Deep learning and computational methods for spatial data analysis

Quentin Blampey

AI & ML in Spatial Single-Cell Transcriptomics Workshop 16 Oct 2025 - Lyon (France)

CentraleSupélec, Paris-Saclay University
Gustave Roussy Institute



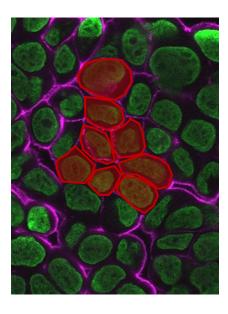


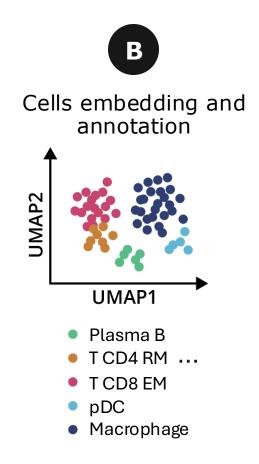


Typical tasks in spatial omics analysis



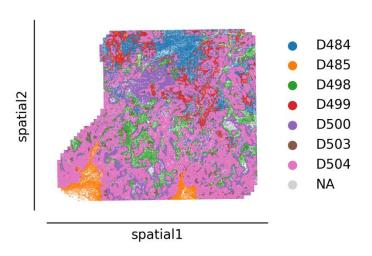
Cell segmentation







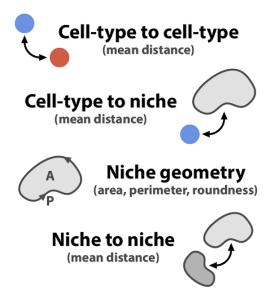
Spatial clustering (also known as spatial domains)



Typical tasks in spatial omics analysis

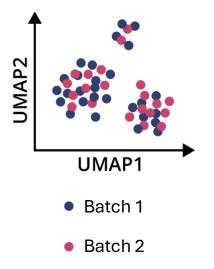


Spatial geometrical statistics and distances





Domains batch-effect correction



Motivation

Building analysis tools for biological research, within the



scverse community



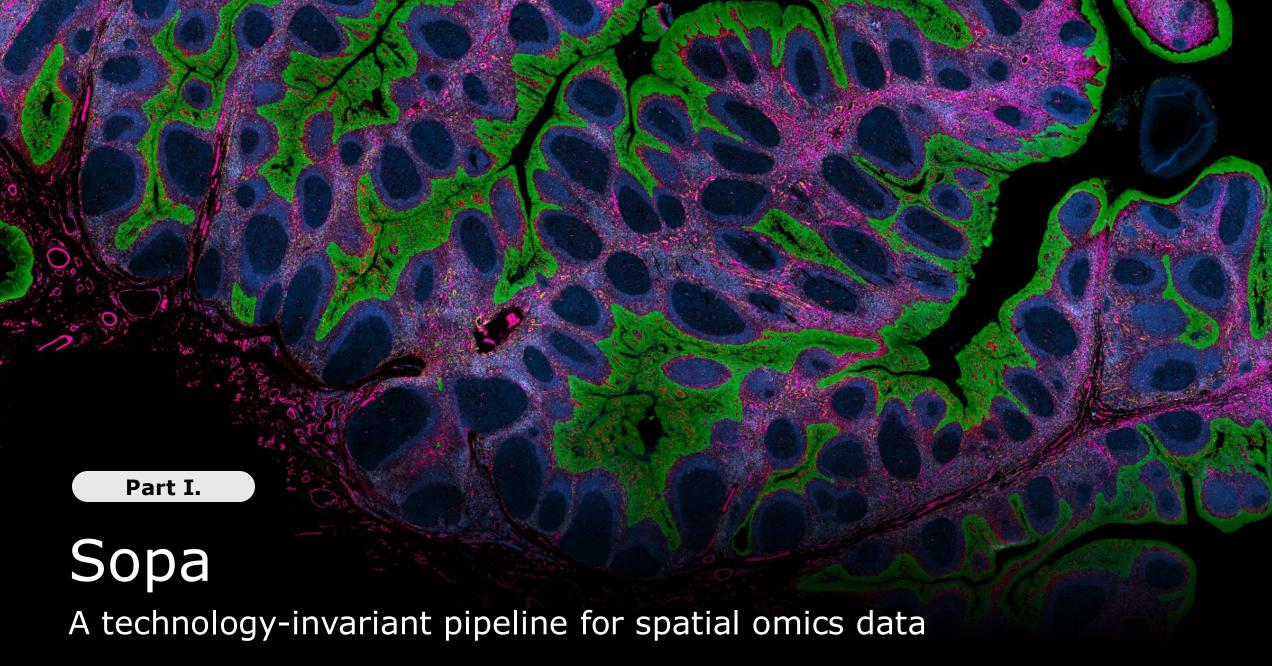
SpatialData is a data structure, not an analysis library

General spatial toolkit

Cell segmentation, tissue segmentation, interactive visualization, efficient aggregation, spatial operations, and many other features

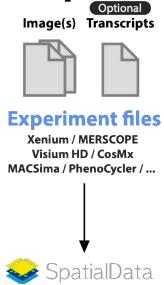
Foundation model

Spatial domains (or niches), hierarchy of domains, domains batch-effect correction, slide architecture analysis, pathway spatial patterns



Blampey et al., 2024, *Nature Communications*

Overview of Sopa



RAM and time efficiency, compared to naive approaches

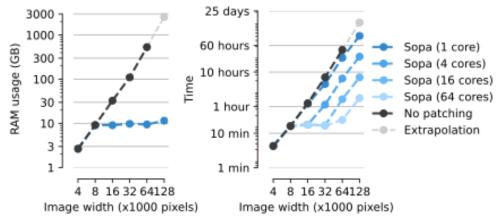
Cellpose segmentation 25 days 10000 Sopa (1 core) RAM usage (GB) 3000 60 hours Sopa (4 cores) 1000 10 hours Sopa (16 cores) 300 Sopa (64 cores) 100 1 hour No patching 30 Extrapolation 10 min 10 3 1 min

4 8 16 32 64128

Larger datasets

Image width (x1000 pixels)

Baysor segmentation 25 days 3000 1000

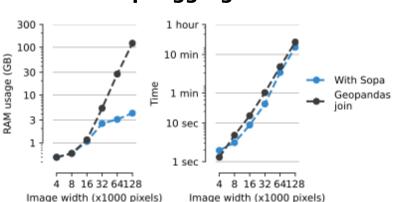


Transcript aggregation

4 8 16 32 64128

Image width (x1000 pixels)

Larger datasets



Channel aggregation

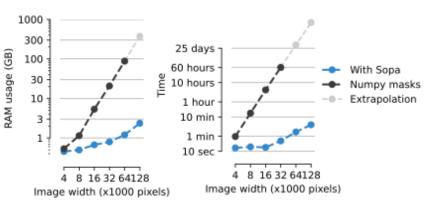
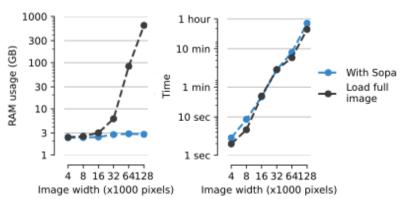
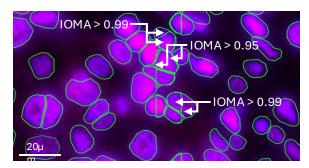


Image on-disk writing

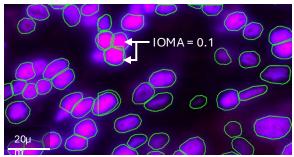


Regarding segmentation conflicts

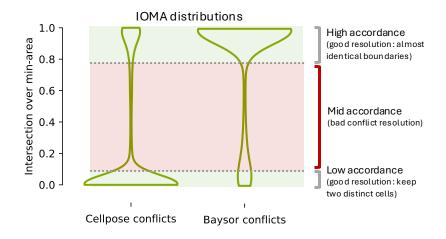
The distribution of the intersection area is close to 0 and 1

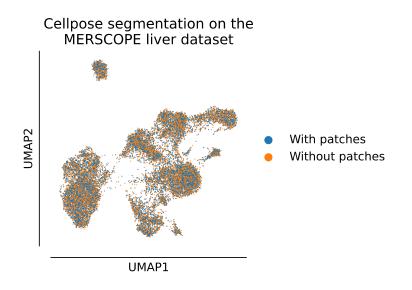


Cells on overlapping patches are segmented twice, the corresponding boundaries are merged



The boundaries of two different cells may slightly overlap (cells kept as distinct)

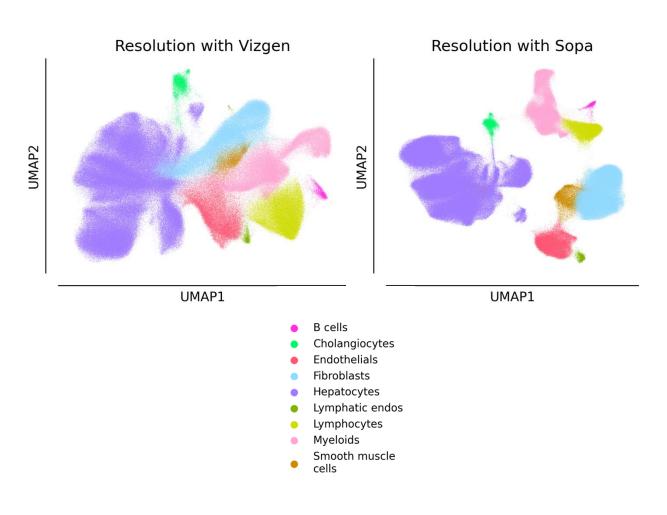


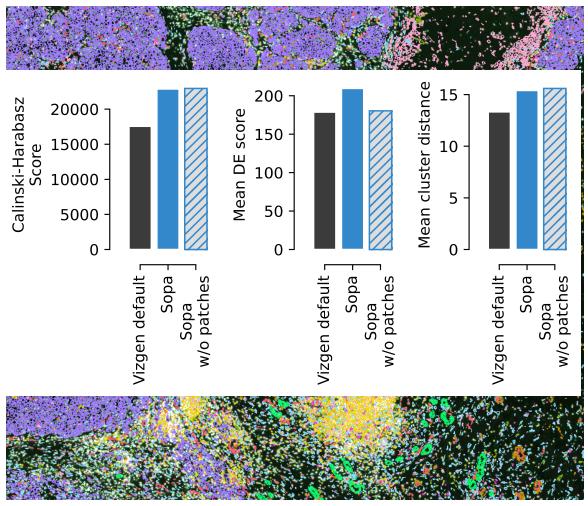


Improved segmentation quality

Comparison: Vizgen segmentation (cellpose-based) and Sopa segmentation (baysor-based)

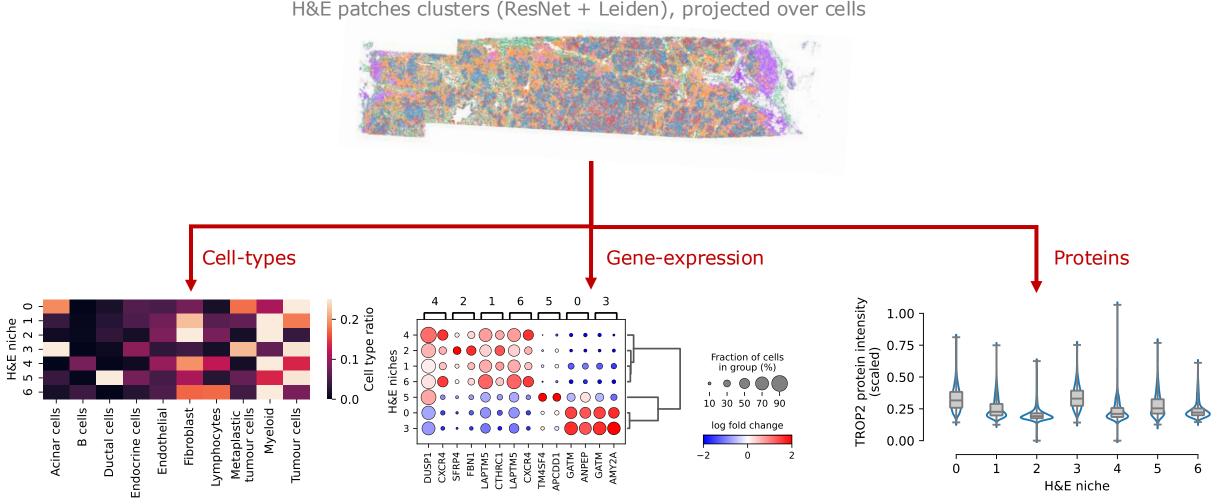
Note: Baysor can't be run without Sopa due to RAM usage





Multi-omics analysis examples

Spatial operations to relate multiple spatial modalities (H&E, RNA, Protein)



Examples using the API

Usage on Xenium data

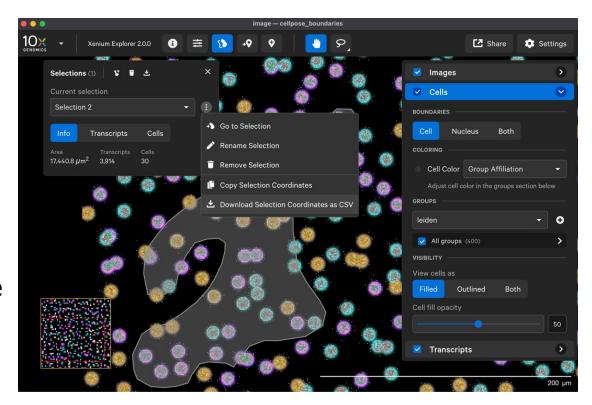
```
import sopa
sdata = sopa.io.xenium("/path/to/data_directory")
sopa.make_image_patches(sdata)
sopa.segmentation.cellpose(sdata, channels=["DAPI"], diameter=35, gpu=True)
sopa.make_transcript_patches(sdata, patch_width=None, prior_shapes_key="cellpose_boundaries")
sopa.segmentation.proseg(sdata)
sopa.aggregate(sdata)
sopa.io.explorer.write(sdata, mode="-it") # update the Xenium Explorer
sdata.pl.render_shapes().pl.show() # with spatialdata-plot
```

Xenium data and Xenium Explorer interoperability



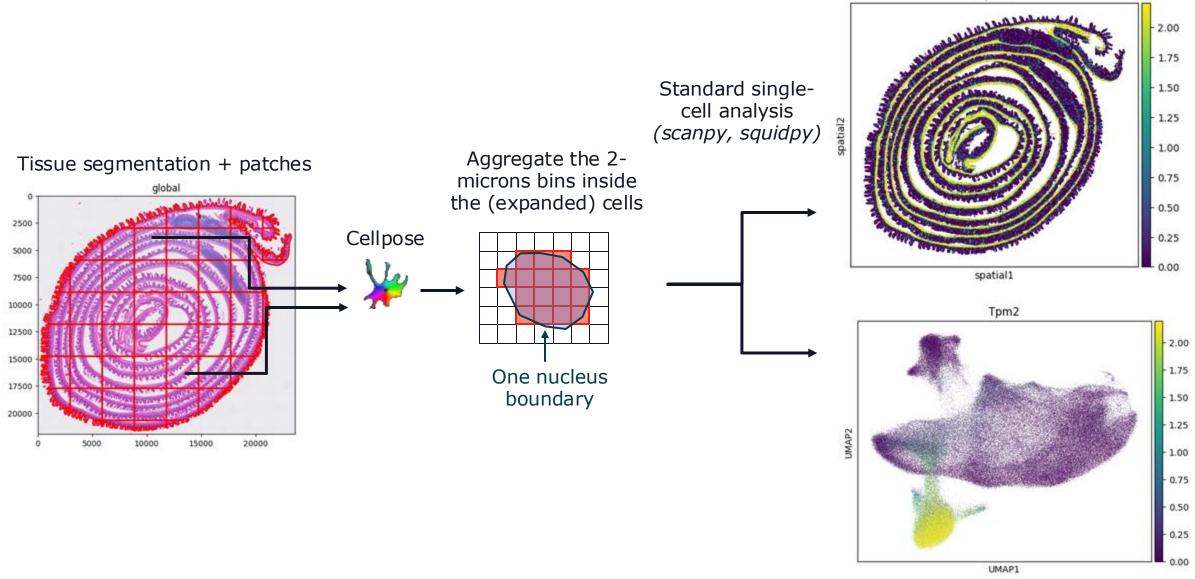
Add the analysis to the explorer

Select cells of interest and analyze them



Tpm2

Approach on Visium HD data



But the API is still very similar

Usage on Visium HD data

```
import sopa
sdata = sopa.io.visium_hd("/path/to/data_directory")
sopa.make_image_patches(sdata)
sopa.segmentation.stardist(sdata)
sopa.segmentation.proseg(sdata, prior_shapes_key="stardist_boundaries")
sopa.aggregate(sdata)
```

Other usage modes

Command line interface (CLI)

```
> sopa --help # show command names and arguments
> sopa convert merscope_directory --technology merscope # read some data
> sopa patchify image merscope_directory.zarr # make patches for low-memory segmentation
> sopa segmentation cellpose merscope_directory.zarr --diameter 60 --channels DAPI # segmentation
> sopa resolve cellpose merscope_directory.zarr # resolve segmentation conflicts at boundaries
> sopa aggregate merscope_directory.zarr --average-intensities # transcripts/channels aggregation
> sopa explorer write merscope_directory.zarr # convert for interactive vizualisation
```







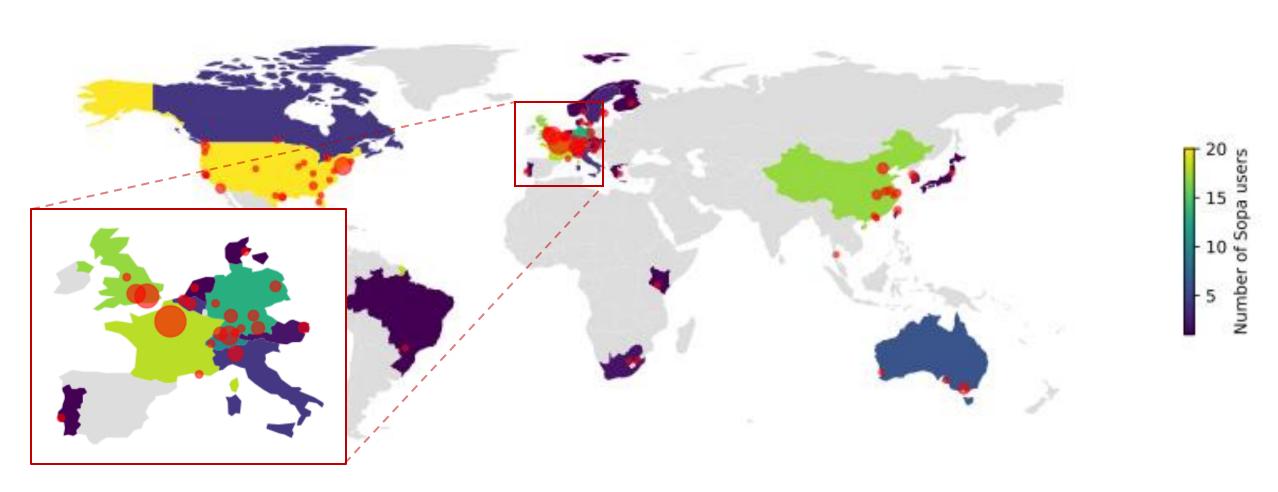






Usage and impact of Sopa (in Nov 2024, i.e., after 12 months)

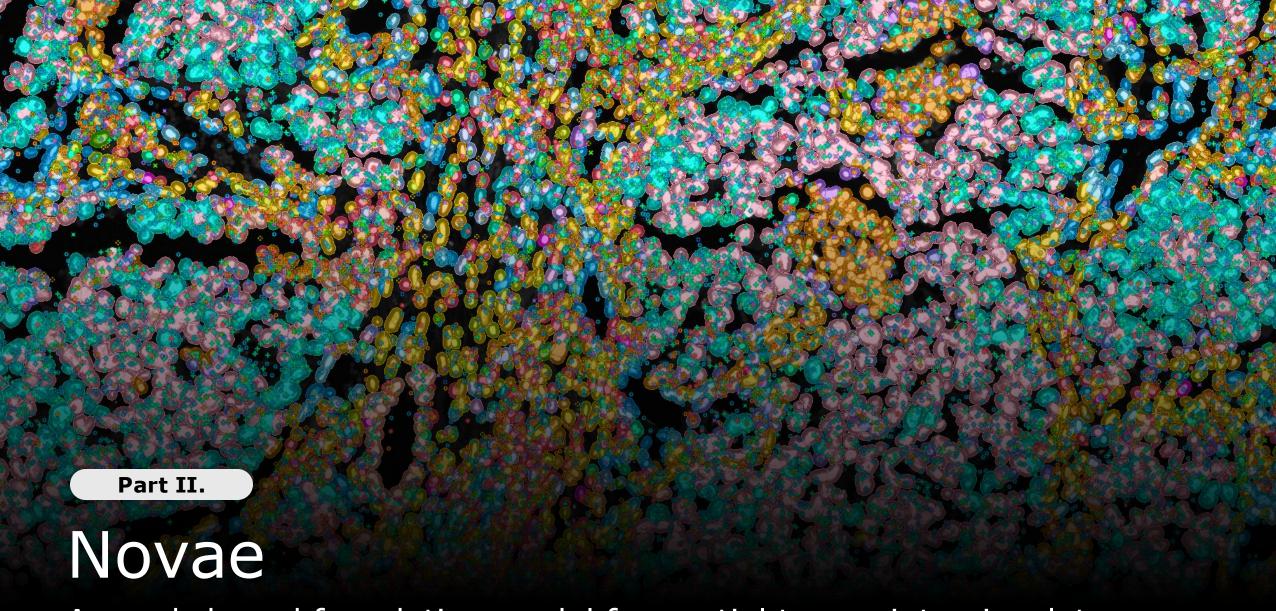
Dots represent researchers who gave a "star" to the repository, or opened an issue/PR



Many other features were not shown, try it out!

https://gustaveroussy.github.io/sopa/

https://github.com/gustaveroussy/sopa



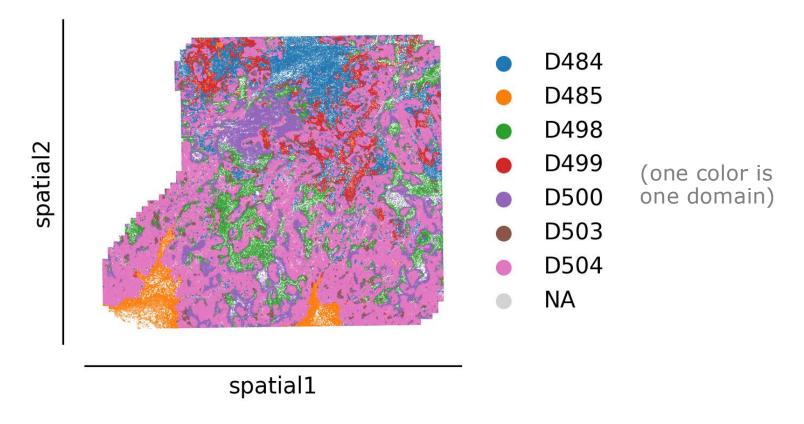
A graph-based foundation model for spatial transcriptomics data

Blampey et al., 2024, accepted in *Nature Methods*

Definition of a spatial domain (or niche)

A spatial domain (a.k.a. niche) is a localized cellular microenvironment within tissues.

> These domains reflect distinct biological functions and interactions



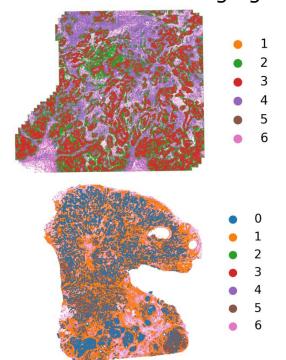
Which implications for tissue function / disease progression?

Main challenges in spatial domains analysis

Existing tools have limitations (STAGATE, SEDR, SpaceFlow, GraphST)

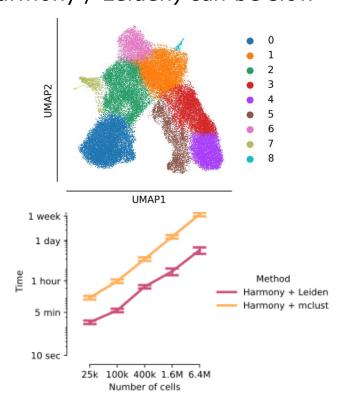


Identifying cross-slide spatial domains is challenging



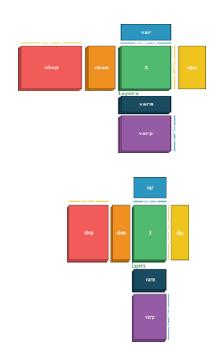
2

Depending on external tools (e.g. Harmony / Leiden) can be slow



3

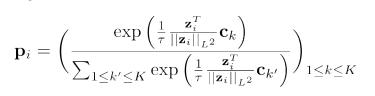
Many studies are composed of different gene panels (technologies evolve)

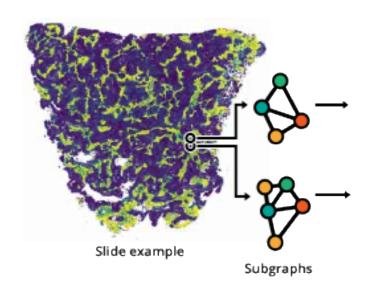


Method overview

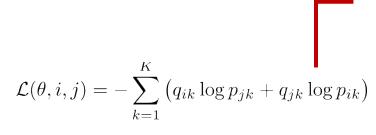
Self-supervision model (inspired from the SwAV framework)

$$\mathbf{Q}^* := \operatorname{argmin}_{\mathbf{Q} \in \mathcal{Q}} \Big(\operatorname{Tr}(\mathbf{Q} \mathbf{C} \mathbf{Z}^\top) - \epsilon H(\mathbf{Q}) \Big)$$



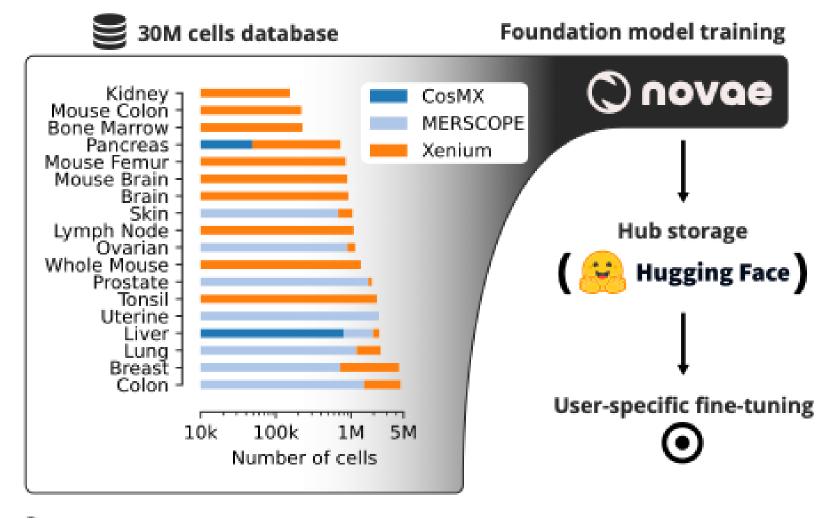


embed(
$$\mathbf{x}_i, \mathcal{P}$$
) = $\frac{\sum_{j=1}^{P} x_{ij}(\mathbf{W} \mathbf{v}_{\mathcal{P}[\mathbf{j}]} + \mathbf{b})}{\sqrt{\sum_{j=1}^{P} (\mathbf{W} \mathbf{v}_{\mathcal{P}[\mathbf{j}]} + \mathbf{b})^2}}$

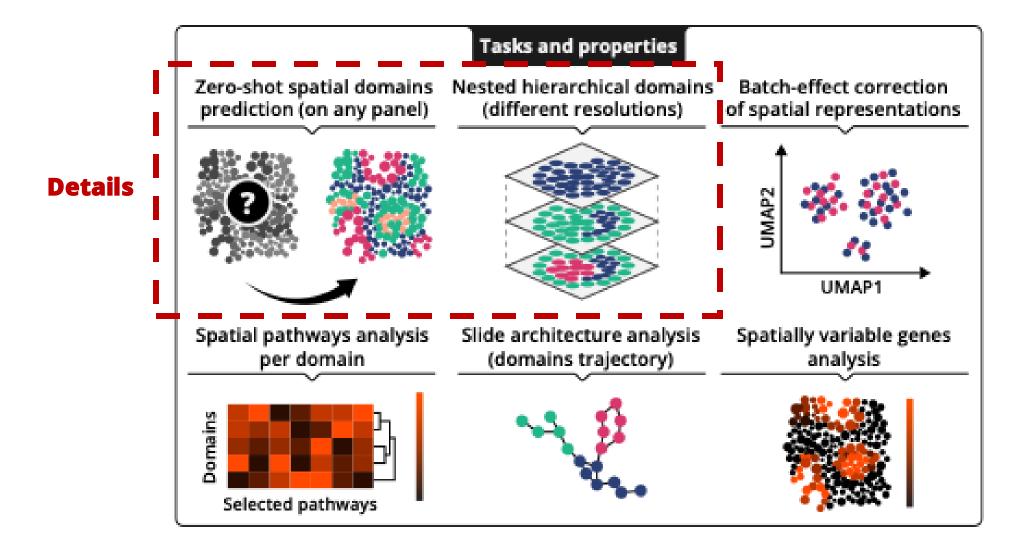


Training and sharing (single-cell resolution dataset)

Public database collected from single-cell resolution technologies

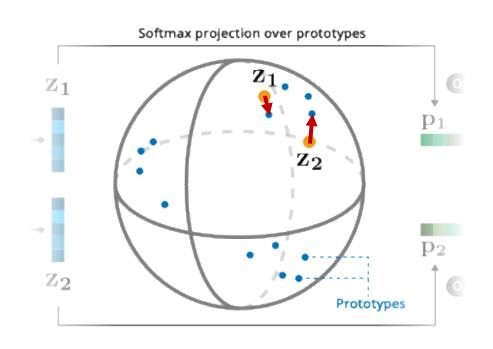


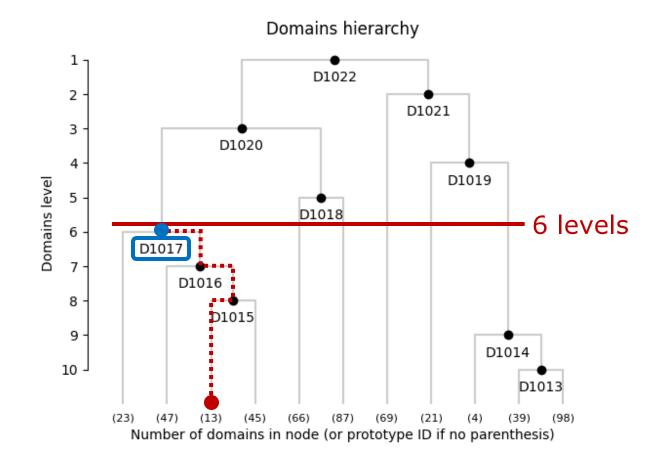
Tasks and properties of Novae



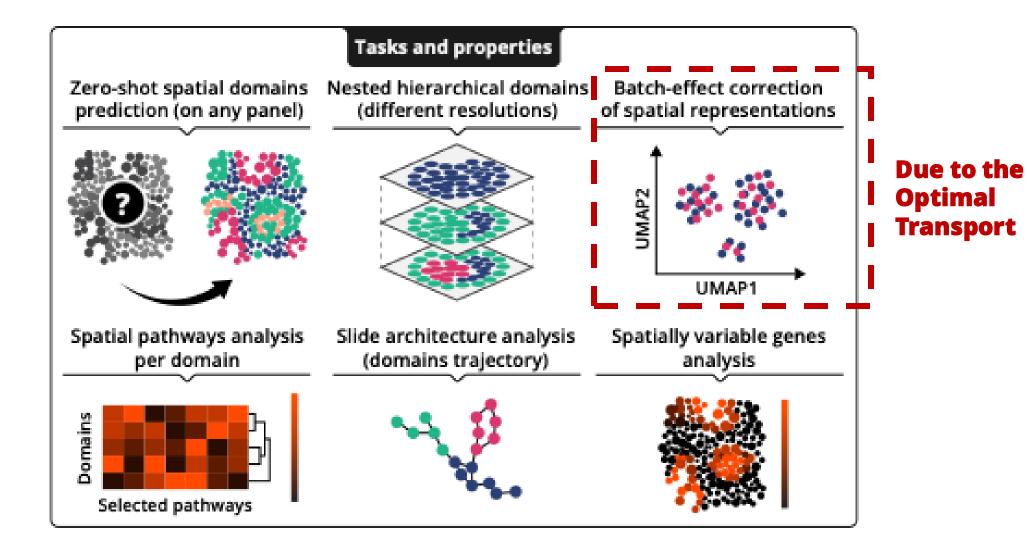
A hierarchy of prototypes

Hierarchical clustering of the prototypes, and assignment to the closest prototype

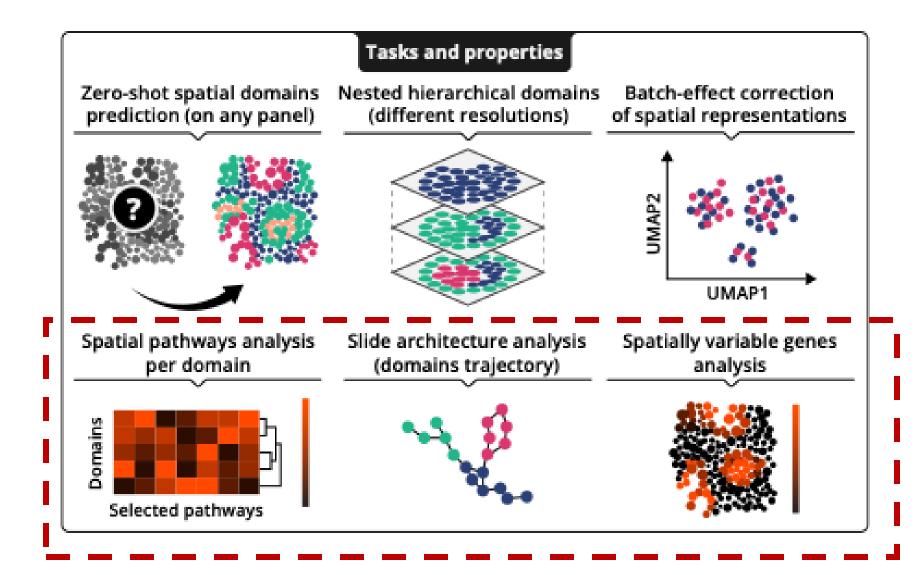




Tasks and properties of Novae



Tasks and properties of Novae

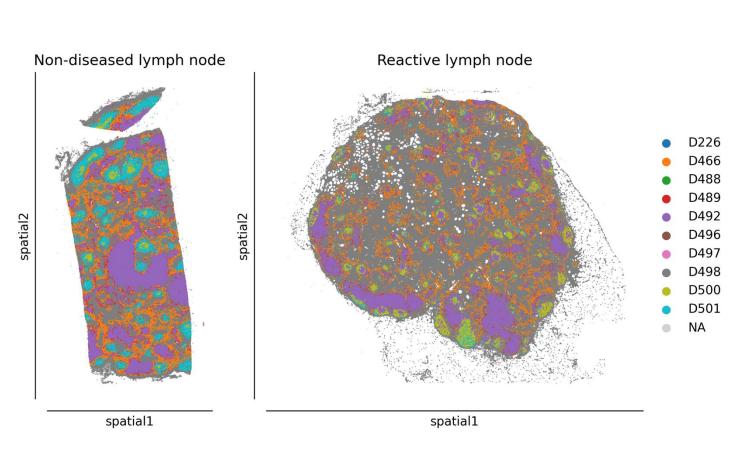


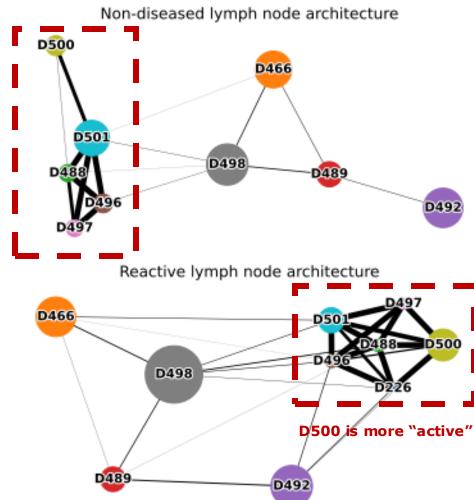
Multiple examples

[Example 1] Architecture of a reactive lymph node

I. Sopa: Technology-invariant pipeline for spatial omics

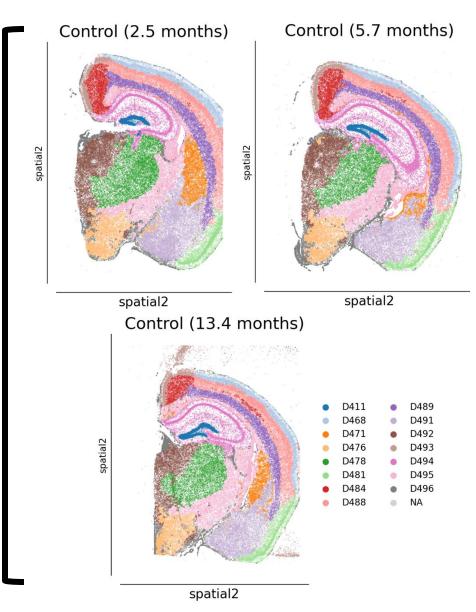
Comparison between a non-diseased and a reactive lymph node



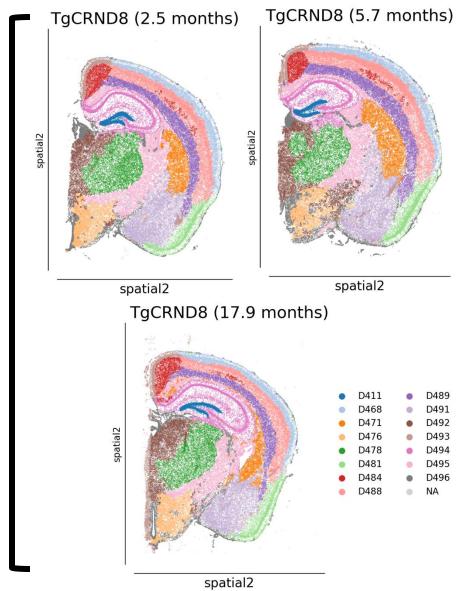


[Example 2] Mice brain slides with Alzheimer-like pathology

Control Mice (3 time points)



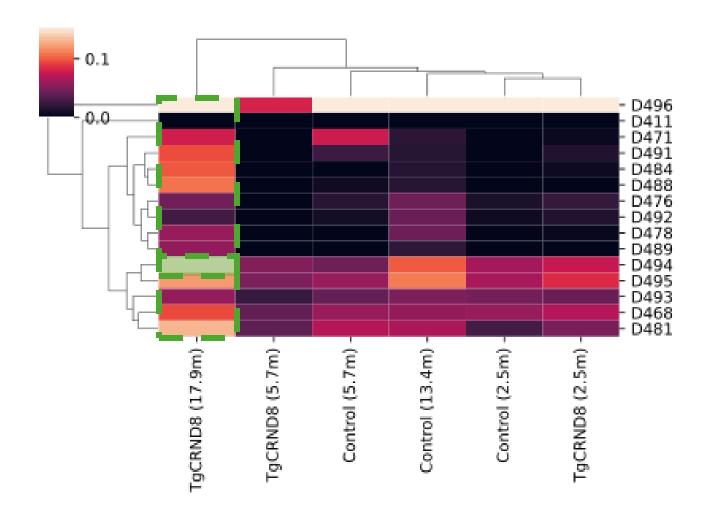
TgCRND8 Mice (3 time points)

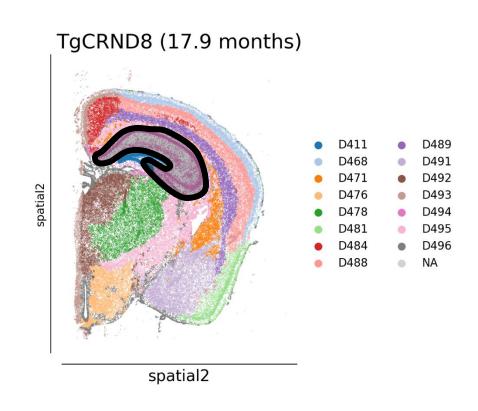


[Example 2] Mice brain slides with Alzheimer-like pathology

I. Sopa: Technology-invariant pipeline for spatial omics

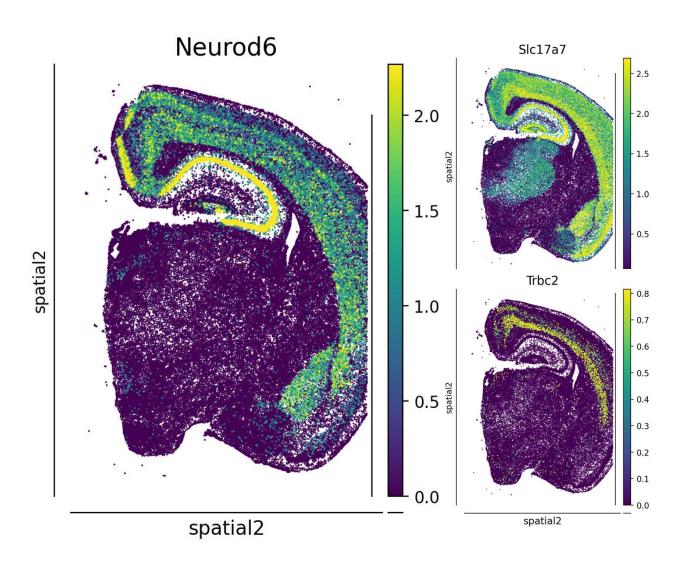
A brain aging pathway is up-regulated in a specific domains of the TgCRND8 17m mouse

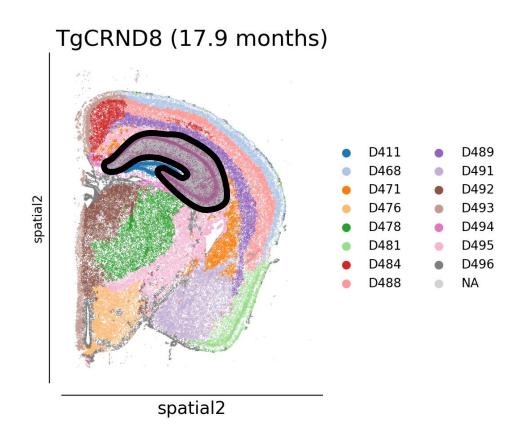




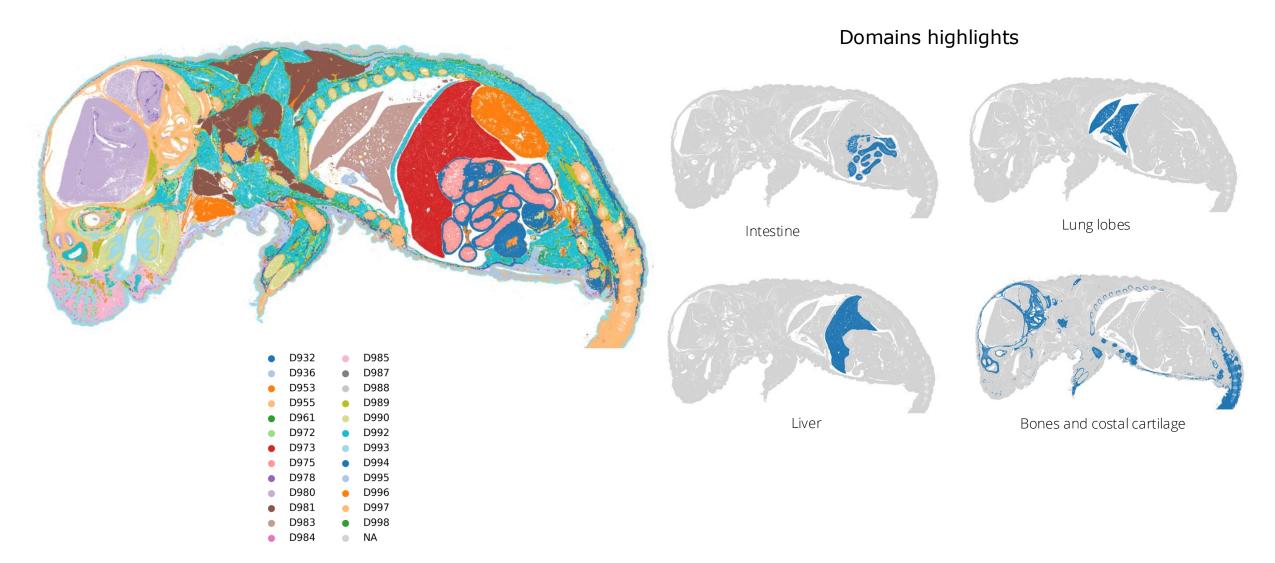
[Example 2] Mice brain slides with Alzheimer-like pathology

Identification of Spatially Variable Genes (DEGs over spatial domains)





[Example 3] Whole mouse spatial domains

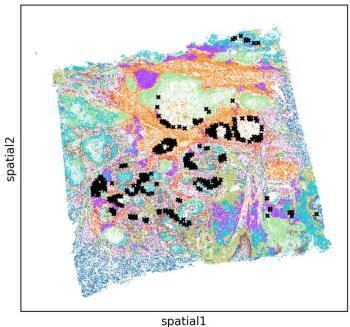


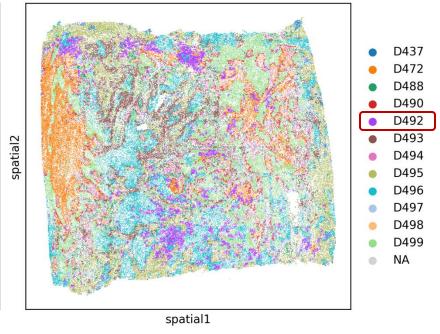
[Example 4] Spatial domains of MGCs in head & neck cancer

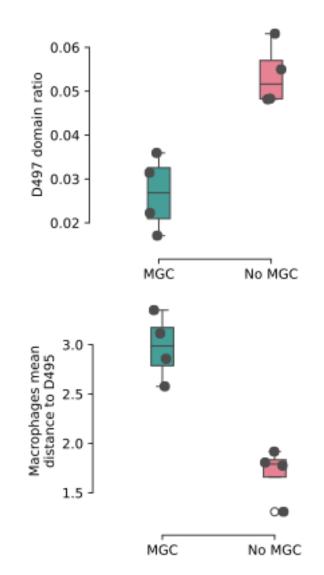
Usage on CosMx Protein Assays

8 patients stratified based on the presence of multinucleated giant cells (MGCs), which are good prognosis

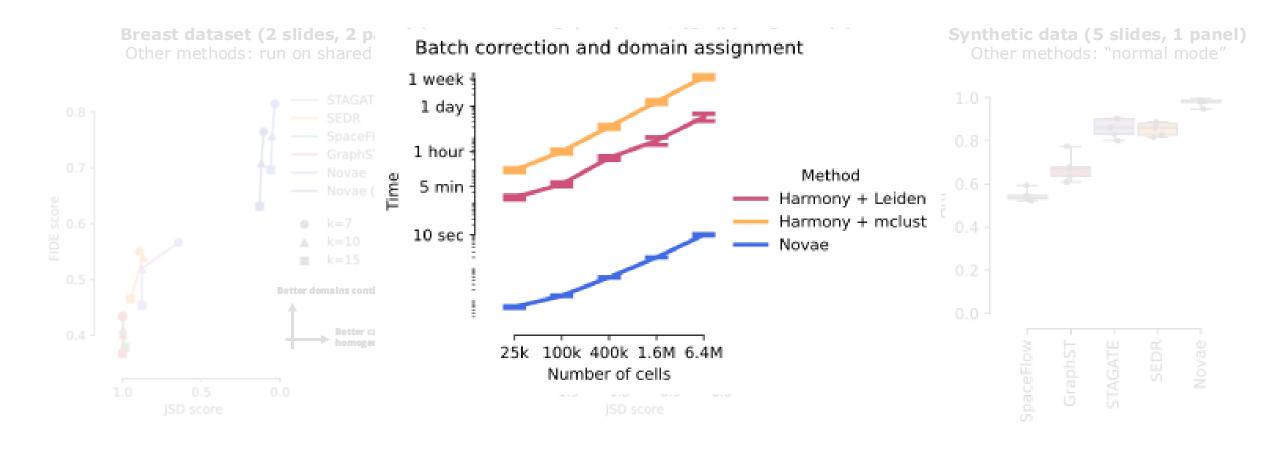
Novae domains (slide with MGCs) Novae domains (slide without MGCs)





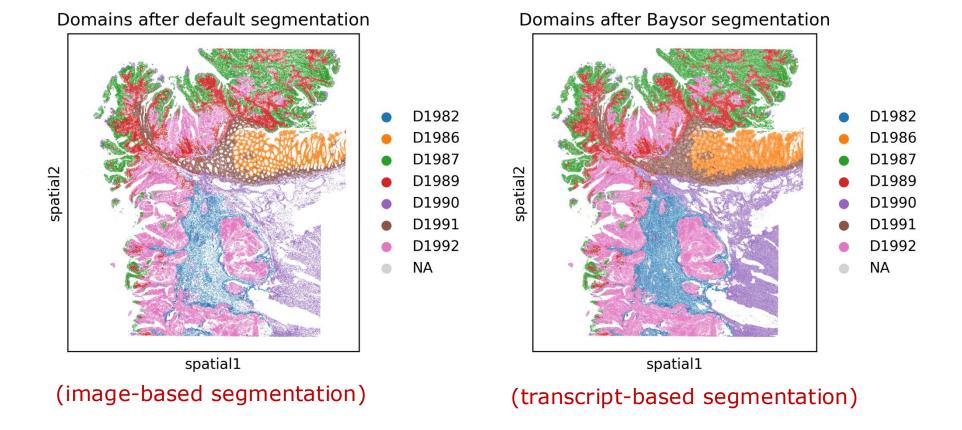


Benchmark - Comparison to existing methods



Robustness analysis: sensitivity to segmentation methods

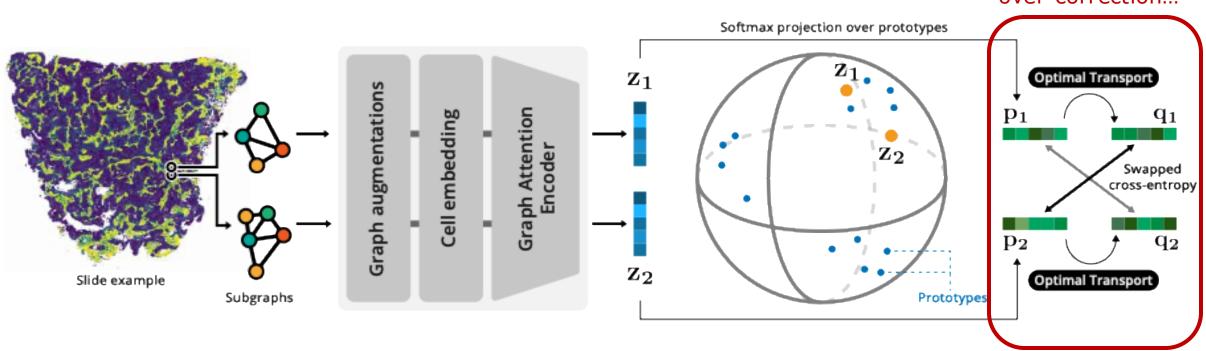
Approach: run two segmentation methods, and show we can retrieve the same domains



Robustness analysis 2: avoiding over-correction

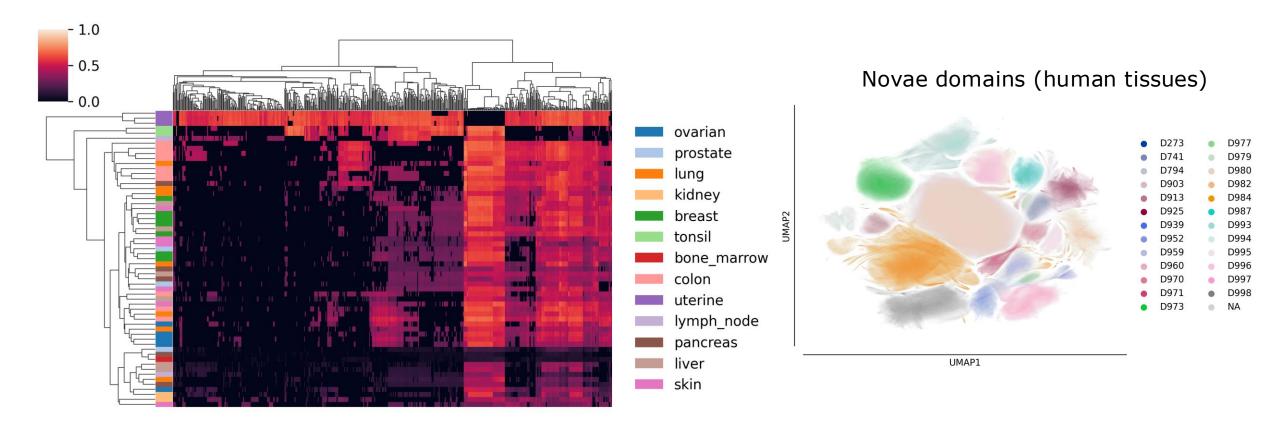
I. Sopa: Technology-invariant pipeline for spatial omics

Reminder: we use optimal transport to correct the batch effect



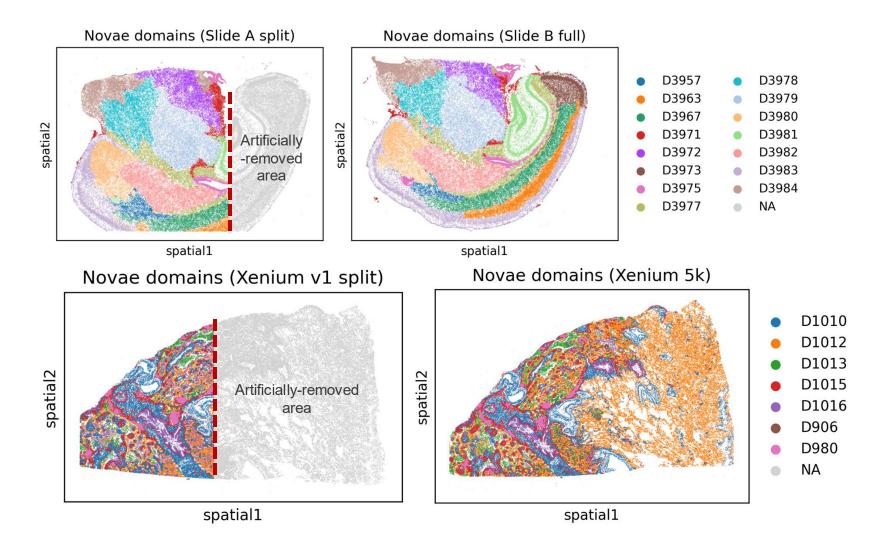
Robustness analysis 2: avoiding over-correction

We allow the prototypes not to be used in every slide



Robustness analysis 2: avoiding over-correction

Removing some parts of the slides to see if it affects the spatial domains assignment



Example usage (API)

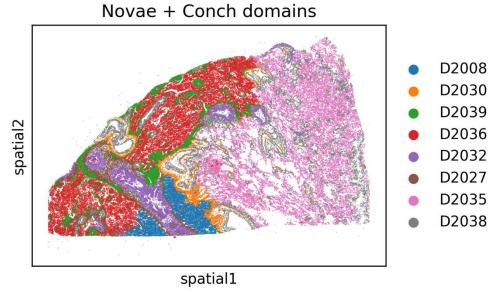
Zero-shot assignment of spatial domains on some slides from our dataset

```
import novae
adatas = novae.load_dataset(tissue="colon", species="human", pattern=".*P2.*")
novae.spatial_neighbors(adata, radius=80)
model = novae.Novae.from_pretrained("MICS-Lab/novae-human-0")
model.compute_representations(adata, zero_shot=True)
model.assign domains(adata)
```

Combining ST information with H&E data

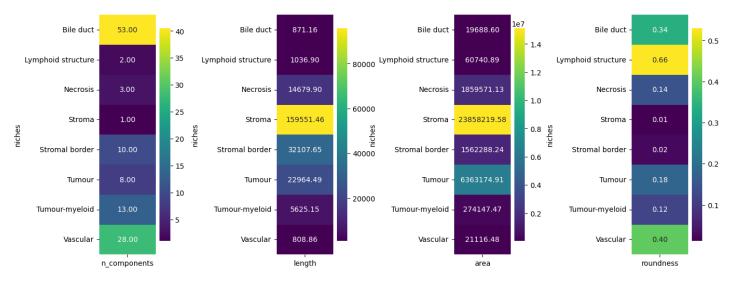
Use case: aligning modalities, then using an H&E foundation model to get embeddings per subgraph, and train a new Novae model in multimodal mode

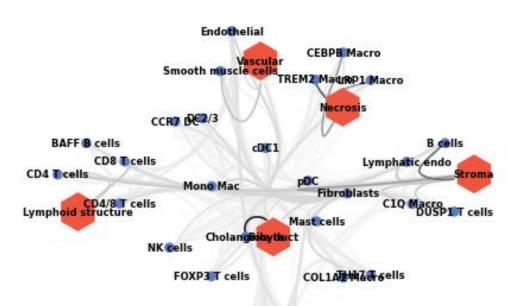






Go back to Sopa to compute some spatial statistics



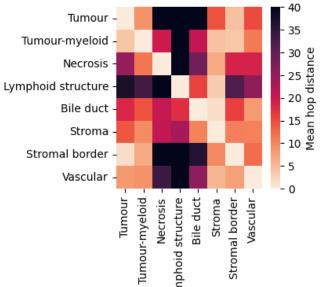


Hepata Wolfr-myeloid

Stromal border

Kupffer cells

SERPINAL MACKS



Or other scverse tools...



Again, feel free to try it!

https://mics-lab.github.io/novae/

https://github.com/MICS-Lab/novae

Acknowledgments

CentraleSupélec

Paul-Henry Cournède

Hakim Benkirane

Karmen Rabar

Stergios Christodoulidis



Institut Gustave Roussy

Fabrice André

Kevin Mulder

Florent Ginhoux

Charles-Antoine Dutertre

Isabelle Pic

Margaux Gardet

Joana Mourato Ribeiro





Institut Cochin [External]

Nadège Bercovici



EMBL / DKFZ [External]

Luca Marconato

Wouter-Michiel Vierdag





Helmholtz Munich [External]

Giovanni Palla

