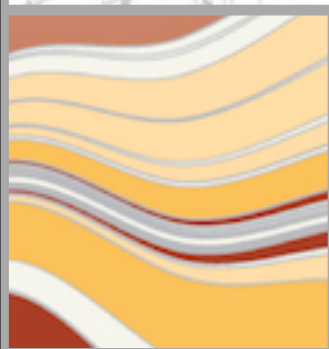
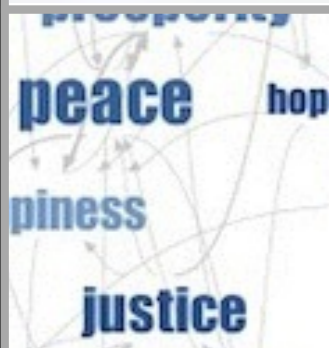
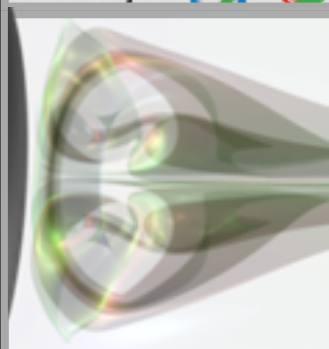
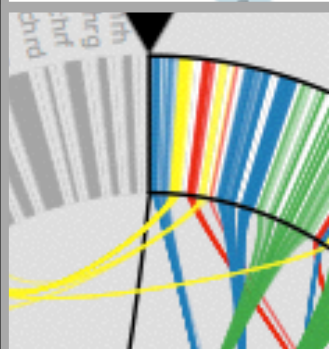
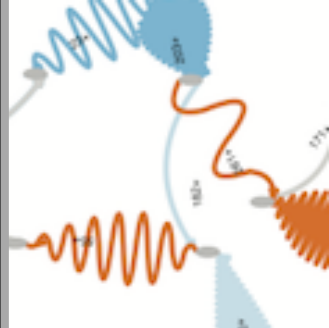


# FILTERING & AGGREGATION

Miriah Meyer  
*University of Utah*



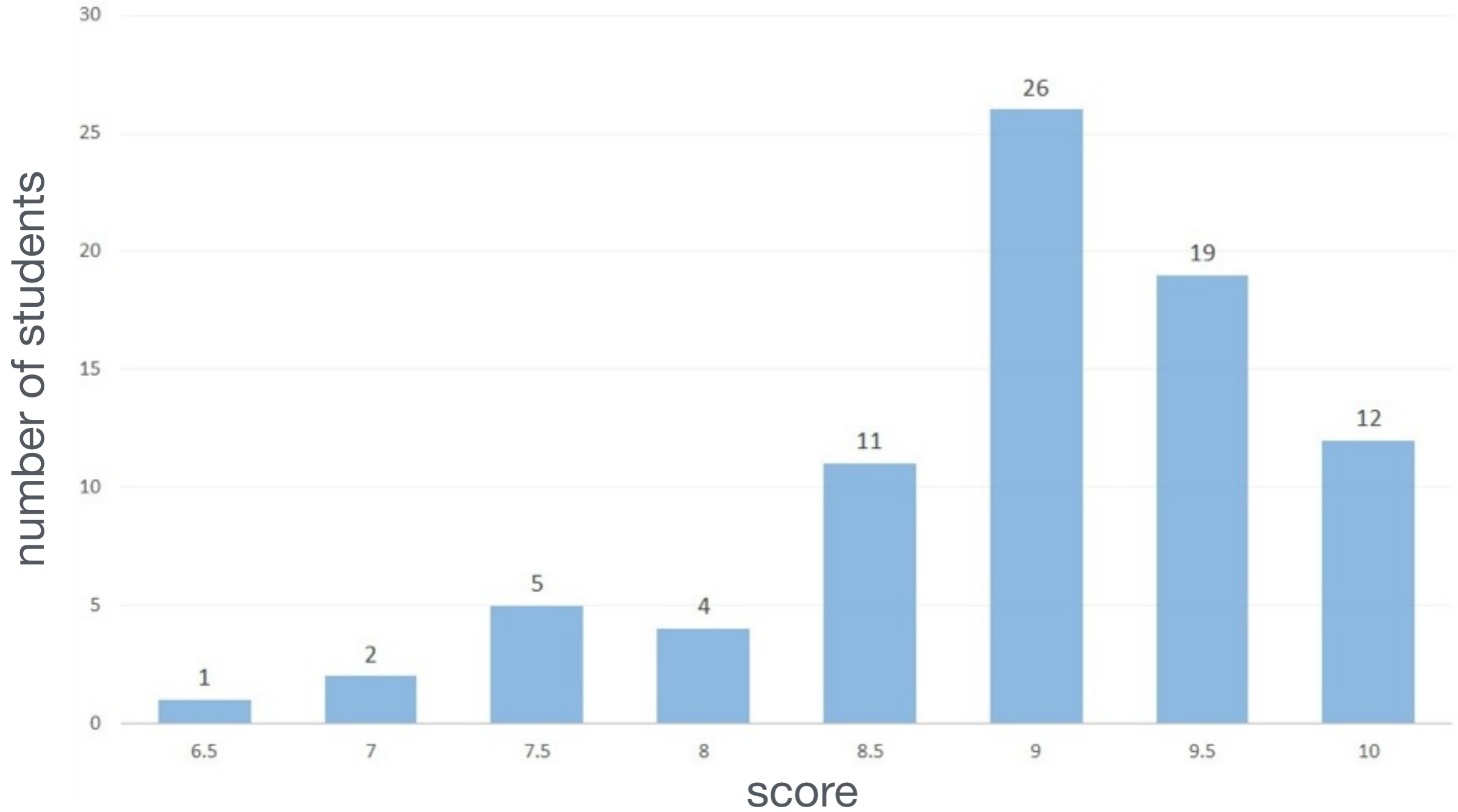
administrivia . . .



- exam on Tuesday

- data exploration grades are in

# data exploration assignment



last time . . .

# FOCUS + CONTEXT

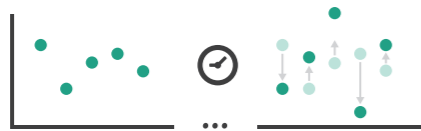
carefully pick what to show

hint at what you are not showing

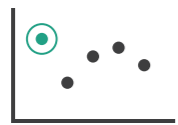
## Manipulate

---

### → Change



### → Select



### → Navigate



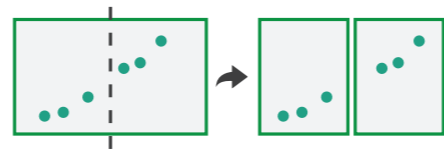
## Facet

---

### → Juxtapose



### → Partition



### → Superimpose



## Reduce

---

### → Embed

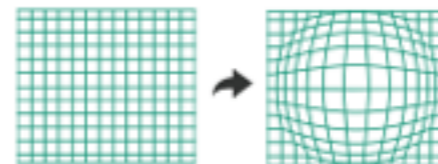
#### → Elide Data



#### → Superimpose Layer



#### → Distort Geometry



what is elision?

# what is elision?

focus items shown in detail, other items summarized for context

superimpose



# superimpose

focus layer limited to a local region of view, instead of stretching across the entire view

distort

# distort

use geometric distortion of the contextual regions to make room for the details in the focus region(s)

# distortion concerns

- unsuitable for relative spatial judgements**
- overhead of tracking distortion**
- visual communication of distortion**
  - gridlines, shading
- target acquisition problem**
  - lens displacing items away from screen location
- mixed results compared to separate views and temporal navigation**
- fish-eye follow-up: concern with enthusiasm over distortion**
  - what* is being shown: selective filtering
  - how* it is being shown: distortion as one possibility

today . . .

**-filtering**

**-aggregation**

# Reducing Items and Attributes

## ① Filter

→ Items



→ Attributes



## ② Aggregate

→ Items



→ Attributes



why reduce?



# filter vs aggregation

# filter

elements are eliminated

→ Items



→ Attributes



# filter

elements are eliminated

→ Items



→ Attributes



# dynamic queries

# filter

elements are eliminated

→ Items



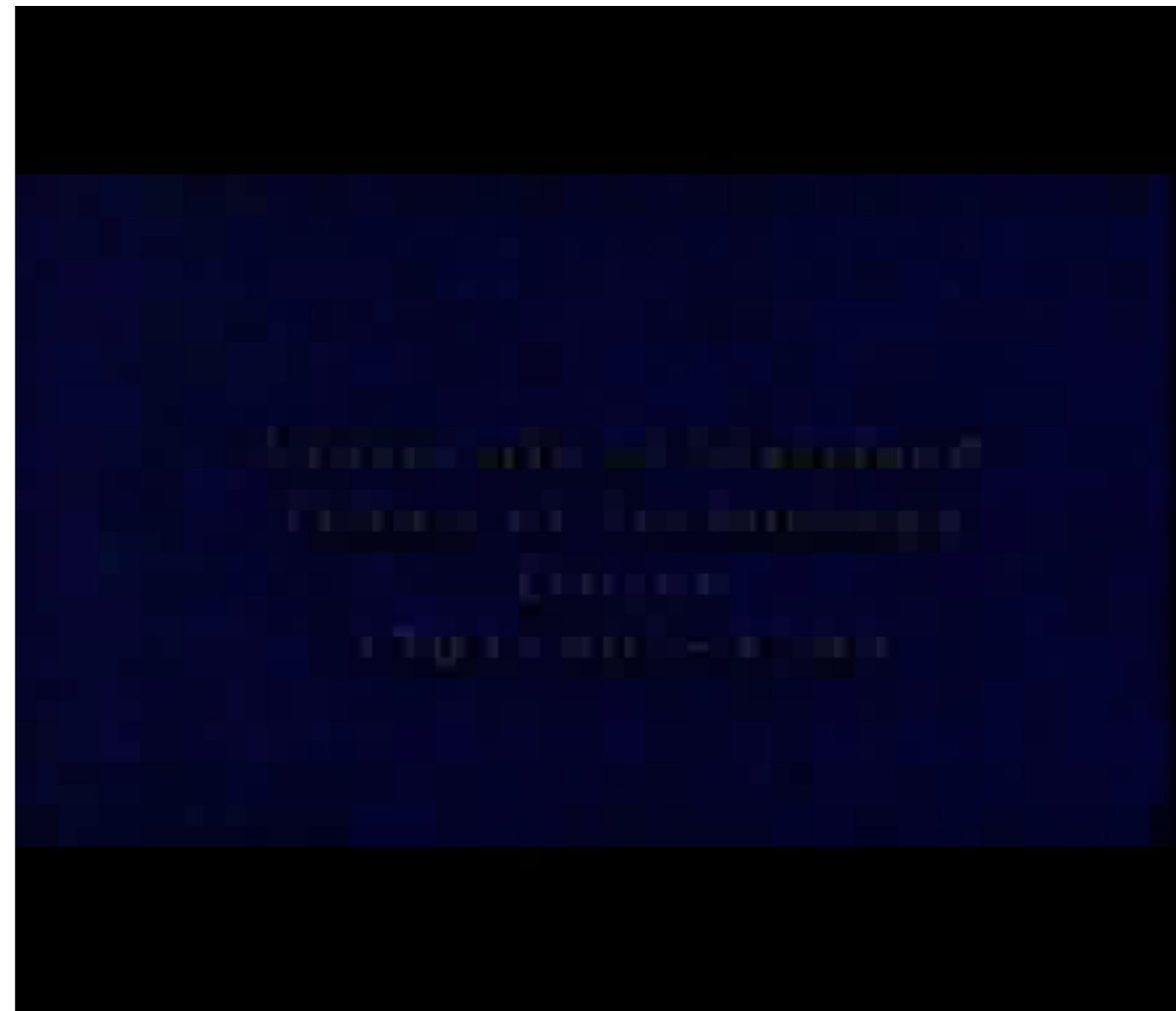
→ Attributes



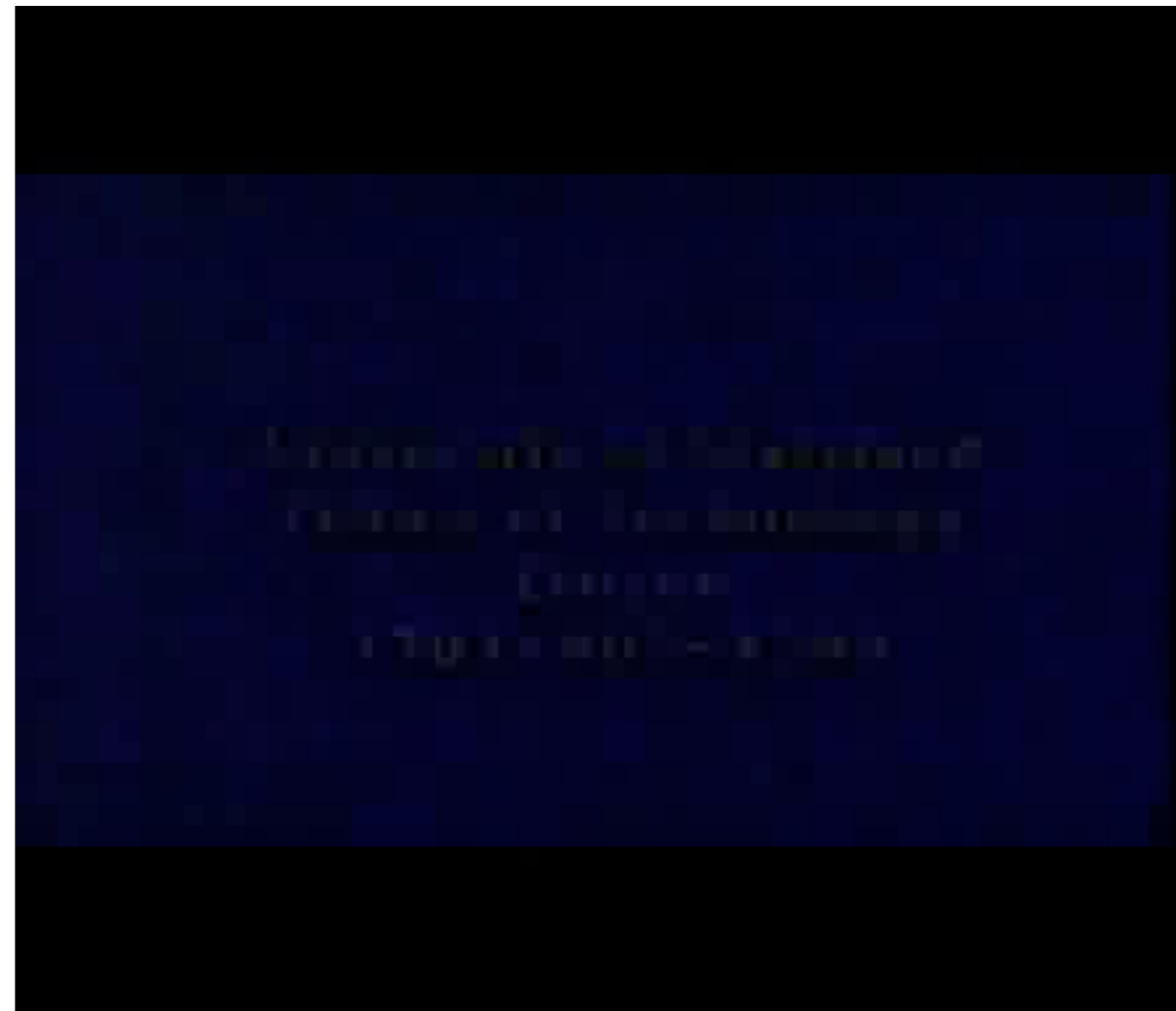
# dynamic queries

coupling between encoding and interaction so that user can immediately see the results of an action

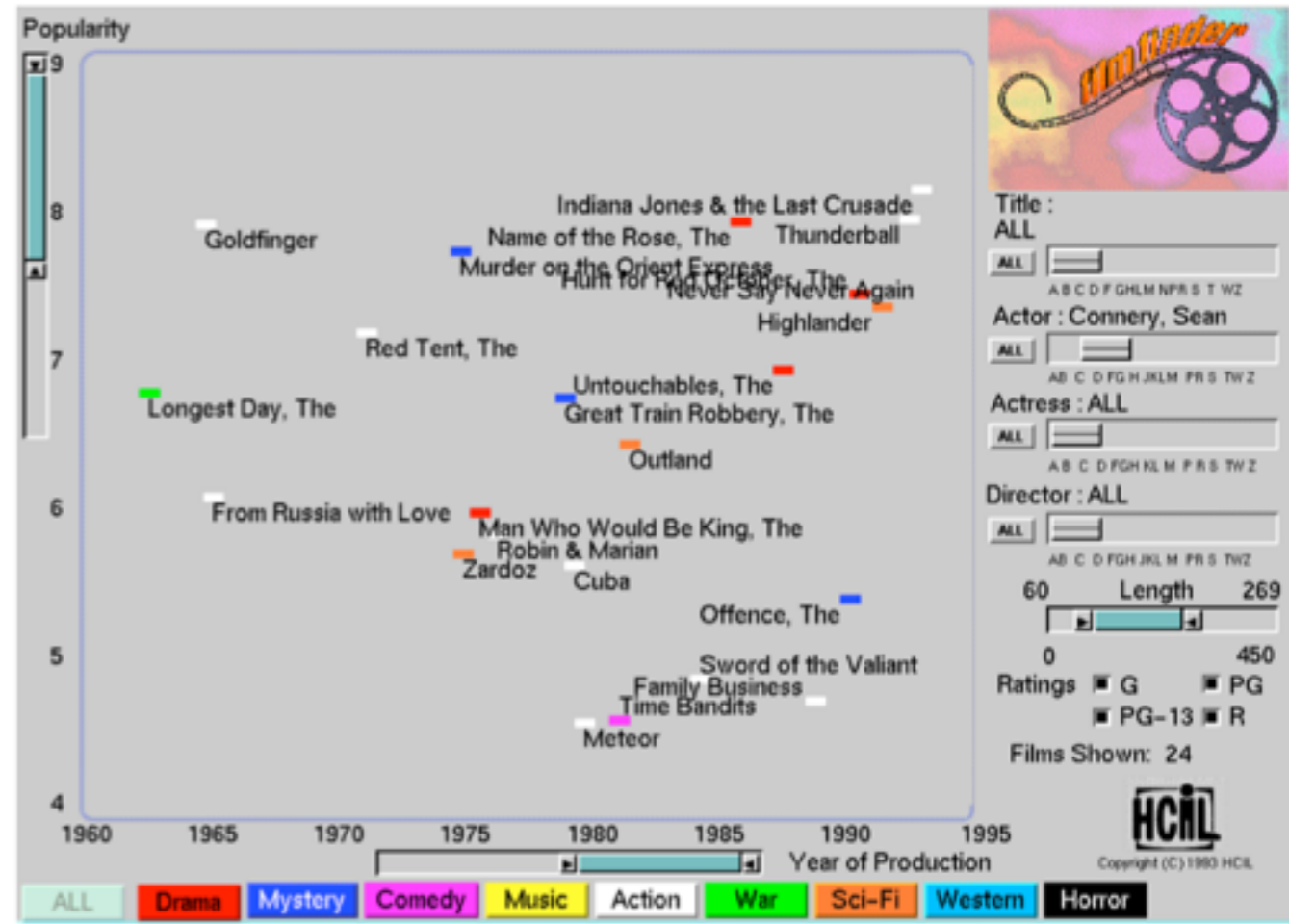
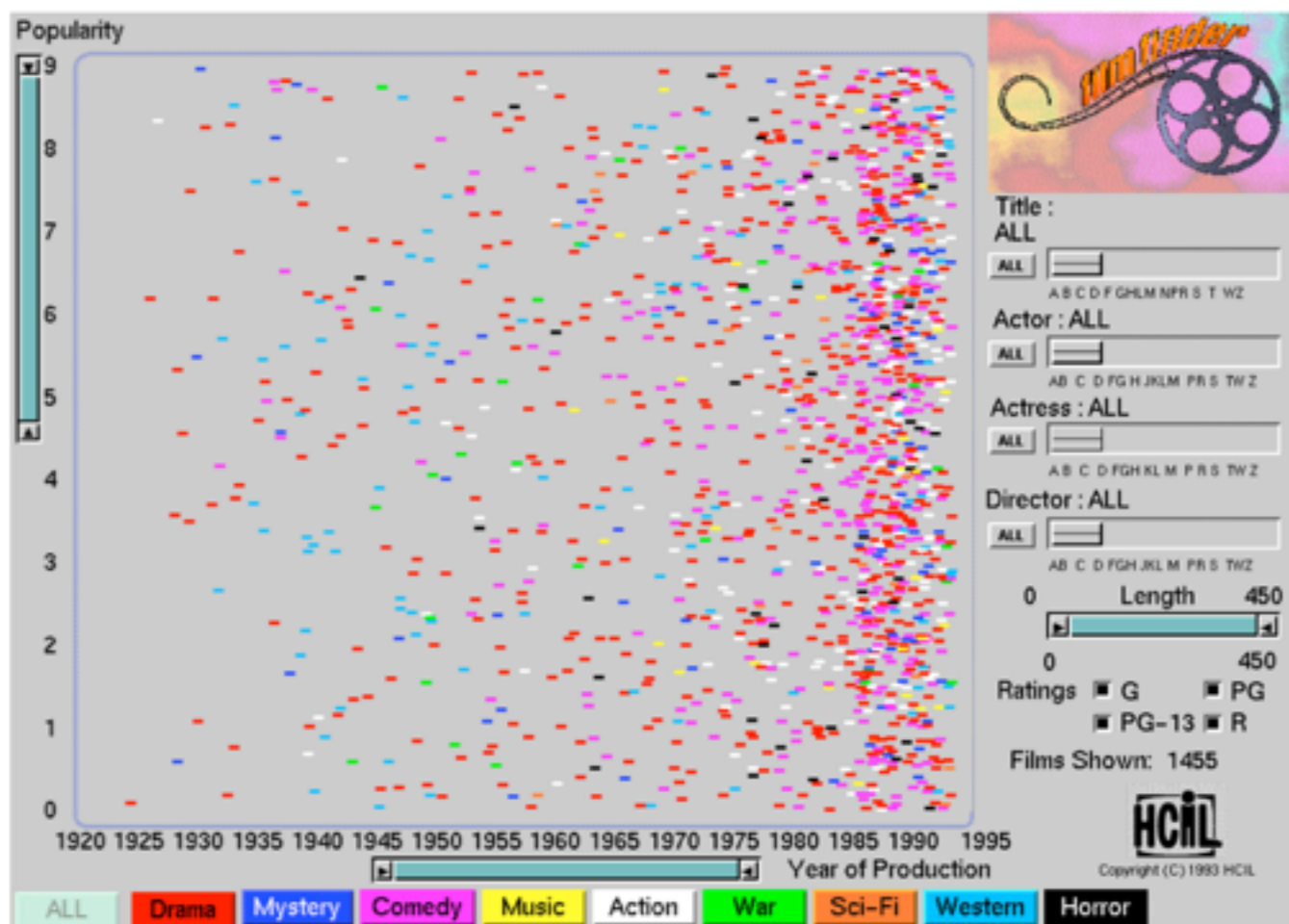
# ITEM FILTERING



# ITEM FILTERING



# ITEM FILTERING





UPDATED June 25, 2012

RECOMMEND TWITTER LINKEDIN SIGN IN TO E-MAIL SHARE

## New York Health Department Restaurant Ratings Map

The New York City Department of Health and Mental Hygiene performs unannounced sanitary inspections of every restaurant at least once per year. Violation points result in a letter grade, which can be explored in the map below, along with violation descriptions. The information on this map will be updated every two weeks. For menus and reviews by New York Times critics, visit [our restaurants guide](#). [Related Article >](#)

FIND A RESTAURANT FIND A LOCATION

FILTER

**A B C** All grades

All violations

All cuisines



Restaurant locations are derived from the New York City Department of Health and Mental Hygiene database. Due to the limitations of the Health Department's database, some restaurants could not be placed.

By JEREMY WHITE | [Send Feedback](#)

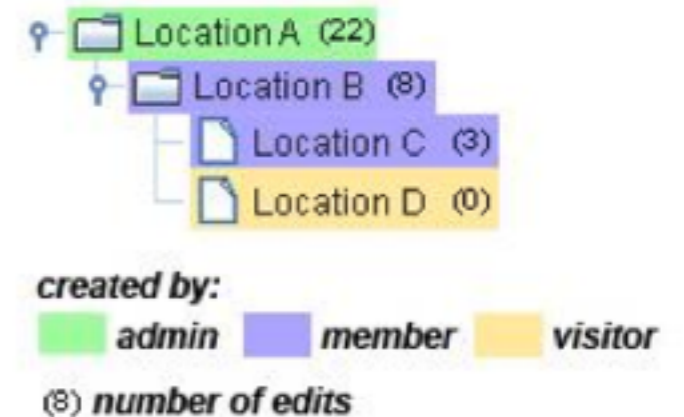
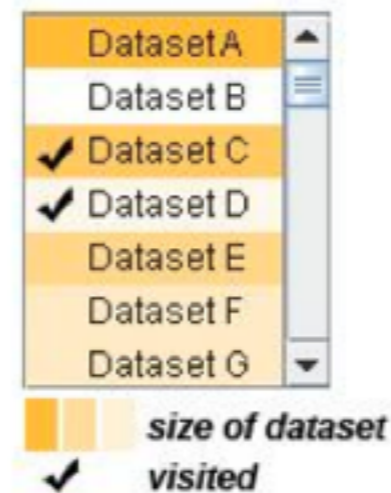
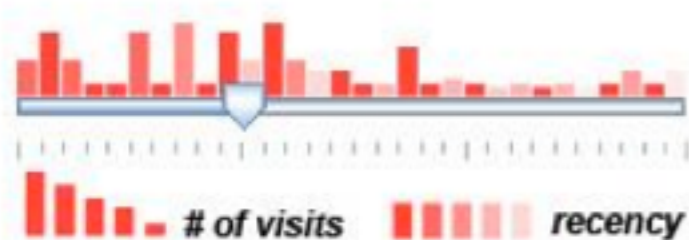
Source: New York City Department of Health and Mental Hygiene



# scented widgets

**information scent:** user's (imperfect) perception of data

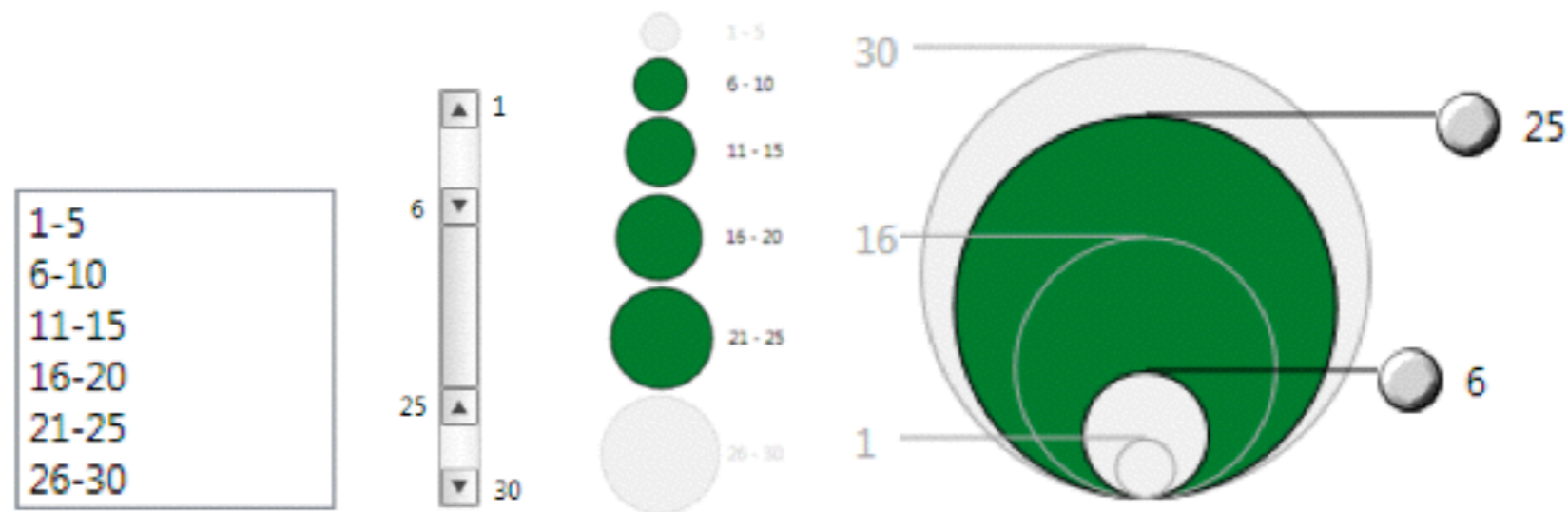
**GOAL: lower the cost of information forging through better cues**



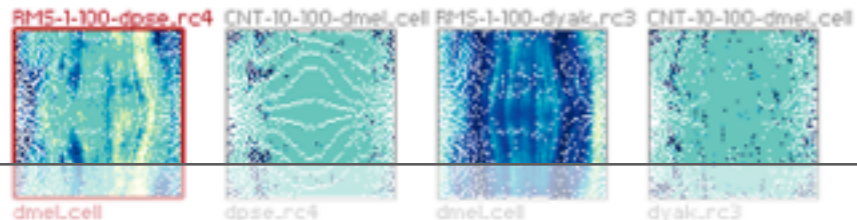
# interactive legends

controls combining the visual representation of static legends with interaction mechanisms of widgets

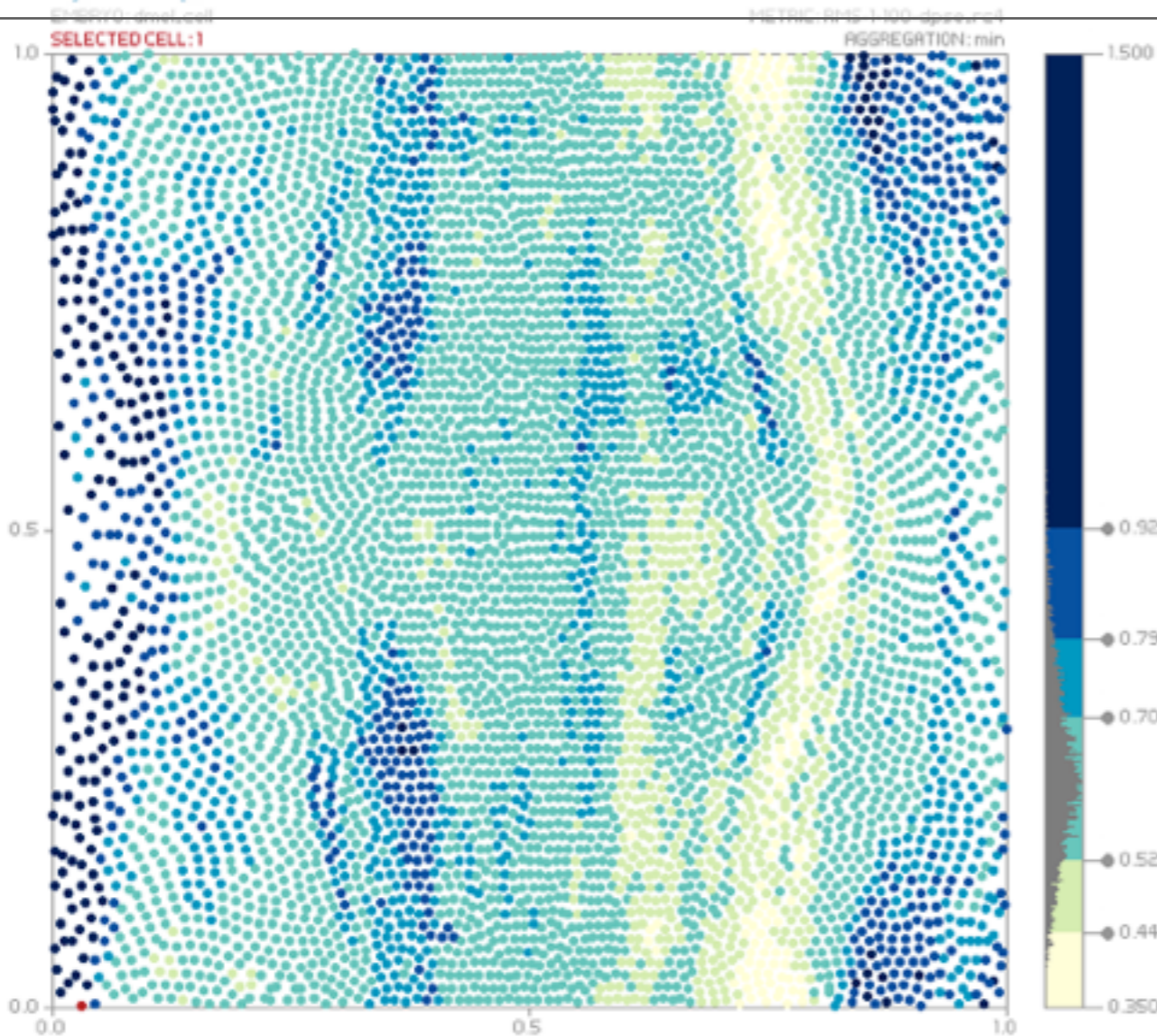
**define and control visual display together**



### Summaries



### Embryo Map



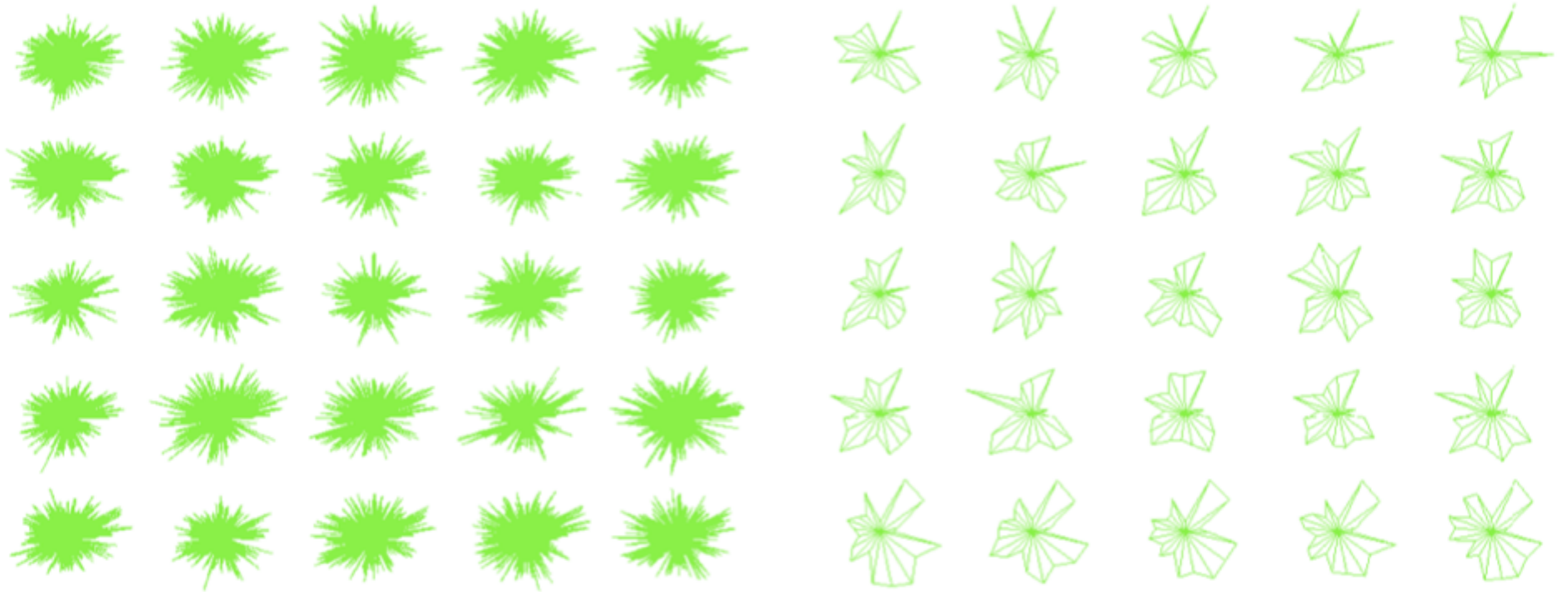
Colormaps ▲

# Demo





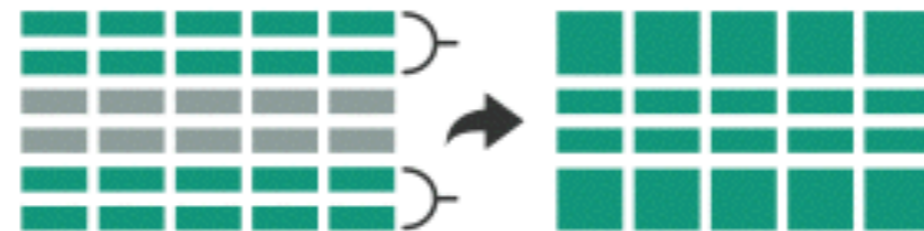
# ATTRIBUTE FILTERING



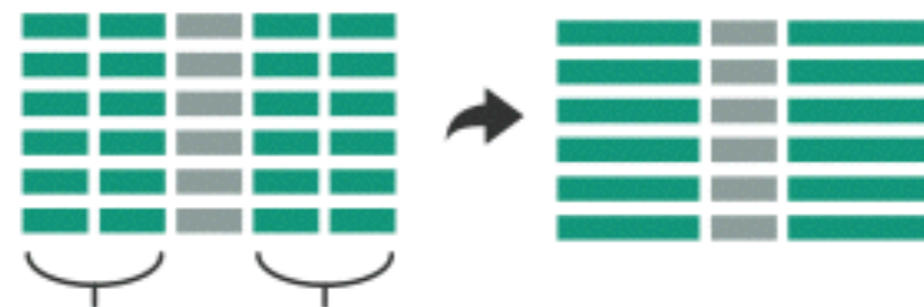
# aggregate

a group of elements is represented by a new derived element that stands in for the entire group

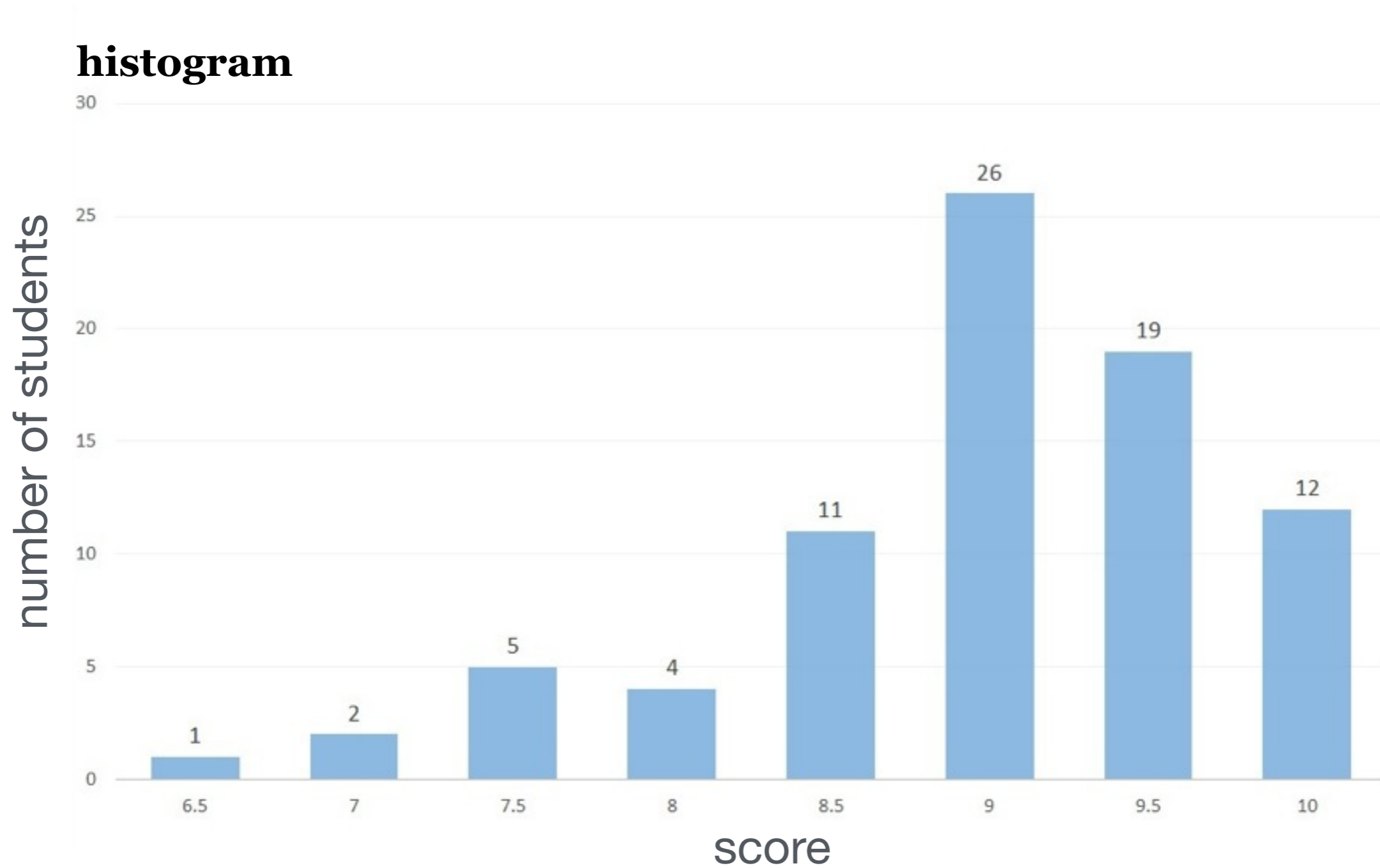
→ Items



→ Attributes

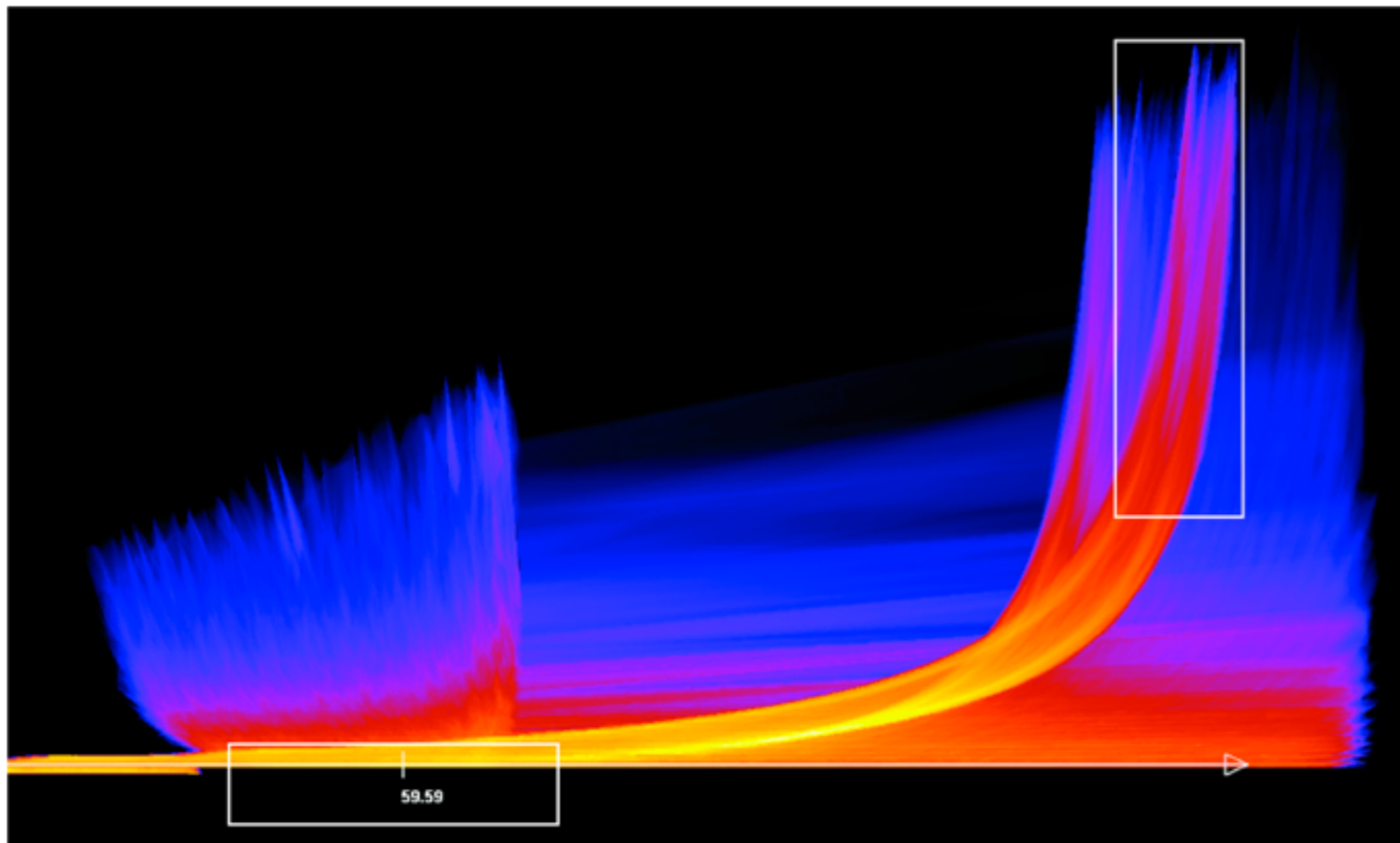


# item aggregation



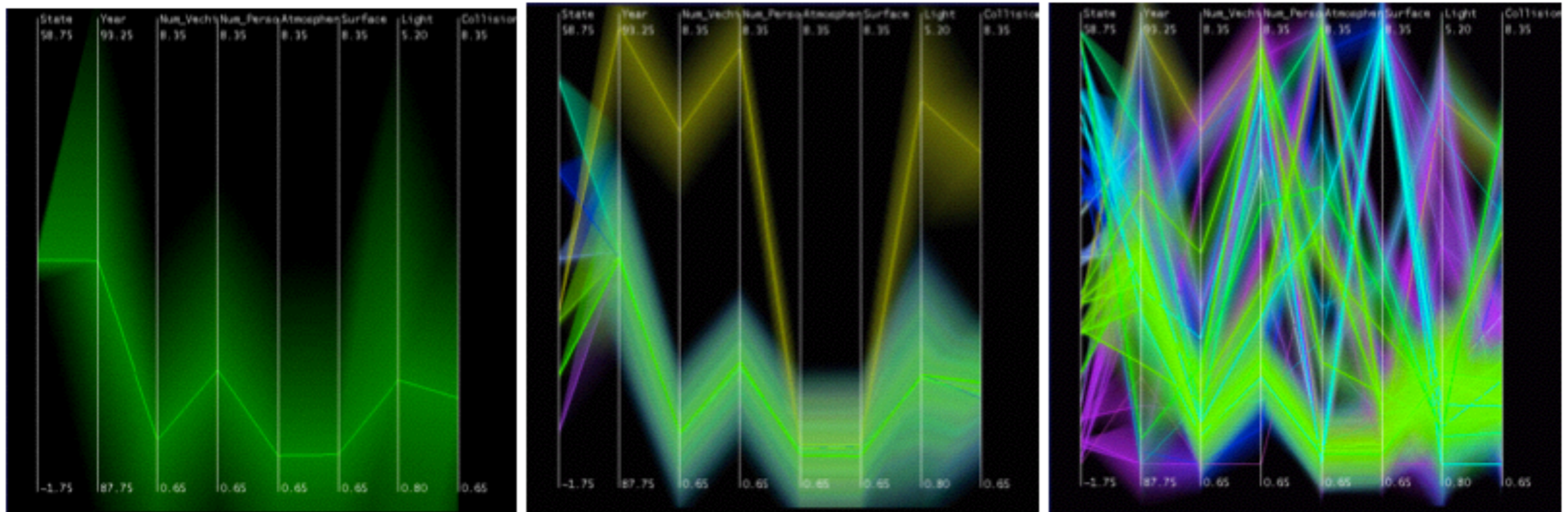
# item aggregation

continuous scatterplot



# item aggregation

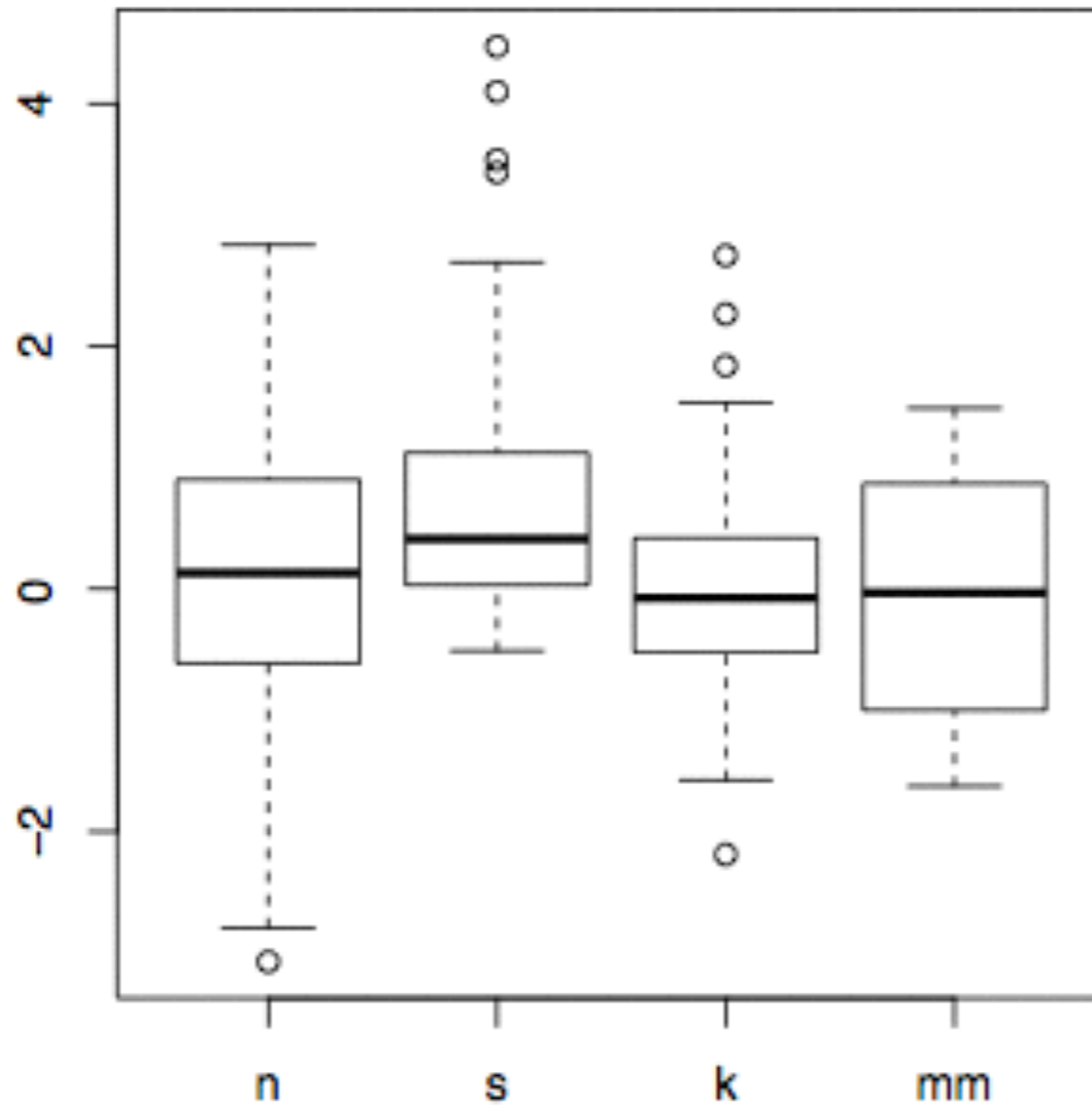
## hierarchical parallel coordinates



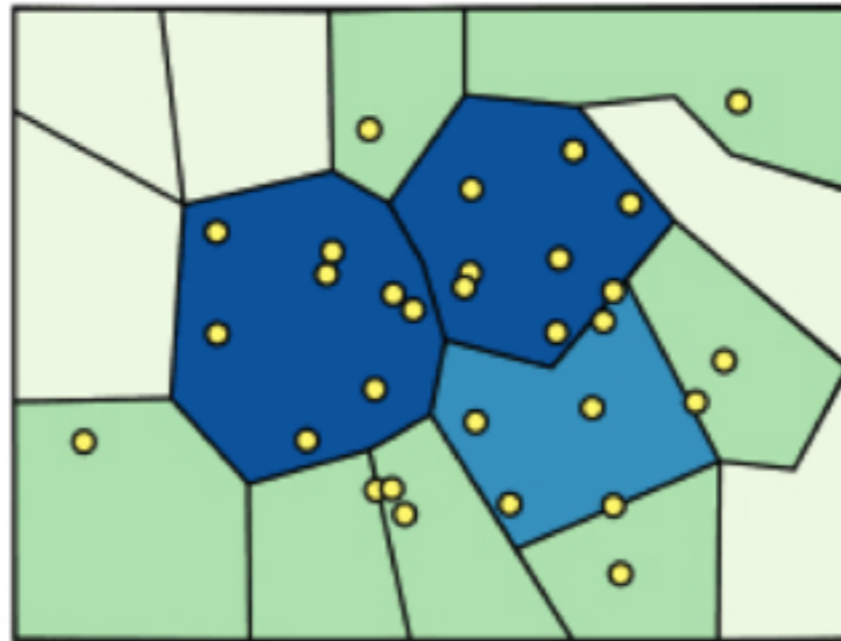


# item aggregation

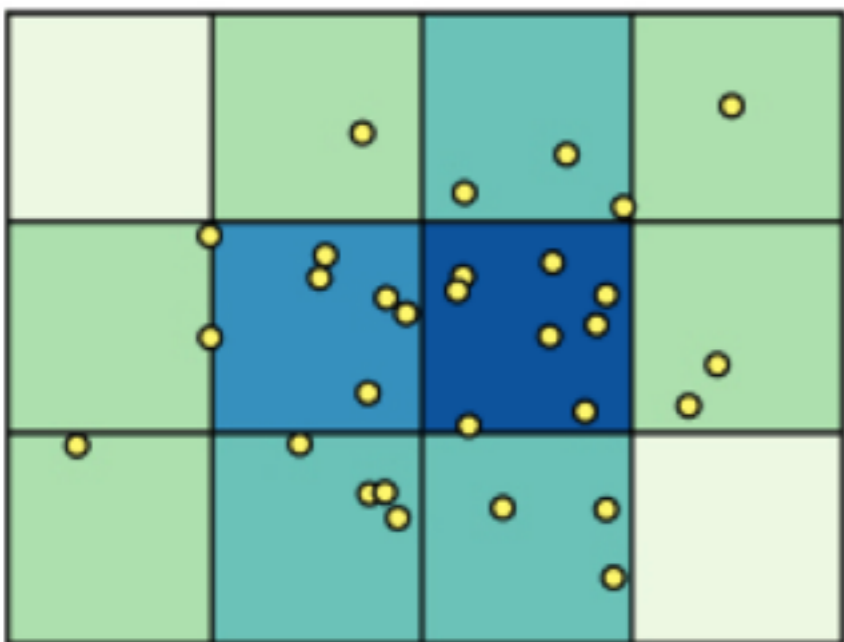
box plot



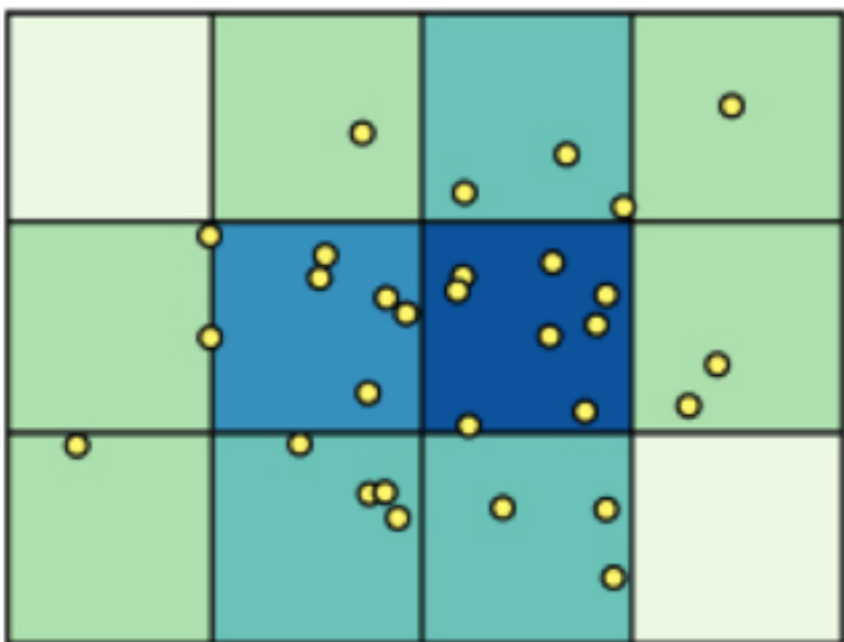
# spatial aggregation



# spatial aggregation



# spatial aggregation



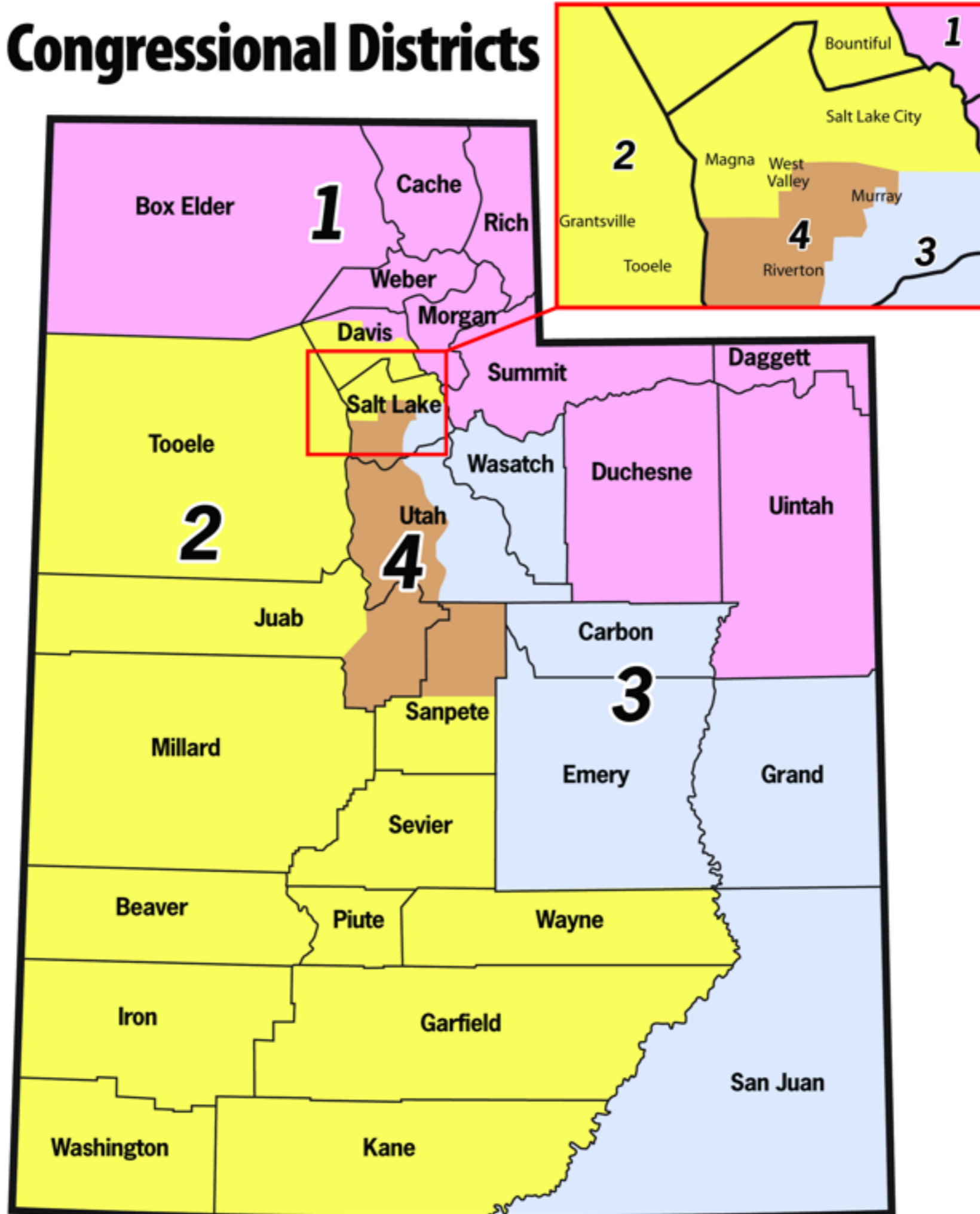
# spatial aggregation



## modifiable areal unit problem

in cartography, changing the boundaries of the regions used to analyze data can yield dramatically different results

# Congressional Districts



# attribute aggregation

- 1) group attributes and compute a similarity score across the set
- 2) dimensionality reduction, to preserve meaningful structure

# attribute aggregation

- 1) group attributes and compute a similarity score across the set**
- 2) dimensionality reduction, to preserve meaningful structure



# MulteeSum

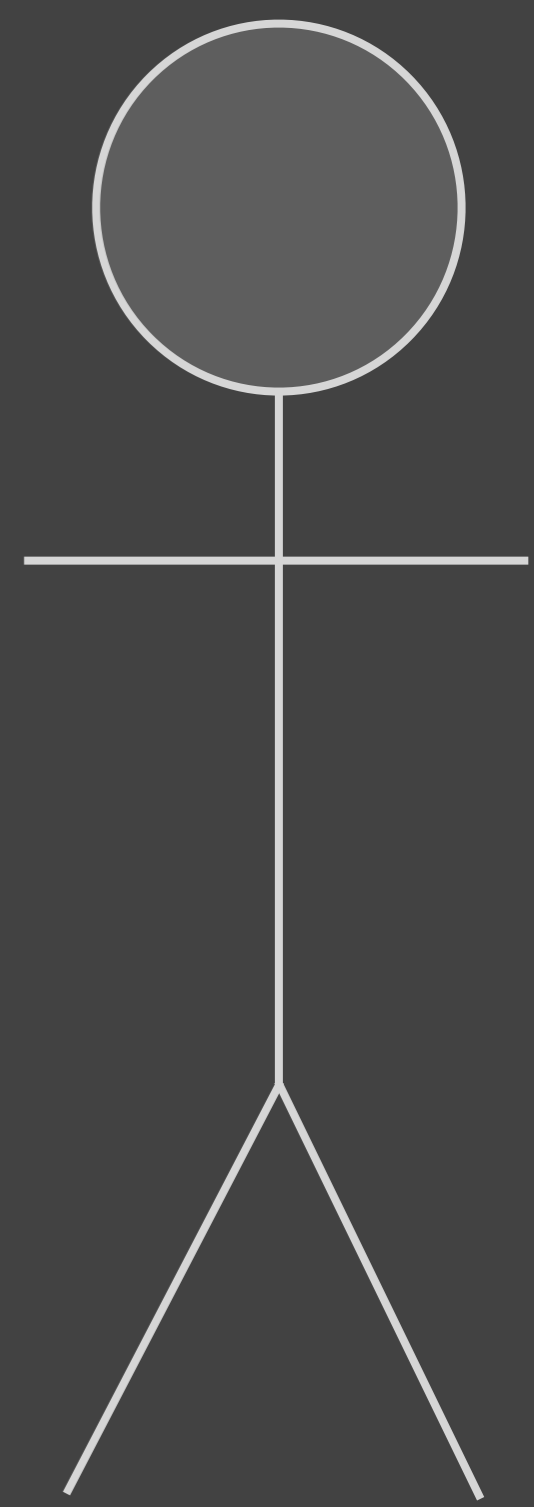
A Tool for Comparative Spatial and Temporal Gene Expression Data

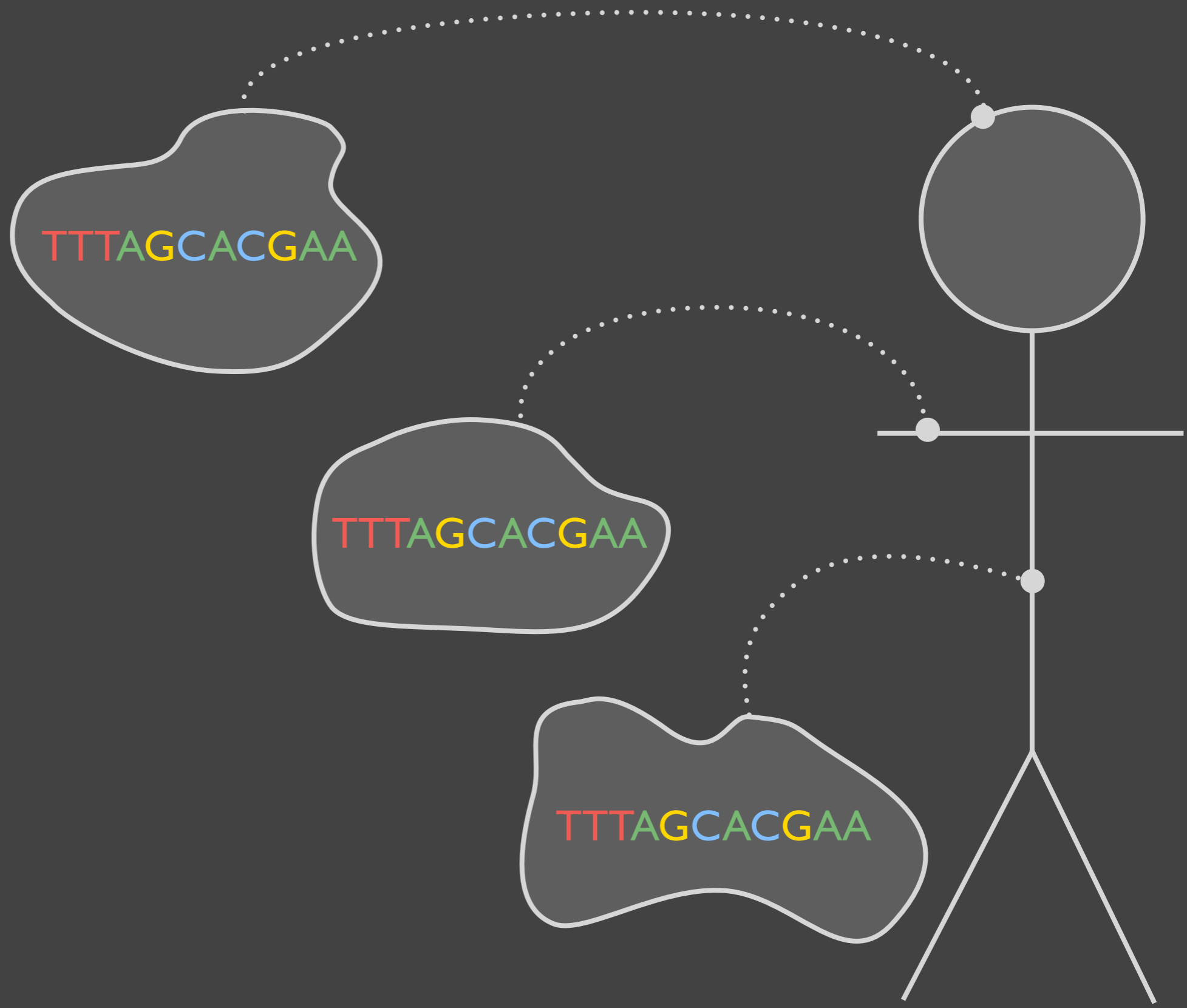
Miriah Meyer,<sup>1</sup> Tamara Munzner,<sup>2</sup> Angela DePace,<sup>3</sup> Hanspeter Pfister<sup>1</sup>

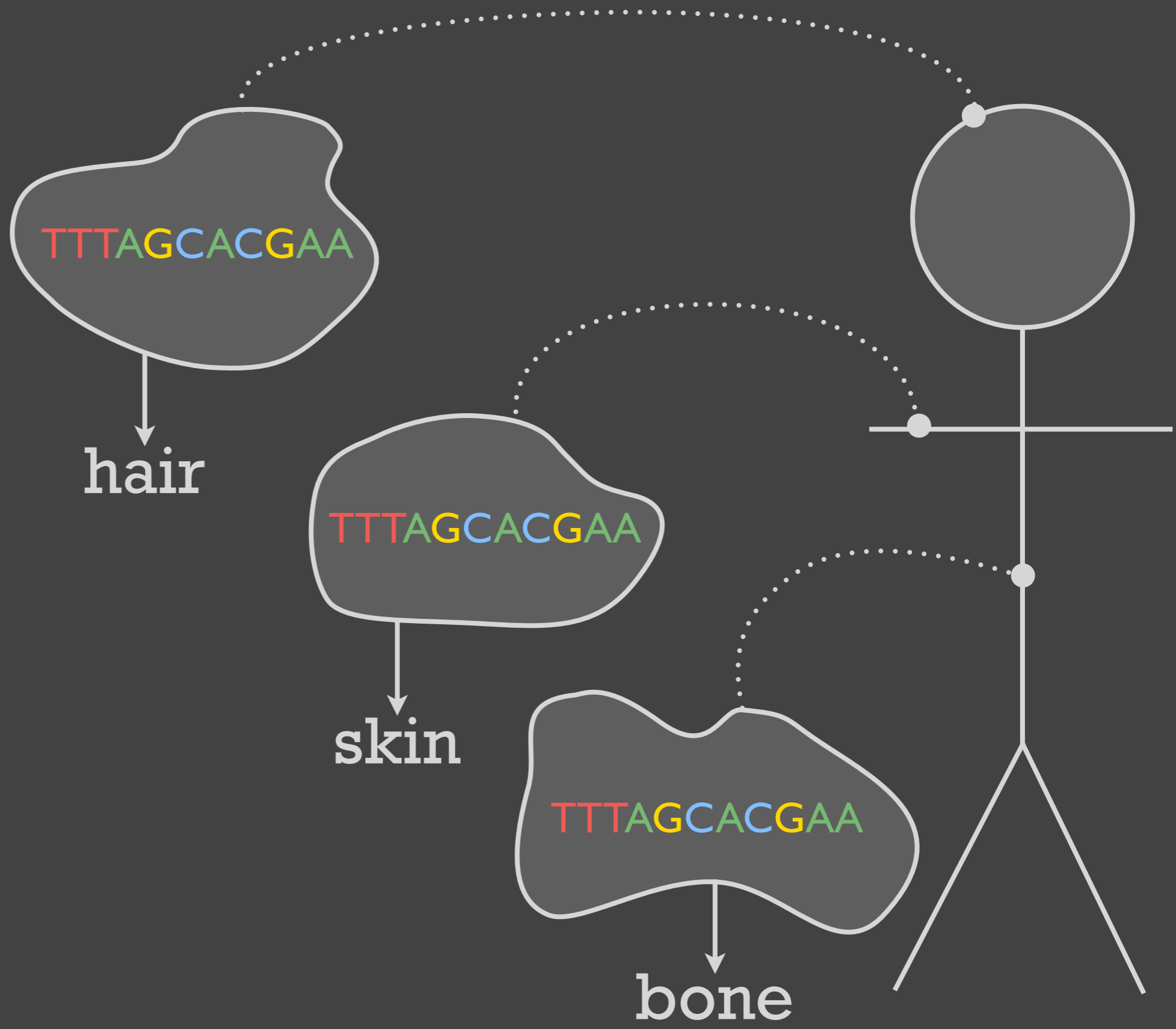
<sup>1</sup> Harvard University

<sup>2</sup> University of British Columbia

<sup>3</sup> Harvard Medical School

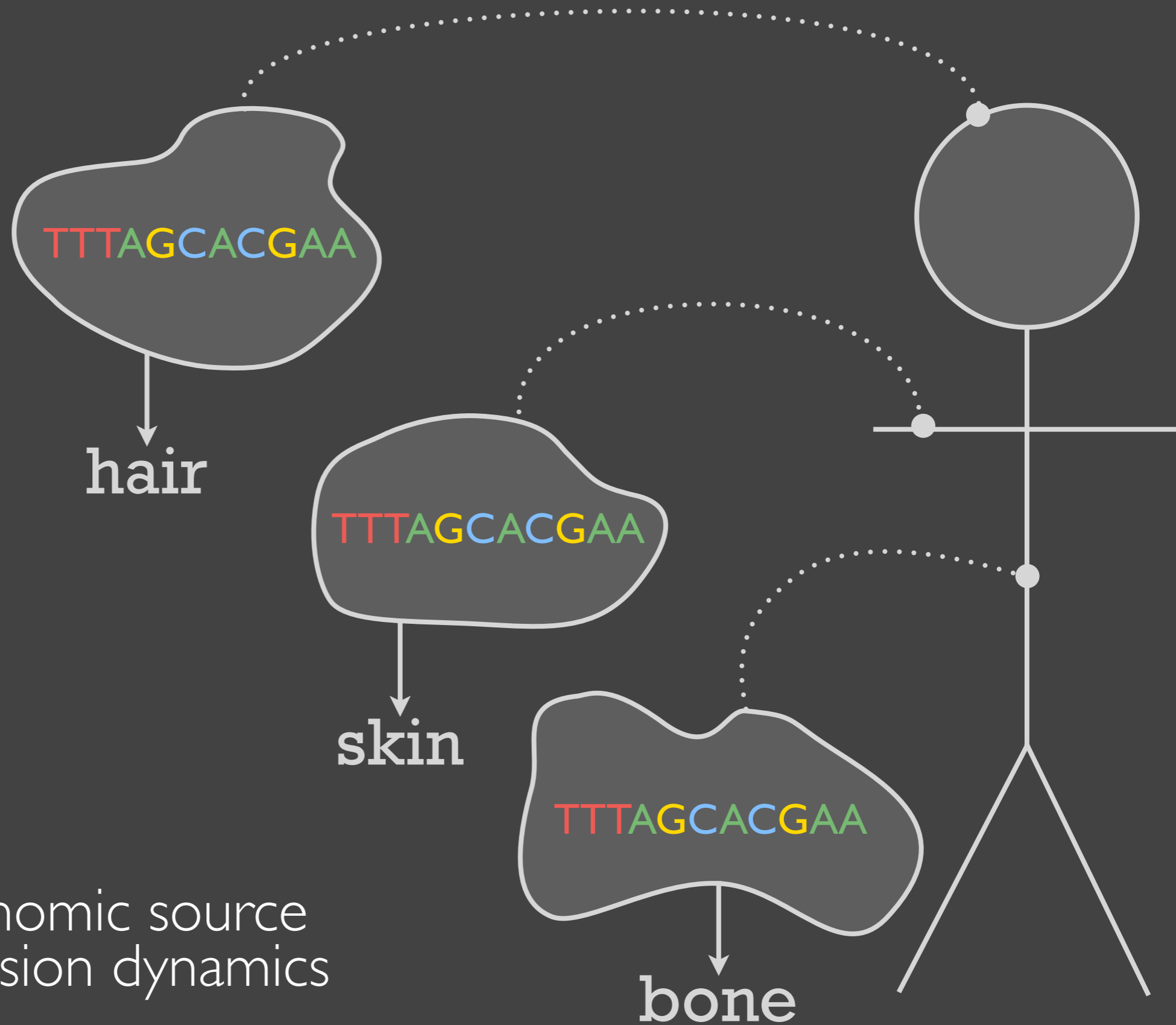






**grand  
challenge  
in biology:**

understand genomic source  
of gene expression dynamics



collaborators: DePace lab at Harvard Medical School

model species: fruit fly





collaborators: DePace lab at Harvard Medical School

model species: fruit fly

scientific goal: link changes in the regulatory part of the genome to species variation



collaborators: DePace lab at Harvard Medical School

model species: fruit fly

scientific goal: link changes in the regulatory part of the genome to species variation

requires: characterize differences in gene expression patterns between species





collaborators: DePace lab at Harvard Medical School

model species: fruit fly

scientific goal: link changes in the regulatory part of the genome to species variation

requires: characterize differences in gene expression patterns between species

**MulteeSum**



# process

two year collaboration

two early prototype systems

feedback from six biologists

*informal interviews, emails*

*one day a week in biology lab*

tool deployed

*currently used several times a week*

data & tool & tasks

summaries & groups

encodings & interaction

conclusions

gene expression is

**gene expression is ...**

*... the measured level of how much a gene is on or off.*

**gene expression is ...**

*... the measured level of how much a gene is on or off.*

*... a single quantitative value.*

0.2

**gene expression is ...**

*... the measured level of how much a gene is on or off.*

*... a single quantitative value.*



**collaborators measure it ...**

*... for multiple time points.*

**gene expression is ...**

*... the measured level of how much a gene is on or off.*

*... a single quantitative value.*

**collaborators measure it ...**

*... for multiple time points.*

*... for multiple genes.*

time →

genes ↓

0.2	0.4	1	1	1	0.8
1	0	0	0	1	1
0.7	0.8	1	1	0.8	0.6
1	0	0.2	0.5	1	1
0.5	0.8	0.5	0.3	0.5	0.8
0.7	0.5	0.8	0.7	1	1
1	0.3	0.4	1	1	1



**gene expression is ...**

*... the measured level of how much a gene is on or off.*

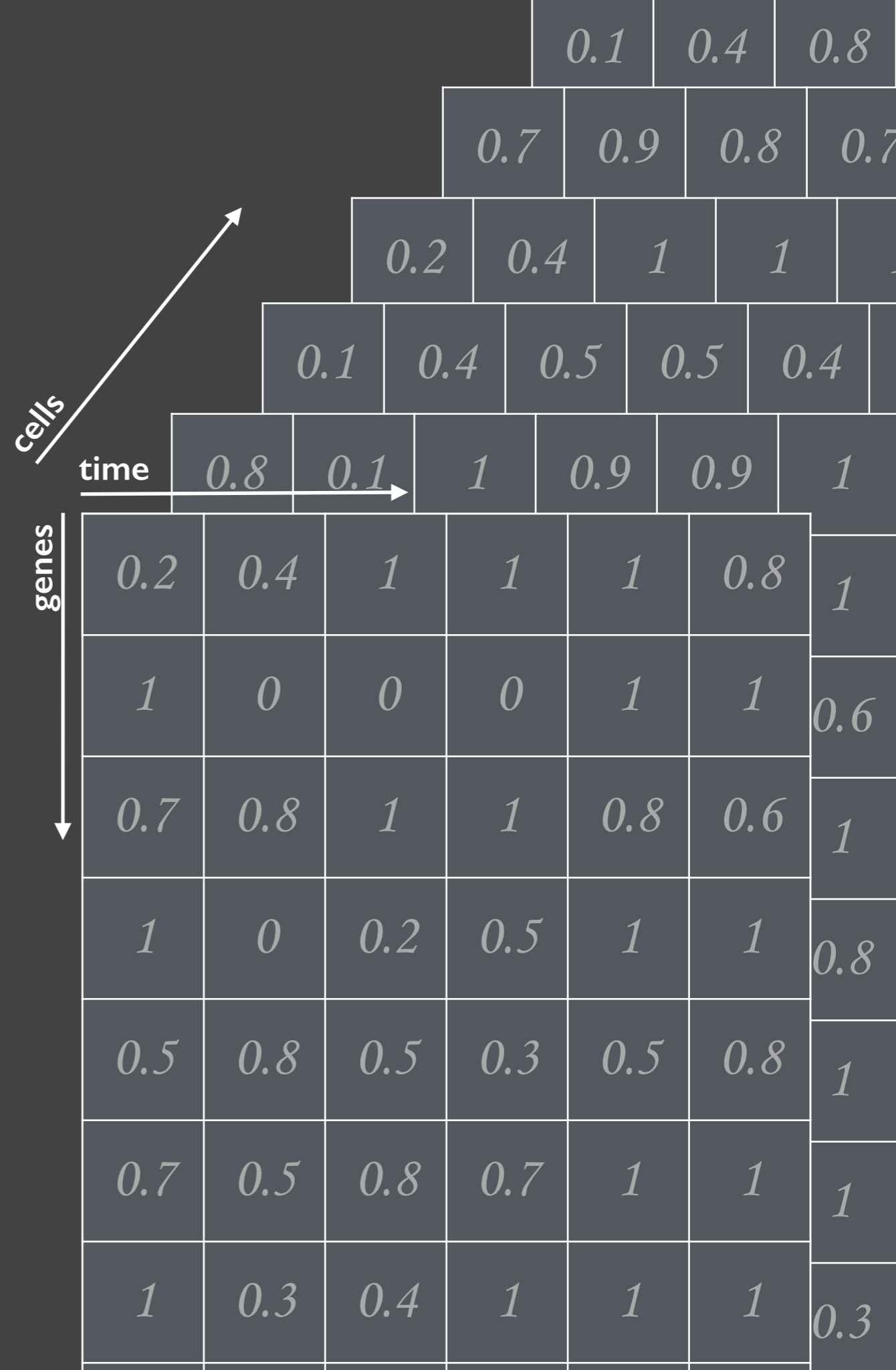
*... a single quantitative value.*

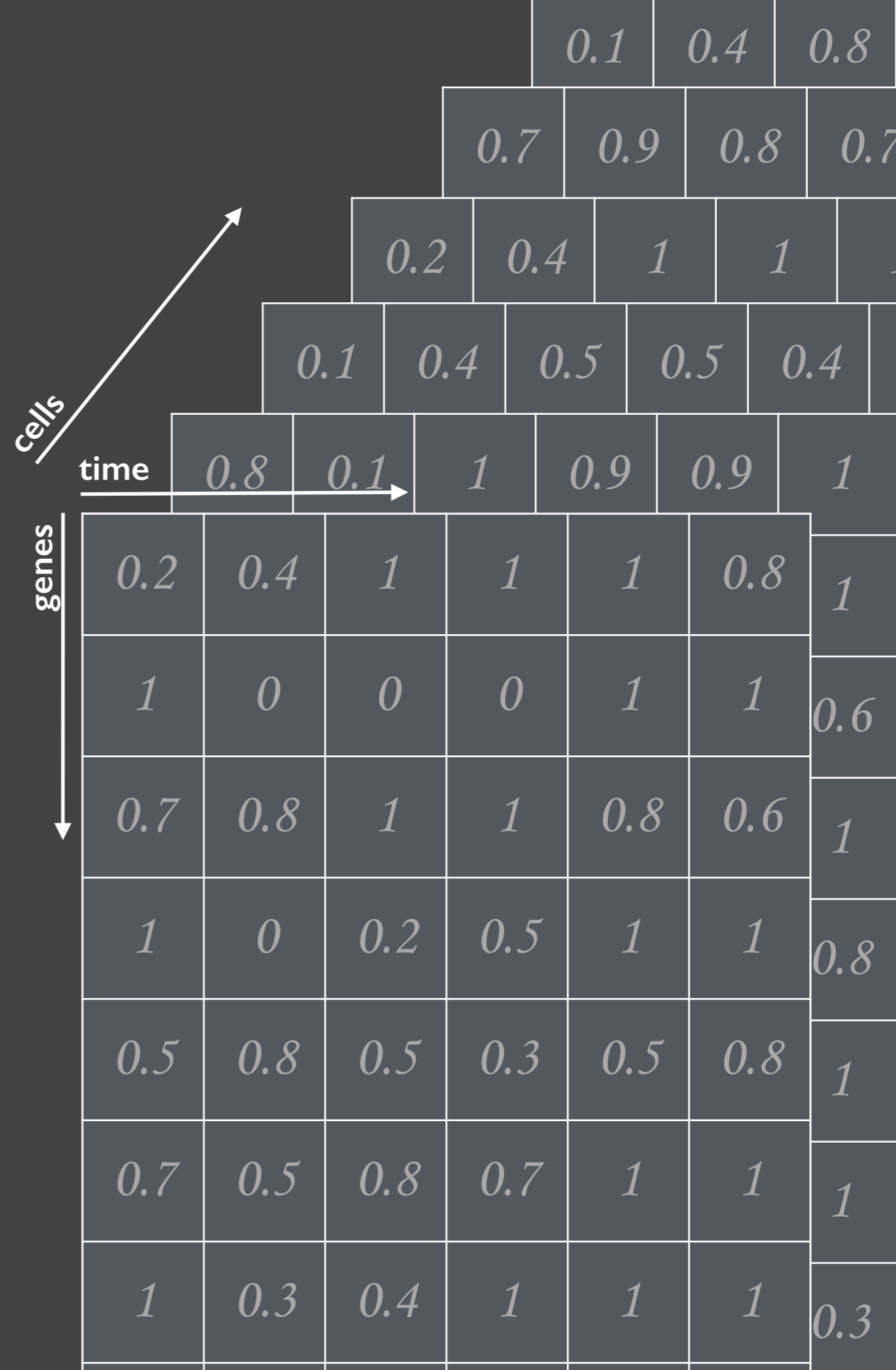
**collaborators measure it ...**

*... for multiple time points.*

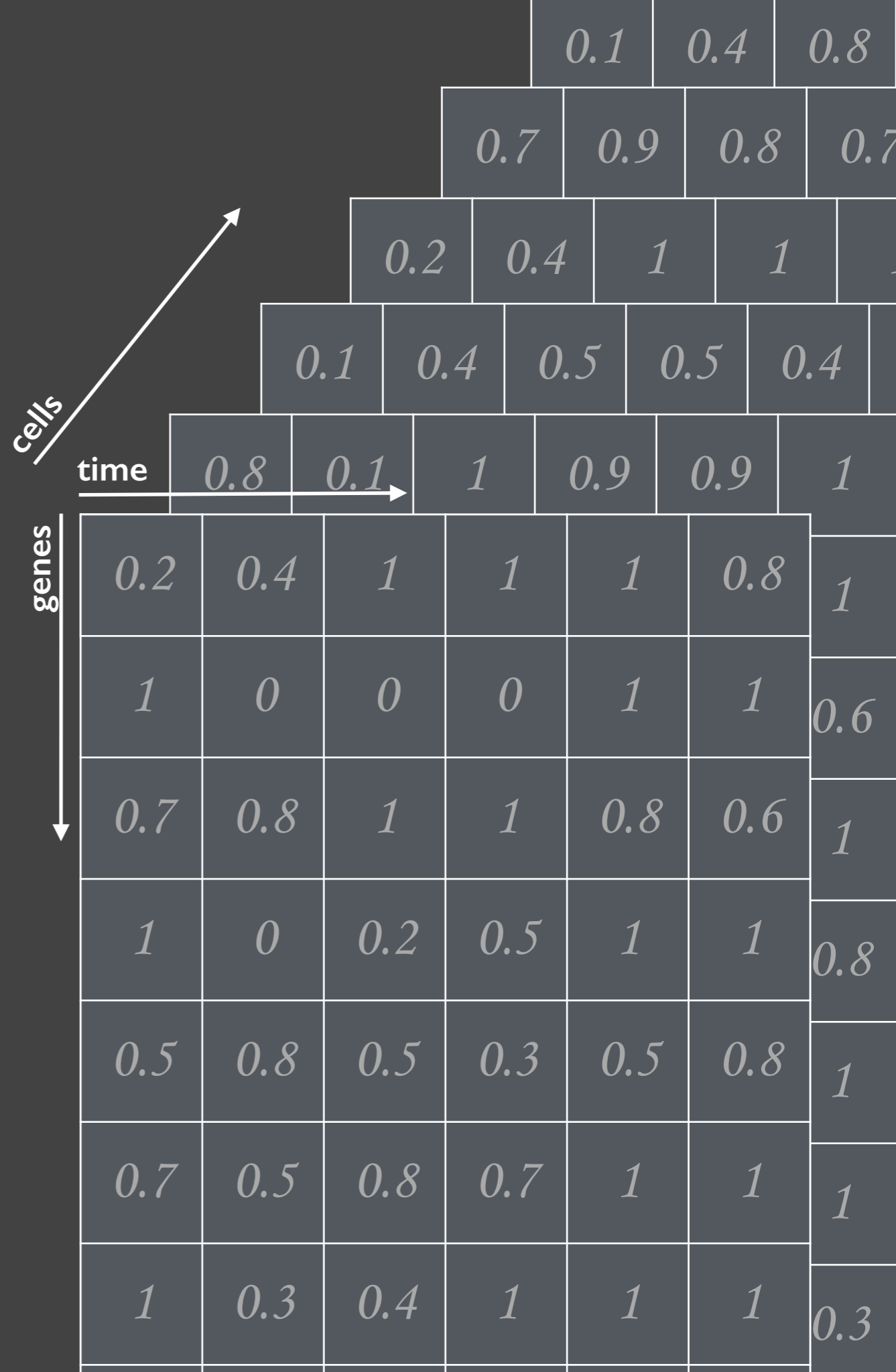
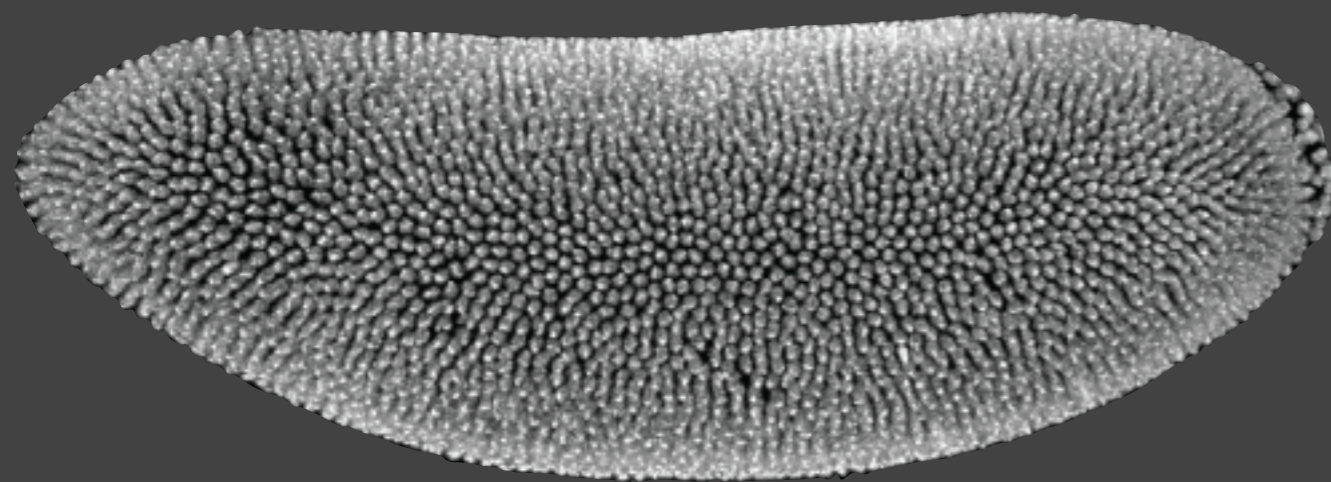
*... for multiple genes.*

*... in many cells.*

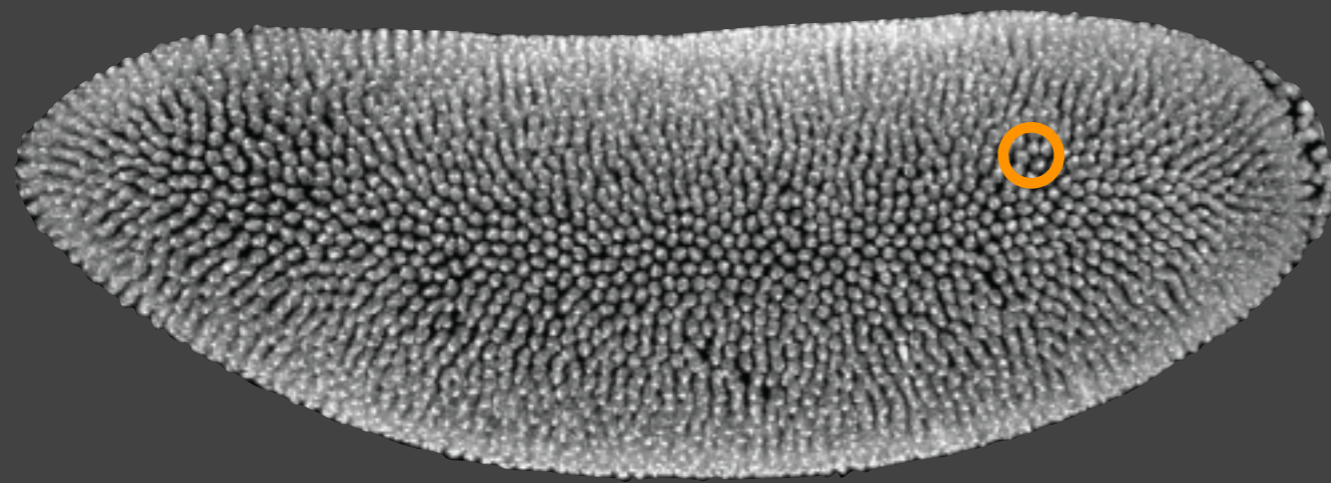




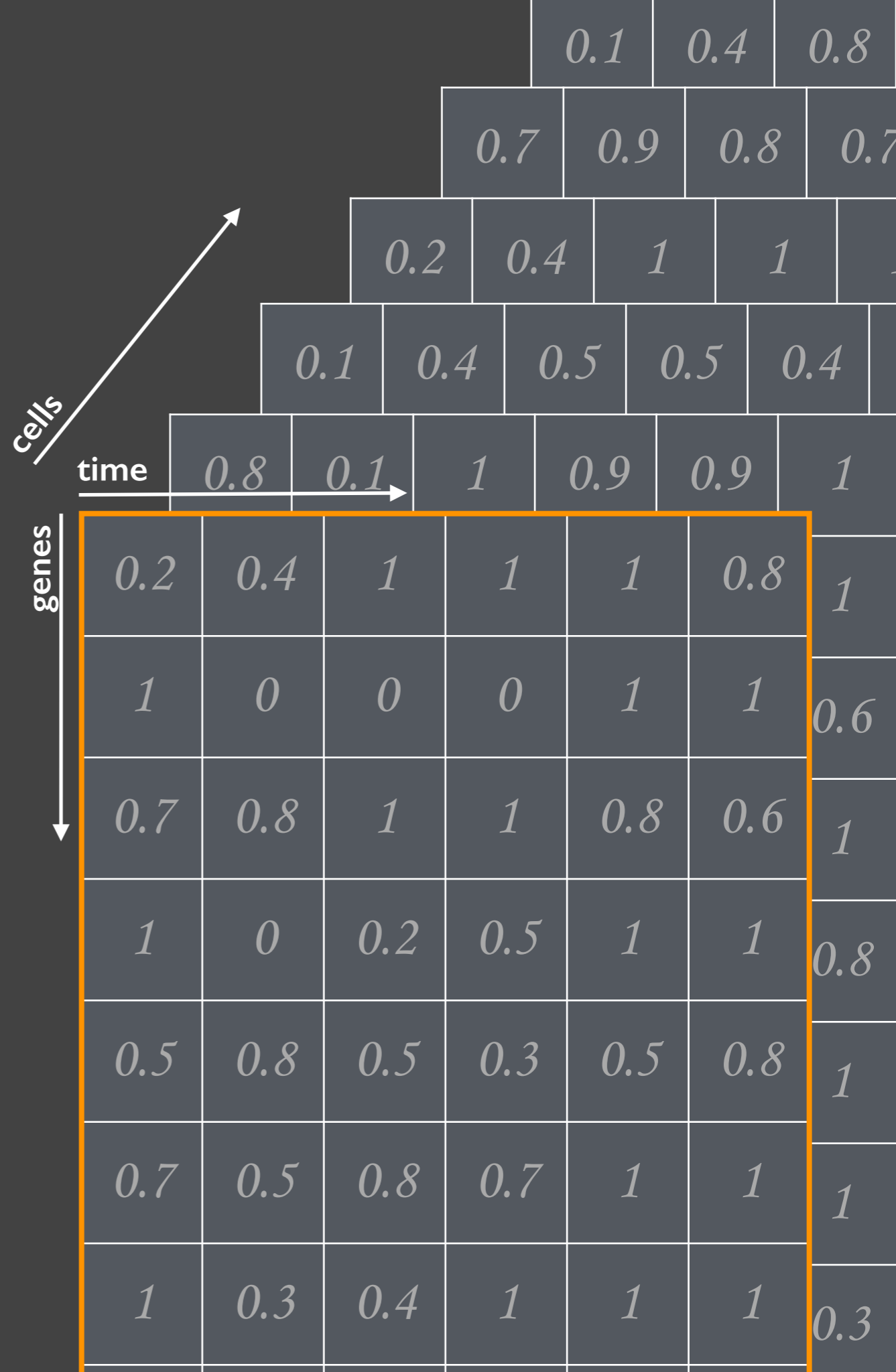
measurements taken for every cell in a single fruit fly embryo.



measurements taken for every cell in a single fruit fly embryo.



correlate gene expression with spatial location.



**virtual embryo**

# virtual embryo

several thousand cells

[5,000  $\pm$  1,000]

# virtual embryo

**several thousand cells**  
[5,000 ± 1,000]

each cell has:

# virtual embryo

several thousand cells  
[5,000  $\pm$  1,000]

each cell has:  
**expression profile**  
[6 time points x 50 genes]

time  $\rightarrow$

genes  $\downarrow$

0.2	0.4	1	1	1	0.8
1	0	0	0	1	1
0.7	0.8	1	1	0.8	0.6
1	0	0.2	0.5	1	1
0.5	0.8	0.5	0.3	0.5	0.8
0.7	0.5	0.8	0.7	1	1
1	0.3	0.4	1	1	1
0.5	0	0	0.7	0.5	0.3



# virtual embryo

**several thousand cells**  
[5,000  $\pm$  1,000]

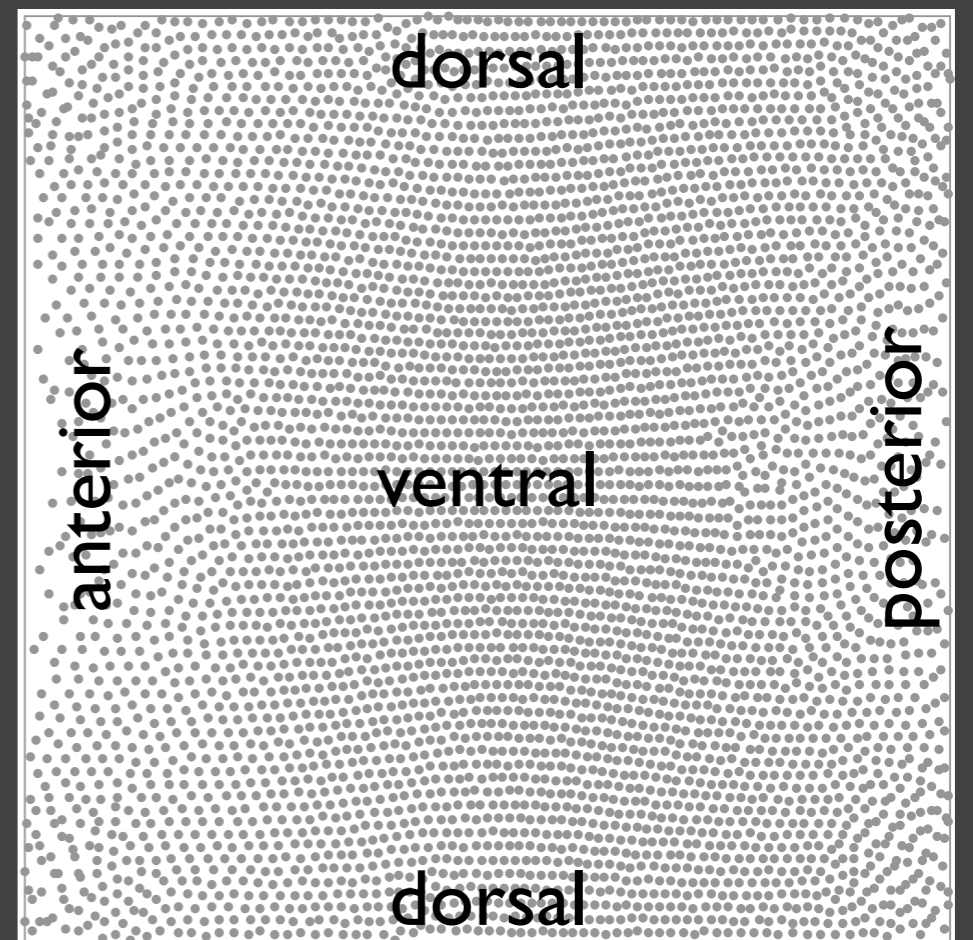
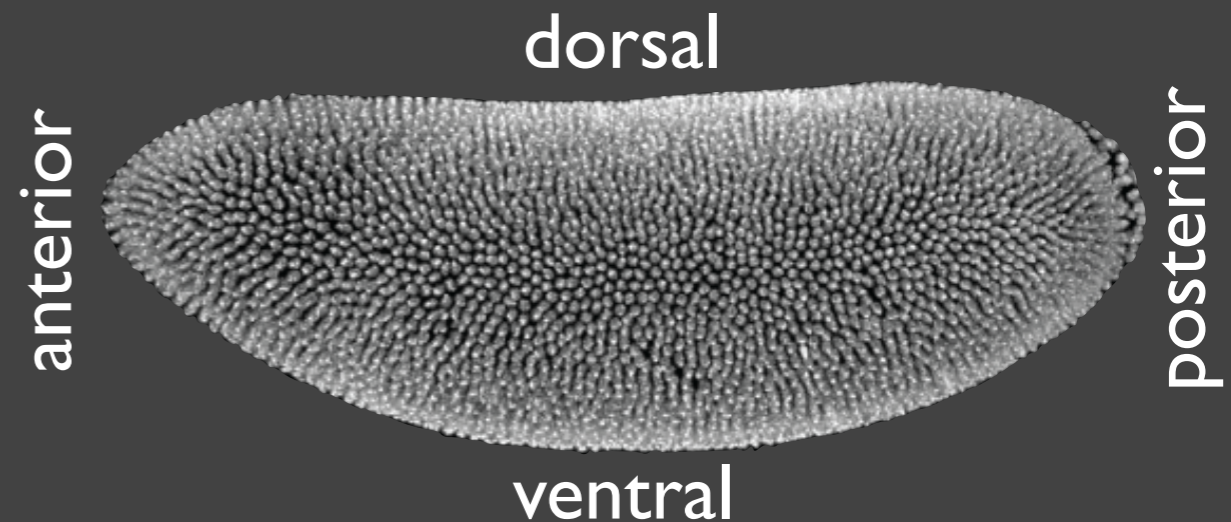
each cell has:

**expression profile**

[6 time points x 50 genes]

**spatial position**

[3D and 2D coordinates]



# virtual embryo

several thousand cells  
[5,000 ± 1,000]

each cell has:

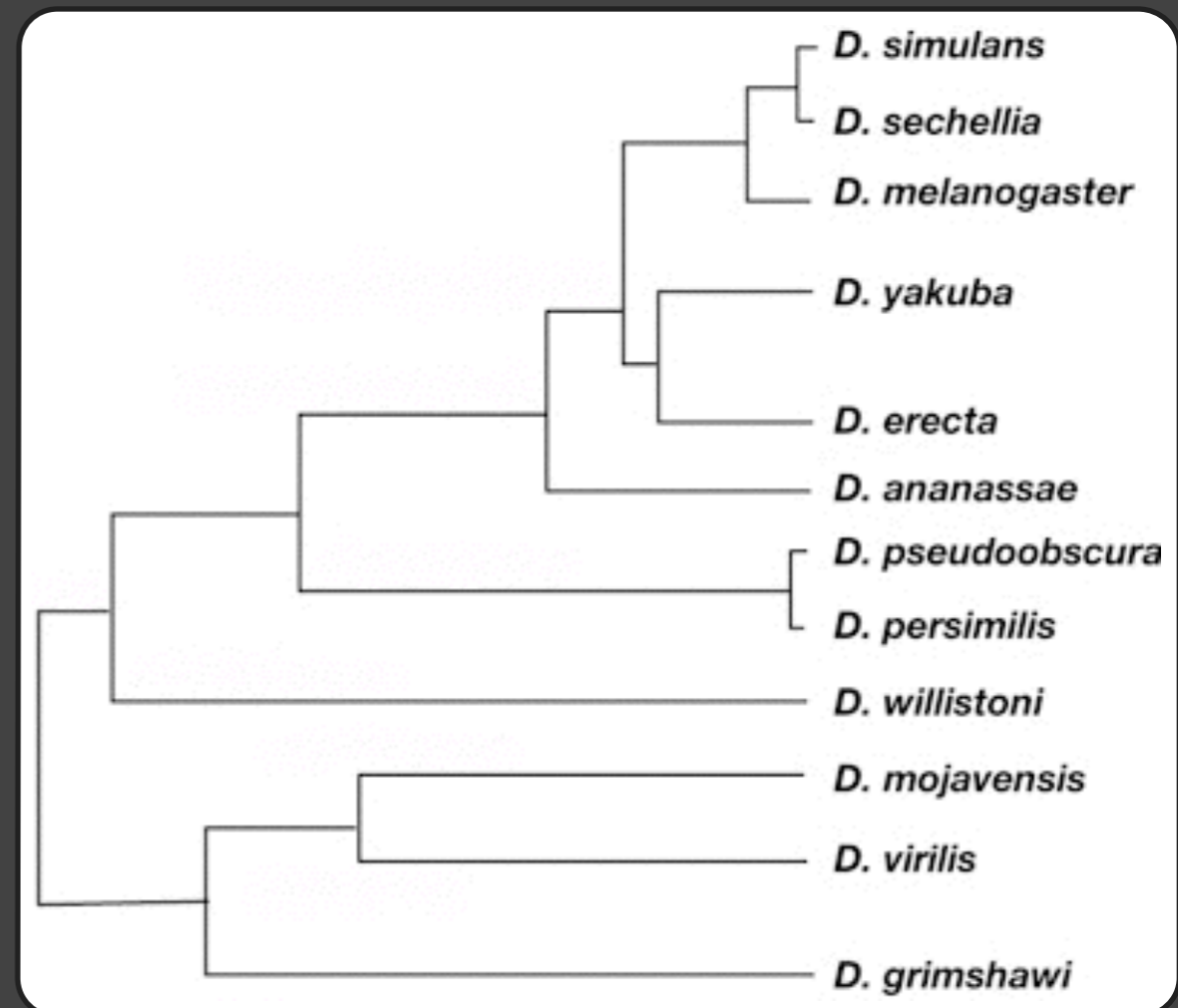
**expression profile**

[6 time points x 50 genes]

**spatial position**

[3D and 2D coordinates]

**12 related species**



# virtual embryo

several thousand cells  
[5,000 ± 1,000]

each cell has:

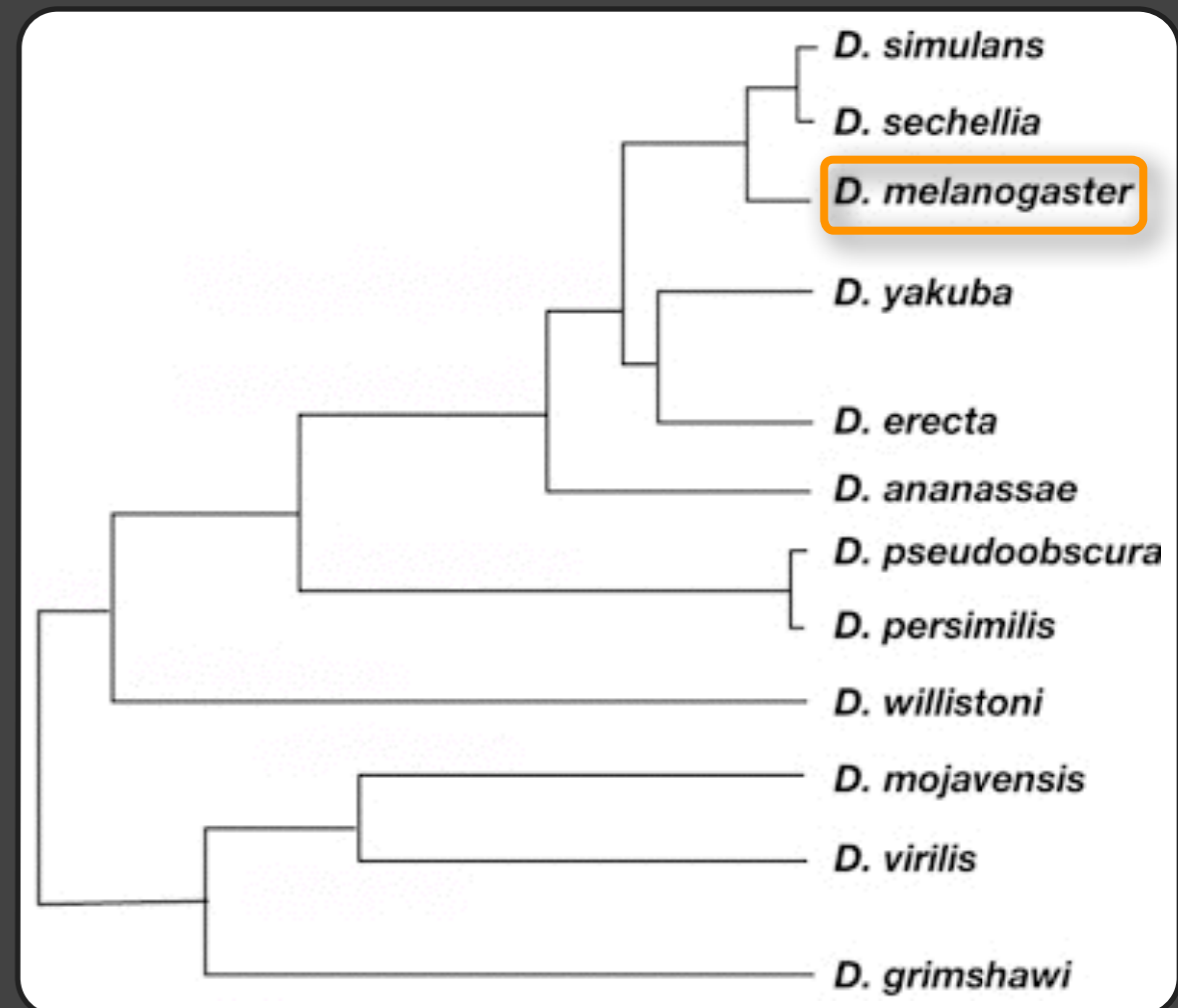
**expression profile**

[6 time points x 50 genes]

**spatial position**

[3D and 2D coordinates]

**12 related species**  
one complete



# virtual embryo

several thousand cells  
[5,000 ± 1,000]

each cell has:

**expression profile**

[6 time points x 50 genes]

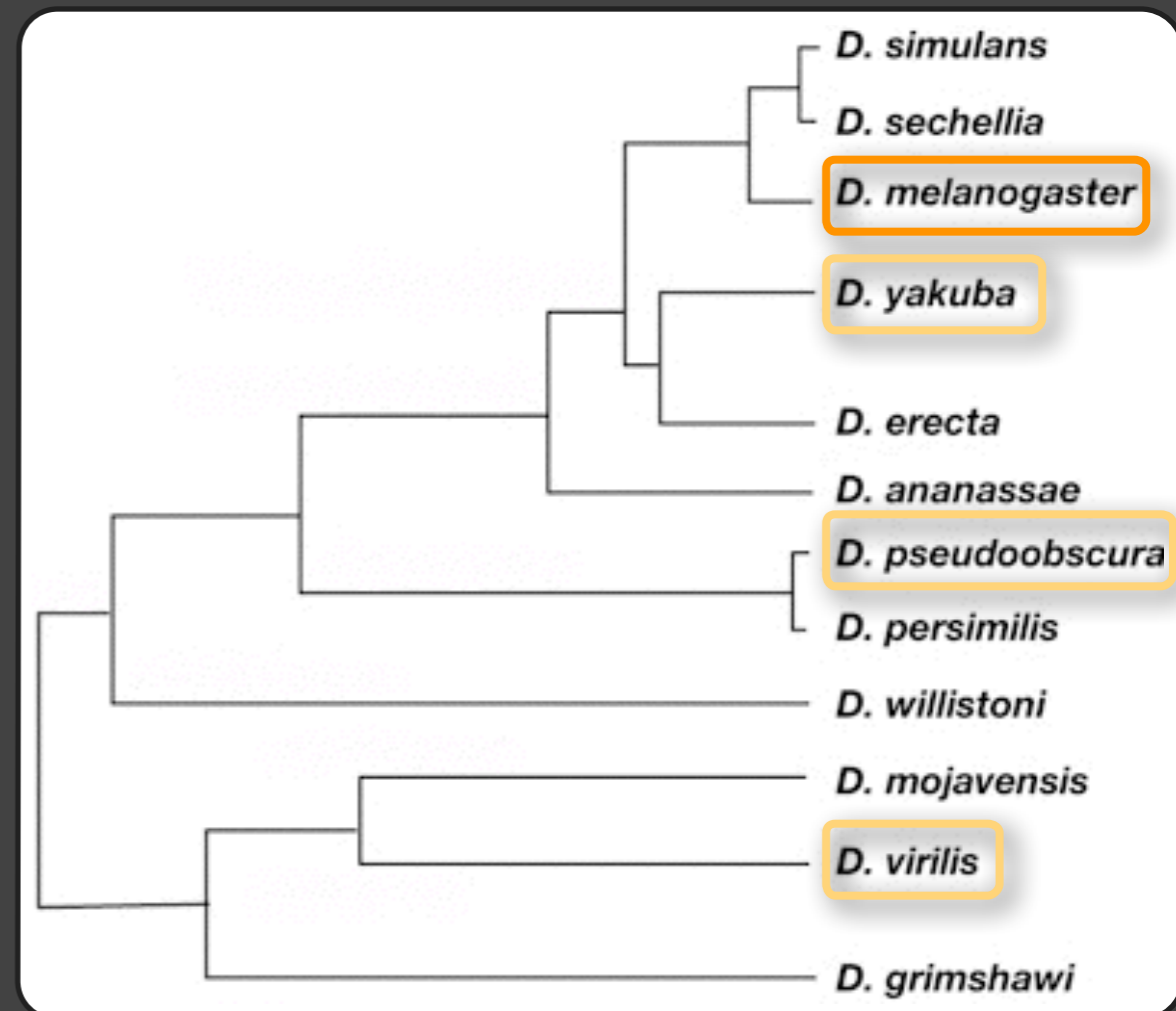
**spatial position**

[3D and 2D coordinates]

**12 related species**

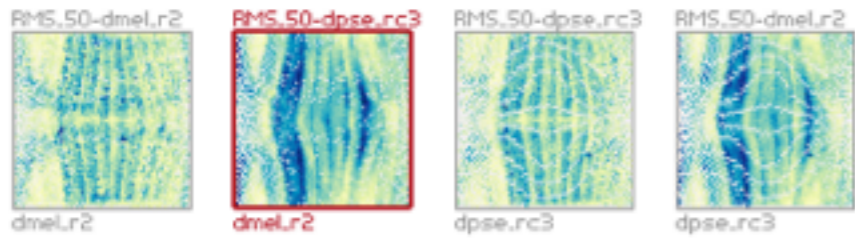
one complete

three preliminary

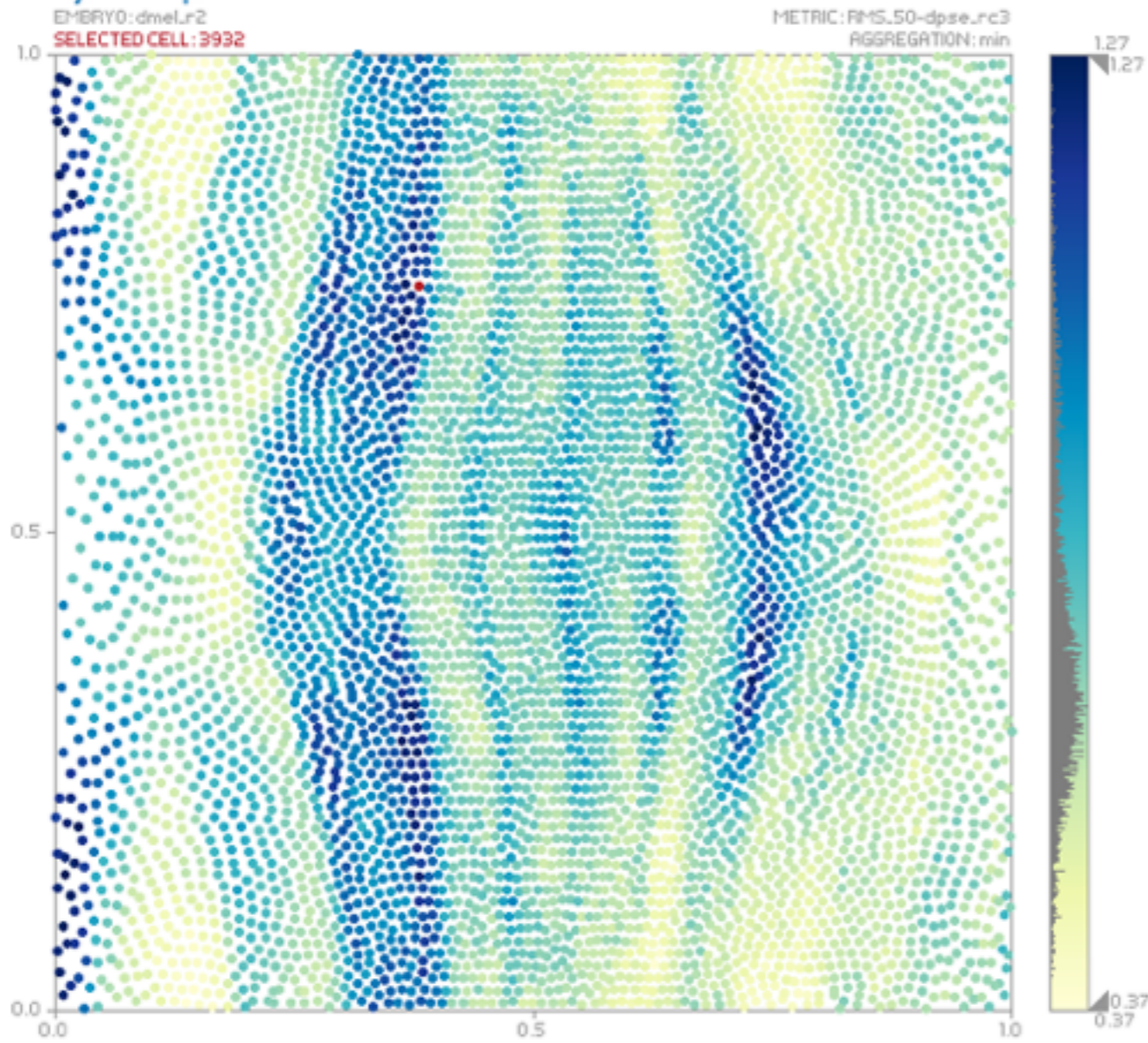




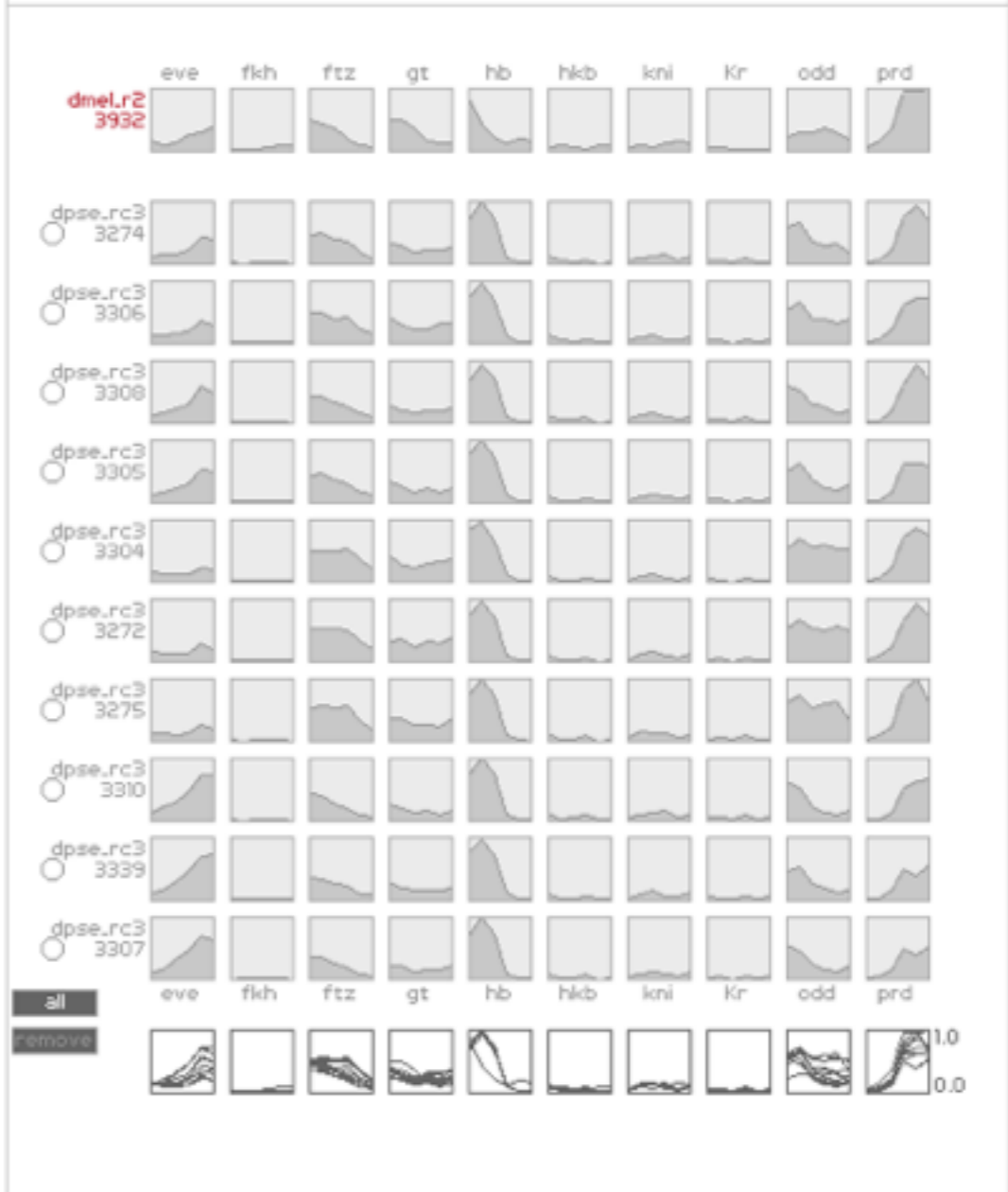
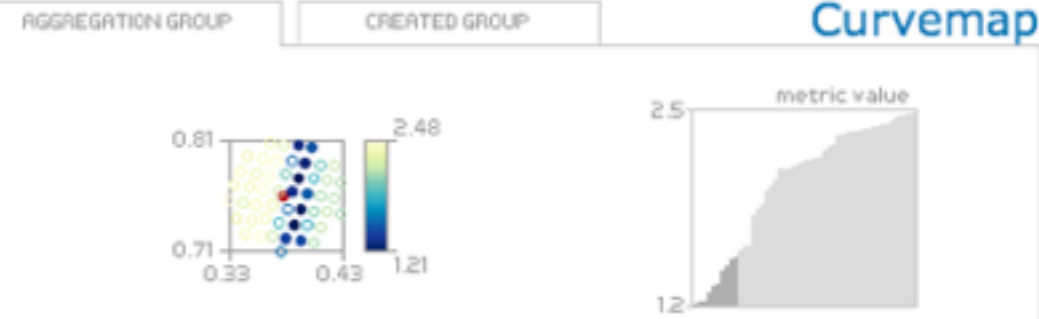
### Summaries



### Embryo Map

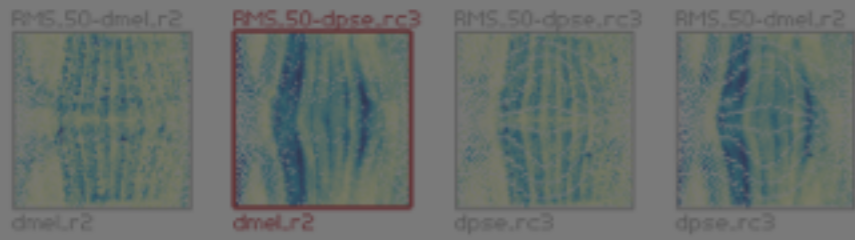


### Curvemap

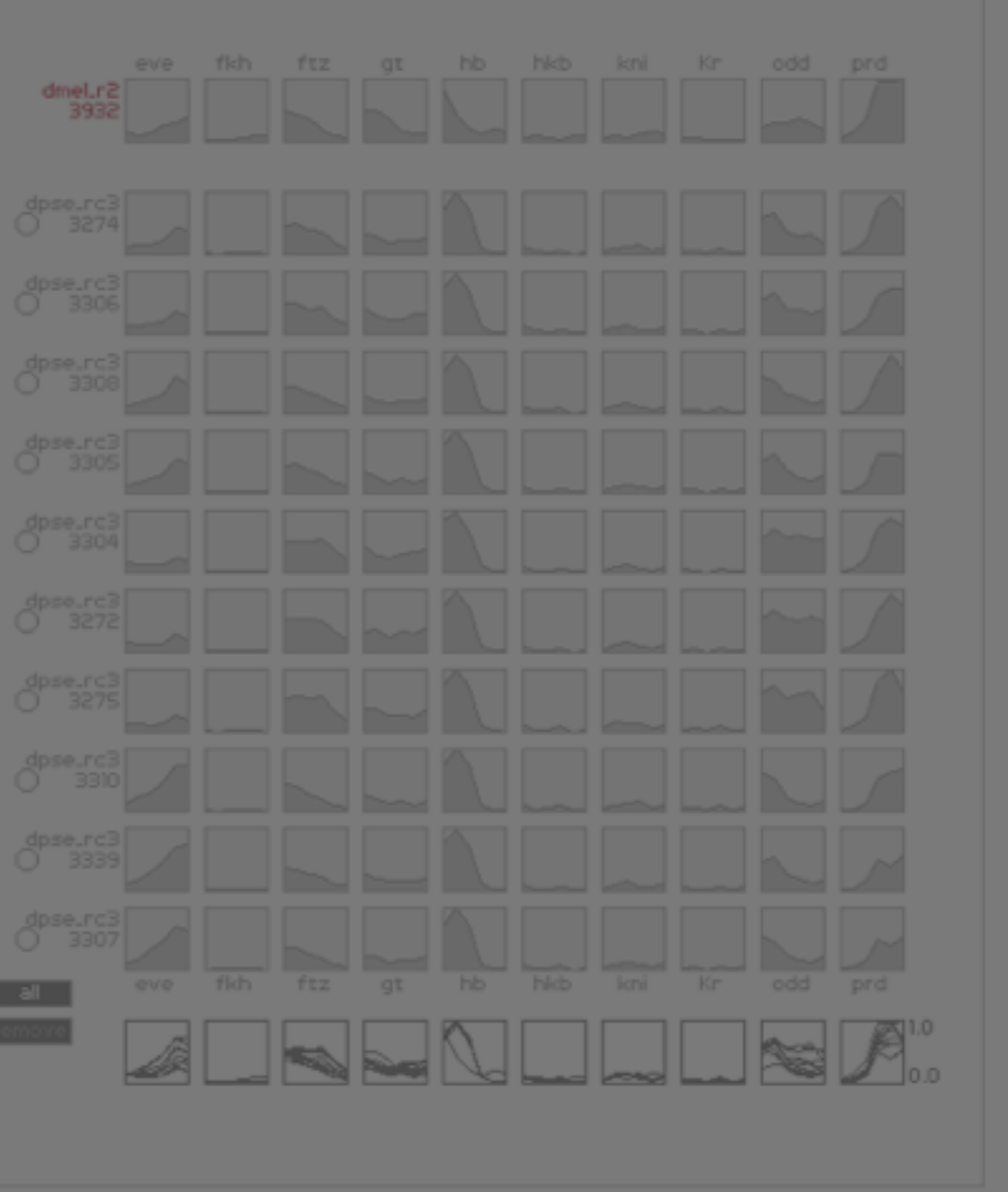
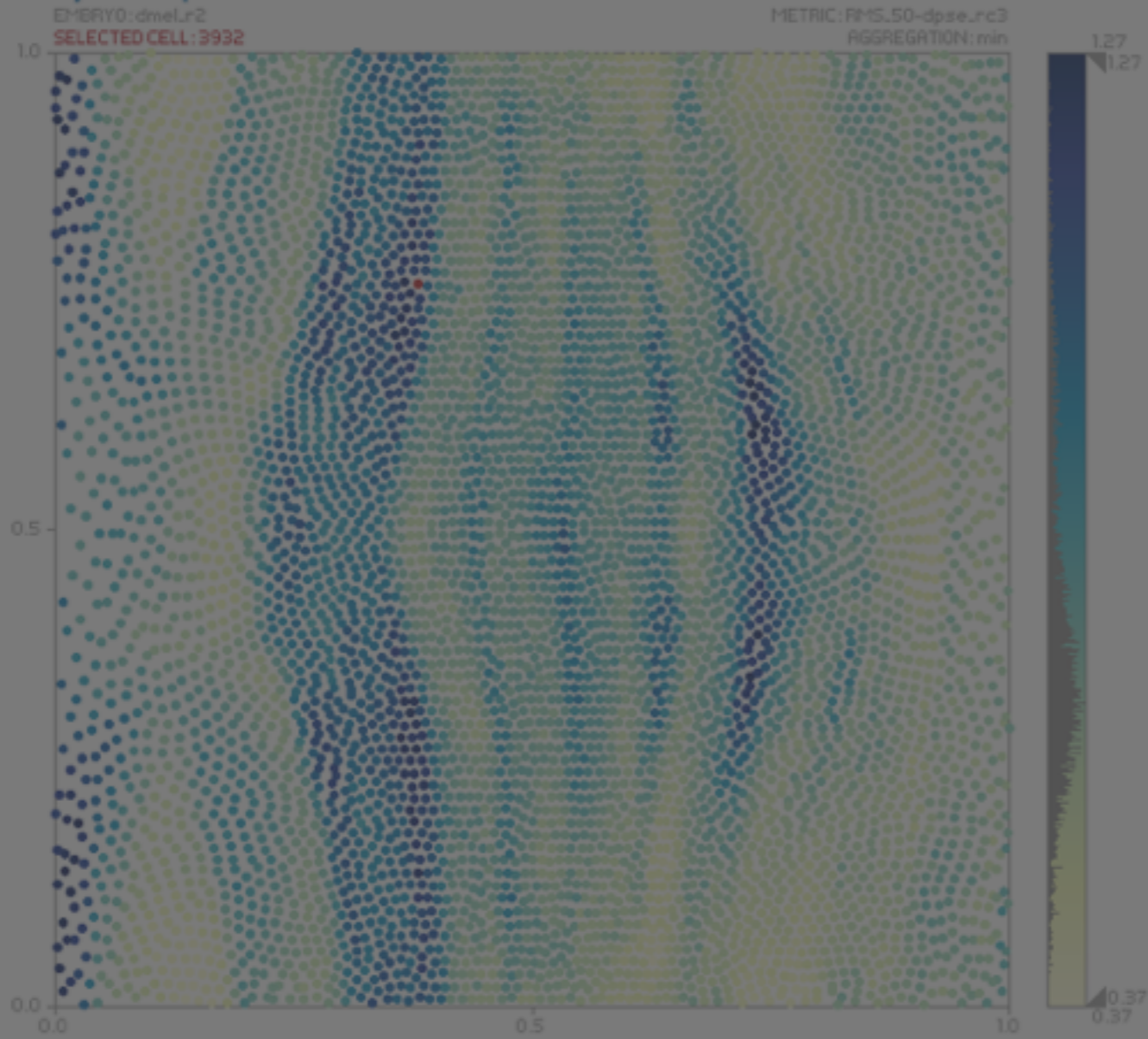




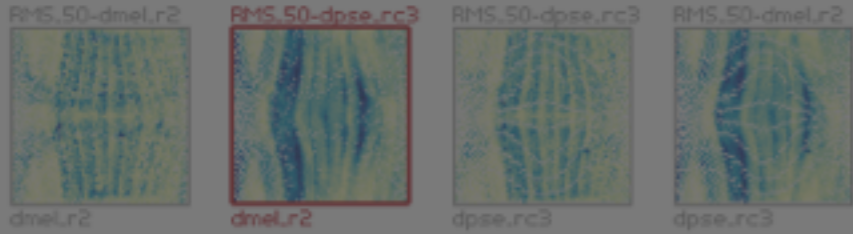
### Summaries



### Embryo Map

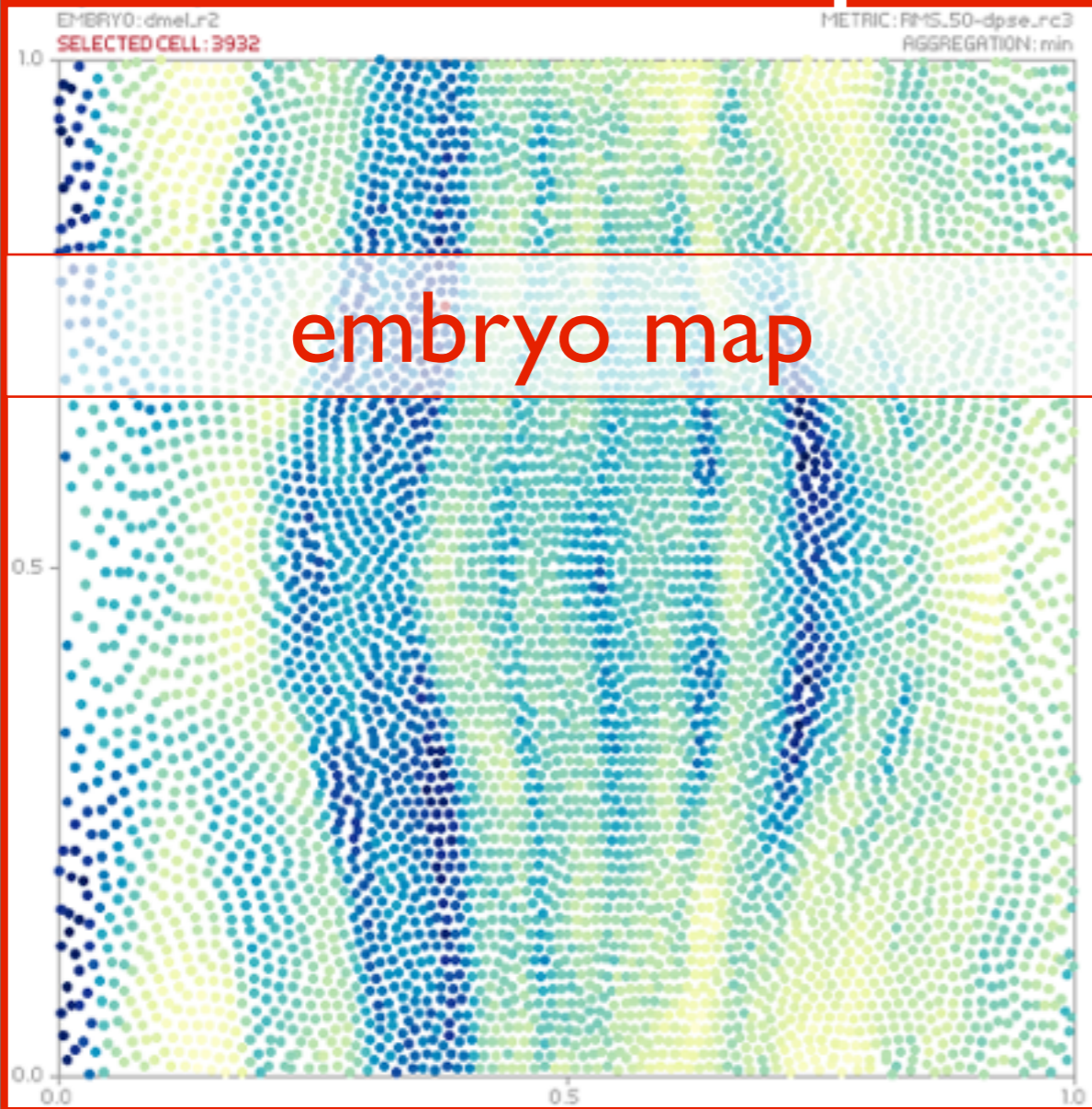


### Summaries



spatial

### Embryo Map



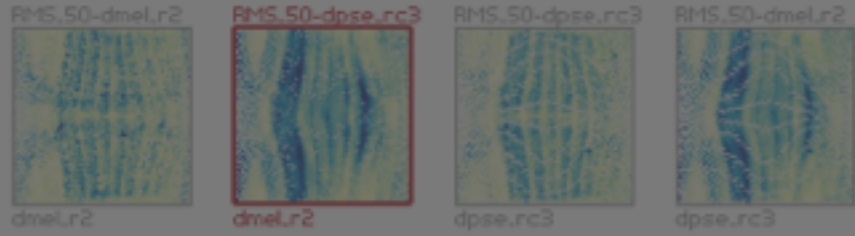
embryo map

### Curvemap

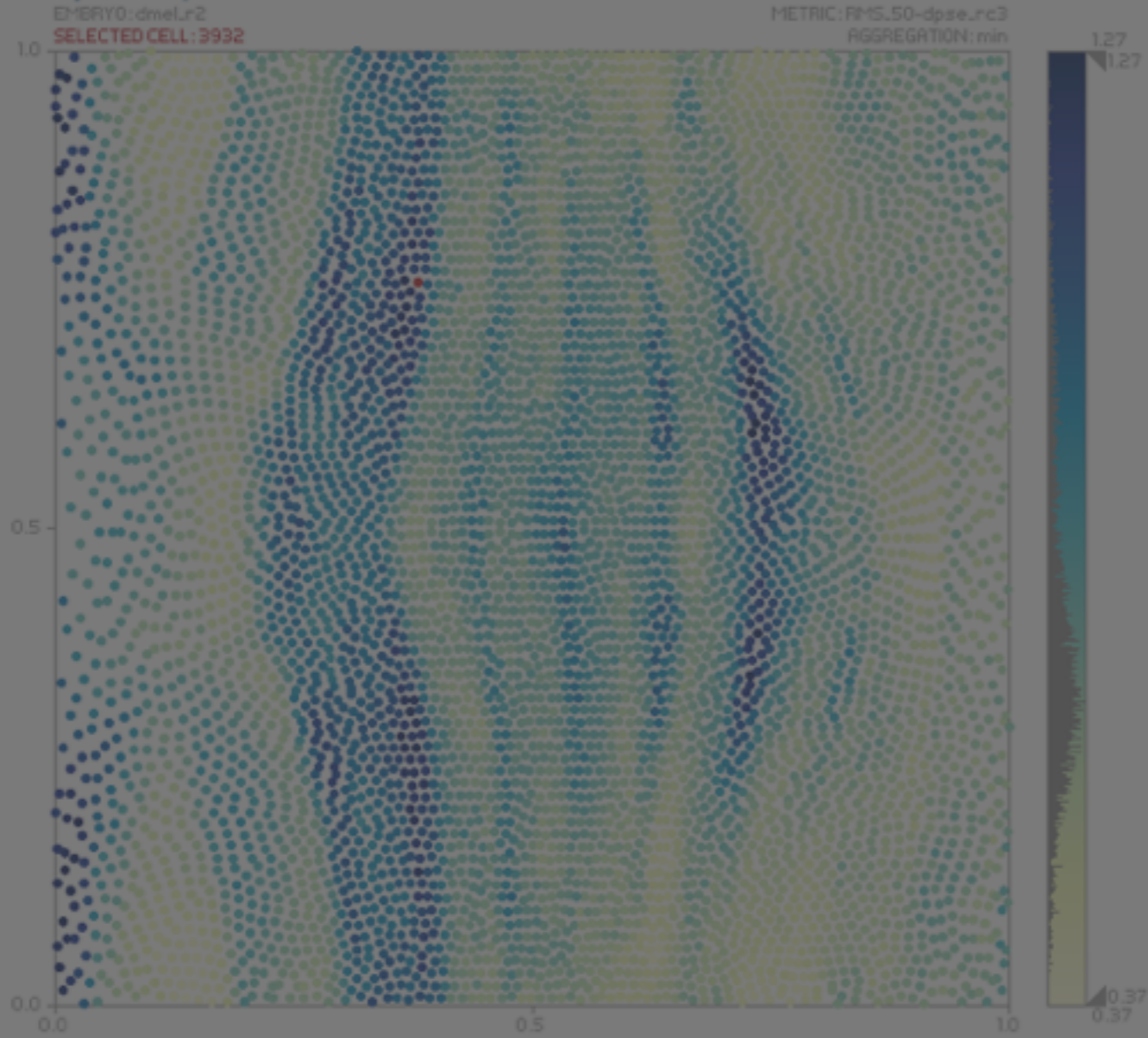




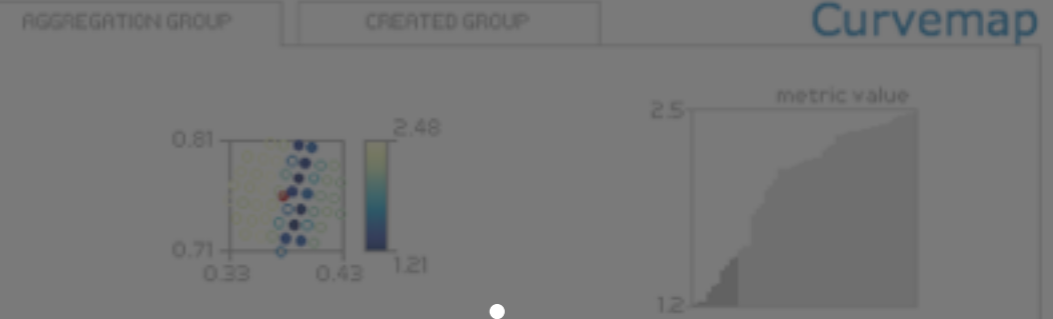
### Summaries



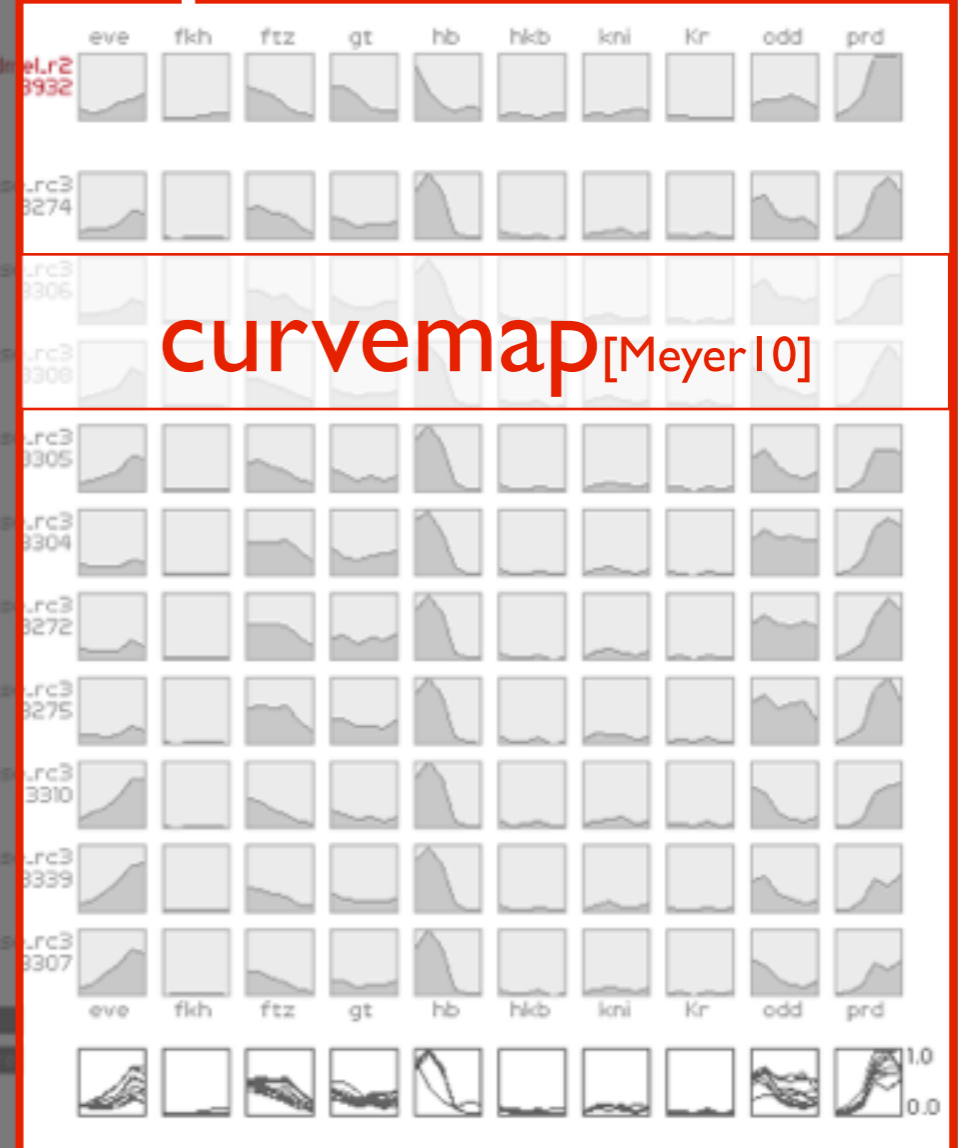
### Embryo Map



### Curvemap

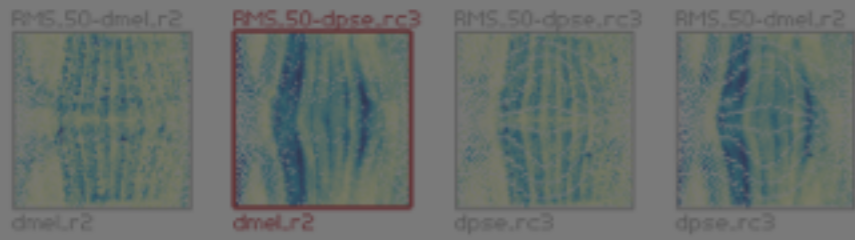


expression

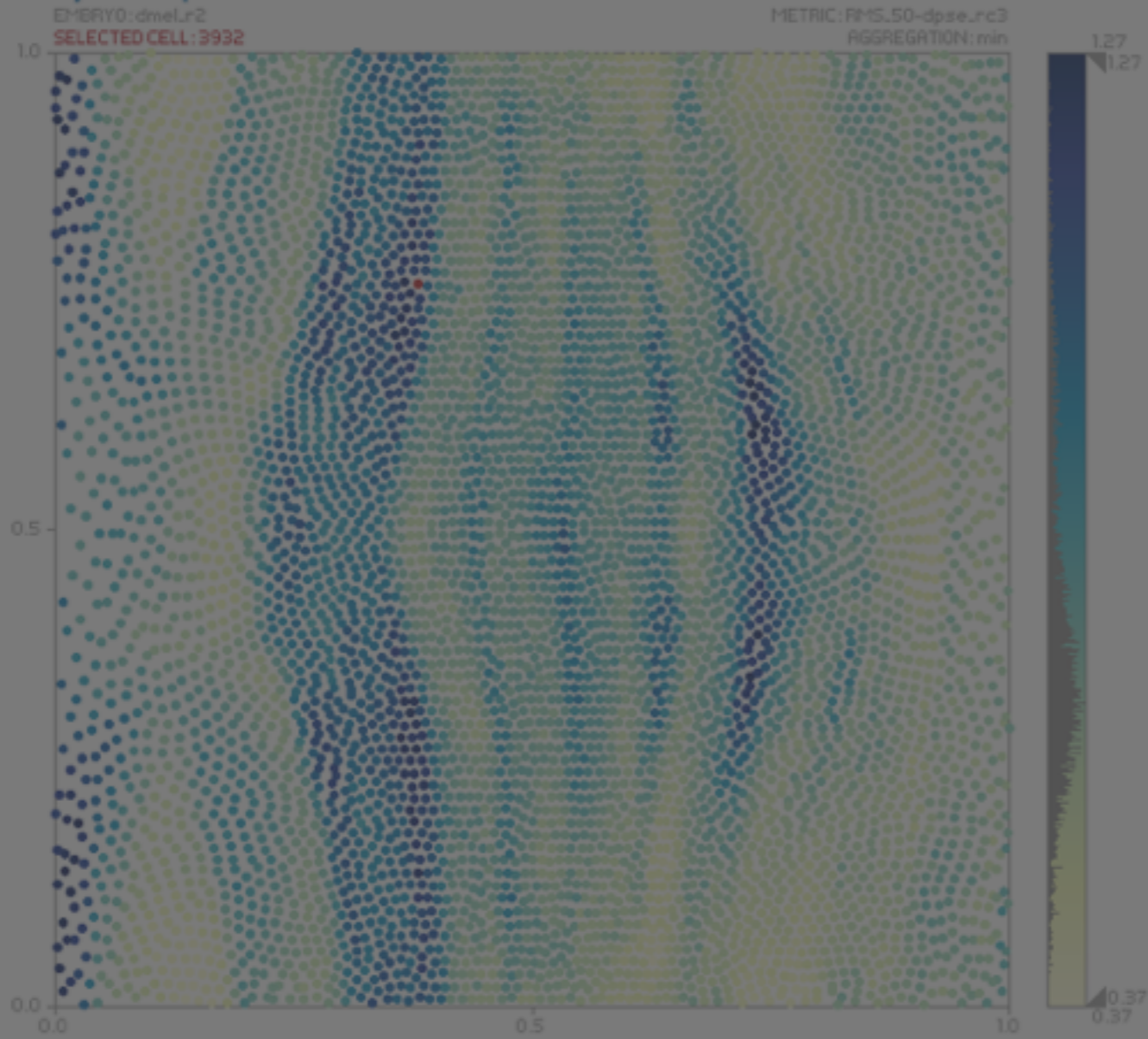




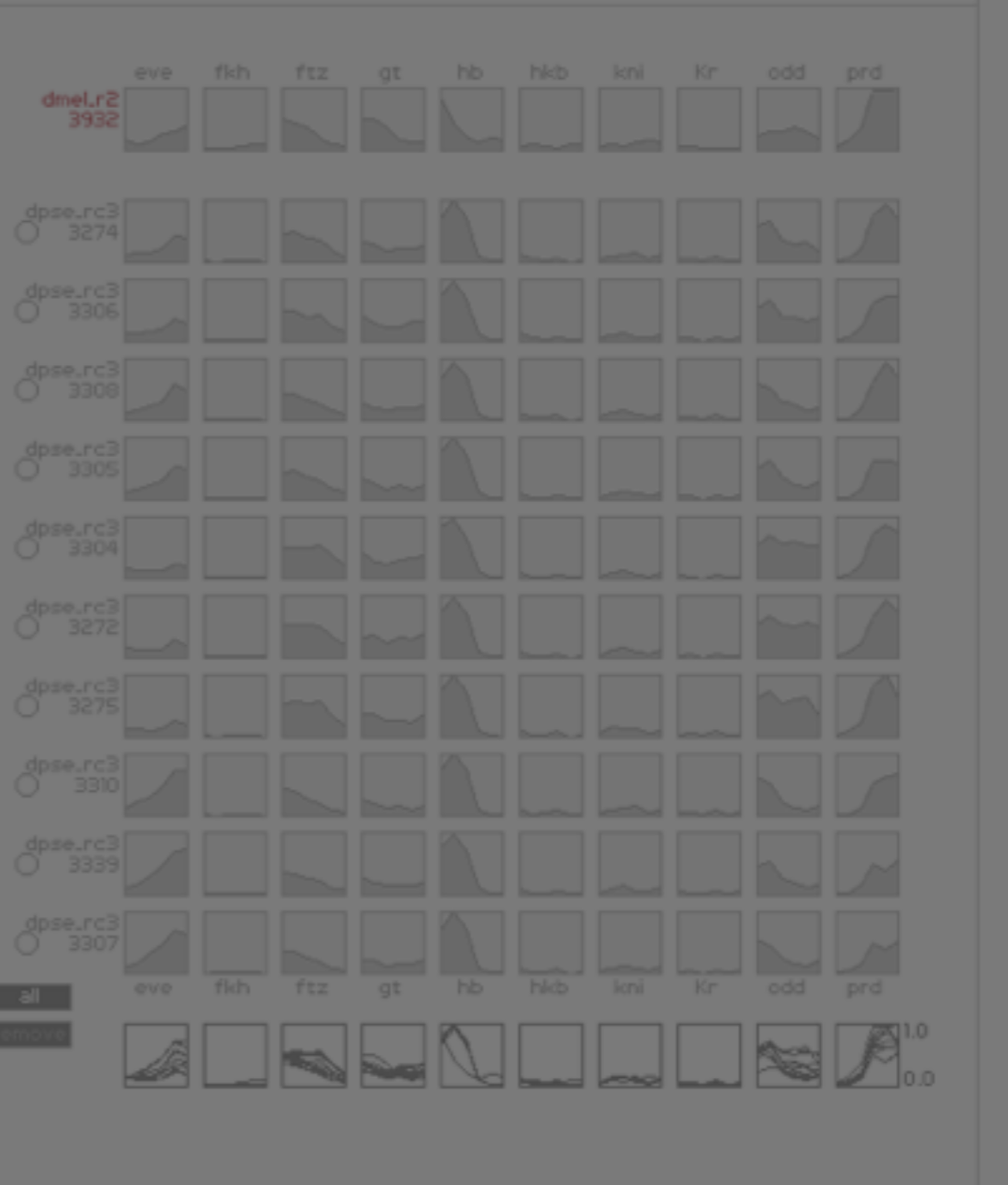
### Summaries



### Embryo Map



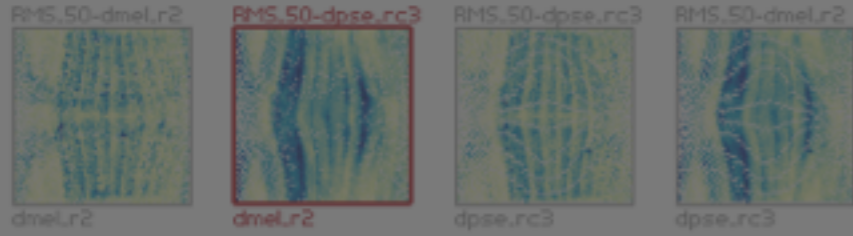
### Curvemap



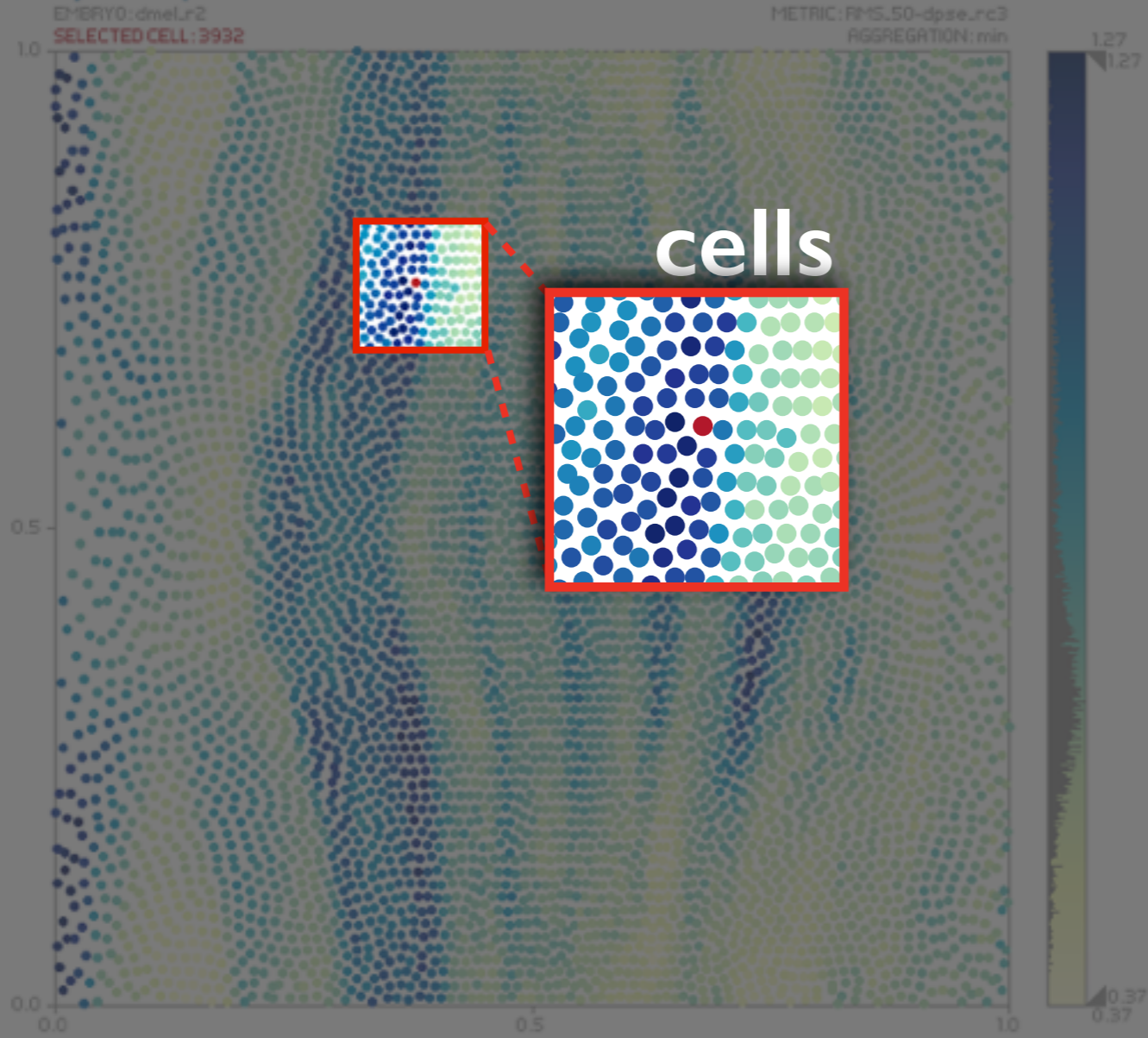




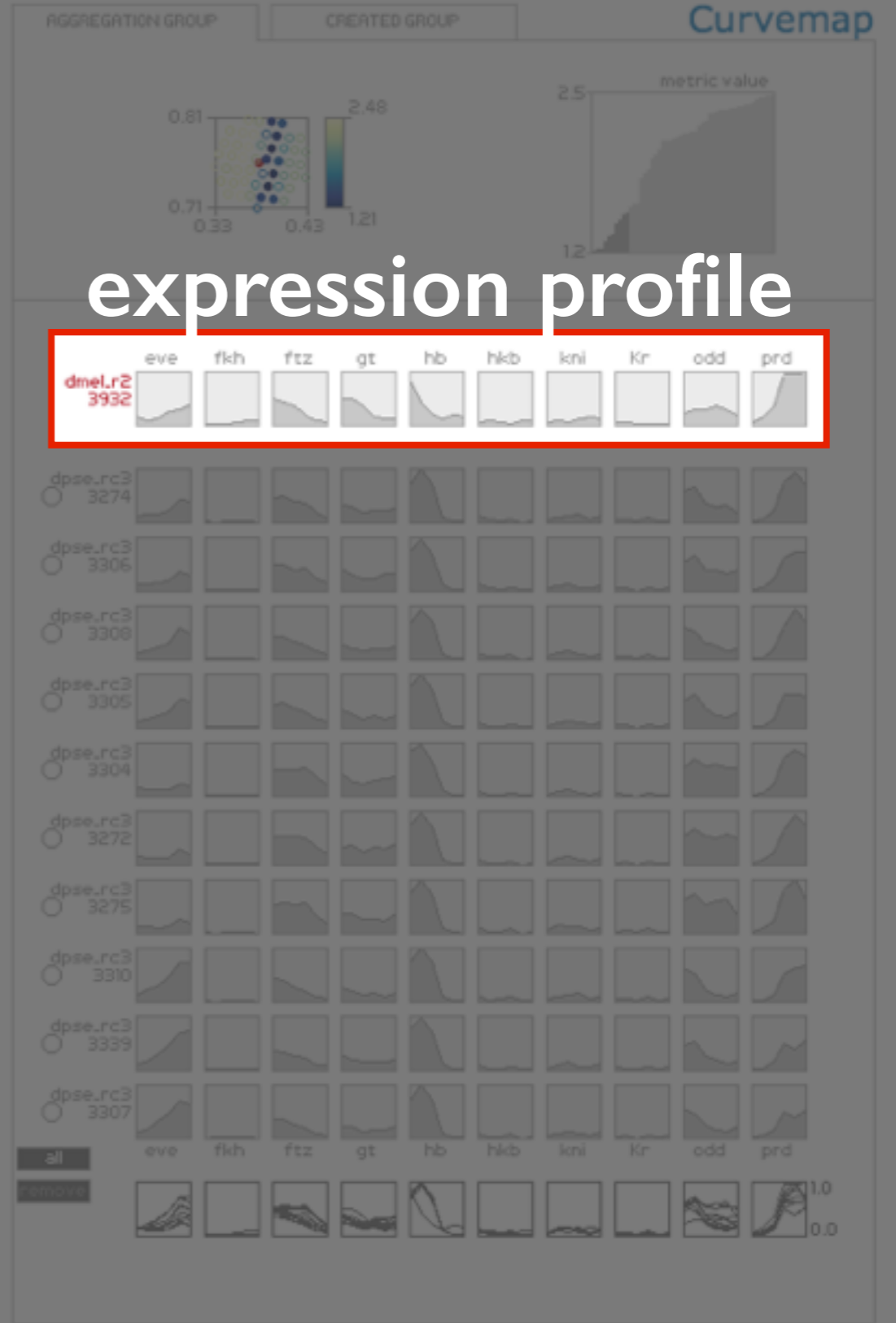
### Summaries



### Embryo Map

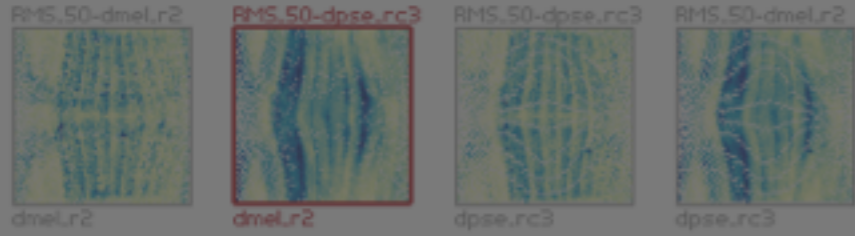


cells

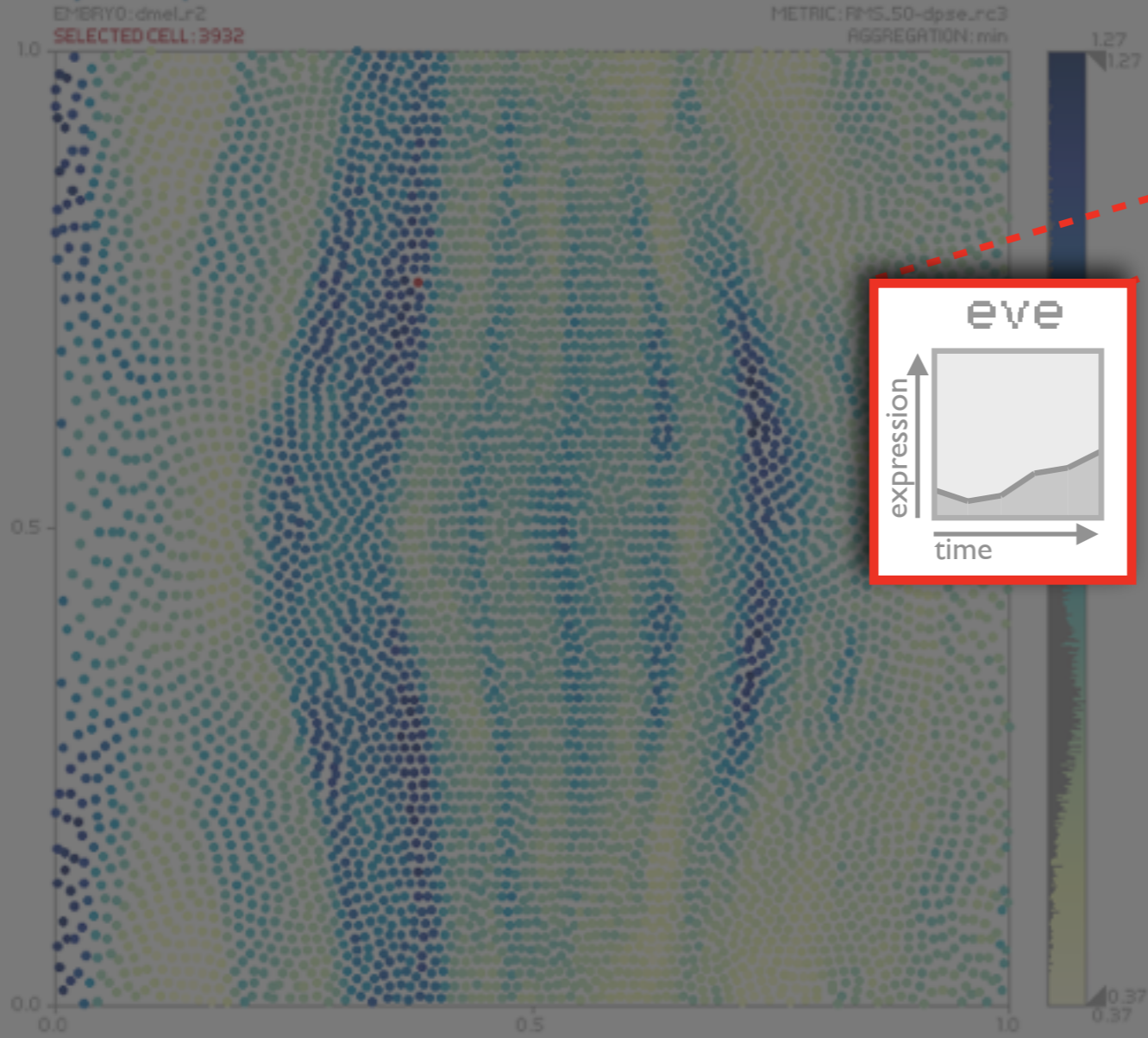


expression profile

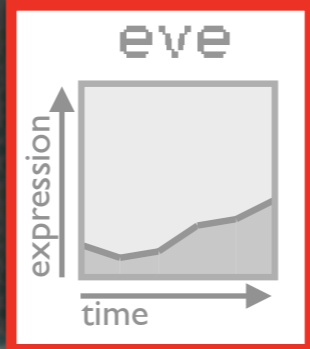
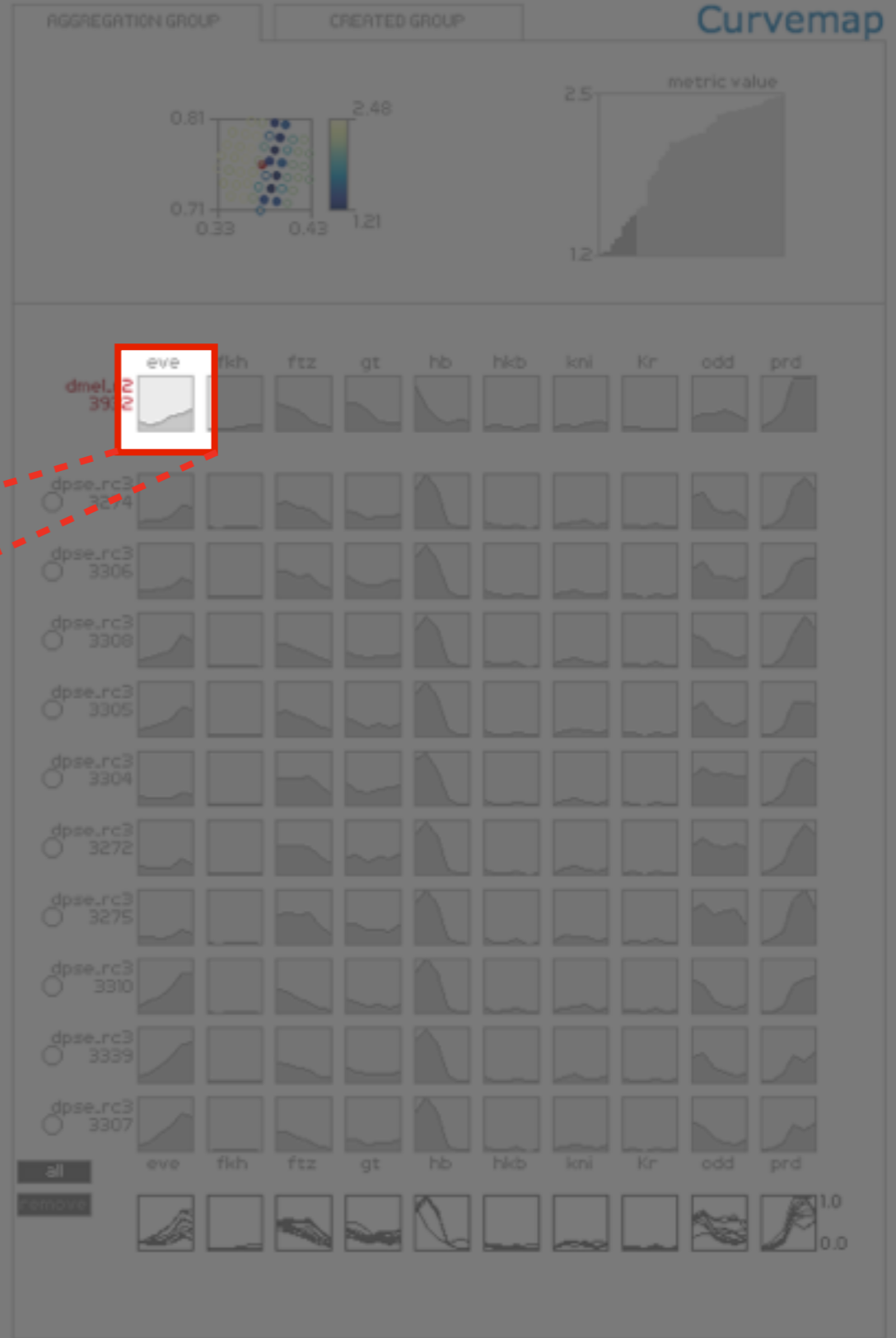
### Summaries



### Embryo Map

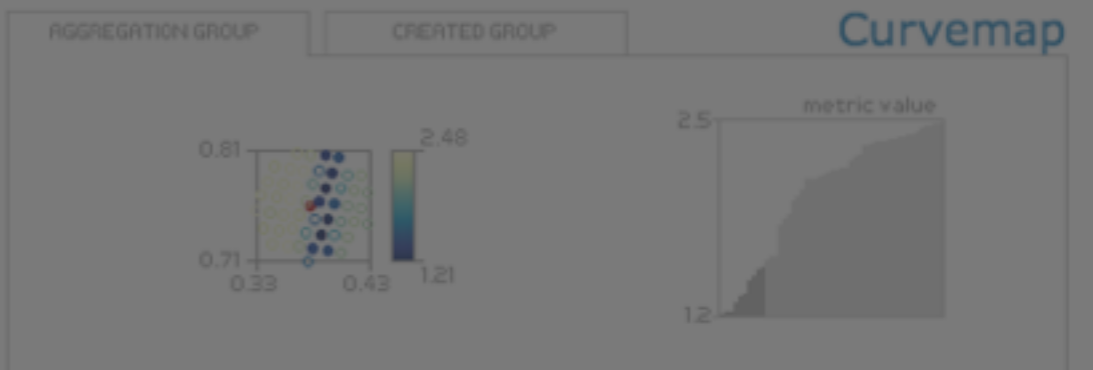
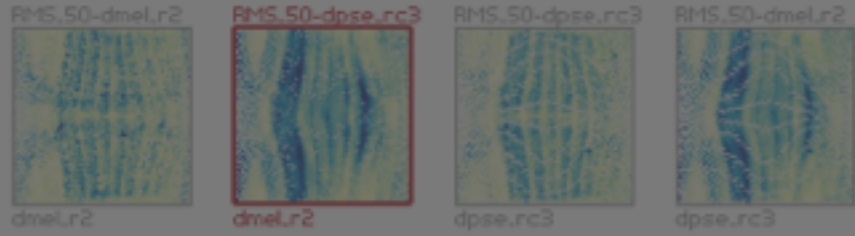


### Curvemap

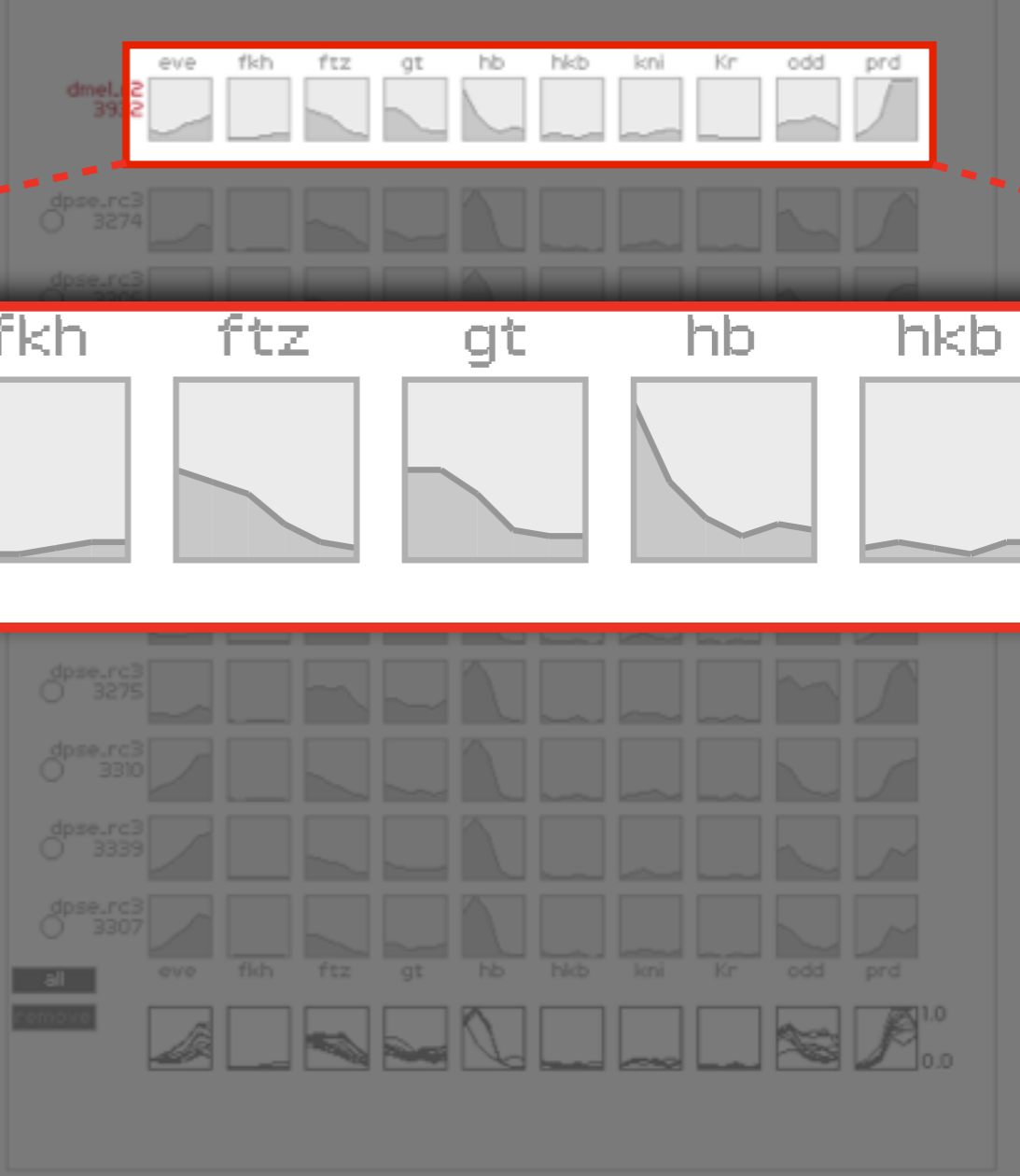
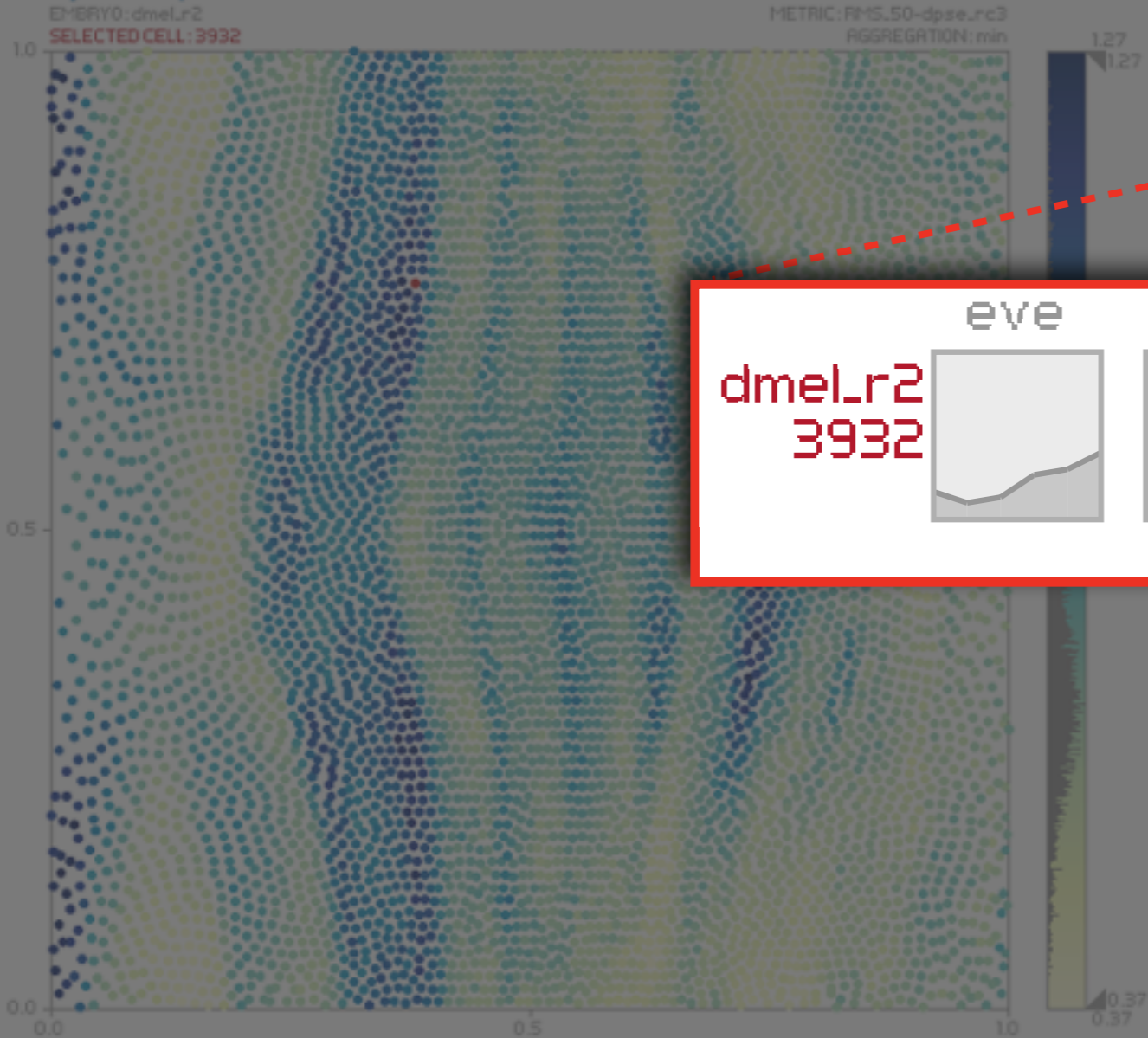




### Summaries



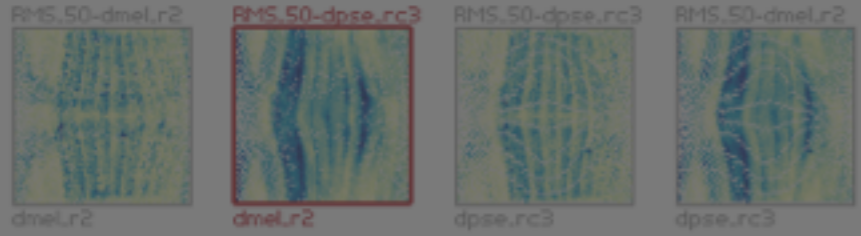
### Embryo Map



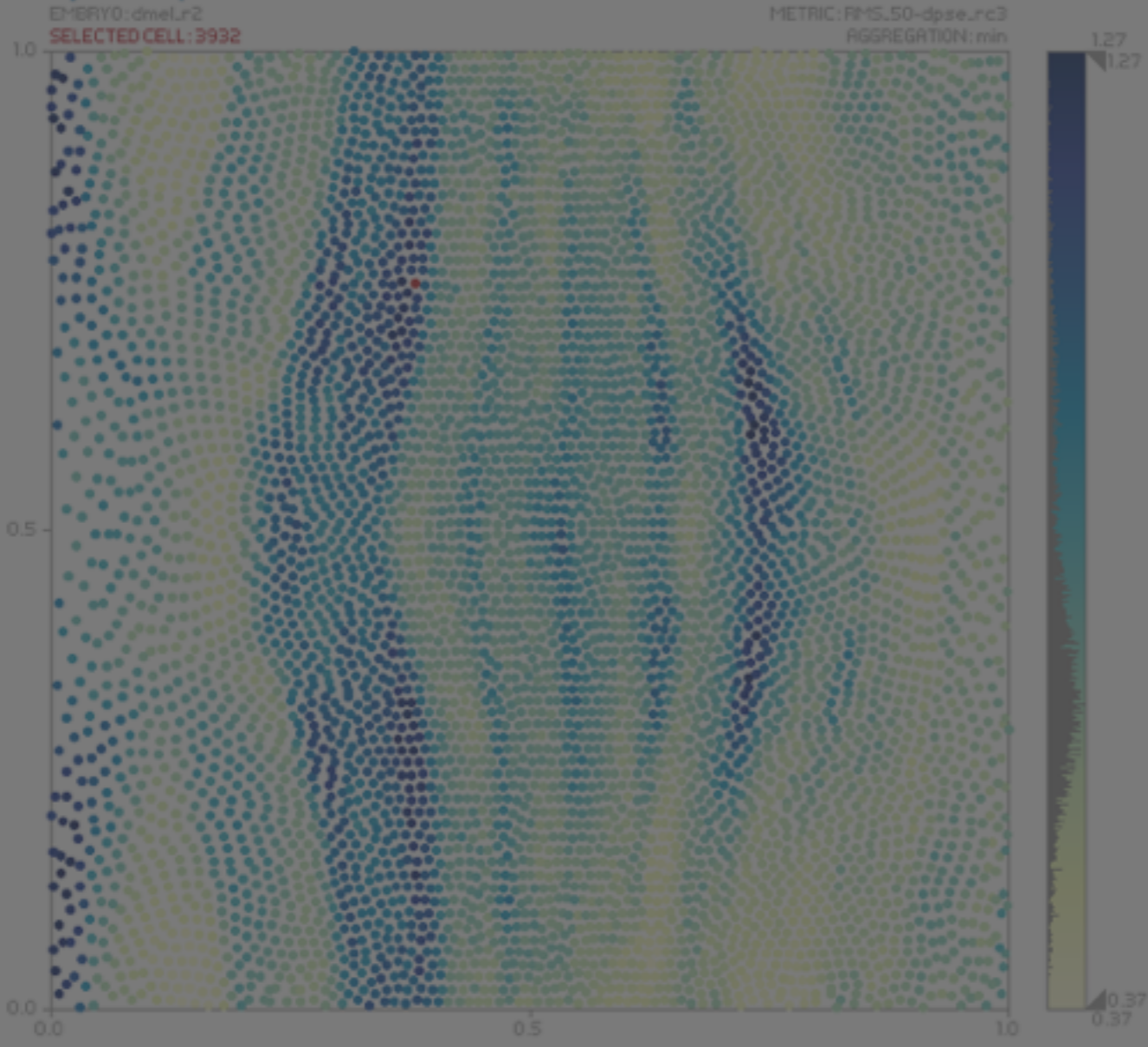




### Summaries

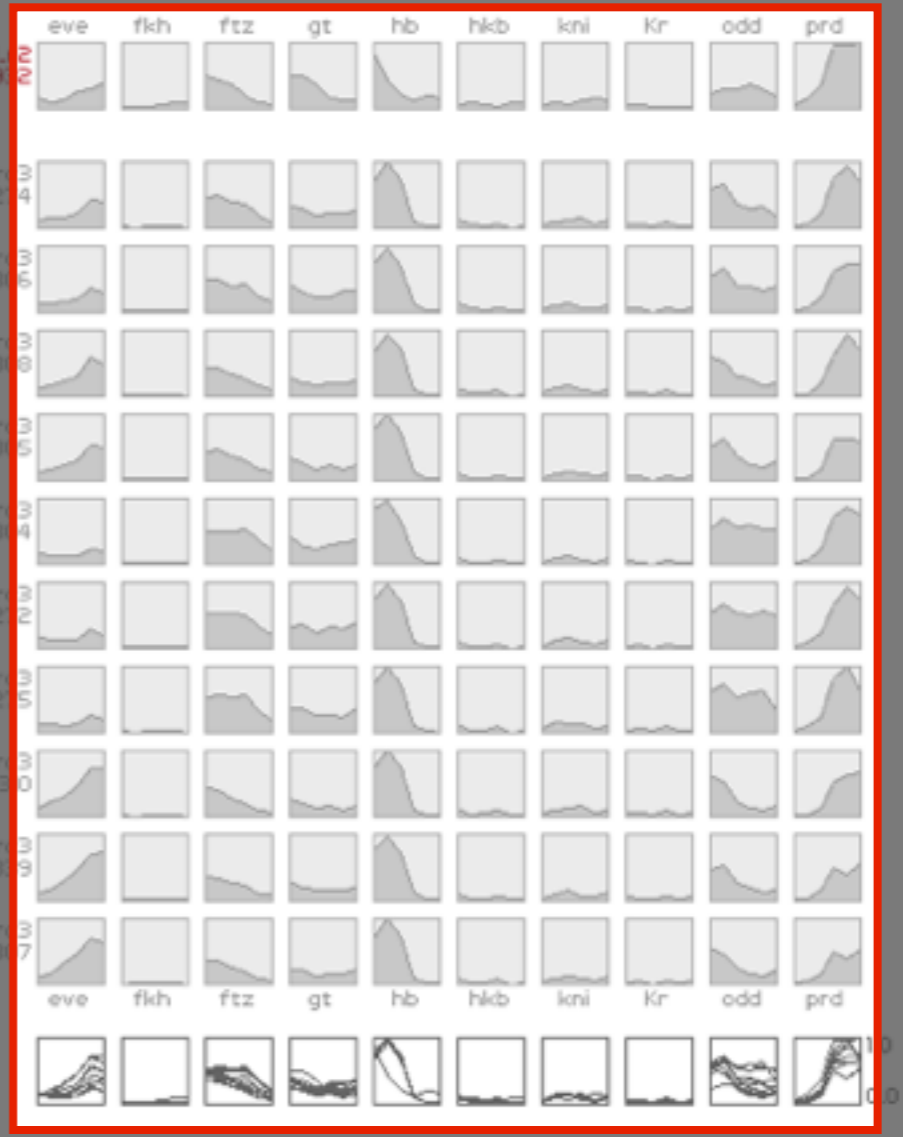


### Embryo Map

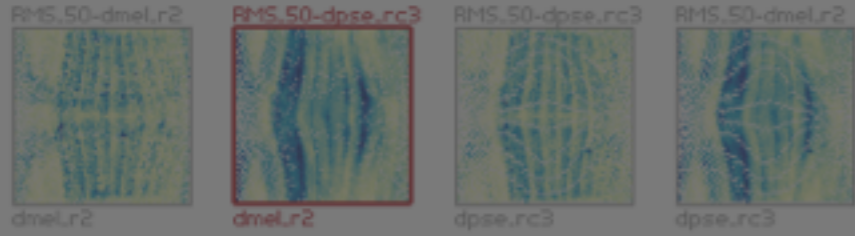


genes

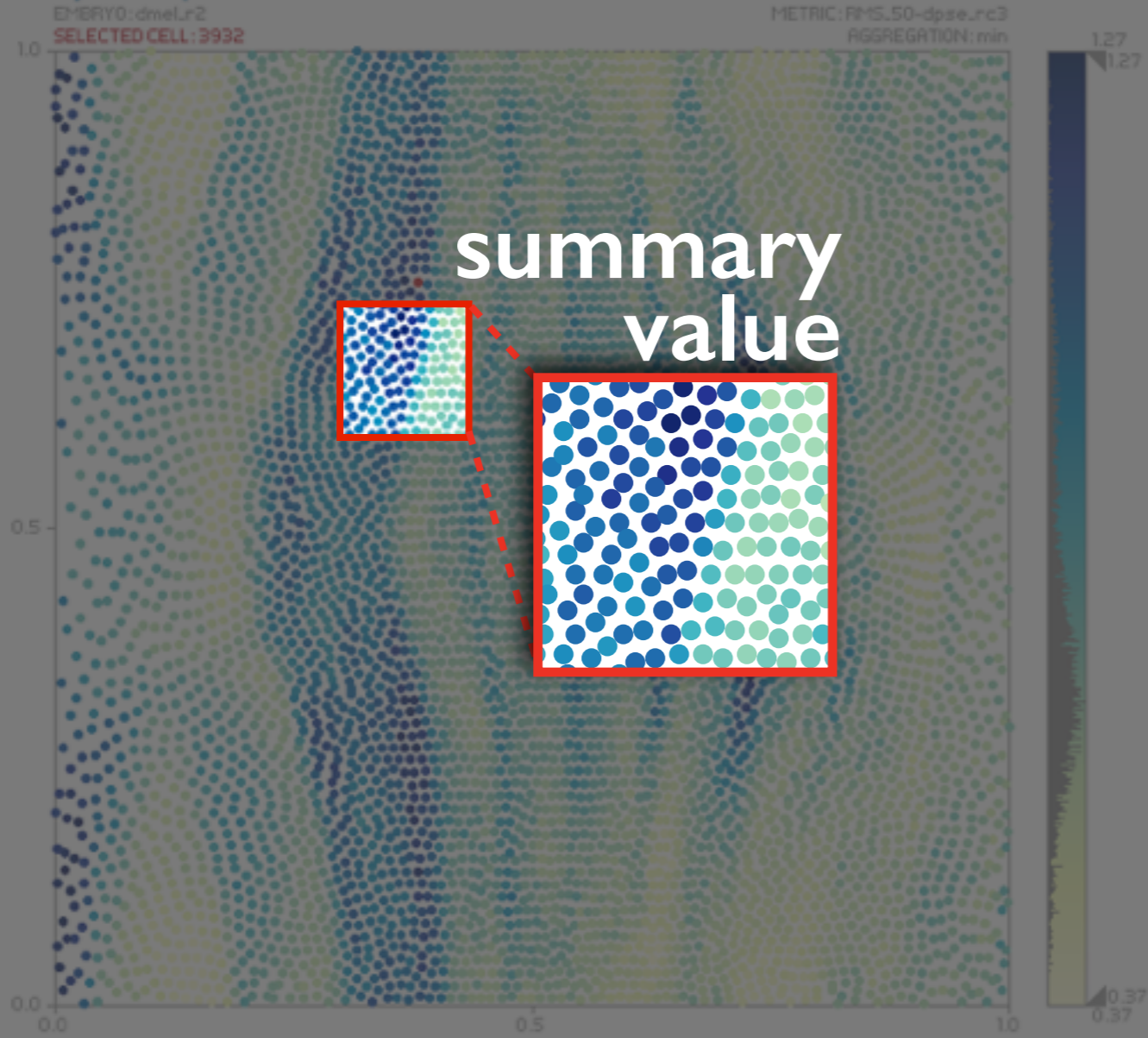
cells



### Summaries



### Embryo Map



### Curvemap

AGGREGATION GROUP | CREATED GROUP

Curvemap plots showing metric value (0.71 to 2.48) and gene expression profiles for various genes (eve, fkh, ftz, gt, hb, hkb, kni, Kr, odd, prd) across different cell lines (dmeLr2, dpse\_rc3).

Cell Line	eve	fkh	ftz	gt	hb	hkb	kni	Kr	odd	prd
dmeLr2 3932	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]
dpse_rc3 3274	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]
dpse_rc3 3306	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]
dpse_rc3 3309	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]
dpse_rc3 3305	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]
dpse_rc3 3304	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]
dpse_rc3 3272	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]
dpse_rc3 3275	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]
dpse_rc3 3310	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]
dpse_rc3 3339	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]
dpse_rc3 3307	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]	[Profile]

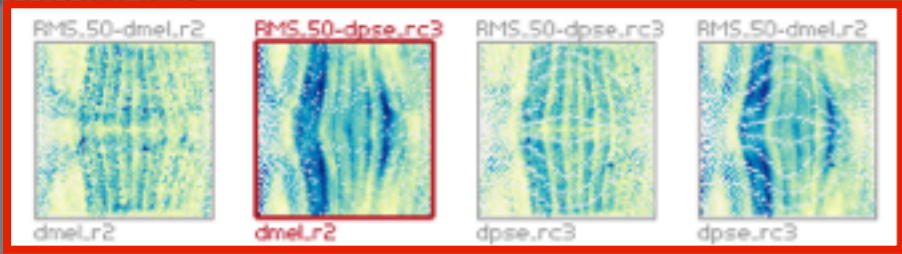
all | remove

0.0 | 1.0

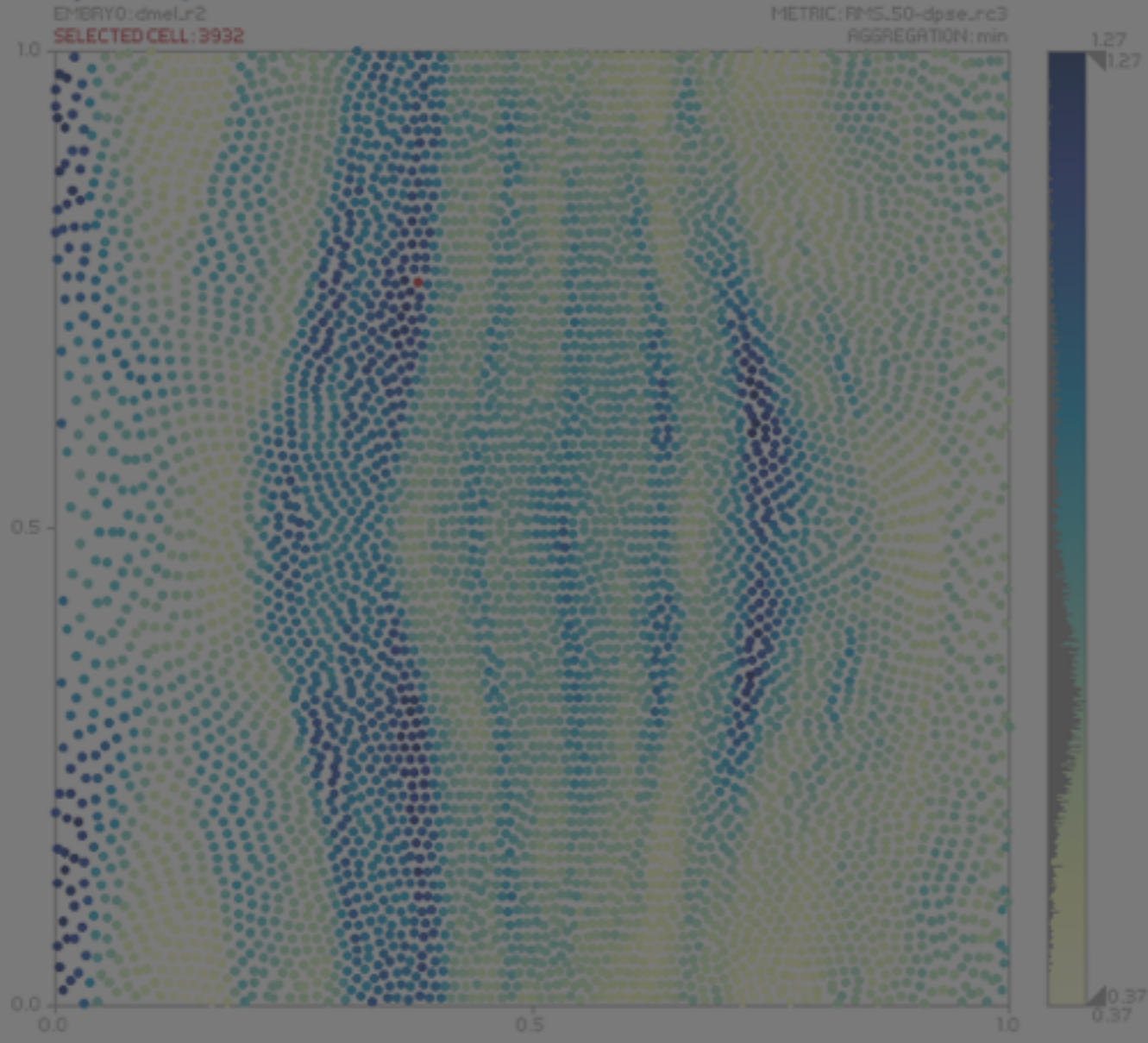


# summaries

## Summaries



## Embryo Map



**Characterize differences in gene expression patterns between species.**

# Characterize differences in gene expression patterns between species.

differences related to:

*spatial position*

*gene expression profiles*

*complex combination*

challenging to characterize manually

support mechanisms:

*summaries, groups*

data & tool & tasks

**summaries & groups**

encodings & interaction

conclusions

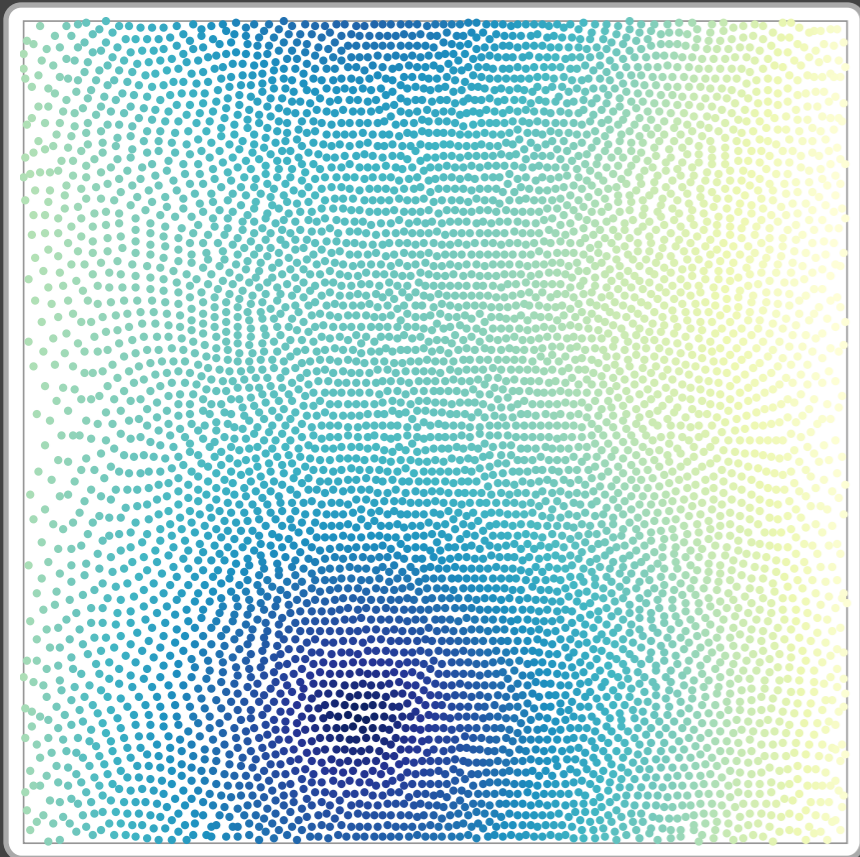
# summary

a value for each cell that summarizes the underlying data

# summary

a value for each cell that summarizes the underlying data

**spatial**



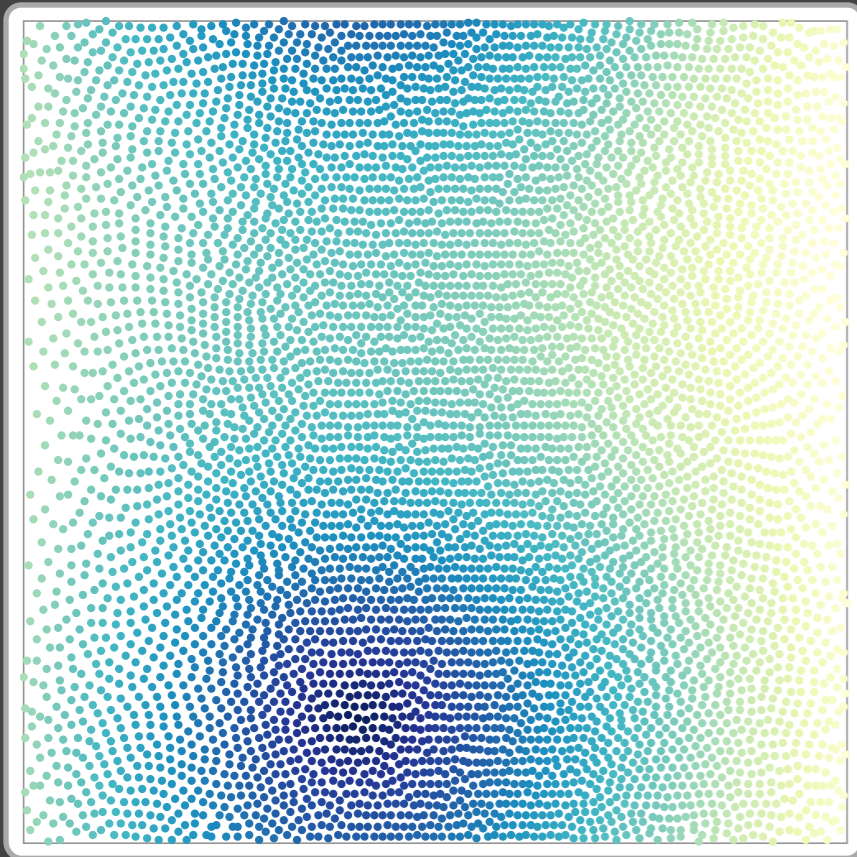
distance to cell  $i$



# summary

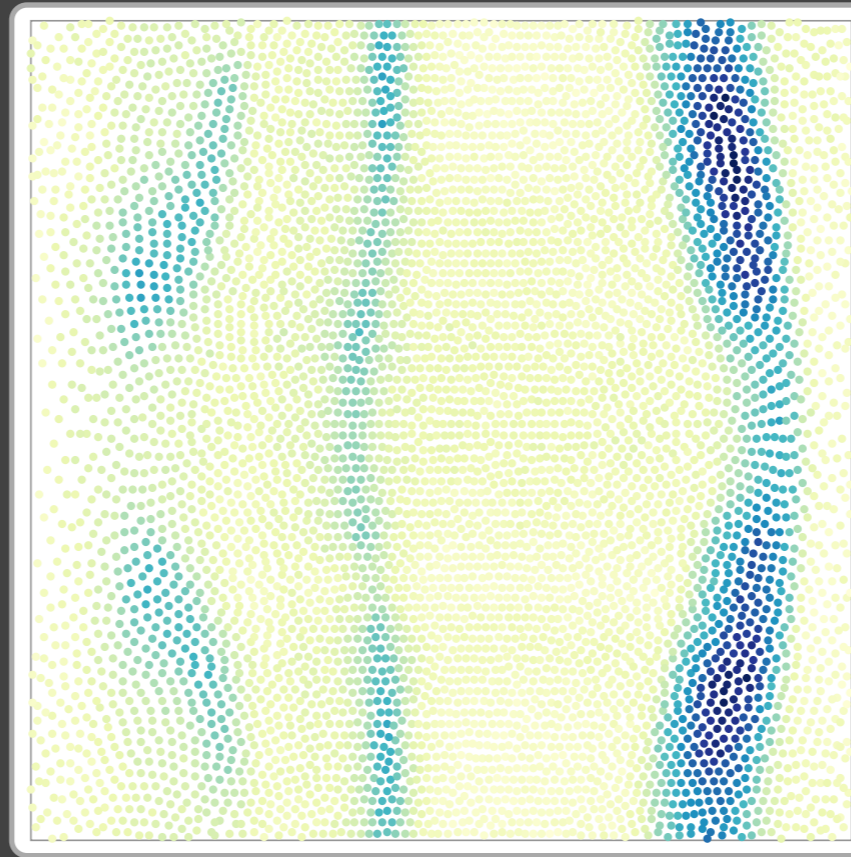
a value for each cell that summarizes the underlying data

**spatial**



distance to cell  $i$

**expression**



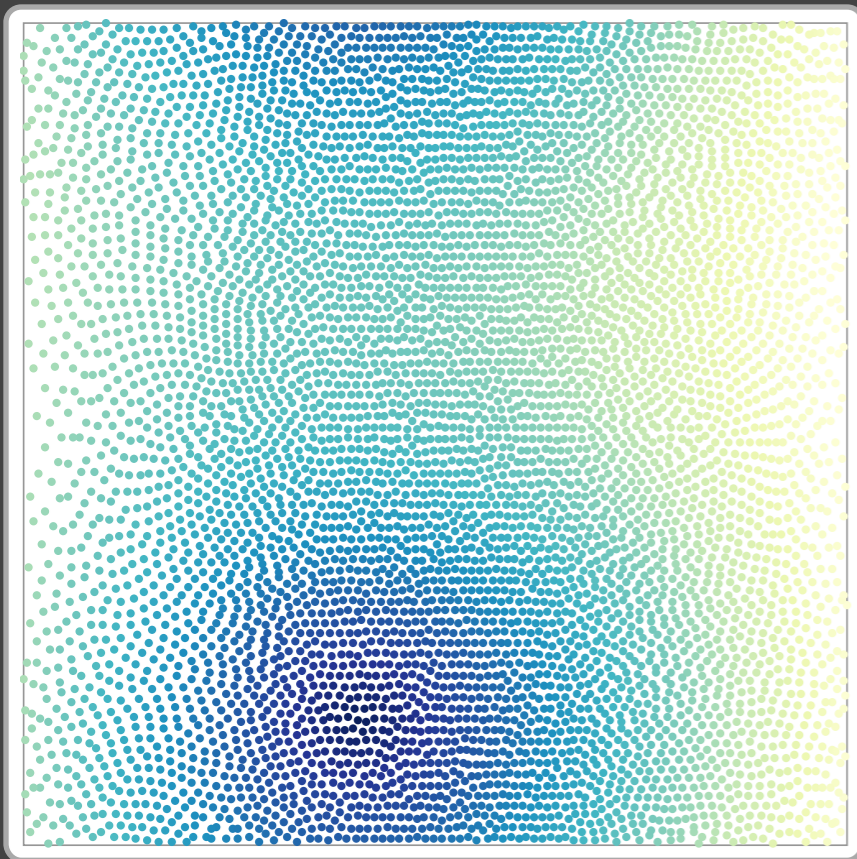
gene  $j$  at timepoint  $k$



# summary

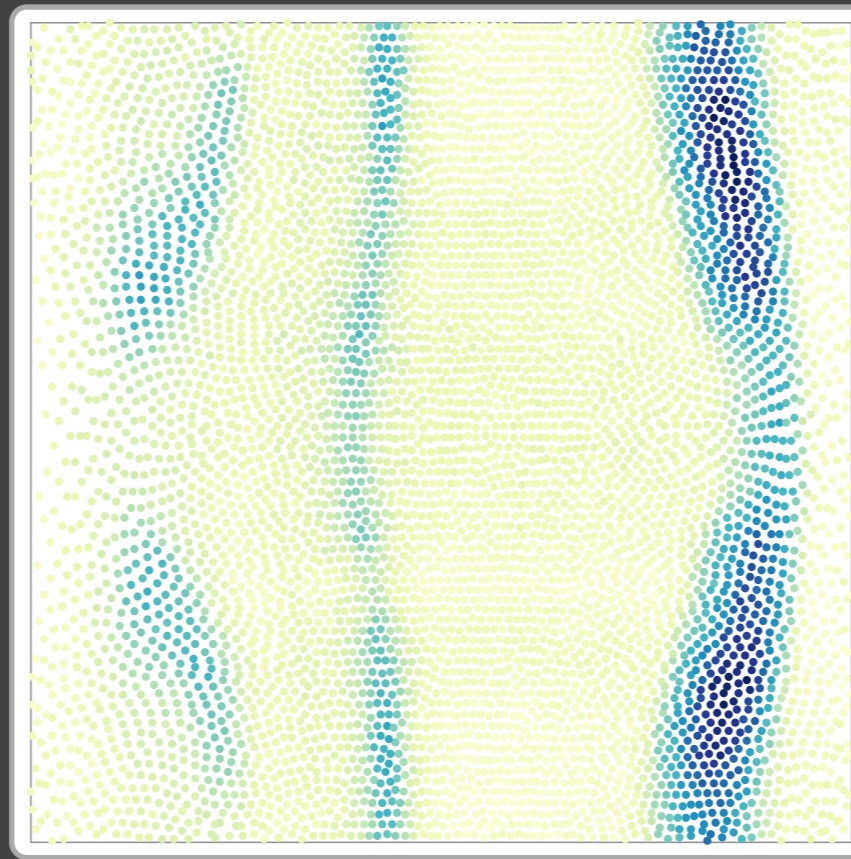
a value for each cell that summarizes the underlying data

**spatial**



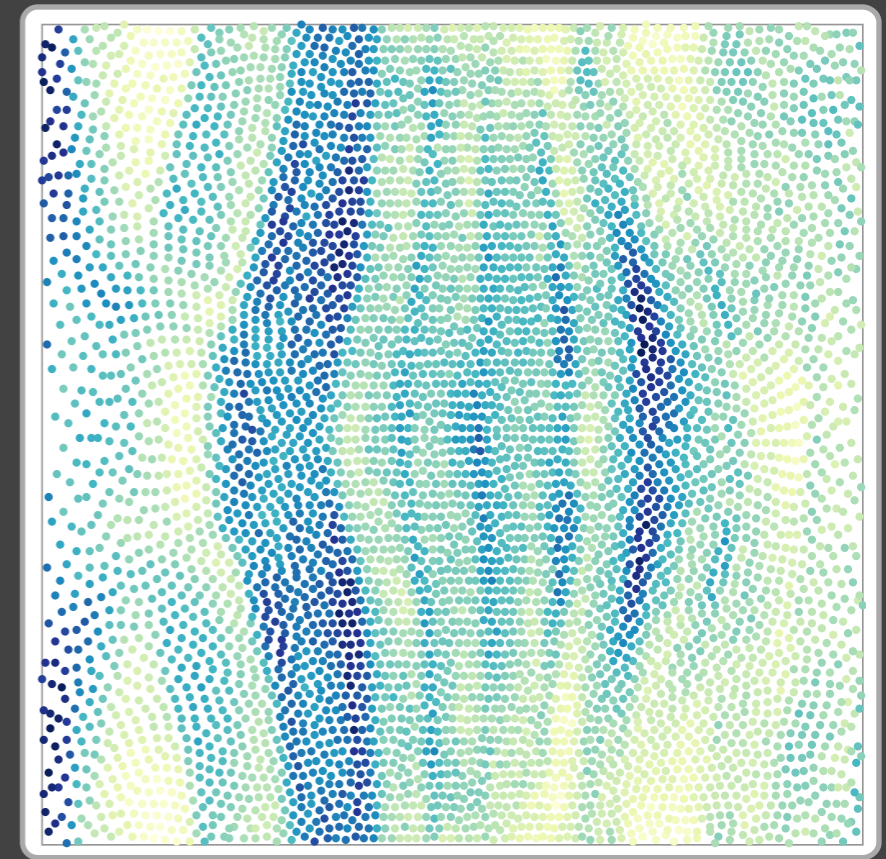
distance to cell  $i$

**expression**



gene  $j$  at timepoint  $k$

**combination**



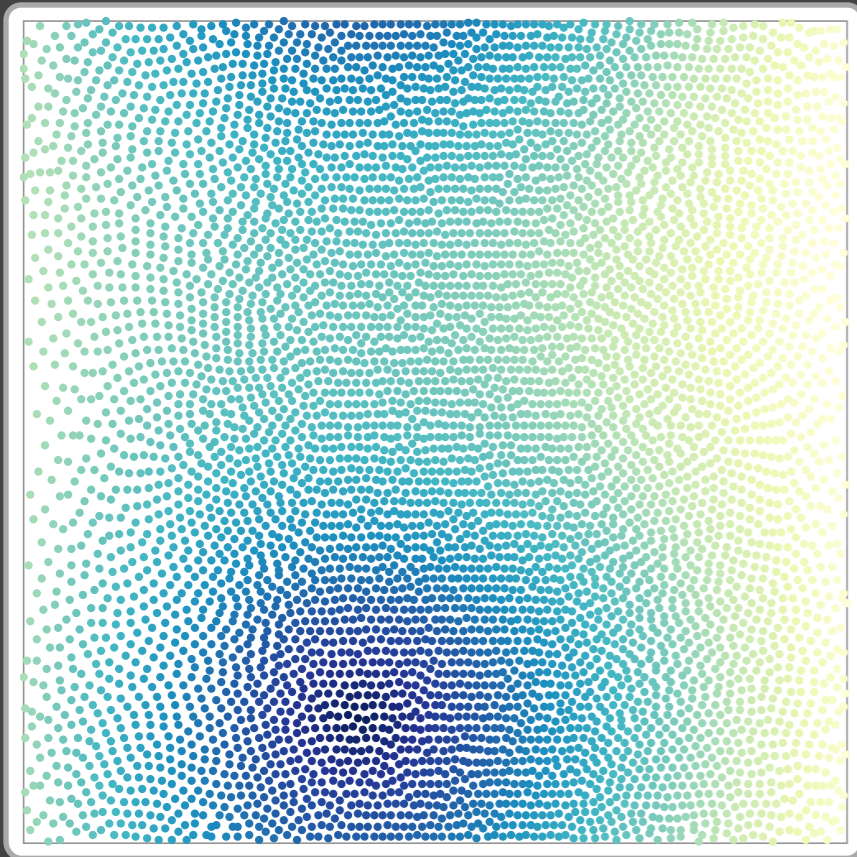
compare embryos  
 $A$  and  $B$



# summary

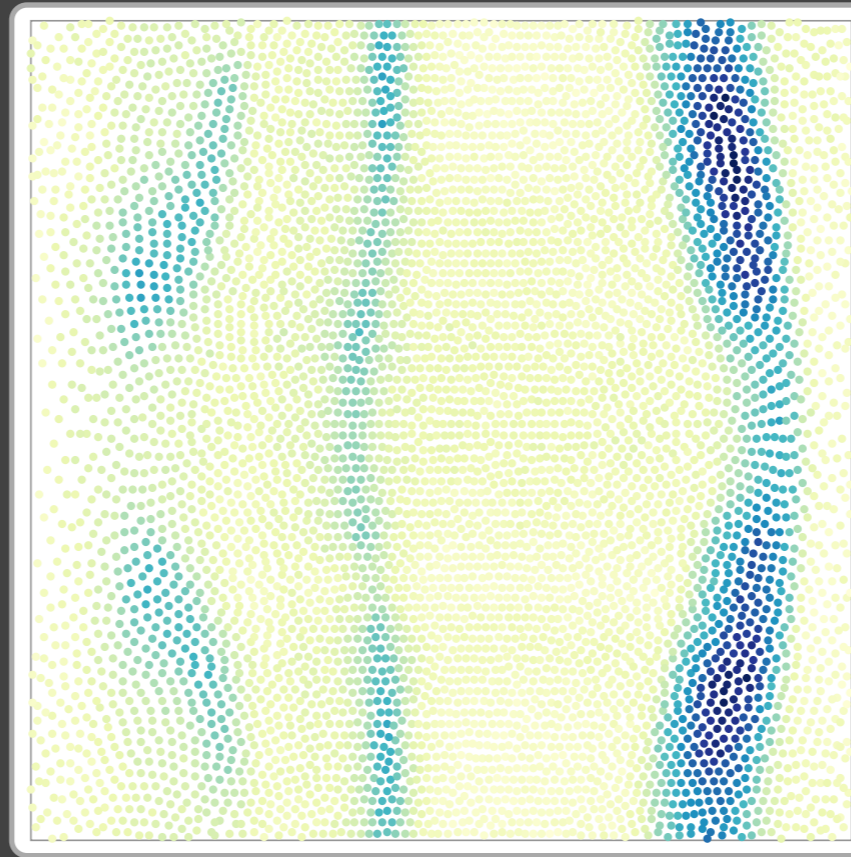
a value for each cell that summarizes the underlying data

**spatial**



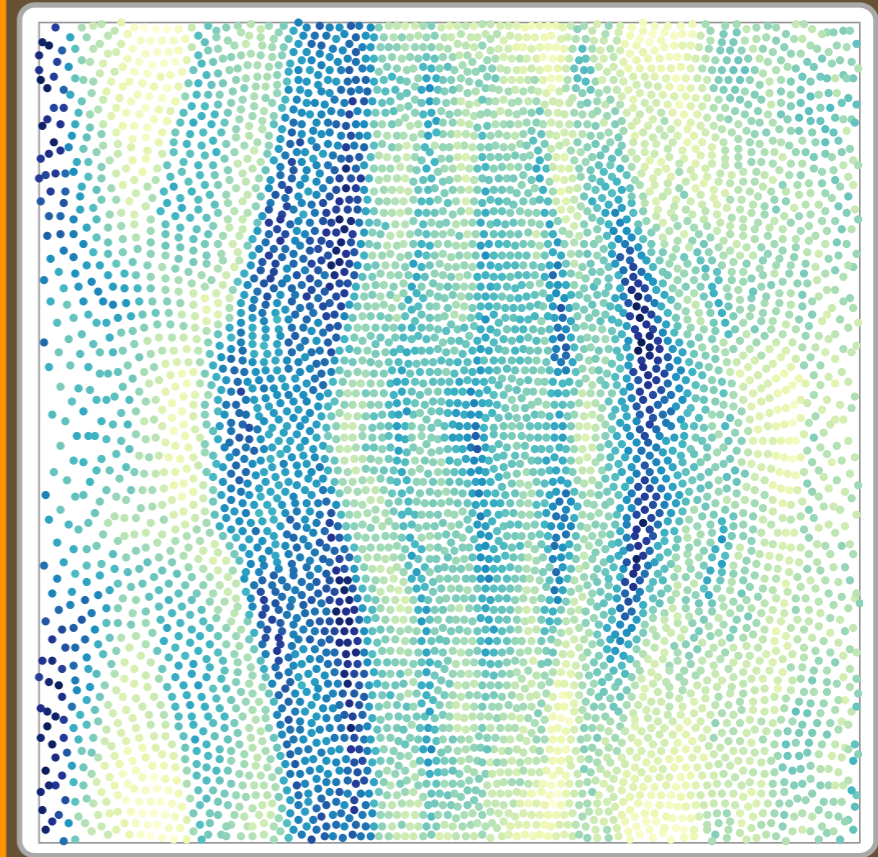
distance to cell  $i$

**expression**



gene  $j$  at timepoint  $k$

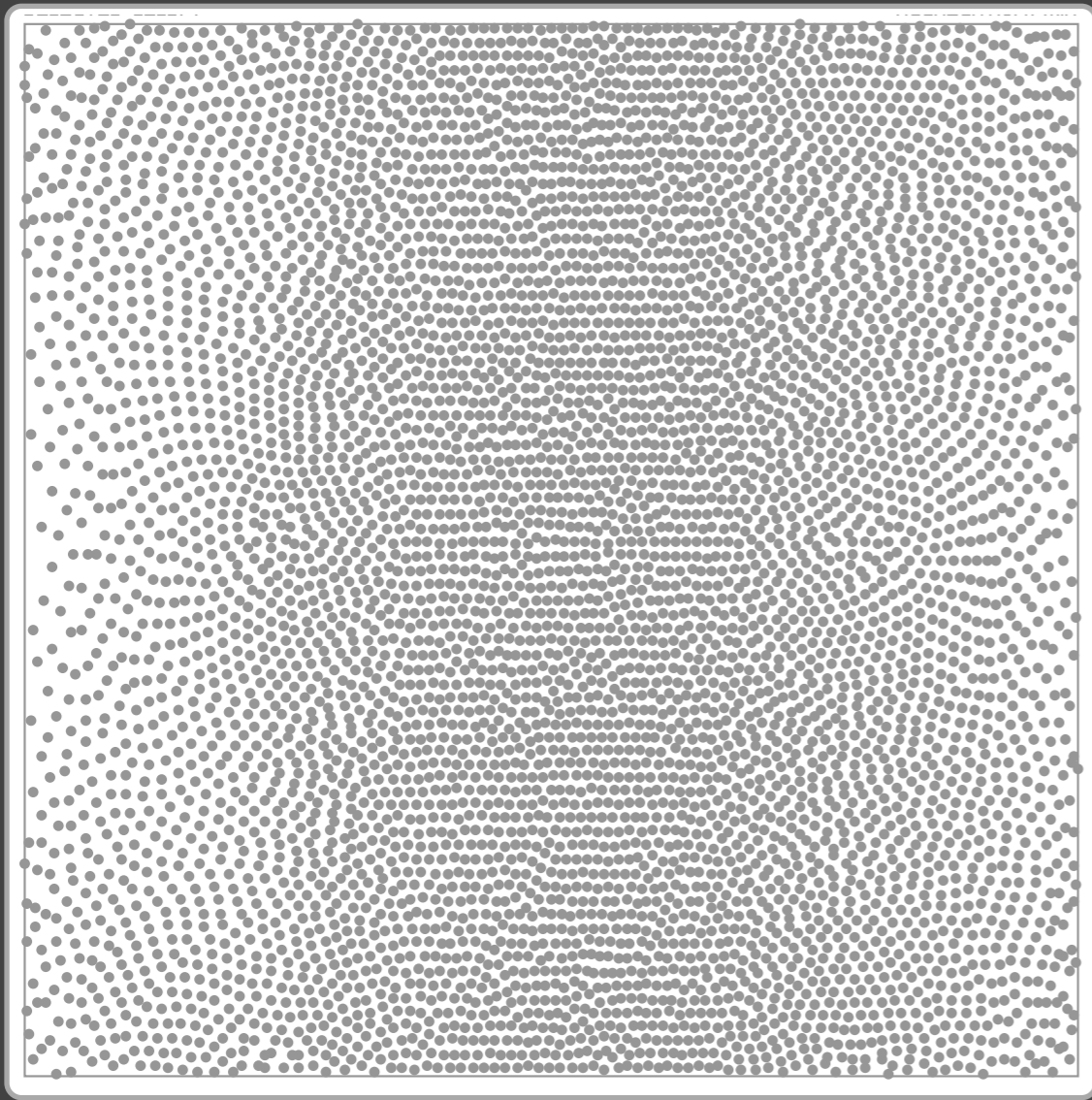
**combination**



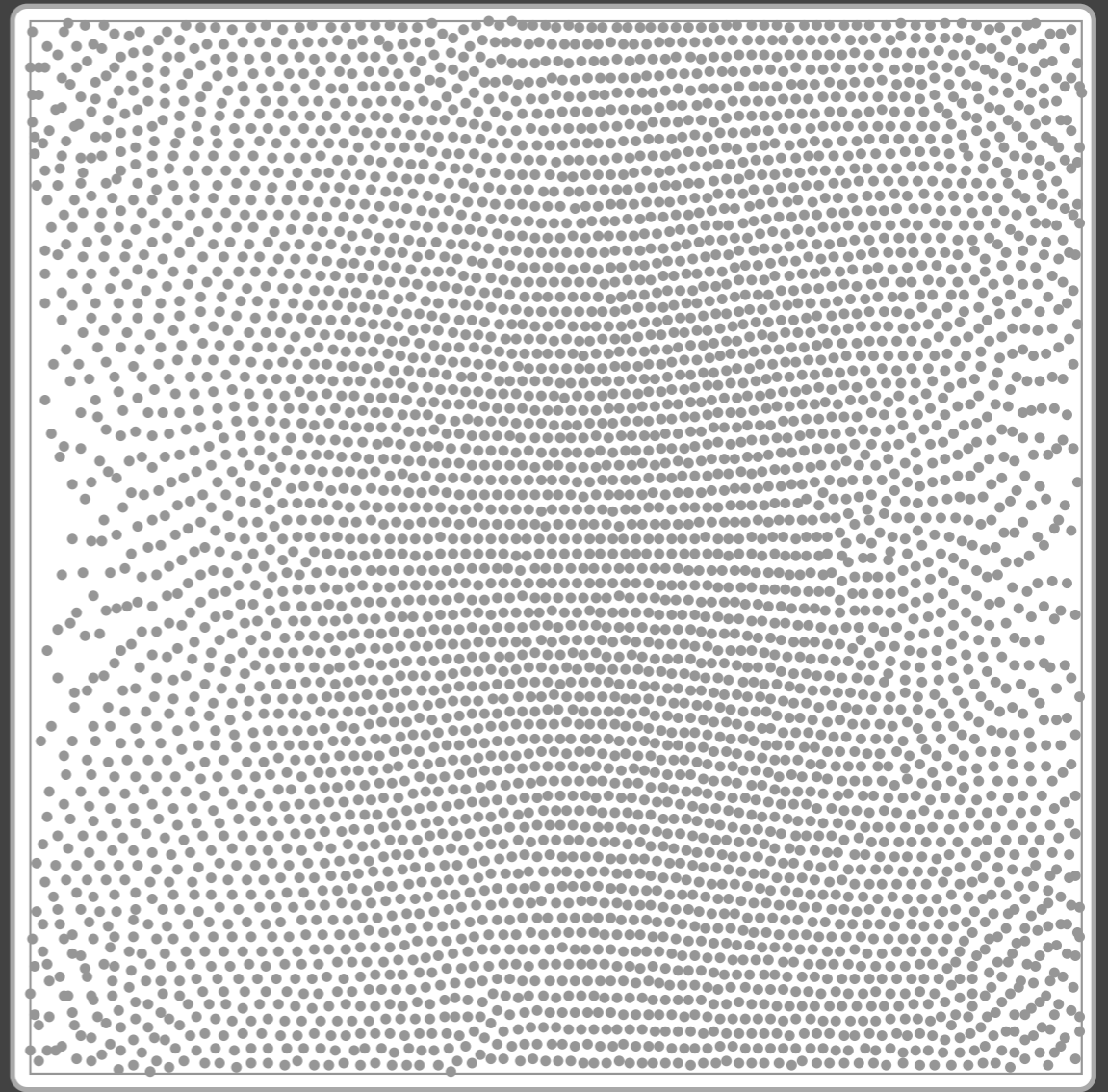
compare embryos  
 $A$  and  $B$

creating a comparison summary

# creating a comparison summary



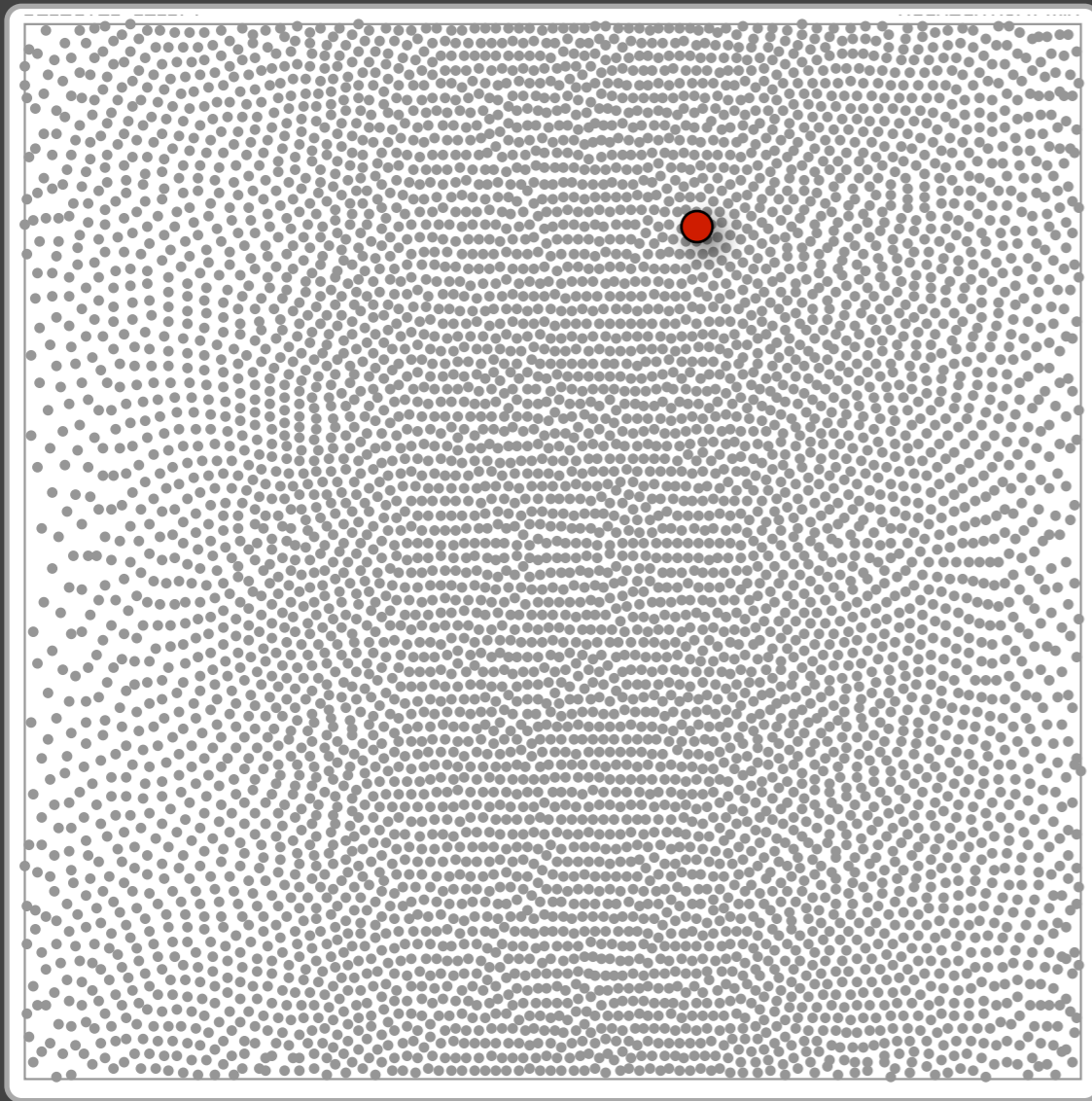
**embryo A**



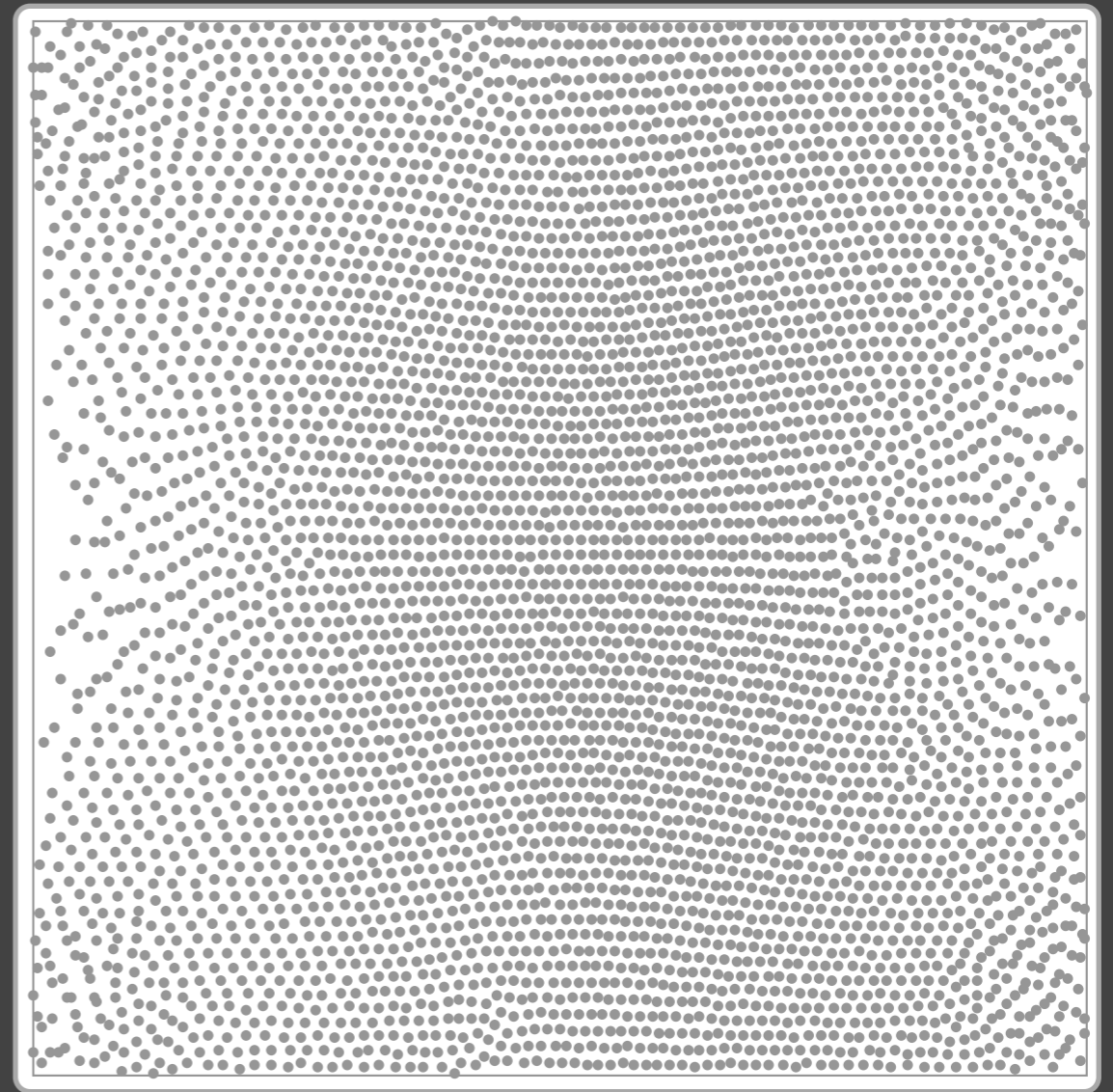
**embryo B**



# creating a comparison summary

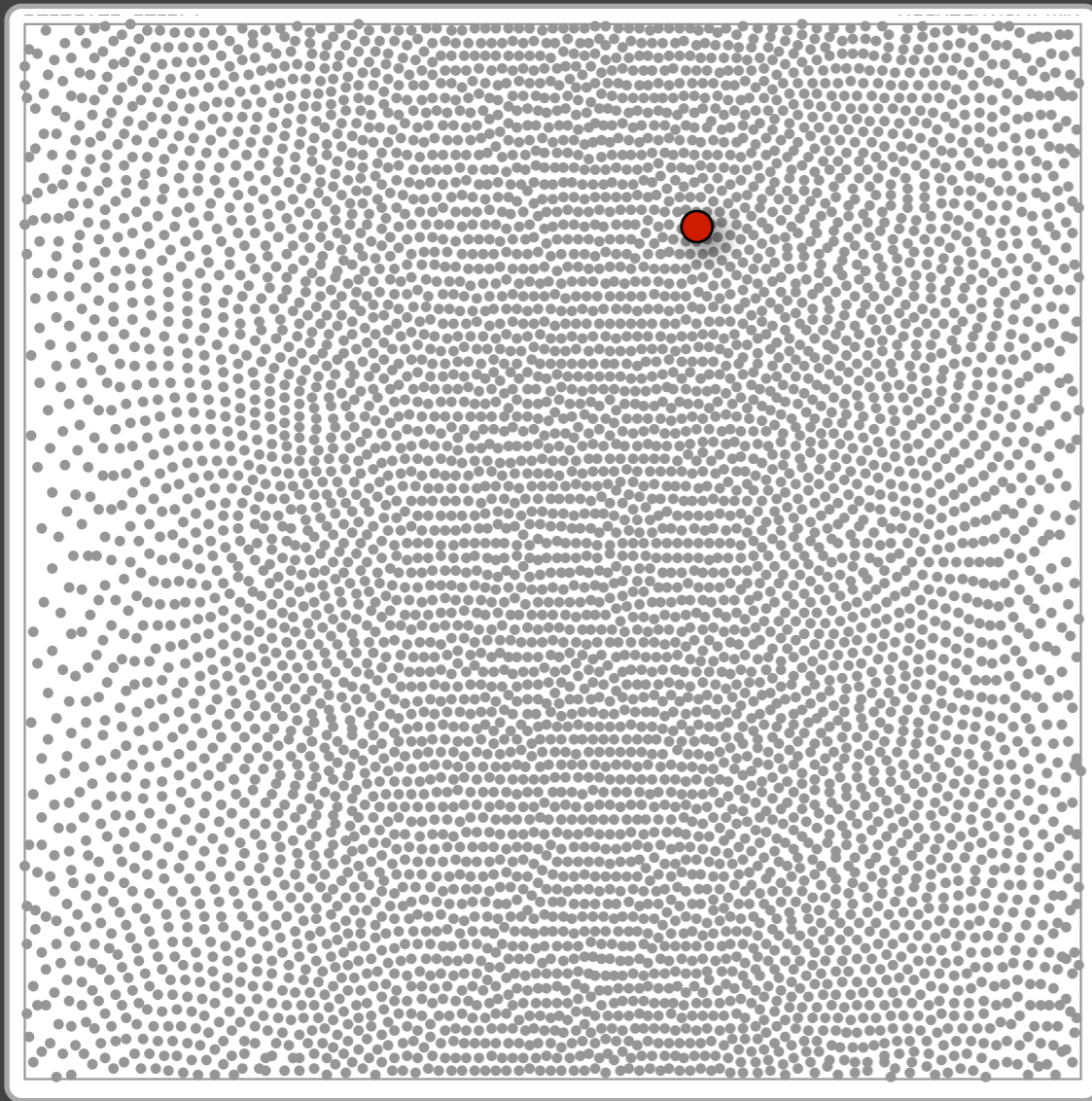


embryo A

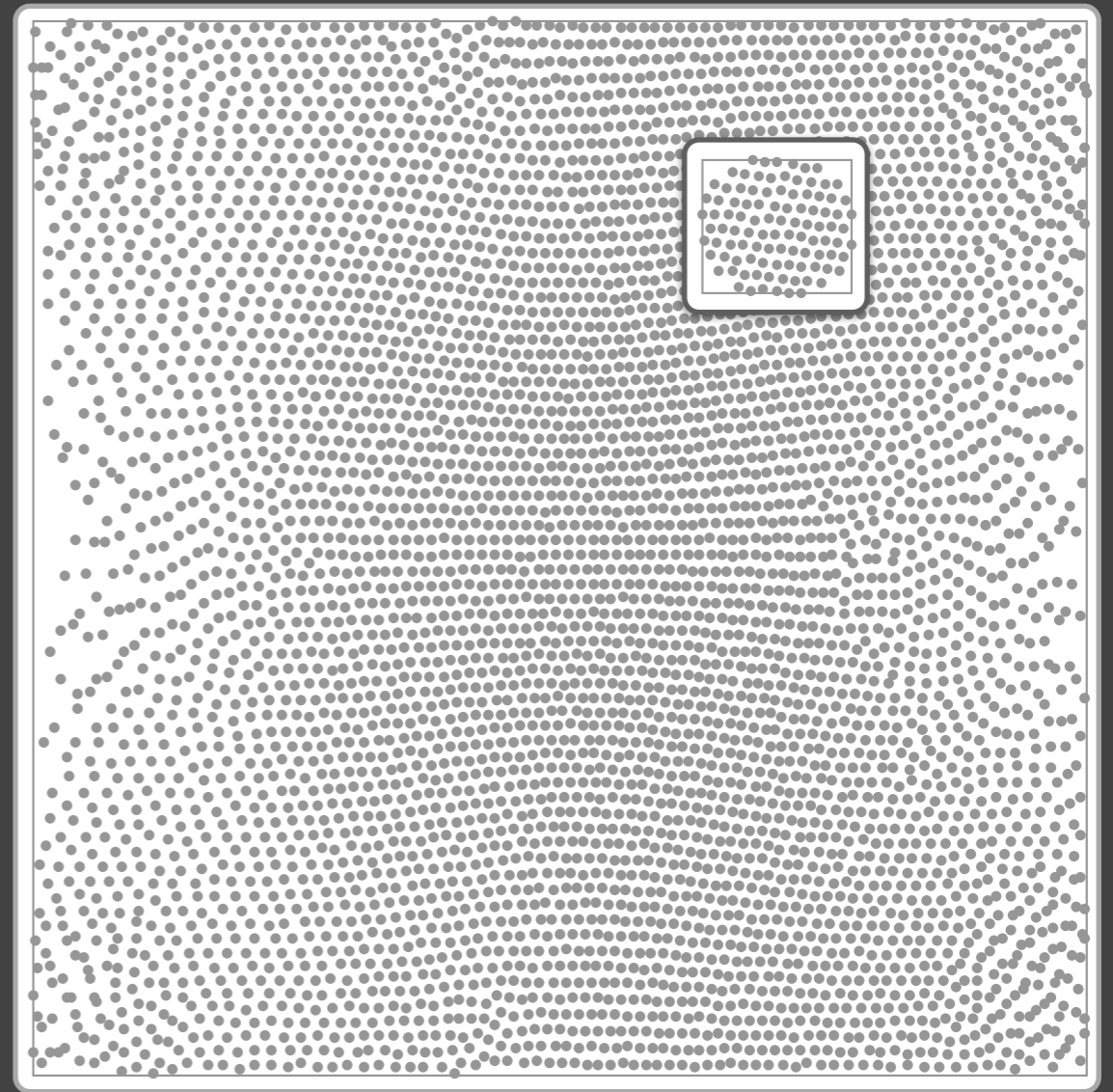


embryo B

# creating a comparison summary



embryo A

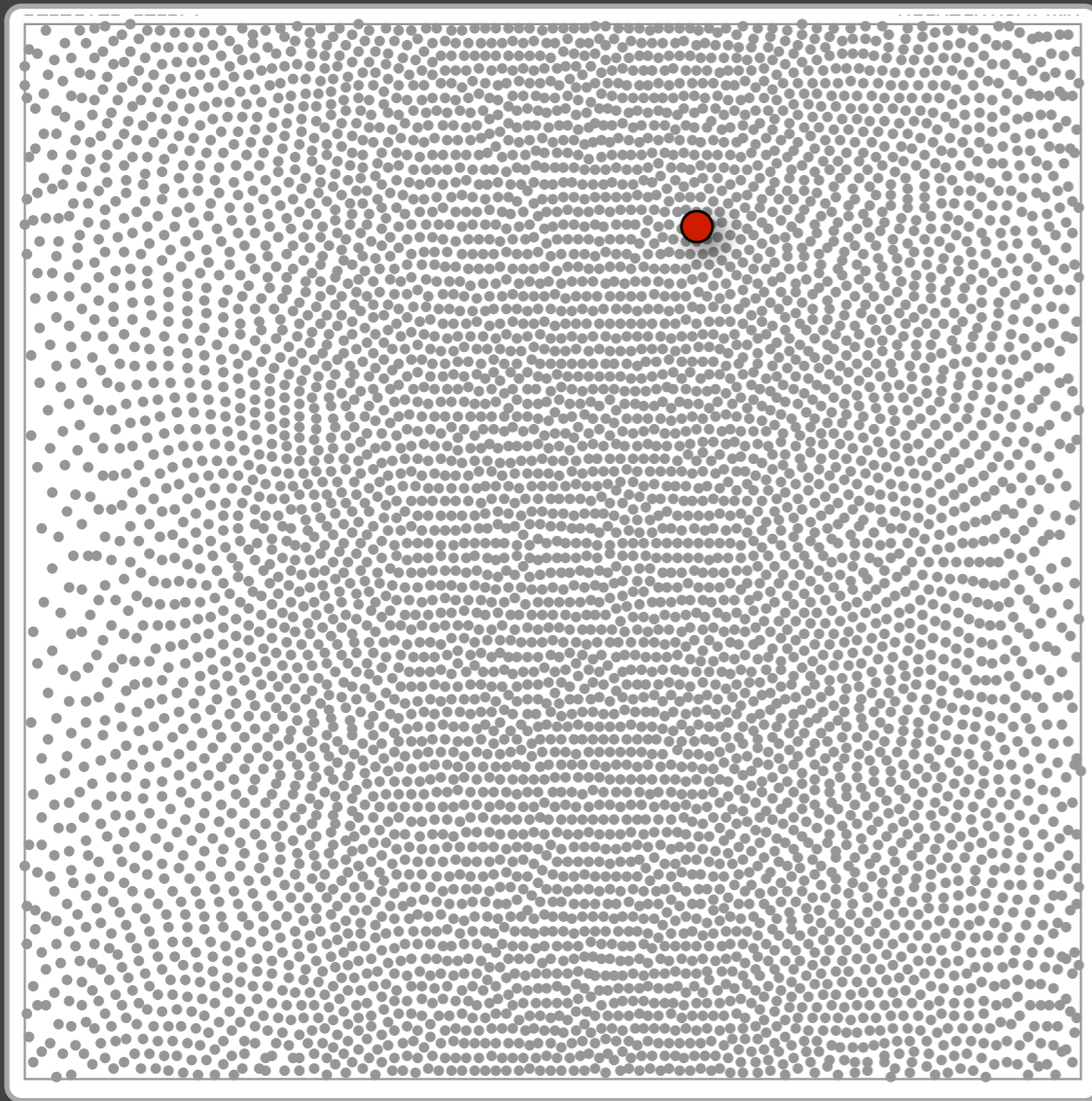


embryo B

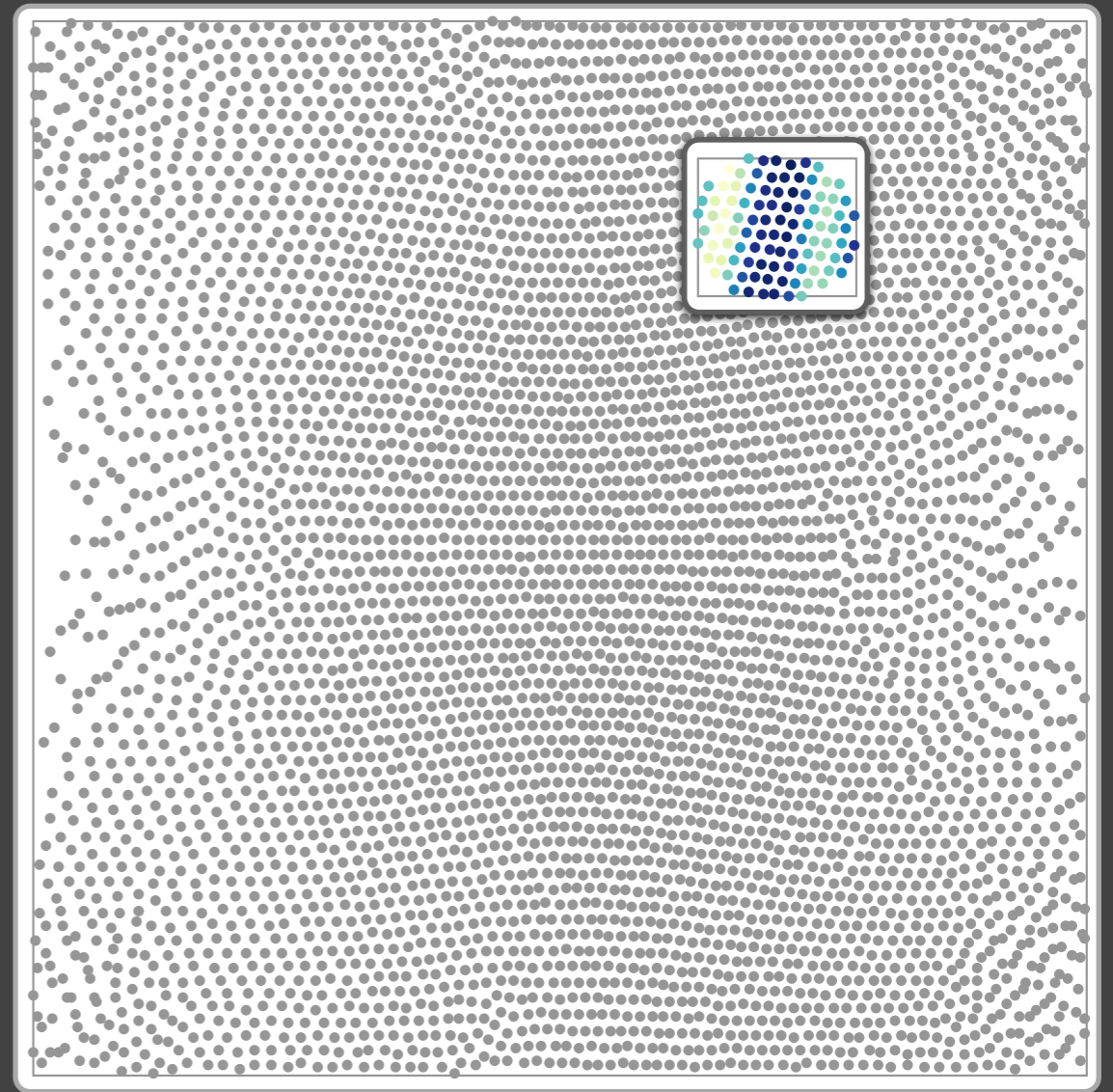


# creating a comparison summary

root-mean-square distance



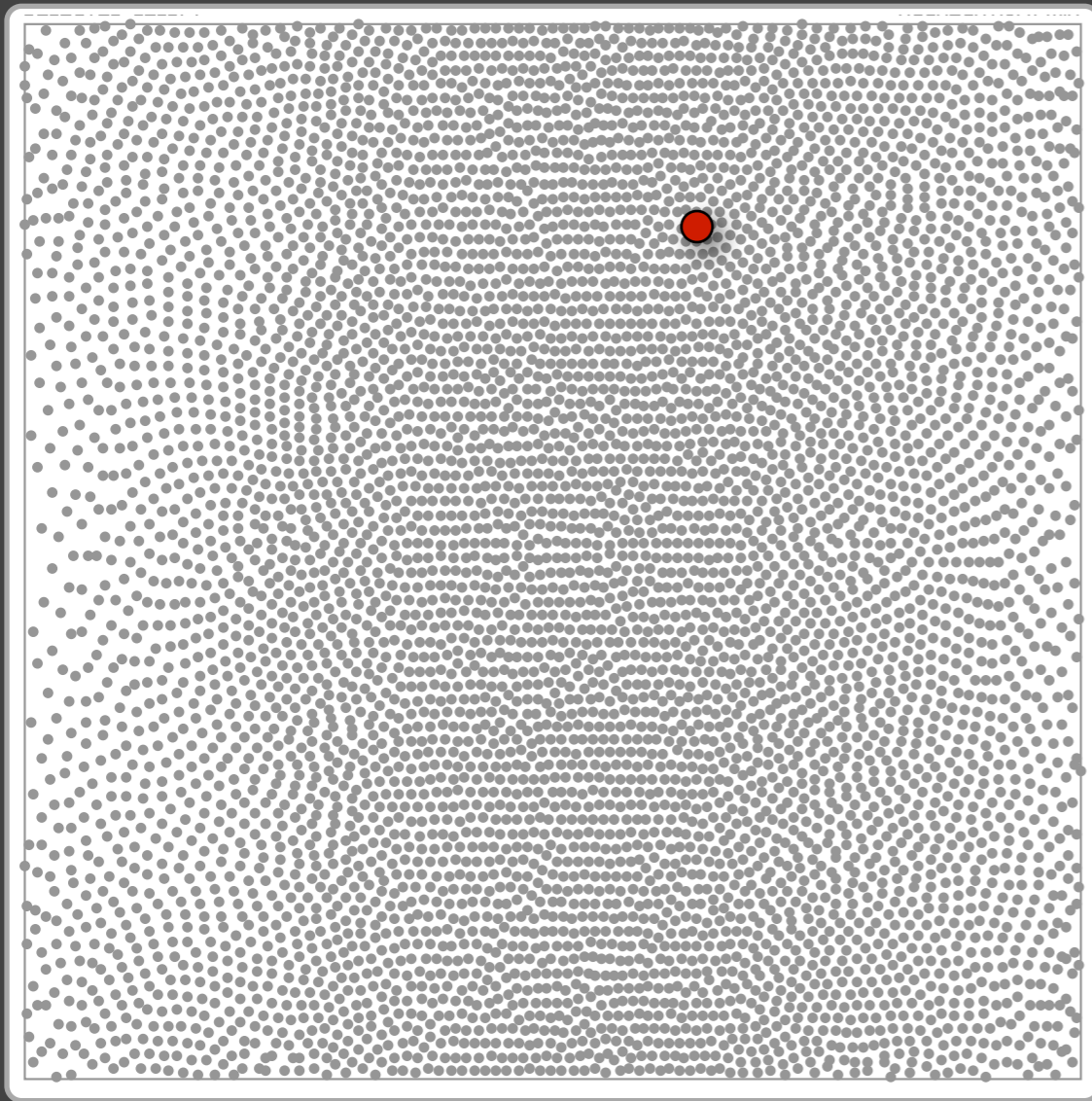
embryo A



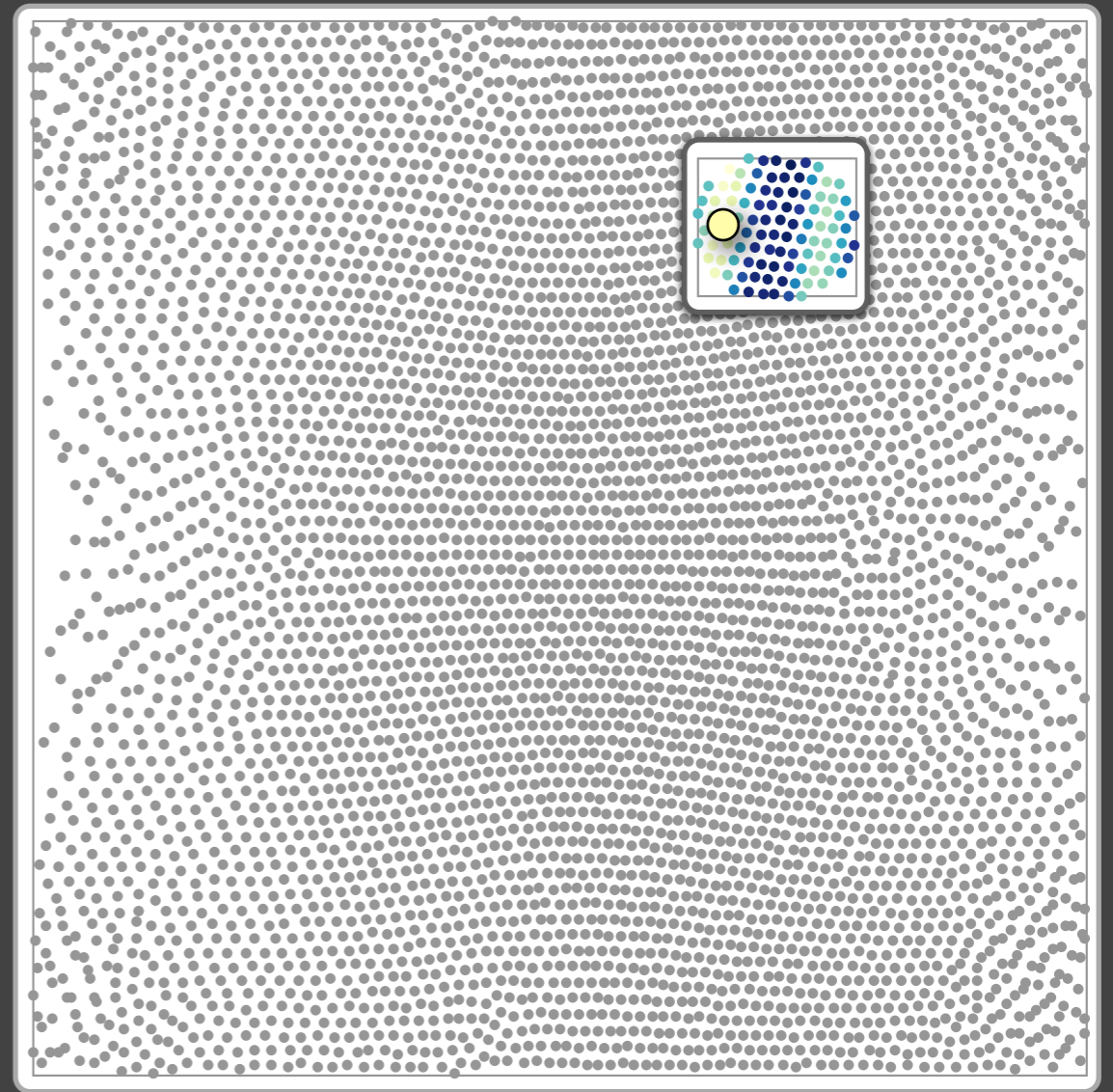
embryo B

# creating a comparison summary

root-mean-square distance



embryo A

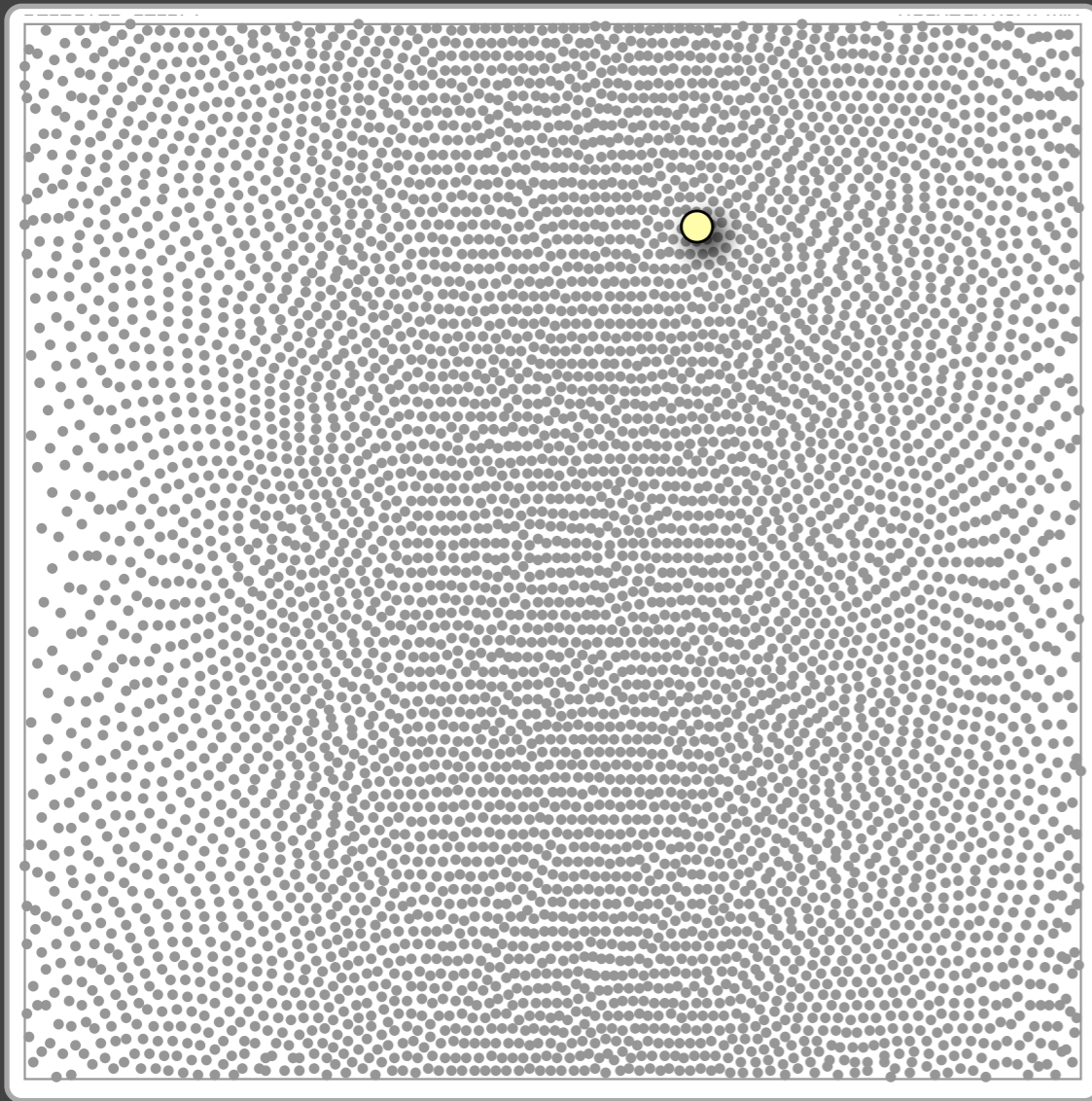


embryo B

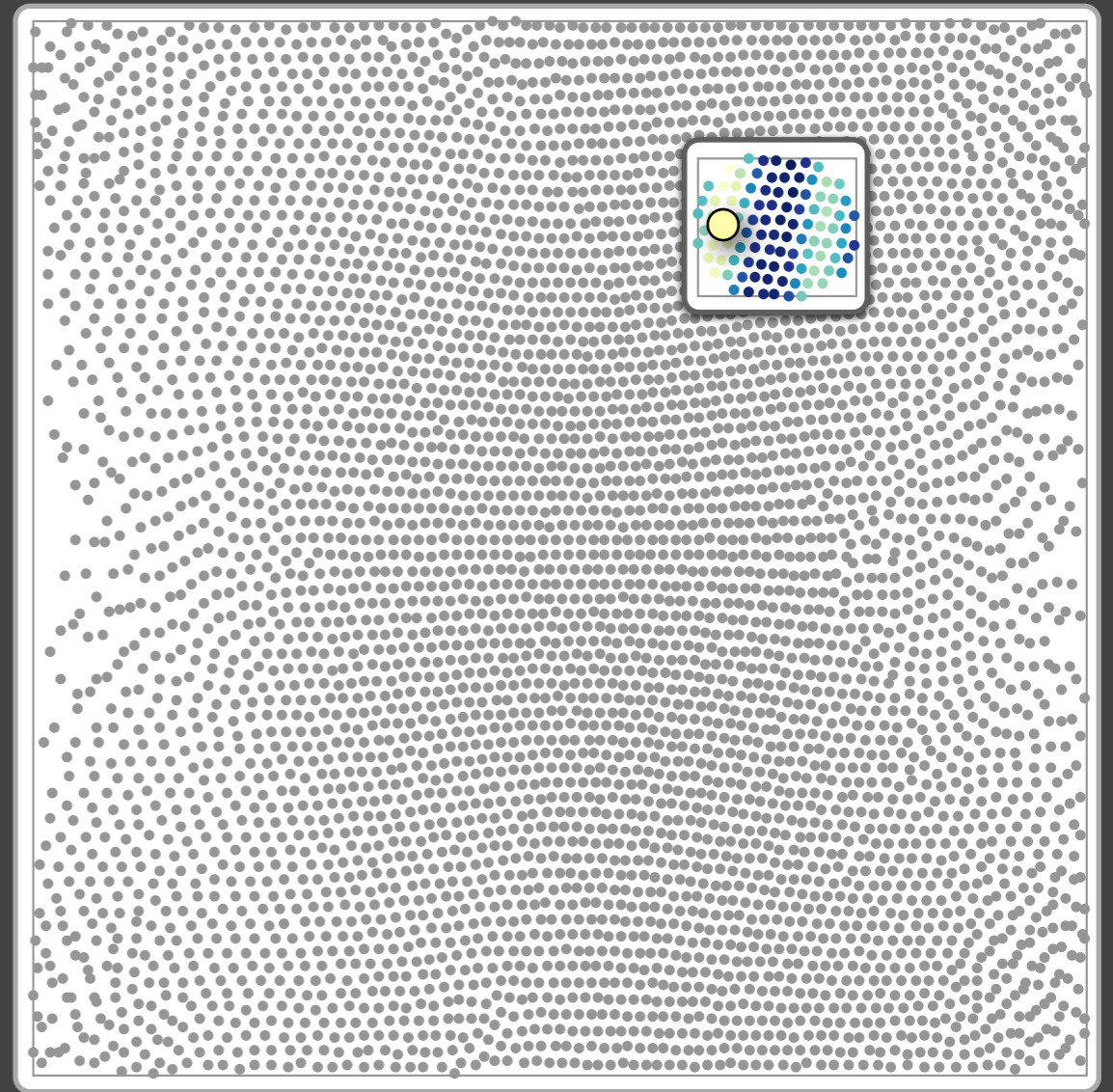


# creating a comparison summary

root-mean-square distance



embryo A

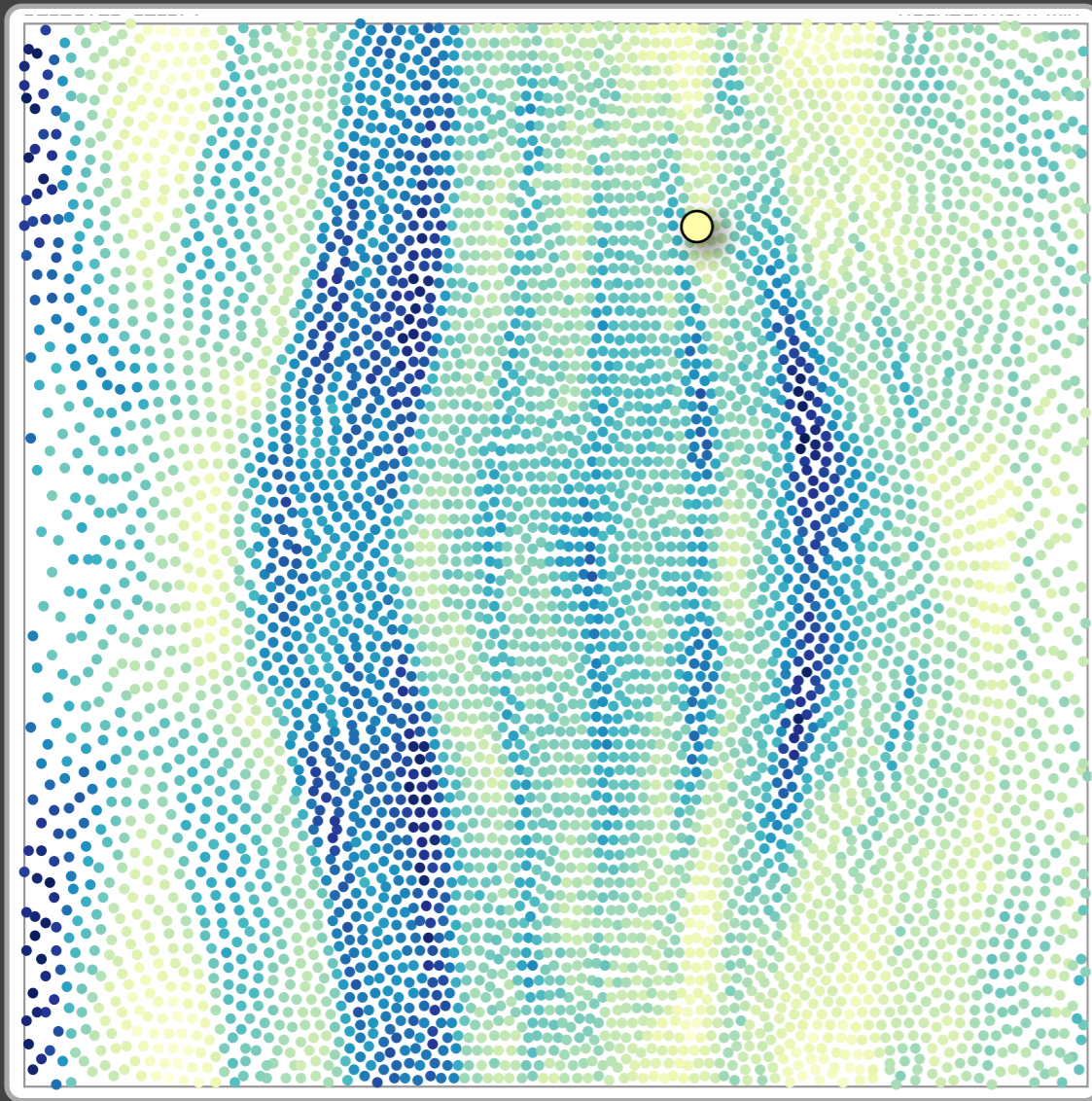


embryo B

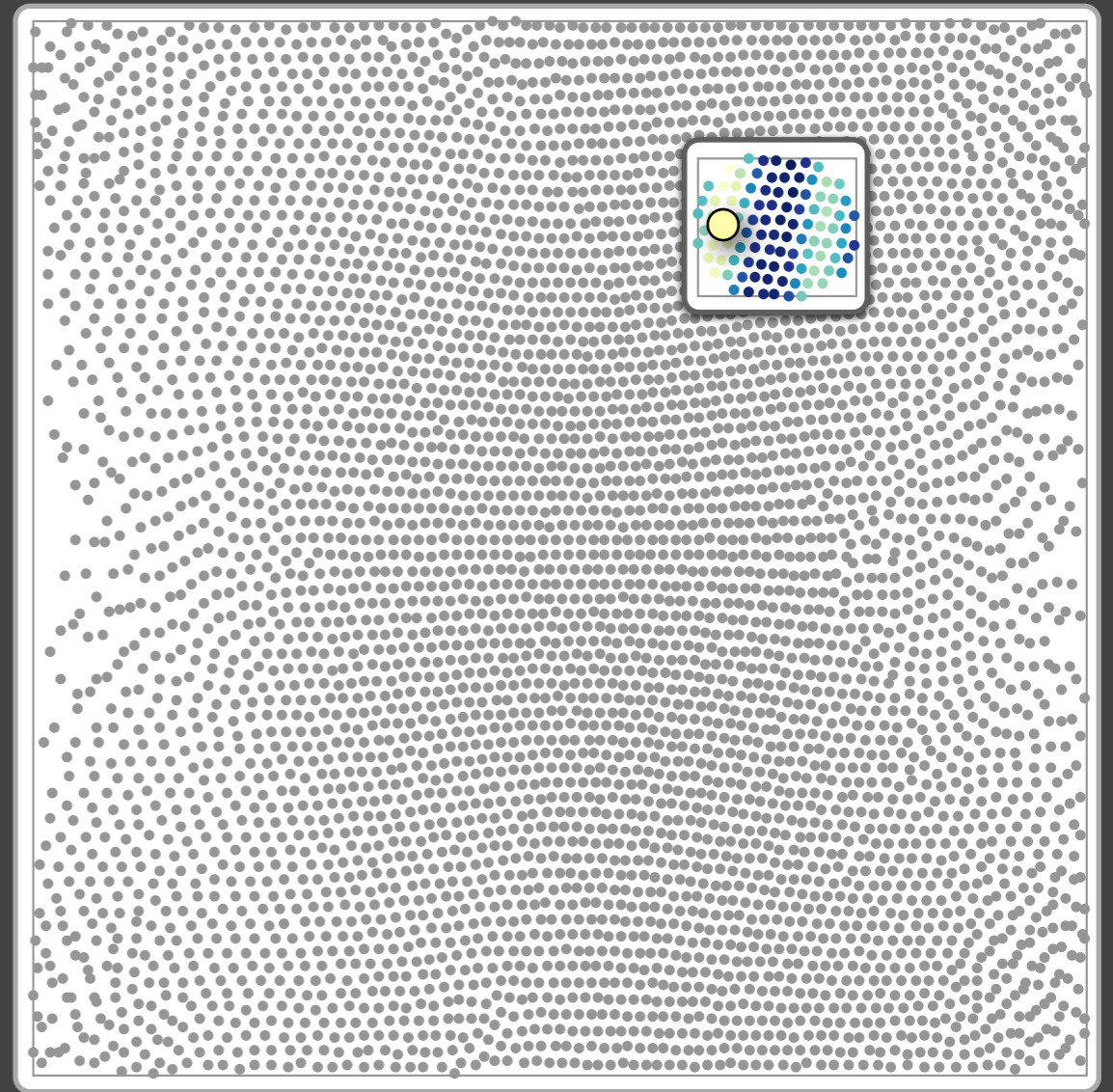


# creating a comparison summary

root-mean-square distance



embryo A

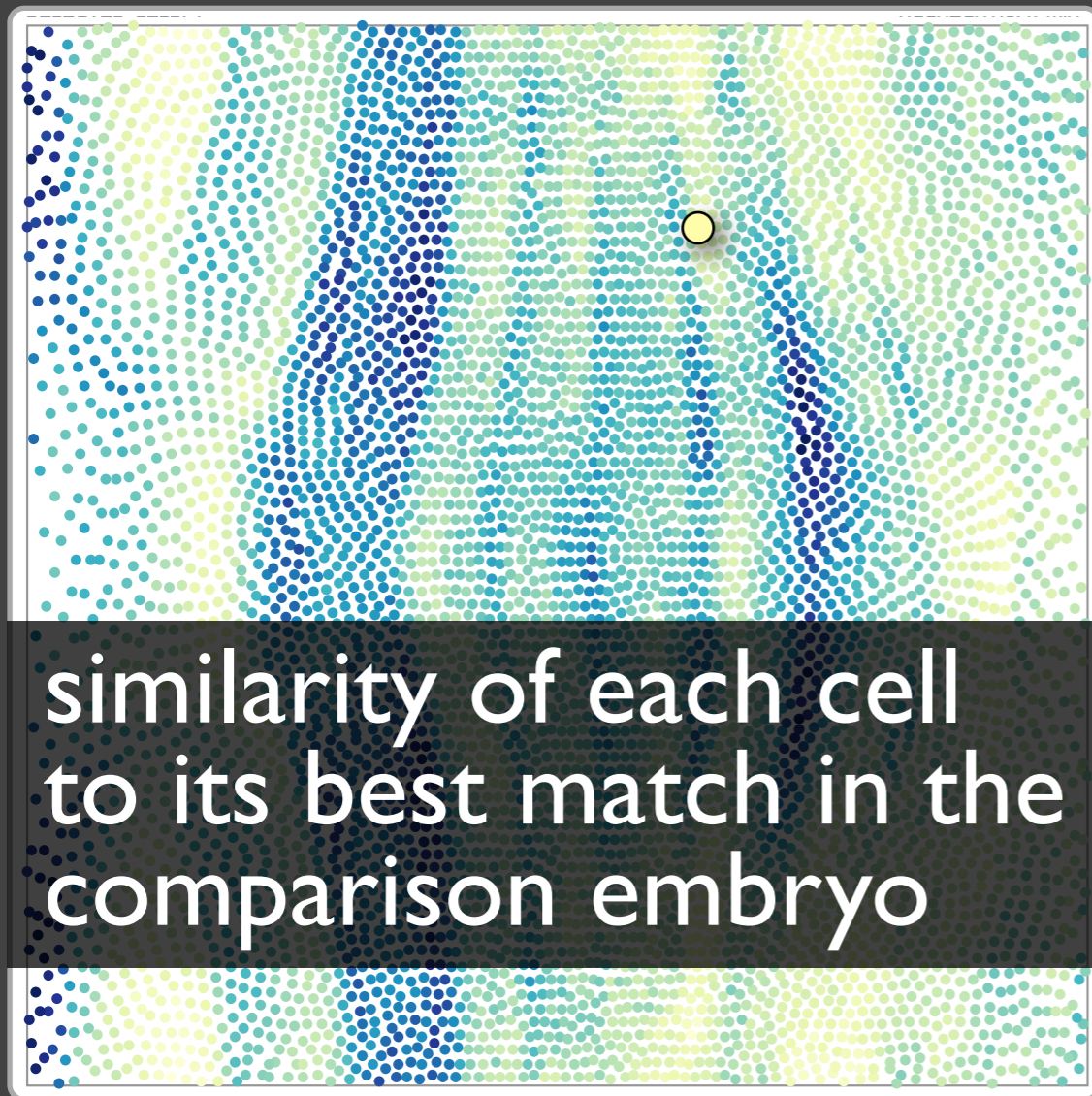


embryo B

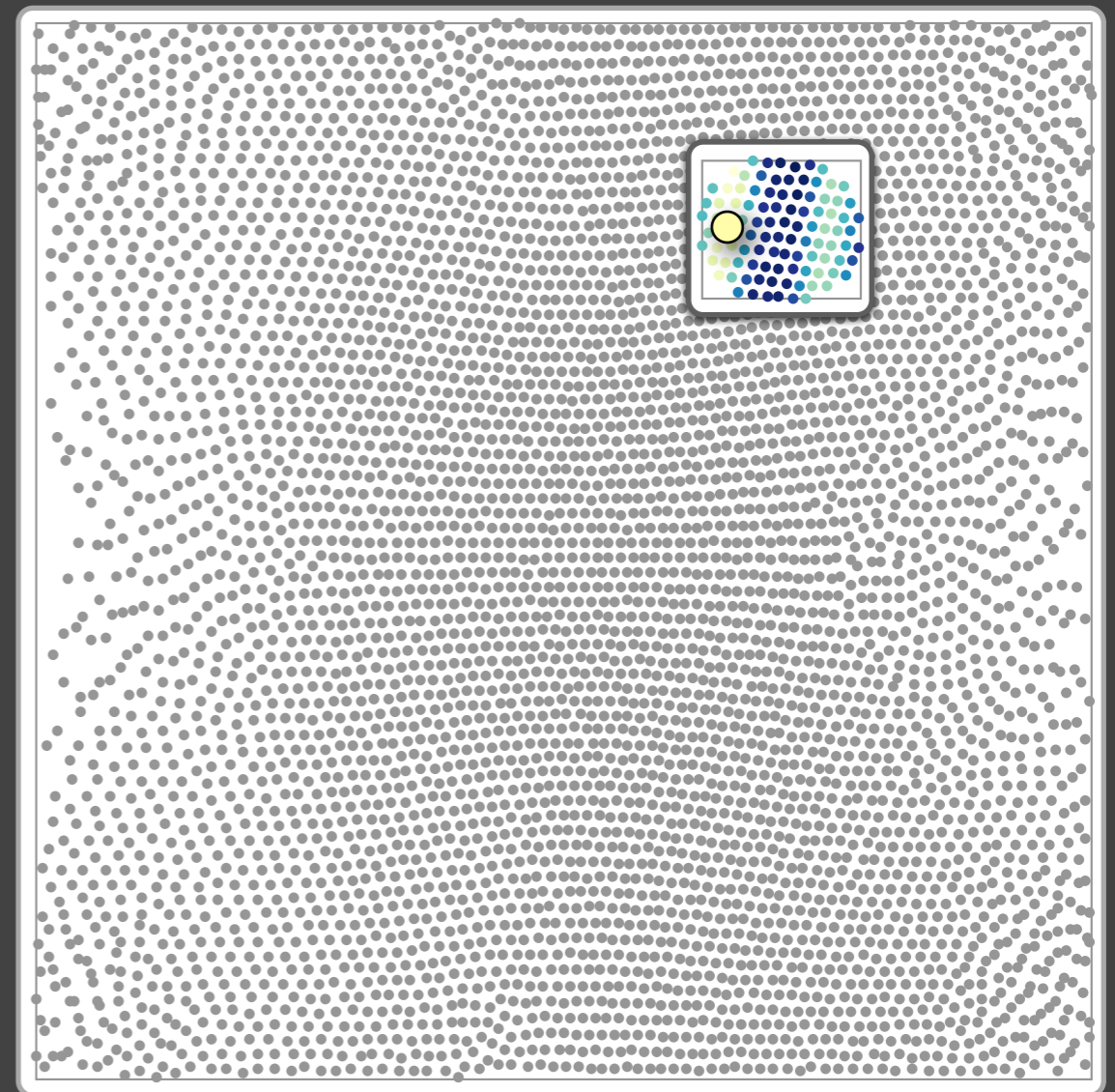


# creating a comparison summary

root-mean-square distance



embryo A



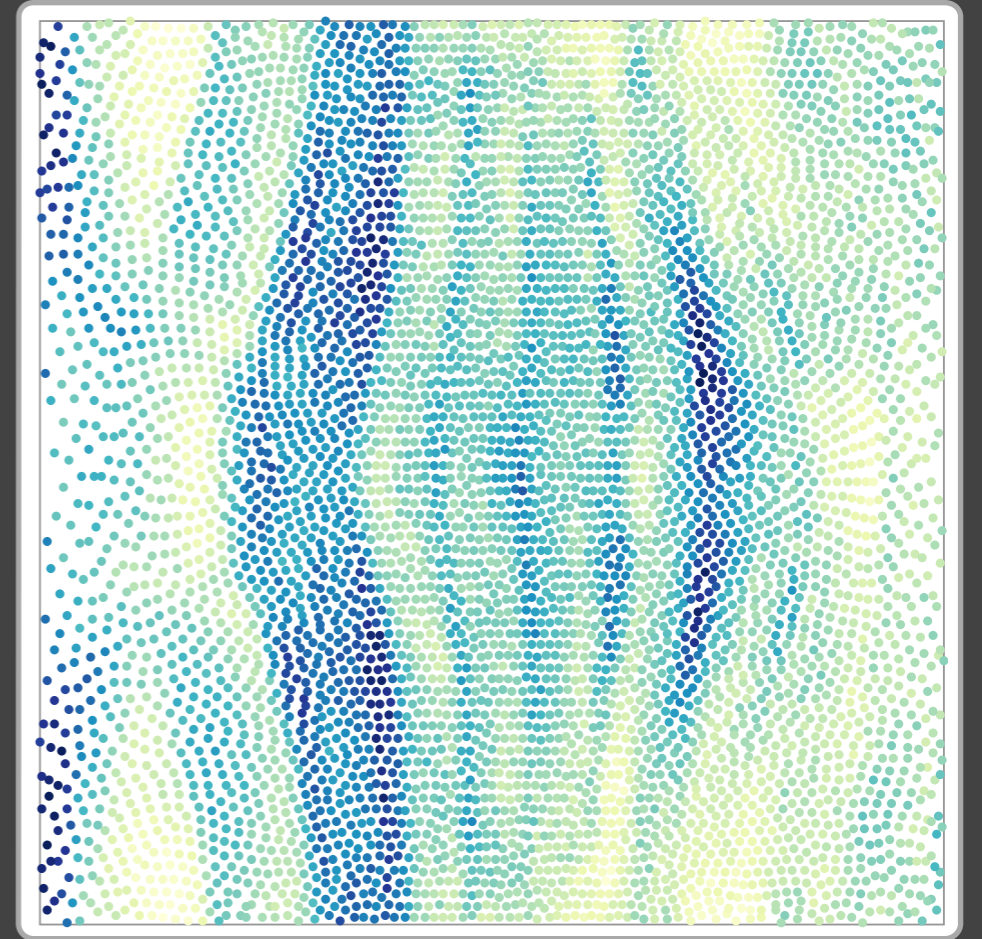
embryo B

# comparative summary components

aggregation group

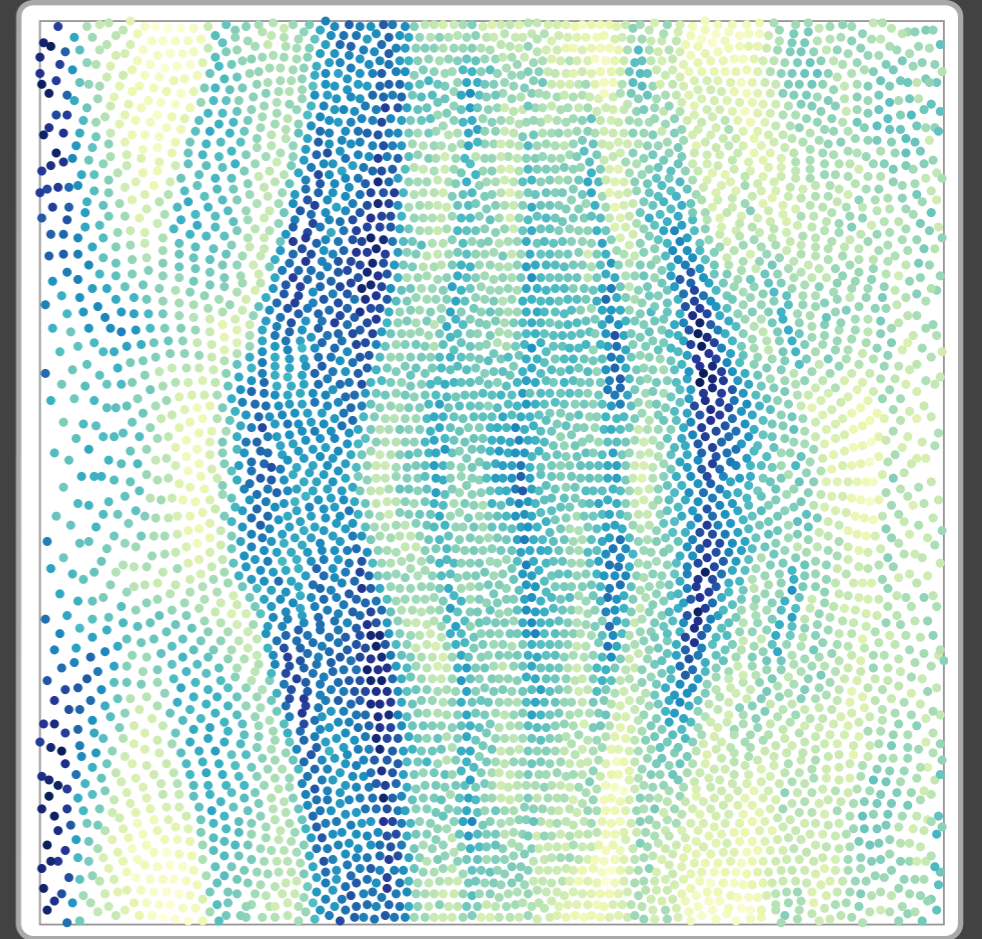
metric

aggregation





# comparative summary components

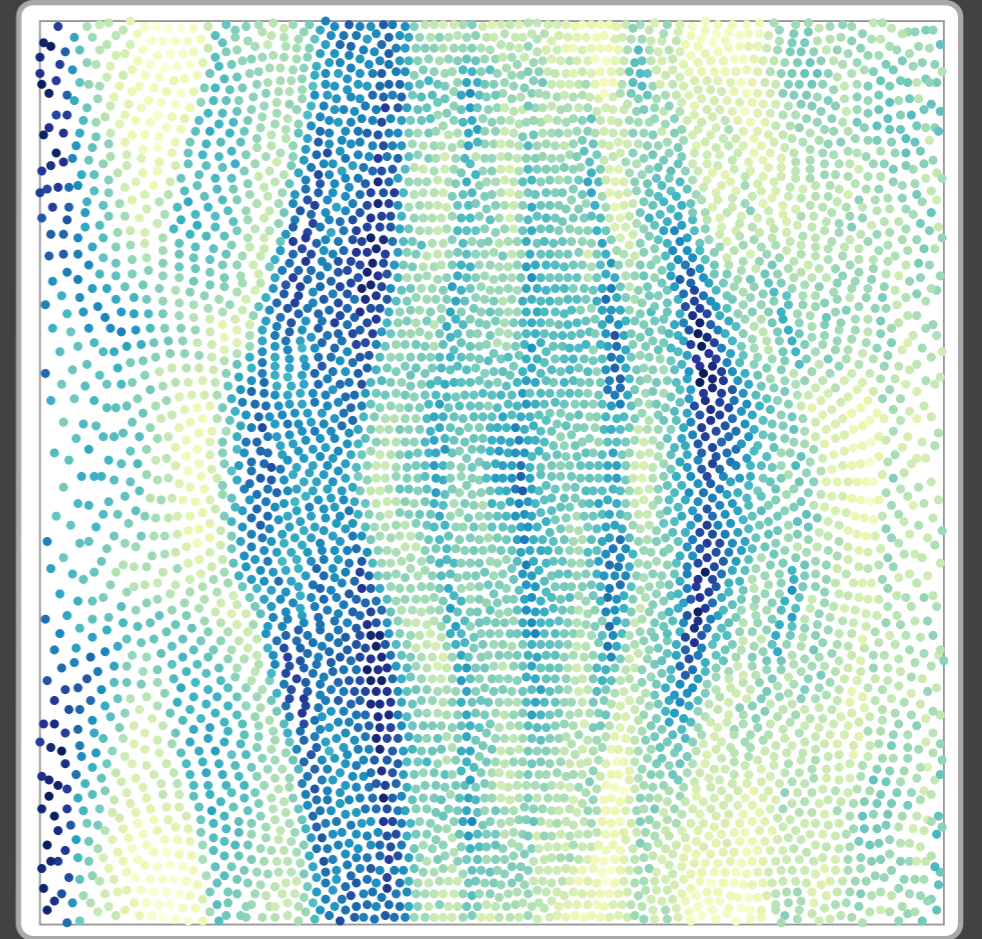


aggregation group: *100 spatially closest cells*

metric

aggregation

# comparative summary components

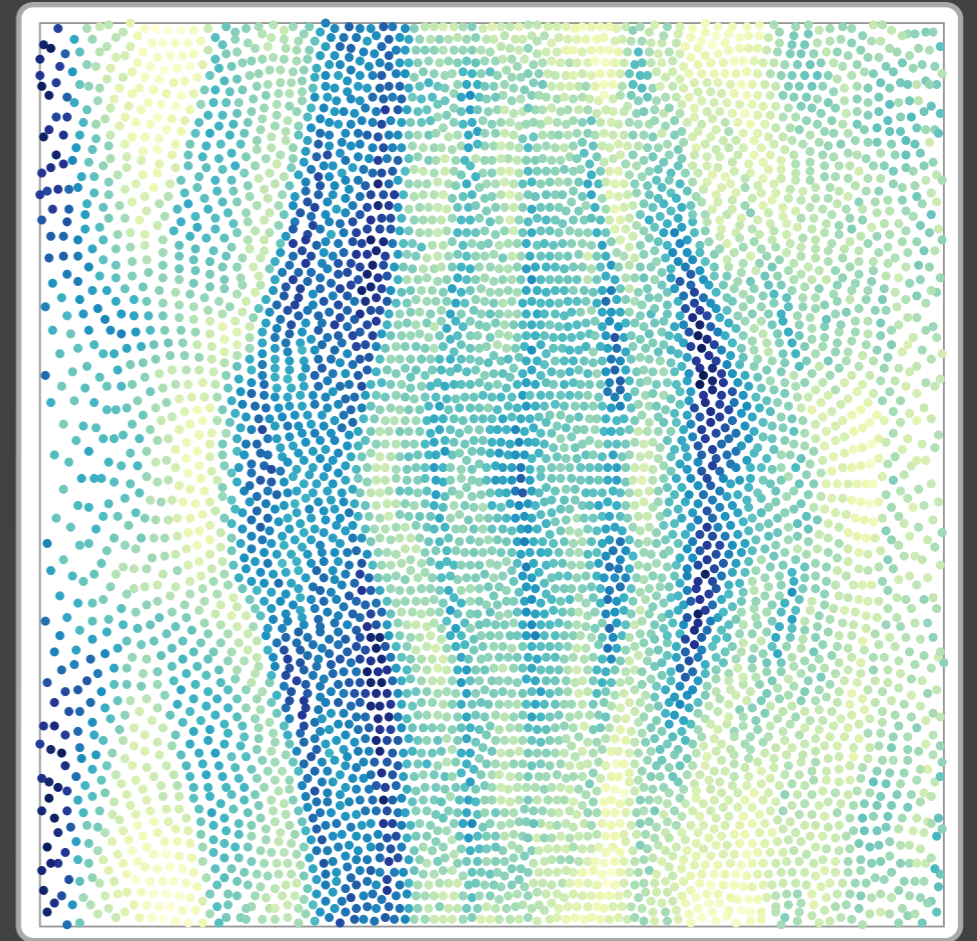


**aggregation group:** *100 spatially closest cells*

**metric:** *root-mean-square distance*

**aggregation**

## comparative summary components



**aggregation group:** *100 spatially closest cells*

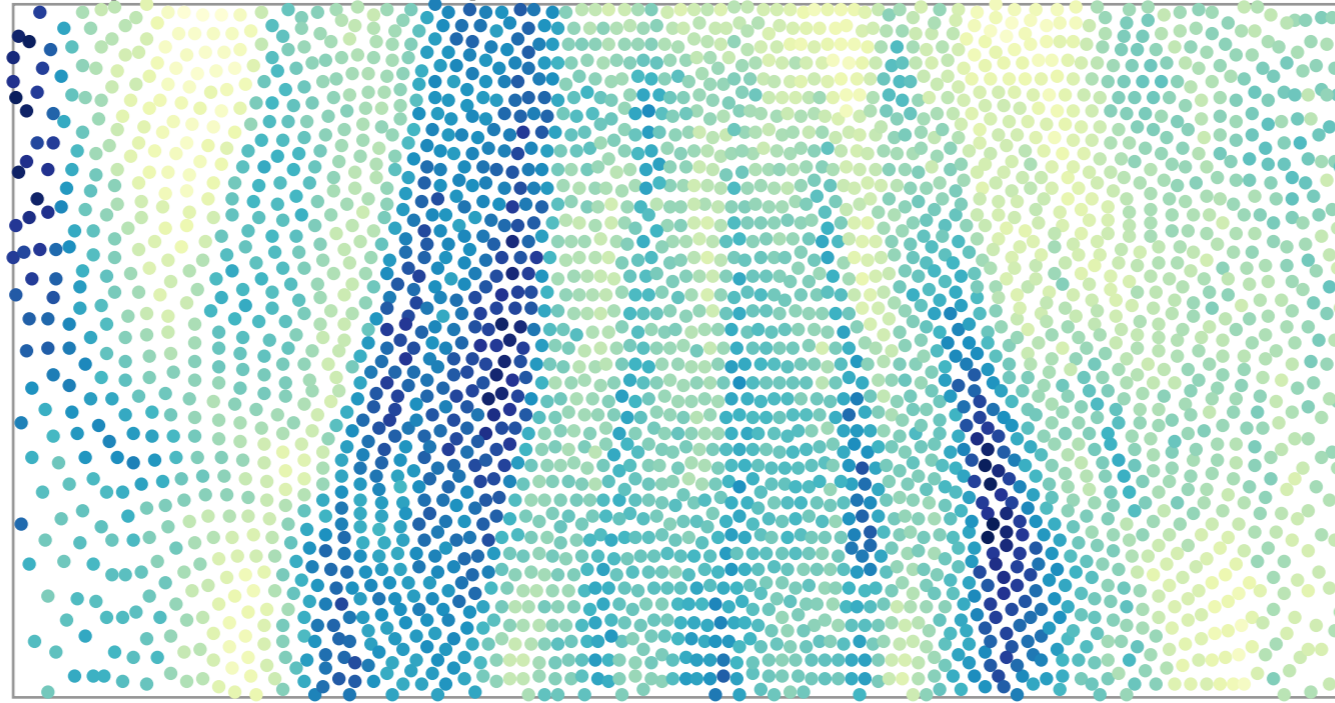
**metric:** *root-mean-square distance*

**aggregation:** *min operator*



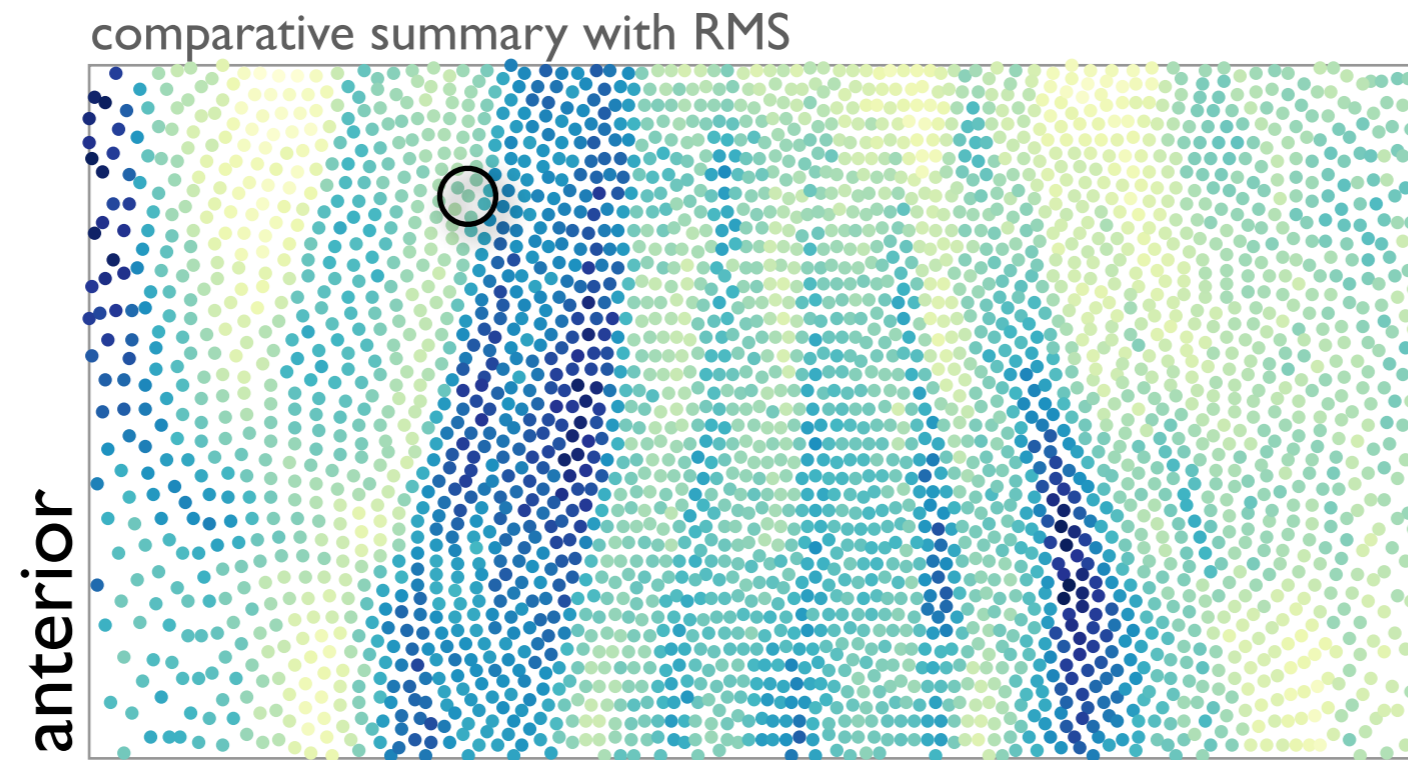
# case study one: characterizing a summary

comparative summary with RMS

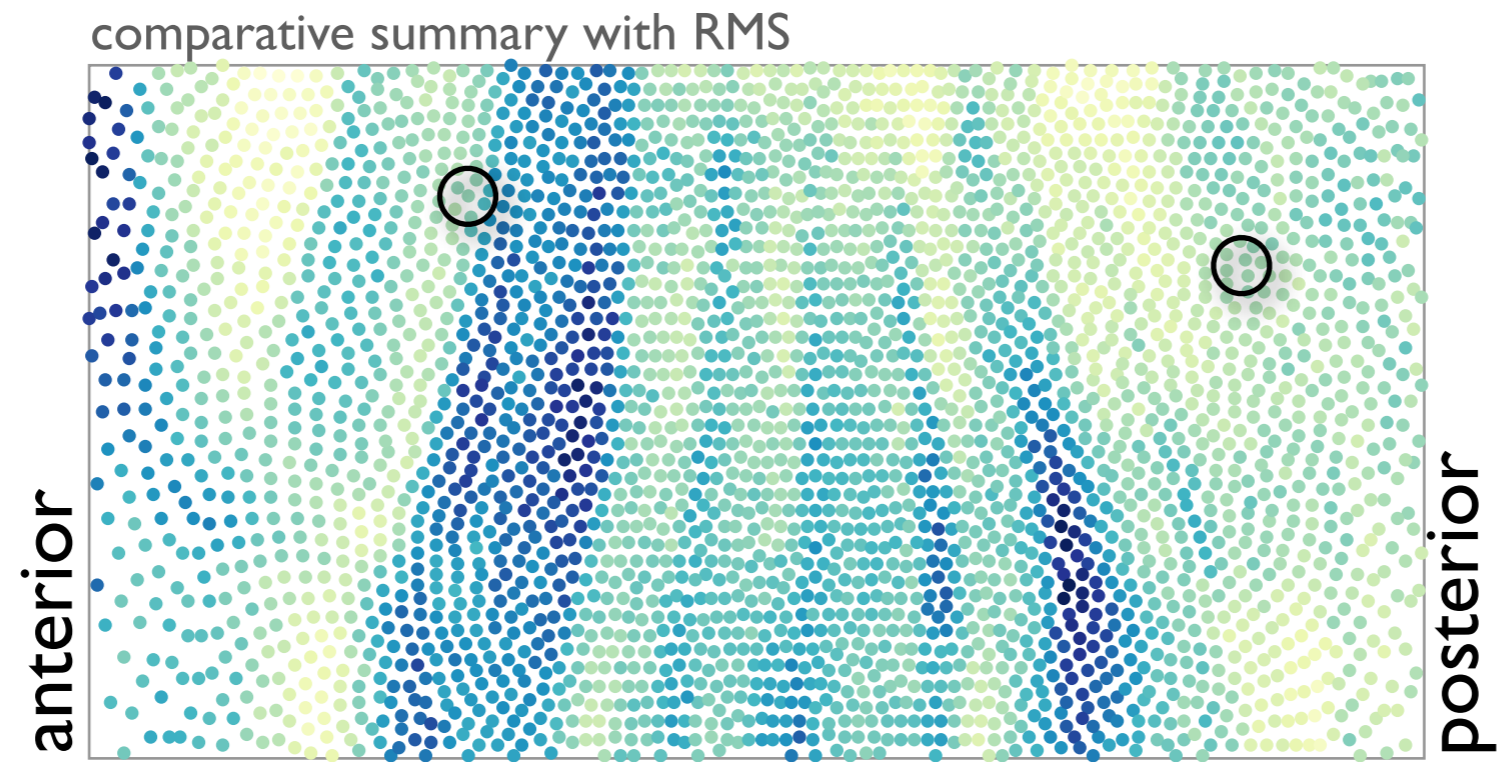




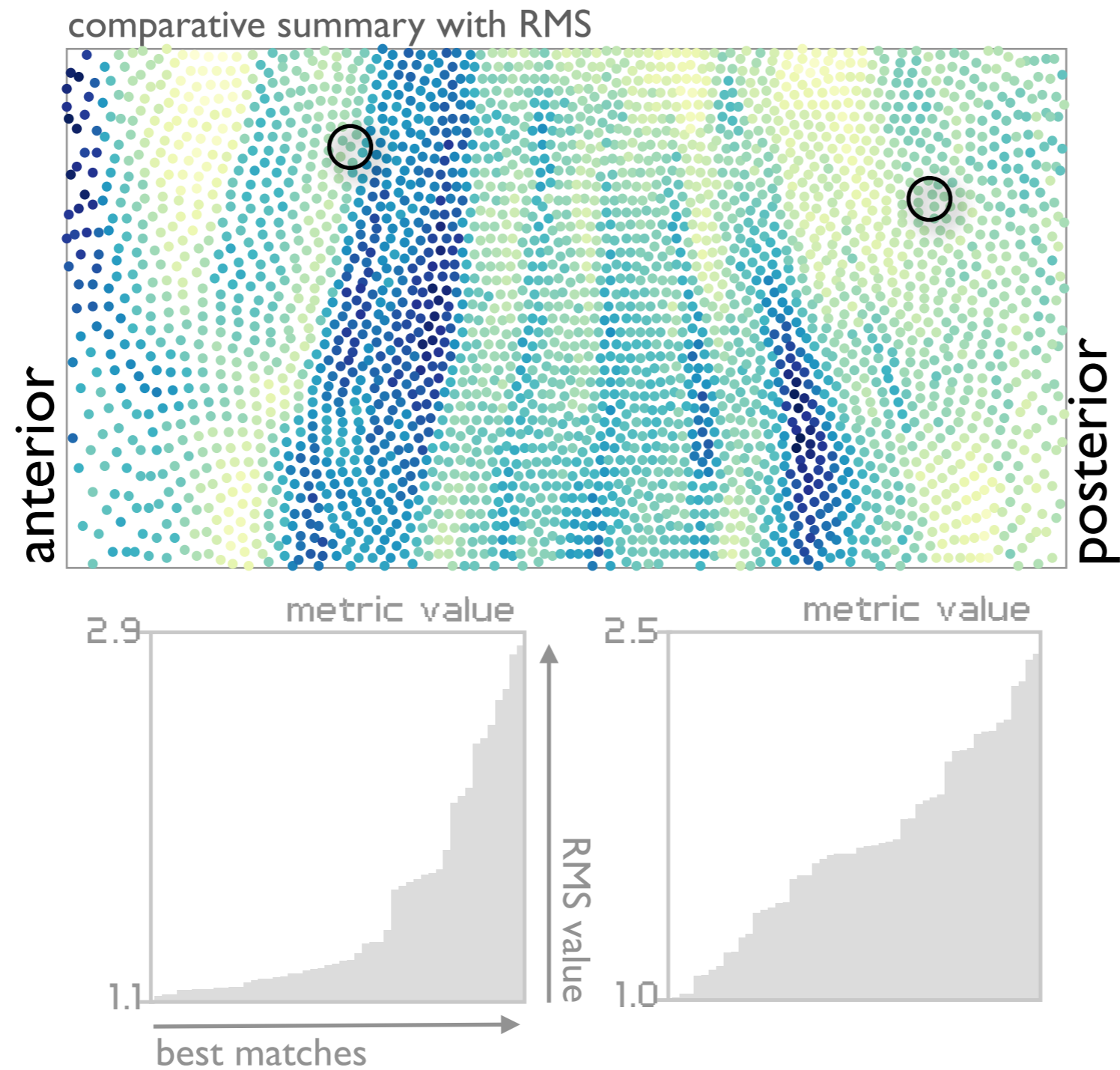
# case study one: characterizing a summary



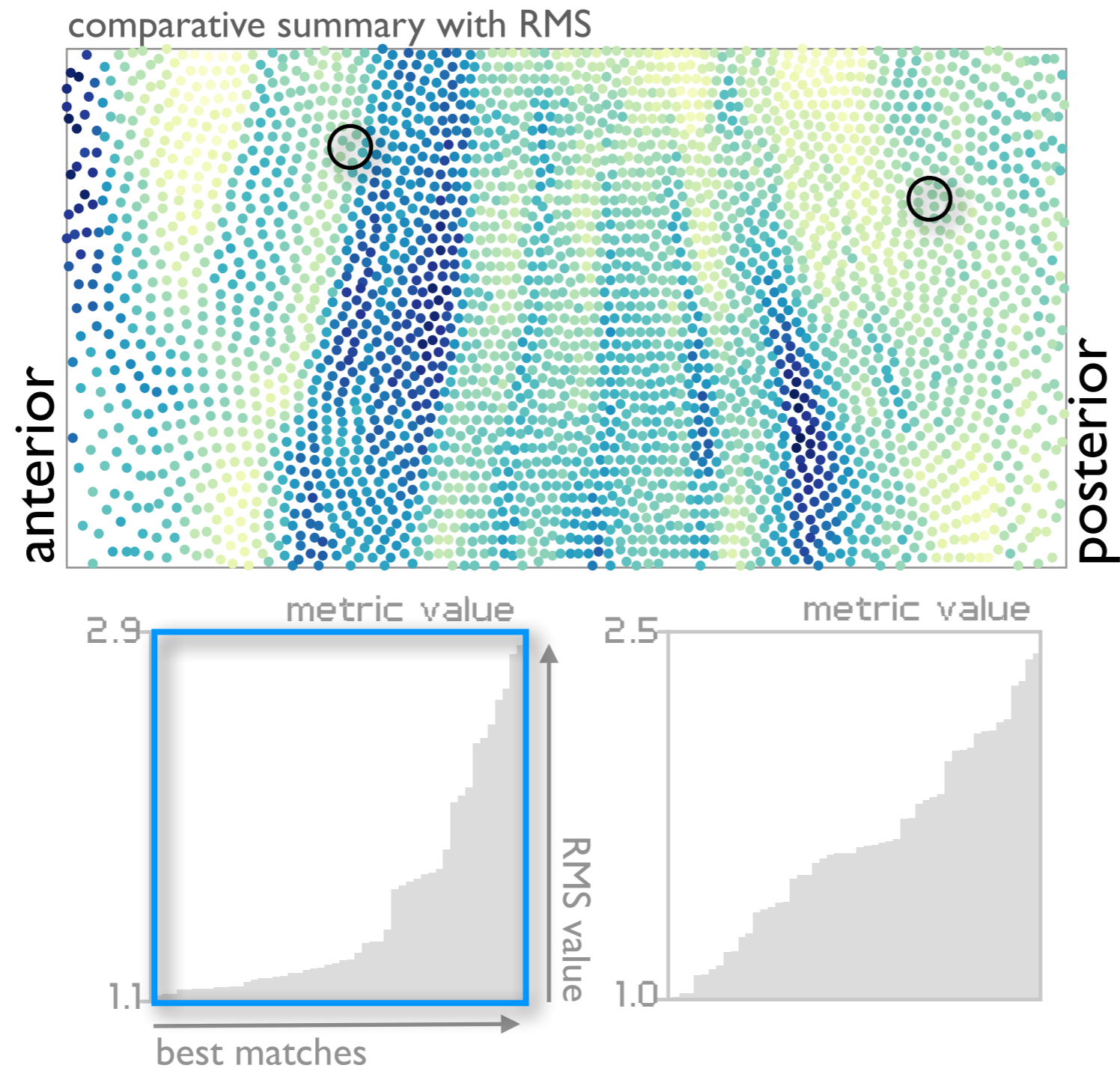
# case study one: characterizing a summary



# case study one: characterizing a summary

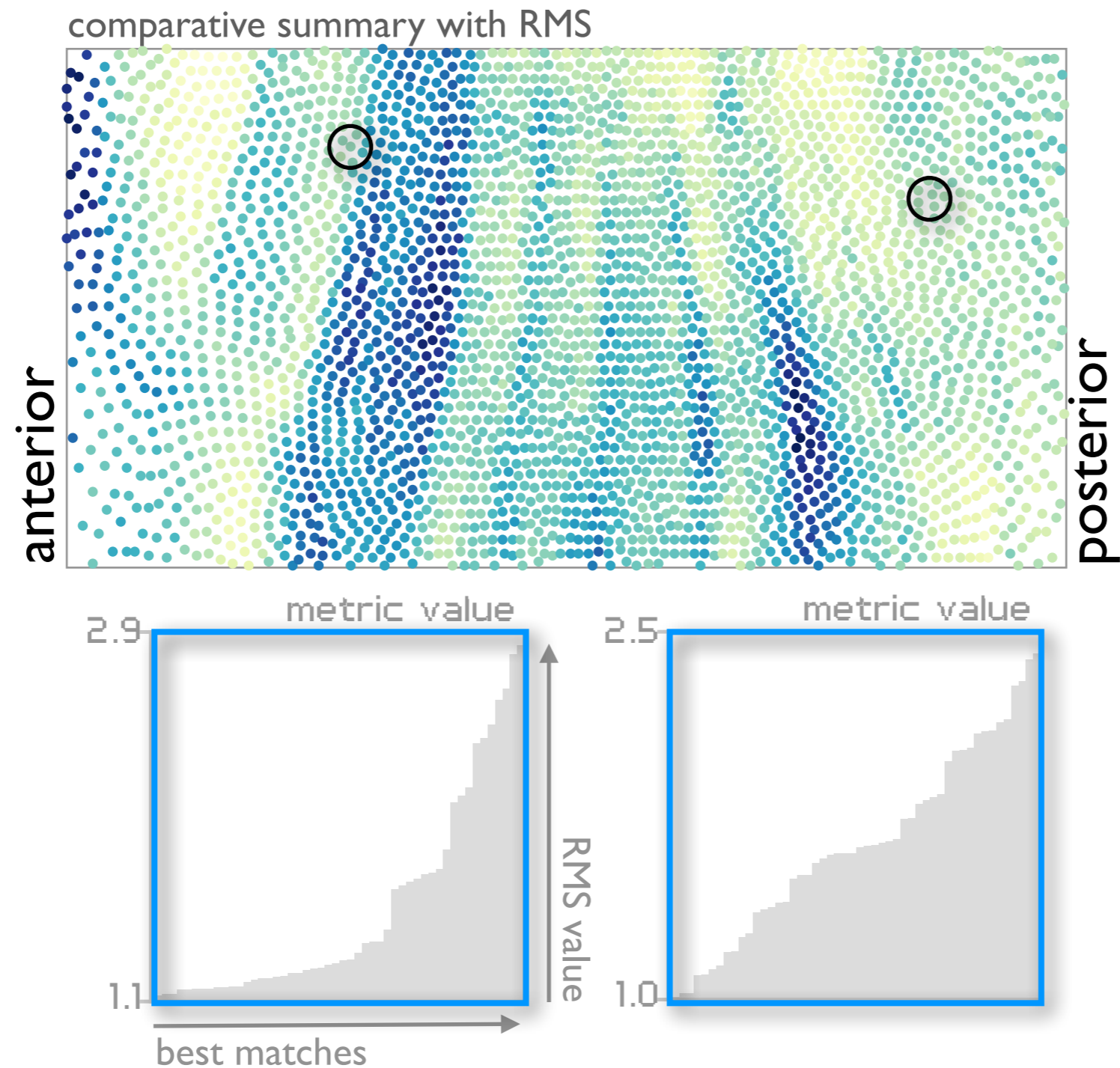


# case study one: characterizing a summary





# case study one: characterizing a summary

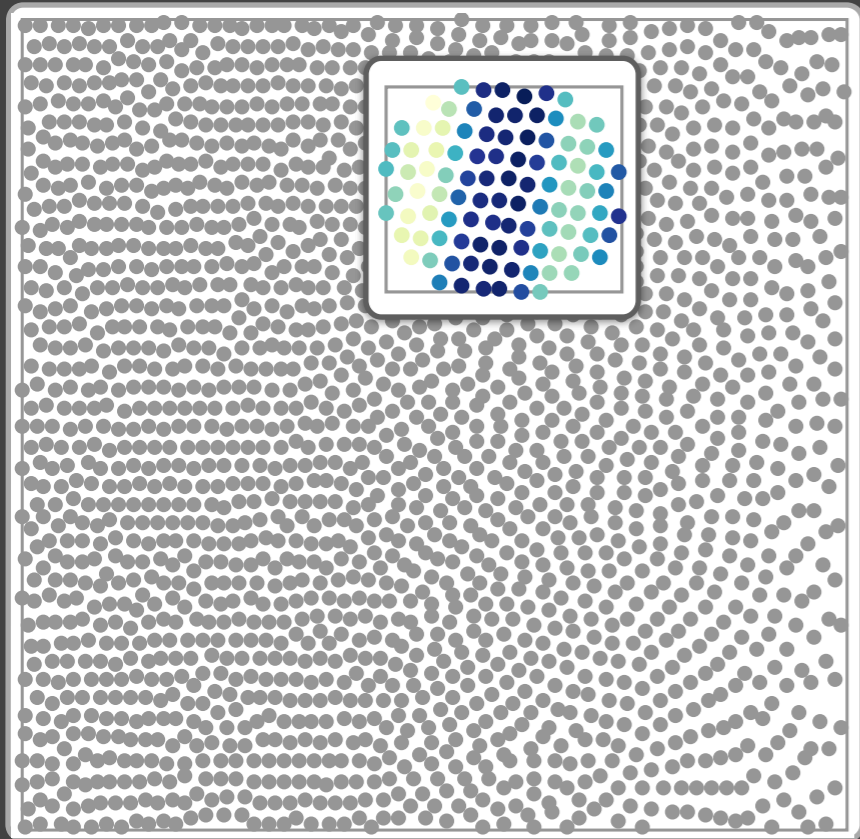


**group**  
a set of cells



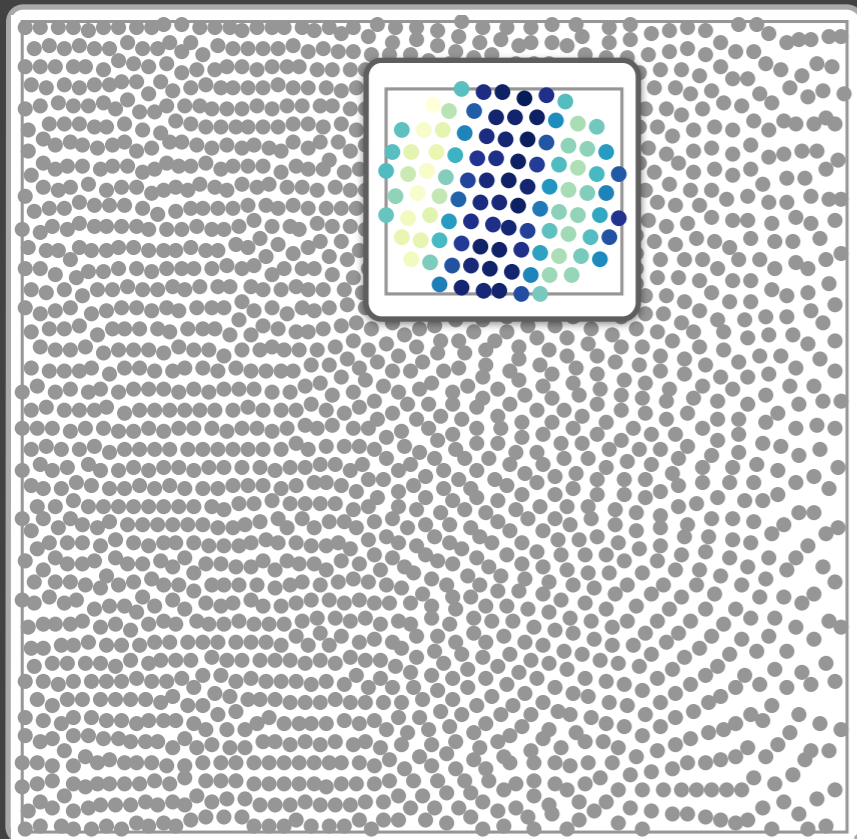
**group**  
a set of cells

**aggregation group**

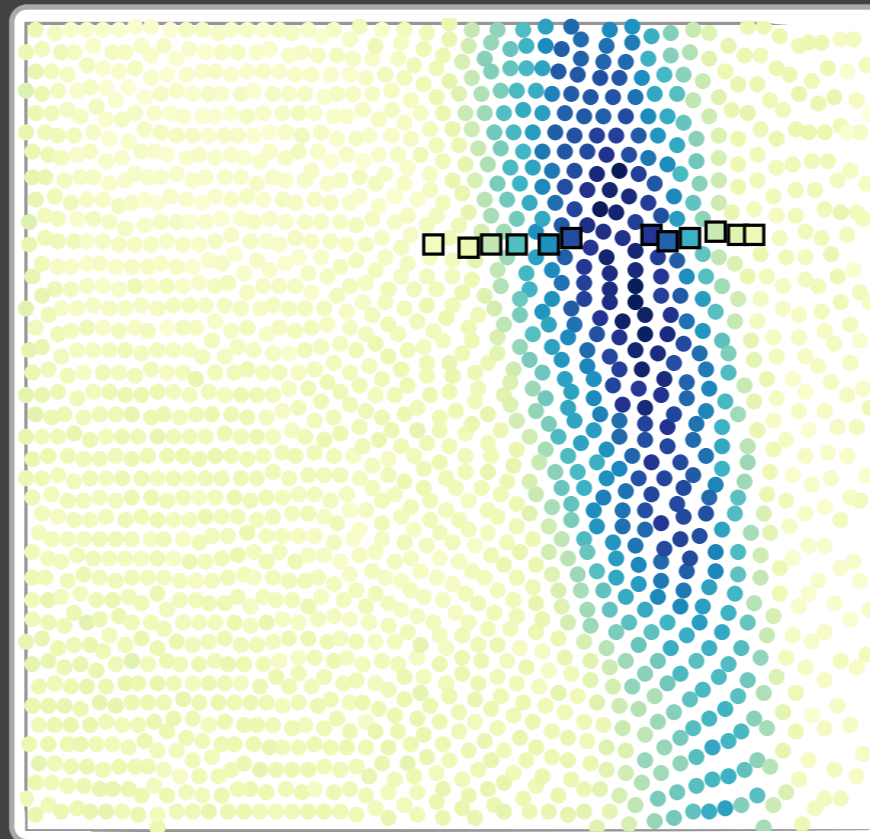


**group**  
a set of cells

**aggregation group**

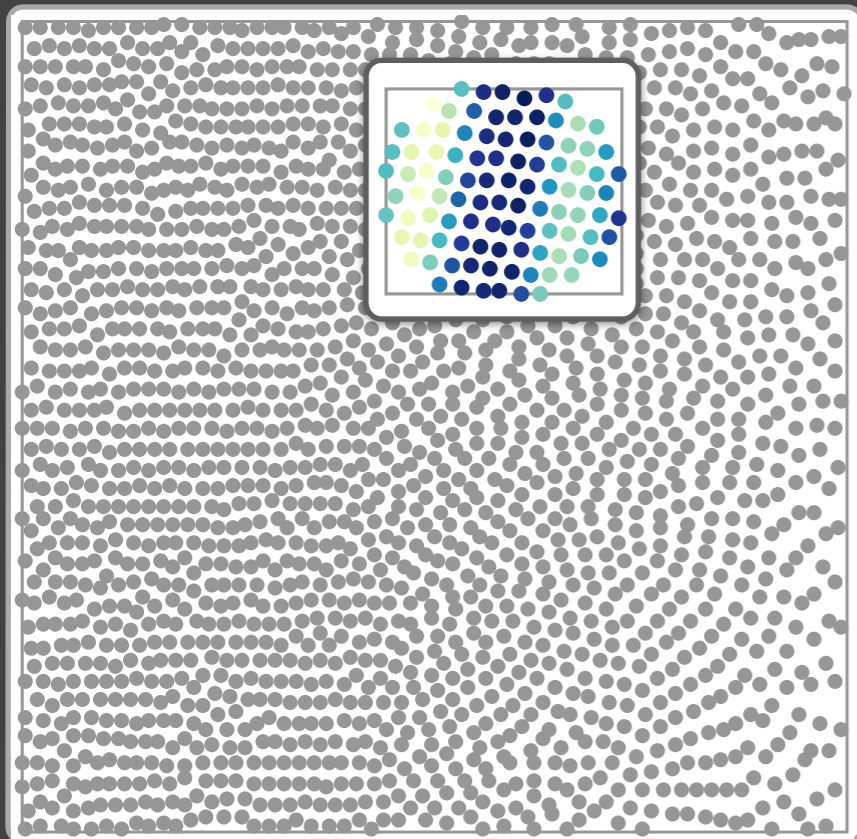


**created group**

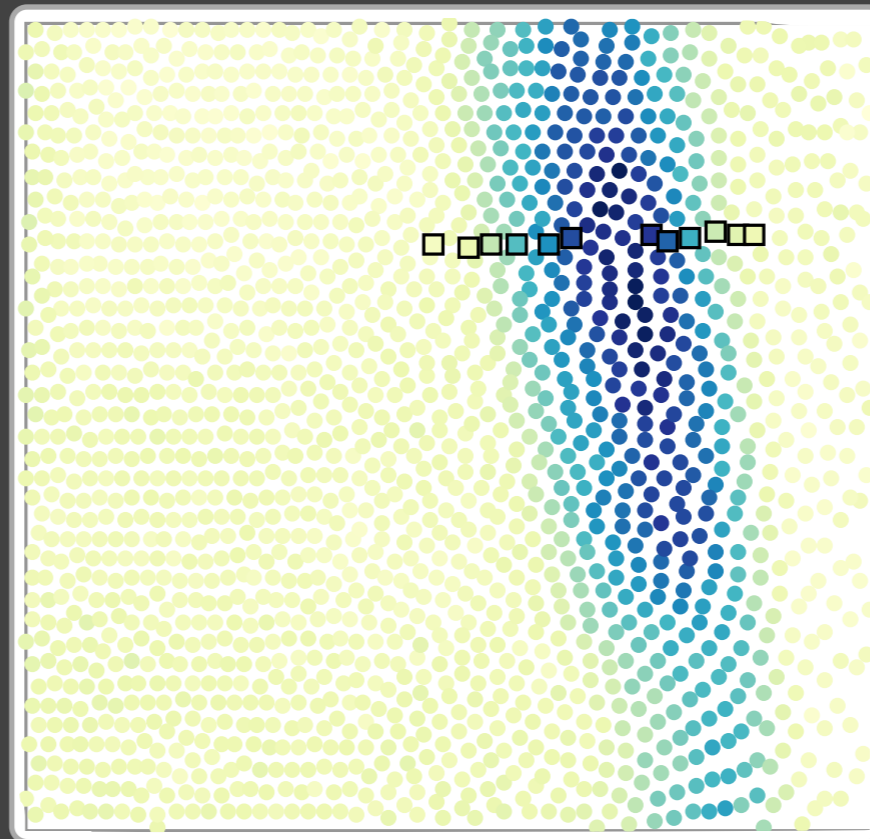


**group**  
a set of cells

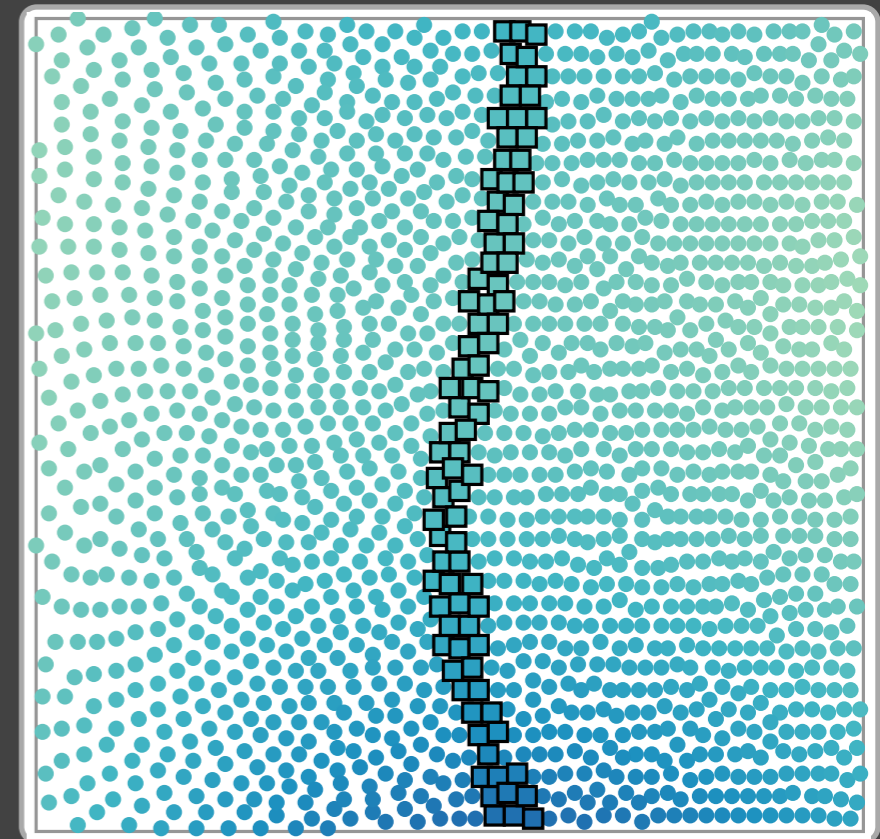
**aggregation group**



**created group**

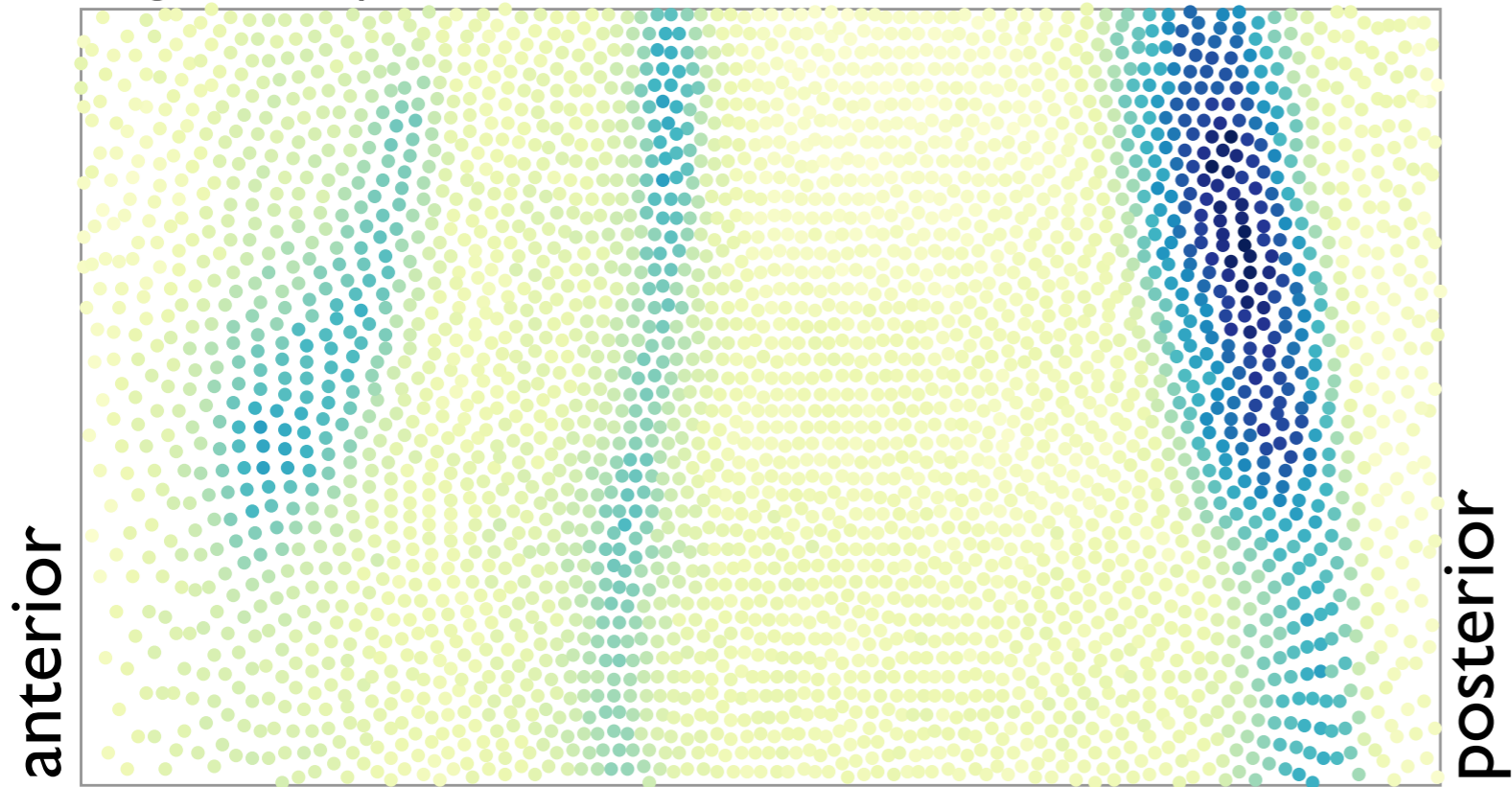


**existing group**



# case study two: exploring groups

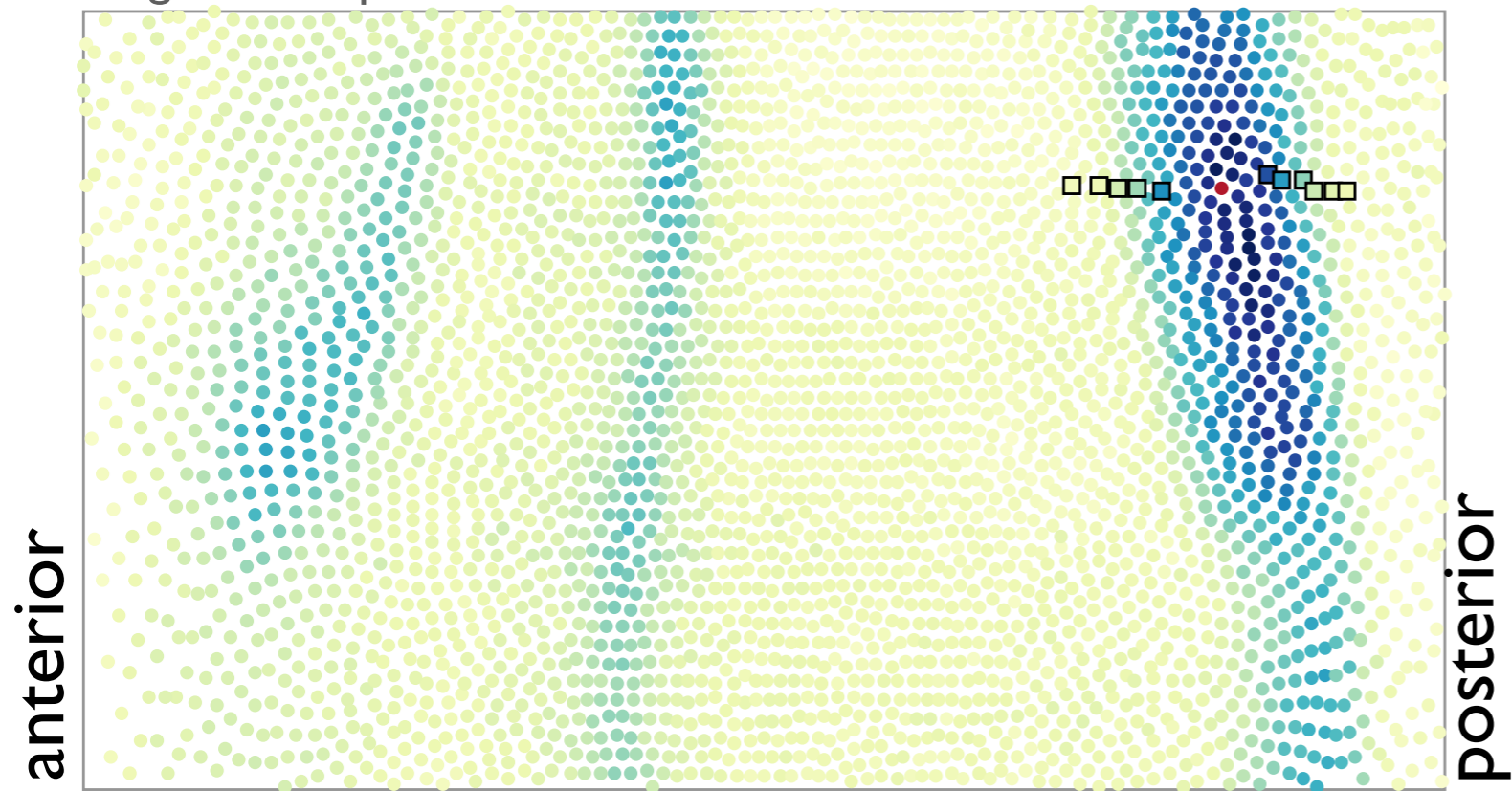
hb gene, timepoint 4





# case study two: exploring groups

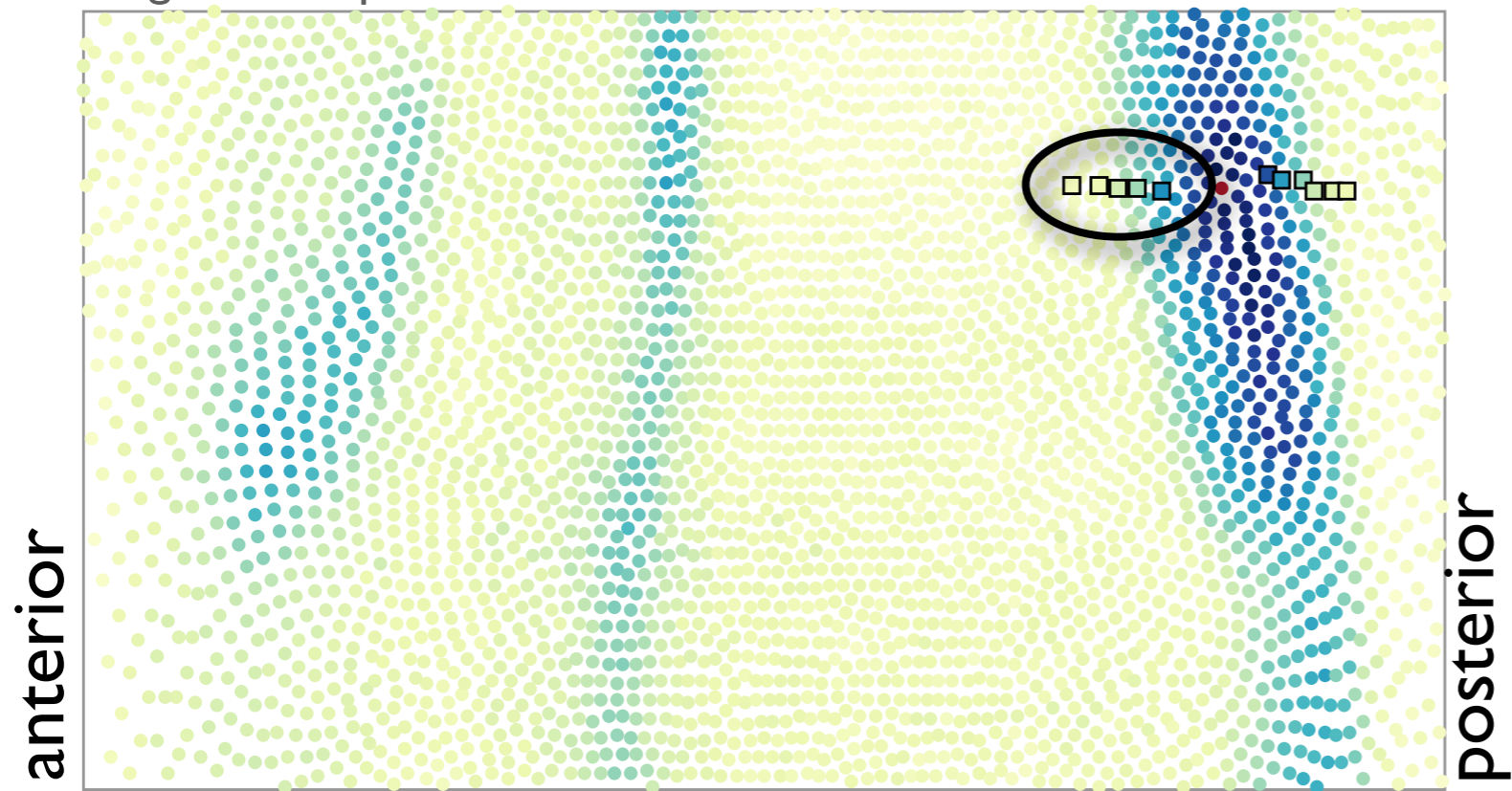
hb gene, timepoint 4



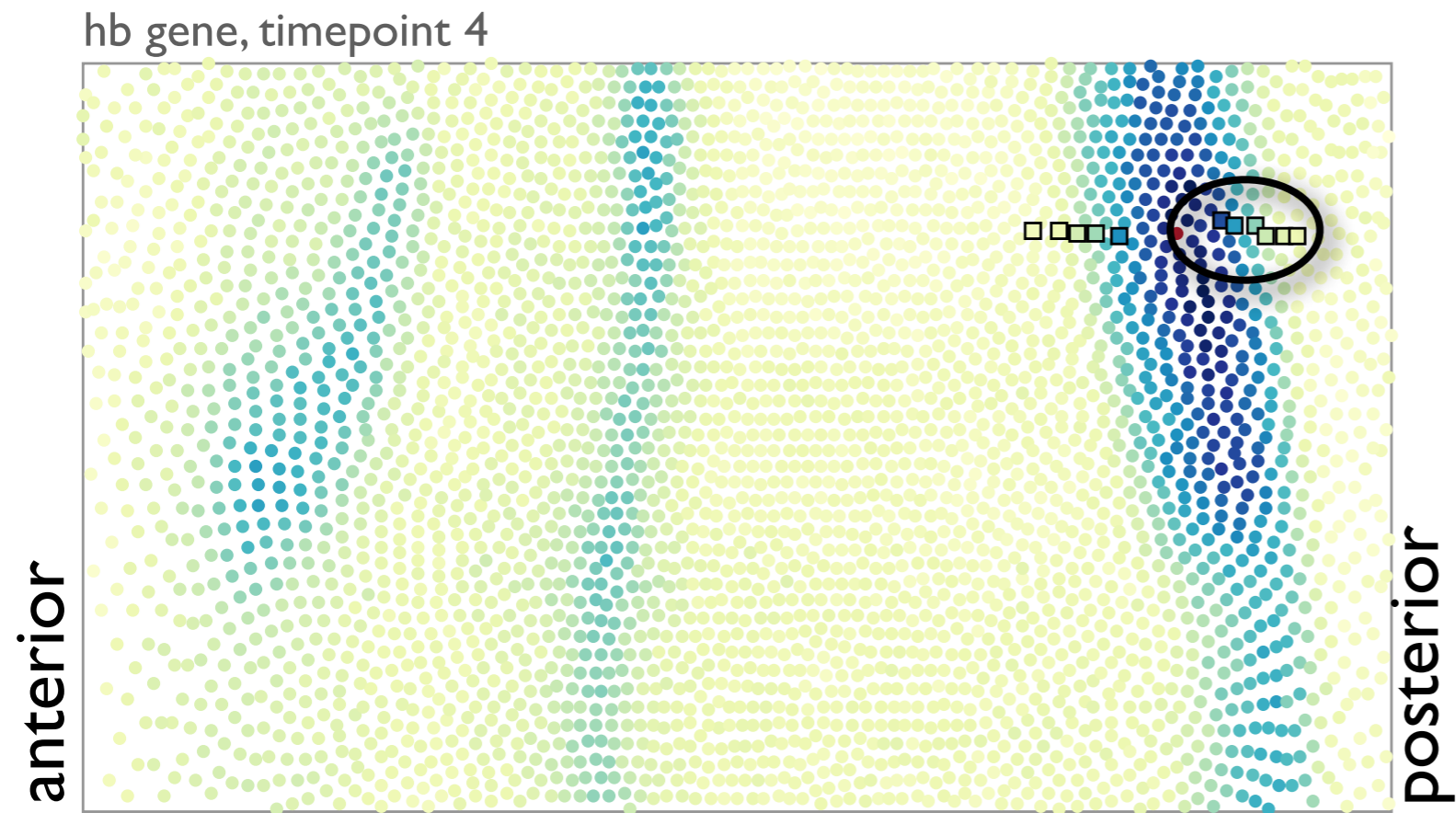


# case study two: exploring groups

hb gene, timepoint 4

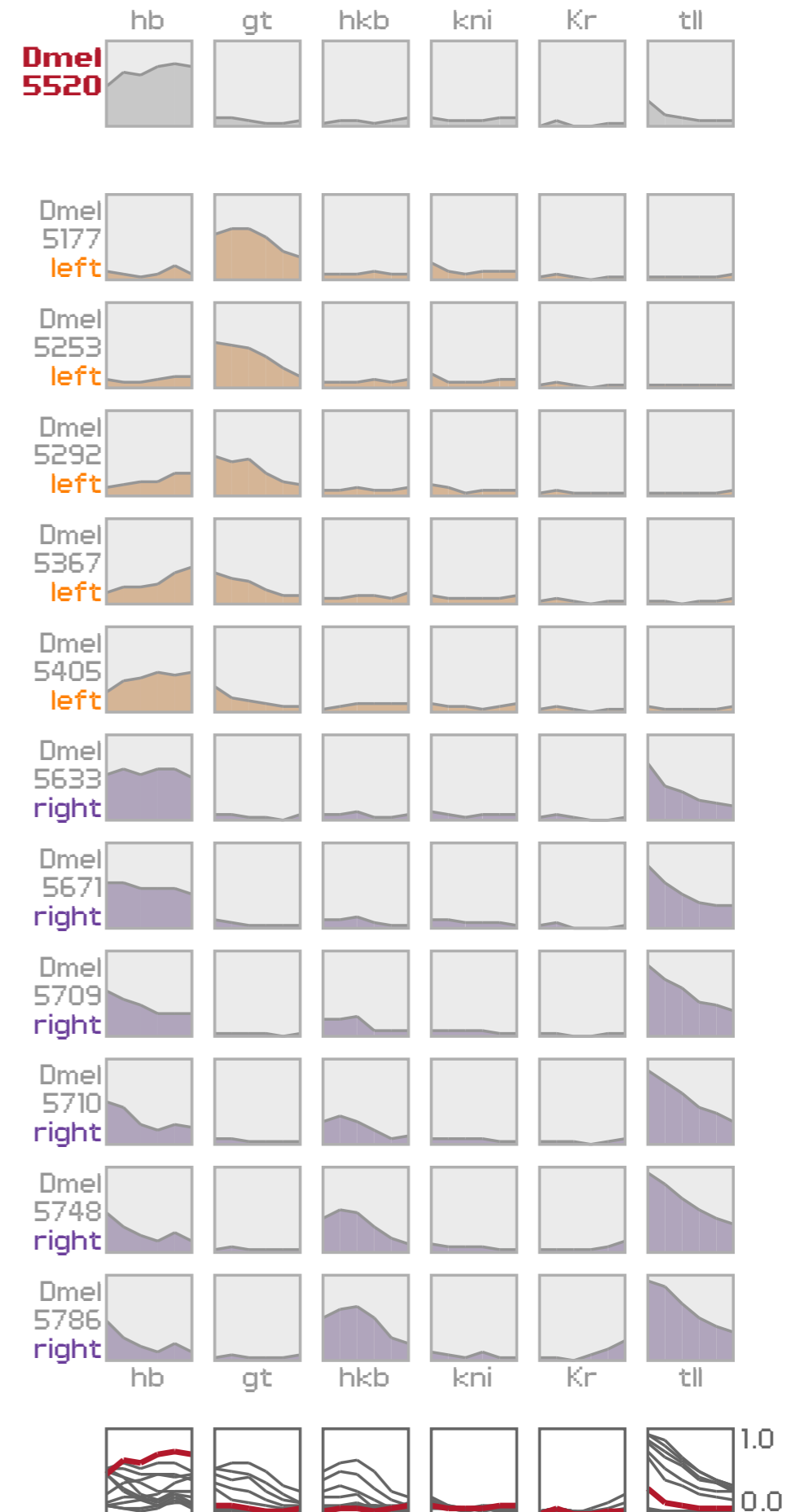
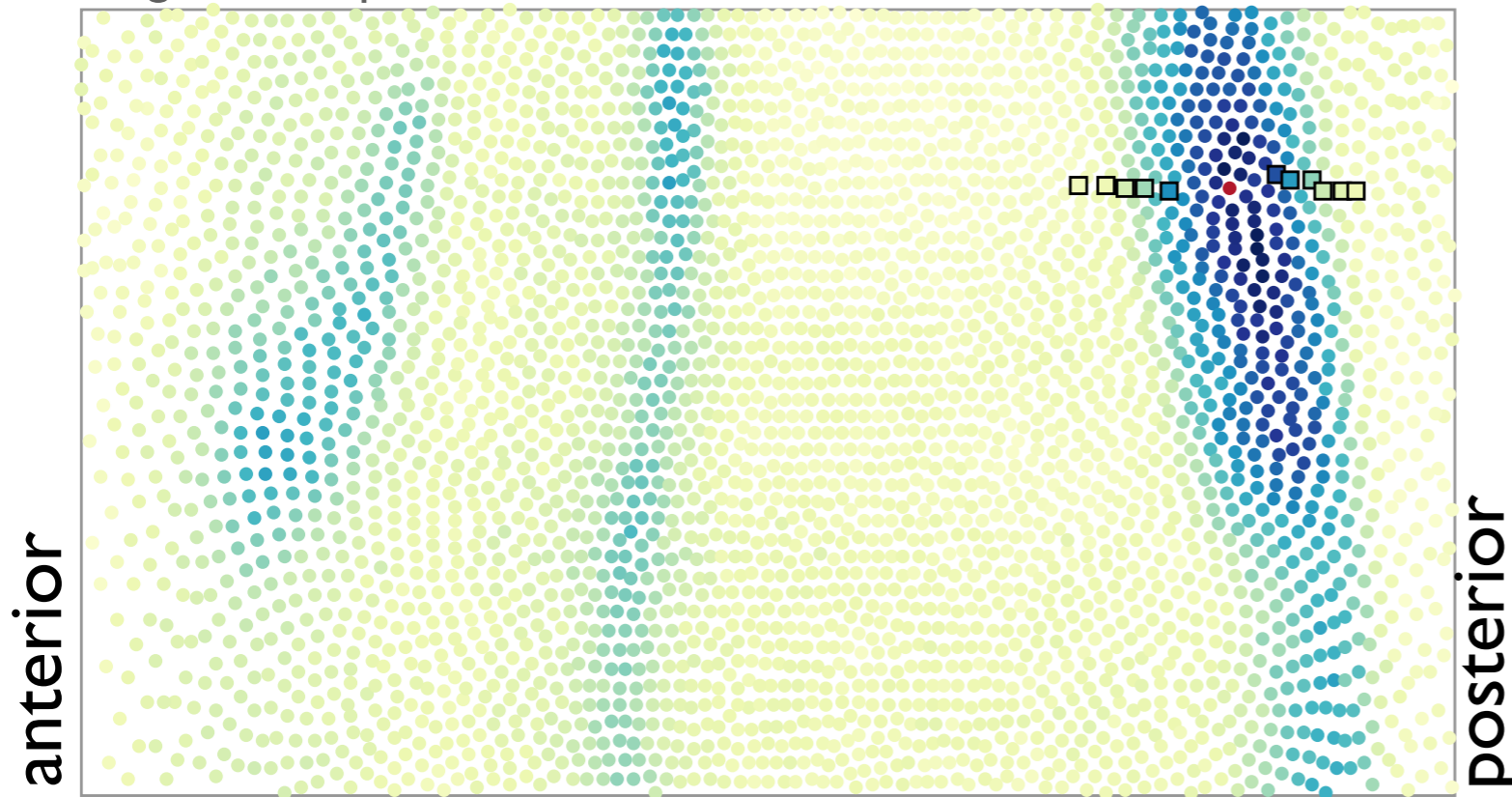


# case study two: exploring groups



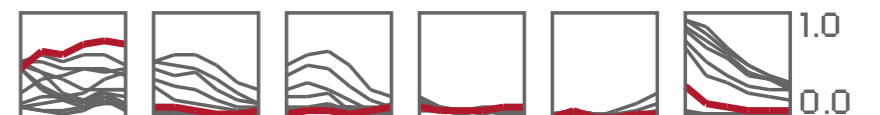
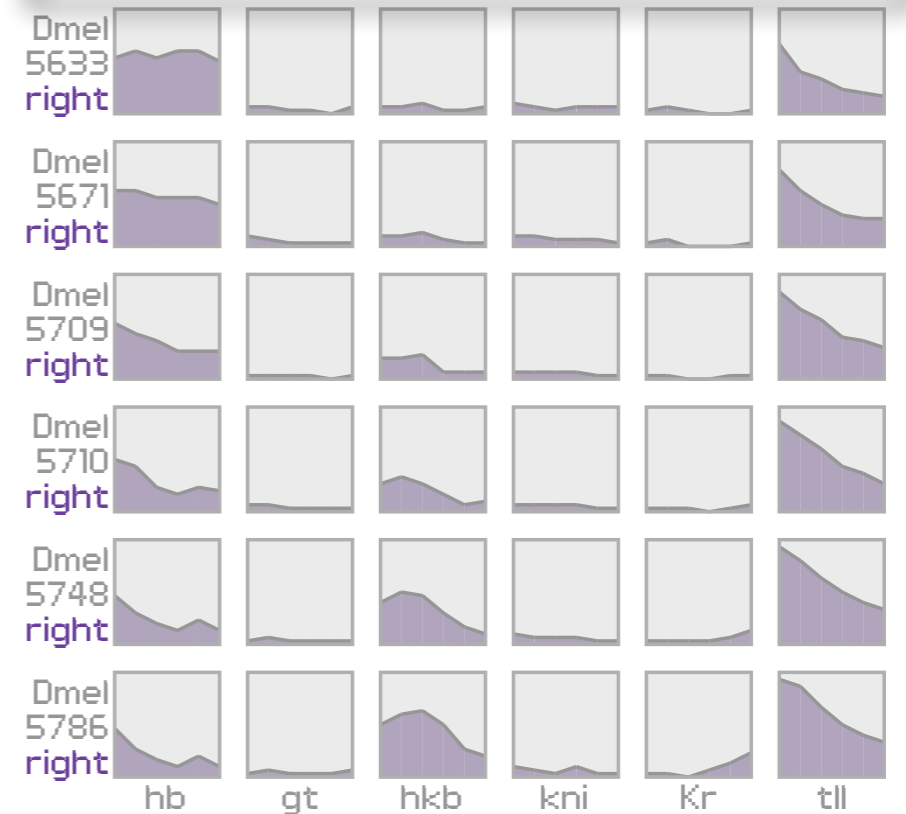
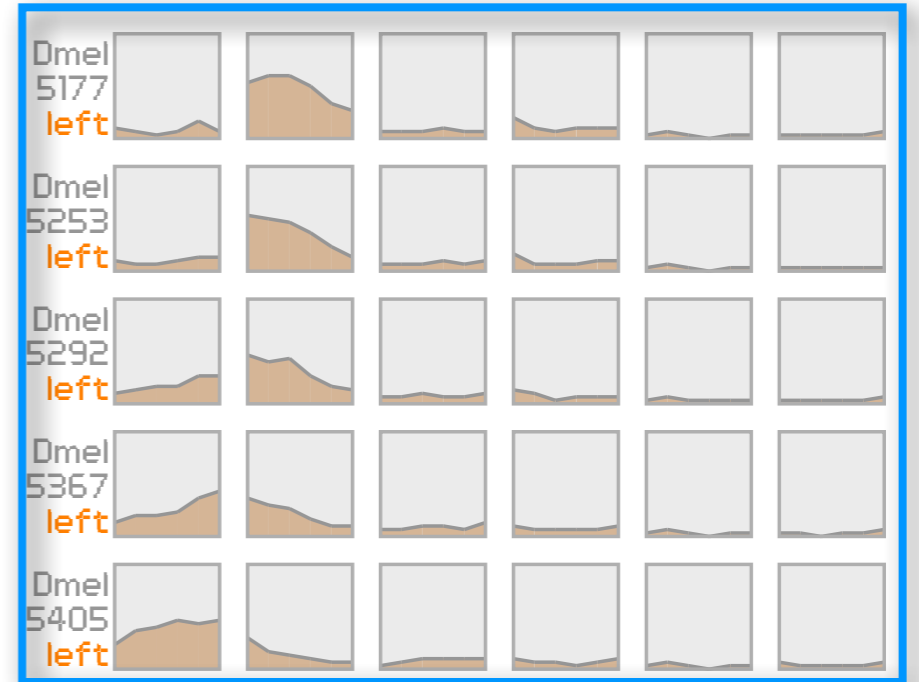
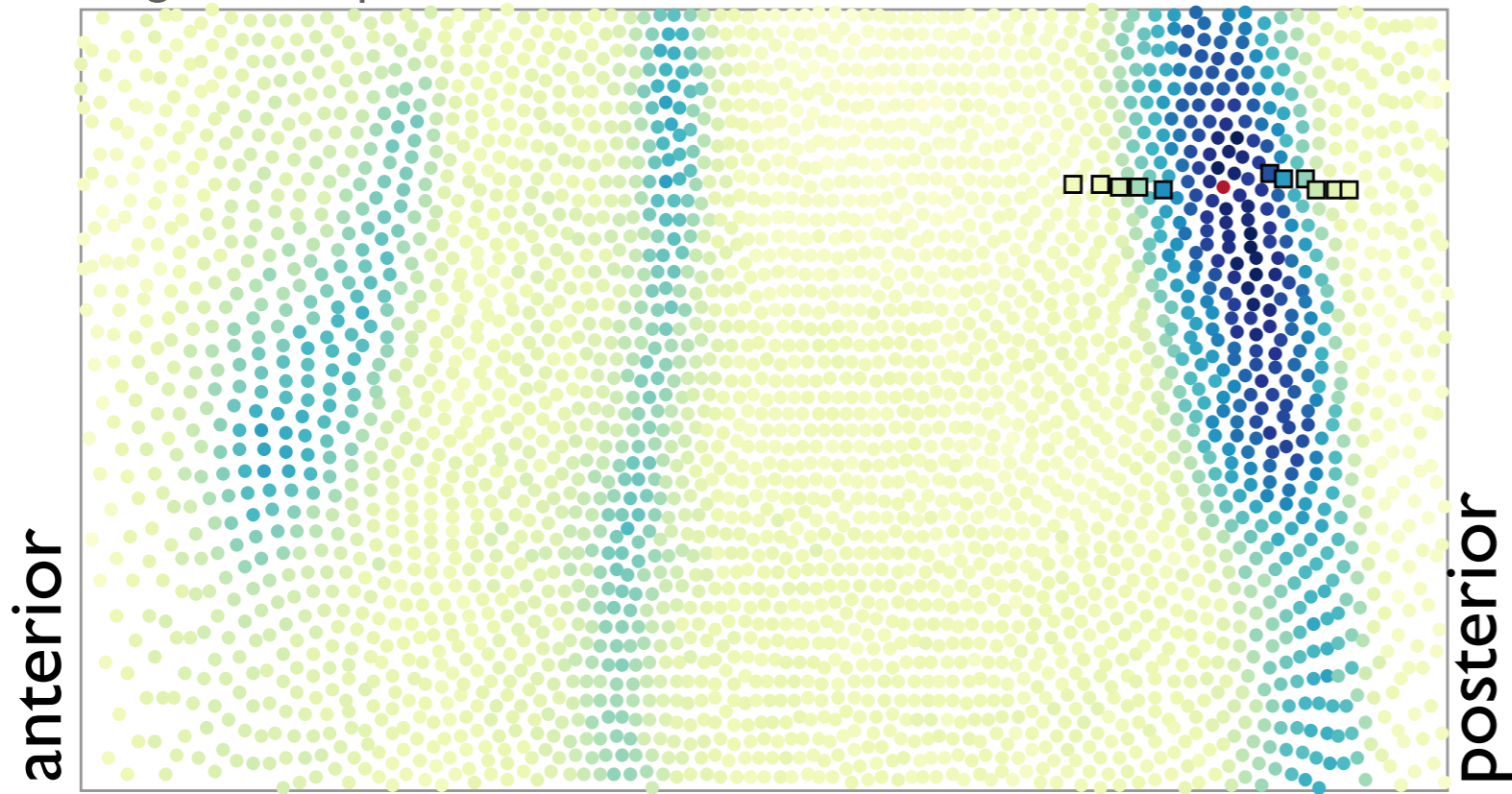
# case study two: exploring groups

hb gene, timepoint 4



# case study two: exploring groups

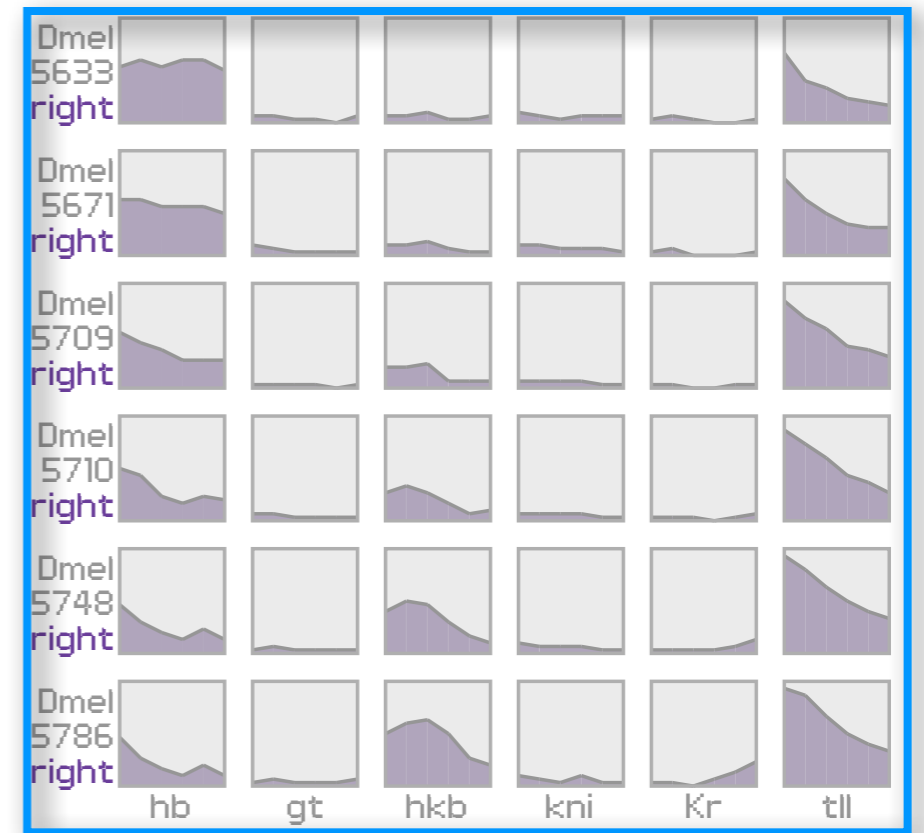
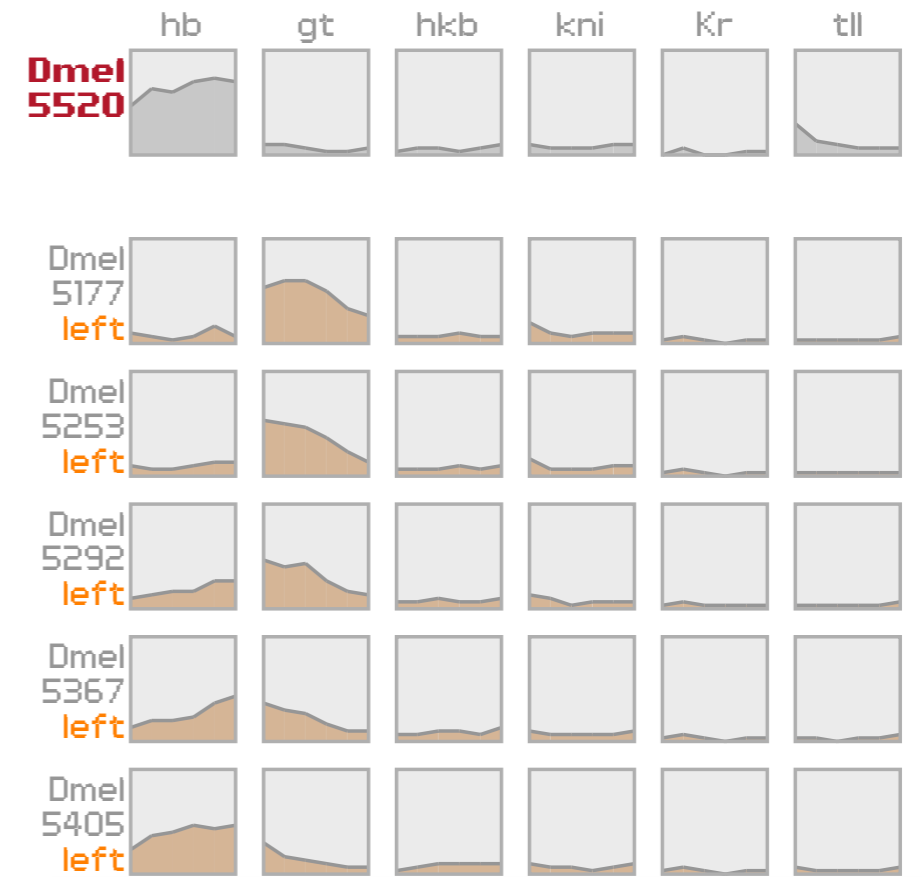
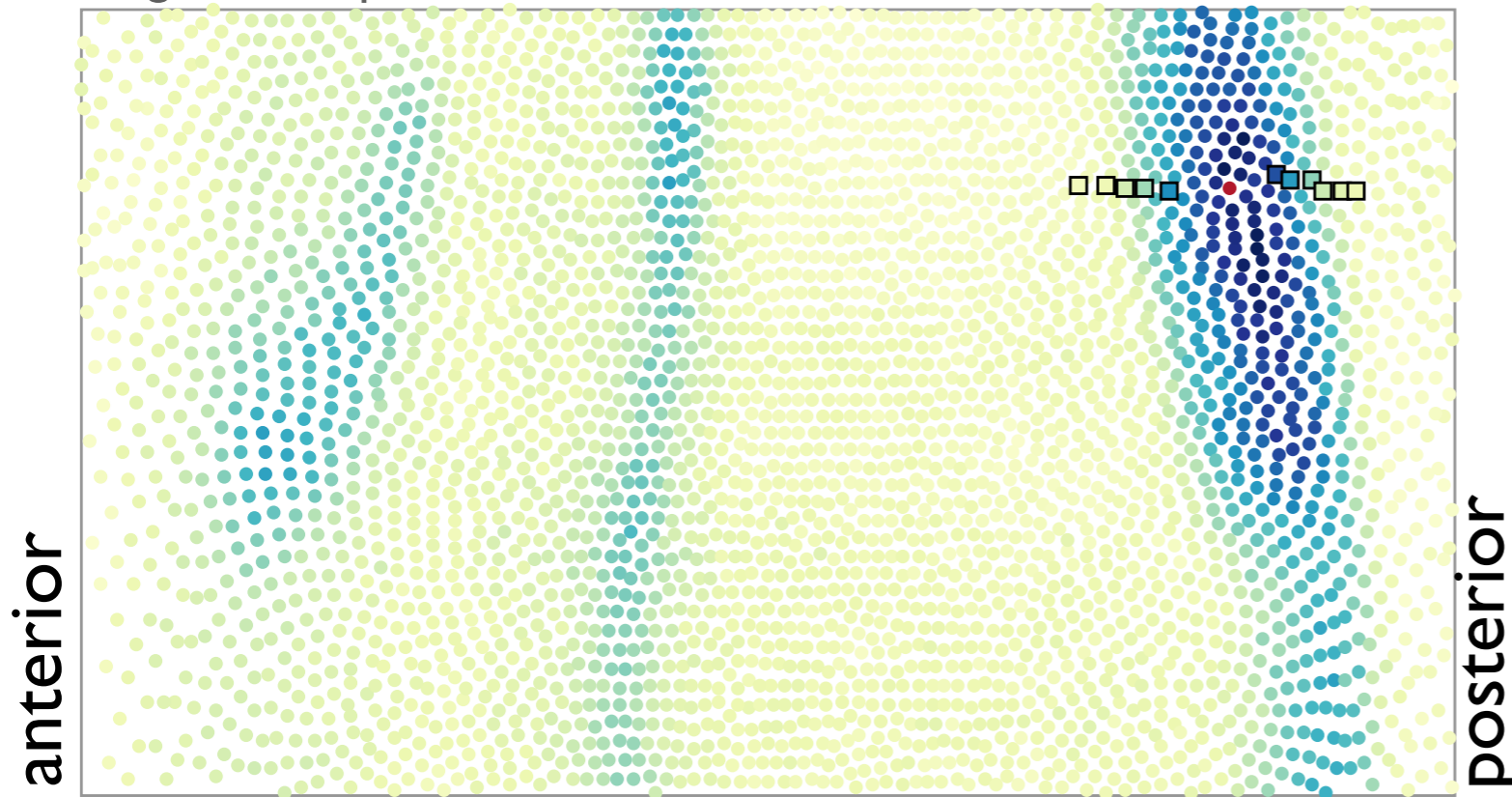
hb gene, timepoint 4





# case study two: exploring groups

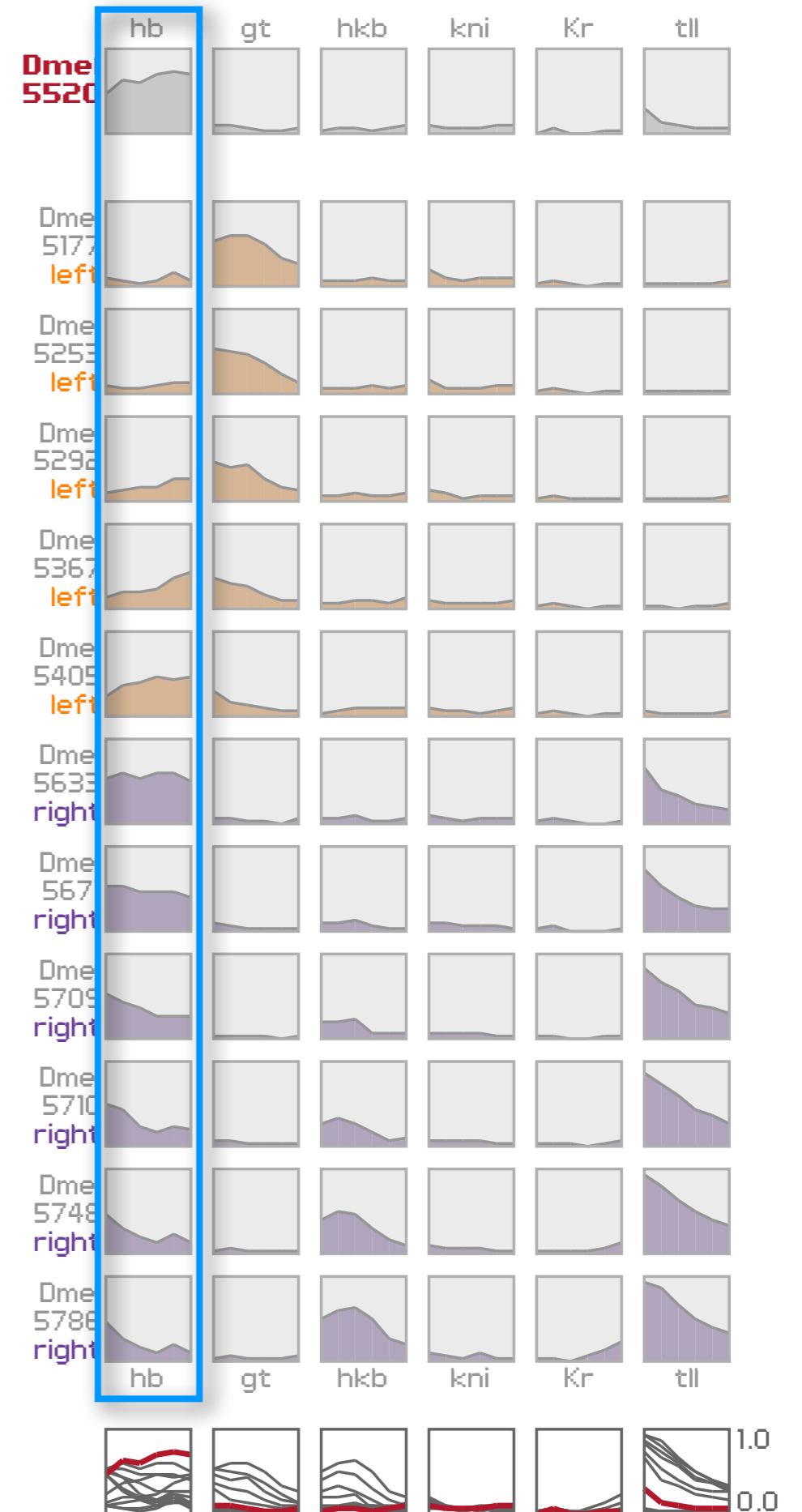
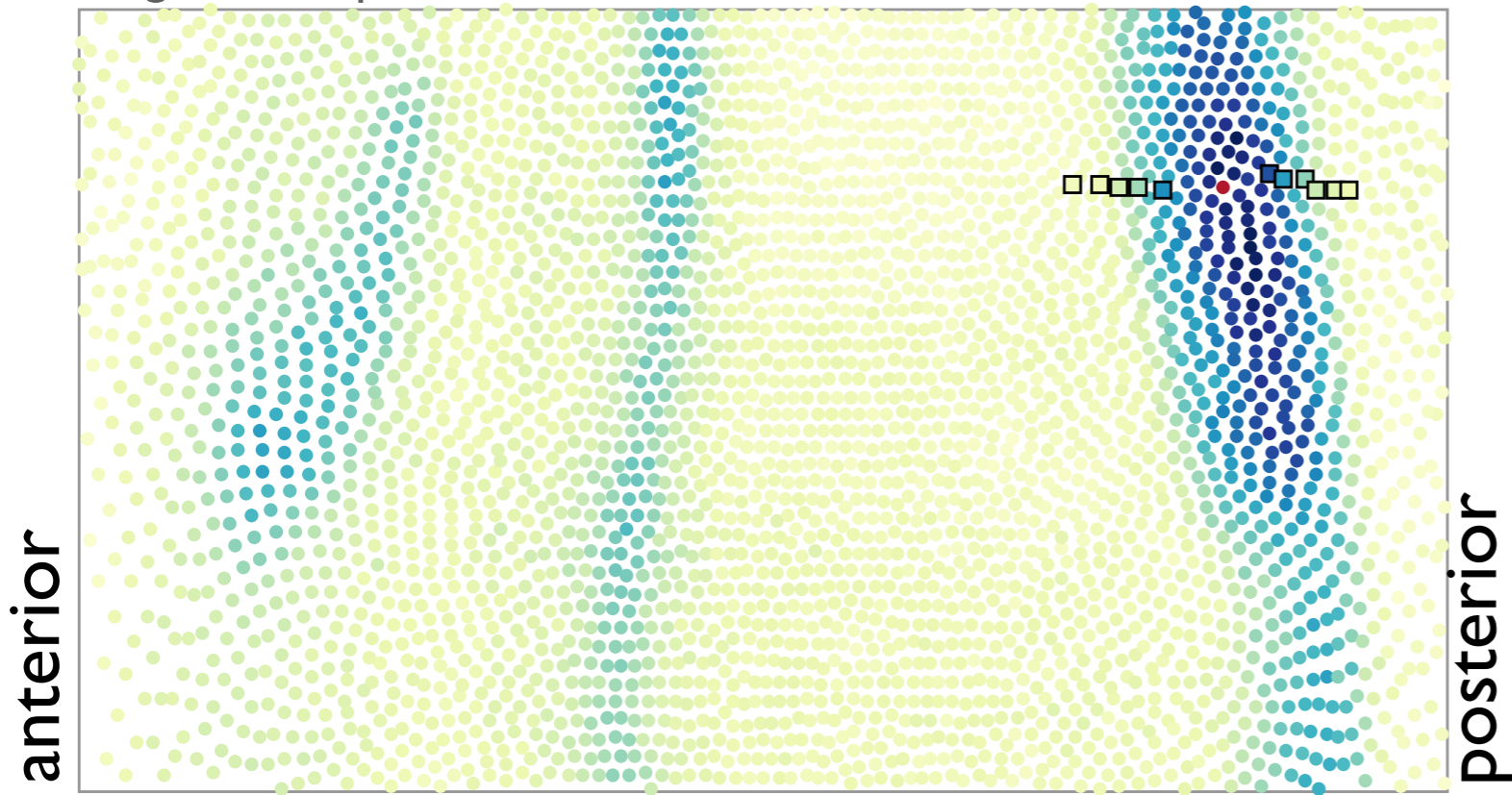
hb gene, timepoint 4





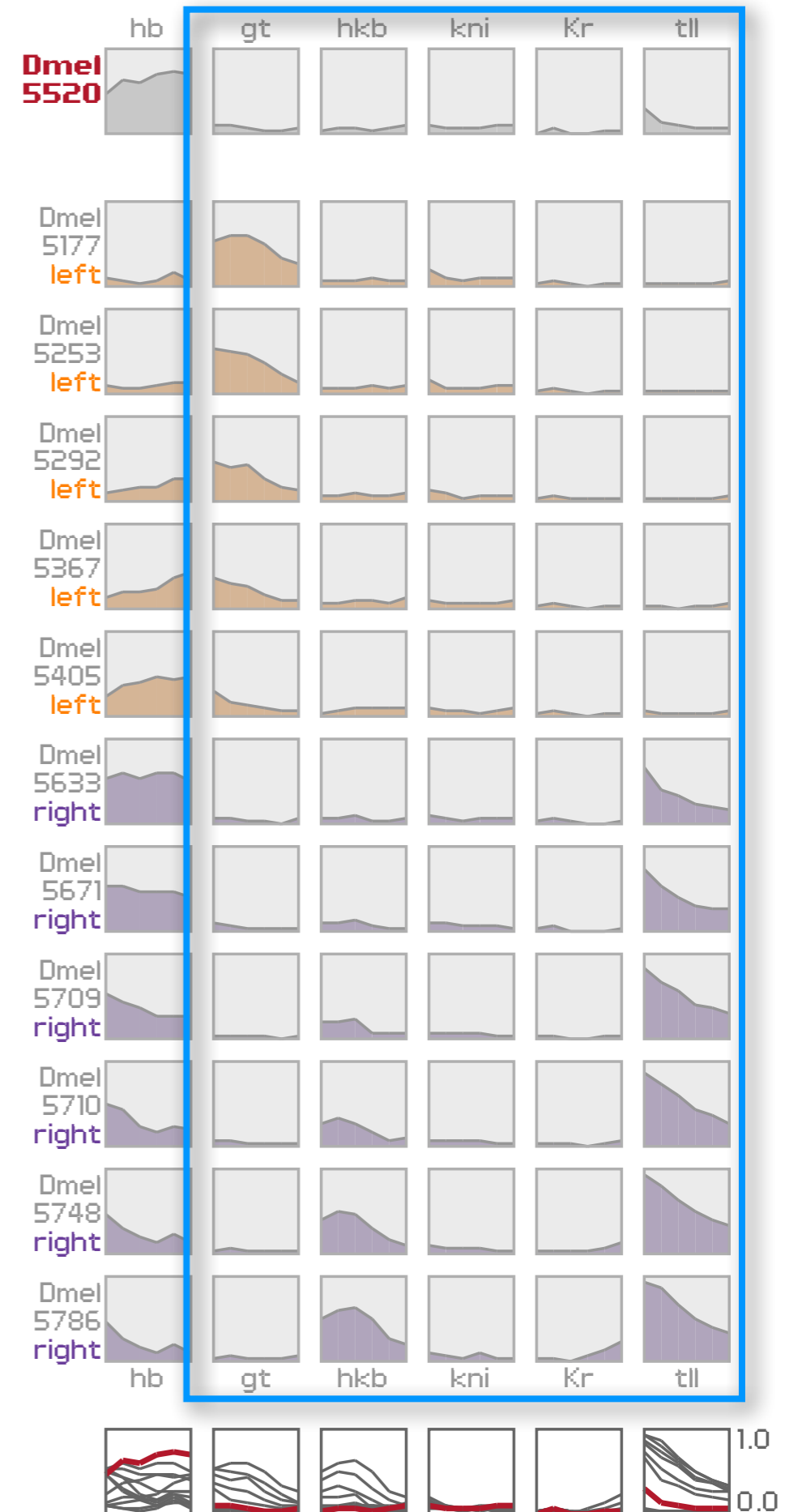
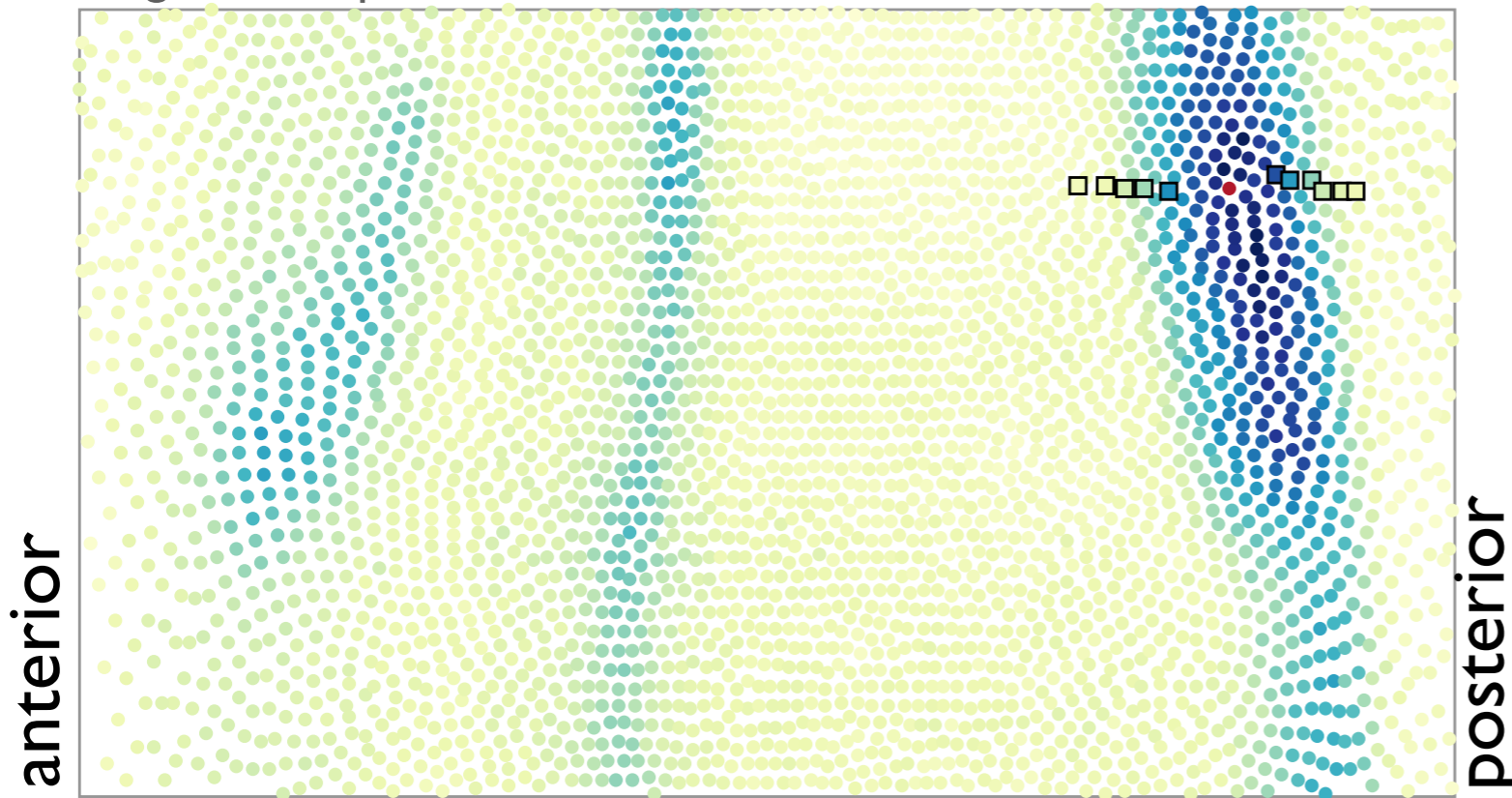
# case study two: exploring groups

hb gene, timepoint 4



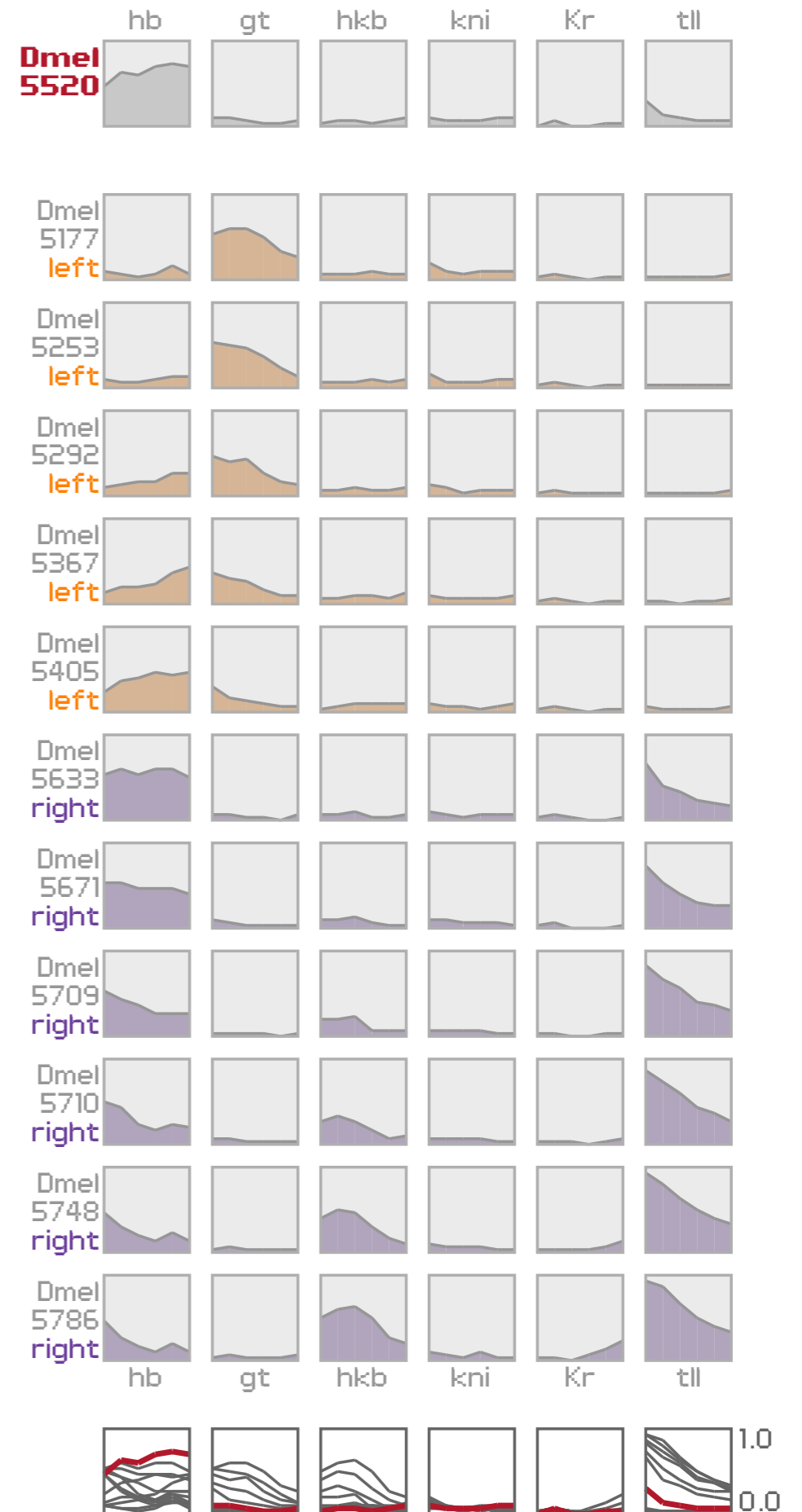
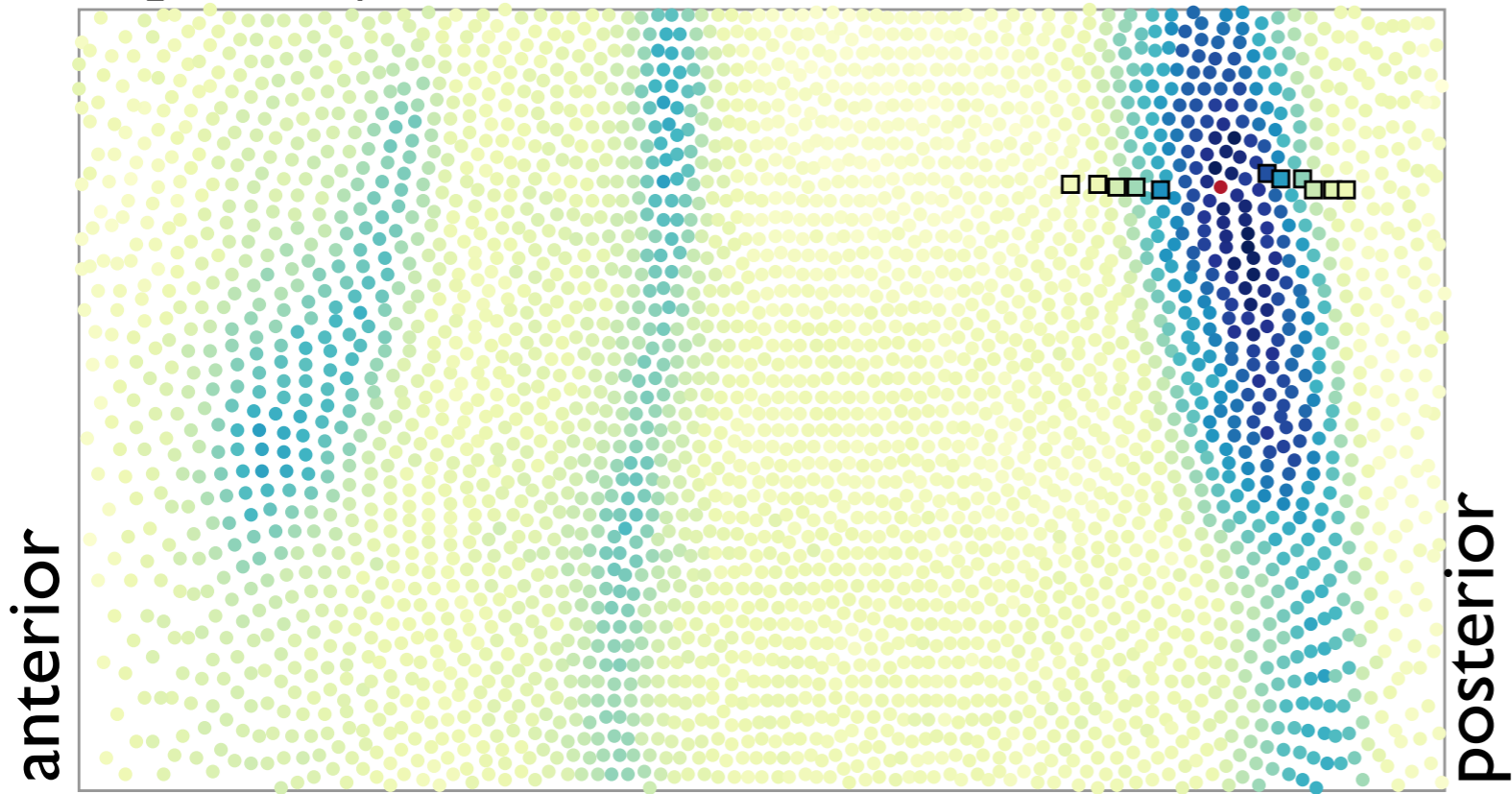
# case study two: exploring groups

hb gene, timepoint 4



# case study two: exploring groups

hb gene, timepoint 4





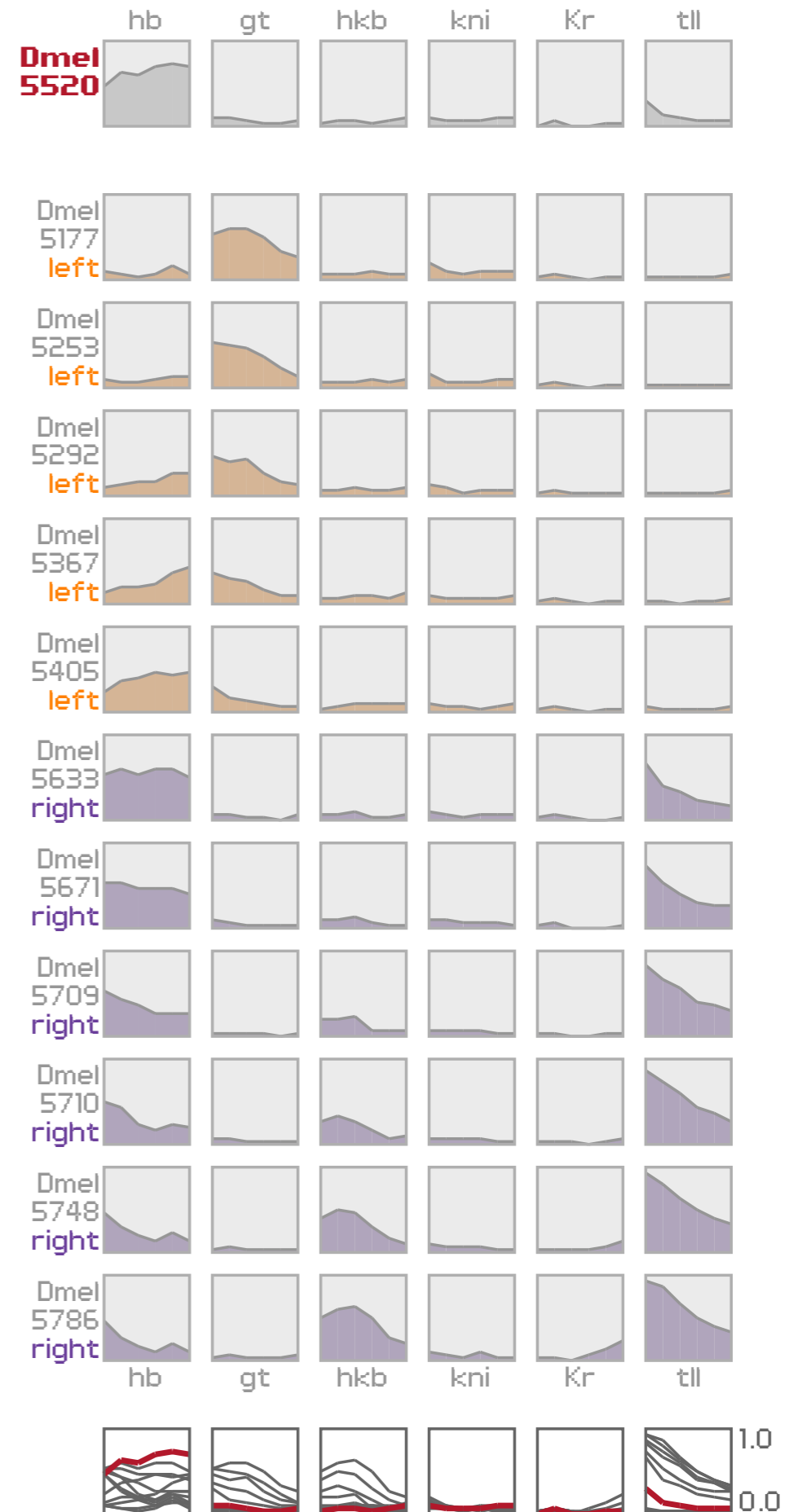
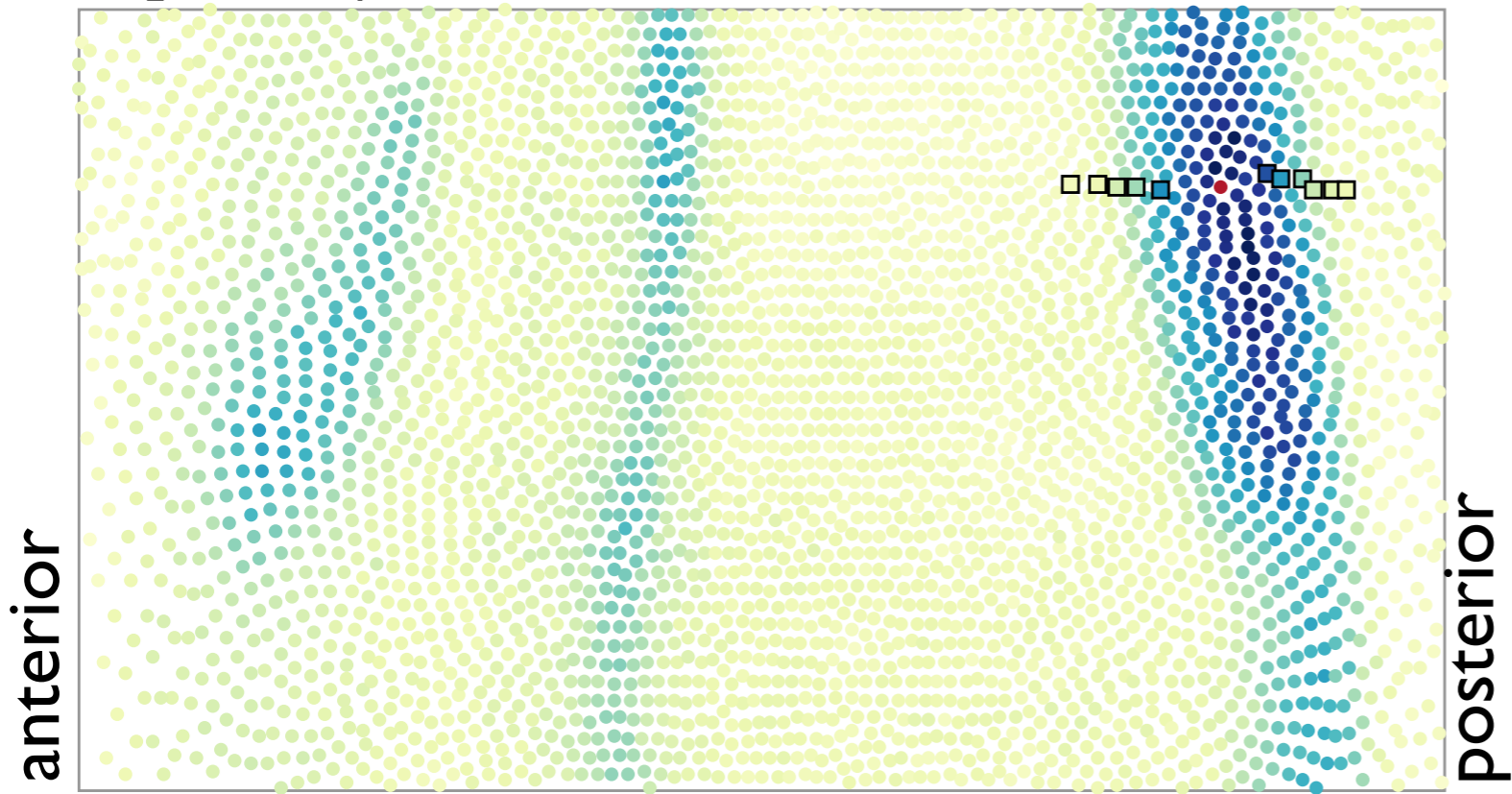






# case study two: exploring groups

hb gene, timepoint 4



data & tool & tasks

summaries & groups

**encodings & interaction**

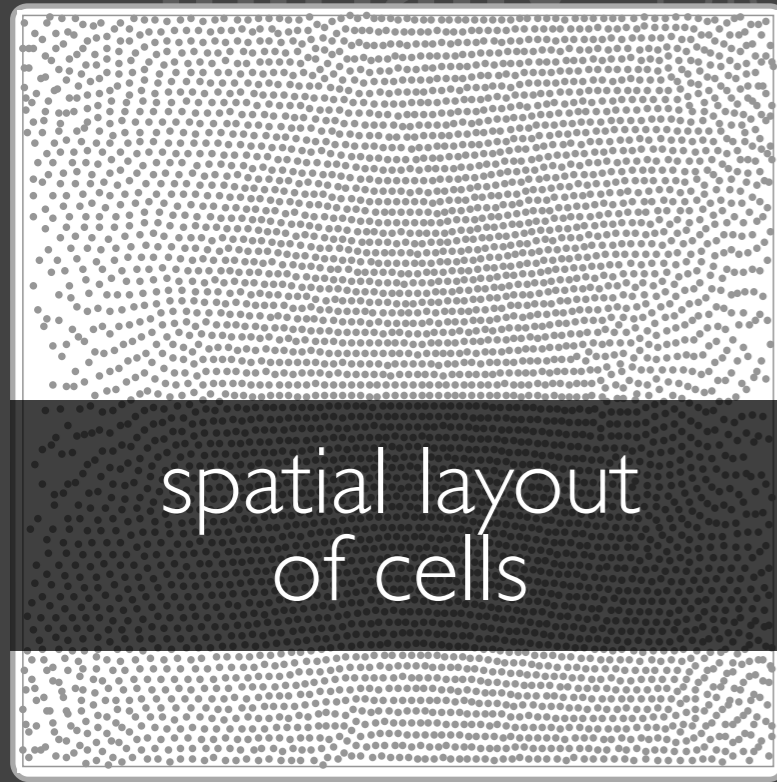
conclusions

**initially:** biologists relied on computational methods to explore data

**goal:** link underlying data with computational results

**visualize:** triad of data

**initially:** biologists relied on  
manual methods to explore data

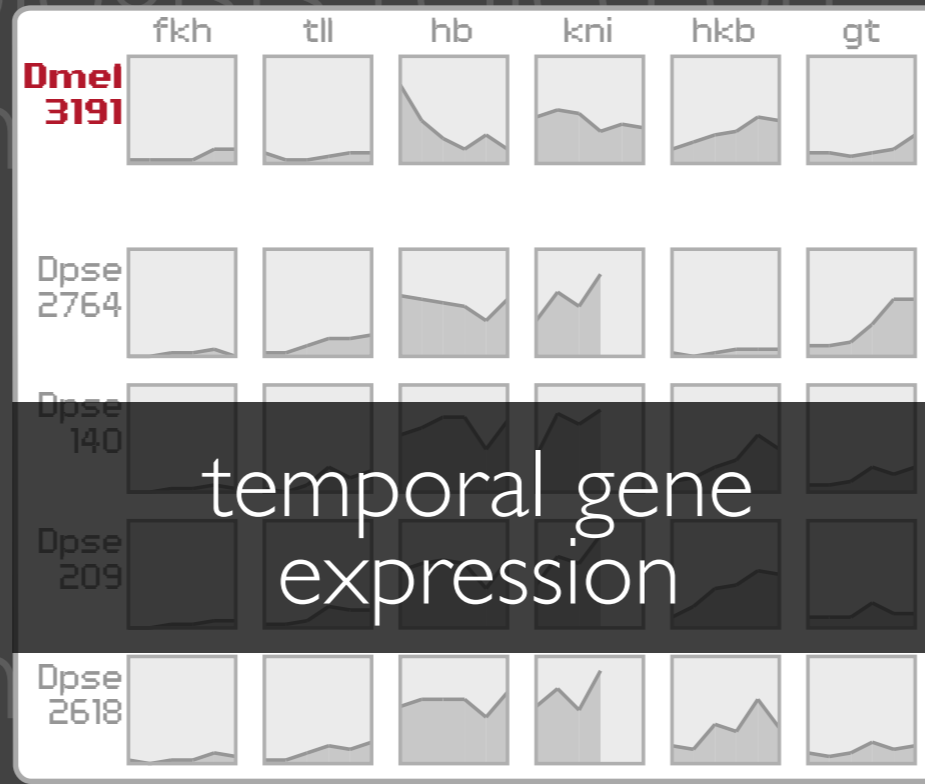
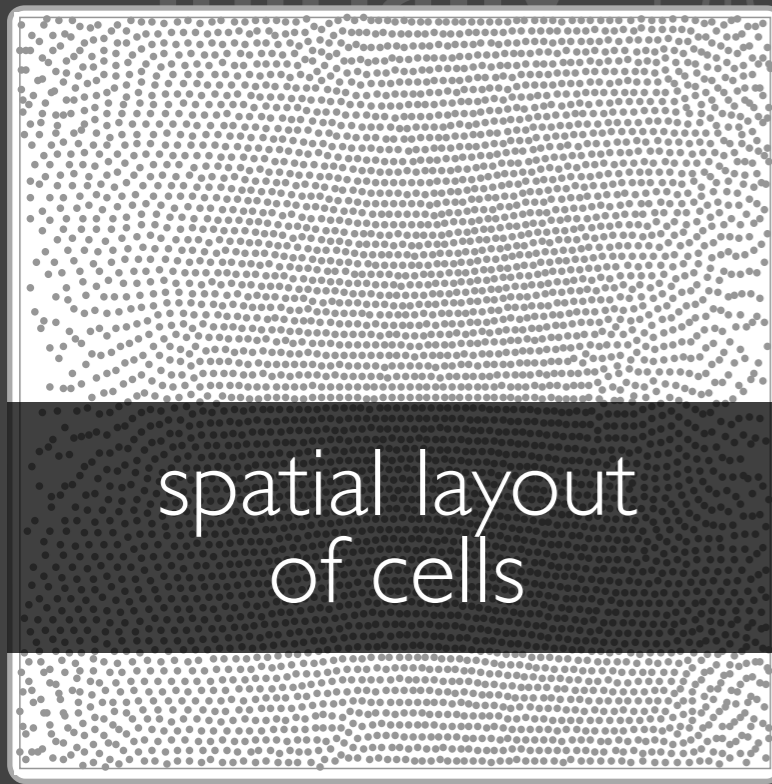


underlying data with  
manual results

embryo map

**visualize:** triad of data

initially: biologists relied on sparse data



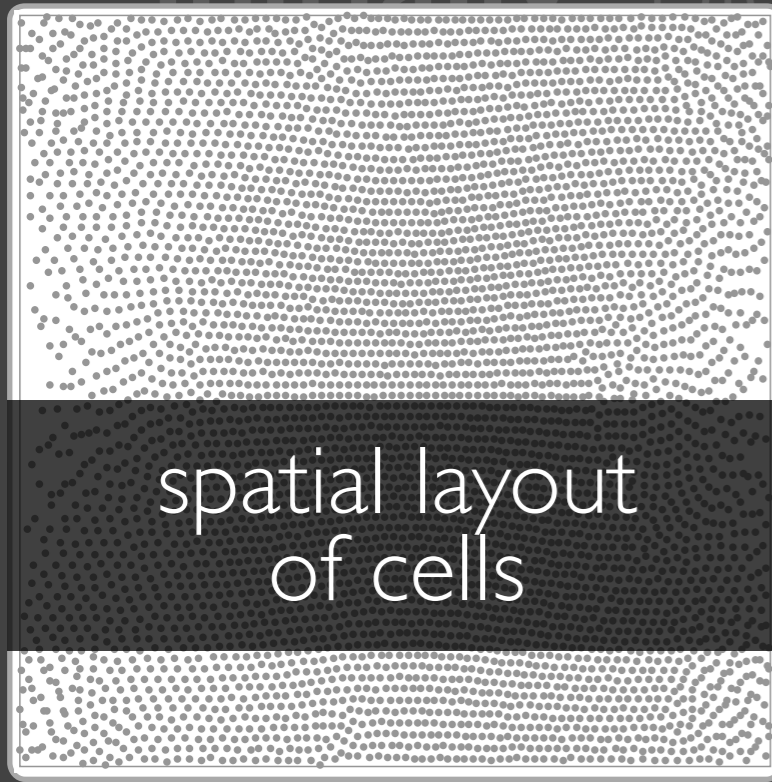
embryo map

curvemaps [Meyer 10]

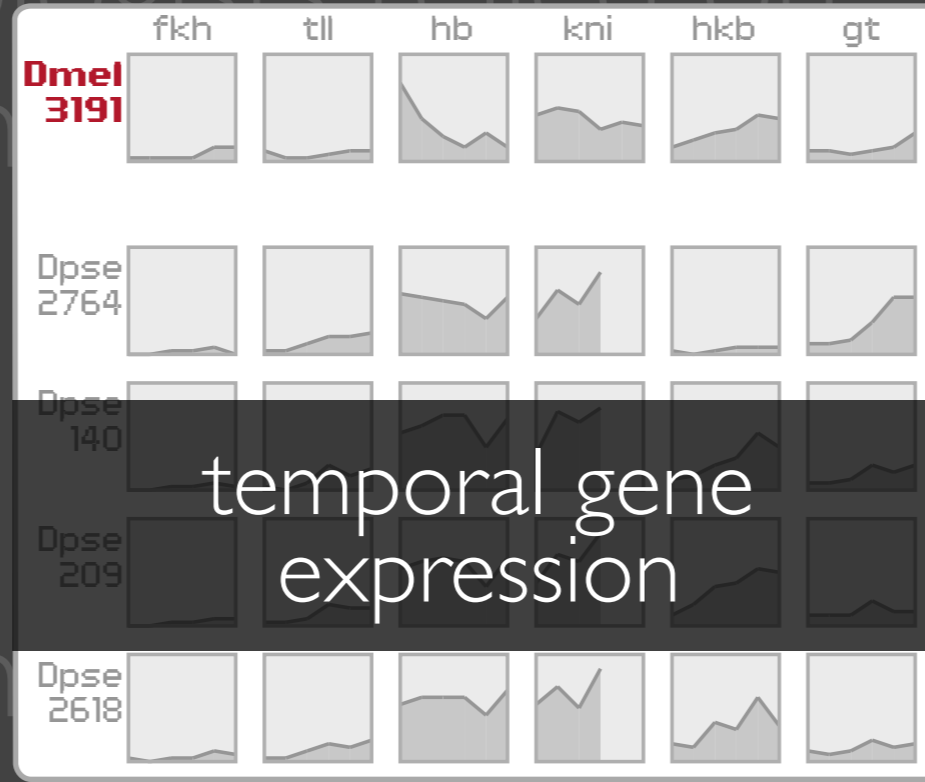
**visualize:** triad of data



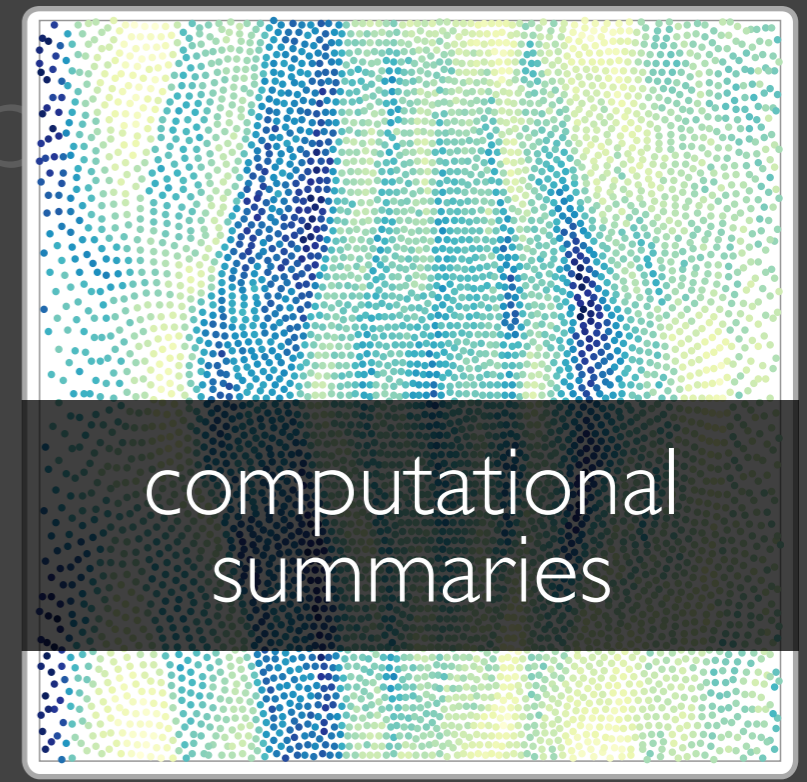
initially: biologists relied on



spatial layout  
of cells



temporal gene  
expression



computational  
summaries

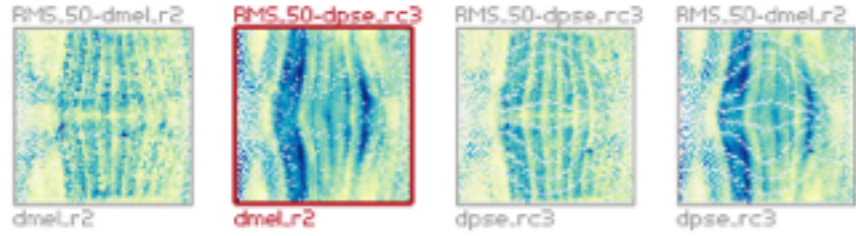
embryo map

curvemaps [Meyer 10]

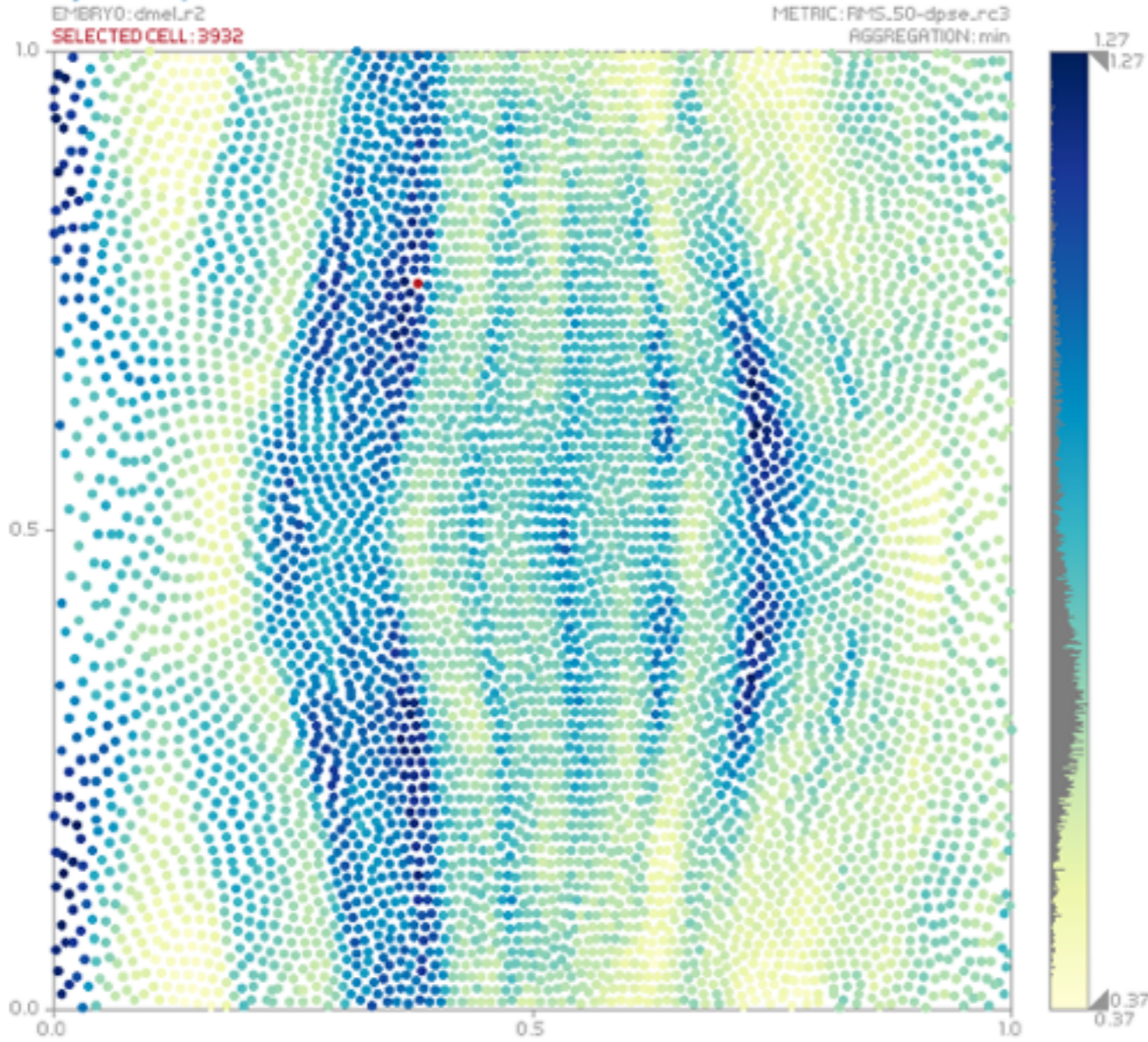
colormap

**visualize:** triad of data

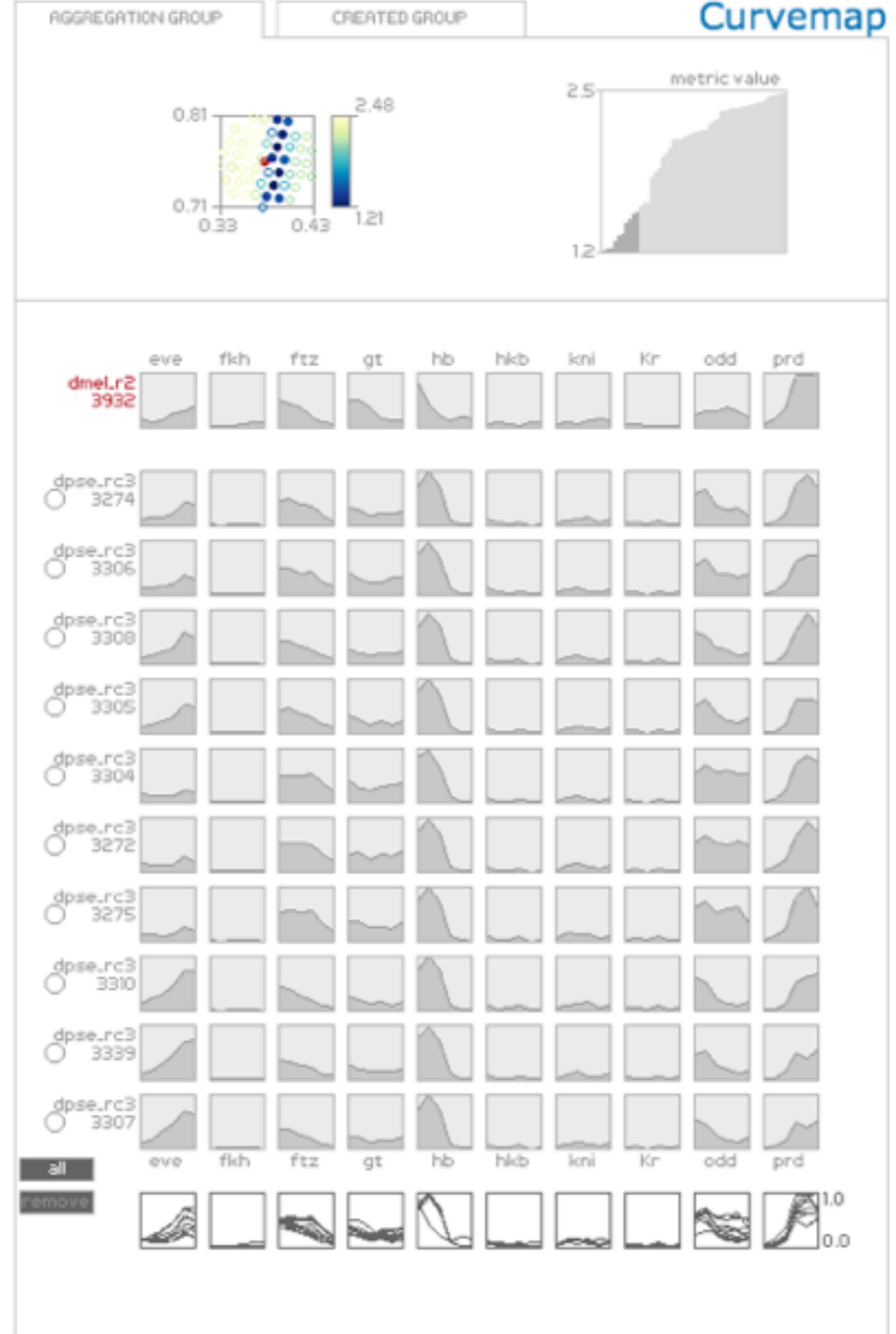
### Summaries



### Embryo Map

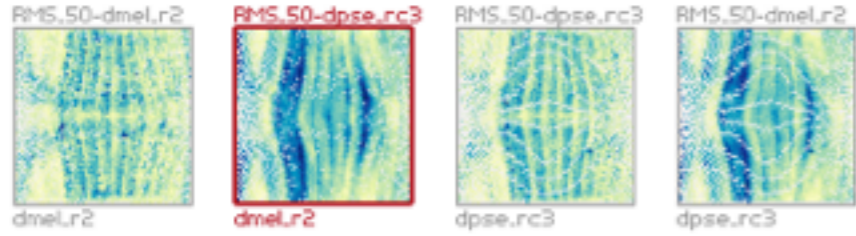


### Curvemap

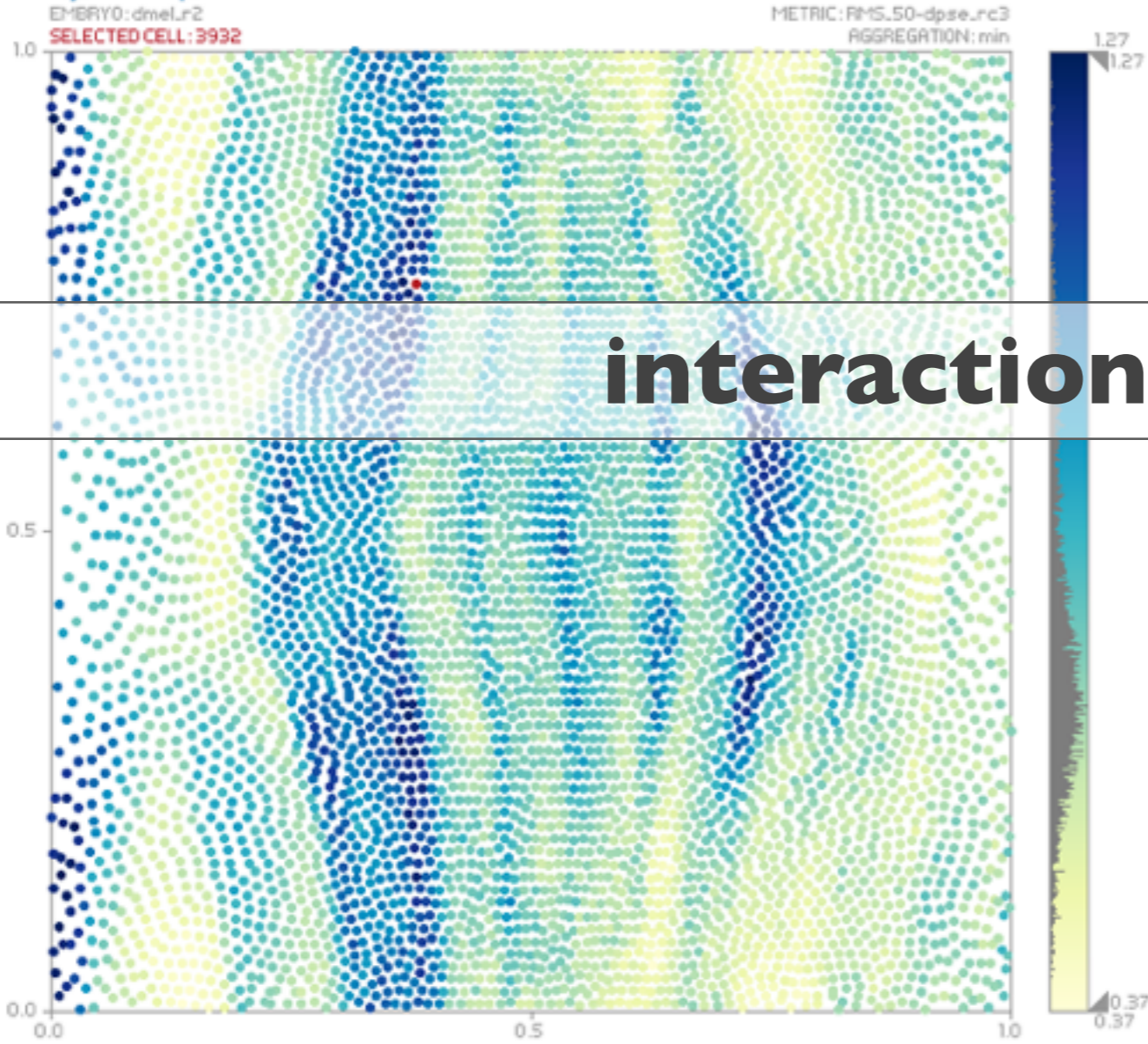




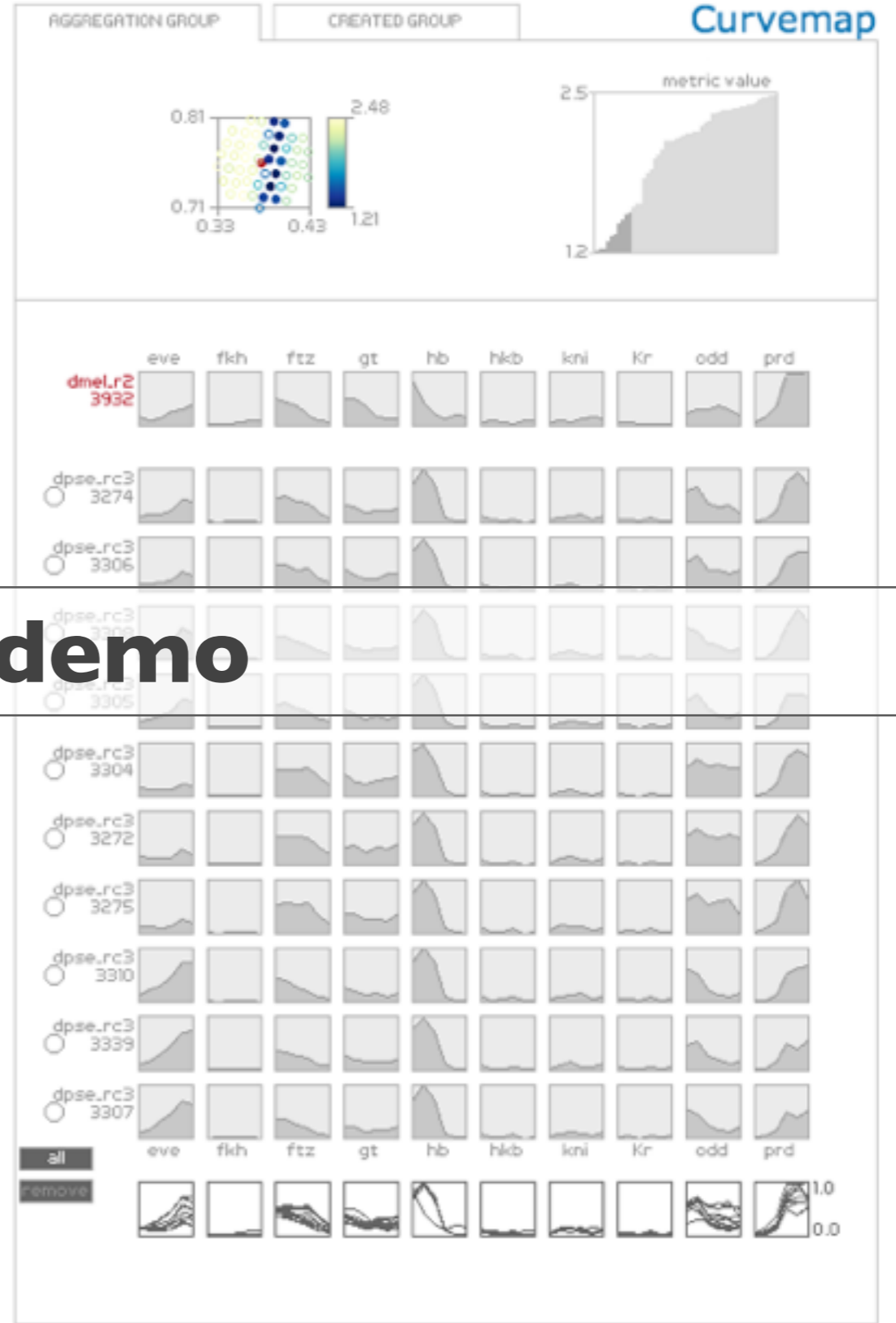
### Summaries



### Embryo Map



interaction demo



data & tool & tasks

summaries & groups

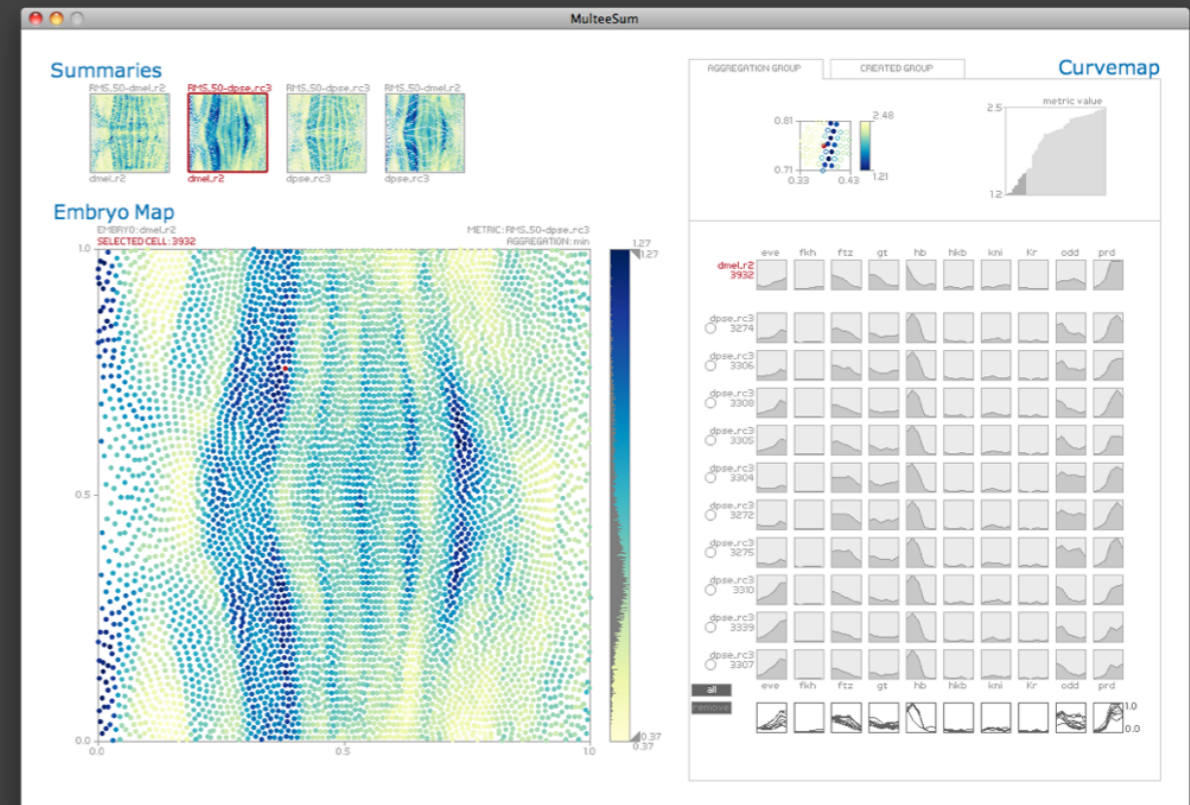
encodings & interaction

**conclusions**

# contributions

## MulteeSum

spatial and temporal gene expression data from multiple species



## workflow

visualization supports upstream computation via summaries

## validation

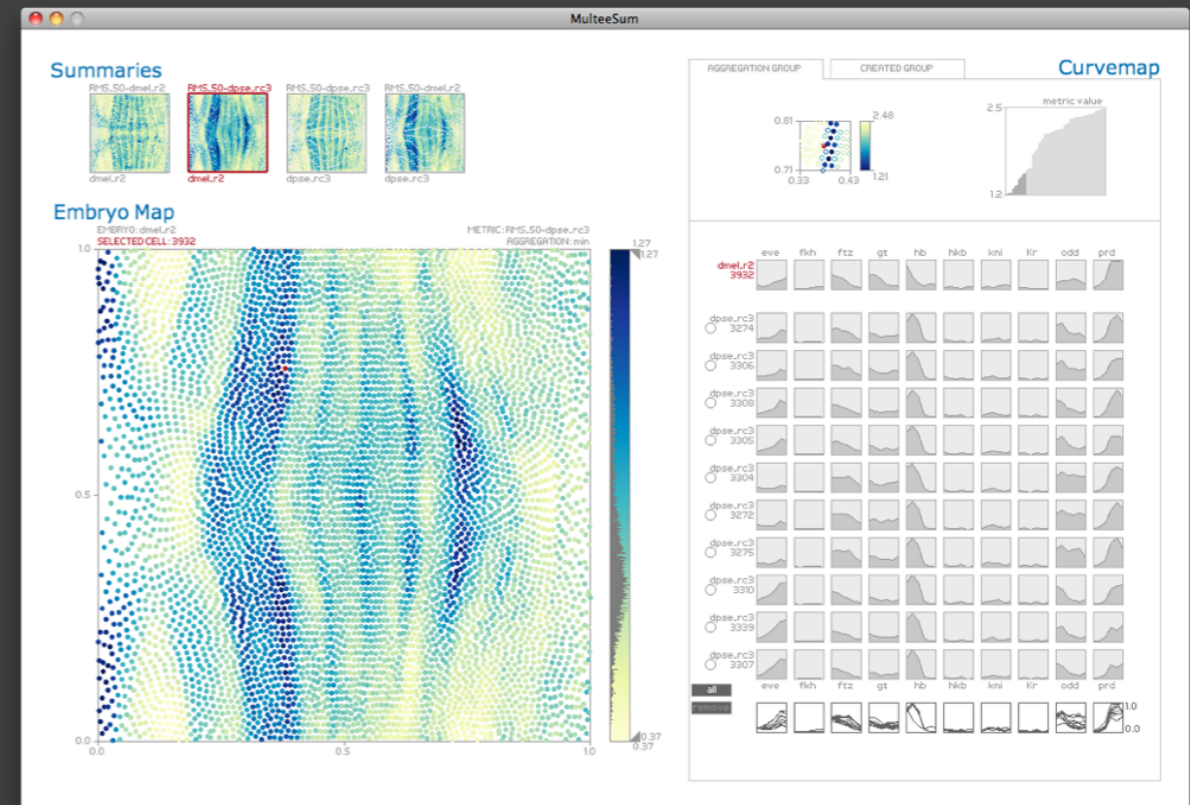
case studies, deployment



# contributions

## MulteeSum

spatial and temporal gene expression data from multiple species



## workflow

visualization supports upstream computation via summaries

## validation

case studies, deployment

Cellular resolution comparison of gene expression in *Drosophila* reveals coordinated shifts in the segmentation network.

*DePace et. al, in preparation.*

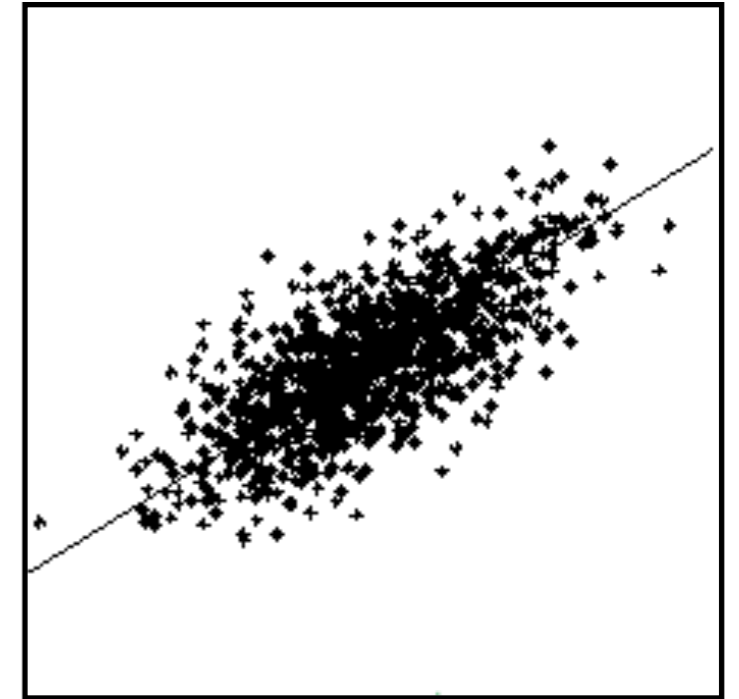
# attribute aggregation

- 1) group attributes and compute a similarity score across the set
- 2) dimensionality reduction, to preserve meaningful structure**

# dimensionality reduction

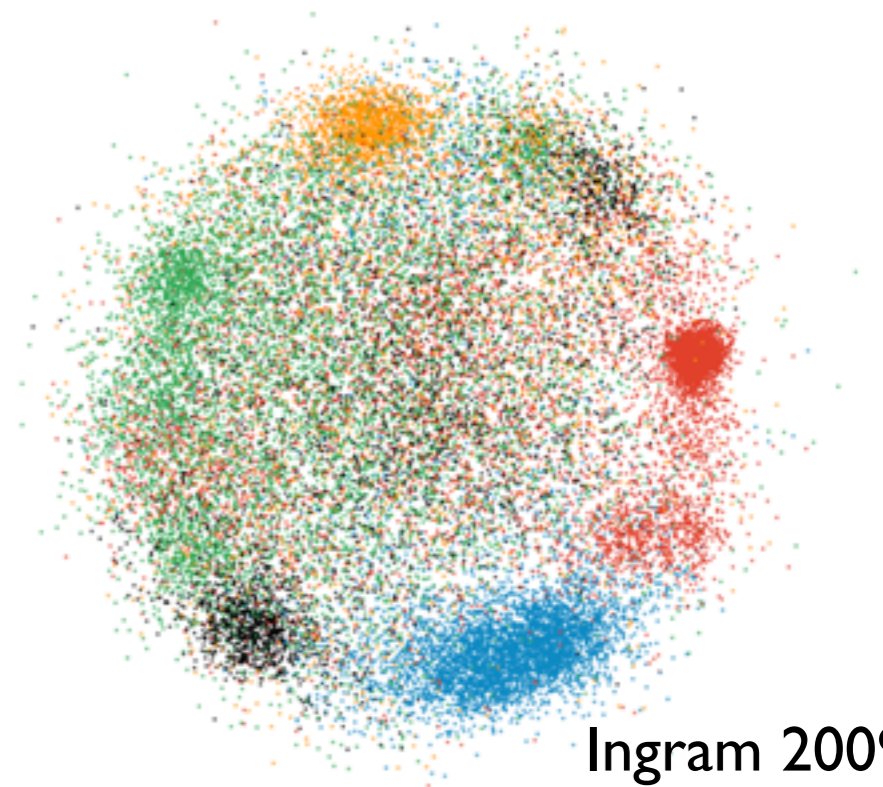
## -PCA

- linear approach
- new dimensions are weighted combinations of original ones
- new dimensions created in order of maximum variance



## -MDS

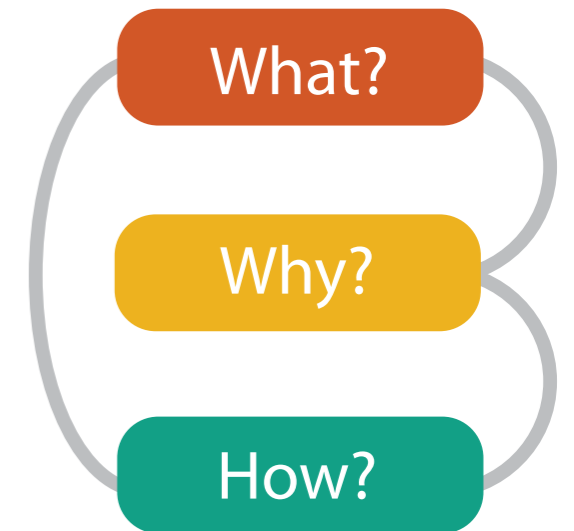
- nonlinear class of approaches
- maximize differences in distances from high dim space in the low dim space



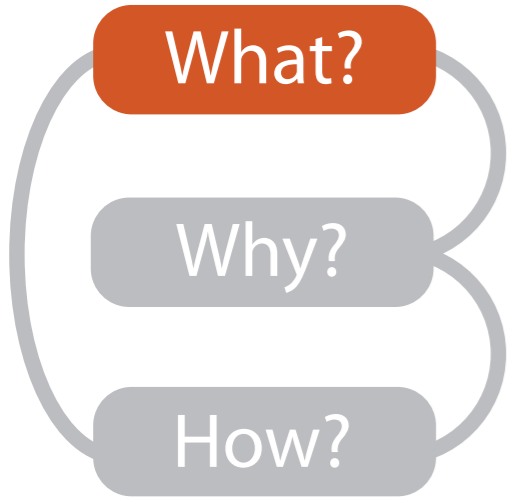
end of foundations...

# Analysis: What, why, and how

- **what** is shown?
  - data abstraction
- **why** is the user looking at it?
  - task abstraction
- **how** is it shown?
  - idiom: visual encoding and interaction
- **abstract vocabulary avoids domain-specific terms**
  - translation process iterative, tricky
- **what-why-how analysis framework as scaffold to think systematically about design space**



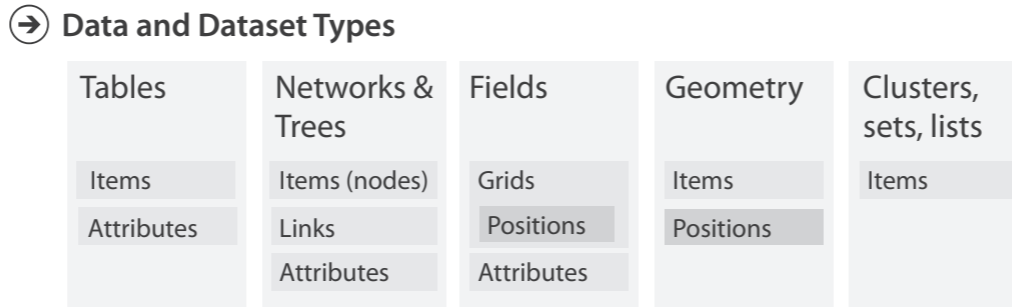




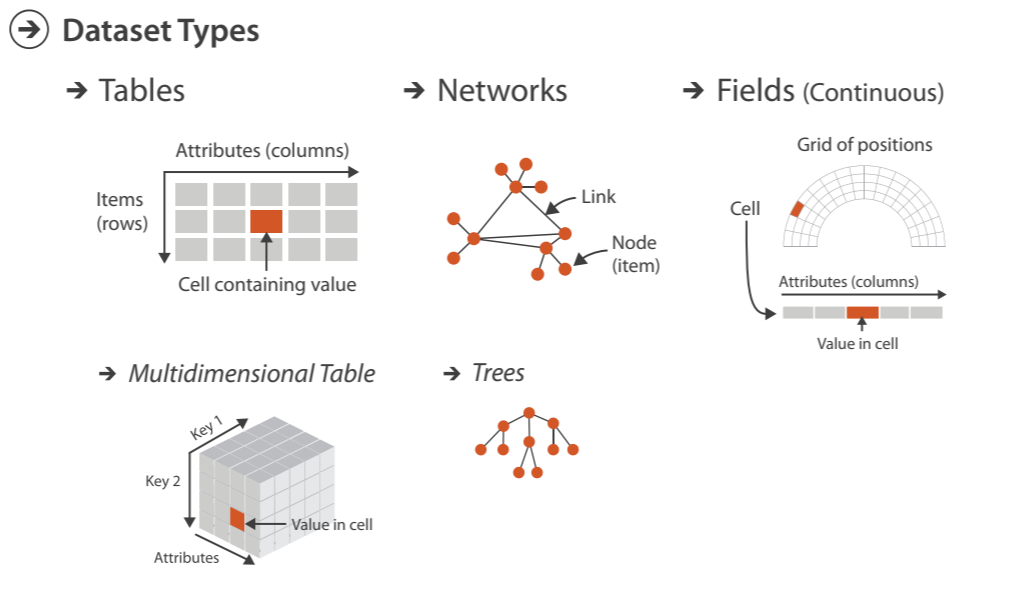
# What?

## Datasets Attributes

→ **Data Types**  
 → Items → Attributes → Links → Positions → Grids



→ **Attribute Types**  
 → Categorical  
 + ● ■ ▲  
 → Ordered  
 → Ordinal  
 ↑ T-shirt icons of increasing size  
 → Quantitative  
 ———|—————|—————

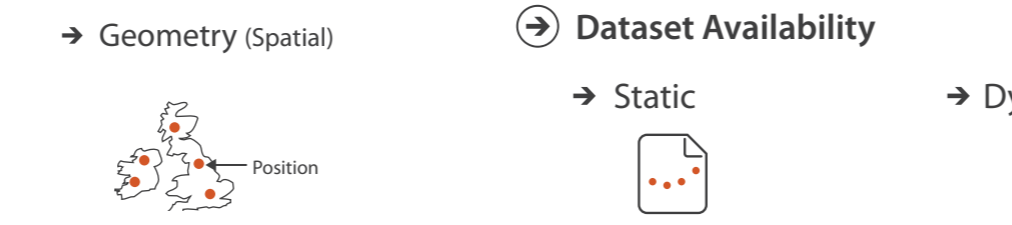


→ **Ordering Direction**

→ Sequential  
 —————→

→ Diverging  
 ←————|————→

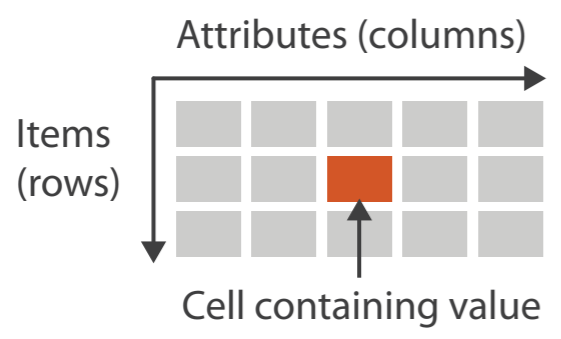
→ Cyclic  
 ↻



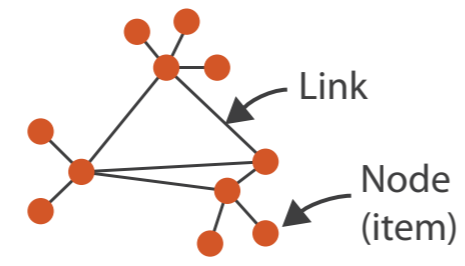
# Dataset types

## → Dataset Types

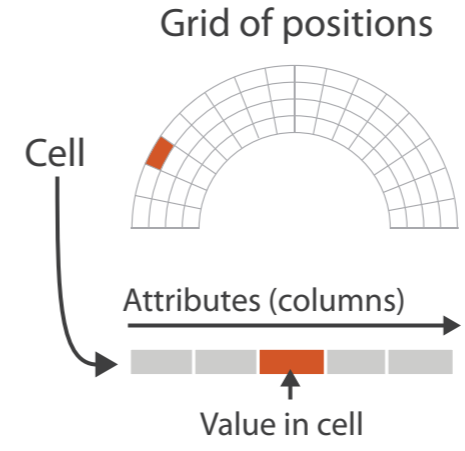
→ Tables



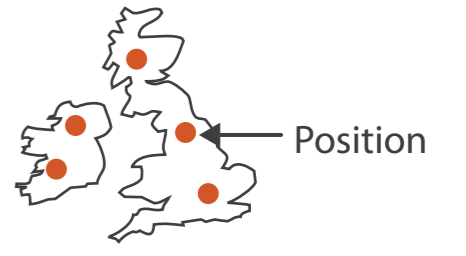
→ Networks



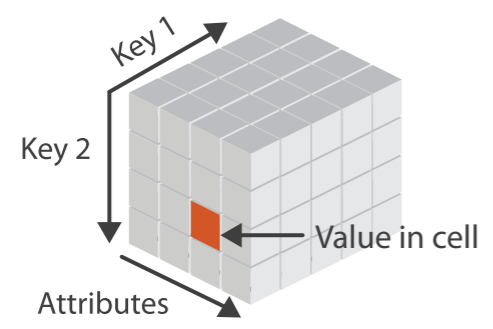
→ Fields (Continuous)



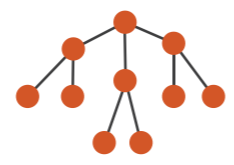
→ Geometry (Spatial)



→ *Multidimensional Table*

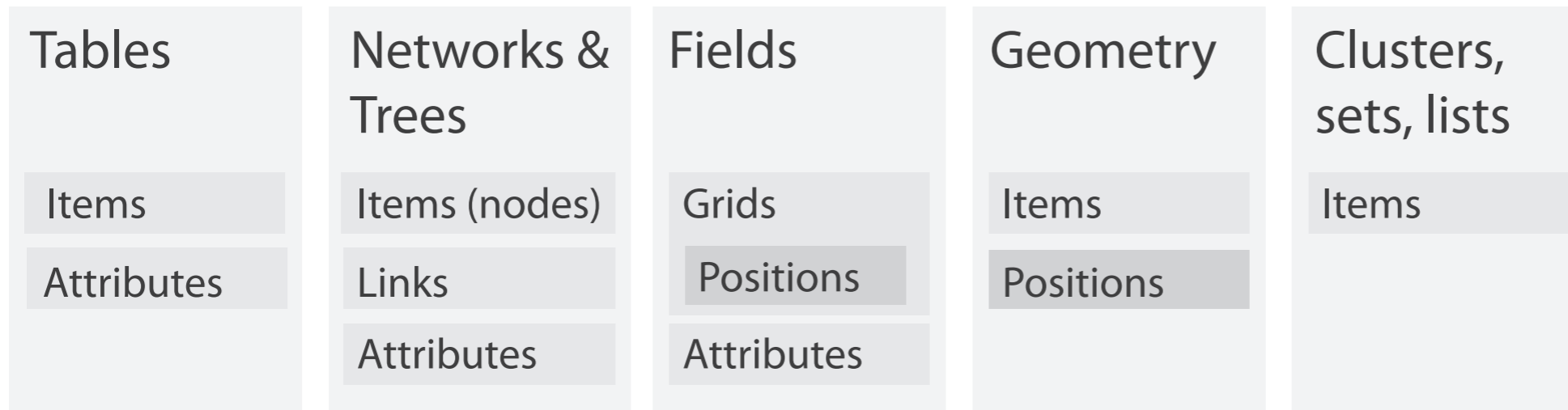


→ *Trees*



# Dataset and data types

## → Data and Dataset Types



## → Data Types

→ Items    → Attributes    → Links    → Positions    → Grids

## → Dataset Availability

→ Static



→ Dynamic



# Attribute types

## ➔ Attribute Types

➔ Categorical



➔ Ordered

➔ *Ordinal*



➔ *Quantitative*



## ➔ Ordering Direction

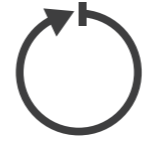
➔ Sequential

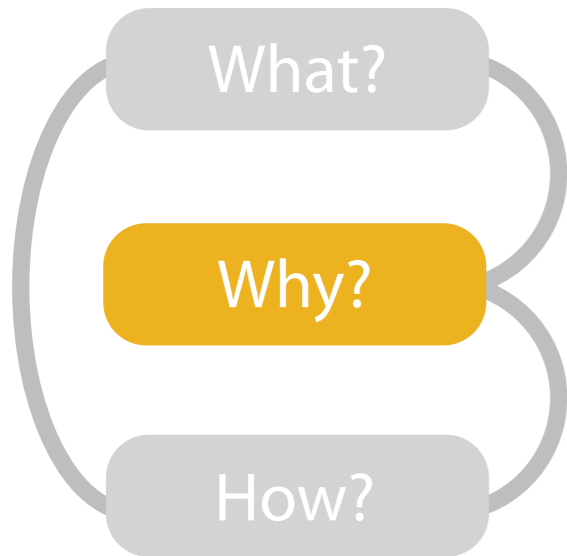


➔ Diverging



➔ Cyclic





# Why?

## 👉 Actions

## 🎯 Targets

### ➔ Analyze

➔ Consume

➔ Discover



➔ Present



➔ Enjoy



➔ Produce

➔ Annotate



➔ Record



➔ Derive

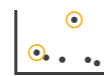


### ➔ Search

	Target known	Target unknown
Location known	••• Lookup	••• Browse
Location unknown	<•••> Locate	<•••> Explore

### ➔ Query

➔ Identify



➔ Compare

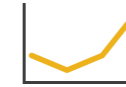


➔ Summarise



### ➔ All Data

➔ Trends



➔ Outliers



➔ Features



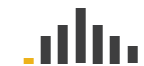
### ➔ Attributes

➔ One

➔ Distribution



➔ Extremes

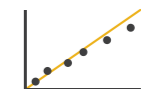


➔ Many

➔ Dependency



➔ Correlation



➔ Similarity



### ➔ Network Data

➔ Topology

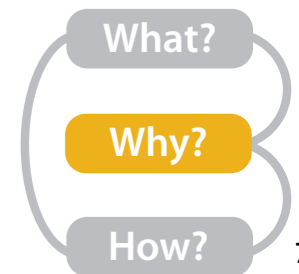


➔ Paths

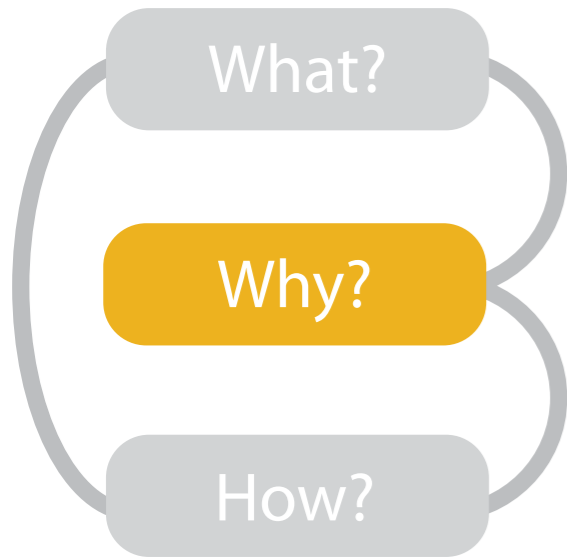


### ➔ Spatial Data

➔ Shape



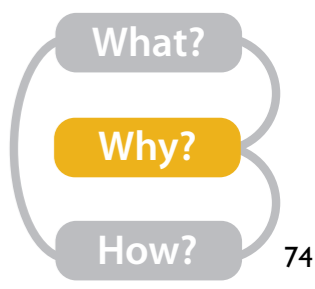




- {action, target} pairs
  - discover distribution
  - compare trends
  - locate outliers
  - browse topology

## Why?

👉 Actions	🎯 Targets									
<p>➔ <b>Analyze</b></p> <p>➔ Consume</p> <p>➔ Discover  ➔ Present  ➔ Enjoy </p> <p>➔ Produce</p> <p>➔ Annotate  ➔ Record  ➔ Derive </p>	<p>➔ <b>All Data</b></p> <p>➔ Trends  ➔ Outliers  ➔ Features </p> <p>➔ <b>Attributes</b></p> <p>➔ One</p> <p>➔ Distribution  ➔ Extremes </p> <p>➔ Many</p> <p>➔ Dependency  ➔ Correlation  ➔ Similarity </p>									
<p>➔ <b>Search</b></p> <table border="1" style="width: 100%; border-collapse: collapse; text-align: center;"> <thead> <tr> <th></th> <th>Target known</th> <th>Target unknown</th> </tr> </thead> <tbody> <tr> <td>Location known</td> <td> Lookup</td> <td> Browse</td> </tr> <tr> <td>Location unknown</td> <td> Locate</td> <td> Explore</td> </tr> </tbody> </table>		Target known	Target unknown	Location known	Lookup	Browse	Location unknown	Locate	Explore	<p>➔ <b>Network Data</b></p> <p>➔ Topology </p> <p>➔ Paths </p>
	Target known	Target unknown								
Location known	Lookup	Browse								
Location unknown	Locate	Explore								
<p>➔ <b>Query</b></p> <p>➔ Identify  ➔ Compare  ➔ Summarise </p>	<p>➔ <b>Spatial Data</b></p> <p>➔ Shape </p>									



# High-level actions: Analyze

- consume

- discover vs present

- classic split
- aka explore vs explain

- enjoy

- newcomer
- aka casual, social

- produce

- annotate, record

- derive

- crucial design choice

## → Analyze

### → Consume

→ Discover



→ Present



→ Enjoy



### → Produce

→ Annotate



→ Record



→ Derive



# Actions: Mid-level search, low-level query





- what does user know?

– target, location

- how much of the data matters?

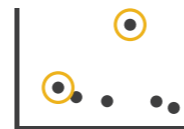
– one, some, all

## → Search

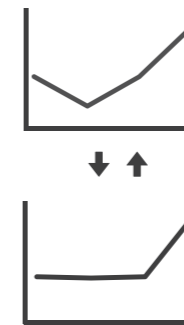
	Target known	Target unknown
Location known	 <i>Lookup</i>	 <i>Browse</i>
Location unknown	 <i>Locate</i>	 <i>Explore</i>

## → Query

→ Identify



→ Compare



→ Summarise



# Why: Targets

## → ALL DATA

→ Trends



→ Outliers



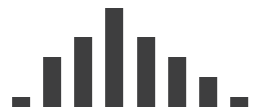
→ Features



## → ATTRIBUTES

→ One

→ *Distribution*



↓ *Extremes*

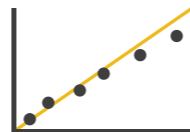


→ Many

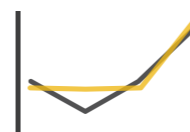
→ *Dependency*



→ *Correlation*

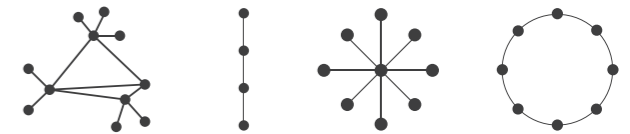


→ *Similarity*

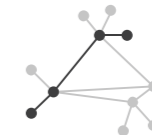


## → NETWORK DATA

→ Topology

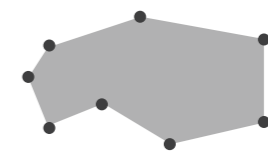


→ *Paths*



## → SPATIAL DATA

→ Shape



# How?

## Encode

### ➔ Arrange

➔ Express



➔ Separate



➔ Order



➔ Align



➔ Use



### ➔ Map

from **categorical** and **ordered** attributes

➔ Color

➔ Hue



➔ Saturation



➔ Luminance



➔ Size, Angle, Curvature, ...



➔ Shape



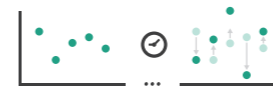
➔ Motion

*Direction, Rate, Frequency, ...*



## Manipulate

### ➔ Change



### ➔ Select



### ➔ Navigate



## Facet

### ➔ Juxtapose



### ➔ Partition



### ➔ Superimpose



## Reduce

### ➔ Filter



### ➔ Aggregate



### ➔ Embed



What?

Why?

How?



+ perception and design guidelines...

L12: Tables

# REQUIRED READING

# Chapter 7

## Arrange Tables

### 7.1 The Big Picture

Figure 7.1 shows the four visual encoding design choices for how to arrange tabular data spatially. One is to express values. The other three are to separate, order, and align regions. The spatial orientation of axes can be rectilinear, parallel, or radial. Spatial layouts may be dense, and they may be space-filling.

► A fifth arrangement choice, to use a given spatial layout, is not an option for nonspatial information; it is covered in Chapter 8.

### 7.2 Why Arrange?

The **arrange** design choice covers all aspects of the use of spatial channels for visual encoding. It is the most crucial visual encoding choice because the use of space dominates the user's mental model of the dataset. The three highest ranked effectiveness channels for quantitative and ordered attributes are all related to spatial position: planar position along an aligned scale, planar position against a common scale, and length. The highest ranked effectiveness channel for categorical attributes, grouping items within the same region, is also about the use of space. Moreover, there are no nonspatial channels that are highly effective for all attribute types: the others are split into being suitable for either ordered or categorical attributes, but not both, because of the principle of expressiveness.

► The primacy of the spatial position channels is discussed at length in Chapter 5, as are the principles of effectiveness and expressiveness.

### 7.3 Arrange by Keys and Values