# The Harvard Chan Bioinformatics Core

Shannan Ho Sui
February 14, 2020
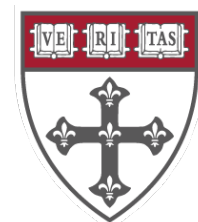
*Webpage:* http://bioinformatics.sph.harvard.edu

*Email:* bioinformatics@hsph.harvard.edu

HARVARD T.H. CHAN
SCHOOL OF PUBLIC HEALTH

**HARVARD**
**T.H. CHAN**
SCHOOL OF PUBLIC HEALTH

Shannan Ho Sui
*Director*

John Hutchinson
*Associate Director*

Victor Barrera

Rory Kirchner

Zhu Zhuo

Peter Kraft
*Faculty Advisor*

Preetida Bhetariya

Meeta Mistry

Mary Piper

Jihe Liu

Radhika Khetani
*Training Director*

Ilya Sytchev

James Billingsley

Sergey Naumenko

Joon Yoon

Maria Simoneau

# Consulting services

- Transcriptomics: RNA-seq, small RNA-seq, single cell RNA-seq

- DNA accessibility and binding: ChIP-seq, ATAC-seq

- Genetic variation: WGS, re-sequencing, exome-seq, structural variation, CNV

- Genome-wide methylation: 450k methylation arrays, RRBS, WGBS

- Data integration

- Functional enrichment analysis

- Experimental design and grant support

# Training

- Basic Data Skills

- Advanced NGS Data Analysis

- In-depth courses (8- to 12-days)

- Free, Monthly, 2 - 3 hour long workshops on various bioinformatics topics
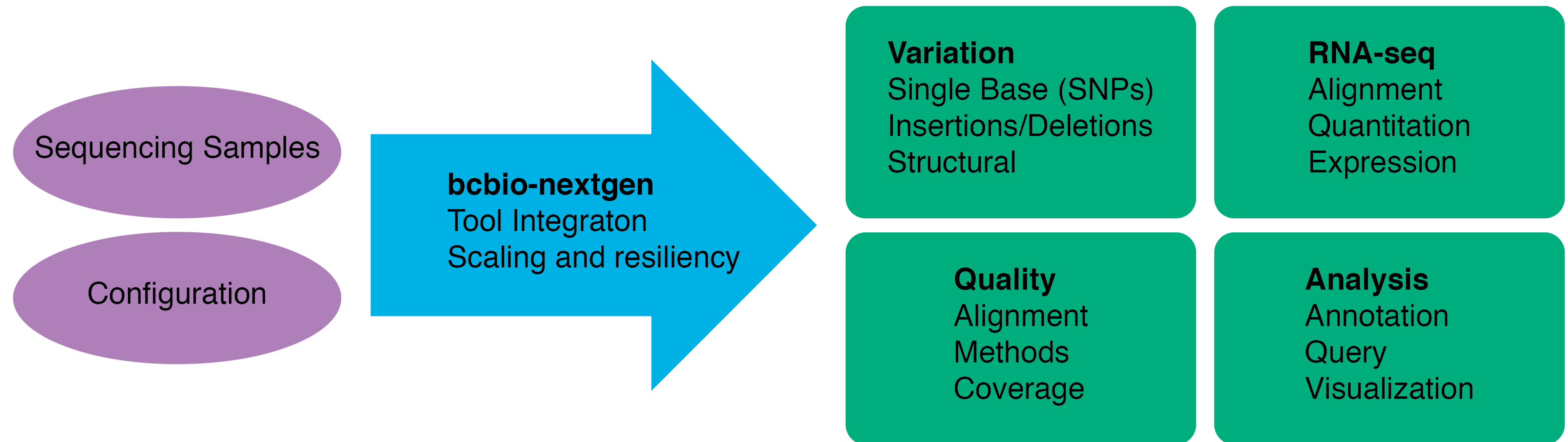
- Lessons: https://hbctraining.github.io/main

Sign up for our mailing list at - https://tinyurl.com/HBC-list-subscribe

**HARVARD T.H. CHAN**
SCHOOL OF PUBLIC HEALTH

# Standardizing using reproducible, scalable, validated best practice workflows

Sequencing Samples

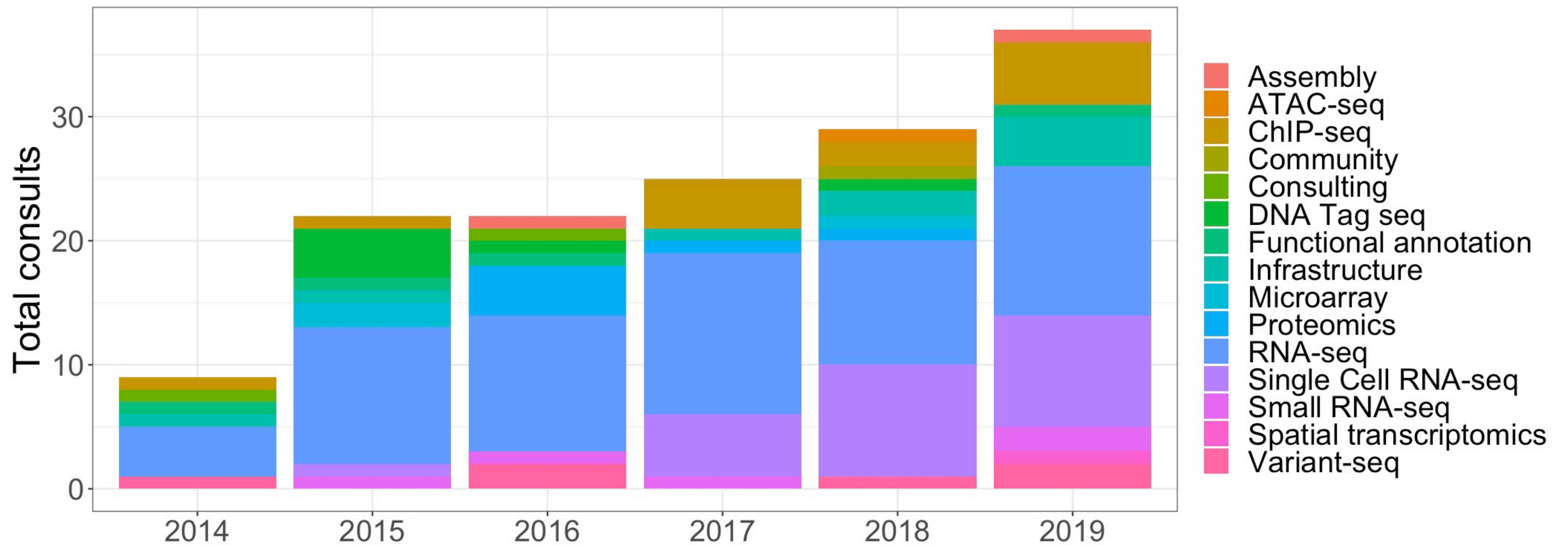Configuration

**bcbio-nextgen**
Tool Integraton
Scaling and resiliency

**Variation**
Single Base (SNPs)
Insertions/Deletions
Structural

**RNA-seq**
Alignment
Quantitation
Expression

**Quality**
Alignment
Methods
Coverage

**Analysis**
Annotation
Query
Visualization

# bcbio-nextgen

Python toolkit to automate best practice NGS pipelines

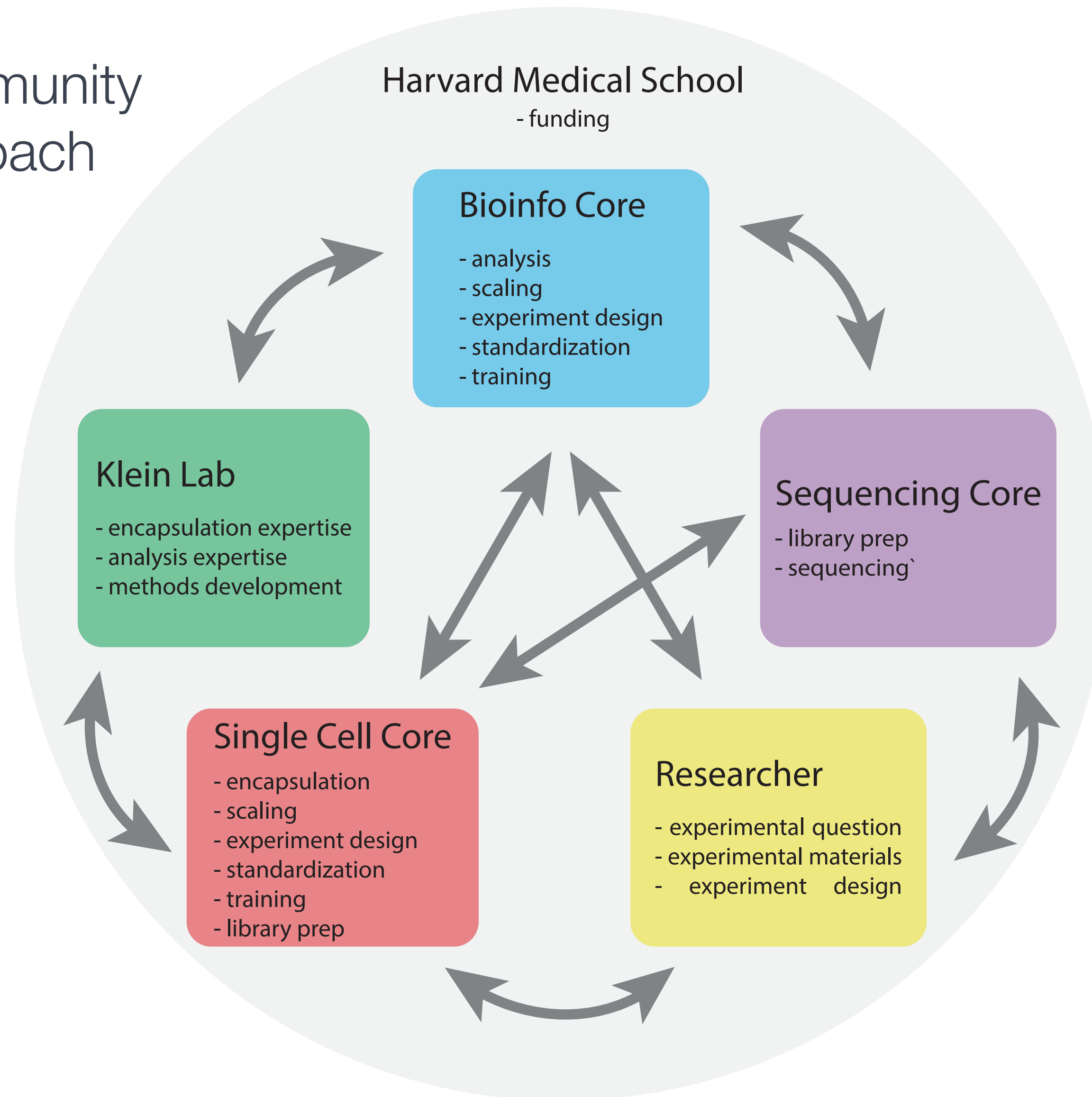Who we support

# Growing demand for service and new techologies
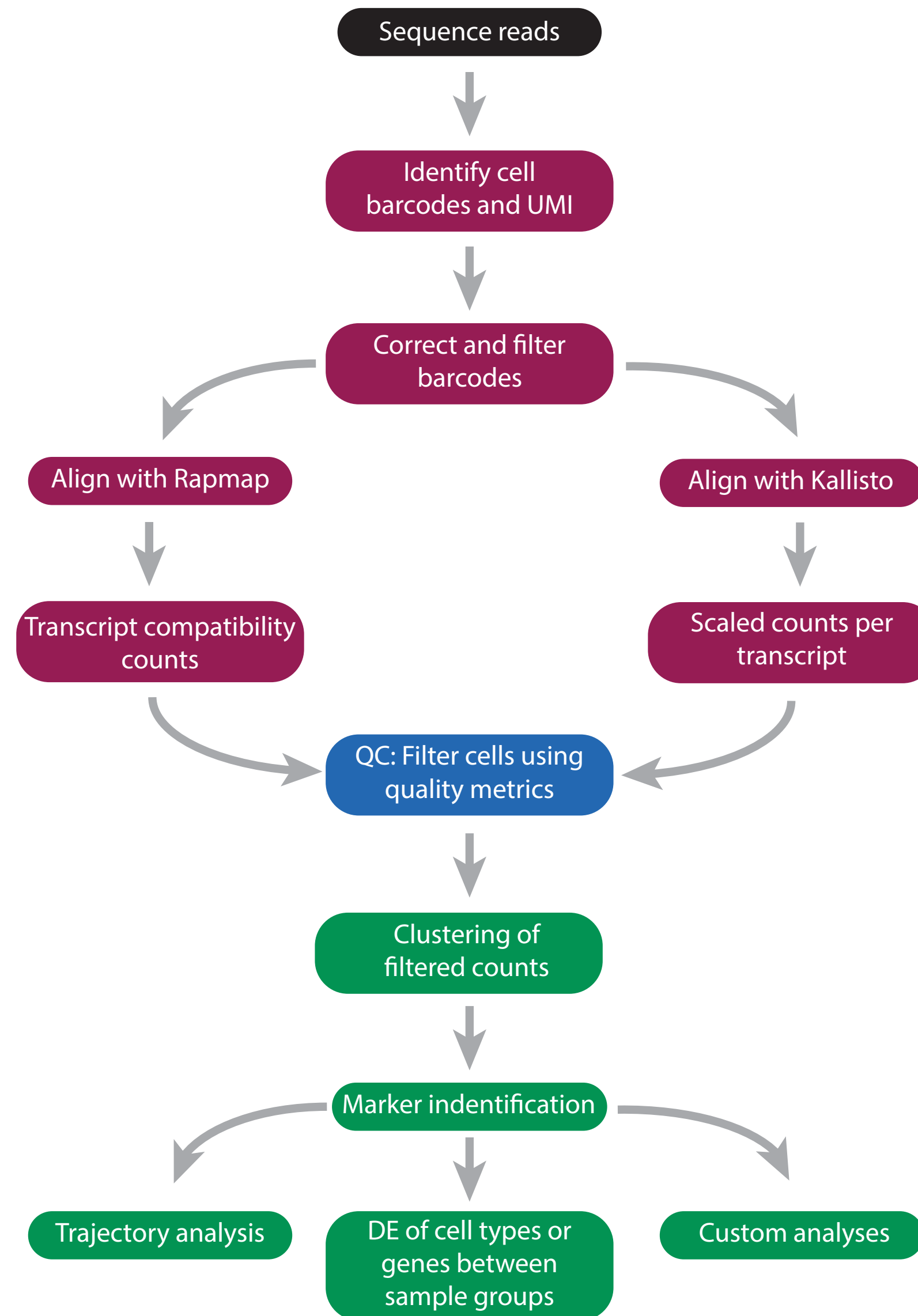
# Common applications of scRNA-seq

▶ Explore which cell types are present in a tissue

▶ Identify unknown/rare cell types or states

▶ Elucidate the changes in gene expression during differentiation or across time or states

▶ Identify genes that are differentially expressed in particular cell types between conditions (e.g. treatment or disease)

▶ Explore changes in expression among a cell type while incorporating spatial, regulatory, and/ or protein information
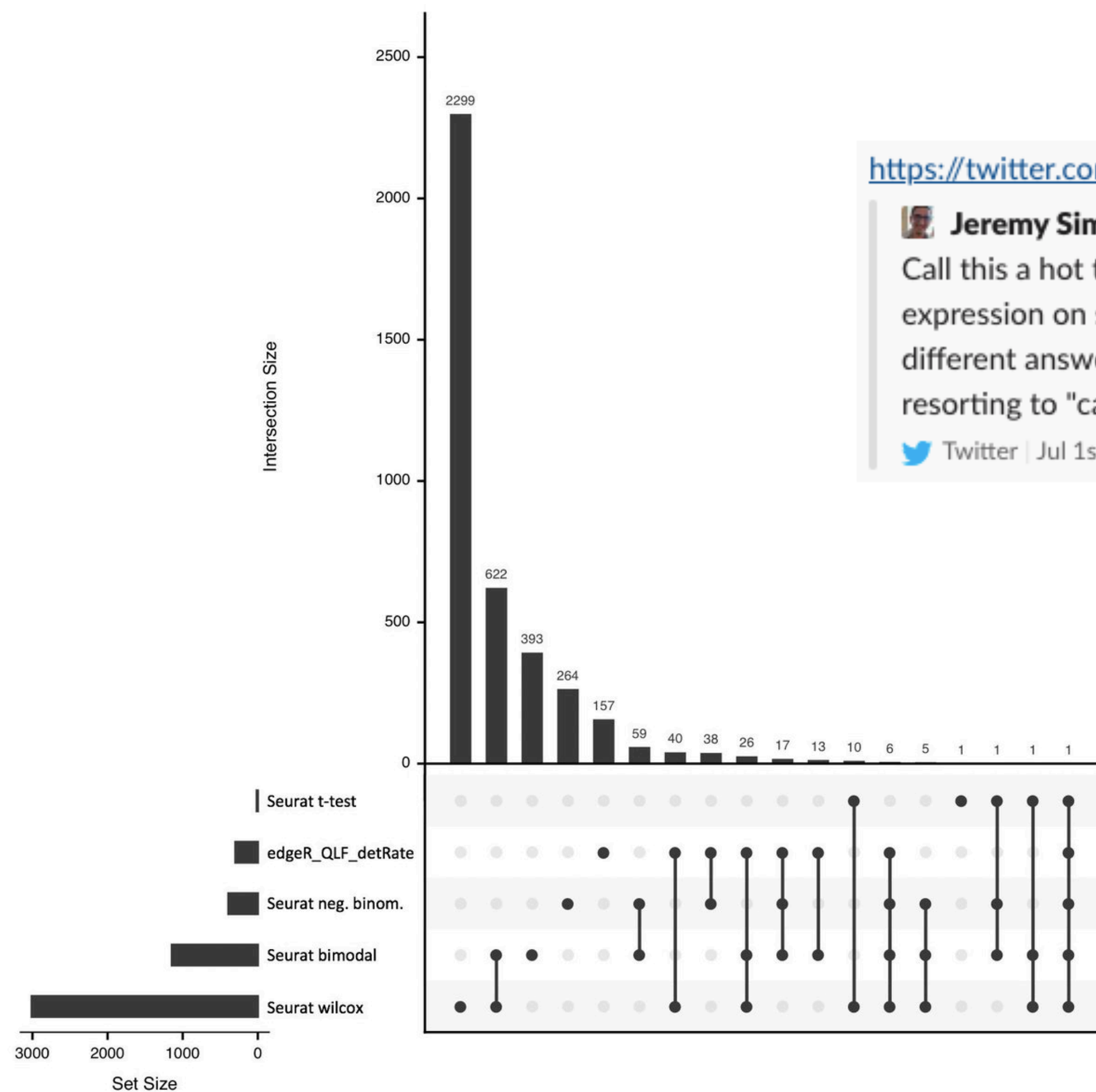
# A Community Approach

Harvard Medical School
- funding

**Bioinfo Core**
- analysis
- scaling
- experiment design
- standardization
- training

**Klein Lab**
- encapsulation expertise
- analysis expertise
- methods development

**Sequencing Core**
- library prep
- sequencing`

**Single Cell Core**
- encapsulation
- scaling
- experiment design
- standardization
- training
- library prep

**Researcher**
- experimental question
- experimental materials
- experiment design

Our single cell workflow

# Quality control

▶ Focus on getting high quality data, messy data is EXTREMELY hard to work with

▶ Cells must be alive, before running the experiment shoot for viability > 95%

▶ When filtering, it is better to start out too strict than too lax

▶ Essential to work closely with the biologists, many, many judgement calls need to be made based on expert knowledge

▶ Before beginning analysis, have a good list of marker genes in hand for each cell type that could be present in the sample

▶ Often markers used for FACS sorting are not good for single-cell RNA-seq

# Challenges and Opportunities

▶ Complex designs - replicates, batches, technologies

▶ Close collaborations to allow for rapid, iterative analyses

▶ Rapidly emerging methods and evolving tools

▶ Which ones to use?

▶ Keeping versions consistent/synchronized (esp. among computing environments)

▶ Different results from different methods

▶ Lots of open questions

12

# What is N in single cell experiments?

▶ Three treated patients, three control patients.

▶ Extract PBMC and want to look at the effect of the treatment on B cells.

▶ Identified B-cell clusters in the treated and non-treated patients via marker genes and found 300 B cells in each patients for a total of 900 B cells in each treatment condition

▶ If I want to ask what the effect of the treatment is on B-cells, what is my N here? Is it 900 for each condition? No. But almost all single-cell papers to date (including ours!) treat it as if it is.

▶ N should be 3, not 900.

# How to get to N=3?

▶ Pseudobulk: sum all B-cells for each sample, and treat it like an *in-silico* FACS sorted experiment

▶ multilevel model: model the patient level data in a multilevel model so you can account for the non-independence of measurements from B-cells of the same patient
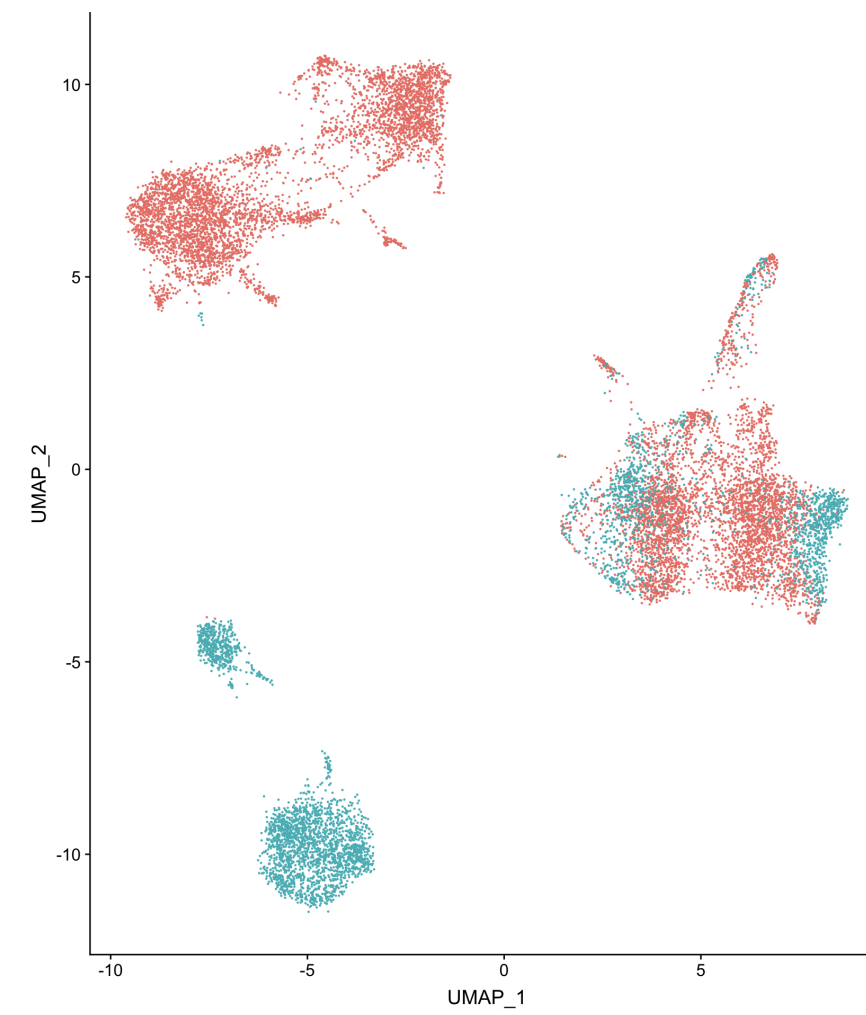
# pseudobulk is simple and works well

# Challenges and Opportunities

▶ Projects take longer to complete

▶ Practical approach to training

▶ Internal training through retreats, development of materials, group discussions

▶ Community training through our Bioinformatics Training Program
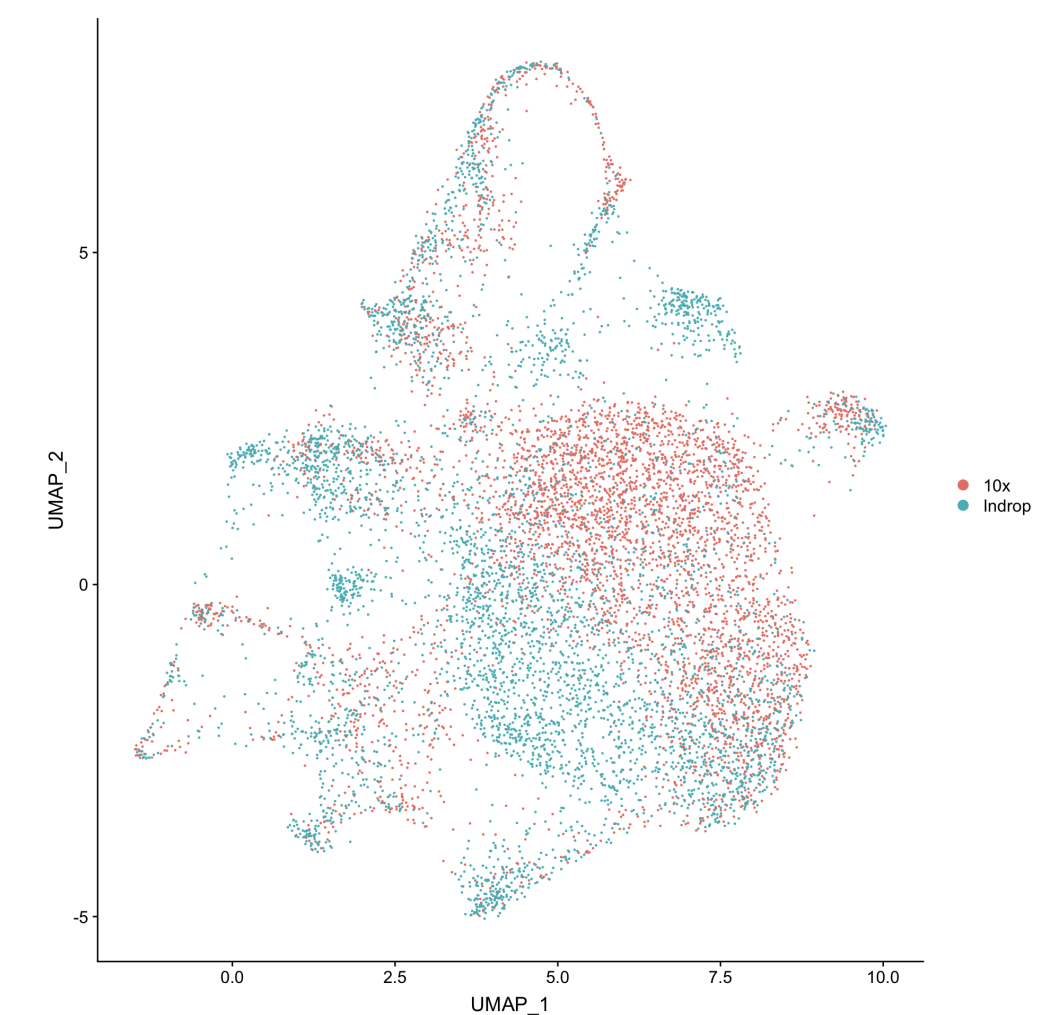
# Seurat v.3.0



State — injured / uninjured
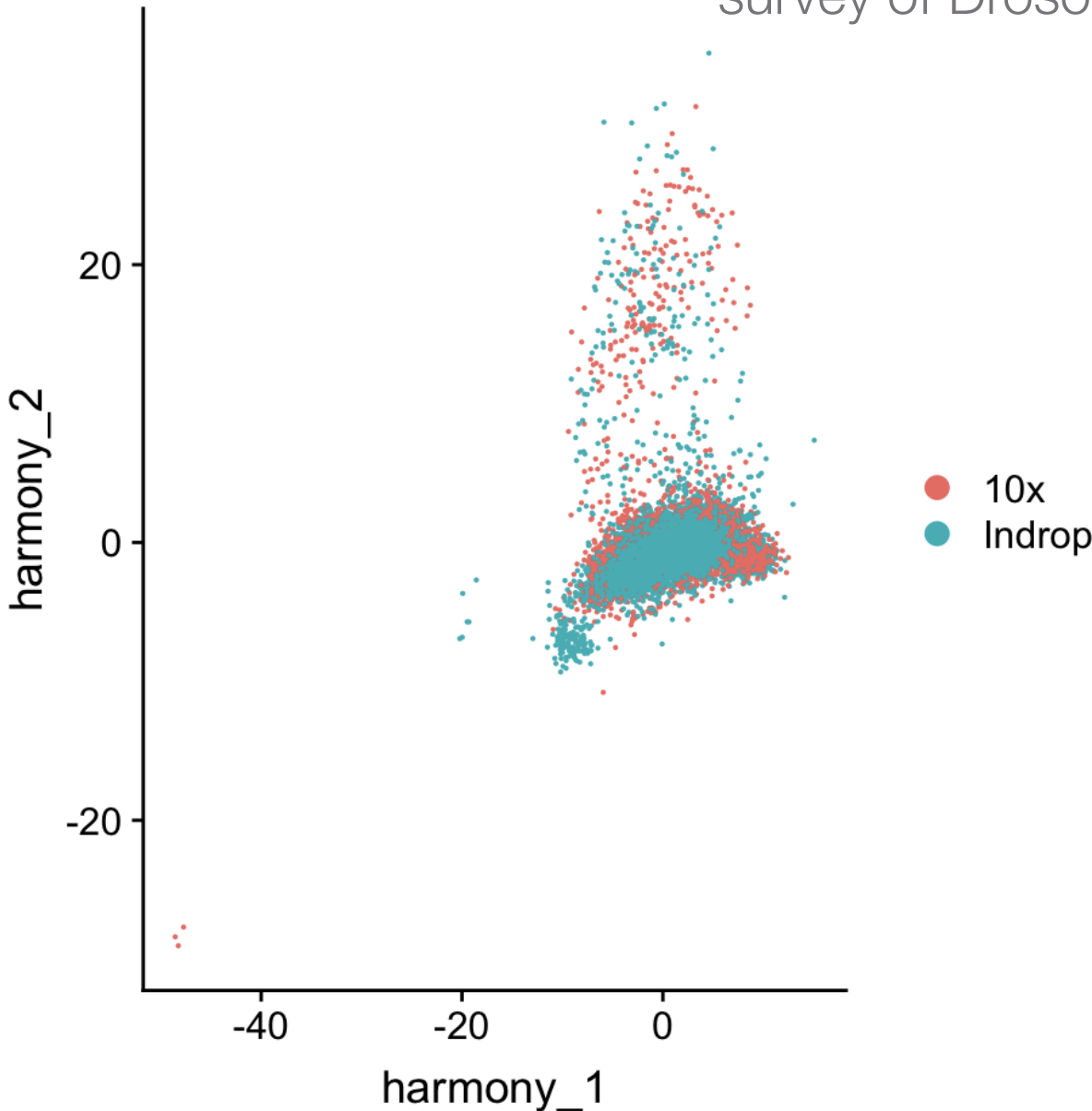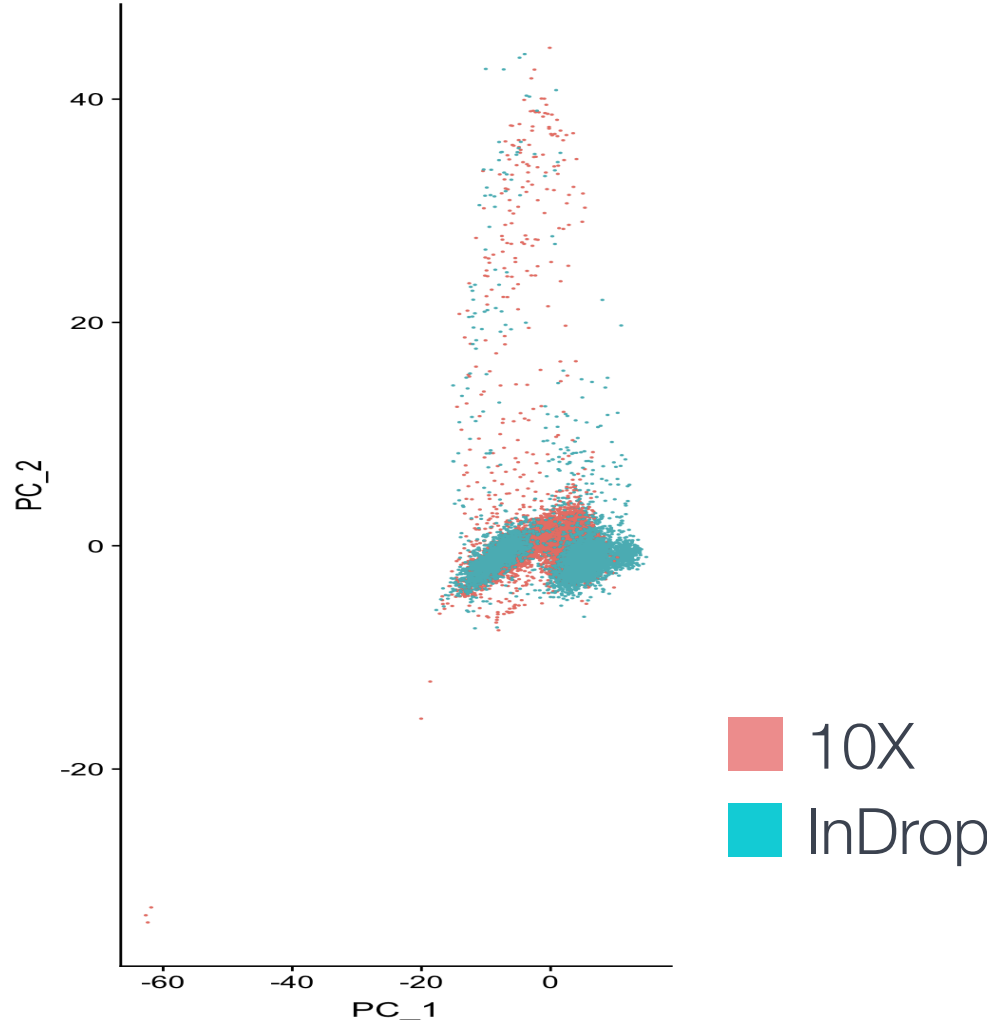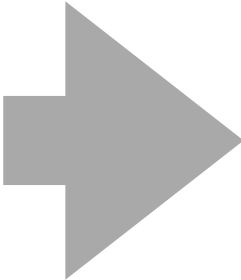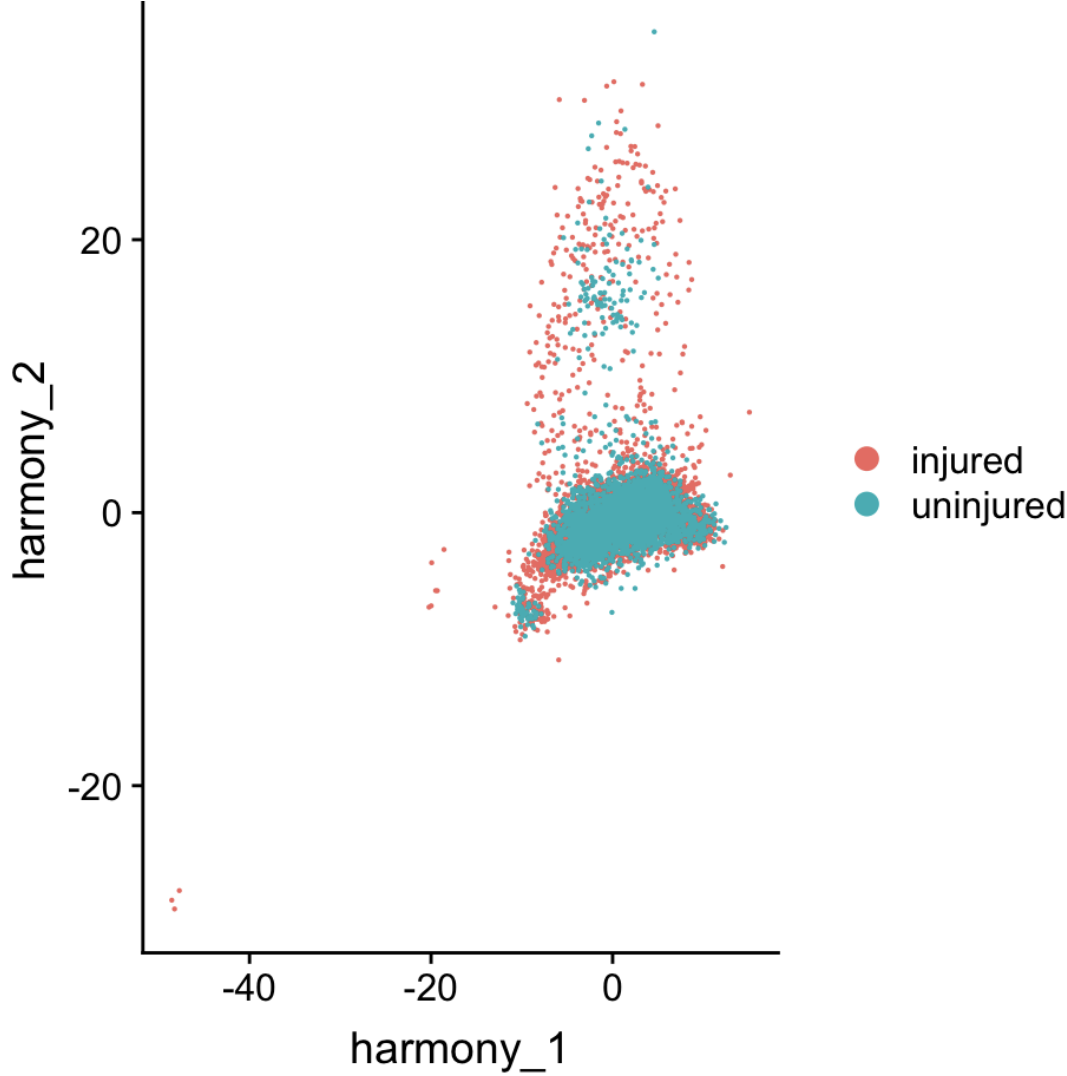
Technology — 10X / InDrop

Multi-CCA — 10X / InDrop

Technology

Injury state

Harmony

19

# Trajectory analysis using Slingshot



A

Cell type
- ISC/EB
- EE
- dEC
- aEC
- mEC
- pEC

B lineage1: ISC/EB -> mEC -> dEC -> aEC

pseudotime
0  5  10  15

C lineage2: ISC/EB -> mEC -> dEC -> pEC

pseudotime
0  5  10  15  20
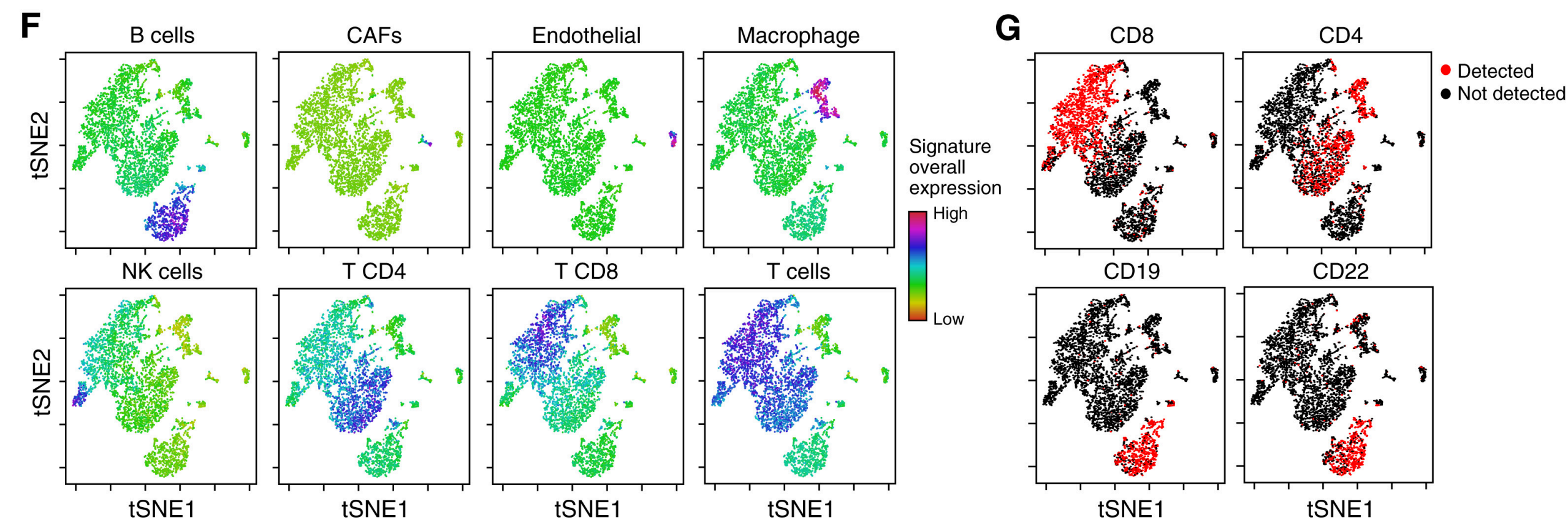
Hung et al., 2019, A cell atlas of the Drosophila midgut, PNAS.

# New technologies we're preparing for

- CITE-seq: Cellular Indexing of Transcriptomes and Epitopes by Sequencing

  - scRNAseq with quantitative and qualitative information on surface proteins with available antibodies

- SLAMseq: High throughput metabolic sequencing of transcripts

  - Analyze transcriptome-wide kinetics of RNA synthesis and turnover

  - Measure nascent RNA expression and transcript stability

  - Enhance the temporal resolution of differential expression

- CyCIF: Highly Multiplexed Cell and Tissue Imaging

  - Spatial profiling using multiplexed antibody staining

# Reproducibility and re-use

**Programs**

```
bamtools,2.4.0
bcbio-nextgen,0.9.8a0-8183767
bcbio-variation,0.2.6
bcftools,1.3
bedtools,2.24.0
biobambam,2.0.42
bioconductor-bubbletree,2.1.5
bowtie2,2.2.8
bwa,0.7.13
chanjo,
cnvkit,0.7.11
cufflinks,2.2.1
cutadapt,1.9.1
fastqc,0.11.5
featurecounts,1.4.4
freebayes,1.0.2
gatk,3.2-2-gec30cee
gatk-framework,3.5.21
gemini,0.18.3
grabix,0.1.6
hisat2,2.0.3beta
htseq,0.6.1p1
lumpy-sv,0.2.12
manta,0.29.6
metasv,0.4.0
mutect,1.1.5
novoalign,3.04.04
novosort,V3.00.02
oncofuse,1.1.0
phylowgs,20150714
picard,1.141
platypus-variant,0.8.1
qualimap,2.1.3
rna-star,2.4.1d
rtg-tools,3.6
sailfish,0.9.0
salmon,0.6.0
sambamba,0.6.1
samblaster,0.1.22
samtools,1.3.1
scalpel,0.5.1
```

**Rmarkdown report with code**



22

# Pipeline development

▶ Updating the code base to meet the Python3 standard

▶ Supporting background inputs for copy number variant (CNV) calling to allow for a pre-computed panel of normals for tumor-only or single sample variant calling

▶ Implementing tumor-only variant calling with duplex barcodes

▶ Whole-genome bisulfite-seq pipeline based on bismark2 for DNA methylation analysis with the Illumina Truseq Methyl Capture platform

▶ Making it multiple orders of magnitudes faster to set up runs with thousands of inputs.

▶ For bulk transcriptomics data, added support for gene fusion calling with Arriba

▶ Set up the hg38 reference in STAR

▶ Enabled 2-pass STAR alignment, which performs sequential alignment, genome indexing and re-alignment to improve quantification of novel splice junctions

▶ Notably, we also resolved more than 500 issues on the bcbio github

▶ Now funded through the CZI

# Pipeline development

- ATAC-seq/ChIP-seq pipeline

  - Stabilizing infrastructure

  - Integration of ATACqv

  - Shifting, removing low quality mapping, subsetting in to NFR regions

  - Still problems with experimental design (lack of replicates)

- 5hmC-seq analysis to detect regions with differential methylation

  - Early stages

- Variant calling analysis

  - Exploring cloud options for a large projects

  - Will be learning more about DRAGEN and Terra in the near term

# Training in 2019

▶ 30 workshops focused on basic data analysis skills and NGS analysis

  ▶ spanning 42 training days and training over 900 researchers

▶ For the third year, members of the training team were involved in the development of the online "Introduction to Omics Research" course offered by Harvard Catalyst

▶ Continued to collaborate with Harvard's FAS Research Computing "RC" group to make high performance computing more accessible to Chan School researchers

▶ Launched a monthly Bioinformatics Breakfast community event

▶ Explored new models for training and collaborating, including embedding core bioinformaticians within research labs to train lab members in private group settings and to oversee their analyses, and mentoring students from our training program to perform bioinformatics consulting within the community

▶ We piloted a new program with more workshops every month and have moved away from the intensive course model

25

# Acknowledgements

Harvard Chan School of Public Health

Harvard Medical School

Harvard Stem Cell Institute

Harvard Catalyst

National Institute of Environmental Health Sciences (NIEHS)

Harvard University Center for AIDS Research (CFAR)

AstraZeneca

Boehringer Ingelheim

Chan Zuckerberg Initiative

Perrimon Lab
  Norbert Perrimon
  Sudhir Tattikota
  Ruie-jiun Hung

HMS RC

HMS Single Cell Core

Biopolyers Facility

Molecular Biology Core Facilities

Bauer Core

BWH Single Cell Core