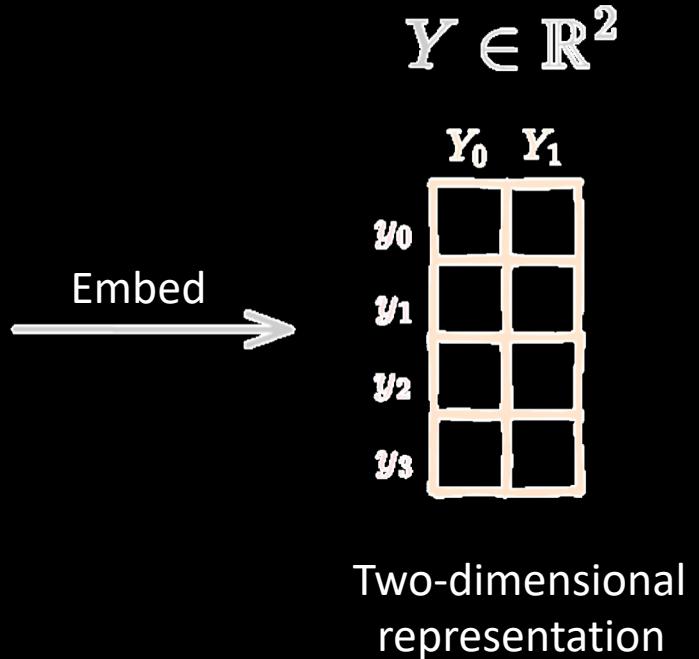
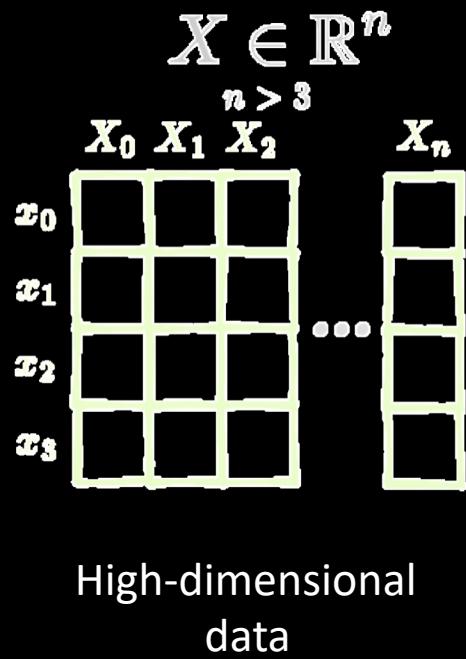


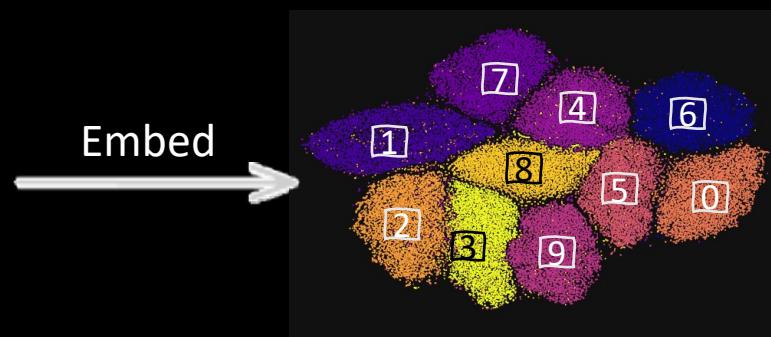
# t-SNE Visualisations With Adaptive Degrees of Freedom parameter

Artemii Shlychkov  
Supervisor: D. Kobak

# Core Idea of SNE



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



Loss function is the Kullback-Leibler divergence between pairwise similarities in the high-dimensional and in the low-dimensional spaces:

$$\mathcal{L} = \sum_{i,j} p_{i,j} \log \frac{p_{i,j}}{q_{i,j}}$$

Where  $p_{i,j}$  represent high dimensional similarities:

$$p_{j|i} = \frac{\exp(-\|\mathbf{x}_i - \mathbf{x}_j\|^2/2\sigma_i^2)}{\sum_{k \neq i} \exp(-\|\mathbf{x}_i - \mathbf{x}_k\|^2/2\sigma_i^2)} \quad p_{ij} = \frac{p_{i|j} + p_{j|i}}{2n}$$

And  $q_{i,j}$  represent low dimensional similarities.

Update rule for the points coordinates is:

$$\boxed{y_i \leftarrow y_i - \eta \cdot \frac{\partial \mathcal{L}}{\partial y_i}}$$

# Low dimensional similarities $q_{i,j}$

$$q_{i,j} = \frac{k(d)}{Z} \quad \text{with } Z \text{ as a normalization term}$$

Gaussian kernel:  $k = e^{(-d_{i,j}^2)}$

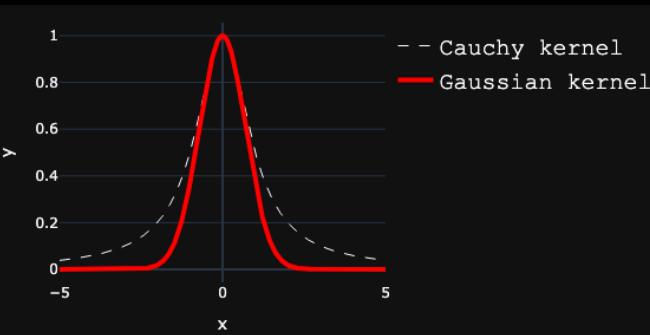
$\alpha \Rightarrow \infty$

?

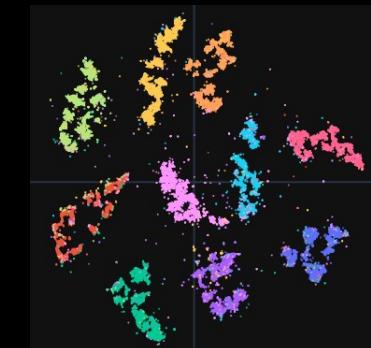
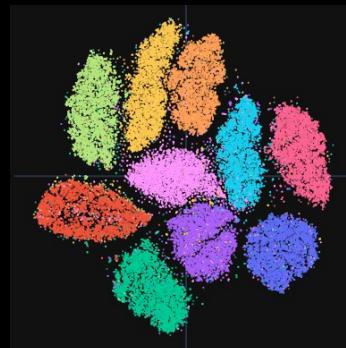
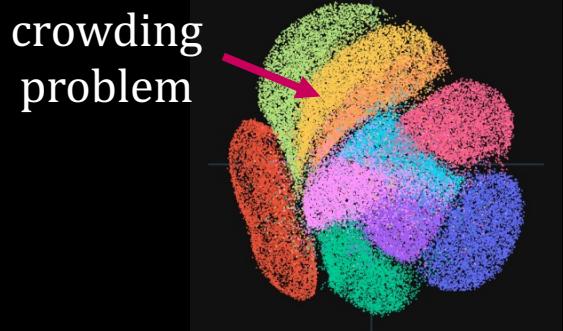
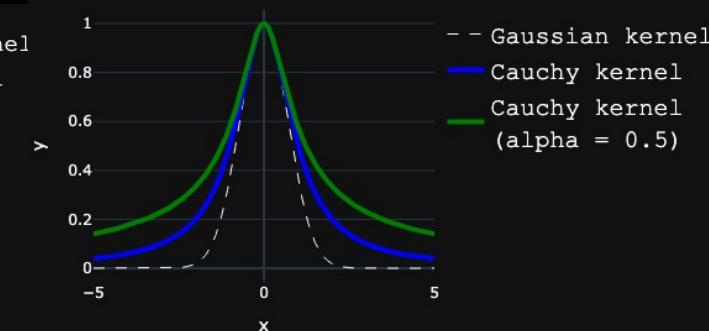
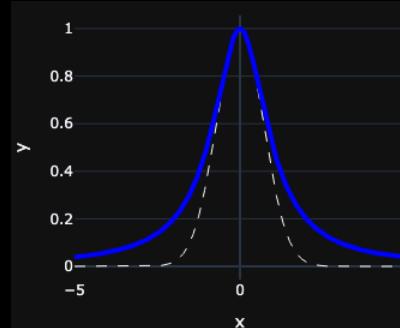
t-distribution kernel:  $k = \left(1 + \frac{d_{i,j}^2}{\alpha}\right)^{-\alpha}$

$\alpha = 1$

$\alpha < 1$



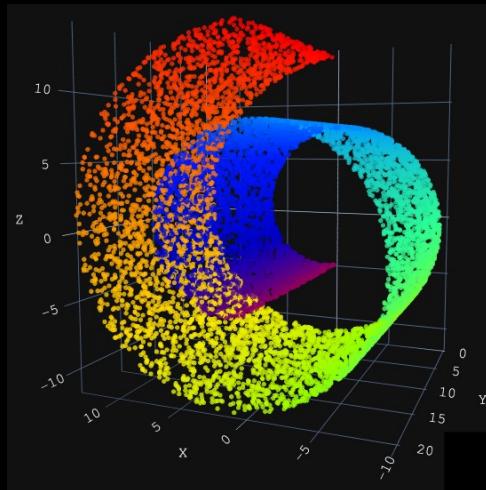
Optimize for alpha ?  
 $\frac{\partial \mathcal{L}}{\partial \alpha}$



# Methods

- openTSNE library:
  - pavlin-policar/openTSNE with standard t-SNE parameters
    - Custom implementation with adaptable alpha parameter
  - Embedding quality metrics: KL divergence and kNN recall
  - Explored datasets:

Swiss Roll

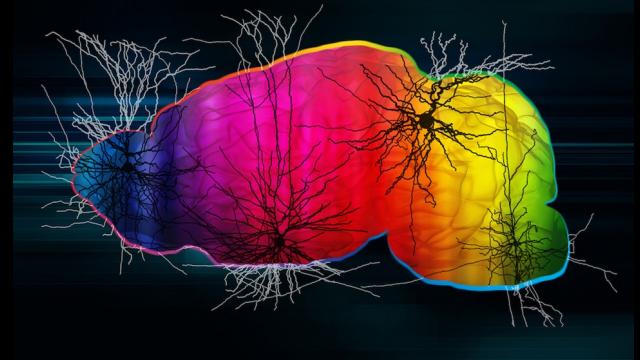


Digits / MNIST

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

Single-cell  
transcriptomics

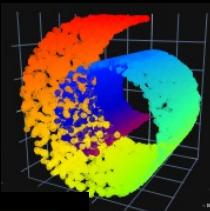
(Tasic et al dataset)\*



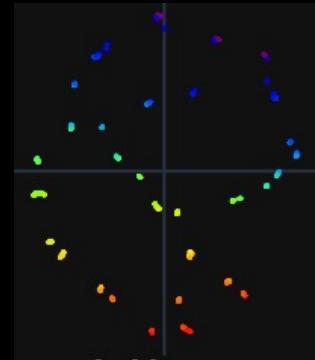
\*Tasic et al, 2018

I. Exploring t-SNE  
embeddings of different  
datasets using an array of  
alphas

# Swiss Roll (3000 samples)

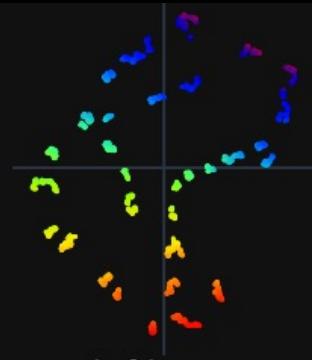


$\alpha = 0.3$



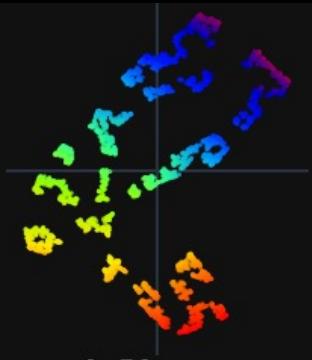
KL: 2.09  
kNN Recall: 0.64

$\alpha = 0.6$



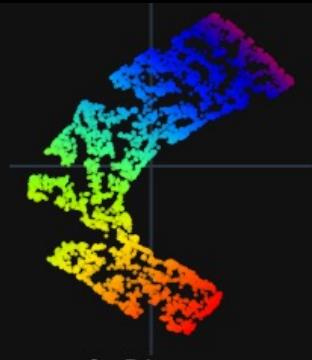
KL: 1.01  
kNN Recall: 0.77

$\alpha = 0.8$



KL: 0.72  
kNN Recall: 0.82

$\alpha = 1$



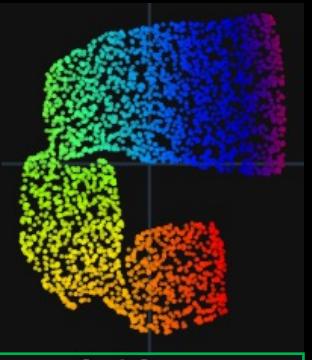
KL: 0.51  
kNN Recall: 0.87

$\alpha = 2$



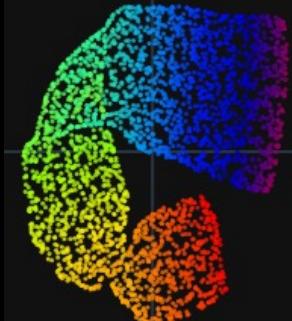
KL: 0.20  
kNN Recall: 0.92

$\alpha = 4$



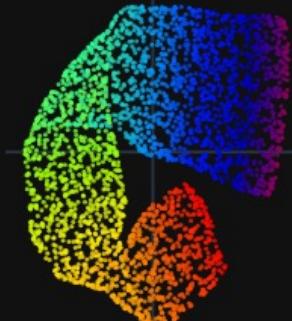
KL: 0.16  
kNN Recall: 0.91

$\alpha = 8$



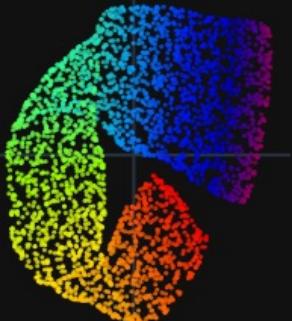
KL: 0.18  
kNN Recall: 0.88

$\alpha = 16$



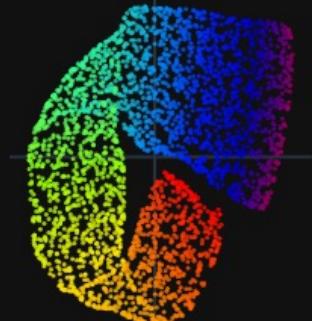
KL: 0.19  
kNN Recall: 0.86

$\alpha = 32$



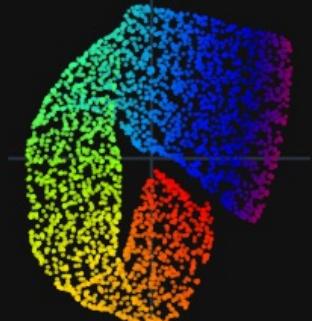
KL: 0.20  
kNN Recall: 0.85

$\alpha = 64$



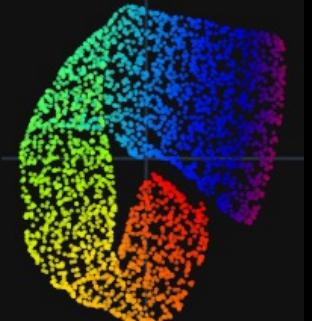
KL: 0.21  
kNN Recall: 0.84

$\alpha = 100$



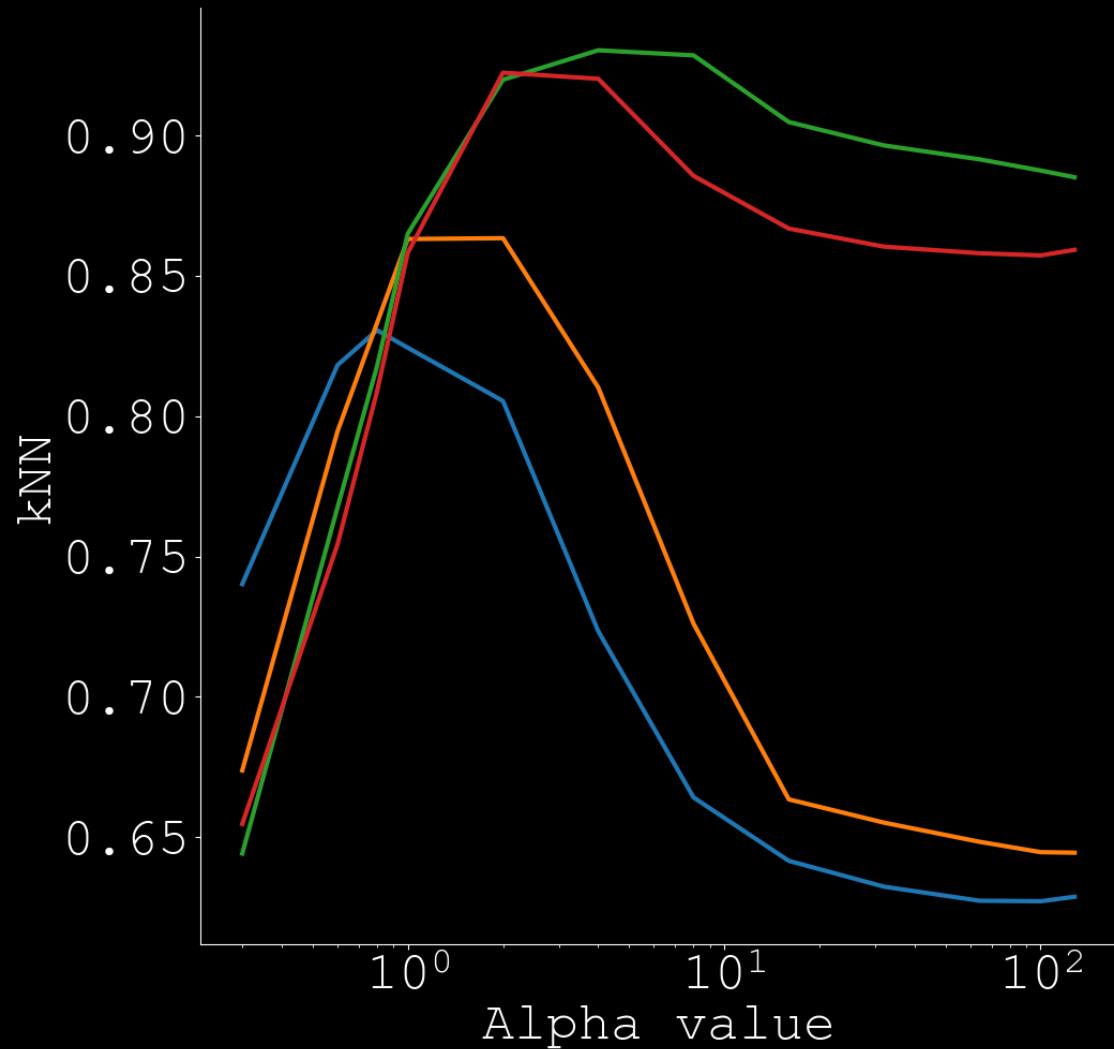
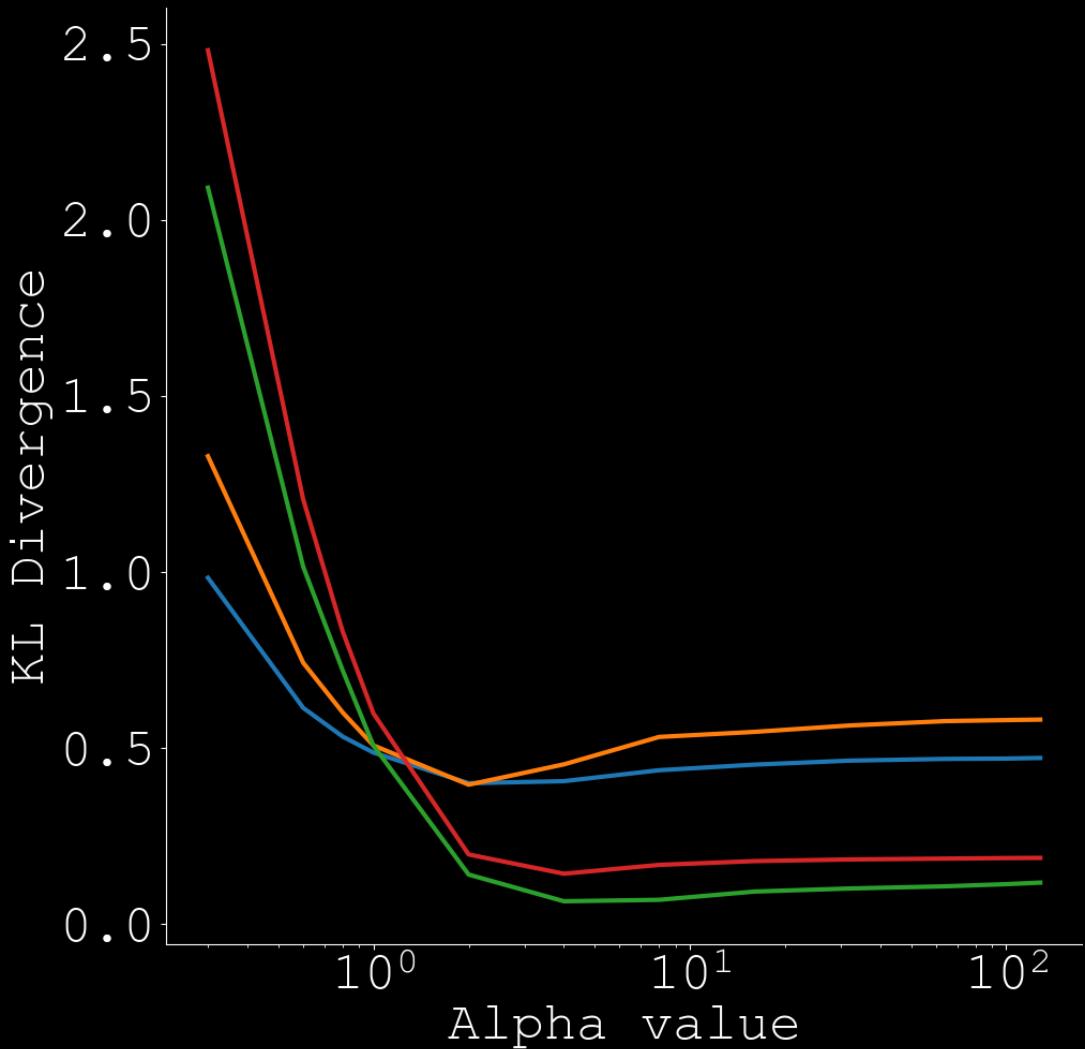
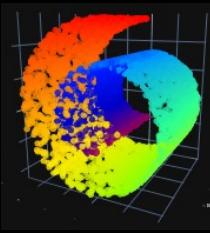
KL: 0.21  
kNN Recall: 0.84

$\alpha = 128$

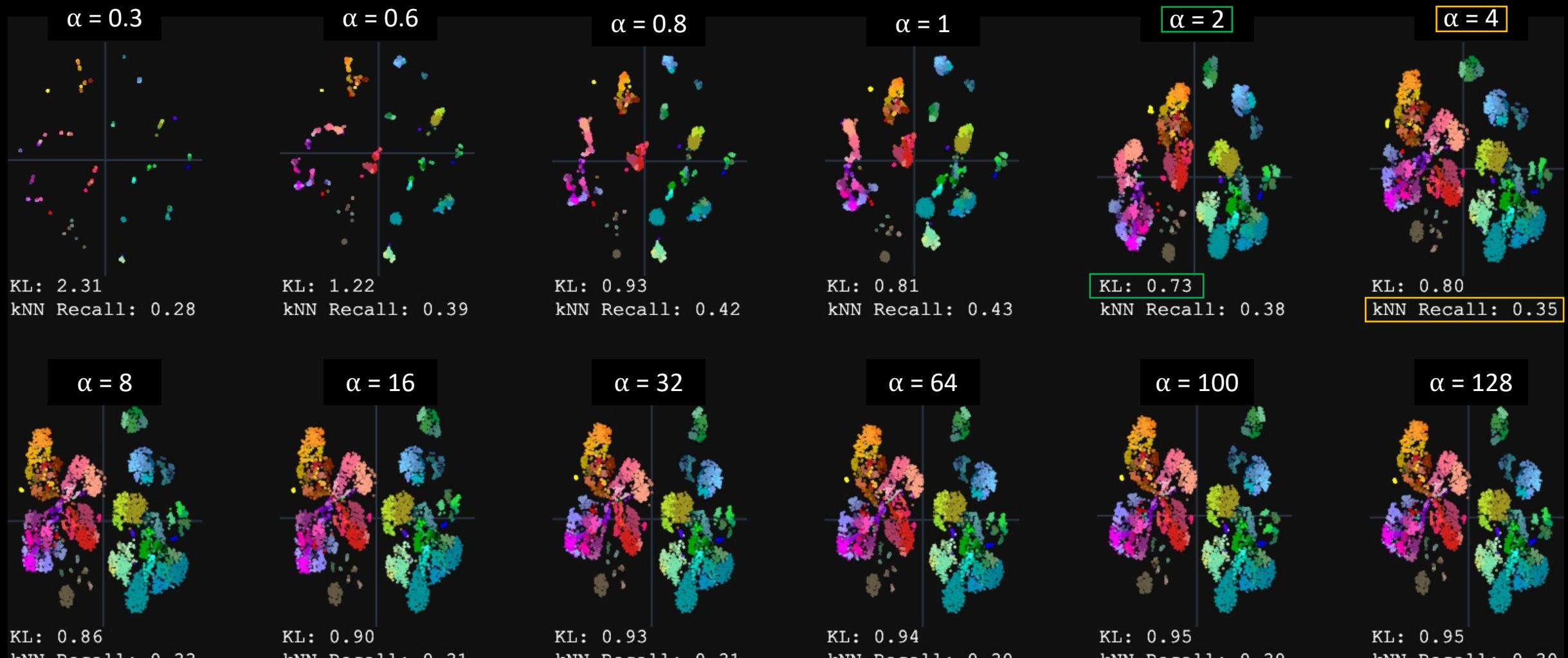


KL: 0.22  
kNN Recall: 0.84

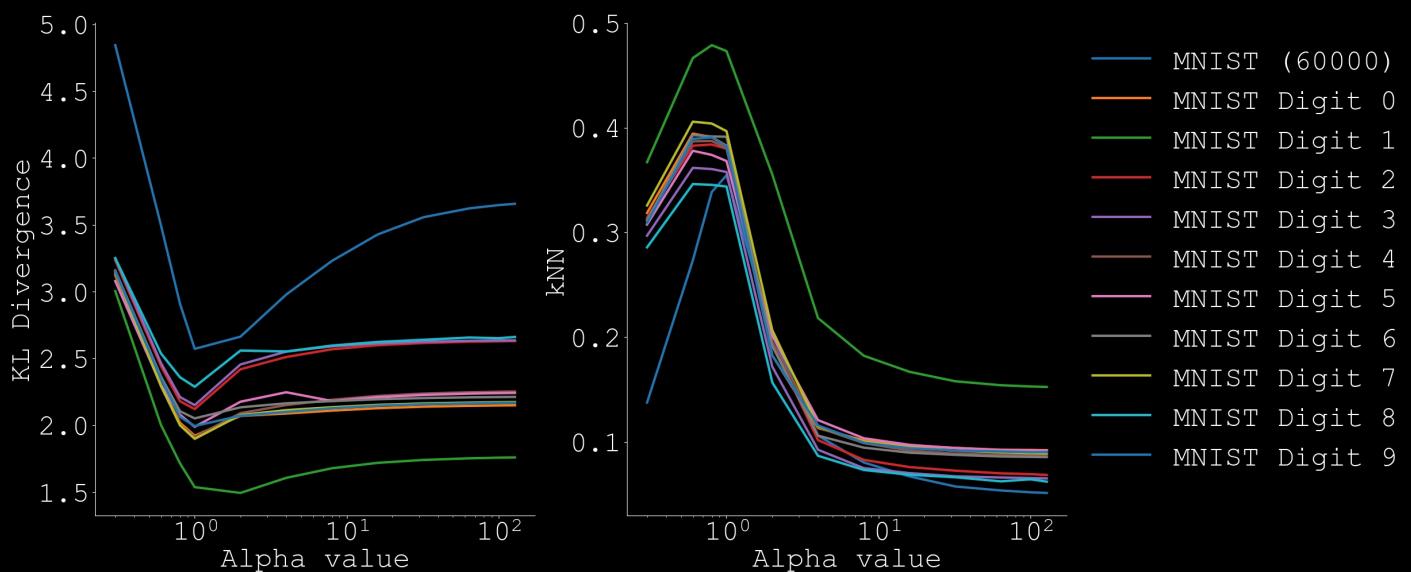
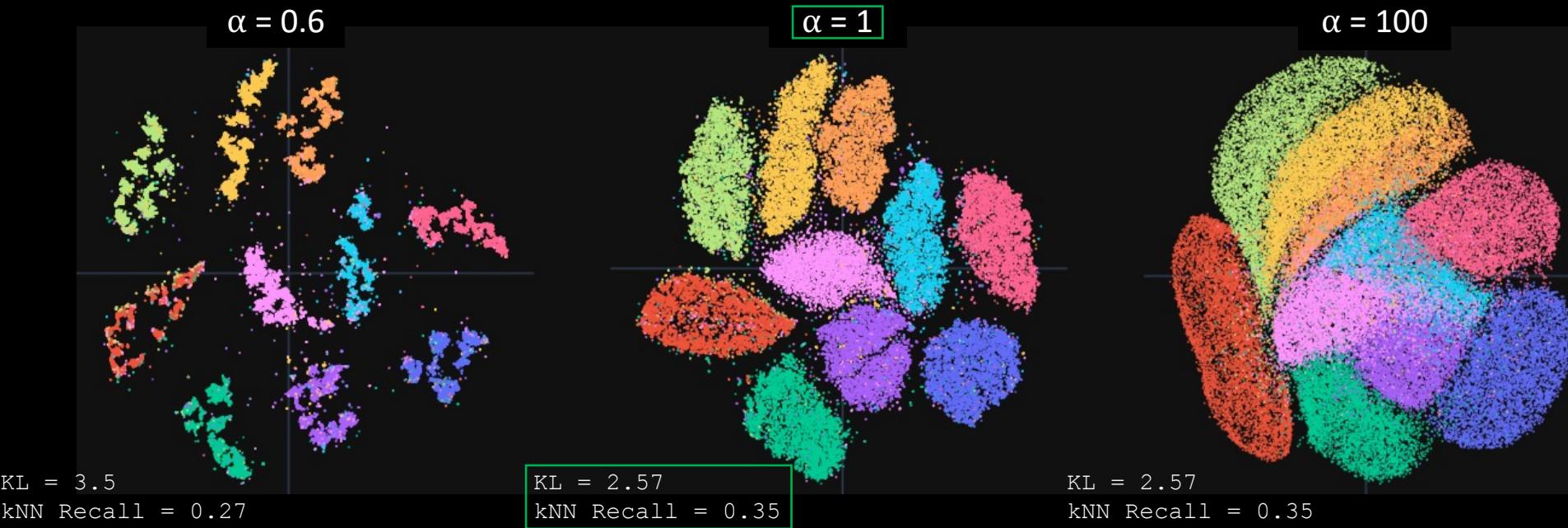
KL divergence and kNN recall for t-SNE embeddings of SwissRoll with different number of samples



# Single cell transcriptomics (5000 samples)



# MNIST dataset (60000 samples)



## II. Optimizing t-SNE embedding for the $\alpha$ parameter

# Derivation of the gradient of alpha

Given the Loss function:

$$\mathcal{L} = \sum_{i,j} p_{i,j} \log \frac{p_{i,j}}{q_{i,j}}$$

Taking partial derivative of the relevant term:

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

where

$$q_{i,j} = \frac{k}{Z}$$

with kernel  $k = \left( 1 + \frac{d_{i,j}^2}{\alpha} \right)^{-\alpha}$

and  $Z = \sum_{k \neq \ell} \left( 1 + \frac{d_{k,\ell}^2}{\alpha} \right)^{-\alpha}$

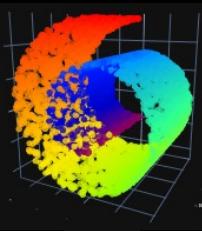
After several transformations finally obtain:

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

With the update rule:

$$\boxed{\alpha \leftarrow \alpha - \eta \cdot \frac{\partial \mathcal{L}}{\partial \alpha}}$$

# Testing alpha optimization on toy data

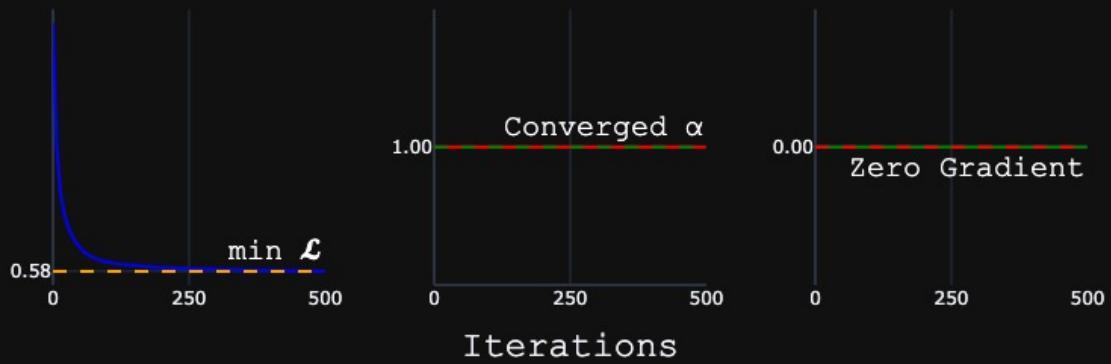


t-SNE embedding with fixed  $\alpha$  parameter

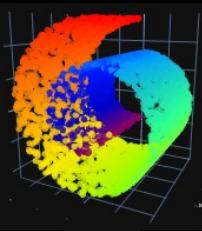
$\alpha = 1.0$

Final Loss = 0.58, kNN recall = 0.84

No optimisation



# Testing alpha optimization on toy data

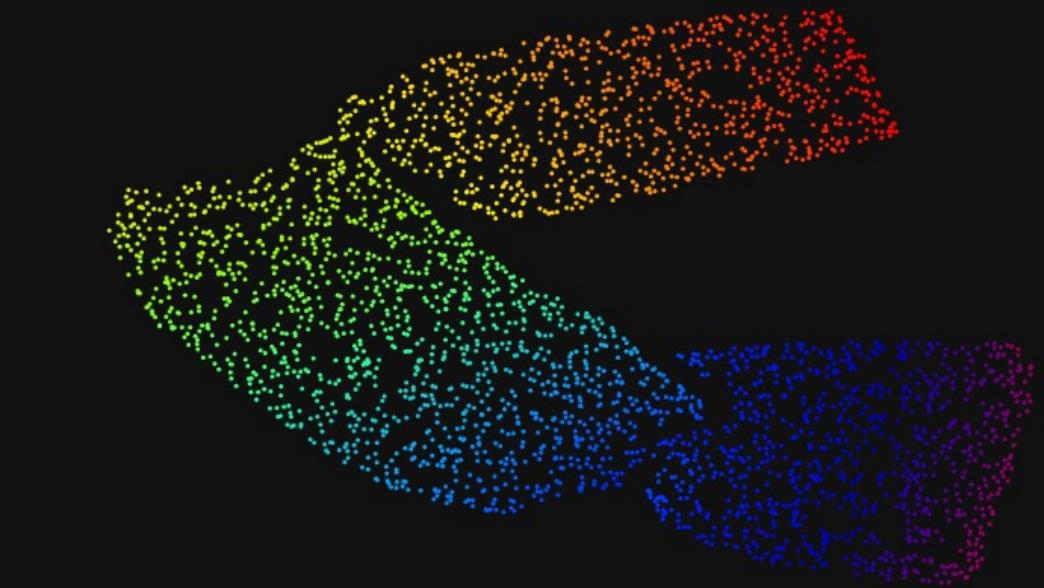


t-SNE embedding with adaptive  $\alpha$  parameter

Initial  $\alpha = 1.00$ ,  $\alpha$  learning rate = 0.9

Final  $\alpha = 3.27$ , Final Loss = 0.16, kNN recall = 0.92

Gradient Descent Optimisation



# Testing alpha optimization on toy data

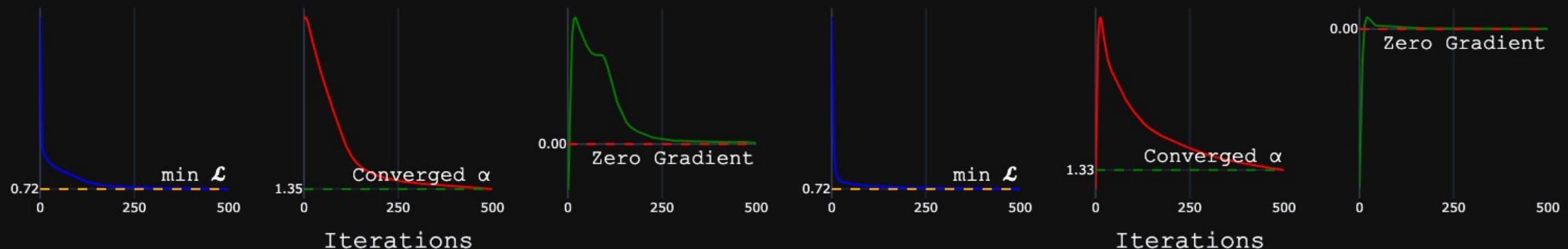
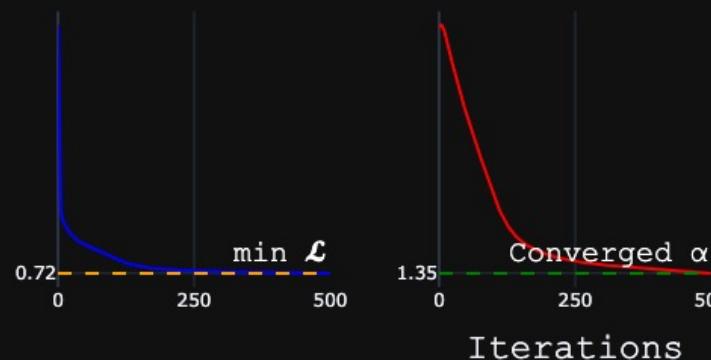
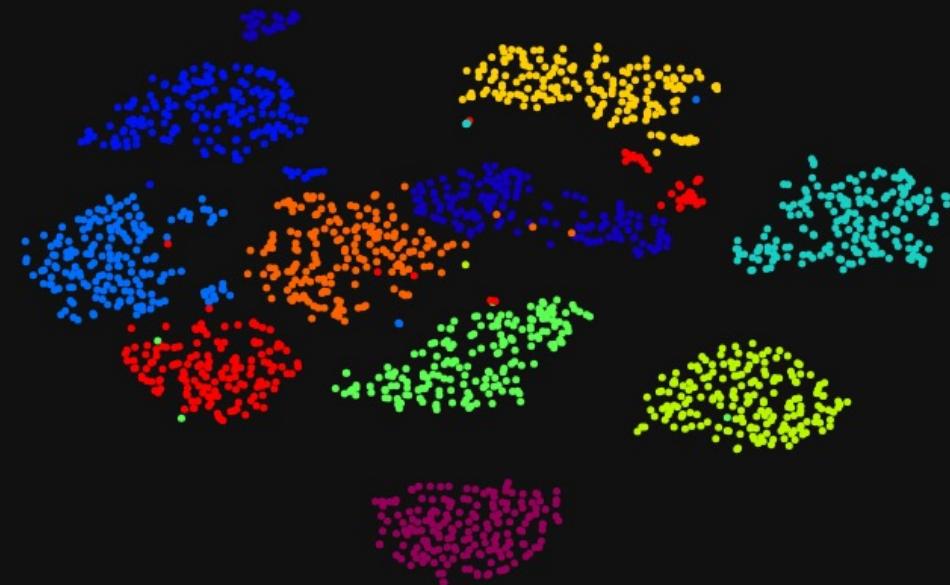
0 8 7 6 4 6 9 7 1 5 / 4 6  
0 1 2 3 4 4 6 2 9 3 0 1 2 3 4  
0 1 2 3 4 5 6 7 0 1 2 3 4 5 6  
7 4 2 0 9 1 2 8 9 4 0 9 5 0  
0 2 7 8 4 8 0 7 7 1 2 9 3 6  
5 3 9 4 3 7 2 3 8 1 2 4 8 8 7  
2 9 1 6 8 0 1 7 1 1 0 3 4 2 6 4  
7 7 6 3 0 7 4 2 7 4 9 1 0 6 8  
2 4 7 8 3 5 5 5 3 5 9 7 4 8 5

t-SNE embedding with adaptive  $\alpha$  parameter

Initial  $\alpha = 5.00$ ,  $\alpha$  learning rate = 0.5

Final  $\alpha = 1.35$ , Final Loss = 0.72, kNN recall = 0.58

(Gradient Descent Optimised)

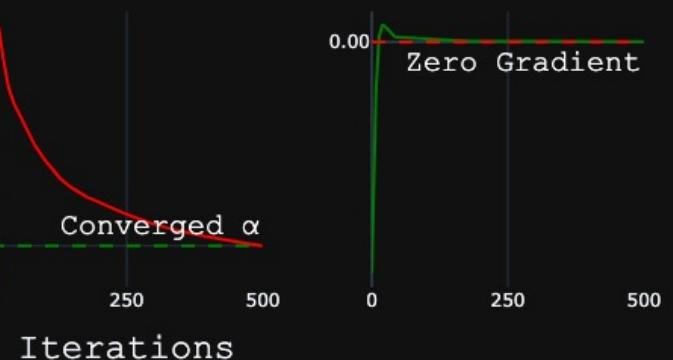
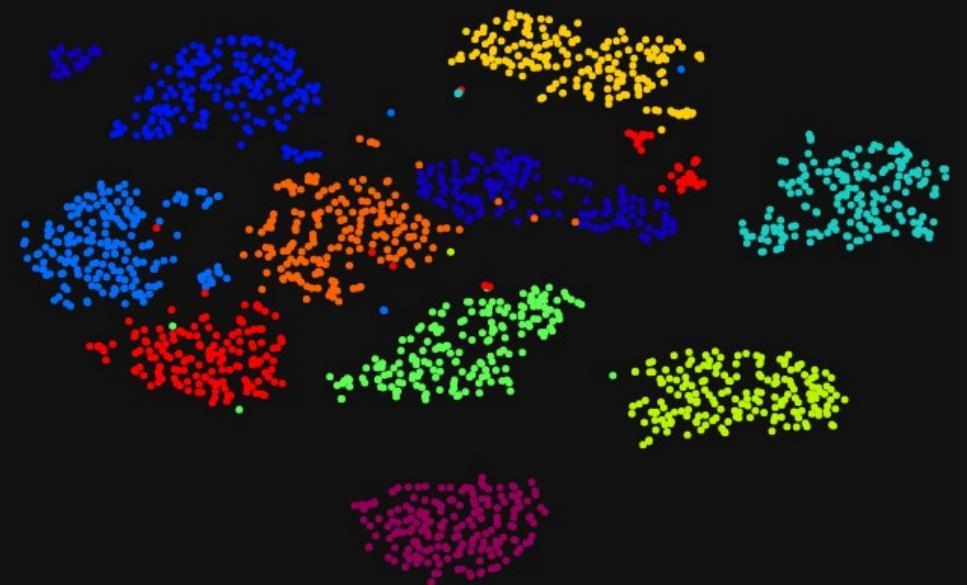


t-SNE embedding with adaptive  $\alpha$  parameter

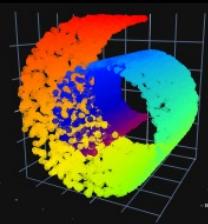
Initial  $\alpha = 1.00$ ,  $\alpha$  learning rate = 0.5

Final  $\alpha = 1.33$ , Final Loss = 0.72, kNN recall = 0.58

(Gradient Descent Optimised)



# Comparing BH-approximation with exact gradient

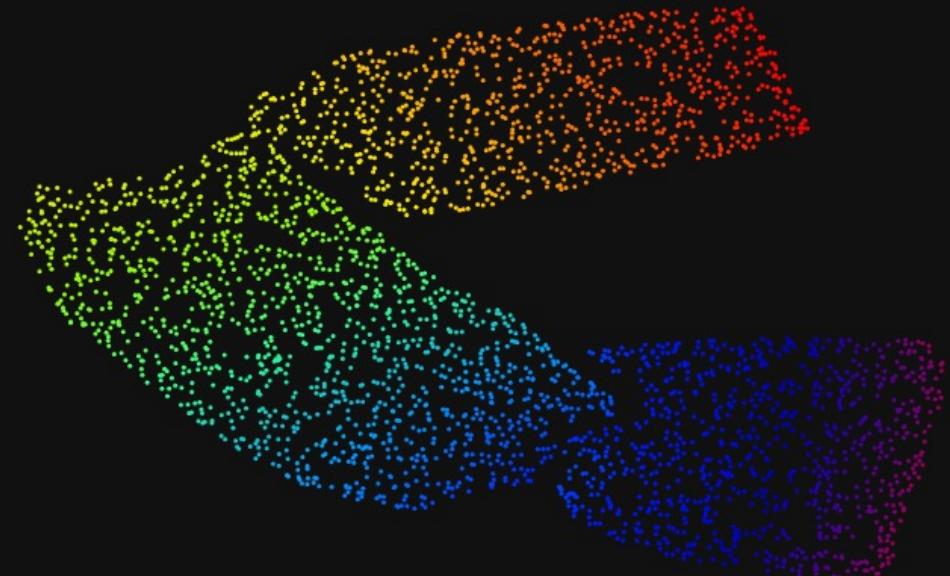


t-SNE embedding with adaptive  $\alpha$  parameter

Initial  $\alpha = 1.00$ ,  $\alpha$  learning rate = 0.9

Final  $\alpha = 3.27$ , Final Loss = 0.16, kNN recall = 0.92

Gradient Descent Optimization

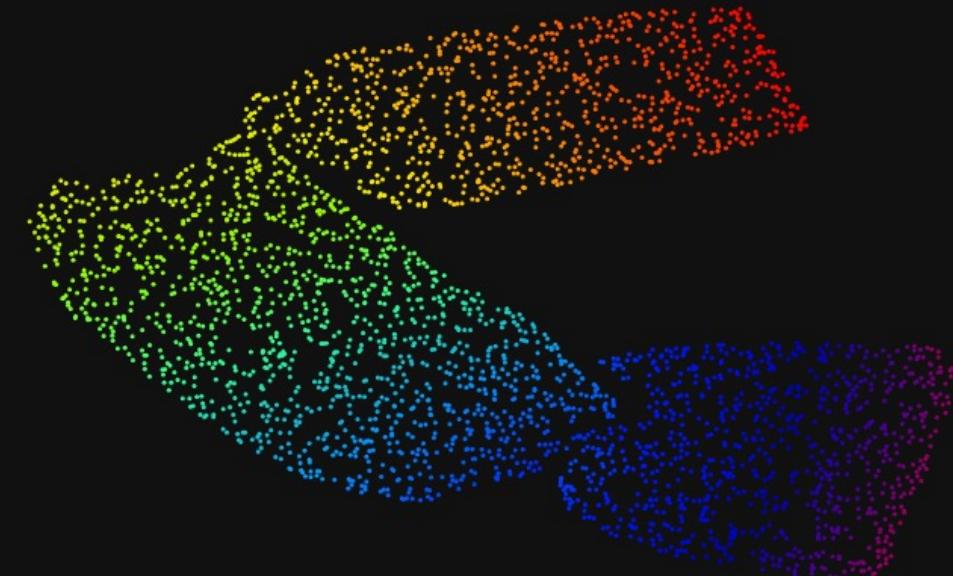


t-SNE embedding with adaptive  $\alpha$  parameter

Initial  $\alpha = 1.00$ ,  $\alpha$  learning rate = 0.9

Final  $\alpha = 3.93$ , Final Loss = 0.16, kNN recall = 0.92

Barnes-Hut Optimization



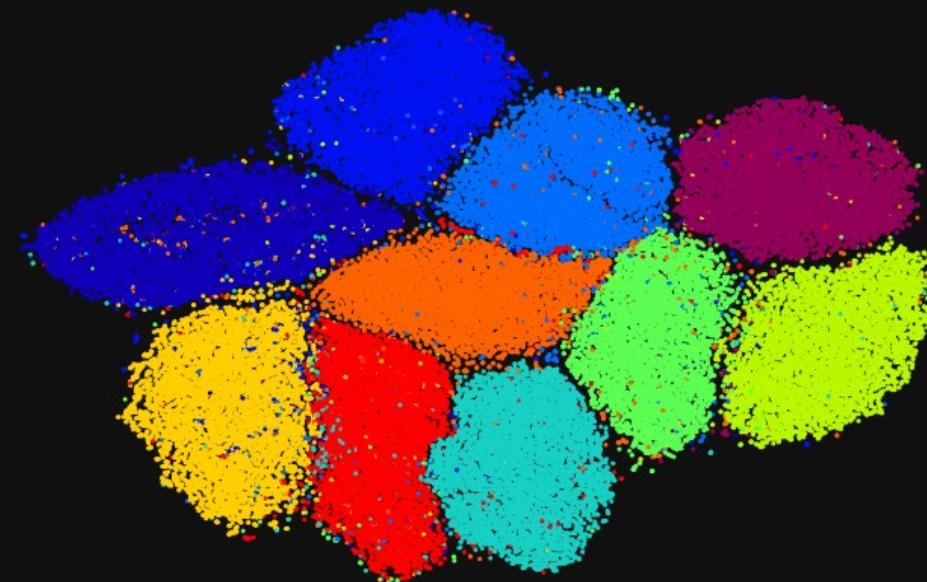
# Optimizing alpha on MNIST dataset

t-SNE embedding with adaptive  $\alpha$  parameter

Initial  $\alpha = 1.00$ ,  $\alpha$  learning rate = 0.5

Final  $\alpha = 1.52$ , Final Loss = 2.61, kNN recall = 0.32

Optimized using BH

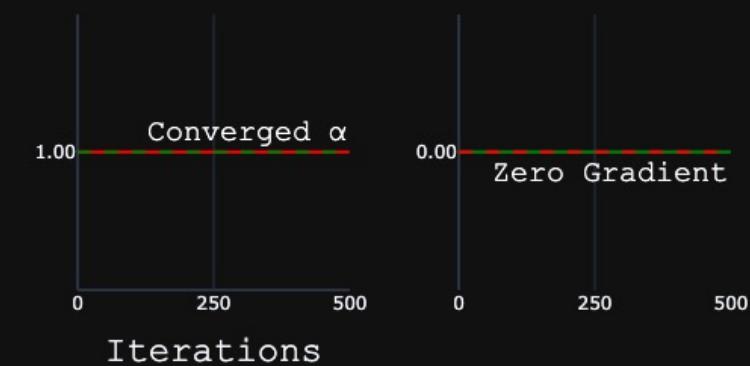
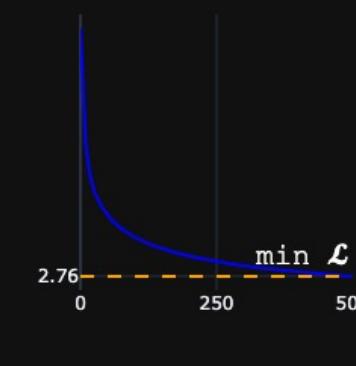
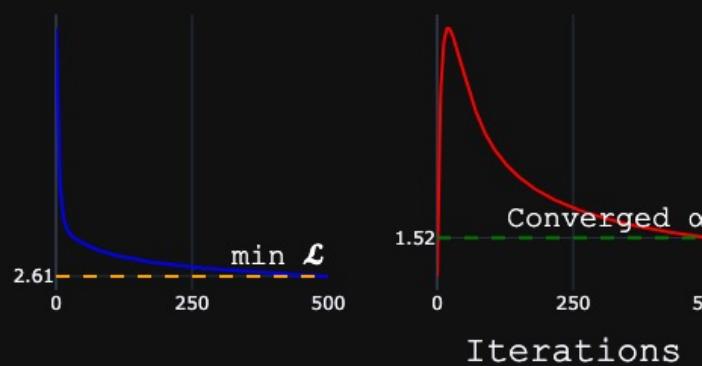
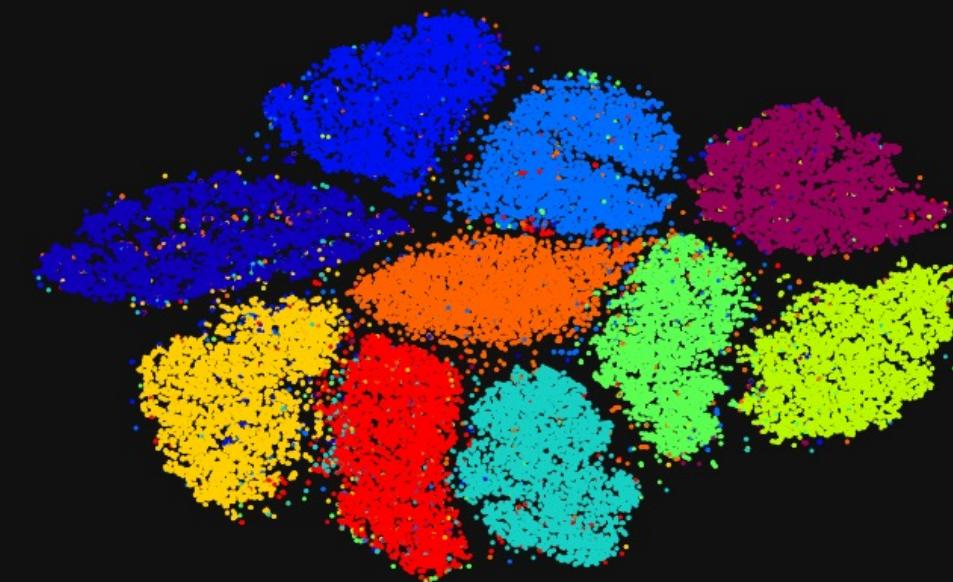


t-SNE embedding with fixed  $\alpha$  parameter

$\alpha = 1.00$

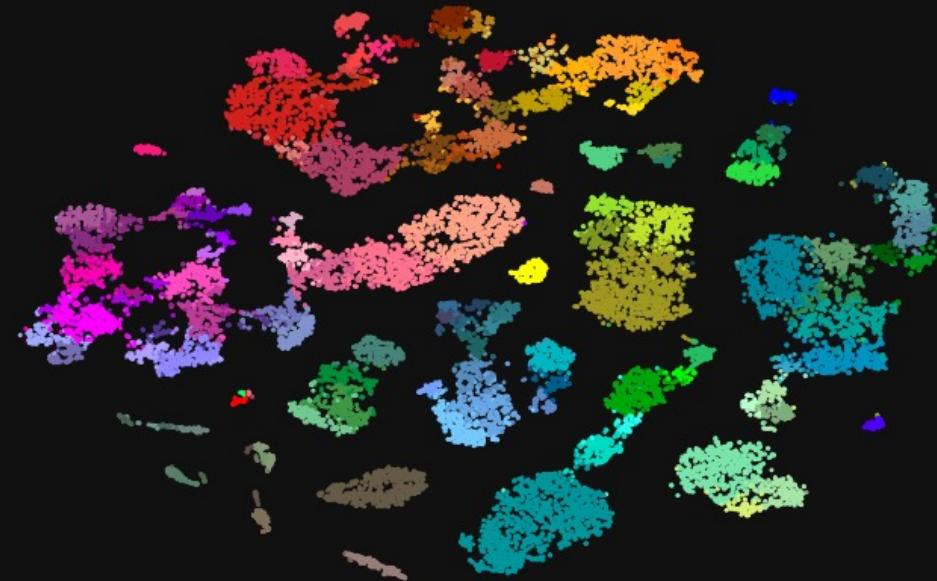
Final Loss = 2.76, kNN recall = 0.37

No optimization

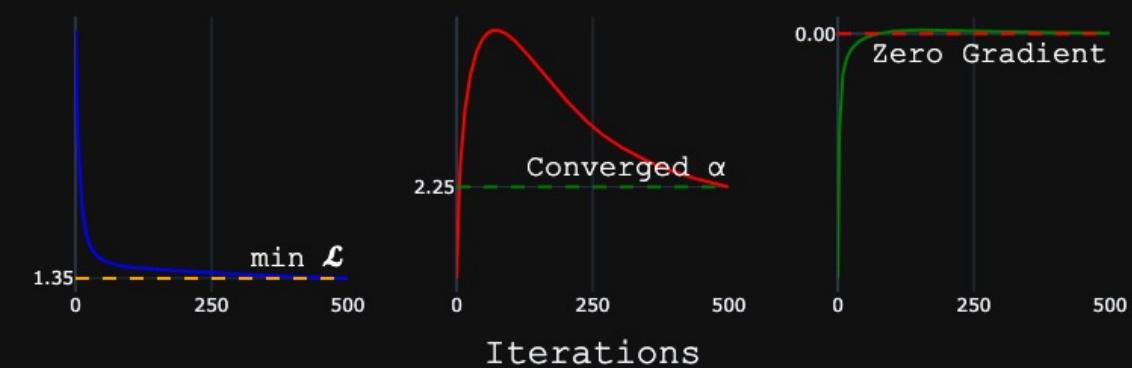
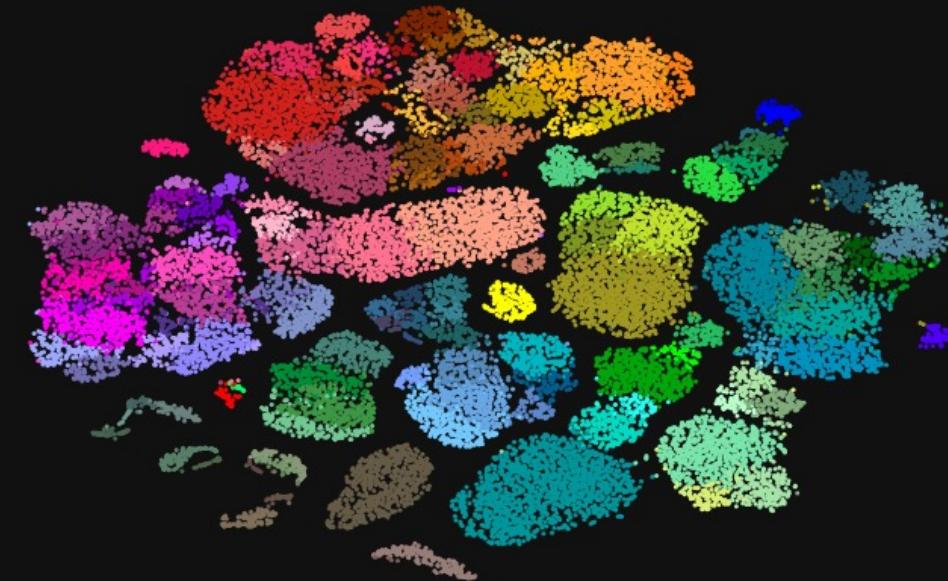


# Optimizing alpha on transcriptomics dataset

t-SNE embedding with fixed  $\alpha$  parameter  
 $\alpha = 1.00$   
Final Loss = 1.59, kNN recall = 0.45  
No optimization



t-SNE embedding with adaptive  $\alpha$  parameter  
Initial  $\alpha = 1.00$ ,  $\alpha$  learning rate = 0.5  
Final  $\alpha = 2.25$ , Final Loss = 1.35, kNN recall = 0.38  
Optimized with Barnes-Hutt



# Optimizing alpha on transcriptomics dataset



## Further work

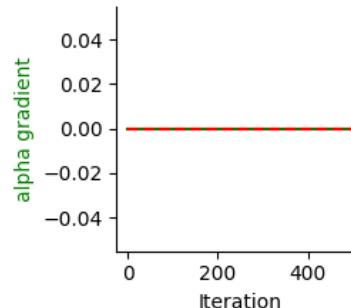
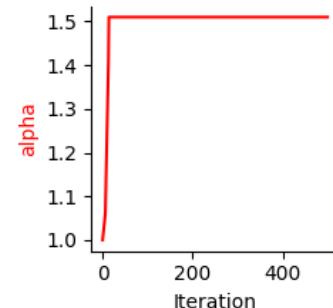
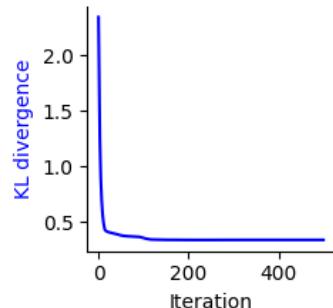
- Exploring larger variety of real-world datasets, employing more objective metrics for the embedding quality
- Exploring interactions of the alpha parameter with other t-SNE parameters (e.g. perplexity)
- Further integration of alpha gradient computation into existing openTSNE framework to achieve faster results at less computational costs

Thank you for your attention!

# 'Drift' optimization for alpha



t-SNE embedding with adaptive degrees of freedom  
SwissRoll with 500 samples  
Initial alpha = 1.0. Alpha learning rate = 0.1  
Final alpha = 1.51. Final Loss = 0.34  
(Drift optimized)



t-SNE embedding with adaptive degrees of freedom  
SwissRoll with 500 samples  
Initial alpha = 1.0. Alpha learning rate = 0.5  
Final alpha = 1.61. Final Loss = 0.33  
(Gradient descent optimized)

