

The Krishnaswamy Laboratory
Yale Genetics and Yale SEAS present

Machine Learning for Single Cell Analysis

Online - May 20-29, 2020

When poll is active, respond at **PollEv.com/yaleml**

Text **YALEML** to **22333** once to join

What is your favorite model organism?

Machine Learning for Single Cell Analysis

Course introduction

Search Krishnaswamy Lab Help

KL Krishnaswamy L... Daniel Burkhardt

Y2 Saved items

Channel browser

People

Apps

Files

Show less

Channels

2020-workshop-byod...

2020-workshop-codin...

2020-workshop-group...

2020-workshop-main

2020-workshop-math-...

2020-workshop-misc-...

2020-workshop-tas

general

magic

meld

phate

random

scprep

workshop

#2020-workshop-main ★

20 | 1 | Add a topic

#2020-workshop-main

You created this channel on May 15th. This is the very beginning of the #2020-workshop-main channel.

Add description Add an app Add people

Friday, May 15th

Daniel Burkhardt 1:45 PM joined #2020-workshop-main along with 19 others.

Today

Pinned by you

Daniel Burkhardt 11:53 AM

Hi everyone! Welcome to the main channel for the 2020 Machine Learning for Single Cell Analysis Workshop! Please join the following channels:

1. #2020-workshop-coding-help
2. #2020-workshop-math-help
3. #2020-workshop-byod-help
4. #2020-workshop-misc-help

Message #2020-workshop-main

Aa @ 😊 🗑

<https://krishnaswamylab.org/get-help>



UNIVERSITY OF
COPENHAGEN

Washington
University in St. Louis

UNIVERSITY OF
CAMBRIDGE

W
UNIVERSITY of
WASHINGTON



FACULTY OF
MEDICINE

novo nordisk



IRELL AND
MANELLA
GRADUATE SCHOOL OF
BIOLOGICAL SCIENCES

JOHNS HOPKINS
UNIVERSITY

McMaster
University



מִצְמָן וַיִּצְמָן לְמִדְעָת
WEIZMANN INSTITUTE OF SCIENCE



COLUMBIA UNIVERSITY
IN THE CITY OF NEW YORK



Northeastern
University



WHITEHEAD
INSTITUTE



Yale

Georgia
Tech



UNIVERSITY
of VIRGINIA

UCLA

McGill



St. Jude Children's
Research Hospital
Finding cures. Saving children.



Fundación Progreso y Salud
CONSEJERÍA DE SALUD



Institut Pasteur



University of
Zurich UZH



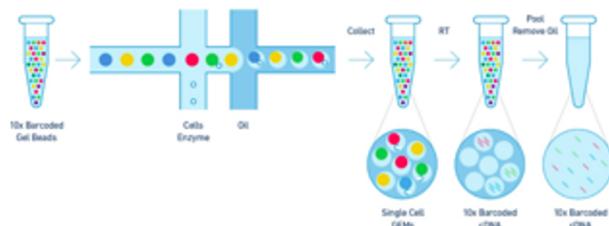
Cold
Spring
Harbor
Laboratory



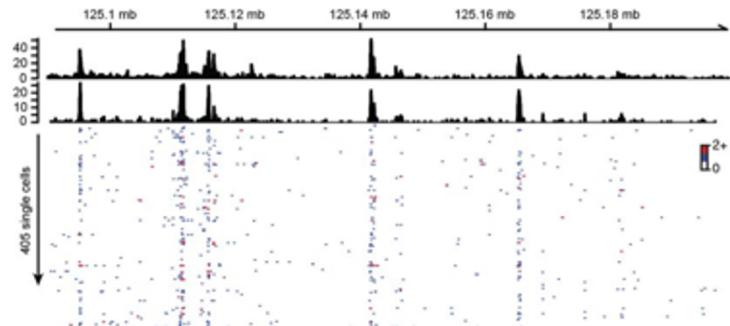
HARVARD
UNIVERSITY

IGIB
INSTITUTE OF GENOMICS & INTEGRATIVE BIOLOGY
Genomics Knowledge Partner

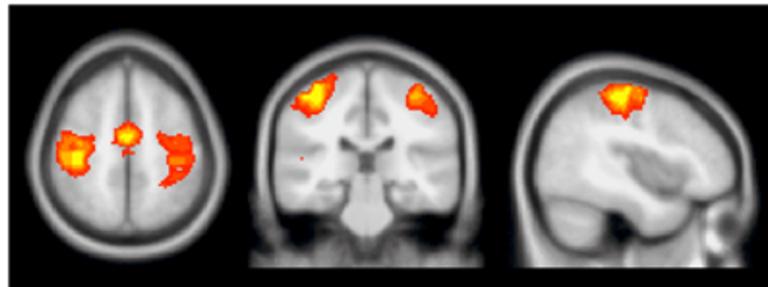
Big biomedical data



ScRNA-seq



ScATAC-seq



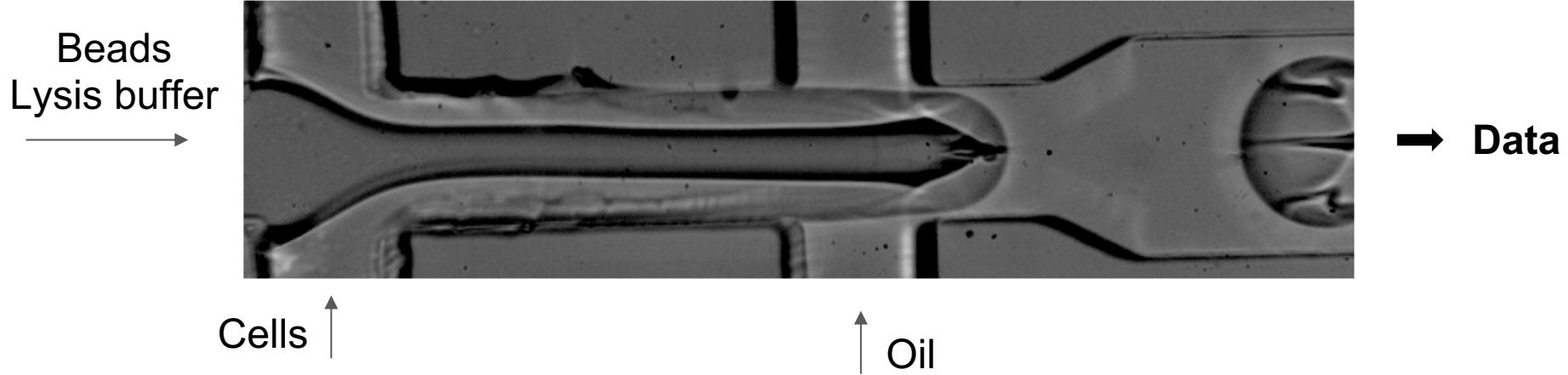
FMRI



Patient Data

Big = Any dataset with many many observations

The single cell revolution



The single cell revolution

Interesting Biological Experiments



Computation



High impact paper

LETTER

<https://doi.org/10.1038/nature08416>

RNA velocity of single cells

Gioele La Manno^{1,3}, Rulan Sudhakar⁴, Amit Zelzer¹, Esterle Braun^{1,3}, Hannah Hochegger^{1,3}, Viktor Pejchal^{1,3}, Katja Lischefski^{1,2}, Maria E. Kastell¹, Peter Lönnerberg^{2,3}, Alessandro Pusztai¹, Jean Fan¹, Lars E. Börre^{1,3}, Zehan Liu¹, Daniel C. Westover¹, Michael A. Hodge¹, Michael S. Hengenberger¹, Gonzalo Caočić¹, Barbara J. Strome⁵, Patrick Cramer^{1,3}, Igor Muravskiy¹, Sten Linnemann^{2,3} & Peter V. Kharchenko^{1,3*}

RNA abundance is a powerful indicator of the state of individual cells. Single-cell RNA sequencing can reveal RNA abundance with high quantitative accuracy, without the need for a reference transcriptome. However, a static snapshot of a transcriptome is not always enough for the analysis of time-resolved phenomena such as gene expression dynamics. The concept of RNA velocity—the time derivative of the gene expression state—can be used to predict the future state of a cell. RNA velocity has been used to predict the fate of cells in the context of cancer^{1–3} and to predict the fate of cells in the context of stem cell differentiation⁴. RNA velocity is also useful for single-cell RNA sequencing protocols. RNA velocity is a high-dimensional vector that predicts the future state of a cell based on its current transcriptomic profile. We used RNA velocity to predict the future state of cells in the context of the neural crest lineage, demonstrate its use on multiple published datasets, and show that RNA velocity can be used to predict the fate of the developing mouse hippocampus, and examine the kinetics of gene expression dynamics in the context of cell differentiation. RNA velocity can greatly aid the analysis of developmental lineages and cellular heterogeneity, particularly in humans.

During development, transcription occurs on a timescale of hours to days, which is comparable to the typical half-life of mRNA. The timescale of RNA velocity is therefore similar to the timescale over which gene expression dynamics can be observed. RNA velocity can be exploited to estimate the rates of gene splicing and degradation, and to predict the future state of a cell based on its current transcriptomic profile. We measured that similar signals may be detectable in single-cell RNA sequencing data, and we used RNA velocity to predict the direction of change of the entire transcriptome during dynamic processes.

All common single-cell RNA-seq protocols rely on oligo-dT primers to select poly(A) mRNA. We analyzed RNA-seq data from 11 studies using single-cell RNA-seq datasets based on the SMART-seq², STRT³, and Dropbead⁴ platforms. In total, we analyzed 1,100 samples. 15–25% of reads contained unspliced intrinsic sequences (Fig. 1a), in contrast to the expected 1–2% of unspliced RNA in a standard 1–20% RNA-seq. Most such reads originated from secondary priming, which is commonly occurring in RNA-seq⁵. Using the Genomics Chromatin library, we also found abundant discordances between the expected and observed RNA-seq profiles after PCR amplification by priming on the first-stranded cDNA. The substantial number of unspliced reads and discordant RNA-seq profiles suggests that these molecules represent unspliced precursor RNA molecules that have been converted to matured spliced mRNA⁶ followed by RNA sequencing using oligo-dT primed RT⁷ (Extended Data Fig. 1a). The balance of unspliced and spliced mRNA abundance, and the future state of the cell, can be predicted based on the equilibrium rate γ , where the opposite is true during regeneration (Fig. 1b). The balance of unspliced and spliced mRNA abundance is determined by the equilibrium rate γ and the transcription rate α . The equilibrium rate γ combines degradation and splicing rates. Capturing gene specific regulatory properties, the rate of gene expression dynamics is often much slower than the timescale of gene expression. Using a recently published compendium of mouse tissues⁸, we found that the equilibrium rate γ was constant across all tissues. The equilibrium rate was consistent with a single fixed slope (Extended Data Fig. 1c). The equilibrium rate γ was also consistent across different mouse tissues (Extended Data Fig. 1d), suggesting tissue specific alternative splicing. The equilibrium rate γ was also consistent across different mouse tissues (Extended Data Fig. 1e), suggesting tissue specific alternative splicing. The equilibrium rate γ was also consistent across different mouse tissues (Extended Data Fig. 1f), suggesting tissue specific alternative splicing.

To validate a dynamic process, an increase in the transcription rate α results in a rapid increase in mRNA abundance (Fig. 1g). Conversely, a drop in the transcription rate α results in a rapid decrease in mRNA abundance (Fig. 1h). We used the equilibrium rate γ and the transcription rate α to extrapolate the future mRNA abundance of the future. We examined a time course of gene expression dynamics in the context of cell differentiation (Fig. 1i). We found that the equilibrium rate γ was consistent across all genes and all tissues (Fig. 1j). Many circadian associated genes showed the expected excess of unspliced mRNA relative to the slope during upregulation, and a deficit of unspliced mRNA relative to the slope during downregulation. The proposed differential equations for each gene allowed us to extrapolate the expected direction of progression of the circadian cycle (Fig. 1k). The equilibrium rate γ was consistent across all genes and all tissues in single-cell measurements, we analysed recently published single-cell mRNA abundance data from the developing mouse hippocampus (Fig. 1l). During development, a substantial proportion of olfactory ensheathing cells, which are neuroectoderm cells of the adrenal medulla, undergo a dynamic process of differentiation. We found that the direction of differentiation can be validated by lineage tracing. Phase portraits of many genes showed the expected deviation

*Correspondence: Michael Hodge (michael.hodge@utoronto.ca); Peter V. Kharchenko (kharchenko@mit.edu). ¹Department of Molecular Biology, University of Toronto, Toronto, Ontario, Canada. ²Department of Molecular Genetics, University of Toronto, Toronto, Ontario, Canada. ³Department of Biochemistry, University of Toronto, Toronto, Ontario, Canada. ⁴Department of Cell Biology, University of Toronto, Toronto, Ontario, Canada. ⁵Department of Cell Biology, Harvard Medical School, Boston, MA, USA. ⁶Department of Cell Biology, Harvard Medical School, Boston, MA, USA. ⁷Department of Cell Biology, Harvard Medical School, Boston, MA, USA. ⁸Department of Cell Biology, Harvard Medical School, Boston, MA, USA.

474 | NATURE | VOL 540 | 21 AUGUST 2016

- Machine learning
- Linear algebra
- Probability theory
- Statistical analysis
- Algorithm design

It's all Greek to me...

Definition 1. The t -step potential distance is defined as $\mathfrak{V}^t(x, y) \triangleq \|U_x^t - U_y^t\|_2$, $x, y \in \mathcal{X}$.

The following proposition shows a relation between the two metrics by expressing the potential distance in embedded diffusion map coordinates¹ for fixed-bandwidth Gaussian-based diffusion (i.e., generated by P_ε from Eq. 2):

Proposition 1. Given a diffusion process defined by a fixed-bandwidth Gaussian kernel, the potential distance from Def 1 can be written as $\mathfrak{V}^t(x, y) = \left(\sum_{z \in \mathcal{X}} \log^2 \left(\frac{1 + \langle \Phi^{t/2}(x), \Phi^{t/2}(z) \rangle}{1 + \langle \Phi^{t/2}(y), \Phi^{t/2}(z) \rangle} \right) \right)^{1/2}$

Proof. According to the spectral theorem, the entries of P_ε^t can be written as

$$[P_\varepsilon^t]_{(x,y)} = \psi_0(y) + \sum_{i=1}^{n-1} \lambda_i^t \phi_i(x) \psi_i(y)$$

since powers of the operator P_ε only affect the eigenvalues, which are taken to the same power, and since the trivial eigenvalue λ_0 is one and the corresponding right eigenvector ϕ_0 only consists of ones. Furthermore, it can be verified that the left and right eigenvectors of P_ε are related by $\psi_i(y) = \phi_i(y) \psi_0(y)$, thus, combined with Eqs. 4 and 6, we get

$$p_{\varepsilon,x}^t(y) = \psi_0(y) \left(1 + \sum_{i=1}^{n-1} \lambda_i^t \phi_i(x) \phi_i(y) \right) = \psi_0(y) (1 + \langle \Phi_\varepsilon^{t/2}(x), \Phi_\varepsilon^{t/2}(y) \rangle) .$$

By applying the logarithm to both ends of this equation we express the entries of the potential representation $U_{\varepsilon,x}^t$ as

$$U_{\varepsilon,x}^t(y) = -\log(1 + \langle \Phi_\varepsilon^{t/2}(x), \Phi_\varepsilon^{t/2}(y) \rangle) - \log(\psi_0(y)) ,$$

and thus for any $j = 1, \dots, N$,

$$\begin{aligned} (U_{\varepsilon,x}^t(x_j) - U_{\varepsilon,y}^t(x_j))^2 &= [\log(1 + \langle \Phi_\varepsilon^{t/2}(x), \Phi_\varepsilon^{t/2}(x_j) \rangle)]^2 \\ &\quad - [\log(1 + \langle \Phi_\varepsilon^{t/2}(y), \Phi_\varepsilon^{t/2}(x_j) \rangle)]^2 \\ &= \log^2 \left(\frac{1 + \langle \Phi_\varepsilon^{t/2}(x), \Phi_\varepsilon^{t/2}(x_j) \rangle}{1 + \langle \Phi_\varepsilon^{t/2}(y), \Phi_\varepsilon^{t/2}(x_j) \rangle} \right) , \end{aligned}$$

which yields the result in the proposition. \square

What reading single cell methods can feel like



What is machine learning?

What is machine learning?

Machine learning is the process of identifying patterns in data.

Two kinds of machine learning

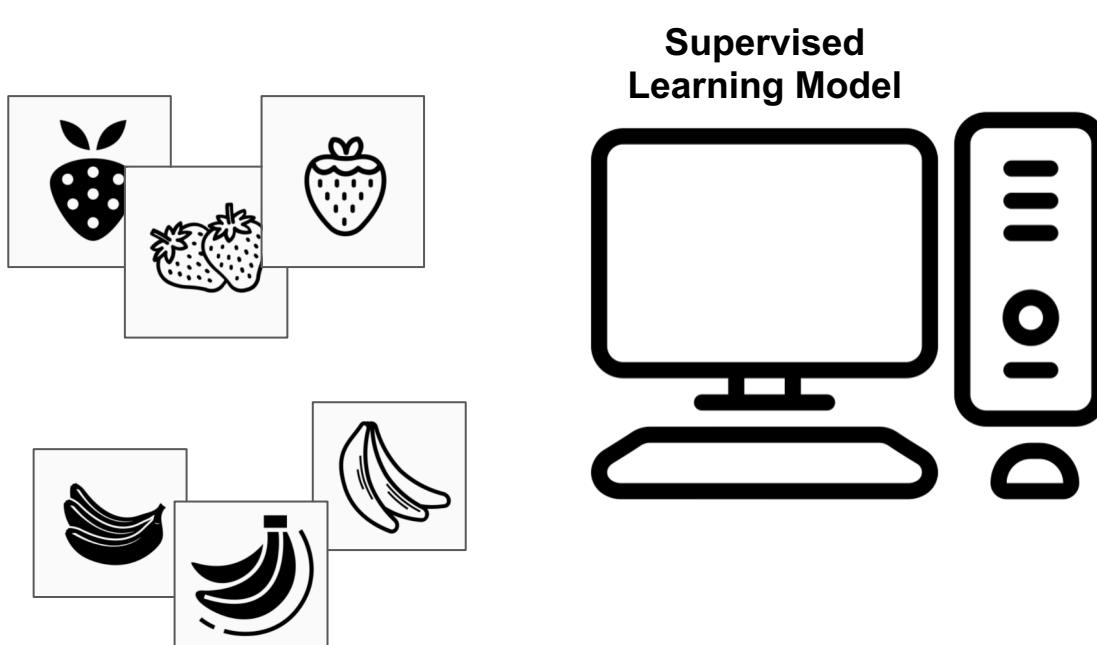
Supervised learning

- Have a bunch of labelled data, want to label new data

Two kinds of machine learning

Supervised learning

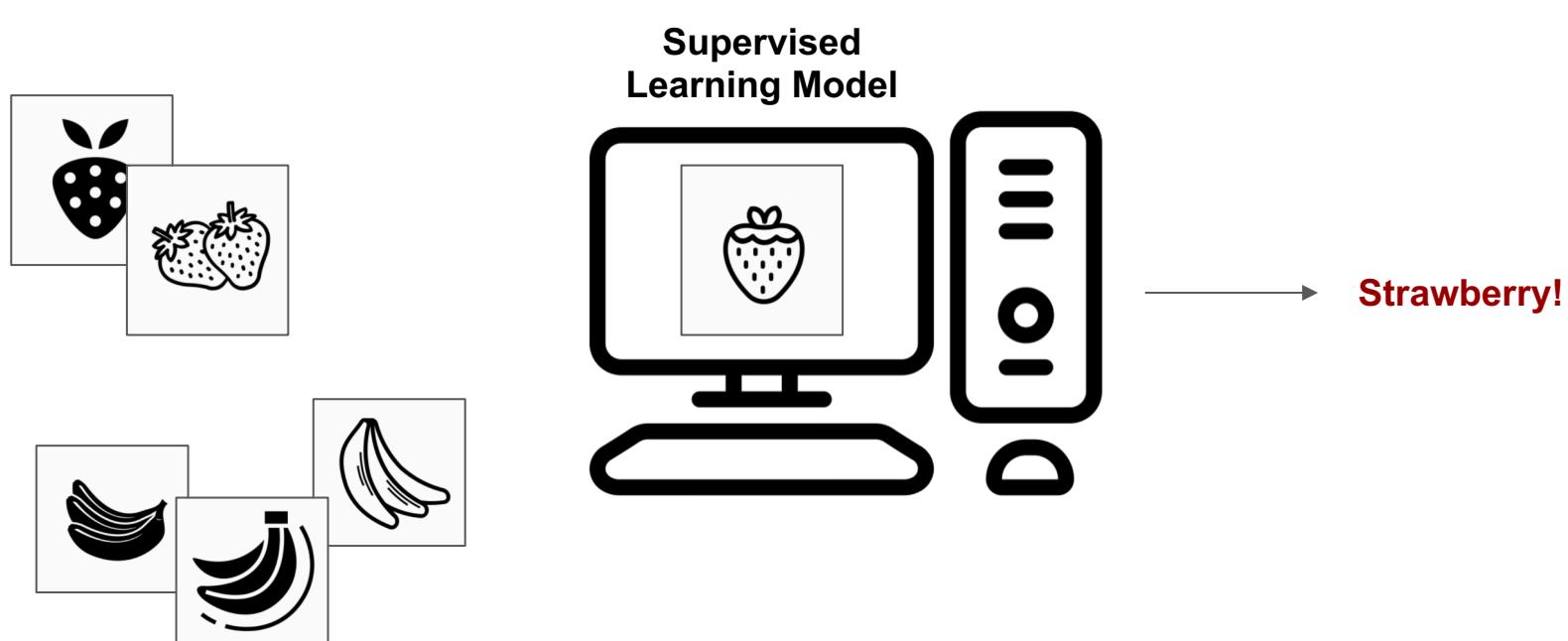
- Have a bunch of labelled data, want to label new data



Two kinds of machine learning

Supervised learning

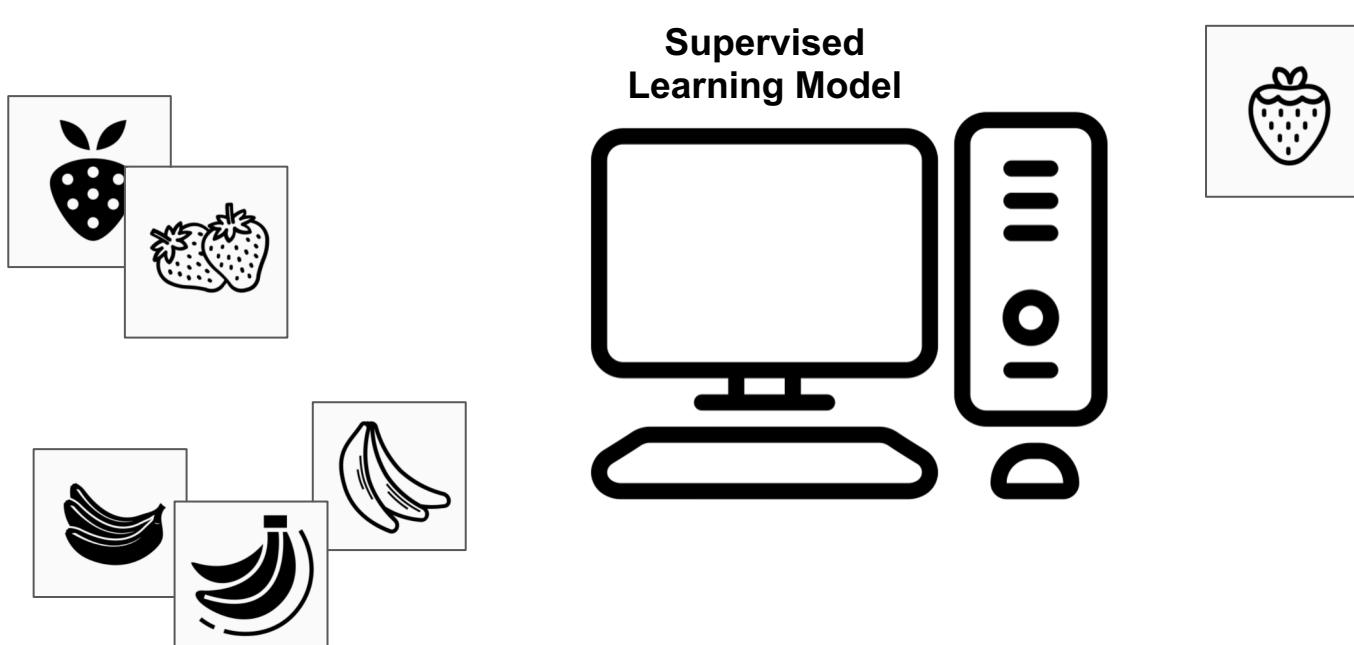
- Have a bunch of labelled data, want to label new data



Two kinds of machine learning

Supervised learning

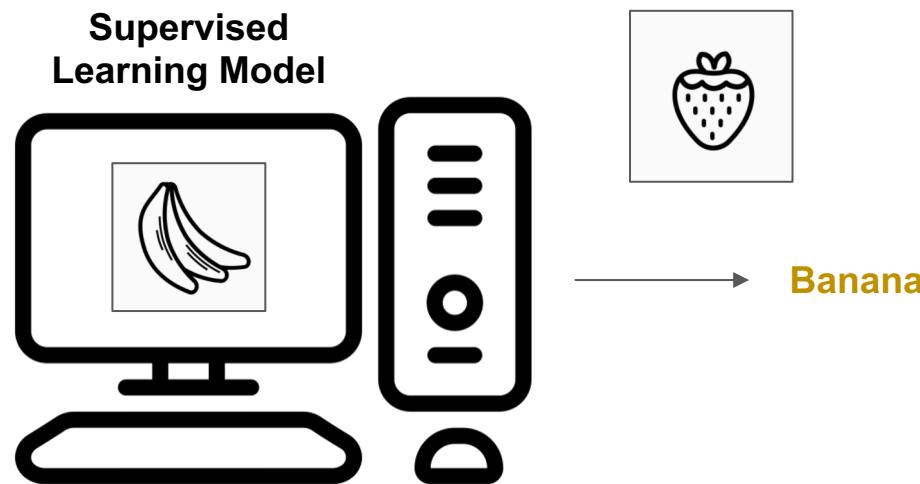
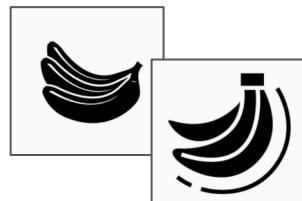
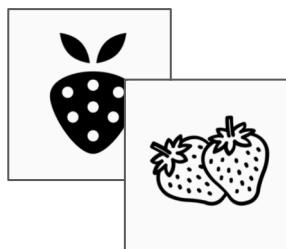
- Have a bunch of labelled data, want to label new data



Two kinds of machine learning

Supervised learning

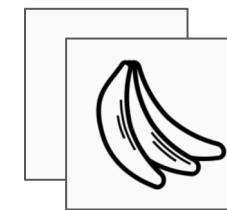
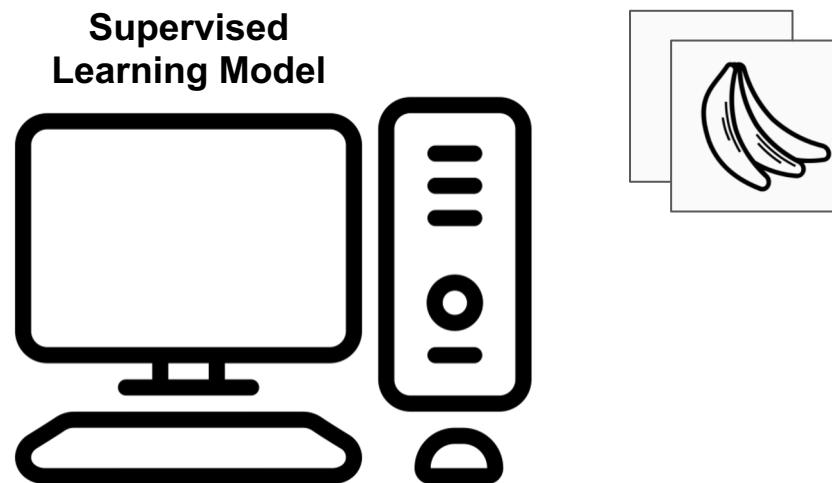
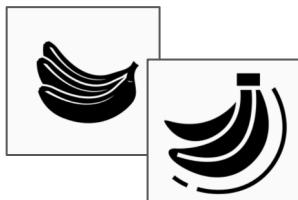
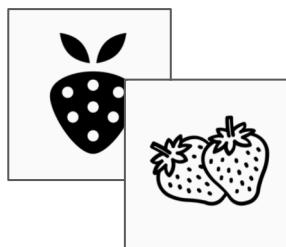
- Have a bunch of labelled data, want to label new data



Two kinds of machine learning

Supervised learning

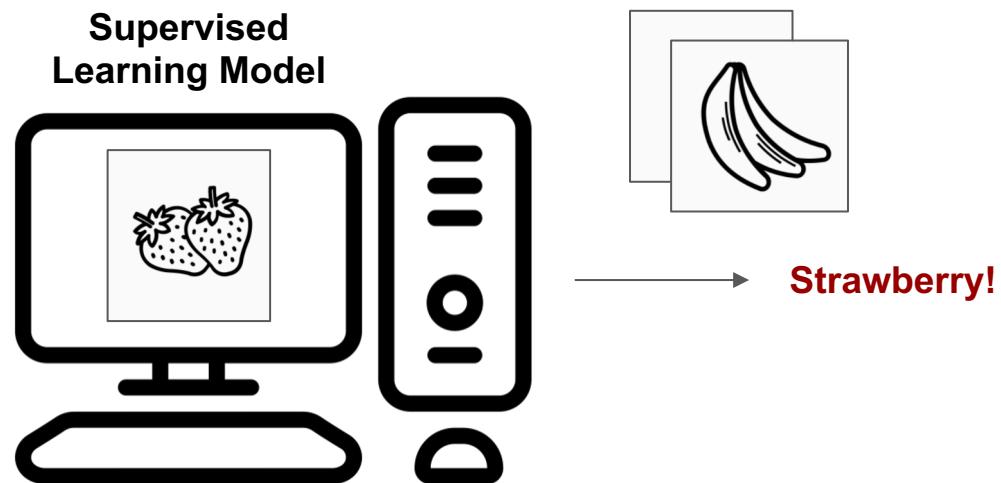
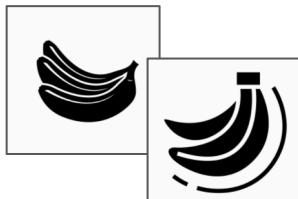
- Have a bunch of labelled data, want to label new data



Two kinds of machine learning

Supervised learning

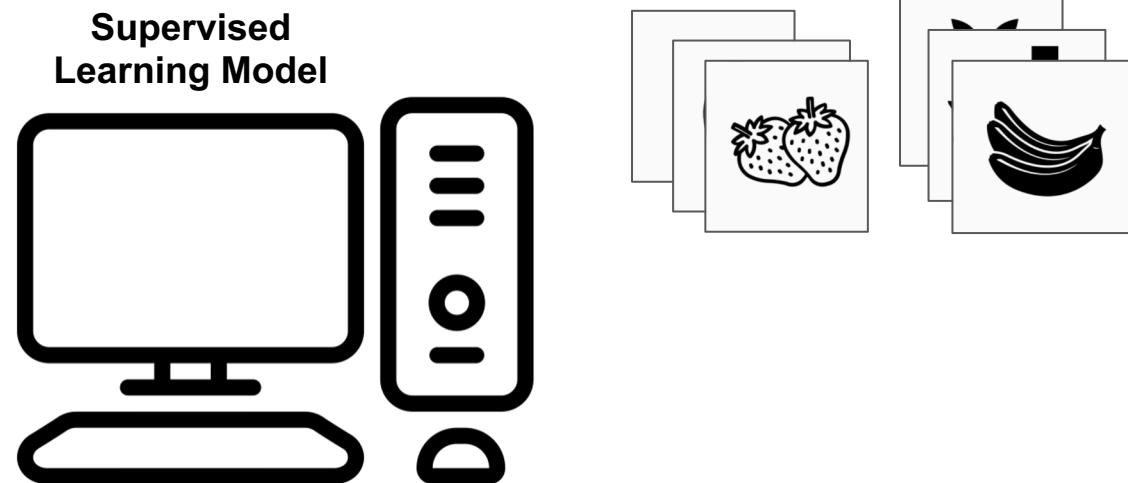
- Have a bunch of labelled data, want to label new data



Two kinds of machine learning

Supervised learning

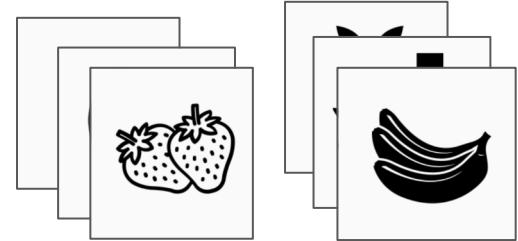
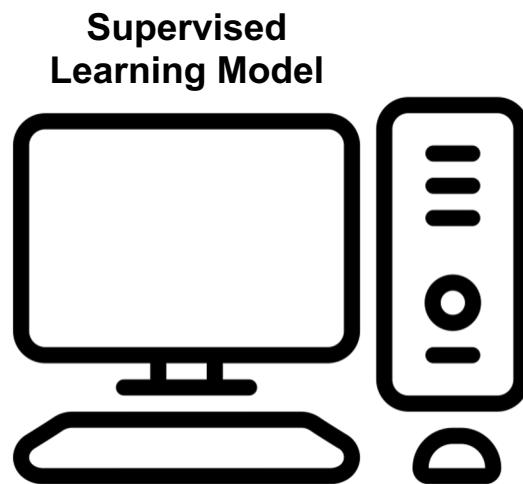
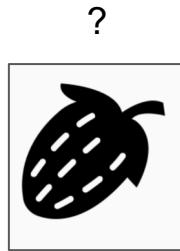
- Have a bunch of labelled data, want to label new data



Two kinds of machine learning

Supervised learning

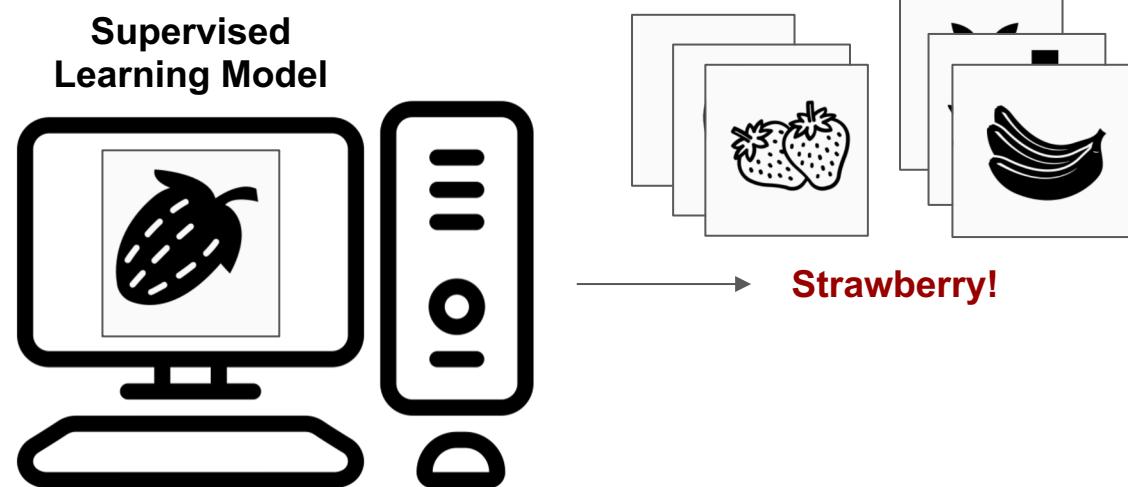
- Have a bunch of labelled data, want to label new data



Two kinds of machine learning

Supervised learning

- Have a bunch of labelled data, want to label new data



Two kinds of machine learning

Supervised learning

- Have a bunch of labelled data, want to label new data

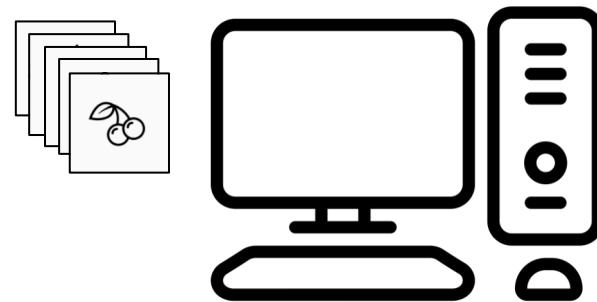
Supervised Learning Model



Unsupervised learning

- Have a bunch of unlabeled data, want to organize it

Unsupervised Learning Model

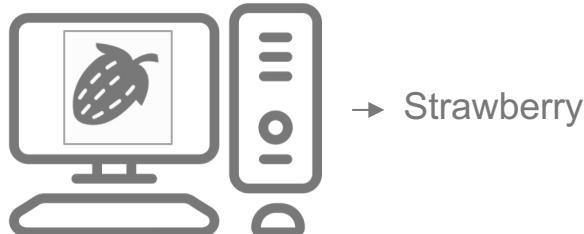


Two kinds of machine learning

Supervised learning

- Have a bunch of labelled data, want to label new data

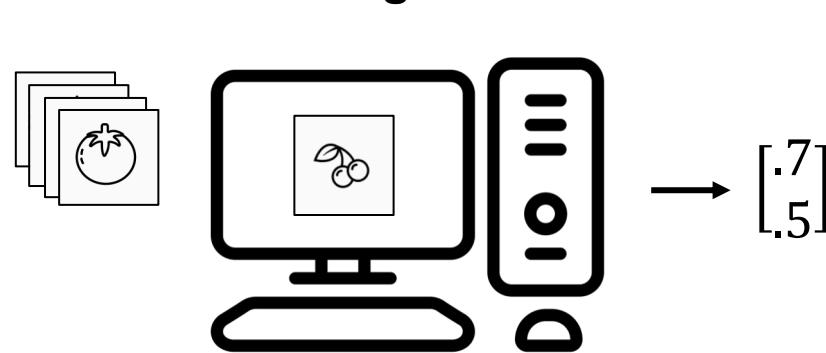
Supervised Learning Model



Unsupervised learning

- Have a bunch of unlabeled data, want to organize it

Unsupervised Learning Model



Two kinds of machine learning

Supervised learning

- Have a bunch of labelled data, want to label new data

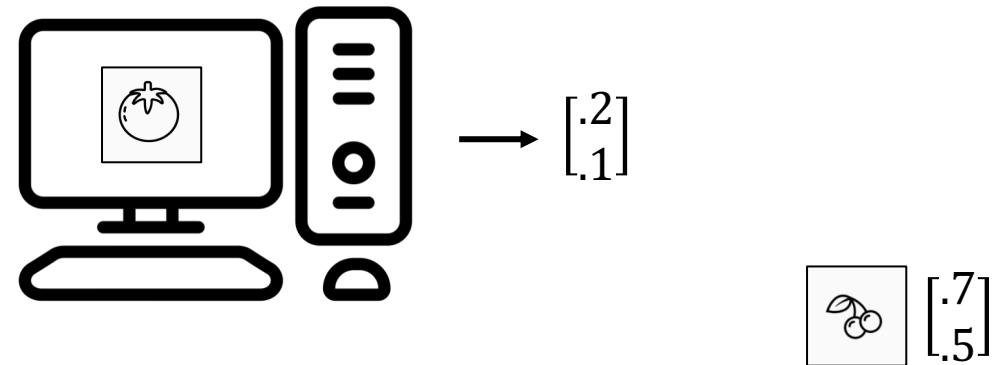
Supervised Learning Model



Unsupervised learning

- Have a bunch of unlabeled data, want to organize it

Unsupervised Learning Model



Two kinds of machine learning

Supervised learning

- Have a bunch of labelled data, want to label new data

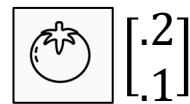
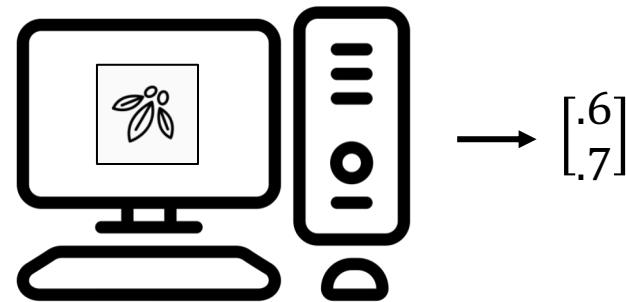
Supervised Learning Model



Unsupervised learning

- Have a bunch of unlabeled data, want to organize it

Unsupervised Learning Model

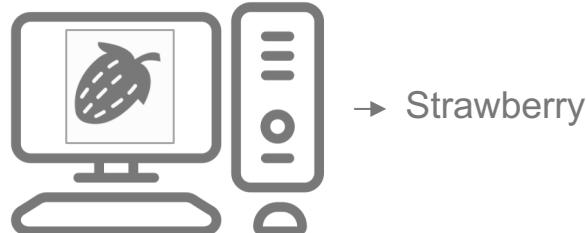


Two kinds of machine learning

Supervised learning

- Have a bunch of labelled data, want to label new data

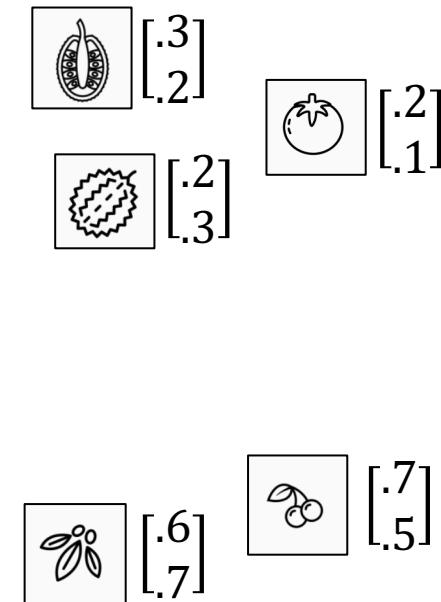
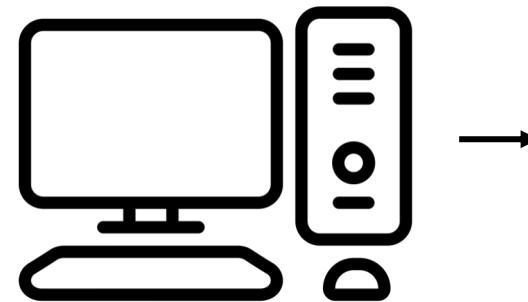
Supervised Learning Model



Unsupervised learning

- Have a bunch of unlabeled data, want to organize it

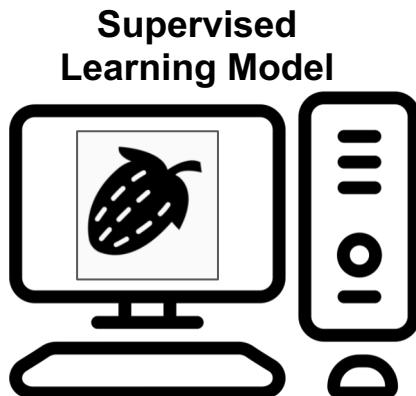
Unsupervised Learning Model



Two kinds of machine learning

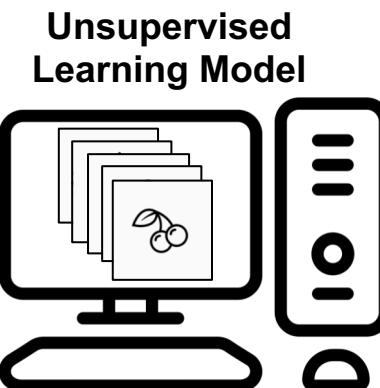
Supervised learning

- Have a bunch of labelled data, want to label new data
- Learn a function $f(X) \rightarrow Y$ where all values of Y are known for some samples of X



Unsupervised learning

- Have a bunch of unlabeled data, want to organize it
- Learn an embedding $f(X) \rightarrow Y, X \in \mathbb{R}^n, Y \in \mathbb{R}^m, n \gg m$
- Lower dimensional, easier to interpret (e.g. as clusters)



Is linear regression an example of supervised or unsupervised machine learning?

Supervised
machine
learning

Unsupervised
machine
learning

Is clustering an example of supervised or unsupervised machine learning?



Supervised
machine learning **A**

Unsupervised
machine learning **B**

Course structure

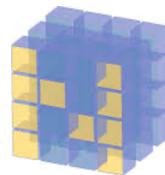
The screenshot shows a web browser window titled "Workshop — Krishnaswamy Lab". The URL is <https://www.krishnaswamylab.org/workshop>. The page content is organized by day:

- Day 1 – Wednesday, May 20th**
 - Lecture [View in Google Slides](#) Introduction
 - [View in Google Slides](#) Preprocessing scRNAseq Data
 - Exercise [Run in Google Colab](#) 1.0. Preprocessing Embryoid Body Data
- Day 2 – Thursday, May 21st**
 - Lecture [View on Google Drive](#) Challenges in Single Cell Analysis
 - [View on Google Drive](#) Thinking about High-Dimensional Data
 - Exercise [Run in Google Colab](#) 2.0. Plotting UCI Wine Data
 - [Run in Google Colab](#) 2.1. Learning Graphs from Data
- Day 3 – Friday, May 22nd**
 - Lecture [View in Google Slides](#) What is Visualization?
 - [View in Google Slides](#) Creating better features with PCA
 - [View in Google Slides](#) Nonlinear dimensionality Reduction
 - Exercise [Run in Google Colab](#) 3.0. Visualizing UCI Wine Data
 - [Run in Google Colab](#) 3.1. PCA on Retinal Bipolar Data
 - [Run in Google Colab](#) 3.2. Visualizing Retinal Bipolar Data
 - [Run in Google Colab](#) 3.3. Visualizing Embryoid Body Data
- Day 4 – Wednesday, May 27th**
 - Lecture [View in Google Slides](#) Clustering and Differential Expression
 - [View on Google Drive](#) Reducing Noise in scRNAseq Measurements
 - [Run in Google Colab](#)

<https://www.krishnaswamylab.org/workshop>

Why Python?

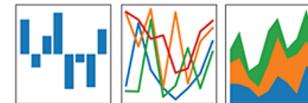




NumPy

pandas

$$y_{it} = \beta' x_{it} + \mu_i + \epsilon_{it}$$



Why Python?



Tensorflow



Pytorch



scanpy



The screenshot shows the Google Colab interface. At the top, there's a browser-like header with tabs, a search bar, and various toolbar icons. Below the header, the main content area has a title "Welcome To Colab" with a "CO" logo. A navigation menu includes File, Edit, View, Insert, Runtime, Tools, Help, and a "Share" button. On the left, a sidebar titled "Table of contents" lists sections like Getting started, Data science, Machine learning, More Resources, and Machine Learning Examples. The main content area features a large "CO" logo and the heading "What is Colab?". It explains that Colab allows writing and executing Python in a browser with zero configuration, free access to GPUs, and easy sharing. It also mentions that Colab can make work easier for students, data scientists, and AI researchers. Below this, a section titled "Getting started" is expanded, explaining that the document is an interactive Colab notebook. It shows a code cell with the Python script `[] seconds_in_a_day = 24 * 60 * 60` and its output `72000`. A note says that to execute the code, one can select it and press Command/Ctrl+Enter or use the keyboard shortcut. It also mentions that variables defined in one cell can be used in others.

CO What is Colab?

Colab allows you to write and execute Python in your browser, with

- Zero configuration required
- Free access to GPUs
- Easy sharing

Whether you're a **student**, a **data scientist** or an **AI researcher**, Colab can make your work easier. Watch [Introduction to Colab](#) to learn more, or just get started below!

▼ Getting started

The document you are reading is not a static web page, but an interactive environment called a **Colab notebook** that lets you write and execute code.

For example, here is a **code cell** with a short Python script that computes a value, stores it in a variable, and prints the result:

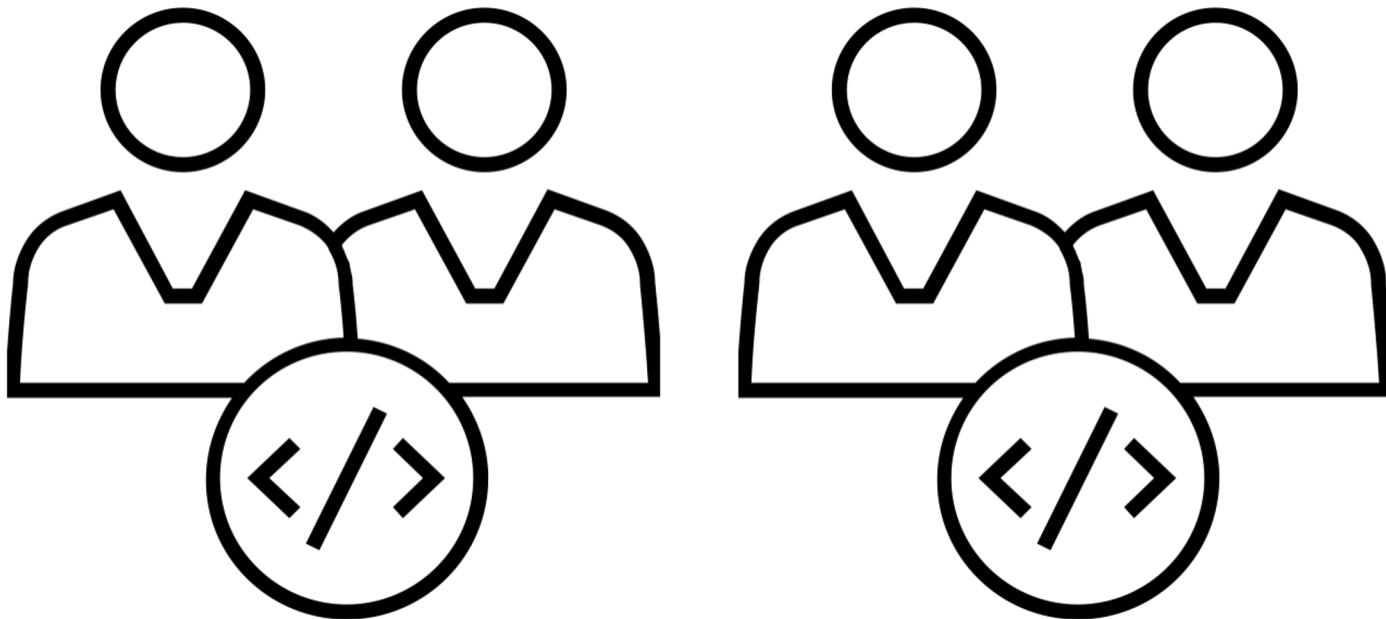
```
[ ] seconds_in_a_day = 24 * 60 * 60
seconds_in_a_day
```

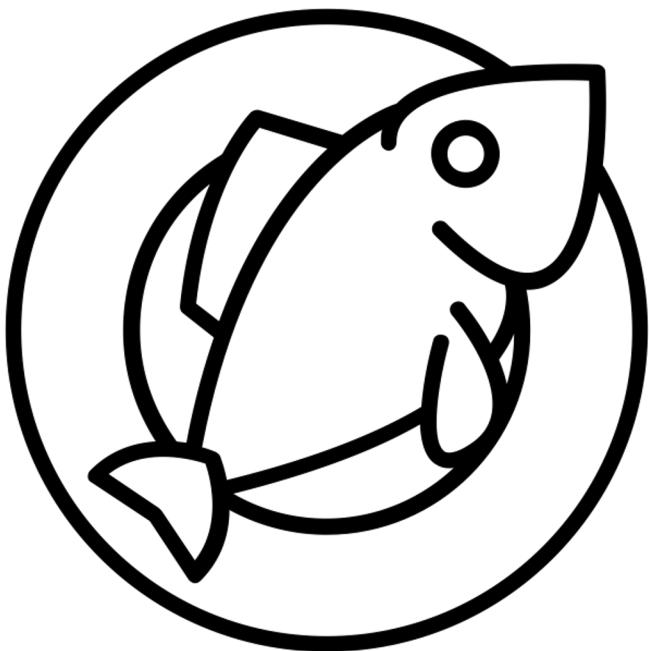
72000

To execute the code in the above cell, select it with a click and then either press the play button to the left of the code, or use the keyboard shortcut "Command/Ctrl+Enter". To edit the code, just click the cell and start editing.

Variables that you define in one cell can later be used in other cells:

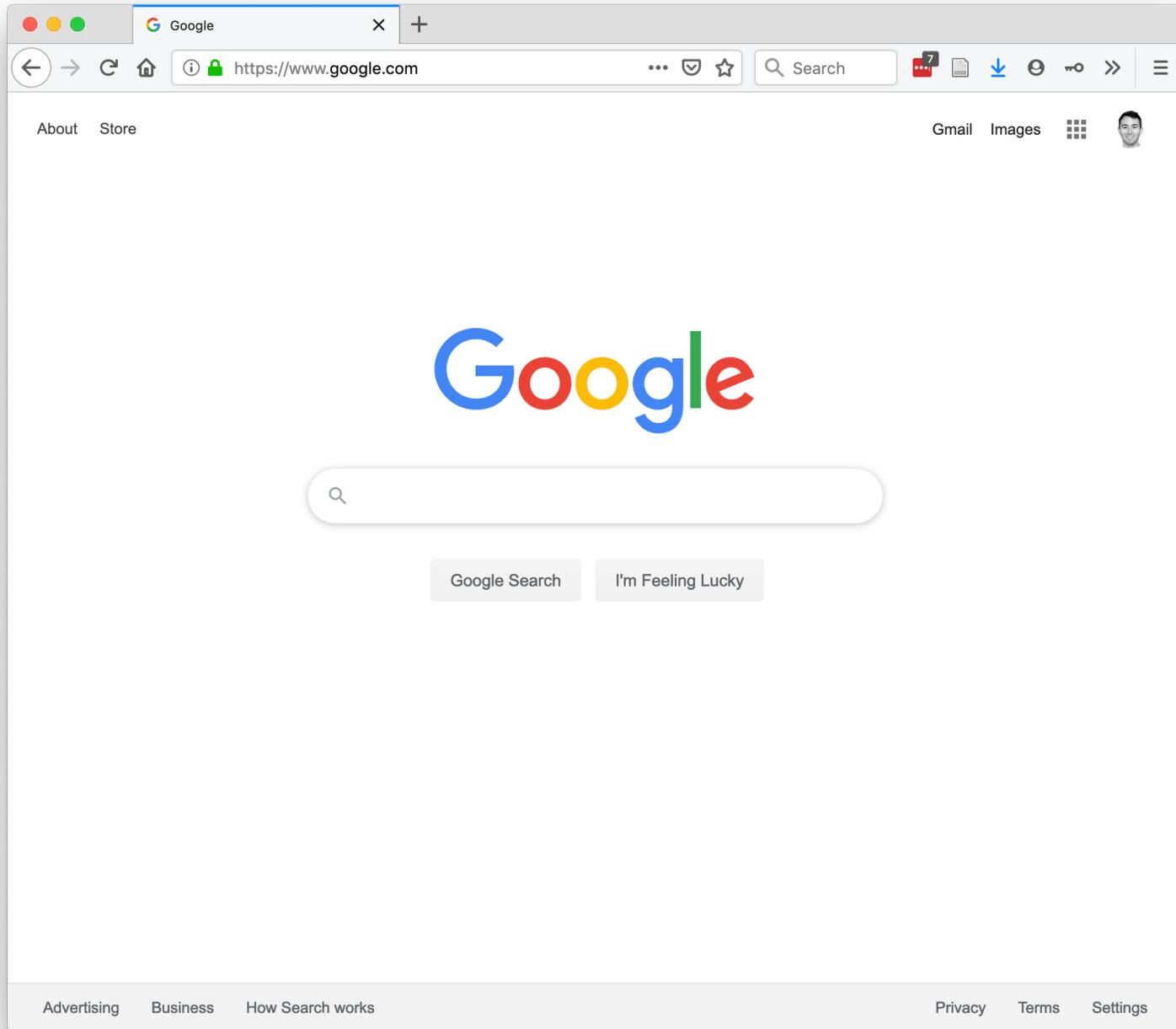
Team programming





vs.





Reference — scprep 1.0.1 documentation

scprep.io.load_10X(*data_dir*, *sparse=True*, *gene_labels='symbol'*, *allow_duplicates=None*) [source]

Basic IO for 10X data produced from the 10X Cellranger pipeline.

A default run of the *cellranger count* command will generate gene-barcode matrices for secondary analysis. For both “raw” and “filtered” output, directories are created containing three files: ‘matrix.mtx’, ‘barcodes.tsv’, ‘genes.tsv’. Running *scprep.io.load_10X(data_dir)* will return a Pandas DataFrame with genes as columns and cells as rows.

Parameters:

- ***data_dir* (string)** – path to input data directory expects ‘matrix.mtx’, ‘genes.tsv’, ‘barcodes.tsv’ to be present and will raise an error otherwise
- ***sparse* (boolean)** – If True, a sparse Pandas DataFrame is returned.
- ***gene_labels* (string, {‘id’, ‘symbol’, ‘both’} optional, default: ‘symbol’)** – Whether the columns of the dataframe should contain gene ids or gene symbols. If ‘both’, returns symbols followed by ids in parentheses.
- ***allow_duplicates* (bool, optional (default: None))** – Whether or not to allow duplicate gene names. If None, duplicates are allowed for dense input but not for sparse input.

Returns:

Return type:

scprep.io.load_10X_HDF5(*filename*, *genome=None*, *sparse=True*, *gene_labels='symbol'*, *allow_duplicates=None*, *backend=None*) [source]

Basic IO for HDF5 10X data produced from the 10X Cellranger pipeline.

Installation Examples Reference Data Input/Output HDF5 Download Filtering Normalization Transformation Measurements Statistics Plotting Dimensionality Reduction Row/Column Selection Utilities External Tools

The POWERFUL PYTHON PLAYBOOK for intermediate+ Python. Download free here

Sponsored · Ads served ethically

Read the Docs v: stable ▾

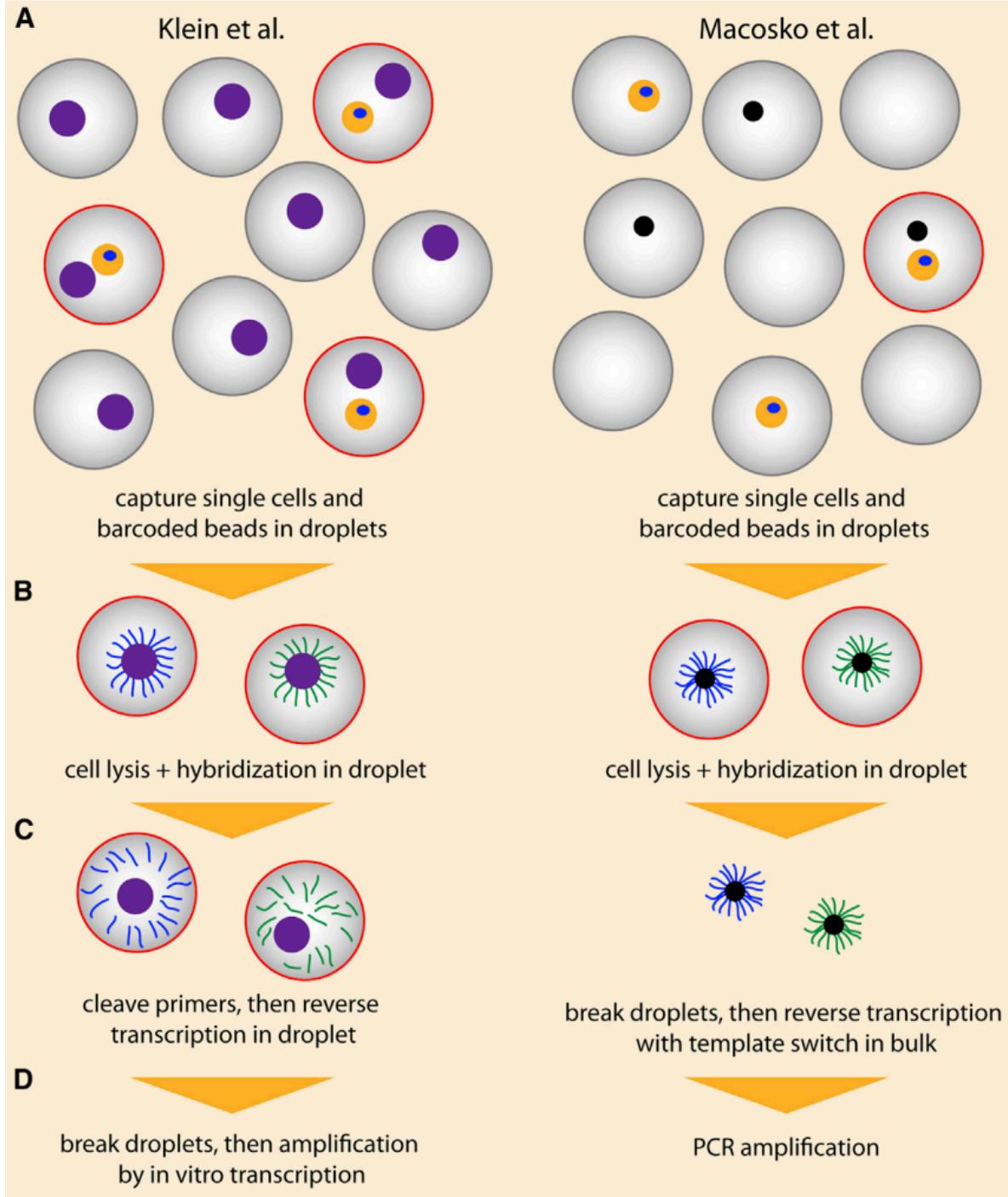
Bring-your-own-data workshop



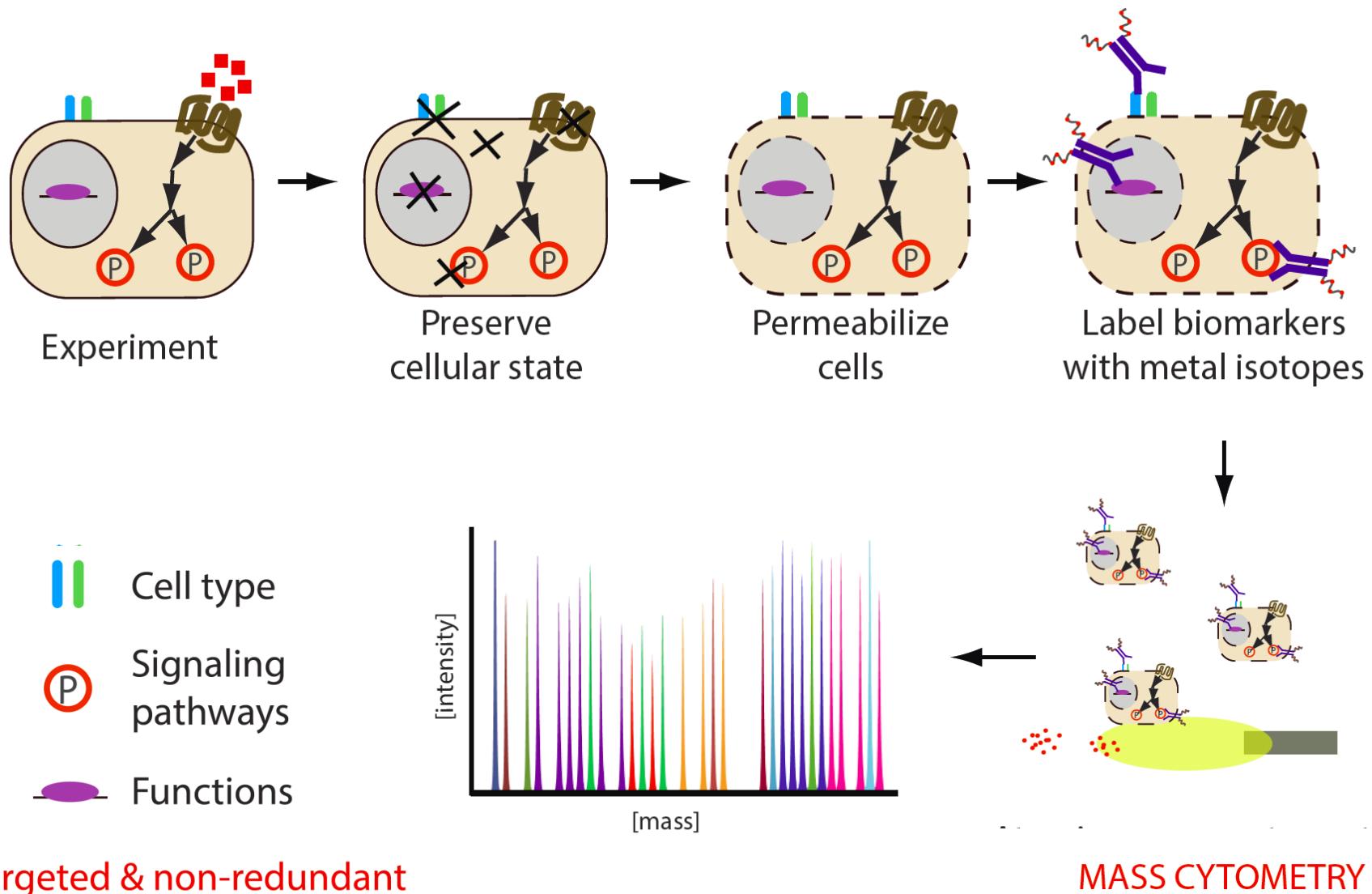
#2020-workshop-byod-help
<https://krishnaswamylab.org/get-help>

Challenges and Opportunities in Single Cell Data

Droplet-based Technologies



Single-Cell Proteomics: Mass Cytometry



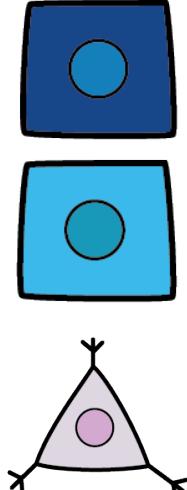
Single Cell Data

- Each cell is a vector of measurements
 - e.g. Cell A = [40 0 20 18 5 0 ...]
- The whole data is a matrix with many observations (cells) and features (proteins, genes)

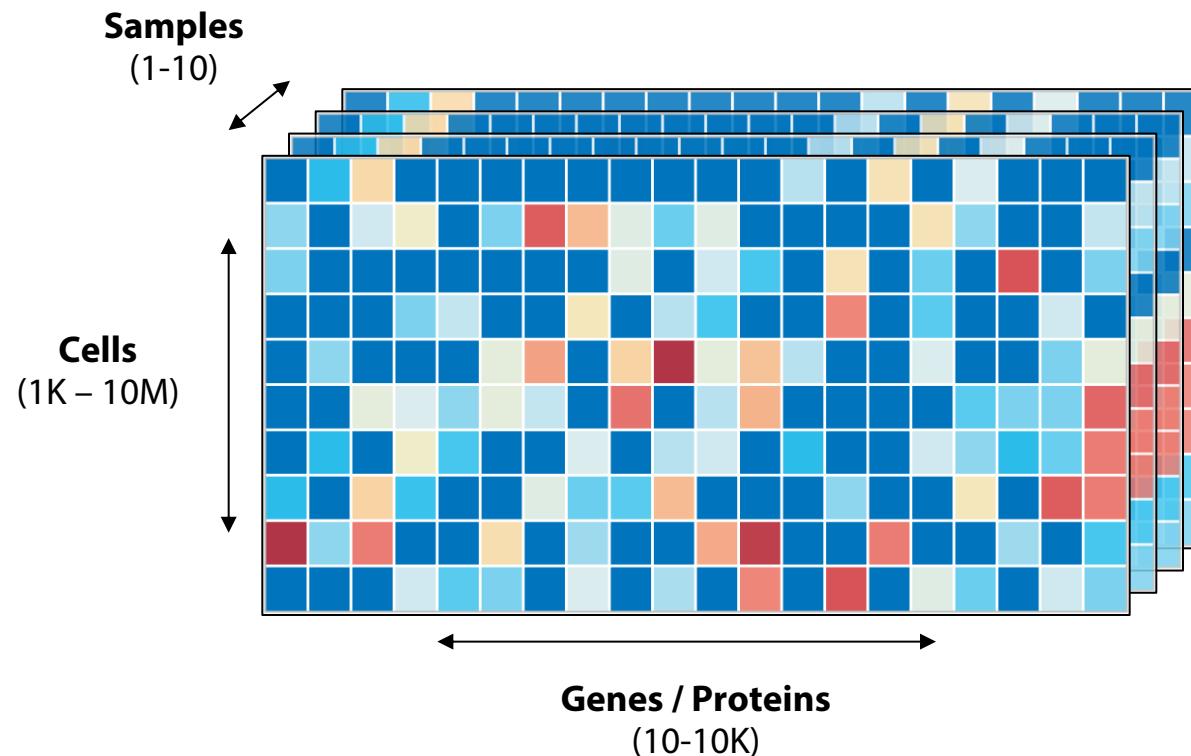
Features
(e.g. genes)

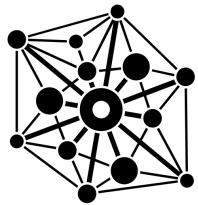
	X	Y	Z
A	10	20	70
B	20	40	140
C	20	0	80

Observations
(e.g. cells)

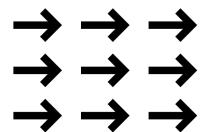


Single Cell Data





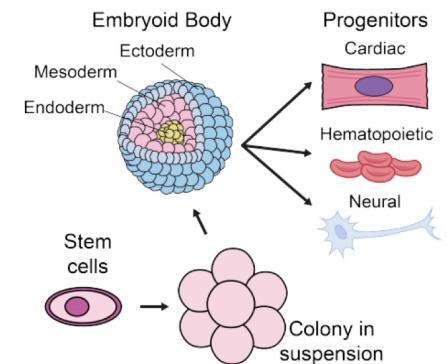
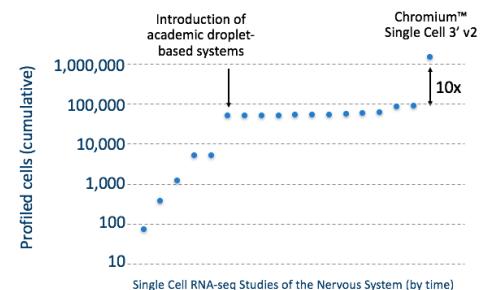
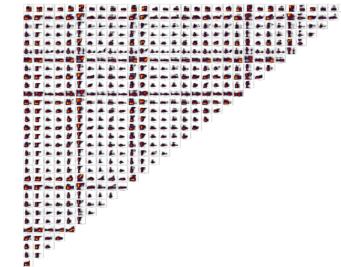
High Dimensional



High Throughput



Heterogeneous



Many dimensions = many measurements



Diagnoses



labs



drug response assays



ECG



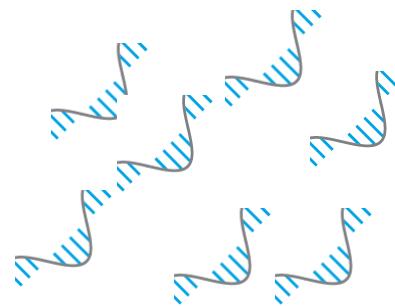
Gene 1



Gene 2



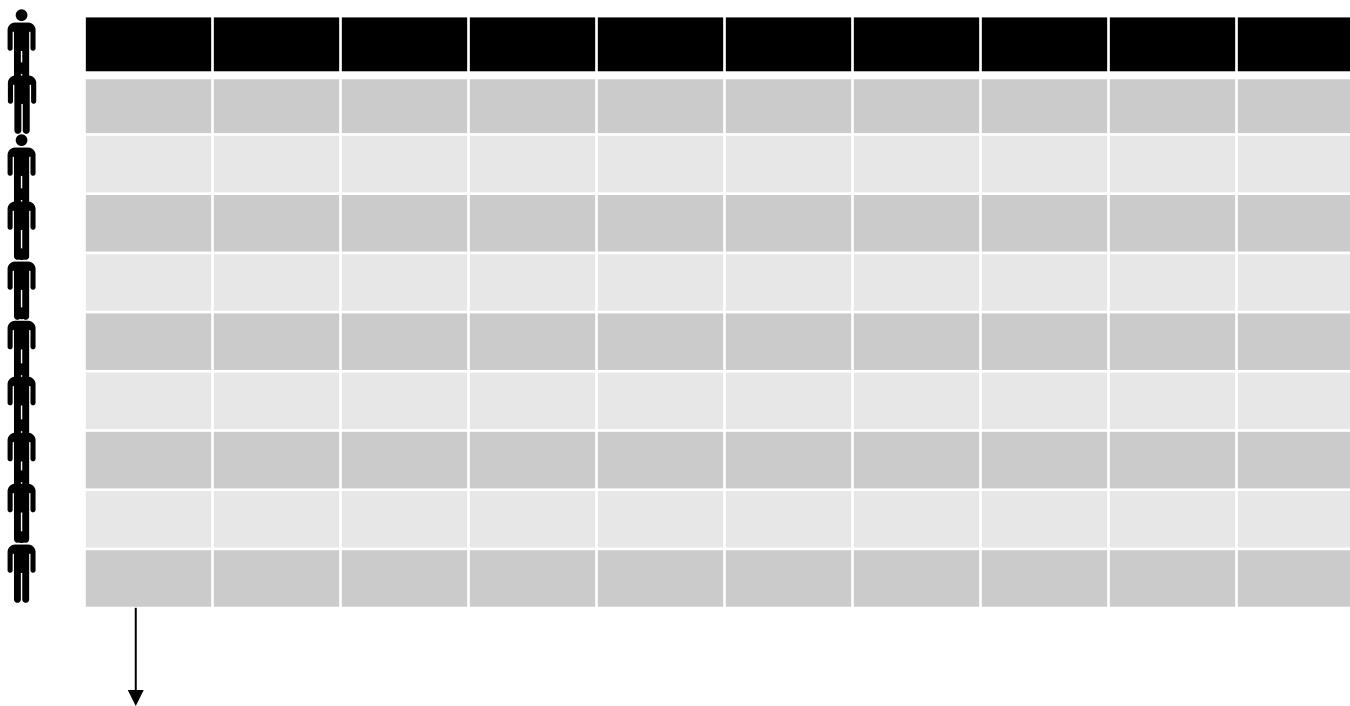
Gene 3



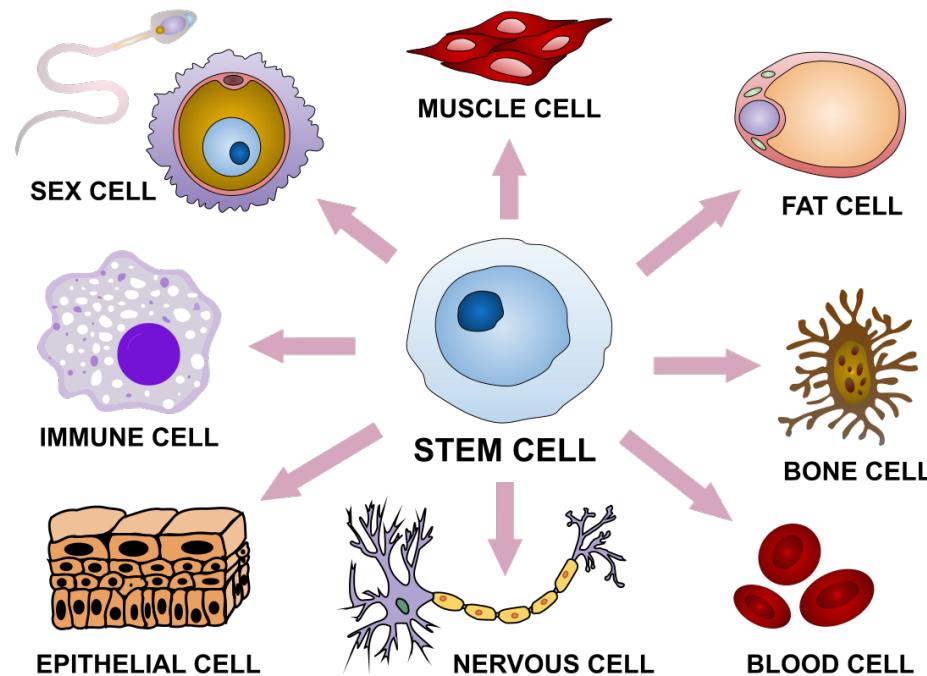
Proteins



High Throughput = Many observations

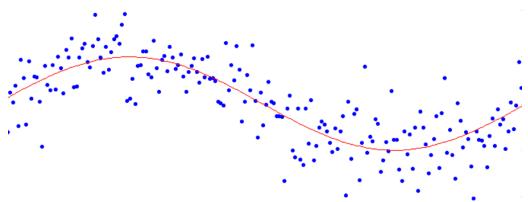


Heterogeneous Observations

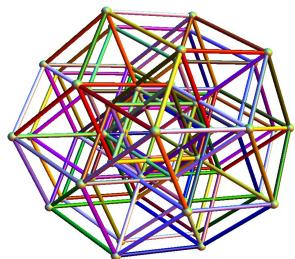
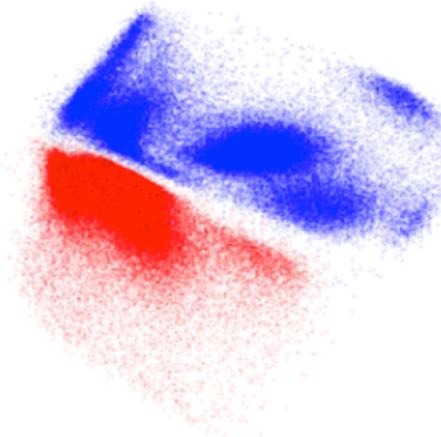


Challenges

Noise



Batch Effects

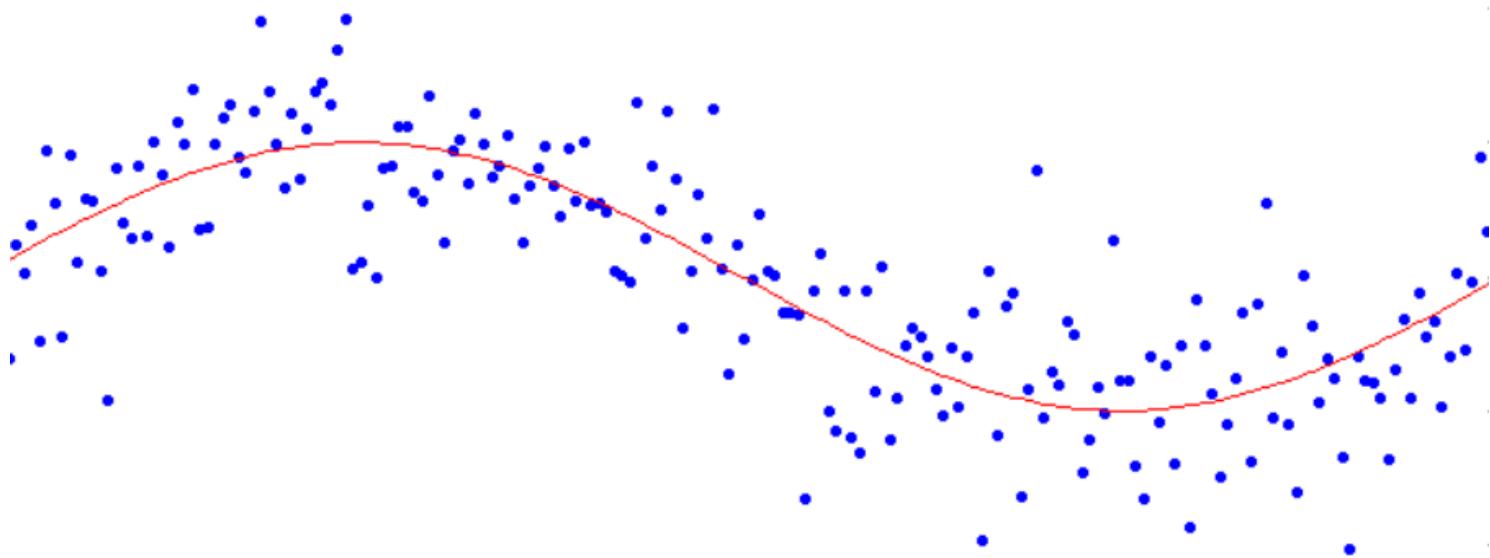


Dimensionality

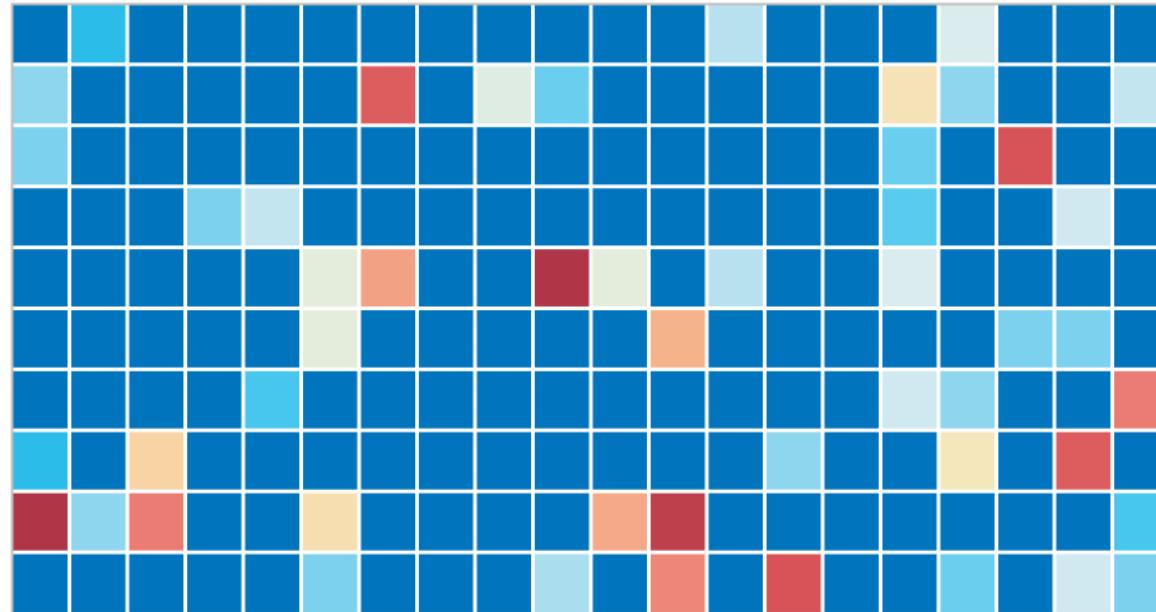


Scale

Noise



Dropout

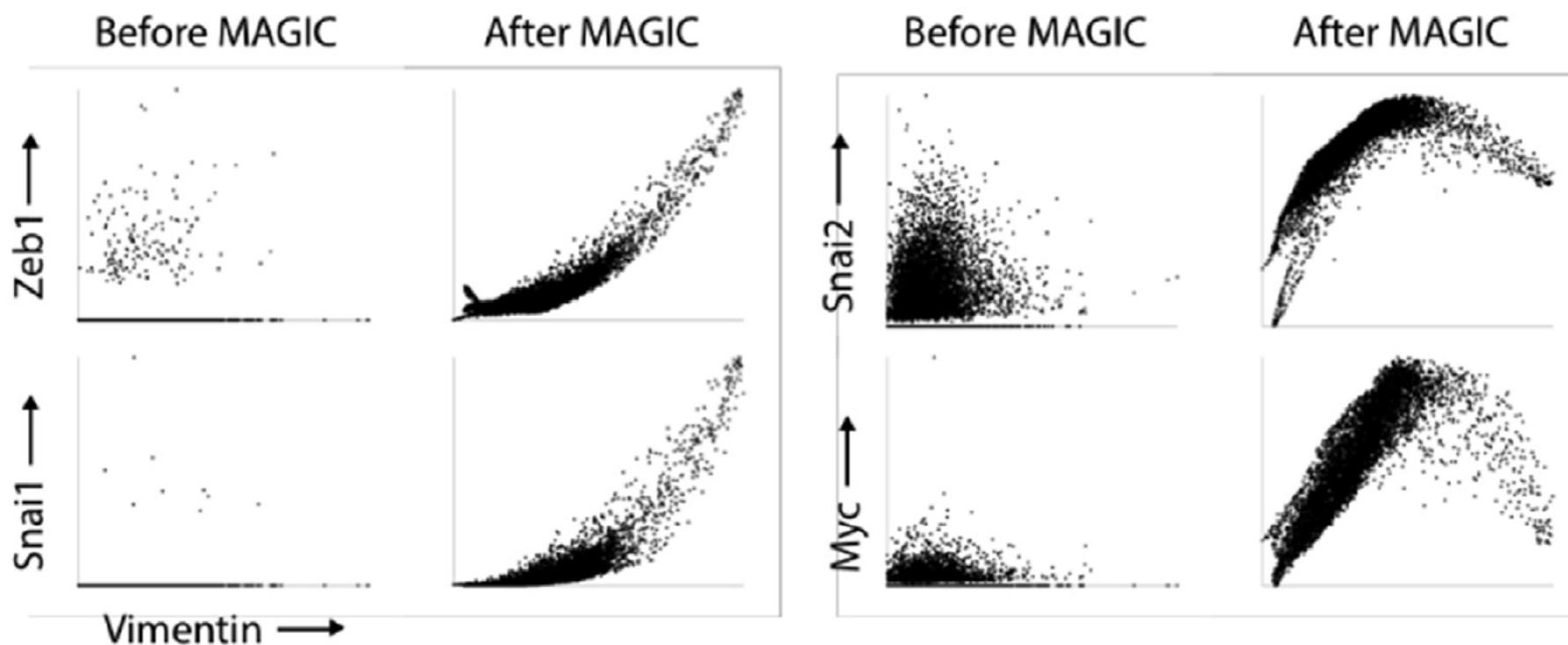


Dropout vs Missing Data

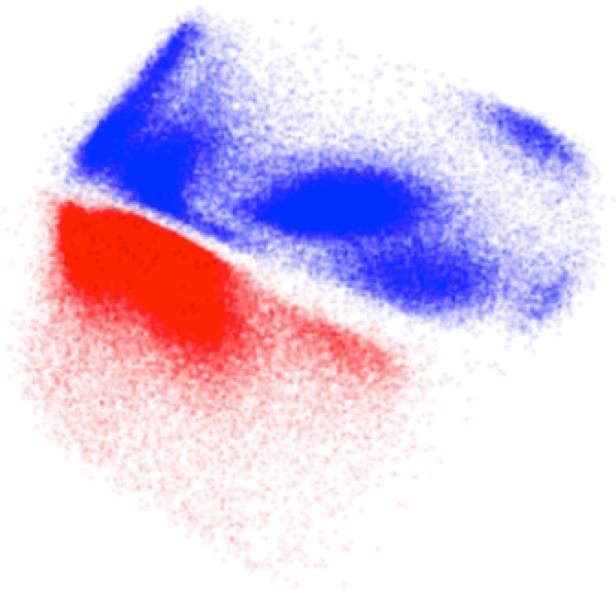
Missing values

PassengerId	Survived	Pclass	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	0	3	male	22	1	0	A/5 21171	7.25		S
2	1	1	female	38	1	0	PC 17599	71.2033	C85	C
3	1	3	female	26	0	0	STON/O2. 3101282	7.925		S
4	1	1	female	35	1	0	113803	53.1	C123	S
5	0	3	male	35	0	0	373450	8.05		S
6	0	3	male		0	0	330877	8.4583		Q

Denoising Data

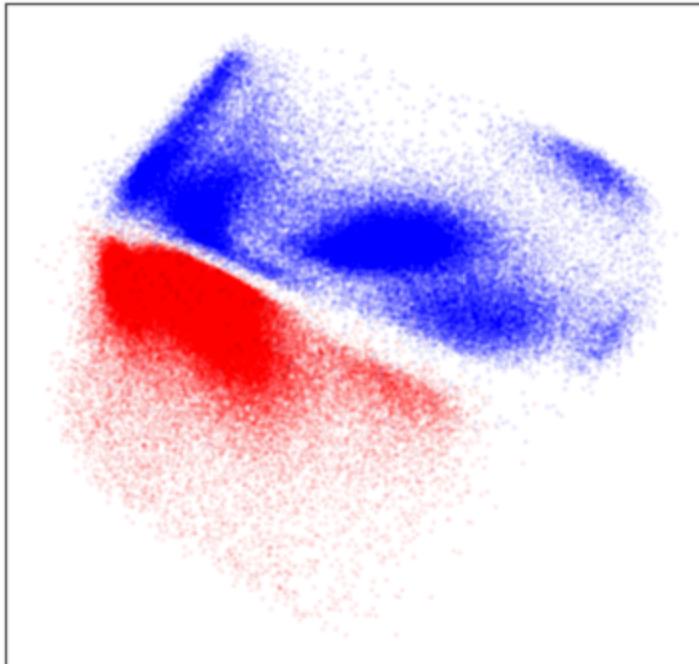


Batch Effects

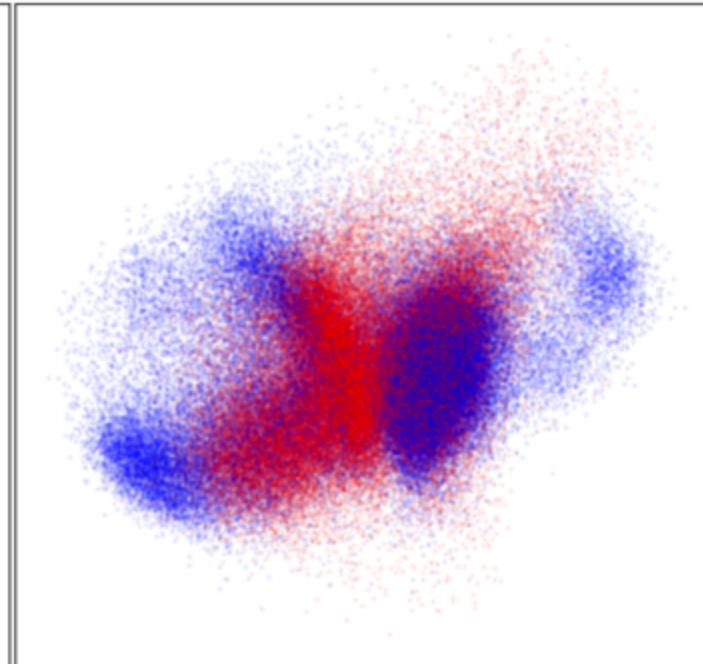


- Systematic differences between samples due to machine calibration, ambient environmental effects
- Variation that is uninteresting to examine and confounds biological variation
- Renders samples uncomparable

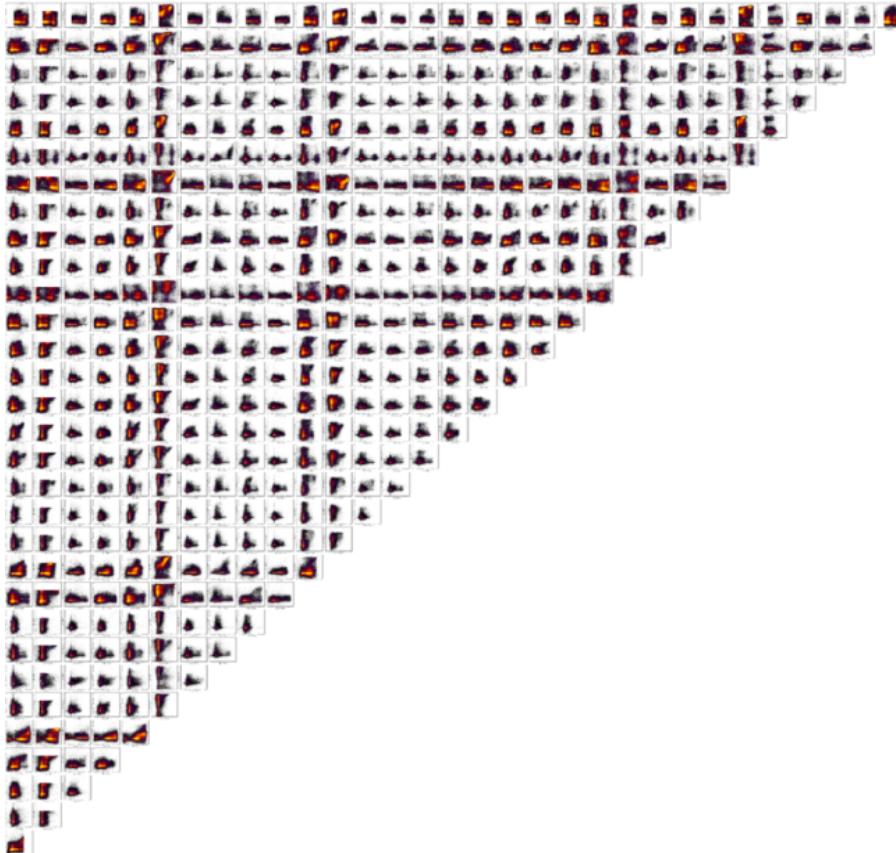
Before MMD



After MMD

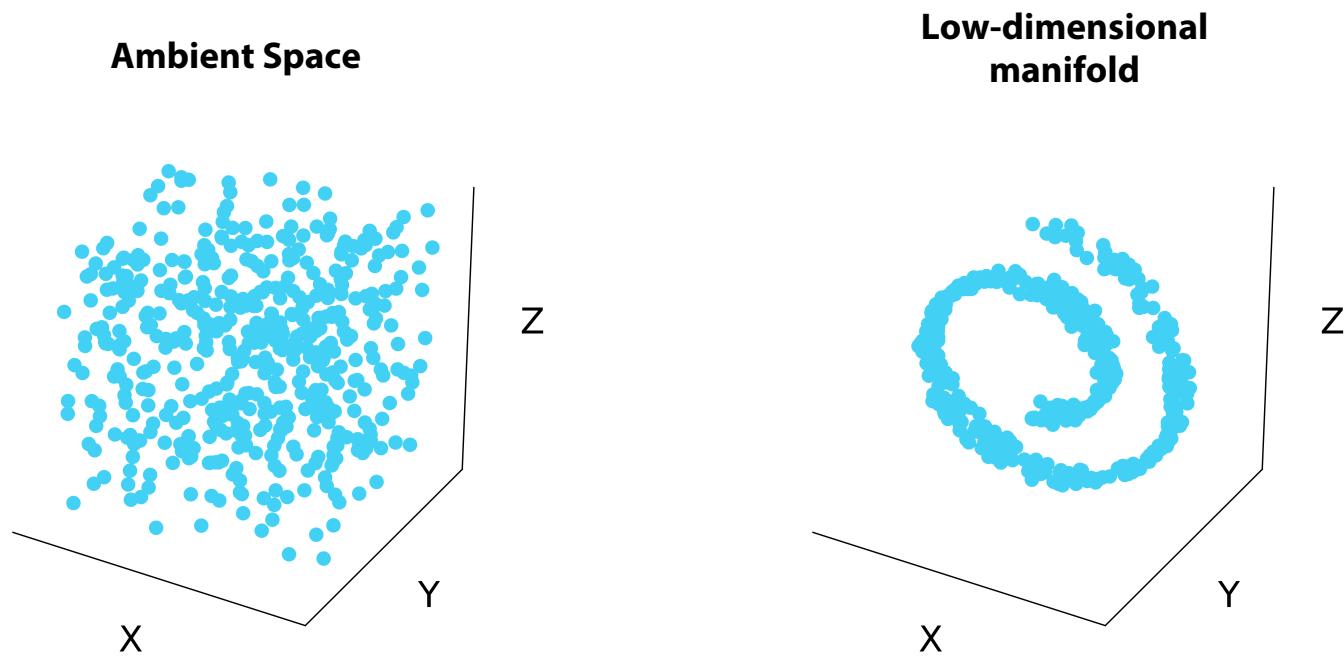


High Dimensionality

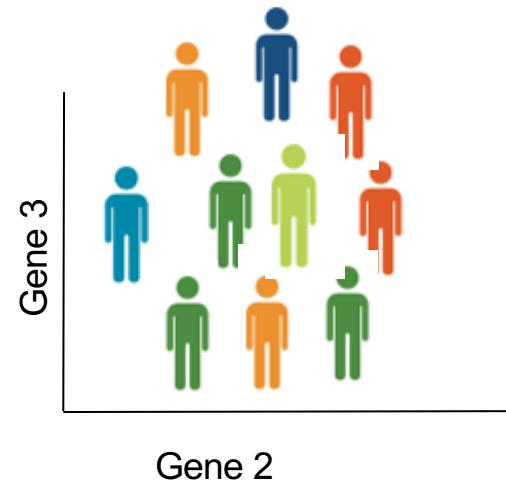
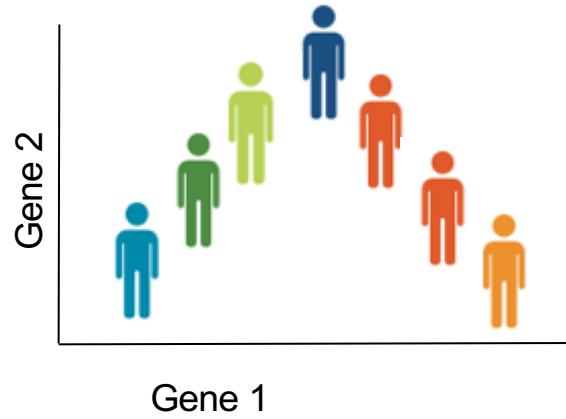


- Systematic differences between samples due to machine calibration, ambient environmental effects
- Variation that is uninteresting to examine and confounds biological variation
- Renders samples incomparable

Latent structure in high dimensional data

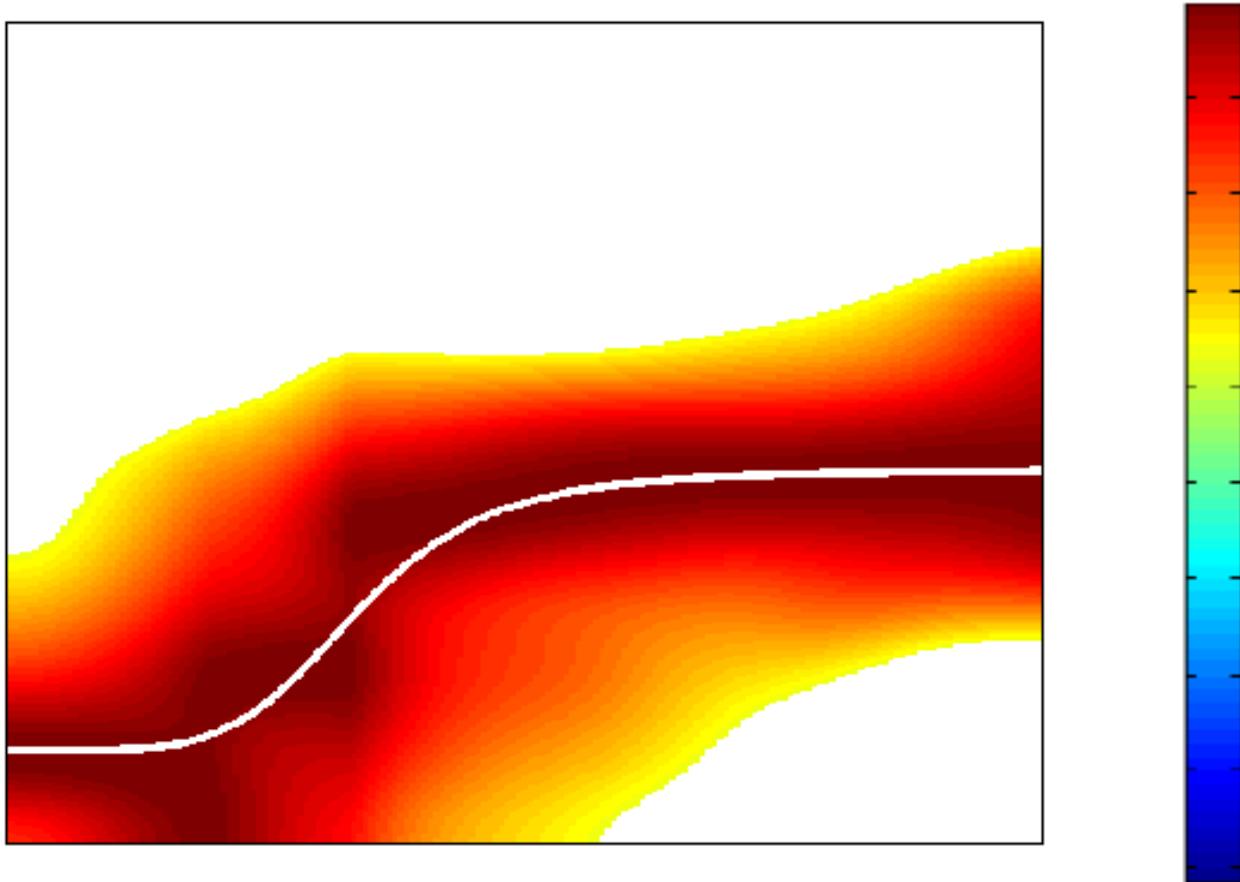
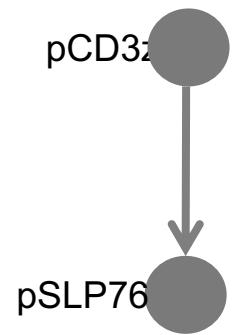


Gene-gene Relationships

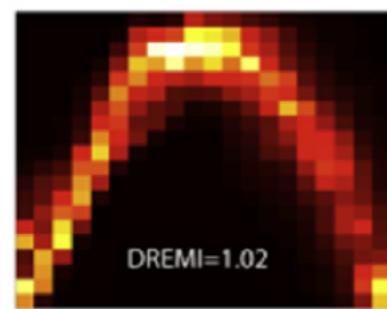
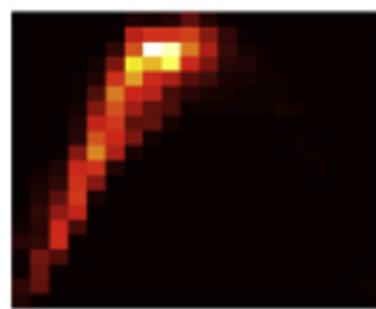
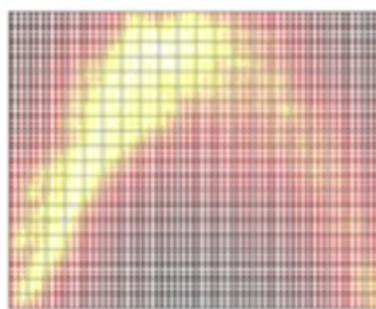
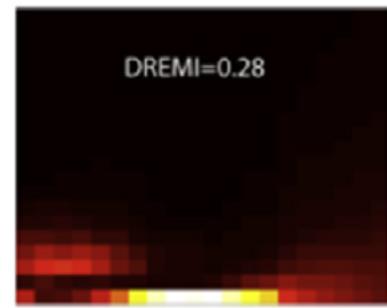
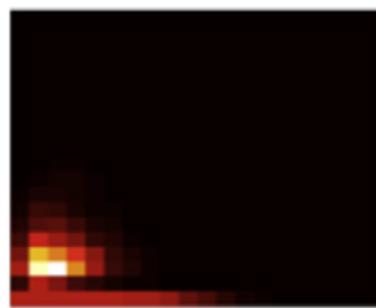
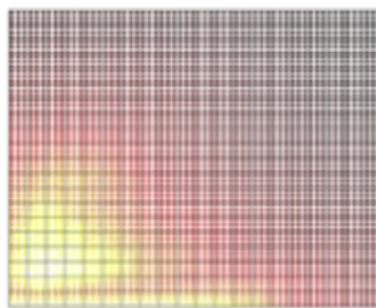
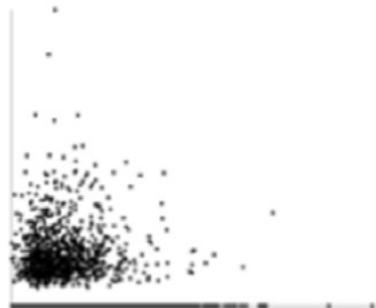


Relationship between features?

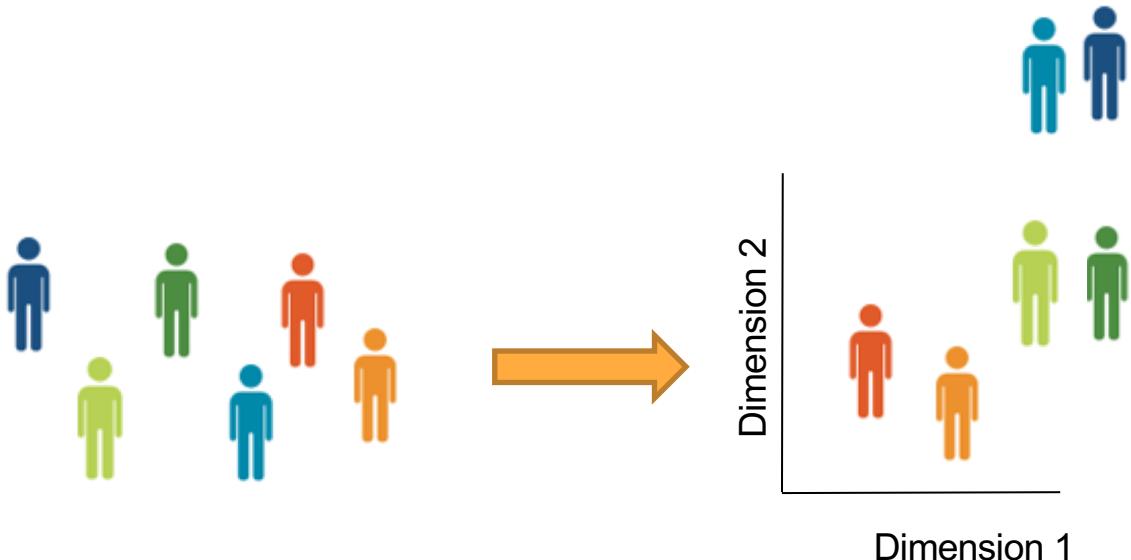
Non-linear Relationships



Mutual Information



Embedding reveals structure

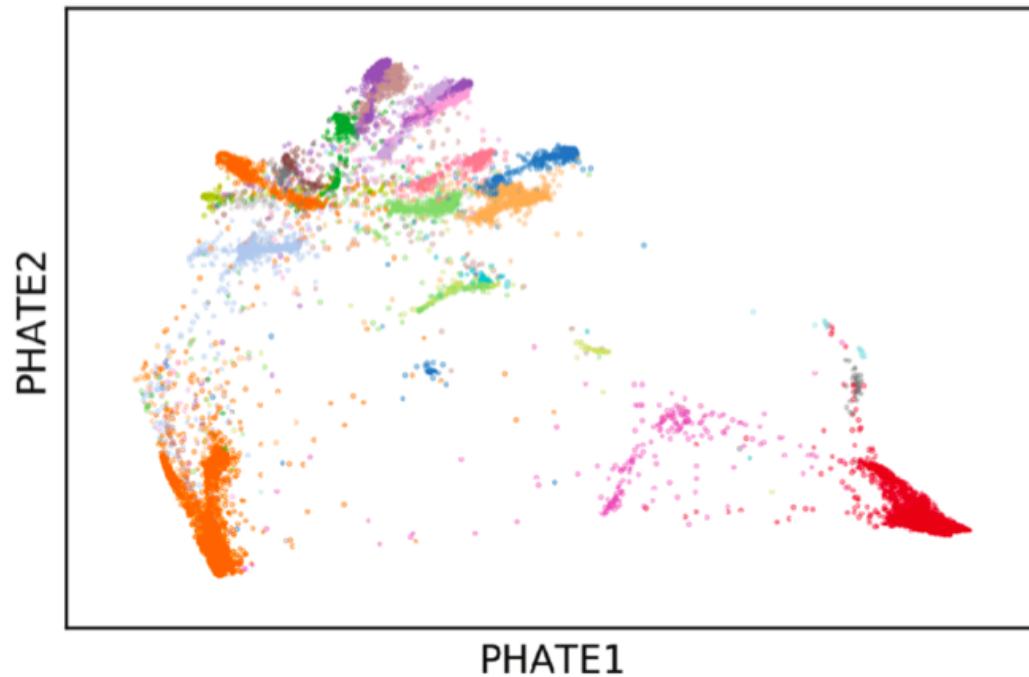


Use high dimensional features and high throughput to understand shape of data

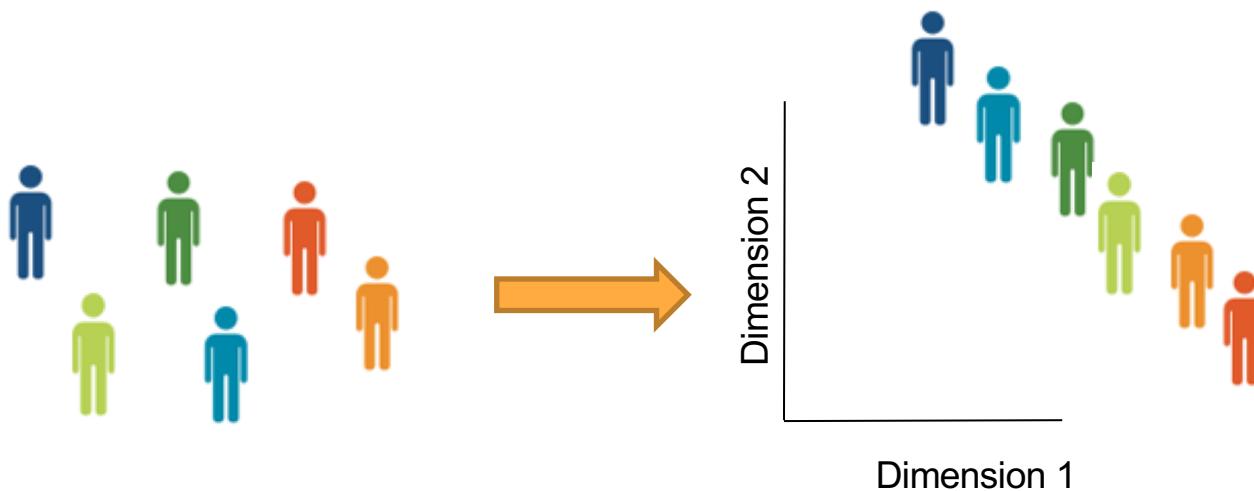
Cluster structure



Retinal Bipolar Cells



Progression continuum

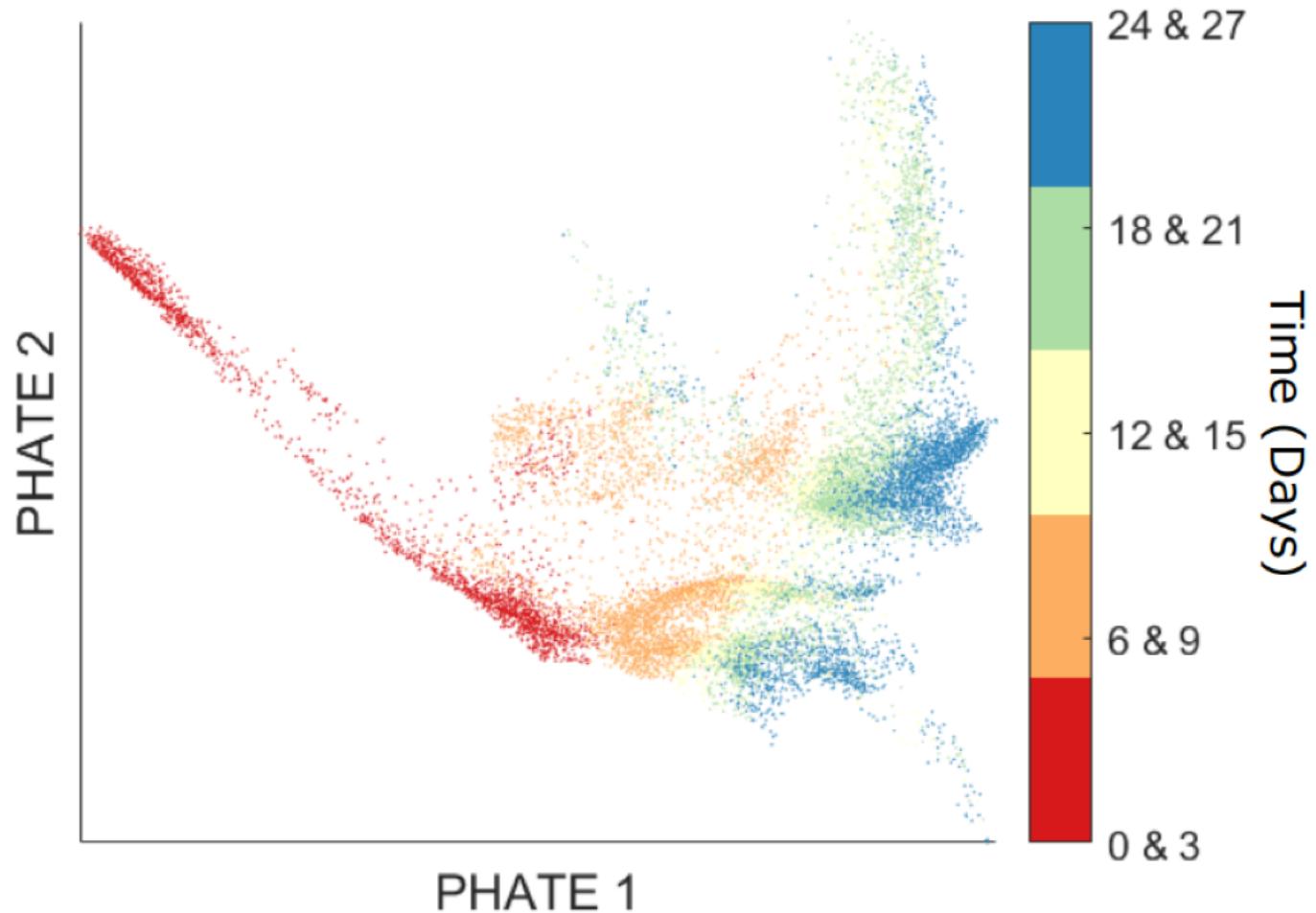


Use high dimensional features and high throughput to understand shape of data

Pseudotime

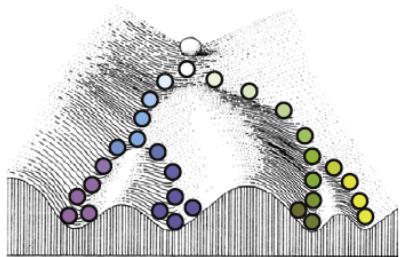


Progressions

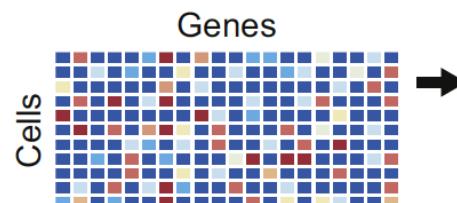


Manifold Learning

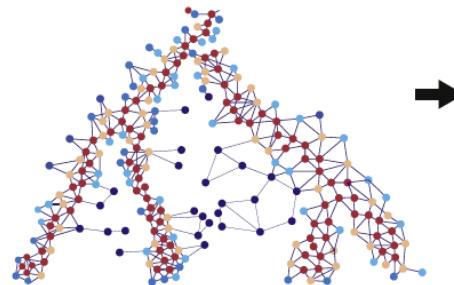
Cells are sampled from an underlying manifold



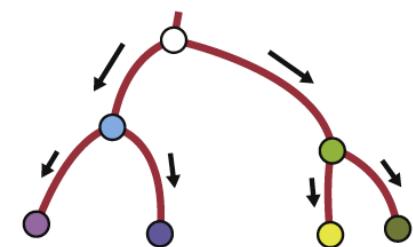
Each cell is represented by a vector of gene expression



Neighborhood structure of the observations is identified

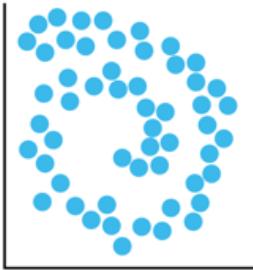


The latent manifold is learned from the data

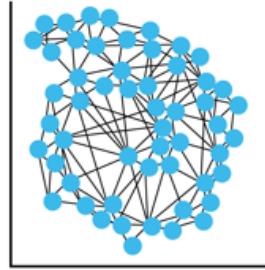


Understanding the shape of data

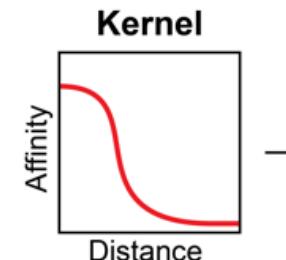
Data in two dimensions



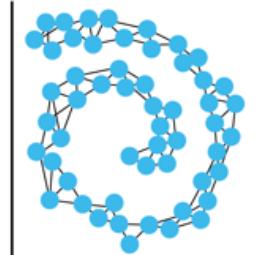
Distances between all points are calculated



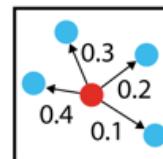
A kernel function calculates affinities from distance



Only local relationships are preserved

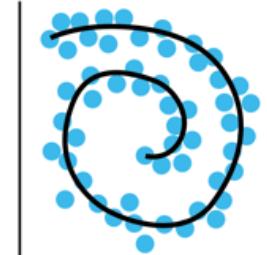


Diffusion shares information between nodes



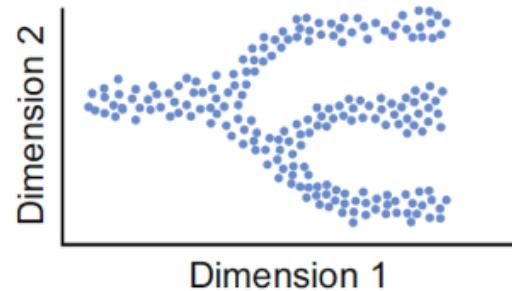
Diffusion distance
≈
Random walk dist.

Underlying manifold is calculated

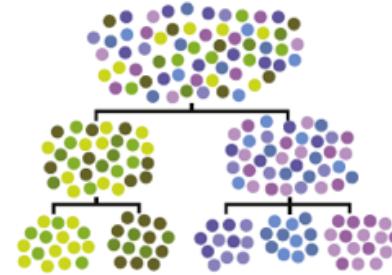


Analysis Tasks

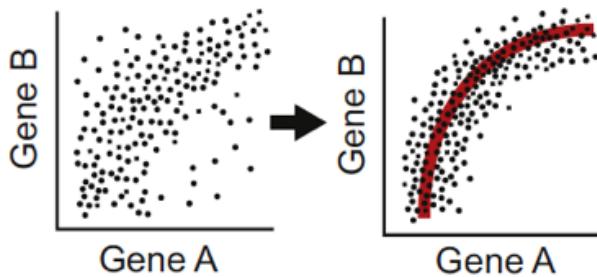
Vizualization



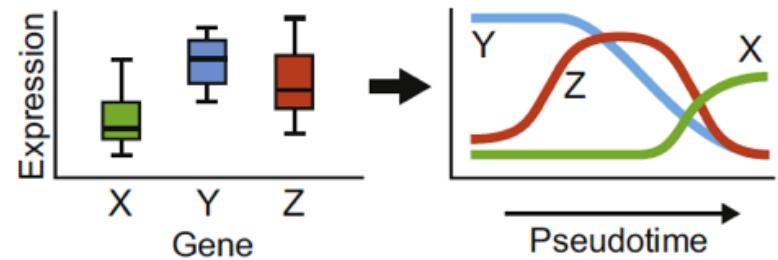
Clustering



Denoising



Pseudotime analysis



Preprocessing single-cell data

Current best practices in single-cell RNA-seq analysis: a tutorial

Malte D Luecken¹  & Fabian J Theis^{1,2,*} 

Abstract

Single-cell RNA-seq has enabled gene expression to be studied at an unprecedented resolution. The promise of this technology is attracting a growing user base for single-cell analysis methods. As more analysis tools are becoming available, it is becoming increasingly difficult to navigate this landscape and produce an up-to-date workflow to analyse one's data. Here, we detail the steps of a typical single-cell RNA-seq analysis, including pre-processing (quality control, normalization, data correction, feature selection, and dimensionality reduction) and cell- and gene-level downstream analysis. We formulate current best-practice recommendations for these steps based on independent comparison studies. We have integrated these best-practice recommendations into a workflow, which we apply to a public dataset to further illustrate how these steps work in practice. Our documented case study can be found at <https://www.github.com/theislab/single-cell-tutorial>. This review will serve as a workflow tutorial for new entrants into the field, and help established users update their analysis pipelines.

Keywords analysis pipeline development; computational biology; data analysis tutorial; single-cell RNA-seq

DOI 10.15252/msb.20188746 | Received 16 November 2018 | Revised 15 March 2019 | Accepted 3 April 2019

Mol Syst Biol. (2019) 15: e8746

Introduction

In recent years, single-cell RNA sequencing (scRNA-seq) has significantly advanced our knowledge of biological systems. We have been able to both study the cellular heterogeneity of zebrafish, frogs

outline current best practices to lay a foundation for future analysis standardization.

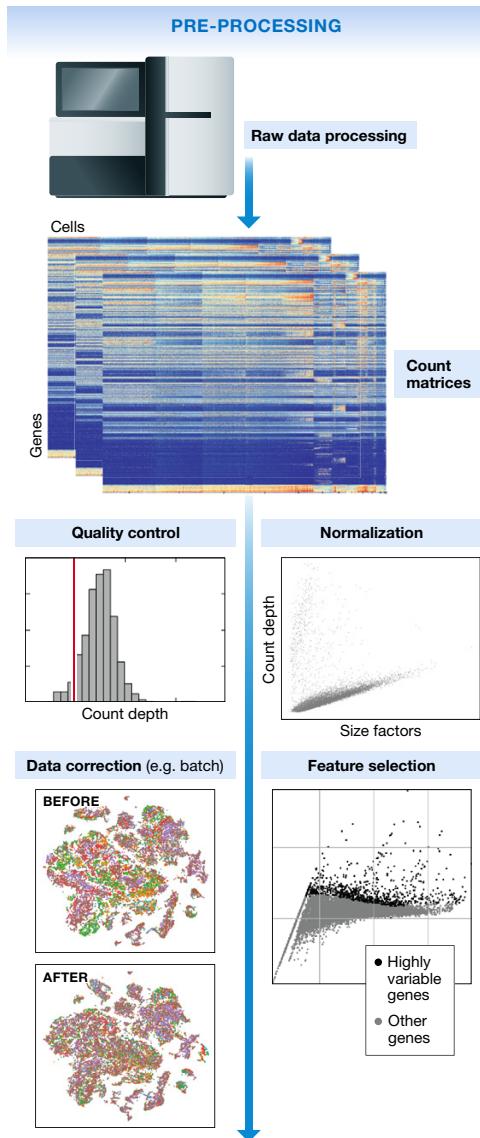
The challenges to standardization include the growing number of analysis methods (385 tools as of 7 March 2019) and exploding dataset sizes (Angerer *et al.*, 2017; Zappia *et al.*, 2018). We are continuously finding new ways to use the data at our disposal. For example, it has recently become possible to predict cell fates in differentiation (La Manno *et al.*, 2018). While the continuous improvement of analysis tools is beneficial for generating new scientific insight, it complicates standardization.

Further challenges for standardization lie in technical aspects. Analysis tools for scRNA-seq data are written in a variety of programming languages—most prominently R and Python (Zappia *et al.*, 2018). Although cross-environment support is growing (preprint: Scholz *et al.*, 2018), the choice of programming language is often also a choice between analysis tools. Popular platforms such as Seurat (Butler *et al.*, 2018), Scater (McCarthy *et al.*, 2017), or Scanpy (Wolf *et al.*, 2018) provide integrated environments to develop pipelines and contain large analysis toolboxes. However, out of necessity these platforms limit themselves to tools developed in their respective programming languages. By extension, language restrictions also hold true for currently available scRNA-seq analysis tutorials, many of which revolve around the above platforms (R and bioconductor tools: <https://github.com/drissi/bioc2016singlecell> and <https://hemberg-lab.github.io/scRNA.seq.course/>; Lun *et al.*, 2016b; Seurat: https://satijalab.org/seurat/get_started.html; Scanpy: <https://scanpy.readthedocs.io/en/stable/tutorials.html>).

Considering the above-mentioned challenges, instead of targeting a standardized analysis pipeline, we outline current best practices and common tools independent of programming language. We guide the reader through the various steps of a scRNA-seq analysis pipeline (Fig 1), present current best practices, and discuss analysis

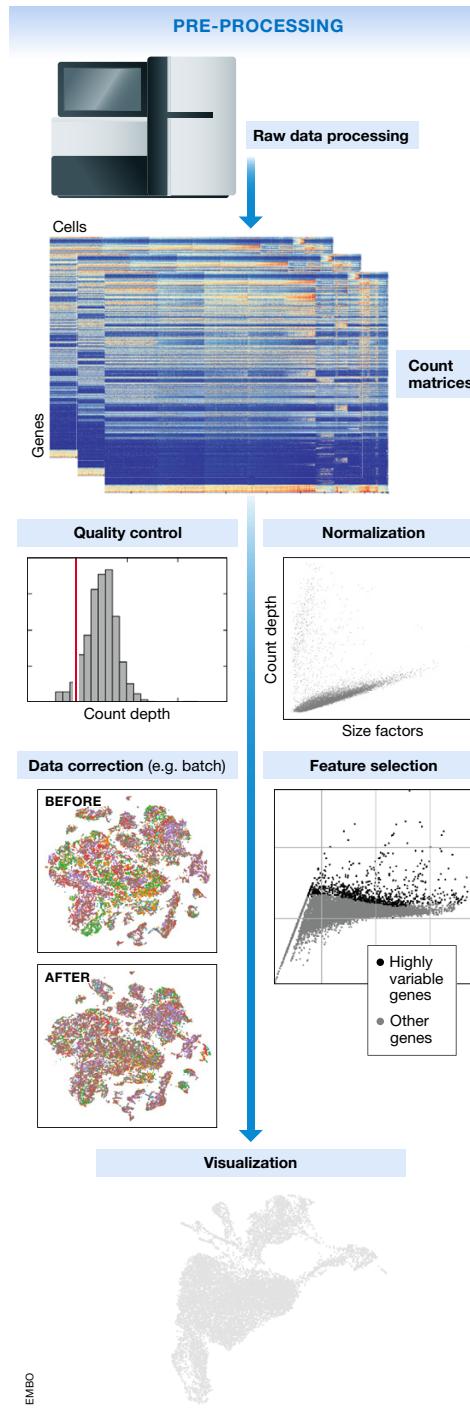
Standard Single-Cell RNA-seq Workflow

1. Sequencing and read mapping
2. Quality control and filtering
3. Normalization
4. Data Correction



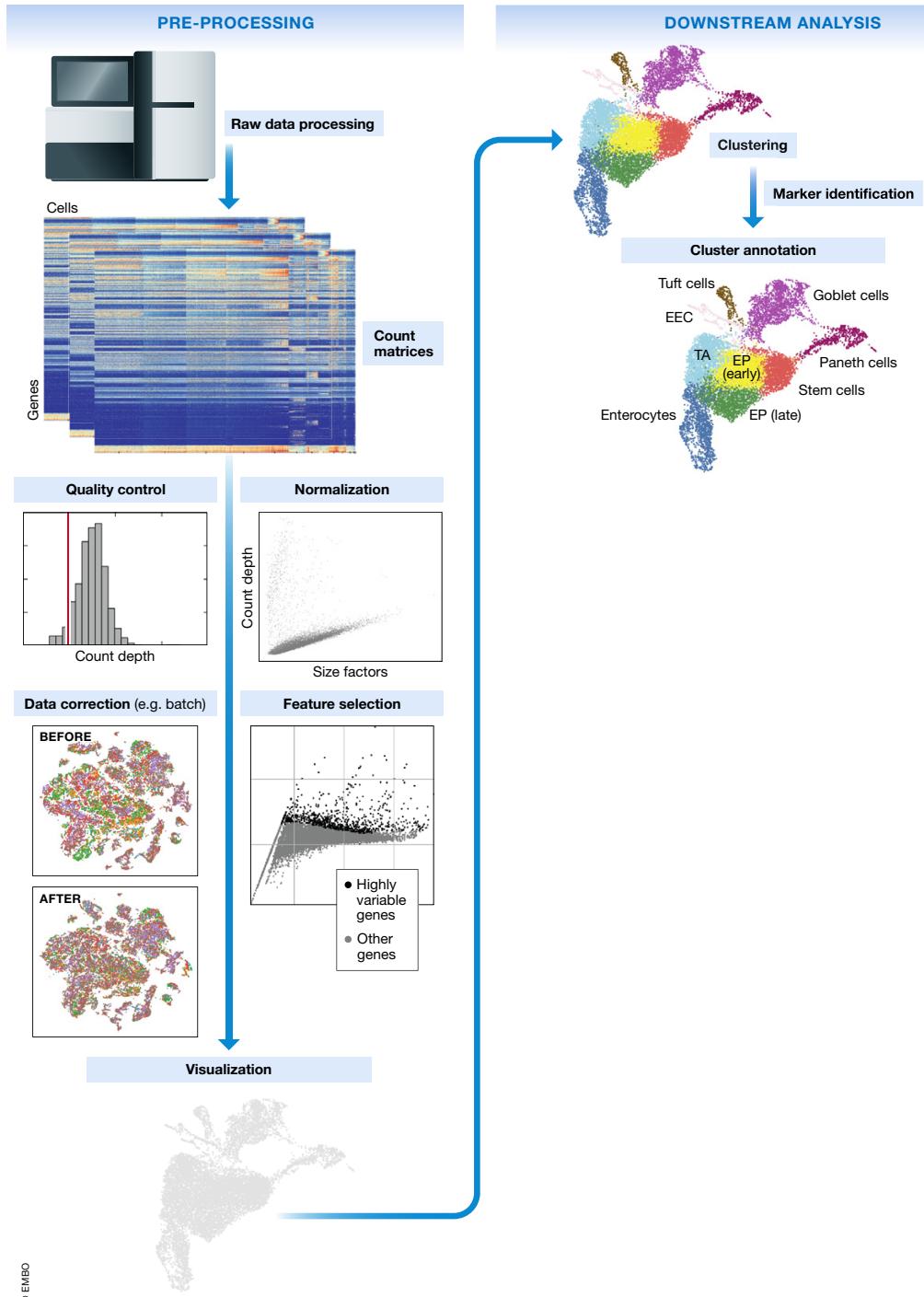
Standard Single-Cell RNA-seq Workflow

1. Sequencing and read mapping
2. Quality control and filtering
3. Normalization
4. Data Correction
5. Dimensionality reduction and visualization
6. Downstream analysis
 1. Clustering
 2. Trajectory inference



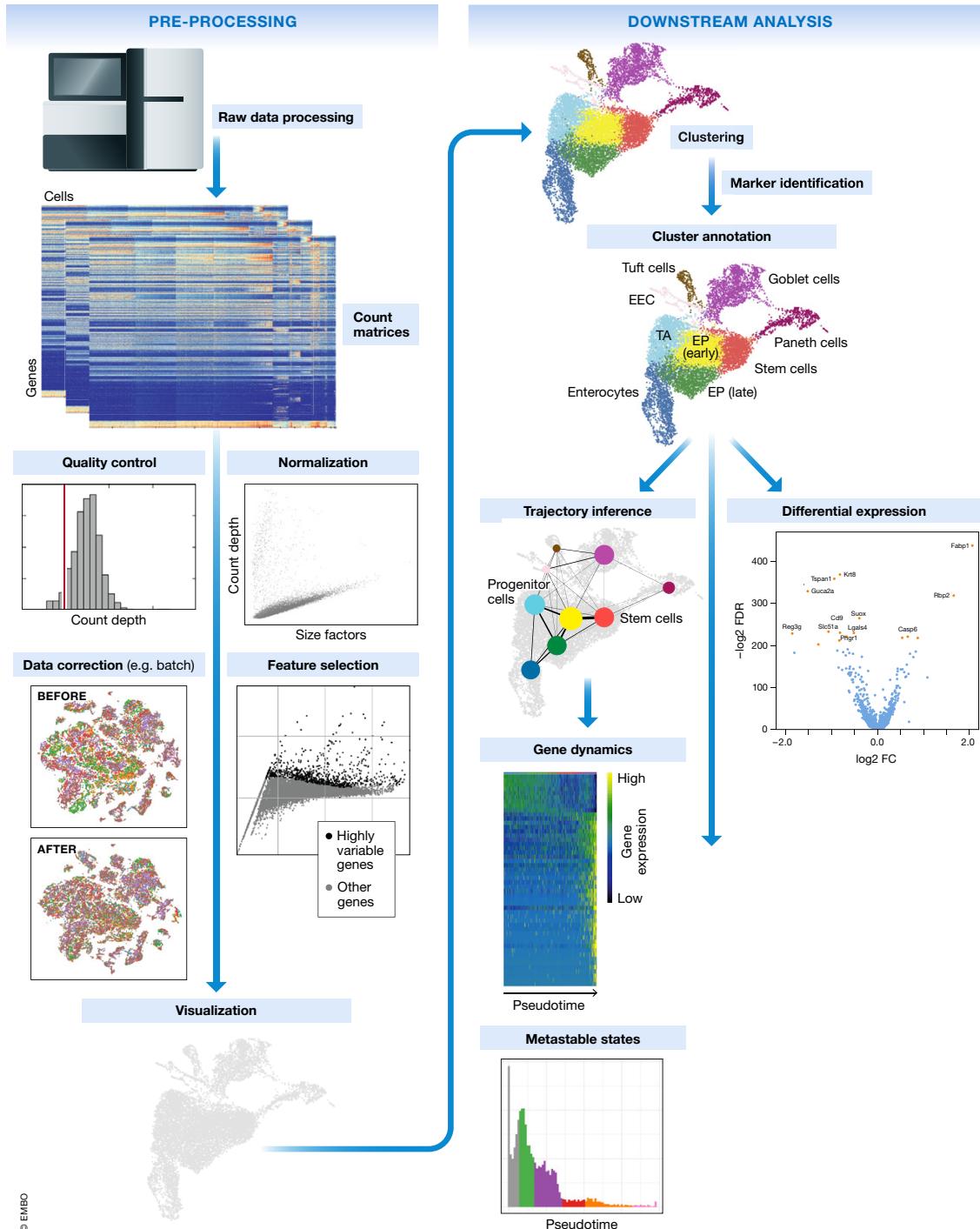
Standard Single-Cell RNA-seq Workflow

1. Sequencing and read mapping
2. Quality control and filtering
3. Normalization
4. Data Correction
5. Dimensionality reduction and visualization
6. Downstream analysis
 1. Clustering
 2. Trajectory inference
 3. Differential expression



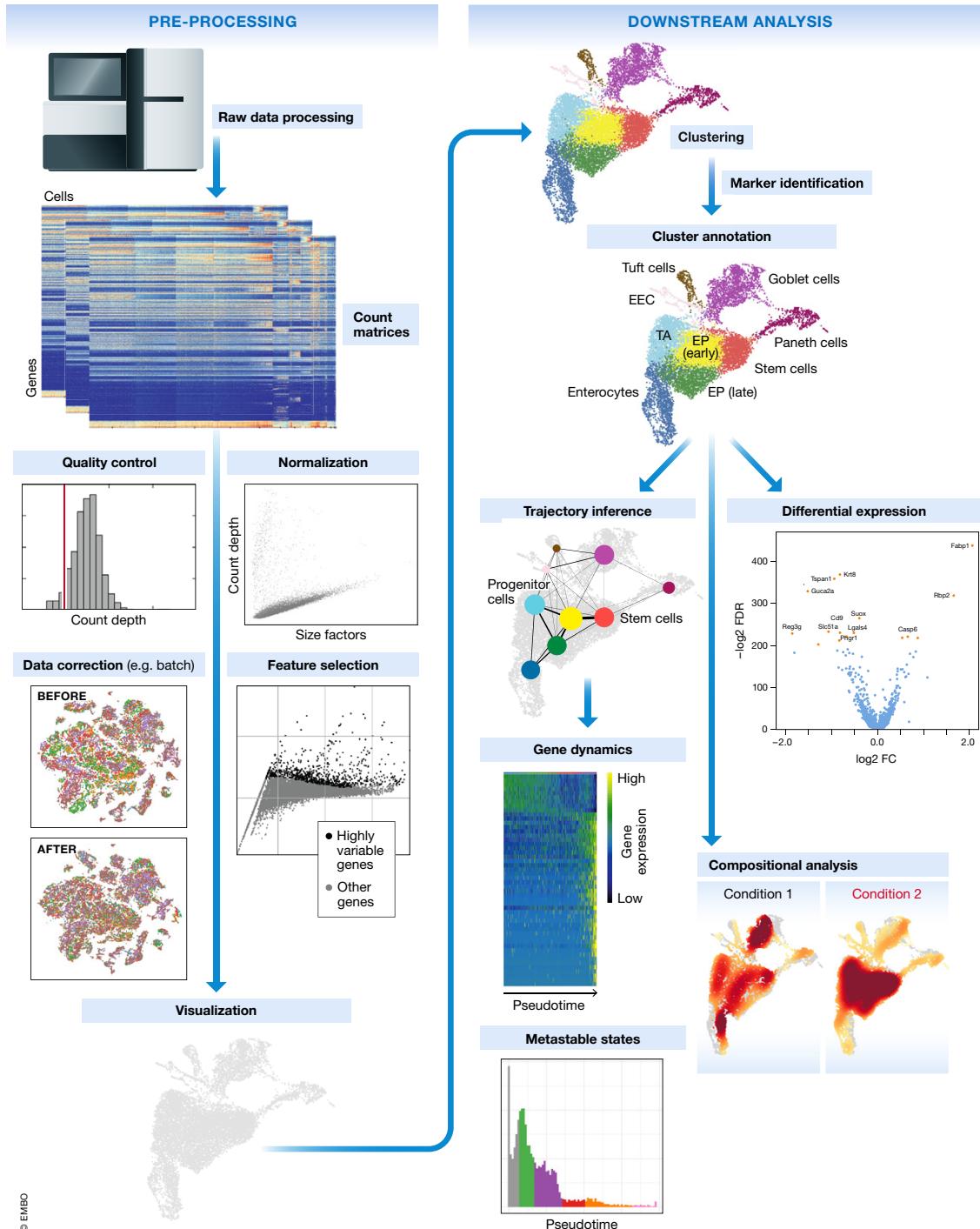
Standard Single-Cell RNA-seq Workflow

1. Sequencing and read mapping
2. Quality control and filtering
3. Normalization
4. Data Correction
5. Dimensionality reduction and visualization
6. Downstream analysis
 1. Clustering
 2. Trajectory inference
 3. Differential expression

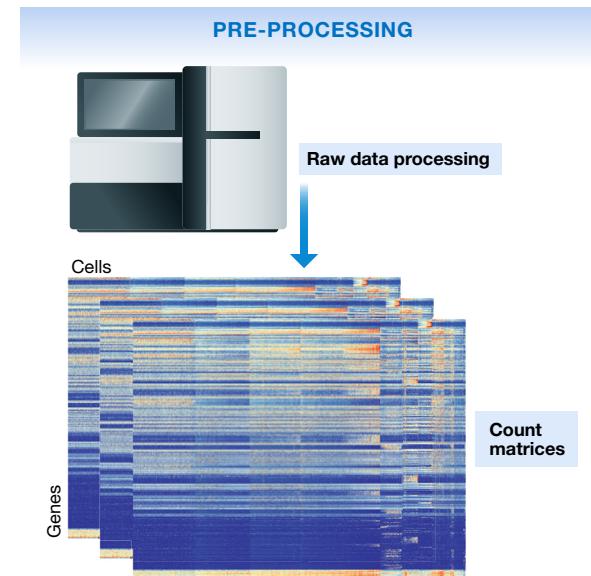
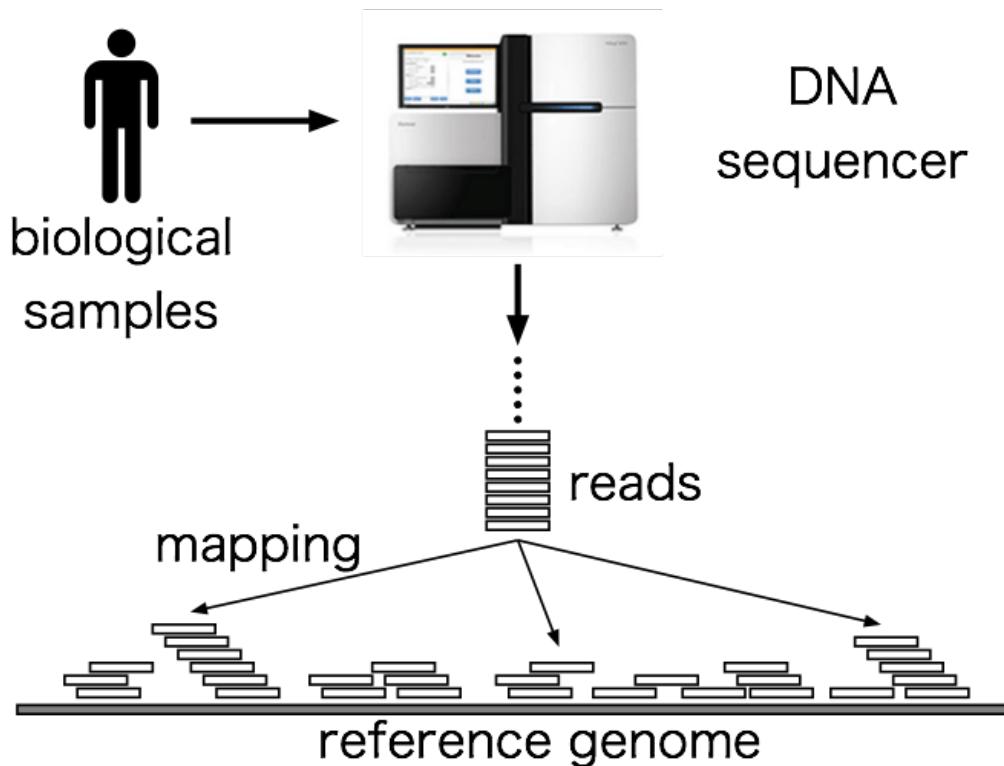


Standard Single-Cell RNA-seq Workflow

1. Sequencing and read mapping
2. Quality control and filtering
3. Normalization
4. Data Correction
5. Dimensionality reduction and visualization
6. Downstream analysis
 1. Clustering
 2. Trajectory inference
 3. Differential expression
7. Comparison of multiple conditions



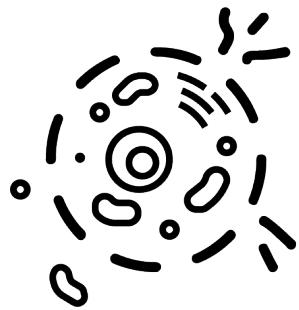
Step 1 - Sequencing & read mapping



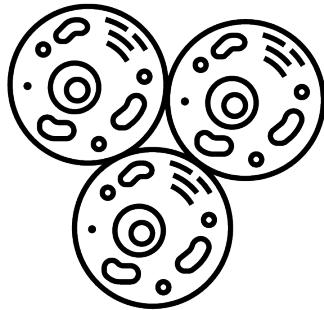
Building Counts Matrix → QC & Filtering → Normalization → Visualization → Analysis

Step 2 – Quality control and filtering

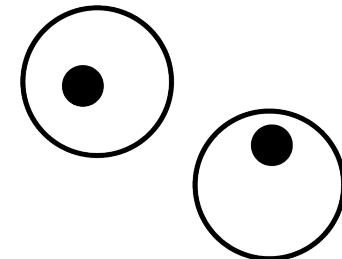
Dying cells



Multiplets



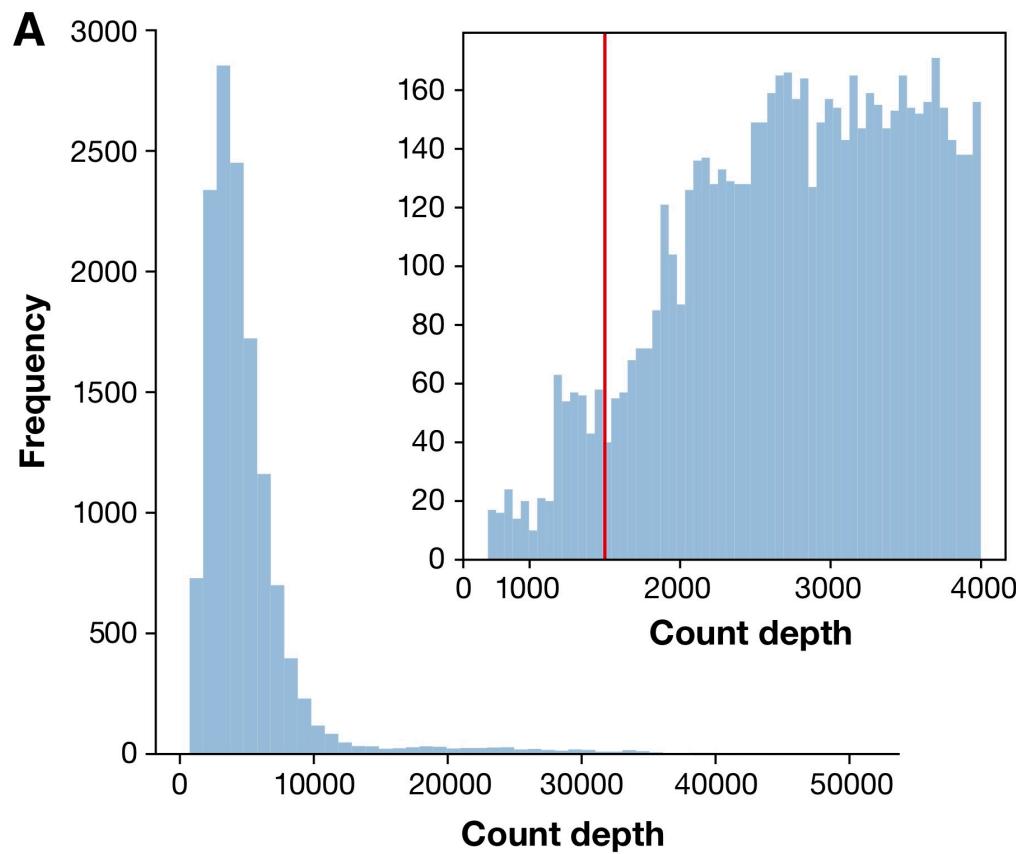
Empty Droplets



What could we look at to discriminate between dying cells, multiplets, or empty droplets and healthy single cells?

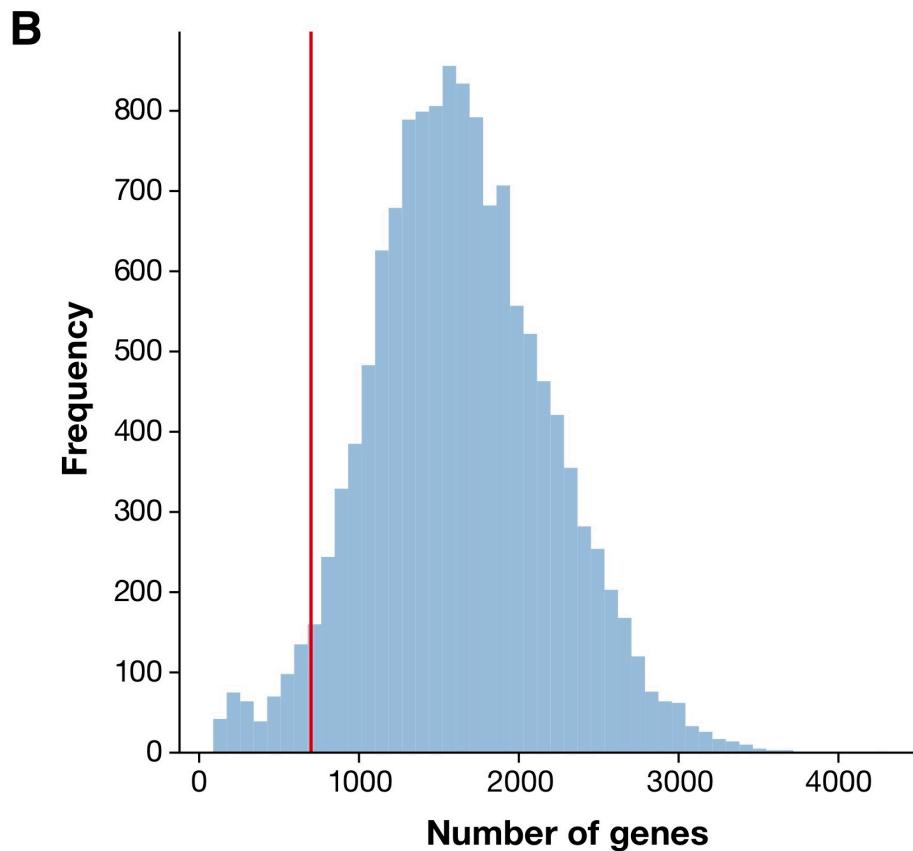
Top

Step 2 – Quality control and filtering



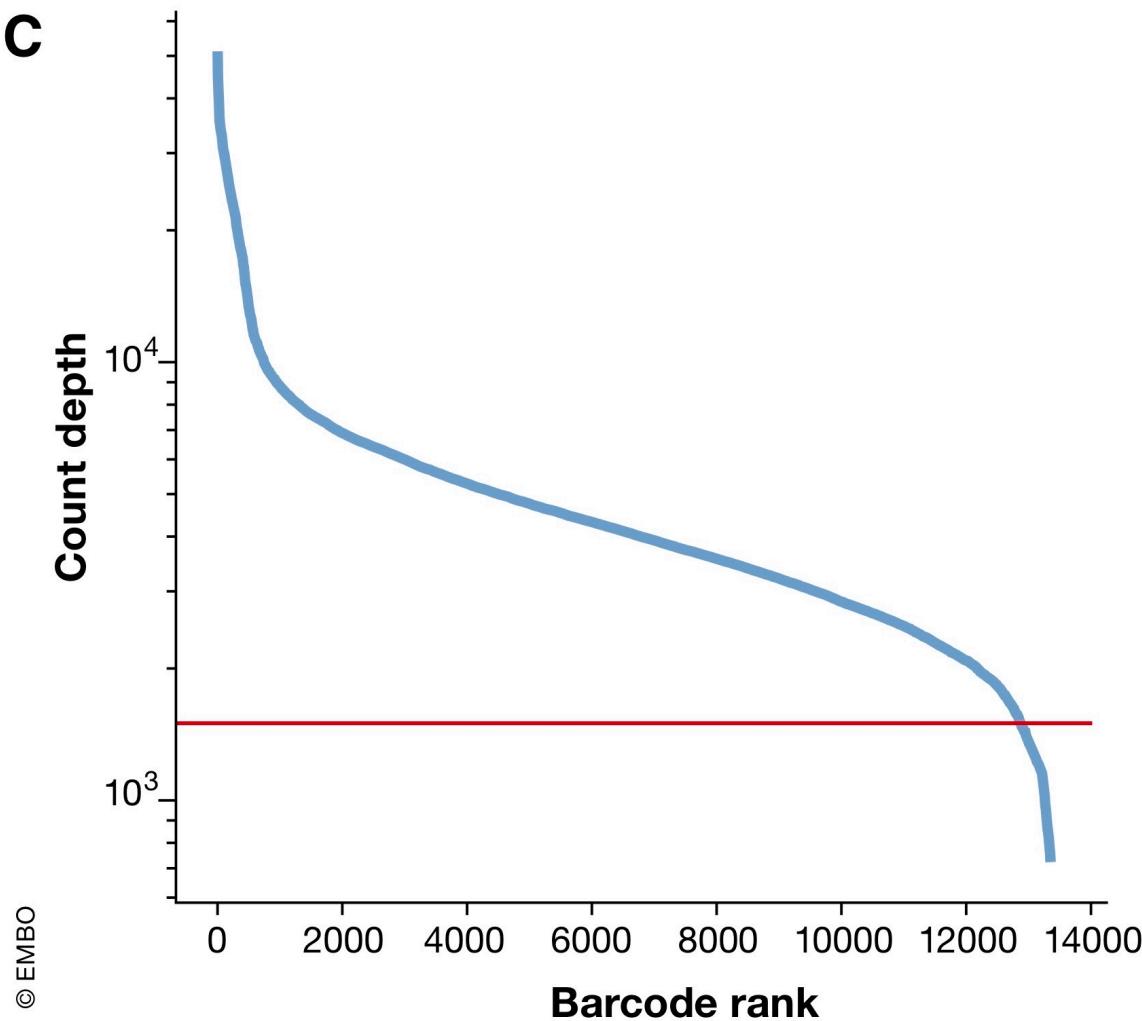
Building Counts Matrix → QC & Filtering → Normalization → Visualization → Analysis

Step 2 – Quality control and filtering



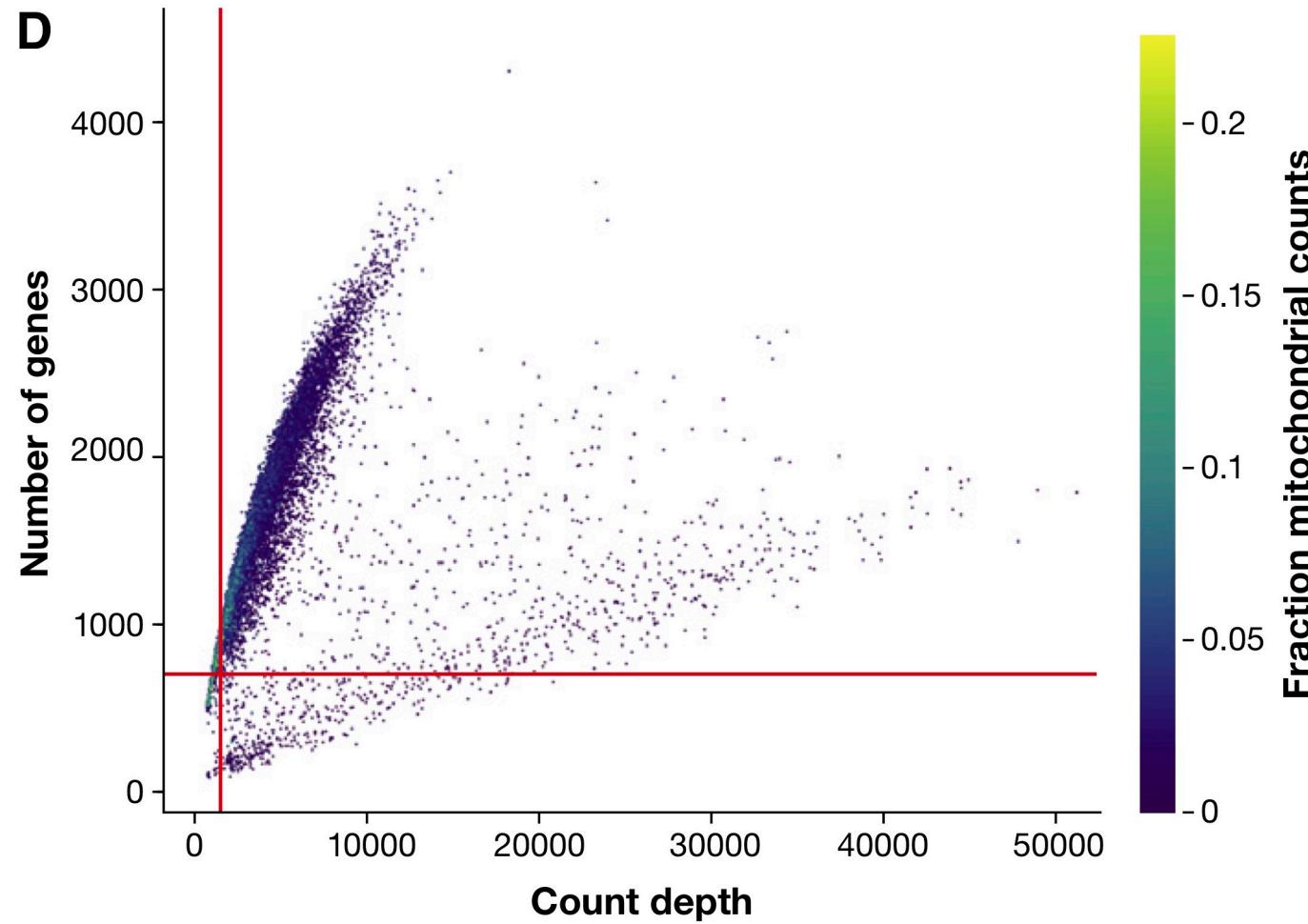
Building Counts Matrix → QC & Filtering → Normalization → Visualization → Analysis

Step 2 – Quality control and filtering



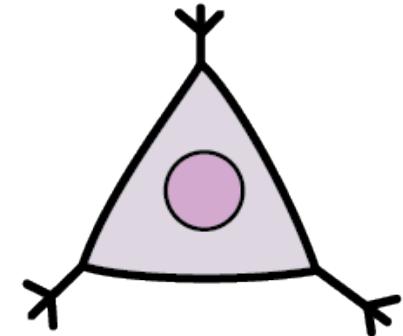
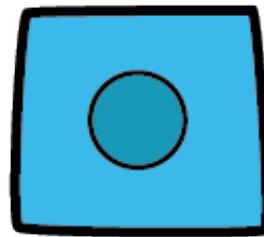
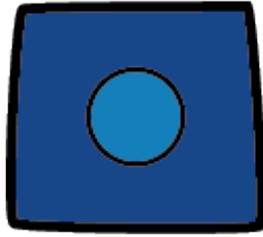
Building Counts Matrix → QC & Filtering → Normalization → Visualization → Analysis

Step 2 – Quality control and filtering



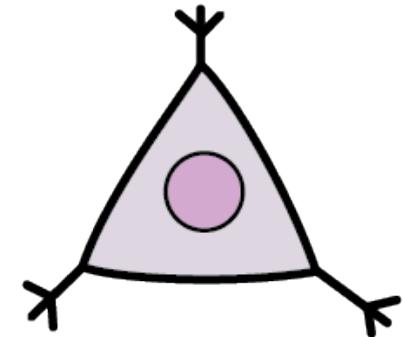
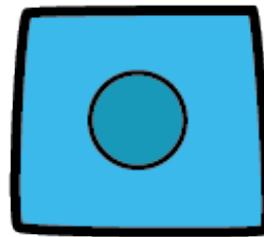
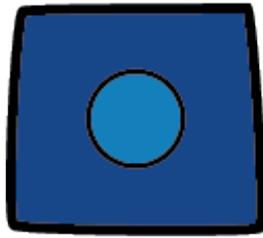
Building Counts Matrix → QC & Filtering → Normalization → Visualization → Analysis

Step 3 - Normalization



If we only have gene expression, how can we determine which cells are similar?

Step 3 - Normalization

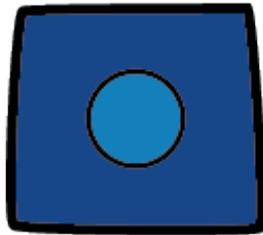


10% Capture Efficiency

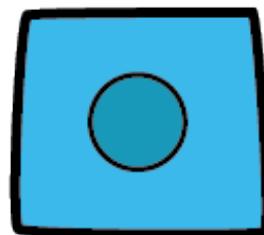
Gene	Cell A
X	10
Y	20
Z	70

Building Counts Matrix → QC & Filtering → Normalization → Visualization → Analysis

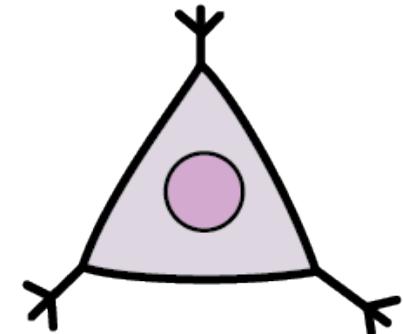
Step 3 - Normalization



10% Capture Efficiency



20% CE

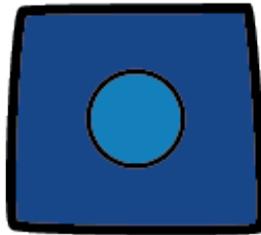


Gene	Cell A
X	10
Y	20
Z	70

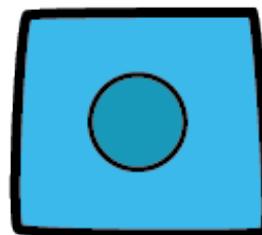
Gene	Cell B
X	20
Y	40
Z	140

Building Counts Matrix → QC & Filtering → Normalization → Visualization → Analysis

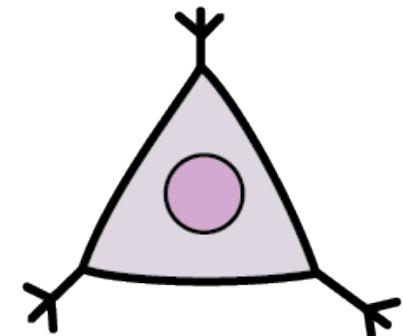
Step 3 - Normalization



10% Capture Efficiency



20% CE



20% CE

Gene	Cell A
X	10
Y	20
Z	70

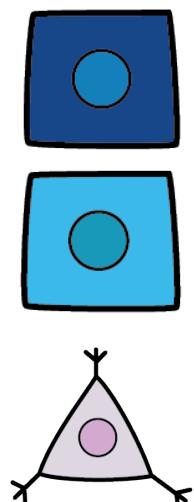
Gene	Cell B
X	20
Y	40
Z	140

Gene	Cell C
X	20
Y	0
Z	80

Building Counts Matrix → QC & Filtering → Normalization → Visualization → Analysis

Step 3 - Normalization

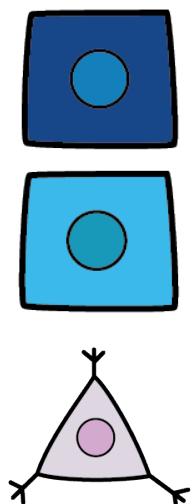
Raw counts



	X	Y	Z
A	10	20	70
B	20	40	140
C	20	0	80

Step 3 - Normalization

Raw counts

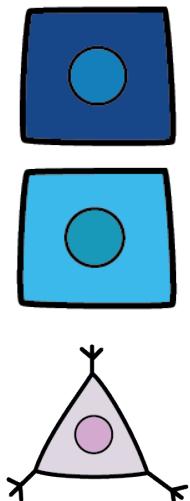


	X	Y	Z
A	10	20	70
B	20	40	140
C	20	0	80

$$\text{dist}(A,B) = \sqrt{(x_A - x_B)^2 + (y_a - y_b)^2 + (z_a - z_b)^2}$$

Step 3 - Normalization

Raw counts



	X	Y	Z
A	10	20	70
B	20	40	140
C	20	0	80

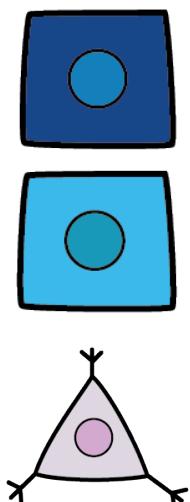
Pairwise distances

$$\text{dist}(A,B) = 71.4$$

$$\text{dist}(A,B) = \sqrt{(x_A - x_B)^2 + (y_a - y_b)^2 + (z_a - z_b)^2}$$

Step 3 - Normalization

Raw counts



	X	Y	Z
A	10	20	70
B	20	40	140
C	20	0	80

Pairwise distances

$$\text{dist}(A,B) = 71.4$$

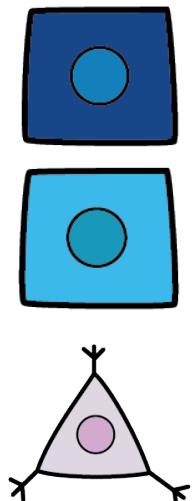
$$\text{dist}(A,C) = 24.5$$

$$\text{dist}(B,C) = 67.1$$

$$\text{dist}(A,B) = \sqrt{(x_A - x_B)^2 + (y_a - y_b)^2 + (z_a - z_b)^2}$$

Step 3 - Normalization

Raw counts

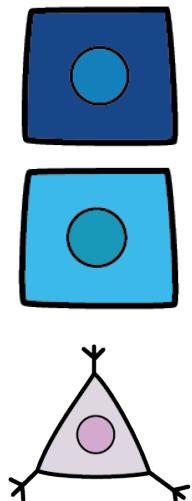


	X	Y	Z	Library Size	Pairwise distances
A	10	20	70	100	$\text{dist}(A,B) = 71.4$
B	20	40	140	200	$\text{dist}(A,C) = 24.5$
C	20	0	80	100	$\text{dist}(B,C) = 67.1$

$$\text{dist}(A,B) = \sqrt{(x_A - x_B)^2 + (y_a - y_b)^2 + (z_a - z_b)^2}$$

Step 3 - Normalization

Normalized counts

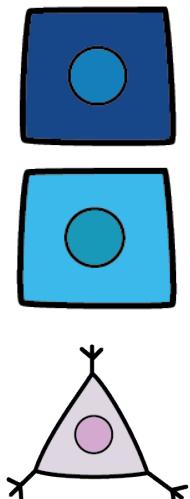


	X	Y	Z	Library Size	Pairwise distances
A	0.1	0.2	0.7	100	$\text{dist}(A,B) = 71.4$
B	0.1	0.2	0.7	200	$\text{dist}(A,C) = 24.5$
C	0.2	0	0.8	100	$\text{dist}(B,C) = 67.1$

$$\text{dist}(A,B) = \sqrt{(x_A - x_B)^2 + (y_a - y_b)^2 + (z_a - z_b)^2}$$

Step 3 - Normalization

Normalized counts

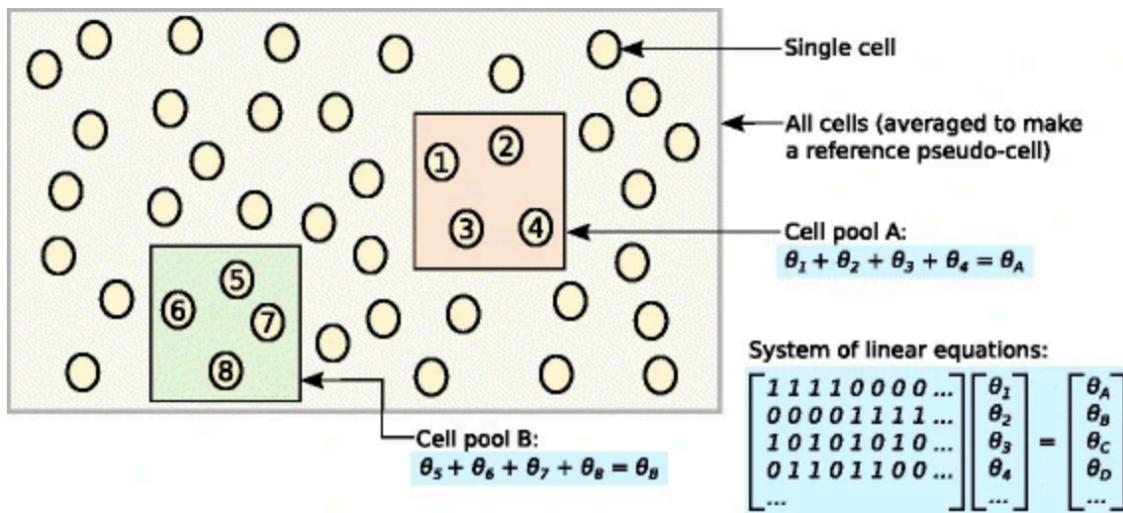


	X	Y	Z	Library Size	Pairwise distances
A	0.1	0.2	0.7	100	$\text{dist}(A,B) = 0$
B	0.1	0.2	0.7	200	$\text{dist}(A,C) = 0.25$
C	0.2	0	0.8	100	$\text{dist}(B,C) = 0.25$

$$\text{dist}(A,B) = \sqrt{(x_A - x_B)^2 + (y_a - y_b)^2 + (z_a - z_b)^2}$$

More complex normalization approaches exist

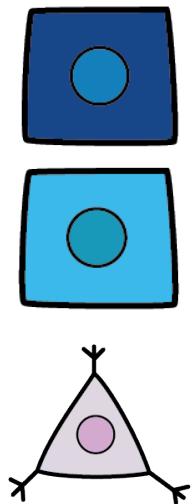
Fig. 3



Schematic of the deconvolution method. All cells in the data set are averaged to make a reference pseudo-cell. Expression values for cells in pool A are summed together and normalized against the reference to yield a pool-based size factor θ_A . This is equal to the sum of the cell-based factors θ_j for cells $j=1-4$ and can be used to formulate a linear equation. (For simplicity, the t_j term is assumed to be unity here.) Repeating this for multiple pools (e.g., pool B) leads to the construction of a linear system that can be solved to estimate θ_j for each cell j

Step 3.5 – Transformation / Scaling

Normalized counts

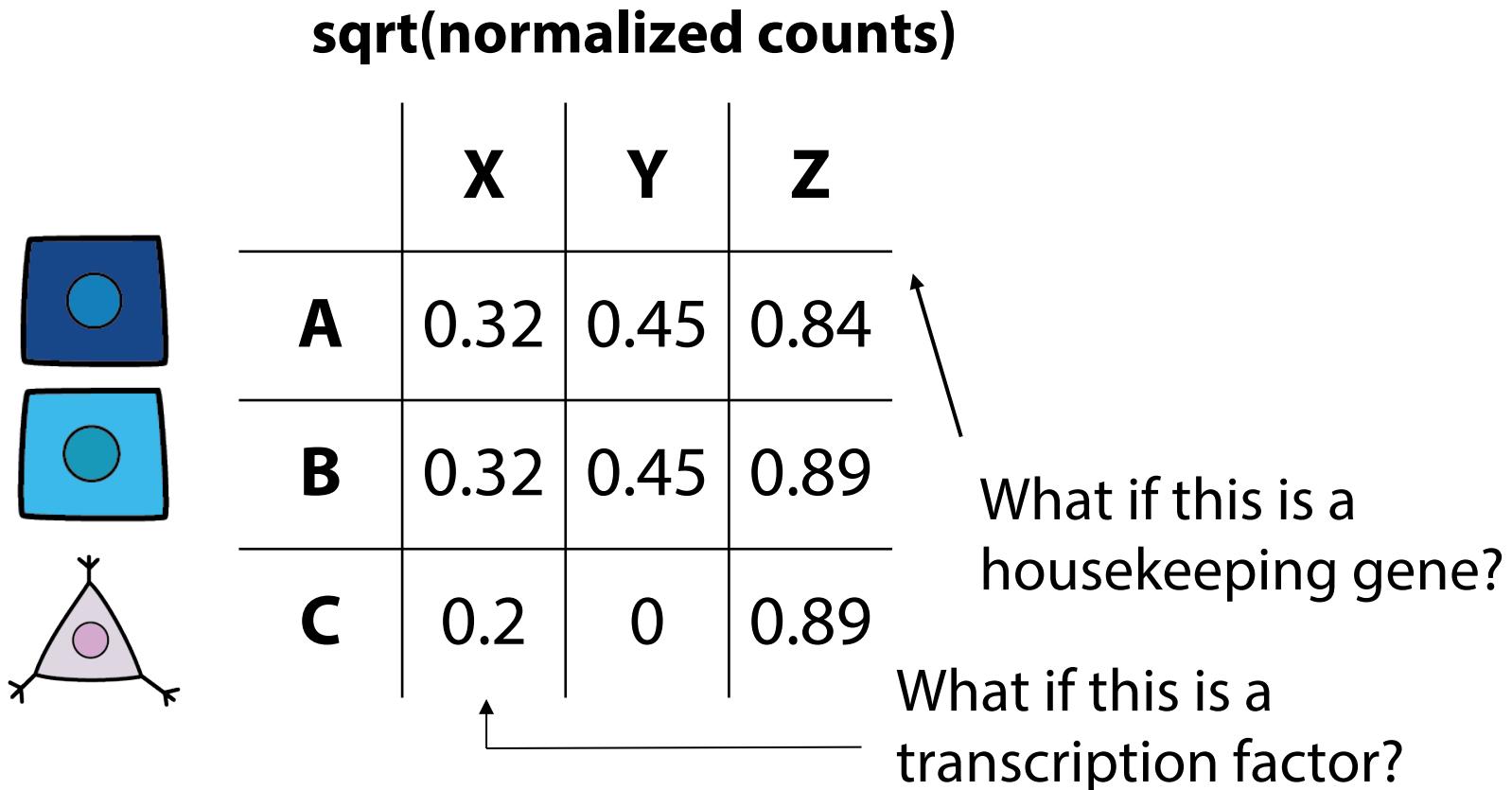


	X	Y	Z
A	0.1	0.2	0.7
B	0.1	0.2	0.7
C	0.2	0	0.8

What if this is a housekeeping gene?

What if this is a transcription factor?

Step 3.5 – Transformation / Scaling



Quick quiz

Other than square-root, what other kinds of transformations might you apply to single cell data?

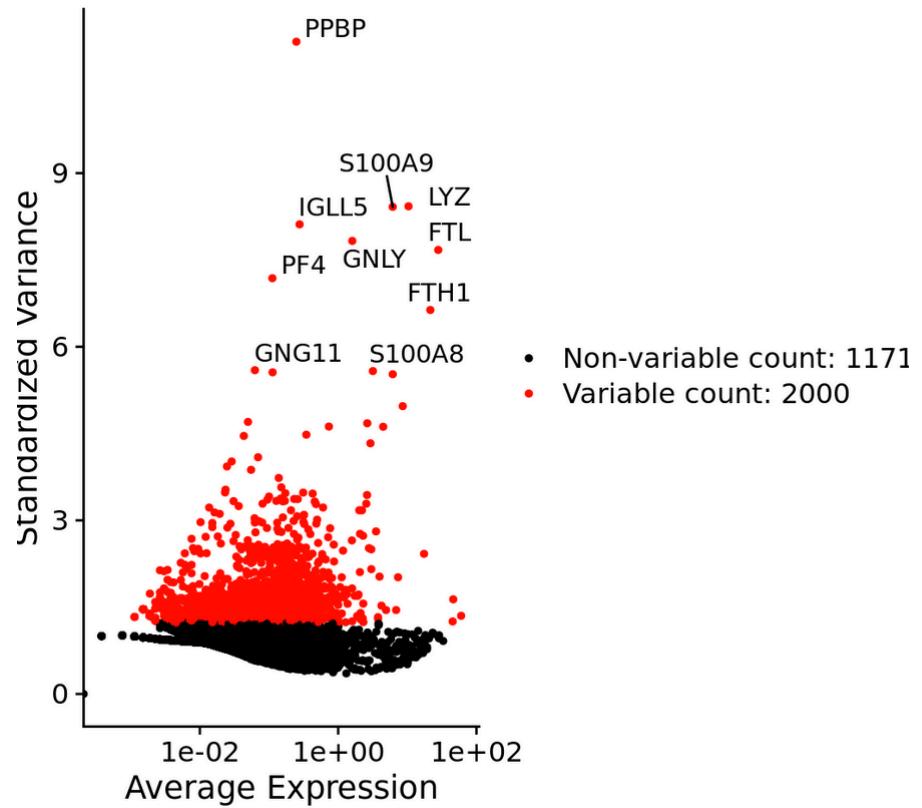
What kind of transformations, other than square-root, could we apply to single cell data?

Top

Step 5 – Dimensionality reduction and visualization

Selecting highly variable genes (HVGs):

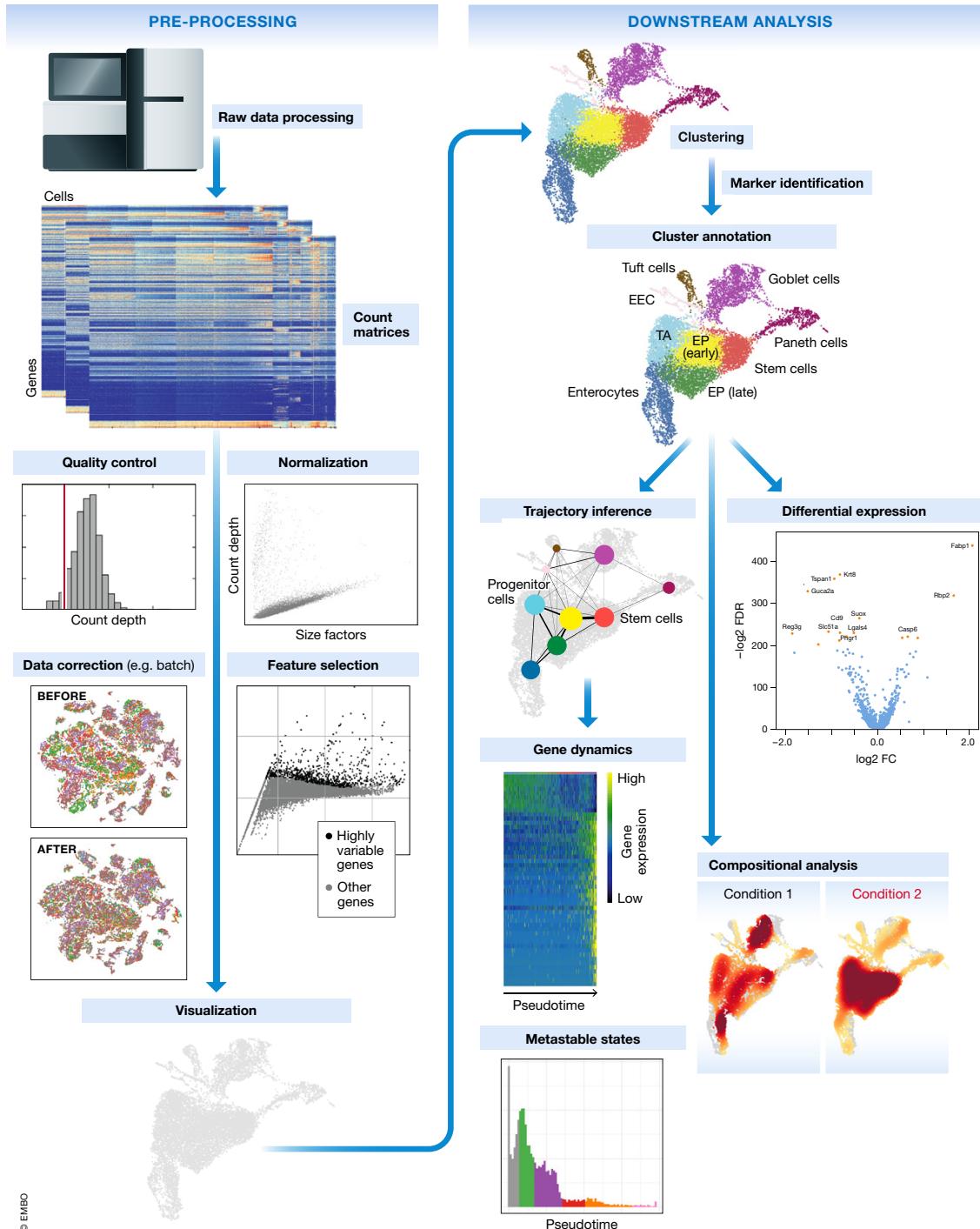
- Calculate log10 mean expression and variance
- Fit a loess curve
- Standardize variance to mean 0 std 1
- Take the top 2000 HVGs



Building Counts Matrix → QC & Filtering → Normalization → Visualization → Analysis

Standard Single-Cell RNA-seq Workflow

1. Sequencing and read mapping
2. Quality control and filtering
3. Normalization
4. Data Correction
5. Dimensionality reduction and visualization
6. Downstream analysis
 1. Clustering
 2. Trajectory inference
 3. Differential expression
7. Comparison of multiple conditions

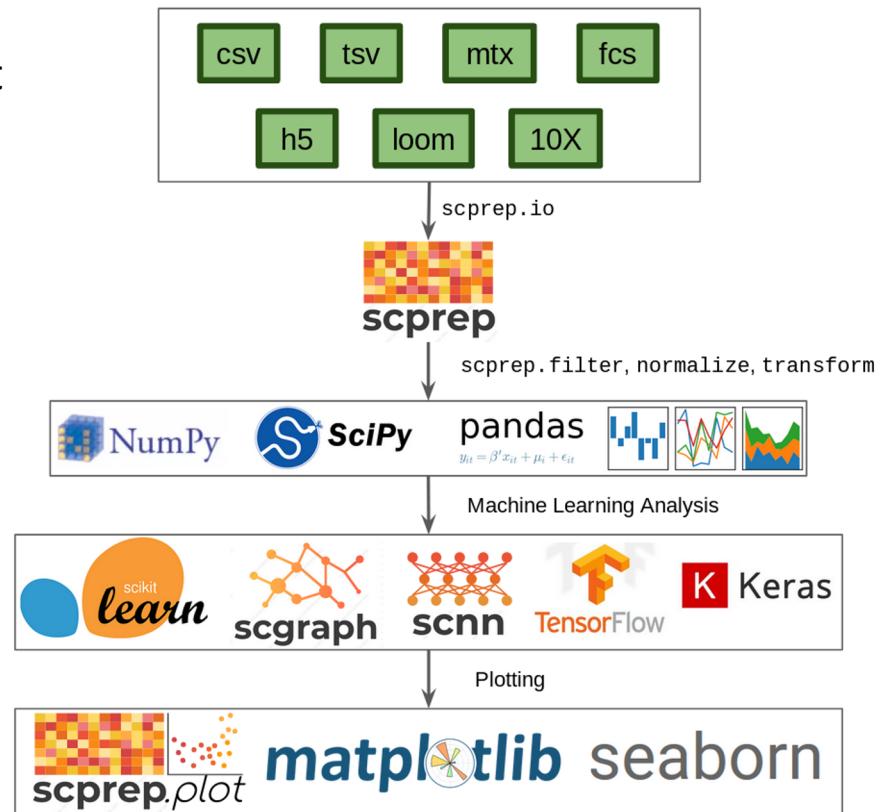


What questions do you have about today's material?

Top

Introducing scprep

- scprep: a lightweight scRNA-seq toolkit for Python Data Scientists
- Native support for ML ecosystem: NumPy, SciPy, pandas and scikit-learn
- Helper functions and wrappers for common preprocessing, analysis and visualization





Exercise!

Load, preprocess, and visualize a scRNAseq dataset generated from a time course of embryoid bodies

