

on differential discovery in scRNA-seq data

...using muscat & miloDE/lemur

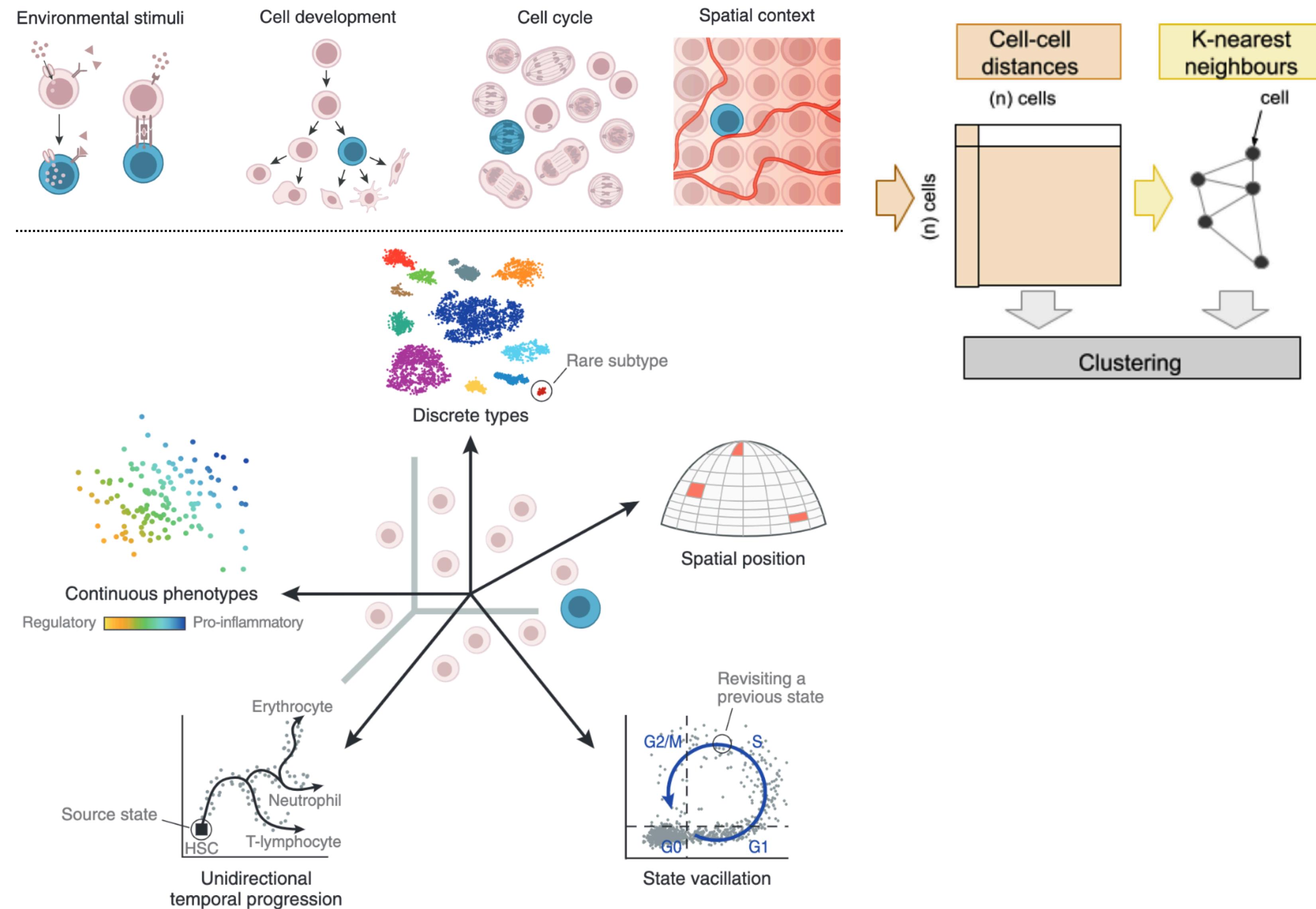
Helena L Crowell, PhD
Autumn School for Single Cell-ers
Oct 21, 2025 · GIMM, Oeiras, Portugal

CNAG – National Center
for Genomic Analysis
Barcelona, Spain

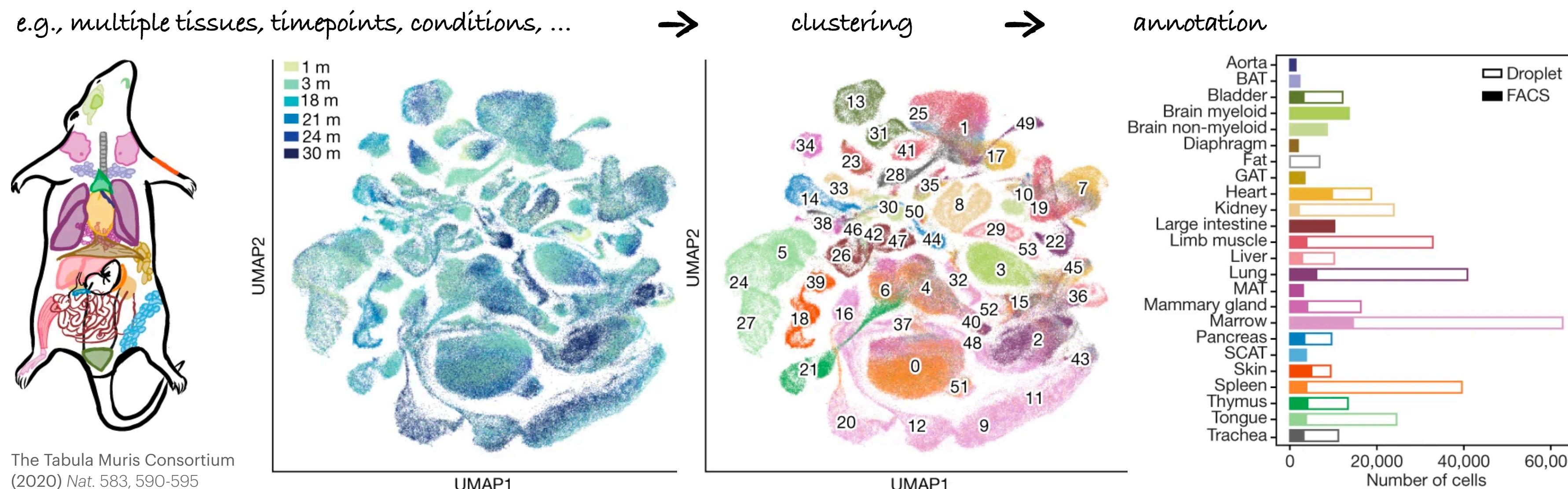
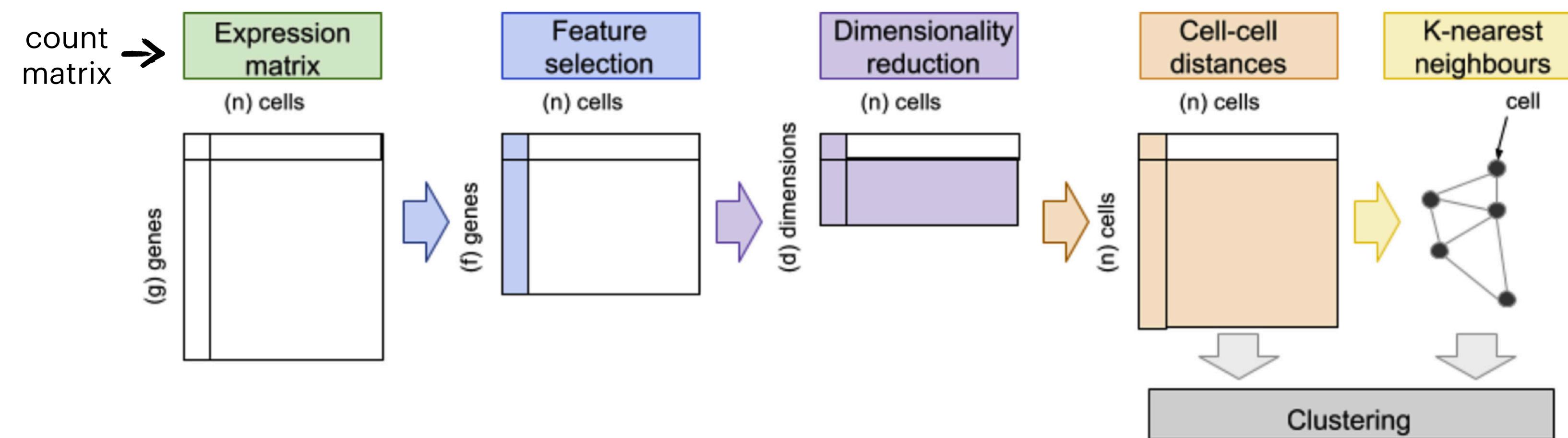


**Swiss National
Science Foundation**

scRNA-seq analysis at a glance



scRNA-seq analysis at a glance



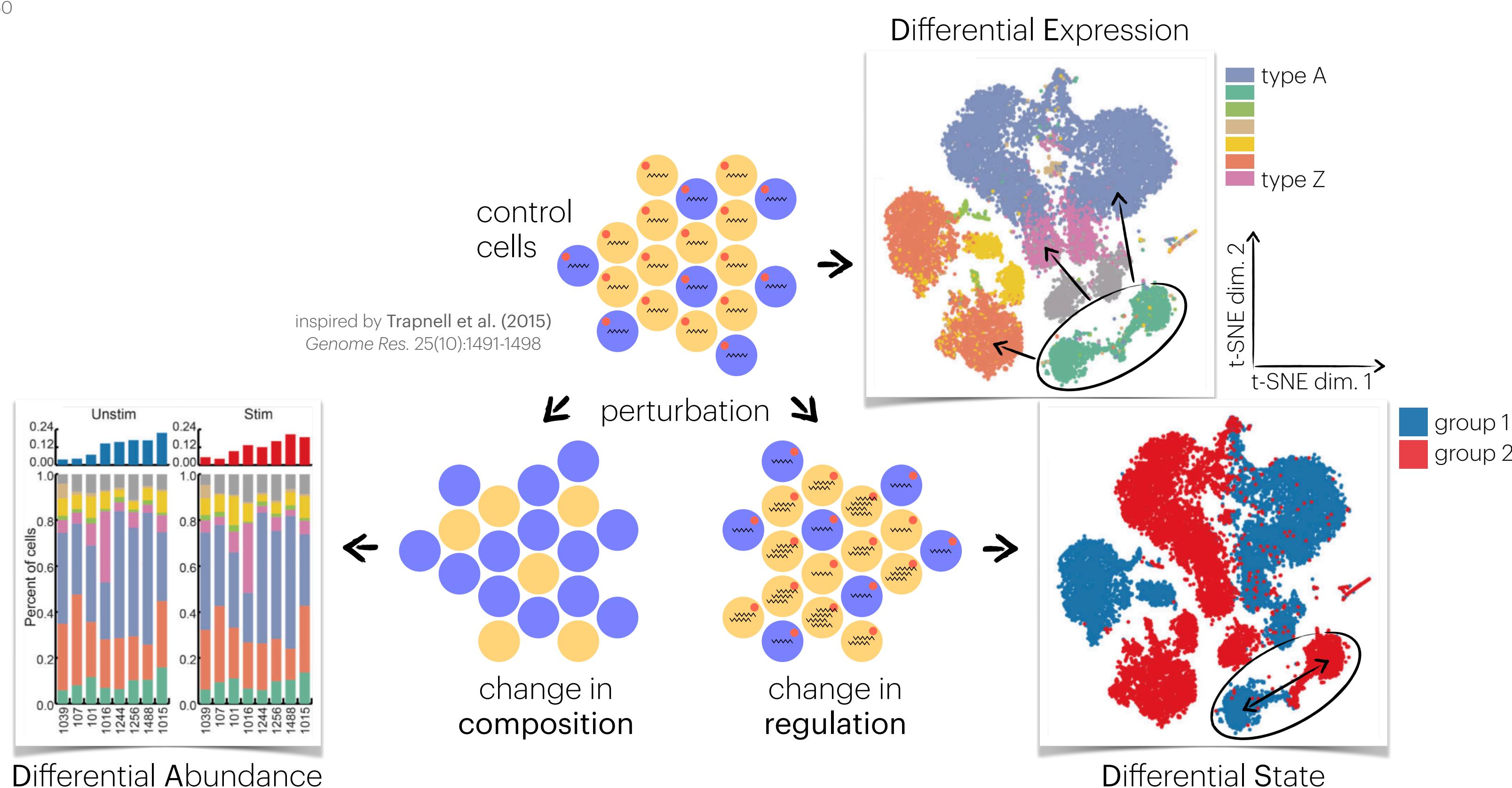
the (wrong but useful) model of cell type & state

The many facets of a cell's identity

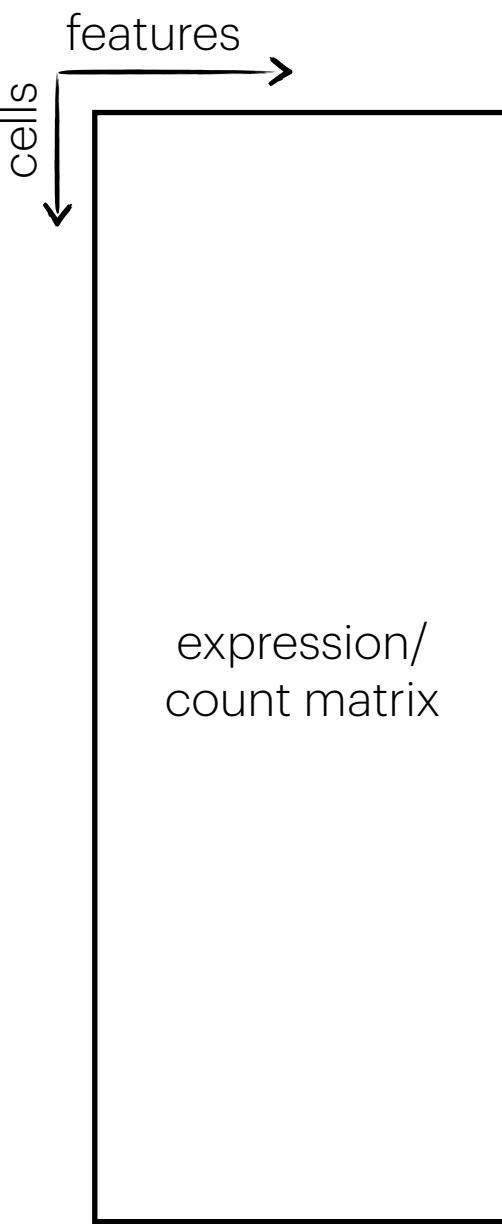
We define a cell's **identity** as the outcome of the instantaneous intersection of all factors that affect it. We refer to the more permanent aspects in a cell's identity as its **type** (e.g., a hepatocyte typically cannot turn into a neuron) and to the more transient elements as its **state**. Cell types are often organized in a hierarchical taxonomy, as types may be further

Wagner et al. (2016) *Nat. Biotechnol.* 34(11):1145-1160

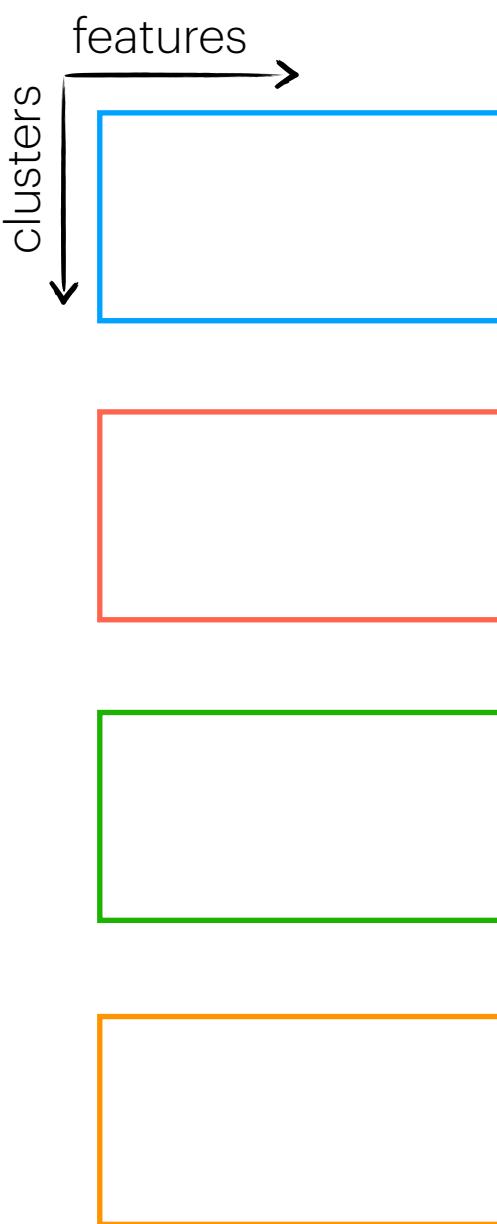
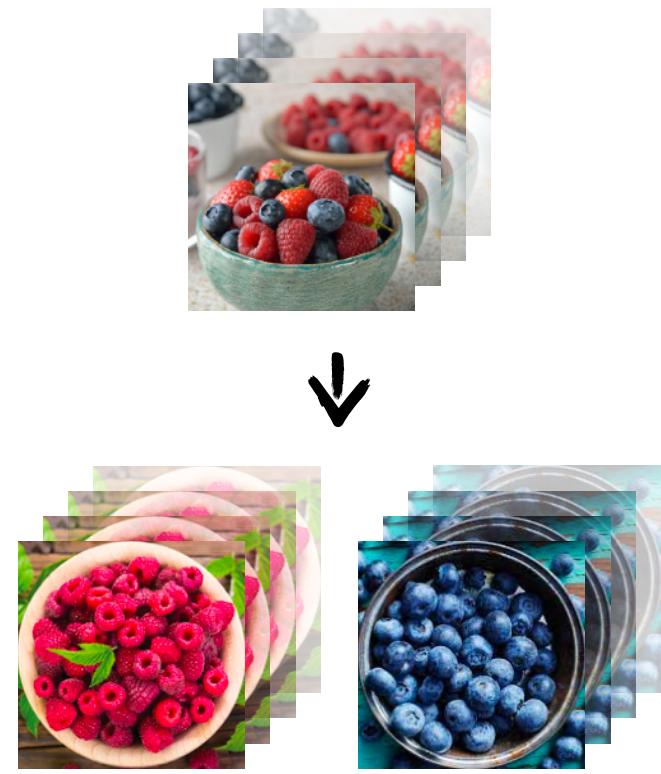
type = permanent/discrete
state = transient/continuous



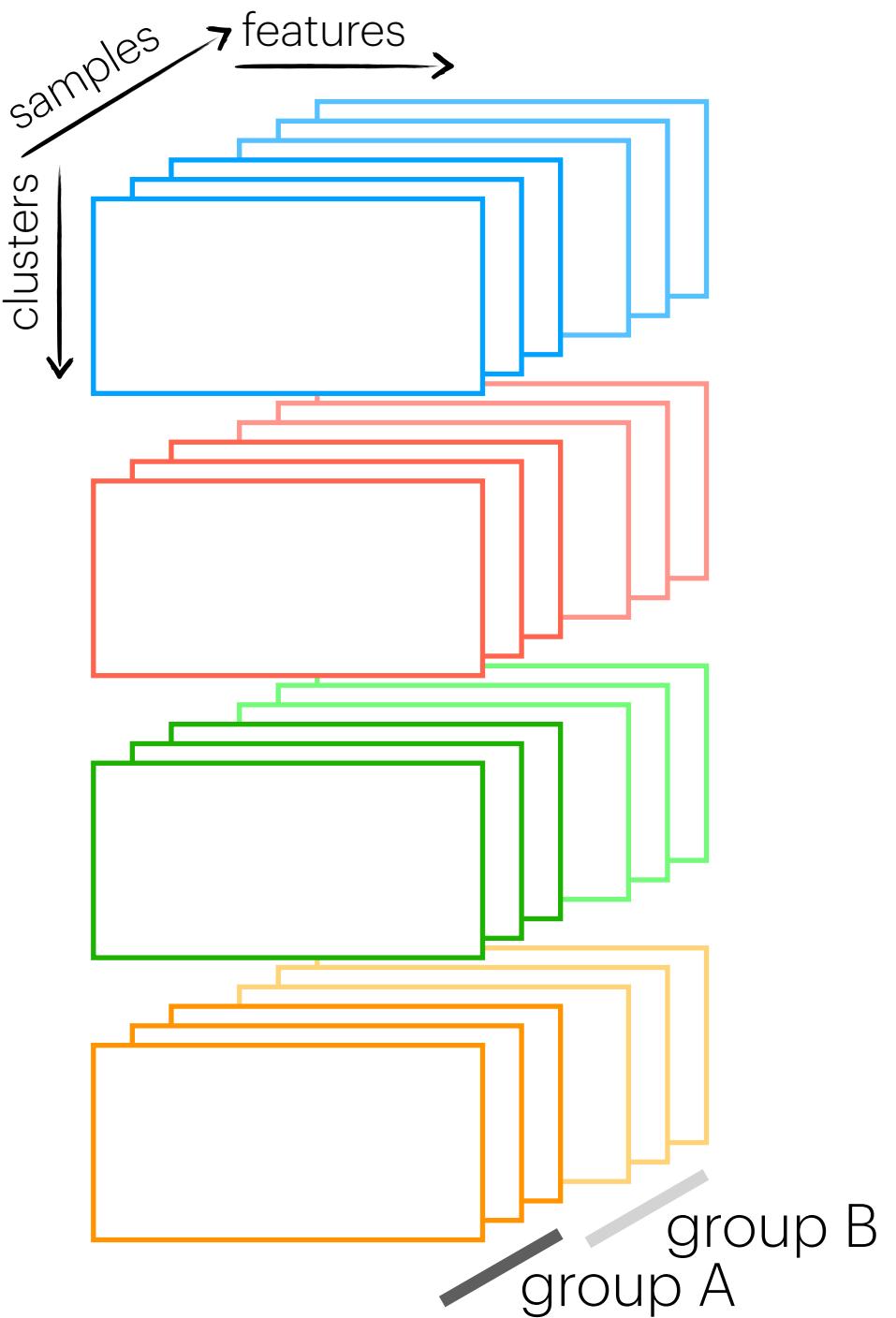
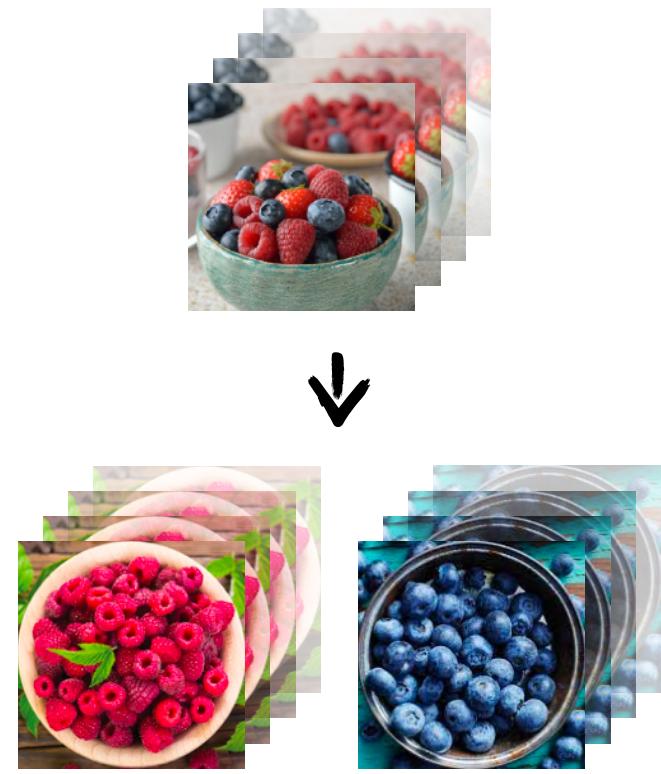
bookkeeping – DS analysis starts with...



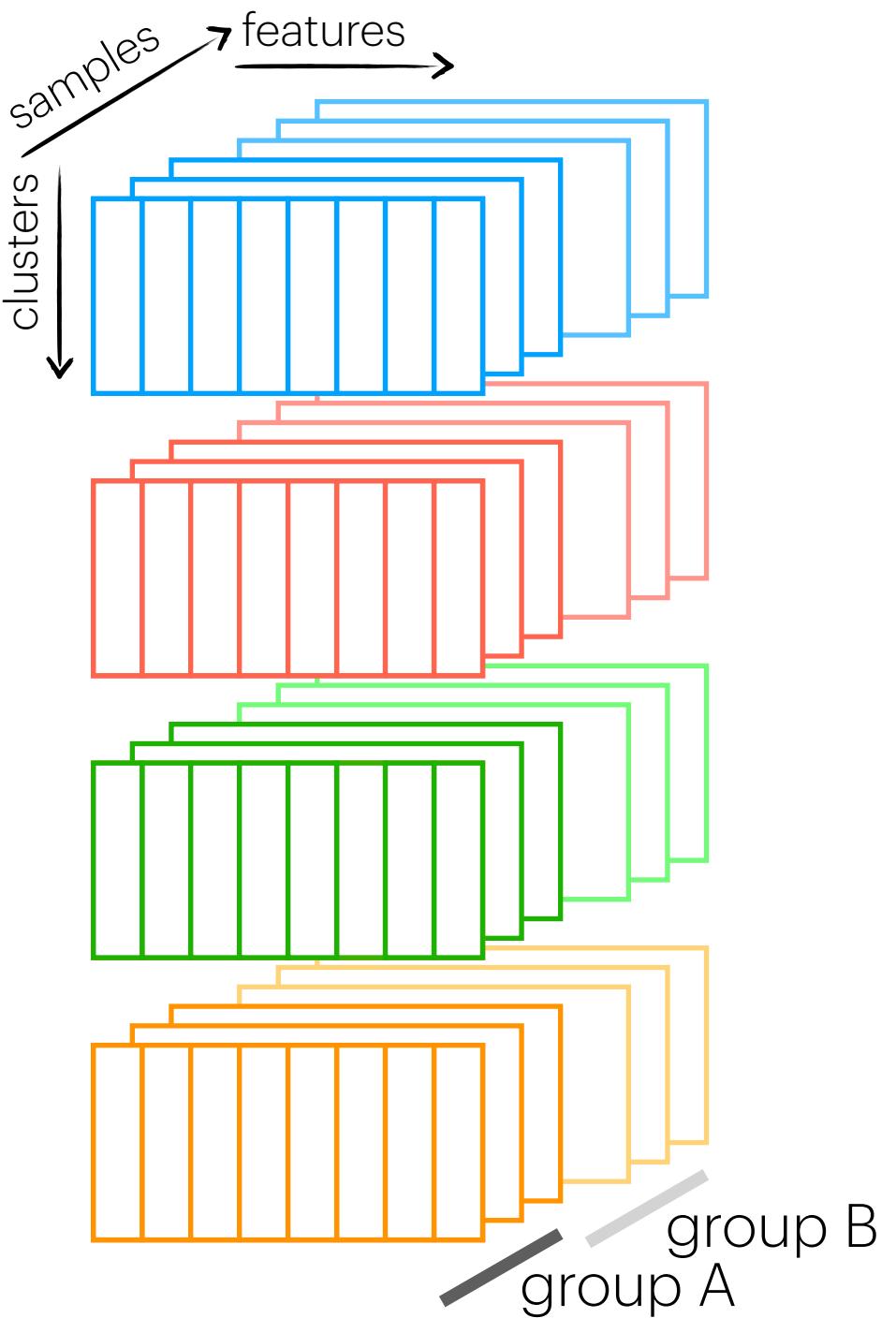
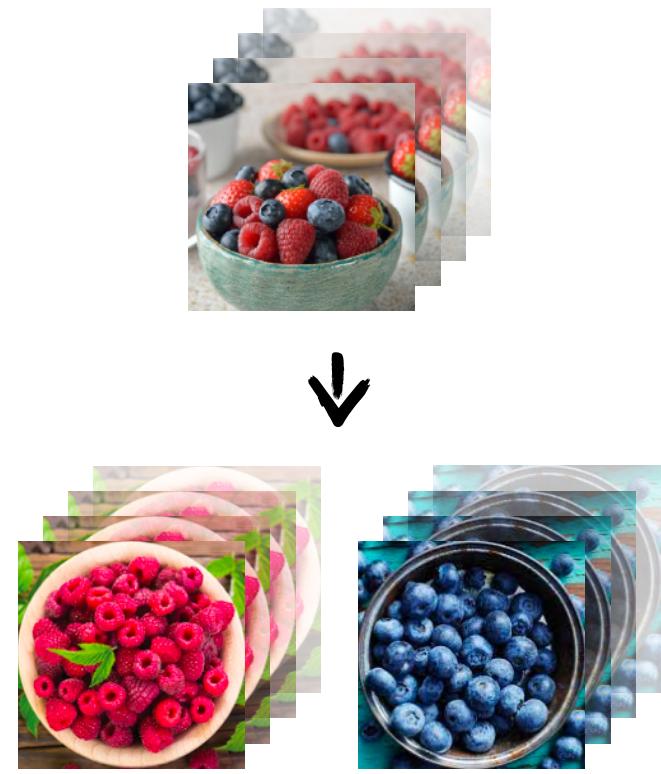
bookkeeping – DS analysis starts with...



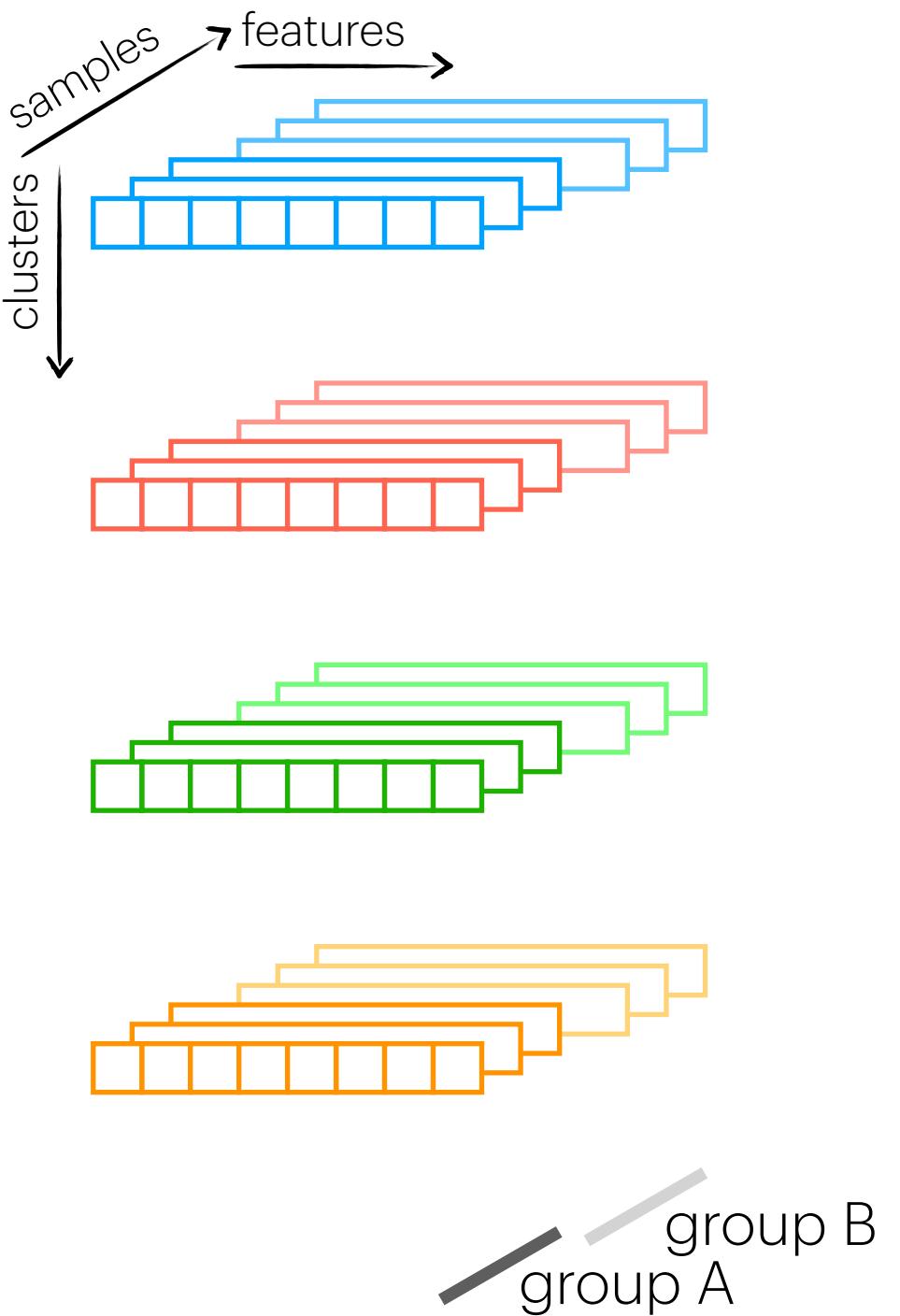
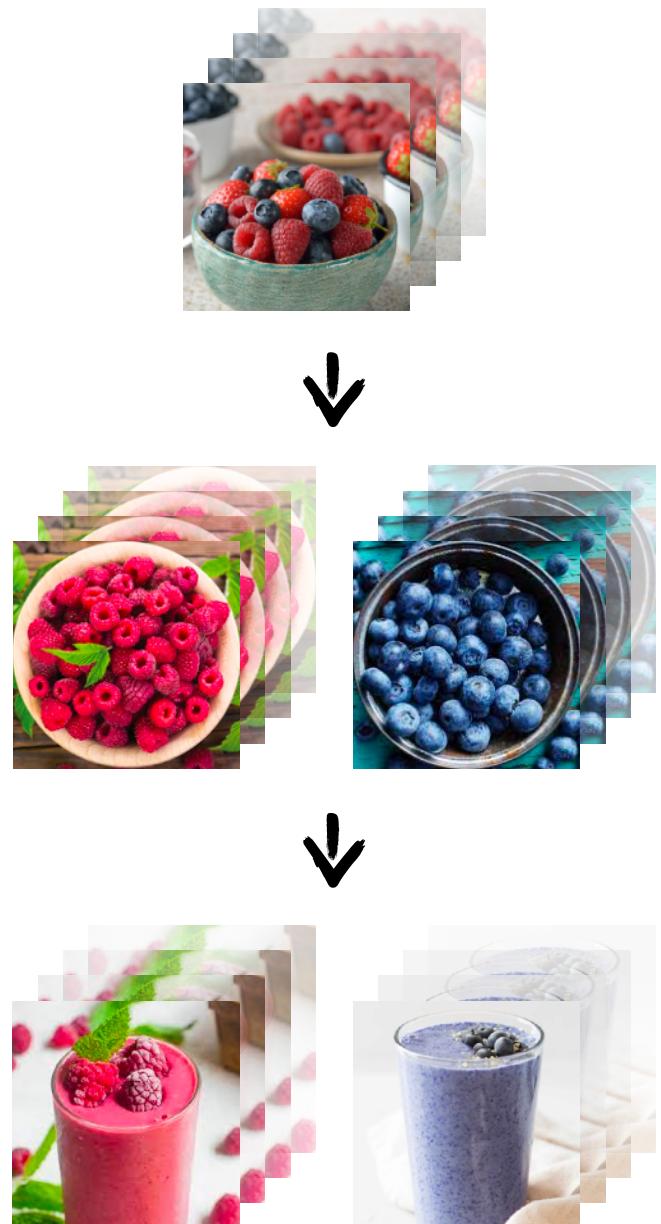
bookkeeping – DS analysis starts with...



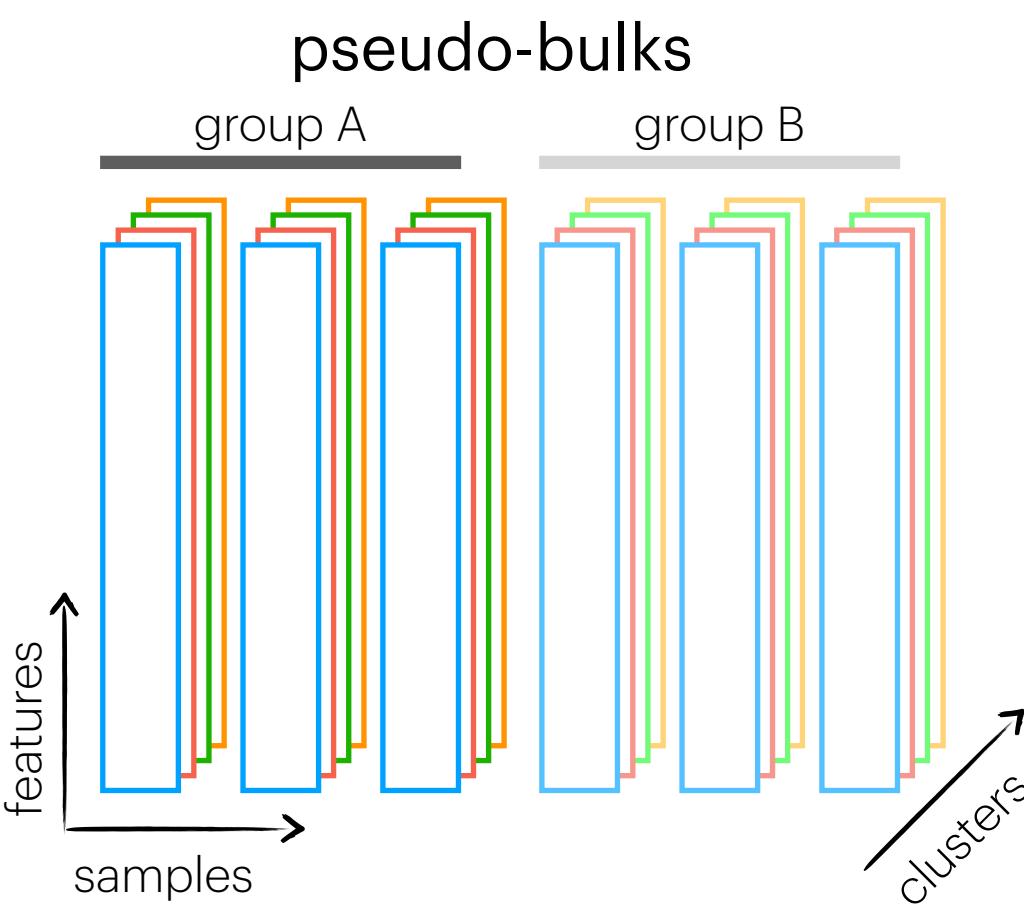
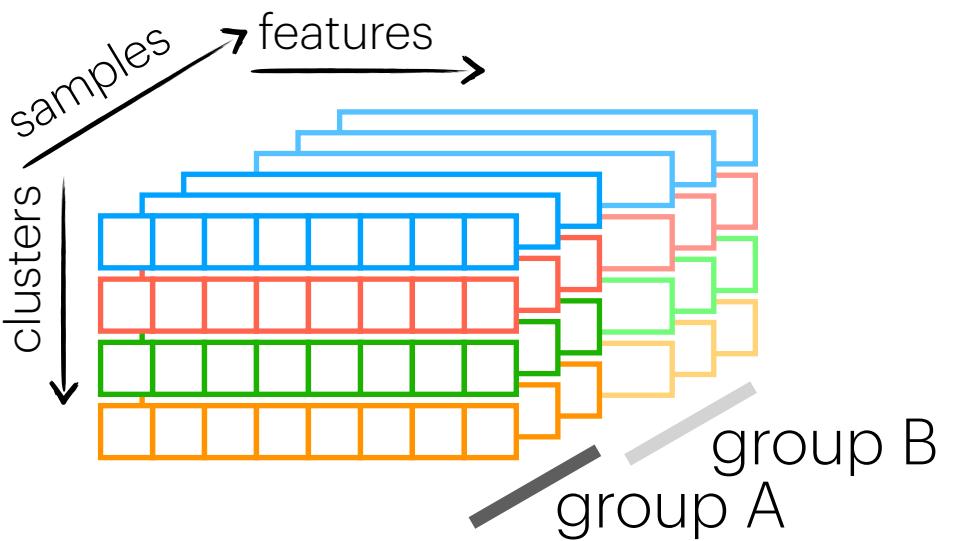
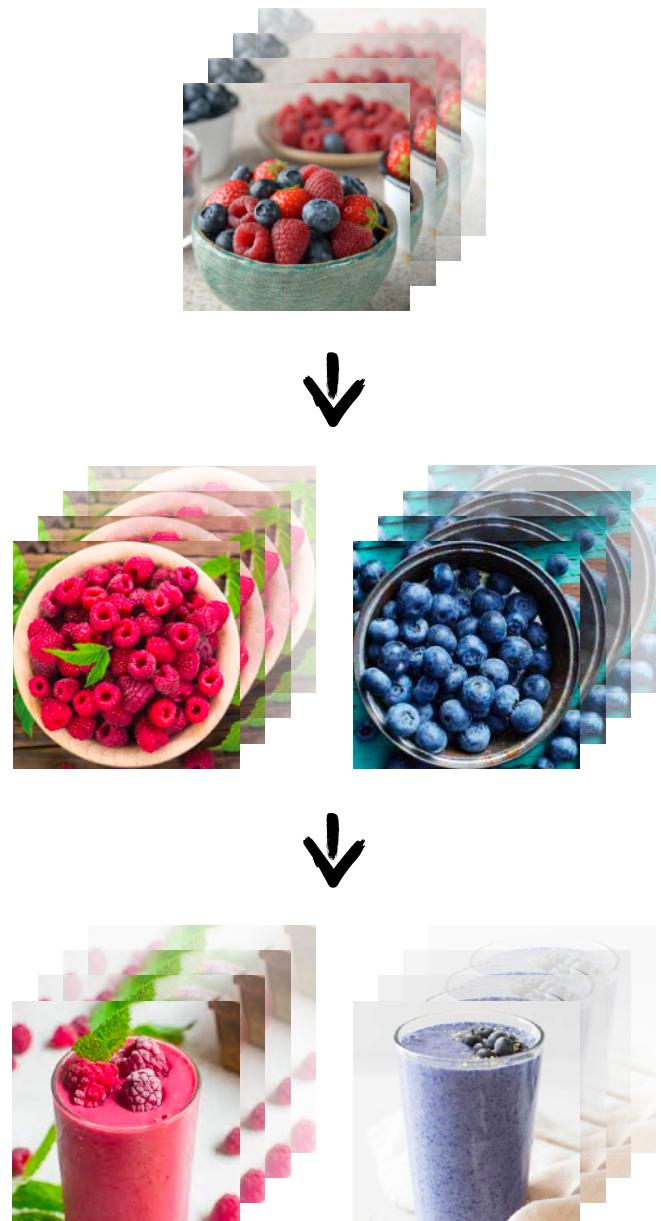
bookkeeping – DS analysis starts with...



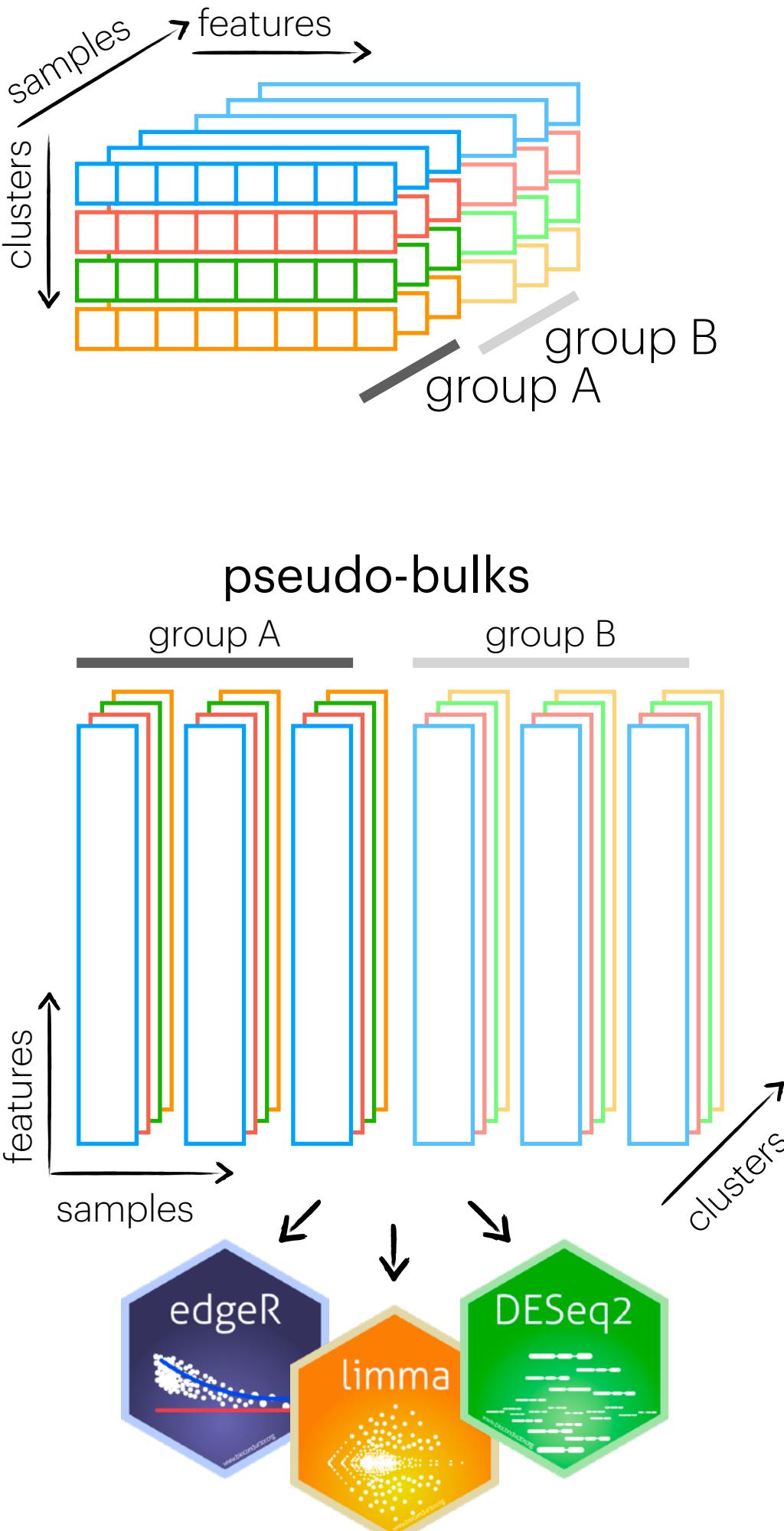
bookkeeping – DS analysis starts with...



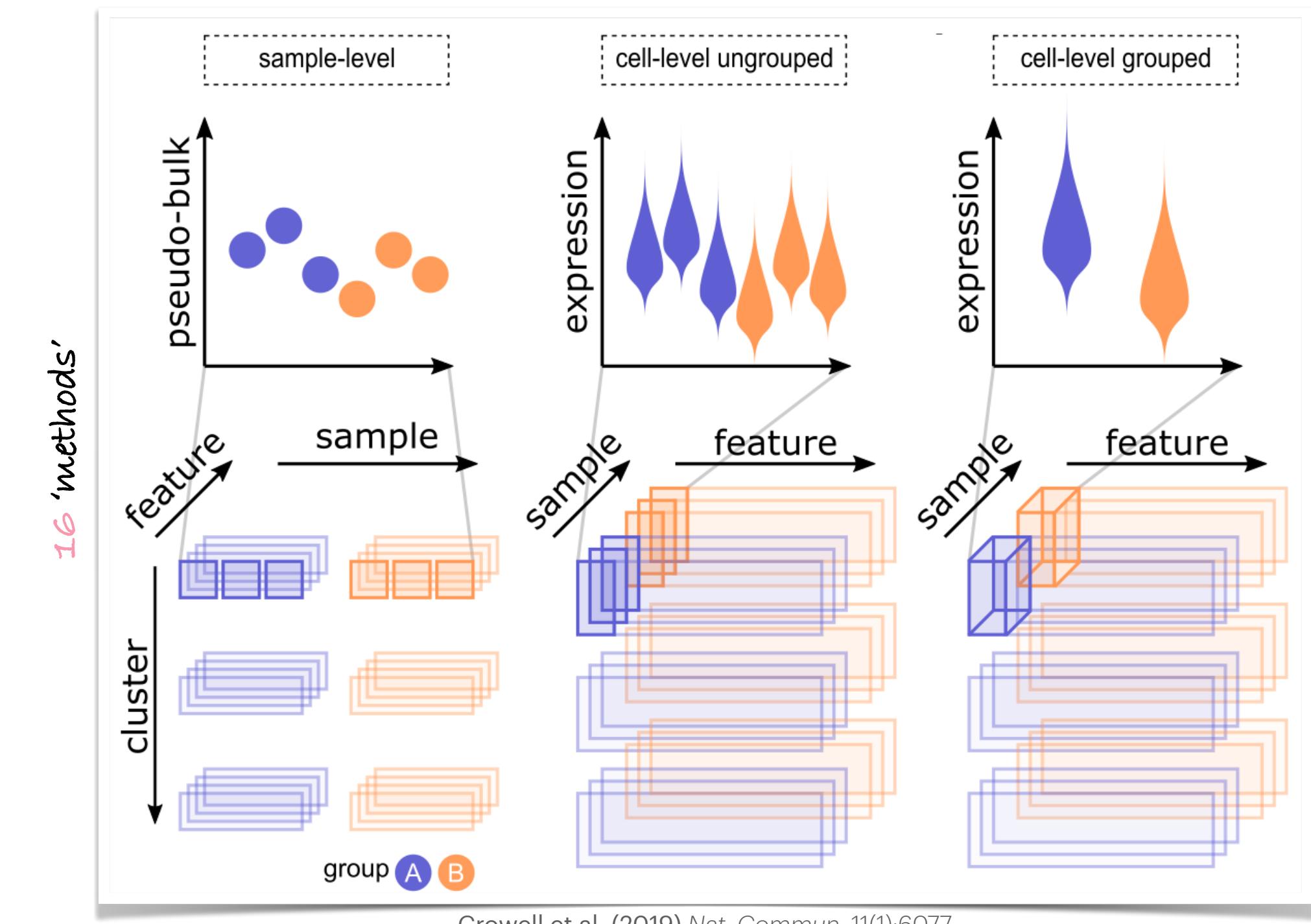
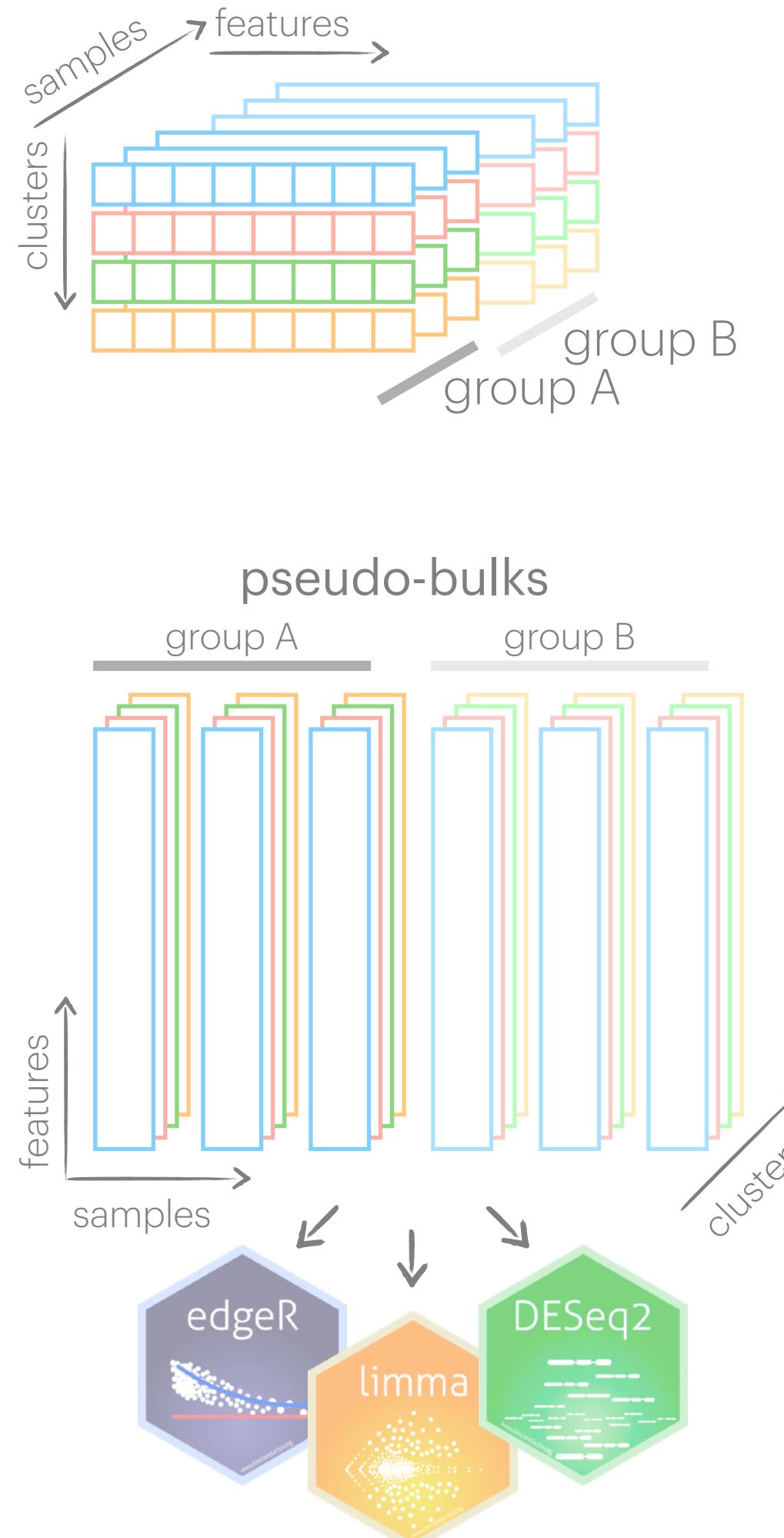
bookkeeping – DS analysis starts with...



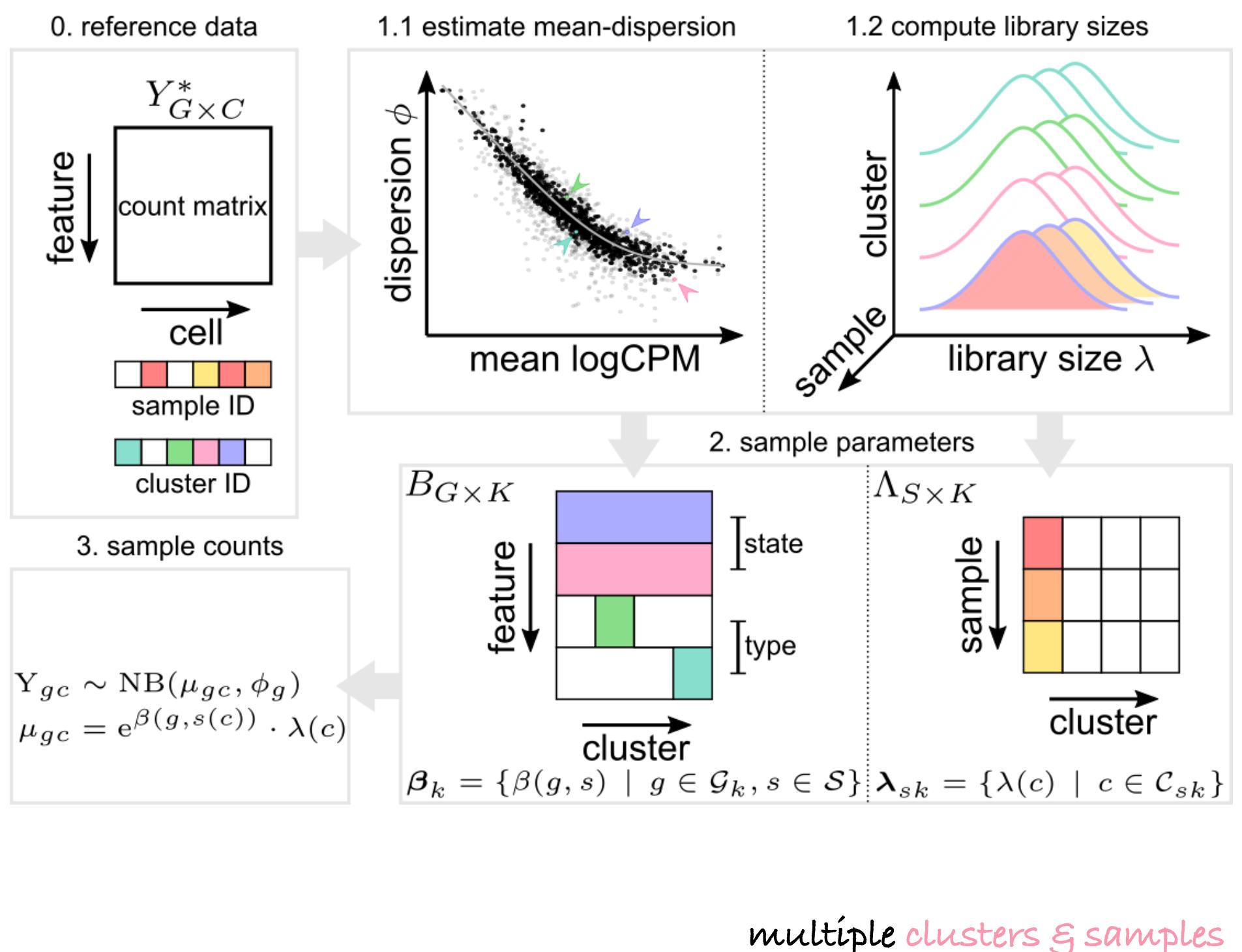
bookkeeping – DS analysis starts with...



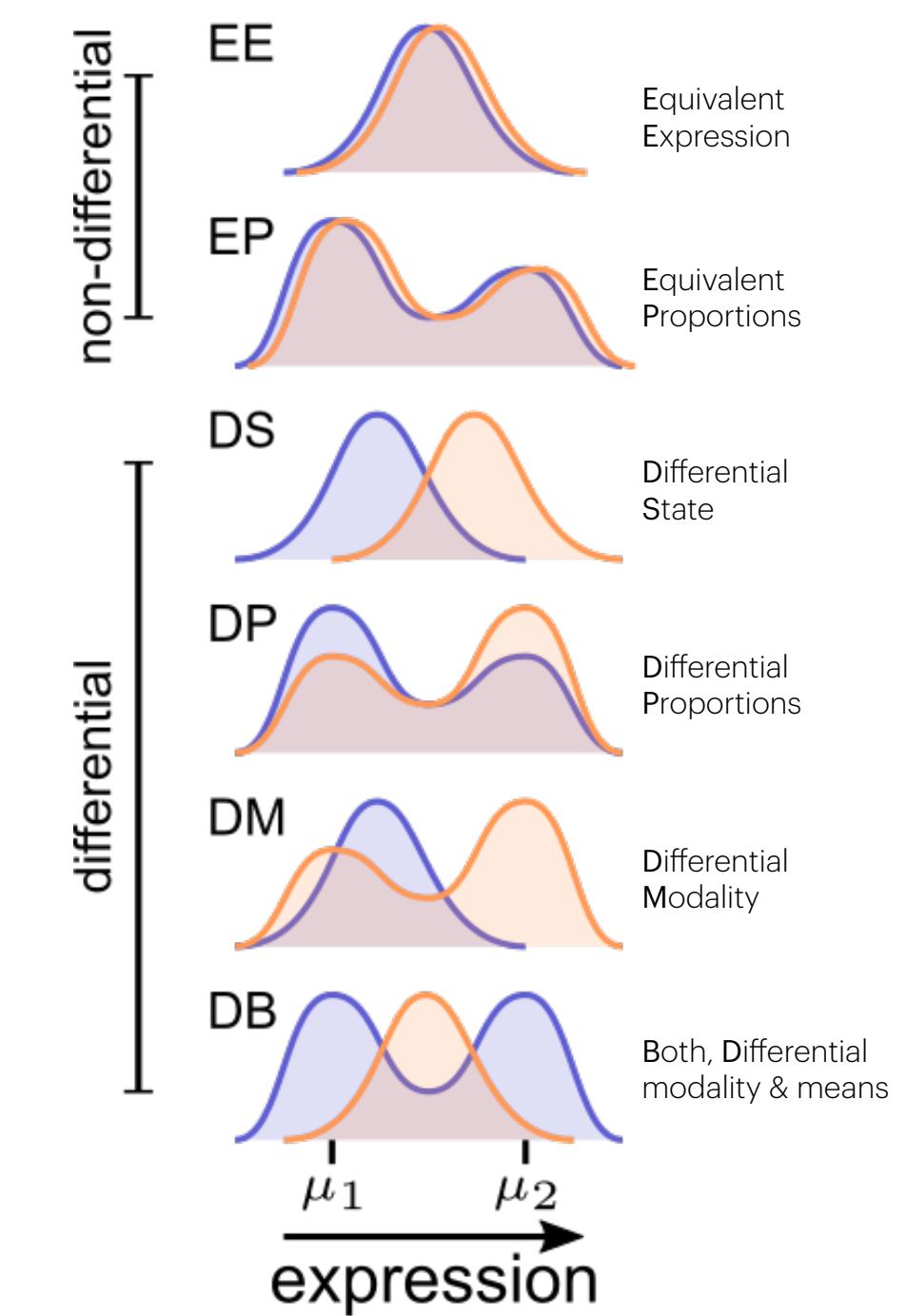
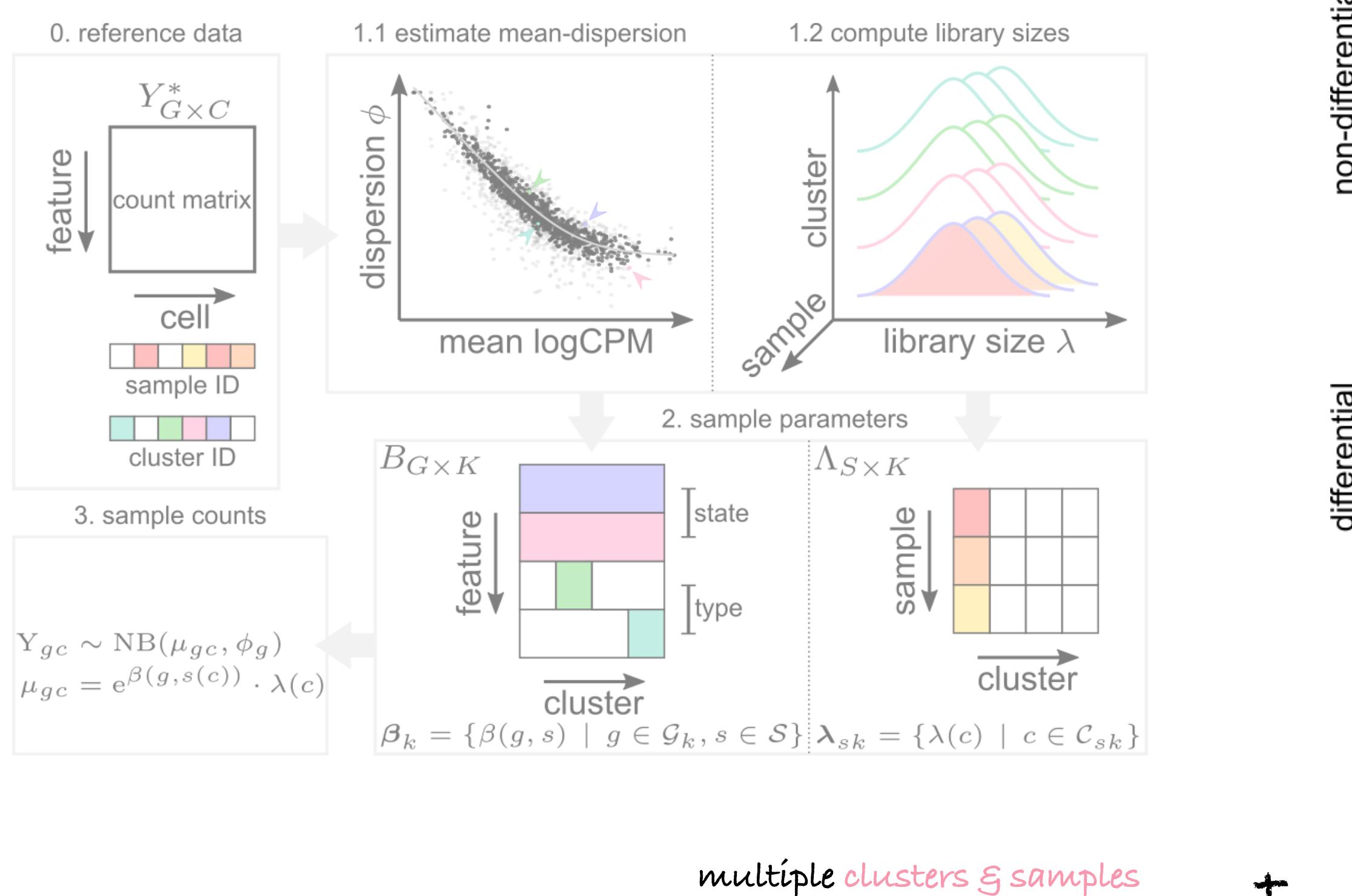
bookkeeping – DS analysis starts with...



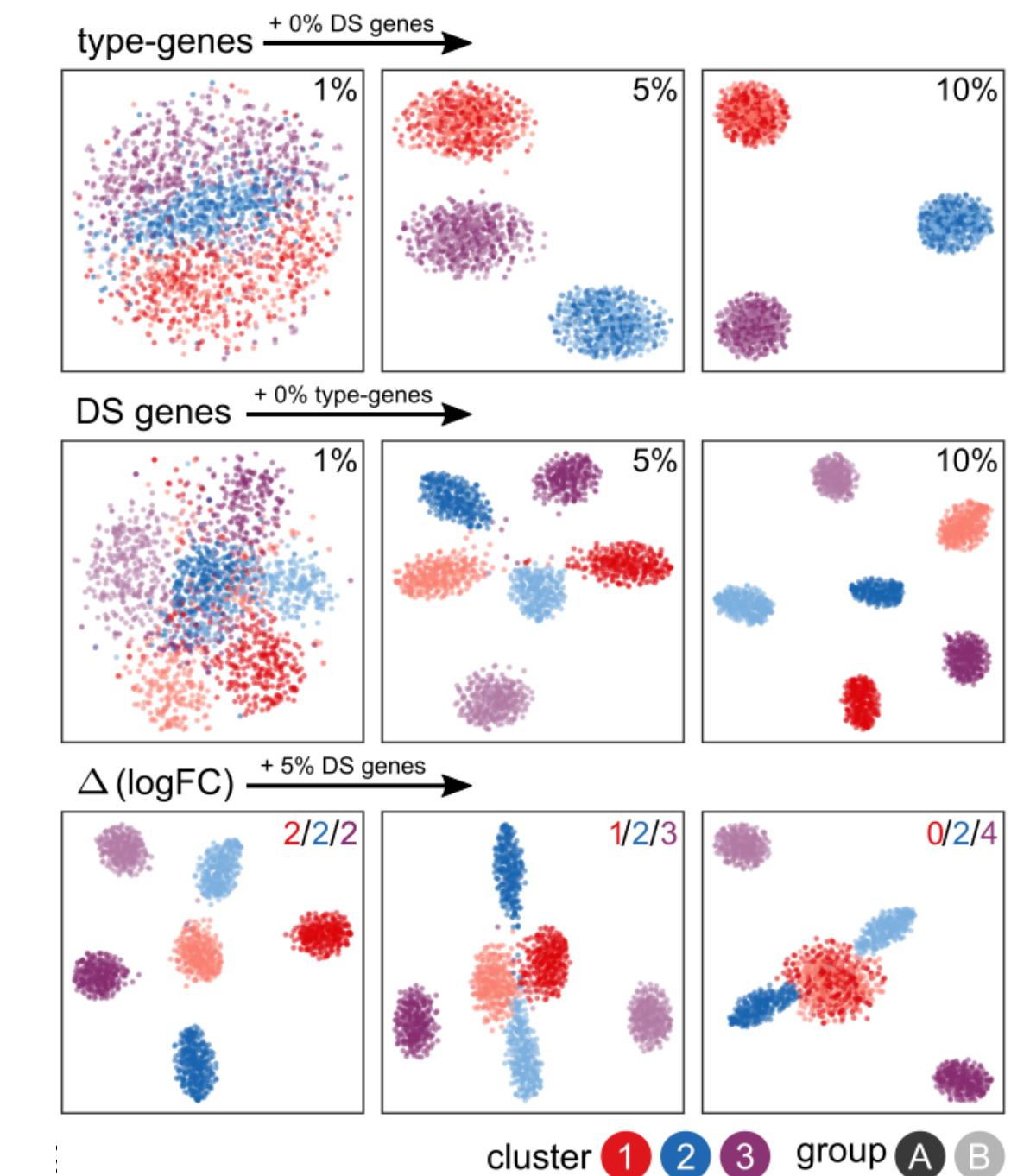
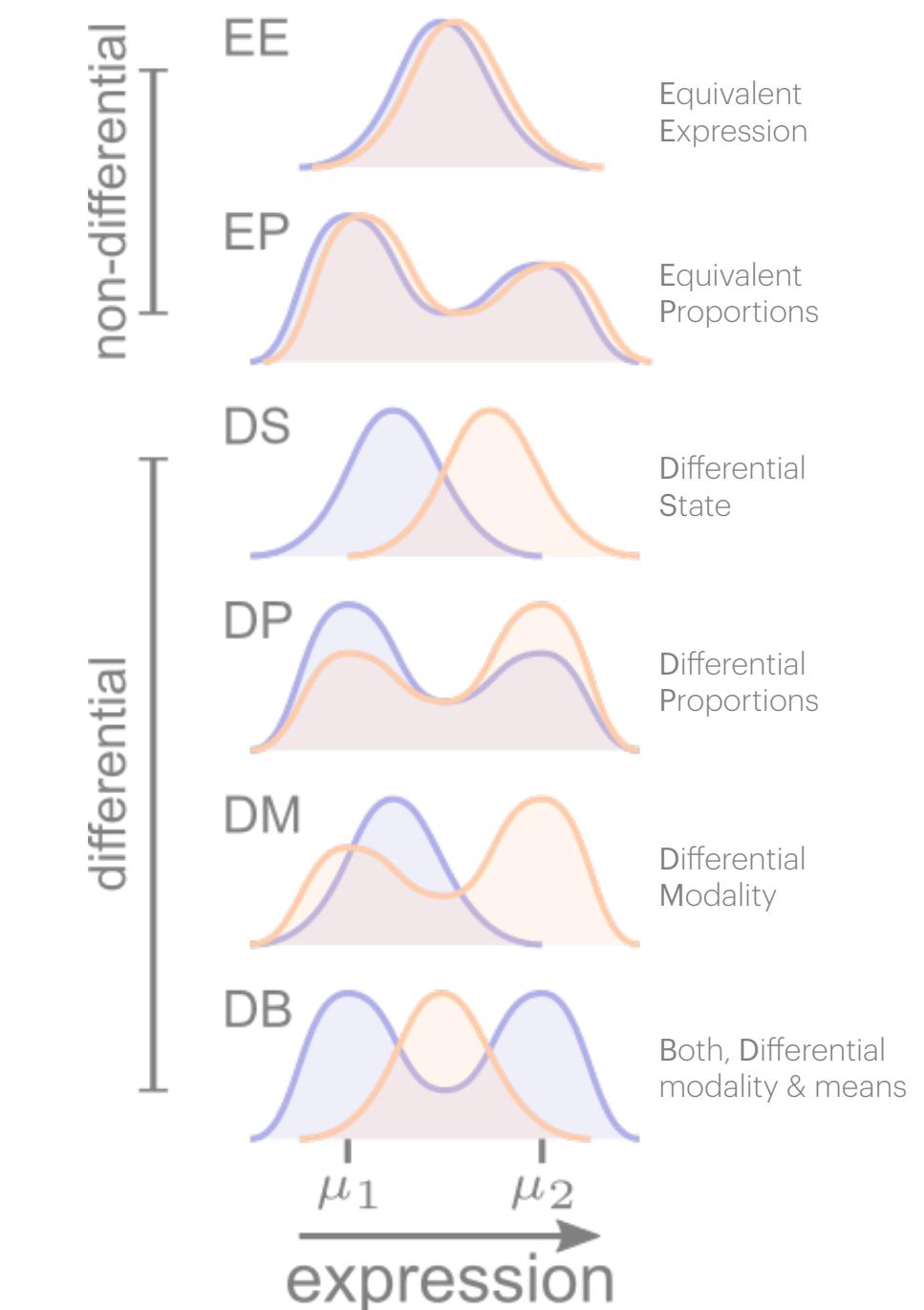
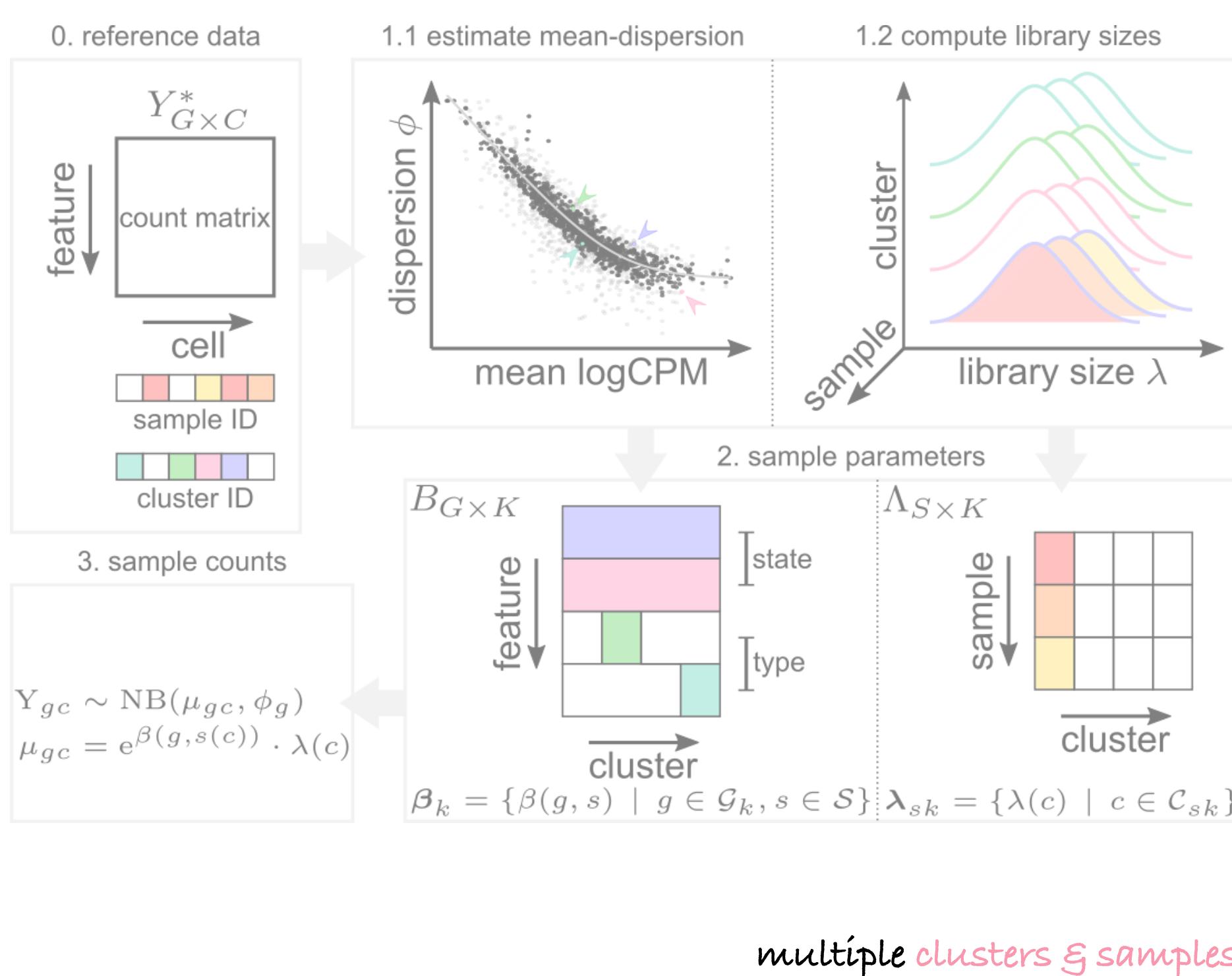
simulation framework to judge method performance



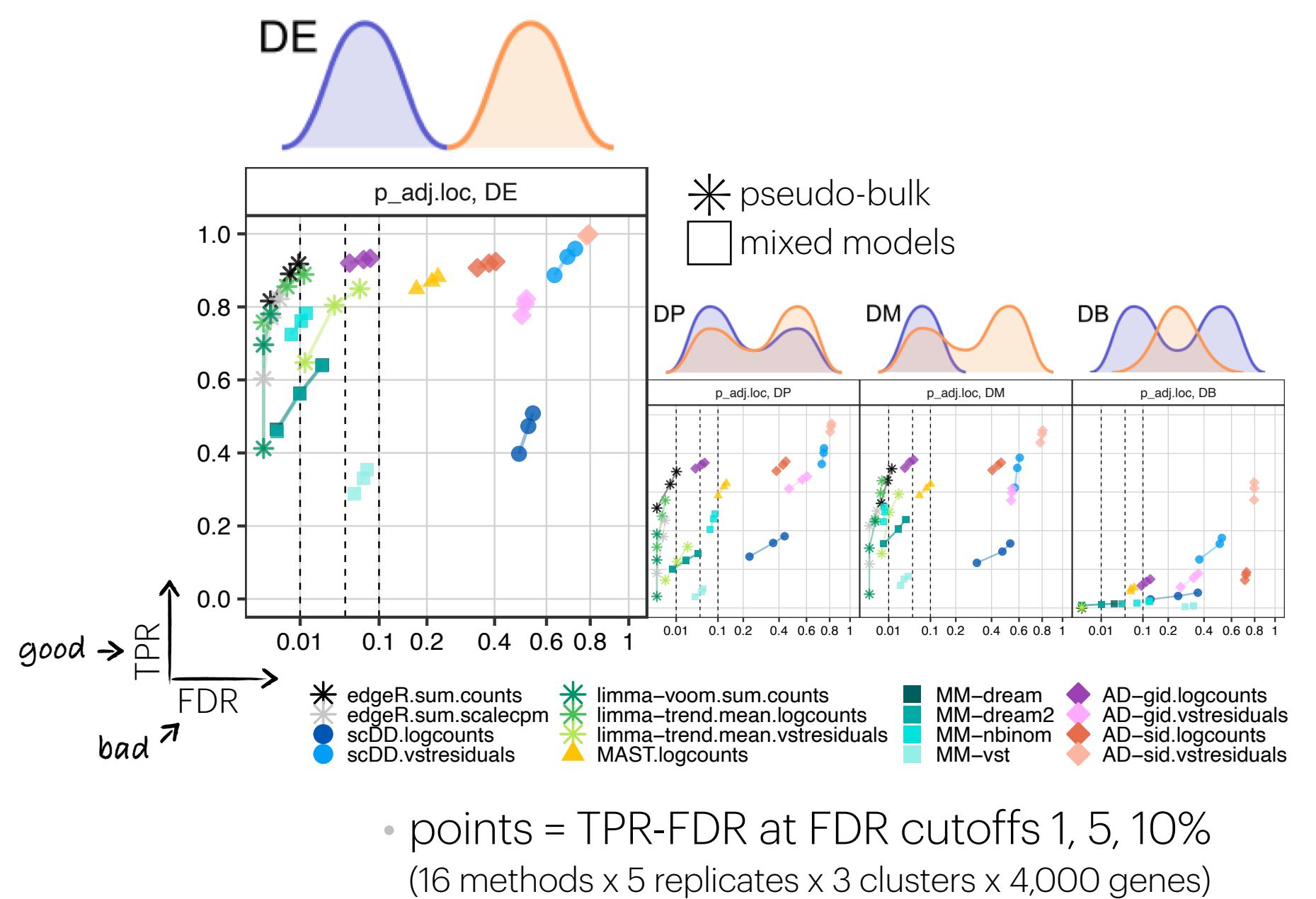
simulation framework to judge method performance



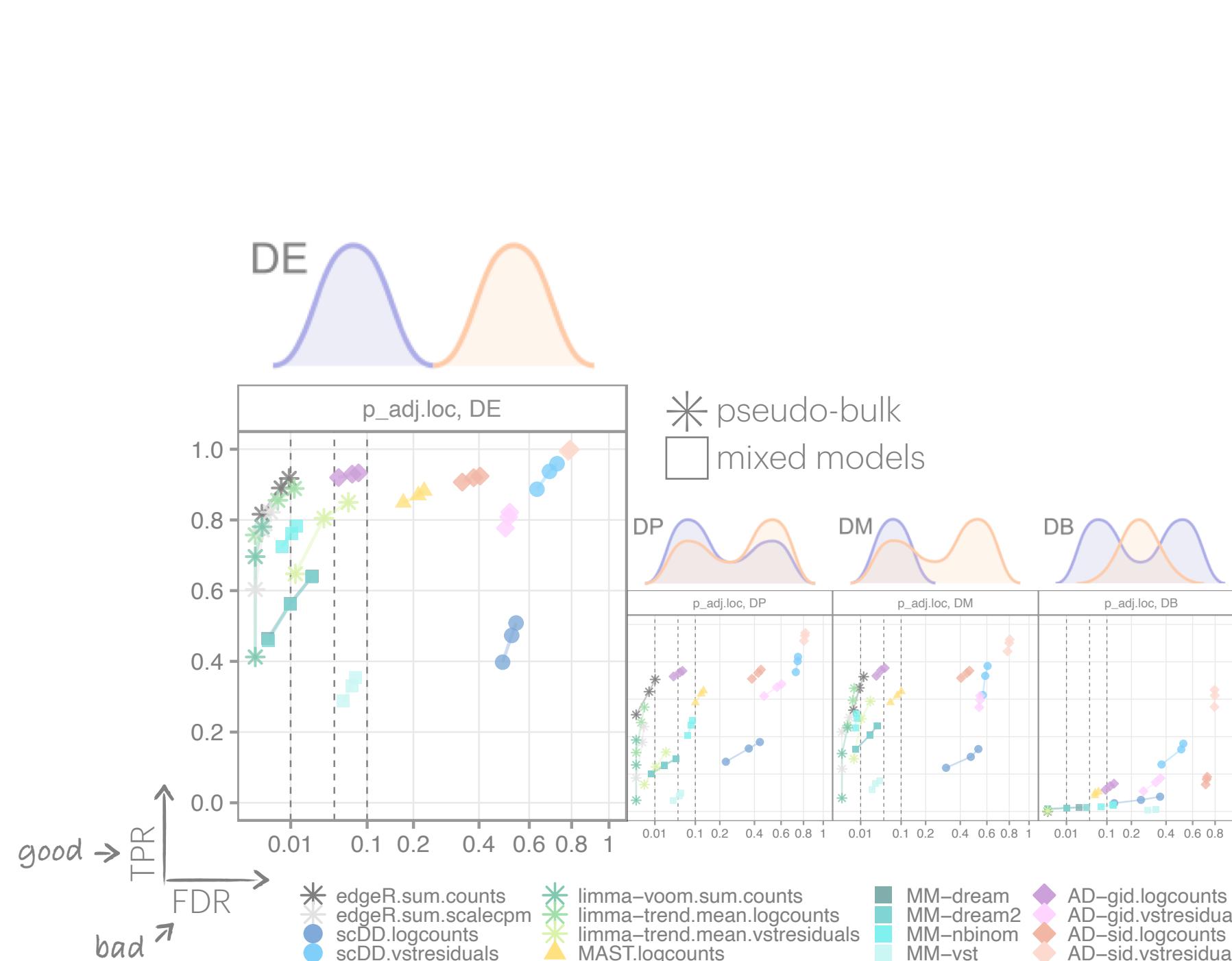
simulation framework to judge method performance



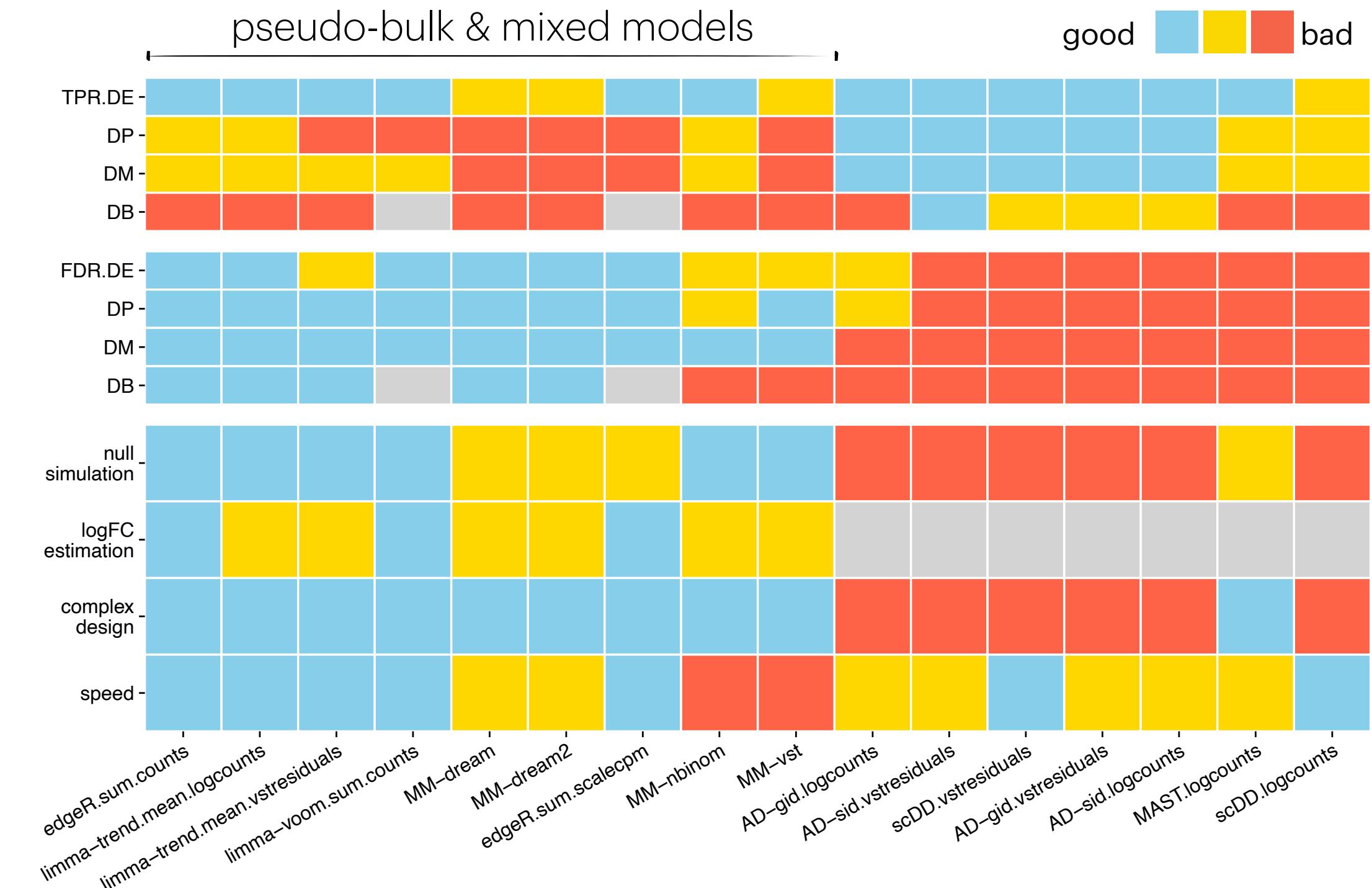
simulation framework to judge method performance



simulation framework to judge method performance



- points = TPR-FDR at FDR cutoffs 1, 5, 10%
(16 methods x 5 replicates x 3 clusters x 4,000 genes)



A balanced measure shows superior performance of pseudobulk methods in single-cell RNA-sequencing analysis

Murphy & Skene (2022) Nat. Commun. 13:785

Confronting false discoveries in single-cell differential expression

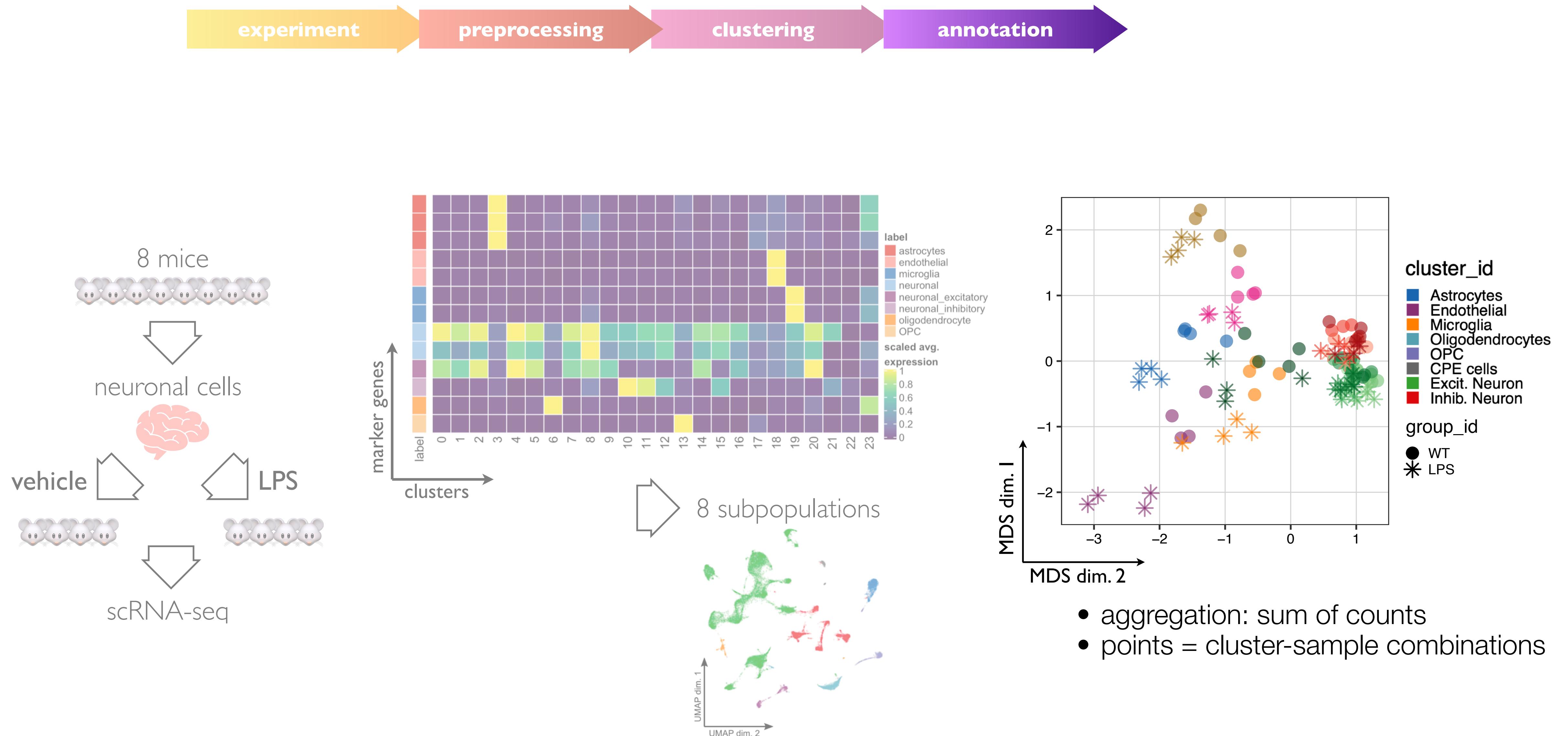
Squair et al. (2021) Nat. Commun. 12:5691

Compare healthy and diabetic samples

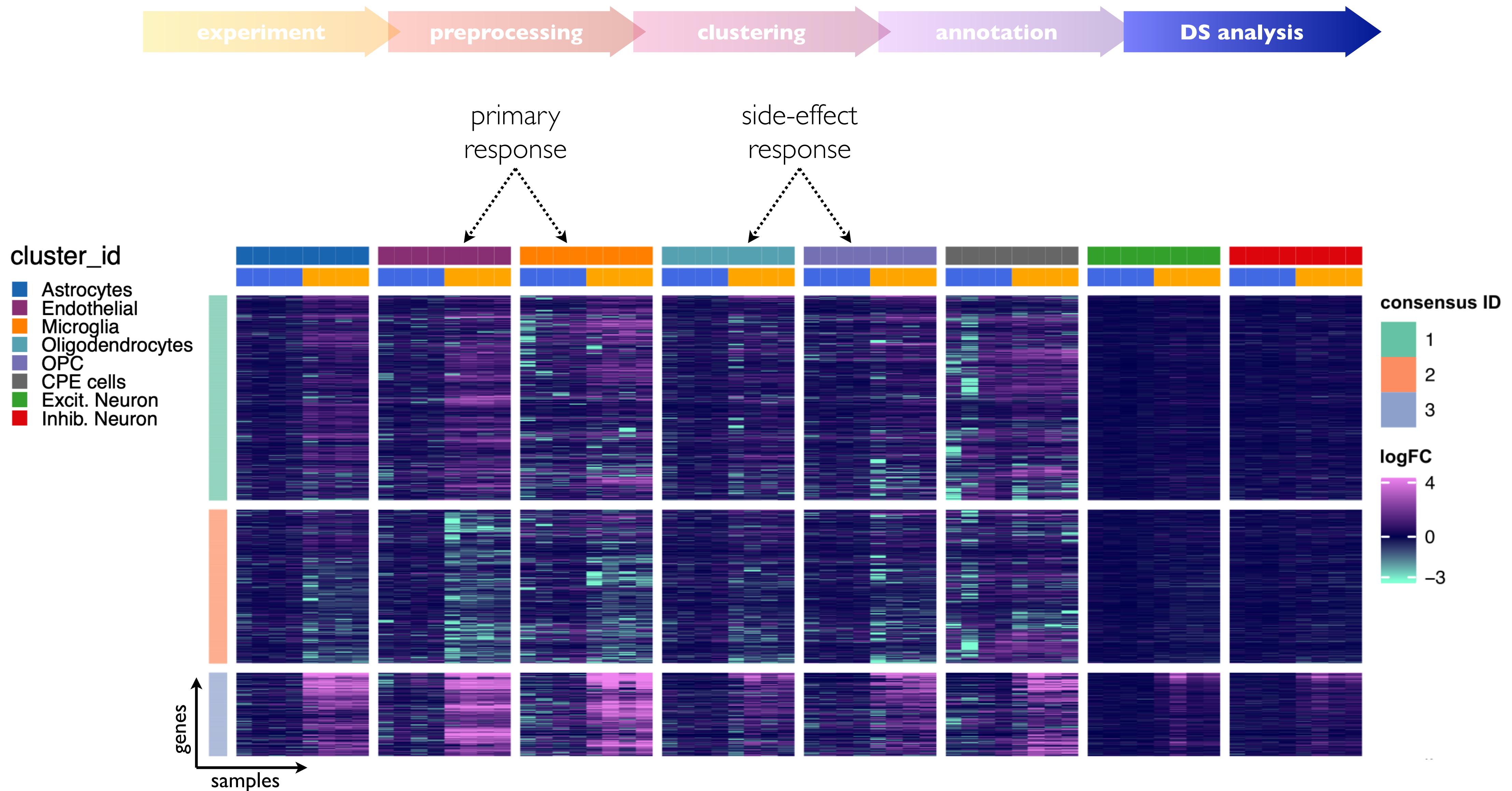
By integrating all samples together, we can now compare healthy and diabetic cells in matched cell states. To maximize statistical power, we want to use all cells - not just the sketched cells - to perform this analysis. As recommended by [Soneson et al.](#) and [Crowell et al.](#), we use an aggregation-based (pseudobulk) workflow. We aggregate all cells within the same cell type and sample using the `AggregateExpression` function. This returns a Seurat object where each 'cell' represents the pseudobulk profile of one cell type in one individual.

Seurat v1

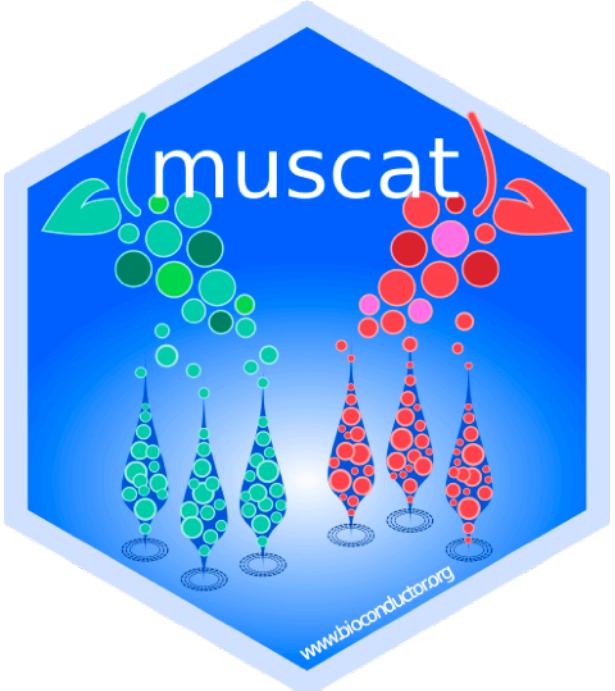
case study – vehicle vs. LPS treated mice



case study – vehicle vs. LPS treated mice



multi-sample multi-group scRNA-seq analysis tools



- Snakemake workflow (simulation study)
- browsable workflowr website (application to mouse cortex data)
- R/Bioconductor package
 - simulation framework
 - pseudobulk- & mixed model-based DS analysis methods
- differential detection (DD)
- bulk-based hypothesis weighting

***muscat* detects subpopulation-specific state transitions from multi-sample multi-condition single-cell transcriptomics data**

Helena L. Crowell^{1,2}, Charlotte Soneson^{1,2,3,6}, Pierre-Luc Germain^{1,4,6}, Daniela Calini⁵, Ludovic Collin⁵, Catarina Raposo⁵, Dheerai Malhotra⁵ & Mark D. Robinson^{1,2*}

Differential detection workflows for multi-sample single-cell RNA-seq data

Jeroen Gilis^{1,2†}, Laura Perin^{3†}, Milan Malfait¹, Helena L. Crowell⁴, Koen Van den Berge⁵, Alemu Takele Assefa⁵, Bie Verbiest⁵, Davide Risso^{3,6} and Lieven Clement^{1*}

The screenshot shows a GitHub repository page for 'muscat-comparison'. At the top, there are buttons for 'Code', 'Issues 1', 'Pull requests 0', 'Actions', 'Projects 0', 'Wiki', 'Security', 'Insights', and 'Settings'. Below this, a message says 'No description, website, or topics provided.' There is a 'Manage topics' button. A summary bar shows '179 commits', '2 branches', '0 packages', '0 releases', and '3 contributors'. A 'Branch: master' dropdown and a 'New pull request' button are also present. The main area displays a list of files and their commit history:

File	Commit Message	Date
CATALYST-project	add fig scripts; rmv old/unused scripts	Latest commit e551c1e on 6 Aug
MAGL	add fig scripts; rmv old/unused scripts	4 months ago
figs	add fig scripts; rmv old/unused scripts	4 months ago
scripts	add fig scripts; rmv old/unused scripts	4 months ago
.Renviron	updtt env	5 months ago
.gitignore	ignore ..	5 months ago
README.md	fix typo	4 months ago
Snakefile	add script for session info	4 months ago
config.yaml	update config	5 months ago

At the bottom, there are navigation tabs for 'Contents', 'Preprocessing', 'Clustering', 'Annotation', 'DS analysis', 'Visualization', 'Geneset analysis', and 'Downstream'.

Contents

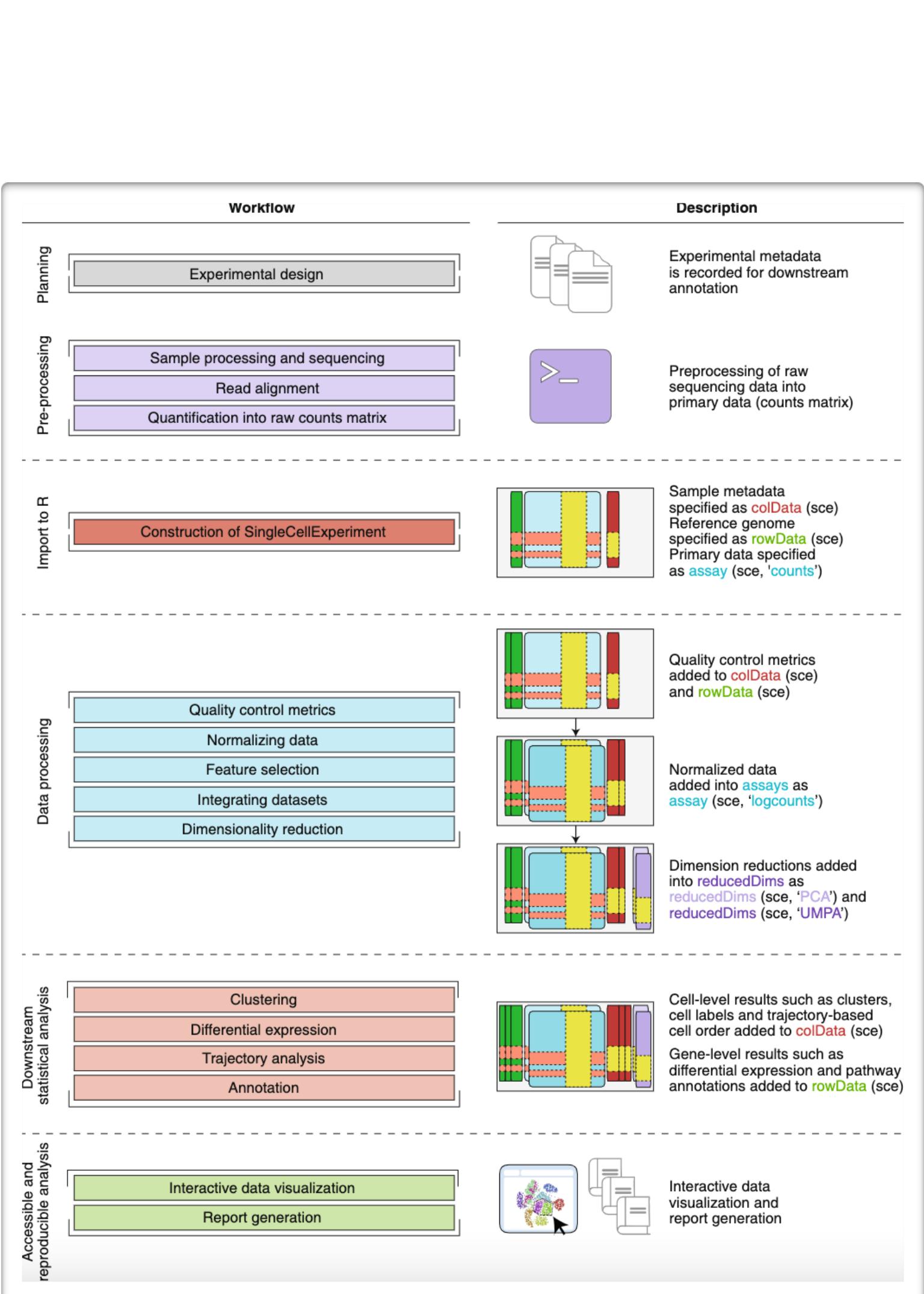
workflow

- Preprocessing:
 - Calculation of QC metrics & identification of outliers using `scater`
 - Filtering of genes and cells
- Clustering:
 - Integration & clustering using a sequence of resolutions
 - t-SNE and UMAP dimension reductions colored by sample, group, and cluster ID
- Annotation:
 - Number of clusters by resolution
 - Number of cells by cluster-sample
 - Relative cluster abundances by sample
 - t-SNE colored by expression of known marker genes
 - Heatmap of mean known marker-gene expressions by cluster
 - Identification of cluster-markers using `scran`
 - Hetamap of mean `scran` marker-gene expressions by cluster
- DS analysis
 - Cluster annotation
 - Aggregation to pseudobulk counts
 - Pseudobulk-level MDS plot
 - Cluster-level DE analysis with `edgeR`
 - Results filtering & overview
 - Dimension reduction: t-SNE & UMAP

OSCA – Orchestrating Single-Cell Analysis with Bioconductor

- comprehensive online book
“covering installation, sources of help, specialized topics pertaining to specific aspects of scRNA-seq analysis and complete workflows [...]”

The screenshot shows the landing page of the online book. It features a sidebar with navigation links like Welcome, What you will learn, and What you won't learn. The main content area has a title 'Orchestrating Single-Cell Analysis with Bioconductor' and a bioRxiv-style abstract. It also includes sections for 'Welcome', 'What you will learn', and 'What you won't learn'.



on differential discovery in scRNA-seq data

...using muscat & miloDE/lemur

Helena L Crowell, PhD
Autumn School for Single Cell-ers
Oct 21, 2025 · GIMM, Oeiras, Portugal

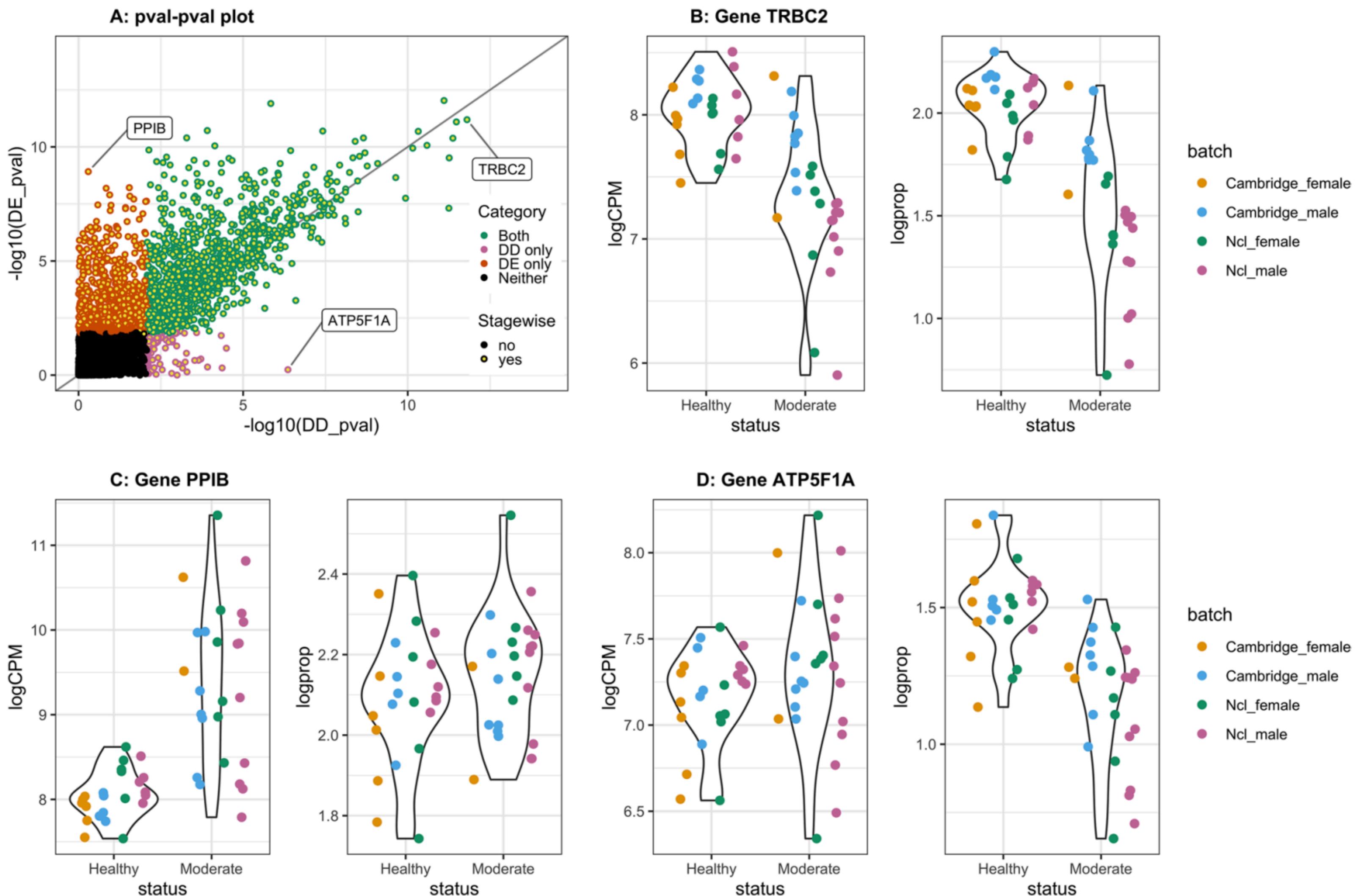
CNAG – National Center
for Genomic Analysis
Barcelona, Spain



**Swiss National
Science Foundation**

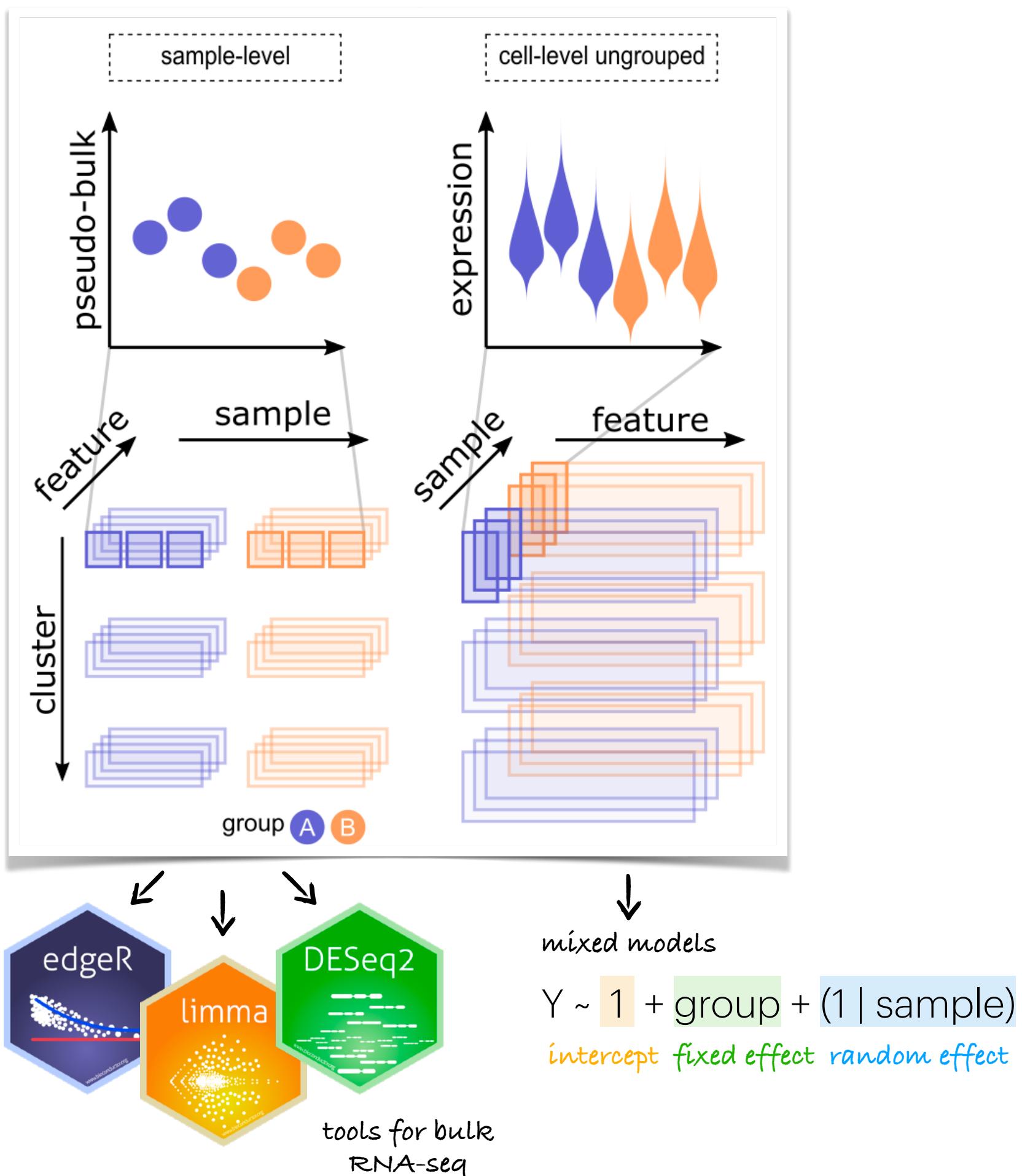
beyond expression – differential detection (DD)

One particularly interesting characteristic of gene expression not explicitly captured by the aforementioned frameworks is differential detection (DD), i.e. finding differences in the fraction of cells in which a gene is detected between groups. It has been reported that gene expression profiles may exhibit characteristic bimodal expression patterns, in which the expression of genes is either strongly positive or undetected within individual cells (e.g., Finak et al. [4]). At the single-cell level, such differences in gene detection may arise from technical artifacts or from the stochastic nature of gene expression [5], a phenomenon commonly referred to as transcriptional bursting [6]. At the tissue level, differential detection can be biologically meaningful as well. For instance, even if the overall expression level of a gene in a tissue is equal between two biological conditions, it is relevant to assess the fraction of cells in the tissue contributing to the observed expression levels. In addition, previous research has shown that assessing differences in detection may provide a more robust means of studying gene expression than assessing differences in the average gene expression.



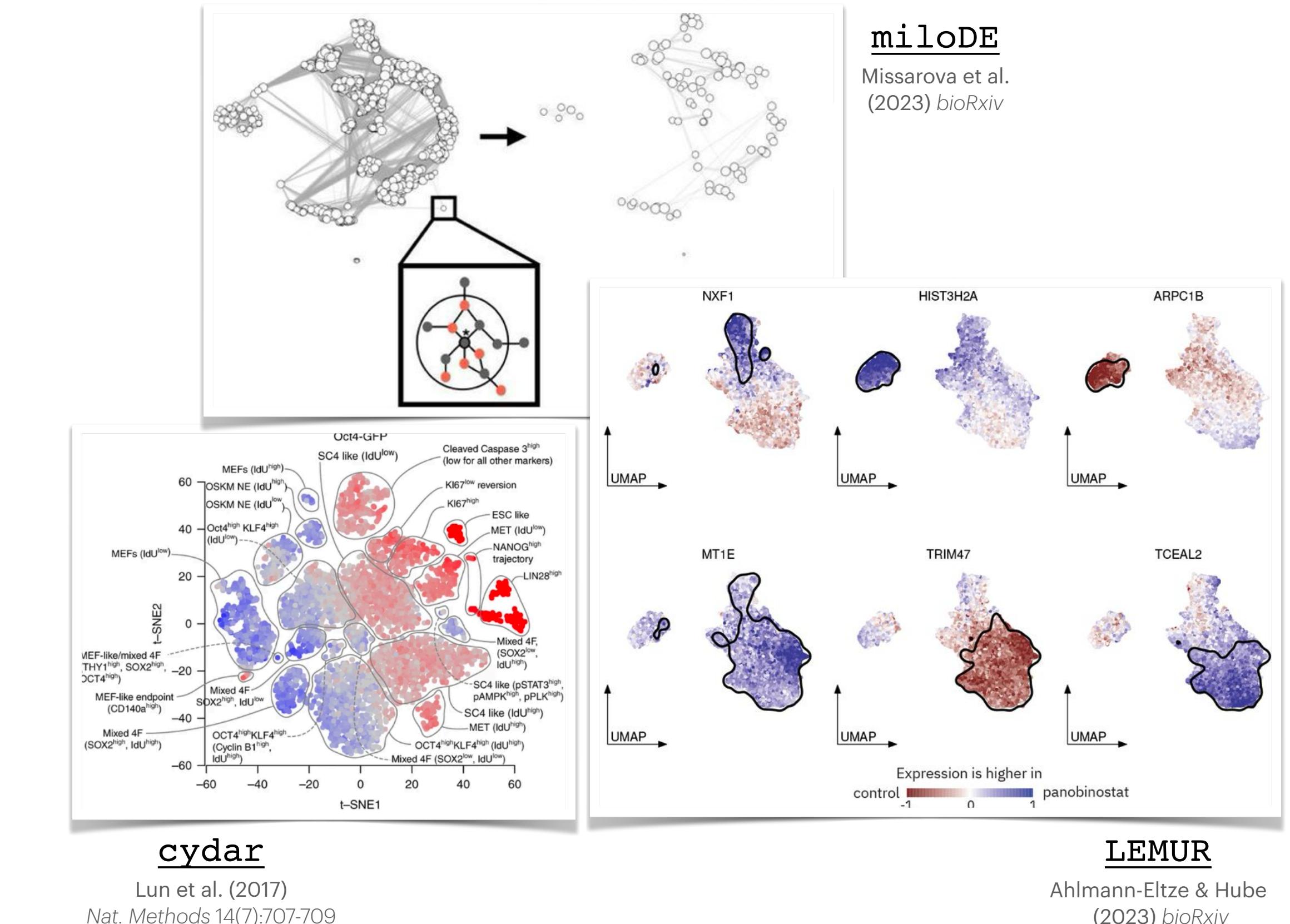
more on DA/S analysis – cluster-based vs. -free approaches

Crowell et al. (2019) *Nat. Commun.* 11(1):6077



a priori grouping of cells

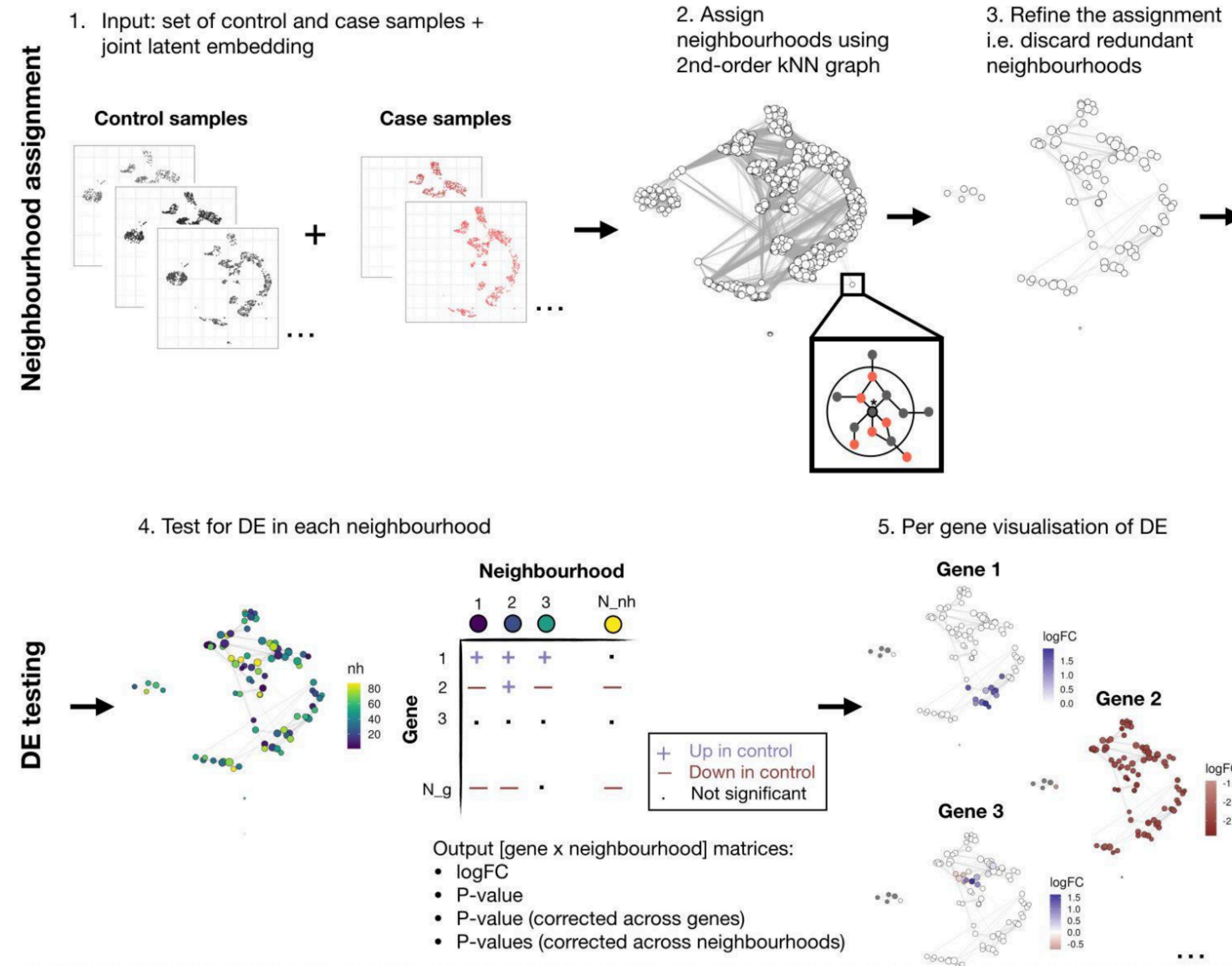
- clustering is not easy...
- ...but more interpretable



interesting neighborhoods are identified in a data-driven way

- no need for clustering...
- ...but need to make (biological) sense of results downstream

a closer look at miloDE

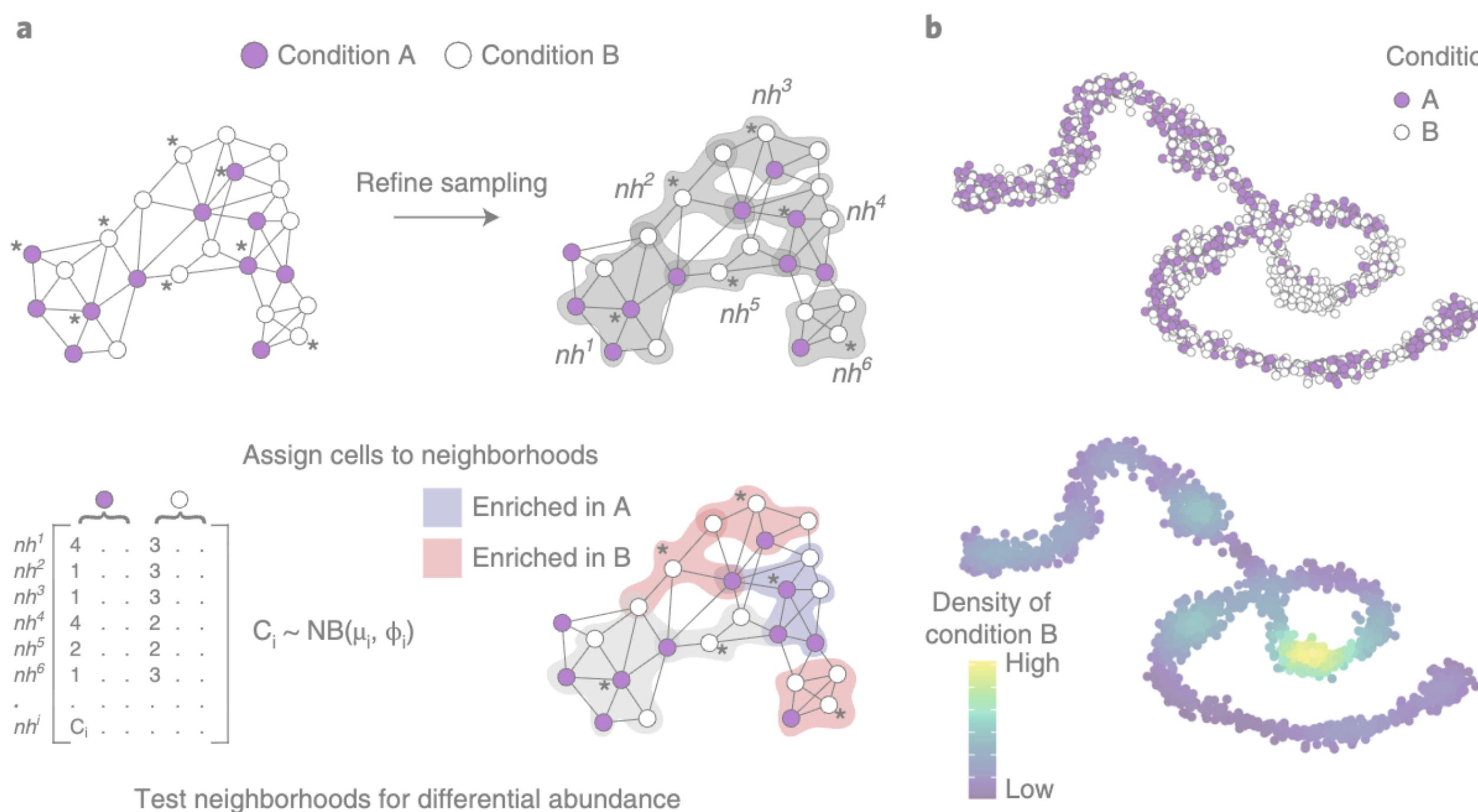


before miloDE, there was milo(DA) – interpretability?

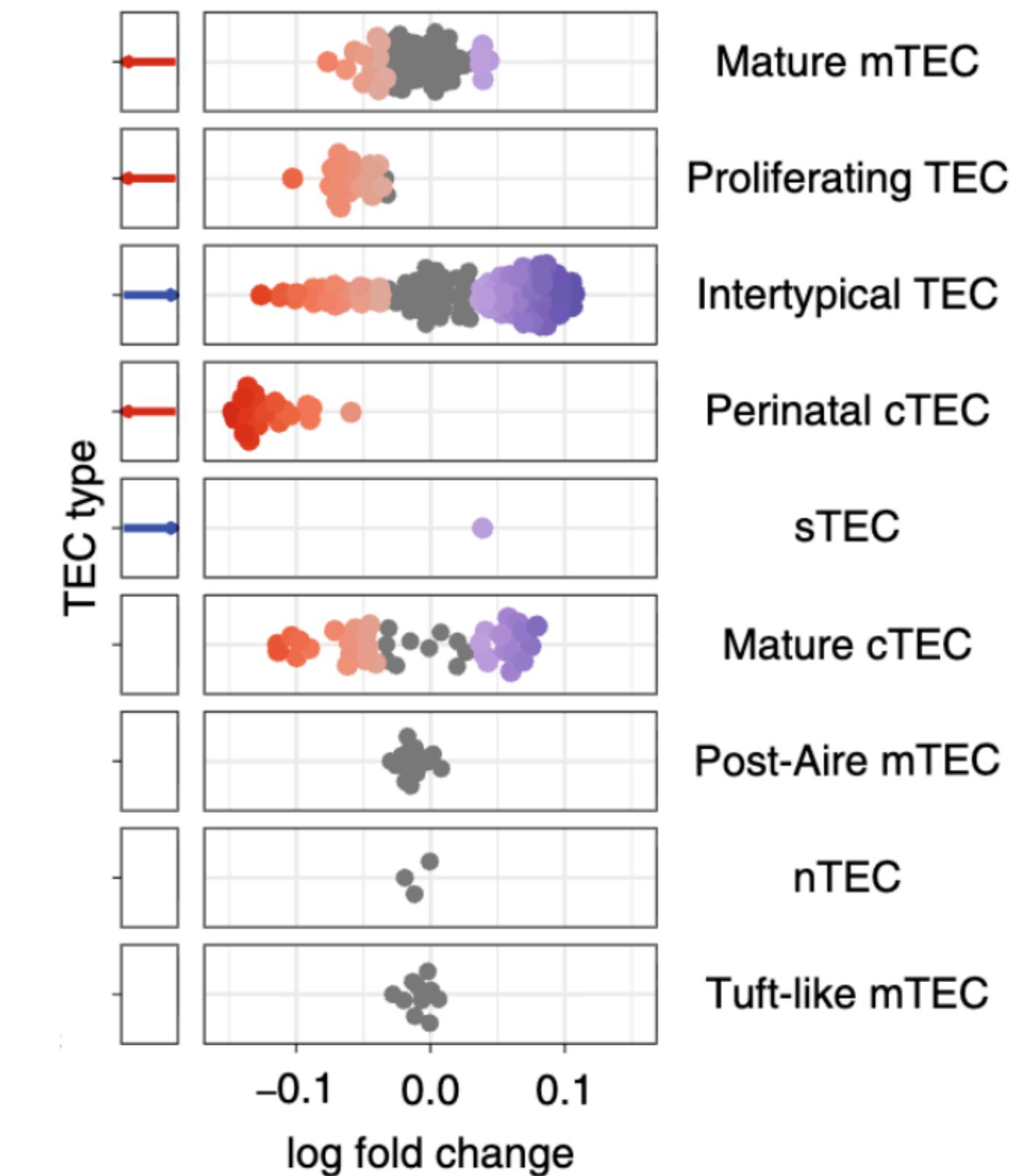
for miloDE: we have this \times the number of genes!

(point = cell, facet = cluster, color = logFC)

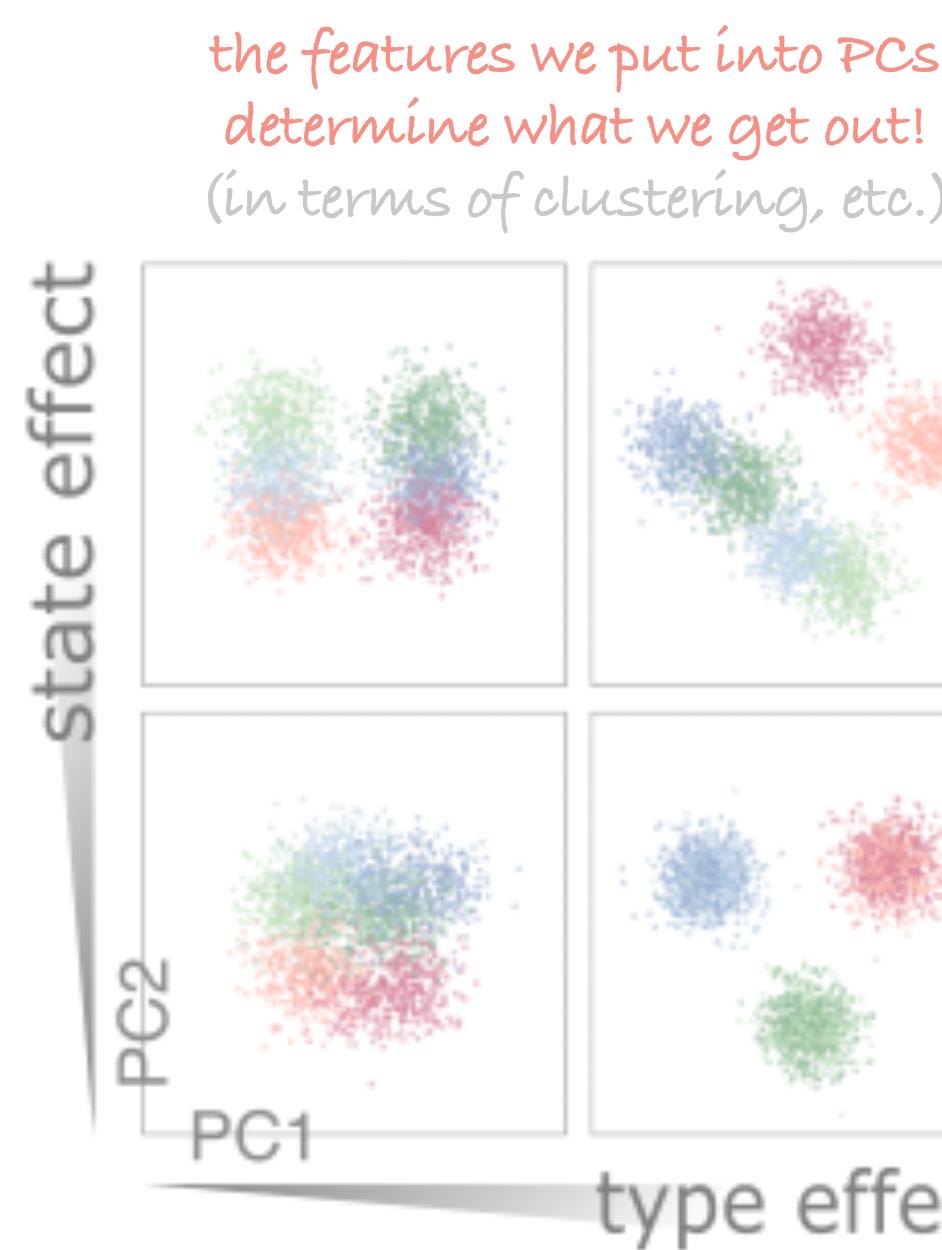
Baran-Gale et al. DA direction
 Enriched with age
 Depleted with age



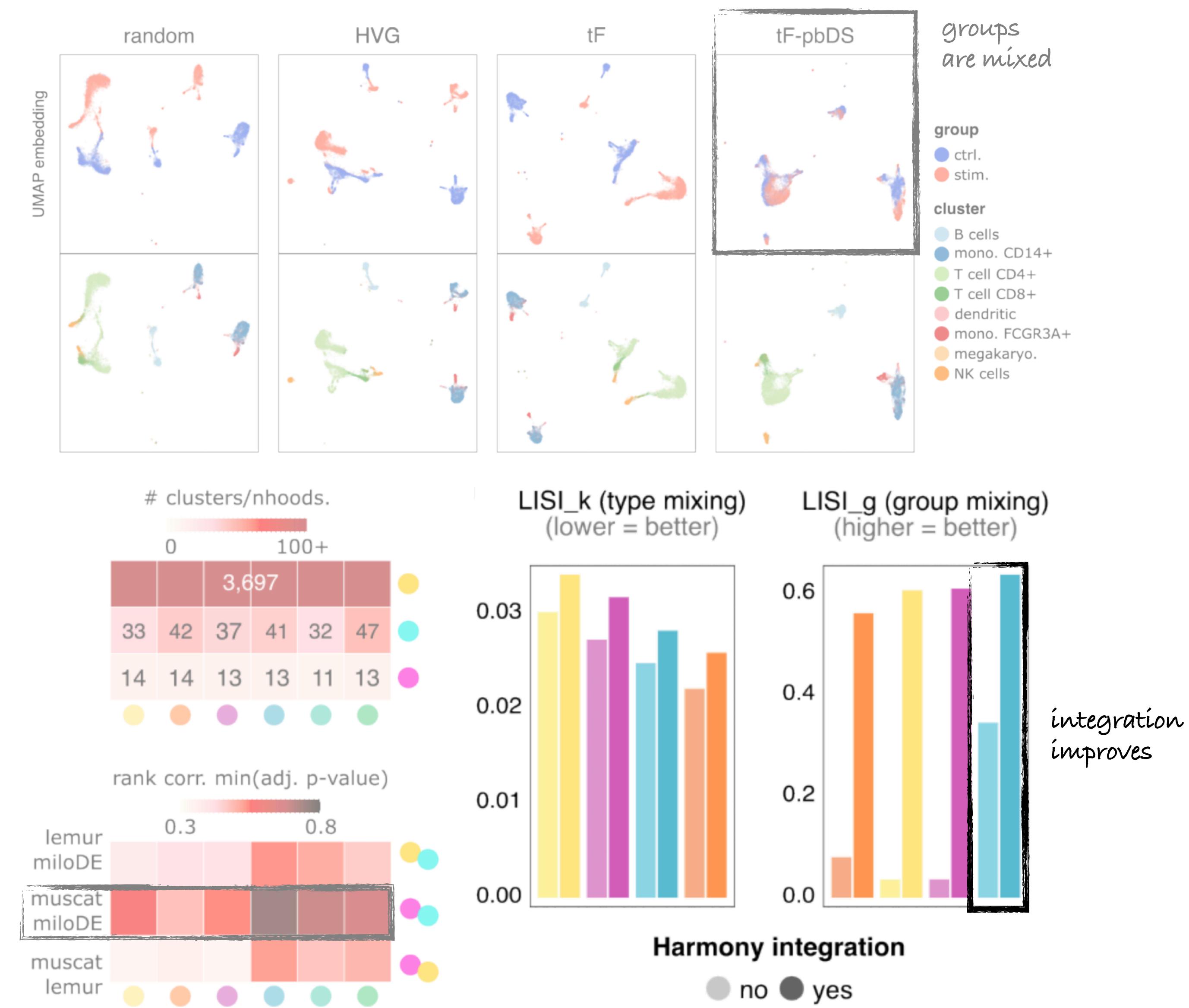
stratification of cell-level results by subpopulation
→



feature selection can disentangle cell type & state



cluster-free δ -based
DS methods comparable



take home message / Gedankenexperiment

There is arguably an interplay between not only the type and state transcriptional programs governing cells, but also between the various types of differential analysis that are applied. While the readily-available DAA and DSA tools can in principle be considered orthogonal, they are often linked by how DE drives subpopulation definition^[17]. As a toy example, consider a substance that strongly increases the expression of a gene γ in all treated (but not control) cells. If γ were used for clustering (or neighborhood generation), γ^- and γ^+ subpopulations should be obtained for each cell type; and, DAA analysis would detect γ^+ and γ^- populations as differentially abundant. In contrast, if γ were *not* contributing to the subpopulation definition, γ^+ and γ^- cells should be present in every cluster (whose markers are expressed independently from γ), and a DSA analysis should detect γ as higher expressed in treated cells (i.e., a differential state) for each cell type.