

Day1 Single cell RNA-seq

NIEHS Workshop--Analyzing NGS data

July 13, 2021

Yunhua Zhu, PhD

Instructor's background

- Bachelor @ National University of Singapore(NUS) in Biochemistry 2000-03
- Ph.D @ NUS in stem cell biology | 2006 -10
- Postdoc @ Hopkins with wet lab & dry lab | 2014 - 2019
 - Neurogenesis w/t single-cell RNA-seq
 - Neurodegeneration w/t single-nucleus RNA-seq
 - Learning R, Bash script statistics through Google and Youtube
- Computational Genomics Specialist @ BCBB | 2019
 - Single-cell RNA-seq, CITE-seq
 - SMARTseq2 pipeline with FASTQC, RSEM, SENIC
 - 10X genomics pipeline with Cellranger, Seurat, Scanpy and RNA Velocity.
 - Bulk RNA-seq
 - Deconvolution Bulk RNA-seq using ABIS and CybersortX.
 - Functional annotation with clusterProfiler.
 - Epigenetics
 - ATAC-seq and Chip-seq



Agenda -- day 1, single cell RNA-seq

- 8-9AM, setting up the system
 - Get to know each other
 - Copy files to scratch and Set up pseudolink to the folders
 - Setup conda environment for installation.
 - Initiate jupyter lab.
- 9-9:30AM, Introducing single cell RNA-seq 30 minutes.
- 9:30-10AM, Getting expression matrix 9:30-10am
 - Convert bcl to fastq files using the tiny-bcl folder 15 minutes
 - Using cellranger count to do get expression matrix
- 10-11AM, Run basic pipelines
 - Using Seurat pipeline to import alignment result do QC, and normalization.
 - (optional) Using Scanpy pipeline to import alignment result do QC, and normalization.
- 11-11:30AM, use RNA velocity to infer trajectory.
 - Run Velocyto CLI to get the reads aligning to non-spliced region.
 - Import into scVelo (python) to do trajectory inference.
- 11:30—1:59AM, Prepare for day2
 - Install IGV and IGM

Outline -- day 2, bulk RNA-seq and ChIP-seq

Part I, Bulk RNA-seq

- 8-8:15AM, introduction to bulk RNA-seq
- 8:15-9AM, fastq and alignment with RSEM, Tophat and visualize the output.
- 9-9:30AM, normalization and differential analysis with DESeq2
- (optional) cellular decomposition using ABIS and Cybersort X

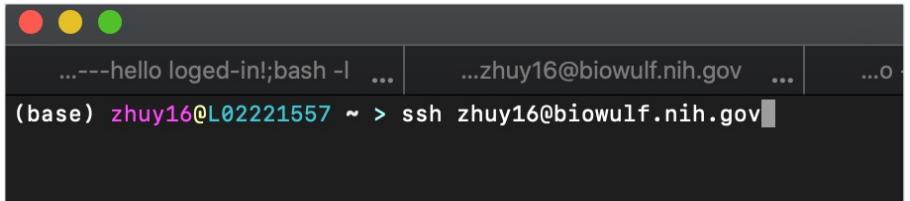
Part II, Chip-seq

- 9:30-10AM, introduction to ChIP-seq analysis
- 10-10:15AM, download fastq, alignment to genome using Bowtie2
- 10:15-10:45AM, use MACS2 to call peaks, and use Homer to call differential peaks
- 10:45-11:15AM, use CEAS to annotate peaks and summarize statistics.
- 11:15am -11:45AM, use Homer to study TF binding site analysis. And use GREAT for gene ontology.

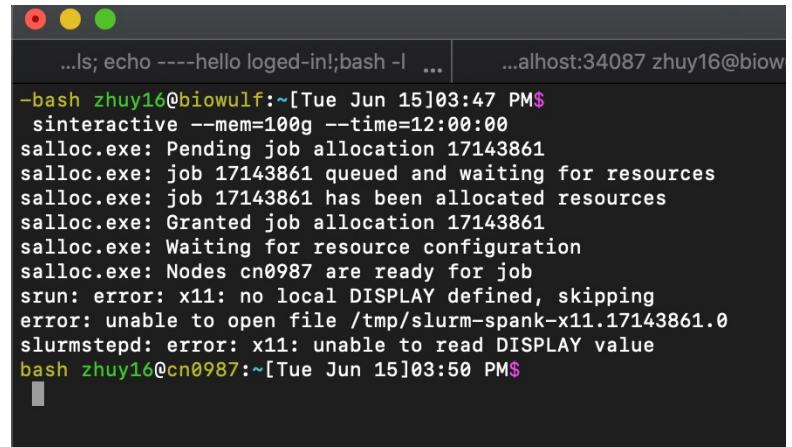
Review & Conclusion

Login to biowulf and go to the folder

- ssh user_id@biowulf.nih.gov
 - # set up sudo links to your storage folder:
 - module load tmux
 - tmux new # get a detachable terminal
 - sinteractive --mem=100g --time=12:00:00 # get a comput node
 - cd
 - ln -s /data/user_id data
 - ln -s /scratch/user_id scratch
- # copy the workshop material to your scratch folder:
- cd scratch; cp -r /spin1/users/classes/NIEHS_NGS/NIEHS_NGS_Workshop /scratch/user_id/



```
...---hello loged-in!;bash -l ... ...zhuy16@biowulf.nih.gov ... ...o ...
(base) zhuy16@L02221557 ~ > ssh zhuy16@biowulf.nih.gov
```



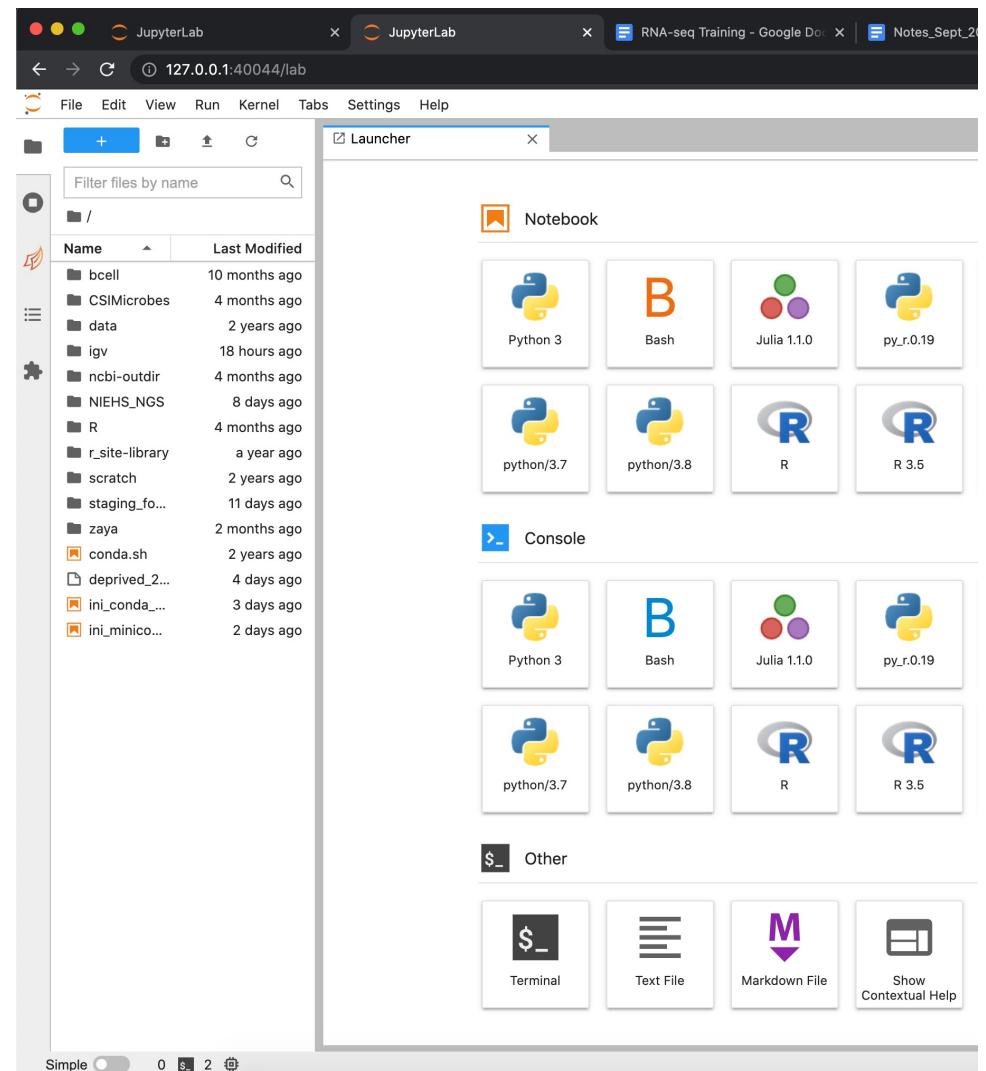
```
...ls; echo ----hello loged-in!;bash -l ... ...alhost:34087 zhuy16@biowulf ...
-bash zhuy16@biowulf:[Tue Jun 15]03:47 PM$ sinteractive --mem=100g --time=12:00:00
salloc.exe: Pending job allocation 17143861
salloc.exe: job 17143861 queued and waiting for resources
salloc.exe: job 17143861 has been allocated resources
salloc.exe: Granted job allocation 17143861
salloc.exe: Waiting for resource configuration
salloc.exe: Nodes cn0987 are ready for job
srun: error: x11: no local DISPLAY defined, skipping
error: unable to open file /tmp/slurm-spank-x11.17143861.0
slurmstepd: error: x11: unable to read DISPLAY value
bash zhuy16@cn0987:[Tue Jun 15]03:50 PM$
```

Start up the jupyter lab,
--our working interface

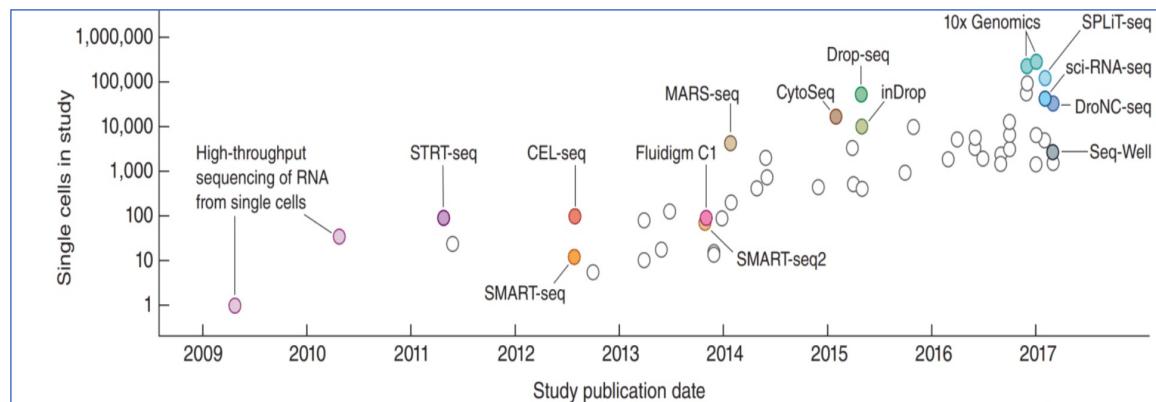
- # Open another terminal
 - ssh user_id@biowulf.nih.gov
 - Enter password
 - module load tmux
 - module load tmux; tmux new -ct 'sinteractive --mem=100g --time=12:00:00 --tunnel'
 - # Copy the tunnel script to another (3rd) terminal, execute and enter password, to establish a ssh tunnel between local computer and the work node.
 - module load jupyter R/4.0.5 && jupyter lab --ip localhost --port \$PORT1 --no-browser
 - # wait until an URL link appears, copy it to your web browser, to get connected to biowulf through a jupyter lab interface.

Jupyter lab

- Navigate files
- Control kernels
- Bash operation
- Monitor activity at biowulf

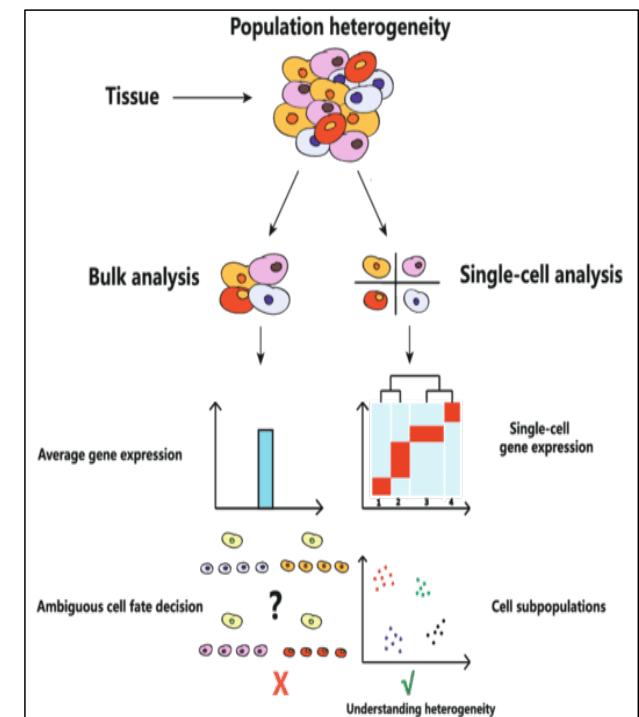


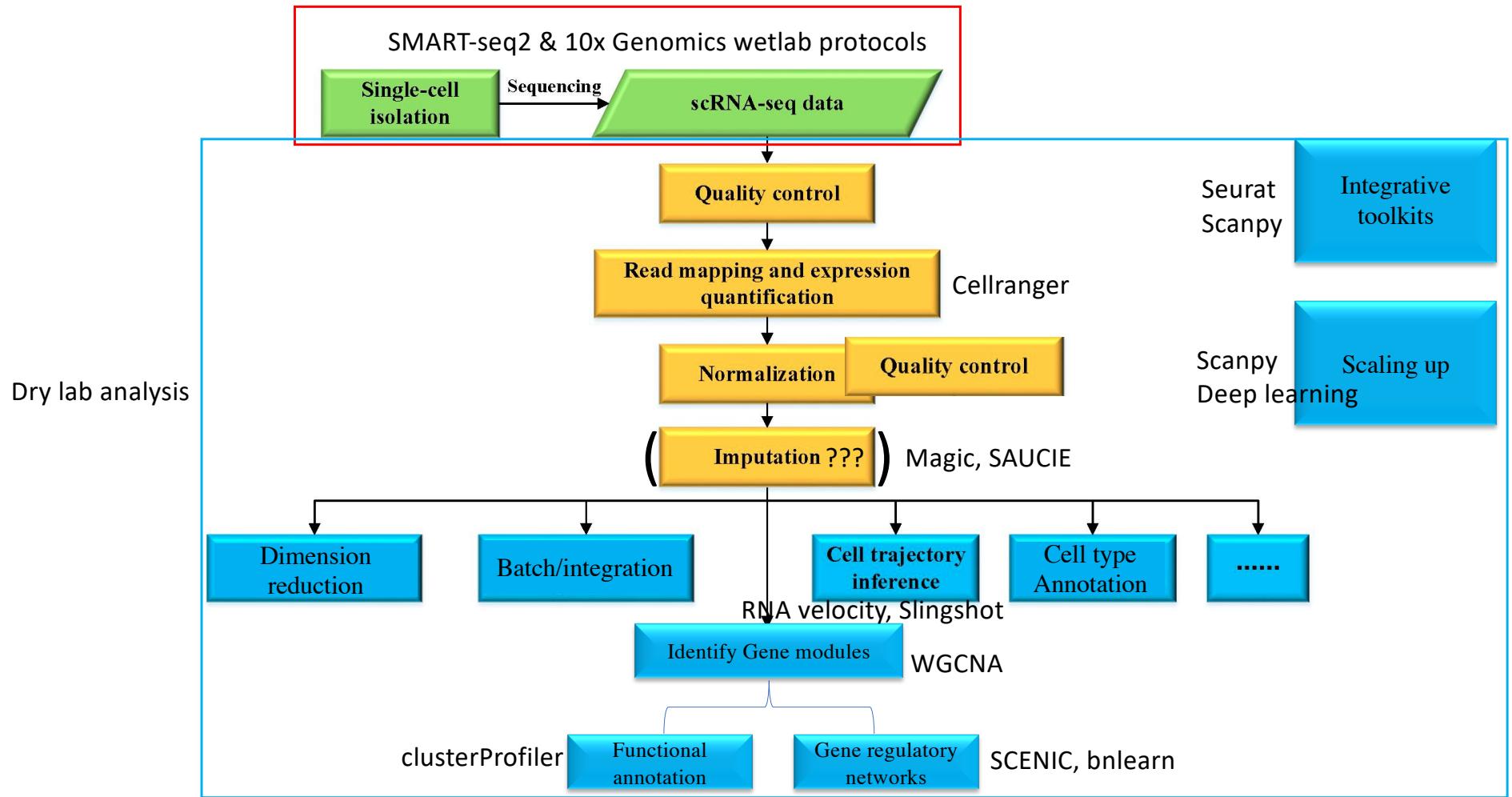
Evolution of single cell technologies



Exponential scaling of single-cell RNA-seq in the past decade.

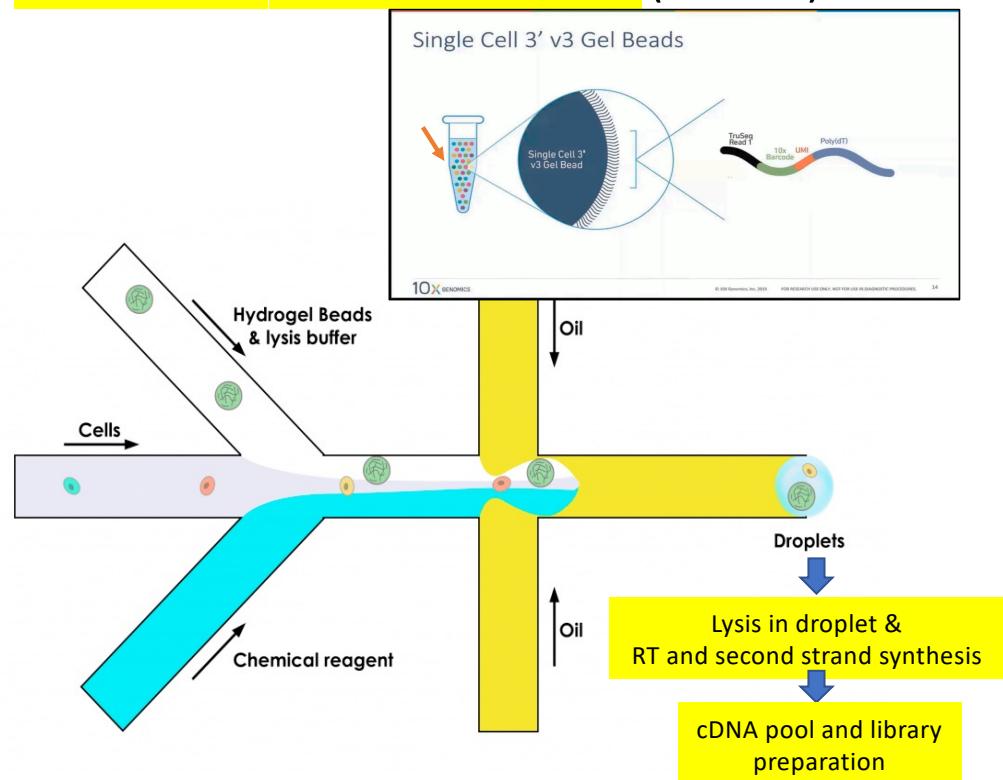
Nature protocols 13;4:599-604)





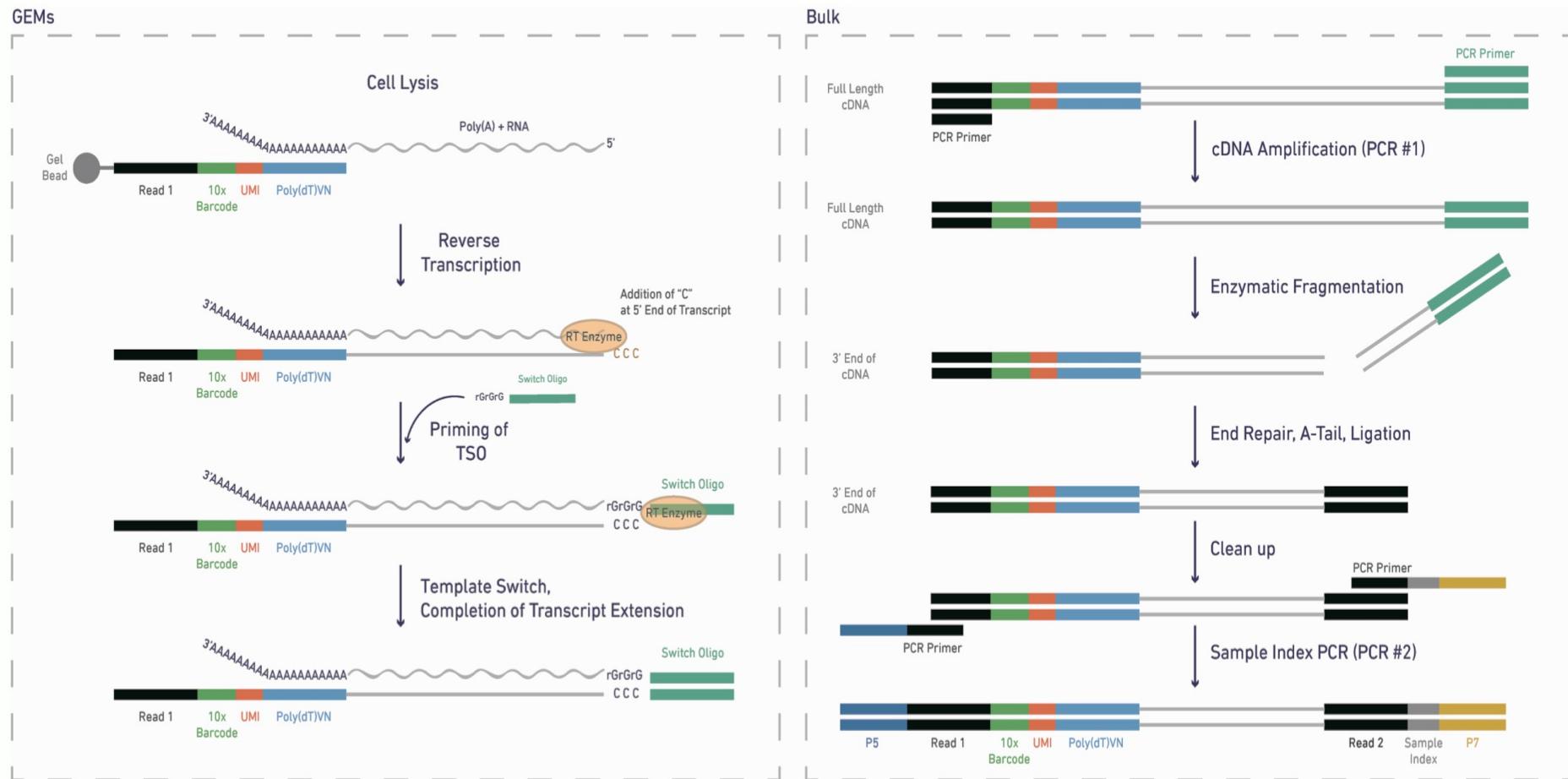
10X Genomics— commercial solution that facilitates automatic generation of Gel Bead-in-Emulsion (GEM)

- In a GEM droplet, one hydrogel bead and one cell were captured
- One **hydrogel bead** is attached with **millions of poly-T primers** with an identical unique barcode.
- cDNA and the second-strand synthesis in the droplet
- Droplets (~1nL) are disrupted to collect all the barcoded samples for highly multiplexed library preparation and sequencing
- **Standardized automation** and reagent has made sequencing library preparation very efficient



Inside individual GEMs (Gel Bead-in-Emulsion)

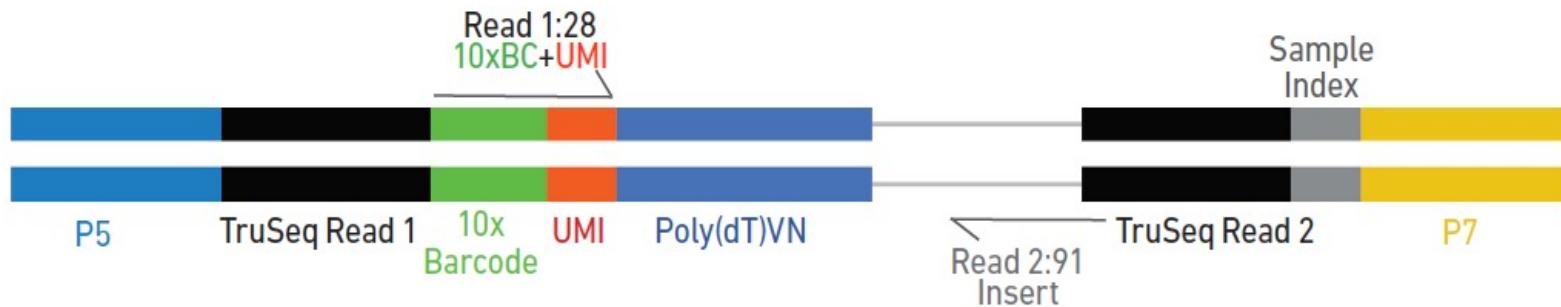
Pooled cDNA processed in bulk

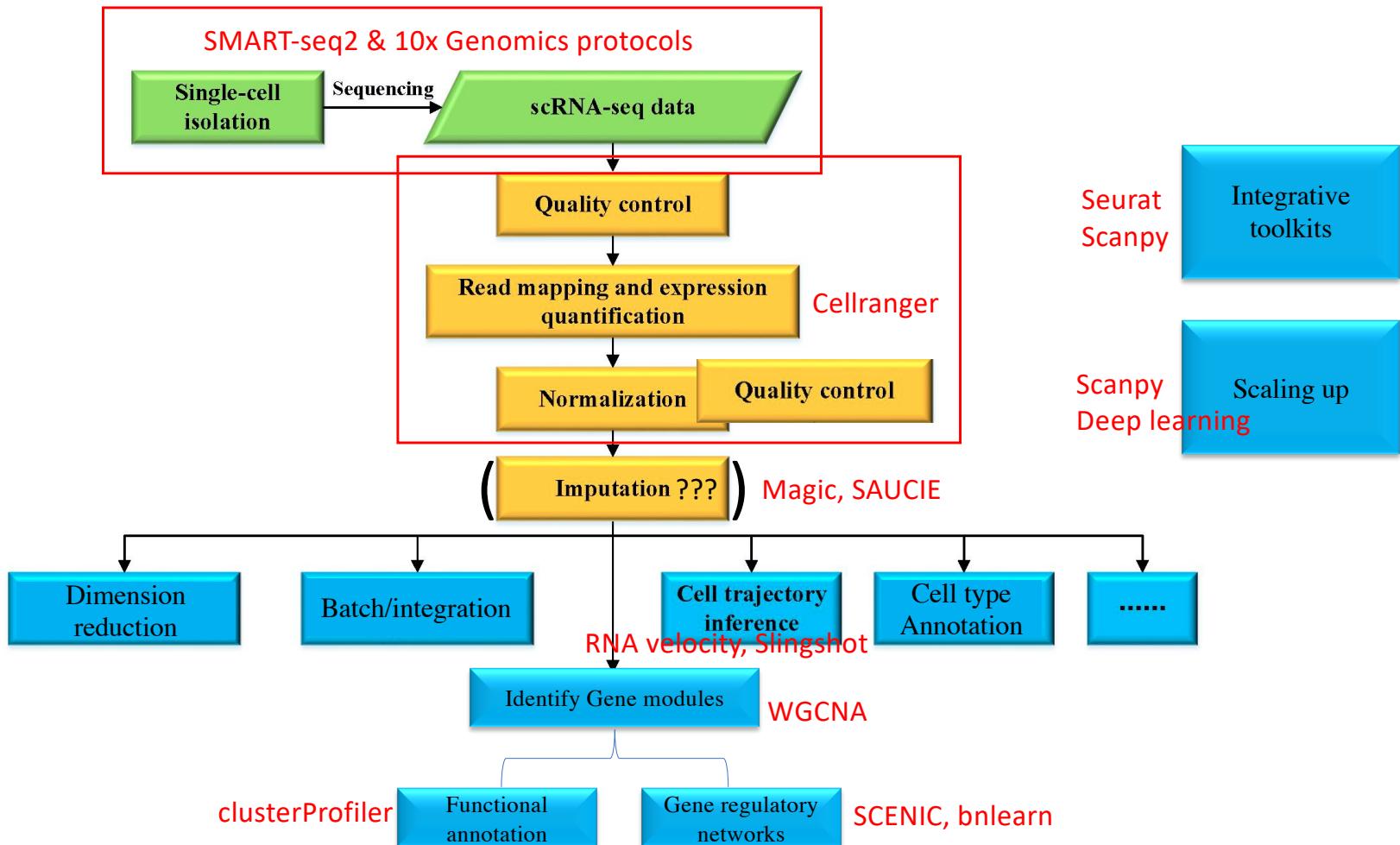


https://assets.ctfassets.net/an68im79xiti/4fly9tr6qQuCWamlloIe/40658acce7a6756e38537584897840e3/CG000108_AssayConfiguration_SC3v2.pdf

Next-seq reading the paired ends

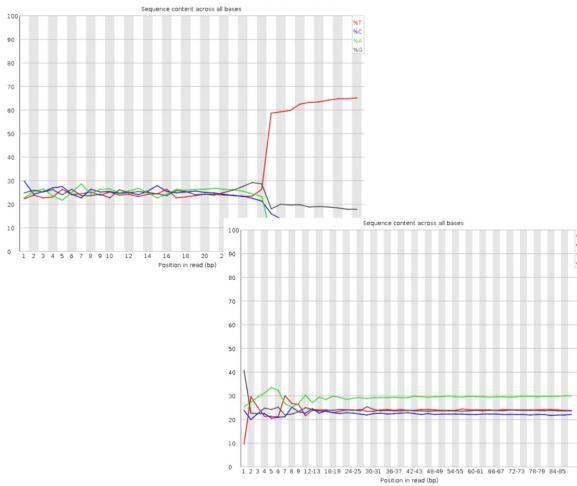
- 10XBC: 16 bp barcodes $2^{16}=65536$ possible unique cells
- UMI, $2^{12}=4096$ unique copies of mRNA can be distinguished for each gene
- Read2 will read into cDNA to identify the identity of the gene
- Sample barcode, identify the batch of your library sample
- All information will be summarized by the Cellranger software





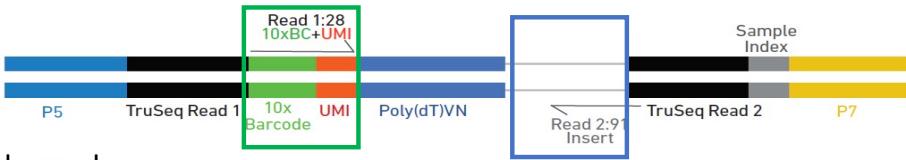
QC on sequencing results

scrALI001_S1_L001_I1_001.fastq.gz
scrALI001_S1_L001_R1_001.fastq.gz
scrALI001_S1_L001_R2_001.fastq.gz



- I1
 - Index file. All identical (or one of 4) at Babraham

- R1
 - Barcode reads
 - 16bp cell level barcode
 - 10bp UMI
- R2
 - 3' RNA-seq read



The screenshot shows a web browser window with several tabs open. The active tab is '001_convert_bcl_to_fastq.ipynb' in a JupyterLab environment. The left sidebar displays a file tree under '/.../scRNA-seq / notebooks /'. The main area contains a code editor with a terminal-like interface showing command-line operations for setting up CellRanger and preparing a samplesheet.

ref:
cellranger documentation on biowulf, <https://hpc.nih.gov/apps/cellranger.html>
davetang's blog on cellranger, <https://davetang.org/muse/2018/08/09/getting-started-with-cell-ranger/>

1. Demo on making fastq from data generated from a flowcell, a 'bcl' files

```
[22]: module load cellranger  
[-] Unloading cellranger 6.0.1  
[+] Loading cellranger 6.0.1  
These files have been downloaded on biowulf
```

```
[23]: # wget -c -N http://s3-us-west-2.amazonaws.com/10x.files/samples/cell-exp/2.1.0/pbmc8k/pbmc8k_f  
# wget -O cellranger-2.2.0.tar.gz "http://cf.10xgenomics.com/releases/cell-exp/cellranger-2.2.0  
# wget http://cf.10xgenomics.com/supp/cell-exp/refdata-cellranger-GRCh38-1.2.0.tar.gz
```

```
[24]: cp ${CELLRANGER_TEST_DATA:-none}/cellranger-tiny-bcl-1.2.0.tar.gz ..../data/  
cp ${CELLRANGER_TEST_DATA:-none}/cellranger-tiny-bcl-samplesheet-1.2.0.csv ..../data/  
cp: cannot create regular file '../data/': Not a directory  
cp: cannot create regular file '../data/': Not a directory  
: 1
```

```
[4]: cat ..../data/cellranger-tiny-bcl-samplesheet-1.2.0.csv
```

[Header],,,,,,,
IEMFileVersion,4,.....
Investigator Name,rir,.....

Alignment results and Quality Controls

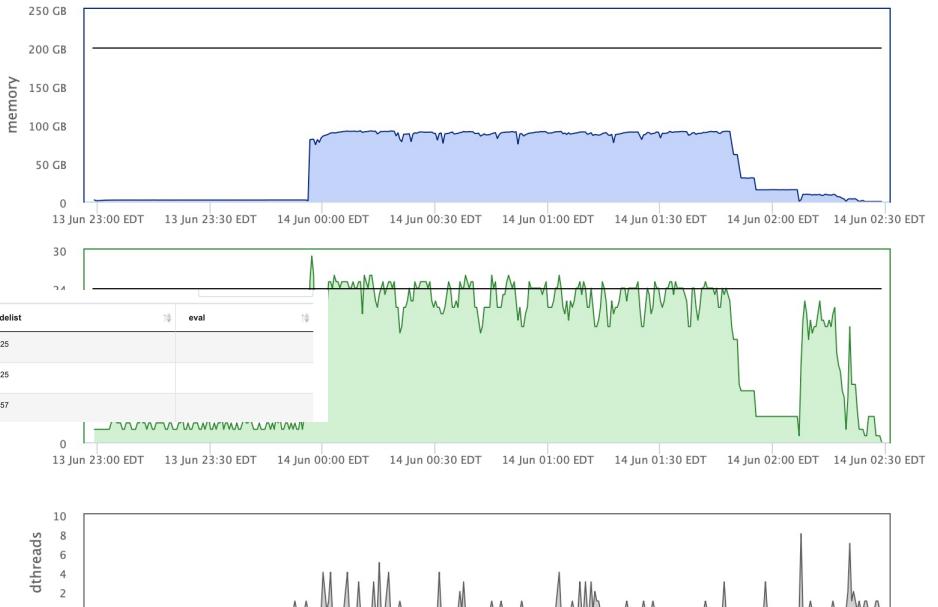
- Cellranger
 - cellranger mkfastq
 - cellranger count
 - cellranger aggr
- Quality controls
 - FASTQC on fastq files
 - Number of cells per experiment
 - Number of UMI per cell
 - Number of genes per cell
 - Percentage of mitochondrial reads
 - Removal of doublets/aggregates

Monitoring jobs on biowulf

Login to the user dashboard to monitor
<https://hpc.nih.gov/dashboard/>

jobid	jobname	state	statetime	nodefile
17006700	cellranger_pbmc3.sh	RUNNING	2021-06-13 22:58:42 EDT	cn0925
17006693	cellranger_pbmc.sh	FAILED	2021-06-13 22:57:30 EDT	cn0925
17006574	sinteractive	RUNNING	2021-06-13 22:45:00 EDT	cn1057

Biowulf Job 17006700



[Export to PNG](#)

Sbatch jobs can be listed on terminal using
sjobs -u user_id

```
sjobs -u zhuy16
User   JobId   JobName   Part    St  Reason  Runtime  Walltime  Nodes  CPUs  Memory  Dependency
===== ====== ====== ====== = = ====== ====== ====== ====== ====== ====== ====== ======
zhuy16 17006003 cellranger norm      R  1:39:28  2:00:00   1    16  200 GB
zhuy16 17006574 sinteracti interactive R  14:30    8:00:00   1    2  100 GB
zhuy16 17006700 cellranger norm      R  0:48     12:00:00  1    24  200 GB
===== ====== ====== ====== = = ====== ====== ====== ====== ====== ====== ====== ======
```

jobname	cellranger_pbmc3.sh
user	zhuy16
submitted	2021-06-13 22:58:41 EDT
state	COMPLETED
submission script	/data/classes/NIEHS_NGS/RNA-workshop-2021/scRNA-seq/data/cellranger_pbmc3.sh
work directory	/data/classes/NIEHS_NGS/RNA-workshop-2021/scRNA-seq/data

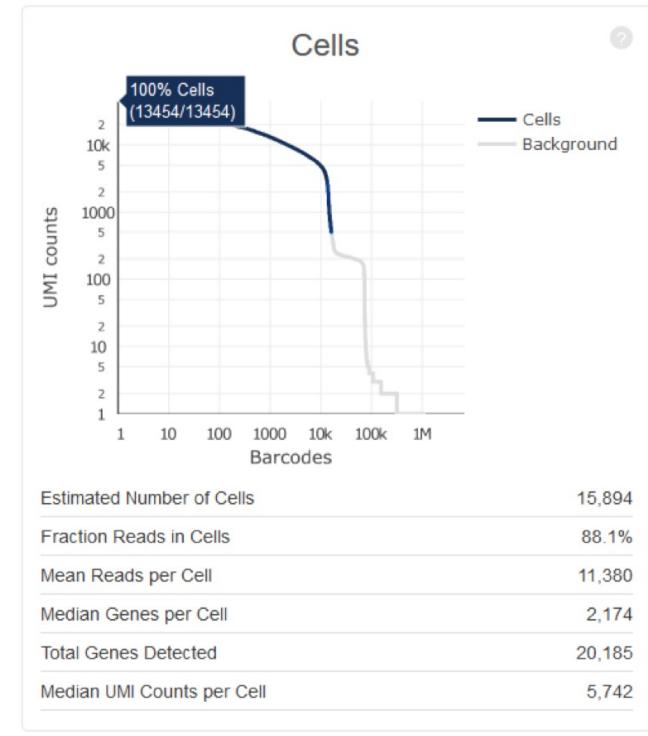
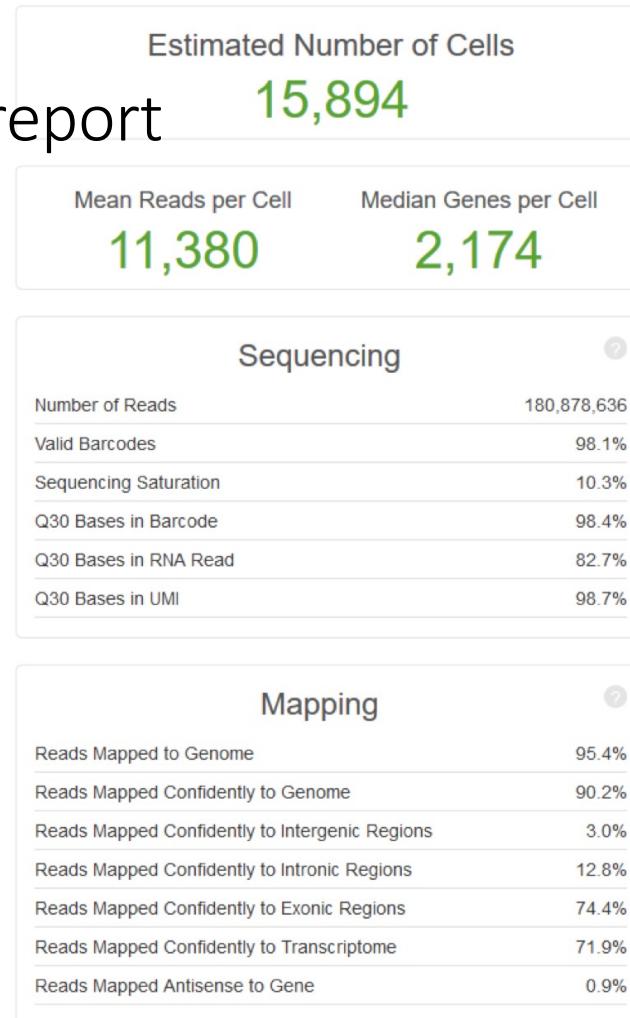
Cellranger count report

FATSTQ

I1.FASTQ
R1.FASTQ
R2.FASTQ

Cellranger →

report summary ...

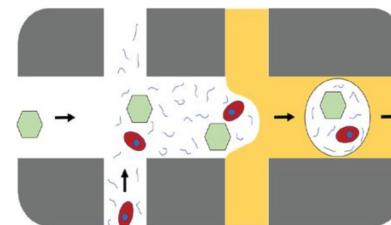


Sample

Name	embryoid_d4
Description	
Transcriptome	mm10
Chemistry	Single Cell 3' v3
Cell Ranger Version	3.0.2

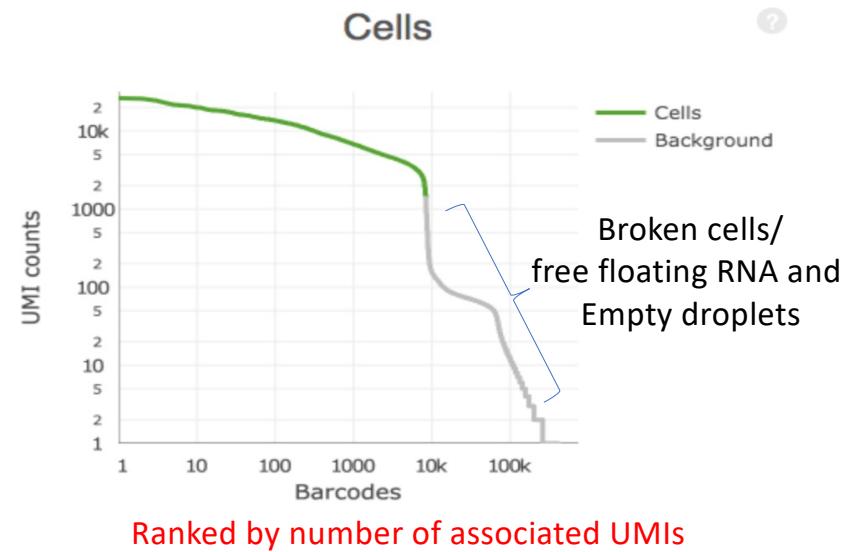
cellranger count output

- Reads level
 - CellRanger
 - cellranger mkfastq
 - Generate fastq files from image “.bcl” files
 - **cellranger count → sparse matrix**
 - umi, unique molecule identifier
 - cellranger aggr
 - Combine count data from multiple batches
 - (For CITE-seq and HASH-tag)
 - Cite-seq-count



Output:

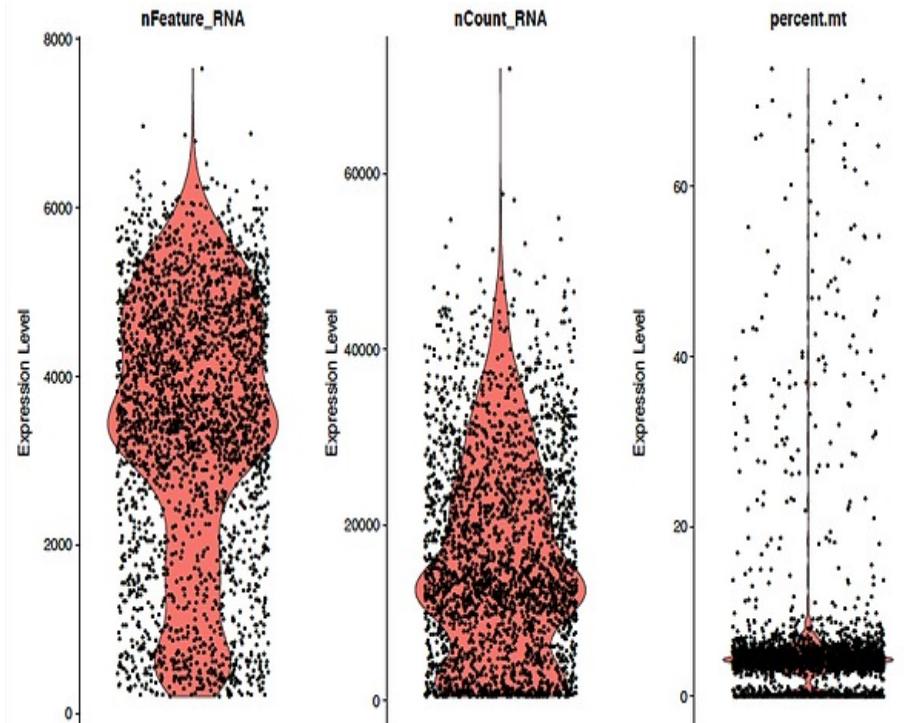
```
$ cd /home/jdoe/runs/sample345/outs
$ tree filtered_feature_bc_matrix
filtered_feature_bc_matrix
├── barcodes.tsv.gz      --cells
└── features.tsv.gz     --genes
    └── matrix.mtx.gz     --sparse matrix
0 directories, 3 files
```



<https://davetang.org/muse/2018/08/09/getting-started-with-cell-ranger/>

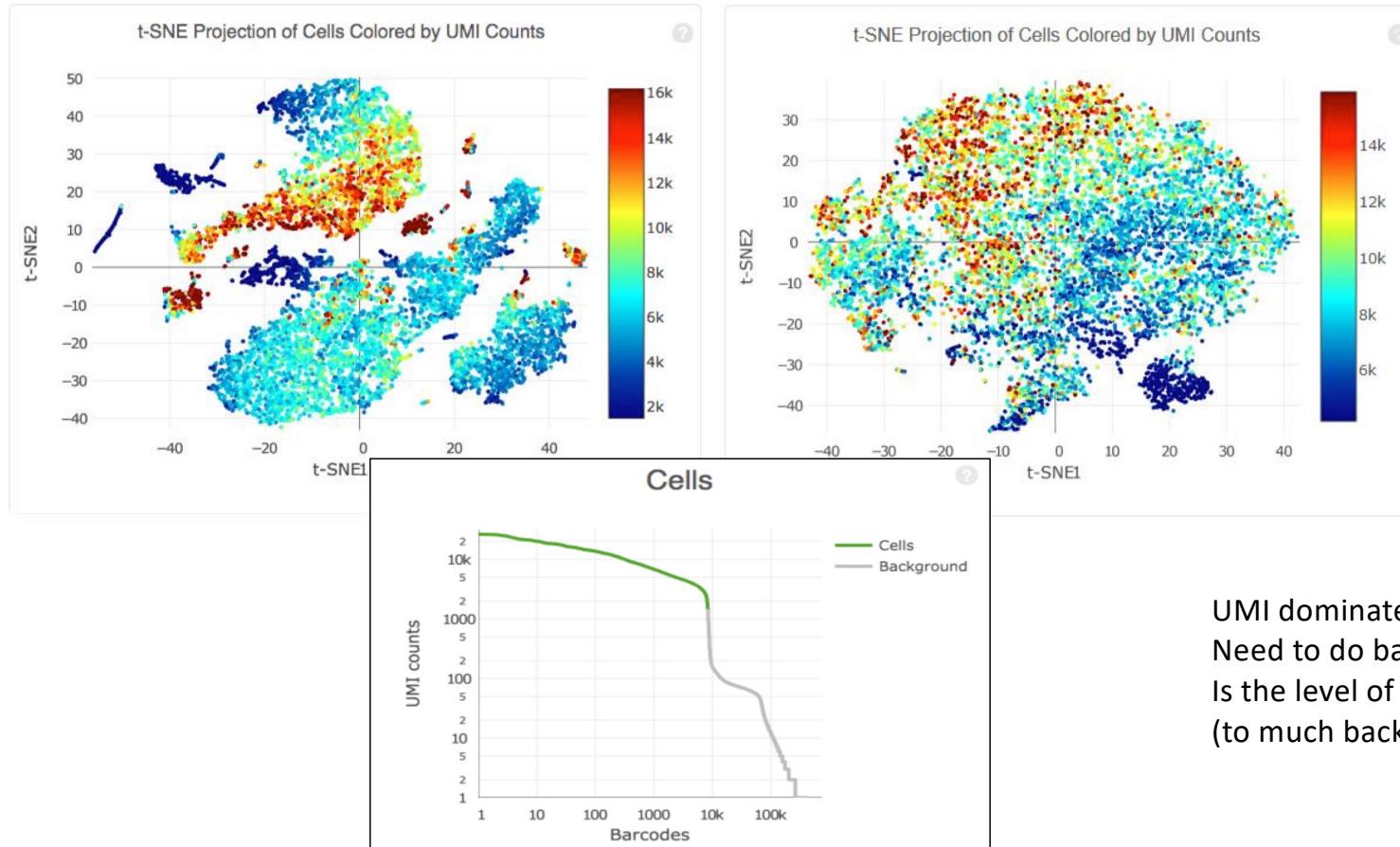
Seurat -- Quality Controls

- Quality of reads – FASTQC
- Percentage of mitochondrial reads
 - Too many mitochondria reads may indicate that cells are dying/dead/broken
 - Case-by-case, the range may vary with method of library prep methods/cell type
- How many cells are you capturing?
 - Typically few thousands in each 10X run
- The sequencing depth
 - Are they acceptable in the field (minimal 2,000 reads /cell?) determined by the Cellranger
- Alignment to the genome and exons
 - Should be 90-100% to the genome
 - A reasonably narrow range 70-80% to the exons
 - (Could be 30% to exons if you use nucleus, which contain lots of introns)
- Expected markers expressed?
 - Highly expressed genes, cell type markers, automatic detection such as scMCA etc
- Be prepared to see differences between RNA (because of the depth and dropouts) and proteins.
- Confounding factors?
 - Batch effect?
 - Is your dimension reduction capturing biological or technical variations?
 - Can be evaluated by WGCNA and visualization in PCA or tSNE



<https://www.biostars.org/p/377422/>

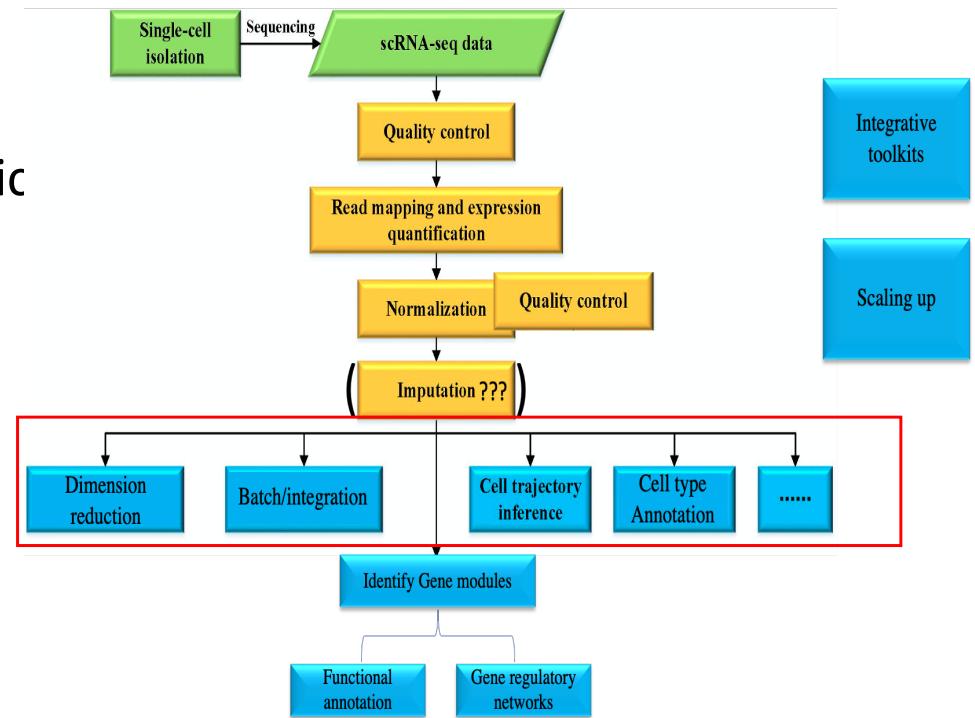
Is coverage variation affecting your data?



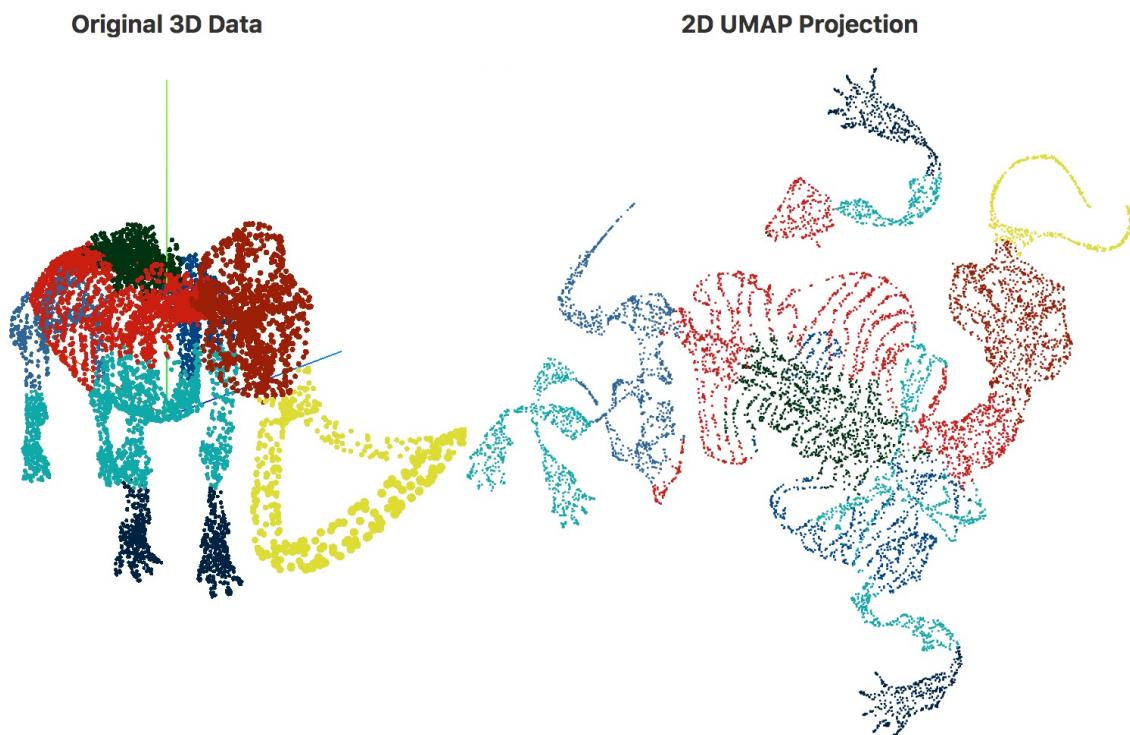
UMI dominate the variance?
Need to do batch correction?
Is the level of separation enough/expected
(to much background RNA?)

Cell-based analysis

- Clustering and annotate the biologic identities of clusters
- Inference of trajectory
- Batch correction/data integration.
- comprehensive workflow/toolkits



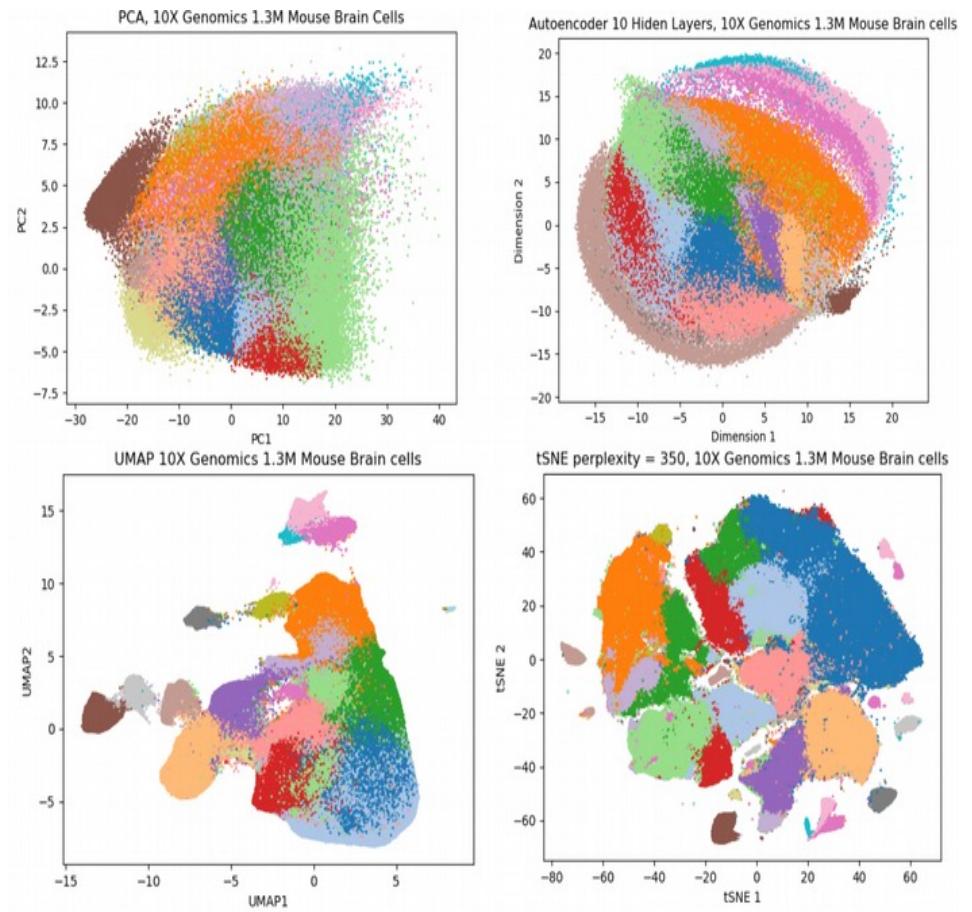
Dimension reduction--an intuitive illustration



Screenshot by Mariam from
<https://pair-code.github.io/understanding-umap/>

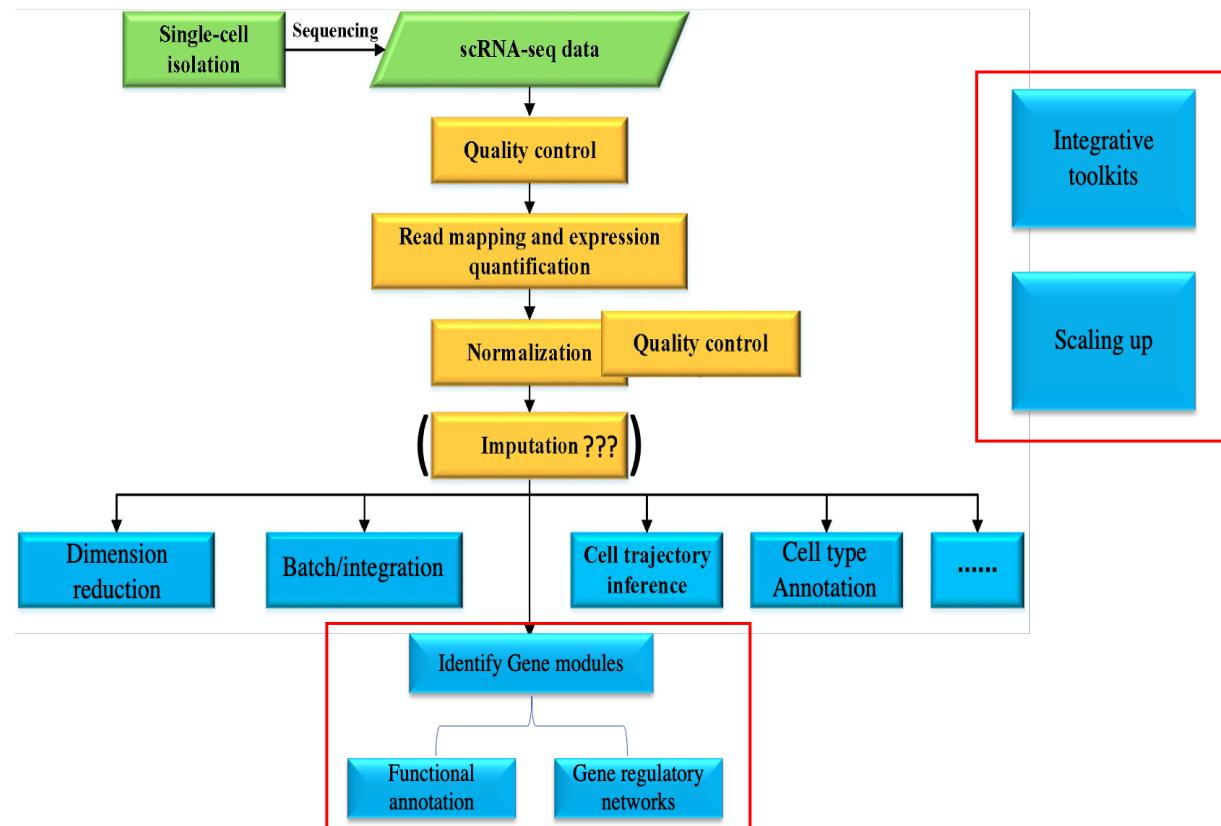
Dimension reduction

- Dimension reduction
 - PCA
 - Linear reduction
 - Based on distances
 - 2D structure in PCA depends on certain observed dominant variations
 - Often not sufficient for large number of cells
 - tSNE
 - Non-linear reduction
 - Attention to **local similarity**
 - Global shape is less meaningful
 - Add new data changes the whole pattern
 - UMAP
 - Consider both global and local structure
 - **Learnt embeddings** can be saved for new batch of data
 - AutoEncoder
 - Fast algorithm to handle up to millions of cells



<https://towardsdatascience.com/deep-learning-for-single-cell-biology-935d45064438>

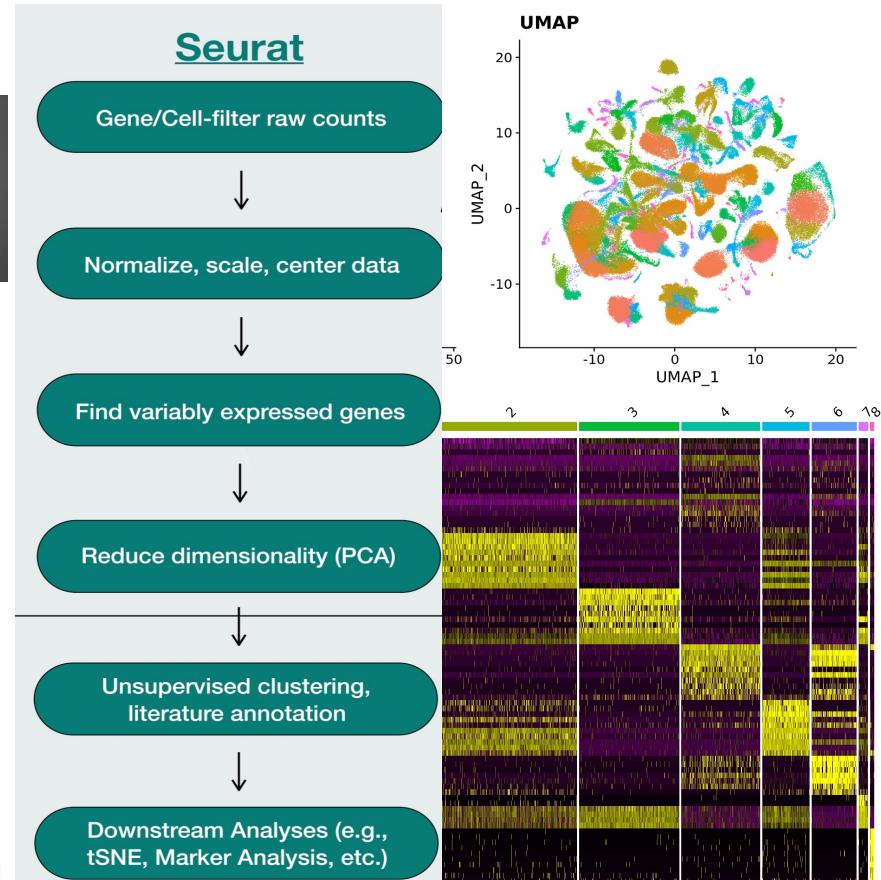
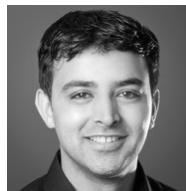
Comprehensive toolkits



Comprehensive pipeline tools for explorative analysis -- Seurat

- **Seurat pipeline**

- General QC assessment
- Cell type annotation
- Batch correction and meta analysis
- Multimodal analysis (for CITE-seq, Hash-tagging, ATAC-seq)
- Comparative analysis across different conditions

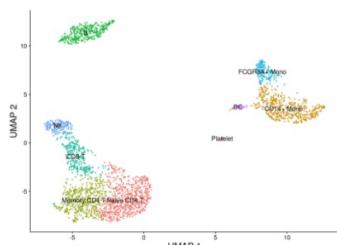


<https://satijalab.org/seurat/vignettes.html>

Multiple vignettes for different tasks

Basic pipeline: QC,
Dimension reduction
Clustering
Marker identifications

Guided tutorial – 2,700 PBMCs



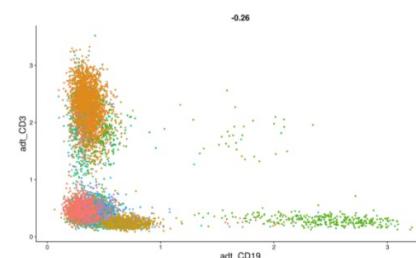
A basic overview of Seurat that includes an introduction to common analytical workflows.

[Start from here](#)

GO

Integrating with CITE-seq, HASH-seq

Multimodal analysis

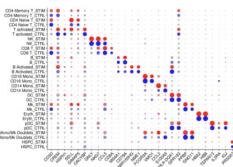


An introduction to working with multimodal datasets in Seurat.

GO

Multiple pipelines for integrating data

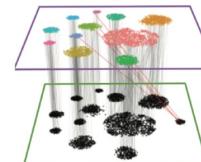
Introduction to scRNA-seq integration



An introduction to integrating scRNA-seq datasets in order to identify and compare shared cell types across experiments

GO

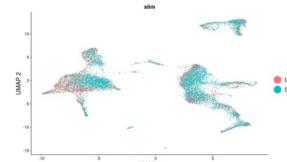
Mapping and annotating query datasets



Learn how to map a query scRNA-seq dataset onto a reference in order to automate the annotation and visualization of query cells

GO

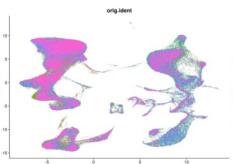
Fast integration using reciprocal PCA (rPCA)



Identify anchors using the reciprocal PCA (rPCA) workflow, which performs a faster and more conservative integration

GO

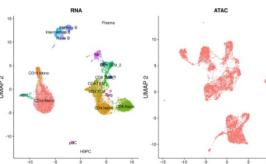
Tips for integrating large datasets



Tips and examples for integrating very large scRNA-seq datasets (including >200,000 cells)

GO

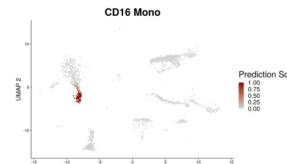
Integrating scRNA-seq and scATAC-seq data



Annotate, visualize, and interpret an scATAC-seq experiment using scRNA-seq data from the same biological system

GO

Multimodal Reference Mapping



Analyze query data in the context of multimodal reference atlases.

GO

Download vignette/tutorial to use on your data

- Linked to github for you to download the code.
- Follow through the tutorial using sample data included in the package.
- Change the input file to your our data to use the analytic pipelines.

The screenshot shows the Seurat 4.0.0 website. The top navigation bar includes links for 'Install', 'Get started', 'Vignettes', 'Extensions', 'FAQ', 'News', 'Reference', and 'Archive'. A home icon is also present. The main content area features a title 'Seurat - Guided Clustering Tutorial'. Below the title, a red circle highlights the text 'Compiled: February 08, 2021' and 'Source: vignettes/pbmc3k_tutorial.Rmd'. To the right, it says '(.Rmd can be converted to ipynb To be handled in Jupyter Notebooks)'. On the far right, there is a 'Contents' sidebar with a list of topics: 'Setup the Seurat Object', 'Standard pre-processing workflow', 'Normalizing the data', 'Identification of highly variable features (feature selection)', 'Scaling the data', 'Perform linear dimensional reduction', 'Determine the 'dimensionality' of the dataset', 'Cluster the cells', 'Run non-linear dimensional reduction (UMAP/TSNE)', 'Finding differentially expressed features (cluster biomarkers)', and 'Assigning cell type identity to clusters'.

Seurat - Guided Clustering Tutorial

Compiled: February 08, 2021
Source: vignettes/pbmc3k_tutorial.Rmd

(.Rmd can be converted to ipynb
To be handled in Jupyter Notebooks)

Setup the Seurat Object

For this tutorial, we will be analyzing the a dataset of Peripheral Blood Mononuclear Cells (PBMC) freely available from 10X Genomics. There are 2,700 single cells that were sequenced on the Illumina NextSeq 500. The raw data can be found [here](#).

We start by reading in the data. The `Read10X()` function reads in the output of the `cellranger` pipeline from 10X, returning a unique molecular identified (UMI) count matrix. The values in this matrix represent the number of molecules for each feature (i.e. gene; row) that are detected in each cell (column).

We next use the count matrix to create a `Seurat` object. The object serves as a container that contains both data (like the count matrix) and analysis (like PCA, or clustering results) for a single-cell dataset. For a technical discussion of the `Seurat` object structure, check out our [GitHub Wiki](#). For example, the count matrix is stored in `pbmc[["RNA"]].@counts`.

```
library(dplyr)
library(Seurat)
library(patchwork)

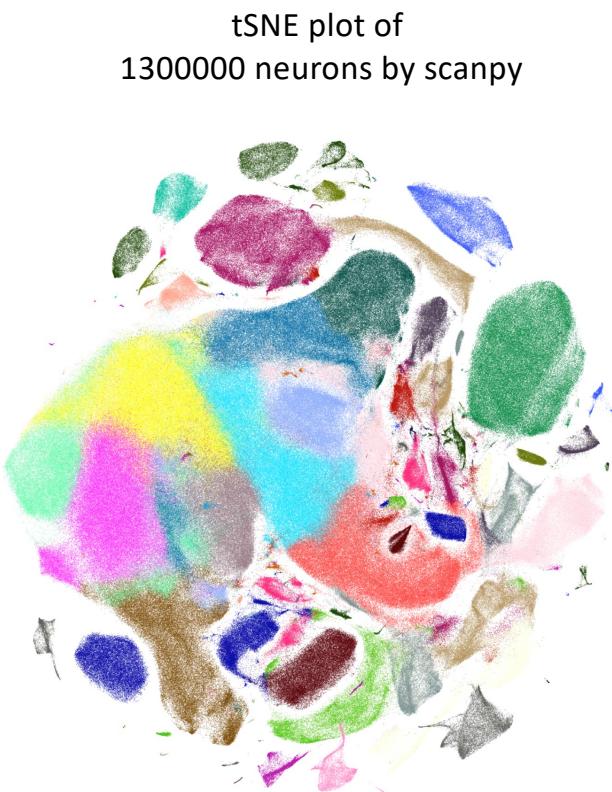
# Load the PBMC dataset
pbmc.data <- Read10X(data.dir = "./data/pbmc3k/filtered_gene_bc_matrices/hg19/")
# Initialize the Seurat object with the raw (non-normalized data).
pbmc <- CreateSeuratObject(counts = pbmc.data, project = "pbmc3k", min.cells = 3, min.features = 20
pbmc
```

Scale up with python implementations

- Python packages/toolkits are increasingly popular
 - scanpy pipeline
 - scVelo pipeline
- Some has a R rapper.



- Use python in R through Reticulate
- Use R in python through rpy2



Scanpy vs. Seurat

Satija et al., Nat. Biotechn. (2015)

Scanpy is benchmarked with Seurat.

- preprocessing: <1 s vs. 14 s
- regressing out unwanted sources of variation: 6 s vs. 129 s
- PCA: <1 s vs. 45 s
- clustering: 1.3 s vs. 65 s
- tSNE: 6 s vs. 96 s
- marker genes (approximation): 0.8 s vs. 96 s

<https://scanpy.readthedocs.io/en/stable/tutorials.html>

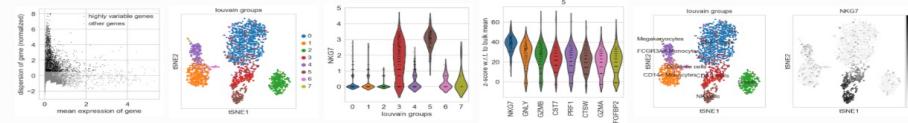
The screenshot shows the Scanpy documentation homepage. At the top is the Scanpy logo with a stylized orange and pink cell icon. Below it is the word "stable". A search bar says "Search docs". On the left, a sidebar titled "Tutorials" lists "Clustering", "Visualization", "Trajectory inference", "Integrating datasets", and "Spatial data". Under "Further Tutorials", there are links for "Usage Principles", "Installation", "API", "External API", "Ecosystem", "Release notes", "News", "Contributing", "Contributors", and "References". At the bottom, it says "Read the Docs" and "v: stable ▾".

» Tutorials

Tutorials

Clustering

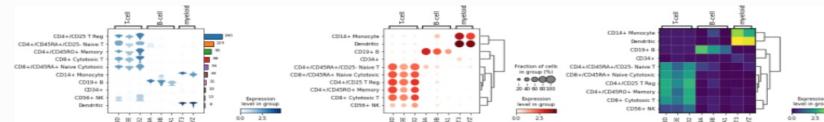
For getting started, we recommend Scanpy's reimplementation [→ tutorial: pbmc3k](#) of Seurat's [\[Satija15\]](#) clustering tutorial for 3k PBMCs containing preprocessing, clustering and the identification of cell types via known marker genes.



(.ipynb format can be converted to .RMD to be run in Rstudio)

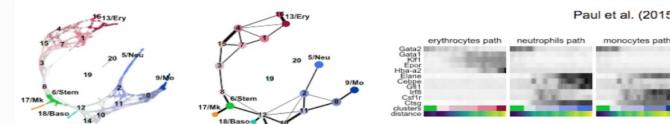
Visualization

This tutorial shows how to visually explore genes using scanpy. [→ tutorial: plotting/core](#)



Trajectory inference

Get started with the following example for hematopoiesis for data of [\[Paul15\]](#): [→ tutorial: paga-paul15](#)

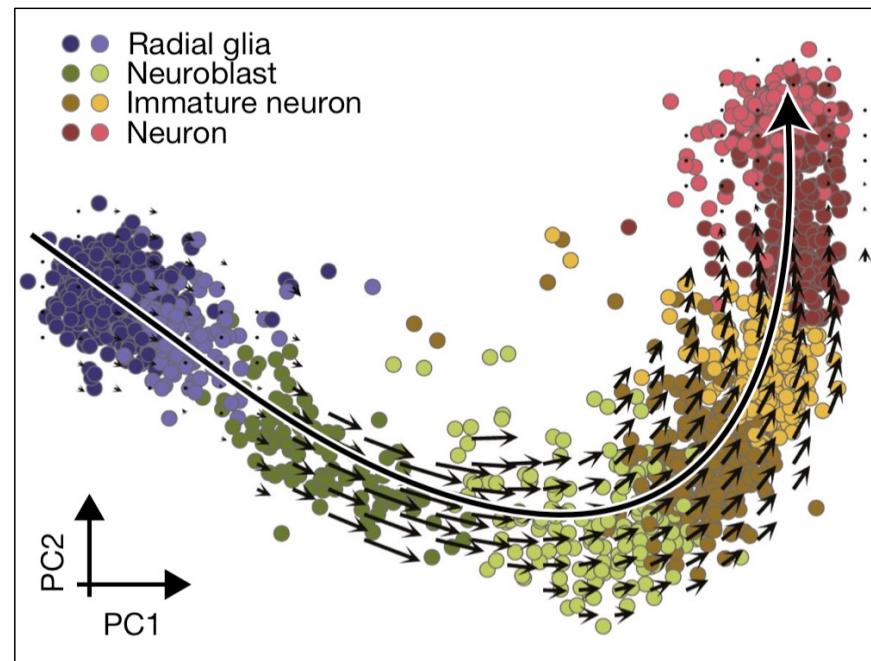


More examples for trajectory inference on complex datasets can be found in the [PAGA](#) repository [\[Wolf19\]](#), for instance, multi-resolut-

- Same data set, similar analysis with Seurat.
- But implemented in python.
- Using [jupyter notebook](#) as a interface.

Trajectory analysis by RNA velocity

- When cells differentiate, **new genes** will start to be expressed
- Transcripts have introns and will be spliced off given time
- Through assessing the present percentage of reads in introns, increase or decrease of expression can be modeled



<https://liorpachter.wordpress.com/tag/velocyto/>
<http://pklab.med.harvard.edu/software.html>

Get information about spliced/unsPLICED transcripts from bam files

```
#Build pipe line and shared scripts and instruction to researchers  
https://github.com/zhu16/single-cell-RNA-seq/tree/master/scvel\_notebooks
```

Step 1, finding reads mapped to un-spliced regions.

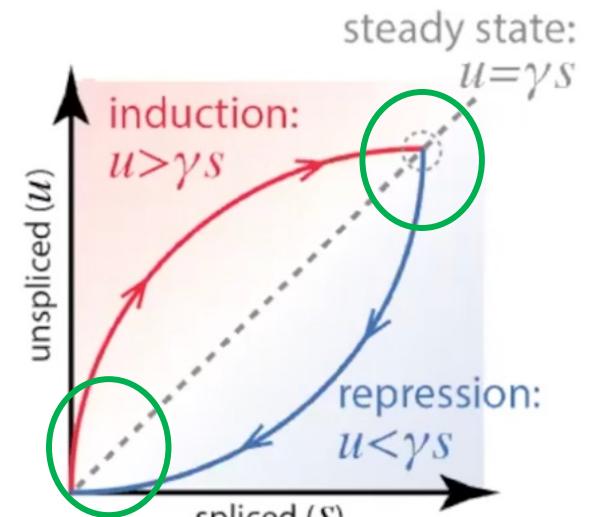
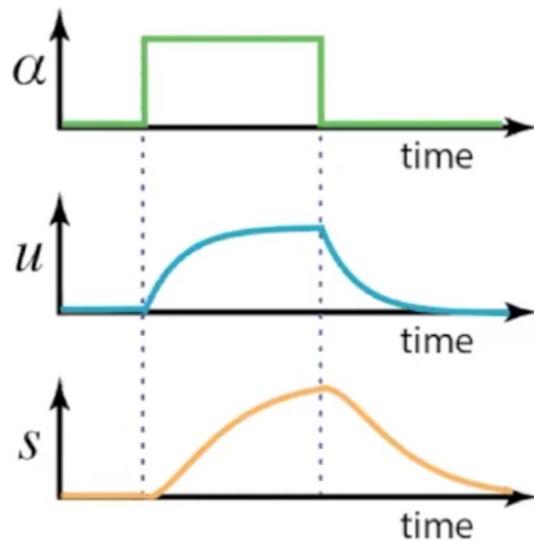
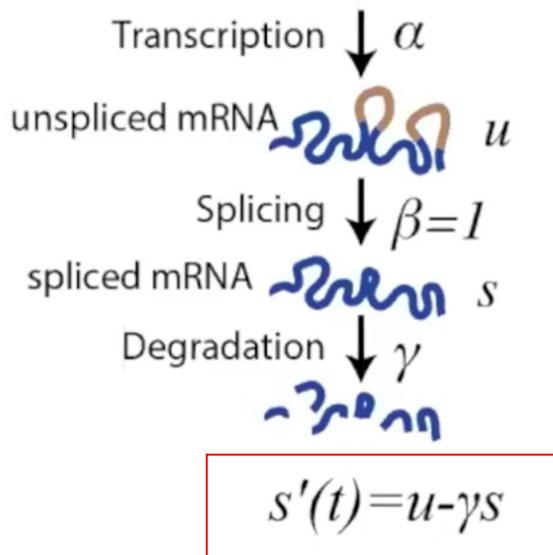
```
#Sort BAM files with Barcodes  
samtools sort -t CB -O BAM -  
o cellsorted_possorted_genome_bam.bam  
possorted_genome_bam.bam
```

```
#Analyze reads from spliced/non-spliced  
transcripts,  
module load python/3.7.3-foss-2016b  
velocyto run10x cellranger_output_folder genes.gtf
```

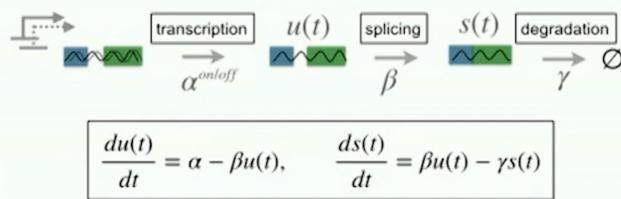
```
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0000.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0001.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0002.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0003.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0004.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0005.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0006.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0007.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0008.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0009.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0010.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0011.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0012.bAM
```

```
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0509.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0510.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0511.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0512.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0513.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0514.bAM  
cellsorTED_possORTed_genOME_bAM.bAM.tMP.0515.bAM
```

Modeling RNA dynamics



Concept of RNA velocity



Steady-state model (velocyto)

- Fit lin.reg. on extreme quantile cells (steady states)
- Estimate velocities as deviation from steady state

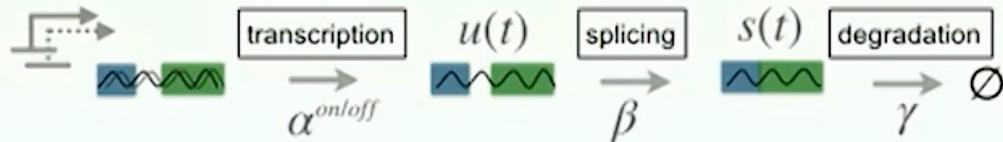
$$u_{\infty} \approx \gamma' s_{\infty} \quad (\beta = 1)$$

$$v_i = u_i - \gamma' s_i$$

2 assumptions:

steady states has been observed
a constant splicing rate β across all RNA

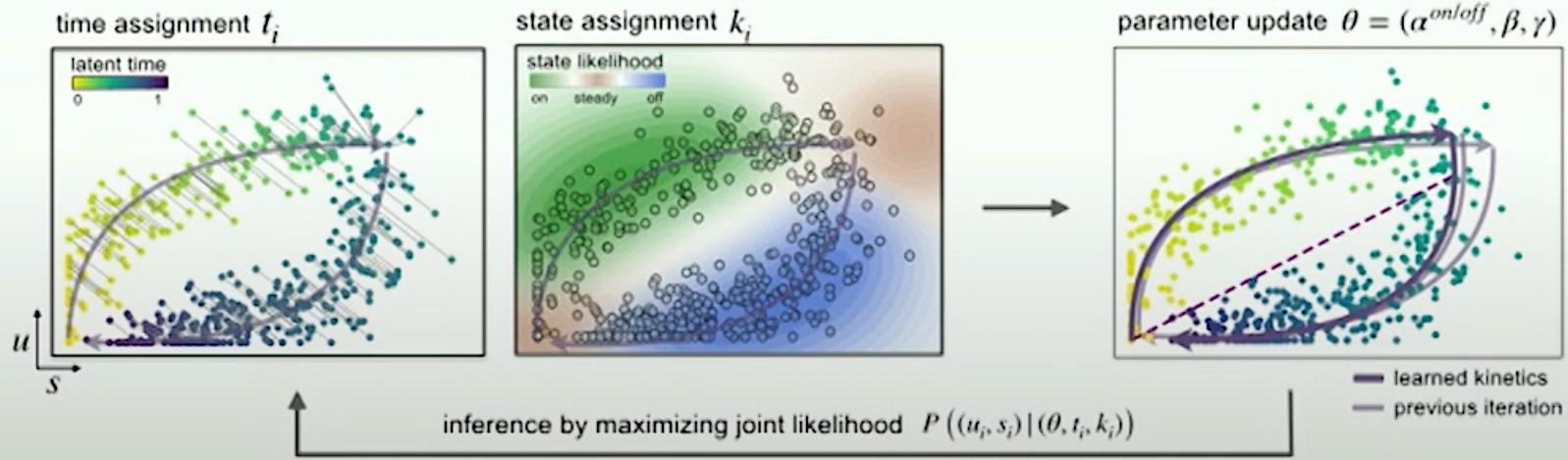
Generalizing RNA velocity to dynamical populations



$$u(t) = u_0 e^{-\beta \tau} + \frac{\alpha}{\beta} (1 - e^{-\beta \tau})$$

$$s(t) = s_0 e^{-\gamma \tau} + \frac{\alpha}{\gamma} (1 - e^{-\gamma \tau}) + \frac{\alpha - \beta u_0}{\gamma - \beta} (e^{-\gamma \tau} - e^{-\beta \tau}) \quad \tau = t - t_0$$

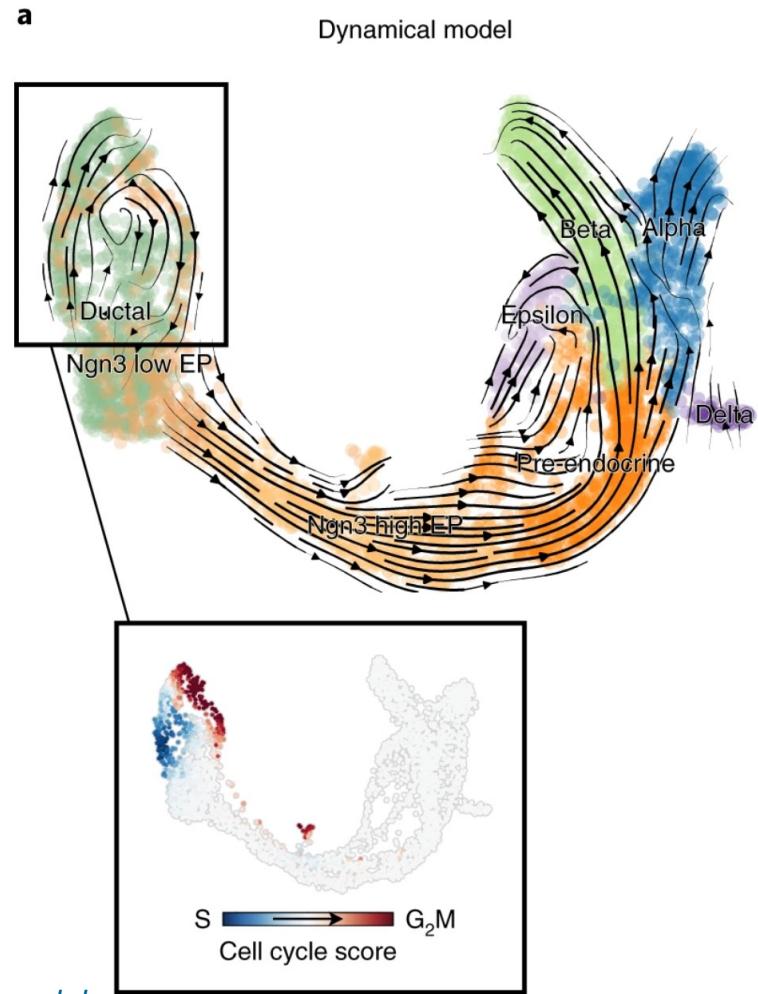
parameters of **reaction rates** $\theta = (\alpha^{off}, \alpha^{on}, \beta, \gamma)$
 cell-specific **latent variables** $\eta_i = (t_0^{(i)}, t_i, k_i)$
 (switch, time, state)



Steps of RNA-velocity

- From cellranger produced bam files
 - Sort bam files
 - Using Velocyto CLI to identify exon/intron reads as loom files.
 - Use scVelo to model the data and recover RNA dynamics
 - Through Markov process to predict which neighbor is the most probable destiny for the cell.
 - Backtracking and forward tracking to get the destiny and ancestor of the cell.

[Nature, 2018. invented the concept](#)
[Nature Biotechnology, 2020. a generalized model](#)



Demo on RNA velocity

Summary of day 1, scRNA-seq

- Set up the working environment
 - Login and pseudolinks
 - Copy files to scratch drive
 - Evoke jupyter lab IDE
- Introduction to scRNA-seq
- Make fastq from bcl files
- Use cellranger count to do alignment
- Use Seurat/scanpy pipeline to do QC and dimension reduction
- Use RNA velocity to infer trajectory
- For tomorrow: please install IGV and IGB for visualization of alignment, on your local machine.