# Single-cell RNA-sequencing (scRNA-seq): knowing the in and outs of the data generated

Ken Lau, Assistant Professor of Cell and Developmental Biology
ken.s.lau@vanderbilt.edu
CQS Summer Academy (8/13/2018)

http://www.mc.vanderbilt.edu/vumcdept/cellbio/laulab/index.html
Twitter: @KenLauLab

**About me and single cells**

- Started lab at Vanderbilt in 2013 with focus on single-cell biology of the gut (IBD and colon cancer)

- Multiplex imaging, CyTOF, scRNA-seq

- Training at Toronto/MIT/Harvard on multivariate analysis, mathematical modeling, and tissue systems

- inDrop in lab since August 2016 (first 1cell customer outside of Boston), > 50 samples ran so far > 500 000 cells sequenced; we have two systems

**Outline**

- Introduction to scRNA-seq techniques

- Discussion on scRNA-seq data issues

- Brief Python introduction

https://github.com/KenLauLab/Discovery_Oriented_Data_Science

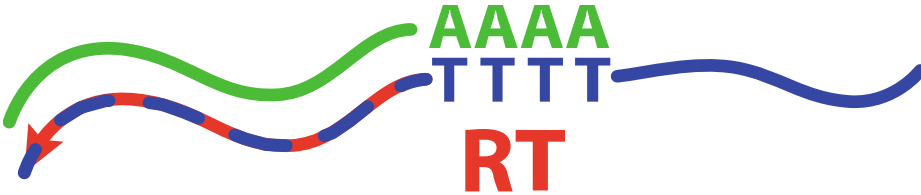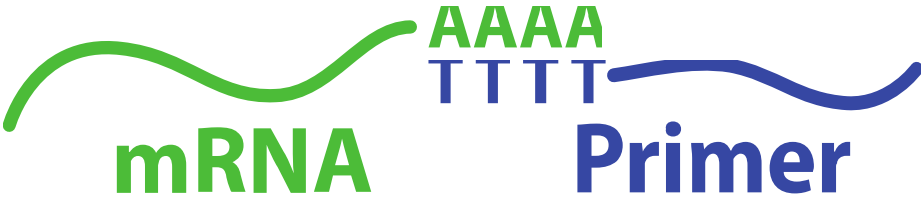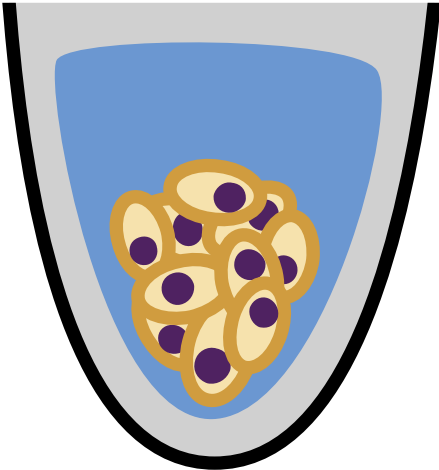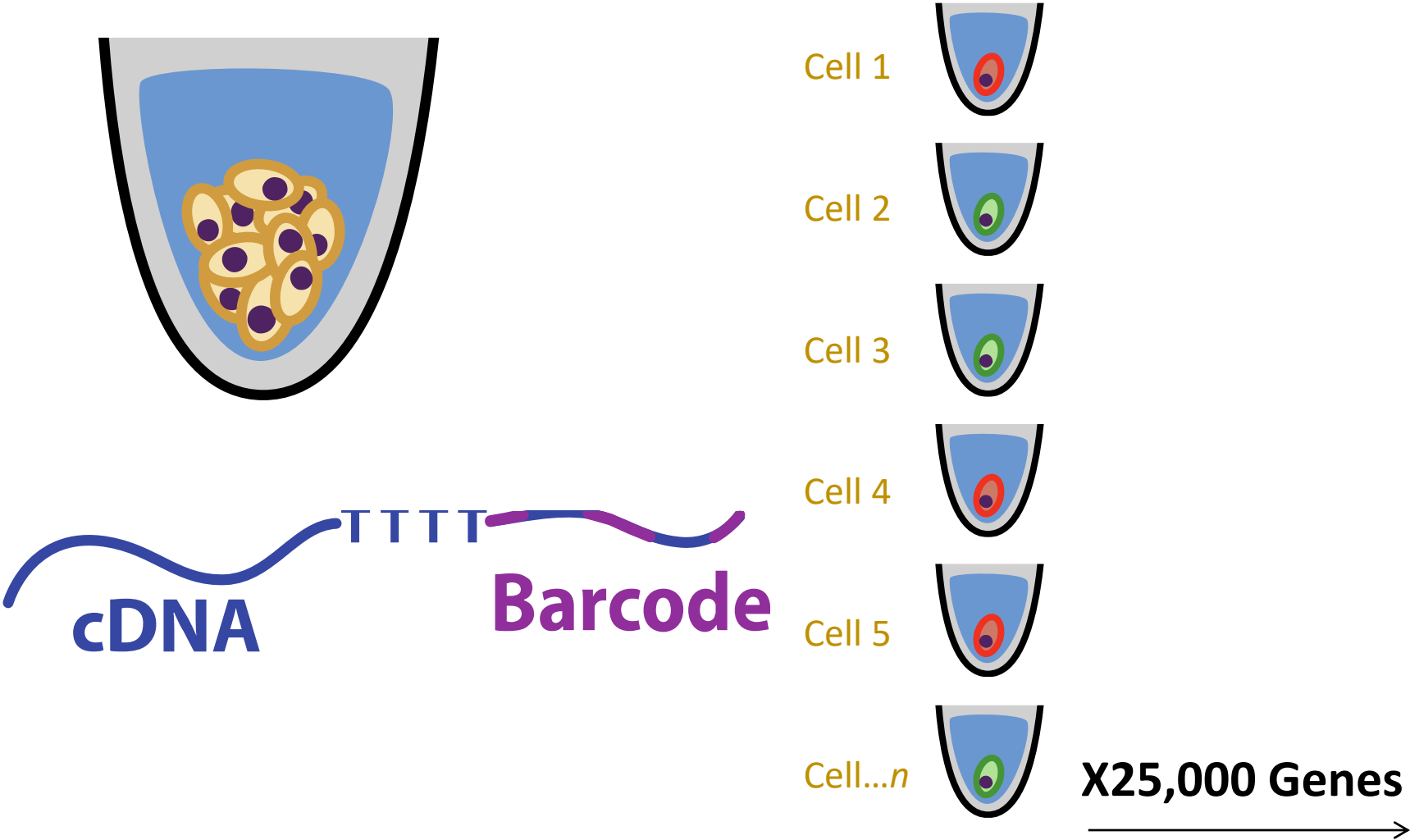# scRNA-seq techniques

# Bulk data versus single-cell data



Bulk

VS.



Single-cell

# Single-cell RNA-seq



cDNA

Barcode

Cell 1

Cell 2

Cell 3

Cell 4

Cell 5

Cell...*n*

**X25,000 Genes**

# scRNAseq protocols

Cell Encapsulation techniques

- Droplet-Based
- Well-Based
- Microfluidic capture (Fluidigm C1)

Lysis and RT

- Coupled – requires balanced mix
- Uncoupled – enables more aggressive lysis
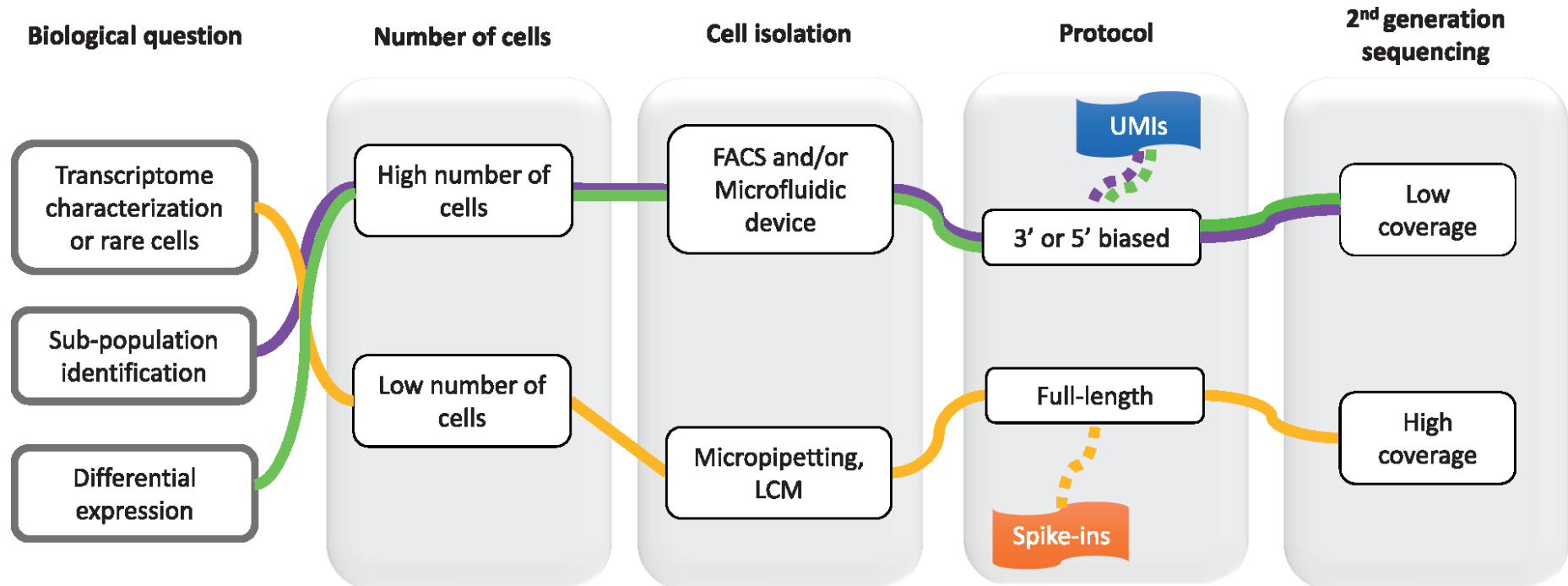
RNA capture strategies

- Poly dT priming
- Targeting / enrichment

Indexing strategies

- During capture/RT    (typically per **cell** indexes - barcodes)
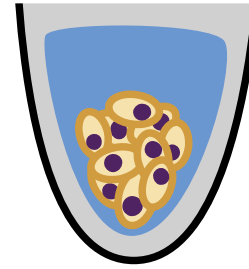- After RT                (typically per **well** indexes)
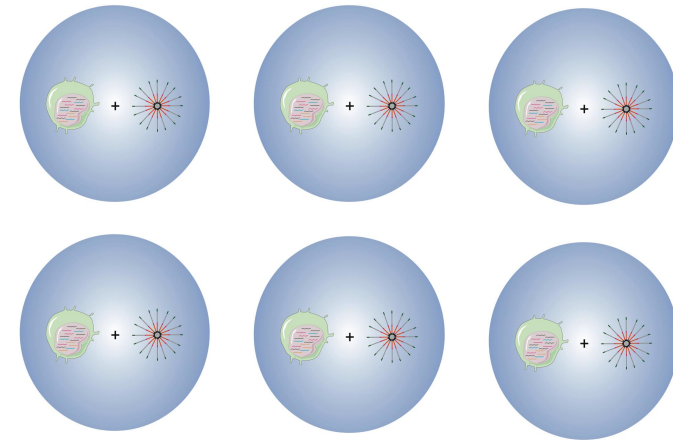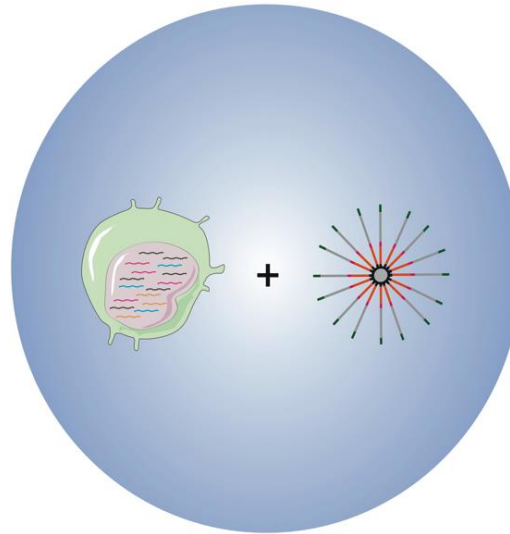
Amplification strategies

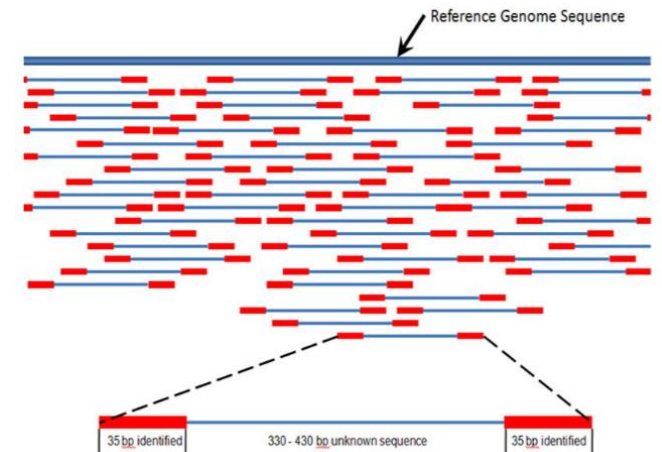- PCR vs IVT

# Typical logic for a scRNA-seq experiment

Alessandra Dal Molin, Barbara Di Camillo; How to design a single-cell RNA-sequencing experiment: pitfalls, challenges and perspectives, *Briefings in Bioinformatics*, , bby007,

**Single cell suspension\*\*\***



**Single-cell encapsulation/
Library preparation**
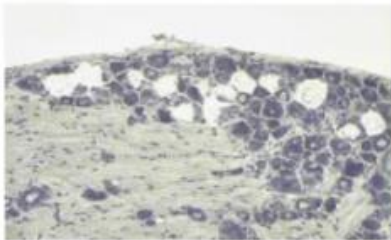


**Sequencing and alignment
(Bioinformatics I)**



Reference Genome Sequence

35 bp identified | 330 - 430 bp unknown sequence | 35 bp identified

# Well-known methods to isolate single cells
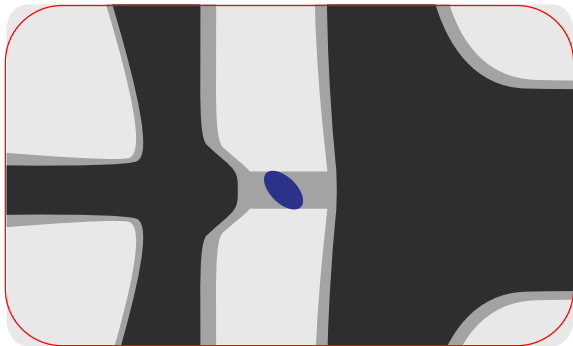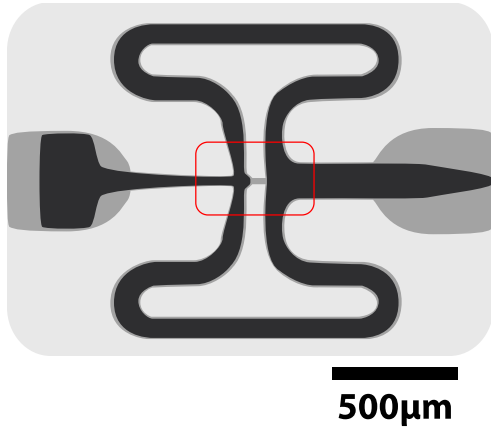
# Plate-Based scRNAseq



20mm

- Isolate RNA, label transcripts using barcoded RT primers (3' seq) or through template switching library prep (enables full length)
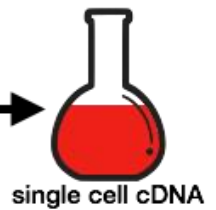
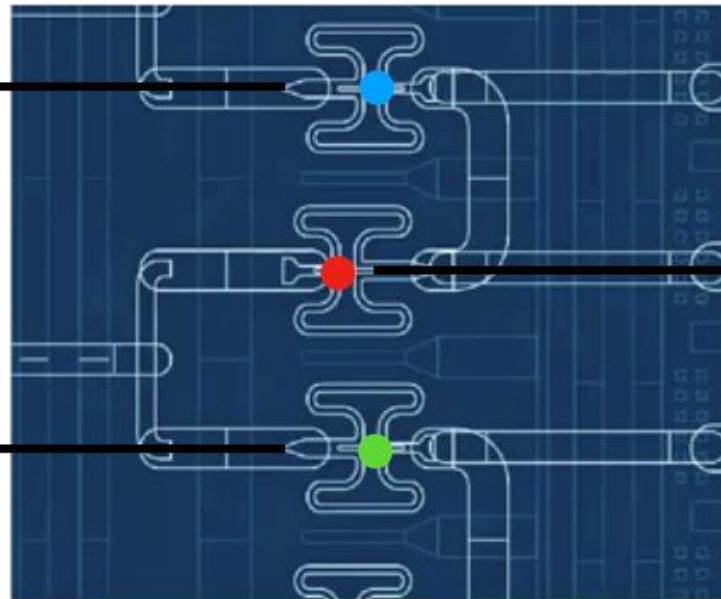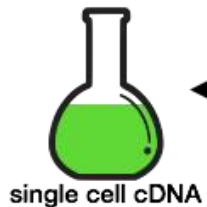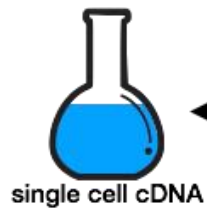Relative to other platforms:
- ~10µl/ cell, <1000 cells (higher volume, lower throughput)

- Deeper sequencing possible - flexibility

- Number of wells become limiting (doublet rate vs. cost)

# Microfluidic capture scRNAseq (Fluidigm C1)



500μm

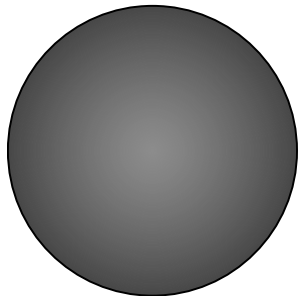Capture chips are built to accommodate specific ranges of cell size

single cell cDNA

single cell cDNA

single cell cDNA

Cells

Stanford course

**Bead-based capture:**
**Immobilized sets of indexed primers**

- Each bead is coated with primers containing a barcode unique to that bead – (index for each cell)
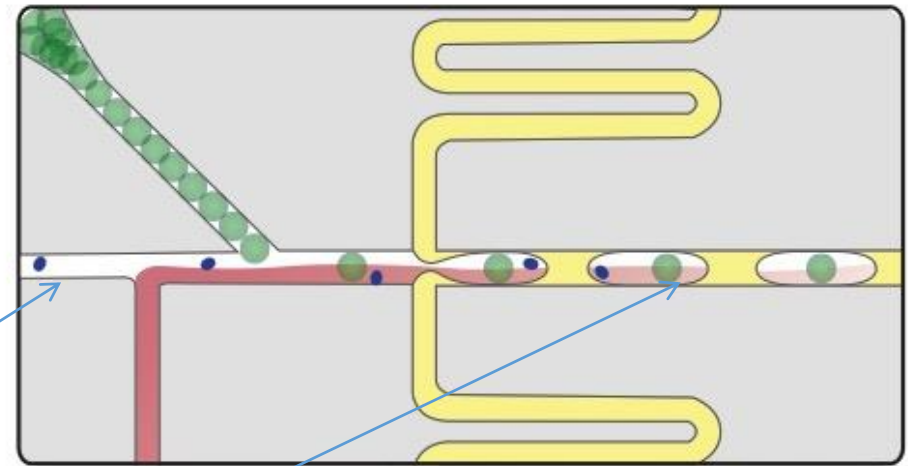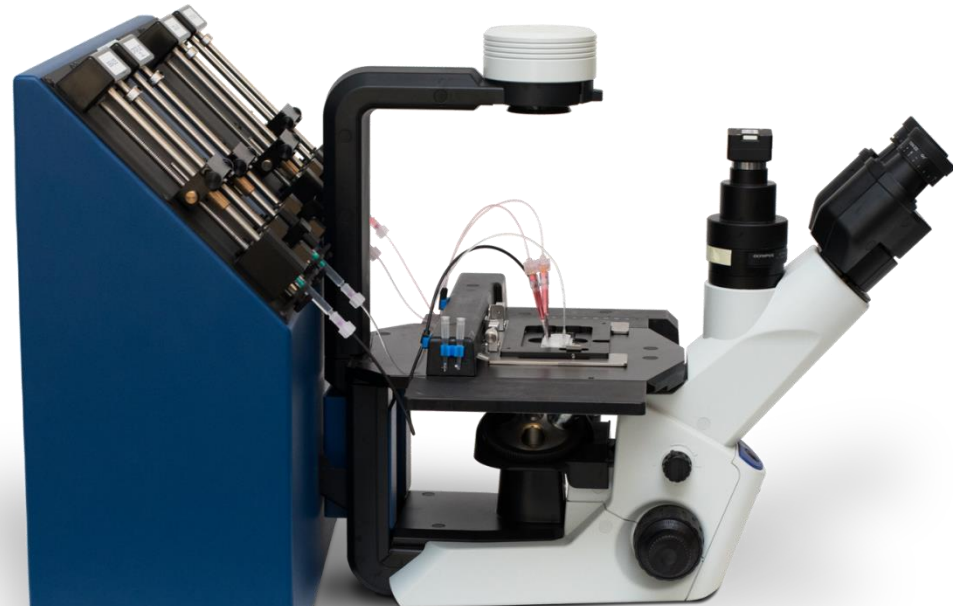
**Barcode**    **UMI**

**TTTT**

- And a Unique Molecular Identifier (UMI) that uniquely tags each primer – (index for each transcript)

# Droplet-based encapsulation

- Co-encapsulating cells and beads in thousands of 1-5nL droplets

- Beads carry barcoded poly-T primers to capture RNA

- Encapsulation rate follows Poisson distribution

- Excess of "vessels" to minimize doublets
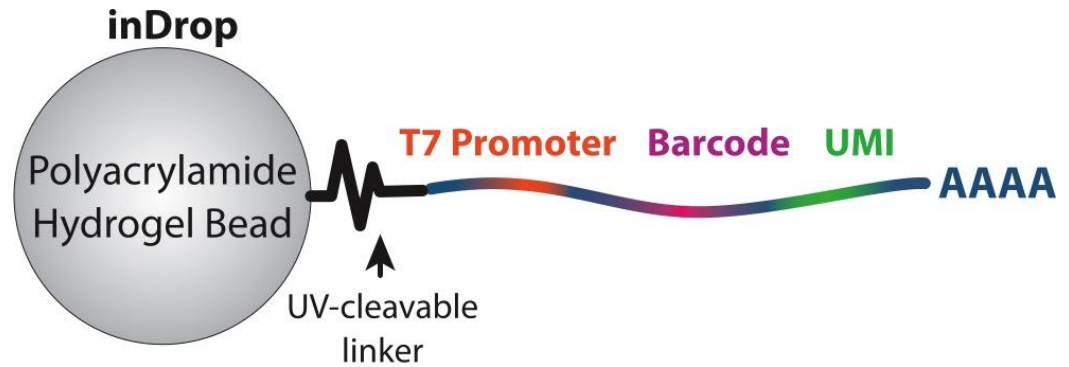
Cells entering/ sec  x          Beads/droplet        =          cells captured/ sec
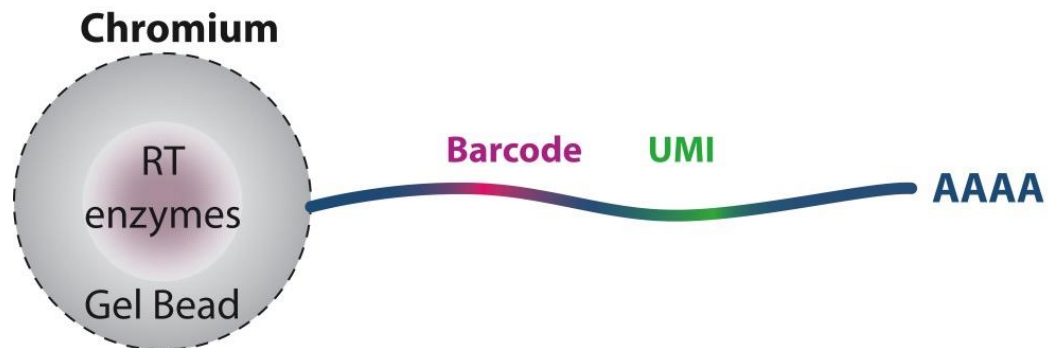
# Beads

- inDrop
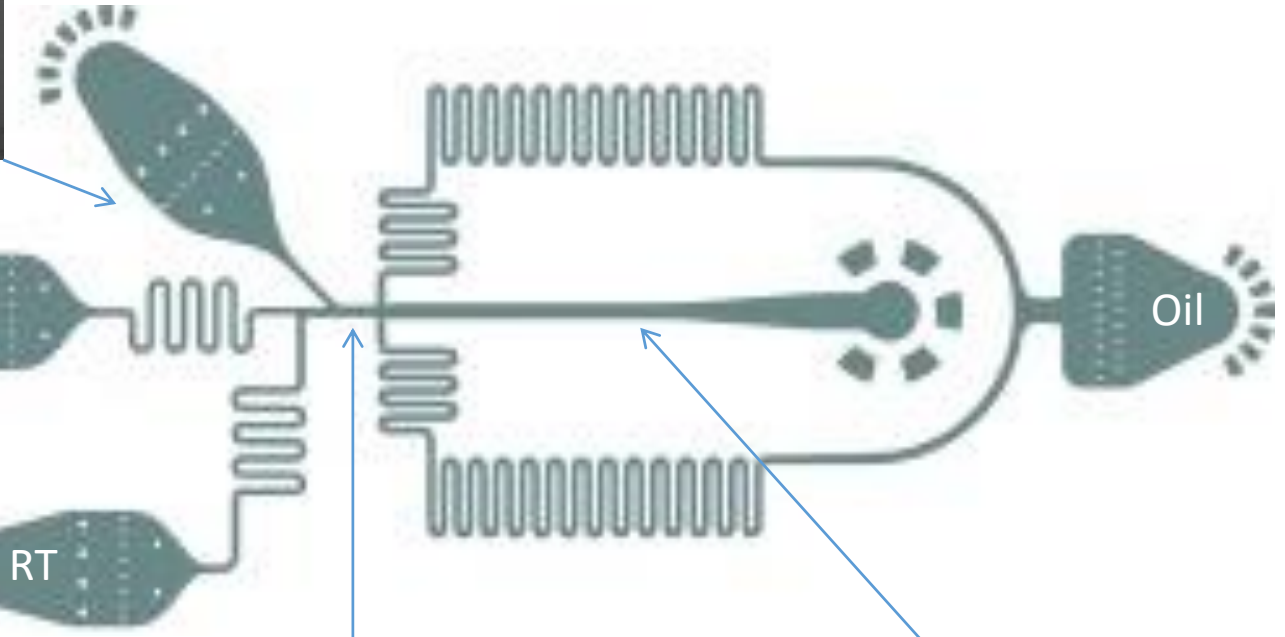  (1cellBio)

- Drop-seq
  (Chemgenes)

- Chromium
  (10x Genomics)



**inDrop**

Polyacrylamide Hydrogel Bead — **T7 Promoter  Barcode  UMI  AAAA**

UV-cleavable linker

**DropSeq**

Polystyrene Bead — **Barcode  UMI  AAAA**

**Chromium**

RT enzymes Gel Bead — **Barcode  UMI  AAAA**

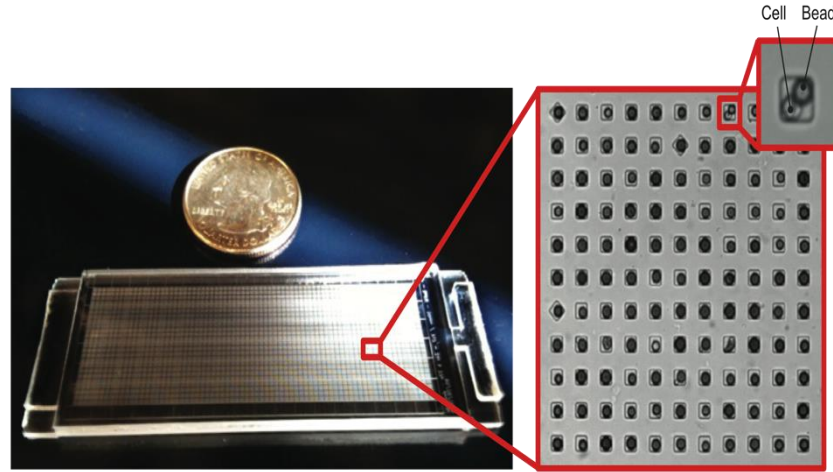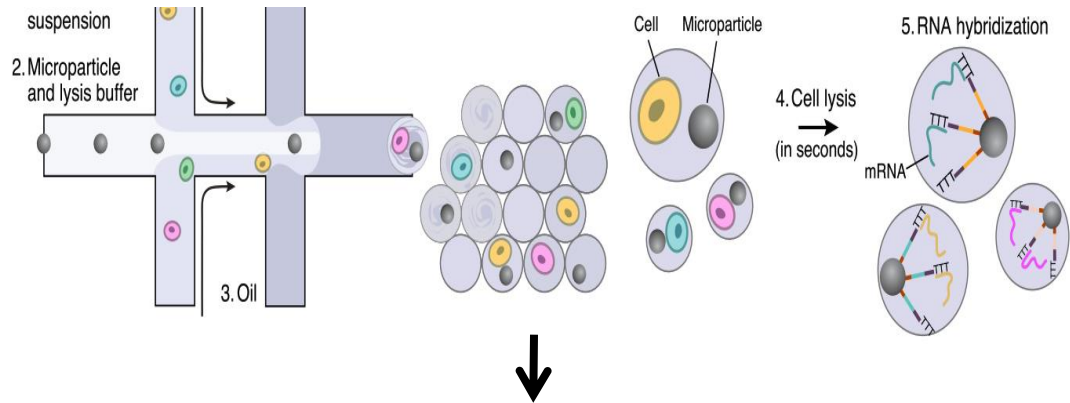# Microfluidic droplet encapsulation chip

Beads

Oil

RT

Cells

Droplet formation

Emulsion output
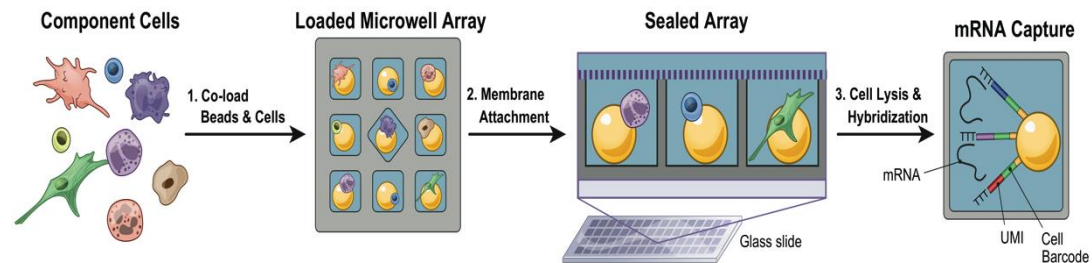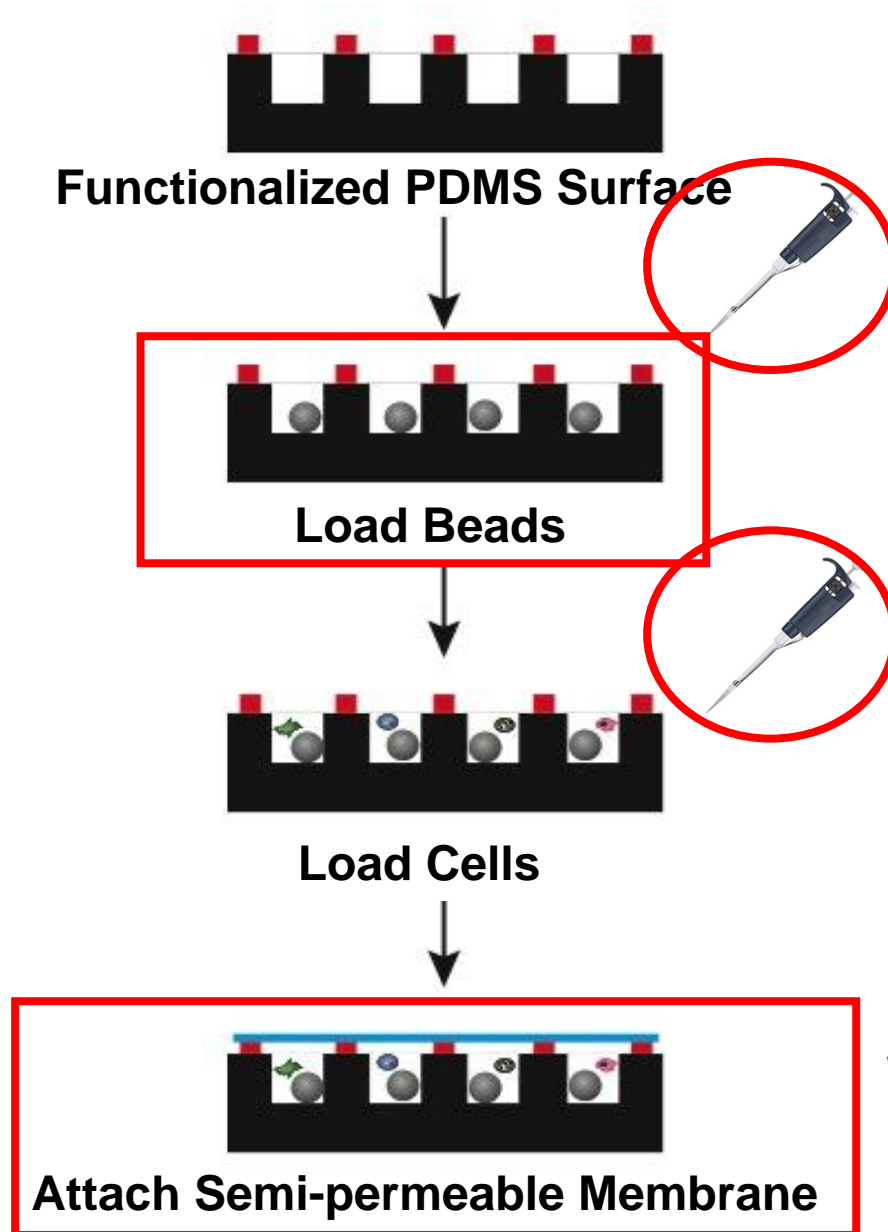
# Seq-Well - microwell sequencing (Shalek lab)

# Seq-Well: Principle

**Functionalized PDMS Surface**

**Load Beads**

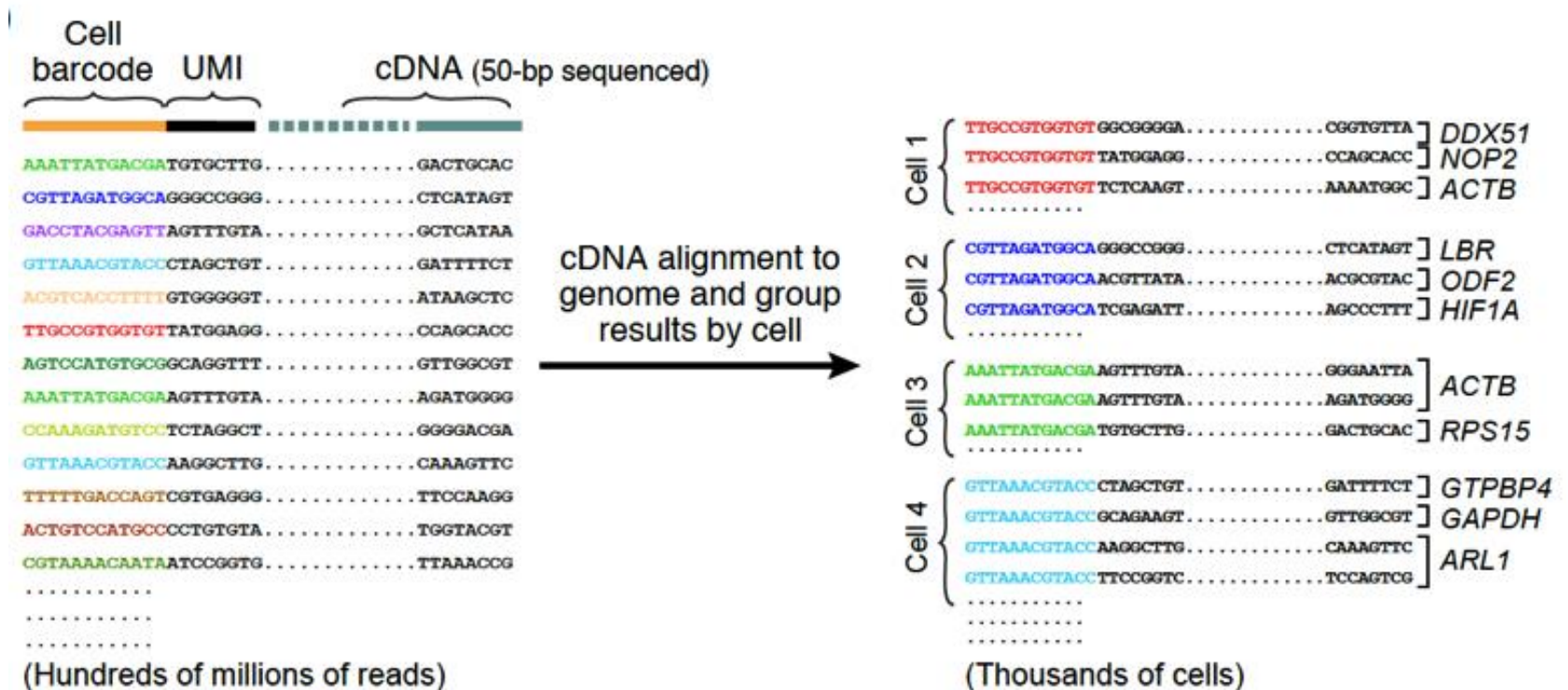**Load Cells**

**Attach Semi-permeable Membrane**

**Nanowell Array**

Bead

Cell

**Size Exclusion ➔ ≤ 1 bead per well**

**Sealing ➔ ~~Cross Contamination~~**

# Deconvolving the data



Reads with same barcodes collapse into cells

Read with same UMIs collapse into transcript counts

# scRNA-seq Data Exploration
## (and problems)

# Table of genes and barcodes

# Inflection point method for identifying barcodes with real cells

# Doublet rate determination



Klein et al, *Cell*, **161**, 2015
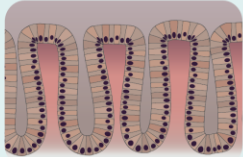Macosko et al, *Cell*, **161**, 2015

# Procedure for the isolation of high viability single-cells from tissues

Keep all reagents on ice and perform all procedures at 4°c, avoid working with overconcentrated tissue/cell solutions
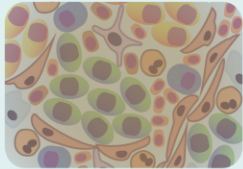
## Tissue

*Keep cold and minimize ischemic time*
*Mince if needed*

Wash in PBS
Decant small debris

Normal colonic epithelium can be isolated through chelation, while other tissues may require direct mechanical processing to achieve 50-500µm fragments. Care should be taken to remove dead cells during washes to maintain viability

Wash in PBS
Decant small debris

The process of isolating tissue fragments should be optimized to the needs of the target cells. This may mean filtering or taking other measures to enrich.

## Chelation

10ml DPBS(-Mg/Ca)
3mM EDTA
0.5mM DTT
10mM NAC

Inversion 20 rpm

~30-60min
Change buffer every 10-15 min

Use forceps to move tissue

Move back to fresh chelation buffer
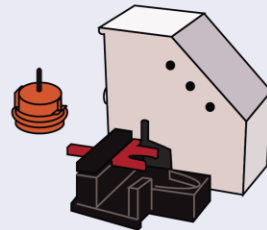
## Medimachine

In cold room
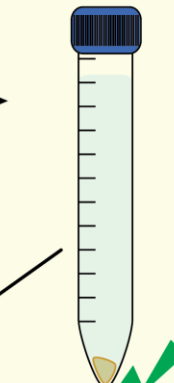Pre-wet 50µm medicon before use
Wash quickly and thoroughly after

*GentleMACS is a potentially less damaging approach we have not yet tested for breaking apart tissue and further dissociating cells.*

## Shaking

10ml DPBS (-Mg/Ca) in a new tube

2-3 Shakes/sec
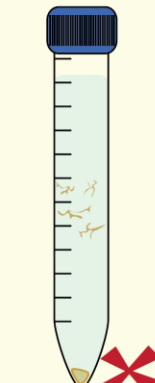~6 inches motion

Wash & inspect crypts

## Washing

DPBS (-Mg/Ca) spin @~300xg for 2.5min

Good pelleting: (fat or pancreas may still float in first wash)

Sticking to tube; may need more chelation (EDTA)*

Floating/clumping typically caused by excess mucin or inviable cells**
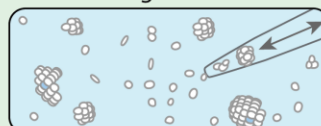
## Cold Protease Dissociation

2-5ml DPBS (-Mg/Ca) + 2.5mg/ml DNase + 5mg/ml Subtilisin
*Use 2ml for every 50-100µl of peleted tissue, dont overconcentrate*
*Add buffer to frozen subtilisin aliquot to thaw immediately before use*
*Pipette with a 1000µl tip every 5-10 min and check for singlets*

4-6°c with gentle motion

*Pellet without EDTA prior to re-suspending in protease*
**Optimize spin conditions so that single cells are decanted*
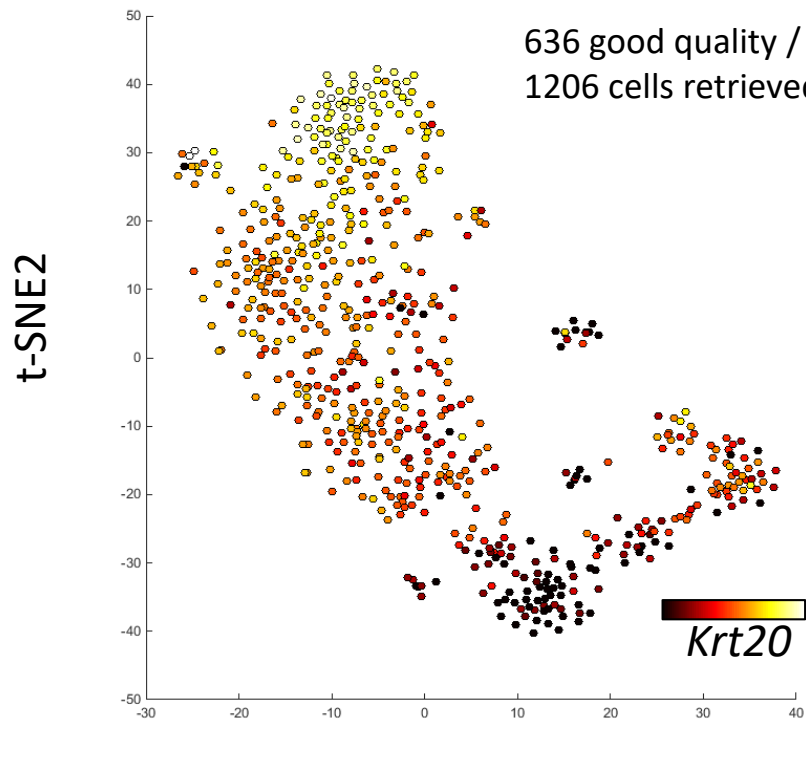**Dont leave tissue in pellet too long, split if too dense*
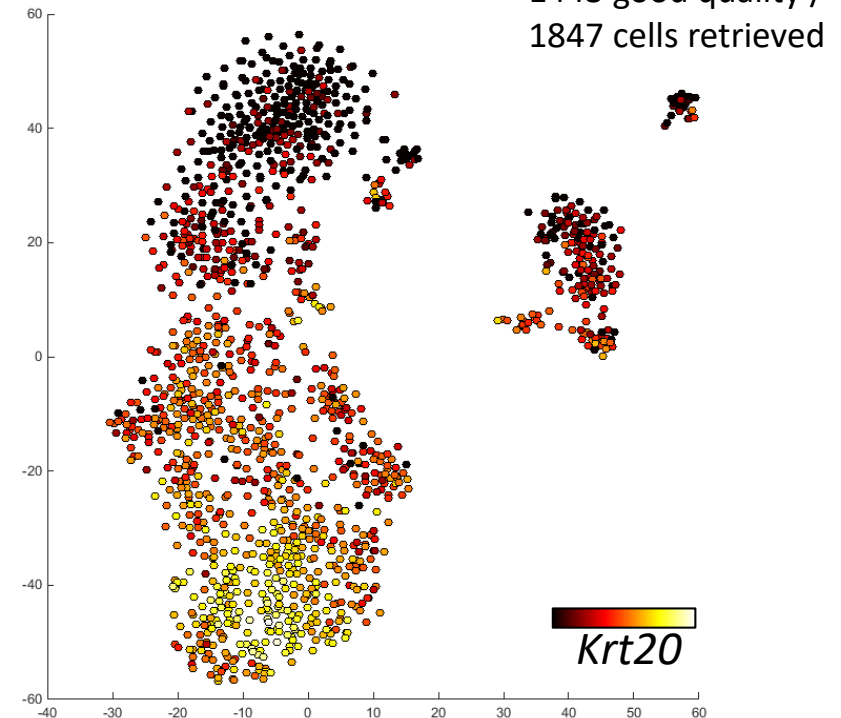**Re-suspend well, but do not over pipette prior to dissociation*

# Single-cell isolation for scRNA-seq to minimize dissociation artifacts

- cold protease from *Bacillus licheniformis,* soil bacteria from Himalayan glaciers

- enables tissue preparation on ice (at 4 degrees)
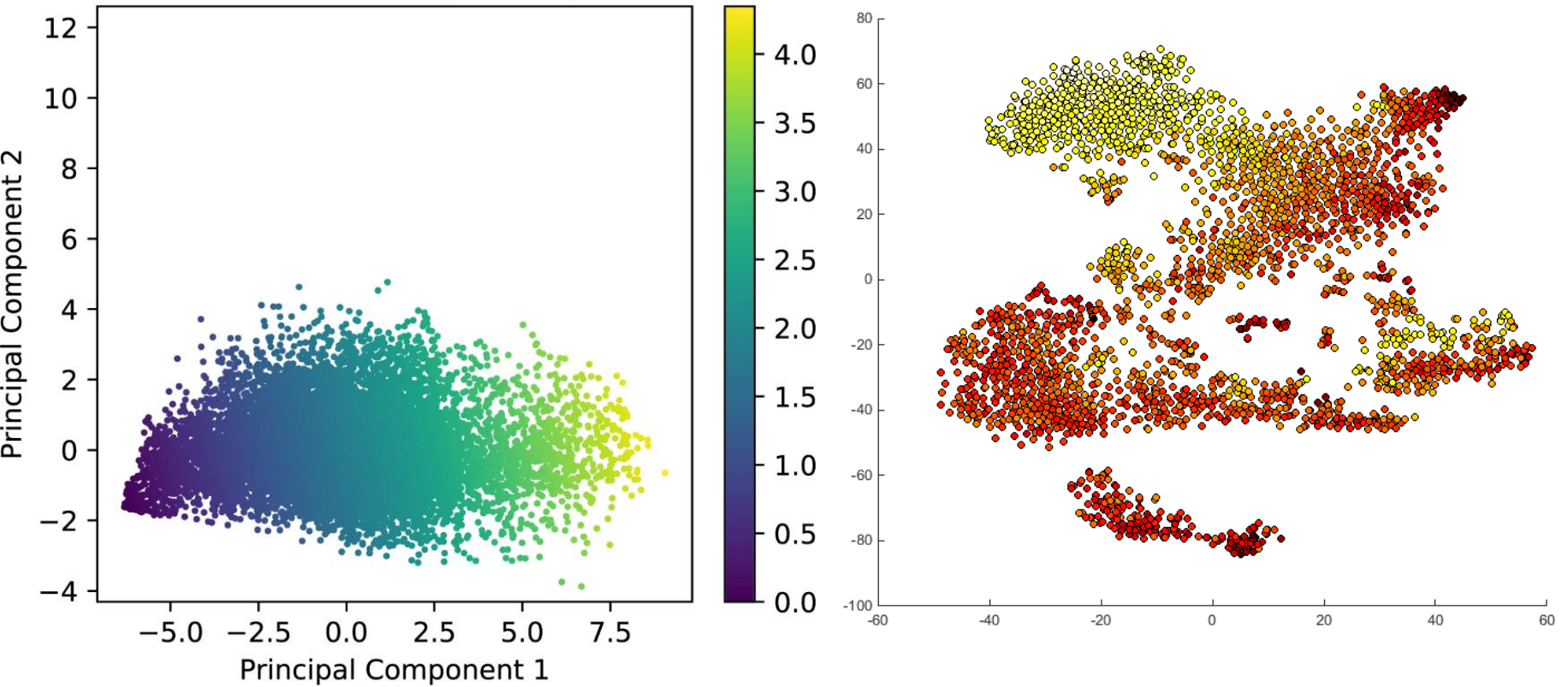
37 degree DNAse/Collagenase

636 good quality /
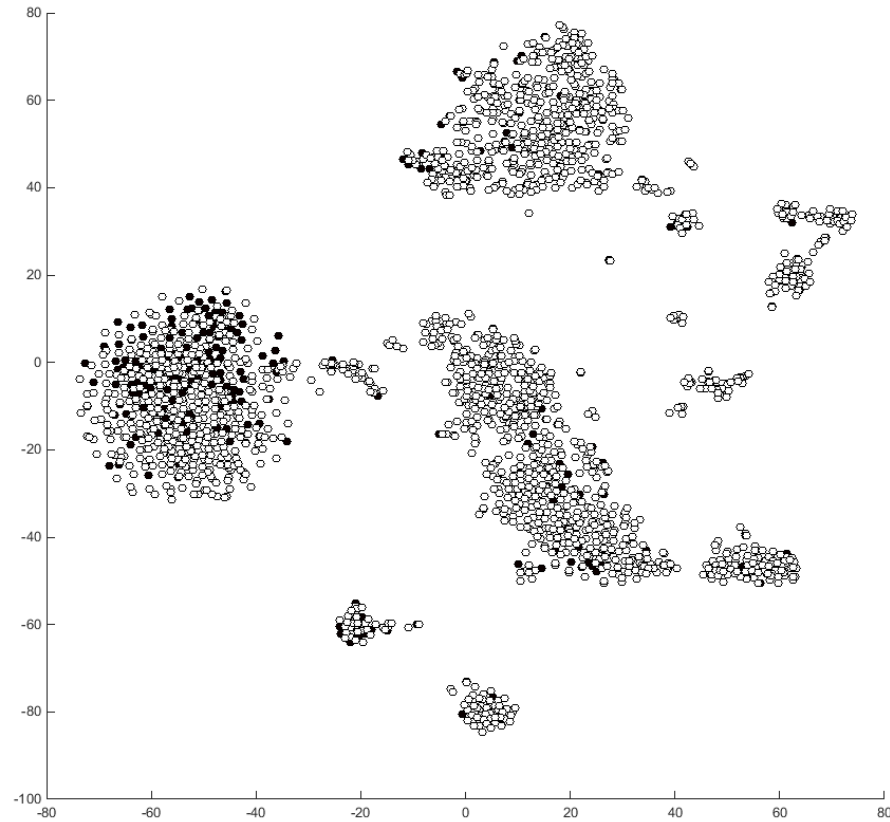1206 cells retrieved

Cold Protease (on ice)

1443 good quality /
1847 cells retrieved

t-SNE2

t-SNE1

*Krt20*

*Krt20*

- caveat – the efficacy for retrieving all cell types from all tissues unknown

(Adam et al., *Development*, 2018)
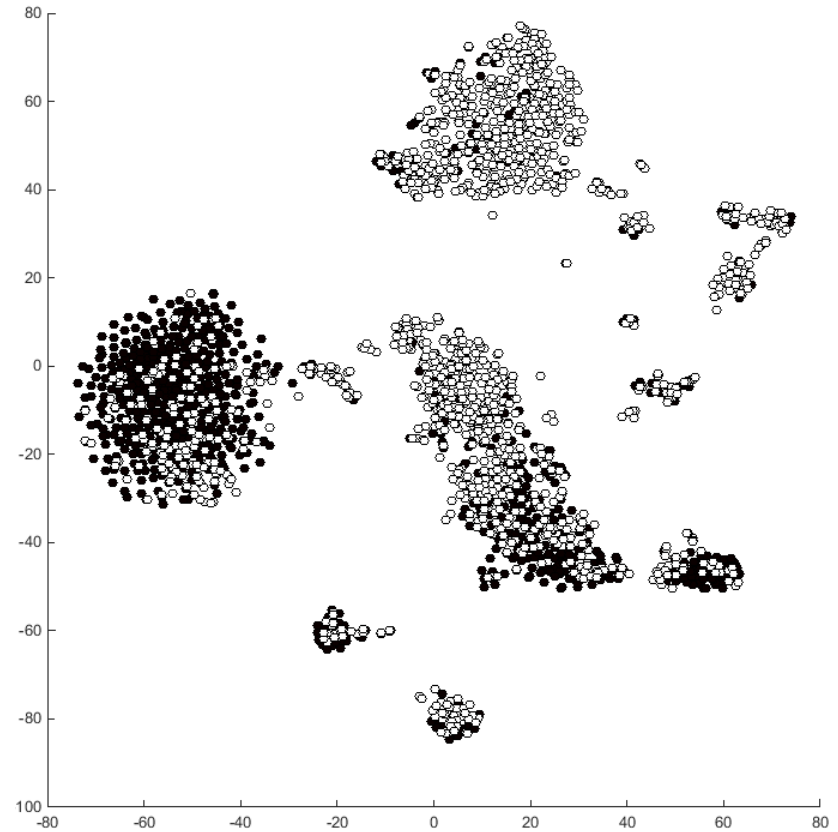
# Stressed/dying cells



PCA of mitochondrial gene expression

# Contaminant from free floating RNA and leaky cell corpses
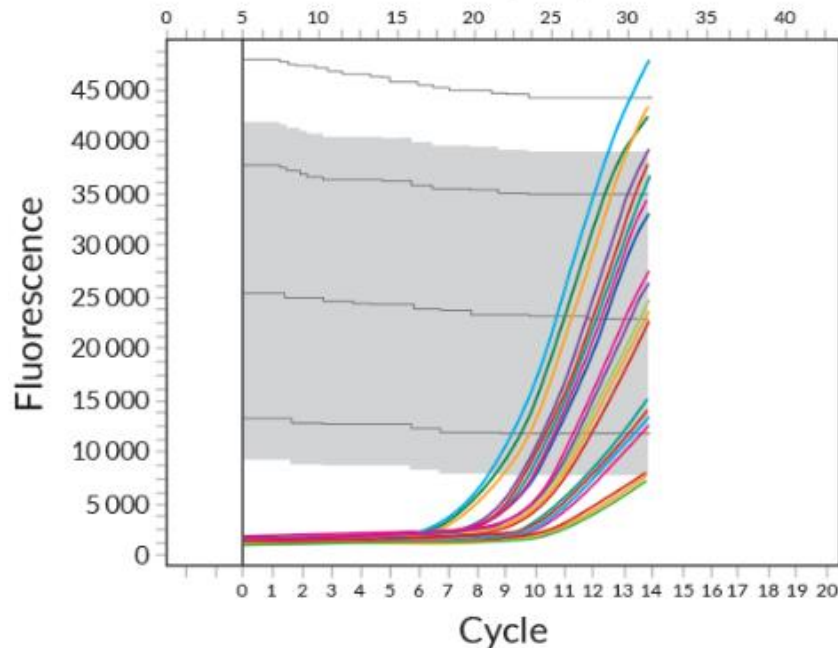

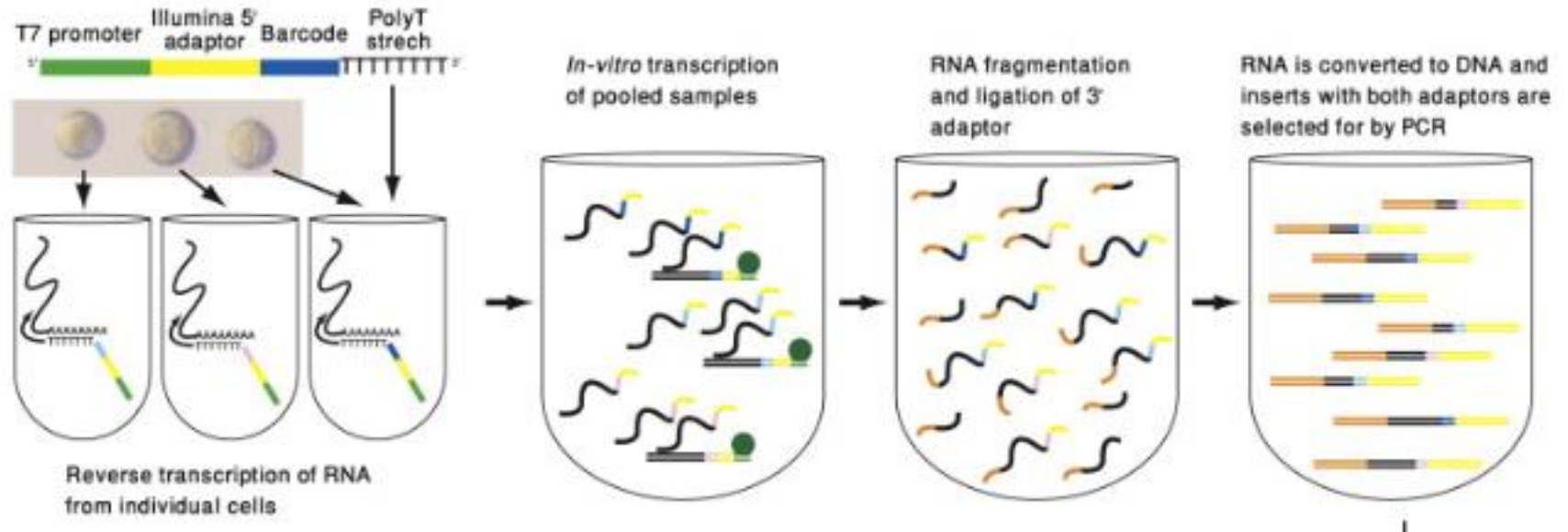
Generous threshold

Strict threshold
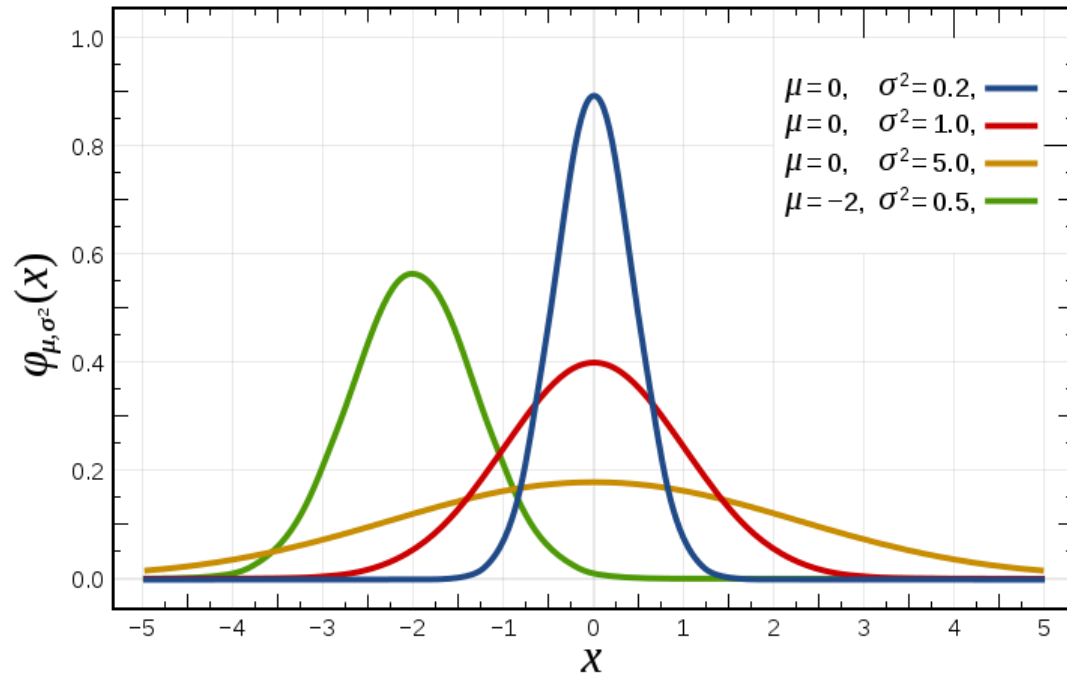
# Amplification bias in scRNA-seq



Non-linear amplification

- Highly expressed transcripts are inappropriately represented and replicated

- Sampling of RNA in a cell 1-10%

- Zero inflated data

- Count data – negative binomial distribution

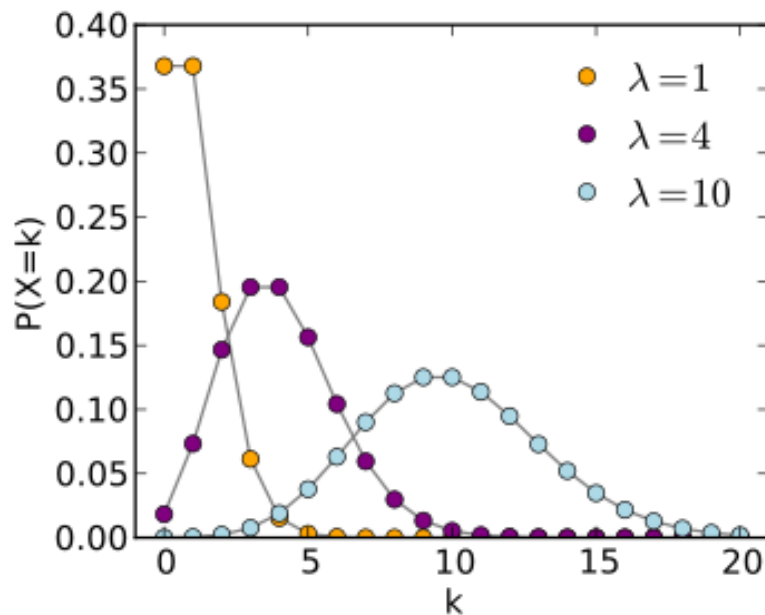# In vitro transcription results in linear amplification of RNA



(Itai lab, CEL-seq)

# Negative binomial versus Gaussian distributions



Gaussian
- Continuous variable
- Symmetric



Negative Binomial
- Discrete variable (counts)
- Asymmetric at small means

# Zero inflation in scRNA-seq



- Only 5% of this table is not 0!

# What does this mean?

- Statistical tests to determine differentially expressed genes with a zero-inflated negative binomial (ZINB)

- False negatives (genes that are supposed to be expressed by appear as 0)

- Many genes (columns) have low counts due to shallow sampling of transcripts – this means the data are noisy

- Unreliable variables that need to be processed/filtered out prior to downstream analysis
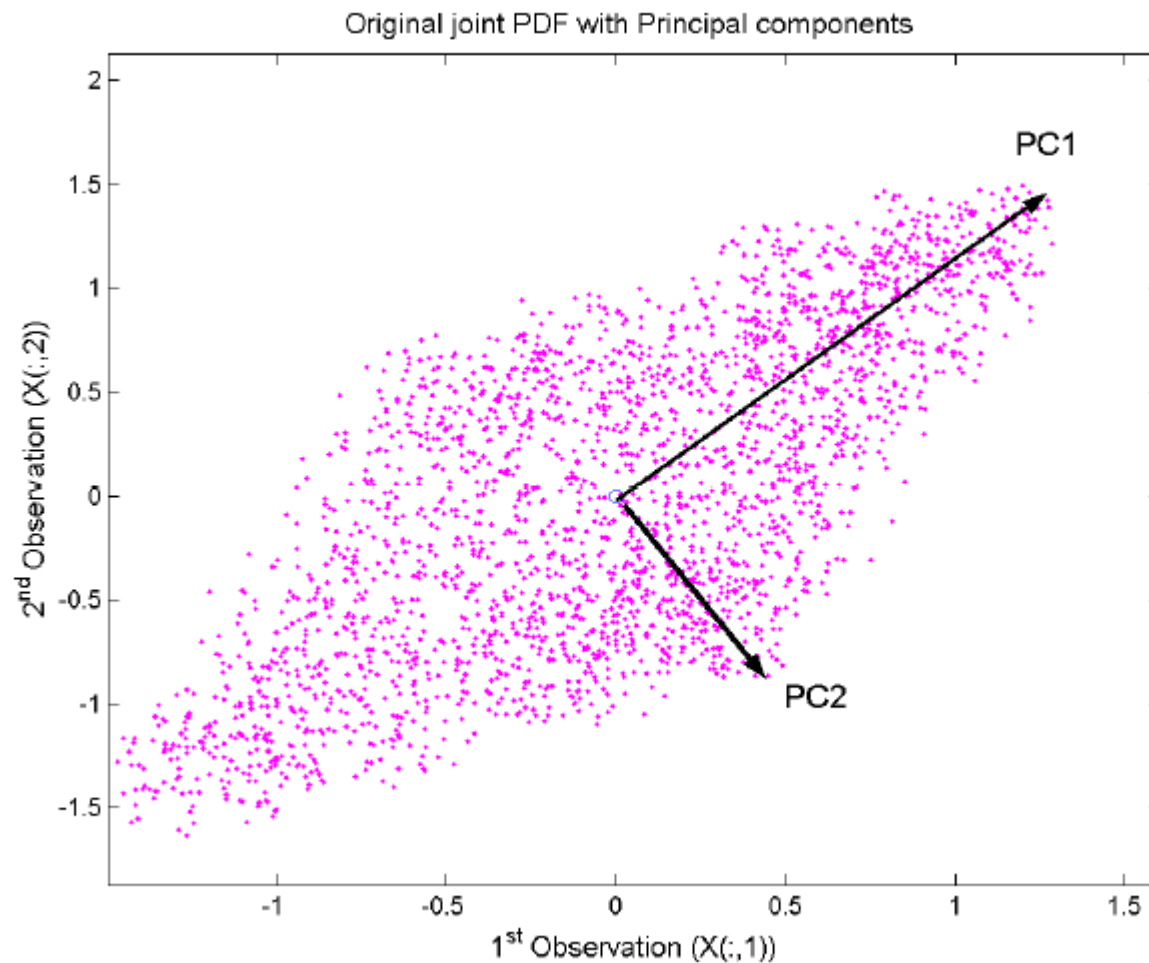
# Feature Selection

# Variance selection

- Easiest – select genes that are the most variable

- Variance = genes that are most different across all cells

$$\sigma^2 = \sum \frac{(X - \mu)^2}{N}$$

- Rank genes by top 500 most variable, for example, discard the rest
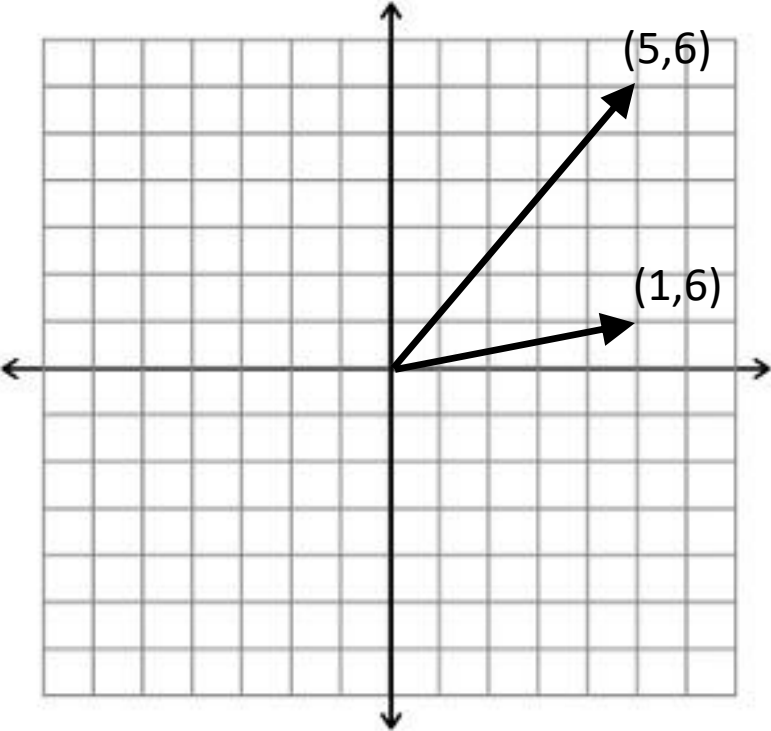
# Highly variable axes = Principal Components



Original joint PDF with Principal components

## PCA (Principal Component Analysis)

- Principle of PCA is to maximize the Variance of X with the least amount of principal components (latent variables)

- What is variance? Spread of the data, information content, change etc.

- Variance is the covariance of a dataset with itself, i.e. Var(X) = Cov(X,X) → Maximize

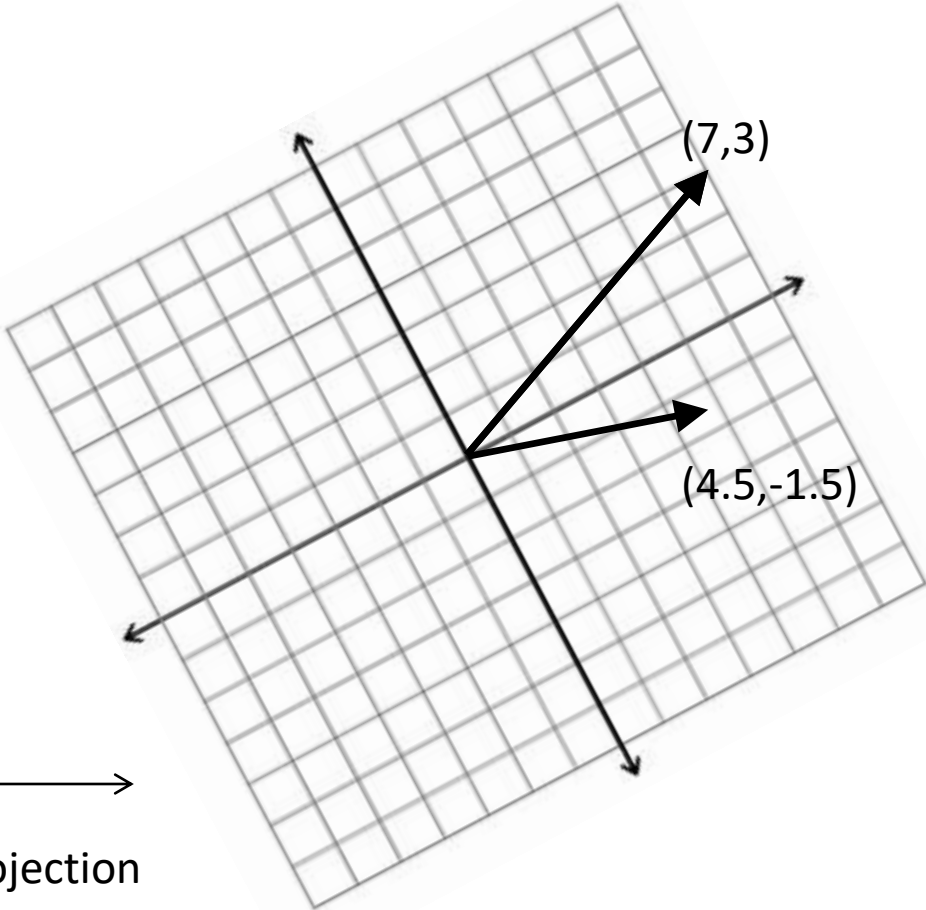- What are principal components? Linear combinations of original variables – linear transformation

# Vectors and projections



Basis Set

(1,0)
(0,1)

Projection

Basis Set

(7,4)
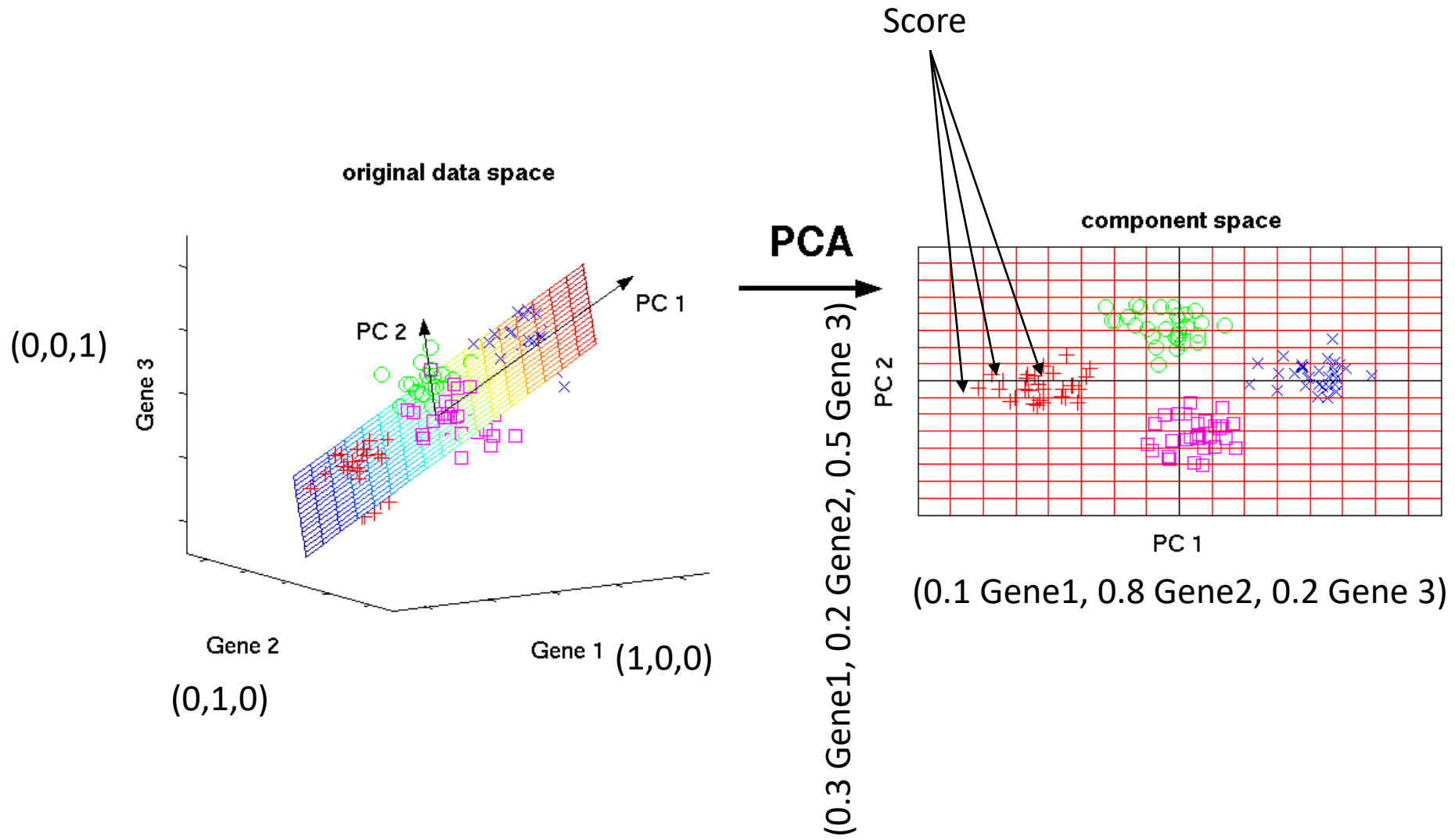(-3.5,7)

# PCA as a dimension reduction tool

# Selection of principal components



$$[\text{stem}(\mathit{diag}(S).{^\wedge}2)]$$

Noise starting

**PRINCIPAL COMPONENTS CAN BE ORDERED BY EIGENVALUES (VARIANCE CAPTURED)**

# Simple summary of this simple feature selection procedure

- Keep most variable genes (over all cells) for downstream analysis, discard rest as noise

- PCA identifies super axes (Principal Components) that are combinations of the original variables that captures the most variance in the data (by eigenvalues)

- Orthogonal – no duplicate or redundant axes – so will only have a few of them