

Lecture 13:

Single cell RNA-seq

You already have a strong
foundation for learning scRNA-seq

Read mapping

Filtering

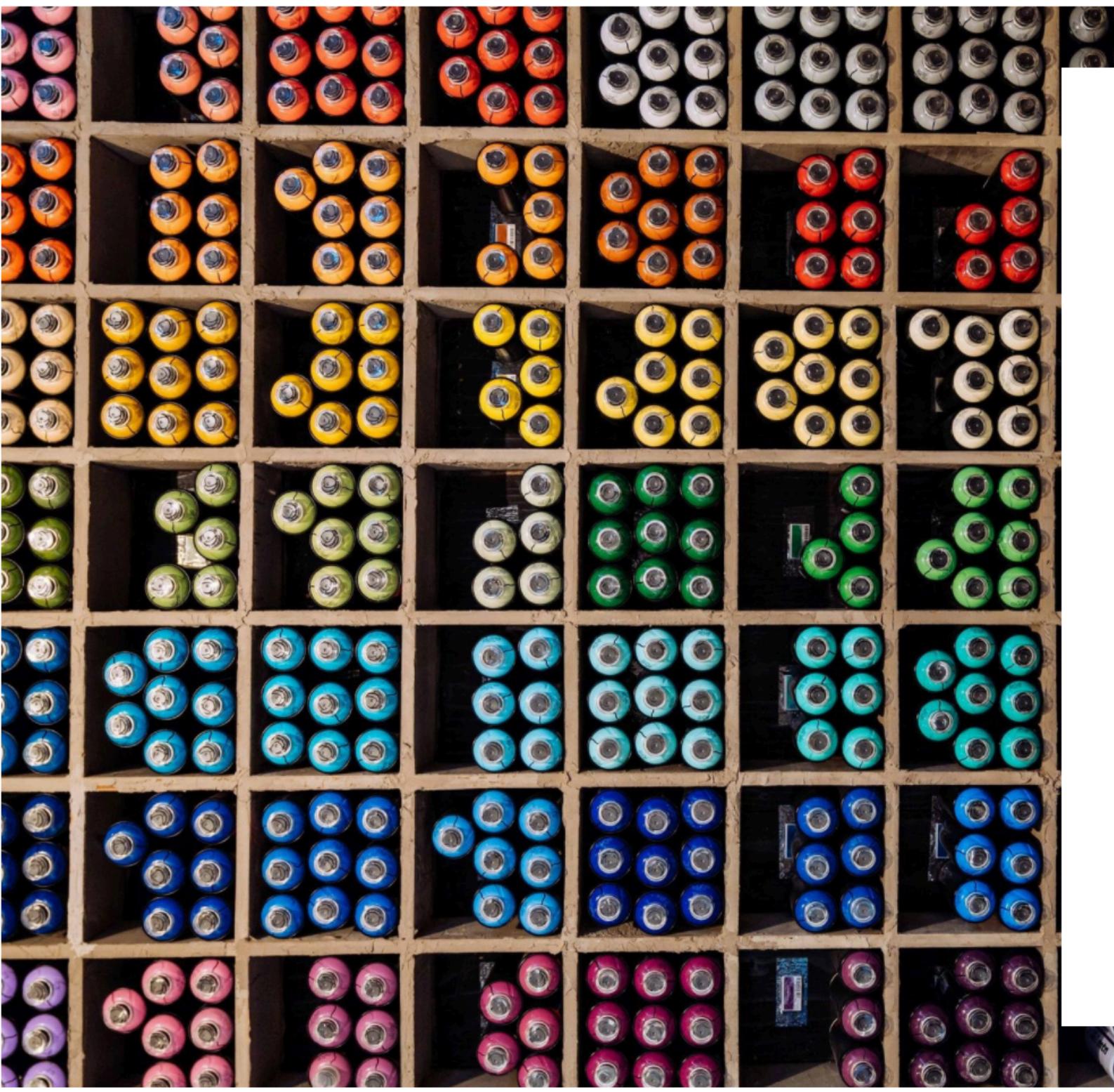
Normalization

Dimensional reduction

Differential gene testing

You already have a strong foundation for learning scRNA-seq

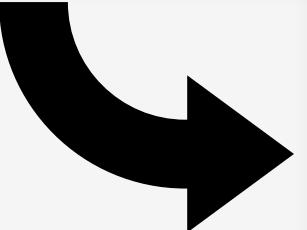




Single cell RNA-seq – principles and processing

Lecture 13 • watch by December 1, 2021

Now that you're comfortable with bulk RNA-seq data analysis, we'll shift our focus to the rapidly developing landscape of single cell RNA-seq (scRNA-seq). In this lecture, you'll learn about the underlying technology and demonstrate how to process raw single cell data directly on your laptop (!) for importing into R/bioconductor.



Single cell RNA-seq – principles and processing

Lecture 13 • watch by December 1, 2021

Overview

Now that you're comfortable with bulk RNA-seq data analysis, we'll shift our focus to the rapidly developing landscape of single cell RNA-seq (scRNA-seq). In this lecture, you'll learn about the underlying technology and we'll demonstrate how to process raw single cell data directly on your laptop (!) for importing into R/bioconductor.

Learning objectives

- Understand droplet-based scRNA-seq technology
- Be able to compare and contrast single cell and bulk RNA-seq methods
- Understand cost and experimental design considerations for scRNA-seq experiments.
- Familiarity with multiplexed single cell assays (CITE-seq, 'multiome', TEA-seq)
- Be able to define common terms and concepts in single cell genomics
- Use Kallisto-BUSTools to preprocess raw scRNA-seq data (via `kb-python`)

What you need to do

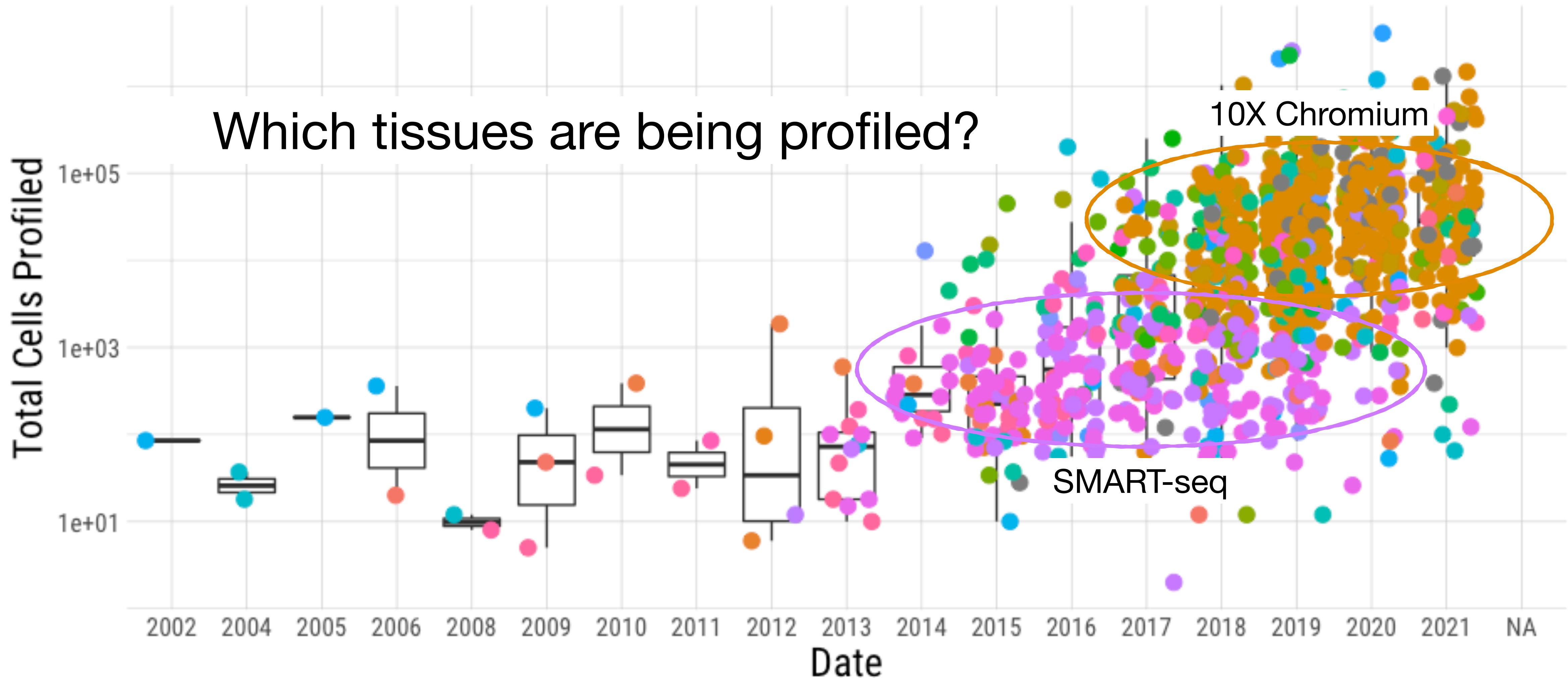
Download raw files. You will need about 5Gb of storage space on your harddrive to accomodate this download. *please do not uncompress these files (leave them as .gz files).* This is data from 1000 peripheral blood mononuclear cells (PBMCs) and is one of the sample datasets provided by 10X Genomics [here](#). I merged the separate lane files to make this simpler to work with for the course.

Human transcriptome reference index file - this is the index you created using Kallisto way back in lecture 2. If you don't have this, remember it's easy to create using `kallisto index`.

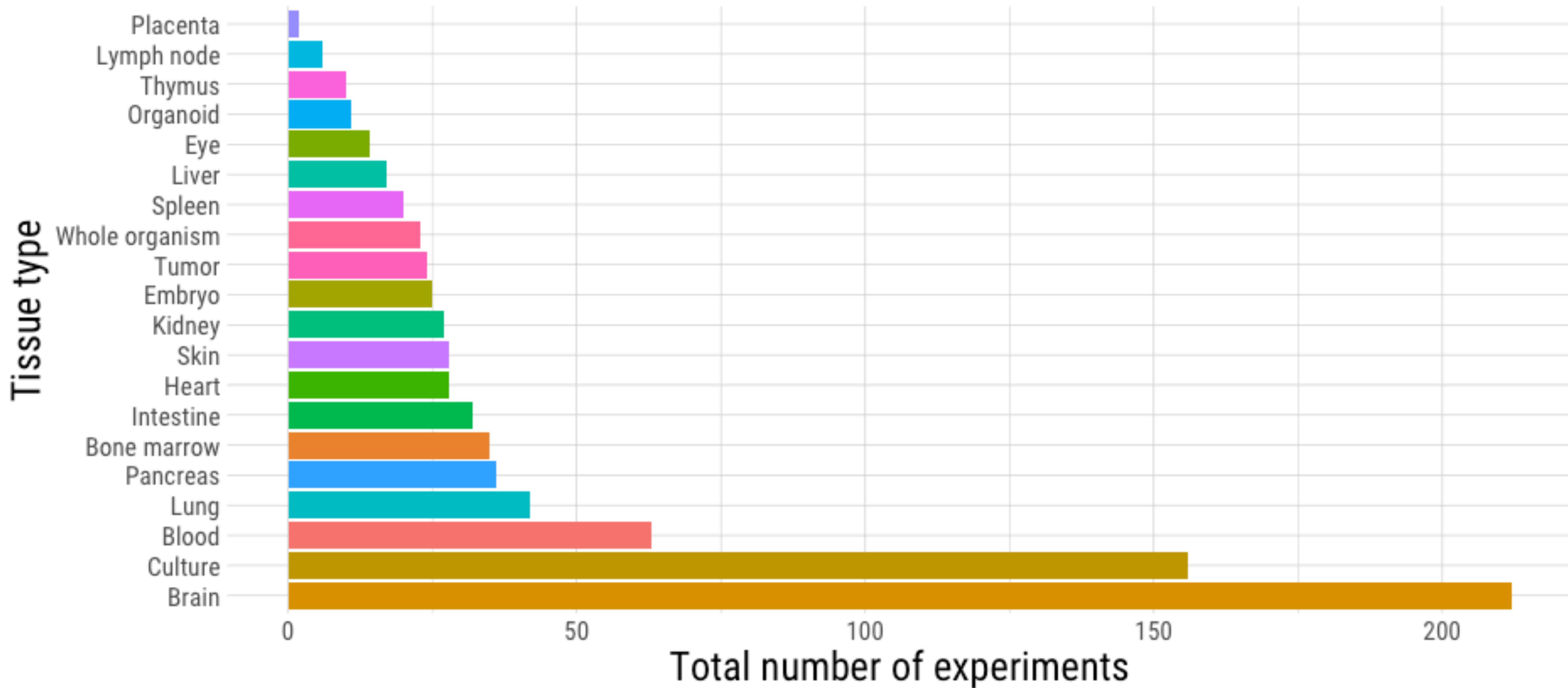
t2g.txt - this is a human transcript-to-gene mapping file that we will use with Kallisto-BUSTools to preprocess our data. This file is easy to generate with `kb ref`, but downloading it now will save you some time.

kb-python - You will need to have this software installed in a Conda environment on your laptop. We did this way back in [lecture 1](#). If you are unable to install or use kb-python, just follow along with the lecture so you understand the concepts.

The landscape of single cell sequencing (scRNA-seq) experiments



scRNAseq experiments by tissue type

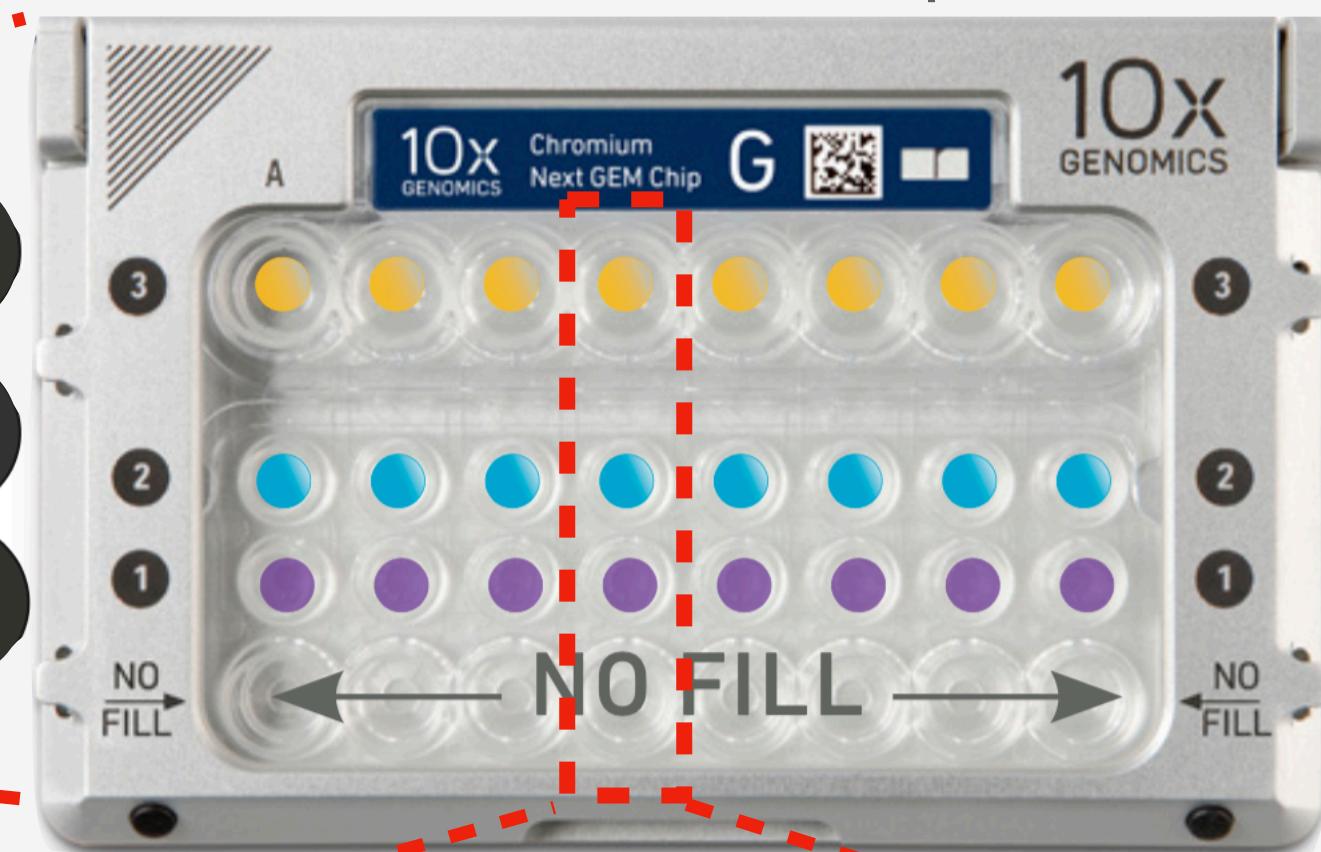


Chromium Controller



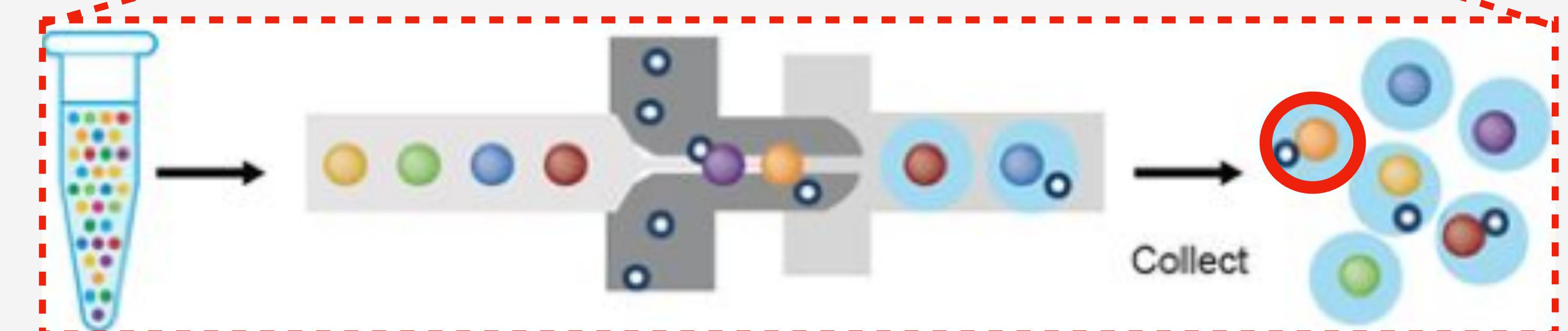
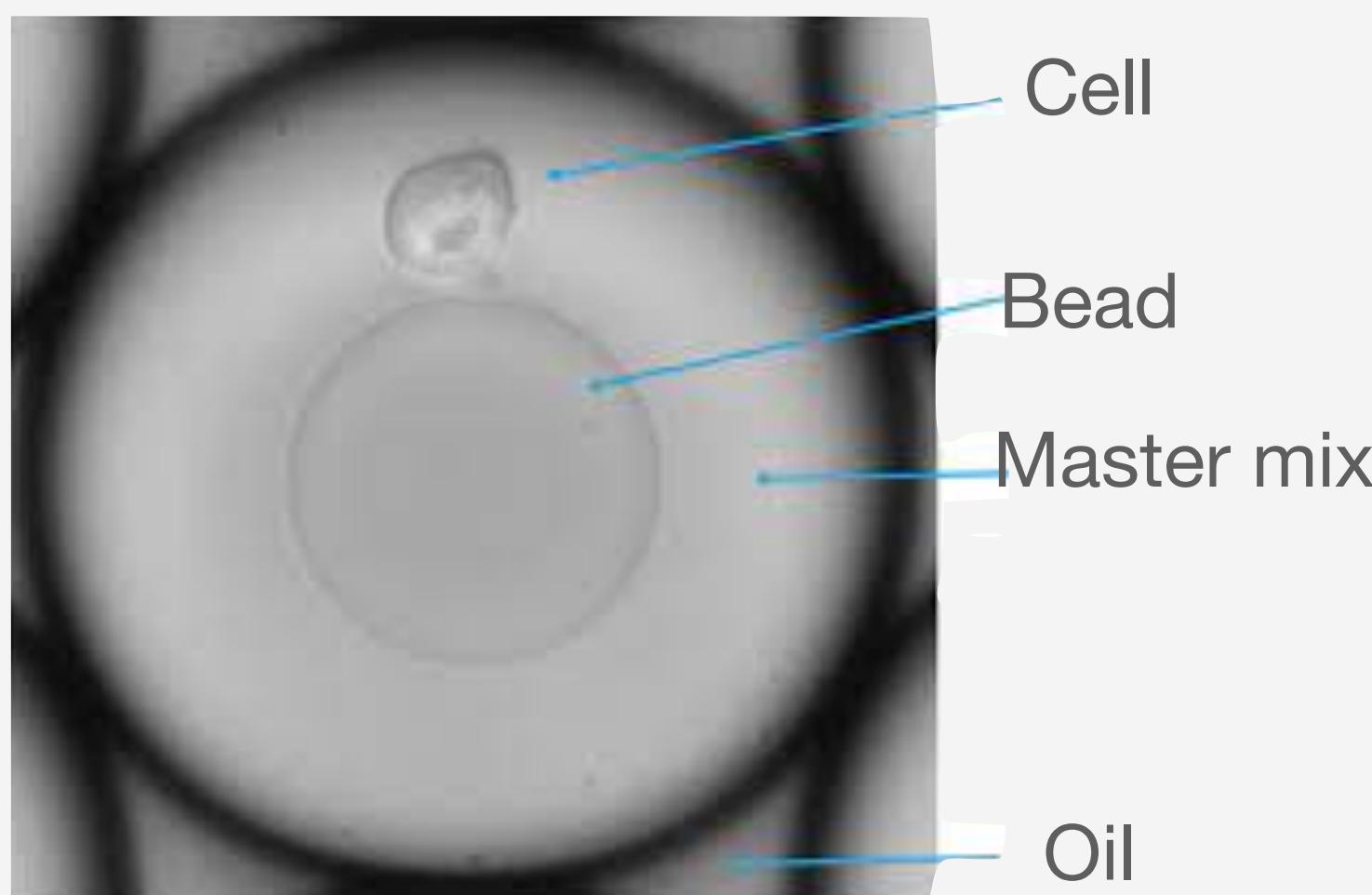
10x workflow for scRNA-seq

Microfluidic chip

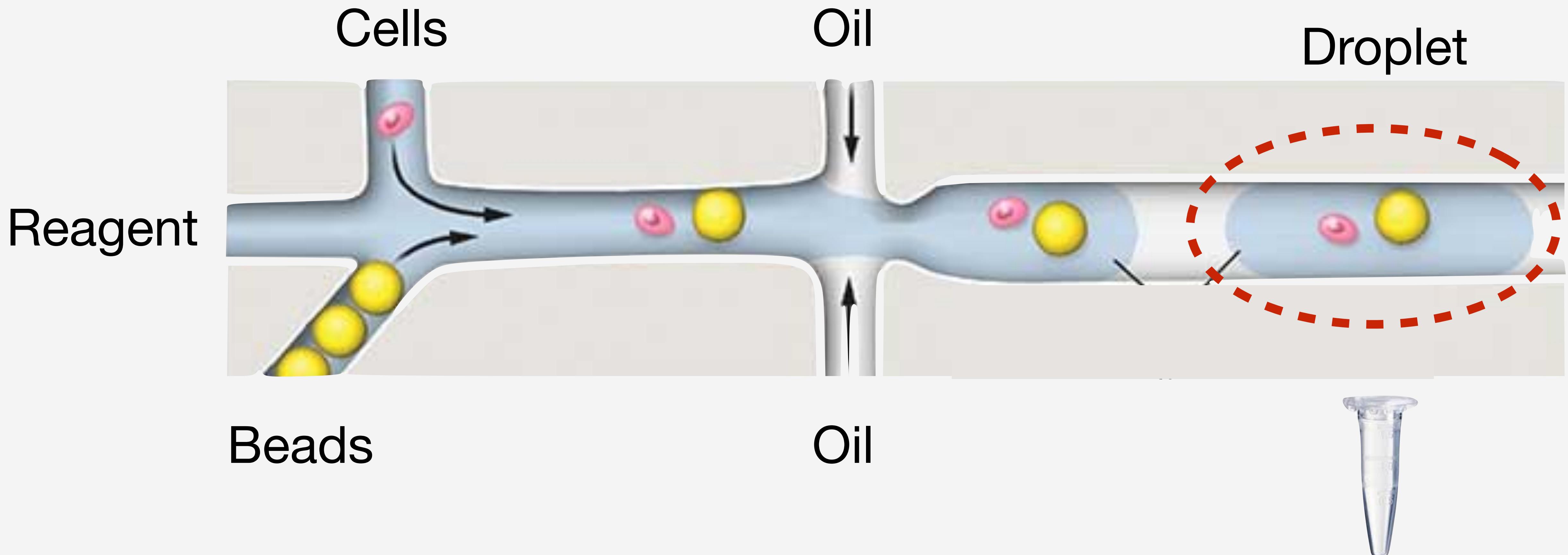


- Oil 3
- Gel Beads 2
- Master mix + Sample 1

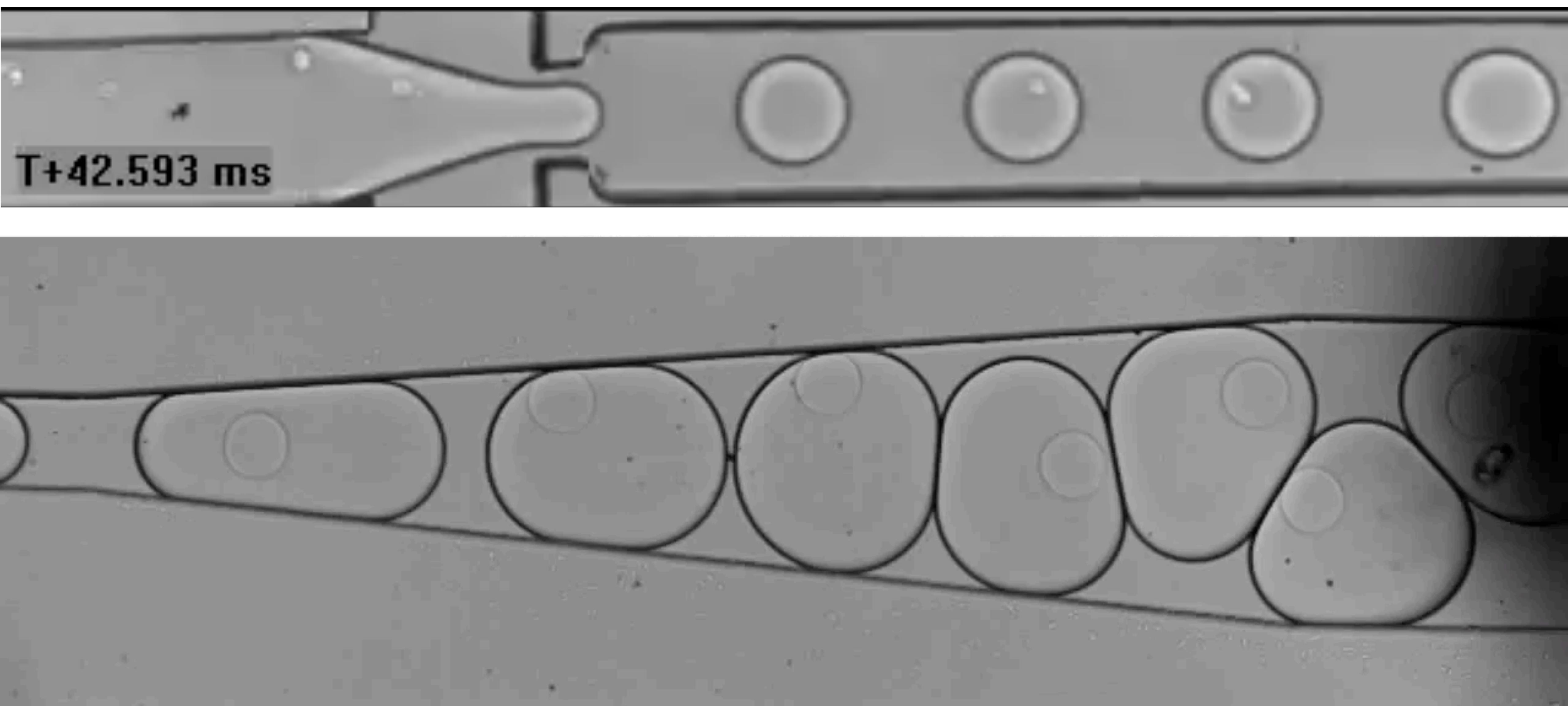
Droplet (a.k.a. GEM)



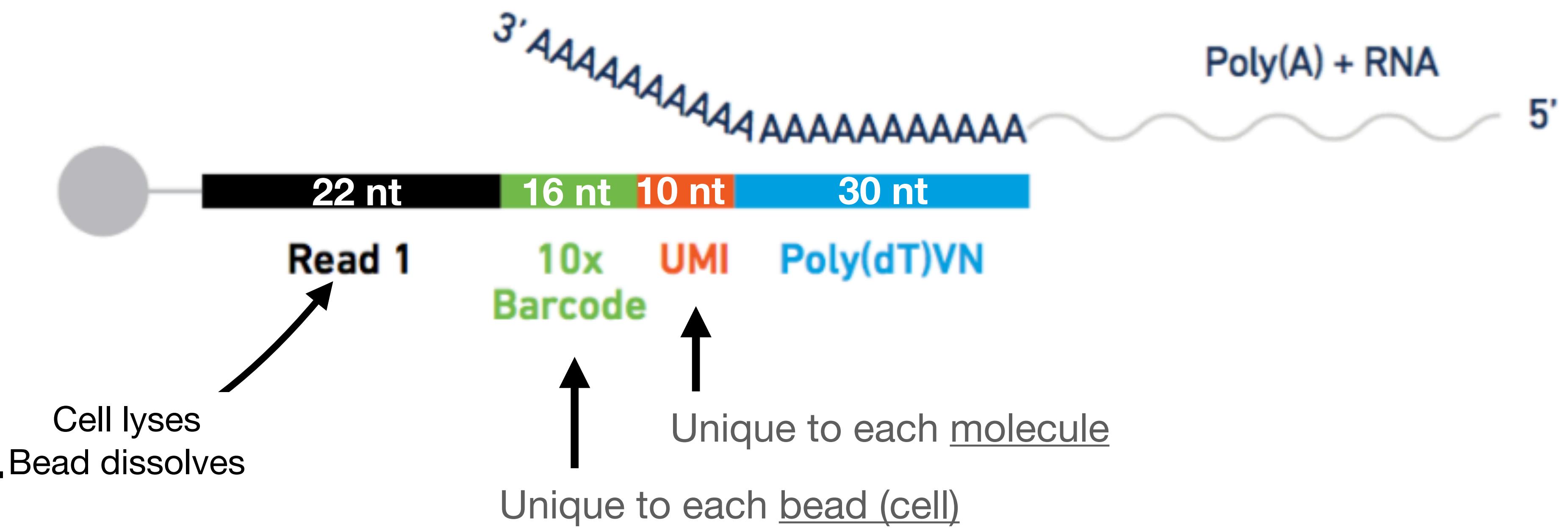
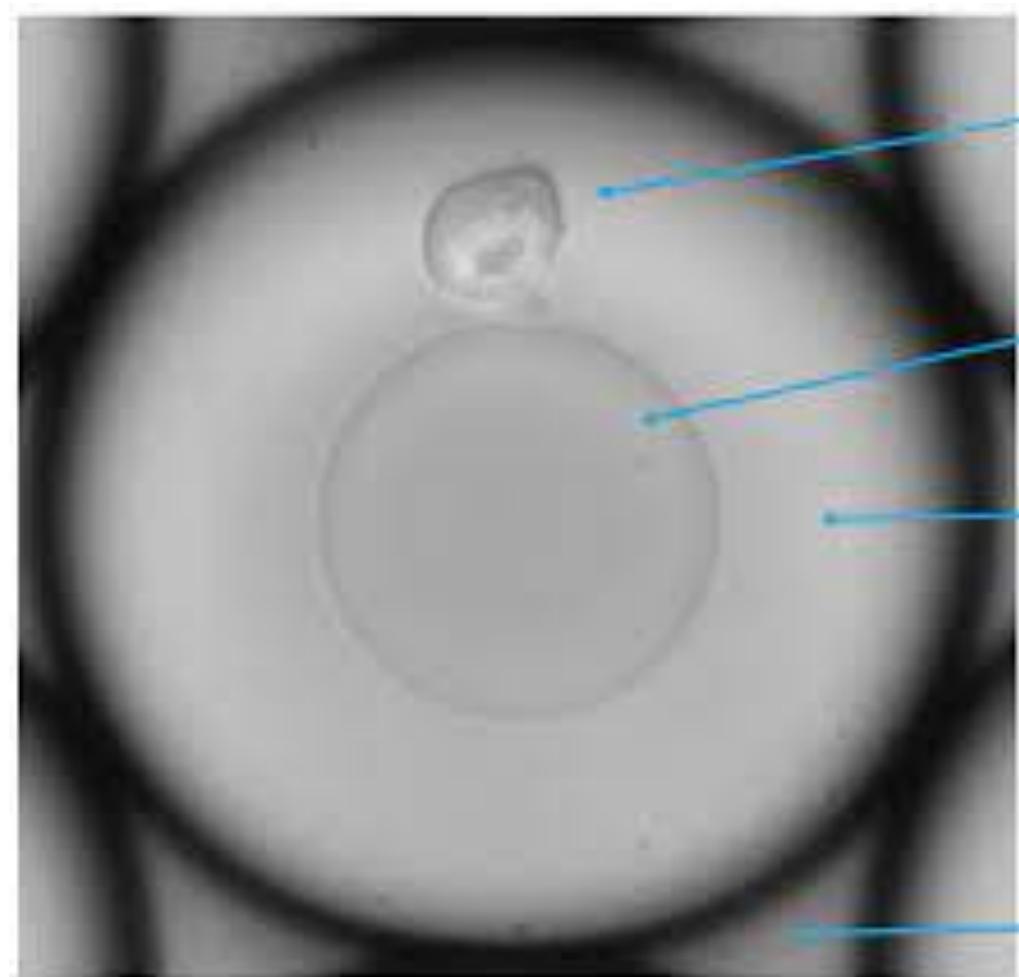
Droplet-based single cell RNAseq (scRNASeq, drop-seq)



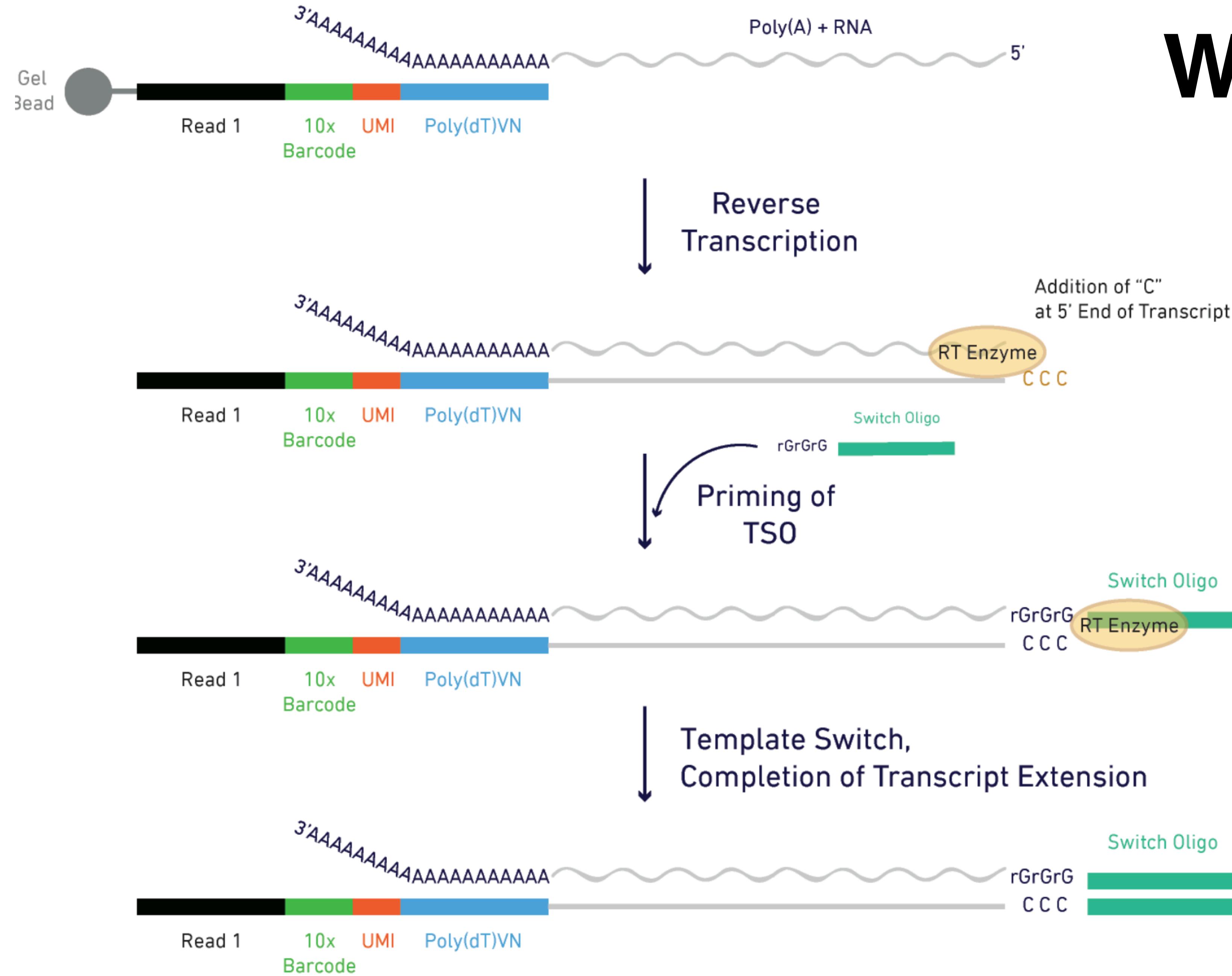
Droplet-based single-cell RNAseq (e.g., drop-seq, 10X)



What happens in a droplet?



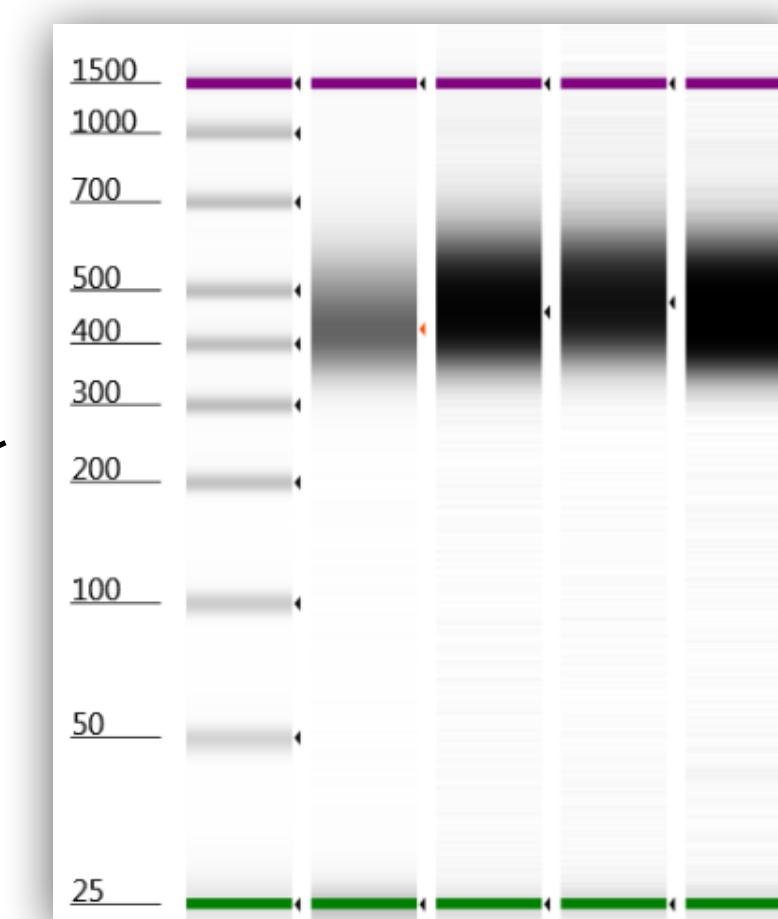
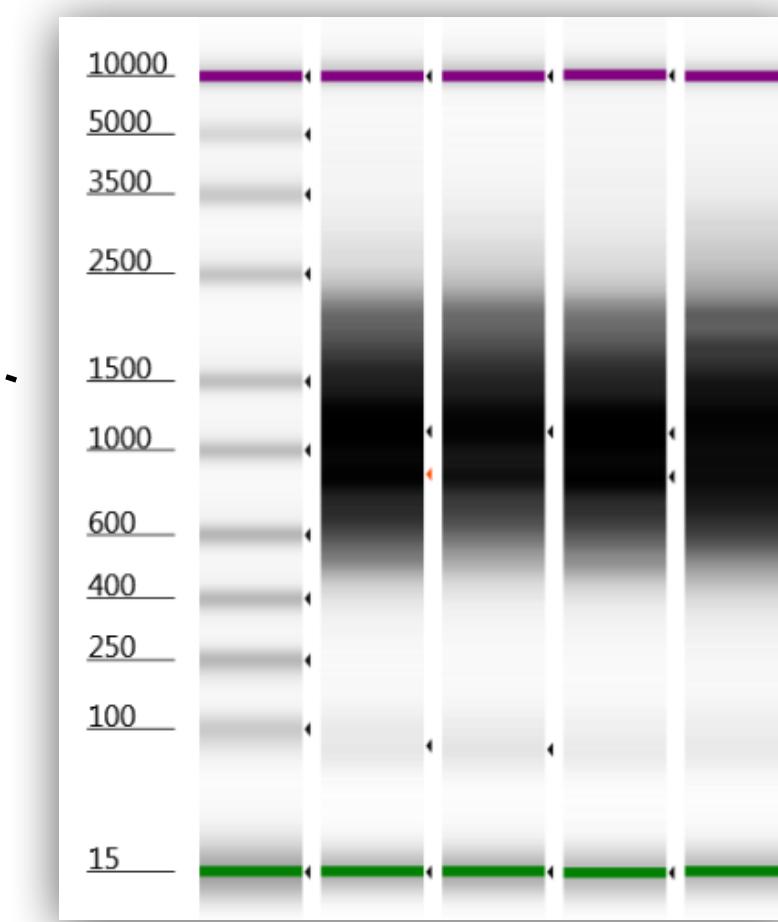
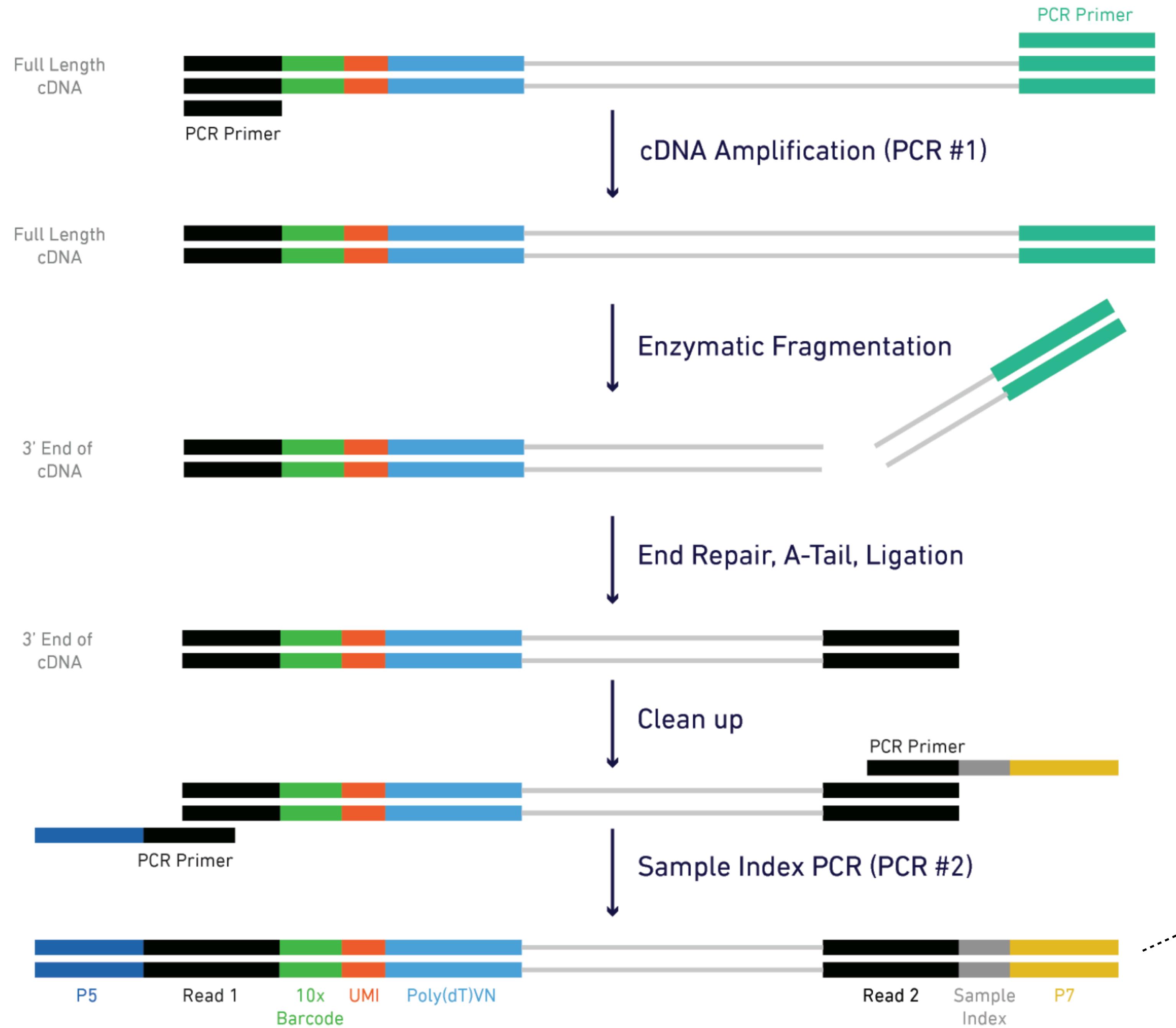
What happens in a droplet?



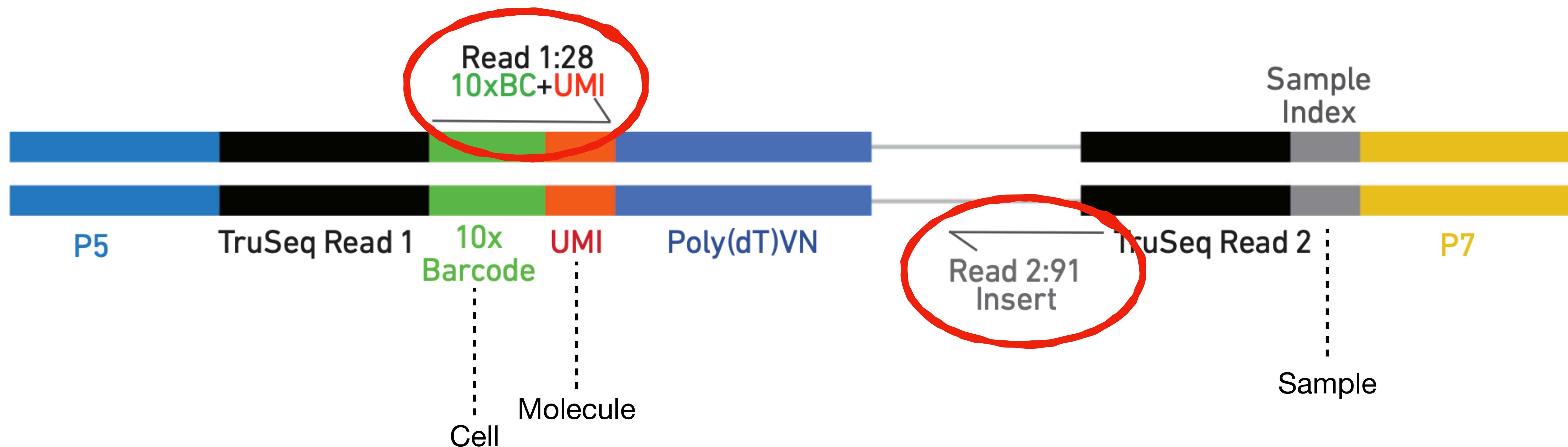
Droplets can now be broken, and remaining steps carried out in bulk

Library prep.

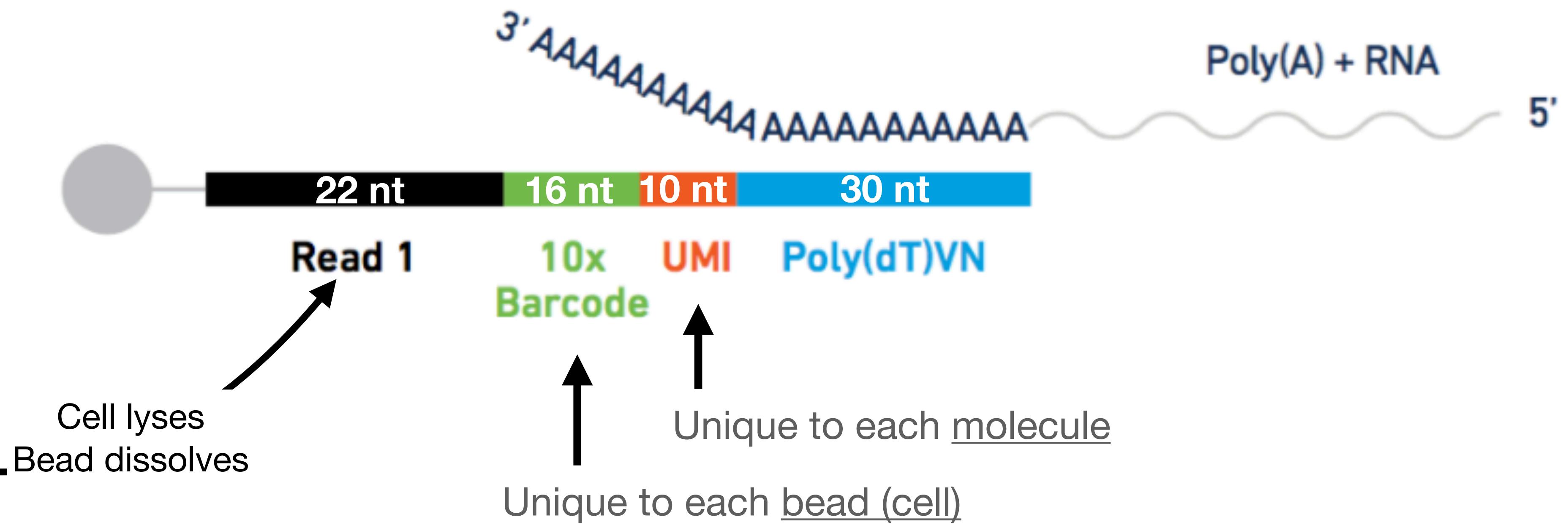
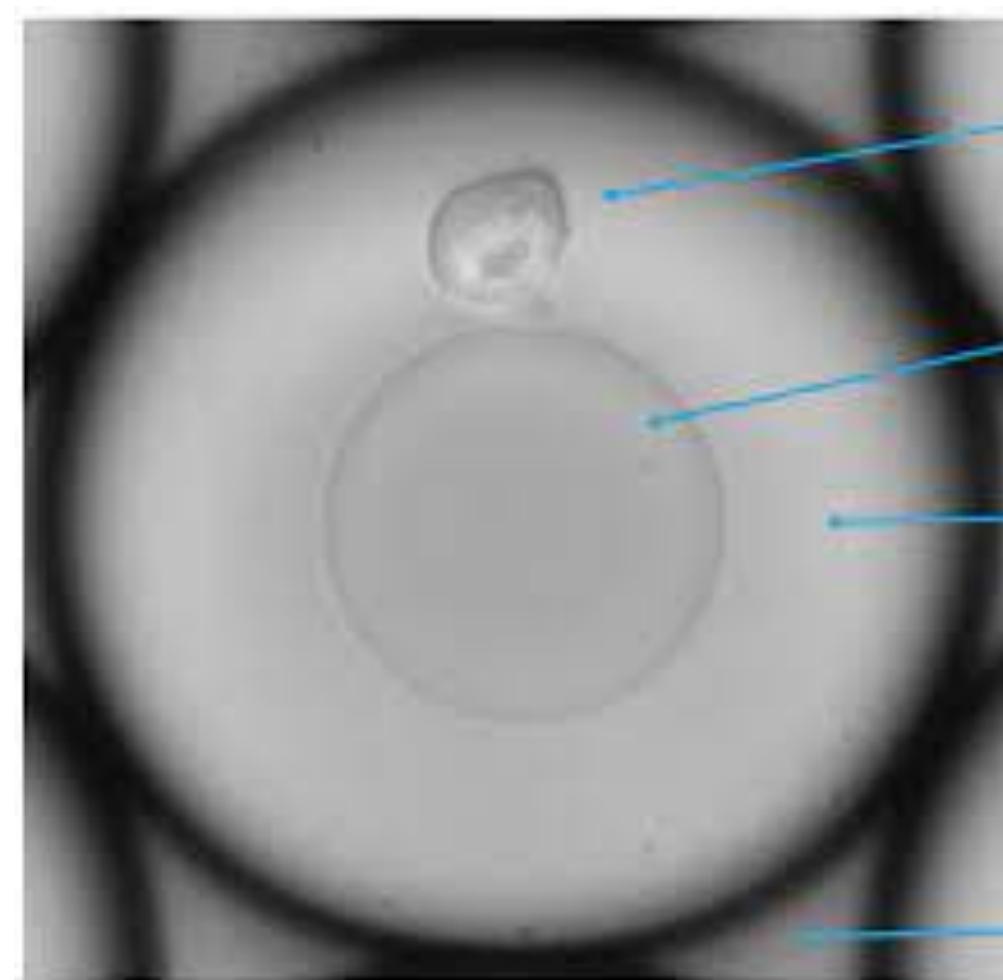
(3' scRNA-seq)



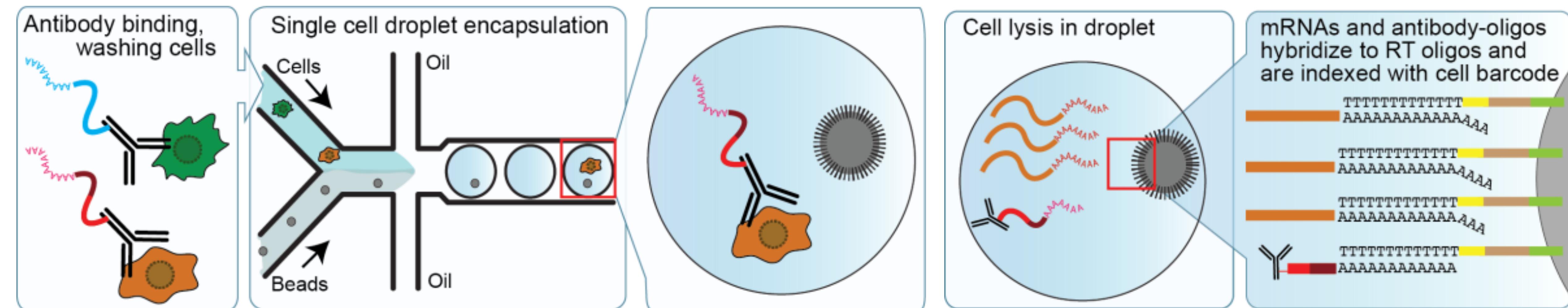
Sequencing



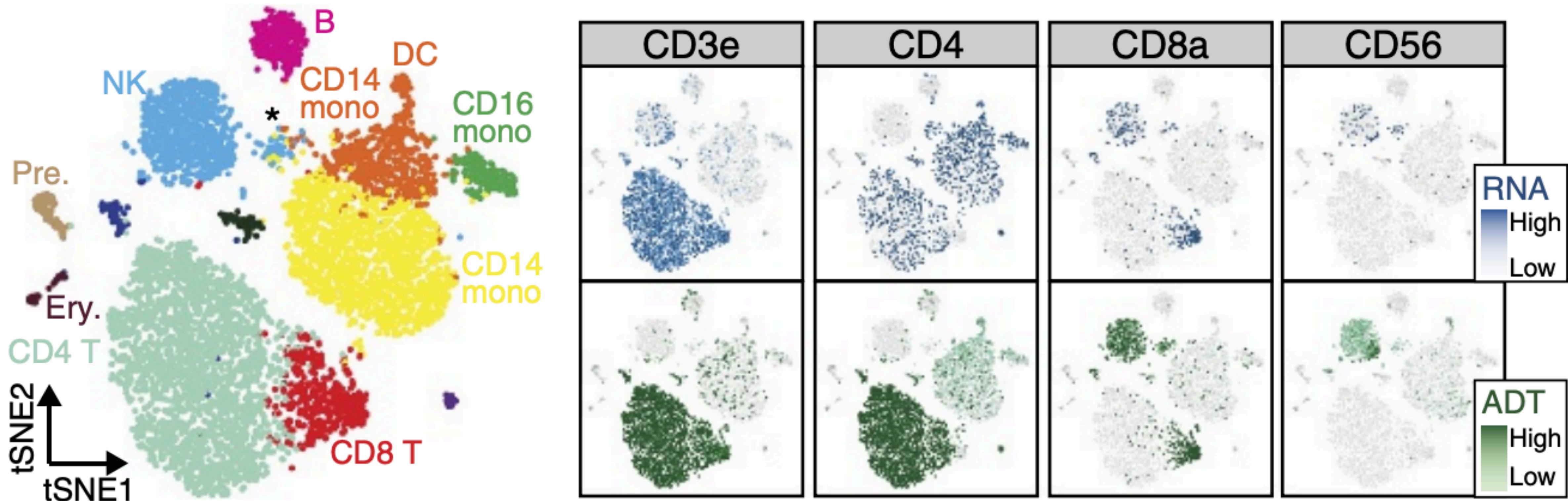
Barcoding is a flexible way to build new single cell assays



Cellular Indexing of Transcriptomes and Epitomes by Sequencing (CITE-seq)



Cellular Indexing of Transcriptomes and Epitomes by Sequencing (CITE-seq)



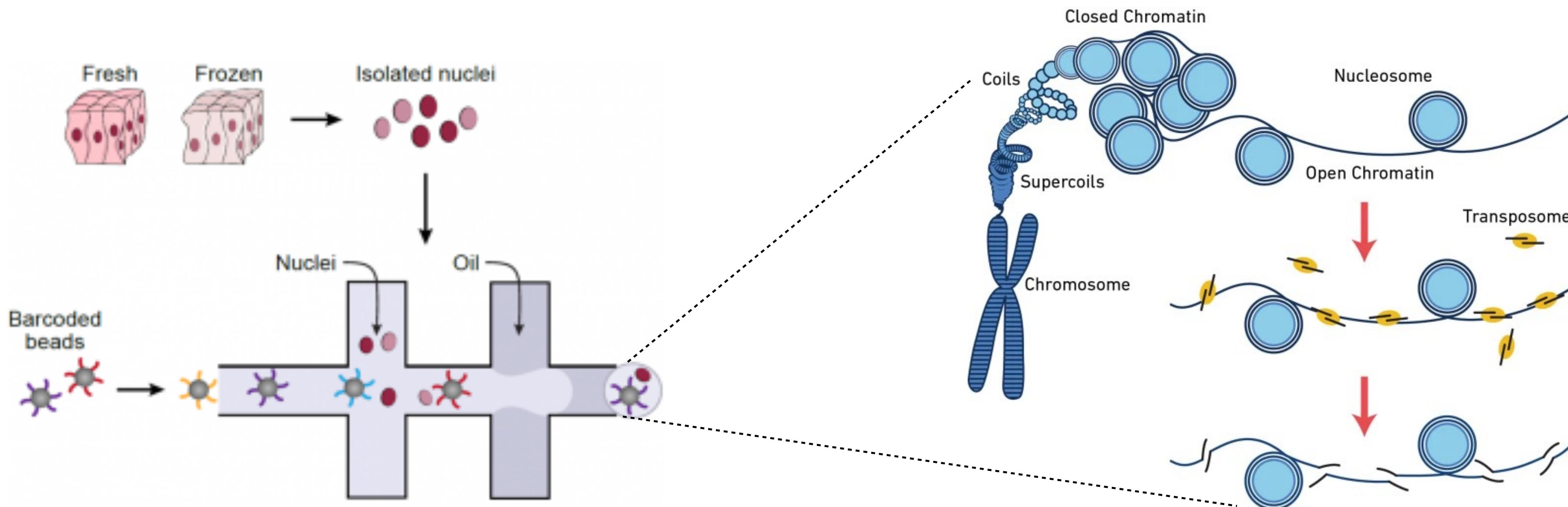
TotalSeq™ and Proteogenomic Analysis



Catalog	Barcode	Description	Clone	Reactivity	Sequence	Ensembl Gene Id
<input type="checkbox"/> 100569	0001	CD4	RM4-5	Mouse	AACAAGACCCTTGAG	ENSMUSG00000023274
<input type="checkbox"/> 100773	0002	CD8a	53-6.7	Mouse	TACCCGTAATAGCGT	ENSMUSG00000053977
<input type="checkbox"/> 119729	0003	CD366 (Tim-3)	RMT3-23	Mouse	ATTGGCACTCAGATG	ENSMUSG00000020399
<input type="checkbox"/> 109123	0004	CD279 (PD-1)	RMP1-30	Mouse	GAAAGTCAAAGCACT	ENSMUSG00000026285
<input type="checkbox"/> 305239	0005	CD80	2D10	Human	ACGAATCAATCTGTG	ENSG00000121594
<input type="checkbox"/> 305443	0006	CD86	IT2.2	Human	GTCTTGTCA GTGCA	ENSG00000114013
<input type="checkbox"/> 329743	0007	CD274 (B7-H1, PD-L1)	29E.2A3	Human	GTTGTCCGACAATAC	ENSG00000120217

Droplet methods are extremely versatile

*Single nuclei encapsulation for frozen or fixed tissue,
or for scATAC-seq*



Droplet methods are extremely versatile

Multiplexed assays

eLife | TOOLS AND RESOURCES | |

Simultaneous trimodal single-cell measurement of transcripts, epitopes, and chromatin accessibility using TEA-seq

Elliott Swanson¹, Cara Lord^{1†}, Julian Reading¹, Alexander T Heubeck¹, Palak C Genge¹, Zachary Thomson¹, Morgan DA Weiss¹, Xiao-jun Li¹, Adam K Savage¹, Richard R Green^{1,2‡}, Troy R Torgerson^{1,3}, Thomas F Bumol¹, Lucas T Graybuck^{1*}, Peter J Skene^{1*}

¹Allen Institute for Immunology, Seattle, United States; ²Department of Biomedical Informatics and Medical Education (BIME), University of Washington, Seattle, United States; ³Department of Pediatrics, University of Washington, Seattle, United States

Abstract Single-cell measurements of cellular characteristics have been instrumental in understanding the heterogeneous pathways that drive differentiation, cellular responses to signals, and human disease. Recent advances have allowed paired capture of protein abundance and transcriptomic state, but a lack of epigenetic information in these assays has left a missing link to gene regulation. Using the heterogeneous mixture of cells in human peripheral blood as a test case, we developed a novel scATAC-seq workflow that increases signal-to-noise and allows paired measurement of cell surface markers and chromatin accessibility: integrated cellular indexing of chromatin landscape and epitopes, called ICICLE-seq. We extended this approach using a droplet-based multiomics platform to develop a trimodal assay that simultaneously measures transcriptomics (scRNA-seq), epitopes, and chromatin accessibility (scATAC-seq) from thousands of single cells, which we term TEA-seq. Together, these multimodal single-cell assays provide a novel toolkit to identify type-specific gene regulation and expression grounded in phenotypically defined cell types.

For correspondence: lucasg@alleninstitute.org (LTG); peter.skene@alleninstitute.org (PJS)

Present address: [†]GlaxoSmithKline, Collegeville, United States; [‡]Department of Biomedical Informatics and Medical Education, University of Washington School of Medicine, Seattle, United States

Competing interest: See page 34

Funding: See page 34

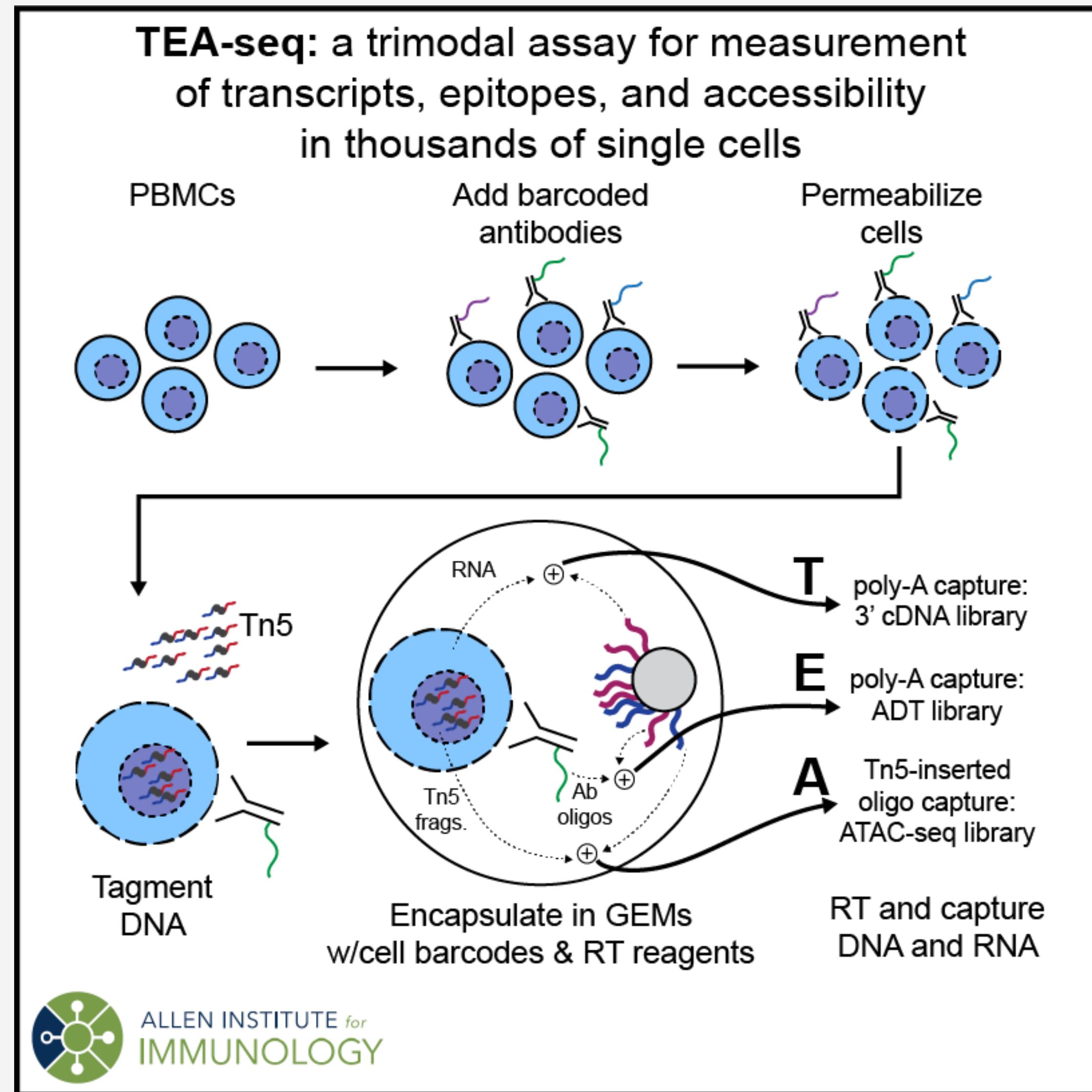
Received: 30 September 2020
Accepted: 11 March 2021
Published: 09 April 2021

Reviewing editor: Howard Y Chang, Stanford University, United States

© Copyright Swanson et al. This article is distributed under the terms of the [Creative Commons Attribution License](#), which permits unrestricted use and redistribution provided that the original author and source are credited.

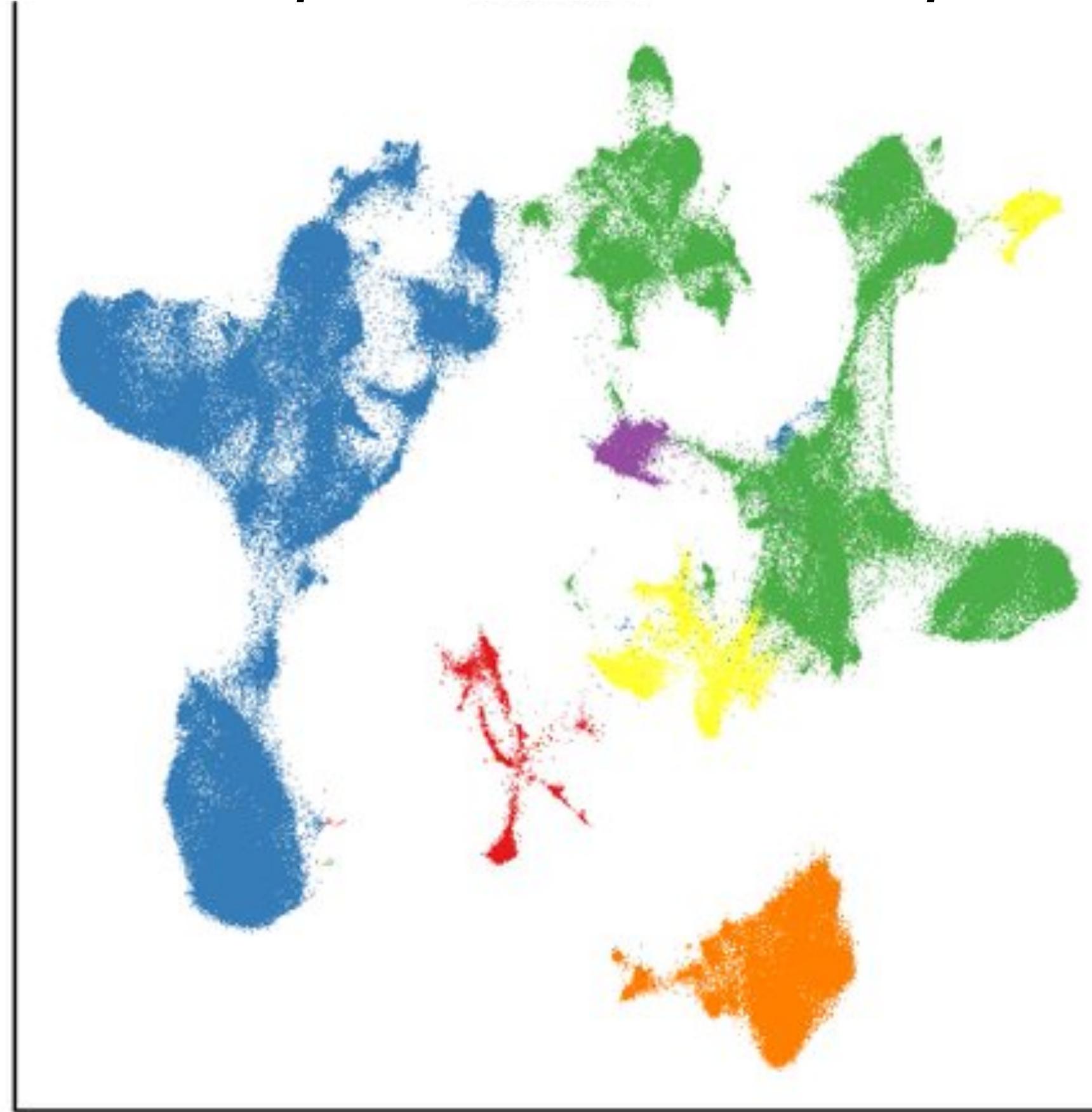
Swanson et al. eLife 2021;10:e63632. DOI: <https://doi.org/10.7554/eLife.63632>

1 of 38

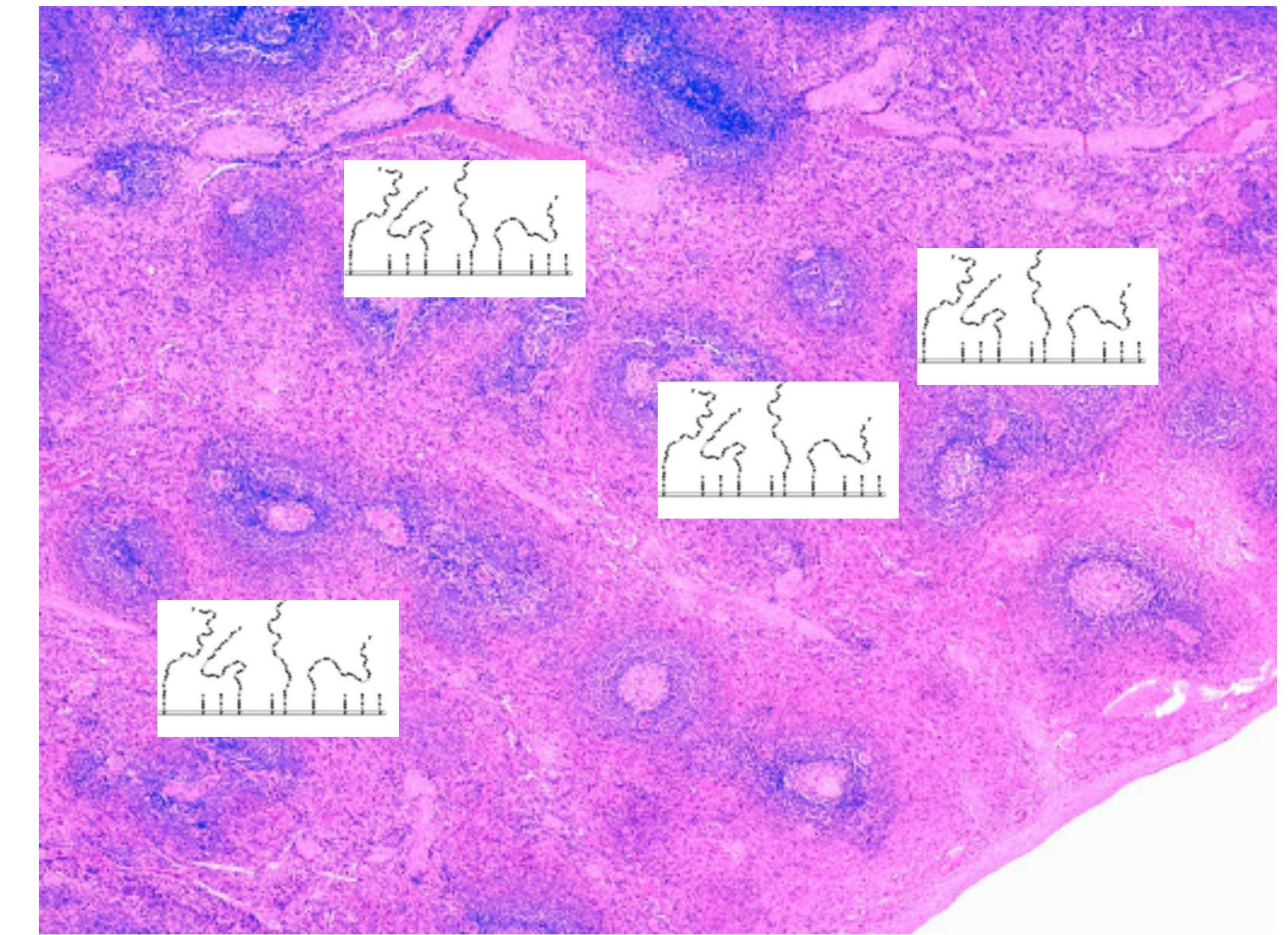


Spatial information is lost in scRNA-seq

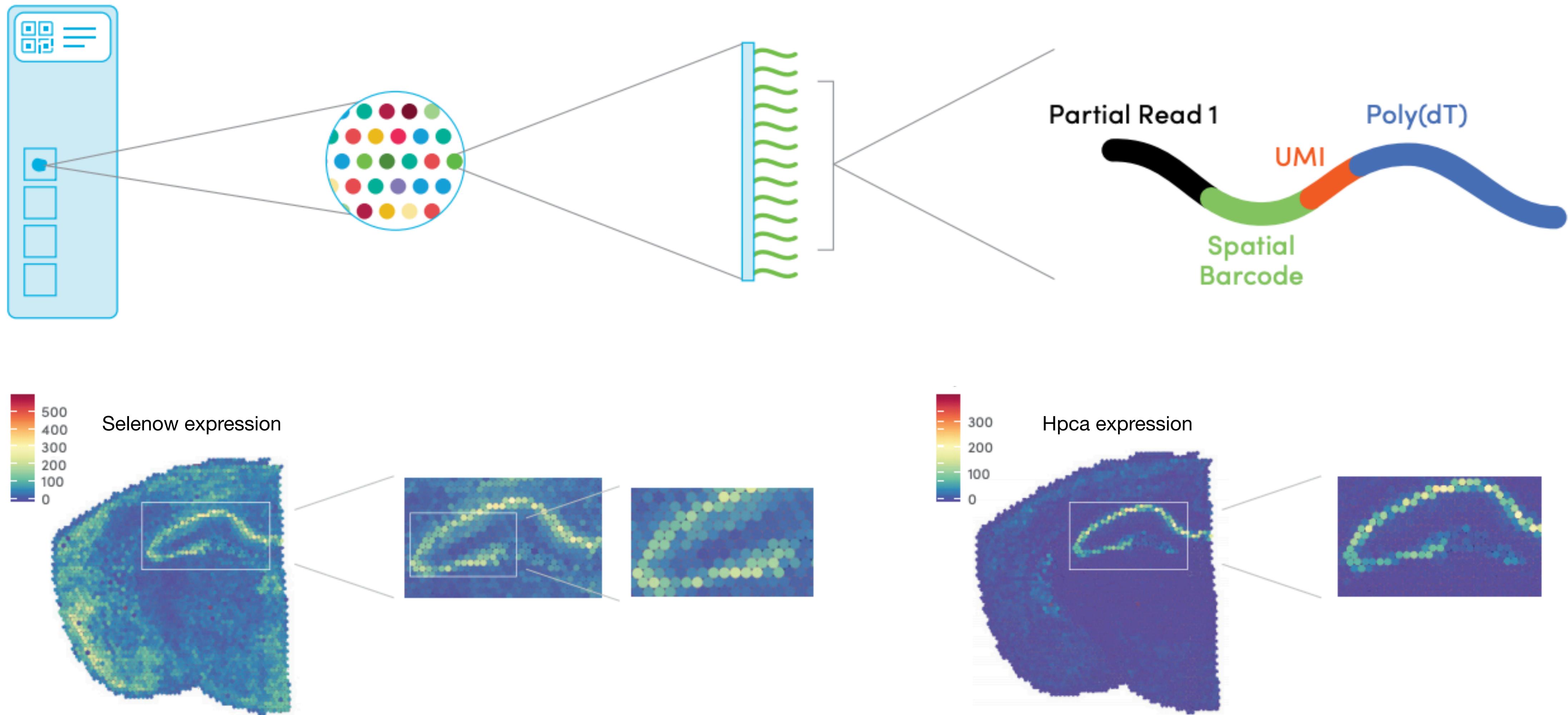
Spleen scRNA-seq



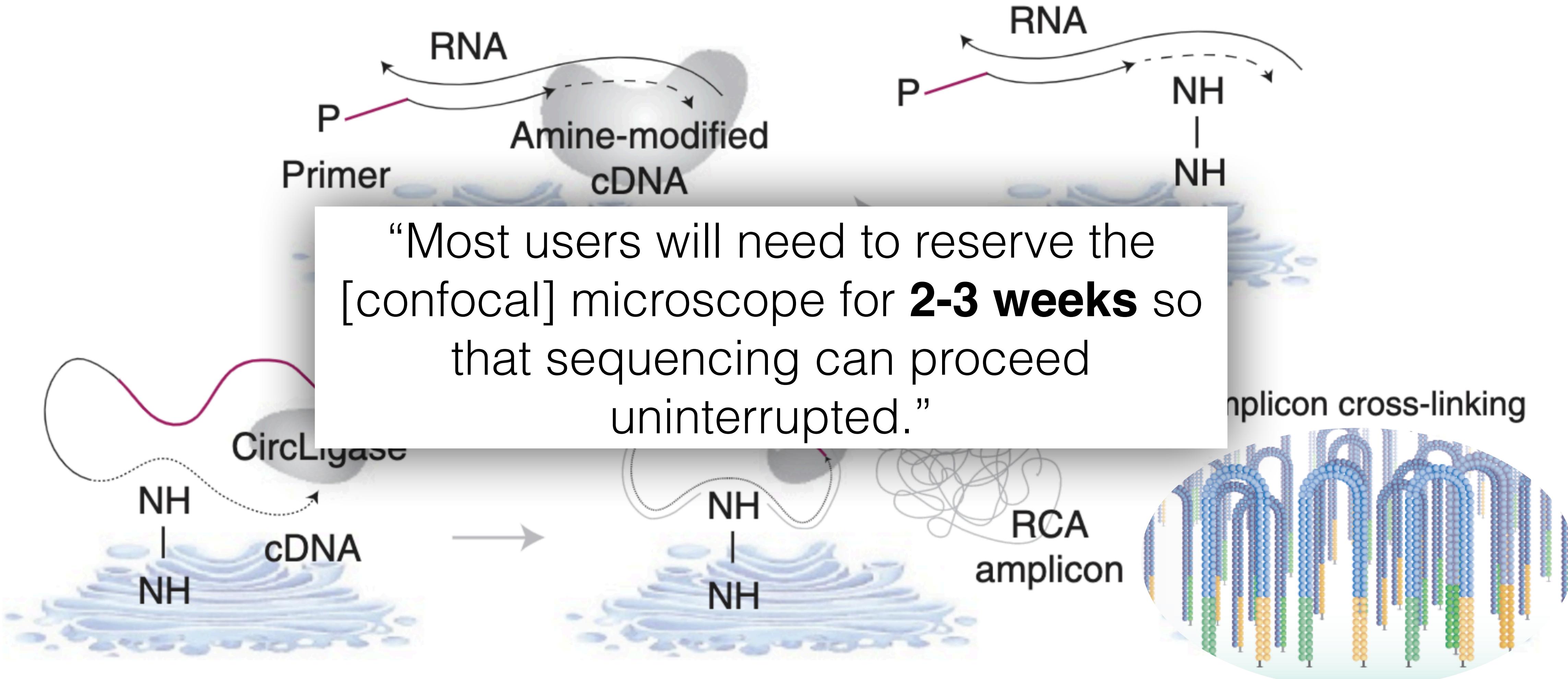
Actual spleen



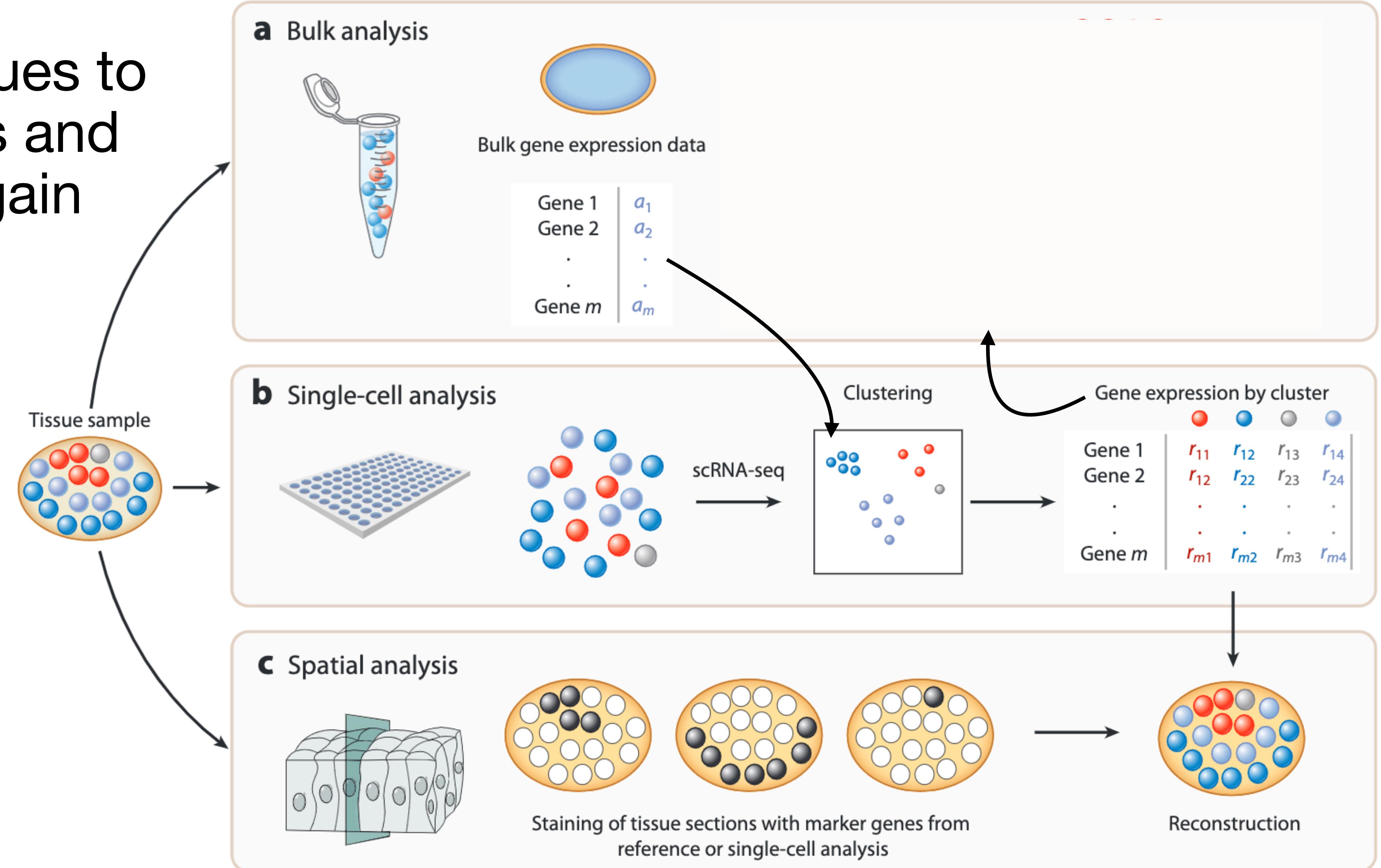
Spatial transcriptomics with ‘zipcoded’ beads on a slide



Fluorescent In Situ Sequencing (FISSEQ)



From Tissues to cell types and back again



End of part 1

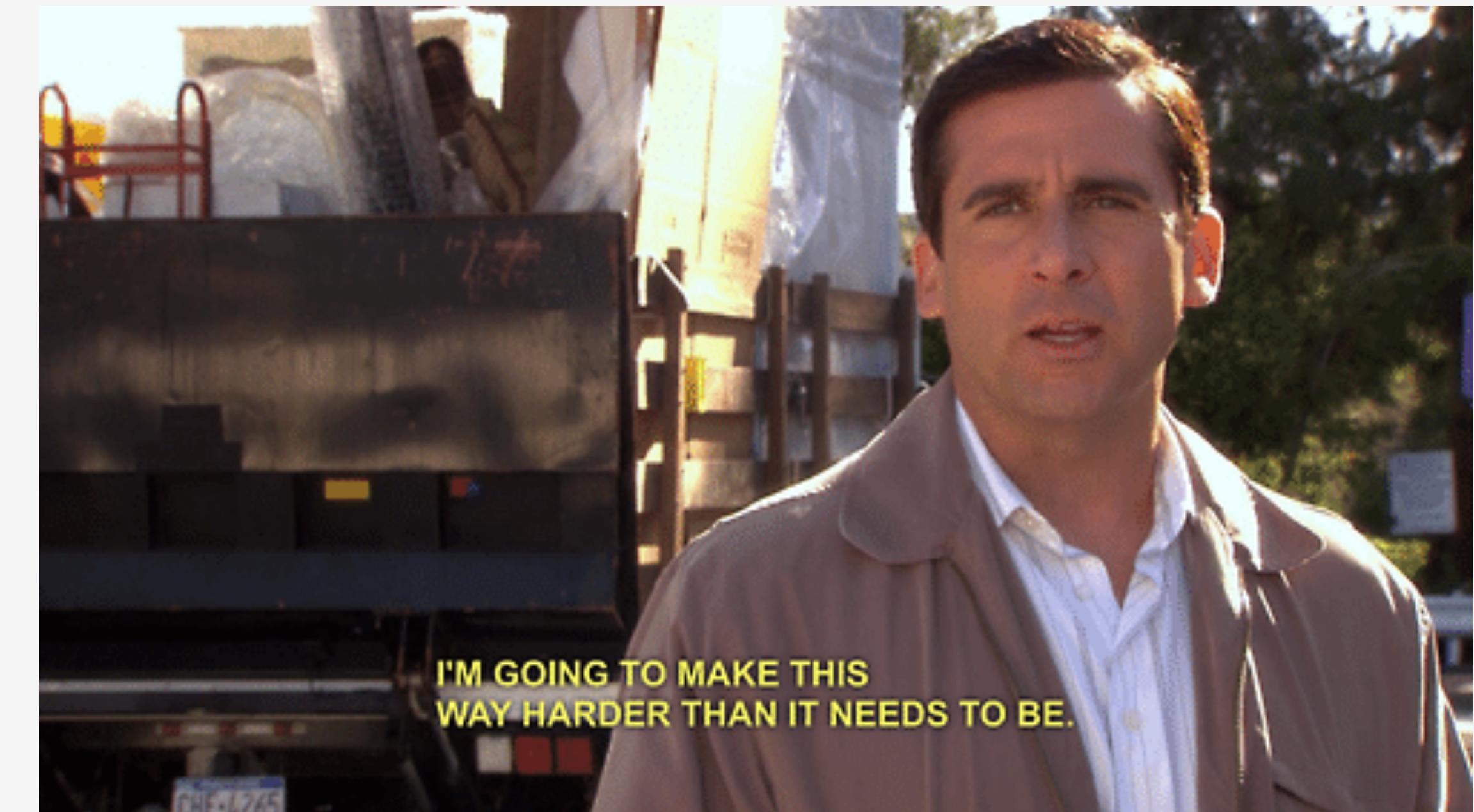
Practical considerations for single cell studies

Practical considerations for single cell studies

Do you *really* need single cell resolution,
or would flow cytometry or bulk RNA-seq be sufficient to answer your question?

scRNA-seq ➤ bulk RNA-seq

	Single cell	Bulk
Heterogeneity	✓	✗
Trajectory	✓	✗
Multiplexing	✓	✗
Depth	+	+++
Cost	++++++	+
Analytic burden	+++	+



Costs associated with bulk* sequencing

From pre-course lecture

Library Preparation



Sequencing



TruSeq kit	cat#	cost	cycles	cat#	cost
stranded total RNA LT (w/ RiboZero)	RS-122-2201	\$5,067	300	FC-404-2004	\$4,222
stranded mRNA LT	RS-122-2101	\$2,430	150	FC-404-2002	\$2,635
			75	FC-404-2005	\$1,374

all kits process 48 samples

12 samples
library prep = \$51/sample
sequencing = \$115/sample
data output = 30M reads/sample

*Costs associated with scRNAseq
are much more complicated

Costs associated with scRNA-seq

Library Preparation



Sequencing



P2 flowcell
400M clusters

P3 flowcell
1.1B clusters

10x Genomics	cat#	cost	Flow cell	cycles	cat#	cost	samples
Chromium Next GEM Single Cell 3' kit v3.1	1000130 (4rxn)	~\$2000/sample	P2	100 (+38)	20046811	\$1420	~2
Chromium Next GEM Single Cell 3' LT kit v3.1	1000322 (4rxn)	~\$900/sample	P3	100 (+38)	20040559	\$3200	~6

2 samples; 10,000 cells/sample;
sequenced on P2

library prep = ~\$2000/sample

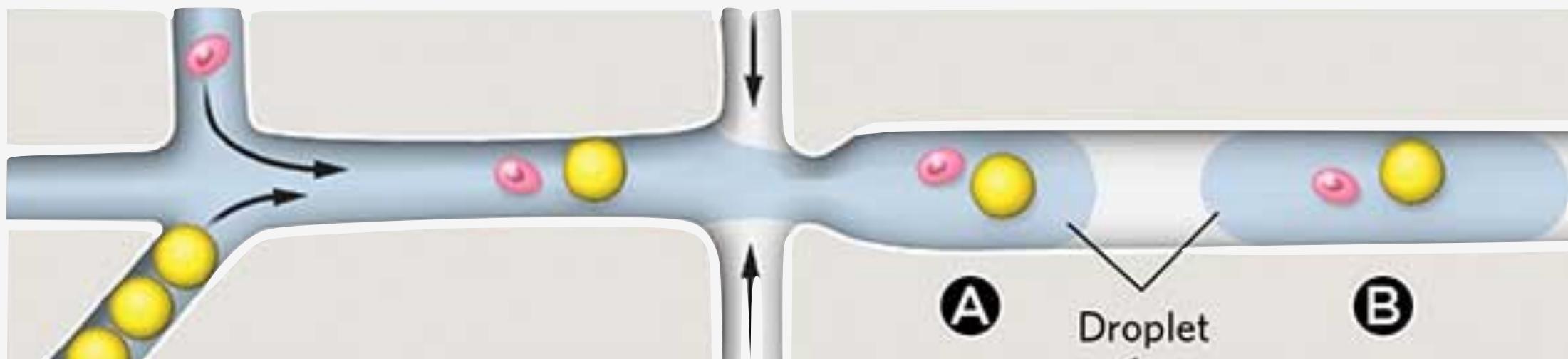
sequencing = ~\$700/sample

data output = ~20,000 reads/cell

Practical considerations for single cell studies

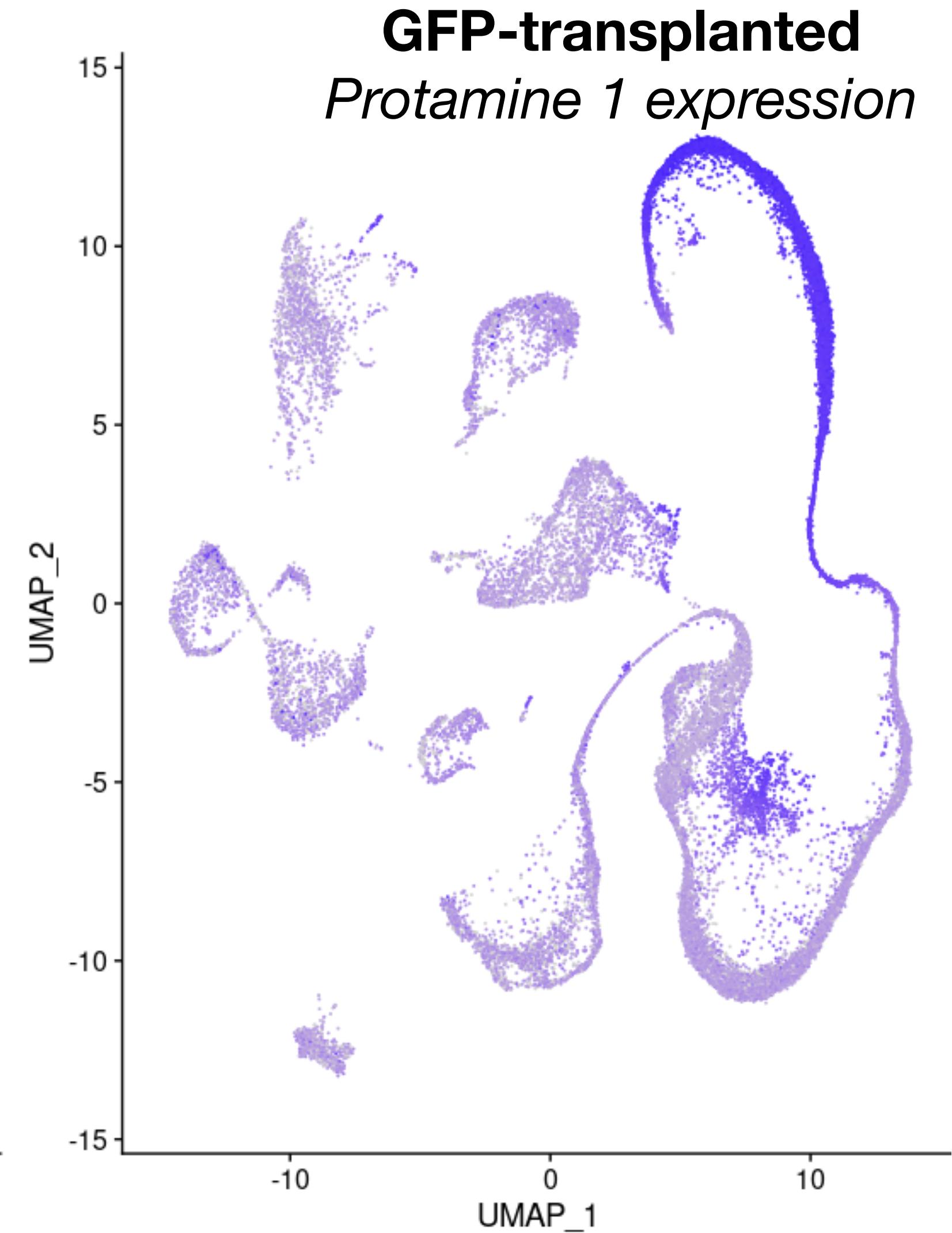
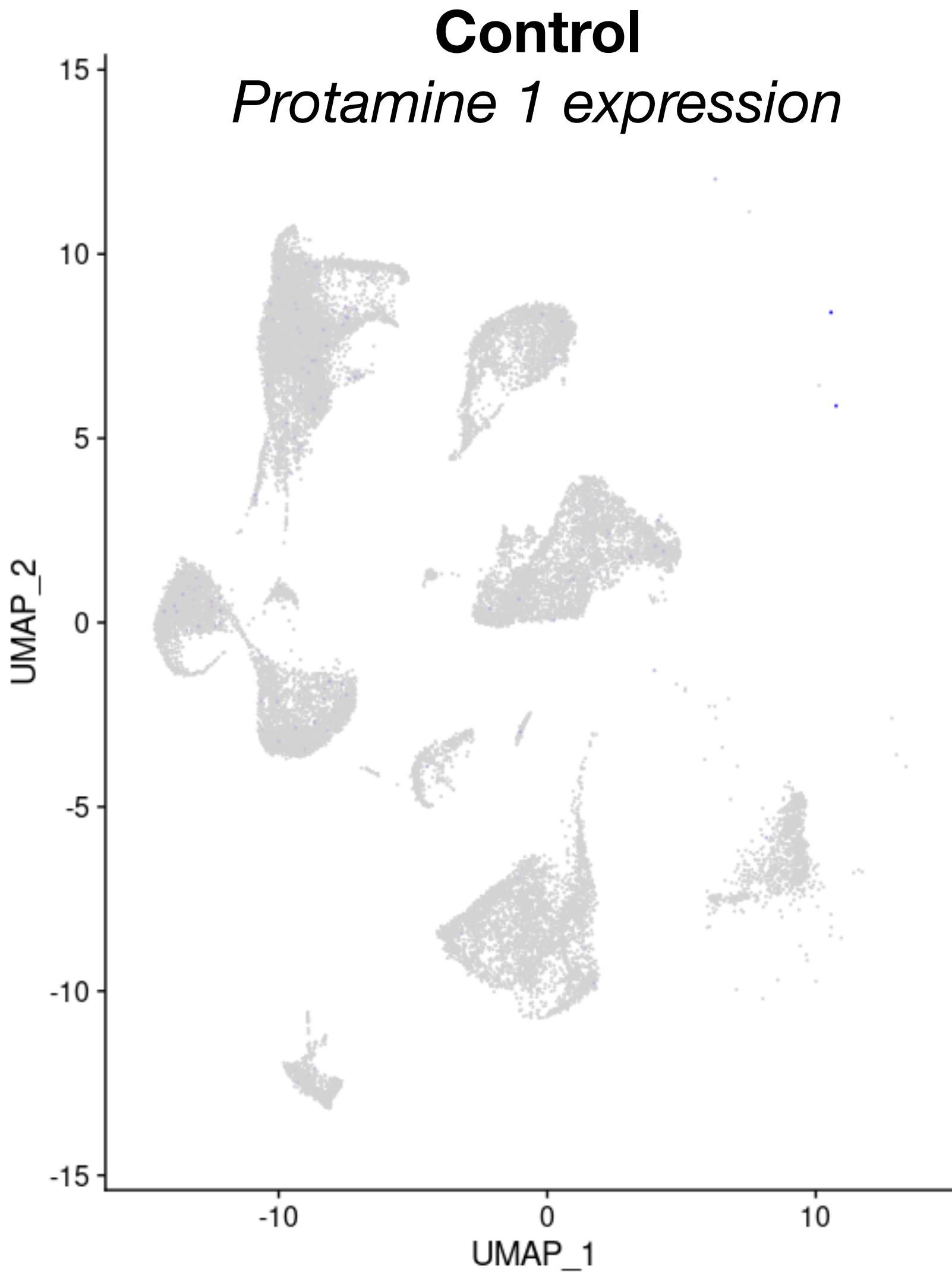
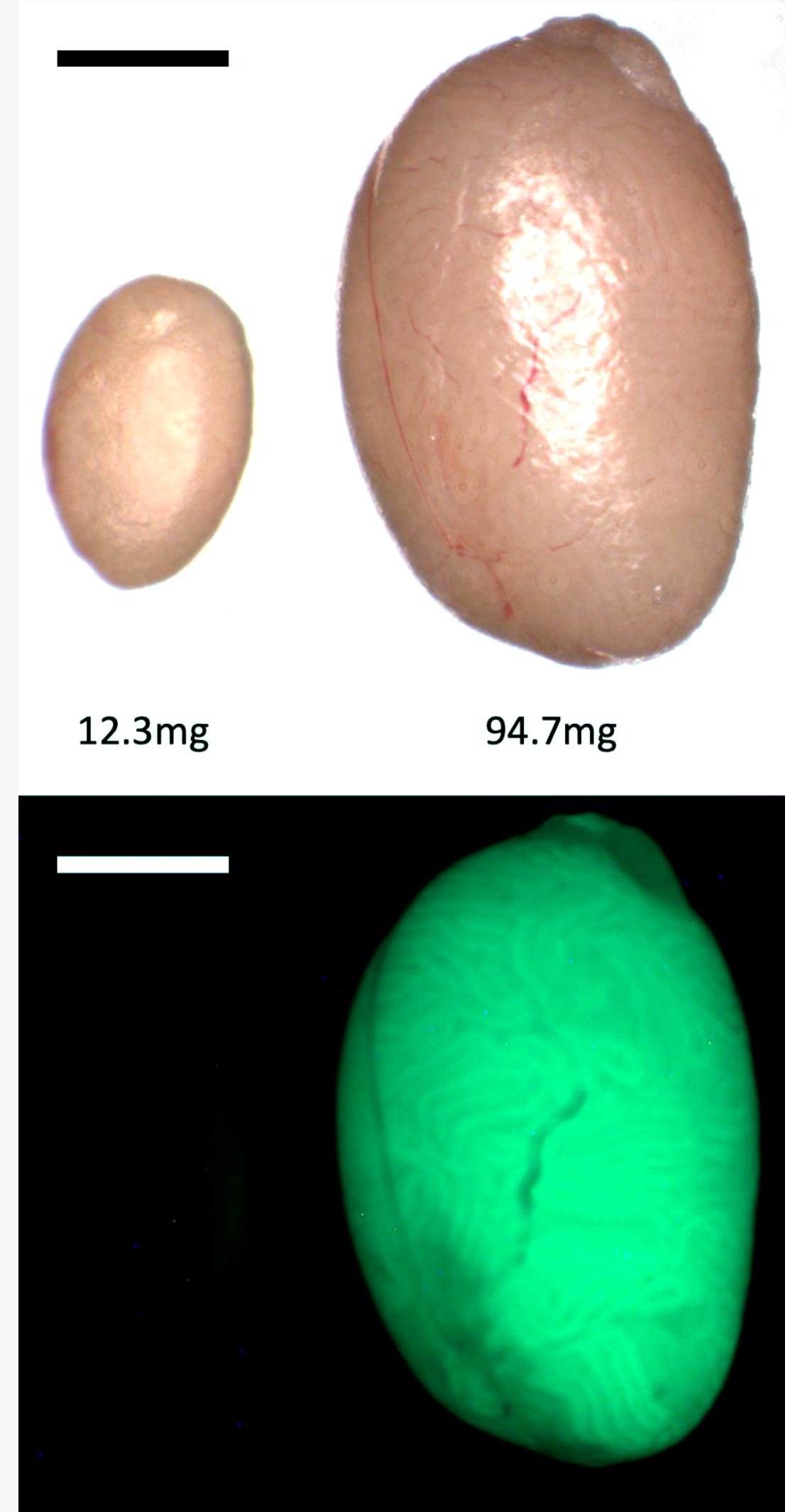
Cell viability is critical to the success of scRNA-seq.

Viability should be >90%

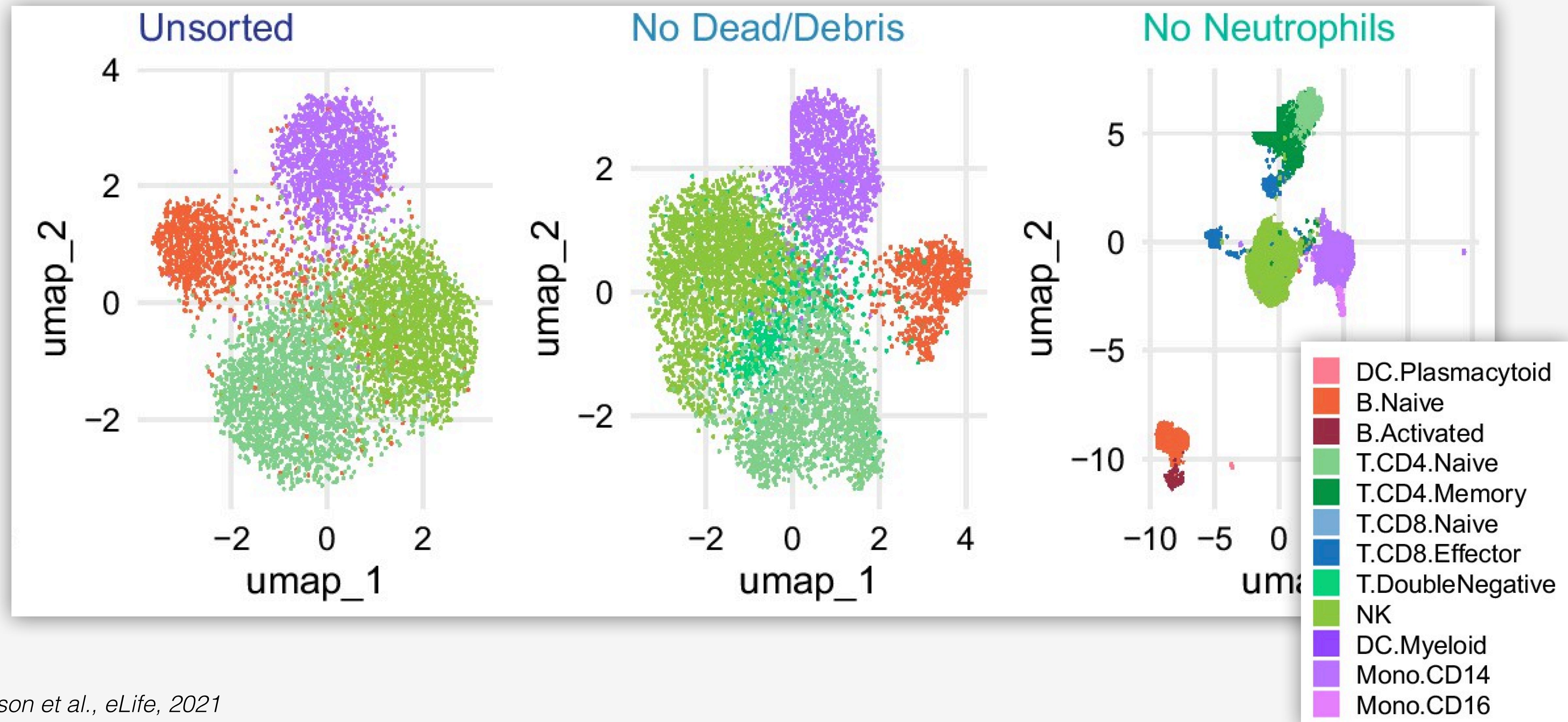


Video source: Instagram account of MicrobialEcology
<https://www.instagram.com/microbialecology/>

Ambient RNA is your enemy!



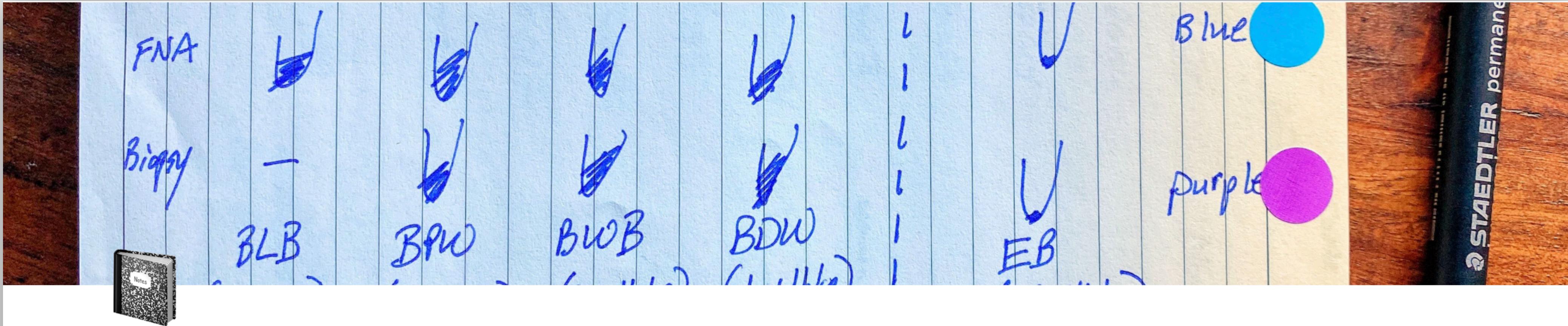
Certain cell types can be problematic in single cell assays



Recommendations for high quality cell preps

- Cell isolation protocols optimized for flow cytometry or *in vitro* functional assays may be insufficient for scRNA-seq experiments.
- If you are accustomed to digesting tissues at 37°C, consider switching to a 30min digestion at 6°C with a cold active protease. See [this paper](#)
- If you need a starting point for tissue dissociation protocols, use the [Worthinton guide](#)
- Avoid mechanical dissociation (e.g. bashing up lymph nodes with a syringe plunger).
- Regardless of whether you use bead-based purification or FACS, you should carry out a dead cell removal using [this kit](#) (or a similar one). Viability should be > 90%.
- Bead-based purification of cells is preferred over flow sorting, because it is usually faster and results in happier cells. Depending on your experimental goals, it may be totally fine if purity isn't great, since you'll get separation of cell populations in the scRNA-seq
- Keep everything cold, cold, cold. Cells should **always** be kept on ice....even when walking from the bench to the centrifuge. Buffers should be pre-chilled on ice. Centrifuges should be pre-chilled to 4deg. This will help to limit cell death and clumping.
- Avoid using EDTA in your tissue dissociation buffers. Excessive amount of EDTA (> 0.1mM) or magnesium (> 3mM) will inhibit reverse transcription reactions and lead to reduced complexity in the scRNA-seq dataset.

scRNA-seq protocols



Protocols

This website provides access to all standard protocols and procedures used at the PennVet Center for Host-Microbial Interactions (CHMI). If you're affiliated with CHMI or a member of the Beiting Lab, then you can collaboratively edit any of these protocols. This wiki-style approach gives us a central place to develop our lab protocols, while always keeping the most up-to-date protocol in one convenient and publicly accessible place. If you're not part of CHMI, no worries, you can still peruse the site to see our protocols. Comments and feedback are welcome!

Bulk RNA-seq and ATAC-seq

- [RNA-seq FAQS](#)
- [Stranded mRNA Prep, Ligation \(Updated TruSeq\)](#)
- [Stranded Total RNA Prep, Ligation with Ribo-Zero Plus \(Updated TruSeq\)](#)
- [Low input RNAseq \(SMART-Seq HT PLUS Kit\)](#)
- [RNA-seq on FFPE tissue \(Takara-Clontech total RNA pico mammalian v2/v3\)](#)
- [DRAFT miRNA Library Prep Protocol](#)

Microbial genomics

- [Microbiome FAQs](#)
- [16S seq](#)
- [NexteraFLEX](#)
- [16S Sanger](#)
- [Bioinformatics for microbial genomics](#)
- [Legacy micro protocols](#)

Practical considerations for single cell studies

Data analysis

System Requirements

Hardware

Cell Ranger pipelines run on Linux systems that meet these minimum requirements:

- 8-core Intel or AMD processor (16 cores recommended)
- 64GB RAM (128GB recommended)
- 1TB free disk space
- 64-bit CentOS/RedHat 6.0 or Ubuntu 12.04

 The minimum requirement of 64GB RAM will allow `cellranger aggr` to aggregate up to 250k cells, more memory will be required beyond that.



“Another challenge in scRNA-seq preprocessing is the amount of data that must be processed. A single-cell experiment can generate 10^6 – 10^{10} reads from 10^3 – 10^6 cells. This is leading to bottlenecks in analysis: for example, the current standard program for preprocessing 10x Genomics Chromium scRNA-seq, the Cell Ranger software, requires approximately 22 h to process 784 million reads using 1.5 Tb of disk space”



Modular, efficient and constant-memory single-cell RNA-seq preprocessing

Páll Melsted^{1,8}, A. Sina Booeshaghi^{2,8}, Lauren Liu³, Fan Gao^{4,5}, Lambda Lu⁴, Kyung Hoi (Joseph) Min^{1,6}, Eduardo da Veiga Beltrame⁴, Kristján Eldjárn Hjörleifsson³, Jase Gehring⁷ and Lior Pachter^{1,3,4}

We describe a workflow for preprocessing of single-cell RNA-sequencing data that balances efficiency and accuracy. Our workflow is based on the kallisto and bustools programs, and is near optimal in speed with a constant memory requirement providing scalability for arbitrarily large datasets. The workflow is modular, and we demonstrate its flexibility by showing how it can be used for RNA velocity analyses.

The quantification of transcript or gene abundances in individual cells from a single-cell RNA-sequencing (scRNA-seq) experiment is a task referred to as preprocessing¹. The preprocessing steps for scRNA-seq bear some resemblance to those used for bulk RNA-seq² and are in principle straightforward: cDNA reads originating from transcripts must be partitioned into groups according to cells of origin and aligned to reference genomes or transcriptomes to determine molecules of origin, and the reads, which originate from PCR-duplicated molecules, must be “collapsed” so that they are counted only once during quantification with unique molecular barcodes that serve as barcodes. Preprocessing scRNA-seq lies in determining choices for the various steps. For example, collapsing UMIs to be performed naively by associating each read with the same gene, with the same UMI, is computationally inefficient procedure is to collapse reads with identical UMIs that are associated with different genes. This is because the CR duplication of a single molecule can be relaxed, resulting in a complete (that is, computation-free) quantification. The amount of data that must be processed by the standard program for preprocessing scRNA-seq, the Cell Ranger software, requires approximately 22 h to process 784 million reads (Methods and Supplementary Table 1). For this reason, the Cell Ranger workflow corrects all barcodes that are one base-pair change away (Hamming distance 1) from barcodes in the whitelist. An examination of a benchmark panel of 20 datasets revealed that this error correction approach can be expected to rescue, on average, 0.8% of the reads in an experiment (Fig. 1a,b), a calculation based on an inferred error rate per base for each dataset (Methods and Supplementary Table 3). Thus, correction

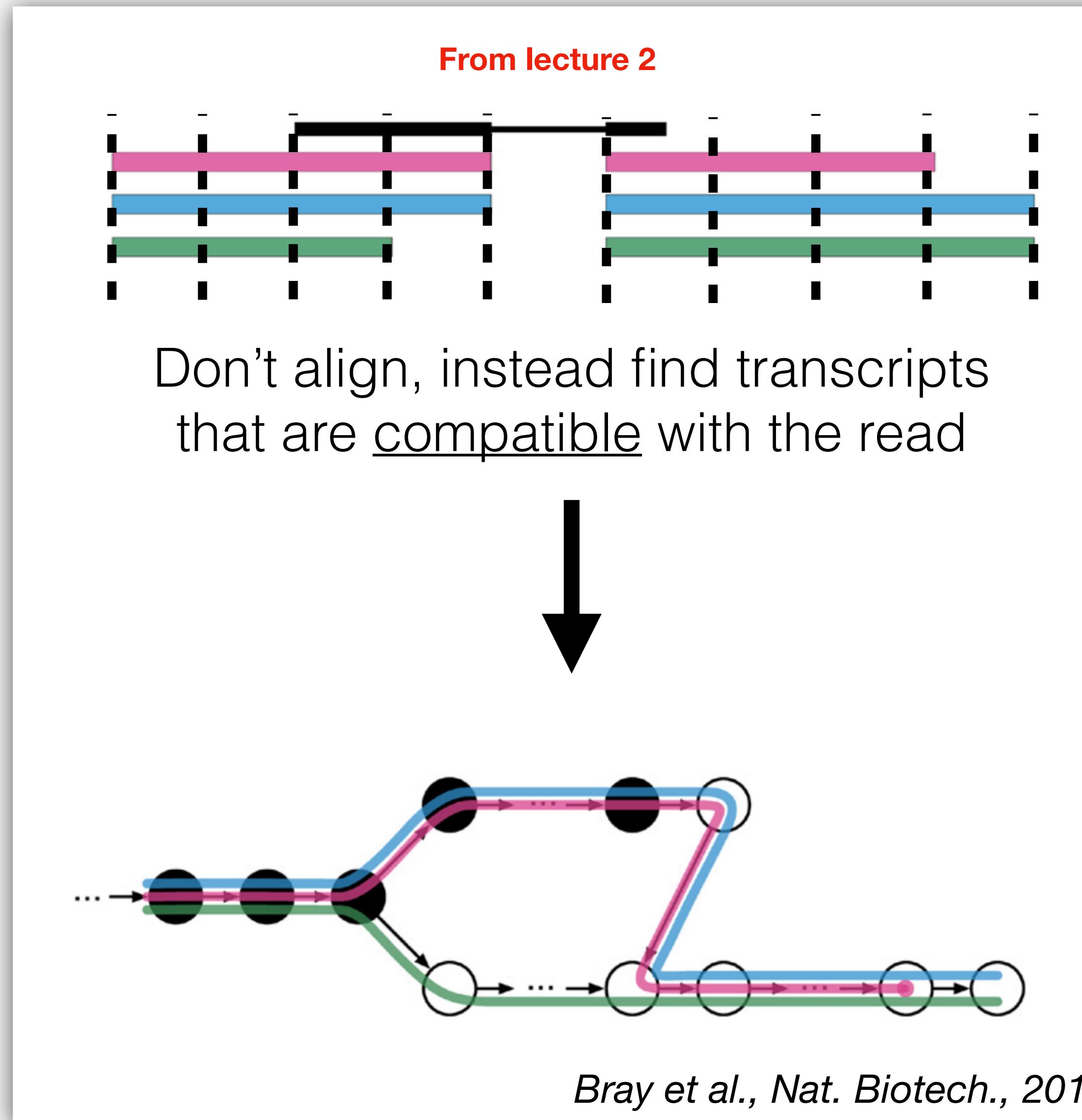
that is untenable given the pace of improvement in technology and the corresponding increase in data volume.

In recent work, we introduced a format for scRNA-seq data that makes possible the development of efficient workflows by virtue of decoupling the computationally demanding step of associating reads to transcripts and genes (alignment) from the other steps required for scRNA-seq preprocessing¹⁰. This format, called BUS (barcode, UMI, set), can be produced by pseudoalignment, and rapidly manipulated by a suite of tools called bustools (Supplementary Table 2). To illustrate the utility, efficiency and flexibility of this approach for scRNA-seq preprocessing, we describe a Chromium preprocessing workflow based on reasoned choices for the key preprocessing steps. While we focus on Chromium, our workflow is general and can be used with other technologies. We show that our preprocessing workflow is faster and has lower memory requirements than existing methods, and we demonstrate the power of modular processing with the BUS format by developing a fast RNA velocity analysis workflow¹¹. We also validate the design decisions underlying the Cell Ranger workflow. Our benchmarking and testing is comprehensive, comprising analysis of almost two dozen datasets and surpassing the scale of testing that has been performed for current workflows. Documentation and tutorials for the kallisto bustools workflow are available at <http://pachterlab.github.io/kallistobustools>.

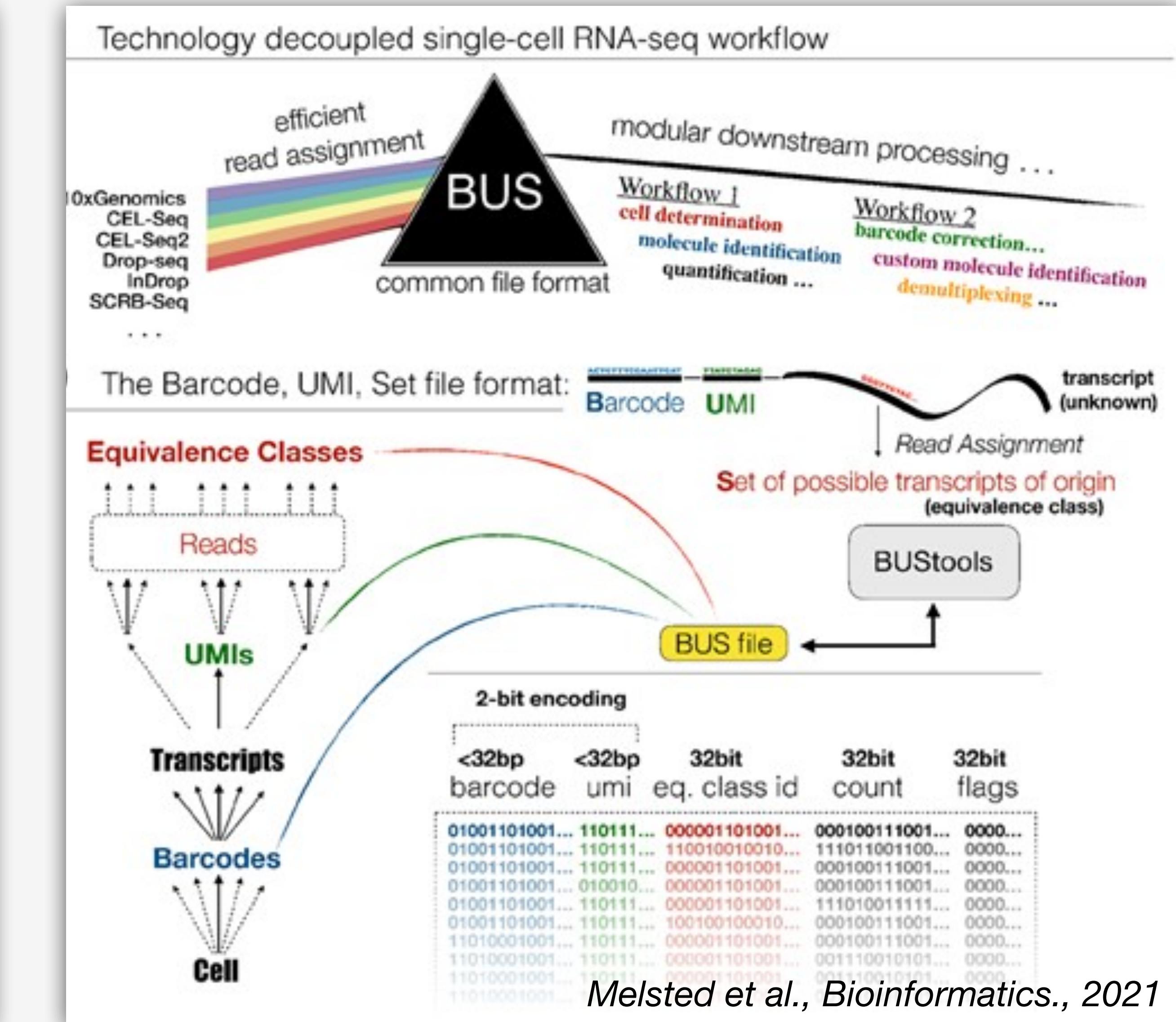
In designing an scRNA-seq preprocessing workflow, we began by investigating each required step: correction of barcodes, collapsing of UMIs and assignment of reads to genes. To achieve single-cell resolution, the Chromium technology produces barcode sequences that are used to associate cDNA reads to individual cells. We began by considering the efficiency–accuracy trade-offs involved in grouping reads with the same, or similar, barcodes to define the contents of individual cells. The Chromium barcodes arise from a ‘whitelist’, a set of pre-defined sequences that are included with the Cell Ranger software. Grouping reads by barcode is therefore straightforward, except for the fact that barcodes may contain sequencing errors. The Cell Ranger workflow corrects all barcodes that are one base-pair change away (Hamming distance 1) from barcodes in the whitelist. An examination of a benchmark panel of 20 datasets revealed that this error correction approach can be expected to rescue, on average, 0.8% of the reads in an experiment (Fig. 1a,b), a calculation based on an inferred error rate per base for each dataset (Methods and Supplementary Table 3). Thus, correction

based on pseudoalignment⁸ have recently been developed^{12,13}. However, despite improvements in running time, current workflows have memory requirements that increase with data size⁵, a situation

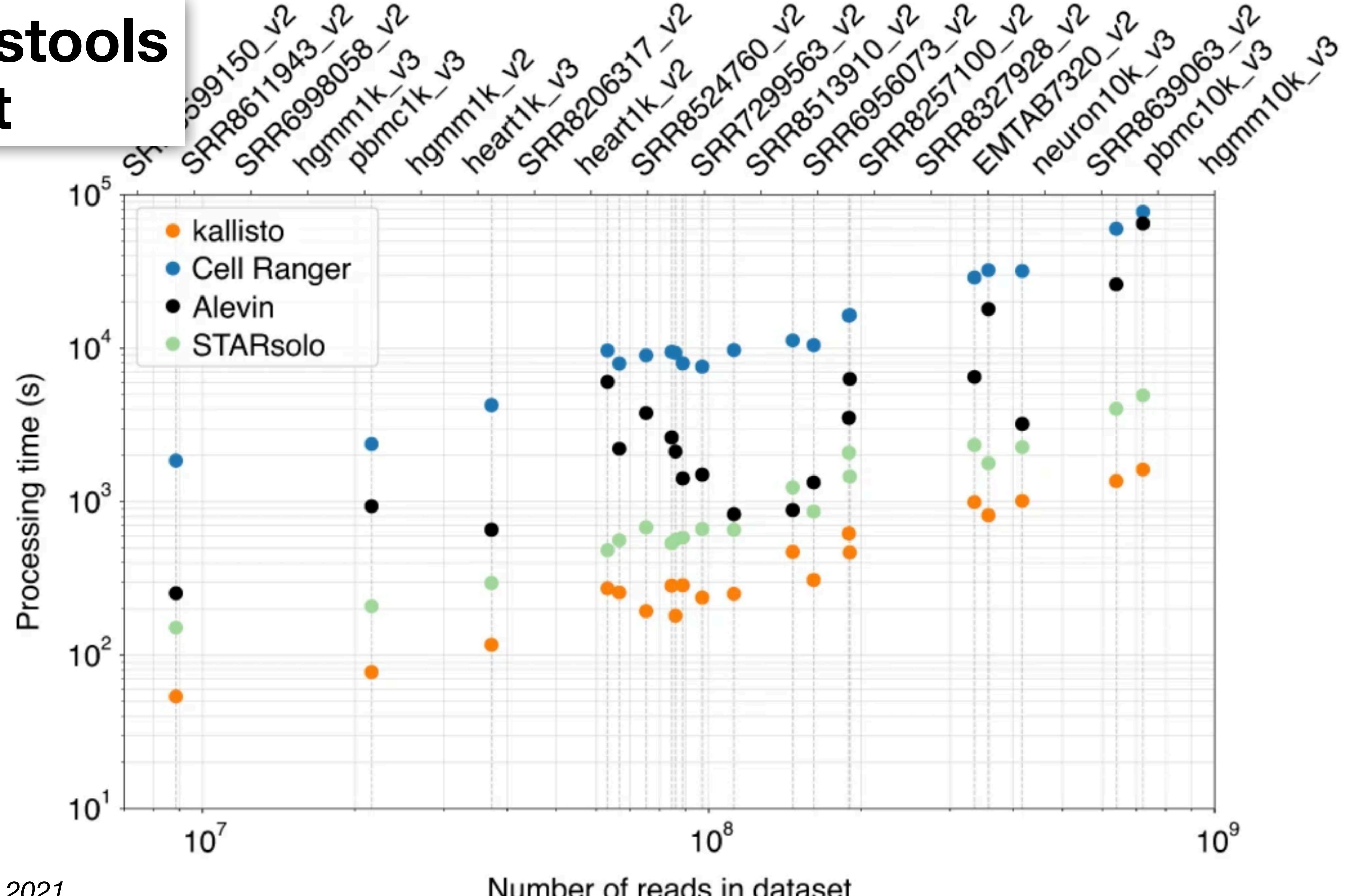
Kallisto pseudoalignment



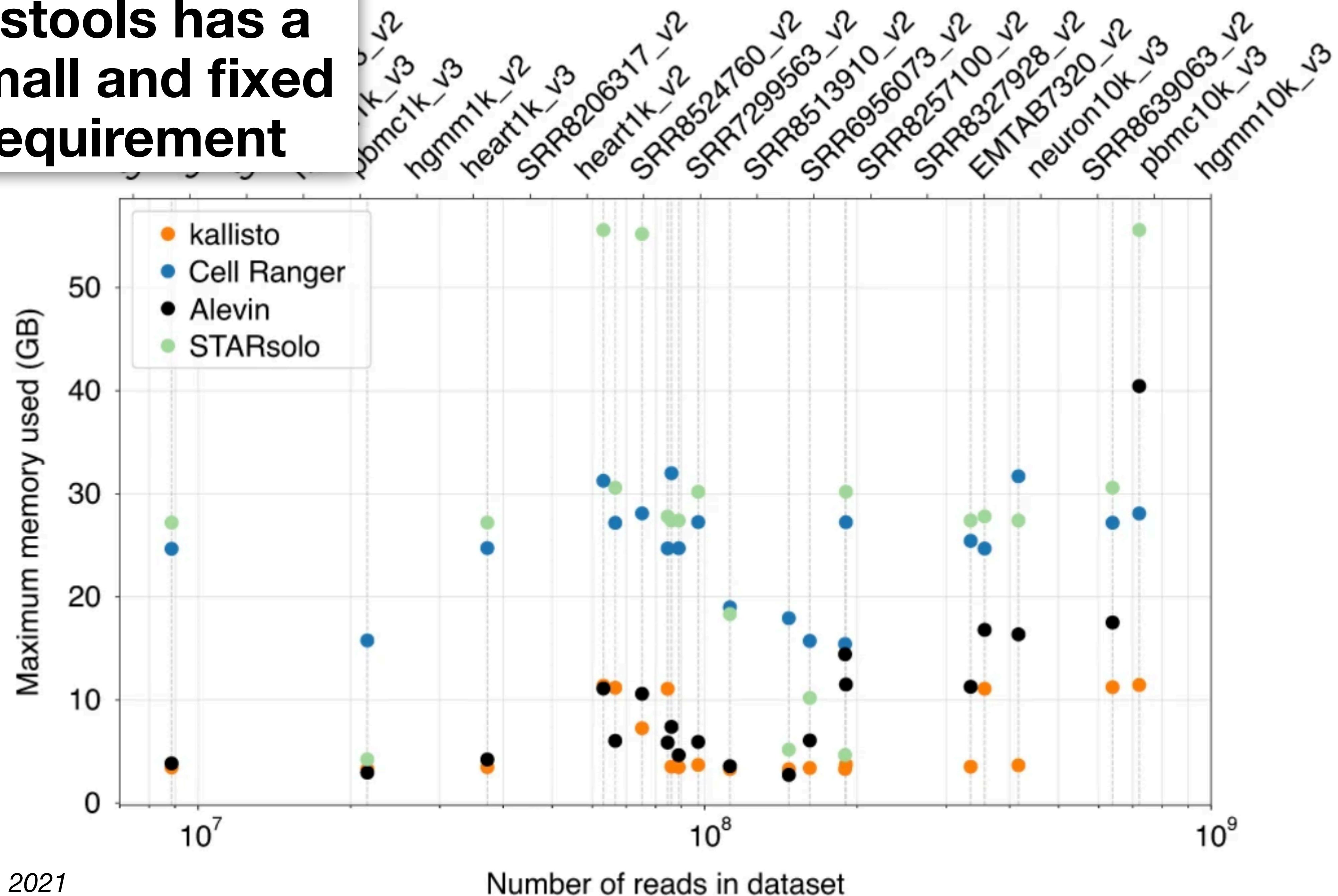
Barcode, UMI, Set (BUS)



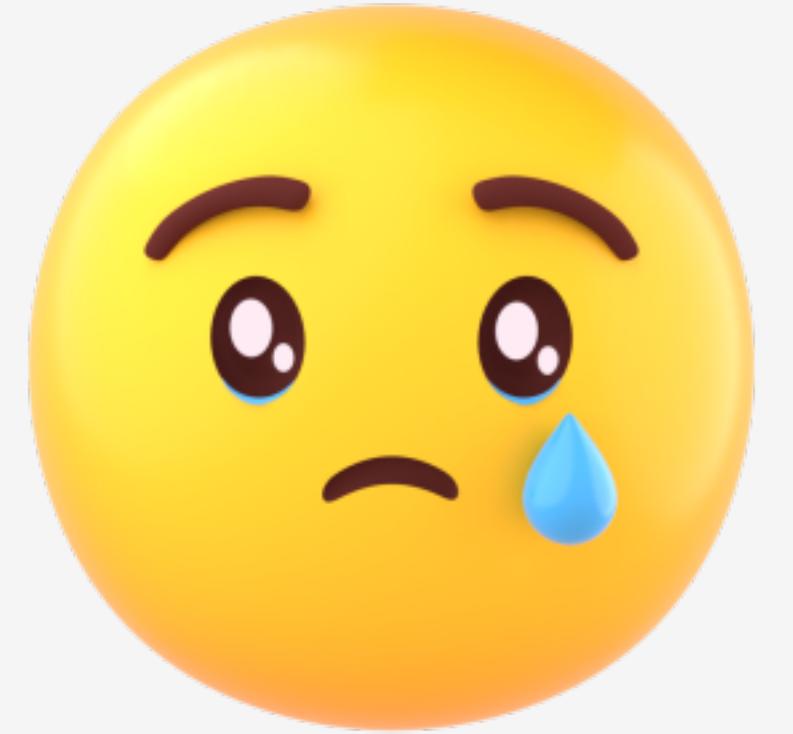
Kallisto Bustools is fast



Kallisto Bustools has a relatively small and fixed memory requirement



**Let's get started using Kallisto bustools to preprocess
raw scRNA-seq data directly on our laptop**



Say goodbye to our course dataset

DIY.transcriptomics

Lectures Labs **Data** Scripts Videos FAQs Help

Data

Access to fastq files and info for the dataset used for the course.

Course dataset

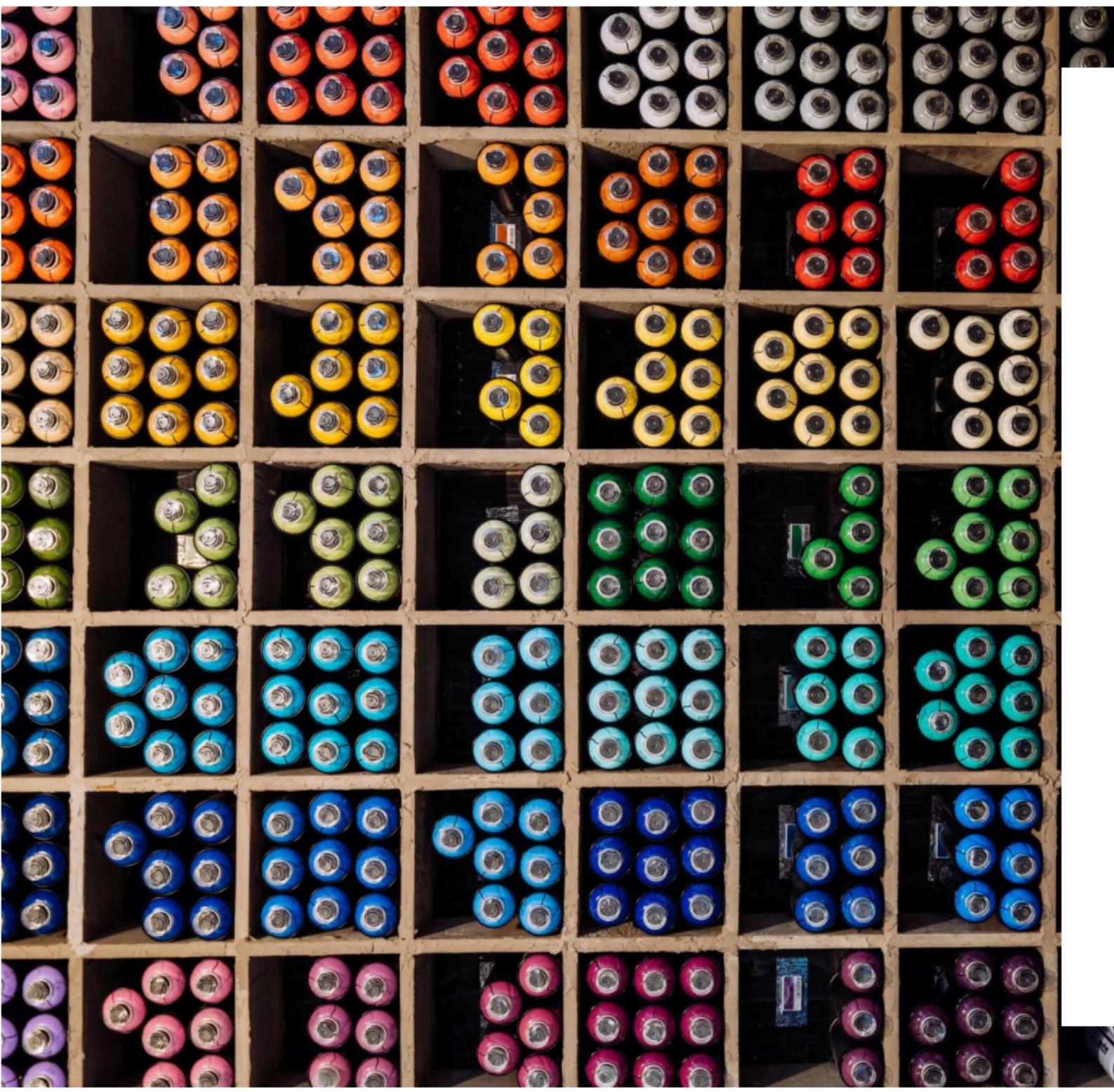
The course dataset comes from a collaboration between my lab, Phil Scott's lab, and our colleagues in Salvador, Brazil ([manuscript here](#)). This is an RNAseq dataset from skin biopsies obtained from patients with cutaneous leishmaniasis, a parasitic disease endemic in Brazil and other areas of South America. You'll be working with data from 5 patients with this disease and 5 healthy endemic controls.

This course offers multiple 'on-ramps', or entry points, for starting an analysis. Below are a few of these options.

Option 1 (preferred)- Download raw data and accessory files

fastq files – You will need about 30Gb of storage space on your harddrive to accomodate this download. *please do not uncompress these files (leave them as .gz files).*

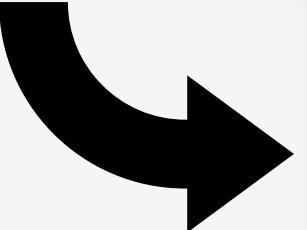
If you're unable to download the files above, you can always get them from [here](#) where these files were deposited during publication. Be sure to select only the specific files we'll be using during the course, which you can see [here](#).



Single cell RNA-seq – principles and processing

Lecture 13 • watch by December 1, 2021

Now that you're comfortable with bulk RNA-seq data analysis, we'll shift our focus to the rapidly developing landscape of single cell RNA-seq (scRNA-seq). In this lecture, you'll learn about the underlying technology and demonstrate how to process raw scRNA-seq data directly on your laptop (!) for importing into R/bioconductor.



Single cell RNA-seq – principles and processing

Lecture 13 • watch by December 1, 2021

Overview

Now that you're comfortable with bulk RNA-seq data analysis, we'll shift our focus to the rapidly developing landscape of single cell RNA-seq (scRNA-seq). In this lecture, you'll learn about the underlying technology and we'll demonstrate how to process raw single cell data directly on your laptop (!) for importing into R/bioconductor.

Learning objectives

- Understand droplet-based scRNA-seq technology
- Be able to compare and contrast single cell and bulk RNA-seq methods
- Understand cost and experimental design considerations for scRNA-seq experiments.
- Familiarity with multiplexed single cell assays (CITE-seq, 'multiome', TEA-seq)
- Be able to define common terms and concepts in single cell genomics
- Use Kallisto-BUSTools to preprocess raw scRNA-seq data (via `kb-python`)

What you need to do

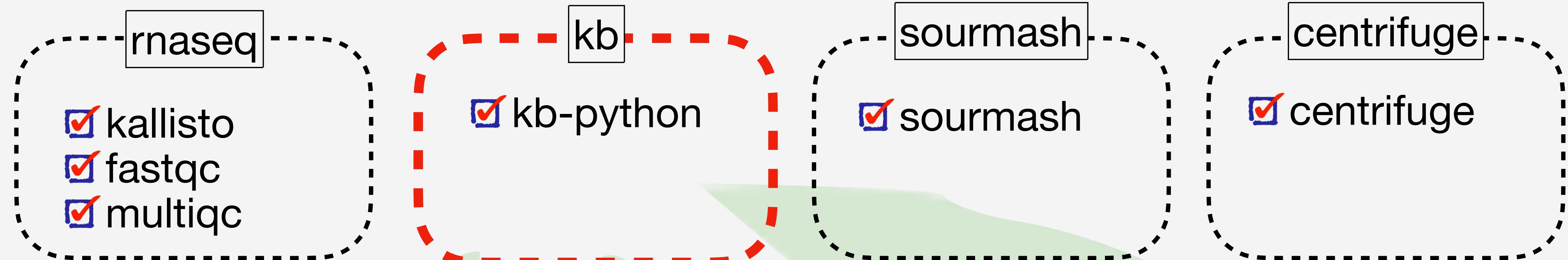
Download raw files. You will need about 5Gb of storage space on your harddrive to accomodate this download. *please do not uncompress these files (leave them as .gz files).* This is data from 1000 peripheral blood mononuclear cells (PBMCs) and is one of the sample datasets provided by 10X Genomics [here](#). I merged the separate lane files to make this simpler to work with for the course.

Human transcriptome reference index file - this is the index you created using Kallisto way back in lecture 2. If you don't have this, remember it's easy to create using `kallisto index`.

t2g.txt - this is a human transcript-to-gene mapping file that we will use with Kallisto-BUSTools to preprocess our data. This file is easy to generate with `kb ref`, but downloading it now will save you some time.

kb-python - You will need to have this software installed in a Conda environment on your laptop. We did this way back in [lecture 1](#). If you are unable to install or use kb-python, just follow along with the lecture so you understand the concepts.

Our Conda software environments



First, activate the environment with *conda activate kb*

Then, call the program with *kb*

Kallisto-Bustools: supported technologies

kb --list

List of supported single-cell technologies					
name	PennVet	whitelist	barcode	umi	cDNA
10XV1		yes	0,0,14	1,0,10	2,None,None
10XV2		yes	0,0,16	0,16,26	1,None,None
10XV3		yes	0,0,16	0,16,28	1,None,None
CELSEQ			0,0,8	0,8,12	1,None,None
CELSEQ2			0,6,12	0,0,6	1,None,None
DROPSEQ			0,0,12	0,12,20	1,None,None
INDROPSV1			0,0,11 0,30,38	0,42,48	1,None,None
INDROPSV2			1,0,11 1,30,38	1,42,48	0,None,None
INDROPSV3		yes	0,0,8 1,0,8	1,8,14	2,None,None
SCRUBSEQ			0,0,6	0,6,16	1,None,None
SURECELL			0,0,6 0,21,27 0,42,48	0,51,59	1,None,None
SMARTSEQ					0,None,None 1,None,None

Preparing reference files scRNA-seq

Calling the program and function

Download a
prebuilt index
(human or mouse)

Path/name of
index file to be
created

Path/name of
transcript-to-gene
mapping file to be
created

Indicates code continues on next line

```
• kb ref \
• -d human \
• -i Homo_sapiens.GRCh38.cdna.all.index \
• -g t2g.txt
```

You do not need to run this code!
I provide the t2g.txt file in the downloads for this lecture

Preprocessing scRNA-seq data

Calling the program and function

Path for Fwd and
Rev fastq files

Path for reference index

Technology used

Path for transcript-to-gene
mapping file

Threads to use

Convert counts to cellranger-
compatible format

Indicates code continues on next line

```
• kb count \
• file1.fastq.gz file2.fastq.gz \
• -i Homo_sapiens.GRCh38.cdna.all.index \
• -x 10XV3 \
• -g t2g.txt \
• -t 8 \
• --cellranger
```