



Projects

Reminders:

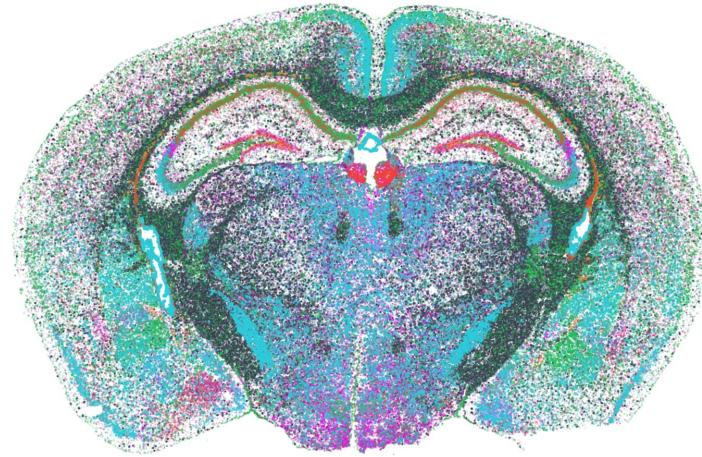
- *project “plan” by **1.12.2025** (a few bullet points)* — done!
- *Private Slack channels (group + Mark, Hubert, optionally collaborators)* — done!
- *Private GitHub repos (of same name)* — done!
- due **9.01.2026** 18.00
- “Office hours” (zoom otherwise Slack): 20.-29.12 Mark is away

bulk single-cell



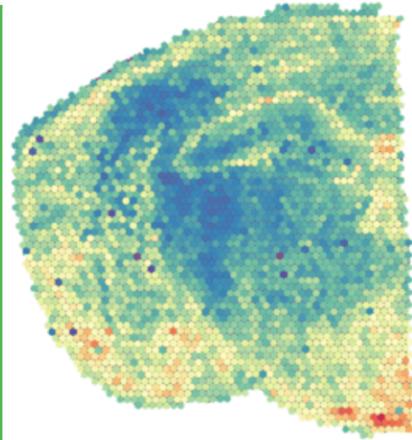
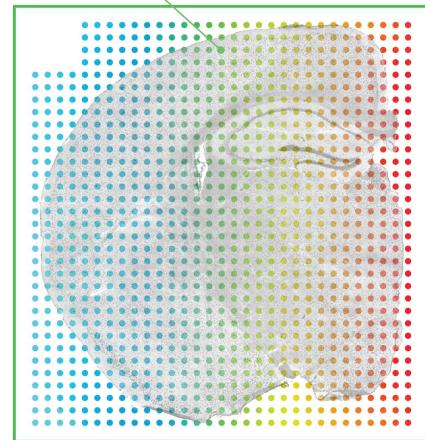
spatial

imaging-based



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

sequencing-based



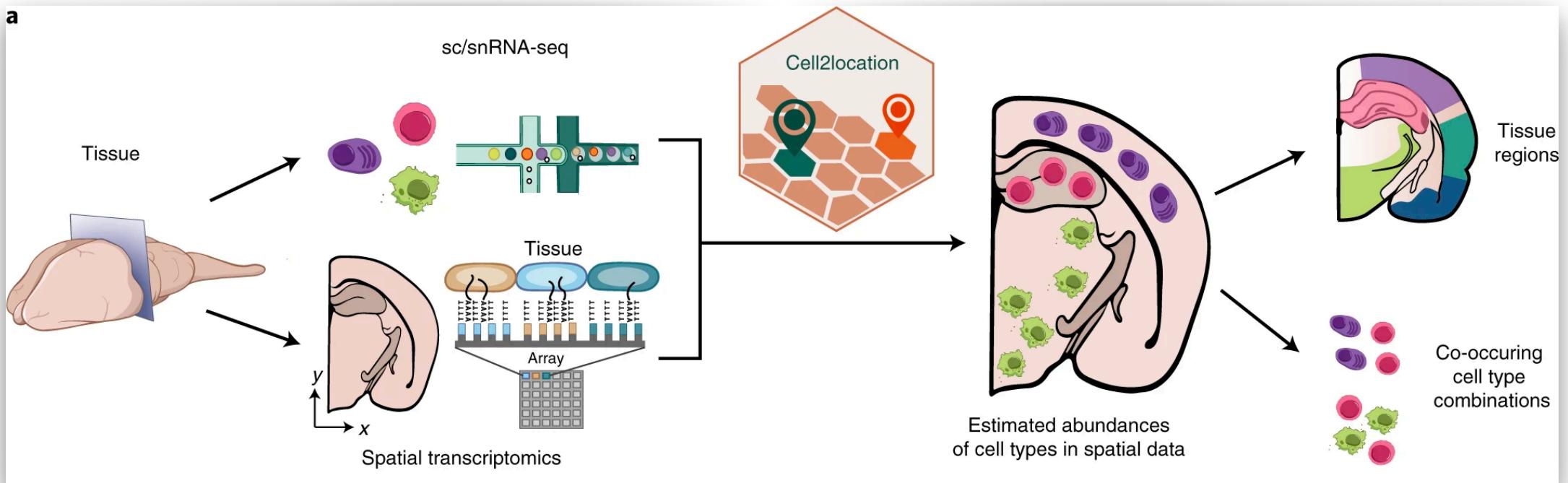
- spot-level data
- whole transcriptome (10,000s of features)
- single-cell resolutions requires aggregation or deconvolution



Slide from
Helena Crowell

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 $g_{f,g}$:

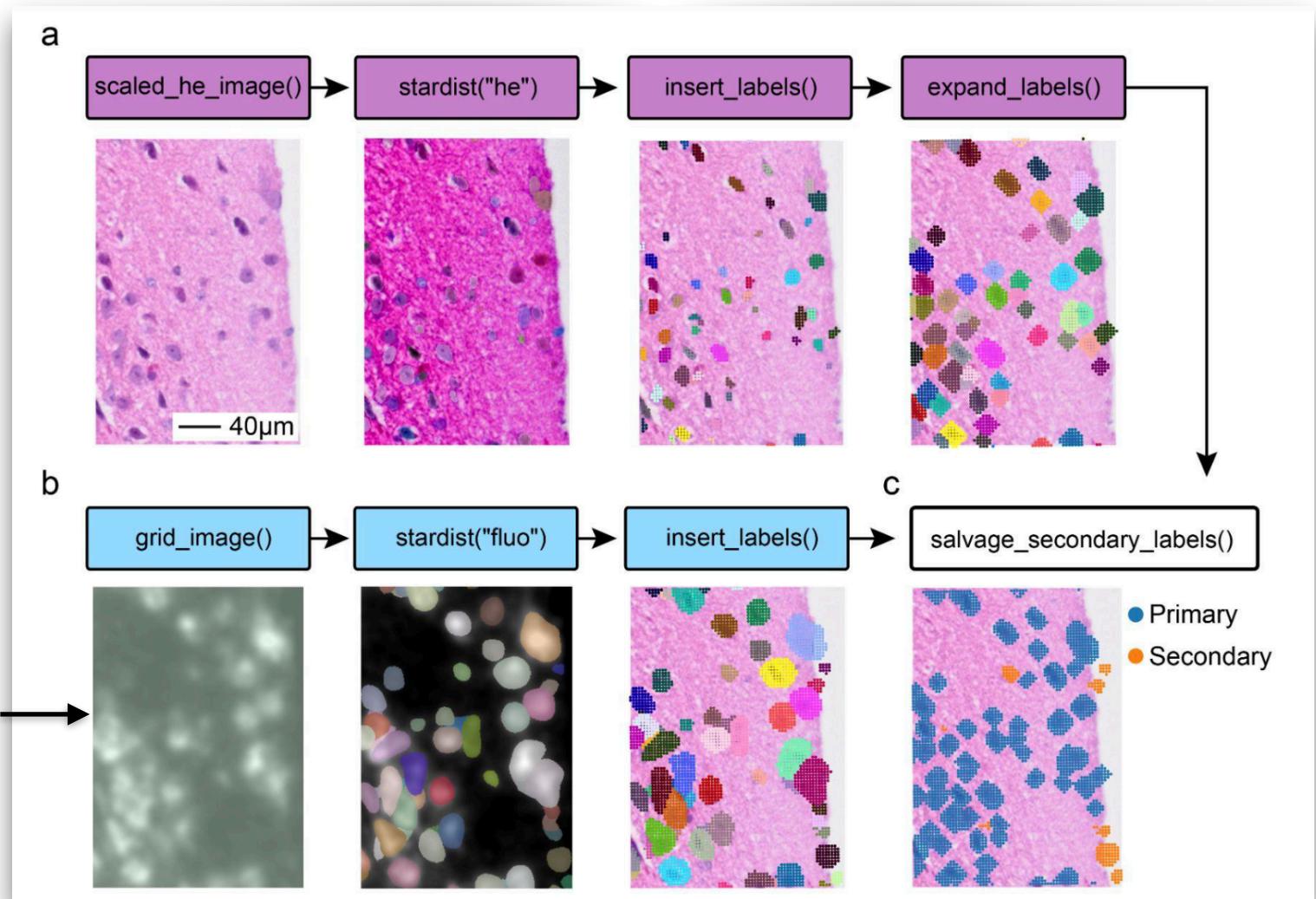
$$\mu_{s,g} = \underbrace{m_g}_{\text{technology sensitivity}} \cdot \underbrace{\sum_f w_{s,f} g_{f,g}}_{\text{cell type contributions}} + \underbrace{s_{e,g}}_{\text{additive shift}} \cdot \underbrace{y_s}_{\text{per-location sensitivity}}.$$



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)

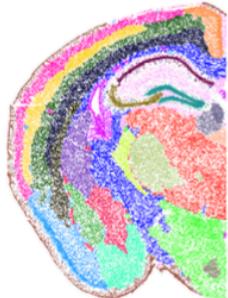


pasta: Data representations determine spatial statistics options

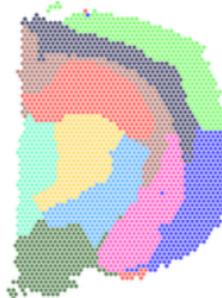
A

Imaging-based

- Targeted
- Higher resolution



STARmap



10X Visium

TECHNOLOGY

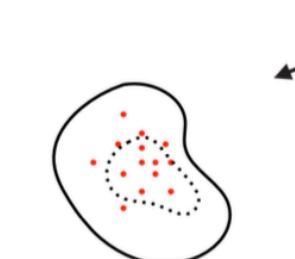
HTS-based

- Untargeted
- Lower resolution

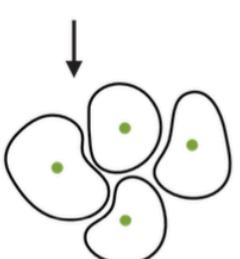


Samuel

B



feature locations



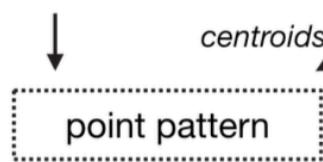
segmentations

depending on
resolution

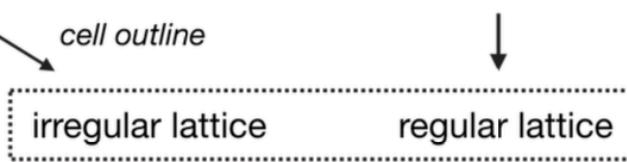


spots / beads / pixels

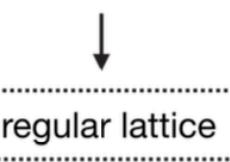
DATA MODALITY



centroids



cell outline



regular lattice

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

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

^①Department of Molecular Life Sciences and SIB Swiss Institute of Bioinformatics, University of Zurich, 8057 Zurich, Switzerland

^②Centro Nacional de Análisis Genómico (CNAG), 08028 Barcelona, Spain

^③Department of Clinical Medicine, Aarhus University, 8200 Aarhus N, Denmark

^④Department of Pathology, Aarhus University Hospital, 8200 Aarhus N, Denmark

^⑤Department of Mathematical Modeling and Machine Learning, University of Zurich, 8057 Zurich, Switzerland

*To whom correspondence should be addressed. Email: mark.robinson@mhs.uzh.ch

[†]The first two authors should be regarded as Joint First Authors.



Martin

pasta: Data representations determine spatial statistics options

Uni-variate

Bi-variate

Multi-variate

Categorical (e.g., cell types)

Colocalization of one cell type and at which scale?

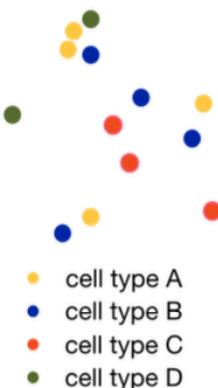
K, L and G functions

Colocalization between two cell types and at which scale?

Cross K, L and G functions

Colocalization of one cell to a set of other cell types?

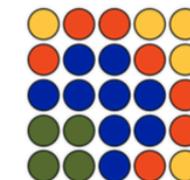
Dot functions



not common

How often are spots of the same cluster neighbouring each other?

Join count statistics



- cluster A
- cluster B
- cluster C
- cluster D

Which clusters are found more frequently neighbouring each other?
Multivariate join count statistics

Numerical (e.g., gene expression)

Uni-variate

Bi-variate

Multi-variate

At which scale is there (spatial) correlation of gene expression?

Mark correlation function

not common



Spatial autocorrelation of a gene?

Moran's I and relatives

Spatial correlation of two genes?
Bivariate Moran's I and relatives



- Gene expression

Spatial correlation of a set of genes?
Multivariate Geary's C



Samuel



Martin



- “a realisation of a spatial point process effectively assumes that the locations of points are not fixed, and that the point pattern is the response or observation of interest.”

Scenario 14.1. *A weather map for Europe displays a symbol for each major city indicating the expected type of weather (e.g., sunny, cloudy, storms).*

Scenario 14.2. *An optical astronomy survey records the sky position and qualitative shape (elliptical, spiral, etc.) of each galaxy in a nearby region of space.*

Scenario 14.3. *Trees in an orchard are examined and their disease status (infected/not infected) is recorded. We are interested in the spatial characteristics of the disease, such as contagion between neighbouring trees.*

Some definitions

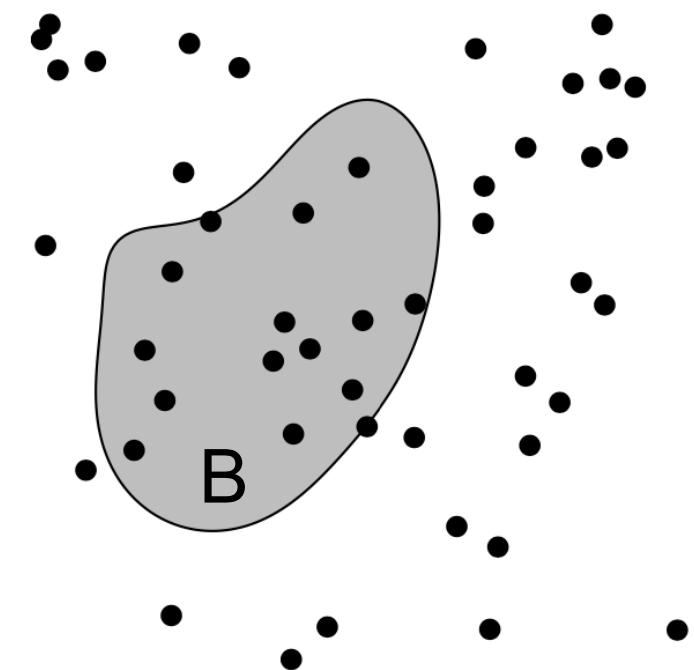
- notation: \mathbf{X} is the point process; \mathbf{x} is the (observed) point pattern
- lambda: intensity function

A point pattern is denoted by a bold lower case letter like \mathbf{x} . It is a set

$$\mathbf{x} = \{x_1, x_2, \dots, x_n\}$$

of points x_i in two-dimensional space \mathbb{R}^2 . The number $n = n(\mathbf{x})$ of points in the pattern is not fixed in advance, and may be any finite nonnegative number *including zero*. In practice, the data points are obviously recorded in some order x_1, \dots, x_n ; but this ordering is artificial, and we treat the pattern \mathbf{x} as an unordered set of points.

$$\mathbb{E}[n(\mathbf{X} \cap B)] = \int_B \lambda(u) du$$





Definitions

- \mathbf{X} is the point process; \mathbf{x} is the (observed) point pattern
- lambda: intensity function
- Complete spatial randomness (CSR) has two properties:

homogeneity: the points have no preference for any spatial location;

independence: information about the outcome in one region of space has no influence on the outcome in other regions.

- More specifically:

homogeneity: the number $n(\mathbf{X} \cap B)$ of random points falling in a test region B has mean value $\mathbb{E}n(\mathbf{X} \cap B) = \lambda |B|$;

independence: for test regions B_1, B_2, \dots, B_m which do not overlap, the counts $n(\mathbf{X} \cap B_1), \dots, n(\mathbf{X} \cap B_m)$ are independent random variables;



Couple more definitions

- Inhomogeneity

The *inhomogeneous Poisson point process* with intensity function $\lambda(u)$ is defined by the following properties:

intensity function: the expected number of points falling in a region B is the integral $\mu = \int_B \lambda(u) du$ of the intensity function $\lambda(u)$ over the region B ;

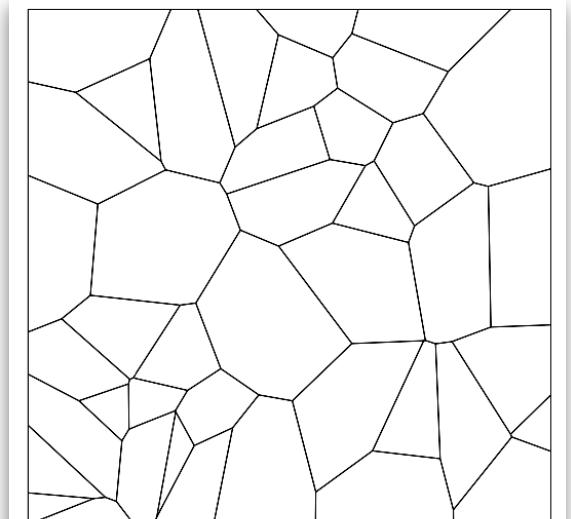
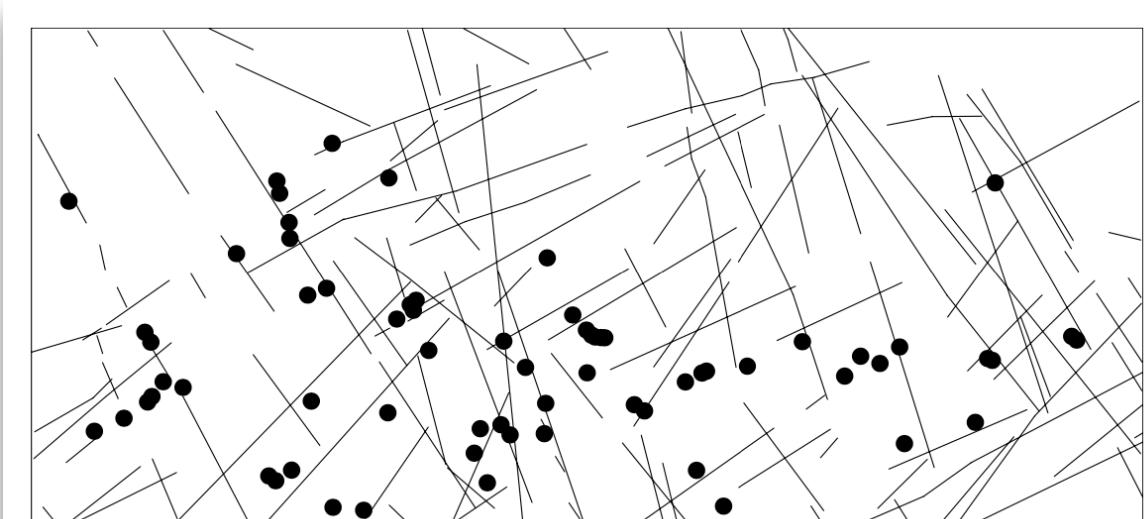
independence: if space is divided into non-overlapping regions, the random patterns inside these regions are independent of each other;

Poisson-distributed counts: the random number of points falling in a given region has a *Poisson* probability distribution;

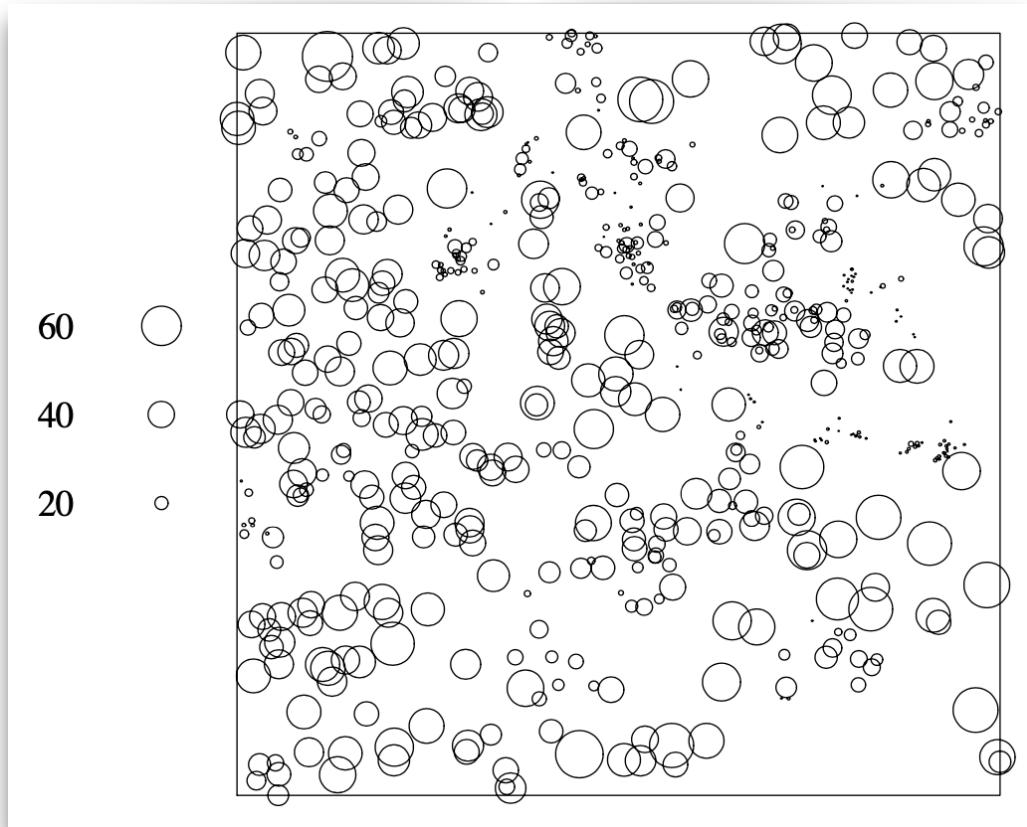


Extensions of point patterns

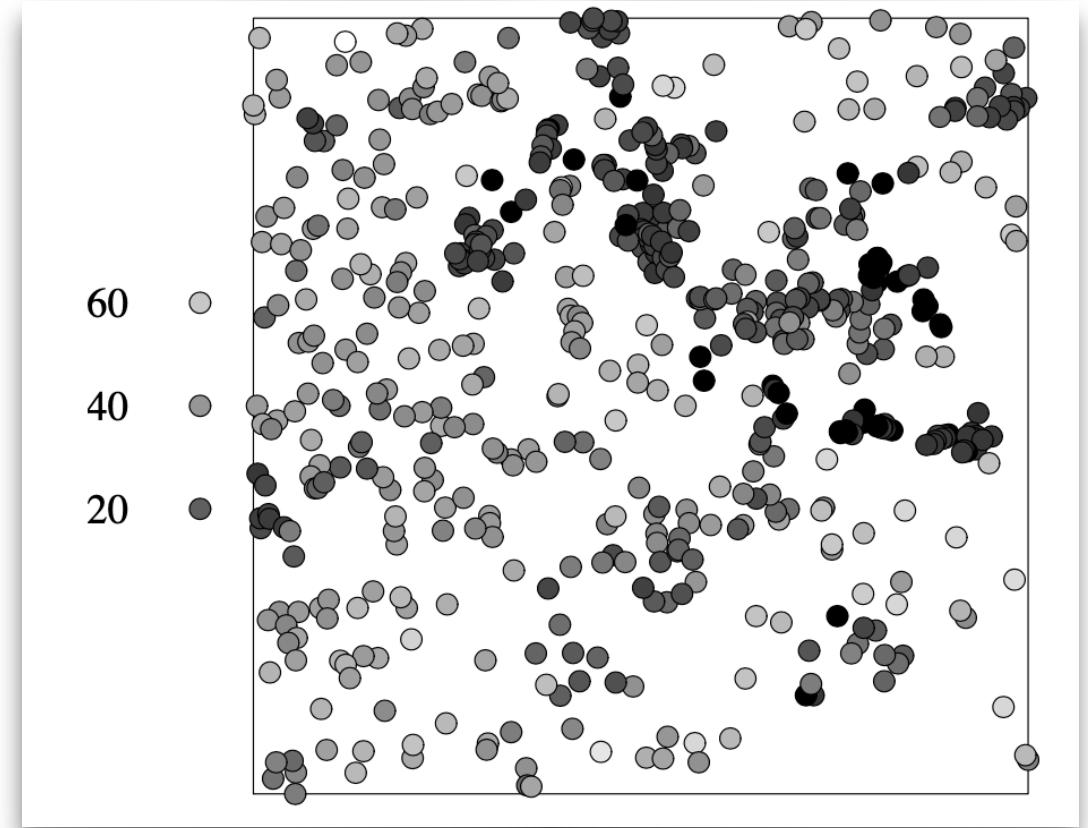
- (patterns can be regions/lines, not just points; can be in higher dimension (e.g., 3D); temporal component .. most of the methods I discuss have extensions; not discussed here)
- marks —> marked point process
- types —> multi-type point process
- covariates



Marks

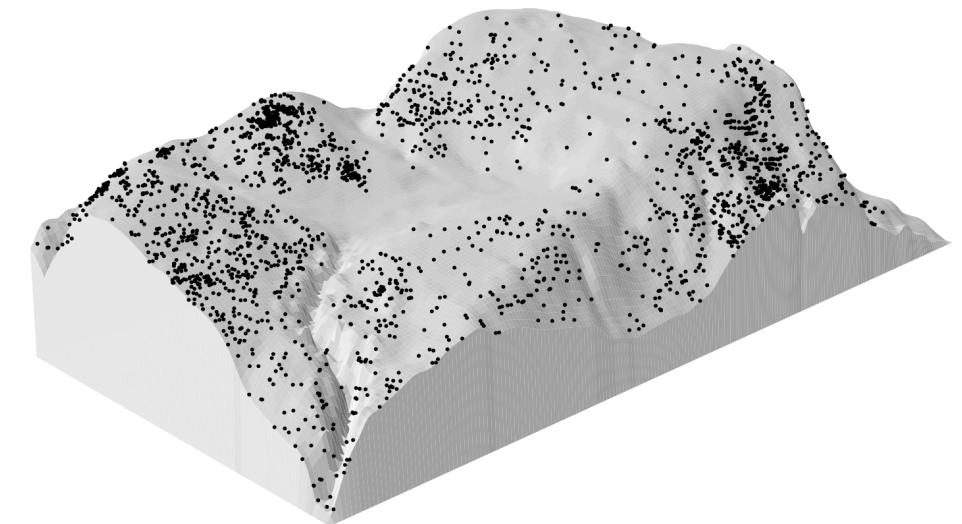
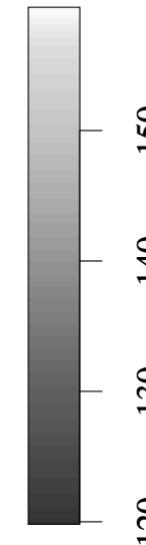
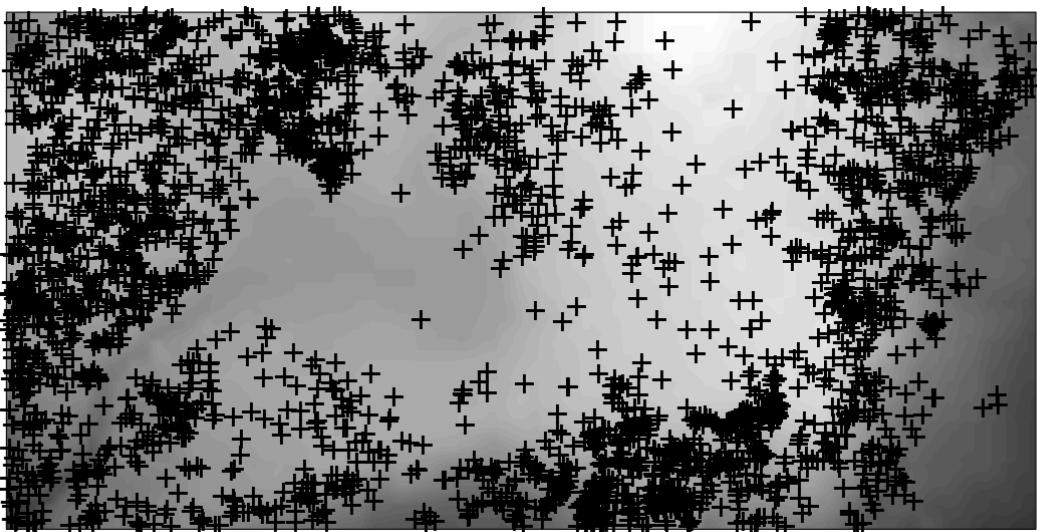


Longleaf pines dataset (location and diameter)





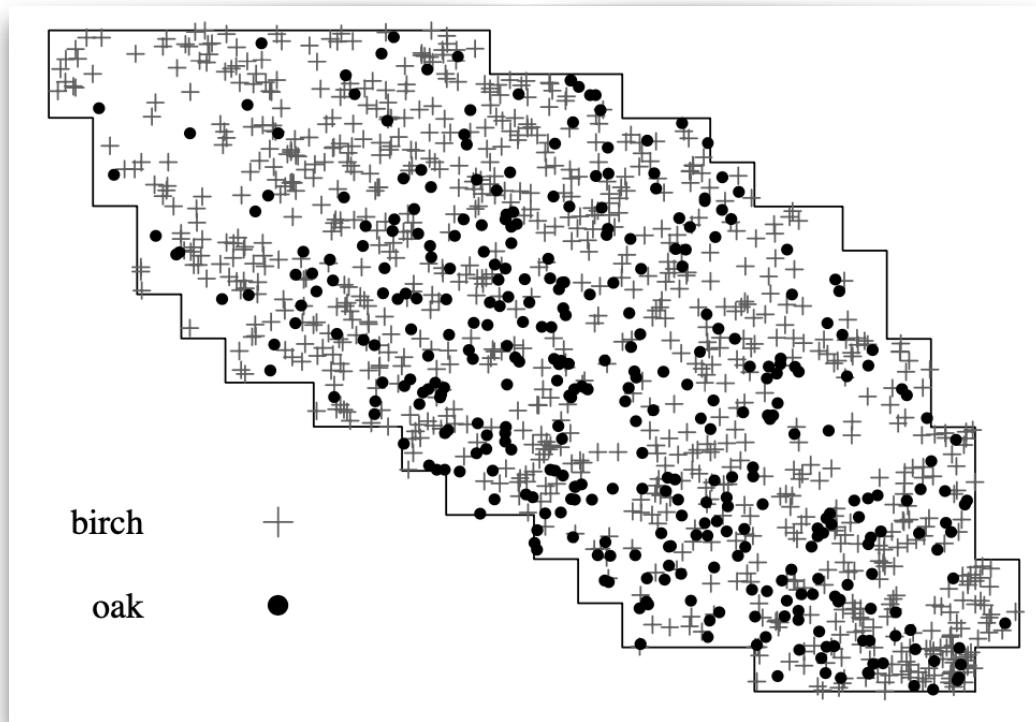
Covariates



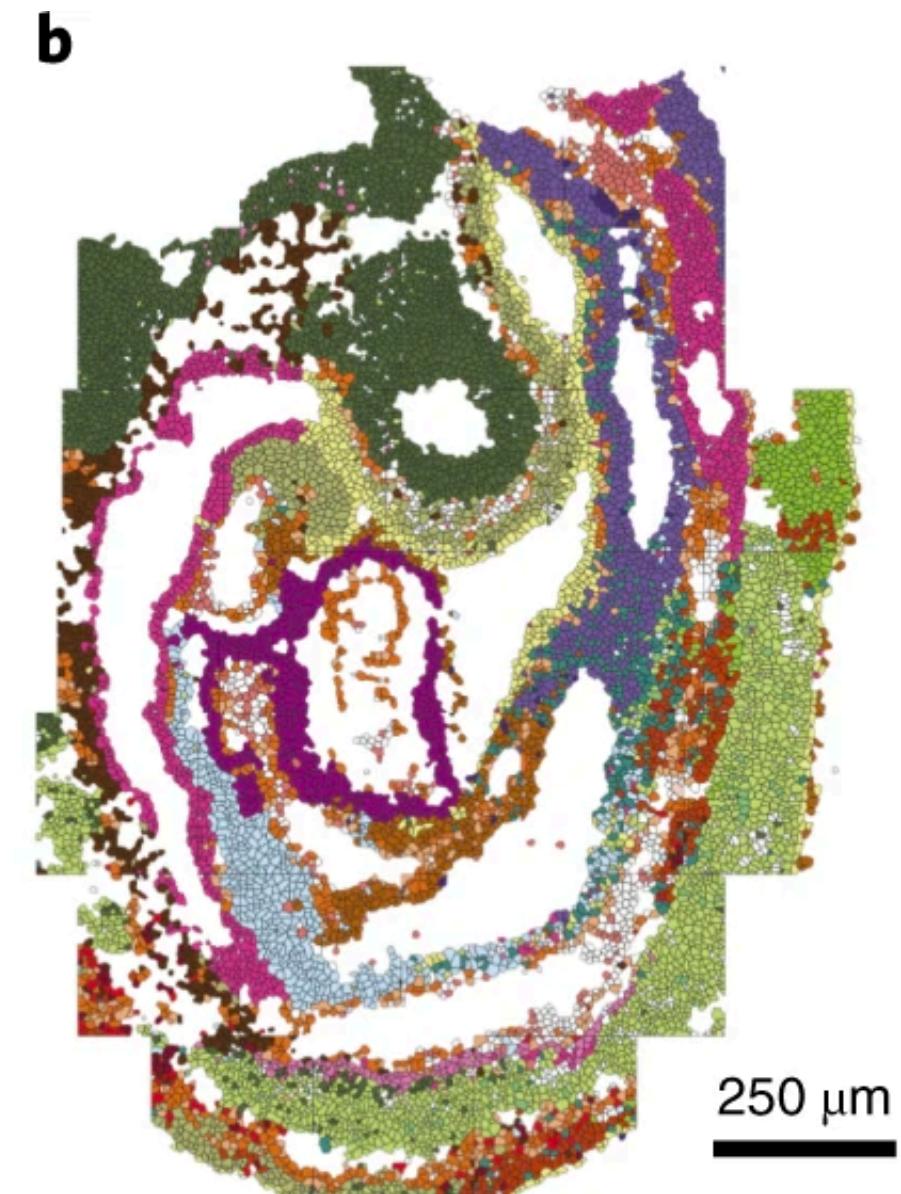
Locations of Beilschmiedia pendula trees (+) and terrain elevation (greyscale) in a 1000×500 metre survey plot in Barro Colorado Island.



Multi-type point patterns



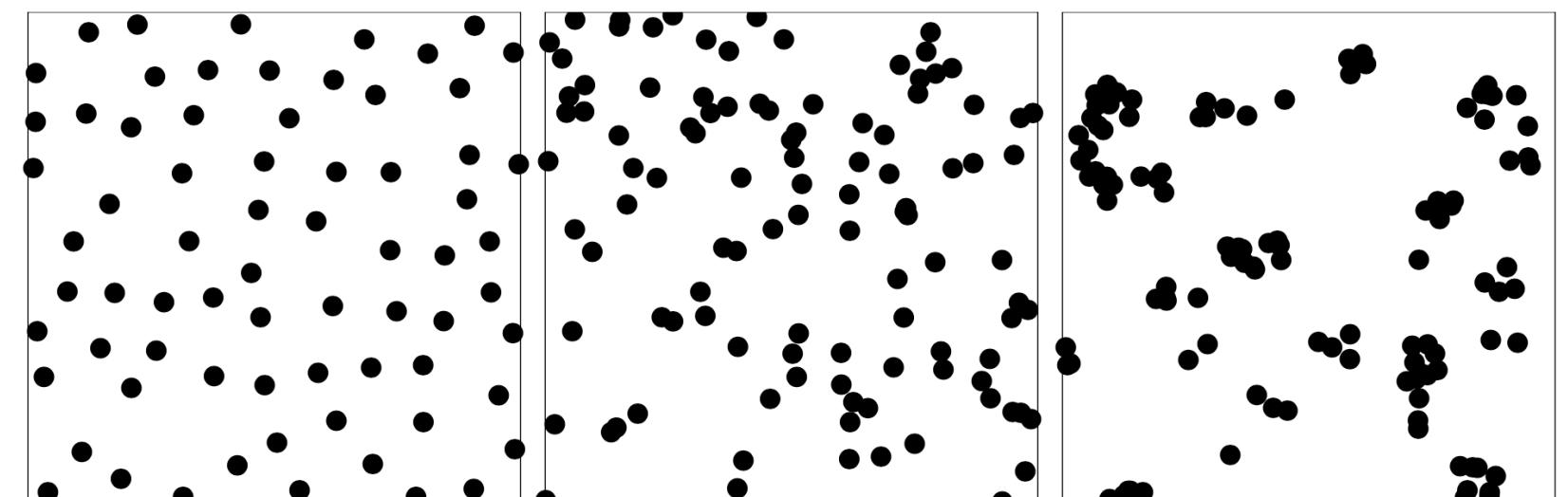
- Allantois
- Anterior somitic tissues
- Blood progenitors
- Cardiomyocytes
- Caudal mesoderm
- Cranial mesoderm
- Definitive endoderm
- Dermomyotome
- Endothelium
- Erythroid
- ExE endoderm
- Forebrain/midbrain/hindbrain
- Gut tube
- Hematoendothelial progenitors
- Intermediate mesoderm
- Lateral plate mesoderm
- Mixed mesenchymal mesoderm
- Neural crest
- NMP
- Presomitic mesoderm
- Sclerotome
- Spinal cord
- Splanchnic mesoderm
- Surface ectoderm



<https://www.nature.com/articles/s41587-021-01006-2>

What is a point pattern?

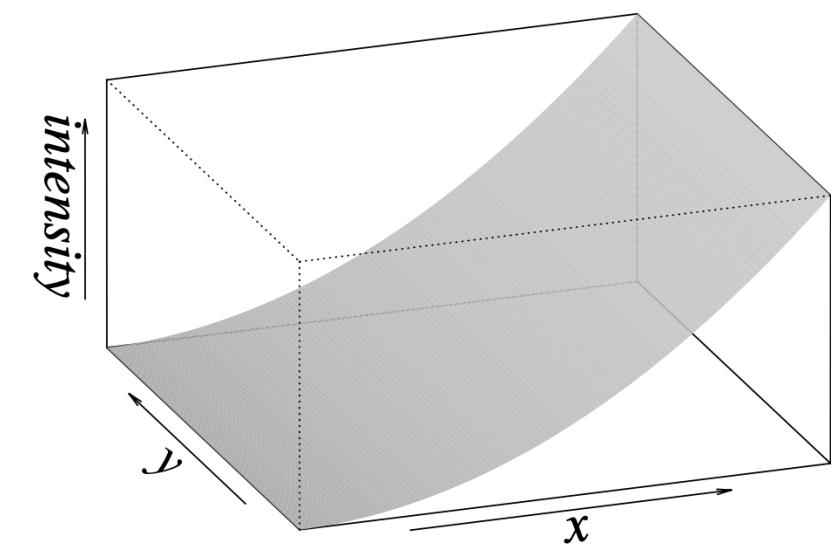
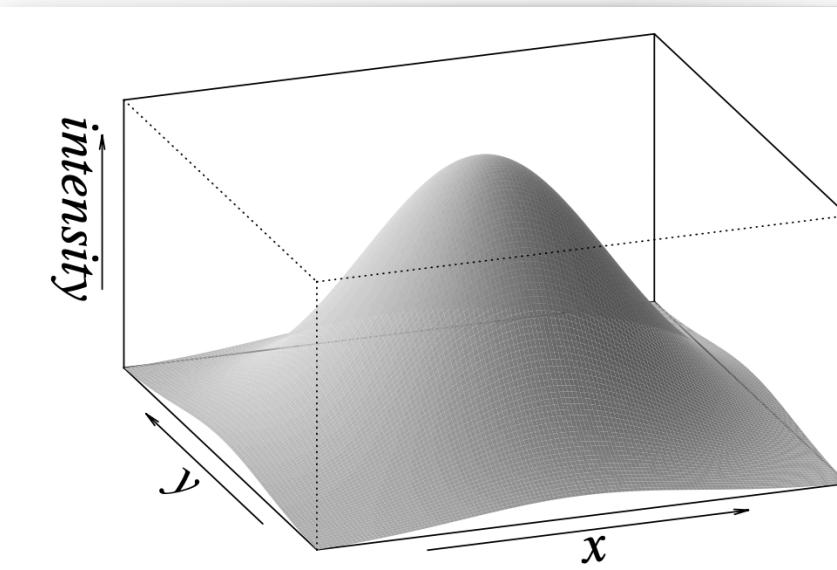
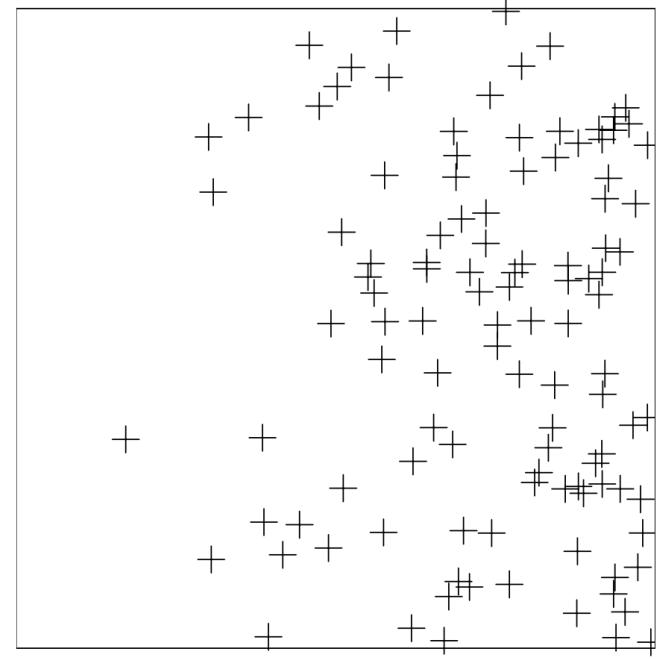
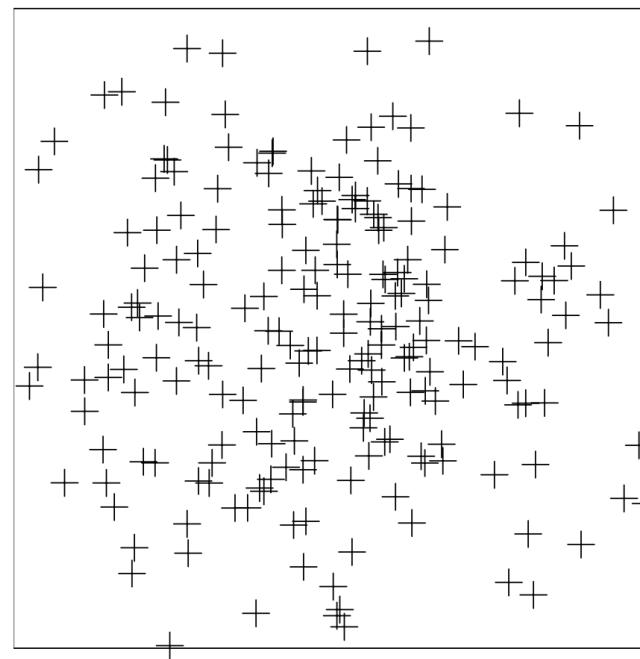
- “a realisation of a spatial point process effectively assumes that the locations of points are not fixed, and that the point pattern is the response or observation of interest.”
- Which of these is homogeneous?
- Which of these is completely spatially random (CSR)?
- Which of these is clustered?
- Which of these is not independent?





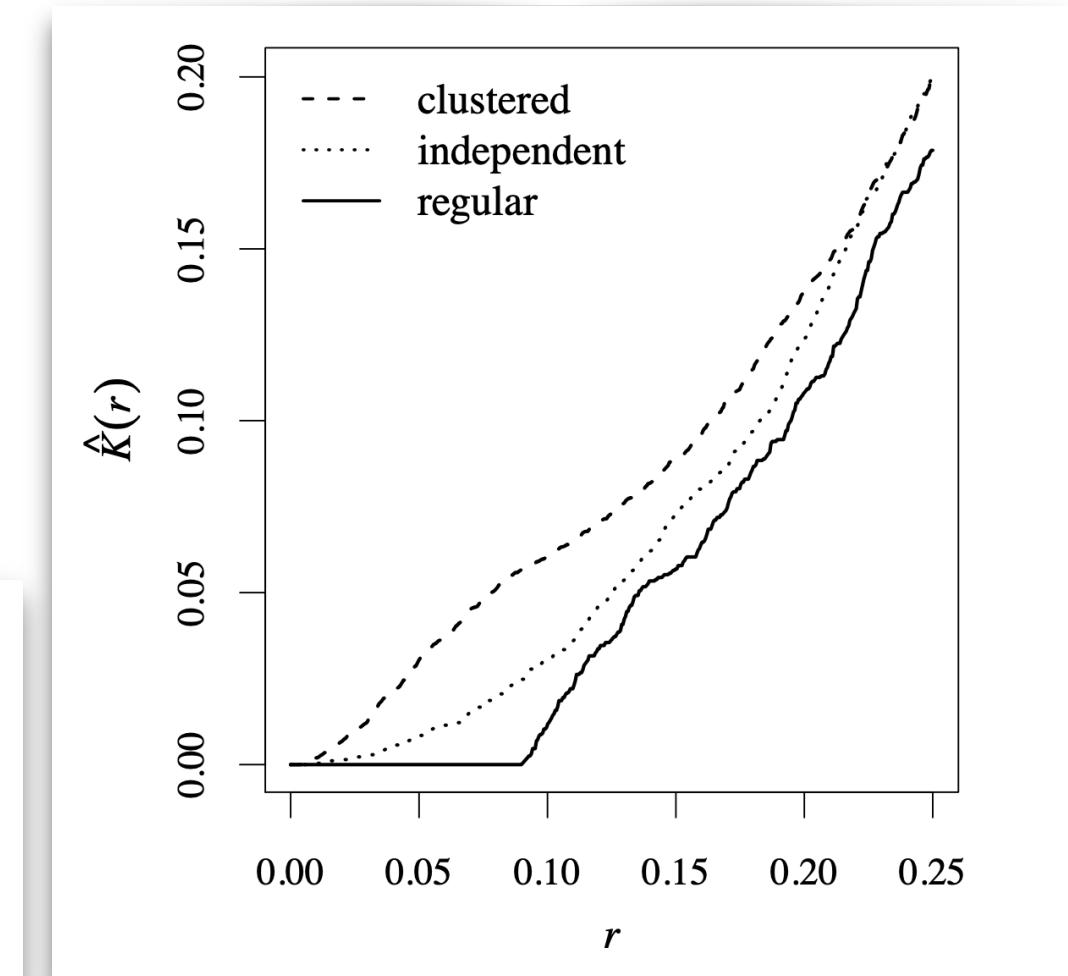
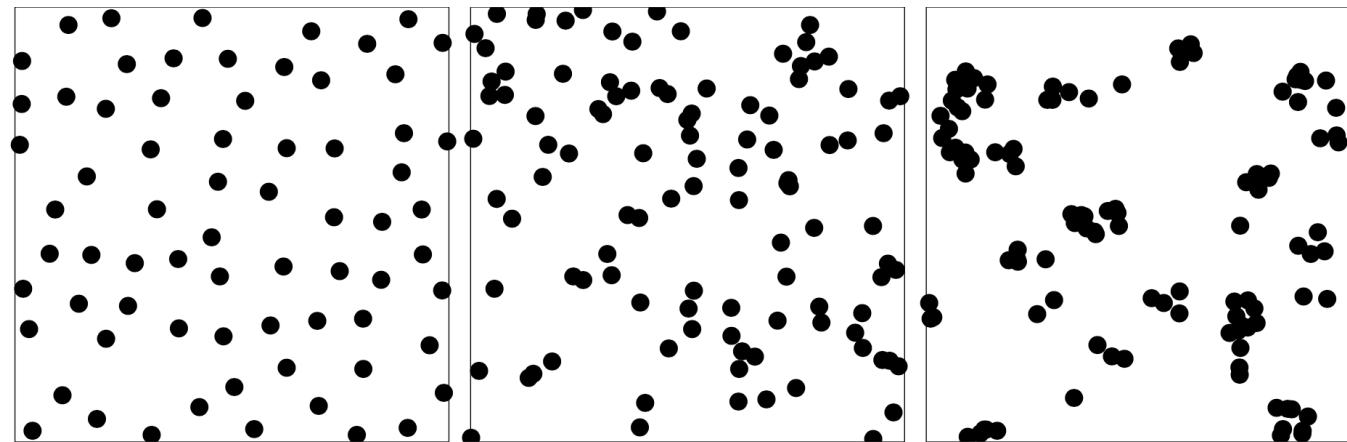
Intensity estimation

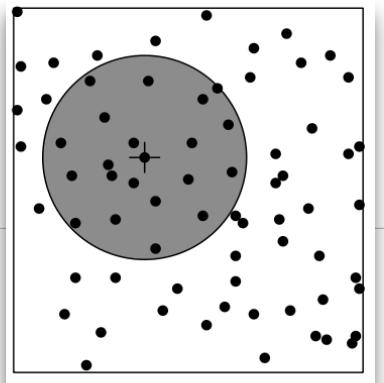
- Methods for intensity estimation



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*





Correlation for **point patterns**

- Ripley's K function
- mathematical definition:

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

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

Definition 7.1. *If \mathbf{X} is a stationary point process, with intensity $\lambda > 0$, then for any $r \geq 0$*

$$K(r) = \frac{1}{\lambda} \mathbb{E} [t(u, r, \mathbf{X}) \mid u \in \mathbf{X}] \tag{7.6}$$

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



What about correlation and intensity together?

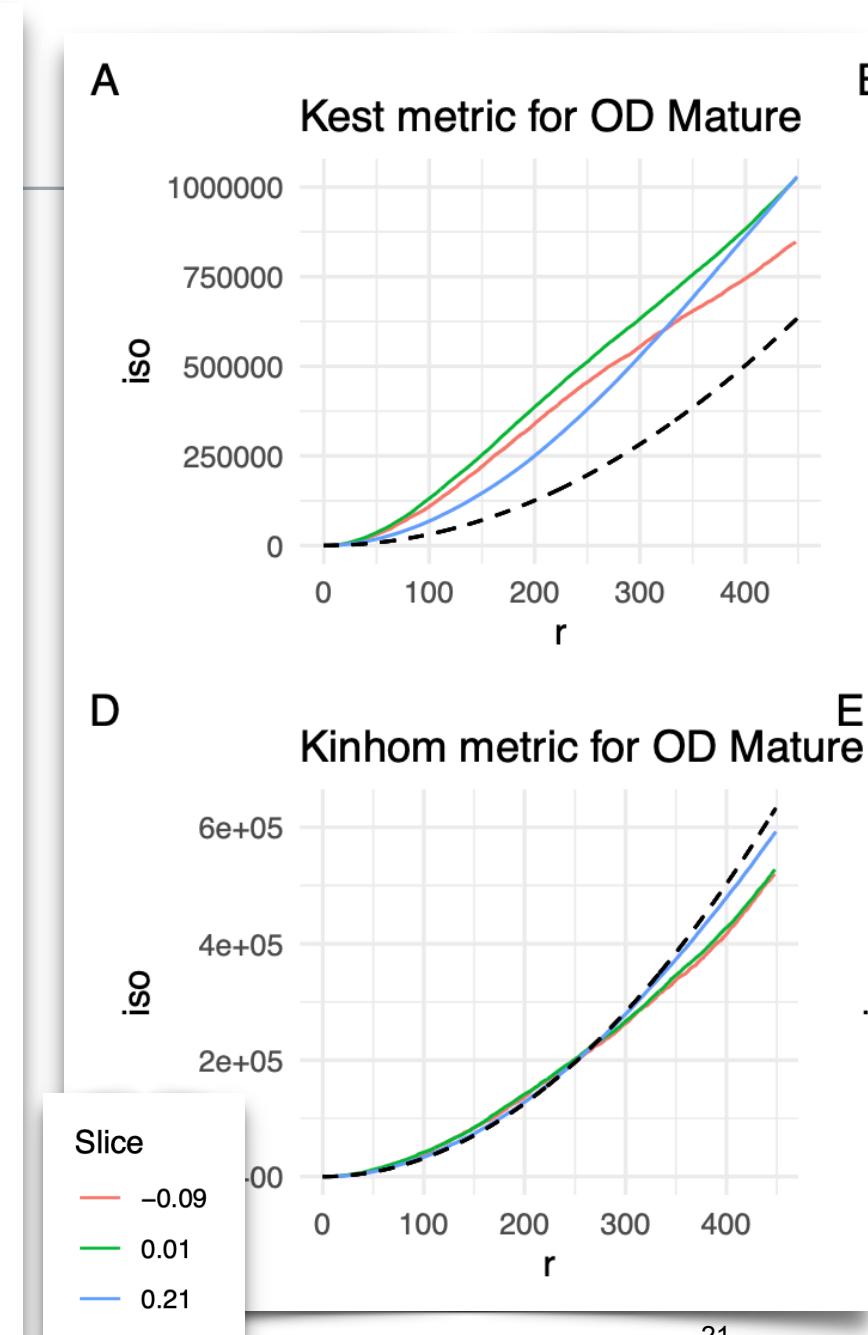
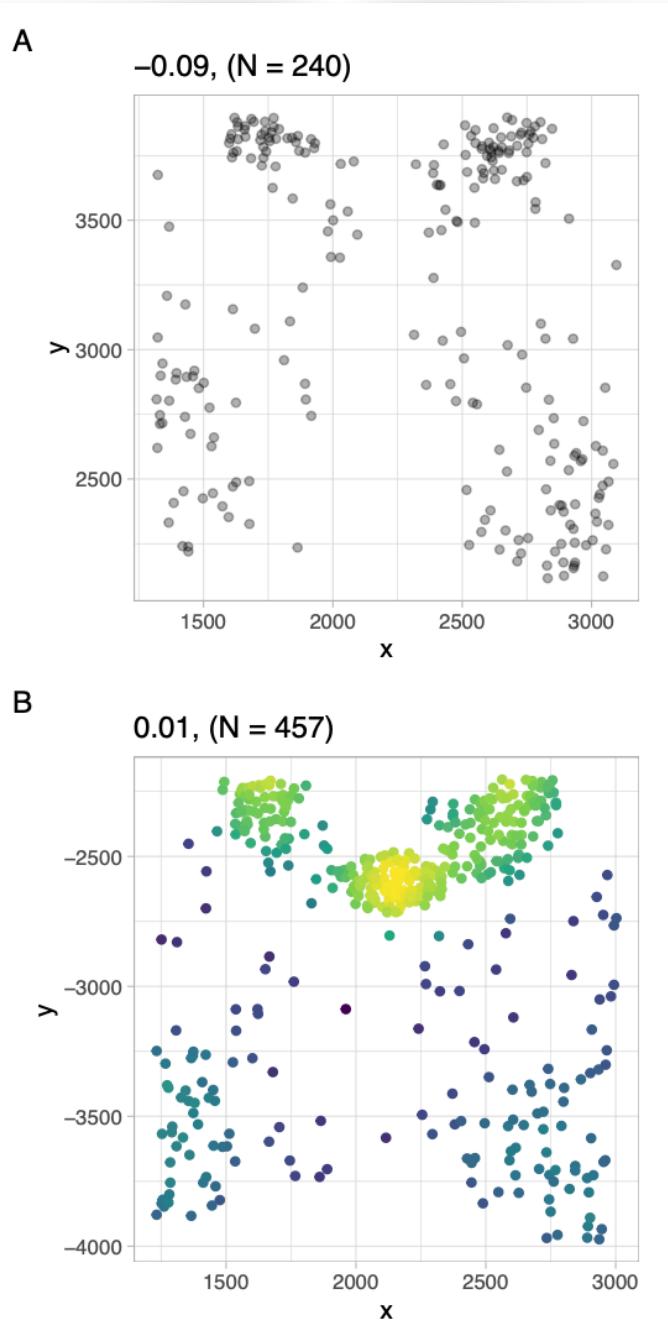
- inhomogeneous correlation functions
- edge correction

$$\widehat{K}_{inhom}(r) = \frac{1}{D^p |W|} \sum_i \sum_{j \neq i} \frac{\mathbf{1}\{||x_i - x_j|| \leq r\}}{\widehat{\lambda}(x_i) \widehat{\lambda}(x_j)} e(x_i, x_j; r)$$

- n.b. confounding of correlation and intensity (next slide)

Confounding between clustering and intensity

- Whether you assume homogeneity or not (in the K-function calculation) can have a big impact on the estimated curves
- Are these cells clustered or have different intensity? Hard to tell.





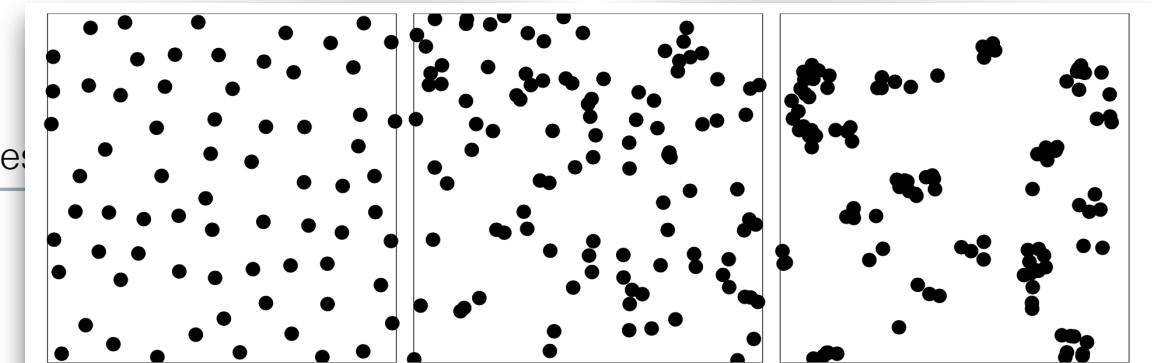
Extensions of the K function (1)

- multitype K-function $K_{ij}(r)$, also called the bivariate or cross-type K-function, is **the expected number of points of type j lying within a distance r of a typical point of type i**, standardised by dividing by the intensity of points of type j.

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

$$K_{ij}(r) = \frac{1}{\lambda_j} \mathbb{E} \left[t(u, r, \mathbf{X}^{(j)}) \mid u \in \mathbf{X}^{(i)} \right]$$

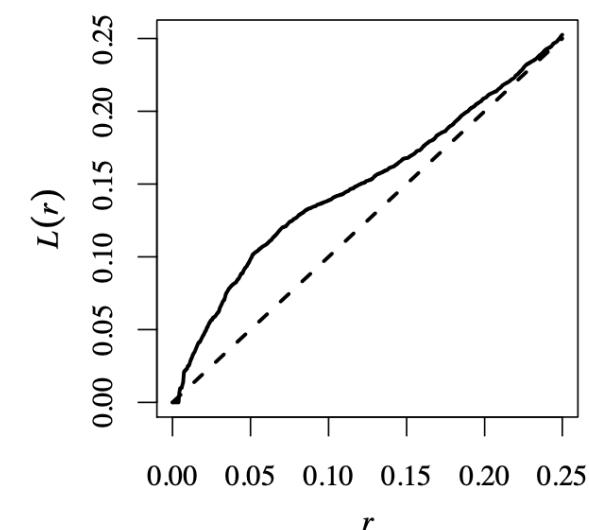
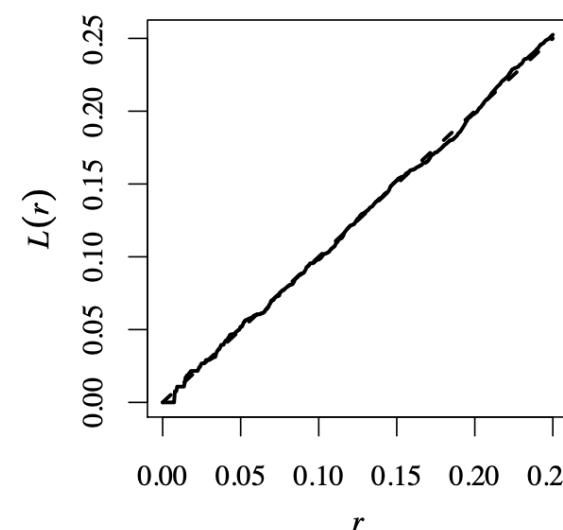
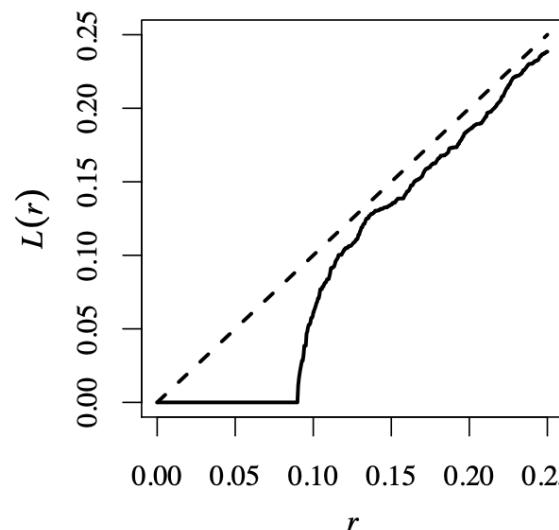
Extensions of the K function (2)



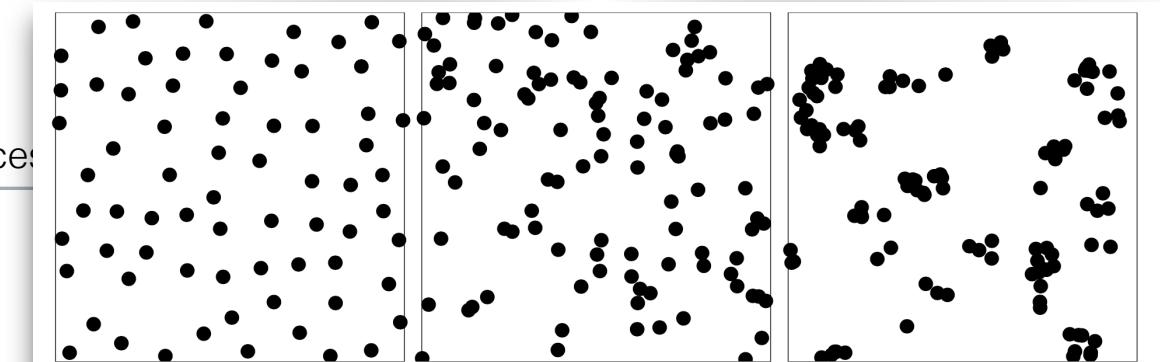
- L and g functions

A commonly used transformation of K proposed by Besag [103] is the **L-function**

$$L(r) = \sqrt{\frac{K(r)}{\pi}}$$



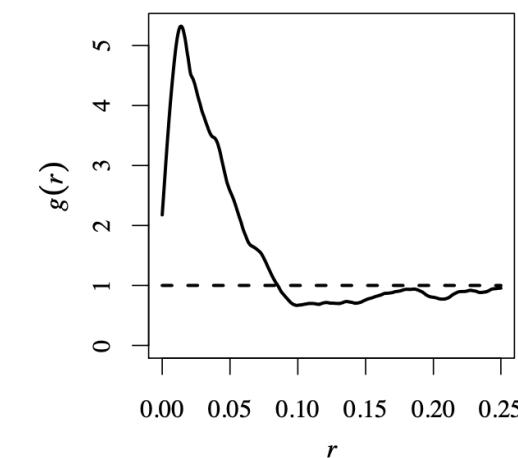
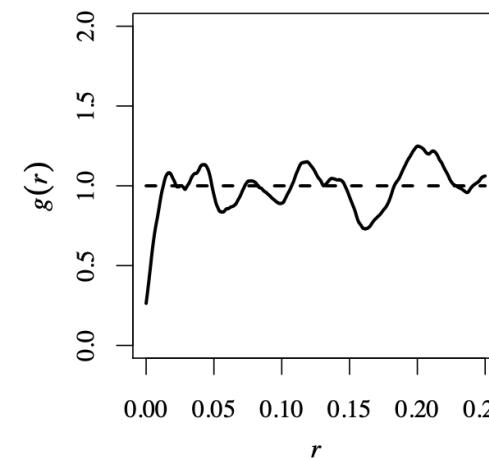
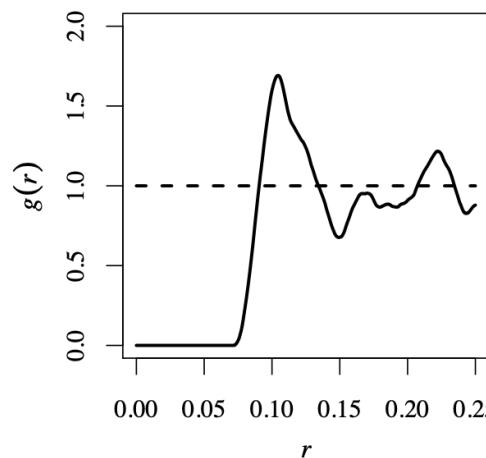
Extensions of the K function (2)



- L and g functions

An alternative tool is the **pair correlation function** $g(r)$ which contains contributions only from interpoint distances *equal to r*. In two dimensions, it can be defined by

$$g(r) = \frac{K'(r)}{2\pi r} \quad (7.22)$$

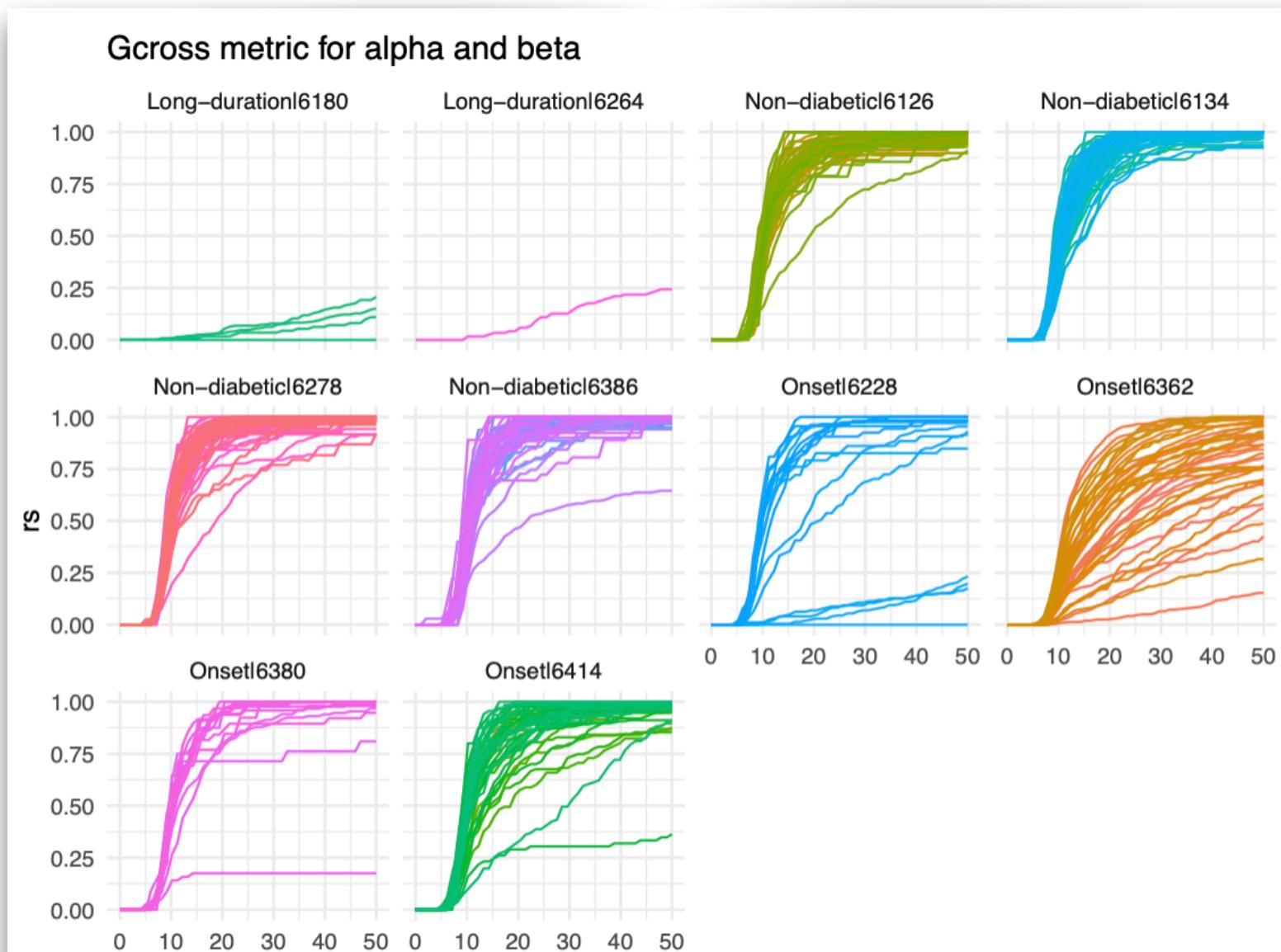


Gross function across 10 samples (multiple FOVs per patient)

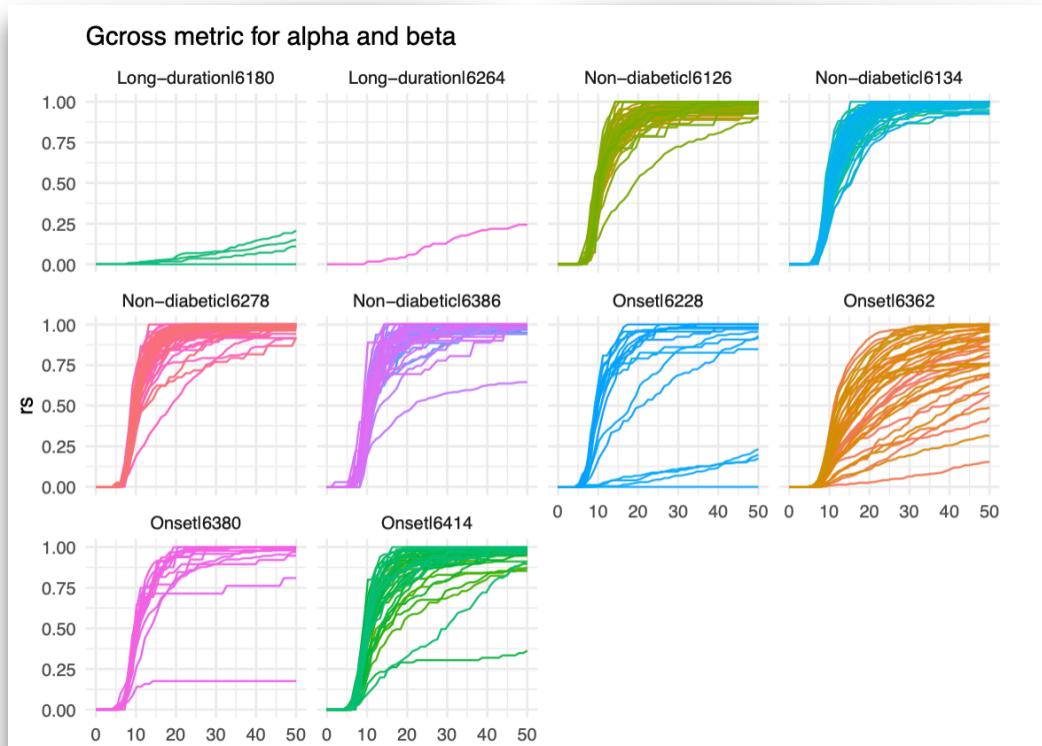


Martin

G-cross:
represents the probability of finding at least one given beta cell within a radius r of any alpha cell



Functional PCA is already useful to decompose FOV-level data



Martin



pasta: Data representations determine spatial statistics options

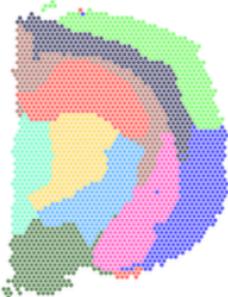
A

Imaging-based

- Targeted
- Higher resolution



STARmap



10X Visium

TECHNOLOGY

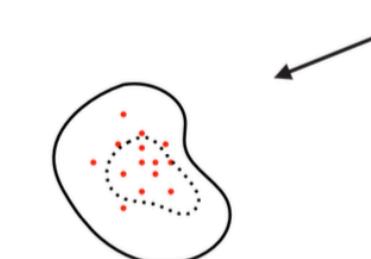
HTS-based

- Untargeted
- Lower resolution

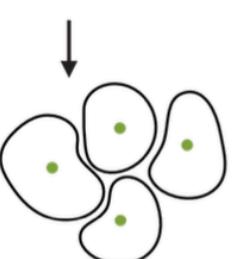


Samuel

B



feature locations



segmentations

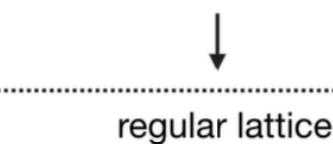
depending on
resolution



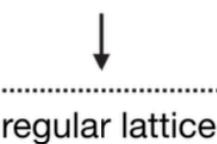
spots / beads / pixels

DATA MODALITY

cell outline



irregular lattice



regular lattice

centroids

point pattern

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, Maite Kuehl ^{3,4}, Reinhard Furrer ⁵, Mark D. Robinson ^{1,*}

¹Department of Molecular Life Sciences and SIB Swiss Institute of Bioinformatics, University of Zurich, 8057 Zurich, Switzerland

²Centro Nacional de Análisis Genómico (CNAG), 08028 Barcelona, Spain

³Department of Clinical Medicine, Aarhus University, 8200 Aarhus N, Denmark

⁴Department of Pathology, Aarhus University Hospital, 8200 Aarhus N, Denmark

⁵Department 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.



Martin



Global autocorrelation measures (Moran's I, Moran's R, Geary's C, etc.)

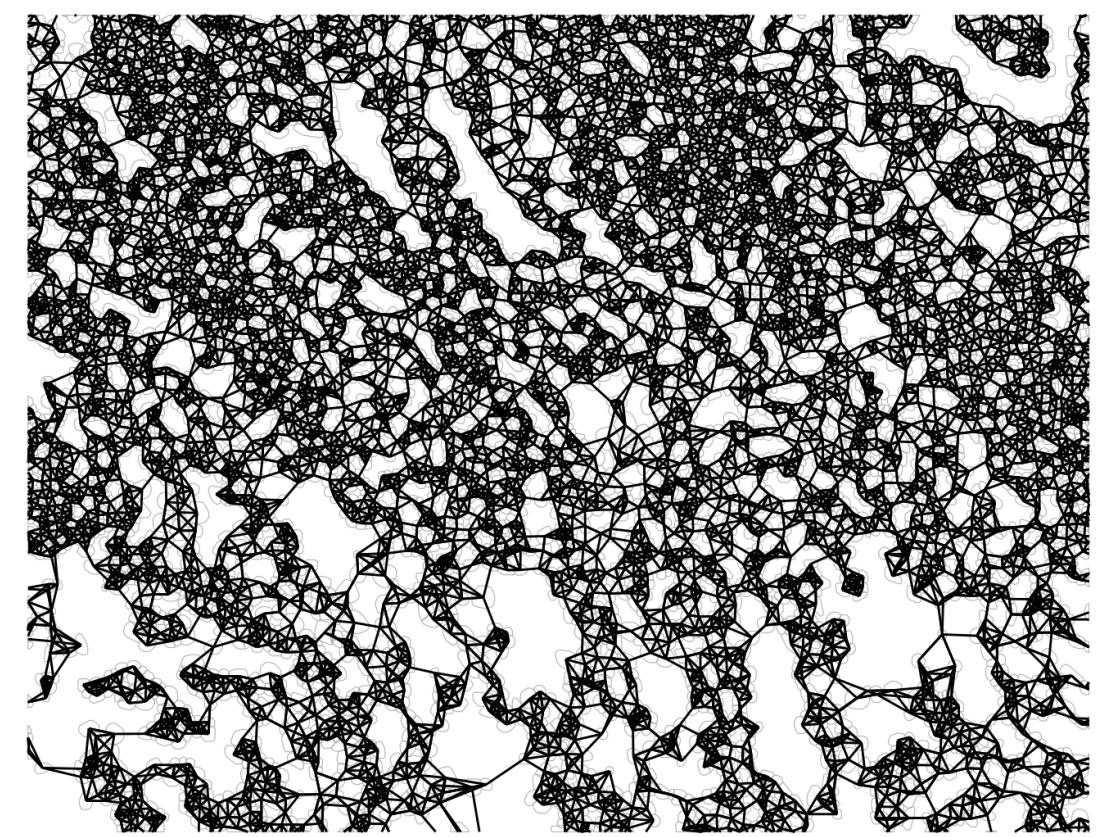
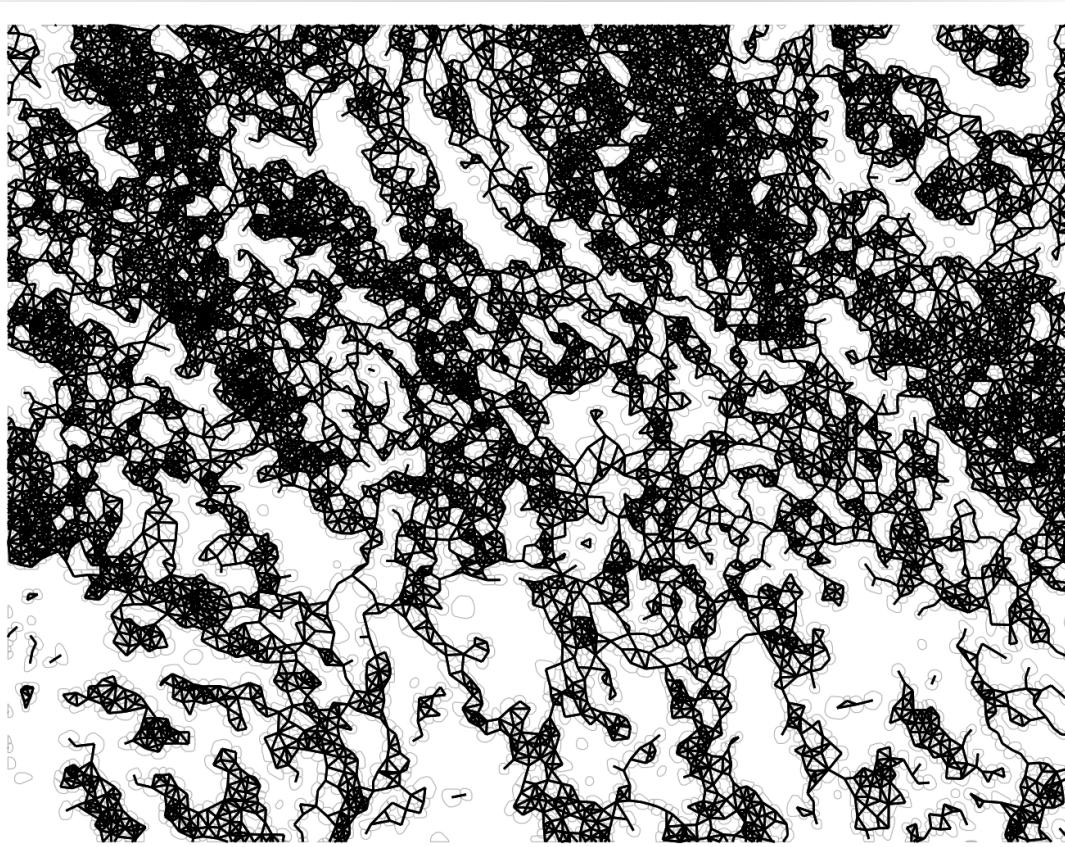
In general, a global spatial autocorrelation measure has the form of a double sum over all locations i, j

$$\sum_i \sum_j f(x_i, x_j) w_{ij}$$

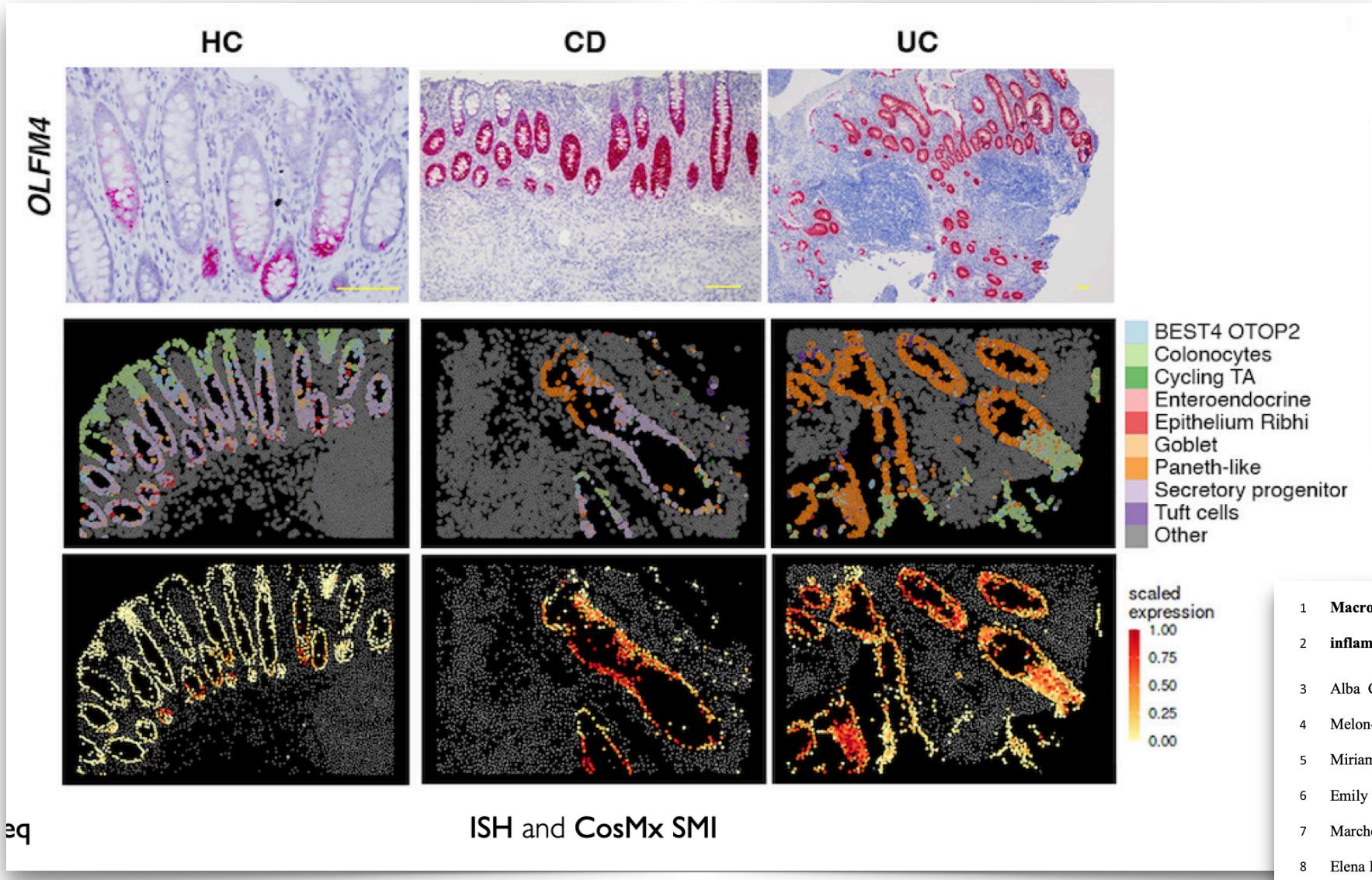
where $f(x_i, x_j)$ is the measure of association between features of interest and w_{ij} scales the relationship by a spatial weight as defined in the weight matrix W . If i and j are not neighbours, i.e. we assume they do not have any spatial association, the corresponding element of the weight matrix is 0 (i.e., $w_{ij} = 0$).



Weight matrix options

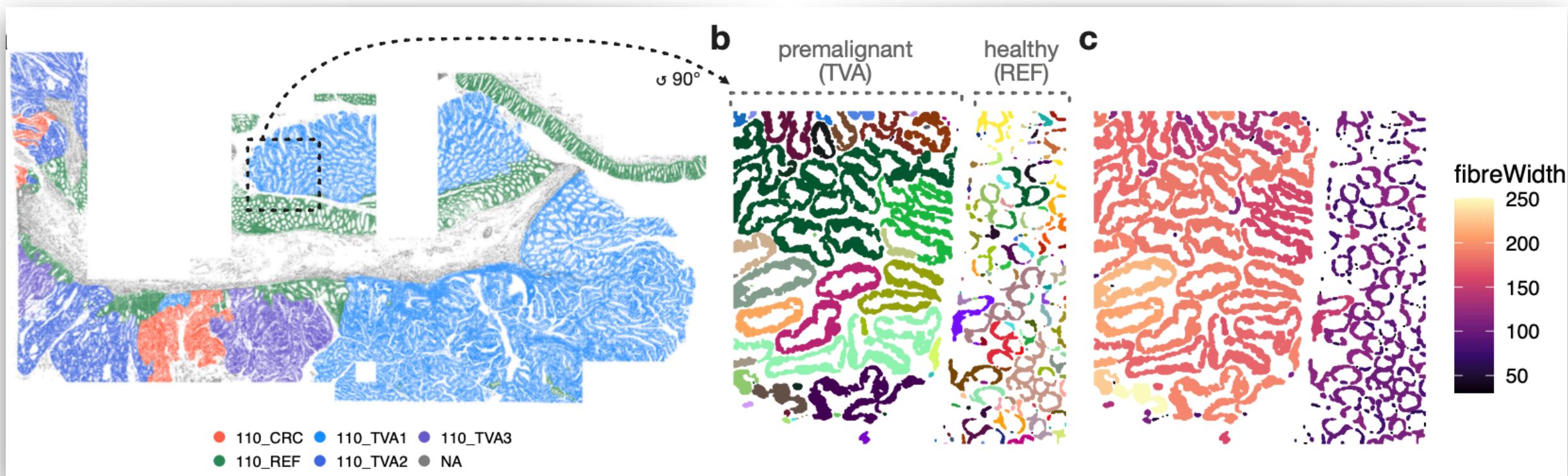


Tissue “structures” are often visible

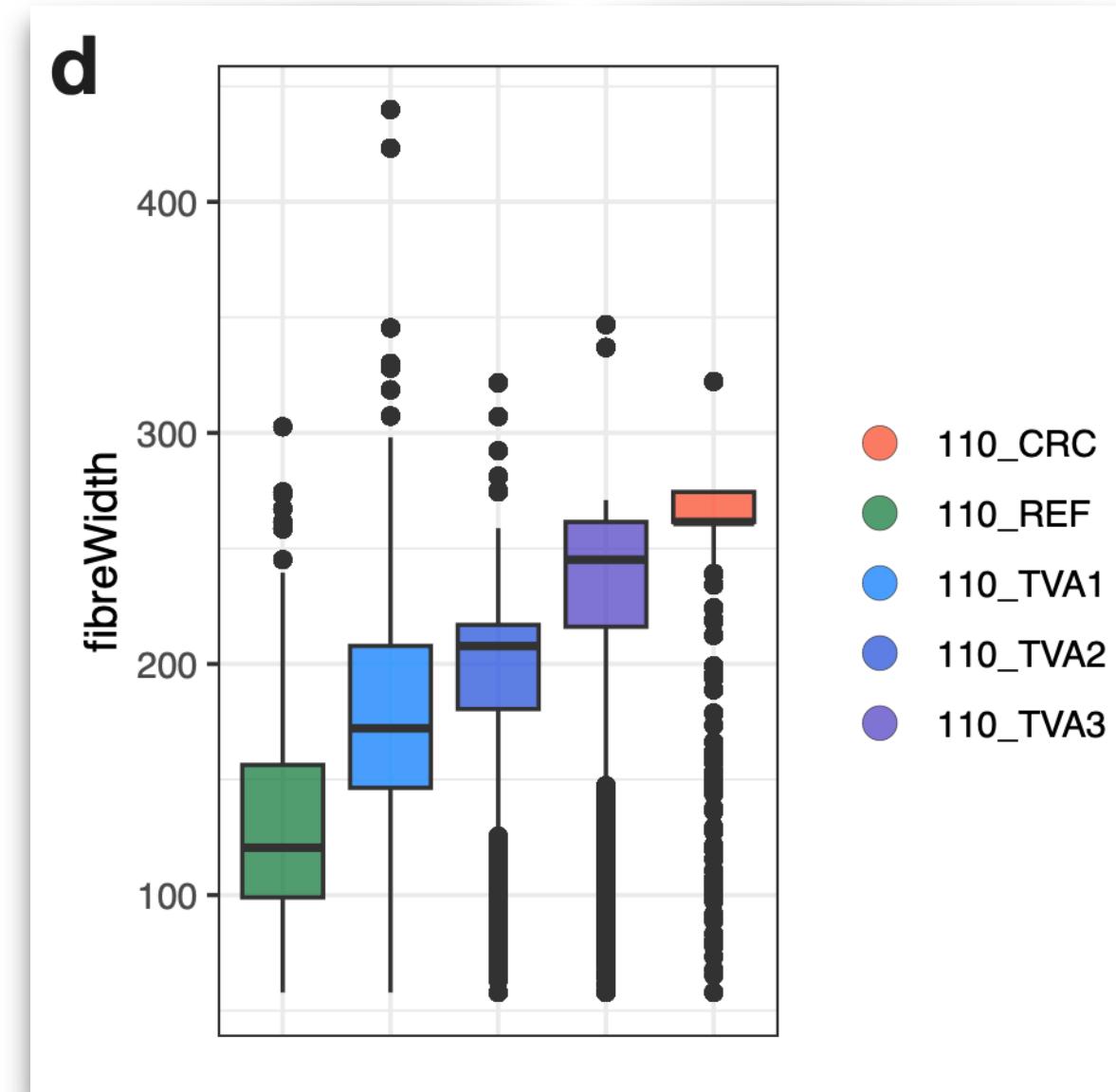
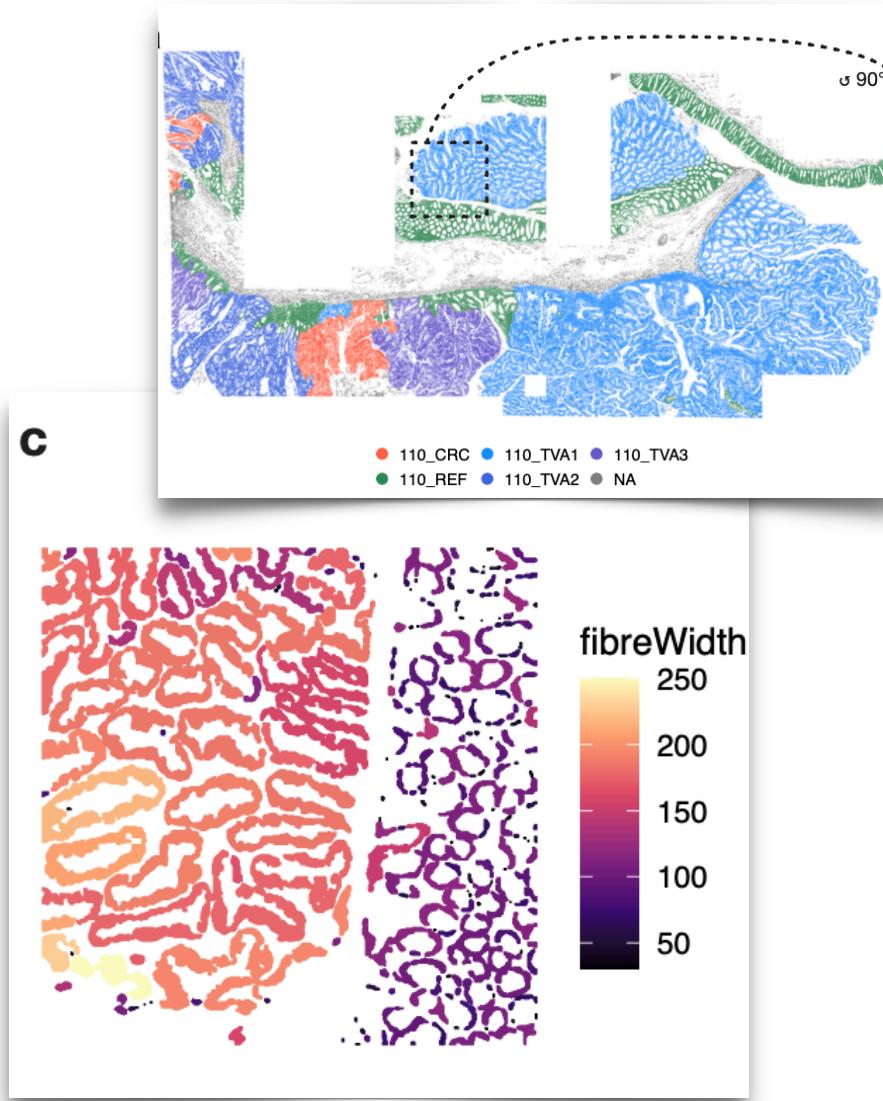


- 1 Macrophage and neutrophil heterogeneity at single-cell spatial resolution in
- 2 inflammatory bowel disease
- 3 Alba Garrido-Trigo^{1,2}, Ana M. Corraliza^{1,2}, Marisol Veny^{1,2}, Isabella Dotti^{1,2}, Elisa
- 4 Melon-Ardanaz^{1,2}, Aina Rill³, Helena L. Crowell⁴, Ángel Corbí⁵, Victoria Gudiño^{1,2},
- 5 Miriam Esteller^{1,2}, Iris Álvarez-Teubel^{1,2}, Daniel Aguilar^{1,2}, M Carme Masamunt^{1,2},
- 6 Emily Killingbeck⁶, Youngmi Kim⁶, Michael Leon⁶, Sudha Visvanathan⁷, Domenica
- 7 Marchese⁸, Ginevra Caratù⁸, 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}

Variation among spatial structures (epithelial example)



Variation among spatial structures (geometric quantifications)



Structures → Reference axis → Expression gradients

