Swiss Institute of Bioinformatics

INTRODUCTION TO SEQUENCING-BASED TRANSCRIPTOMICS DATA ANALYSIS

# Clustering

**Joana Carlevaro Fita**

December 9-10, 2025
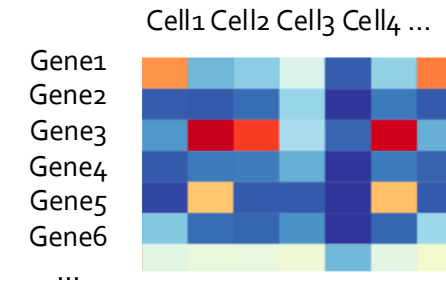
ELIXIR SWITZERLAND

# Standard scRNAseq methods applied

Non-spatially aware clustering

Measure

Cell1 Cell2 Cell3 Cell4 ...

Gene1
Gene2
Gene3
Gene4
Gene5
Gene6
...
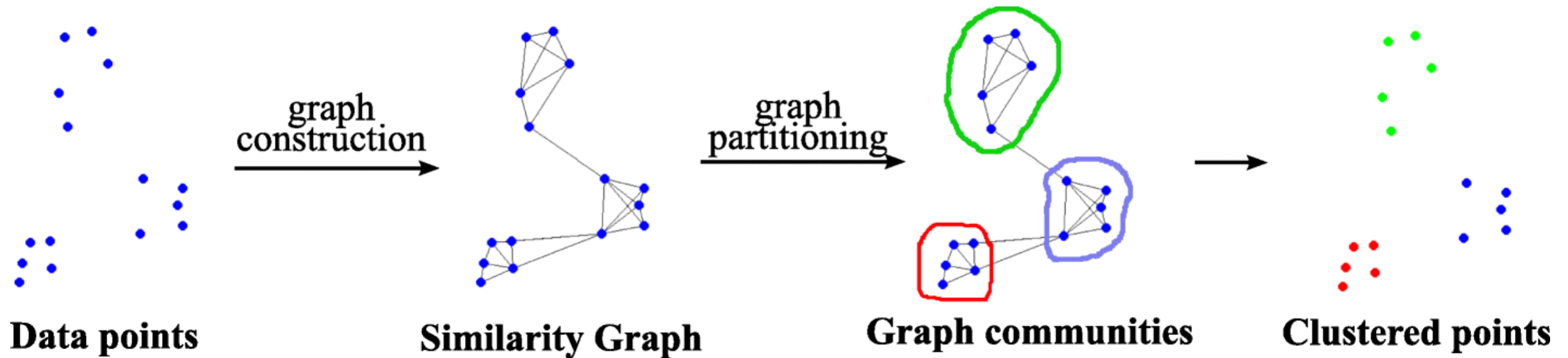
# Standard scRNAseq methods applied

**Graph-based clustering**  Cells within the same community are assigned to the same cluster



Spatially unaware graph-based clustering: based on a shared nearest neighbor (SNN) graph and the Leiden or Louvain algorithm for community detection
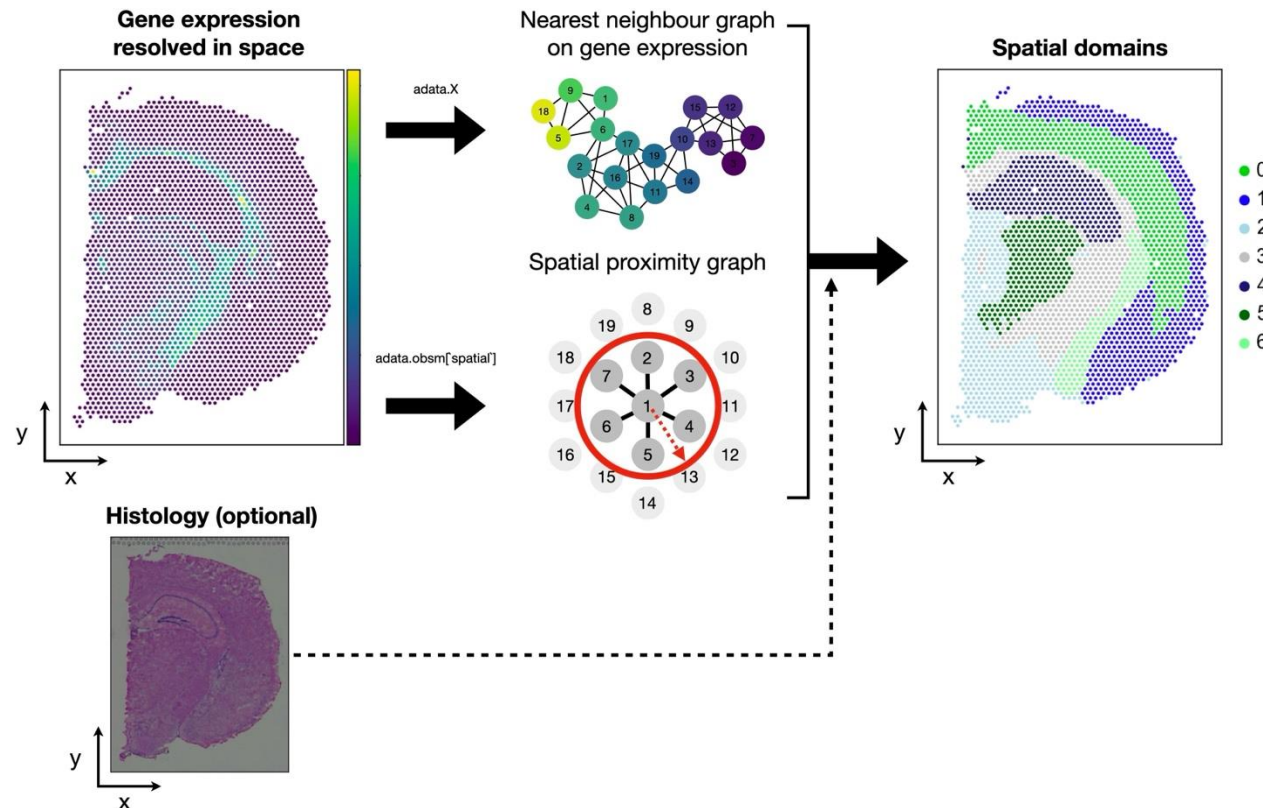
# Clustering

 Non-spatially aware clustering ➡️
- Use only gene expression
- Can give discontinous results, they don't use information about neighbours
- Can be coupled to spatially-aware pre-processing methods

 Spatially aware clustering



Cell type clustering vs tissue domain identification
Different algorithmic problem

Heumos, L. et. Al. Nat Rev Genet, 2023

# Spatially-aware clustering

Spatial domains can be identified based on morphology/biological knowladge or using clustering (spatially aware methods).
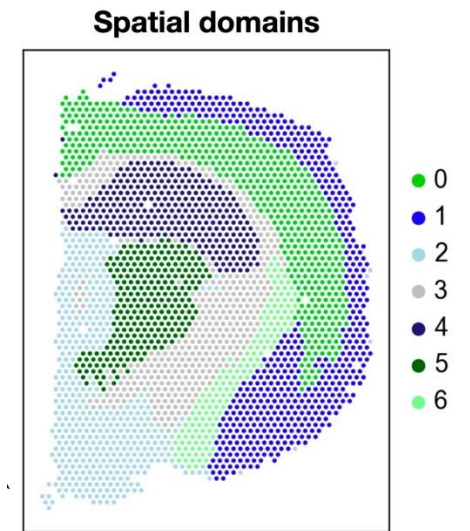
Several methods developed for **spatial domain identification**

Goals
- Use gene expression **and** spatial information
- Define spatial regions with similar spatial gene expression / cell composition --> downstream analysis (ie. Annotation, DE...)
- Understand biological processes in terms of gene expression across spatial localisation

Assumptions
- Cell type / cell state is influenced by interactions with **neighbouring** cells.
- Spatial domains may be composed of different cell types
- Cells from same type/state may be spatially far a part (ie. Cerebral hemispheres, epithelial layers, blood vessels...)

**Spatial domains**

- 0
- 1
- 2
- 3
- 4
- 5
- 6

# Spatially-aware clustering methods

Methods have different assumptions, methodologies and computational trade-offs

Following OSTA book classification (Crowell et al. Biorxiv, 2025) :

- **Probabilistic** (HMRF, BayesSpace)
Examine each cell and its surrounding cell's expression to define domains
Encourages neighbouring spots to have the same label

- **Encoder-based** (CellCharter, STAGATE)
Uses encoder architectures to generate latent embedding
Model spatial neighbourhoods as graphs, or jointly model gene expression and spatial coordinates.
May be a good fit for imaging-based data.

- **Neighbourhood-based** (BANKSY)
Use augmented features and embed cells in a product space containing information on the own cell and its local microenvironment. Clustered using standard algorithms
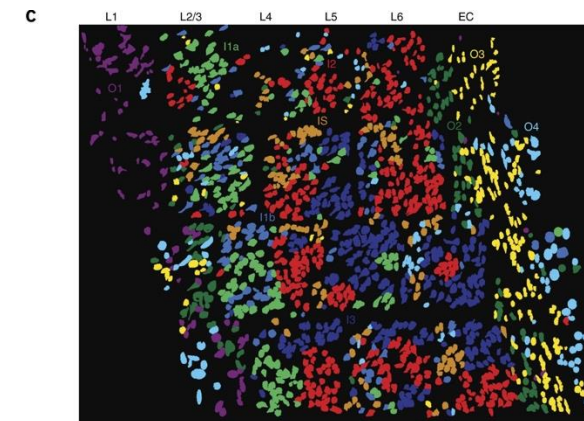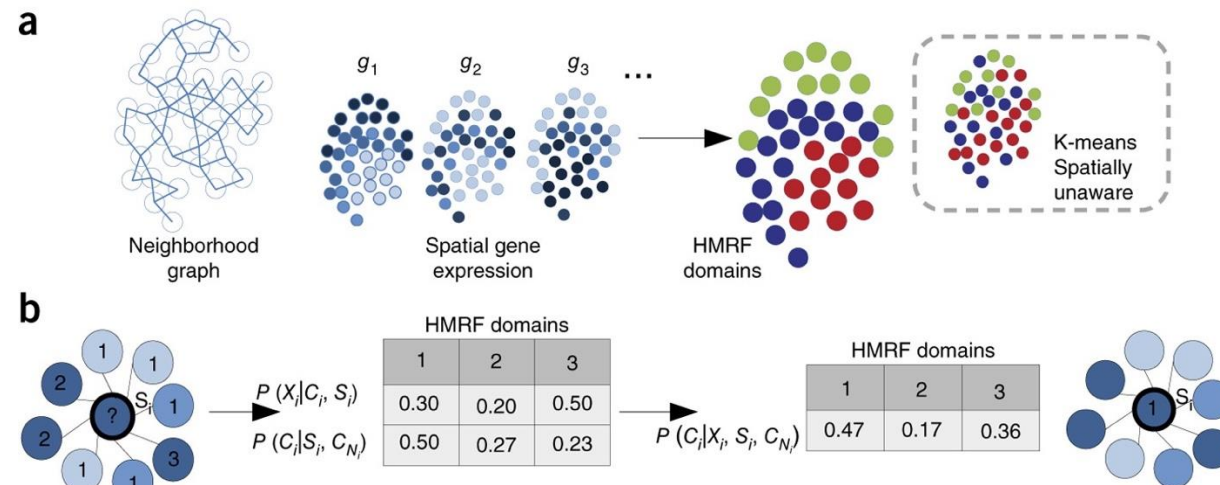
# Spatially-aware clustering methods

## Hidden Markov Random Field (HMRF)

Method: Graph-based model commonly used for pattern recognition in image data analyses

Assumption: Cell type / cell state is influenced by (interactions with) **neighbouring** cells.

Intuition: Examine each cell and its surrounding cell's expression to define domains

# BayesSpace

Intuition:

"BayesSpace enables spatial clustering by modelling a low-dimensional representation of the gene expression matrix and encouraging neighbouring spots to belong to the same cluster via a spatial prior"
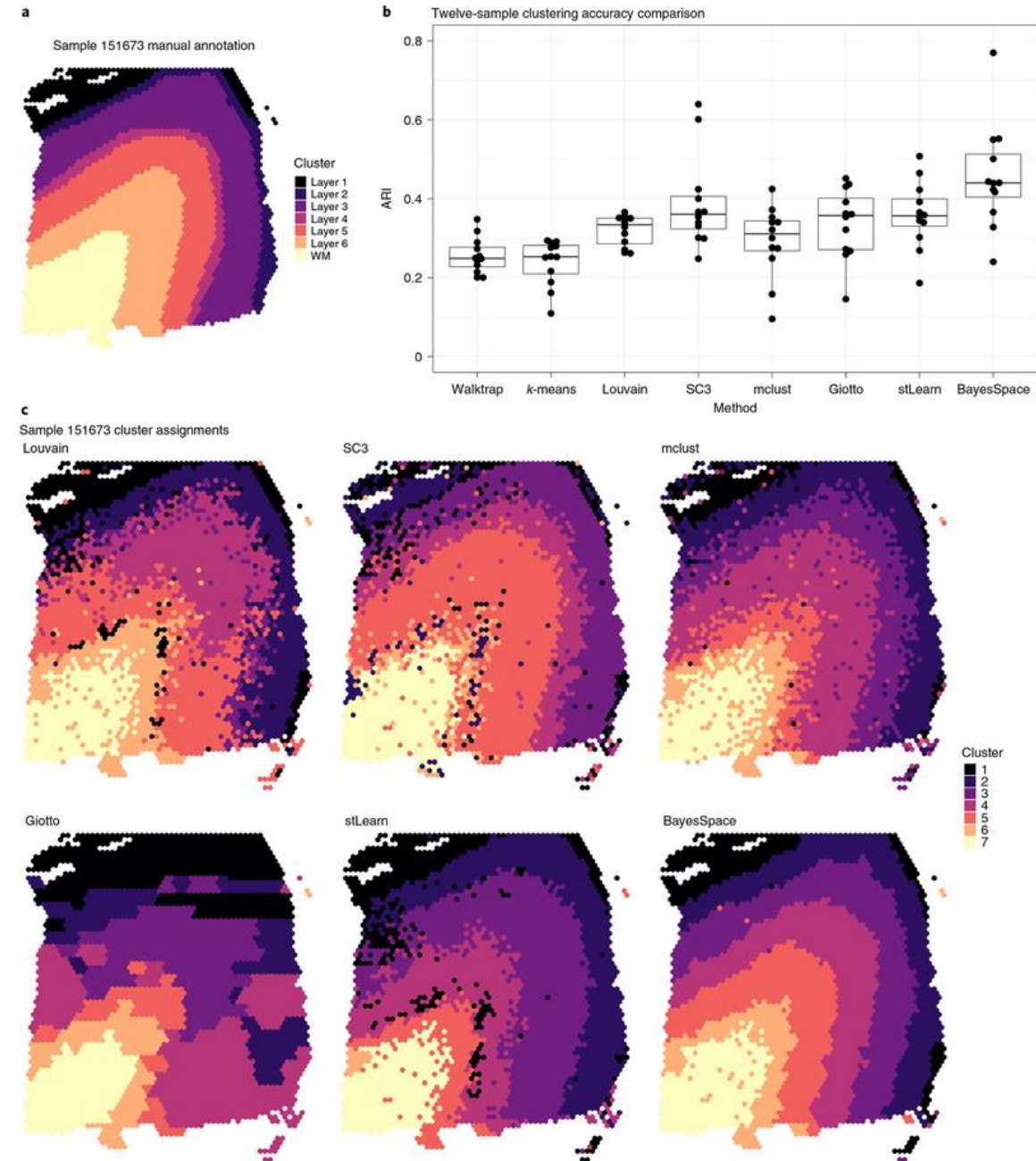
Assumptions:

Nieghbouring cells are more likely to have similar transcriptomes

Doesn't require pre-selection of marker genes

Goals:

Doesn't require independent scRNAseq datasets to perform deconvolution

Adress noise and sparcity getting a smoother separation (good domain segmentation)

# BANKSY

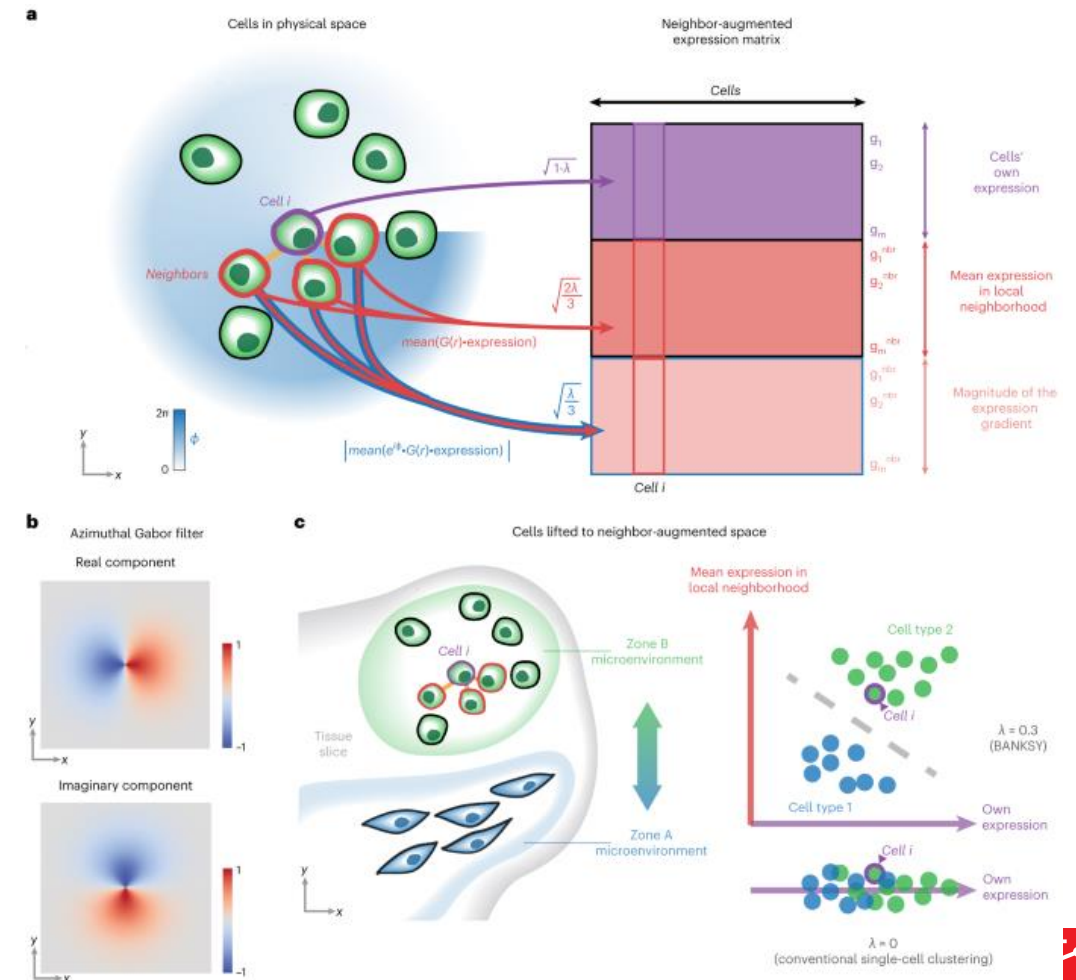## Building Aggregates with a Neighborhood Kernel and Spatial Yardstick (BANKSY)

Intuition:

"BANKSY uses a pair of spatial kernels to encode the transcriptomic texture of the microenvironment, one constructed using the weighted mean of gene expression in each cell's neighborhood and the other using an azimuthal Gabor filter (AGF)

Assumes

- A cell's transcriptome doesn't necessarily resemble the average transcriptome of its domain
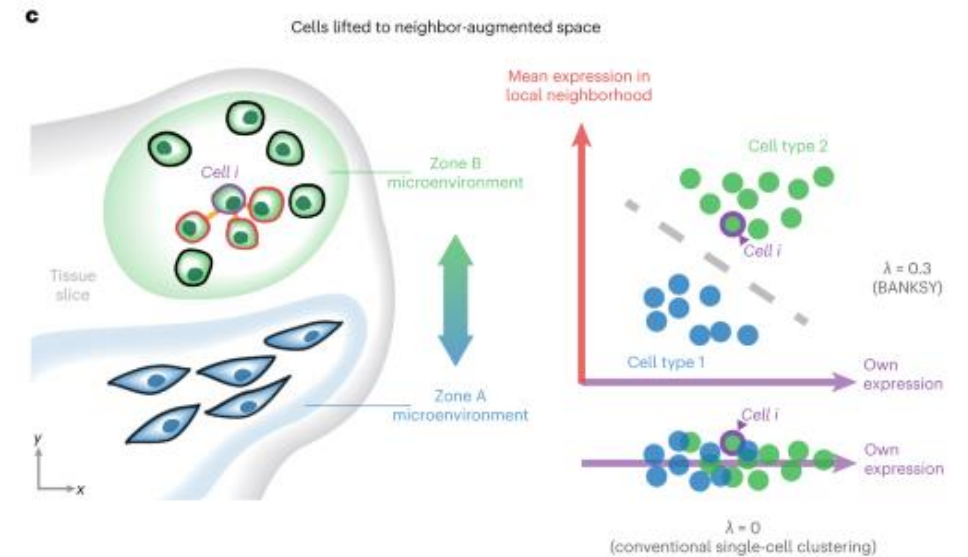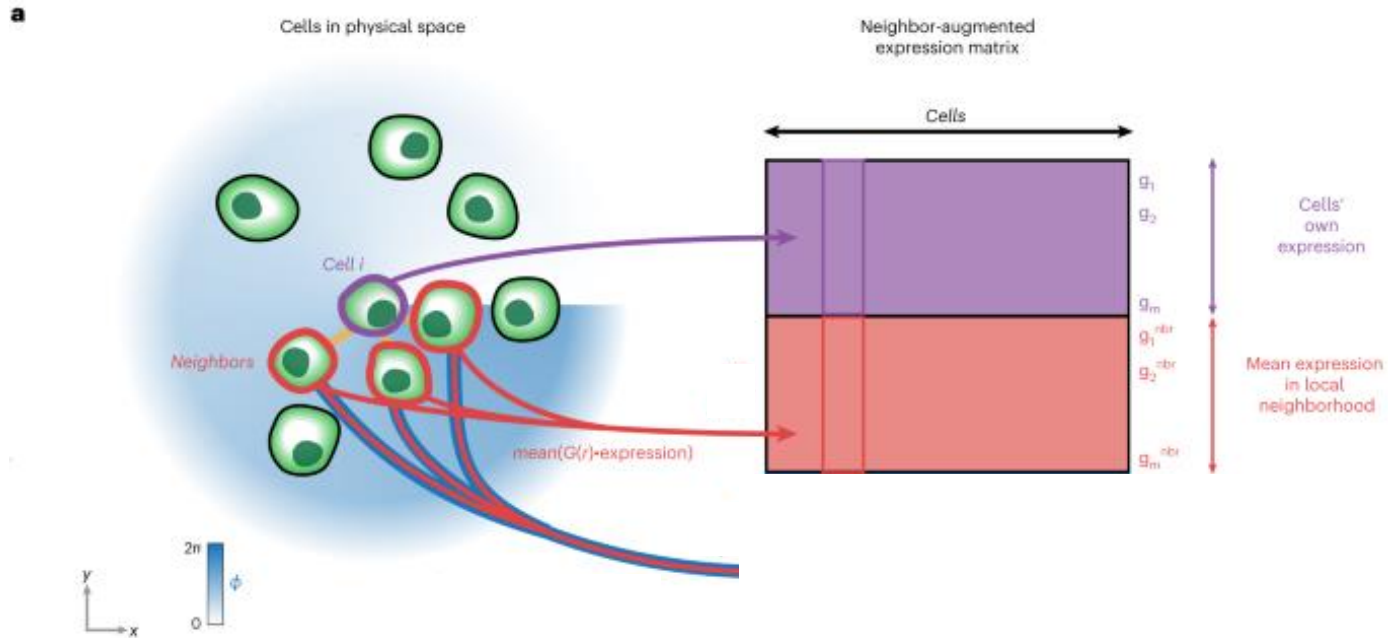
Goals:

Take into account that similar cells may be far apart (intermingled, repeated patterns..) → solve both cell typing and domain segmentation
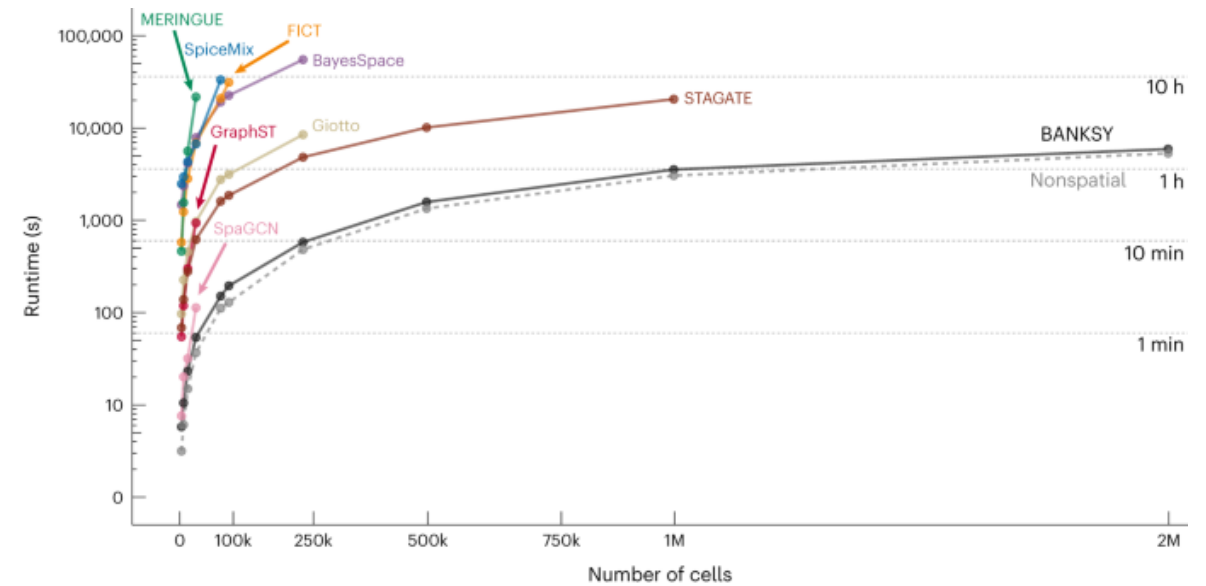
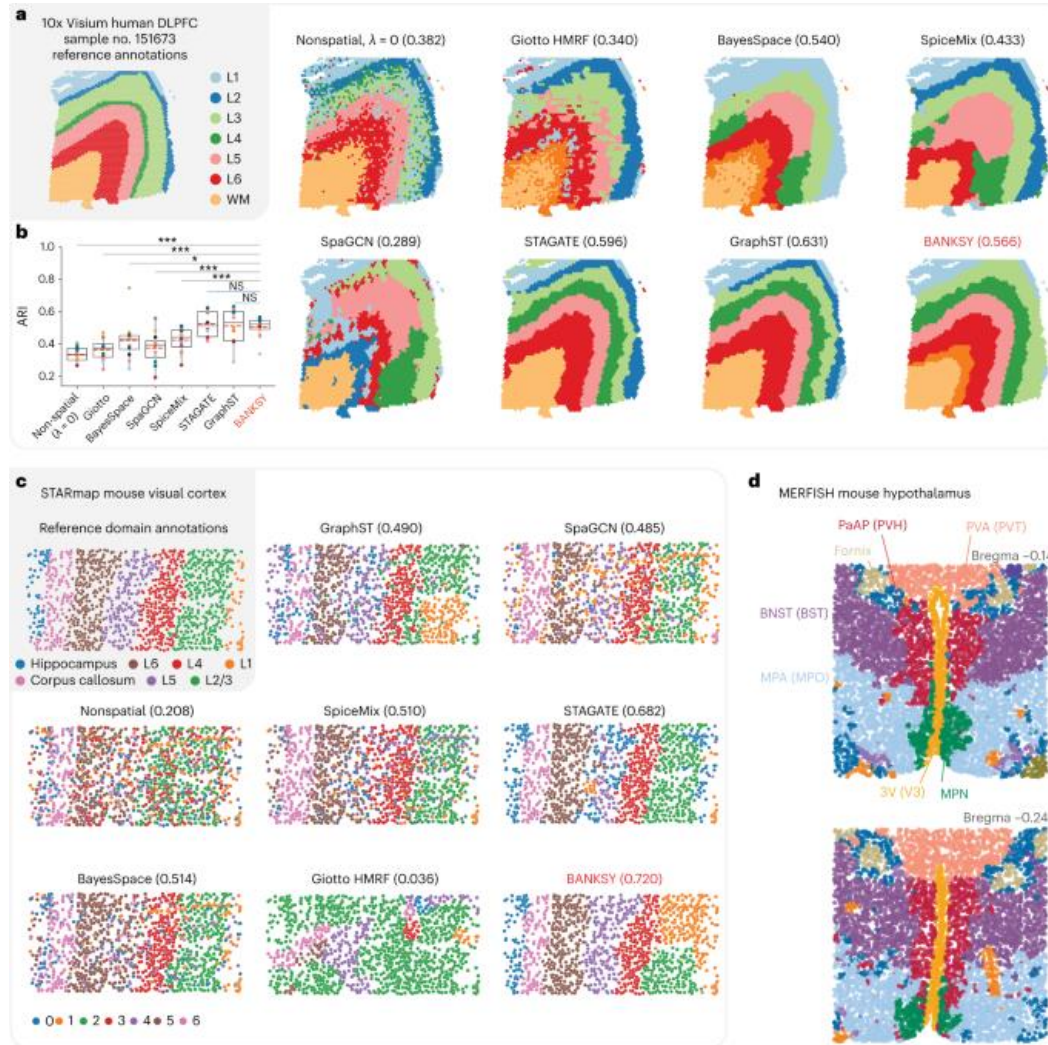

Singhal et al. Nat Gen, 2024

# BANKSY

## Building Aggregates with a Neighborhood Kernel and Spatial Yardstick (BANKSY)

# BANKSY



Singhal et al. Nat Gen, 2024

# Benchmark on methods



Yuan et al. Nature Methods, 2024

single cells

Single-cell transcriptomic data

gene expression matrix

HVGs — gene expression

non-HVGs

cells

Highly variable genes (HVGs)

cell clusters

Cell clustering

gene expression

Cluster-marker gene

spatial spots

Spatially resolved transcriptomic data

SVG   non-SVG   gene expression

Spatially variable genes (SVGs) "Overall SVGs"

spatial domains

Spatial domain identification

gene expression
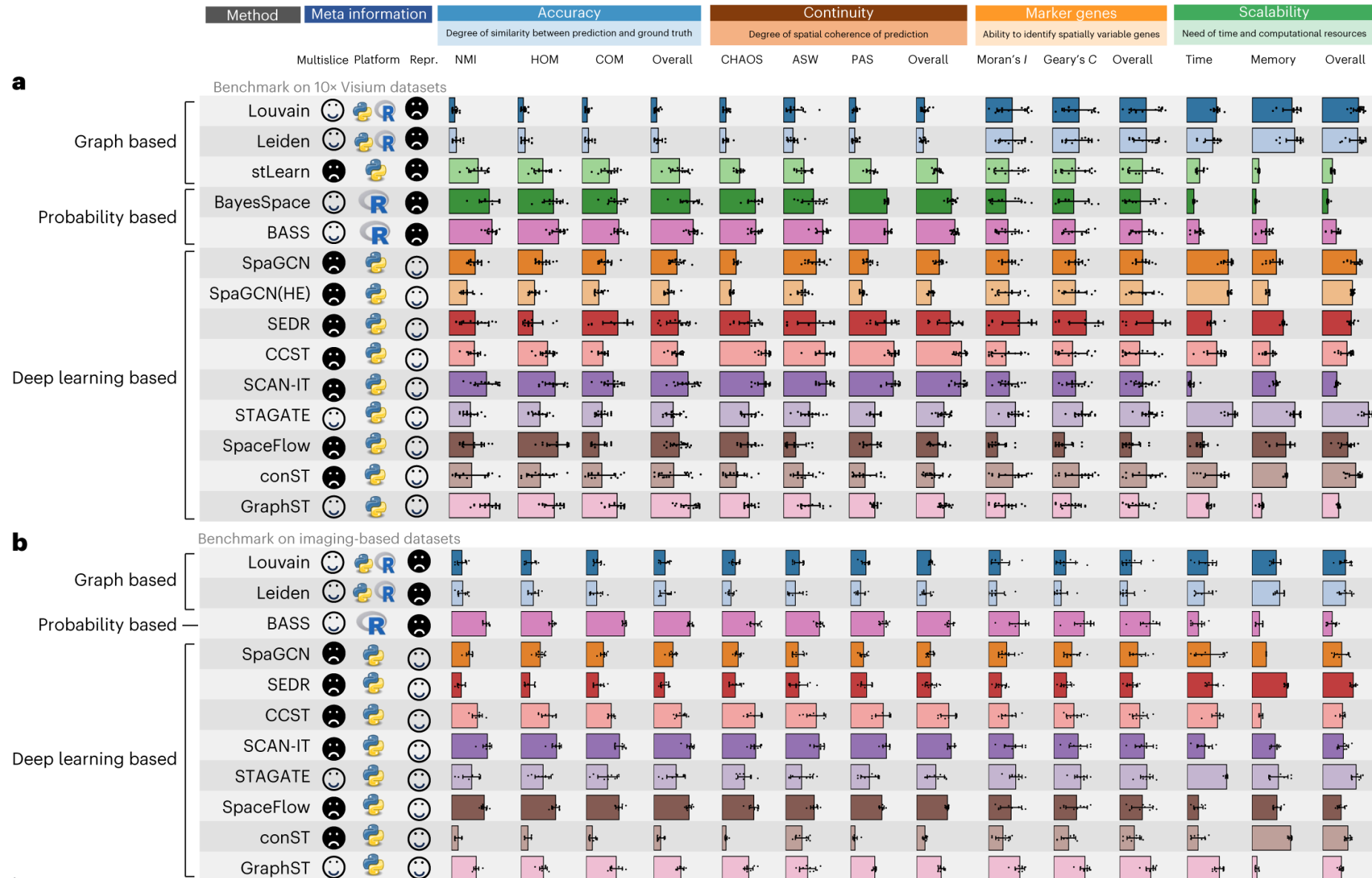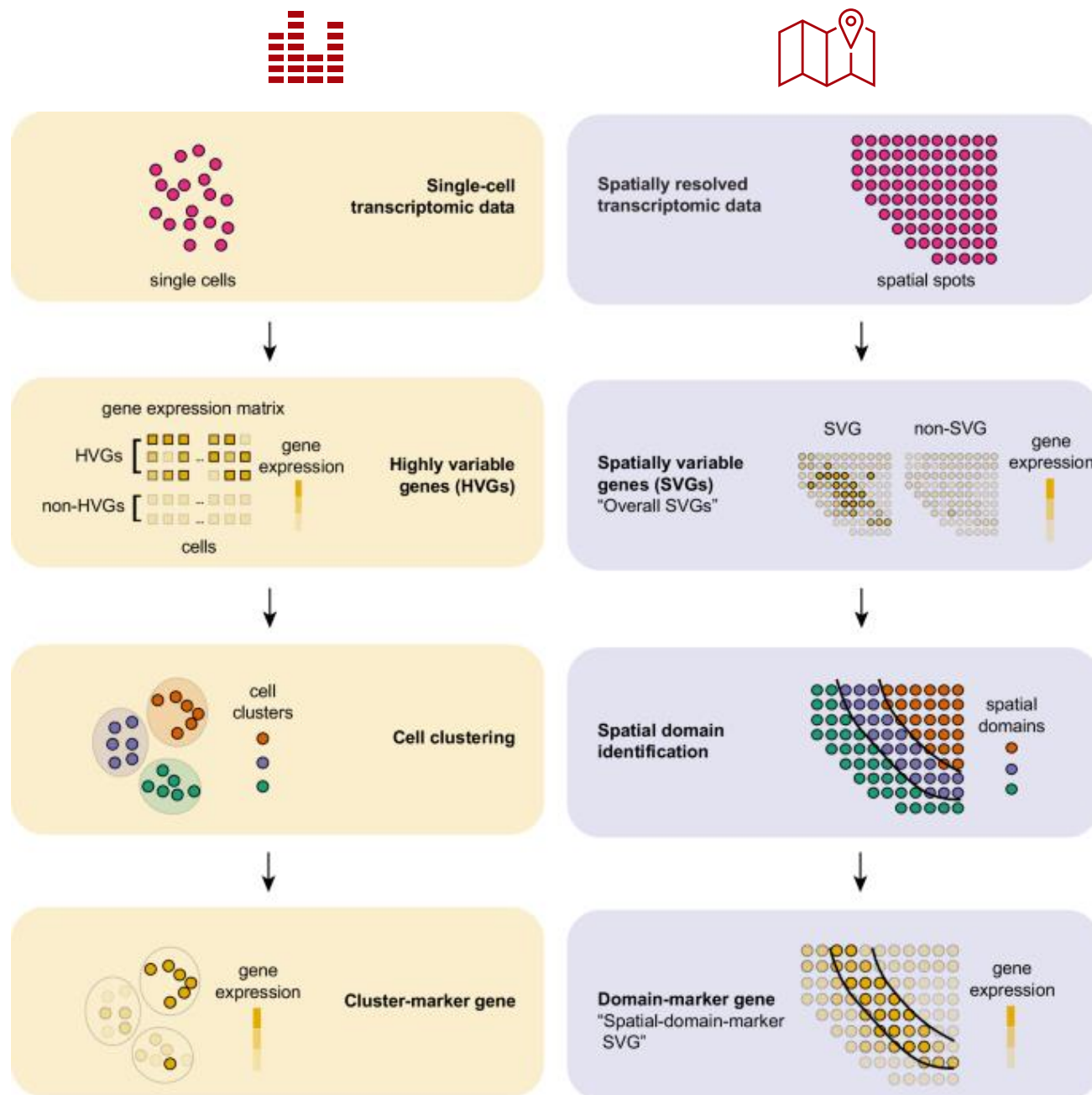
Domain-marker gene "Spatial-domain-marker SVG"

Yan, Hua, Li. Nat Comm,, 2025
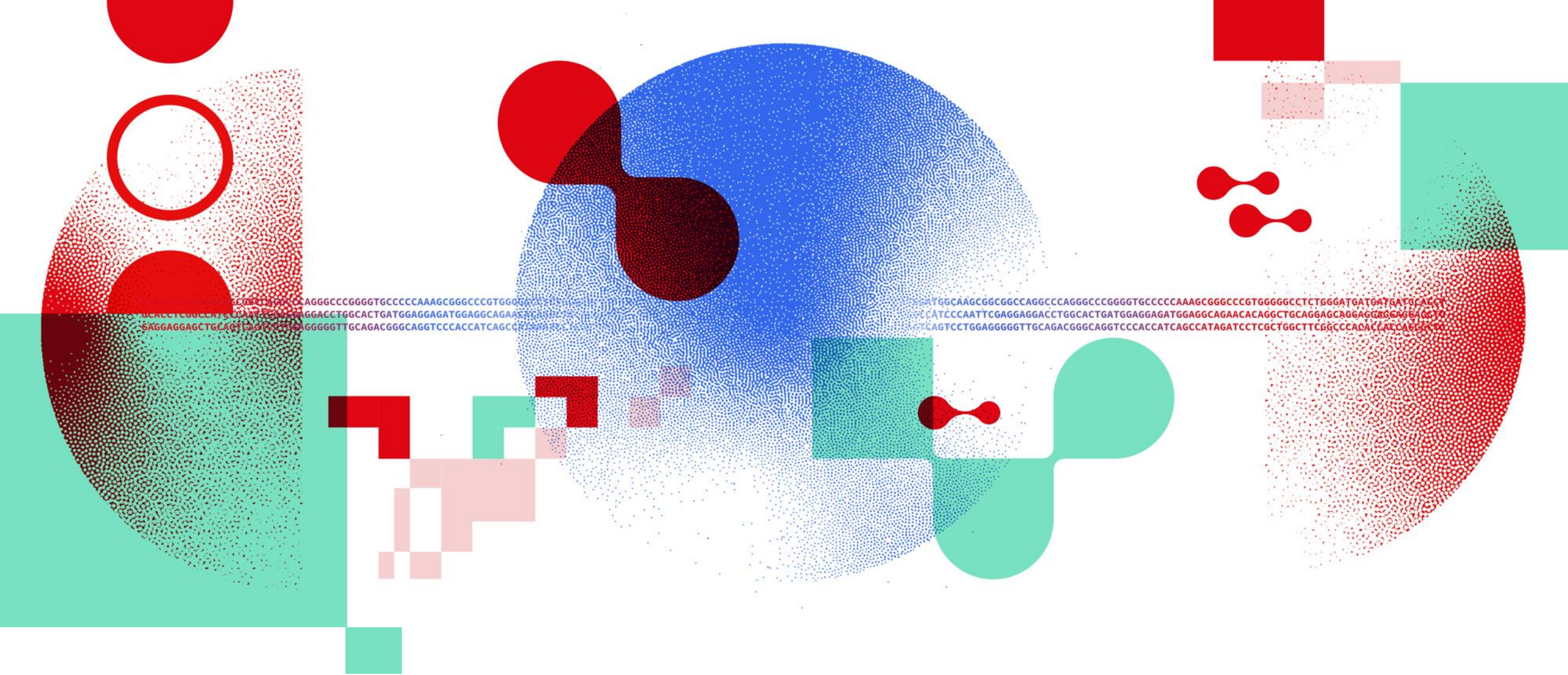
# Conclusions

- Algorithms have been developed to incorporate spatial information in the process of clustering cells into spatial domains.

- Graph neural networks has integrated tissue images to improve the performance using anatomical information.


- Challenges remain
    - Methods tested in few datasets, mostly human dorsolateral prefrontal cortex (DLPFC) → problem with over-fitting
    - Scalability (tested in small datasets)

# Thank you

DATA SCIENTISTS FOR LIFE

sib.swiss