# Statistical methods for spatial omics data

- Overview on the technologies (review)

- Finding spatially-variable genes

- Deconvoluting low-resolution (or aggregating high-resolution) spatial omics data

- Spatially-aware dimension reduction / clustering

- Cell-cell communication —> co-localization

- Classical spatial statistics

  ‣ Point patterns: random, clustered, intensity/correlation

  ‣ Lattice data: useful summaries / functions

  ‣ models with spatially correlated errors
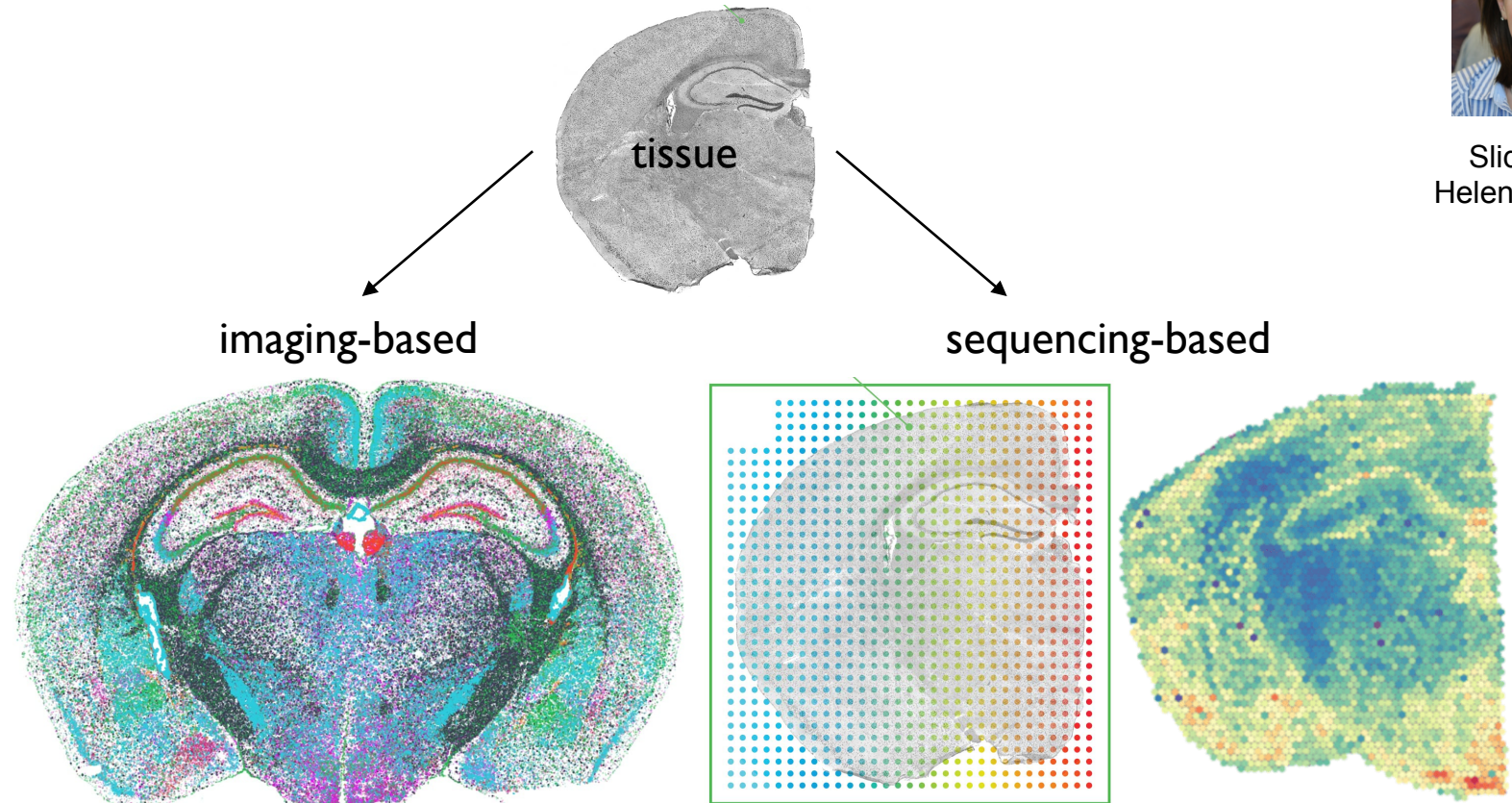
Slide from
Helena Crowell

bulk

single-cell

spatial

tissue

imaging-based

sequencing-based

- molecule-level data
- targeted panel (100s of features; >2024: 1000s)
- single-cell resolution requires segmentation

- spot-level data
- whole transcriptome (10,000s of features)
- single-cell resolutions requires aggregation or deconvolution

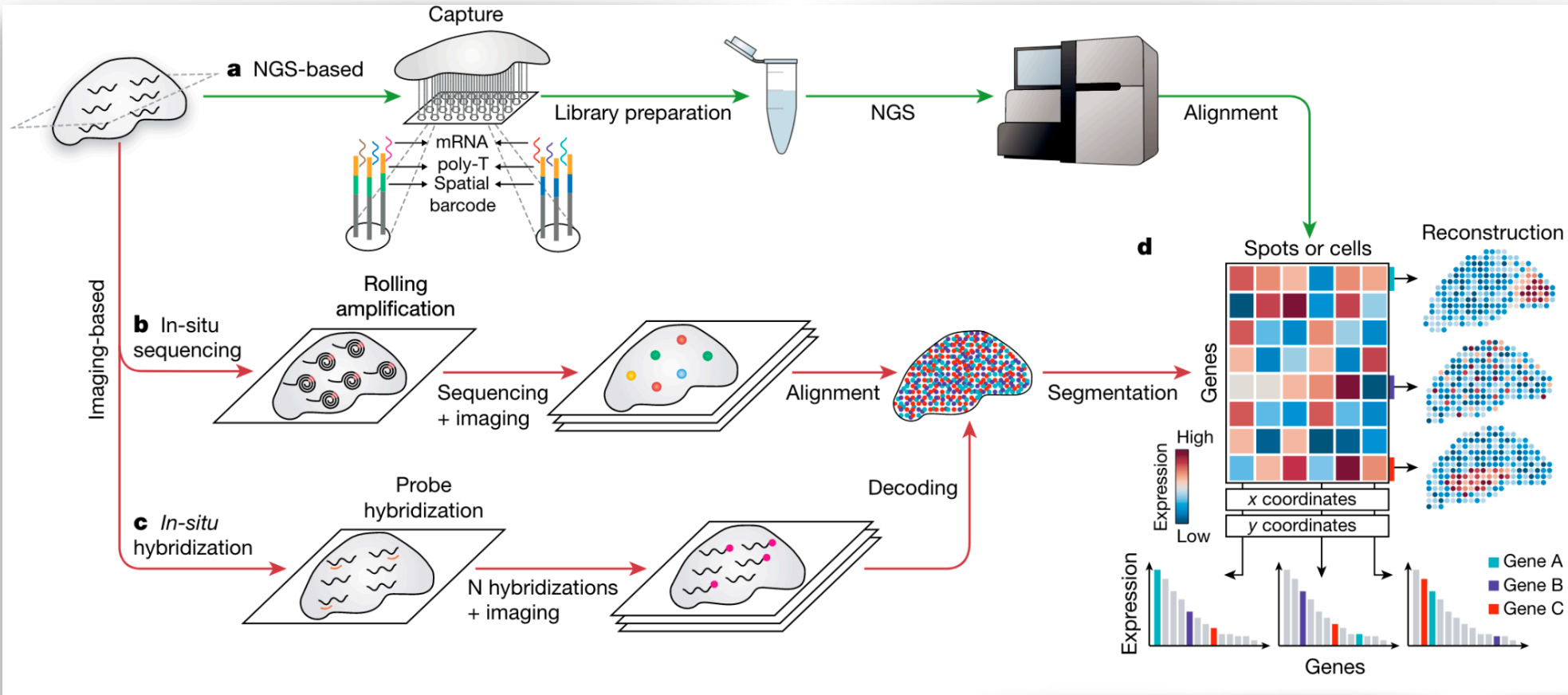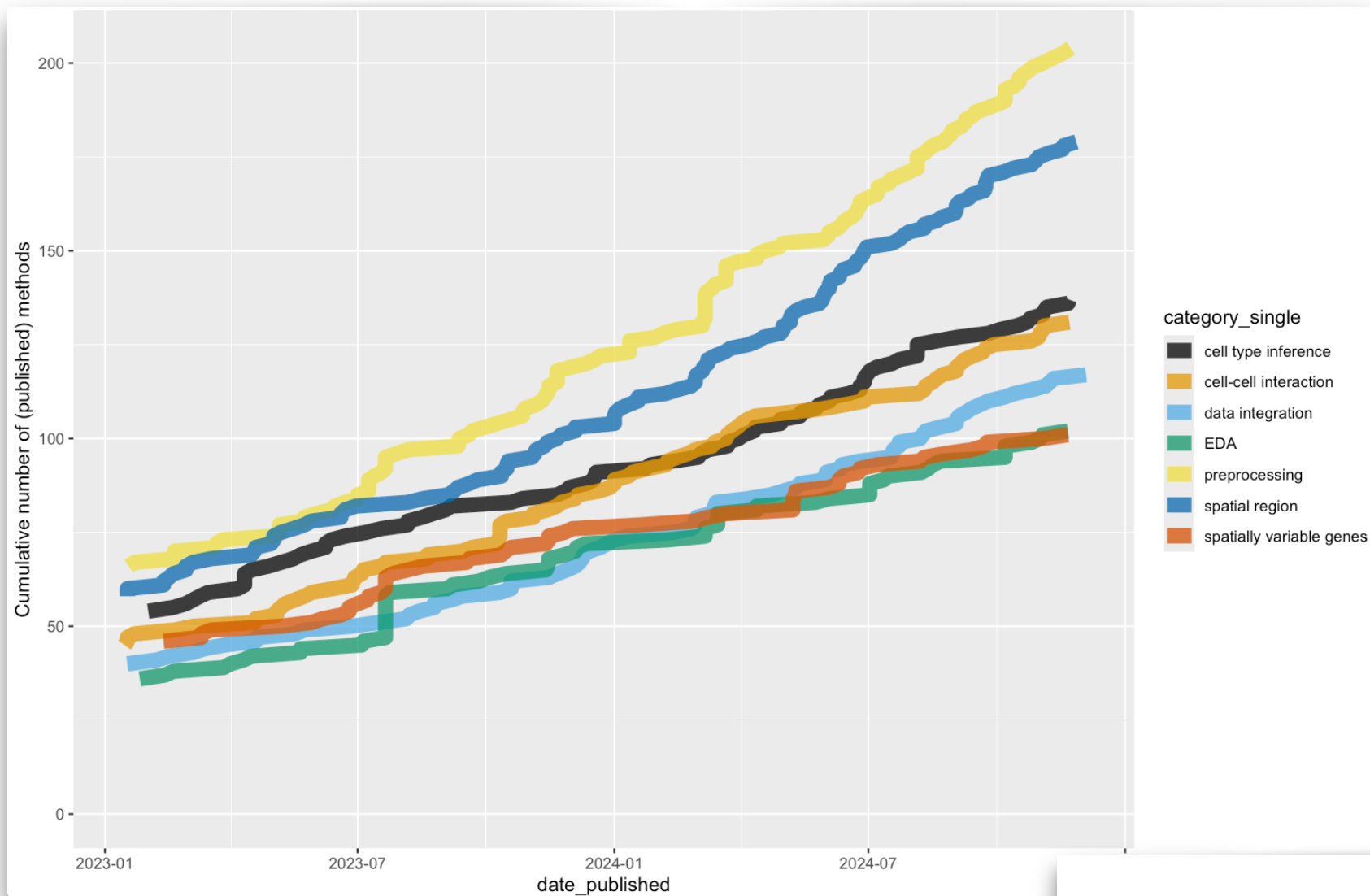# Technology choices: expression table + coordinates



**Fig. 1 | The technologies of spatial transcriptomics provide a gene-expression matrix. a**, NGS-based spatial transcriptomic methods barcode transcripts according to their location in a lattice of spots. **b**, ISS approaches directly read out the transcript sequence within the tissue. **c**, ISH methods detect ta... fluorescent probe... gene-expression m... genes and location...

(Spatial omics) computational method explosion

Museum of spatial transcriptomics

Lambda Moses [iD][1] and Lior Pachter [iD][1,2] ✉

## Finding spatially-variable genes: SpatialDE

– SpatialDE: response = normal distribution with covariance with two components: i) based on distance b/w points - exponential decay; ii) constant non-spatial variance

– Null model: fit just the non-spatial variance (i.e., without sigma)

– Fit 2 models, likelihood ratio test

**SpatialDE model.** SpatialDE models gene expression profiles $y = (y_1, \ldots, y_N)$ for a given gene across spatial coordinates $X = (x_1, \ldots, x_N)$, using a multivariate normal model of the form

$$P(y \mid \mu, \sigma_s^2, \delta, \Sigma) = N(y \mid \mu \cdot 1, \sigma_s^2 \cdot (\Sigma + \delta \cdot I)) \qquad (1)$$

The fixed effect $\mu_g \cdot 1$ accounts for the mean expression level, and $\Sigma$ denotes a spatial covariance matrix defined on the basis of the input coordinates of pairs of cells. SpatialDE uses the so-called squared exponential covariance function to define $\Sigma$:

$$\Sigma_{i,j} = k(x_i, x_j) = \exp\left( -\frac{|x_i - x_j|^2}{2 \cdot l^2} \right) \qquad (2)$$

# Spatially variable genes

Lukas M. Weber [1], Arkajyoti Saha[2], Abhirup Datta [1], Kasper D. Hansen [1] & Stephanie C. Hicks [1]

• different types (senses?)
  of spatially variable genes

# University of Zurich UZH

Statistical Bioinformatics // Department of Molecular Life Sciences

Lukas M. Weber [1], Arkajyoti Saha[2], Abhirup Datta [1], Kasper D. Hansen [1] & Stephanie C. Hicks [1]

## Spatially variable genes

$$C_{ij}(\boldsymbol{\theta}) = \sigma^2 \exp\left(\frac{-\|\mathbf{s_i} - \mathbf{s_j}\|}{l}\right)$$

LIBD

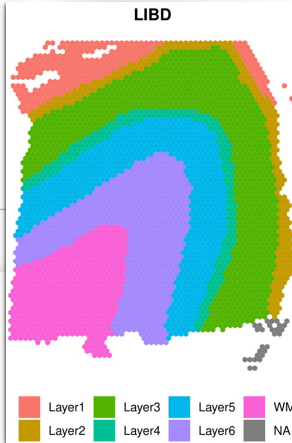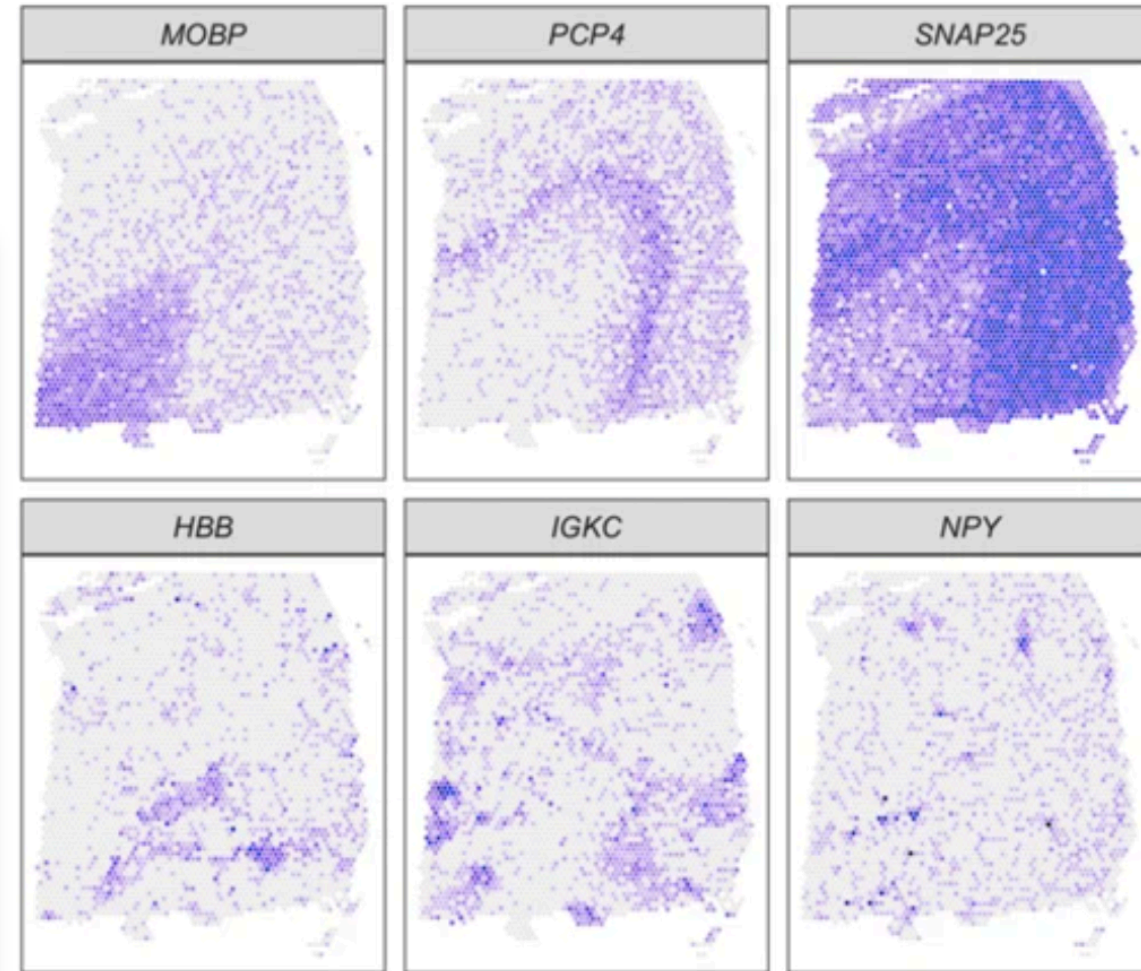Layer1  Layer3  Layer5  WM
Layer2  Layer4  Layer6  NA

Selected SVGs: human DLPFC

MOBP  PCP4  SNAP25

HBB  IGKC  NPY

counts
600

0

**b** nnSVG length scales: human DLPFC

HBB = 0.035
IGKC = 0.046
NPY = 0.015
PCP4 = 0.148
SNAP25 = 0.278
MOBP = 0.876

density — estimated length scale

# Spatially variable versus highly variable



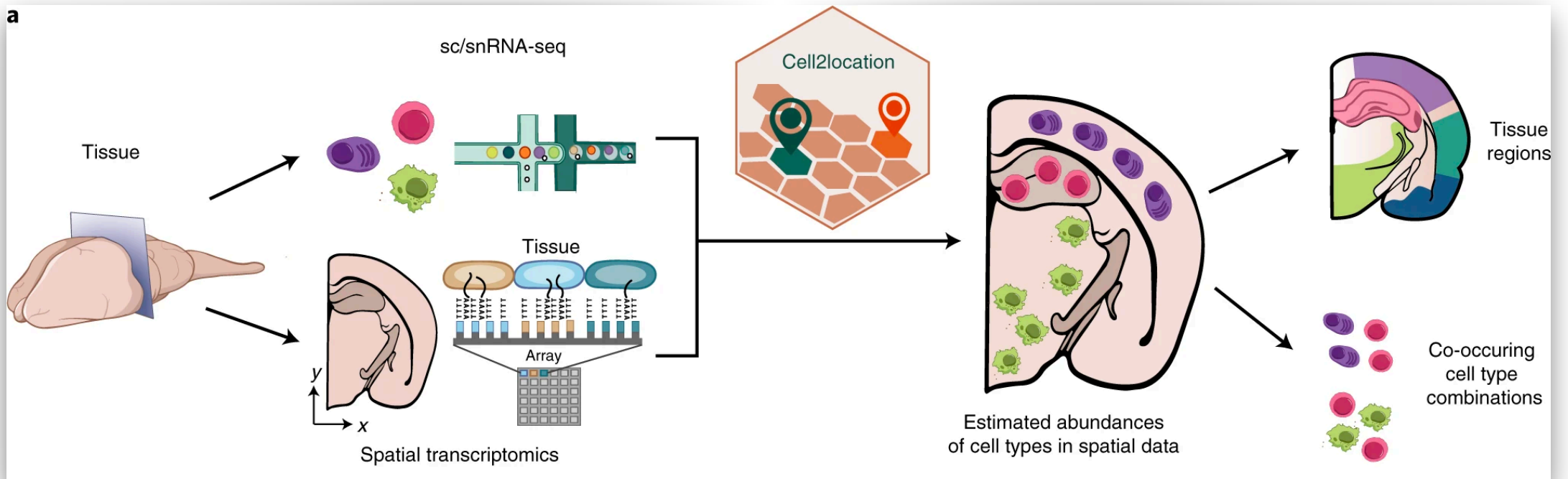(More mathematical details on Moran's I below)

8

# Statistical methods for spatial omics data

- Overview on the technologies (review)

- Finding spatially-variable genes

- Deconvoluting low-resolution (or aggregating high-resolution) spatial omics data

- Spatially-aware dimension reduction / clustering

- Cell-cell communication —> co-localization

- Classical spatial statistics

  ‣ Point patterns: random, clustered, intensity/correlation

  ‣ Lattice data: useful summaries / functions

  ‣ models with spatially correlated errors

# Deconvoluting low-resolution spatial omics (sequencing) data

– Cell2location: negative binomial regression for reference cell type signatures; decompose spot-level mRNA counts into reference cell types

# Deconvoluting low-resolution spatial omics data

– Cell2location: negative binomial regression for reference cell type signatures; decompose spot-level mRNA counts into reference cell types

*Cell2location model.* Cell2location models the elements of the spatial expression count matrix $d_{s,g}$ as negative binomial distributed, given an unobserved gene expression level (rate) $\mu_{s,g}$ and gene- and batch-specific over-dispersion $\alpha_{e,g}$:
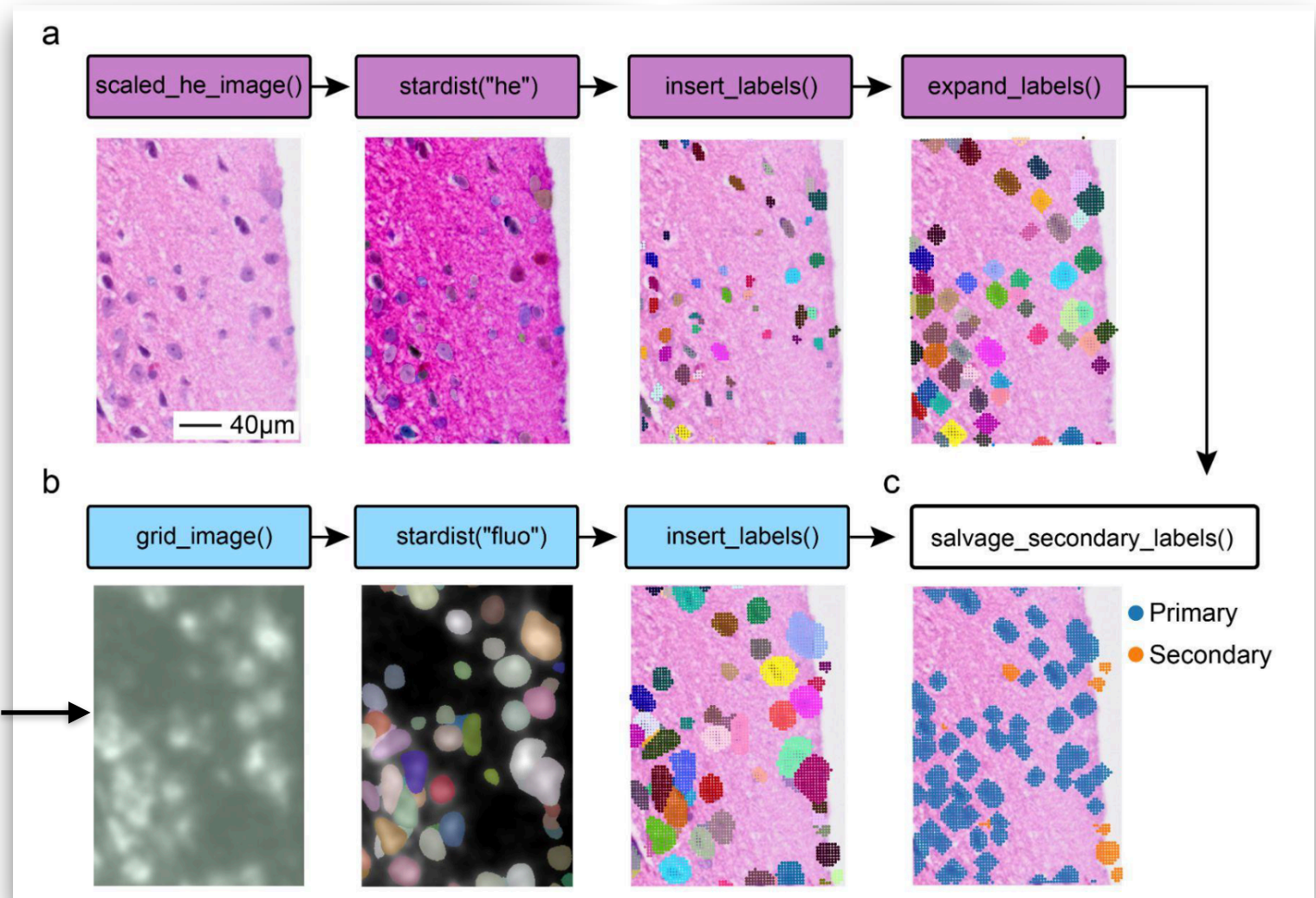
$$d_{s,g} \sim NB\left(\mu_{s,g}, \alpha_{e,g}\right).$$

The expression rate of genes $g$ at location $s$, $\mu_{s,g}$ in the mRNA count space is modeled as a linear function of reference cell types signatures $g_{f,g}$:

$$\mu_{s,g} = \left( \underbrace{m_g}_{\text{technology sensitivity}} \cdot \underbrace{\sum_f w_{s,f}\, g_{f,g}}_{\text{cell type contributions}} + \underbrace{s_{e,g}}_{\text{additive shift}} \right) \cdot \underbrace{y_s}_{\text{per-location sensitivity}} \cdot$$

# Aggregating high-resolution spatial omics (sequencing) data

– bin2cell: combines segmentation on H&E/IF and segmentation on gene expression counts

Image of counts per spot (smoothed)

# Statistical methods for spatial omics data

- Overview on the technologies (review)

- Finding spatially-variable genes

- Deconvoluting low-resolution (or aggregating high-resolution) spatial omics data

- Spatially-aware dimension reduction / clustering

- Cell-cell communication —> co-localization

- Classical spatial statistics

  ‣ Point patterns: random, clustered, intensity/correlation

  ‣ Lattice data: useful summaries / functions

  ‣ models with spatially correlated errors

Vipul Singhal [1,13], Nigel Chou [1,13], Joseph Lee [2], Yifei Yue[3], Jinyue Liu [1], Wan Kee Chock [1], Li Lin[4], Yun-Ching Chang[5], Erica Mei Ling Teo[5], Jonathan Aow [1], Hwee Kuan Lee[4,6,7,8,9,10], Kok Hao Chen [1] & Shyam Prabhakar [1,11,12]

Spatial clustering / domain detection (BANKSY)
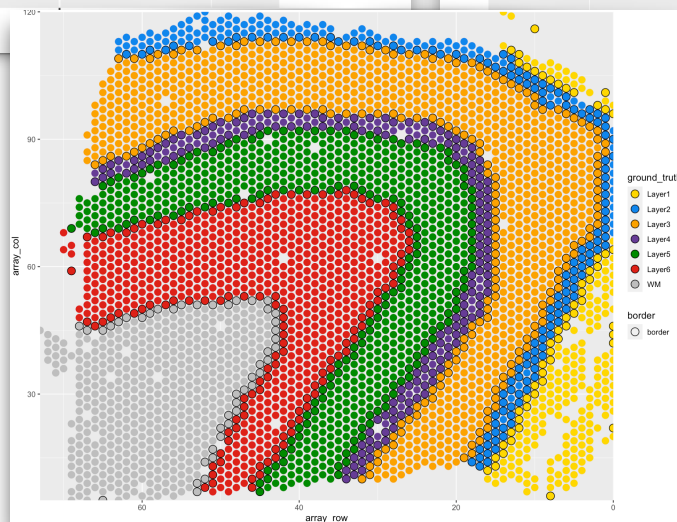—>  combine transcription and spatial information

Non-spatial PCA

Banksy PCA

Sample 151673

# Spatially aware dimension reduction for spatial transcriptomics

Lulu Shang [1,2] & Xiang Zhou [1,2] ✉

## Spatial domain detection ~ spatially homogeneous regions ~ spatial niches



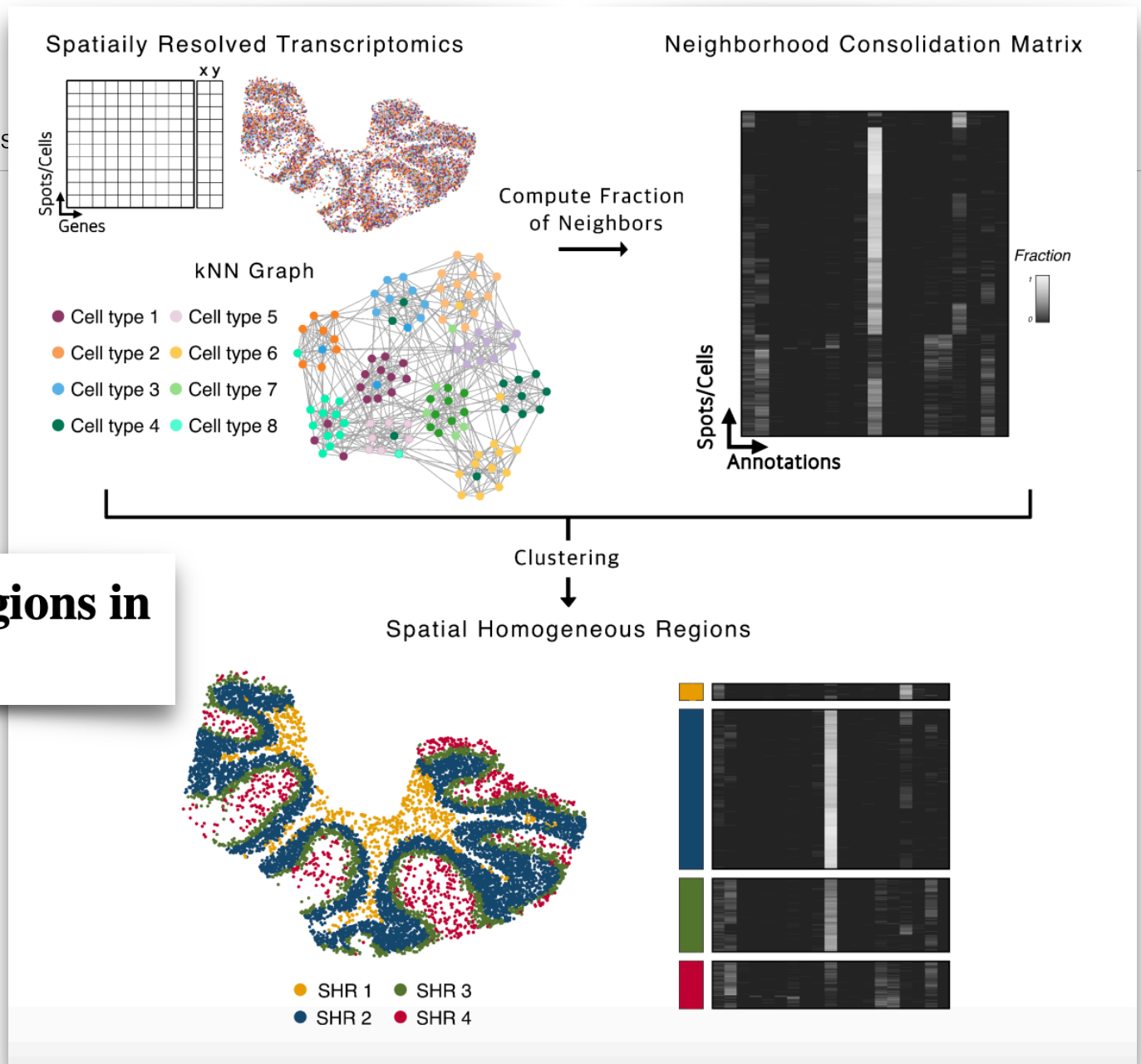|  | 🟡 | 🩷 | 🟢 | 🔵 |
|---|---|---|---|---|
| Scenario 1 | 70% | 10% | 10% | 10% |
| Scenario 2 | 45% | 45% | 5% | 5% |
| Scenario 3 | 60% | 30% | 5% | 5% |
| Scenario 4 | 35% | 30% | 30% | 5% |

https://www.nature.com/articles/s41467-022-34879-1

# Spatial domain detection ~ spatially homogeneous regions

**Identification of spatial homogeneous regions in tissues with concordex**



https://www.biorxiv.org/content/10.1101/2023.06.28.546949v2

# Alternatively, spatially variable features = DE between domains



Simone Tiberi

Peiying Cai

Figure 2. Three examples of simulated SVGs, from the *LIBD* data, following *bottom/right*, *circular*, and *annotations* patterns. Examples of SVGs from *mixture* and *inverted mixture* patterns are presented in Supplementary Fig. S1.

To find spatially variable genes (SVGs); spatial clustering + classical statistical method works quite well

BayesSpace_DESpace    SPARK-X    MERINGUE    BayesSpace_findMarkers    StLearn_FindAllMarkers
StLearn_DESpace    SpatialDE    SpaGCN    StLearn_findMarkers
SPARK    SpatialDE2    nnSVG    BayesSpace_FindAllMarkers

JOURNAL ARTICLE

*DESpace*: spatially variable gene detection via differential expression testing of spatial clusters

Peiying Cai, Mark D Robinson, Simone Tiberi ✉

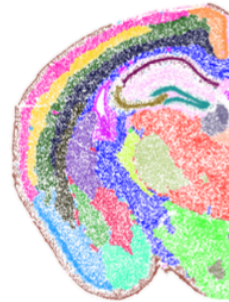# Statistical methods for spatial omics data

- Overview on the technologies (review)

- Finding spatially-variable genes

- Deconvoluting low-resolution (or aggregating high-resolution) spatial omics data

- Spatially-aware dimension reduction / clustering

- Cell-cell communication —> co-localization

- Classical spatial statistics

  ‣ Point patterns: random, clustered, intensity/correlation

  ‣ Lattice data: useful summaries / functions

  ‣ models with spatially correlated errors

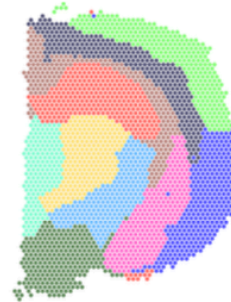# `pasta`: Data representations determine spatial statistics options



Samuel

Martin

**Harnessing the potential of spatial statistics for spatial omics data with *pasta***

**Martin Emons** [1,†]**, Samuel Gunz** [1,†]**, Helena L. Crowell** [2]**, Izaskun Mallona** [1]**, Malte Kuehl** [3,4]**,**
**Reinhard Furrer** [5]**, Mark D. Robinson** [1,*]

[1]Department of Molecular Life Sciences and SIB Swiss Institute of Bioinformatics, University of Zurich, 8057 Zurich, Switzerland
[2]Centro Nacional de Análisis Genómico (CNAG), 08028 Barcelona, Spain
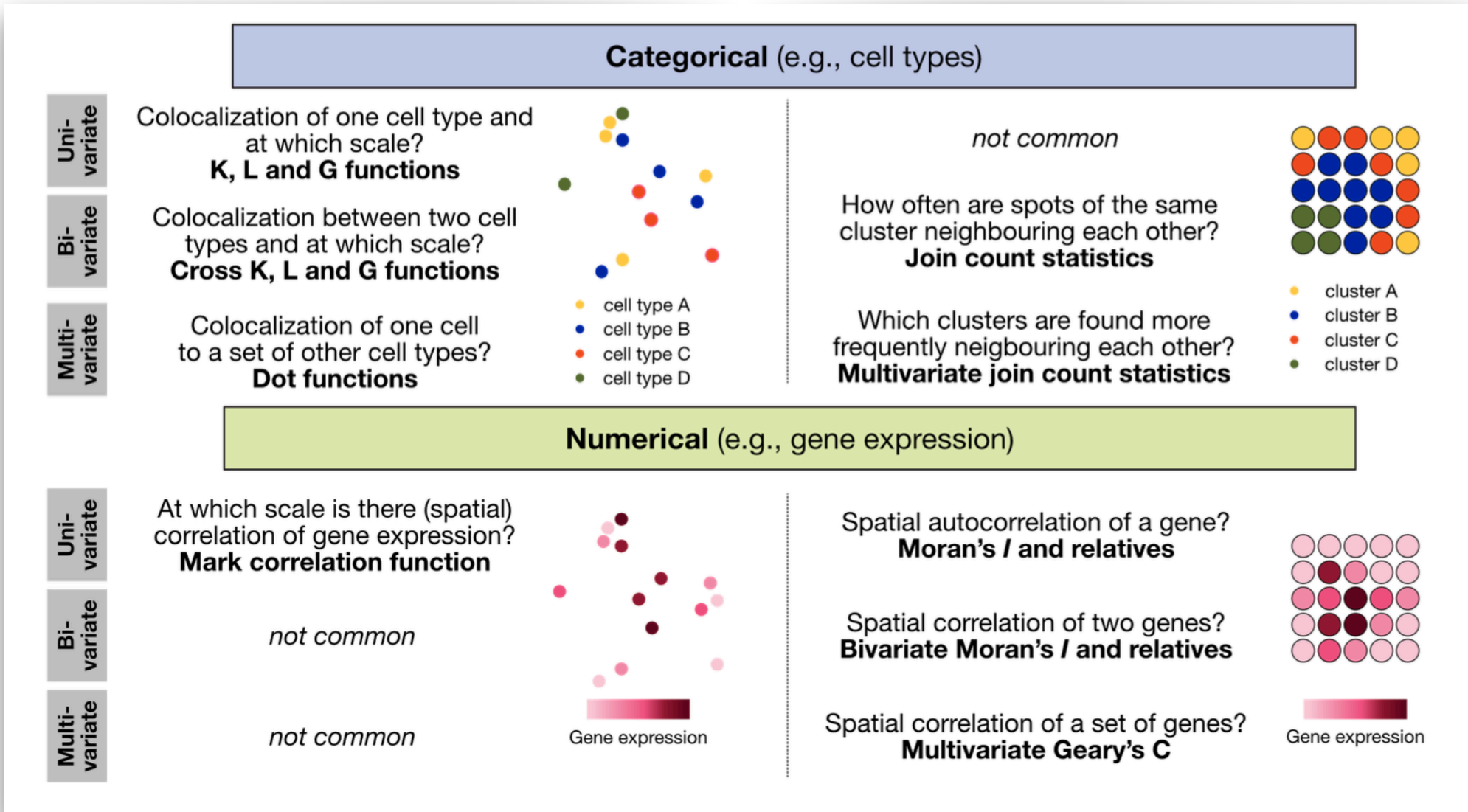[3]Department of Clinical Medicine, Aarhus University, 8200 Aarhus N, Denmark
[4]Department of Pathology, Aarhus University Hospital, 8200 Aarhus N, Denmark
[5]Department of Mathematical Modeling and Machine Learning, University of Zurich, 8057 Zurich, Switzerland

[*]To whom correspondence should be addressed. Email: mark.robinson@mls.uzh.ch
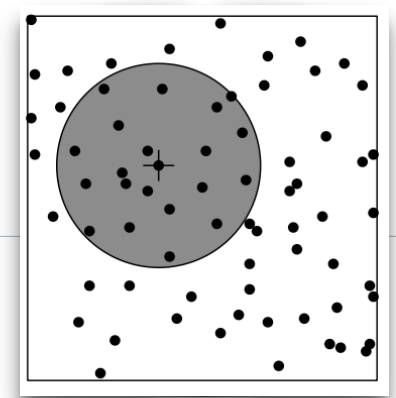[†]The first two authors should be regarded as Joint First Authors.

# `pasta`: Data representations determine spatial statistics options



Samuel

Martin

# Correlation for **point patterns**

– Ripley's K function
– mathematical definition:

$$K(r) = \frac{1}{\lambda} \mathbb{E}\left[\text{number of } r\text{-neighbours of u} \mid \mathbf{X} \text{ has a point at location } u\right]$$

$$t(u, r, \mathbf{x}) = \sum_{j=1}^{n(\mathbf{x})} \mathbf{1}\left\{0 < \|u - x_j\| \le r\right\}$$

**Definition 7.1.** *If $\mathbf{X}$ is a stationary point process, with intensity $\lambda > 0$, then for any $r \ge 0$*
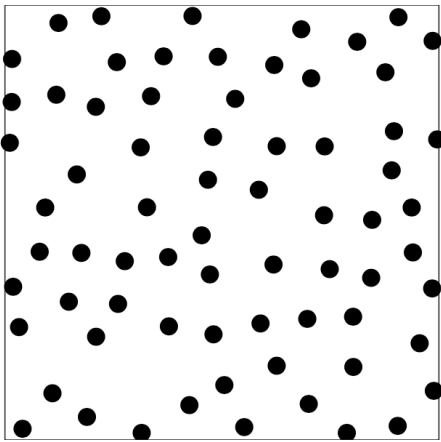
$$K(r) = \frac{1}{\lambda}\mathbb{E}\left[t(u, r, \mathbf{X}) \mid u \in \mathbf{X}\right] \tag{7.6}$$

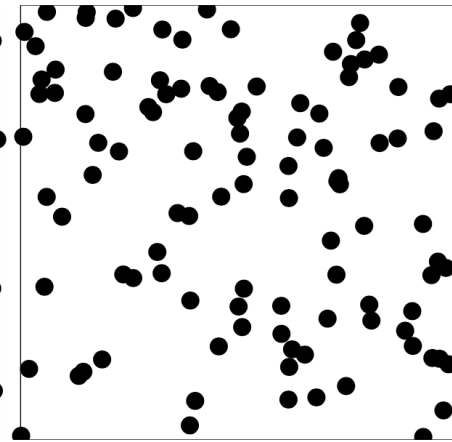*does not depend on the location u, and is called the K-function of $\mathbf{X}$.*

# Correlation for **point patterns**

– Ripley's K function

– words definition: *the empirical K-function K(r) is the cumulative average number of data points lying within a distance r of a typical data point*

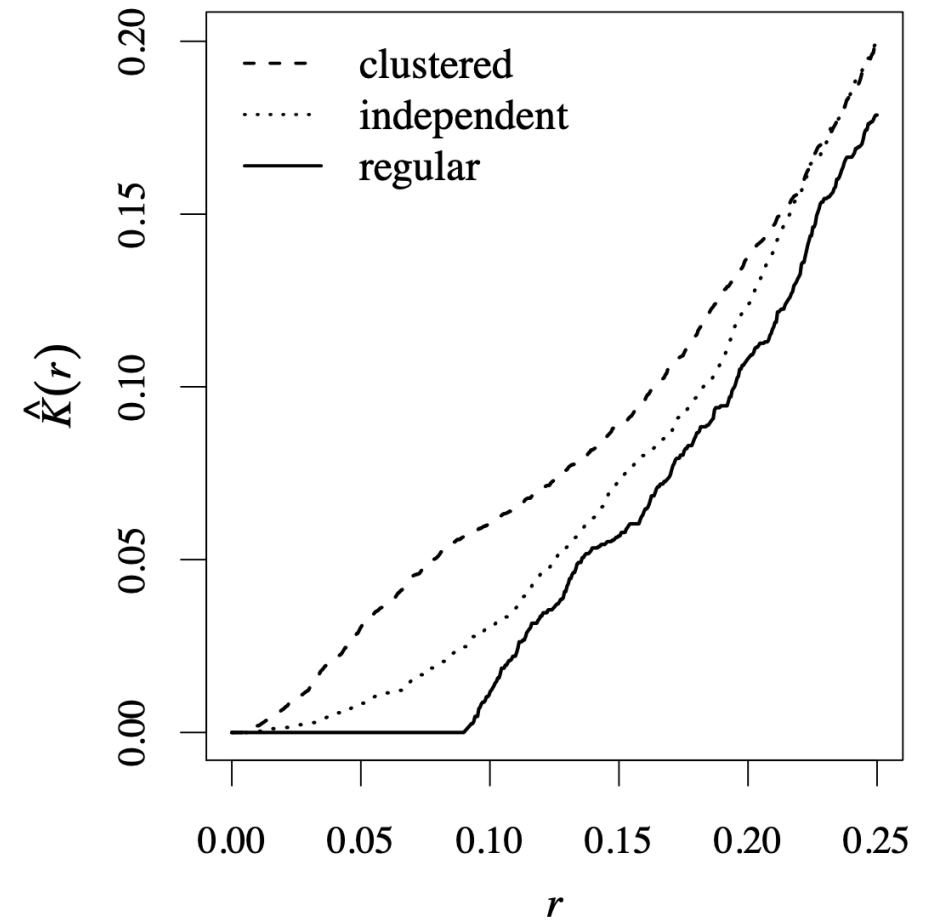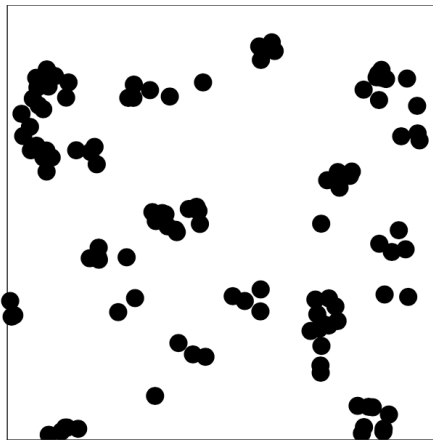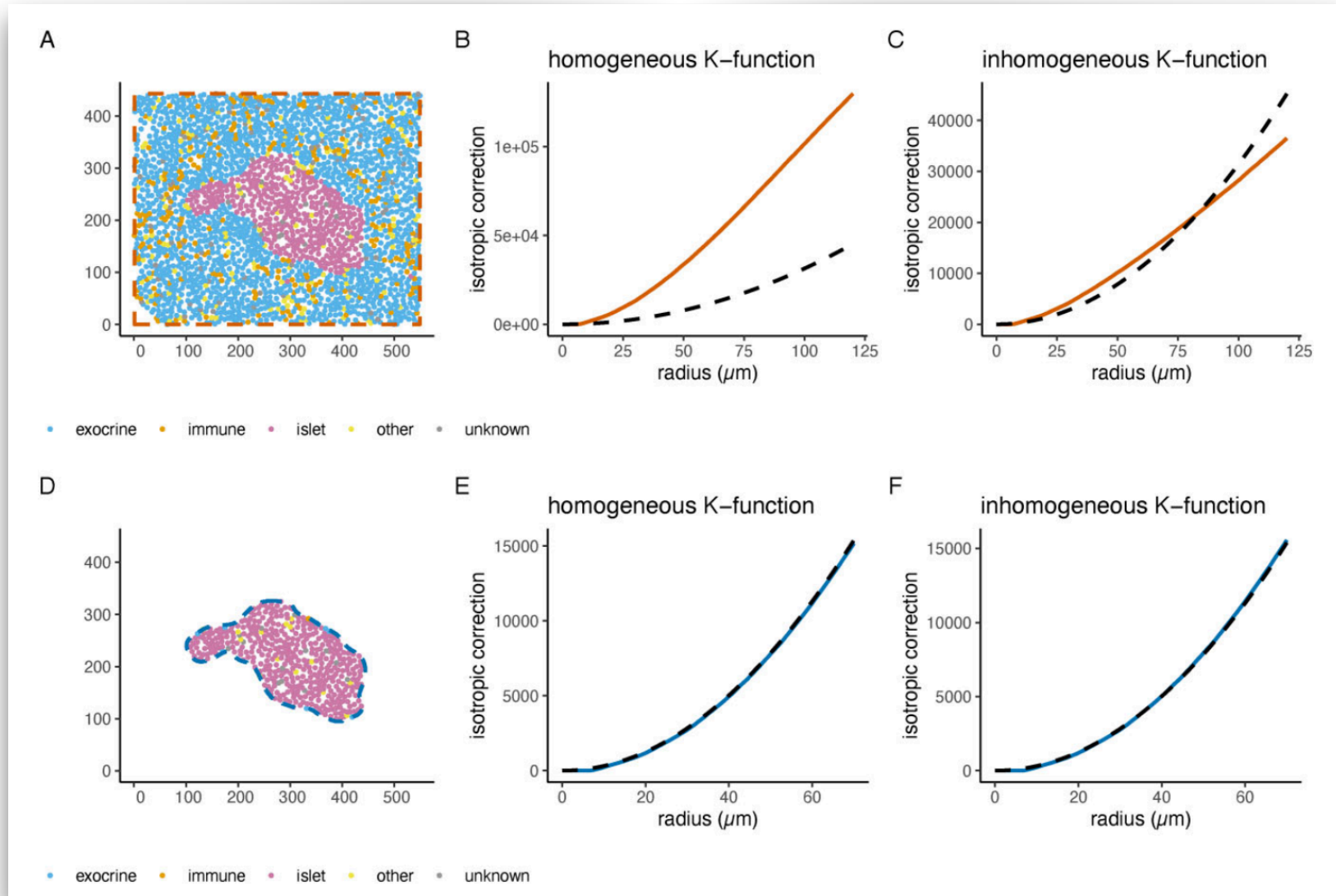regular            independent          clustered

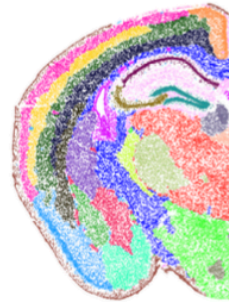# `pasta`: the 'gotcha' of spatial statistics — is it clustering or intensity?



Samuel

Martin

K-functions here: clustering / intensity of pink cells (islets).

# `pasta`: Data representations determine spatial statistics options



**A** TECHNOLOGY
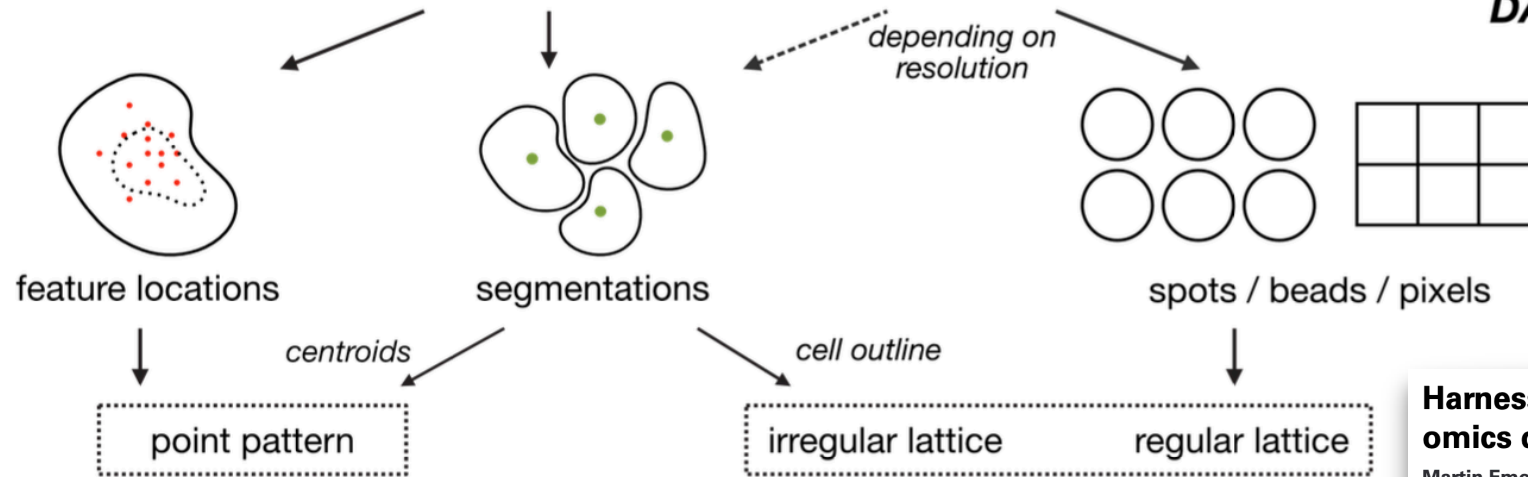
Imaging-based
- Targeted
- Higher resolution

STARmap

HTS-based
- Untargeted
- Lower resolution

10X Visium

**B** DATA MODALITY

depending on resolution

feature locations

segmentations

spots / beads / pixels

centroids

cell outline

point pattern

irregular lattice    regular lattice

Samuel

Martin

**Harnessing the potential of spatial statistics for spatial omics data with *pasta***

Martin Emons [1,†], Samuel Gunz [1,†], Helena L. Crowell [2], Izaskun Mallona [1], Malte Kuehl [3,4], Reinhard Furrer [5], Mark D. Robinson [1,*]

[1]Department of Molecular Life Sciences and SIB Swiss Institute of Bioinformatics, University of Zurich, 8057 Zurich, Switzerland
[2]Centro Nacional de Análisis Genómico (CNAG), 08028 Barcelona, Spain
[3]Department of Clinical Medicine, Aarhus University, 8200 Aarhus N, Denmark
[4]Department of Pathology, Aarhus University Hospital, 8200 Aarhus N, Denmark
[5]Department of Mathematical Modeling and Machine Learning, University of Zurich, 8057 Zurich, Switzerland

[*]To whom correspondence should be addressed. Email: mark.robinson@mls.uzh.ch
[†]The first two authors should be regarded as Joint First Authors.

Statistical Bioinformatics // Department of Molecular Life Sciences

# Spatial autocorrelation: Global Moran's I

– Global measure of auto-correlation (correlation to signal nearby in space); assume homogeneity!
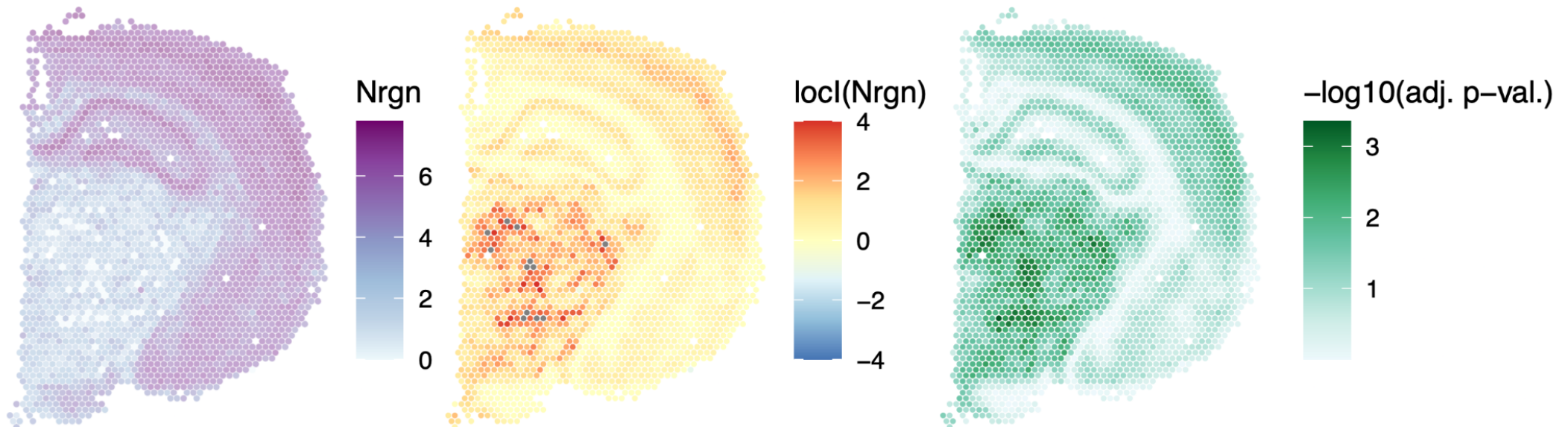
– Alternative: Geary's C

$$I = \frac{1}{\sum_{ij} w_{ij}} \frac{\sum_{ij} w_{ij}(X_i - \overline{X})(X_j - \overline{X})}{N^{-1} \sum_i (X_i - \overline{X})^2}$$

$$C = \frac{(N-1)\sum_i \sum_j w_{ij}(x_i - x_j)^2}{2W \sum_i (x_i - \bar{x})^2}$$

# Spatial autocorrelation: Local Moran's I

– Local measure of auto-correlation (correlation to signal nearby in space)

$$I_i = \frac{x_i - \bar{x}}{\sum_{k=1}^{n}(x_k - \bar{x})^2/(n-1)} \sum_{j=1}^{n} w_{ij}(x_j - \bar{x})$$

$$\text{Global Moran's } R = \frac{\sum_i \sum_j w_{ij}(x_i - \bar{x})(y_j - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2}\sqrt{\sum_i (y_i - \bar{y})^2}},$$

## Cell-cell communication

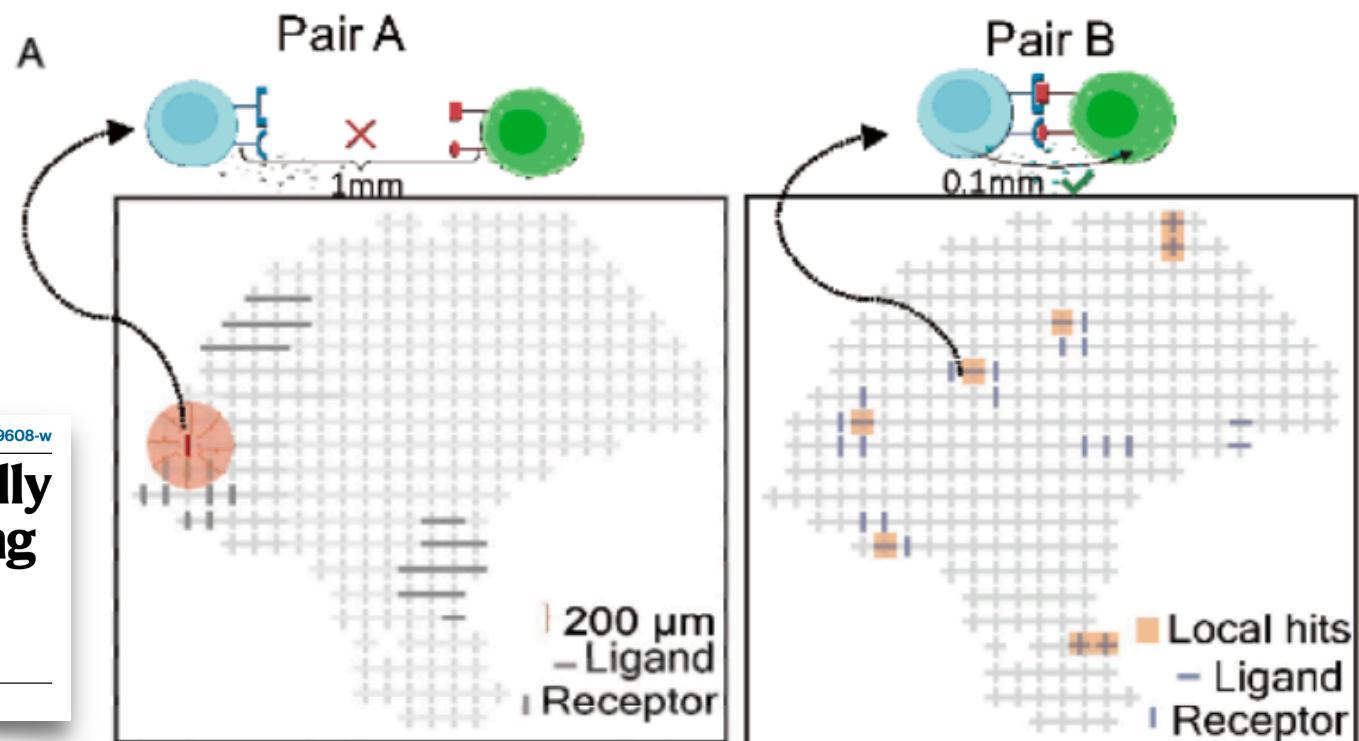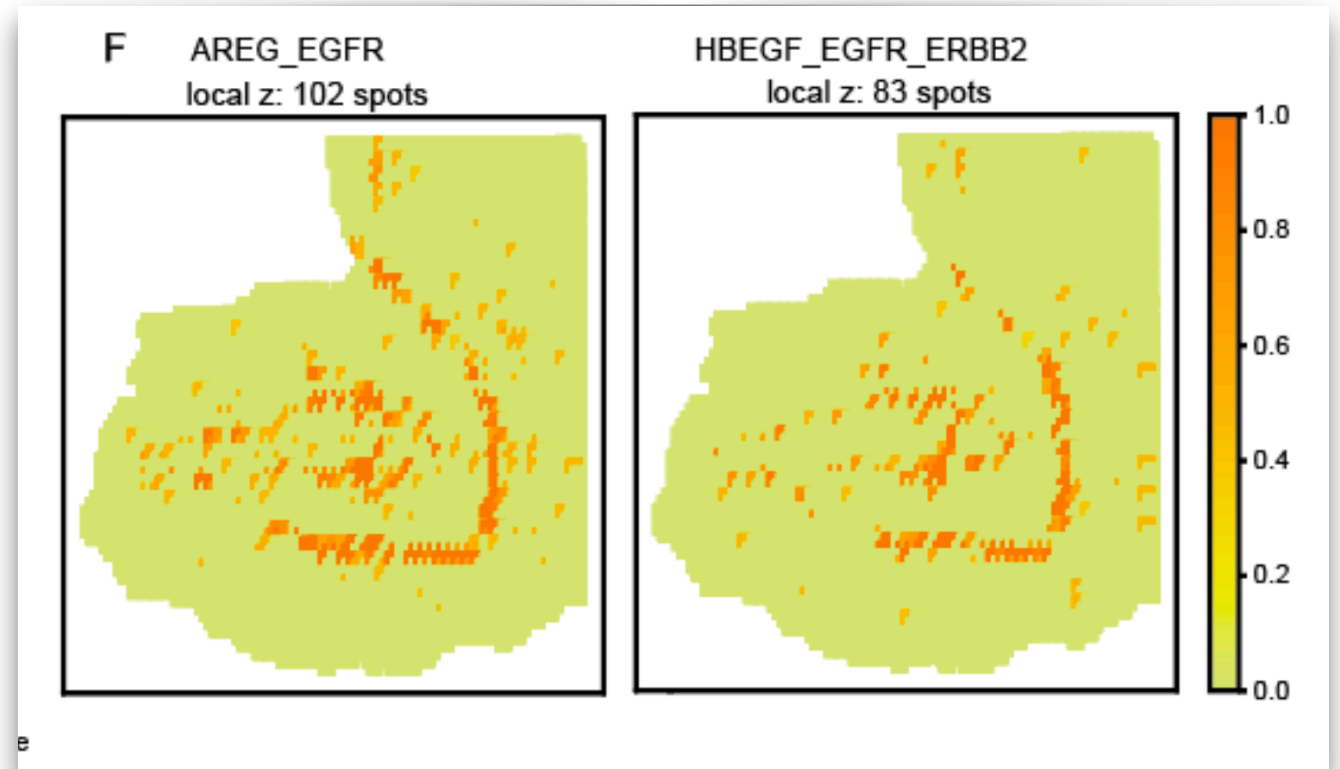– SpatialDM: Global Moran's R, which is a bivariate version of Moran's I

$$\text{Global Moran's } R = \frac{\sum_i \sum_j w_{ij}(x_i - \bar{x})(y_j - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2}\sqrt{\sum_i (y_i - \bar{y})^2}},$$

## Cell-cell communication

– SpatialDM: Global Moran's R, which is a bivariate version of Moran's I



SpatialDM

Non-significant due to spatial range

$$\text{Global Moran's R} = \frac{N}{W}\sum_i \sum_j w_{ij}\, \tilde{x}_i\, \tilde{y}_j$$
$$\text{Local Moran's R} = \sum_j w_{ij}\, \tilde{x}_i\, \tilde{y}_j + \sum_j w_{ij}\, \tilde{y}_i\, \tilde{x}_j$$

F  AREG_EGFR
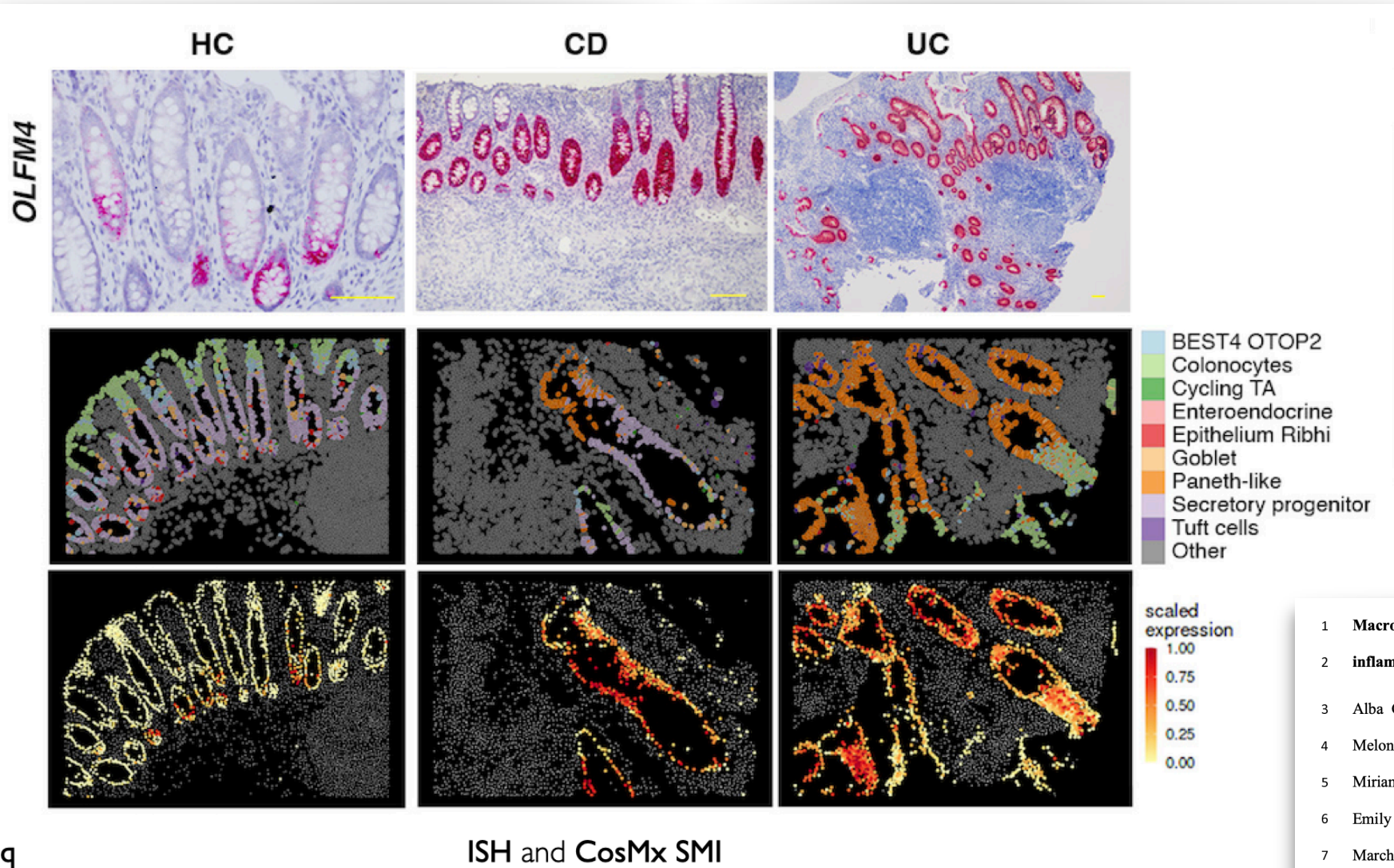local z: 102 spots

HBEGF_EGFR_ERBB2
local z: 83 spots

# Research

- `spatialFDA`: Flexible modeling of point pattern summaries —> Martin

- `DESpace2`: DE beyond markers/SVGs: "differential spatial patterns" —> Peiying

- **`sosta`: "Spatial structure"-focused analyses**

- OSTA: Orchestrating spatial transcriptomics analysis with Bioconductor

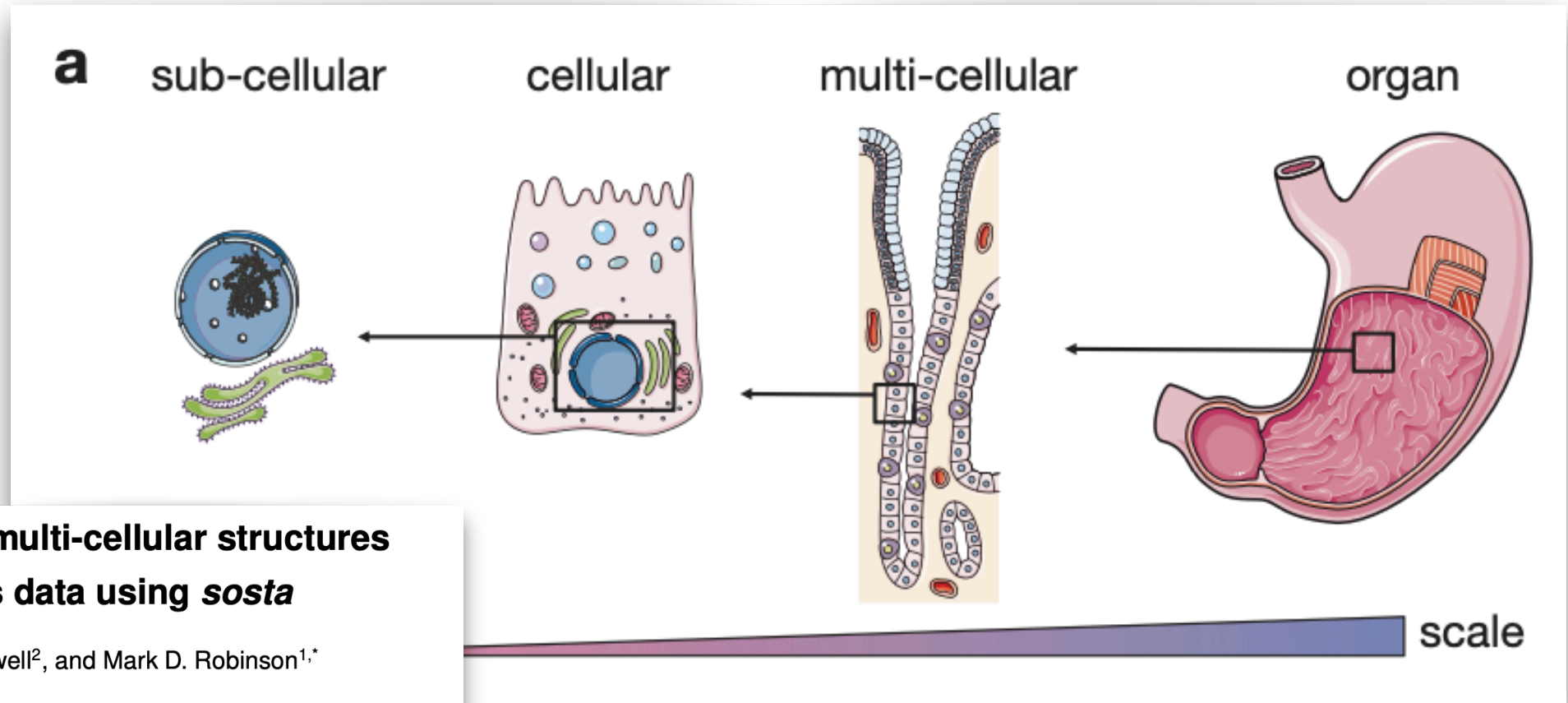- SpaceHack: using consensus clustering to consolidate domain detection

# Tissue "structures" are often visible



- healthy control (HC)
- Crohn's disease (CD)
- ulcerative colitis (UC)

1  **Macrophage and neutrophil heterogeneity at single-cell spatial resolution in**
2  **inflammatory bowel disease**

3  Alba Garrido-Trigo[1,2], Ana M. Corraliza[1,2], Marisol Veny[1,2], Isabella Dotti[1,2], Elisa
4  Melon-Ardanaz[1,2], Aina Rill[3], Helena L. Crowell[4], Ángel Corbí[5], Victoria Gudiño[1,2],
5  Miriam Esteller[1,2], Iris Álvarez-Teubel[1,2], Daniel Aguilar[1,2], M Carme Masamunt[1,2],
6  Emily Killingbeck[6], Youngmi Kim[6], Michael Leon[6], Sudha Visvanathan[7], Domenica
7  Marchese[8], Ginevra Caratù[8], Albert Martin-Cardona[2,9], Maria Esteve[2,9], Julian Panés[1,2],
8  Elena Ricart[1,2], Elisabetta Mereu[3,*], Holger Heyn[8,10,*], Azucena Salas[1,2]

# Tissue "structures" occur at different scales



**Analysis of anatomical multi-cellular structures from spatial omics data using *sosta***

Samuel Gunz[1], Helena L. Crowell[2], and Mark D. Robinson[1,*]

[1]Department of Molecular Life Sciences and
SIB Swiss Institute of Bioinformatics, University of Zurich, Zurich, Switzerland
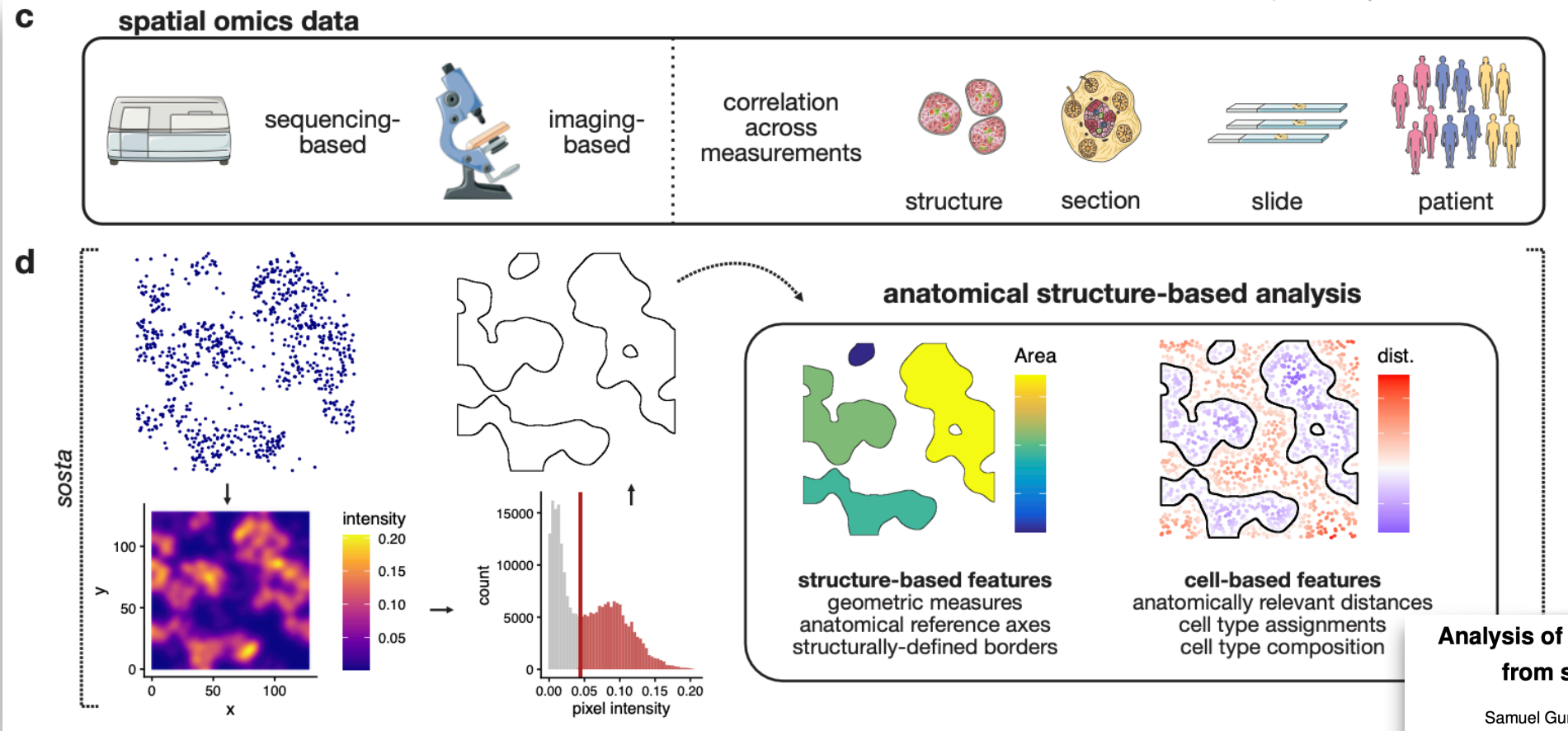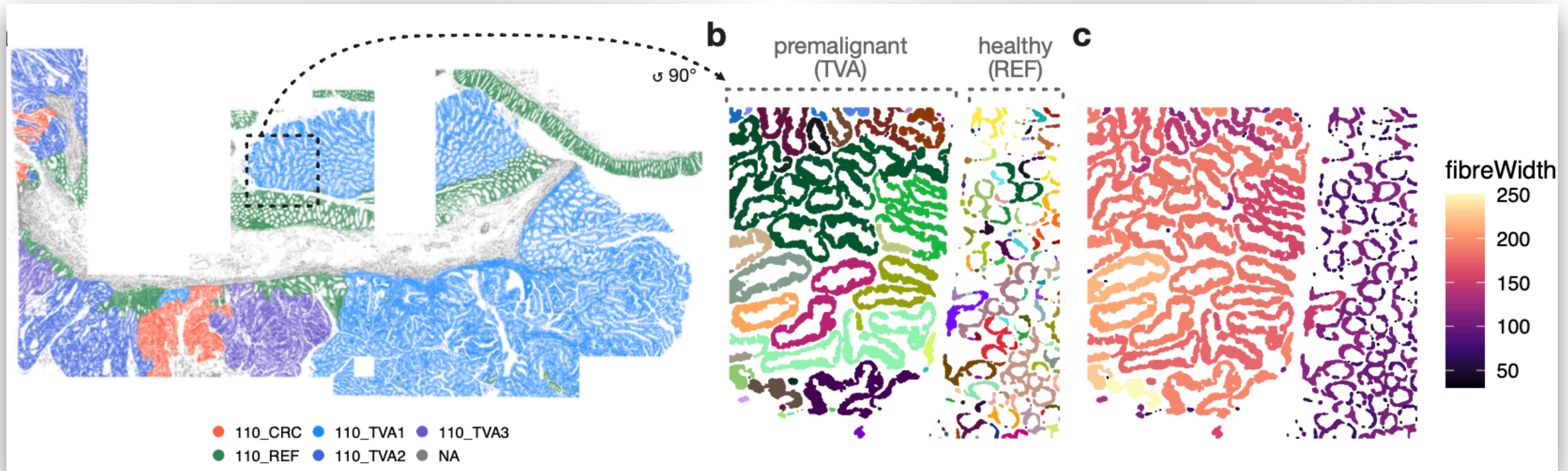[2]Centro Nacional de Análisis Genómico, Barcelona, Spain
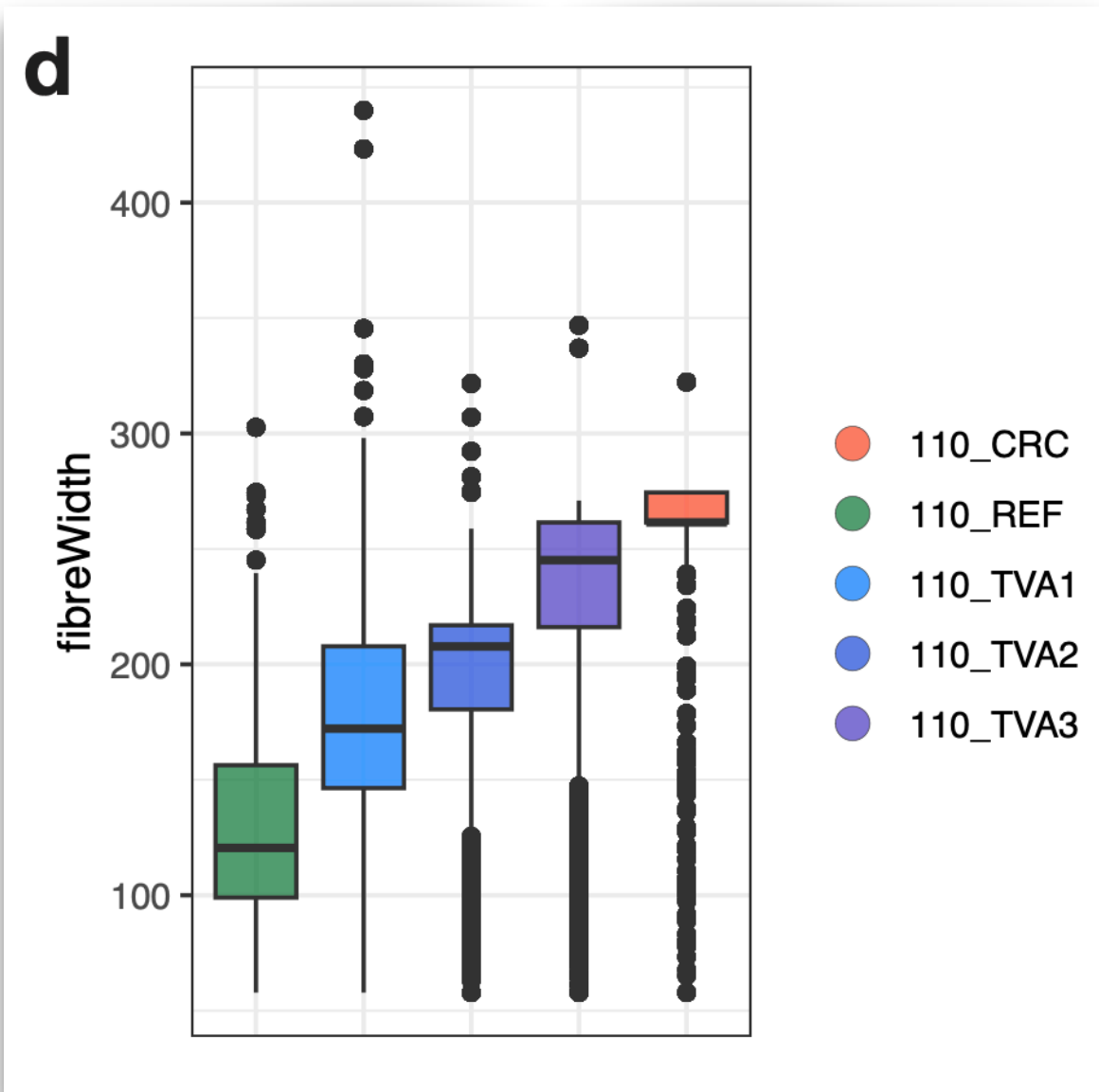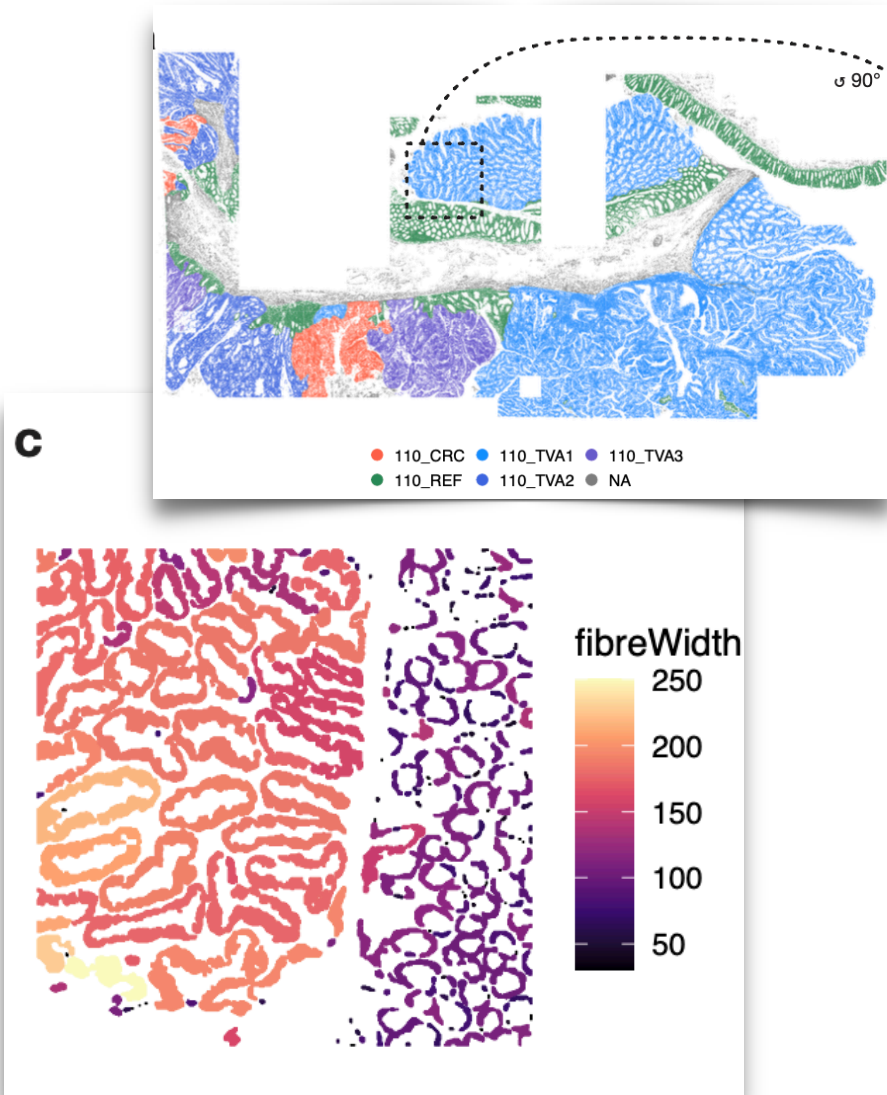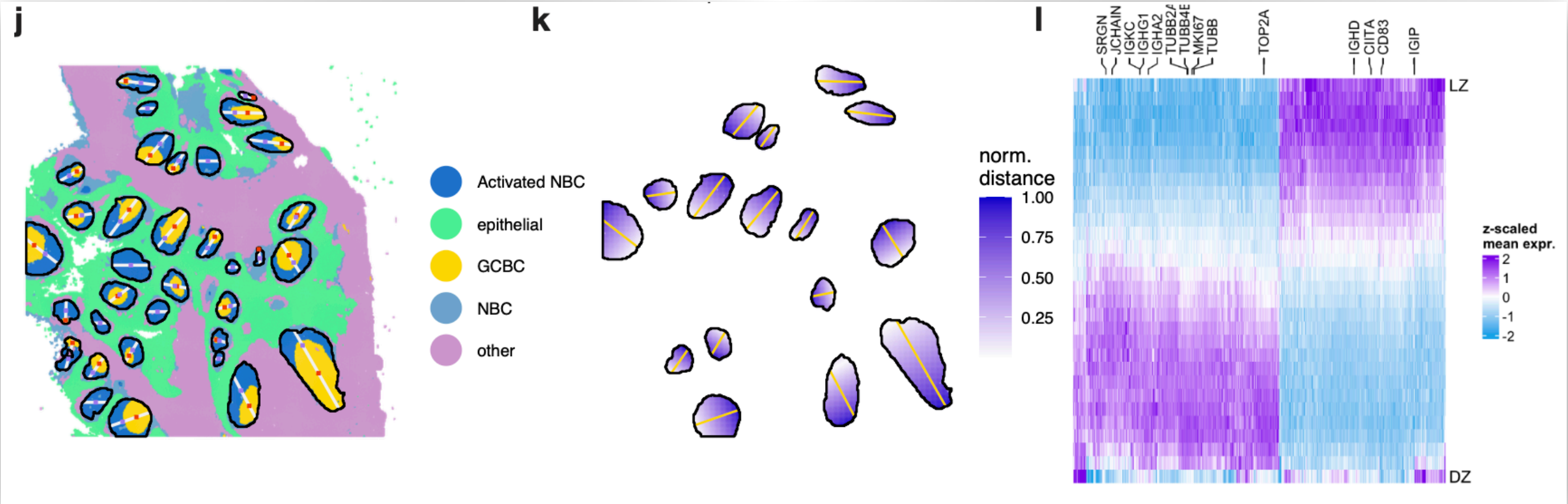[*]Correspondence to: mark.robinson@mls.uzh.ch

October 29, 2025

# sosta: extracting spatial "structures" + quantifying metrics + modelling (differential discovery)

Samuel

**Analysis of anatomical multi-cellular structures from spatial omics data using *sosta***

Samuel Gunz[1], Helena L. Crowell[2], and Mark D. Robinson[1,*]

[1]Department of Molecular Life Sciences and
SIB Swiss Institute of Bioinformatics, University of Zurich, Zurich, Switzerland
[2]Centro Nacional de Análisis Genómico, Barcelona, Spain
*Correspondence to: mark.robinson@mls.uzh.ch

October 29, 2025

Preprint available.

# Variation among spatial structures (epithelial example)

# Variation among spatial structures (geometric quantifications)

# Structures —> Reference axis —> Expression gradients

# Modeling requires accounting for repeated measurements



Potentially i) multiple structures per tissue slice; ii) multiple slices per patient; iii) replication ——> **multiple levels of variability** ——> mixed models generally most appropriate

# Orchestrating Spatial Transcriptomics Analysis with Bioconductor

- https://bioconductor.org/books/OSTA

**Orchestrating Spatial Transcriptomics Analysis with Bioconductor**

Helena L. Crowell[1,*,✉], Yixing Dong[2,3,*], Ilaria Billato[4], Peiying Cai[5,6], Martin Emons[5,6], Samuel Gunz[5,6], Boyi Guo[7], Mengbo Li[8,9,10], Alexandru Mahmoud[11], Artür Manukyan[12], Hervé Pagès[13], Pratibha Panwar[14,15,16], Shreya Rao[14,15,17], Callum J. Sargeant[8], Lori Shepherd Kern[18], Marcel Ramos[19,20], Jieran Sun[2,3], Michael Totty[21], Vincent J. Carey[11], Yunshun Chen[8,9,10], Leonardo Collado-Torres[21,22,23], Shila Ghazanfar[14,15,16], Kasper D. Hansen[21,24,25], Keri Martinowich[22,26,27,28], Kristen R. Maynard[22,26,27], Ellis Patrick[14,15,16,17], Dario Righelli[29], Davide Risso[30,31], Simone Tiberi[32], Levi Waldron[19,20], Raphael Gottardo[2,3,33,†,✉], Mark D. Robinson[5,6,†,✉], Stephanie C. Hicks[21,25,34,35,†,✉], and Lukas M. Weber[36,†,✉]

Book is available. Preprint on bioRxiv.
(Successor of the OSCA book: https://bioconductor.org/books/OSCA/)

# Meta-benchmark

Reported method performances are inconsistent across studies

Jieran Sun[1†], Kirti Biharie[2,3†], Peiying Cai[4†], Niklas Müller-Bötticher[5†], Paul Kiessling[6†], Meghan A. Turner[7†], Søren H. Dam[8,9†], Florian Heyl[10,11†], Sarusan Kathirchelvan[4], Martin Emons[4], Samuel Gunz[4], Sven Twardziok[5], Amin El-Heliebi[12], Martin Zacharias[13], SpaceHack 2.0 participants, Roland Eils[5], Marcel Reinders[3], Raphael Gottardo[1], Christoph Kuppe[6], Brian Long[7*], Ahmed Mahfouz[2,3*], Mark D. Robinson[4*], Naveed Ishaque[5*]

Peiying Cai



A — Reported ARI for BayesSpace

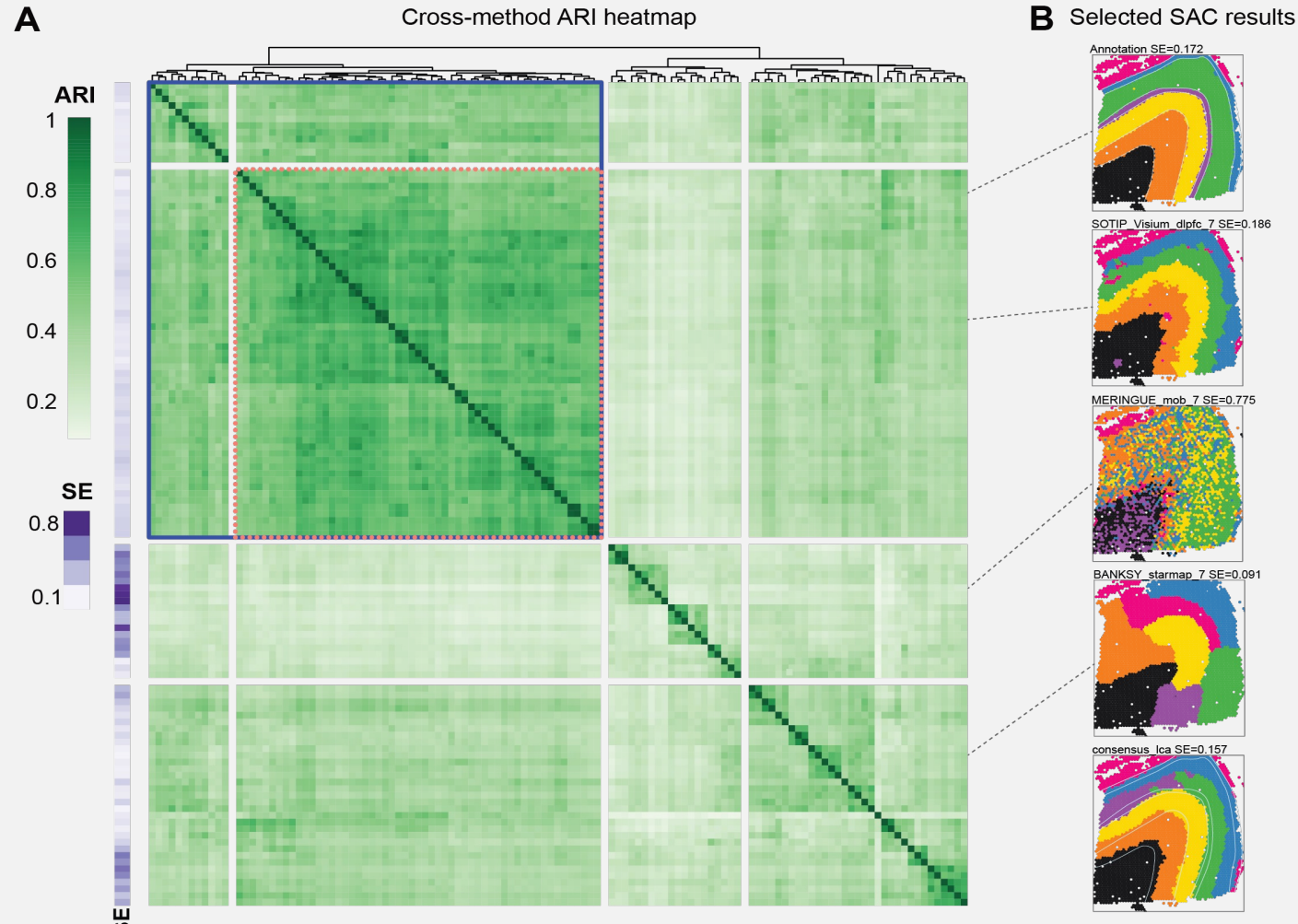B — Self reported ARI vs. ARI from other studies

# Ensemble clustering

Methods are often more similar to each other than to the ground truth.

Beyond benchmarking: an expert-guided consensus approach to spatially aware clustering
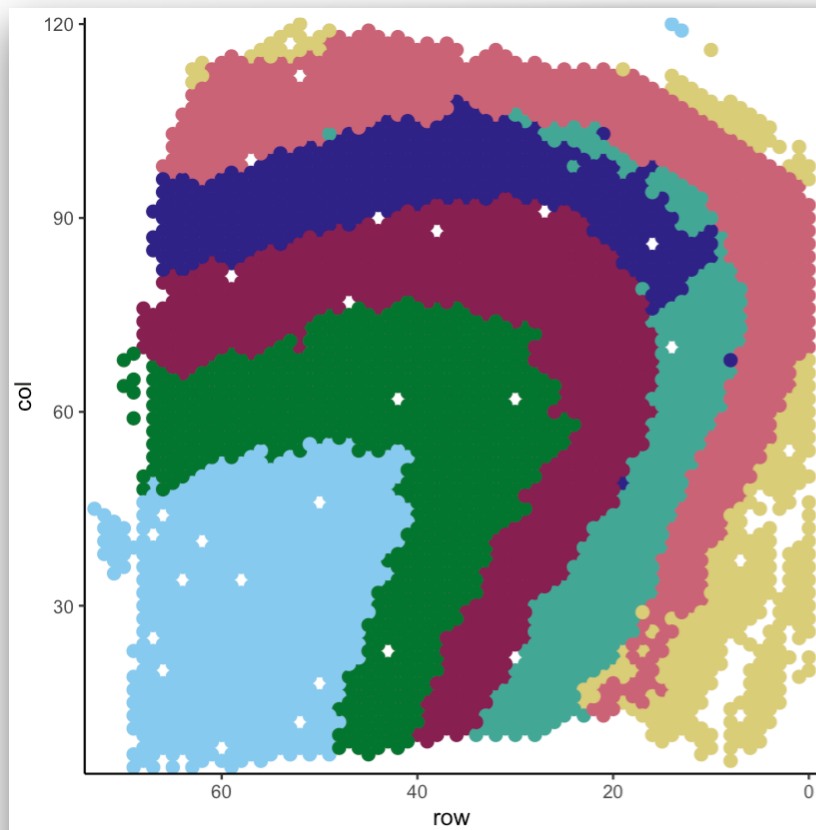
Jieran Sun[1†], Kirti Biharie[2,3†], Peiying Cai[4†], Niklas Müller-Bötticher[5†], Paul Kiessling[6†], Meghan A. Turner[7†], Søren H. Dam[8,9†], Florian Heyl[10,11†], Sarusan Kathirchelvan[4], Martin Emons[4], Samuel Gunz[4], Sven Twardziok[5], Amin El-Heliebi[12], Martin Zacharias[13], SpaceHack 2.0 participants, Roland Eils[5], Marcel Reinders[3], Raphael Gottardo[1], Christoph Kuppe[6], Brian Long[7*], Ahmed Mahfouz[2,3*], Mark D. Robinson[4*], Naveed Ishaque[5*]

**A** Cross-method ARI heatmap

**B** Selected SAC results



Smoothness Entropy (Low = smooth)

Annotation SE=0.172

SOTIP_Visium_dlpfc_7 SE=0.186

MERINGUE_mob_7 SE=0.775

BANKSY_starmap_7 SE=0.091

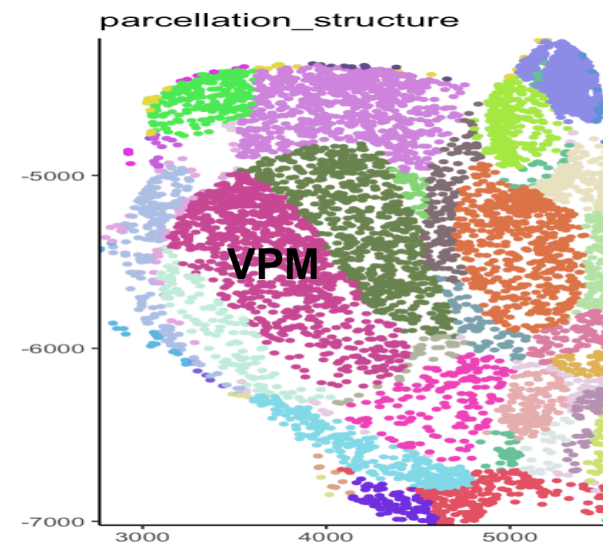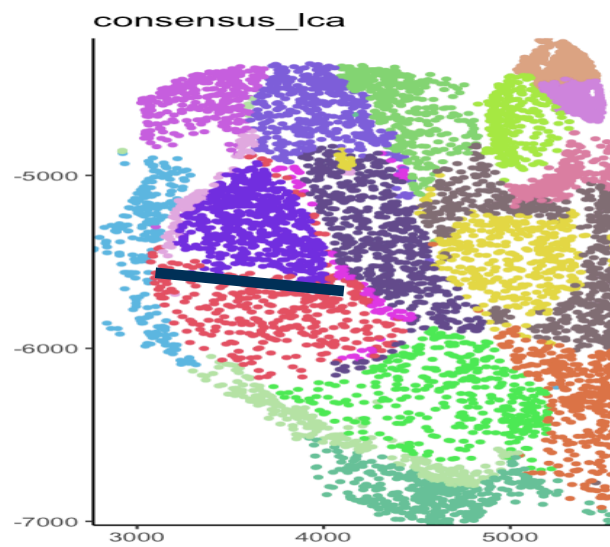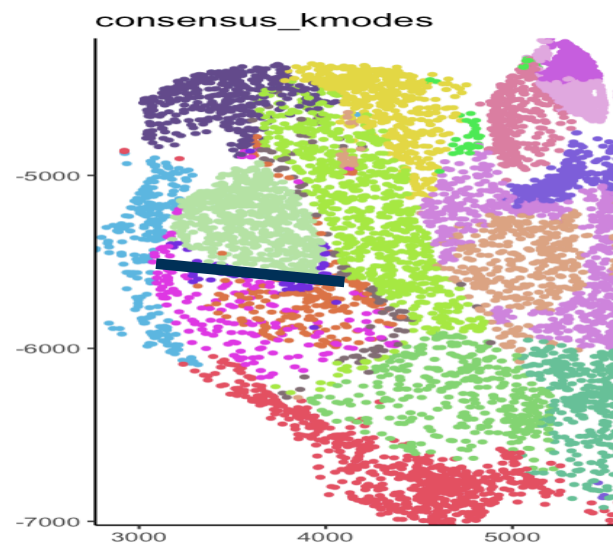consensus_lca SE=0.157

# Consensuses

# Entropy #2: Understanding spot-level uncertainty (across methods)

Entropy in the sense of <u>how stable across algorithms</u>

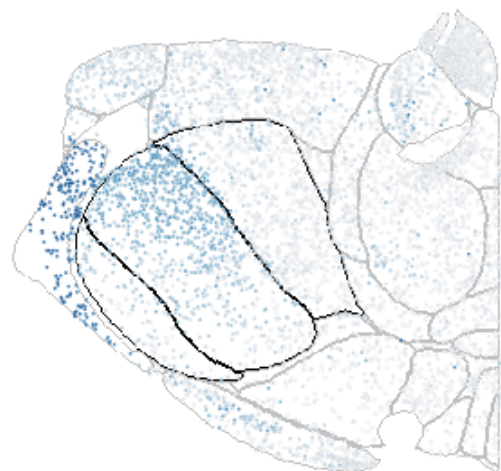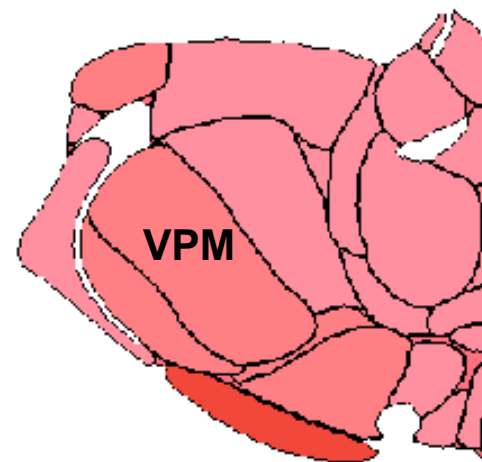(align the spot-wise cluster labels across methods, entropy across label proportions)

# VPM

# Concluding remarks

- You are collecting/analyzing spatial data: what **spatial features** do you want to quantify?

- A few places where (classical) spatial statistics might be useful; data determines: point patterns versus lattice

- Functional data analysis, multi-cellular structure-based analyses, caveats of benchmarking