

# **Natively Interpretable *t*-SNE**

# Neighbor Embedding

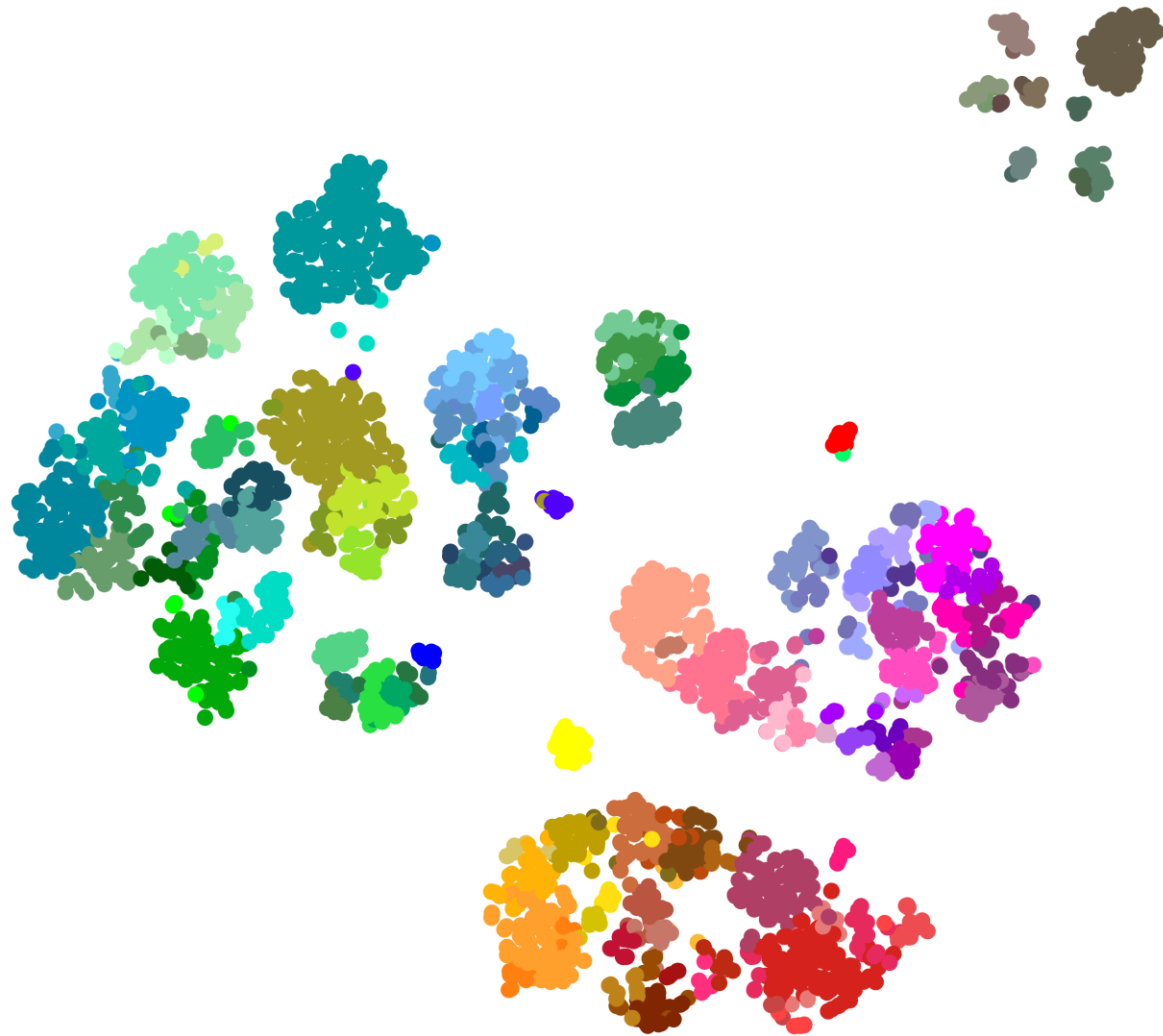


... and many others !

Mouse cortex dataset [1]

Single cell RNA sequencing

Dim = ~ 9500



[1] Tasic, Bosiljka, et al. "Shared and distinct transcriptomic cell types across neocortical areas." *Nature* 563.7729 (2018)

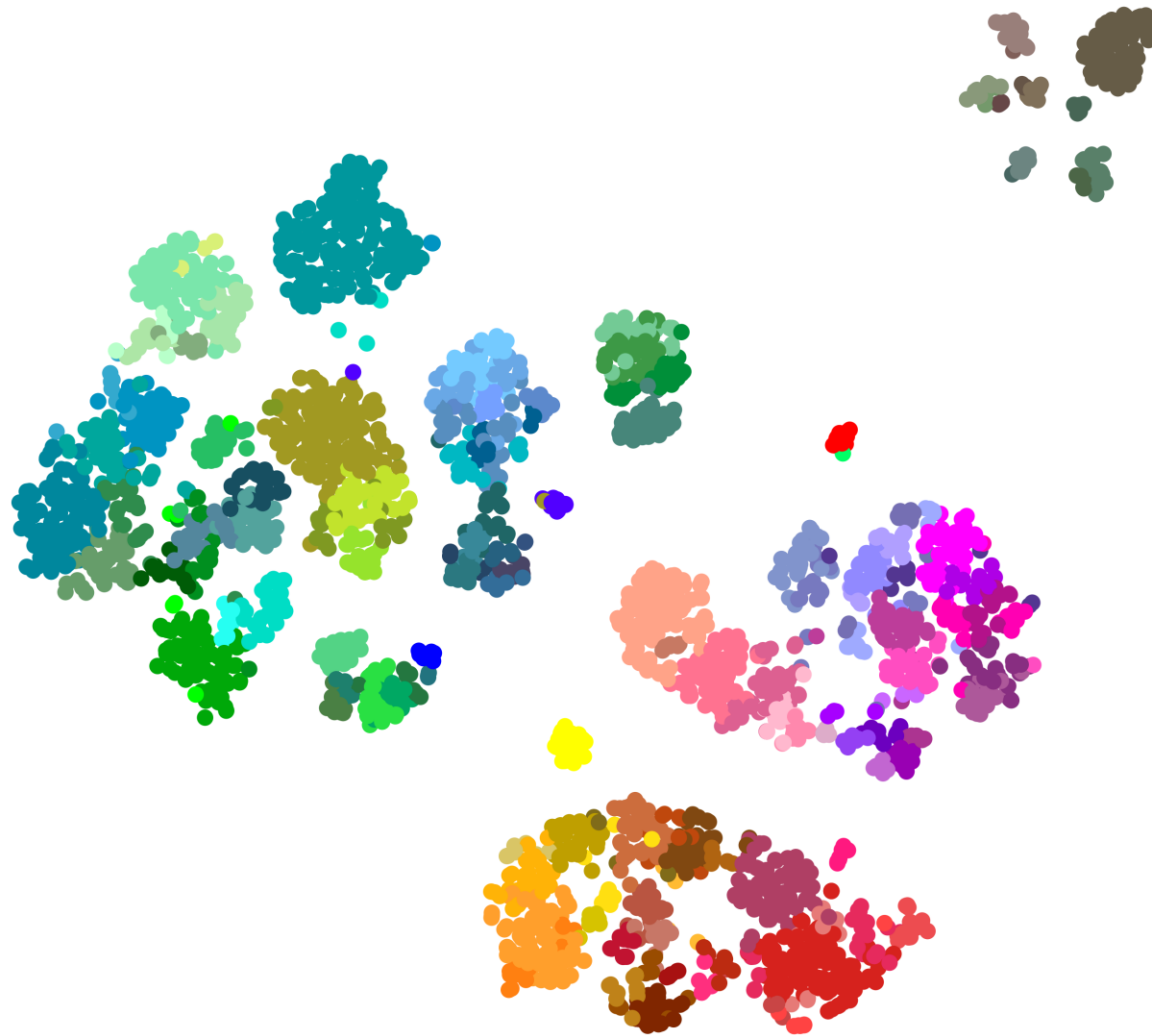
Mouse cortex dataset [1]

Single cell RNA sequencing

Dim = ~ 9500

Preprocessing [2]

Dim = 50



[2] Kobak, Dmitry, and Philipp Berens. "The art of using t-SNE for single-cell transcriptomics." *Nature communications* 10.1 (2019)

Mouse cortex dataset [1]

Single cell RNA sequencing

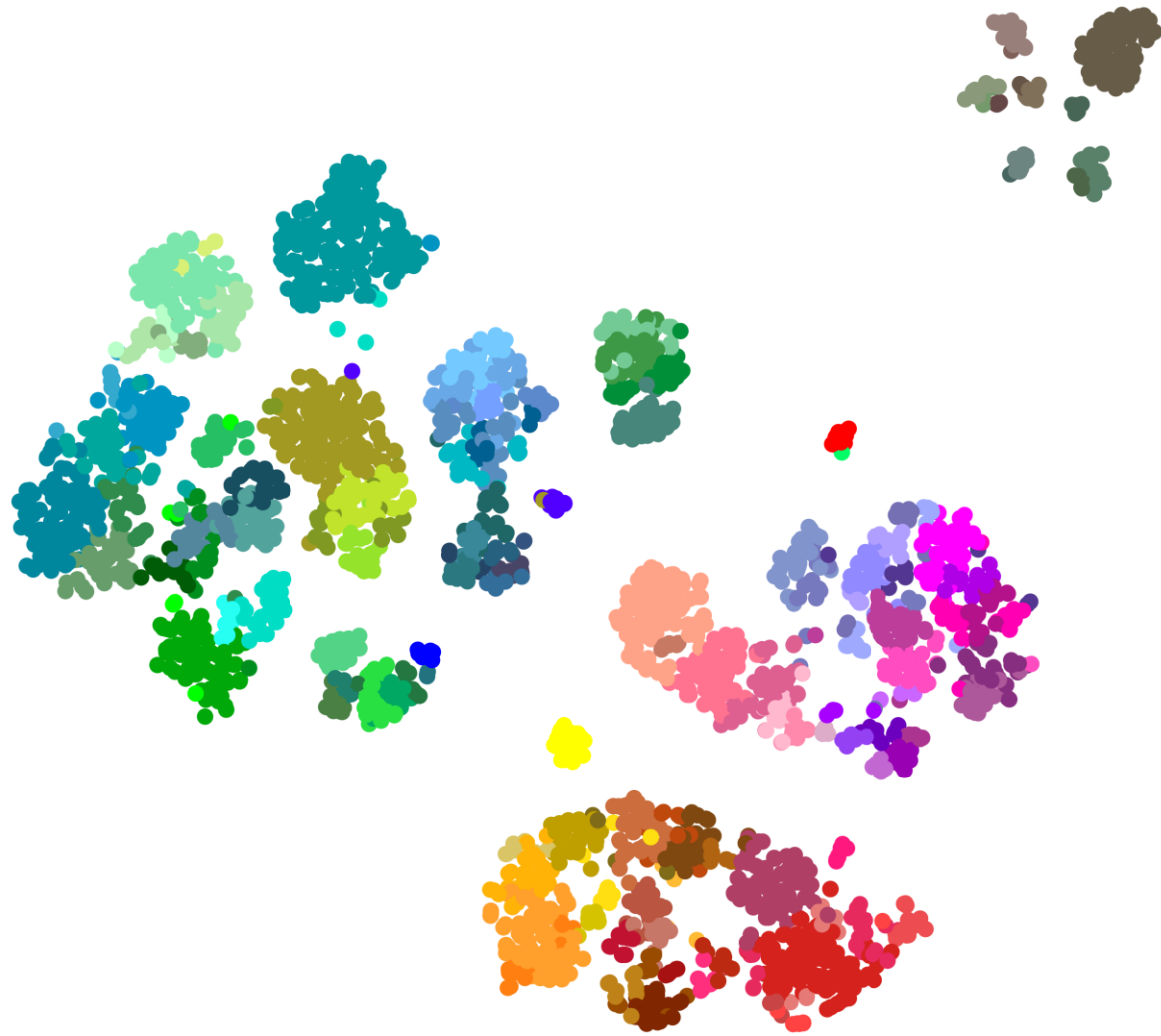
Dim = ~ 9500

Preprocessing [2]

Dim = 50 (HD)

Ms.t-SNE

Dim = 2 (LD)



Mouse cortex dataset [1]

Single cell RNA sequencing

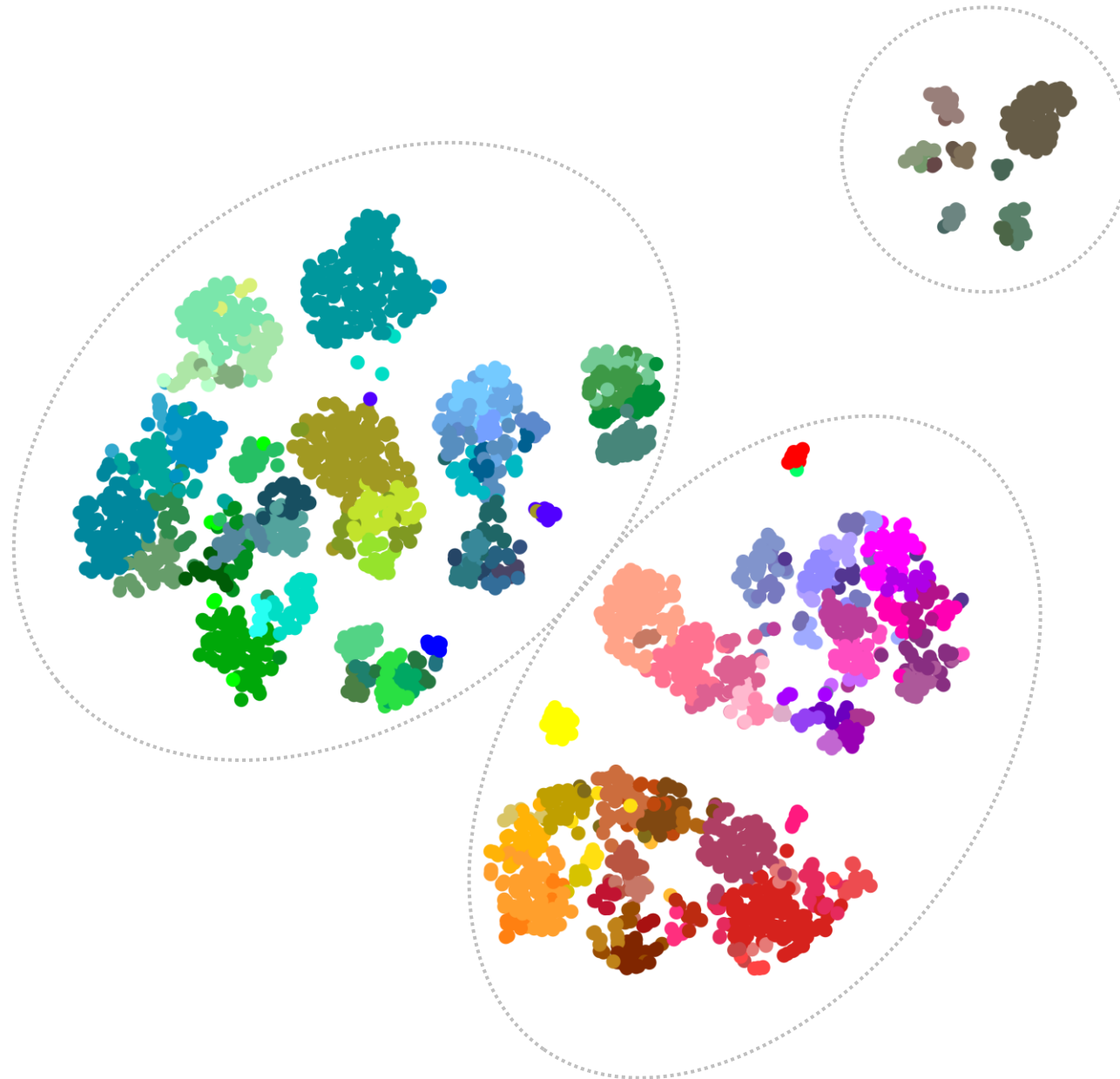
Dim = ~ 9500

Preprocessing [2]

Dim = 50 (HD)

Ms.t-SNE

Dim = 2 (LD)



Excitatory neurons

Inhibitory neurons

Non neural cells

Mouse cortex dataset [1]

Single cell RNA sequencing

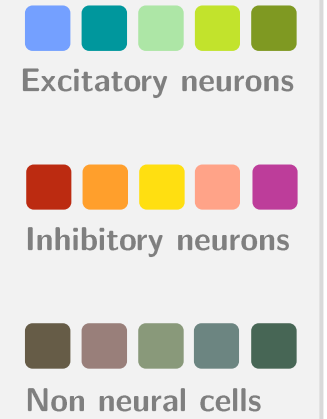
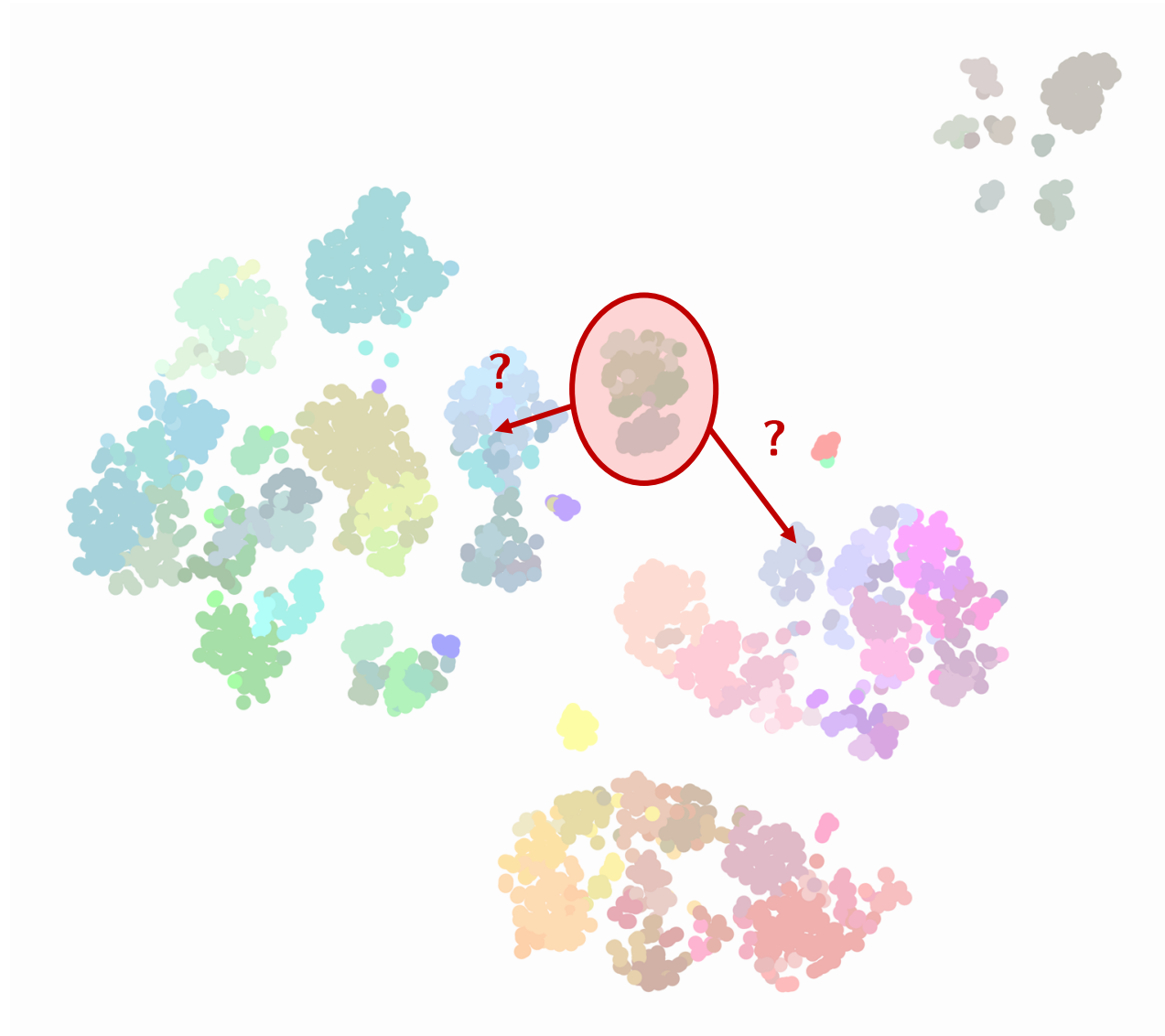
Dim = ~ 9500

Preprocessing [2]

Dim = 50 (HD)

Ms.t-SNE

Dim = 2 (LD)



Mouse cortex dataset [1]

Single cell RNA sequencing

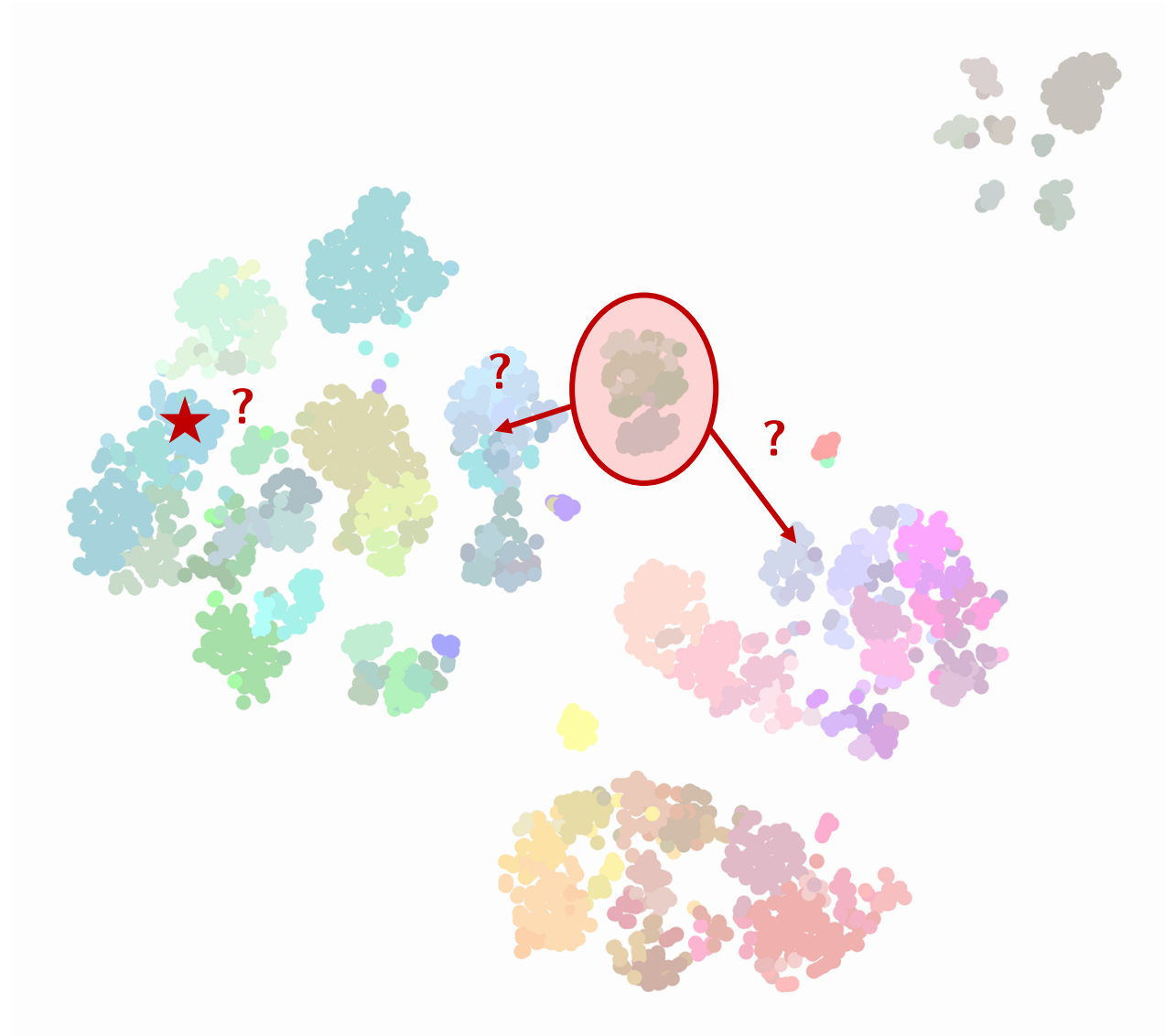
Dim = ~ 9500

Preprocessing [2]

Dim = 50 (HD)

Ms.t-SNE

Dim = 2 (LD)



Excitatory neurons

Inhibitory neurons

Non neural cells



Mouse cortex dataset [1]

Single cell RNA sequencing

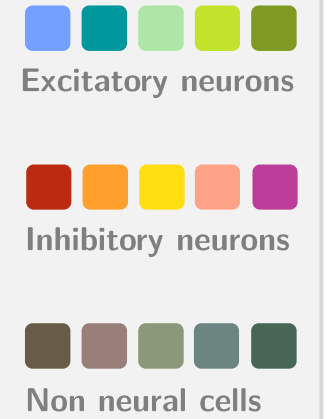
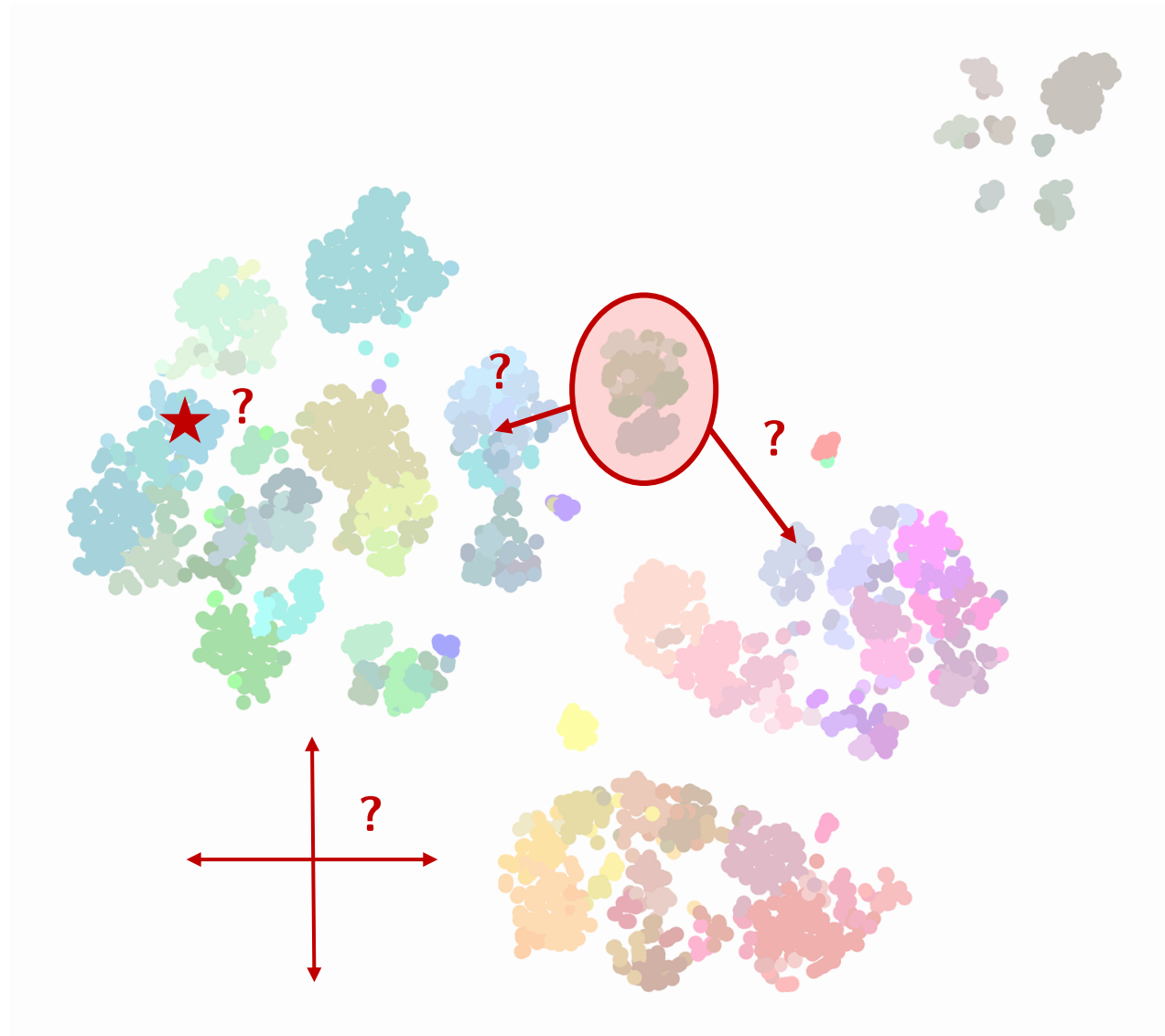
Dim = ~ 9500

Preprocessing [2]

Dim = 50 (HD)

Ms.t-SNE

Dim = 2 (LD)



## *M*-dim space

Data

$$\{\xi_i\}_{i=1}^N$$

## ***M*-dim space**

Data

$$\{\xi_i\}_{i=1}^N$$

## **2-dim space**

Initialization

$$\{x_i\}_{i=1}^N$$

## ***M*-dim space**

Data

$$\{\xi_i\}_{i=1}^N$$

Similarities

$$\tau_{ij}$$

## **2-dim space**

Initialization

$$\{x_i\}_{i=1}^N$$

Similarities

$$t_{ij}$$

## M-dim space

Data

$$\{\xi_i\}_{i=1}^N$$

Similarities

$$\tau_{ij}$$

## 2-dim space

Initialization

$$\{x_i\}_{i=1}^N$$

Similarities

$$t_{ij}$$

Loss function

$$C = \sum_{i \in \mathcal{I}} d(\tau_i, t_i)$$

## M-dim space

Data

$$\{\xi_i\}_{i=1}^N$$

Similarities

$$\tau_{ij}$$

## 2-dim space

Initialization

$$\{x_i\}_{i=1}^N$$

Similarities

$$t_{ij}$$

Updates

$$x_i^k = x_i^{k-1} - r \cdot \nabla_{x_i} C$$

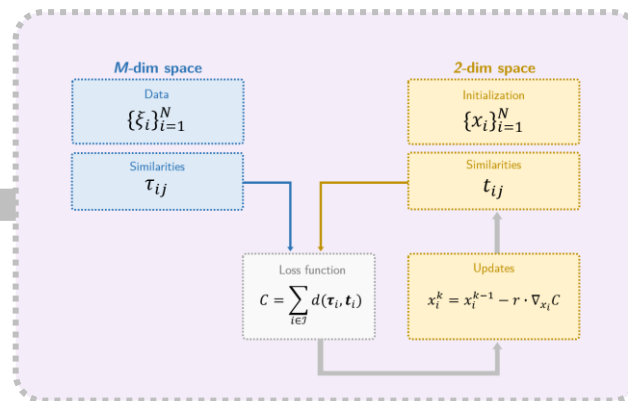
Loss function

$$C = \sum_{i \in \mathcal{I}} d(\tau_i, t_i)$$

## HD Features

$$\{\xi_i\}_{i=1}^N$$

## Neighbor Embedding



## LD Coordinates

$$\{x_i\}_{i=1}^N$$

HD Features

$$\{\xi_i\}_{i=1}^N$$

Neighbor Embedding

$$f(\cdot) ?$$

LD Coordinates

$$\{x_i\}_{i=1}^N$$





Natively Interpretable *t*-SNE

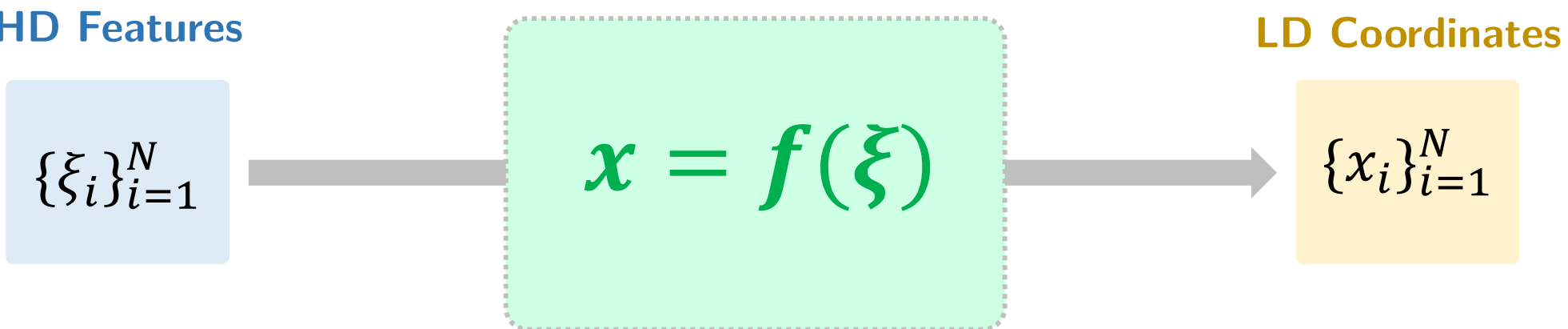
HD Features

$$\{\xi_i\}_{i=1}^N$$

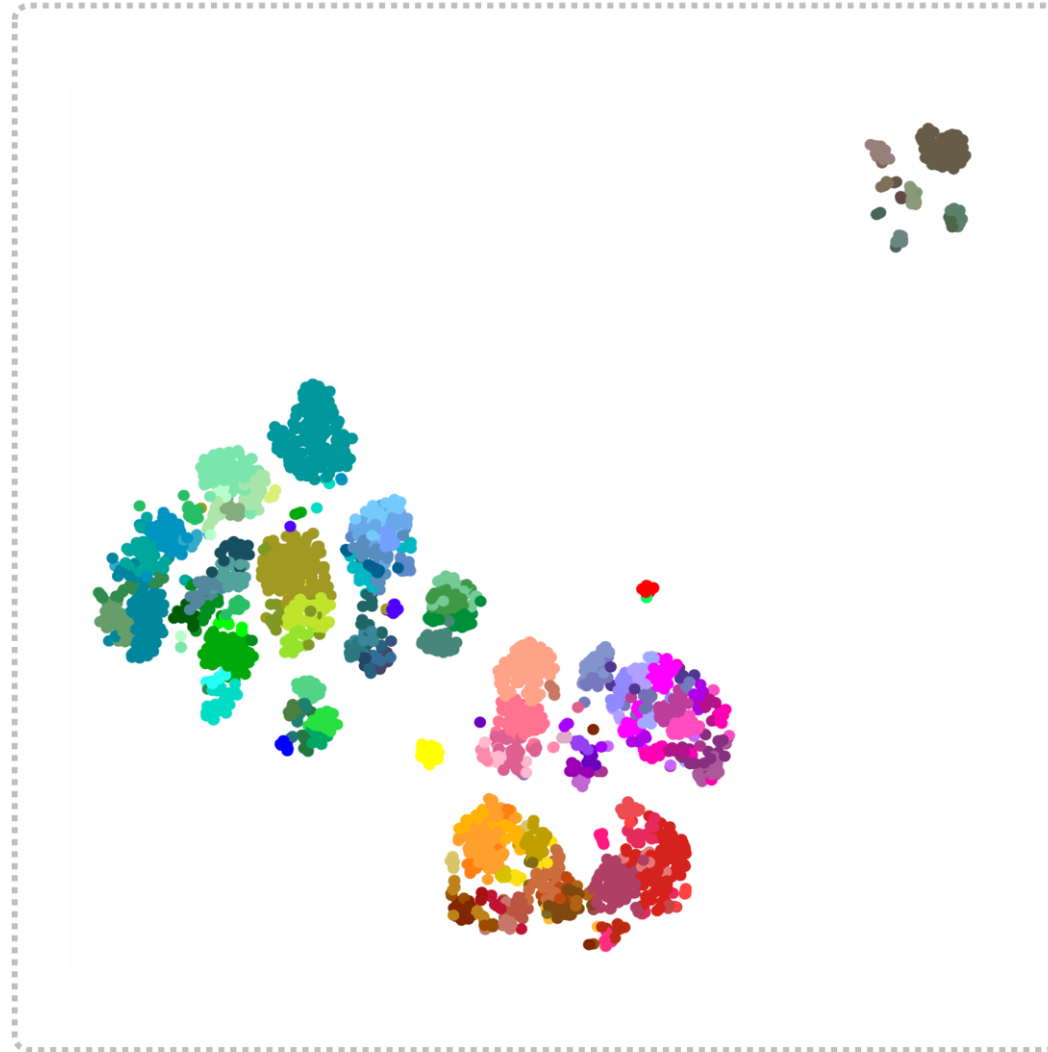
$$x = f(\xi)$$

LD Coordinates

$$\{x_i\}_{i=1}^N$$

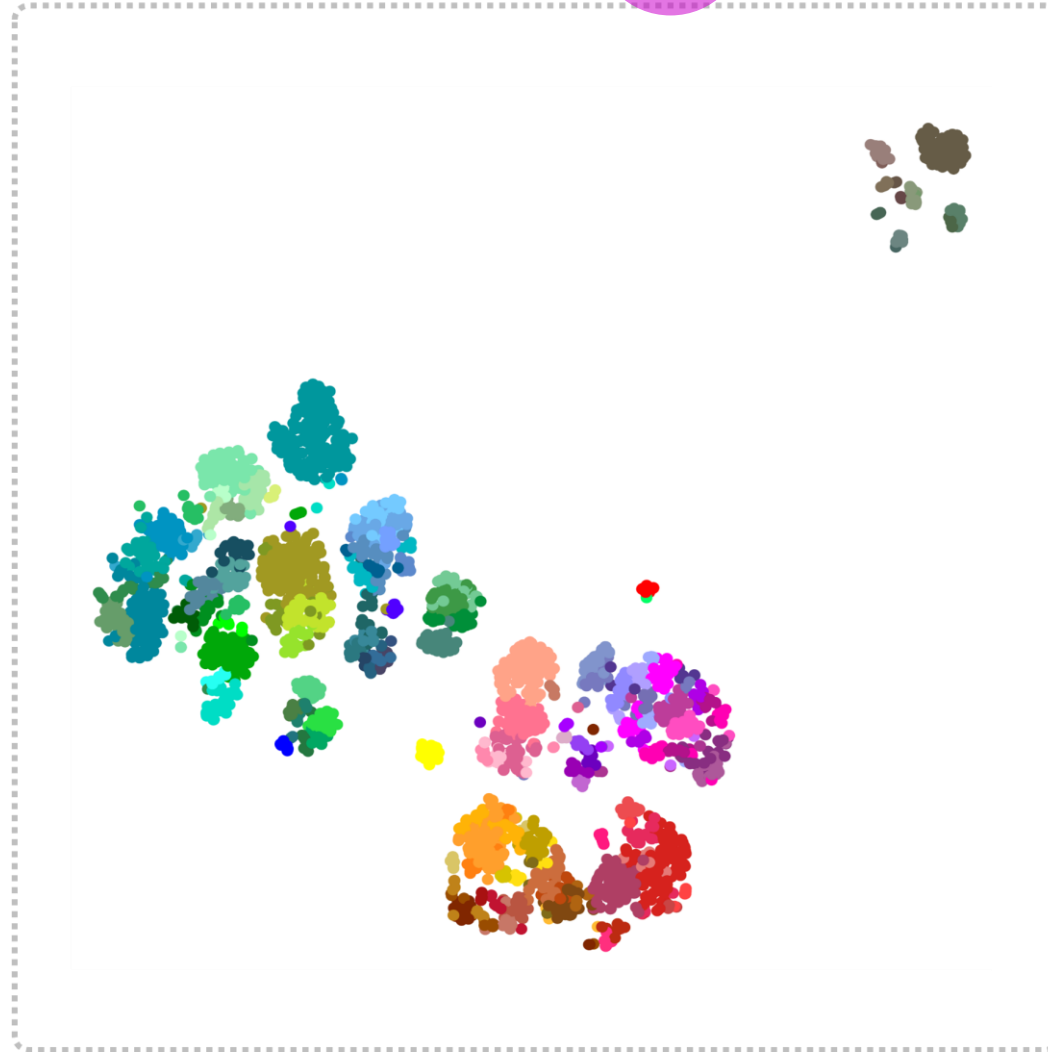


Ms. *t*-SNE [-]

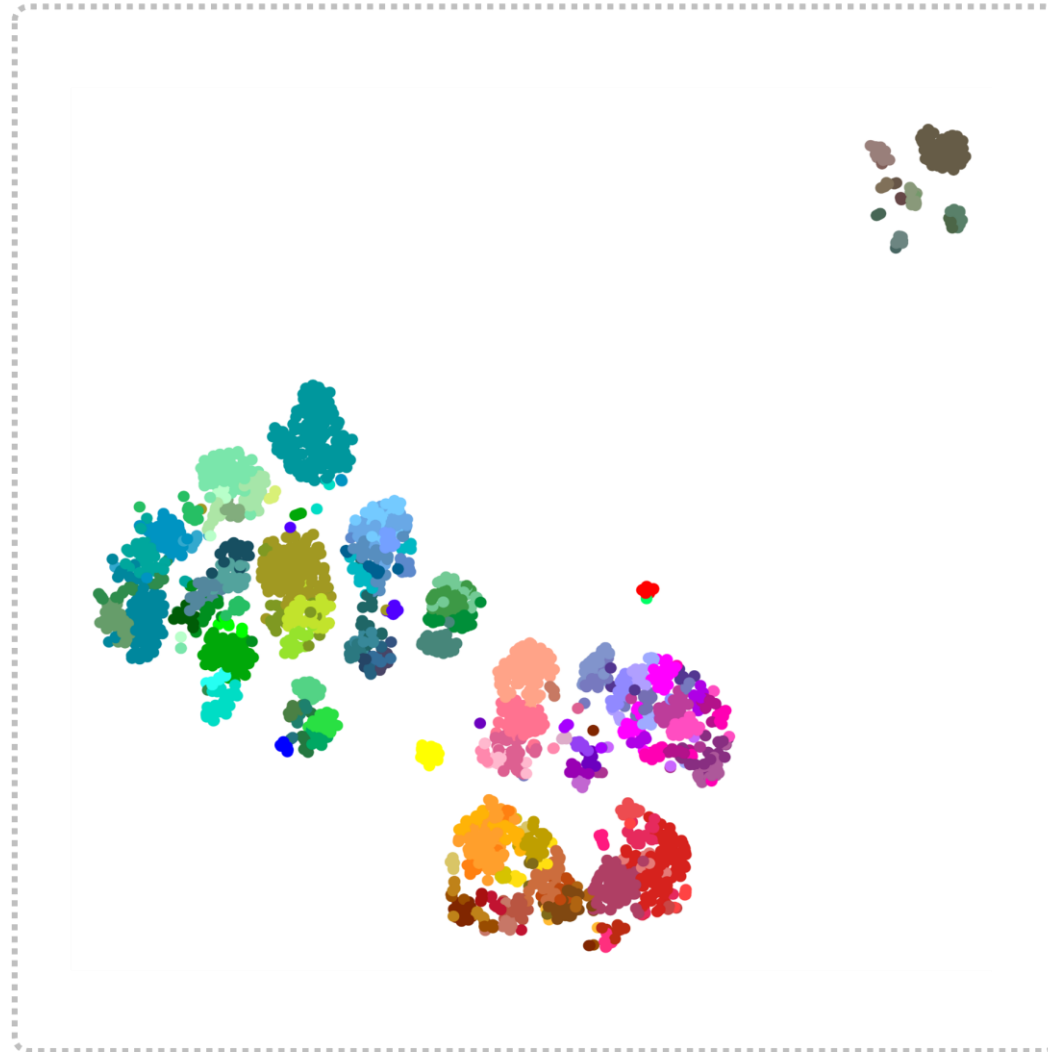


Ms. *t*-SNE [-]

Nonparametric

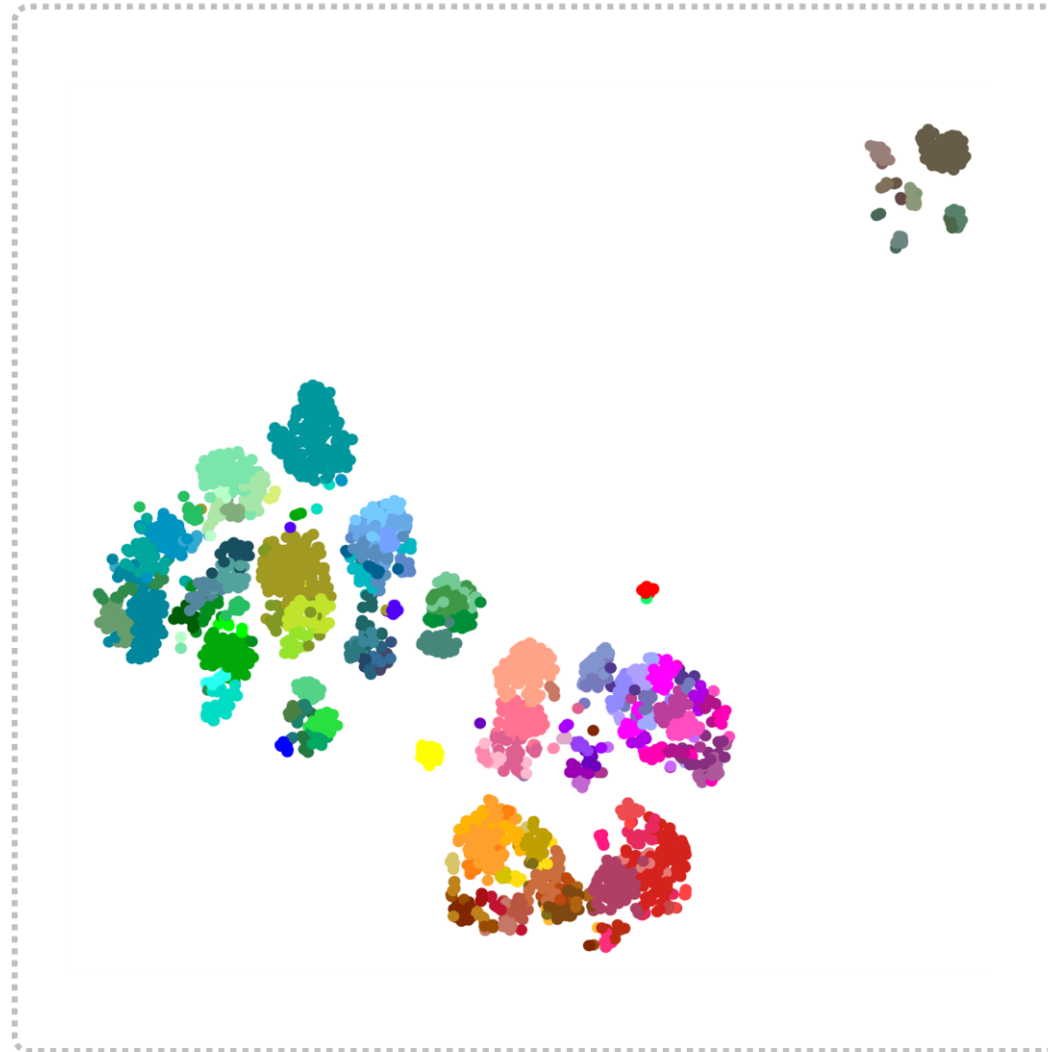


Ms. *t*-SNE [-]



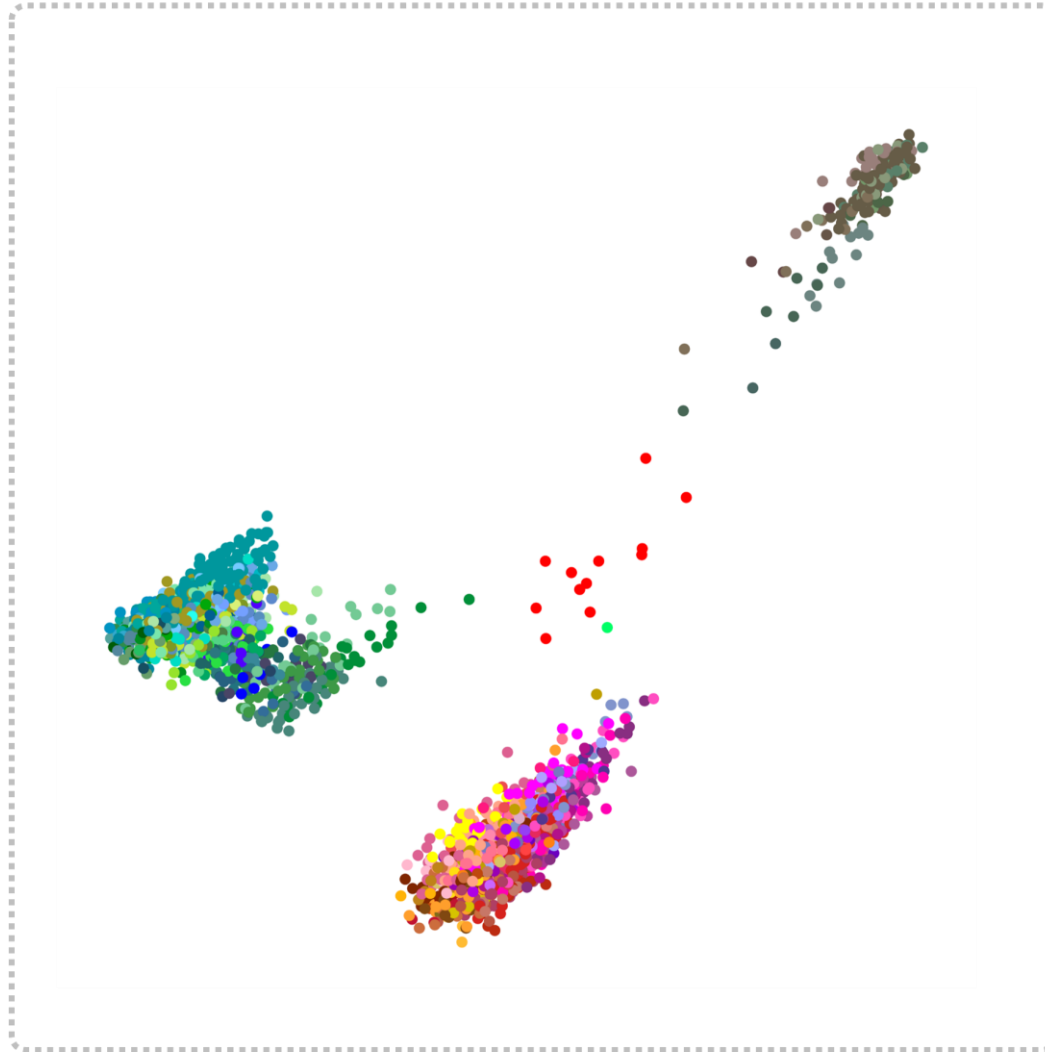
Interpretability  

Ms. *t*-SNE [-]



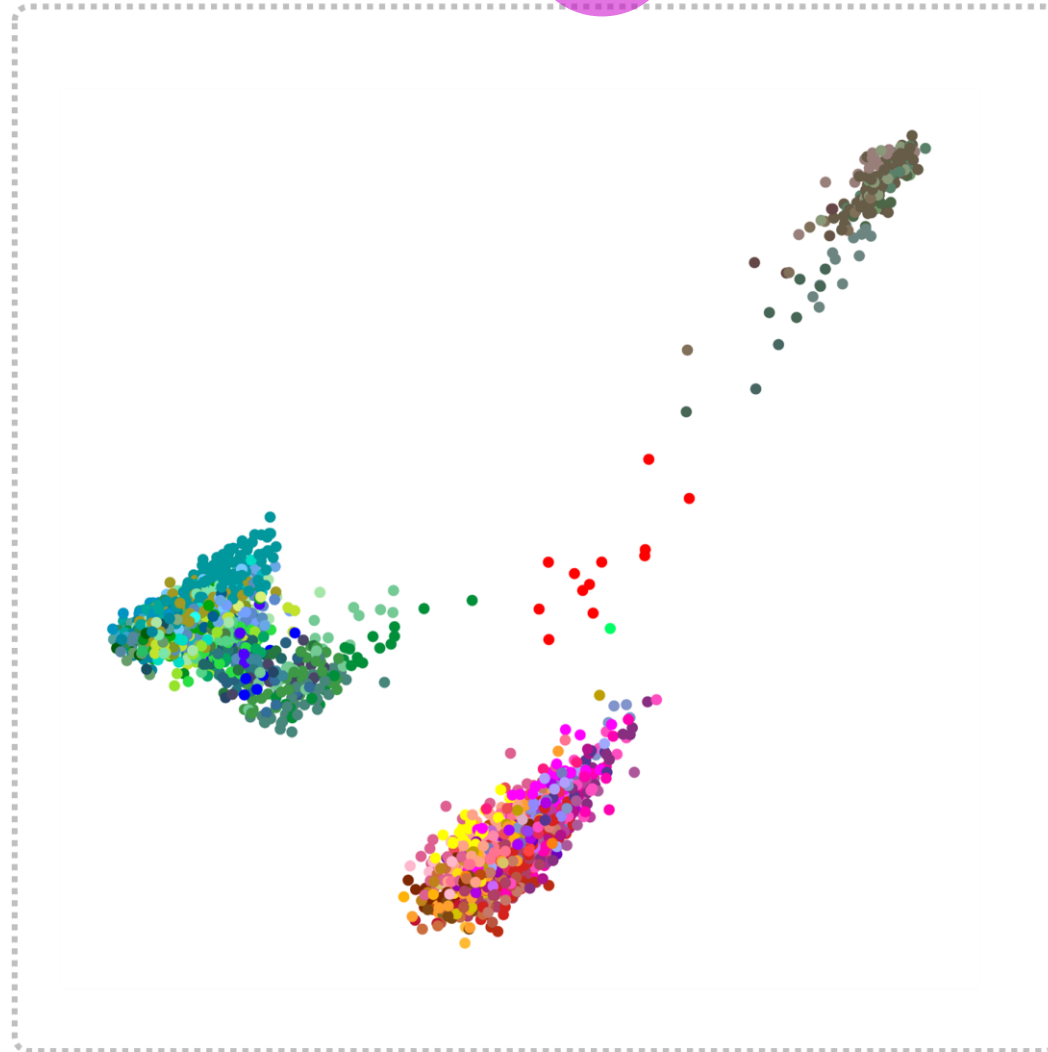
Interpretability **- -** Quality **+ +**

PCA [ W ]

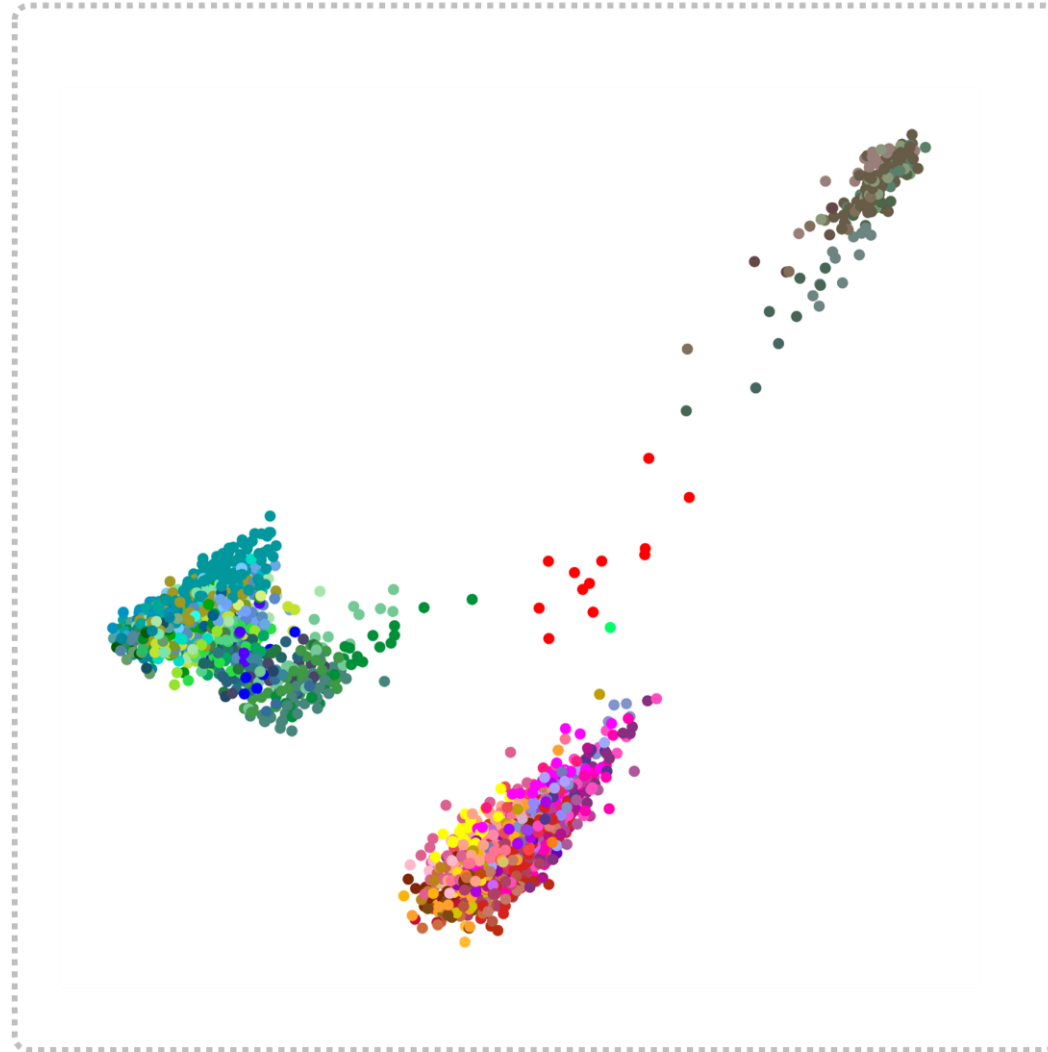


PCA [W]

← Parametric



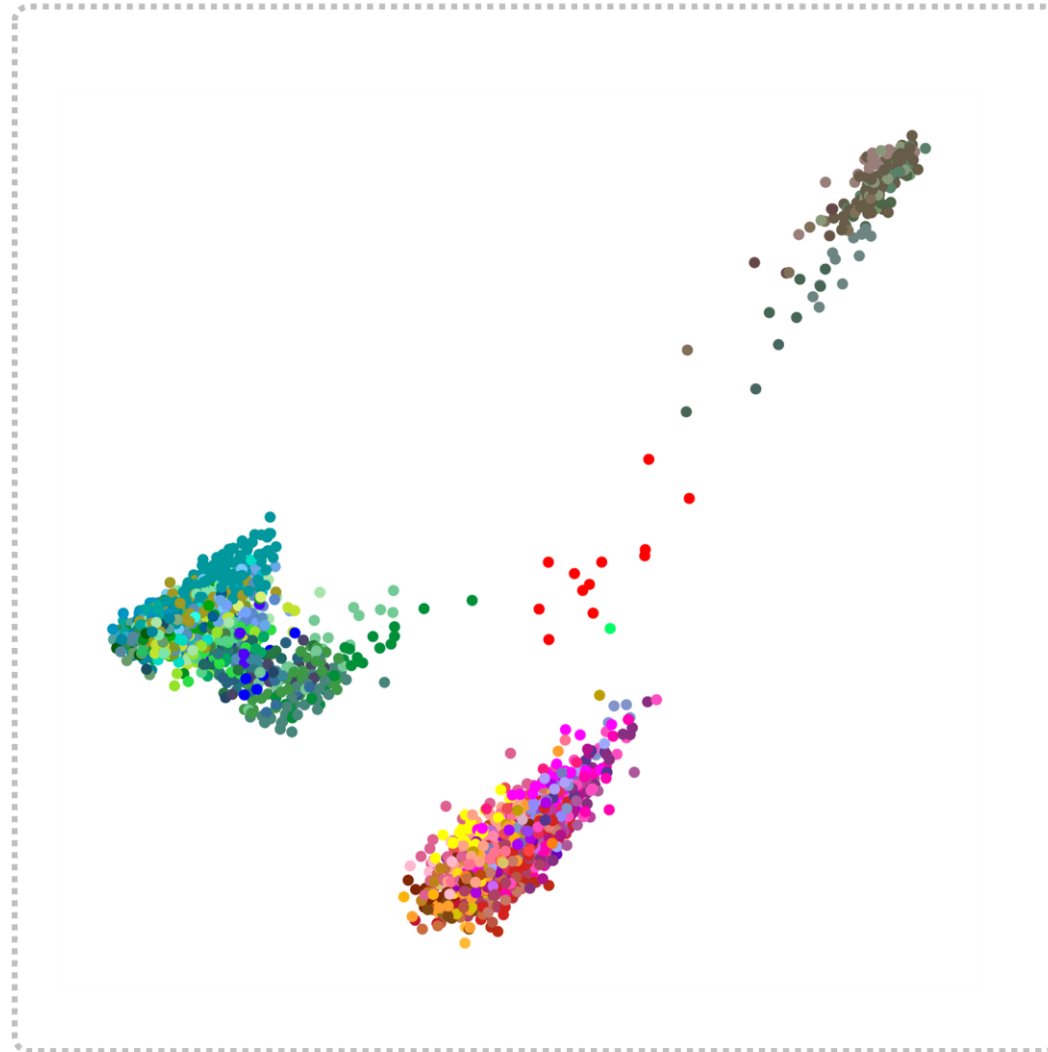
PCA [ W ]



Interpretability **+** **+**

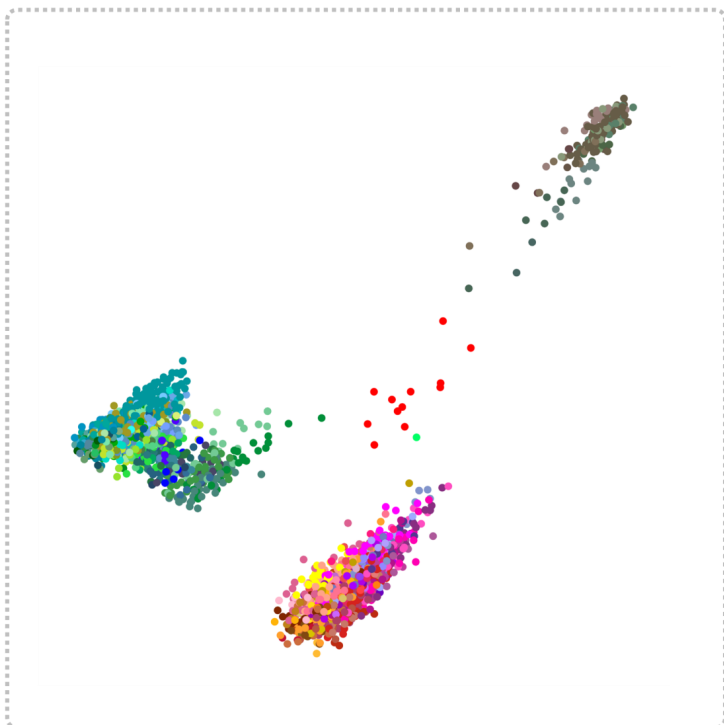


PCA [ W ]



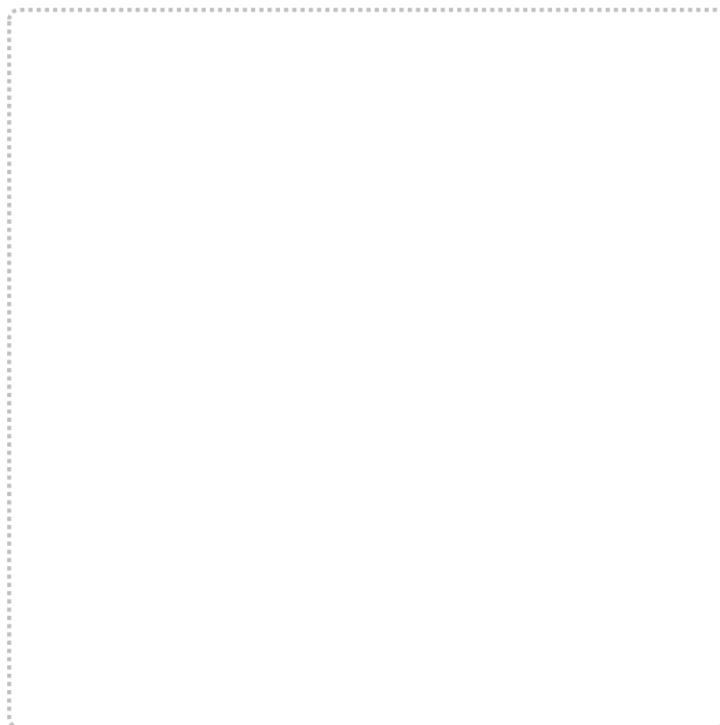
Interpretability **+** **+**      Quality **-** **-**

PCA [ W ]

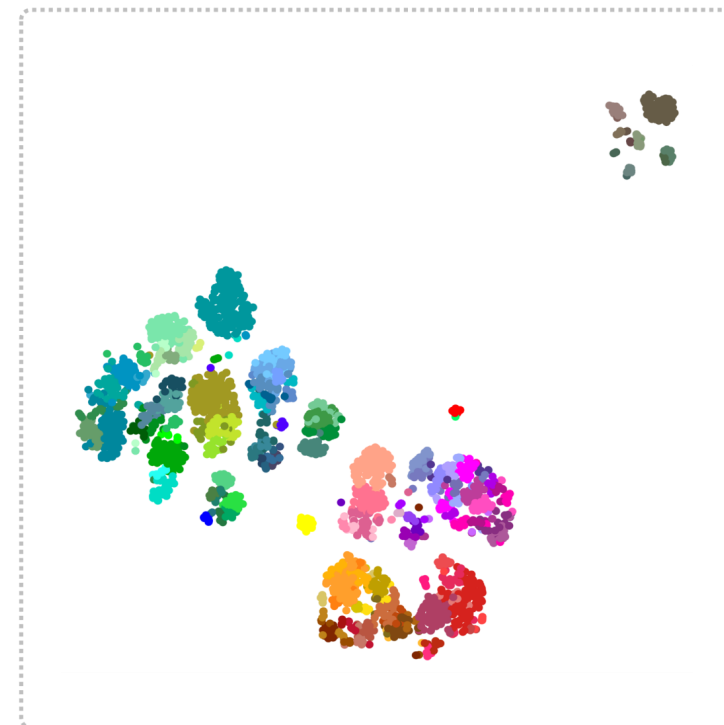


Interpretability + + Quality - -

?

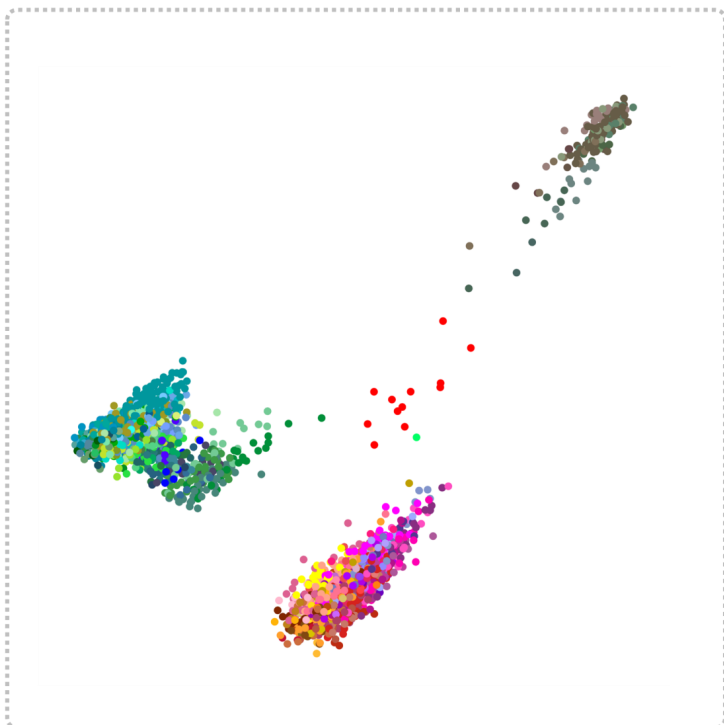


Ms. *t*-SNE [ - ]



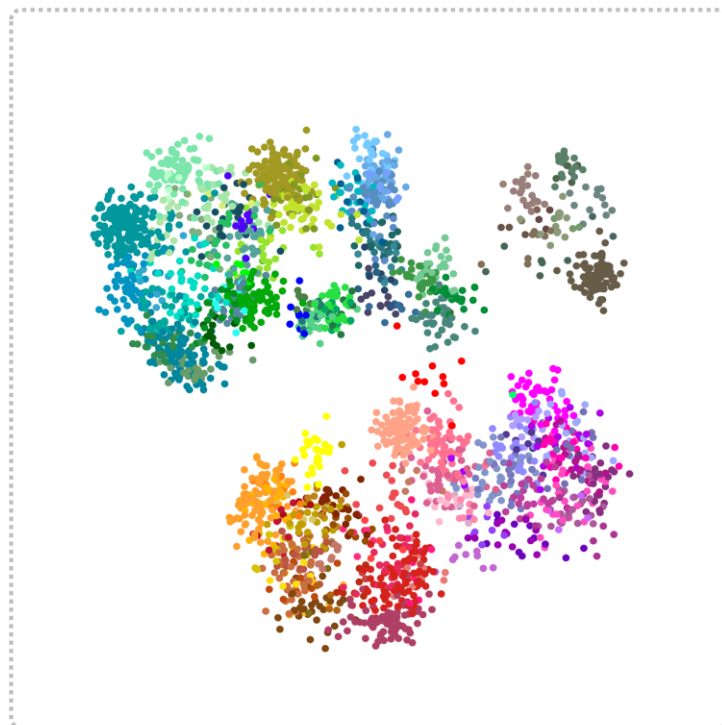
Interpretability - - Quality + +

PCA [ W ]

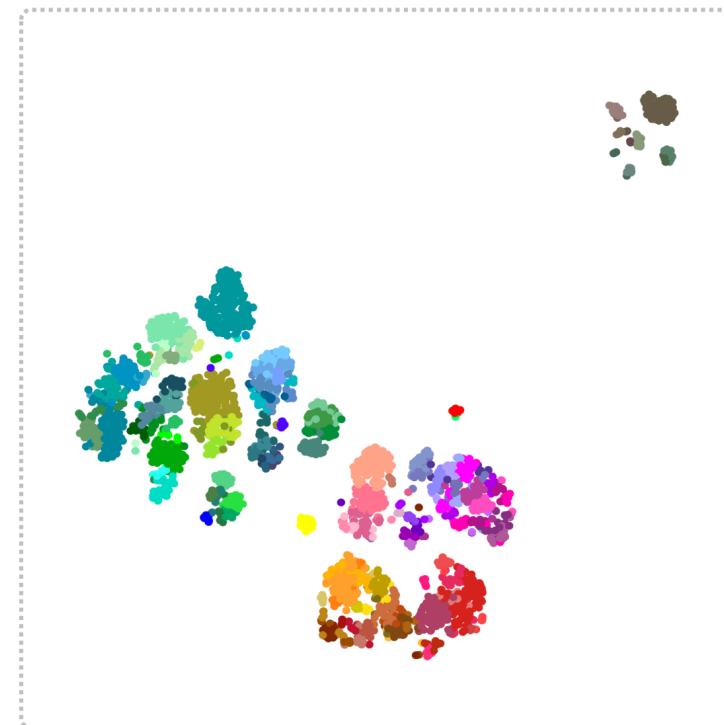


Interpretability + + Quality - -

Ms. *t*-SNE [ W ]

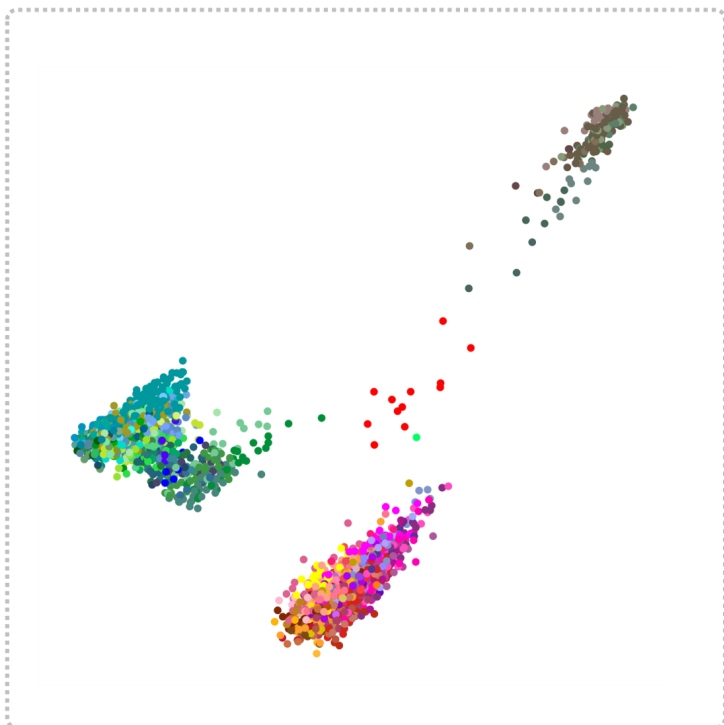


Ms. *t*-SNE [ - ]



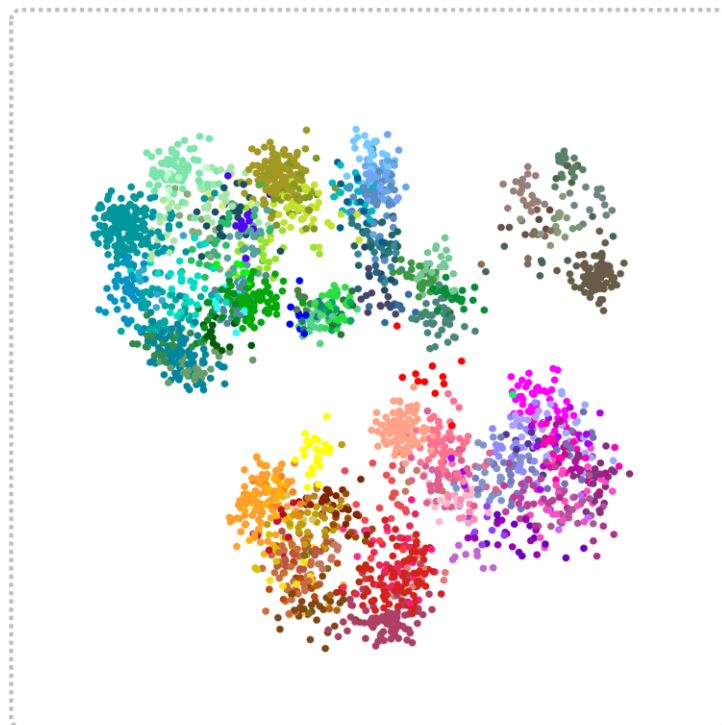
Interpretability - - Quality + +

PCA [ W ]



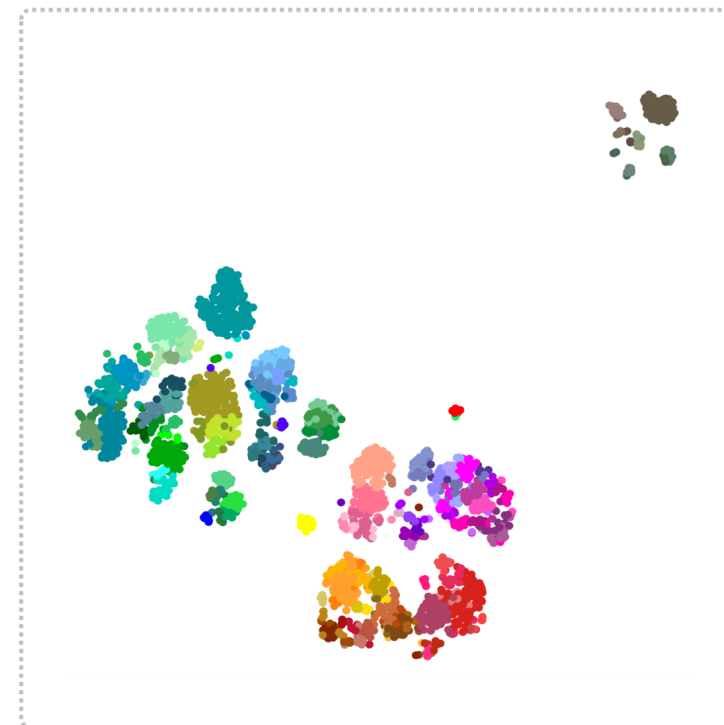
Interpretability + + Quality - -

Ms. *t*-SNE [ W ]



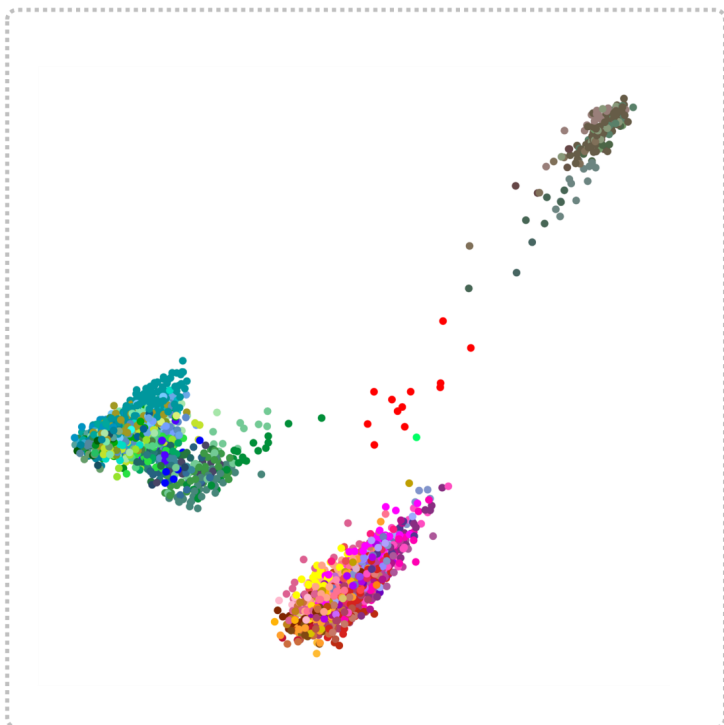
Interpretability + +

Ms. *t*-SNE [ - ]



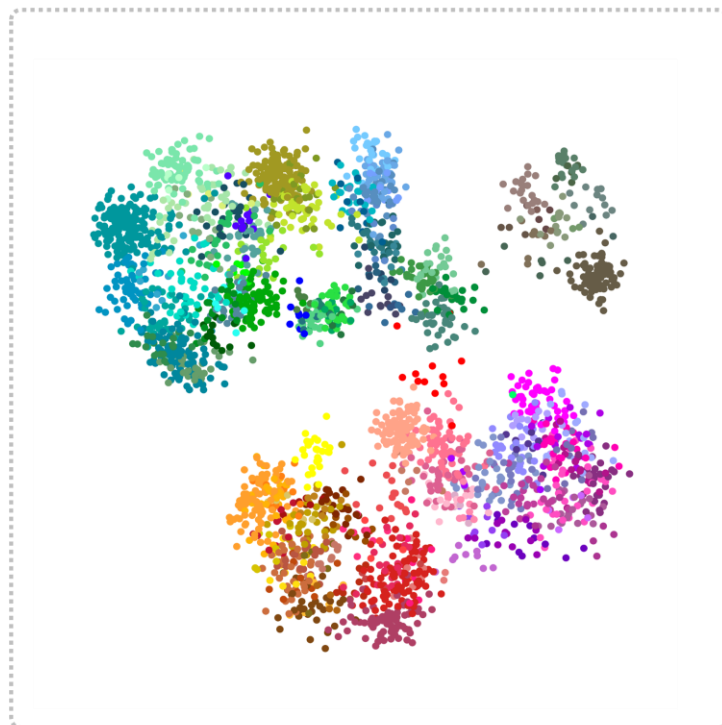
Interpretability - - Quality + +

PCA [ W ]



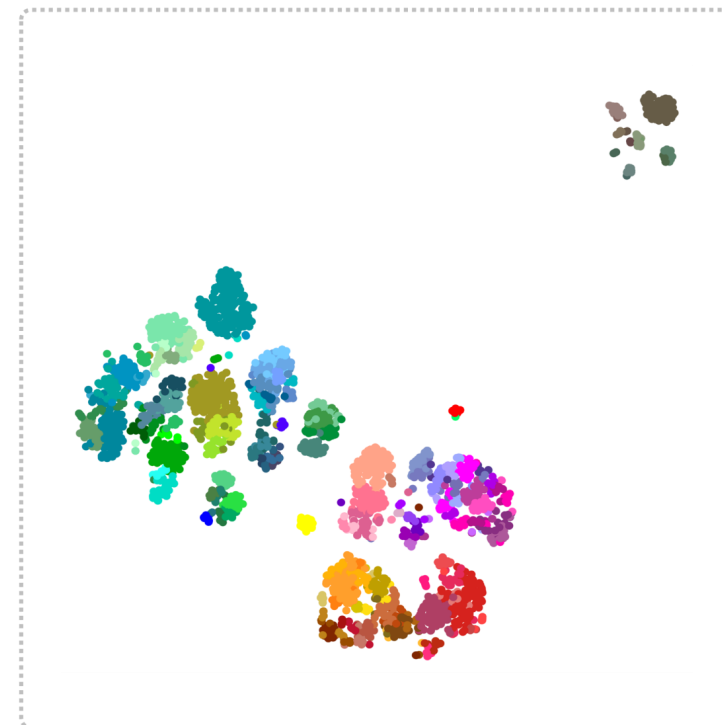
Interpretability **+** **+** Quality **-** **-**

Ms. *t*-SNE [ W ]



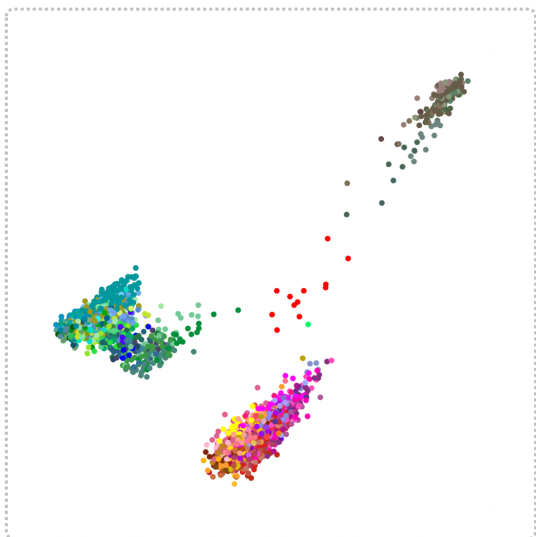
Interpretability **+** **+** Quality **-**

Ms. *t*-SNE [ - ]



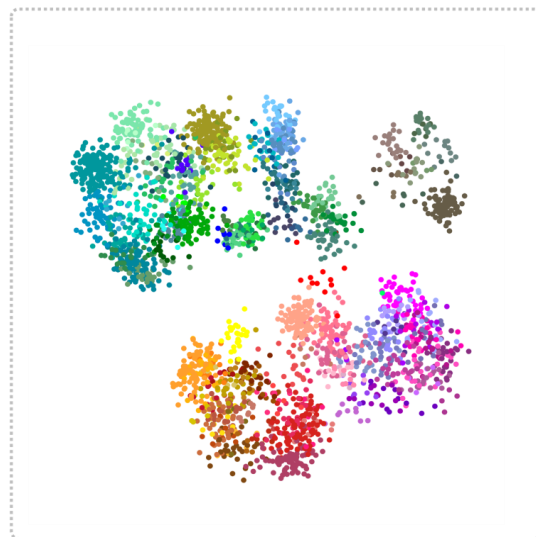
Interpretability **-** **-** Quality **+** **+**

PCA [W]



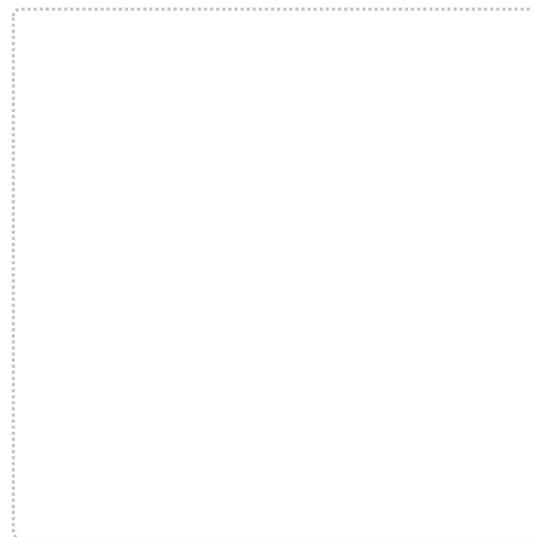
Interpretability **+** **+** Quality **-** **-**

Ms. *t*-SNE [W]

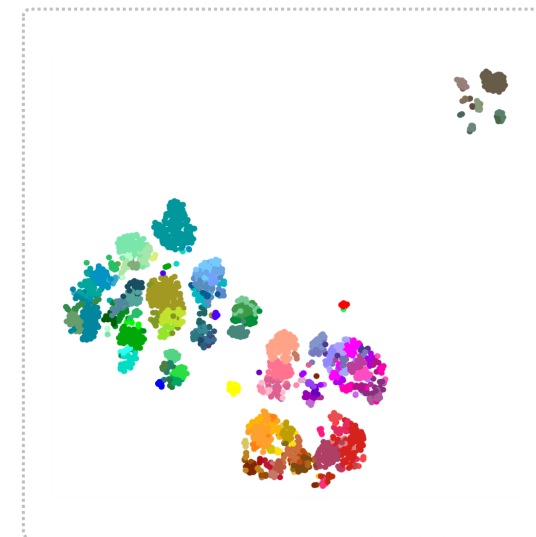


Interpretability **+** **+** Quality **-**

?

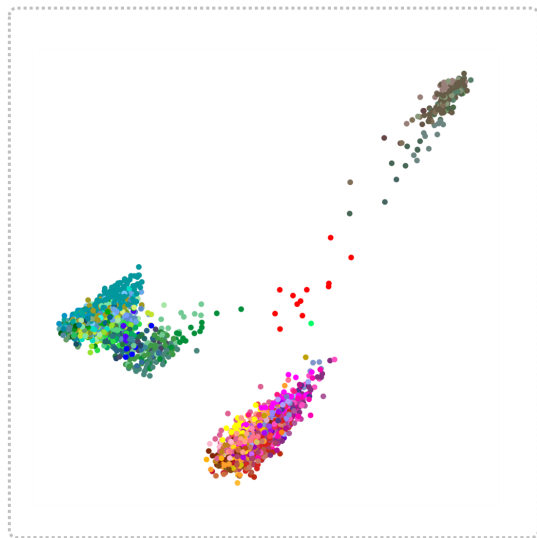


Ms. *t*-SNE [-]



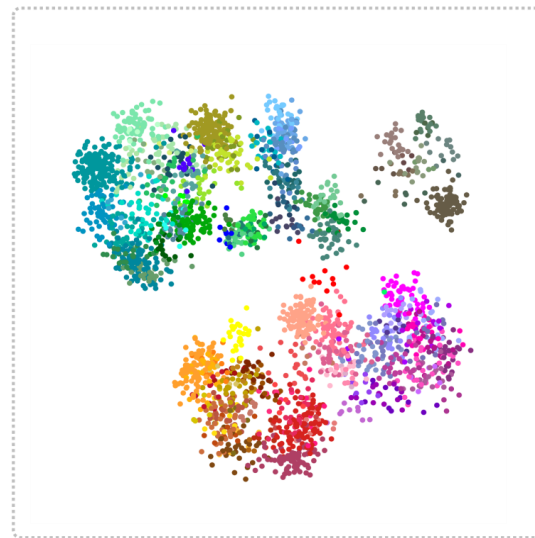
Interpretability **-** **-** Quality **+** **+**

PCA [W]



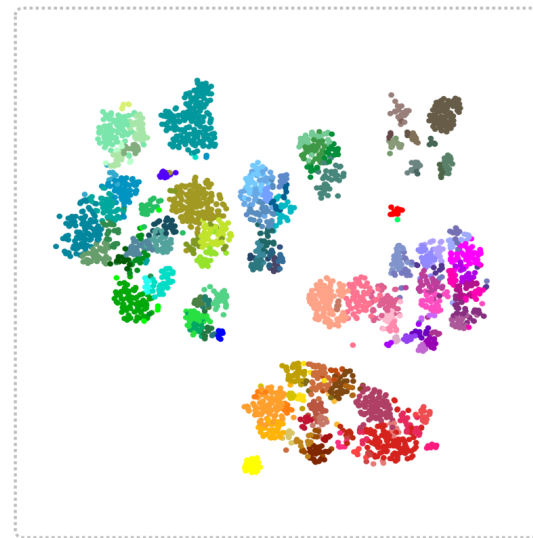
Interpretability **+** **+** Quality **-** **-**

Ms. *t*-SNE [W]

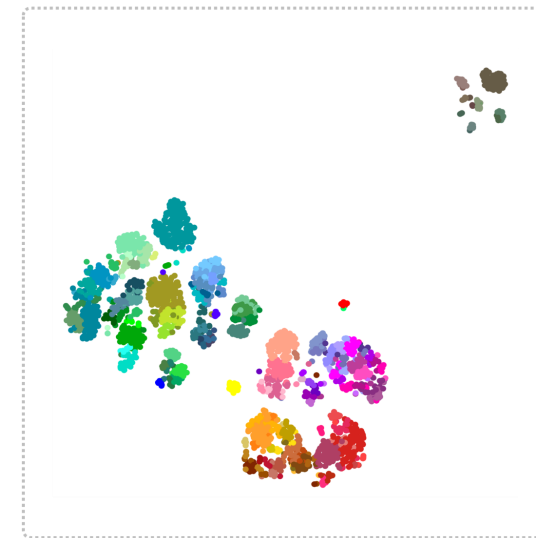


Interpretability **+** **+** Quality **-**

Ms. *t*-SNE [ $W_i$ ]

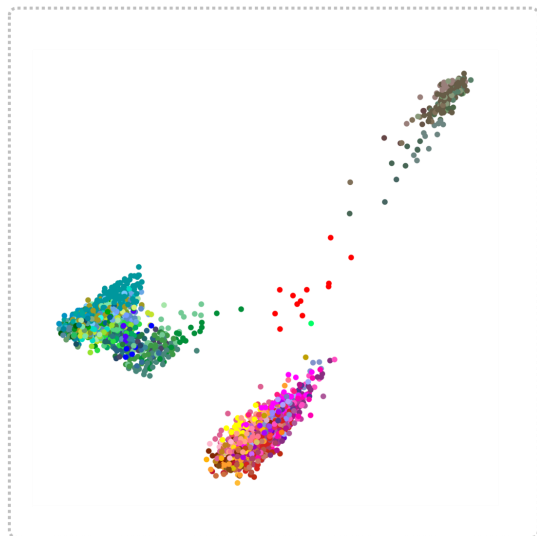


Ms. *t*-SNE [-]



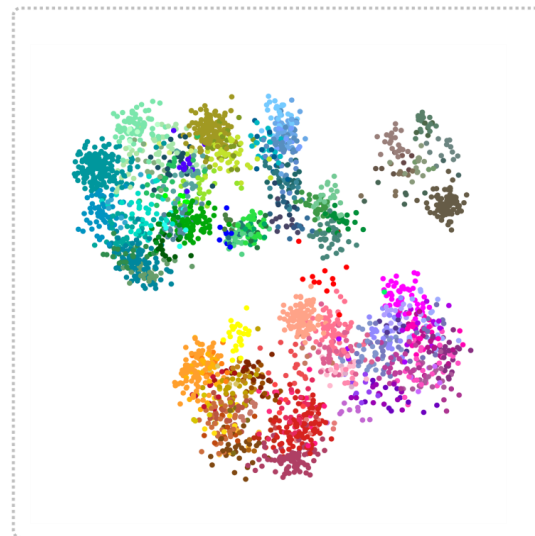
Interpretability **-** **-** Quality **+** **+**

PCA [W]



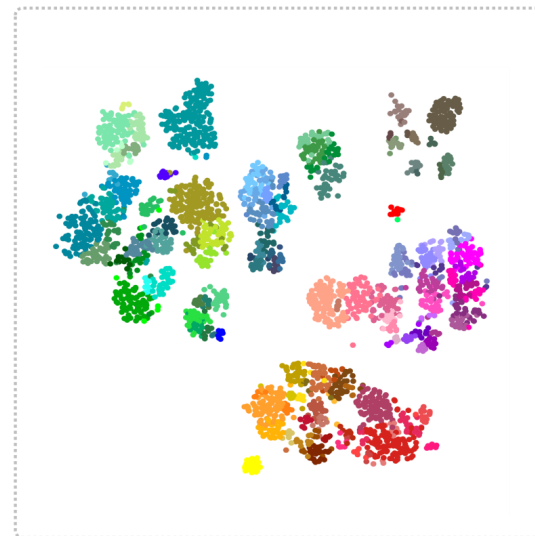
Interpretability **+** **+** Quality **-** **-**

Ms. *t*-SNE [W]



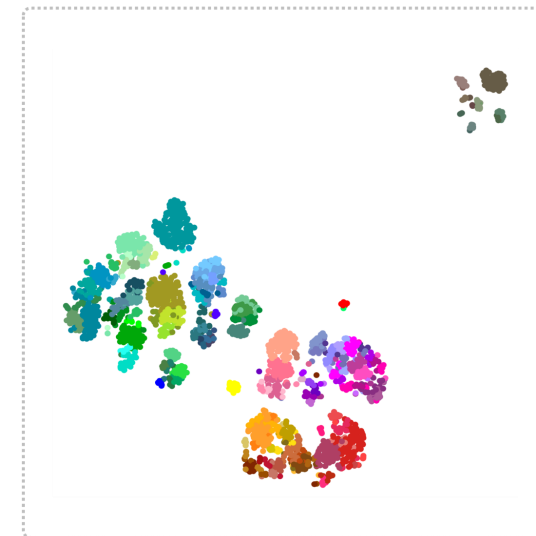
Interpretability **+** **+** Quality **-**

Ms. *t*-SNE [ $W_i$ ]



Interpretability **+**

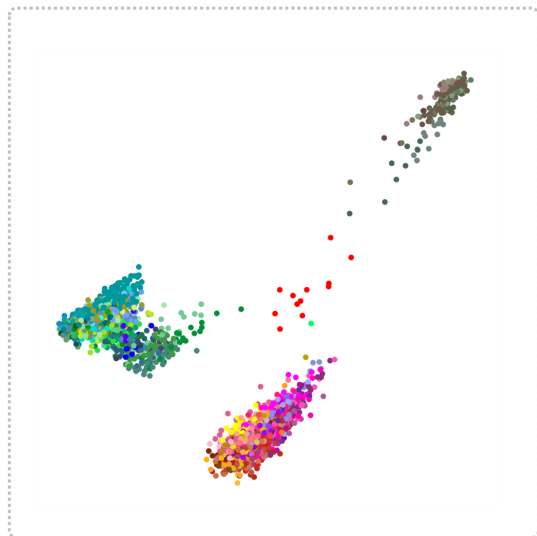
Ms. *t*-SNE [-]



Interpretability **-** **-** Quality **+** **+**

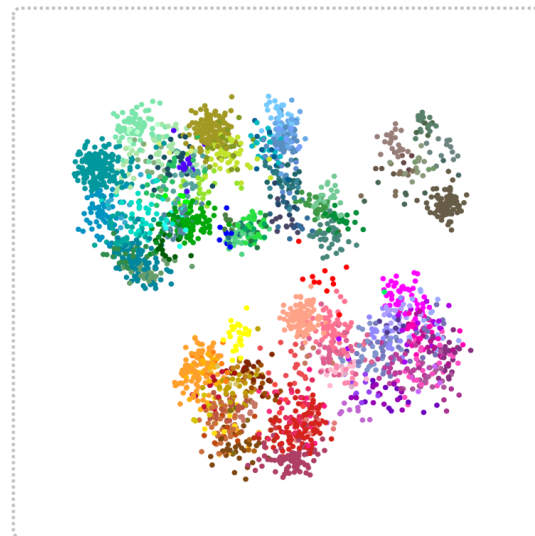


PCA [W]



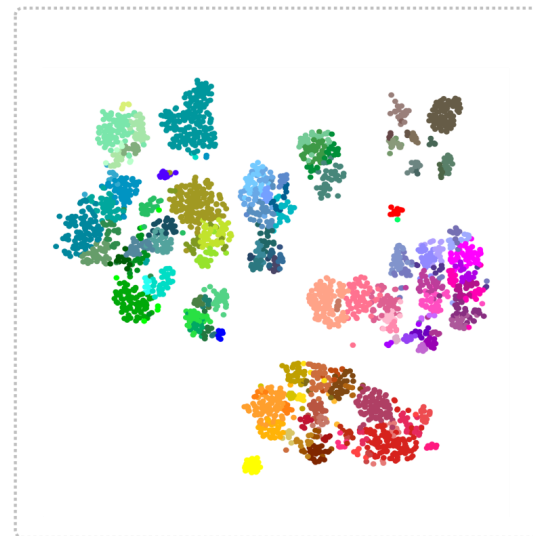
Interpretability **+** **+** Quality **-** **-**

Ms. *t*-SNE [W]



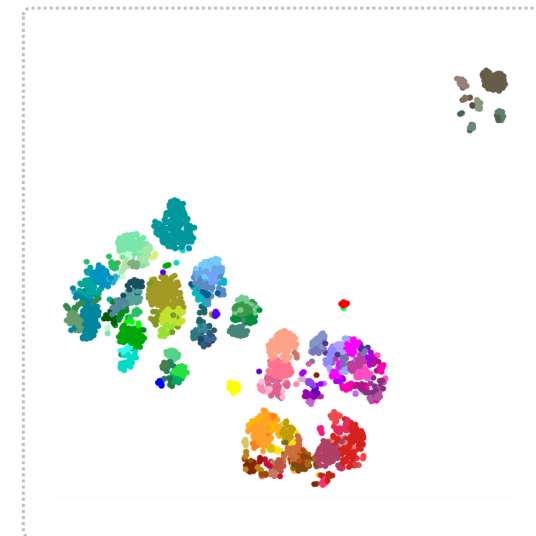
Interpretability **+** **+** Quality **-**

Ms. *t*-SNE [ $W_i$ ]



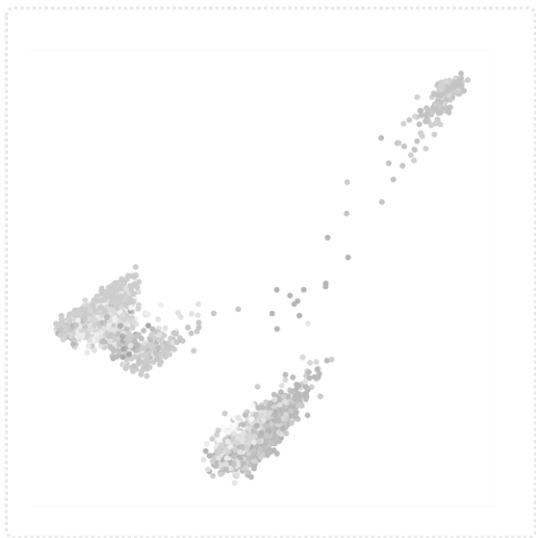
Interpretability **+** Quality **+**

Ms. *t*-SNE [-]



Interpretability **-** **-** Quality **+** **+**

PCA [W]



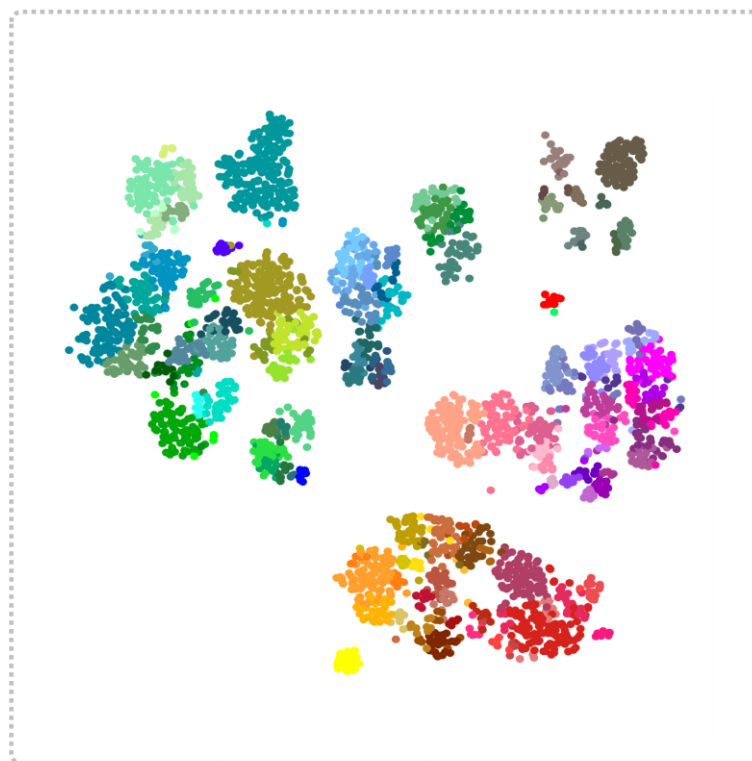
Interpretability ++ Quality --

Ms. *t*-SNE [W]



Interpretability ++ Qu

Ms. *t*-SNE [ $W_i$ ]



Interpretability + Quality +

*t*-SNE [-]



Interpretability - Quality ++

## M-dim space

Data

$$\{\xi_i\}_{i=1}^N$$

Similarities

$$\tau_{ij}$$

## 2-dim space

Initialization

$$\{x_i\}_{i=1}^N$$

Similarities

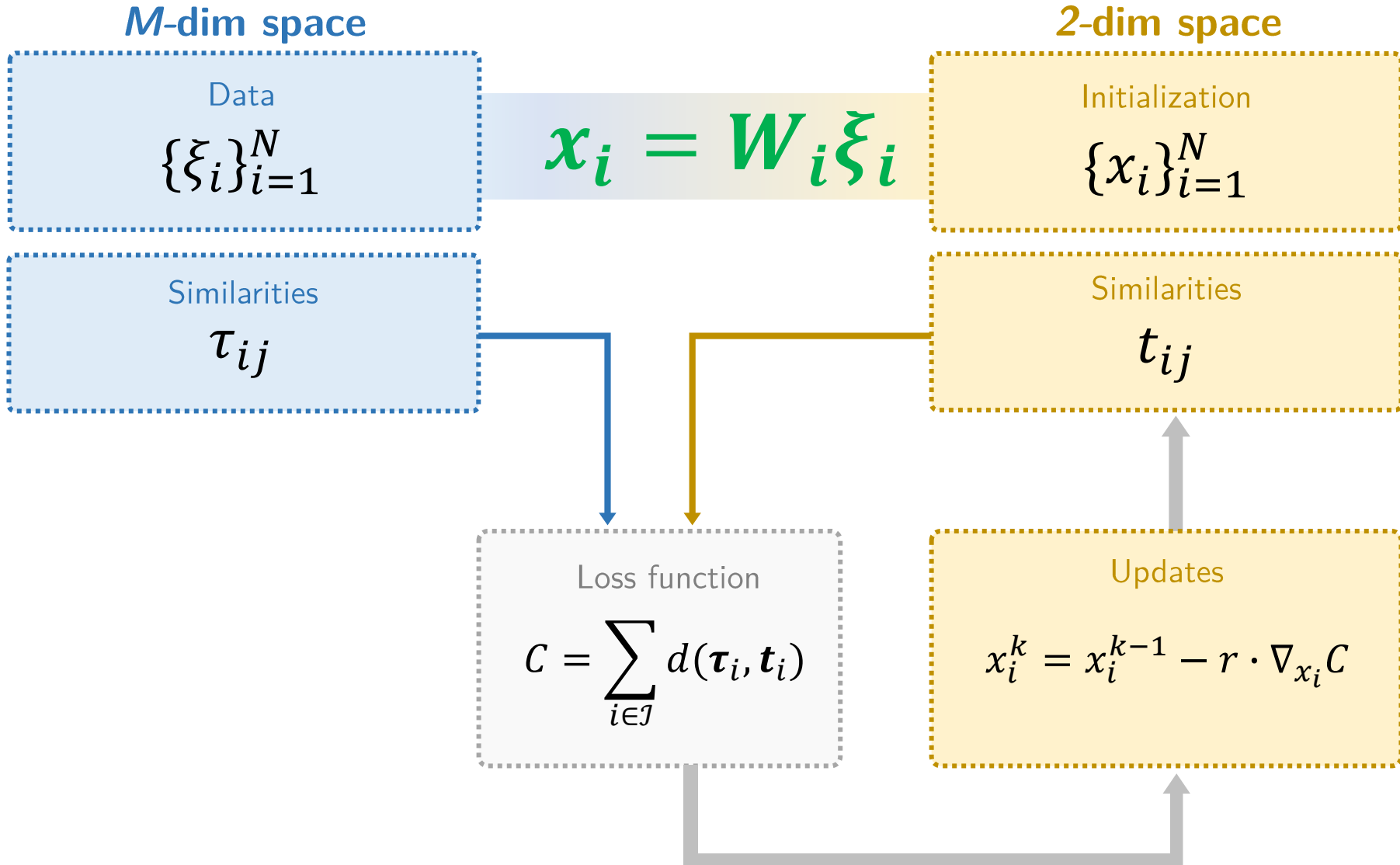
$$t_{ij}$$

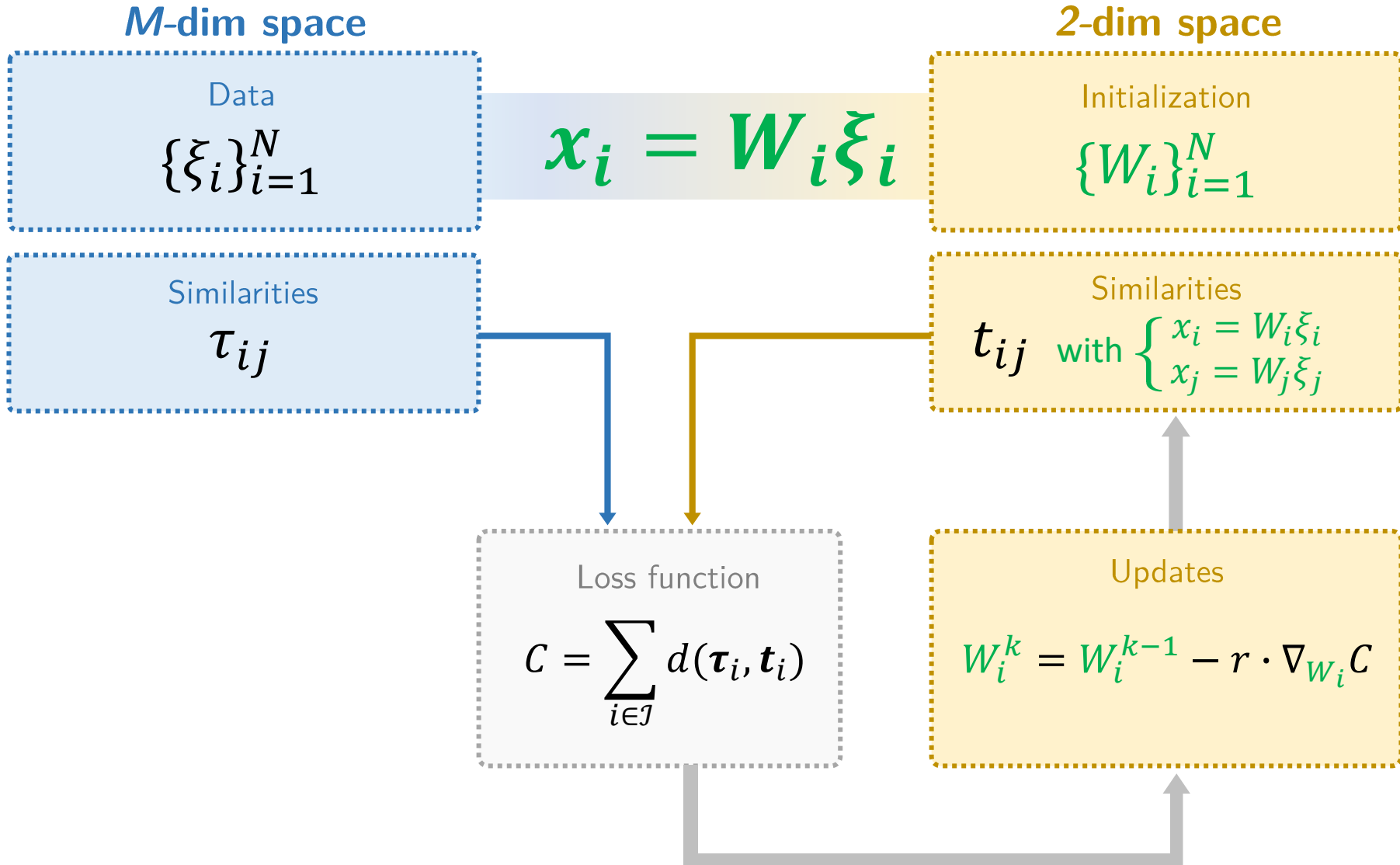
Updates

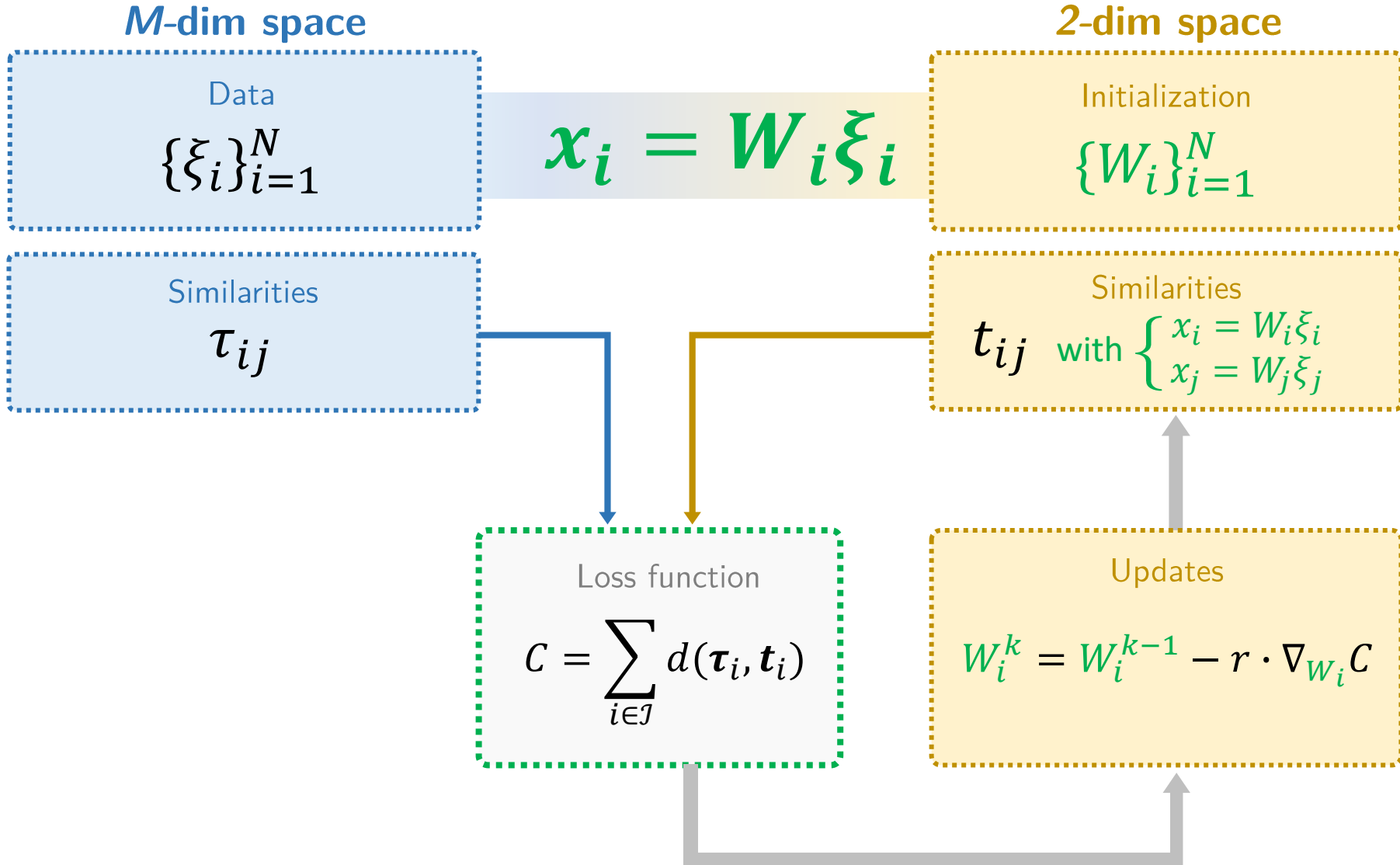
$$x_i^k = x_i^{k-1} - r \cdot \nabla_{x_i} C$$

Loss function

$$C = \sum_{i \in \mathcal{I}} d(\tau_i, t_i)$$







$$C = \sum_{i \in \mathcal{I}} \left( d(\boldsymbol{\tau}_i, \mathbf{t}_i) \right)$$

$$C = \sum_{i \in \mathcal{I}} \left( d(\boldsymbol{\tau}_i, \mathbf{t}_i) + \alpha \cdot \sum_{j \in \mathcal{J}} t_{ij} \|W_i - W_j\|_F \right)$$

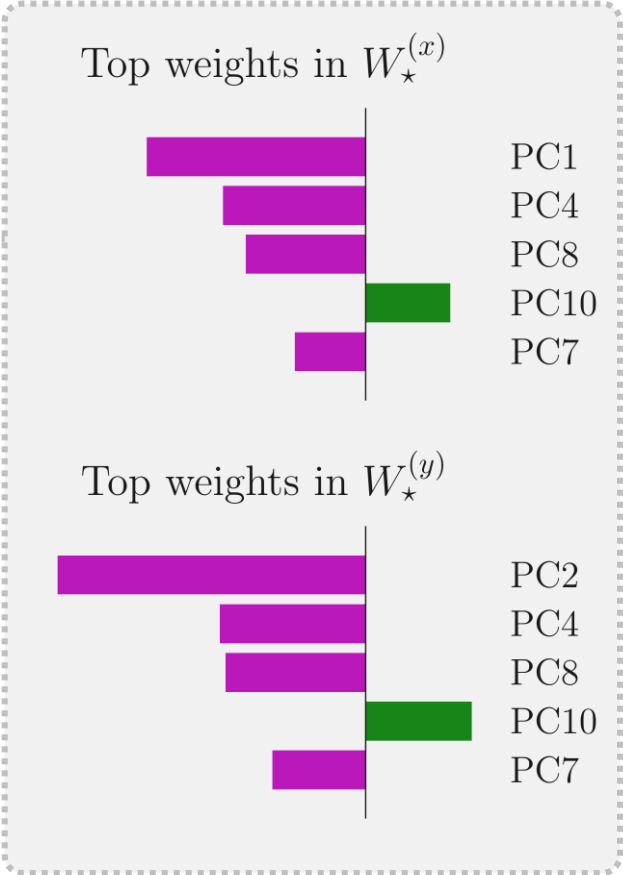
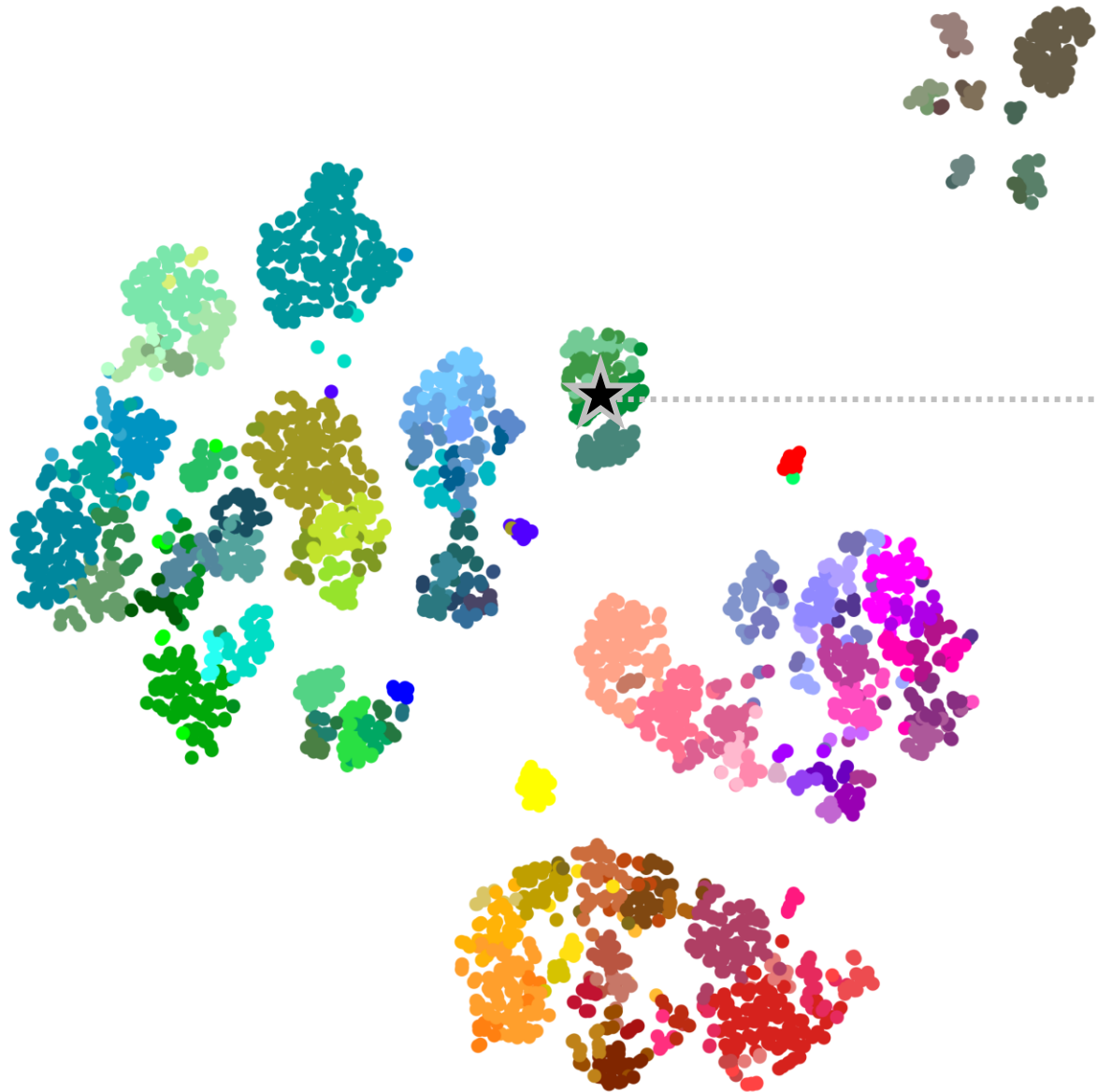
Coherence

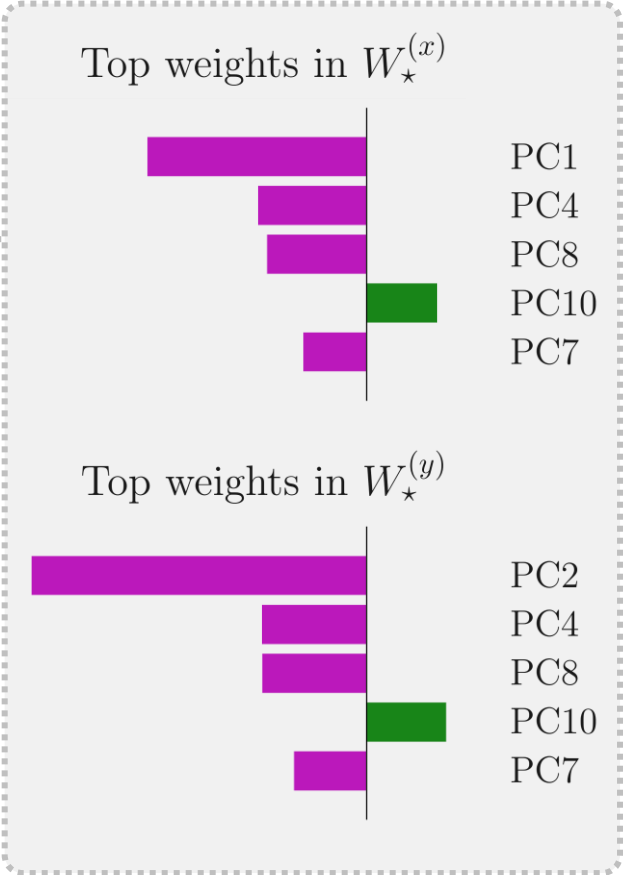
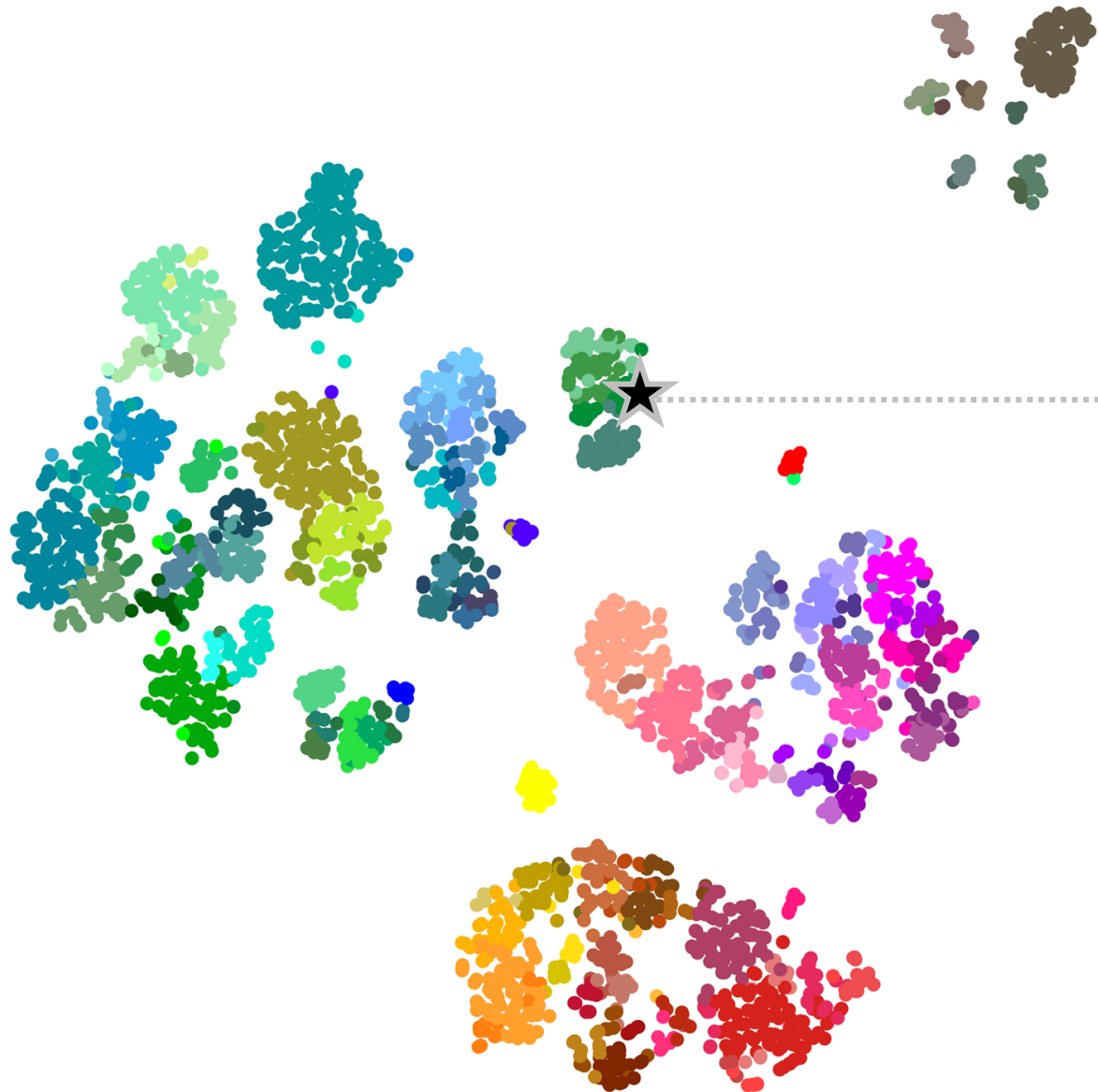


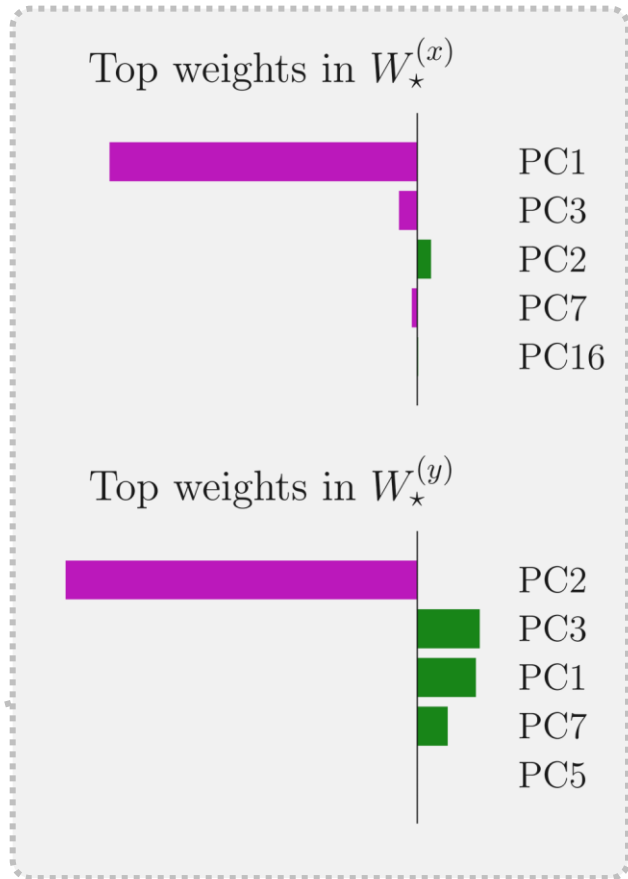
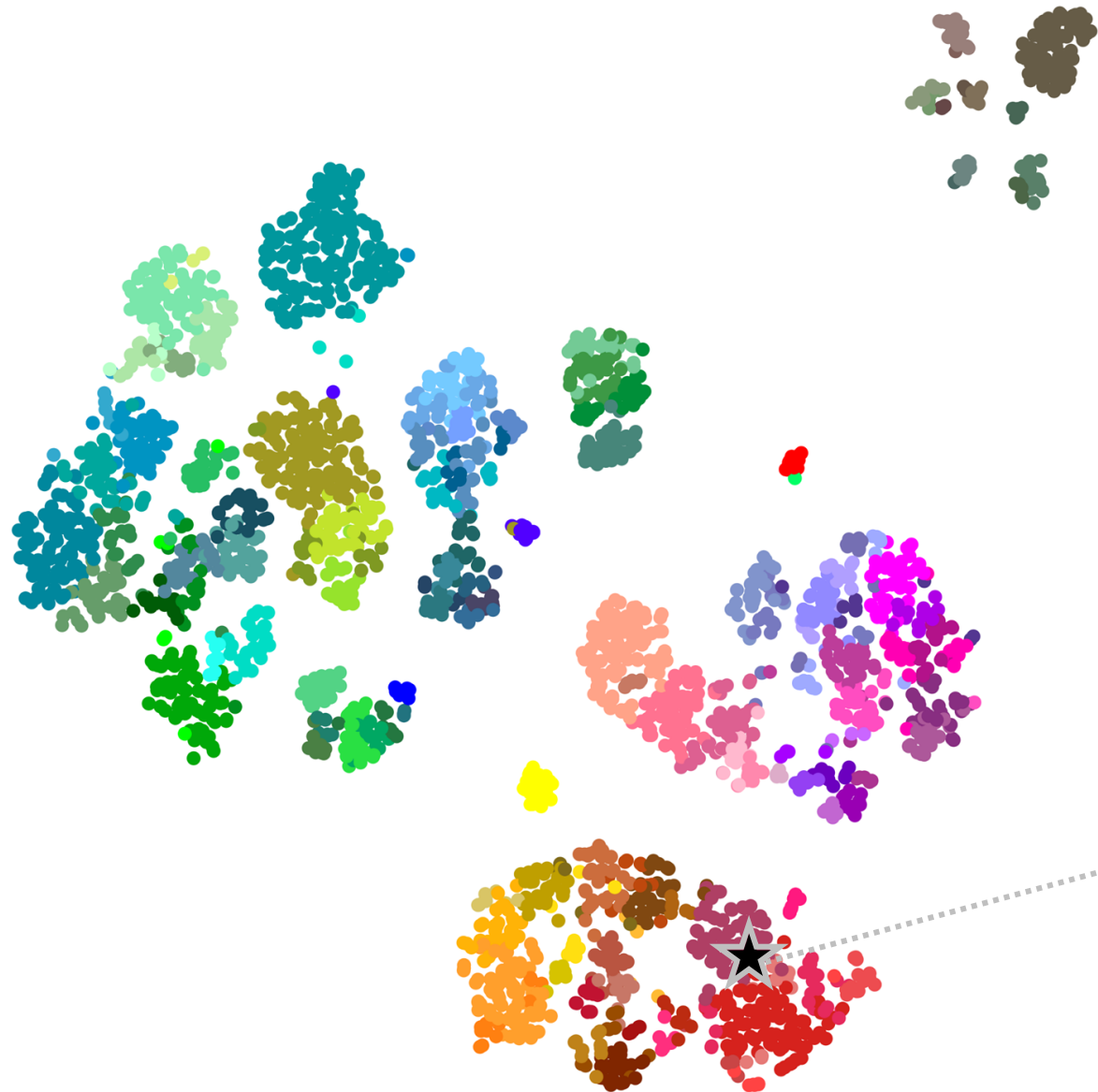
$$C = \sum_{i \in \mathcal{I}} \left( d(\boldsymbol{\tau}_i, \mathbf{t}_i) + \alpha \cdot \sum_{j \in \mathcal{J}} t_{ij} \|W_i - W_j\|_F + \beta \cdot \|W_i\|_1 \right)$$

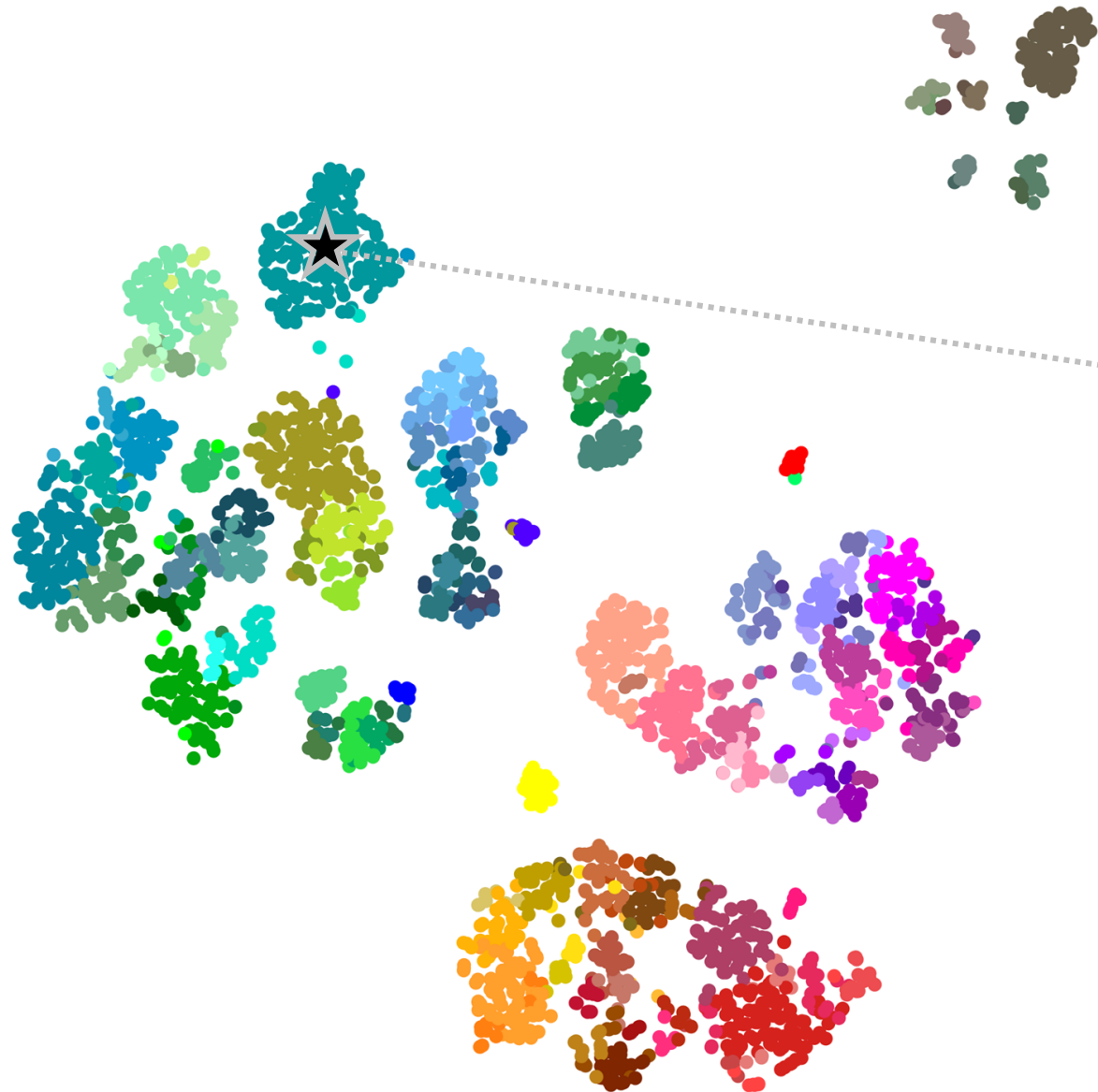
Coherence

Sparsity









Top weights in  $W_{\star}^{(x)}$



PC1  
PC44  
PC42  
PC17  
PC29

Top weights in  $W_{\star}^{(y)}$



PC2  
PC1  
PC4  
PC8  
PC9

# Main takeaways

# Main takeaways

- A natively **interpretable** version of  $t$ -SNE

# Main takeaways

- A natively **interpretable** version of  $t$ -SNE
- Explicit locally linear mapping  $\mathbf{x}_i = \mathbf{W}_i \tilde{\xi}_i$



# Main takeaways

- A natively **interpretable** version of  $t$ -SNE
- Explicit locally linear mapping  $x_i = W_i \xi_i$
- Navigable trade-off: **interpretability/quality**

Don't hesitate to reach out if you'd like to learn more!



[edouard.couplet@uclouvain.be](mailto:edouard.couplet@uclouvain.be)



Edouard Couplet



Pierre Lambert



Michel Verleysen



Dounia Mulders

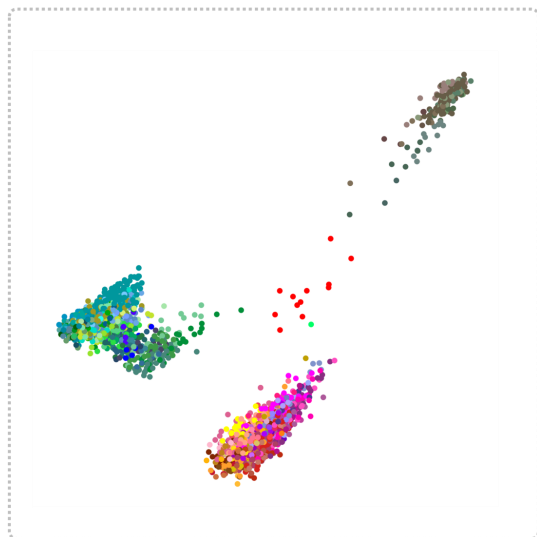


John Lee



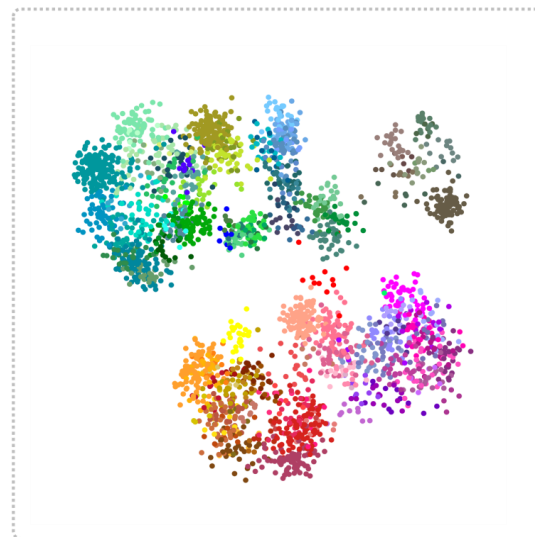
Cyril de Bodt

PCA [W]



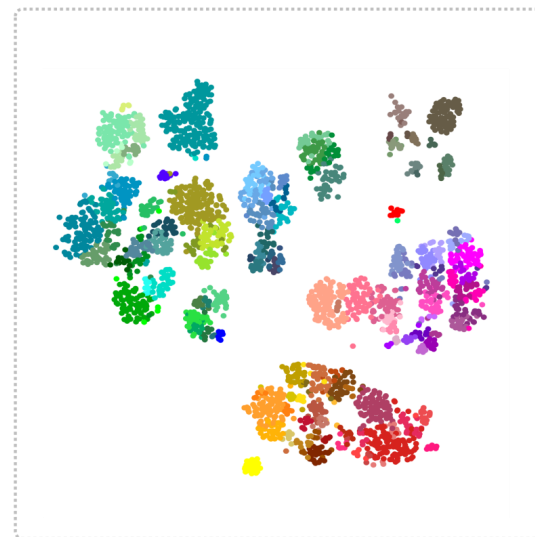
Interpretability **+** **+** Quality **-** **-**

Ms. *t*-SNE [W]



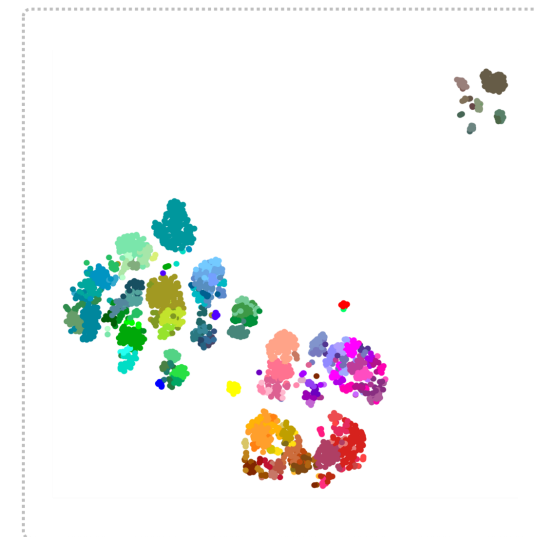
Interpretability **+** **+** Quality **-**

Ms. *t*-SNE [ $W_i$ ]

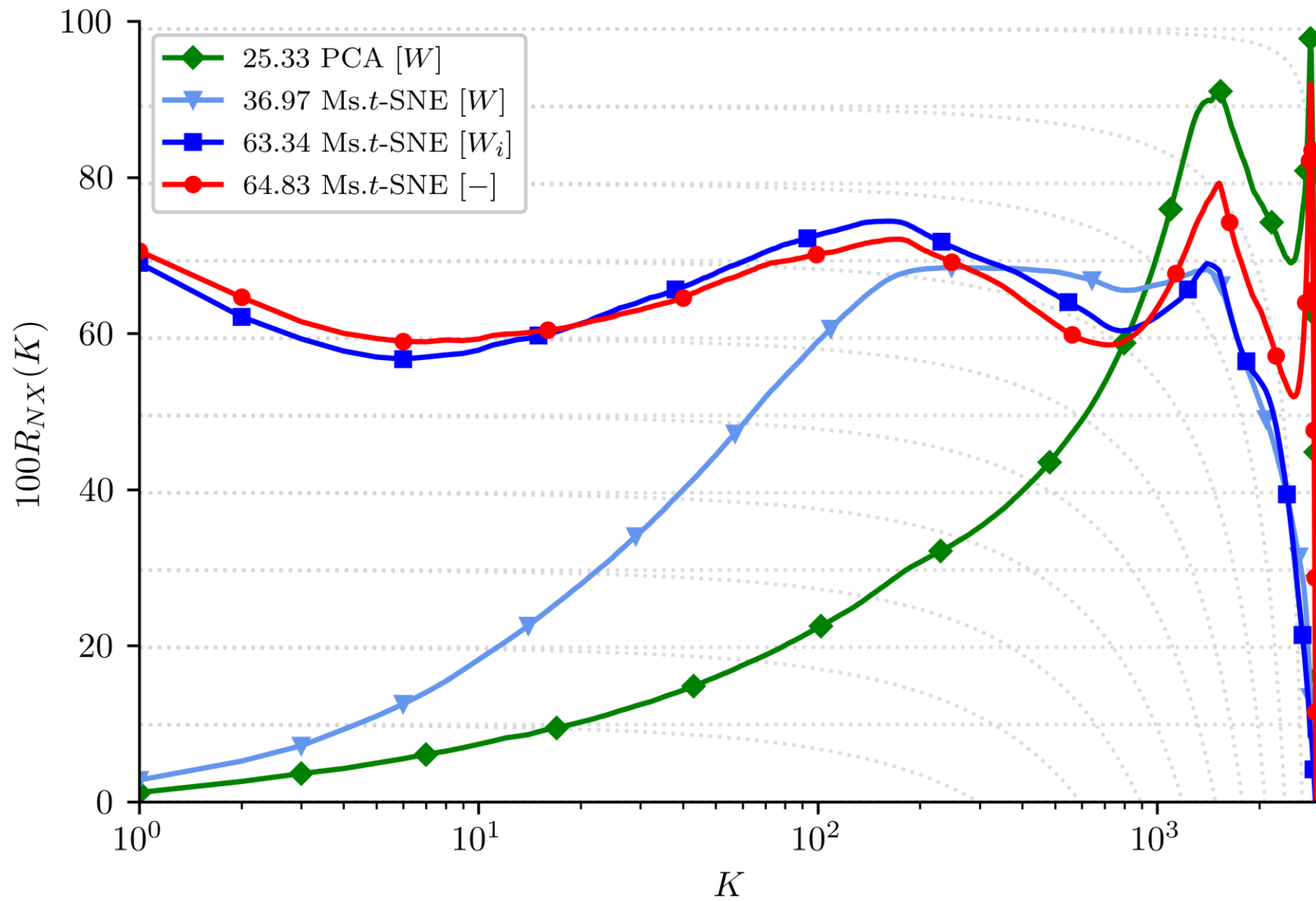


Interpretability **+** Quality **+**

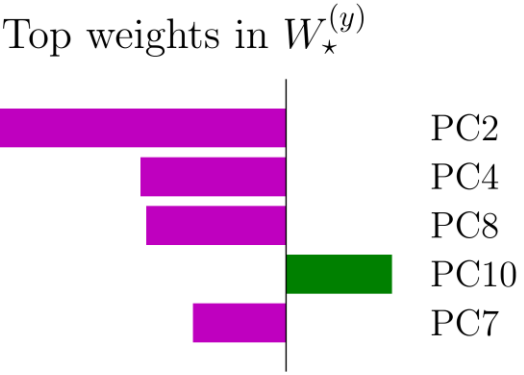
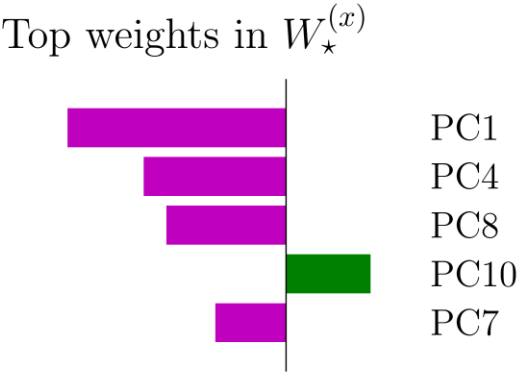
Ms. *t*-SNE [-]



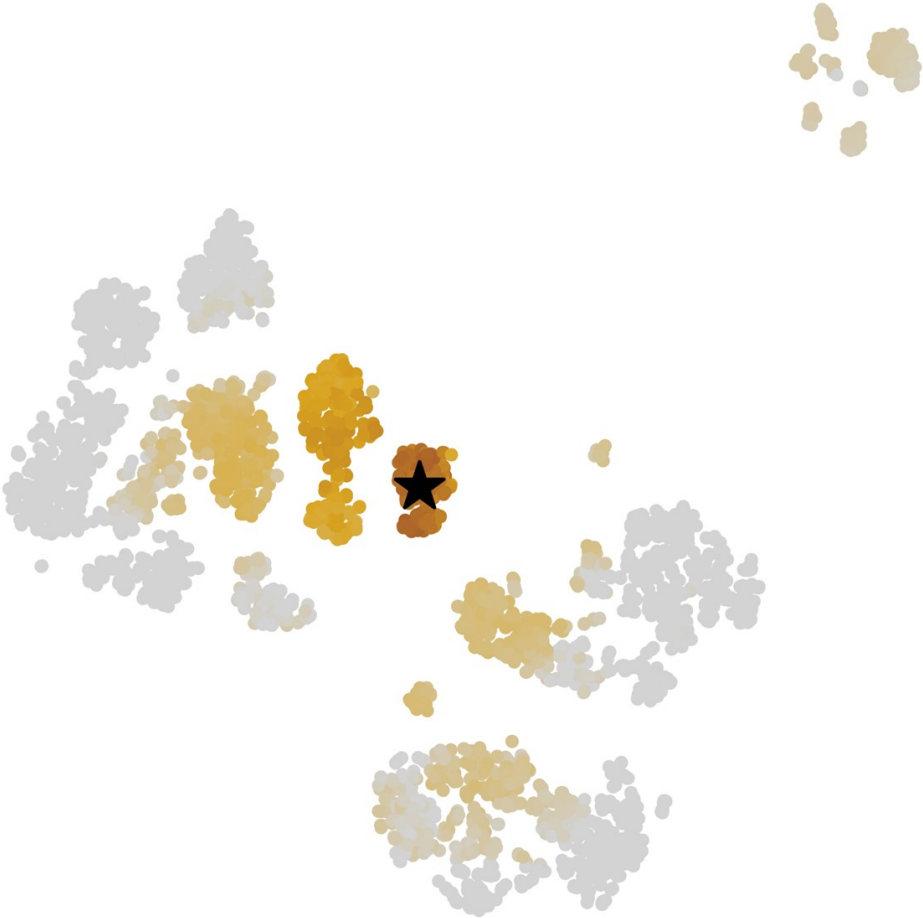
Interpretability **-** **-** Quality **+** **+**



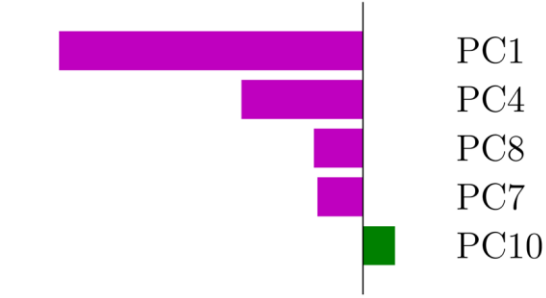
$\alpha = 10^6$     $\beta = 10^2$



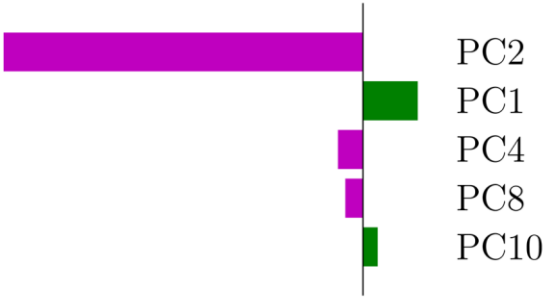
$\alpha = 10^7$     $\beta = 10^2$



Top weights in  $W_{\star}^{(x)}$



Top weights in  $W_{\star}^{(y)}$

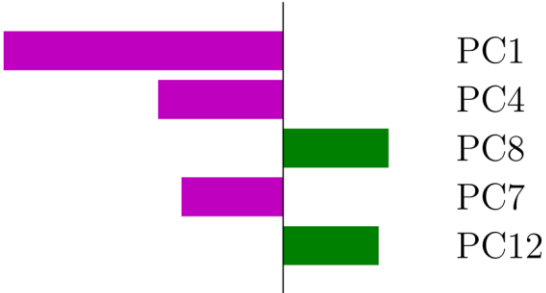


$$\alpha = 10^8$$

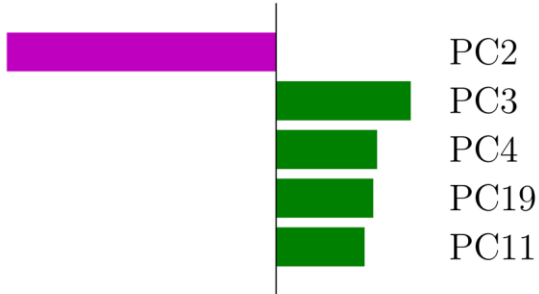
$$\beta = 10^2$$



Top weights in  $W_{\star}^{(x)}$



Top weights in  $W_{\star}^{(y)}$

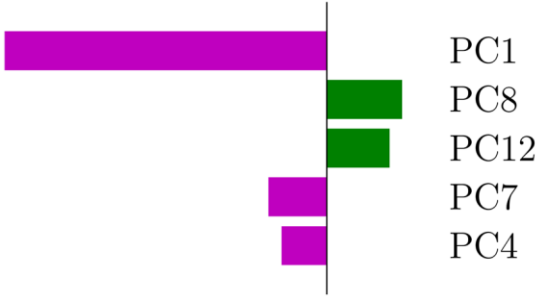


$$\alpha = 10^9$$

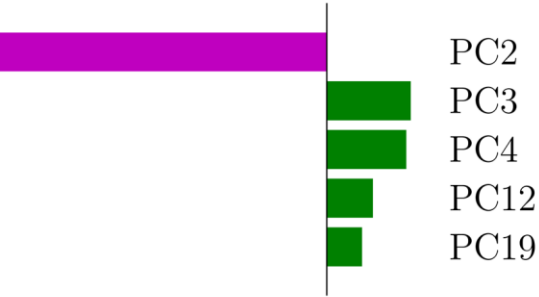
$$\beta = 10^2$$



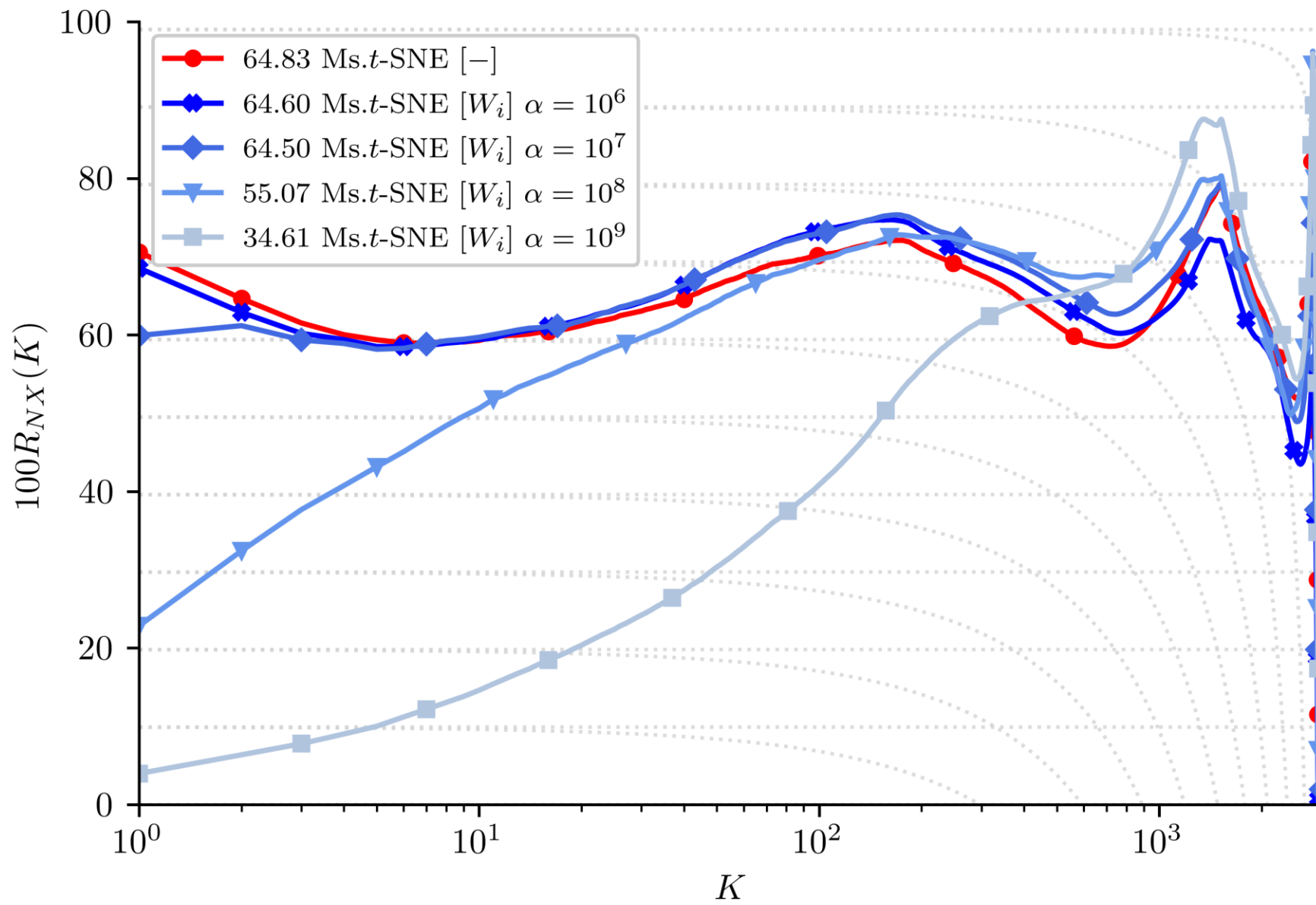
Top weights in  $W_{\star}^{(x)}$



Top weights in  $W_{\star}^{(y)}$







$\alpha = 10^6$

$\beta = 10^1$

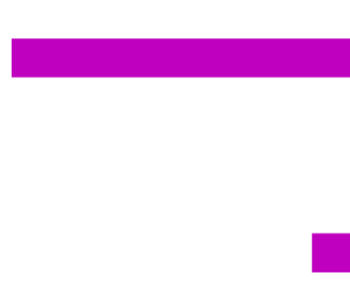


Top weights in  $W_{\star}^{(x)}$



PC1  
PC8  
PC3  
PC7  
PC4

Top weights in  $W_{\star}^{(y)}$



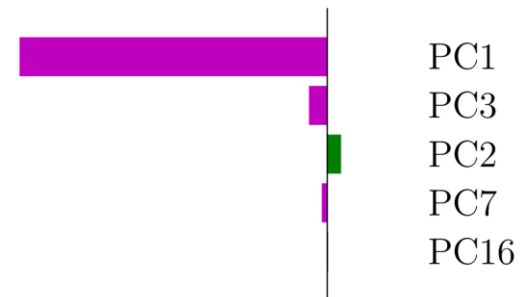
PC2  
PC3  
PC1  
PC7  
PC4

$$\alpha = 10^6$$

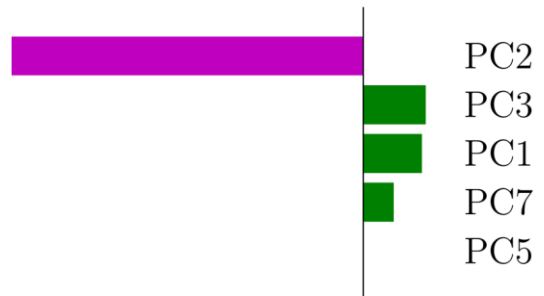
$$\beta = 10^2$$



Top weights in  $W_{\star}^{(x)}$

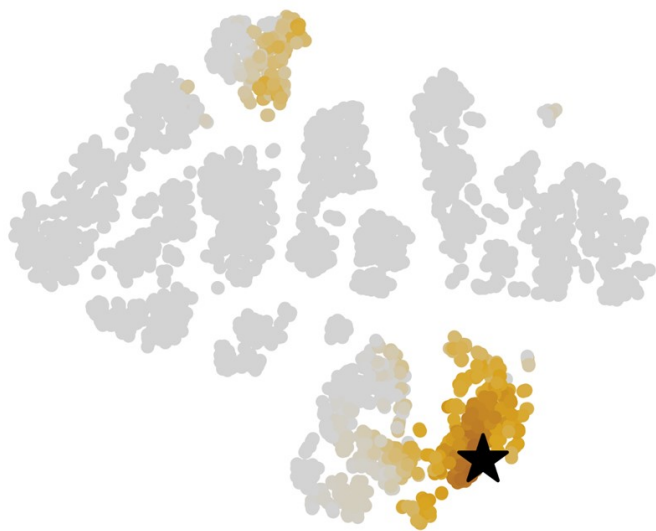


Top weights in  $W_{\star}^{(y)}$

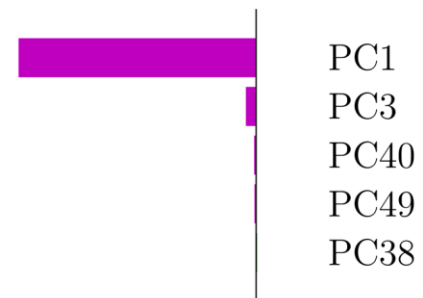


$$\alpha = 10^6$$

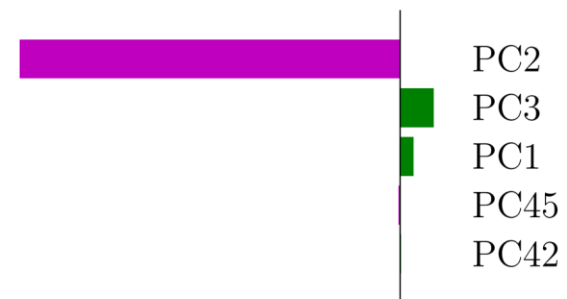
$$\beta = 10^3$$



Top weights in  $W_{\star}^{(x)}$

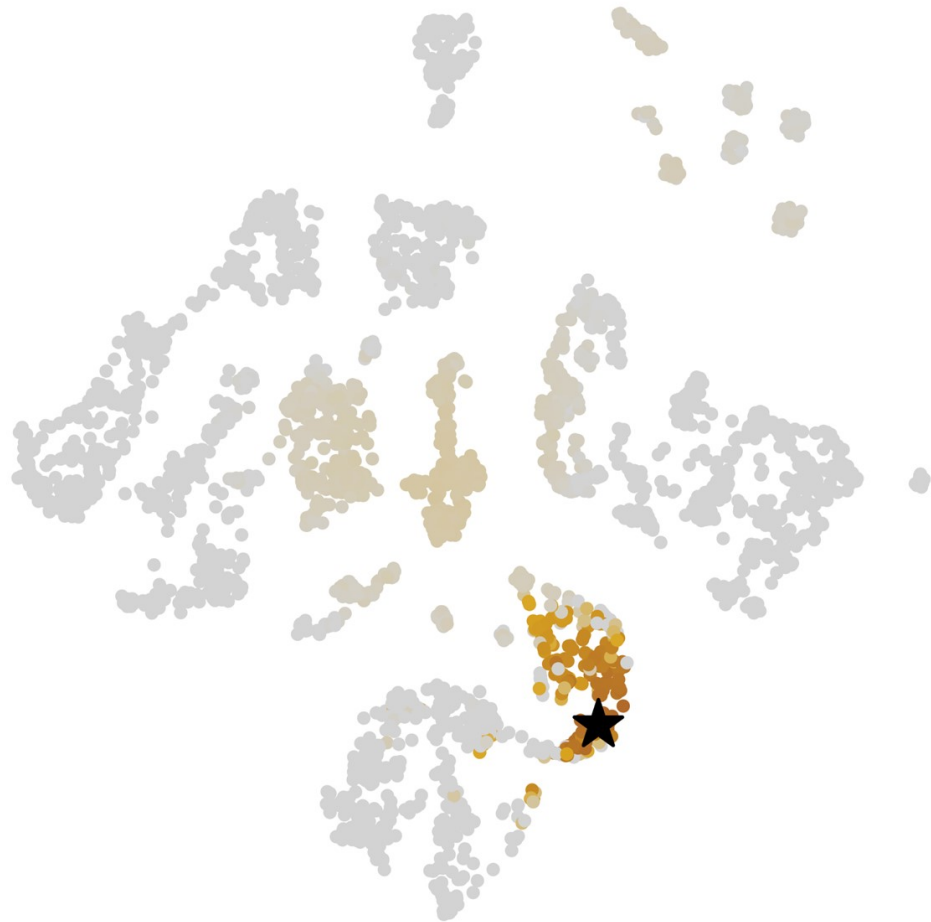


Top weights in  $W_{\star}^{(y)}$

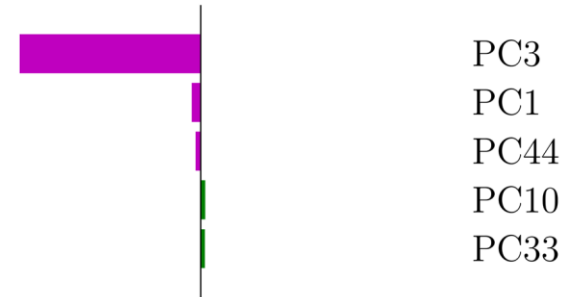


$\alpha = 10^6$

$\beta = 10^4$



Top weights in  $W_{\star}^{(x)}$



Top weights in  $W_{\star}^{(y)}$

