# Models, Inference & Algorithms (MIA)

**Primer: Permutation Enhances the Rigor of Genomics Data Analysis**

Jingyi Jessica Li

**Meeting: mcRigor: a statistical method to enhance the rigor of metacell partitioning in single-cell data analysis**

Pan Liu

November 5, 2025

# Genomics is a "liberal" discipline

1. Interdisciplinary nature

2. Data-driven focus

3. Rapid evolution of methods

4. Flexible analytical approaches



We currently track **1837** tools...

Submit an article   Journal homepage

≡  ◀)) Listen  ▶

Reviews

# What are the Most Important Statistical Ideas of the Past 50 Years?

Andrew Gelman ✉ 🔵 & Aki Vehtari

https://medium.com/bitgrit-data-science-publication/the-8-most-important-statistical-ideas-of-the-past-50-years-11220e46736f

# Teaser: bulk RNA-seq DE analysis

# Teaser: bulk RNA-seq DE analysis

Q: Why are many genes identified as DE genes from permuted data?

A: The negative binomial assumption does not hold on this dataset.

*MXD1* expression

Pre–therapy      On–therapy

Empirical quantiles
(edgeR fitted NB CDF)

Theoretical quantiles
(Uniform[0,1])

# Teaser: bulk RNA-seq DE analysis

**False discoveries may mislead scientific conclusions**

[Li et al., *Genome Biology*, 2022]

# How to permute data?

**Supervised learning**

X    Y

features

samples

**Unsupervised learning**

X

features

samples

# Teaser: bulk RNA-seq DE analysis

## False discoveries may mislead scientific conclusions

[Li et al., *Genome Biology*, 2022]

# Teaser: bulk RNA-seq DE analysis

**Q:** Why are many genes identified as DE genes from permuted data?

**A:** The negative binomial assumption does not hold on this dataset.



*MXD1* expression

# Teaser: bulk RNA-seq DE analysis

**False discoveries may mislead scientific conclusions**



[Li et al., *Genome Biology*, 2022]

# Teaser: bulk RNA-seq DE analysis

**Q: Why are many genes identified as DE genes from permuted data?**

**A: The negative binomial assumption does not hold on this dataset.**

*MXD1* expression

# Teaser: bulk RNA-seq DE analysis

Exaggerated false positives by popular differential expression methods when analyzing human population samples

[Li*, Ge* et al., *Genome Biology*, 2022]

X/Twitter: @jsb_ucla

# Teaser: bulk RNA-seq DE analysis
## False discoveries may mislead scientific conclusions

[Li et al., *Genome Biology*, 2022]

# How to permute data?

**Supervised learning**

X          Y

features

samples

**Bulk RNA-seq:**

features = genes

Y = sample condition labels

**Unsupervised learning**

X

features

samples

**Single-cell RNA-seq:**

samples = cells;

features = genes

# Two examples where permutation helps

## 1. Single-cell data **visualization**

> **Statistical method scDEED for detecting dubious 2D single-cell embeddings and optimizing t-SNE and UMAP hyperparameters**
>
> Lucy Xia, Christy Lee & Jingyi Jessica Li ✉
>
> *Nature Communications* **15**, Article number: 1753 (2024) | Cite this article

## 2. Aggregating single cells into **metacells**

> **mcRigor: a statistical method to enhance the rigor of metacell partitioning in single-cell data analysis**
>
> Pan Liu & Jingyi Jessica Li ✉
>
> *Nature Communications* **16**, Article number: 8602 (2025) | Cite this article

# Example 1: dubious t-SNE/UMAP embeddings?

## How to Use t-SNE Effectively



- **Hyperparameters** really matter

- **Distances between clusters** might not mean anything

- ...

Source: https://distill.pub/2016/misread-tsne/

# Example 1: dubious t-SNE/UMAP embeddings?

## Seeing data as t-SNE and UMAP do

Vivien Marx ✉

Dimension reduction helps to visualize high-dimensional datasets.
These tools should be used thoughtfully and with tuned parameters.
Sometimes, these methods take a second thought.

# Example 1: dubious t-SNE/UMAP embeddings?

**Permuted cells are exchangeable → A cell's neighbors are random**

permuted data

each cell

vs.

null distribution of reliability scores

density

dubious    reliability score    trustworthy

reliability score = cor( )

all *n* cells

# Example 1: dubious t-SNE/UMAP embeddings?

**scDEED** detects dubious embeddings



Original (perplexity 40)

*Hydra* single-cell RNA-seq data [Siebert et al., *Science*, 2019]

# Example 1: dubious t-SNE/UMAP embeddings?

**scDEED** optimizes hyperparameters by minimizing dubious embeddings

Original (perplexity 40)

scDEED optimized (perplexity 230)

- neuron ec1
- neuron ec3
- ecEP_sc

scaled gene expression

cells

genes

# Example 1: dubious t-SNE/UMAP embeddings?
## scDEED enhances the consistency between t-SNE and UMAP

# Two examples where permutation helps

1. Single-cell data **visualization**

   ## Statistical method scDEED for detecting dubious 2D single-cell embeddings and optimizing t-SNE and UMAP hyperparameters

   Lucy Xia, Christy Lee & Jingyi Jessica Li ✉

   *Nature Communications* **15**, Article number: 1753 (2024) | Cite this article

2. Aggregating single cells into **metacells**

   ## mcRigor: a statistical method to enhance the rigor of metacell partitioning in single-cell data analysis

   Pan Liu & Jingyi Jessica Li ✉

   *Nature Communications* **16**, Article number: 8602 (2025) | Cite this article

# Example 2: aggregating single cells into metacells

Metacell: a heuristic solution to the sparsity issue in single-cell data



Bilous, M., et al. "Building and analyzing metacells in single-cell genomics data." Molecular Systems Biology (2024): 1-23.

# Example 1: dubious t-SNE/UMAP embeddings?

**scDEED** enhances the consistency between t-SNE and UMAP

1. Single-cell data visualization

Statistical method scDEED for detecting dubious 2D single-cell embeddings and optimizing t-SNE and UMAP hyperparameters

Lucy Xia, Christy Lee & Jingyi Jessica Li

*Nature Communications* 15, Article number: 1753 (2024) | Cite this article

2. Aggregating single cells into metacells

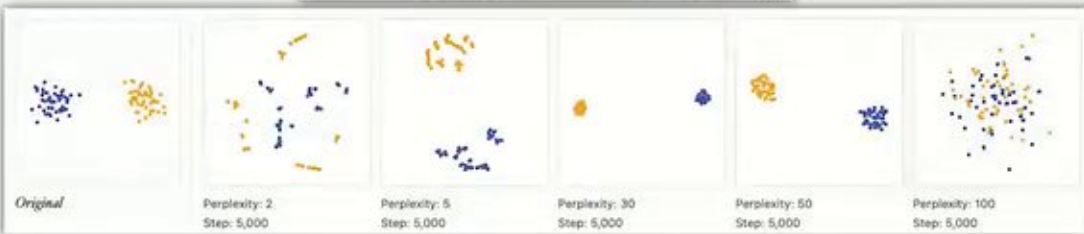mcRigor: a statistical method to enhance the rigor of metacell partitioning in single-cell data analysis

*Nature Communications* 16, Article number: 8602 (2025) | Cite this article

# Examples where permutation helps

## Two examples where scDEED/UMAP embeddings?
**scDEED** enhances the consistency between t-SNE and UMAP

1. Single-cell data **visualization**



Statistical method scDEED for detecting dubious 2D single-cell embeddings and optimizing t-SNE and UMAP hyperparameters

Lucy Xia, Christy Lee & Jingyi Jessica Li

*Nature Communications* **15**, Article number: 1753 (2024) | Cite this article

2. Aggregating single cells into **metacells**

mcRigor: a statistical method to enhance the rigor of metacell partitioning in single-cell data analysis

Pan Liu & Jingyi Jessica Li

*Nature Communications* **16**, Article number: 8602 (2025) | Cite this article

B cell    CD14 + monocyte    CD4+ T cell    Cytotoxic T cell

# Example 2: aggregating single cells into metacells

## Metacell: a heuristic solution to the sparsity issue in single-cell data



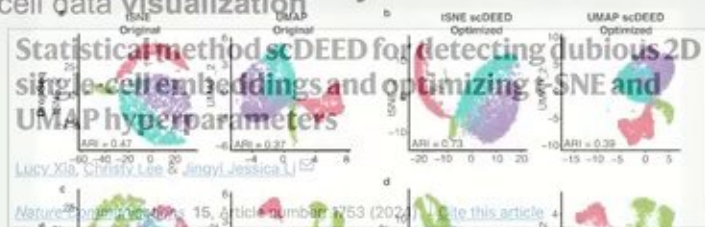Bilous, M., et al. "Building and analyzing metacells in single-cell genomics data." Molecular Systems Biology (2024): 1-23.

# Example 2: aggregating single cells into metacells

## A statistical definition of "metacell"

"A *homogeneous* collection of single-cell profiles that could have been resampled from the *same* original cell."

$\implies$ Variation within a metacell is attributed exclusively to measurement error

**Two-layer observation model:**
Cell (observation) $i = 1, \ldots, n$
Feature $j = 1, \ldots, p$

**Expression model:** $\lambda_i \sim \mathcal{F}(\cdot | \mathbf{x}_i)$

**Measurement model:** $y_{ij} \sim \mathcal{G}(y_{i+} \lambda_{ij})$

$\Downarrow$

**Statistical definition:** A **metacell** is a group of single cells that share the same $\lambda$

$\Downarrow$

Satisfying this definition?

Yes: **trustworthy metacells**        No: **dubious metacells**

A statistical problem

# Example 2: aggregating single cells into metacells

**Our proposal: mcRigor**

**Goals:** a statistical criterion to
- Identify dubious metacells consisting of single cells from different cell states
- Nominate the top-performing metacell method and optimize its hyperparameter

$$\text{granularity level } \gamma = \frac{\#\text{single cells}}{\#\text{metacells}}$$

in a data-specific way

corr

trustworthy    dubious

Pan Liu
(JSB)

# Example 2: aggregating single cells into metacells

Q: Is within-gene permutation enough?

genes

genes

cells

permuted cells

within-gene permutation

# Example 2: aggregating single cells into metacells

## Our proposal: mcRigor

**Goals:** a statistical criterion to
- Identify dubious metacells consisting of single cells from different cell states
- Nominate the top-performing metacell method and optimize its hyperparameter

$$\text{granularity level } \gamma = \frac{\#\text{single cells}}{\#\text{metacells}}$$

in a data-specific way

**Intuition:**
- Within a trustworthy metacell, features are approximately uncorrelated

Pan Liu
(JSB)

# Example 2: aggregating single cells into metacells

## A statistical definition of "metacell"

"A homogeneous collection of single-cell profiles that could have been resampled from the same original cell."

$\implies$ Variation within a metacell is attributed exclusively to measurement error

**Two-layer observation model:**

Cell (observation) $i = 1, \ldots, n$

Feature $j = 1, \ldots, p$

**Expression model:** $\lambda_i \sim \mathcal{F}(\cdot | \mathbf{x}_i)$

**Measurement model:** $y_{ij} \sim \mathcal{G}(y_{i+}, \lambda_{ij})$

$\Downarrow$

**Statistical definition:** A **metacell** is a group of single cells that share the same $\lambda$

$\Downarrow$

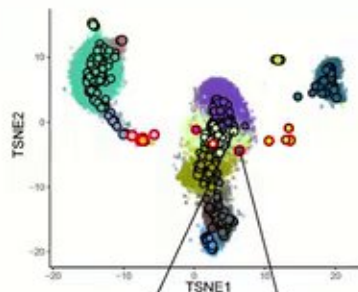Satisfying this definition?

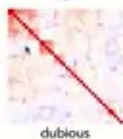Yes: **trustworthy metacells**        No: **dubious metacells**

A statistical problem

# Example 2: aggregating single cells into metacells

**Our proposal: mcRigor**

**Goals:** a statistical criterion to

- Identify dubious metacells consisting of single cells from different cell states
- Nominate the top-performing metacell method and optimize its hyperparameter

$$\text{granularity level } \gamma = \frac{\#\text{single cells}}{\#\text{metacells}}$$
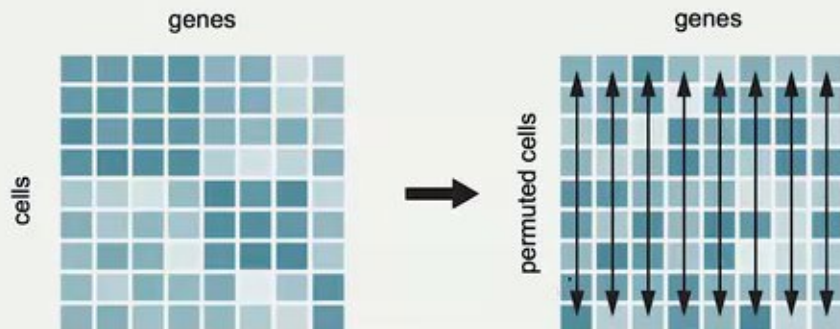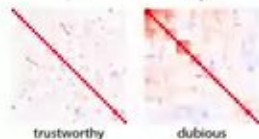
in a data-specific way

**Intuition:**

- Within a trustworthy metacell, features are approximately uncorrelated

Pan Liu
(JSB)

# Example 2: aggregating single cells into metacells

**Q:** Is within-gene permutation enough?

**A:** Genes become uncorrelated, but cell library sizes are gone.



within-gene permutation

# Example 2: aggregating single cells into metacells

**Double permutation**

Within-gene permutation:

- **preserves genes marginal distributions**
- **removes gene correlations**
- **removes cell library sizes**

genes

cells

mcDiv

within-gene permutation →

genes

permuted cells

mcDiv^null

X

Pan Liu
(JSB)

# Example 2: aggregating single cells into metacells

## mcRigor function 1: detecting dubious metacells

# Example 2: aggregating single cells into metacells

## mcRigor function 1: detecting dubious metacells

# Example 2: aggregating single cells into metacells

## mcRigor function 1: detecting dubious metacells

Summary

More generally,

Synthetic null data

# Single-cell post-clustering DE analysis

Q: How to control false discoveries using synthetic null data?
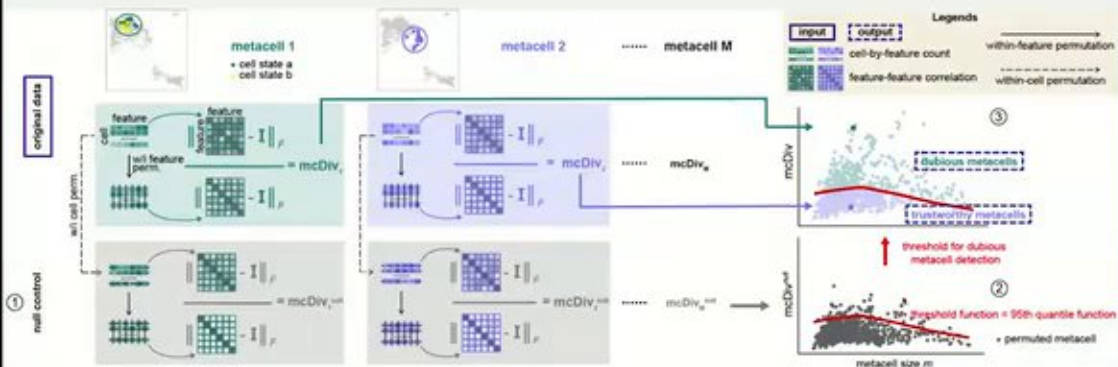
# Single-cell post-clustering DE analysis

**Expectation 1:** Cell-type marker genes should be found as top DE genes.

**Expectation 2:** Housekeeping genes should NOT be found as top DE genes.

— CD14⁺/CD16⁺ Monocyte Markers — Housekeeping Genes

# Single-cell post-clustering DE analysis

Q: Why does **ClusterDE** NOT identify housekeeping genes as top DE genes?

A: **ClusterDE** uses contrast scores (= target DE score – null DE score).

# Single-cell post-clustering DE analysis

ClusterDE guides the merging of spurious clusters.



Drosophila
visual system
developmental
atlas

Ozel et al.
*Nature* (2021)

# Single-cell post-clustering DE analysis

## ClusterDE identifies cell-type markers in a cell cluster hierarchy.

# Spatial post-clustering DE analysis

Preprocessed data → Spatial clustering → Enhanced clustering → Enhanced gene expression → Differential expression

Source: https://www.nature.com/articles/s41587-021-00935-2

# Spatial post-clustering DE analysis

**(1) Synthetic null generation**

**(2) Clustering**

**(3) DE analysis Clusters 1 vs. 2**

**(4) FDR control**
Contrast scores

$c_j = s_j - \tilde{s}_j$

Target FDR (e.g., 0.05)

DE genes

scDesign3 — Spatial analysis pipeline — Clipper

**Spatial post-clustering DE analysis**

(1) Synthetic null generation — scDesign3

(2) Clustering — Spatial analysis pipeline

(3) DE analysis Clusters 1 vs. 2

Target DE scores: $S_1$ $S_2$ $S_3$ ...... $S_m$

Null DE scores: $\tilde{S}_1$ $\tilde{S}_2$ $\tilde{S}_3$ ...... $\tilde{S}_m$

(4) FDR control — Clipper

Contrast scores: $C_1$ $C_2$ $C_3$ ...... $C_m$

$C_j = S_j - \tilde{S}_j$

Target FDR (e.g., 0.05)

DE genes

# Spatial post-clustering DE analysis

**(1) Synthetic null generation**

**(2) Clustering**

**(3) DE analysis Clusters 1 vs. 2**

**(4) FDR control**

# Spatial post-clustering DE analysis



(a) Manual annotation (sample 151673) | Target data (spatial domain detected by BayesSpace) | Synthetic null data (spatial domain detected by BayesSpace)

Domains
- L5
- L6
- WM

(b) The top genes identified by ClusterDE (L5 vs L6) | The top genes identified by Seurat (L5 vs L6)

The top genes identified by ClusterDE (L6 vs WM) | The top genes identified by Seurat (L6 vs WM)

expression
1.00
0.75
0.50
0.25
0.00

# More Generally: **Nullstrap**

- A general framework for statistical inference that:

- Constructs **synthetic null data** using a model trained under the null

- Applies **the same selection and inference procedure** to both the observed and synthetic null datasets

- Calibrates the results from the two datasets to control the **FDR**

- Avoids data alteration and preserves **power**

arXiv > stat > arXiv:2501.05012

Statistics > Methodology

[Submitted on 9 Jan 2025 (v1), last revised 15 Jul 2025 (this version, v2)]

**Nullstrap: A Simple, High-Power, and Fast Framework for FDR Control in Variable Selection for Diverse High-Dimensional Models**

Changhu Wang, Ziheng Zhang, Jingyi Jessica Li

# Real data application: the triple-omic dataset ($n = 150$, $p = 6331$)

# An independent benchmark: performance comparison result

# An independent benchmark: runtime comparison result

| Method | $p = 1000$ | $p = 3000$ | $p = 10\,000$ |
| --- | --- | --- | --- |
| Mutual Information | $1.96 \pm 0.02$ | $5.87 \pm 0.05$ | $19.52 \pm 0.19$ |
| Random Forest | $41.03 \pm 0.77$ | $122.16 \pm 1.90$ | $406.27 \pm 7.78$ |
| Stabl | $92.60 \pm 2.13$ | $348.66 \pm 29.11$ | $1644.52 \pm 151.70$ |
| Nullstrap | $1.17 \pm 0.16$ | $3.08 \pm 0.40$ | $8.50 \pm 0.75$ |

https://doi.org/10.1101/2025.09.09.675248

# Metacell: a heuristic solution to the sparsity issue in single-cell data

Bilous, M., et al. "Building and analyzing metacells in single-cell genomics data." *Molecular Systems Biology* (2024): 1-23.

# Metacell: a heuristic solution to the sparsity issue in single-cell data



Bilous, M., et al. "Building and analyzing metacells in single-cell genomics data." Molecular Systems Biology (2024): 1-23.

Metacell: a heuristic solution to the sparsity issue in single-cell data

Bilous, M., et al. "Building and analyzing metacells in single-cell genomics data." Molecular Systems Biology (2024): 1-23.

# Metacell methods and applications

**Metacell Methods**

- MetaCell: analysis of single-cell RNA-seq data using K-nn graph partitions — MetaCell
- Metacell-2: a divide-and-conquer metacell algorithm for scalable scRNA-seq analysis — MetaCell-2
- Metacells untangle large and complex single-cell transcriptome networks — SuperCell
- SEACells infers transcriptional and epigenomic cellular states from single-cell genomics data — SEACells
- MetaQ: fast, scalable and accurate metacell inference via single-cell quantization — MetaQ

**Metacell Applications**

- Single-cell characterization of haematopoietic progenitors and their trajectories in homeostasis and perturbed haematopoiesis
- Single-cell mapping of the thymic stroma identifies IL-25-producing tuft epithelial cells
- NASH limits anti-tumour surveillance in immunotherapy-treated HCC
- Temporal single-cell tracing reveals clonal revival and expansion of precursor exhausted T cells during anti-PD-1 therapy in lung cancer

## Questions we wish to answer

Q: How to define a "metacell"?

Q: How to detect dubious metacells?

Q: How to optimize metacell partitioning?

Article | Open access | Published: 29 September 2025

## mcRigor: a statistical method to enhance the rigor of metacell partitioning in single-cell data analysis

Pan Liu & Jingyi Jessica Li ✉

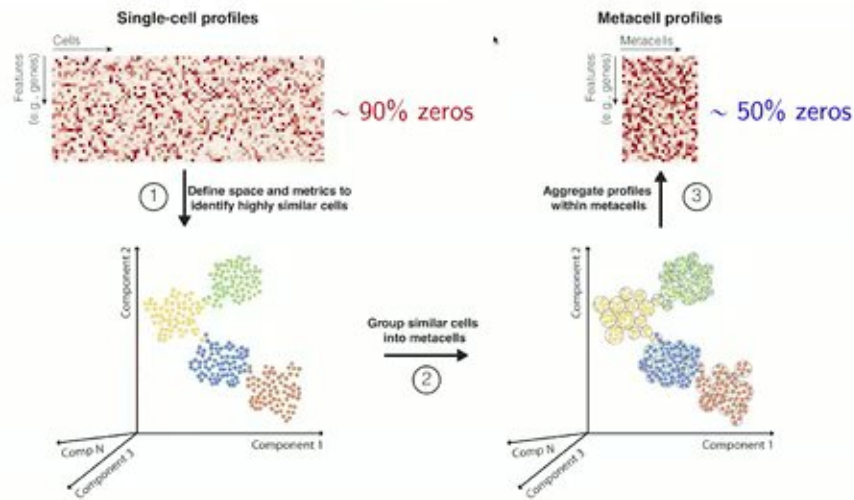*Nature Communications* **16**, Article number: 8602 (2025) | Cite this article

# How to define "metacell" in a statistical way?

*"A homogeneous collection of single-cell profiles that could have been resampled from the same original cell."*

$\implies$ Variation within a metacell is attributed exclusively to measurement error

**Two-layer observation model:**

$$\text{Expression model:} \quad \boldsymbol{\lambda}_i \sim \mathcal{F}(\cdot | \mathbf{x}_i)$$

$$\text{Measurement model:} \quad y_{ij} | \boldsymbol{\lambda}_i \overset{\text{ind}}{\sim} \text{Poisson}(c_i \lambda_{ij})$$

Cell (observation) $i = 1, \ldots, n$, Feature $j = 1, \ldots, p$

- $\lambda_{ij}$: the relative expression level of feature $j$ in cell $i$; $\boldsymbol{\lambda}_i = [\lambda_{i1}, \ldots, \lambda_{ip}]^{\top}$
- $\mathbf{x}_i \in \mathbb{R}^q$: the covariates of cell $i$, e.g., cell types, pseudotimes in lineage trajectories
- $\mathcal{F}$: the $p$-dim expression distribution that captures meaningful biological information
- $y_{ij}$: the observed count of feature $j$ in cell $i$; $y_{i+} = \sum_{j=1}^{p} y_{ij}$; $c_i = \mathsf{E}[y_{i+} | \boldsymbol{\lambda}_i]$

# Dubious metacells can bias analysis

Bilous, M., et al. "Building and analyzing metacells in single-cell genomics data." Molecular Systems Biology (2024): 1-23.

# Dubious metacells introduce bias and artifacts

Single cells

Metacells

Single-cell profiles
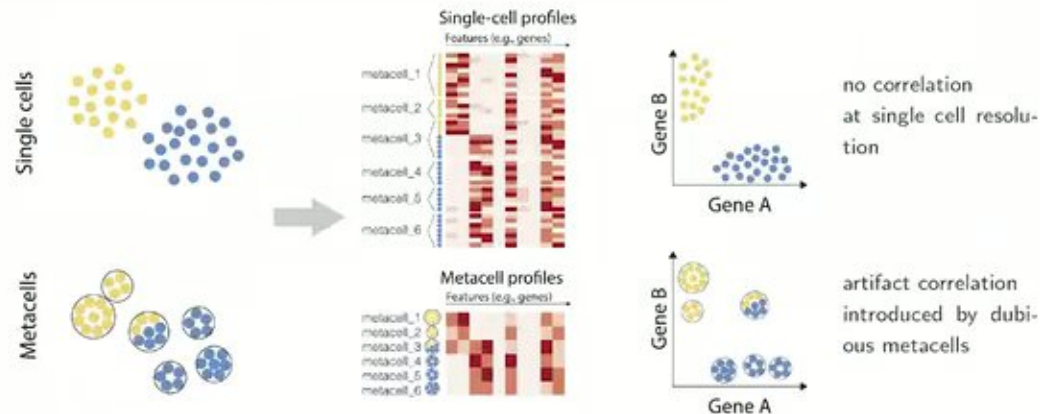Features (e.g. genes)

metacell_1
metacell_2
metacell_3
metacell_4
metacell_5
metacell_6

Gene B / Gene A

no correlation at single cell resolution

Metacell profiles
Features (e.g. genes)

metacell_1
metacell_2
metacell_3
metacell_4
metacell_5
metacell_6

Gene B / Gene A

artifact correlation introduced by dubious metacells

Bilous, Mariia, et al. "Building and analyzing metacells in single-cell genomics data." Molecular Systems Biology (2024): 1-23.
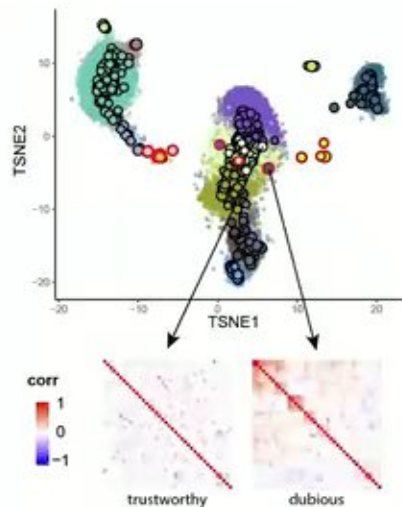
# Our proposal: mcRigor



**Goals:** a statistical criterion to

- Identify dubious metacells consisting of single cells from different cell states
- Nominate the top-performing metacell method and optimize its hyperparameter

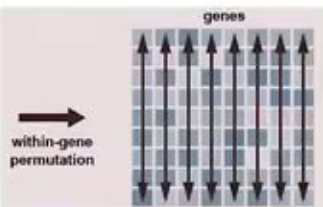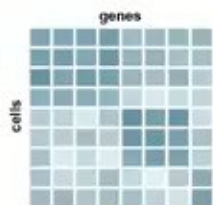$$\text{granularity level } \gamma = \frac{\#\text{single cells}}{\#\text{metacells}}$$

in a data-specific way

# Double permutation for null construction

Within-gene permutation:

- removes genes correlations
- removes cell library sizes
- preserves genes marginal distributions

Within-cell permutation:

- preserves cell library sizes
- removes genes correlations
- removes genes marginal distributions

genes

cells

within-gene permutation

mcDiv

within-cell permutation

permuted genes

cells

within-gene permutation

permuted genes

mcDiv$^{null}$

Normalization

Metacell: a heuristic solution to the sparsity issue in single-cell data

Bilous. M.. et al. "Building and analyzing metacells in single-cell genomics data." Molecular Systems Biology (2024): 1-23.

mcRigor method—mcDiv and mcDiv$^{null}$

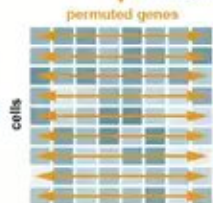# Double permutation for null construction

Within-gene permutation:

- removes genes correlations
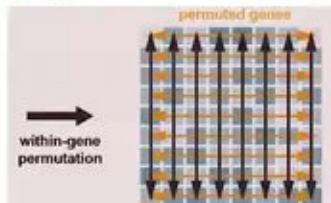- removes cell library sizes
- preserves genes marginal distributions

Within-cell permutation:

- preserves cell library sizes
- removes genes correlations
- removes genes marginal distributions

genes

cells

within-cell permutation

permuted genes

cells

genes

within-gene permutation

mcDiv

permuted genes

within-gene permutation

mcDiv$^{null}$

Normalization

## mcRigor method—mcDiv and mcDiv$^{null}$

# mcRigor method

# mcRigor effectively detects dubious metacells with high accuracy

Simulator: Song, D., et al. "scDesign3 generates realistic in silico data for multimodal single-cell and spatial omics." Nat Biotech (2024)
MetaCell: Baran, Y., et al. "MetaCell: analysis of single-cell RNA-seq data using K-nn graph partitions." Genome Biology (2019)

# Test of mcRigor on barcode multiplets

**Barcode multiplet:** a set of cell-like observations in which each observation is assigned a unique cell barcode but actually originates from the same physical cell.

Lareau, C. A., et al. "Droplet-based combinatorial indexing for massive-scale single-cell chromatin accessibility." Nat Biotech (2019)
**SuperCell**: Bilous, M., et al. "Metacells untangle large and complex single-cell transcriptome networks." BMC Bioinformatics (2022)

# mcRigor enhances gene co-expression analysis

## Gene-gene correlation matrix

by single cells / by trustworthy metacells / by all metacells

Healthy

COVID-19

p-value = 0.00043 / p-value = 7.6e-19 / p-value = 0.54632

enriched / not enriched

metacells / single cells

IGKV4−1 / IGLV2−14

IGHV4−59 / IGLV2−8

IGHV4−59 / CCR2

corr 1 0 −1 ■ antigen processing via MHC Class II ■ adaptive immune response ■ response to interferon-alpha

**SuperCell**: Bilous, M., et al. "Metacells untangle large and complex single-cell transcriptome networks." BMC Bioinformatics (2022)
**Data**: Wilk, A. J., et al. "A single-cell atlas of the peripheral immune response in patients with severe covid-19." Nat Med (2020)

## mcRigor improves the reliability of gene regulatory inference

very low signal — validated enhancer (LOC117038772)

*GATA2*

single-cell level

$\rho = 0.531$
p-value = 0.003

peak1

$\rho = 0.457$
p-value = 0.089

peak2

peak accessibility

gene expression

SEACells: Persad, S., et al. "SEACells infers transcriptional and epigenomic cellular states from single-cell genomics data." Nat Biotech (2023).

mcRigor optimizes hyperparameter $\gamma$ by balancing sparsity and dubiousness

mcRigor helps distinguish biological zeros from technical zeros

Data: Torre, Eduardo, et al. "Rare cell detection by single-cell RNA sequencing as guided by single-molecule RNA FISH." Cell systems (2018).

# mcRigor selects optimal metacell method and $\gamma$ for DEG detection

Optimal method-hyparameter configuration:

SEACells with $\gamma = 13$

F score for DE genes derived from different metacell partitions (q-val<0.05)

v.s.

DE genes derived from bulk RNA data (q-val<0.05)

Data: Chu, L.-F. et al. "Single-cell RNA-seq reveals novel regulators of human embryonic stem cell differentiation to definitive endoderm." Genome Biol. (2016).

# Q: How could dubious metacells be handled appropriately?

## mcRigor two-step: an extension of mcRigor

**Step 1:** A method–hyperparameter configuration (i.e., a metacell partitioning method with a granularity level $\gamma_1$) is either specified by the user or selected by mcRigor. This configuration is then applied to partition single cells into metacells. If mcRigor detects dubious metacells within the partition, it is re-applied to the same partition using a lower divergence score threshold (as below) to label more metacells as dubious.

$$\theta(m_k) = q_{0.85}\left(\left\{\mathrm{mcDiv}_{k'}^{\mathrm{null}} : m_{k'} \in [m_k - h, m_k + h], k' = 1, \ldots, M\right\}\right)$$

**Step 2:** The selected metacell partitioning method is re-applied to the subset of single cells that belong to the metacells now marked as dubious. This yields a refined metacell partition under a new granularity level $\gamma_2 < \gamma_1$, which can be selected by mcRigor from the candidate set of granularity levels $2, \ldots, \gamma_1 - 1$.

# mcRigor two-step effectively resolves rare cell types



SEACells: Persad, S., et al. "SEACells infers transcriptional and epigenomic cellular states from single-cell genomics data." Nat Biotech (2023).
Data: Stuart, T., et al. "Comprehensive integration of single-cell data." Cell (2019).

# mcRigor two-step effectively resolves rare cell types



MetaCell *Baran, Y., et al. "MetaCell: analysis of single-cell RNA-seq data using K-nn graph partitions." Genome Biology (2019).*
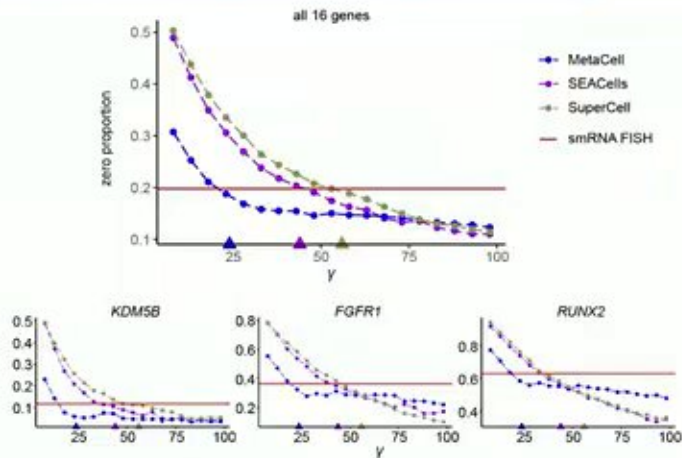SuperCell *Bilous, M., et al. "Metacells untangle large and complex single-cell transcriptome networks." BMC Bioinformatics (2022).*
Data: *Stuart, T., et al. "Comprehensive integration of single-cell data." Cell (2019).*

# Conclusion

A: We give a statisical definition of "metacell" based on the two-layer model.

A: We detect dubious metacells using per-metacell mcDiv statistics and null contructed through double permutation.

A: We optimize metacell partitioning by balancing sparsity and dubiousness.

Article | Open access | Published: 29 September 2025

## mcRigor: a statistical method to enhance the rigor of metacell partitioning in single-cell data analysis

Pan Liu & Jingyi Jessica Li ✉

*Nature Communications* **16**, Article number: 8602 (2025) | Cite this article

# R package and tutorial on Github

mcRigor 1.0    Reference    Articles ▾

# mcRigor

Functionality 1: detect dubious metacells for a given metacell partition
Functionality 2: optimize metacell partitioning
Implementing metacell partitioning methods
Extension: mcRigor two-step

The R package **mcRigor** is a statistical method to enhance the rigor of metacell partitioning in single-cell data analysis. It can be used as an add-on for any existing metacell partitioning methods for obtaining more reliable metacells. mcRigor has two main functionalities: 1) detecting dubious metacells, which are composed of heterogeneous single cells, for a given metacell partition, and 2) optimizing the hyperparameter of a metacell partitioning method. The core of mcRigor is a feature-correlation-based statistic that measures the heterogeneity of a metacell, with its null distribution derived from a double permutation scheme. The following figure illustrates the schematics of mcRigor for dubious metacell detection (a) and hyperparameter optimization (b).

**License**

MIT + file LICENSE

**Citation**

Citing mcRigor

**Developers**

Pan Liu
Maintainer

a



**R package:** https://github.com/JSB-UCLA/mcRigor    **Tutorial:** https://jsb-ucla.github.io/mcRigor