# Single-cell RNA-sequencing

- **Differentially expressed gene analysis between two groups**

- Wilcoxon rank sum test between two groups for each cell type (or cluster)
→ Nonparametric approach (does not require a specific distribution of data)
→ Adjusted p-value, average_Log2FC + expression cell ratio (pct.1, pct.2)

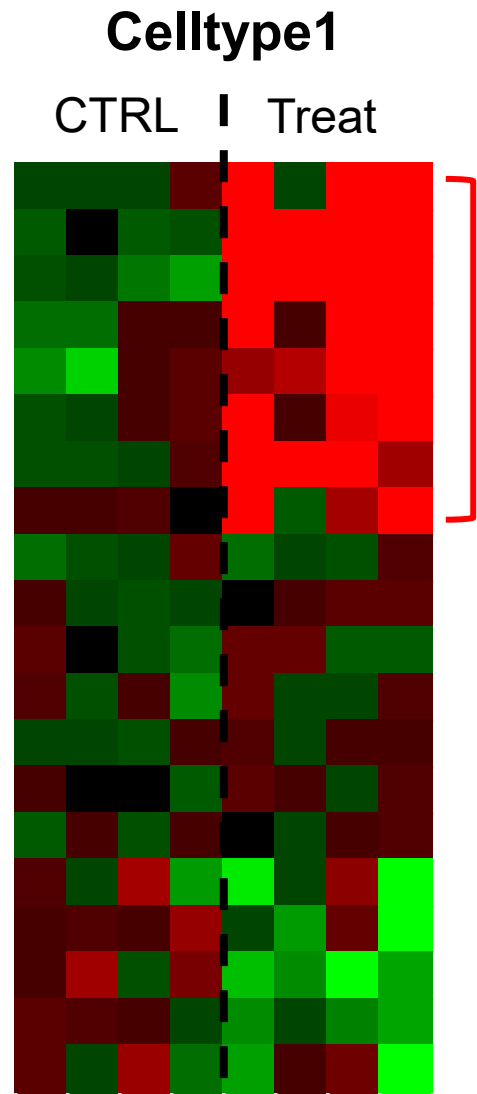| | p_val | avg_log2FC | pct.1 | pct.2 | p_val_adj |
|---|---|---|---|---|---|
| TMSB4X | 1.180041e-123 | 0.8131724 | 1.000 | 1.000 | 3.863218e-119 |
| CD74 | 1.336753e-110 | 2.4380974 | 0.841 | 0.460 | 4.376260e-106 |
| HLA-DRA | 6.658395e-84 | 2.6055562 | 0.616 | 0.121 | 2.179825e-79 |
| RPS29 | 9.853846e-84 | -0.6238238 | 0.997 | 1.000 | 3.225952e-79 |
| RPS14 | 5.520321e-81 | -0.6245685 | 0.993 | 1.000 | 1.807243e-76 |
| RGS1 | 8.877664e-78 | 1.3781774 | 0.872 | 0.485 | 2.906370e-73 |
| RPS27 | 2.936975e-74 | -0.4452566 | 1.000 | 1.000 | 9.615067e-70 |
| RPLP2 | 8.556250e-74 | -0.5271463 | 0.998 | 1.000 | 2.801145e-69 |
| DUSP4 | 1.420529e-72 | 2.1971333 | 0.553 | 0.075 | 4.650527e-68 |
| RPL3 | 5.104196e-72 | -0.6507585 | 0.982 | 0.992 | 1.671012e-67 |
| EEF1A1 | 3.198898e-69 | -0.5291642 | 0.997 | 1.000 | 1.047255e-64 |
| RPS3 | 7.501267e-69 | -0.5543631 | 0.994 | 0.998 | 2.455765e-64 |
| HLA-DRB1 | 3.762663e-62 | 1.7883327 | 0.662 | 0.273 | 1.231821e-57 |
| HLA-DPB1 | 1.903542e-61 | 1.7512534 | 0.649 | 0.275 | 6.231815e-57 |

GZMA: Tumor_CD8 T cell > Normal_CD8 T cell
→ Cytotoxic

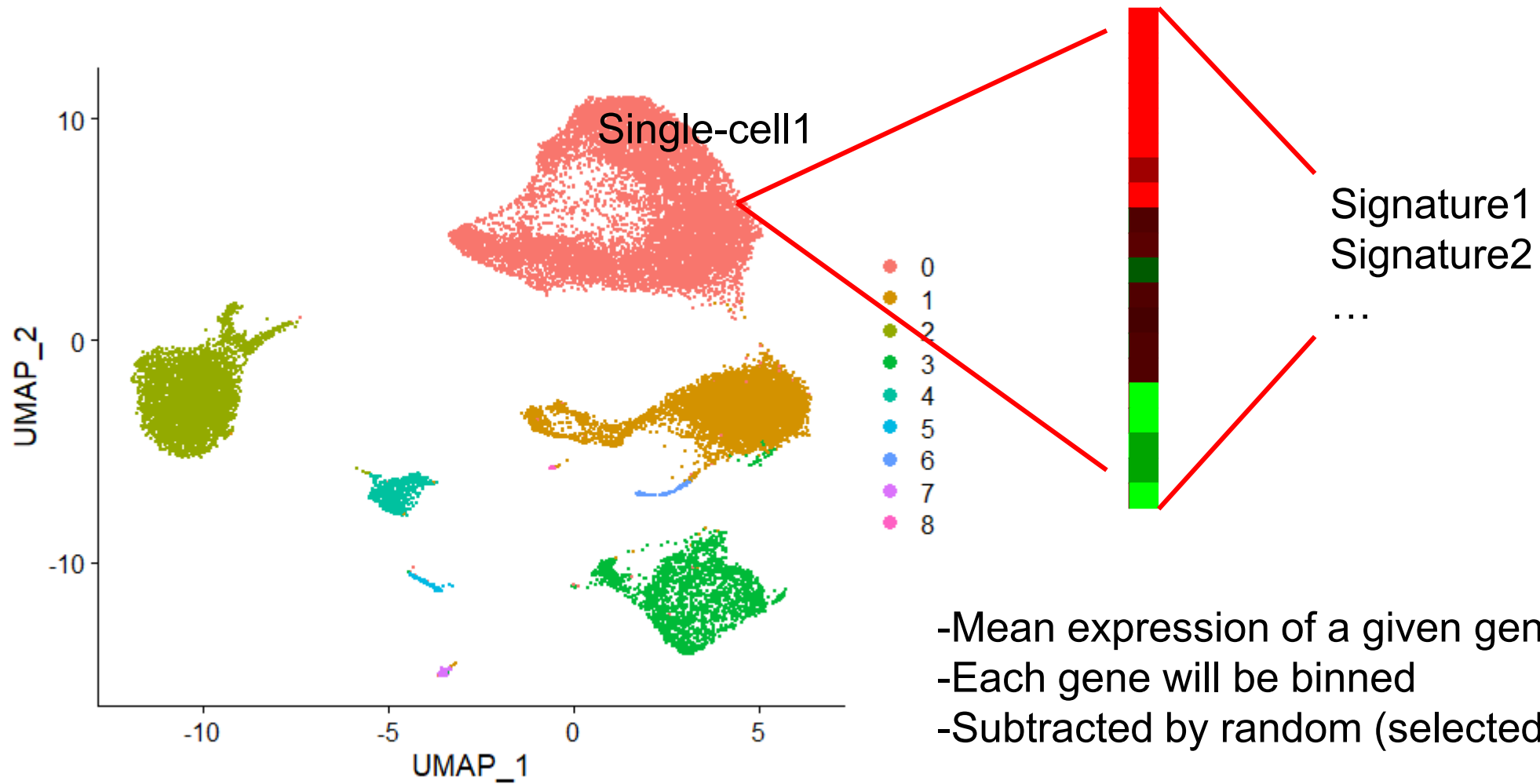HLA-DRA, HLA-DRB1, HLA-DPB1: Tumor_CD8 T cell >
Normal_CD8 T cell
→ activated

- Geneset analysis

**Celltype1**          **Celltype2 …**

CTRL    Treat

→ Fisher's exact test  or GSEA

→ Or measuring "signature score" for each cell
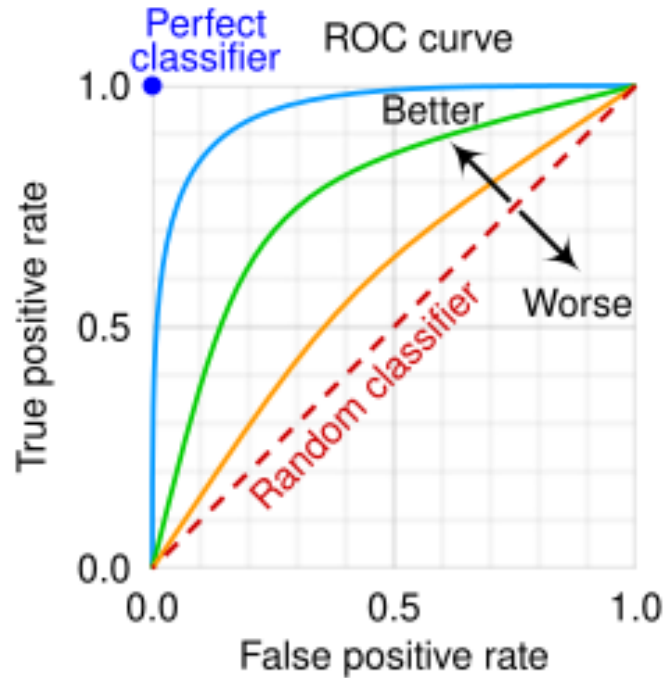→ Kind of supervised dimension reduction (?)

# AddModuleScore



- Mean expression of a given geneset
- Each gene will be binned
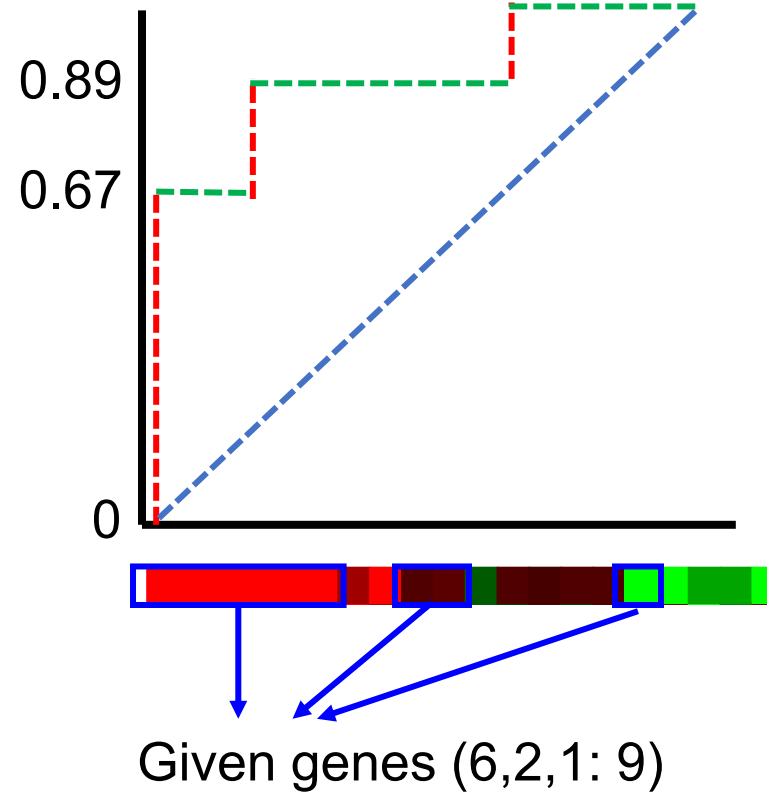- Subtracted by random (selected by each bin) noise

# AUCell

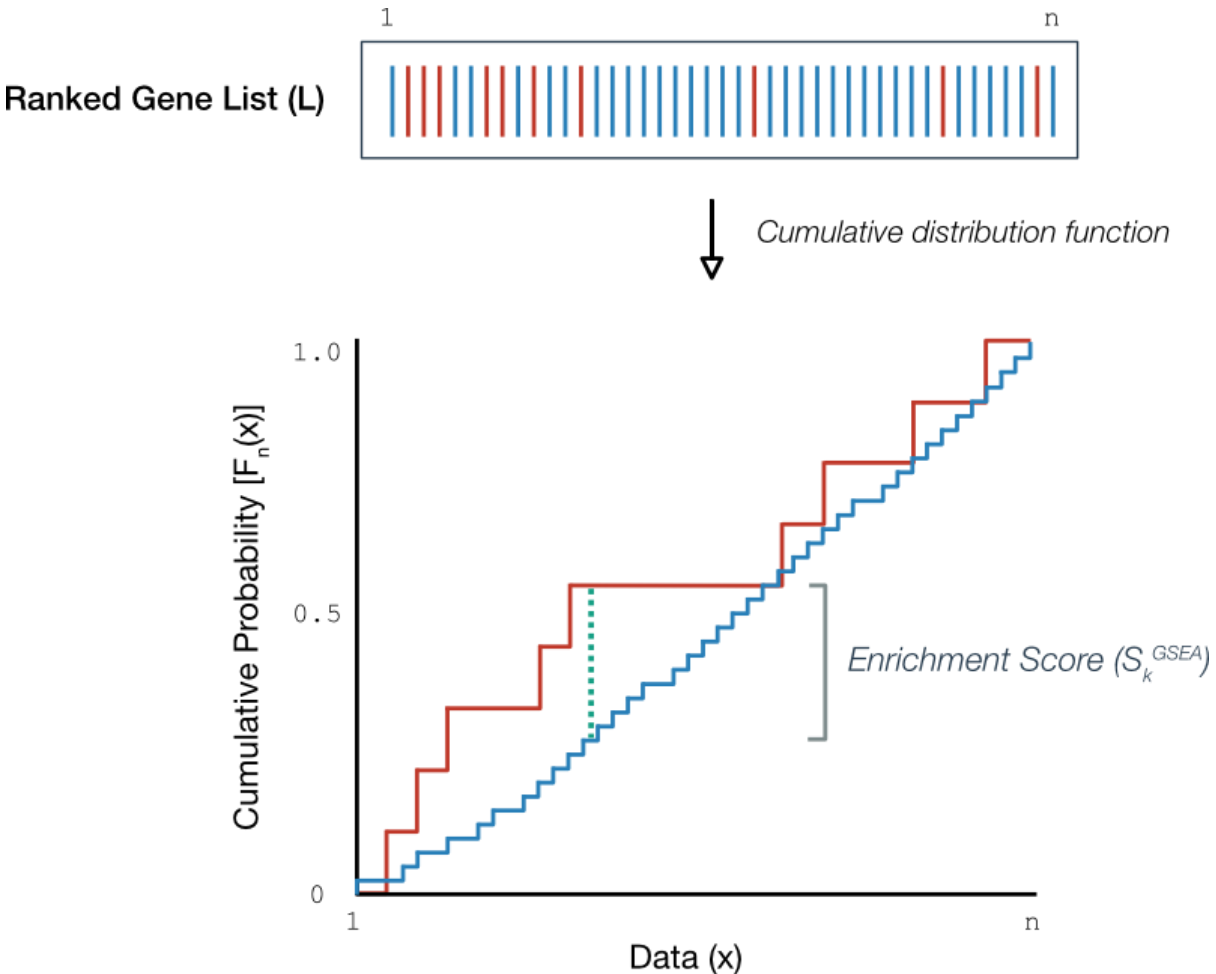-Order genes by expression for each cell
-Measure AUC for a given geneset



-Receiver operating characteristic (ROC) curve: performance measurement
-AUC: area under curve
→ Quantification

Given genes (6,2,1: 9)

- # ssGSEA & GSVA

*ssGSEA
-Order genes by expression for each cell
-Make a ECDF (empirical cumulative distribution function)
for a given geneset and remaining genes, respectively
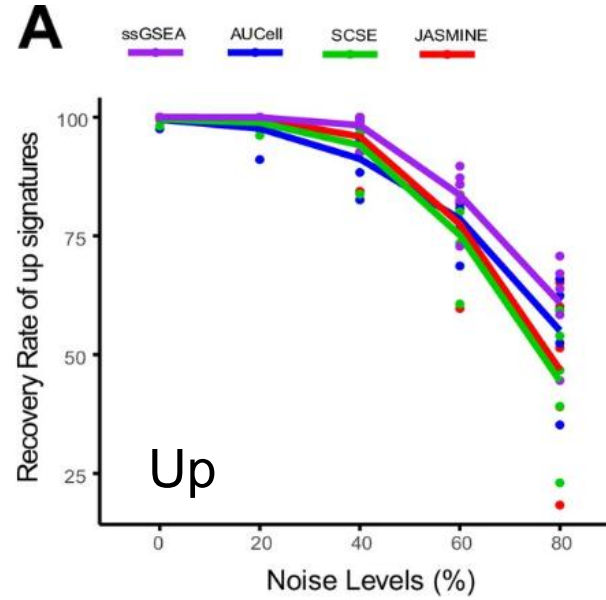-Integration of difference between two ECDFs

*GSVA
-Order genes by expression for each cell
-Make a ECDF
-KS-test (Kolmogorov-Smirnov test)
Statistics by maximum difference

# Is not always good



-Biased to cancer (upregulation >> down regulation)
-Robustness (against noise): Up > Down
-Robustness (against down-sampling): ssGSEA (worst)

Signature-scoring methods developed for bulk samples are not adequate for cancer single-cell RNA sequencing data

- ## cNMF

(consensus non-negative matrix factorization)

*Unsupervised approach
NMF: Decomposition method
Make W,H be close to X
Gradient descending

$$X = WH$$

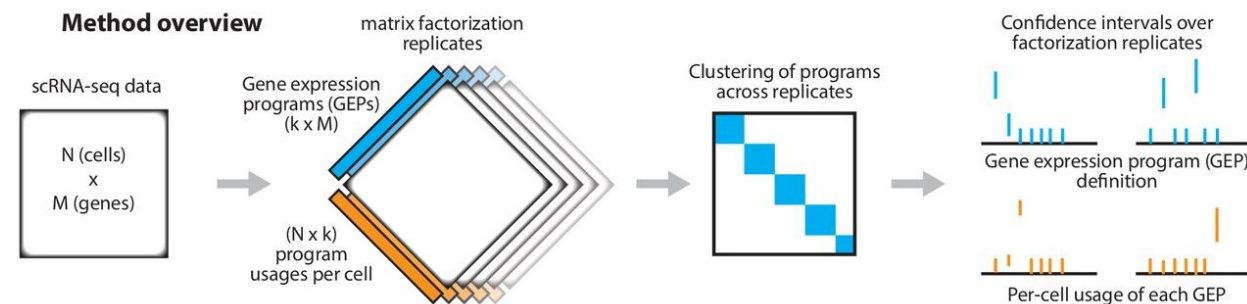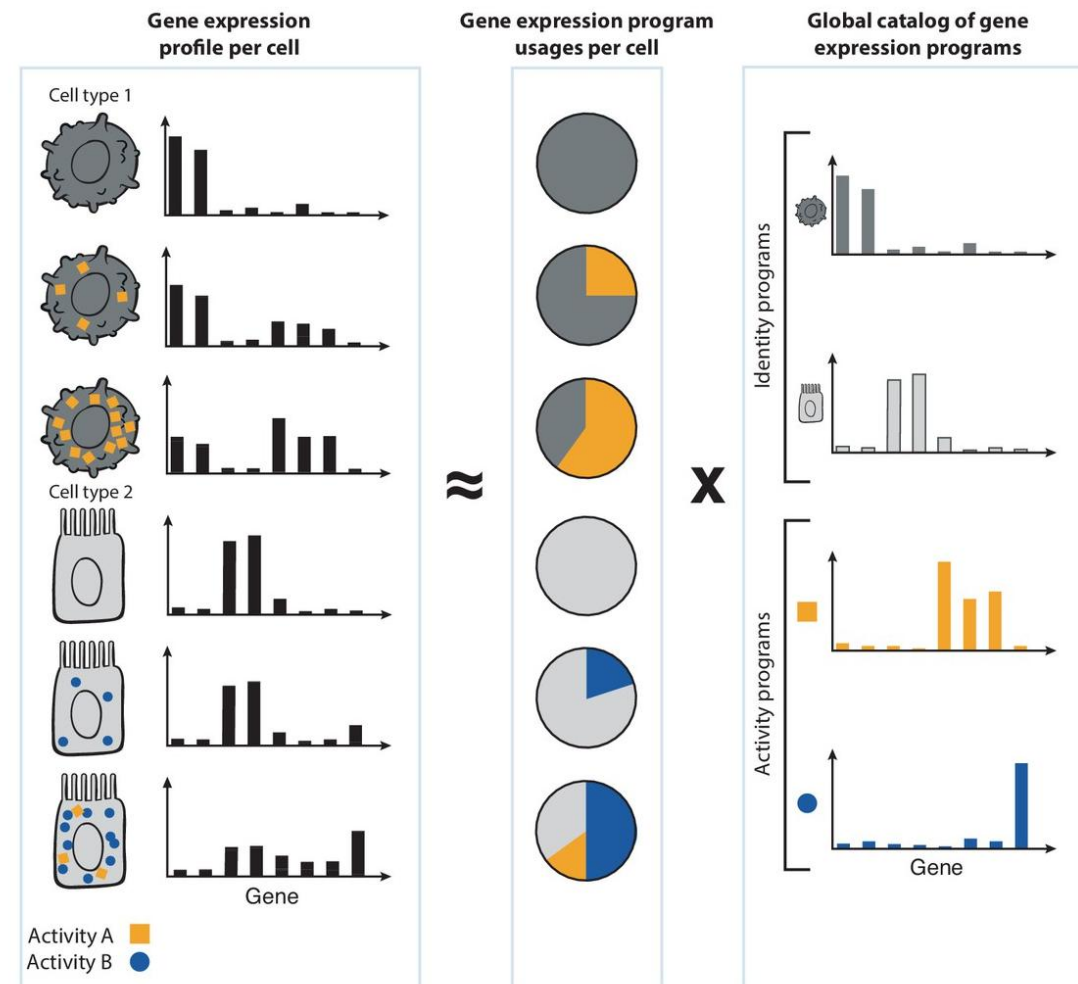$$H := H - \eta_H \circ \boxed{\nabla_H \|X - WH\|_F^2}$$

$$W := W - \eta_W \circ \boxed{\nabla_W \|X - WH\|_F^2}$$

$$\therefore H := H \circ \frac{W^T X}{W^T WH}$$

$$식 (36) \Rightarrow W := W + \frac{W}{WHH^T} \circ (XH^T - WHH^T)$$

$$= W + W \circ \frac{XH^T}{WHH^T} - W \circ \frac{WHH^T}{WHH^T} = W \circ \frac{XH^T}{WHH^T}$$

# cNMF

X = WH

-X: gene expression (N x M)
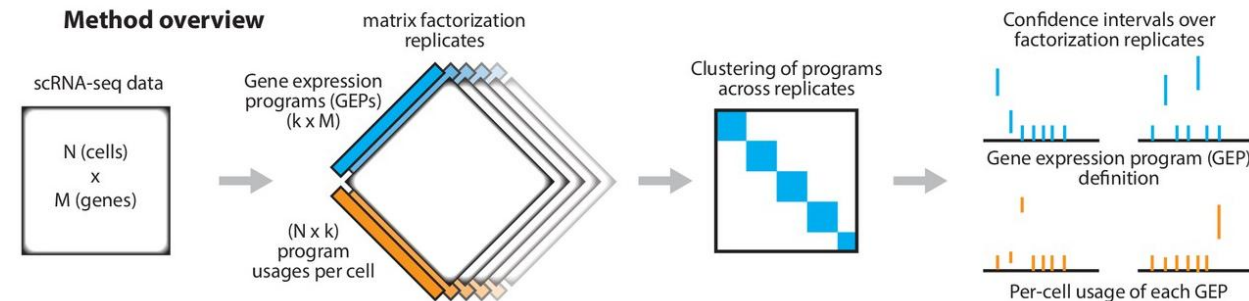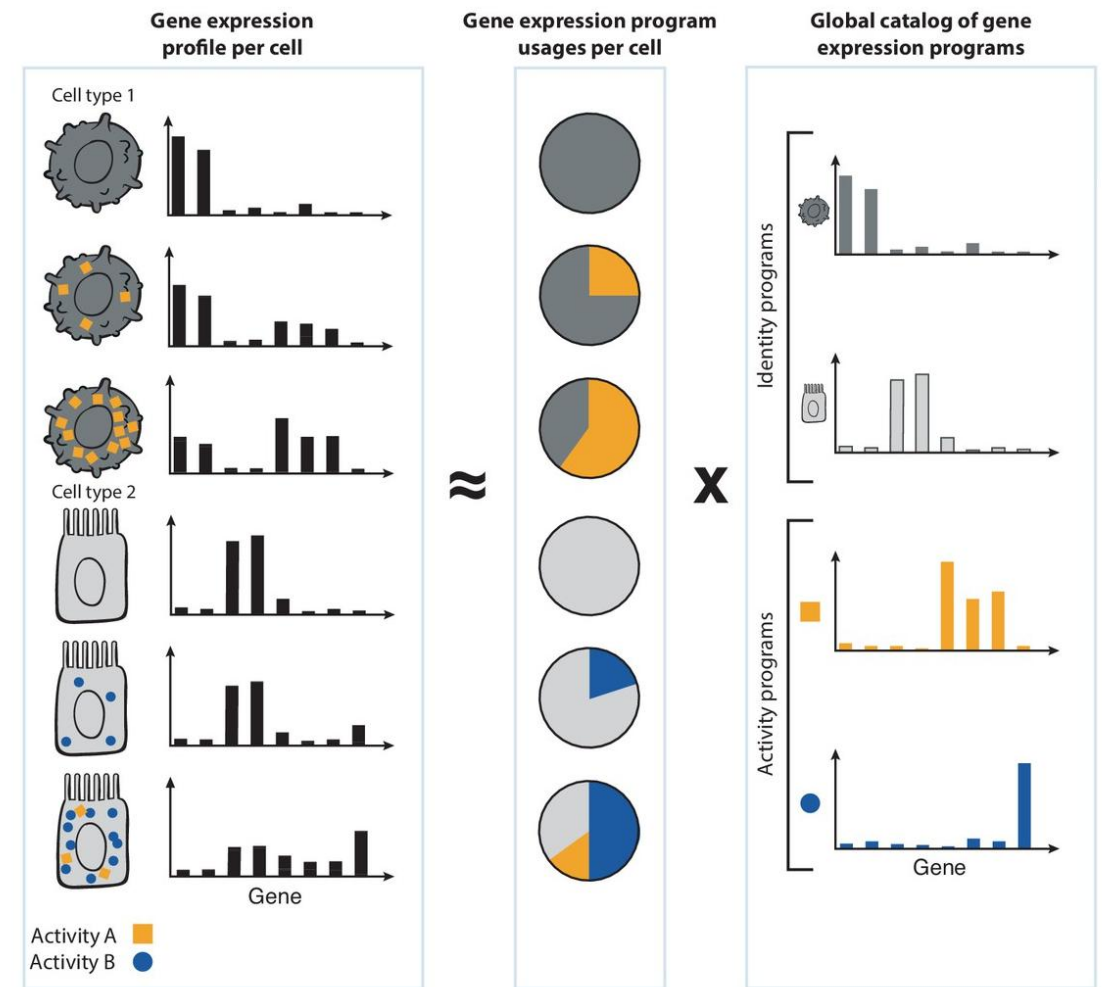N: cell, M: gene
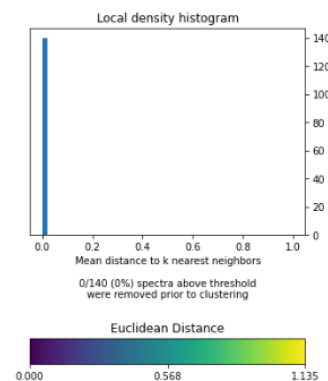
-W: program usage (activity): N x k
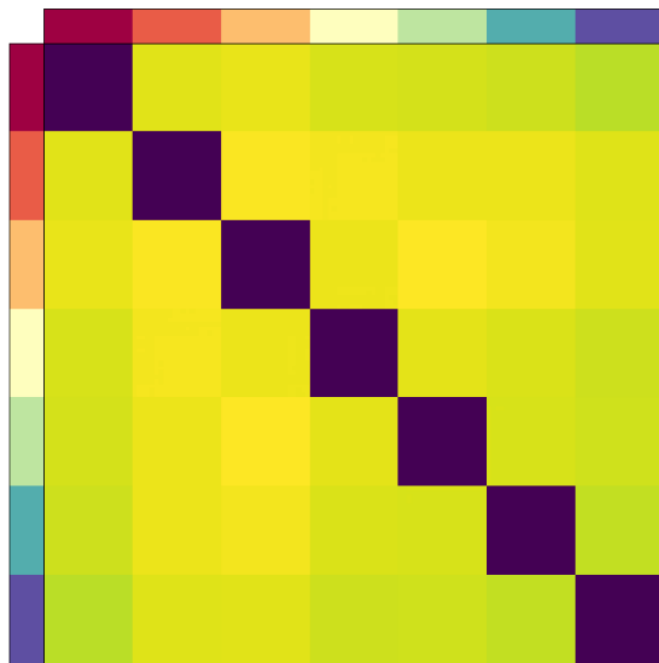k: number of program

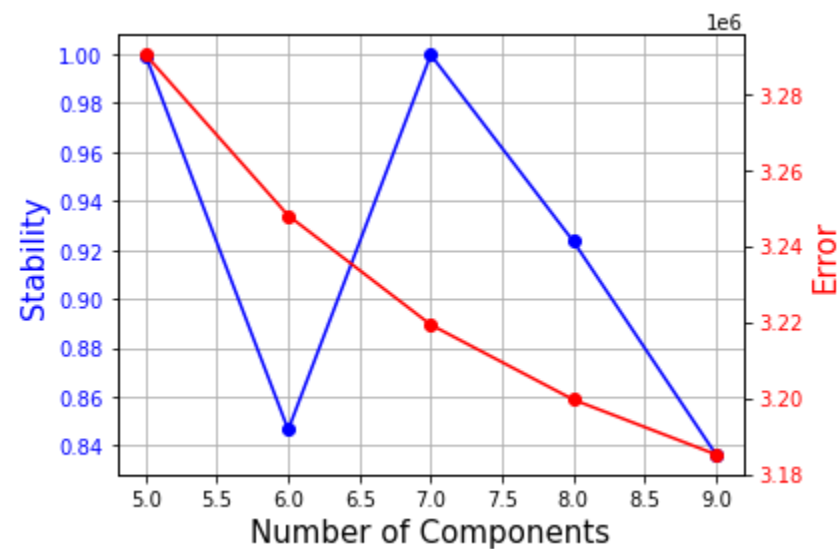-H: gene expression program: weight of each gene
k x M

Consensus → robustness
Take median value of each gene

- cNMF



-K selection → stability: high, error: low

-Define density_threshold by
KNN distance distribution

-Batch correction for input count matrix (harmony)
moe_correct_ridge (same algorithm in harmony)

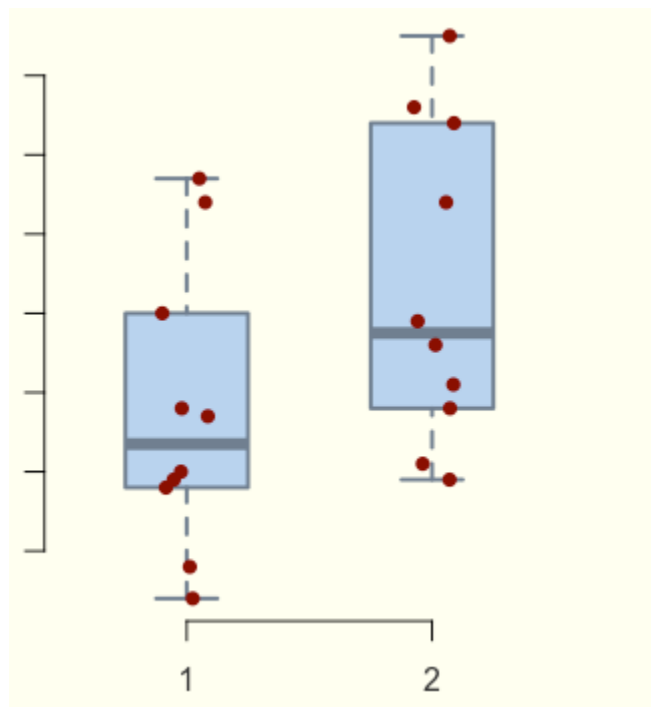# • Signature analysis

Tumor infiltrating T cell → might be exhausted

Exhaustion signature: PDCD1, CTLA4, HAVCR2, LAG3, TOX

CD8 T cells

P value: 2.595e-05

→ T cells in the tumor-microenvironment are exhausted

- ## Cell abundance

*T-test, Wilcoxon



Always! Relative abundance
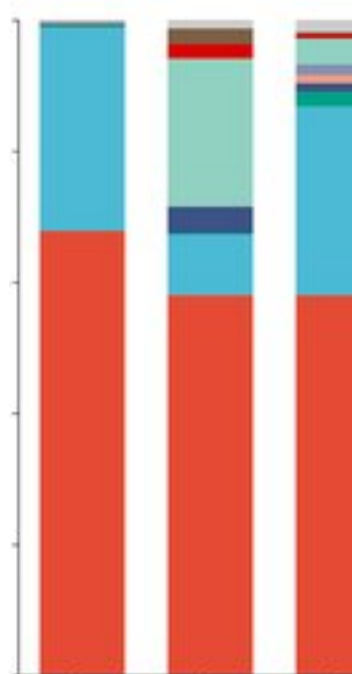Why? The cell counts for each sample is always different

Sample size is usually very small for scRNA-seq
→ Poor power analysis (less significant)

*Fisher's exact test
-Comparing by group-level
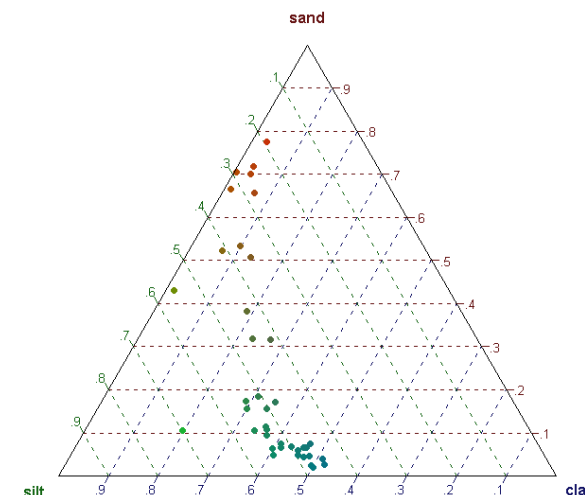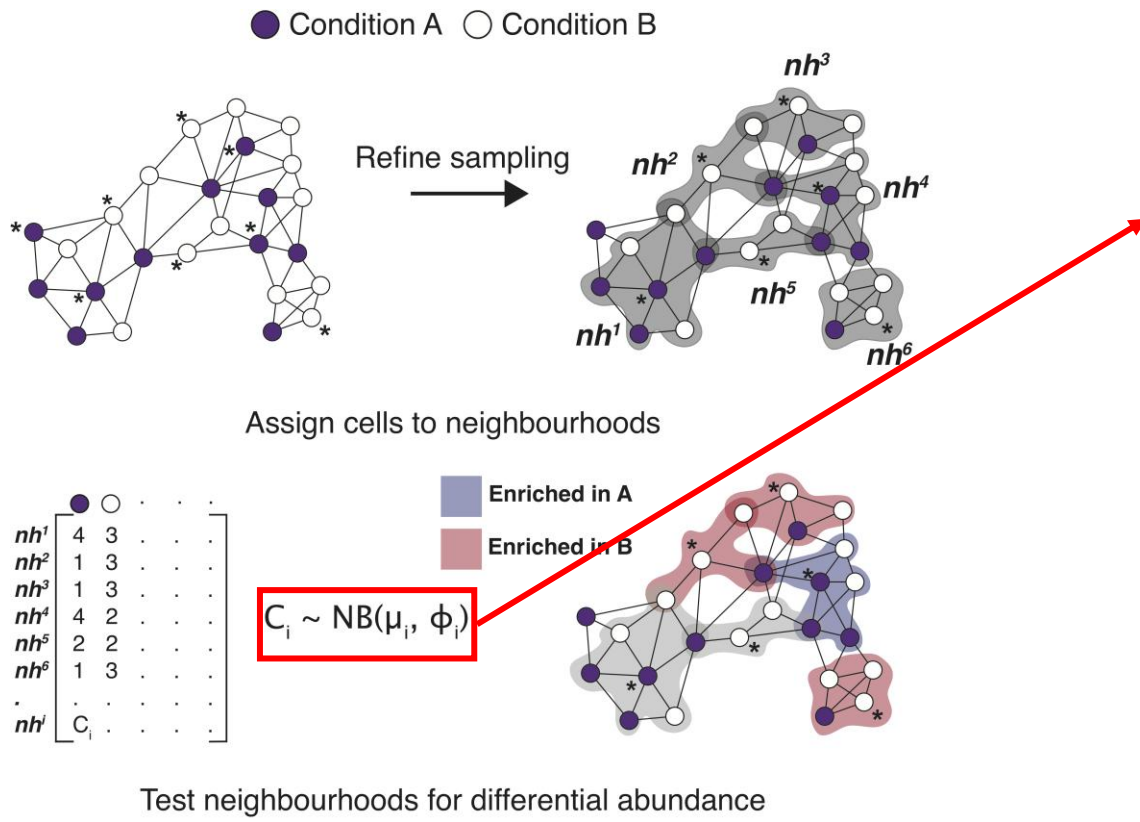-Very sensitive; high false-positive



*Dirichlet Regression
-one celltype ↑→ one celtype ↓
-**prior** reference
celltype selection

# Cell abundance (MILO)



Condition A    Condition B

Refine sampling

$nh^1$, $nh^2$, $nh^3$, $nh^4$, $nh^5$, $nh^6$

Assign cells to neighbourhoods

|        |   |   |   |   |   |
|--------|---|---|---|---|---|
| $nh^1$ | 4 | 3 | . | . | . |
| $nh^2$ | 1 | 3 | . | . | . |
| $nh^3$ | 1 | 3 | . | . | . |
| $nh^4$ | 4 | 2 | . | . | . |
| $nh^5$ | 2 | 2 | . | . | . |
| $nh^6$ | 1 | 3 | . | . | . |
| .      | . | . | . | . | . |
| $nh^i$ | $C_i$ | . | . | . | . |

$$C_i \sim NB(\mu_i, \phi_i)$$

Enriched in A

Enriched in B

Test neighbourhoods for differential abundance
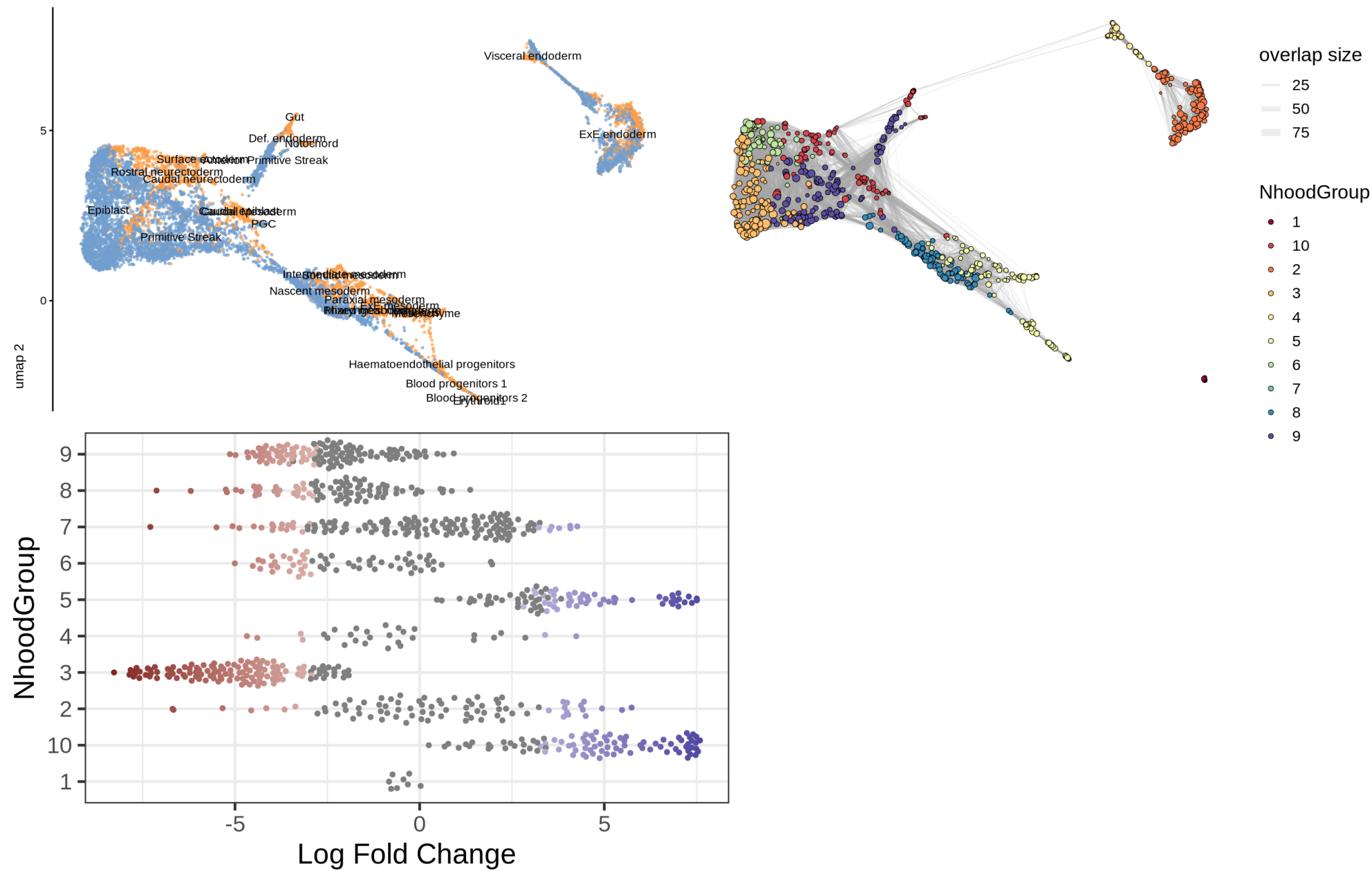
-KNN graph of cells
-Sampling to increase statistical power
-Perform enrichment test for each sampling
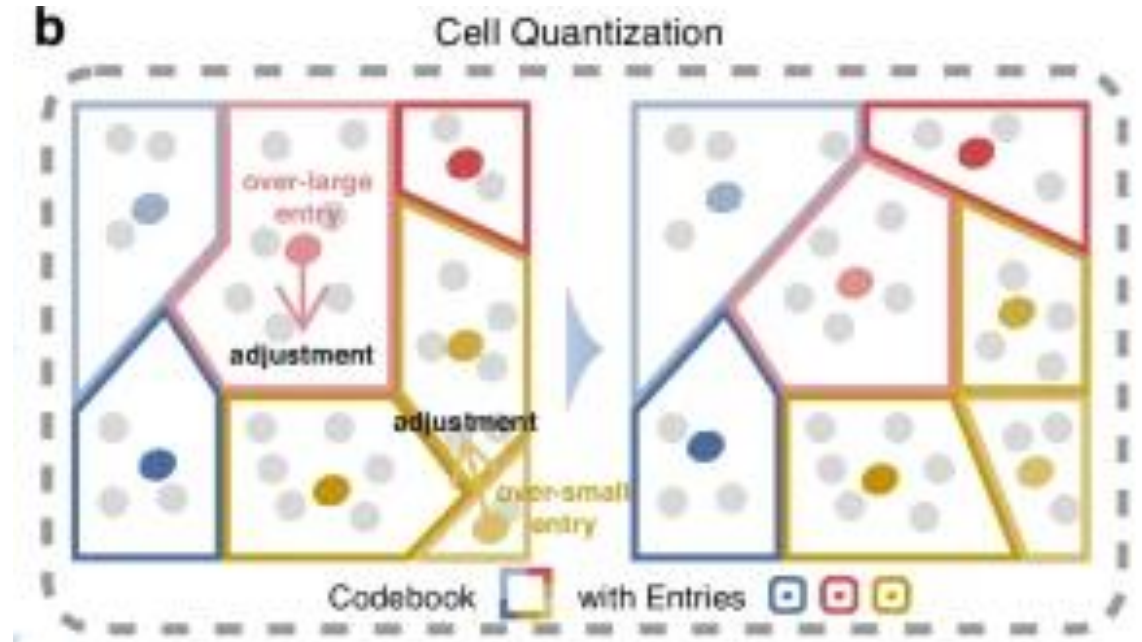Which **condition** has more in the neighborhood

# Cell abundance (MILO)

# • Cell-pooling



-High drop-out rate: zero count ↑
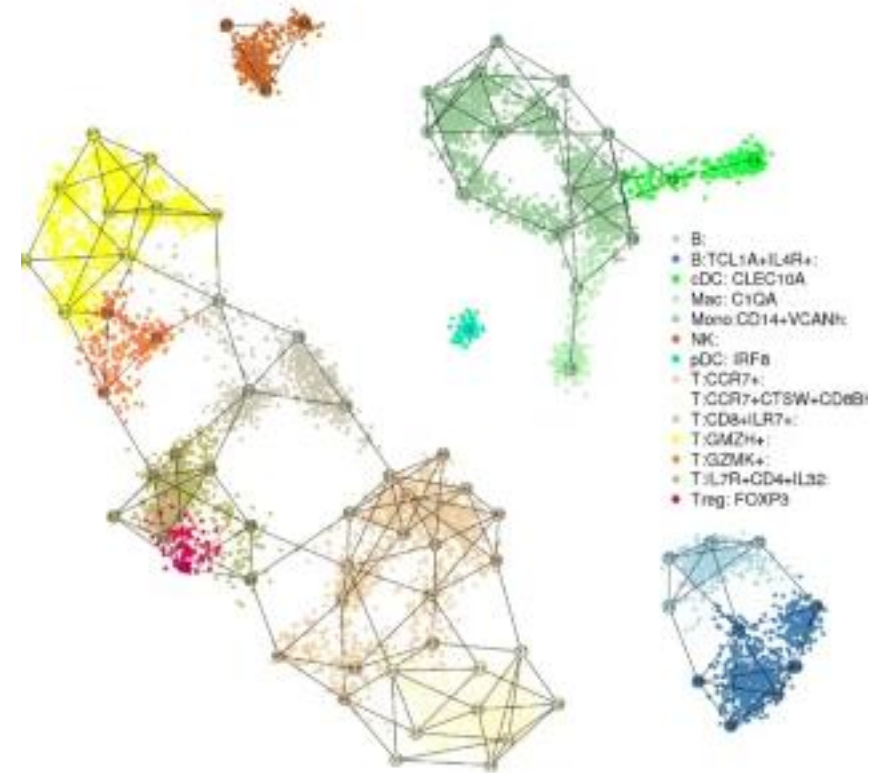-merge cells → pseudo cell → averaging →
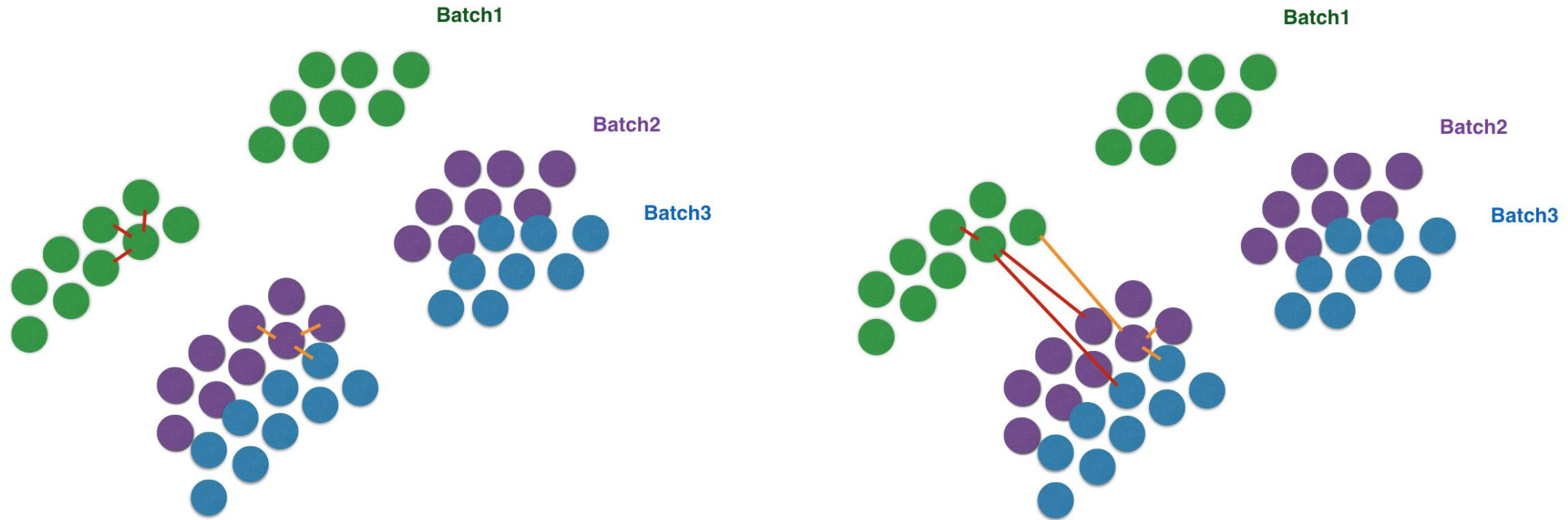overcome drop-out!

- ## Metacell



-balanced KNN graph construction
-resampling → consensus-based partitioning
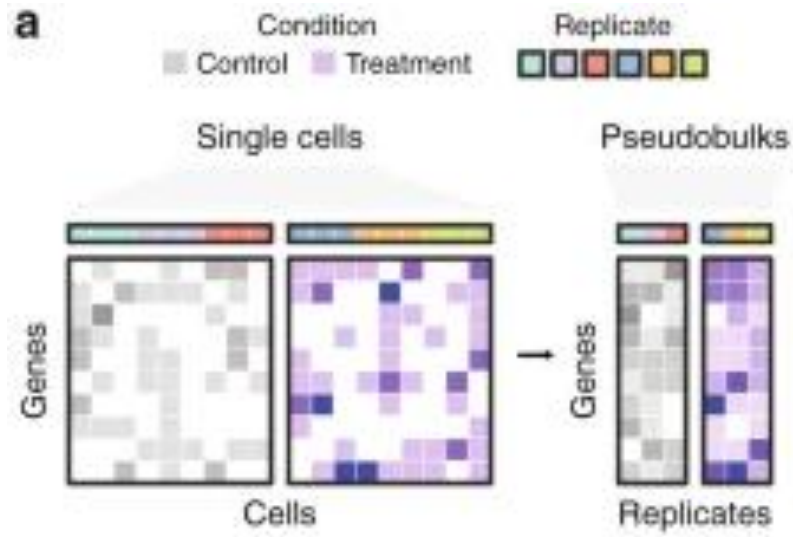-remove outliers

- **BBKNN**

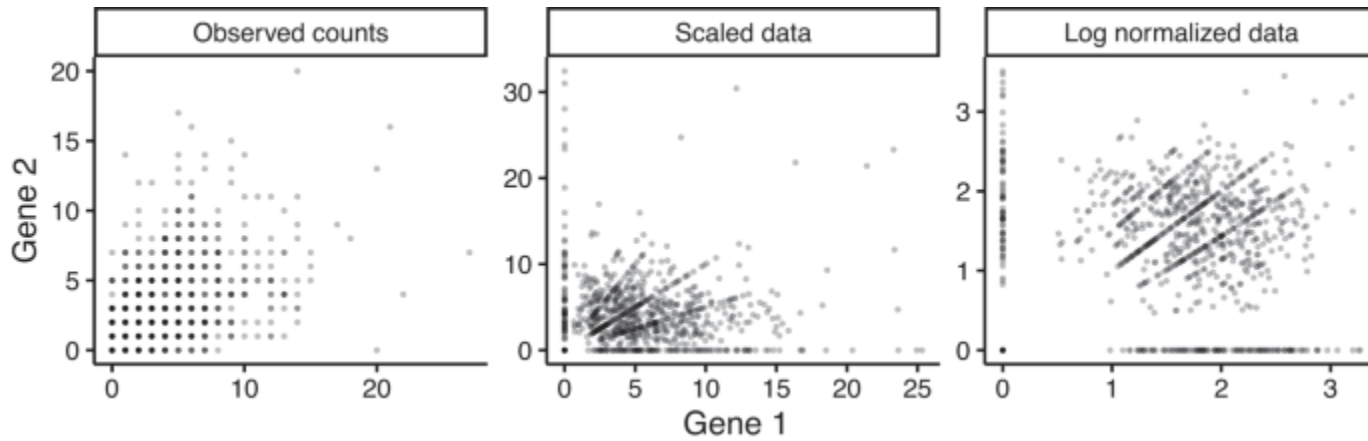Difference by cell type > difference by batch in a given cell type



- KNN for whole data
- KNN across batch with smaller k
→ Batch corrected pooling

# • Pseudobulk DE analysis



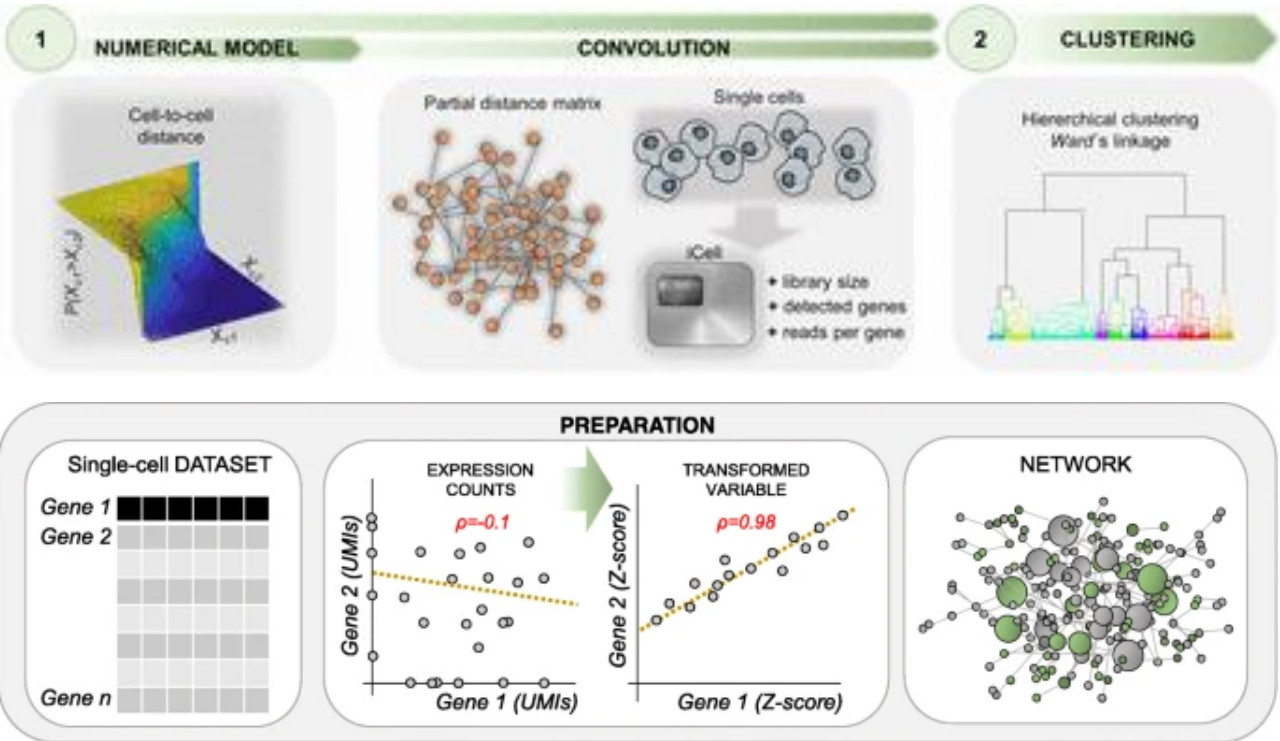-Extreme case
-Generate a pseudobulk for each sample (each cell type)
→ Perform bulk DE analysis (DESeq2, Limma, edgeR …)

-Overcome high dropout
-susceptible to outlier cells
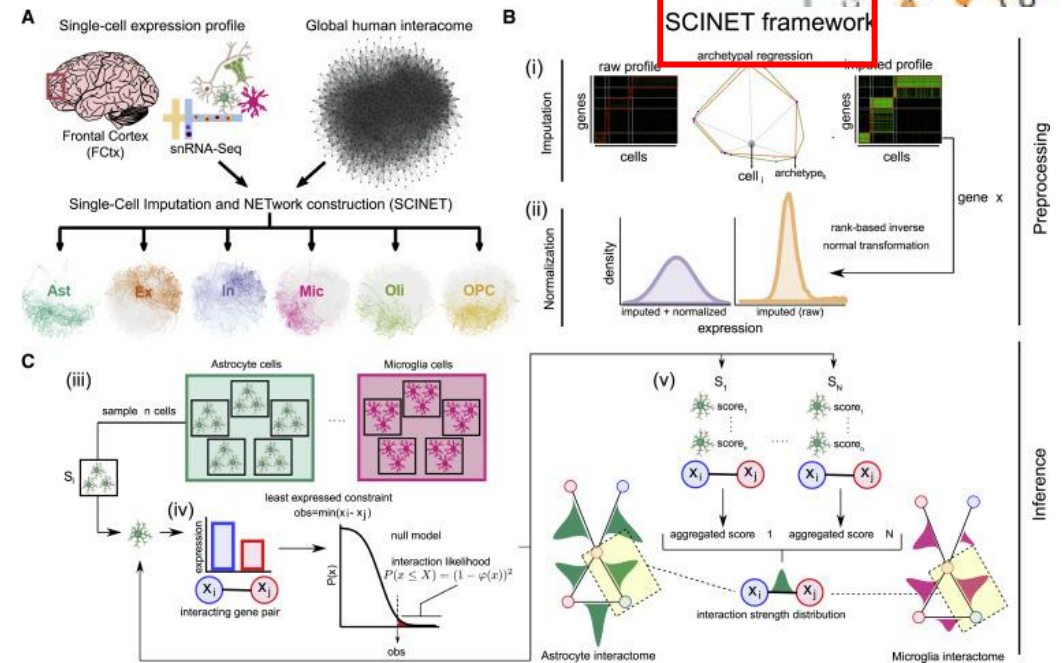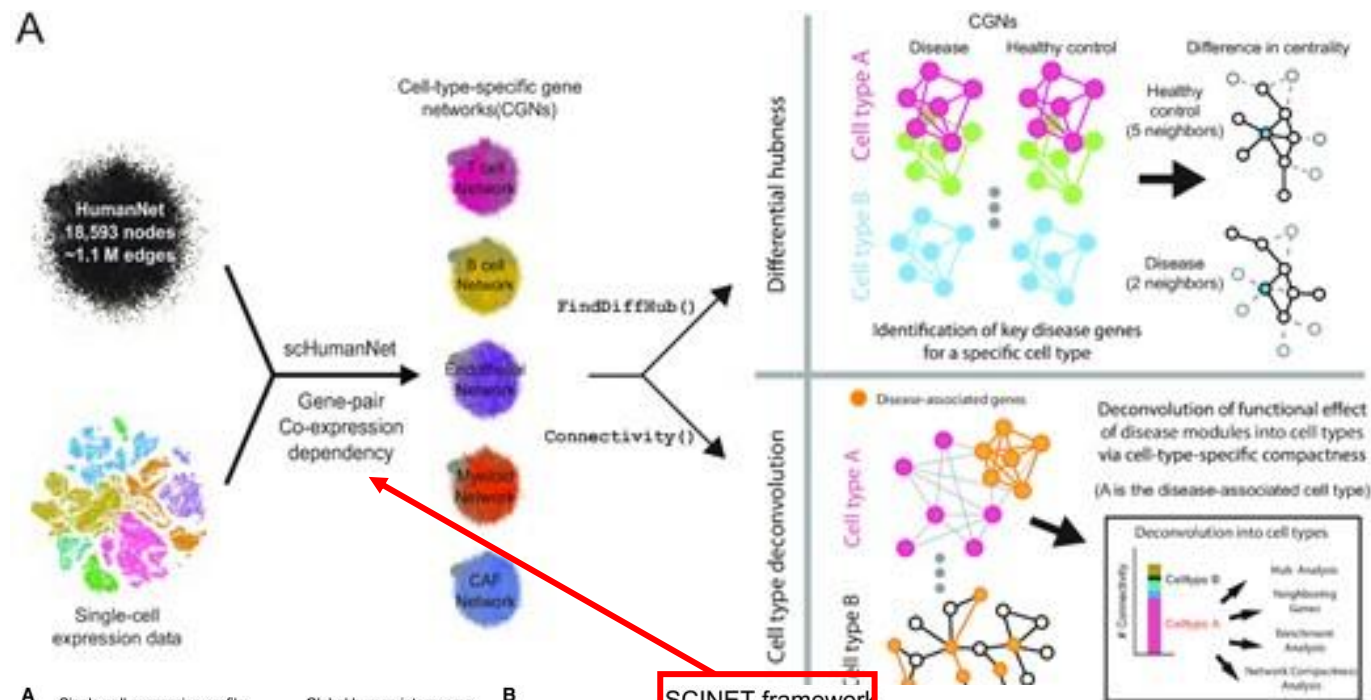-cannot account for expressing cell ratio

- Network analysis in scRNA-seq



-Skewed too much to zero-counts
→ Hard to obtain a suitable correlation

Cell-type-specific co-expression inference from single cell RNA-sequencing data

- BigScale2



-High granularity clustering: Recursive clustering (Hierarchical clustering)
-all pairwise comparison → DE → measure Z-score for each gene
→ Correlation (similar effect of cell-pooling)

# scHumanNet



-Filtering out the "cell-type-specific" network from the reference network
-HumanNetv3, String
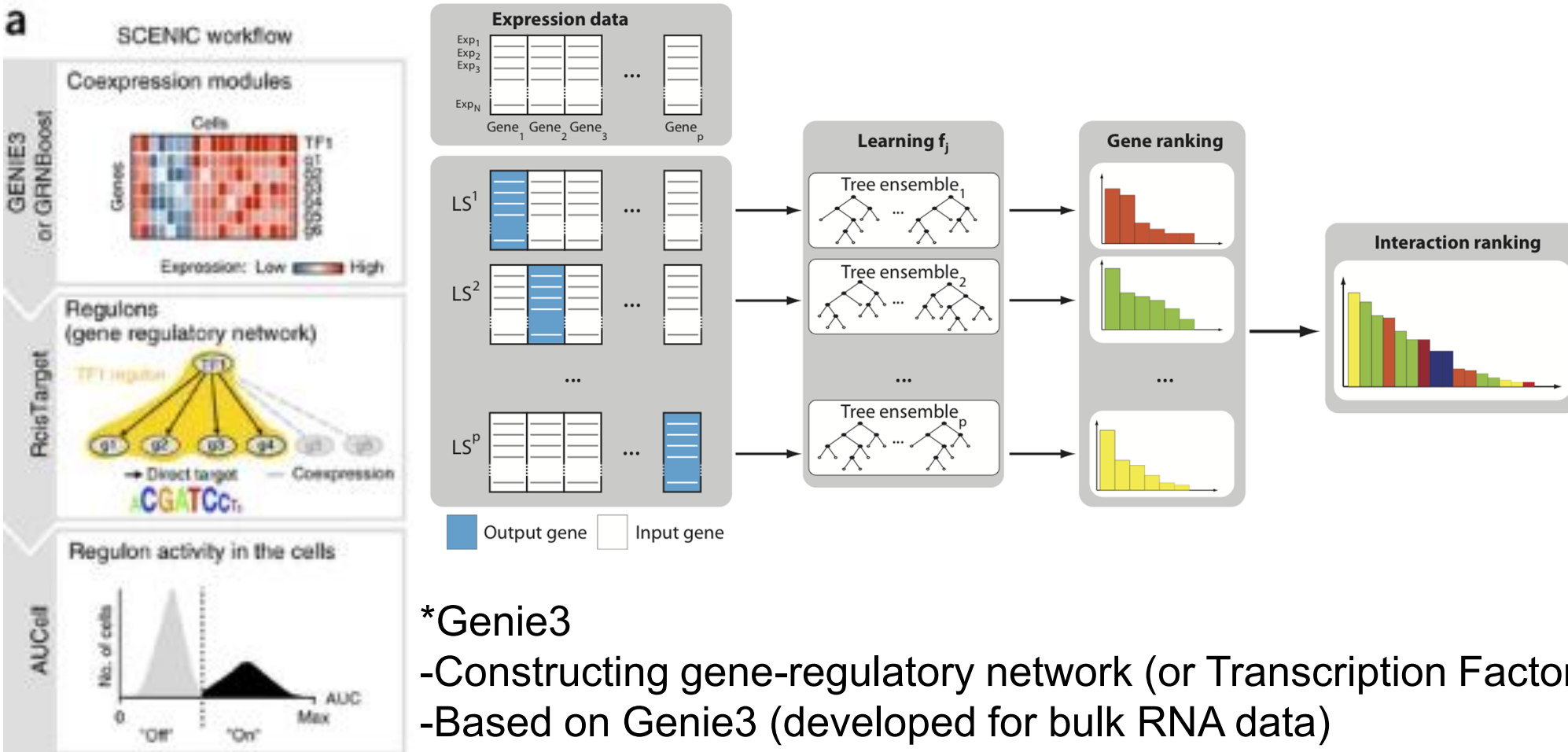
*SCINET framework
-clustering → Archetype → transcriptome interpolation (smoothing)
-gene expression transformation → Better distribution

-subsampling (per cell type)
-p-value for interacting gene-pair vs Null
-aggregate p-values by Fisher's method

Or

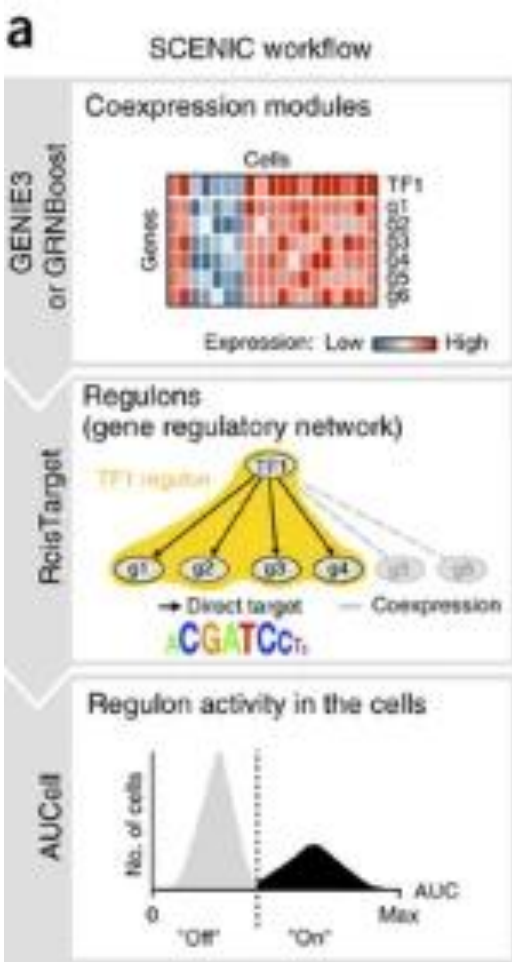Just take the edge and use the original edge score

# SCENIC (genie3)



*Genie3
-Constructing gene-regulatory network (or Transcription Factor regulatory NW: TRN)
-Based on Genie3 (developed for bulk RNA data)
-output gene exp ← explained by input genes (random forest) [coexpression pattern]
→ TF filtering
-output gene (i) ← input gene (j1, j2, j3) score → ranking (by importance)
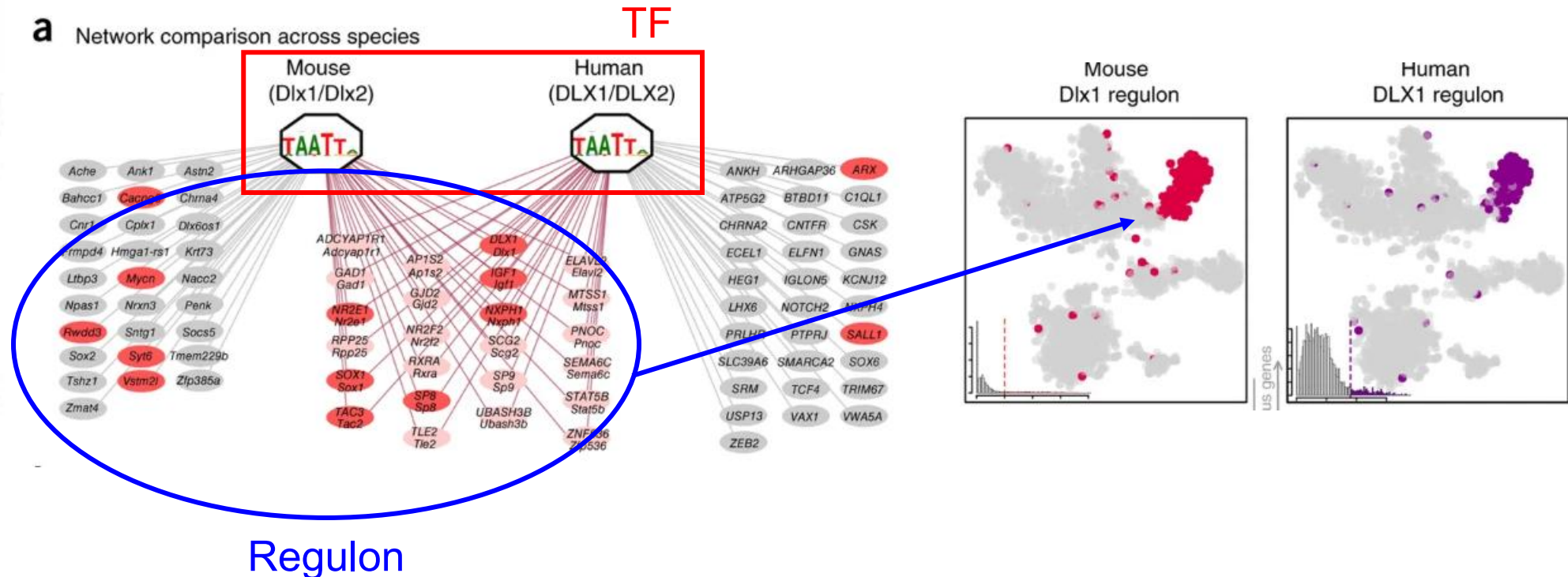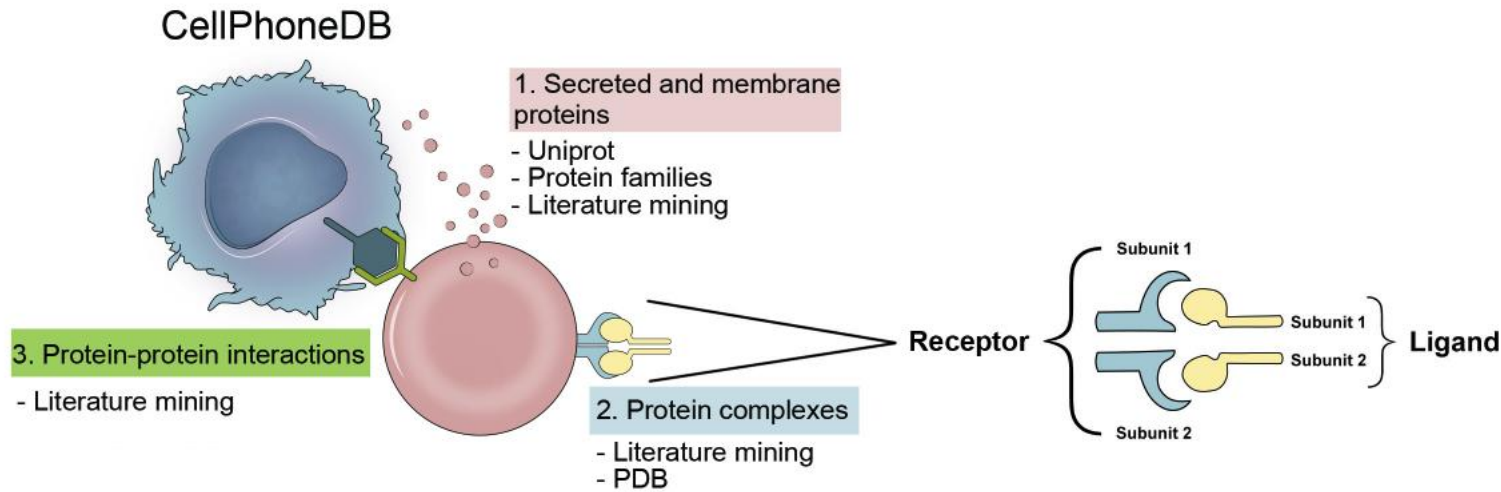
→ GRNBoost for speed

# SCENIC (genie3)



-We obtained TF- target gene (correlation or association)
→ Not all the TF binds to the target gene
→ RcisTarget: cis-regulatory motif analysis
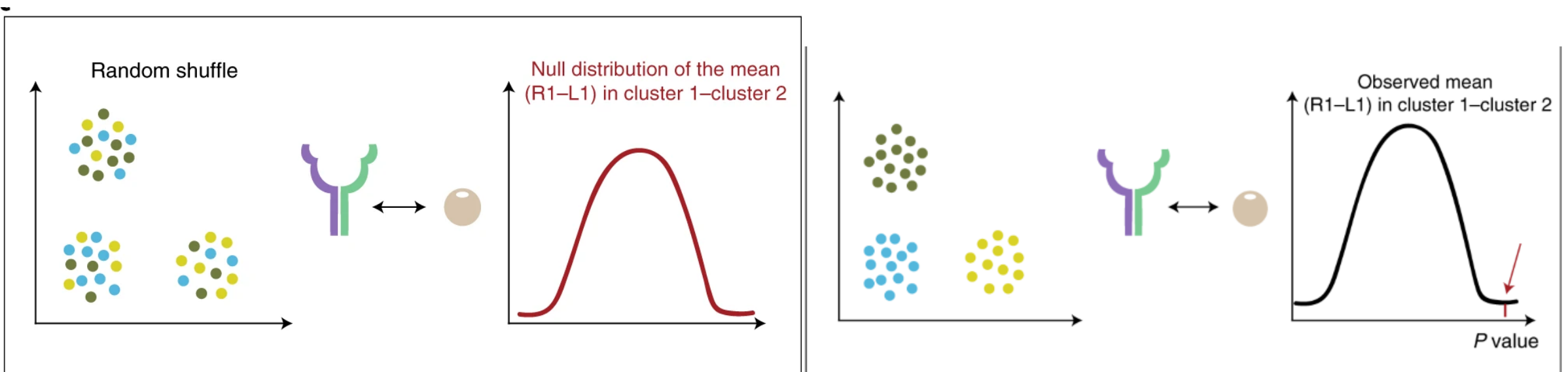→ Only enriched TF can bind to the promoter of a given genes

*Regulon activity: AUCell

TF
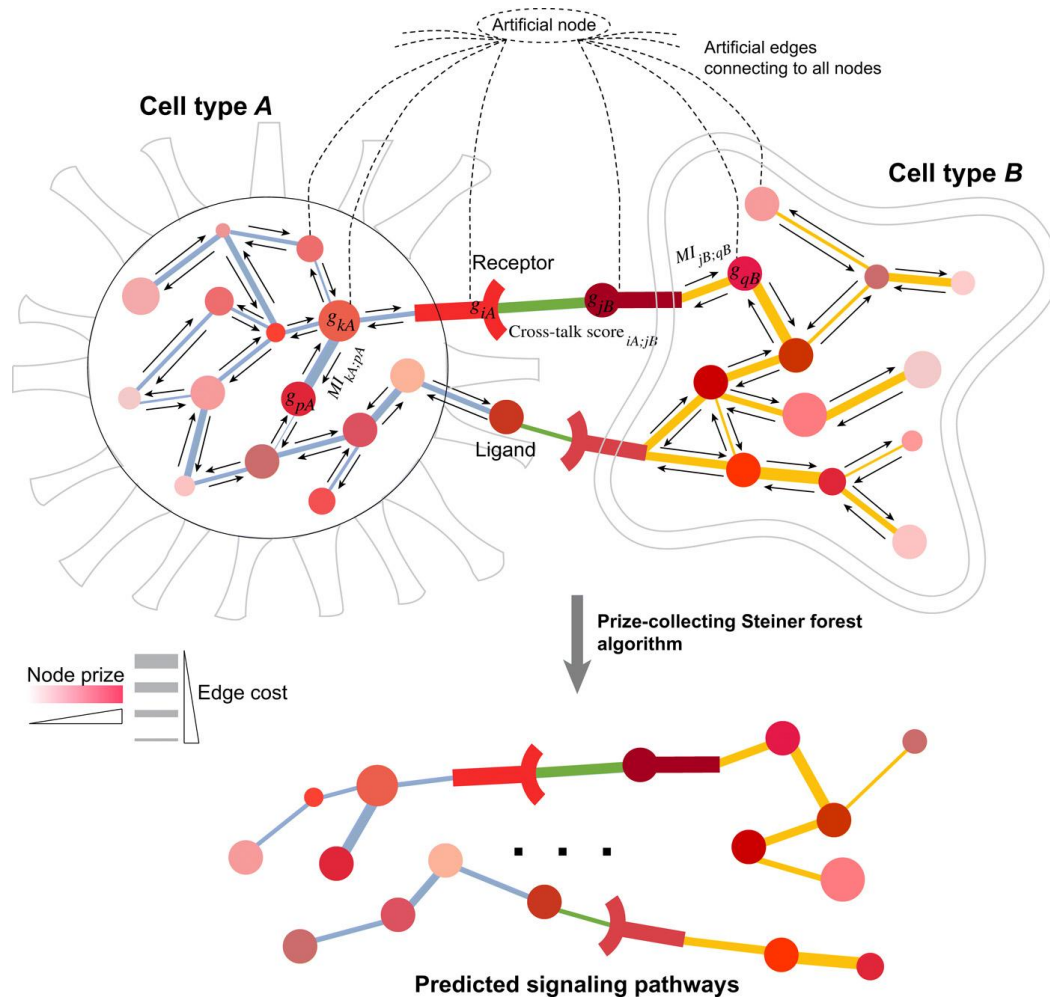
Regulon

# Cell-Cell interaction (CellPhoneDB, CellChat)

-Cluster-to-Cluster interaction
→ Measure the mean expression of receptor-ligand pair
→ Shuffle the cells: Null distribution
→ P-value measurement for each pair
→ Strength: mean expression

+ complex: mean expression



CellPhoneDB

1. Secreted and membrane proteins
- Uniprot
- Protein families
- Literature mining

3. Protein-protein interactions
- Literature mining

2. Protein complexes
- Literature mining
- PDB

Receptor { Subunit 1 / Subunit 2 } + { Subunit 1 / Subunit 2 } Ligand



Random shuffle

Null distribution of the mean (R1–L1) in cluster 1–cluster 2
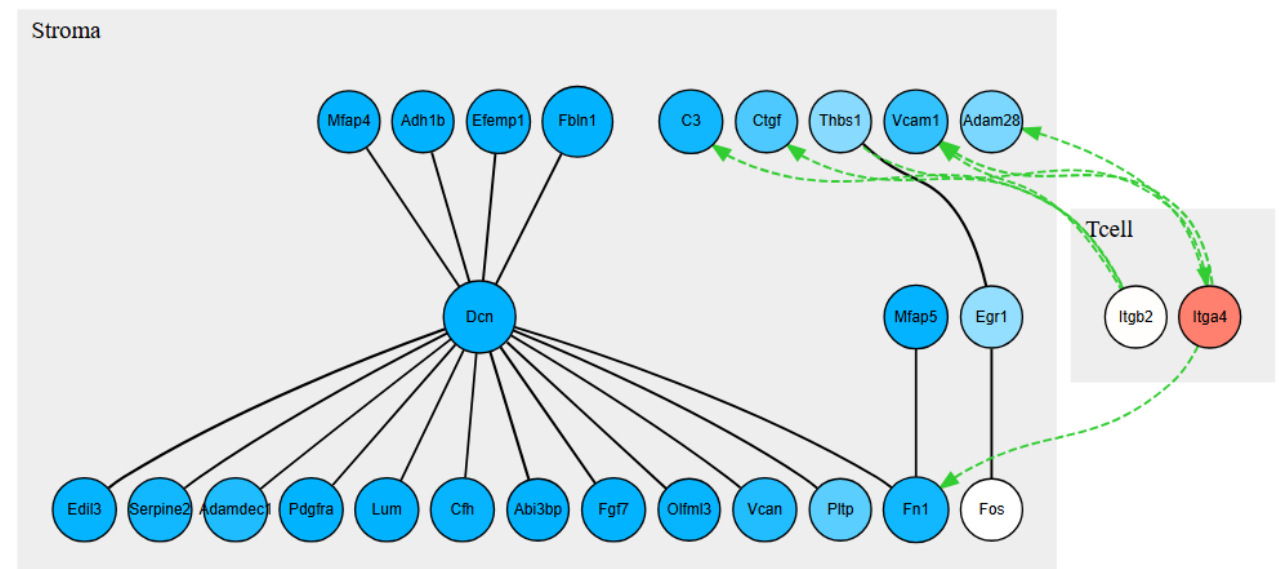
Observed mean (R1–L1) in cluster 1–cluster 2
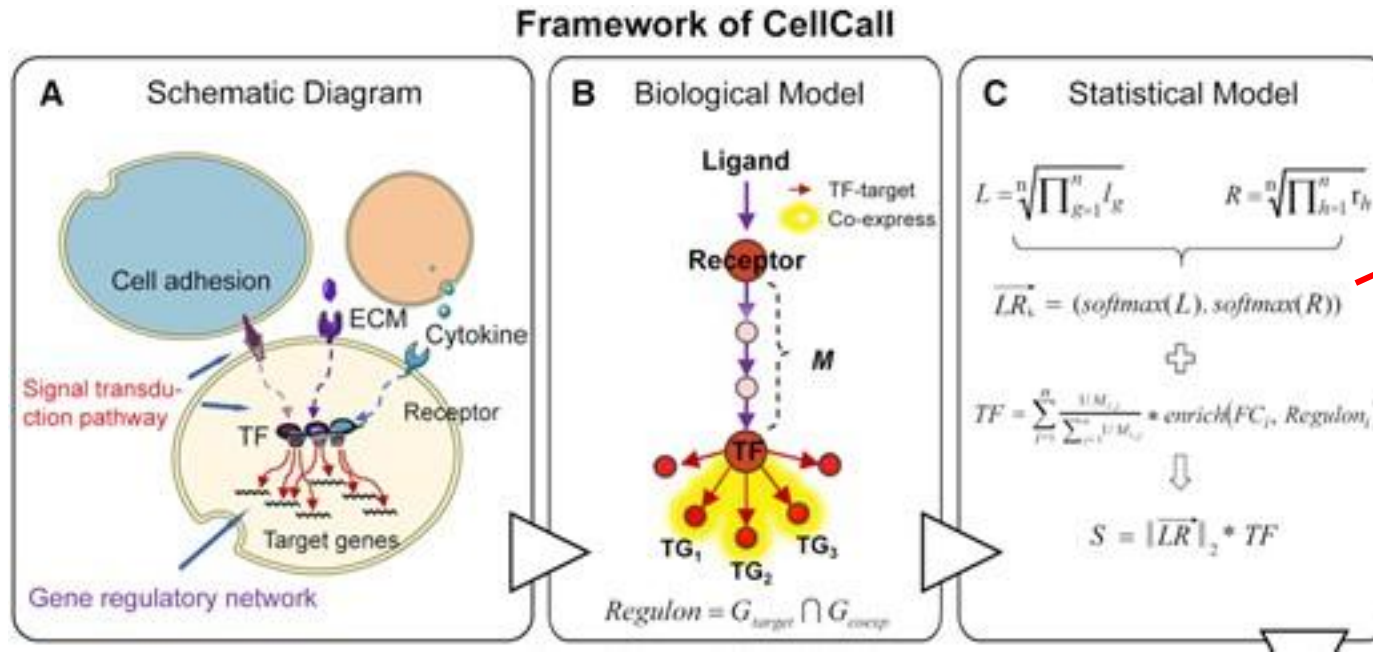
P value

# Cell-Cell interaction (CytoTalk)



What is the downstream of interaction?

-Cell-cell interaction: intercellular interaction
-Intra-cellular network: mutual information (co-occurrence of gene expression across cells)

-Network propagation algorithm
→ Remain only the significant edges

- # Cell-Cell interaction (CellCall)



Framework of CellCall

LR: conceptually mean expression

Only looking at the gene expression from the ligand-receptor is not enough
It should show some **perturbation of target genes** due to cell-cell interaction!

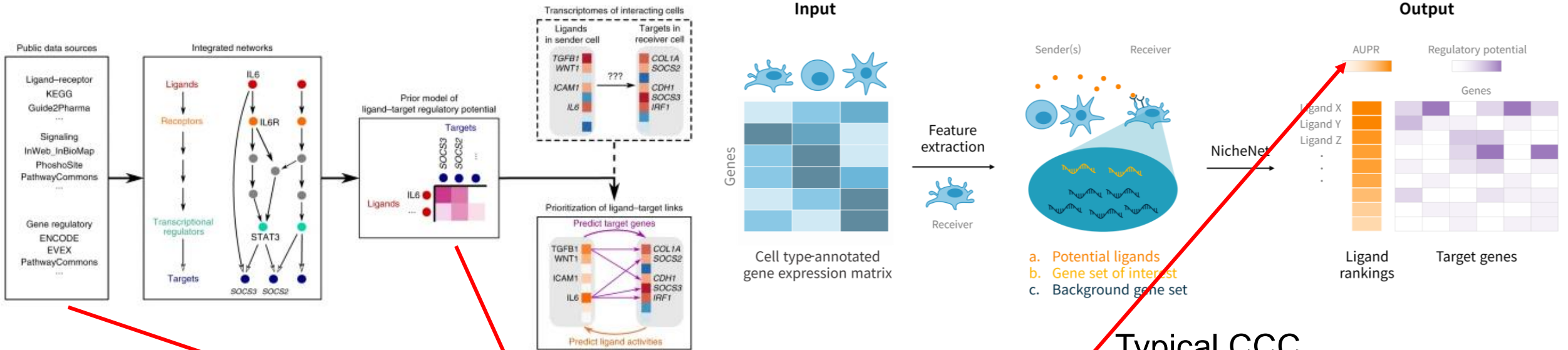Interaction score = LR * $TF_k$
$TF_k$: regulon activity of TF
-LR → TF: KEGG, etc
-TF → regulon: known DB (TRRUST …) & coexpressed with TF
-$TF_k$: GSEA for those regulon
-Multiple $TF_k$ : weight sum (number of node; LR → TF)

# Cell-Cell interaction (NicheNet)



Input

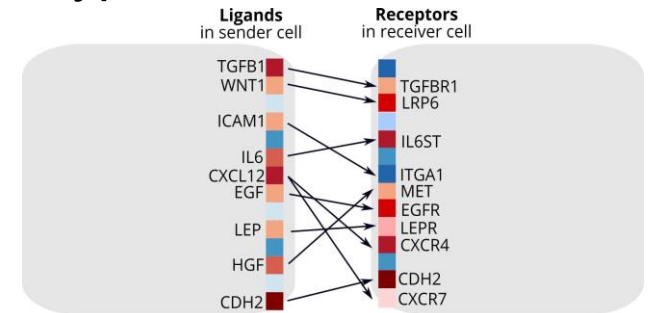Output

Typical CCC

NicheNet

-Merge each path from DB: weighted network (prior model)
Ligand – receptor – TF – target genes
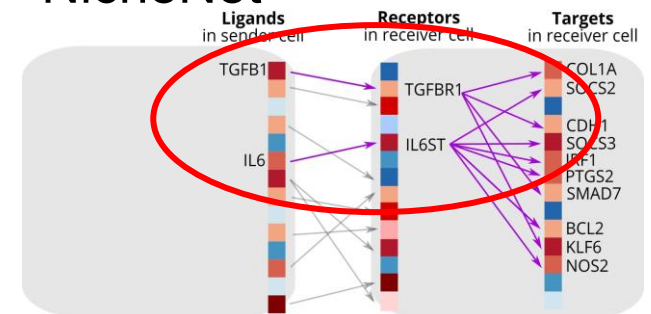
-ligand ~ potential target genes vs bg genes
-ligand ranking: ligand [exp] ~ predefined targets [gene exp]
(how much ligand expression can differentially express target genes)
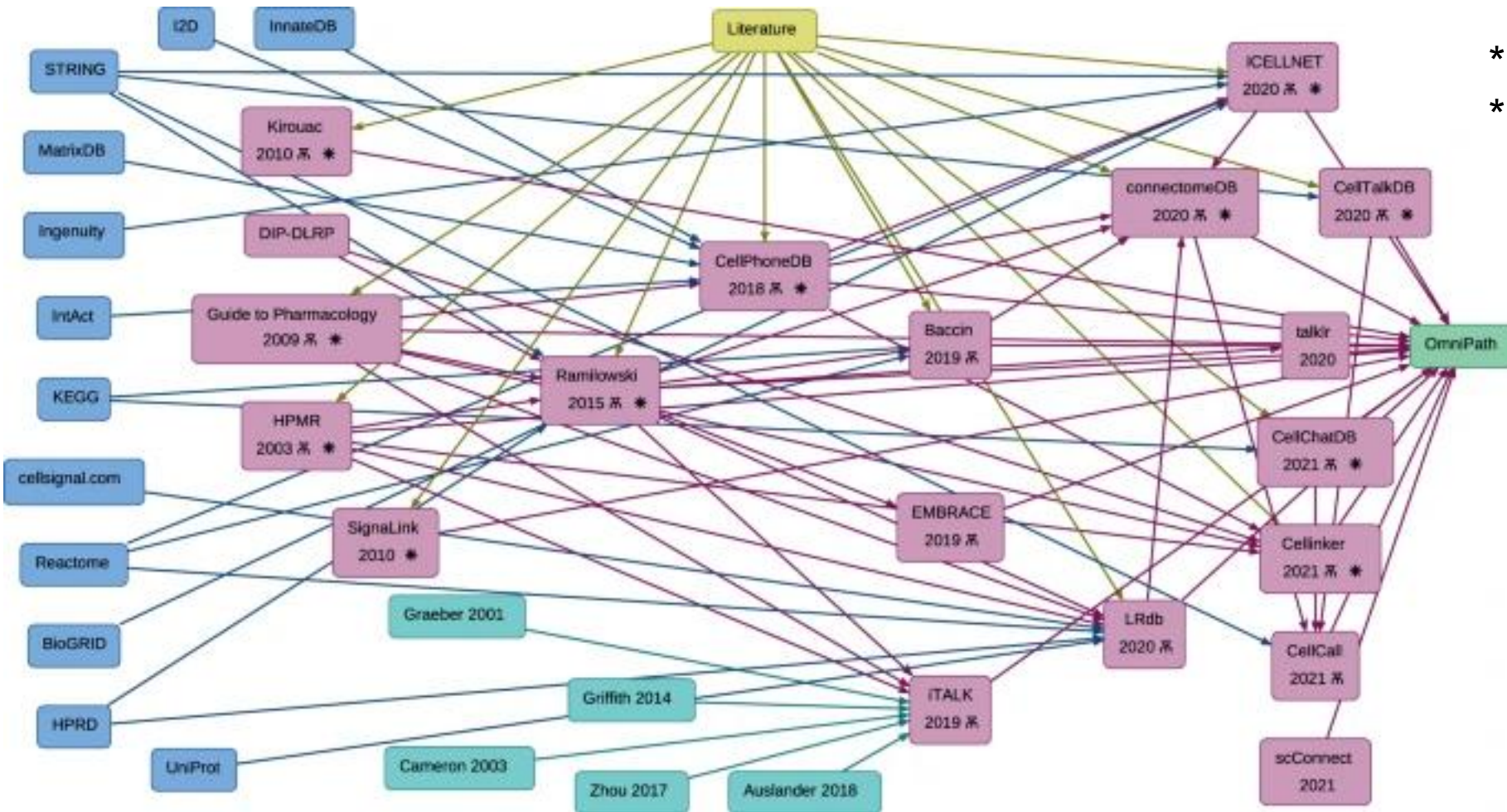-target genes are selected by a predefined ligand-target link

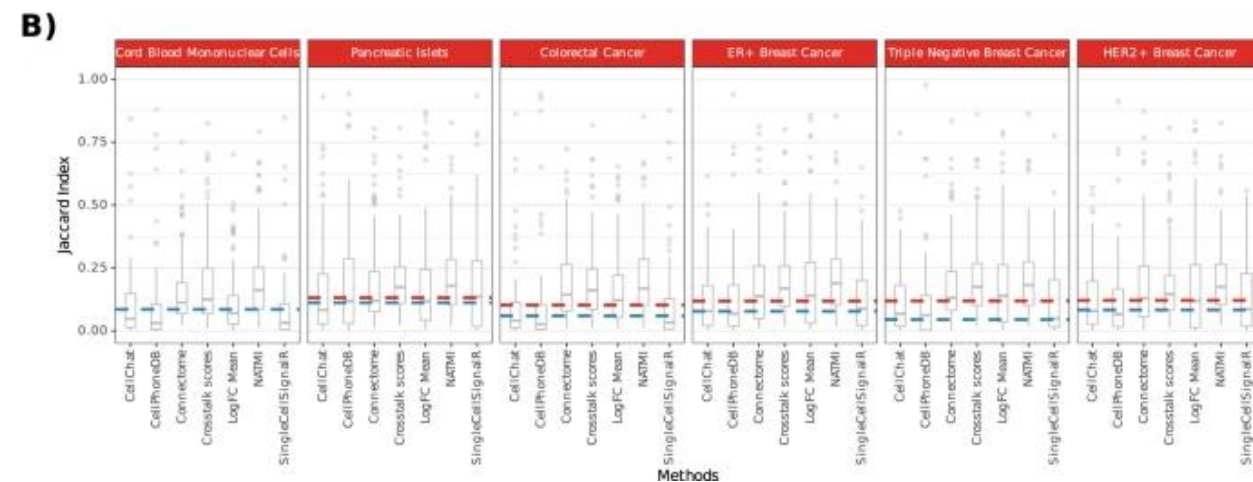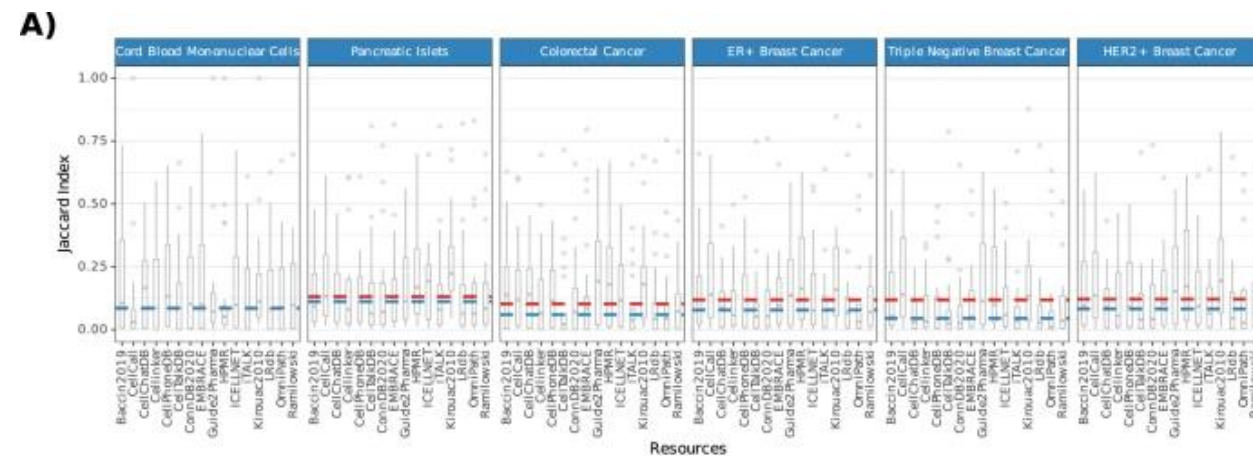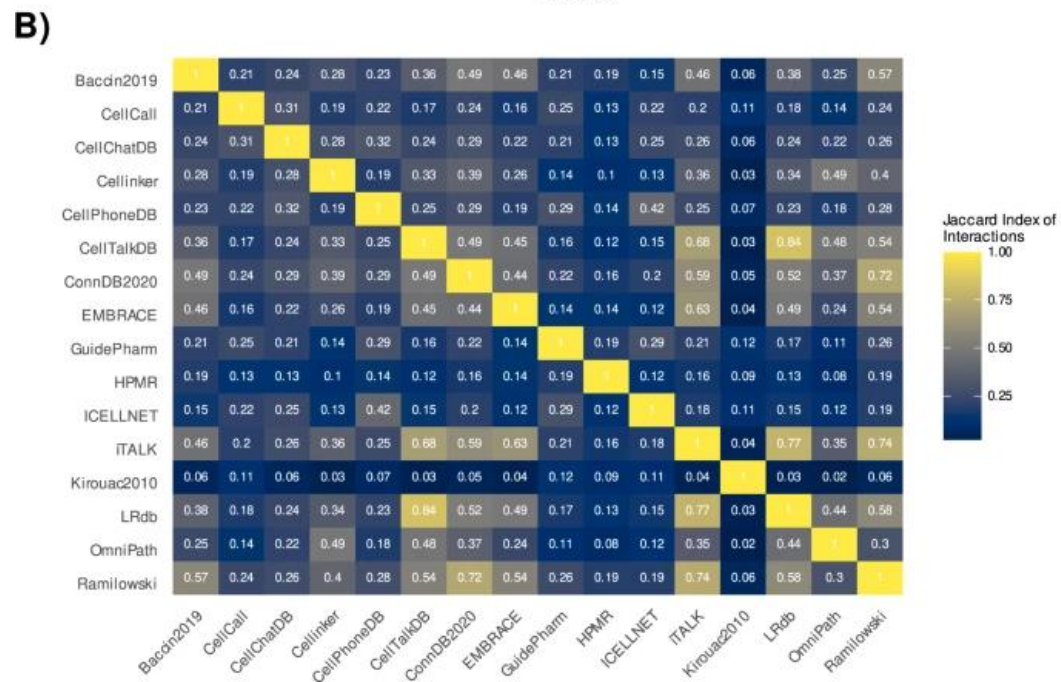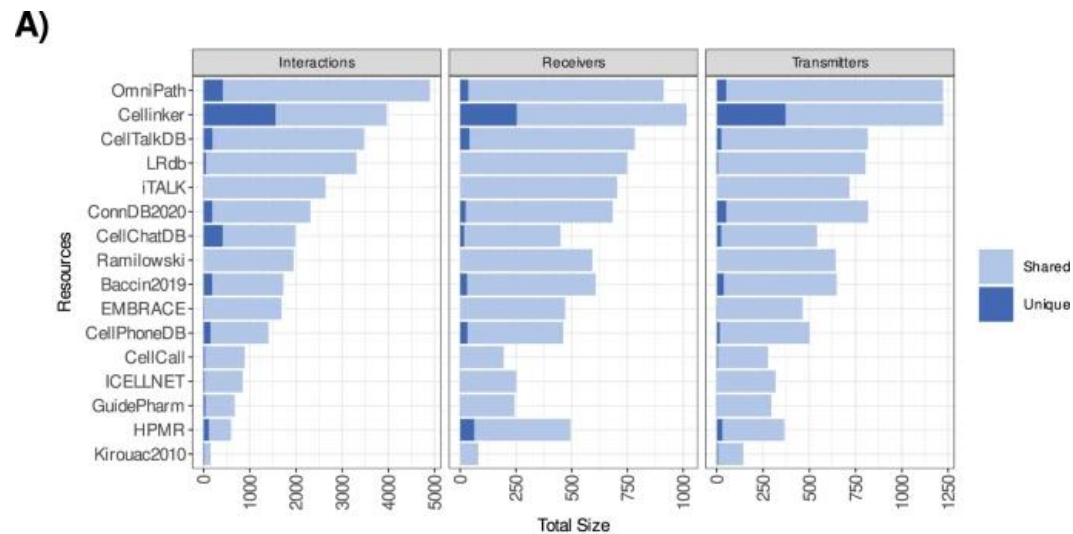- Cell-Cell interaction comparison
-Heterogeneous DB



*KEGG, Reactome, STRING
*Published literature

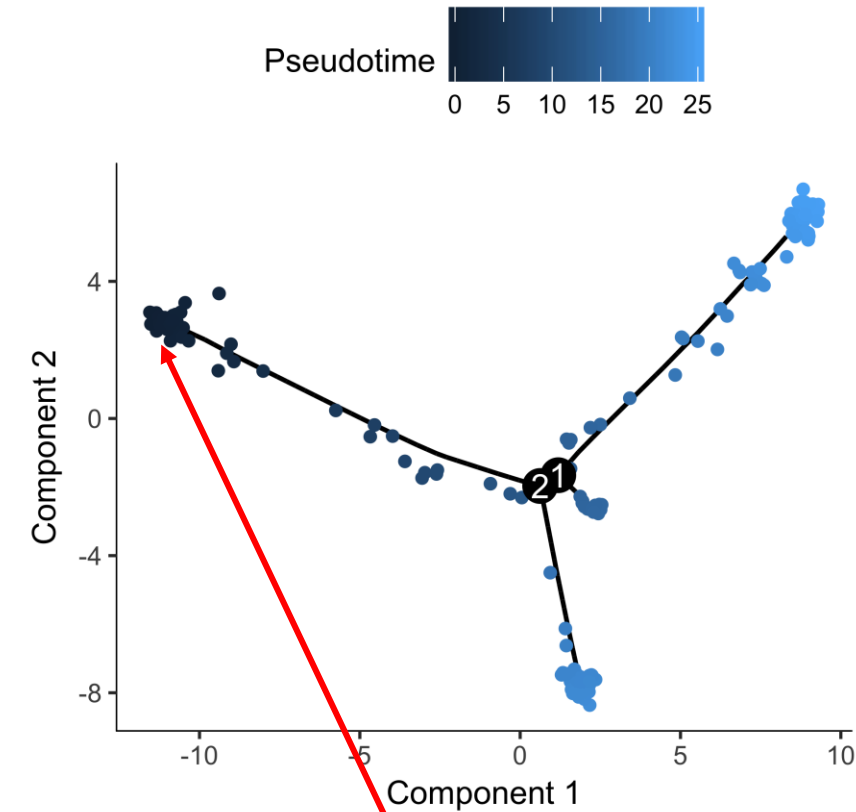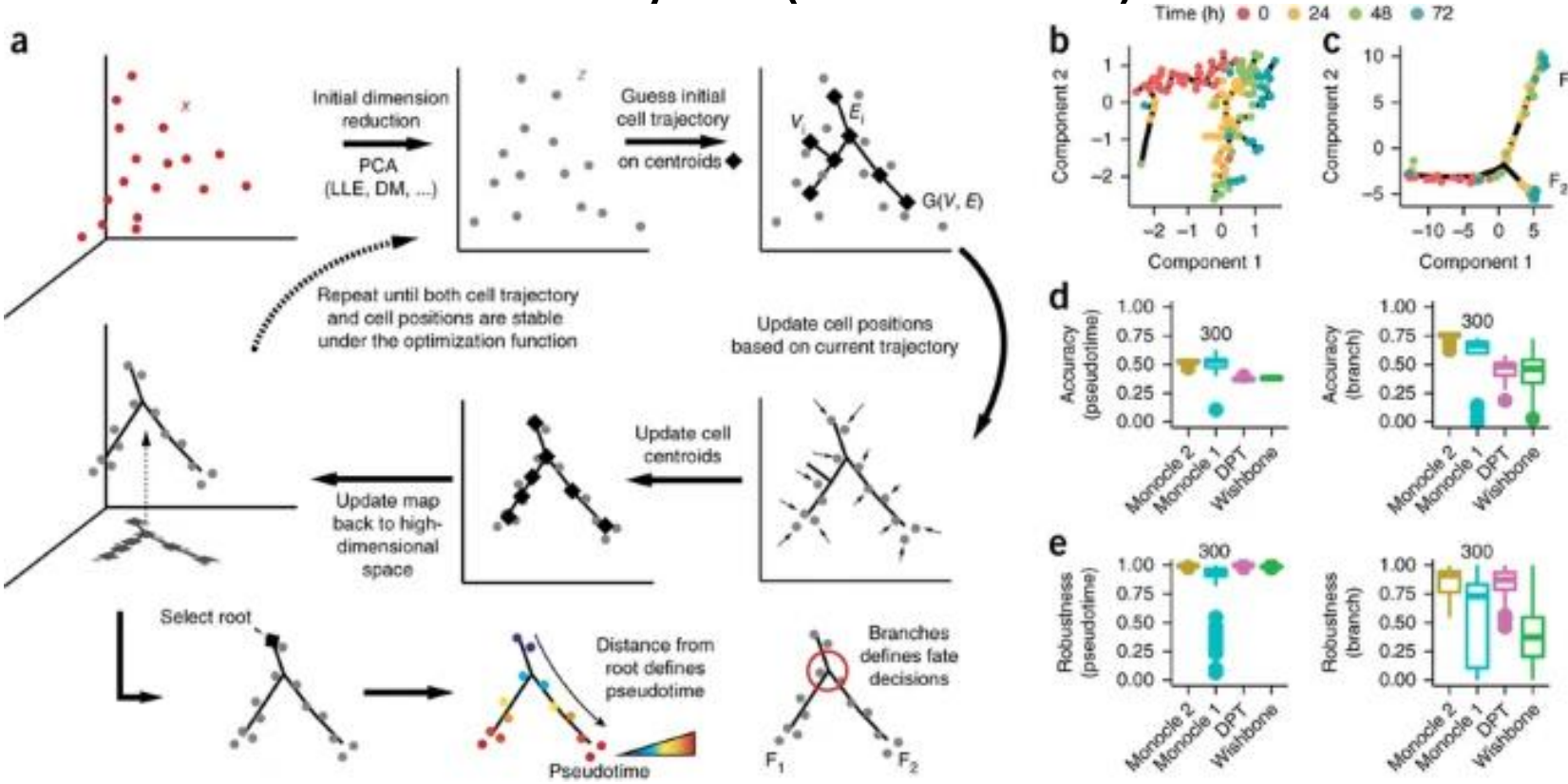- Cell-Cell interaction comparison
-Different CCC methods are too different



Low Jaccard index (low overlap)
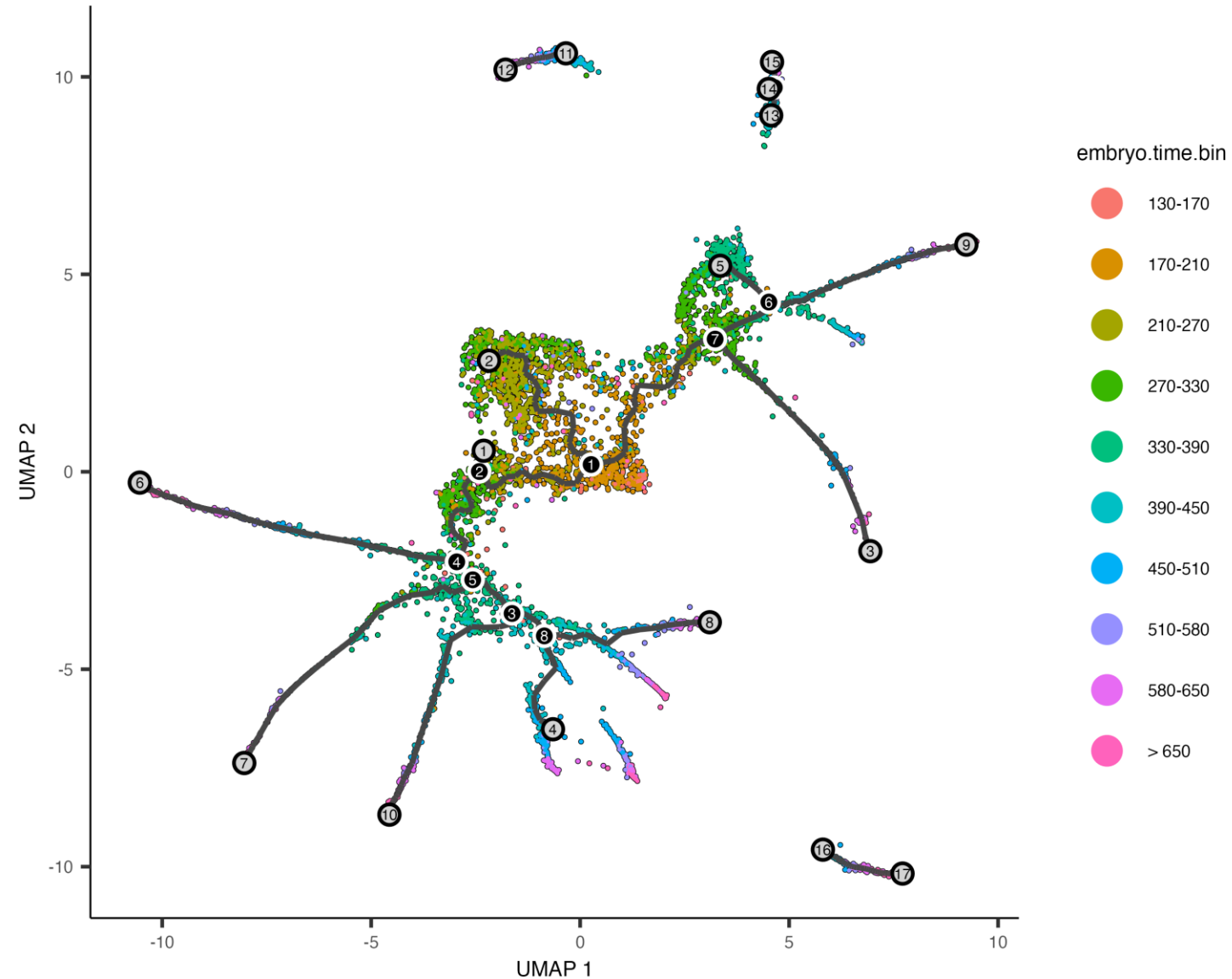-across DB
-across method

# Pseudotime analysis (Monocle2)



Aligning cells into a virtual embedding
-Dimension reduction → initial centroid (k-mean)
→ update cell position
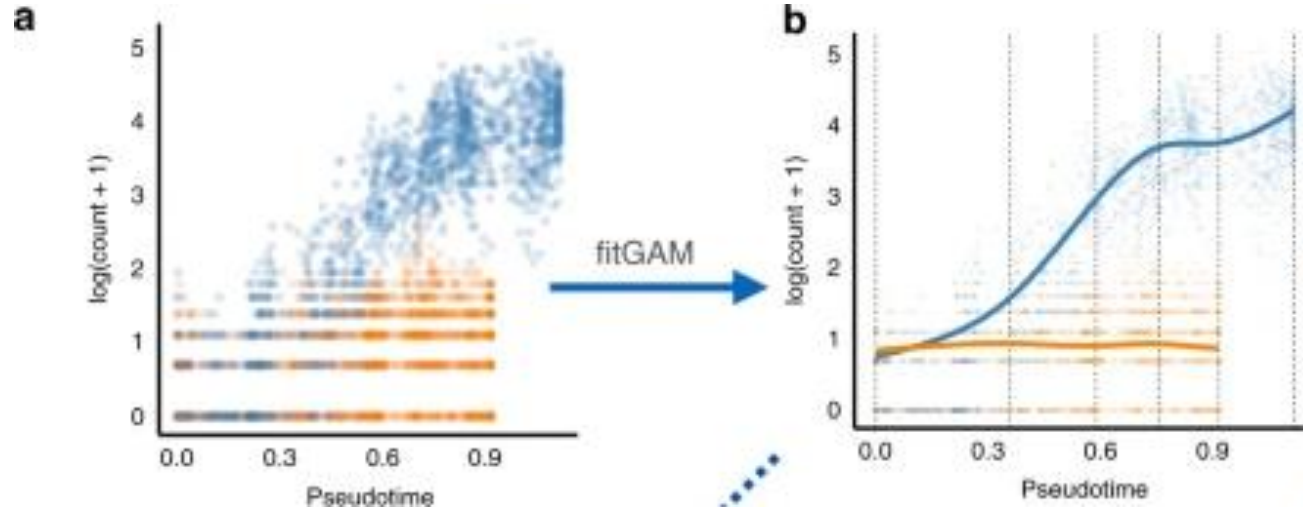→ High dimension → re-do until convergence

Conceptually: clustering + network construction
Gene: DEG (across differentiation), HVG …

*Require root cell or cluster

# Pseudotime analysis (Monocle3)



Adapt dimension reduction space into a familiar UMAP projection

- # TradeSeq


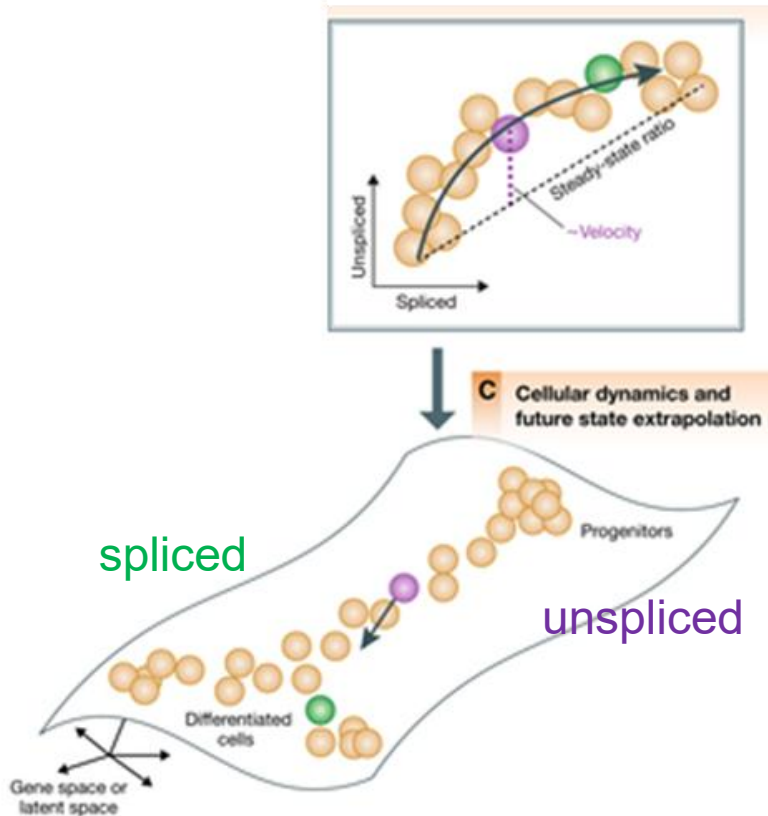
What kind of genes are associated with a given trajectory
-Pseudotime ~ Gene expression (correlation) →  likely to be poor
-Negative binomial generalized additive model (nonlinear approach)
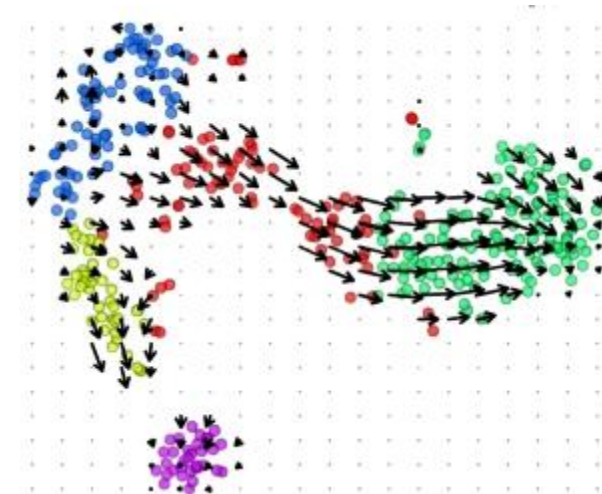
→ but, super slow …
→ Maybe, binning??

- # Velocyto, scVelo



spliced

unspliced

Assumption: during the differentiation, progenitor may have more unspliced RNA while differentiated cell may have fully spliced form

-Compare between unspliced/spliced ratio! (RNA velocity) for each gene
-Merge all the velocities (from all the gene) → project on the user embedding space (UMAP)
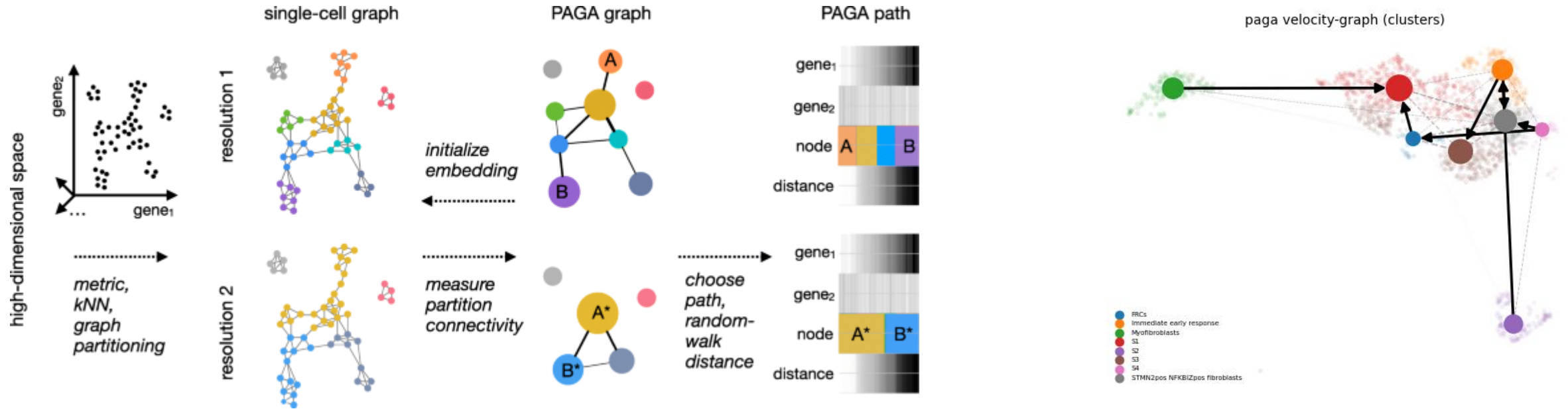-Sum of transition P = 1 for each cell

*RNA velocity: steady-state approximation
spliced RNA deg speed = splicing rate * unspliced RNA
- Spliced RNA * degradation rate = 0

$$\frac{ds}{dt} = \beta u - \gamma s$$

*scVelo: those parameters change across time and cell states

# PAGA (Representation)



-Abstraction of trajectory (scvelo) result
Using graph-based cell-cell similarity (→ connectivity calculated by trajectory)

-PAGA transition confidence score: actual / random expected model

The confidence should be interpreted as the ratio of the actual versus the expected value of connections under the null model of randomly connecting partitions.

- ImmuneDictionary