



Nuts and bolts of neighbour embeddings

Dmitry Kobak
28 November 2025



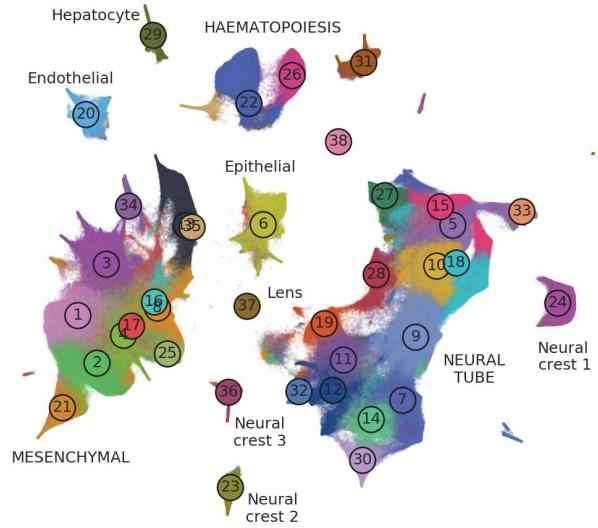
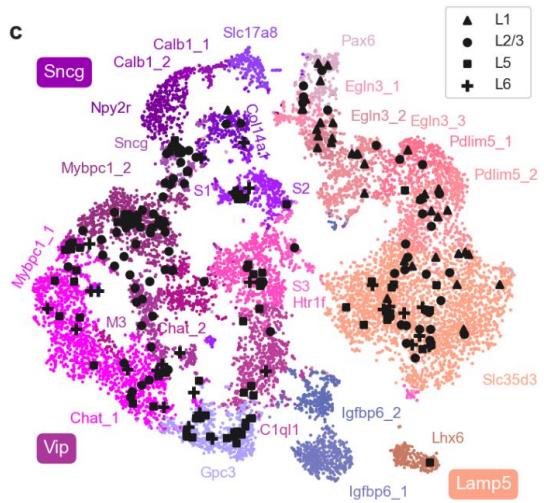
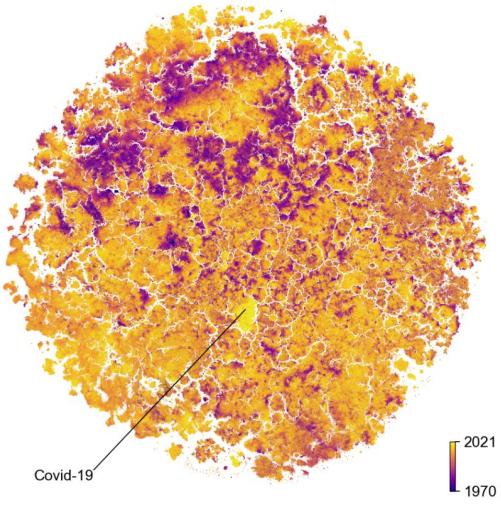
Everything you always wanted to know about t -SNE but were afraid to ask

Dmitry Kobak
28 November 2025



*Much more than you ever
Everything you ~~always~~ wanted to know
about t -SNE but were afraid to ask*

Dmitry Kobak
28 November 2025





Neighbor embeddings

*identify similar objects and
embed them close to each other*

Neighbor embeddings

identify similar objects and embed them close to each other

[PDF] Stochastic neighbor embedding

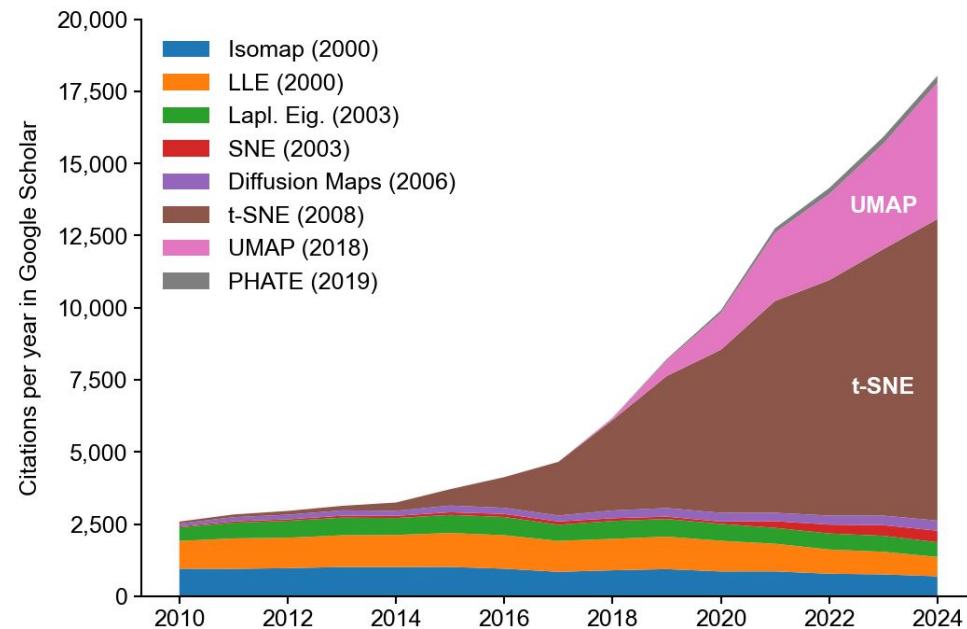
G Hinton, ST Roweis - NIPS, 2002 - Citeseer

We describe a probabilistic approach to the task of placing objects, described by high-dimensional vectors or by pairwise dissimilarities, in a low-dimensional space in a way that preserves neighbor identities. A Gaussian is centered on each object in the high ...

[PDF] Visualizing data using t-SNE.

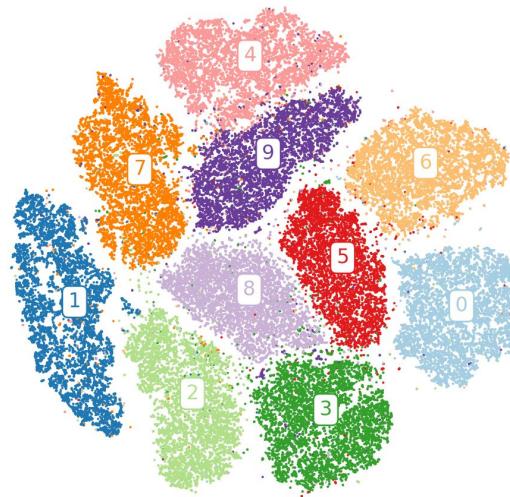
L. Van der Maaten, G. Hinton - Journal of machine learning research, 2008 - jmlr.org

We present a new technique called “**t-SNE**” that visualizes high-dimensional data by giving each datapoint a location in a two or three-dimensional map. The technique is a variation of Stochastic Neighbor Embedding (Hinton and Roweis, 2002) that is much easier to optimize ...

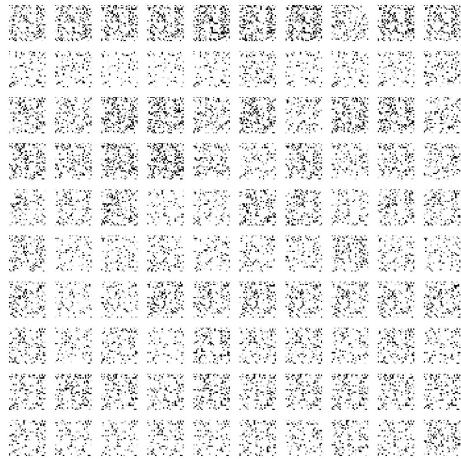


0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9

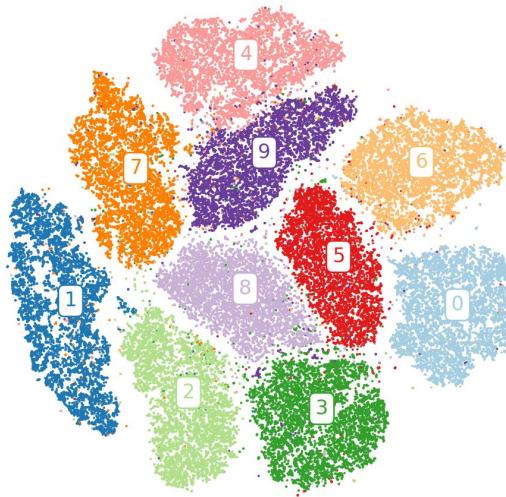
MNIST



t -SNE

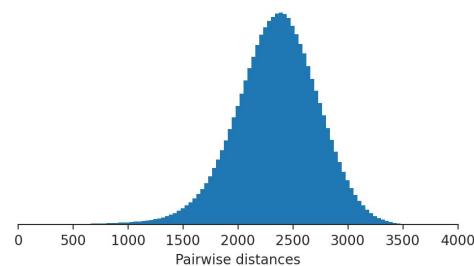
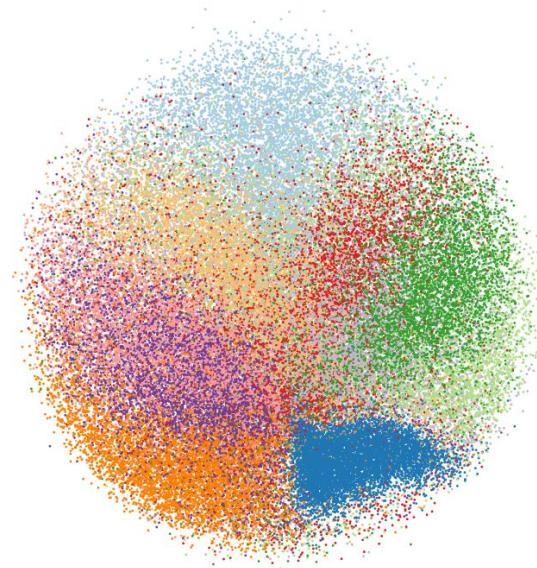
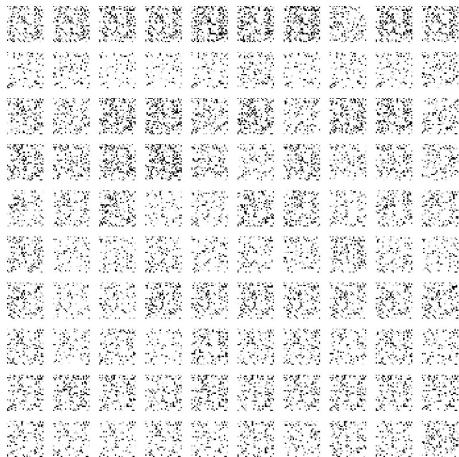


MNIST



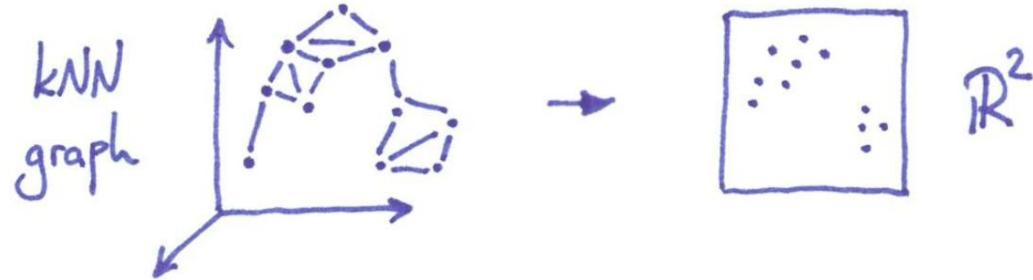
t -SNE

$$\mathcal{L} = \sum_{i < j} (d_{ij} - \|\mathbf{y}_i - \mathbf{y}_j\|)^2$$

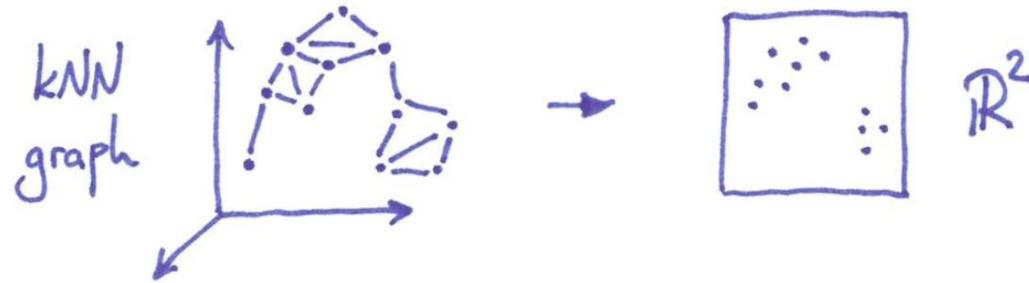


MDS

Neighbor
embeddings



t -SNE

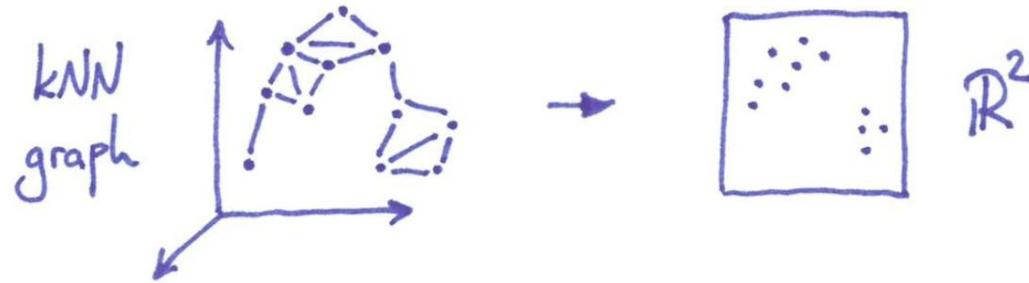


$$\mathcal{L}_{t\text{-SNE}} = \sum_{ij} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$

p_{ij} can be normalized symmetrized adjacency matrix of the k NN graph.

$$q_{ij} = \frac{w_{ij}}{Z} = \frac{w_{ij}}{\sum w_{ij}} \text{ with } w_{ij} = \frac{1}{1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2}.$$

t -SNE



$$\mathcal{L}_{t\text{-SNE}} = \sum_{ij} p_{ij} \log \frac{p_{ij}}{q_{ij}} \sim - \sum_{ij} p_{ij} \log q_{ij} = - \sum_{ij} p_{ij} \log \frac{w_{ij}}{Z} = - \underbrace{\sum_{ij} p_{ij} \log w_{ij}}_{\text{attraction}} + \underbrace{\log \sum_{ij} w_{ij}}_{\text{repulsion}}.$$

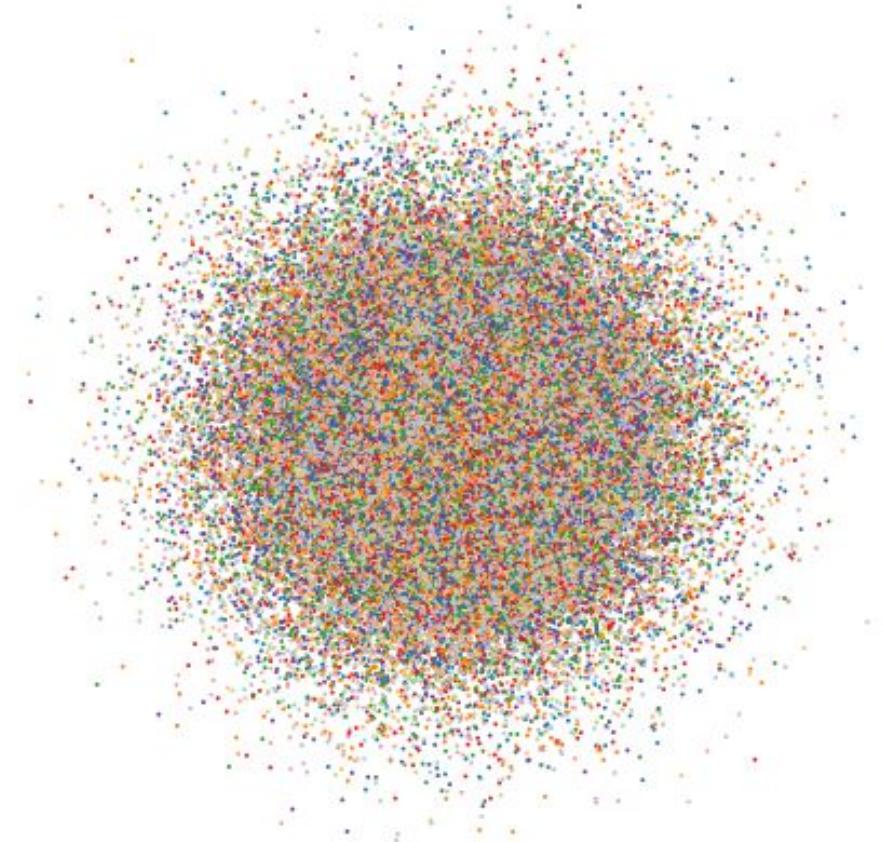
p_{ij} can be normalized symmetrized adjacency matrix of the k NN graph.

$$q_{ij} = \frac{w_{ij}}{Z} = \frac{w_{ij}}{\sum w_{ij}} \text{ with } w_{ij} = \frac{1}{1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2}.$$

Gradient descent

Gradient descent starting from a random configuration.

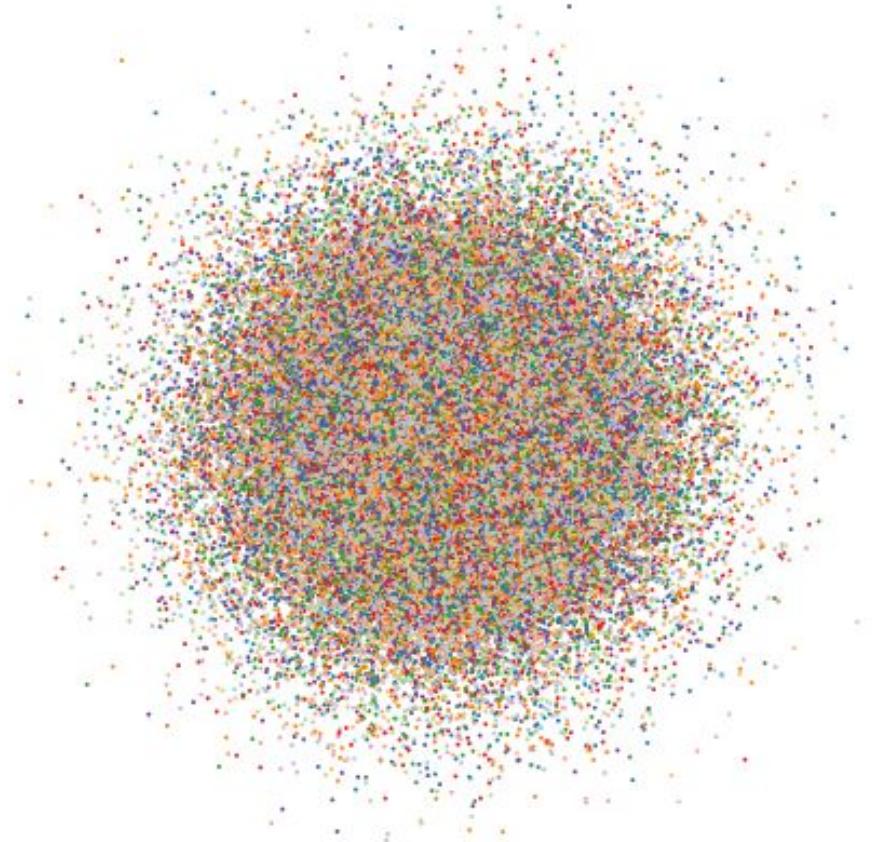
Each frame is scaled
(in reality embedding is initialized small and slowly grows in size).



Gradient descent

Optimization issues:

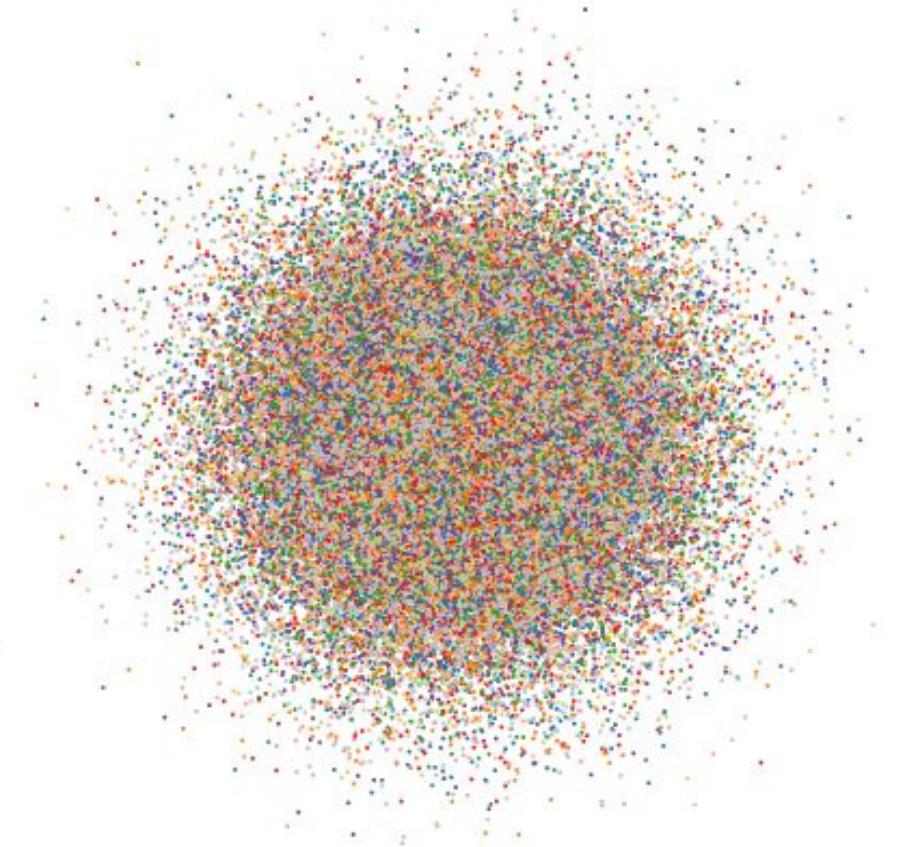
- 1) Initialization
- 2) Learning rate
- 3) Local minima
- 4) Runtime



Early exaggeration

Multiply all attractive forces
by 12 for 250 iterations.

$$-\underbrace{\sum_{ij} p_{ij} \log w_{ij}}_{\text{attraction}} + \underbrace{\log \sum_{ij} w_{ij}}_{\text{repulsion}}$$



Early exaggeration

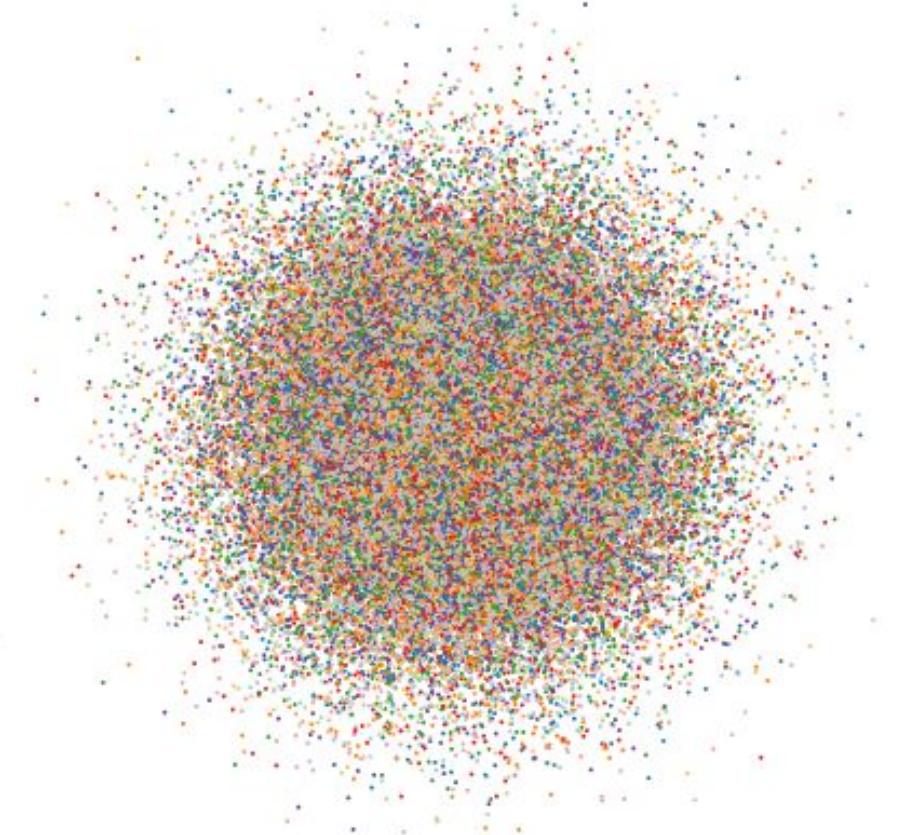
Multiply all attractive forces
by 12 for 250 iterations.

$$-\underbrace{\sum_{ij} p_{ij} \log w_{ij}}_{\text{attraction}} + \underbrace{\log \sum_{ij} w_{ij}}_{\text{repulsion}}$$

The learning rate should grow
with sample size:

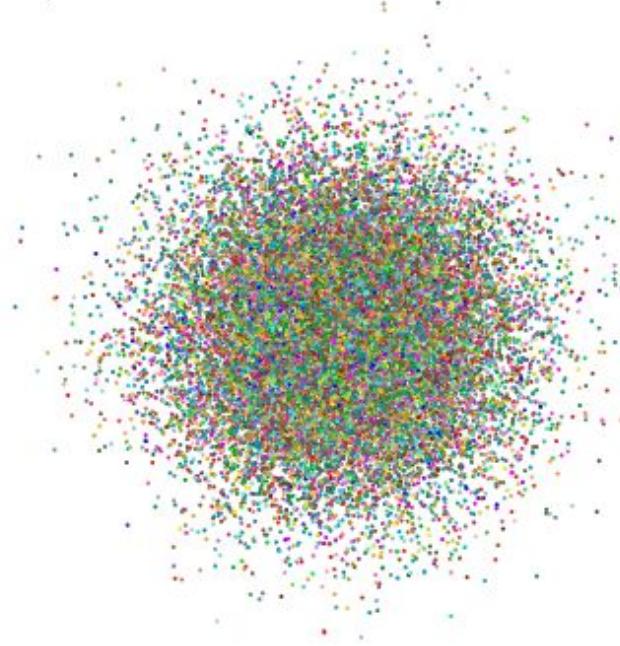
$$\eta = n / \text{exagg.}$$

*Belkina et al. (2019) based on
Linderman & Steinerberger (2019)*

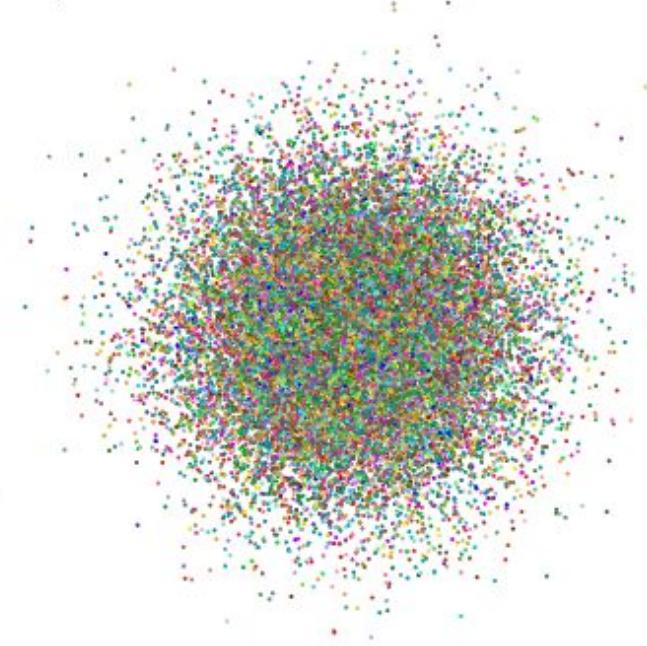


Early exaggeration: scRNA-seq data

Iter 1, KL=6.52



Iter 1, KL=6.52

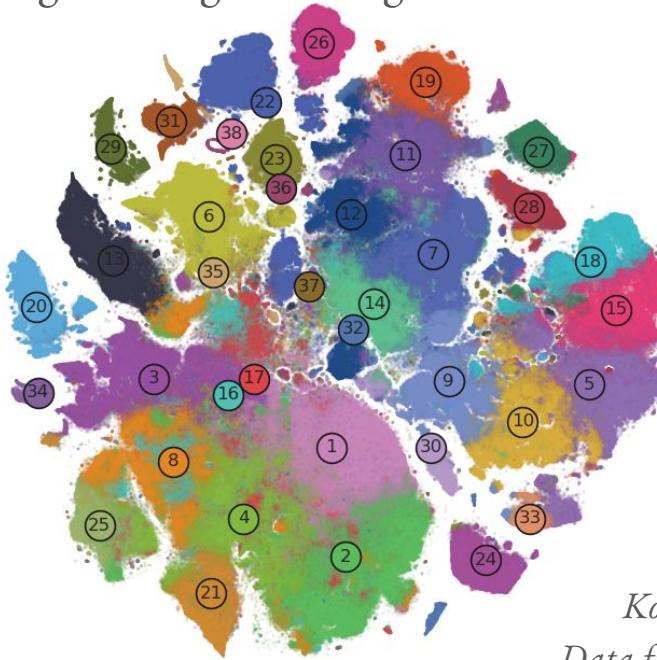


Data from Tasic et al. (2018)

Early exaggeration and the learning rate: real-world example

Single-cell transcriptomic study of mouse embryogenesis ($n \approx 2,000,000$).

Left — t-SNE published in the original paper. Right — high learning rate.

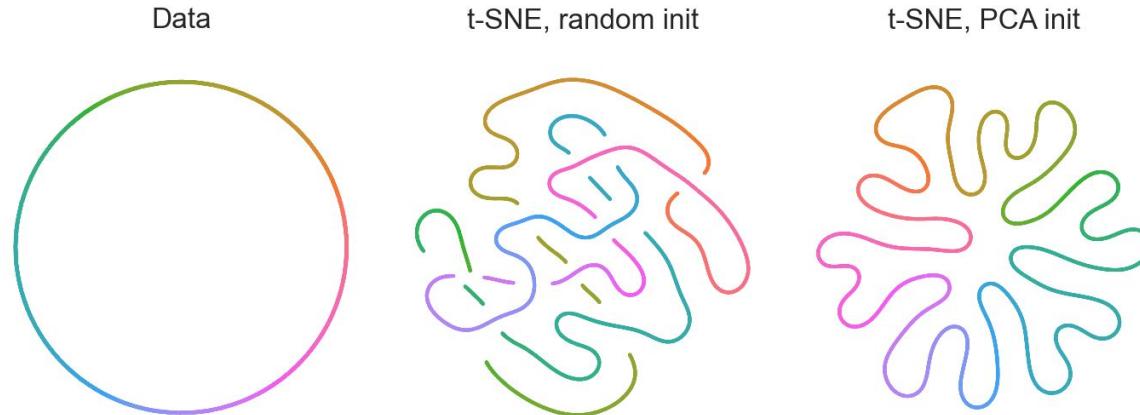


Kobak & Berens (2019)

Data from Cao et al. (2019)

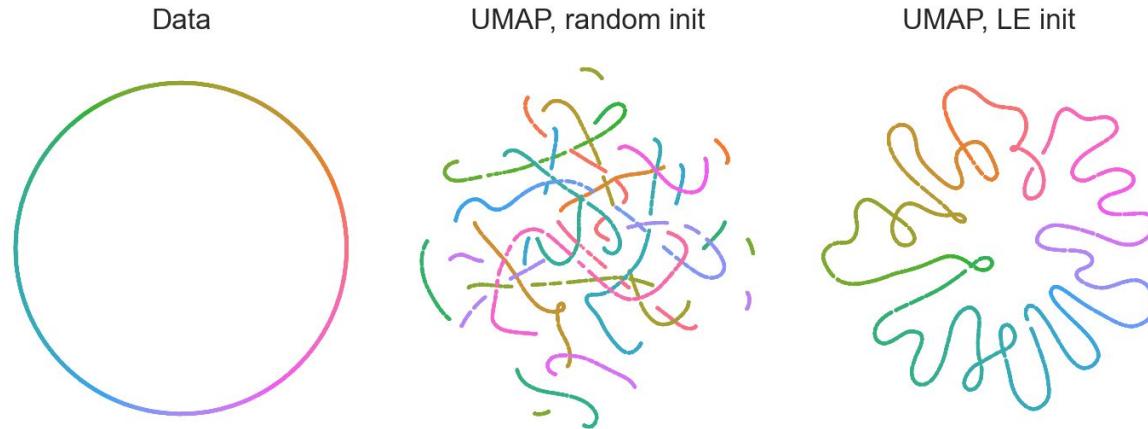
Initialization

The loss function has many local minima. Even with early exaggeration, initialization can play a big role. With random initialization, t-SNE will often struggle to preserve global structure.



Initialization

The loss function has many local minima. Even with early exaggeration, initialization can play a big role. With random initialization, t-SNE will often struggle to preserve global structure.



Initialization: real-world example

t-SNE (random initialization)



PCA

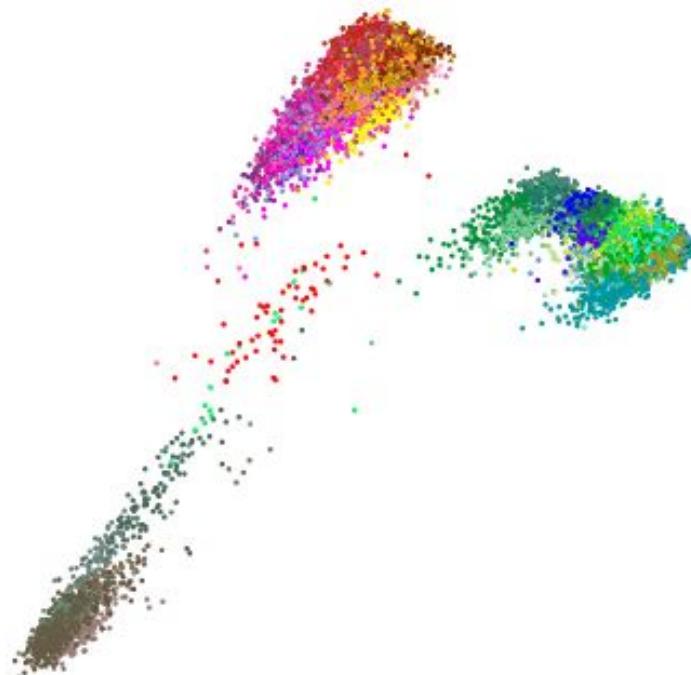


Kobak & Berens, *Nature Comms* 2019)

Data from Tasic et al. (2018)

Single-cell RNA-seq data from mouse cortex

Initialization: real-world example



Single-cell RNA-seq data from mouse cortex

Kobak & Berens, *Nature Comms* 2019)

Data from Tasic et al. (2018)

Approximations

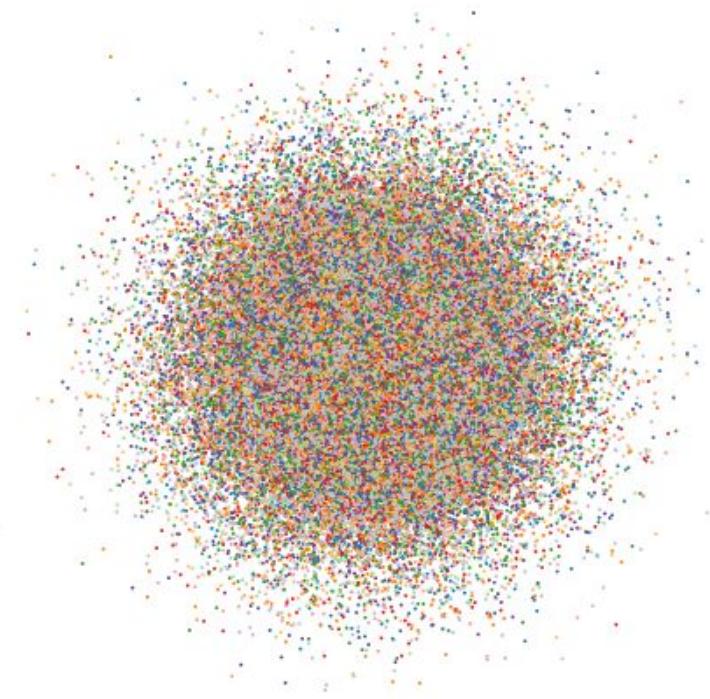
$$-\sum_{ij} p_{ij} \log w_{ij} + \log \sum_{ij} w_{ij}$$

$\underbrace{}_{\text{attraction}}$ $\underbrace{}_{\text{repulsion}}$

Approximate k NN graph
with small k

Approximations:

- * Barnes–Hut
- * Fourier interpolation
- * Sampling



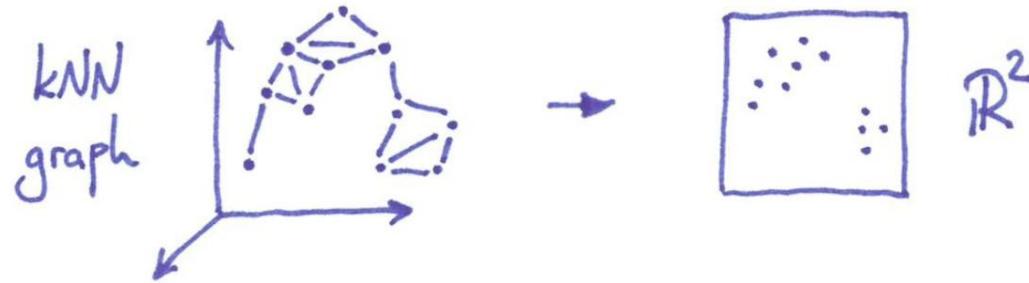
Optimization issues

A good t-SNE implementation should:

1. use informative initialization by default;
2. set appropriate optimization parameters (learning rate set to n);
3. support approximate nearest neighbors and $O(n)$ approximation of repulsive forces.

=> In Python, use openTSNE (*Poličar et al., 2024*).

t -SNE

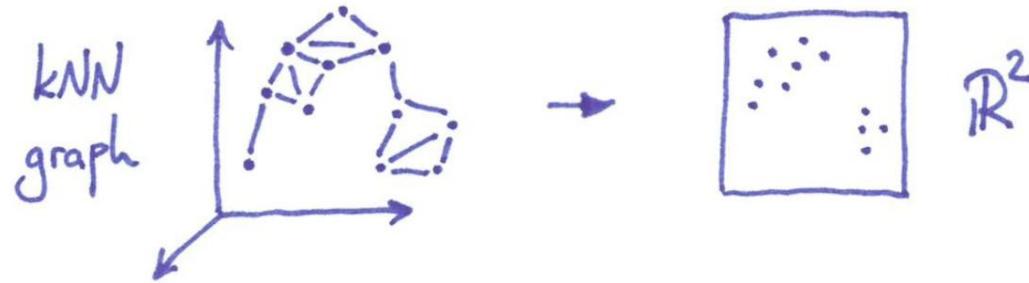


$$\mathcal{L}_{t\text{-SNE}} = \sum_{ij} p_{ij} \log \frac{p_{ij}}{q_{ij}} \sim - \sum_{ij} p_{ij} \log q_{ij} = - \sum_{ij} p_{ij} \log \frac{w_{ij}}{Z} = - \underbrace{\sum_{ij} p_{ij} \log w_{ij}}_{\text{attraction}} + \underbrace{\log \sum_{ij} w_{ij}}_{\text{repulsion}}.$$

p_{ij} can be normalized symmetrized adjacency matrix of the k NN graph.

$$q_{ij} = \frac{w_{ij}}{Z} = \frac{w_{ij}}{\sum w_{ij}} \text{ with } w_{ij} = \frac{1}{1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2}.$$

t-SNE



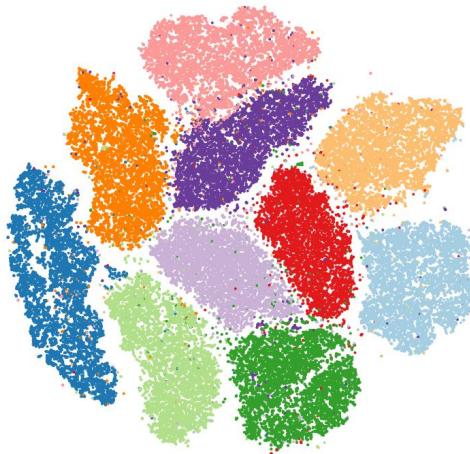
$$\mathcal{L}_{\text{t-SNE}} = \sum_{ij} p_{ij} \log \frac{p_{ij}}{q_{ij}} \sim - \sum_{ij} p_{ij} \log q_{ij} = - \sum_{ij} p_{ij} \log \frac{w_{ij}}{Z} = - \underbrace{\sum_{ij} p_{ij} \log w_{ij}}_{\text{attraction}} + \underbrace{\log \sum_{ij} w_{ij}}_{\text{repulsion}}.$$

p_{ij} can be normalized symmetrized adjacency matrix of the k NN graph.

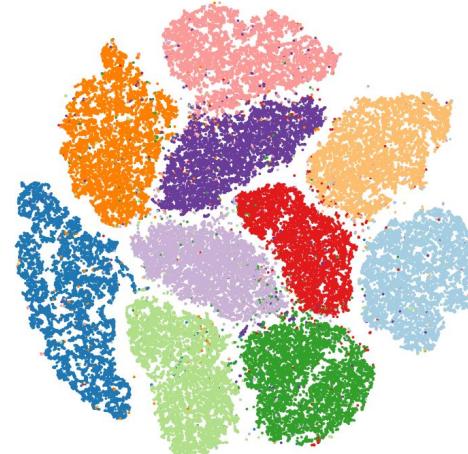
Usually high-dim affinities
are obtained with a
Gaussian kernel via
perplexity parameter

$$q_{ij} = \frac{w_{ij}}{Z} = \frac{w_{ij}}{\sum w_{ij}} \text{ with } w_{ij} = \frac{1}{1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2}.$$

t-SNE with perplexity 30



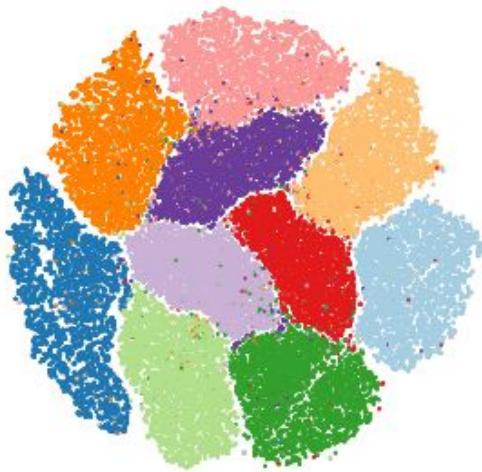
t-SNE with a kNN kernel (k=10)



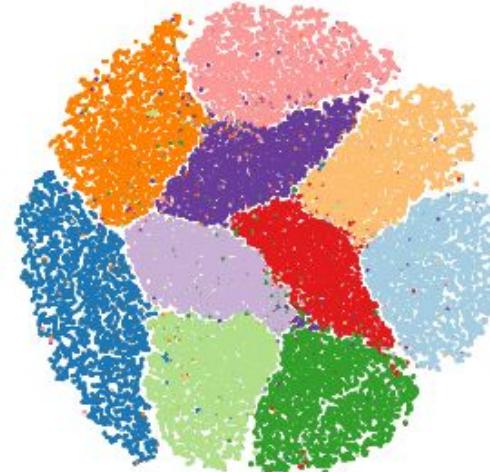
p_{ij} can be normalized symmetrized adjacency matrix of the k NN graph.

Usually high-dim affinities
are obtained with a
Gaussian kernel via
perplexity parameter

t-SNE with perplexity 9



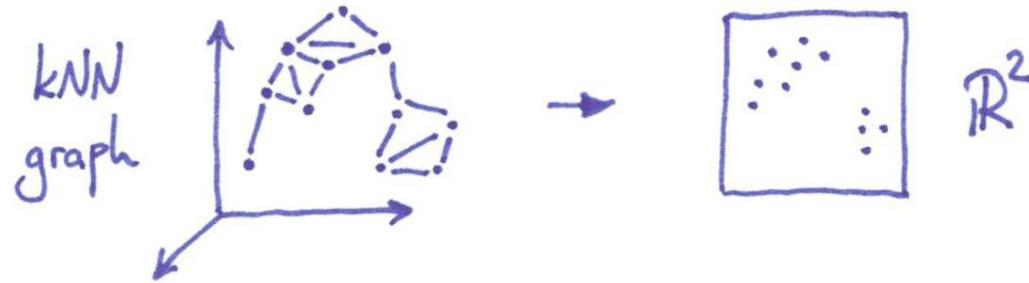
t-SNE with a kNN kernel, k=3



p_{ij} can be normalized symmetrized adjacency matrix of the k NN graph.

Usually high-dim affinities
are obtained with a
Gaussian kernel via
perplexity parameter

t -SNE

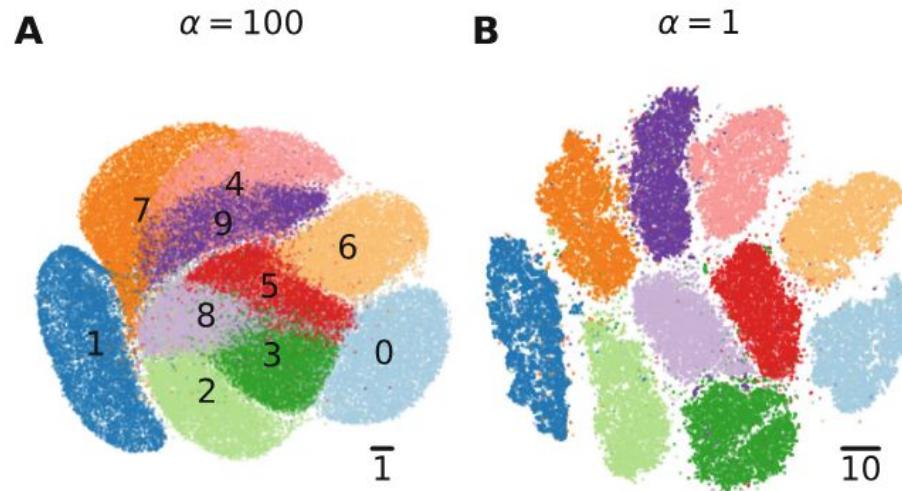


$$\mathcal{L}_{t\text{-SNE}} = \sum_{ij} p_{ij} \log \frac{p_{ij}}{q_{ij}} \sim - \sum_{ij} p_{ij} \log q_{ij} = - \sum_{ij} p_{ij} \log \frac{w_{ij}}{Z} = - \underbrace{\sum_{ij} p_{ij} \log w_{ij}}_{\text{attraction}} + \underbrace{\log \sum_{ij} w_{ij}}_{\text{repulsion}}.$$

p_{ij} can be normalized symmetrized adjacency matrix of the k NN graph.

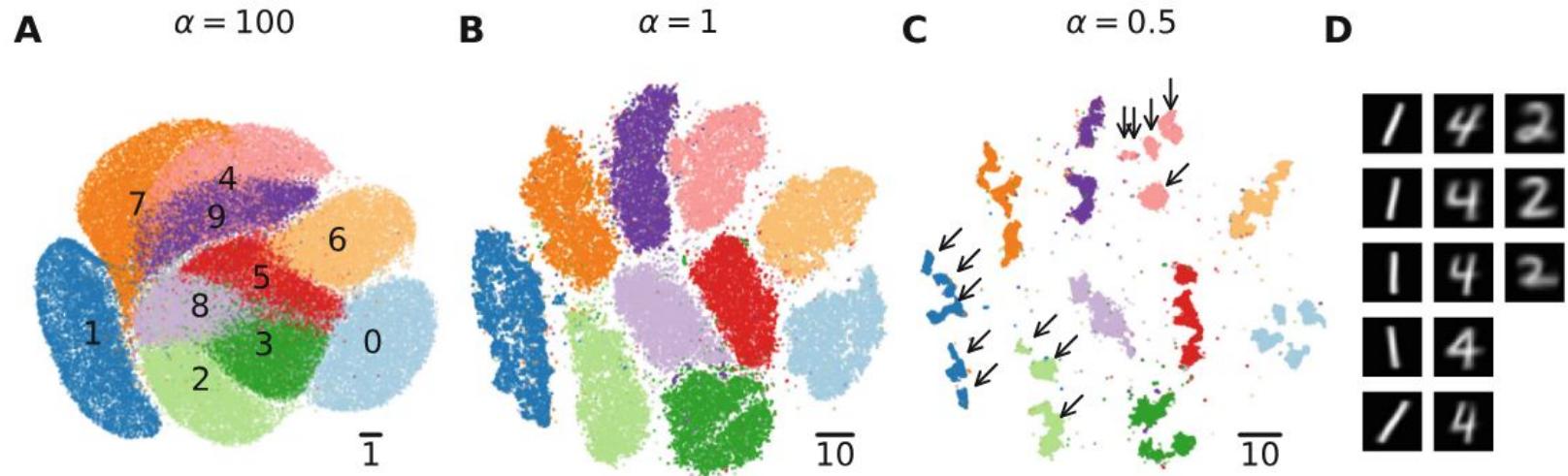
$$q_{ij} = \frac{w_{ij}}{Z} = \frac{w_{ij}}{\sum w_{ij}} \text{ with } w_{ij} = \frac{1}{1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2}.$$

The t -distribution kernel: from SNE to t -SNE



$$k(d) = \frac{1}{(1 + d^2/\alpha)^\alpha}$$

The t -distribution kernel: beyond t -SNE

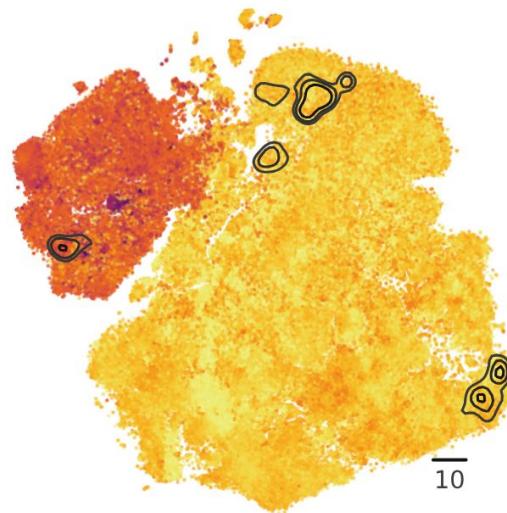


$$k(d) = \frac{1}{(1 + d^2/\alpha)^\alpha}$$

The t -distribution kernel: beyond t -SNE

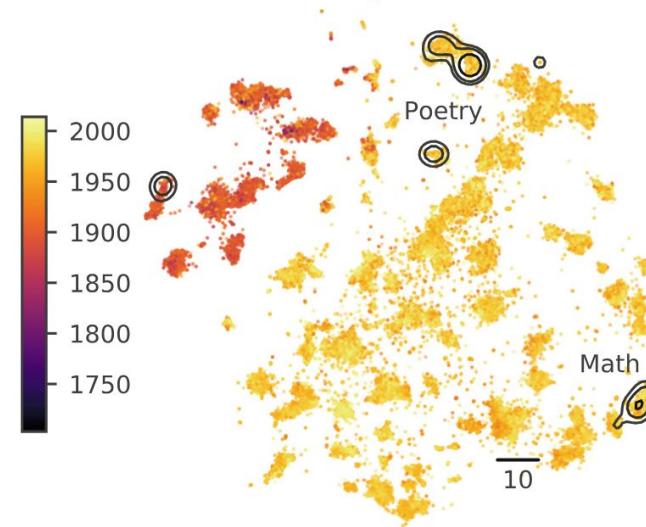
A

$\alpha = 1$



B

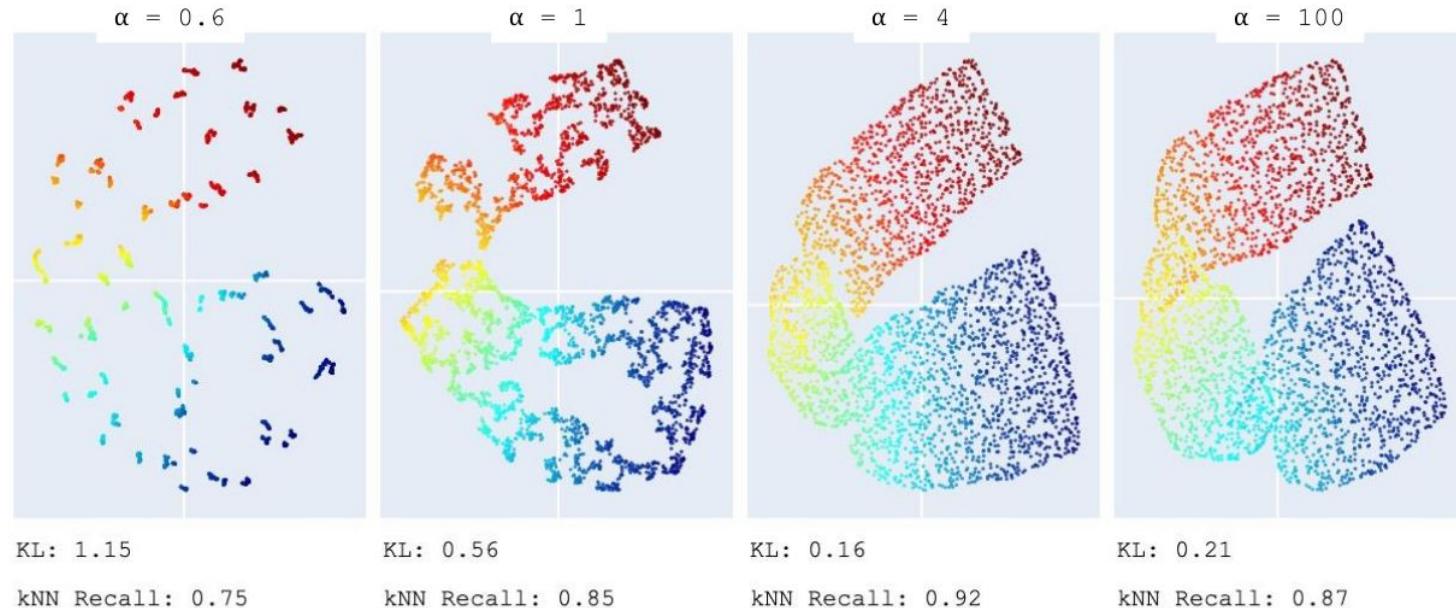
$\alpha = 0.5$



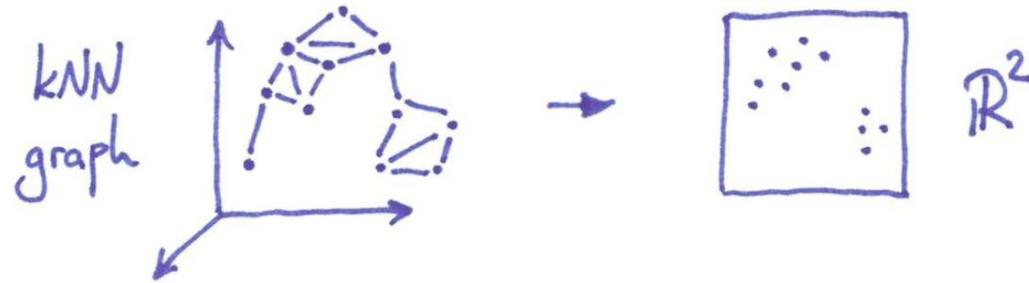
Russian language part of the HathiTrust library ($n = 408\,291$).

The t -distribution kernel: optimizing the α

$$\frac{\partial L}{\partial \alpha} = \sum_{i,j} (p_{i,j} - q_{i,j}) \left(\log \left(1 + \frac{d_{i,j}^2}{\alpha} \right) - \frac{d_{i,j}^2}{\alpha + d_{i,j}^2} \right)$$



t -SNE



$$\mathcal{L}_{t\text{-SNE}} = \sum_{ij} p_{ij} \log \frac{p_{ij}}{q_{ij}} \sim - \sum_{ij} p_{ij} \log q_{ij} = - \sum_{ij} p_{ij} \log \frac{w_{ij}}{Z} = - \underbrace{\sum_{ij} p_{ij} \log w_{ij}}_{\text{attraction}} + \underbrace{\log \sum_{ij} w_{ij}}_{\text{repulsion}}.$$

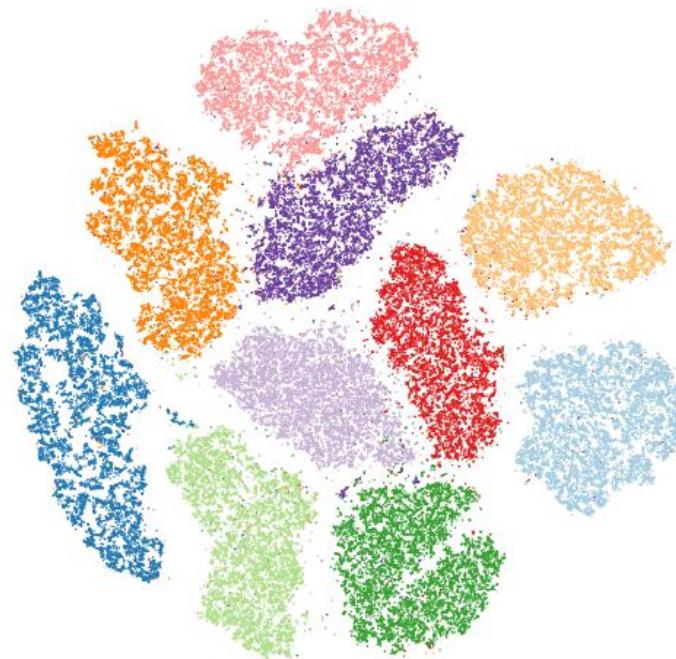
p_{ij} can be normalized symmetrized adjacency matrix of the k NN graph.

$$q_{ij} = \frac{w_{ij}}{Z} = \frac{w_{ij}}{\sum w_{ij}} \text{ with } w_{ij} = \frac{1}{1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2}.$$

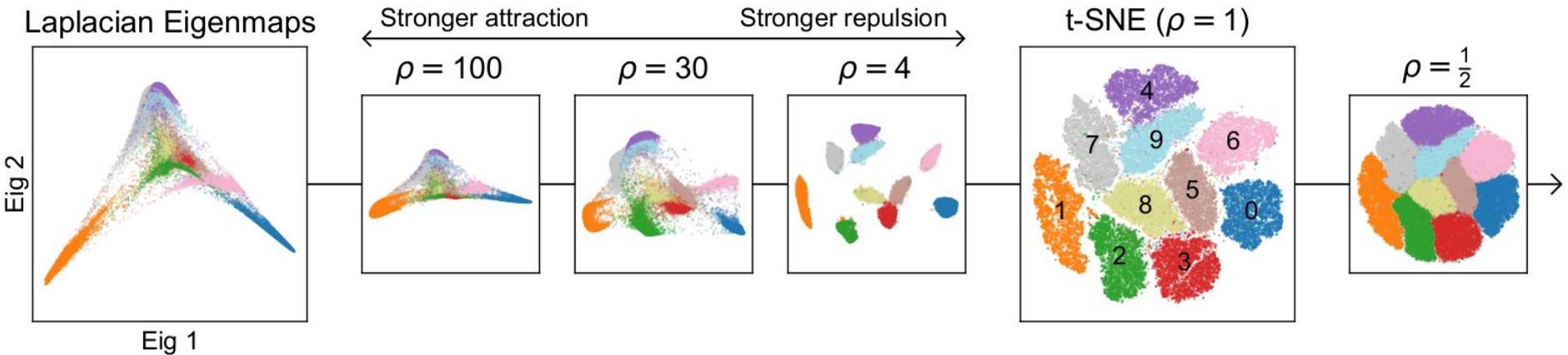


Attraction-repulsion spectrum

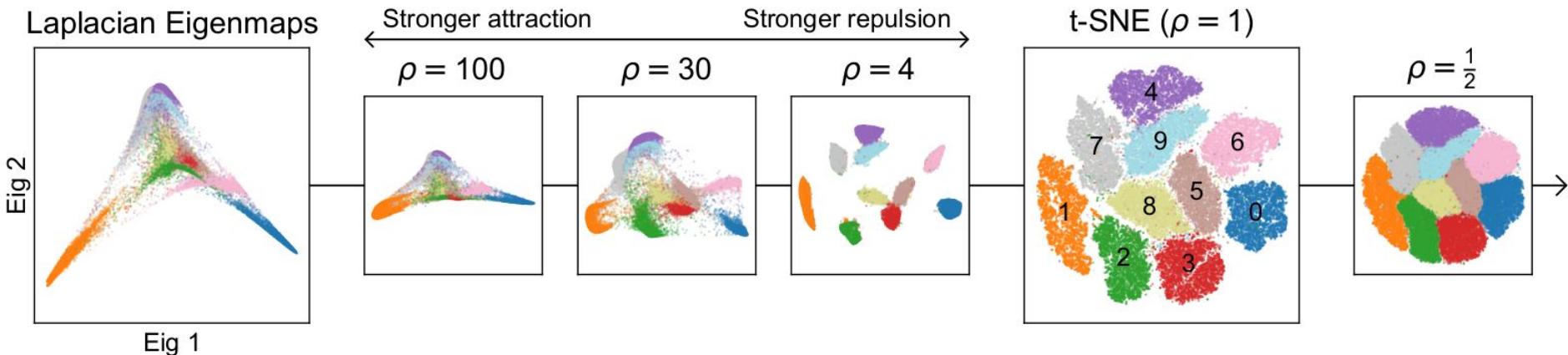
Niklas Böhm



Attraction-repulsion spectrum

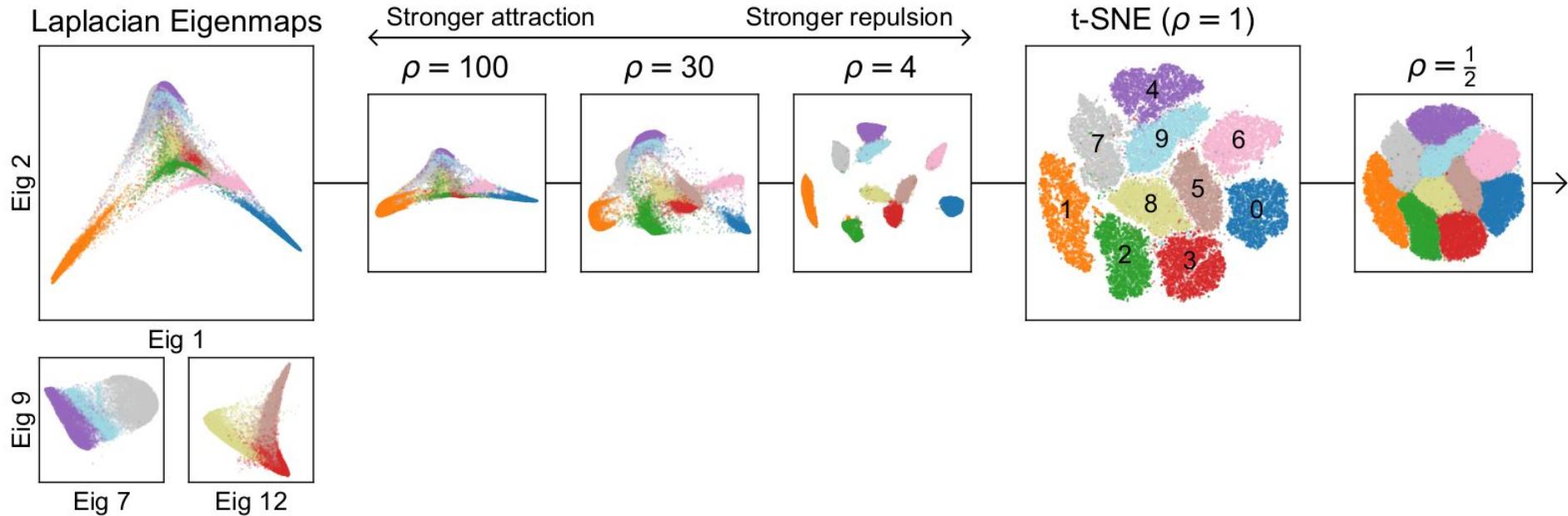


Attraction-repulsion spectrum



For very strong exaggeration, t-SNE approximates Laplacian Eigenmaps of the affinity matrix (*Linderman & Steinerberger, 2020*).

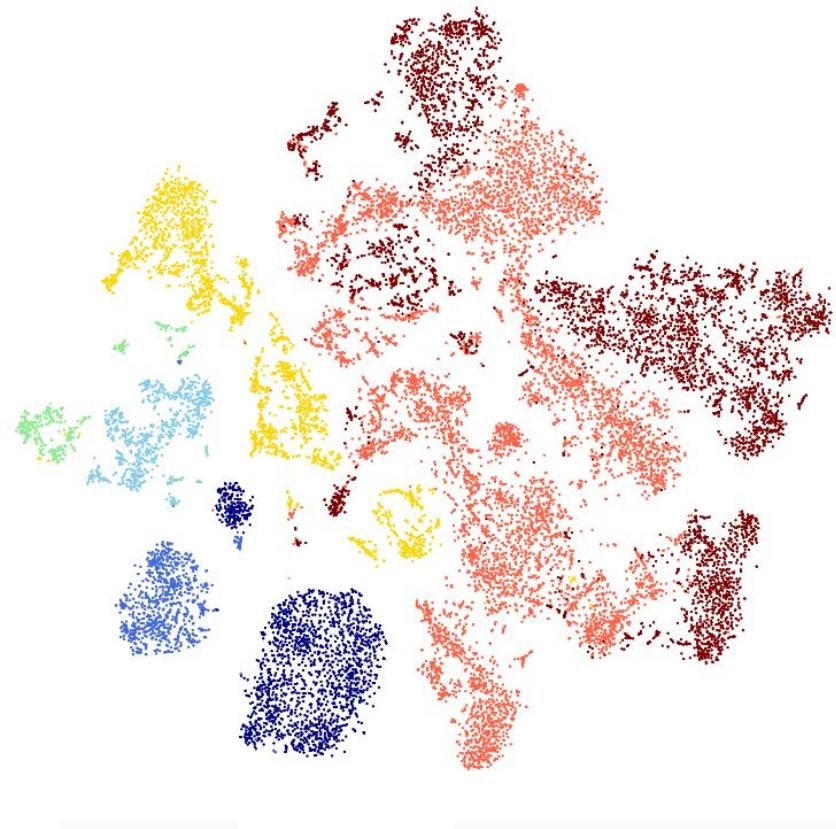
Attraction-repulsion spectrum



The ‘data-agnostic’ repulsion term brings our information from higher Laplacian eigenvectors.



Attraction-repulsion spectrum



Brain organoid scRNA-seq data
from *Kanton et al., Nature 2019*

Böhm et al., JMLR 2022

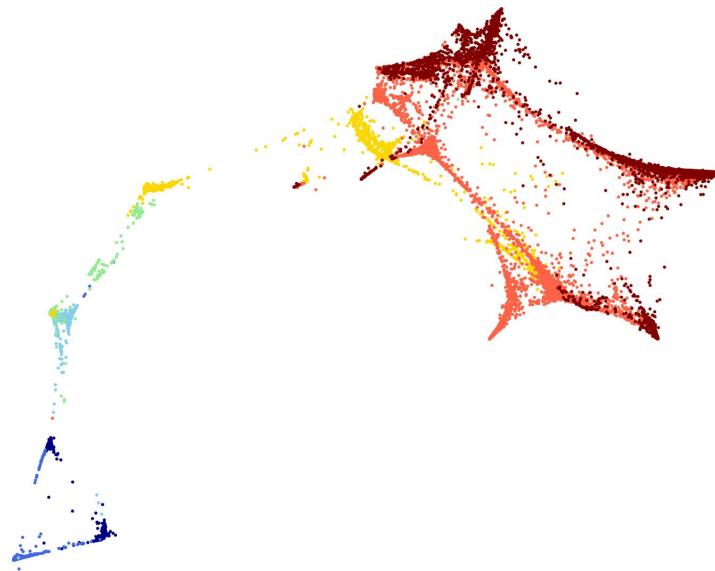
Niklas Böhm



Attraction-repulsion spectrum

Niklas Böhm

- 0 days
- 4 days
- 10 days
- 15 days
- 1 month
- 2 months
- 4 months

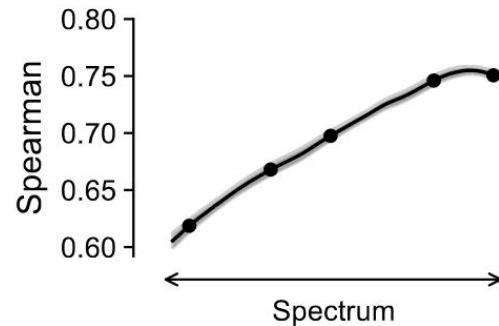
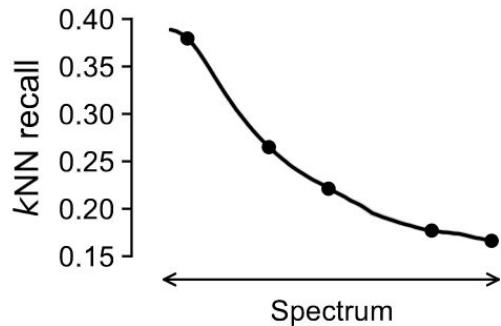


Brain organoid scRNA-seq data
from Kanton et al., *Nature* 2019

Böhm et al., *JMLR* 2022



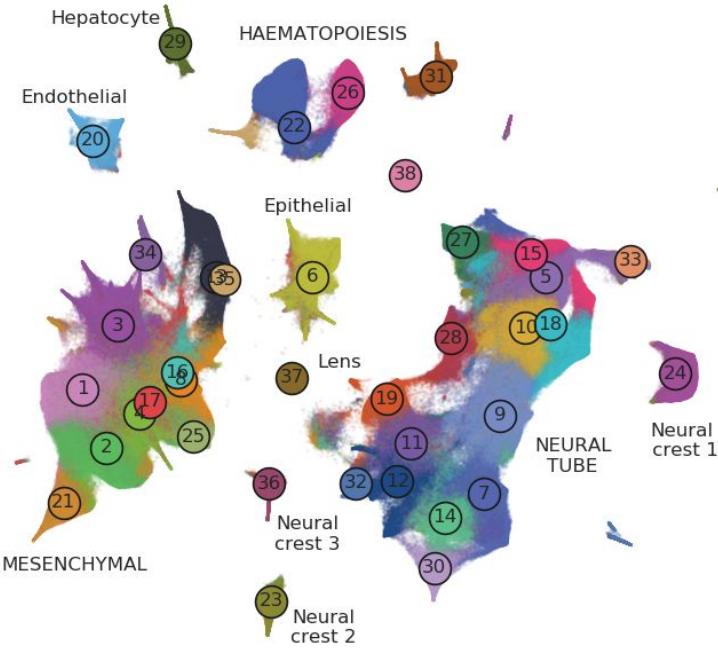
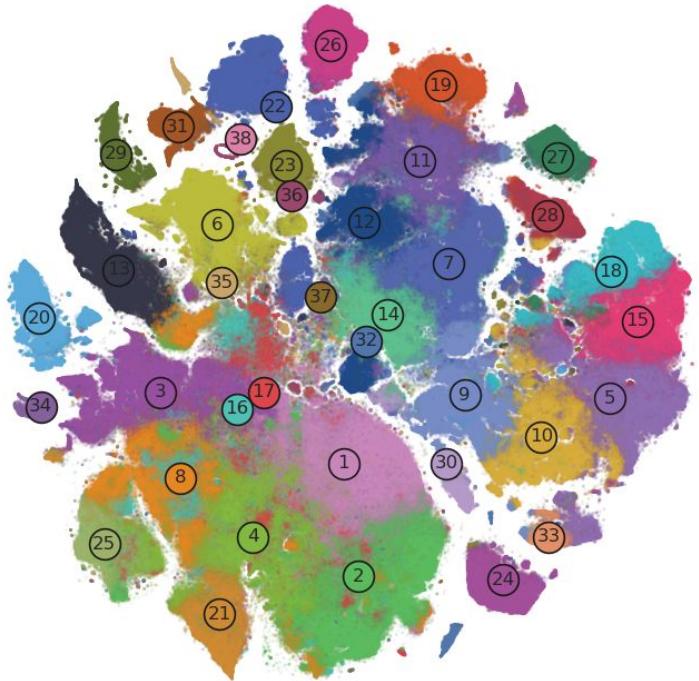
Local-global trade-off



Sebastian Damrich



Local-global trade-off ($n = 2$ mln)

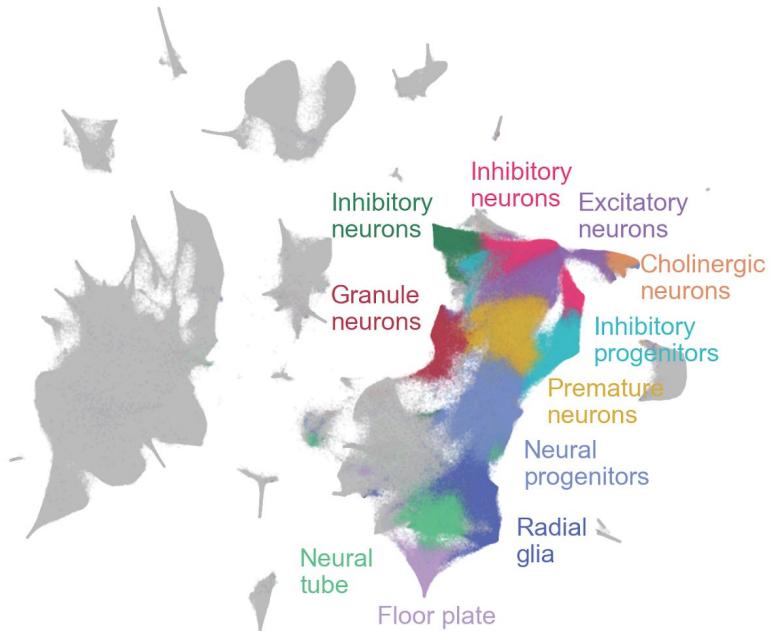
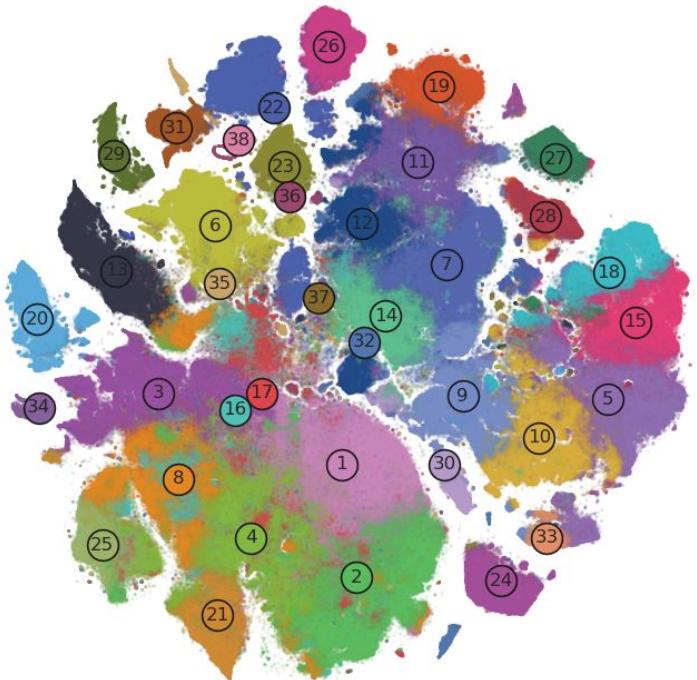


Data from

Cao et al., *Nature* 2019

Kobak & Berens, *Nature Comms* 2019

Local-global trade-off ($n = 2$ mln)



Data from

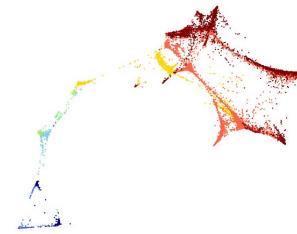
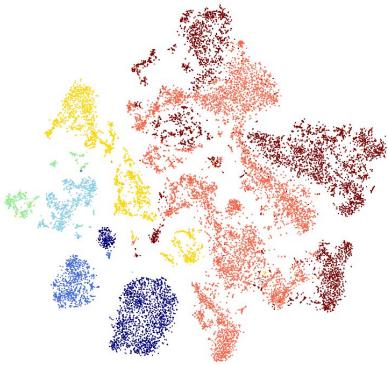
Cao et al., *Nature* 2019

Kobak & Berens, *Nature Comms* 2019



UMAP is equivalent to t -SNE with increased attraction

Sebastian Damrich



UMAP

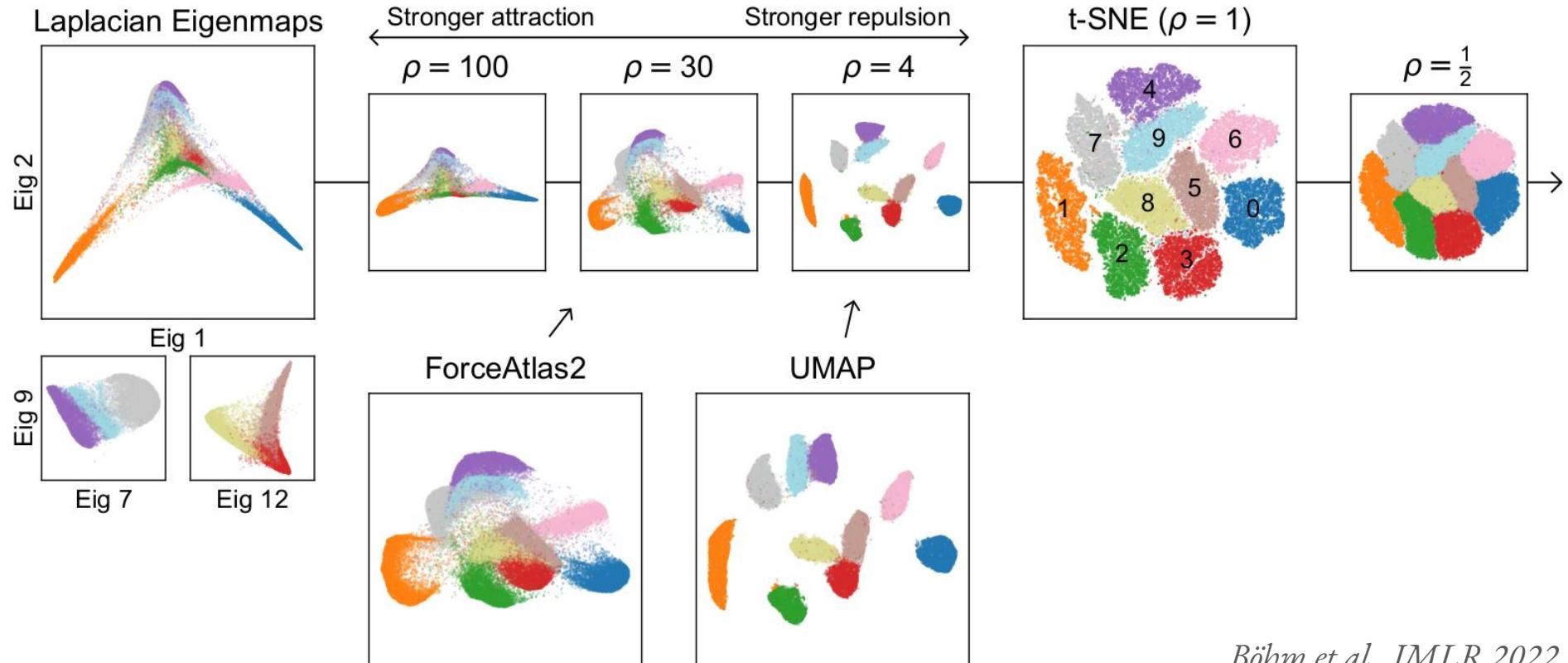


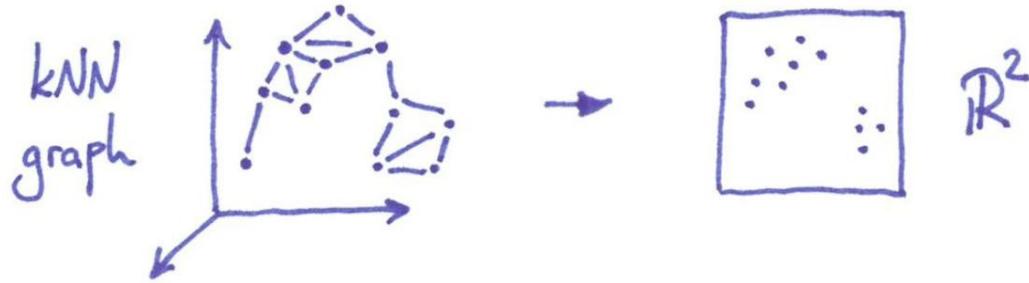
ForceAtlas2



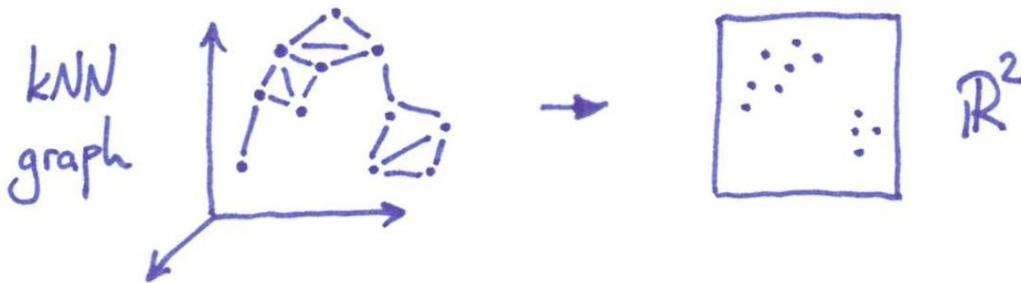
Damrich et al., ICLR 2023

Other neighbour embedding algorithms





1. What if my graph is not a k NN graph?
2. What if I want higher embedding dimensionality?



1. What if my graph is not a k NN graph?
2. What if I want higher embedding dimensionality?



$$-\sum_{ij} p_{ij} \log \frac{w_{ij}}{Z} \quad -\sum_{ij} p_{ij} \log \left(\frac{w_{ij}}{w_{ij} + \sum_{k \in \mathcal{I}_m} w_{ik}} \right)$$

InfoNC-t-SNE: sampling-based approximation

Can be optimized by SGD in batches for any output dimensionality.

$$\mathcal{L}_{\text{InfoNCE}} = - \sum_{ij} p_{ij} \log \left(\frac{w_{ij}}{w_{ij} + \sum_{k \in \mathcal{I}_m} w_{ik}} \right).$$

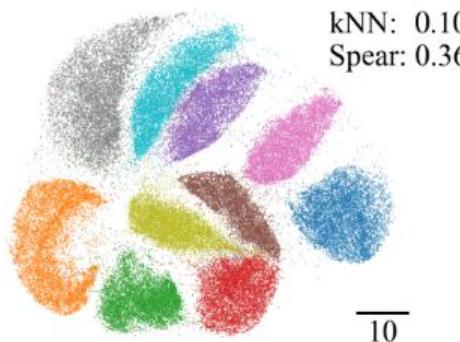
InfoNC-t-SNE

$$\mathcal{L}_{\text{InfoNCE}} = - \sum_{ij} p_{ij} \log \left(\frac{w_{ij}}{w_{ij} + \sum_{k \in \mathcal{I}_m} w_{ik}} \right).$$



InfoNC-t-SNE

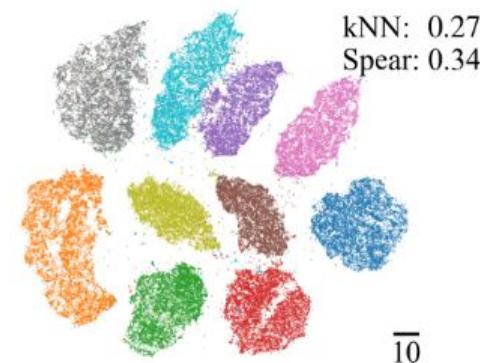
$$\mathcal{L}_{\text{InfoNCE}} = - \sum_{ij} p_{ij} \log \left(\frac{w_{ij}}{w_{ij} + \sum_{k \in \mathcal{I}_m} w_{ik}} \right).$$



(i) $m = 5$



(j) $m = 50$

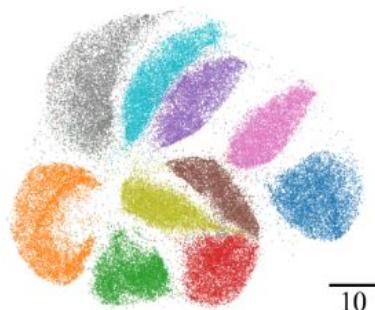


(k) $m = 500$

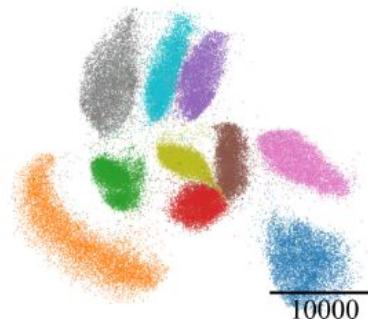
Parametric InfoNC-t-SNE

$$\mathcal{L}_{\text{InfoNCE}} = - \sum_{ij} p_{ij} \log \left(\frac{w_{ij}}{w_{ij} + \sum_{k \in \mathcal{I}_m} w_{ik}} \right).$$

Can optimize a parametric mapping (e.g. MLP) with the same loss!



(c) InfoNC-*t*-SNE

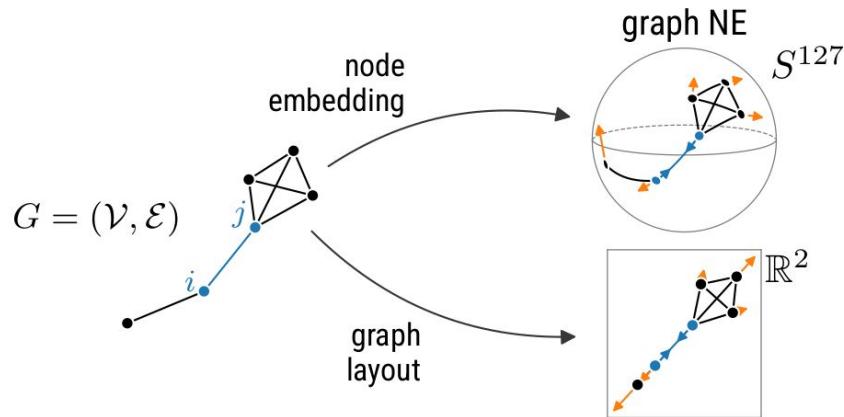


(f) Param. InfoNC-*t*-SNE



Niklas Böhm

Graph layouts & node embeddings via NE algorithms



InfoNCE

$$\ell_{ij} = -\log \frac{\exp(\mathbf{y}_i^\top \mathbf{y}_j / \tau)}{\sum_k \exp(\mathbf{y}_i^\top \mathbf{y}_k / \tau)}$$

KL divergence

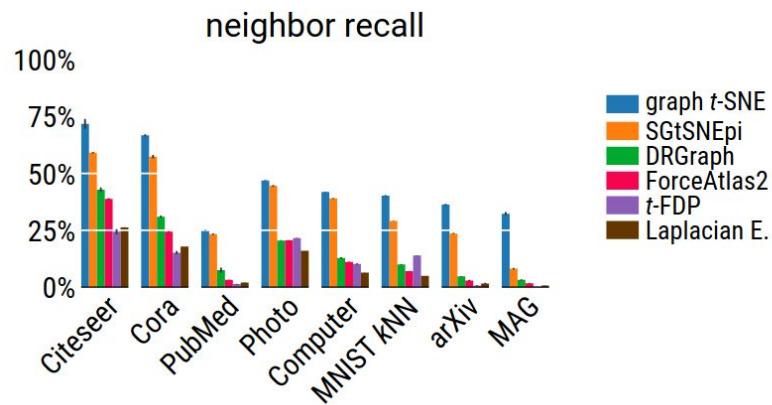
$$\ell_{ij} = -\log \frac{(1 + \|\mathbf{y}_i - \mathbf{y}_j\|^2)^{-1}}{\sum_{kl} (1 + \|\mathbf{y}_k - \mathbf{y}_l\|^2)^{-1}}$$



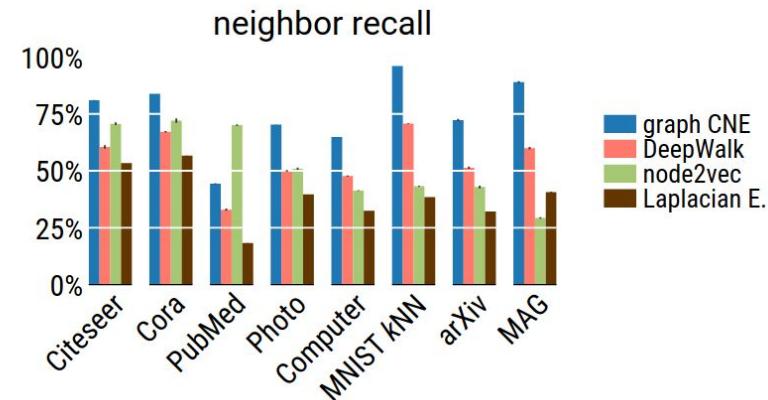
Marius Keute

Graph layouts & node embeddings via NE algorithms

2D embeddings

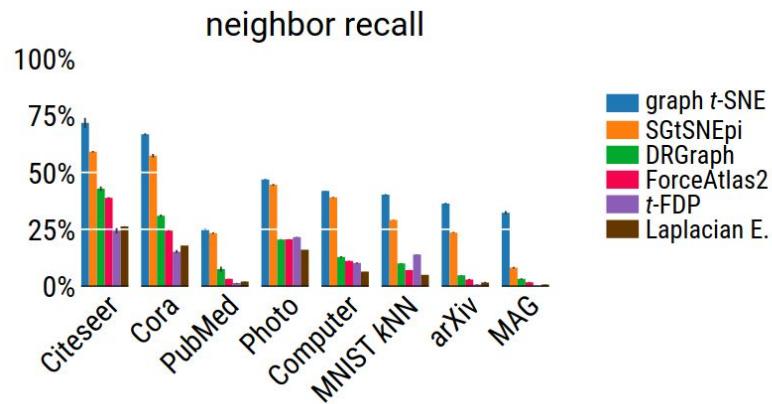


128D embeddings

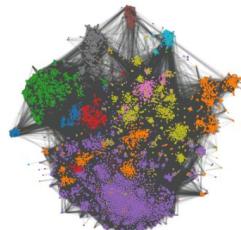


Graph layouts & node embeddings via NE algorithms

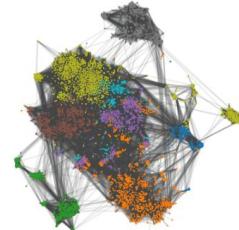
2D embeddings



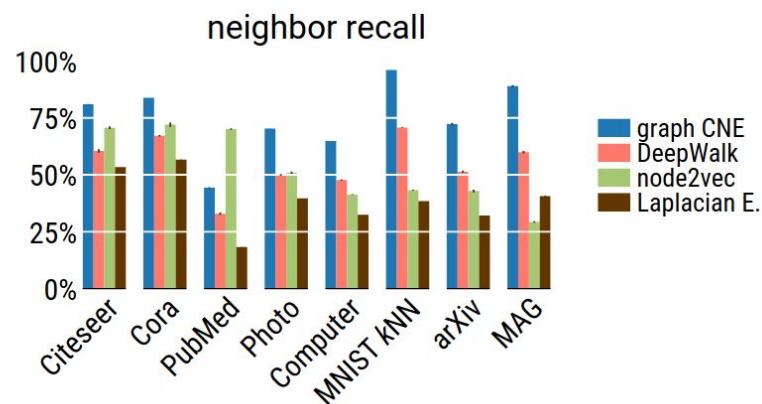
Computer



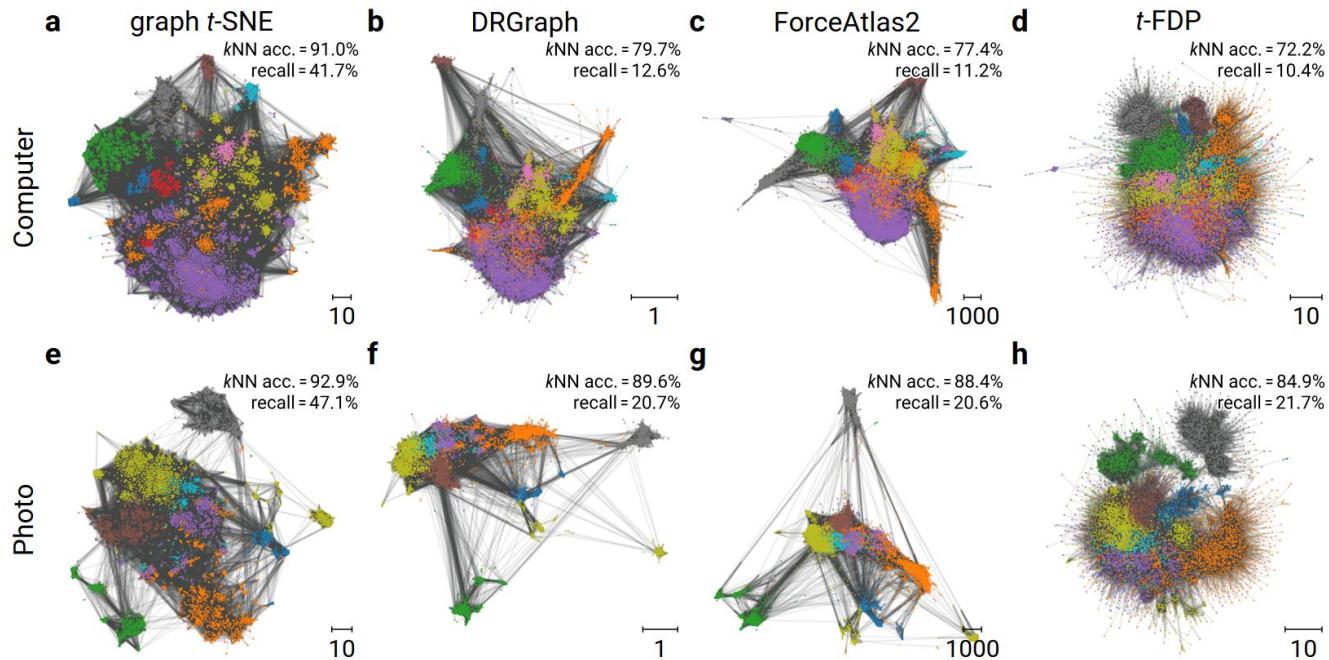
Photo



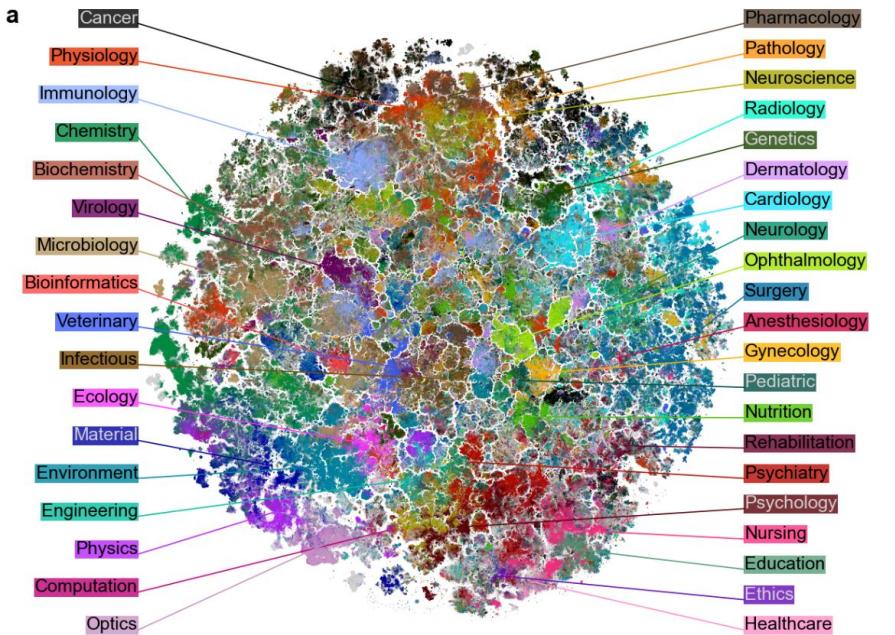
128D embeddings



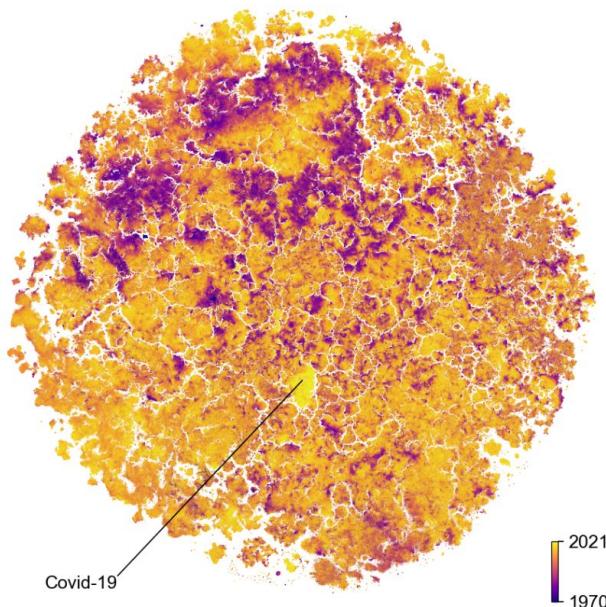
Graph layouts & node embeddings via NE algorithms



LLM-based embeddings of text data



b

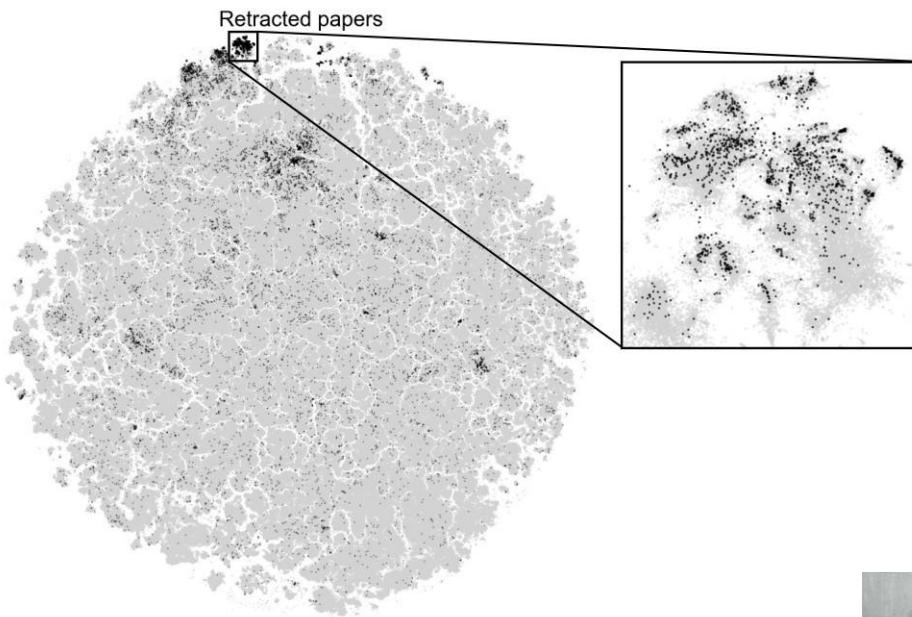
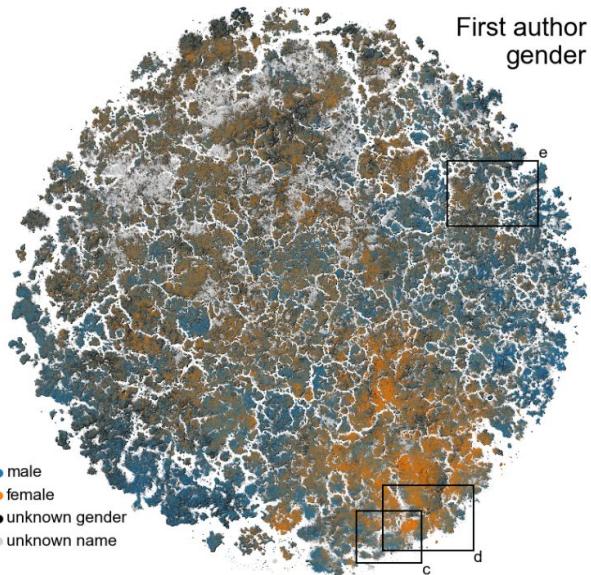


All 21 million English abstracts from PubMed

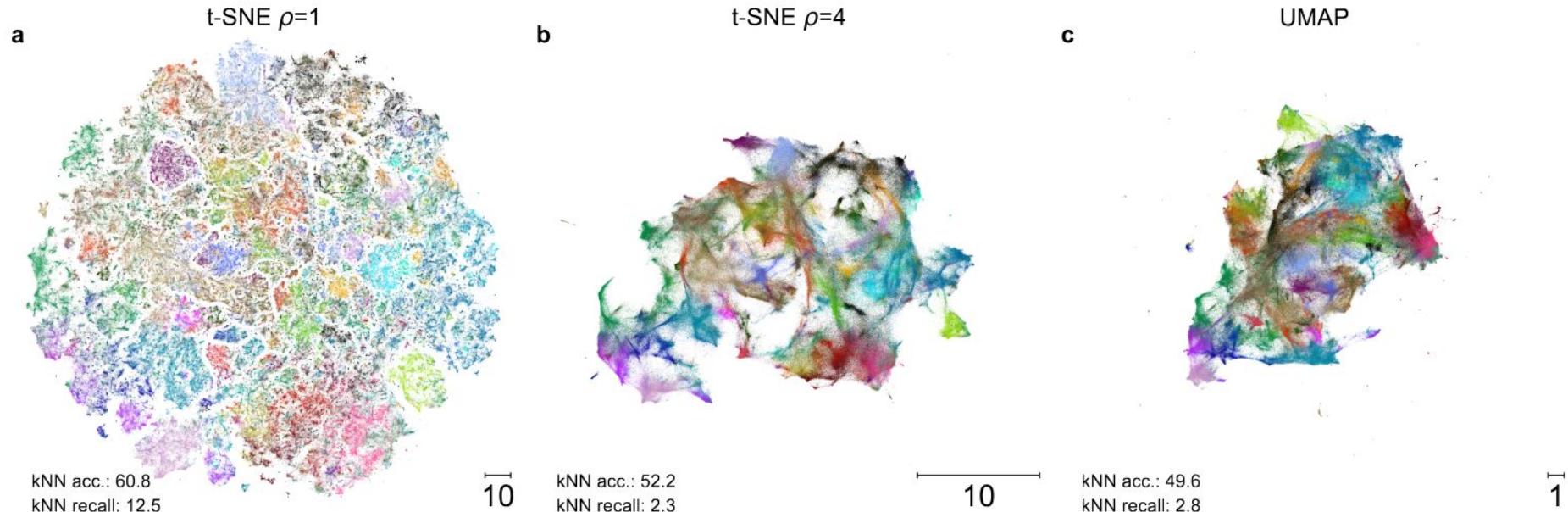


<https://static.nomic.ai/pubmed.html>

LLM-based embeddings of text data

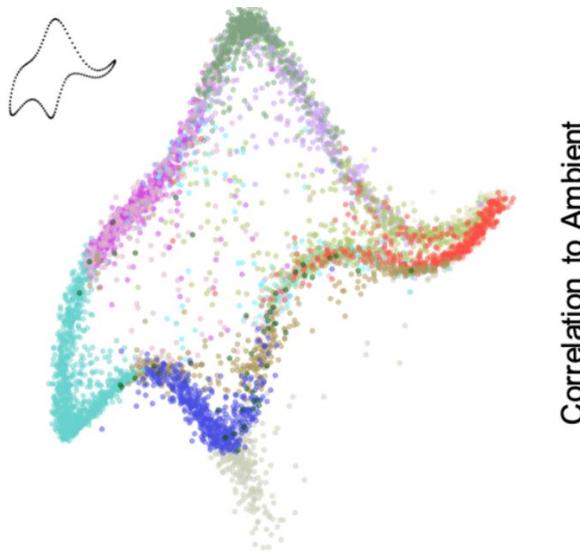


Attraction-repulsion spectrum (*t*-SNE/UMAP) for PubMed

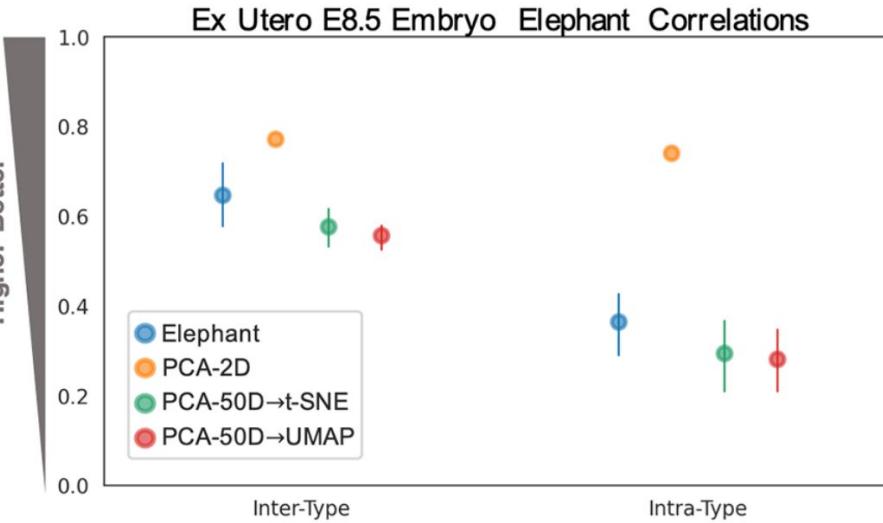


Embeddings of a subset of the PubMed dataset using different neighbor embedding methods. Subset size: 1,000,000 labeled papers.

“All embeddings are wrong ...”



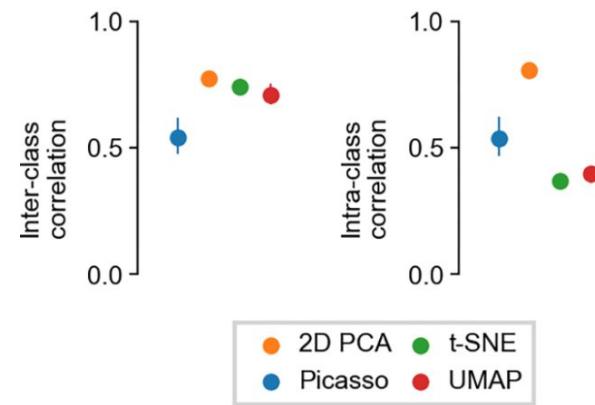
Correlation to Ambient





“... but some are useful”

Jan Lause





Jan Lause

“... but some are useful”

