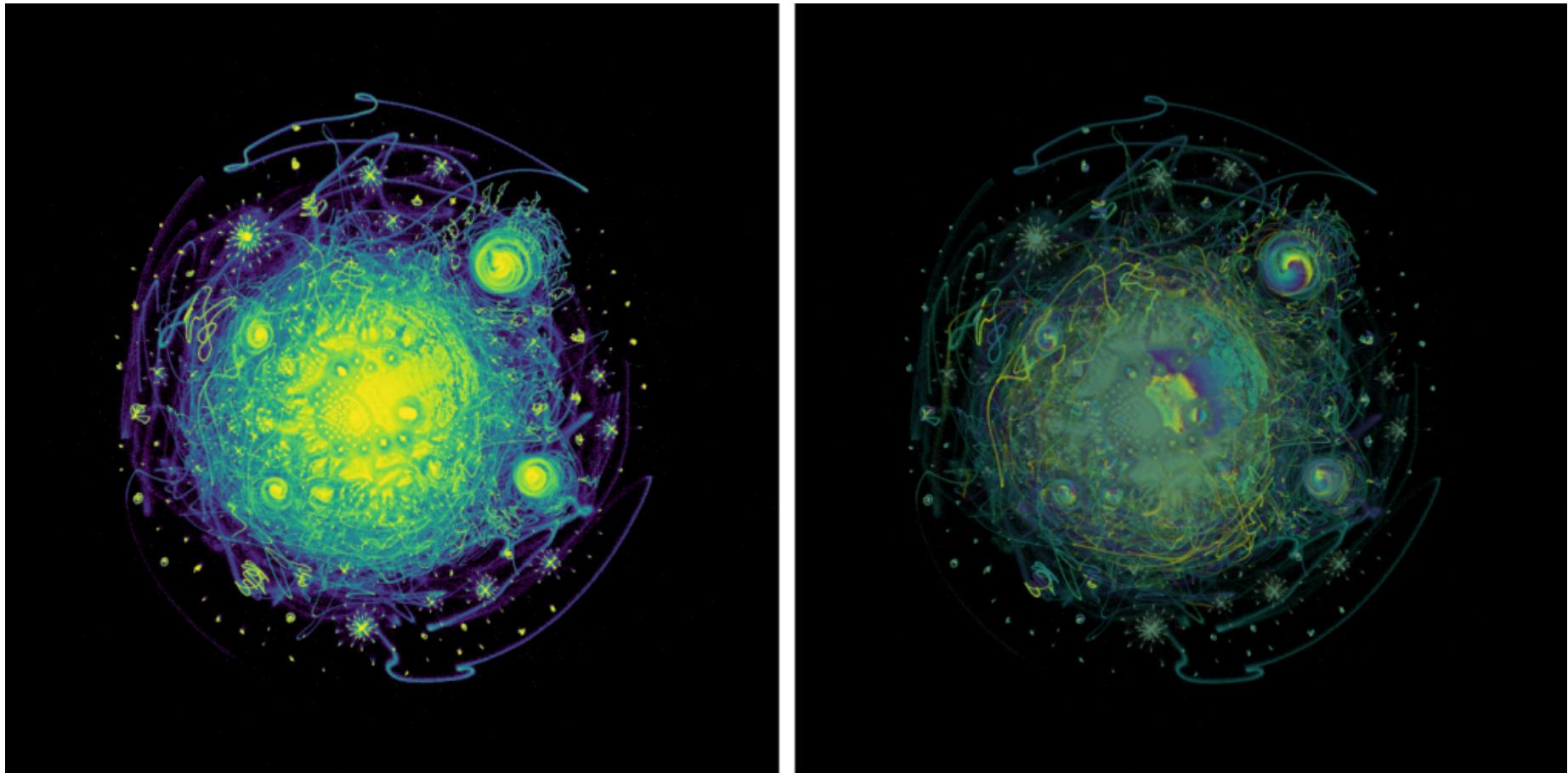


Dimension Reduction for OMICs Integration

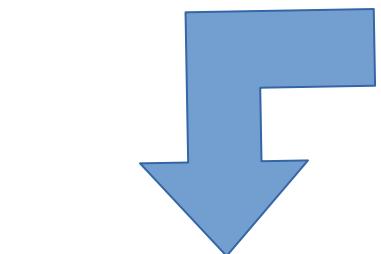
OMICs Integration and Systems Biology course

Nikolay Oskolkov, NBIS SciLifeLab

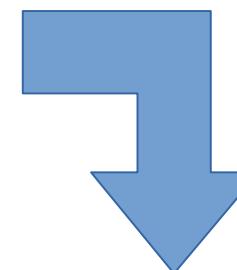
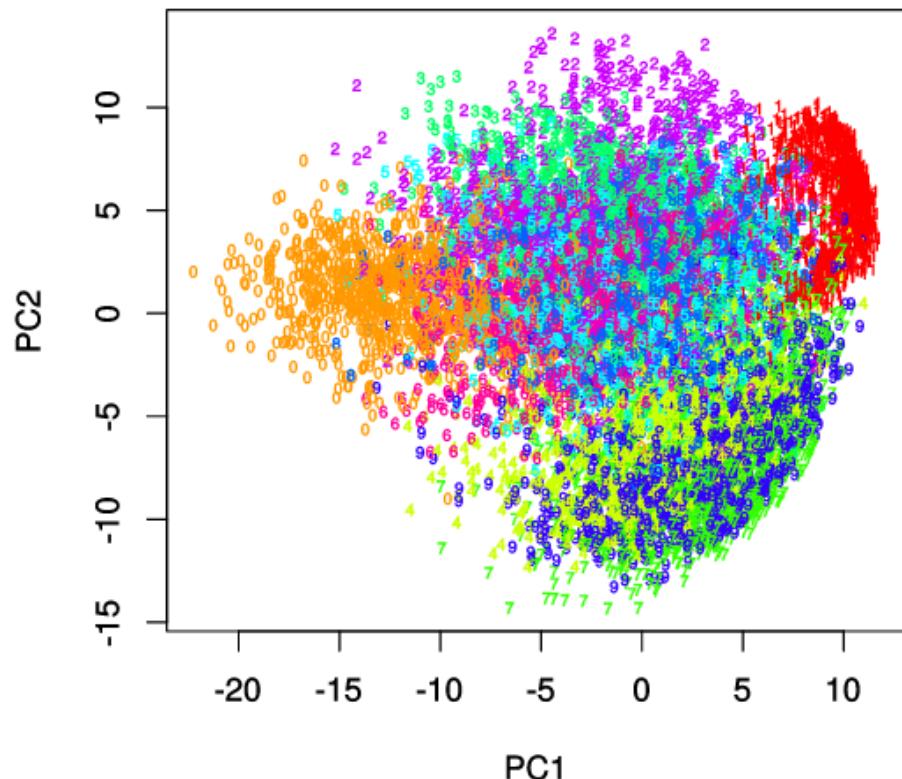
Lund, 5.10.2020



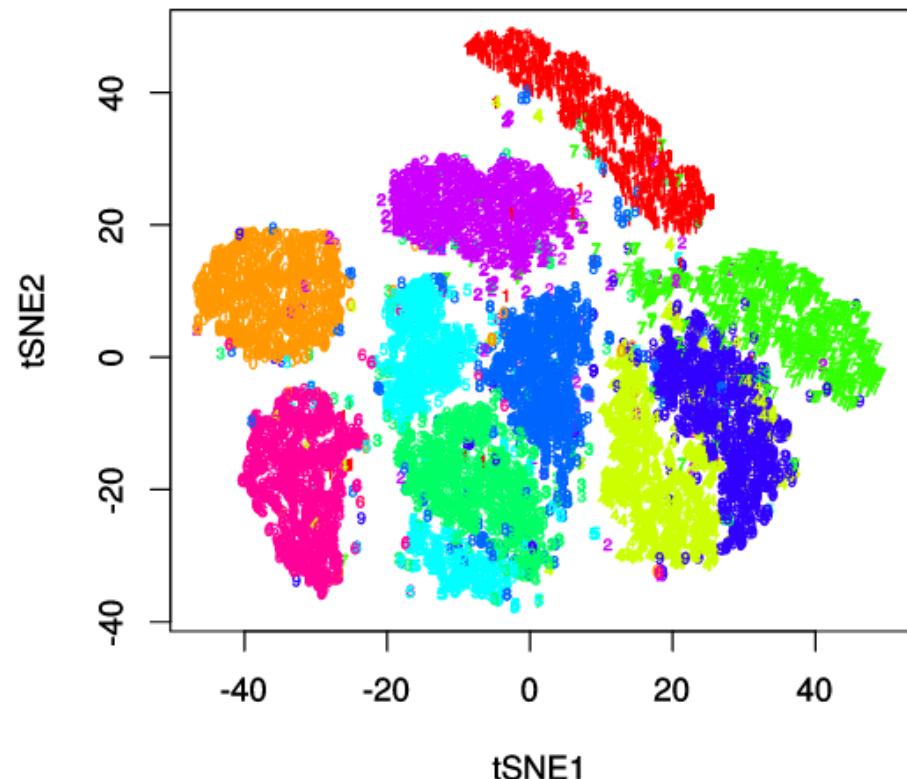
0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9



PCA PLOT WITH PRCOMP



tSNE MNIST



Dimension reduction is not for visualization but overcoming the Curse of Dimensionality

P is the number of features (genes, proteins, genetic variants etc.)
 N is the number of observations (samples, cells, nucleotides etc.)

Biomedicine

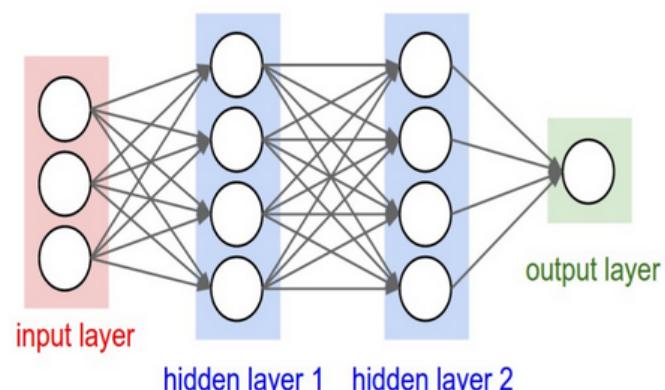
Bayesianism

 $P \gg N$

Frequentism

 $P \sim N$

Deep Learning

 $P \ll N$ 

Amount of Data



Ex.1

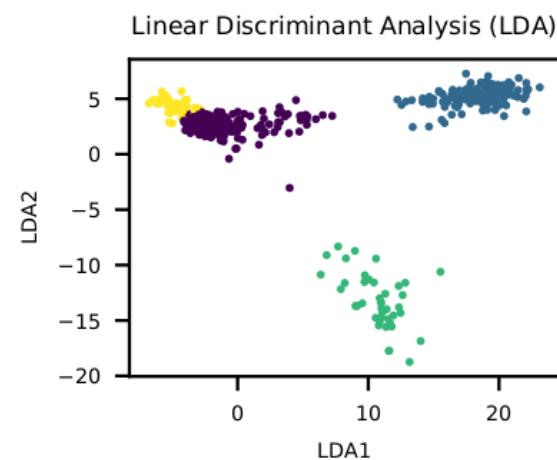
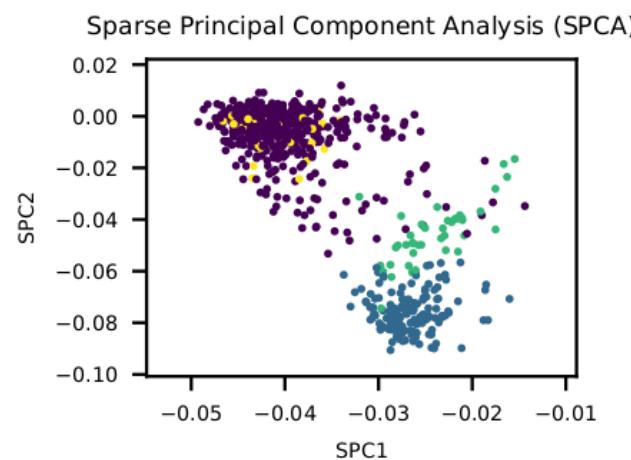
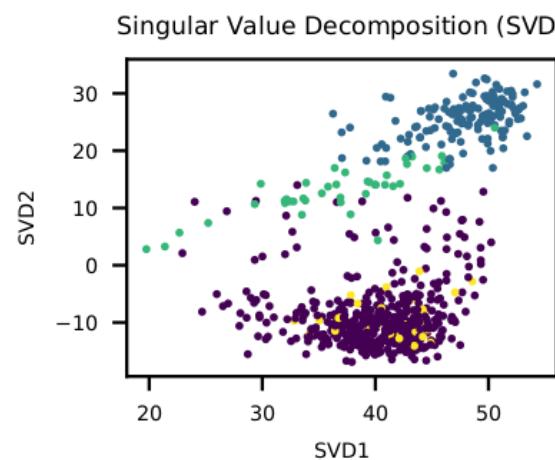
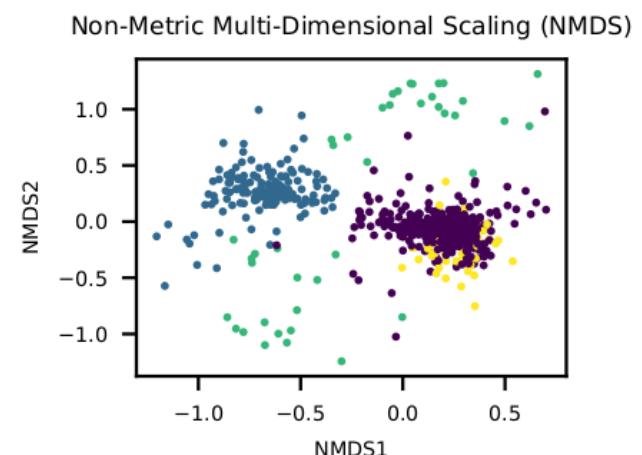
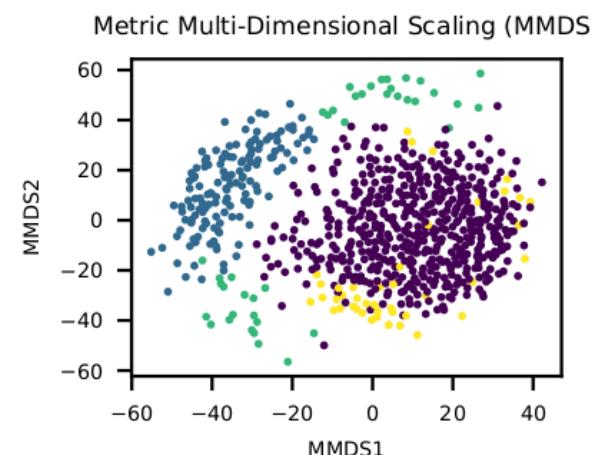
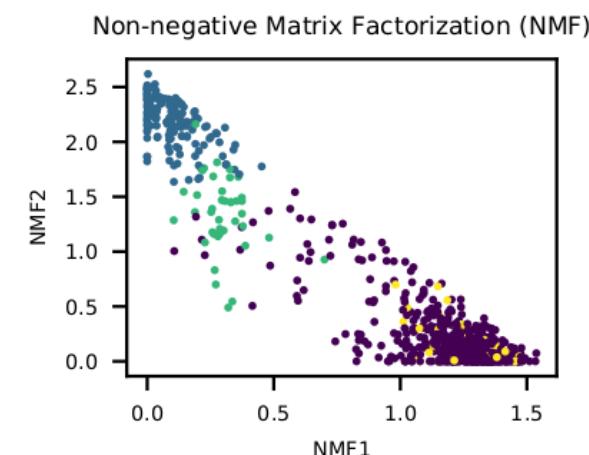
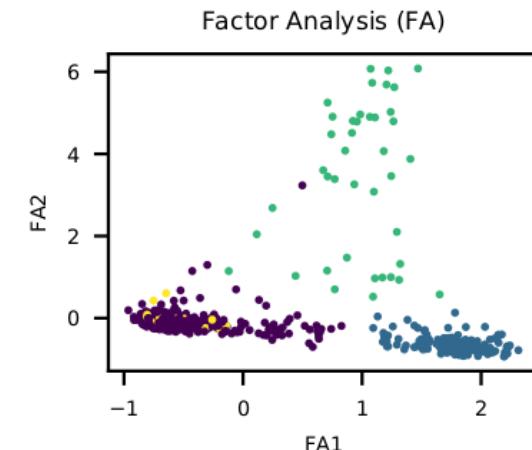
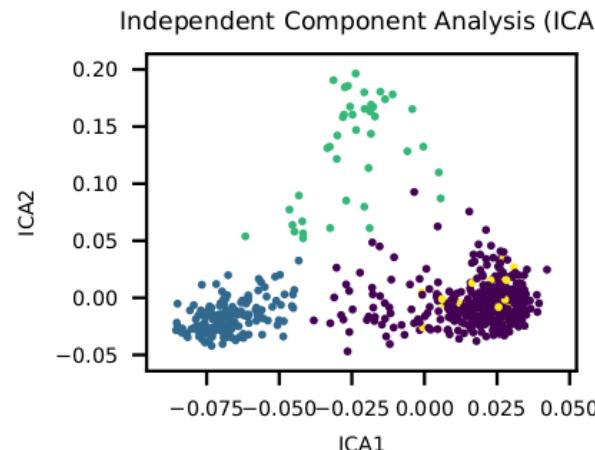
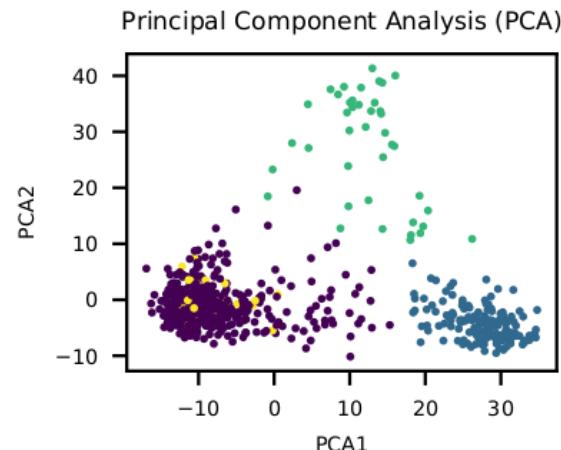
$$Y = \alpha + \beta X$$

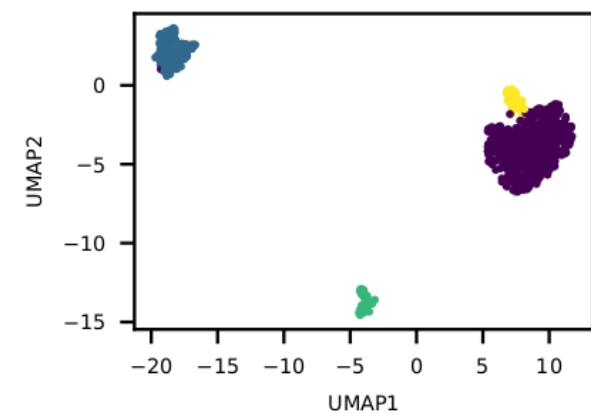
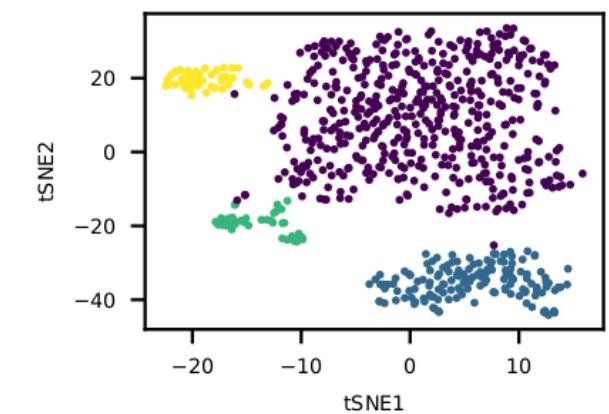
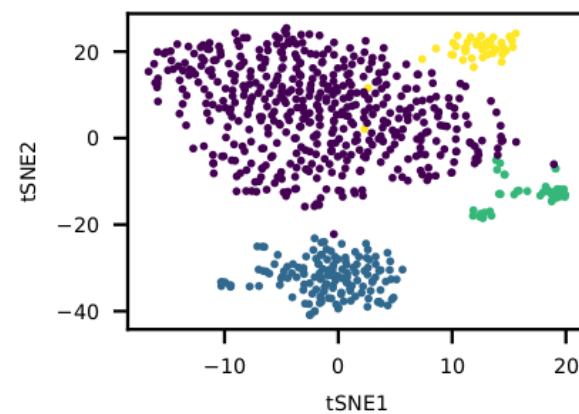
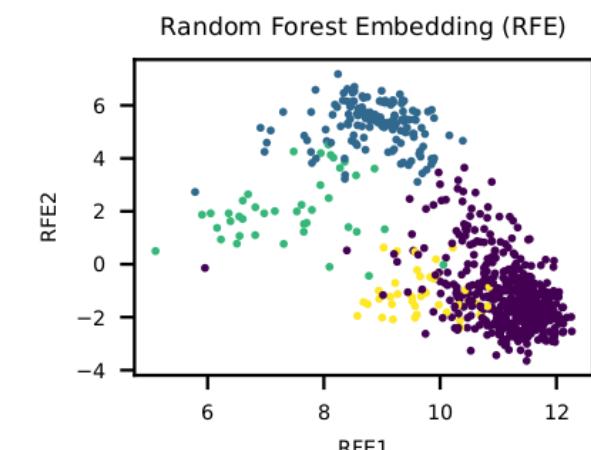
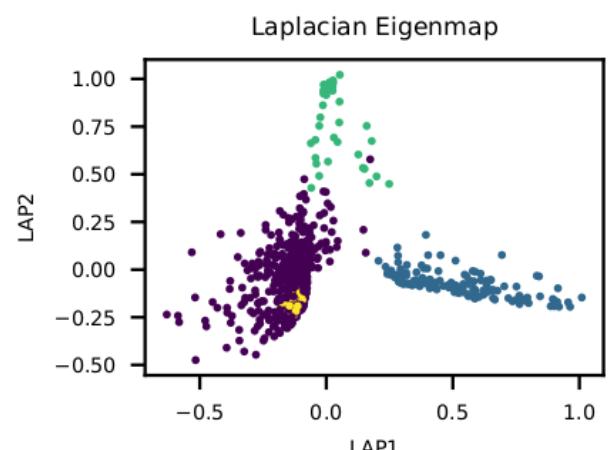
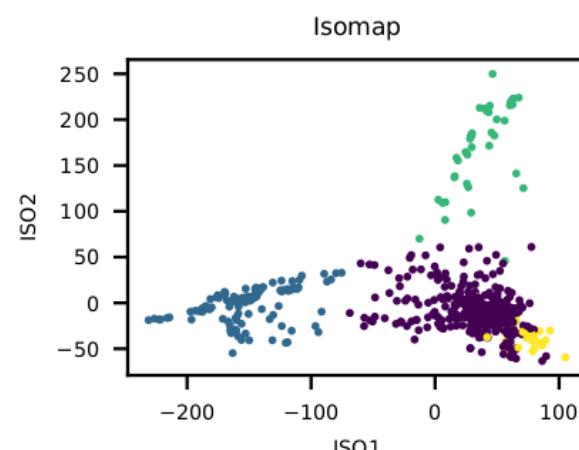
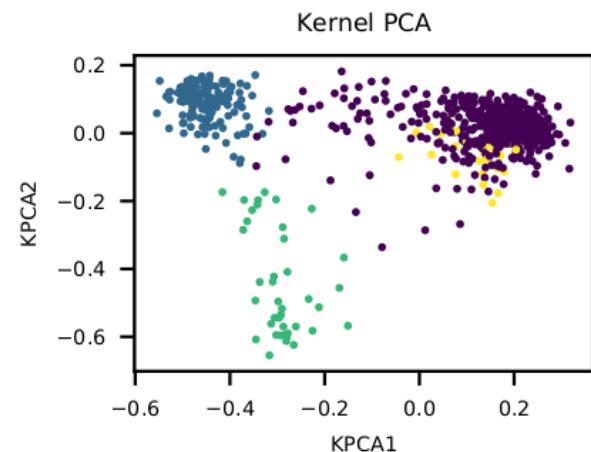
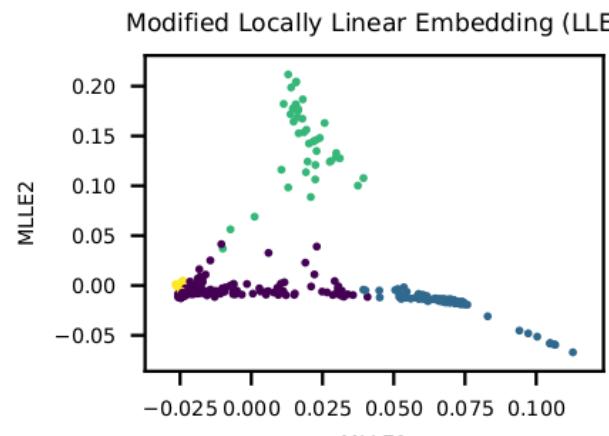
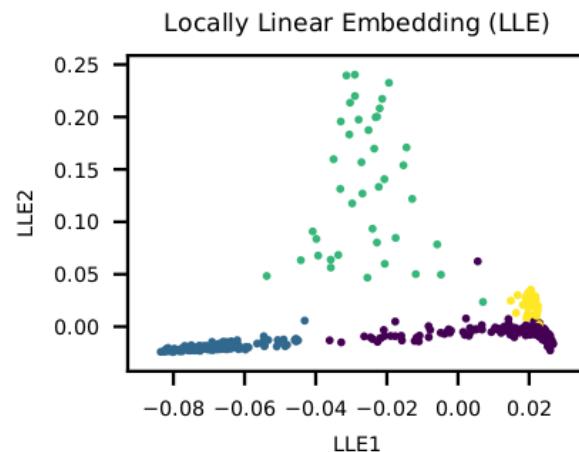
$$\beta = (X^T X)^{-1} X^T Y$$

$$(X^T X)^{-1} \sim \frac{1}{\det(X^T X)} \dots \rightarrow \infty, \quad n \ll p$$

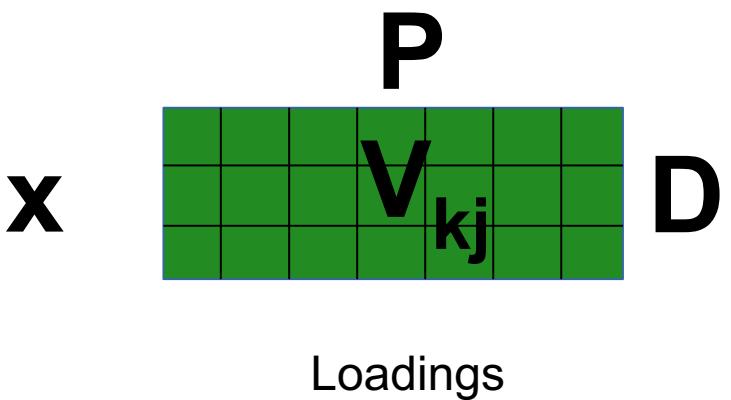
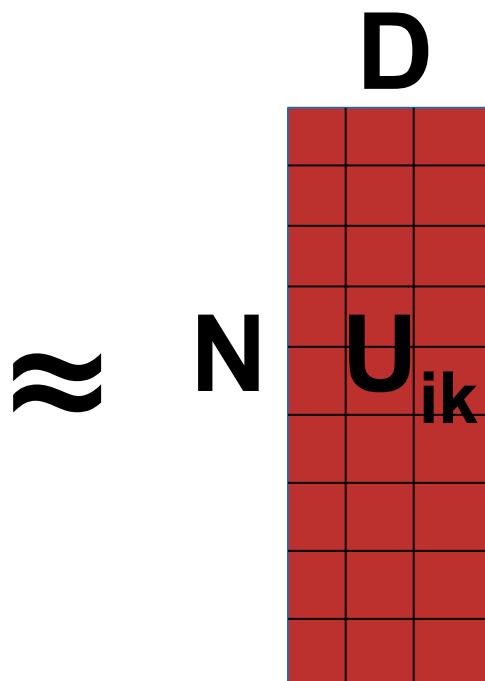
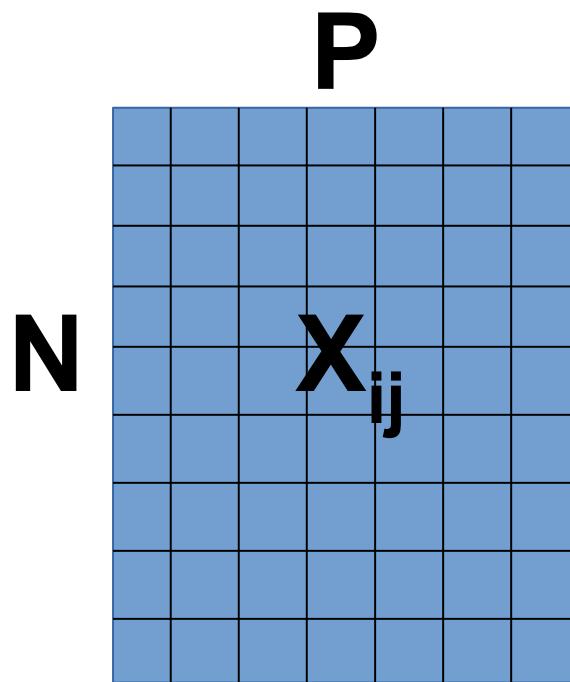
$$\text{Ex.2} \quad E[\hat{\sigma}^2] = \frac{n-p}{n} \sigma^2$$

Biased ML variance estimator in HD-space



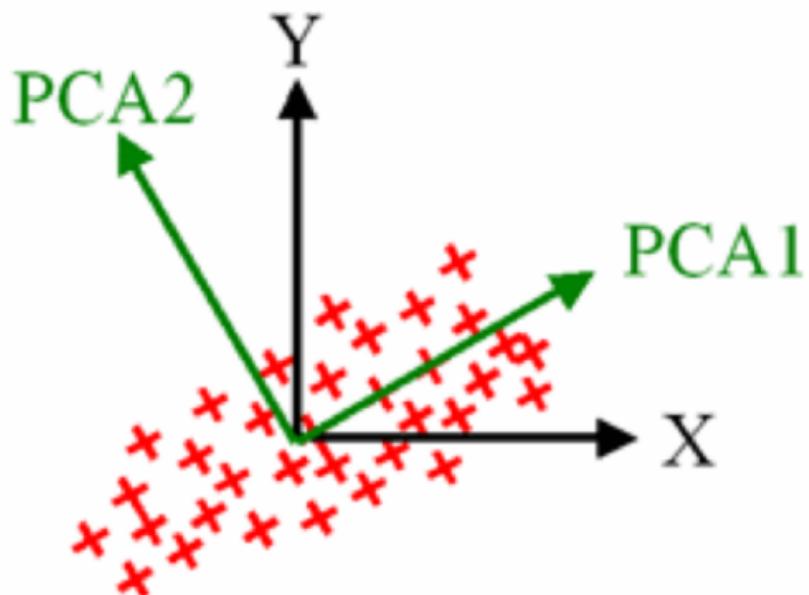


$$\mathbf{X}_{ij} \approx \mathbf{U}_{ik} \mathbf{V}_{kj}$$



$$\text{Loss} = \sum_{i=1}^N \sum_{j=1}^P (\mathbf{X}_{ij} - \mathbf{U}_{ik} \mathbf{V}_{kj})^2$$

- Collapse p features ($p \gg n$) to few latent features and keep variation
- Rotation and shift of coordinate system toward maximal variance
- PCA is an **eigen matrix decomposition** problem



$$PC = u^T X = X^T u$$

X is mean centered $\implies \langle PC \rangle = 0$

$$\langle (PC - \langle PC \rangle)^2 \rangle = \langle PC^2 \rangle = u^T X X^T u$$

$$X X^T = A$$

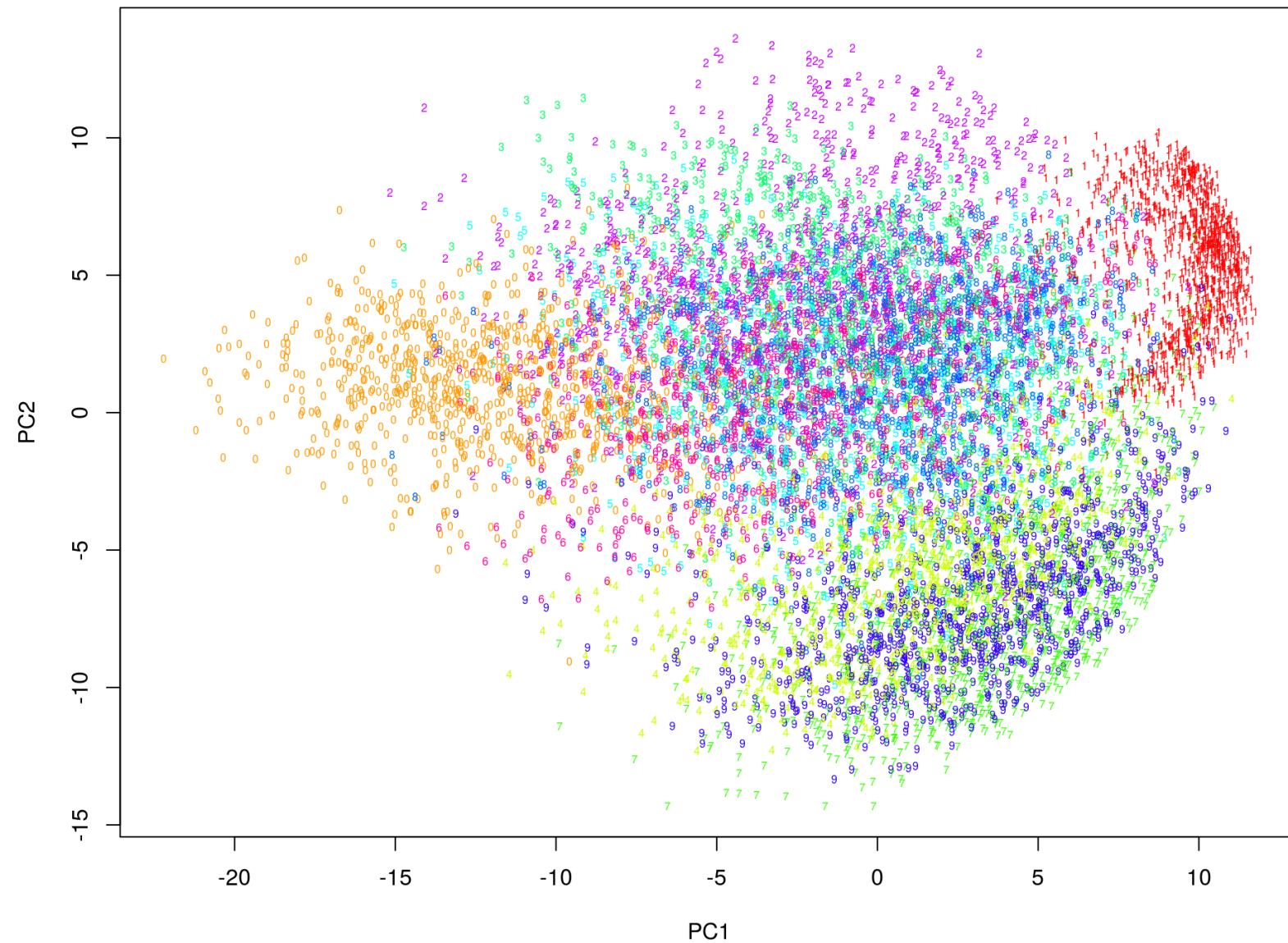
$$\langle PC^2 \rangle = u^T A u$$

A is **variance-covariance** of X

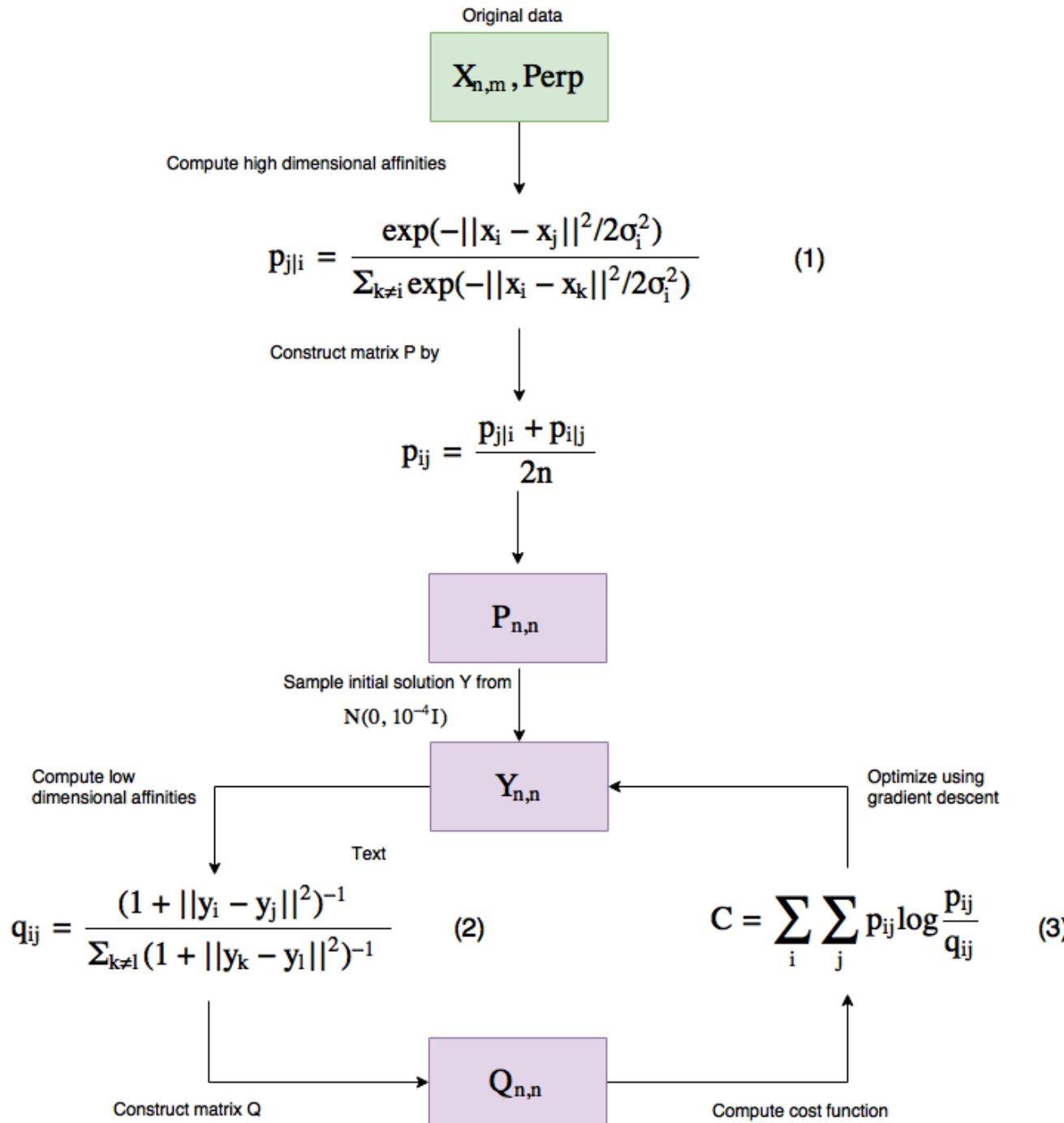
$$\max(u^T A u + \lambda(1 - u^T u)) = 0$$

$$A u = \lambda u$$

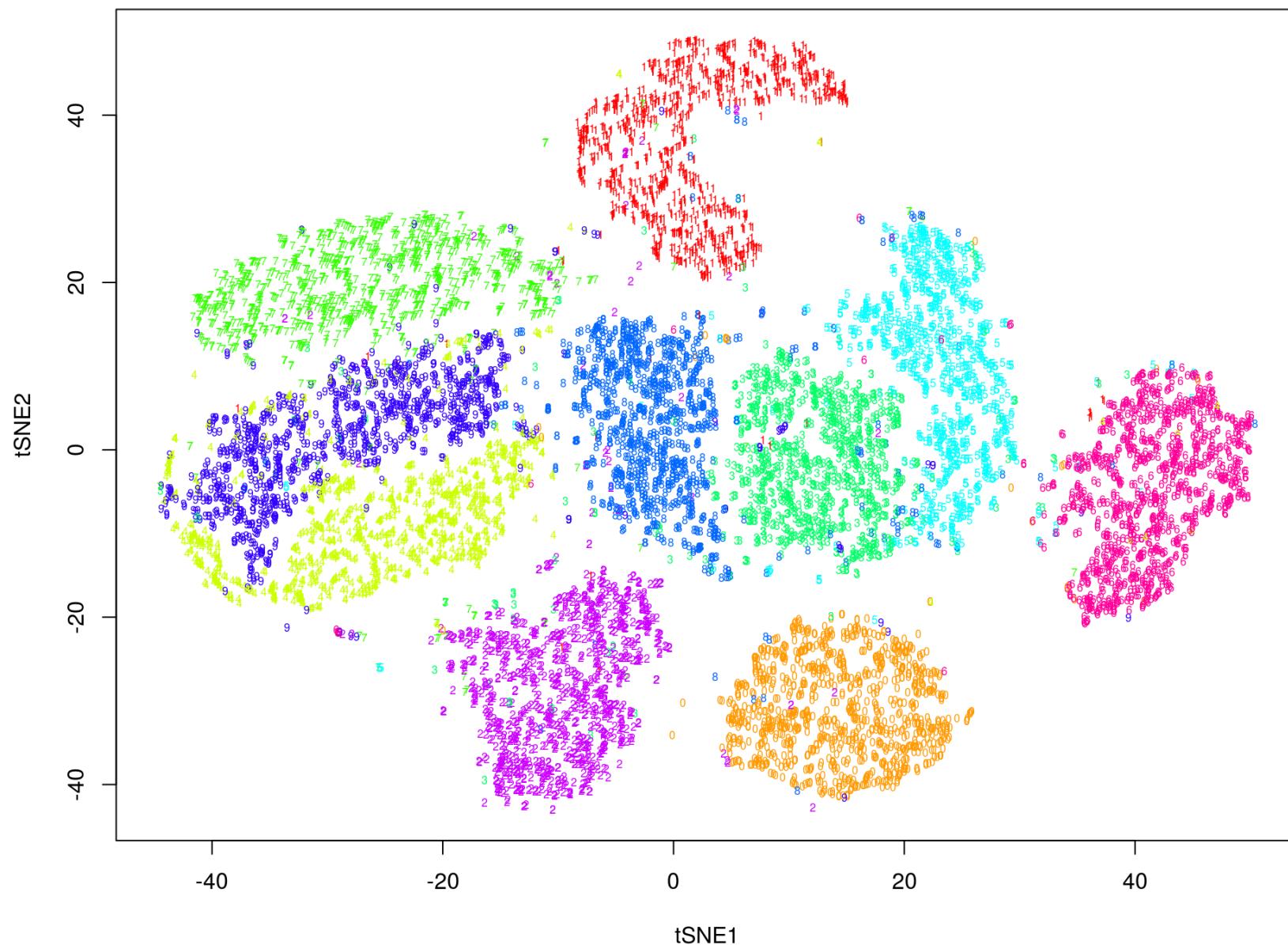
PCA PLOT WITH PRCOMP

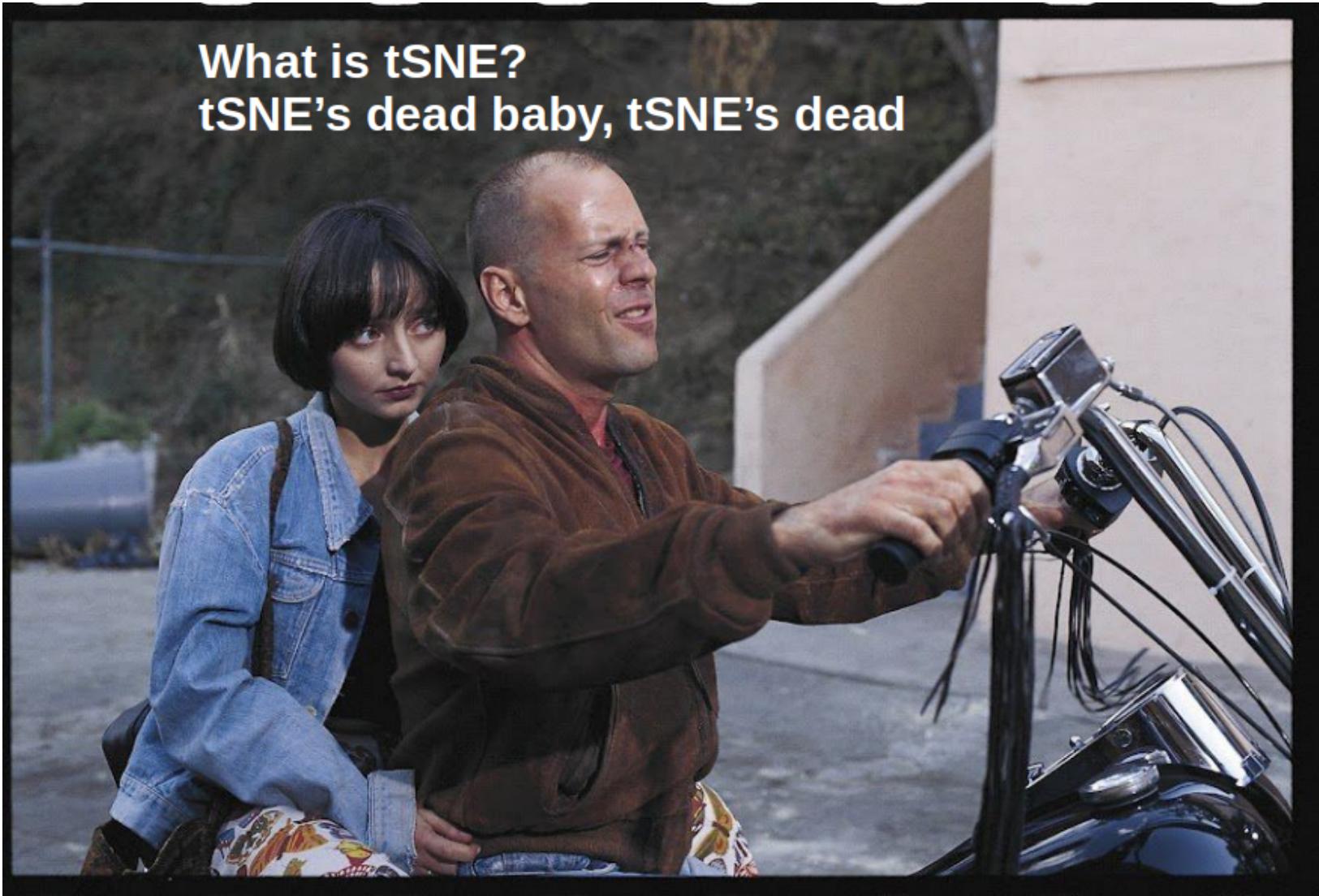


t-distributed Stochastic Neighbor Embedding (tSNE)



tSNE MNIST





Is UMAP really superior over tSNE?

tSNE does not scale for large data sets?

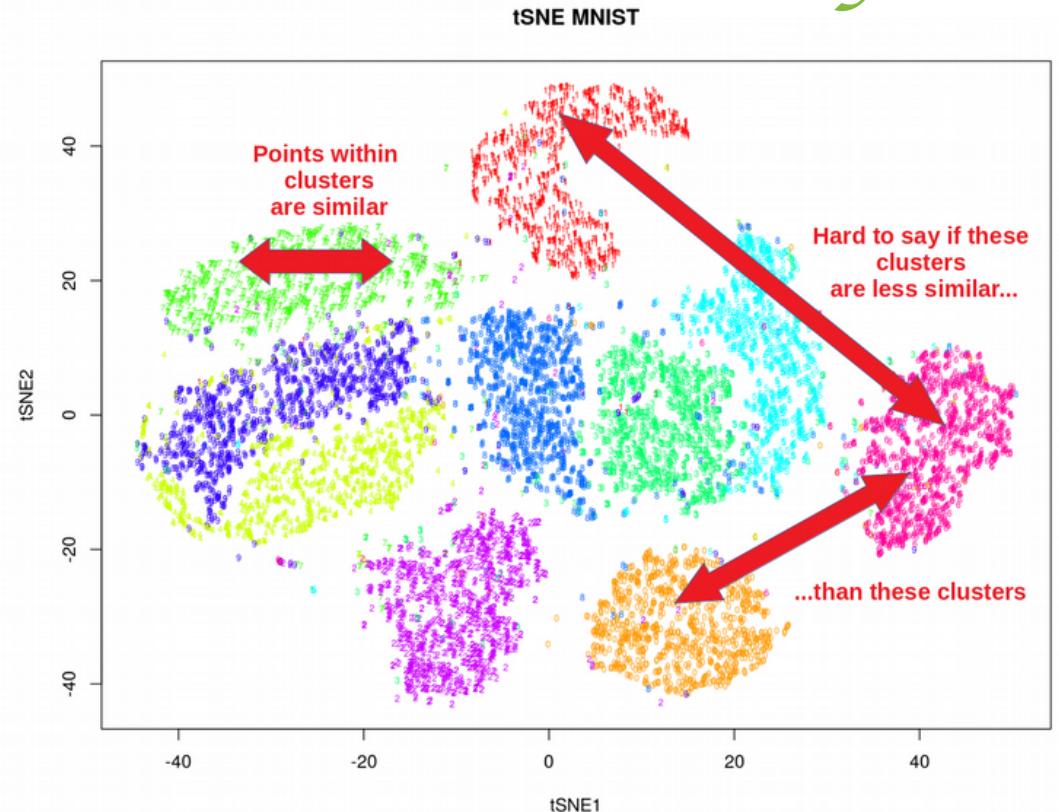
tSNE does not preserve global structure?

tSNE can only embed into 2-3 dims?

tSNE performs non-parametric mapping?

tSNE can not work with high-dimensional data directly (PCA needed)?

tSNE uses too much memory at large perplexities (FitSNE does not solve it)?

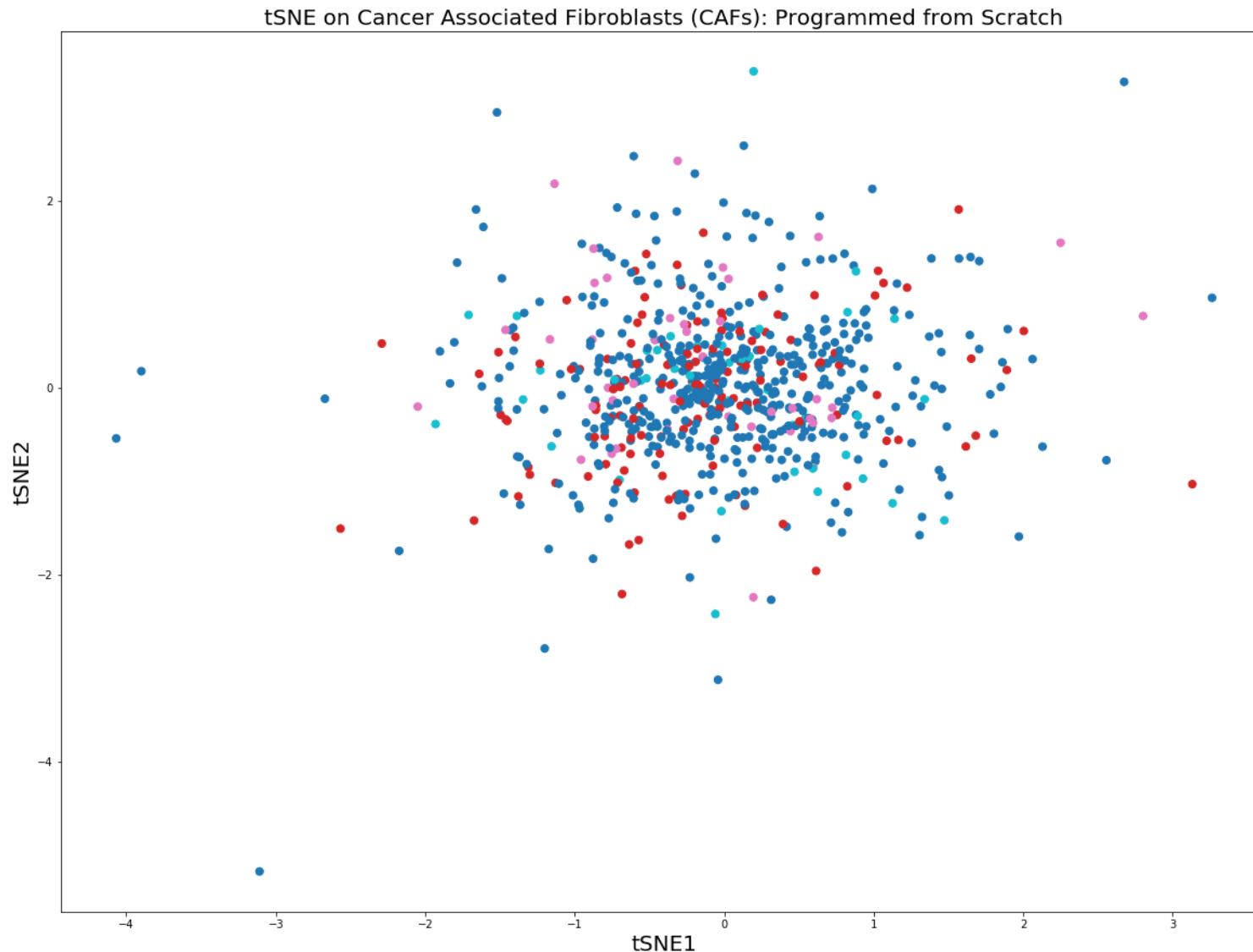


$$p_{j|i} = \frac{\exp\left(-\|x_i - x_j\|^2 / 2\sigma_i^2\right)}{\sum_{k \neq i} \exp\left(-\|x_i - x_k\|^2 / 2\sigma_i^2\right)}, \quad p_{ij} = \frac{p_{i|j} + p_{j|i}}{2N} \quad (1)$$

$$\text{Perplexity} = 2^{-\sum_j p_{j|i} \log_2 p_{j|i}} \quad (2)$$

$$q_{ij} = \frac{\left(1 + \|y_i - y_j\|^2\right)^{-1}}{\sum_{k \neq l} \left(1 + \|y_k - y_l\|^2\right)^{-1}} \quad (3)$$

$$KL(P_i || Q_i) = \sum_i \sum_j p_{j|i} \log \frac{p_{j|i}}{q_{j|i}}, \quad \frac{\partial KL}{\partial y_i} = 4 \sum_j (p_{ij} - q_{ij})(y_i - y_j) \left(1 + \|y_i - y_j\|^2\right)^{-1} \quad (4)$$



UMAP uses local connectivity for high-dim probabilities

$$p_{i|j} = e^{-\frac{d(x_i, x_j) - \rho_i}{\sigma_i}}$$

UMAP does not normalize probabilities (speed-up)

UMAP uses slightly different expression for nearest neighbors

$$k = 2 \sum_i p_{ij}$$

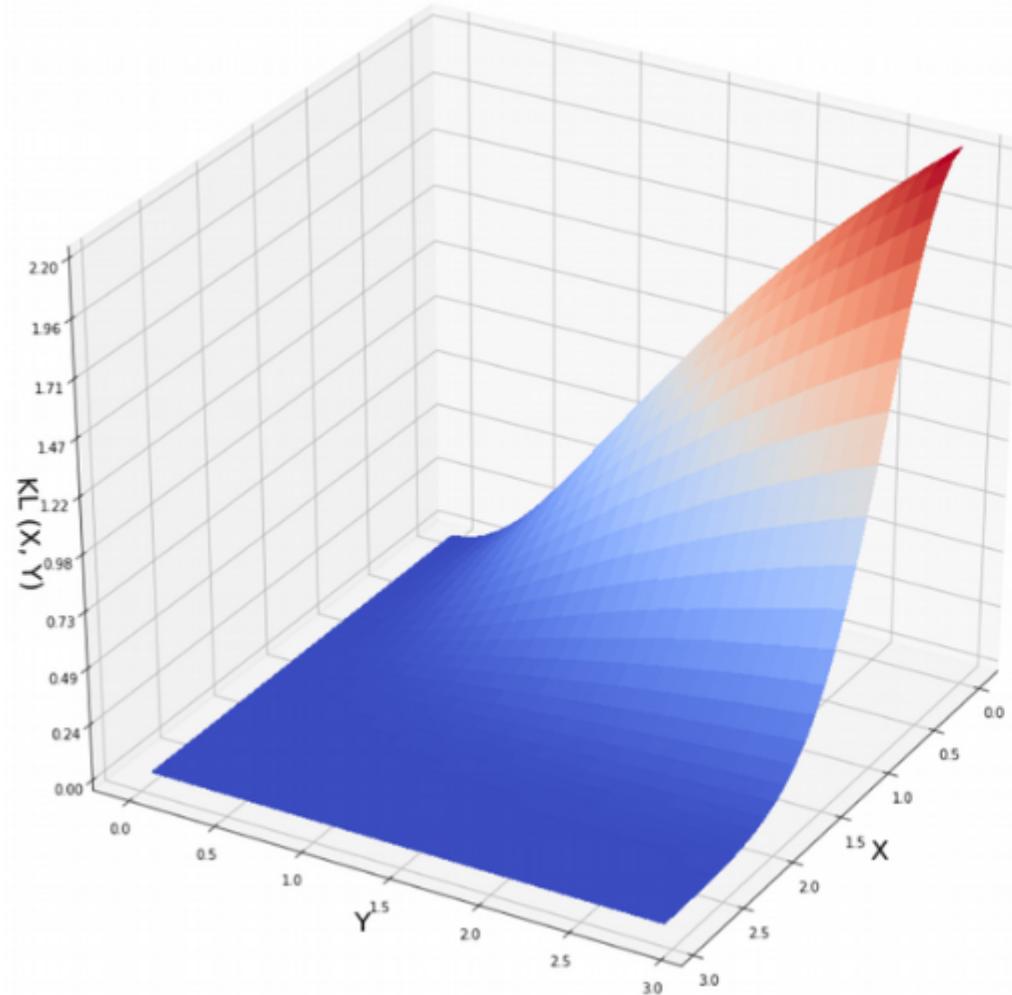
UMAP uses Laplacian Eigenmap for initialization

UMAP uses Cross-Entropy (not KL) as cost function

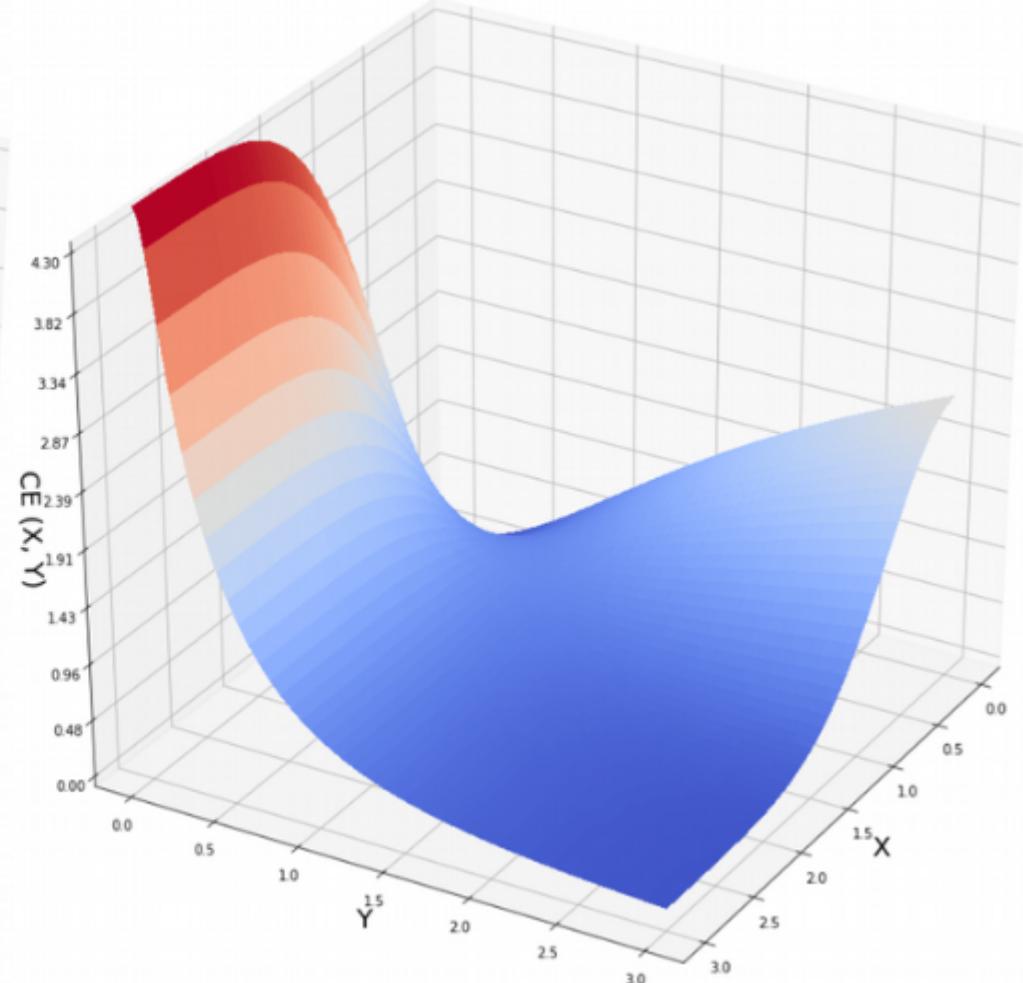
$$CE(X, Y) = \sum_i \sum_j \left[p_{ij}(X) \log \left(\frac{p_{ij}(X)}{q_{ij}(Y)} \right) + (1 - p_{ij}(X)) \log \left(\frac{1 - p_{ij}(X)}{1 - q_{ij}(Y)} \right) \right]$$

This is similar to tSNE cost function

This term is UMAP specific

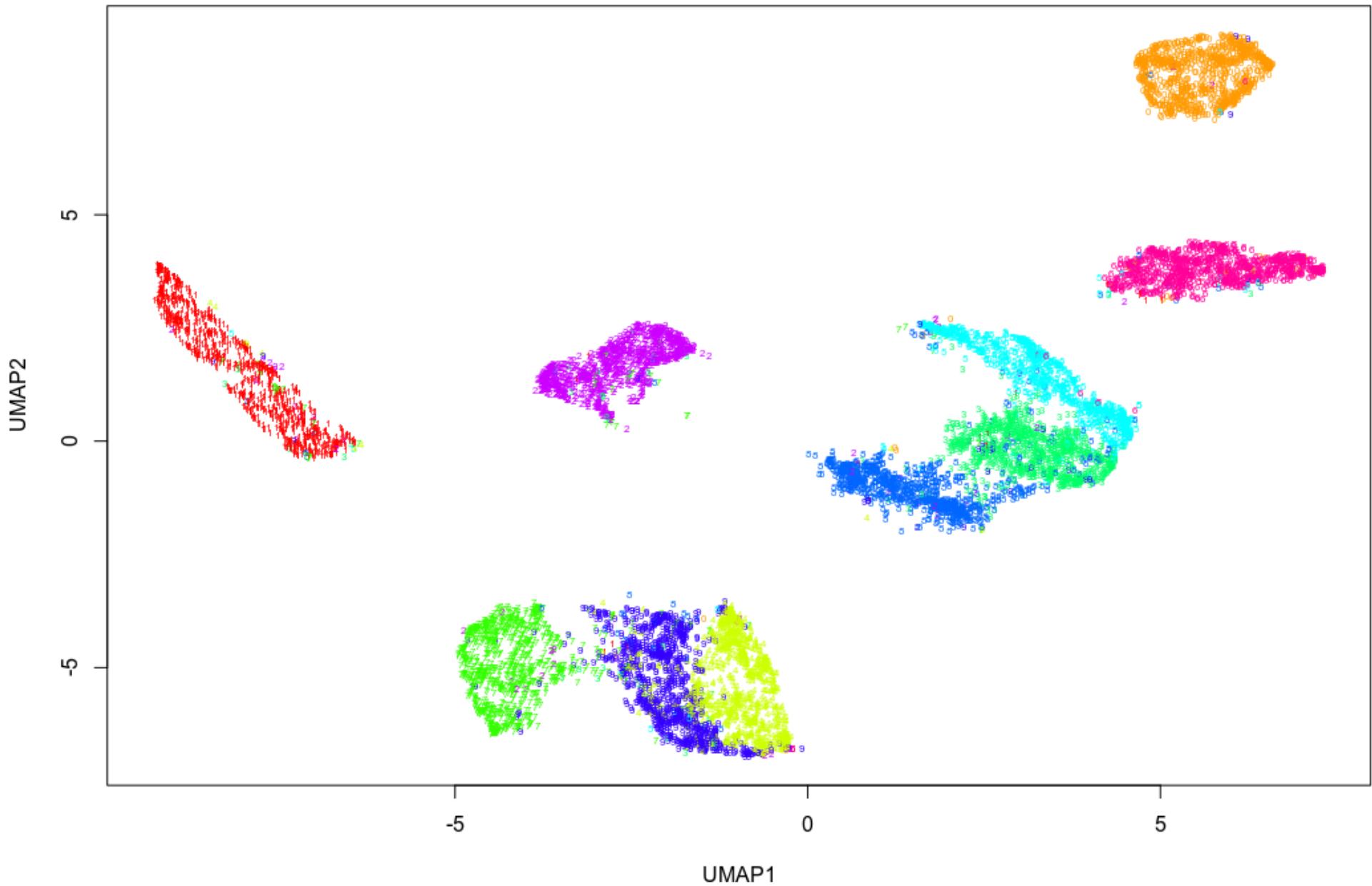


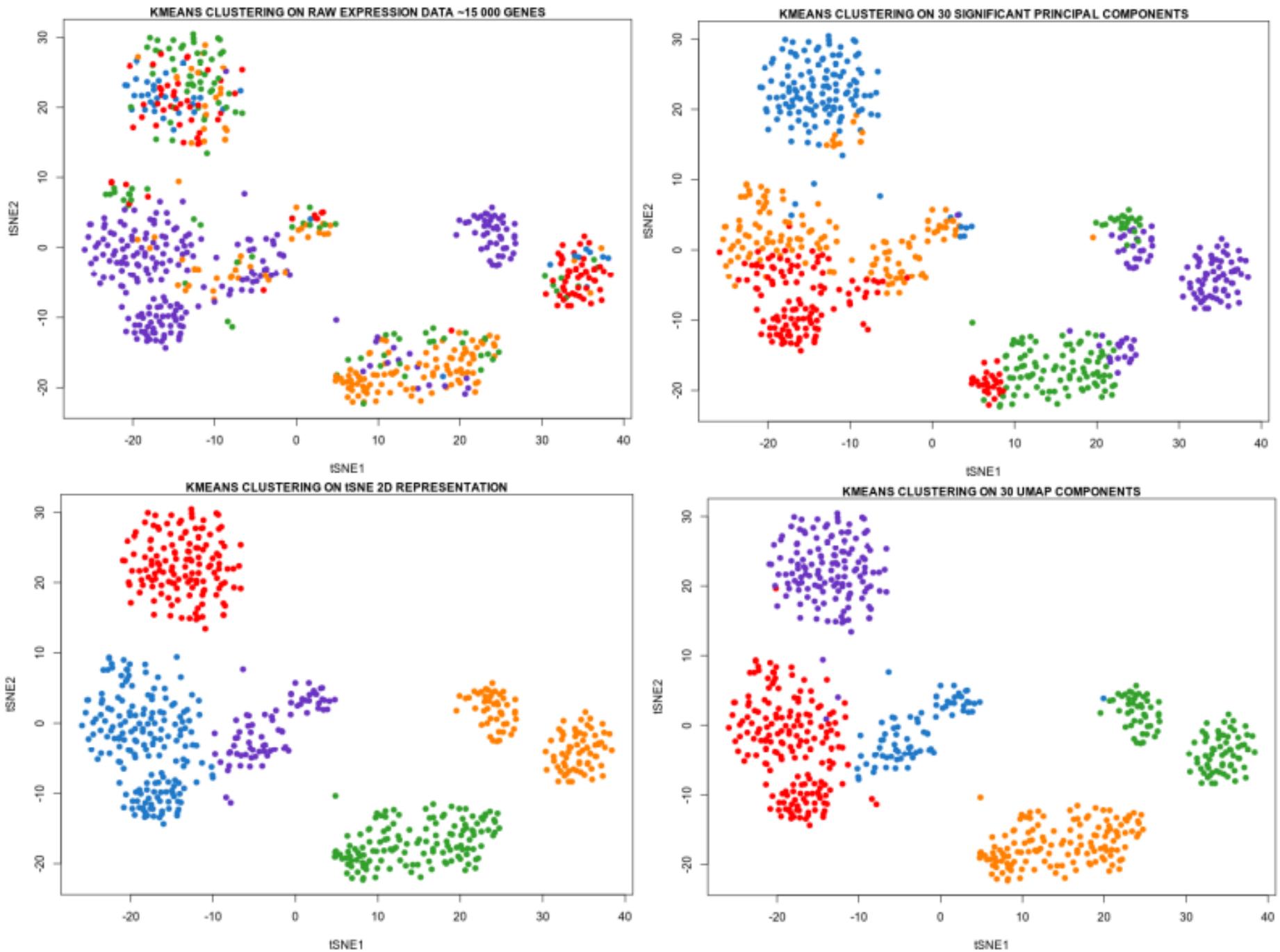
$X \rightarrow \infty$, Y can be any

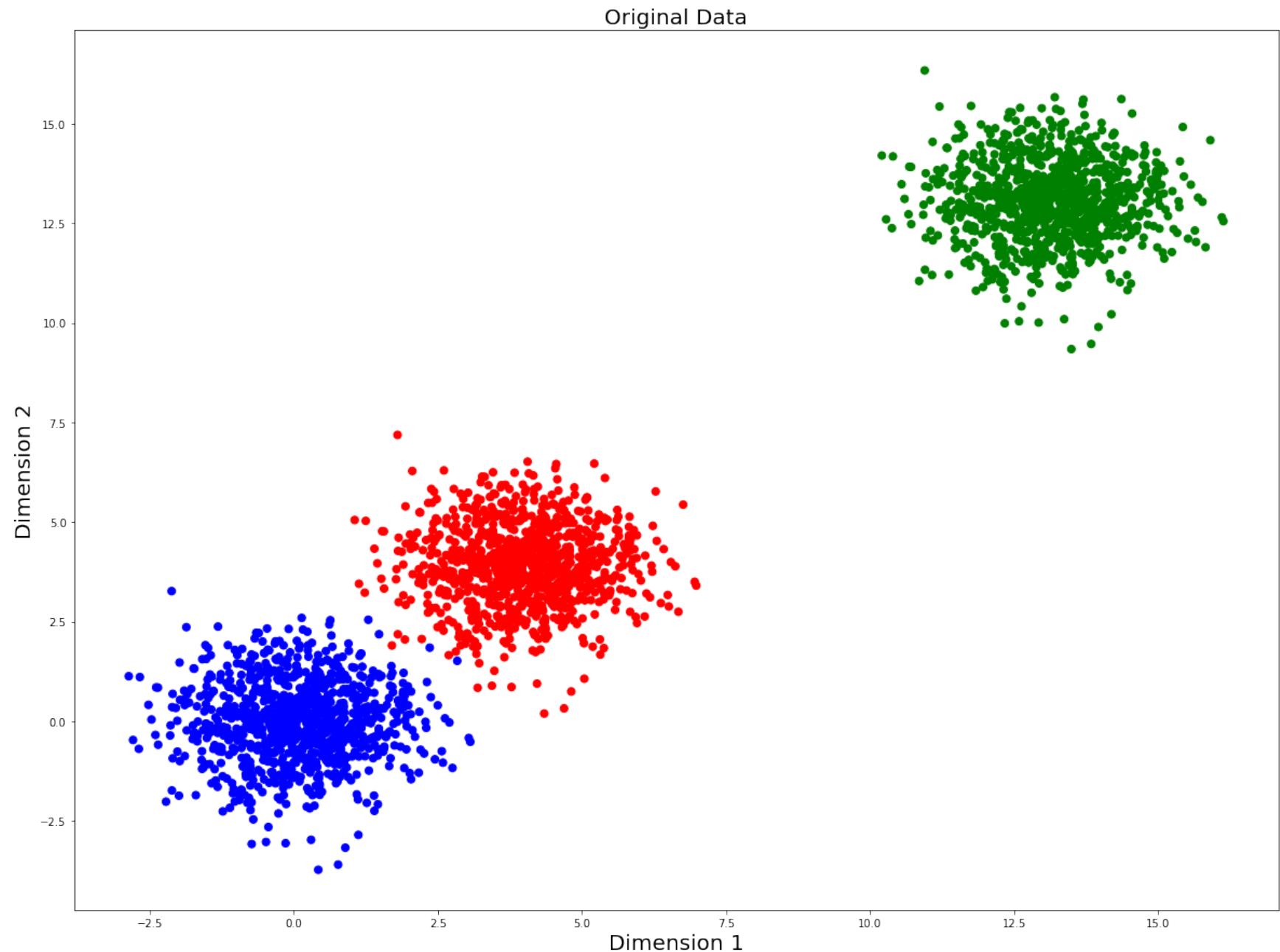


$X \rightarrow \infty$, $Y \rightarrow \infty$

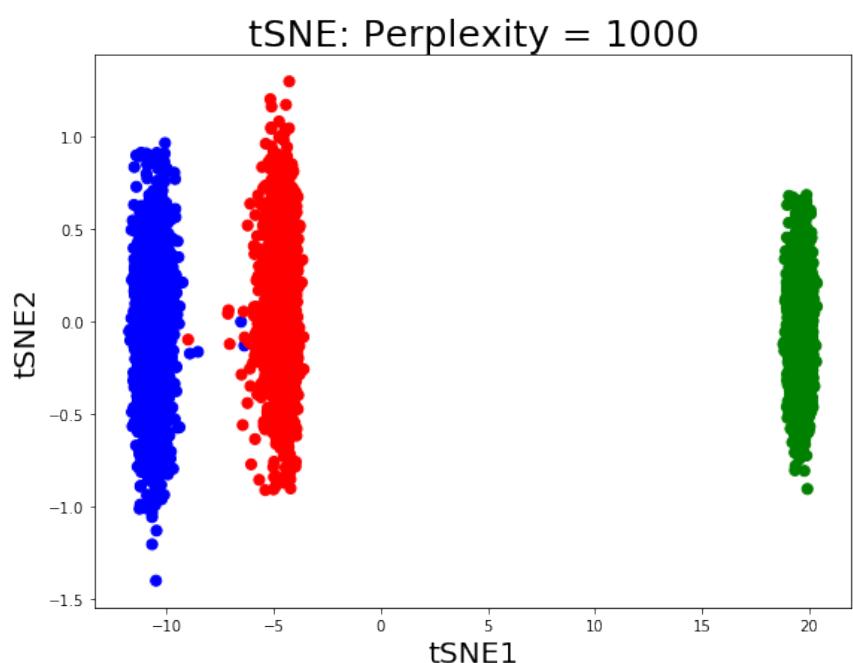
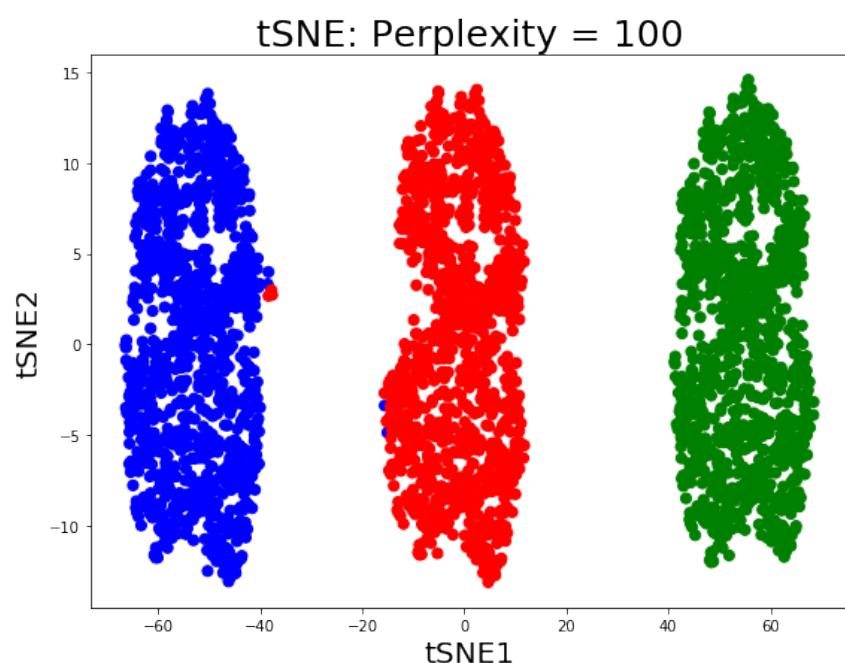
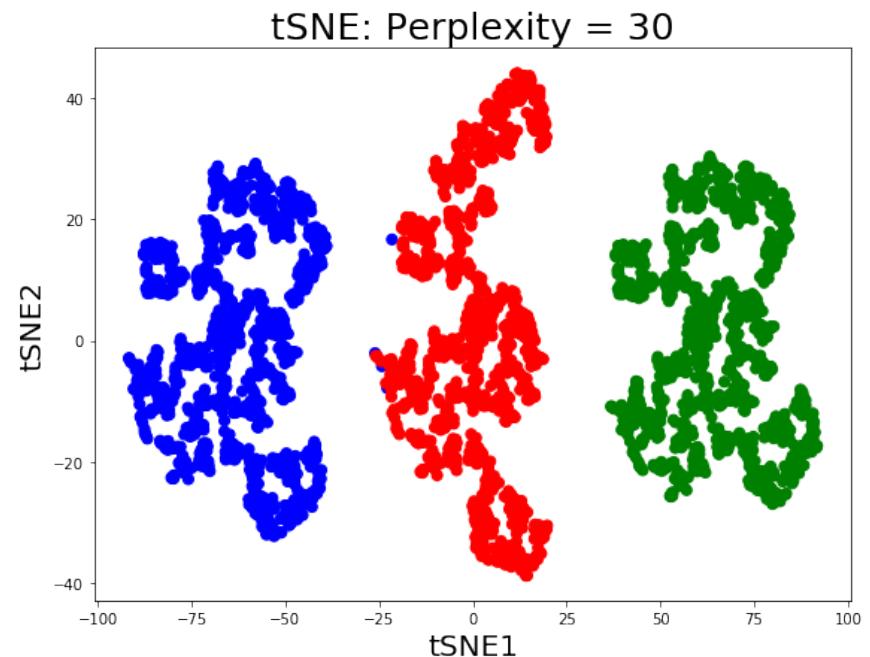
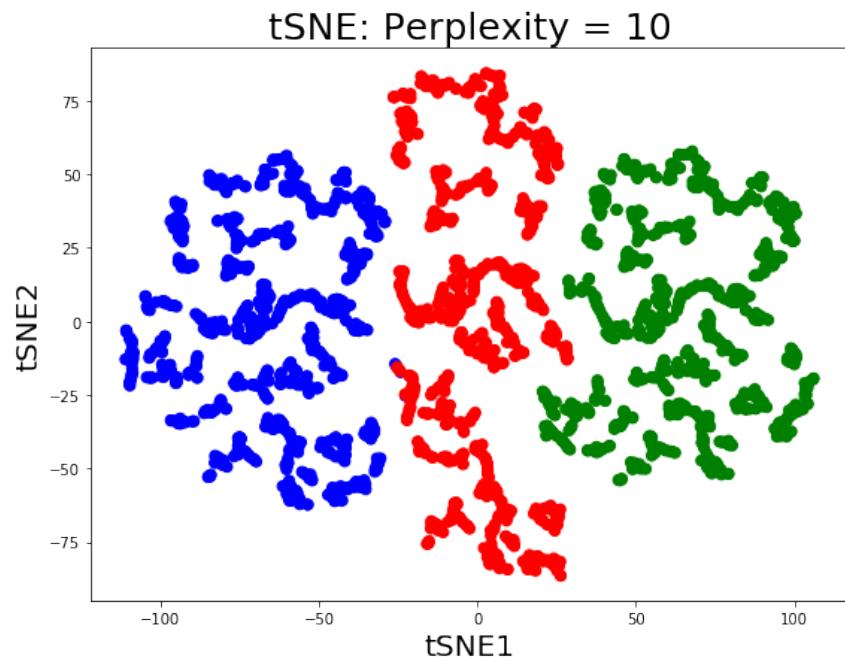
UMAP MNIST

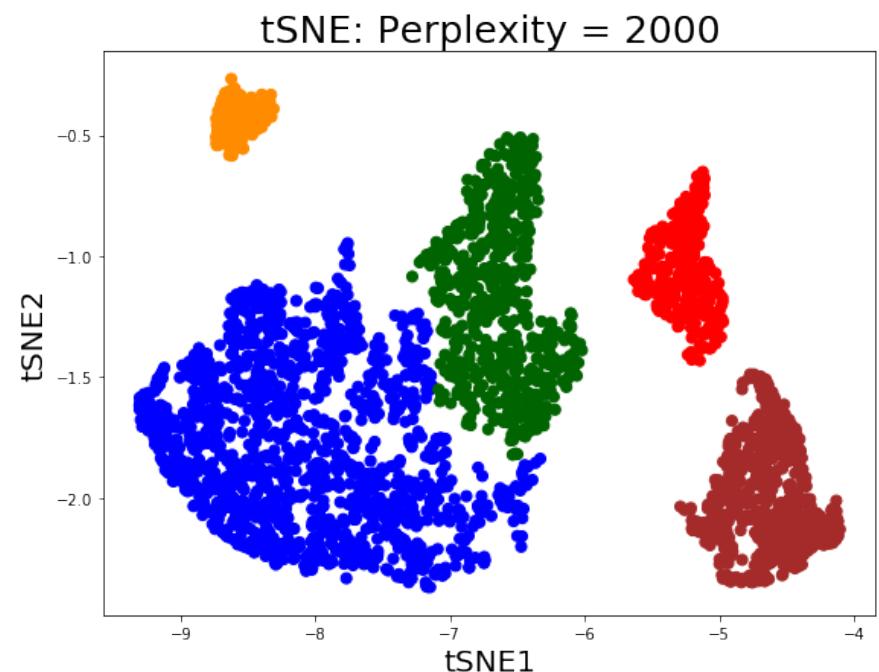
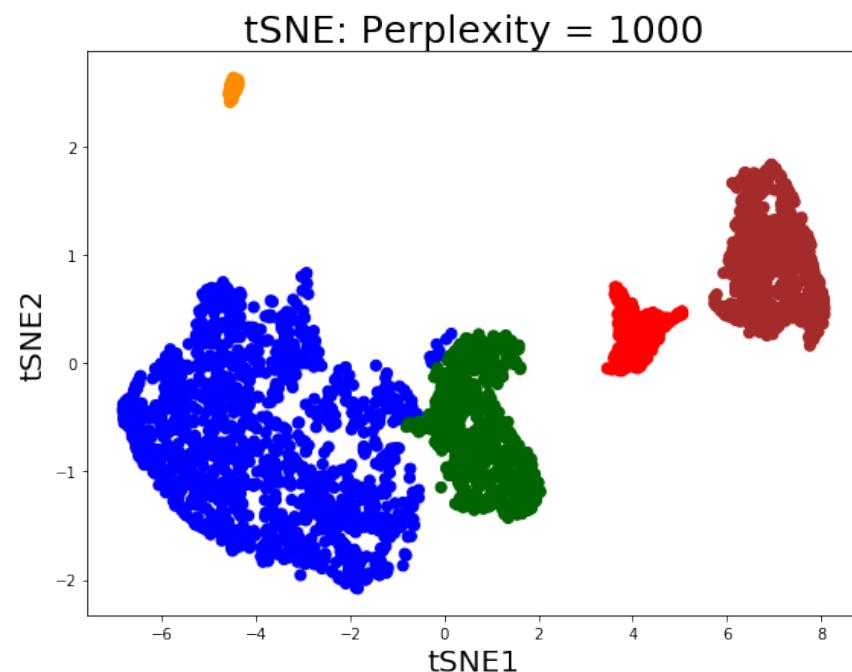
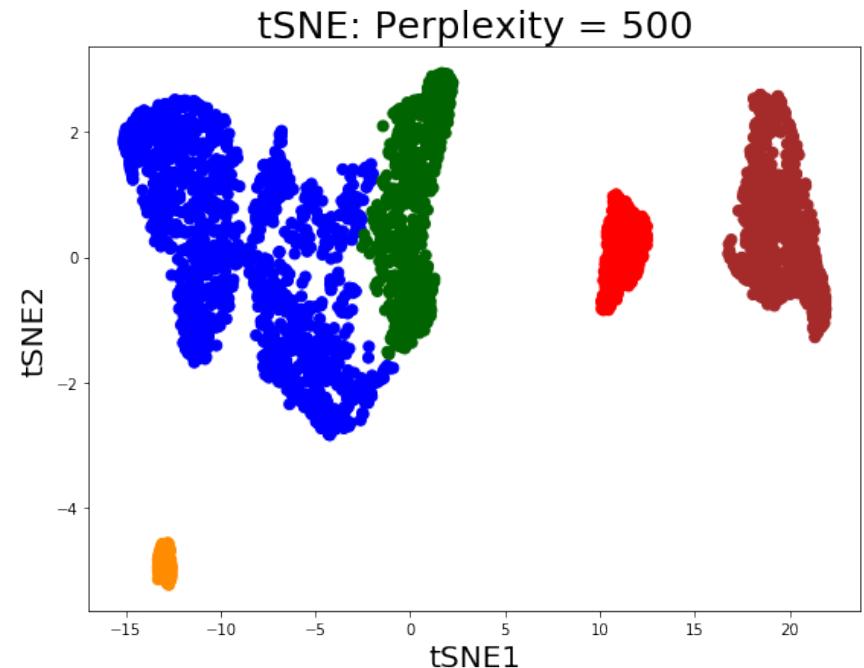
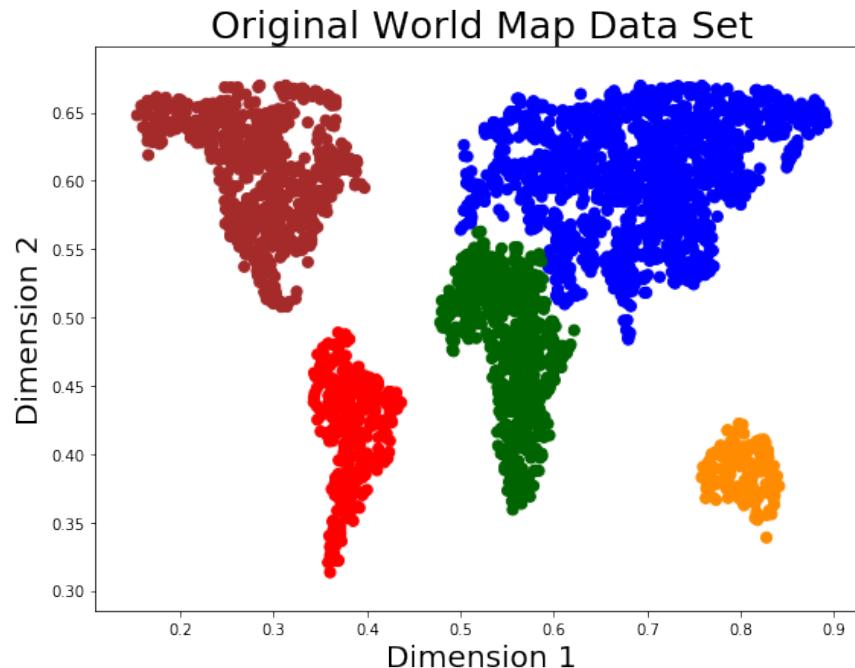






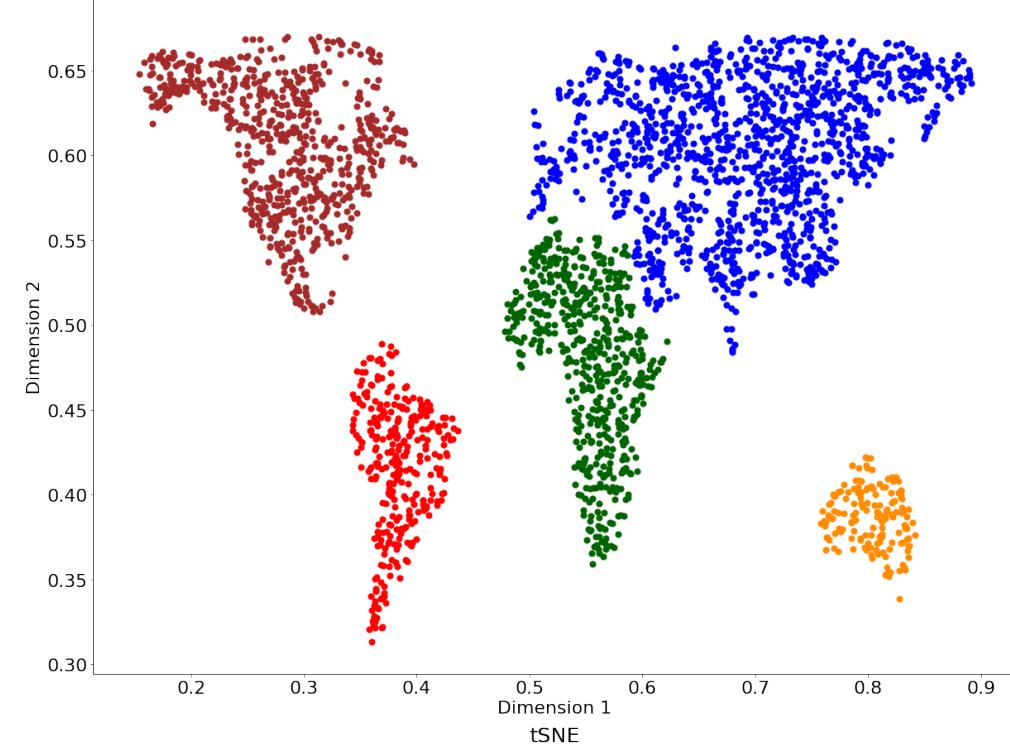
Why Global Structure Preservation is Important



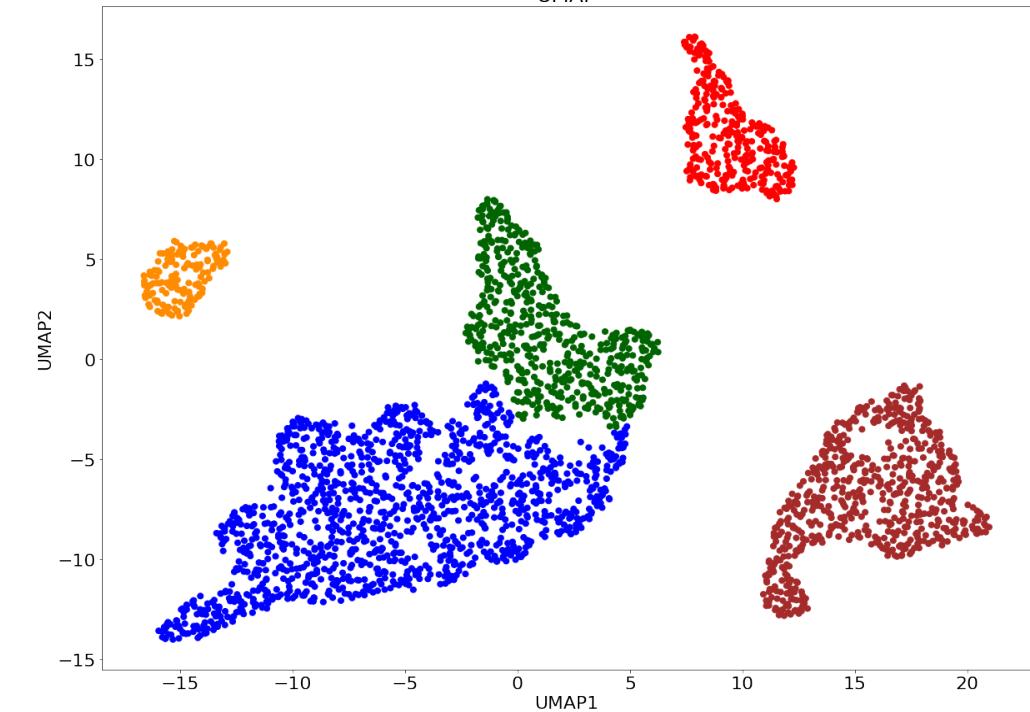
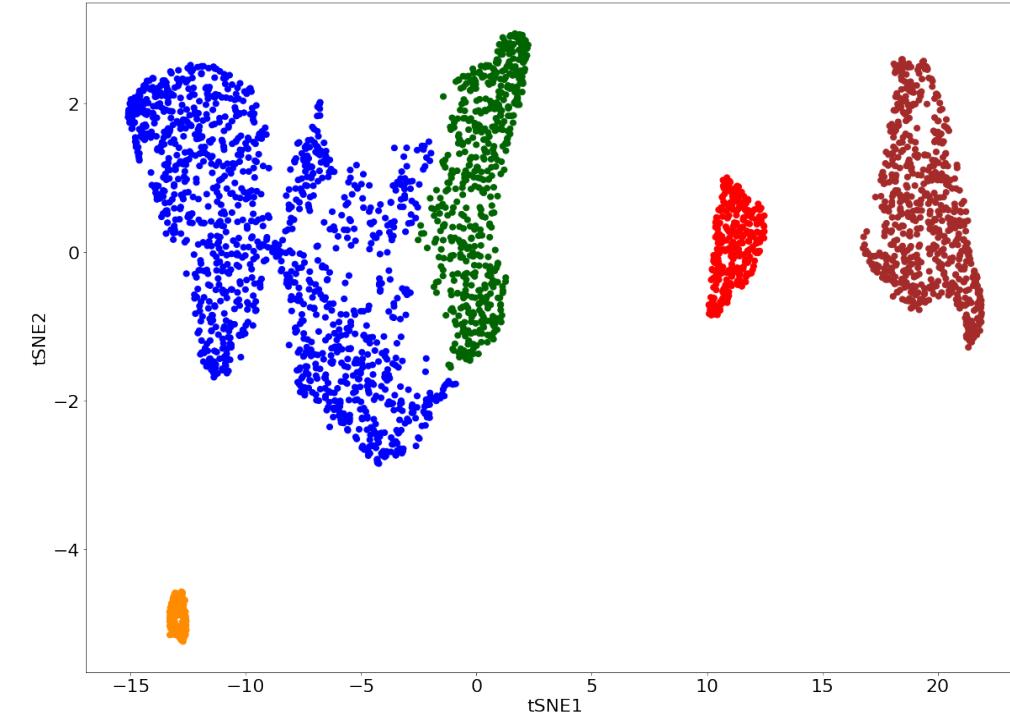
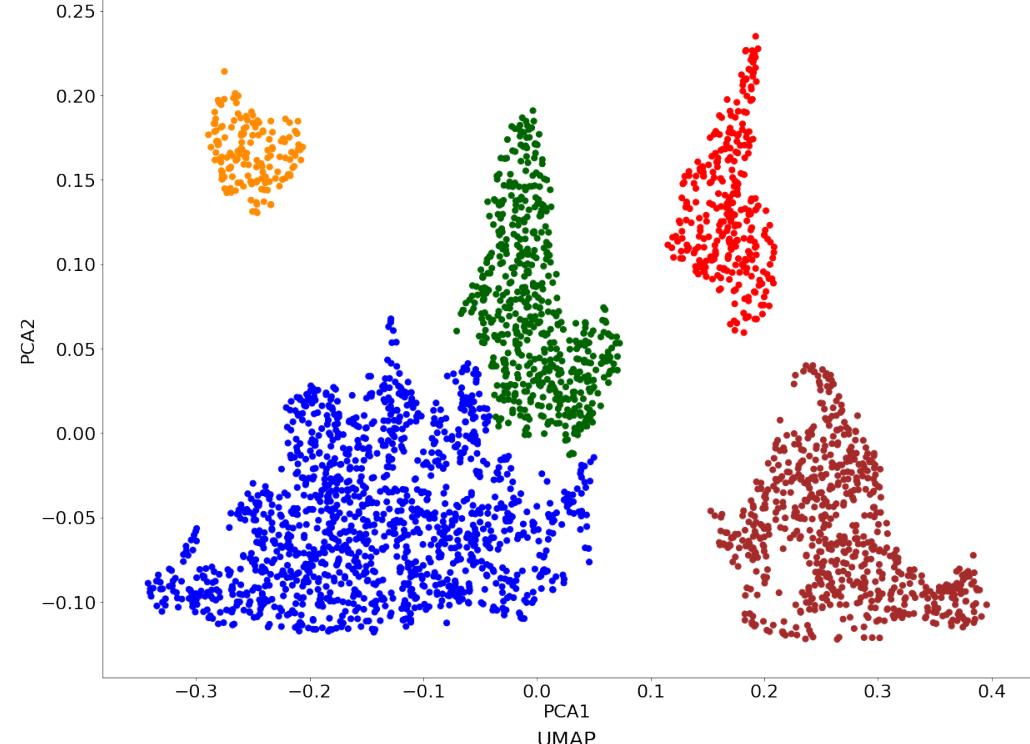


Global Structure Preservation for Linear Manifold

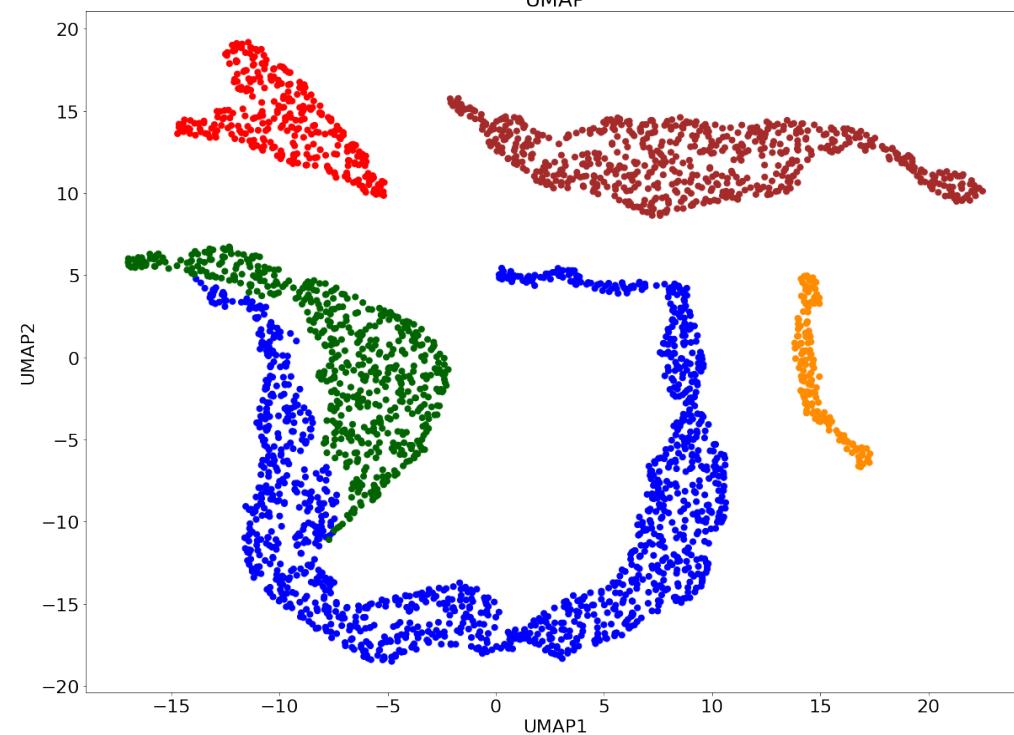
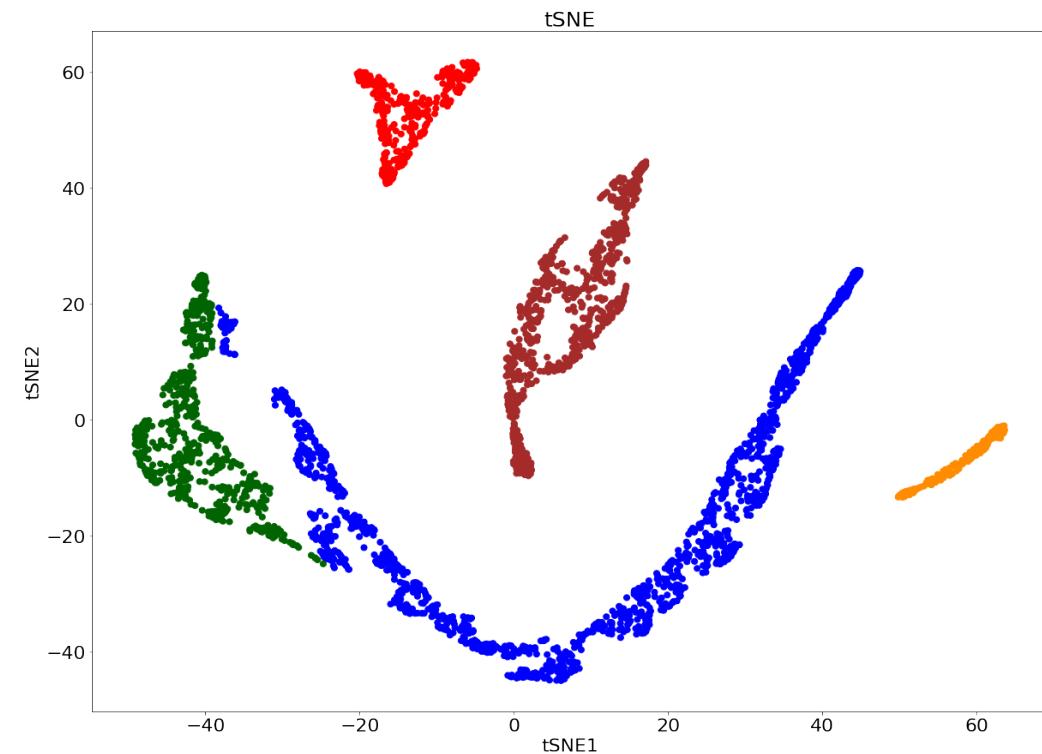
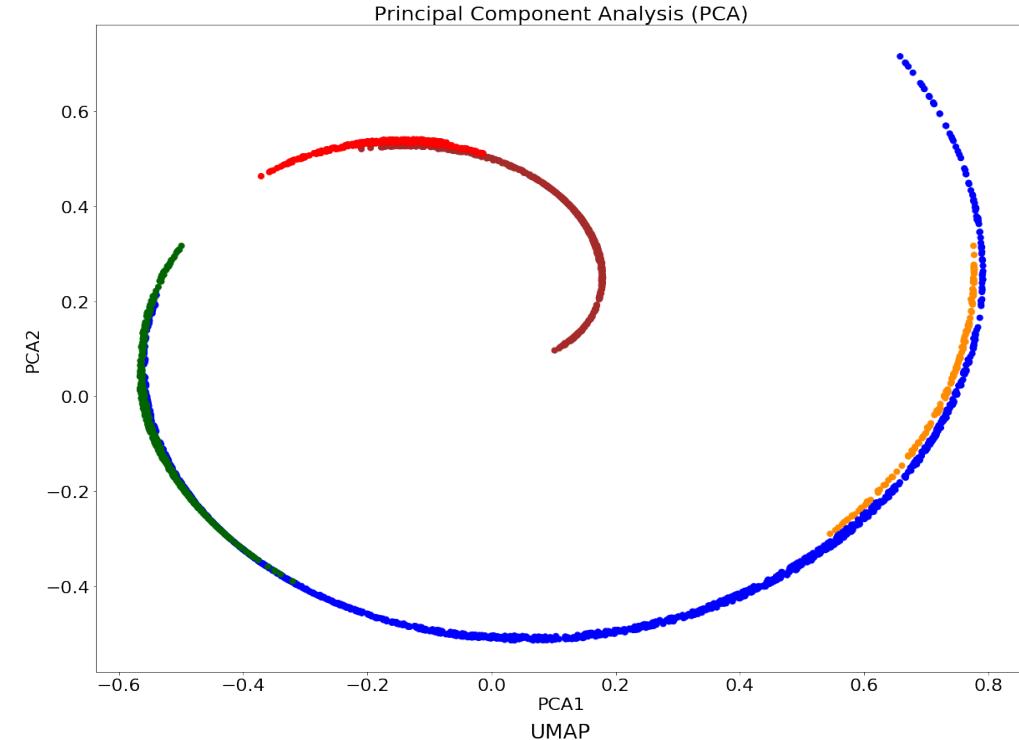
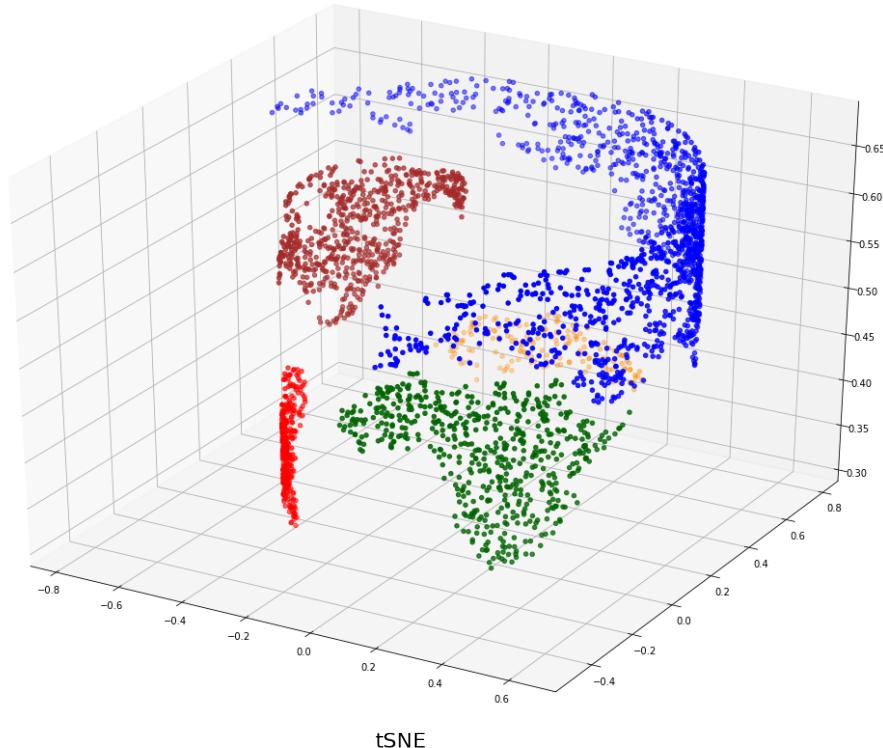
Original World Map Data Set



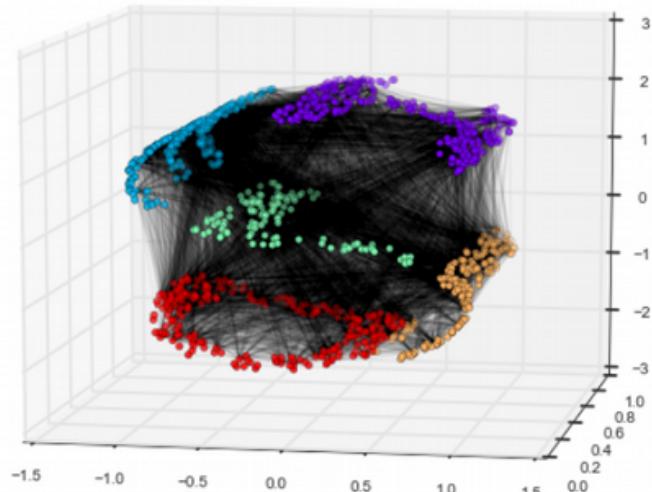
Principal Component Analysis (PCA)



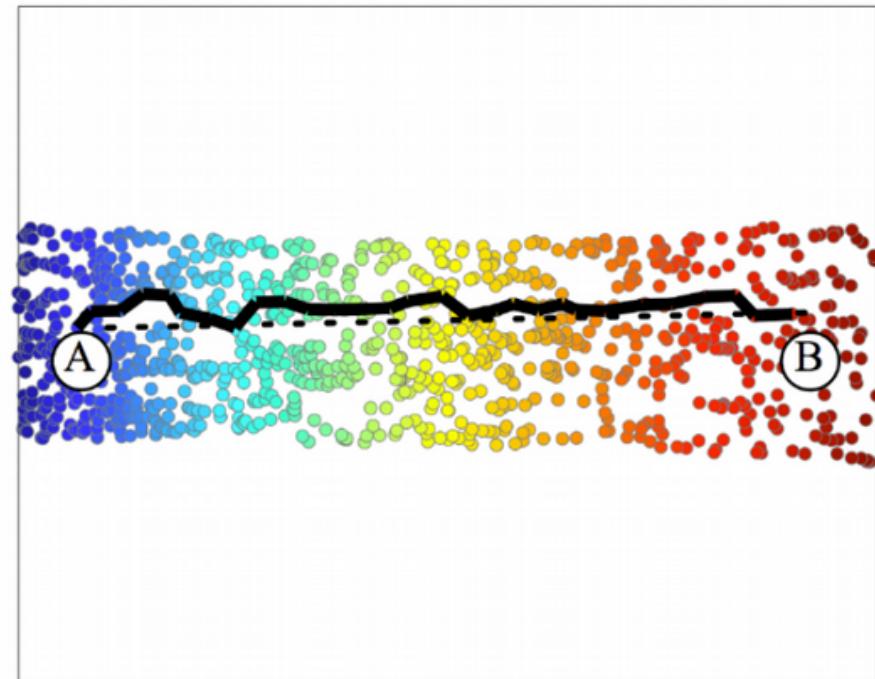
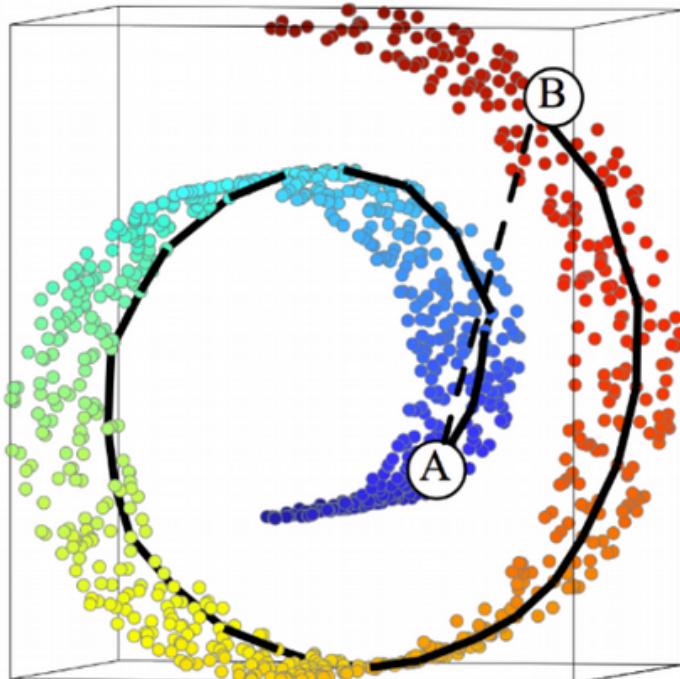
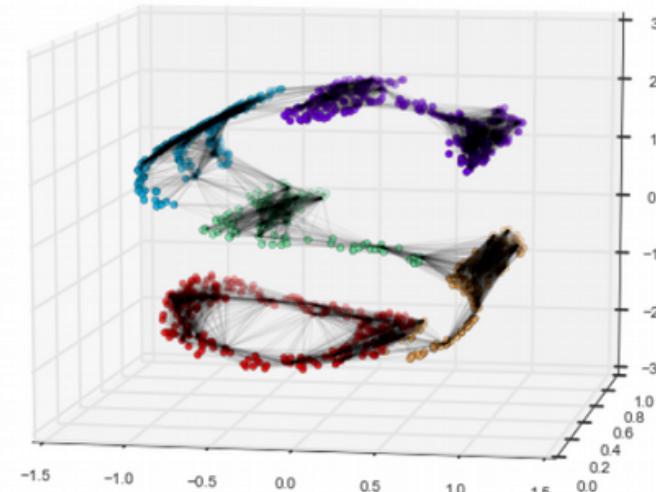
Global Structure Preservation for Non-Linear Manifold



MDS Linkages



LLE Linkages (100 NN)





- Watson, initialization makes UMAP superior over tSNE

- But Holmes ... did you think about KL-divergence?

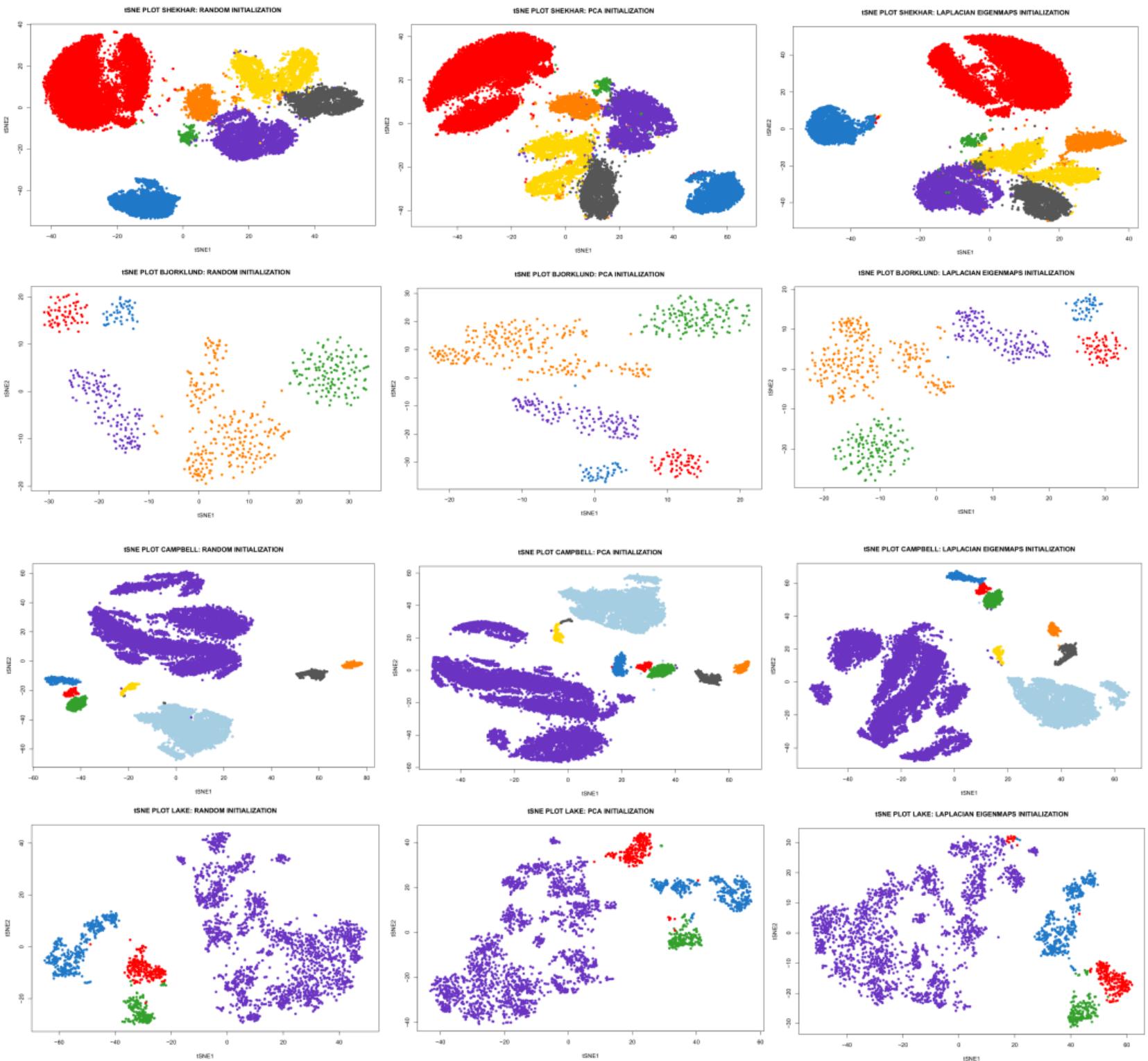
Initialization: PCA

Cost function: KL / CE

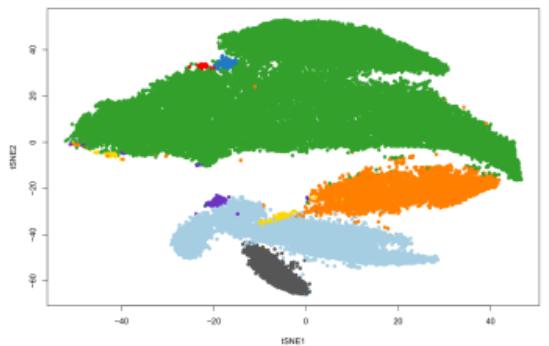
$$y_i = y_i - \mu \frac{\partial \text{Cost}}{\partial y_i}$$

Final embedding

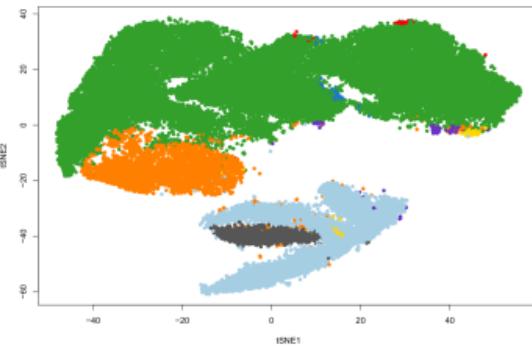
Learning rate



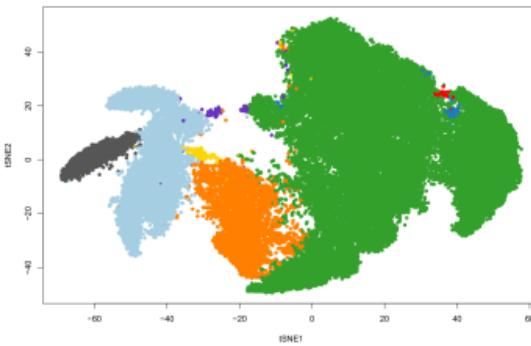
tSNE PLOT MACOSKO: RANDOM INITIALIZATION



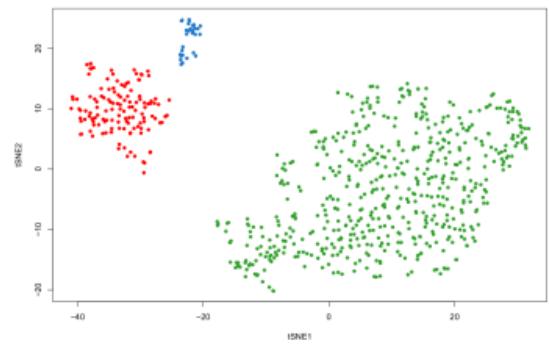
tSNE PLOT MACOSKO: PCA INITIALIZATION



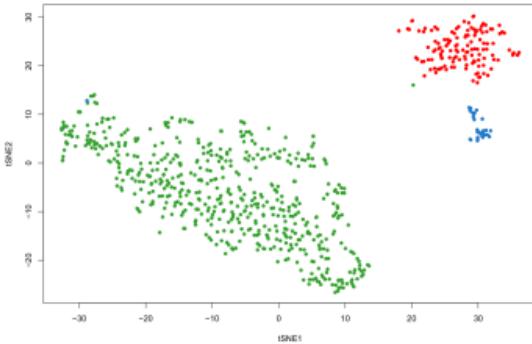
tSNE PLOT MACOSKO: LAPLACIAN EIGENMAPS INITIALIZATION



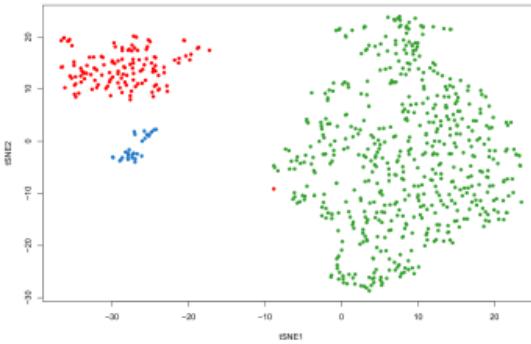
tSNE PLOT BARTOSHEK: RANDOM INITIALIZATION



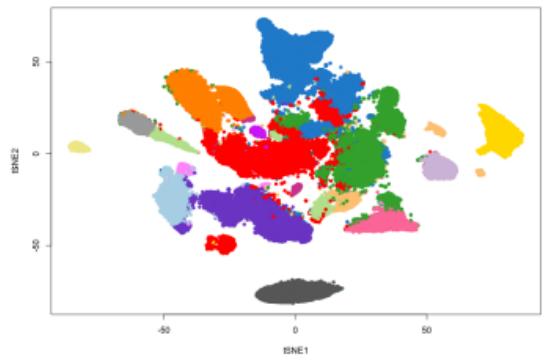
tSNE PLOT BARTOSHEK: PCA INITIALIZATION



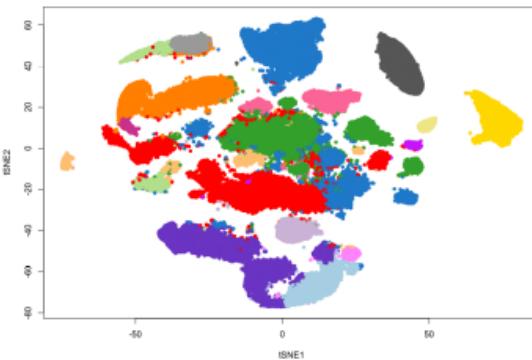
tSNE PLOT BARTOSHEK: LAPLACIAN EIGENMAPS INITIALIZATION



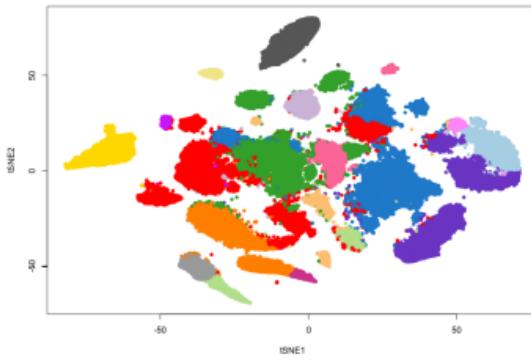
tSNE PLOT MCA 250K CELLS: RANDOM INITIALIZATION



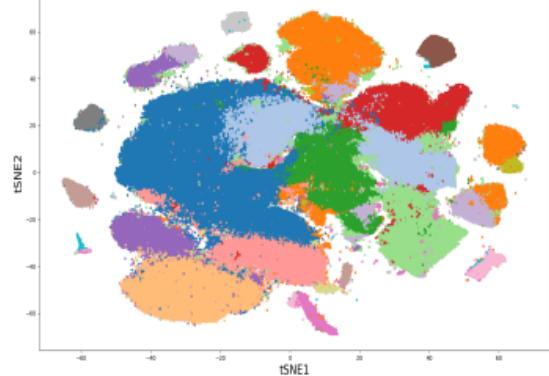
tSNE PLOT MCA 250K CELLS: PCA INITIALIZATION



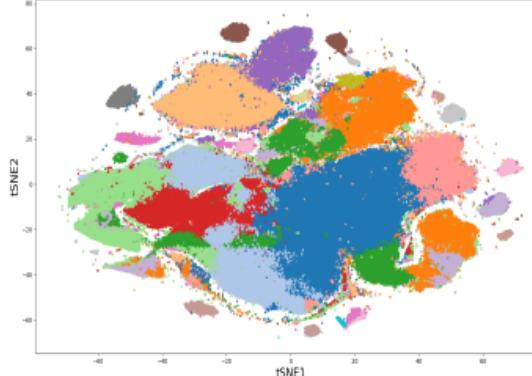
tSNE PLOT MCA 250K CELLS: LAPLACIAN EIGENMAPS INITIALIZATION



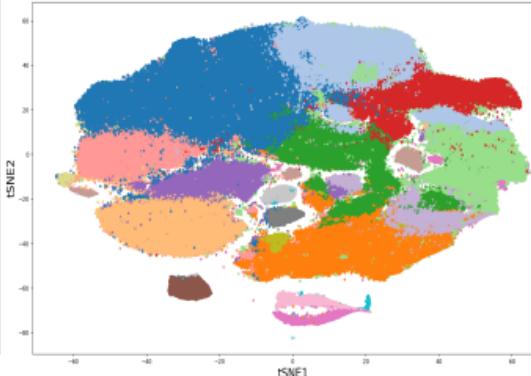
tSNE PLOT 10X 1.3M CELLS MOUSE BRAIN: RANDOM INITIALIZATION



tSNE PLOT 10X 1.3M CELLS MOUSE BRAIN: PCA INITIALIZATION

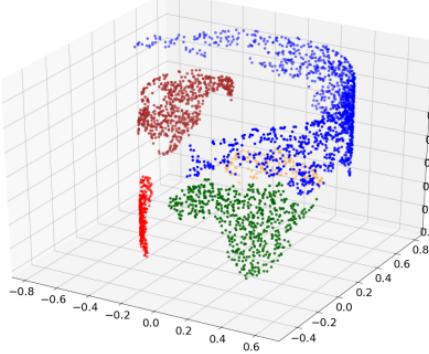


tSNE PLOT 10X 1.3M CELLS MOUSE BRAIN: LAPLACIAN EIGENMAPS INITIALIZATION

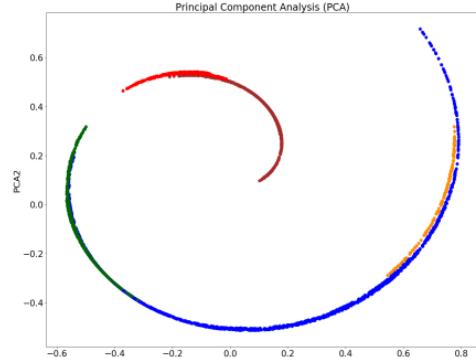


Non-Linear Manifold: PCA vs. tSNE vs. UMAP

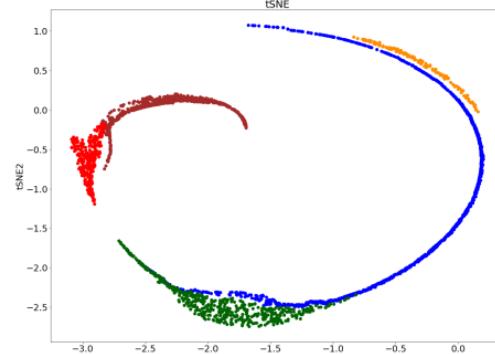
Swiss Roll: 3023 points



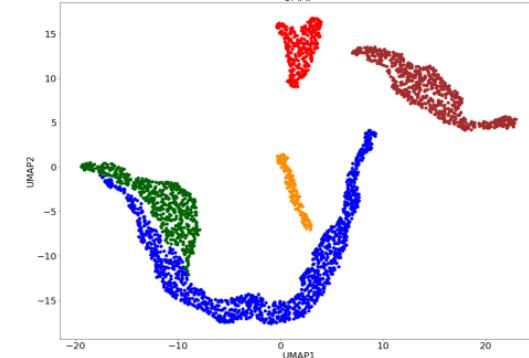
Principal Component Analysis



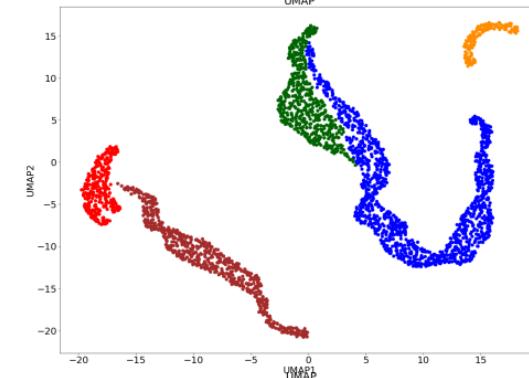
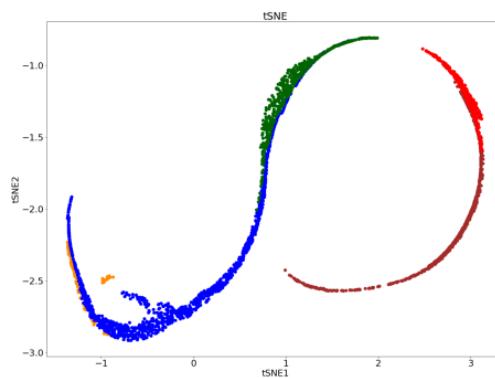
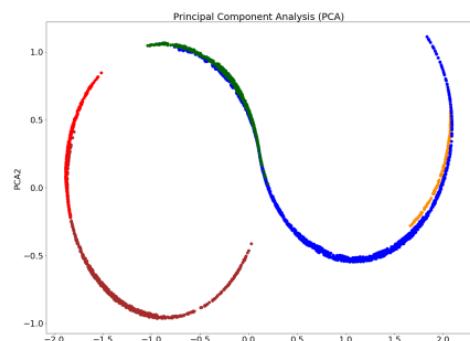
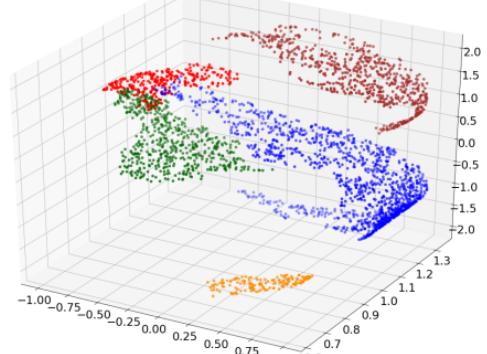
tSNE: perplexity = 2000



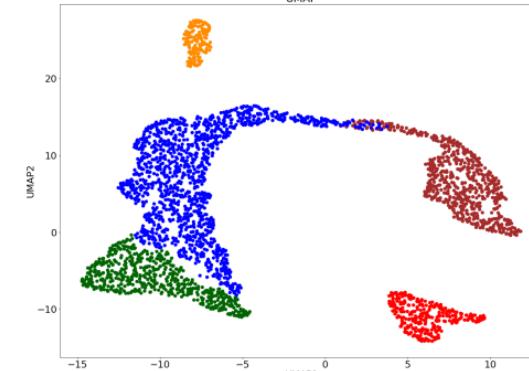
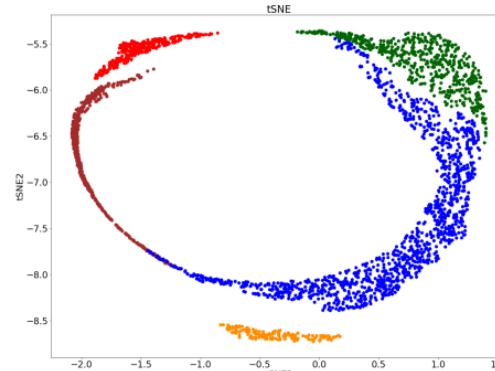
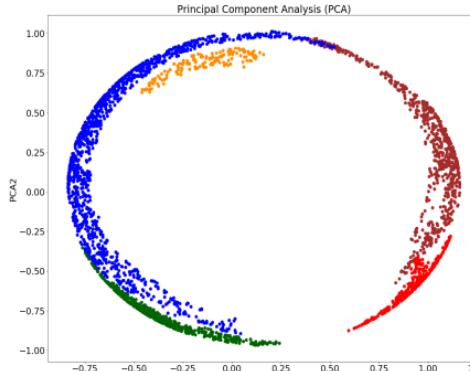
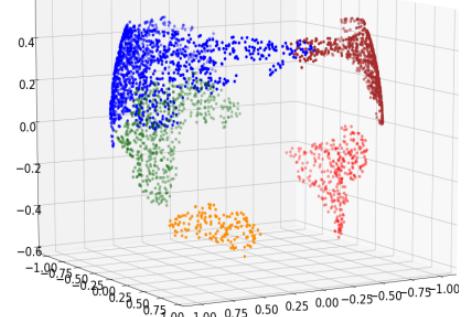
UMAP: n_neighbor = 2000

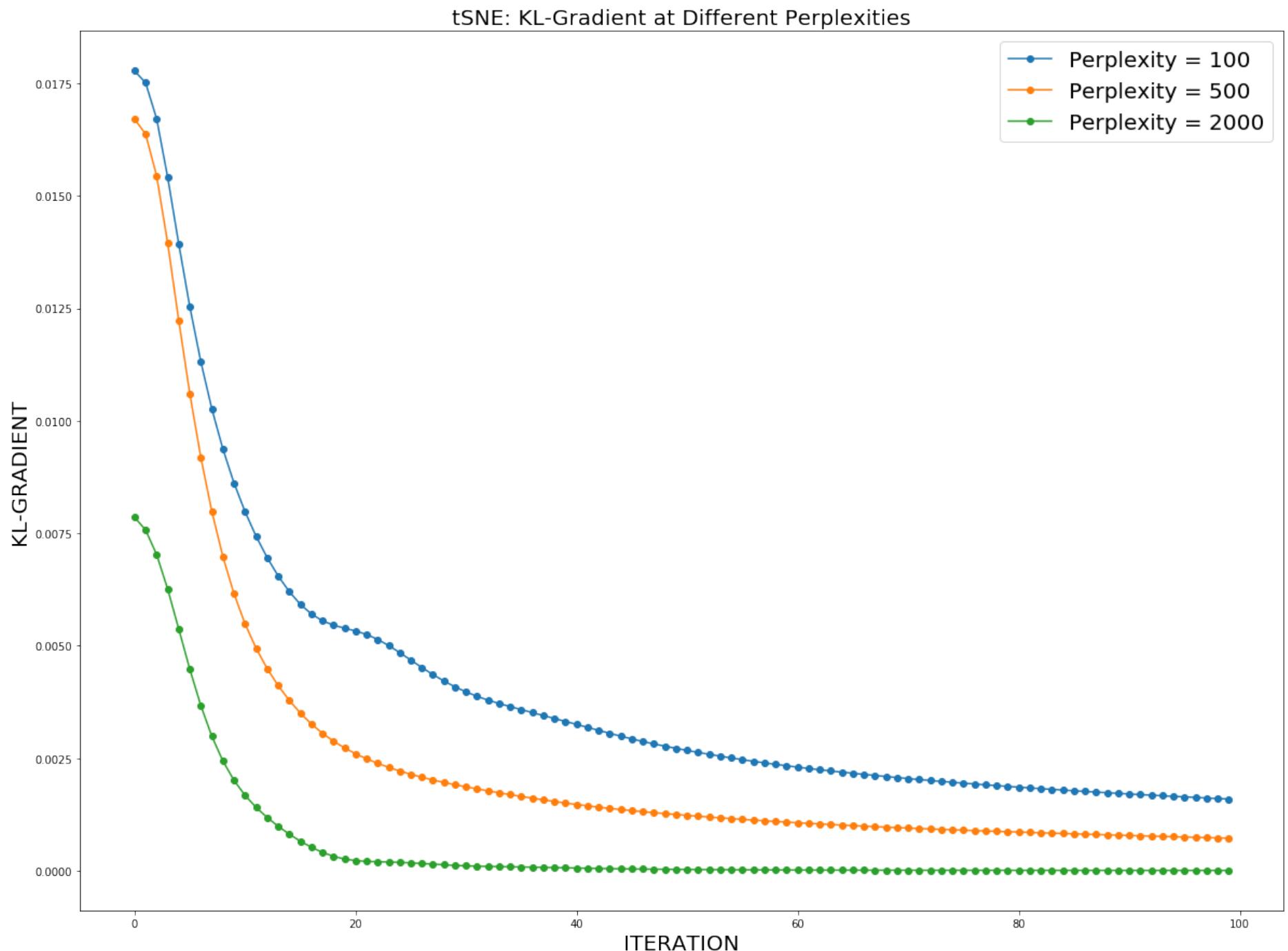


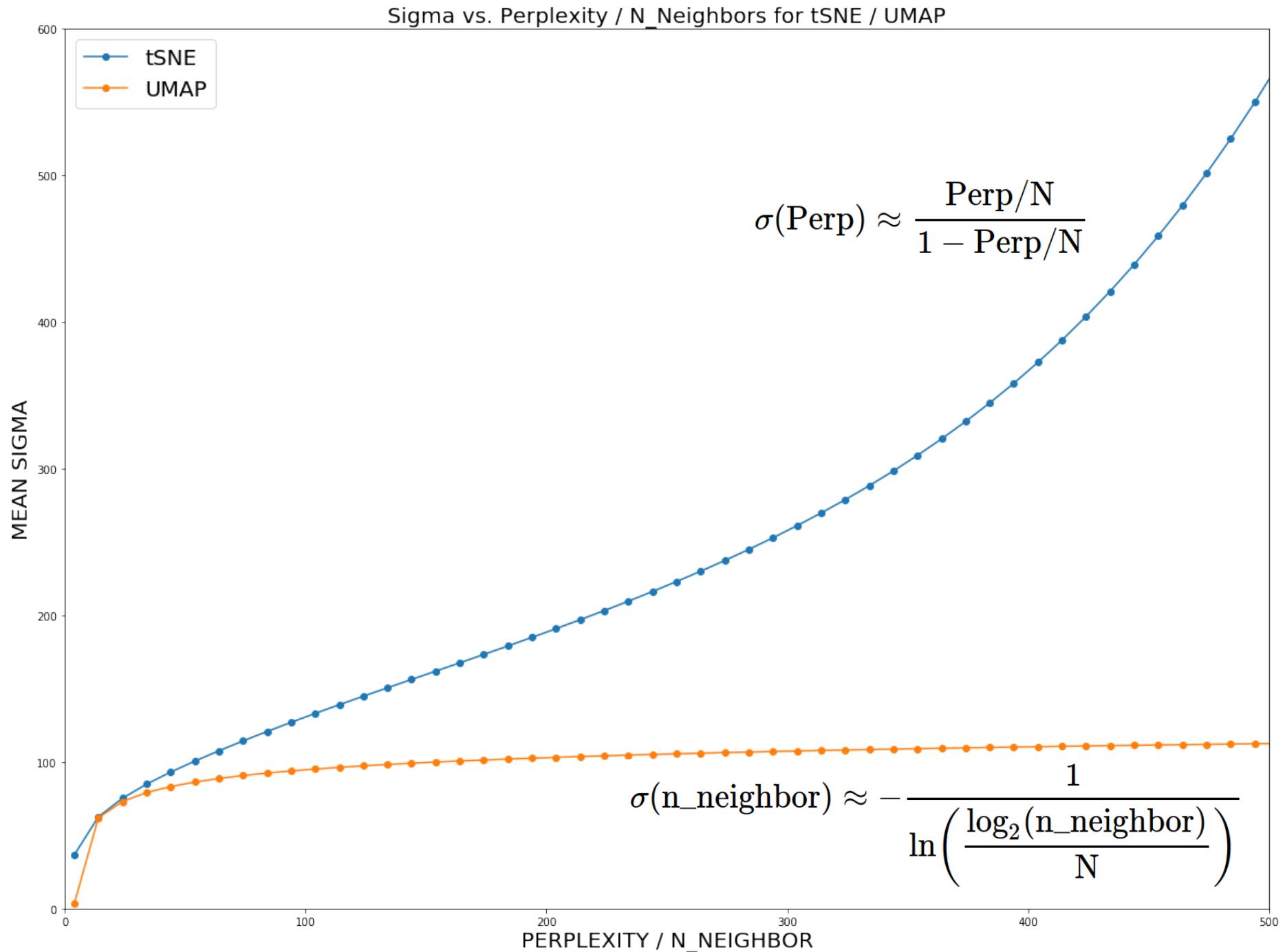
S-shape: 3023 points



Sphere: 3023 points







$$y_i = y_i - \mu \frac{\partial KL}{\partial y_i}; \quad \frac{\partial KL}{\partial y_i} = 4 \sum_j (p_{ij} - q_{ij})(y_i - y_j) \frac{1}{1 + (y_i - y_j)^2}; \quad q_{ij} \approx \frac{1}{1 + (y_i - y_j)^2}; \quad p_{ij} \approx e^{-\frac{(x_i - x_j)^2}{2\sigma_i^2}}$$

In the limit $\sigma_i \rightarrow \infty$ the probability to observe points at a distance in high-dimensional space becomes $p_{ij} \rightarrow 1$. Therefore:

$$\frac{\partial KL}{\partial y_i} \approx 4 \sum_j \left(1 - \frac{1}{1 + (y_i - y_j)^2}\right) (y_i - y_j) \frac{1}{1 + (y_i - y_j)^2} = 4 \sum_j \frac{(y_i - y_j)^3}{(1 + (y_i - y_j)^2)^2}$$

In the limit of close embedding points:

$$y_i - y_j \rightarrow 0 : \quad \frac{\partial KL}{\partial y_i} \approx 4 \sum_i (y_i - y_j)^3 \rightarrow 0$$

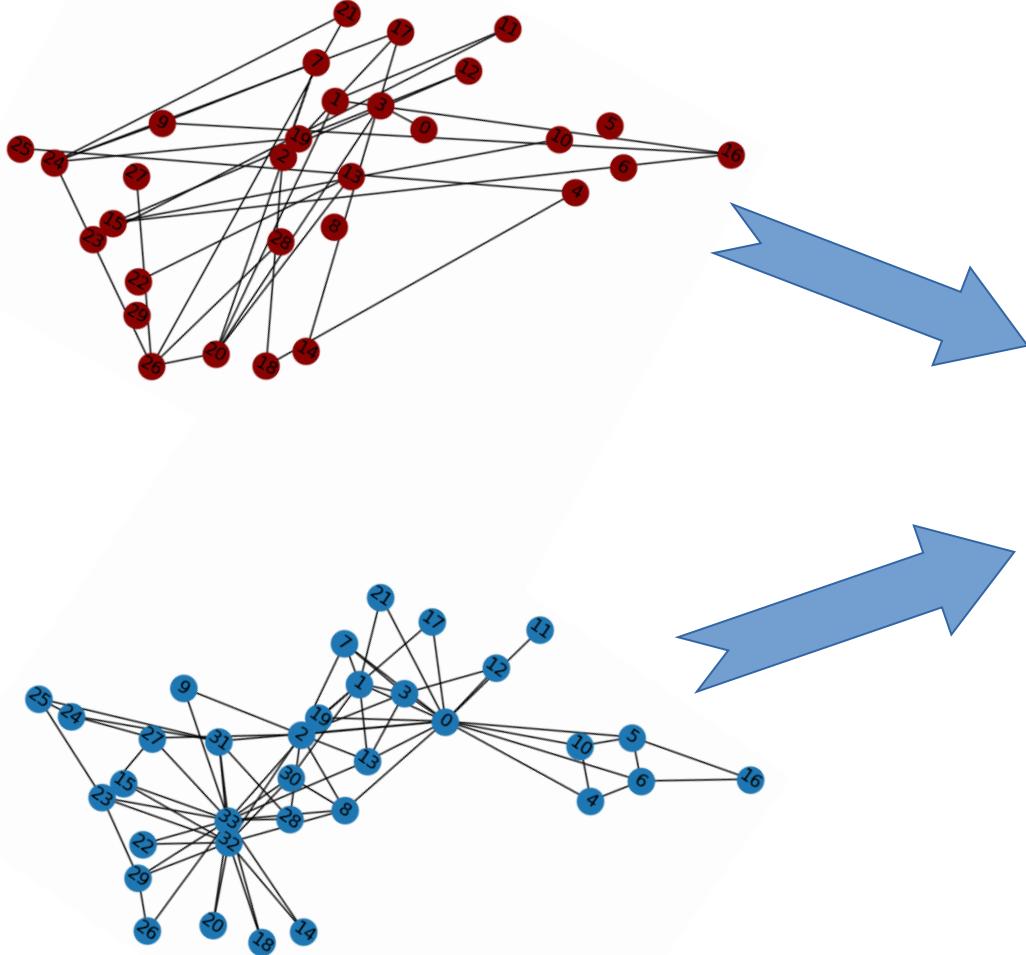
In the limit of distant embedding points:

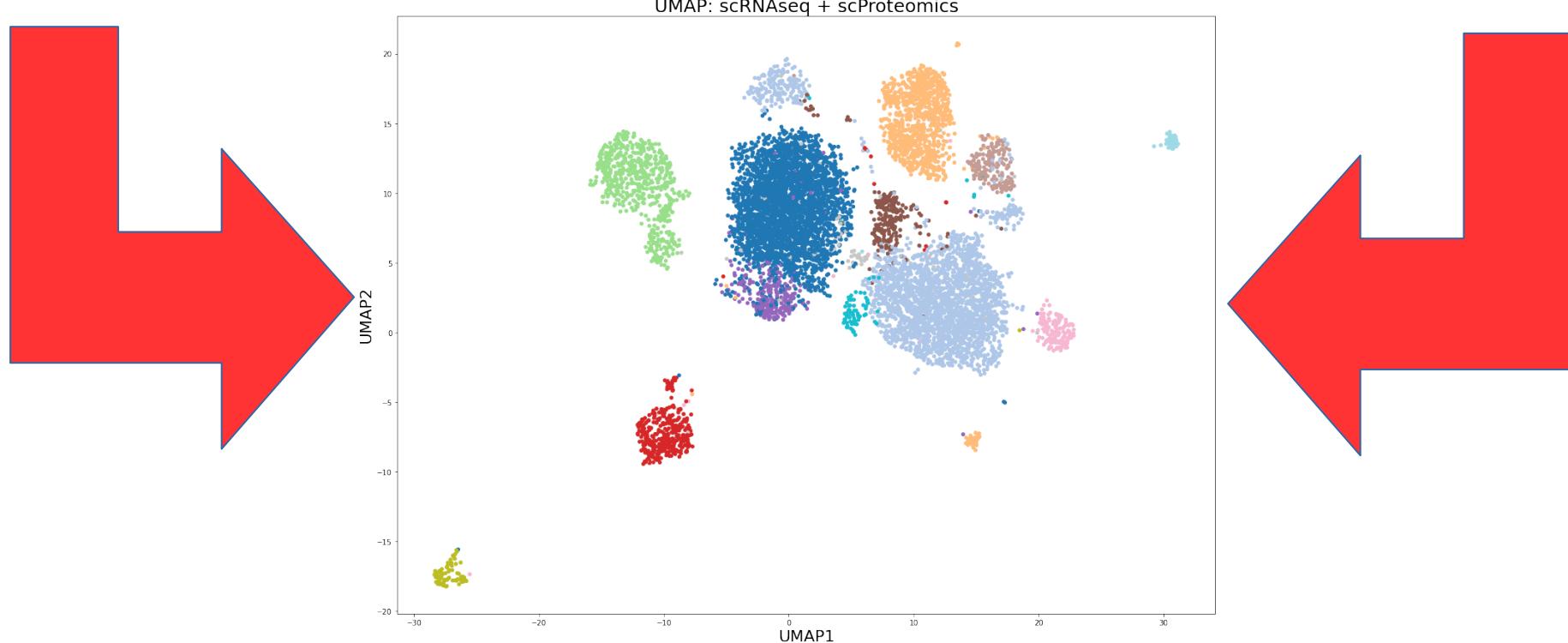
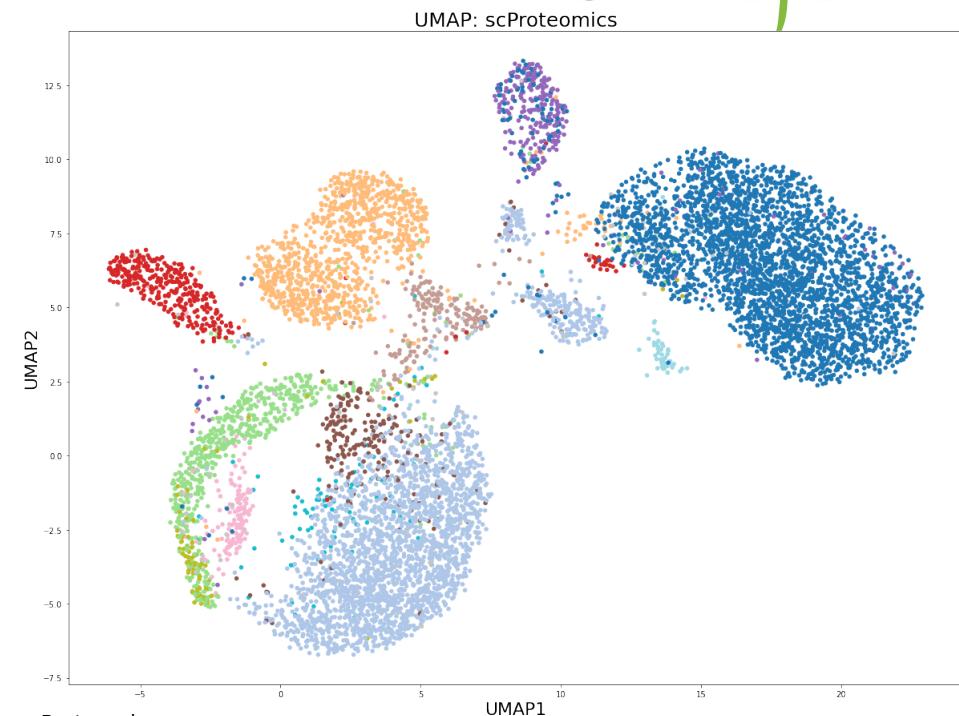
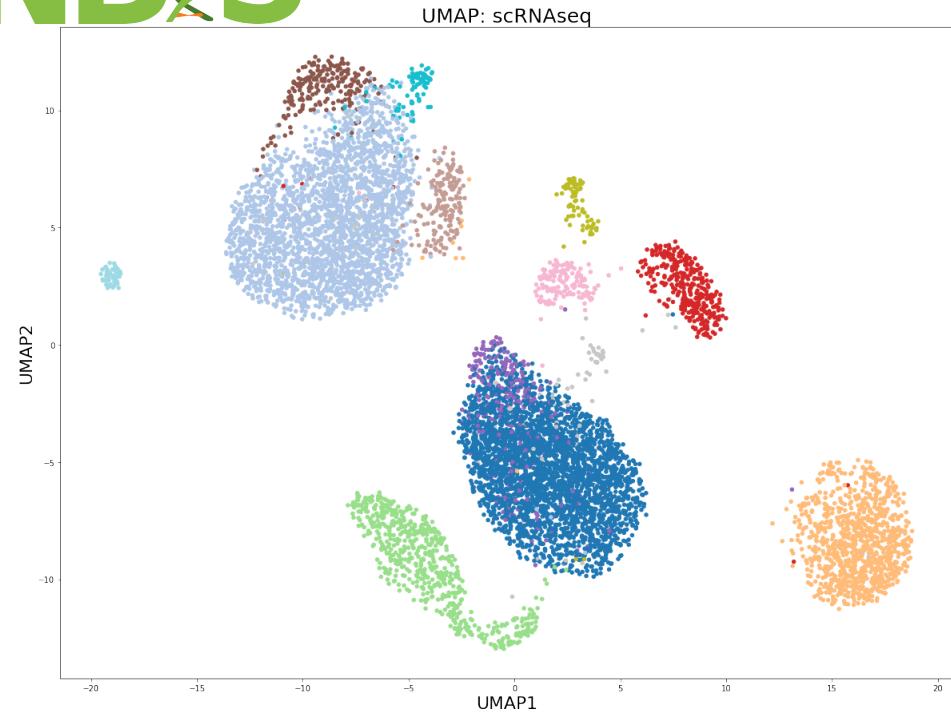
$$y_i - y_j \rightarrow \infty : \quad \frac{\partial KL}{\partial y_i} \approx 4 \sum_i \frac{1}{y_i - y_j} \rightarrow 0$$

Large Perplexity might lead to disappearance of the KL-gradient,

therefore tSNE might degrade to something linear,

so if you start with PCA initialization and increase Perplexity, you end up with PCA







National Bioinformatics Infrastructure Sweden (NBIS)

SciLifeLab



*Knut och Alice
Wallenbergs
Stiftelse*



LUNDS
UNIVERSITET