# Non-Linear Dimensionality Reduction in R

Nikolay Oskolkov, MRG Group Leader, LIOS, Riga, Latvia
R course, 09.02.2026



@NikolayOskolkov

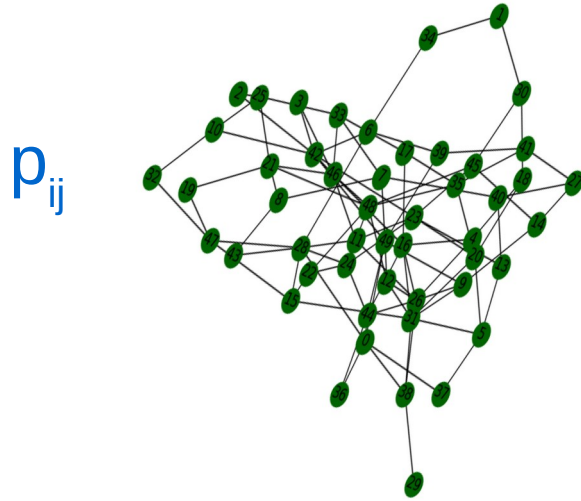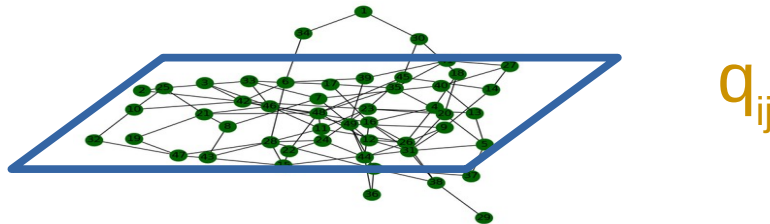@oskolkov.bsky.social

Personal homepage:
https://nikolay-oskolkov.com

Image generated by ChatGPT

Topics we'll cover in this session:

1) Neighborhood graph principle behind non-linear dimension reduction

2) Overview of tSNE and UMAP algorithms

3) Limitations of tSNE and UMAP algorithms

4) Comparing pros and cons of PCA and UMAP applications to cell biology and population genetics studies
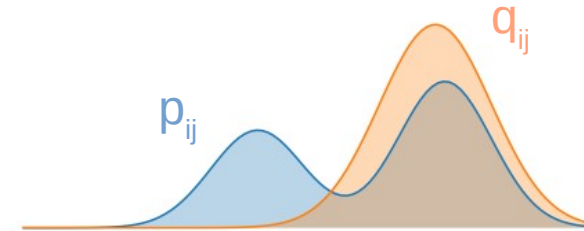
# Non-linear dimension reduction: neighborhood graph
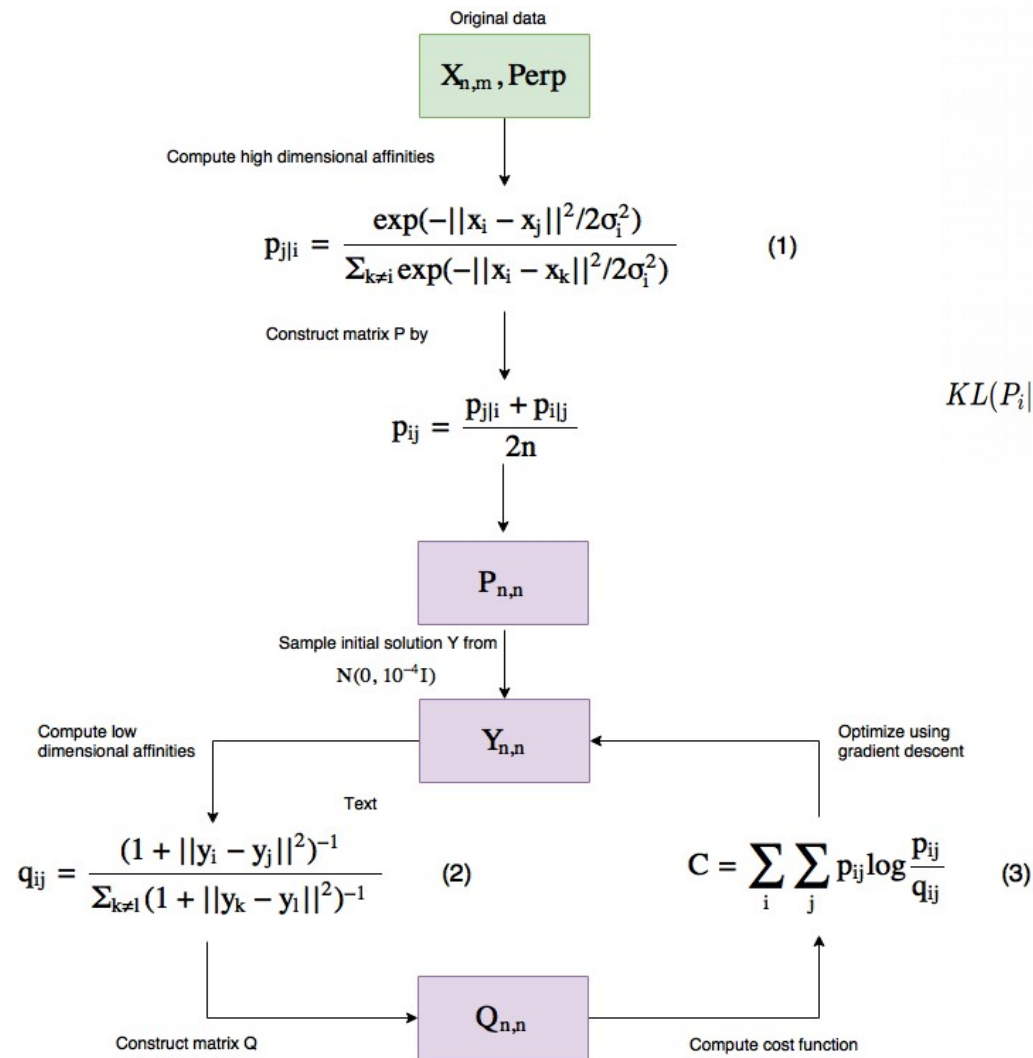
1) Construct high-dimensional graph

$p_{ij}$

3) Collapse the graphs together

$q_{ij}$

$p_{ij}$

Kullback-Leibler divergence

2) Construct low-dimensional graph

$q_{ij}$

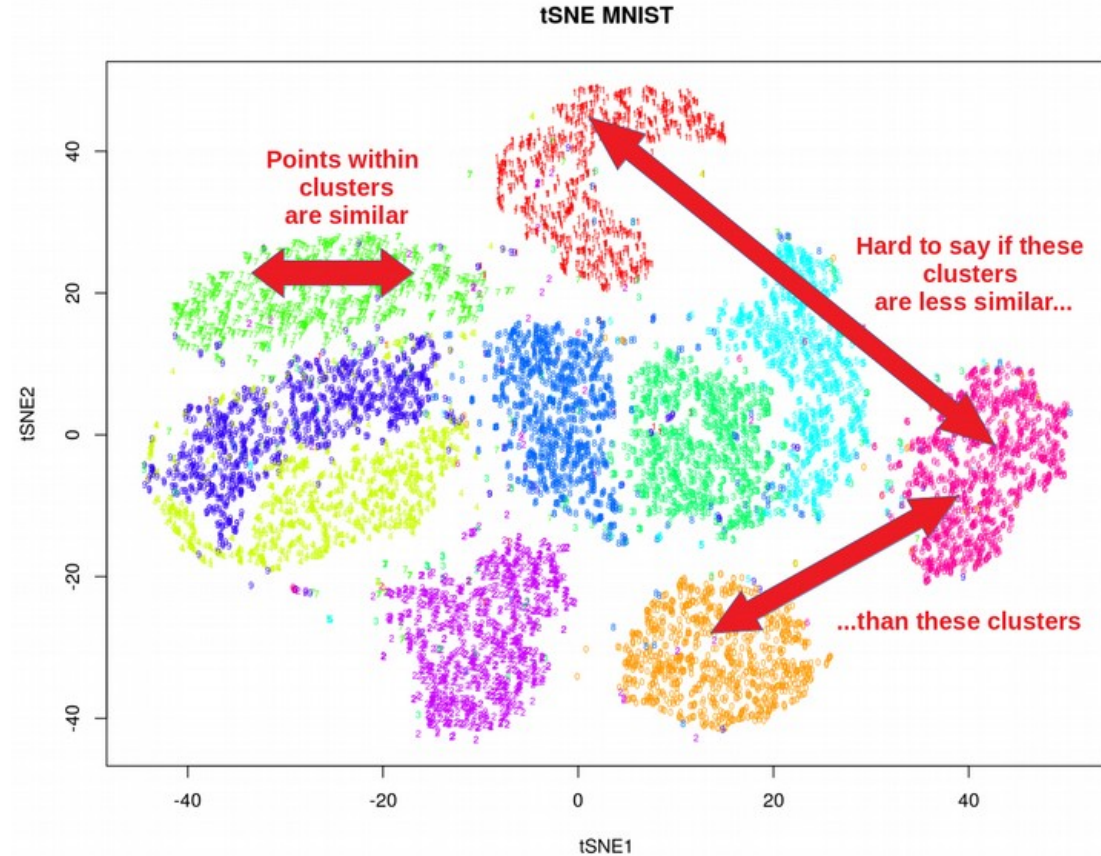tSNE does not scale for large data sets?
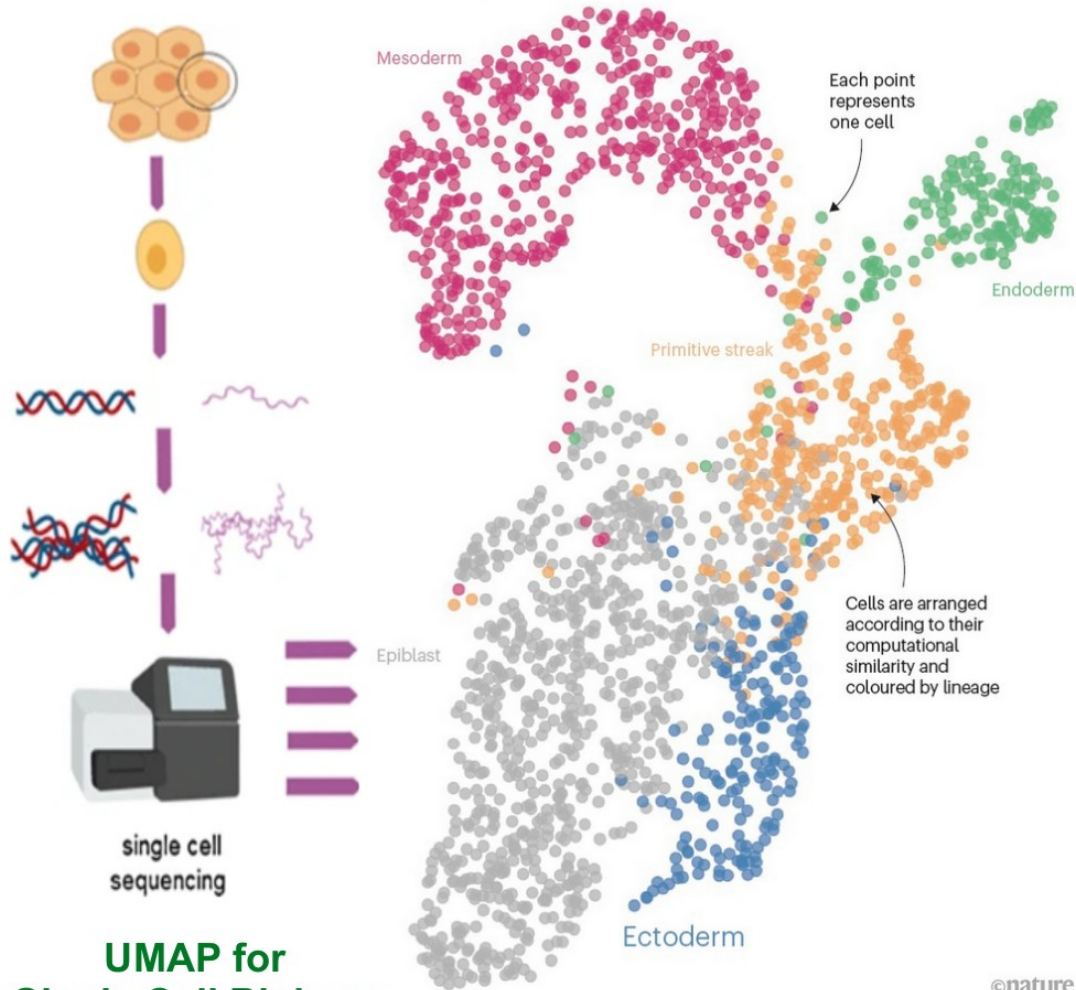
tSNE does not preserve global structure?

tSNE can only embed into 2-3 dims?

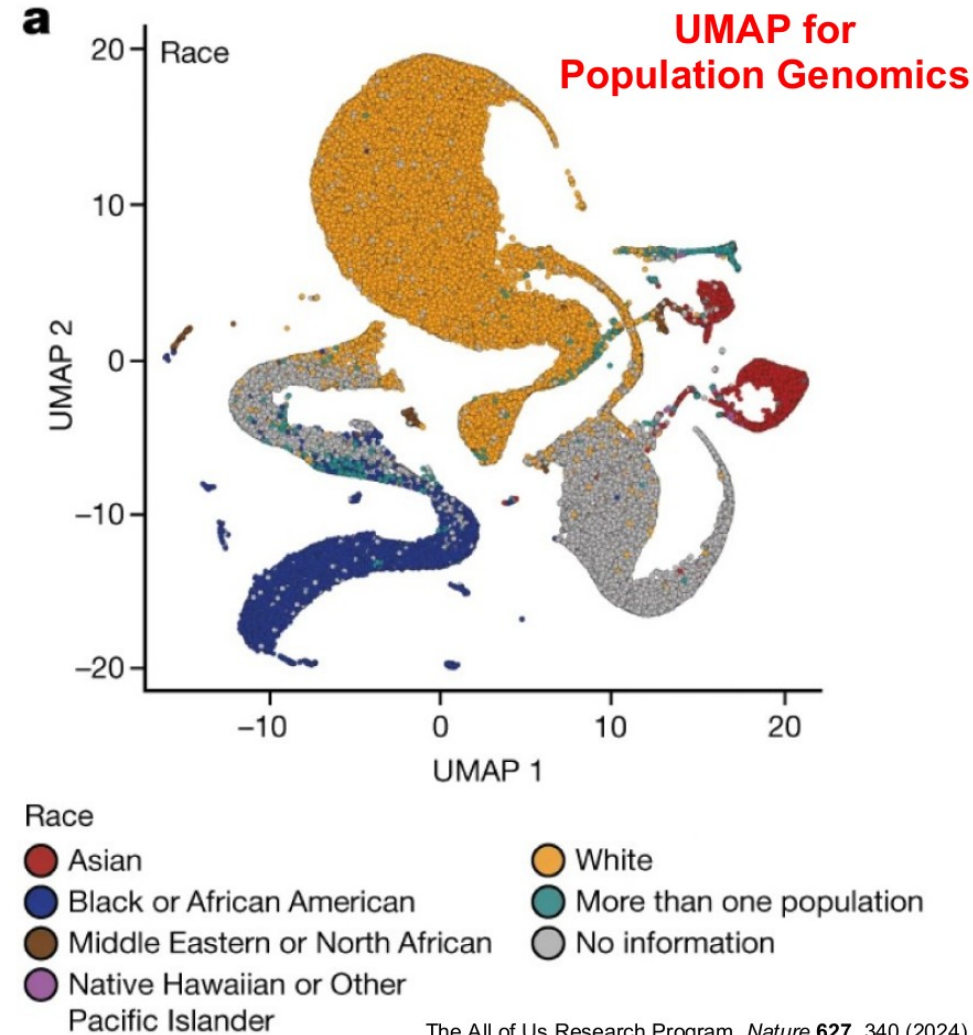tSNE performs non-parametric mapping (no variance explained statistics)?

tSNE can not work with high-dimensional data directly (PCA needed)?
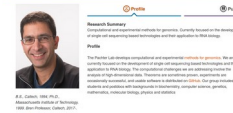
tSNE uses too much RAM at large perp?

# UMAP: Single Cell vs. PopGen



**UMAP for Single Cell Biology**

Mesoderm

Each point represents one cell

Endoderm

Primitive streak

Cells are arranged according to their computational similarity and coloured by lineage

Epiblast

Ectoderm

©nature

https://www.nature.com/articles/d41586-021-01994-w

**UMAP for Population Genomics**

a    Race

Race
- 🔴 Asian
- 🔵 Black or African American
- 🟤 Middle Eastern or North African
- 🟣 Native Hawaiian or Other Pacific Islander
- 🟠 White
- 🟢 More than one population
- ⚪ No information

The All of Us Research Program, *Nature* **627**, 340 (2024)

# UMAP (and Single Cell) Criticism

**UMAP**

**PCA**

Are PUR genetically closer to GBR + IBS than ASW + YRI?

**ABSOLUTELY NOT!**

- Because of their meaningless inter-cluster distances tSNE / UMAP are less useful for population genomics than PCA

- The goal of tSNE / UMAP is to **discover clusters**, which is sufficient for Single Cell Biology but not for PopGen.

- In PopGen we generally do not discover clusters, we have an idea about e.g. human populations, and the aim is often to explore the **genetic relatedness** between the populations, a task UMAP can absolutely not solve!

Take home messages of the session:

1) Neighborhood graph is the key of non-linear dimension reduction

2) Kullback-Leibler divergence is the objective function of tSNE

3) Inter-cluster distances in tSNE and UMAP are meaningless

4) tSNE and UMAP are appropriate for cell biology but not for human genetics applications where PCA is more accurate and informative

# Acknowledgments: LIOS + TARGETWISE