

eggplant - building common coordinate frameworks for spatial transcriptomics

data

Alma Andersson
SCOG Workshop | May 24th 2022



@aalmaander



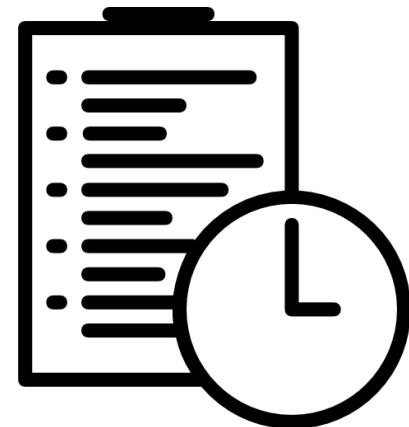
alma.andersson@differentiable.net



differentiable.net

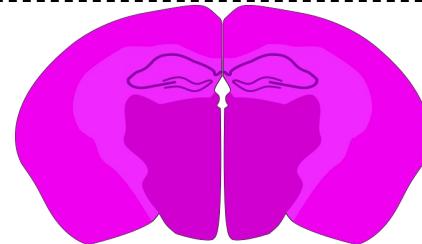
Agenda

- ⇒ **Introducing the problem** | explaining why we need a method like this
- ⇒ **Presenting a solution** | outlining how our approach works
- ⇒ **Results** | 3 examples highlighting different aspects of *eggplant*
- ⇒ **Questions** | addressing questions relating to the presentation
- ⇒ **Notebook walkthrough** | getting familiar with *eggplant*

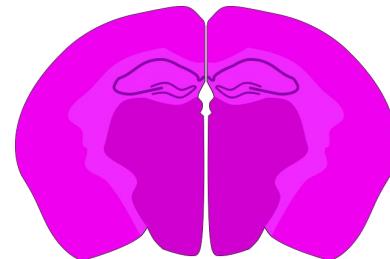


Setting the scene

- **Scenario:** we have collected tissue samples **representing the same structure** from multiple sources:
 - different conditions
 - time points
 - replicates
- **Objective:** **compare and relate** features between our samples
- **Example projects:** Atlas efforts (e.g., HCA)

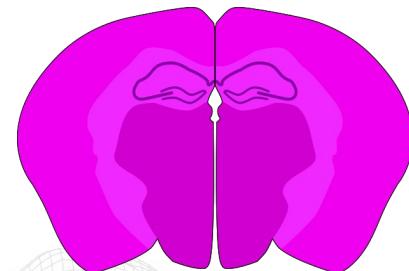


Individual 1



Individual 2

⋮



Individual N



The strategy

non-spatial data

Employ strategies for **batch correction** or
data integration.

Create a **low-dimensional embedding**
based on gene expression.

Relate observations in the low-dimensional
gene expression space.



The strategy

non-spatial data

Employ strategies for **batch correction** or **data integration**.

Create a **low-dimensional embedding** based on gene expression.

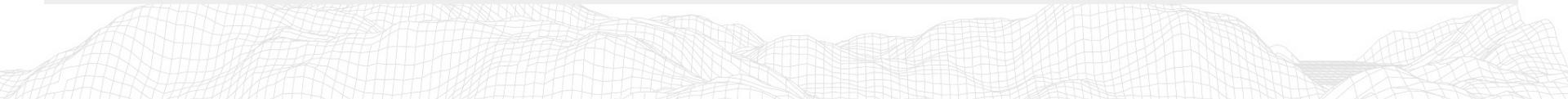
Relate observations in the low-dimensional **gene expression space**.

spatial data

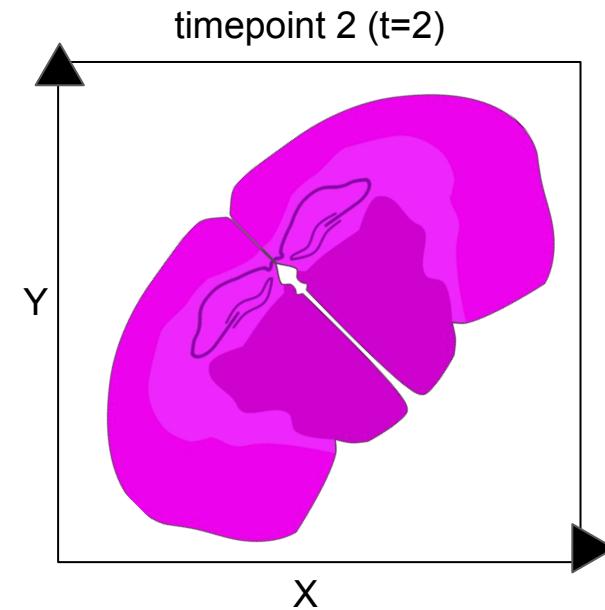
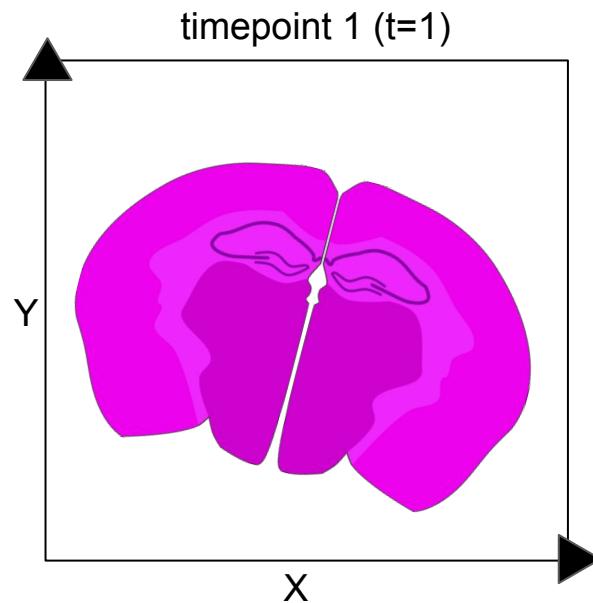
Already possess an **inherent low-dimensional representation**: the spatial domain.

Extend idea of integration beyond removal of batch effects, to also **relate observations in the physical domain**.

Certain **challenges** to this, biggest one being the issue of **independent coordinate systems**



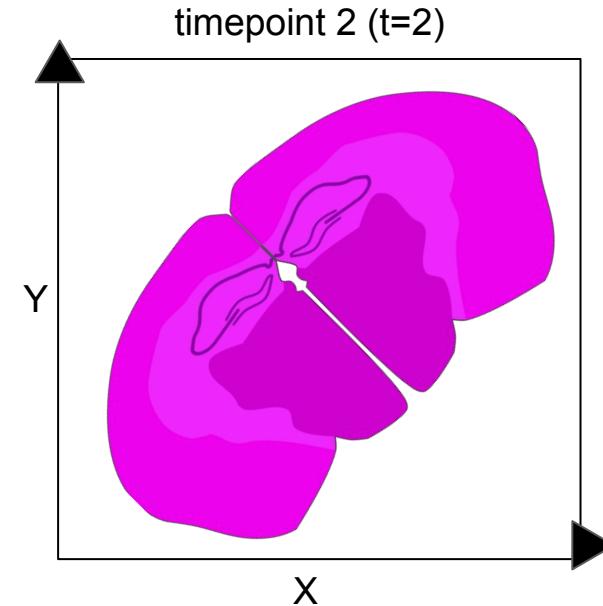
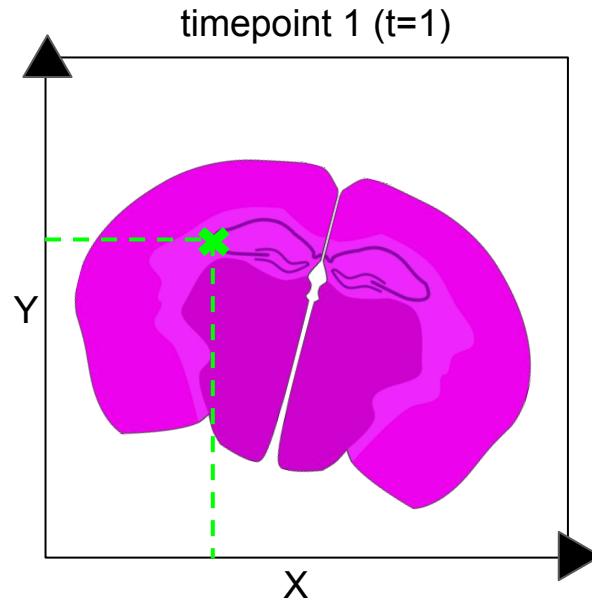
The challenge :: independent coordinate systems



Example: we have two tissue samples from different timepoints, t=1 and t=2

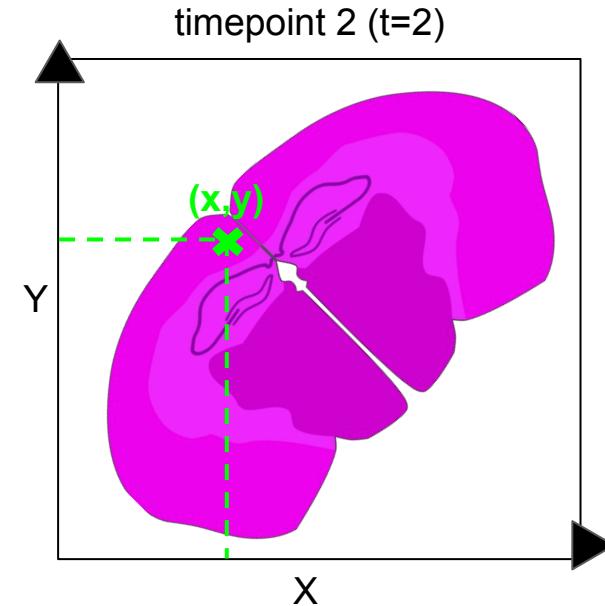
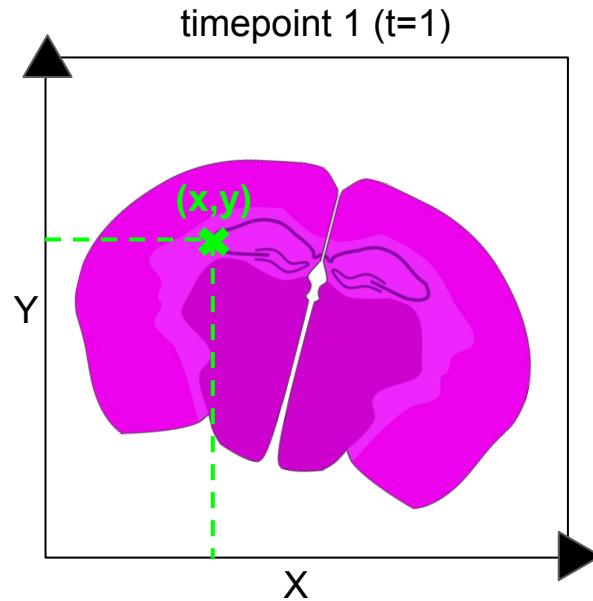


The challenge :: independent coordinate systems



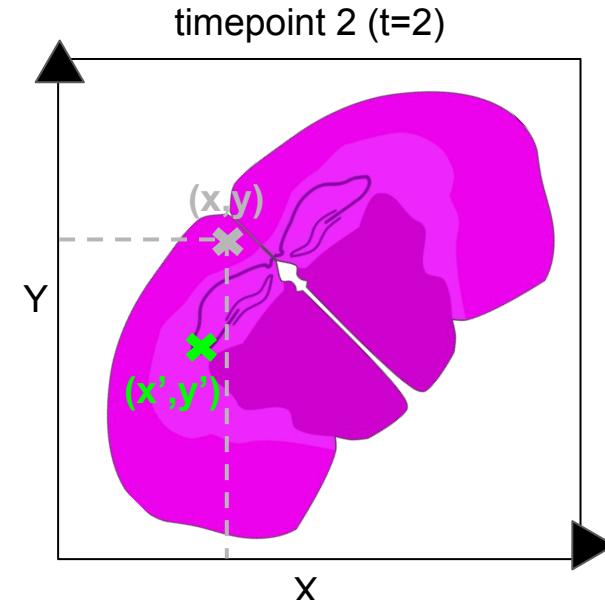
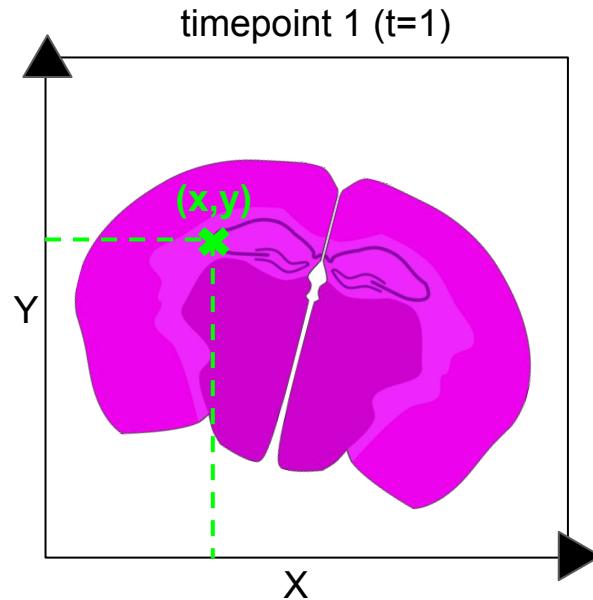
Objective: examine how the gene expression at a certain **position** changes over time

The challenge :: independent coordinate systems



Issue: coordinates (x, y) do not represent the same position in t=1 and t=2; independent coordinate systems

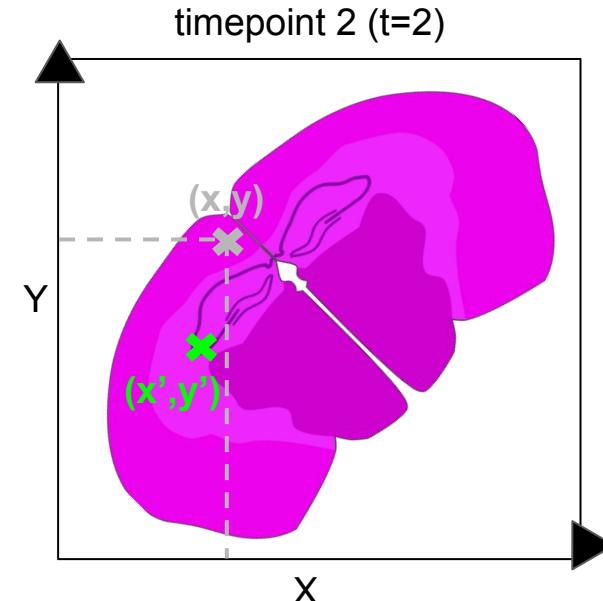
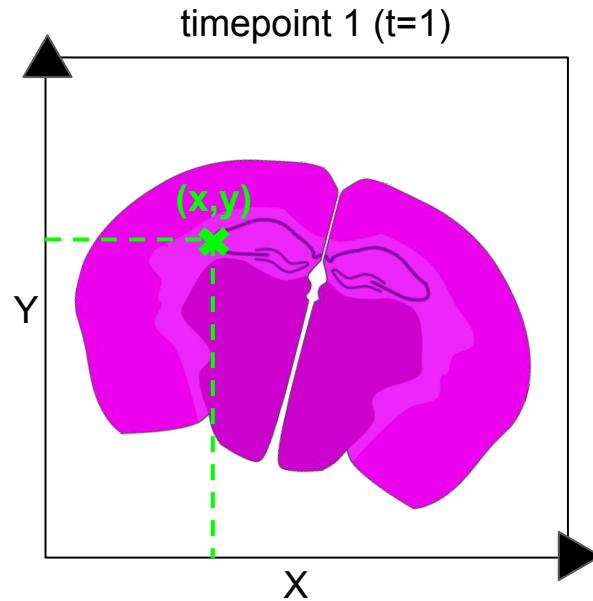
The challenge :: independent coordinate systems



Issue: coordinates (x, y) do not represent the same position in t=1 and t=2; independent coordinate systems

Correct approach: compare expression at coordinates (x, y) at t=1 with expression at (x', y') at t=2

The challenge :: independent coordinate systems



Issue: coordinates (x, y) do not represent the same position in t=1 and t=2; independent coordinate systems

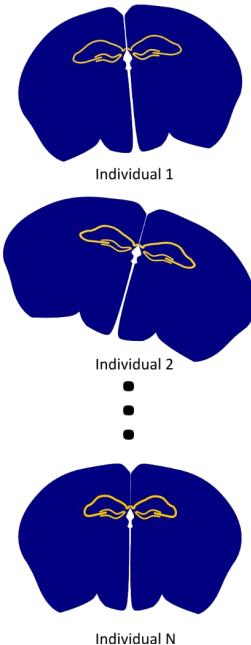
Correct approach: compare expression at coordinates (x, y) at t=1 with expression at (x', y') at t=2

Solution: transfer the expression to a **common coordinate framework** where (x, y) represents same position in all samples

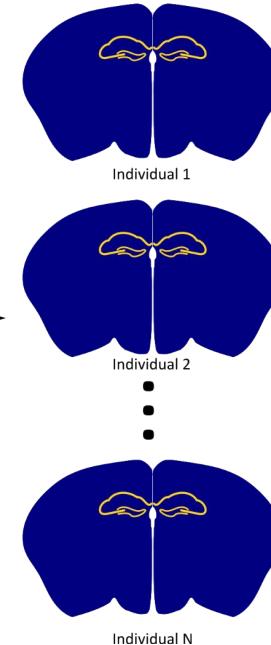


Common Coordinate Frameworks (CCFs)

Observed data



Transferred data



Reference

- Transfer information from observed data to reference/CCF
- Information “inhabits” the same spatial domain.
- **More than an alignment**
- **Important:** corrects for non-linear distortions.

Feature value

Low

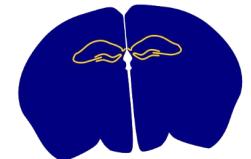
High





Common Coordinate Frameworks (CCFs)

Observed data



Individual 1



Individual 2

⋮



Individual N

Feature value

Low

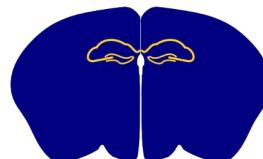
High

Reference

Transferred data



Individual 1



Individual 2

⋮



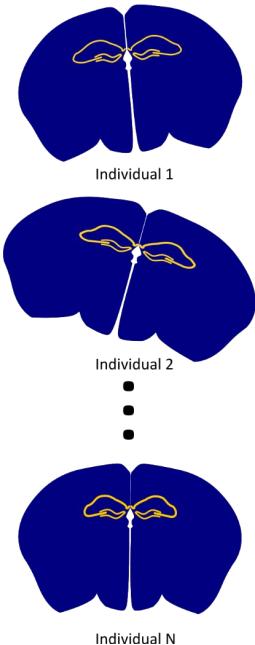
Individual N

- Transfer information from observed data to reference/CCF
- Information “inhabits” the same spatial domain.
- **More than an alignment**
- **Important:** corrects for non-linear distortions.

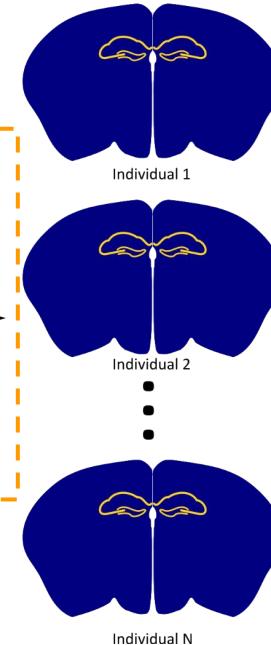


Common Coordinate Frameworks (CCFs)

Observed data



Transferred data



Reference

Our focus: transfer process

- Transfer information from observed data to reference/CCF
- Information “inhabits” the same spatial domain.
- More than an alignment**
- Important:** corrects for non-linear distortions.

Feature value

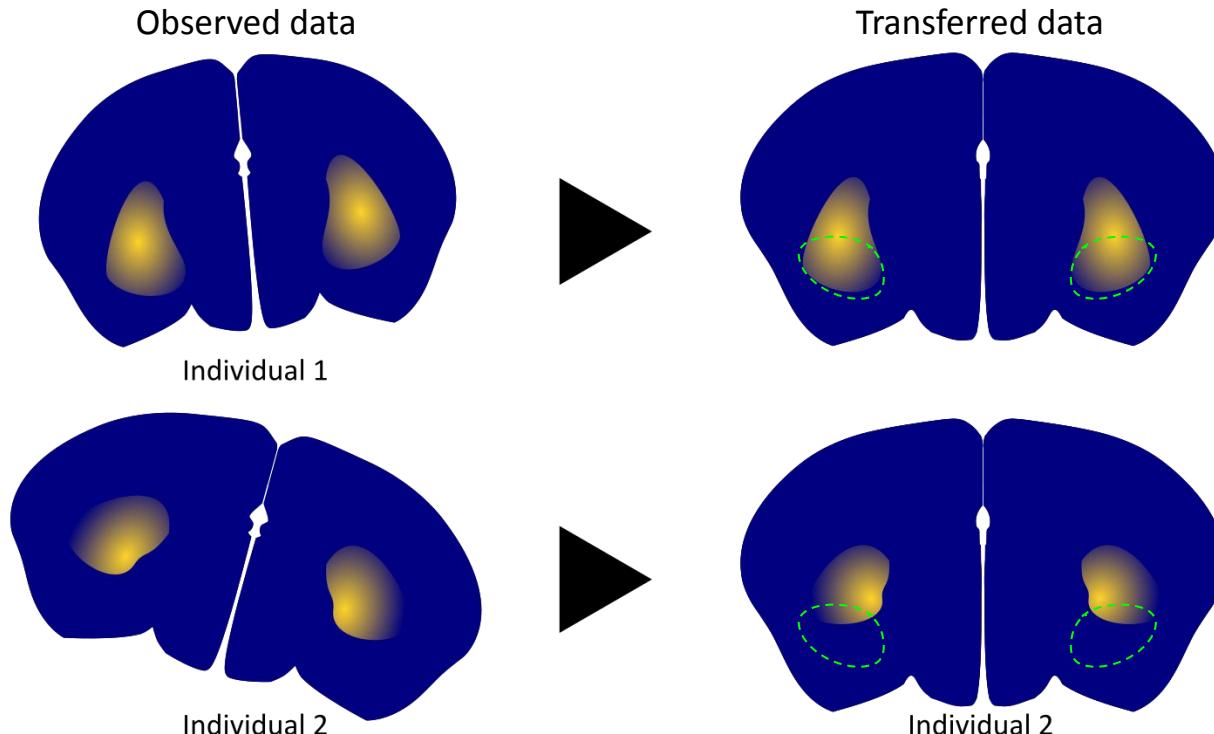
Low

High





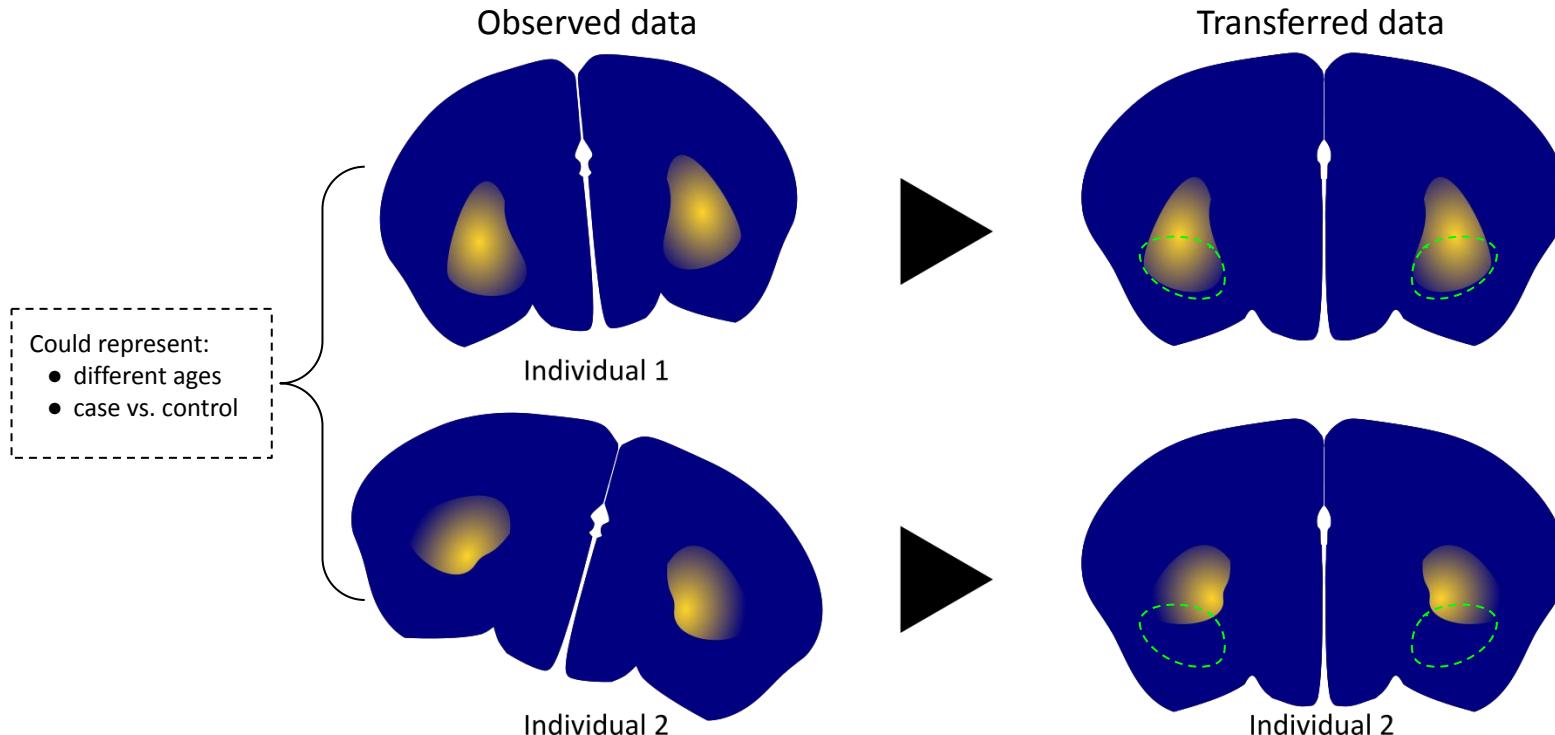
Common Coordinate Frameworks (CCFs)



We can **compare** expression **location-by-location** between different samples
Differences can be assessed both visually and computationally!



Common Coordinate Frameworks (CCFs)



We can **compare** expression **location-by-location** between different samples
Differences can be assessed both visually and computationally!

So what's the purpose?

- Transfer of information to a reference/CCF is **usually not the end goal**.
- Enables plenty of downstream analyses:
 - ***Spatiotemporal modeling:***
 - Analyze and characterize how the expression changes over **time and space**
 - ***Local expression changes:***
 - Assess whether expression is up-or downregulated in certain regions
 - ***Aggregation of data:***
 - Jointly represent information from multiple samples by **composite representations**
 - Identify **general** spatial expression patterns found in the data
- etc ...





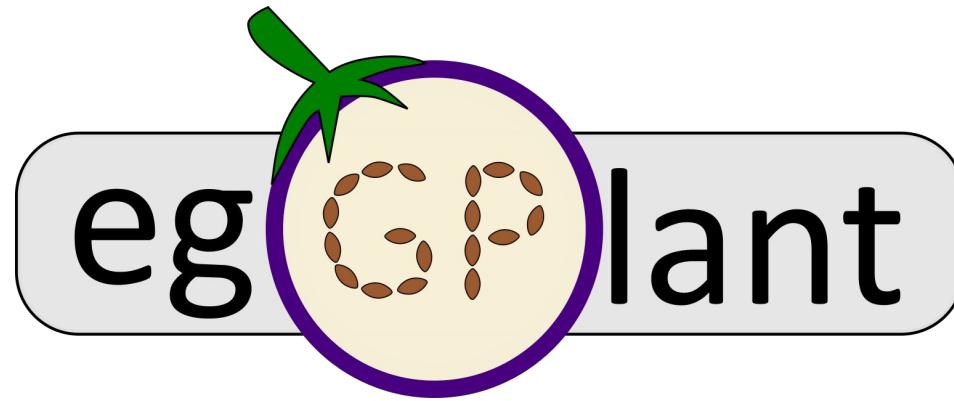
The solution?





The Our solution?

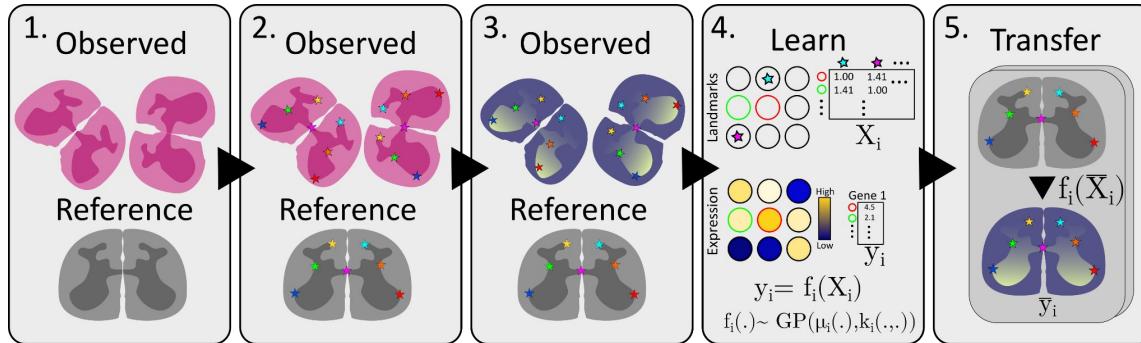




effortless generic GP landmark transfer

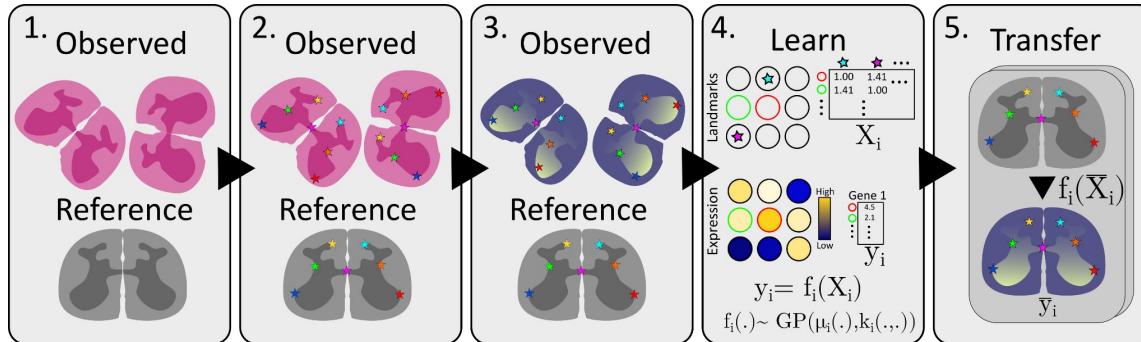


eggplant :: Method overview



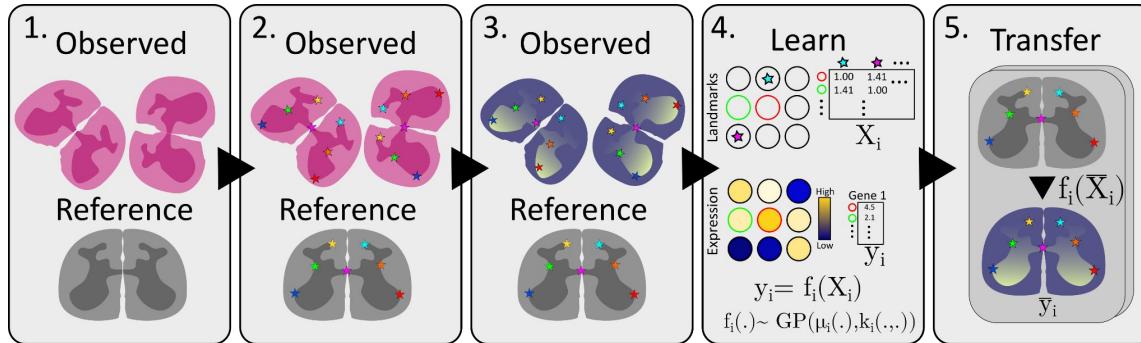
1. Construct/choose a reference and select samples to transfer information from
2. *Chart the landmarks*, i.e., annotate common landmarks
3. Select feature of interest, e.g., expression of your favorite gene
4. Learn a **transfer function** relating landmark distances to expression
5. **Transfer** information to the reference, **apply transfer function** to each location in the reference

eggplant :: Method overview



1. Construct/choose a reference and select samples to transfer information from
 2. Chart the **landmarks**, i.e., annotate common landmarks
 3. Select feature of interest, e.g., expression of your favorite gene
 4. Learn a **transfer function** relating landmark distances to expression
 5. **Transfer** information to the reference, **apply transfer function** to each location in the reference
- Will explain in more detail

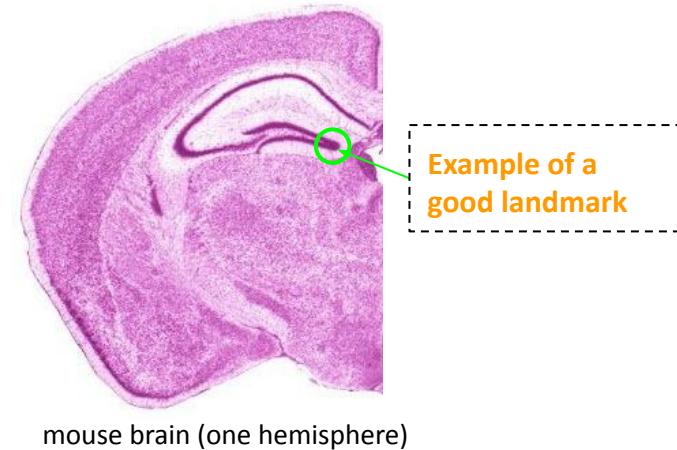
eggplant :: Method overview



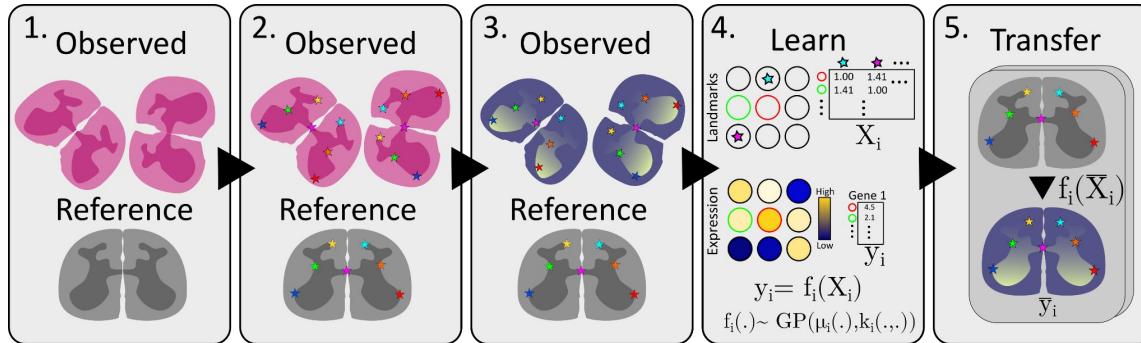
1. Construct/choose a reference and select samples to transfer information from
2. *Chart the landmarks*, i.e., annotate common landmarks
3. Select feature of interest, e.g., expression of your favorite gene
4. Learn a **transfer function** relating landmark distances to expression
5. **Transfer** information to the reference, **apply transfer function** to each location in the reference

■ ■ | What is a landmark?

- **Landmarks:** structures, signals or any features **easily identified across samples**
 - definition largely based on “*Toward a Common Coordinate Framework for the Human Body*” (Rood et al.)
- **Key:** a landmark should always represent the **same approximate position**
- Landmarks can, for example, be based on:
 - Morphology (Image)
 - Gene expression
 - Protein signals
- **Important:** *eggplant* only focuses on the **transfer process**
 - Agnostic to annotation method
 - Manual landmark annotation so far, automatic on the way



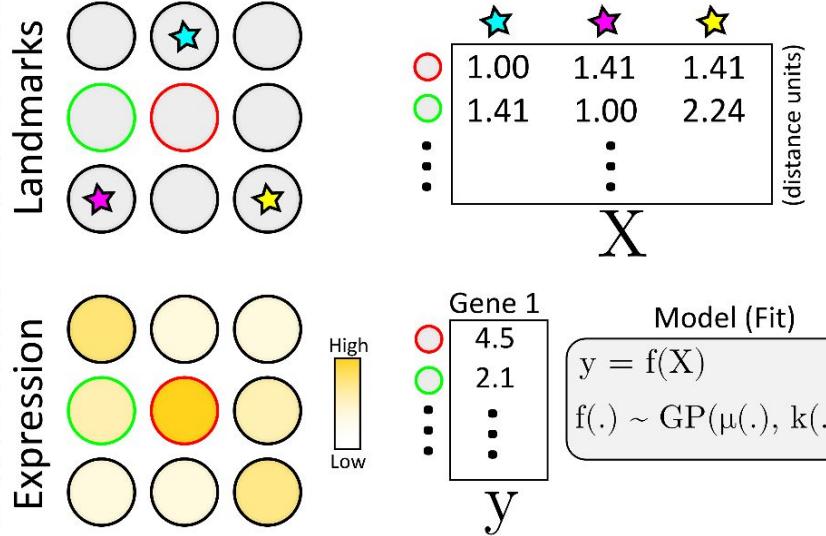
eggplant :: Method overview



1. Construct/choose a reference and select samples to transfer information from
2. *Chart the landmarks*, i.e., annotate common landmarks
3. Select feature of interest, e.g., expression of your favorite gene
4. Learn a **transfer function** relating landmark distances to expression
5. **Transfer** information to the reference, **apply transfer function** to each location in the reference

Learning the transfer function

Observed



- For each location: measure distance to all landmarks → **distance matrix** (X)
- Assume that gene expression (y) is **related** to landmark distances **via a function**: $y = f(X)$
- Assumes that the transfer function (f) follows a **Gaussian Process** (GP) distribution
- Use **Gaussian Process Regression** to determine character of the transfer function

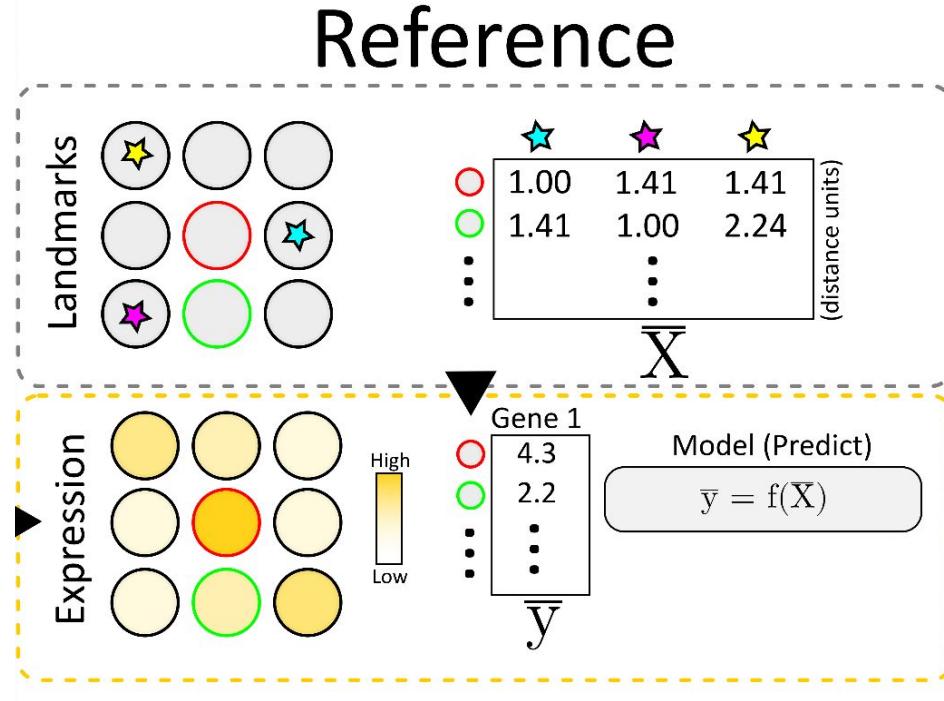
X : observed landmark distance matrix (given)

y : observed gene expression vector (given)



Transferring information to the reference

- Create **distance matrix** (\bar{X}) for the reference
- **Transfer** information by **applying** transfer function to \bar{X} (prediction)

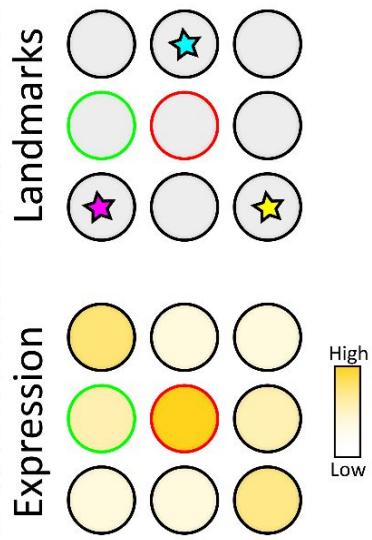


\bar{X} : reference landmark distance matrix (given)

\bar{y} : **predicted** gene expression vector

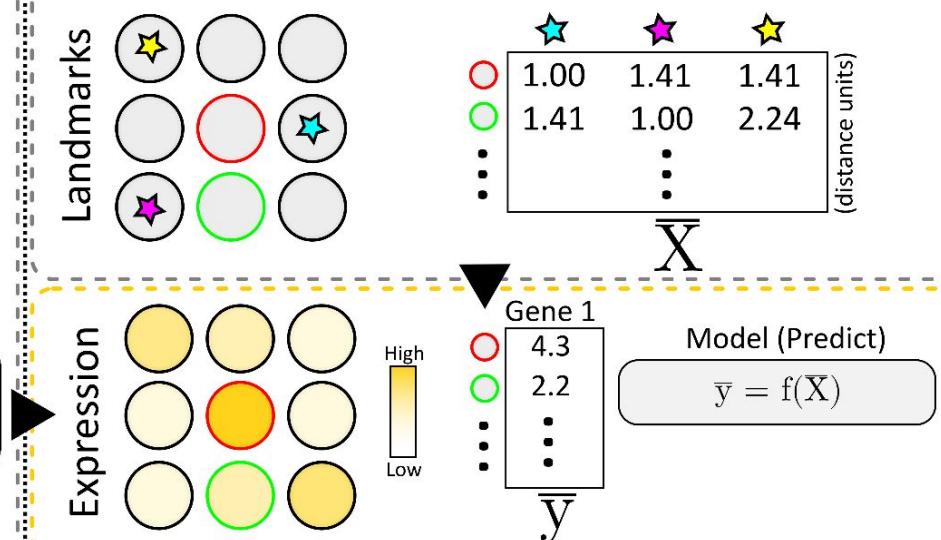
Transferring information to the reference

Observed



X : observed landmark distance matrix (given)
 y : observed gene expression vector (given)

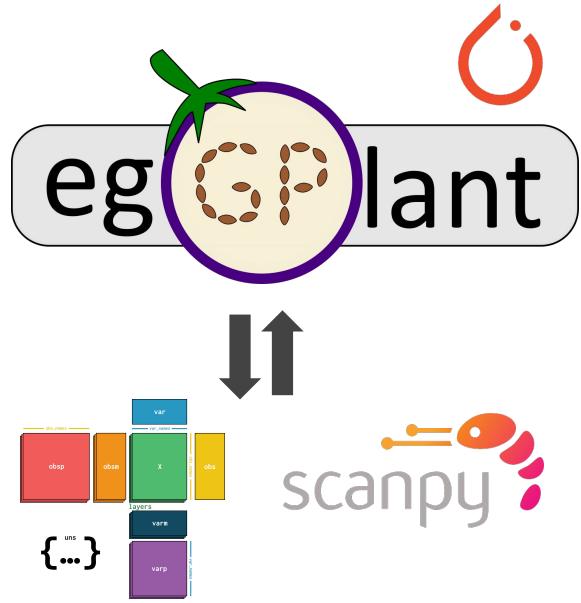
Reference



\bar{X} : reference landmark distance matrix (given)
 \bar{y} : **predicted** gene expression vector

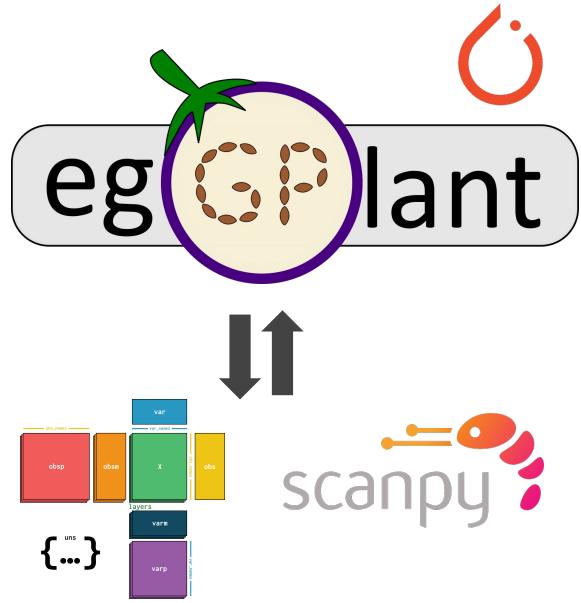
Implementation

- **Python package:**
 - written to be compatible with AnnData objects
 - scanpy-like API
 - can be installed via PyPi and GitHub (recommended)
 - documented on: Read the Docs
- **Backend:** PyTorch/GPyTorch



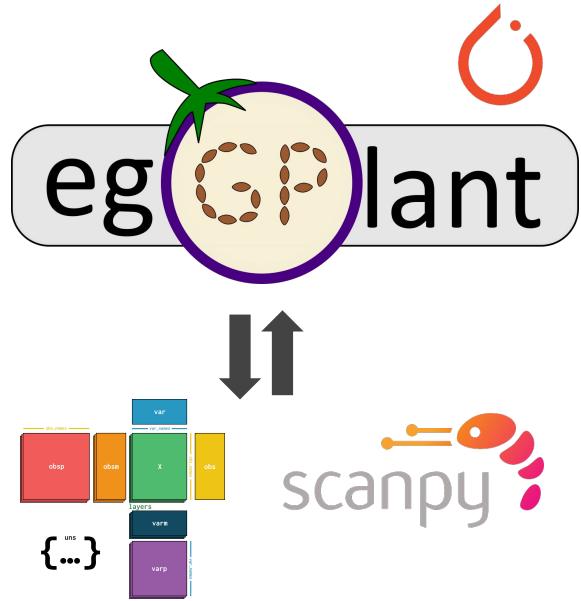
■ ■ ■ Implementation

- **Python package:**
 - written to be compatible with AnnData objects
 - scanpy-like API
 - can be installed via PyPi and GitHub (recommended)
 - documented on: Read the Docs
- **Backend:** PyTorch/GPyTorch
- **Challenge:**
 - GP's have notoriously bad scaling: $O(n^3)$ where n = number of observations
 - We have plenty of observations, and many features.



■ ■ ■ Implementation

- **Python package:**
 - written to be compatible with AnnData objects
 - scanpy-like API
 - can be installed via PyPi and GitHub (recommended)
 - documented on: Read the Docs
- **Backend:** PyTorch/GPyTorch
- **Challenge:**
 - GP's have notoriously bad scaling: $O(n^3)$ where n = number of observations
 - We have plenty of observations, and many features.



■ ■ | Implementation : the challenge

Strategies for scaling the method to large omics data sets:

- **GPU acceleration:**
 - General acceleration of the inference process
 - Handled by PyTorch/GPyTorch
- **Variational Inference:**
 - Uses inducing points
 - reduces complexity w.r.t. Observations : $O(n^3) \rightarrow O(m^2n)$ ($m = \#$ inducing points)
- **Fast approximate transfer:**
 - Uses a data decomposition step
 - Reduces the number of features that has to be transferred
 - Steps:
 1. Use PCA to project full-dimensional data into a K-dimensional subspace
 2. Transfer loadings of each principal component to the reference
 3. Reconstruct full-dimensional data in the reference

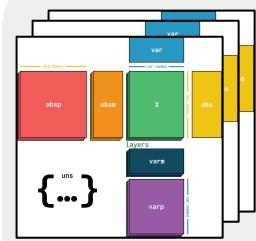




Implementation : input and output

Input

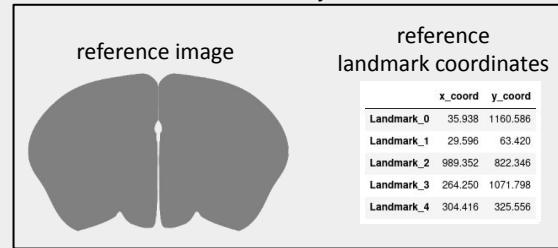
anndata objects



- *.X* : features of interest (e.g., gene expression)
- *.obsm["spatial"]* : coordinates of observations (e.g., spots)
- *.uns["curated_landmarks"]* : position of landmarks
 - *.obsm["landmark_distances"]* : distances to landmarks



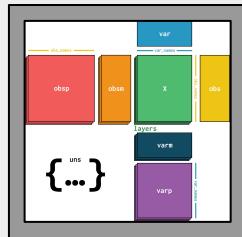
reference object



eggplant

Output

reference object
(anndata object)



Individual 1



Individual 2



Individual N



■ ■ | Example analyses

➡ Human developmental heart (❖)

- Basic example
- Variance in structure between individuals
- Comparison with alignment

➡ Mouse brain (sagittal)

- Complex structures

➡ Perturbation in mouse brain (❖)

- spatial differential expression analysis (SDEA)

■ ■ ■ Example analyses

➡ Human developmental heart (❖)

- Basic example
- Variance in structure between individuals
- Comparison with alignment

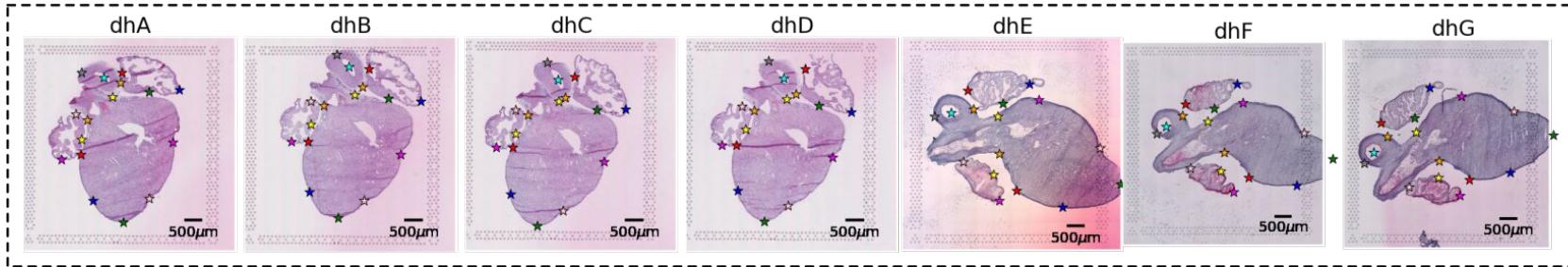
➡ Mouse brain (sagittal)

- Complex structures

➡ Perturbation in mouse brain (❖)

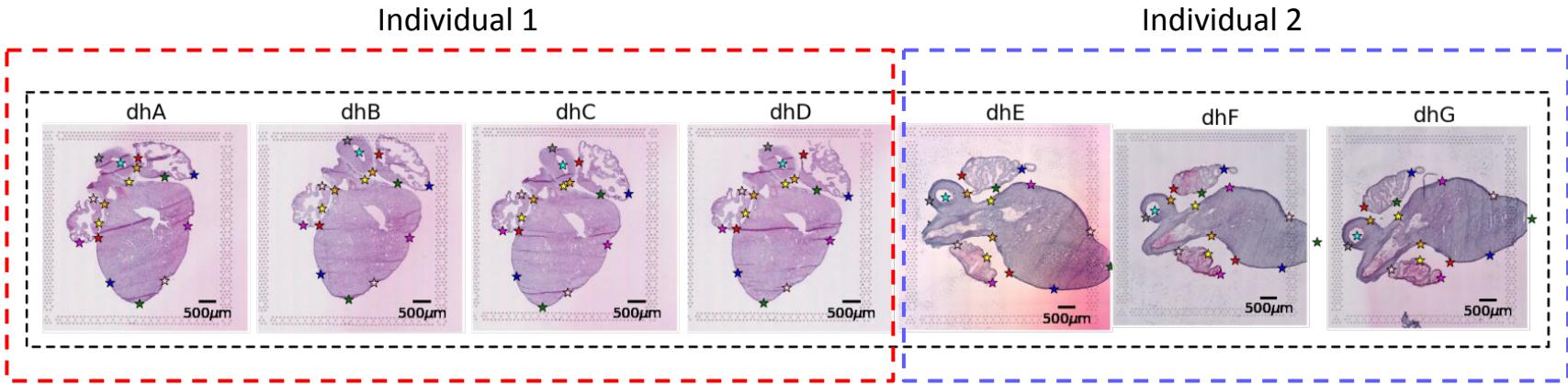
- spatial differential expression analysis (SDEA)

eggplant :: developmental heart analysis



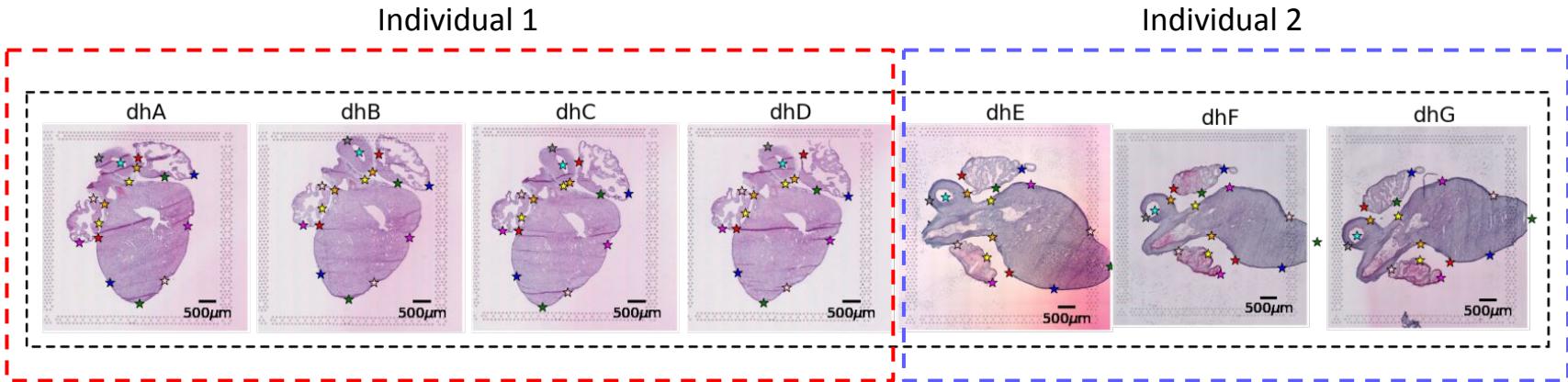
★ = Landmark

eggplant :: developmental heart analysis



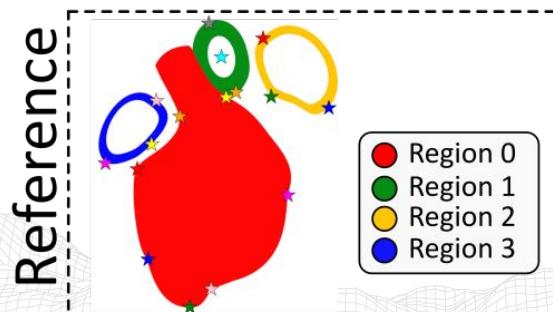
★ = Landmark

eggplant :: developmental heart analysis

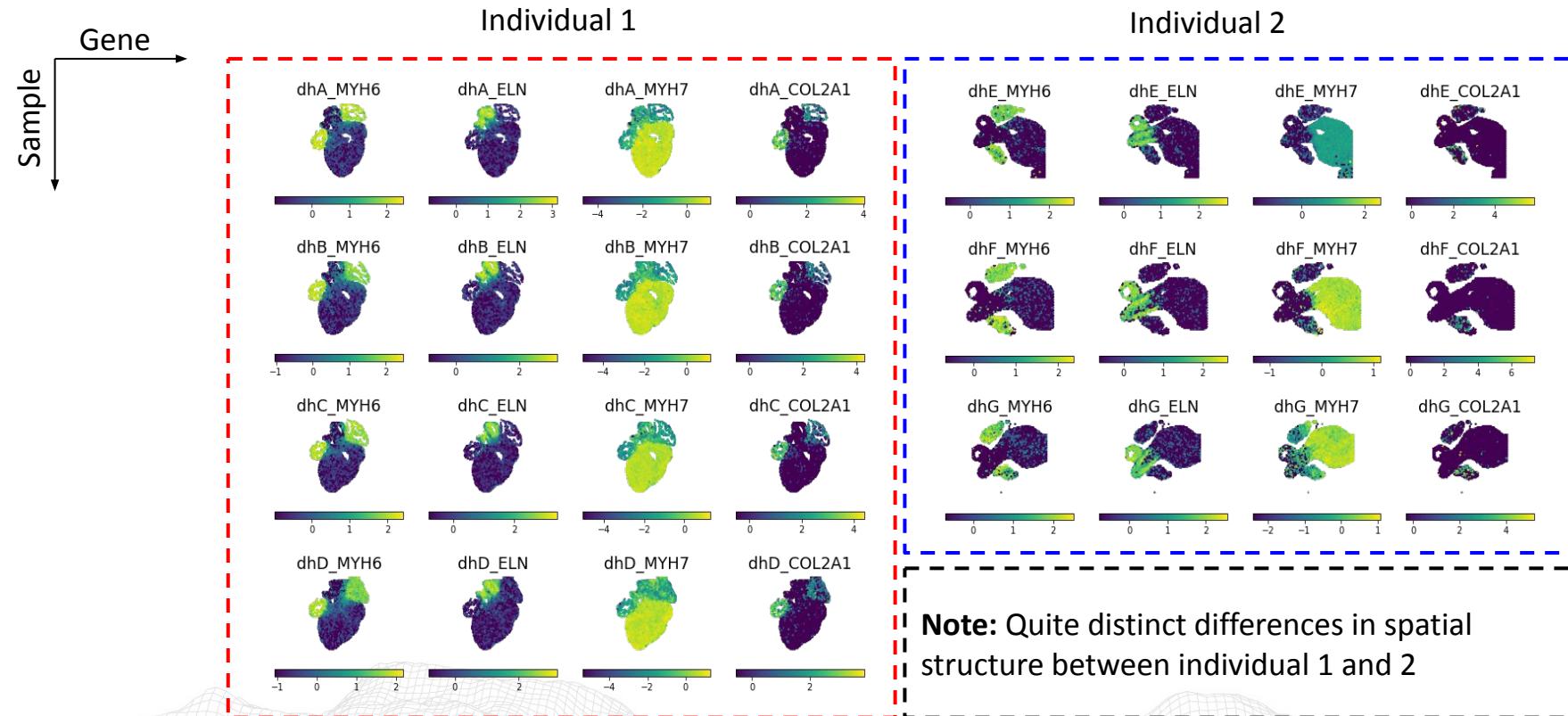


Transfer (eggplant)

★ = Landmark



■ ■ II Observed expression (developmental heart)



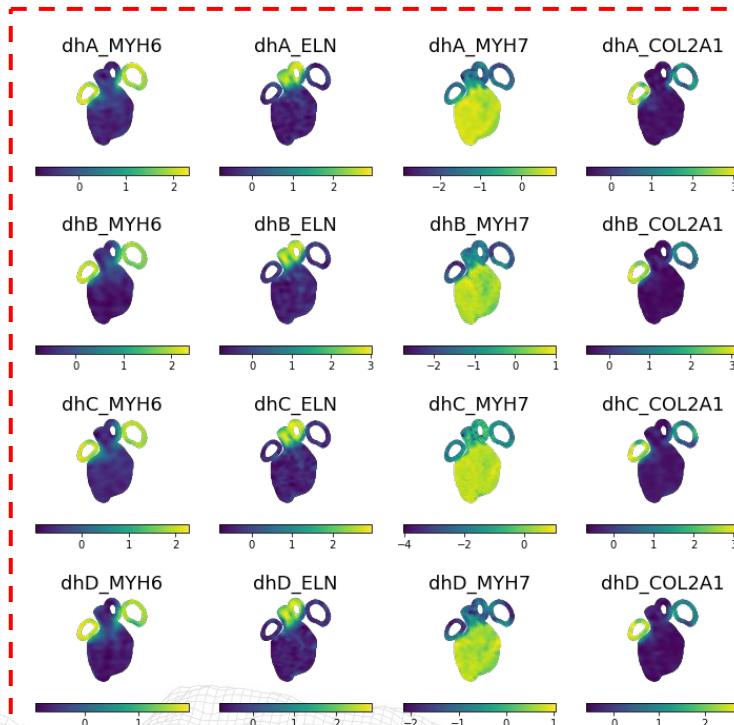


Transferred expression (developmental heart)

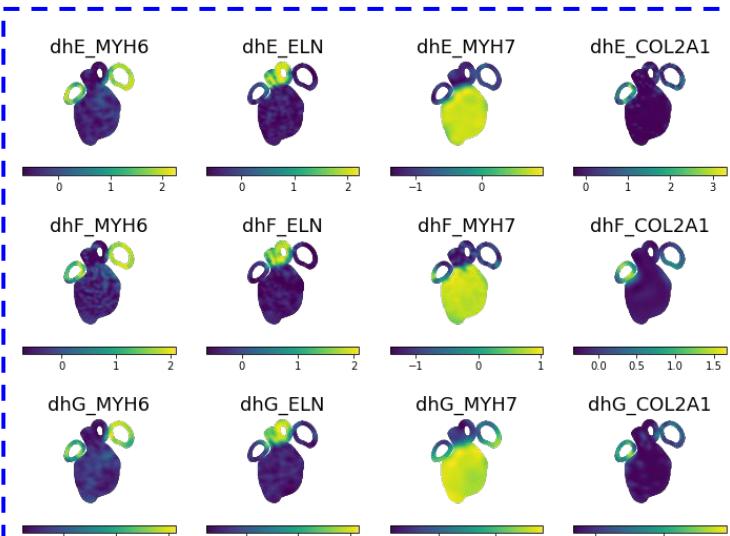
Sample

Gene

Individual 1



Individual 2



Good agreement between samples from different individuals!

eggplant :: Average representations and regional enrichment

Composite

Composite Profile:
COL2A1



Composite Profile:
ELN



Composite Profile:
MYH6

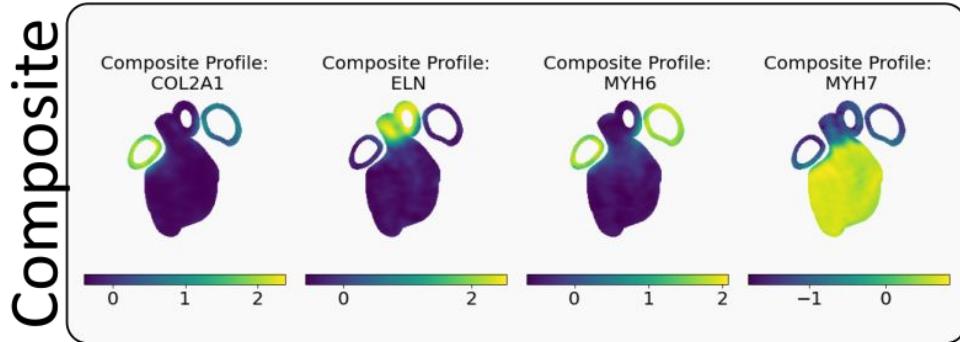


Composite Profile:
MYH7

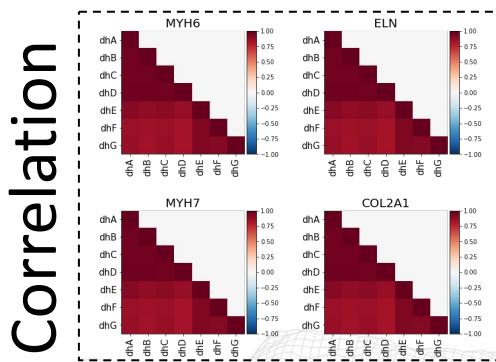


- Composite representations to **summarize information**

 | eggplant :: Average representations and regional enrichment

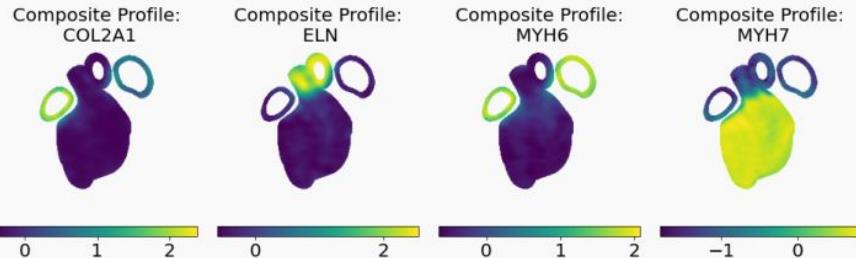


- Composite representations to **summarize information**
 - Location-wise correlation across the 7 samples
 - **high correlation** within and **between individuals**



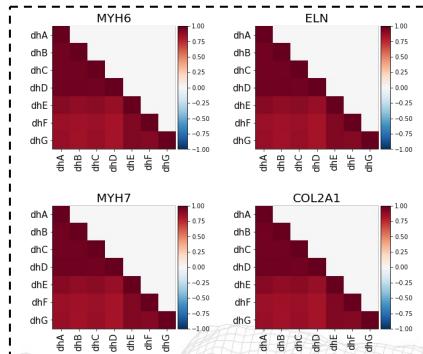
eggplant :: Average representations and regional enrichment

Composite

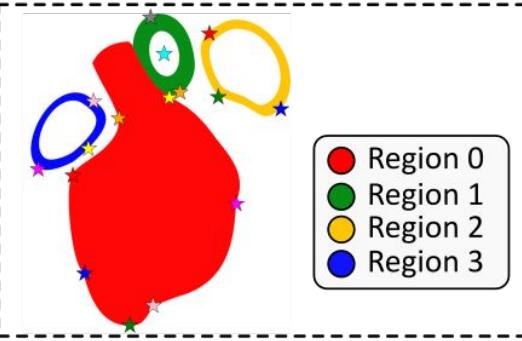


- Composite representations to **summarize information**
- Location-wise correlation across the 7 samples
 - **high correlation** within and **between individuals**
- Use annotated regions to for **enrichment** analysis
 - Assess up/downregulation in regions

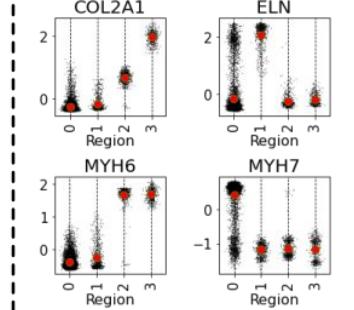
Correlation



Reference

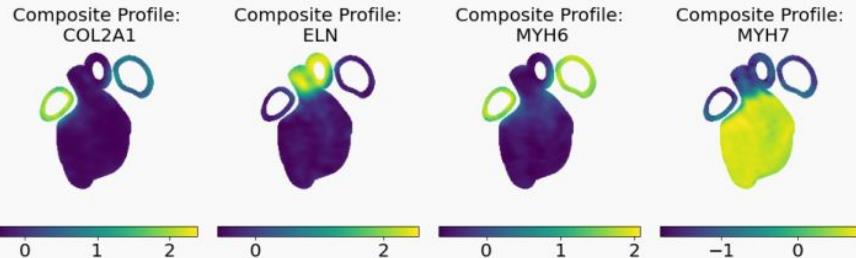


Enrichment



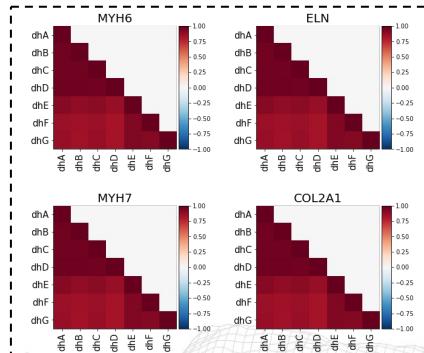
eggplant :: Average representations and regional enrichment

Composite

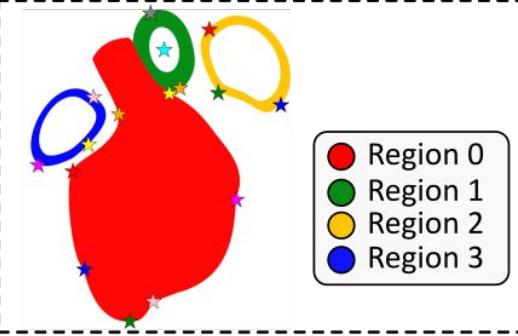


- Composite representations to **summarize information**
- Location-wise correlation across the 7 samples
 - **high correlation** within and **between individuals**
- Use annotated regions to for **enrichment** analysis
 - Assess up/downregulation in regions

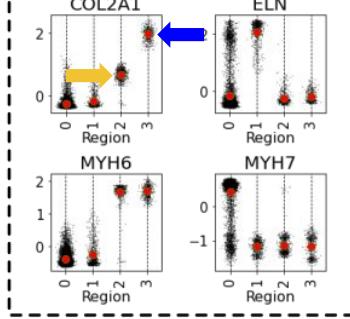
Correlation



Reference



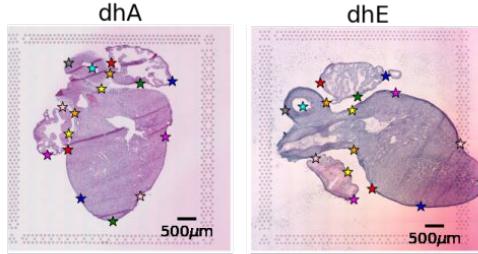
Enrichment





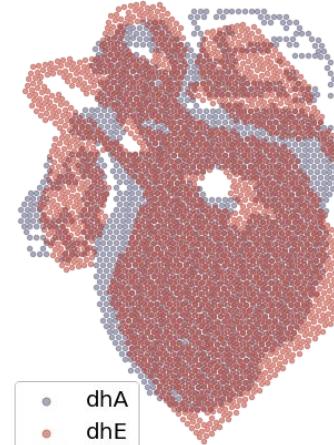
Alignment vs. information transfer

Alignment-based



► PASTE ◀

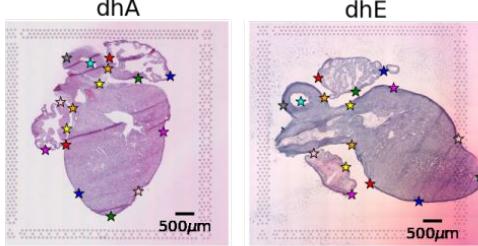
Aligned dhA and dhE



Coordinates still
don't correspond
to the same
position.

PASTE: "Alignment and integration of spatial transcriptomics data", Zeira et al. | doi.org/10.1038/s41592-022-01459-6

CCF-based



► egGPLant ◀

dhA_MYH6

dhE_MYH6

dhA_ELN

dhE_ELN

dhA_MYH7

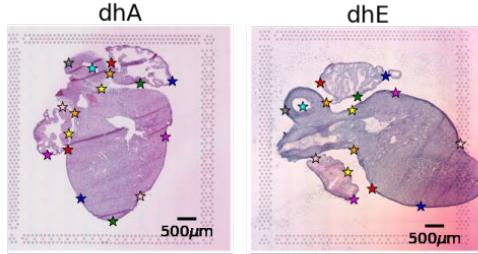
dhE_MYH7

The data inhabits the same space, corresponding coordinates.



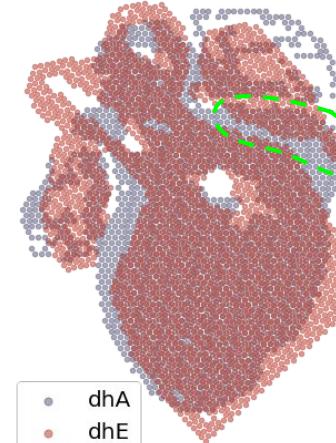
Alignment vs. information transfer

Alignment-based



► PASTE ◀

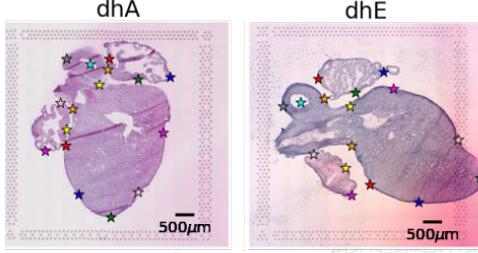
Aligned dhA and dhE



Coordinates still
don't correspond
to the same
position.

PASTE: "Alignment and integration of spatial transcriptomics data", Zeira et al. | doi.org/10.1038/s41592-022-01459-6

CCF-based



► egGPLant ◀

dhA_MYH6

dhE_MYH6

dhA_ELN

dhE_ELN

dhA_MYH7

dhE_MYH7

The data inhabits the same space, corresponding coordinates.

■ ■ | Example analyses

➡ Human developmental heart (❖)

- Basic example
- Variance in structure between individuals
- Comparison with alignment

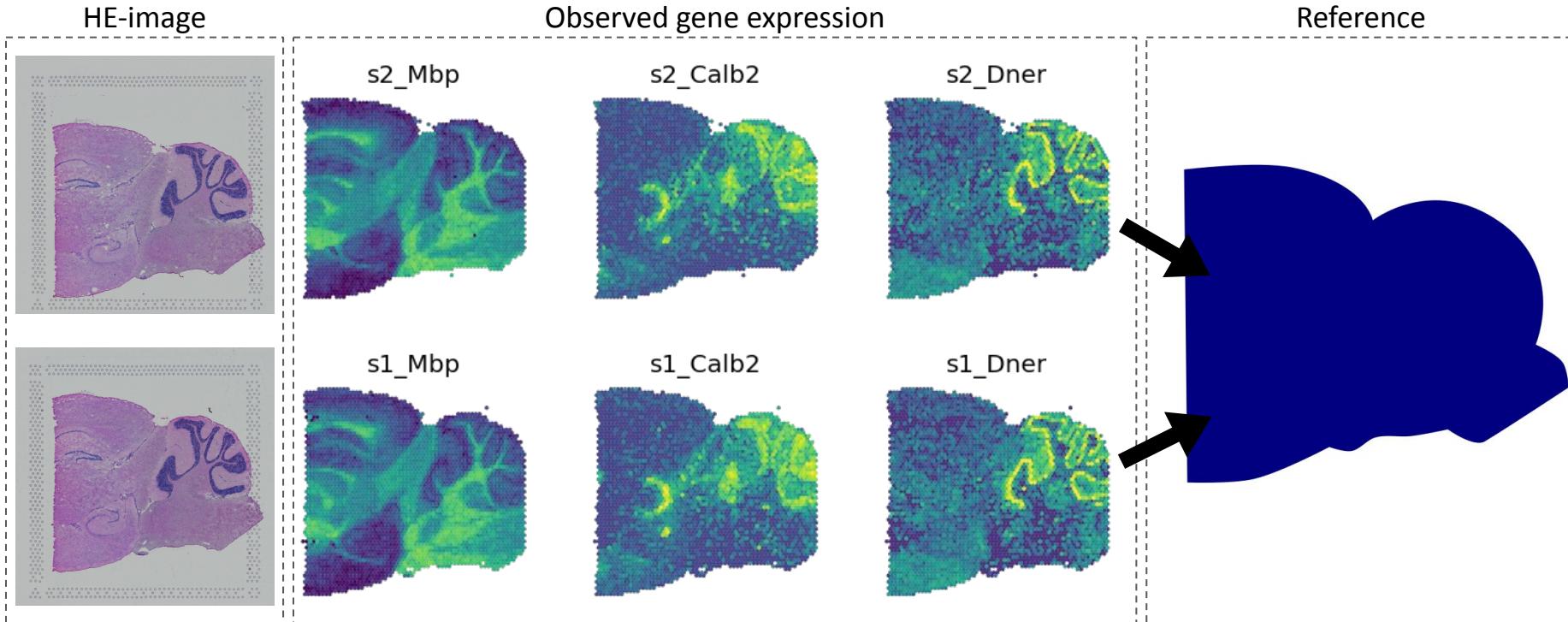
➡ Mouse brain (sagittal)

- Complex structures

➡ Perturbation in mouse brain (❖)

- spatial differential expression analysis (SDEA)

■■■ Mouse brain (sagittal)

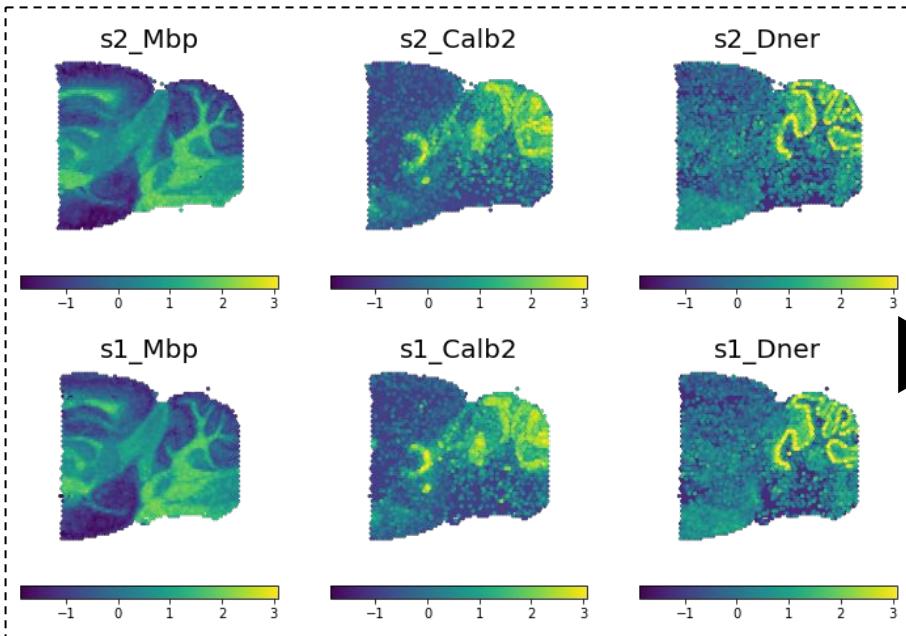


- 2 replicates of mouse brain (sagittal sections)
- Want to transfer: expression of 3 genes with **complex structure**

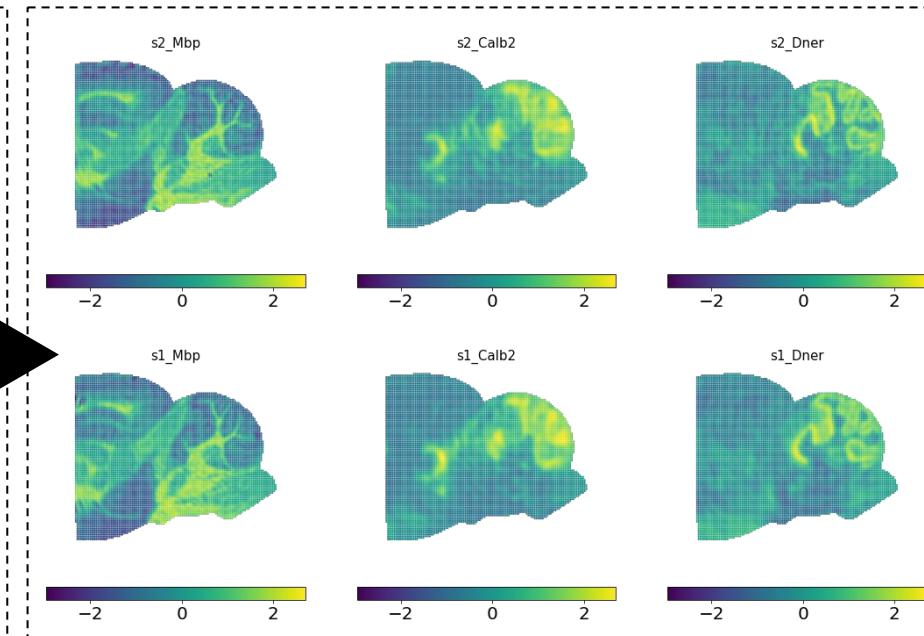


Mouse brain (sagittal)

Observed gene expression



Transferred gene expression



■ ■ | Example analyses

➡ Human developmental heart (❖)

- Basic example
- Variance in structure between individuals
- Comparison with alignment

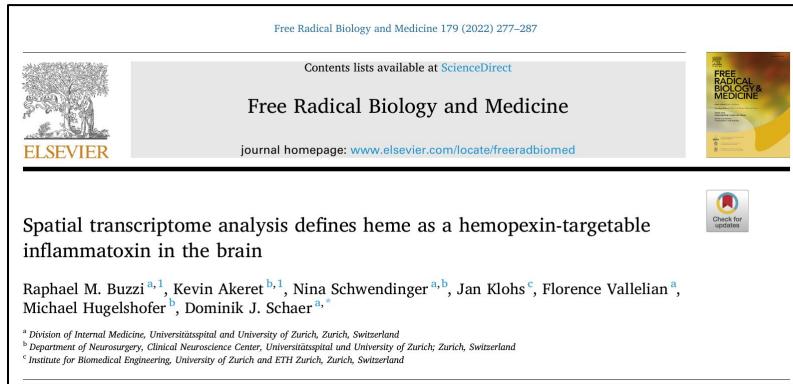
➡ Mouse brain (sagittal)

- Complex structures

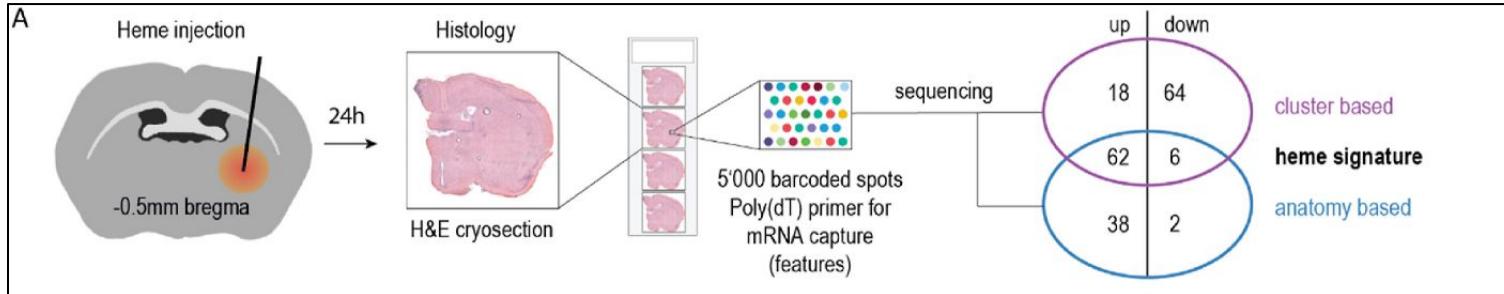
➡ Perturbation in mouse brain (❖)

- spatial differential expression analysis (SDEA)

Perturbation in mouse brain

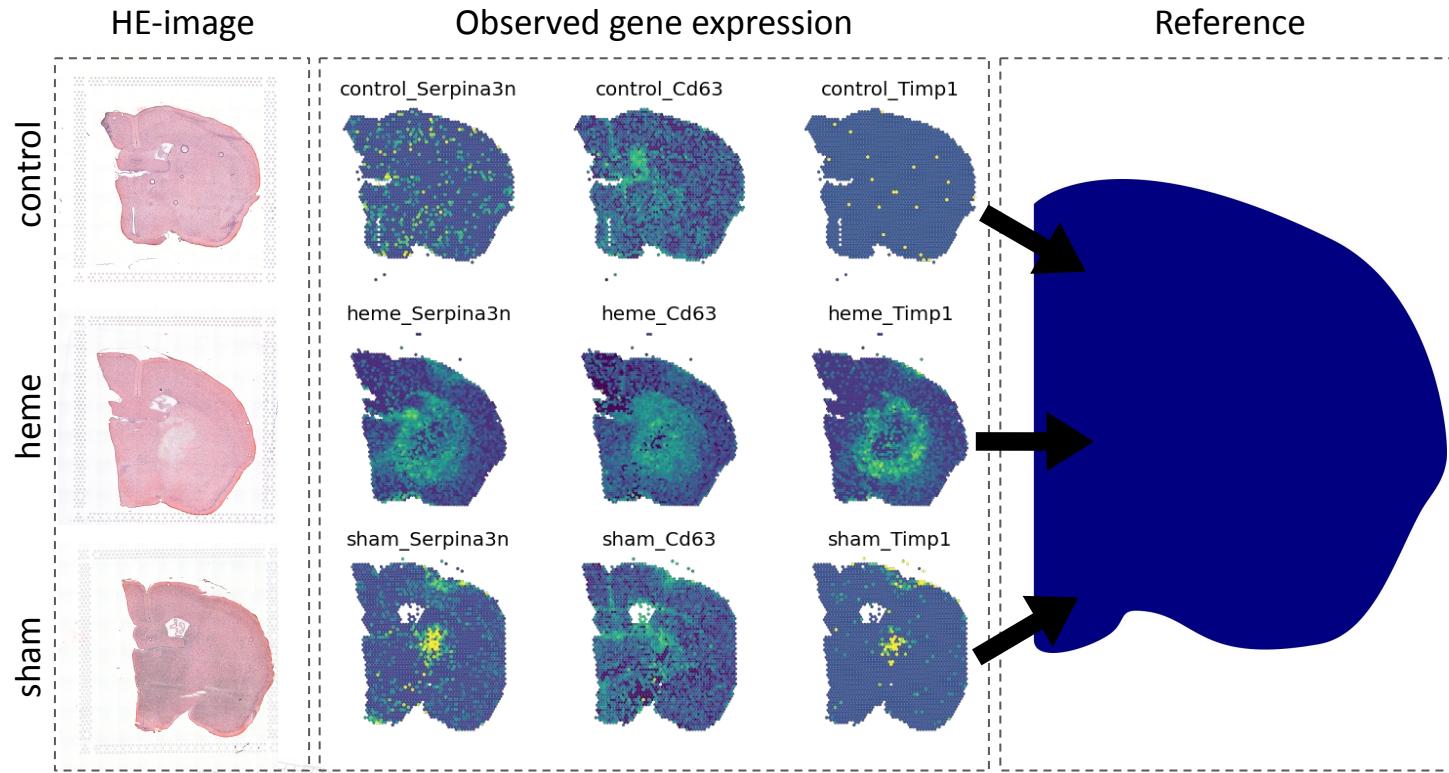


- Already published data
- The authors **injected mice with:**
 - nothing (control)
 - saline (sham)
 - heme
- Then **produced Visium data** from said mice



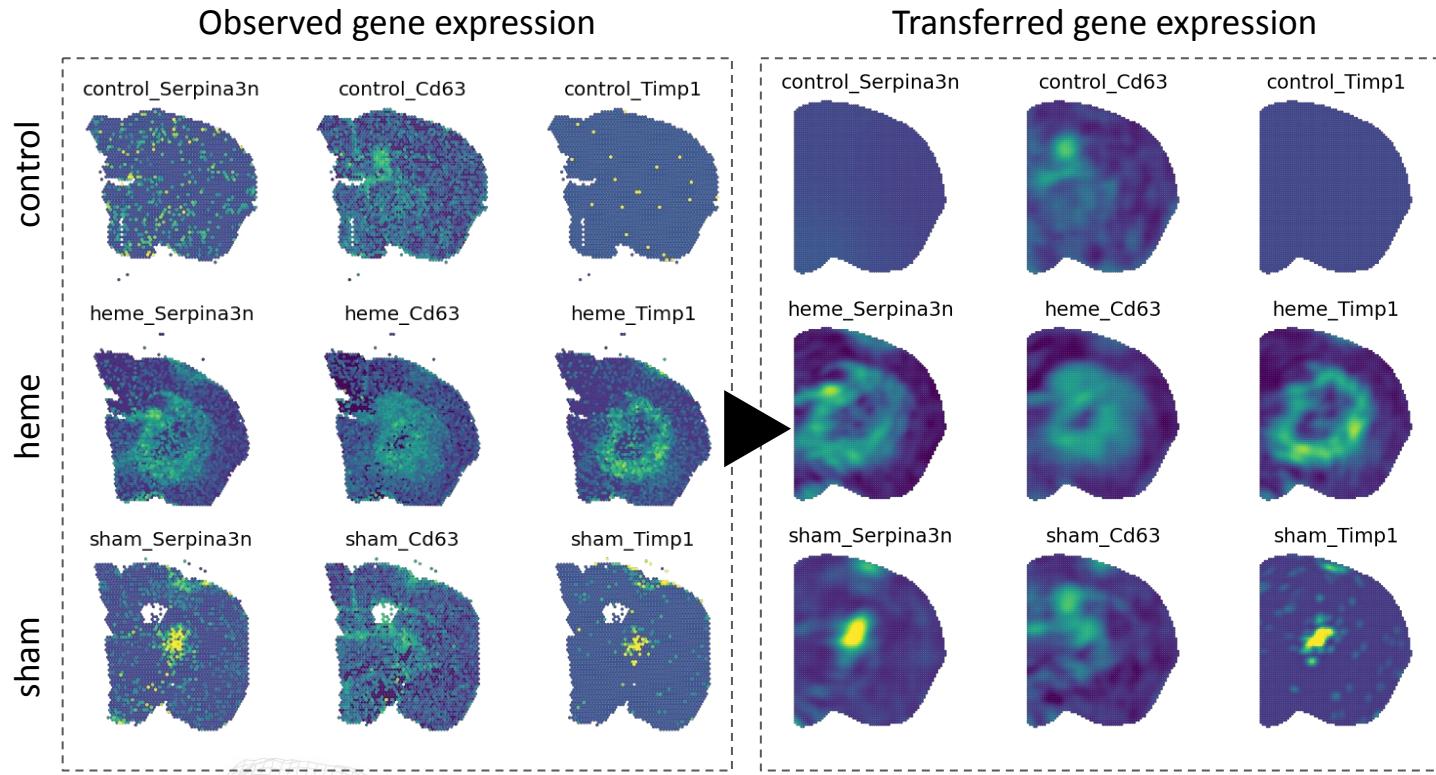


Perturbation in mouse brain



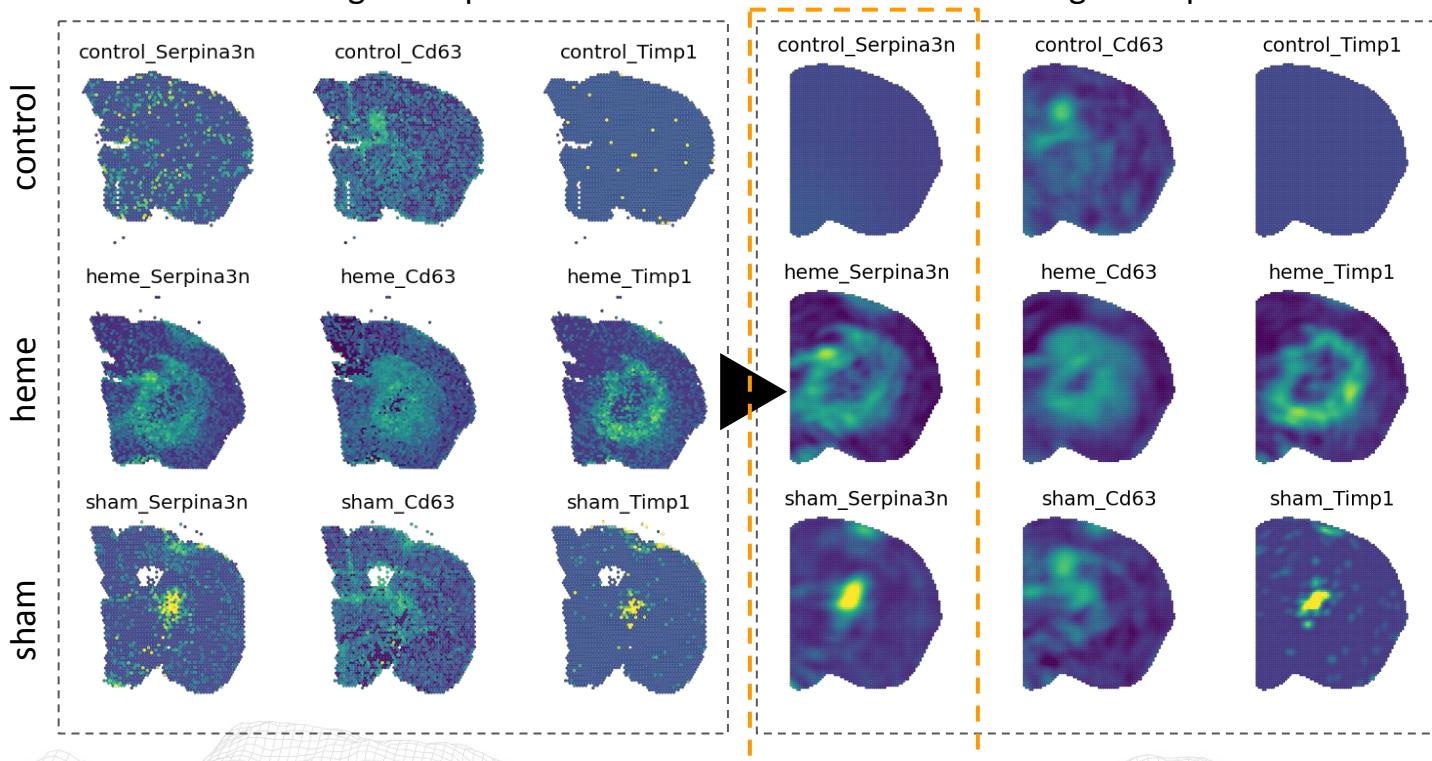


Perturbation in mouse brain



Perturbation in mouse brain

Observed gene expression

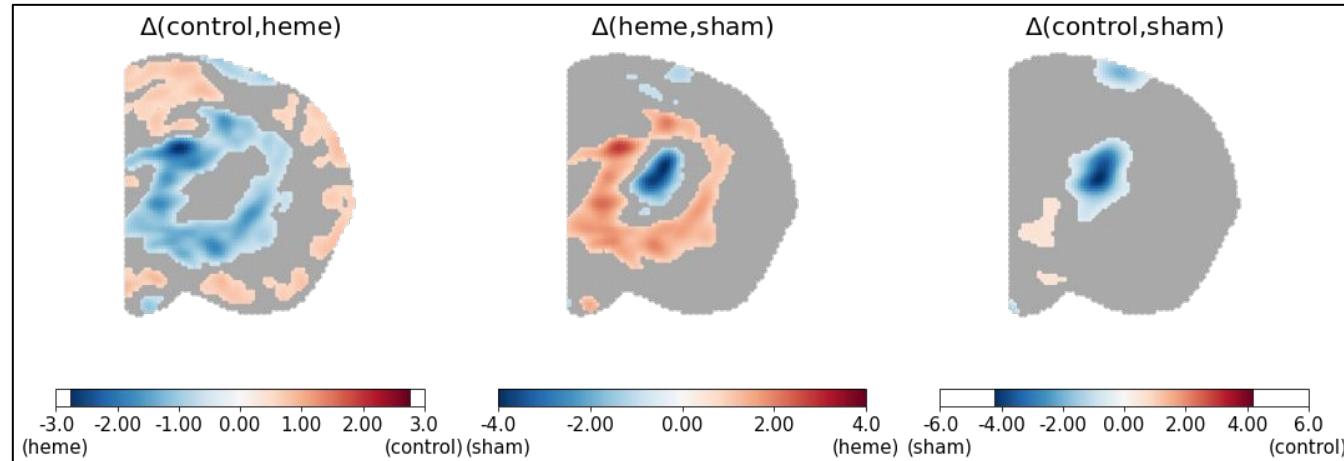


Example of SDEA with *Serpina3n*



Perturbation in mouse brain

Spatial Differential Expression Analysis (SDEA) of *Serpina3n*

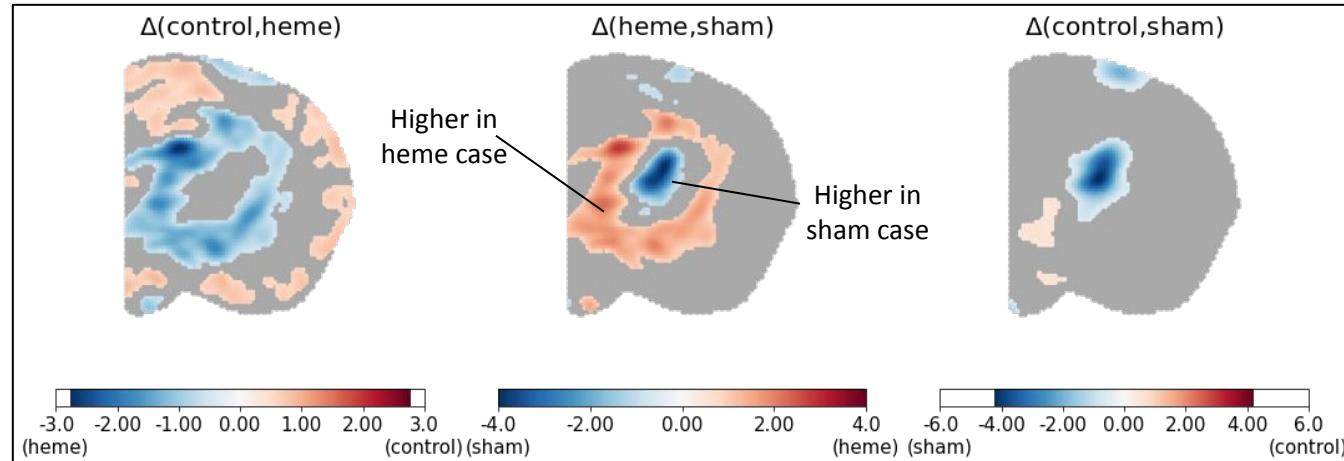


- Compare spatial gene expression between conditions
- Gray indicates non-DE regions
- Only possible when the information inhabits the same reference/CCF



Perturbation in mouse brain

Spatial Differential Expression Analysis (SDEA) of *Serpina3n*

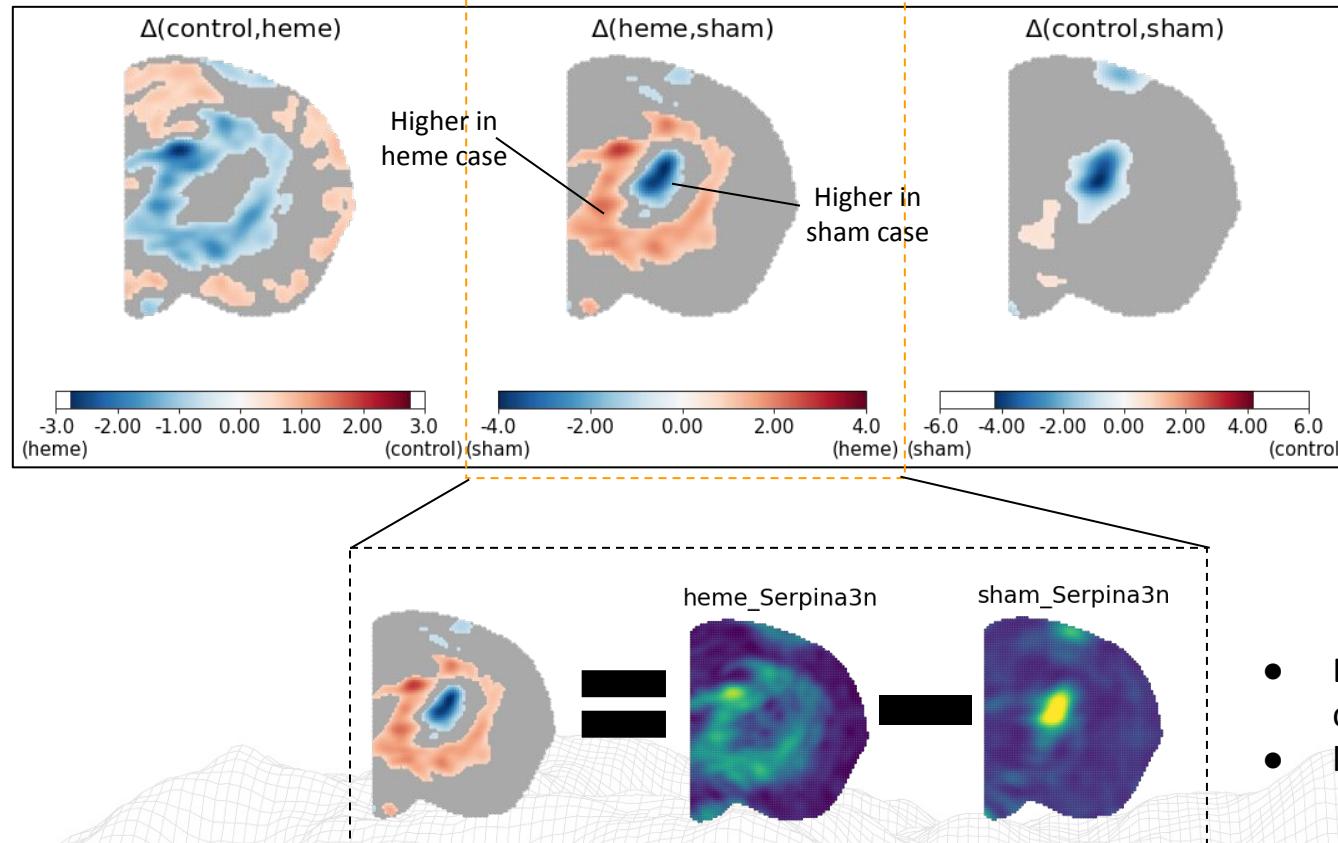


- Compare spatial gene expression between conditions
- Gray indicates non-DE regions
- Only possible when the information inhabits the same reference/CCF



Perturbation in mouse brain

Spatial Differential Expression Analysis (SDEA) of *Serpina3n*



- Compare spatial gene expression between conditions
- Gray indicates non-DE regions
- Only possible when the information inhabits the same reference/CCF

- Difference between the two conditions (spatial arithmetics)
- location-by-location comparison



Summary

- Common coordinate frameworks (**CCFs**) represents:
 - a unified framework that allows us to relate observations across samples
 - implicitly, this also means comparisons across **time points and conditions**
- Transferring data to a CCF is **not the same as aligning** samples
- For spatial data, **spatially aware CCFs** can be highly useful
- Our method (*eggplant*) uses:
 - **Landmark annotation**
 - Gaussian Process Regression
- Spatial Differential Expression Analysis (SDEA) is:
 - enabled by having the data in the same reference/CCF
 - a way to compare expression location-by-location



More information



THE PREPRINT SERVER FOR BIOLOGY

A Landmark-based Common Coordinate Framework for Spatial Transcriptomics

Data

Alma Andersson, Žaneta Andrusiová, Paulo Czarnecki, Xiaofei Li, Erik Sundström,
Joakim Lundeberg

doi: <https://doi.org/10.1101/2021.11.11.468178>

This screenshot shows the bioRxiv preprint page for the manuscript "A Landmark-based Common Coordinate Framework for Spatial Transcriptomics Data". The page features a large image of a purple eggplant with a green sprout. The title "eggplant" is overlaid on the image. Below the image, a text block explains the purpose of the package: "This repository contains the source code for the package `eggplant` presented in the manuscript "A Landmark-based Common Coordinate Framework for Spatial Transcriptomics Data"; which - in short - is a method designed to transfer information from multiple spatial-transcriptomics data sets to a single reference representing a **Common Coordinate Framework (CCF)**." At the bottom of the page, there is a GitHub link: <https://github.com/almaan/eggplant>.

<https://github.com/almaan/eggplant>

This screenshot shows the documentation for the `eggplant` package. It includes a search bar, a "GETTING STARTED" section, and a sidebar with links to various modules: `eggplant package`, `eggplant.methods module`, `eggplant.models module`, `eggplant.plot module`, `eggplant.preprocess module`, and `eggplant.sdea module`. A specific class, `PoissonDiscSampler`, is highlighted in the `eggplant.methods module` documentation. The class definition is shown in a code block:

```
class eggplant.methods.PoissonDiscSampler(crd: numpy.ndarray, min_dist: float, seed: Optional[int] = None) [source]
```

The `Bases` field indicates that it inherits from `object`. The `Poisson Disc Sampler` is described as "Designed according to the principles outlined in Bridson [Bri07] but for d=2." A method, `add_k_in_annulus`, is also defined in the class:

```
add_k_in_annulus(point: Union[numpy.ndarray, Tuple[float, float]], k: int = 5) → numpy.ndarray [source]
```

The method adds `k` points randomly in an annulus around a given point.

<https://spatial-eggplant.readthedocs.io>

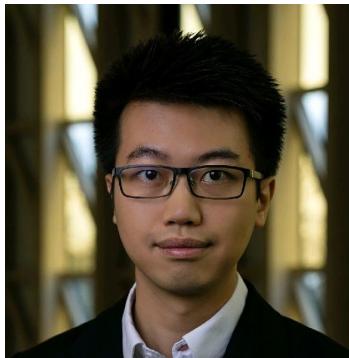


Acknowledgments

Acknowledgements



Žaneta Andrusiová



Xiaofei Li



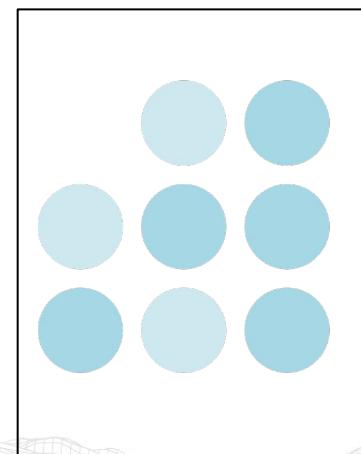
Paulo Czarnewski



Erik Sundström



Joakim Lundeberg



Spatial Research Group

Developmental Heart Data:
Zaneta, Xiaofei and Erik

Code beta-testing and discussions:
Paulo Czarnewski

Supervisor:
Joakim Lundeberg

Special Shoutout: Giovanni Palla for
excellent feedback!!!



Thank *you* for the attention!
Questions?



@aalmaander



alma.andersson@differentiable.net



differentiable.net

Notebook walkthrough

Included:

- Load data and reference
- Prepare data and reference for transfer of information
- Transfer of information to reference, using two strategies:
 - exact (but slightly slower) approach
 - approximate (but *significantly* faster) approach
- Creating composite representations
- Regional enrichment analysis
- Spatial differential gene expression analysis (SDEA)

Datasets:

- Human developmental heart
- Perturbed mouse brain

Comment:

- Run notebook in **DEMO** mode to load pre-analyzed data

eggplant : an introduction

- Author: [Alma Andersson](#)
- Purpose: originally designed for the [SCOG workshop 2022-05-24](#)
- Created: 2022-05-17

Welcome to the `eggplant` session, here we'll explore some of the most basic features of *eggplant*. A method and python package to transfer information from different spatial-omics sections to a single reference (common coordinate framework).

How does the method work? (brief)

The method is more thoroughly described in it's associated [manuscript](#), but to briefly outline the five main steps, we will use the same flowchart as is presented in Figure 1A (of the above referenced manuscript).

https://github.com/theislab/spatial_scog_workshop_2022/tree/main/eggplant/notebooks