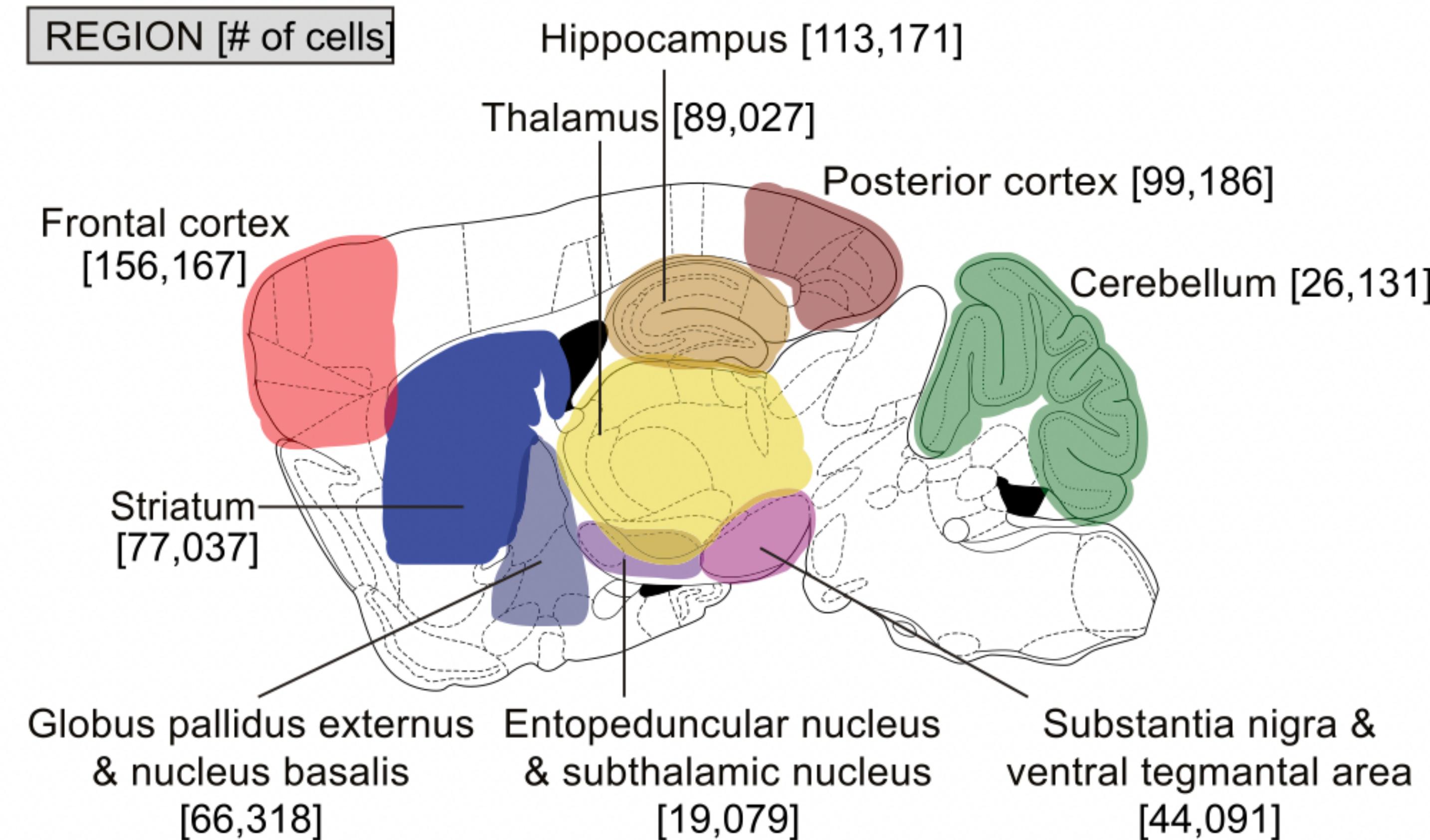


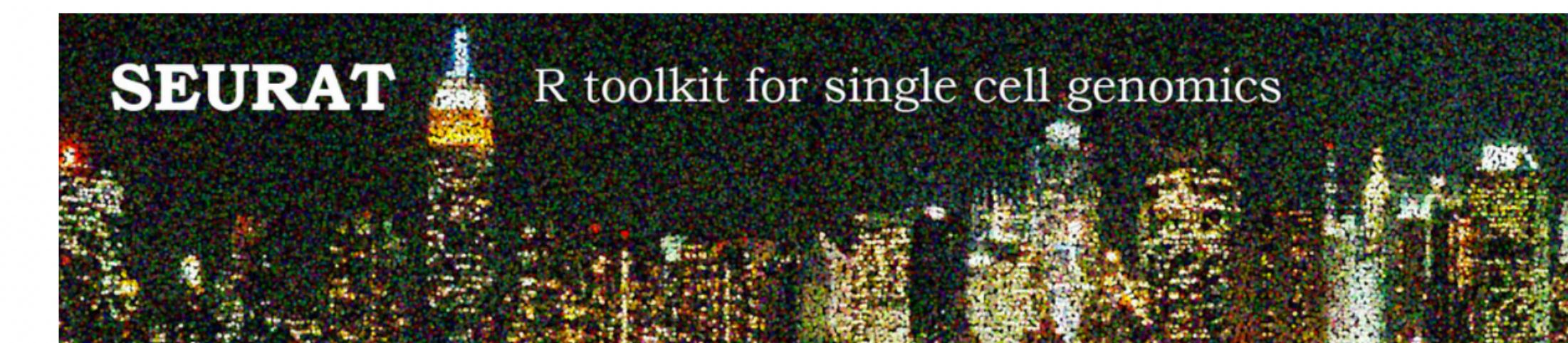
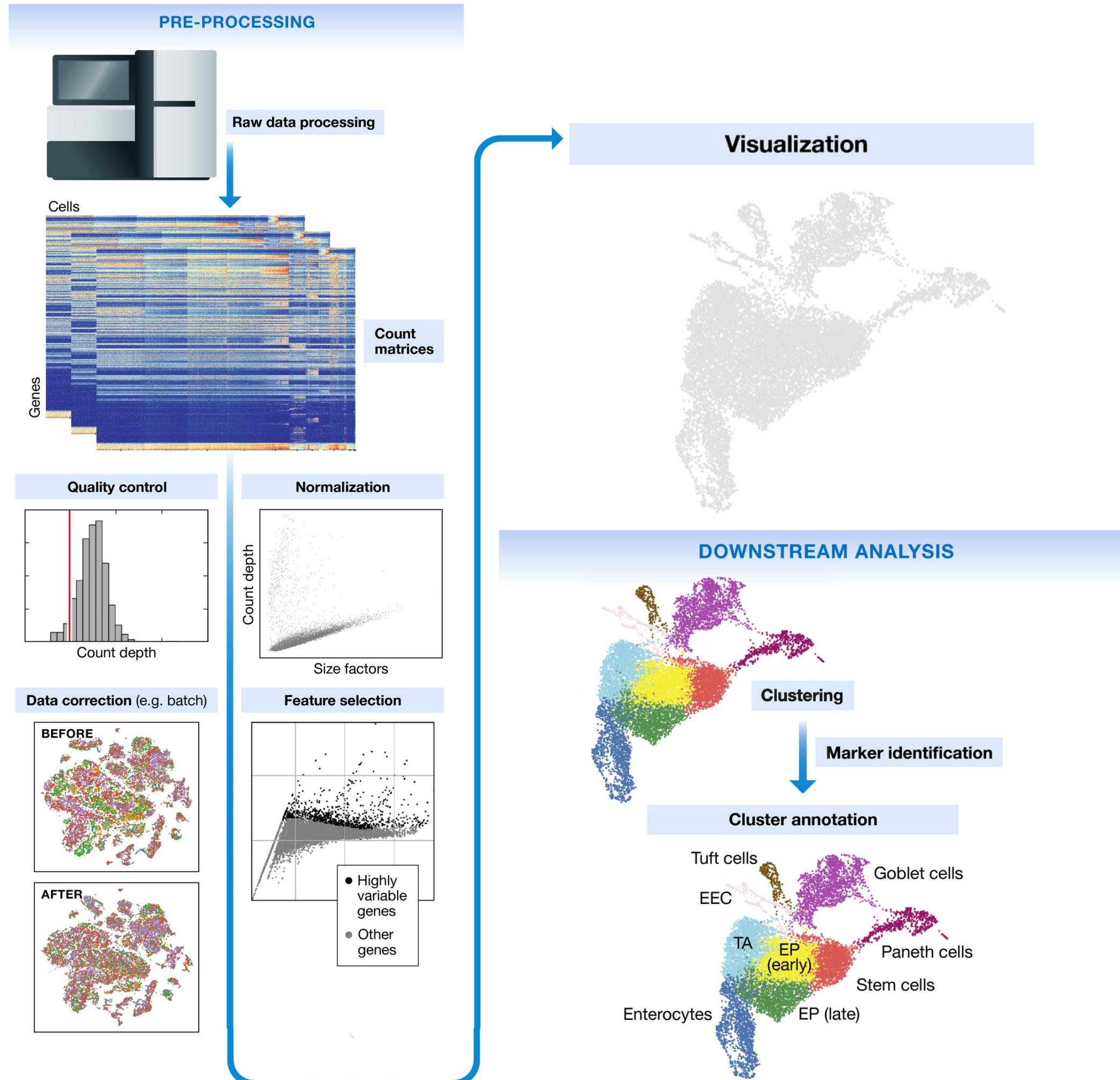
Computational Methods for Phenotyping a Perturbation in Single Cell Data

Stanley Center Primer
11/27/2022

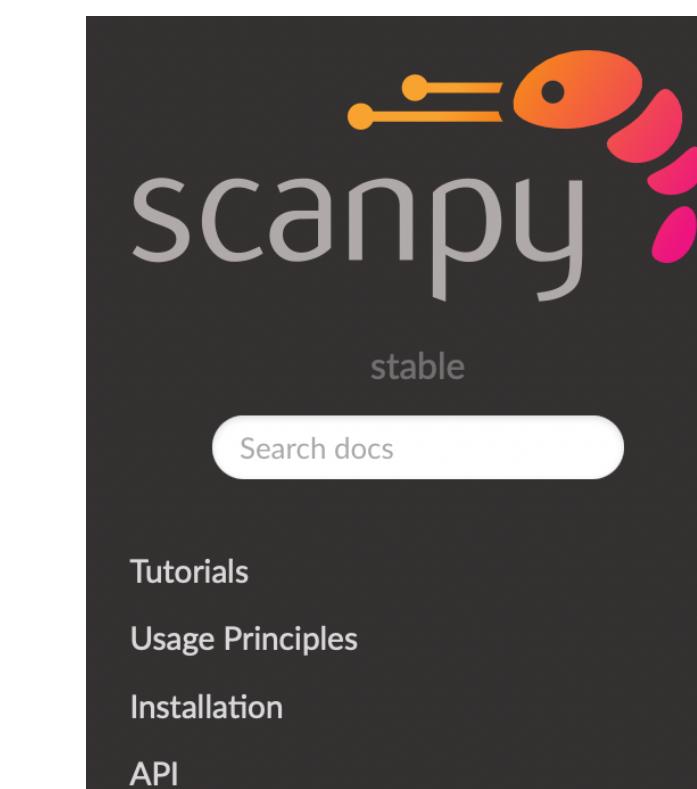
Single cell atlases create comprehensive maps of transcriptional heterogeneity in complex tissues



We have an established pipeline for interrogating the heterogeneity in a sample



Official release of Seurat 4.0



» Scanpy – Single-Cell Analysis in Python

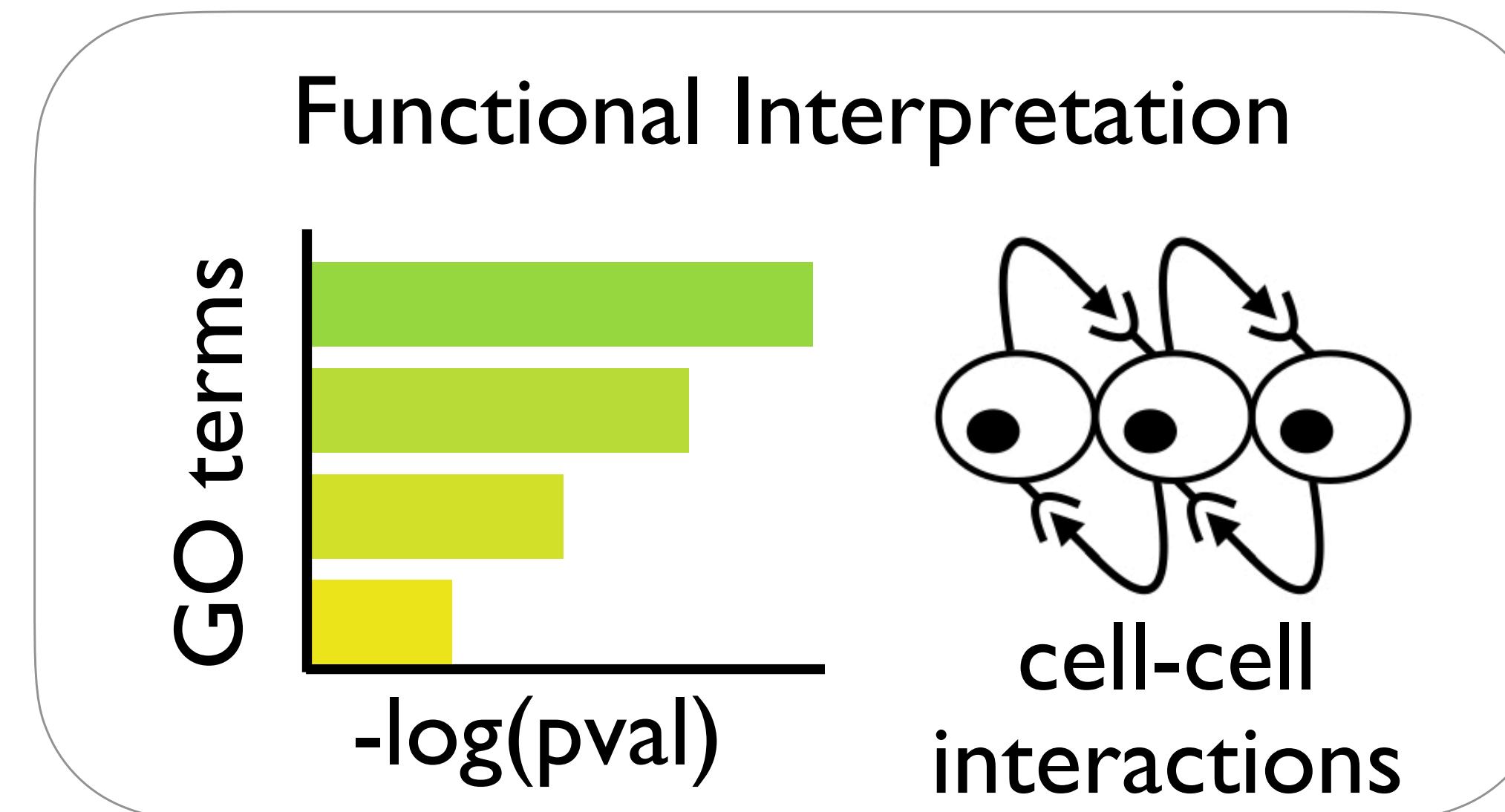
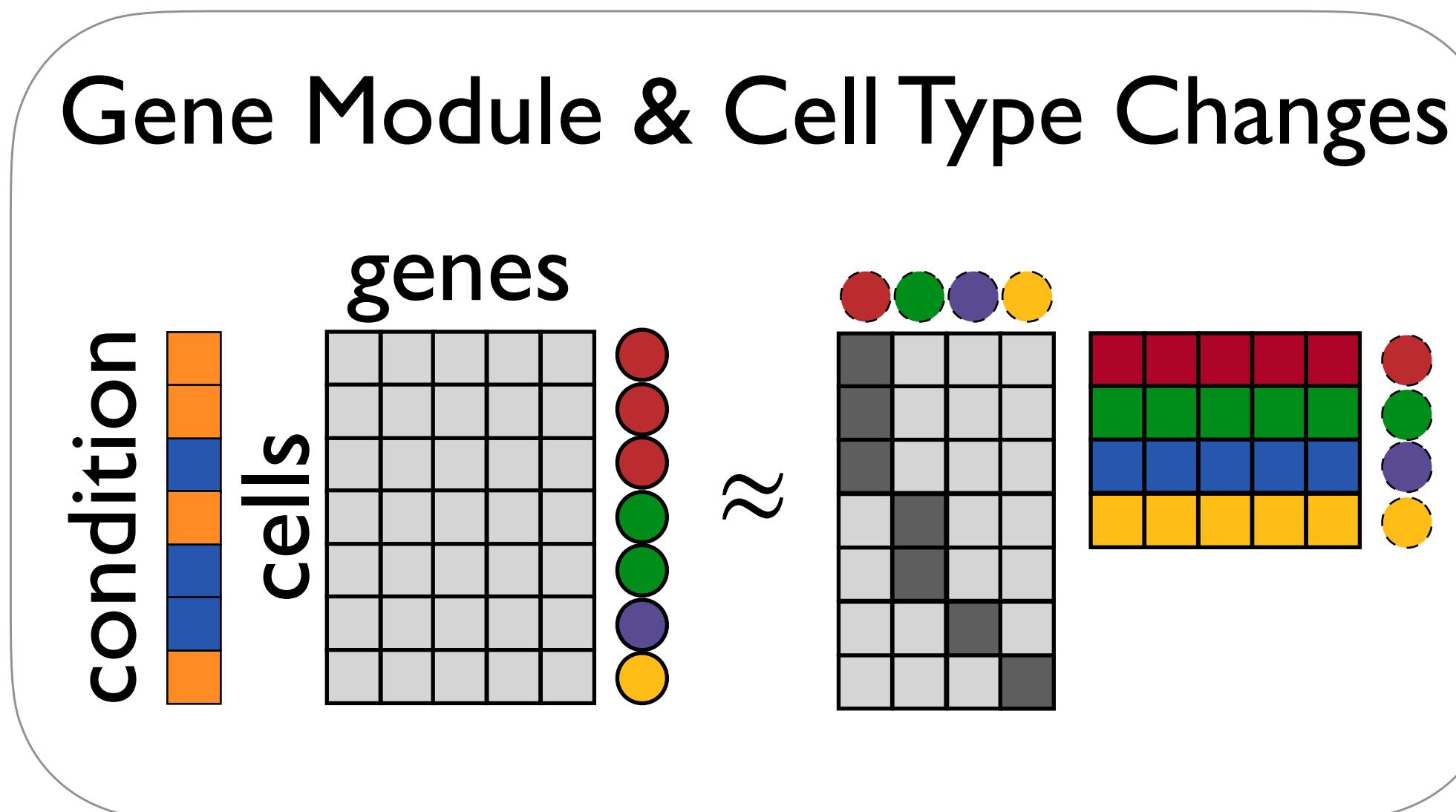
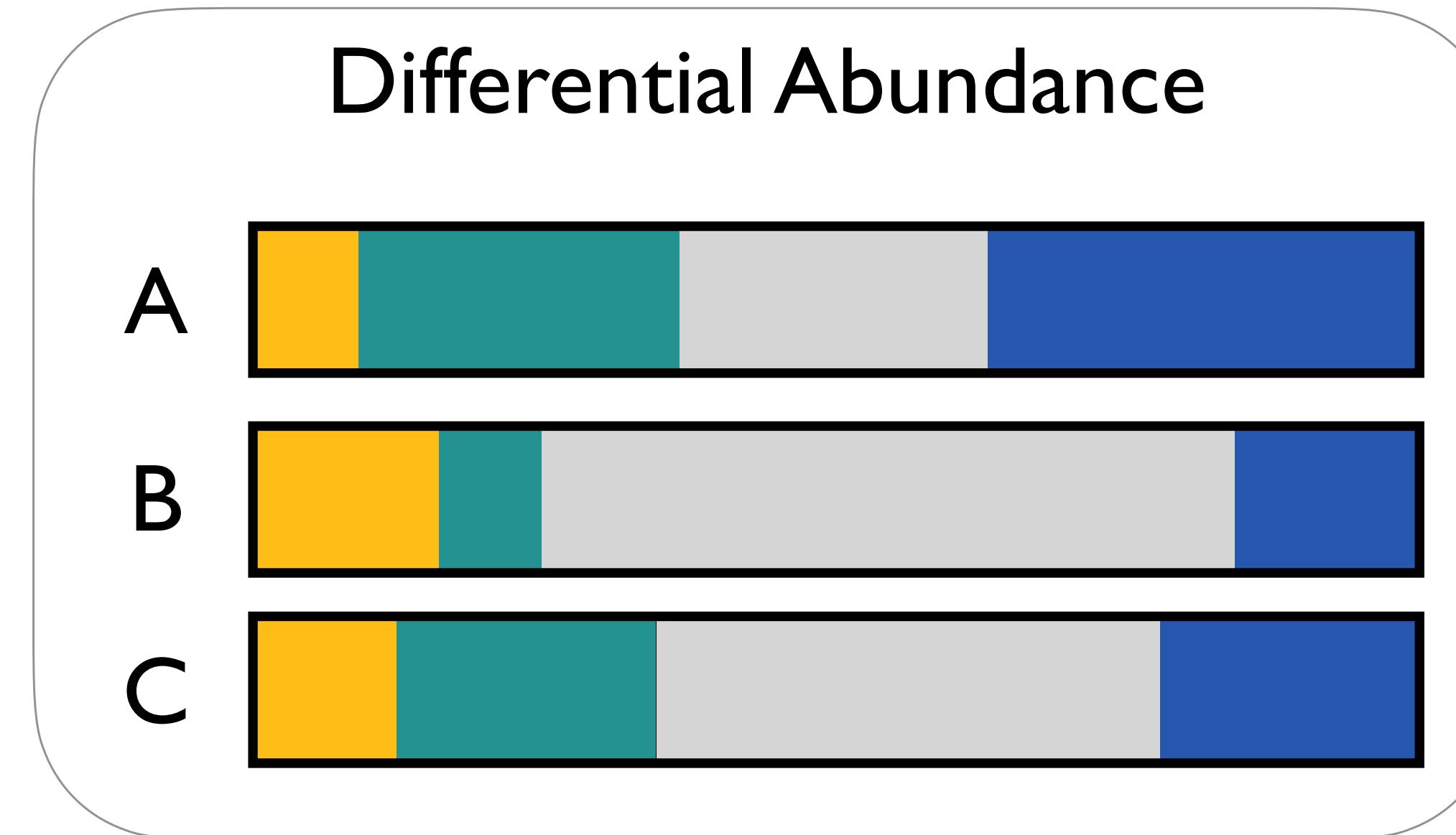
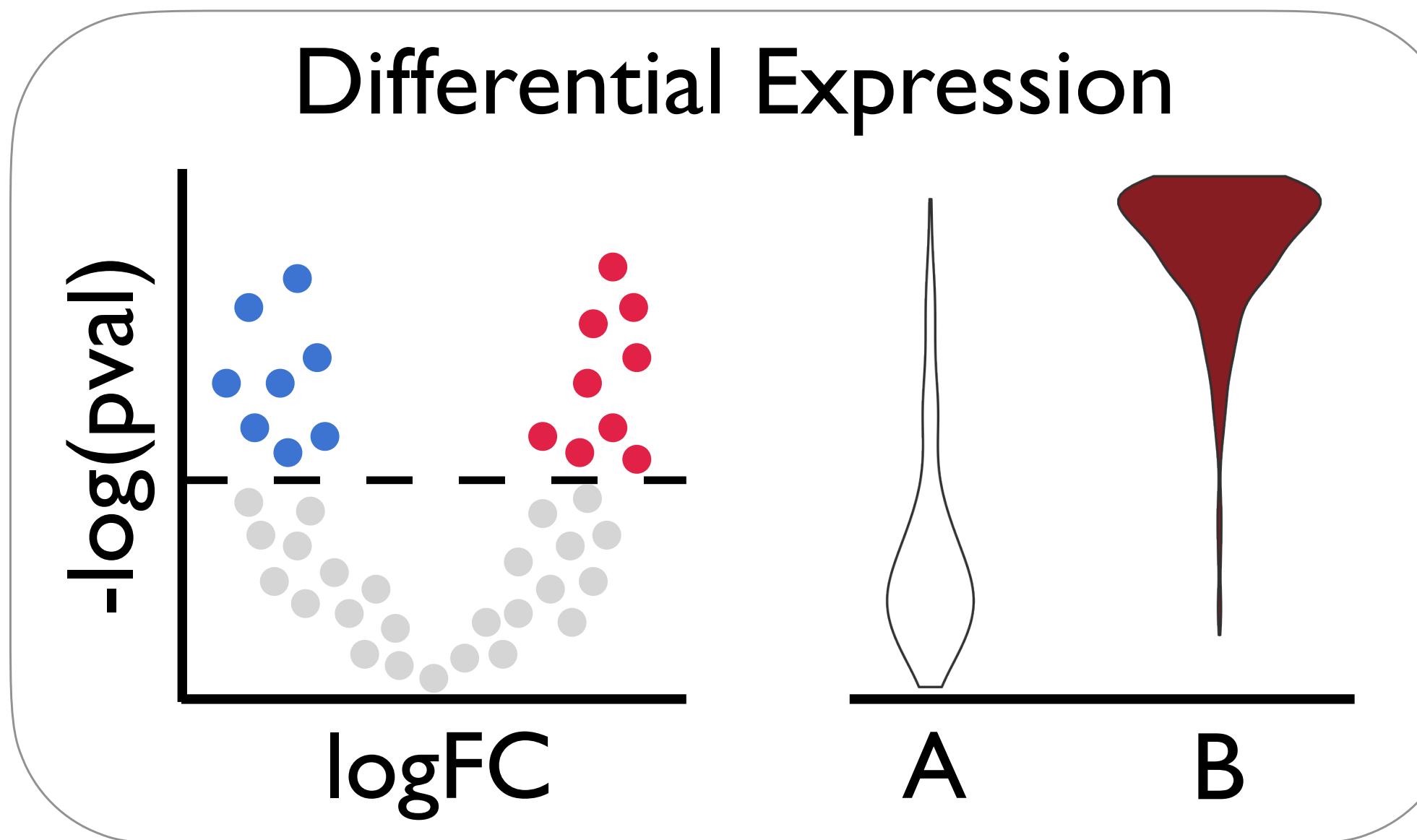
Edit on GitHub

stars 1.3k pypi v1.9.1 downloads 1M downloads 29k docs passing
Azure Pipelines succeeded discourse 1.8k posts zulip join chat

Scanpy – Single-Cell Analysis in Python

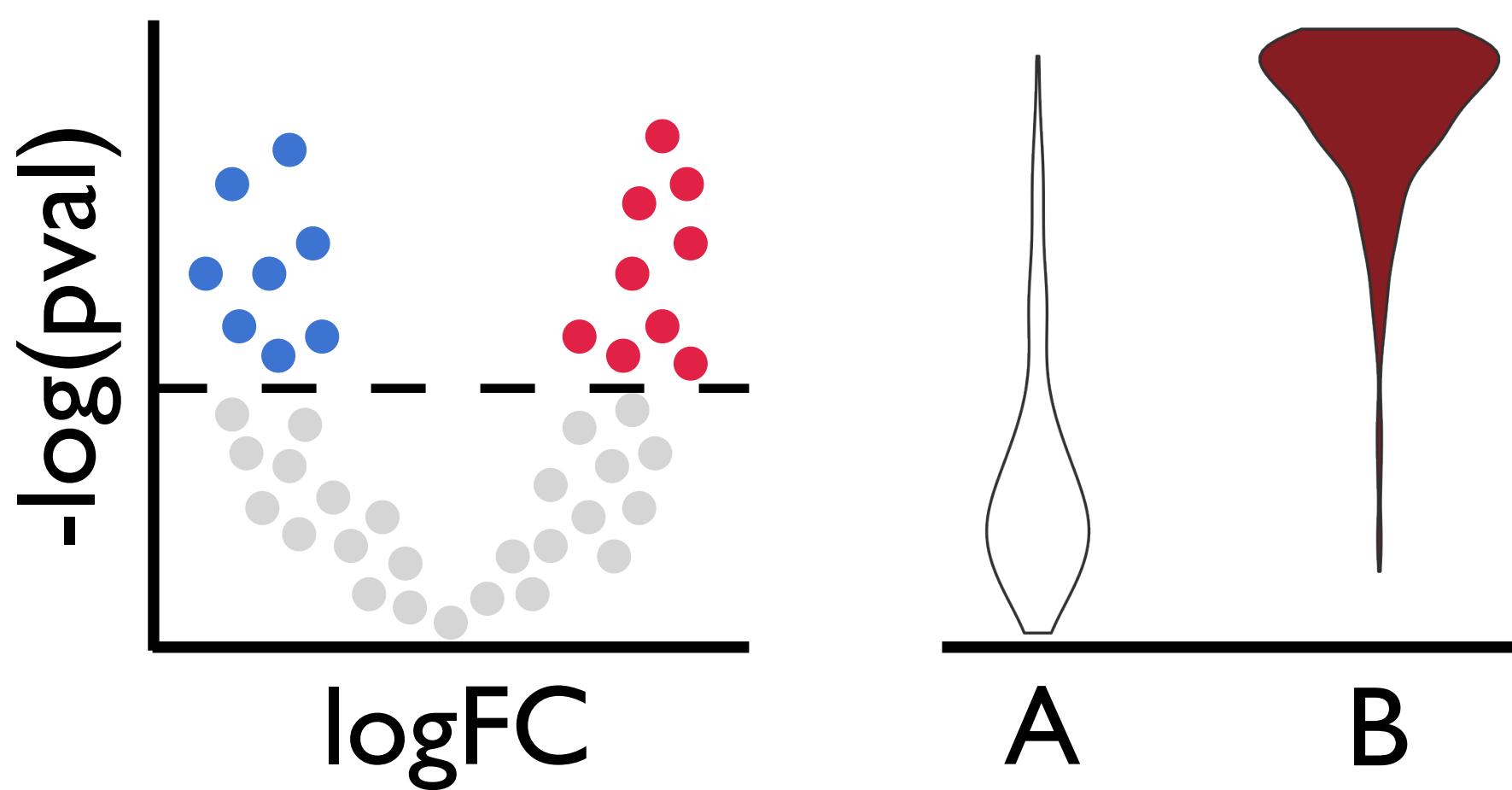
Scanpy is a scalable toolkit for analyzing single-cell gene expression data built jointly with [anndata](#). It includes preprocessing, visualization, clustering, trajectory inference and differential expression testing. The Python-based implementation efficiently deals with datasets of more than one million cells.

Today we will focus on methods for finding differences across conditions



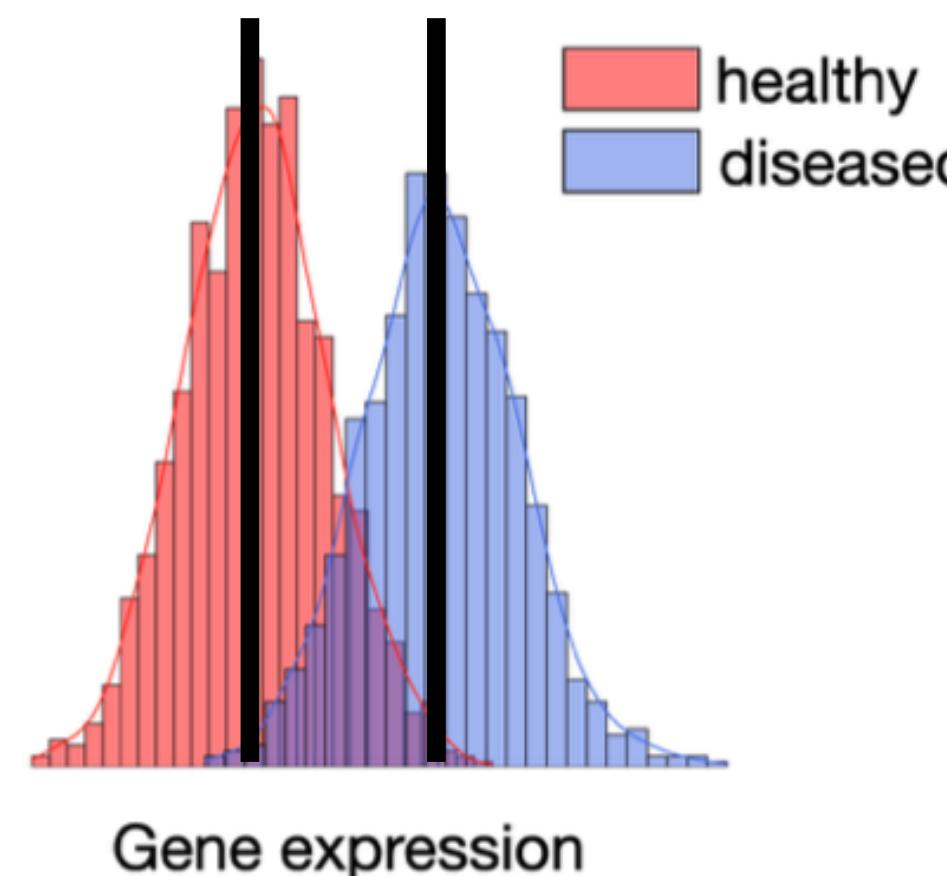
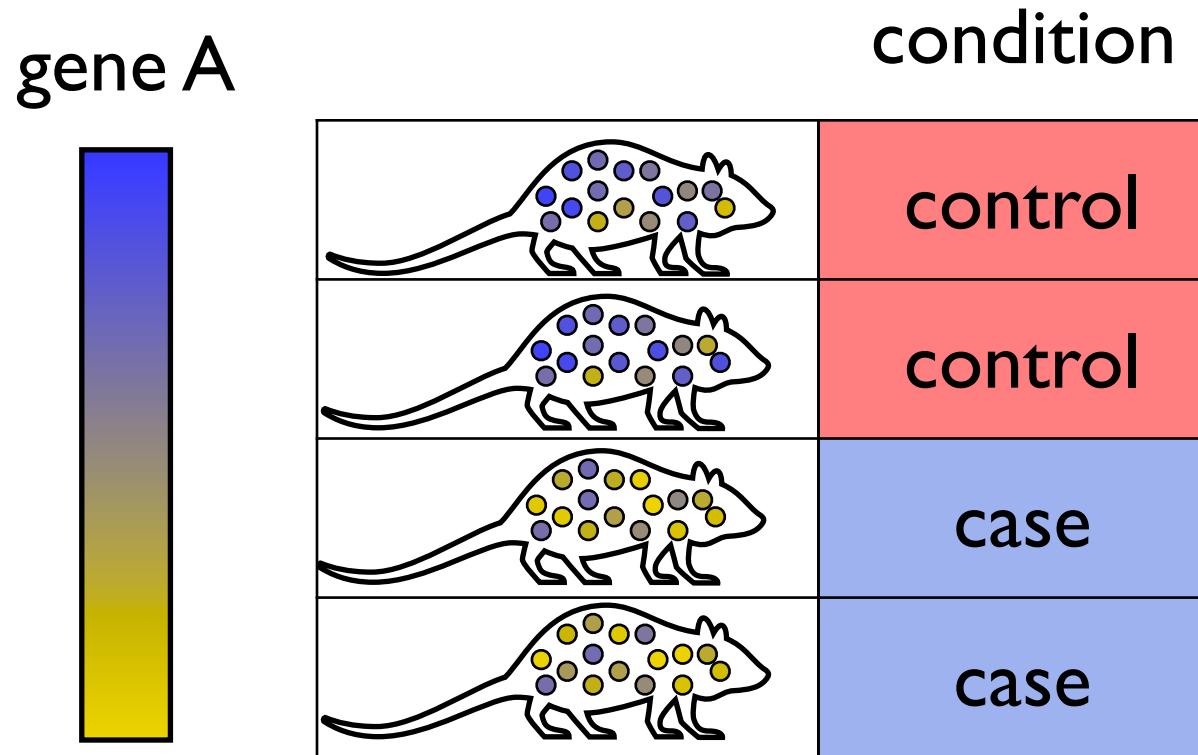
Differential Expression

Identification of genes associated with a perturbation



Condition-level aggregation

The default method



Seurat - Guided Clustering Tutorial

Compiled: January 11, 2022

Finding differentially expressed features (cluster biomarkers)

```
# find all markers of cluster 2
cluster2.markers <- FindMarkers(pbmc, ident.1 = 2, min.pct = 0.25)
head(cluster2.markers, n = 5)
```

Gene expression markers of identity classes

Finds markers (differentially expressed genes) for identity classes

```
FindMarkers(object, ...)

# S3 method for default
FindMarkers(
  object,
  slot = "data",
  counts = numeric(),
  cells.1 = NULL,
  cells.2 = NULL,
  features = NULL,
  logfc.threshold = 0.25,
  test.use = "wilcox",
```

scipy.tl.rank_genes_groups

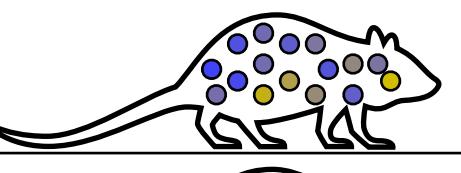
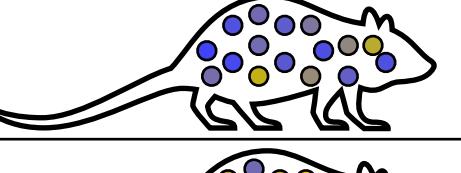
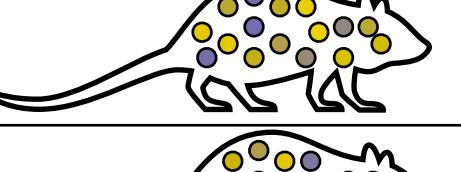
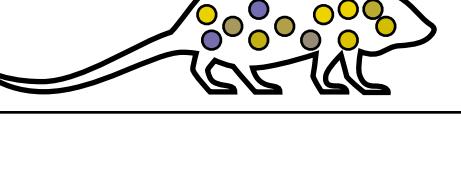
method : `Optional[Literal['logreg', 't-test', 'wilcoxon', 't-test_overestim_var']]` (default: `None`)

The default method is `'t-test'`, `'t-test_overestim_var'` overestimates variance of each group, `'wilcoxon'` uses Wilcoxon rank-sum, `'logreg'` uses logistic regression. See [Ntranos18], [here](#) and [here](#), for why this is meaningful.

Sample-level aggregation

Allows for accounting for covariates

gene A

	condition	sex	weight	animal ID
	control	F	53	1
	control	M	61	2
	case	F	55	3
	case	M	58	4

[Bioinformatics](#). 2010 Jan 1; 26(1): 139–140.

Published online 2009 Nov 11. doi: [10.1093/bioinformatics/btp616](https://doi.org/10.1093/bioinformatics/btp616)

PMCID: PMC2796818

PMID: [19910308](#)

edgeR: a Bioconductor package for differential expression analysis of digital gene expression data

[Mark D. Robinson](#),^{1,2,*†} [Davis J. McCarthy](#),^{1,2,†} and [Gordon K. Smyth](#)²

Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2

[Michael I Love](#), [Wolfgang Huber](#) & [Simon Anders](#)✉

[Genome Biology](#) 15, Article number: 550 (2014) | [Cite this article](#)

395k Accesses | 28144 Citations | 110 Altmetric | [Metrics](#)

PMCID: PMC4402510

PMID: [25605792](#)

$$\begin{matrix} \text{color bar} \\ = \end{matrix} \text{NB}\left(\beta_0 + \beta_1 \begin{matrix} \text{color bar} \\ \text{red} \end{matrix} + \beta_2 \begin{matrix} \text{color bar} \\ \text{yellow} \end{matrix} + \beta_2 \begin{matrix} \text{color bar} \\ \text{grey} \end{matrix}, \sigma^2\right)$$

[Nucleic Acids Res.](#). 2015 Apr 20; 43(7): e47.

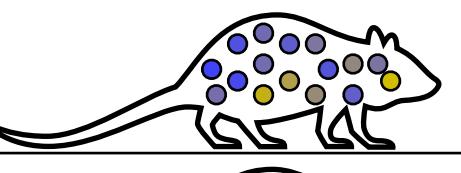
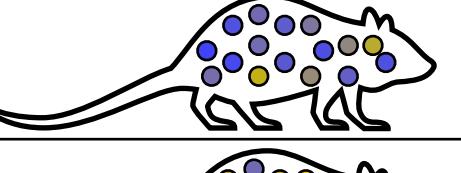
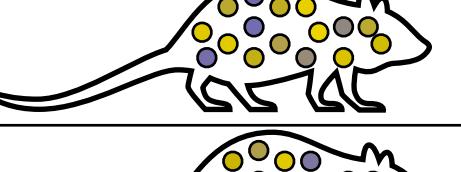
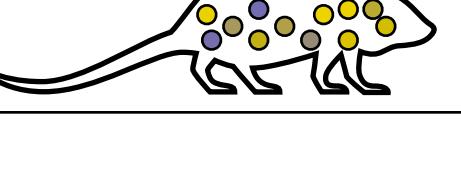
Published online 2015 Jan 20. doi: [10.1093/nar/gkv007](https://doi.org/10.1093/nar/gkv007)

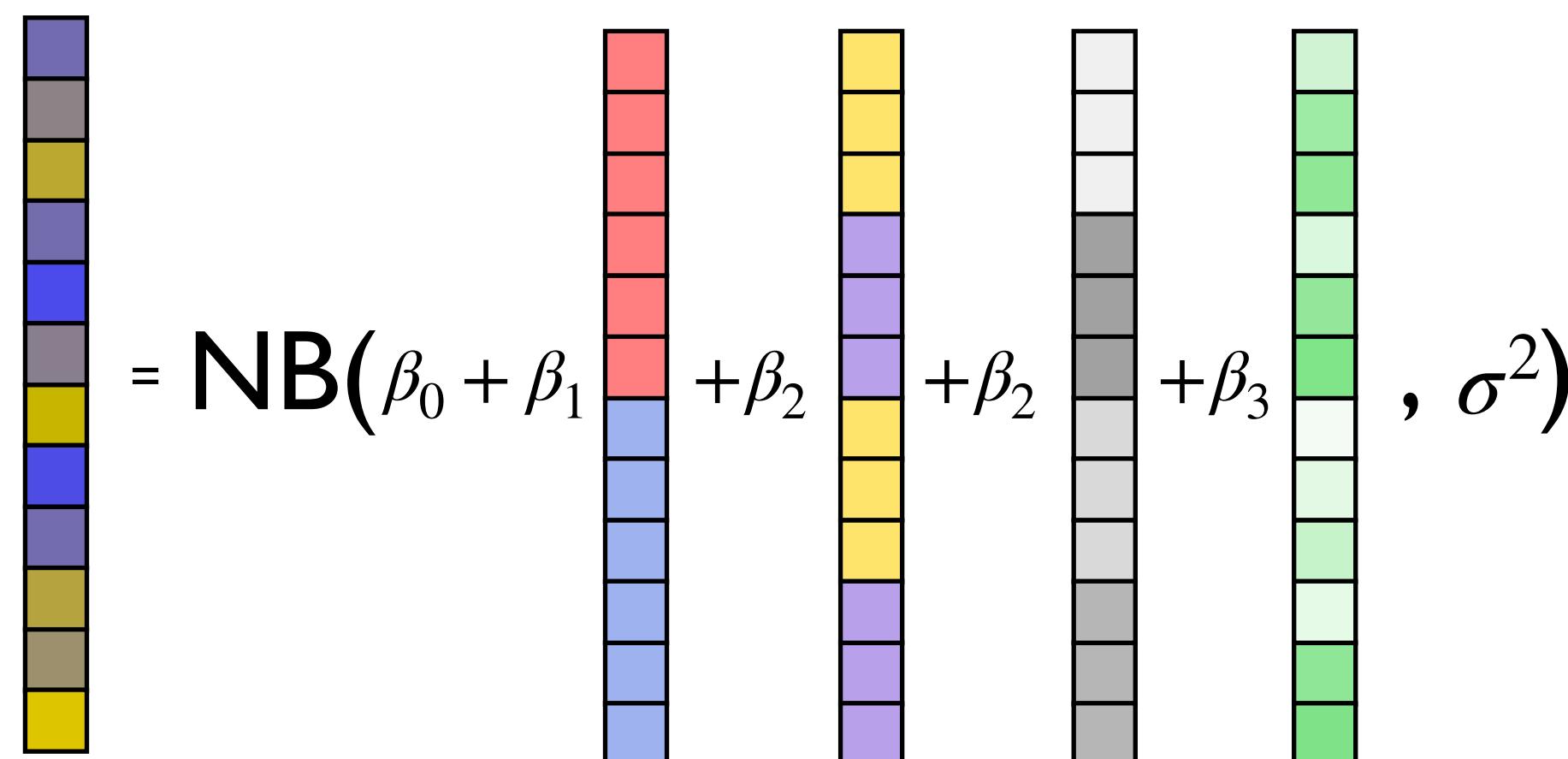
limma powers differential expression analyses for RNA-sequencing and microarray studies

[Matthew E. Ritchie](#),^{1,2} [Belinda Phipson](#),³ [Di Wu](#),⁴ [Yifang Hu](#),⁵ [Charity W. Law](#),⁶ [Wei Shi](#),^{5,7} and [Gordon K. Smyth](#)^{2,5,*}

Bespoke single-cell methods

Can account for cell-level covariates

gene A	condition	sex	weight	animal ID
	control	F	53	1
	control	M	61	2
	case	F	55	3
	case	M	58	4



glmGamPoi: fitting Gamma-Poisson generalized linear models on single cell count data 

Constantin Ahlmann-Eltze , Wolfgang Huber

Bioinformatics, Volume 36, Issue 24, 15 December 2020, Pages 5701–5702,
<https://doi.org/10.1093/bioinformatics/btaa1009>

MAST: a flexible statistical framework for assessing transcriptional changes and characterizing heterogeneity in single-cell RNA sequencing data

[Greg Finak](#), [Andrew McDavid](#), [Masanao Yajima](#), [Jingyuan Deng](#), [Vivian Gersuk](#), [Alex K. Shalek](#), [Chloe K. Slichter](#), [Hannah W. Miller](#), [M. Juliana McElrath](#), [Martin Prlic](#), [Peter S. Linsley](#) & [Raphael Gottardo](#) 

Genome Biology **16**, Article number: 278 (2015) | [Cite this article](#)

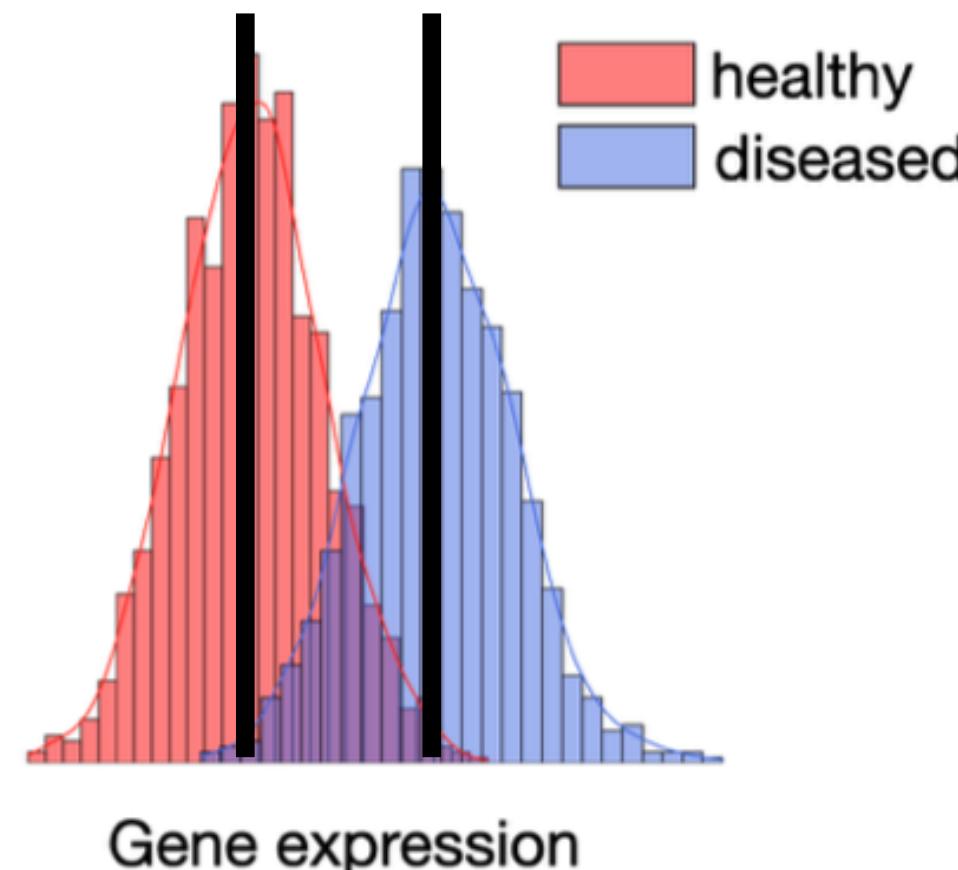
Differential Expression

Summary

gene A	condition	sex	weight	animal ID
	control	F	53	1
	control	M	61	2
	case	F	55	3
	case	M	58	4

t-test, Wilcoxon

- one gene at a time
- average per condition



edgeR, DEseq2, limma

- regress one gene at a time
- count model
- typically used in pseudo-bulk mode
- can be used per cell but are slow
- can account for covariates

$$= \text{NB}(\beta_0 + \beta_1 + \beta_2 + \beta_3, \sigma^2)$$

glmGamPoi

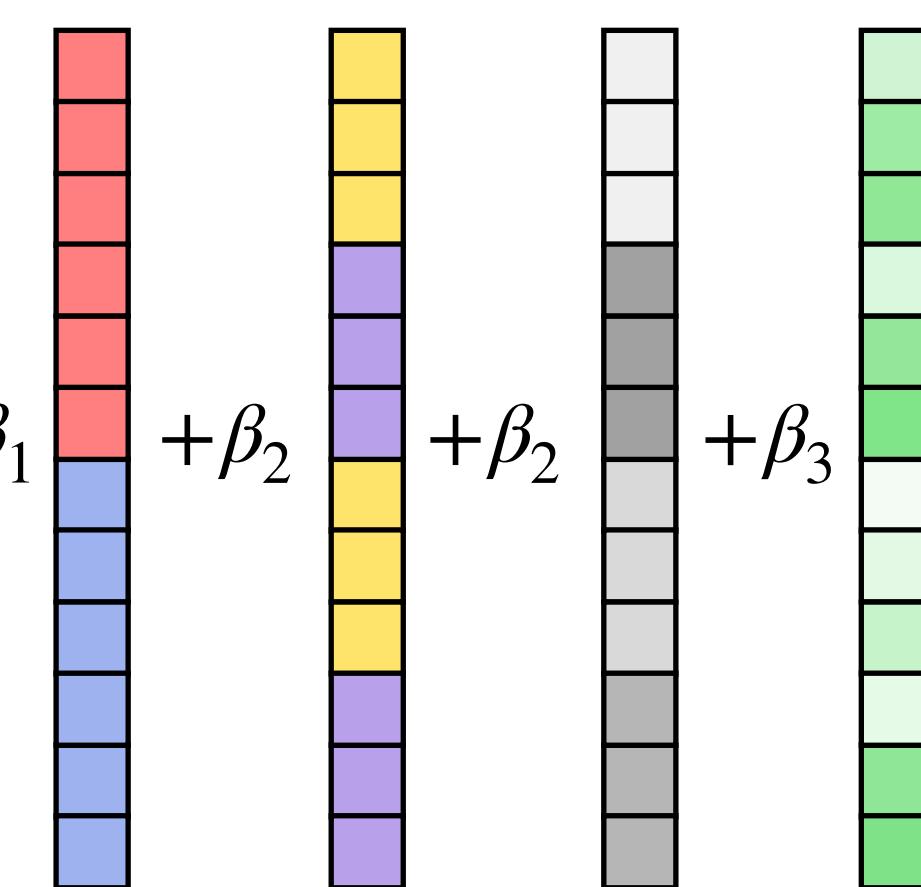
- fast alternative to edgeR
- ran per cell
- fixed effects only



$$= \text{NB}(\beta_0 + \beta_1 + \beta_2 + \beta_3, \sigma^2)$$

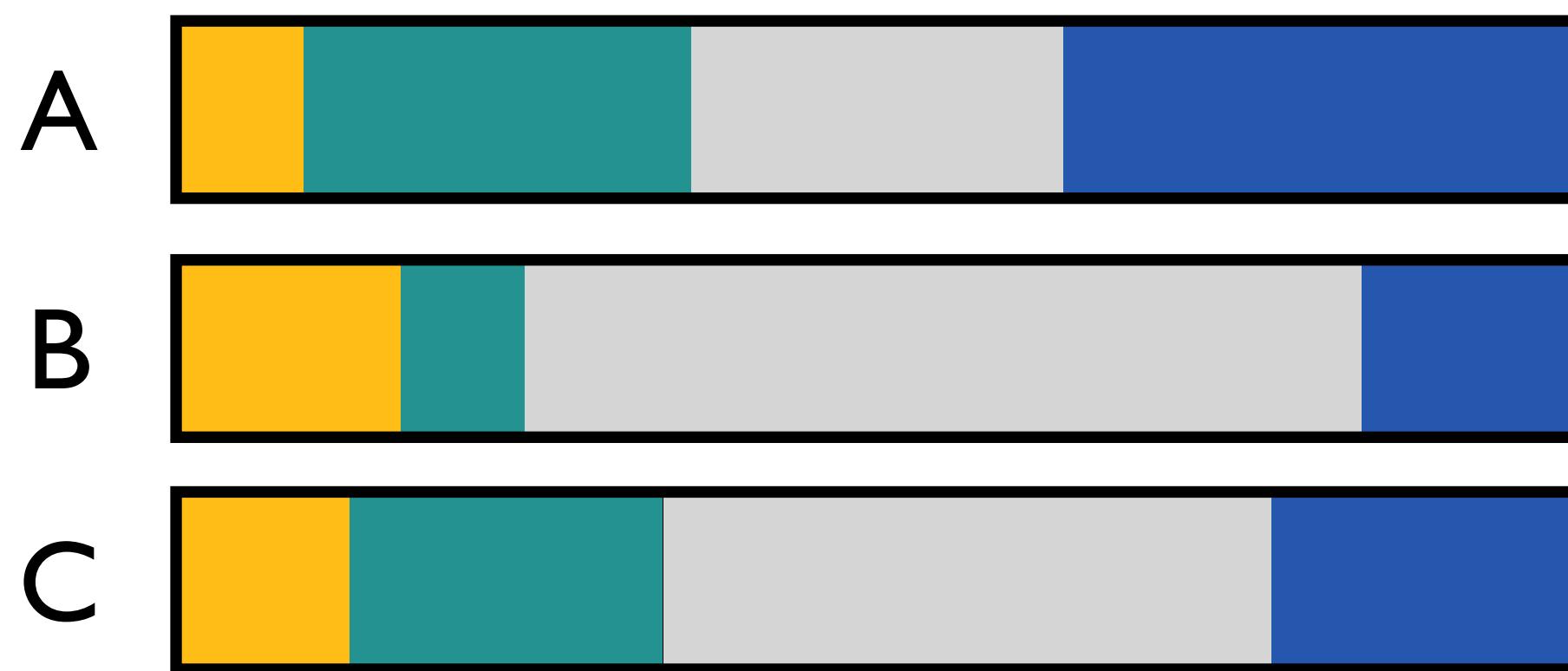
MAST

- ran per cell
- can include random effects



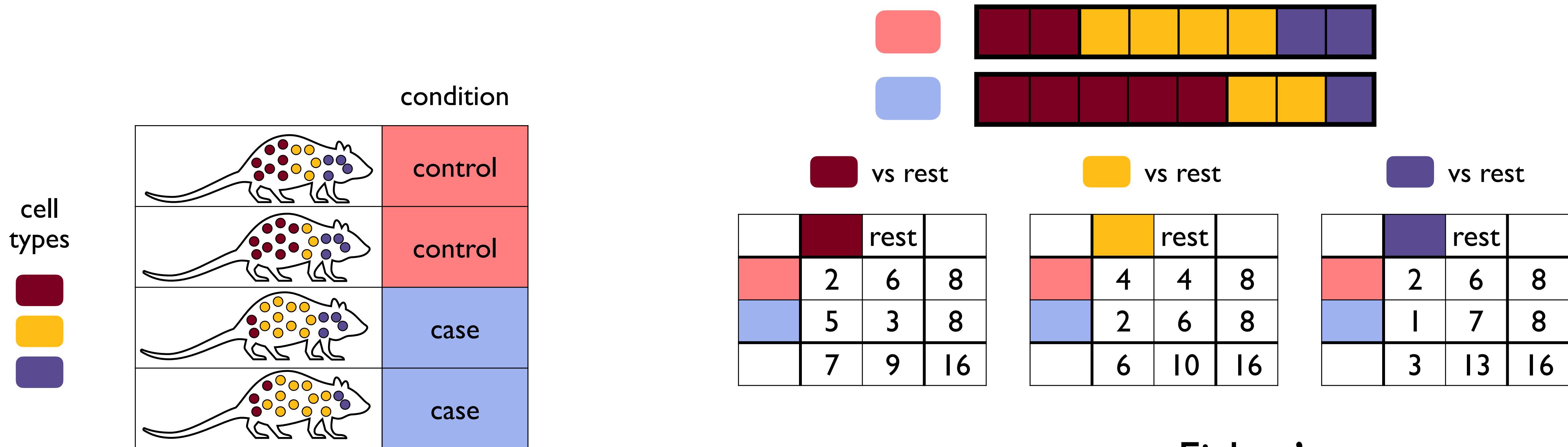
Differential Abundance

Identification of cell type proportion changes across conditions

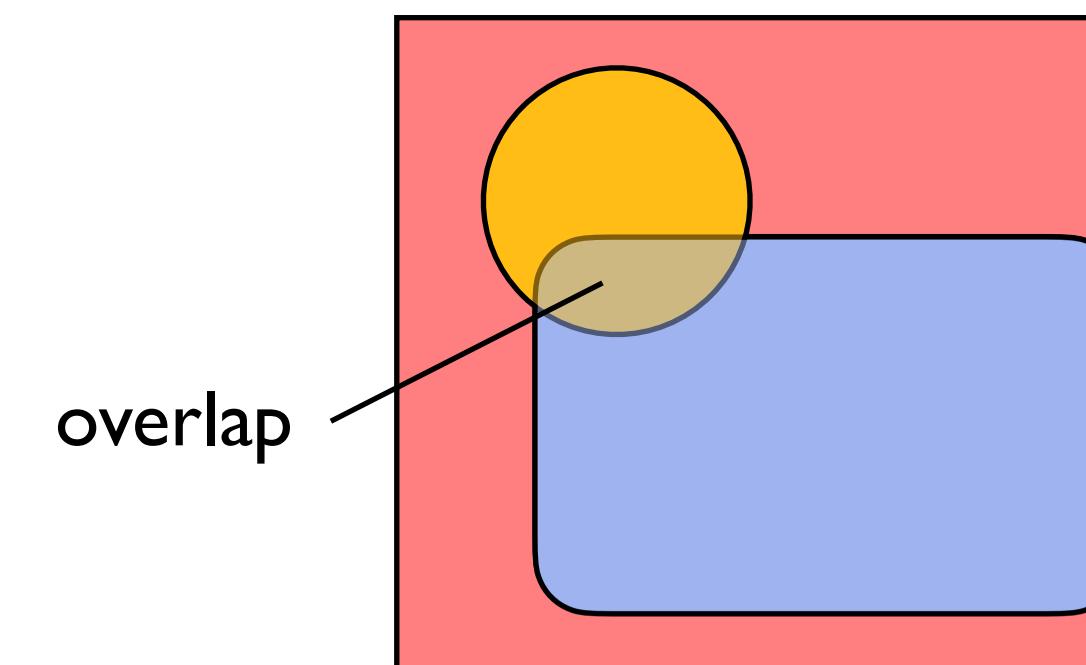


Univariate tests

Reflects how enriched each cell type is across conditions, but does not account for how changes in other cell types might affect that



$$\frac{\binom{\text{total case}}{\text{yellow case}} \binom{\text{total control}}{\text{yellow control}}}{\binom{\text{all cells}}{\text{total yellow}}}$$

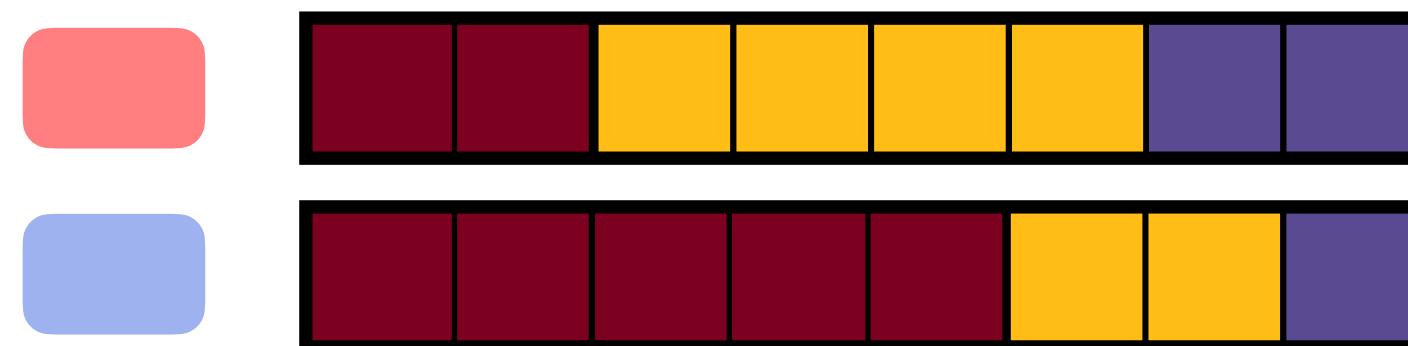


Fisher's exact test

- test of independence of row and column variable in the 2x2 table
- assumes fixed row and col totals
- equivalent to hypergeometric
- relax fixed marginals assumptions
- elevated FDR

Cell type proportions are not independent of each other

Since all proportions sum to 1, an increase in the proportion of one cell subset will necessarily lead to a decrease in the proportions of other cell subsets



DirichletReg: Dirichlet Regression for Compositional Data in R

Marco Maier

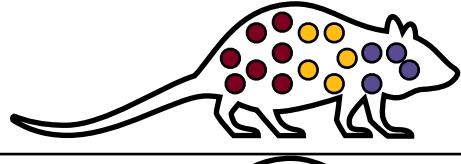
Institute for Statistics and Mathematics

[nature](#) > [nature communications](#) > [articles](#) > article

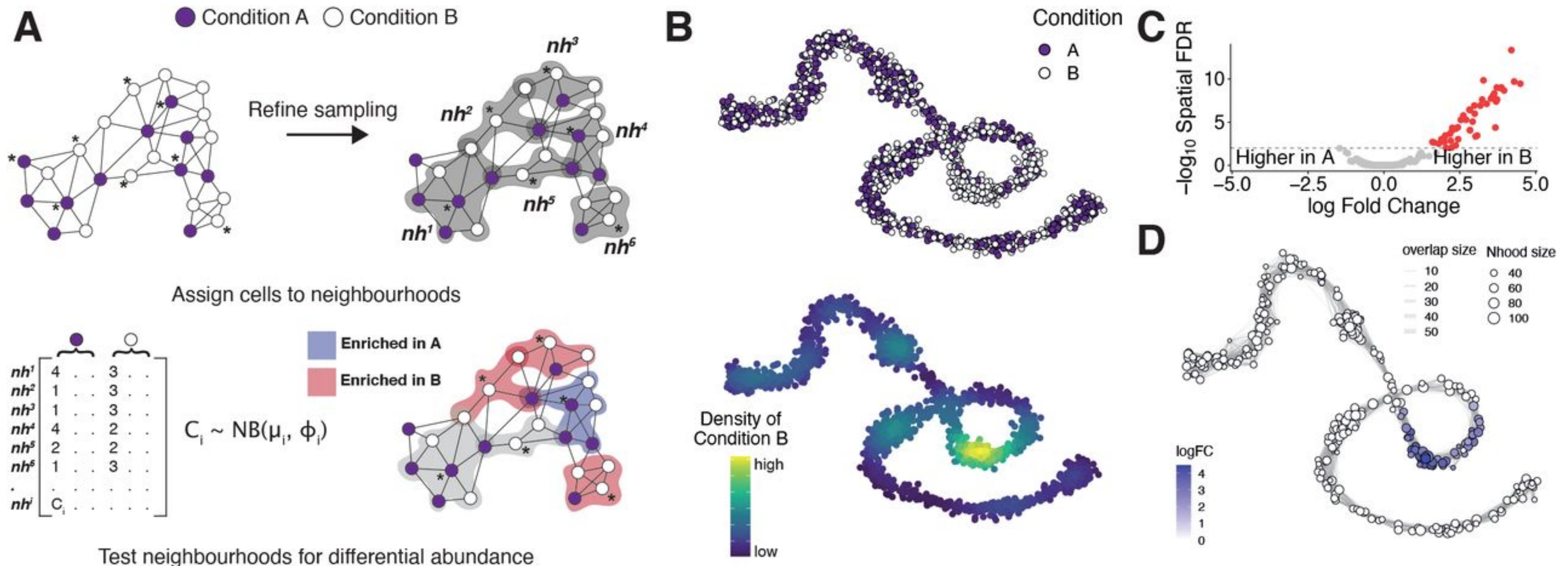
Article | [Open Access](#) | Published: 25 November 2021

scCODA is a Bayesian model for compositional single-cell data analysis

[M. Büttner](#), [J. Ostner](#), [C. L. Müller](#) , [F. J. Theis](#) & [B. Schubert](#) 

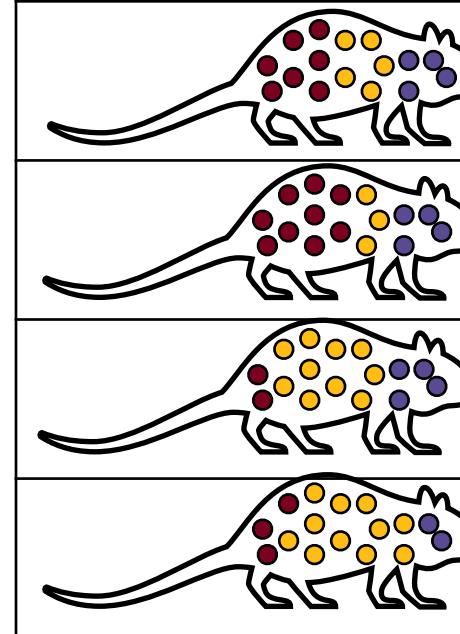
cell types	condition	sex	weight	animal ID
	control	F	53	1
	control	M	61	2
	case	F	55	3
	case	M	58	4

Compositional changes when cells lie on a continuum



Differential Abundance

Summary

cell types	condition	sex	weight	animal ID
	control	F	53	1
	control	M	61	2
	case	F	55	3
	case	M	58	4

Univariate tests

- one vs all
- Fisher's exact test
- elevated FDR

Dirichlet-multinomial regression

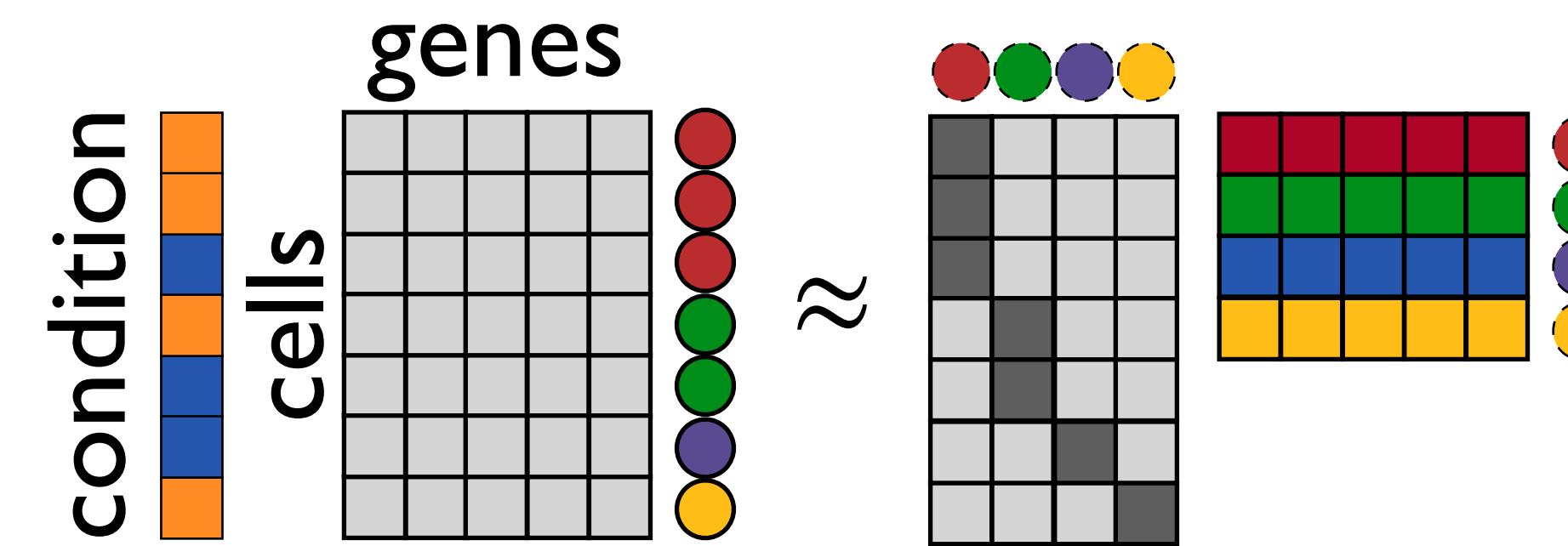
- accounts for proportions of all cell types
- can account for covariates

Milo

- differential abundance testing when cells are on a continuum rather than belonging to discrete clusters

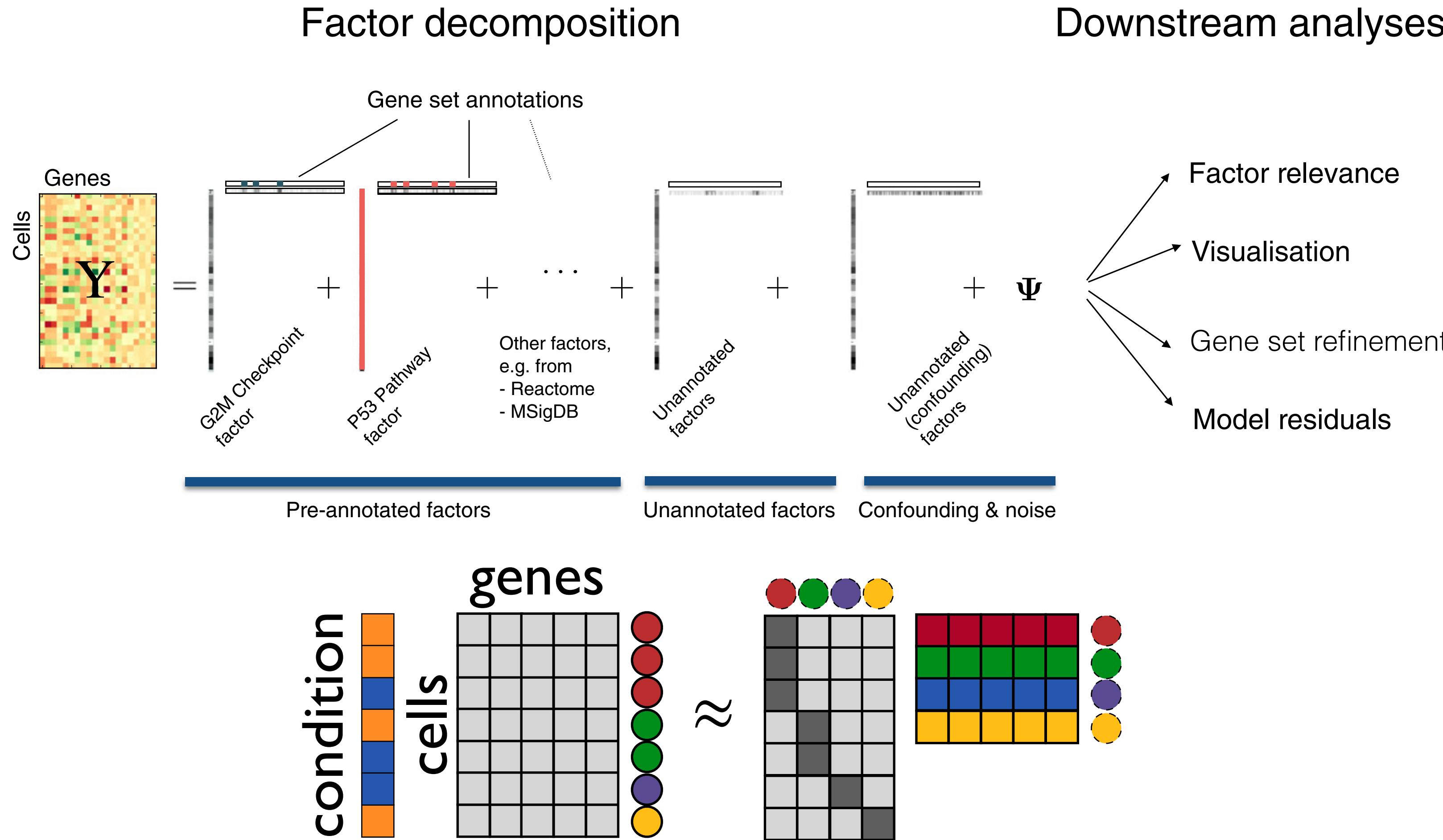
Gene Module and Cell Type Changes

Detection of groups of genes and groups of cells correlated with the perturbation

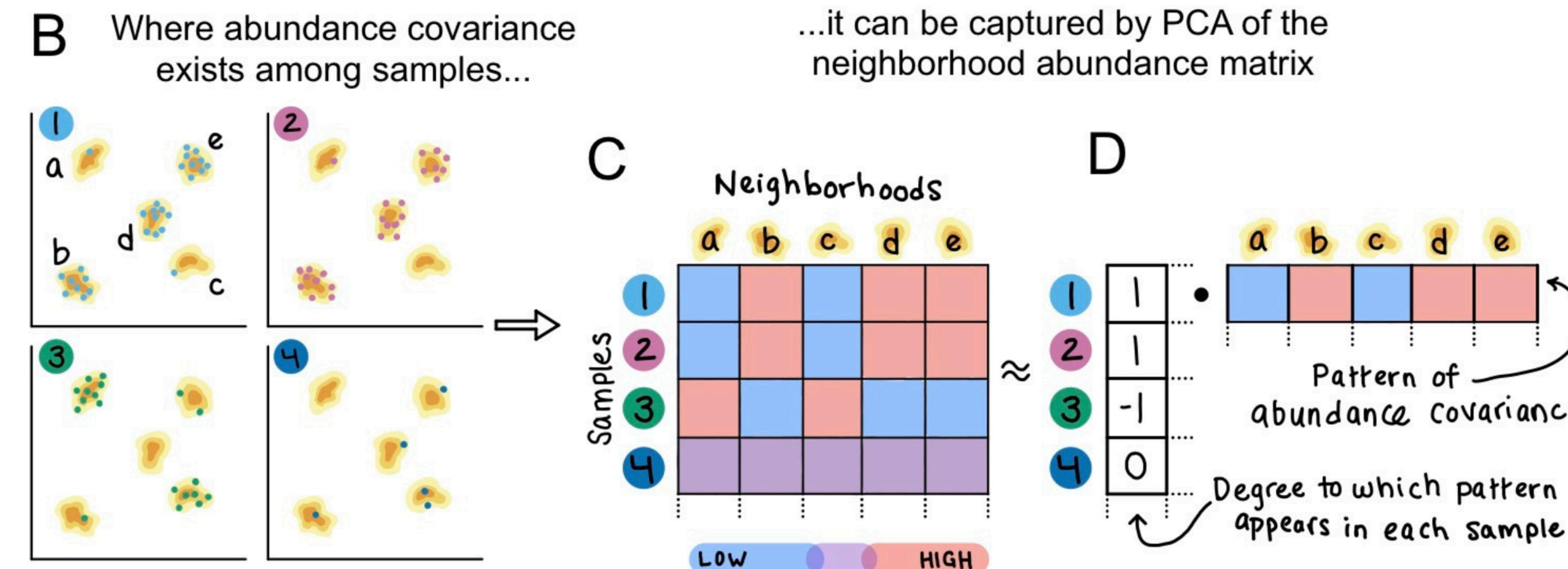
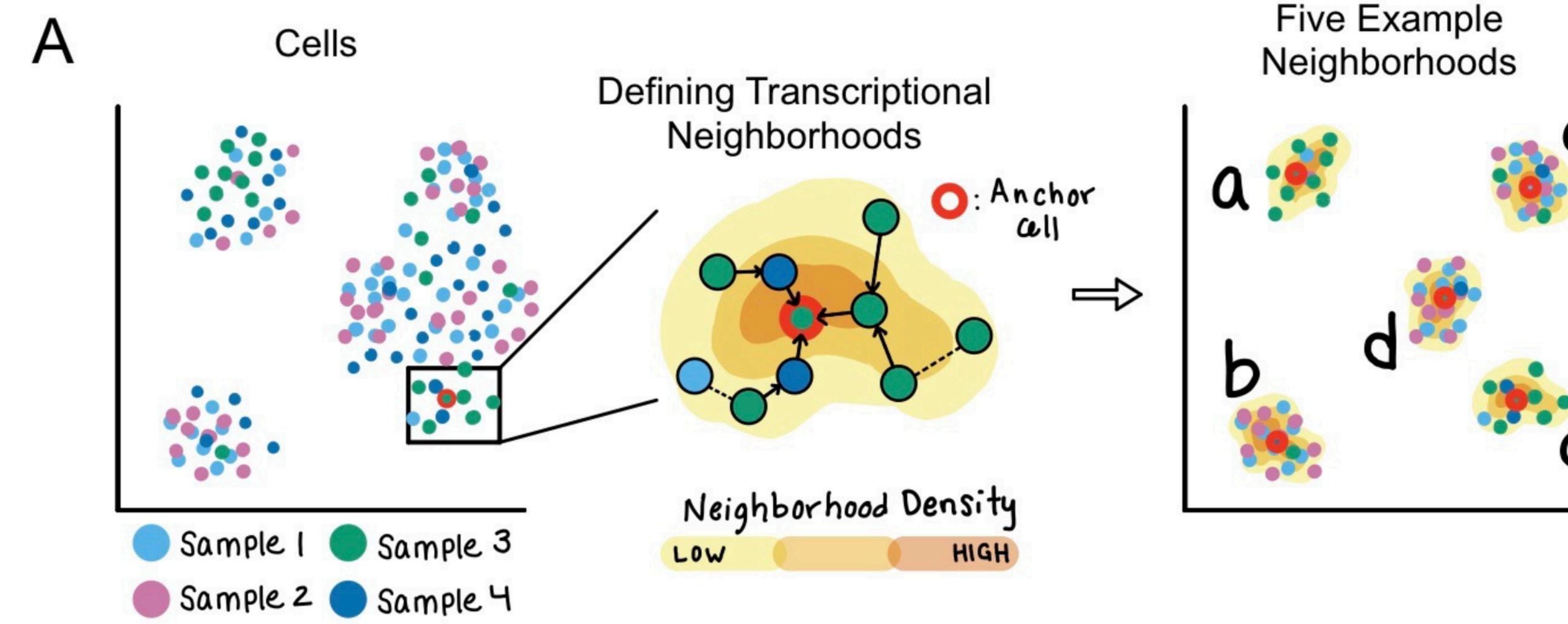


Factor analysis approach

Correlating groups of genes with the perturbation effect



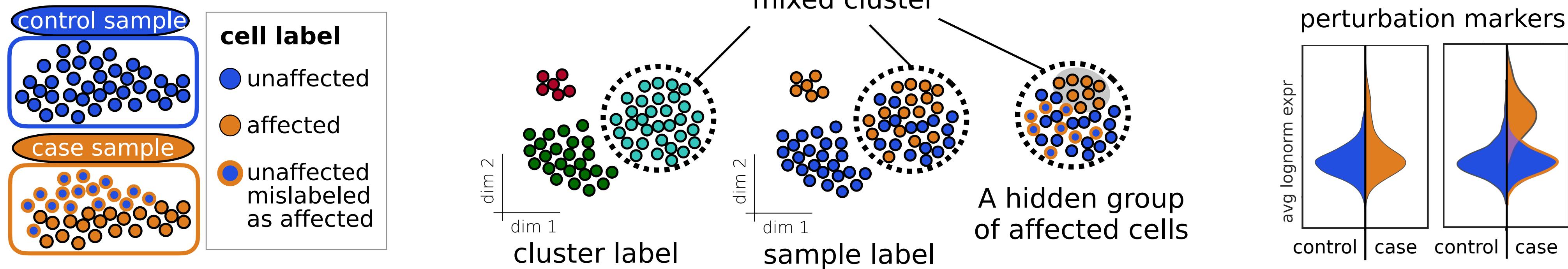
Cluster-free approach for identifying cell states associated with a perturbation



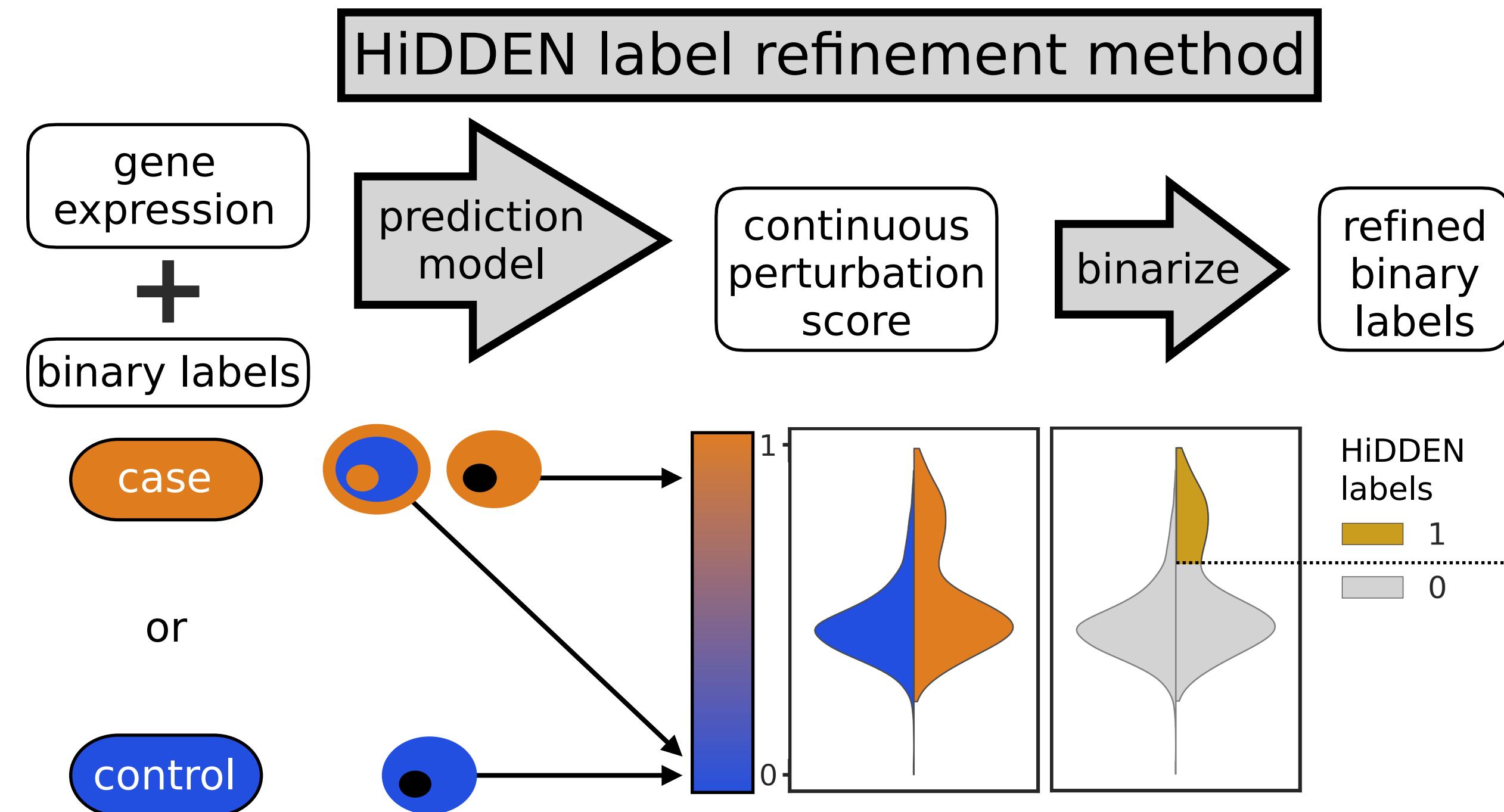
Key assumptions of all prior methods

- The latent space is a complete and accurate representation of variation in the data
- Unsupervised analyses are sufficient to identify groups of cells affected by the perturbation
- Perturbation labels correctly represent the presence or absence of an effect in individual cells

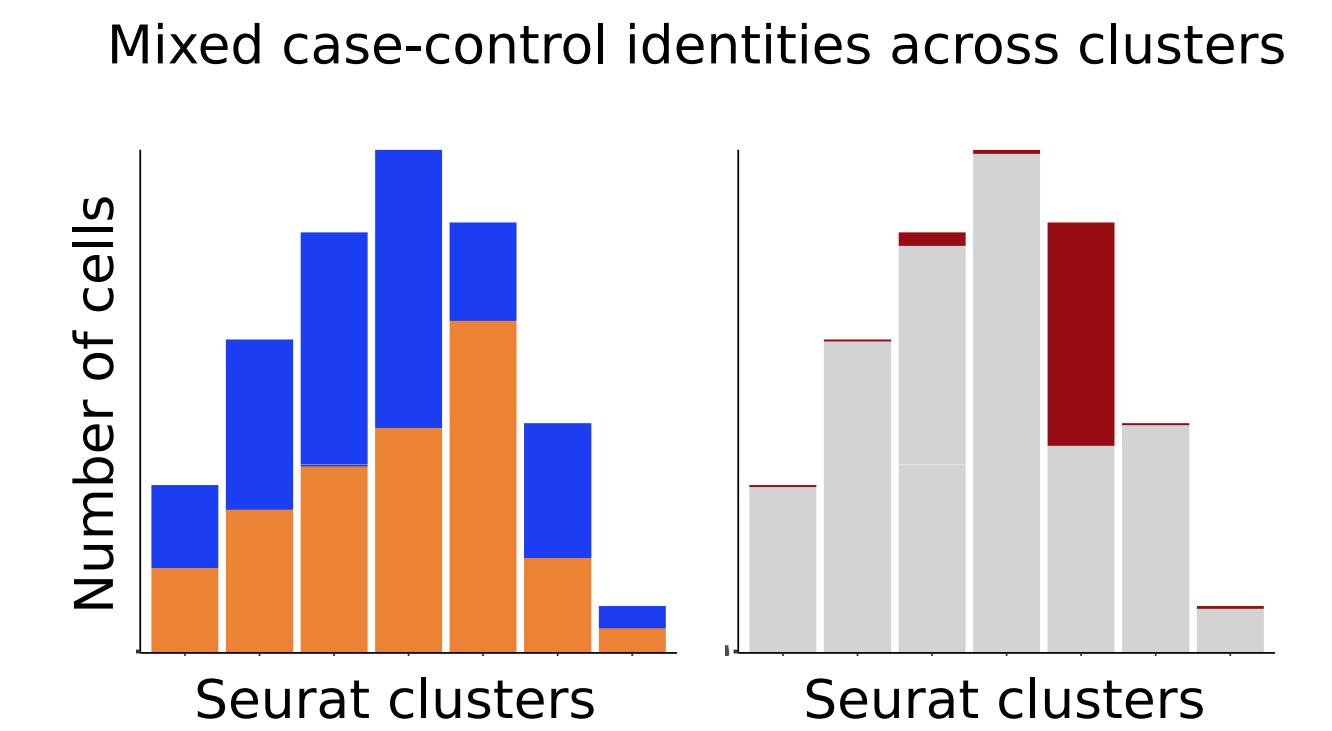
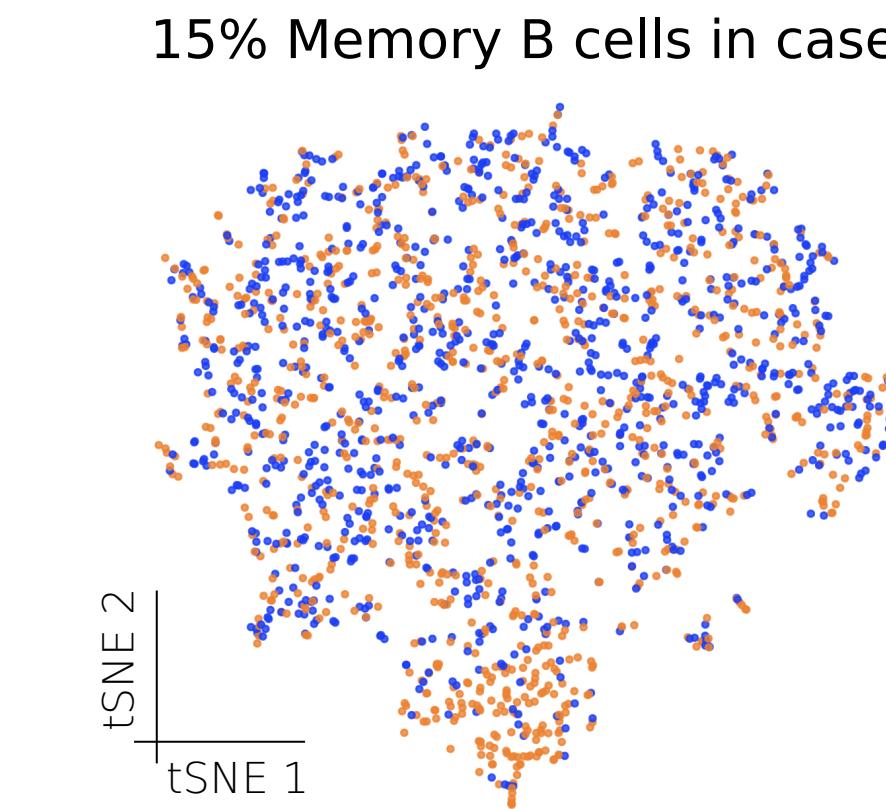
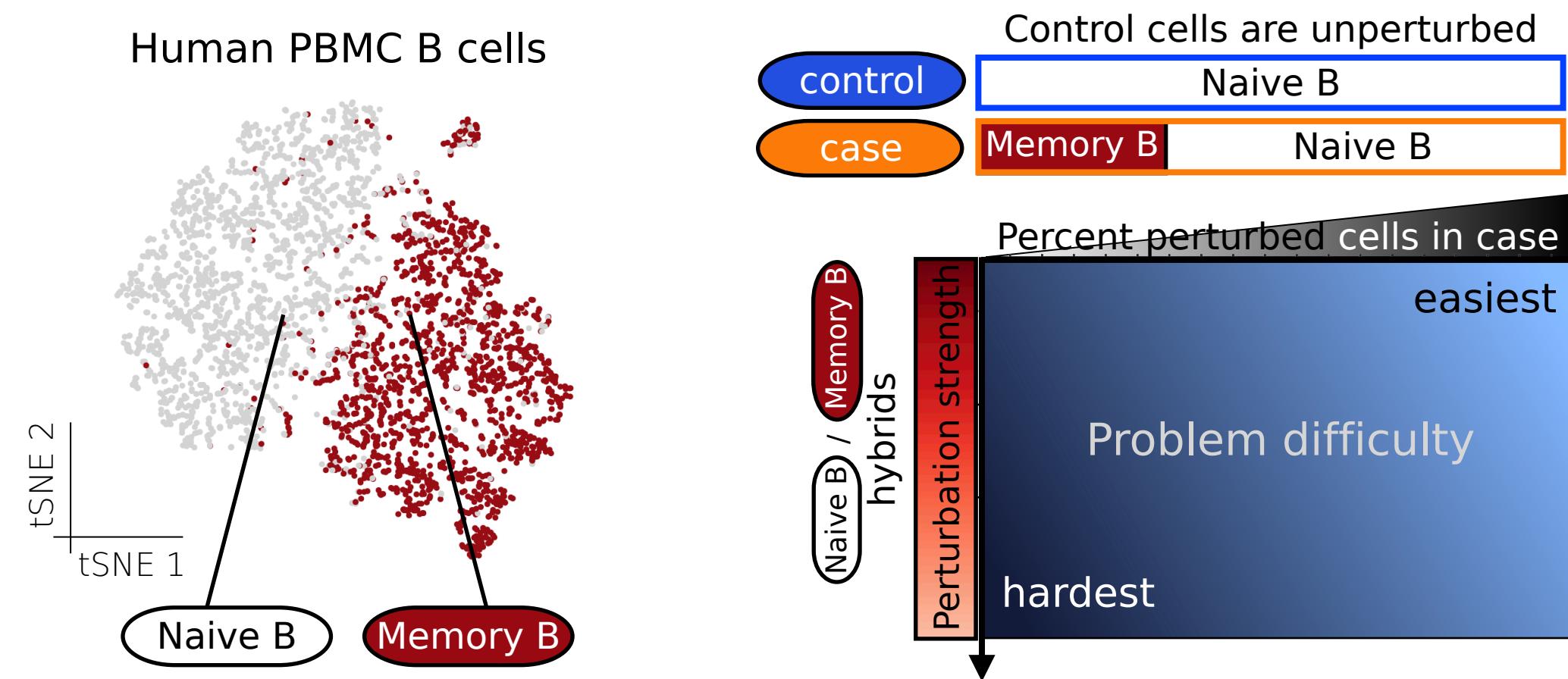
HiDDEN: A machine learning label refinement method for detection of disease-relevant populations in case-control single-cell transcriptomics



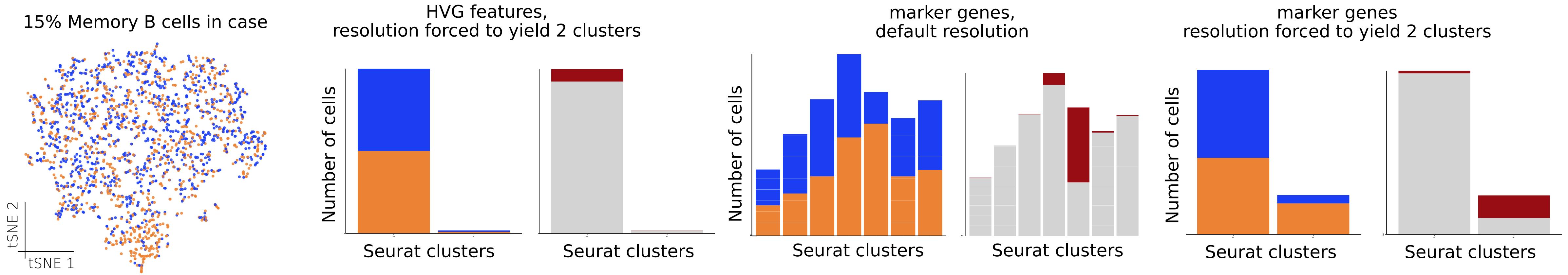
Method overview



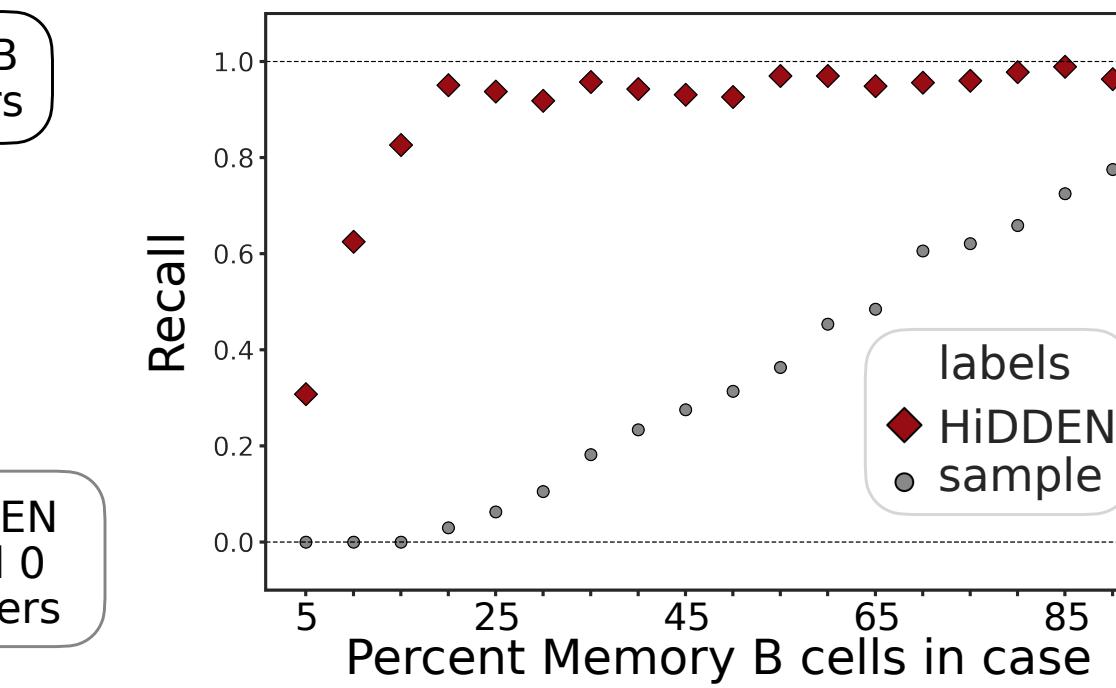
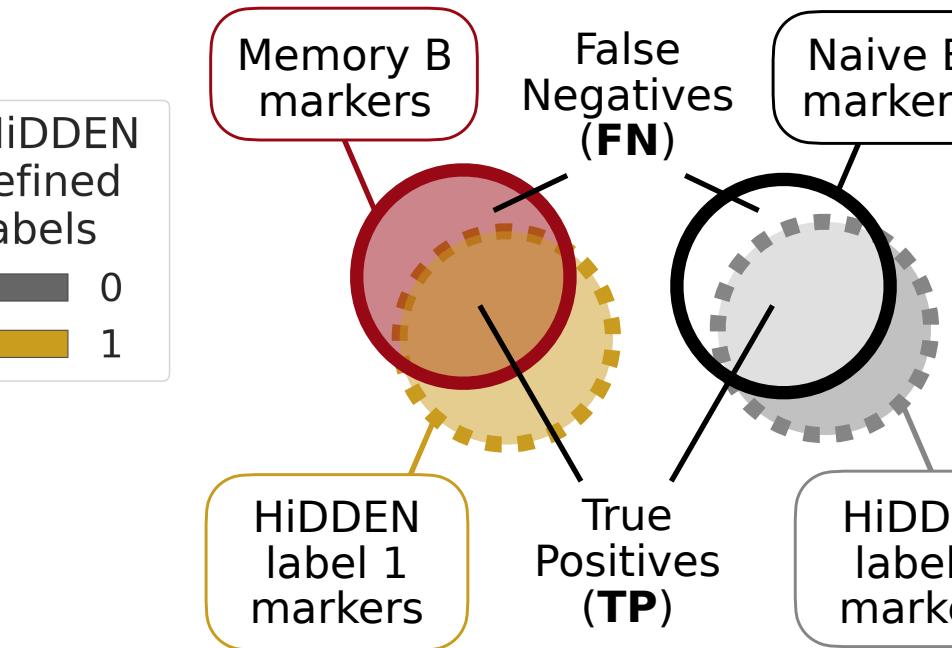
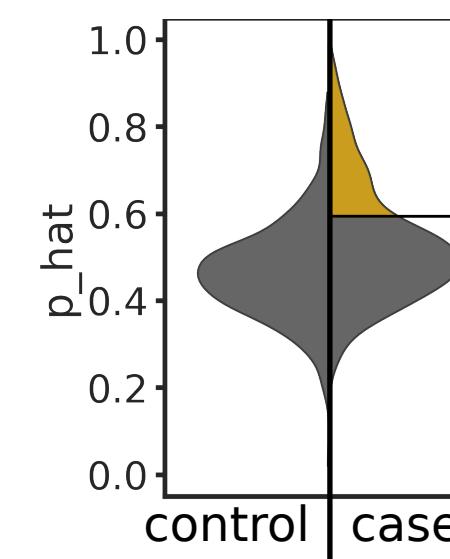
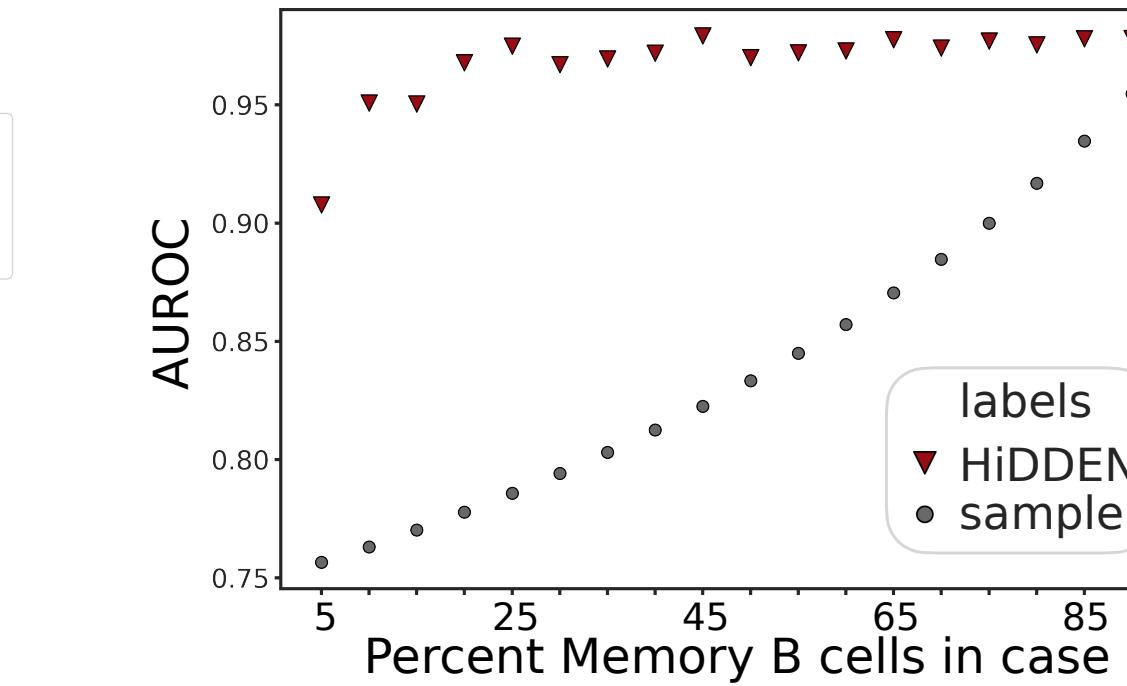
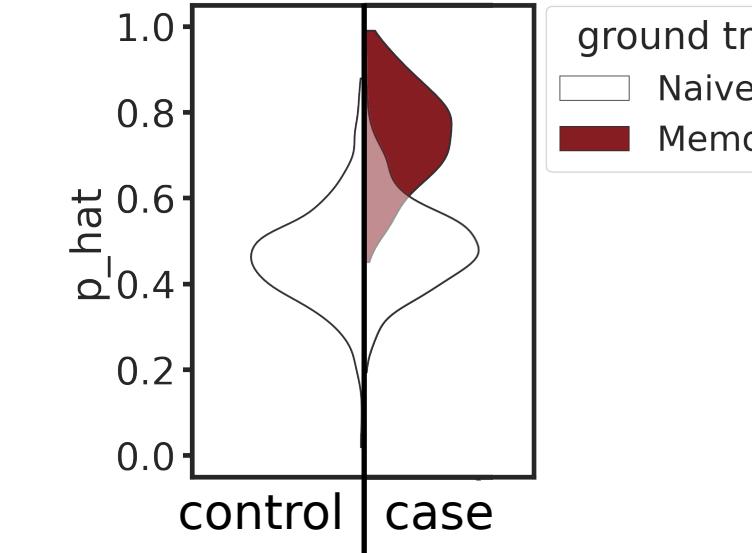
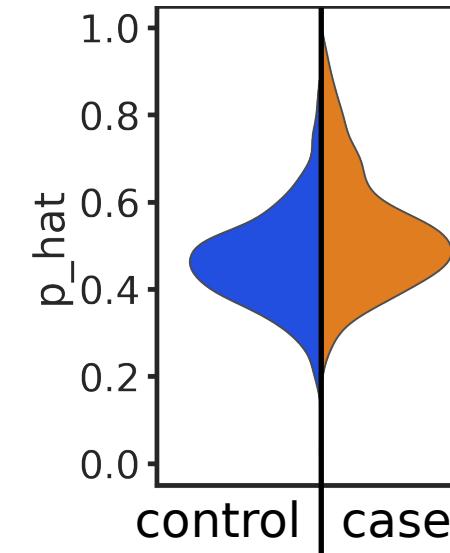
The standard analysis workflow can miss the biological signal in simulated ground truth datasets of cell type mixtures



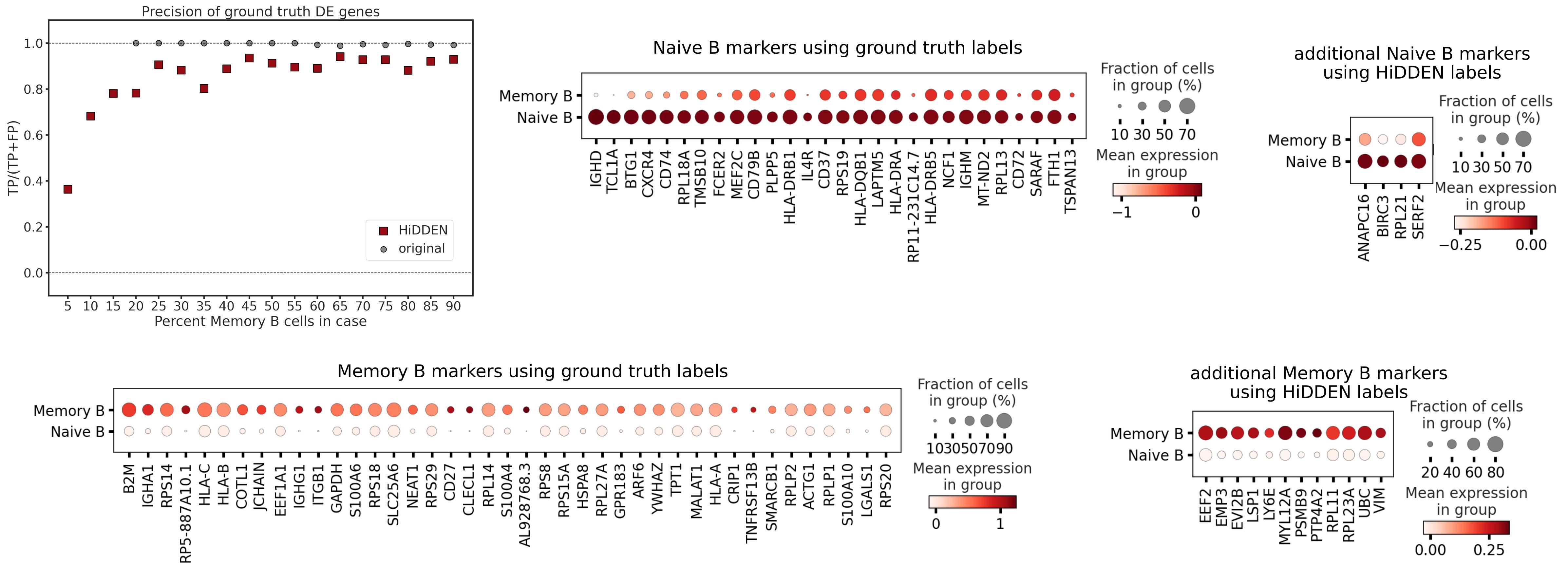
The recovery of the Memory B cluster could not be improved by gene selection or adjusting the resolution parameter alone



HiDDEN detects biological signal missed by the standard analysis workflow

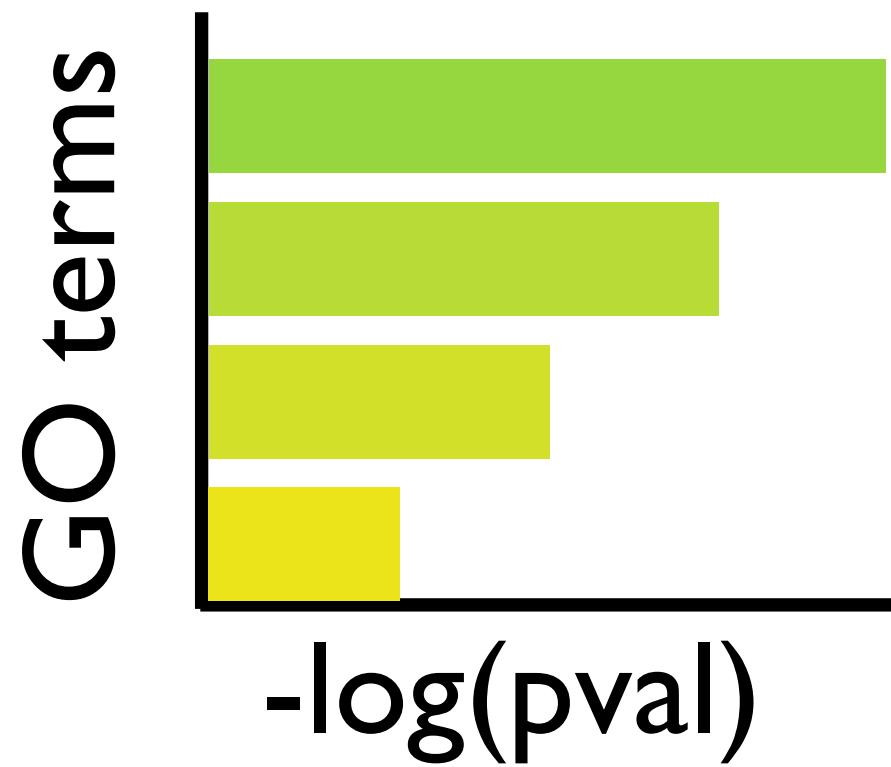


HiDDEN identifies additional marker genes indicating a slight amount of misclassification might have occurred in the original annotation



Functional Interpretation

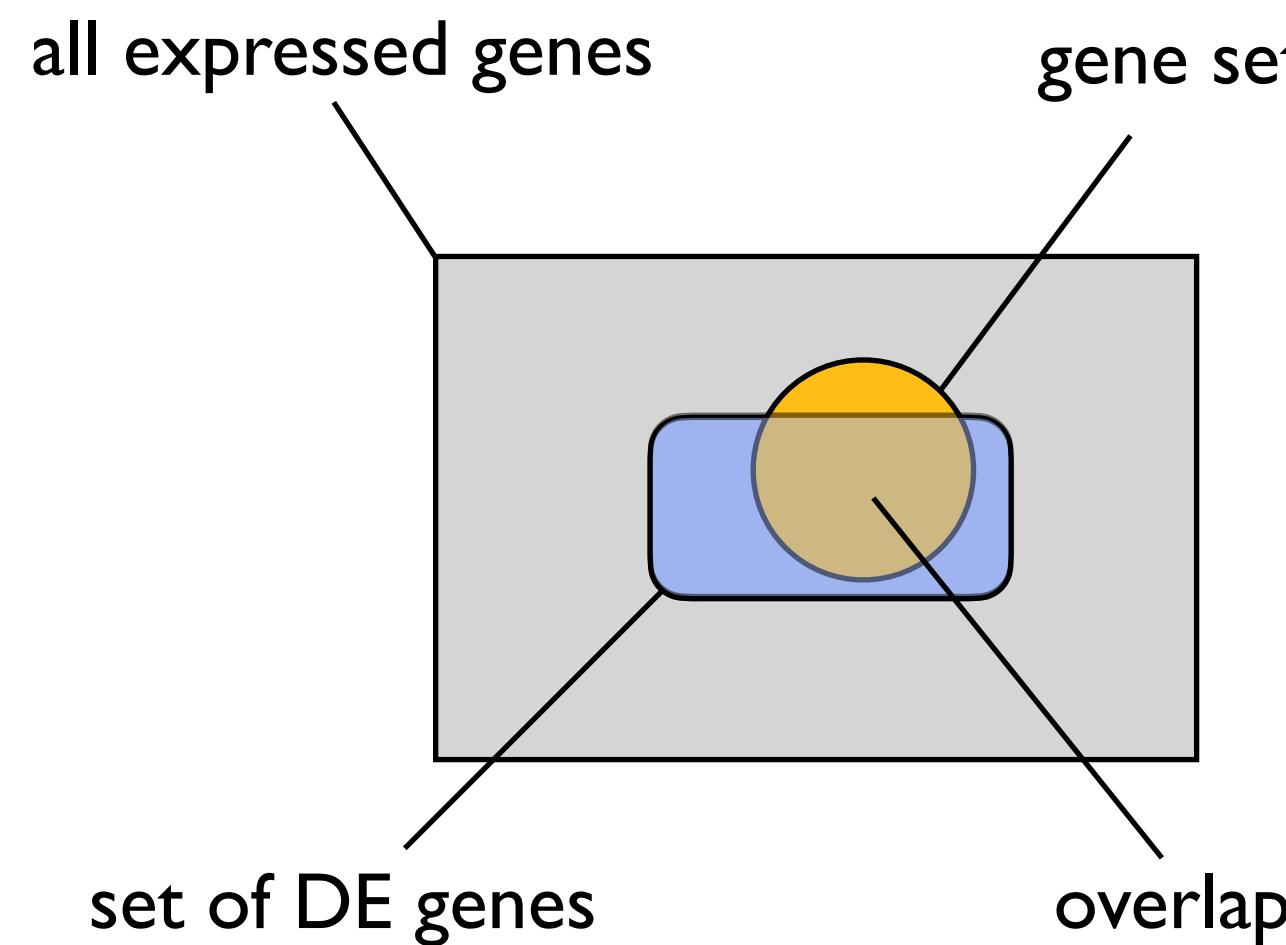
Extracting the biological meaning of genes affected by the perturbation



Differential expression analyses can return many genes

Endow genes affected by the perturbation with a biological interpretation

Overrepresentation test
set of DE genes

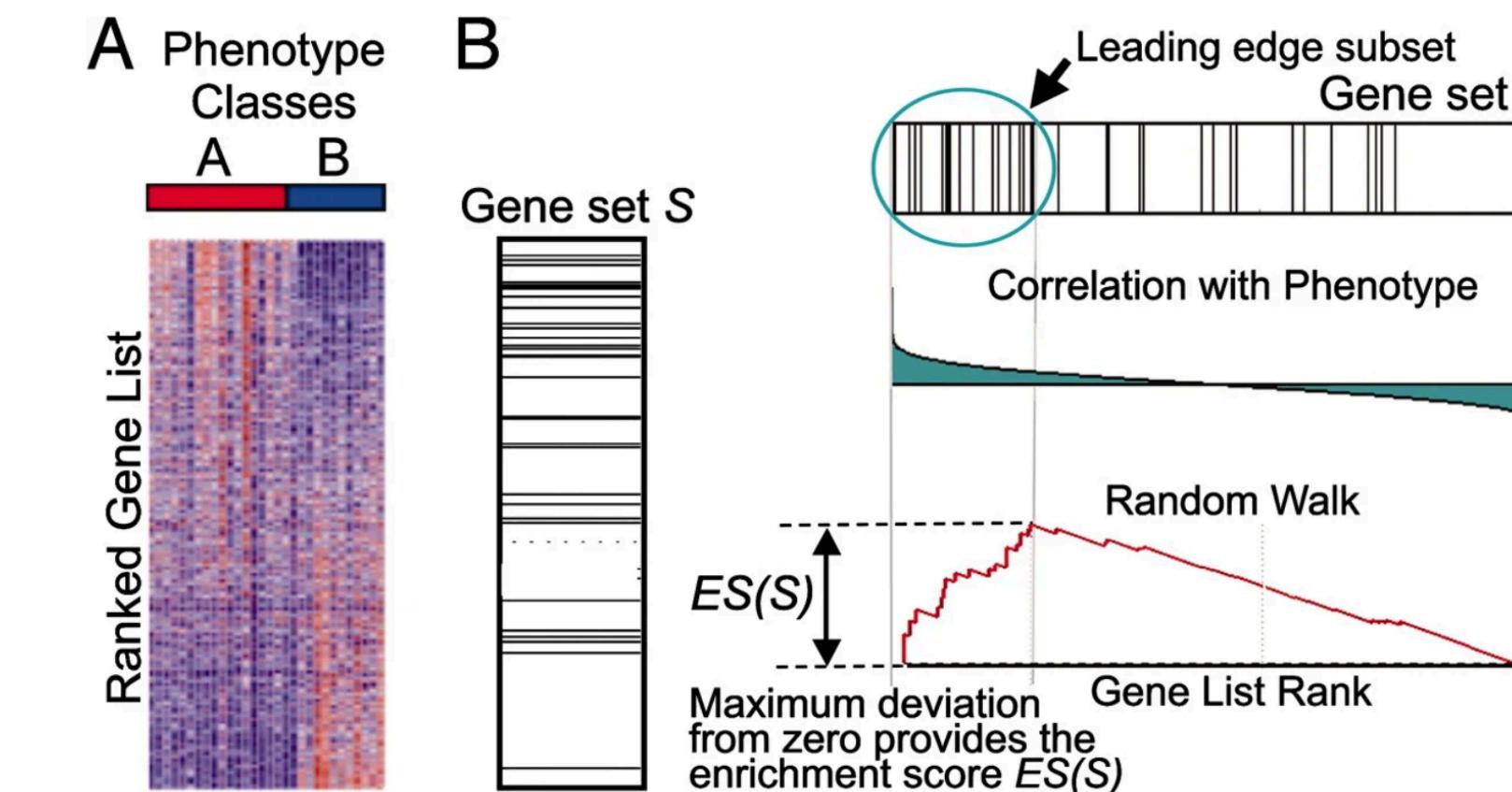


GSEA
ranked list of all genes

RESEARCH ARTICLE | BIOLOGICAL SCIENCES | ✓
Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles

Aravind Subramanian, Pablo Tamayo, Vamsi K. Mootha, Sayan Mukherjee, Benjamin L. Ebert, Michael A. Gillette, Amanda Paulovich, Scott L. Pomeroy, Todd R. Golub, Eric S. Lander, and Jill P. Mesirov [7] [Authors Info & Affiliations](#)

September 30, 2005 | 102 (43) 15545-15550 | <https://doi.org/10.1073/pnas.0506580102>



Gene set enrichment analyses often return many terms

Revigo can summarize them by removing redundant GO terms

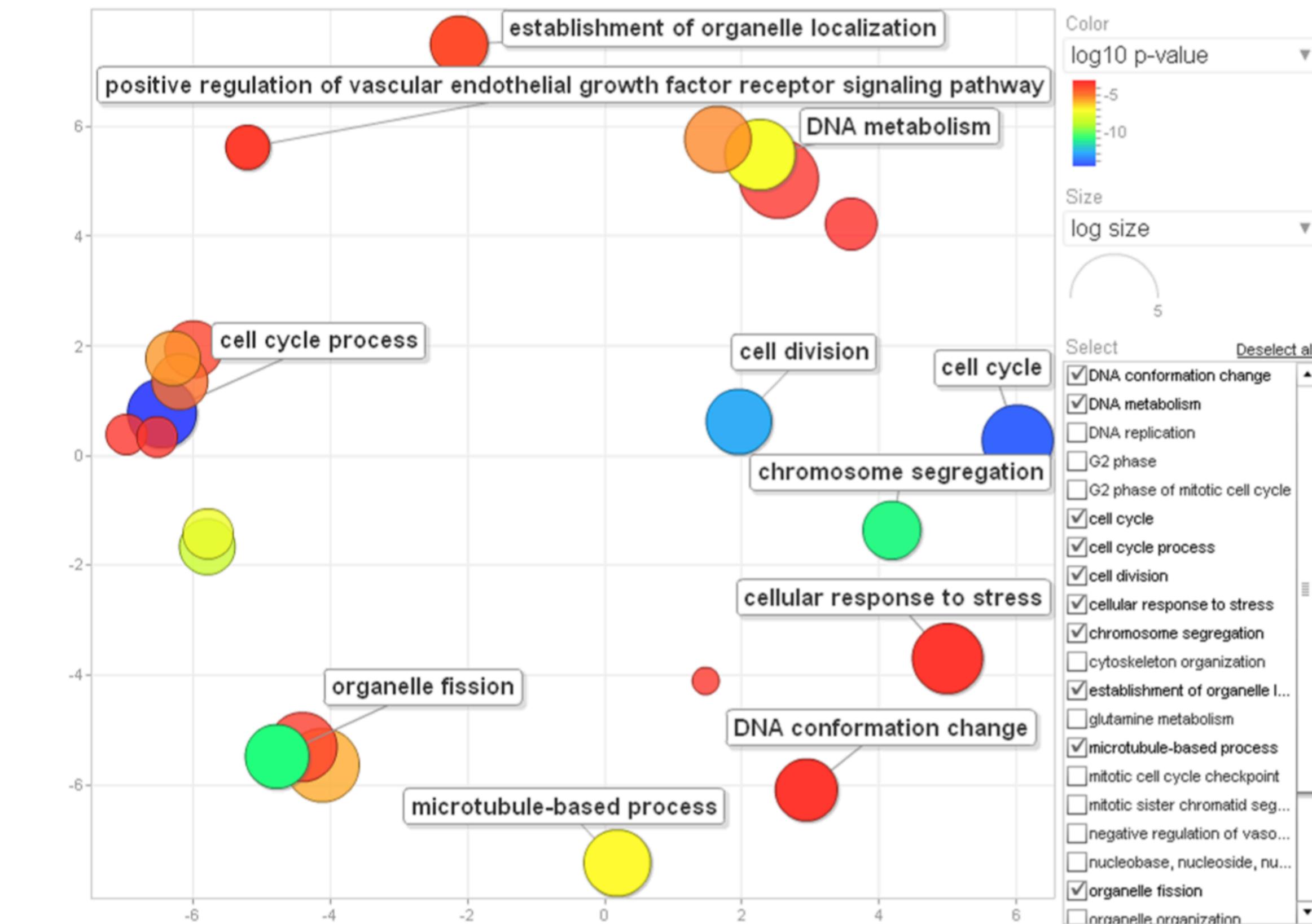


reduce + visualize Gene Ontology

REVIGO Summarizes and Visualizes Long Lists of Gene Ontology Terms

Fran Supek Matko Bošnjak, Nives Škunca, Tomislav Šmuc

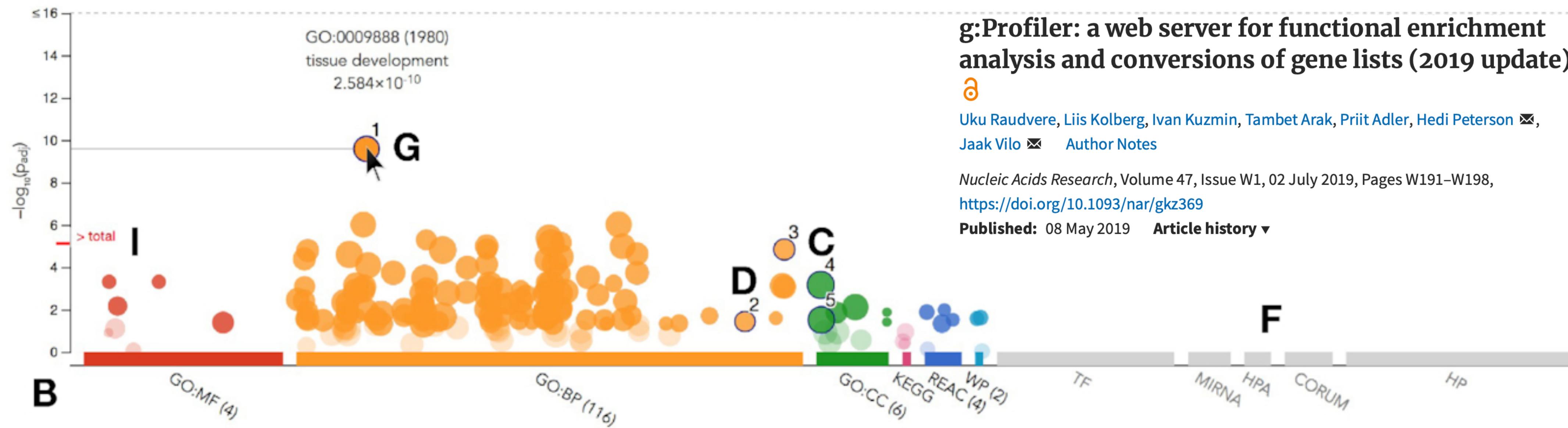
Published: July 18, 2011 • <https://doi.org/10.1371/journal.pone.0021800>



Gene set enrichment

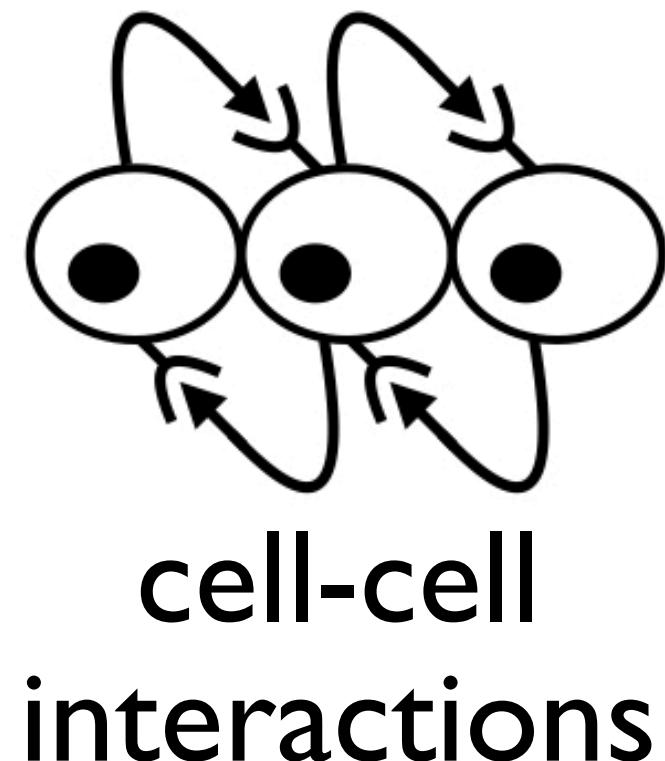
Summary

- The results can be more sensitive to the choice of gene set database rather than statistical methods



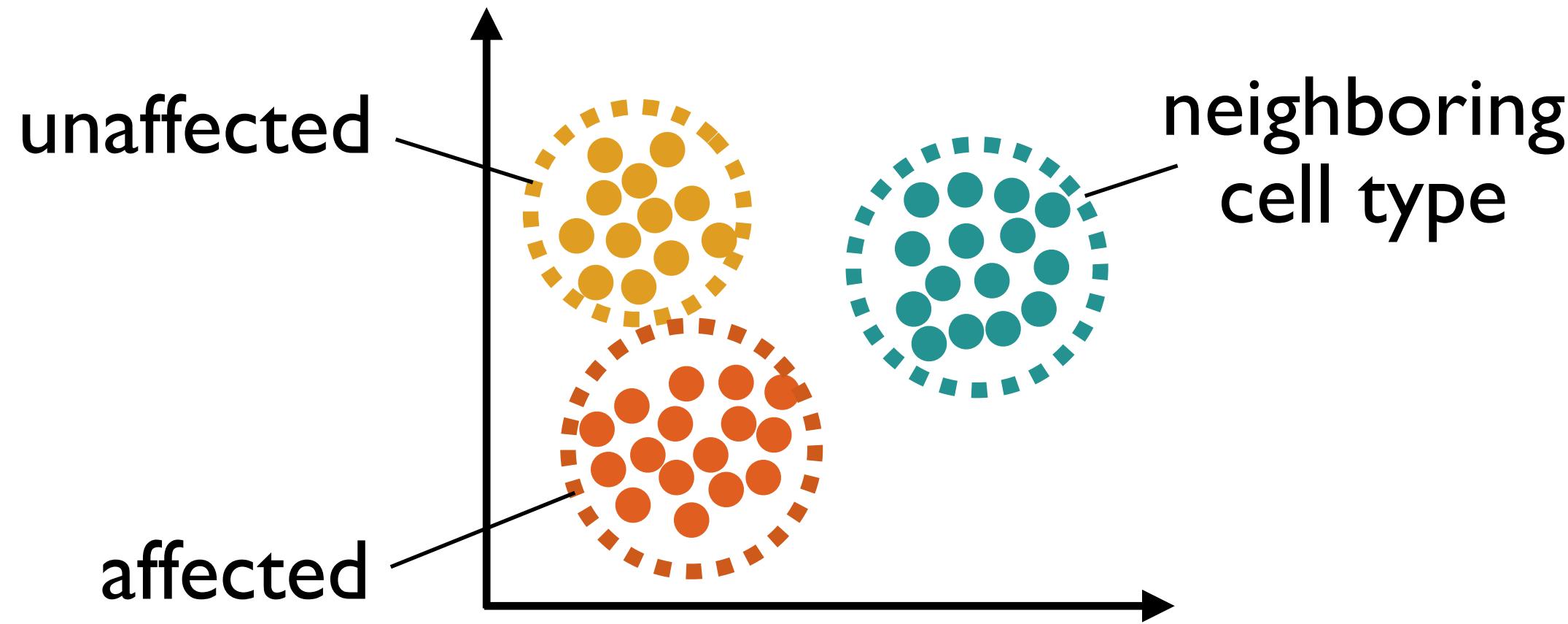
Functional Interpretation

Testing for perturbation-associated cell-cell communication changes

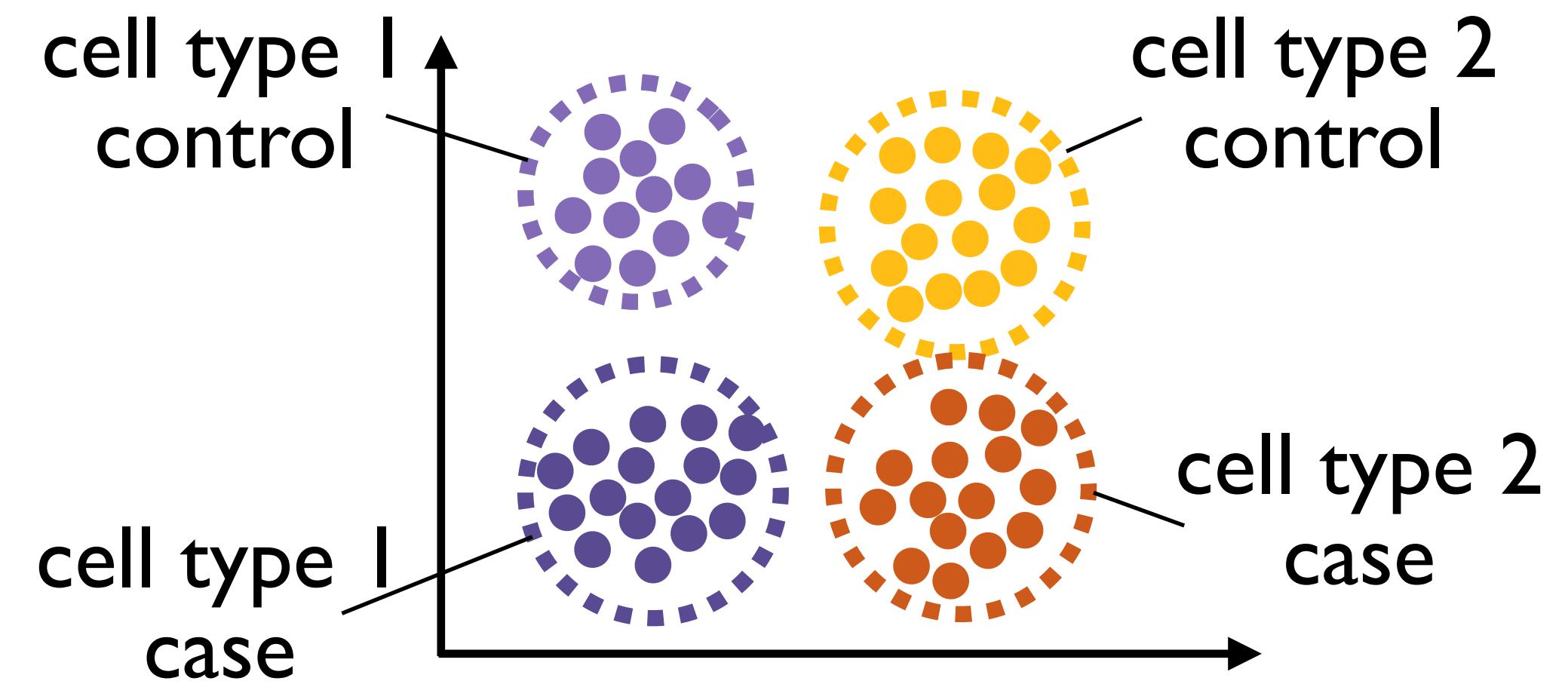


Contrast cell-cell interactions

of affected and unaffected cells
with a neighboring cell type

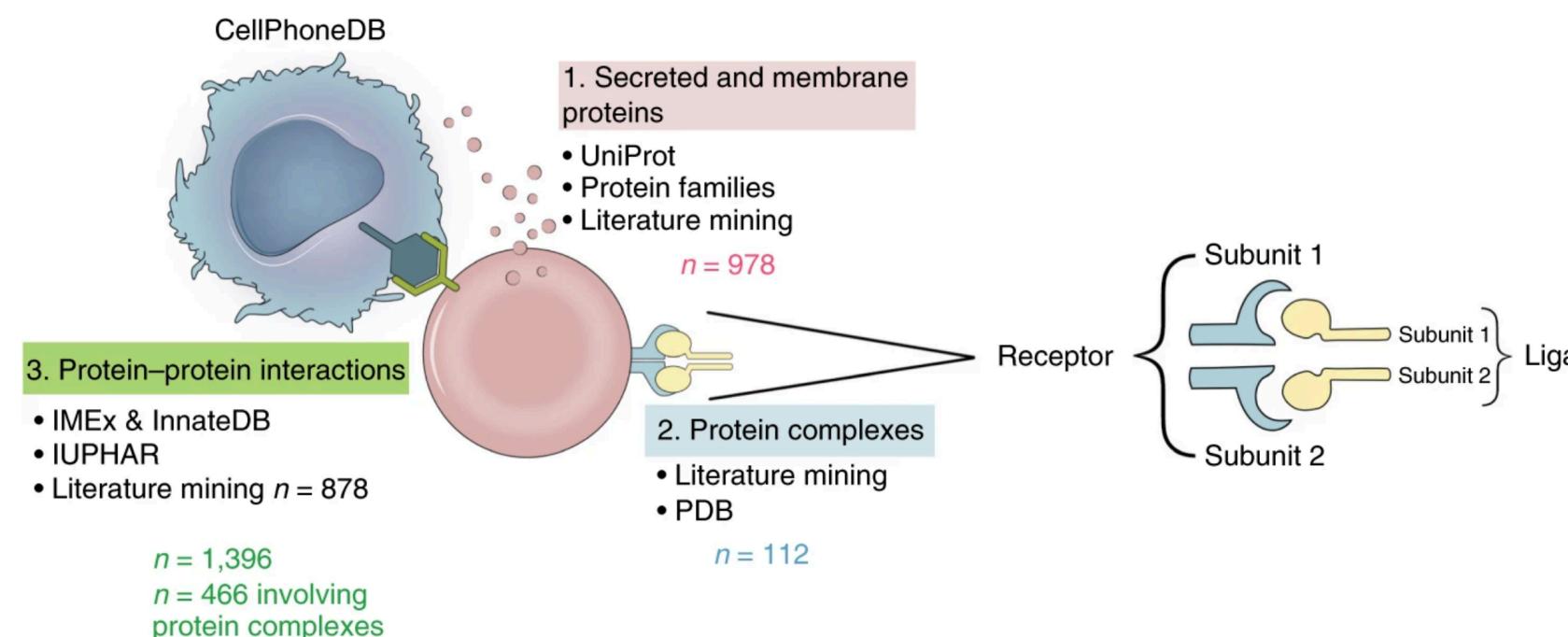


between two cell types
across conditions

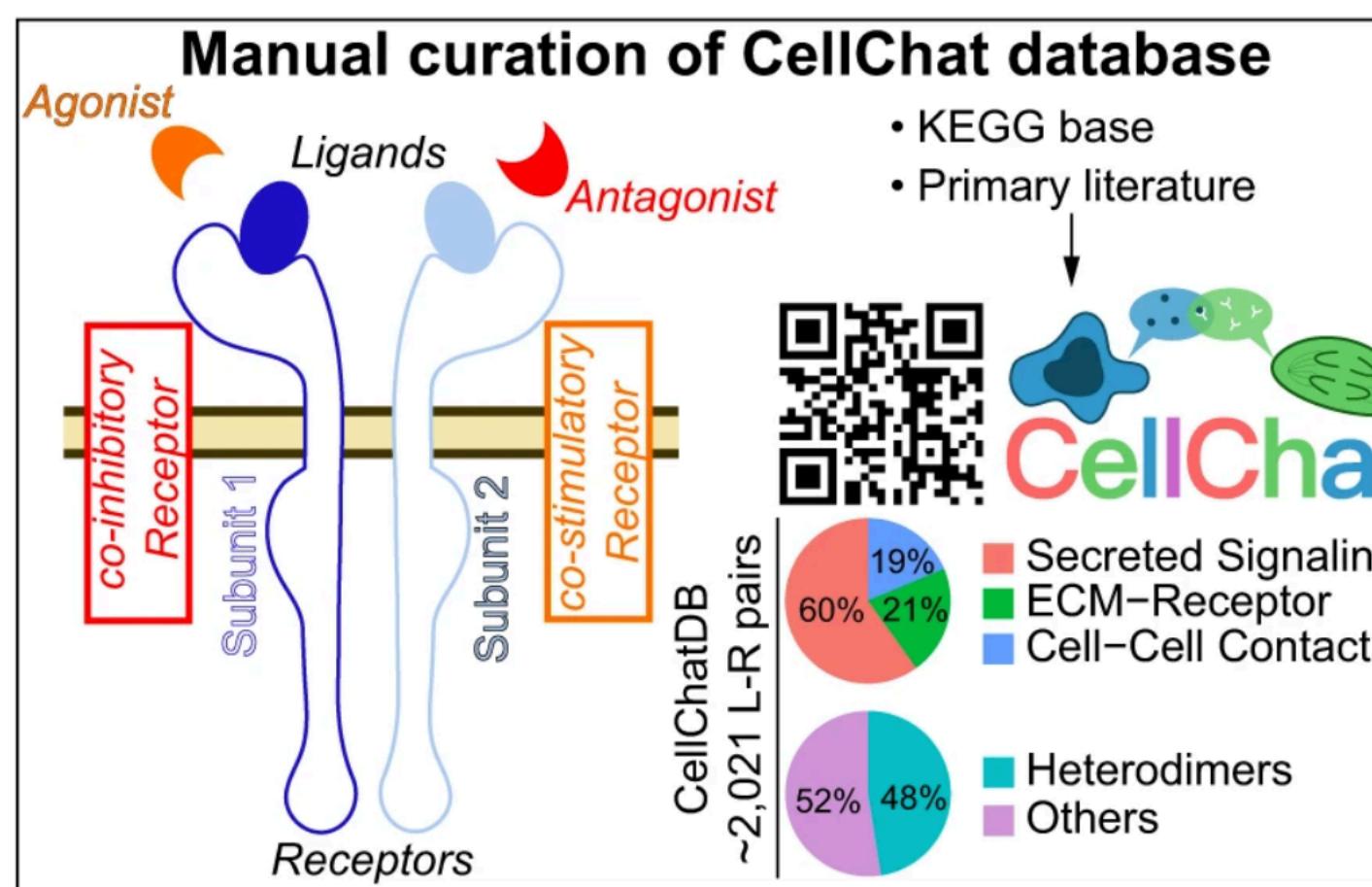


Ligand-receptor interaction repositories

CellphoneDB, CellChat heteromeric complexes

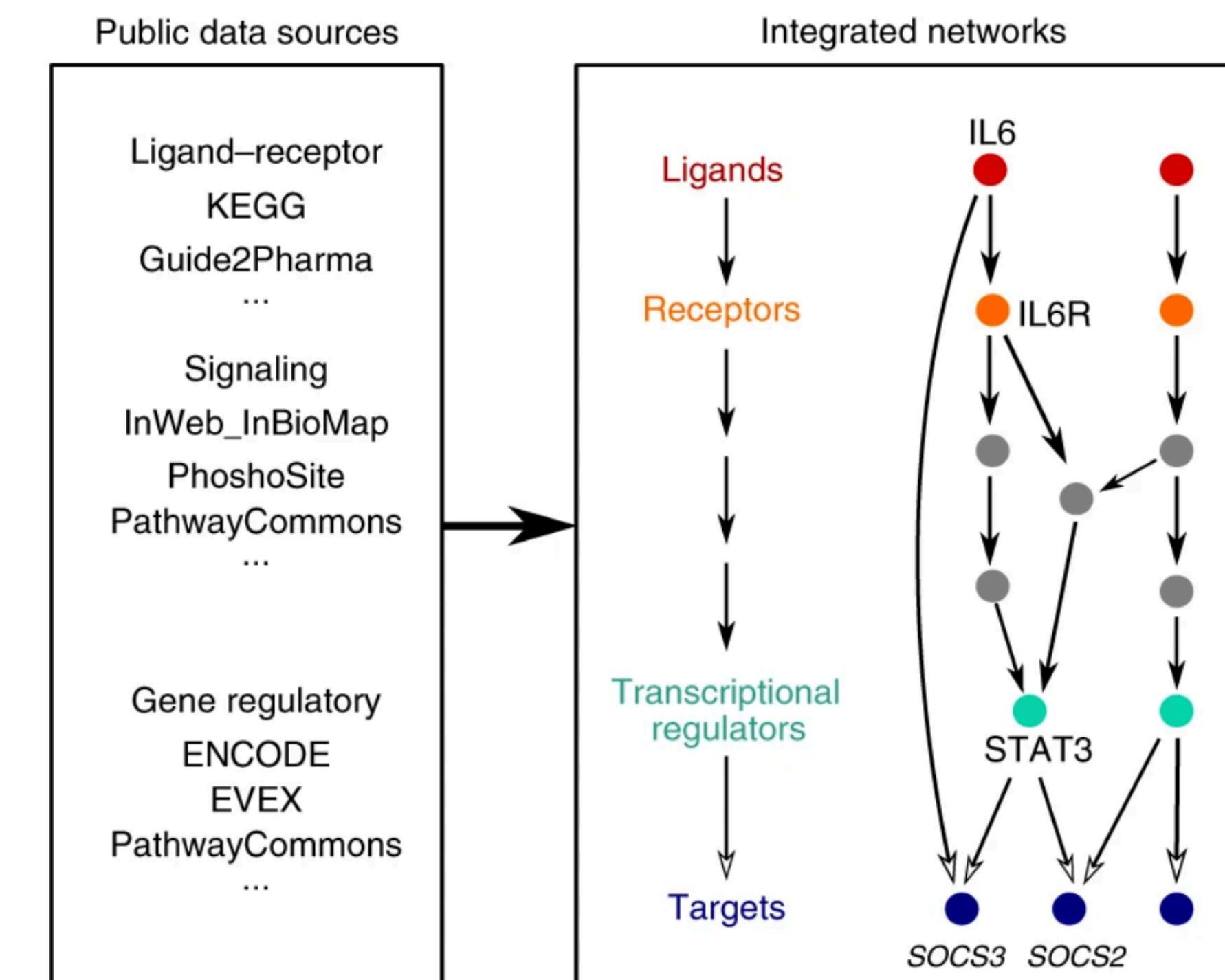


[Efremova, M., et al. Nature Protocols \(2020\)](#)



[Jin, S., et al. Nature Communications \(2021\)](#)

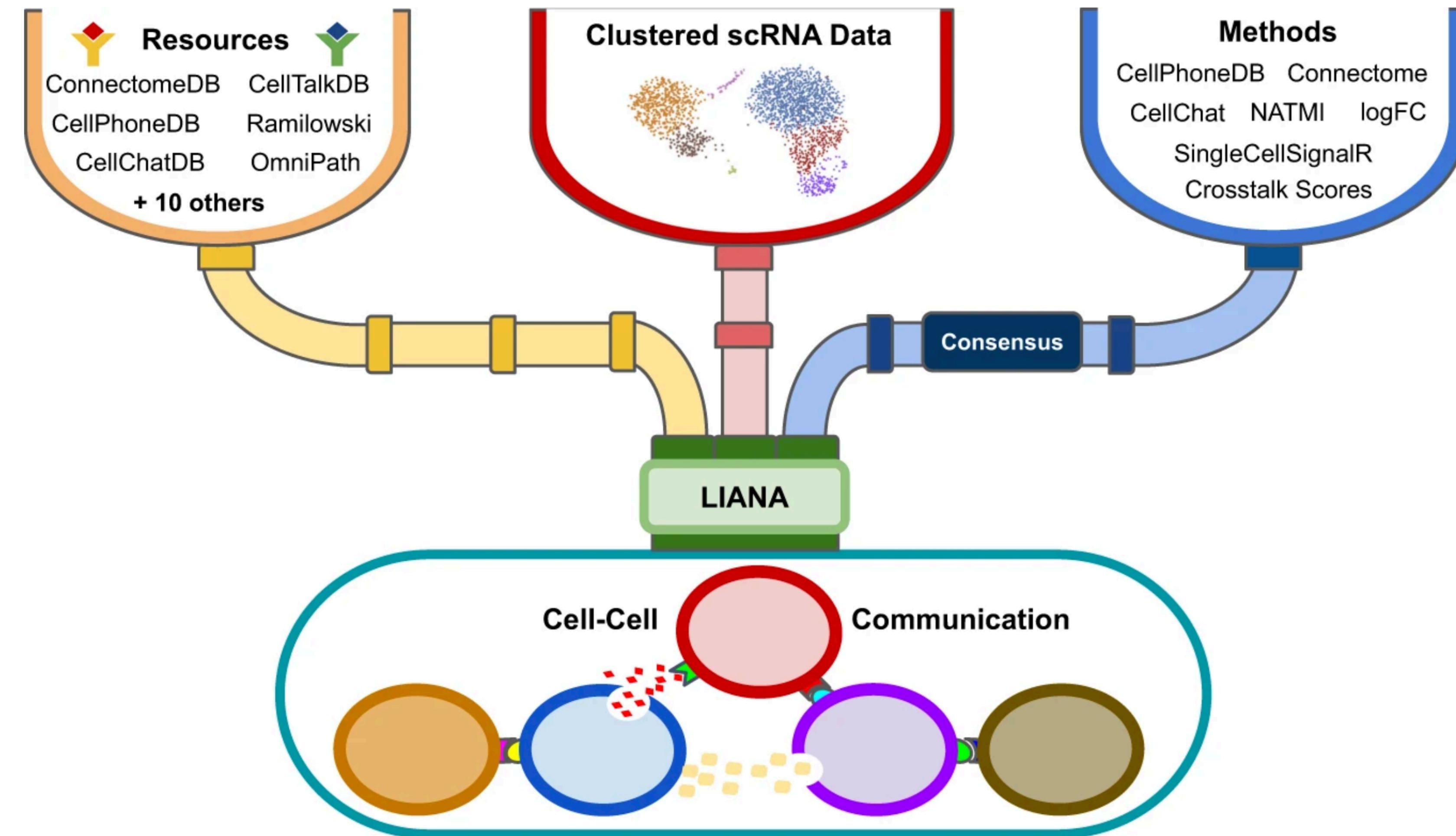
Nichenet linking ligands to target genes



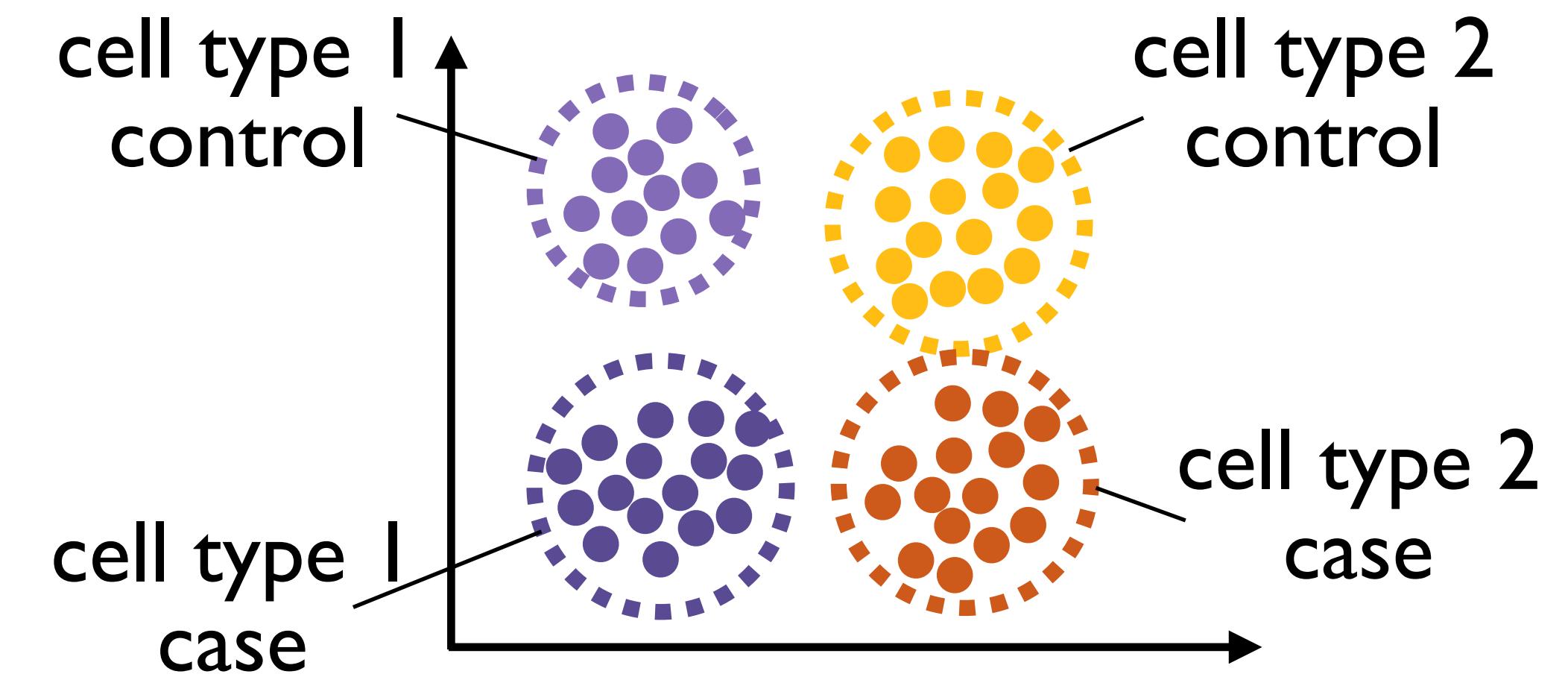
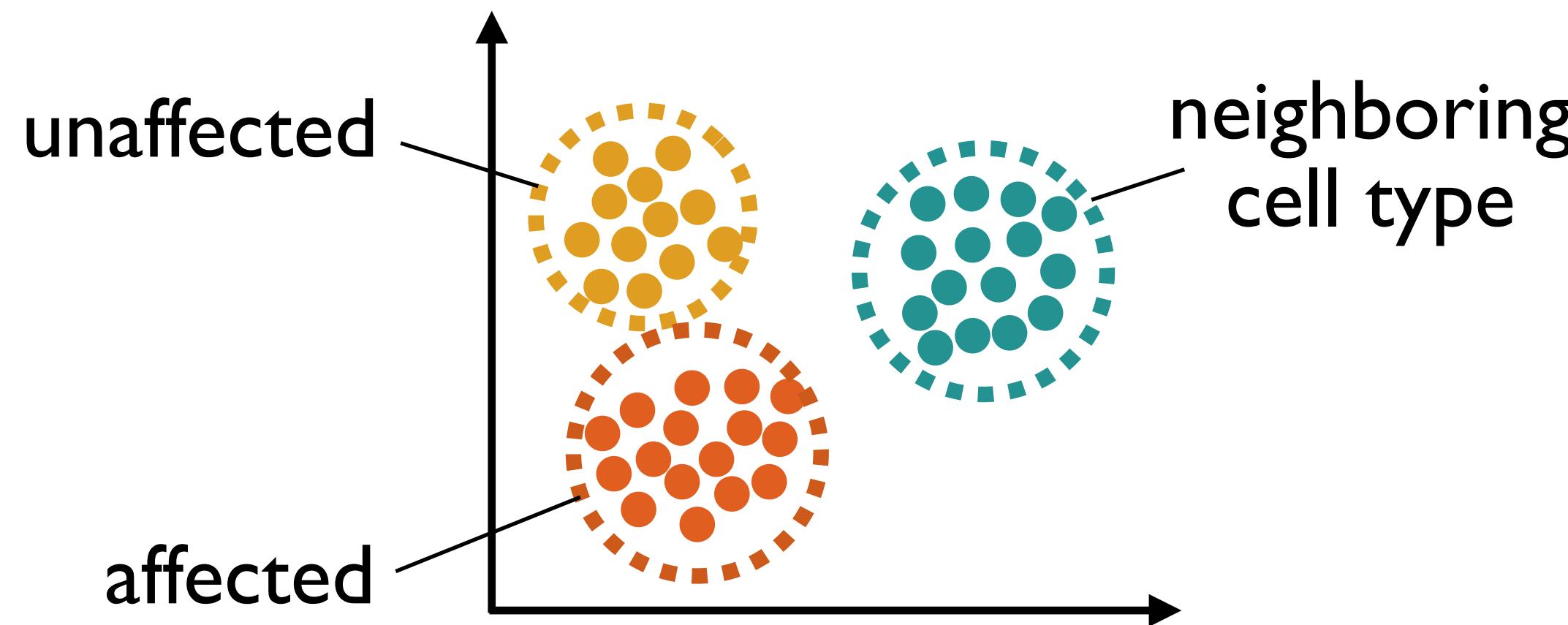
[Browaeys, R., et al. Nature Methods \(2020\)](#)

Lack of consensus between repositories and computational tools

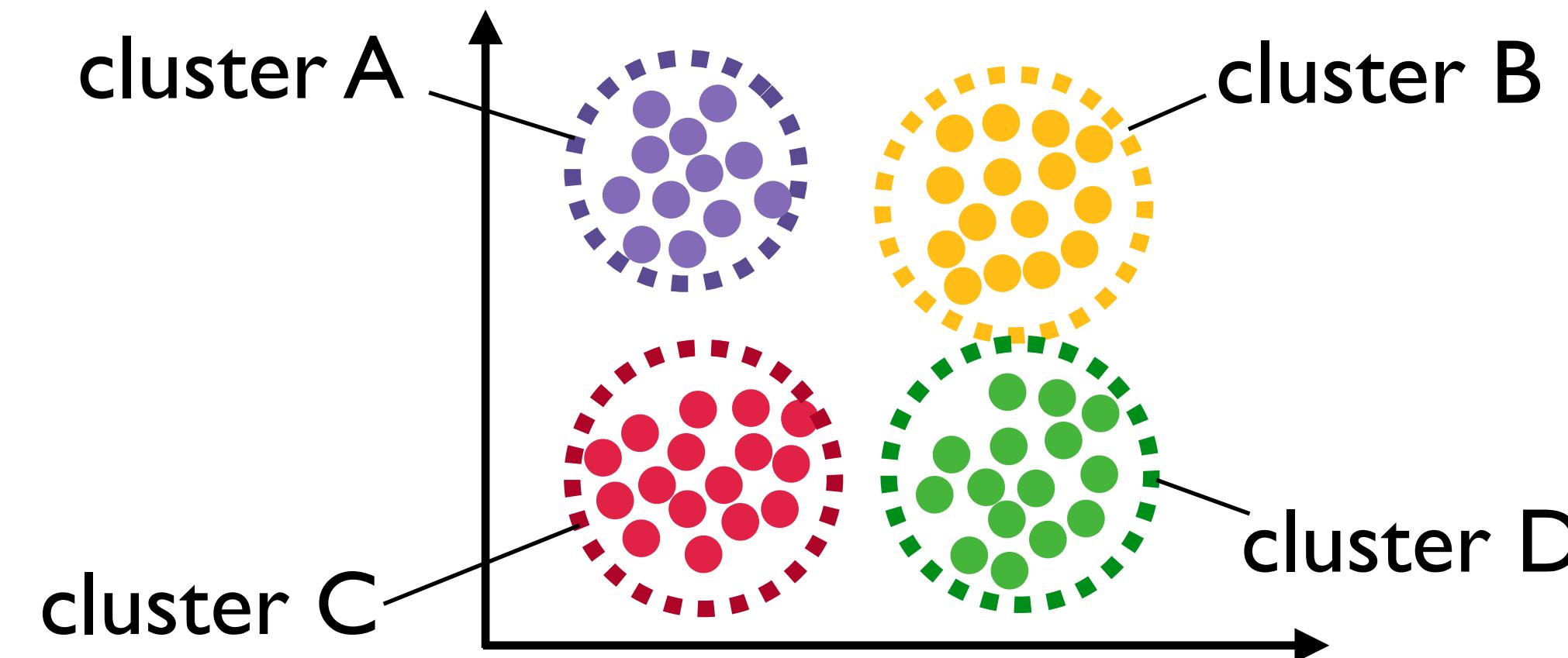
LIANA allows all combinations of resources and methods



Cell-cell communication methods are not built with our targeted contrasts in mind

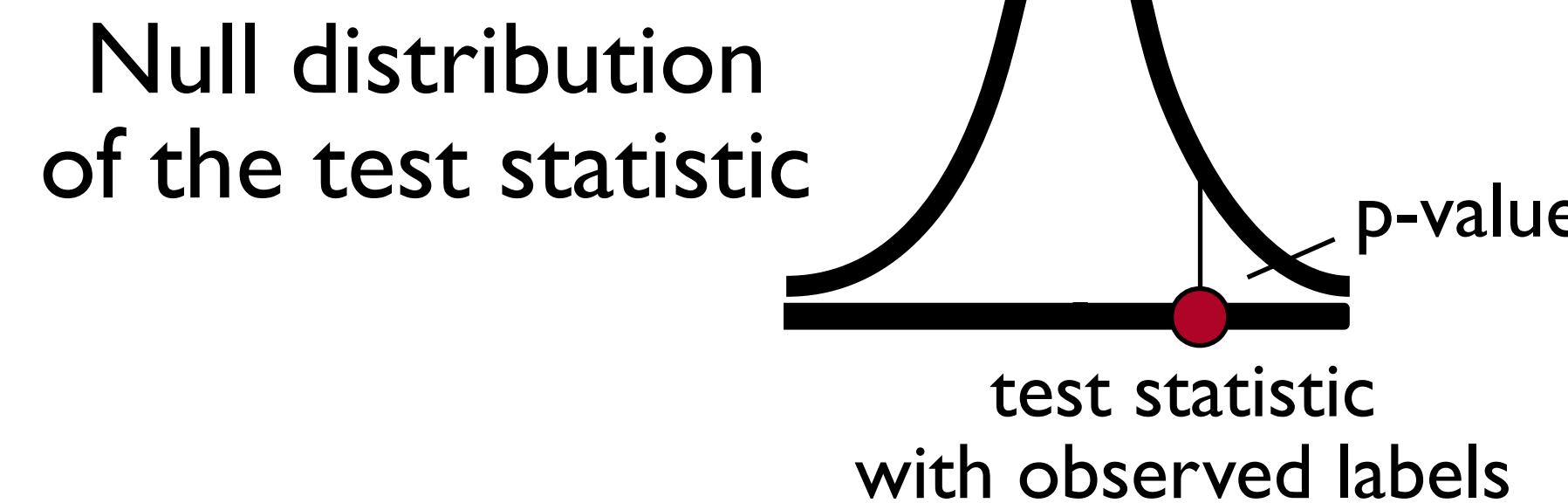


Take CellphoneDB as an example

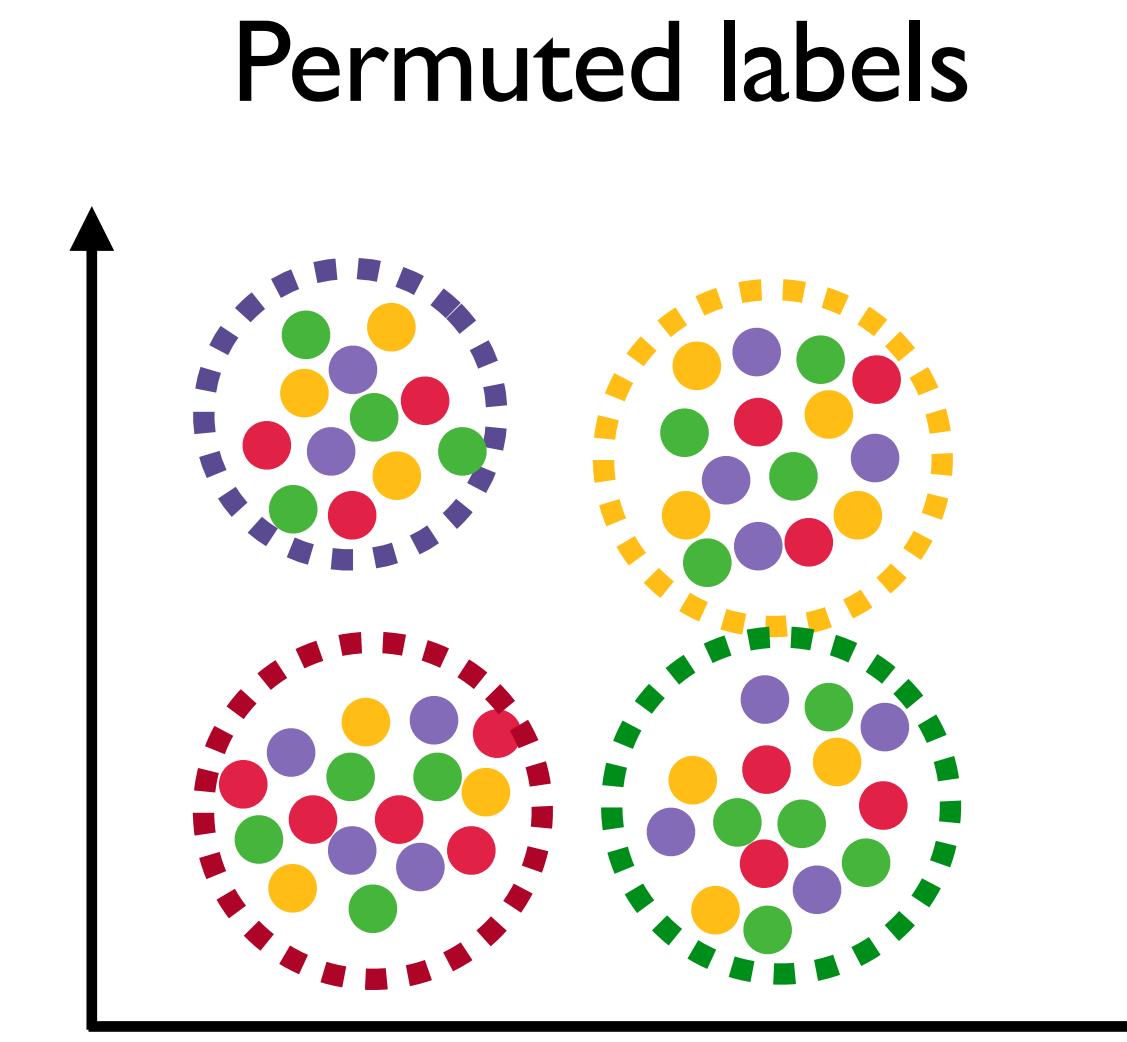
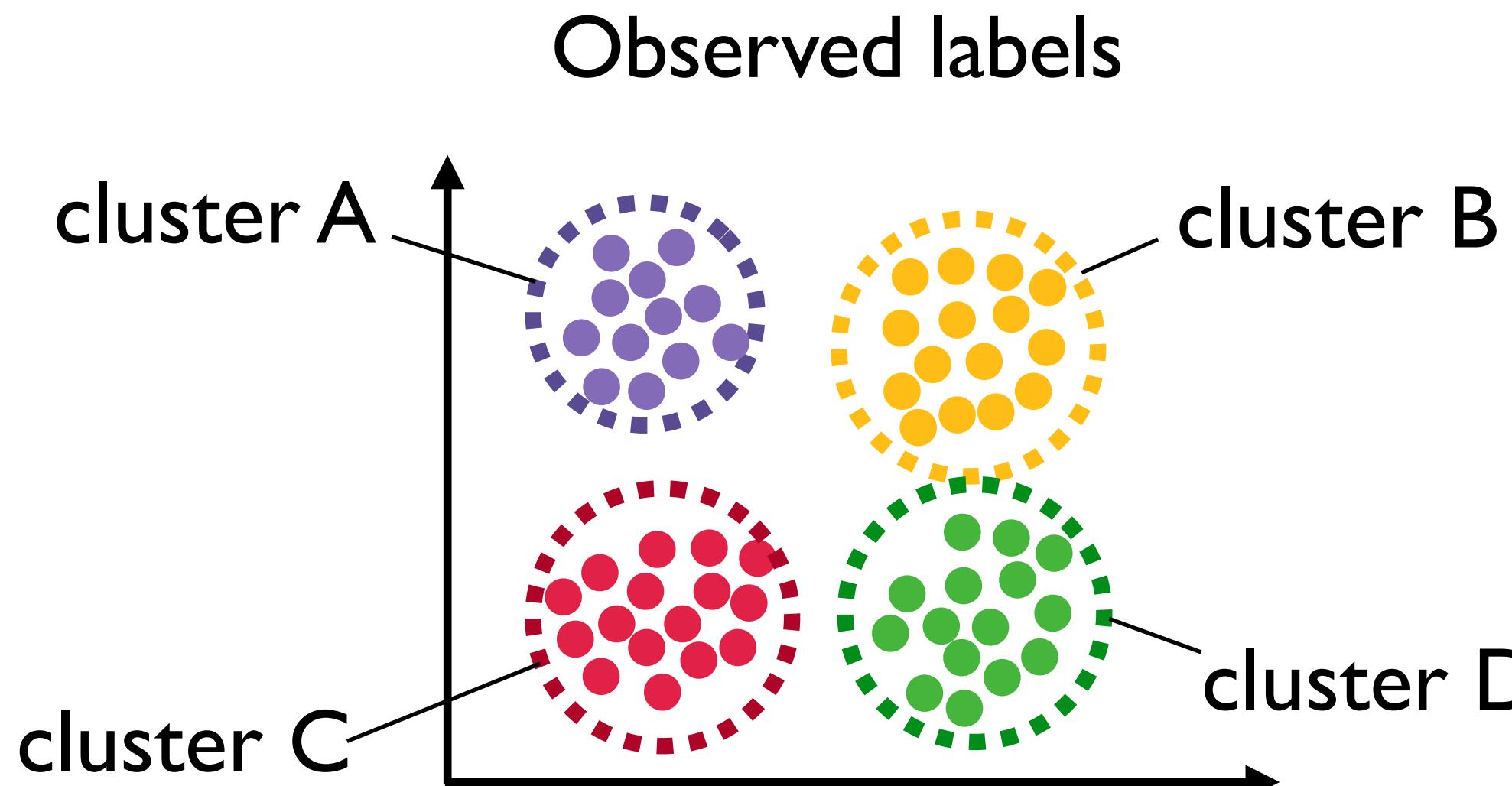


For any given ordered ligand-receptor pair (L, R) and an ordered pair of clusters (clA, clB) :

$$\text{Test Statistic} = \text{mean}(\bar{L}_{clA}, \bar{R}_{clB})$$

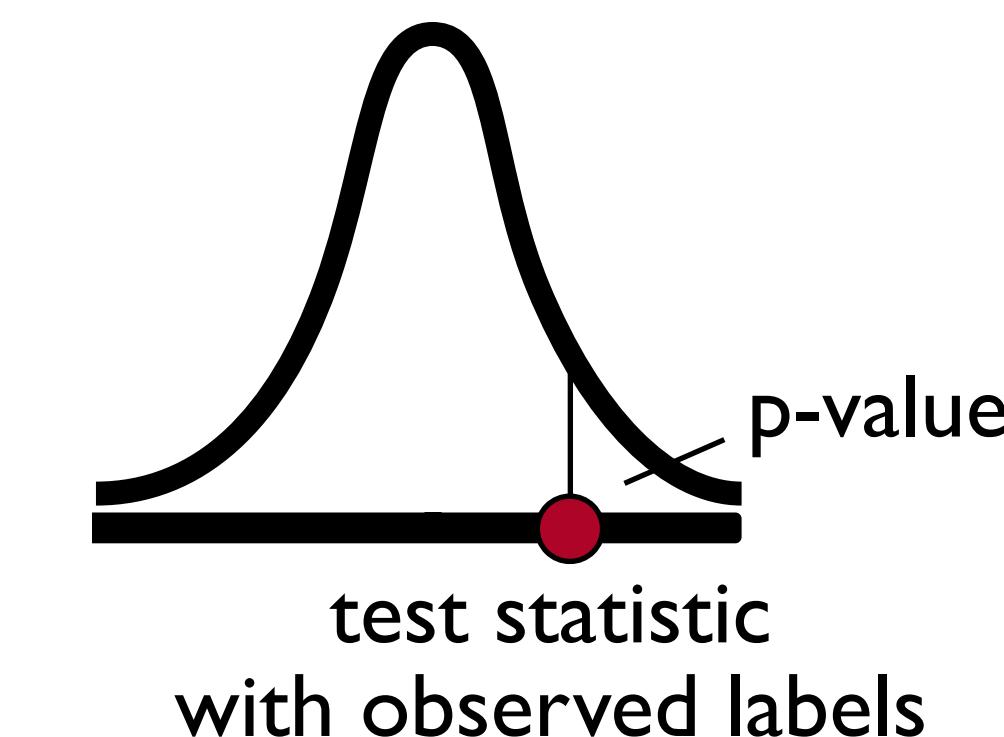


Permute labels to generate a null distribution



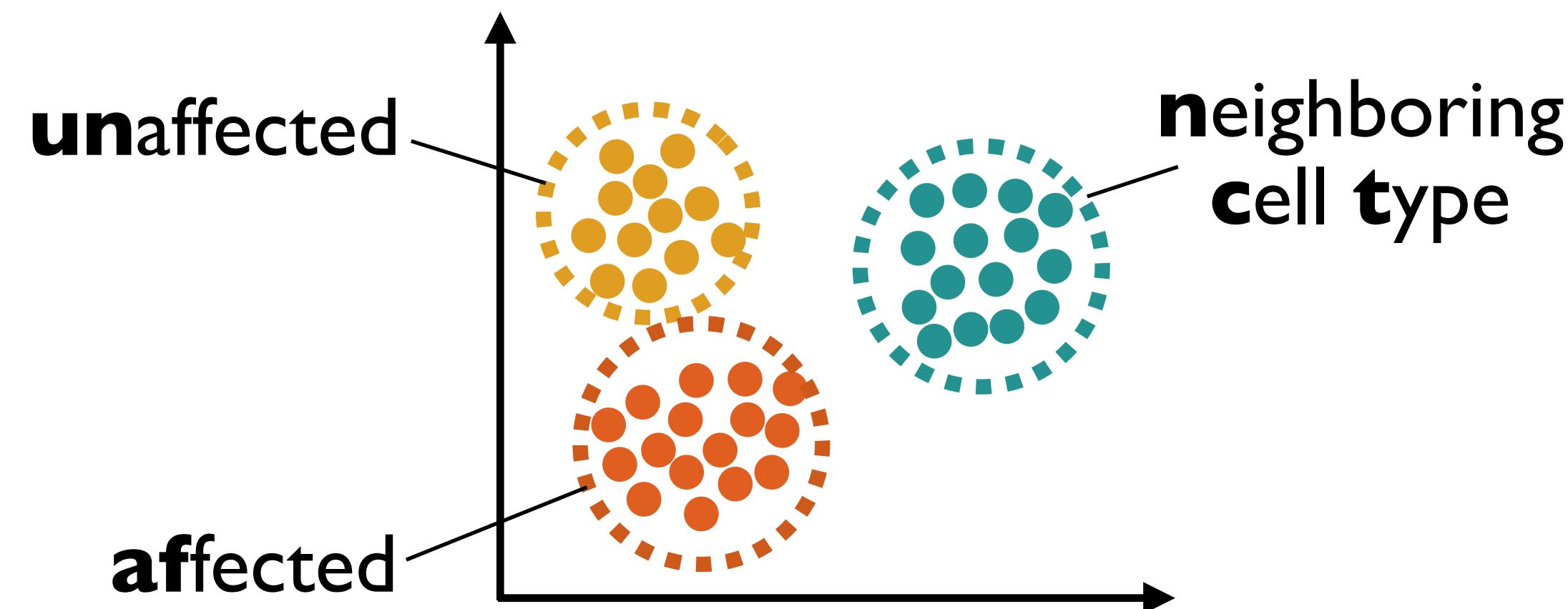
This tests whether this (L, R) interaction is significantly enriched in (clA, clB) relative to any other ordered pair of clusters.

$$\text{Test Statistic} = \text{mean}(\bar{L}_{clA}, \bar{R}_{clB})$$



Empirical distribution of the test statistic under permutation of all labels

Contrast cell-cell interactions of **affected** and **unaffected** cells with a neighboring **cell type**



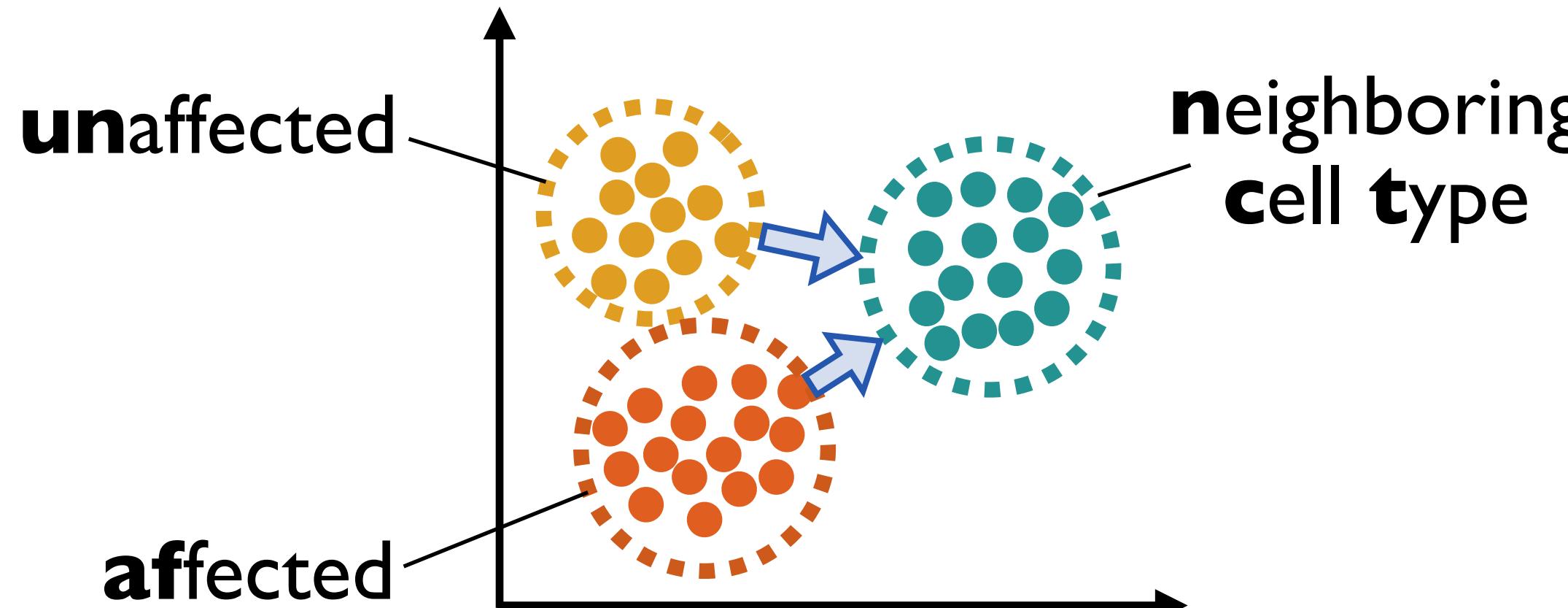
Test for each ordered (L, R) pair if
the interaction between
(un with nct) and **(af with nct)**
is significantly different.

$$H_0: \text{un} \rightarrow \text{nct} \equiv \text{af} \rightarrow \text{nct}$$

$$H_0: \text{nct} \rightarrow \text{un} \equiv \text{nct} \rightarrow \text{af}$$

CellphoneDB builds a general, not targeted, comparison

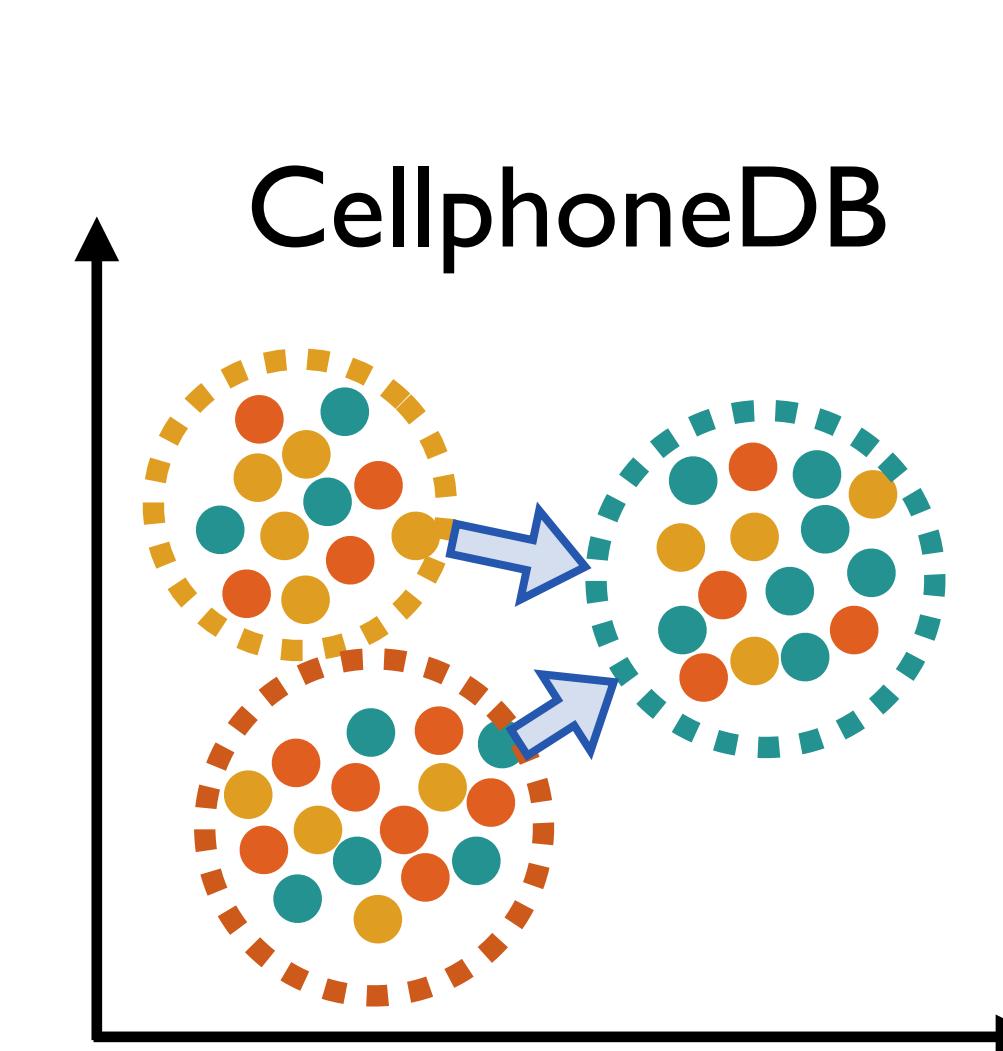
Testing for both directions and comparing to all pairwise cluster combinations



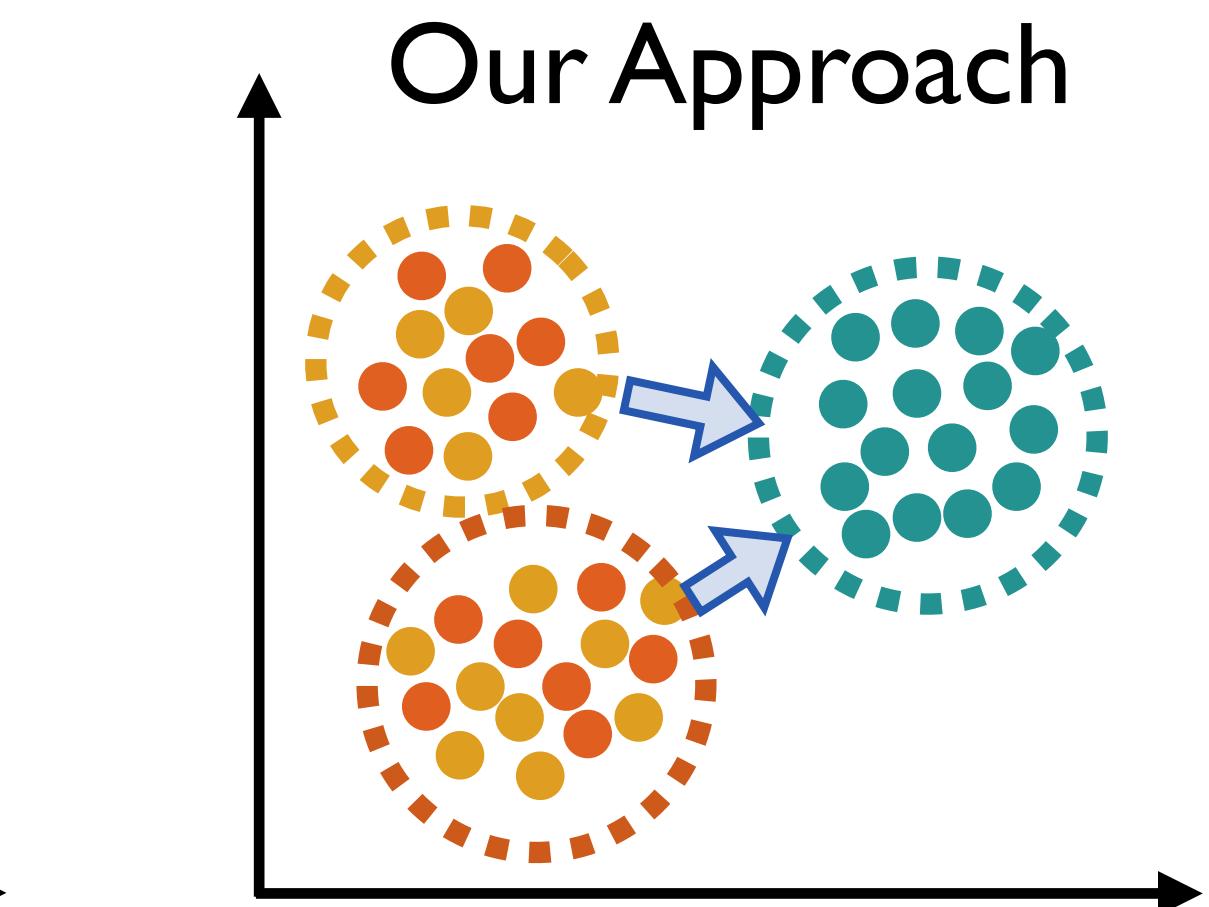
Test for each ordered (L, R) pair if
the interaction between
(un with nct) and (af with nct)
is significantly different.

$$H_0: \text{un} \rightarrow \text{nct} \equiv \text{af} \rightarrow \text{nct}$$

$$H_0: \text{nct} \rightarrow \text{un} \equiv \text{nct} \rightarrow \text{af}$$



	un	af	nct
un			
af			
nct			



Test Statistic =

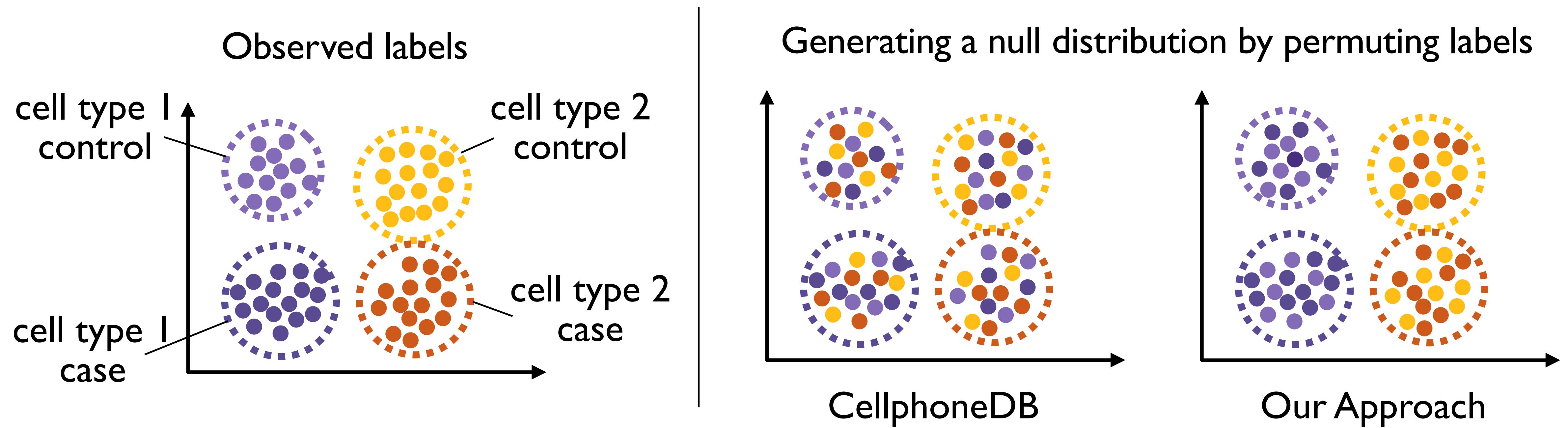
$$\text{mean}(\bar{L}_{un}, \bar{R}_{nct}) - \text{mean}(\bar{L}_{af}, \bar{R}_{nct})$$

Our approach generalizes to various cross-condition contrasts

For example, case-control comparisons

$$H_0: ct1 \rightarrow ct2 \equiv ct1 \rightarrow ct2$$

$$H_0: ct2 \rightarrow ct1 \equiv ct2 \rightarrow ct1$$



Outlook

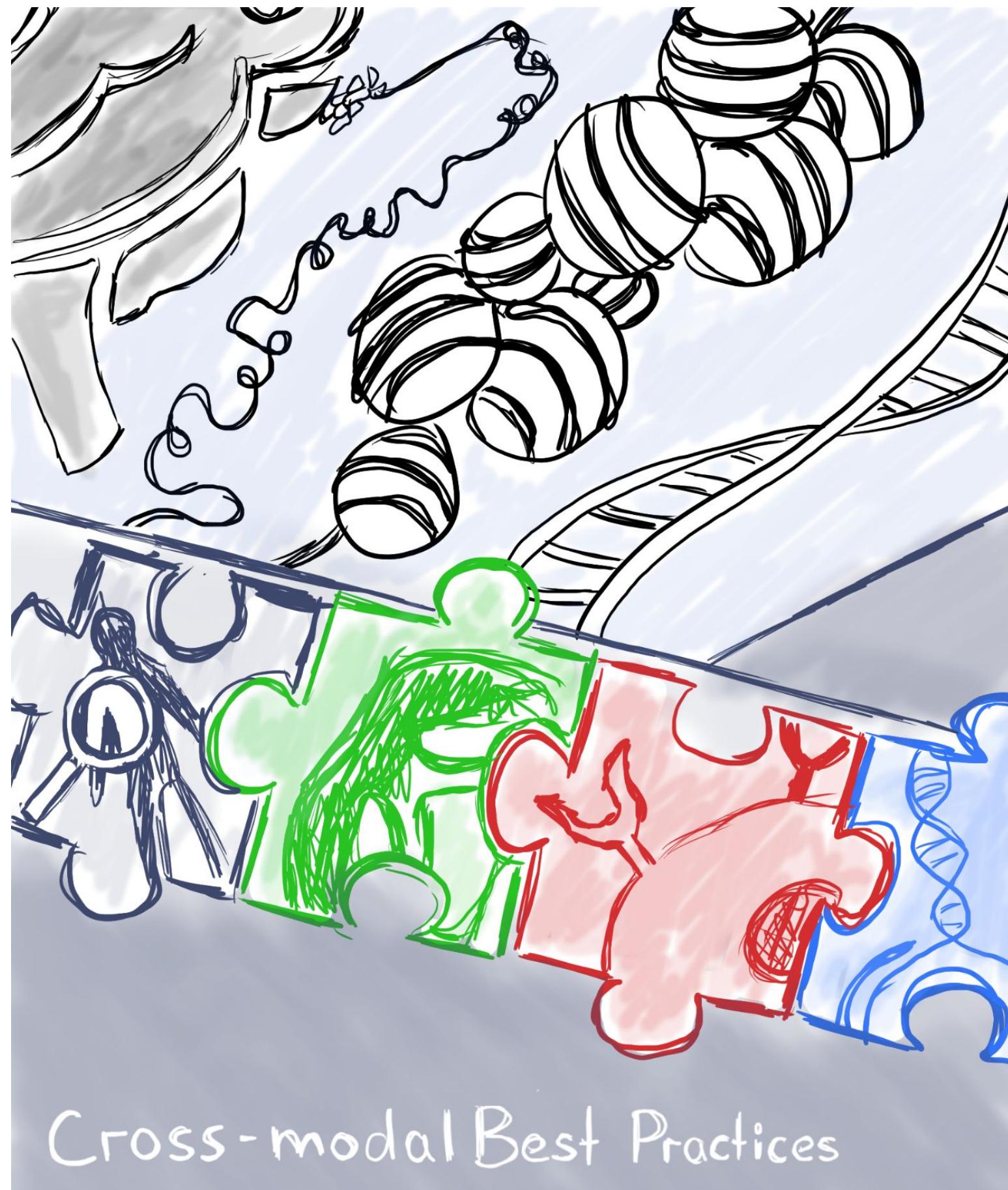
An inventory of perturbations



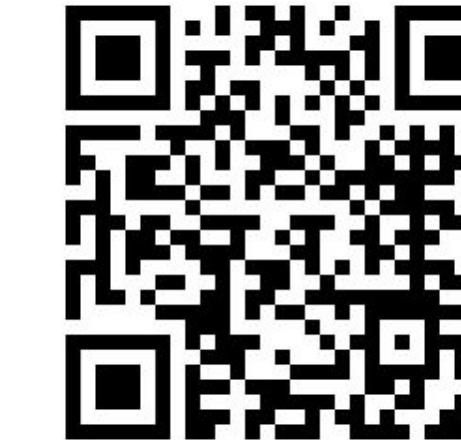
We need large-scale systematic projects that measure and characterize multi-modal cellular changes in response to a comprehensive collection of perturbations so we can relate our datasets to it.



Single cell best practices book



> 40 contributors
> 50 chapters

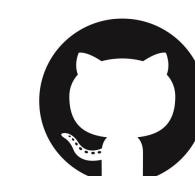


sc-best-practices.org Lukas Heumos



Anna Schaar

- Online book with tutorials and best practices for single cell analysis steps
 - Unimodal and multimodal single cell analysis of RNA, ATAC, ADT, TCR/BCR and spatial omics
 - Preprocessing to advanced downstream analysis
 - Suggestions based on independent benchmarks
 - All chapters downloadable as Jupyter Notebooks
- Living resource for and with the community
 - Contribute content
 - Provide feedback
 - Work in progress...



github.com/theislab/single-cell-best-practices



lukas.heumos@helmholtz-munich.de
anna.schaar@helmholtz-munich.de