

AI for Genomics: from CNNs and LSTMs to Transformers

Nikolay Oskolkov, Group Leader (PI) at LIOS, Riga, Latvia

Physalia course, 09.09.2025

Session 2b: Autoencoder neural network applications for Cell Biology and Data Integration



@NikolayOskolkov



@oskolkov.bsky.social



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

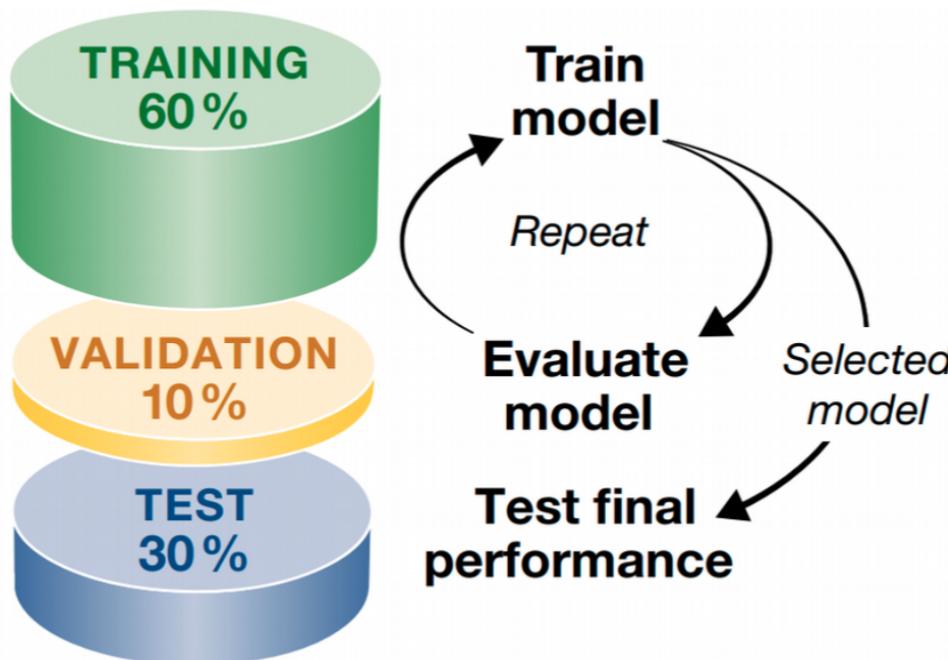
Topics we'll cover in this session:

- 1) Introduction to unsupervised machine learning
- 2) Autoencoder neural network architecture and principles
- 3) Introduction to advances in single cell biology
- 4) Autoencoders for dimension reduction in single cell biology
- 5) Autoencoders for integration of heterogeneous data

$Y = f(X)$, where X is input (data) and Y is output (response)

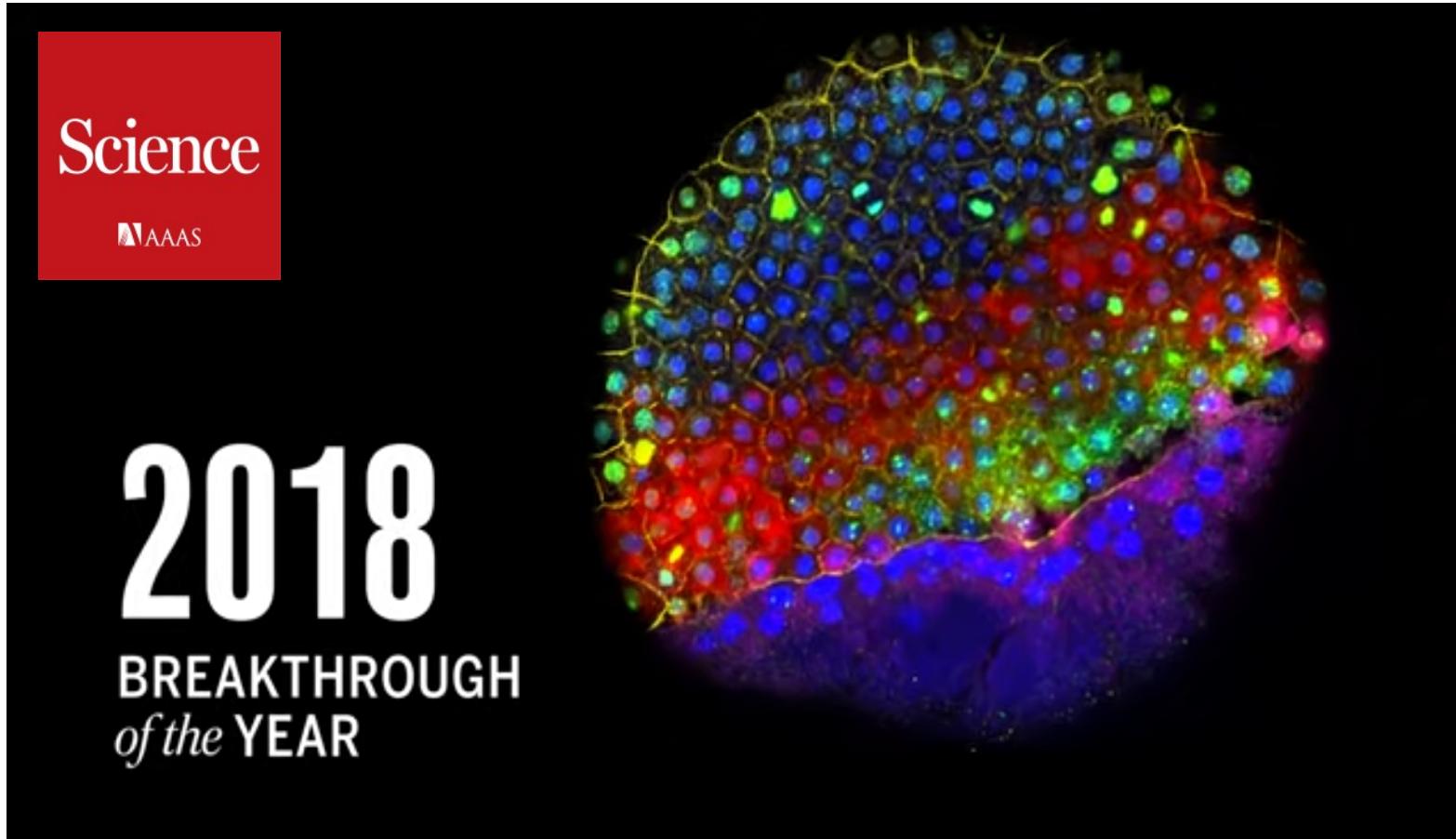
Y is present – supervised machine learning

Y is absent – unsupervised machine learning



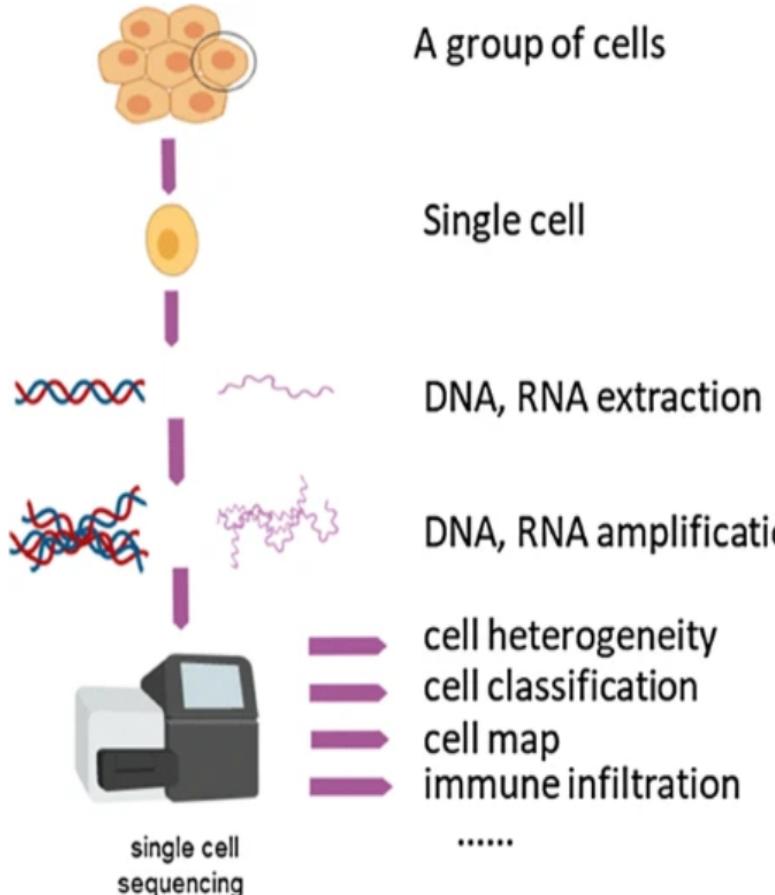
Machine Learning typically involves five basic steps:

1. Split data set into train, validation and test subsets
2. Fit the model on the train subset
3. Validate your model on the validation subset
4. Repeat train - validation split many times and tune hyperparameters
5. Test the accuracy of the optimized model on the test subset.

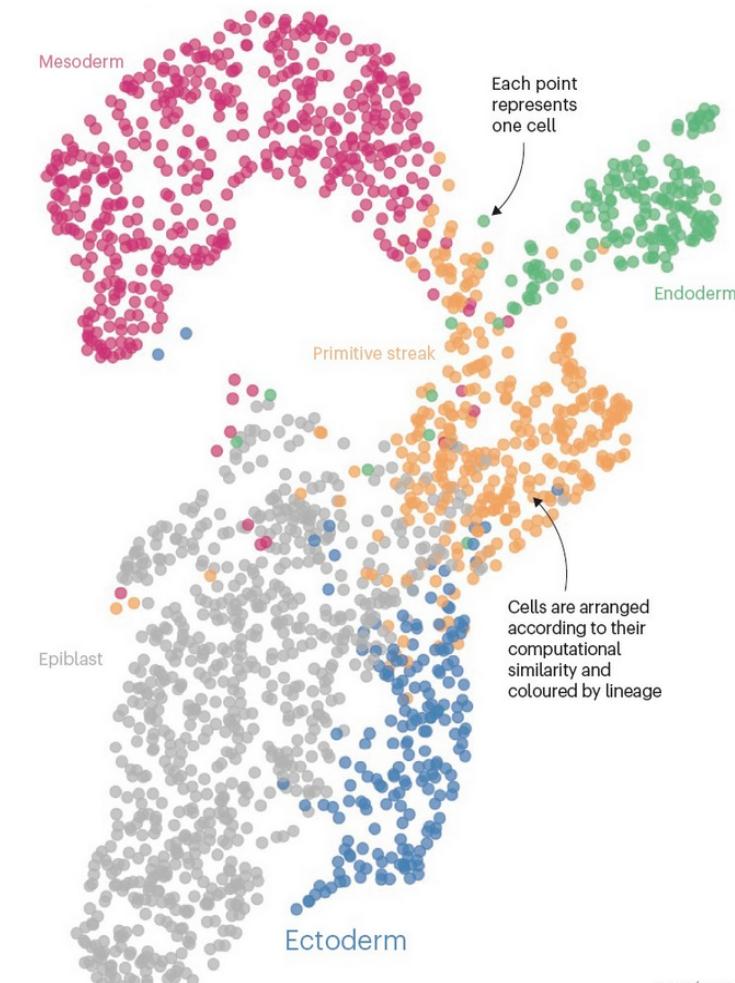


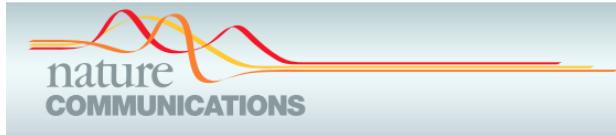
Nature method of the year: 2019 (single cell multi-omics), 2020 (spatial transcriptomics)

Image adapted from Science, 362, 6421, 2018



The principle of single-cell sequencing. It is a process of isolating a single cell for sequencing and studying cell heterogeneity, molecular mapping, immune infiltration and epigenetic changes





ARTICLE

DOI: 10.1038/s41467-018-07582-3

OPEN

Spatially and functionally distinct subclasses of breast cancer-associated fibroblasts revealed by single cell RNA sequencing

Michael Bartoschek¹, Nikolay Oskolkov², Matteo Bocci¹, John Lövrot¹, Christer Larsson¹, Mikael Sommarin⁴, Chris D. Madsen¹, David Lindgren¹, Gyula Pekar⁵, Göran Karlsson⁴, Markus Ringnér¹, Jonas Bergh³, Åsa Björklund¹ & Kristian Pietras¹

Cancer-associated fibroblasts (CAFs) are a major constituent of the tumor microenvironment, although their origin and roles in shaping disease initiation, progression and treatment response remain unclear due to significant heterogeneity. Here, following a negative selection strategy combined with single-cell RNA sequencing of 768 transcriptomes of mesenchymal cells from a genetically engineered mouse model of breast cancer, we define three distinct subpopulations of CAFs. Validation at the transcriptional and protein level in several experimental models of cancer and human tumors reveal spatial separation of the CAF subclasses attributable to different origins, including the peri-vascular niche, the mammary fat pad and the transformed epithelium. Gene profiles for each CAF subtype correlate to distinctive functional programs and hold independent prognostic capability in clinical cohorts by association to metastatic disease. In conclusion, the improved resolution of the widely defined CAF population opens the possibility for biomarker-driven development of drugs for precision targeting of CAFs.

Oskolkov et al. *Skeletal Muscle* (2022) 12:16
<https://doi.org/10.1186/s1395-022-00299-4>

Skeletal Muscle

RESEARCH

Open Access



High-throughput muscle fiber typing from RNA sequencing data

Nikolay Oskolkov^{1,2}, Małgorzata Santel³, Hemang M. Parikh⁴, Ola Ekström¹, Gray J. Camp³, Eri Miyamoto-Mikami⁵, Kristoffer Ström^{1,6}, Bilal Ahmad Mir¹, Dmytro Kryvokhyza¹, Mikko Lehtovirta^{1,7}, Hiroyuki Kobayashi⁸, Ryo Kakigi⁹, Hisashi Naito⁵, Karl-Fredrik Eriksson¹, Björn Nystedt¹⁰, Noriyuki Fuku⁵, Barbara Treutlein³, Svante Pääbo^{3,11} and Ola Hansson^{1,2*}

Abstract

Background: Skeletal muscle fiber type distribution has implications for human health, muscle function, and performance. This knowledge has been gathered using labor-intensive and costly methodology that limited these studies. Here, we present a method based on muscle tissue RNA sequencing data (totRNAseq) to estimate the distribution of skeletal muscle fiber types from frozen human samples, allowing for a larger number of individuals to be tested.

Methods: By using single-nuclei RNA sequencing (snRNAseq) data as a reference, cluster expression signatures were produced by averaging gene expression of cluster gene markers and then applying these to totRNAseq data and inferring muscle fiber nuclei type via linear matrix decomposition. This estimate was then compared with fiber type distribution measured by ATPase staining or myosin heavy chain protein isoform distribution of 62 muscle samples in two independent cohorts ($n = 39$ and 22).

Results: The correlation between the sequencing-based method and the other two were $r_{ATPase} = 0.44$ [0.13–0.67], [95% CI], and $r_{myosin} = 0.83$ [0.61–0.93], with $p = 5.70 \times 10^{-3}$ and 2.00×10^{-6} , respectively. The deconvolution inference of fiber type composition was accurate even for very low totRNAseq sequencing depths, i.e., down to an average of $\sim 10,000$ paired-end reads.

Conclusions: This new method (<https://github.com/OlaHanssonLab/PredictFiberType>) consequently allows for measurement of fiber type distribution of a larger number of samples using totRNAseq in a cost and labor-efficient way. It is now feasible to study the association between fiber type distribution and e.g. health outcomes in large well-powered studies.

Introduction

Our bodies constitute to ~ 30 –40% of the skeletal muscle, and it is the most abundant form of the three types of muscle, the others being smooth and cardiac. The skeletal muscle is composed of different fiber types (i.e., muscle cell types), and the relative proportions of these types vary among the muscles, locations within the

muscles, individuals, and the sex of individuals [1–4]. The oxidative and glycolytic potential and the contractile properties differ considerably between fiber types, with the mitochondria-rich slow-twitch fibers (type I) having higher oxidative capacity, and fast-twitch fibers (type IIa and type IIx) having higher glycolytic capacity [1]. The proportions also change as people age, with type II fibers being preferentially affected by sarcopenia [5]. Exercising the skeletal muscle is a major site for catabolic metabolism of the blood glucose and lipids and the metabolic characteristics of this tissue influence both the

*Correspondence: Ola.Hansson@med.lu.se

¹Department of Clinical Sciences, Lund University, Malmö, Sweden
 Full list of author information is available at the end of the article



© The Author(s) 2022. **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

PERSPECTIVE

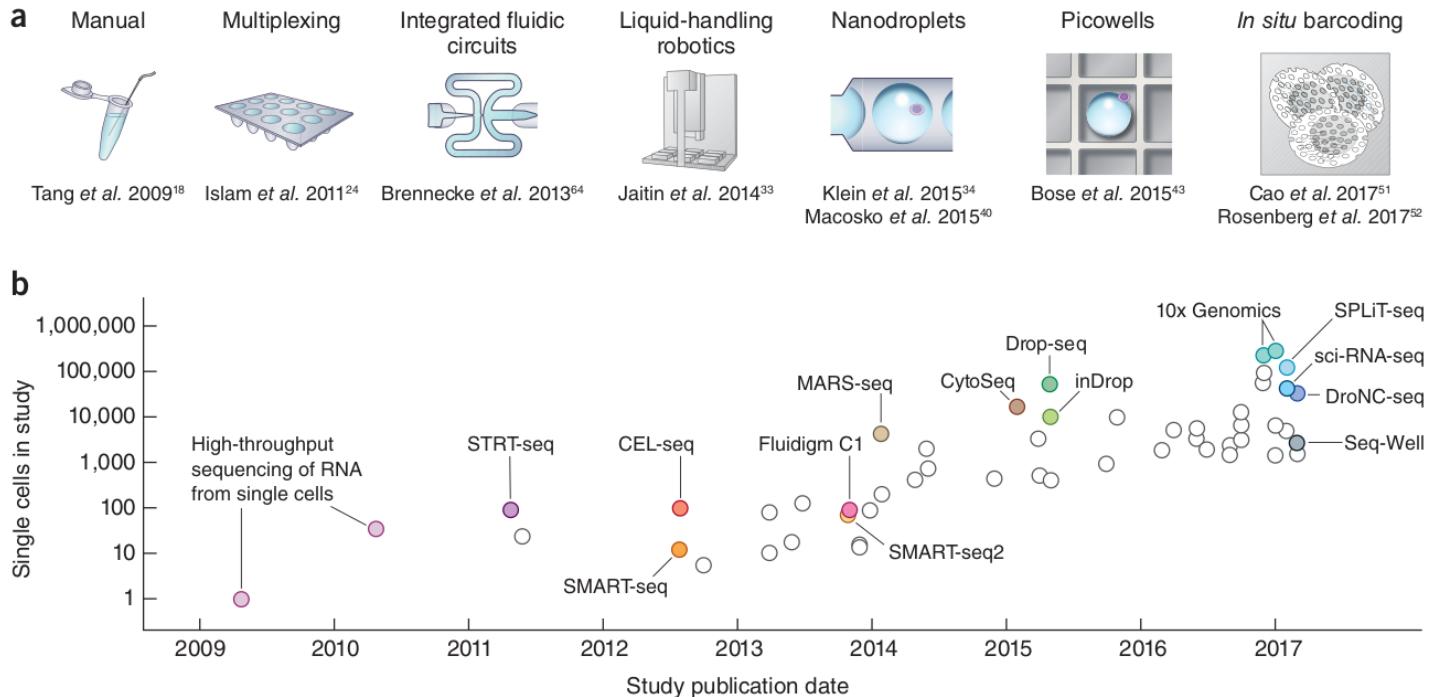


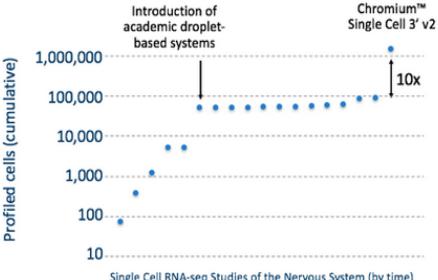
Figure 1 | Scaling of scRNA-seq experiments. (a) Key technologies that have allowed jumps in experimental scale. A jump to ~100 cells was enabled by sample multiplexing, and then a jump to ~1,000 cells was achieved by large-scale studies using integrated fluidic circuits, followed by a jump to several thousands of cells with liquid-handling robotics. Further orders-of-magnitude increases bringing the number of cells assayed into the tens of thousands were enabled by random capture technologies using nanodroplets and picowell technologies. Recent studies have used *in situ* barcoding to inexpensively reach the next order of magnitude of hundreds of thousands of cells. (b) Cell numbers reported in representative publications by publication date. Key technologies are indicated. A full table with corresponding numbers is available as [Supplementary Table 1](#).

Available online at www.sciencedirect.com

ScienceDirect

[« Back to Blog](#)

< Newer Article Older Article >



Our 1.3 million single cell dataset is ready 0 KUDOS



POSTED BY grace-10x on Feb 21, 2017 at 2:28 PM

At ASHG last year, we announced our 1.3 Million Brain Cell Dataset, which is, to date, the largest dataset published in the single cell RNA-sequencing (scRNA-seq) field. Using the Chromium™ Single Cell 3' Solution (v2 Chemistry), we were able to sequence and profile 1,308,421 individual cells from embryonic mice brains. Read more in our application note [Transcriptional Profiling of 1.3 Million Brain Cells with the Chromium™ Single Cell 3' Solution](#).

Single cells make big data: New challenges and opportunities in transcriptomics

Philipp Angerer¹, Lukas Simon¹, Sophie Tritschler¹, F. Alexander Wolf¹, David Fischer¹ and Fabian J. Theis^{1,2}

Abstract

Recent technological advances have enabled unprecedented insight into transcriptomics at the level of single cells. Single cell transcriptomics enables the measurement of transcriptomic information of thousands of single cells in a single experiment. The volume and complexity of resulting data make it a paradigm of big data. Consequently, the field is presented with new scientific and, in particular, analytical challenges where currently no scalable solutions exist. At the same time, exciting opportunities arise from increased resolution of single-cell RNA sequencing data and improved statistical power of ever growing datasets. Big single cell RNA sequencing data promises valuable insights into cellular heterogeneity which may significantly improve our understanding of biology and human disease. This review focuses on single cell transcriptomics and highlights the inherent opportunities and challenges in the context of big data analytics.

Addresses

¹ Institute of Computational Biology, Helmholtz Center Munich, German Research Center for Environmental Health, Neuherberg, Germany

² Department of Mathematics, Technical University of Munich, Garching, Germany

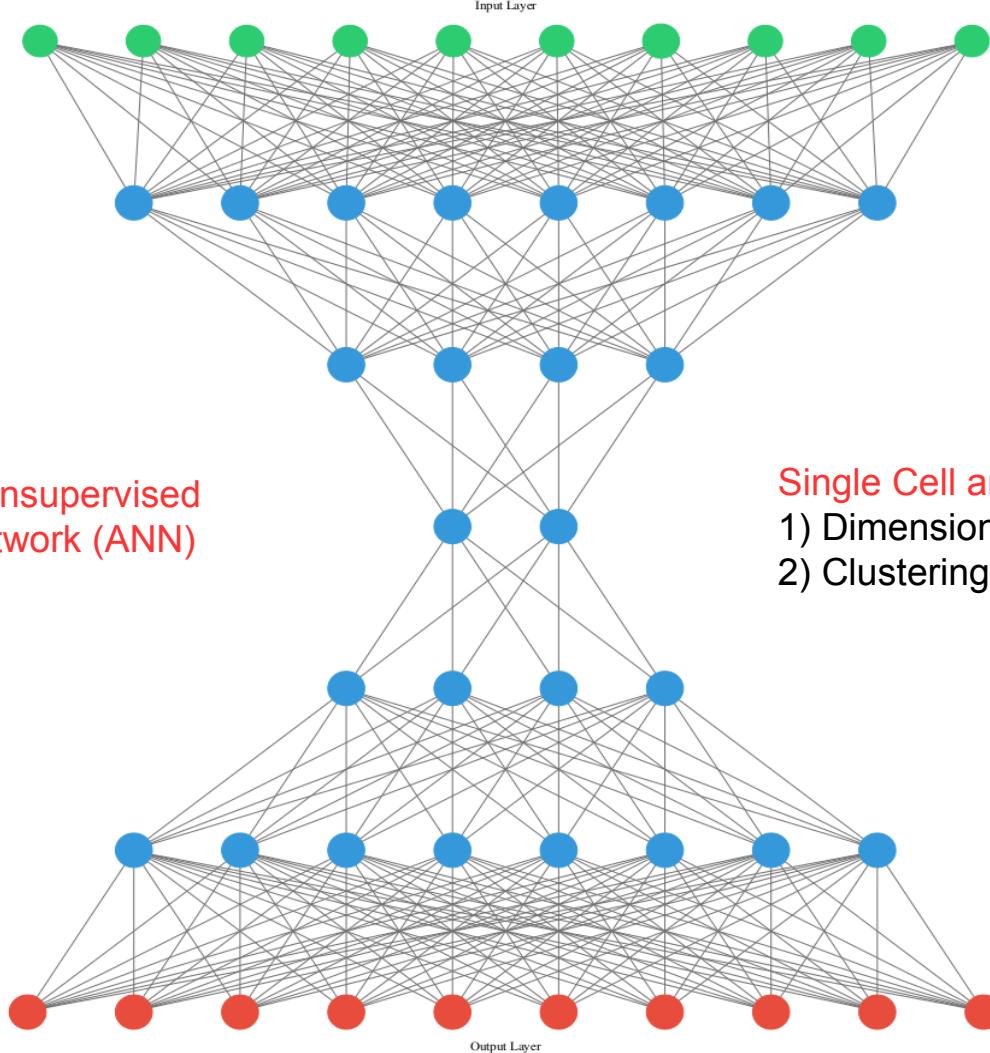
Corresponding author: Theis, Fabian J. Institute of Computational Biology, Helmholtz Center Munich, German Research Center for Environmental Health, Neuherberg, Germany. (fabian.theis@helmholtz-muenchen.de)

Current Opinion in Systems Biology 2017, 4:85–91

This review comes from a themed issue on **Big data acquisition and**

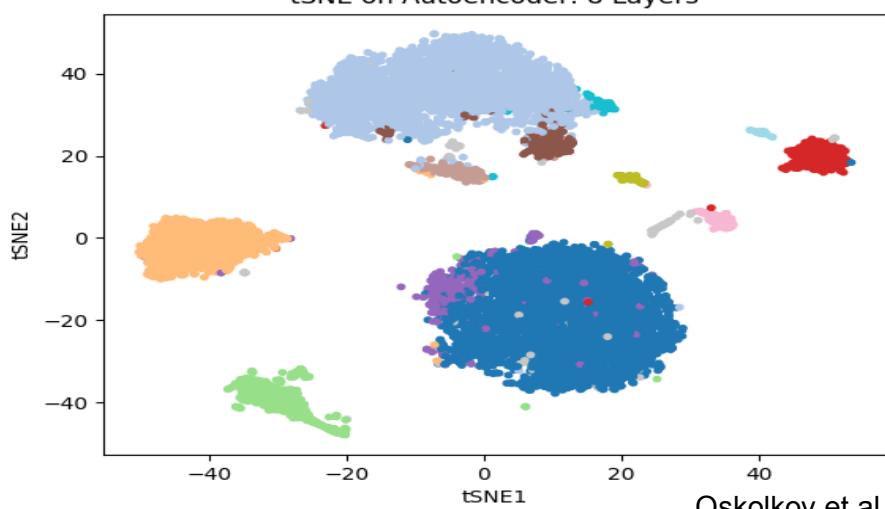
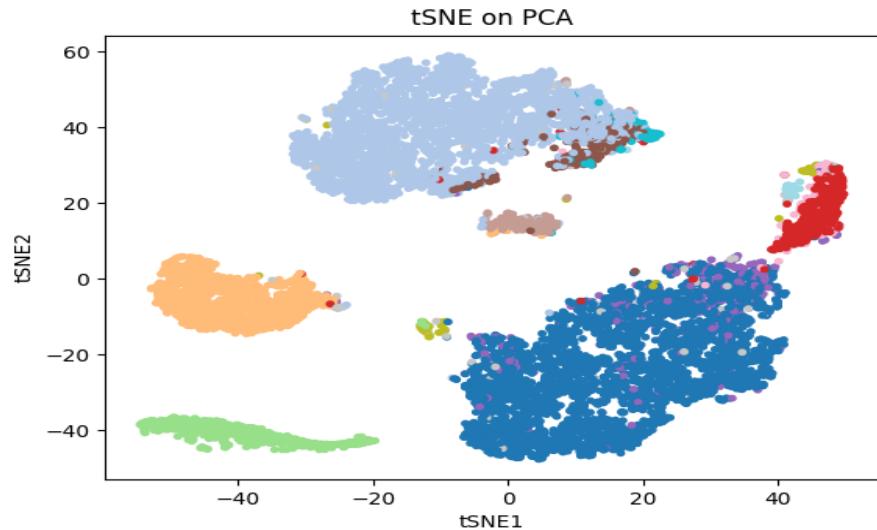
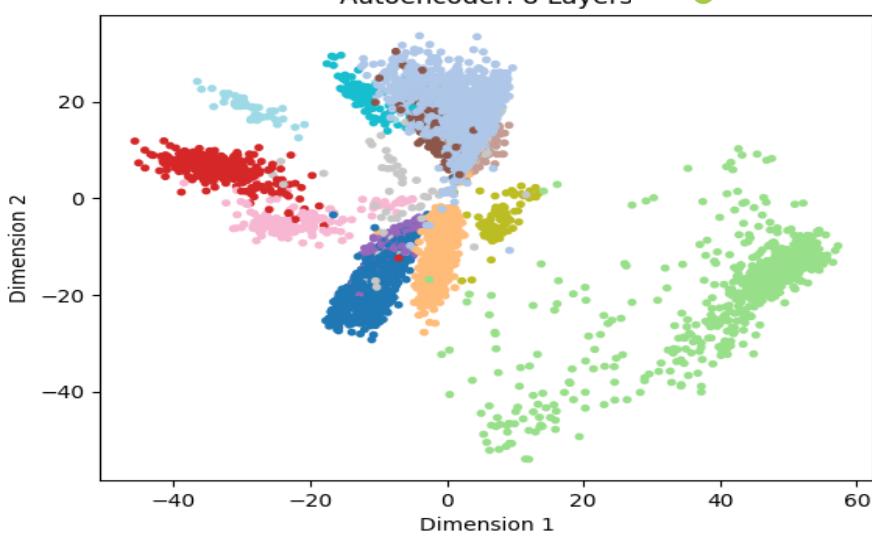
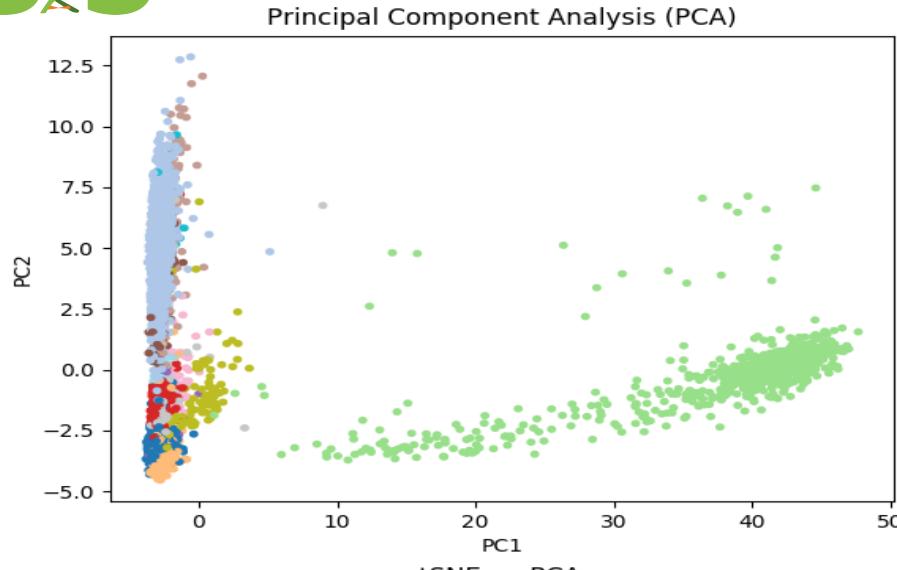
volume – the amount of data, velocity – the required processing speed, veracity – trustworthiness and availability, and variety – necessary model complexity [2]. The traditional scientific big data field is astronomy because of the huge *volume* of image data produced by telescopes with a high daily *velocity* [3]. Big data has also reached biology, mainly driven through the advent of next generation sequencing technology. For biologists, assessing *veracity* through statistical means is nothing new.

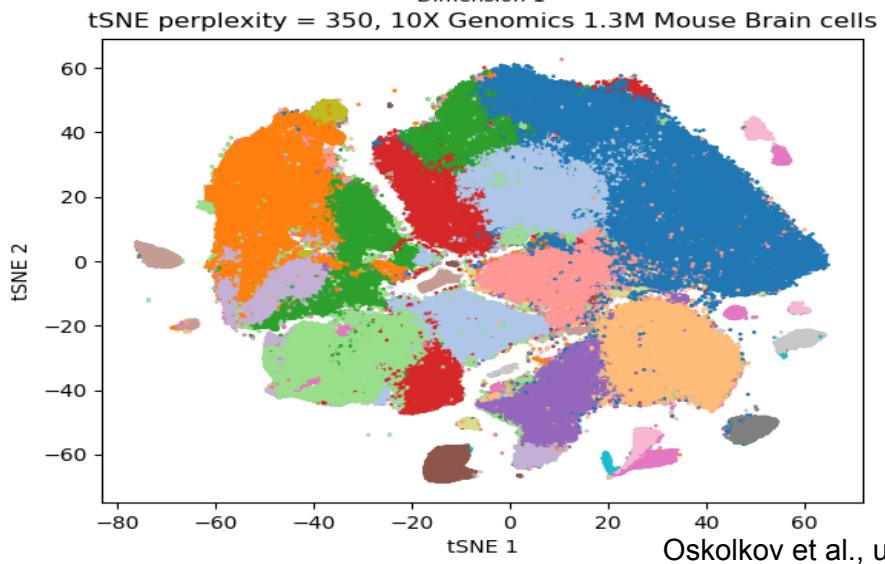
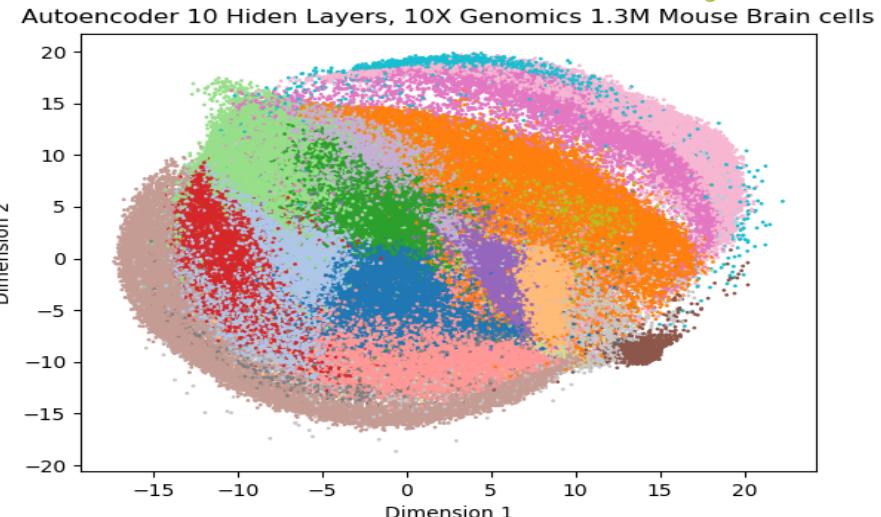
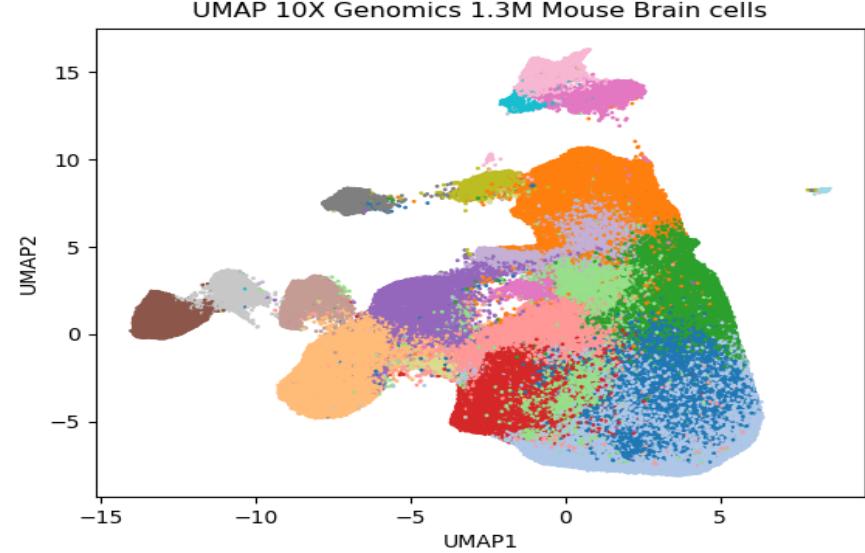
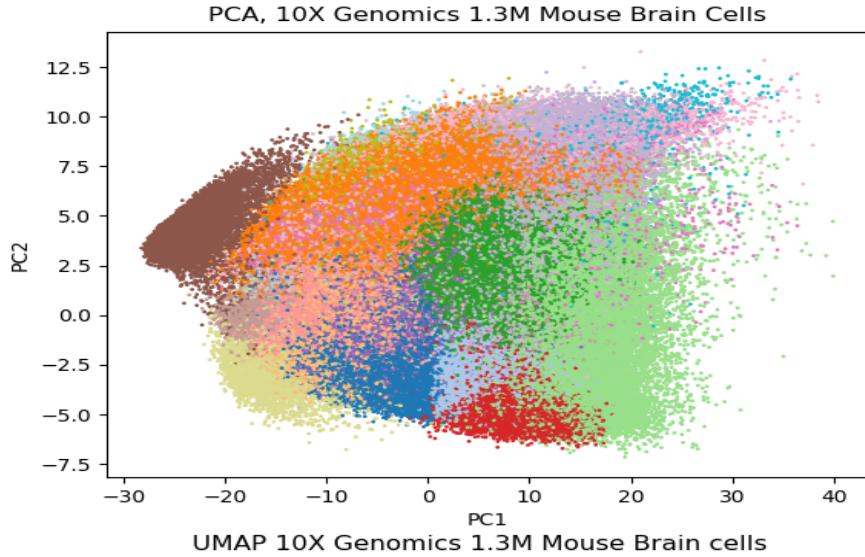
Recent technological advances now allow the profiling of single cells at a *variety* of omic layers (genomes, epigenomes, transcriptomes and proteomes) at an unprecedented level of resolution [4]. Single cell transcriptomics (SCT) entails the profiling of all messenger RNAs present in a single cell and constitutes the most widely-used sc profiling technology [4]. Unlike bulk RNA-seq profiling where sequencing libraries are generated from thousands of cells, scRNA-seq technologies isolate single cells and generate cell-specific sequencing libraries (e.g. Fluidigm [5]) mark RNA content with a cell-specific molecular barcode [6–9]. Both approaches generate gene expression estimates at the single cell level [10]. SCT enables, for the first time, the measurement of the transcriptomic information of thousands, and up to millions of single cells, in a single experiment [7]. The complexity of SCT data coupled with the massive volume inherent to next generation sequencing data makes it a paradigm of big data.



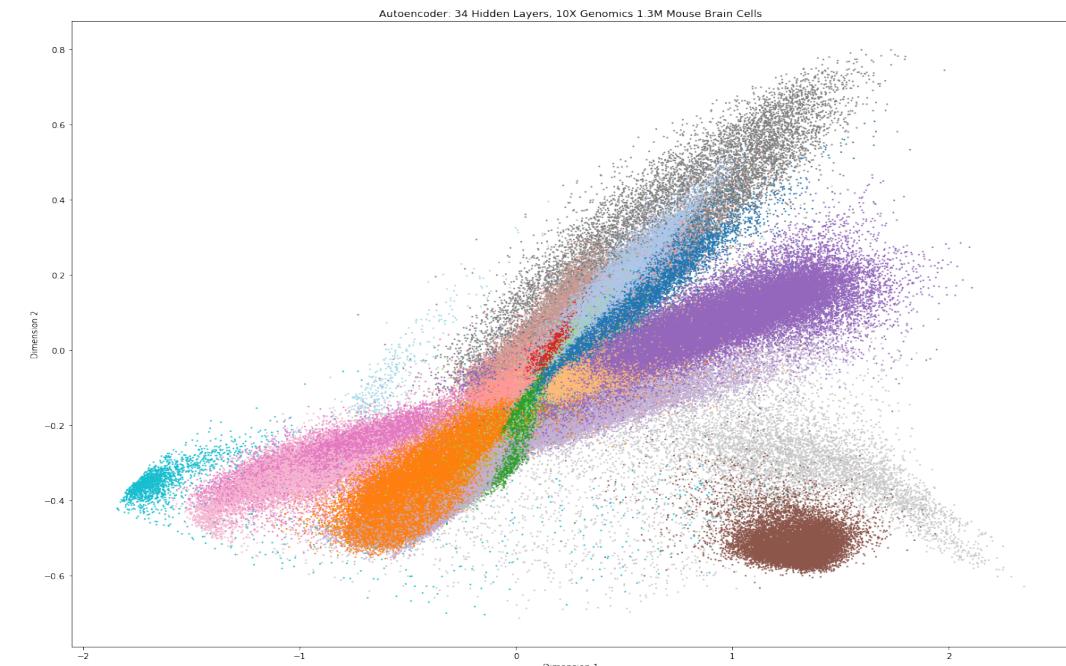
Autoencoder is an unsupervised
Artificial Neural Network (ANN)

Single Cell analysis is unsupervised
1) Dimensionality reduction: visualization
2) Clustering of cells: discover cell populations



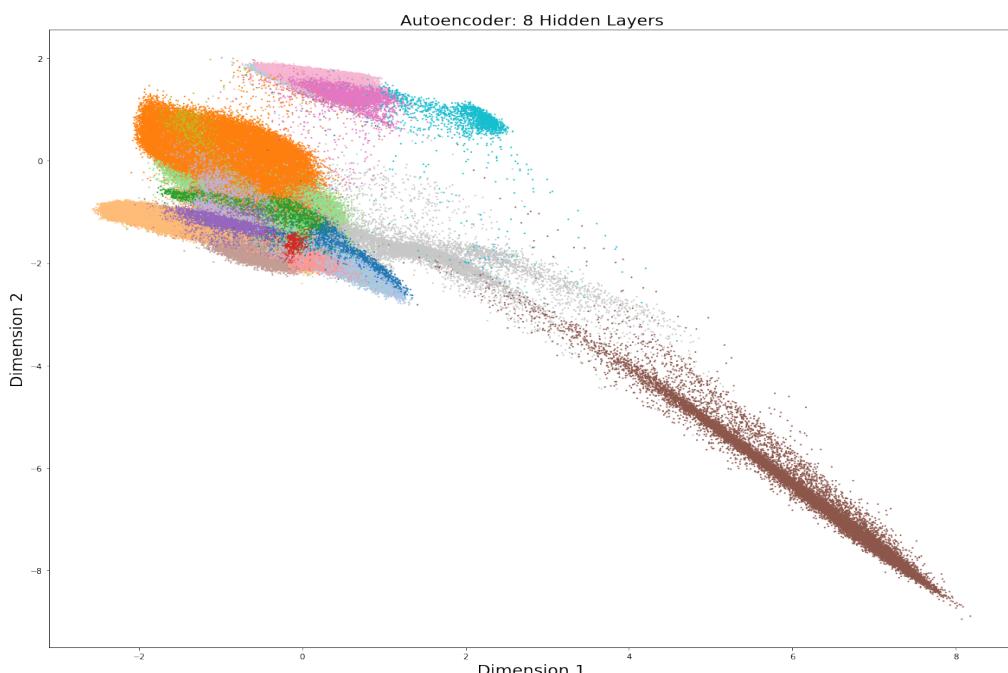


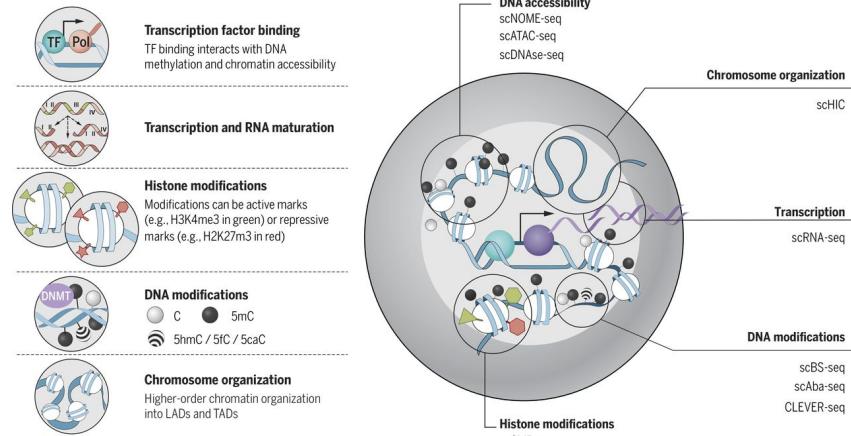
10X Genomics mouse brain 1.3M cells: applying Autoencoders with Tensorflow and Keras



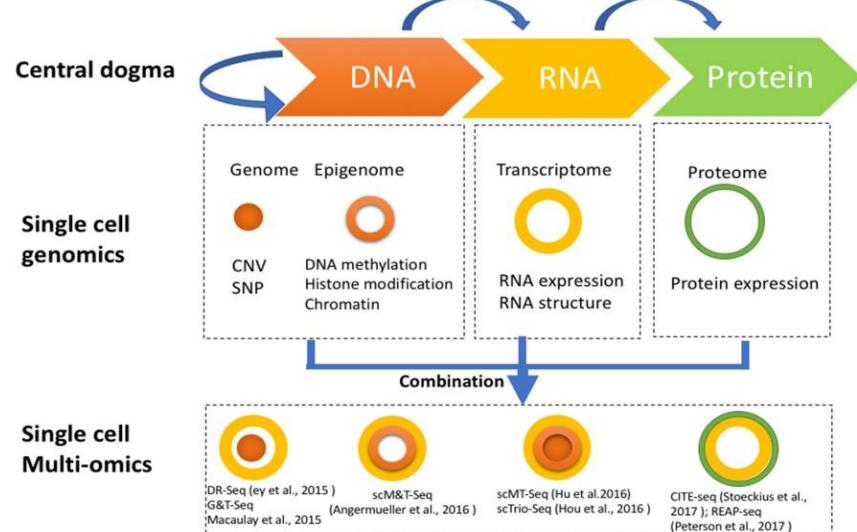
Autoencoders themselves are perhaps
not optimal for visualization of scOmics

Autoencoders can be promising for non-linear
data pre-processing, the bottleneck can potentially
be fed to tSNE / UMAP

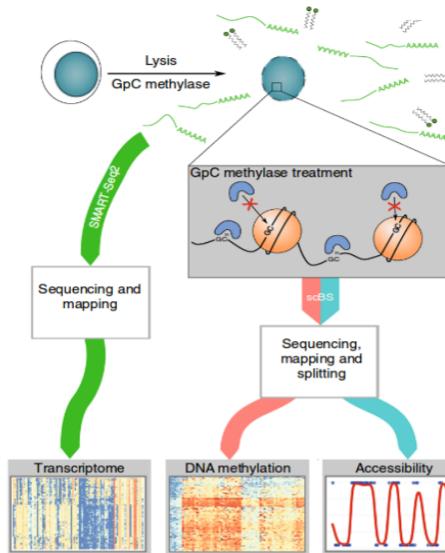




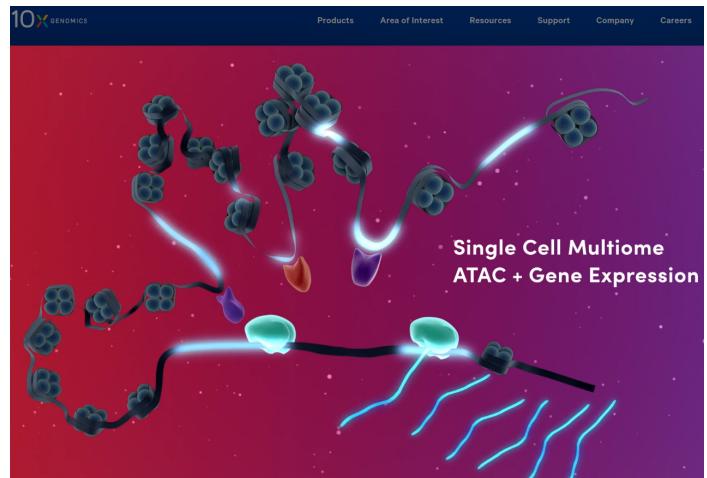
Kelsey et al., 2017, Science 358, 69-75

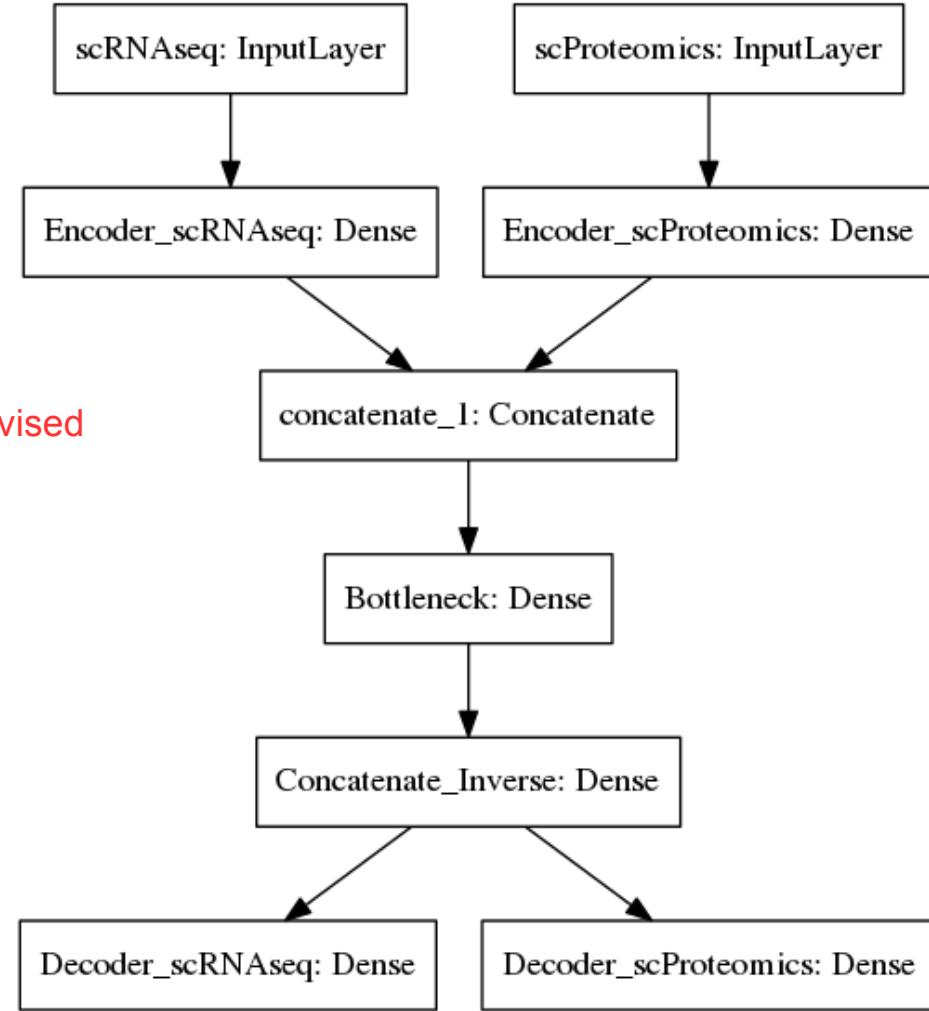
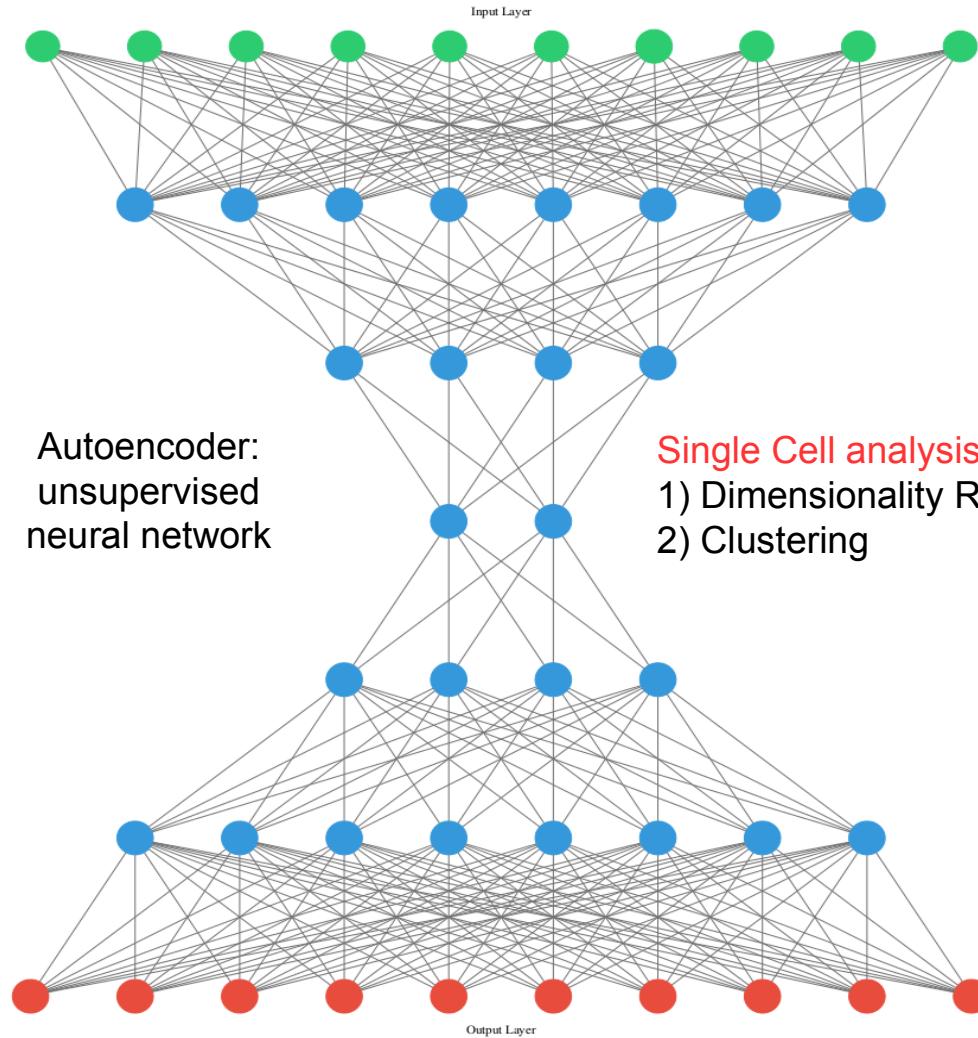


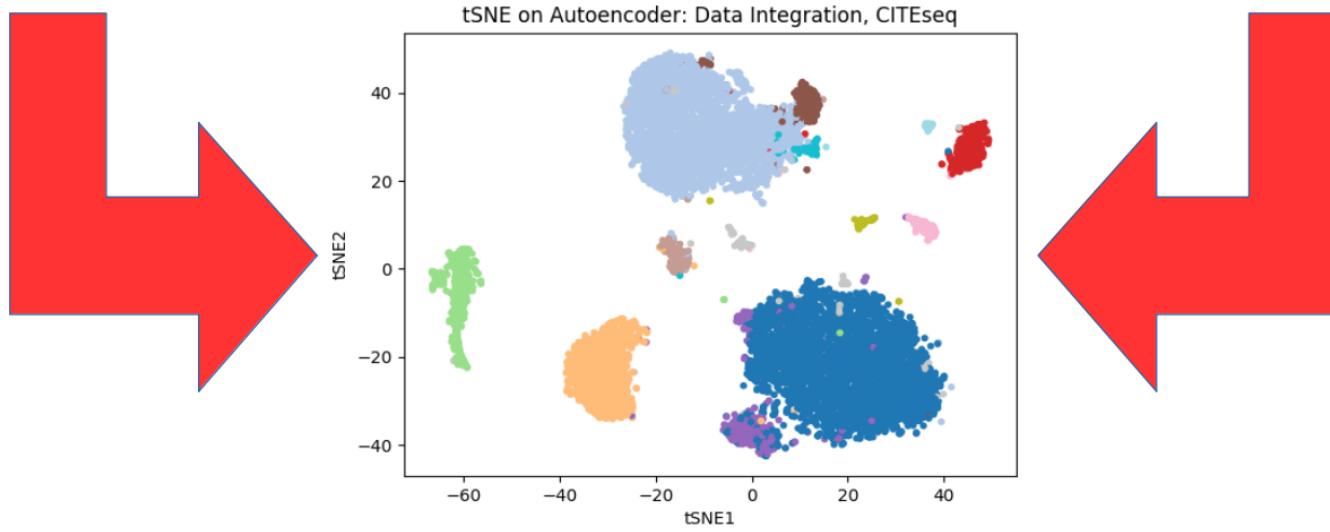
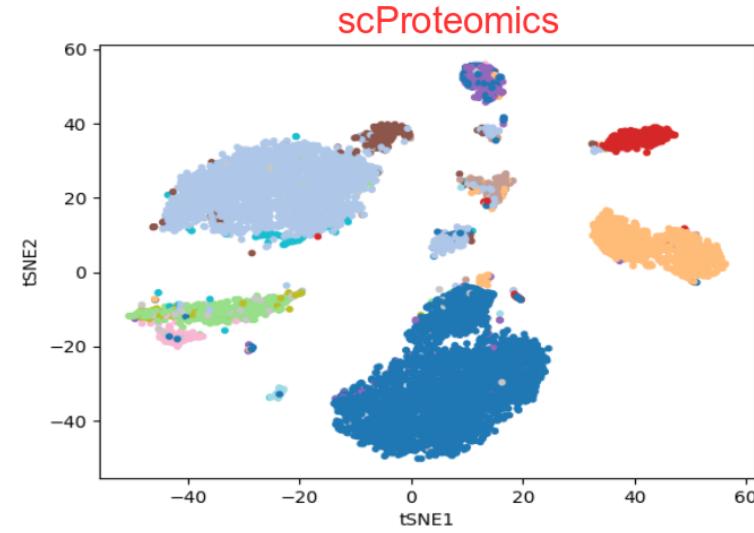
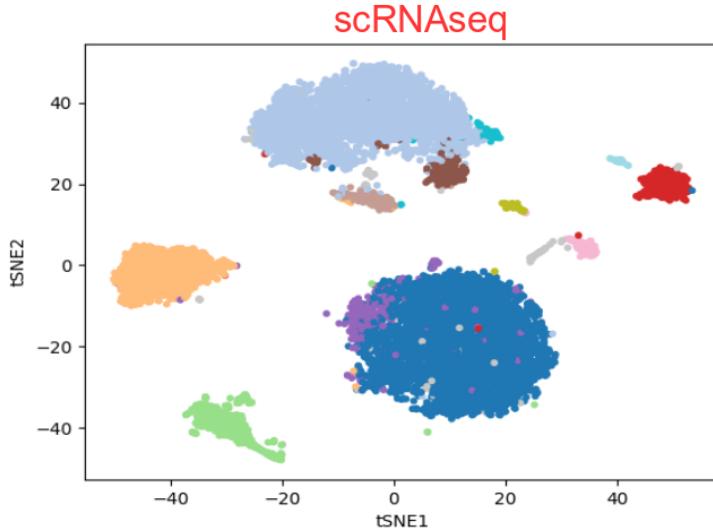
Hu et al., 2018, Frontier in Cell and Developmental Biology 6, 1-13



Clark et al., 2018, Nature Communications 9, 781

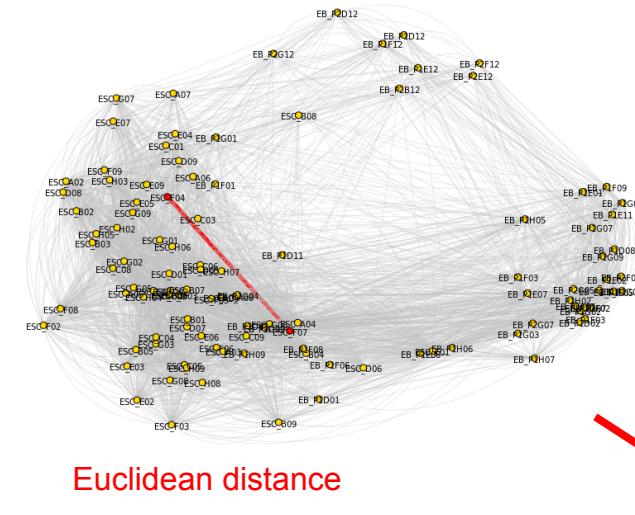






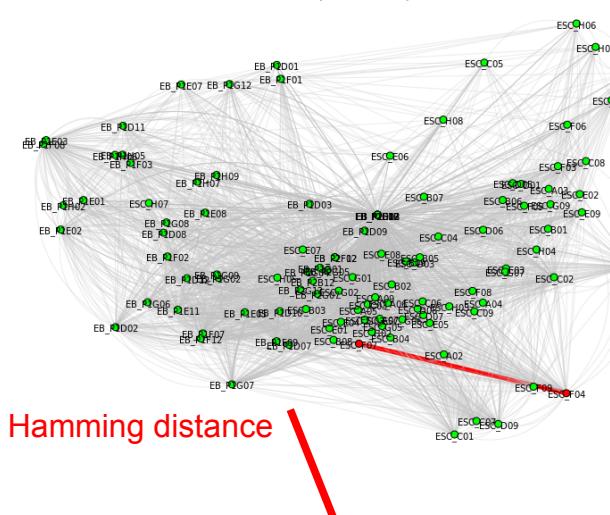
Graph Intersection Method

scRNAseq KNN Graph



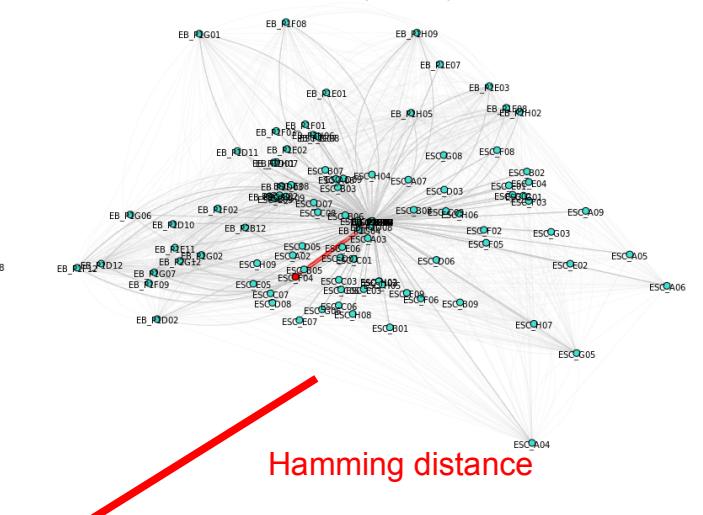
Euclidean distance

scBSseq KNN Graph



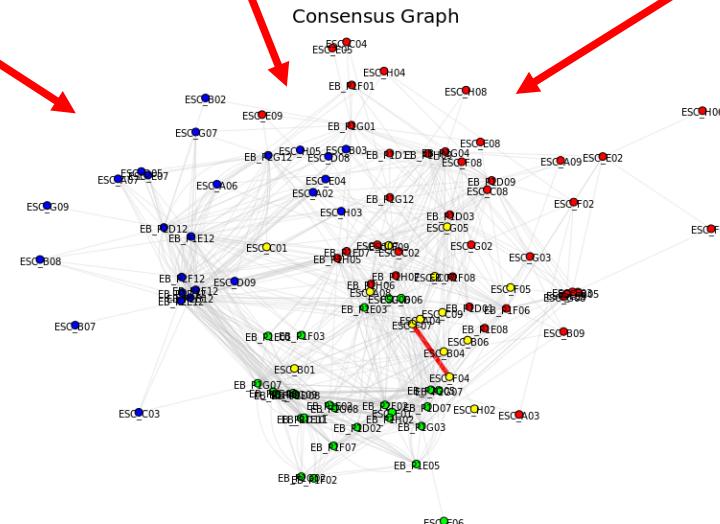
Hamming distance

scATACseq KNN Graph



Hamming distance

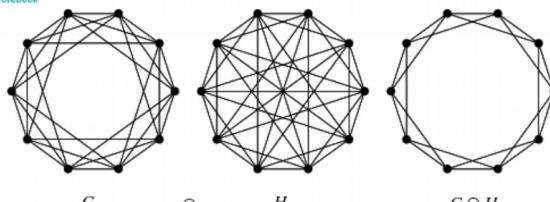
Consensus Graph



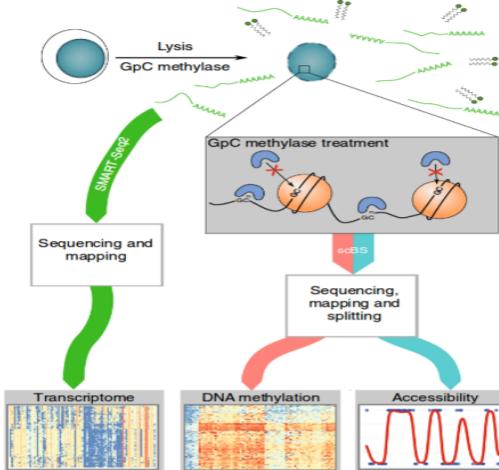
Keep edges consistently
present across the Omics

Graph Intersection

[DOWNLOAD](#)
Wolfram Notebook

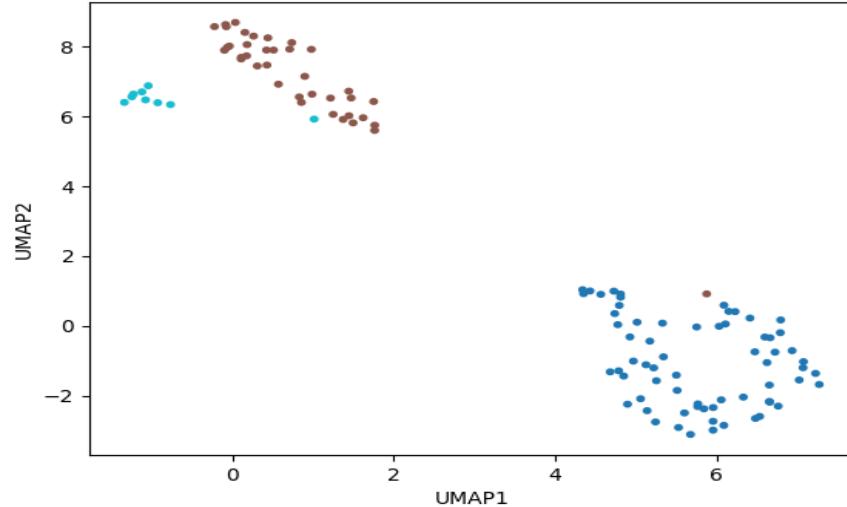


Let S be a set and $F = \{S_1, \dots, S_p\}$ a nonempty family of distinct nonempty subsets of S whose union is $\bigcup_{i=1}^p S_i = S$. The intersection graph of F is denoted $\Omega(F)$ and defined by $V(\Omega(F)) = F$, with S_i and S_j adjacent whenever $i \neq j$ and $S_i \cap S_j \neq \emptyset$. Then a graph G is an intersection graph on S if there exists a family F of subsets for which G and $\Omega(F)$ are isomorphic graphs (Harary 1994, p. 19). Graph intersections can be computed in the Wolfram Language using `GraphIntersection[g, h]`.



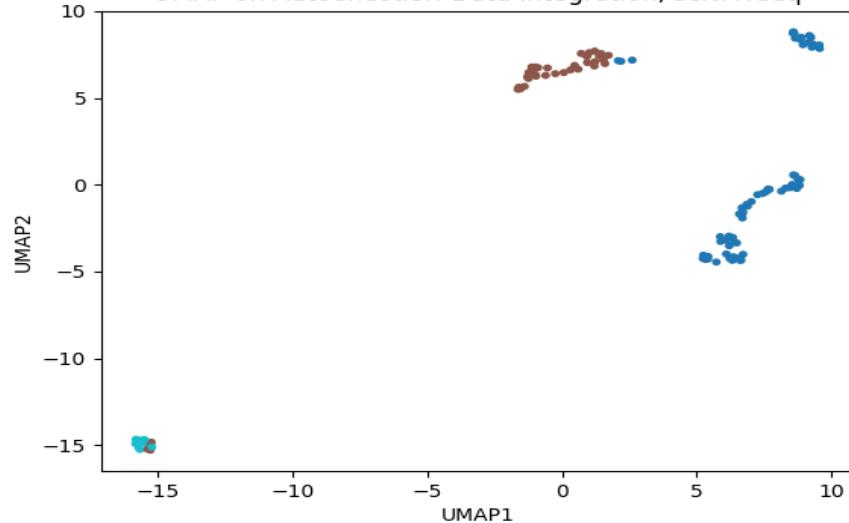
scNMTseq: Clark et al., 2018, Nature Communications 9, 781

UMAP on PCA: scNMTseq, scRNASeq



From Single
To
multi-Omics

UMAP on Autoencoder: Data Integration, scNMTseq



Take home messages of the session:

- 1) Unsupervised ANNs are useful for nonlinear data exploration
- 2) Autoencoders are data-compression-type neural networks
- 3) Single cell data is Life Science-specific Big Data type
- 4) Autoencoders are used for dimension reduction in single cell
- 5) Autoencoders are common for single cell data integration



National Bioinformatics Infrastructure Sweden (NBIS)



*Knut och Alice
Wallenbergs
Stiftelse*



**LUNDS
UNIVERSITET**