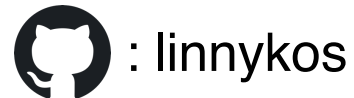




Tilted-CCA: Quantifying common and distinct information in multi-modal single-cell data via matrix factorization

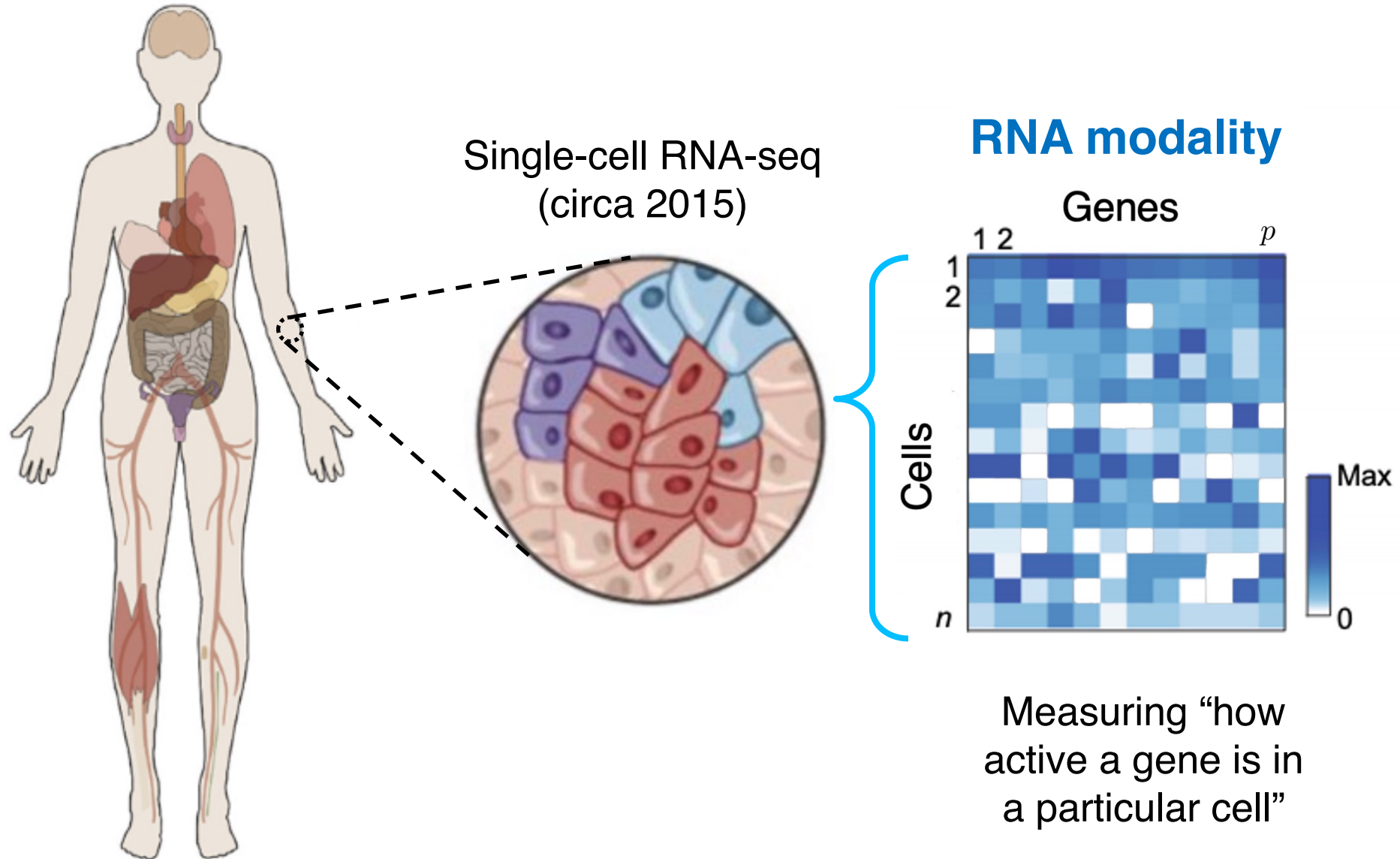
Kevin Lin



January 3,
2023



My body of work focuses on advancing statistical ideas that are inspired by advancements in biomedical technology.



Previous statistical methods & theory I've developed for single-cell RNA-seq:

- Relevant for denoising single-cell data (Matrix factorization)
 - K. Lin, H. Liu, K. Roeder (JASA 2021)
 - K. Lin, J. Lei, K. Roeder (JASA 2022)
- Relevant for answering biological questions (Network methods and changepoint detection)
 - K. Lin, J. Sharpnack, A. Rinaldo, R.J. Tibshirani (Neurips 2017)
 - S. Hyun, K. Lin, M. G'Sell, R.J. Tibshirani (Biometrics 2021)
 - J. Lei, K. Lin (JASA 2022)

Previous statistical methods & theory I've developed for single-cell RNA-seq:

- Relevant for denoising single-cell data (Matrix factorization)
 - K. Lin, H. Liu, K. Roeder (JASA 2021)
 - K. Lin, J. Lei, K. Roeder (JASA 2022)
- Relevant for answering biological questions (Network methods and changepoint detection)
 - K. Lin, J. Sharpnack, A. Rinaldo, R.J. Tibshirani (Neurips 2017)
 - S. Hyun, K. Lin, M. G'Sell, R.J. Tibshirani (Biometrics 2021)
 - J. Lei, K. Lin (JASA 2022)

However, our talk today is inspired by newer biomedical technology.

Our story today starts with a newer technology:

FOCUS | EDITORIAL

Method of the Year 2019: Single-cell multimodal omics

Multimodal omics measurement offers opportunities for gaining holistic views of cells one by one.

Teichmann, Efremova. Nature Methods (2020)

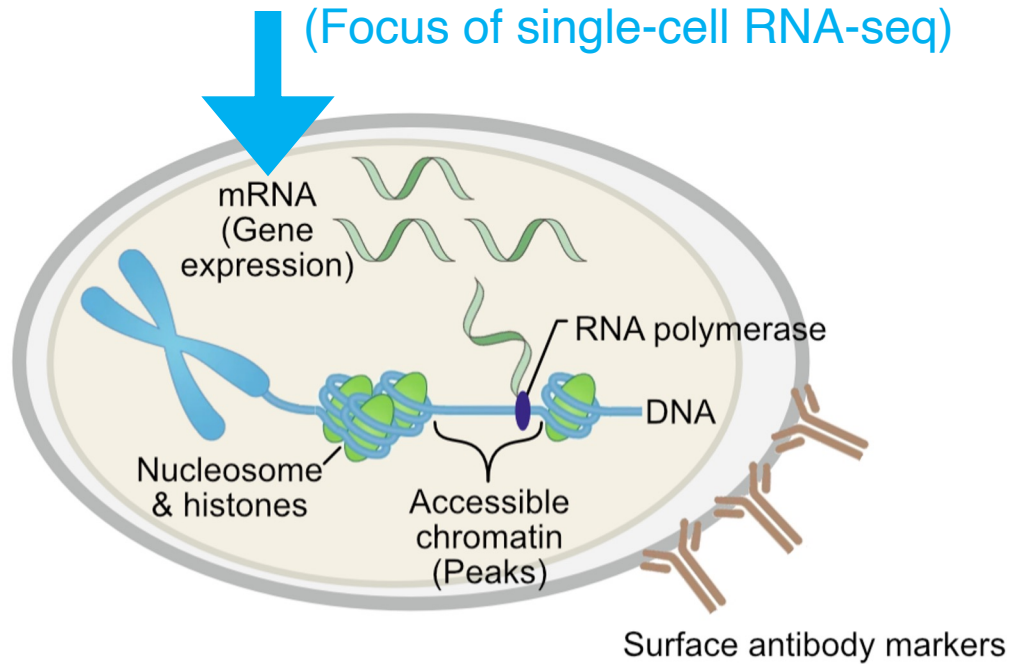
FOCUS | COMMENT

Single-cell multimodal omics: the power of many

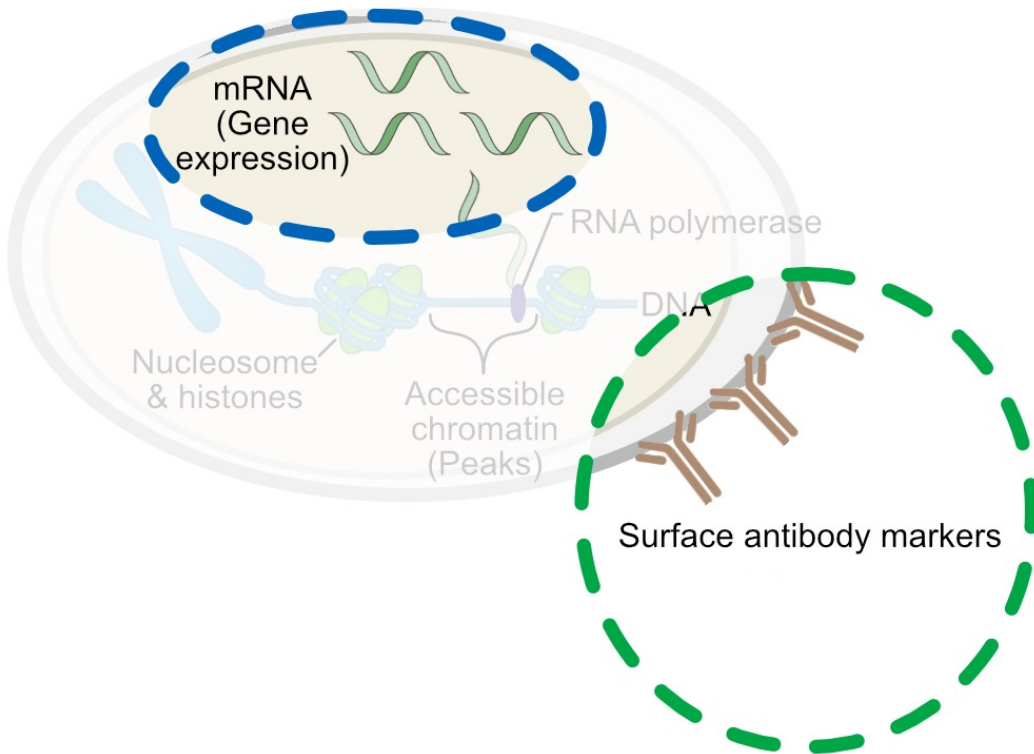
Advances in single-cell genomics technologies have enabled investigation of the gene regulation programs of multicellular organisms at unprecedented resolution and scale. Development of single-cell multimodal omics tools is another major step toward understanding the inner workings of biological systems.

Zhu, Preissl, Ren. Nature Methods (2020)

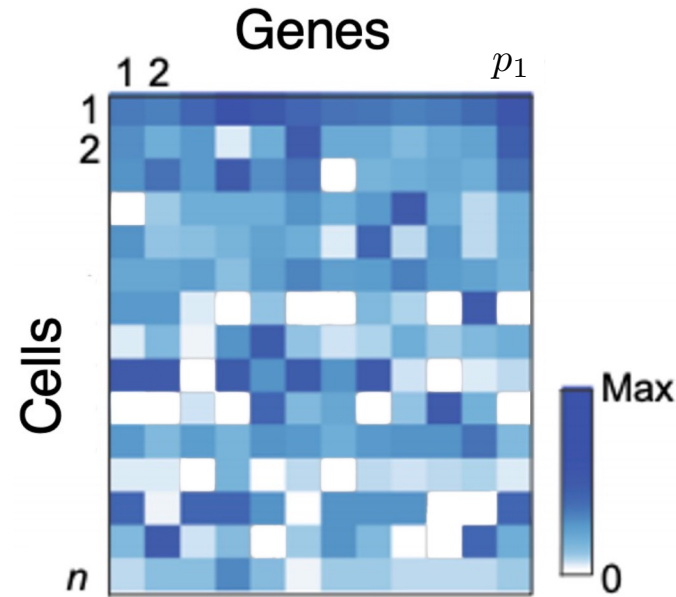
Our story today starts with a newer technology: Single-cell multi-modal (i.e., “multi-view”) sequencing (circa 2020).



Our story today starts with a newer technology: Single-cell multi-modal (i.e., “multi-view”) sequencing (circa 2020).

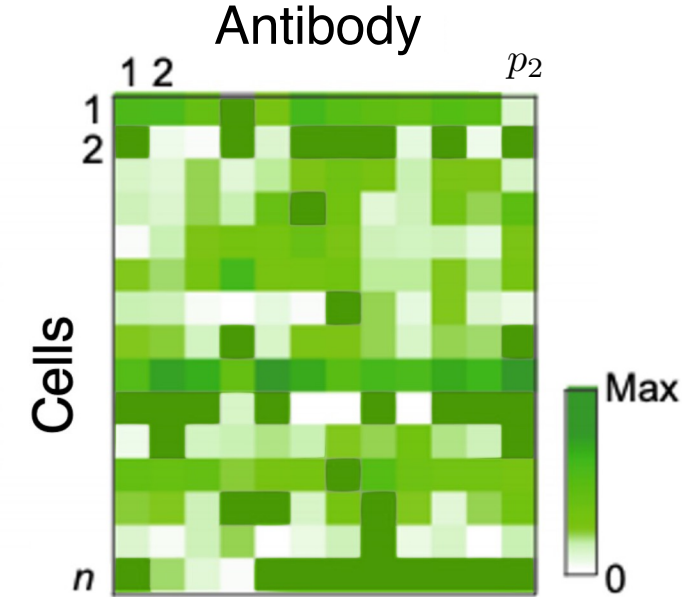


RNA modality



$$X^{(1)} \in \mathbb{R}^{n \times p_1}$$

Protein modality

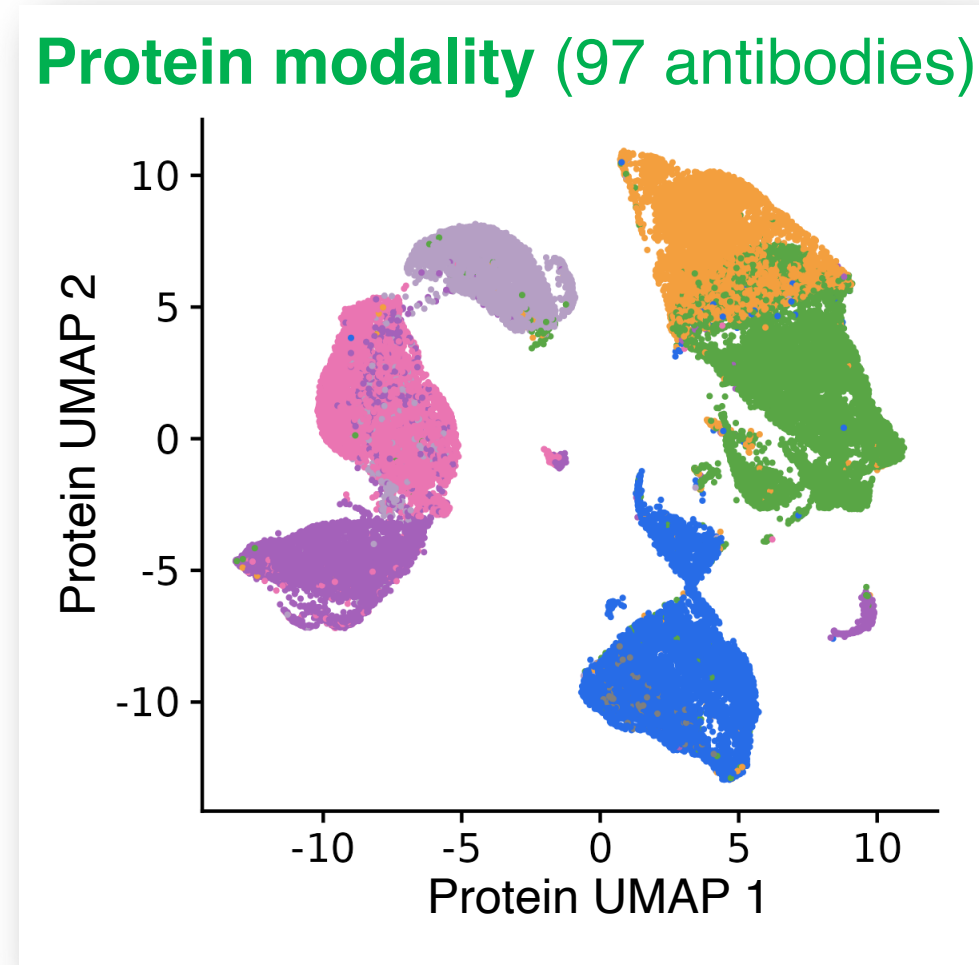
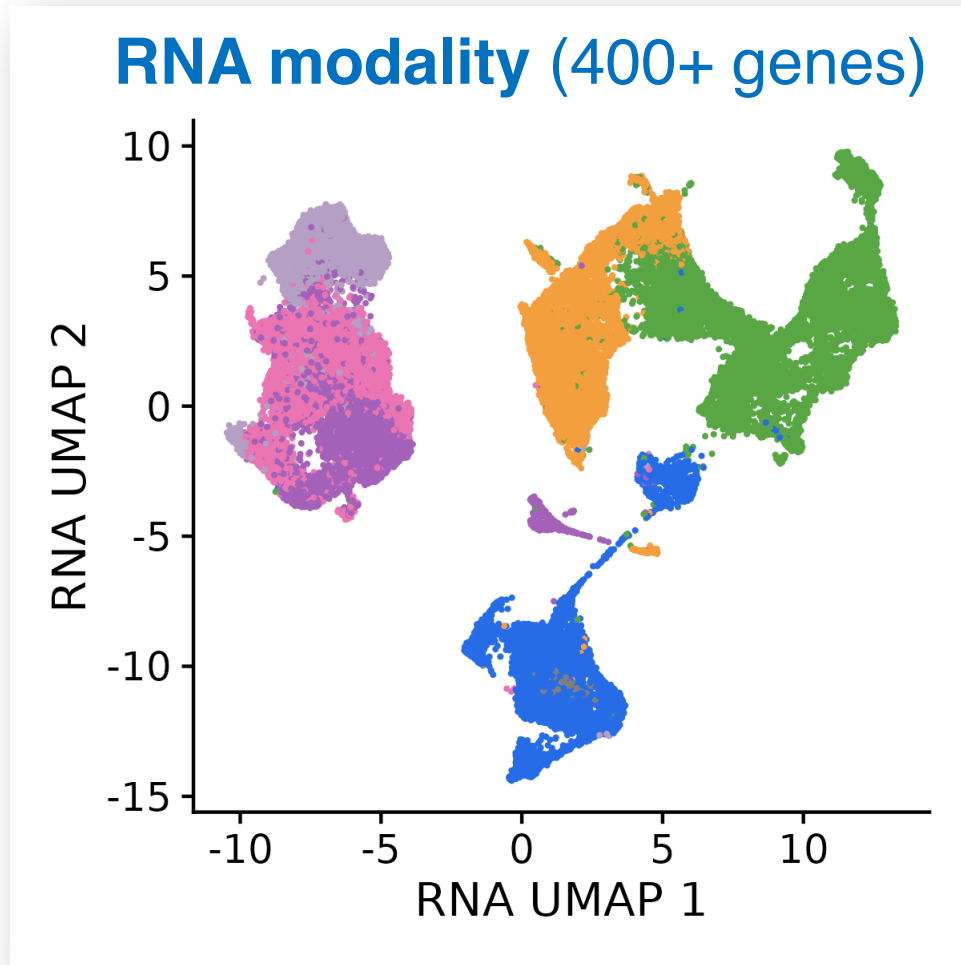


$$X^{(2)} \in \mathbb{R}^{n \times p_2}$$

Central question: What “information” is unique to a modality, or represents the coordination between both modalities?

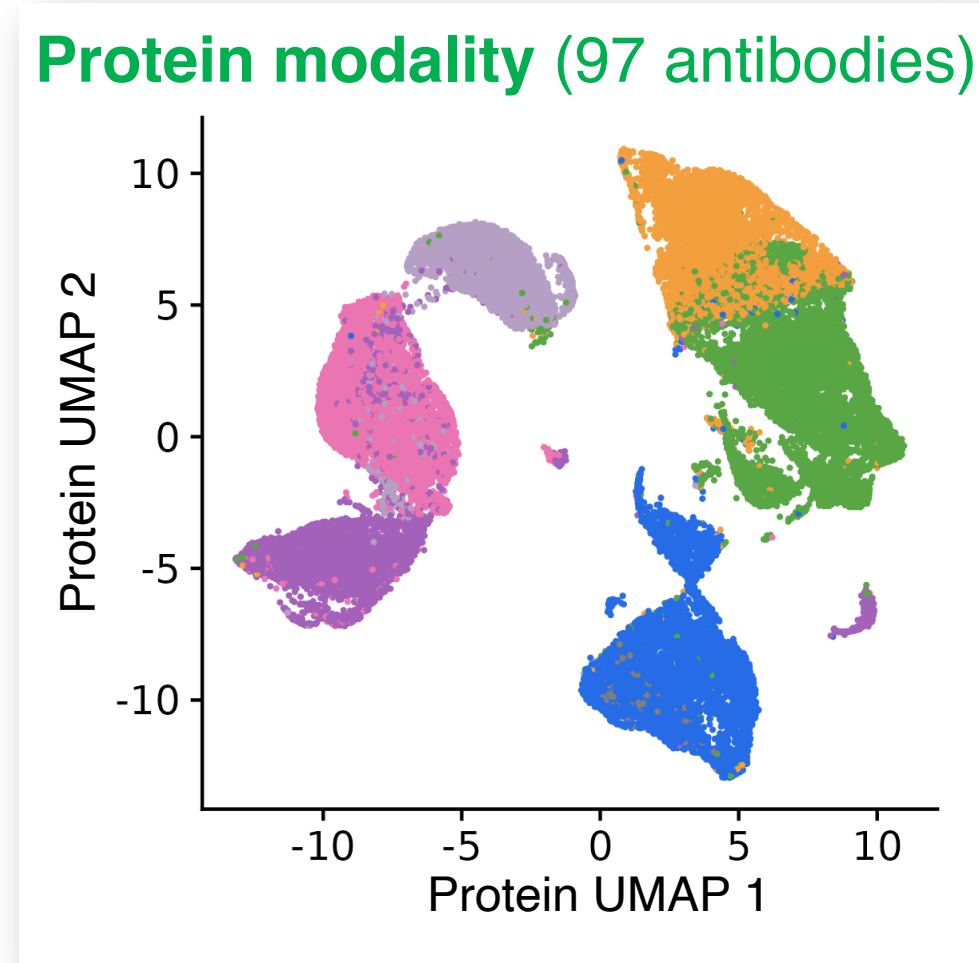
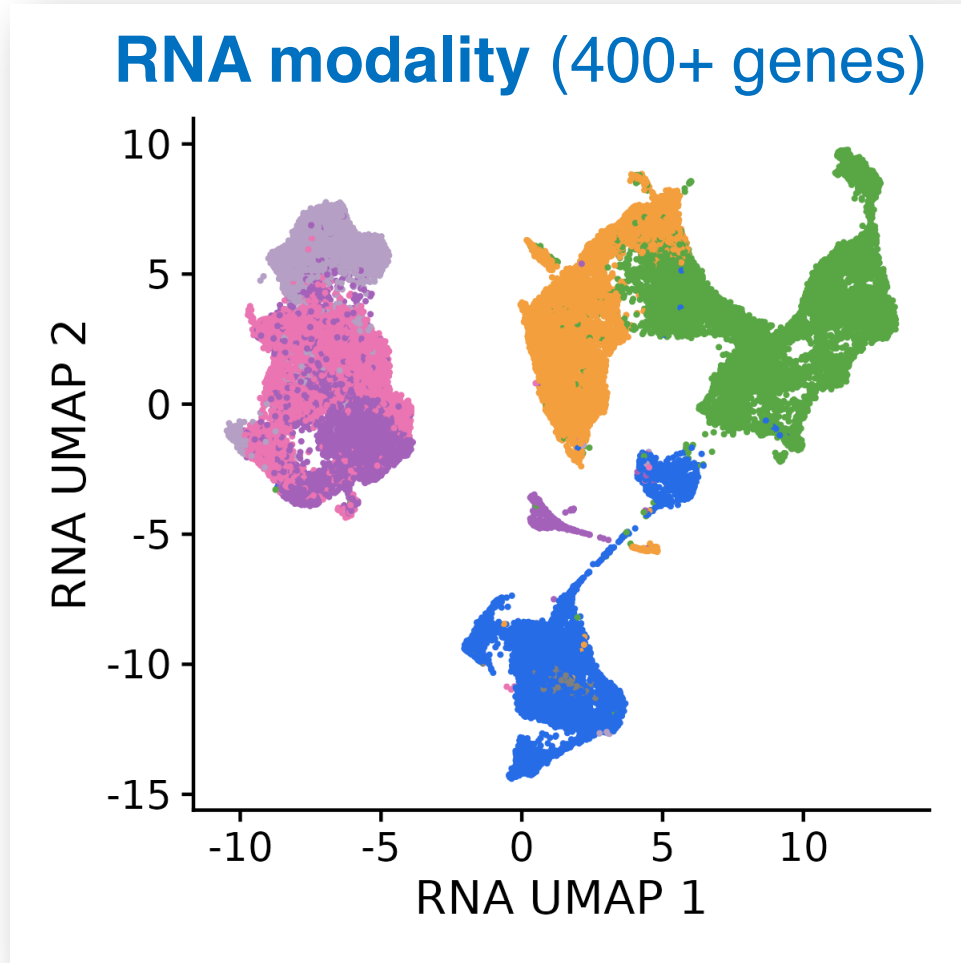
What “information” is unique to a modality, or represents the coordination between both modalities?

What “information” is **unique** to a modality, or represents the **coordination** between both modalities?

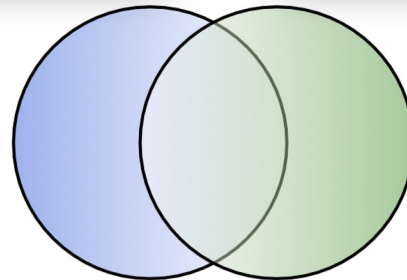


Human bone marrow (AbSeq, Triana et al., 2021), 49000+ cells (colored by annotated cell type)

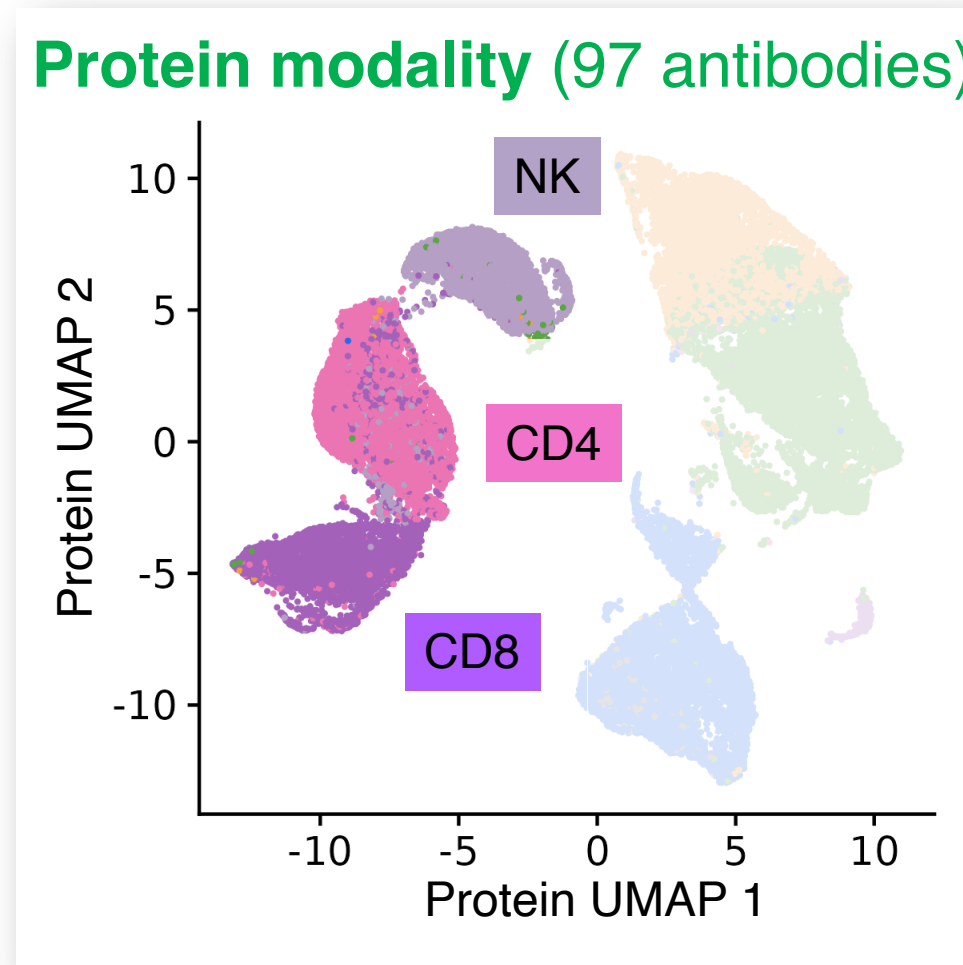
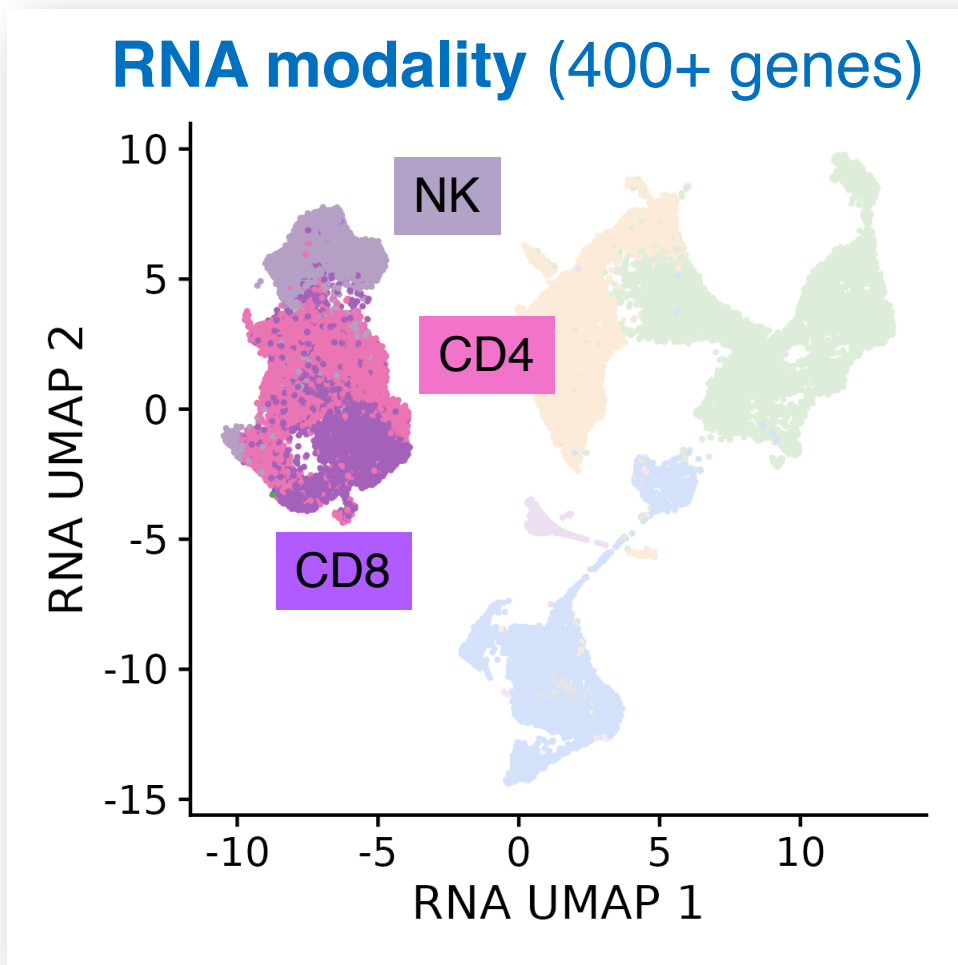
What “information” is **unique** to a modality, or represents the **coordination** between both modalities?



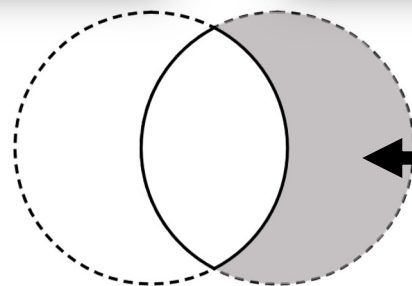
Venn diagram of geometry (“information”):



What “information” is **unique** to a modality, or represents the **coordination** between both modalities?

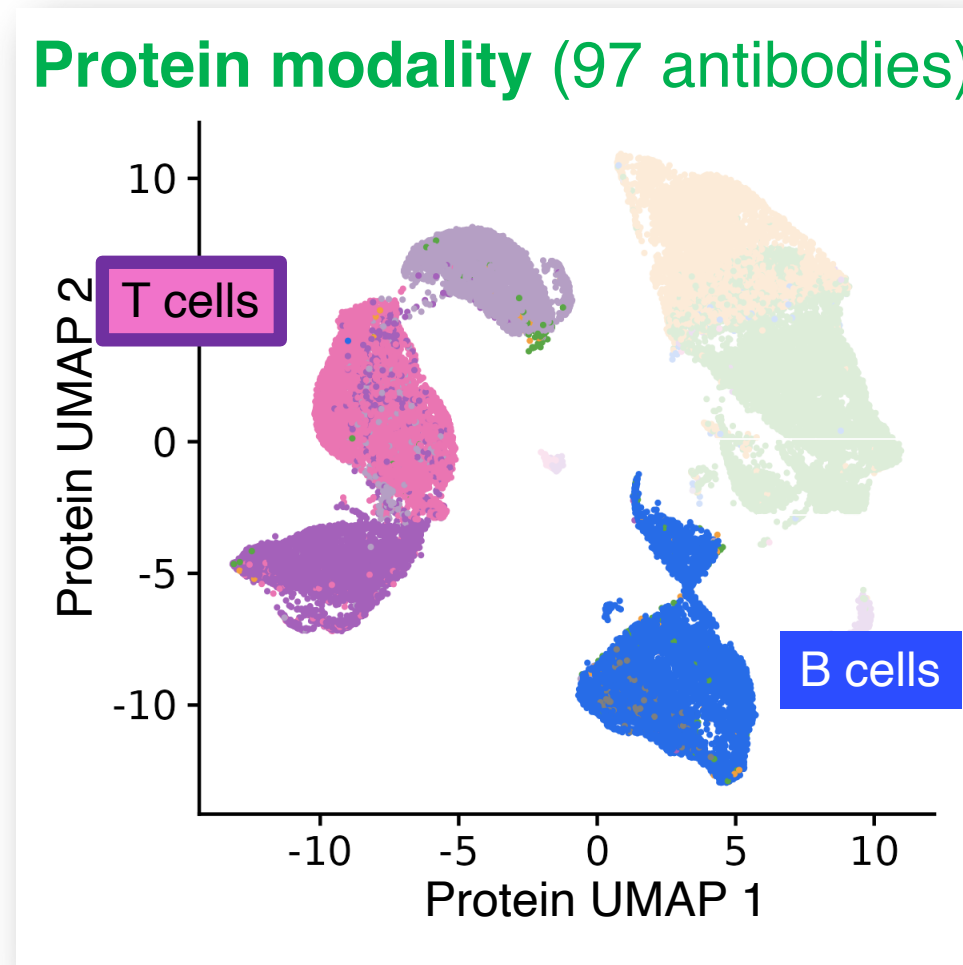
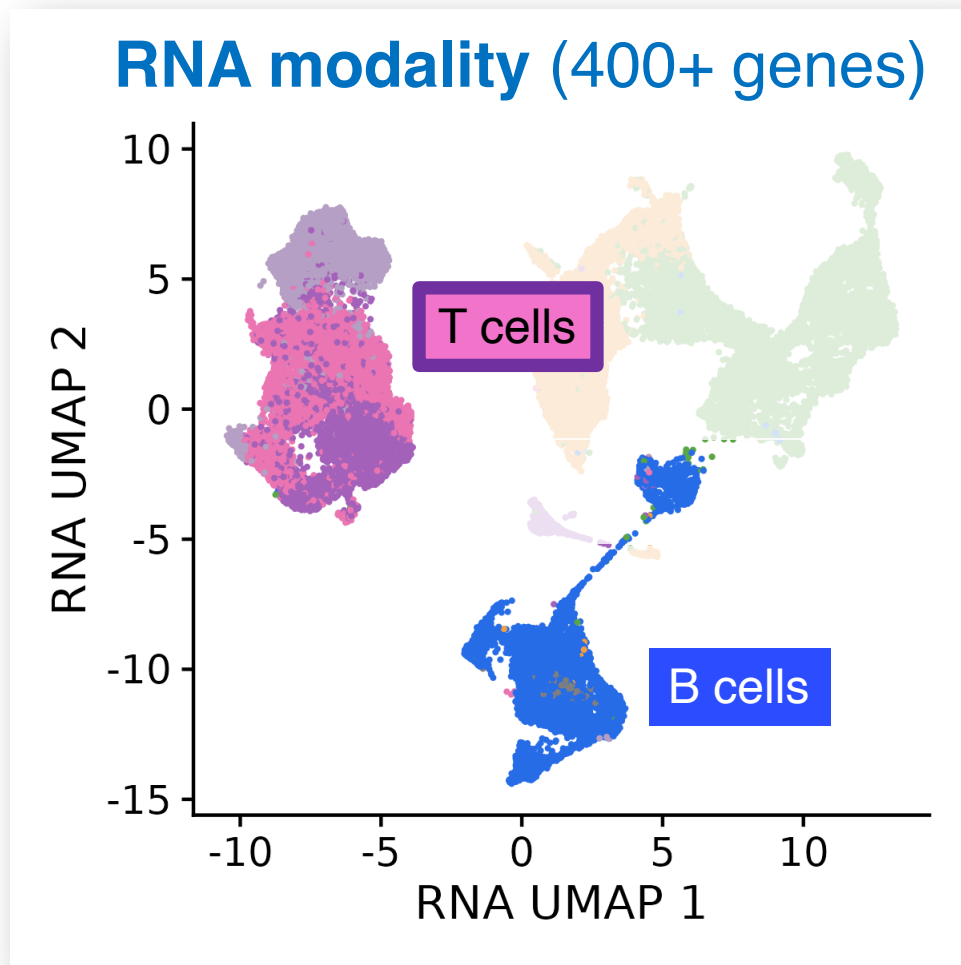


Venn diagram of geometry (“information”):

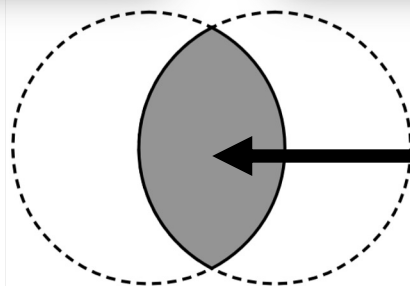


Cell-type separation unique to proteins

What “information” is **unique** to a modality, or represents the **coordination** between both modalities?

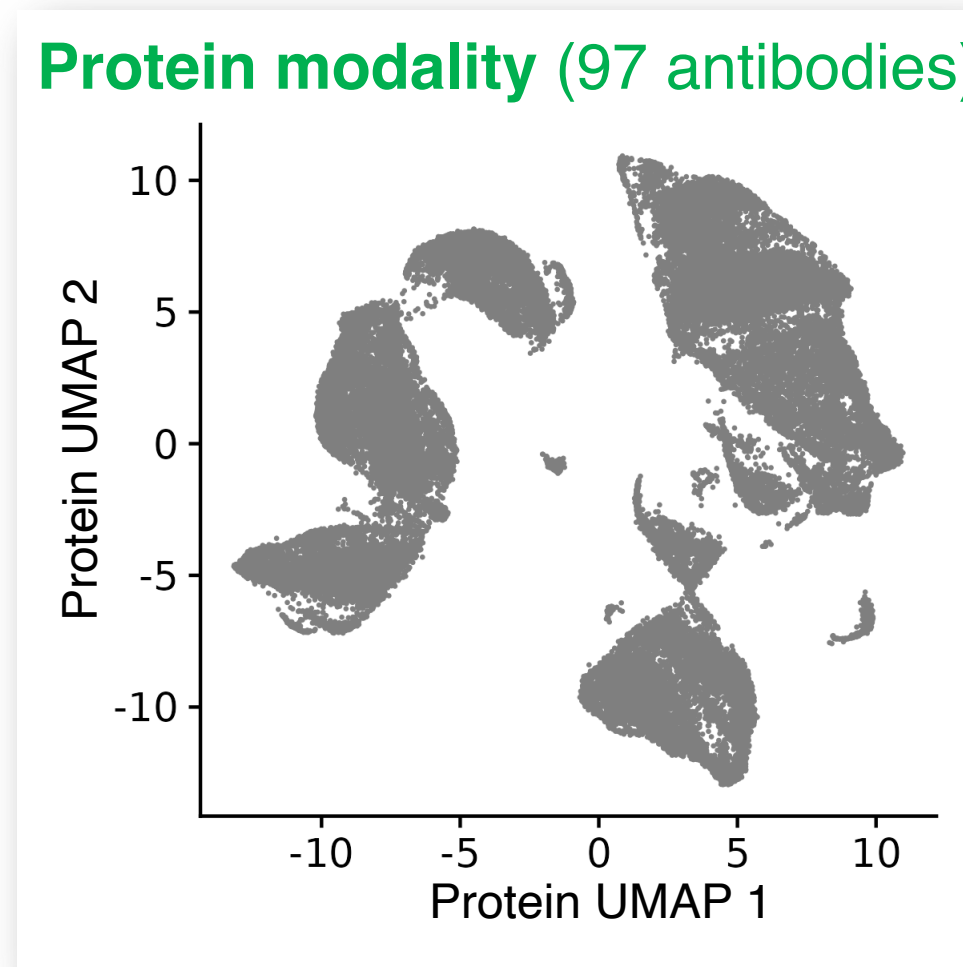
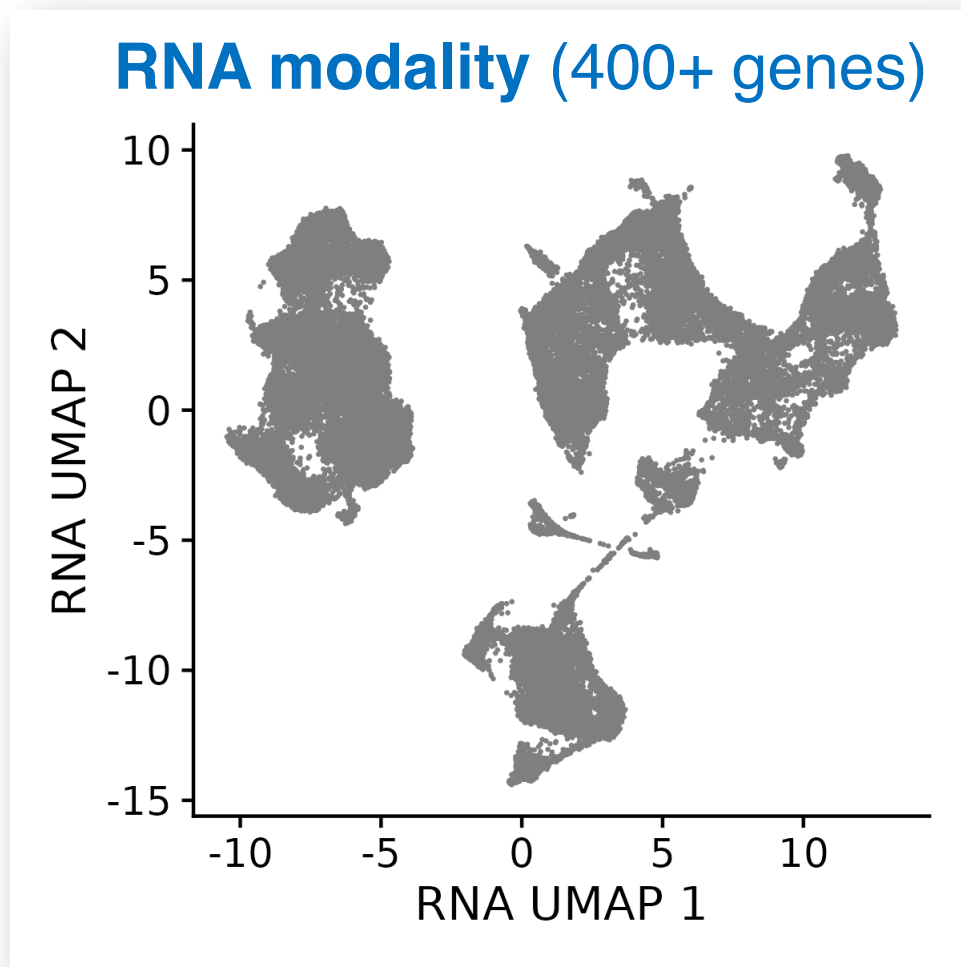


Venn diagram of geometry (“information”):

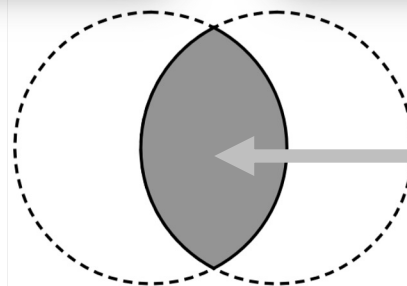


Between-modality coordination

What “information” is **unique** to a modality, or represents the **coordination** between both modalities?



Venn diagram of geometry (“information”):
(No cell-type information)



Between-modality
coordination

Statistical goal: Develop a new matrix factorization framework for multi-modal data based on shared/unique geometry to answer the following biological questions.

1. **(Experimental design):** Which pair of modalities should biologist sequence to have the most comprehensive understanding?

Statistical goal: Develop a new matrix factorization framework for multi-modal data based on shared/unique geometry to answer the following biological questions.

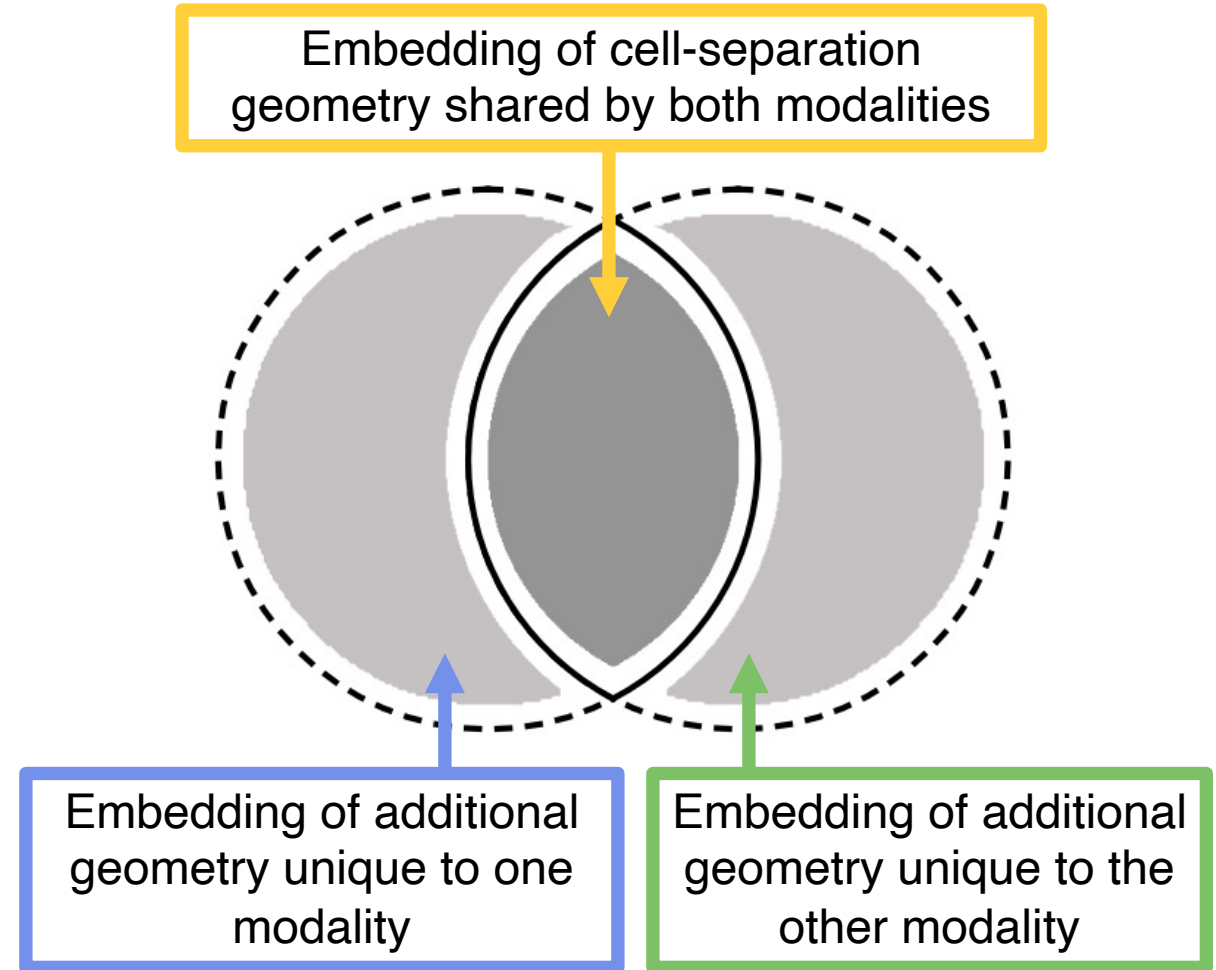
1. **(Experimental design):** Which pair of modalities should biologist sequence to have the most comprehensive understanding?
2. **(Variable selection):** For RNA-Protein data, how can we pick the antibodies that contribute the most additional information to the RNA modality?

Statistical goal: Develop a new matrix factorization framework for multi-modal data based on shared/unique geometry to answer the following biological questions.

1. **(Experimental design):** Which pair of modalities should biologist sequence to have the most comprehensive understanding?
2. **(Variable selection):** For RNA-Protein data, how can we pick the antibodies that contribute the most additional information to the RNA modality?
3. **(Developmental biology):** Can the amount of coordination between two modalities tell us if a cell is in a steady-state or is undergoing development?

Statistical goal: Develop a new matrix factorization framework for multi-modal data based on shared/unique geometry to answer the following biological questions.

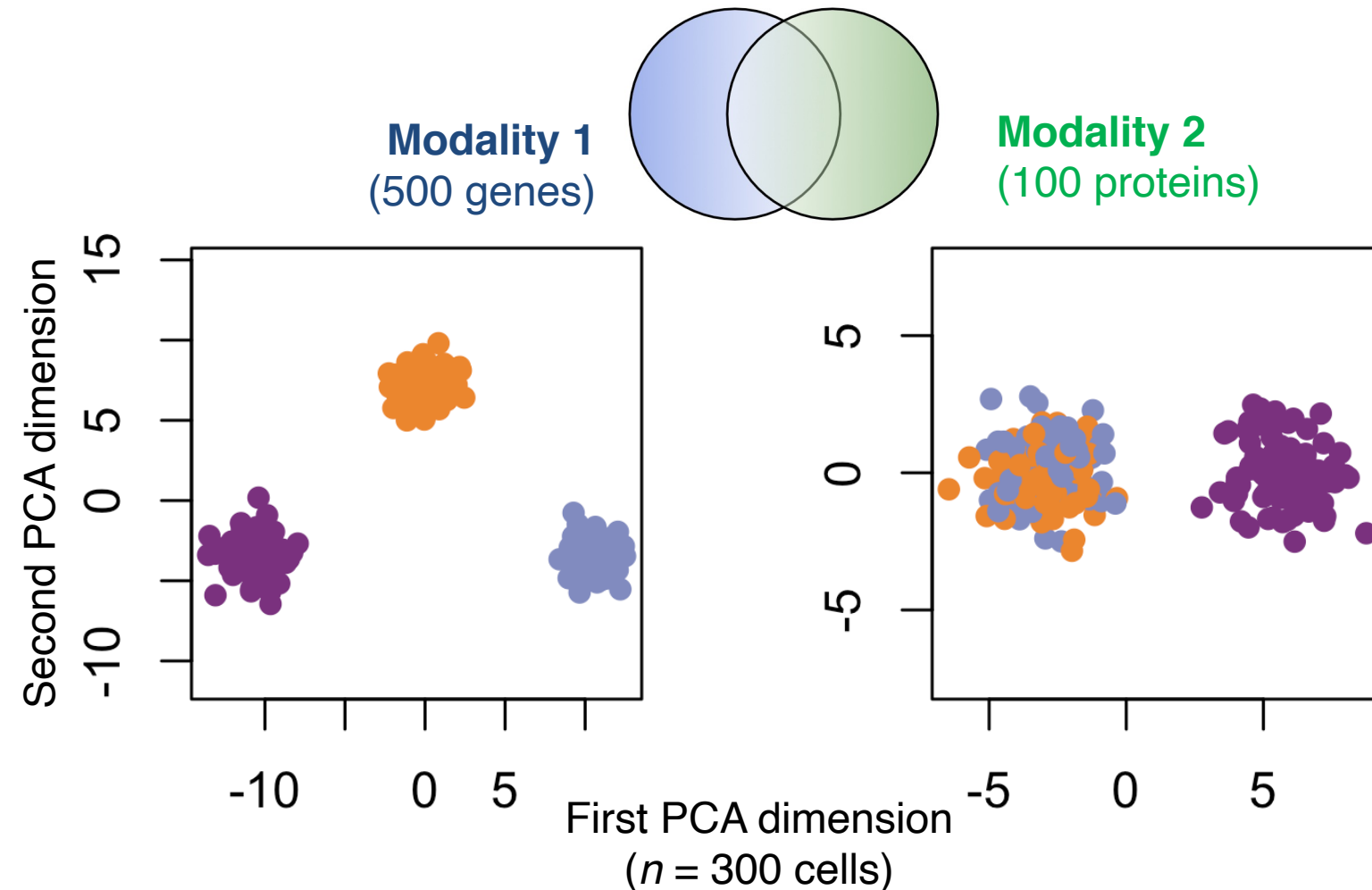
1. **(Experimental design):** Which pair of modalities should biologist sequence to have the most comprehensive understanding?
2. **(Variable selection):** For RNA-Protein data, how can we pick the antibodies that contribute the most additional information to the RNA modality?
3. **(Developmental biology):** Can the amount of coordination between two modalities tell us if a cell in a steady-state or is undergoing development?



Our method inspires new theoretical questions and is applicable to any multi-modal dataset.

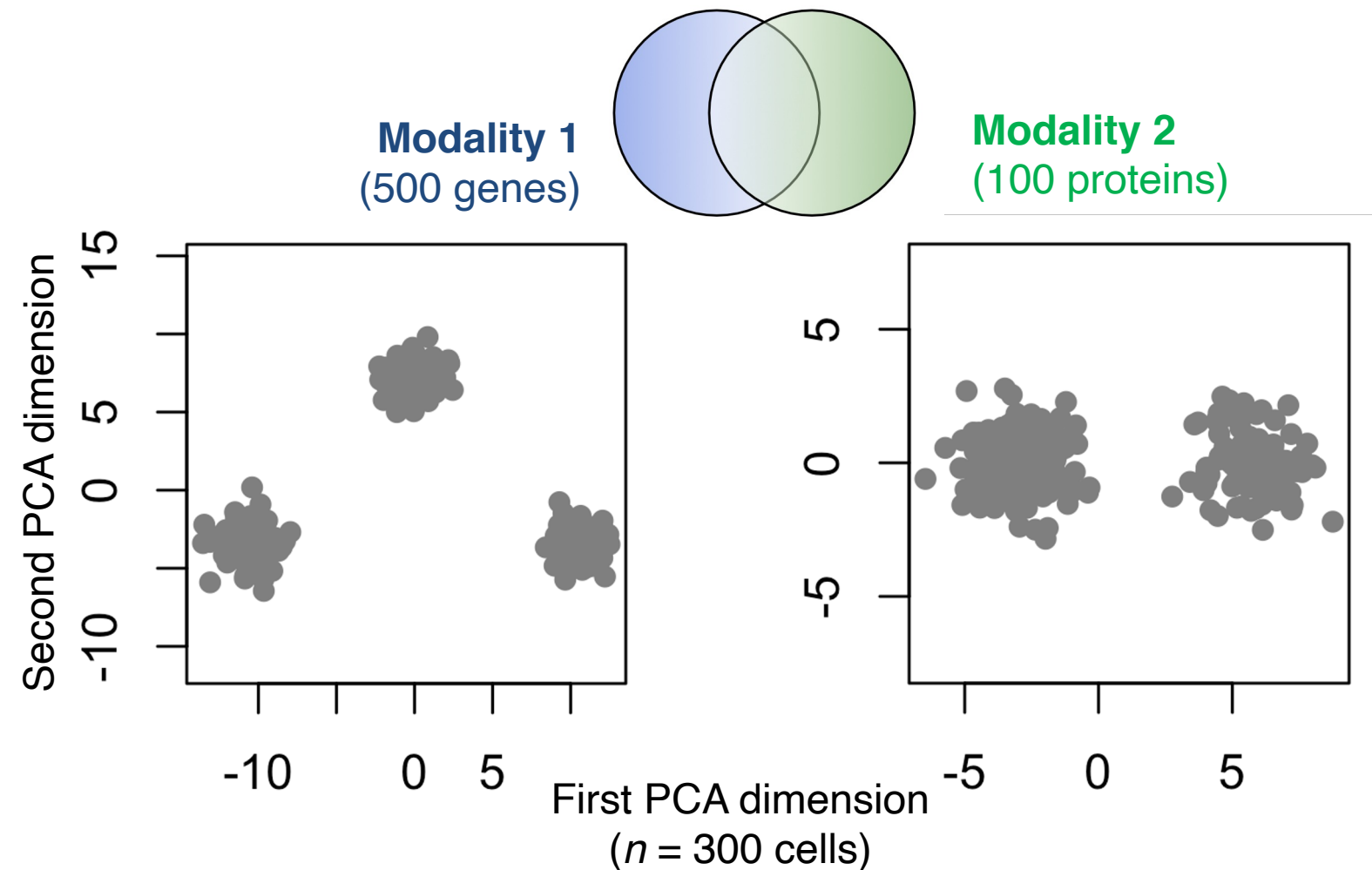
To start thinking about how to define geometry shared/unique to each modality, consider a toy example where one modality has more “information” than the other.

- Geometry = “Information” \approx “Density clustering”



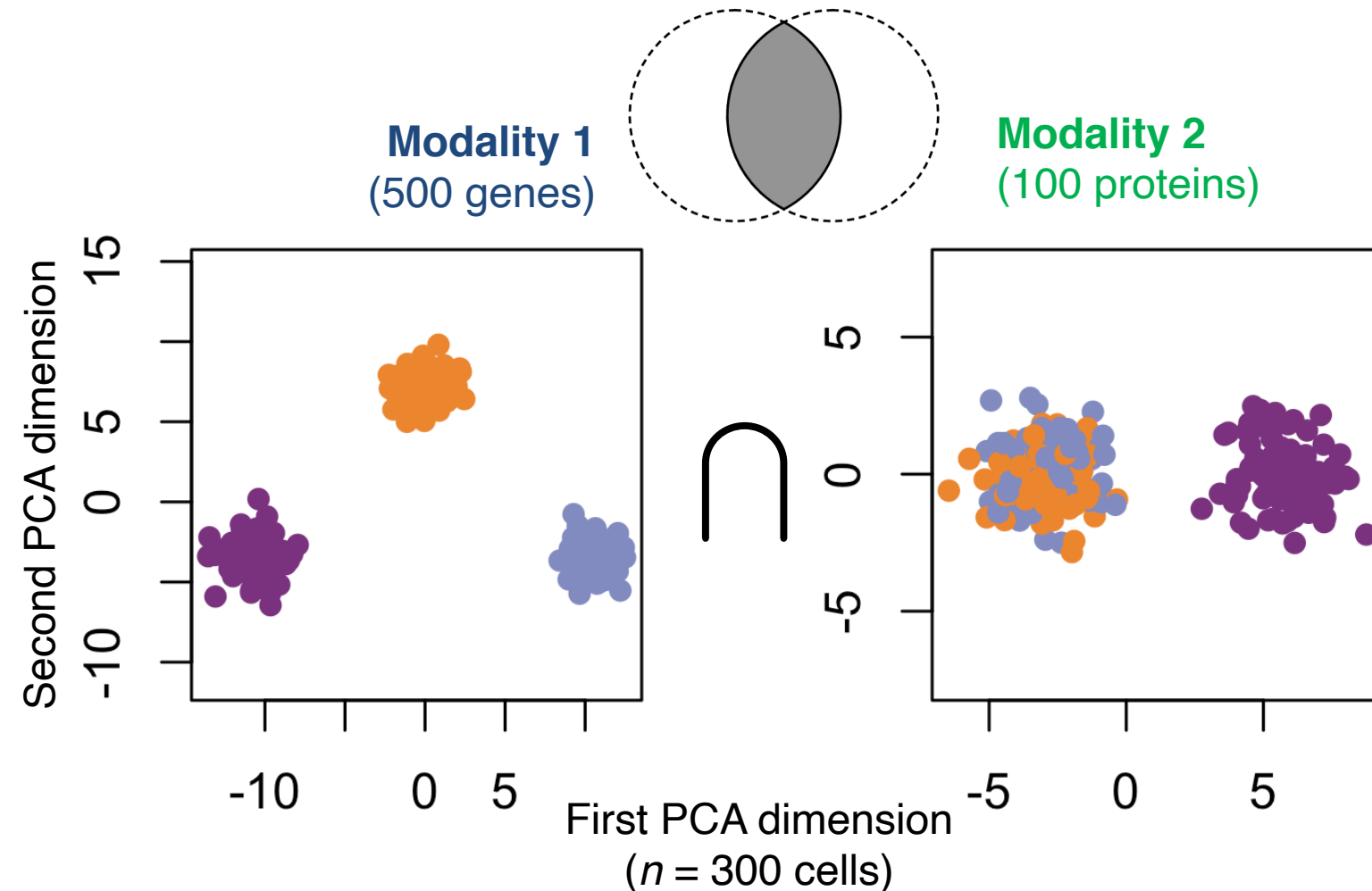
To start thinking about how to define geometry shared/unique to each modality, consider a toy example where one modality has more “information” than the other.

- Geometry = “Information” \approx “Density clustering”



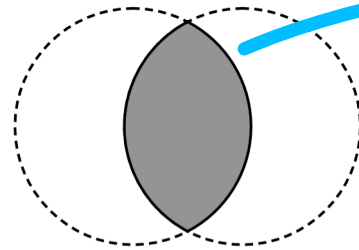
To start thinking about how to define geometry shared/unique to each modality, consider a toy example where one modality has more “information” than the other.

- Geometry \approx “Information” \approx “Density clustering”



To start thinking about how to define geometry shared/unique to each modality, consider a toy example where one modality has more “information” than the other.

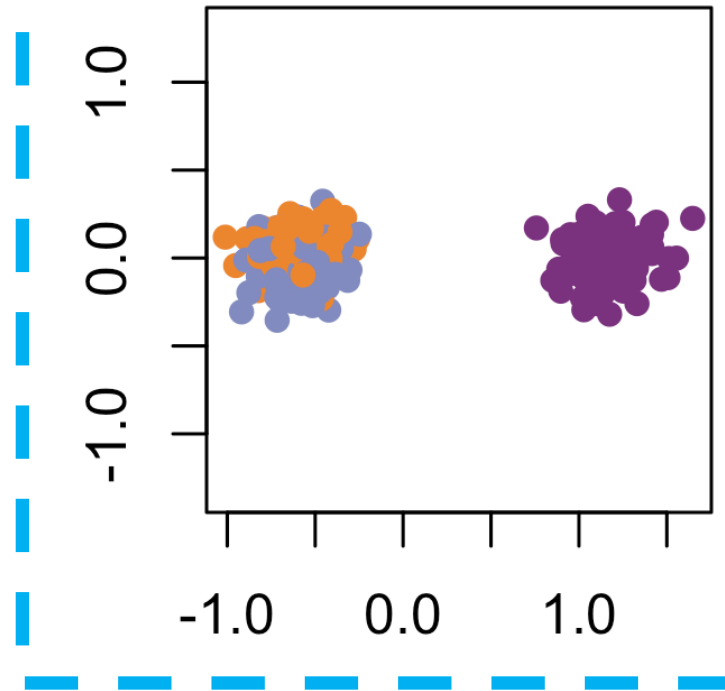
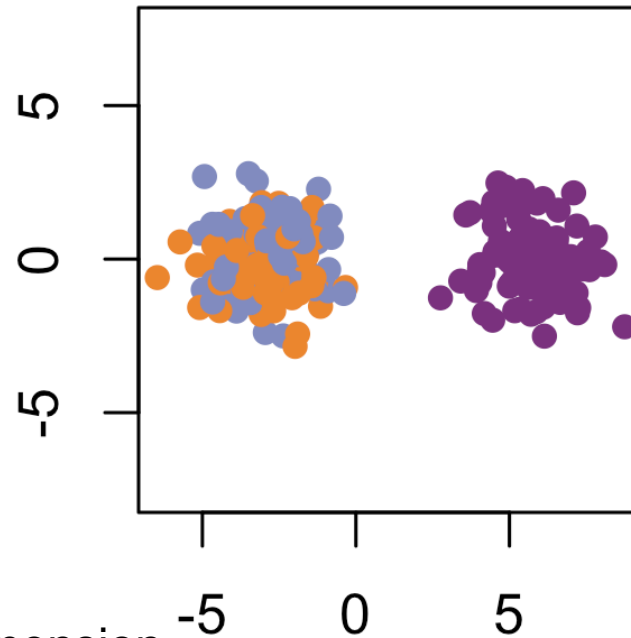
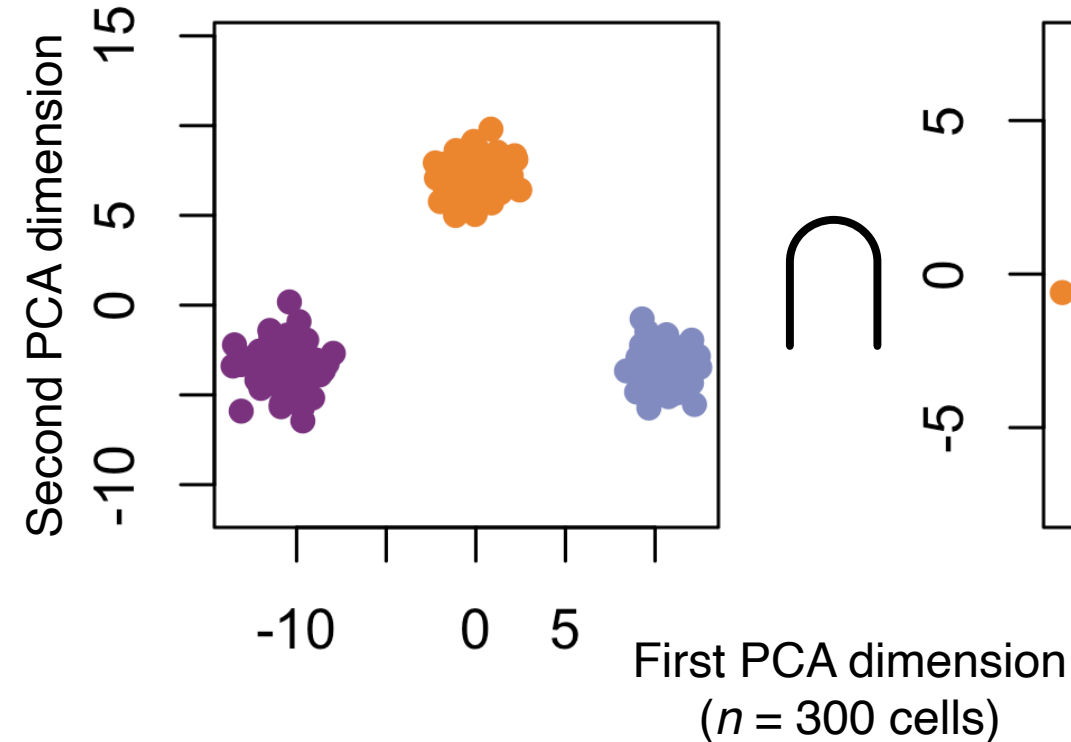
- Geometry \approx “Information” \approx “Density clustering”



Modality 1
(500 genes)

Modality 2
(100 proteins)

Desired 2D embedding



Many existing multi-modal dimension reduction methods focus minimizing reconstruction error, but these encapsulate the “union” of information.

Many existing multi-modal dimension reduction methods focus minimizing reconstruction error, but these encapsulate the “union” of information.

(a phrase we coined for
multimodal matrix factorization)

Many existing multi-modal dimension reduction methods focus minimizing reconstruction error, but these encapsulate the “union” of information.

Consider Consensus PCA (Wold et al., 1987):

$$\hat{L} = \arg \min_{L \in \mathbb{R}^{(p_1 + p_2) \times r}} \left\| \begin{bmatrix} X^{(1)} & ; & X^{(2)} \end{bmatrix} - \begin{bmatrix} X^{(1)} & ; & X^{(2)} \end{bmatrix} L L^\top \right\|_F^2$$

Resulting embedding via SVD:

$$\begin{bmatrix} X^{(1)} & ; & X^{(2)} \end{bmatrix} = U D V^\top$$

Many existing multi-modal dimension reduction methods focus minimizing reconstruction error, but these encapsulate the “union” of information.

Consider Consensus PCA (Wold et al., 1987):

$$\hat{L} = \arg \min_{L \in \mathbb{R}^{(p_1 + p_2) \times r}} \left\| \begin{bmatrix} X^{(1)} & ; & X^{(2)} \end{bmatrix} - \begin{bmatrix} X^{(1)} & ; & X^{(2)} \end{bmatrix} L L^\top \right\|_F^2$$

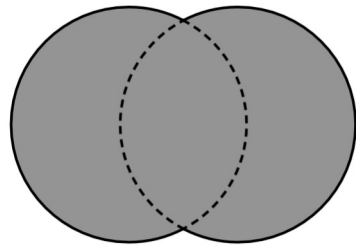
Resulting embedding via SVD:

$$\begin{bmatrix} X^{(1)} & ; & X^{(2)} \end{bmatrix} = \underbrace{U D V^\top}$$

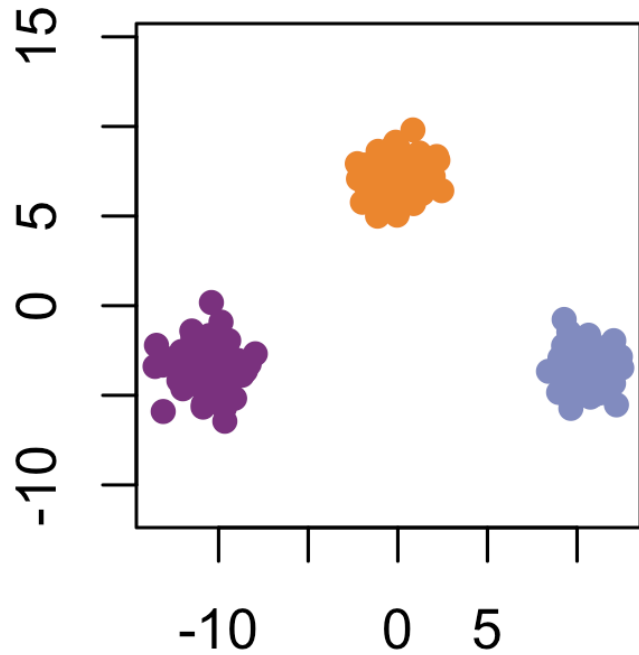
Low-dimensional embedding

Observation: The embedding combines the axes of variation from each modality (i.e. the “union”).

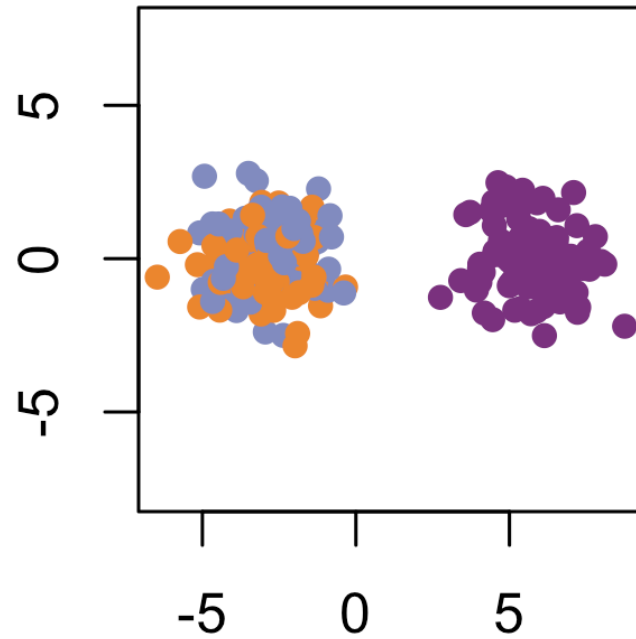
Let's see the "union of information" in toy example:



Modality 1
(500 genes)



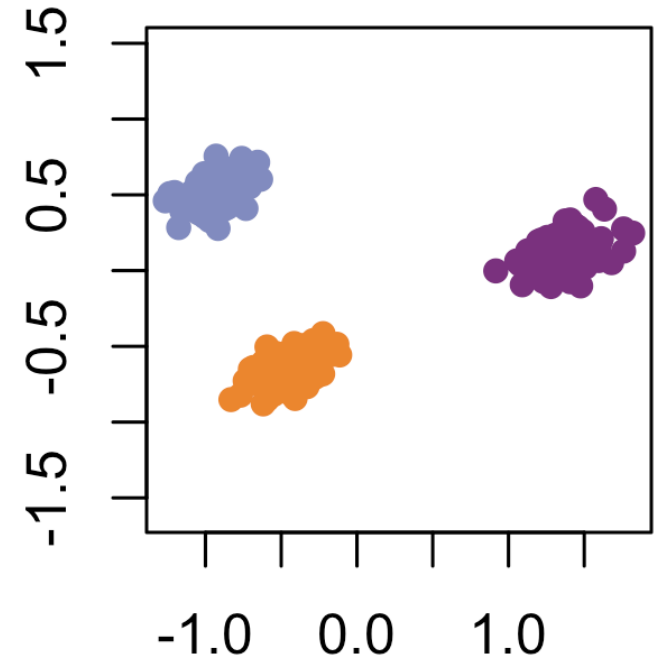
Modality 2
(100 proteins)



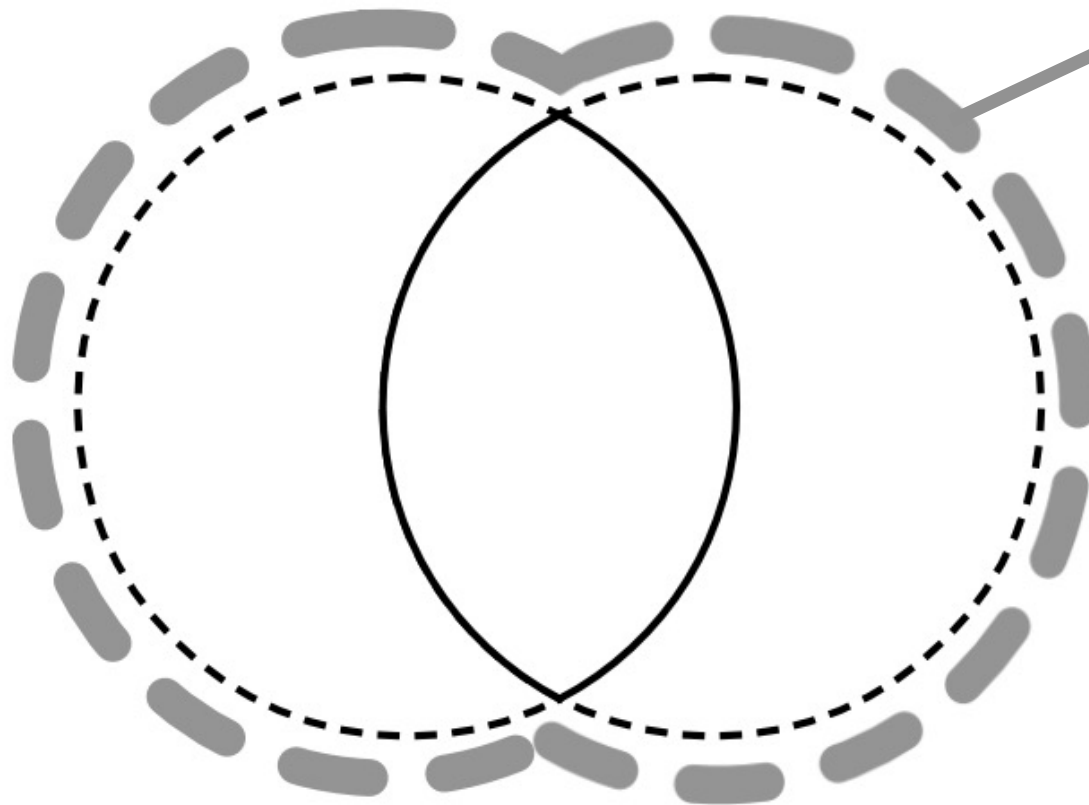
\cup

\parallel

**Consensus
PCA's 2-dim.
embedding**



Our contribution: A matrix factorization for multi-modal data that separates of “intersection” and “unique” of information, as opposed many existing methods that represent the “union” of information.

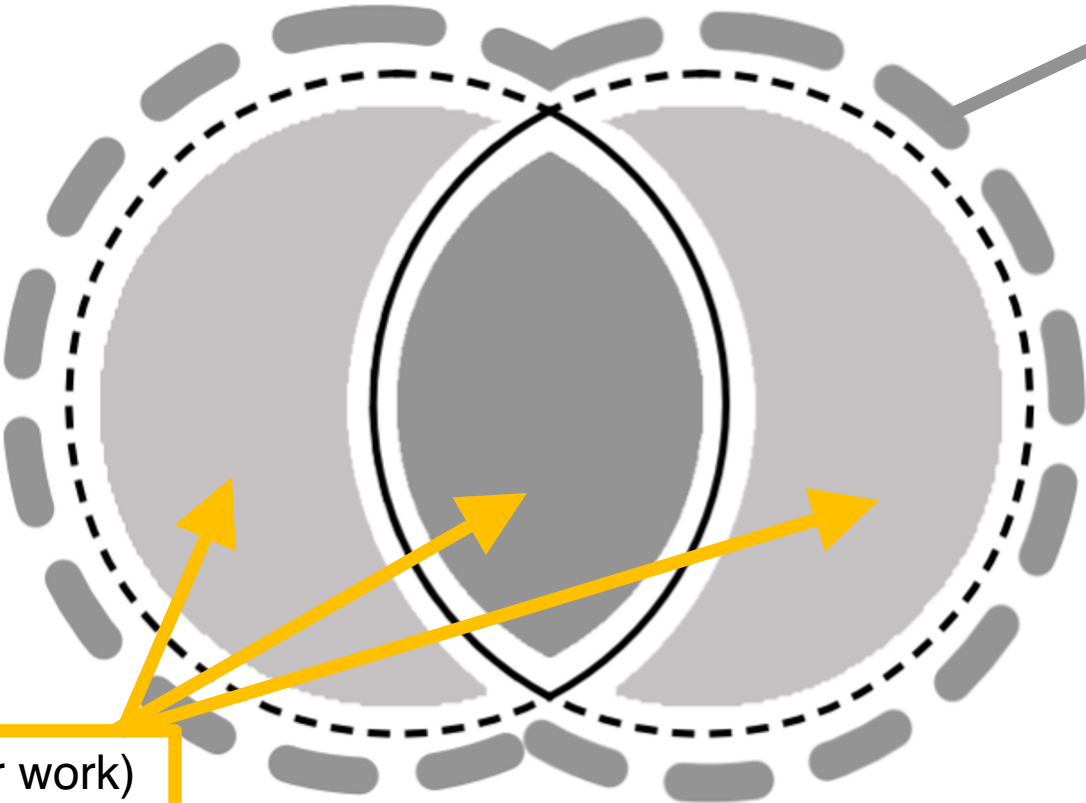


- [Consensus PCA](#) (Word et al., 1987)
- [scAI](#) (Nie et al., 2020)
- [MOFA+](#) (Stegle et al., 2020)
- [JSNMF](#) (Ma et al., 2022)

- [JIVE](#) (Nobel et al., 2013)
- [WNN](#) (Satija et al., 2021)
- [Cobalt](#) (Purdom et al., 2021)

Cell embeddings for the “union” geometry: Useful for making an “atlas” across both modalities

Our contribution: A matrix factorization for multi-modal data that separates of “intersection” and “unique” of information, as opposed many existing methods that represent the “union” of information.



- **Consensus PCA** (Word et al., 1987)
- **scAI** (Nie et al., 2020)
- **MOFA+** (Stegle et al., 2020)
- **JSNMF** (Ma et al., 2022)

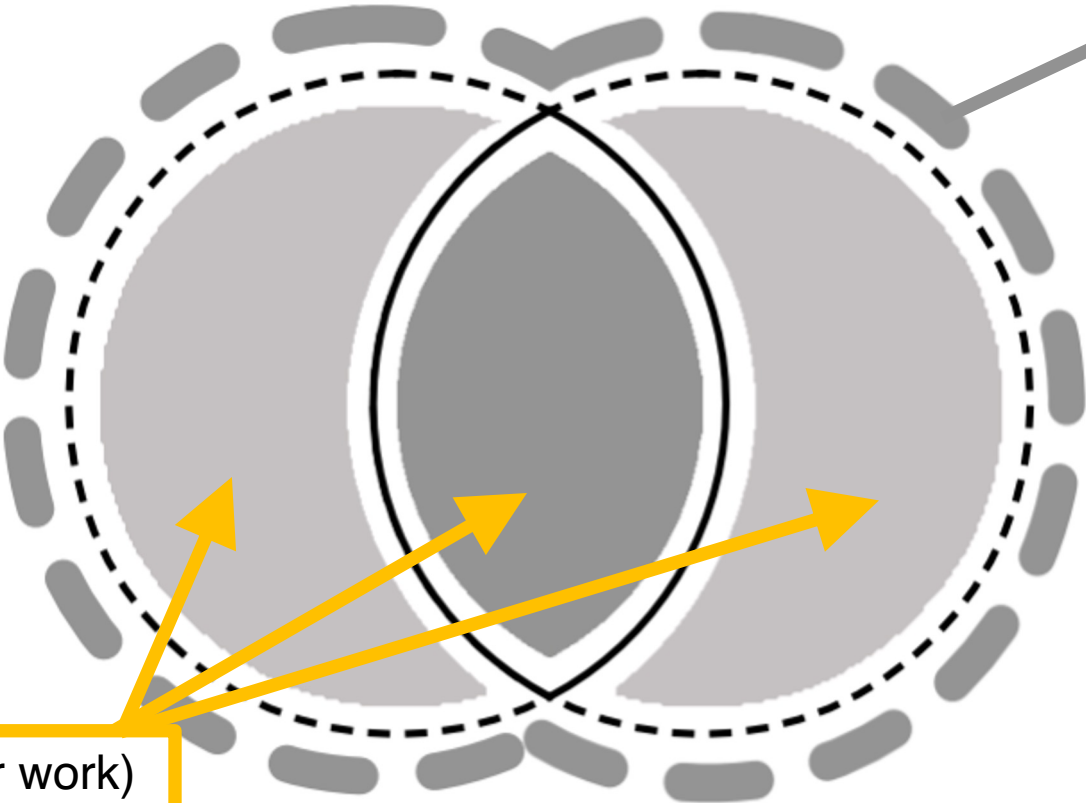
- **JIVE** (Nobel et al., 2013)
- **WNN** (Satija et al., 2021)
- **Cobalt** (Purdom et al., 2021)

Cell embeddings for the “union” geometry: Useful for making an “atlas” across both modalities

Tilted-CCA (Our work)

Cell embedding for the “intersection”/“unique” geometry:
Useful to understand the coordination between modalities

Our contribution: A matrix factorization for multi-modal data that separates of “intersection” and “unique” of information, as opposed many existing methods that represent the “union” of information.



- Consensus PCA (Word et al., 1987)
- scAI (Nie et al., 2020)
- MOFA+ (Stegle et al., 2020)
- JSNMF (Ma et al., 2022)

- JIVE (Nobel et al., 2013)*
- WNN (Satija et al., 2021)*
- Cobalt (Purdom et al., 2021)*

Cell embeddings for the “union” geometry: Useful for making an “atlas” across both modalities

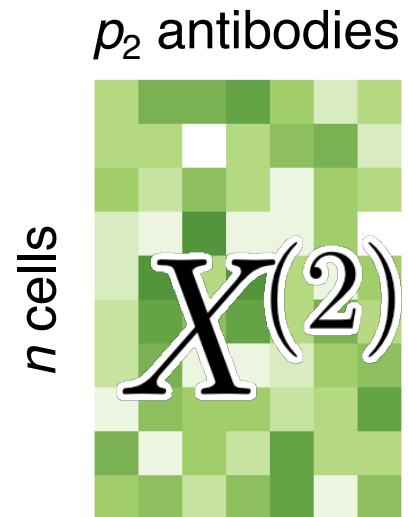
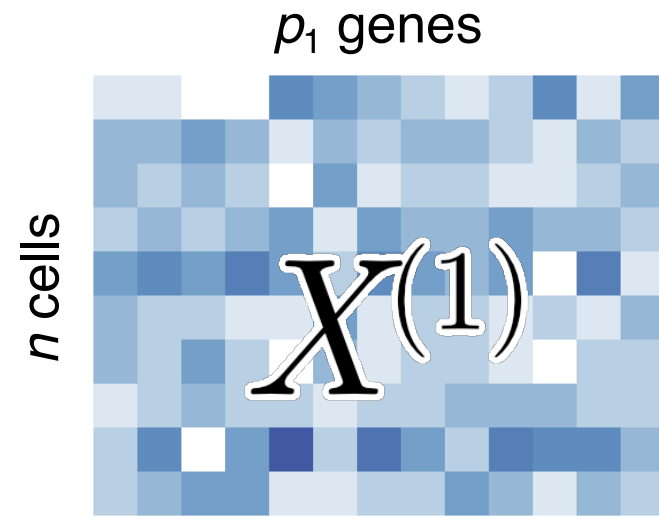
Tilted-CCA (Our work)

Cell embedding for the “intersection”/“unique” geometry:
Useful to understand the coordination between modalities

Statistical method:

New dimension-reduction framework for multi-modal data,
not about denoising but instead the geometry

Proposed model for multi-modal data: Capturing the “shared” geometry



(Assumed to be sufficiently preprocessed)

Proposed model for multi-modal data: Capturing the “shared” geometry

p_1 genes

n cells

$X^{(1)}$

$\approx \left[C + D^{(1)} \right] \times L^{(1)}$

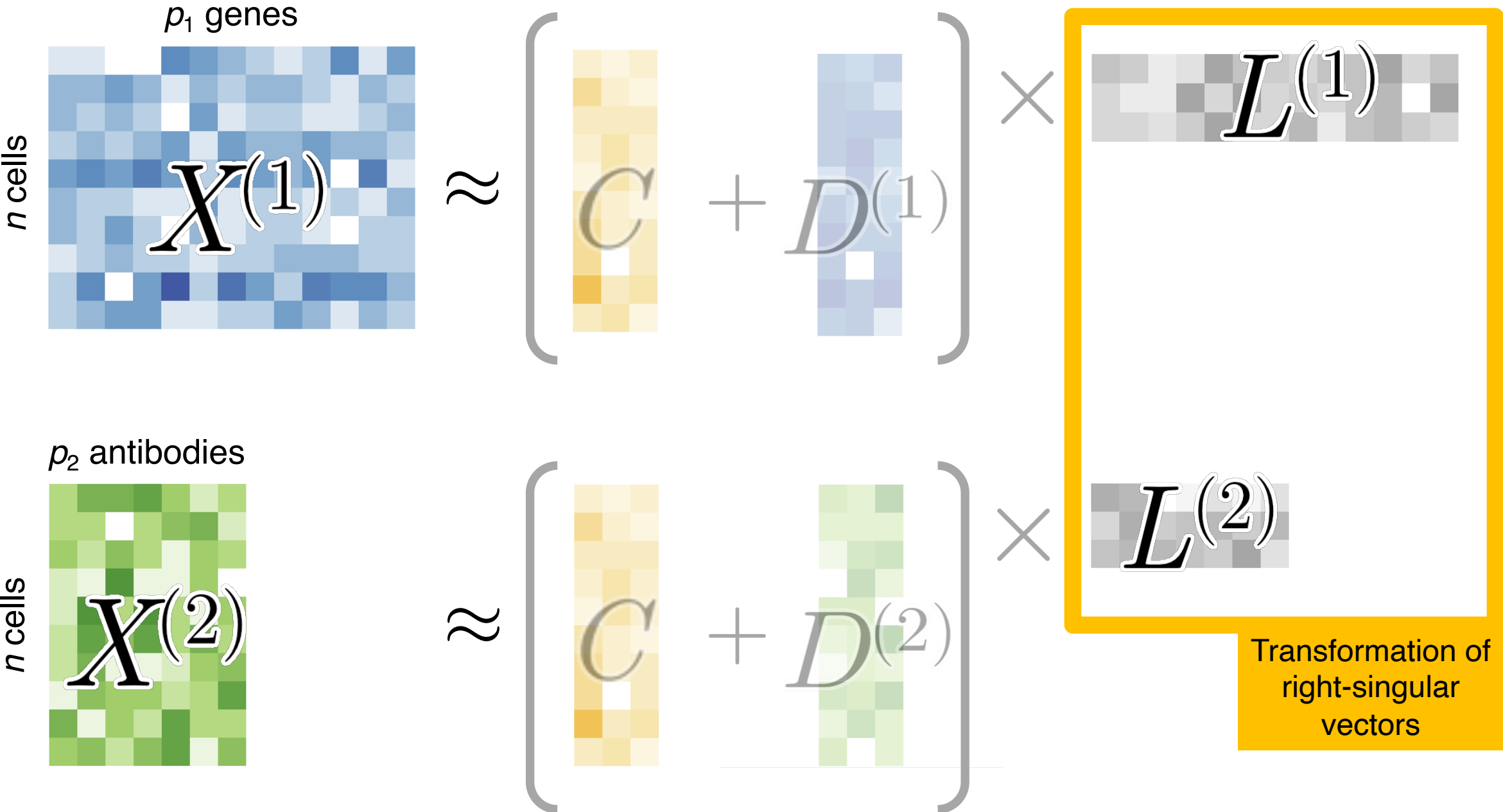
p_2 antibodies

n cells

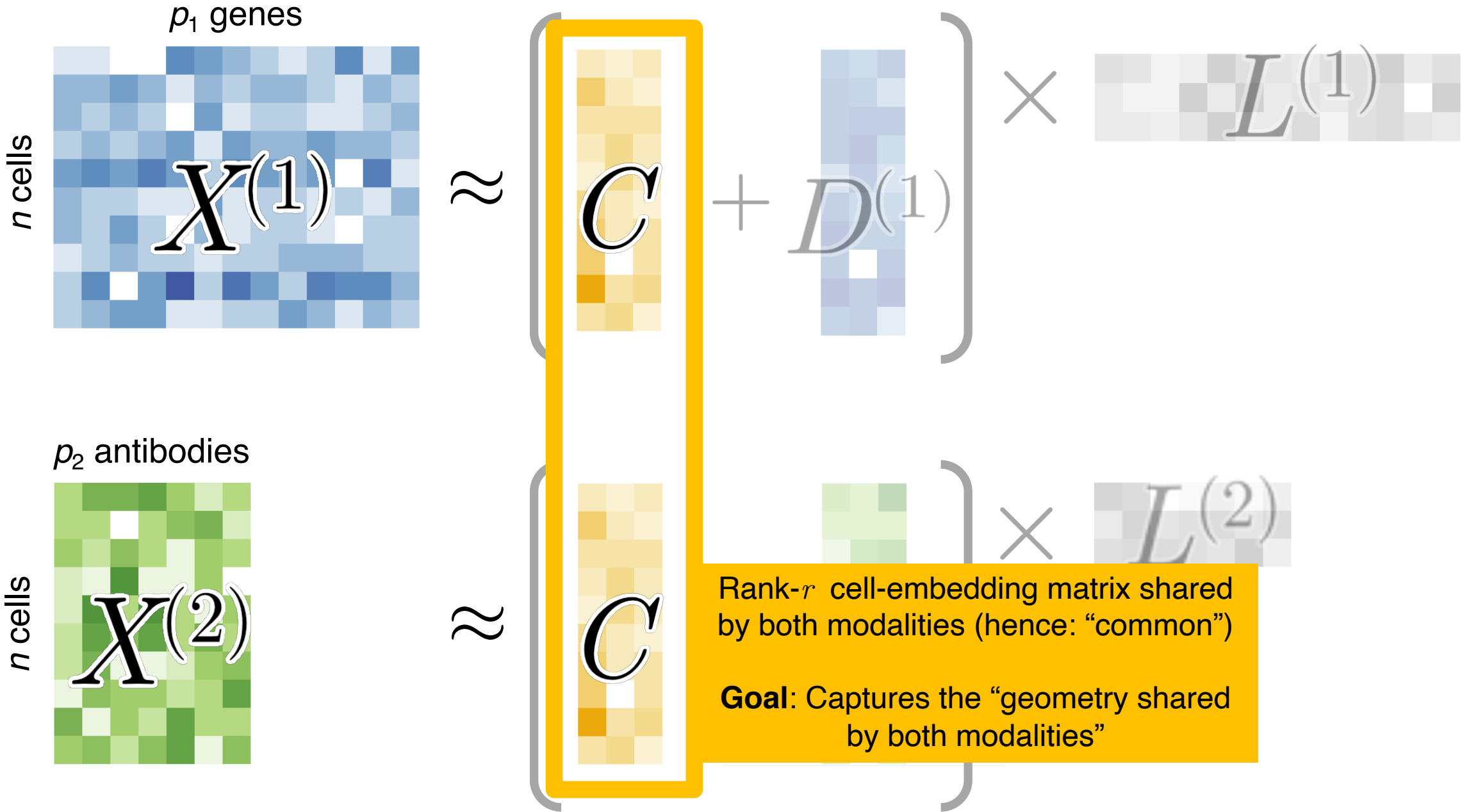
$X^{(2)}$

$\approx \left[C + D^{(2)} \right] \times L^{(2)}$

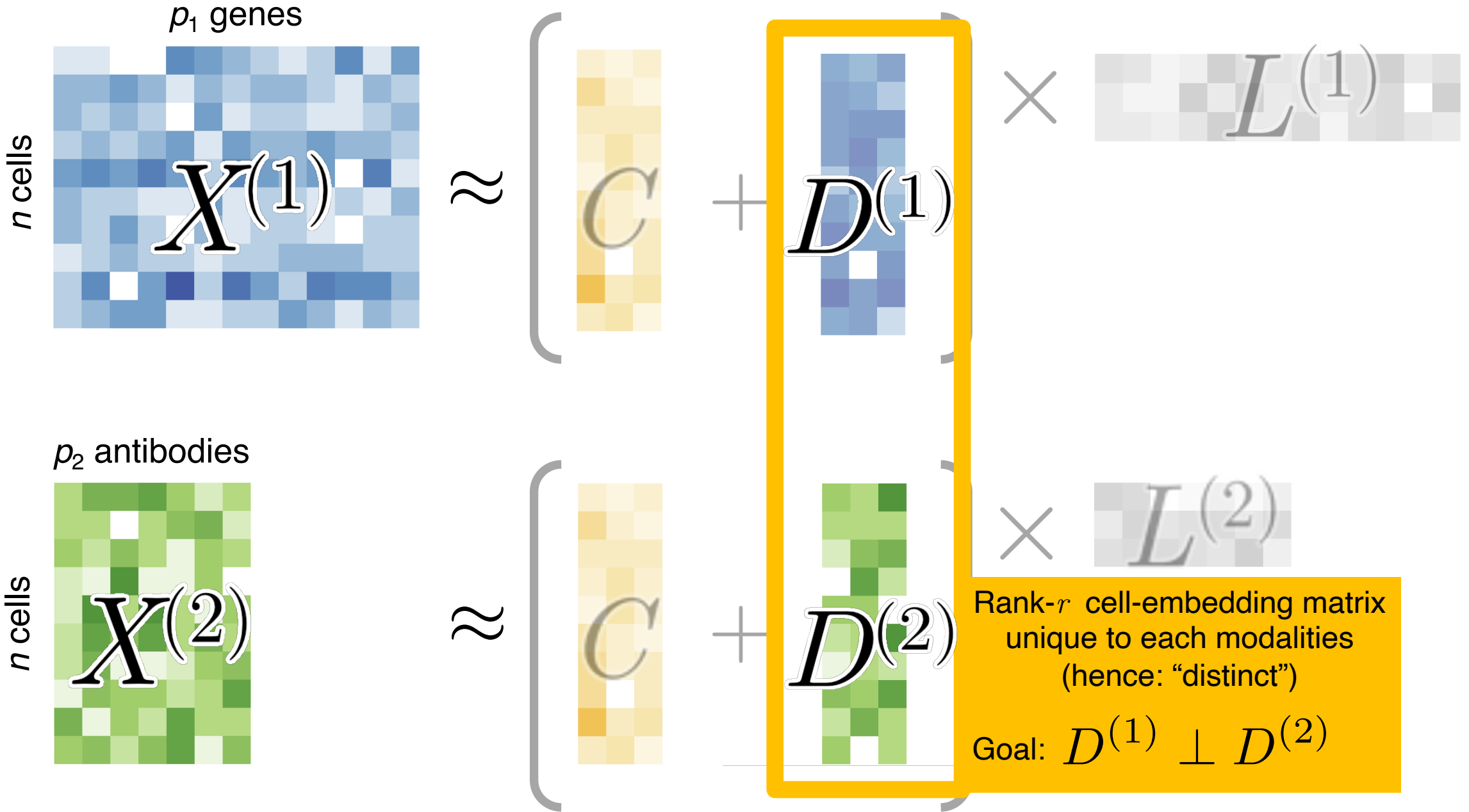
Proposed model for multi-modal data: Capturing the “shared” geometry



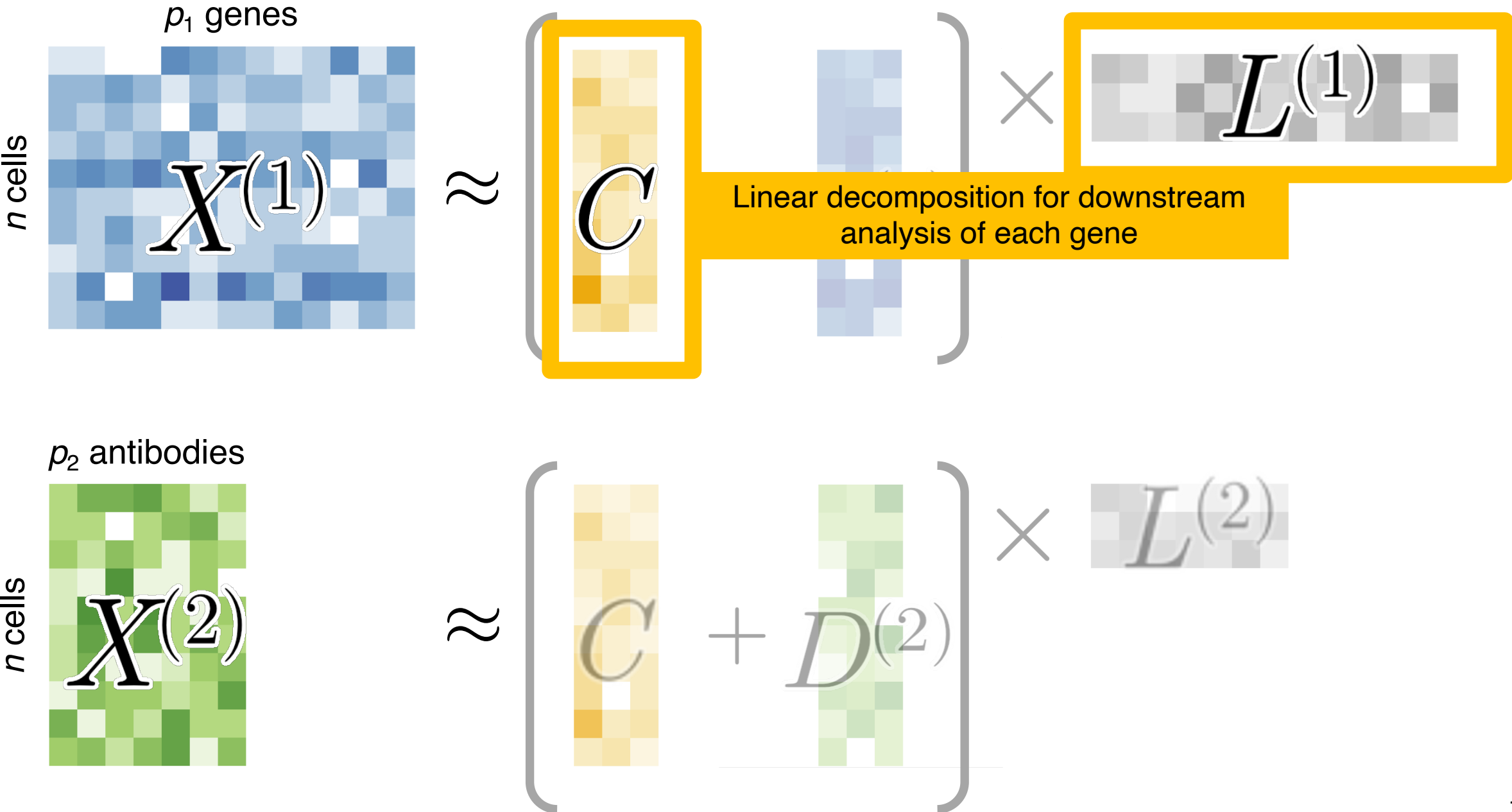
Proposed model for multi-modal data: Capturing the “shared” geometry



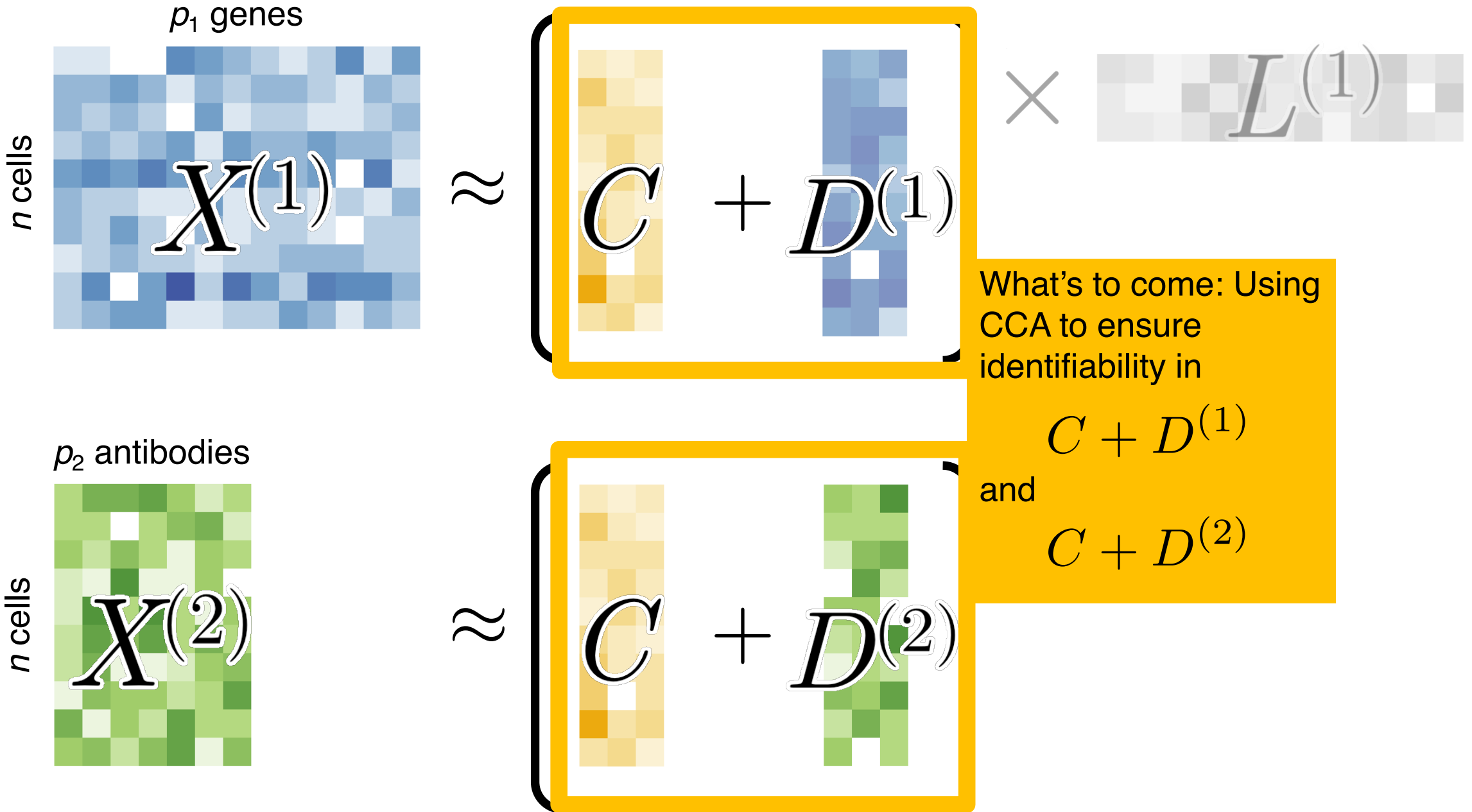
Proposed model for multi-modal data: Capturing the “shared” geometry



Proposed model for multi-modal data: Capturing the “shared” geometry



Proposed model for multi-modal data: Capturing the “shared” geometry



Brief aside: Review of Canonical Correlation Analysis

CCA for the first latent dimension:

Given two modalities, find linear combination of variables that are highly correlated:

$$\{\hat{a}, \hat{b}\} = \arg \max_{\substack{a \in \mathbb{R}^{p_1} \\ b \in \mathbb{R}^{p_2}}} \text{Corr}\left(X^{(1)}a, X^{(2)}b\right)$$

General multi-dimensional CCA:

$$\{\hat{A}, \hat{B}\} = \arg \max_{\substack{A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r \\ B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r}} \text{Tr}\left(A^\top \left(X^{(1)}\right)^\top X^{(2)} B\right)$$

Brief aside: Review of Canonical Correlation Analysis

CCA for the first latent dimension:

Given two modalities, find linear combination of variables that are highly correlated:

$$\{\hat{a}, \hat{b}\} = \arg \max_{\substack{a \in \mathbb{R}^{p_1} \\ b \in \mathbb{R}^{p_2}}} \text{Corr} \left(\underbrace{X^{(1)} a, X^{(2)} b}_{\text{Pair of canonical score vectors}} \right)$$

General multi-dimensional CCA:

$$\{\hat{A}, \hat{B}\} = \arg \max \text{Tr} \left(A^\top (X^{(1)})^\top X^{(2)} B \right)$$
$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$
$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

Brief aside: Review of Canonical Correlation Analysis

CCA for the first latent dimension:

Given two modalities, find linear combination of variables that are highly correlated:

$$\{\hat{a}, \hat{b}\} = \arg \max_{\substack{a \in \mathbb{R}^{p_1} \\ b \in \mathbb{R}^{p_2}}} \text{Corr}\left(X^{(1)} a, X^{(2)} b\right)$$

General multi-dimensional CCA:

$$\{\hat{A}, \hat{B}\} = \arg \max \text{Tr}\left(A^\top (X^{(1)})^\top X^{(2)} B\right)$$

$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$

$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

Orthogonality of
canonical scores

Brief aside: Review of Canonical Correlation Analysis

CCA for the first latent dimension:

Given two modalities, find linear combination of variables that are highly correlated:

$$\{\hat{a}, \hat{b}\} = \arg \max_{\substack{a \in \mathbb{R}^{p_1} \\ b \in \mathbb{R}^{p_2}}} \text{Corr}\left(X^{(1)}a, X^{(2)}b\right)$$

General multi-dimensional CCA:

$$\{\hat{A}, \hat{B}\} = \arg \max_{A, B} \text{Tr}\left(A^\top (X^{(1)})^\top X^{(2)} B\right)$$

$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$

$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

At optimality:
Orthogonality of
canonical scores
across modalities

Decomposition based on CCA:

Decomposition based on CCA:

$$\{\hat{A}, \hat{B}\} = \arg \max \operatorname{Tr} \left(A^\top (X^{(1)})^\top X^{(2)} B \right)$$

$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$

$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

Let: $Z^{(1)} = X^{(1)} \hat{A}$

$Z^{(2)} = X^{(2)} \hat{B}$
 $\mathbb{R}^{n \times r}$

Decomposition based on CCA:

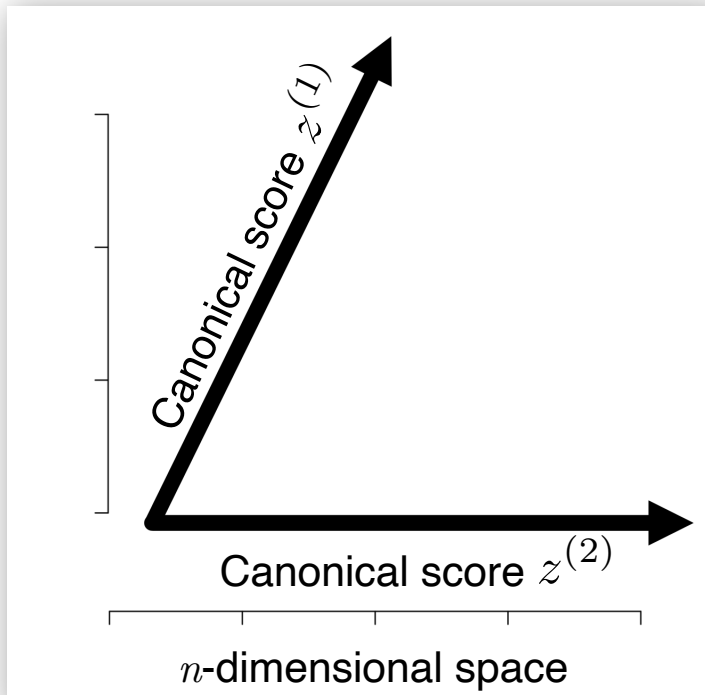
$$\{\hat{A}, \hat{B}\} = \arg \max \operatorname{Tr} \left(A^\top (X^{(1)})^\top X^{(2)} B \right)$$

$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$

$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

$$\text{Let: } Z^{(1)} = X^{(1)} \hat{A}$$

$$Z^{(2)} = X^{(2)} \hat{B}$$



The 2D hyperplane in
 n -dimensional space

Decomposition based on CCA:

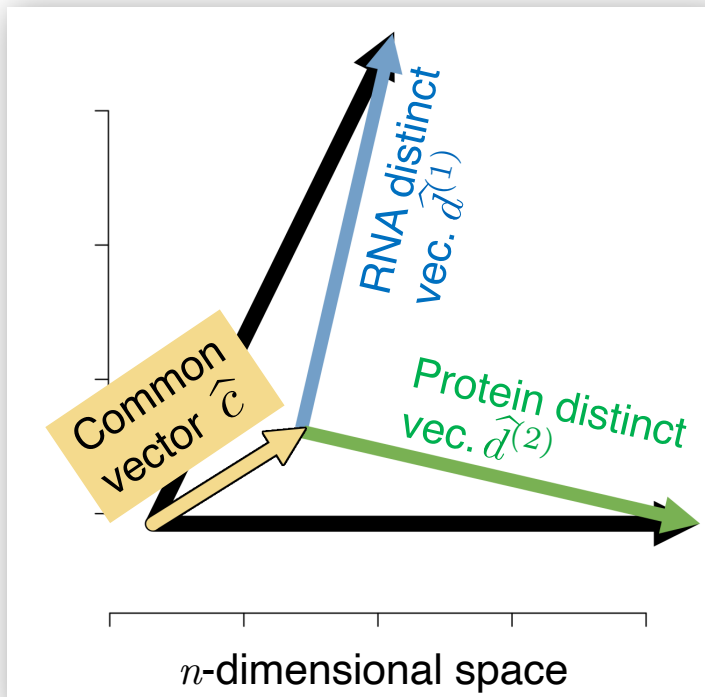
$$\{\hat{A}, \hat{B}\} = \arg \max \text{Tr} \left(A^\top (X^{(1)})^\top X^{(2)} B \right)$$

$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$

$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

$$\text{Let: } Z^{(1)} = X^{(1)} \hat{A}$$

$$Z^{(2)} = X^{(2)} \hat{B}$$



Decomposition:

$$z^{(1)} = \hat{c} + \hat{d}^{(1)}$$

$$z^{(2)} = \hat{c} + \hat{d}^{(2)}$$

(quantifying what
"common" means)

Decomposition based on CCA:

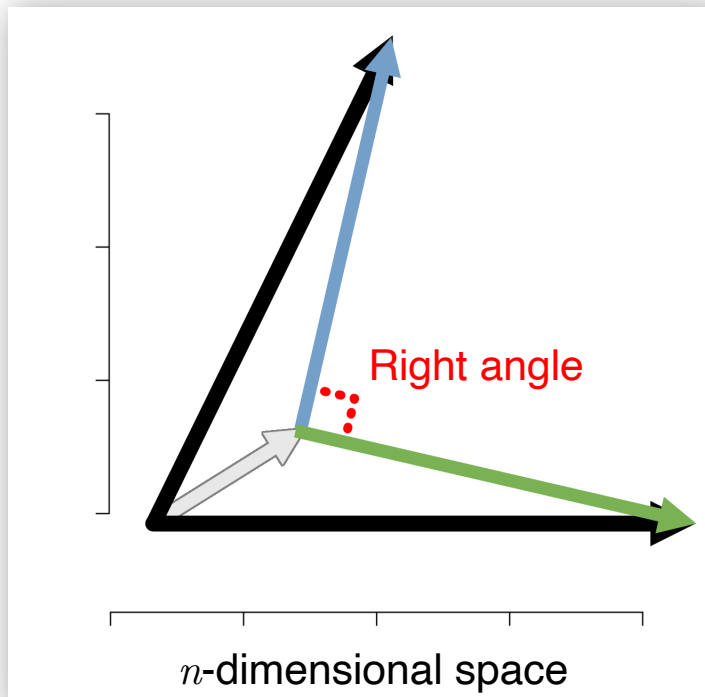
$$\{\hat{A}, \hat{B}\} = \arg \max \operatorname{Tr} \left(A^\top (X^{(1)})^\top X^{(2)} B \right)$$

$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$

$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

$$\text{Let: } Z^{(1)} = X^{(1)} \hat{A}$$

$$Z^{(2)} = X^{(2)} \hat{B}$$



Constrained to:

$$\hat{d}^{(1)} \perp \hat{d}^{(2)}$$

(quantifying what
"distinct" means)

Decomposition based on CCA:

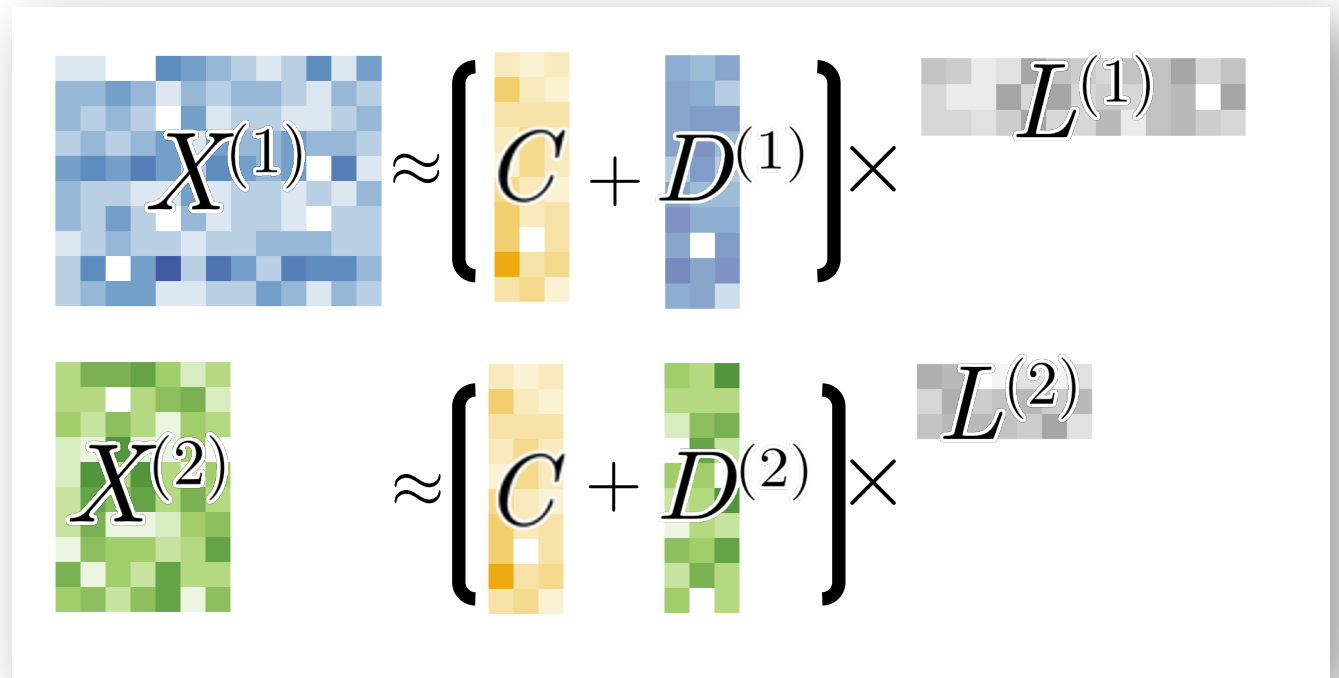
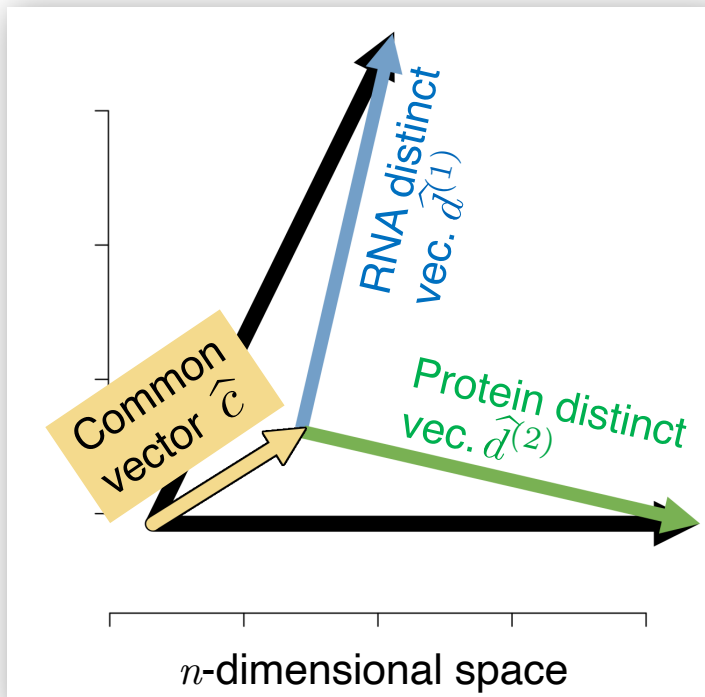
$$\{\hat{A}, \hat{B}\} = \arg \max \text{Tr} \left(A^\top (X^{(1)})^\top X^{(2)} B \right)$$

$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$

$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

$$\text{Let: } Z^{(1)} = X^{(1)} \hat{A}$$

$$Z^{(2)} = X^{(2)} \hat{B}$$



Decomposition based on CCA:

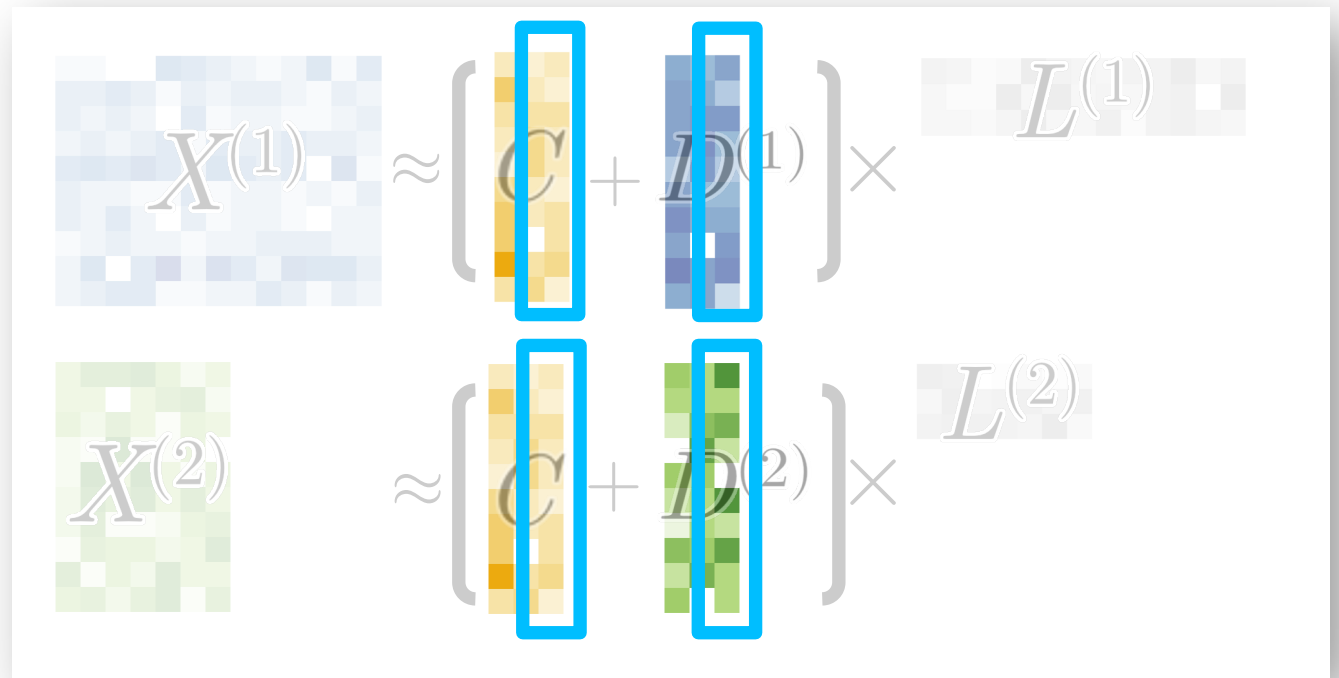
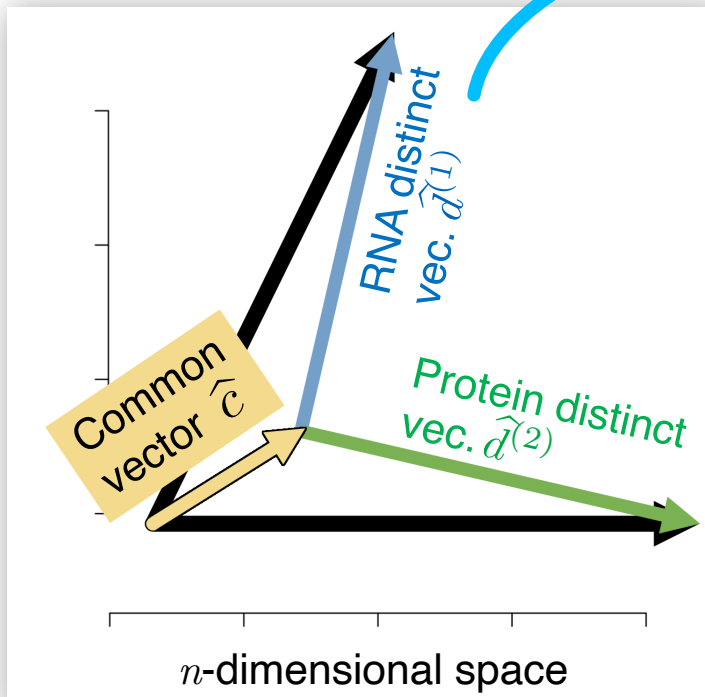
$$\{\hat{A}, \hat{B}\} = \arg \max \text{Tr} \left(A^\top (X^{(1)})^\top X^{(2)} B \right)$$

$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$

$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

$$\text{Let: } Z^{(1)} = X^{(1)} \hat{A}$$

$$Z^{(2)} = X^{(2)} \hat{B}$$



Decomposition based on CCA:

$$\{\hat{A}, \hat{B}\} = \arg \max \text{Tr} \left(A^\top (X^{(1)})^\top X^{(2)} B \right)$$

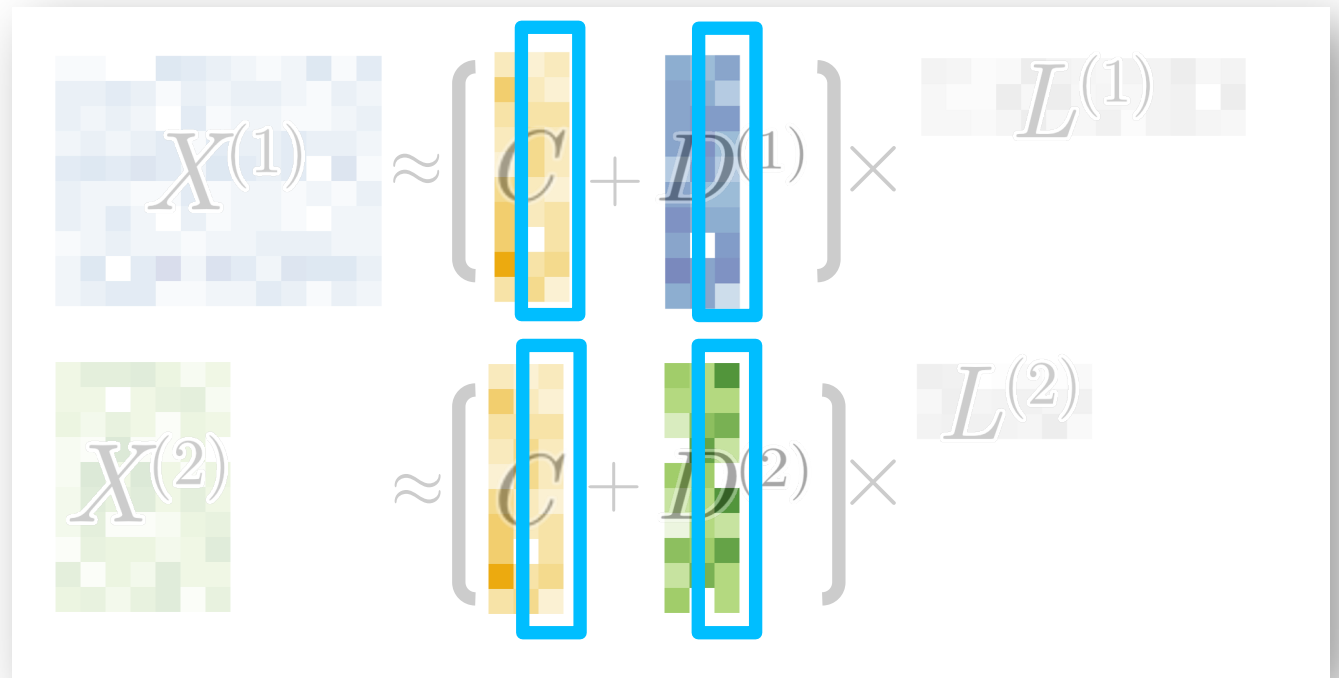
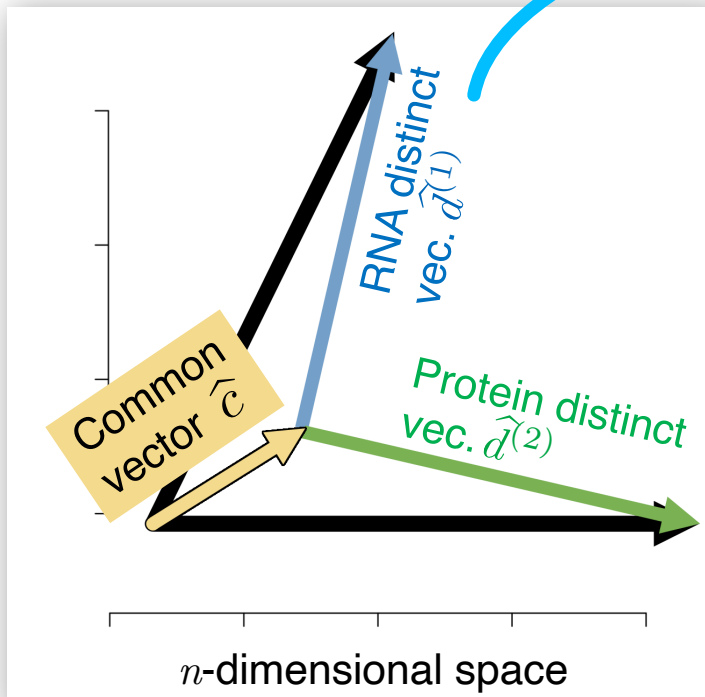
$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$

$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

Identifiability ensured thanks to CCA's properties

$$\text{Let: } Z^{(1)} = X^{(1)} \hat{A}$$

$$Z^{(2)} = X^{(2)} \hat{B}$$



Decomposition based on CCA

$$\{\hat{A}, \hat{B}\} = \arg \max \text{Tr} \left(A^\top (X^{(1)})^\top X^{(2)} B \right)$$

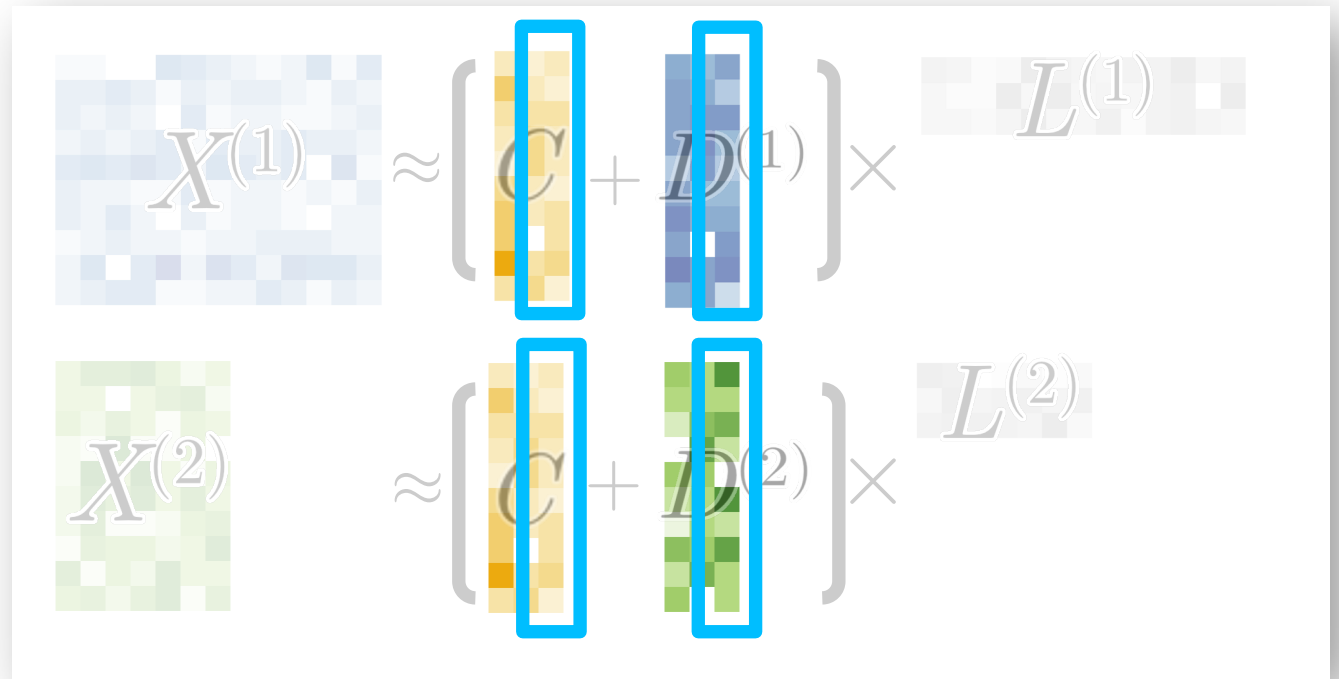
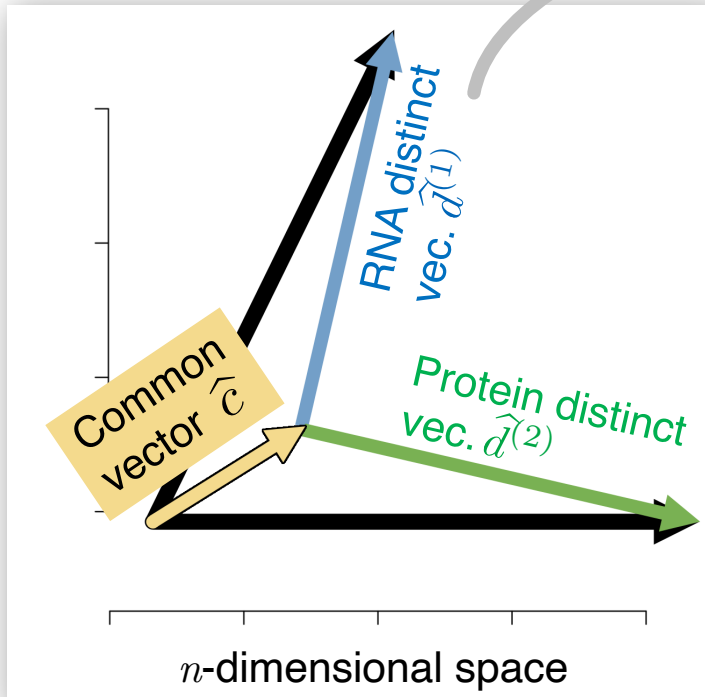
$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$

$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

Identifiability ensured thanks to CCA's properties

$$\text{Let: } Z^{(1)} = X^{(1)} \hat{A}$$

$$Z^{(2)} = X^{(2)} \hat{B}$$



✓ Decomp. for each latent dim.

✓ Ortho. distinct vectors

? Common vector shared geometry

Decomposition based on CCA

$$\{\hat{A}, \hat{B}\} = \arg \max \operatorname{Tr} \left(A^\top (X^{(1)})^\top X^{(2)} B \right)$$

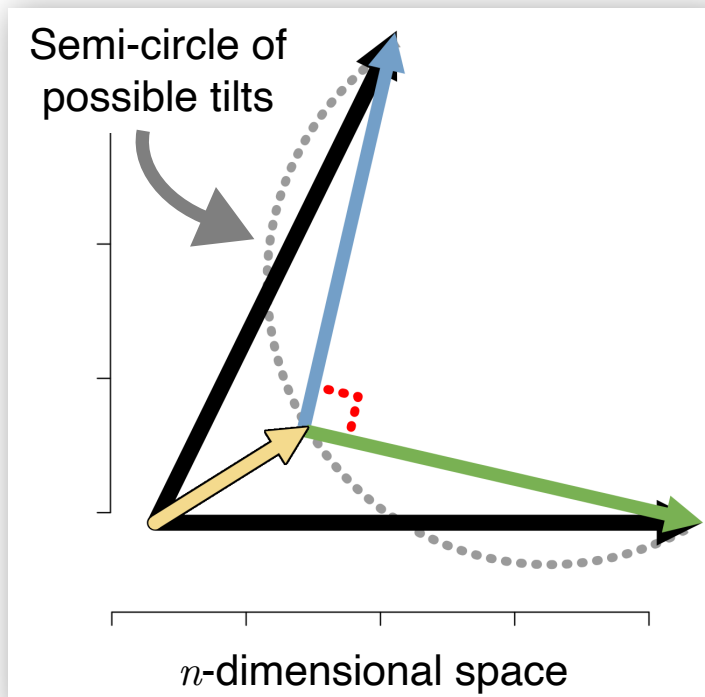
$$A \in \mathbb{R}^{p_1 \times r} : A^\top \Sigma^{(1)} A = I_r$$

$$B \in \mathbb{R}^{p_2 \times r} : B^\top \Sigma^{(2)} B = I_r$$

Identifiability ensured thanks to CCA's properties

$$\text{Let: } Z^{(1)} = X^{(1)} \hat{A}$$

$$Z^{(2)} = X^{(2)} \hat{B}$$

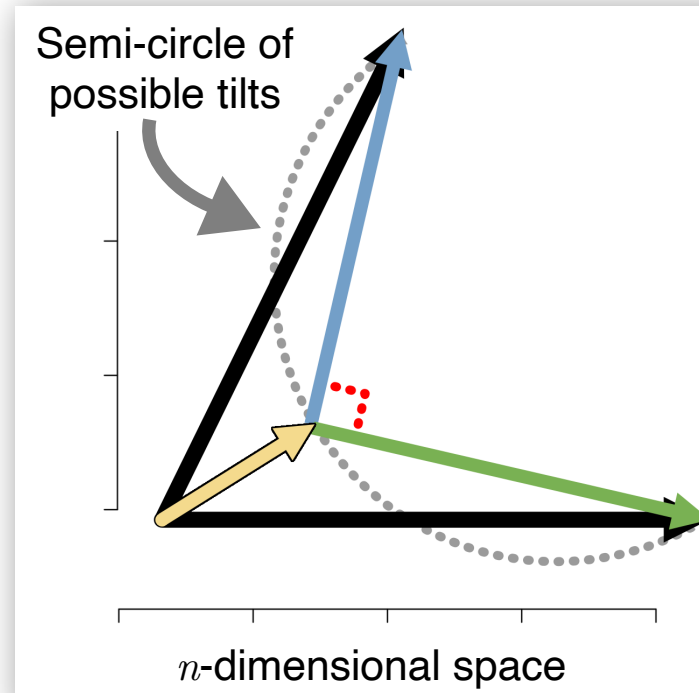


Our insight:

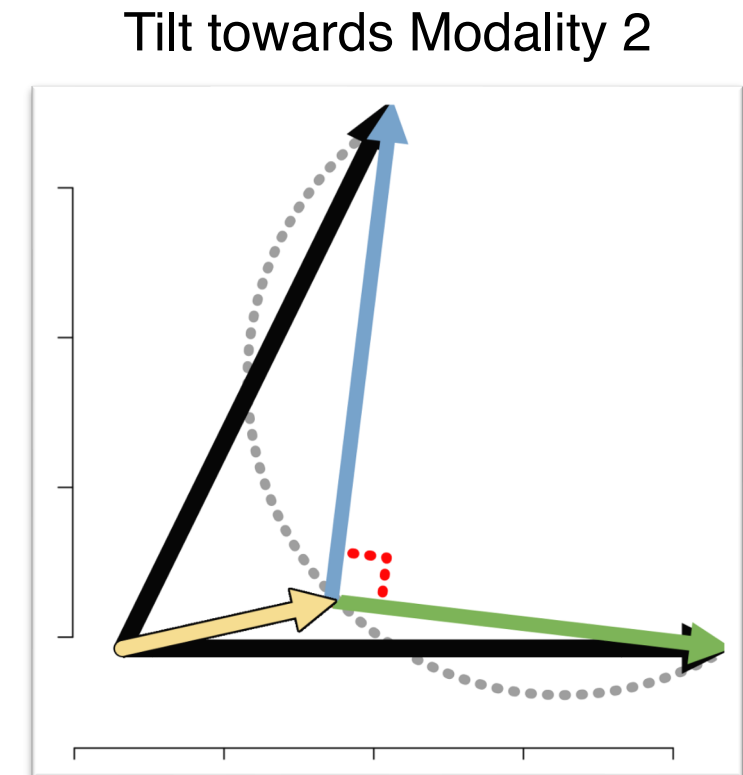
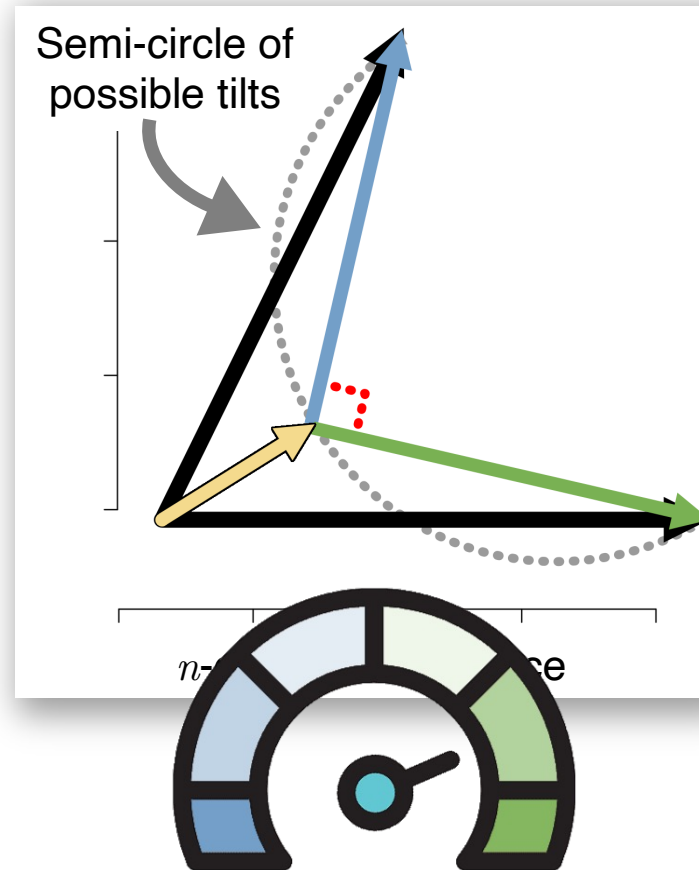
The common vector can fall anywhere along this **semi-circle** and still retain orthogonal distinct vectors.

(One extra degree of freedom per dimension)

Our novelty (combining ideas in CCA with geometry): The tilt of the common vector (if chosen appropriately) allows the common embedding to reflect the shared geometry between both modalities.



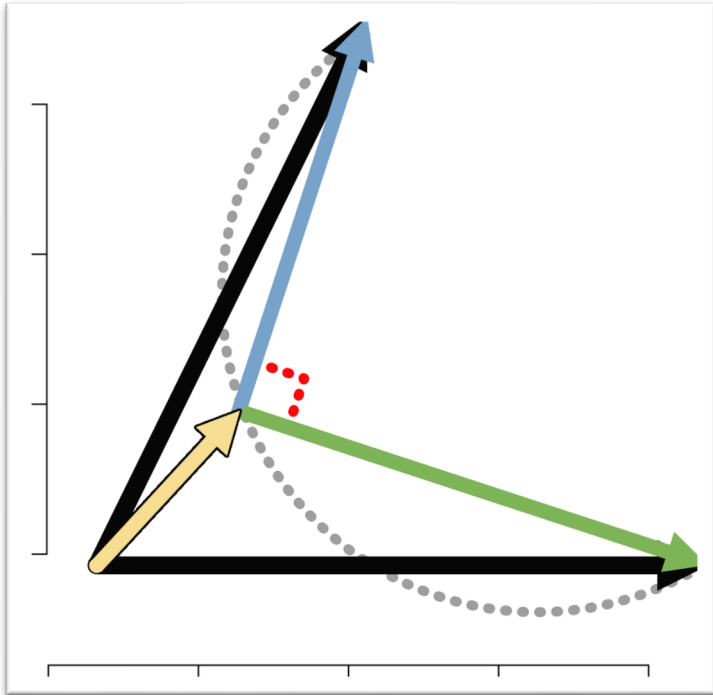
Our novelty (combining ideas in CCA with geometry): The tilt of the common vector (if chosen appropriately) allows the common embedding to reflect the shared geometry between both modalities.



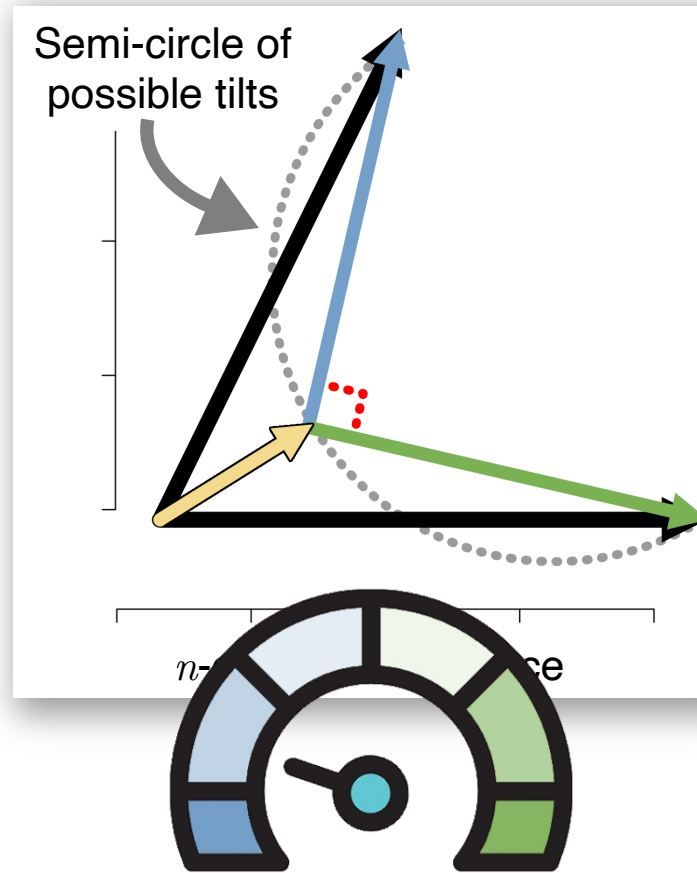
Common embedding has geometry similar to Modality 2

Our novelty (combining ideas in CCA with geometry): The tilt of the common vector (if chosen appropriately) allows the common embedding to reflect the shared geometry between both modalities.

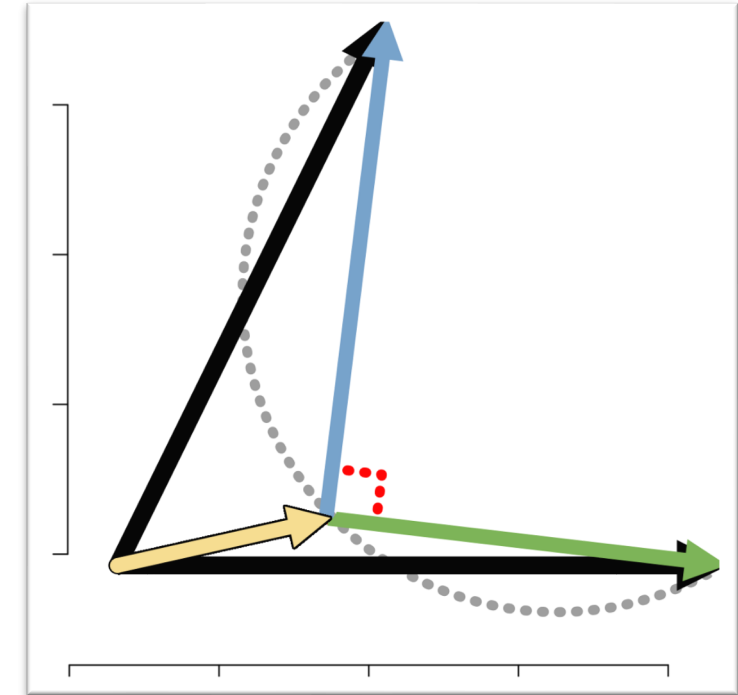
Tilt towards Modality 1



Common embedding has geometry similar to Modality 1

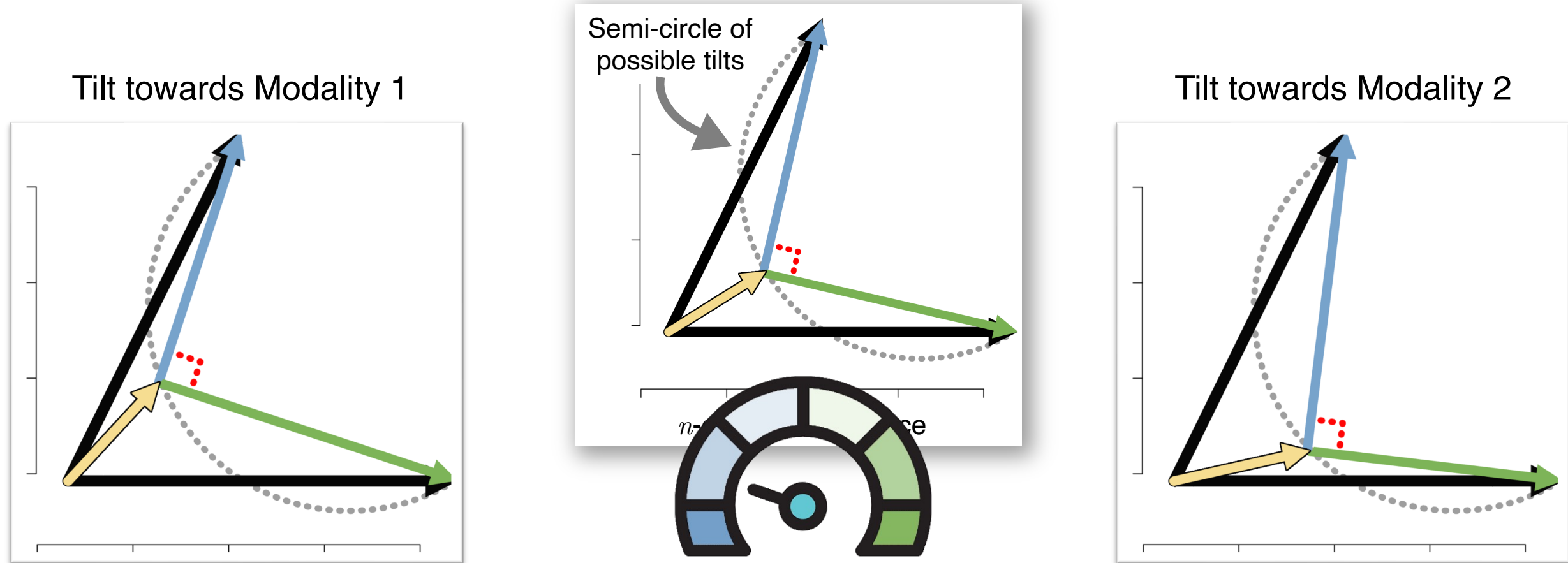


Tilt towards Modality 2



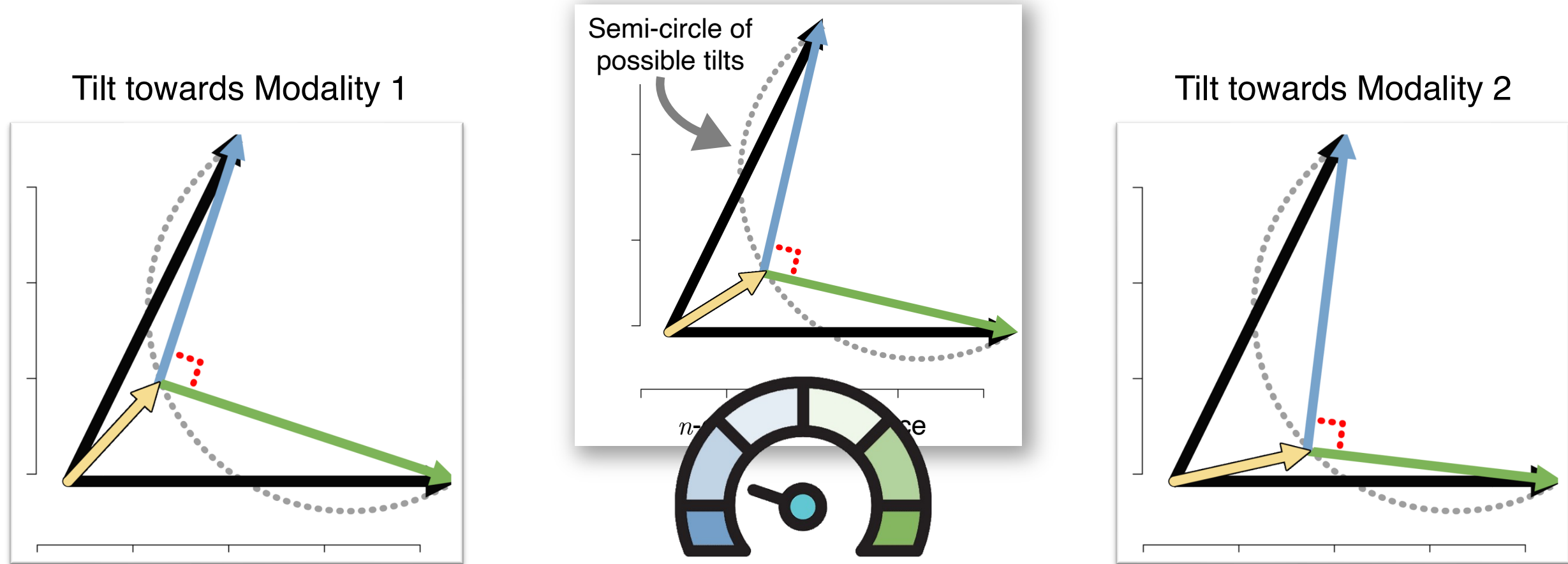
Common embedding has geometry similar to Modality 2

Our novelty (combining ideas in CCA with geometry): The tilt of the common vector (if chosen appropriately) allows the common embedding to reflect the shared geometry between both modalities.



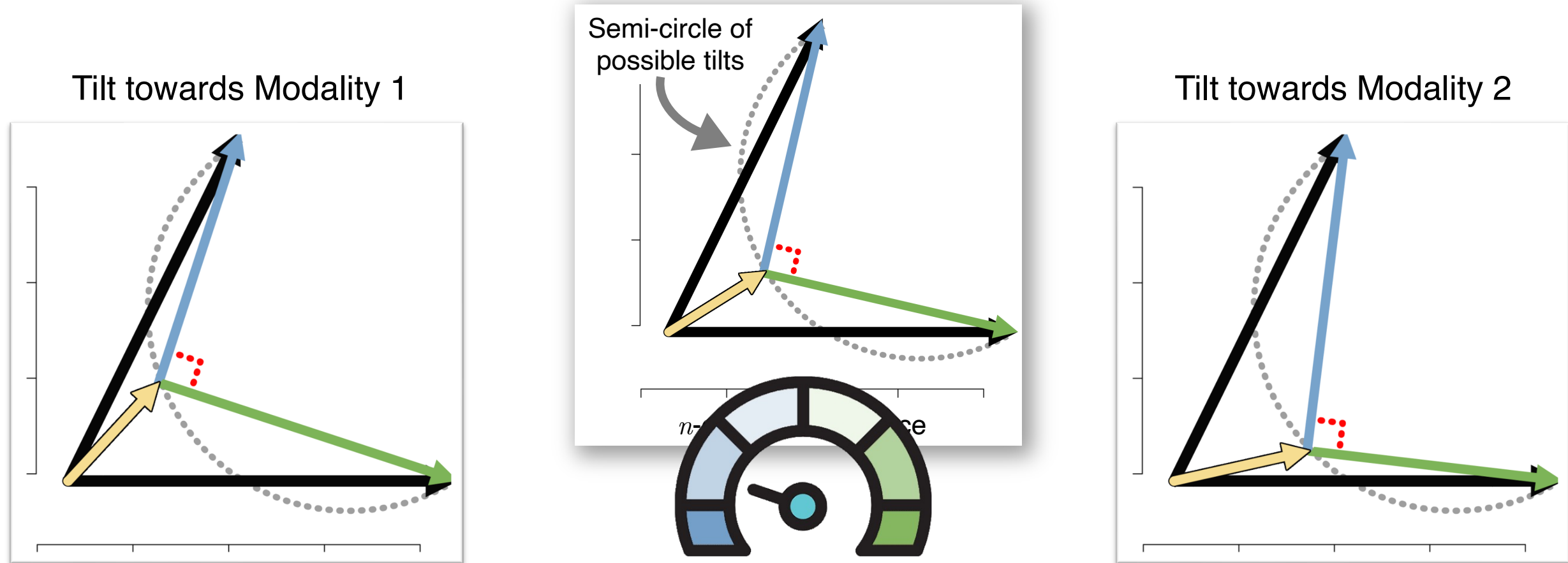
Main Goal: To estimate the tilt $\tau_j \in [1, 2]$ in each latent dimension j to capture the desired shared geometry

Our novelty (combining ideas in CCA with geometry): The tilt of the common vector (if chosen appropriately) allows the common embedding to reflect the shared geometry between both modalities.



Main Goal: To estimate the tilt $\tau_j \in [1, 2]$ in each latent dimension j to capture the **desired shared geometry**

Our novelty (combining ideas in CCA with geometry): The tilt of the common vector (if chosen appropriately) allows the common embedding to reflect the shared geometry between both modalities.





Main Goal: To **estimate** the tilt $\tau_j \in [1, 2]$ in each latent dimension j to capture the **desired shared geometry**

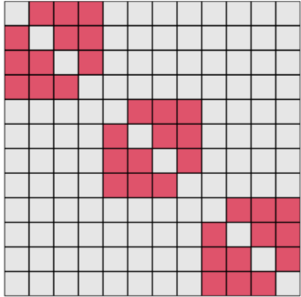
(Unsupervised) Estimation: Part 1 – Defining the shared geometry

(Unsupervised) Estimation: Part 1 – Defining the shared geometry

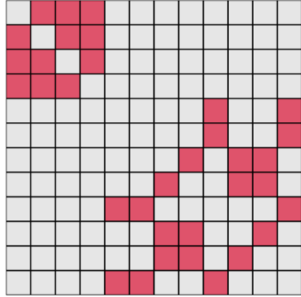
1

	Edge between two cells
	Otherwise

Cell-by-cell
nearest-neighbor
graph
(RNA modality)



“3 clusters”



“2 clusters”

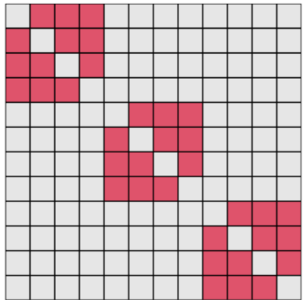
Cell-by-cell
nearest-neighbor
graph
(Protein modality)

(Unsupervised) Estimation: Part 1 – Defining the shared geometry

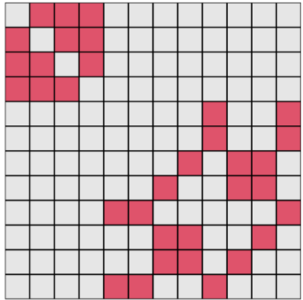
1

■ Edge between two cells
□ Otherwise

Cell-by-cell nearest-neighbor graph (RNA modality)

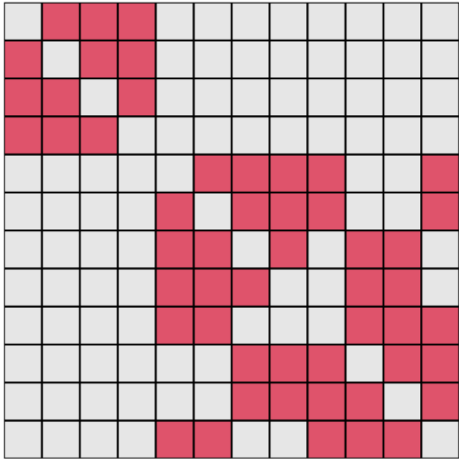


Cell-by-cell nearest-neighbor graph (Protein modality)



2

Desired shared geometry



“2 clusters”

Entry-wise operation:
Aggregating the edges in both graphs

Intuition: Cell clusters unique to one modality should not appear in the shared geometry.

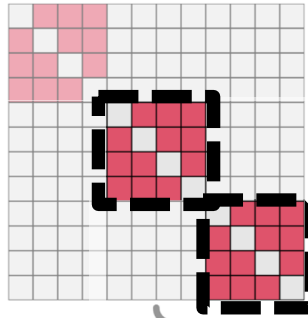
(Unsupervised) Estimation: Part 1 – Defining the shared geometry

1

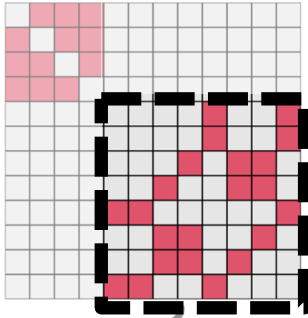
Legend:

- Red square: Edge between two cells
- Orange square: Otherwise

Cell-by-cell nearest-neighbor graph (RNA modality)

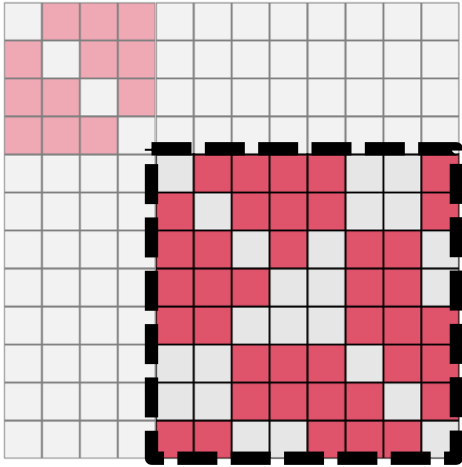


Cell-by-cell nearest-neighbor graph (Protein modality)



2

Desired shared geometry



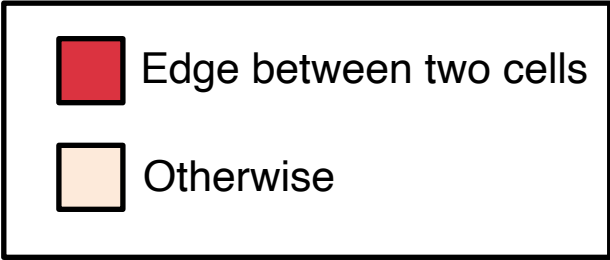
“2 clusters”

Entry-wise operation:
Aggregating the edges in both graphs

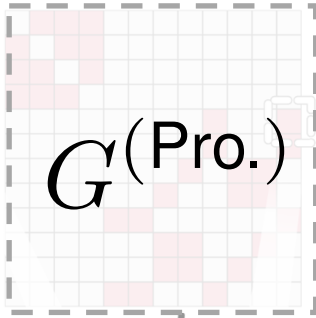
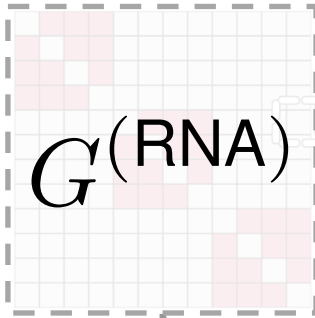
Intuition: Cell clusters unique to one modality should not appear in the shared geometry.

(Unsupervised) Estimation: Part 1 – Defining the shared geometry

1



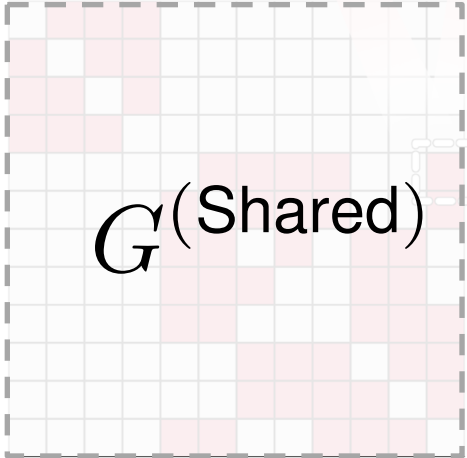
Cell-by-cell nearest-neighbor graph (RNA modality)



Cell-by-cell nearest-neighbor graph (Protein modality)

2

Desired shared geometry



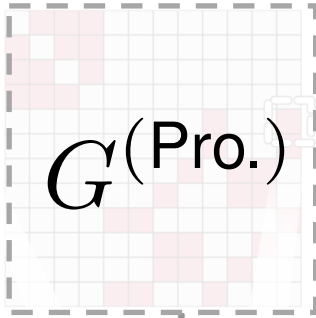
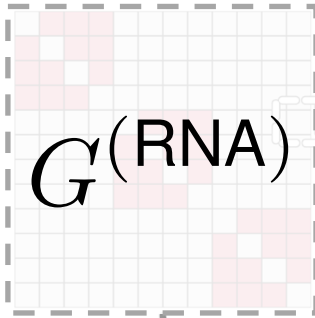
Entry-wise operation:
Aggregating the edges in both graphs

Intuition: Cell clusters unique to one modality should not be appear in the shared geometry.

(Unsupervised) Estimation: Part 2 – Estimating the corresponding matrix factorization

1

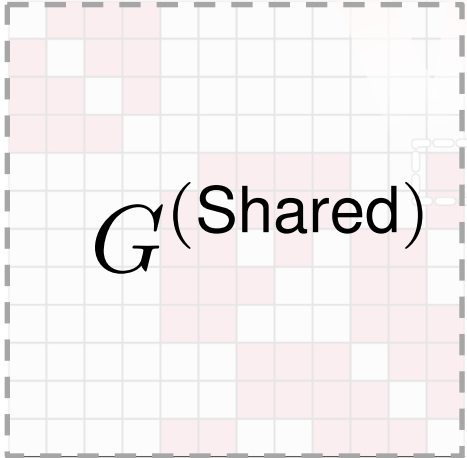
Cell-by-cell
nearest-neighbor
graph
(RNA modality)



Cell-by-cell
nearest-neighbor
graph
(Protein modality)

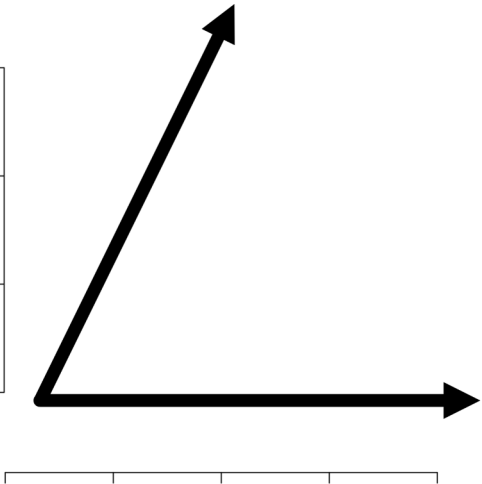
2

Desired shared geometry



3

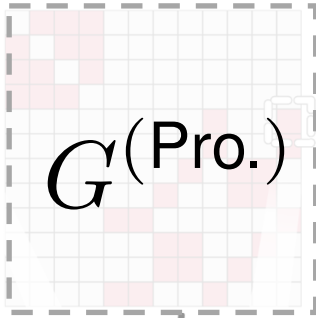
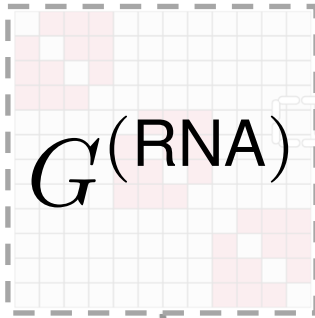
Perform
CCA



(Unsupervised) Estimation: Part 2 – Estimating the corresponding matrix factorization

1

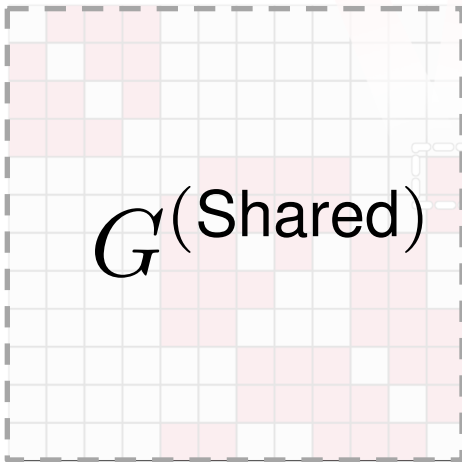
Cell-by-cell nearest-neighbor graph (RNA modality)



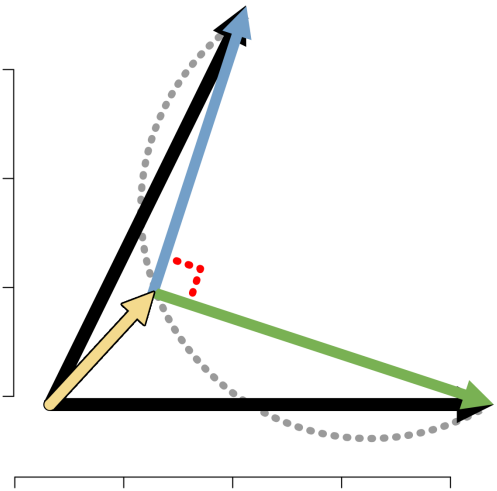
Cell-by-cell nearest-neighbor graph (Protein modality)

2

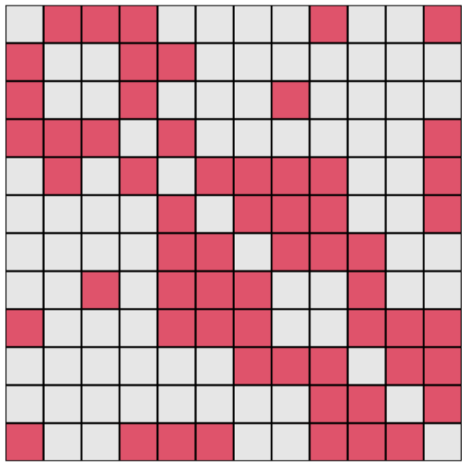
Desired shared geometry



For a latent dim, posit a tilt of the common vector...



... and its corresponding cell-by-cell nearest-neighbor graph for the common embedding



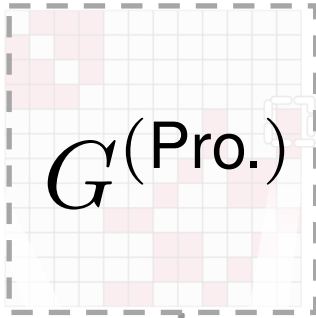
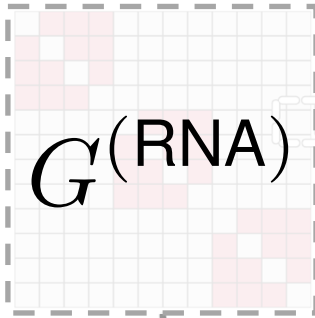
3

Perform CCA

(Unsupervised) Estimation: Part 2 – Estimating the corresponding matrix factorization

1

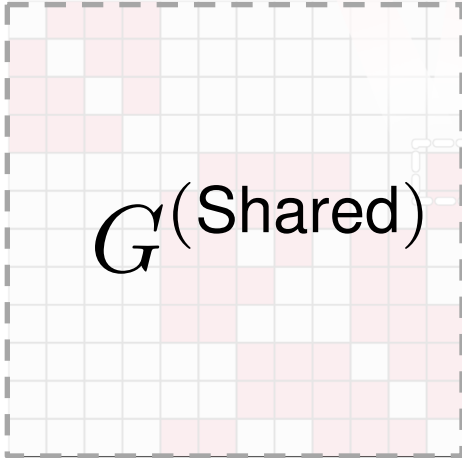
Cell-by-cell nearest-neighbor graph (RNA modality)



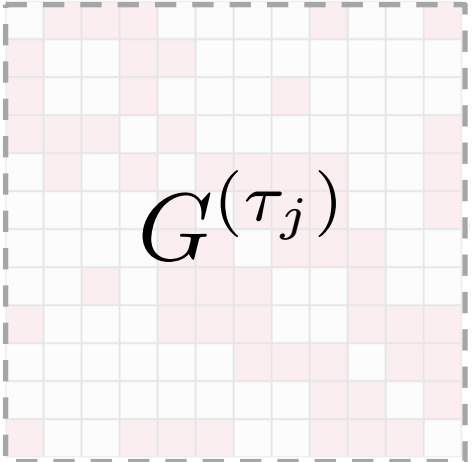
Cell-by-cell nearest-neighbor graph (Protein modality)

2

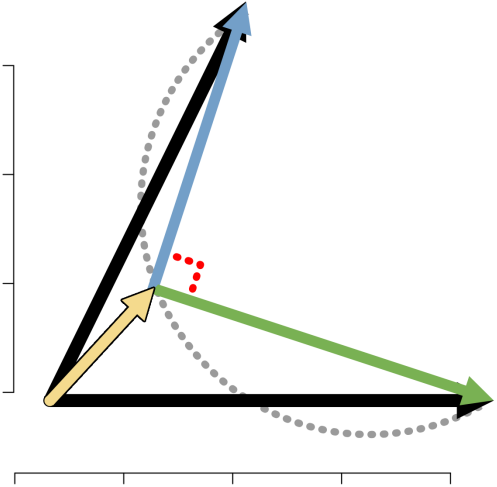
Desired shared geometry



... and its corresponding cell-by-cell nearest-neighbor graph for the common embedding



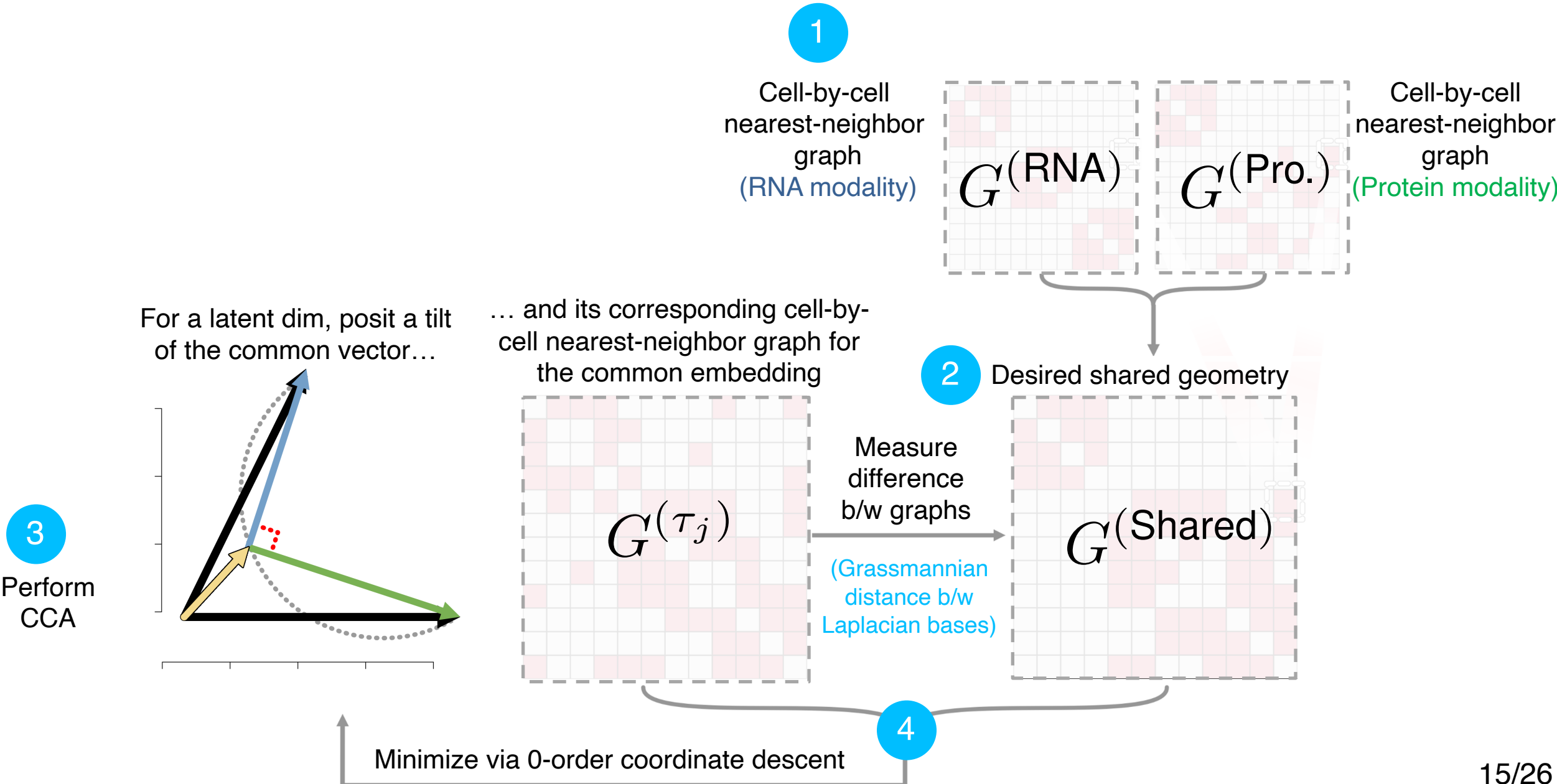
For a latent dim, posit a tilt of the common vector...



3

Perform CCA

(Unsupervised) Estimation: Part 2 – Estimating the corresponding matrix factorization



(Unsupervised) Estimation: Part 2 – Estimating the corresponding matrix factorization

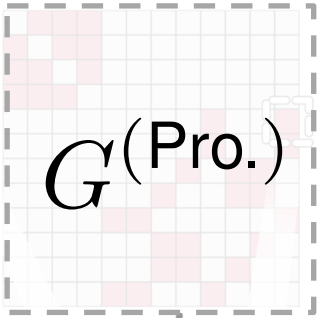
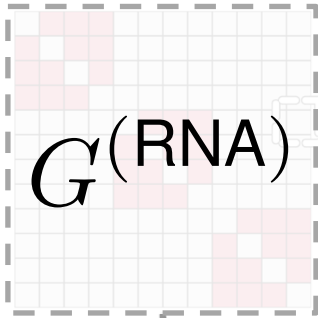
$$X^{(1)} \approx \left[C + D^{(1)} \right] \times L^{(1)}$$

$$X^{(2)} \approx \left[C + D^{(2)} \right] \times L^{(2)}$$

$C + D^{(1)} = Z^{(1)}, C + D^{(2)} = Z^{(2)}$

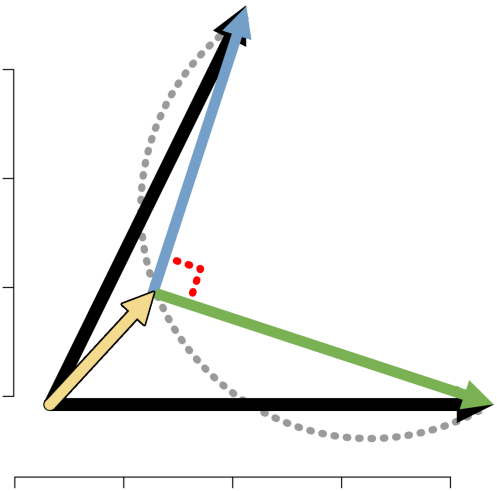
1

Cell-by-cell nearest-neighbor graph (RNA modality)

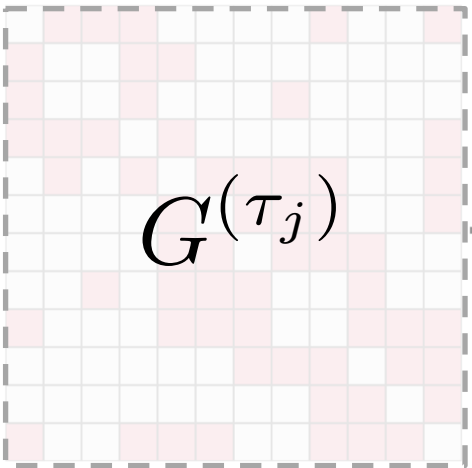


Cell-by-cell nearest-neighbor graph (Protein modality)

For a latent dim, posit a tilt of the common vector...

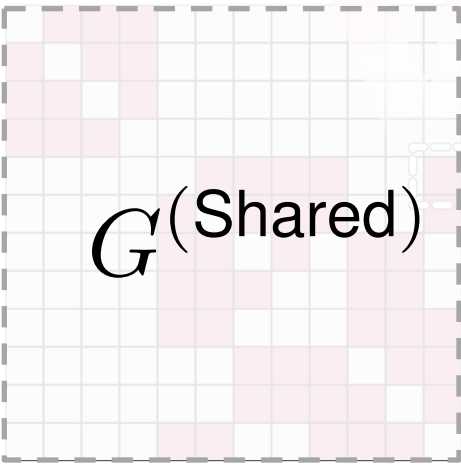


... and its corresponding cell-by-cell nearest-neighbor graph for the common embedding



2

Desired shared geometry



Measure difference b/w graphs (Grassmannian distance b/w Laplacian bases)

3

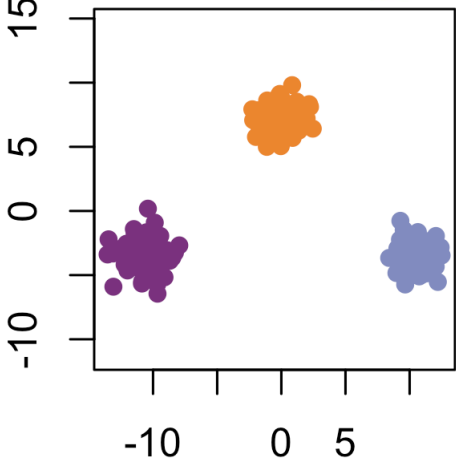
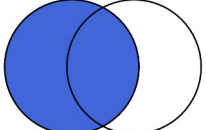
Perform CCA

4

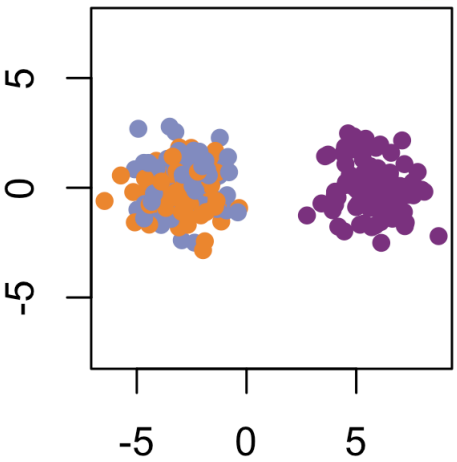
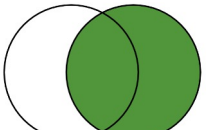
Minimize via 0-order coordinate descent

Simulation 1: Toy

Modality 1
(Leading PC's)

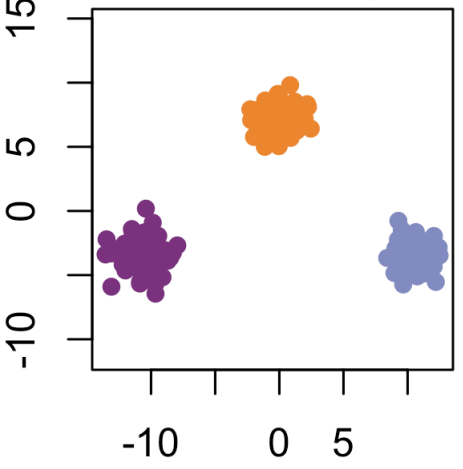


Modality 2
(Leading PC's)

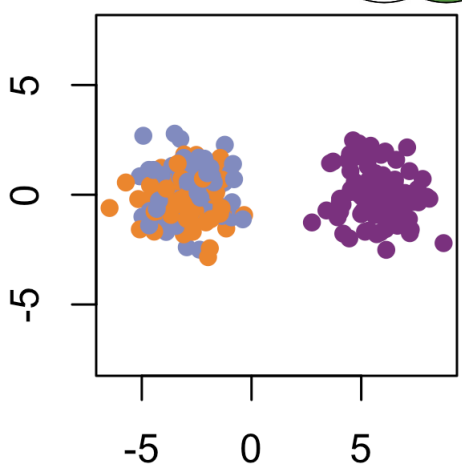


Simulation 1: Toy

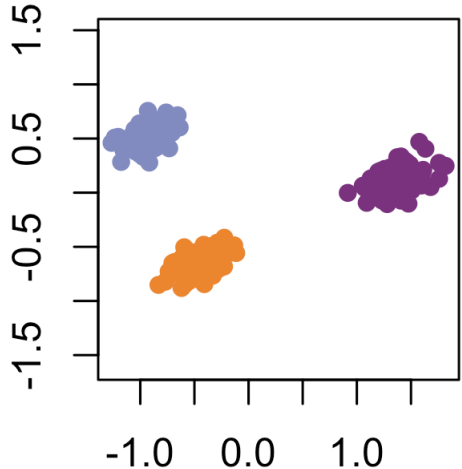
Modality 1
(Leading PC's)



Modality 2
(Leading PC's)

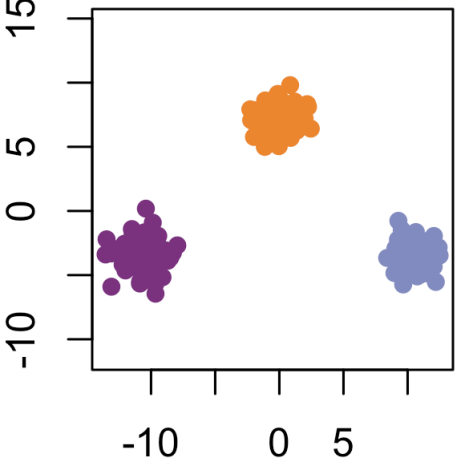


Consensus PCA
(2D: Union)

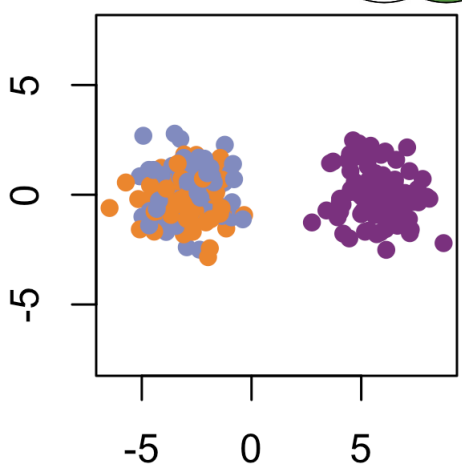


Simulation 1: Toy

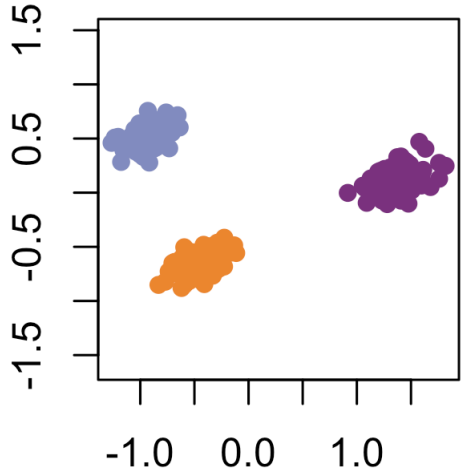
Modality 1
(Leading PC's)



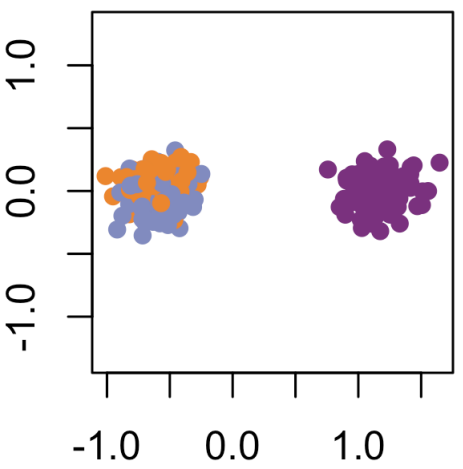
Modality 2
(Leading PC's)



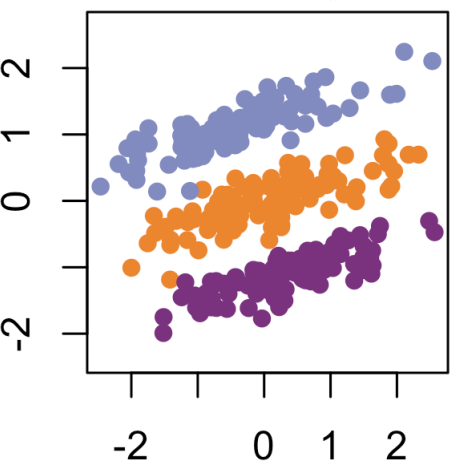
Consensus PCA
(2D: Union)



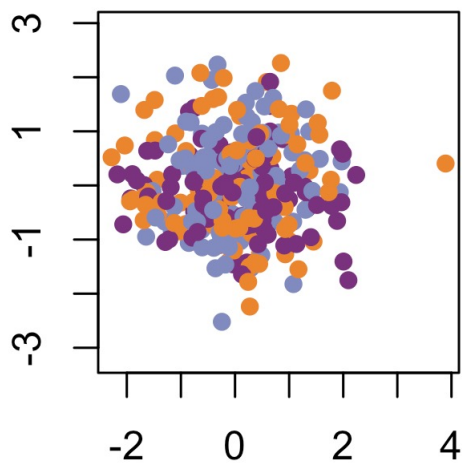
Tilted-CCA: Common
(2D: Intersection)



Distinct 1

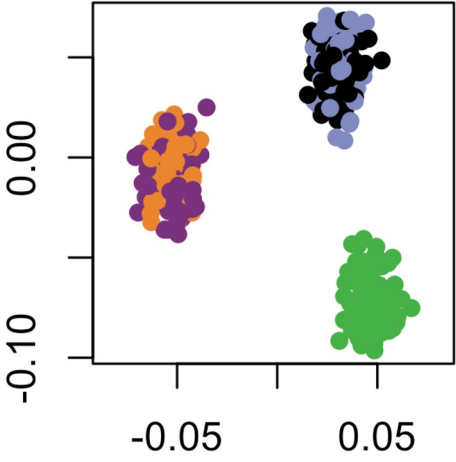
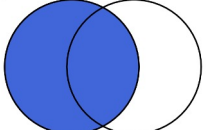


Distinct 2

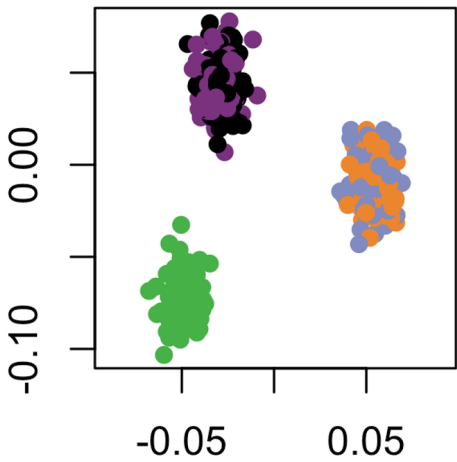
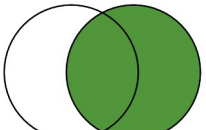


Simulation 2: Criss-cross

Modality 1
(Leading PC's)

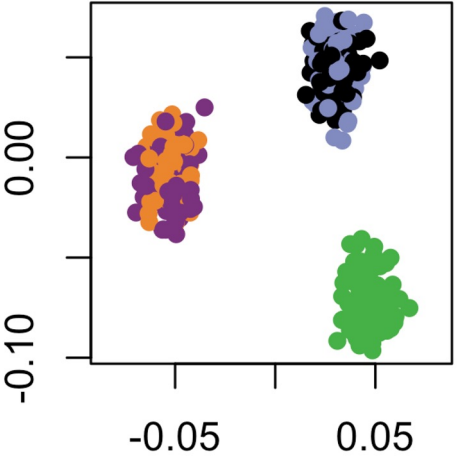
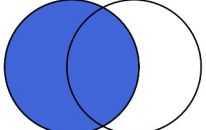


Modality 2
(Leading PC's)

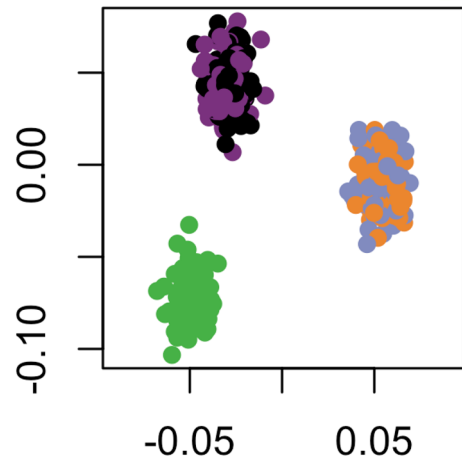
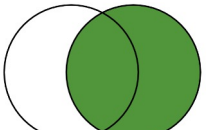


Simulation 2: Criss-cross

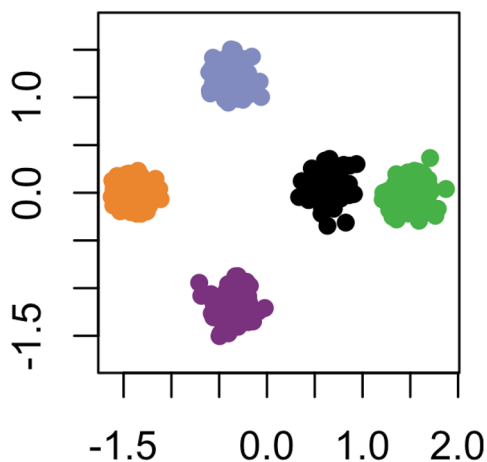
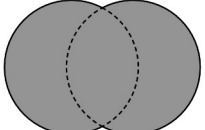
Modality 1
(Leading PC's)



Modality 2
(Leading PC's)

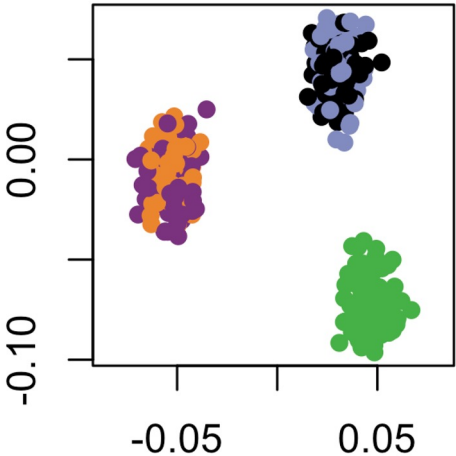


Consensus PCA
(2D: Union)

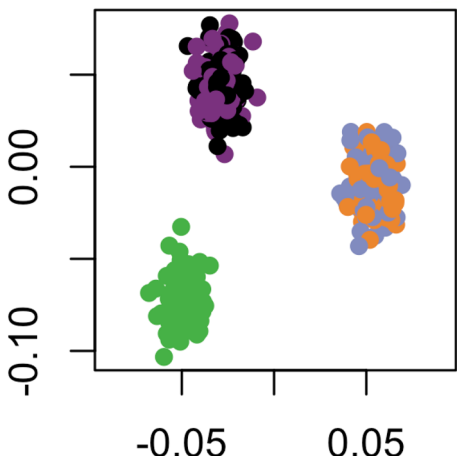


Simulation 2: Criss-cross

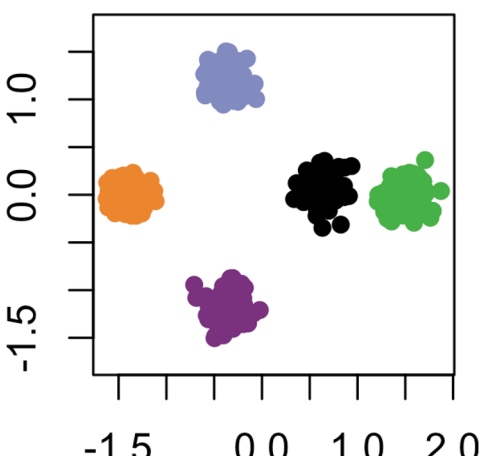
Modality 1
(Leading PC's)



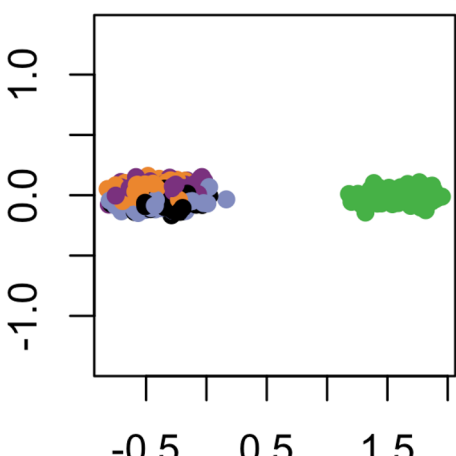
Modality 2
(Leading PC's)



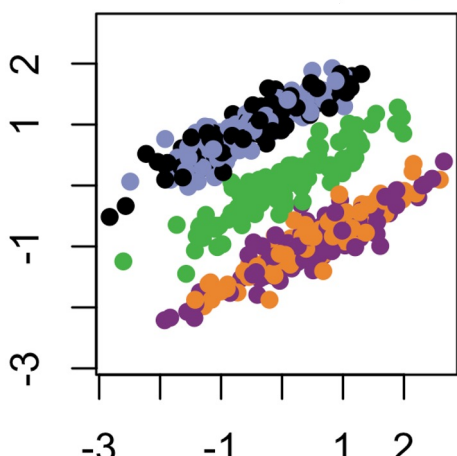
Consensus PCA
(2D: Union)



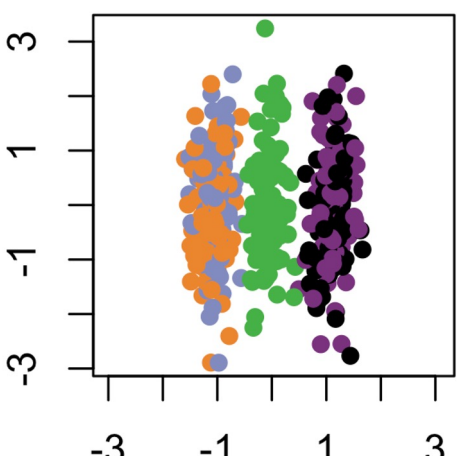
Tilted-CCA: Common
(2D: Intersection)



Distinct 1



Distinct 2

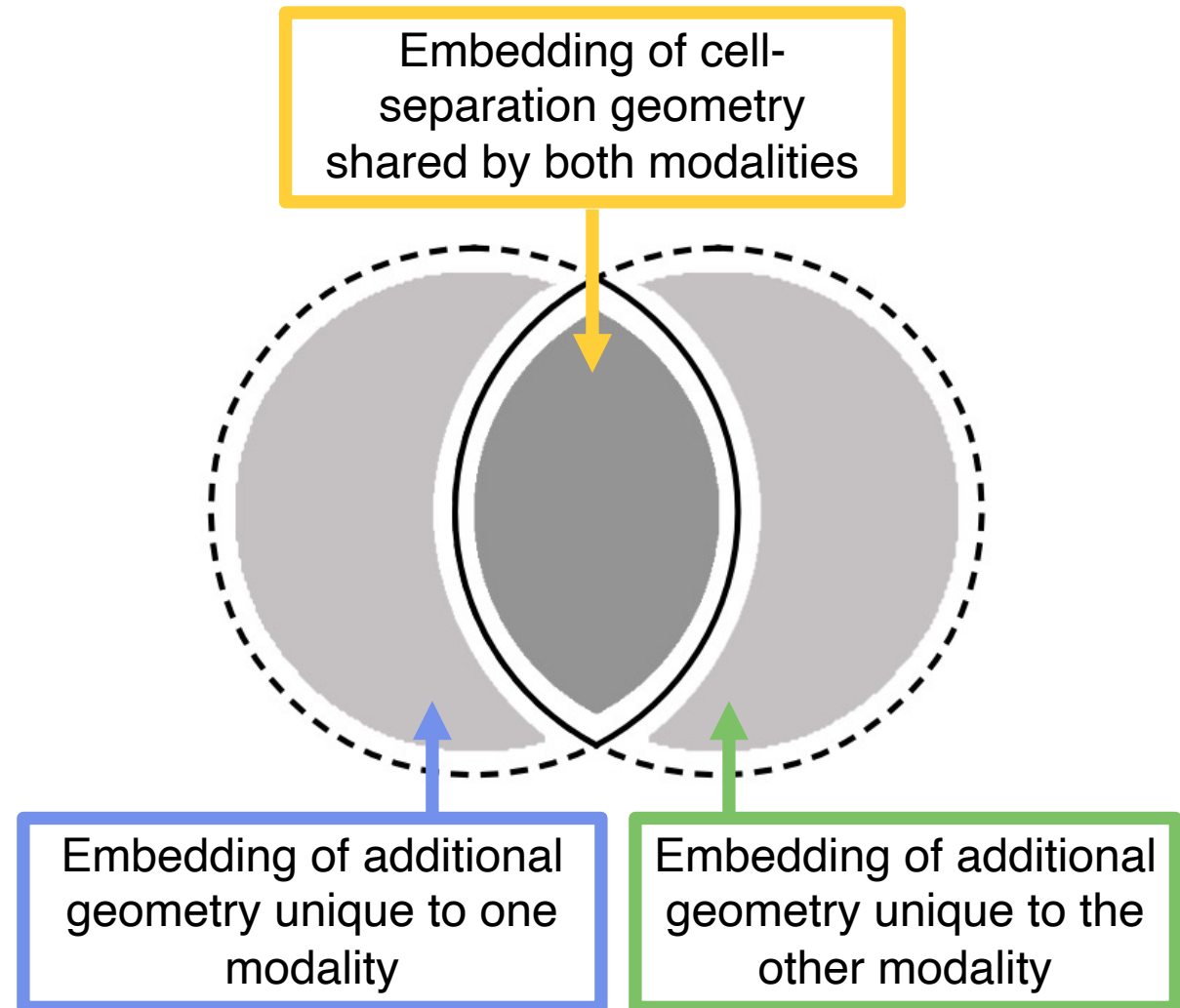


Single-cell investigation:

New perspective of developmental biology, now that we have estimated the shared/unique geometry

Recap of the biological goals:

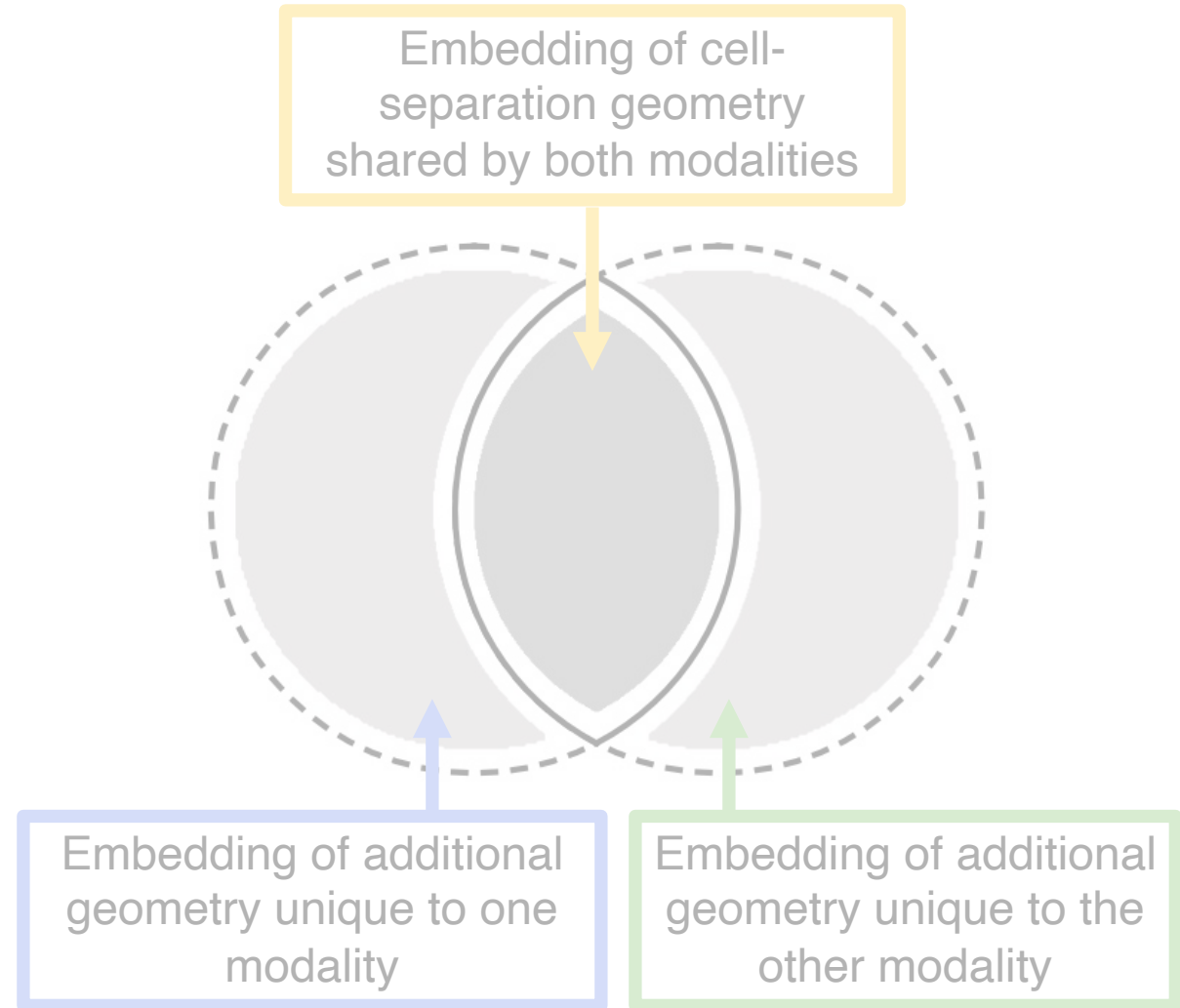
1. **(Experimental design):** Which pair of modalities should biologist sequence to have the most comprehensive understanding?
2. **(Variable selection):** For RNA-Protein data, how can we pick the antibodies that contribute the most additional information to the RNA modality?
3. **(Developmental biology):** Can the amount of coordination between two modalities tell us if a cell in a steady-state or is undergoing development?



Recap of the biological goals:

1. **(Experimental design):** Which pair of modalities should biologist sequence to have the most comprehensive understanding?
2. **(Variable selection):** For RNA-Protein data, how can we pick the antibodies that contribute the most additional information to the RNA modality?
3. **(Developmental biology):** Can the amount of coordination between two modalities tell us if a cell in a steady-state or is undergoing development?

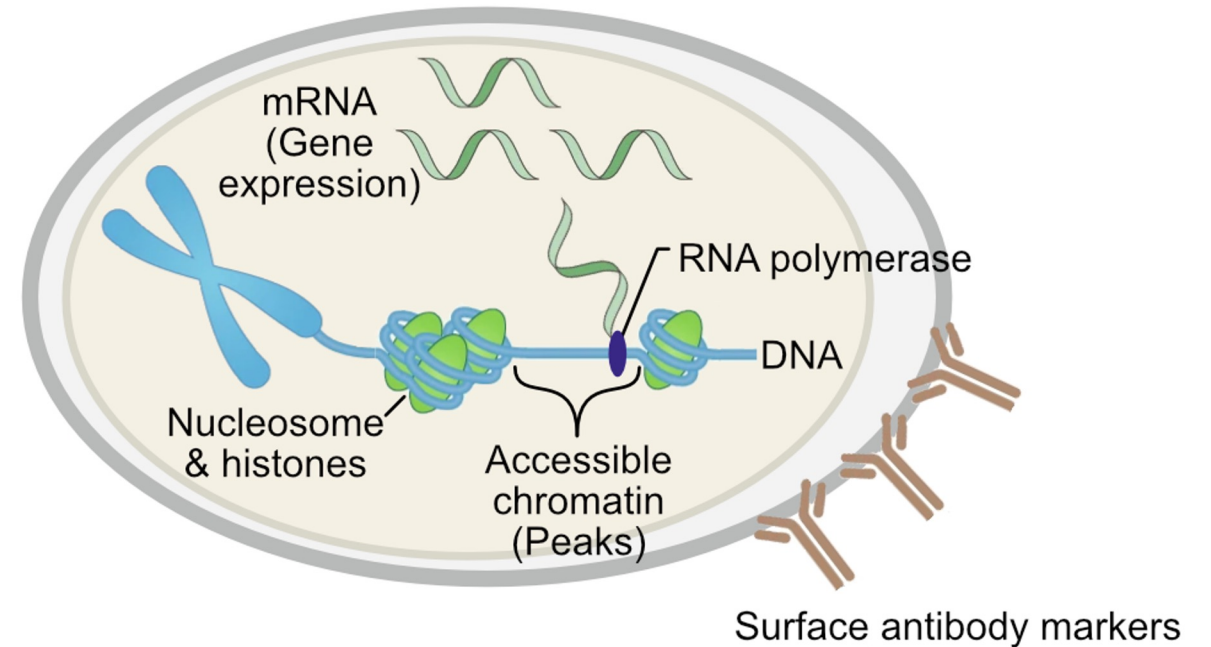
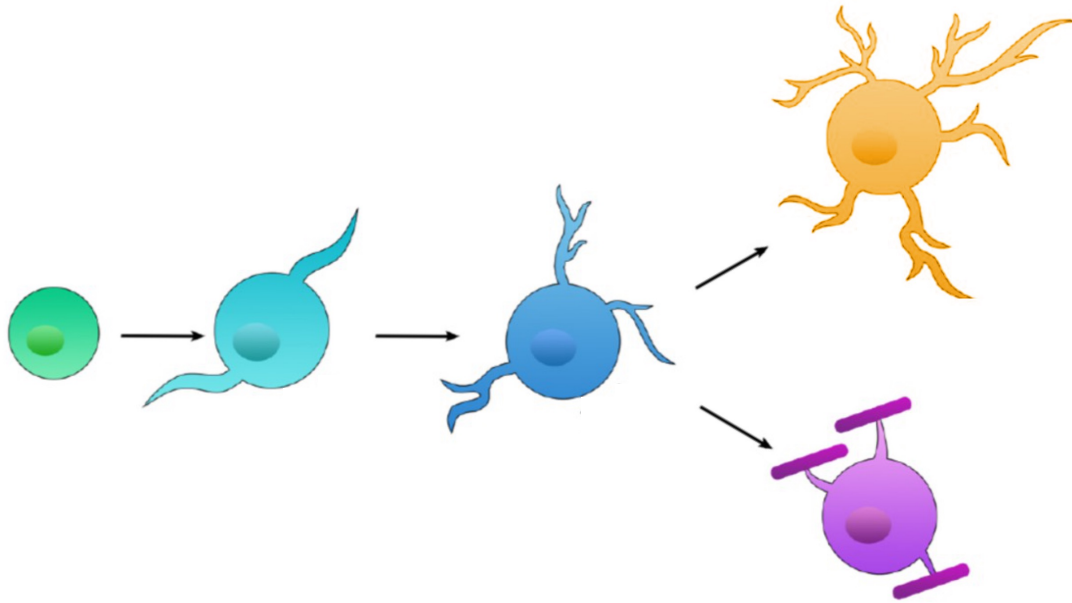
Focus of the remainder of the talk



In developmental biology, we're interested in studying how cells continually specialize over time despite given static snapshots.

Progenitor cells
(Youngest)

Mature neurons
(Oldest)



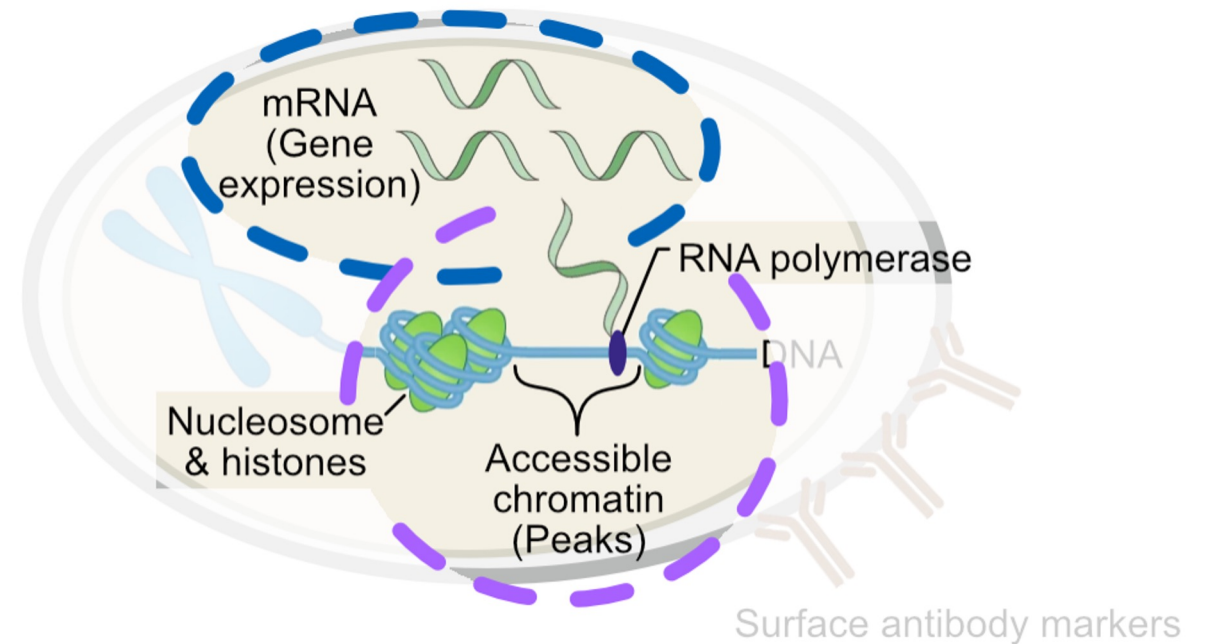
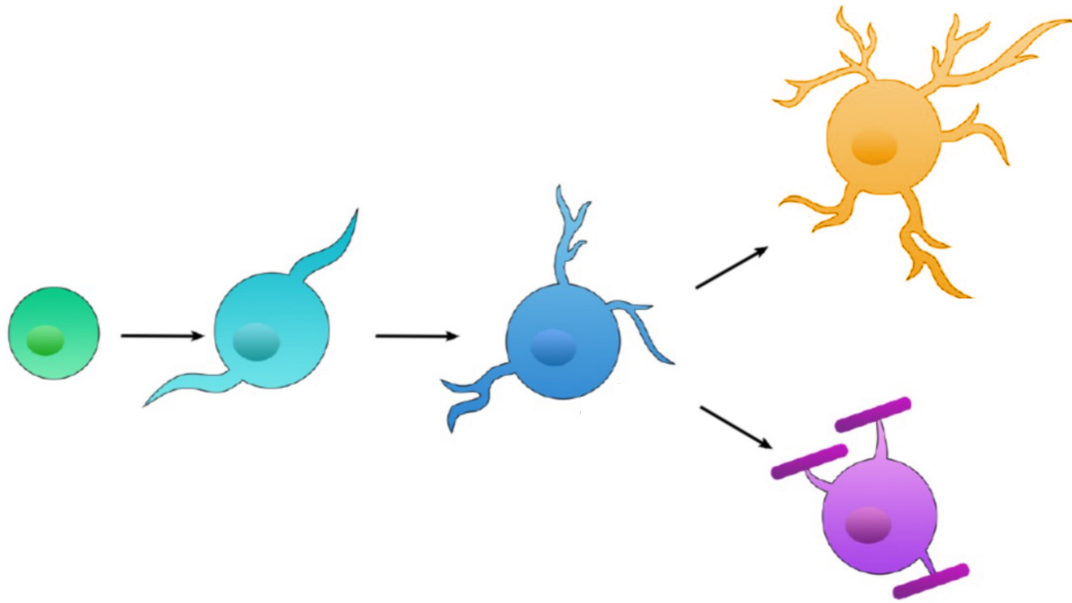
Existing methods:

Monocle (Trapnell et al., Nature Biotech., 2014),
Slingshot (Dudoit et al., BMC Genomics, 2018),
RNA velocity (Kharchenko et al., Nature, 2018)

In developmental biology, we're interested in studying how cells continually specialize over time despite given static snapshots.

Progenitor cells
(Youngest)

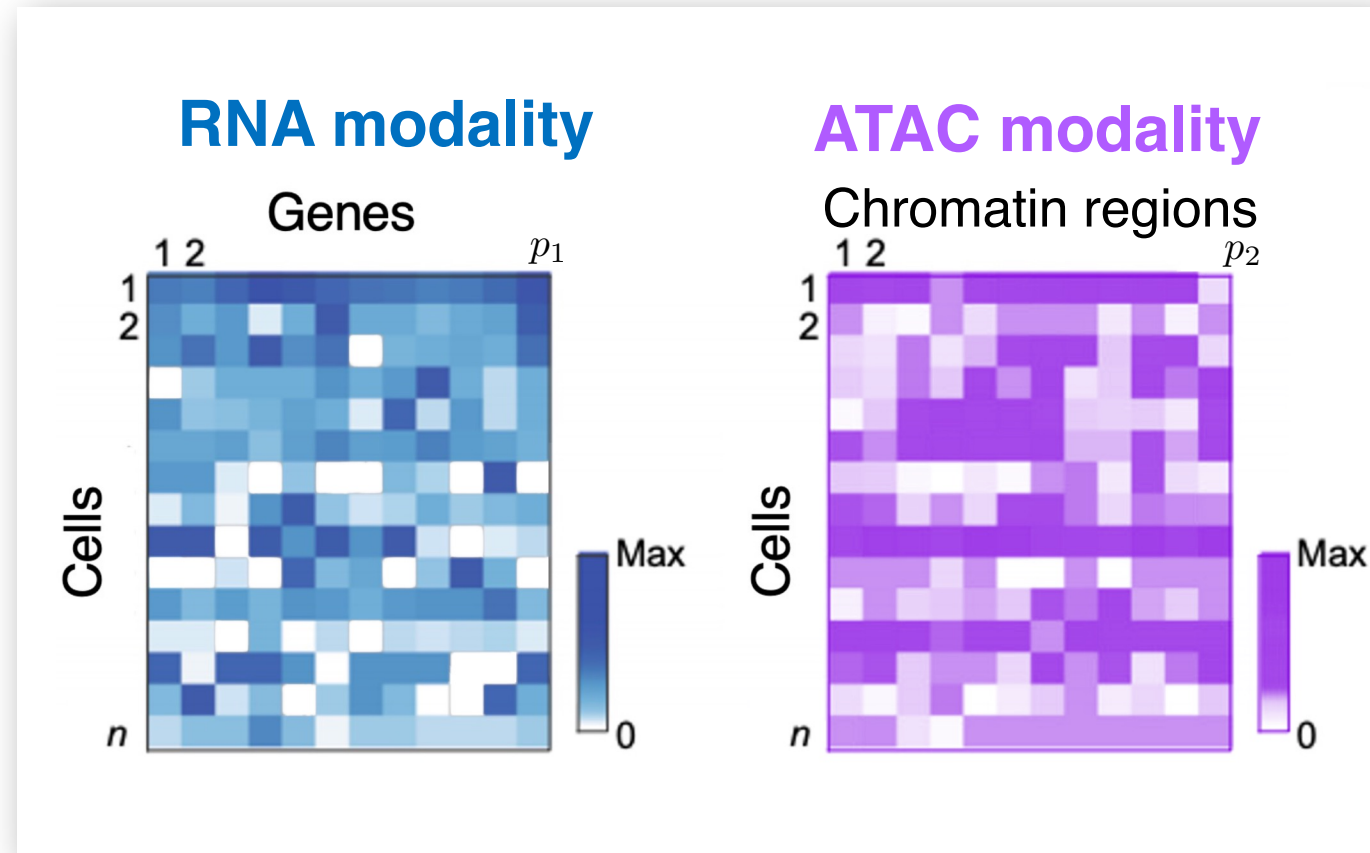
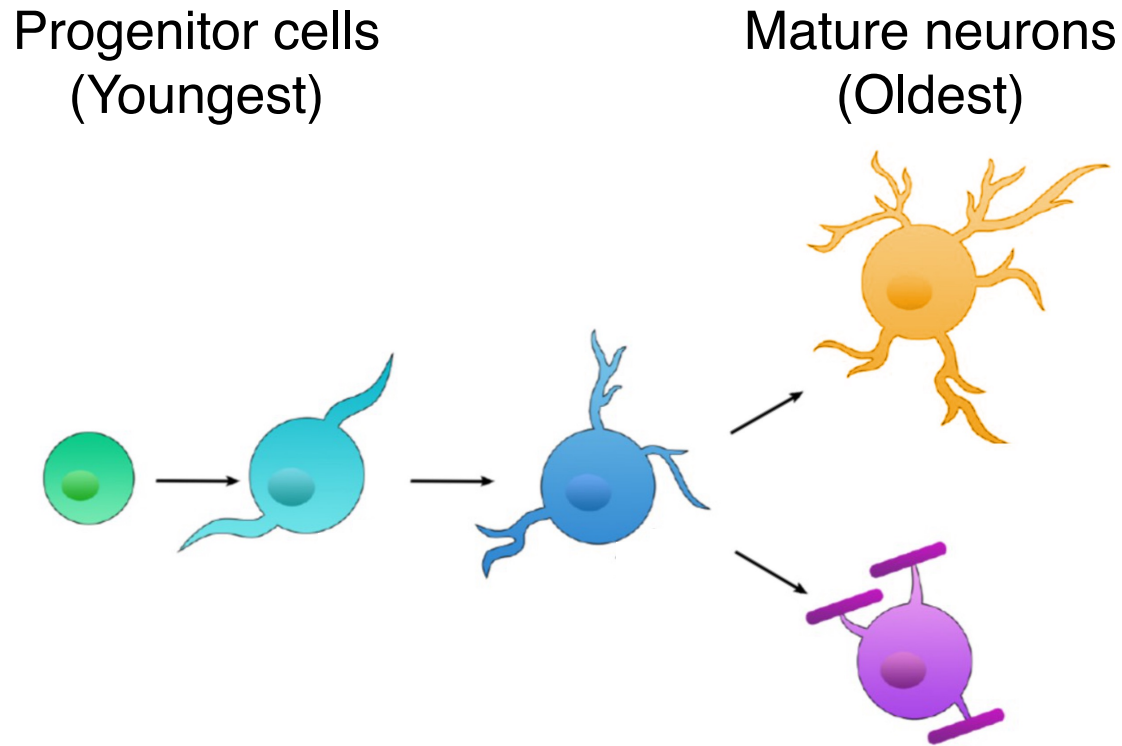
Mature neurons
(Oldest)



Existing methods:

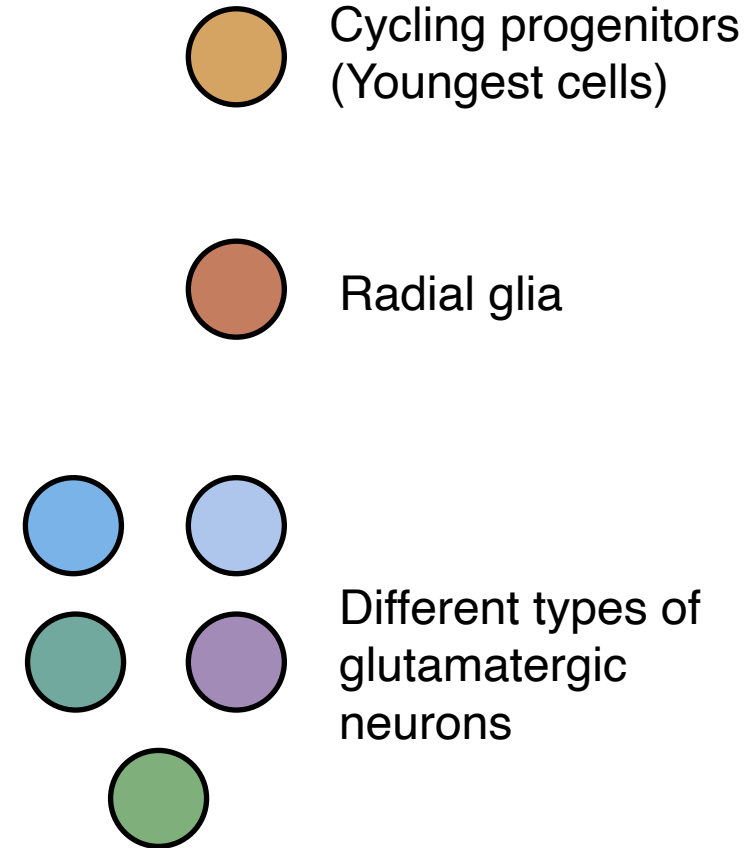
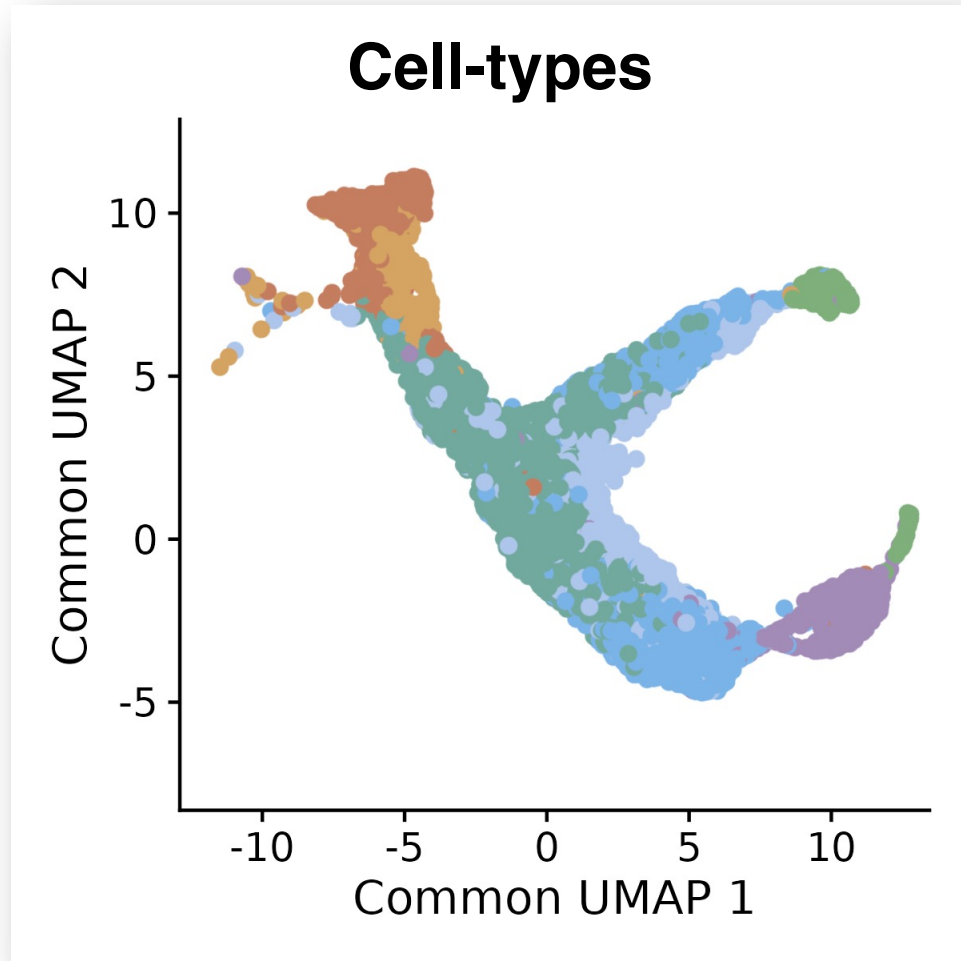
Monocle (Trapnell et al., Nature Biotech., 2014),
Slingshot (Dudoit et al., BMC Genomics, 2018),
RNA velocity (Kharchenko et al., Nature, 2018)

In developmental biology, we're interested in studying how cells continually specialize over time despite given static snapshots.



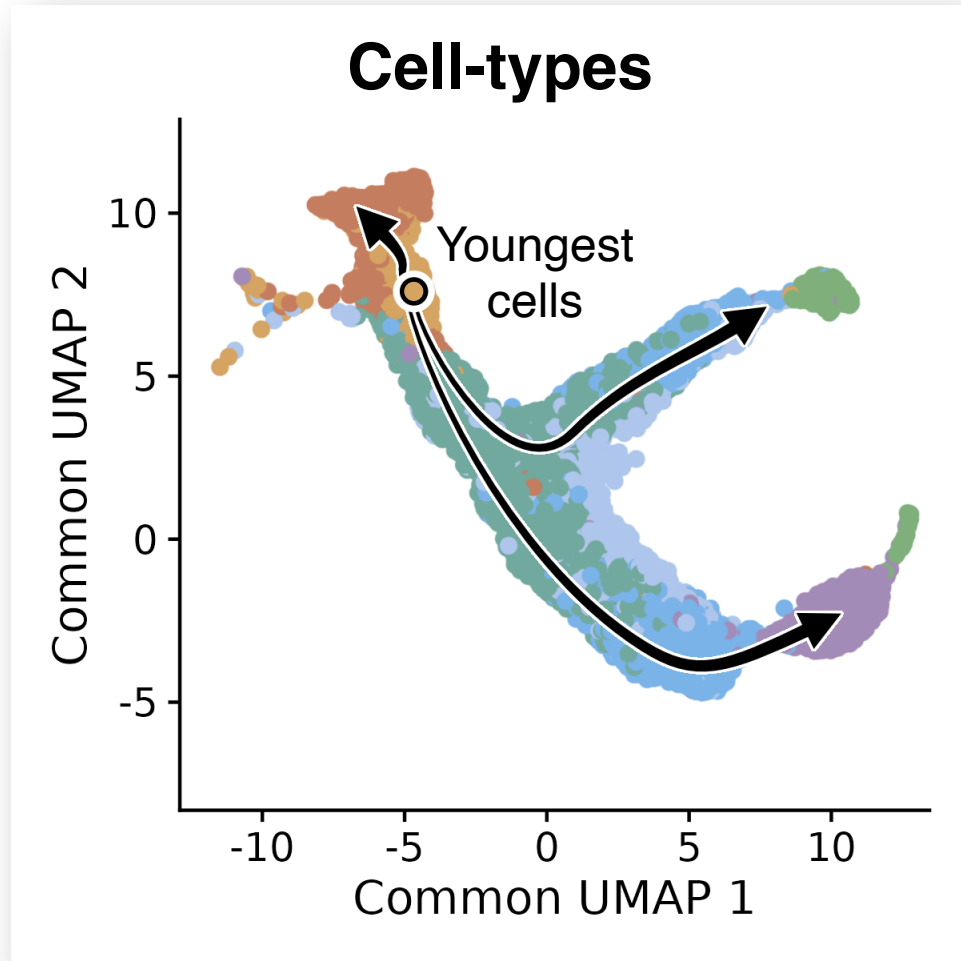
While existing methods focusing on one modality tell us **how** cells develop (“layout of the highway”), we want to know **what** status the cell is in (“speed of traffic”).

Knowing the developmental trajectories is not the same as knowing if a cell is undergoing development or is in steady-state.

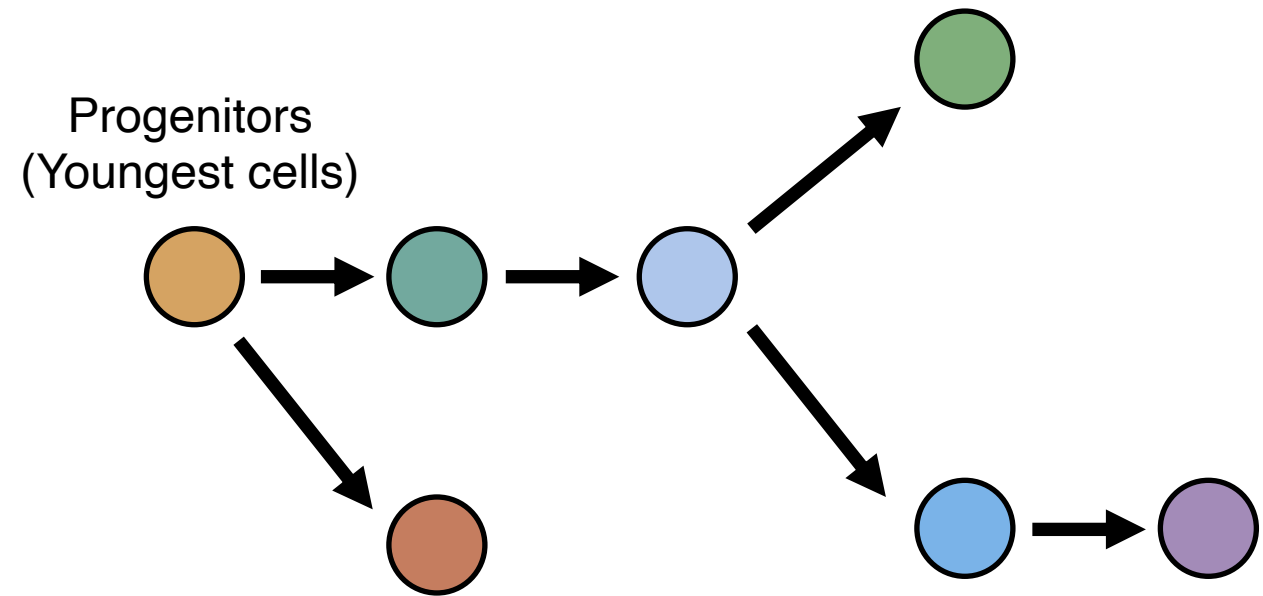


Human brain development (10x. RNA & ATAC. Greenleaf et al., Cell, 2021): 6000+ cells

Knowing the developmental trajectories is not the same as knowing if a cell is undergoing development or is in steady-state.

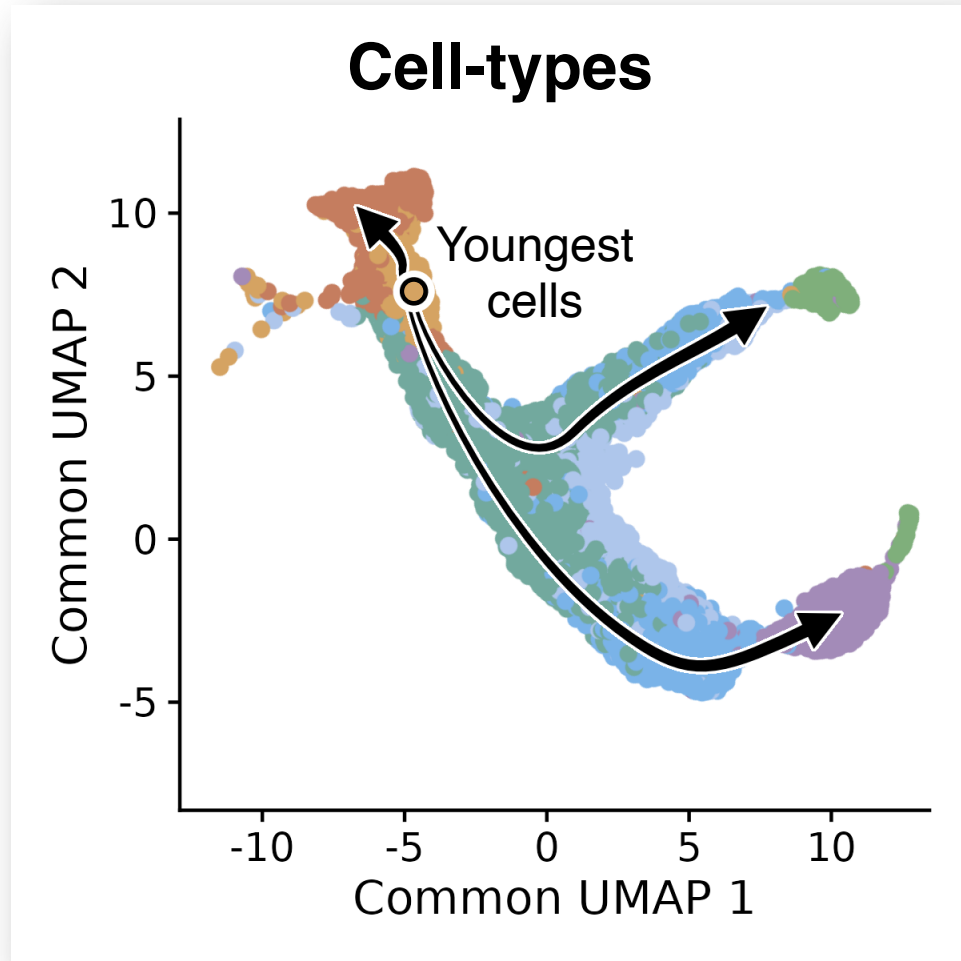


Human brain development (10x. RNA & ATAC. Greenleaf et al., Cell, 2021): 6000+ cells

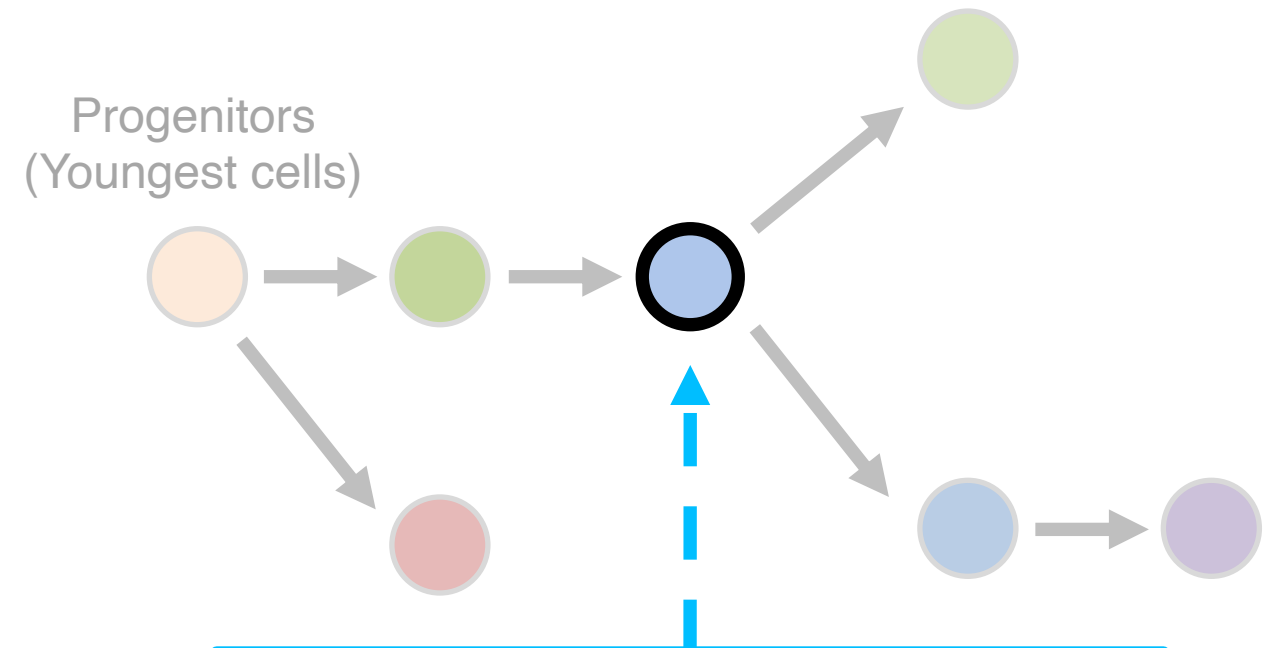


Existing methods:
“Layout of the highway”

Knowing the developmental trajectories is not the same as knowing if a cell undergoing development or is in steady-state.

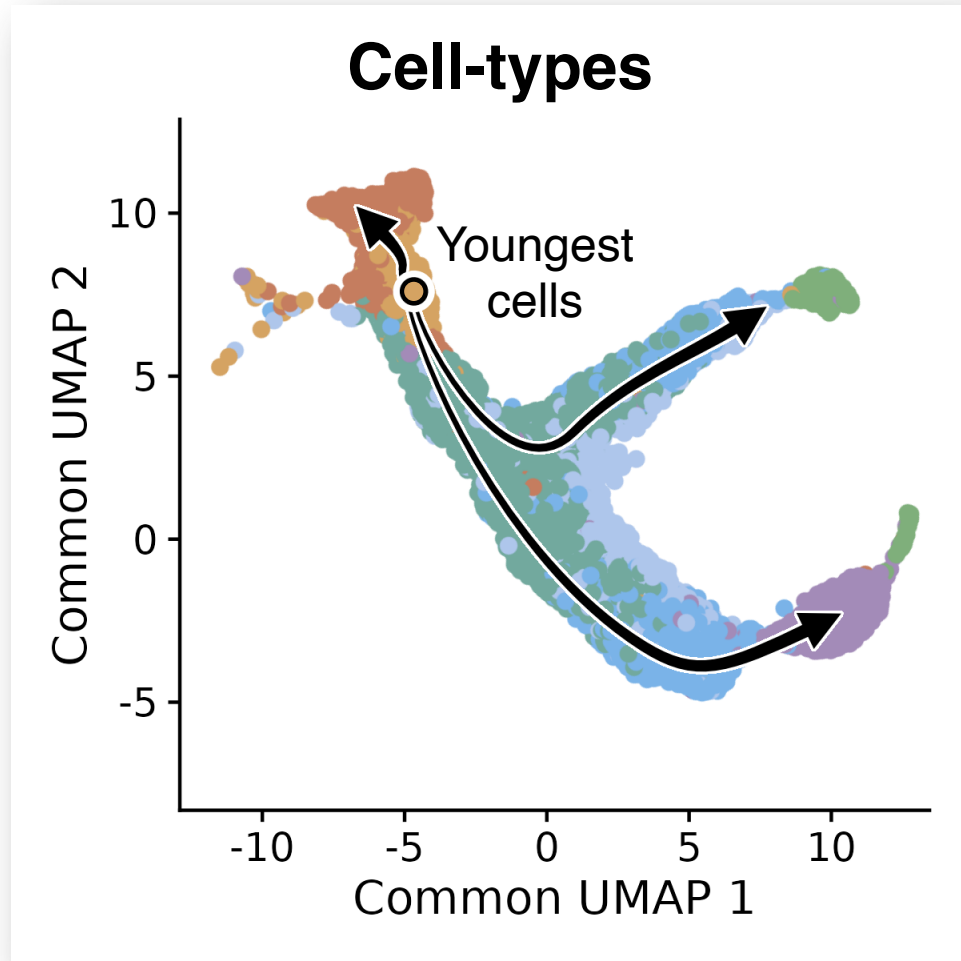


Human brain development (10x. RNA & ATAC. Greenleaf et al., Cell, 2021): 6000+ cells

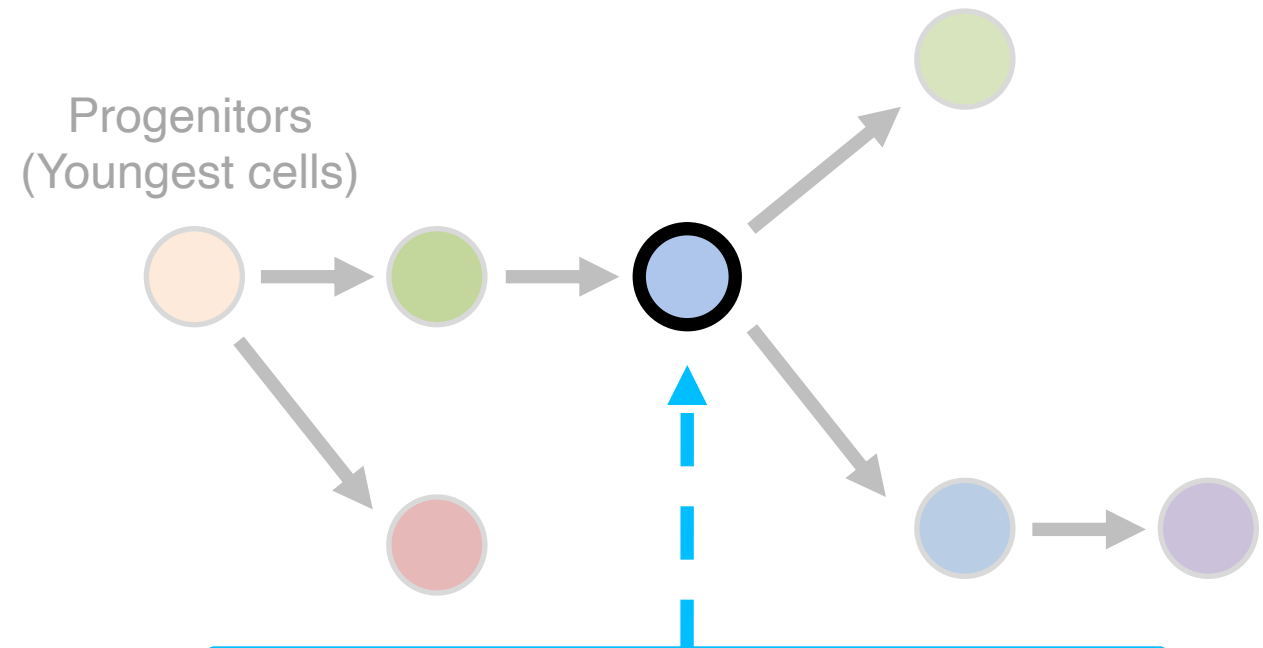


“Speed of traffic”:
Are the neurons rapidly passing through **OR** idling around/waiting for an external stimulus?

Knowing the developmental trajectories is not the same as knowing if a cell undergoing development or is in steady-state.



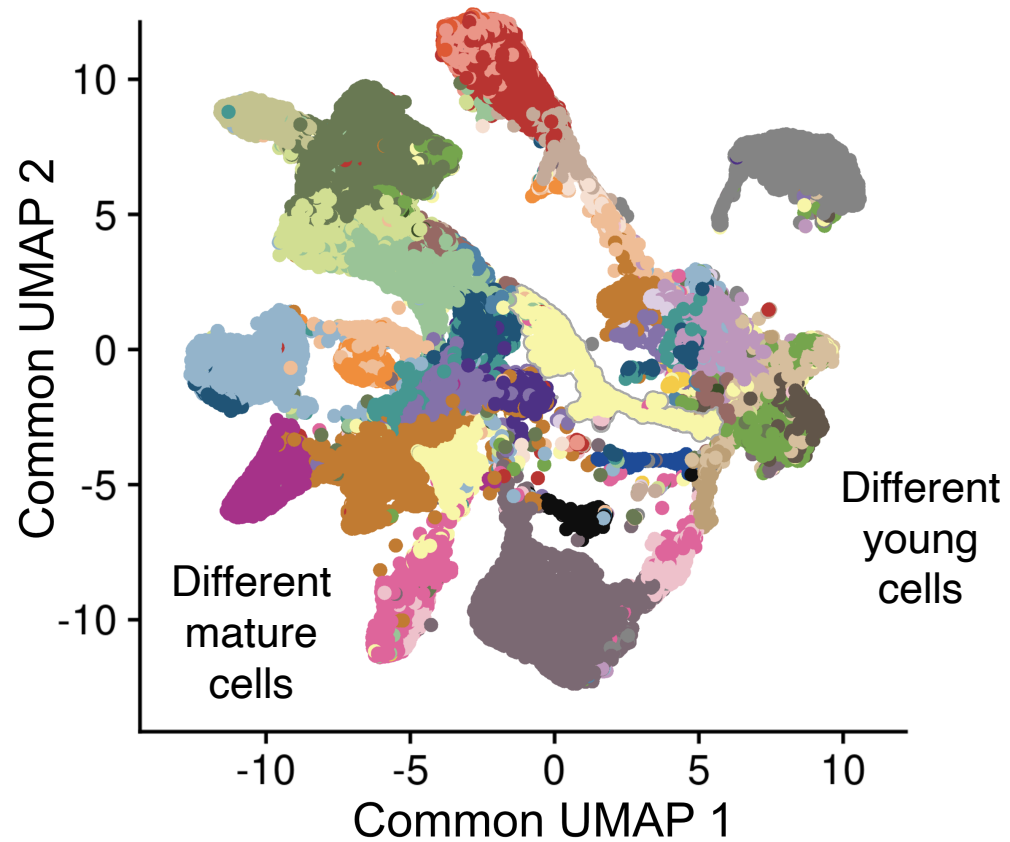
Human brain development (10x. RNA & ATAC. Greenleaf et al., Cell, 2021): 6000+ cells



“Speed of traffic”:
Are the neurons in-development
OR in steady-state?

Brief aside: More complex systems

Mouse embryonic development
(10x. Reik et al., bioRxiv, 2022)

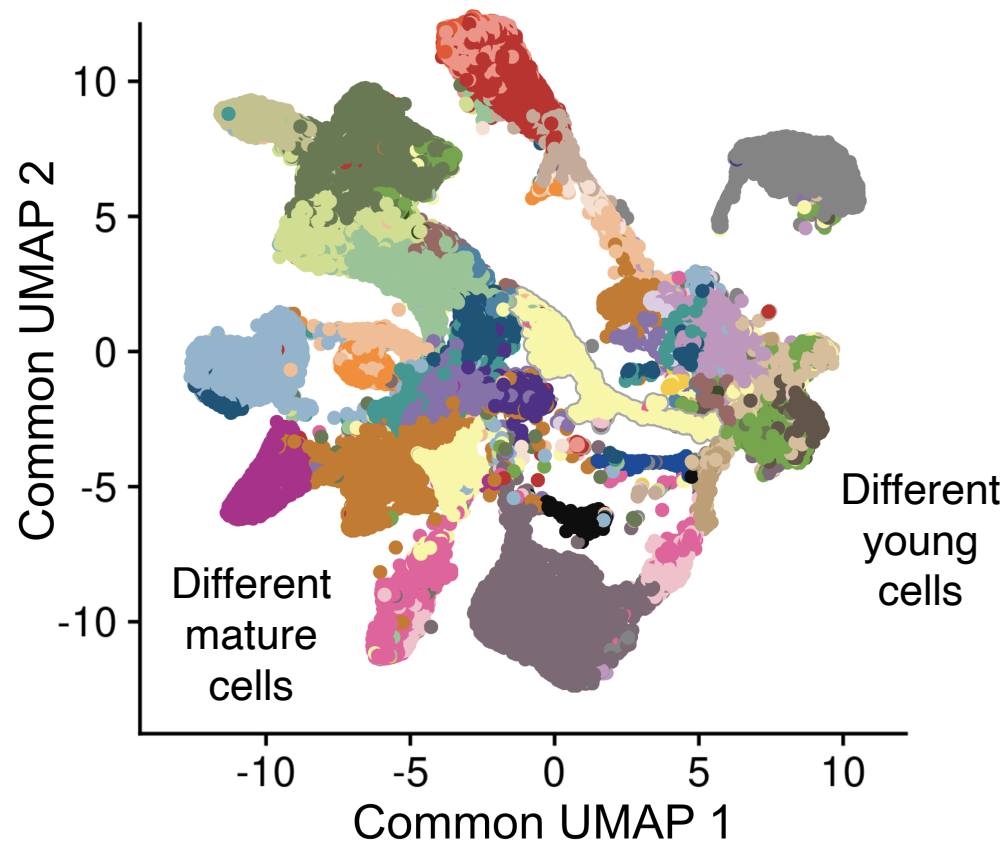


Difficult due to multiple trajectories

Brief aside: More complex systems

Mouse embryonic development

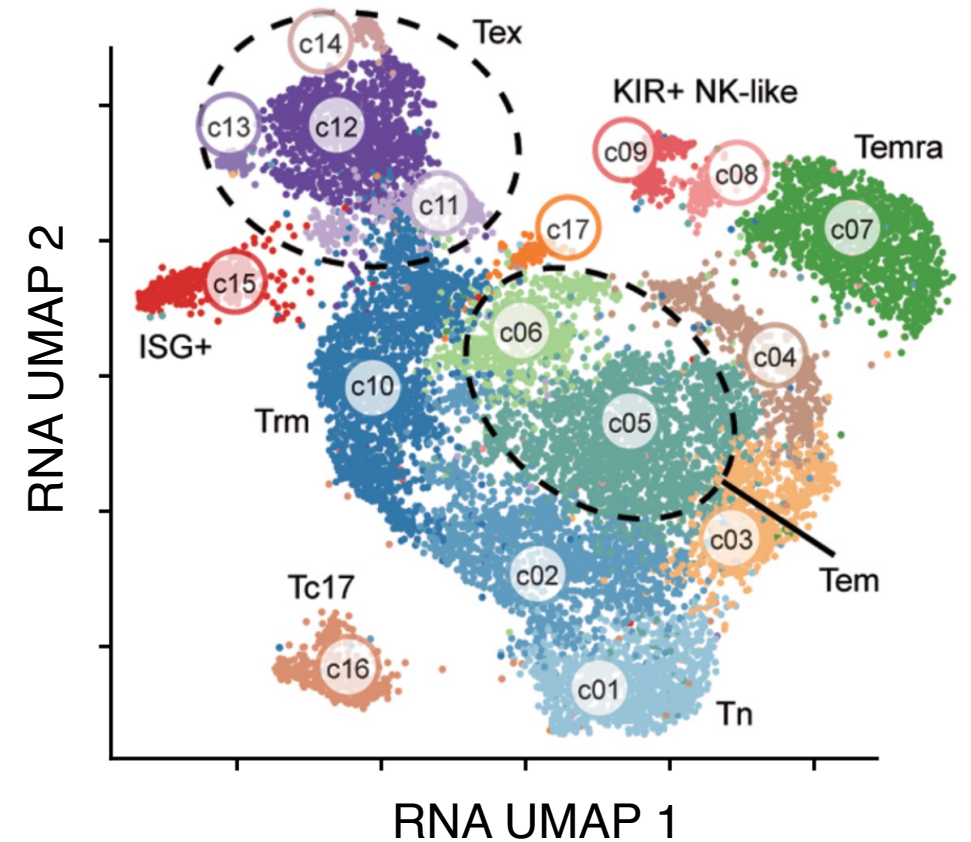
(10x. Reik et al., bioRxiv, 2022)



Difficult due to multiple trajectories

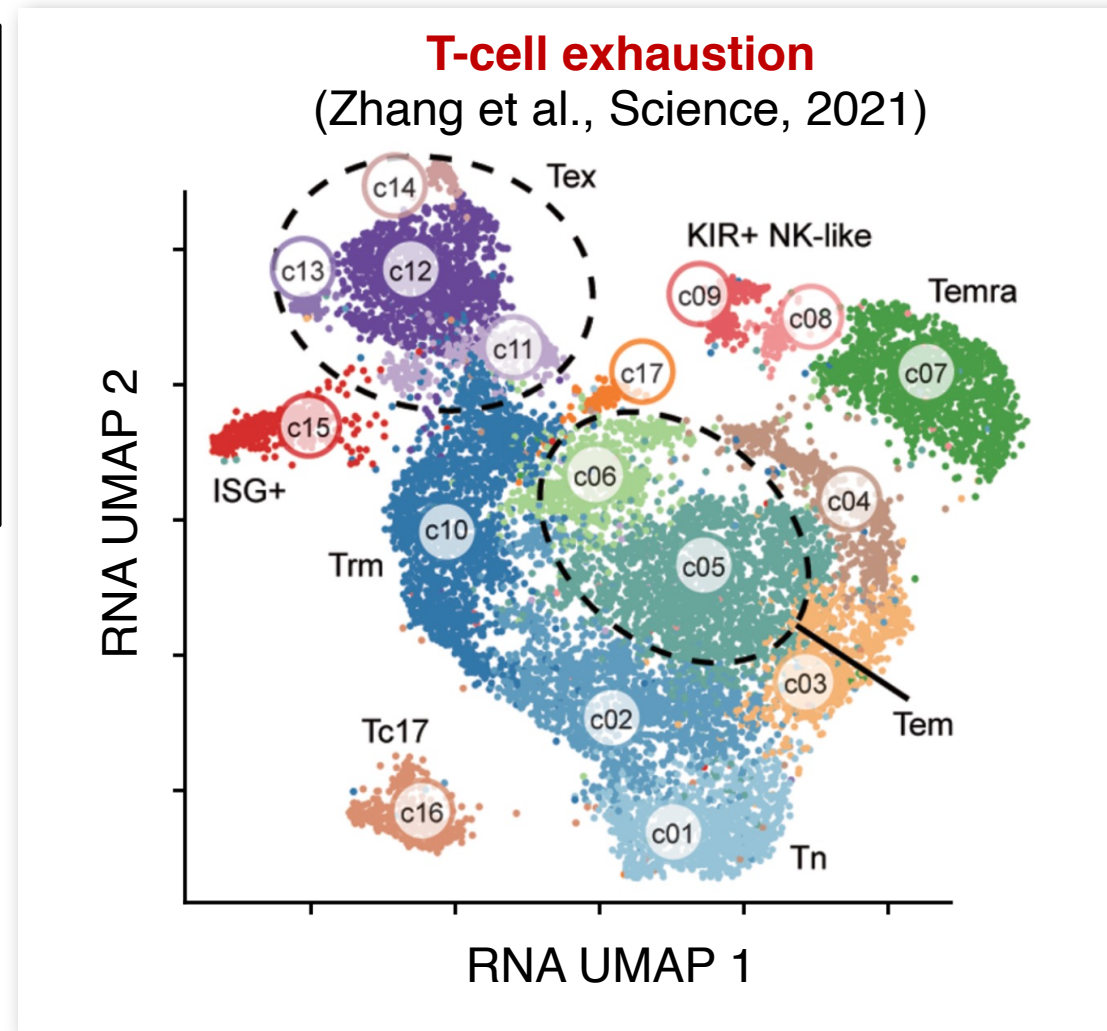
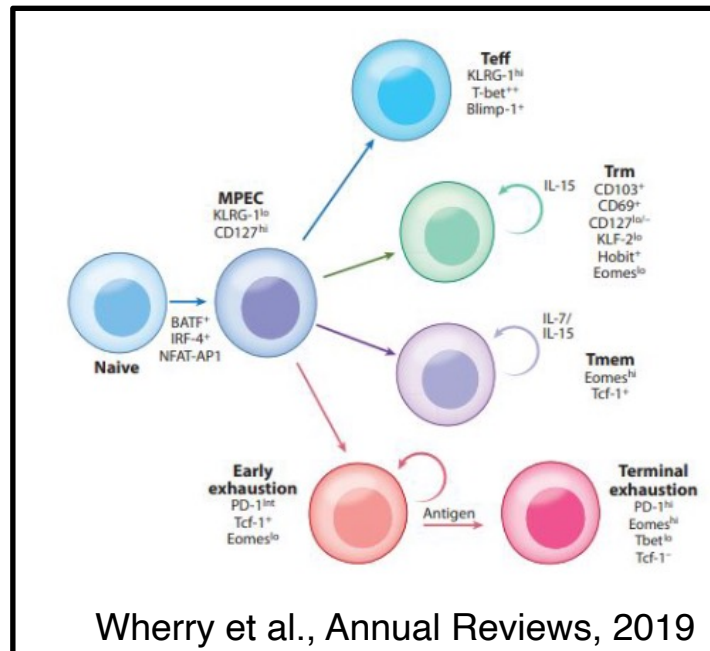
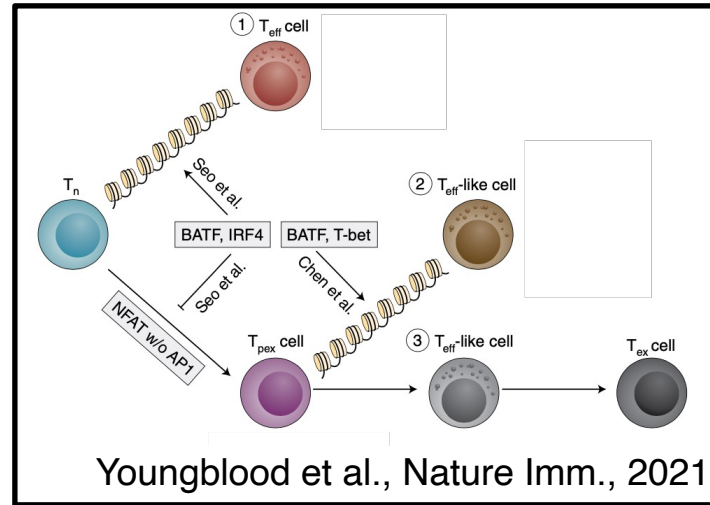
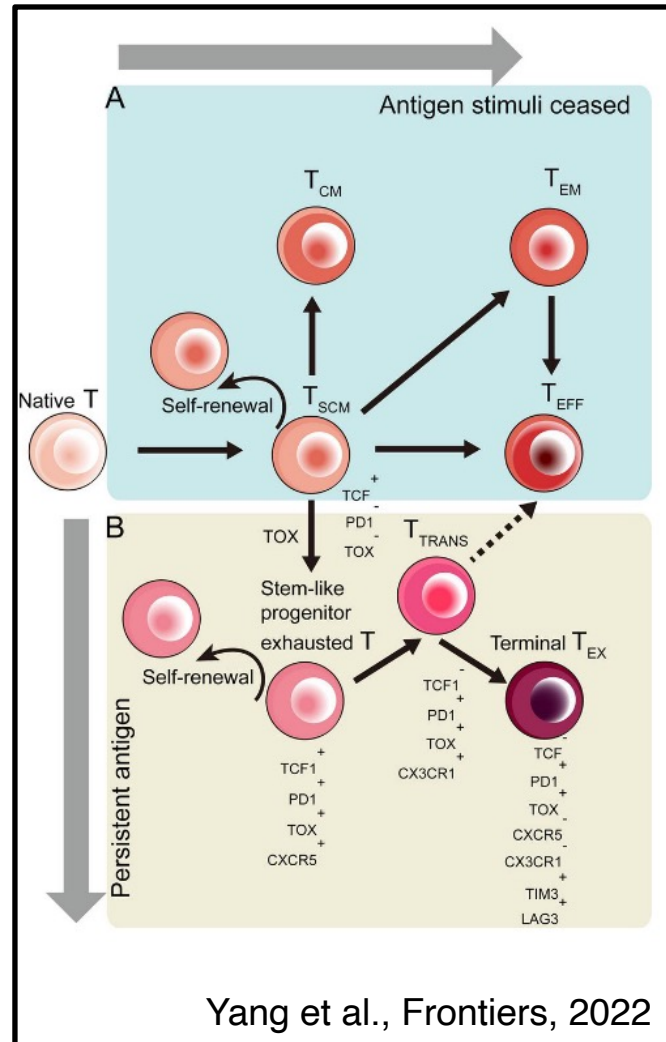
T-cell exhaustion

(Zhang et al., Science, 2021)



Even more difficult since the biology is still anyone's guess

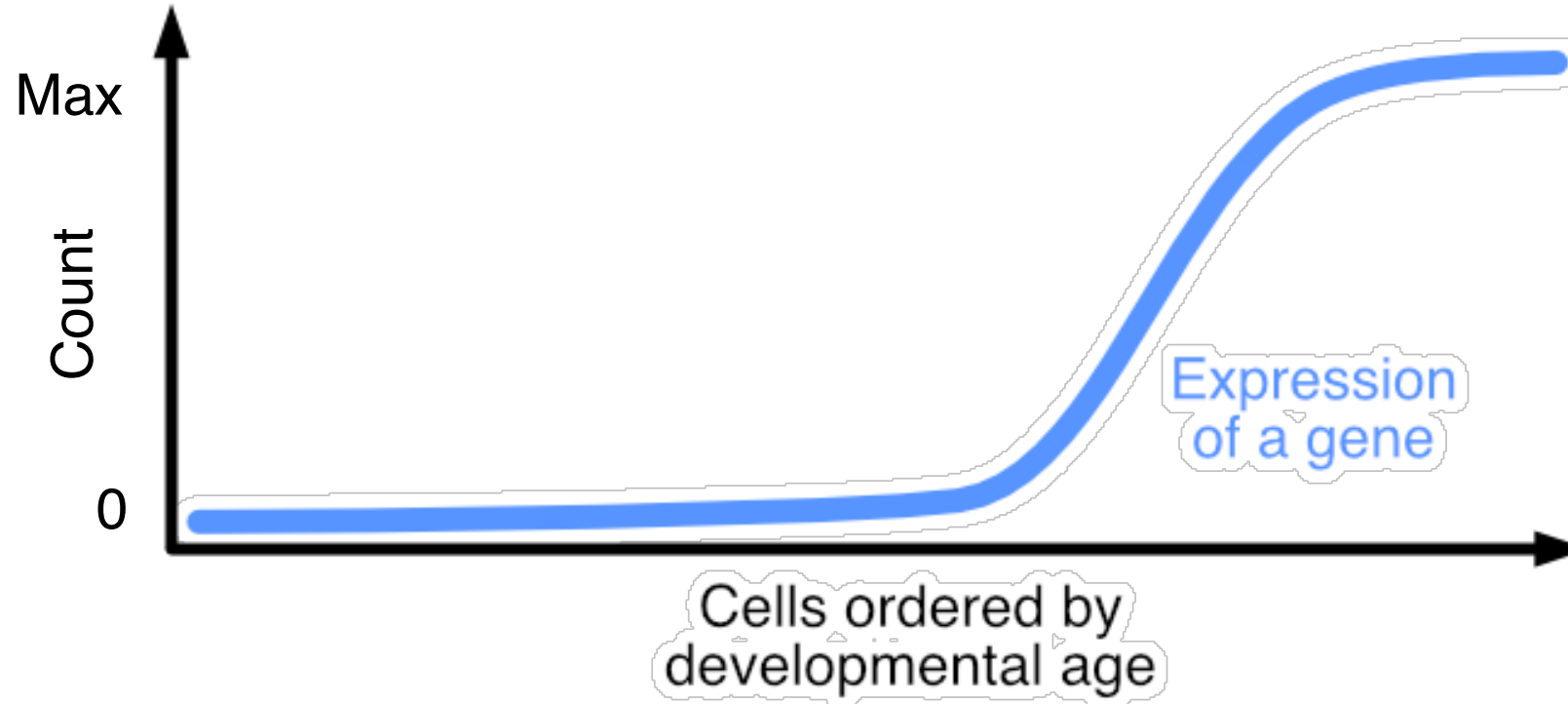
Brief aside: More complex systems



Even more difficult since the biology is still anyone's guess

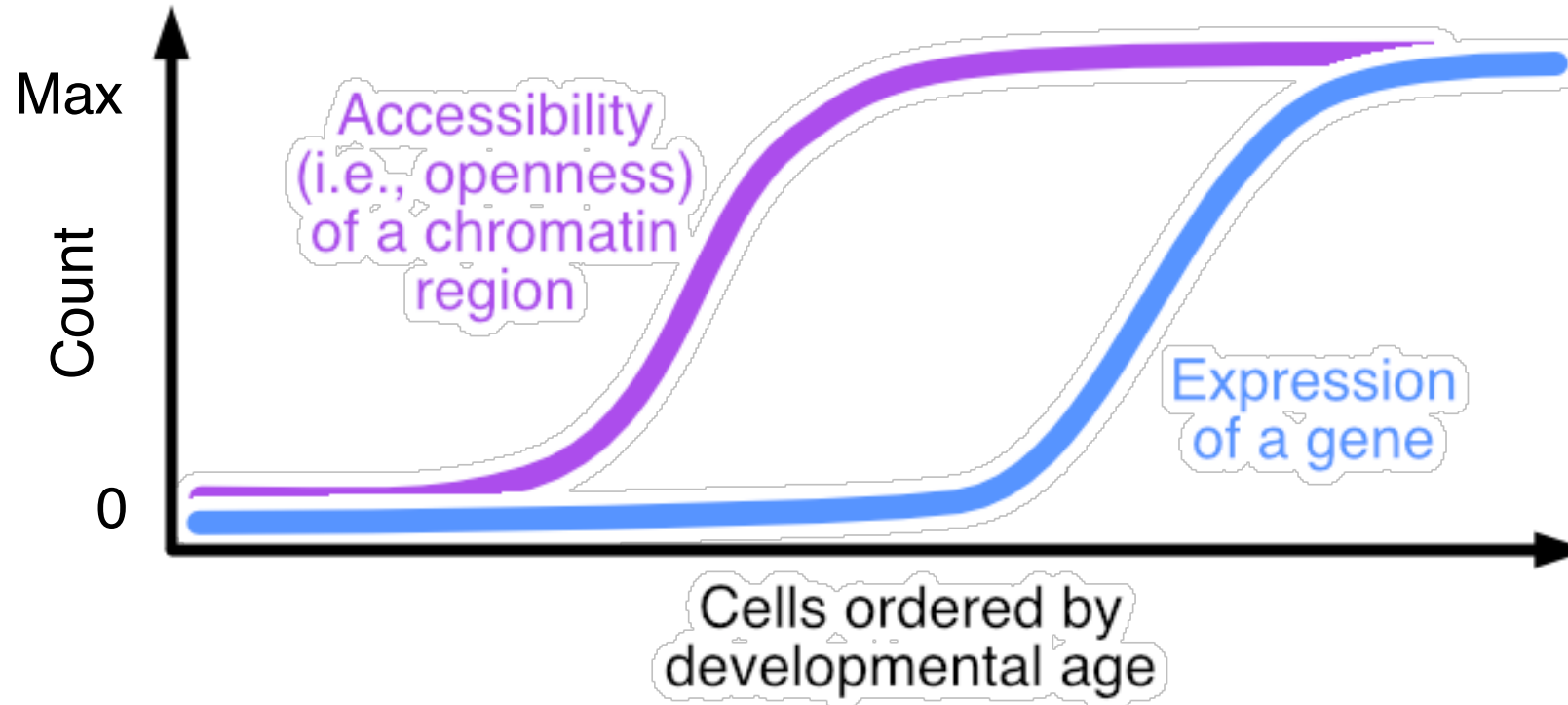
Tilted-CCA can learn which cells are in steady-state via **RNA-ATAC** relations thanks to principles discovered in previous work.

Schematic from SHARE-seq (Buenrostro et al., Cell, 2020)



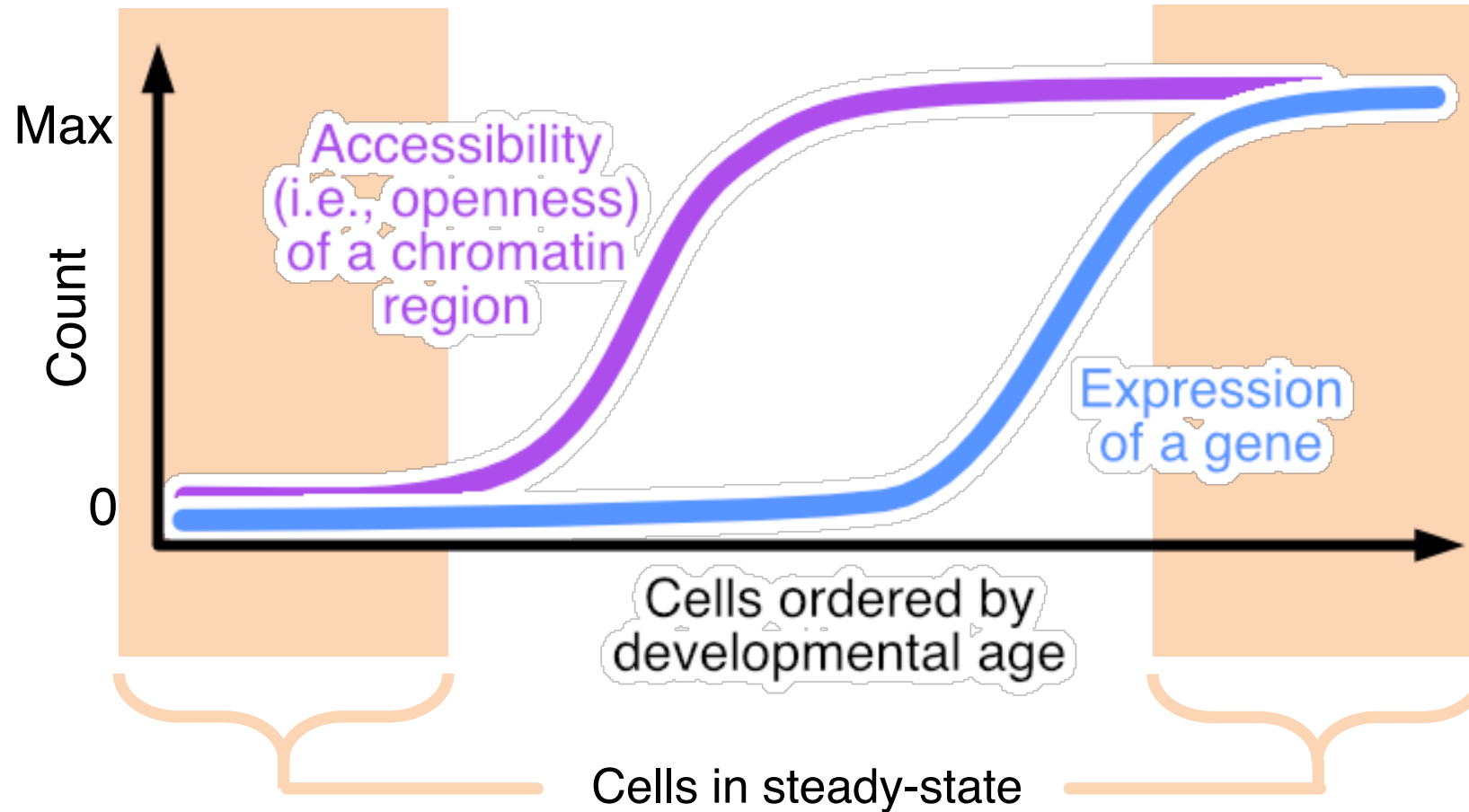
Tilted-CCA can learn which cells are in steady-state via **RNA-ATAC** relations thanks to principles discovered in previous work.

Schematic from SHARE-seq (Buenrostro et al., Cell, 2020)



Tilted-CCA can learn which cells are in steady-state via **RNA-ATAC** relations thanks to principles discovered in previous work.

Schematic from SHARE-seq (Buenrostro et al., Cell, 2020)



We want to learn which cells are in steady-state (w/o temporal ordering):

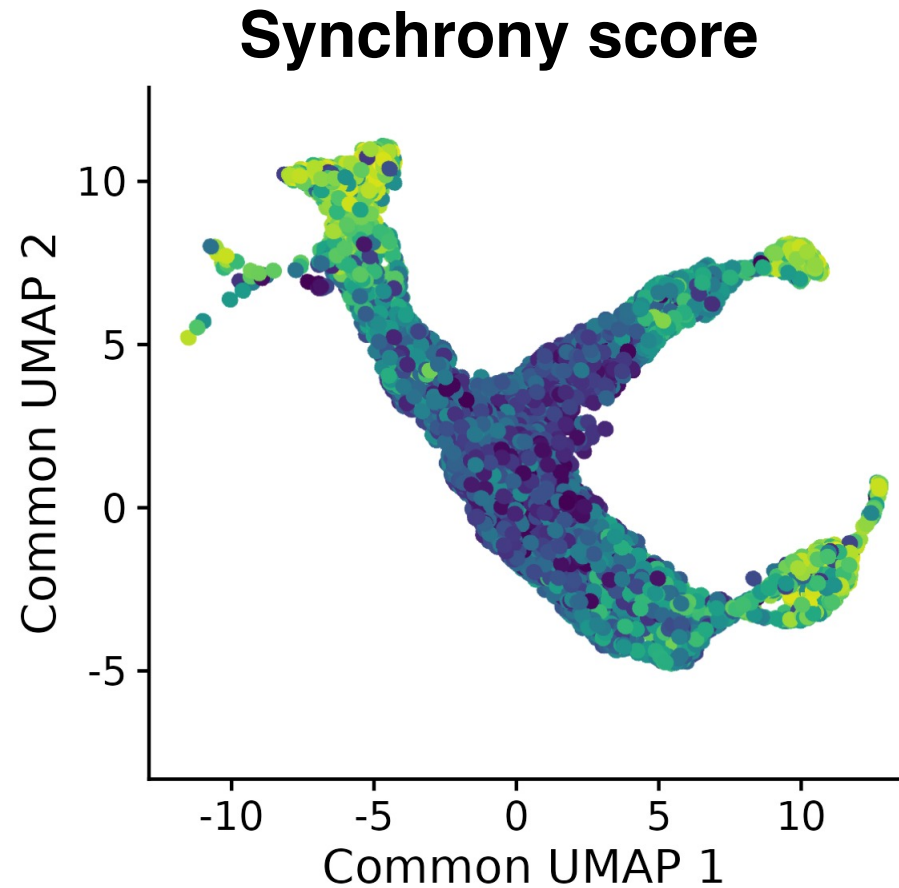
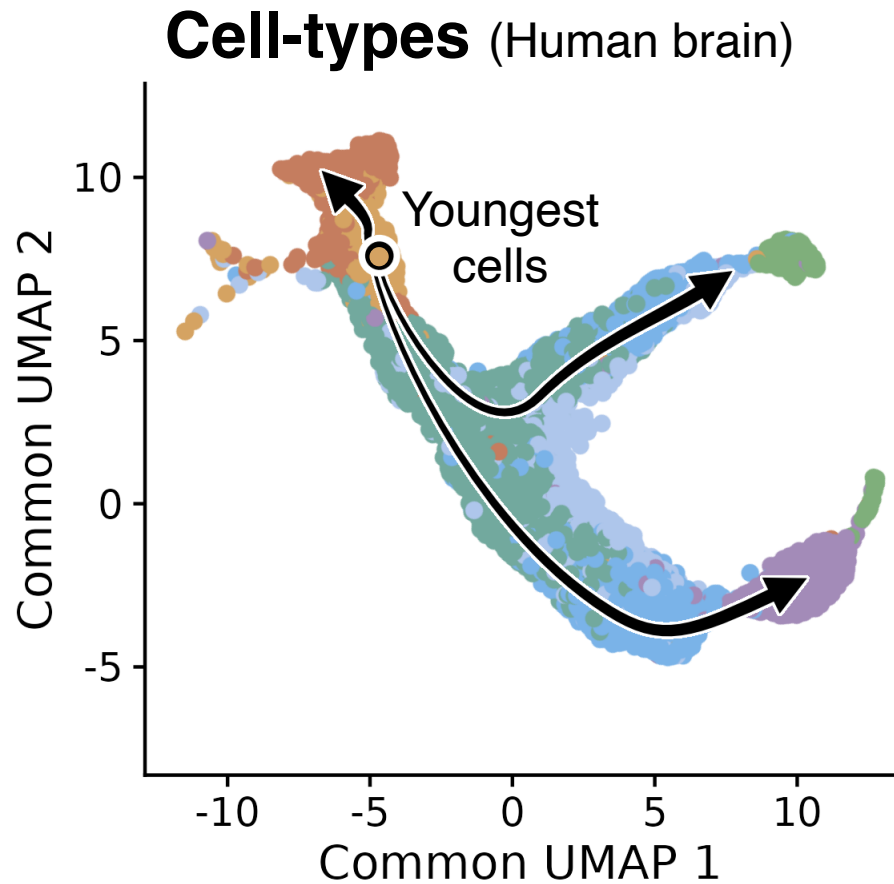
Local geometries similar in both modalities → Small distinct components

Our contribution: The synchrony score, which measures the coordination b/w genes and chromatin accessibility within a cell, to reveal if a cell is in steady-state

After fitting Tilted-CCA, for each cell i : Correlation $\left(\begin{array}{l} \text{common} \\ \text{component} \\ \text{among genes} \end{array} , \begin{array}{l} \text{common + distinct} \\ \text{components among} \\ \text{genes} \end{array} \right)$

Our contribution: The synchrony score, which measures the coordination b/w genes and chromatin accessibility within a cell, to reveal if a cell is in steady-state

After fitting Tilted-CCA, for each cell i : Correlation $\left(\begin{array}{l} \text{common} \\ \text{component} \\ \text{among genes} \end{array}, \begin{array}{l} \text{common + distinct} \\ \text{components among} \\ \text{genes} \end{array} \right)$



High synchrony score
(i.e., steady-state)

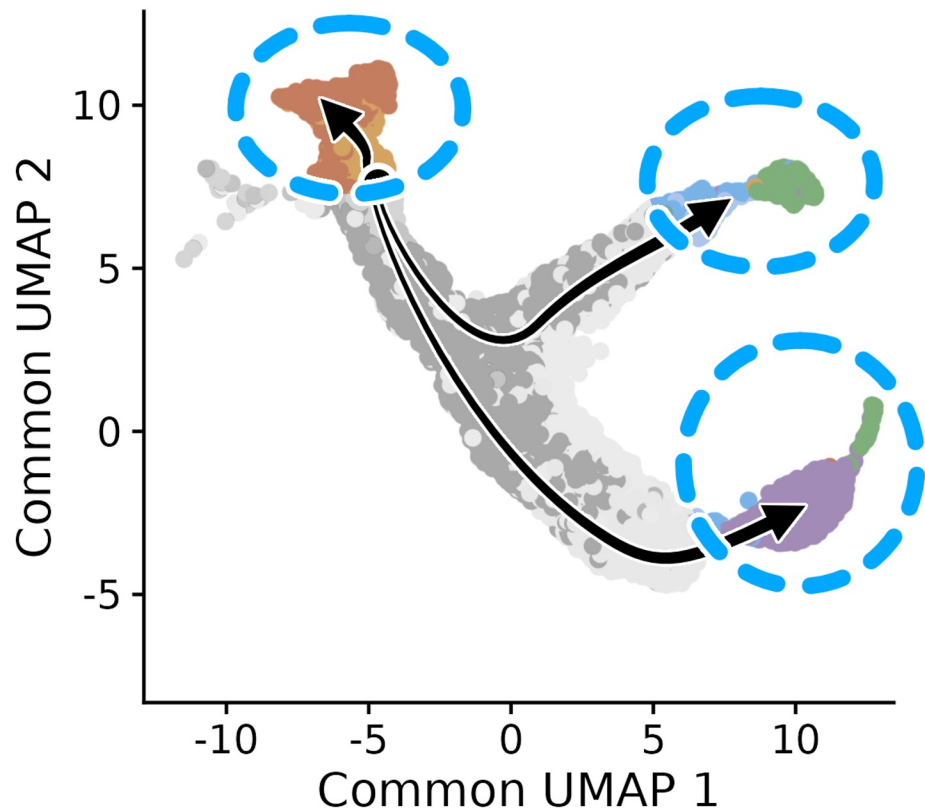


Low synchrony score
(i.e., undergoing
development)

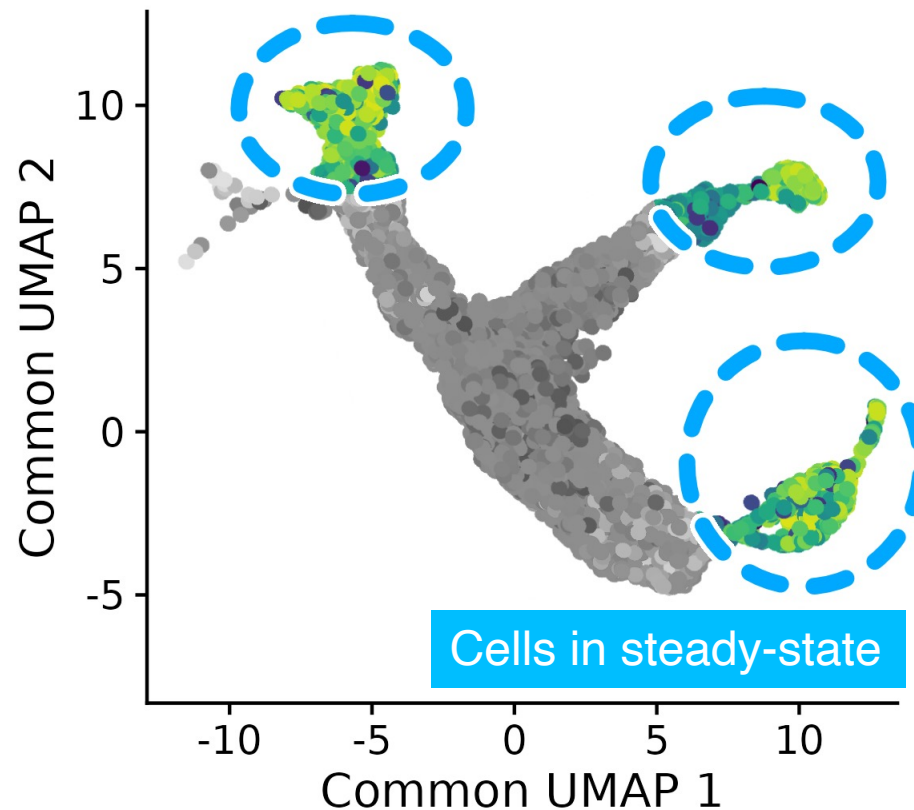
Our contribution: The synchrony score, which measures the coordination b/w genes and chromatin accessibility within a cell, to reveal if a cell is in steady-state

After fitting Tilted-CCA, for each cell i : Correlation $\left(\begin{array}{l} \text{common} \\ \text{component} \\ \text{among genes} \end{array}, \begin{array}{l} \text{common + distinct} \\ \text{components among} \\ \text{genes} \end{array} \right)$

Cell-types (Human brain)



Synchrony score



High synchrony score
(i.e., steady-state)

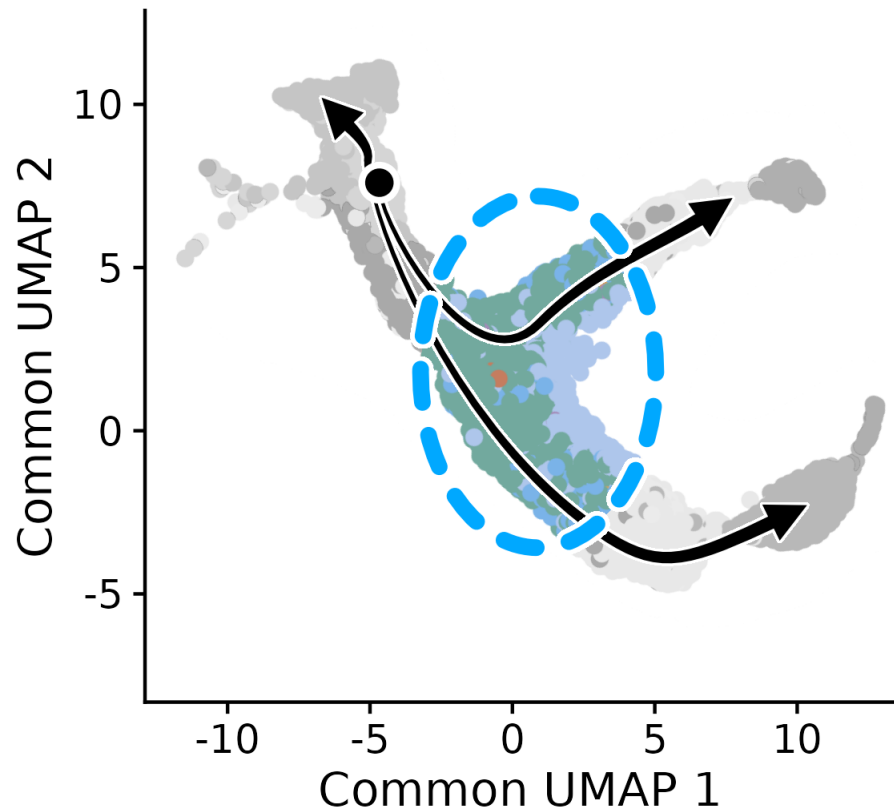


Low synchrony score
(i.e., undergoing
development)

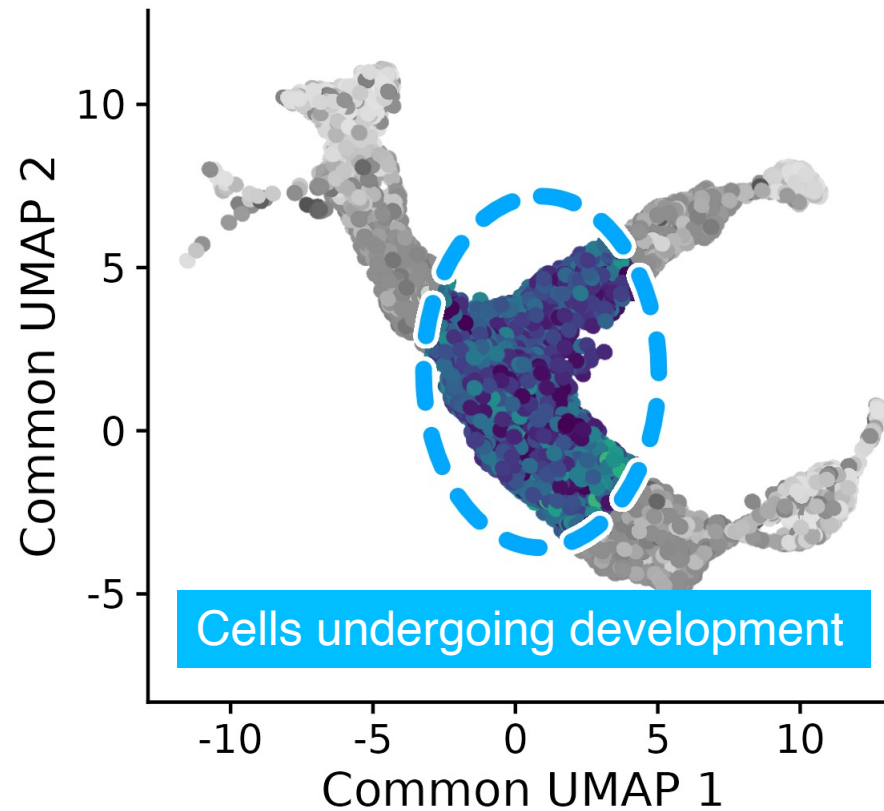
Our contribution: The synchrony score, which measures the coordination b/w genes and chromatin accessibility within a cell, to reveal if a cell is in steady-state

After fitting Tilted-CCA, for each cell i : Correlation $\left(\begin{array}{l} \text{common} \\ \text{component} \\ \text{among genes} \end{array}, \begin{array}{l} \text{common + distinct} \\ \text{components among} \\ \text{genes} \end{array} \right)$

Cell-types (Human brain)



Synchrony score



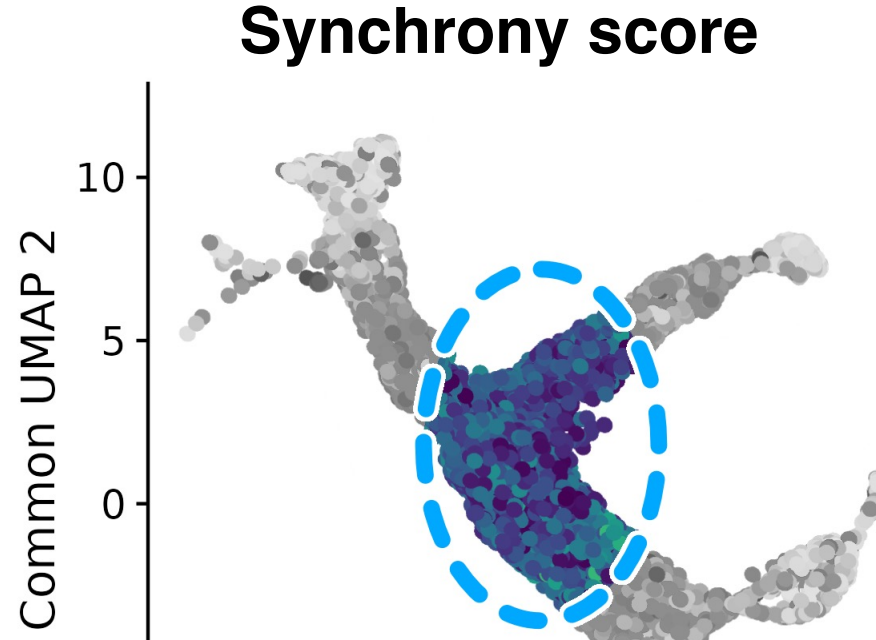
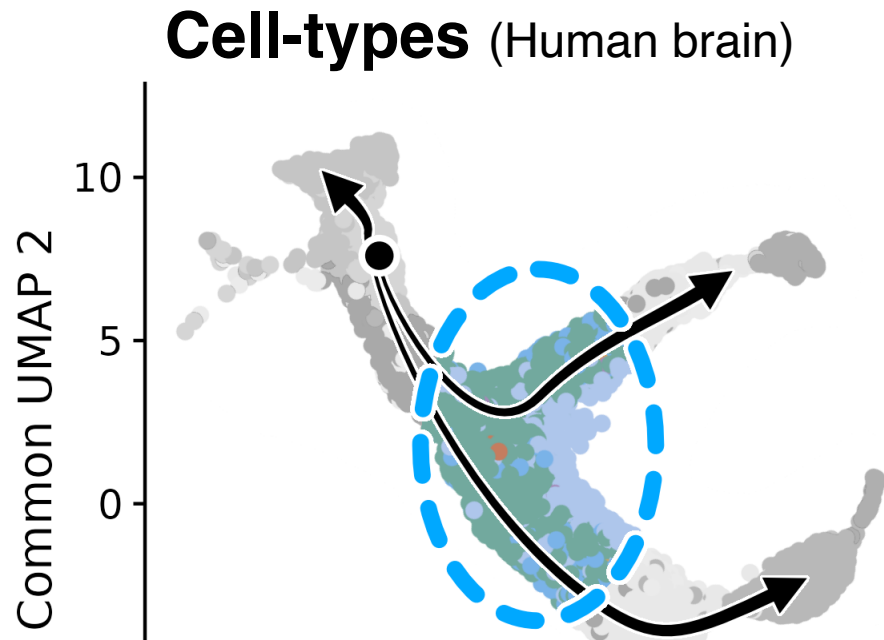
High synchrony score
(i.e., steady-state)



Low synchrony score
(i.e., undergoing
development)

Our contribution: The synchrony score, which measures the coordination b/w genes and chromatin accessibility within a cell, to reveal if a cell is in steady-state

After fitting Tilted-CCA, for each cell i : Correlation $\left(\begin{array}{l} \text{common} \\ \text{component} \\ \text{among genes} \end{array}, \begin{array}{l} \text{common + distinct} \\ \text{components among} \\ \text{genes} \end{array} \right)$



High synchrony score
(i.e., steady-state)



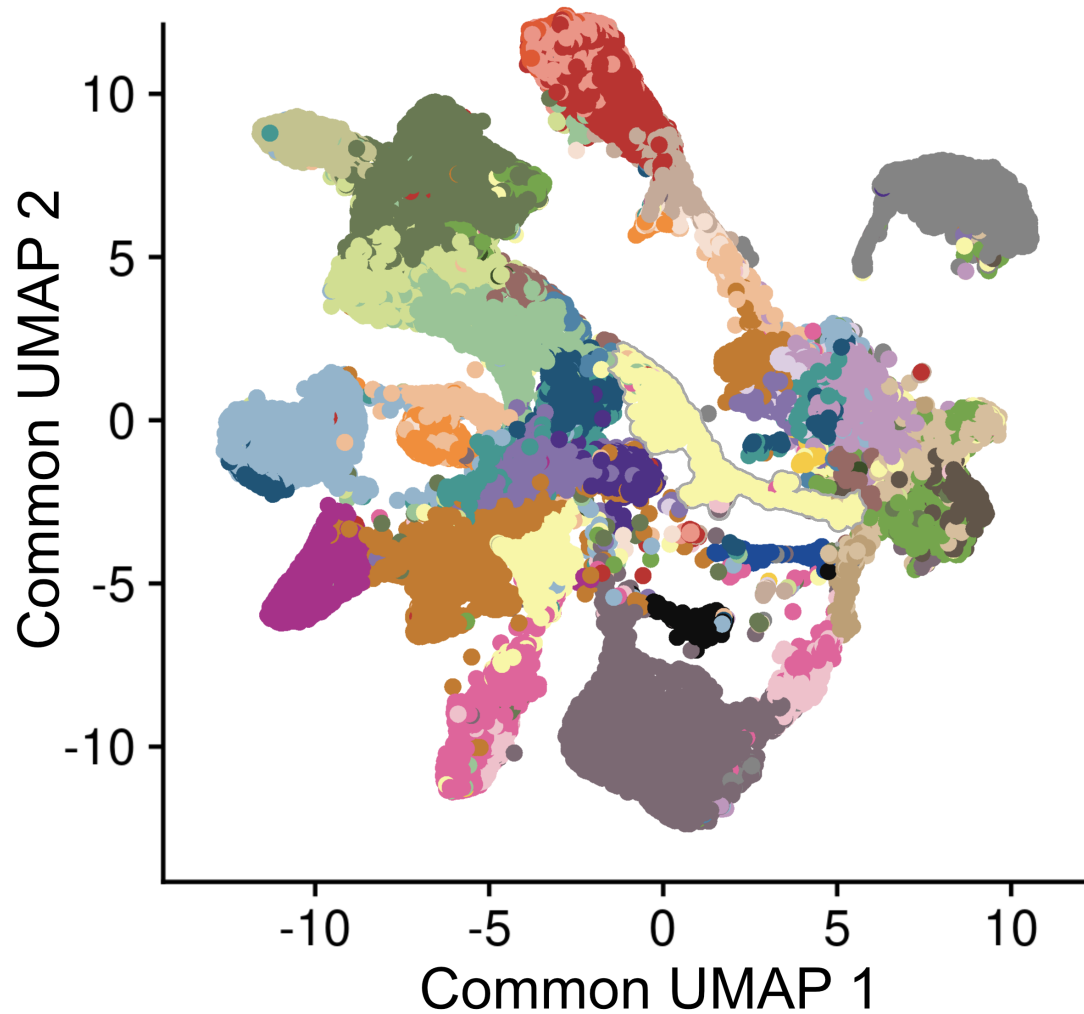
Low synchrony score
(i.e., undergoing
development)

No trajectory estimation needed for the synchrony scores!
Learning the “highways” from “the speed” rely on different principles.

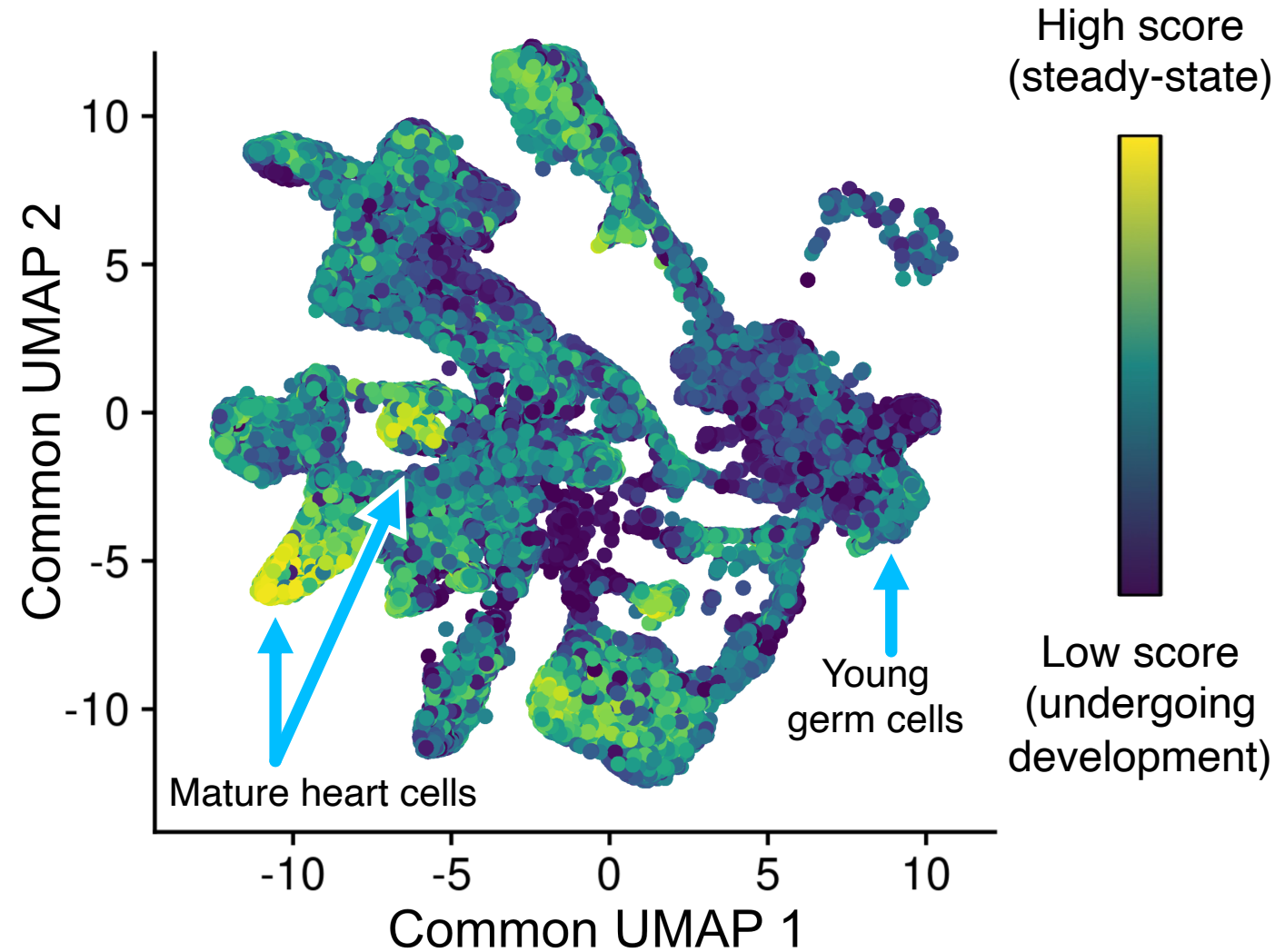
We demonstrate the utility of the synchrony scores on other harder systems.

We demonstrate the utility of the synchrony scores on other harder systems.

Cell-types (Mouse embryo)

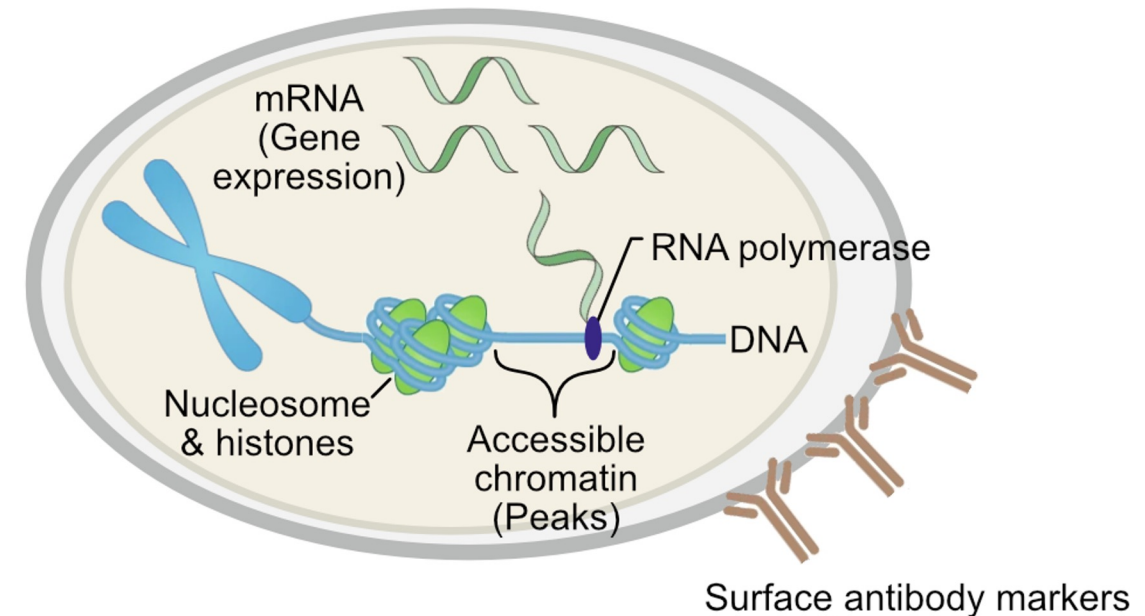


Synchrony score



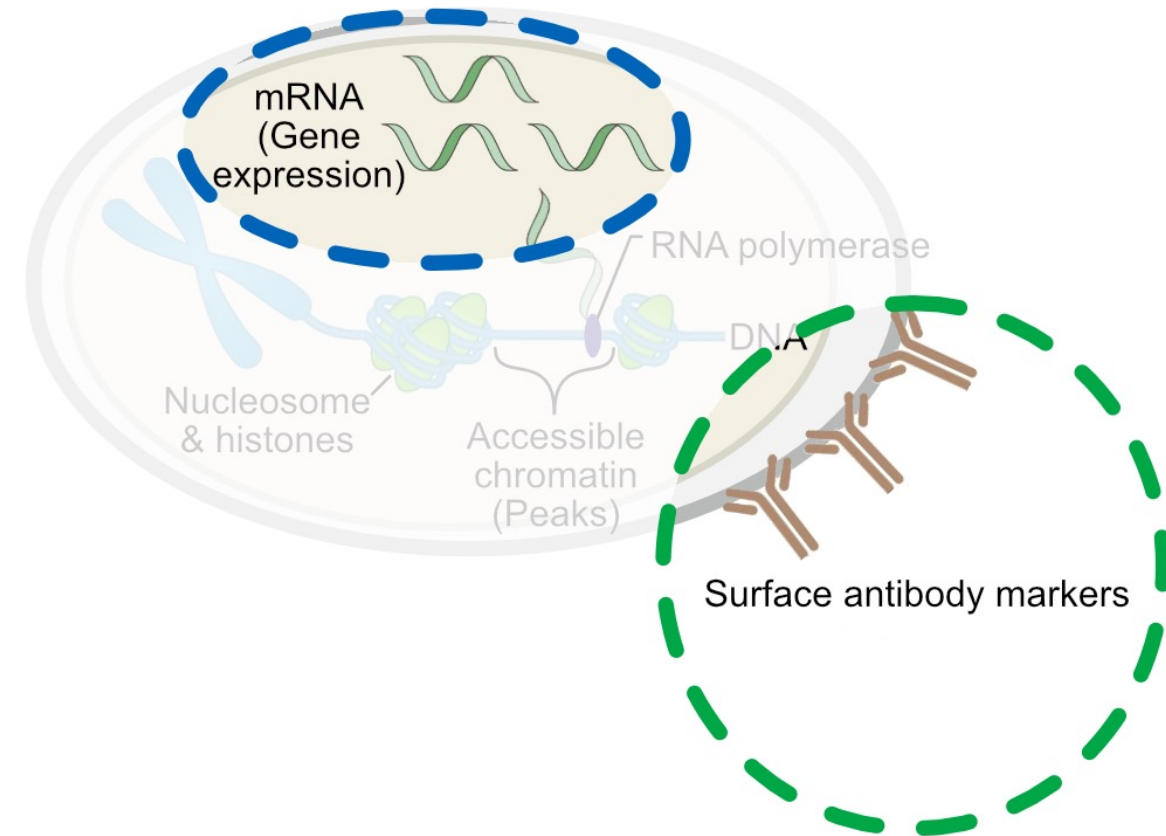
Recap: By matrix factorizing multi-modal data based on the common/distinct geometry, Tilted-CCA enables new perspectives to understand cell biology.

1. **(Experimental design):** Which pair of modalities should biologist sequence to have the most comprehensive understanding?
2. **(Variable selection):** For RNA-Protein data, how can we pick the antibodies that contribute the most additional information to the RNA modality?
3. **(Developmental biology):** Can the amount of coordination between two modalities tell us if a cell in a steady-state or is undergoing development?



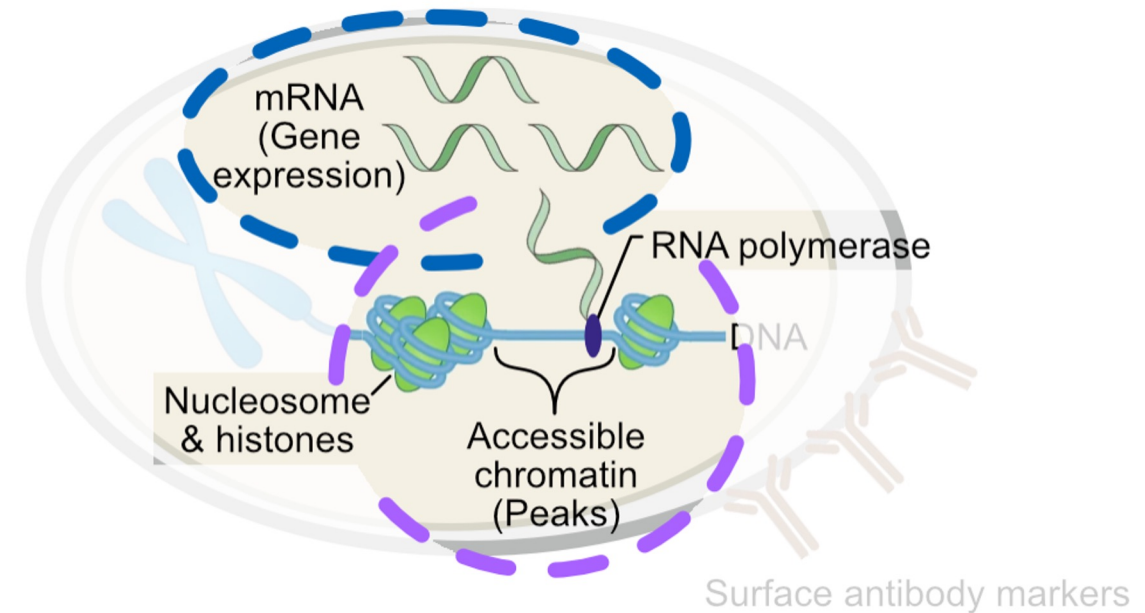
Recap: By matrix factorizing multi-modal data based on the common/distinct geometry, Tilted-CCA enables new perspectives to understand cell biology.

1. **(Experimental design):** Which pair of modalities should biologist sequence to have the most comprehensive understanding?
2. **(Variable selection):** For RNA-Protein data, how can we pick the antibodies that contribute the most additional information to the RNA modality?
3. **(Developmental biology):** Can the amount of coordination between two modalities tell us if a cell in a steady-state or is undergoing development?



Recap: By matrix factorizing multi-modal data based on the common/distinct geometry, Tilted-CCA enables new perspectives to understand cell biology.

1. **(Experimental design):** Which pair of modalities should biologist sequence to have the most comprehensive understanding?
2. **(Variable selection):** For RNA-Protein data, how can we pick the antibodies that contribute the most additional information to the RNA modality?
3. **(Developmental biology):** Can the amount of coordination between two modalities tell us if a cell in a steady-state or is undergoing development?



Recap: By matrix factorizing multi-modal data based on the common/distinct geometry, Tilted-CCA enables new perspectives to understand cell biology.

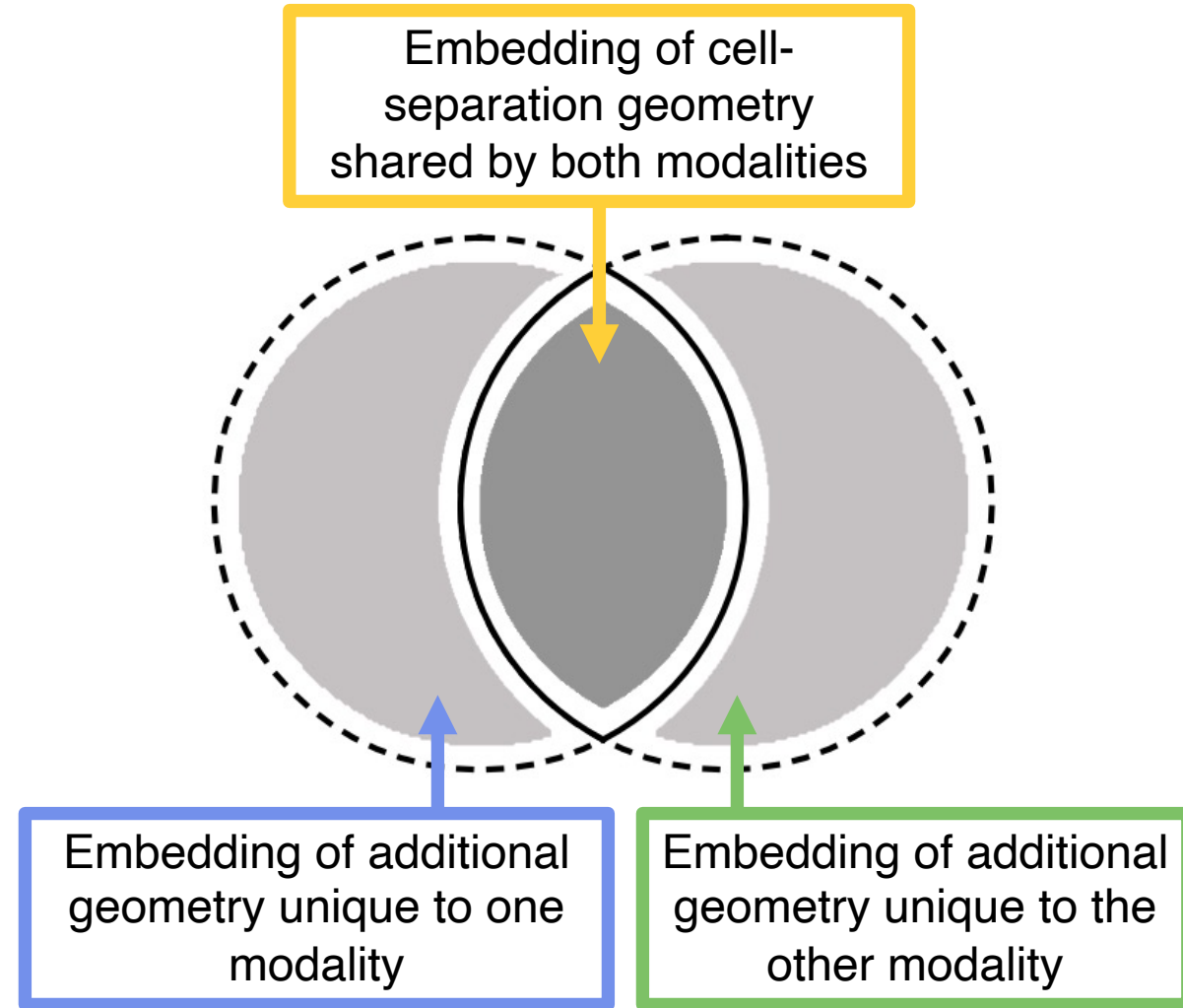
1. **(Experimental design):** Which pair of modalities should biologist sequence to have the most comprehensive understanding?
2. **(Variable selection):** For RNA-Protein data, how can we pick the antibodies that contribute the most additional information to the RNA modality?
3. **(Developmental biology):** Can the amount of coordination between two modalities tell us if a cell is in a steady-state or is undergoing development?

Multi-modal analyses will only become more prevalent in biomedical research, and this will raise new biological theories that will require new matrix factorization approaches to study.

Takeaways:

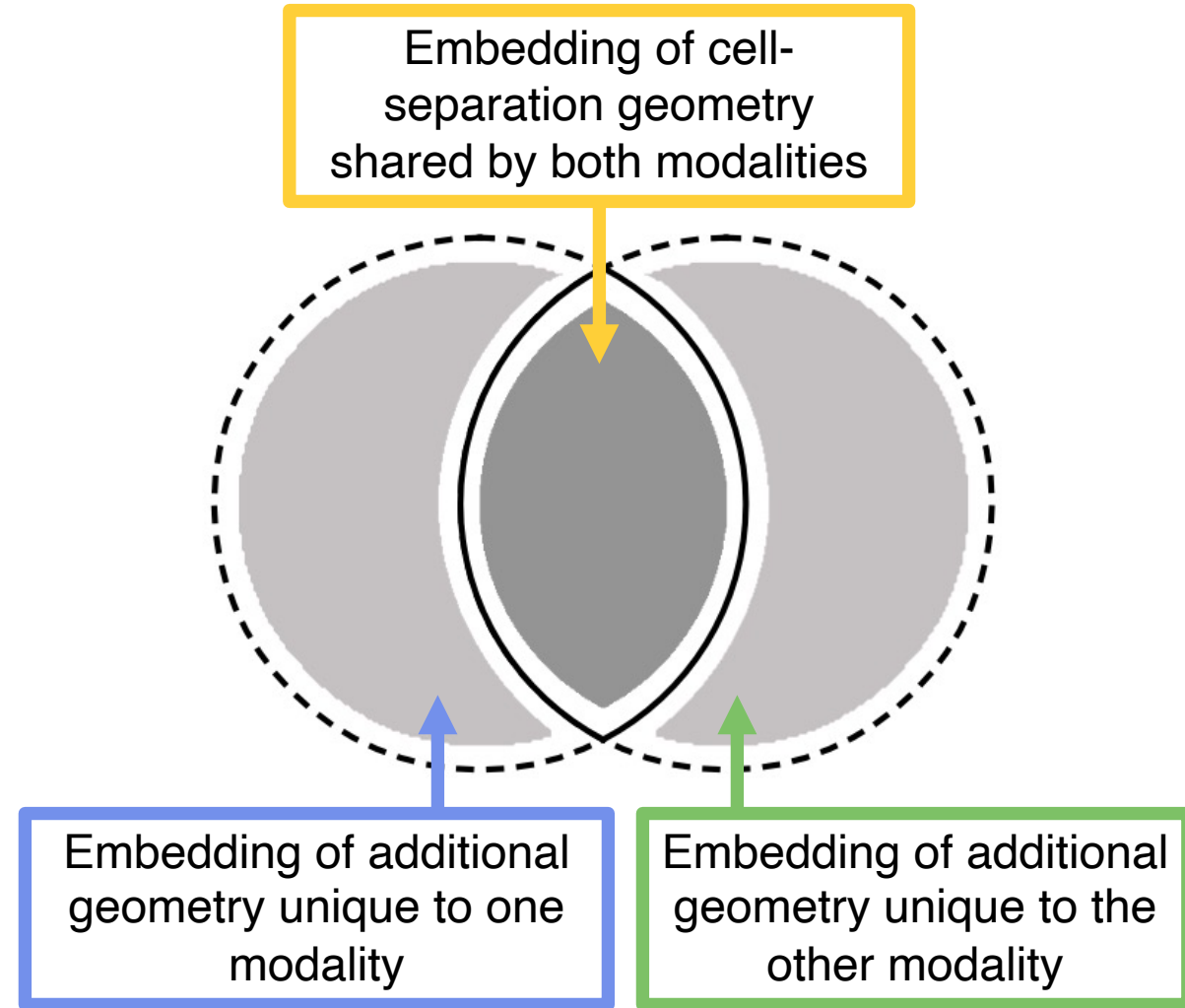
Takeaways:

- **Novelty:** Estimating the “intersection of information” in multi-modal data, as opposed many existing methods for “union of information”
- **Strategy:** Matrix-factorization to estimate embeddings of the common and distinct geometries
- **Specifically:** Build upon CCA by “tilting” of the common vector to approximate the desired shared geometry



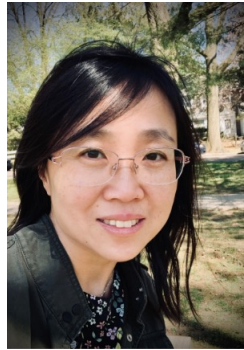
Takeaways:

- **Novelty:** Estimating the “intersection of information” in multi-modal data, as opposed many existing methods for “union of information”
- **Strategy:** Matrix-factorization to estimate embeddings of the common and distinct geometries
- **Specifically:** Build upon CCA by “tilting” of the common vector to approximate the desired shared geometry
- **Future theory:** What is the “population” shared geometry?
- **Future biology:** Cancer biology, where we need to understand the specifics of cellular mechanisms



Thank you!

Tilted-CCA on Biorxiv: 2022.10.07.511320



Nancy
Zhang



Sydney
Shaffer



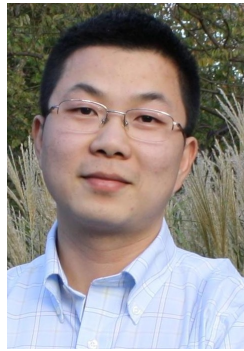
Andy
Minn



E. John
Wherry



Kathryn
Roeder



Jing
Lei



Ryan
Tibshirani



Alessandro
Rinaldo



Han
Liu



Max
G'Sell



James
Sharpnack



Sangwon
Hyun