



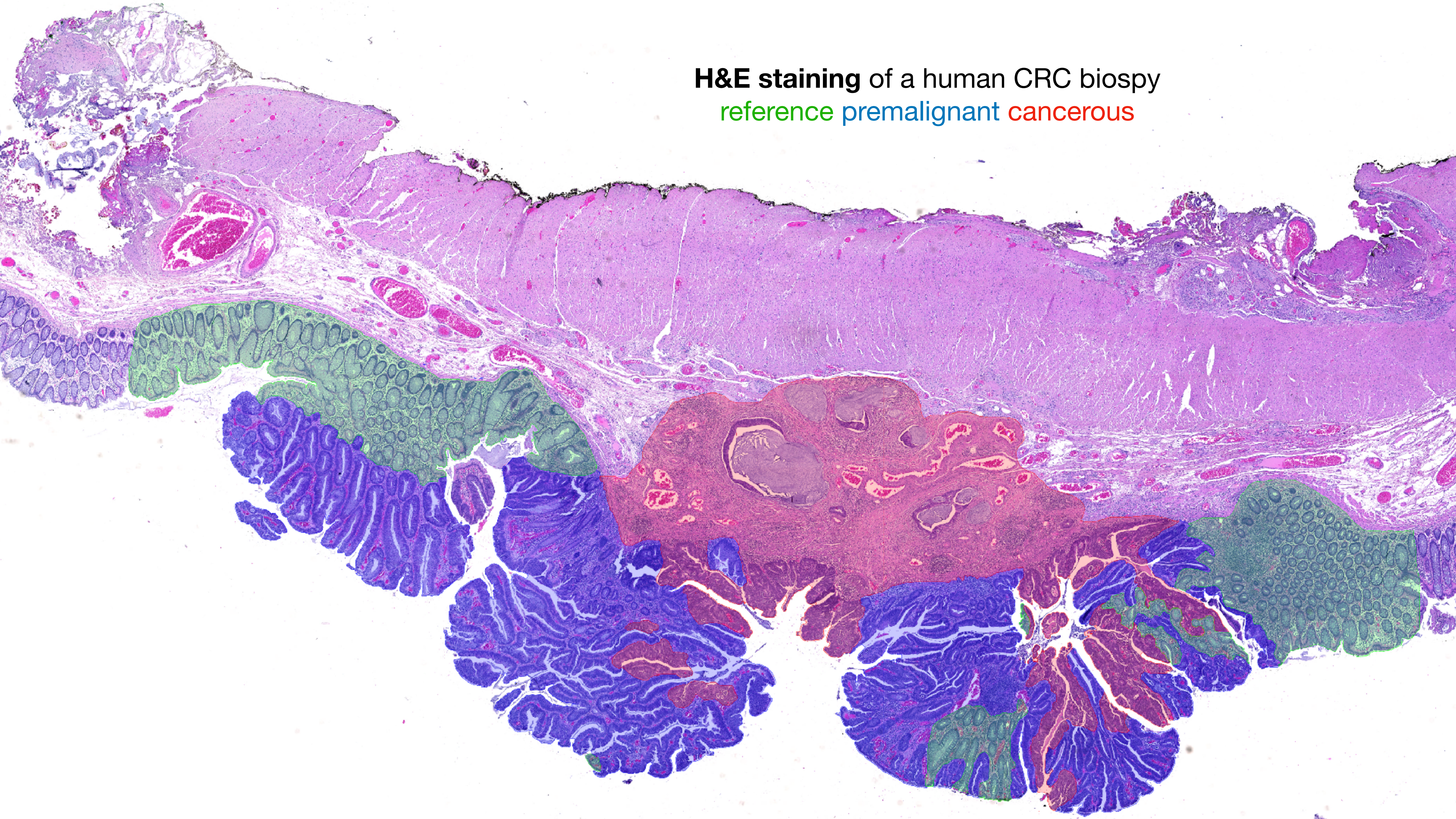
# **(pre)preprocessing** of imaging- based spatial transcriptomics data

Swiss Institute of Bioinformatics  
Spatial Omics Data Analysis

Helena L. Crowell — January 21<sup>st</sup>, 2025  
in Lausanne, Switzerland

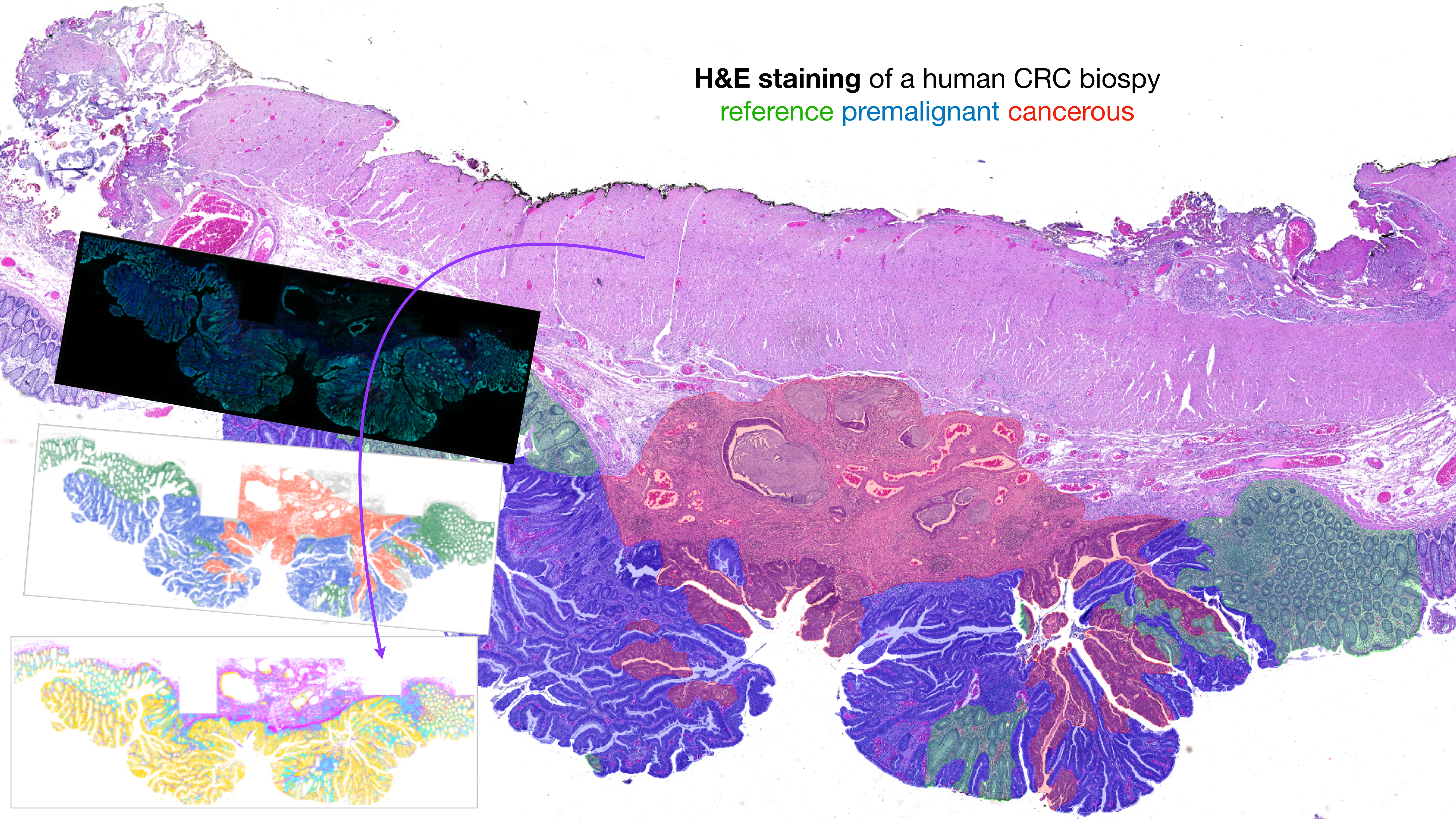


**H&E staining** of a human CRC biopsy  
reference premalignant cancerous





**H&E staining** of a human CRC biopsy  
reference premalignant cancerous

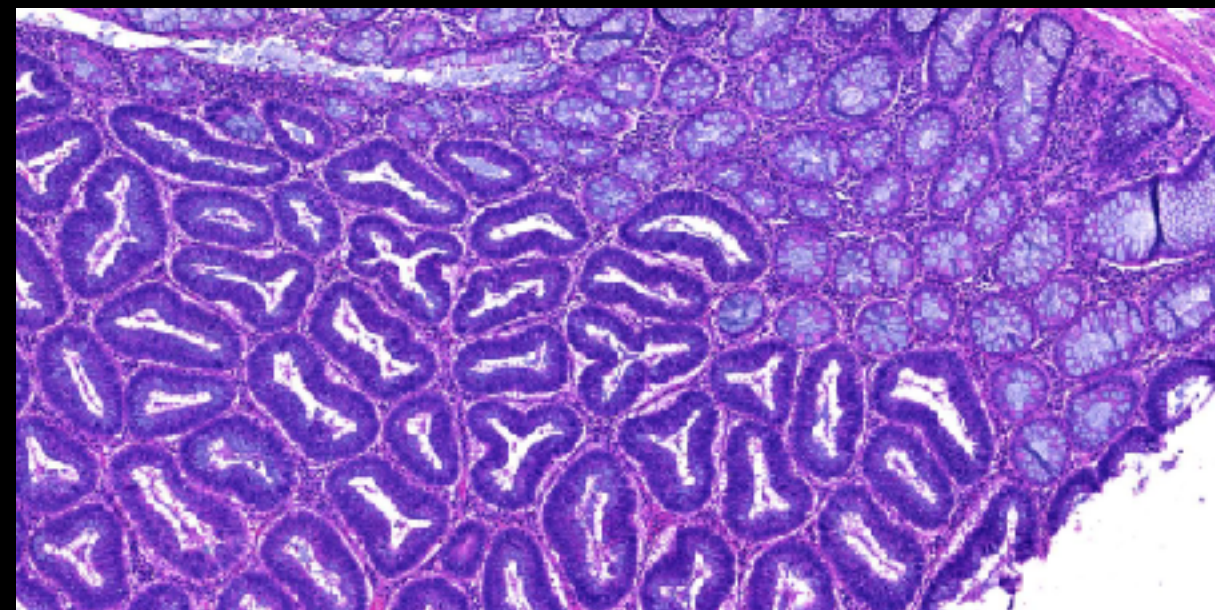




# the path from tissue to pathology ain't easy

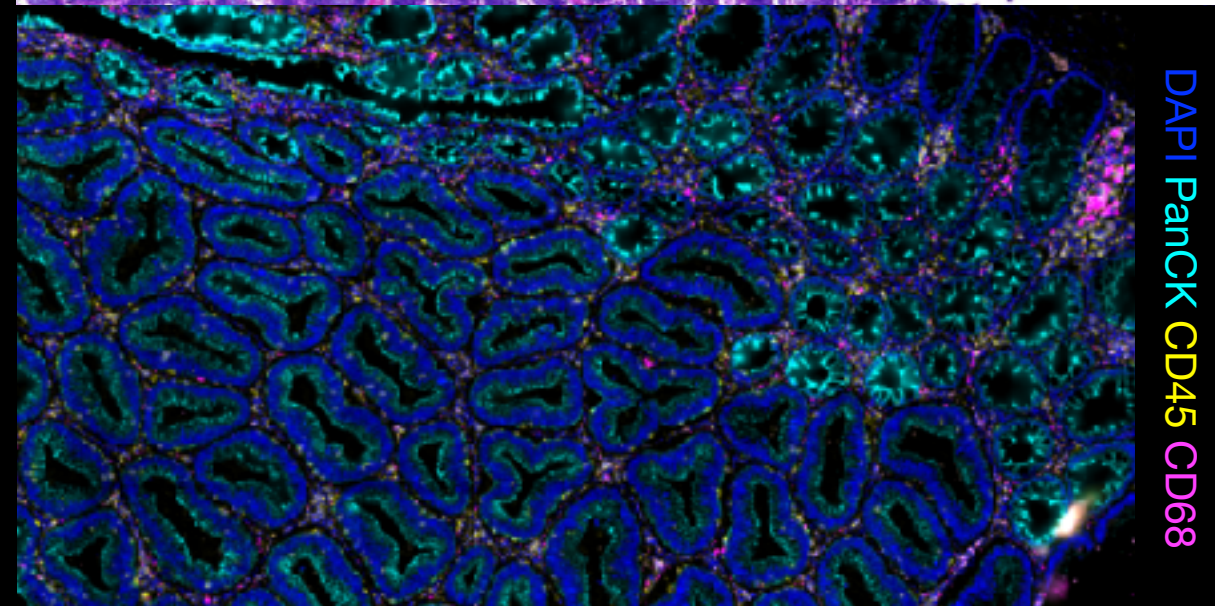
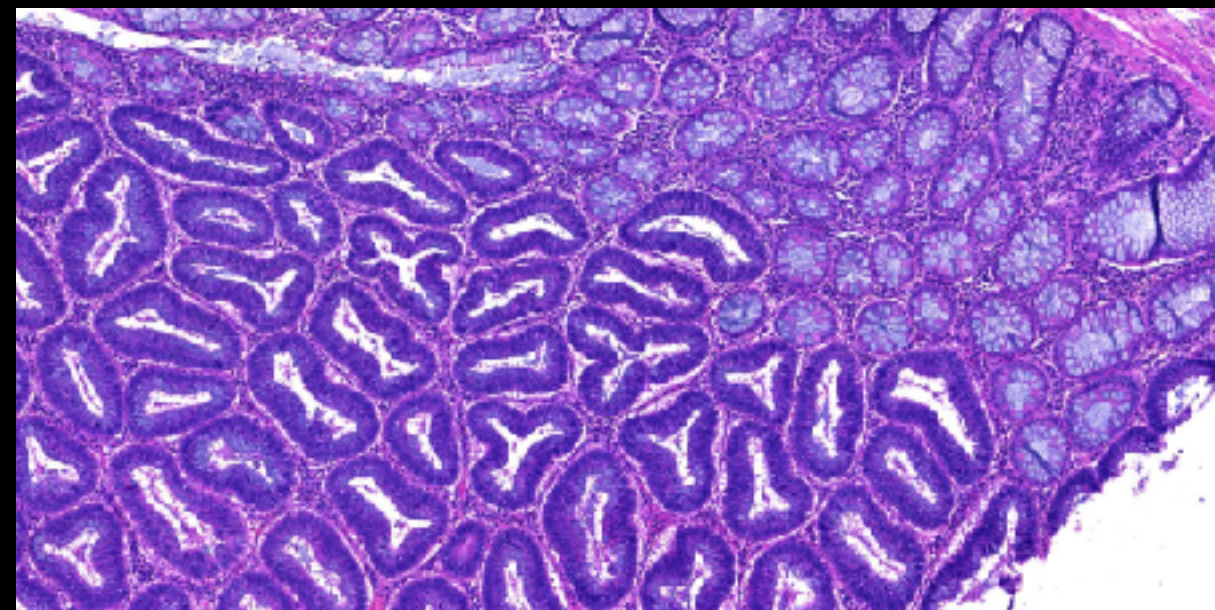
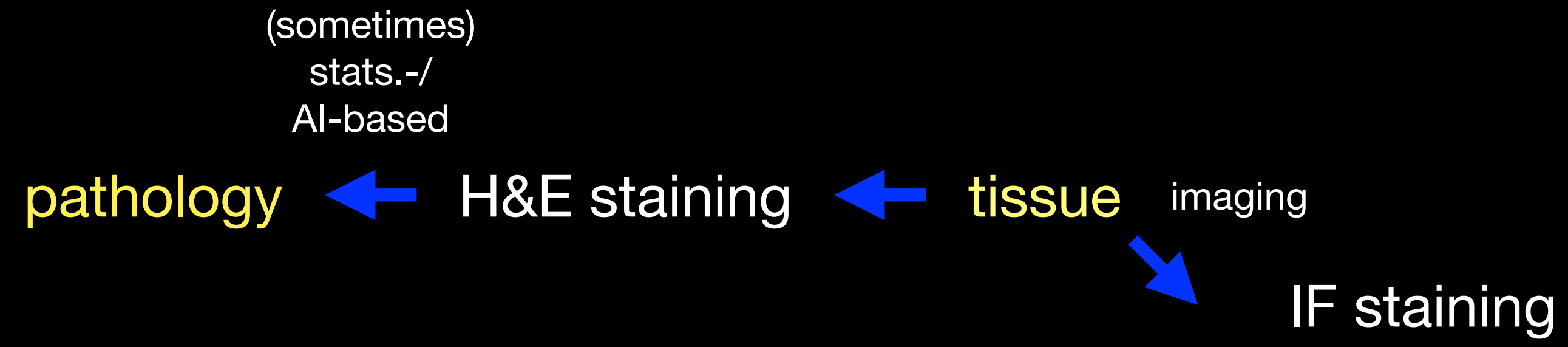
(sometimes)  
stats.-/  
AI-based

pathology ← H&E staining ← tissue





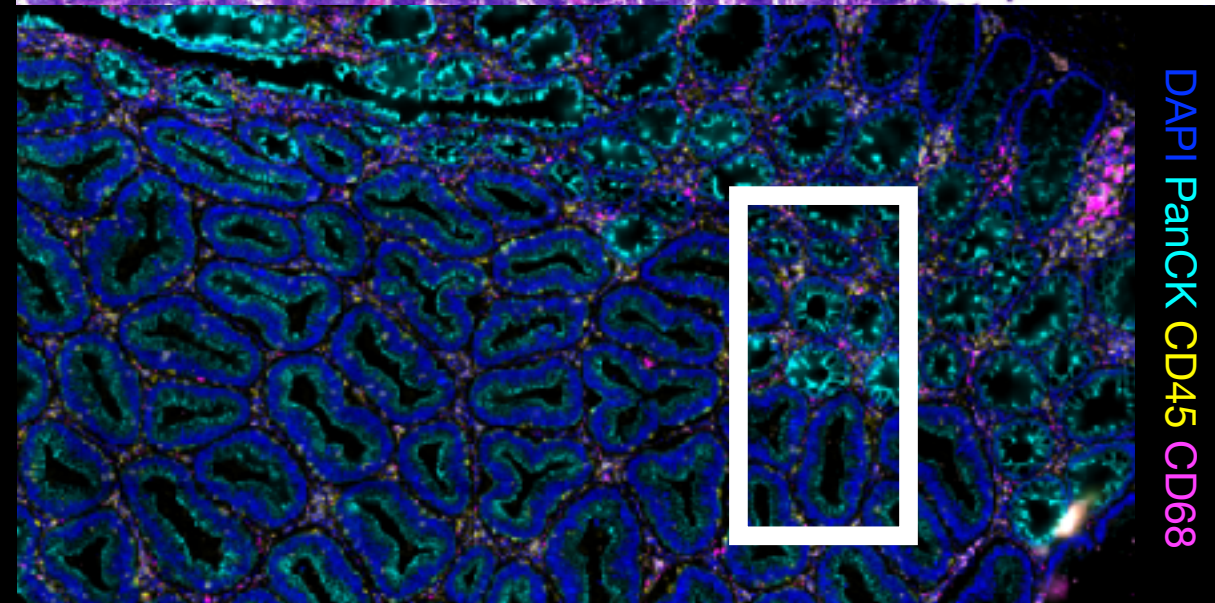
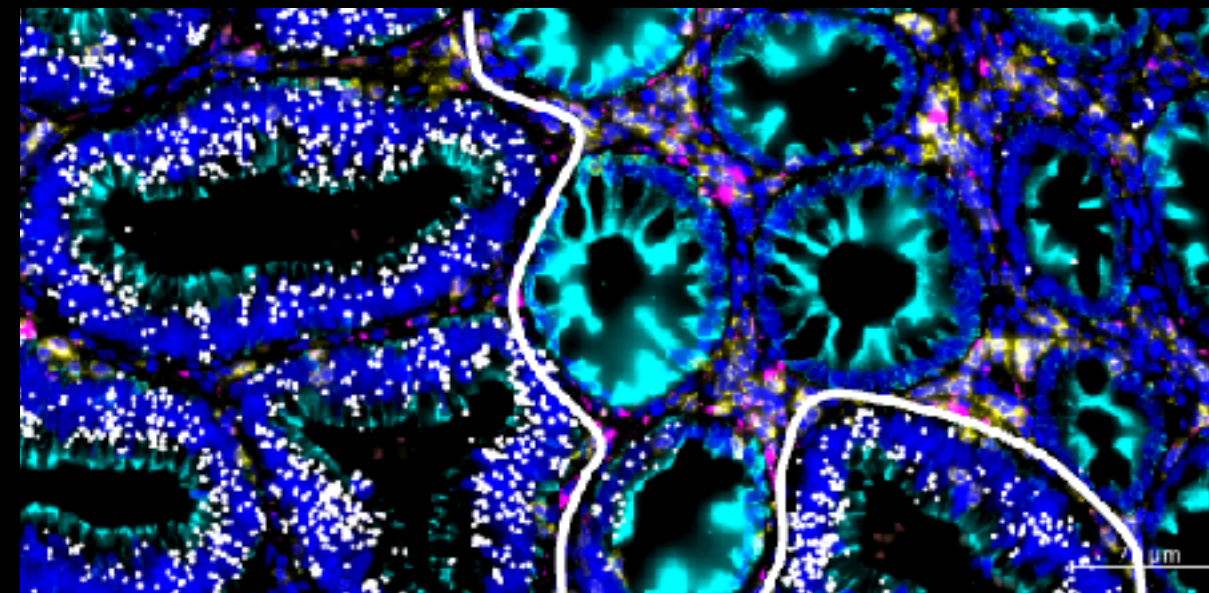
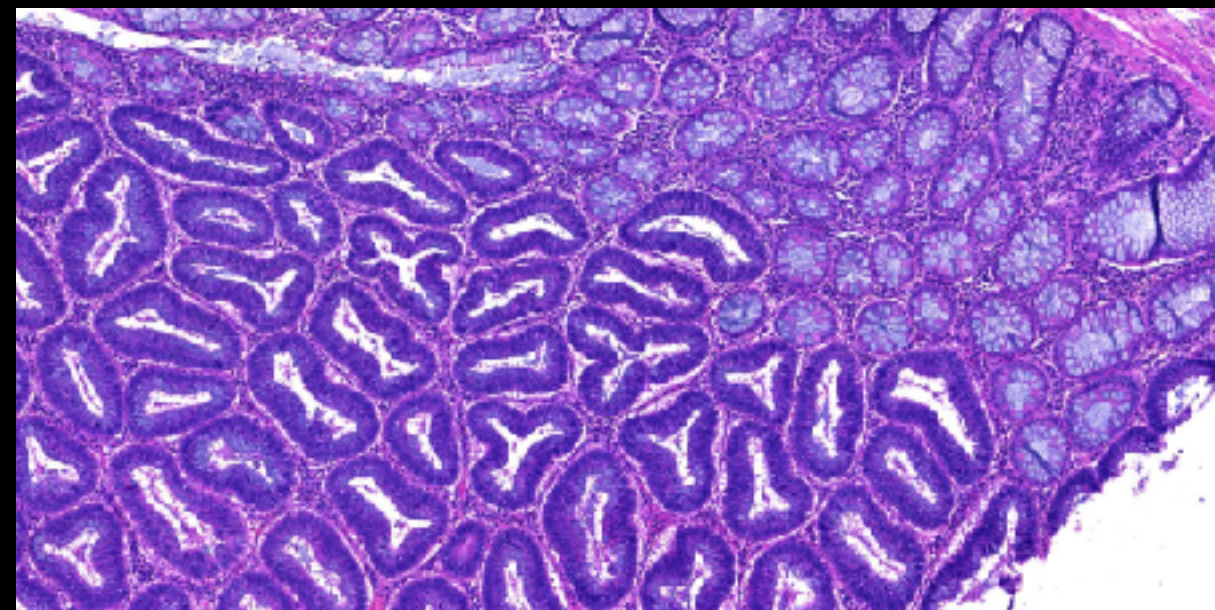
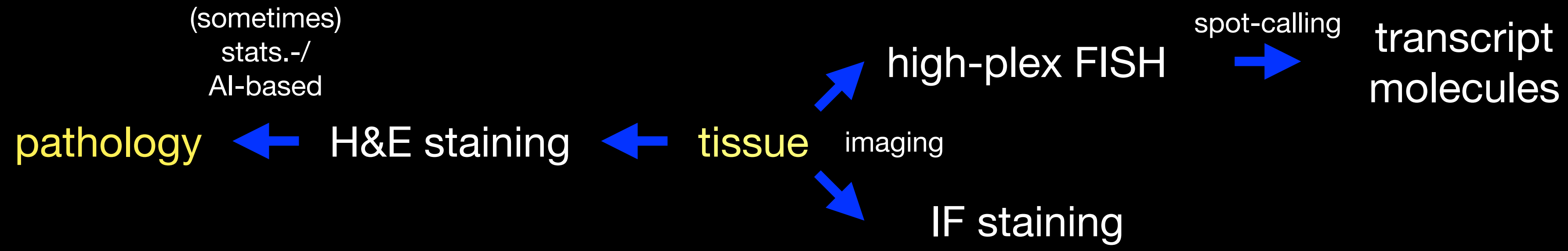
# the path from tissue to pathology ain't easy



DAPI PanCK CD45 CD68



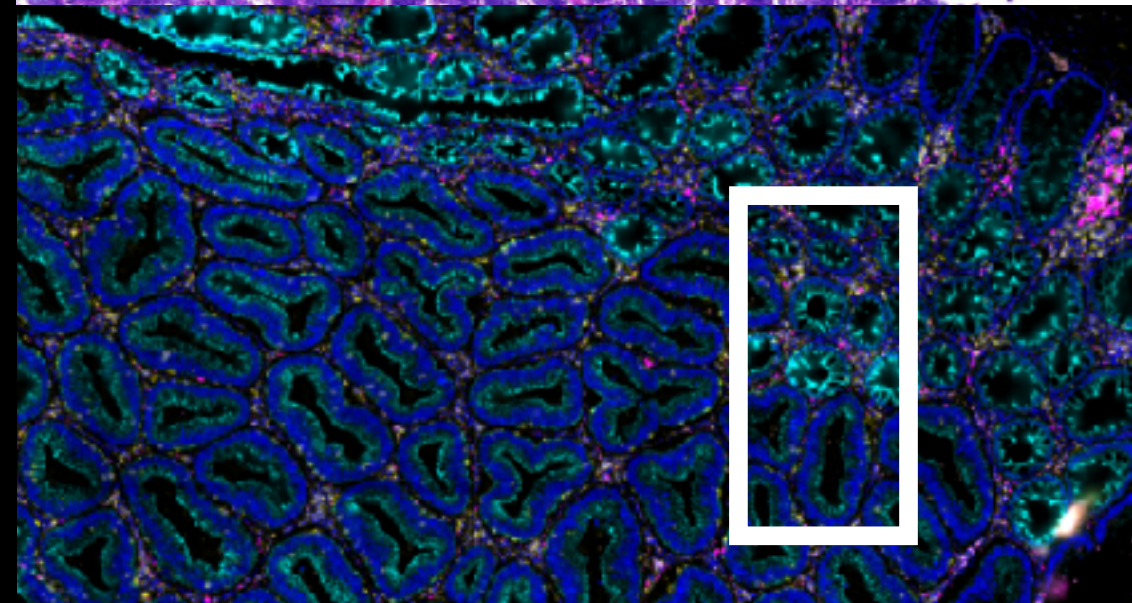
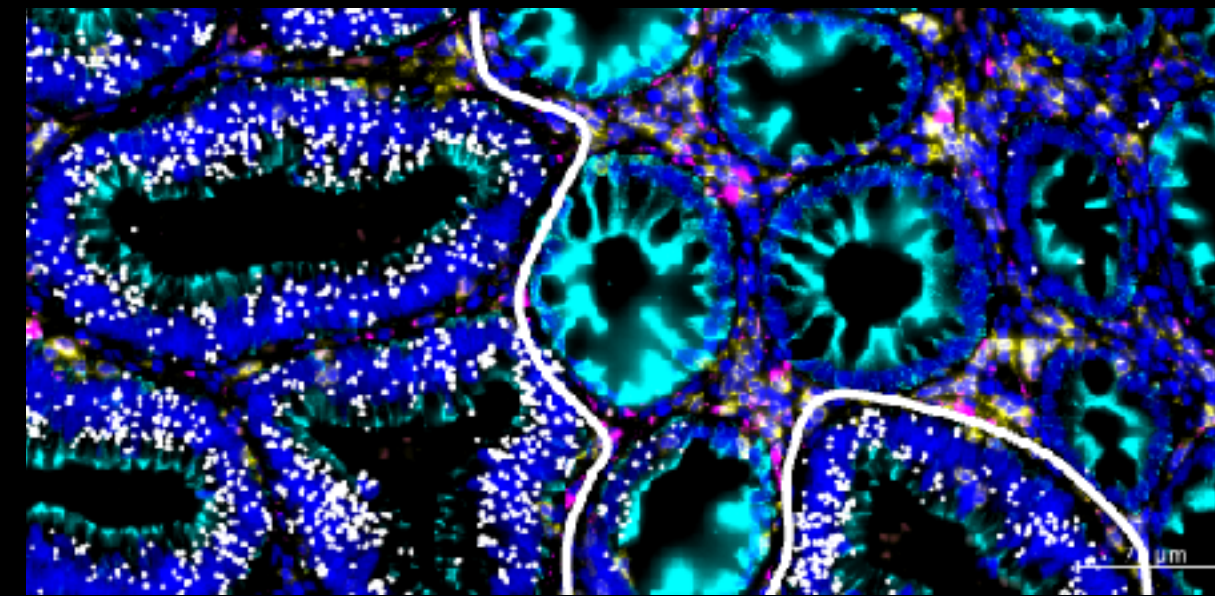
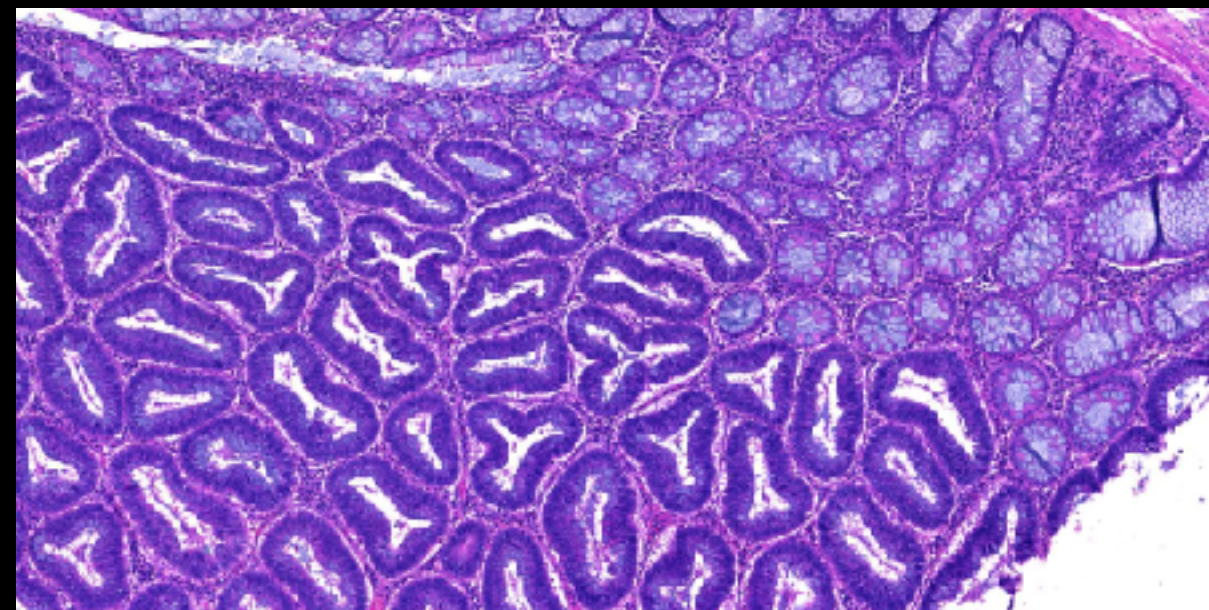
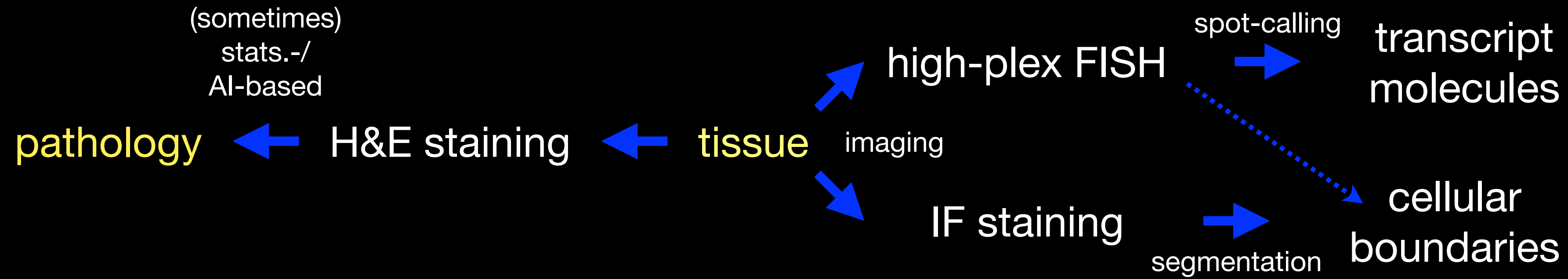
# the path from tissue to pathology ain't easy



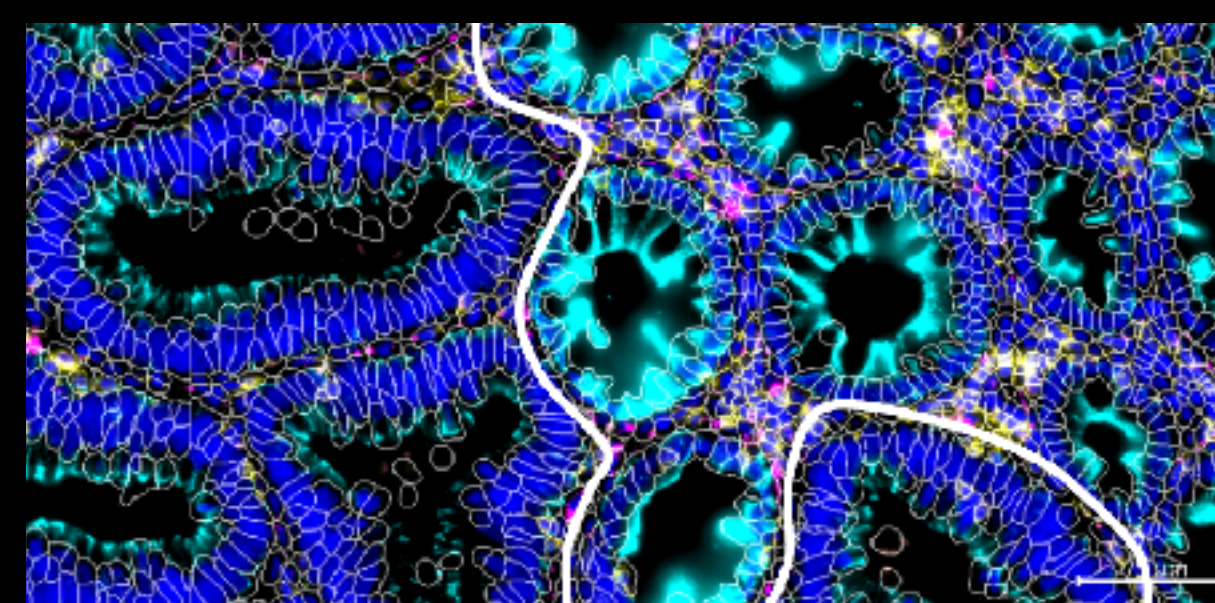
DAPI PanCK CD45 CD68



# the path from tissue to pathology ain't easy

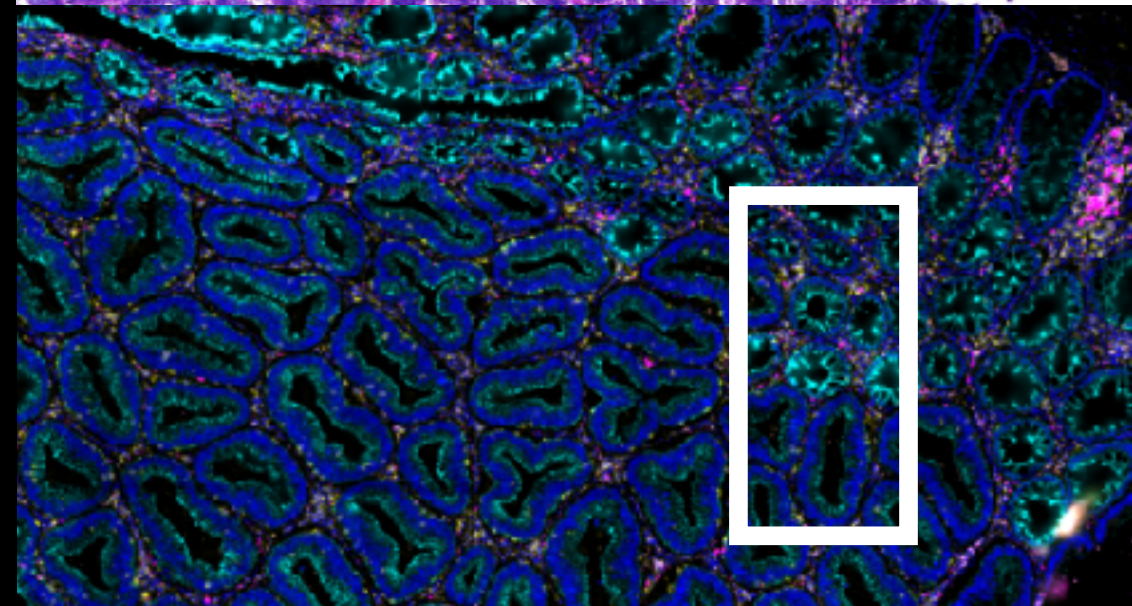
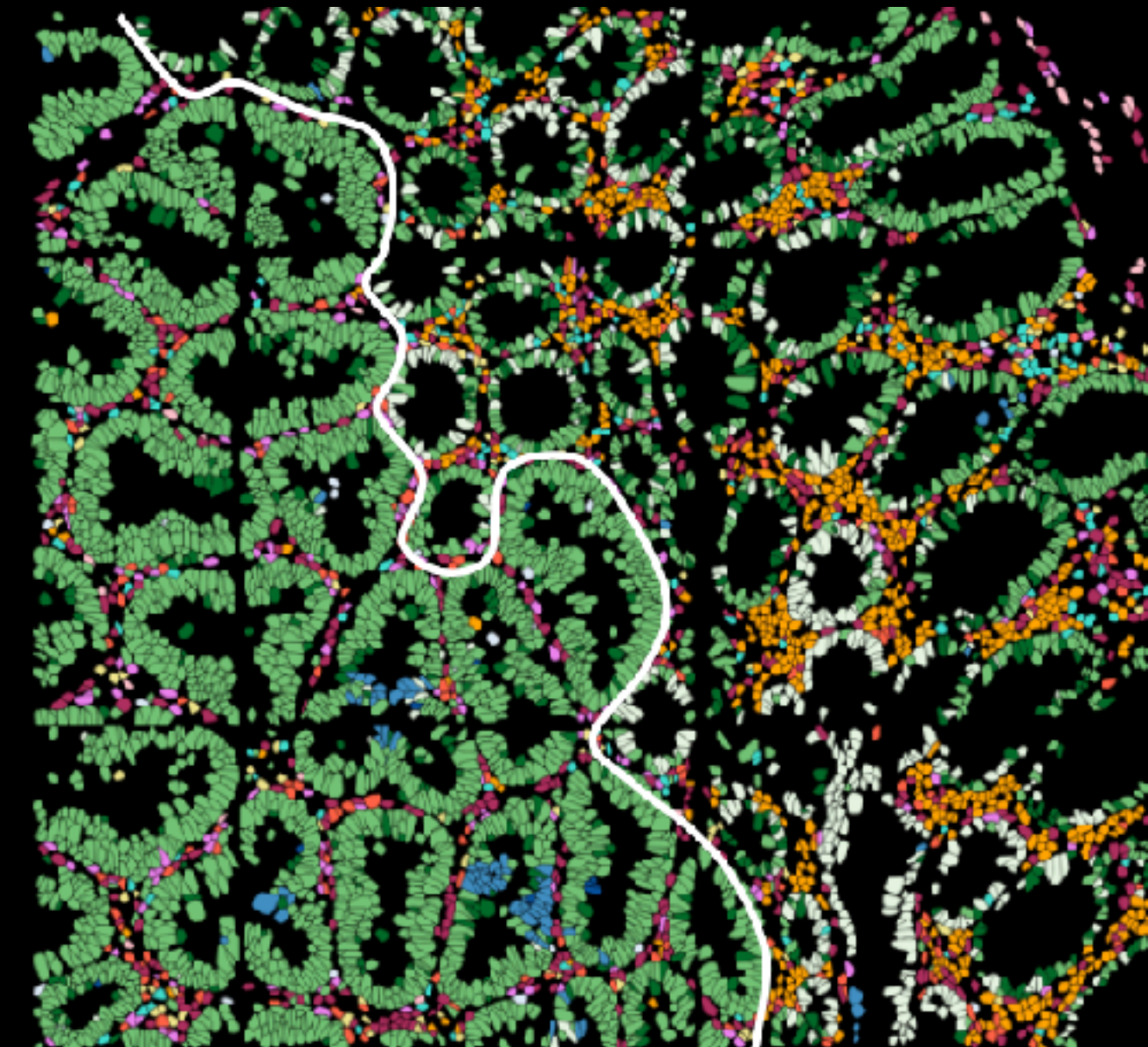
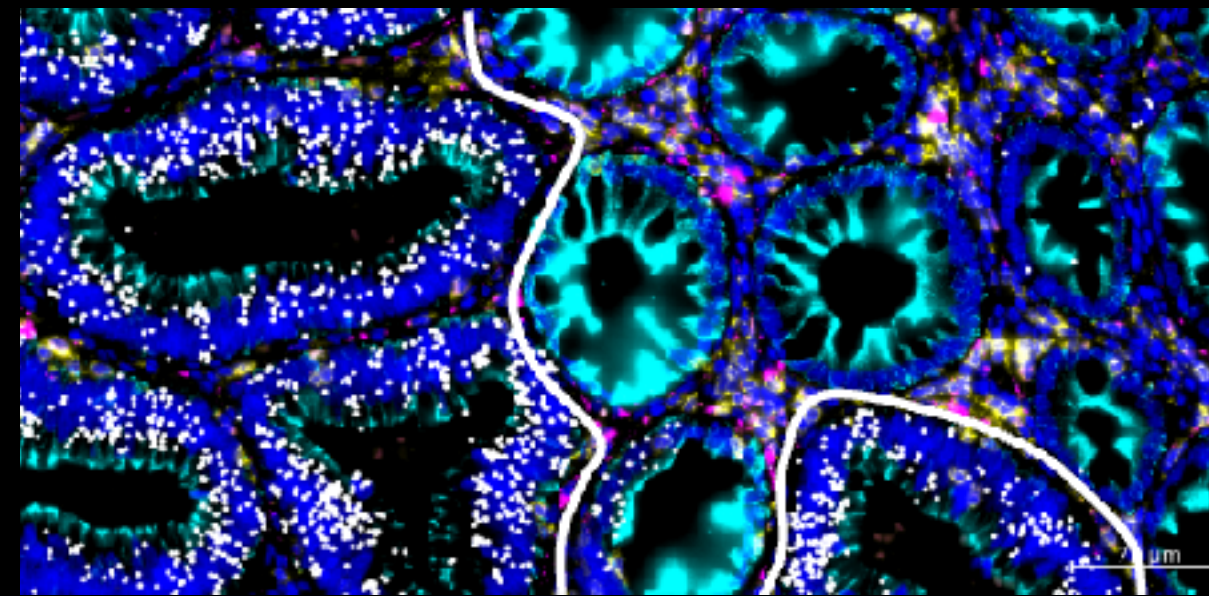
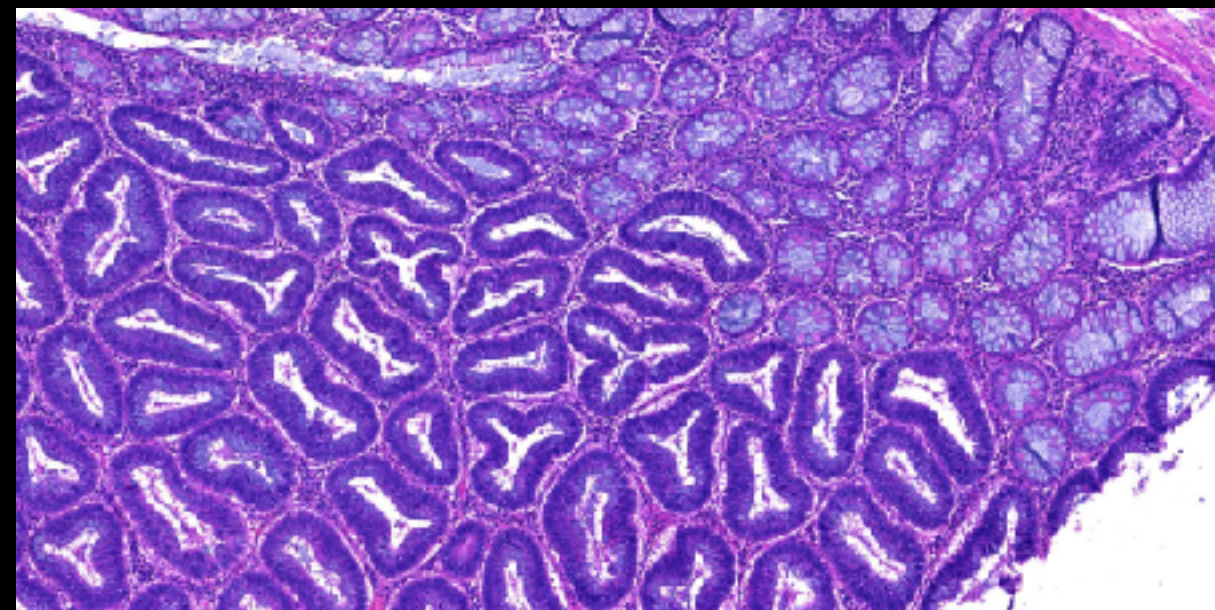
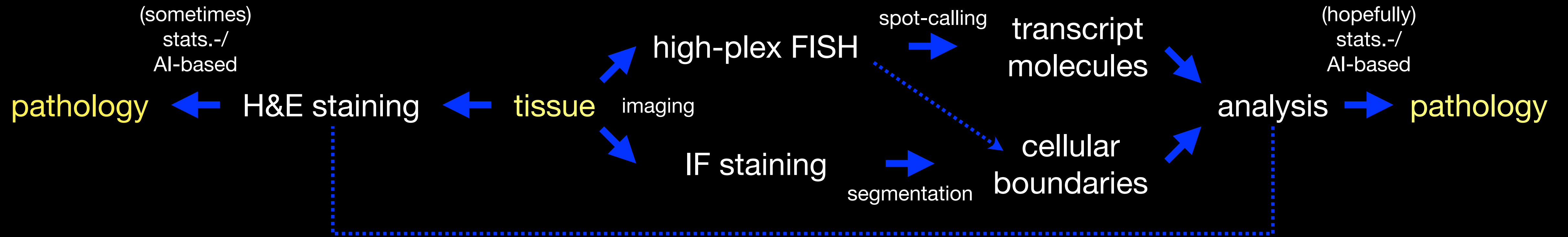


DAPI PanCK CD45 CD68

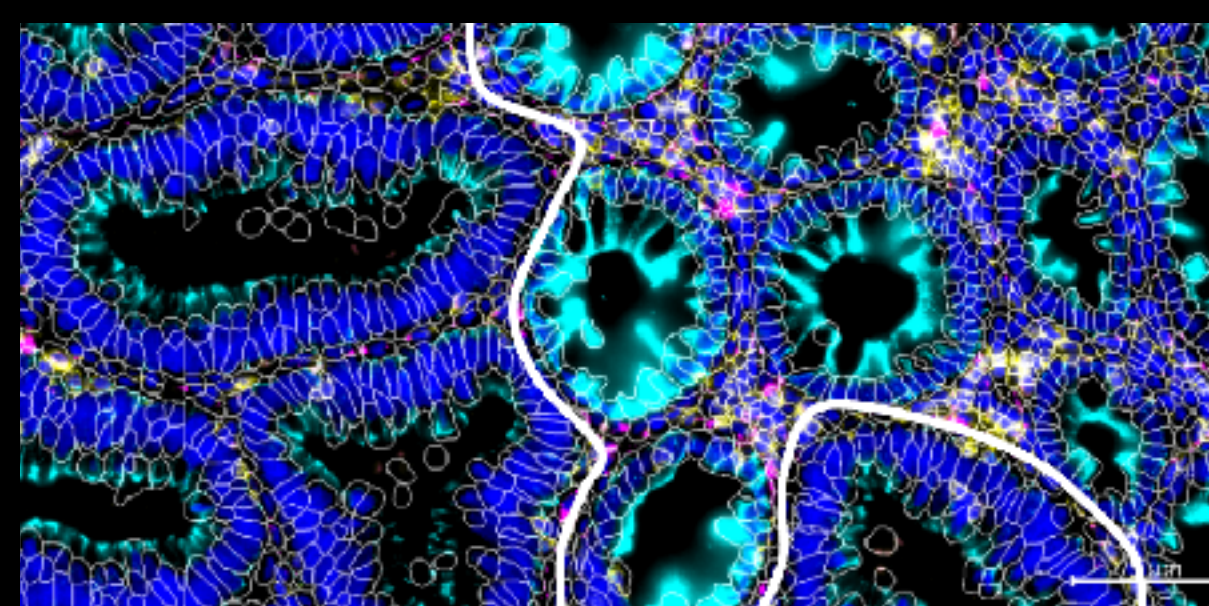




# the path from tissue to pathology ain't easy

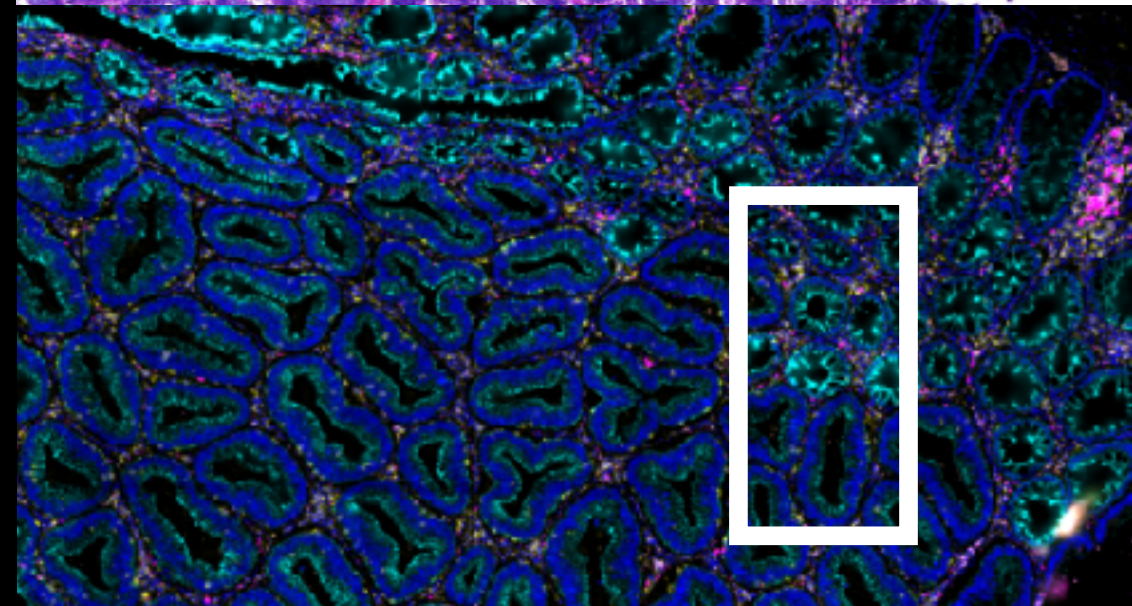
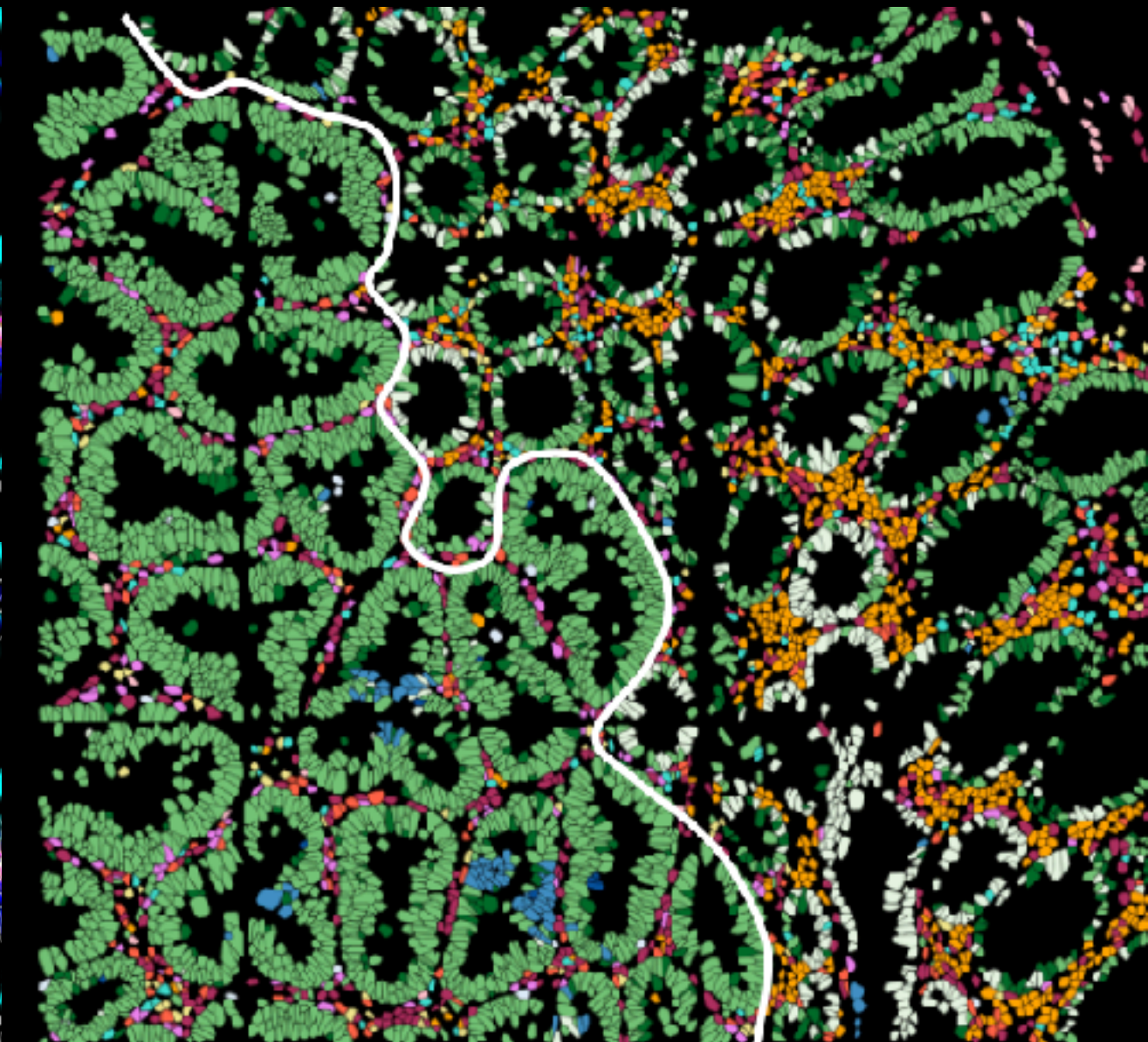
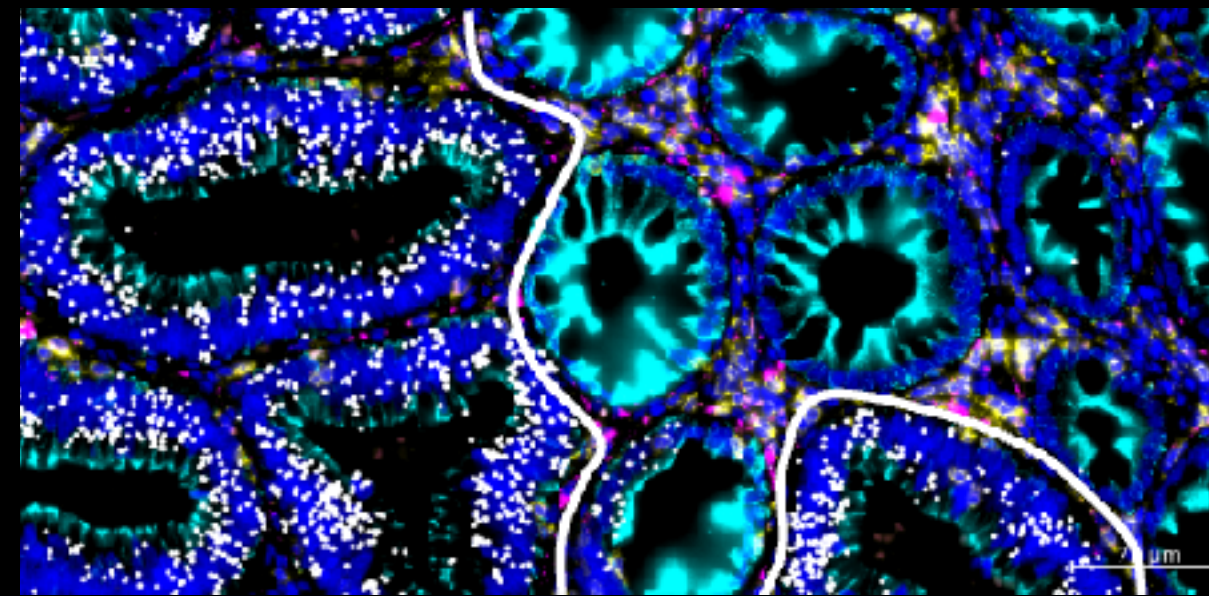
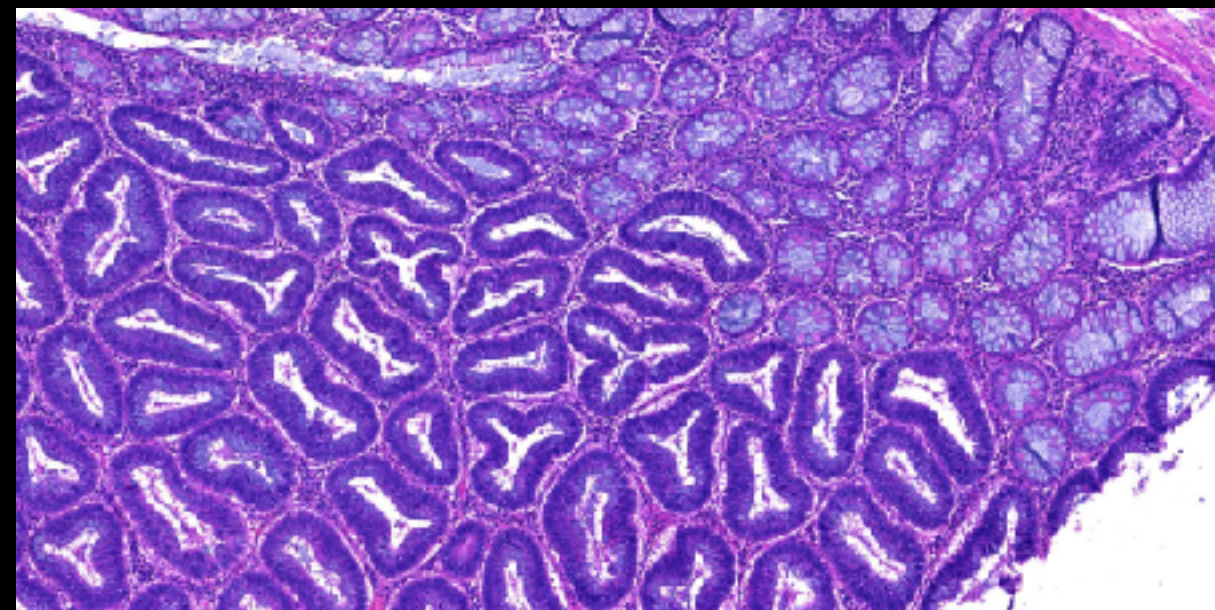
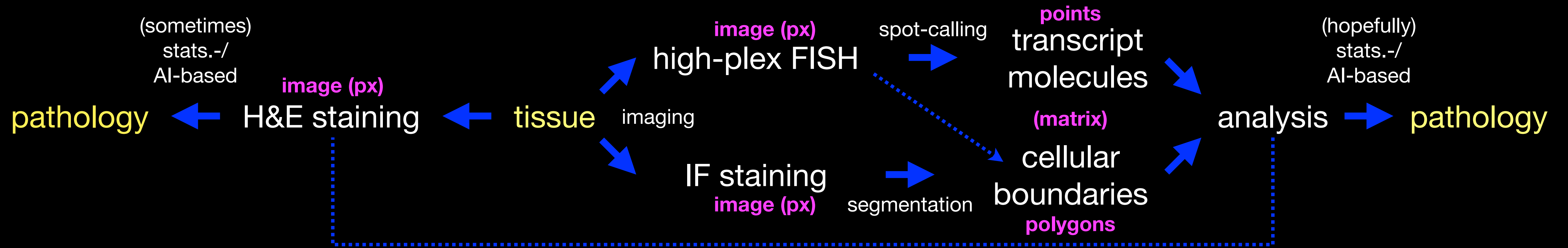


DAPI PanCK CD45 CD68

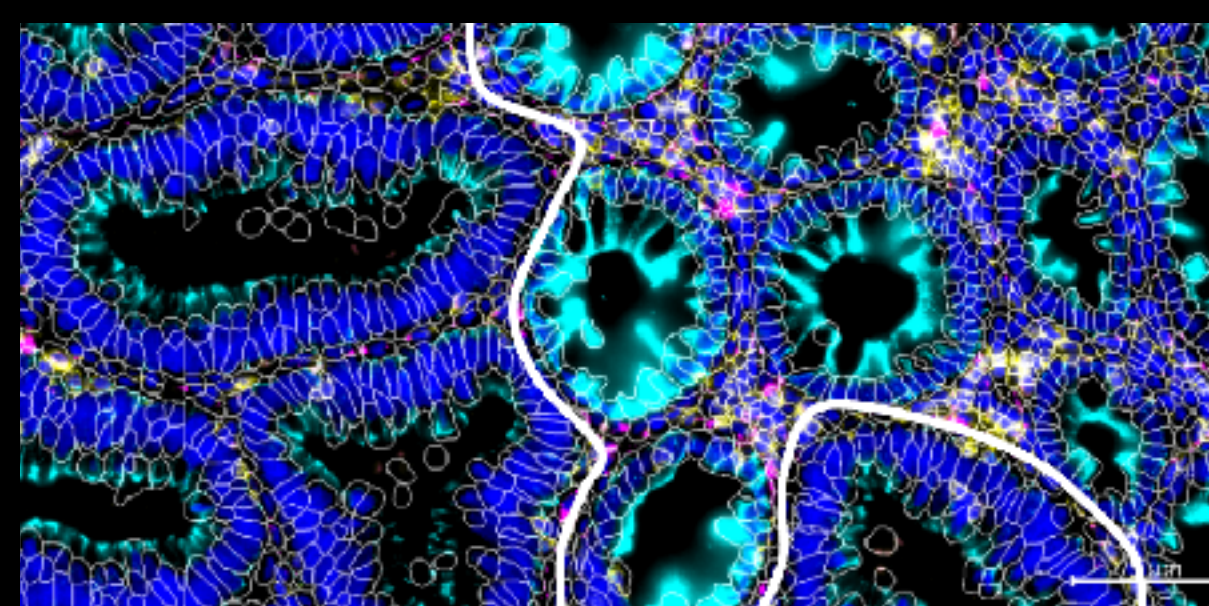




# the path from tissue to pathology ain't easy



DAPI PanCK CD45 CD68





# quality control for img-ST is ways from (sc)RNA-seq

## high-throughput RNA sequencing

fragmentation, reverse transcription, mapping



# of raw counts per transcript varies with transcript length, GC content, sequencing depth



normalization strategies aim to minimize these effects & *“there’s awareness that misinterpretation of results where biological & technical effects are correlated”*



# quality control for img-ST is ways from (sc)RNA-seq

## high-throughput RNA sequencing

fragmentation, reverse transcription, mapping



# of raw counts per transcript varies with transcript length, GC content, sequencing depth



normalization strategies aim to minimize these effects & *“there’s awareness that misinterpretation of results where biological & technical effects are correlated”*

## imaging-based spatial transcriptomics

tissue preparation, chemistry, imaging



tissue damage/detachment, image/transcript loss, varying detection across space & experiments



# quality control for img-ST is ways from (sc)RNA-seq

## high-throughput RNA sequencing

fragmentation, reverse transcription, mapping



# of raw counts per transcript varies with transcript length, GC content, sequencing depth



normalization strategies aim to minimize these effects & *“there’s awareness that misinterpretation of results where biological & technical effects are correlated”*

## imaging-based spatial transcriptomics

tissue preparation, chemistry, imaging



tissue damage/detachment, image/transcript loss, varying detection across space & experiments



*“sources of these [...] are known [but] it’s often unclear how often errors occur, how to best detect & describe [them] & how [they] impact downstream analyses [...]”*



# quality control for img-ST is ways from (sc)RNA-seq

*DISCLAIMER: We don't really know what's best (yet) — I'll try to summarize some recent ideas, and personal pains.*

## high-throughput RNA sequencing

fragmentation, reverse transcription, mapping



# of raw counts per transcript varies with transcript length, GC content, sequencing depth



normalization strategies aim to minimize these effects & *“there's awareness that misinterpretation of results where biological & technical effects are correlated”*

## imaging-based spatial transcriptomics

tissue preparation, chemistry, imaging



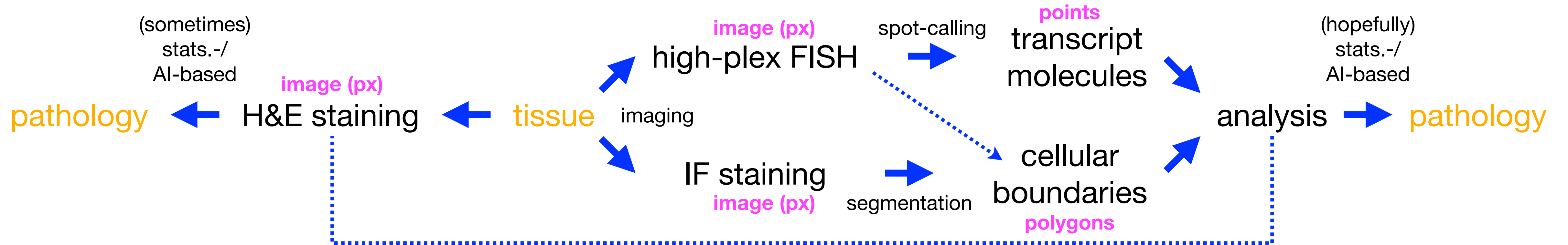
tissue damage/detachment, image/transcript loss, varying detection across space & experiments



*“sources of these [...] are known [but] it's often unclear how often errors occur, how to best detect & describe [them] & how [they] impact downstream analyses [...]”*

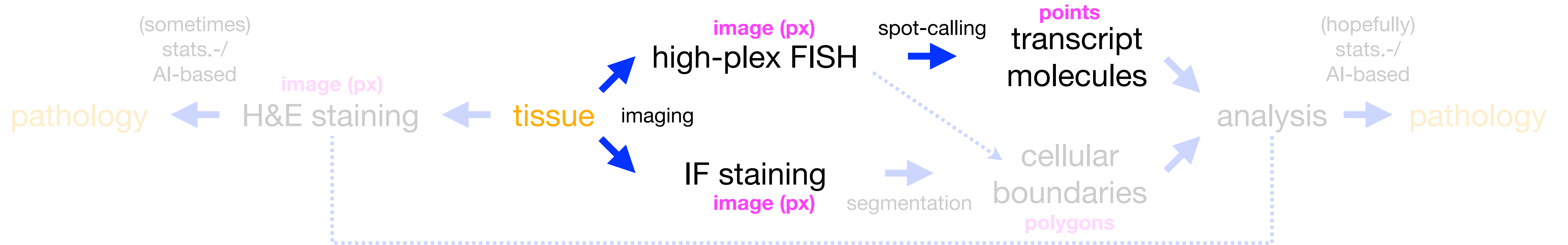


# MerQuaCo proposes a **pixel classifier** & quality metrics





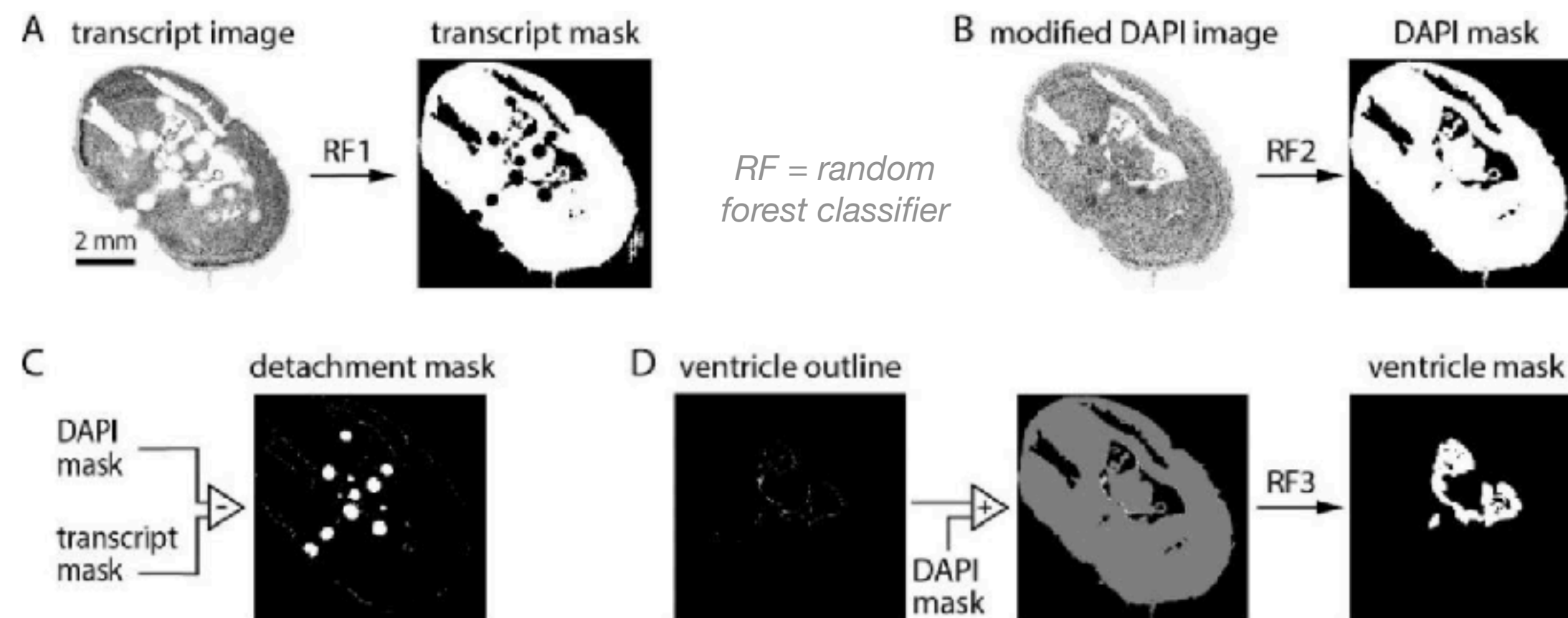
# MerQuaCo proposes a **pixel classifier** & quality metrics





# MerQuaCo proposes a **pixel classifier** & quality metrics

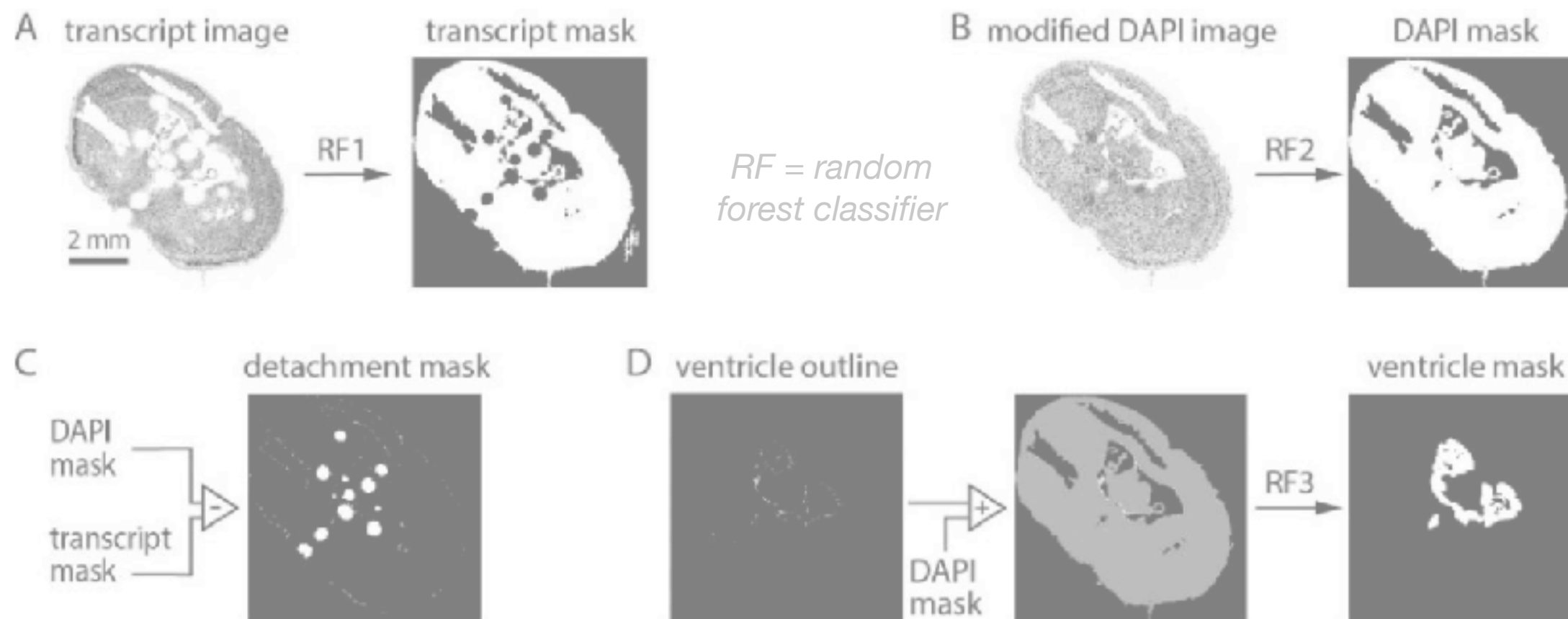
- input: transcript locations + DAPI staining
- series of **binary masks** (& combinations thereof), **trained on few manually annotated sections**



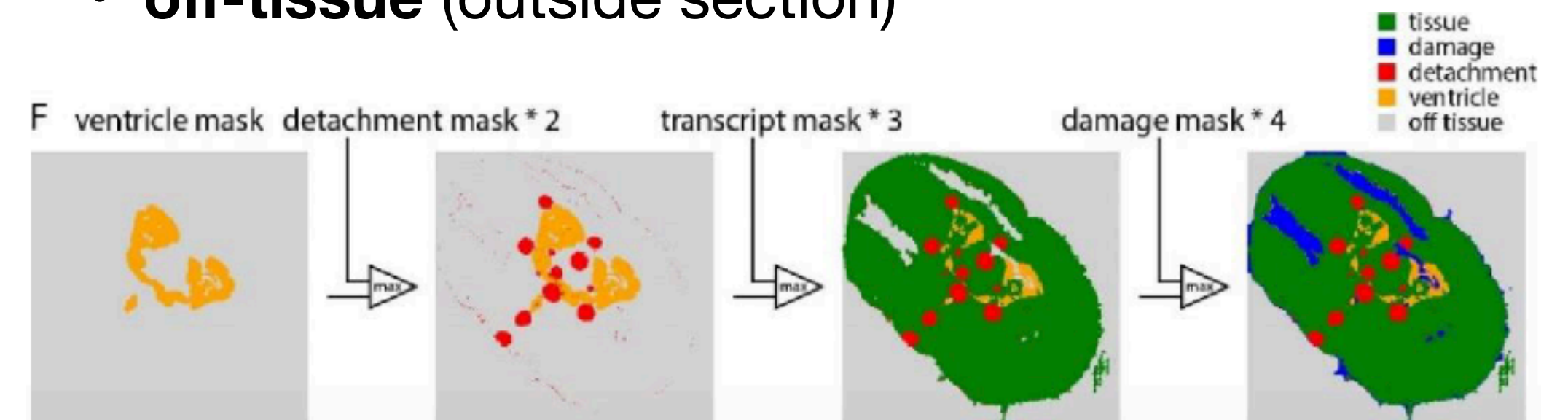


# MerQuaCo proposes a **pixel classifier** & quality metrics

- input: transcript locations + DAPI staining
- series of **binary masks** (& combinations thereof), **trained on few manually annotated sections**



- each location is classified into one of **five categories**
  - **tissue** within image volumen
  - **detachment** (tissue not imaged)
  - **ventricle** (no tissue but no loss)
  - **damage** (no tissue due to loss)
  - **off-tissue** (outside section)





# MerQuaCo proposes a pixel classifier & quality metrics

## **perfusion rate**

- log files can reveal inconsistencies (e.g., blockage) of volume per time during solution exchange

## **data loss**

- iterative comparison of transcript counts between neighboring fields of view (FOVs)

## **detection efficiency**

- across section: periodicity,  
through section: p6/p0 ratio

## **transcript density**

- should vary across section, but little between sections of comparable quality



# MerQuaCo proposes a pixel classifier & quality metrics

## perfusion rate

- log files can reveal inconsistencies (e.g., blockage) of volume per time during solution exchange

## data loss

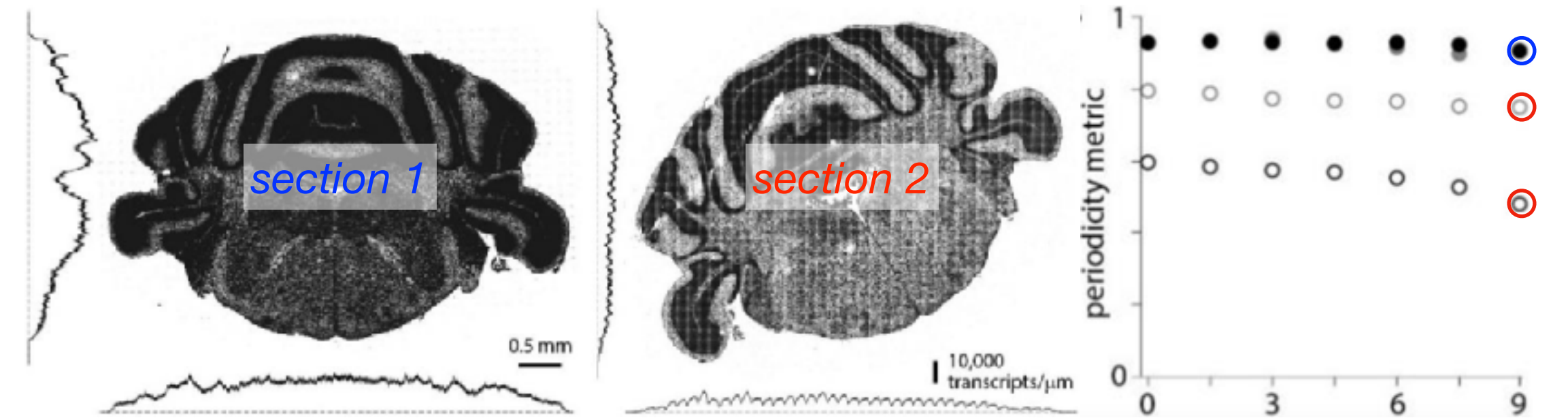
- iterative comparison of transcript counts between neighboring fields of view (FOVs)

## detection efficiency

- across section: periodicity,  
through section: p6/p0 ratio

## transcript density

- should vary across section, but little between sections of comparable quality





# MerQuaCo proposes a pixel classifier & quality metrics

## perfusion rate

- log files can reveal inconsistencies (e.g., blockage) of volume per time during solution exchange

## data loss

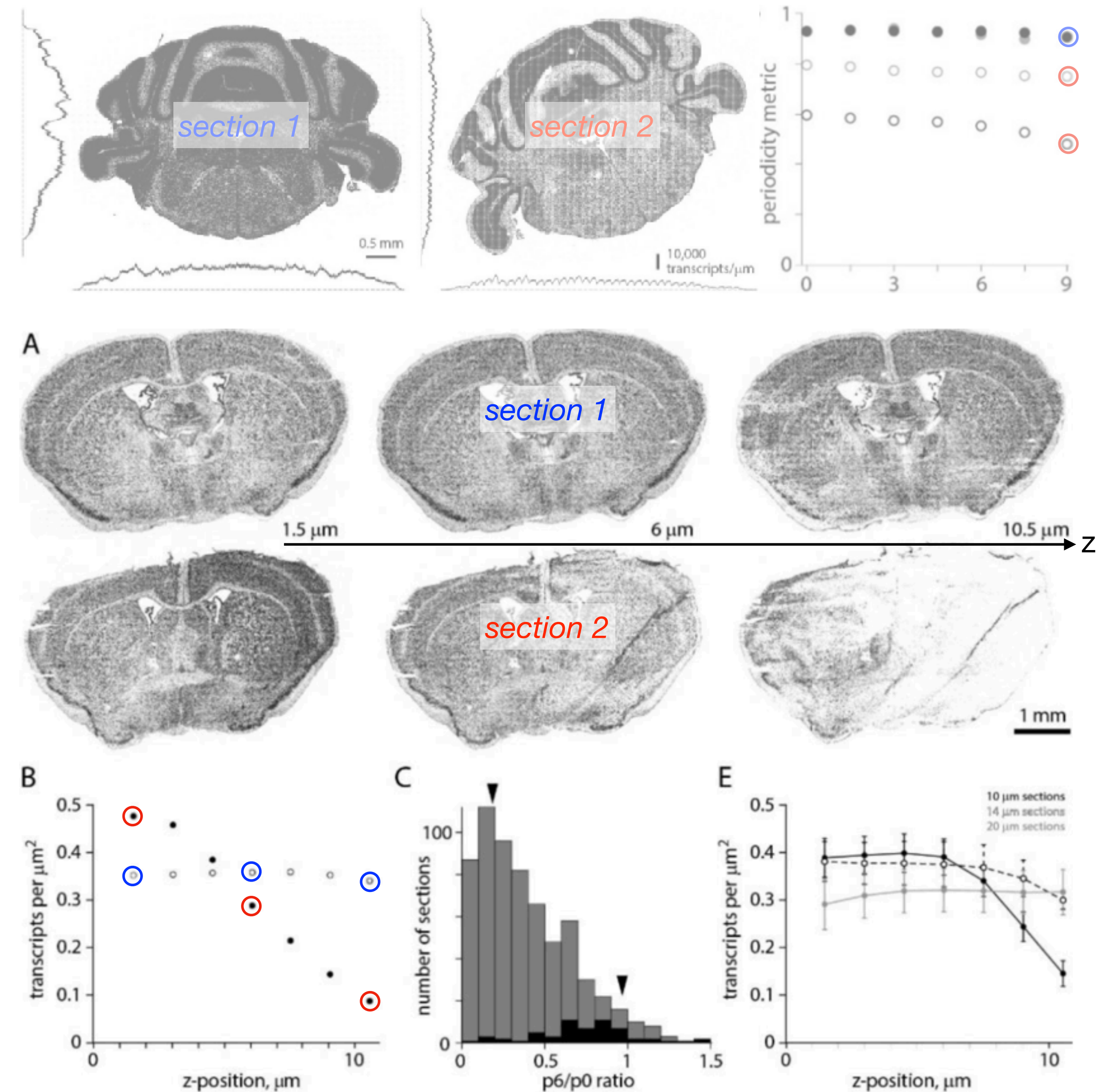
- iterative comparison of transcript counts between neighboring fields of view (FOVs)

## detection efficiency

- across section: periodicity, through section: p6/p0 ratio

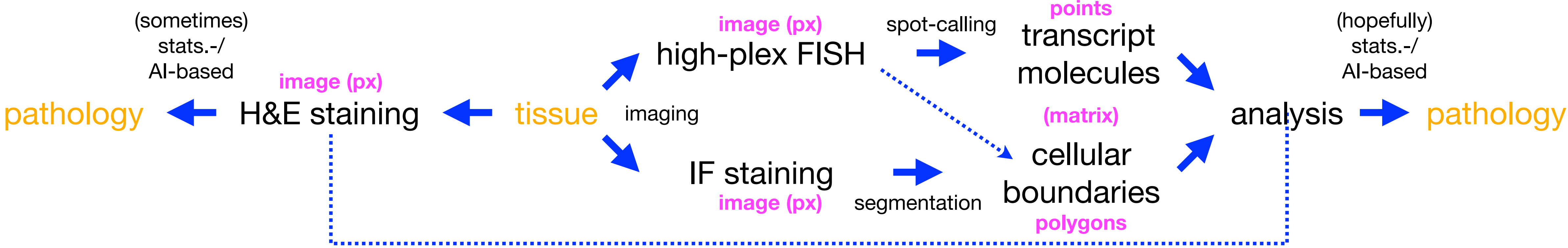
## transcript density

- should vary across section, but little between sections of comparable quality



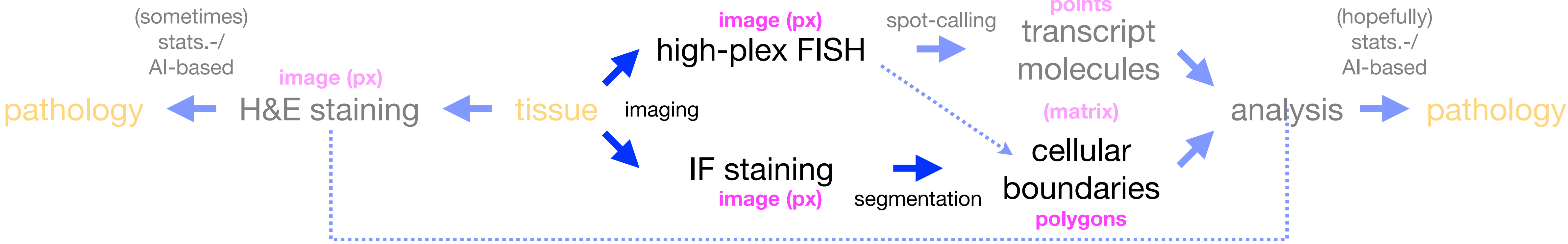


# for (most) analyses, we need to go **from images to counts**



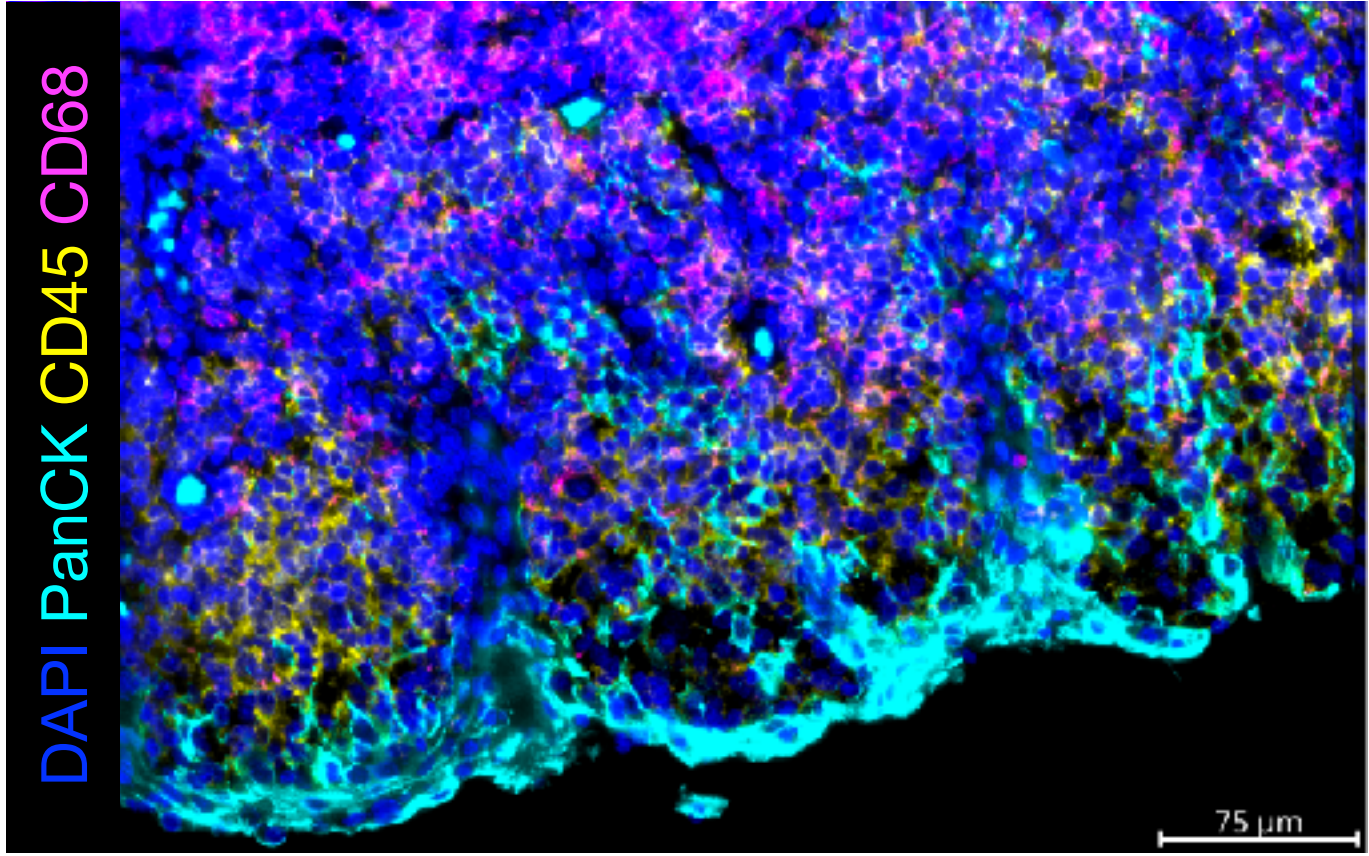
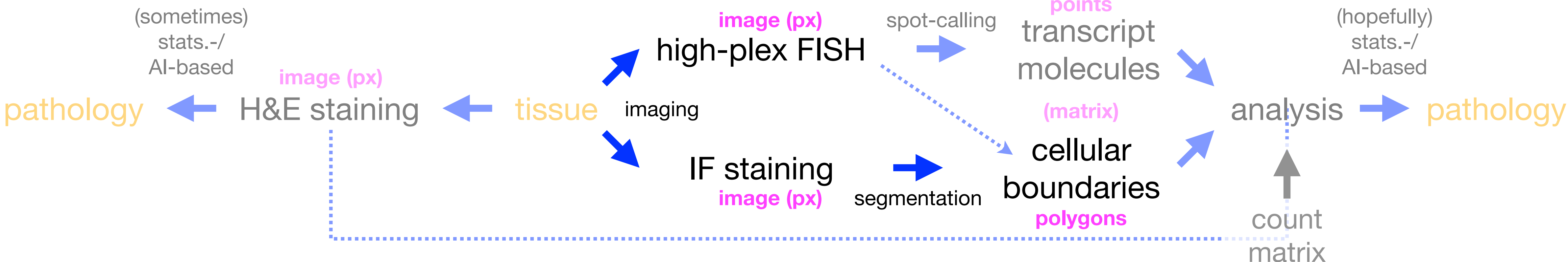


# for (most) analyses, we need to go **from images to counts**

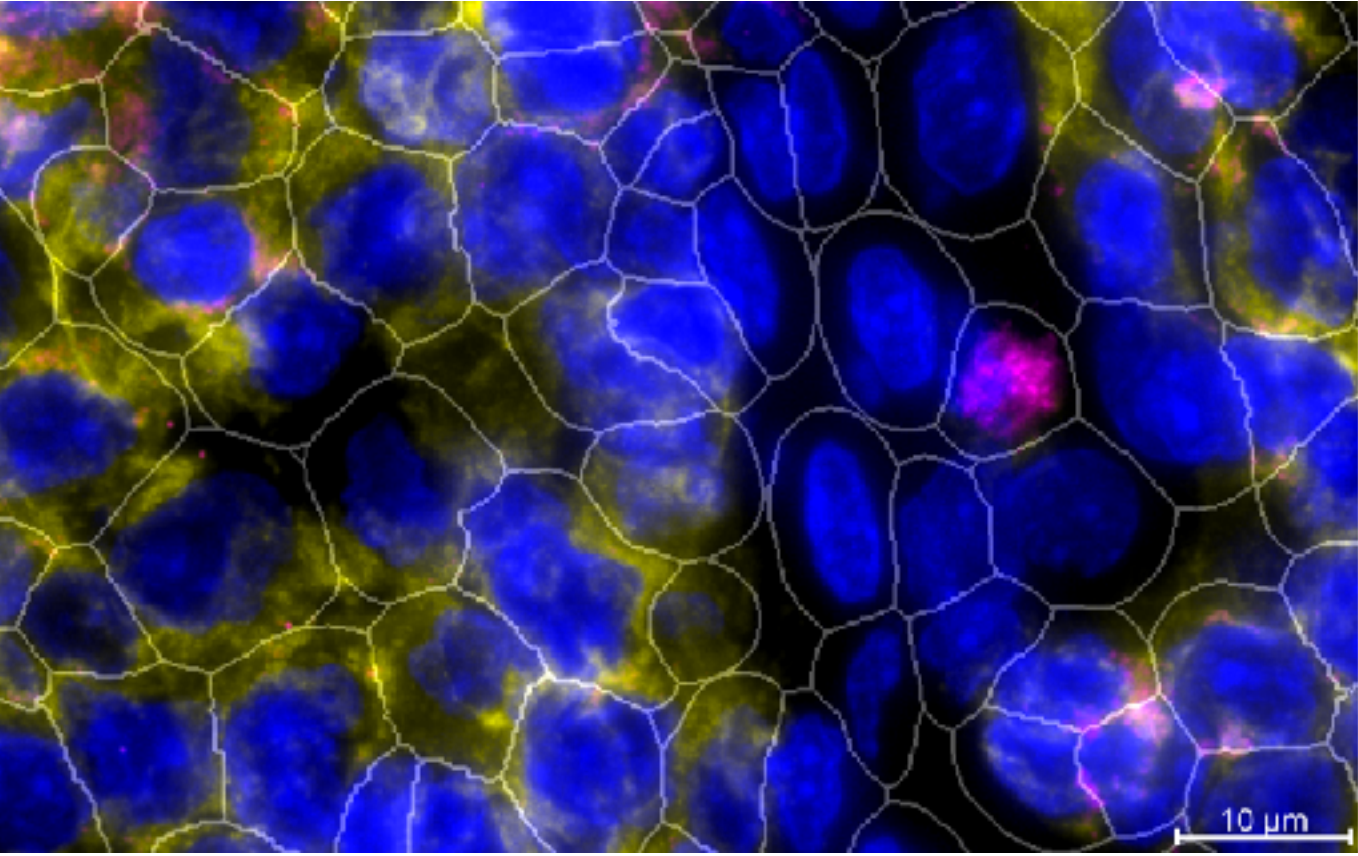




# for (most) analyses, we need to go from images to counts



segmentation →



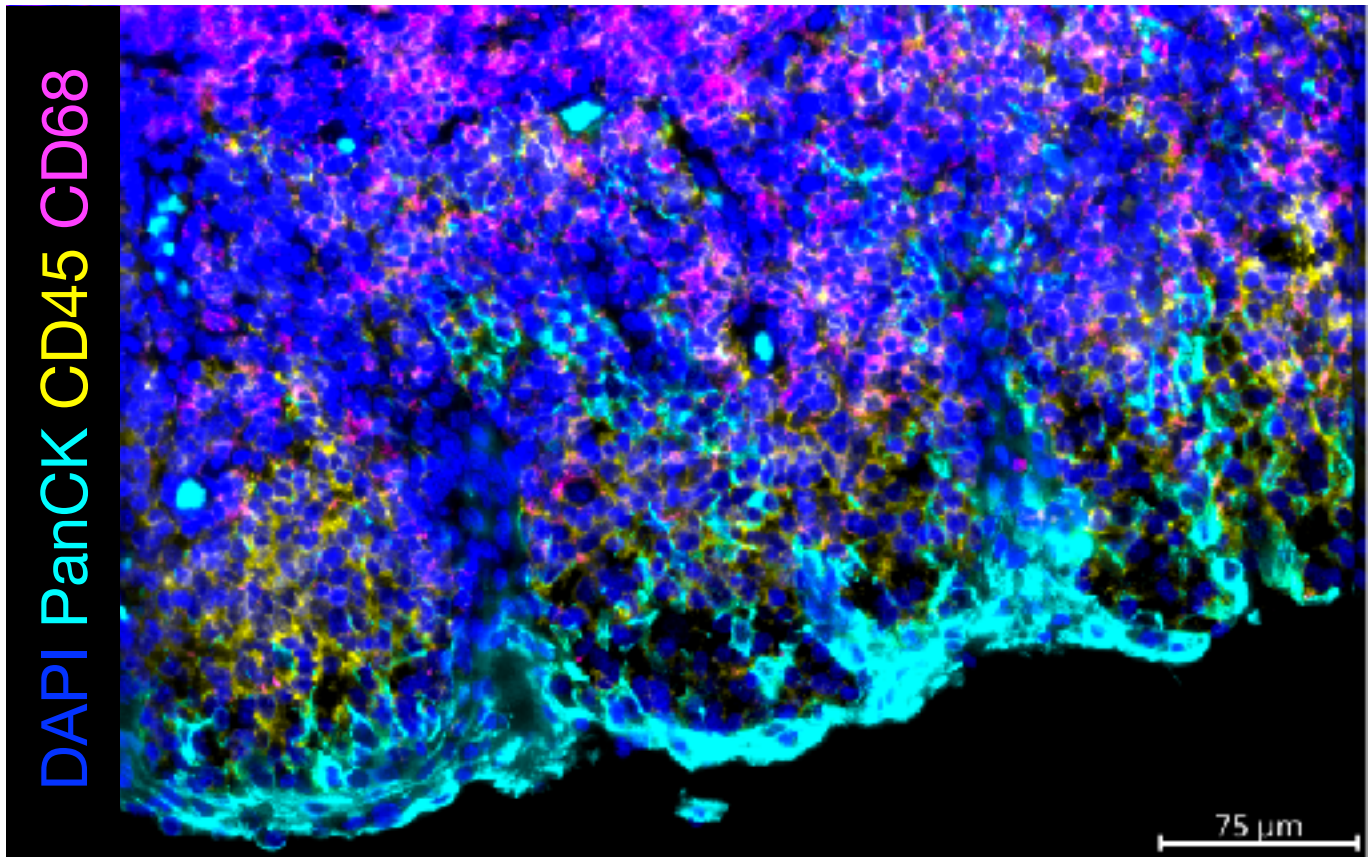
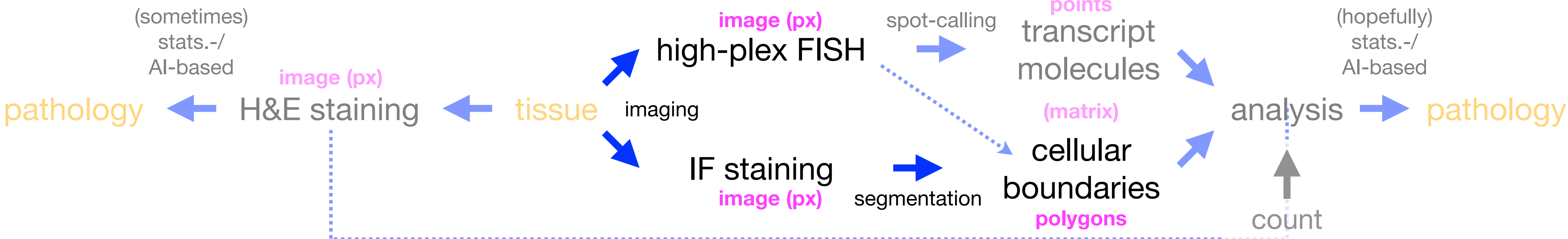
transcripts + boundaries

count matrix

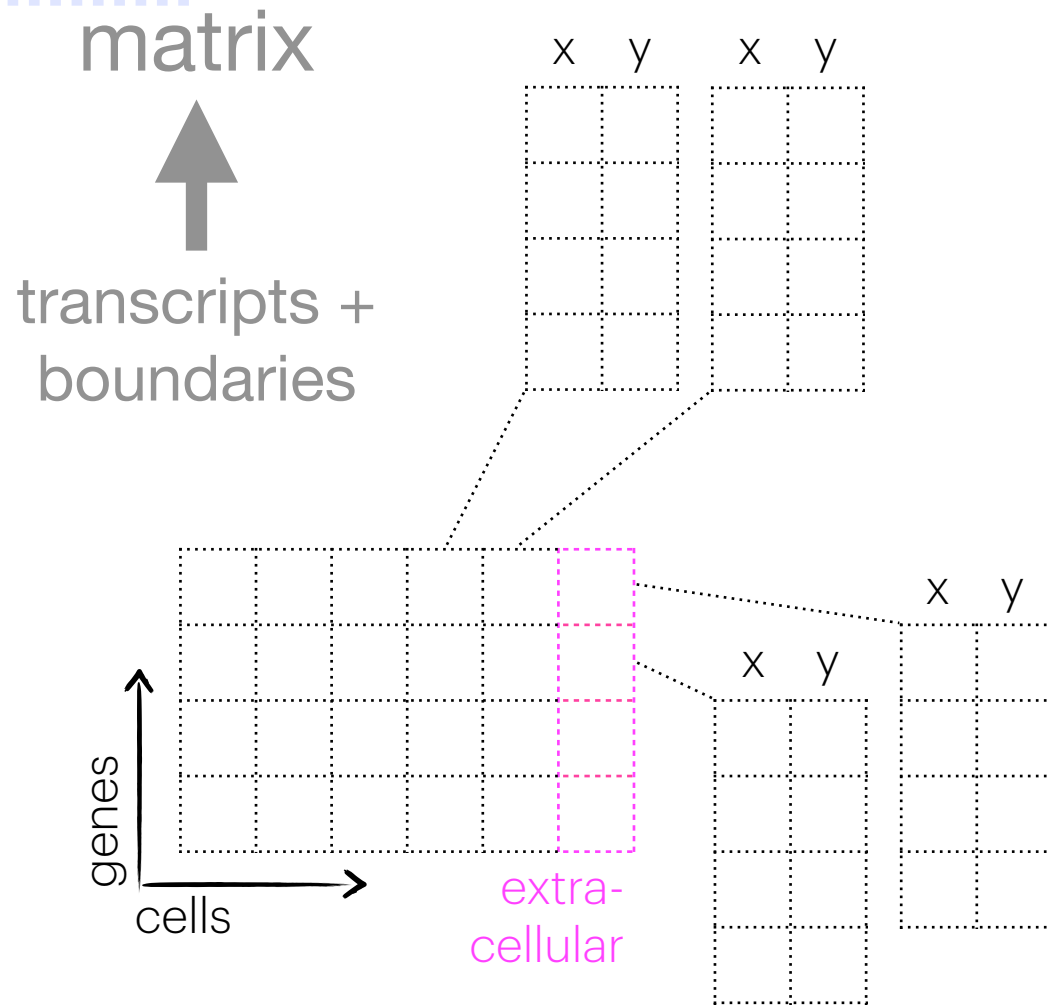
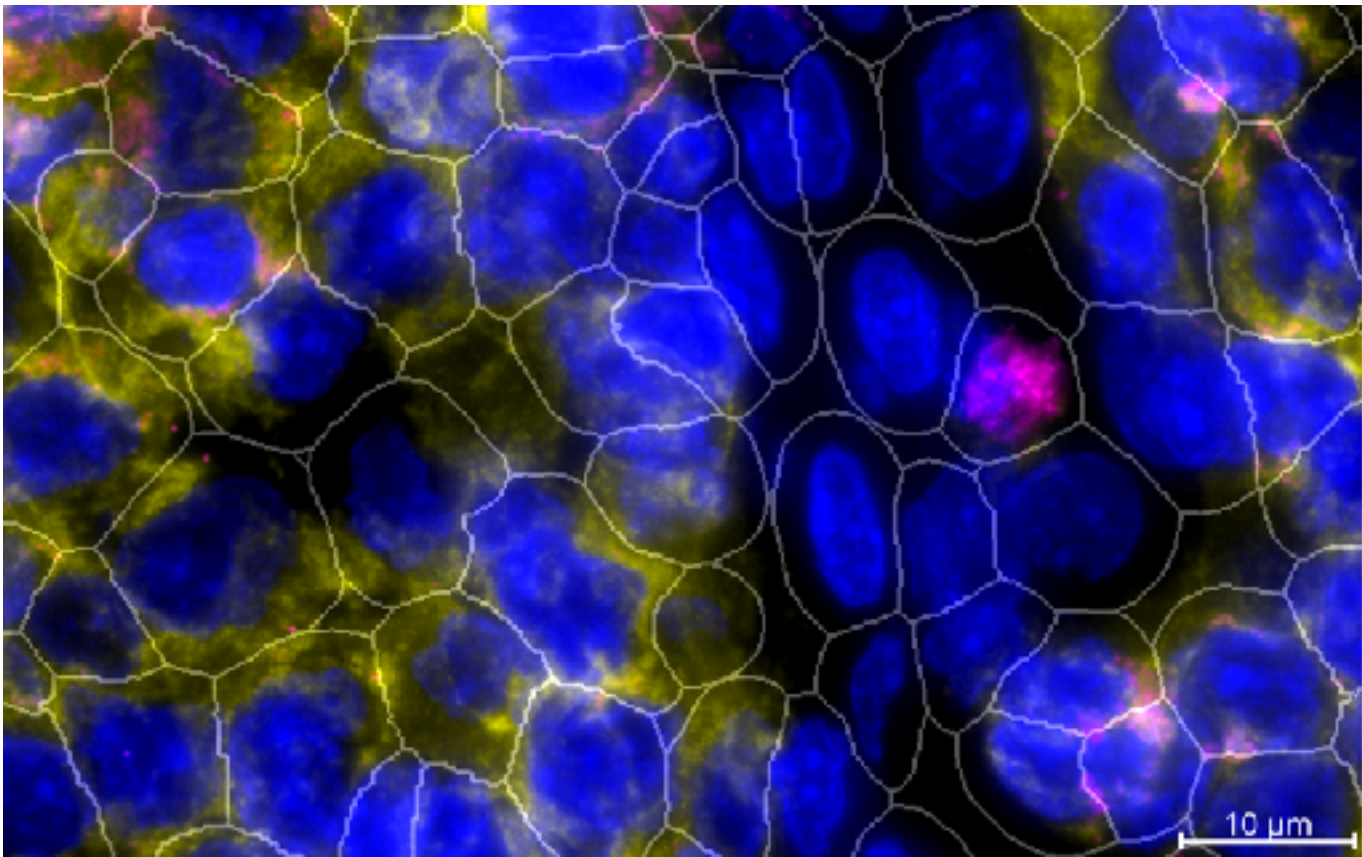
analysis → pathology



# for (most) analyses, we need to go from images to counts

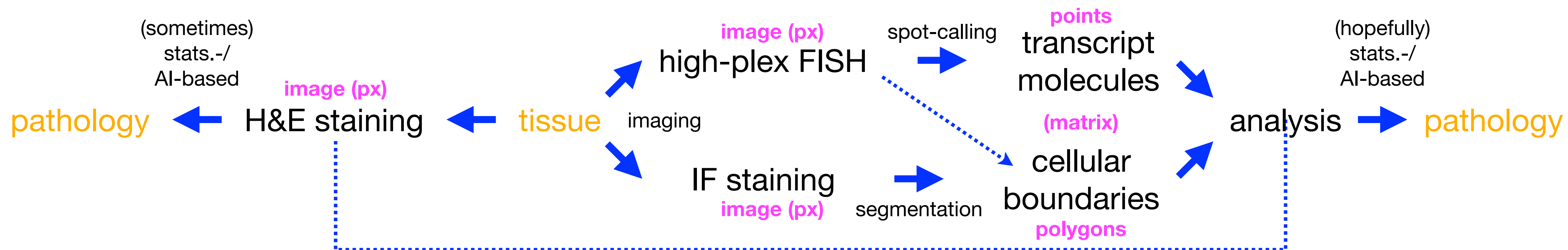


segmentation





# beware standard QC metrics — counts relate area AND biology





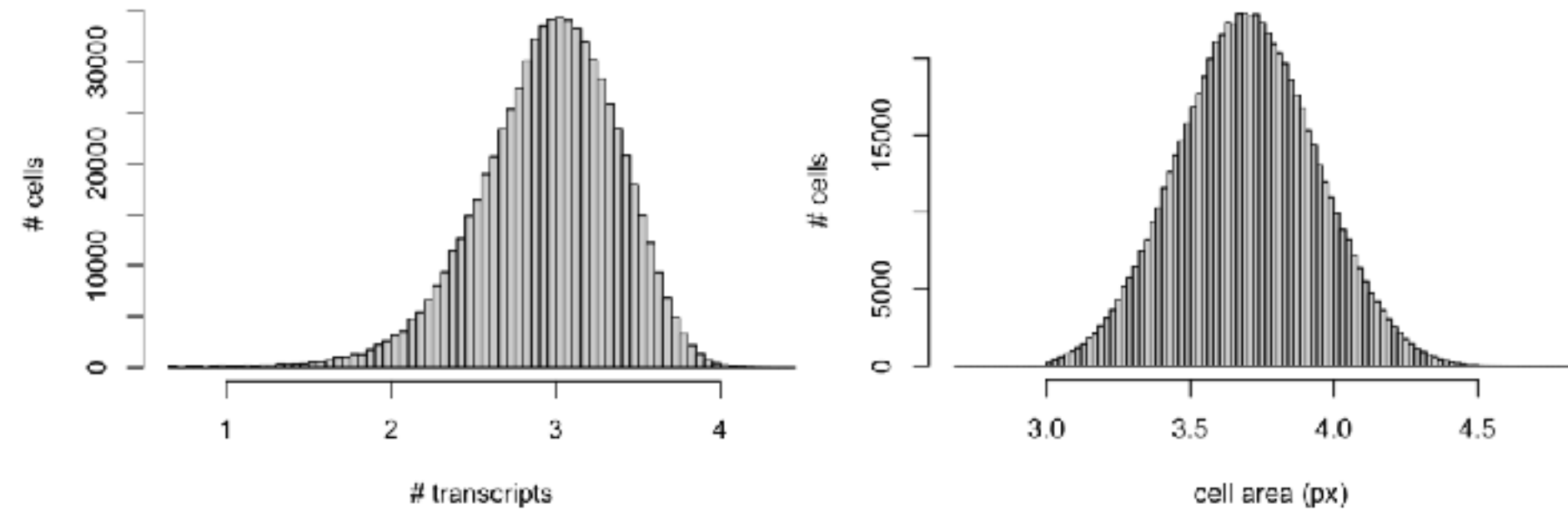
# beware standard QC metrics — counts relate area AND biology





# beware standard QC metrics — counts relate area AND biology

all axes  $\log_{10}$ -transformed (...what we'd usually filter on)

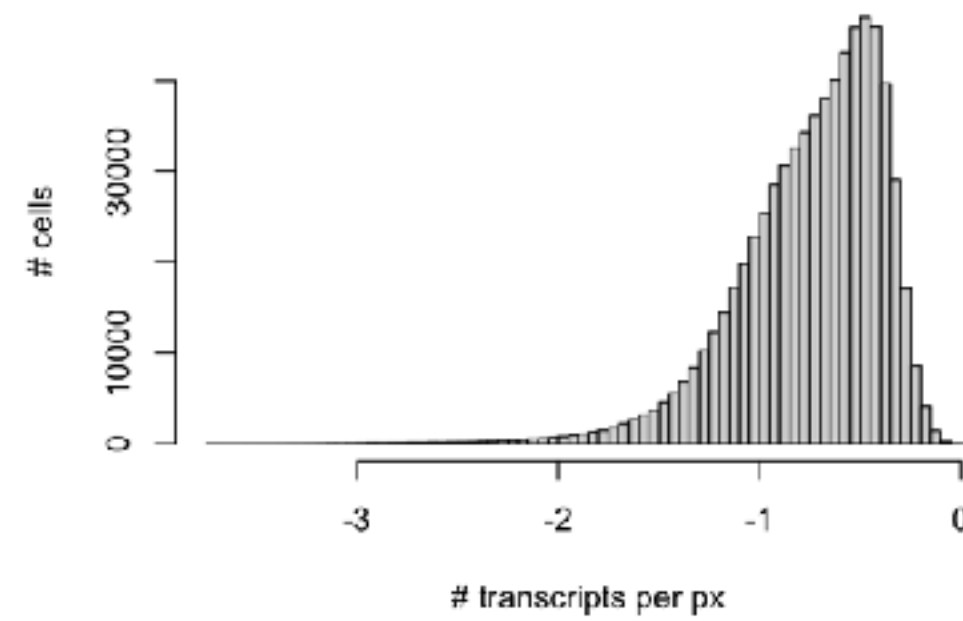
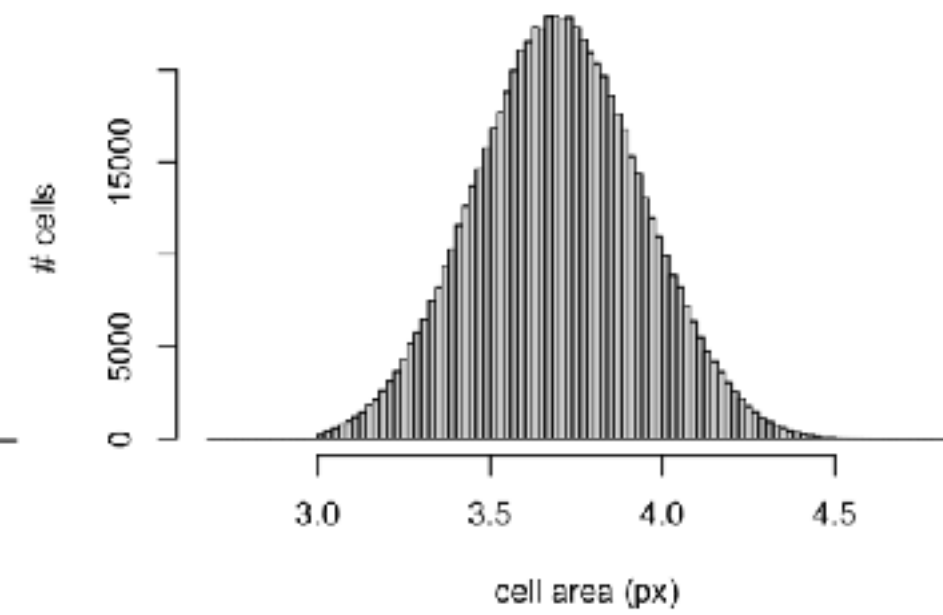
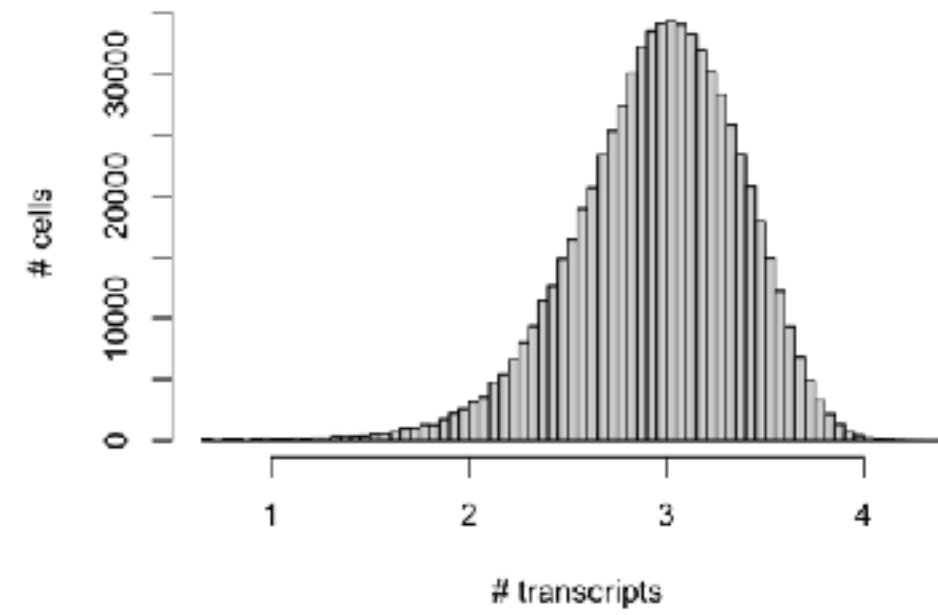


- filtering based on standard scRNA-seq QC metrics may bias against smaller & transcriptionally less complex cells

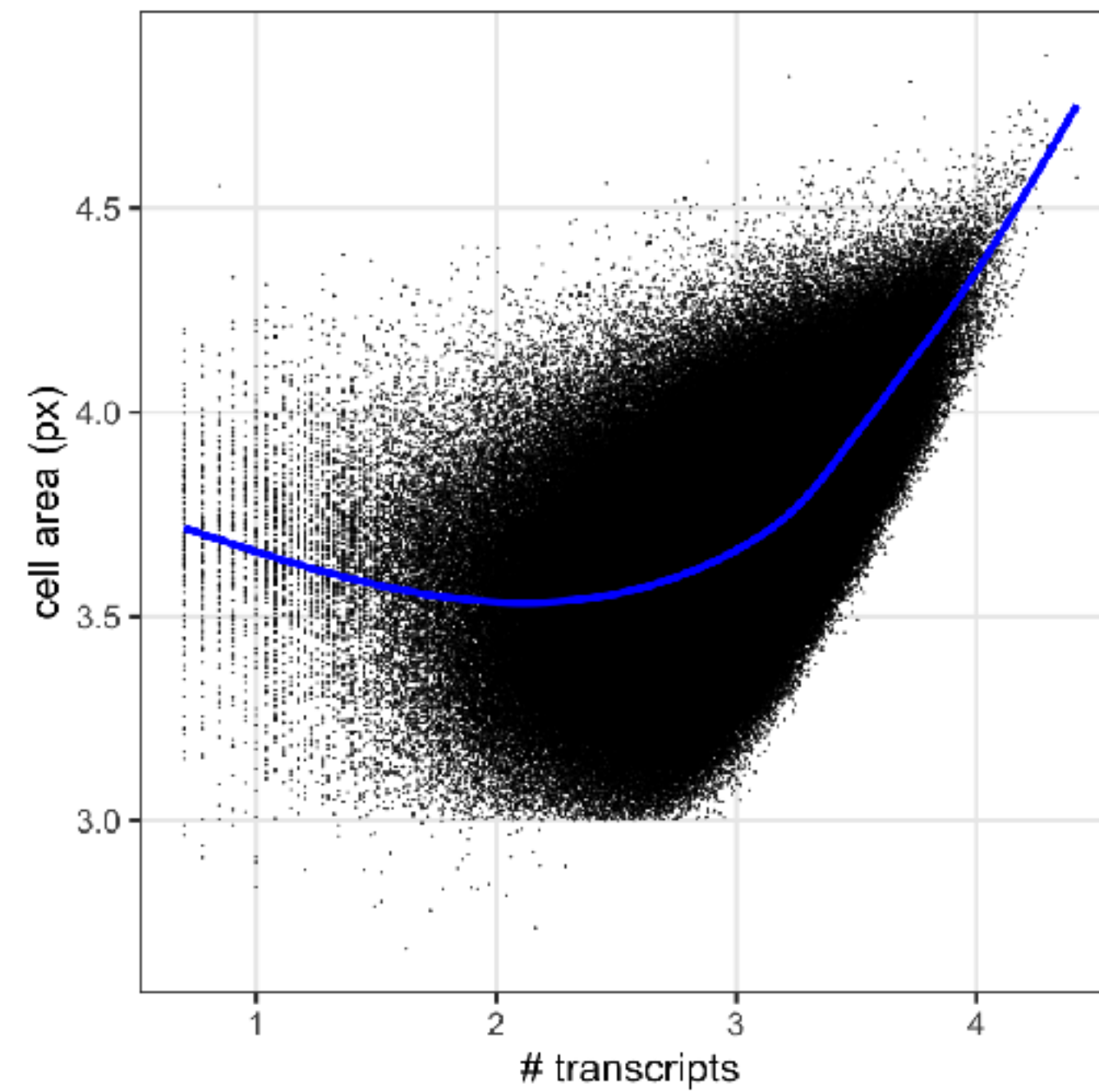


# beware standard QC metrics — counts relate area AND biology

all axes  $\log_{10}$ -transformed (...what we'd usually filter on)



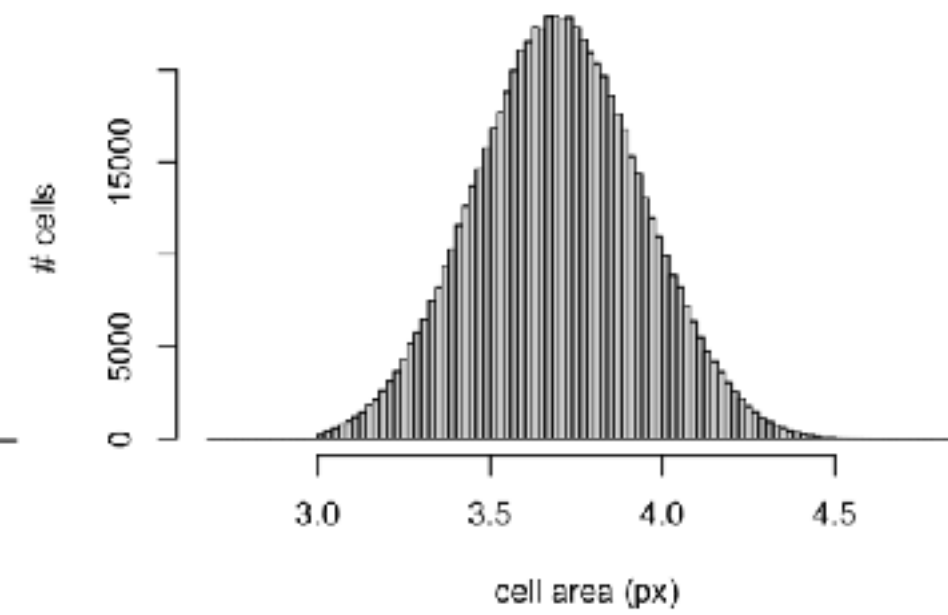
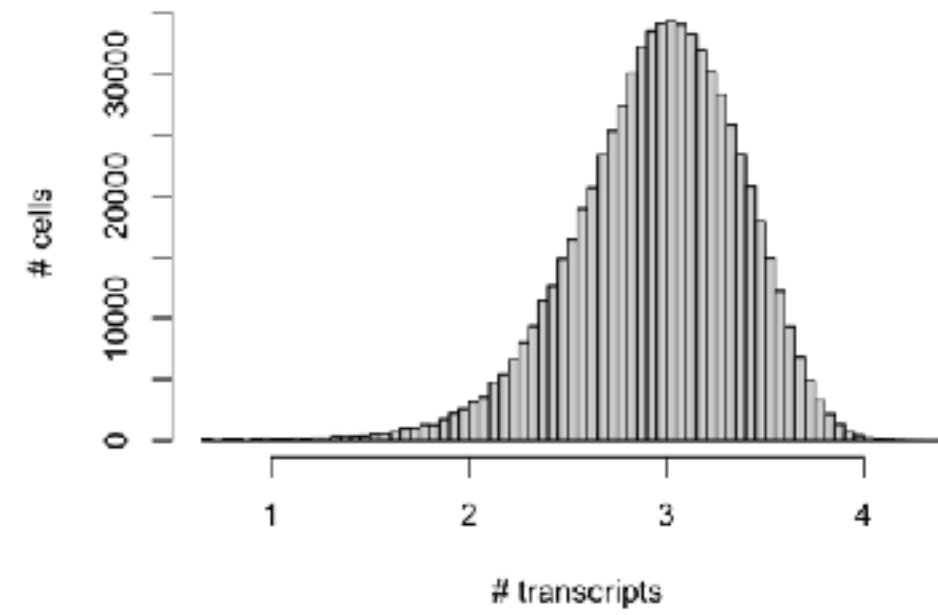
- filtering based on standard scRNA-seq QC metrics may bias against smaller & transcriptionally less complex cells



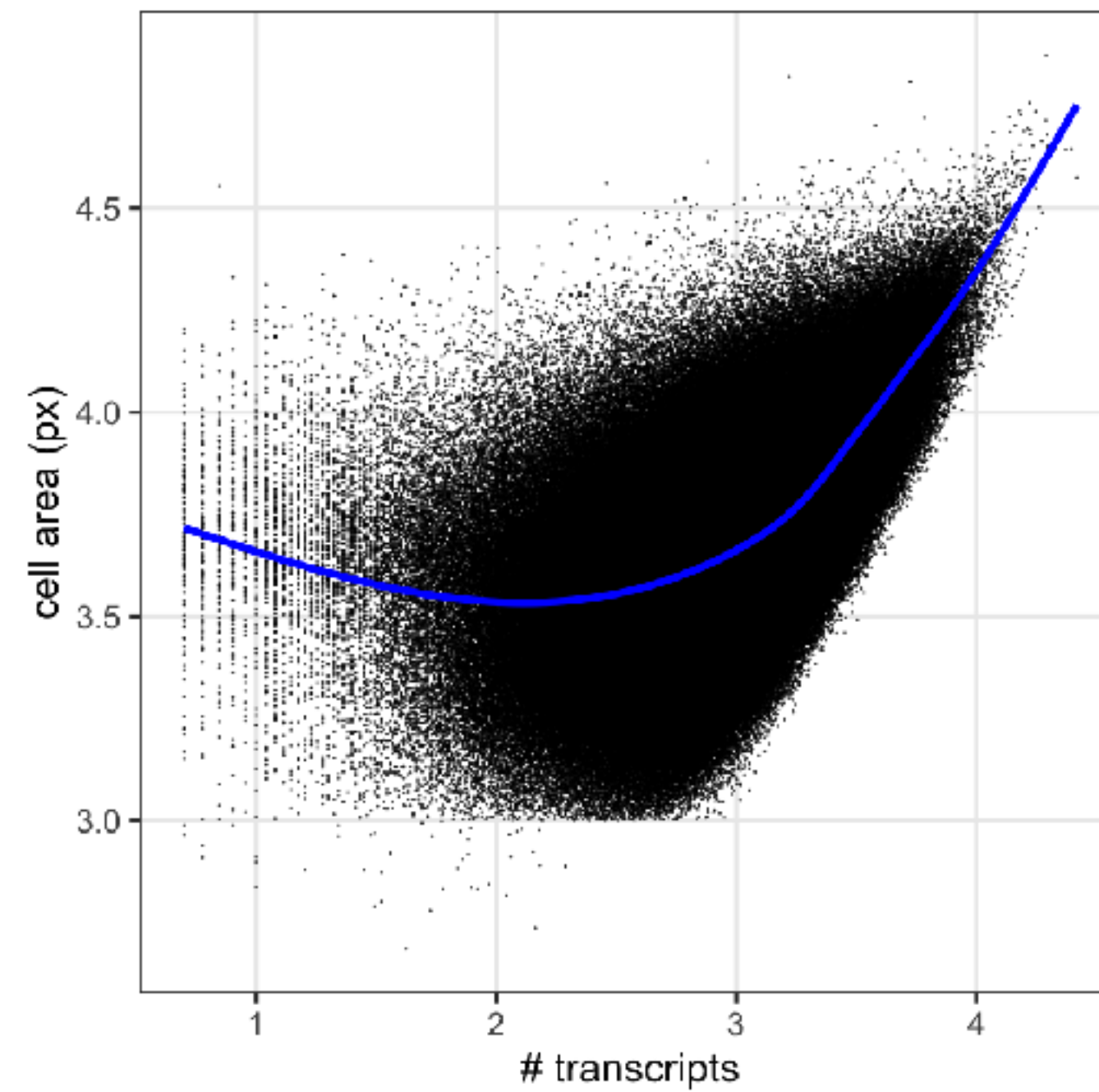
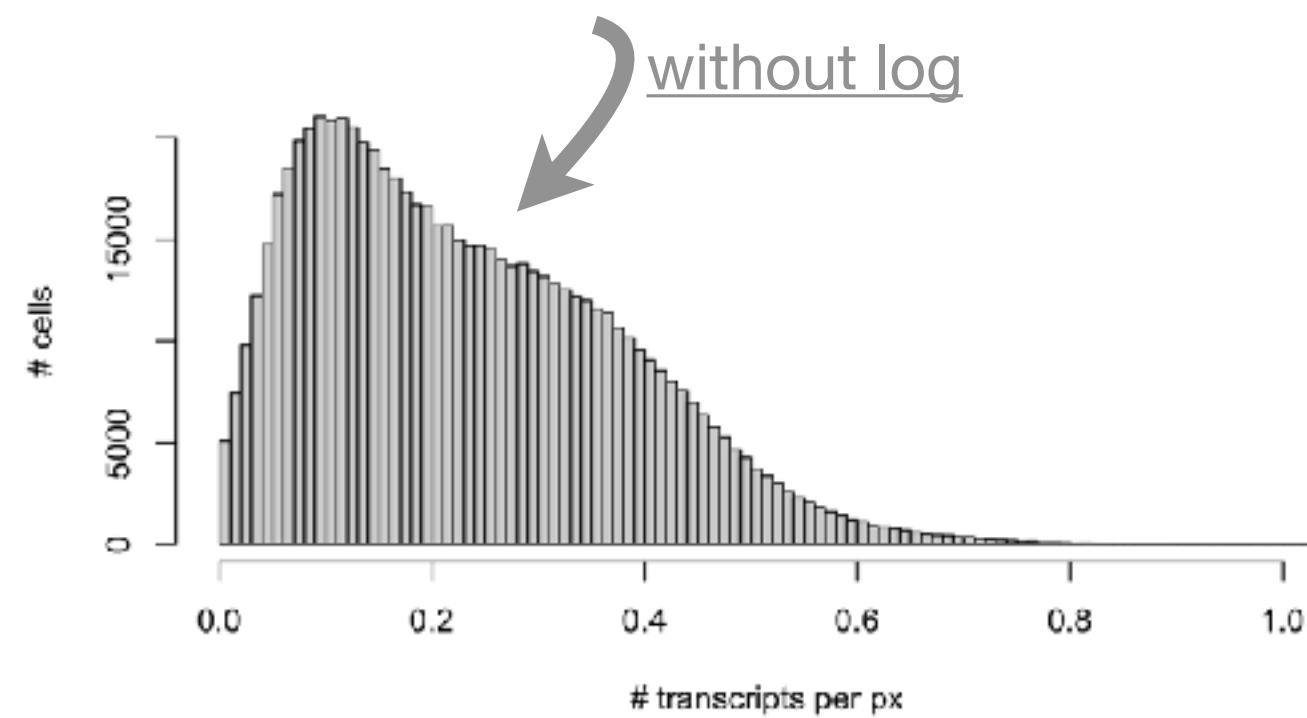
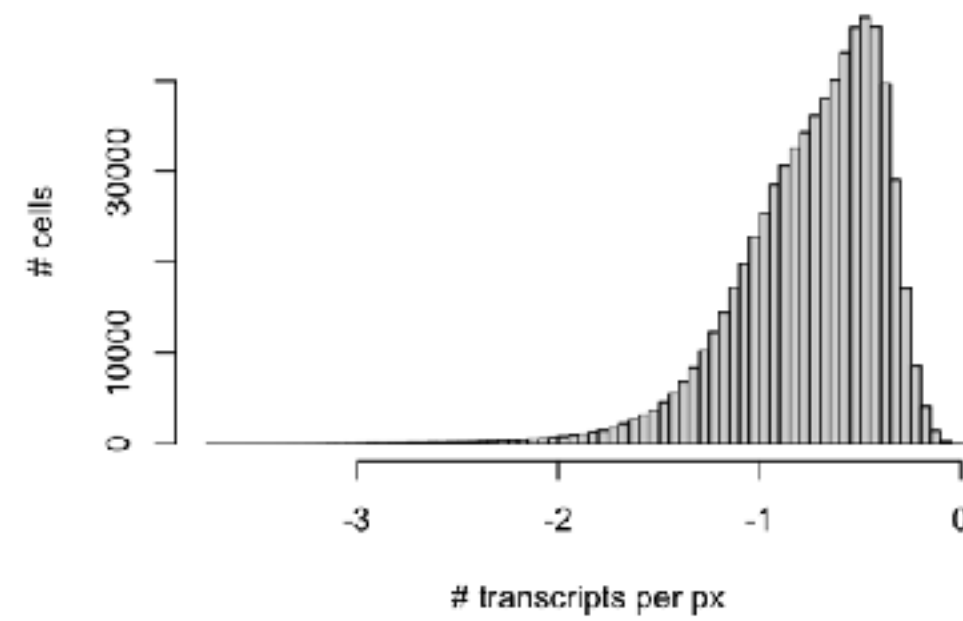


# beware standard QC metrics — counts relate area AND biology

all axes  $\log_{10}$ -transformed (...what we'd usually filter on)



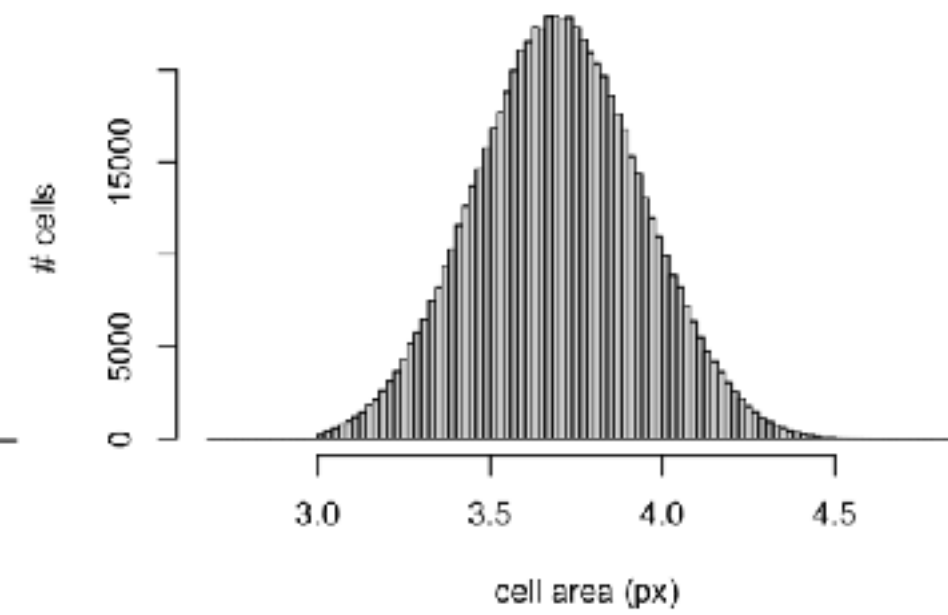
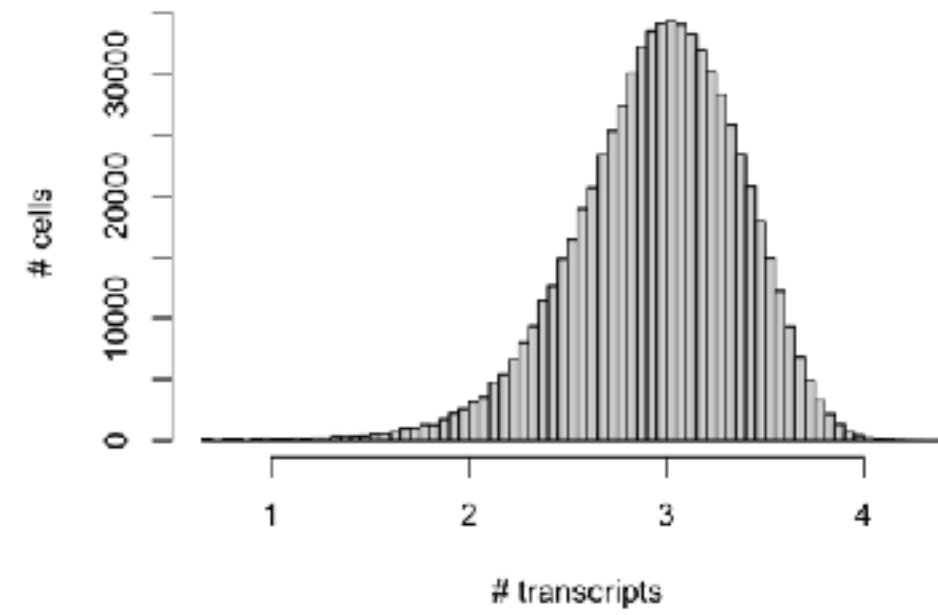
- filtering based on standard scRNA-seq QC metrics may bias against smaller & transcriptionally less complex cells



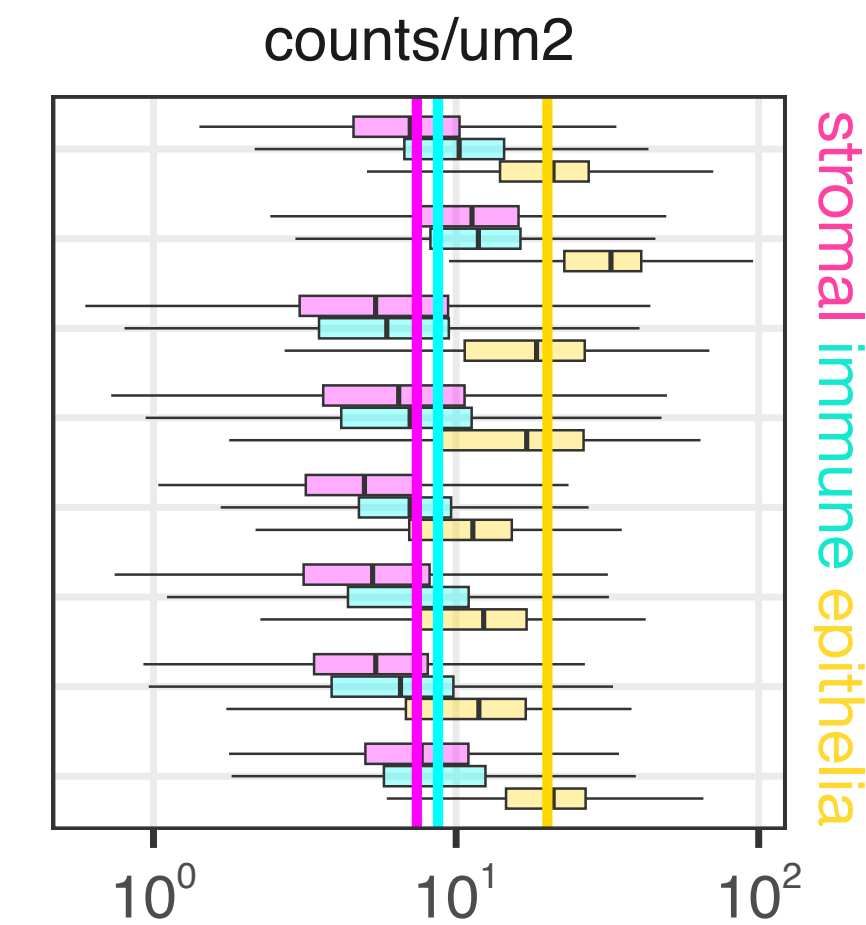
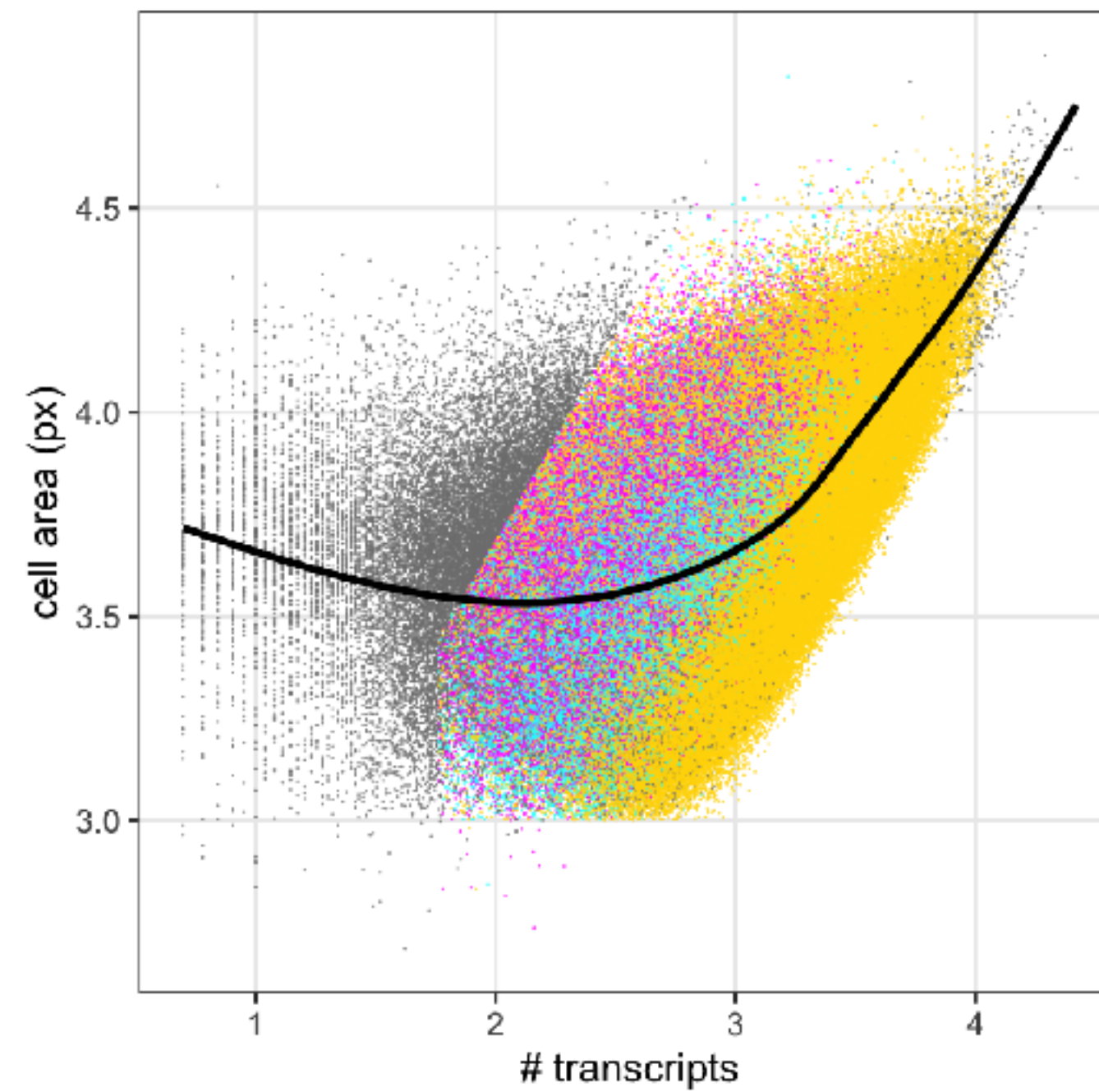
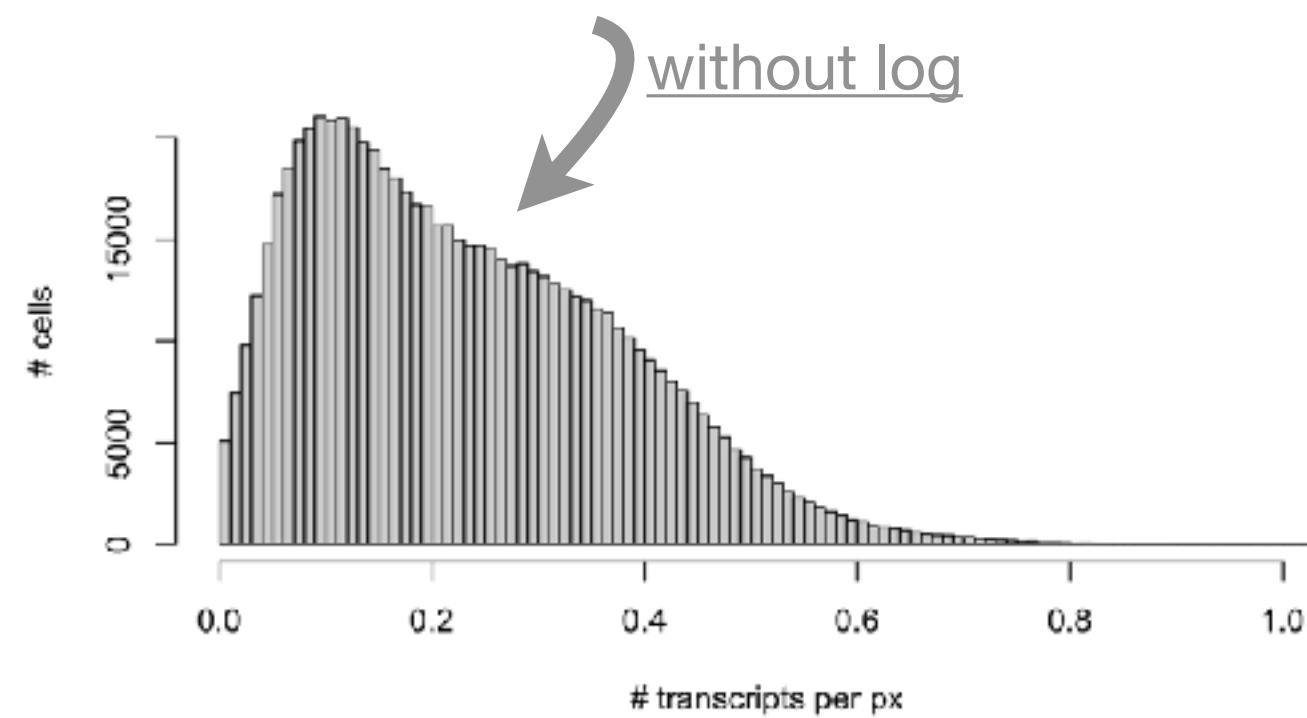
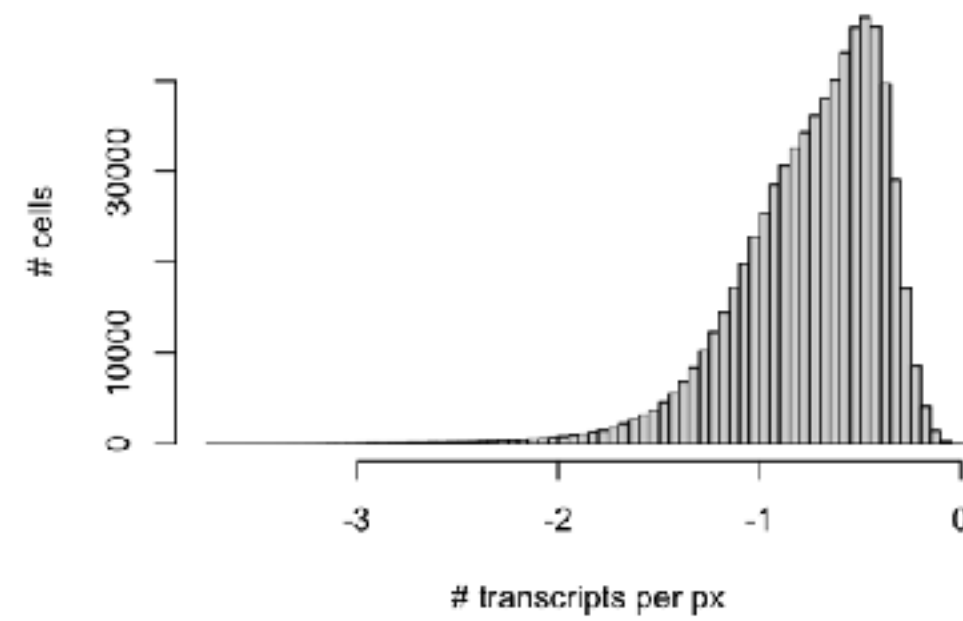


# beware standard QC metrics — counts relate area AND biology

all axes  $\log_{10}$ -transformed (...what we'd usually filter on)



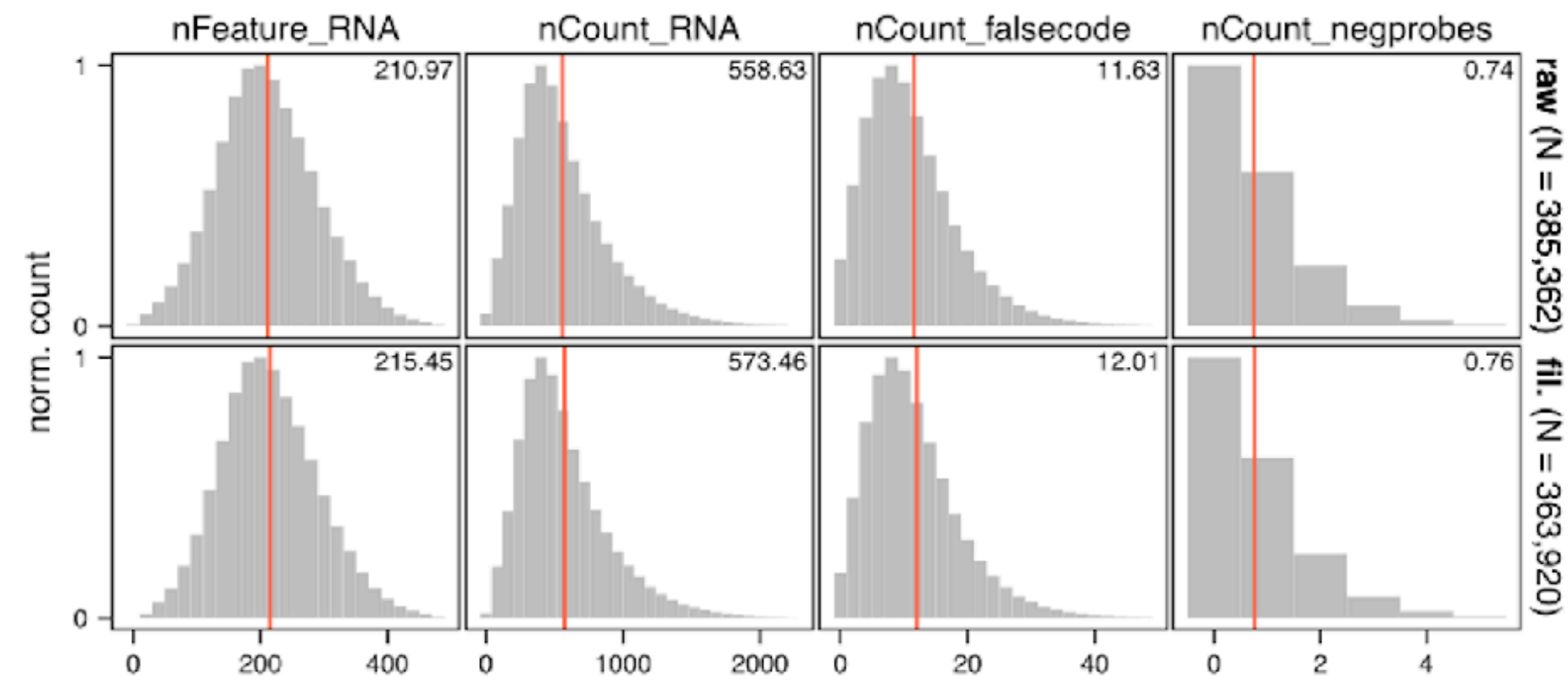
- filtering based on standard scRNA-seq QC metrics may bias against smaller & transcriptionally less complex cells



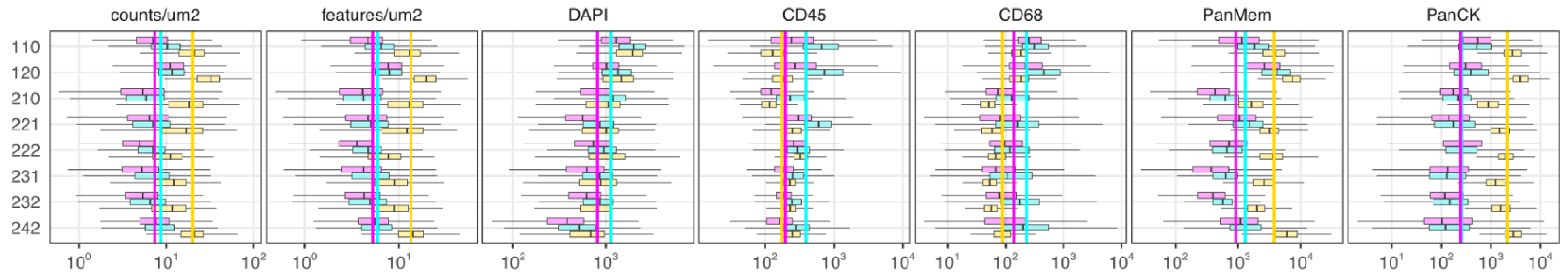
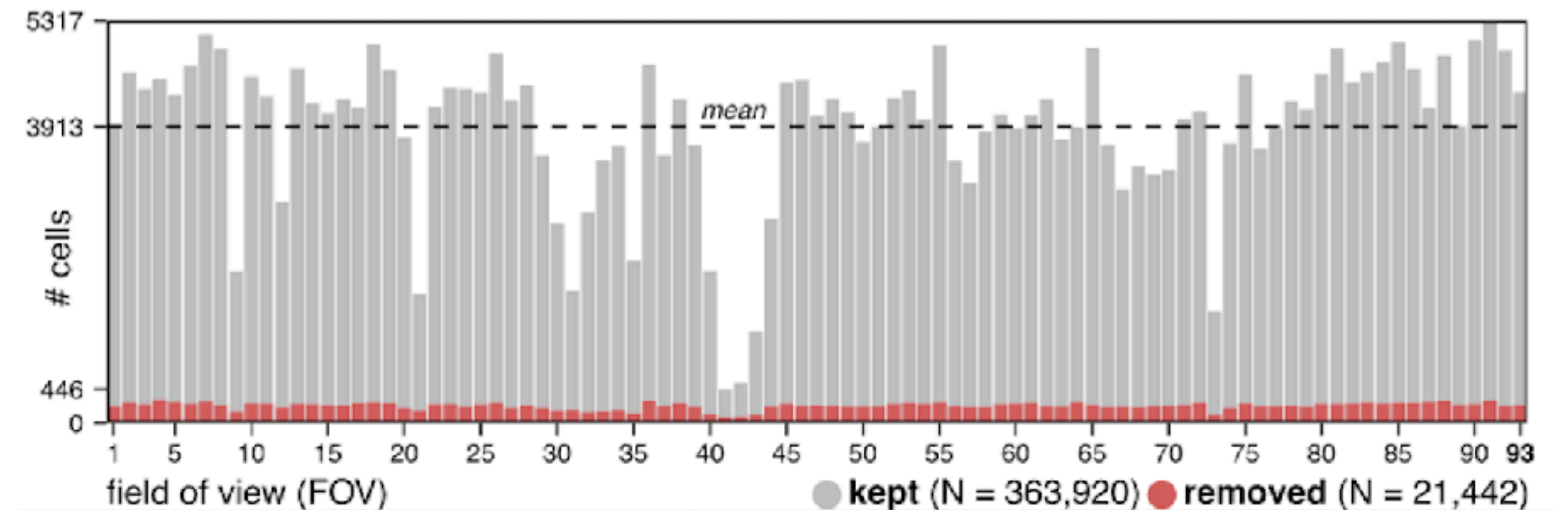


# but there's more we could be looking at... some **QC** examples

*standard QC metrics*



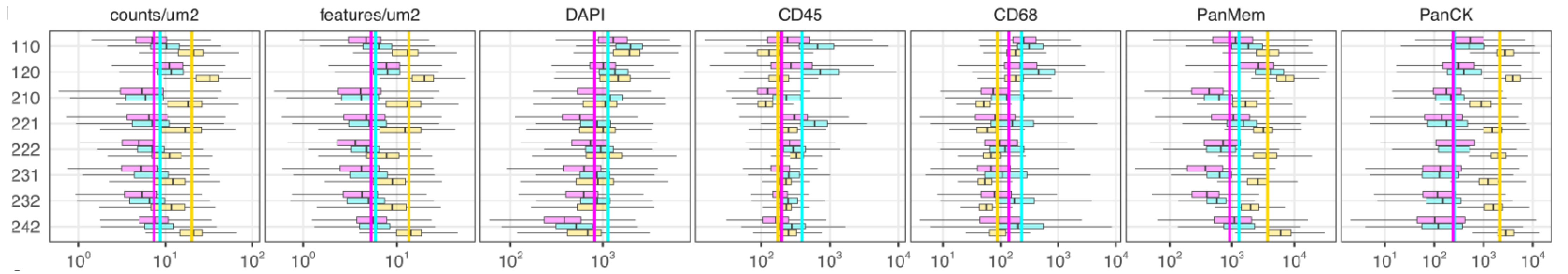
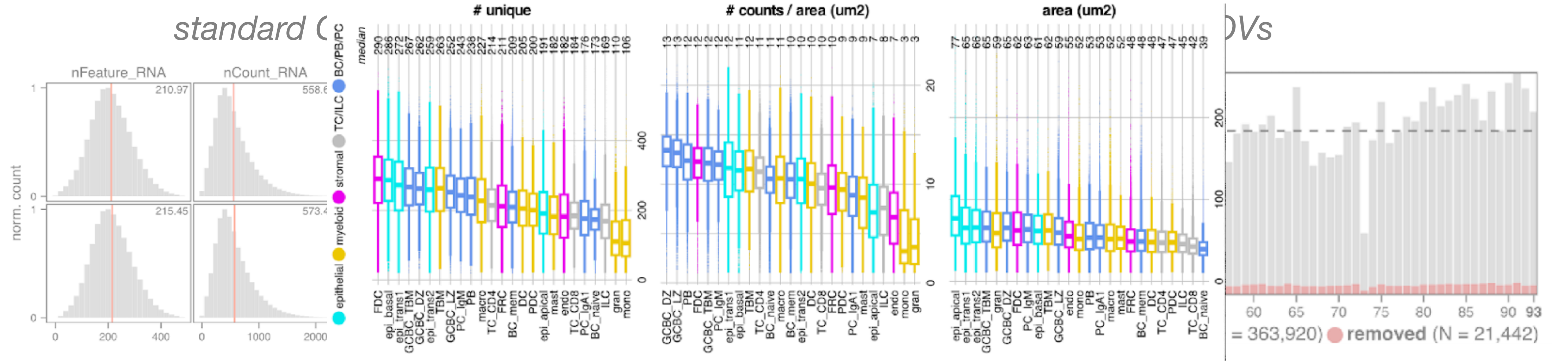
*# cells across FOVs*



*RNA per area & IF markers*



but there's more we could be looking at... some **QC examples**



*RNA per area & IF markers*



spatial organization of excluded cells might indicate a **bias against specific types** (but it depends!)





# img-ST data stem from serial imaging of **fields of view (FOVs)**

H&E staining

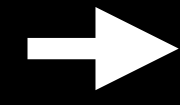


REF TVA CRC

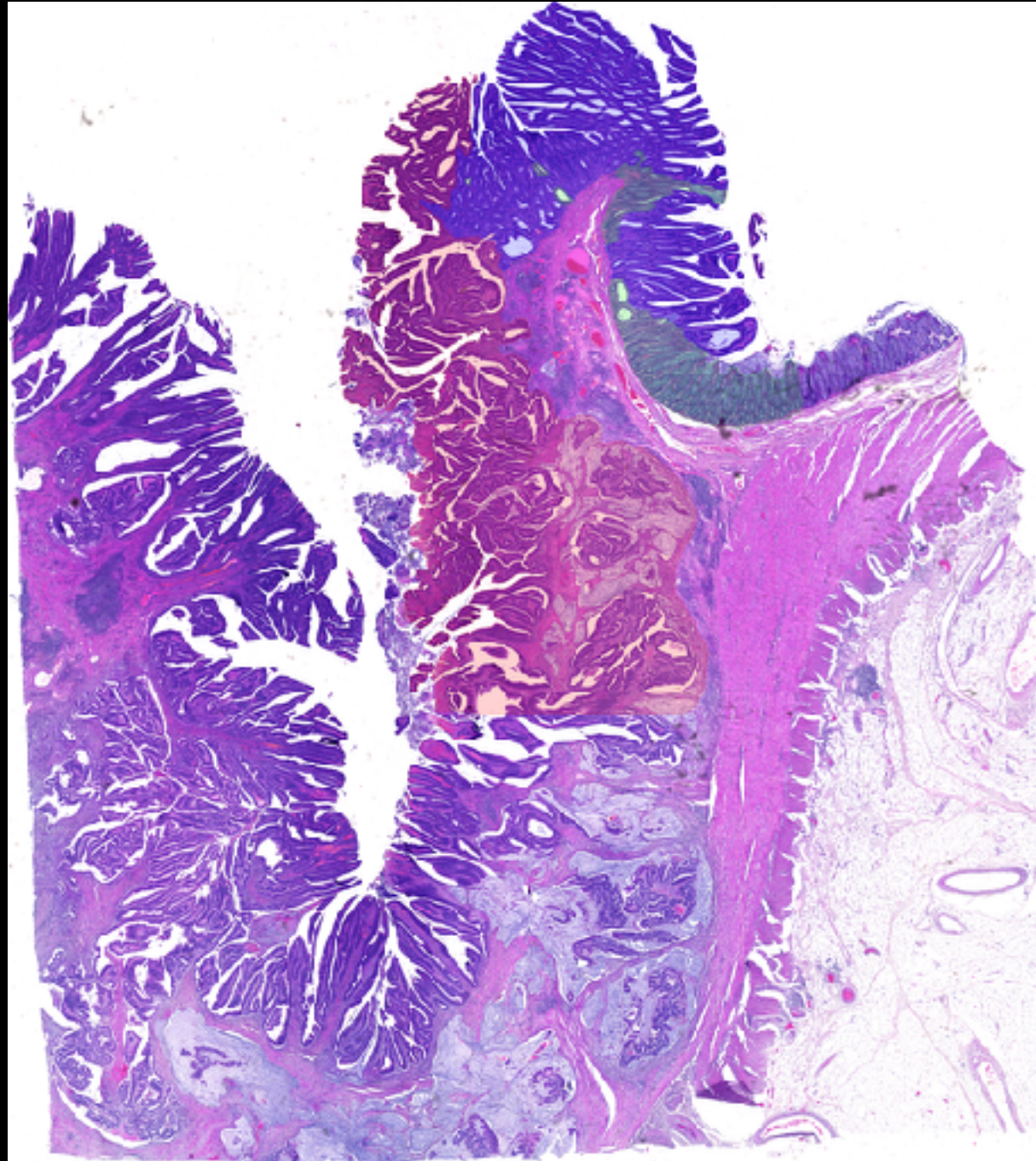


# img-ST data stem from serial imaging of **fields of view (FOVs)**

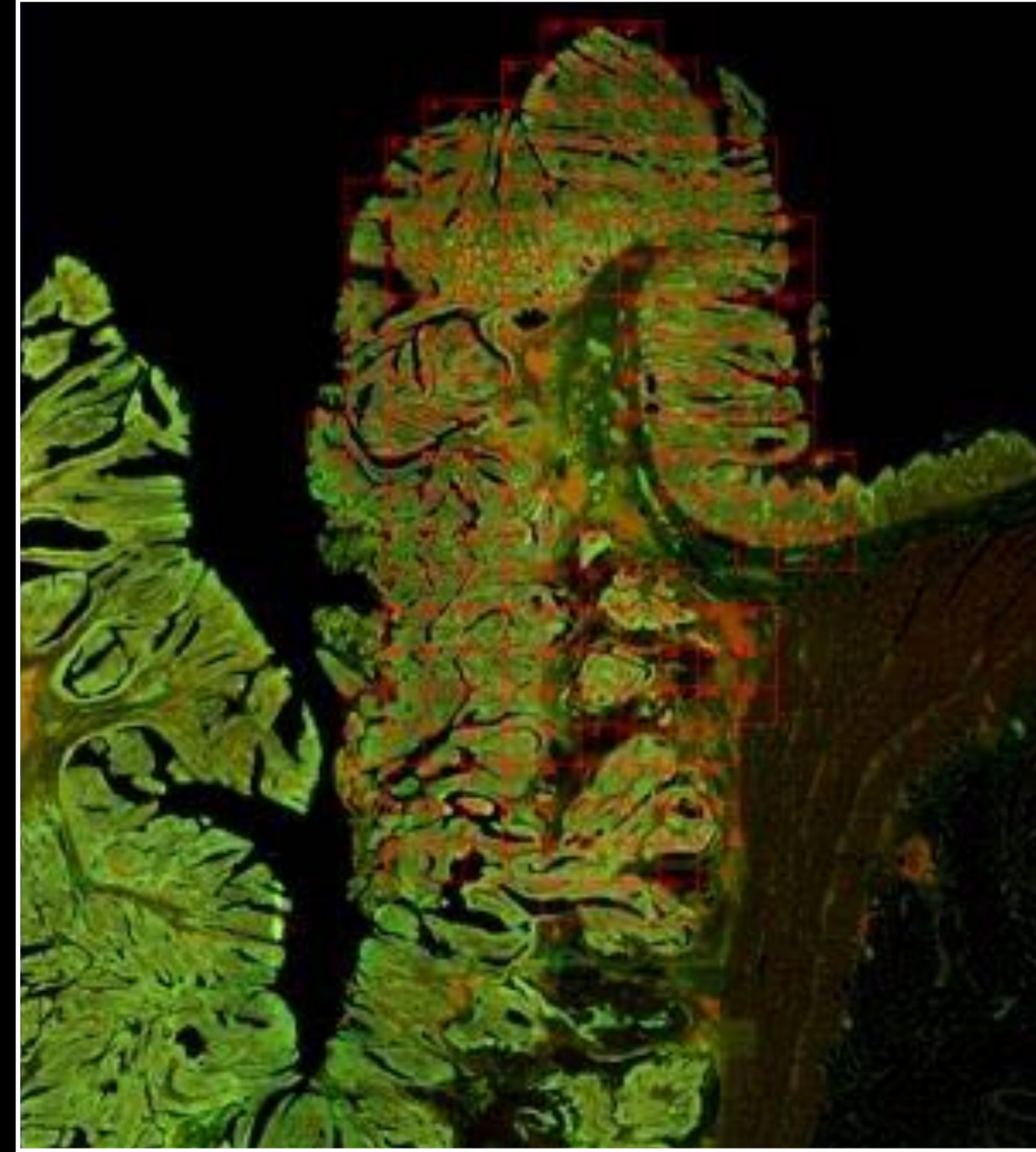
H&E staining



FOV placement



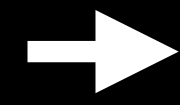
REF TVA CRC



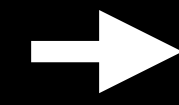


# img-ST data stem from serial imaging of **fields of view (FOVs)**

H&E staining



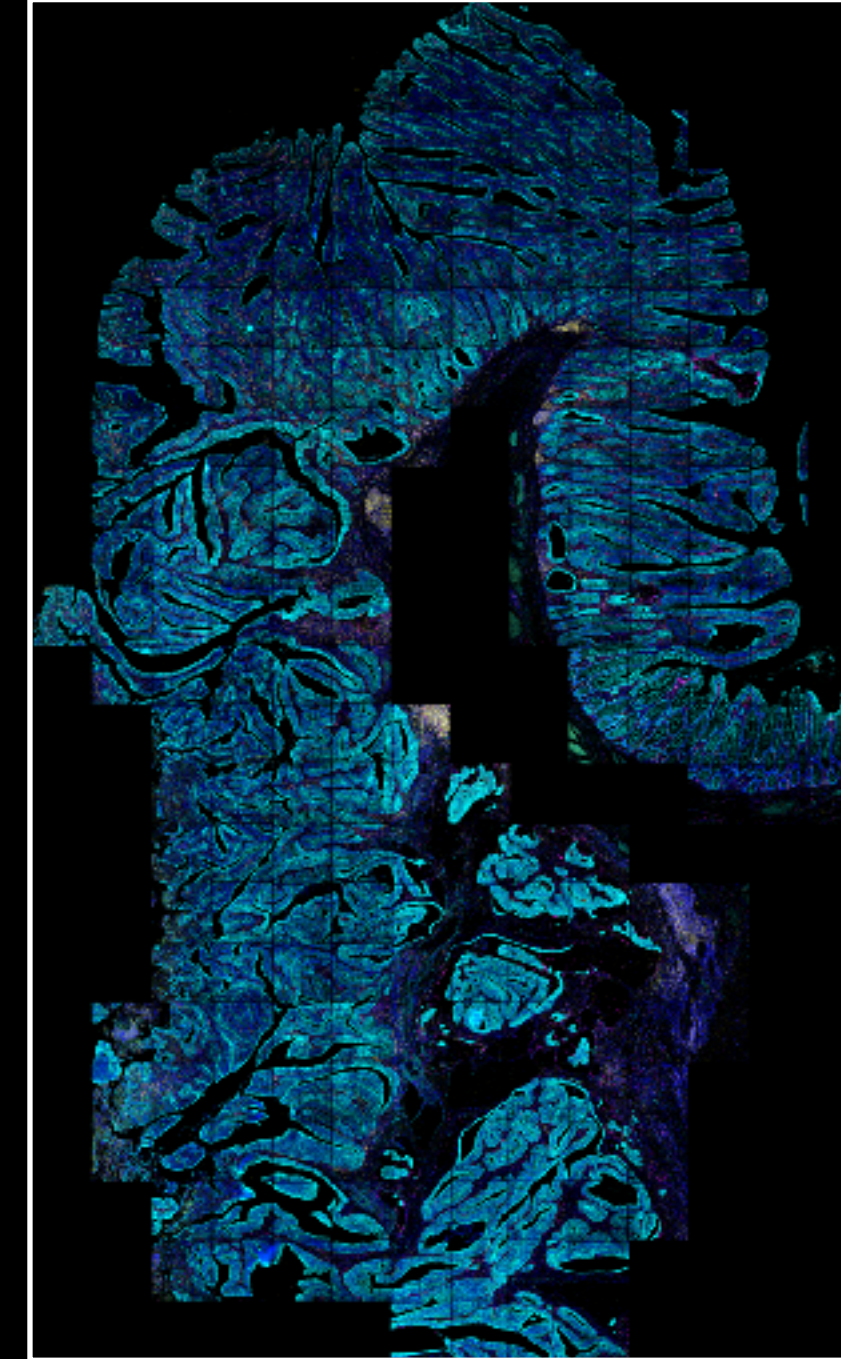
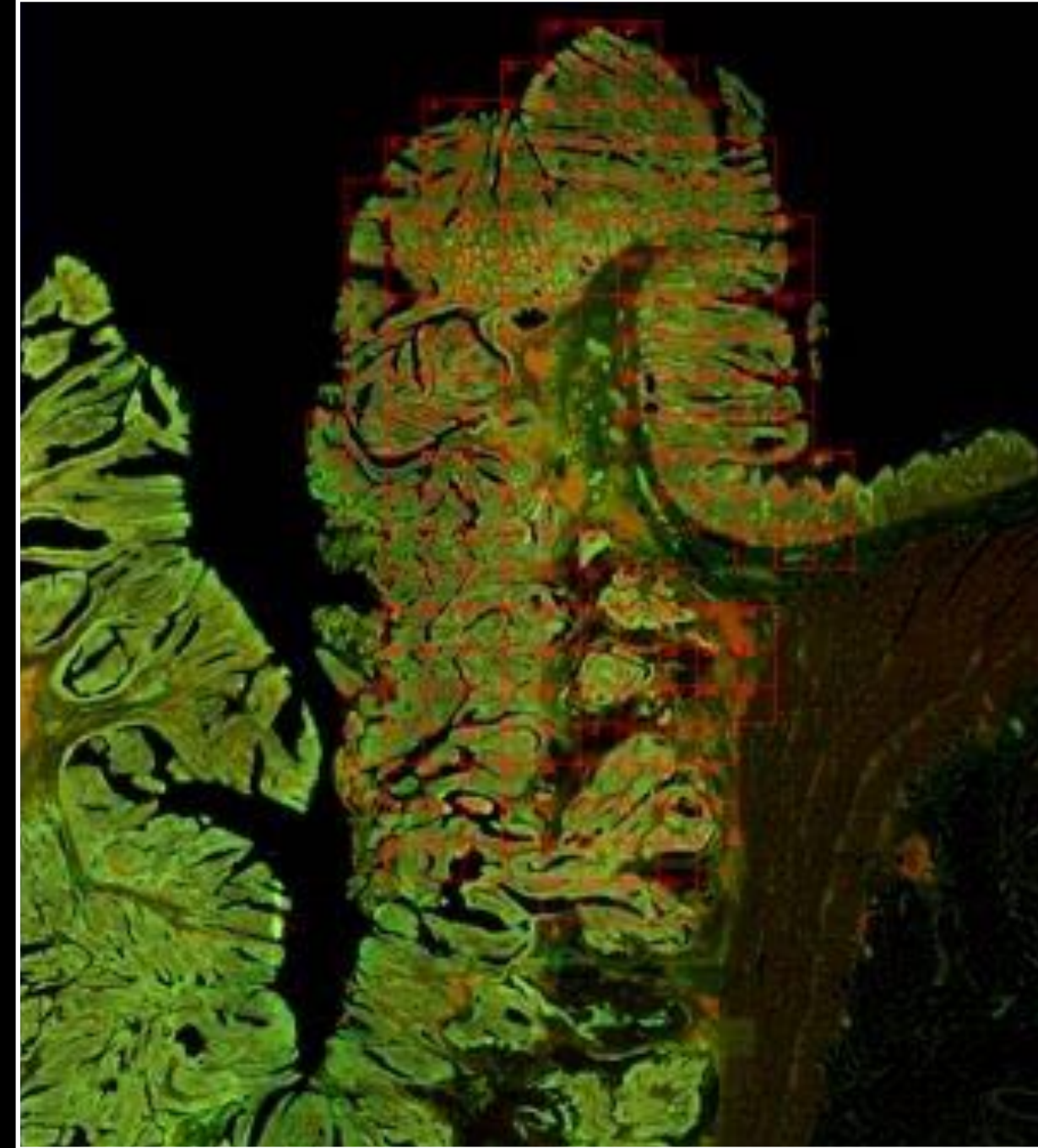
FOV placement



imaging



REF TVA CRC

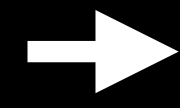


DAPI PanCK CD45 CD68

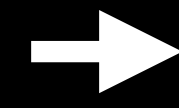


# img-ST data stem from serial imaging of **fields of view (FOVs)**

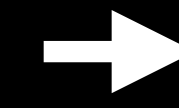
H&E staining



FOV placement



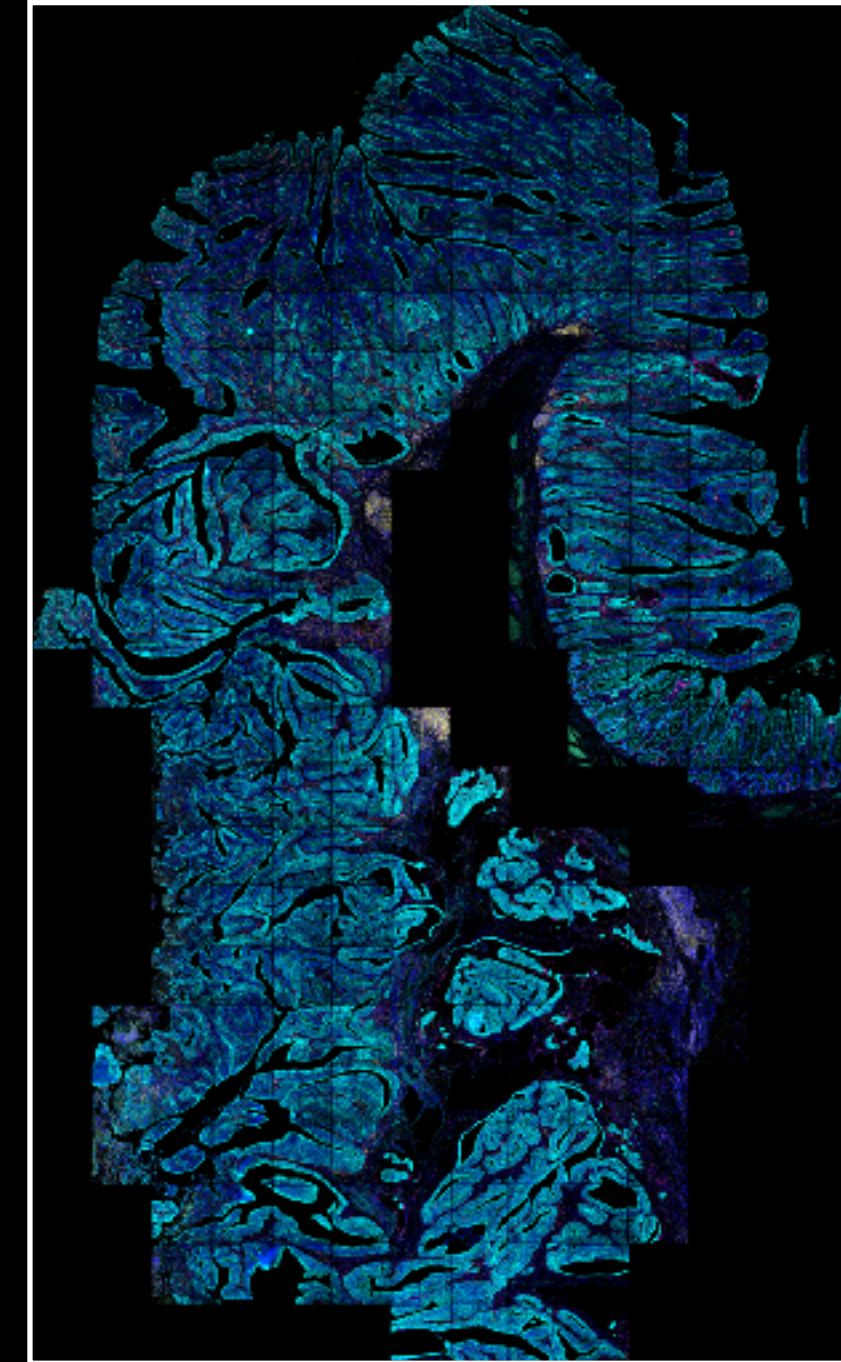
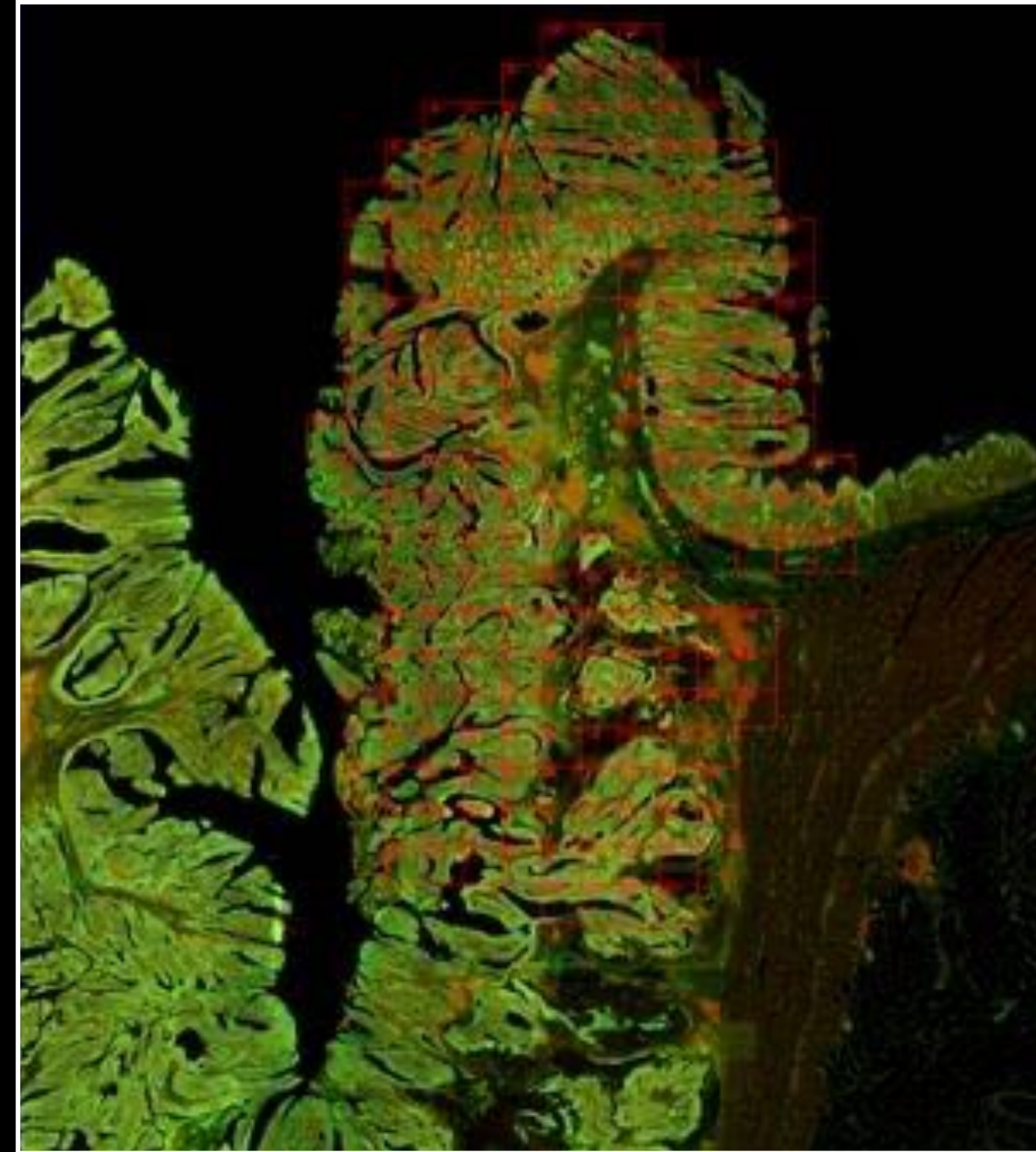
imaging



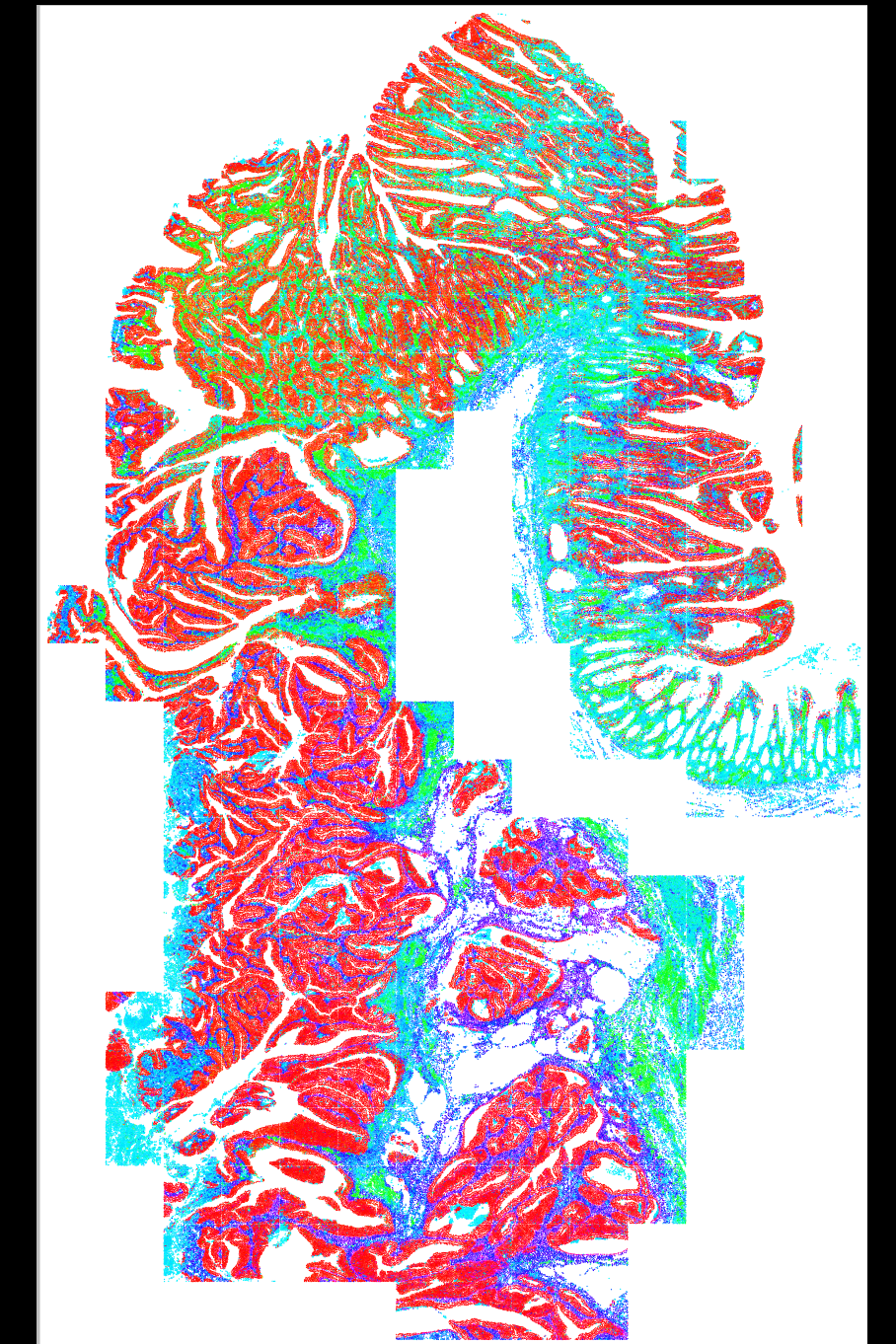
analysis



REF TVA CRC



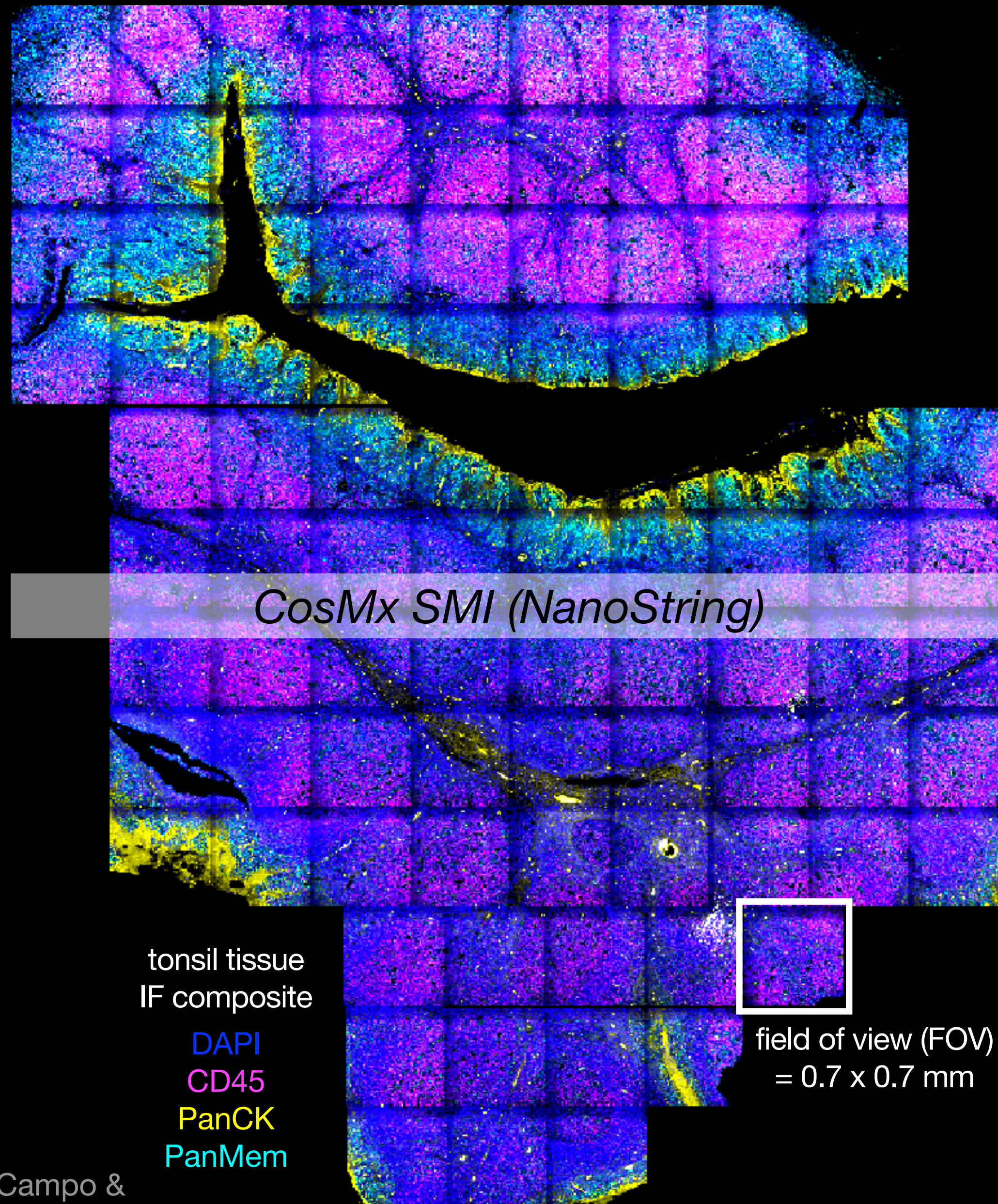
DAPI PanCK CD45 CD68



PC1 PC2 PC3



# changes in optical performance lead to **lower optical resolution at FOV edges**

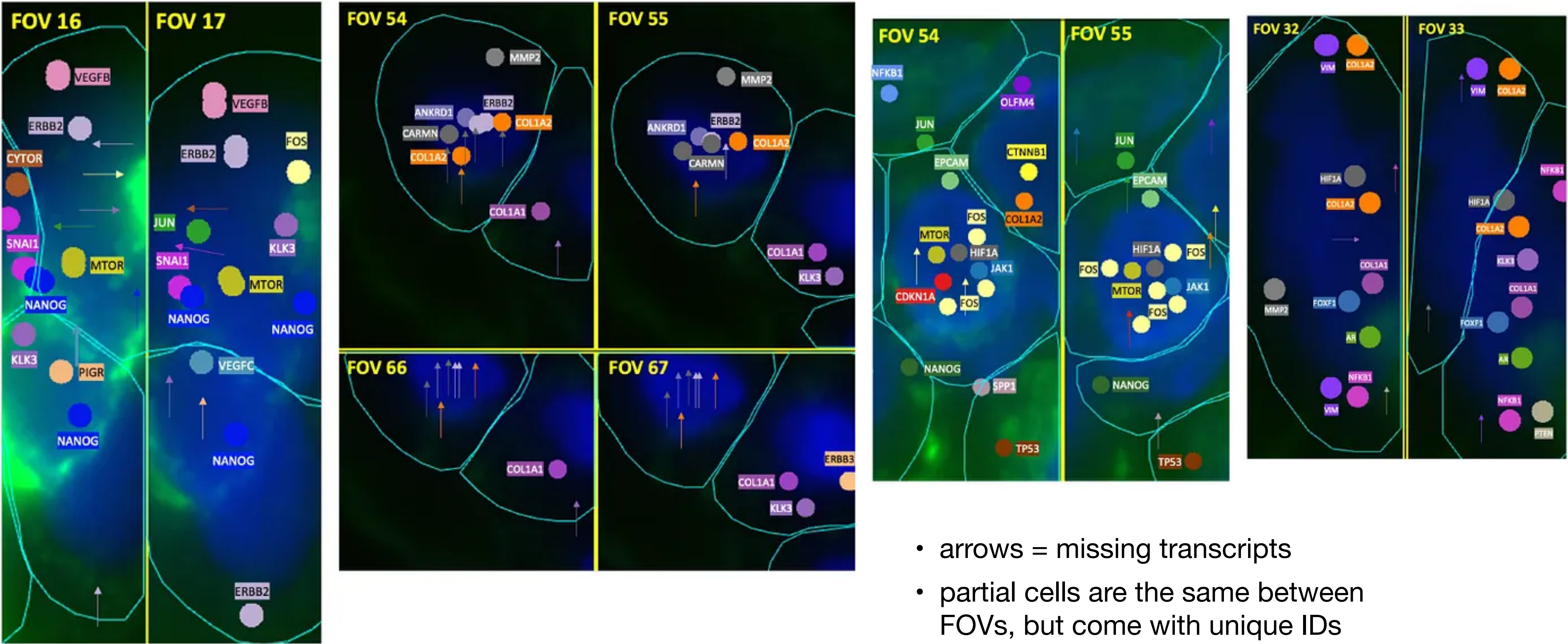






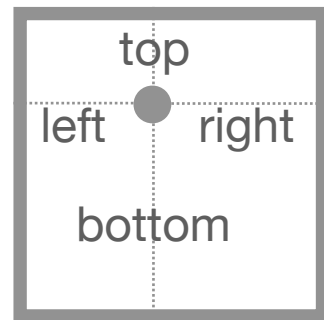


# lack of stitching leads to cell fragmentation, duplication & inconsistencies



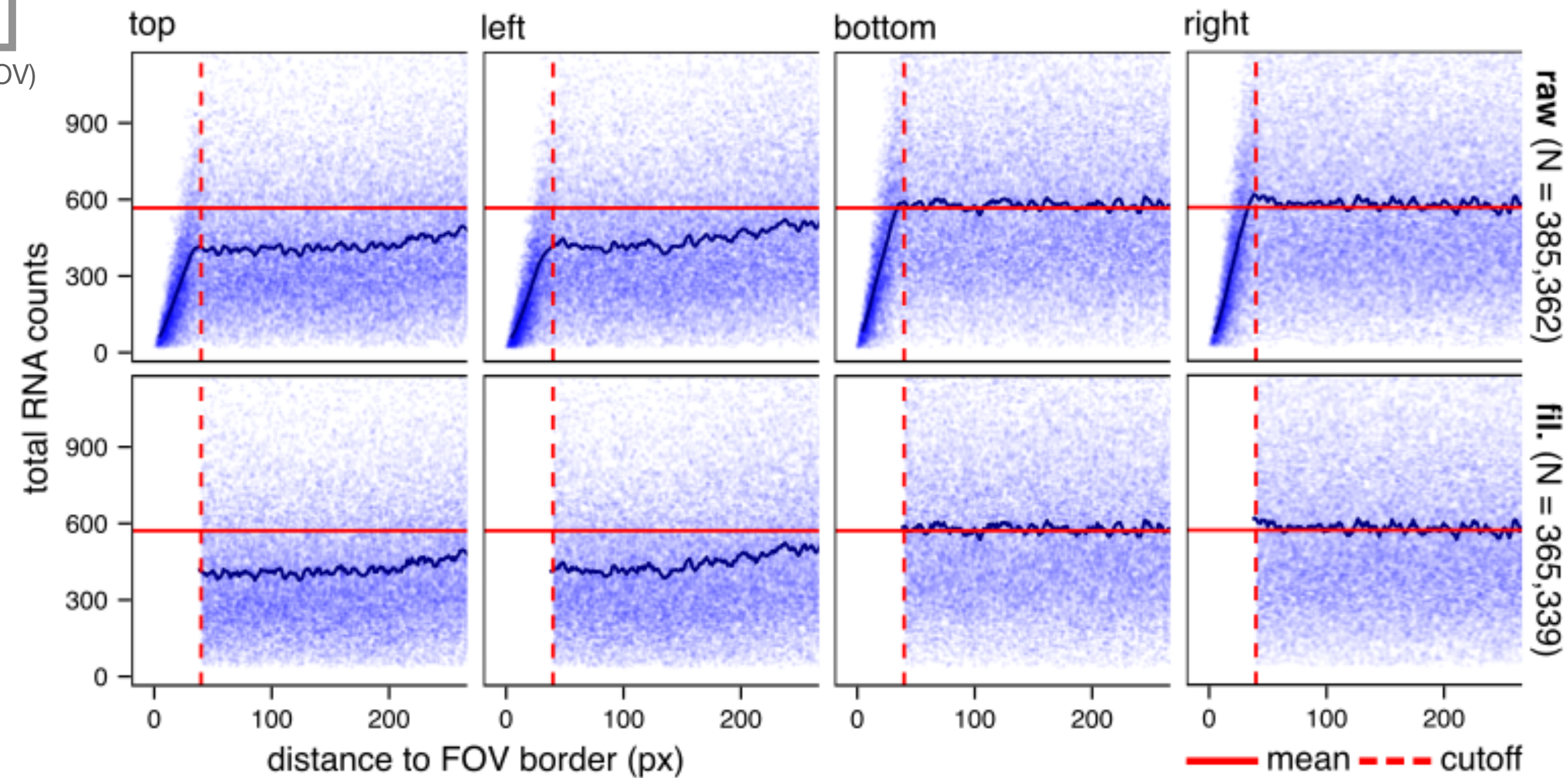


# border cells have fewer counts & can highlight other artefacts



field of view (FOV)

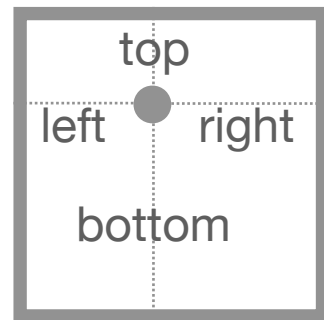
- compute each cell's distance to each FOV border
- plot counts vs. distance, stratified by direction



points = cells (across all FOVs), lines = running median

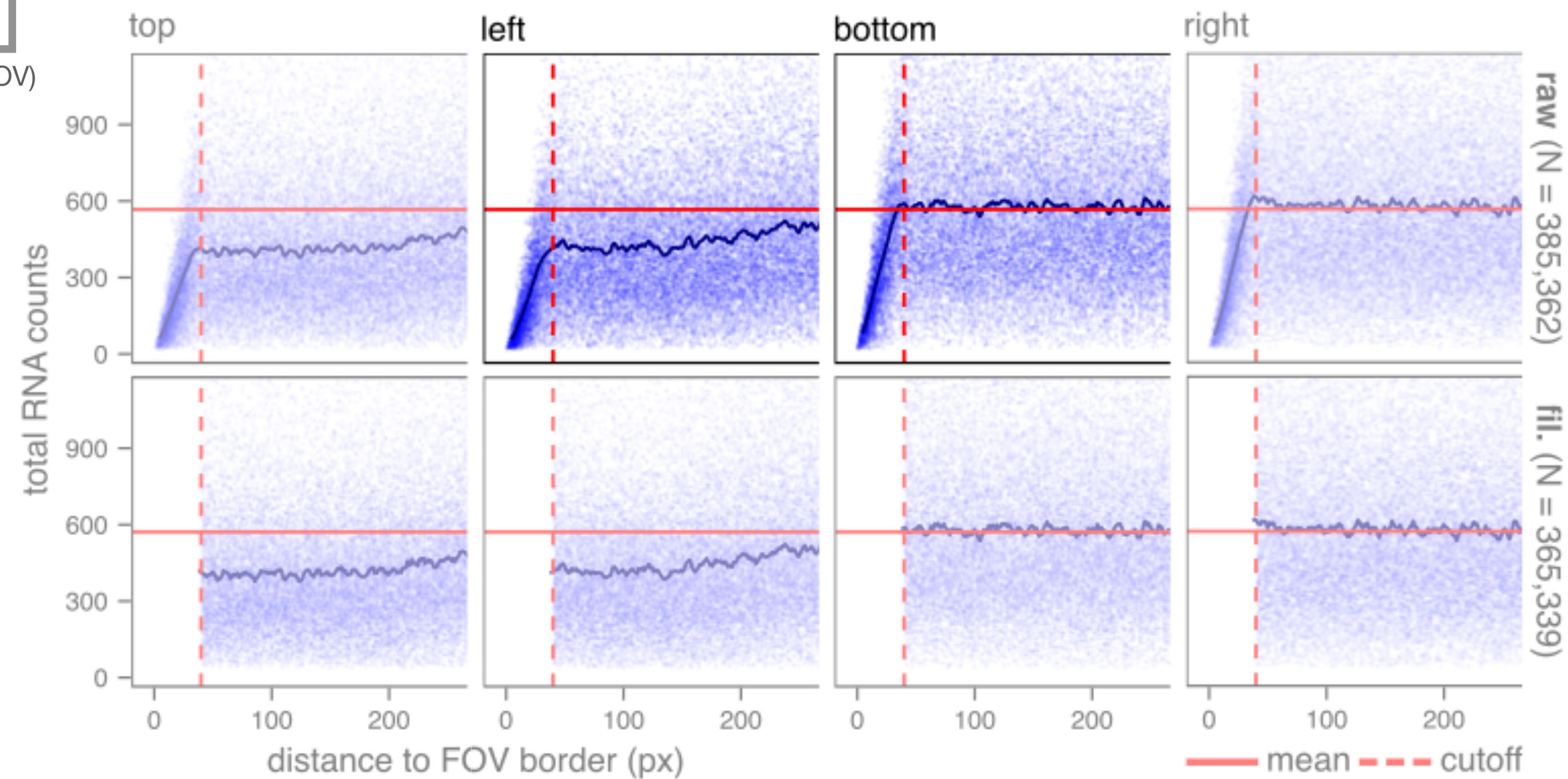


# border cells have fewer counts & can highlight other artefacts

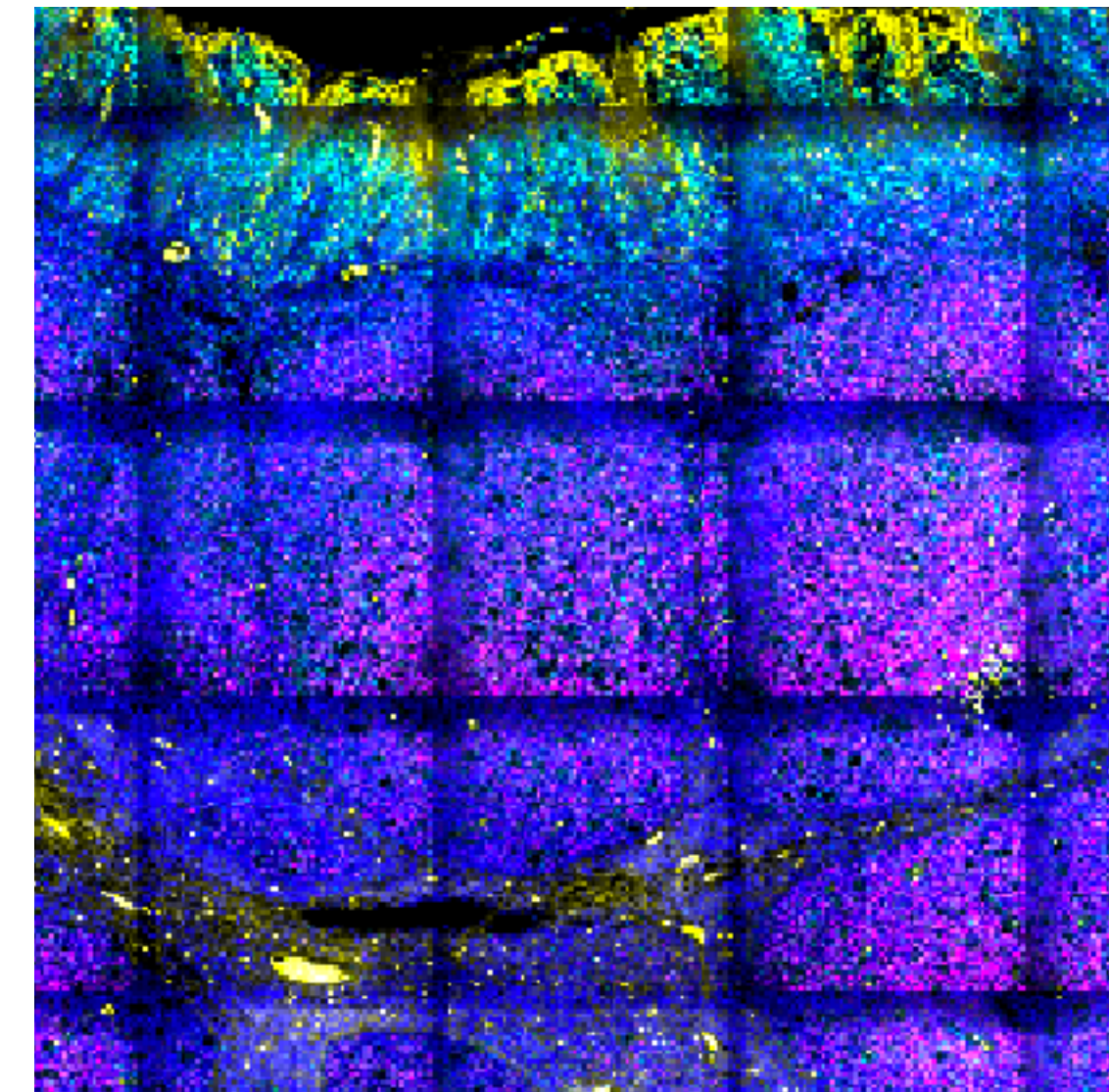


field of view (FOV)

- compute each cell's distance to each FOV border
- plot counts vs. distance, stratified by direction



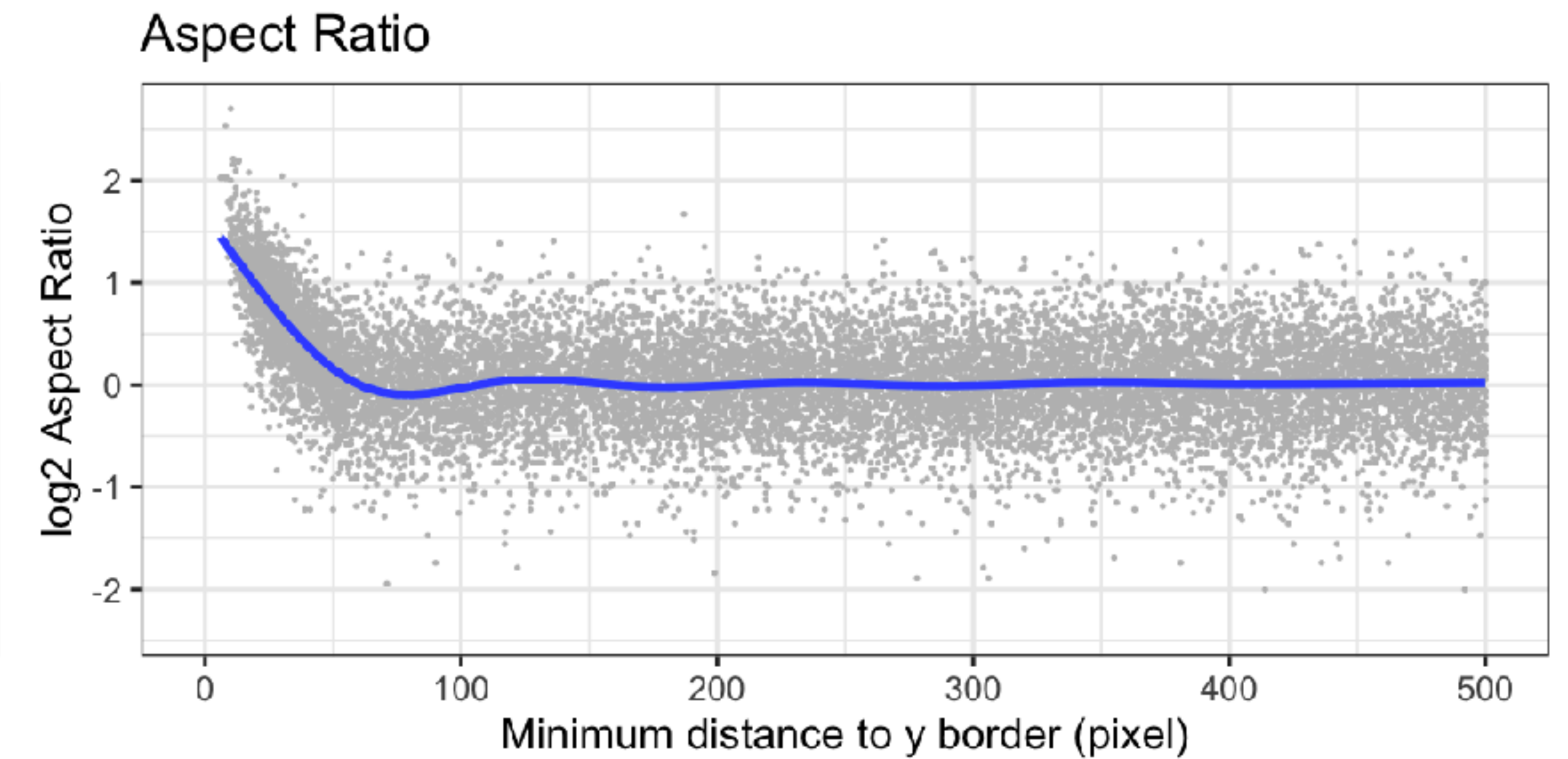
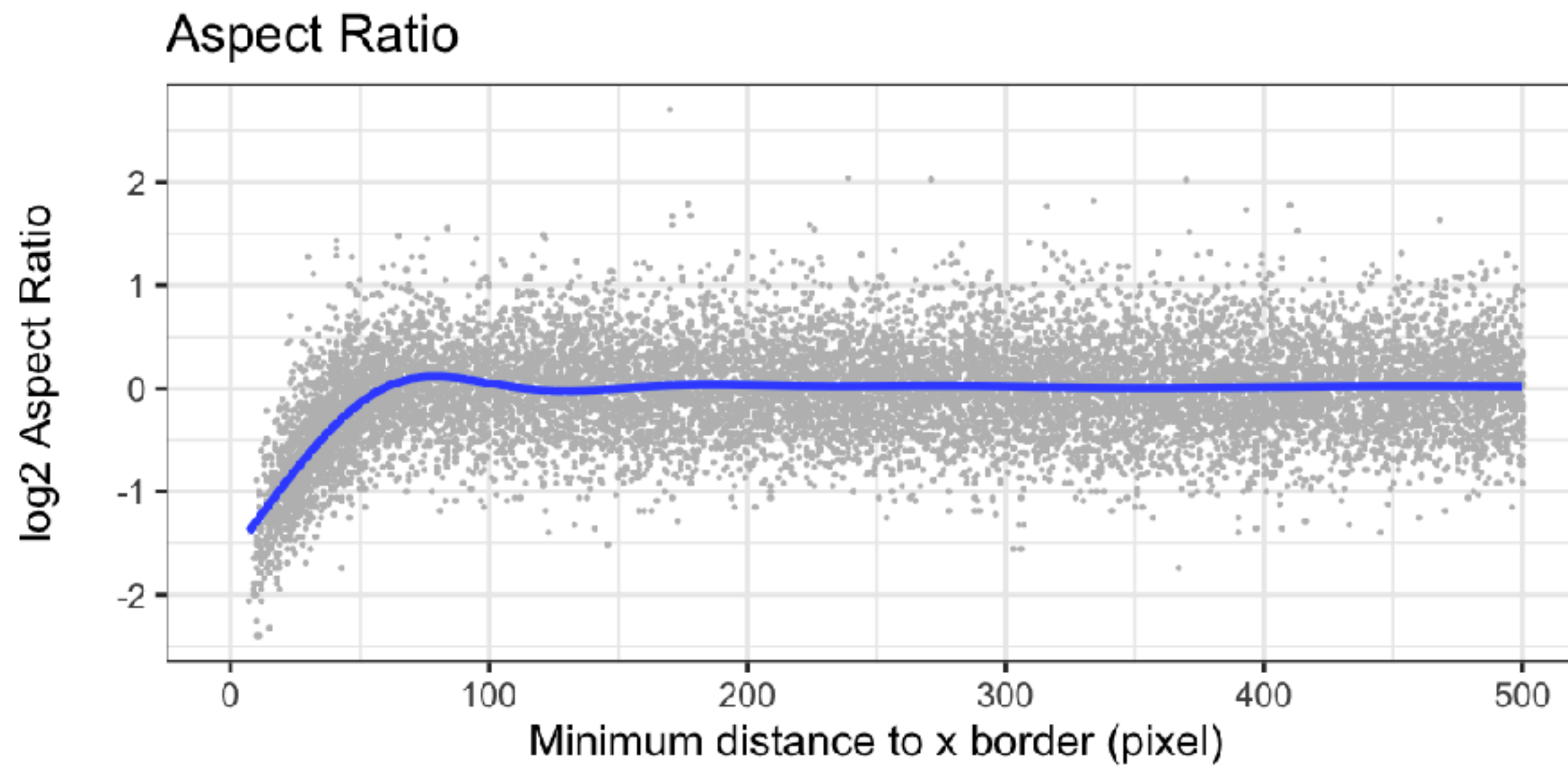
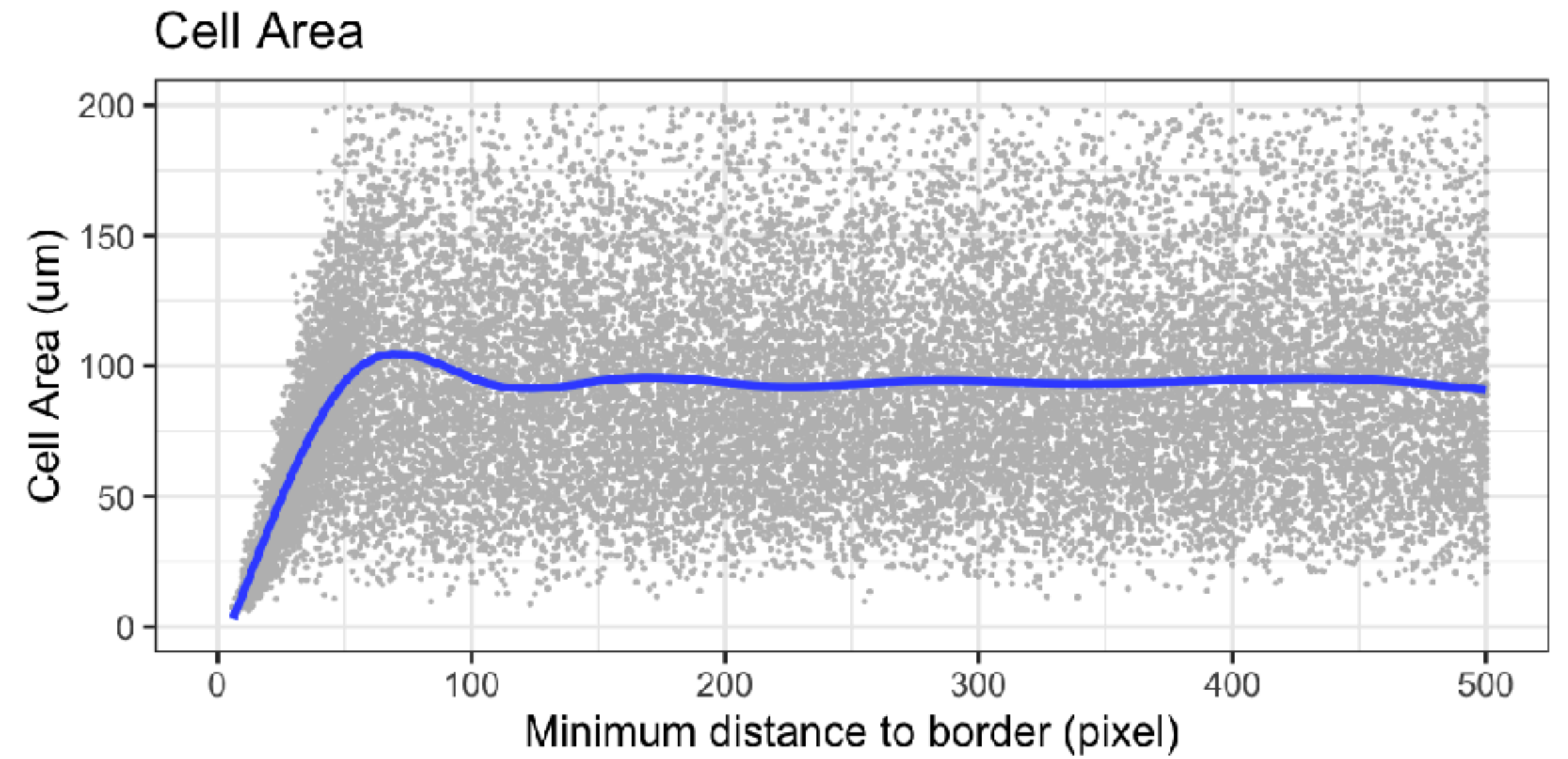
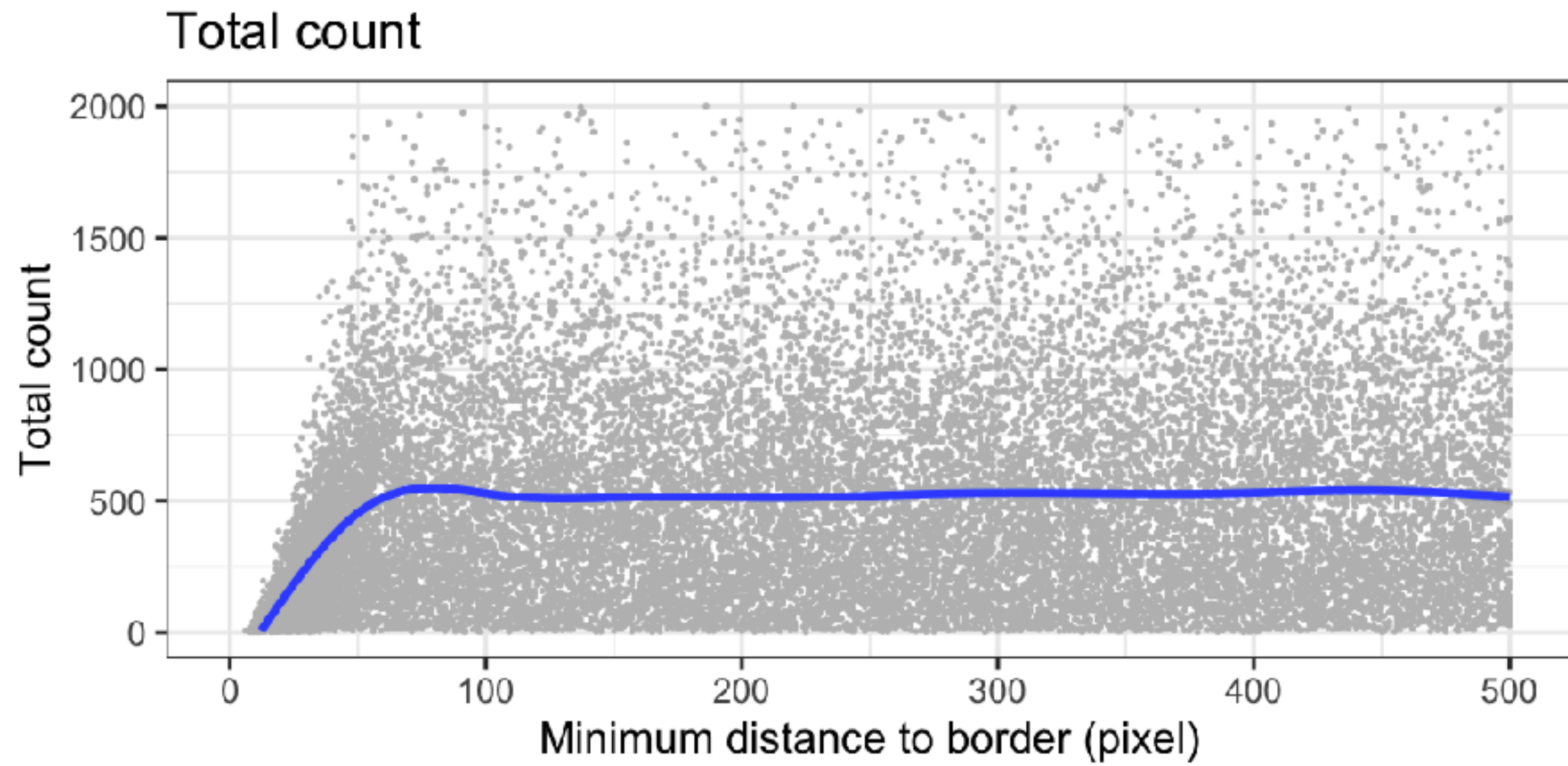
points = cells (across all FOVs), lines = running median



top-left corners exhibit dimmer IF signals

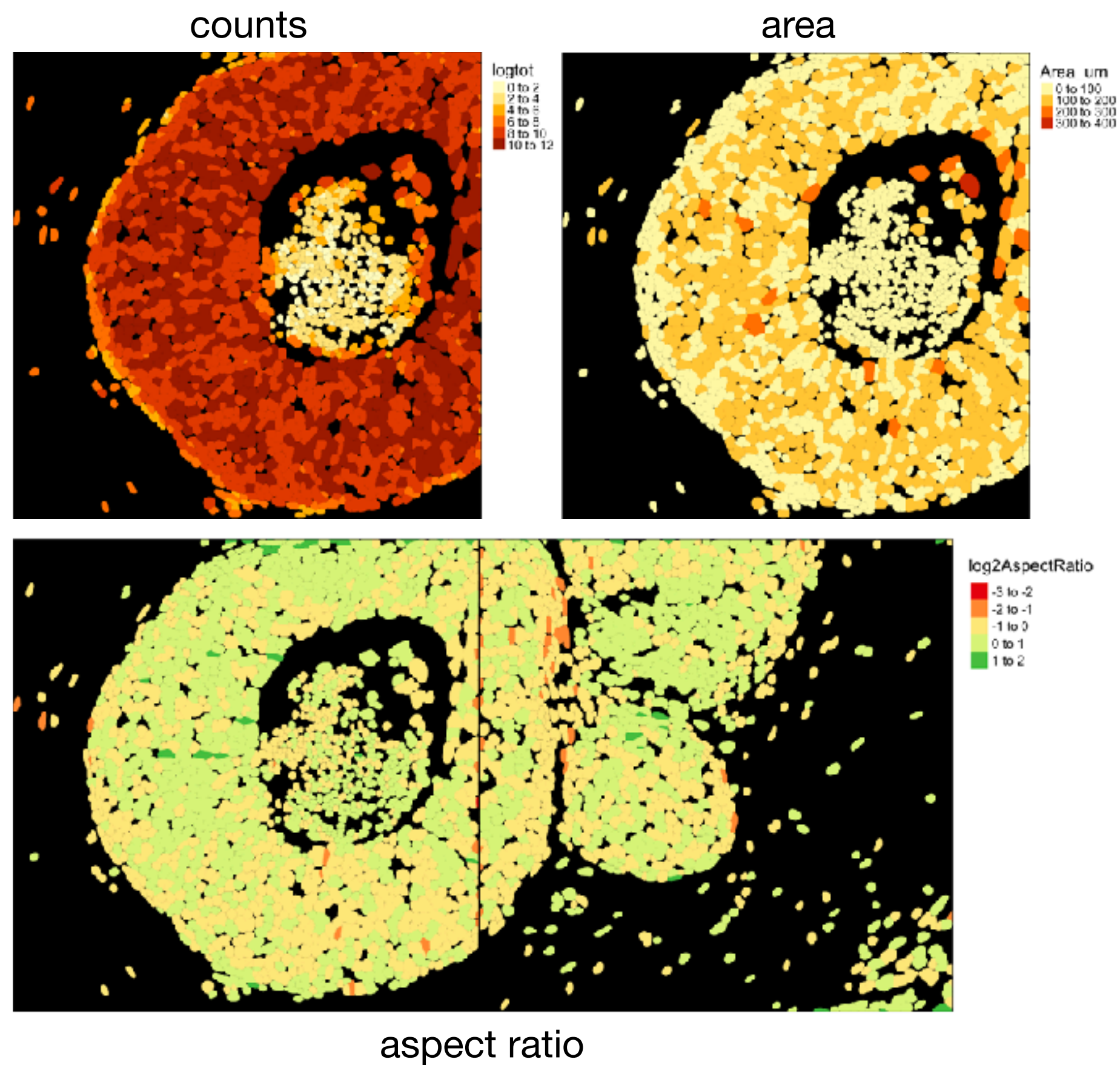


besides fewer counts, **border cells are smaller & slimmer**



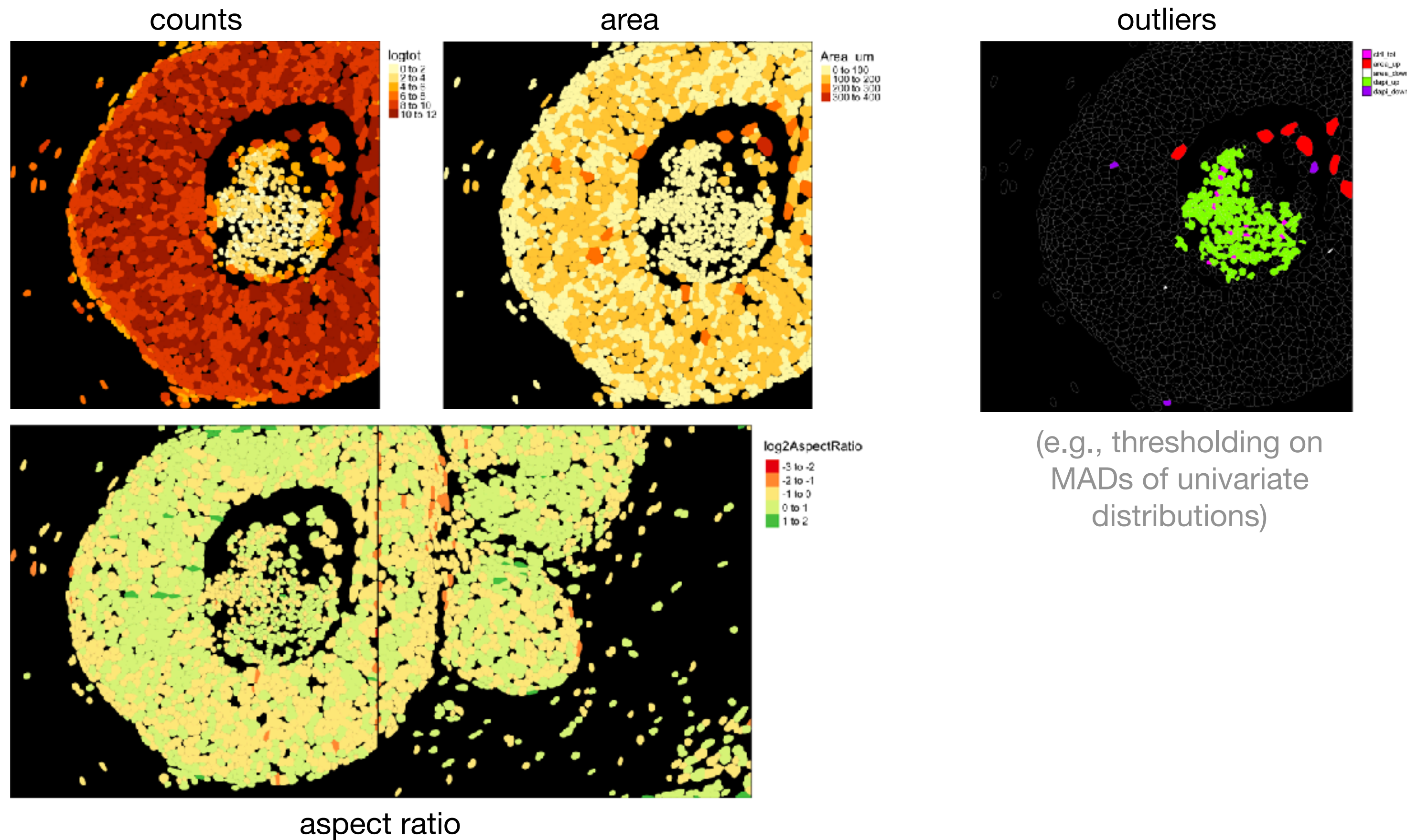


# SpaceTrooper proposes a *flag score* combining several metrics





# SpaceTrooper proposes a *flag score* combining several metrics



(e.g., thresholding on MADs of univariate distributions)

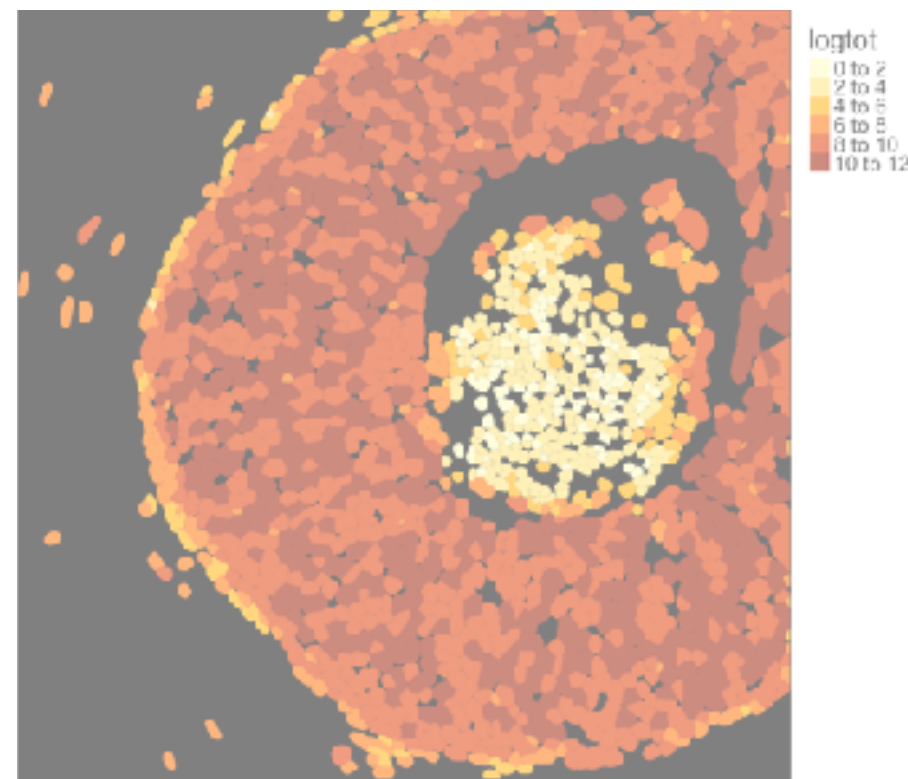


# SpaceTrooper proposes a *flag* score combining several metrics

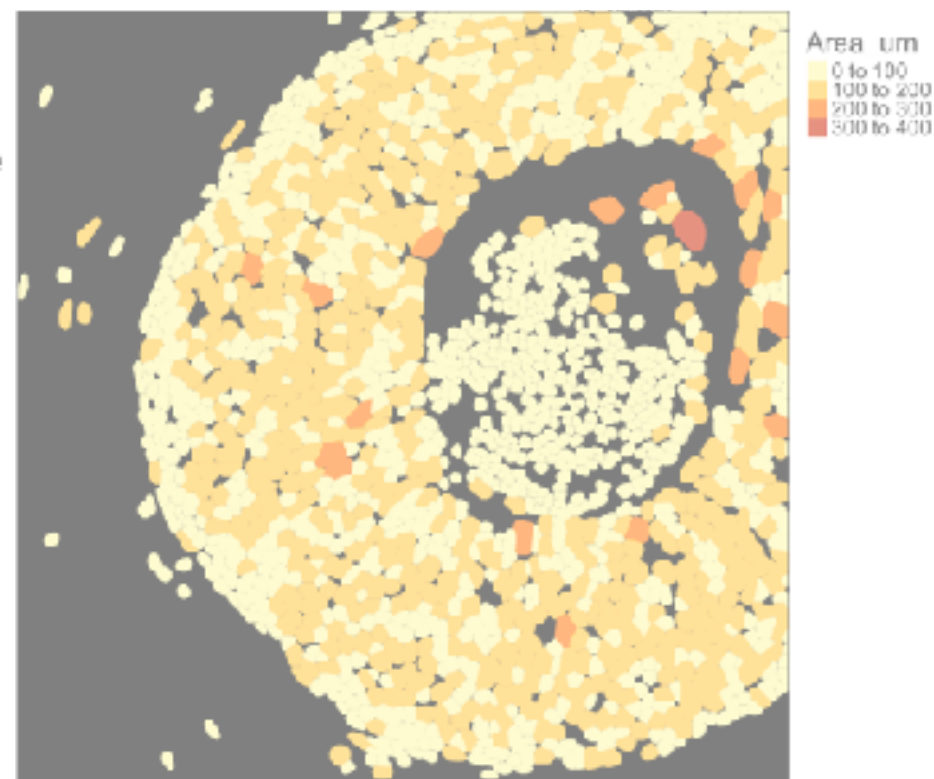


$$\text{logit}(F) = \beta_0 + \beta_1 \log \frac{\text{count}}{\text{area}} + \beta_2 |\log(\text{aspect ratio})| \cdot I_{\{d < \text{threshold}\}}$$

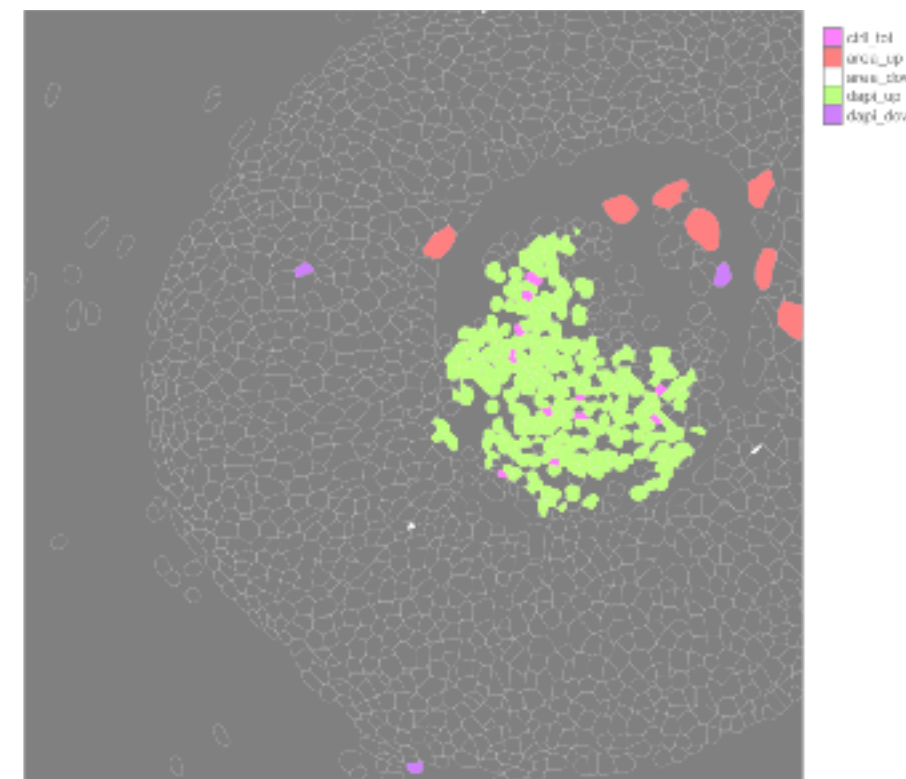
counts



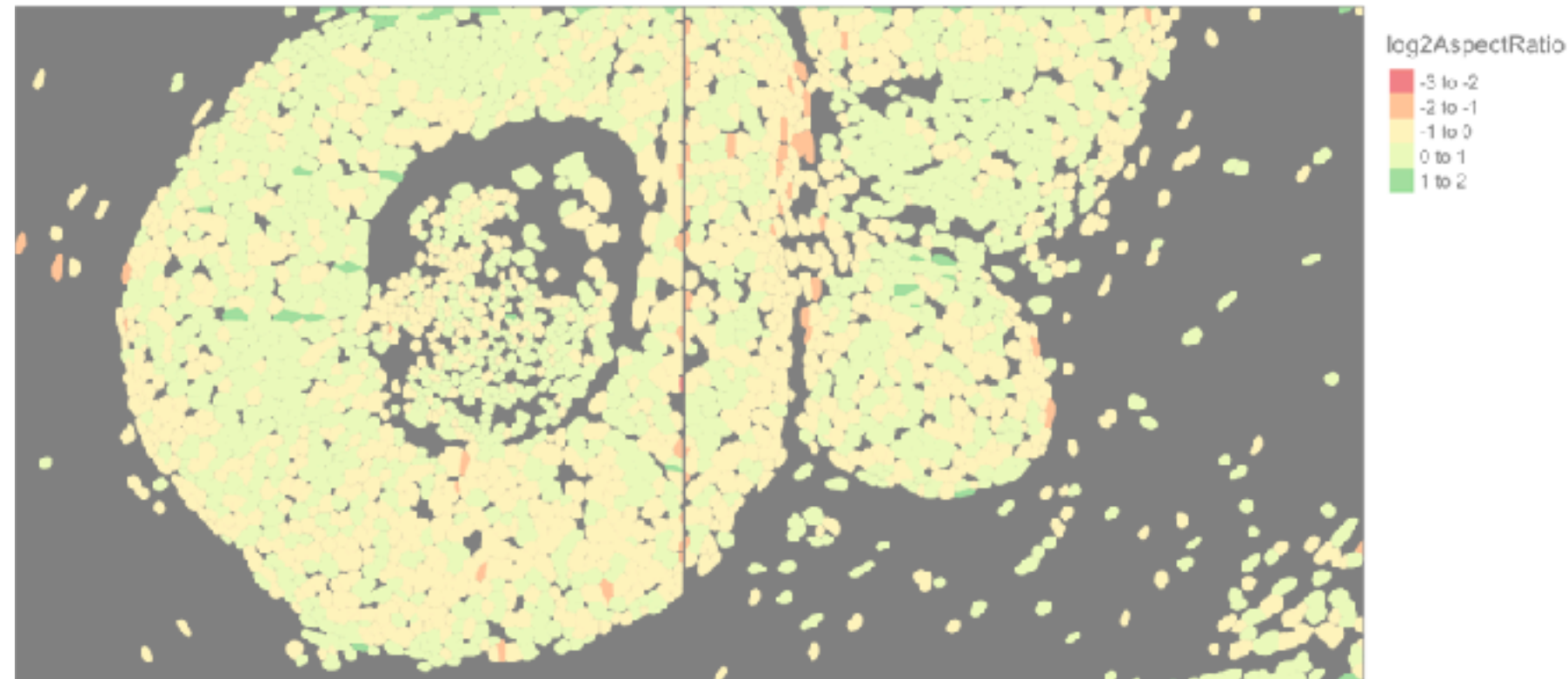
area



outliers



(e.g., thresholding on  
MADs of univariate  
distributions)



aspect ratio

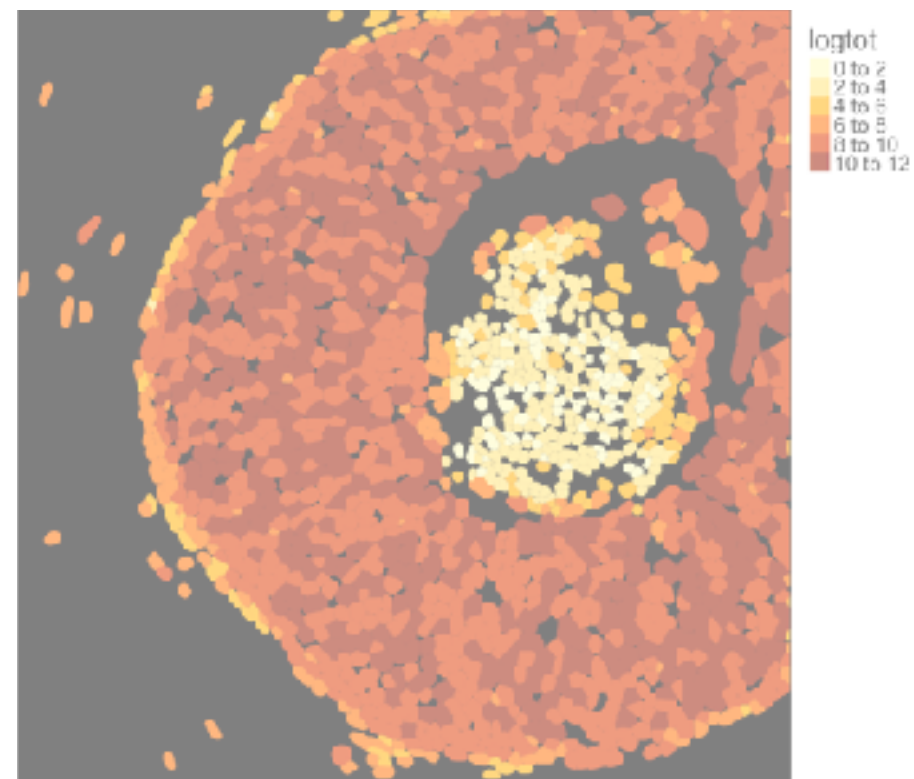


# SpaceTrooper proposes a *flag score* combining several metrics

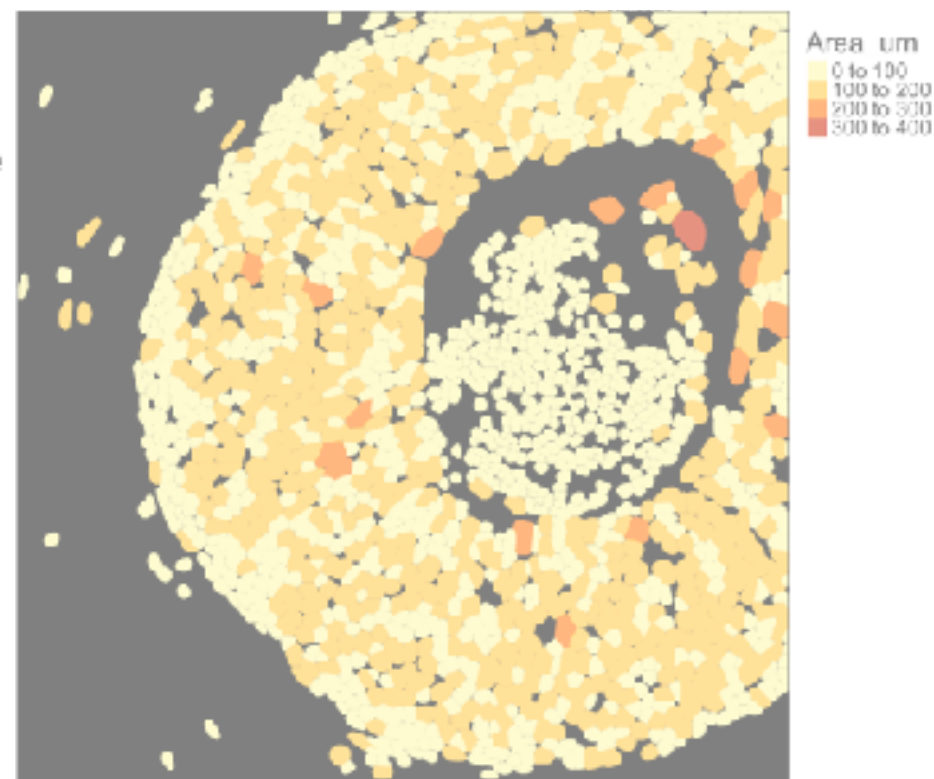


$$\text{logit}(F) = \beta_0 + \beta_1 \log \frac{\text{count}}{\text{area}} + \beta_2 |\log(\text{aspect ratio})| \cdot I_{\{d < \text{threshold}\}}$$

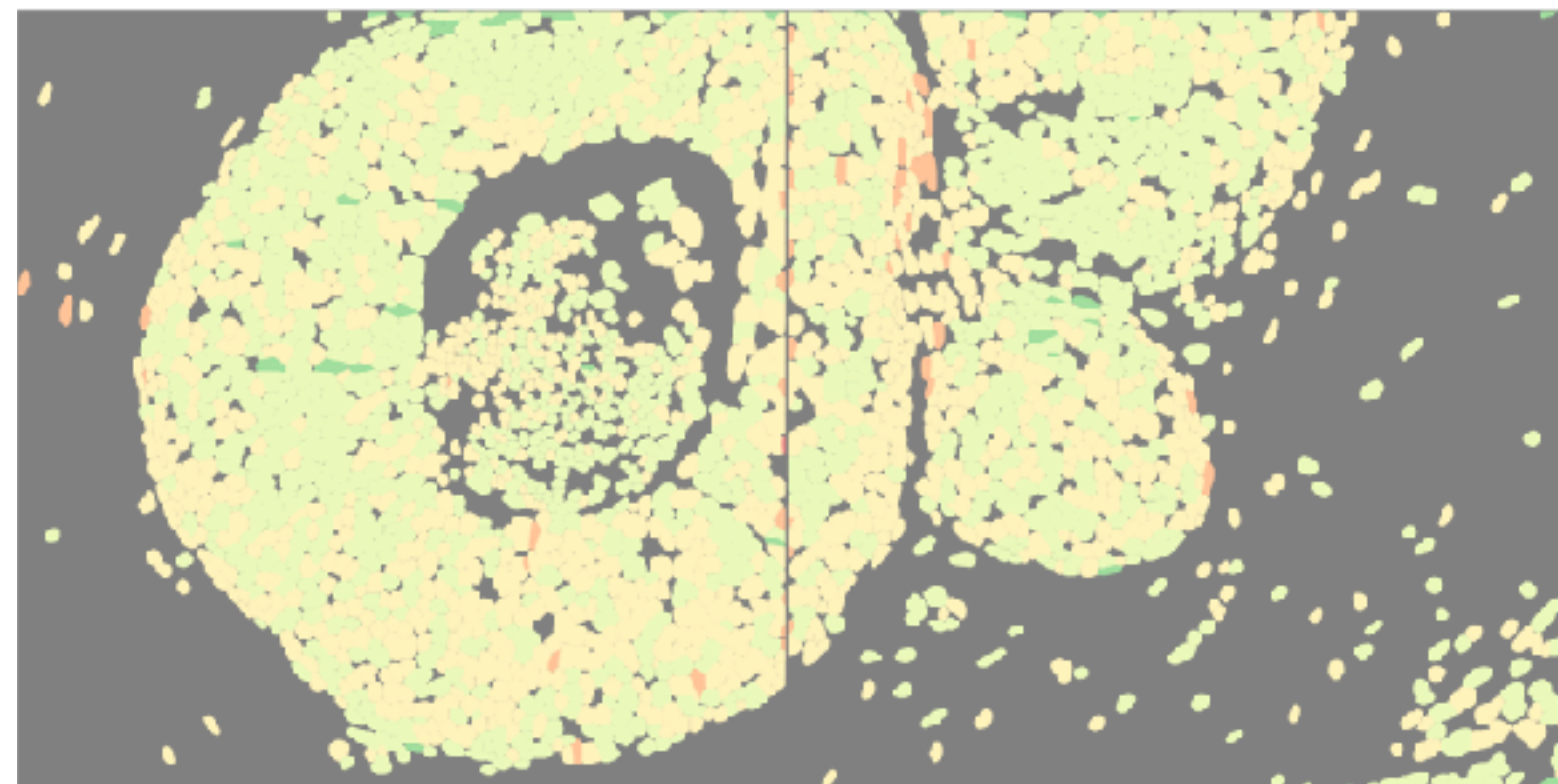
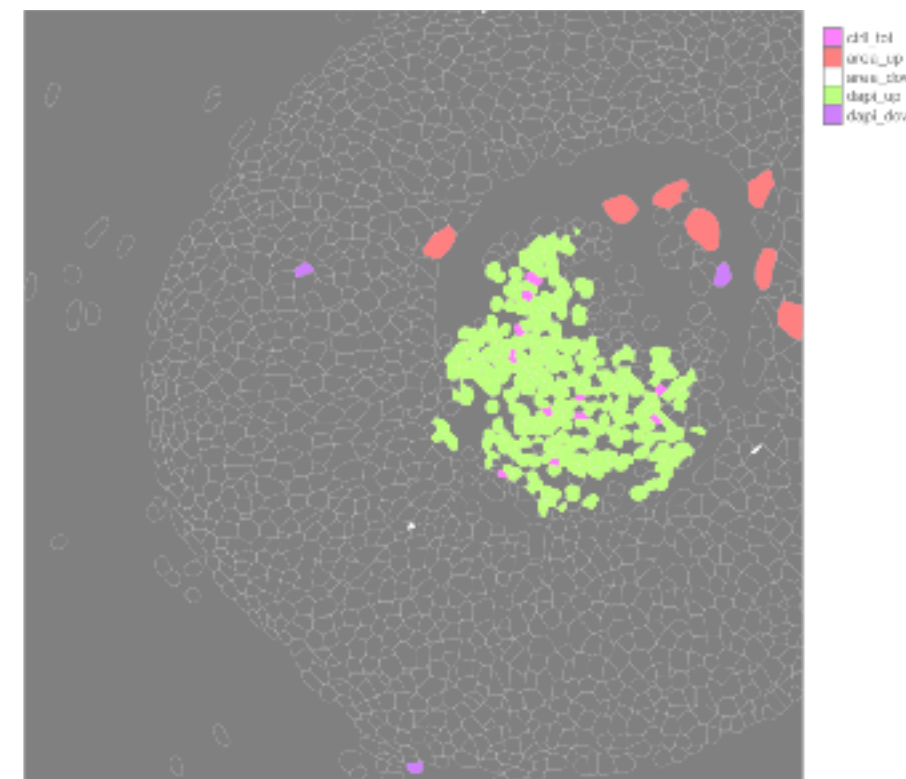
counts



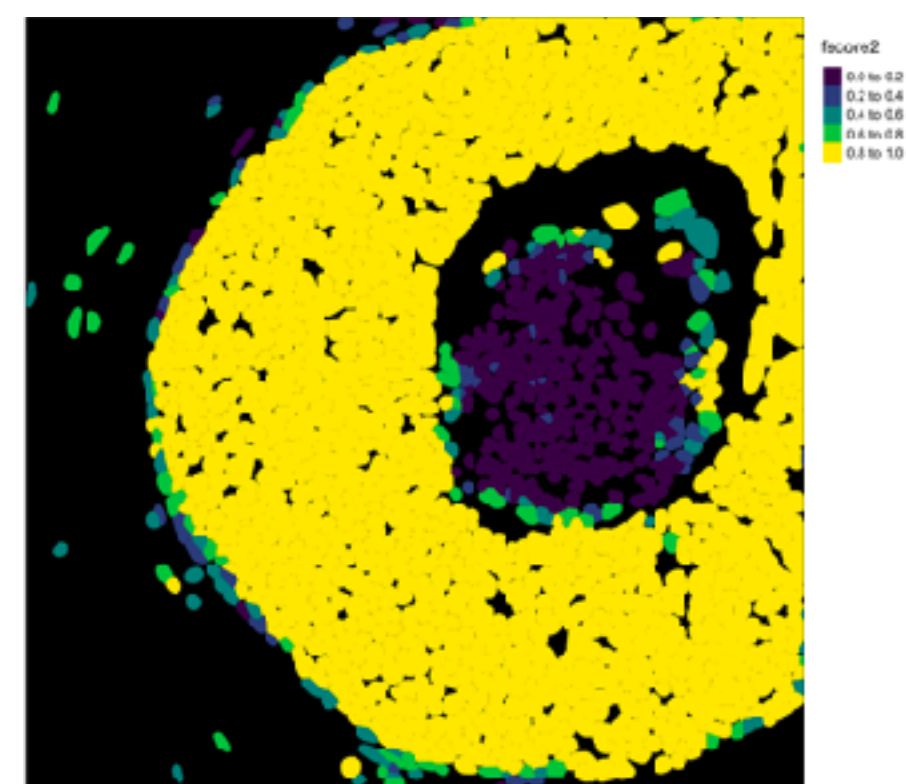
area



outliers



aspect ratio



flag score

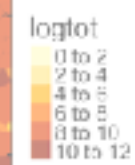
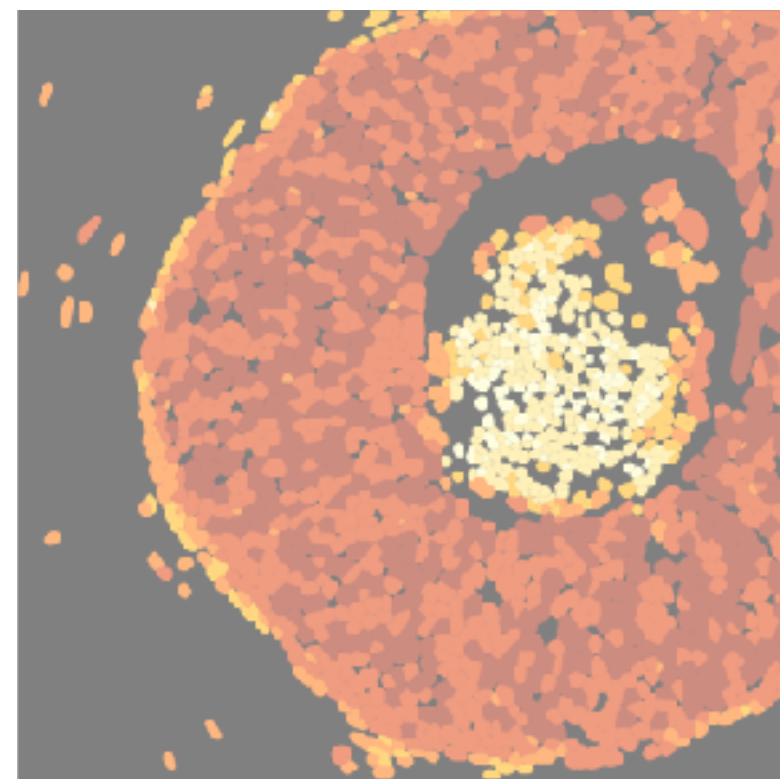


# SpaceTrooper proposes a *flag score* combining several metrics

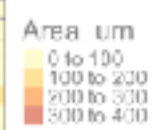
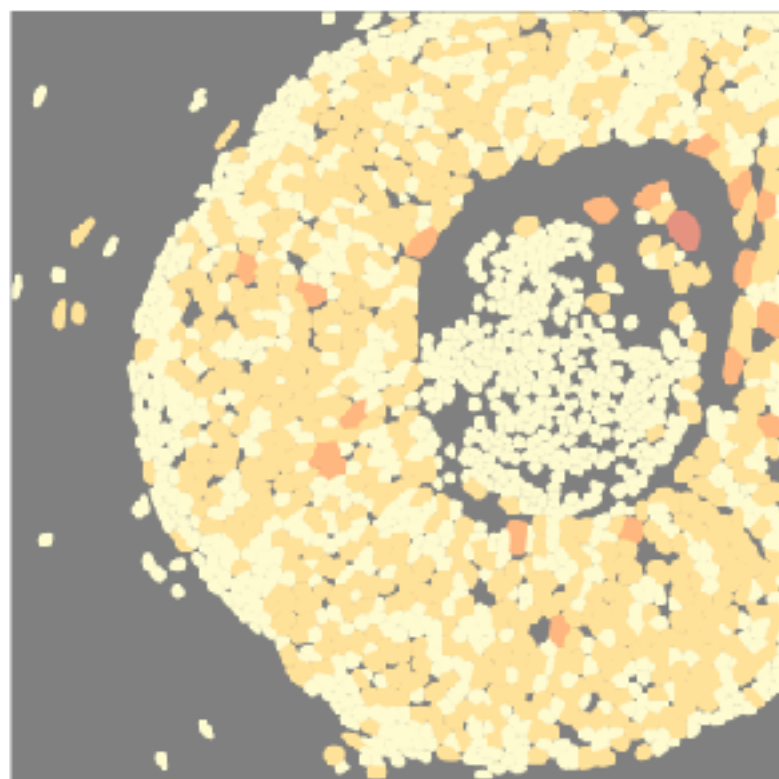


$$\text{logit}(F) = \beta_0 + \beta_1 \log \frac{\text{count}}{\text{area}} + \beta_2 |\log(\text{aspect ratio})| \cdot I_{\{d < \text{threshold}\}}$$

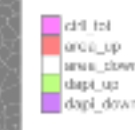
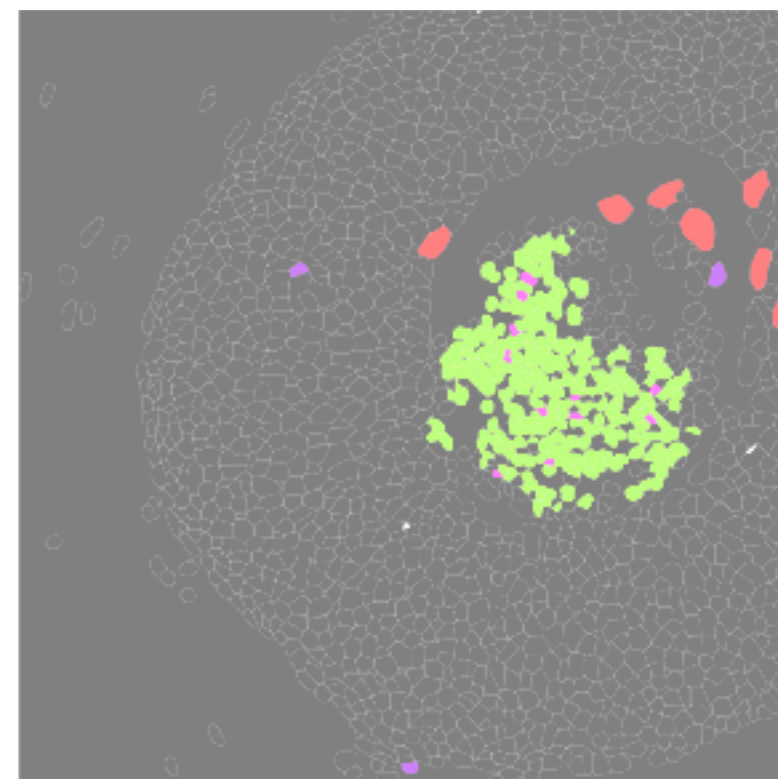
counts



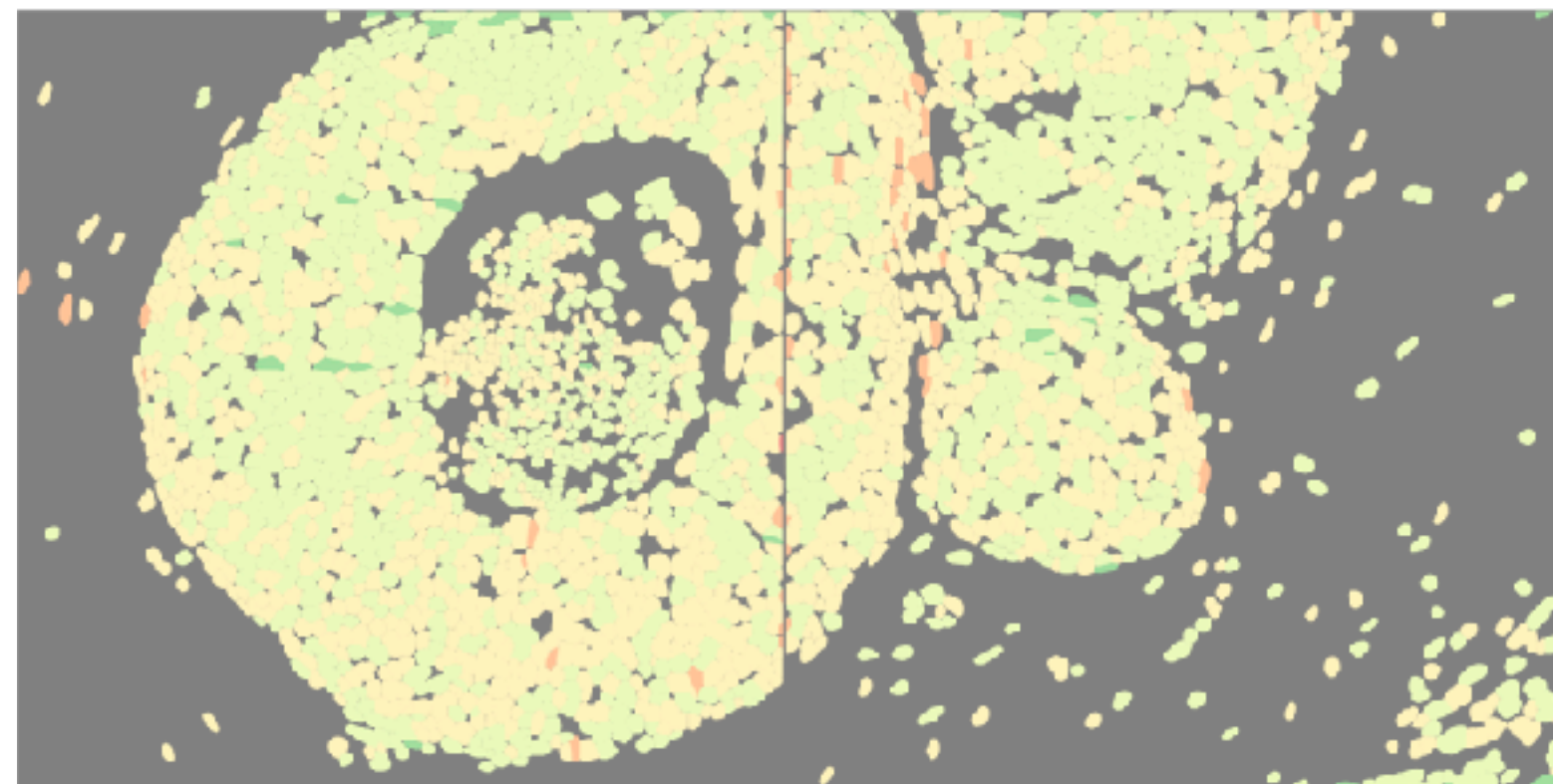
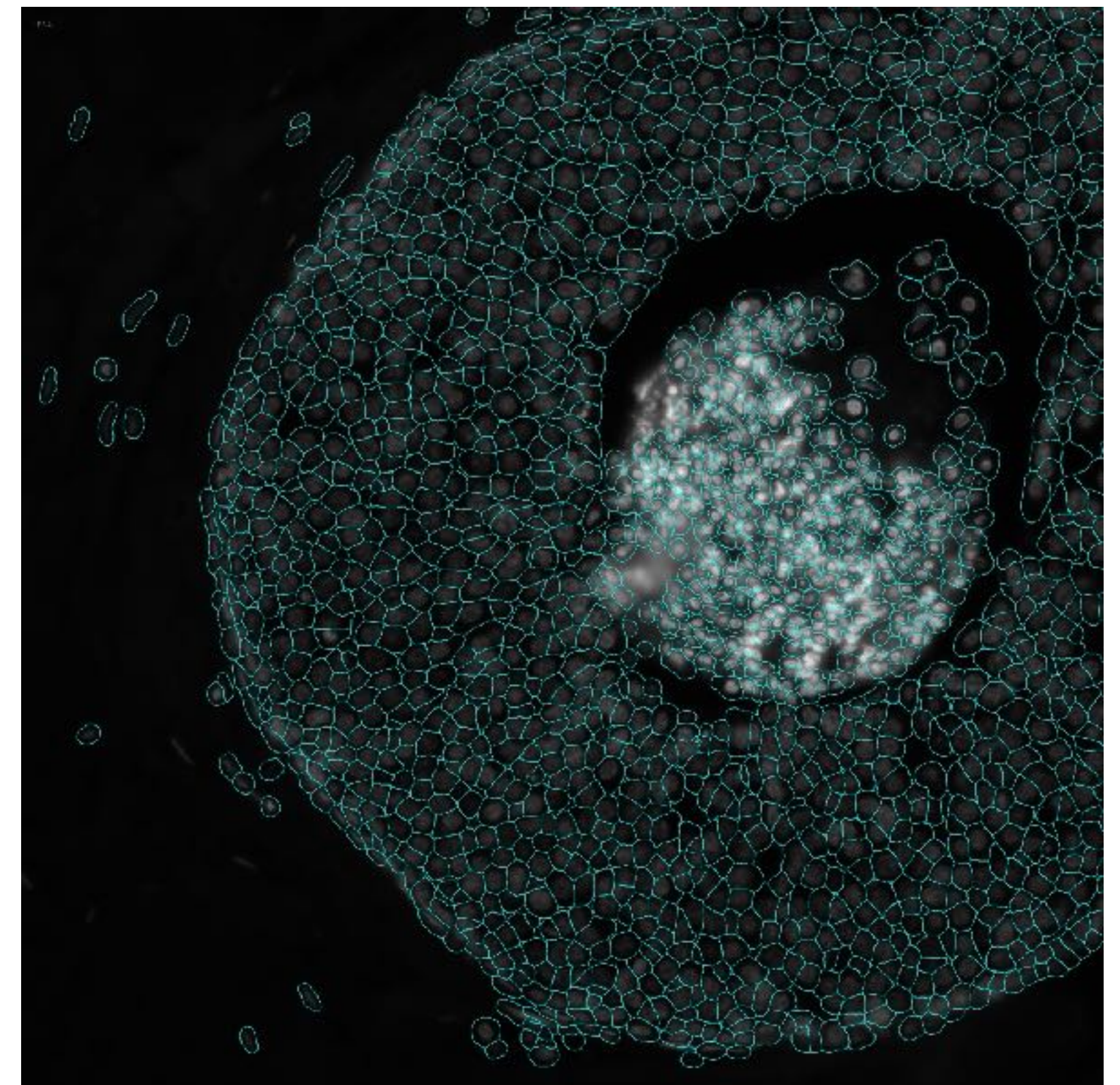
area



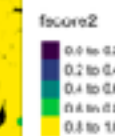
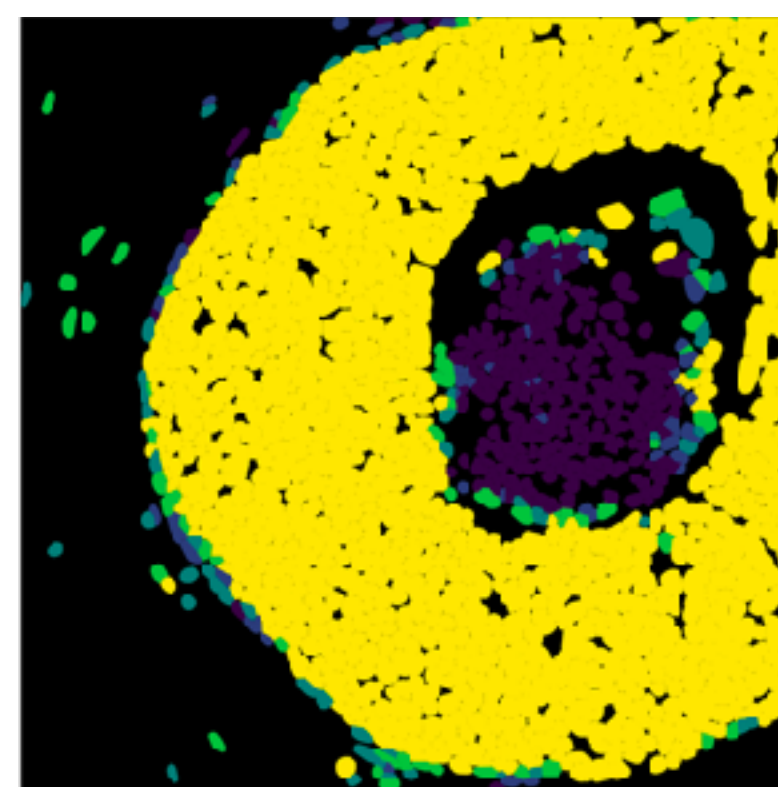
outliers



DAPI + segmentation



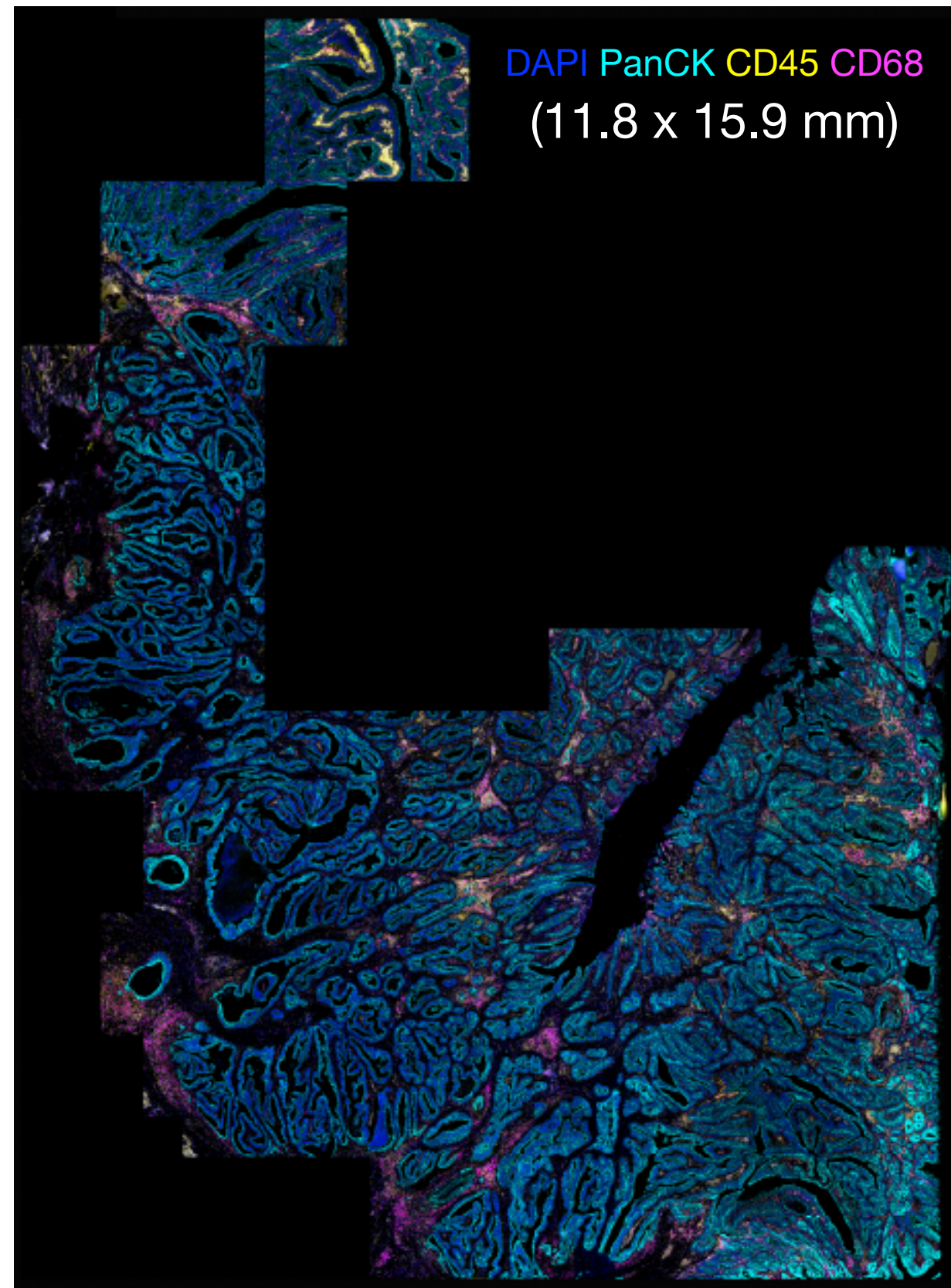
aspect ratio



flag score

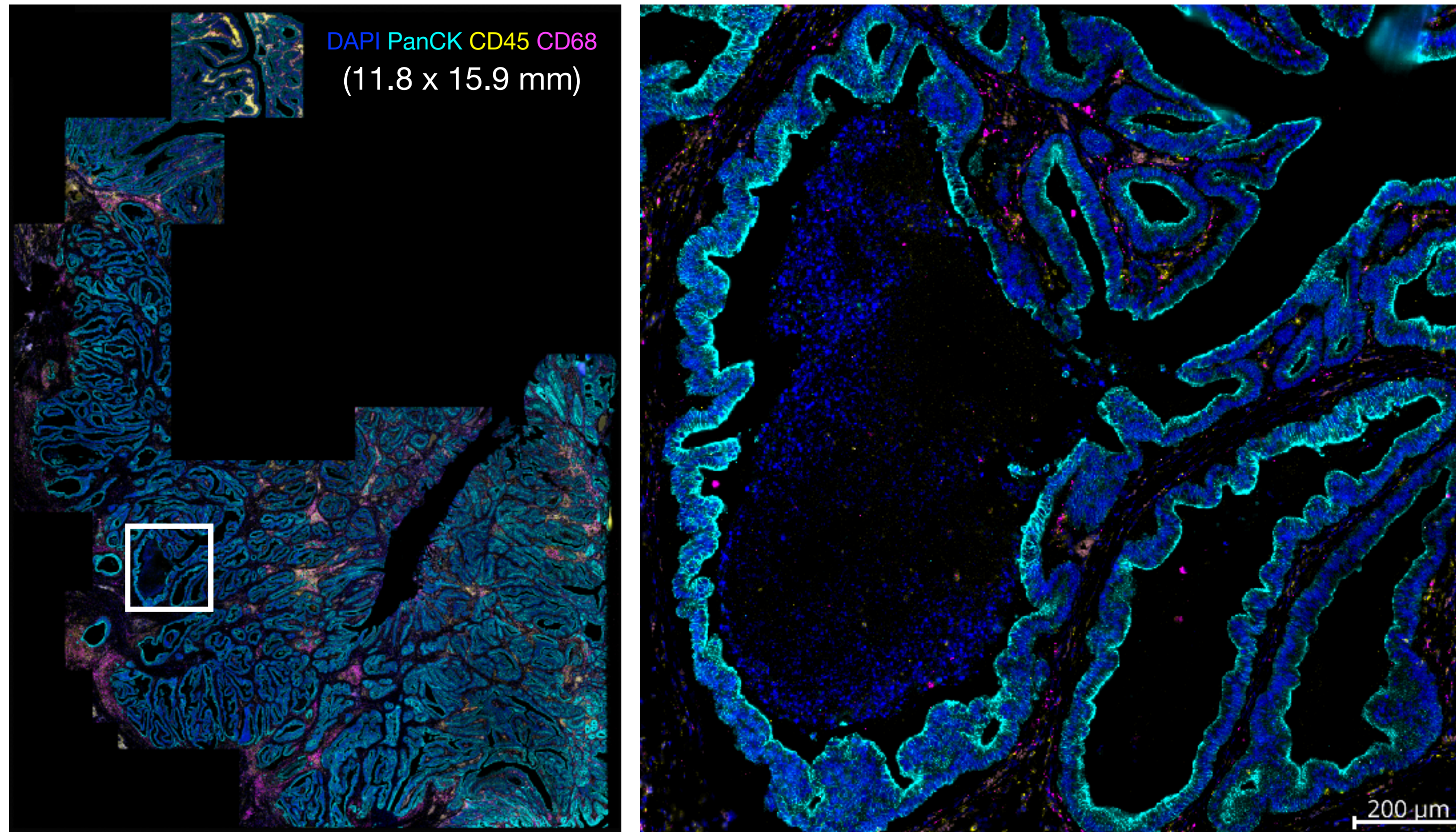


# debris is usually segmented, yet hard to distinguish



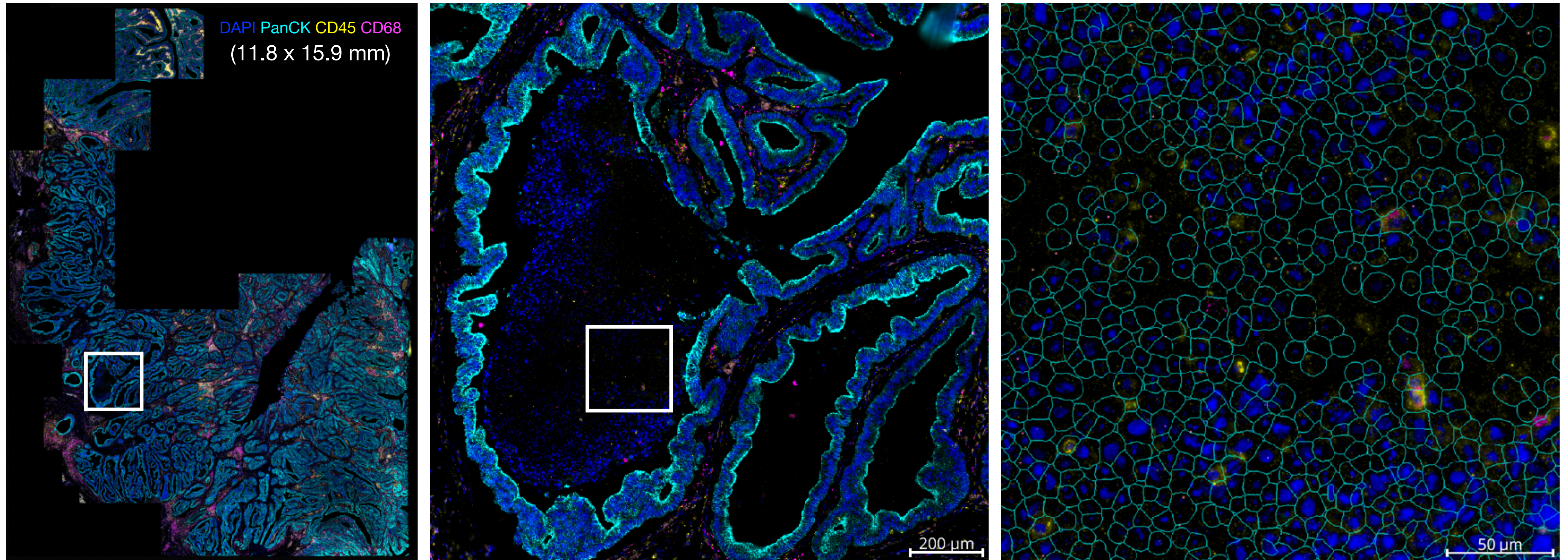


# debris is usually segmented, yet hard to distinguish



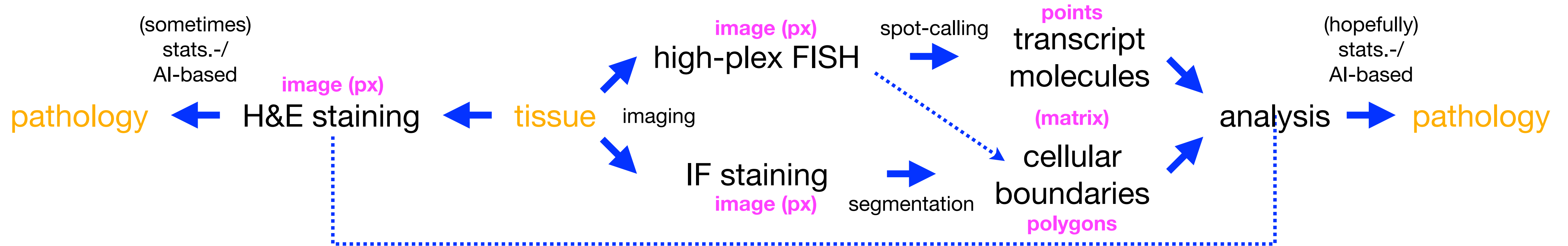


# debris is usually segmented, yet hard to distinguish



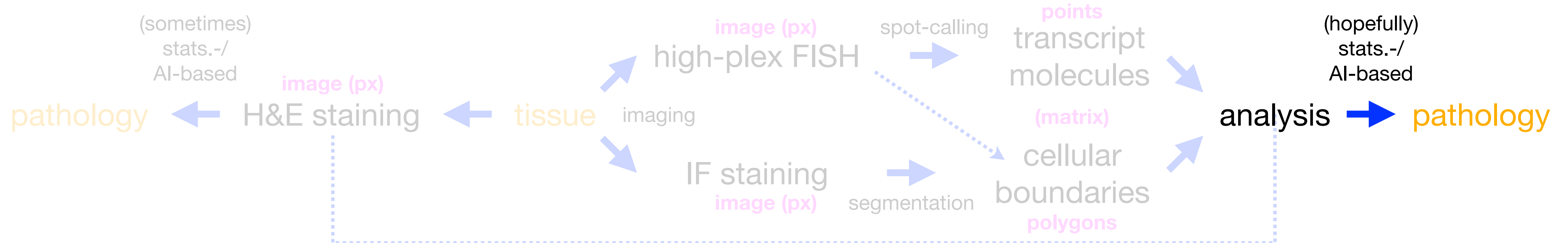


# library size normalization on non-representative panels introduces biases





# library size normalization on non-representative panels introduces biases



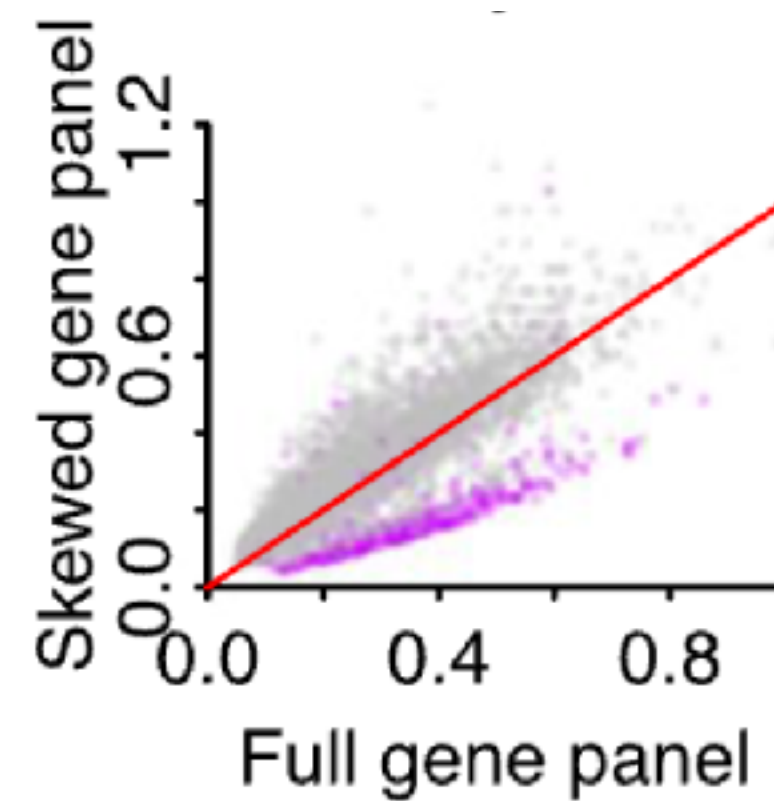
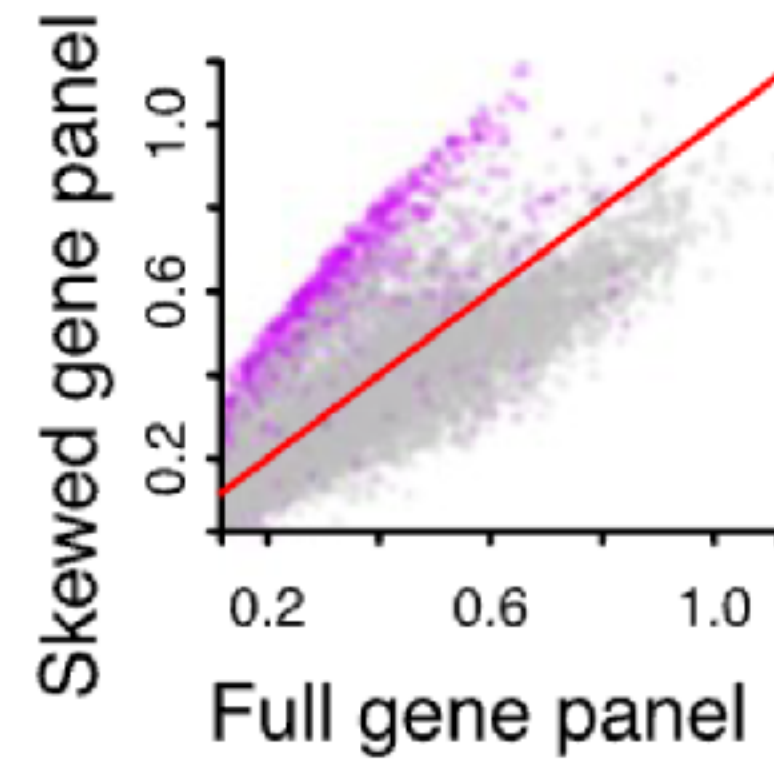


# library size normalization on non-representative panels introduces biases

- **scaling factors** for cells in R are systematically **larger**

100-gene panel skewed towards some region R

- **gene expression** for cells in R are systematically **smaller**



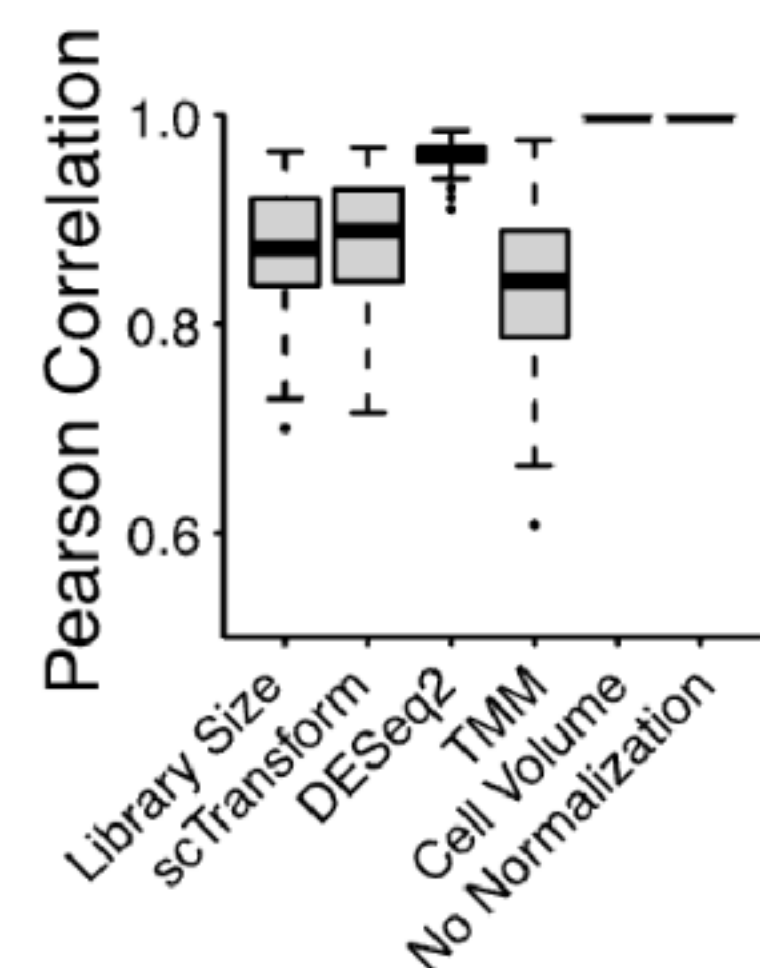
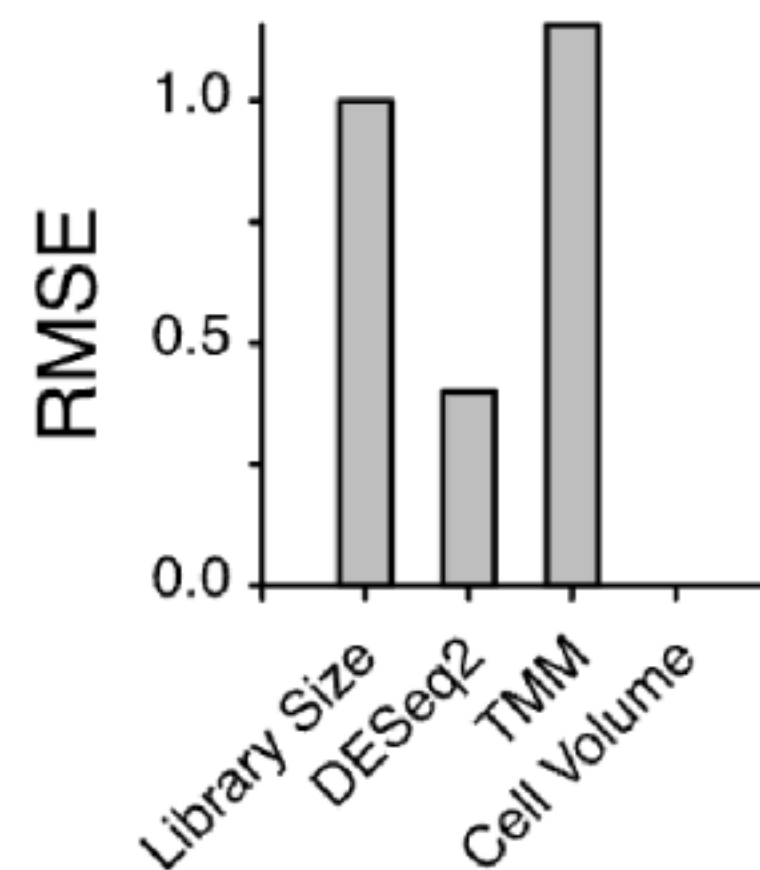
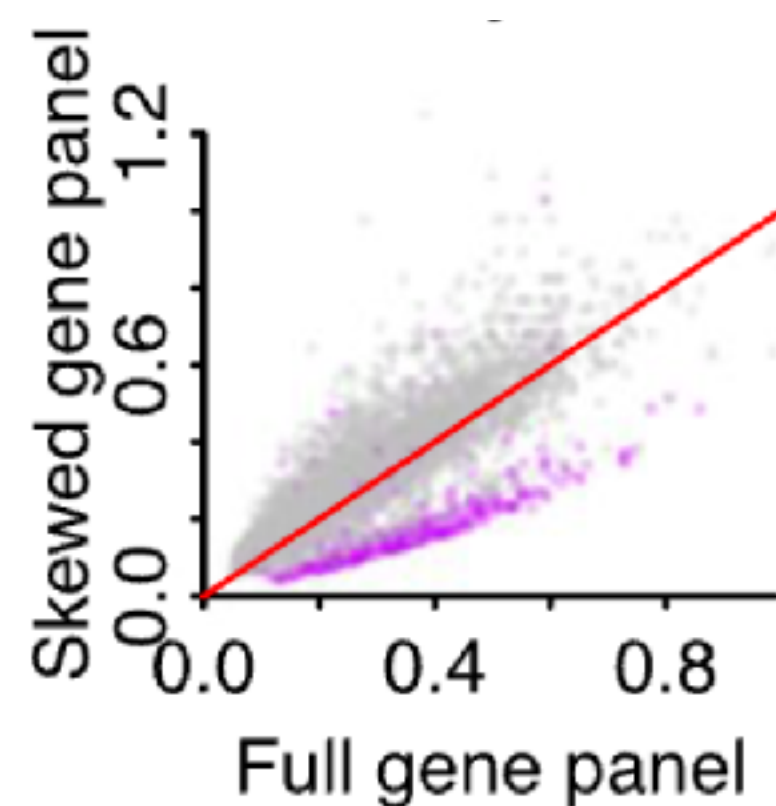
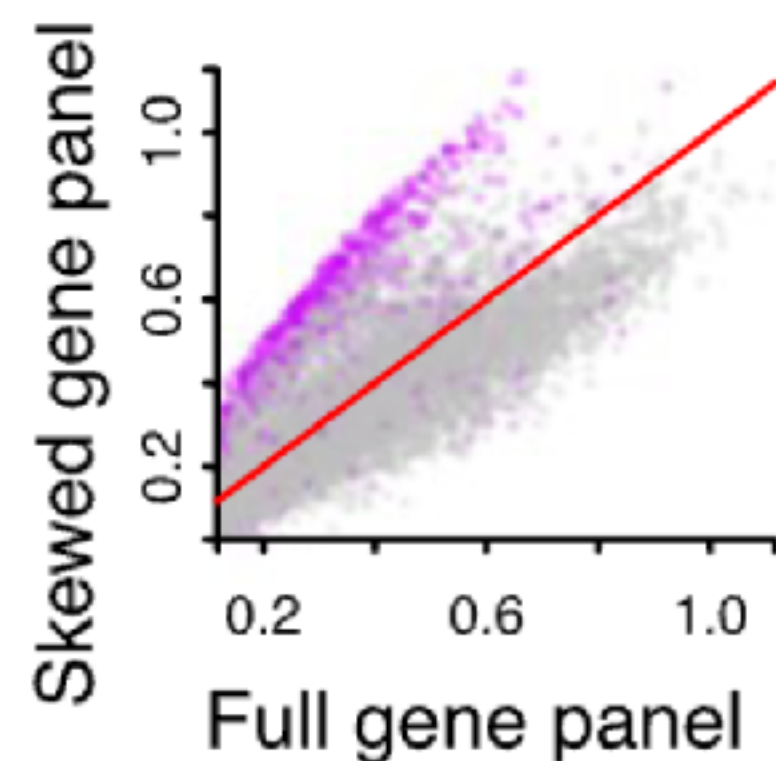


# library size normalization on non-representative panels introduces biases

- **scaling factors** for cells in R are systematically **larger**

100-gene panel skewed towards some region R

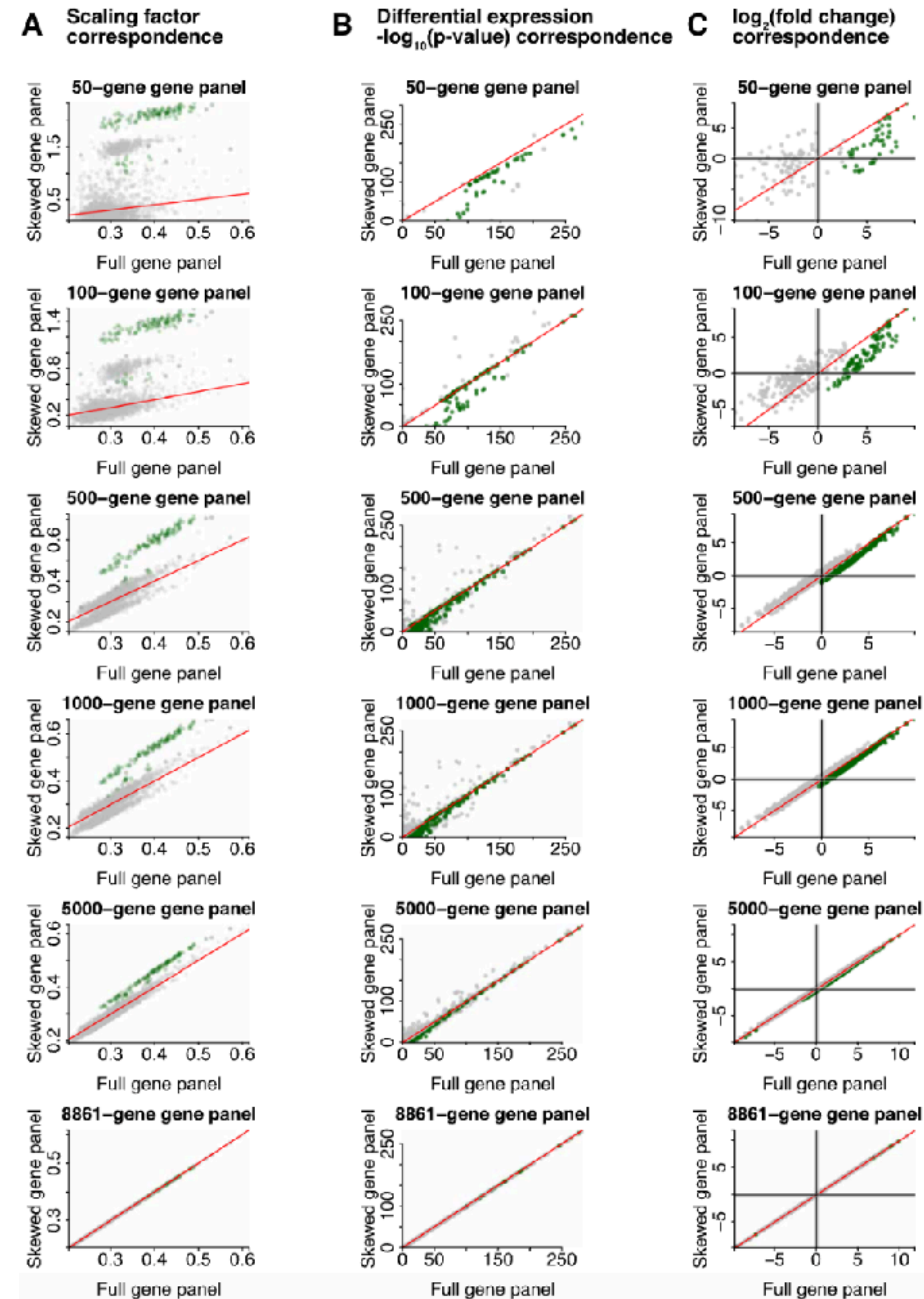
- **gene expression** for cells in R are systematically **smaller**



- systematic biases affect analyses to evaluate **differential gene expression & spatially variable genes**



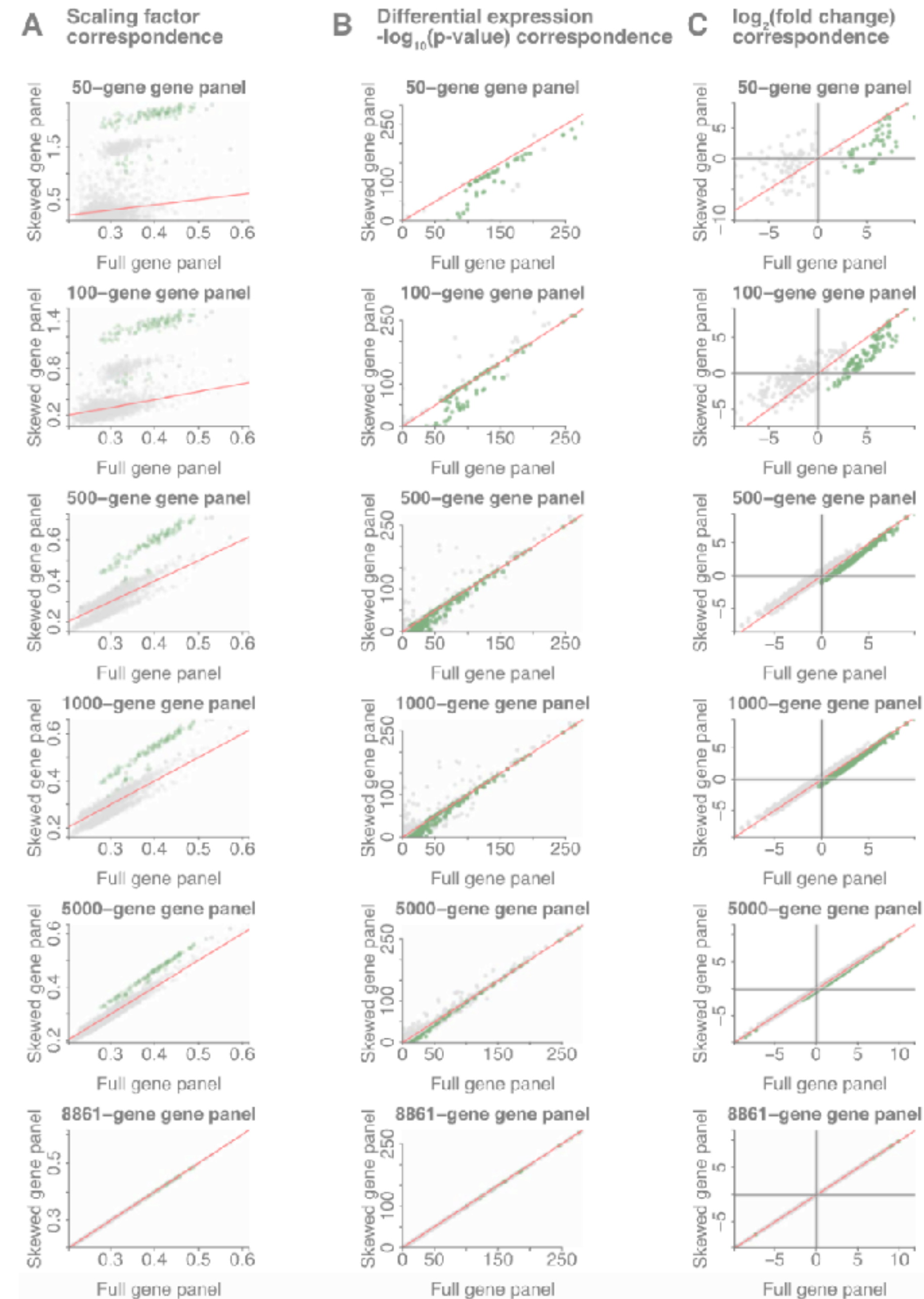
# larger/more representative panels help mitigate region-specific effects



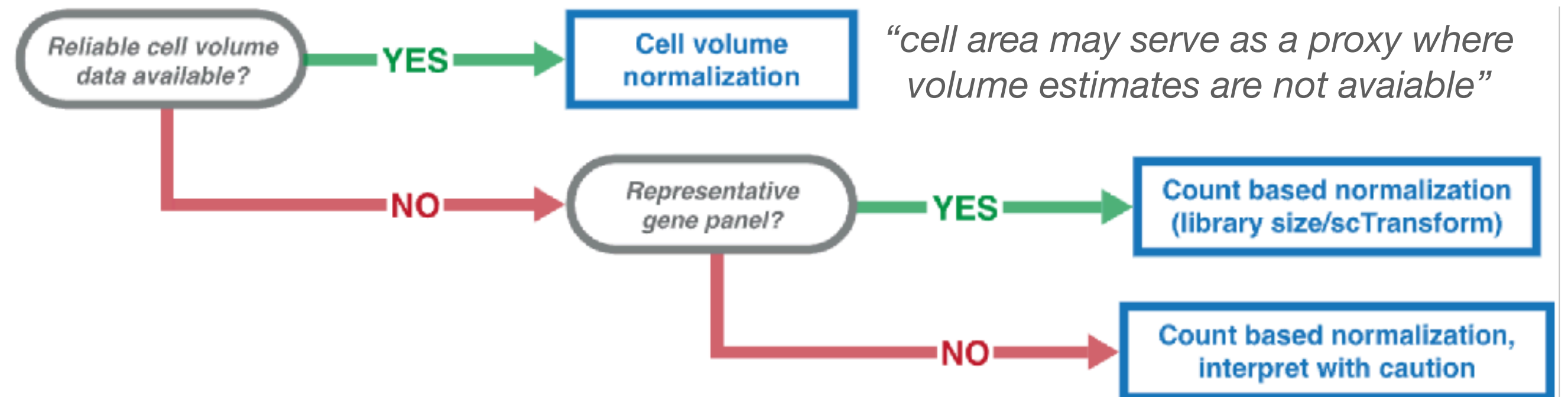
- skewed panels of 50...5,000 vs. all genes (simulated based on scRNA-seq data)
- differences are observed for skewed panels of all sizes, but their extent decreases as panel size increases



# larger/more representative panels help mitigate region-specific effects

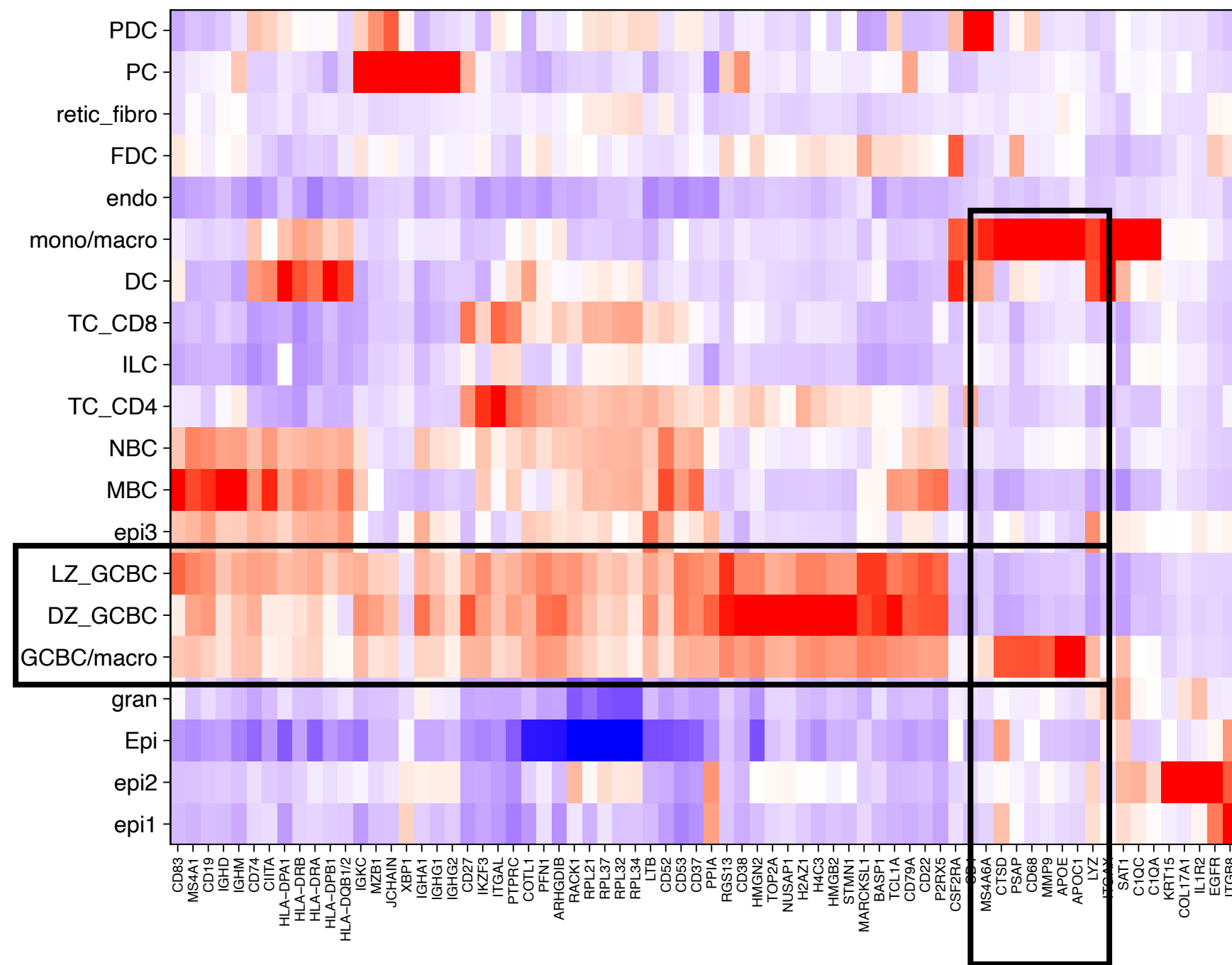


- skewed panels of 50...5,000 vs. all genes (simulated based on scRNA-seq data)
- differences are observed for skewed panels of all sizes, but their extent decreases as panel size increases



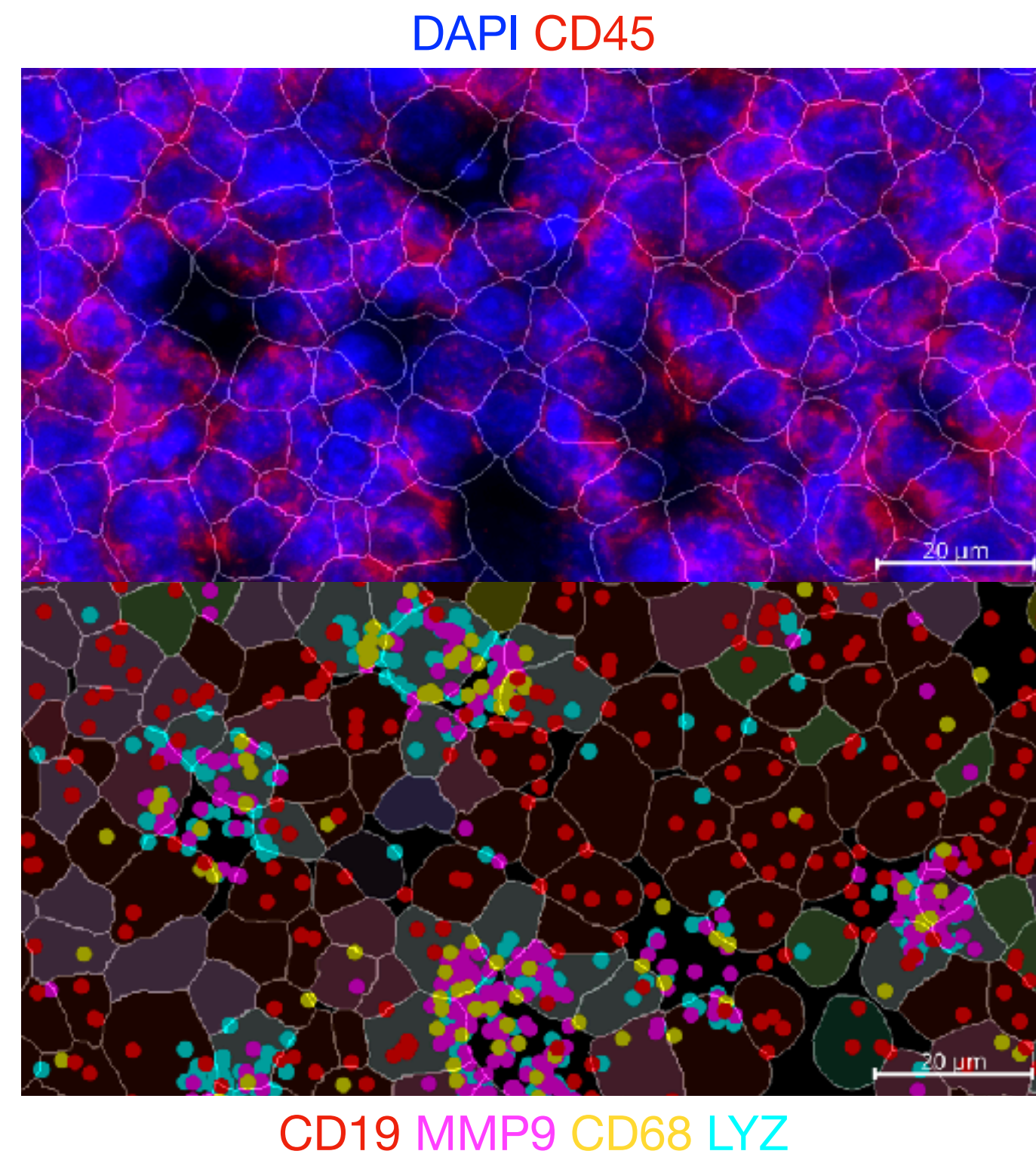
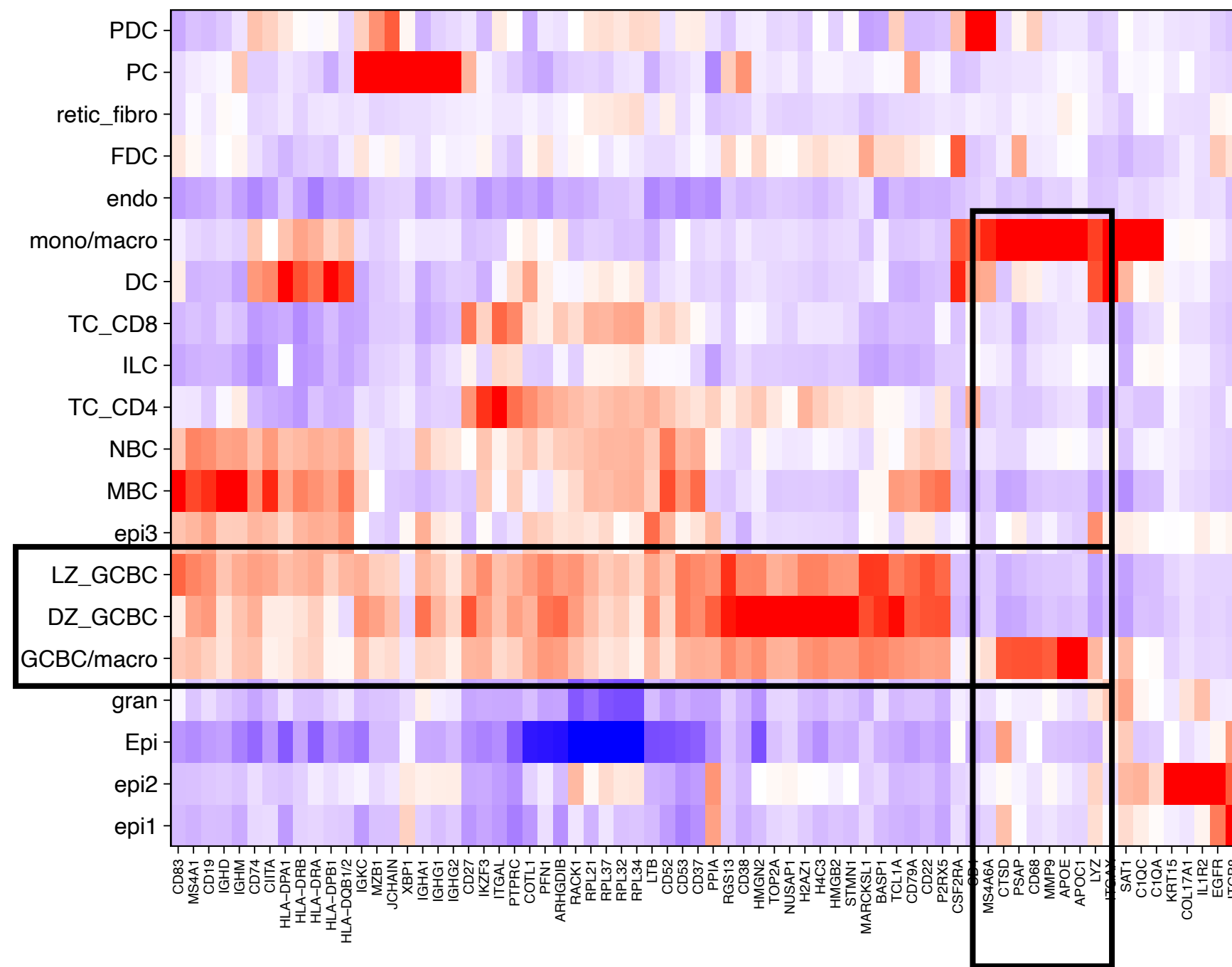


# CD68+ B cells — what's going on here?



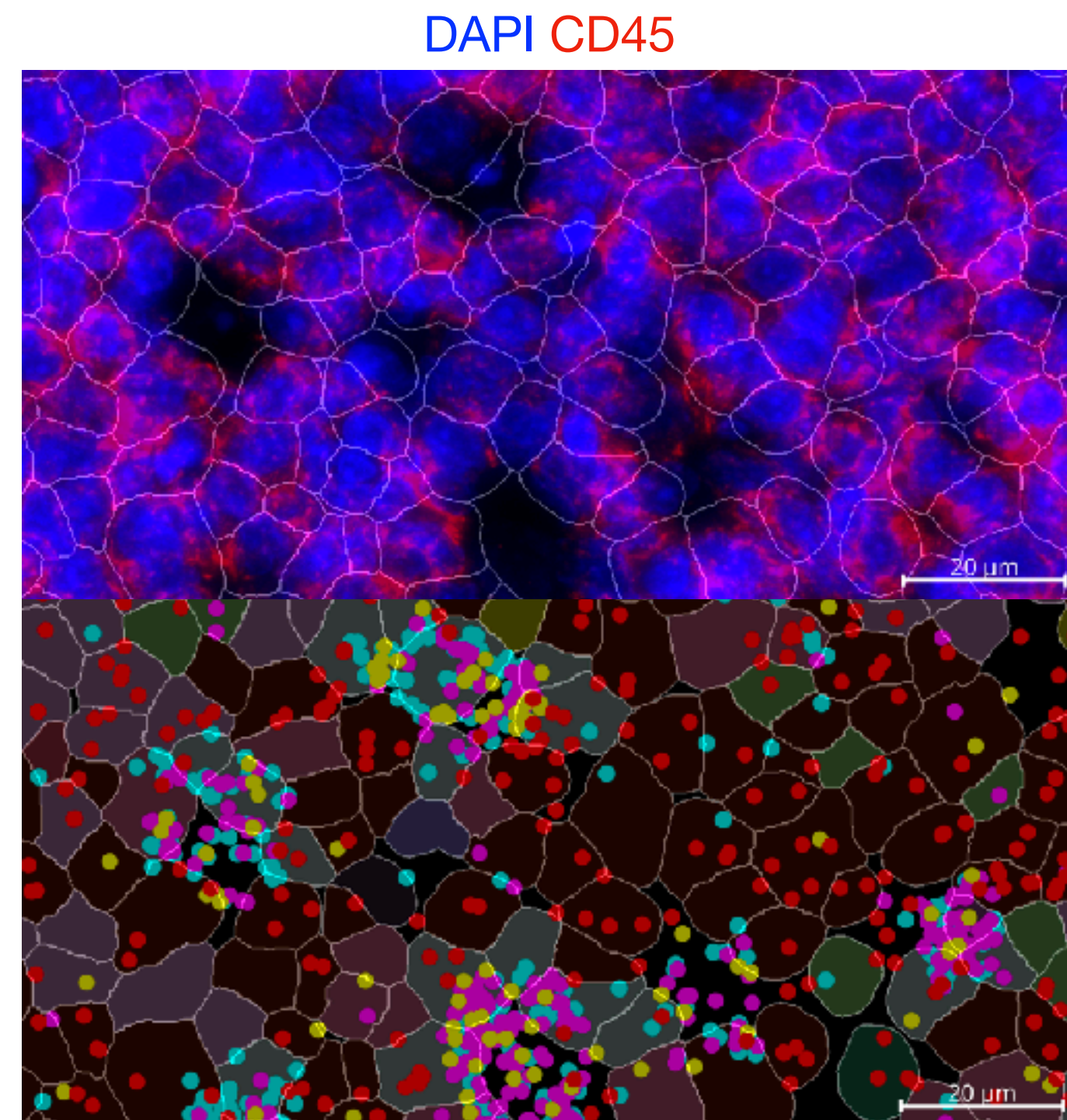
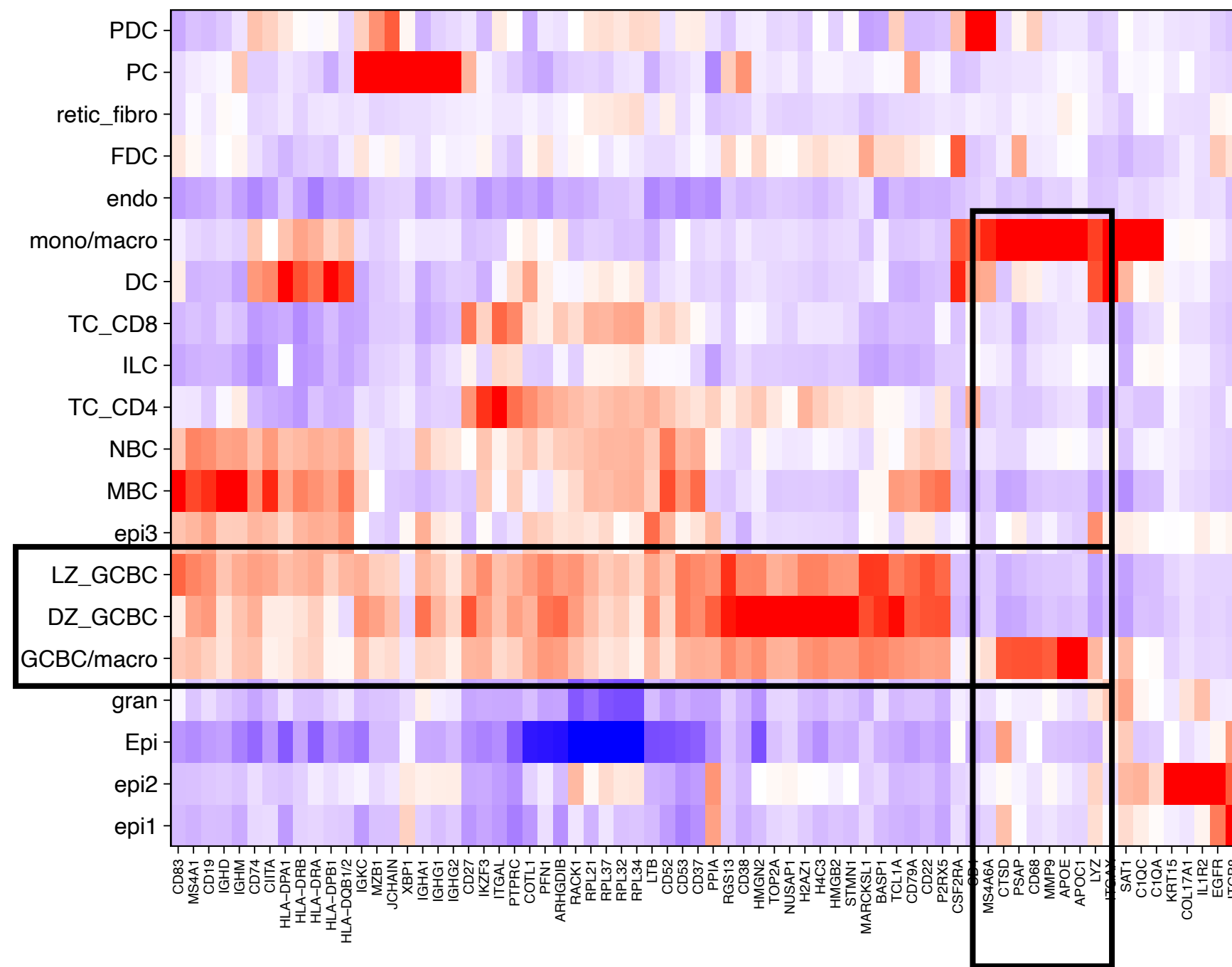


# CD68+ B cells — what's going on here?



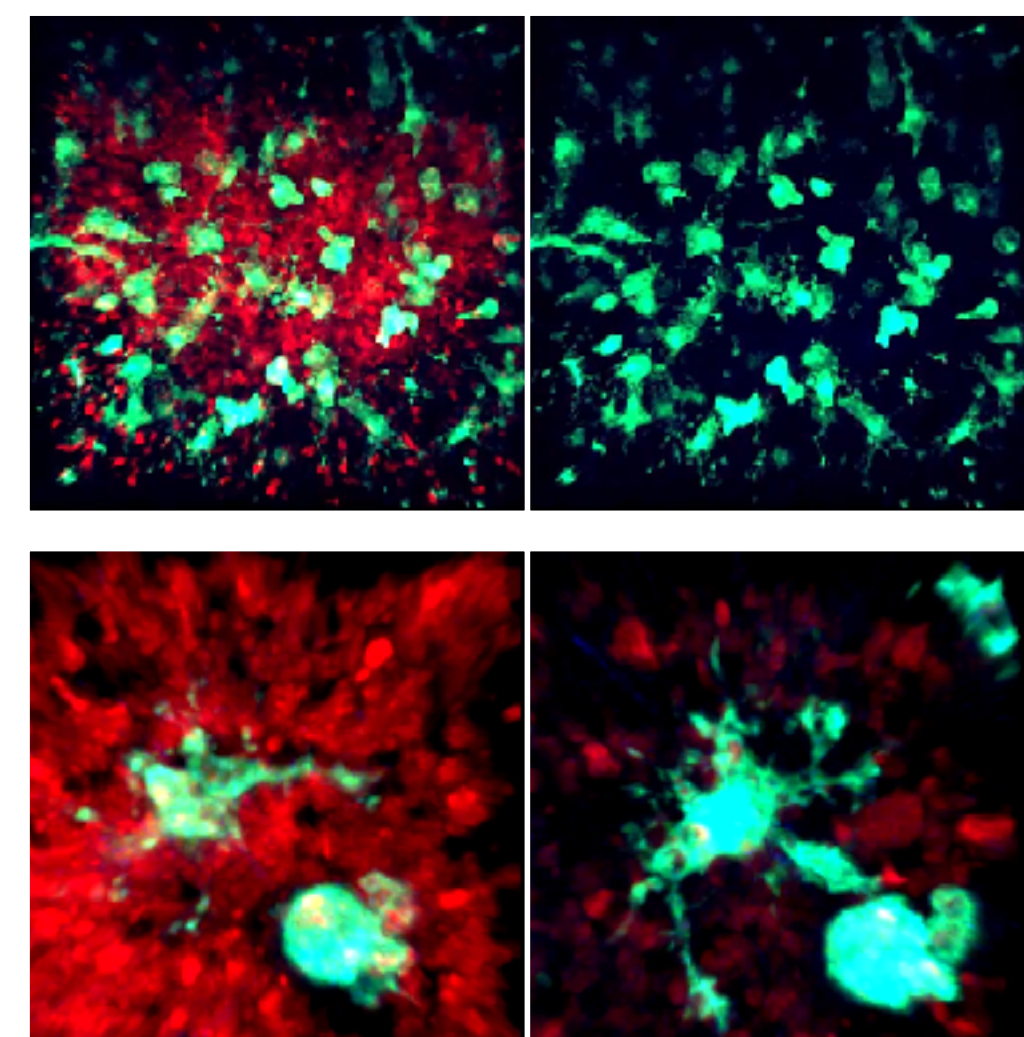


# CD68+ B cells — what's going on here?



CD19 MMP9 CD68 LYZ

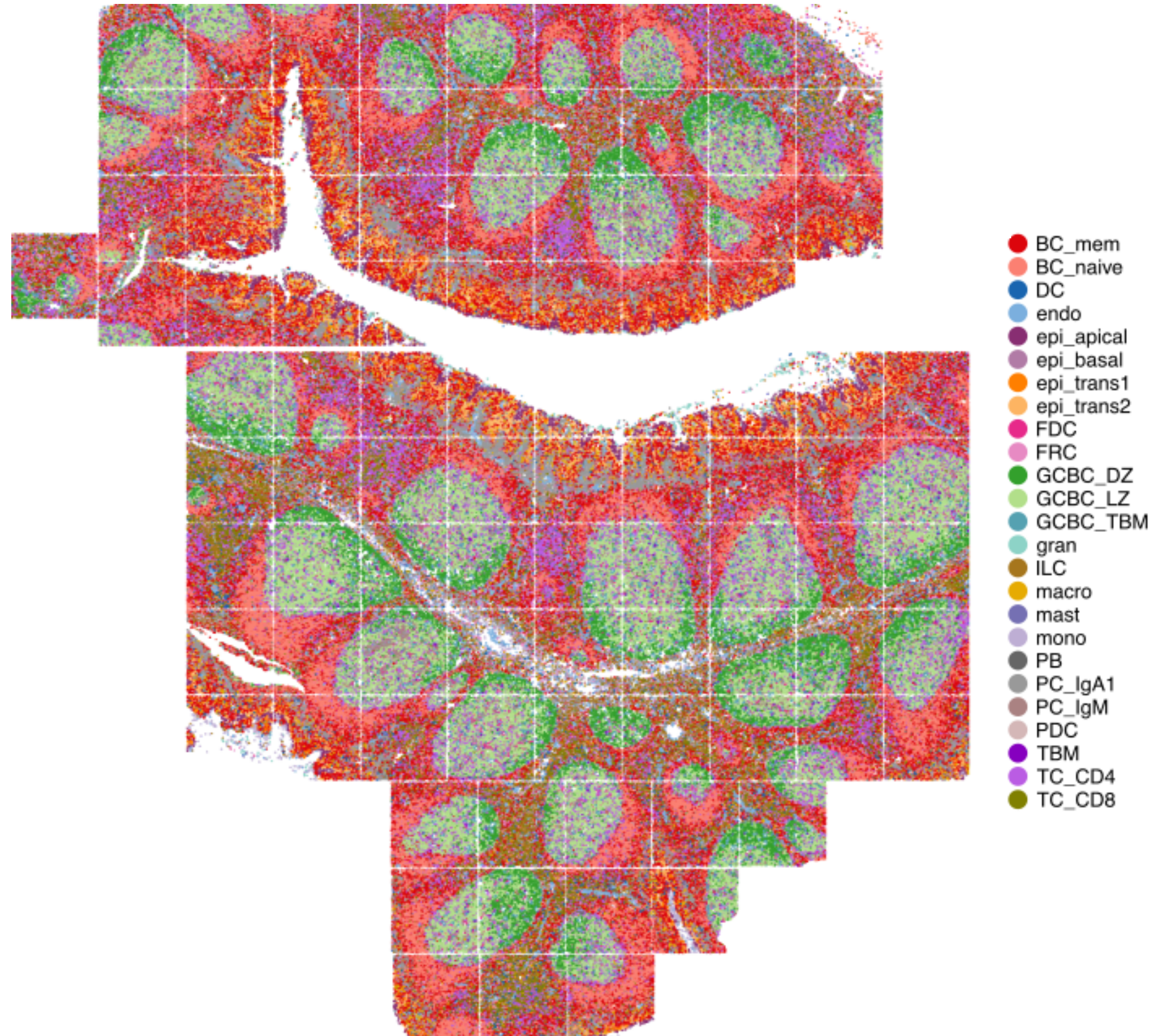
Article | January 27, 2023  
**Tingible body macrophages arise from lymph node-resident precursors and uptake B cells by dendrites**  
 Nela Staritz, Liel Stalor-Bark, Niklas Schwarz, Anshu Bandhyopadhyay, Michael Meyer-Hermann, Ziv Shulman  
 + Author and Article Information | Check for updates  
 J Exp Med (2023) 220 (4): e20222173 | <https://doi.org/10.1084/jem.20222173> | Article history



B cells macrophages

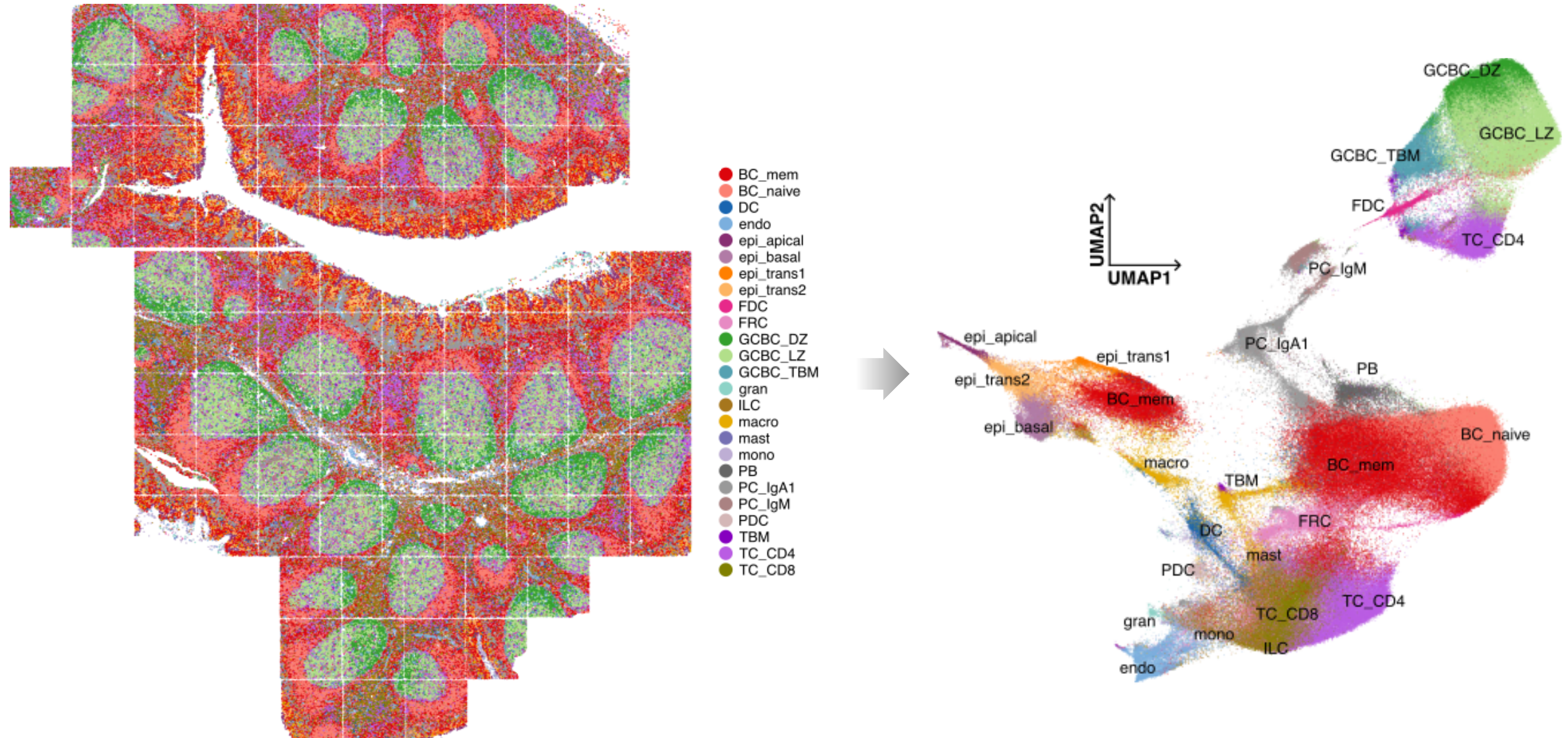


# spatial bleeding manifests in RNA counts, hence PCs, and UMAPs



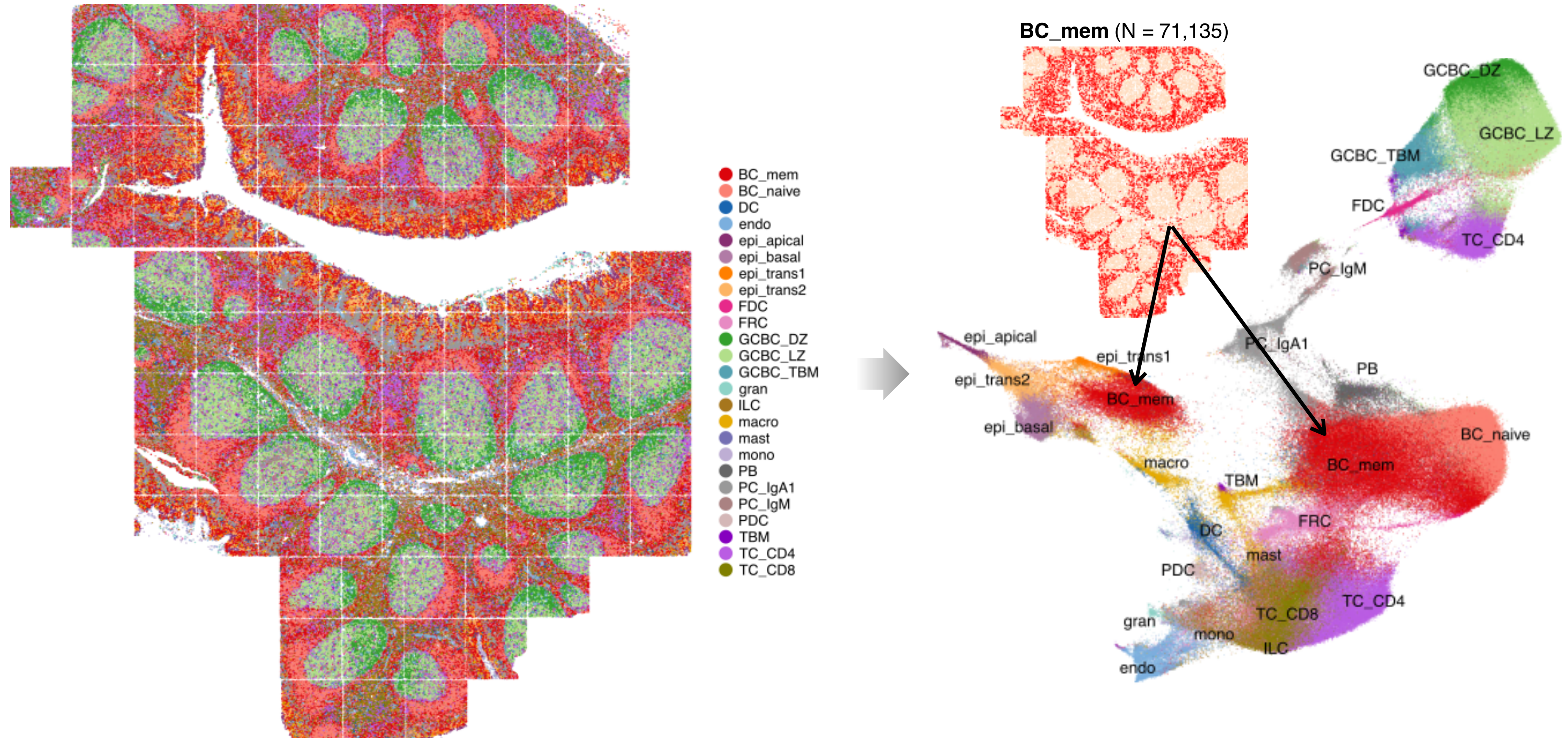


# spatial bleeding manifests in RNA counts, hence PCs, and UMAPs



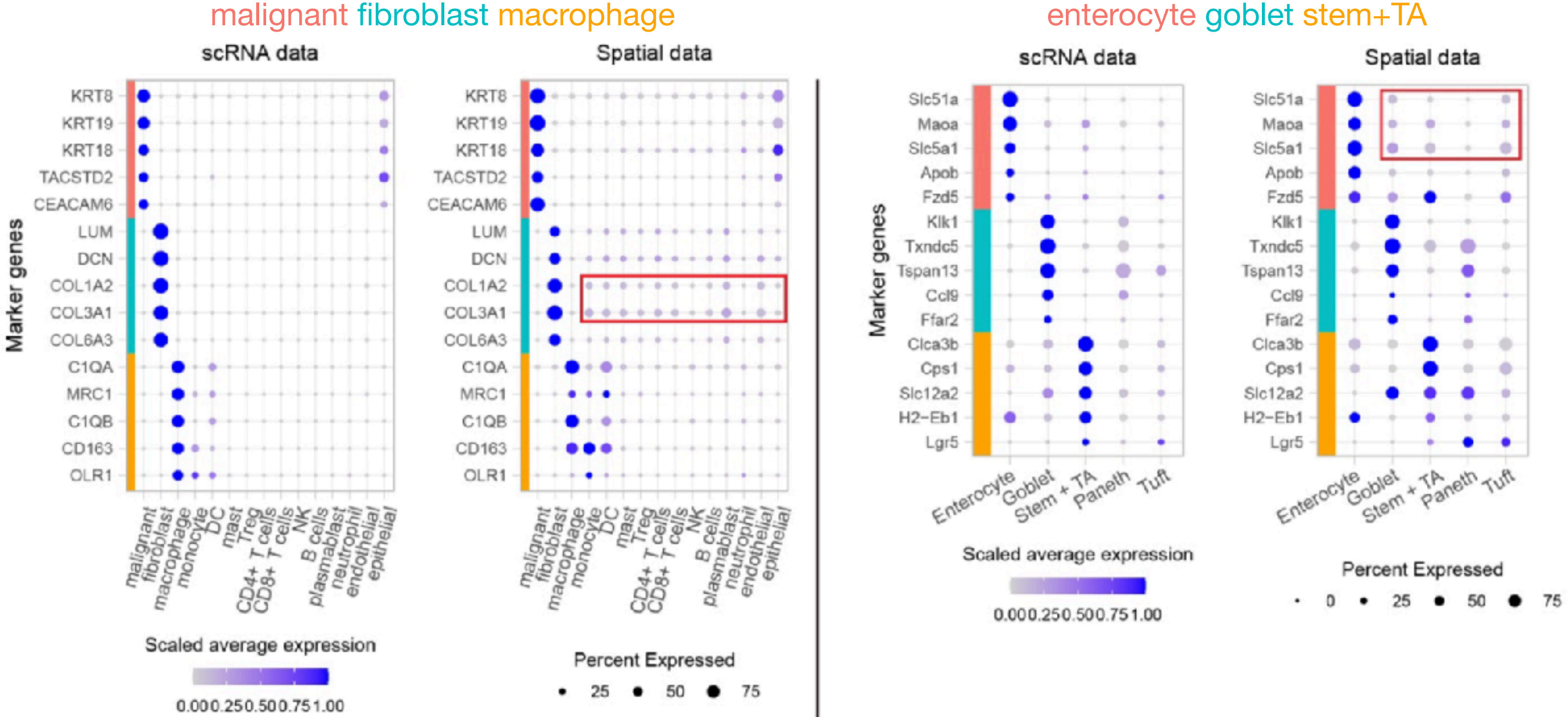


# spatial bleeding manifests in RNA counts, hence PCs, and UMAPs



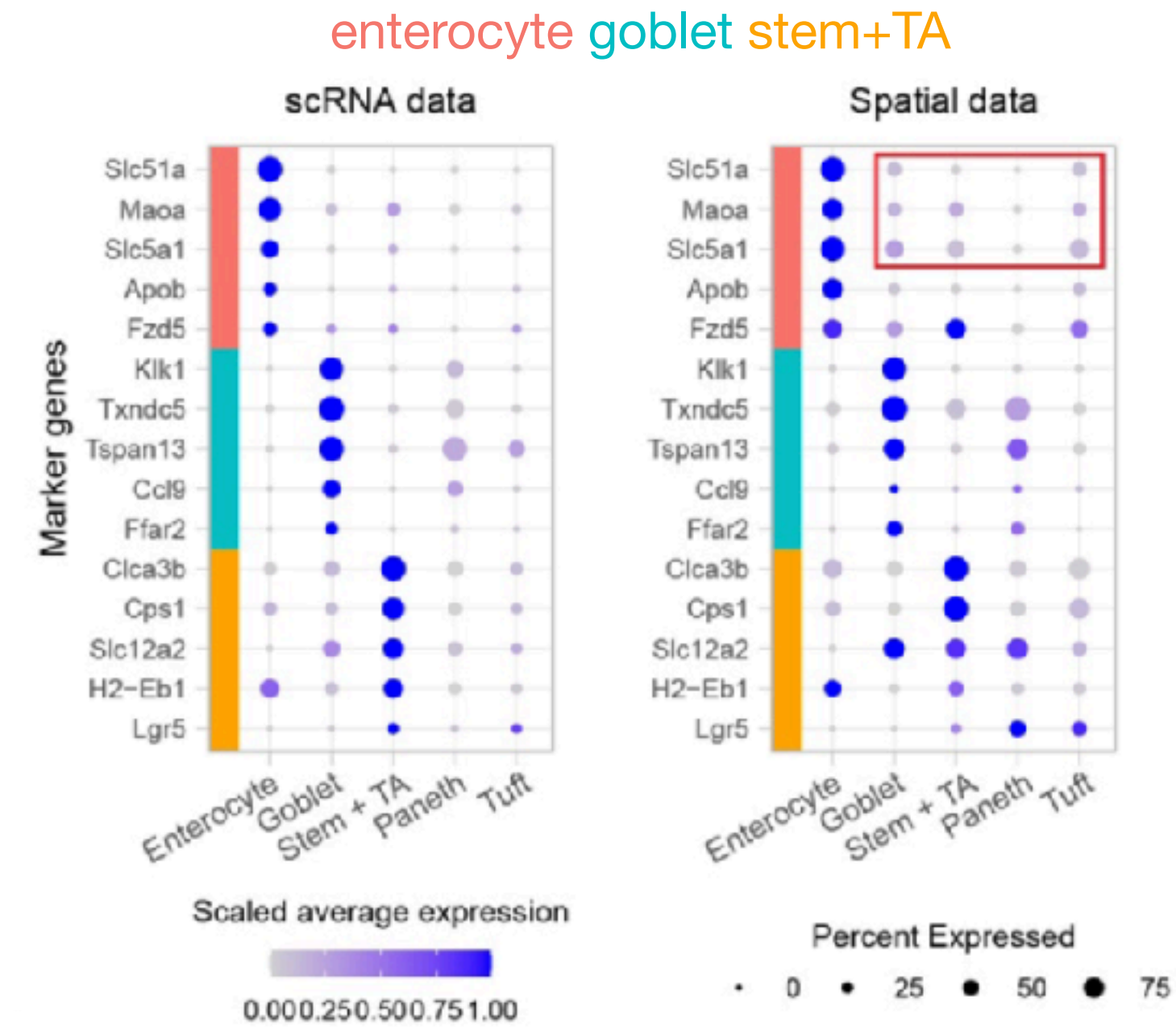
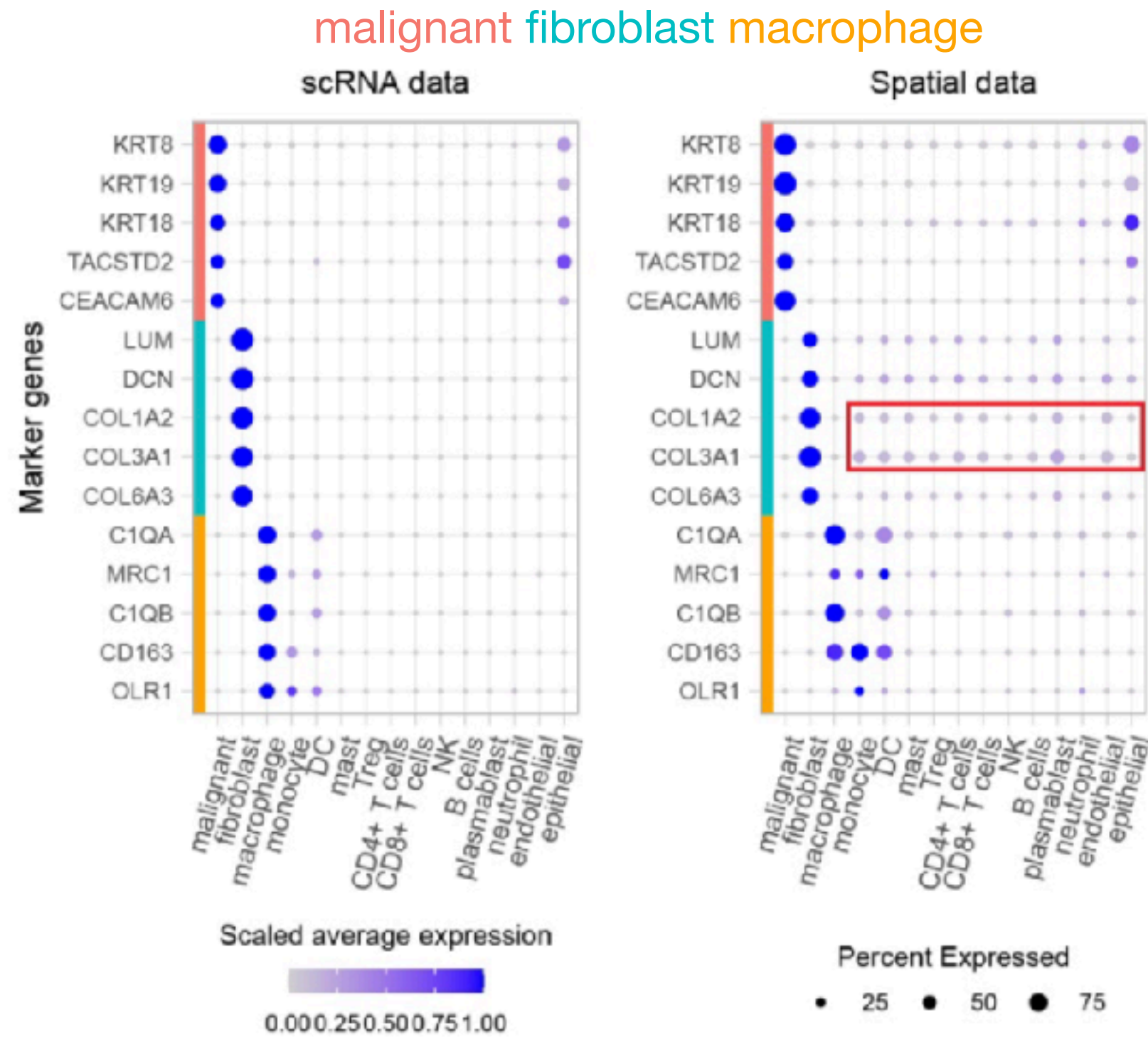


# bleeding occurs in 3D – around, above & below cells

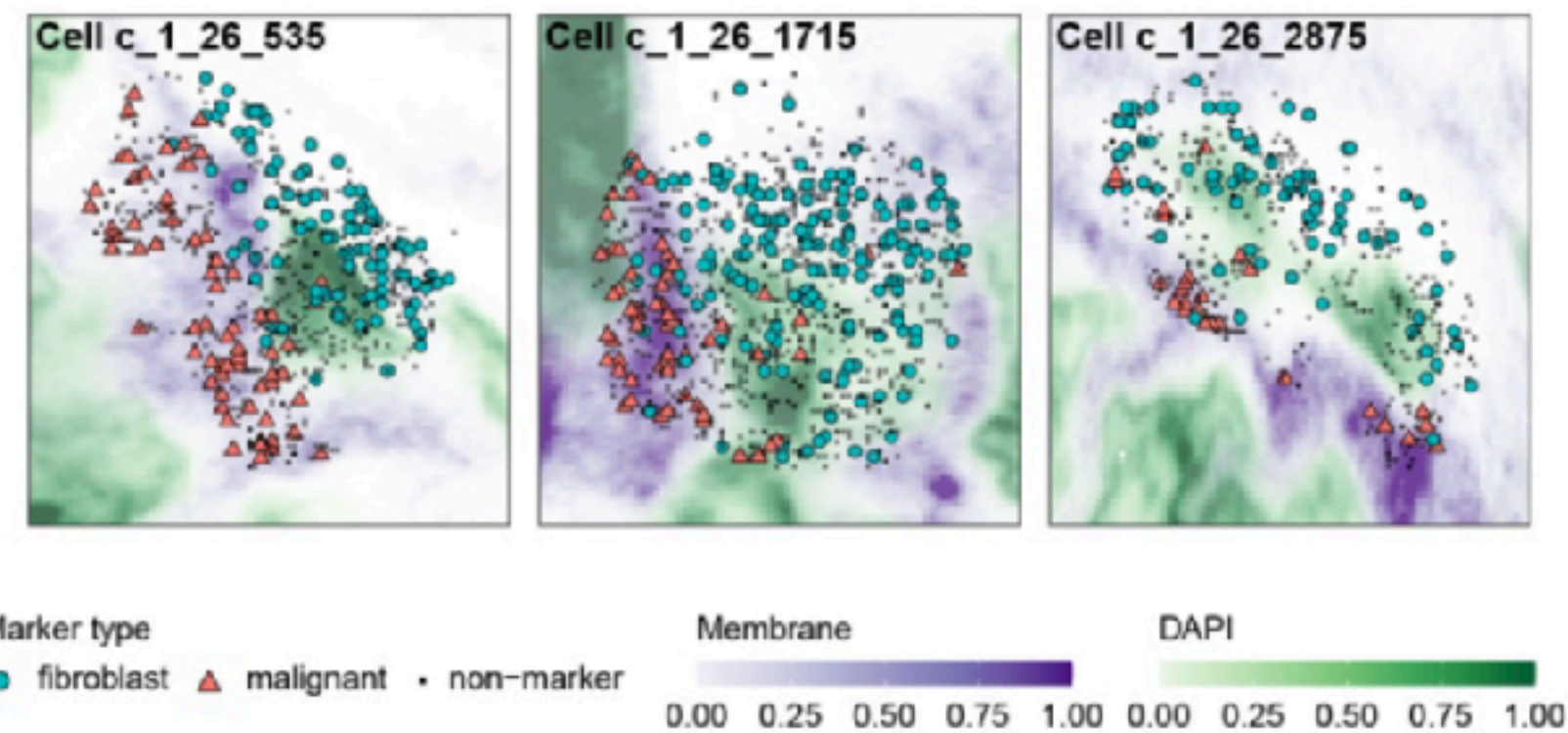




# bleeding occurs in 3D – around, above & below cells

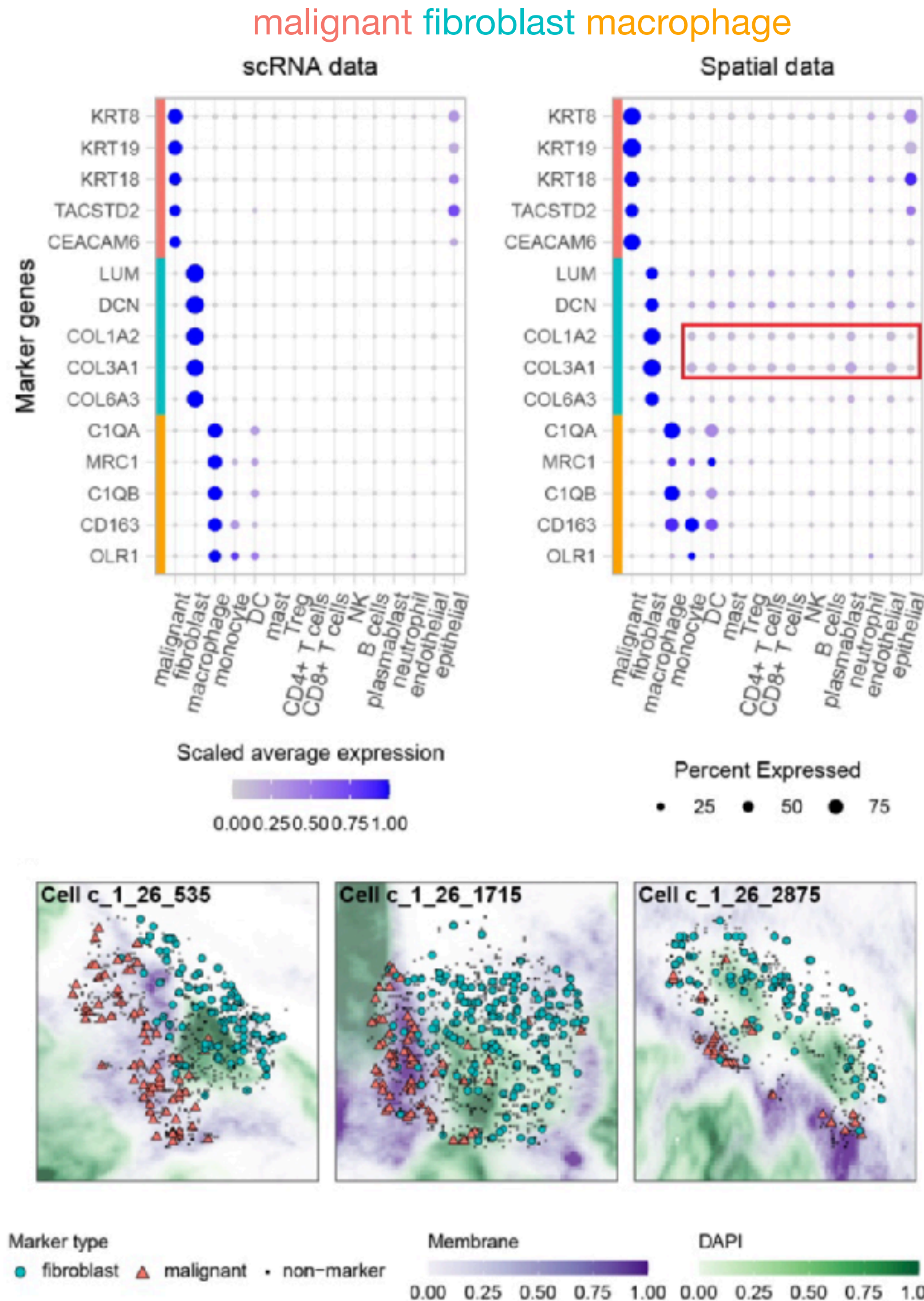


admixture occurs at cellular periphery

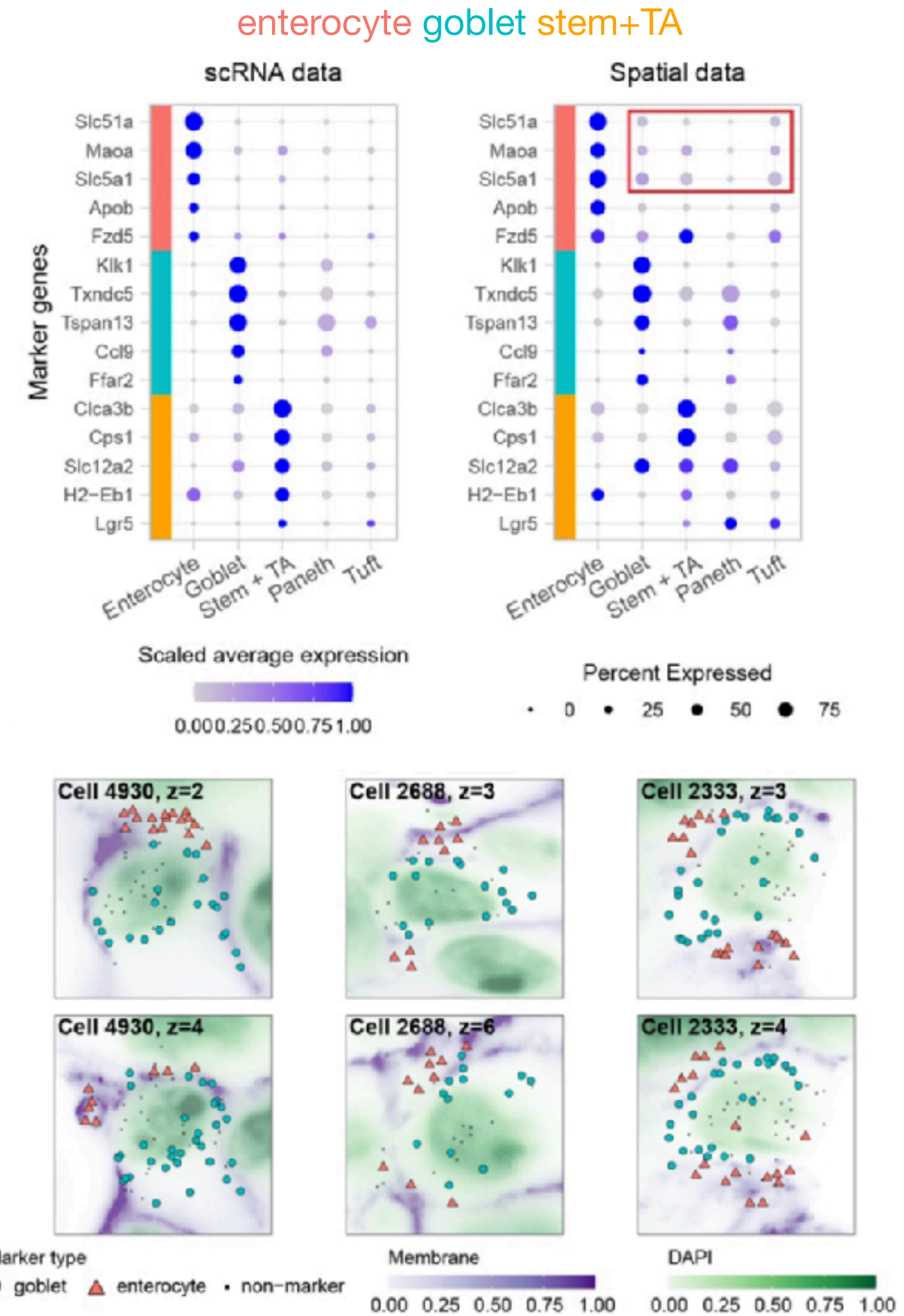




# bleeding occurs in 3D – around, above & below cells



admixtures occur at cellular periphery

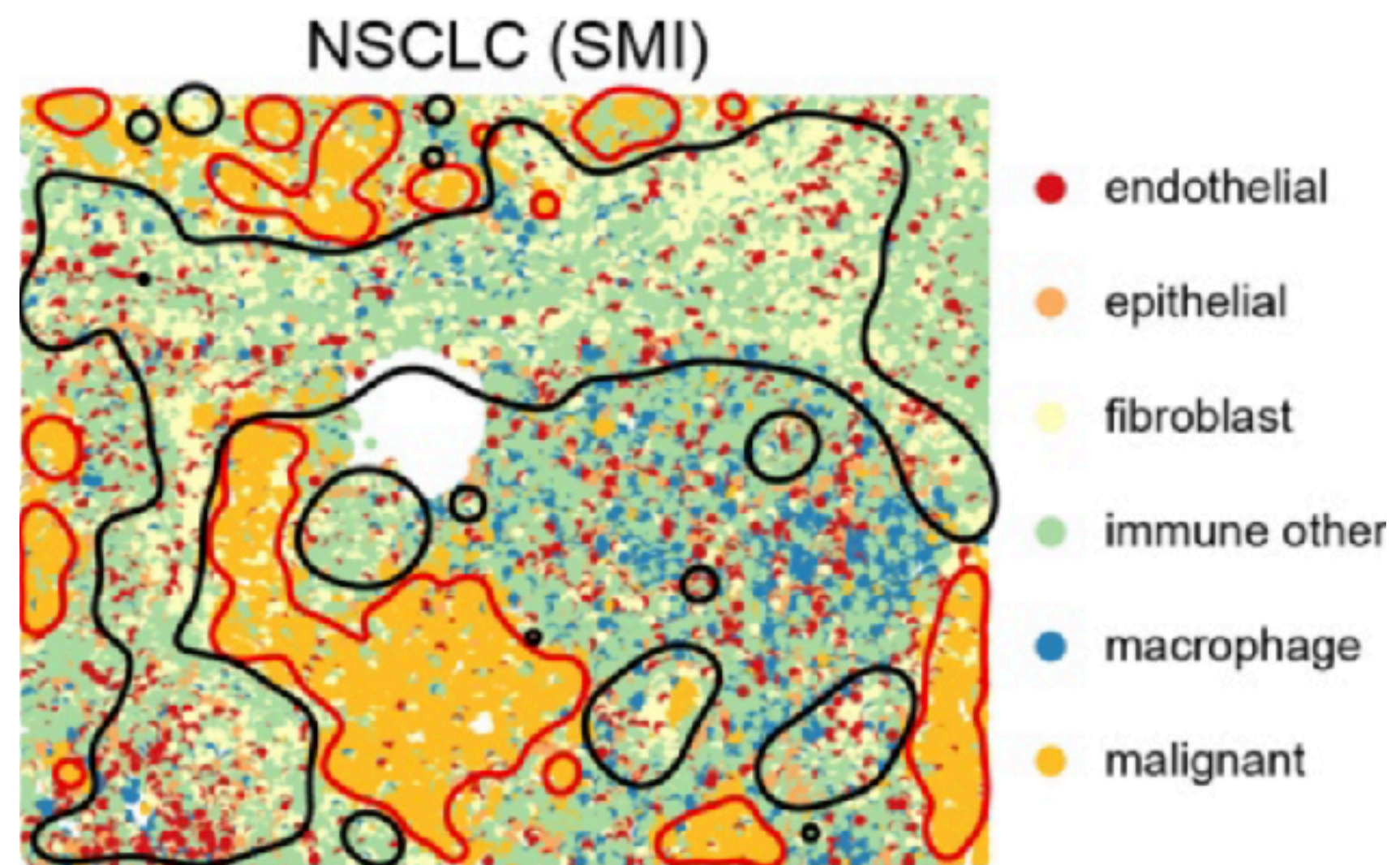


admixtures occur at upper/lower z-plane



# DE genes between regions reflect compositional differences (not differences in state)

- manual annotation into **tumor** & **stromal** regions

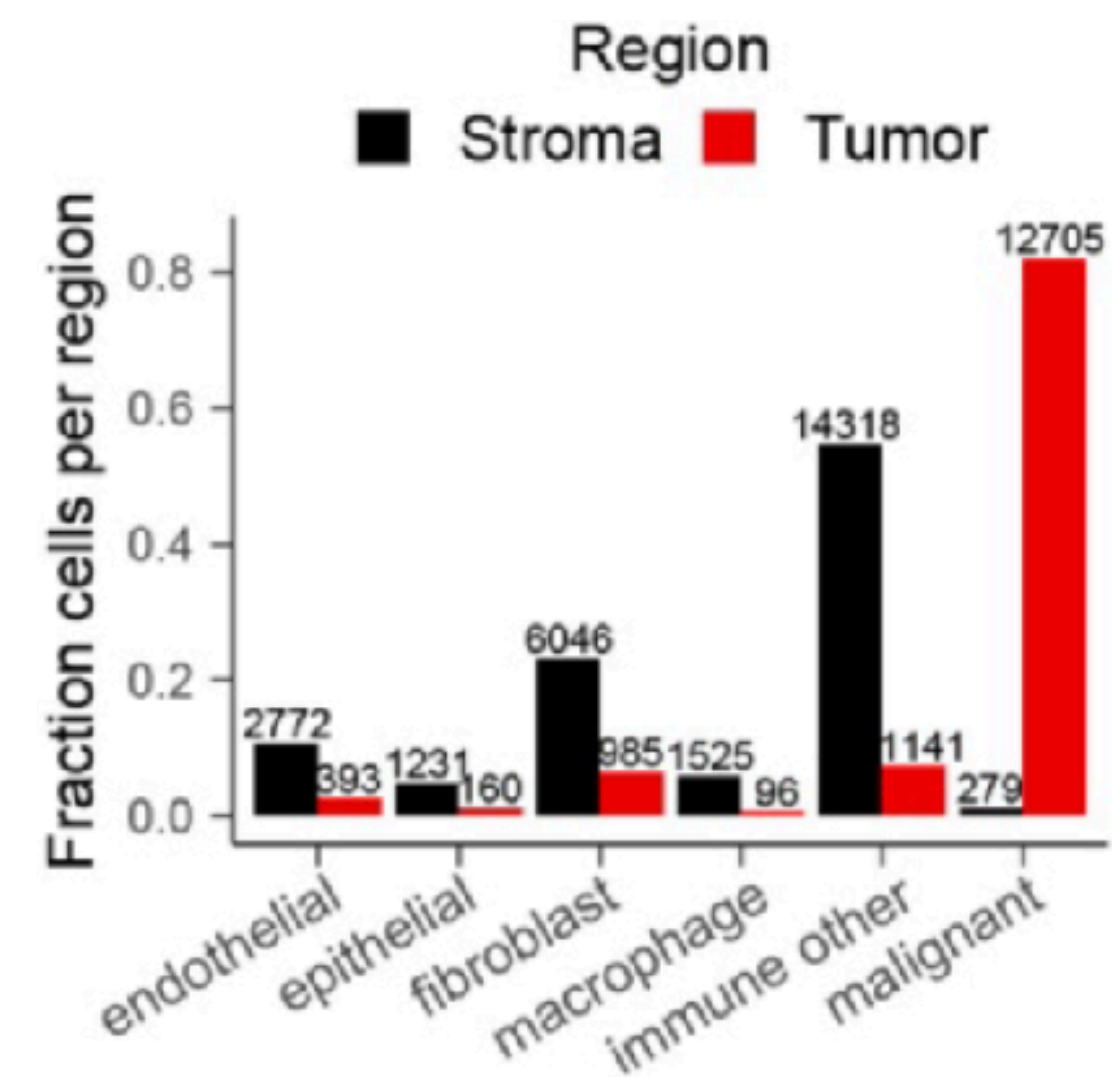
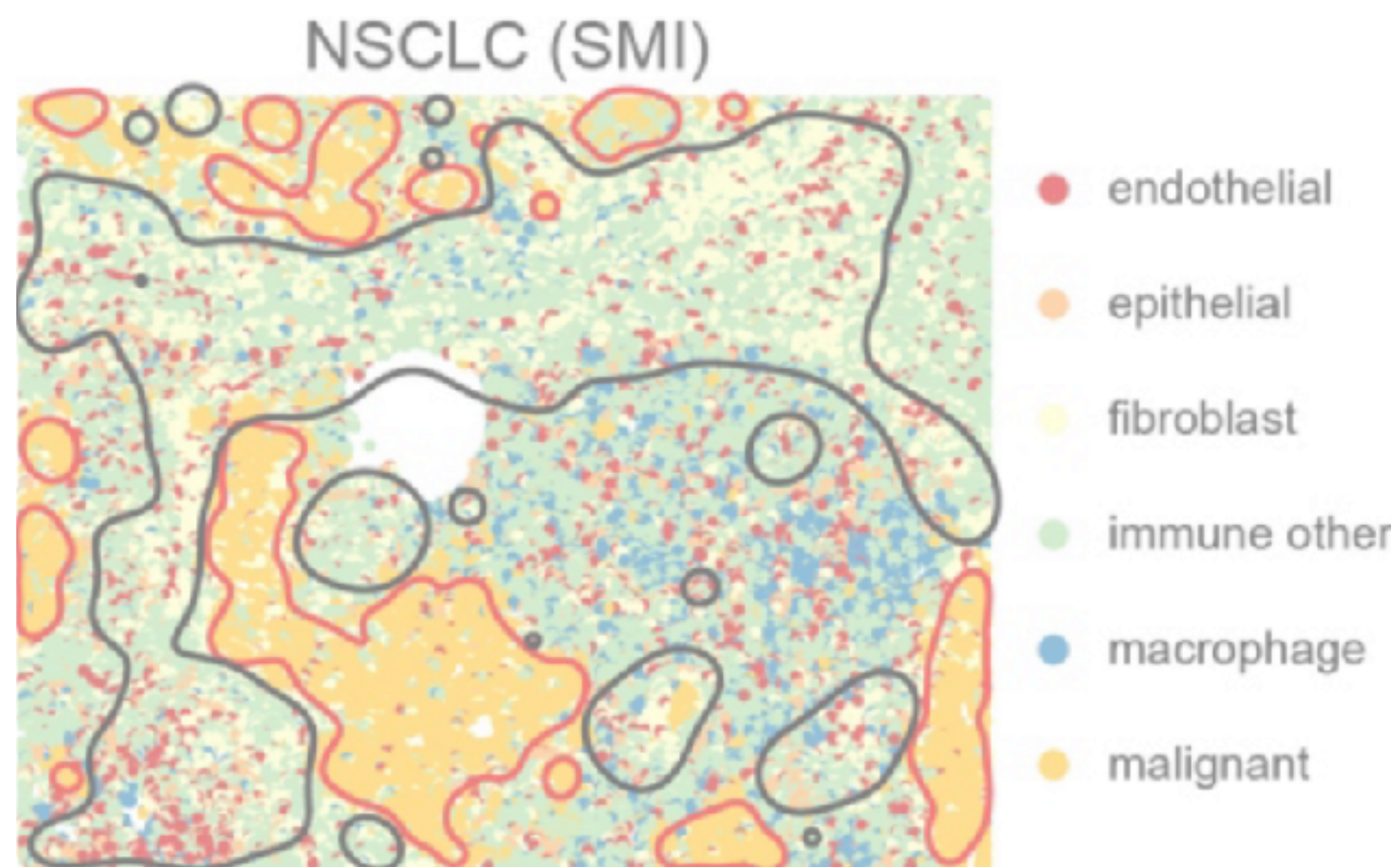




# DE genes between regions reflect compositional differences (not differences in state)

- manual annotation into **tumor** & **stromal** regions

- tumor is dominated by malignant, stroma is dominated by other cells



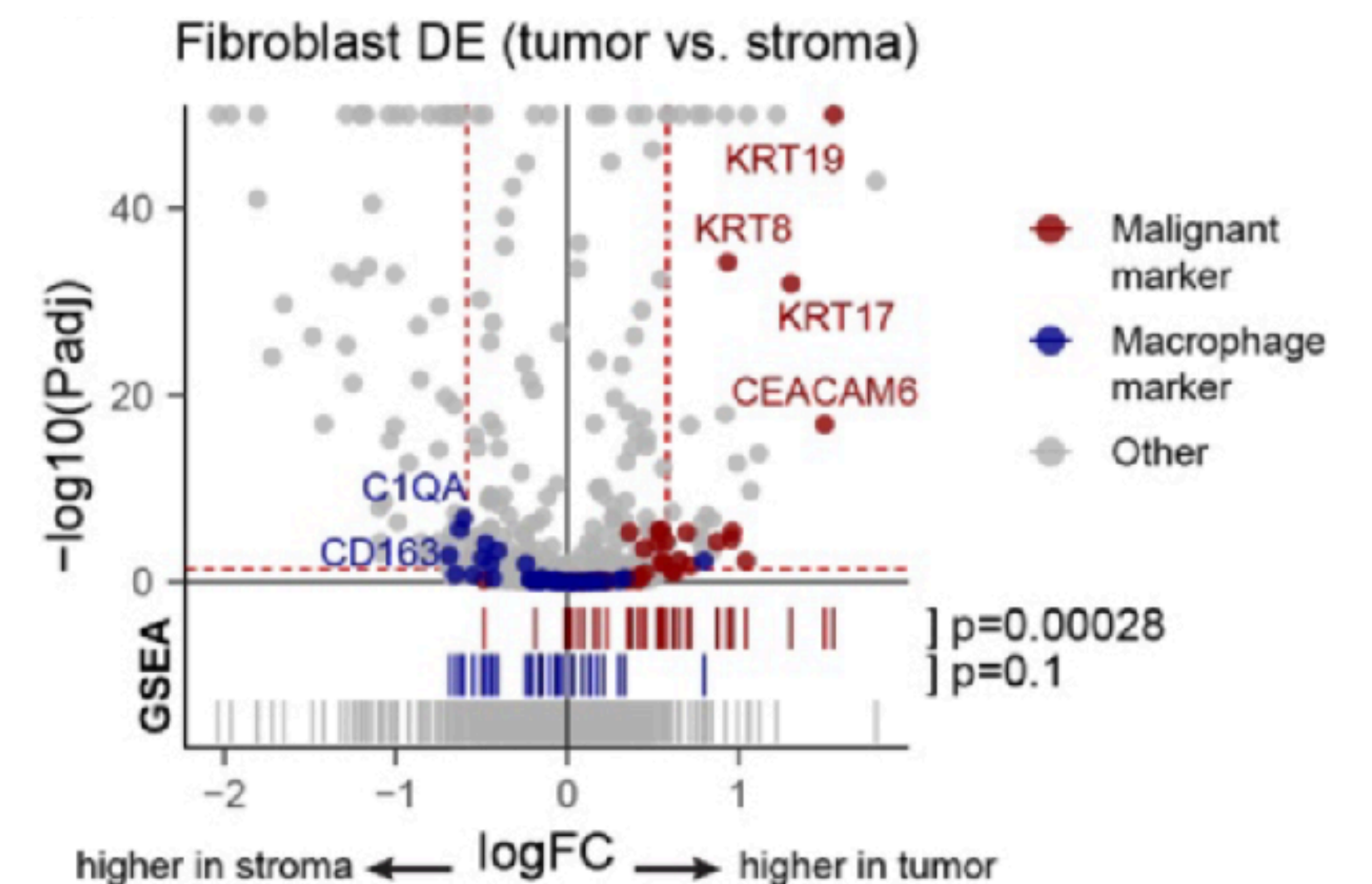
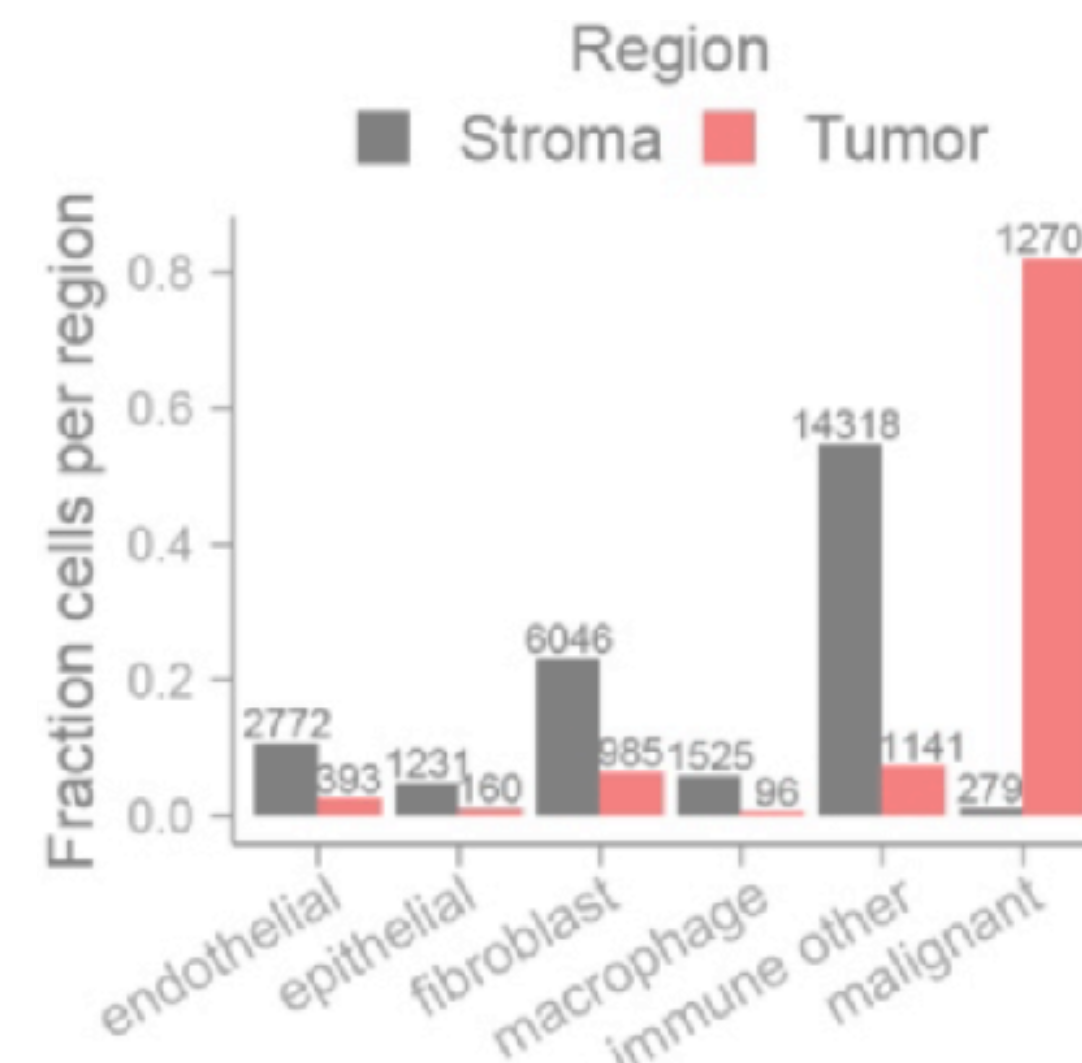
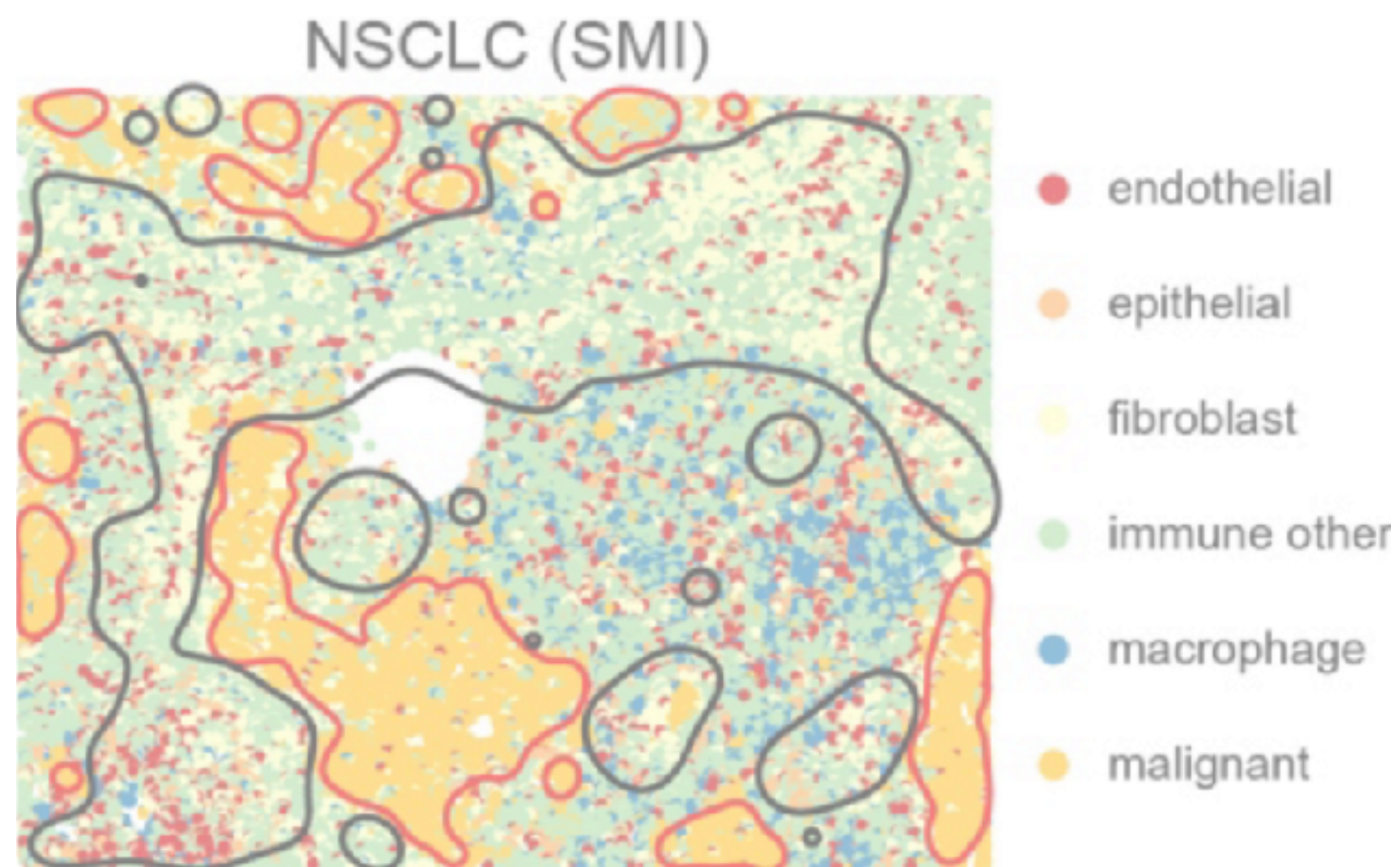


# DE genes between regions reflect compositional differences (not differences in state)

- manual annotation into **tumor** & **stromal** regions

- tumor is dominated by malignant, stroma is dominated by other cells

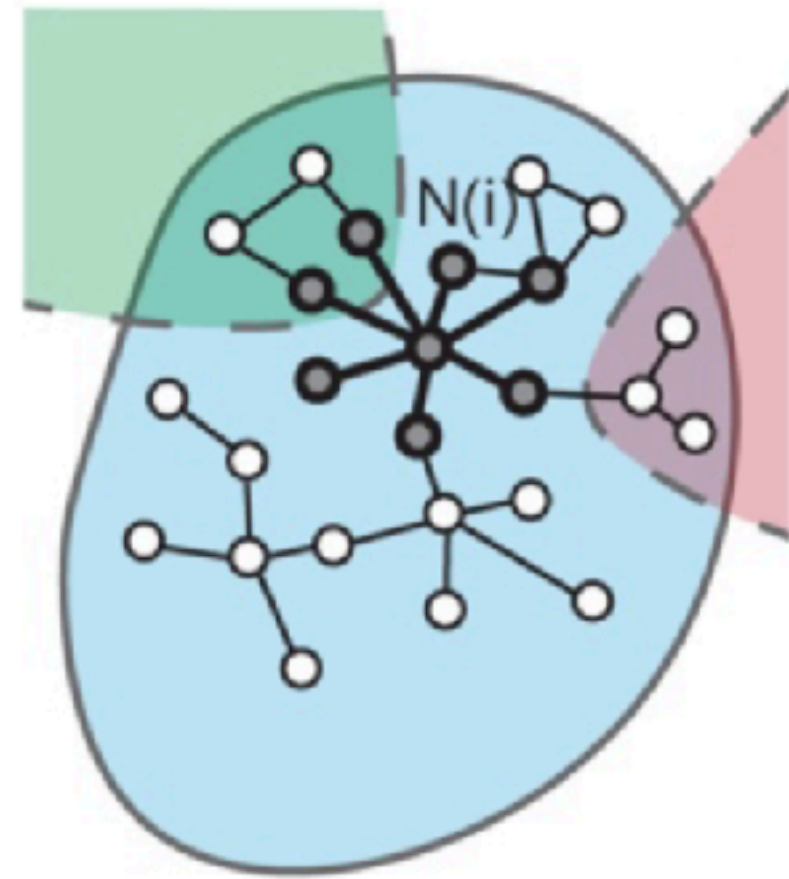
- comparing regions, genes upregulated in fibroblasts are **epithelial markers**



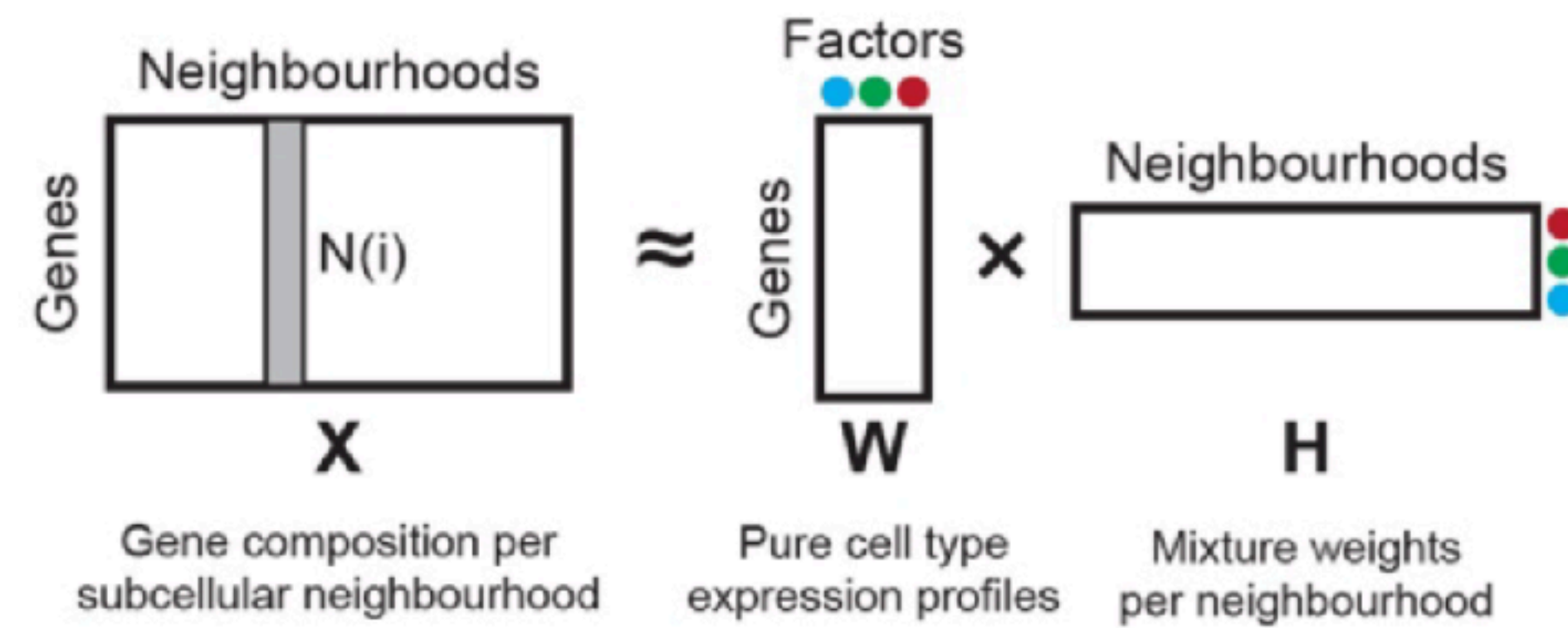


Mitchel *et al.* propose **NMF + CRF clean-up** to mitigate spatial bleeding

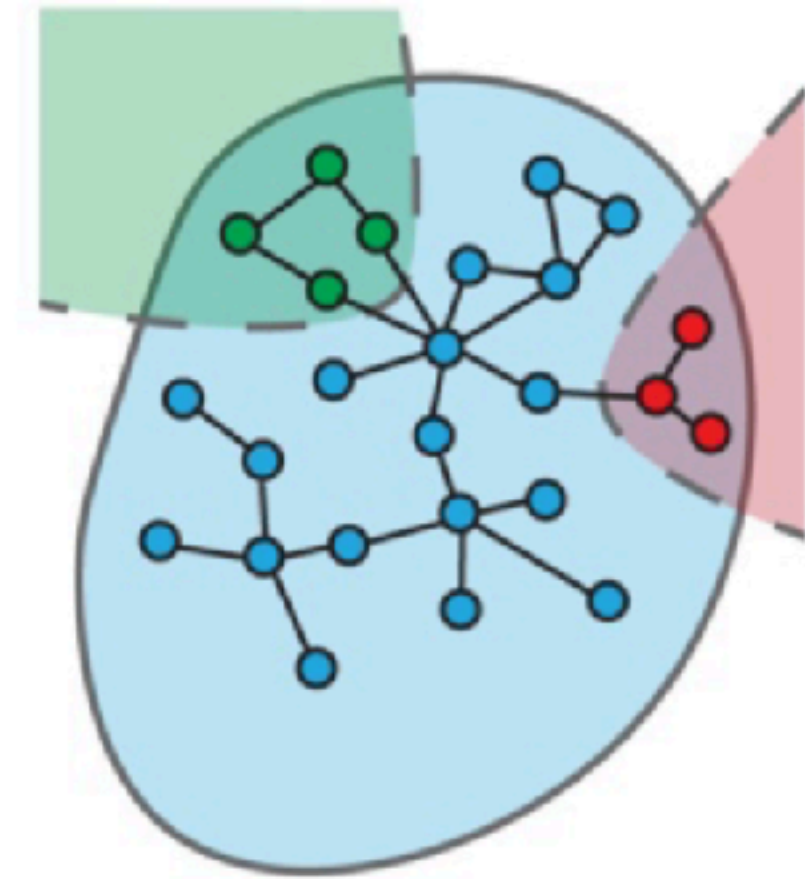
1. Construct KNN graph per cell



2. Recover pure expression profiles using weighted NMF



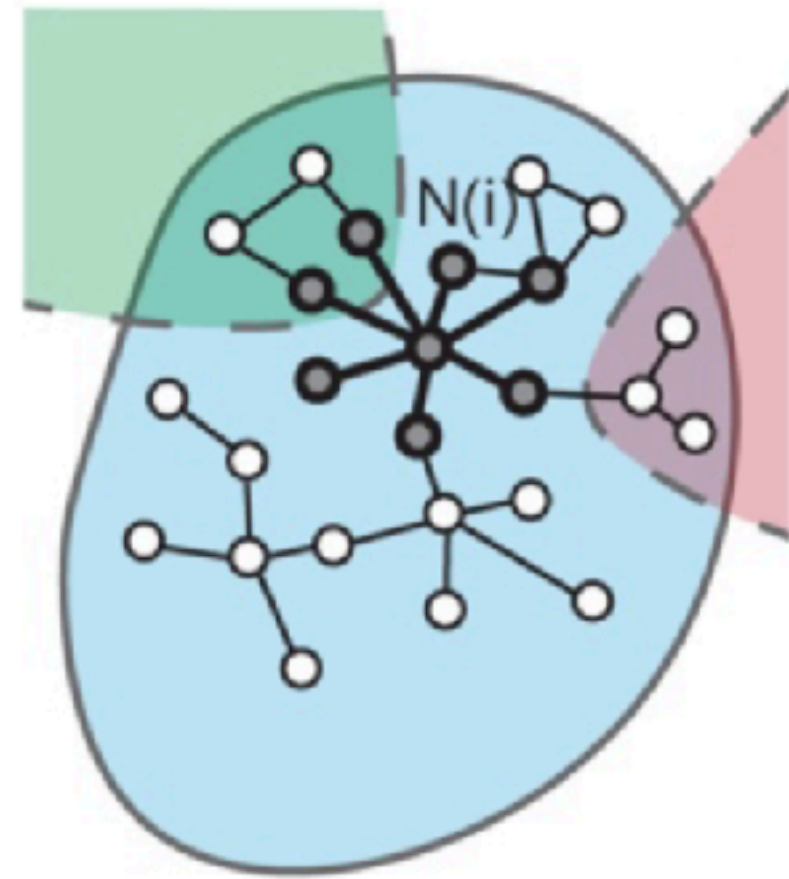
3. Label admixture molecules using CRF



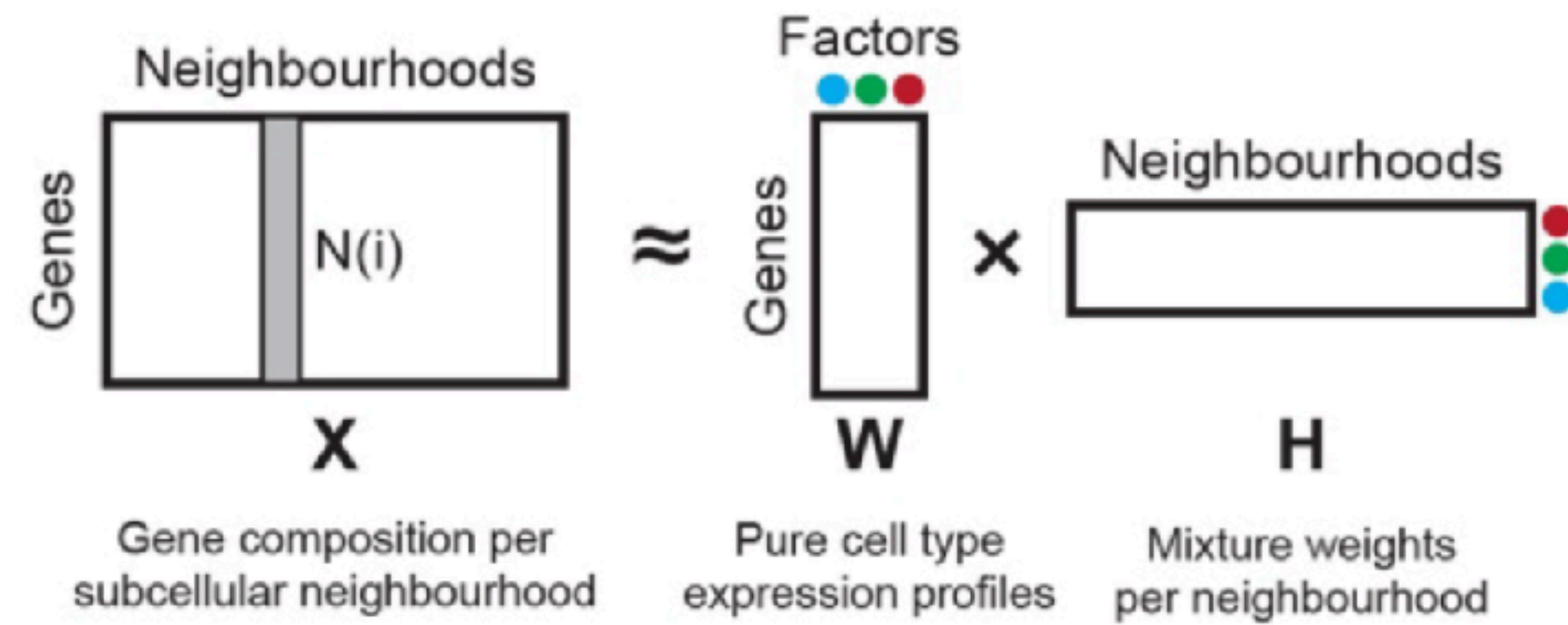


Mitchel *et al.* propose **NMF + CRF clean-up** to mitigate spatial bleeding

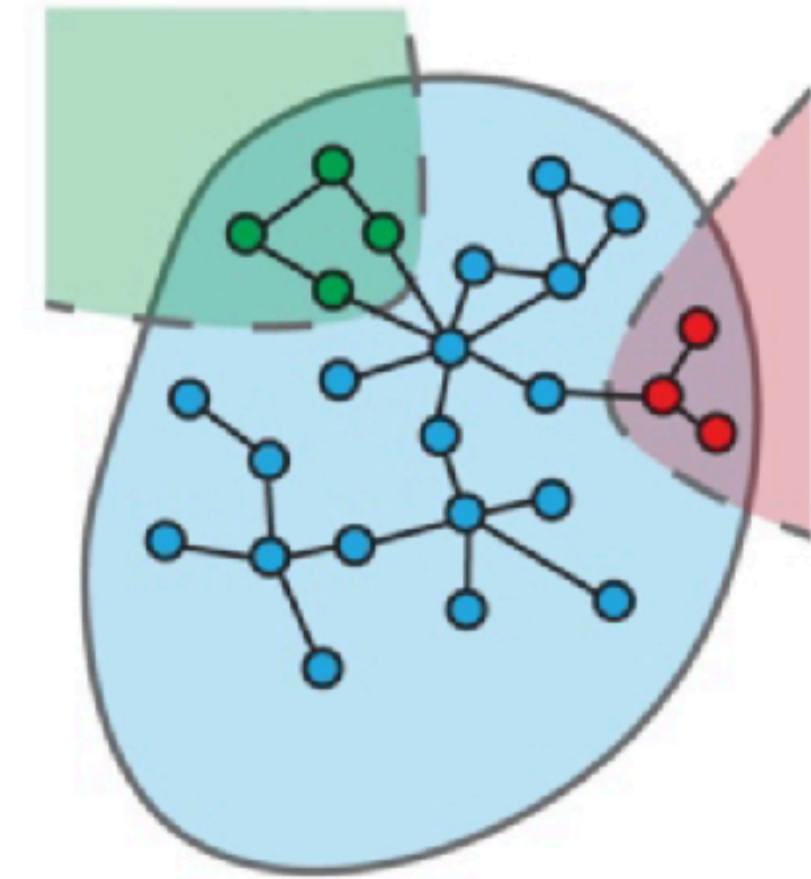
1. Construct KNN graph per cell



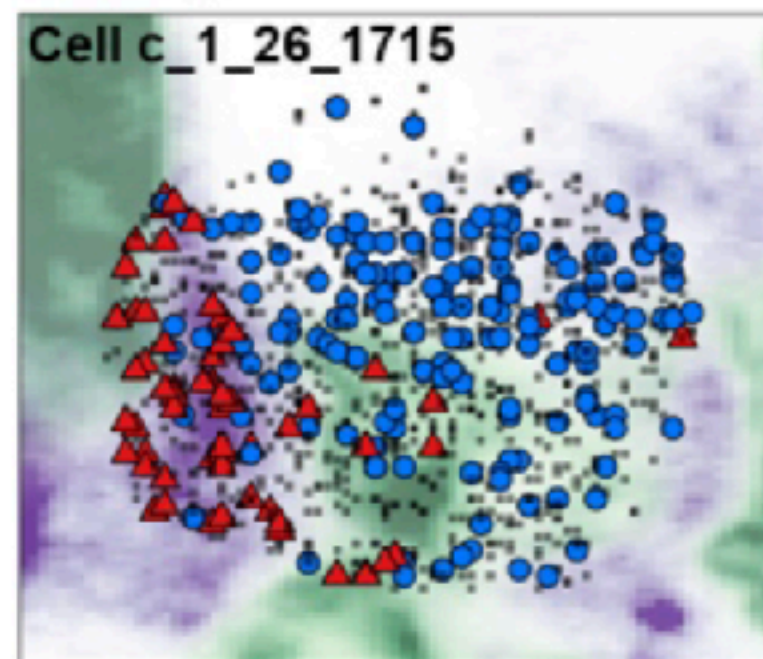
2. Recover pure expression profiles using weighted NMF



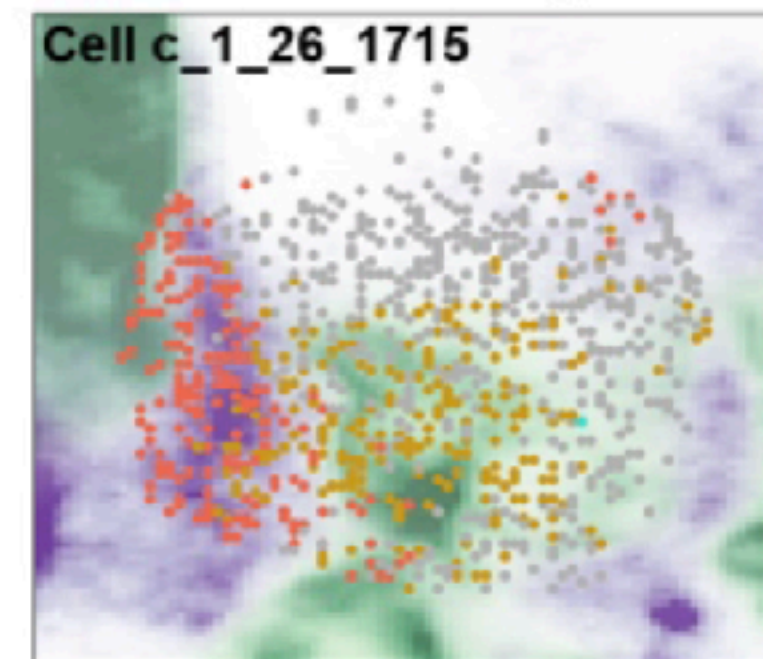
3. Label admixture molecules using CRF



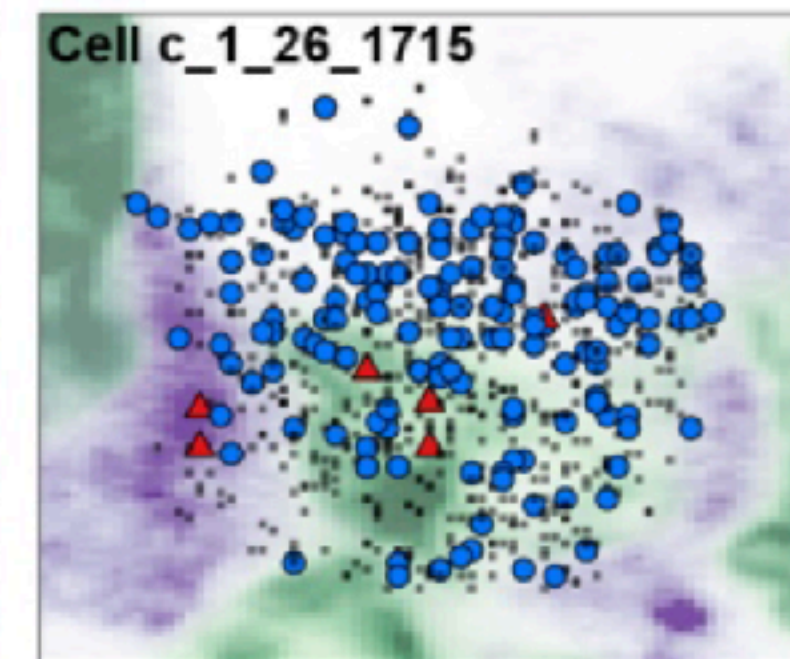
Original fibroblast cells



CRF molecule assignments



Cleaned fibroblast cells



Marker type

● fibroblast ▲ malignant • non-marker

Factor

● 1 ● 2 ● 3 ● 4 ● 5

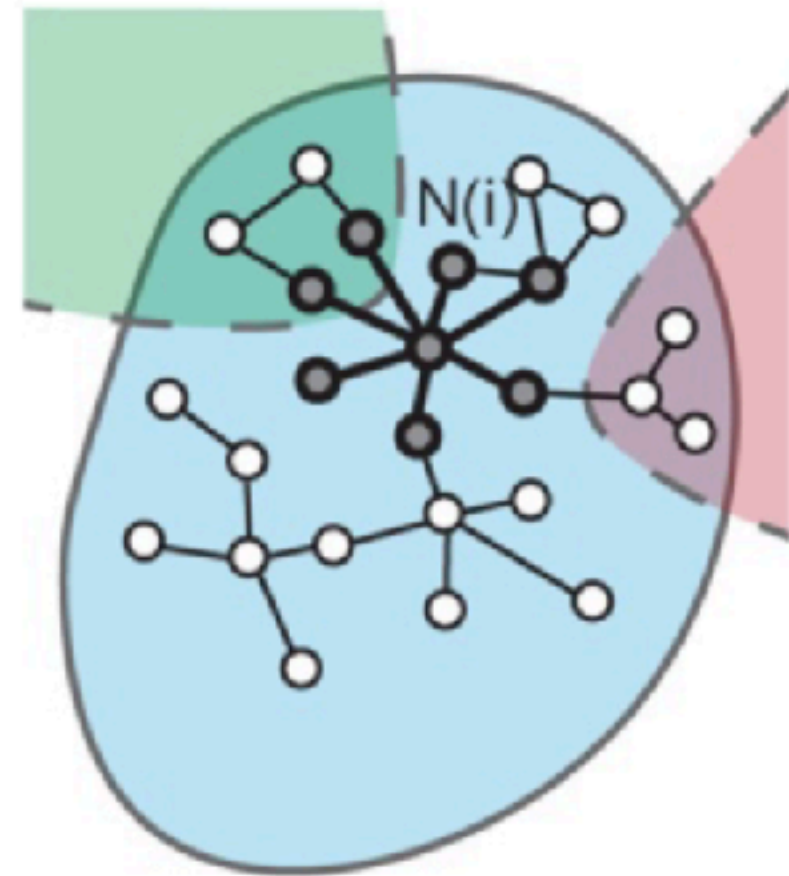
Marker type

● fibroblast ▲ malignant • non-marker

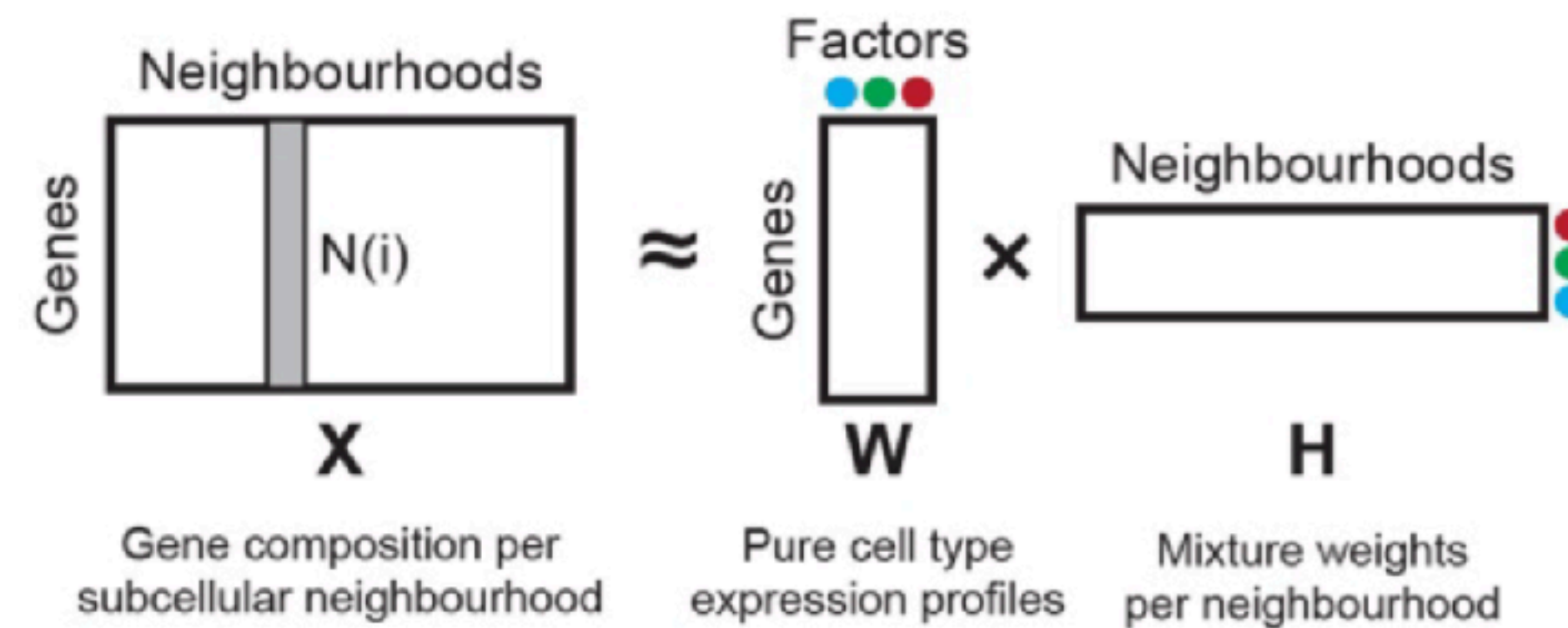


# Mitchel *et al.* propose **NMF + CRF clean-up** to mitigate spatial bleeding

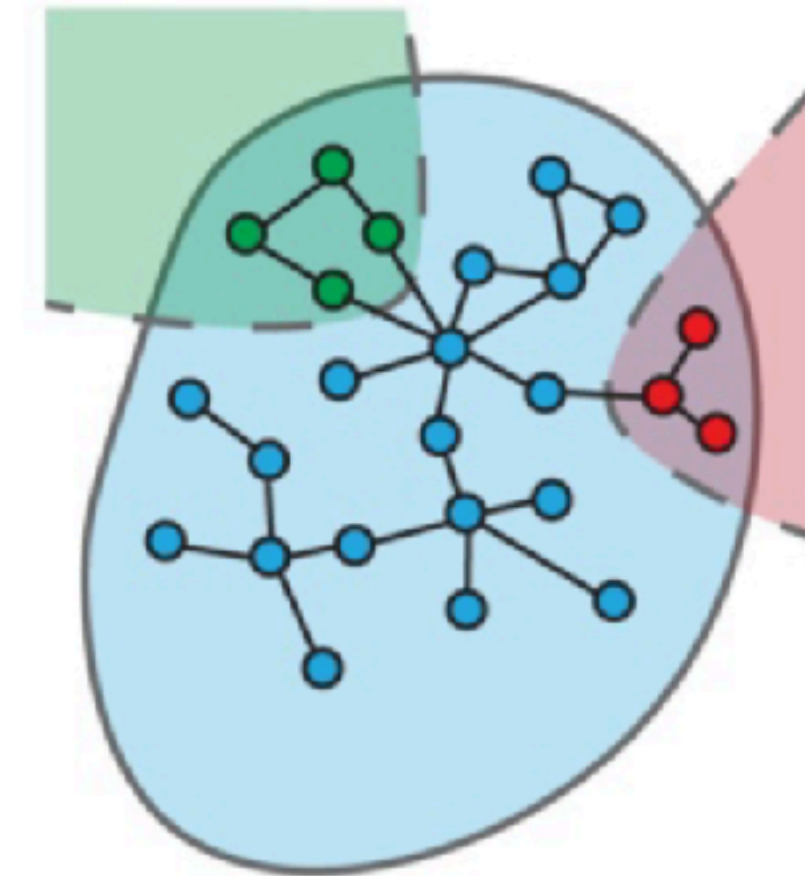
1. Construct KNN graph per cell



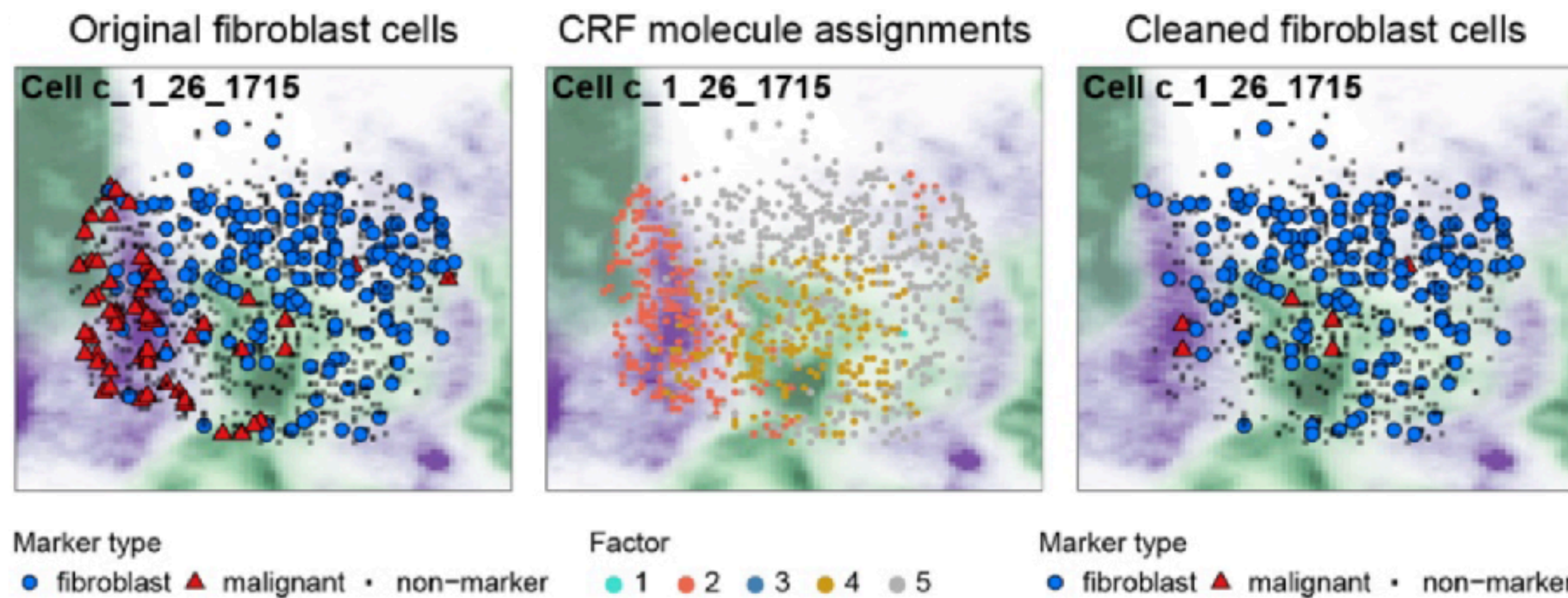
2. Recover pure expression profiles using weighted NMF



3. Label admixture molecules using CRF



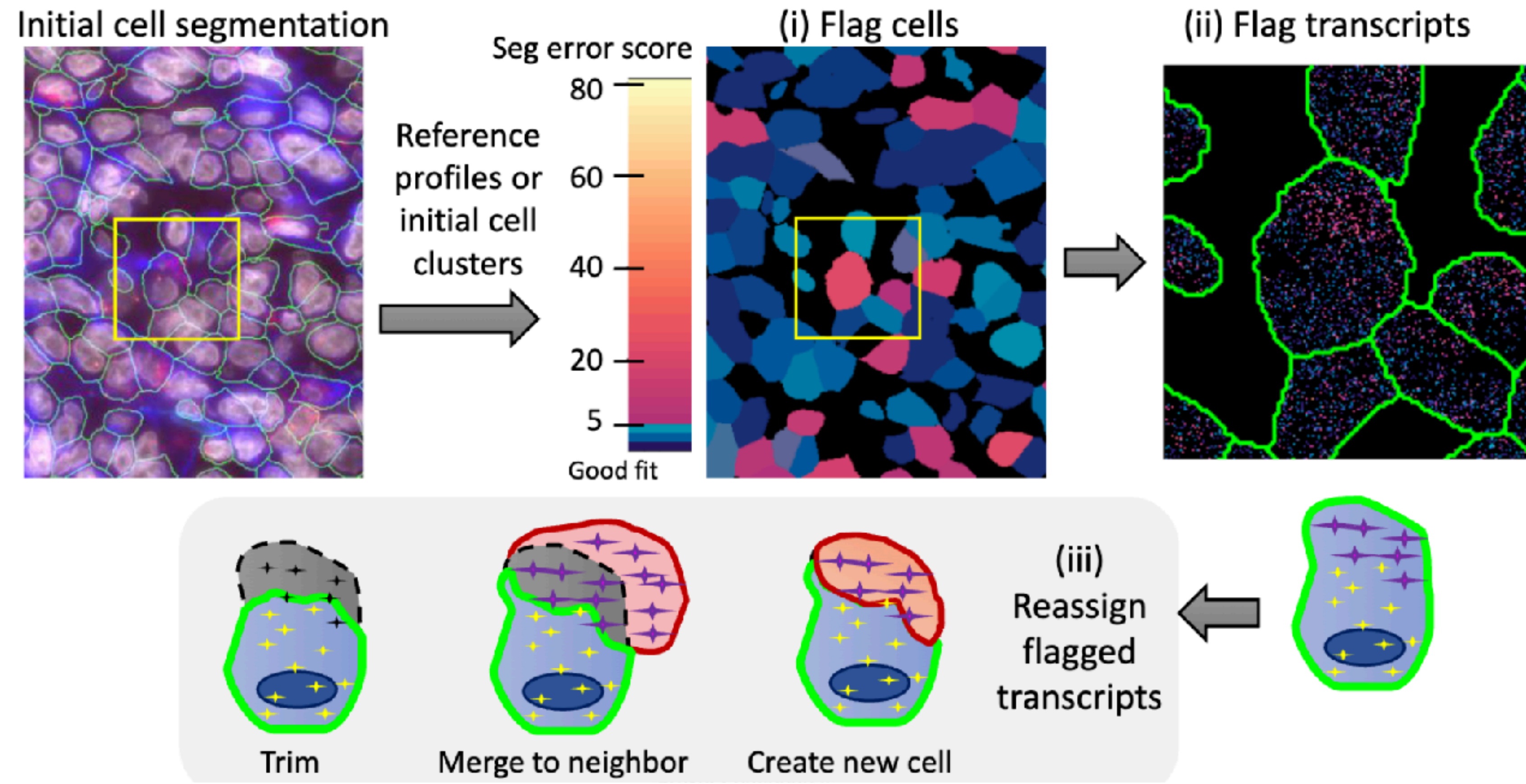
- **subcellular features** can include
  - recurrent **admixture patterns** (e.g., between frequently co-occurring cell types)
  - true **cellular structures** (e.g., ER, nuclei, polarization)





# FastReseg uses transcript locations to refine img-based segmentation

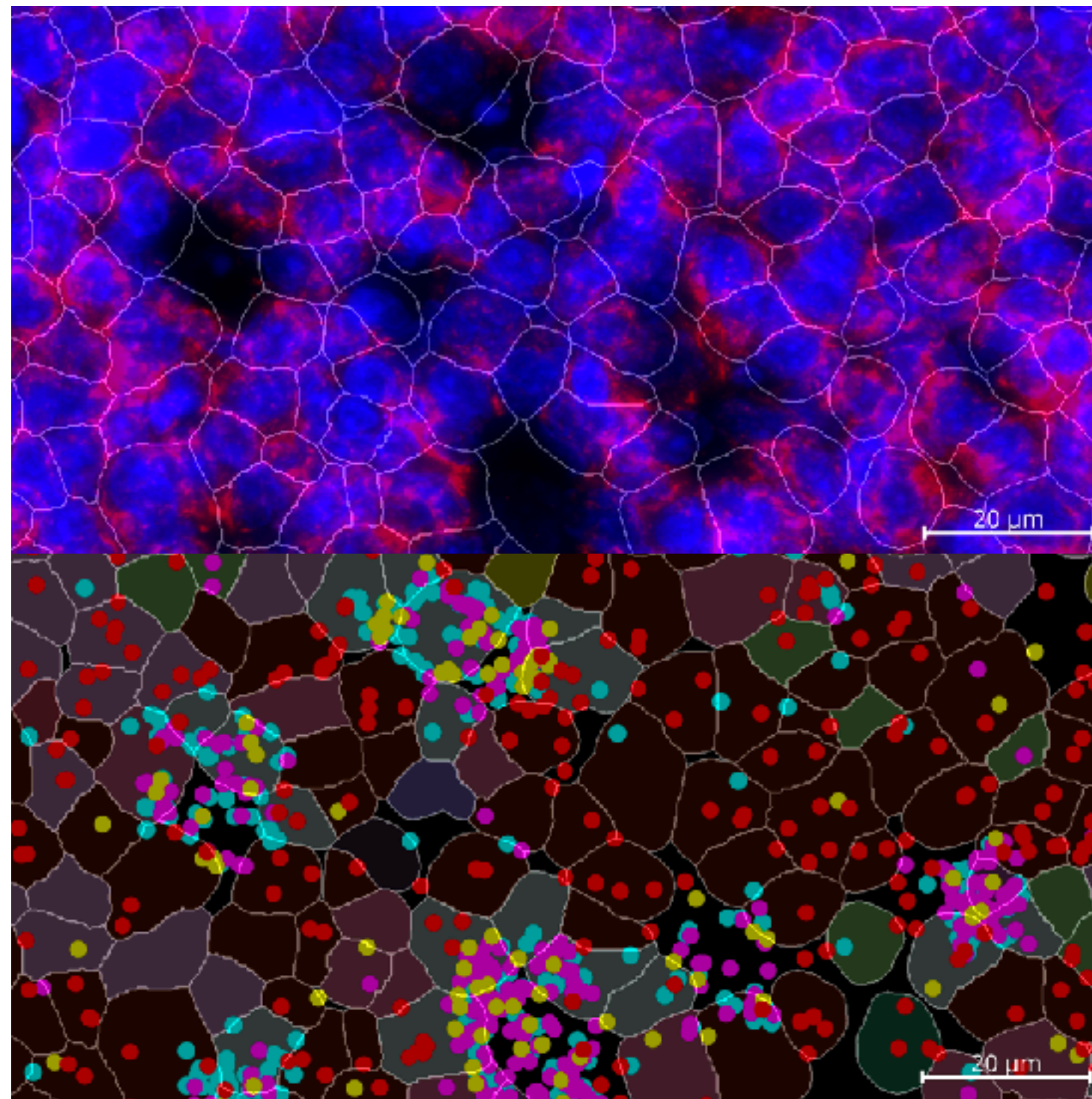
- **transcript scoring** based on initial host cell
- **flag spatial doublets** as putative segmentation errors
- **flag misassigned transcripts** within flagged cells only
- **correct counts** (but not segmentation boundaries)



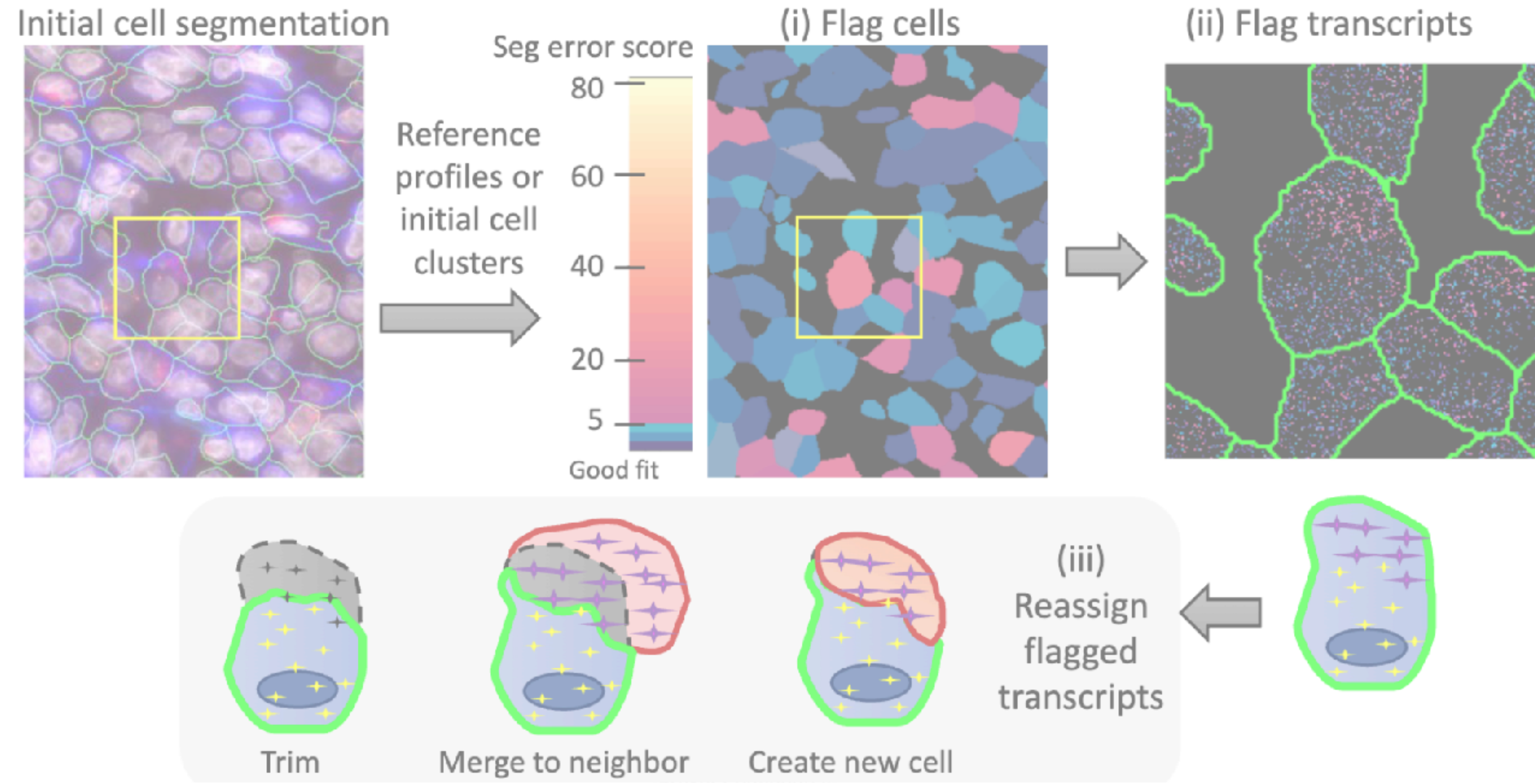


# FastReseg uses transcript locations to refine img-based segmentation

- flag B cells surrounding black holes
- flag macrophage-related genes
- correct B cell counts & create new macrophages

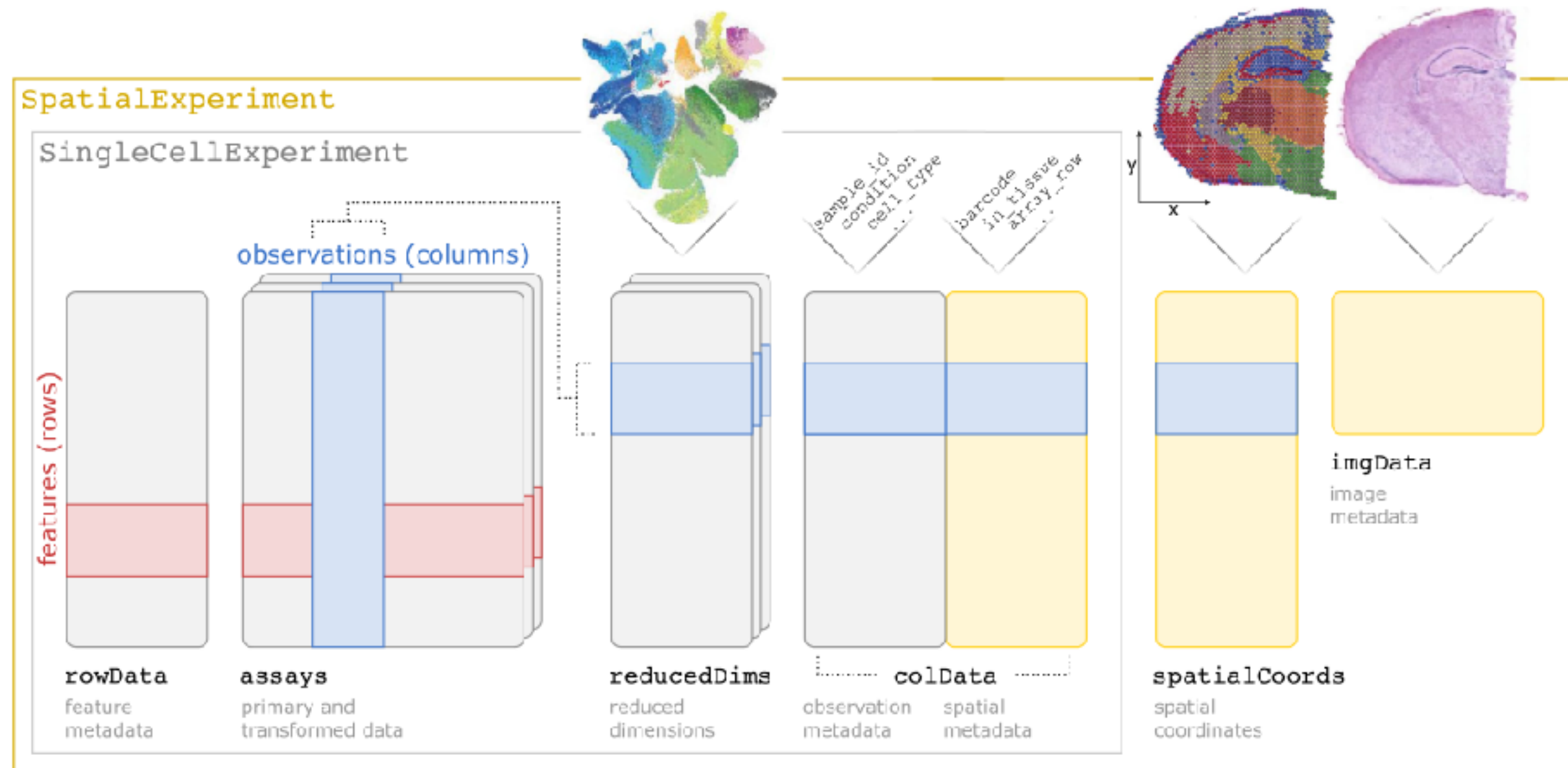


CD19 MMP9 CD68 LYZ



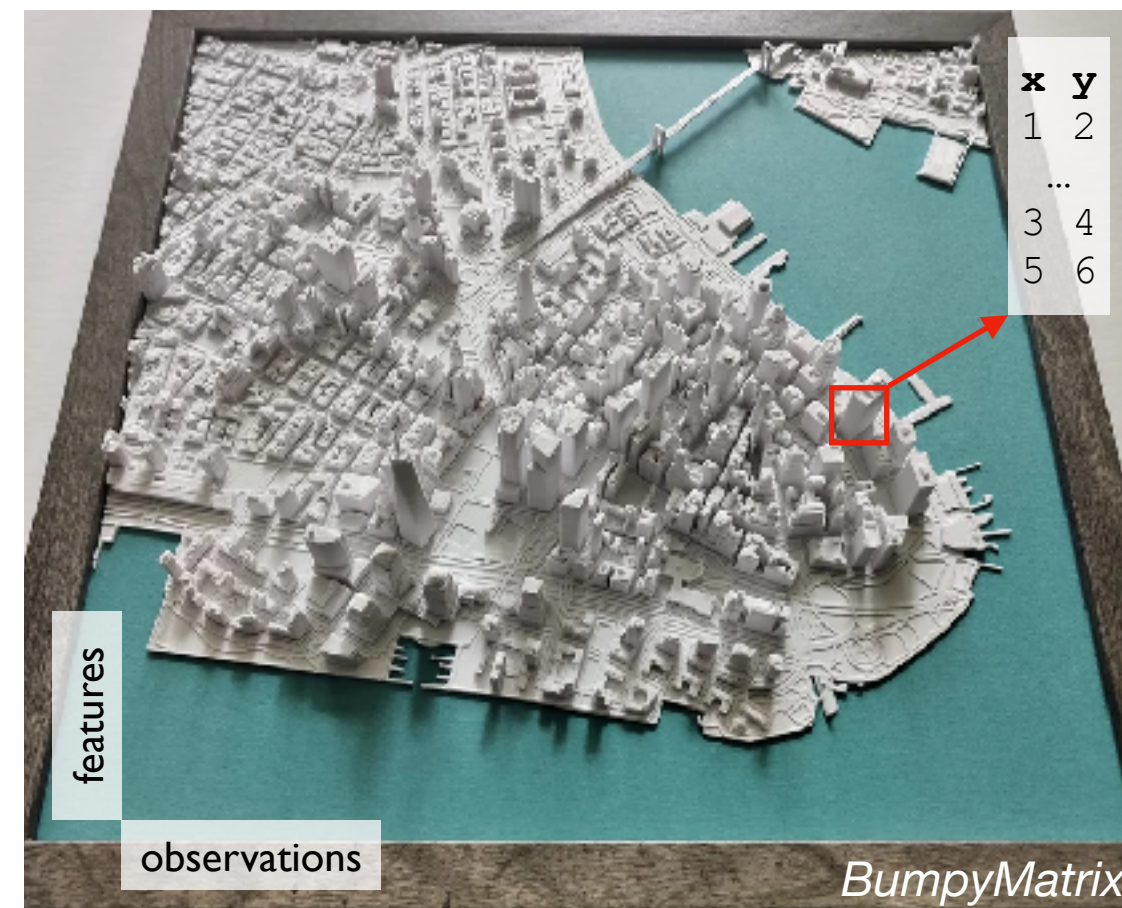


# infrastructure for handling img-ST data in R



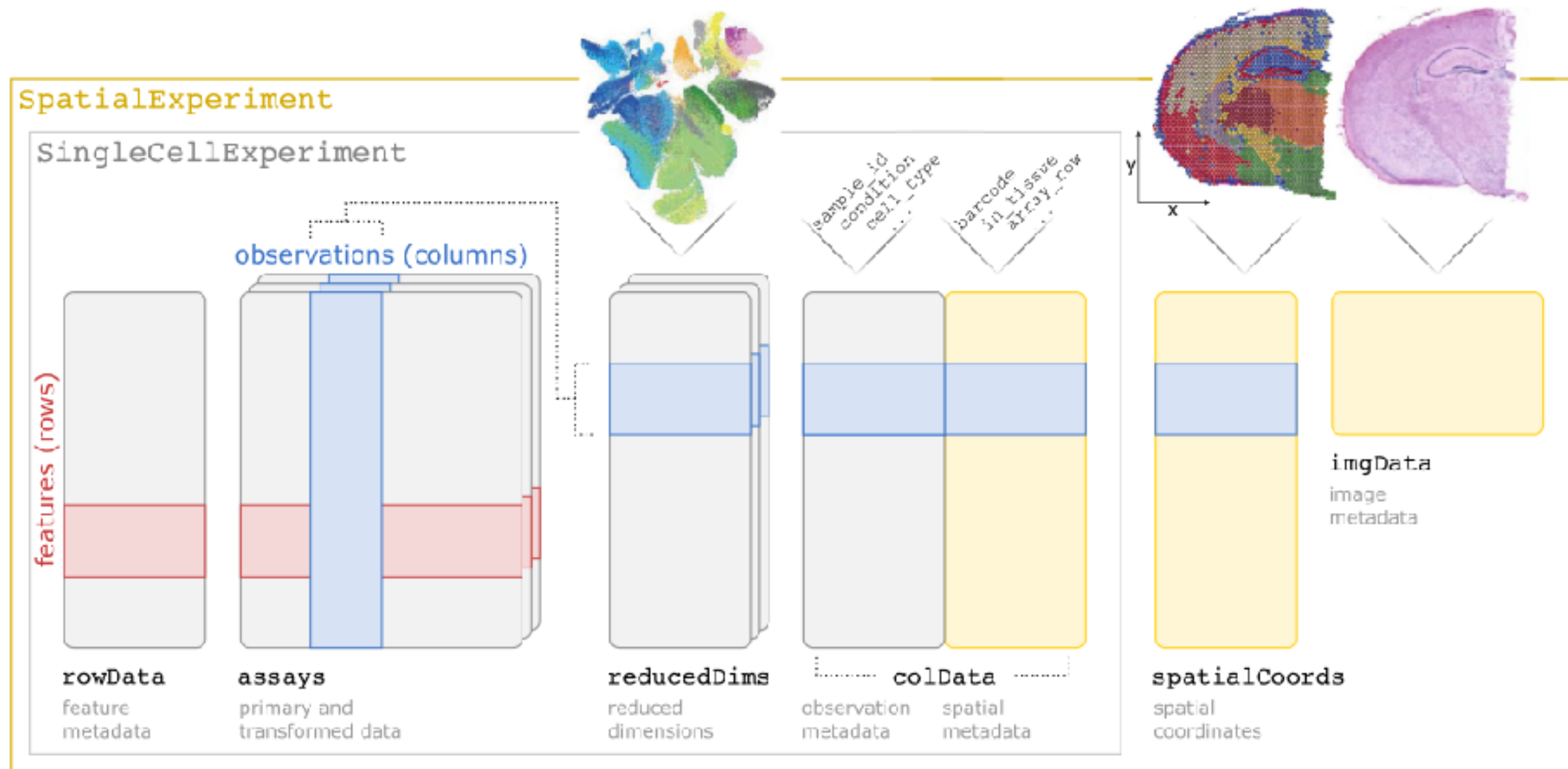
Righelli, Weber, Crowell et al. (2022)  
*Bioinformatics* 38(11):3128-3131

## SpatialExperiment



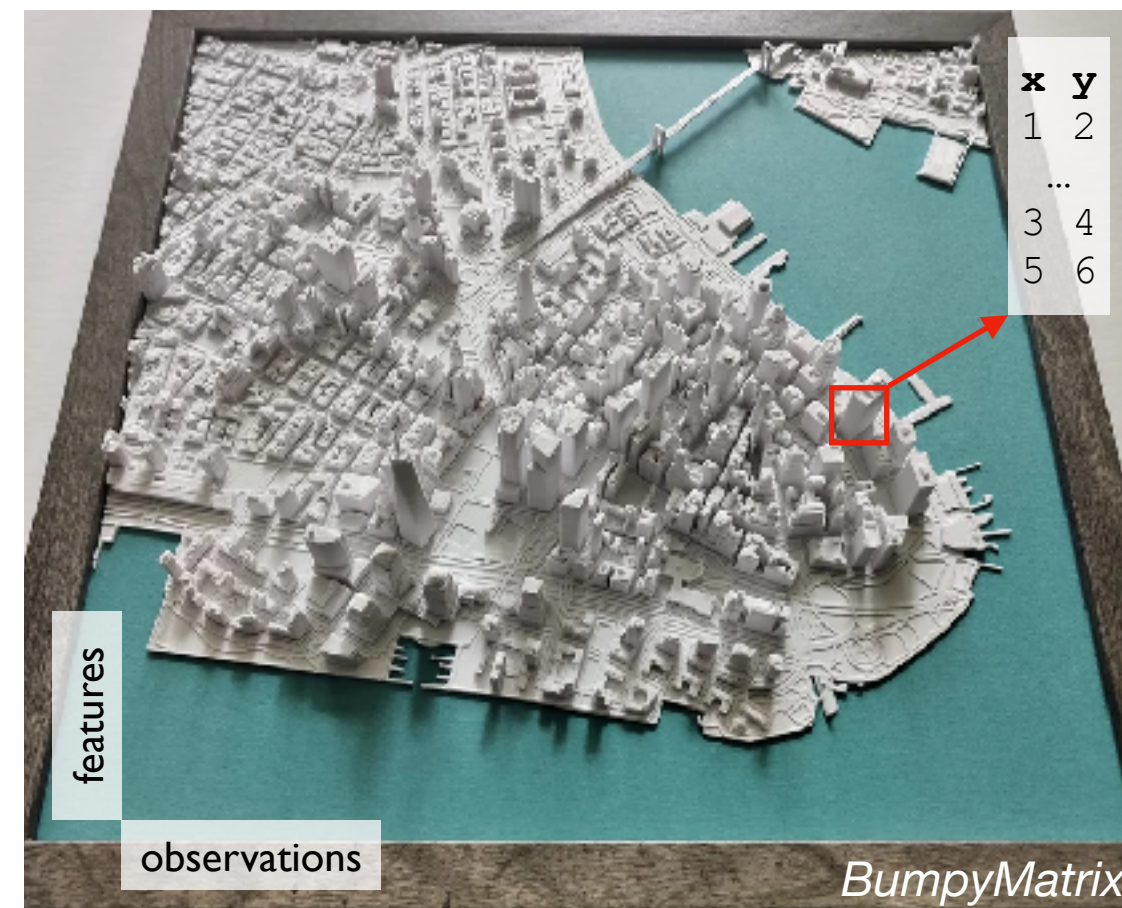


# infrastructure for handling img-ST data in R

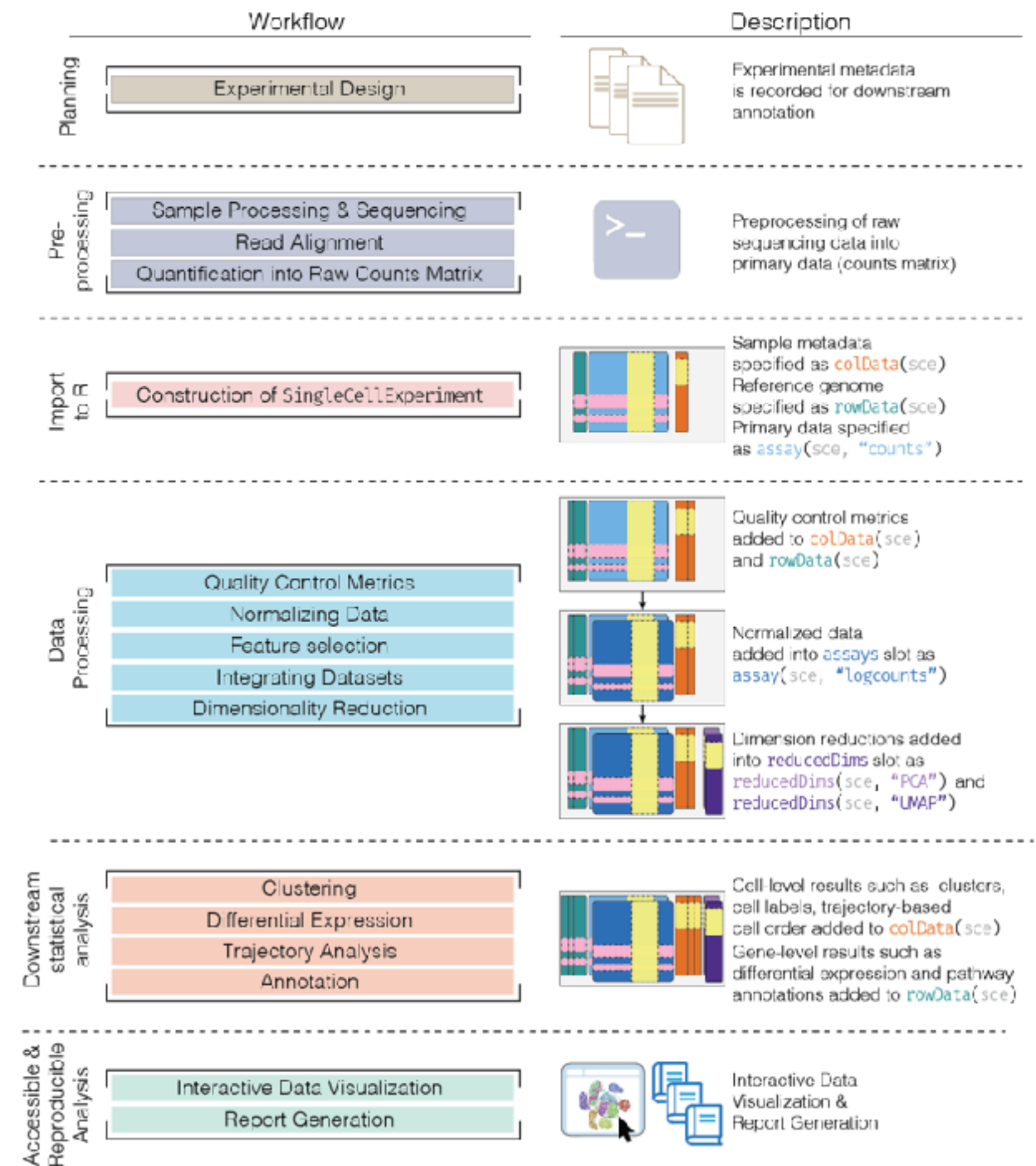


Righelli, Weber, Crowell et al. (2022)  
*Bioinformatics* 38(11):3128-3131

## SpatialExperiment

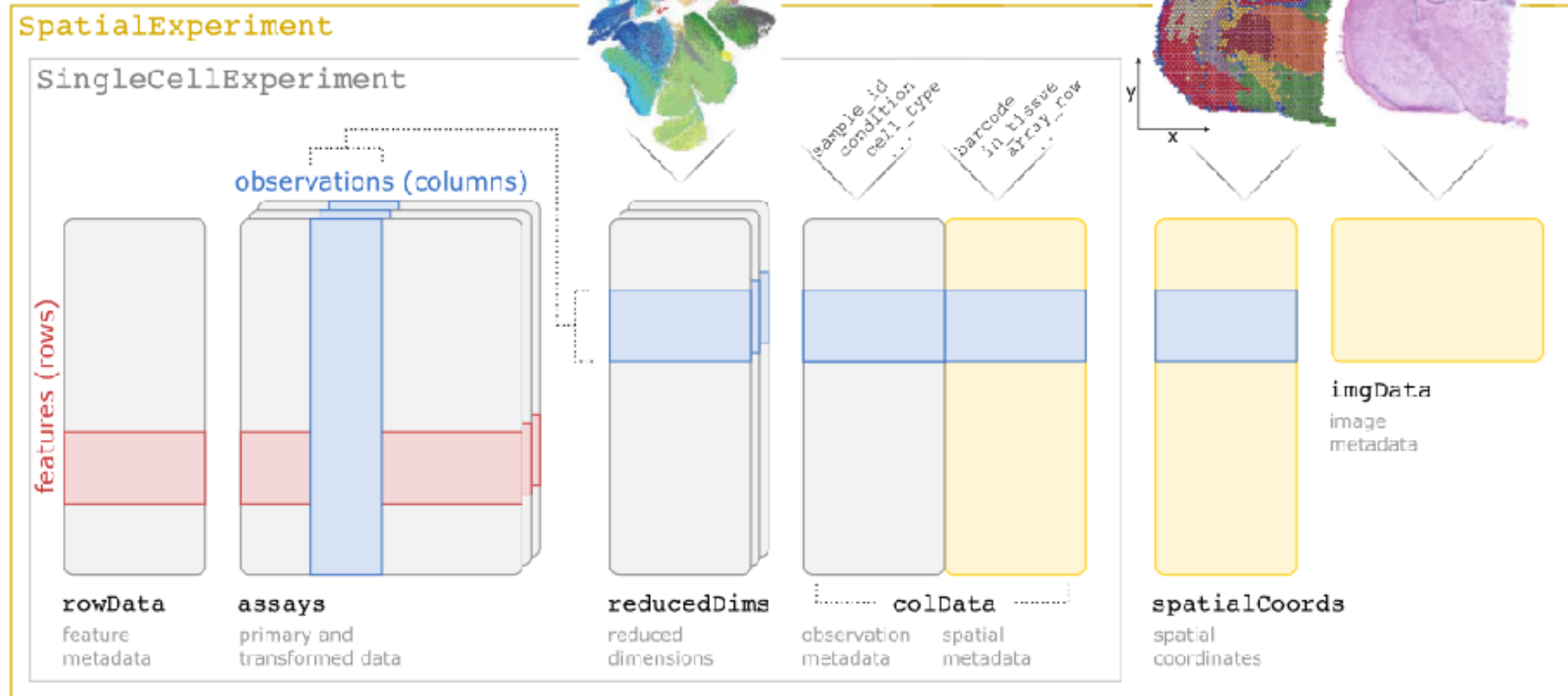


*Why? Because infrastructure around single-cell data is large!*



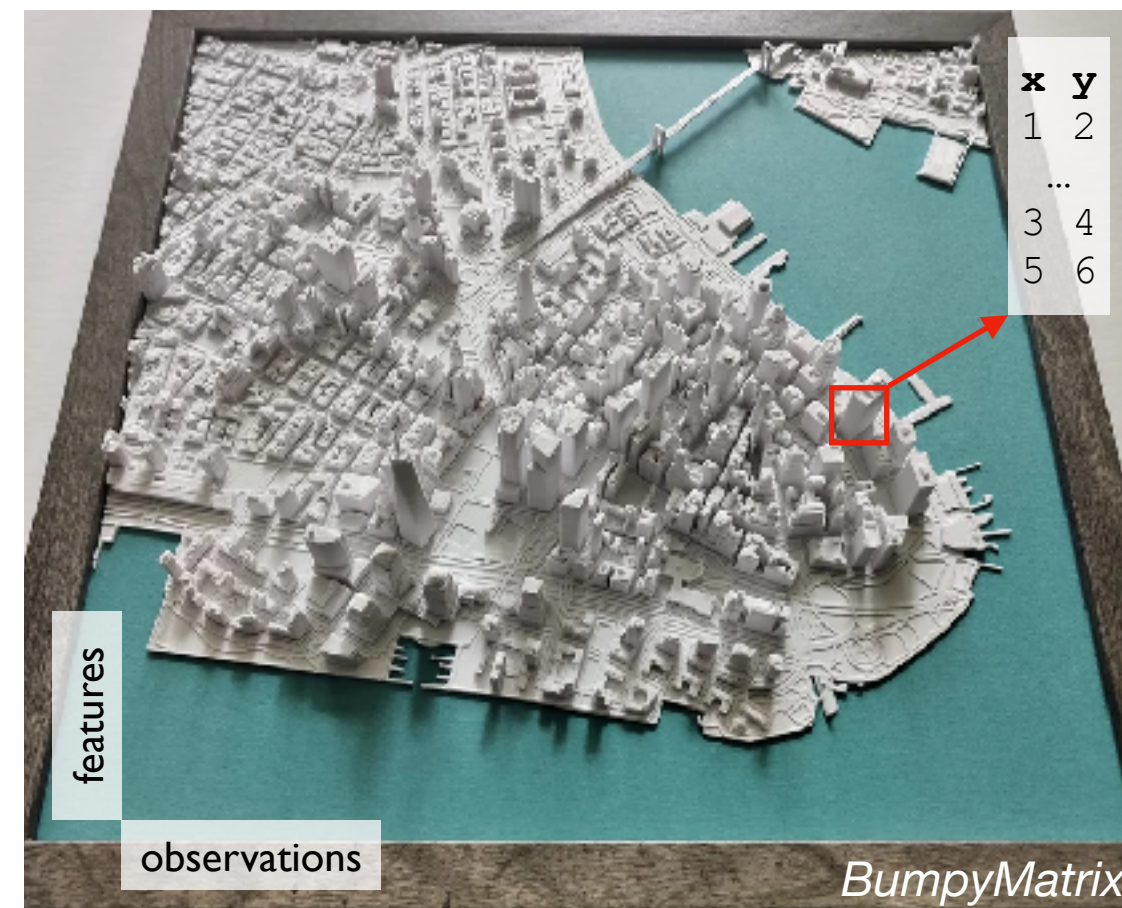


# infrastructure for handling img-ST data in R



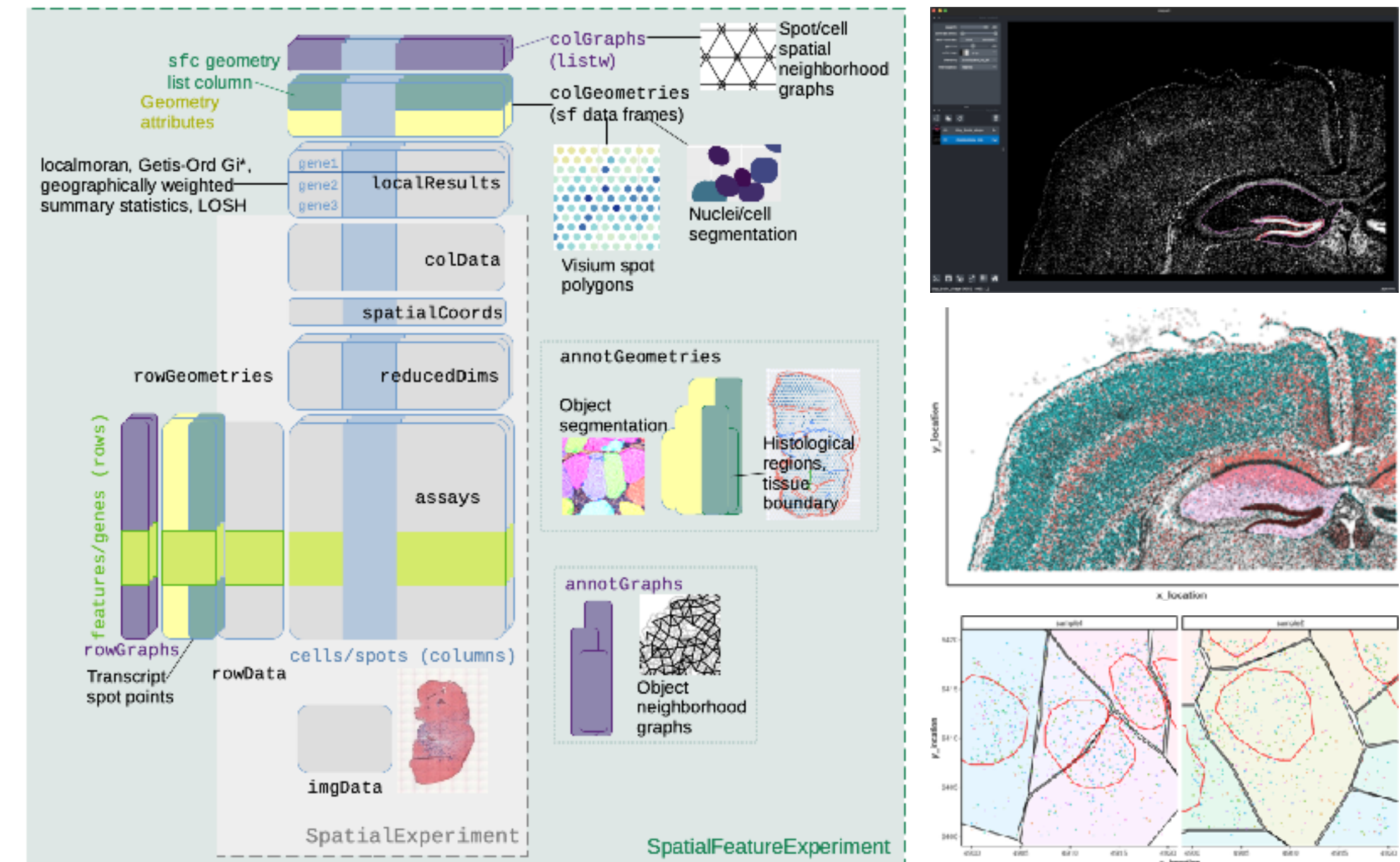
Righelli, Weber, Crowell et al. (2022)  
*Bioinformatics* 38(11):3128-3131

## SpatialExperiment



Moses et al. (2023) *bioRxiv* 2023.07.20.549945

## SpatialFeatureExperiment

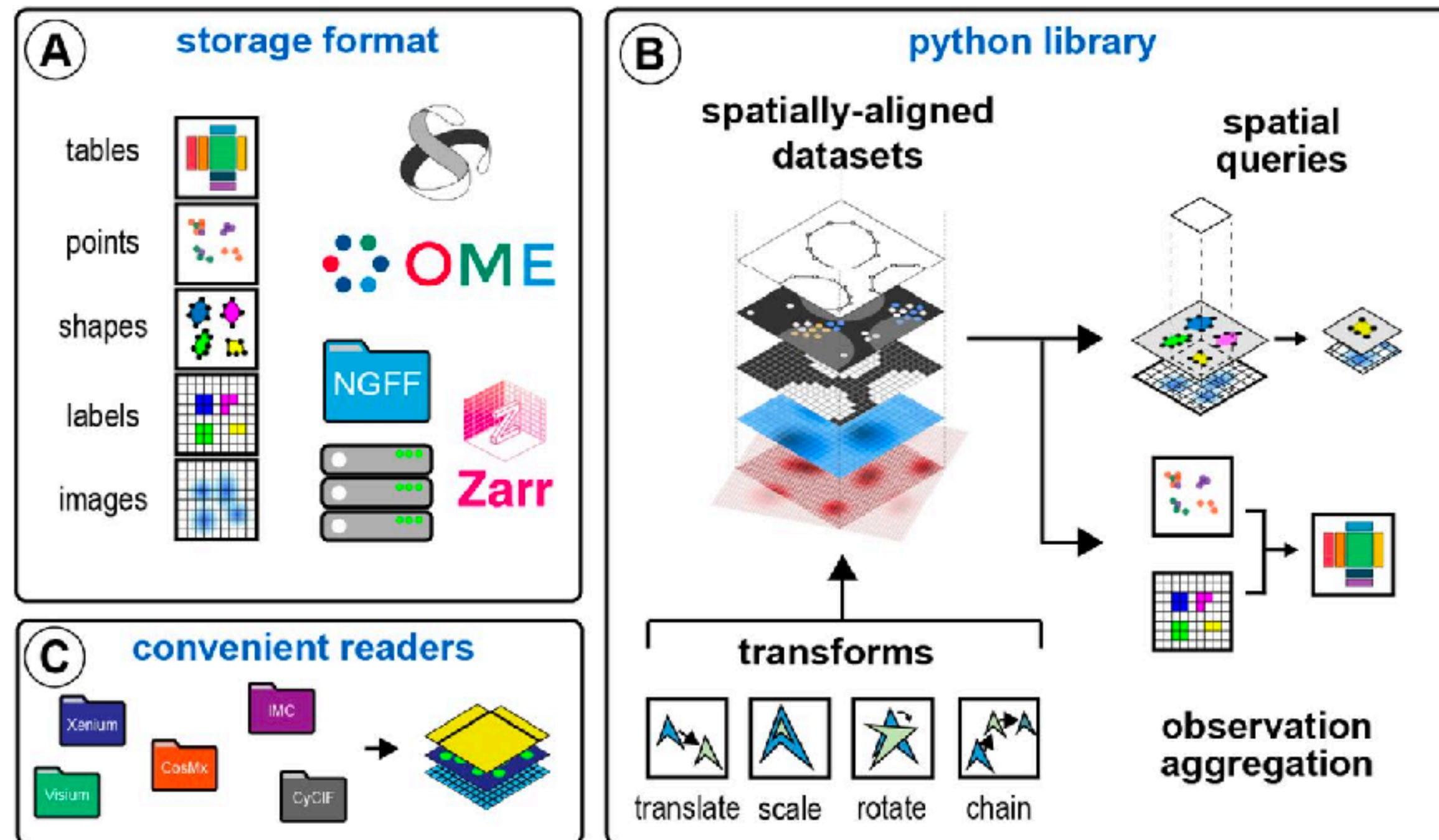


## MoleculeExperiment

Couto et al. (2023) *bioRxiv*  
 2023.05.16.541040

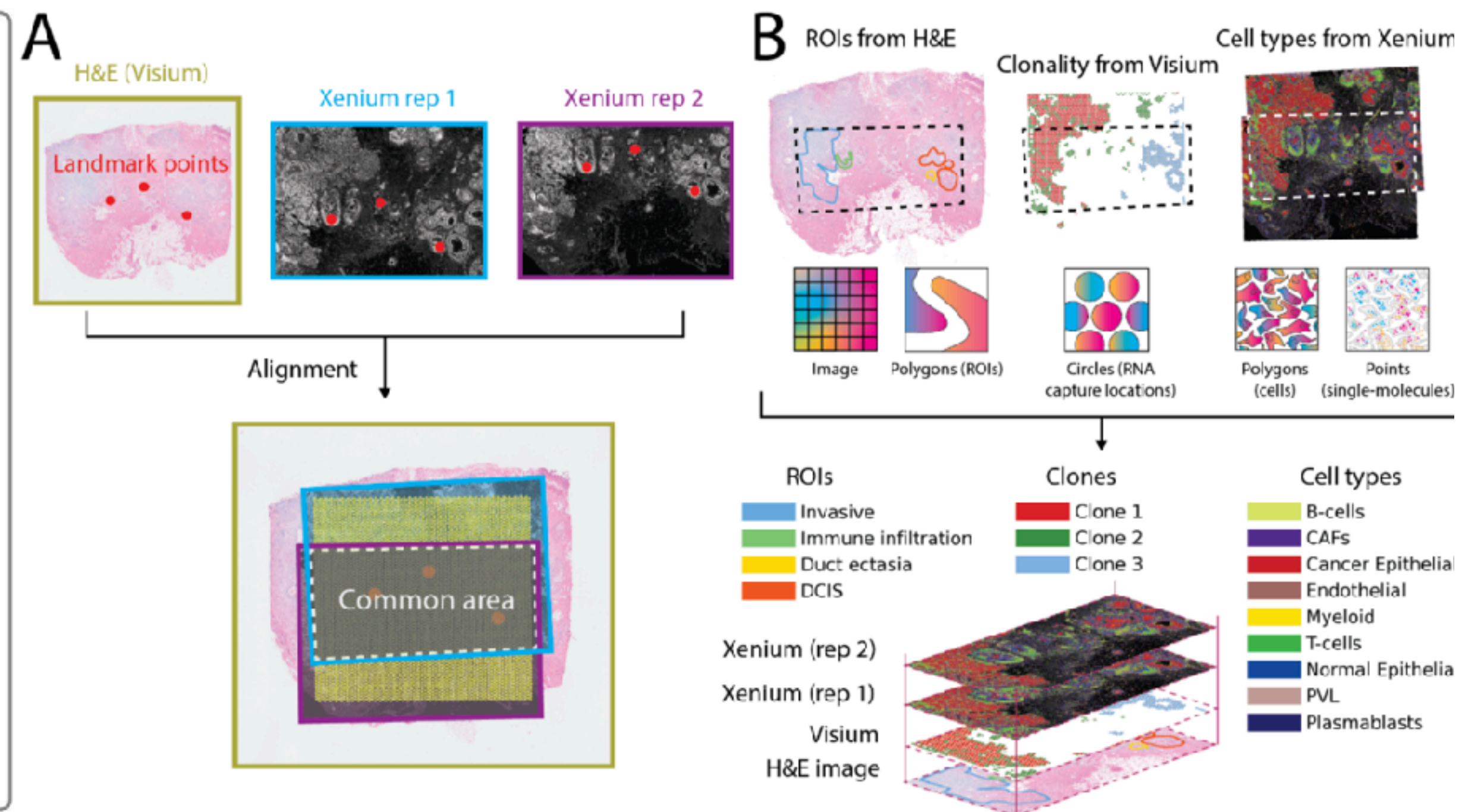
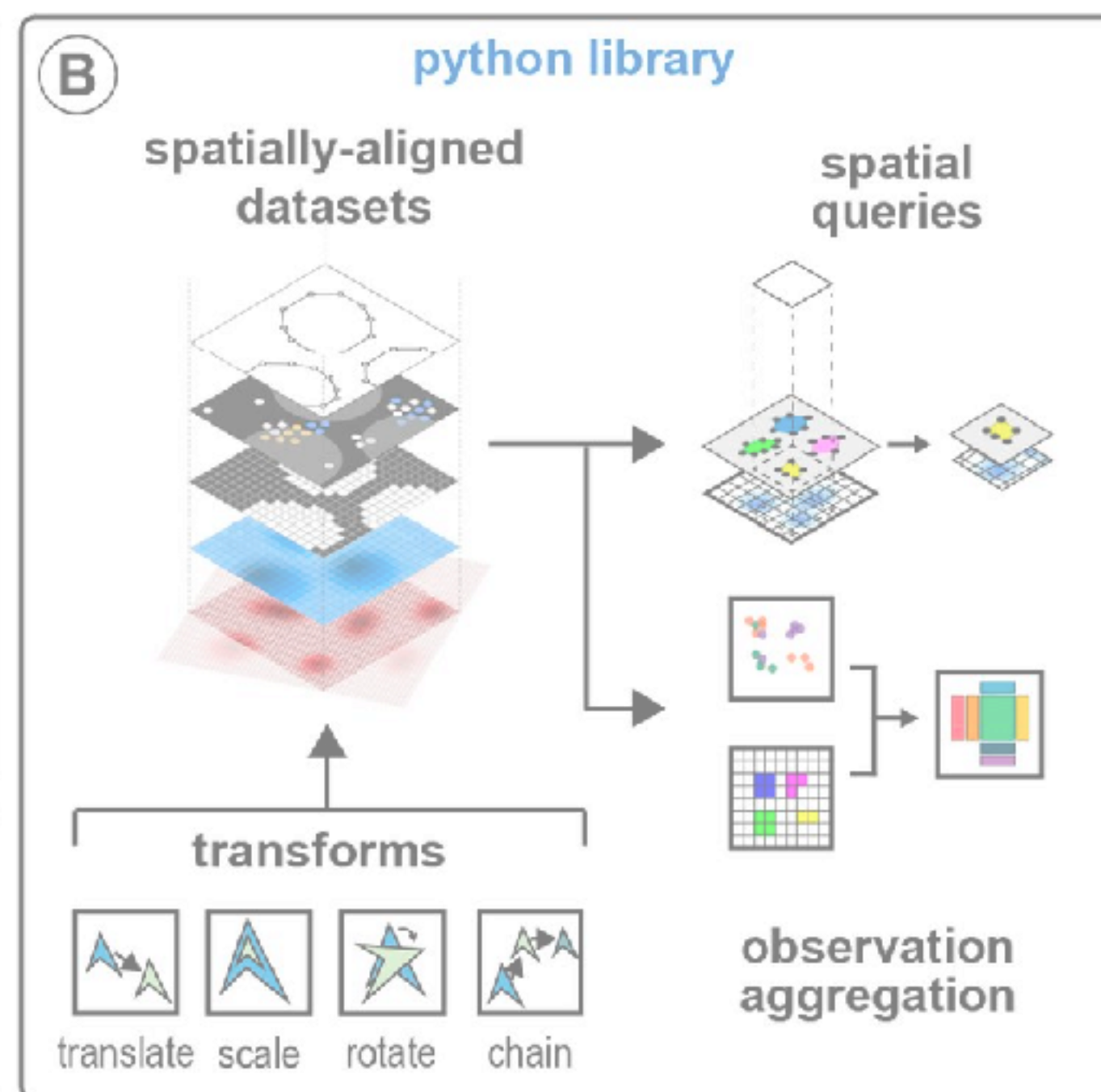


# infrastructure for handling img-ST data in Python





# infrastructure for handling img-ST data in Python





# appendix — references & resources

## technology

- He *et al.* (2022). High-plex Multiomic Analysis in FFPE at Subcellular Level by Spatial Molecular Imaging. *bioRxiv* [2021.11.03.467020](https://doi.org/10.1101/2021.11.03.467020)
- Khafizov *et al.* (2024). Sub-cellular imaging of the entire protein-coding human transcriptome (18933-plex) on FFPE tissue using SMI. *bioRxiv* [2024.11.27.625536](https://doi.org/10.1101/2024.11.27.625536)

## segmentation

- Stringer *et al.* (2021). Cellpose: a generalist algorithm for cellular segmentation. *Nature Methods* [18\(1\):100-106](https://doi.org/10.1038/s41592-017-0000-8)
- Petukhov *et al.* (2021). Cell segmentation in imaging-based spatial transcriptomics. *Nature Biotechnology* [40:345-354](https://doi.org/10.1038/s41587-021-0088-4)
- Wu *et al.* (2024). FastReseg: using transcript locations to refine image-based cell segmentation results in spatial transcriptomics. *bioRxiv* [2024.12.05.627051](https://doi.org/10.1101/2024.12.05.627051)
- Mitchel *et al.* (2024). Impact of Segmentation Errors in Analysis of Spatial Transcriptomics Data. *bioRxiv* [2025.01.02.631135](https://doi.org/10.1101/2025.01.02.631135)

## miscellaneous

- Martin *et al.* (2024). MerQuaCo: a computational tool for quality control in image-based spatial transcriptomics. *bioRxiv* [2024.12.04.626766](https://doi.org/10.1101/2024.12.04.626766)
- NanoString scratch space: <https://nanosting-biostats.github.io/CosMx-Analysis-Scratch-Space>

## normalization

- Bhuva *et al.* (2024): Library size confounds biology in spatial transcriptomics data. *Genome Biology* [25:99](https://doi.org/10.1186/s12864-024-10999-4)
- Atta *et al.* (2024). Gene count normalization in single-cell imaging-based spatially resolved transcriptomics. *Genome Biology* [25:153](https://doi.org/10.1186/s12864-024-10999-4)

## infrastructure

- Righelli, Weber, Crowell *et al.* (2022). SpatialExperiment: infrastructure for spatially resolved transcriptomics data in R using Bioconductor. *Bioinformatics* [38\(11\):3128-3131](https://doi.org/10.1093/bioinformatics/btad311)
- Couto *et al.* (2023). MoleculeExperiment enables consistent infrastructure for molecule-resolved spatial transcriptomics data in Bioconductor. *bioRxiv* [2023.05.16.541040](https://doi.org/10.1101/2023.05.16.541040)
- Moses *et al.* (2023). Voyager: exploratory single-cell genomics data analysis with geospatial statistics. *bioRxiv* [2023.07.20.549945](https://doi.org/10.1101/2023.07.20.549945)
- Marconato *et al.* (2024). SpatialData: an open and universal data framework for spatial omics. *Nature Methods* [s41592-024-02212-x](https://doi.org/10.1038/s41592-024-02212-x)

## benchmarks

- Wang *et al.* (2023). Systematic benchmarking of imaging spatial transcriptomics platforms in FFPE tissues. *bioRxiv* [2023.12.07.570603](https://doi.org/10.1101/2023.12.07.570603)
- Cook *et al.* (2023): A comparative analysis of imaging-based spatial transcriptomics platforms. *bioRxiv* [2023.12.13.571385](https://doi.org/10.1101/2023.12.13.571385)
- Rademacher *et al.* (2024): Comparison of spatial transcriptomics technologies using tumor cryosections. *bioRxiv* [2024.04.03.586404](https://doi.org/10.1101/2024.04.03.586404)
- Ren *et al.* (2024). Systematic Benchmarking of High-Throughput Subcellular Spatial Transcriptomics Platforms. *bioRxiv* [2024.12.23.630033](https://doi.org/10.1101/2024.12.23.630033)