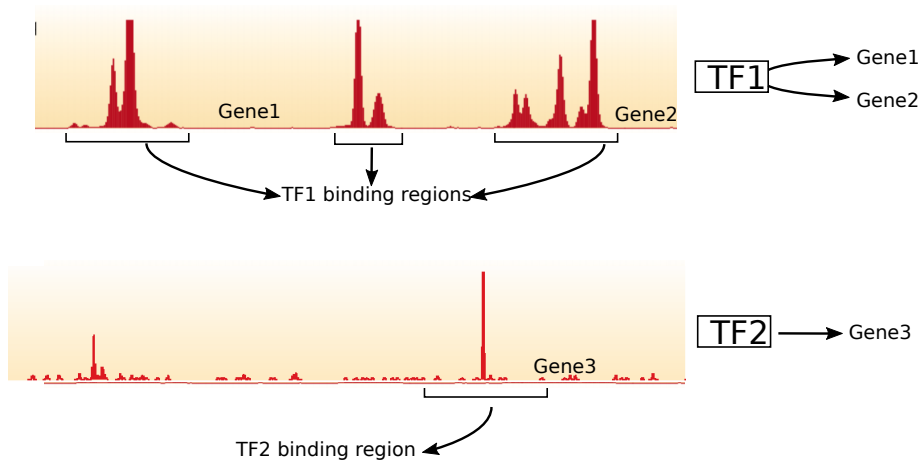
# ChIP-seq

**Chromatin:** any protein interacting with DNA, e.g., TF

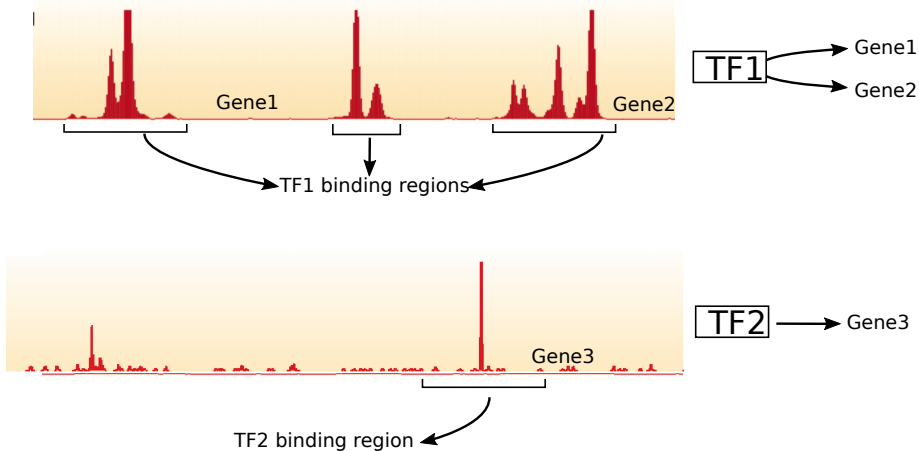
**ImmunoPrecipitation:** enrichment of DNA bound to the protein of interest



# ChIP-seq



# Cell-type specific ChIP-seq network



**Very few TFs tested for each cell-type**

# Non-specific ChIP-seq network

- Collected curated TF-gene interactions from
  - 1 RegNetwork<sup>1</sup>
  - 2 TRRUST<sup>2</sup>
  - 3 DoRothEA<sup>3</sup>



1. Liu et al. (2015) "RegNetwork: an integrated ..." Database, 2015
2. Han et al. (2018) "TRRUSTv2 ..." NAR, 46(D1):D380–D386
3. Garcia-Alonso et al. (2019) "Benchmark ..." Gen. Res., 29:1363–1375



# Gene Expression Datasets

Name	#Cells	#Nodes	#Edges	# TFs
mESC <sup>1</sup>	471	896	6,893	516
mHSC <sup>2</sup>	3,175	4,158	17,309	445
mMac <sup>3</sup>	6,283	7,428	35,347	747
hESC <sup>4</sup>	758	1,142	4,597	292

---

<sup>1</sup>Hayashi *et al.* (2018) "Single-cell full-length..." *Nat Comm*, 9, 619

<sup>2</sup>Nestorowa *et al.* (2016) "A Single-Cell Resolution..." *Blood*, 128(8):e20-31

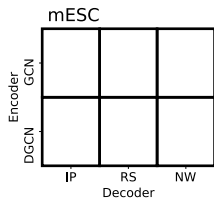
<sup>3</sup>Alavi *et al.* (2018) "A web server for..." *Nat Comm*, 9, 4768

<sup>4</sup>Chu *et al.* "Single-cell RNA-seq reveals ... *Genome Biology*, 17(1), 173

# Best GCN-based autoencoder architecture

- Median test early precision from 10-fold evaluations

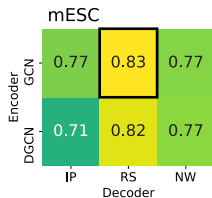
(a) 10-fold edge CV



# Best GCN-based autoencoder architecture

- Median test early precision from 10-fold evaluations
- **GCN-RS** performs the best for 10-fold edge cross-validation

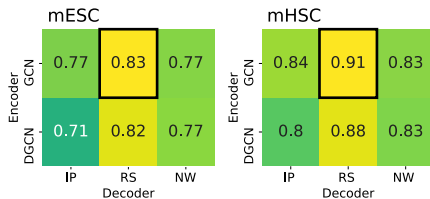
(a) 10-fold edge CV



# Best GCN-based autoencoder architecture

- Median test early precision from 10-fold evaluations
- **GCN-RS** performs the best for 10-fold edge cross-validation

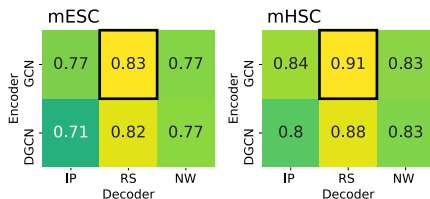
(a) 10-fold edge CV



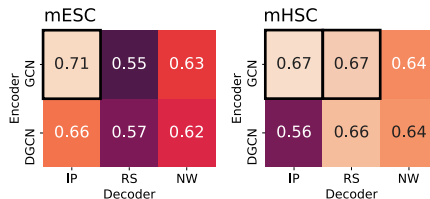
# Best GCN-based autoencoder architecture

- Median test early precision from 10-fold evaluations
- **GCN-RS** performs the best for 10-fold edge cross-validation
- **GCN-IP** performs the best for 10-fold TF cross-validation

(a) 10-fold edge CV



(b) 10-fold TFCV



## Methods evaluated

- For 10-fold edge CV: **GCN-RS** autoencoder
- For 10-fold TF CV: **GCN-IP** autoencoder

---

<sup>1</sup>Yuan (2020) “Deep learning for inferring . . .” *PNAS*, 116 (52) 27151-27158

## Methods evaluated

- For 10-fold edge CV: **GCN-RS** autoencoder
- For 10-fold TF CV: **GCN-IP** autoencoder
- **CNNC**<sup>1</sup>: CNN-based method that uses normalized empirical probability function (NEPDF) as features for every pair of genes
- **MLP-C**: a multi-layer perceptron with as many hidden layers as in the GCN and with concatenated expression vectors as input features
- **SVM-C**: Linear SVM with concatenated expression vectors as input features

---

<sup>1</sup>Yuan (2020) “Deep learning for inferring . . .” *PNAS*, 116 (52) 27151-27158

## Methods evaluated

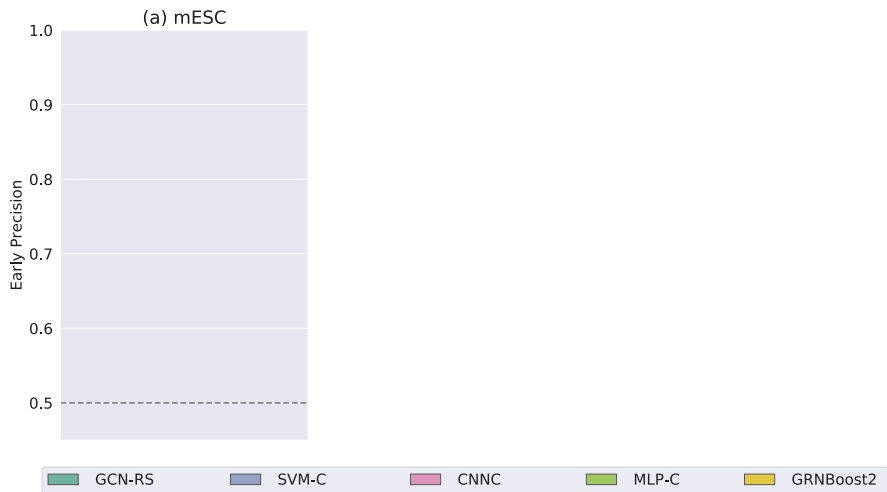
- For 10-fold edge CV: **GCN-RS** autoencoder
- For 10-fold TF CV: **GCN-IP** autoencoder
- **CNNC**<sup>1</sup>: CNN-based method that uses normalized empirical probability function (NEPDF) as features for every pair of genes
- **MLP-C**: a multi-layer perceptron with as many hidden layers as in the GCN and with concatenated expression vectors as input features
- **SVM-C**: Linear SVM with concatenated expression vectors as input features
- **GRNBoost2**: One of the top performing unsupervised learning methods from BEELINE (baseline)

---

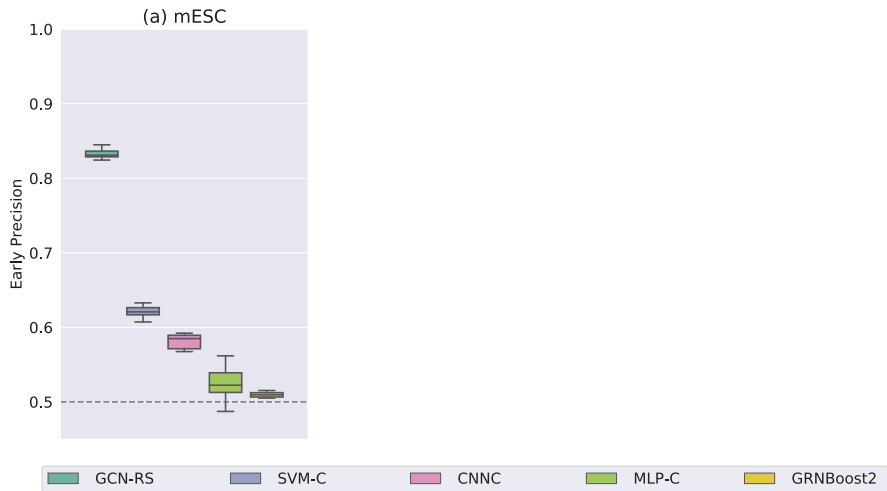
<sup>1</sup>Yuan (2020) “Deep learning for inferring . . .” *PNAS*, 116 (52) 27151-27158



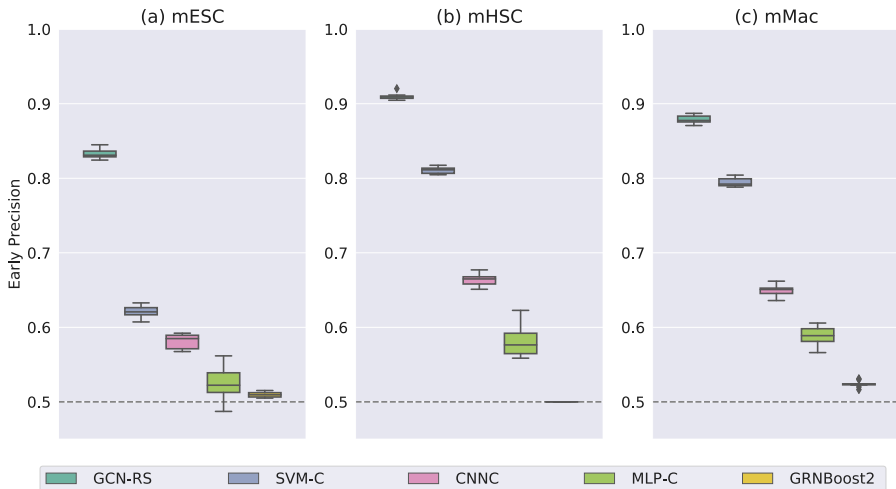
# Evaluation: 10-fold edge cross-validation



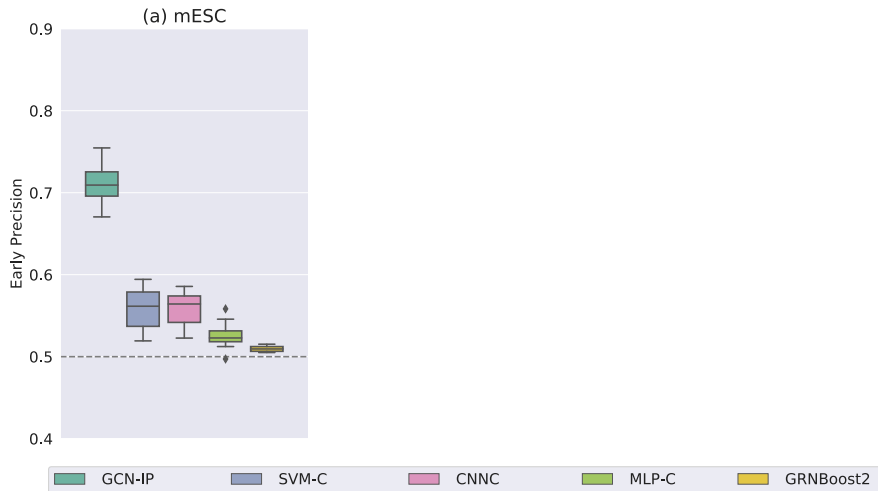
# Evaluation: 10-fold edge cross-validation



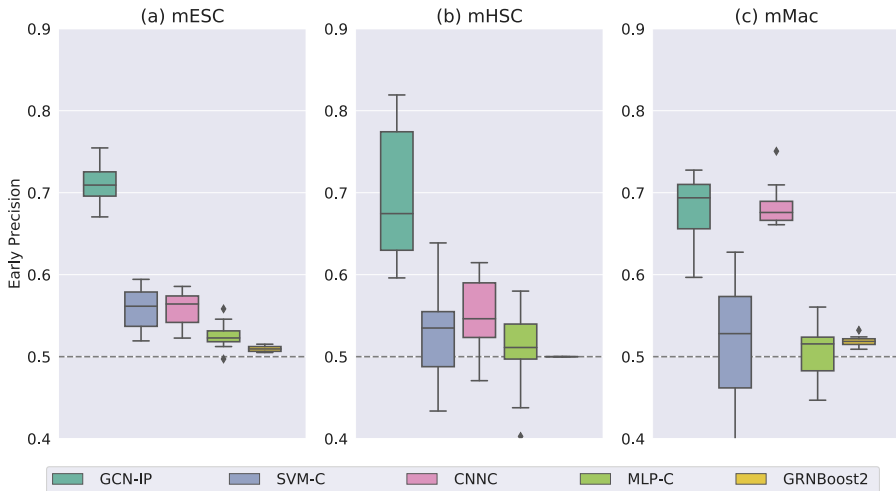
# Evaluation: 10-fold edge cross-validation



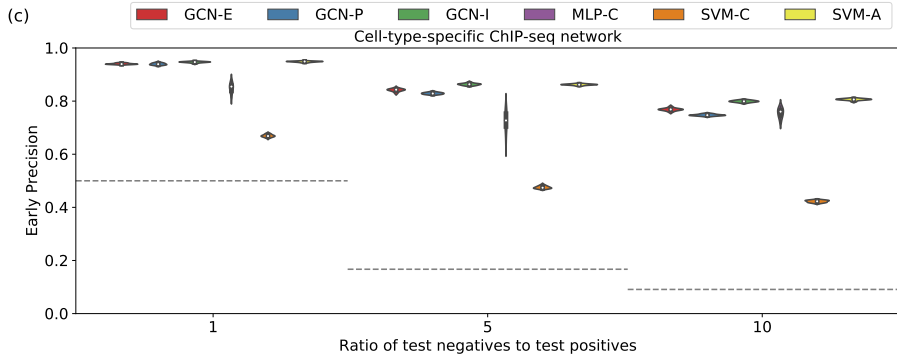
# Evaluation: 10-fold TF-holdout cross-validation



# Evaluation: 10-fold TF-holdout cross-validation



# Evaluation: mESC ChIP-seq network



# Case Study: hESC scRNA-seq dataset <sup>3</sup>

- Human embryonic stem cell dataset (hESC)

Dataset	#Cells	#Nodes	#Edges	# TFs
hESC	758	1,142	4,597	292

- Ground-truth network: Non-cell-type specific ChIP-seq network<sup>1 2</sup>
- Training set-up:
  - ▶ **GCN-RS-E**
  - ▶ Positives: Edges in the human non-specific ChIP-seq network
  - ▶ Negatives: All possible TF-gene edges that not in the human non-specific ChIP-seq network

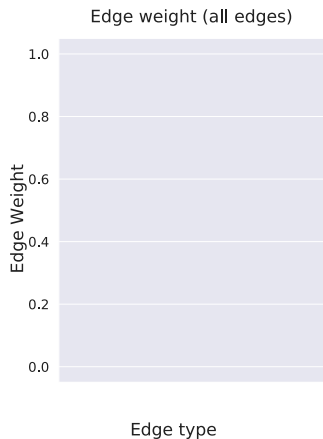
---

<sup>1</sup>Liu et al. (2015) "RegNetwork: an integrated ..." Database, 2015

<sup>2</sup>Han et al. (2108) "TRRUSTv2 ..." Nucleic Acids Res., 46(D1):D380–D386

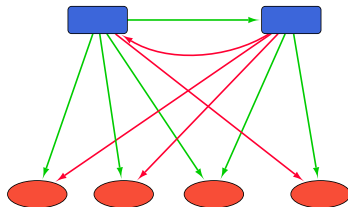
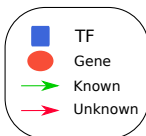
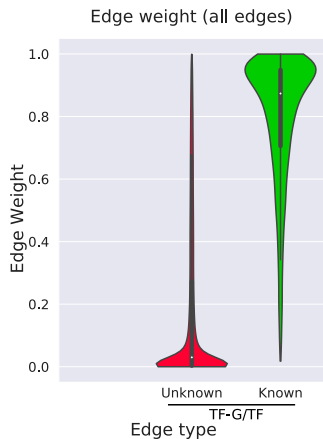
<sup>3</sup>hesc-single-cell-genbio-december-2016

# Predicted network: Edge weight distribution

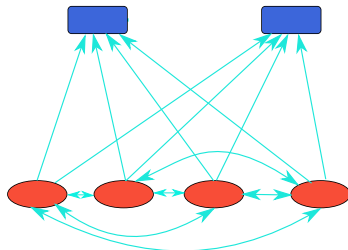
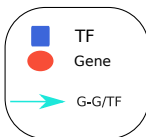
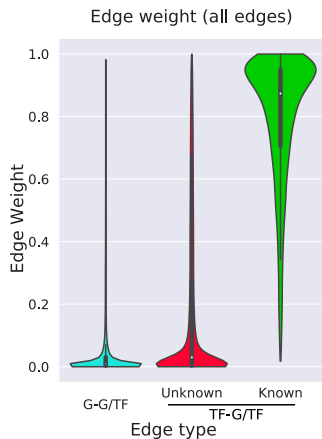




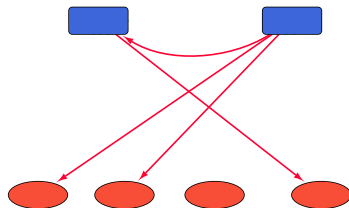
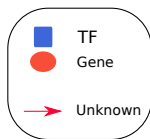
# Predicted network: Edge weight distribution



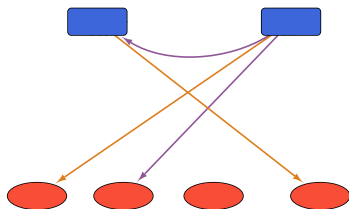
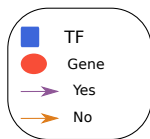
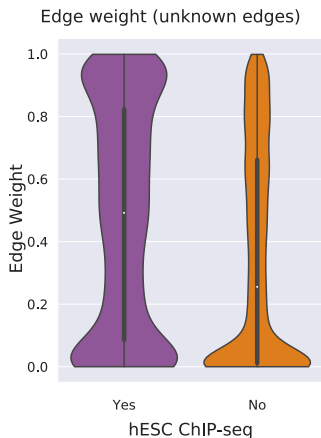
# Predicted network: Edge weight distribution



# Unknown Edges: hESC cell-type specific network



# Unknown Edges: hESC cell-type specific network



# Summary

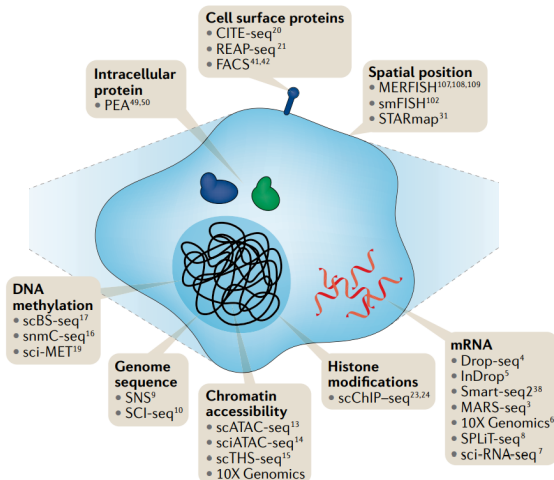
- GCN-based autoencoders are useful for denoising and reducing dimensions
- We use GCN-based autoencoders to train model for supervised GRN inference

# Summary

- GCN-based autoencoders are useful for denoising and reducing dimensions
- We use GCN-based autoencoders to train model for supervised GRN inference
- GCN-autoencoder outperforms other methods for supervised GRN inference
- Can identify cell-type specific regulatory interactions even when trained on non-cell type specific GRN

# Future Research

- Integrative single-cell analysis



Stuart et al. (2019) "Integrative single-cell..." *Nat Rev Genet* 20, 257–272.