

# Natural Language Processing

Lecture 11

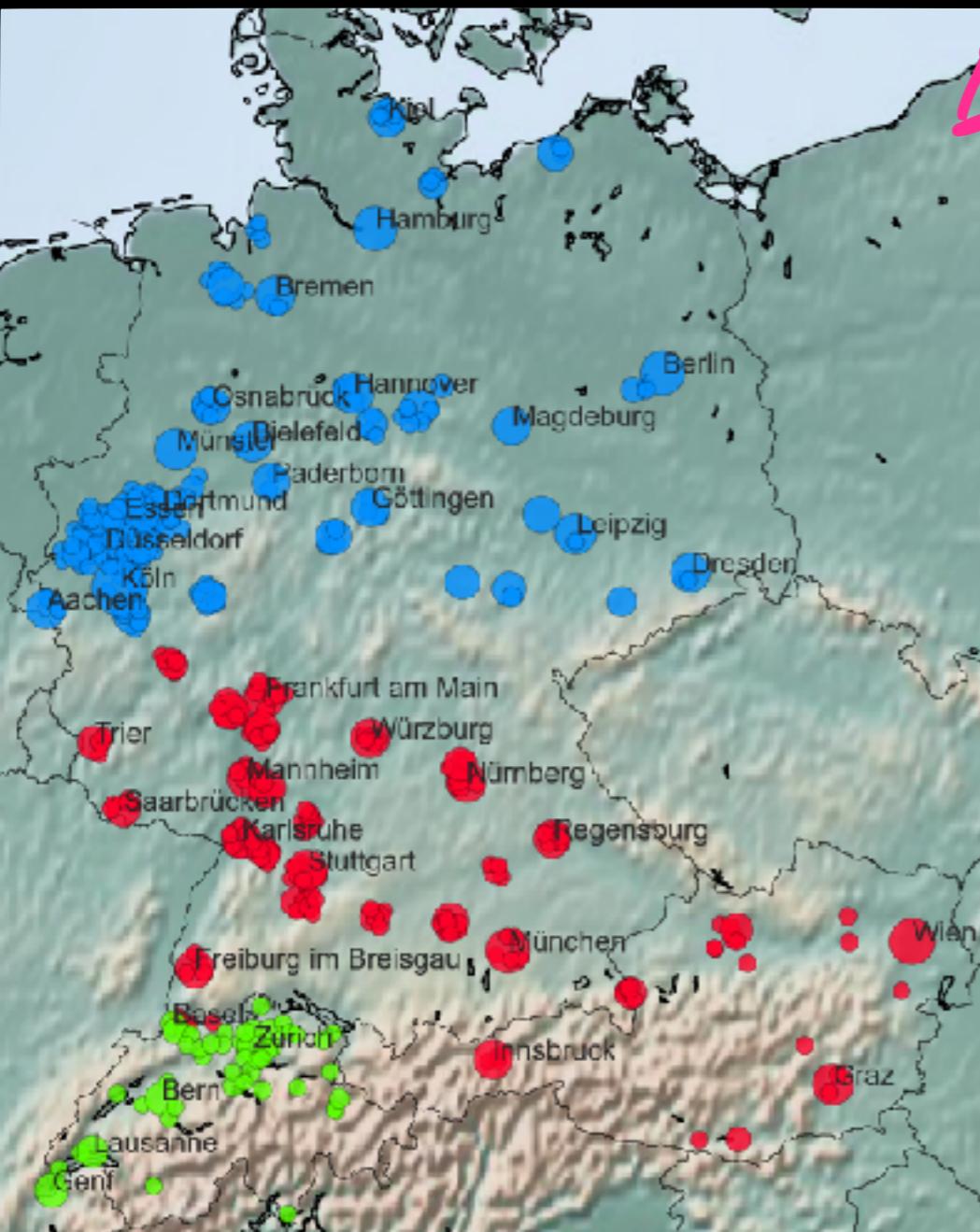
Dirk Hovy

[dirk.hovy@unibocconi.it](mailto:dirk.hovy@unibocconi.it)

 @dirk\_hovy

Bocconi

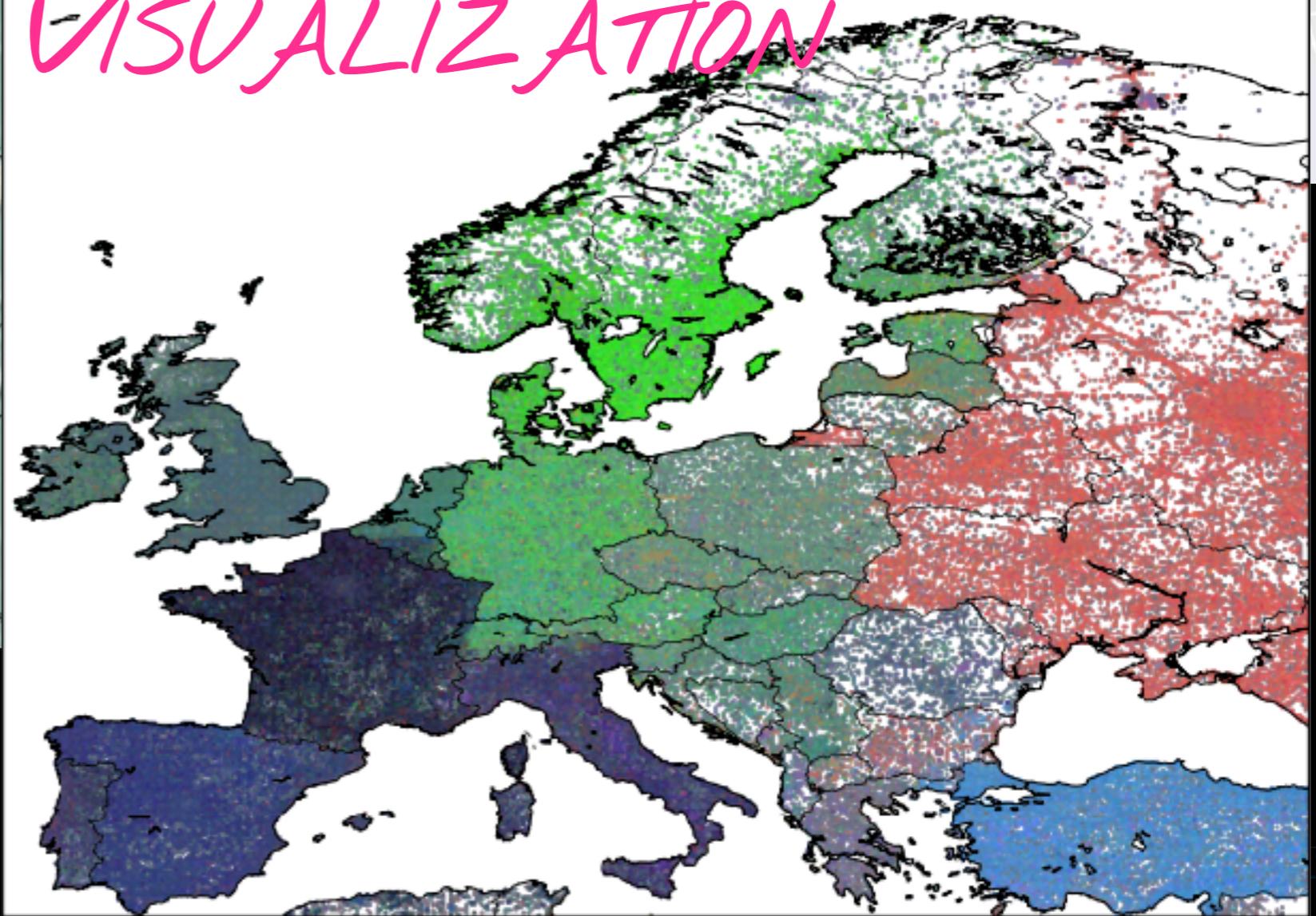
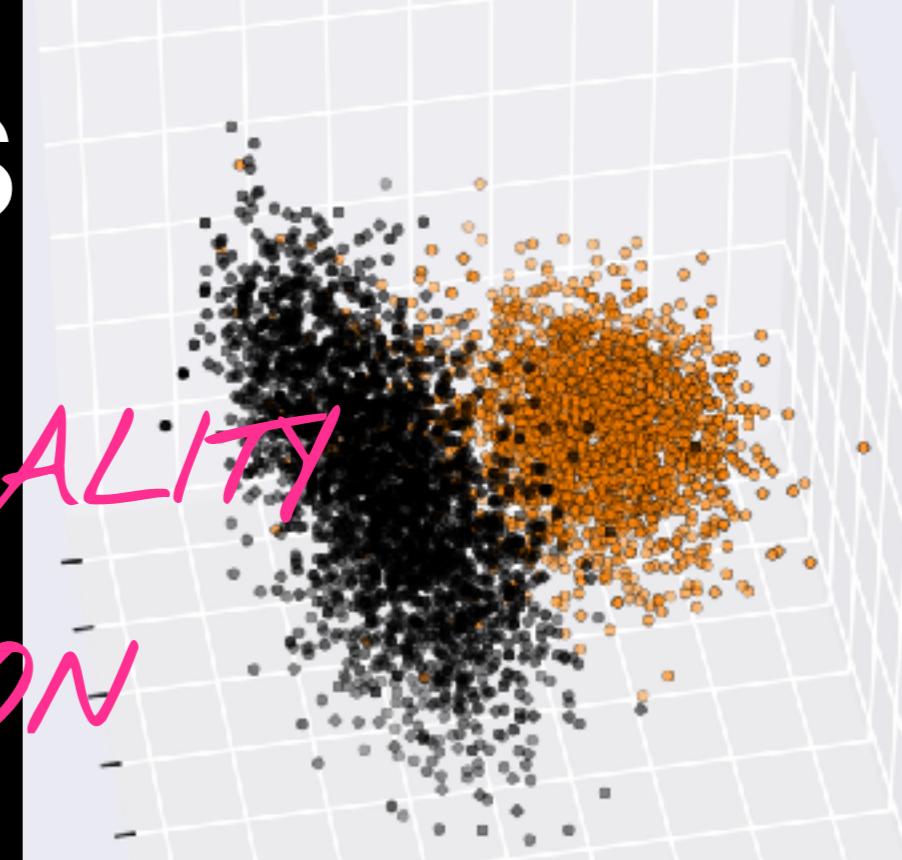
# Examples



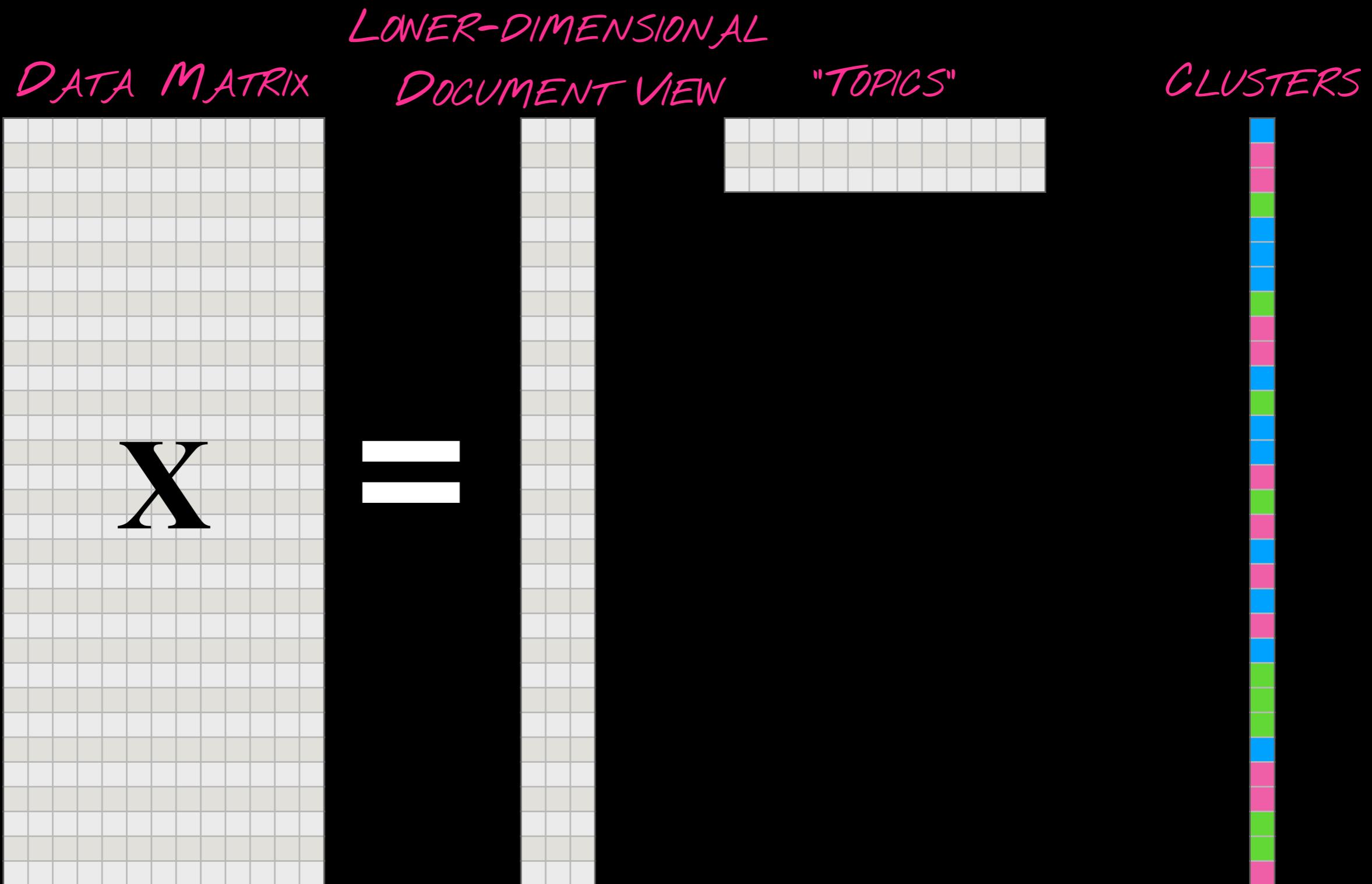
CLUSTERING

DIMENSIONALITY  
REDUCTION

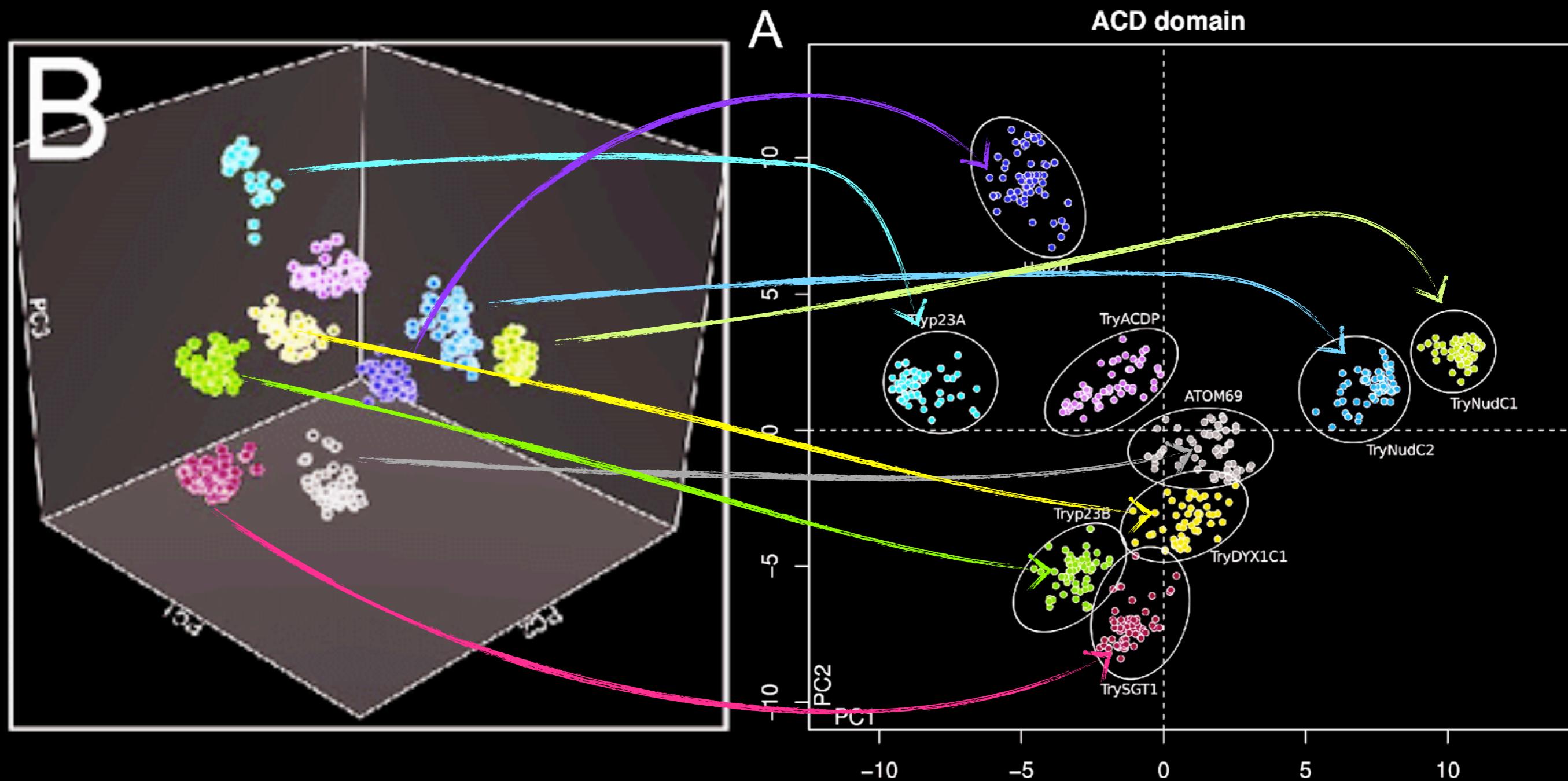
VISUALIZATION



# Latent Dimensions



# Latent Dimensions

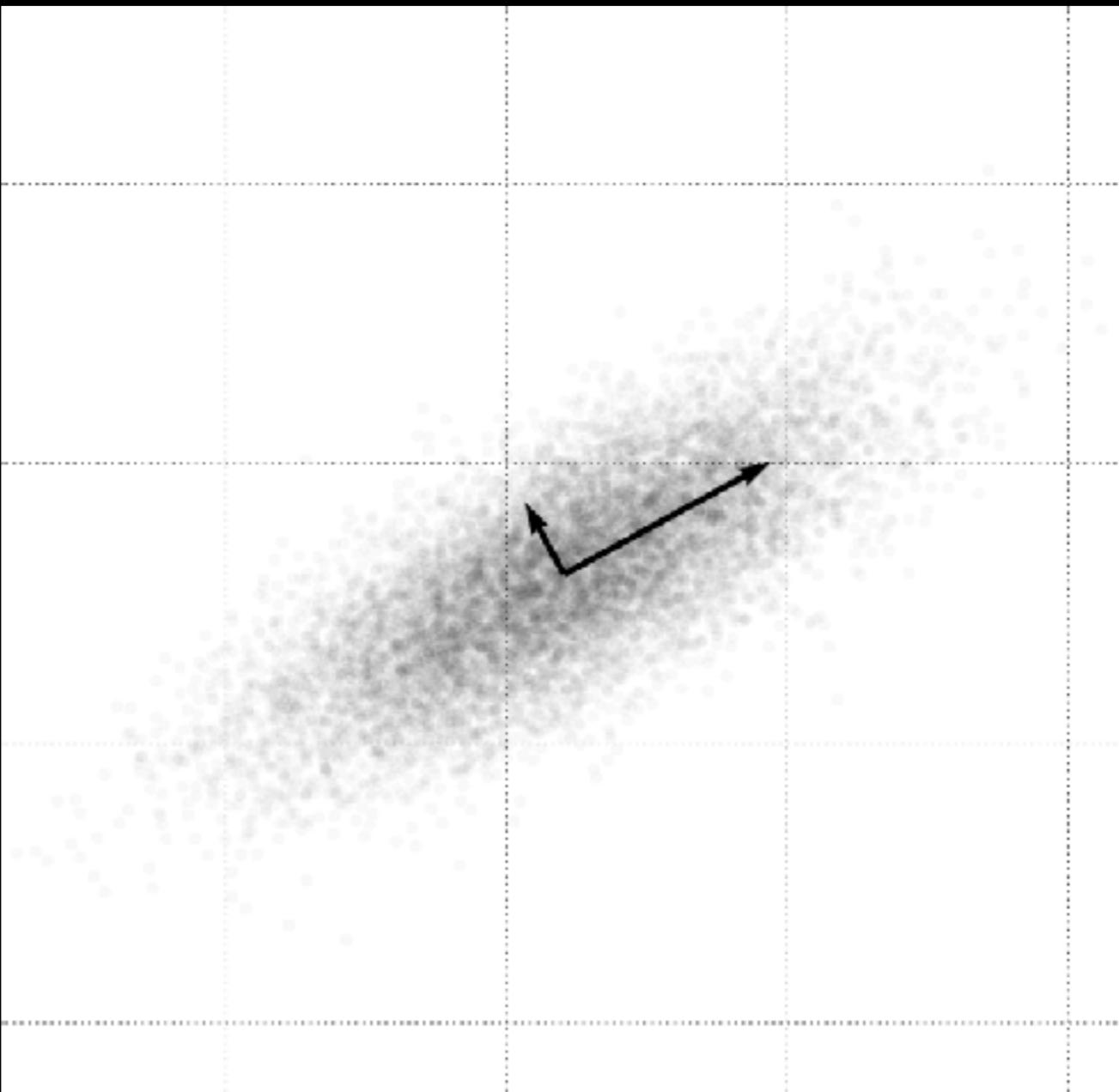


# Goals for Today

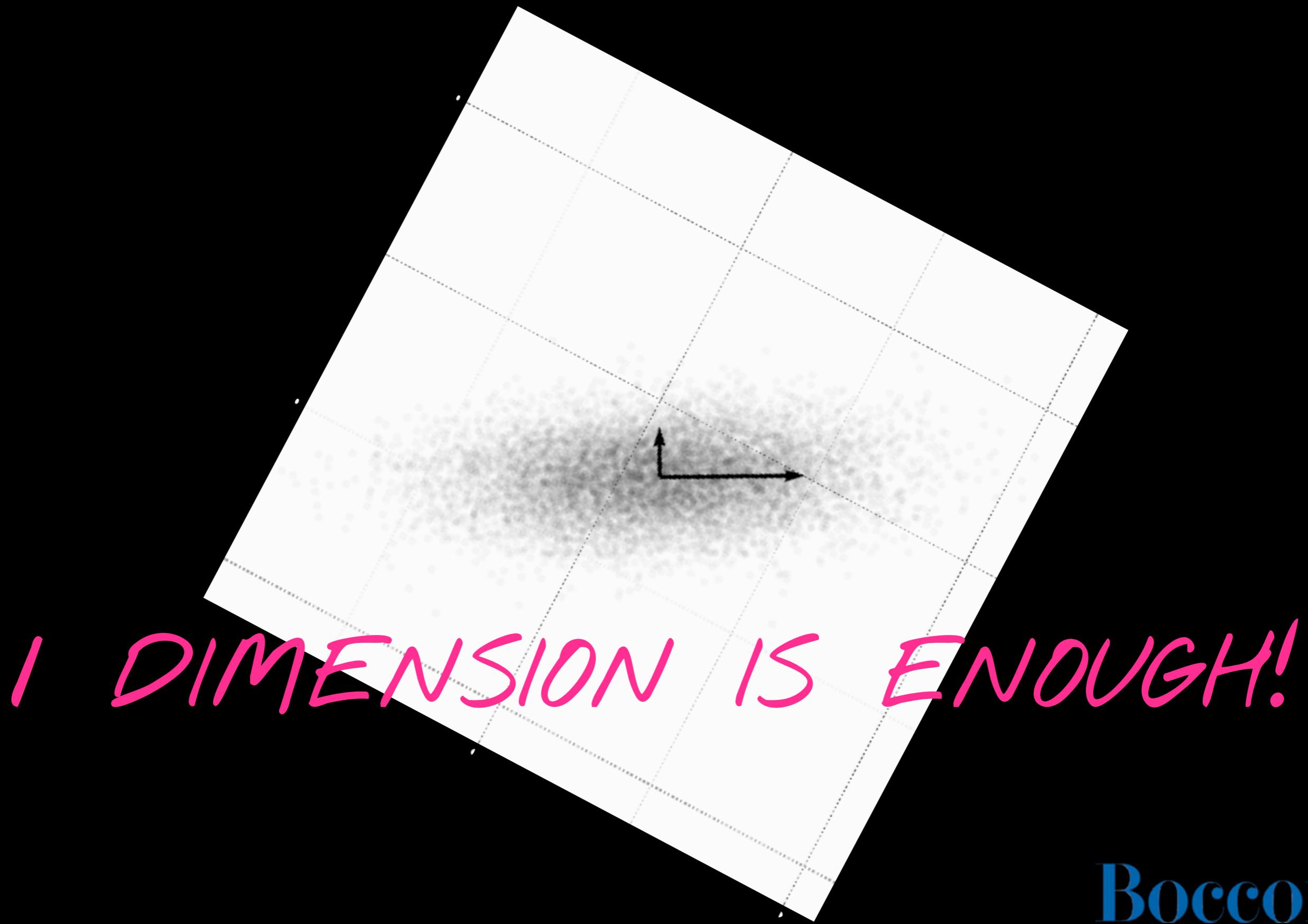
- Learn about **matrix factorization** and its use for **semantic similarity** and **visualization**
- Learn about **k-means** and **agglomerative clustering**
- Learn about **evaluation** criteria

# Matrix Factorization

# Singular Value Decomposition



# Singular Value Decomposition



*1 DIMENSION IS ENOUGH!*

# Singular Value Decomposition

- “principal component analysis”: discover the dimensions that matter
- idea: matrix is made up of few hidden dimensions
- Dimensions correspond to **documents**, **terms**, and **latent concepts**

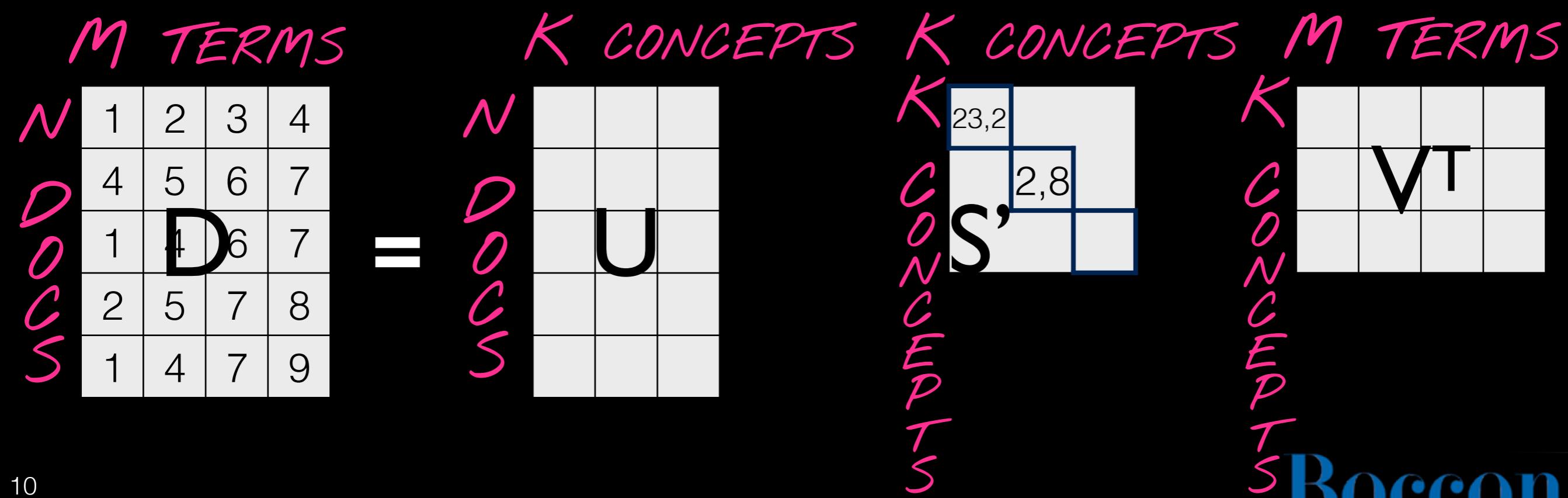
$$\begin{matrix} M \text{ TERMS} \\ N \\ D \\ O \\ C \\ S \end{matrix} \quad \begin{matrix} K \text{ CONCEPTS} \\ N \\ D \\ O \\ C \\ S \end{matrix} \quad = \quad \begin{matrix} K \text{ CONCEPTS} \\ K \\ C \\ O \\ N \\ C \\ E \\ P \\ T \\ S \end{matrix} \quad \begin{matrix} M \text{ TERMS} \\ K \\ C \\ O \\ N \\ C \\ E \\ P \\ T \\ S \end{matrix}$$

The diagram illustrates the Singular Value Decomposition (SVD) of a document-term matrix  $D$ . The matrix  $D$  has dimensions  $M \times N$  (5 terms by 4 documents). It is decomposed into three matrices:  $U$  (orthogonal matrix of size  $N \times N$ ),  $S$  (diagonal matrix of size  $K \times K$  containing singular values), and  $V^T$  (orthogonal matrix of size  $M \times M$ ). The matrix  $S$  is shown with its top-left element highlighted as 23,2.

The diagram shows the decomposition of a 5x4 document-term matrix  $D$  into  $U$ ,  $S$ , and  $V^T$ . The labels are color-coded: pink for rows and columns, black for the matrices themselves. The matrix  $S$  is highlighted with a blue border around its top-left corner, which contains the value 23,2.

# Singular Value Decomposition

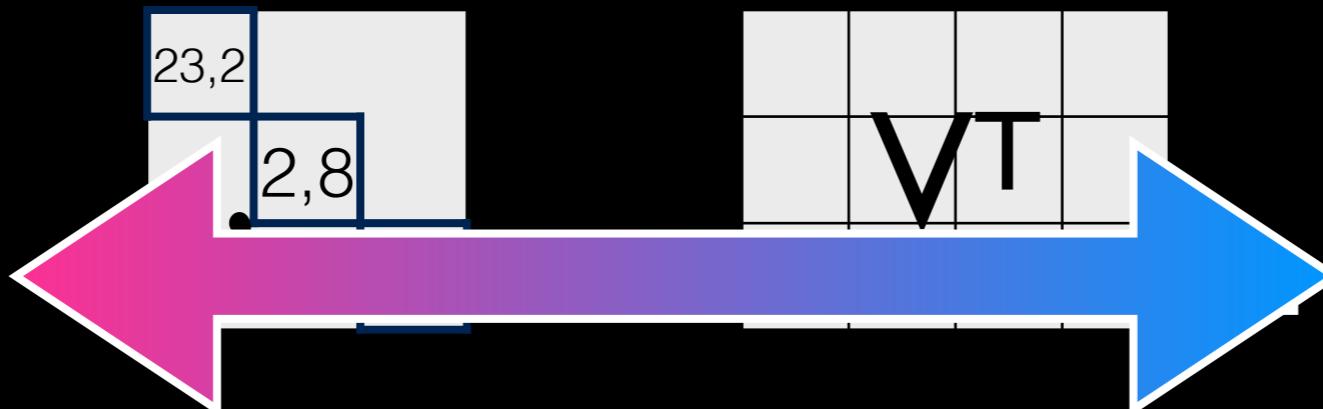
- reduce principal components/concepts to smaller number



# Singular Value Decomposition

- reconstruct original matrix in new concept space:  
**Latent Semantic Analysis**

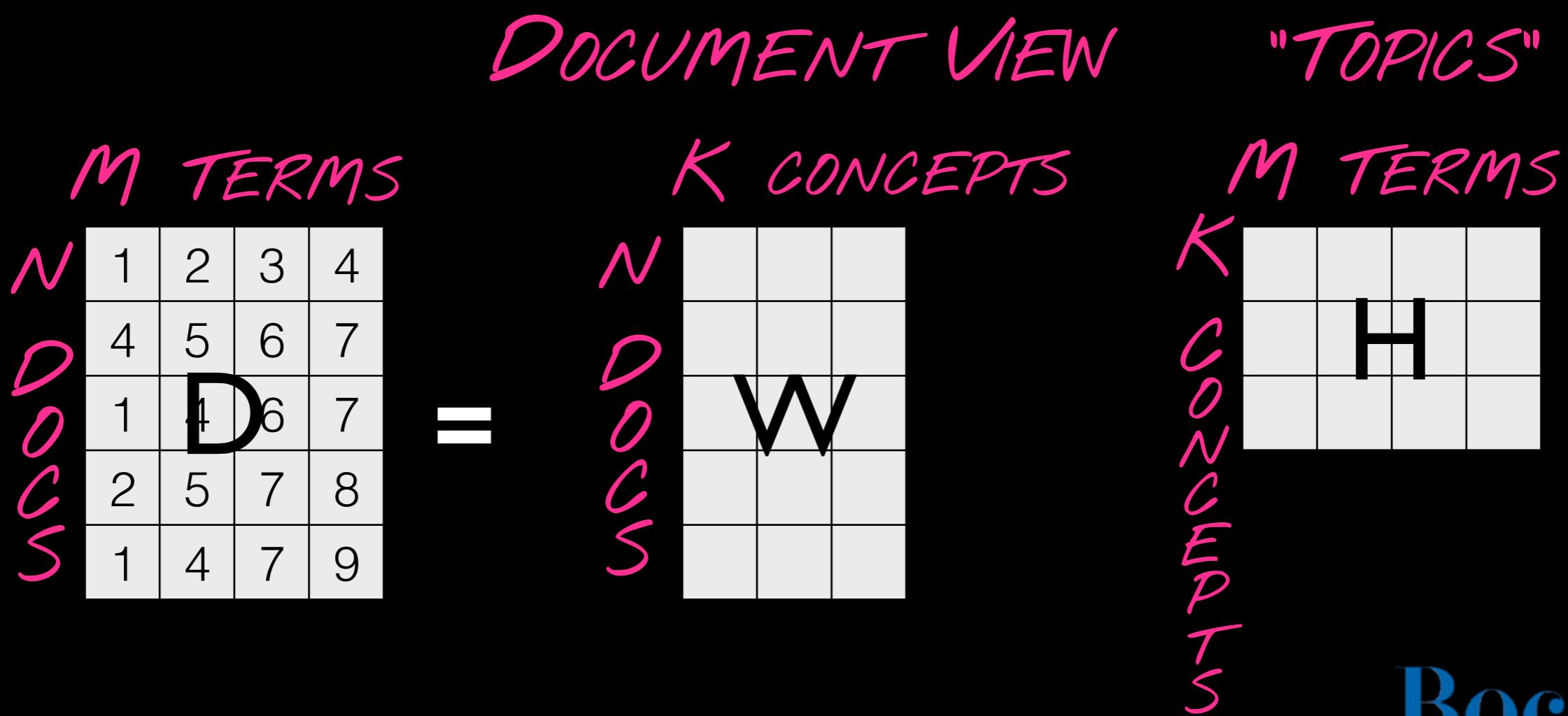
1	2	3	4
4	5	6	7
1	4	6	7
2	5	7	8
1	4	7	9



0,9	2,1	3,2	3,8
3,9	5,1	6	6,9
1,1	3,8	4,9	7,2
2,2	4,7	6,9	8,2
0,8	4,3	7,1	8,8

# Non-negative Matrix Factorization

- Use only positive values
- Find approximation of two components



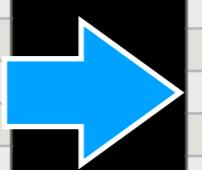
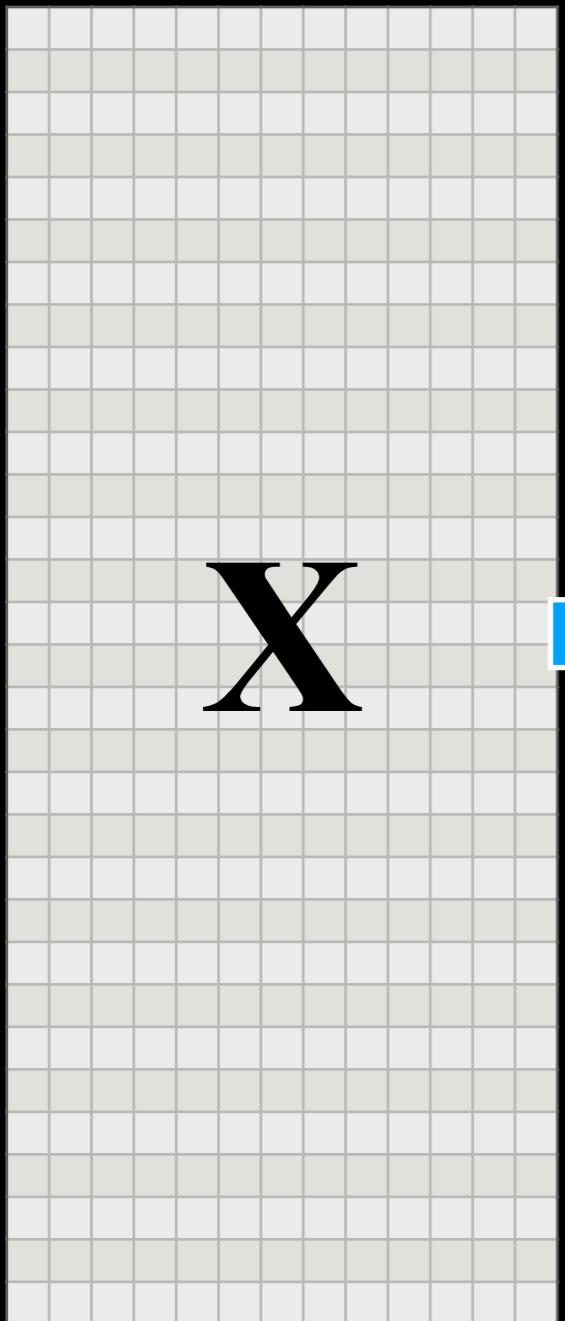
# Comparison

	SVD	NMF
<b>Negative values (embeddings) as input?</b>	yes	no
<b>#components</b>	$3: U, S, V$	$2: W, H$
<b>document view?</b>	yes: $U$	yes: $W$
<b>term view?</b>	yes: $V$	yes: $H$
<b>strength ranking?</b>	yes: $S$	no
<b>exact?</b>	yes	no
<b>"topic" quality</b>	mixed	better
<b>sparsity</b>	low	medium

# Yes, but: What is it Good for?

- Find latent **topic** dimensions (alternative: LDA)
- Find **word similarity** in latent space (alternative: Word2Vec)
- Find **document similarity** in latent space (alternative: Doc2Vec)
- Reduce dimensionality for **visualization**

# Latent Word Dimension Topics



EVT

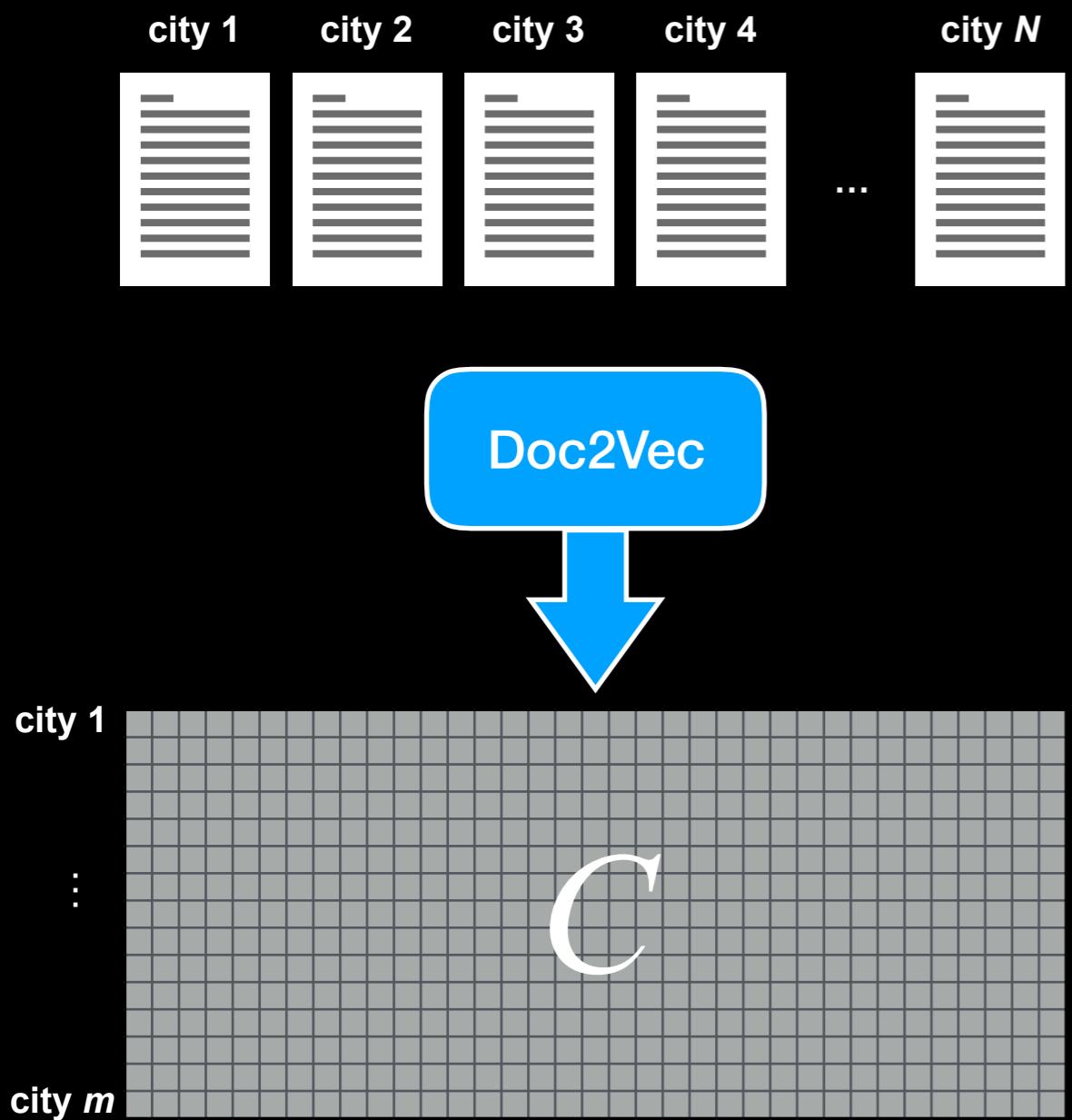
NMF(10)

FOR MOBY DICK

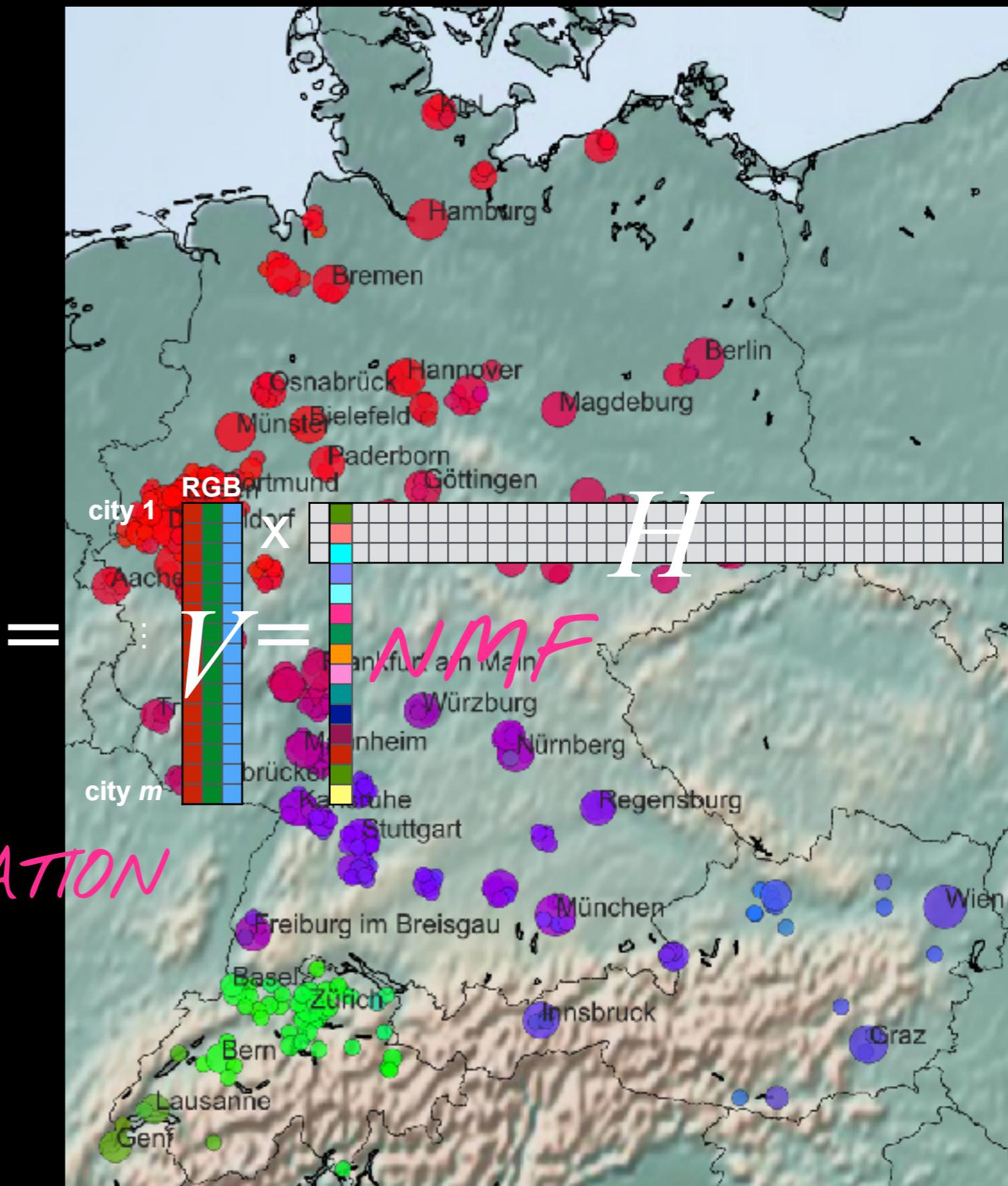
AHAB	ahab, captain, cried, captain ahab, cried ahab
STRUCTURE	chapter, folio, octavo, ii, iii
???	like, ship, sea, time, way
MEN	man, old, old man, look, young man
???	oh, life, starbuck, sweet, god
CHARACTERS	said, stubb, queequeg, don, starbuck
???	sir, aye, let, shall, think
OLD-TIMEY	thou, thee, thy, st, god
WHALES	whale, sperm, sperm whale, white, white whale
???	ye, look, say, ye ye, men

# Dimensionality Reduction for Visualizations

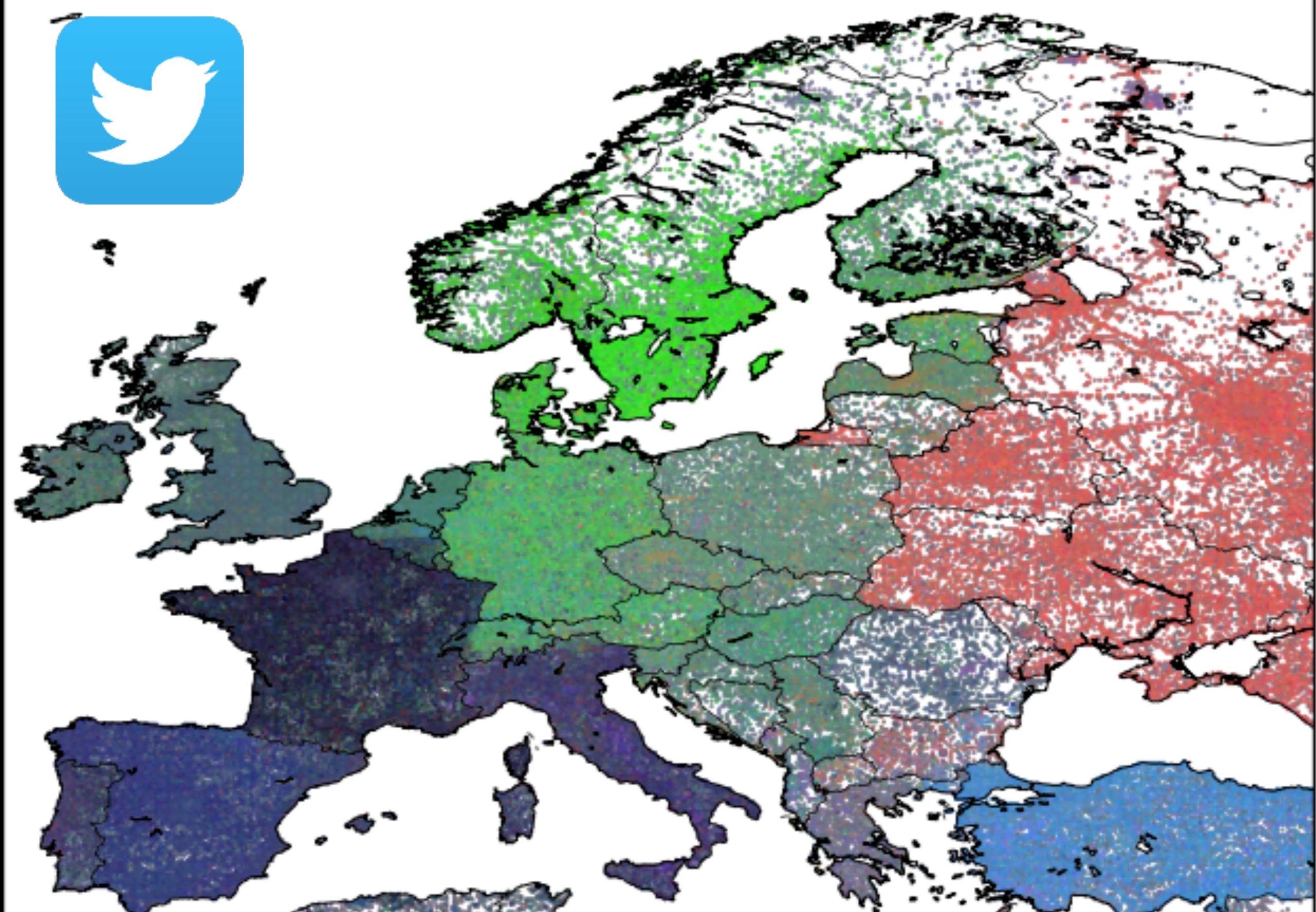
# Dimensions as RGB



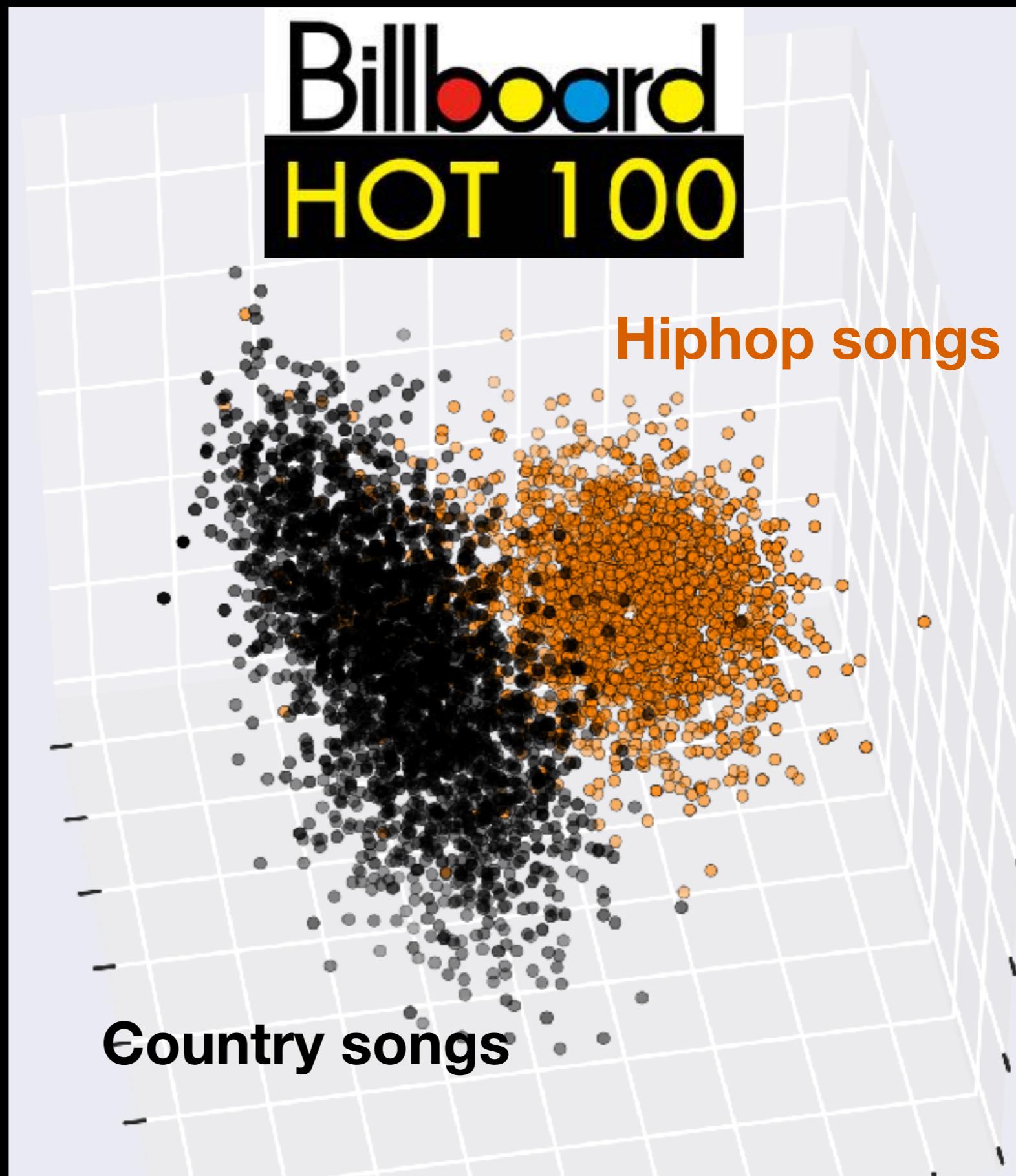
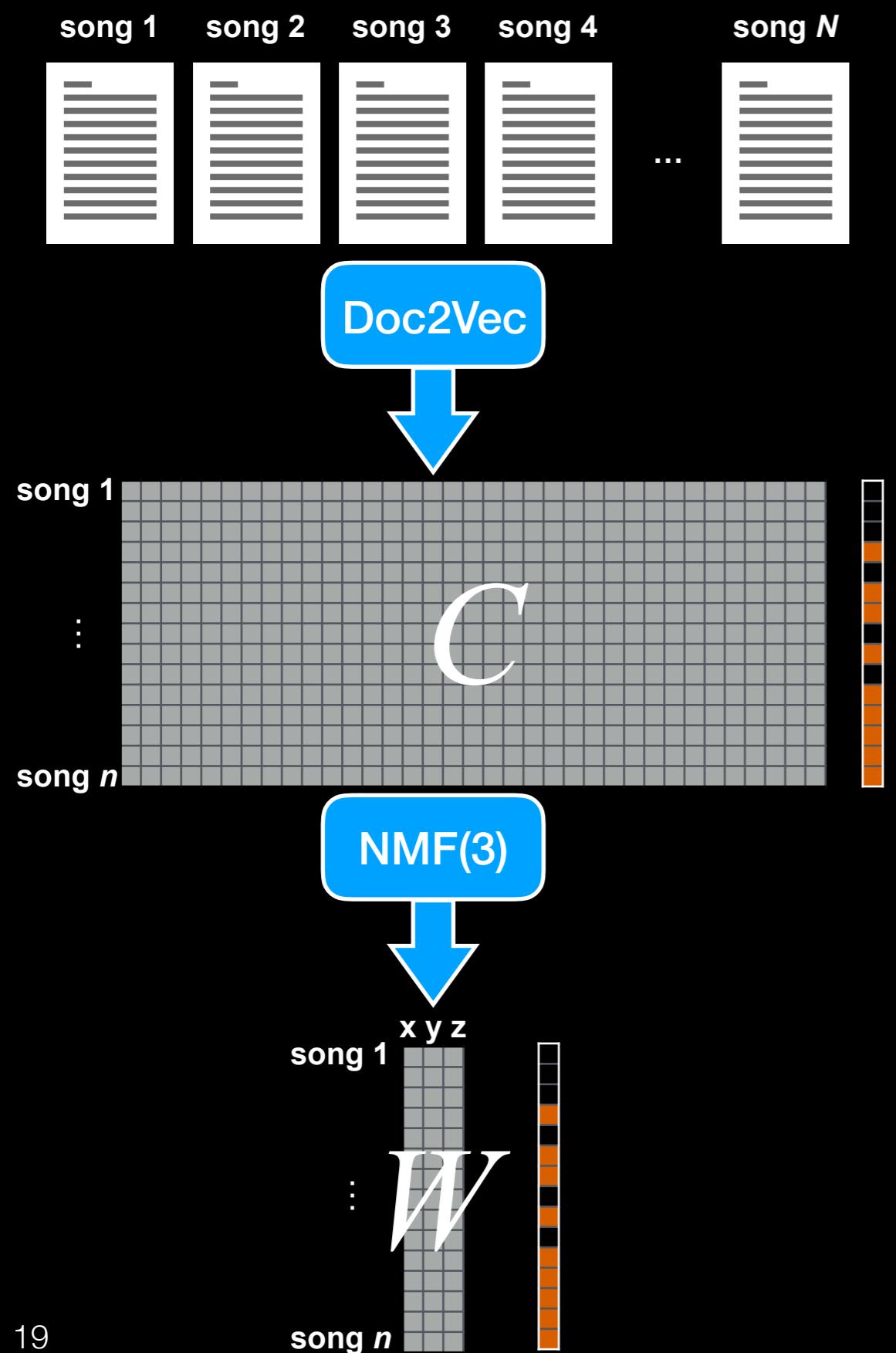
*DENSE REPRESENTATION*



# Dimensions as RGB

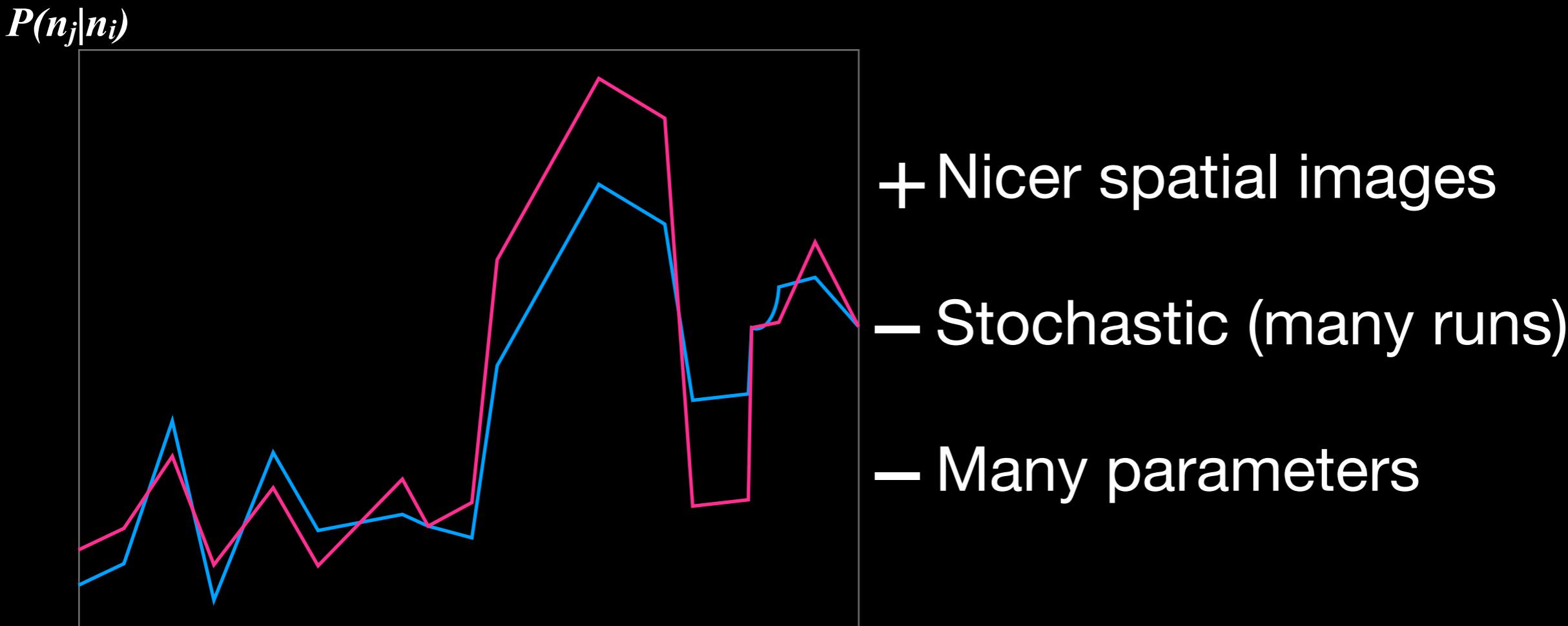


# Dimensions as Coordinates



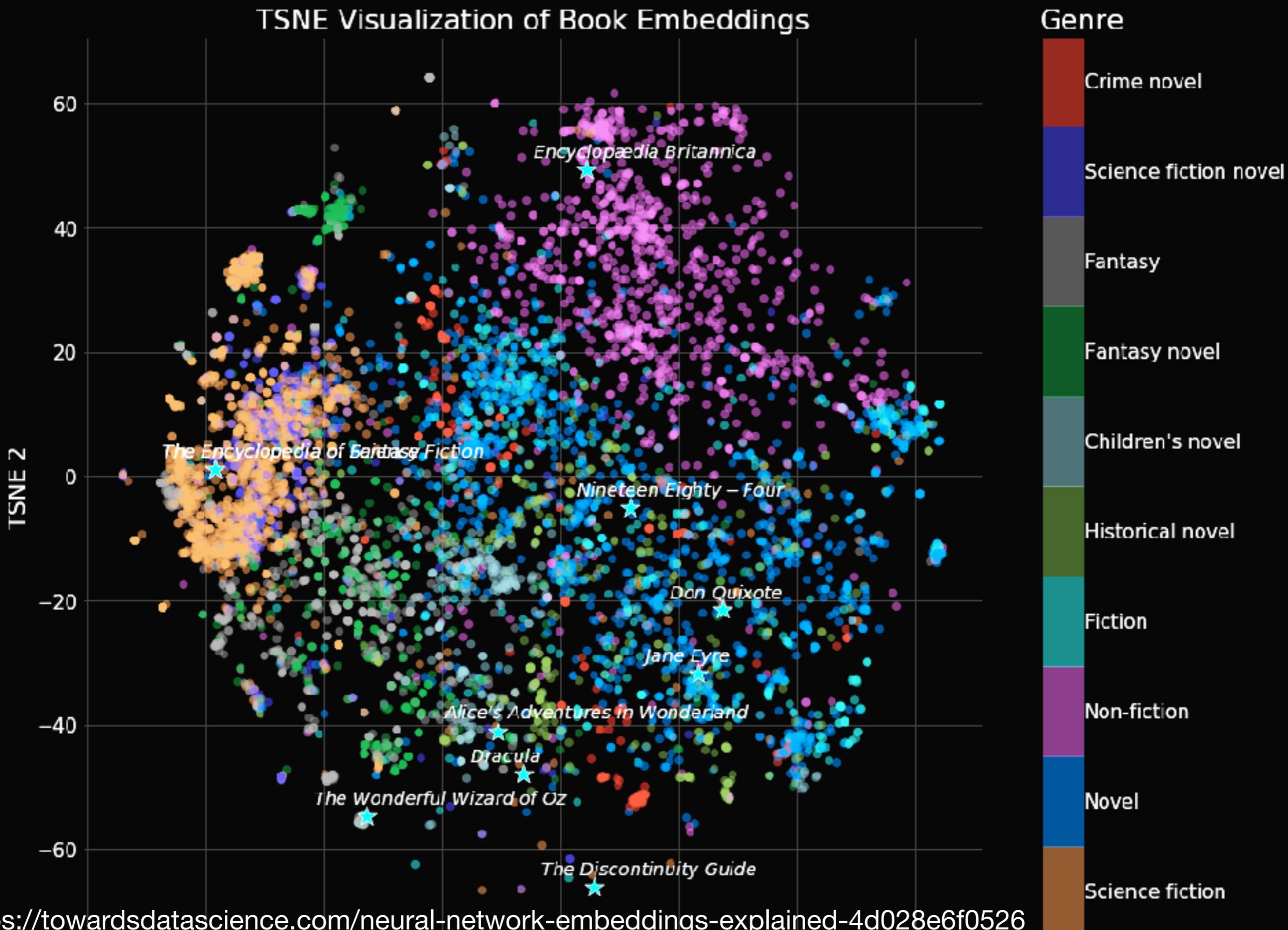
# t-SNE

- Map/preserve neighborhood structure of high-dimensional space in lower-dimensional space
- Minimize difference between probability distributions over neighbors in both dimensions



# t-SNE

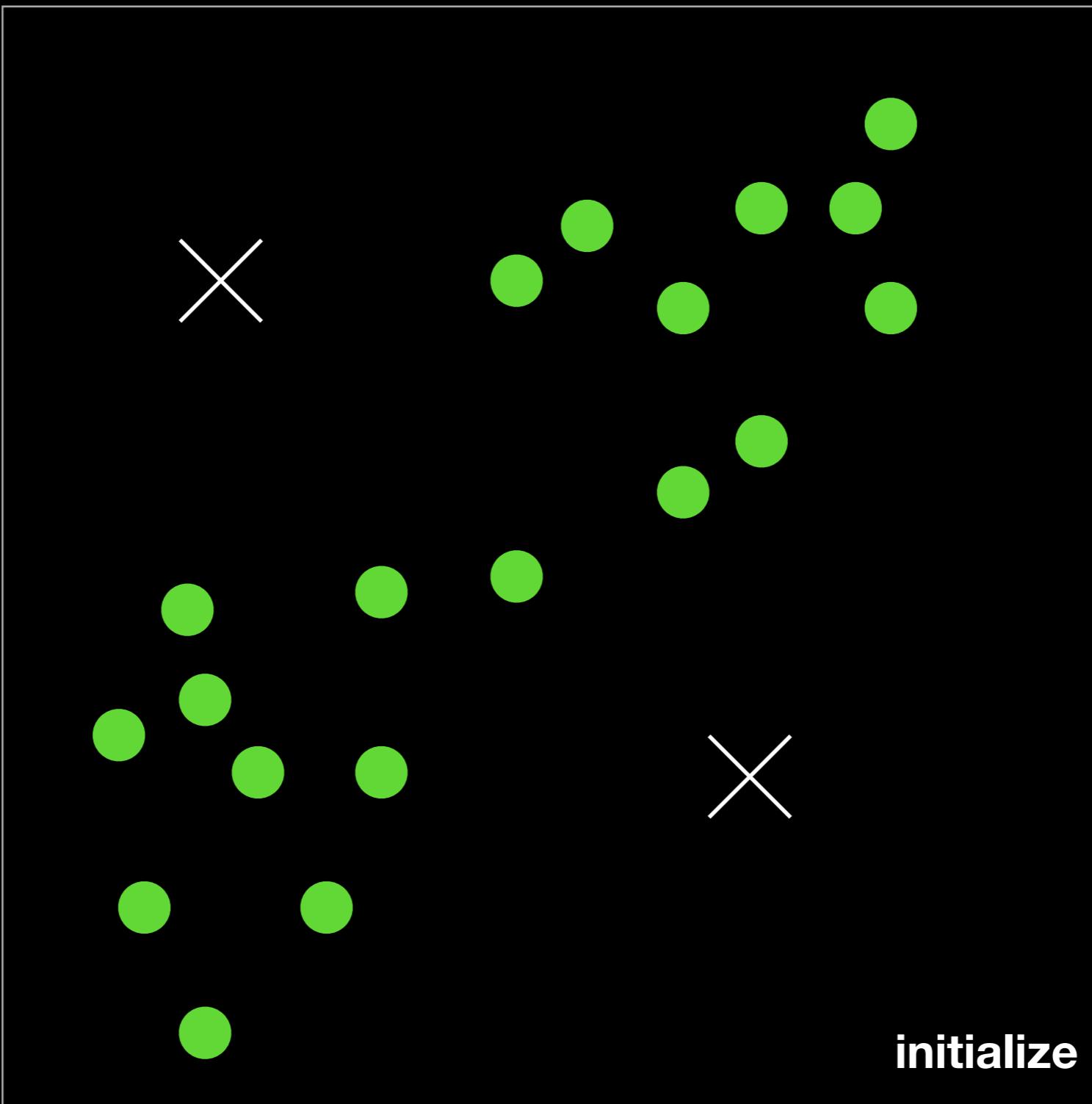
TSNE Visualization of Book Embeddings



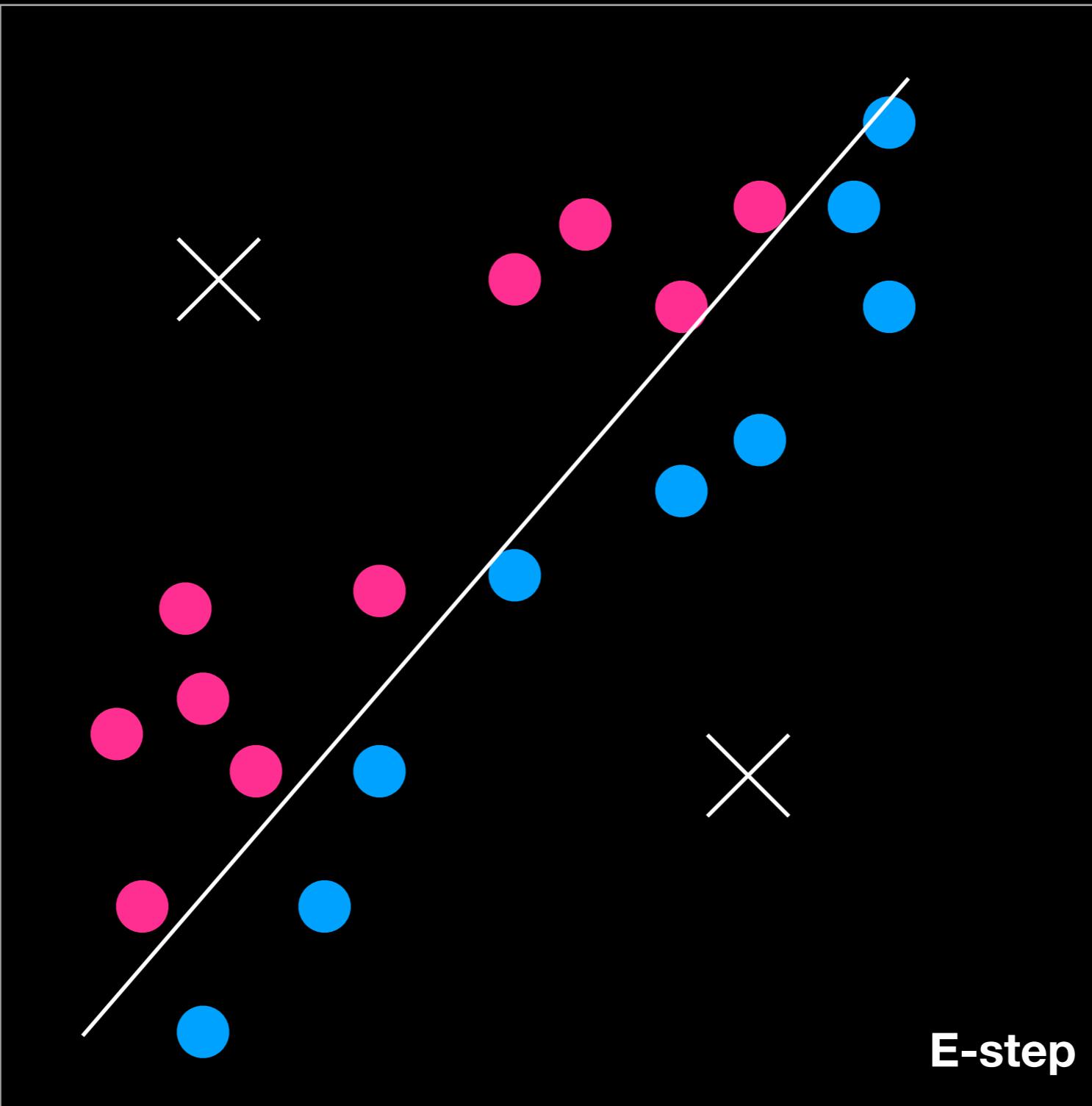
# Clustering

# $k$ -Means Clustering

# $k$ -Means

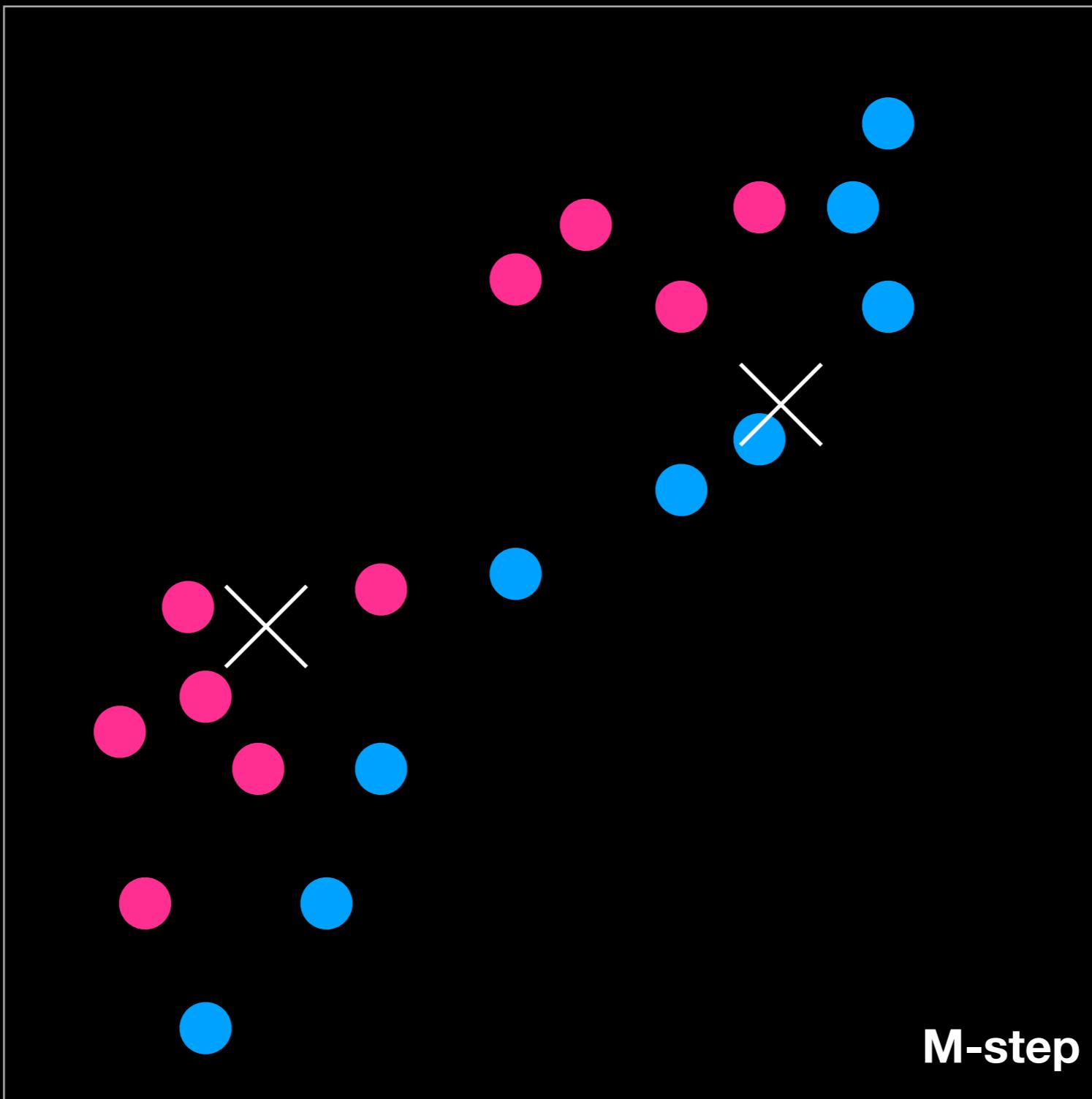


# $k$ -Means



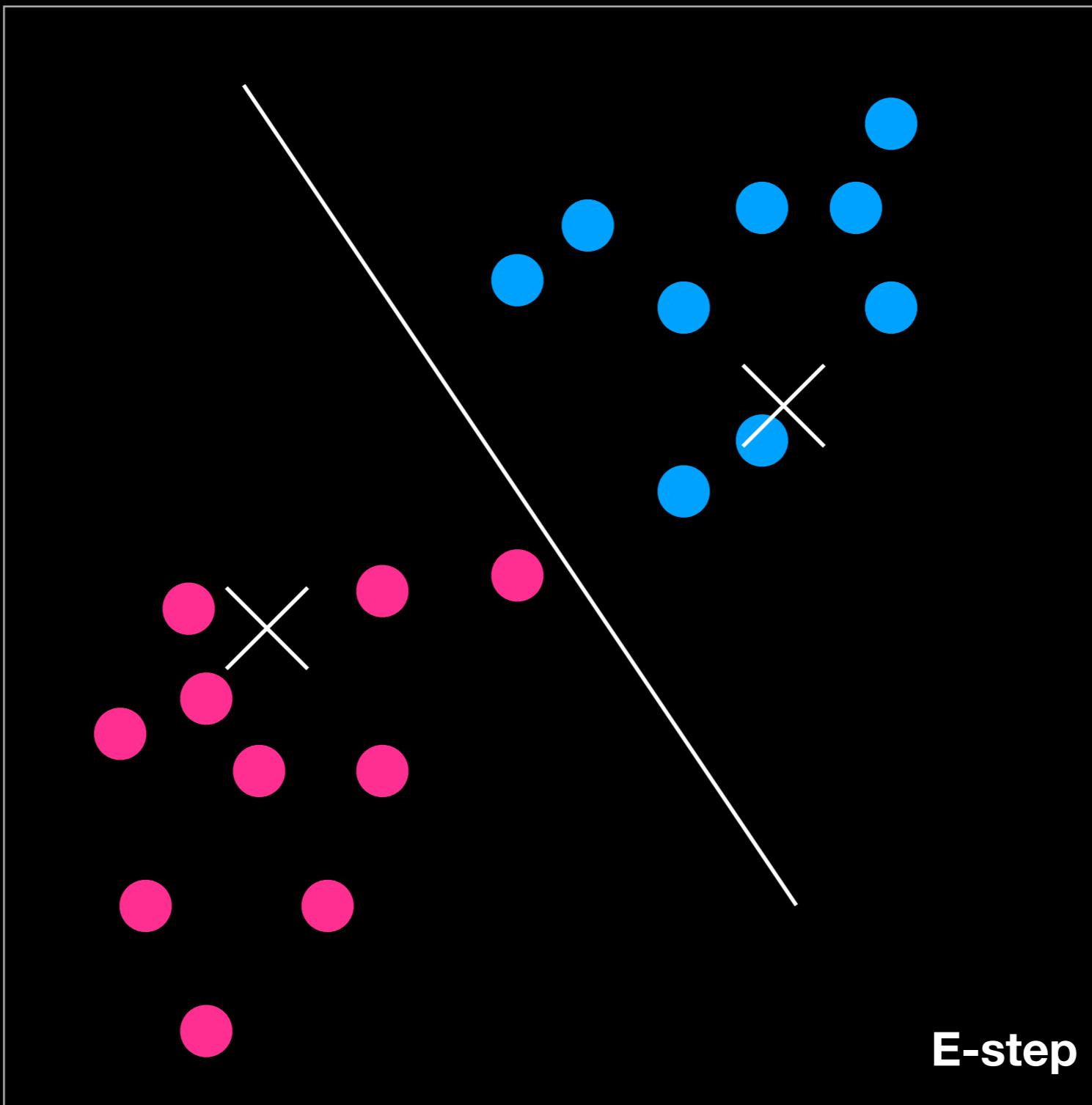
*ASSIGN POINTS TO CENTROIDS*

# $k$ -Means

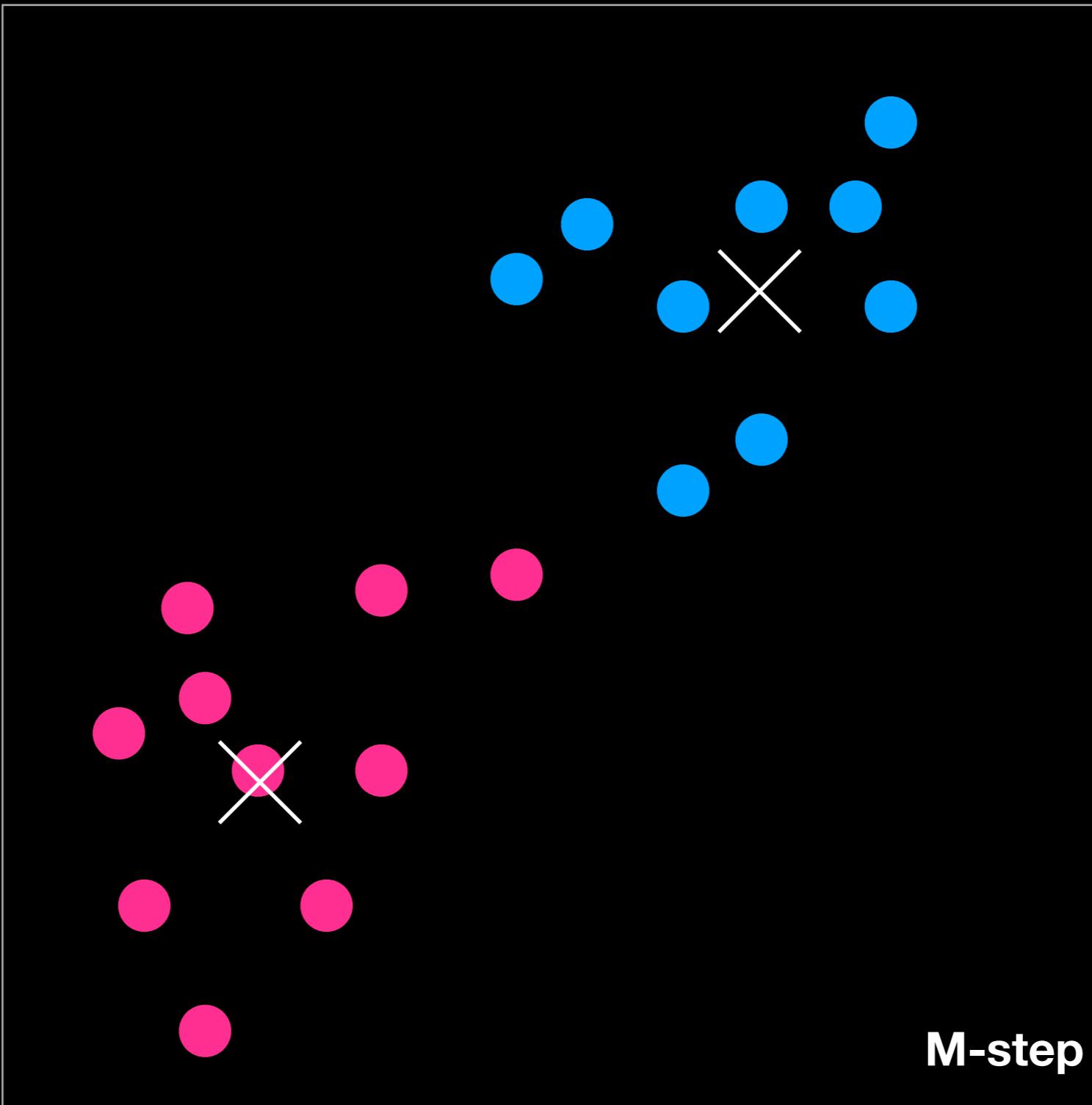


*RECOMPUTE CENTROIDS*

# $k$ -Means

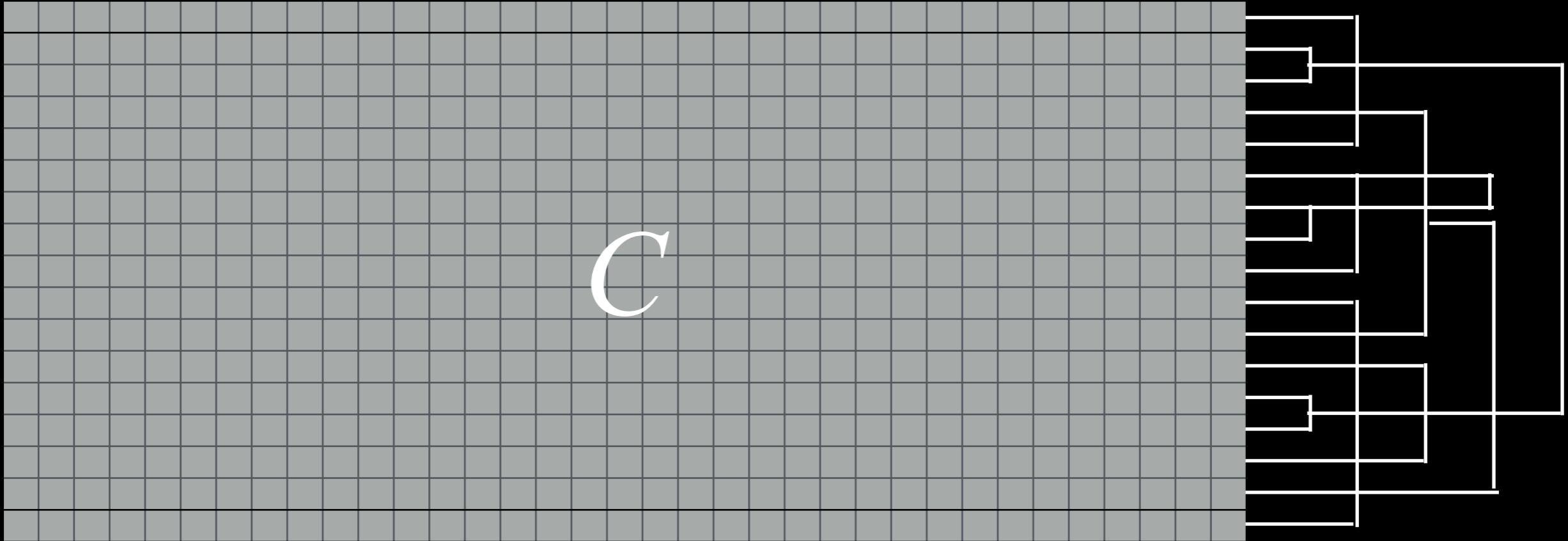


# $k$ -Means

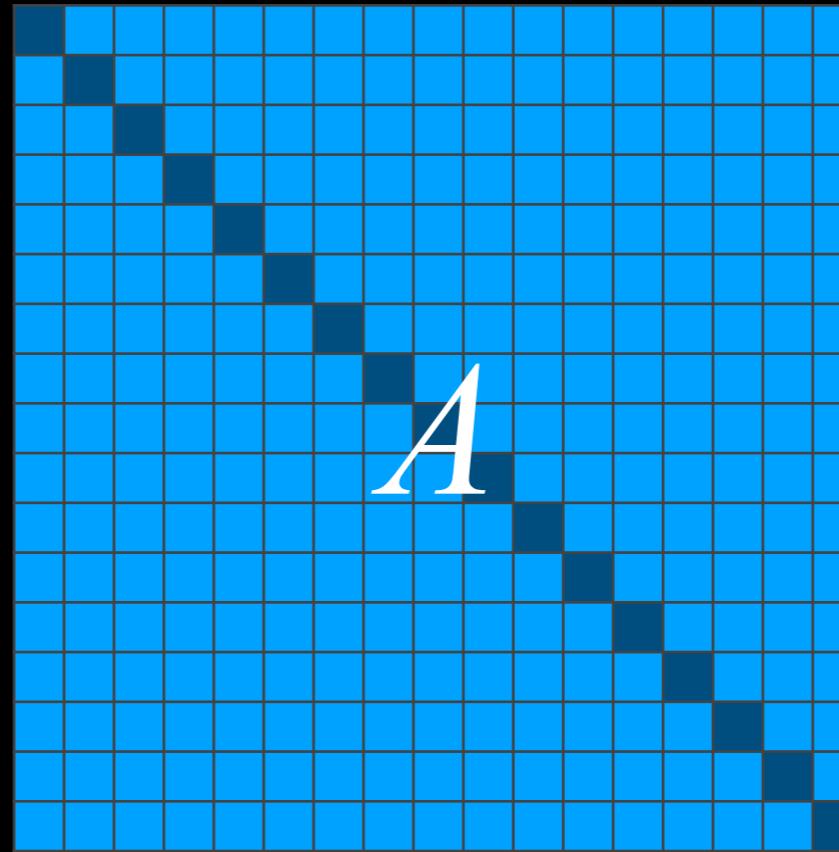


# Agglomerative Clustering

# Building Up

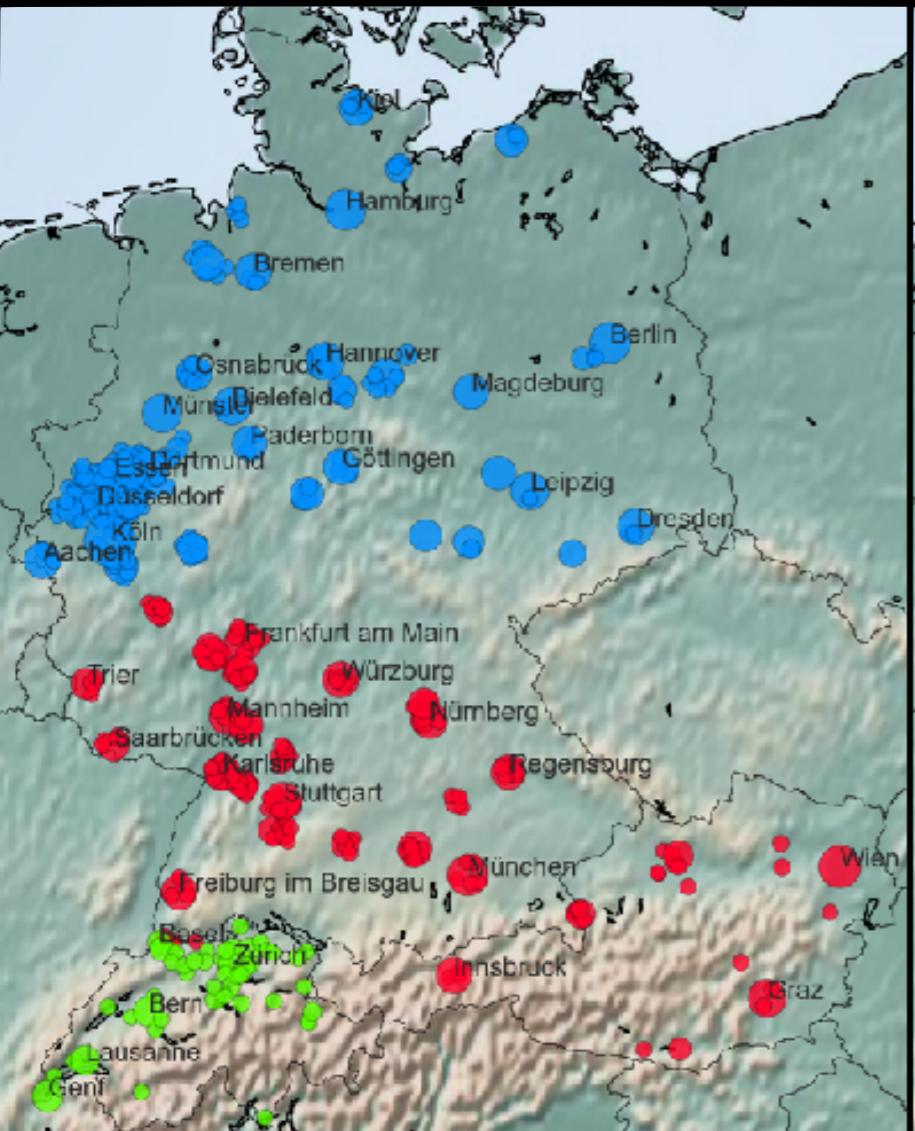


*ADJACENCY  
MATRIX*

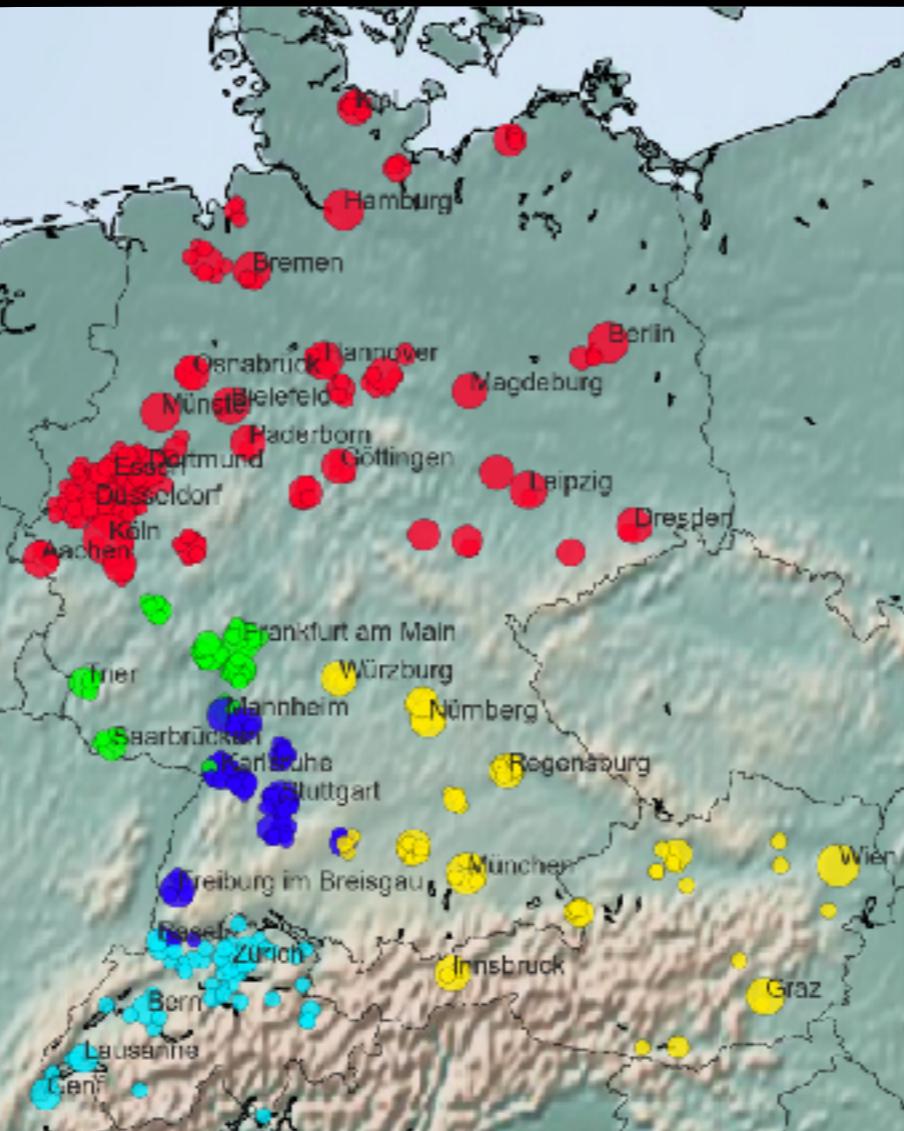


# Dialect Clusters

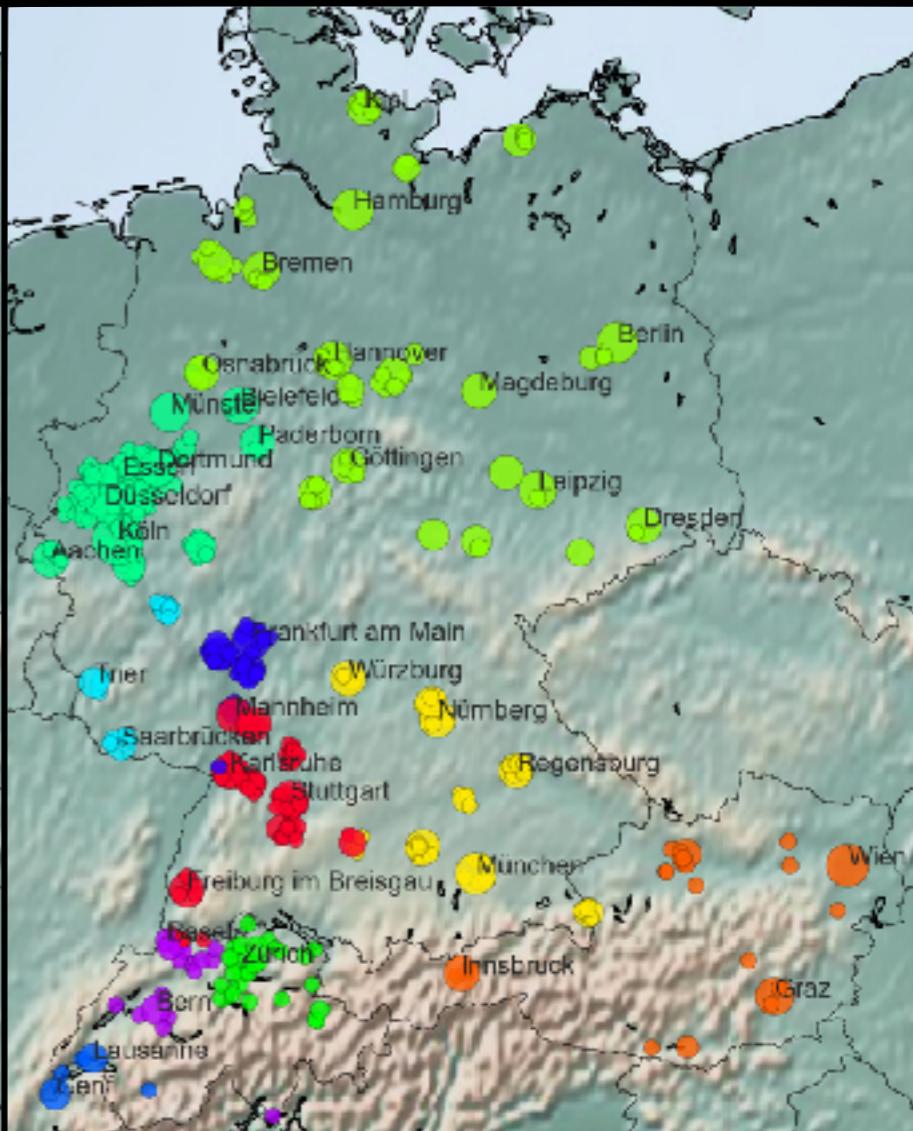
3



5

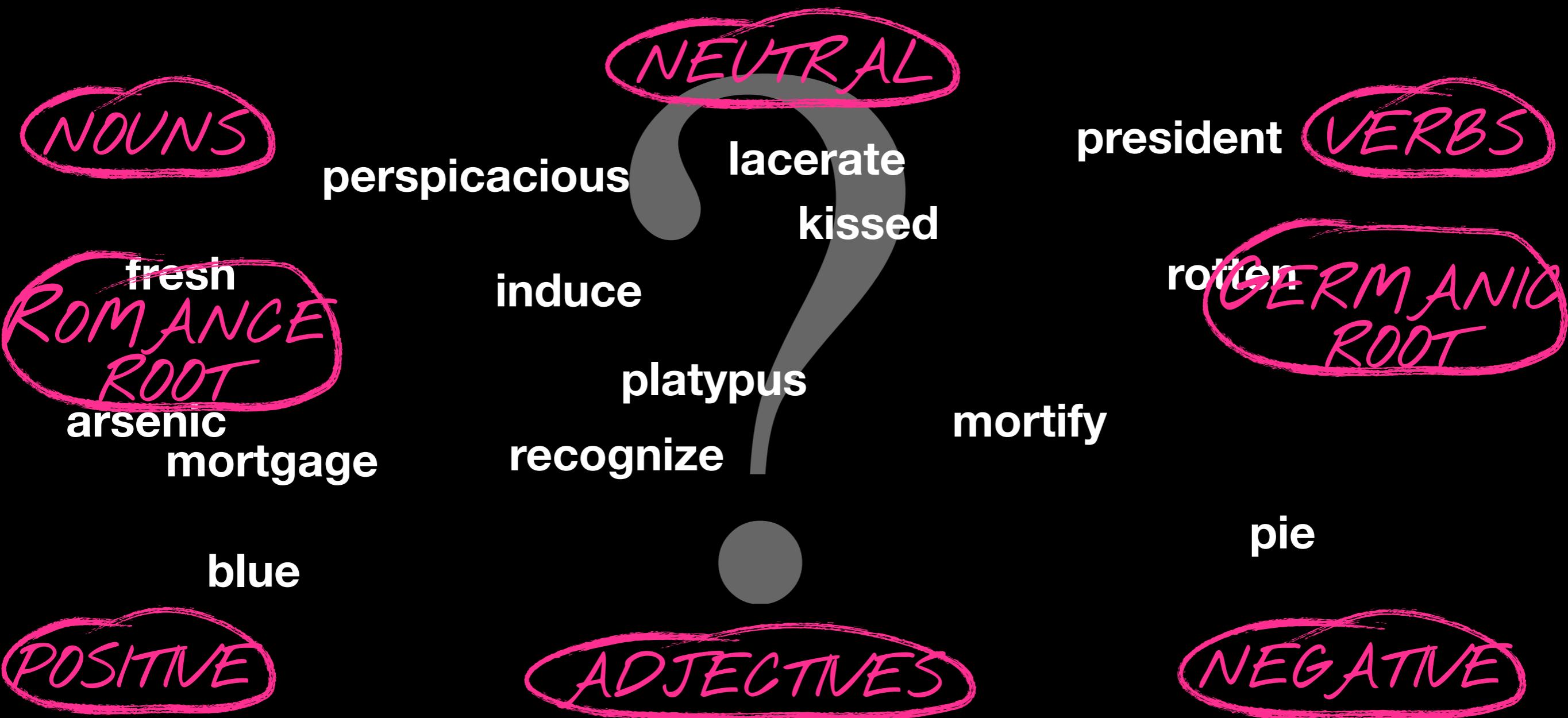


10



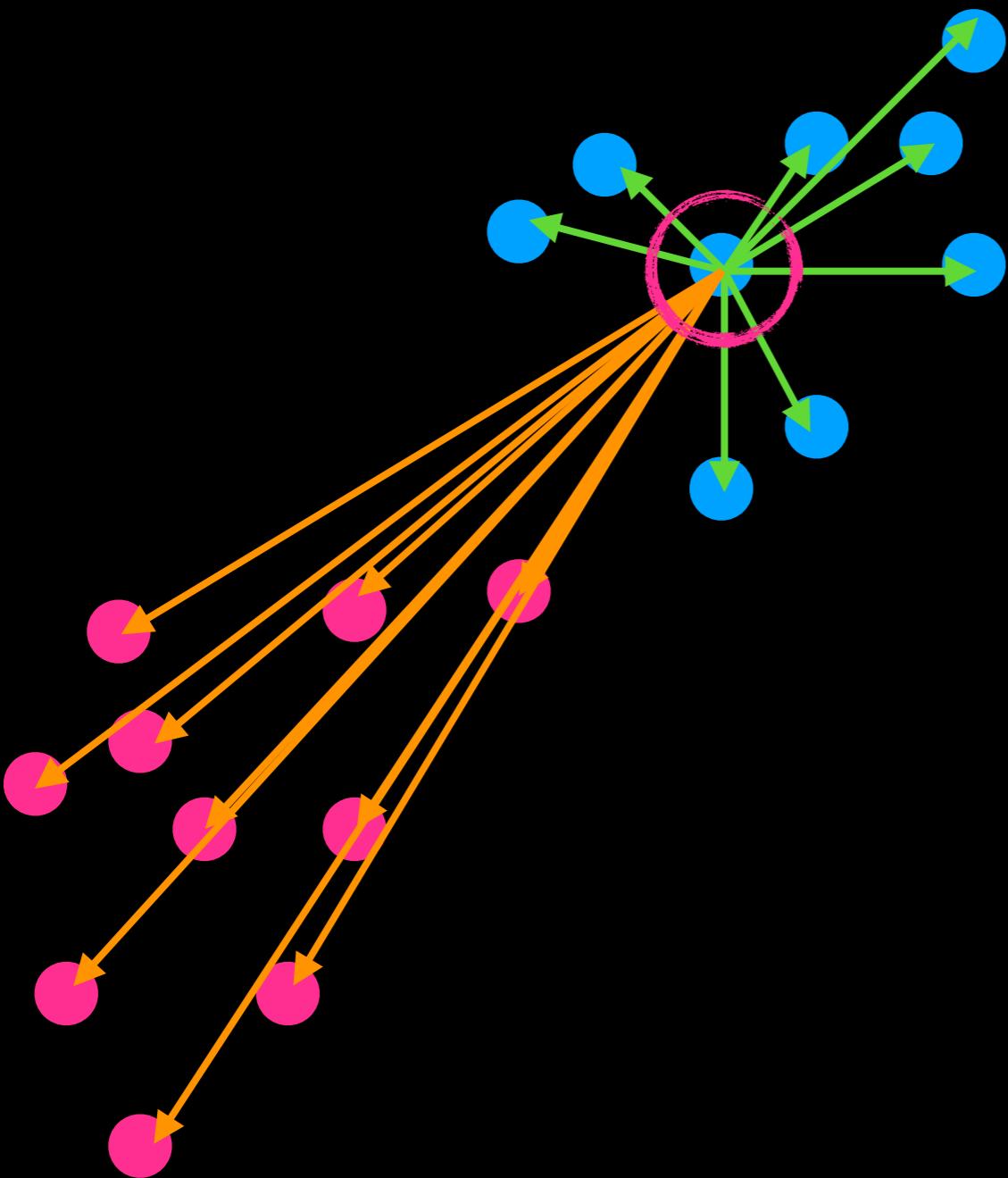
# Evaluating Clusters

# Making Sense of Clusters



# How Many Clusters?

Silhouette Score



$a = \text{mean intra-cluster distance}$

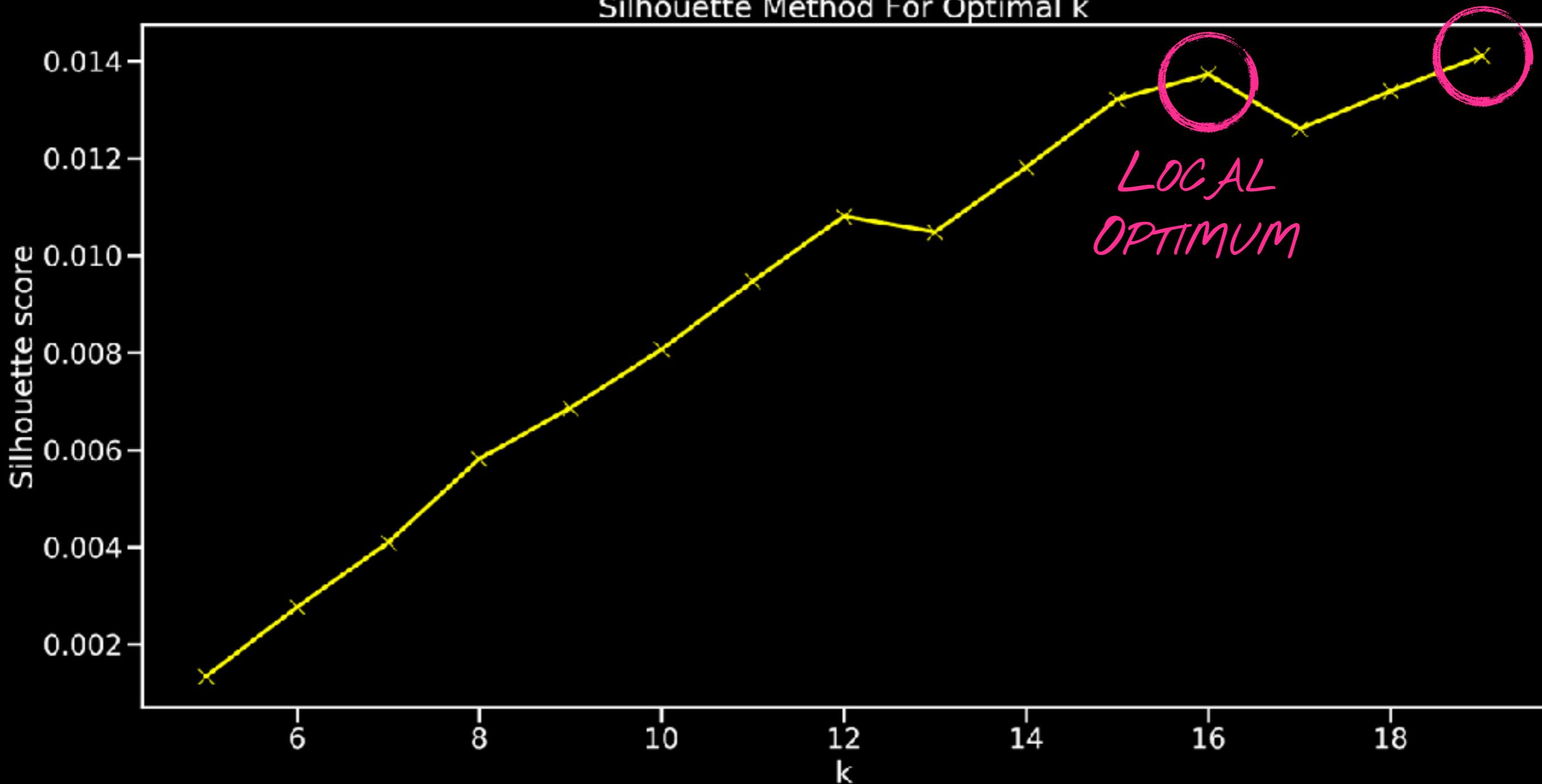
$$S = \frac{(b - a)}{\max(a, b)}$$

$b = \text{mean dist. nearest cluster}$

# Silhouette Scores

*DEPENDS ON PATIENCE/COMPUTE POWER*

Silhouette Method For Optimal k



# Supervised Evaluation Metrics

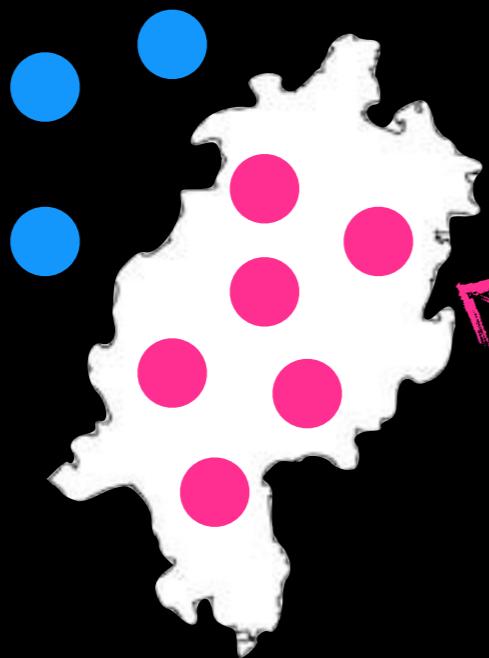
## Homogeneity

cluster has only 1 gold label

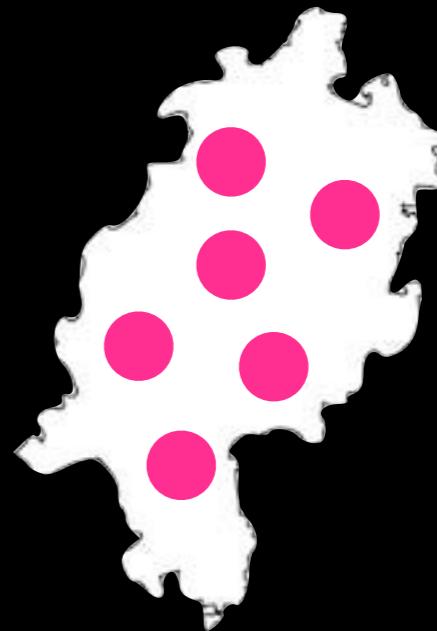
## Completeness

gold label has only 1 cluster

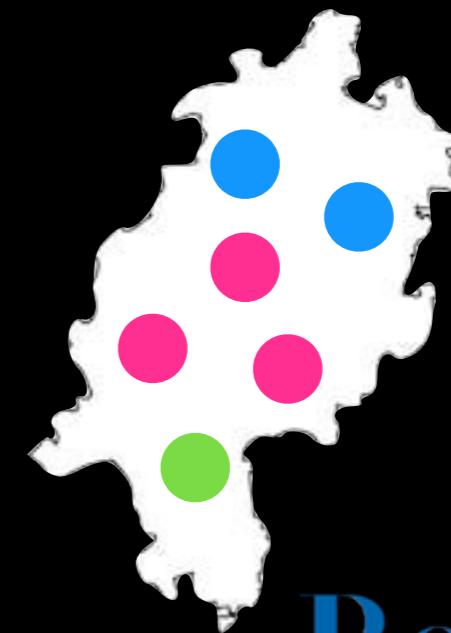
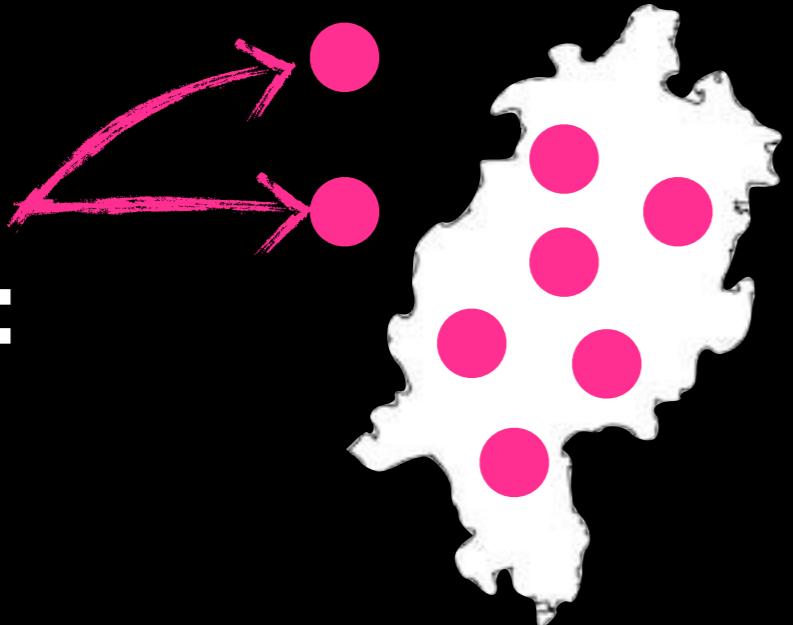
Good:



GOLD LABEL  
(REGION)



Bad:



# Comparison

	<b><i>k-means</i></b>	<b>Agg</b>
<b>scalable</b>	yes	no (up to ~20k)
<b>repeatable result</b>	no	yes
<b>include external info</b>	no	yes
<b>Good on dense clusters?</b>	no	yes

# Wrapping Up

# When to Use What

	Discrete Features	Embeddings
Latent topics	NMF	<i>Not applicable</i>
RGB translation	NMF	SVD + scaling
Plotting	SVD	t-SNE
Clustering	<i>Reduce dimensions</i>	<i>Use as-is</i>

# Take-Home Points

- **Matrix factorization** assumes latent concept dimensions
  - Can be used for semantic similarity (**LSA**)
  - Reduced components can be **visualized** in **graphs** or as **RGB** colors
- **Clusters** can group input in new ways
- Trade-off between speed and interpretability