# A visual interrogation of dimension reduction tools for single-cell analysis

Tyler J Burns, PhD

AG Mei, DRFZ Berlin
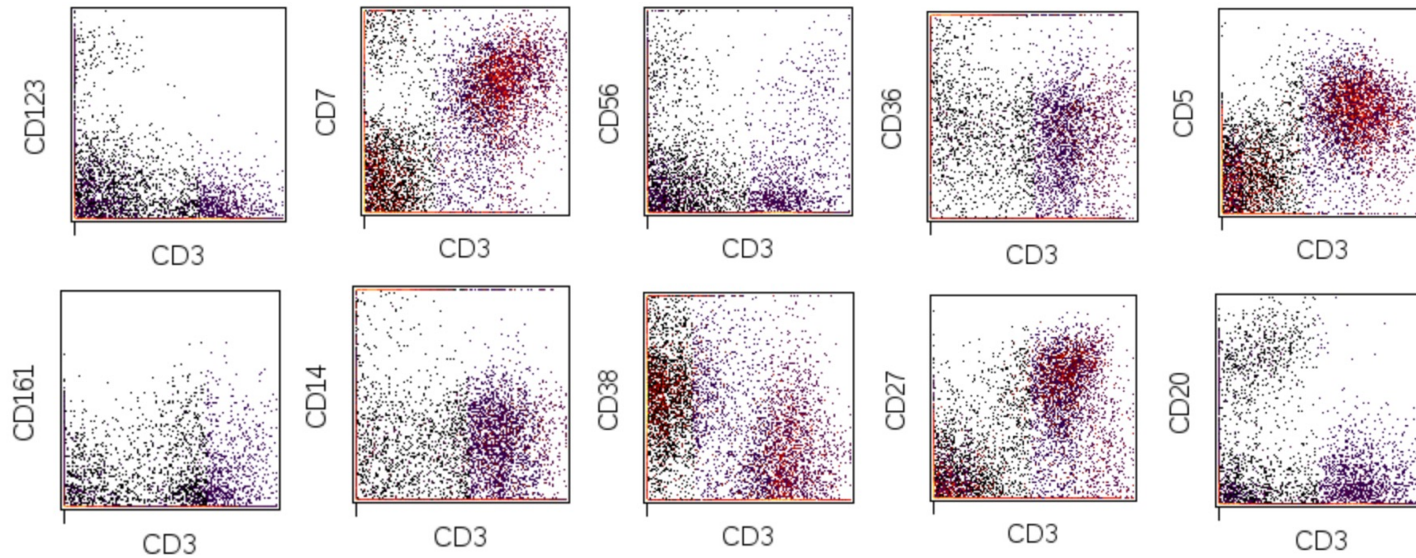
# Outline

- Part 1: Introduction
- Part 2: Preservation of local structure
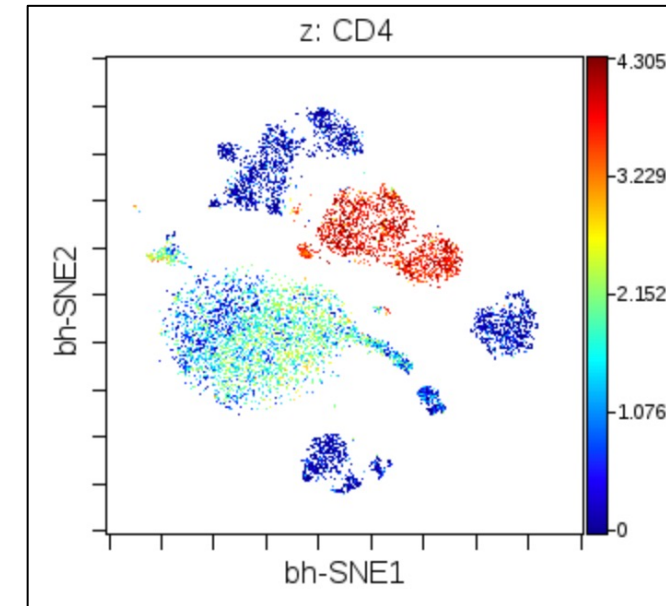- Part 3: Preservation of global structure

# Outline

- **Part 1: Introduction**
- Part 2: Preservation of local structure
- Part 3: Preservation of global structure

# Dimension reduction makes large amounts of information human-readable without too much human work



t-SNE(cells)

Amir et al, *Nat Biotechnology* 2013

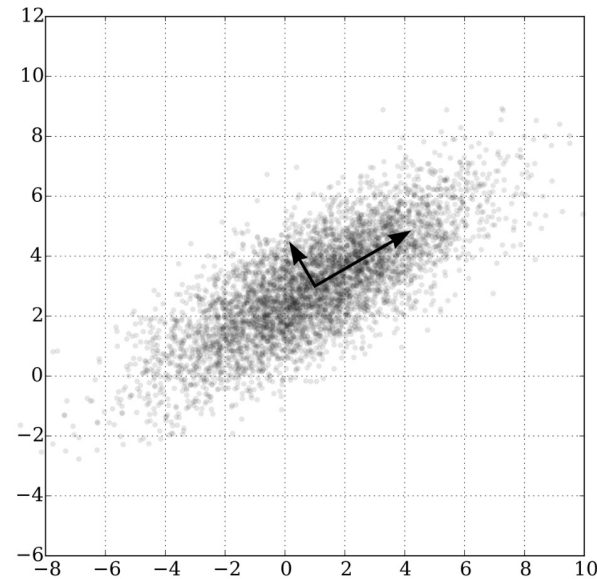# Early dimension reduction tools: Principal Component Analysis (PCA)

[ 559 ]

LIII. *On Lines and Planes of Closest Fit to Systems of Points in Space. By* KARL PEARSON, *F.R.S., University College, London* *. (1901)

(1) IN many physical, statistical, and biological investigations it is desirable to represent a system of points in plane, three, or higher dimensioned space by the "best-fitting" straight line or plane. Analytically this consists in taking

$$y = a_0 + a_1 x, \quad \text{or} \quad z = a_0 + a_1 x + b_1 y,$$
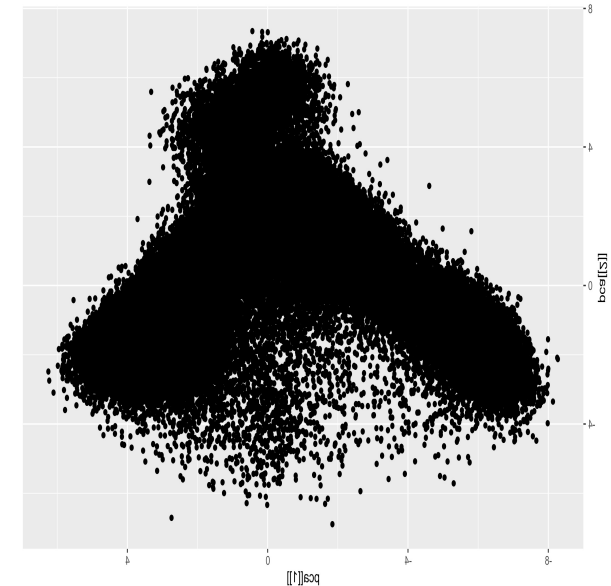$$\text{or} \quad z = a_0 + a_1 x_1 + a_2 x_2 + a_3 x_3 + \ldots + a_n x_n,$$

where $y$, $x$, $z$, $x_1$, $x_2$, … $x_n$ are variables, and determining the "best" values for the constants $a_0$, $a_1$, $b_1$, $a_0$, $a_1$, $a_2$, $a_3$, … $a_n$ in relation to the observed corresponding values of the variables. In nearly all the cases dealt with in the text-books of least squares, the variables on the right of our equations are treated as the independent, those on the left as the dependent variables. The result of this treatment is that we get one straight line or plane if we treat some one variable as independent, and a quite different one if we treat another variable as the independent variable. There is no paradox about this; it is, in fact, an easily understood and most im-portant feature of the theory of a system of correlated.

Axes span the direction with highest variance



https://en.wikipedia.org/wiki/Principal_component_analysis

Samusik_01 bone marrow CyTOF dataset

# t-SNE preserves local information, produces more well clustered maps

**Laurens van der Maaten**
*TiCC*
*Tilburg University*
*P.O. Box 90153, 5000 LE Tilburg, The Netherlands*

LVDMAATEN@GMAIL.COM

**Geoffrey Hinton**
*Department of Computer Science*
*University of Toronto*
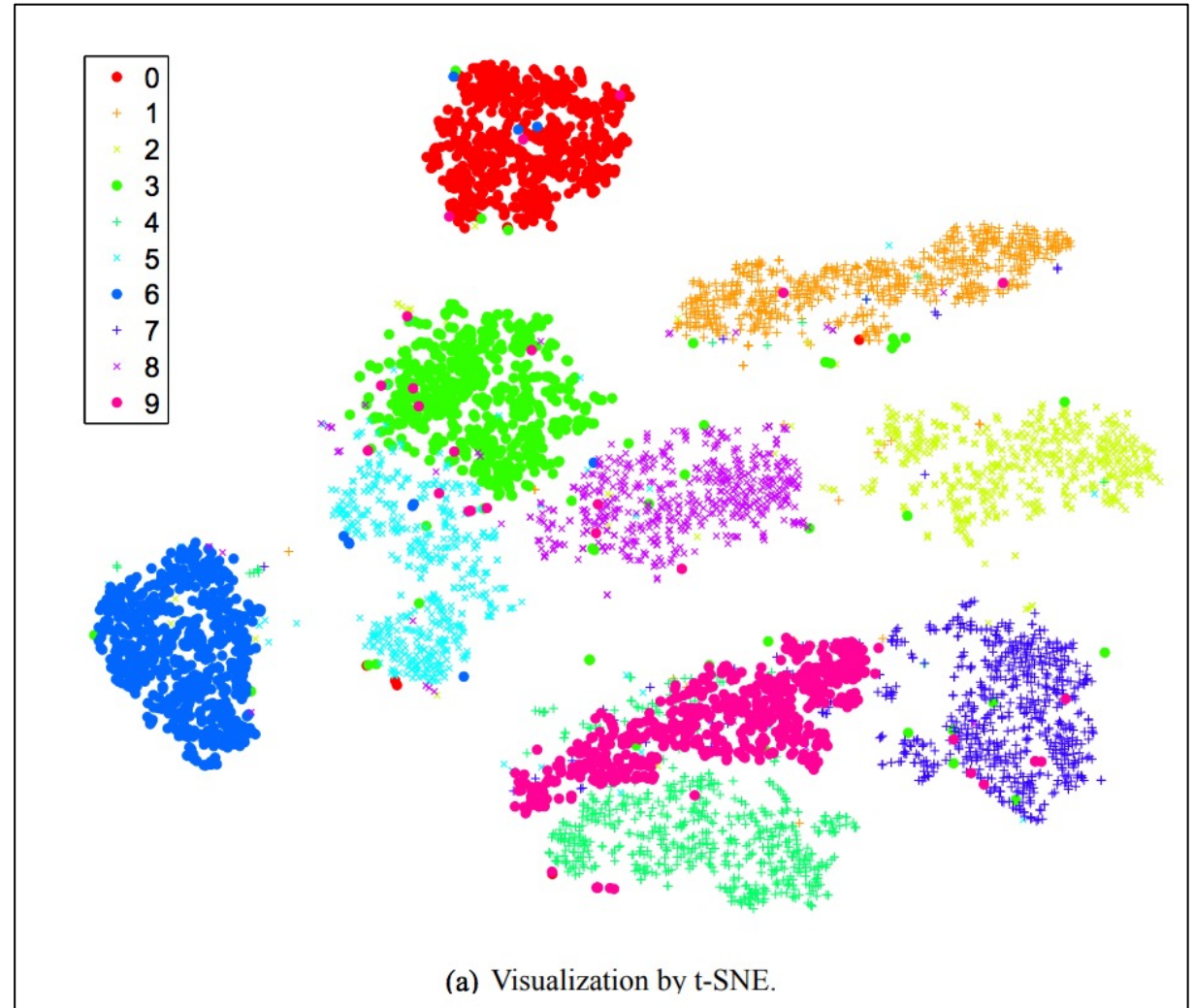*6 King's College Road, M5S 3G4 Toronto, ON, Canada*

HINTON@CS.TORONTO.EDU

**Editor:** Yoshua Bengio

**Visualizing Data using t-SNE**

## Abstract

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, and produces significantly better visualizations by reducing the tendency to crowd points together in the center of the map. t-SNE is better than existing techniques at creating a single map that reveals structure at many different scales. This is particularly important for high-dimensional data that lie on several different, but related, low-dimensional manifolds, such as images of objects from multiple classes seen from multiple viewpoints. For visualizing the structure of very large data sets, we show how t-SNE can use random walks on neighborhood graphs to allow the implicit structure of all of the data to influence the way in which a subset of the data is displayed. We illustrate the performance of t-SNE on a wide variety of data sets and compare it with many other non-parametric visualization techniques, including Sammon mapping, Isomap, and Locally Linear Embedding. The visualizations produced by t-SNE are significantly better than those produced by the other techniques on almost all of the data sets.

**Keywords:** visualization, dimensionality reduction, manifold learning, embedding algorithms, multidimensional scaling

(a) Visualization by t-SNE.

# viSNE: the adaptation of t-SNE to CyTOF

## viSNE enables visualization of high dimensional single-cell data and reveals phenotypic heterogeneity of leukemia

El-ad David Amir,[1] Kara L Davis,[2,3] Michelle D Tadmor,[1,3] Erin F Simonds,[2,3] Jacob H Levine,[1,3] Sean C Bendall,[2,3] Daniel K Shenfeld,[1,3] Smita Krishnaswamy,[1] Garry P Nolan,[2,4] and Dana Pe'er[1,4,*]

Author information ► Copyright and License information ► Disclaimer

### Abstract

Go to: ☑

High-dimensional single-cell technologies are revolutionizing the way we understand biological systems. Technologies such as mass cytometry measure dozens of parameters simultaneously in individual cells, making interpretation daunting. We developed viSNE, a tool to map high-dimensional cytometry data onto 2D while conserving high-dimensional structure. We integrated mass cytometry with viSNE to map healthy and cancerous bone marrow samples. Healthy bone marrow maps into a canonical shape that separates between immune subtypes. In leukemia, however, the shape is malformed: the maps of cancer samples are distinct from the healthy map and from each other. viSNE highlights structure in the heterogeneity of surface phenotype expression in cancer, traverses the progression from diagnosis to relapse, and identifies a rare leukemia population in minimal residual disease settings. As several new technologies raise the number of simultaneously measured parameters in each cell to the hundreds, viSNE will become a mainstay in analyzing and interpreting such experiments.

# Emergence of UMAP as an alternative to t-SNE for single-cell analysis



## UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction

Leland McInnes
Tutte Institute for Mathematics and Computing
leland.mcinnes@gmail.com

John Healy
Tutte Institute for Mathematics and Computing
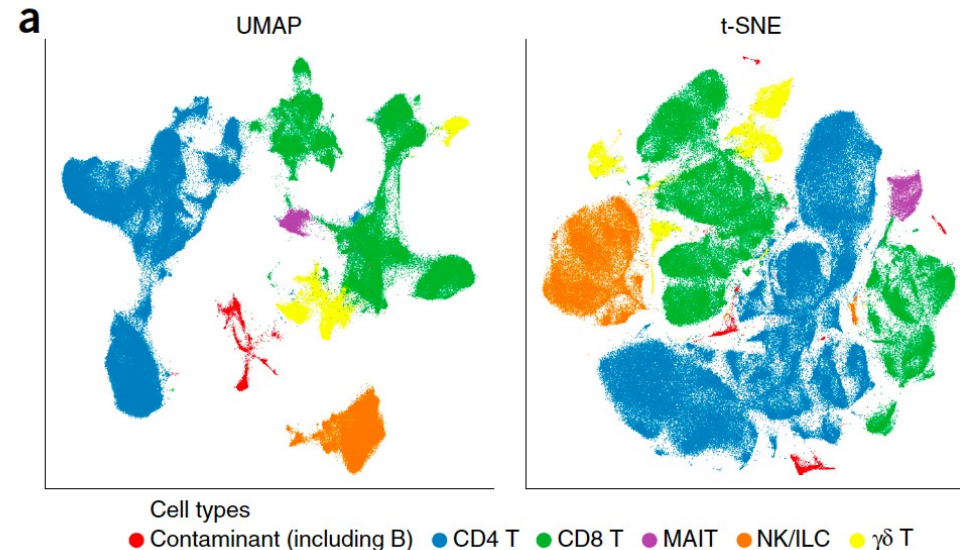jchealy@gmail.com

James Melville
jlmelville@gmail.com

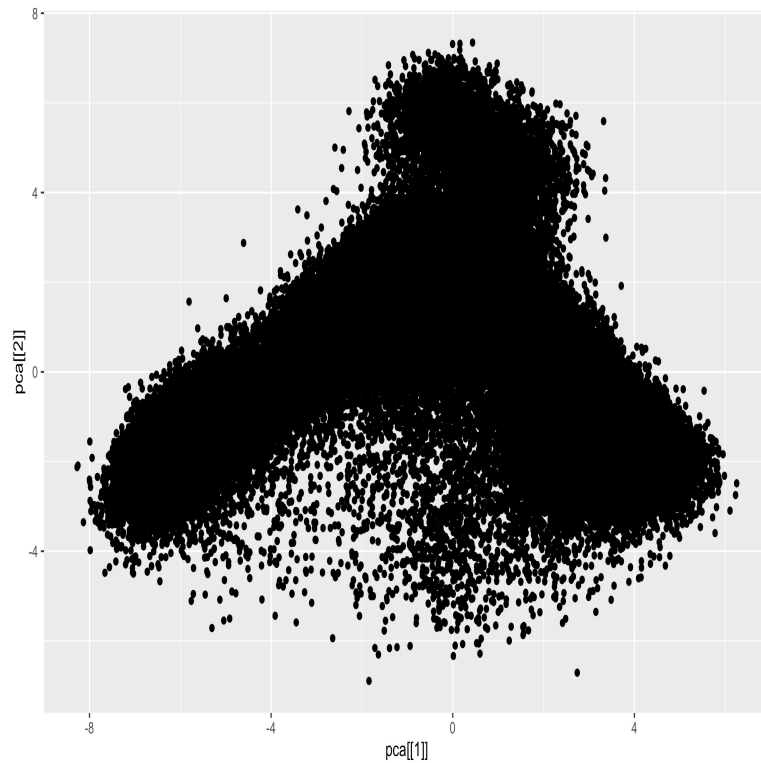December 7, 2018

ANALYSIS

nature biotechnology

Dimensionality reduction for visualizing single-cell data using UMAP

Etienne Becht[1], Leland McInnes[2], John Healy[2], Charles-Antoine Dutertre[1], Immanuel W H Kwok[1], Lai Guan Ng[1], Florent Ginhoux[1] & Evan W Newell[1,3]

# What PCA, t-SNE and UMAP look like on a bone marrow CyTOF dataset



PCA

t-SNE

UMAP

Dataset: Samusik *et al, Nature Methods* 2016

# t-SNE and UMAP are accessible from single-cell analysis user interfaces

## Cytobank

## FlowJo



Others: Astrolabe, OMIQ, Tercen

What additional information about dimension reduction maps should we know for their proper use?

# Credit for the following t-SNE and UMAP explanations

## UMAP Uniform Manifold Approximation and Projection for Dimension Reduction | SciPy 2018

29,591 views • Jul 13, 2018

👍 762   👎 4   ➤ SHARE

**Enthought**
41K subscribers

This talk will present a new approach to dimension reduction called UMAP. UMAP is grounded in manifold learning and topology, making an effort to preserve the topological structure of the data. The resulting algorithm can provide both 2D visualisations of data of comparable quality to t-SNE,

SHOW MORE

## StatQuest: t-SNE, Clearly Explained

17,938 views

👍 489   👎 10

**StatQuest with Josh Starmer**
Published on Sep 18, 2017

t-SNE is a popular method for making an easy to read graph from a complex dataset, but not many people know how it works. Here's the dope! Also, if you'd like to see a code example in R, here's one:

SHOW MORE

# The goal of t-SNE and UMAP is to reduce dimensions while preserving specific information about each cell's neighbors

Higher dimensional space

Low dimensional embedding

# t-SNE and UMAP start with a low-dimensional embedding of randomly placed points

# t-SNE weights its neighbors based on distance fitted to a distribution



First, measure the distance between two points...

Then plot that distance on a normal curve that is centered on the point of interest...

...lastly, draw a line from the point to the curve. The length of that line is the "unscaled similarity".

# UMAP weights its neighbors based on topology

Find the diameter of a neighborhood. Think of it like a ball.

Find the probability that a 1-simplex exists between two points in a neighborhood

Resulting neighbor graph is a bunch of simplexes glued together (simplicial complex). Simpler structure but preserves topological information.

0-simplex    1-simplex    2-simplex    3-simplex

# The weighted neighborhood graphs can be represented as similarity matrices



Image source: YouTube: StatQuest with Josh Starmer: *t-SNE, clearly explained*

# Make these similarity matrices as similar to each other as possible, and then you're done



t-SNE uses a t-distribution here

...without it the clusters would all clump up in the middle and be harder to see.

# Make these similarity matrices as similar to each other as possible, and then you're done



UMAP also makes a 2-D simplicial complex



UMAP only moves one or a few cells at a time, rather than all of them as t-SNE does.

A cell will move toward one cell and away from another

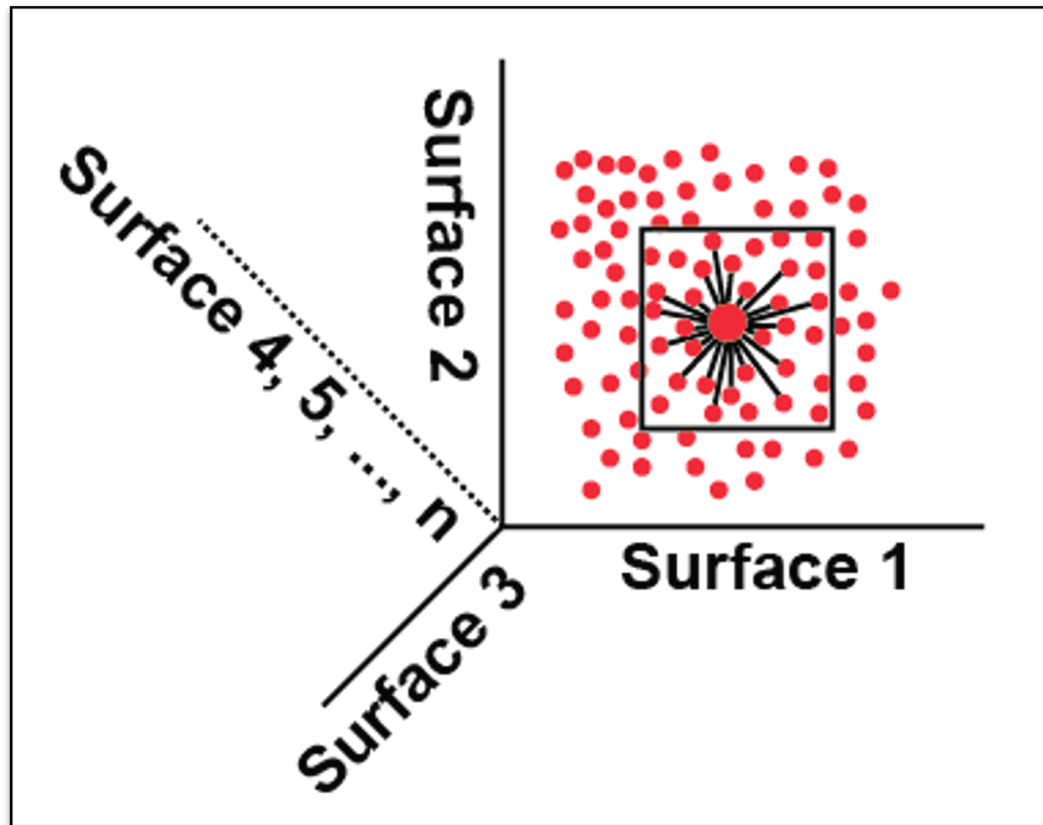# Dimension reduction maps group similar cells near each other



Samusik Bone marrow

# Outline

- Part 1: Introduction
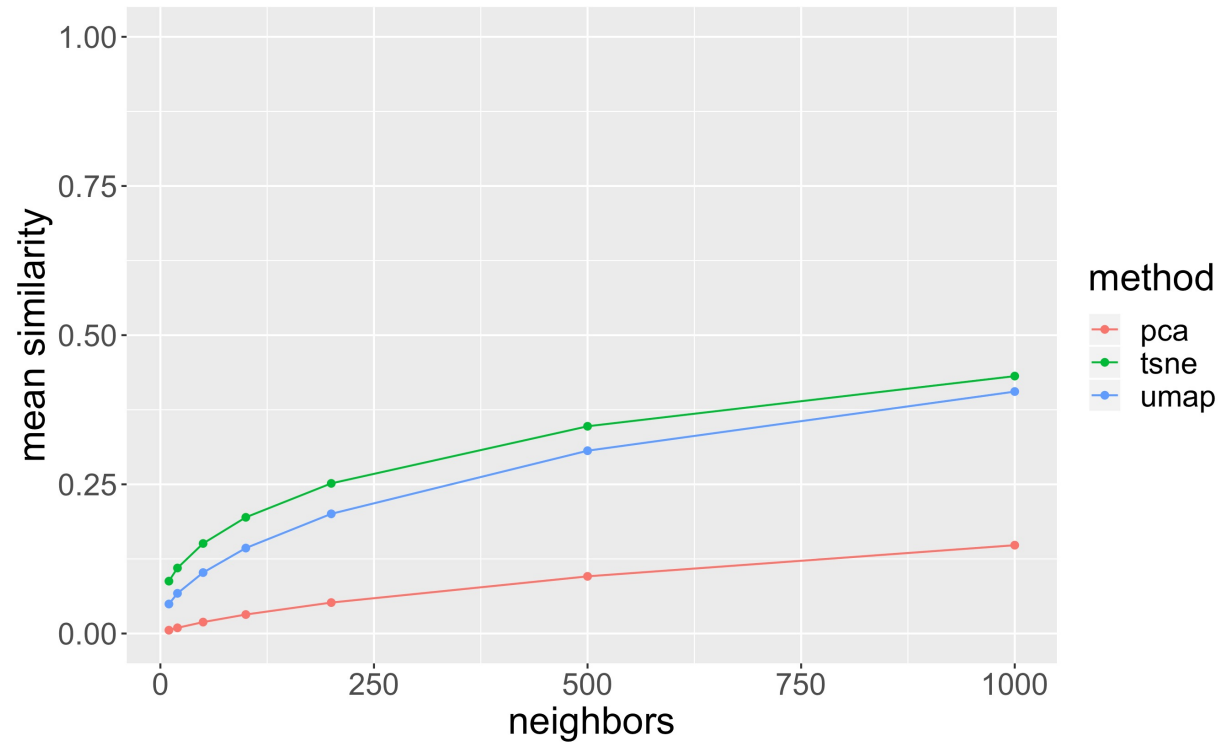- **Part 2: Preservation of local structure**
- Part 3: Preservation of global structure

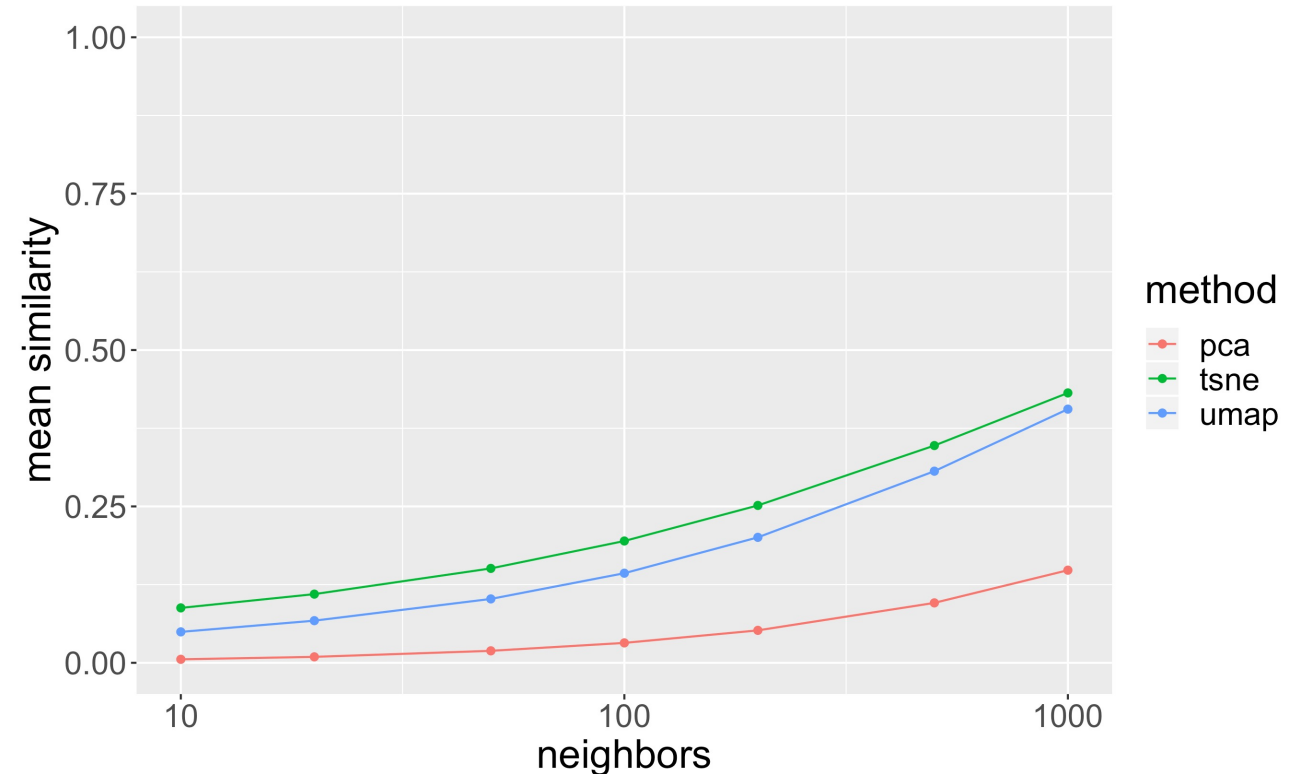# KNN to determine preservation of lower dimensional embeddings

KNN orig    KNN low-D

Find KNN for each cell from high-dim space → Find KNN for each cell from a 2-D embedding (eg t-SNE) → Compare KNN identities from the 2-D embedding and high-dim space

Repeat across a wide range of values for K

Surface 4, 5, ..., n

Surface 2

Surface 1

Surface 3

Bioconductor package: Sconify

# Global KNN comparison between t-SNE, UMAP, and PCA

Dataset: Samusik Bone marrow (public)
Num. cells: 100k

X axis is on a log scale

# t-SNE outperforms UMAP in KNN preservation, has been observed in scRNA seq data

## The art of using t-SNE for single-cell transcriptomics (2019)

Dmitry Kobak [1*] & Philipp Berens [1,2,3,4*]
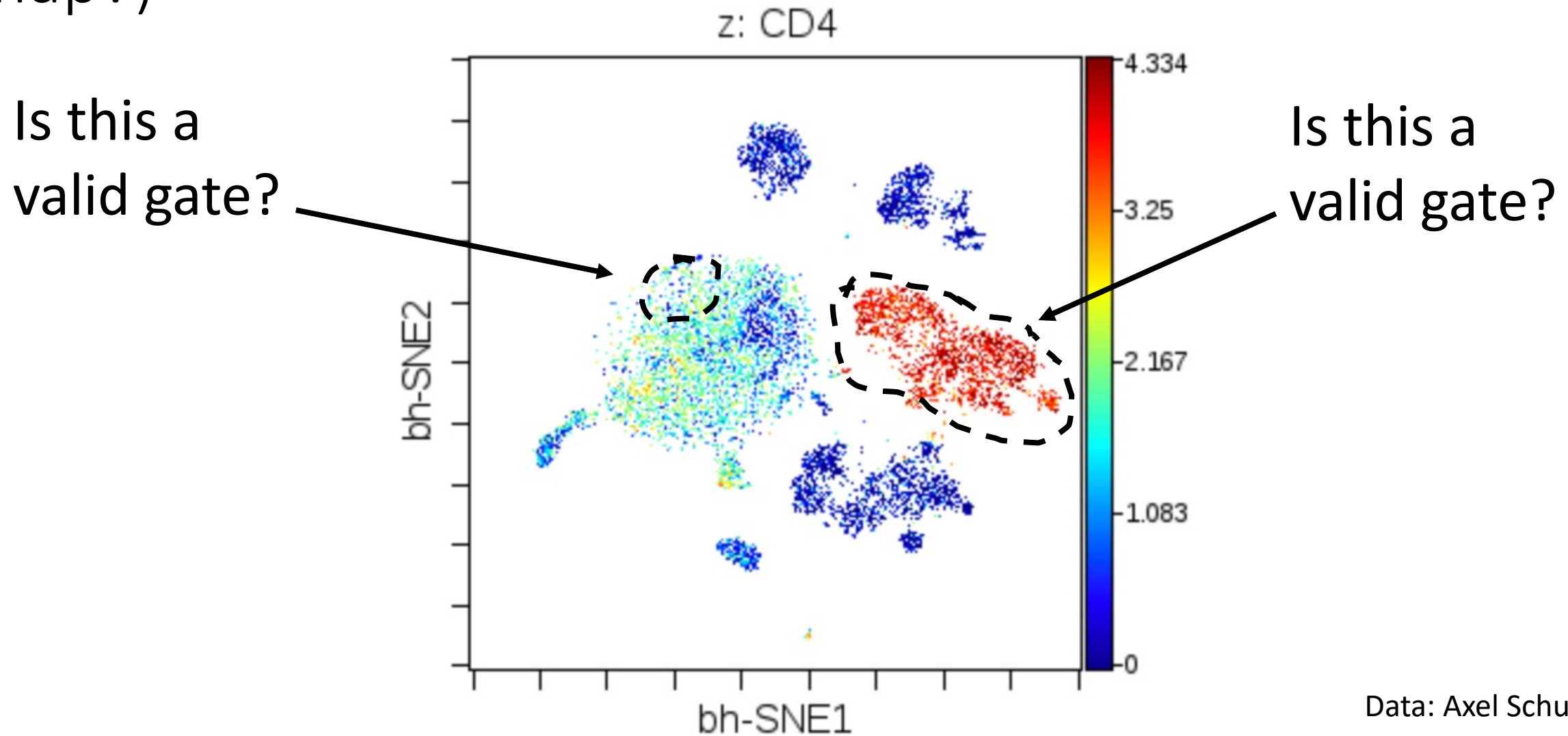
(Also did KNN preservation, K = 10 only)

To compare UMAP with our t-SNE approach in terms of preservation of global structure, we first ran UMAP on the synthetic and on the Tasic et al.[3] data sets (Supplementary Fig. 2). We used the default UMAP parameters, and also modified the two key parameters (number of neighbours and tightness of the embedding) to produce a more t-SNE-like embedding. In both cases and for both data sets, all three metrics (KNN, KNC, and CPD) were considerably lower than with our t-SNE approach.

X axis is on a log scale



method
- pca
- tsne
- umap

# Results confirmed across 3 datasets, but with very large standard deviation

# Does dimension reduction maps preserve some regions better than others (should and/or how should we gate the map?)



Is this a valid gate?

Is this a valid gate?

z: CD4

bh-SNE2

bh-SNE1

Data: Axel Schulz, PhD

# People are already gating and clustering dimension reduction maps. Guidelines are needed!
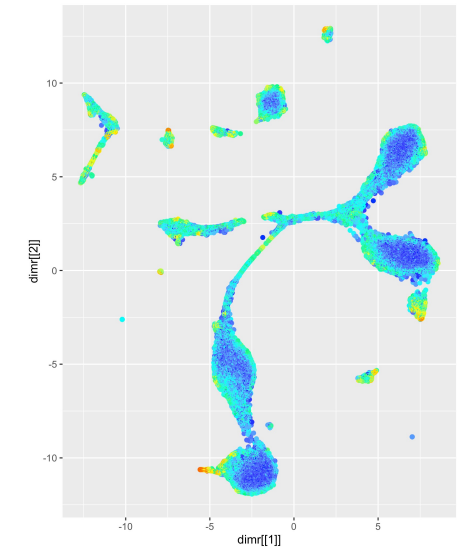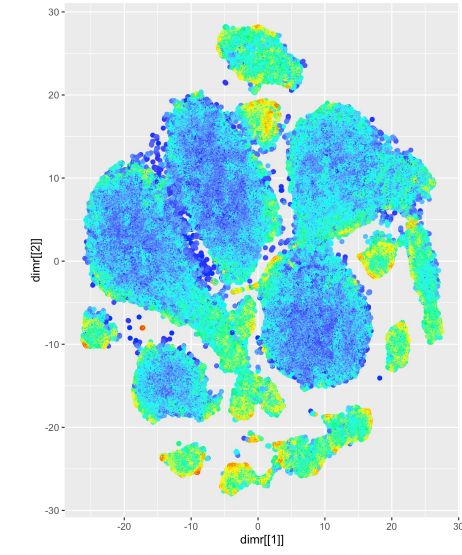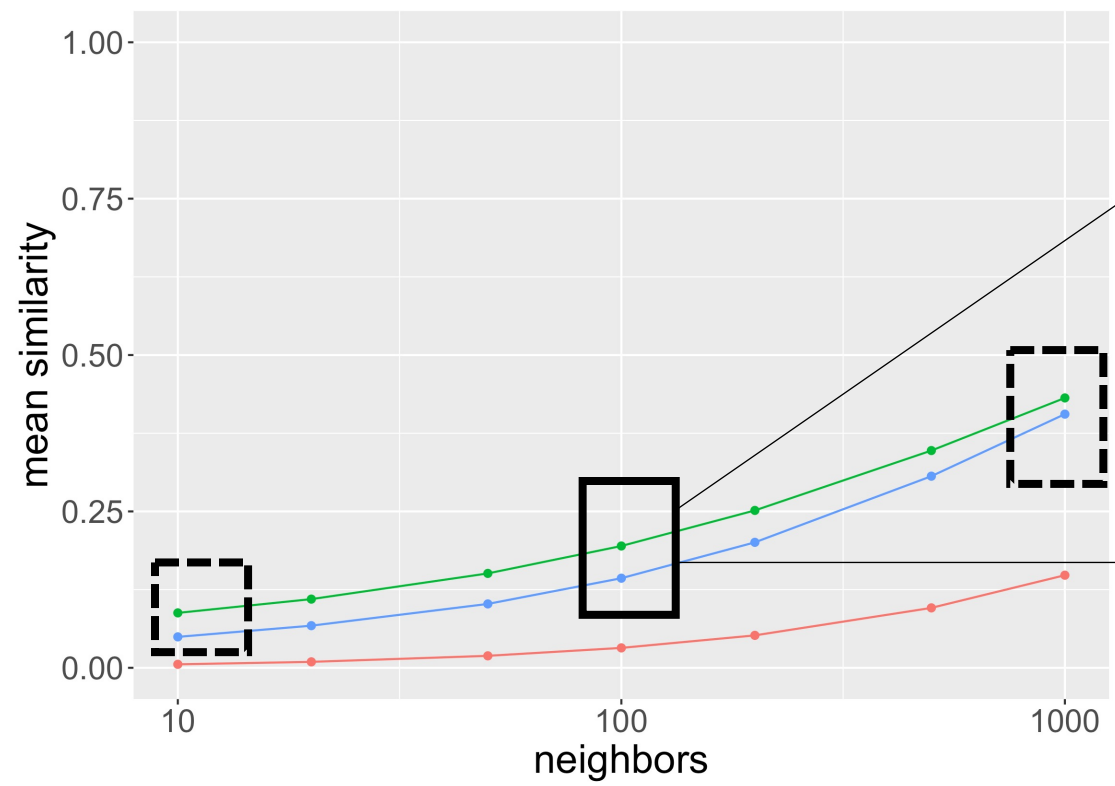
Michael Wong and Evan Newell: Manually gating a t-SNE map



Wong *et al*,
*Cell* 2016

Accense (Petter Brodin): Clustering a t-SNE map



Shekar *et al*,
*PNAS* 2014

# Color a dimension reduction map by it's own neighborhood preservation, given k

# Local comparison for t-SNE

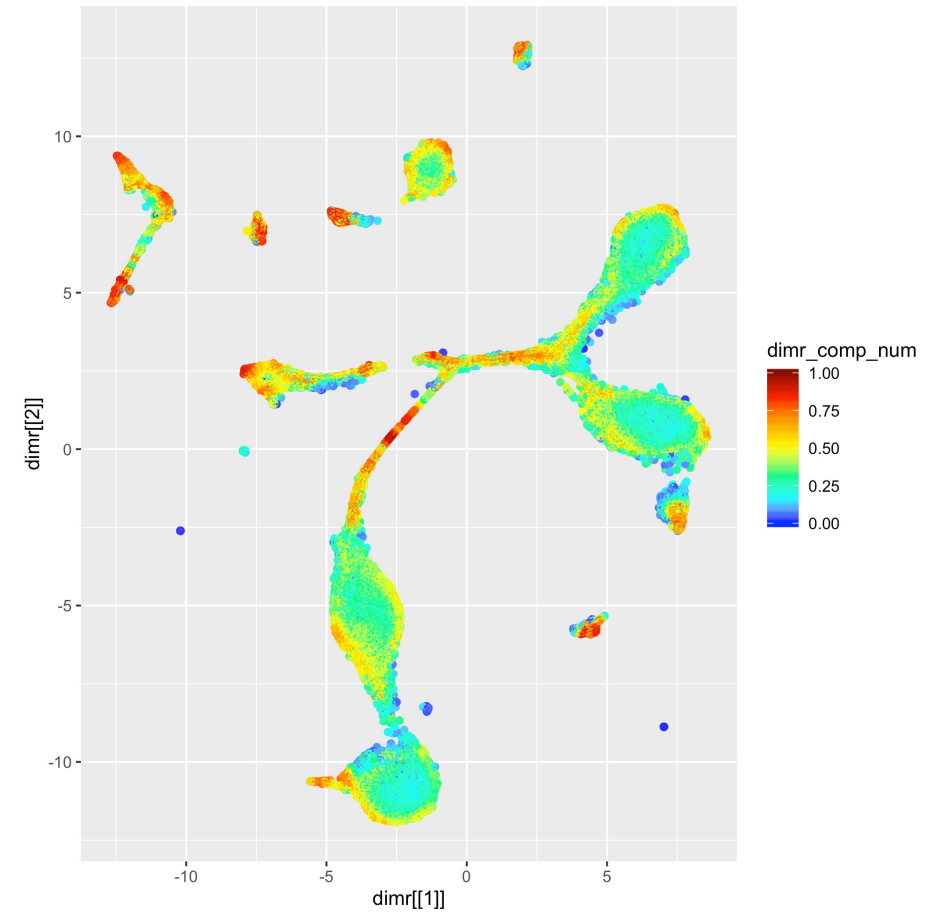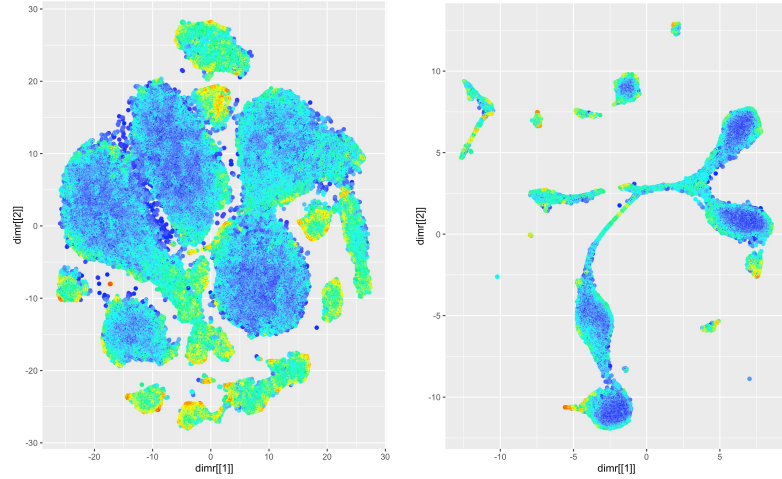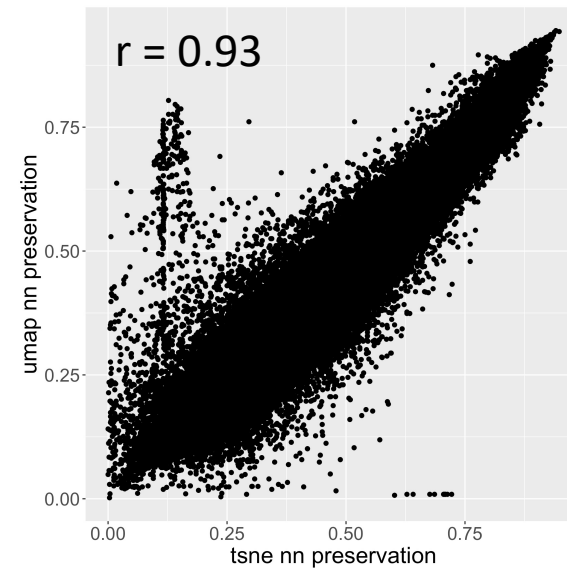K = 10

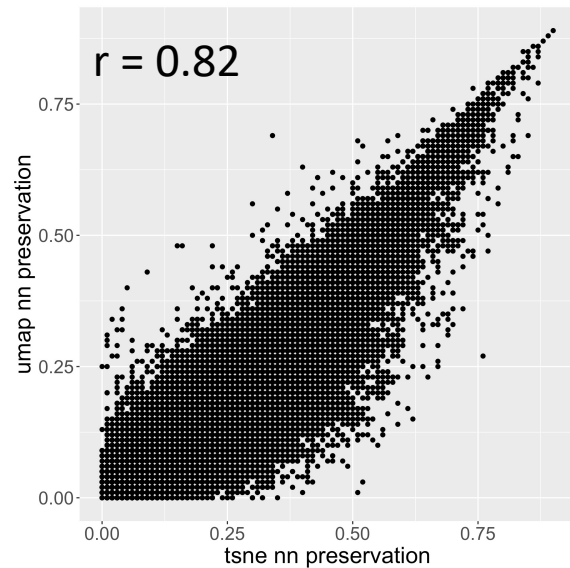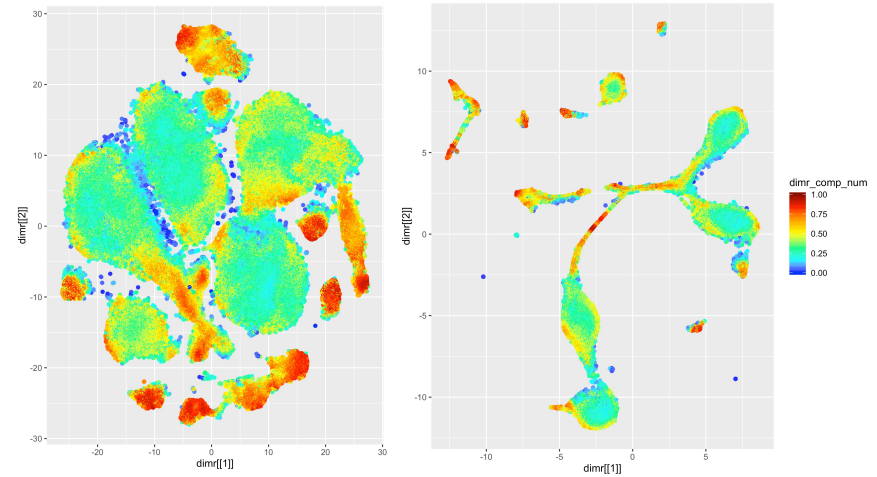K = 100

K = 1000

# Local comparison for UMAP

# t-SNE and UMAP are preserving the data in a similar manner

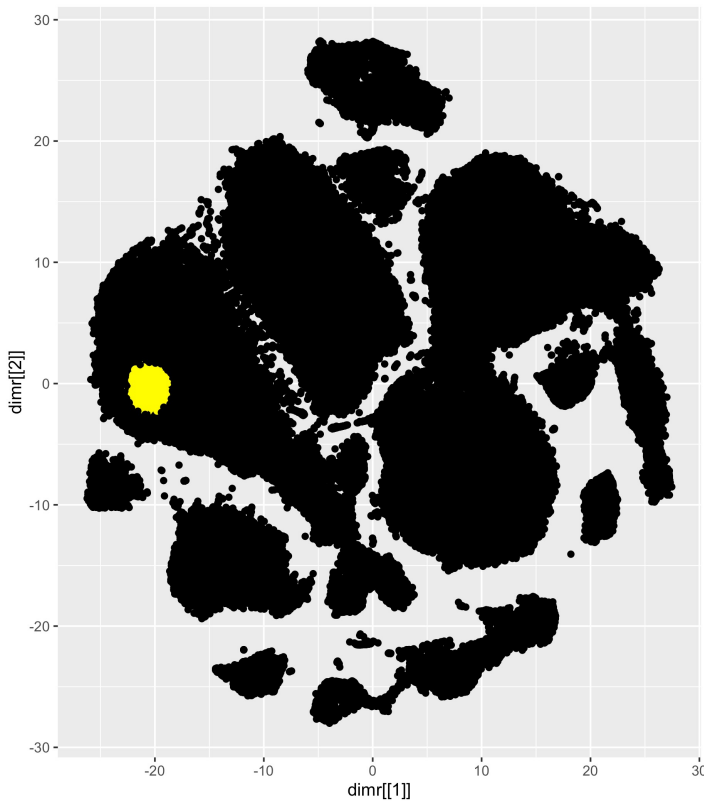K = 100

K = 1000

# Part 1 conclusions

- t-SNE outperforms UMAP (though only slightly) in KNN preservation
- Both t-SNE and UMAP outperform PCA in KNN preservation
- KNN preservation performance varies in specific patterns across both t-SNE and UMAP
- t-SNE and UMAP have better KNN preservation in smaller islands/corridors in the data. Implications on how to gate the maps
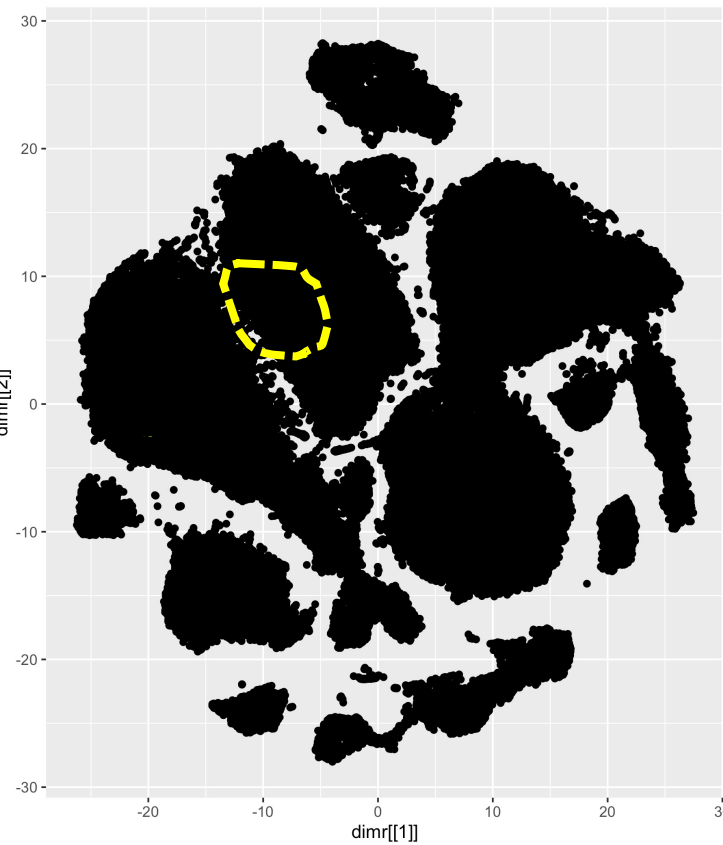
# Outline

- Part 1: Introduction
- Part 2: Preservation of local structure
- **Part 3: Preservation of global structure**

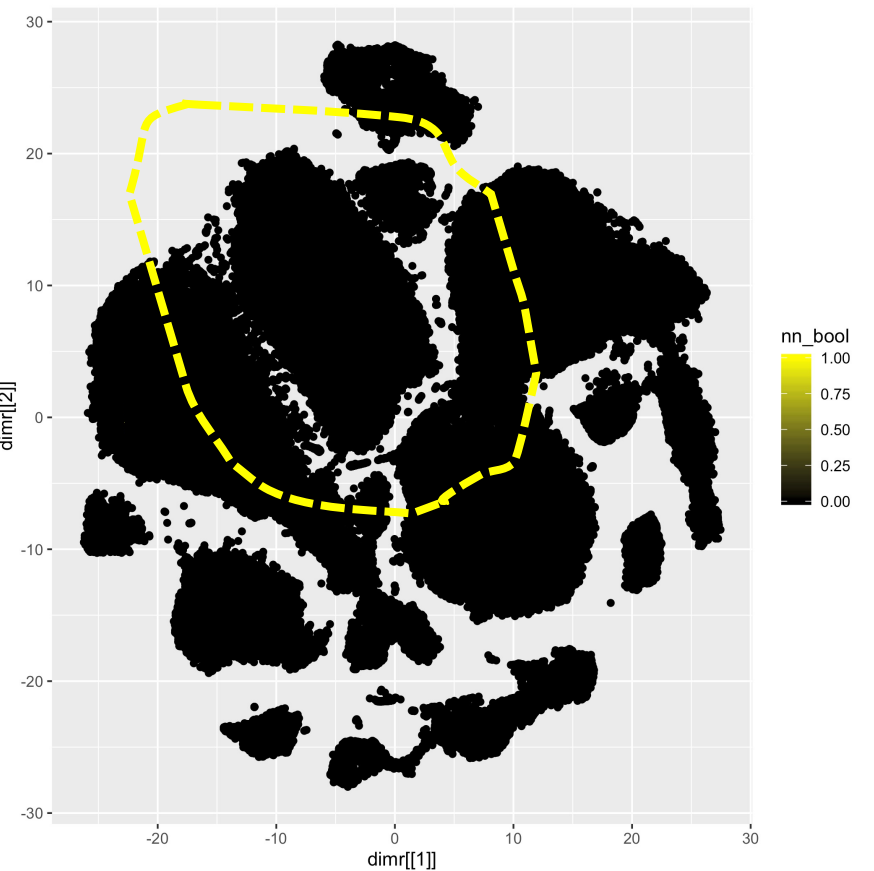# If a "gate" on the map has 30% KNN preservation, where are the other cells?
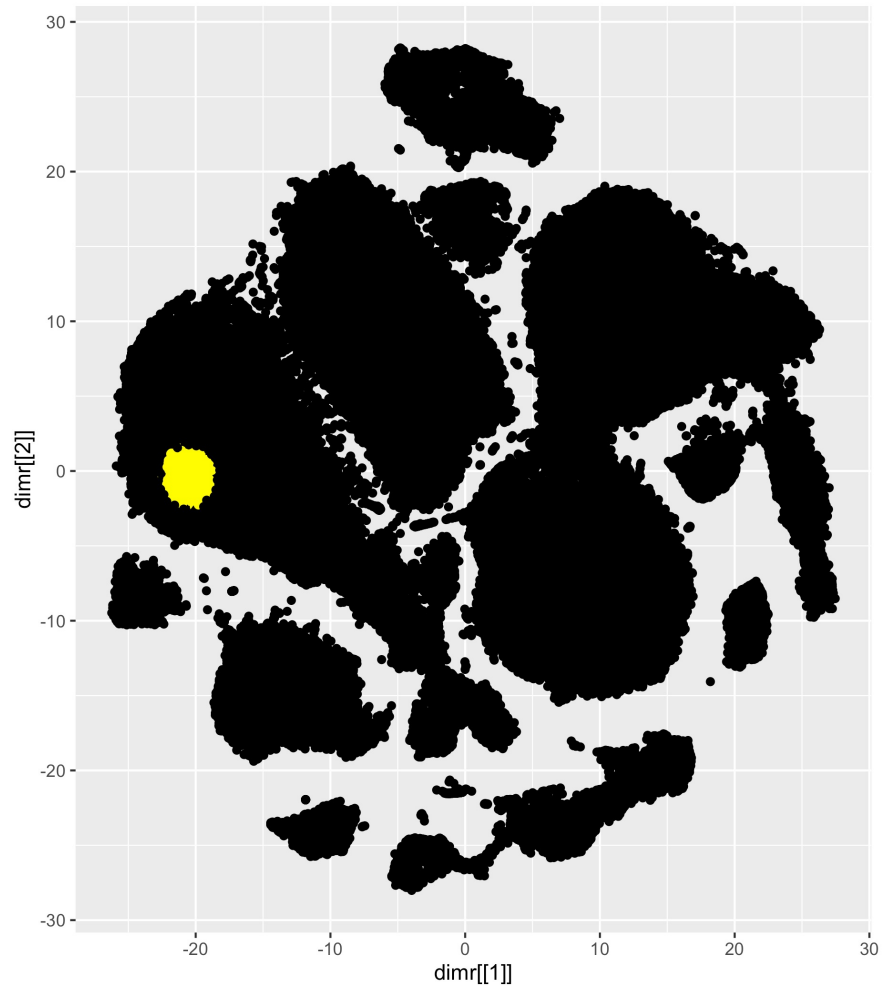
KNN ID on t-SNE map

KNN from hi-D, locations
Hypothesis 1
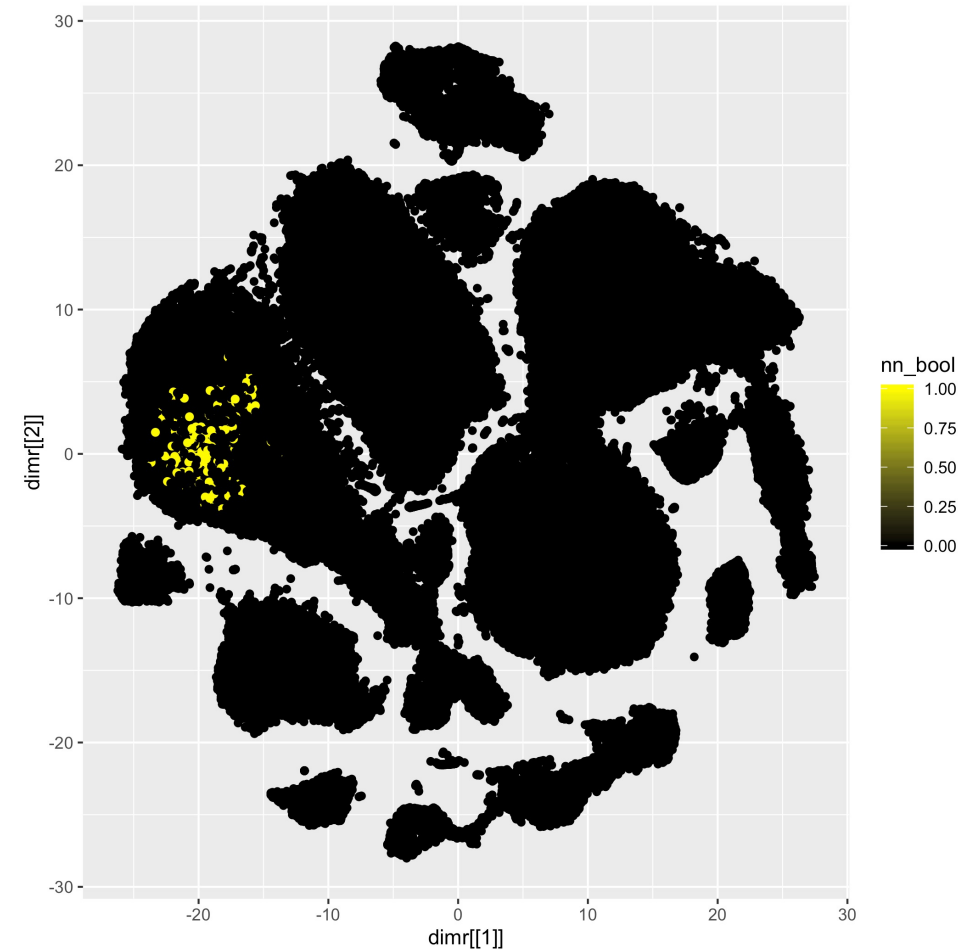
KNN from hi-D, locations
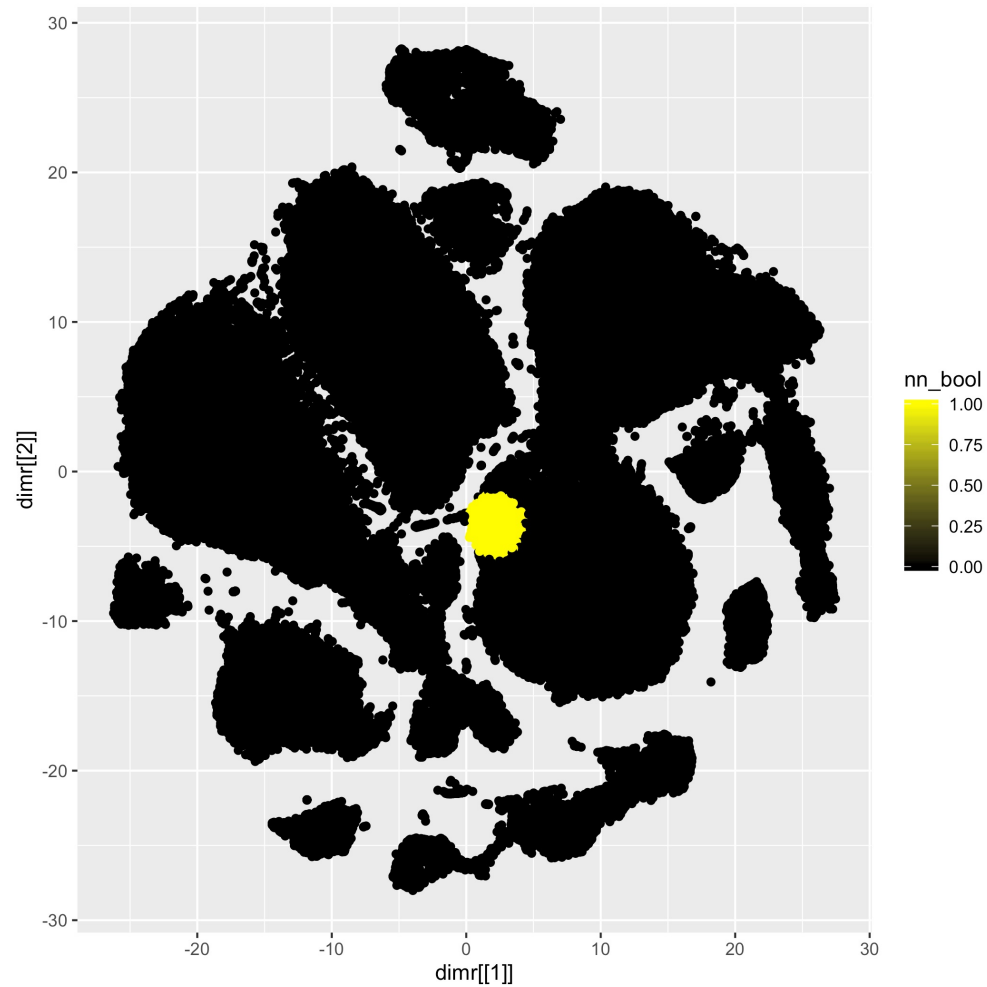Hypothesis 2

# KNN identity for t-SNE, k = 1000, cell 1

# KNN identity for t-SNE, k = 1000, cell 4

# KNN identity for t-SNE, k = 1000, cell 6

# KNN identity for UMAP, k = 1000, cell 1

# KNN identity for UMAP, k = 1000, cell 6

Use my tool knn_sleepwalk
see the feature space KNN
for your own data

# Global preservation, measured by pairwise distances



Dimensionality reduction for visualizing single-cell data using UMAP

Etienne Becht[1], Leland McInnes[2], John Healy[2], Charles-Antoine Dutertre[1], Immanuel W H Kwok[1], Lai Guan Ng[1], Florent Ginhoux[1] & Evan W Newell[1,3]

# K-farthest neighbors (KFN) to determine global preservation of dimension reduction maps
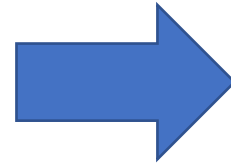
KFN orig    KFN low-D



Find KFN for each cell from high-dim space → Find KFN for each cell from a 2-D embedding (eg tSNE) → Compare KFN identities from the 2-D embedding and high-dim space

Repeat across a wide range of values for K

Bioconductor package: Sconify

# Global KFN comparison between PCA, t-SNE and UMAP (10k cell subsample)

X axis is on a log scale



therefore concerns itself primarily with accurately representing local structure. While we believe that UMAP can capture more global structure than these other techniques, it remains true that if global structure is of primary interest then UMAP may not be the best choice for dimension reduction.

Mcinness *et al, Arxiv* 2018 (the UMAP paper)

# Across 3 datasets, bar plots with error bars and p values

# Part 3 conclusions

- Nearest neighborhoods computed from high-D space and dimension reduction space occupy similar regions

- Positioning of the islands relative to each other could be arbitrary

- K-farthest neighborhood (KFN) preservation reveals global structure preservation: PCA > UMAP > t-SNE

# Next steps: initialization matters

UMAP does not preserve global structure any better than t-SNE when using the same initialization
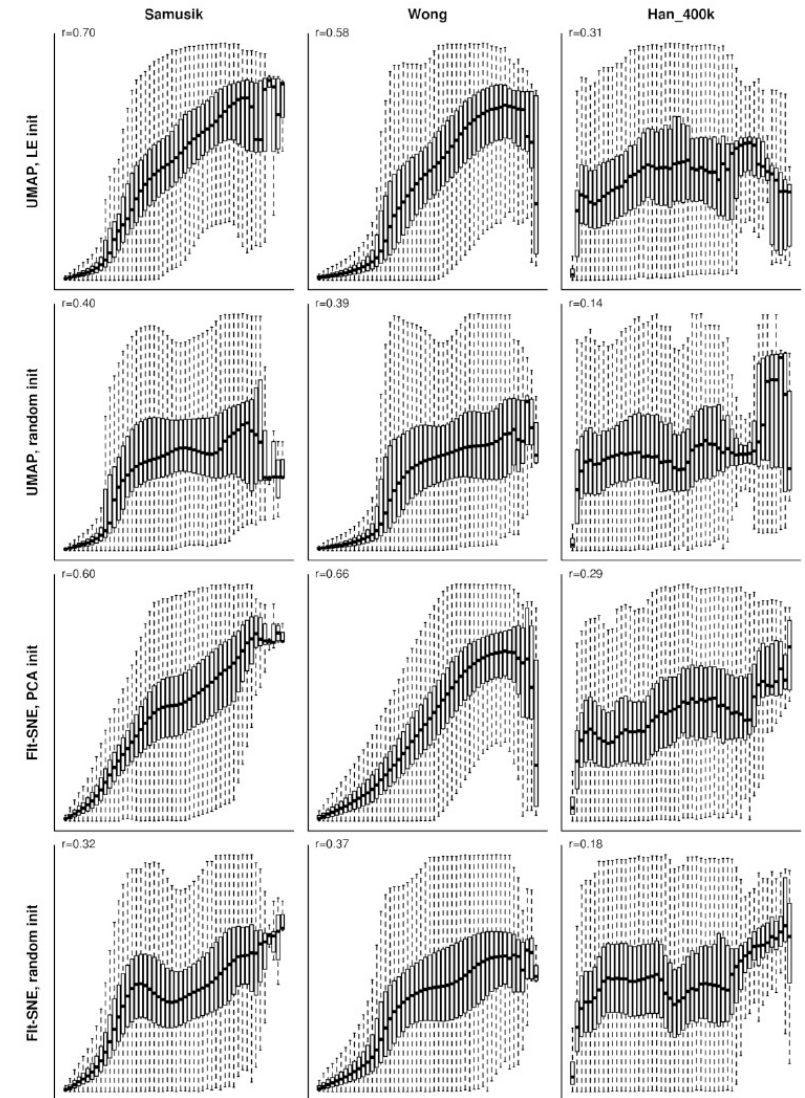
Dmitry Kobak[1] and George C. Linderman[2]

[1] *Institute for Ophthalmic Research, University of Tübingen, Germany*
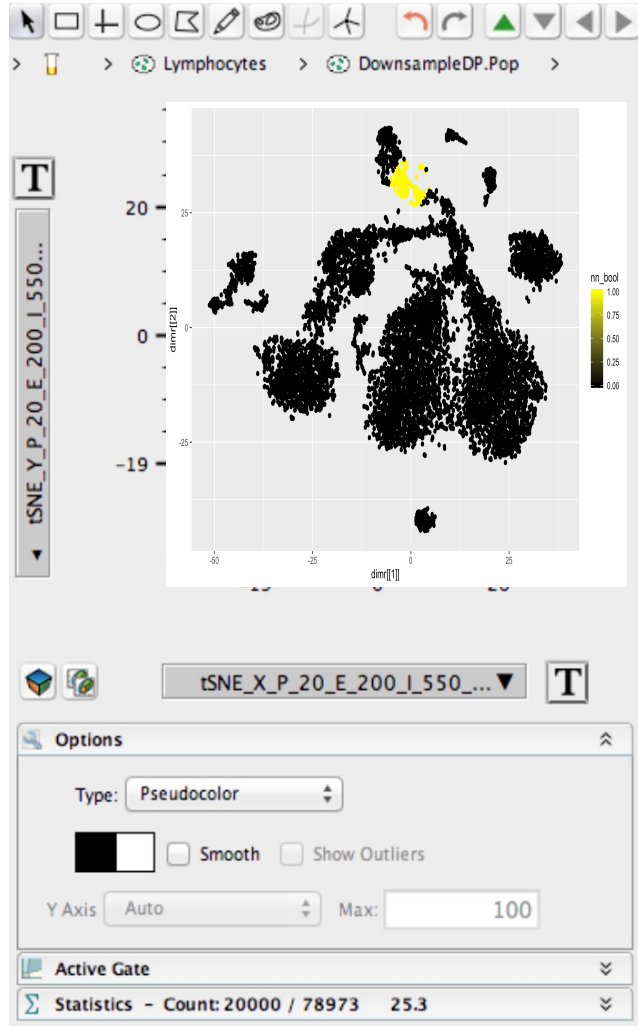[2] *Applied Mathematics Program, Yale University, New Haven, CT, USA*
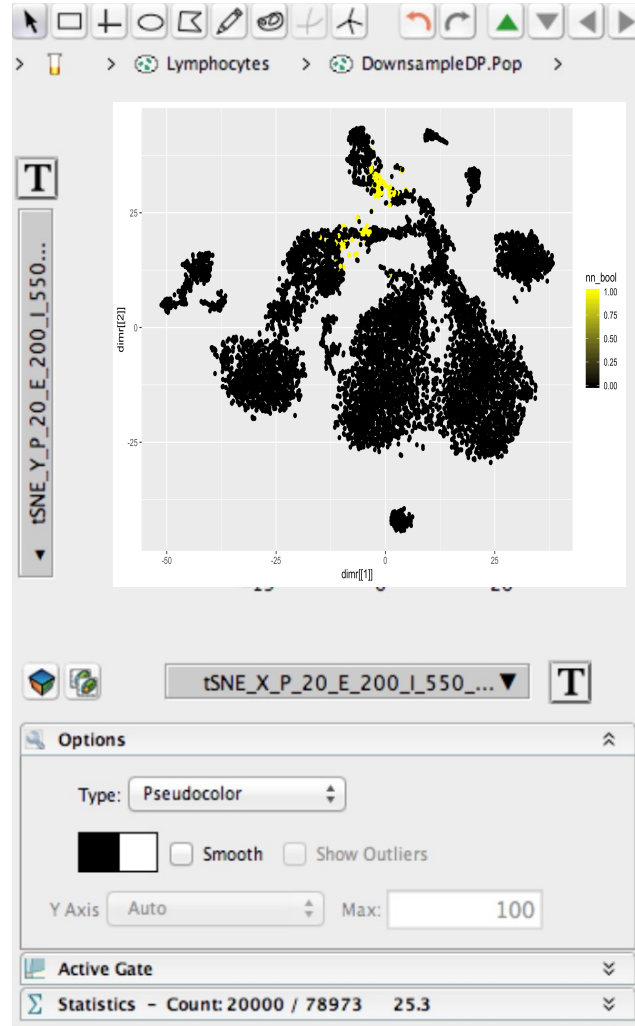dmitry.kobak@uni-tuebingen.de, george.linderman@yale.edu

(BioRxiv)

# Toward a "safe" manual gating interface for dimension reduction maps
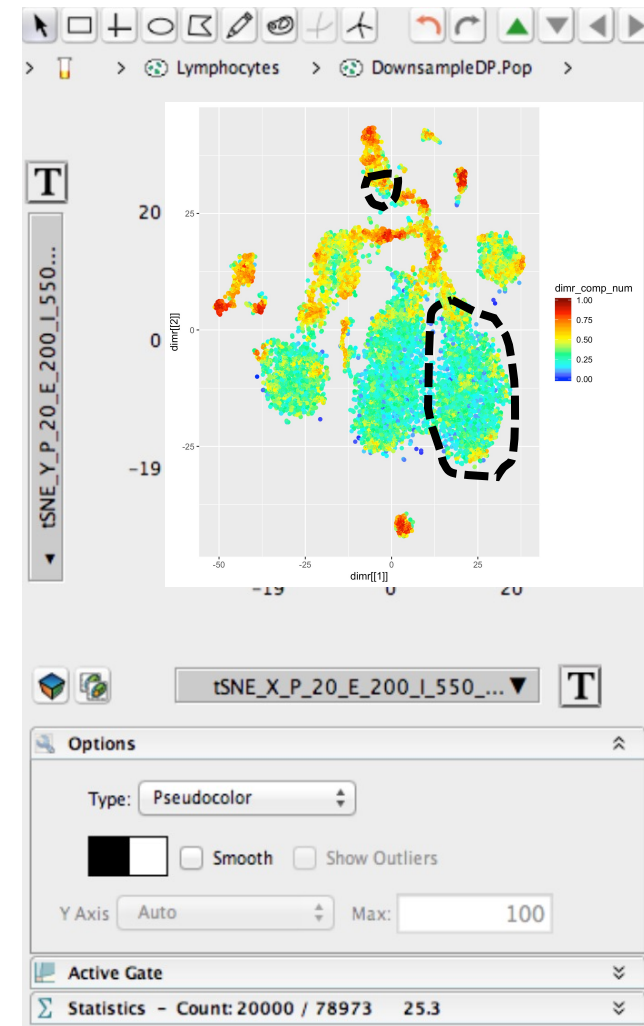
Identity color (dimr)

Identity color (high-d)

Identity comparison

# nnvis: an R package to do neighbor-based preservation analysis on your dimension reductions

## tjburns08/nnvis: Make KNN-based identity comparisons between different manifolds (eg. original space vs t-SNE space)

This package examines the quality of a low-dimensional embedding by comparing the membership of each cell's k-nearest neighbors (KNN) in original high dimensional marker space with this cell's KNN in the low-dimensional space. Comparisons can be visualized with average fidelity plots for different values of K, or the t-SNE maps themselves can be colored by their own fidelity. The package also provides wrappers for popular low dimensional embeddings.

### Getting started

README.md

### Browse package contents

📄 Vignettes

📄 Man pages

</> API and functions

📂 Files

Search within the tjburns08/nnvis package 🔍

# Acknowledgments

Heike Hirseland    Antonia Niedobitek    René Riedel

Sarah Gräßle    Marie Urbicht    Silke Stanislawiak



Sabine Baumgart    **Andreas Grützkau**    Pawel Durek

**Henrik Mei**    **Axel Schulz**

DRFZ BERLIN
Deutsches Rheuma-Forschungszentrum
Ein Institut der Leibniz-Gemeinschaft