



# **Applying Single Cell genomics to your research: Discussion of experimental and computational frameworks**

**Alexandra-Chloé Villani, PhD**

**Society of Immunotherapy of Cancer Meeting  
Workshop on Single Cell Techniques in Immunology and Cancer  
Immunotherapy**

**November 9<sup>th</sup> 2017**

- **No disclosure**



# Outline

1. Introduction
  - A. Relevance & advances in single cell sequencing
  - B. Overview single cell assay
2. Cell isolation for single cell readout
3. scRNAseq protocols
4. Other single cell readouts & multi-omics
5. Analysis overview
6. Technical challenges
7. Experimental design & common questions
8. Application & validation

# ARTICLE

doi:10.1038/nature13437

## Single-cell RNA-seq reveals dynamic paracrine control of cellular variation

Alex K. Shalek<sup>1,2,3\*</sup>, Rahul Satija<sup>3\*</sup>, Joe Shuga<sup>4\*</sup>, John J. Trombetta<sup>3</sup>, Dave Gennert<sup>3</sup>, Diana Lu<sup>3</sup>, Peilin Chen<sup>4</sup>, Rona S. Gertner<sup>1,2</sup>, Jellert T. Gaublotte<sup>1,2</sup>, Nir Yosef<sup>2</sup>, Schraga Schwartz<sup>2</sup>, Brian Fowler<sup>4</sup>, Suzanne Weaver<sup>4</sup>, Jing Wang<sup>4</sup>, Xiaohui Wang<sup>4</sup>, Ruihua Ding<sup>1,2</sup>, Raktima Raychowdhury<sup>3</sup>, Nir Friedman<sup>5</sup>, Nir Hacohen<sup>3,6</sup>, Hongkun Park<sup>1,2,3</sup>, Andrew P. May<sup>4</sup> & Aviv Regev<sup>3,7</sup>

## Distinct myeloid progenitor–differentiation pathways identified through single-cell RNA sequencing

Roy Drissen<sup>1,2</sup>, Natalija Buza-Vidas<sup>2</sup>, Petter Woll<sup>1,3</sup>, Supat Thongjuea<sup>1</sup>, Adriana Gambardella<sup>1,2</sup>, Alice Giustacchini<sup>1,3</sup>, Elena Mancini<sup>4</sup>, Alya Zriwil<sup>5</sup>, Michael Lutteropp<sup>1,3</sup>, Amit Grover<sup>1,2,4</sup>, Adam Mead<sup>1,3</sup>, Ewa Sitnicka<sup>5</sup>, Sten Eirik W Jacobsen<sup>1,3,6</sup> & Claus Nerlov<sup>1,2,4,6</sup>

Cell

## An Immune Atlas of Clear Cell Renal Cell Carcinoma

Stéphane Chevrier<sup>1,15</sup>, Jacob Harrison Levine<sup>2,15</sup>, Vito Riccardo Tomaso Zanotelli<sup>1,3</sup>, Karina Silina<sup>4</sup>, Daniel Schulz<sup>1</sup>, Marina Bacac<sup>5</sup>, Carola Hermine Ries<sup>6</sup>, Laurie Ailles<sup>7,8</sup>, Michael Alexander Spencer Jewett<sup>8</sup>, Holger Moch<sup>9</sup>, Maries van den Broek<sup>4</sup>, Christian Beisel<sup>10</sup>, Michael Beda Stadler<sup>11,12</sup>, Craig Gedye<sup>13</sup>, Bernhard Reis<sup>14</sup>, Dana Pe'er<sup>2</sup> and Bernd Bodenmiller<sup>1,16,\*</sup>

Cell

## Innate Immune Landscape in Early Lung Adenocarcinoma by Paired Single-Cell Analyses

Yonit Lavin<sup>1,2,3</sup>, Soma Kobayashi<sup>1,2,3,14</sup>, Andrew Leader<sup>1,2,3,14</sup>, El-ad David Amir<sup>2,3,9</sup>, Naama Elefant<sup>10</sup>, Camille Bigenwald<sup>1,2,3</sup>, Romain Remark<sup>1,2,3,13</sup>, Robert Sweeney<sup>6,7</sup>, Christian D. Becker<sup>4</sup>, Jacob H. Levine<sup>11</sup>, Klaus Meinhof<sup>4</sup>, Andrew Chow<sup>1,2,3</sup>, Seunghee Kim-Shulze<sup>2,3,9</sup>, Andrea Wolf<sup>6</sup>, Chiara Medaglia<sup>10</sup>, Hanjie Li<sup>10</sup>, Julie A. Rytlewski<sup>12</sup>, Ryan O. Emerson<sup>12</sup>, Alexander Solovoyov<sup>1,3,5,8</sup>, Benjamin D. Greenbaum<sup>1,3,5,8</sup>, Catherine Sanders<sup>12</sup>, Marissa Vignali<sup>12</sup>, Mary Beth Beasley<sup>8</sup>, Raja Flores<sup>6</sup>, Sacha Gnjatovic<sup>2,3,5,9</sup>, Dana Pe'er<sup>11</sup>, Adeeb Rahman<sup>2,3,7,9</sup>, Ido Amit<sup>10</sup> and Miriam Merad<sup>1,2,3,9,15,\*</sup>

# LETTER

doi:10.1038/nature20105

## Single-cell RNA-seq identifies a PD-1<sup>hi</sup> ILC progenitor and defines its development pathway

Yong Yu<sup>1\*</sup>, Jason C. H. Tsang<sup>1,2,3\*</sup>, Cui Wang<sup>1,4\*</sup>, Simon Clare<sup>1</sup>, Juexuan Wang<sup>1</sup>, Xi Chen<sup>1</sup>, Cordelia Brandt<sup>1</sup>, Leanne Kane<sup>1</sup>, Lia S. Campos<sup>1</sup>, Liming Lu<sup>5</sup>, Gabrielle T. Belz<sup>6,7</sup>, Andrew N. J. McKenzie<sup>8</sup>, Sarah A. Teichmann<sup>1,9</sup>, Gordon Dougan<sup>1,10</sup> & Pentao Liu<sup>1</sup>

# Science

## Massively Parallel Single-Cell RNA-Seq for Marker-Free Decomposition of Tissues into Cell Types

Diego Adhemar Jaitin<sup>1\*</sup>, Ephraim Kenigsberg<sup>2,3\*</sup>, Hadas Keren-Shaul<sup>1\*</sup>, Naama Elefant<sup>1</sup>, Franziska Paul<sup>1</sup>, Irina Zaretsky<sup>1</sup>, Alexander Mildner<sup>1</sup>, Nadav Cohen<sup>2,3</sup>, Steffen Jung<sup>1</sup>, Amos Tanay<sup>2,3</sup>††, Ido Amit<sup>1</sup>††

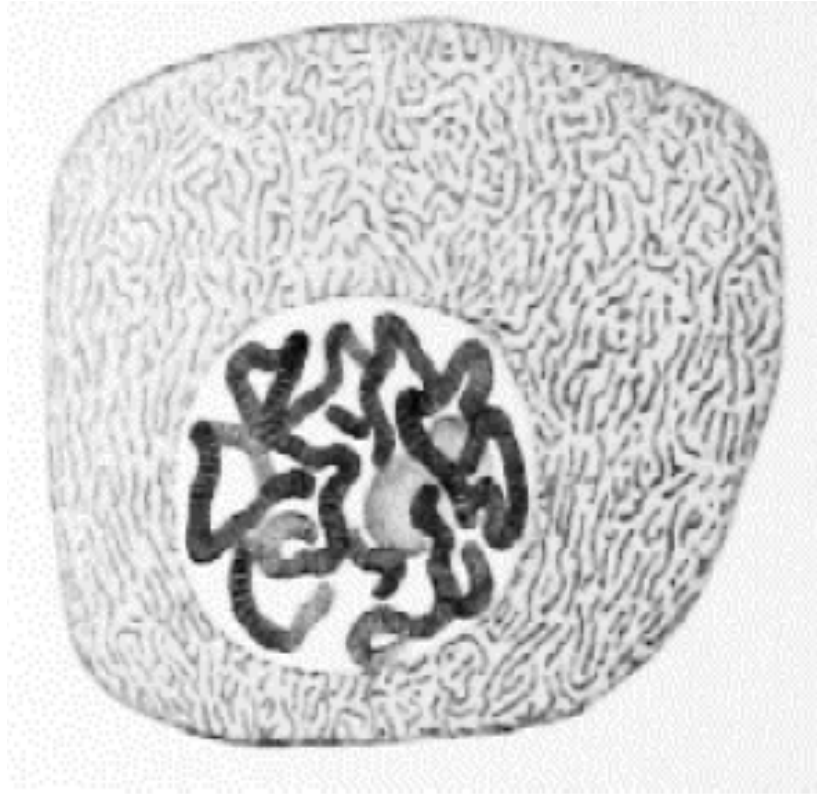
### RESEARCH ARTICLES

#### CANCER GENOMICS

## Dissecting the multicellular ecosystem of metastatic melanoma by single-cell RNA-seq

Itay Tirosh<sup>1\*</sup>, Benjamin Izar<sup>1,2,3\*</sup>††, Sanjay M. Prakadan<sup>1,4,5,6</sup>, Marc H. Wadsworth II<sup>1,4,5,6</sup>, Daniel Treacy<sup>1</sup>, John J. Trombetta<sup>1</sup>, Asaf Rotem<sup>1,2,3</sup>, Christopher Rodman<sup>1</sup>, Christine Lian<sup>7</sup>, George Murphy<sup>7</sup>, Mohammad Fallahi-Sichani<sup>8</sup>, Ken Dutton-Regester<sup>1,2,9</sup>, Jia-Ren Lin<sup>10</sup>, Ofir Cohen<sup>1</sup>, Parin Shah<sup>2</sup>, Diana Lu<sup>1</sup>, Alex S. Genshaft<sup>1,4,5,6</sup>, Travis K. Hughes<sup>1,4,6,11</sup>, Carly G. K. Ziegler<sup>1,4,6,11</sup>, Samuel W. Kazer<sup>1,4,5,6</sup>, Aleth Gaillard<sup>1,4,5,6</sup>, Kellie E. Kolb<sup>1,4,5,6</sup>, Alexandra-Chloé Villani<sup>1</sup>, Cory M. Johannessen<sup>1</sup>, Aleksandr Y. Andreev<sup>1</sup>, Eliezer M. Van Allen<sup>1,2,3</sup>, Monica Bertagnolli<sup>12,13</sup>, Peter K. Sorger<sup>8,10,14</sup>, Ryan J. Sullivan<sup>15</sup>, Keith T. Flaherty<sup>15</sup>, Dennie T. Frederick<sup>15</sup>, Judit Jané-Valbuena<sup>1</sup>, Charles H. Yoon<sup>12,13</sup>†, Orít Rozenblatt-Rosen<sup>1</sup>†, Alex K. Shalek<sup>1,4,5,6,11,16</sup>†, Aviv Regev<sup>1,17,18</sup>††, Levi A. Garraway<sup>1,2,3,14</sup>††

# Analysis at single cell level is an old concept!



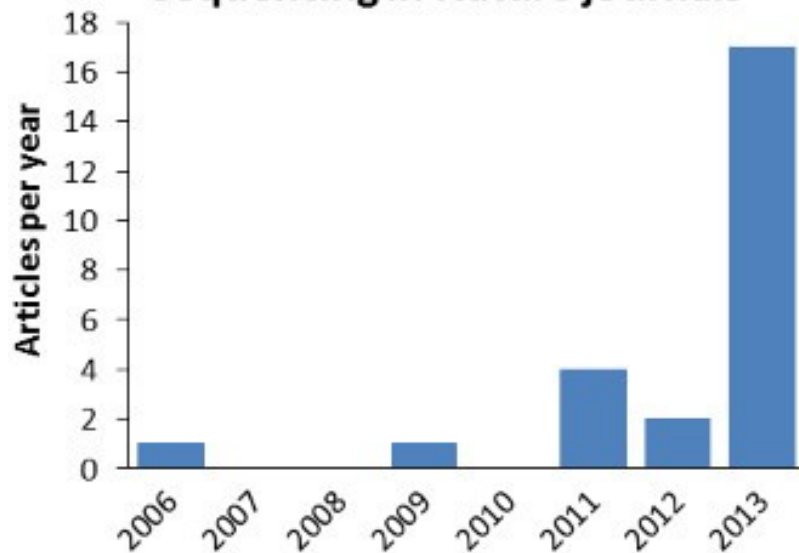
A single-cell genome image of polytene chromosomes from insects  
from 1882 monograph by Flemming

# 2013 METHOD OF THE YEAR

Methods to sequence the DNA and RNA of single cells are poised to transform many areas of biology and medicine.

--- Nature Methods

Research articles using single-cell sequencing in Nature journals



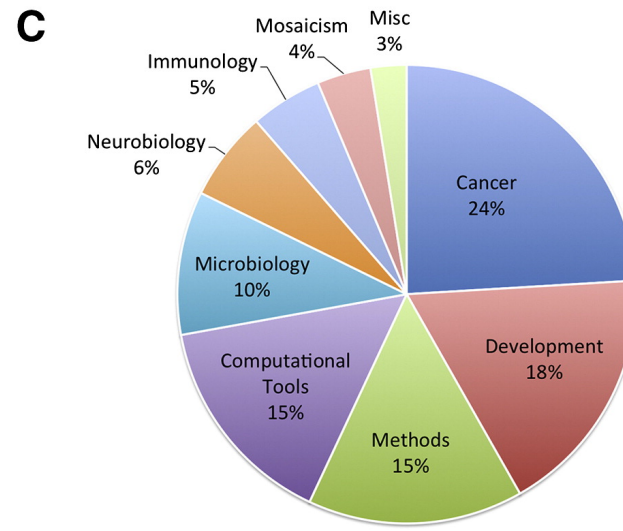
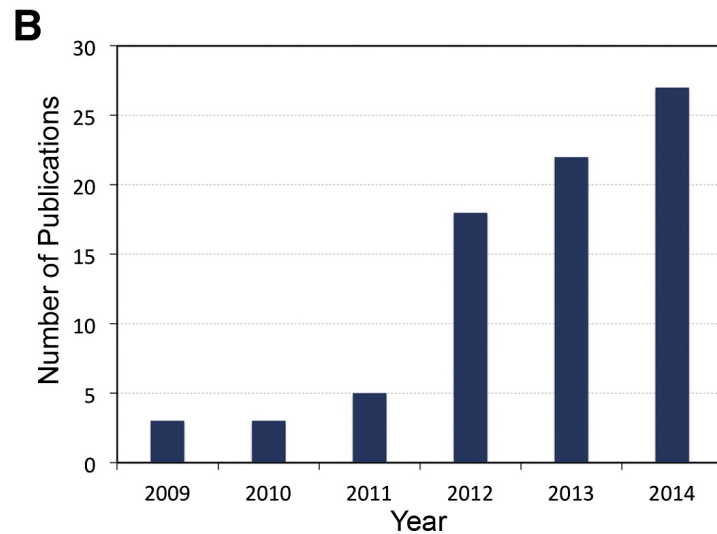
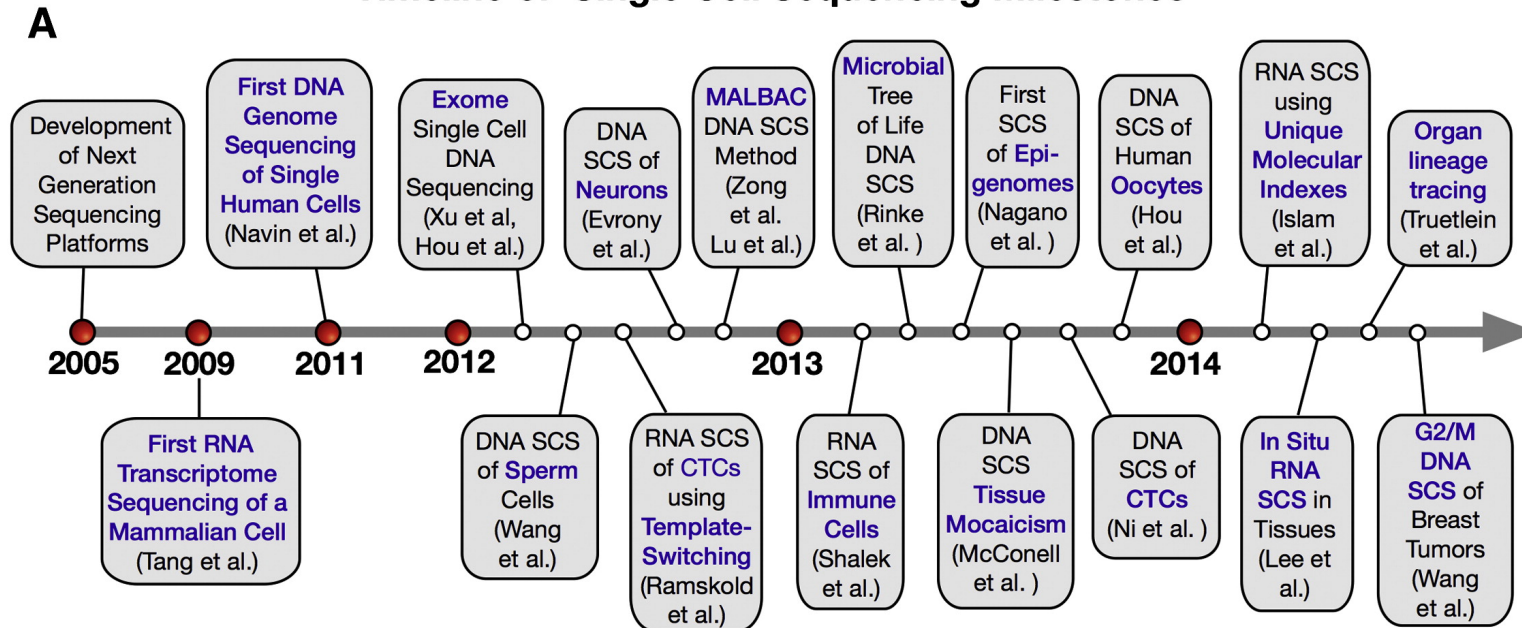
← Significant increase in publications and data in the last 2 years



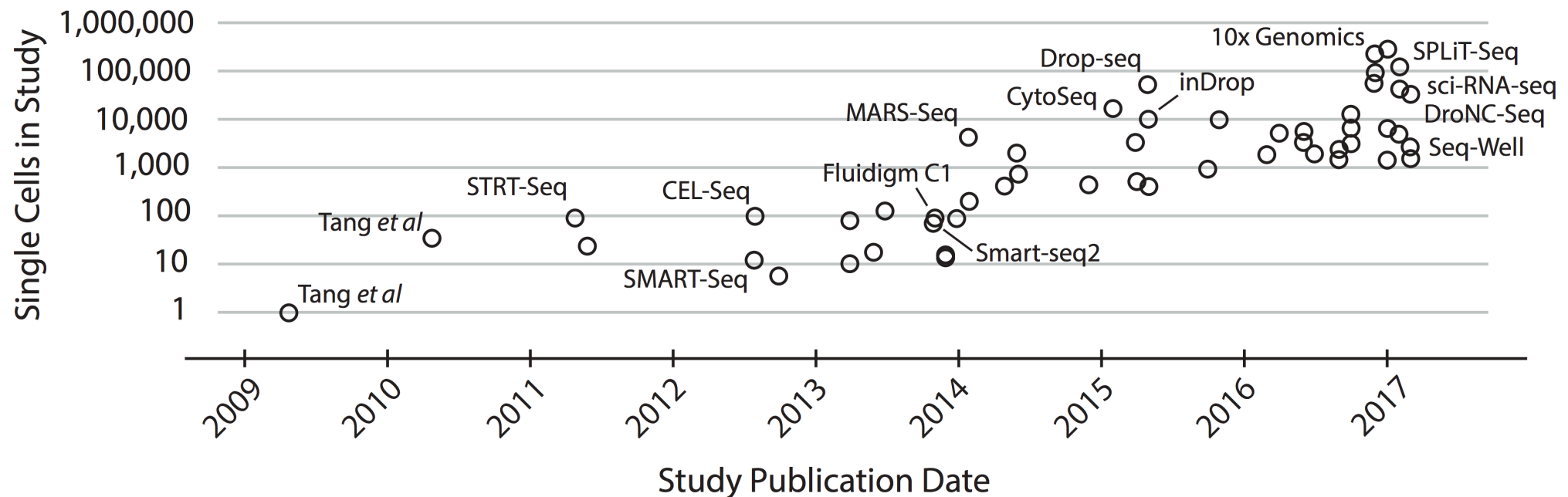
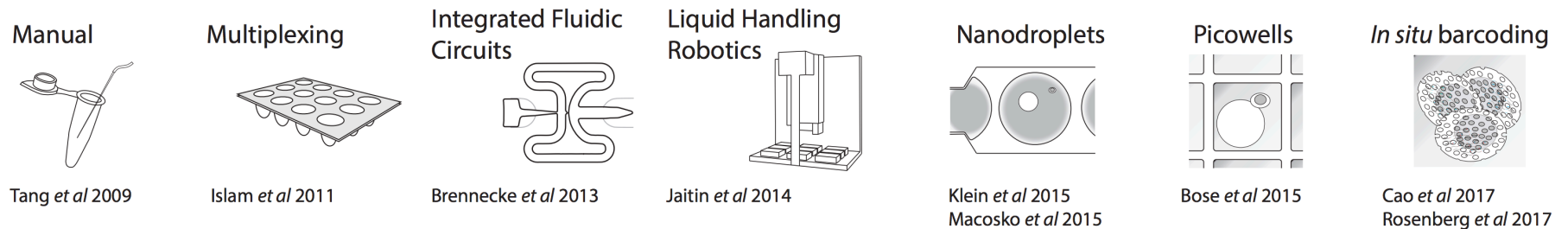


# Advances & Application of single cell sequencing

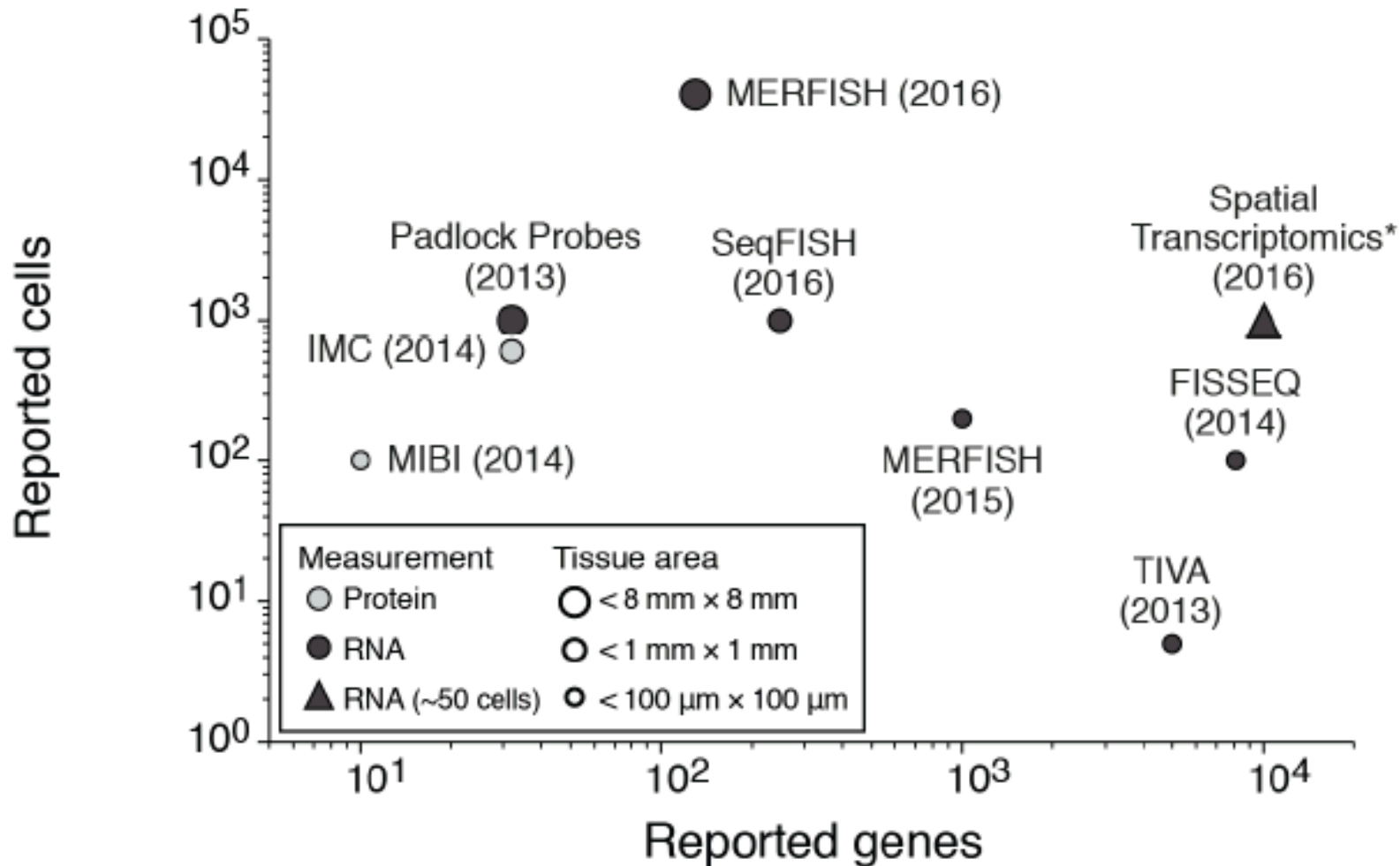
## Timeline of Single Cell Sequencing Milestones

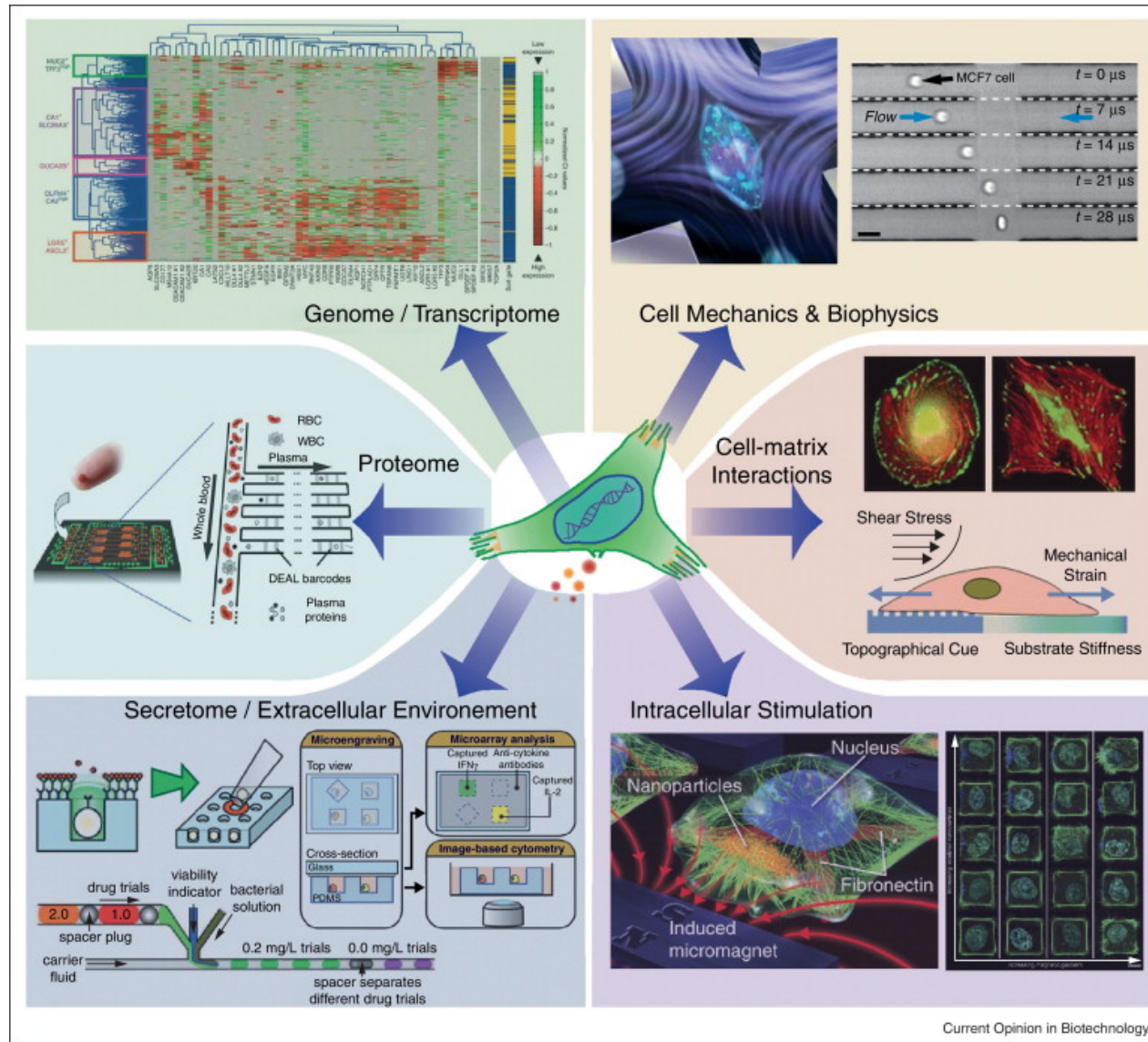


# Technological advances are empowering scalability & additional dimensionalities



# Technological advances are empowering scalability & additional dimensionalities





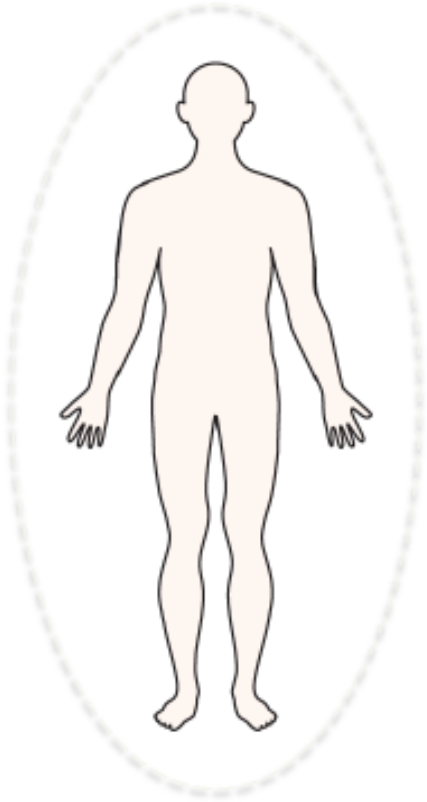
*“Single-cell approaches stand poised to revolutionize our capacity to understand the scale of genomic, epigenomic, and transcriptomic diversity that occurs during the lifetime of an individual organism.”*

Machaulay & Voet 2014

Weaver, 2014

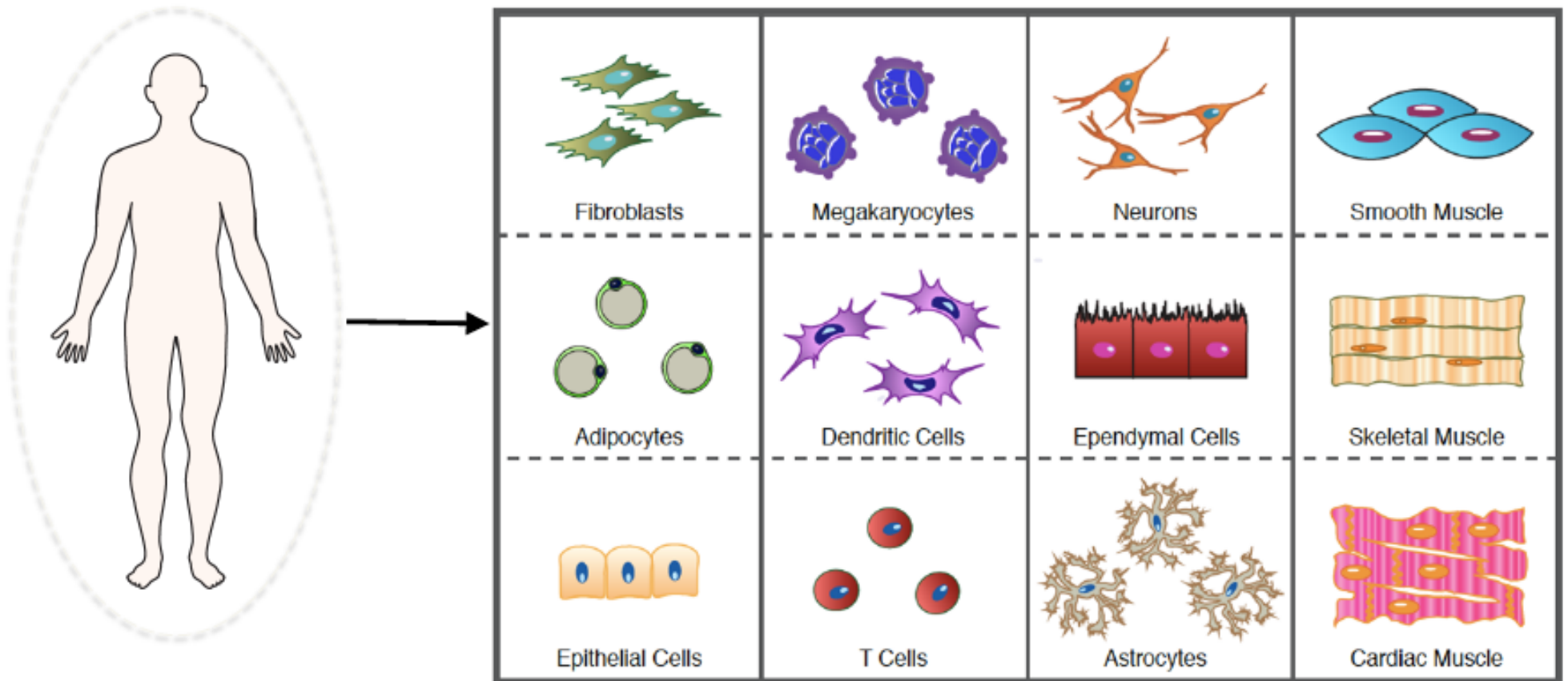


# Do we really know cells defining the human system?



- ~ 30 trillion cells
- Text book → ~ 300 'major' cell types?
- Science → ~ 100 subtypes of immune cells!

# Cells define our core constituents



**How do we define and classify cell type?**

# How do we define and classify cell types?

molecular markers

morphology

spatial localization

physical properties

functions

developmental origins

transcription factor dependency

growth factor dependency

chromatin states

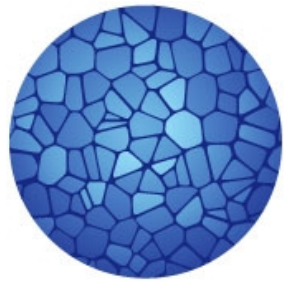
biochemical states

...

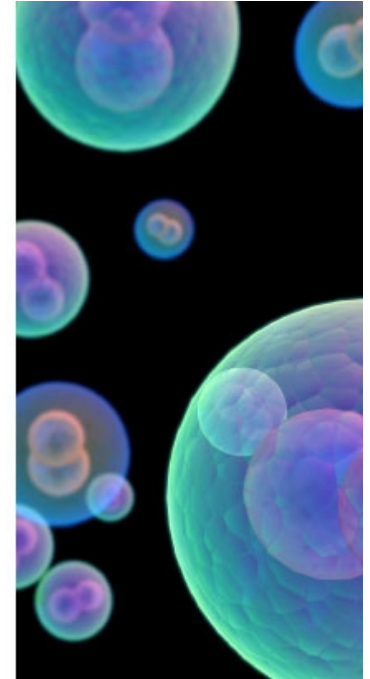
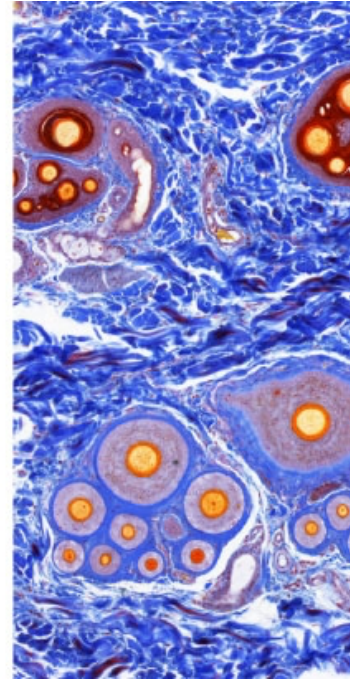
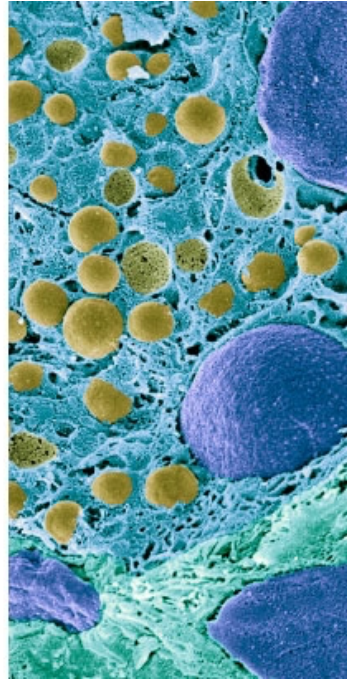
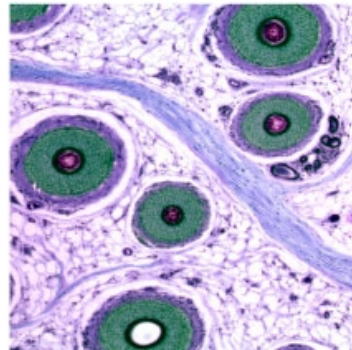
# Limitations of current cell type/state definitions

- **Purity:** Defined cell types may not be pure using the historically defined markers
- **Species:** The more well-defined mouse cell types may not directly translate to human
- **Variations:** An immune response induces new and unexpected states
  - Do existing 'standard' set of surface markers truly define distinct immune cell types?
  - Are there more cell subsets that are not currently appreciated?

**Solution: Leveraging the power of single cell profiling to generate map *de novo* & integrate legacy knowledge**

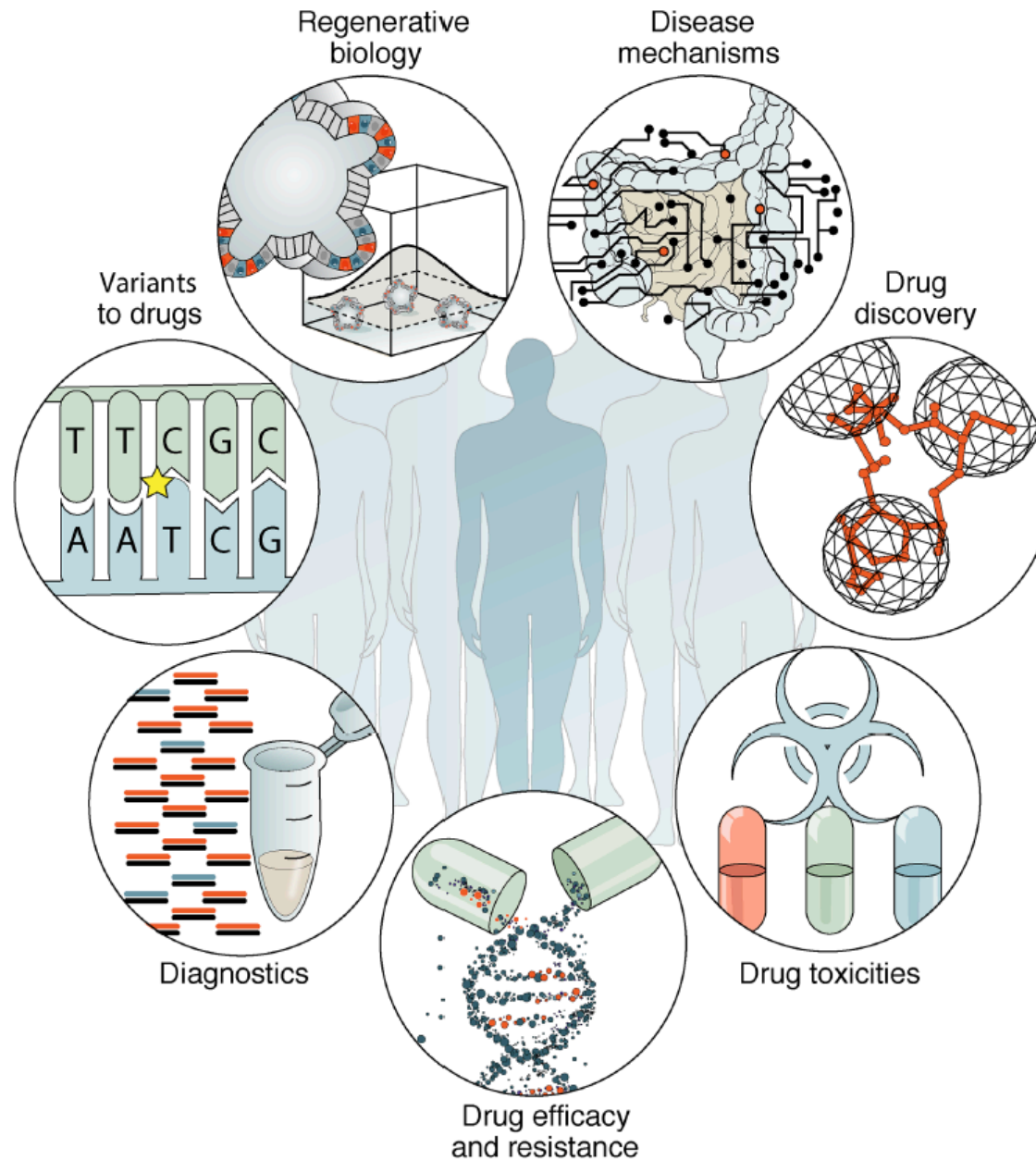


# HUMAN CELL ATLAS



**Mission:** To create comprehensive reference maps of all human cells—the fundamental units of life—as a basis for both understanding human health and diagnosing, monitoring, and treating disease

# Redefining the human system at single cell resolution has tremendous potential for biology & medicine



# What can we learn from single cell

- **Taxonomy & Census** → data-driven molecular definition of cell types & dissection of tissue heterogeneity
- **Anatomy & Physiology** → spatial structure of tissue
- **Pathology** → defining disease cells and associated ecosystem
- **Physiology** → dissection of temporal changes, responses to challenges (e.g. drug treatment)
- **Developmental biology** → cell fate / lineage mapping
- **Molecular mechanisms** → cellular circuitry

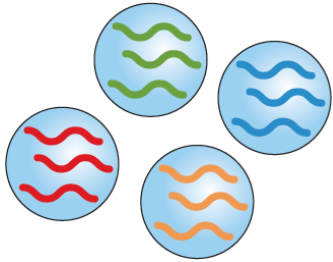


**First critical step → cell isolation**

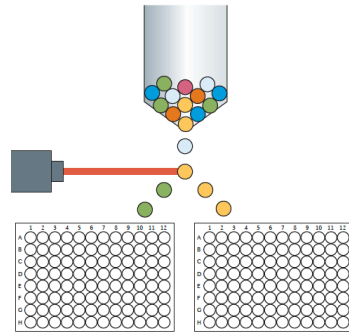


Cell Isolation

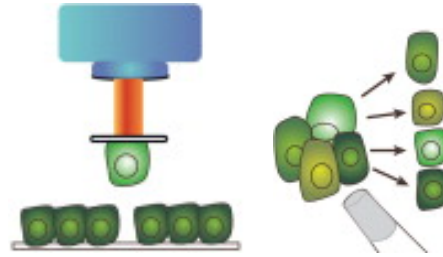
Trapping cells in droplets



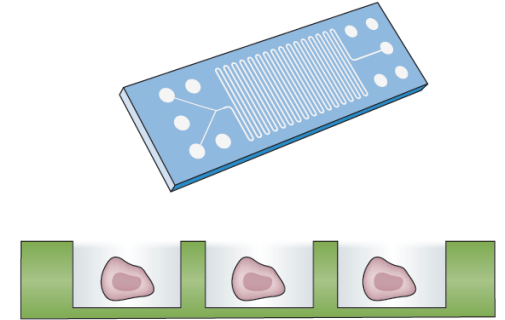
FACS / cell sorter



Microdissection & micromanipulation



Microfluidics & microwells

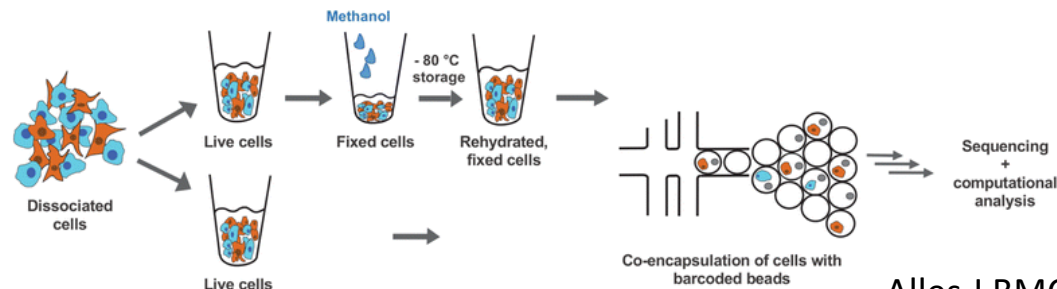


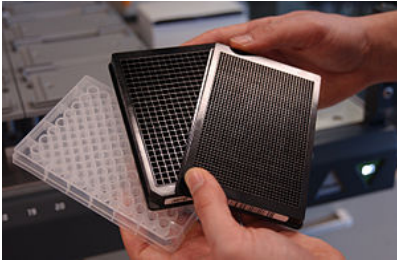
Applications  
Sequencing  
method

Amplification  
method

# Common considerations for sample collection & dissociation

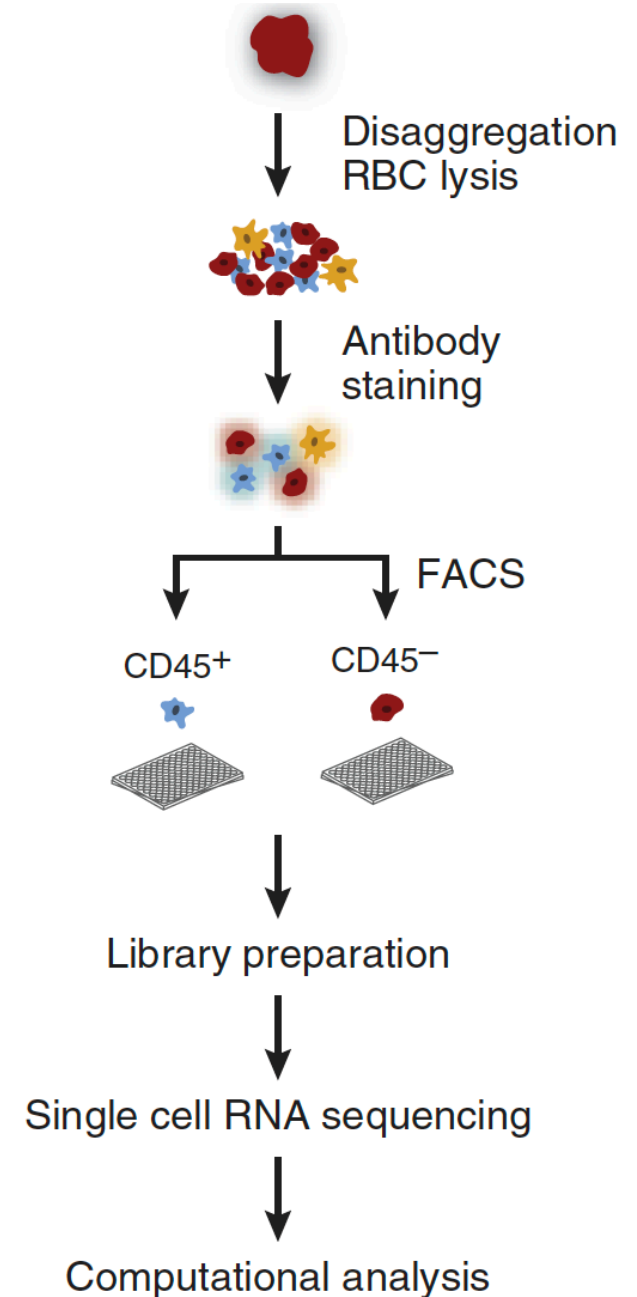
- **Fresh vs. Frozen** → cells vs. nuclei (e.g. considering multi-sites study?)
- **Cell dissociation optimization**
  - Minimizing leakage and RNA degradation
  - Need to optimize for every tissue → e.g readouts: FACS & bulk sequencing
  - Challenging dissociation? Consider LCM & nuclei sequencing
- **Enrichment strategy**
  - Even the sampling to enrich for rare cells (e.g. profiling human blood)
  - Separate immune from non-immune cells (sorting or bead/column)
  - Profiling uniquely T and B cells for TCR & BCR
- **Cell death & RBC removal**
  - Live/death & CD235a marker-based depletion by FACS
  - Magnetic bead depletion-based
  - Column-based (e.g. MACS) depletion → some cell types get caught in columns
- **Work to limit RNA degradation** (fixation protocol work in some case)





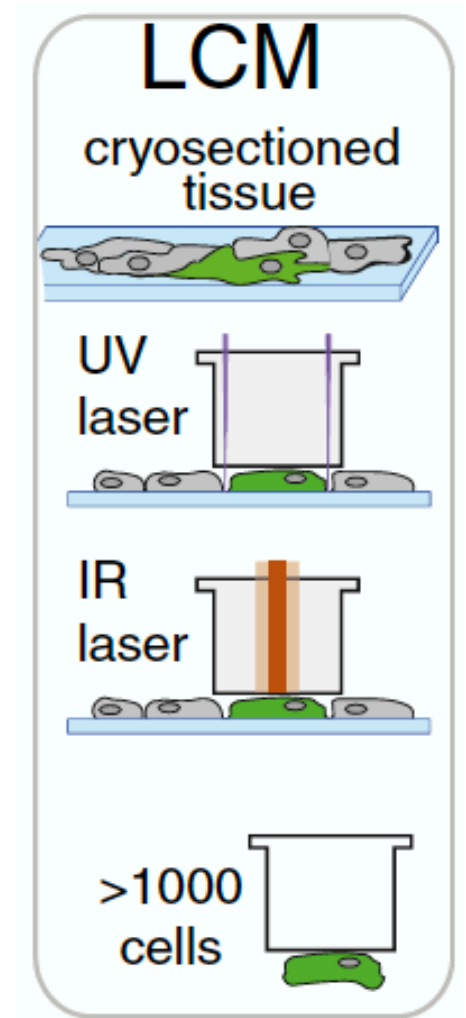
# FACS isolation

- **Advantages:**
  - Sorting based on specific cell phenotype
  - Archiving potential
  - Full-length cDNA readout possible
- **Disadvantages:**
  - Larger amount of cell required
  - Occasional isolation of more than one cells
  - Putative damage of cells (epithelial cells)
  - Labor intensive & more costly
- **Know your cells, are they sticky, are they big?**
  - Select an appropriate sized nozzle
- **Don't sort too quickly (1-2k cells per second or lower)**
  - The slower the more time cells sit in lysis after sorting
  - 10 minutes max in lysis (some say 30 minutes)
- **Calibrate speed of instrument with beads**
  - Check alignment every 5-6 plates
- **Afterwards spin down to make sure cells are in lysis buffer**
  - Flash freeze on dry ice and move to -80C (use very adherent seals for archiving)

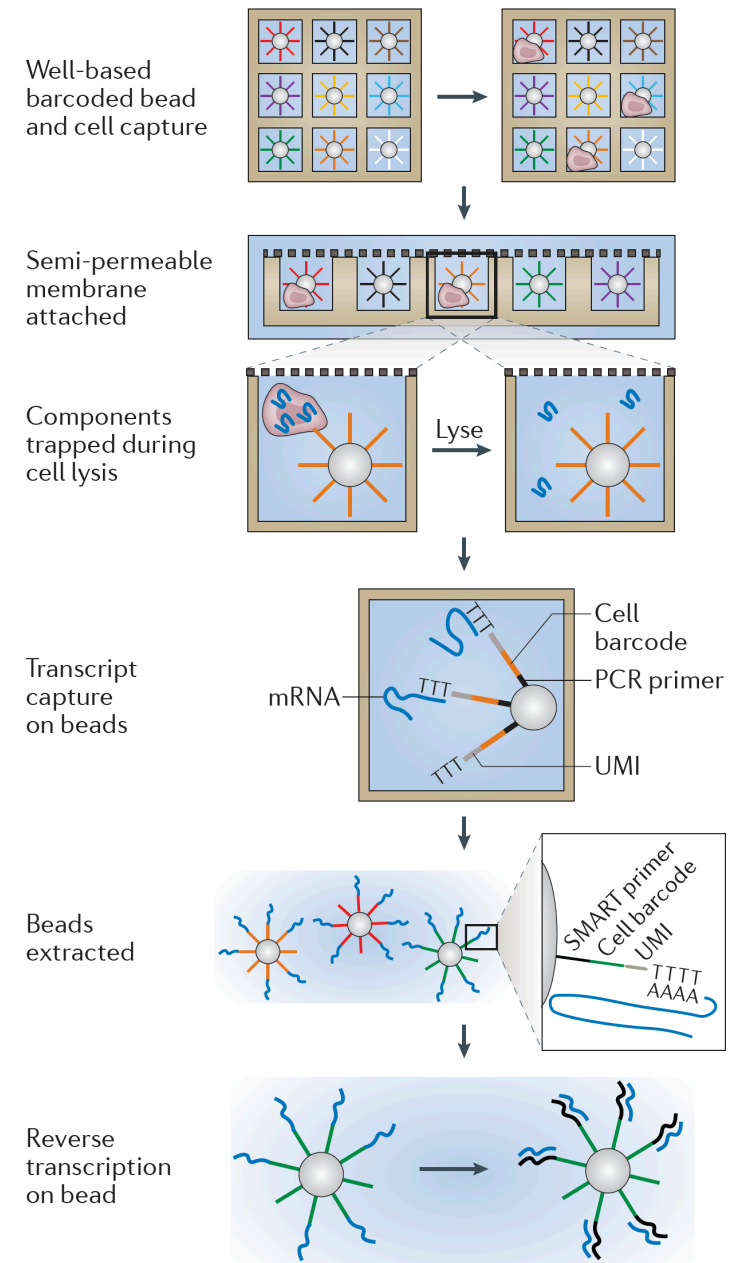
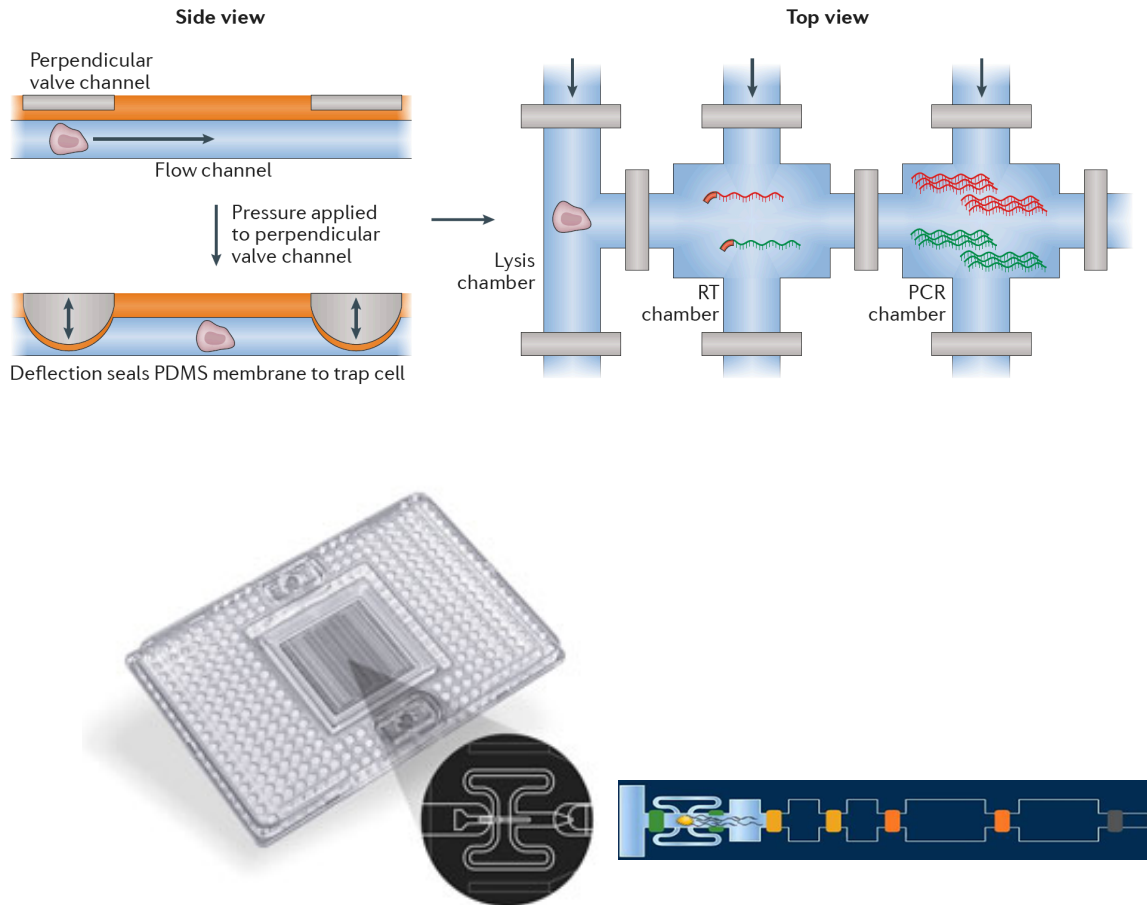


# Micromanipulation & LCM

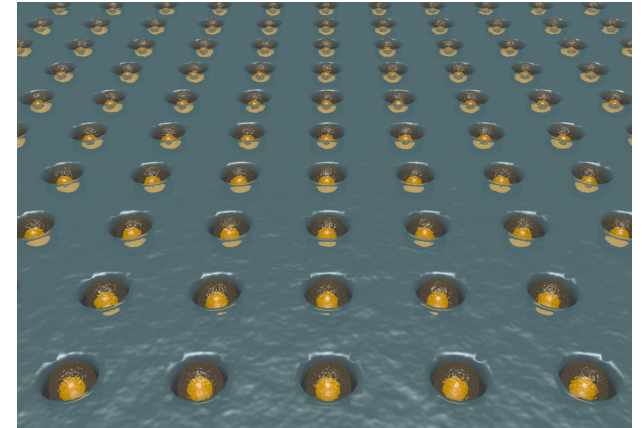
- **Advantages:**
  - Visual confirmation
  - Applicable when only few cells are available
  - Retain topological information of the cell
  - Permits isolation of a cell from fixed tissue or cryosection
- **Disadvantages:**
  - Low throughput
  - Lengthy process → RNA degradation
  - Operator bias
  - Contamination of other cells
  - Potential loss of cellular material (LCM)



# Microfluidics & Micro-wells



# Microfluidics & Micro-wells



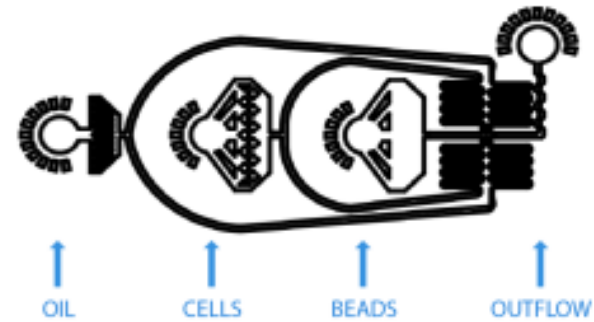
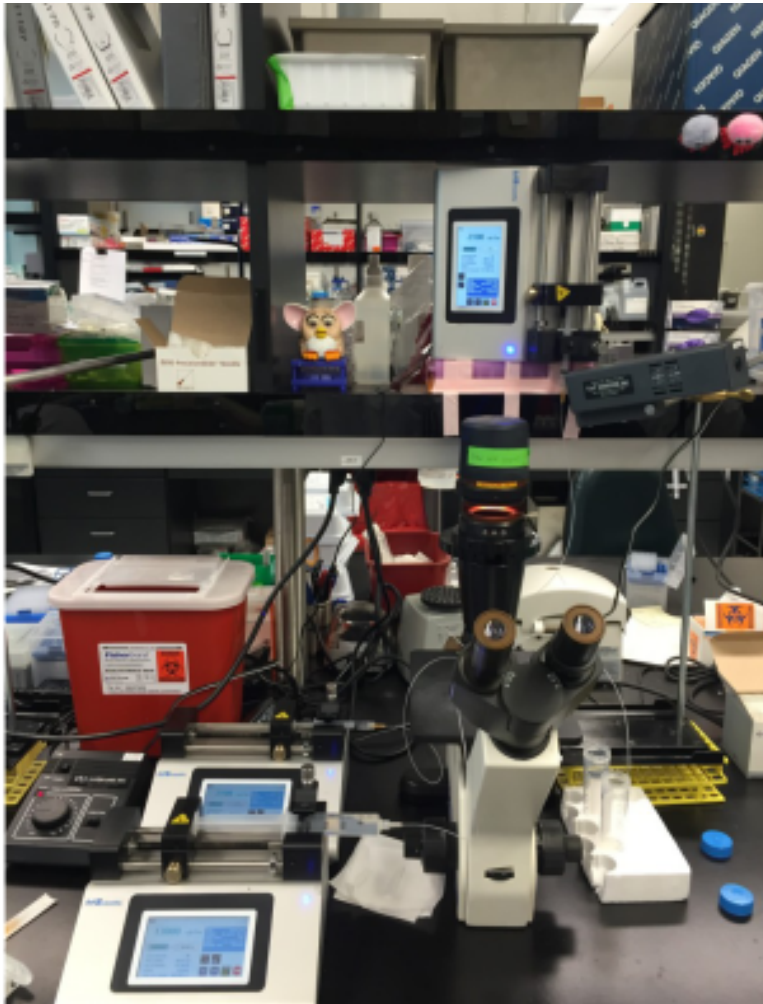
Credit: David Wood

- **Advantages:**
  - Highly standardized nanoliter reaction (lower reagent cost/cell)
  - Less operator bias in cell isolation and enzymatic reactions
  - Automated higher throughput cell isolation with visual confirmation
- **Disadvantages:**
  - Putative loss of cells → capture efficiency lower than if sorting in plates in some cases
  - Cannot select specific cells (unlike cell sorting)
  - Bias driven by cell size and adherence (fixed size devices)
  - Bias driven by cell type frequency (will capture mostly abundant types)
  - In some case still need to enrich first and cells sit around longer before lysis

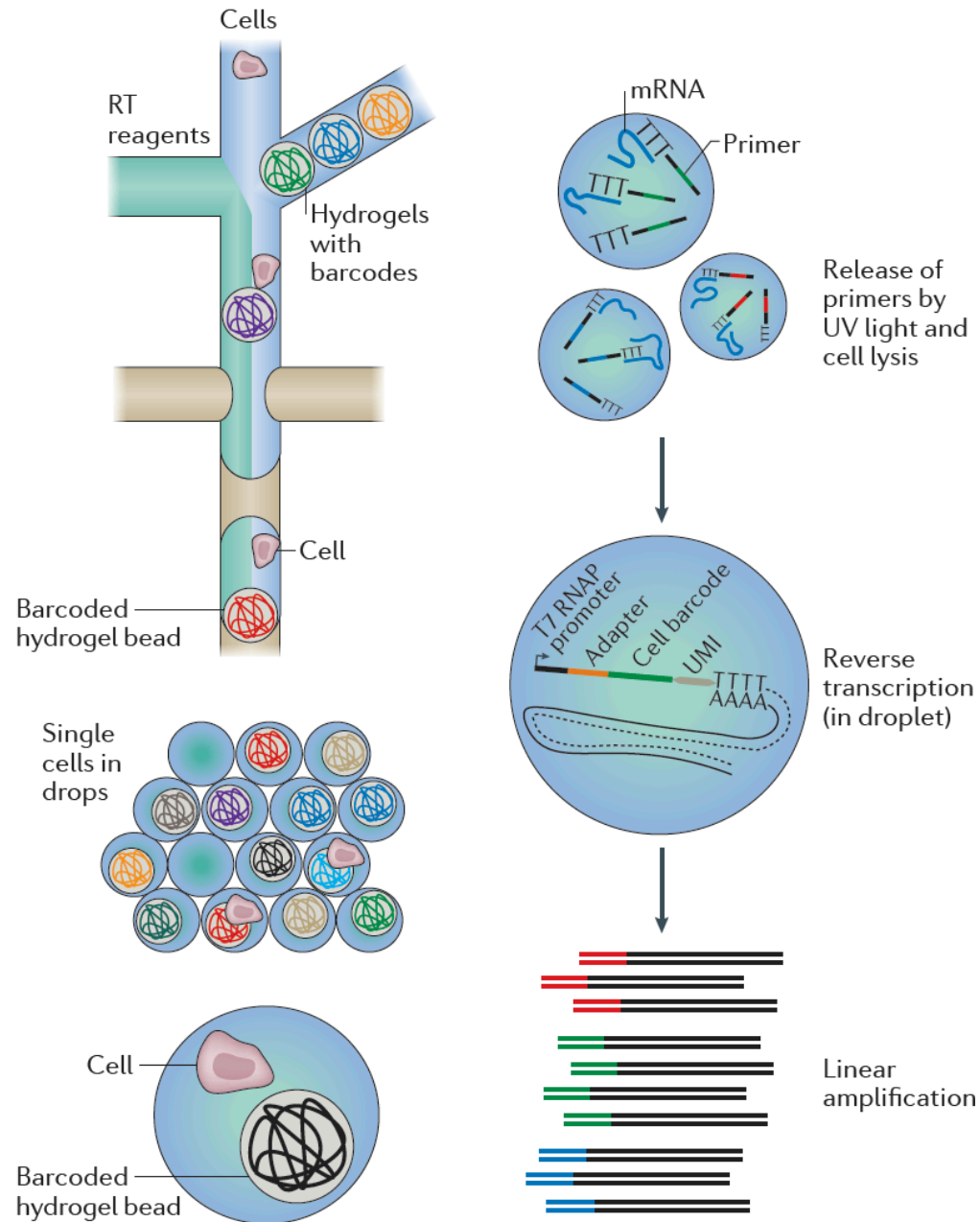


# Emulsion-based / Droplets

## DropSeq setup

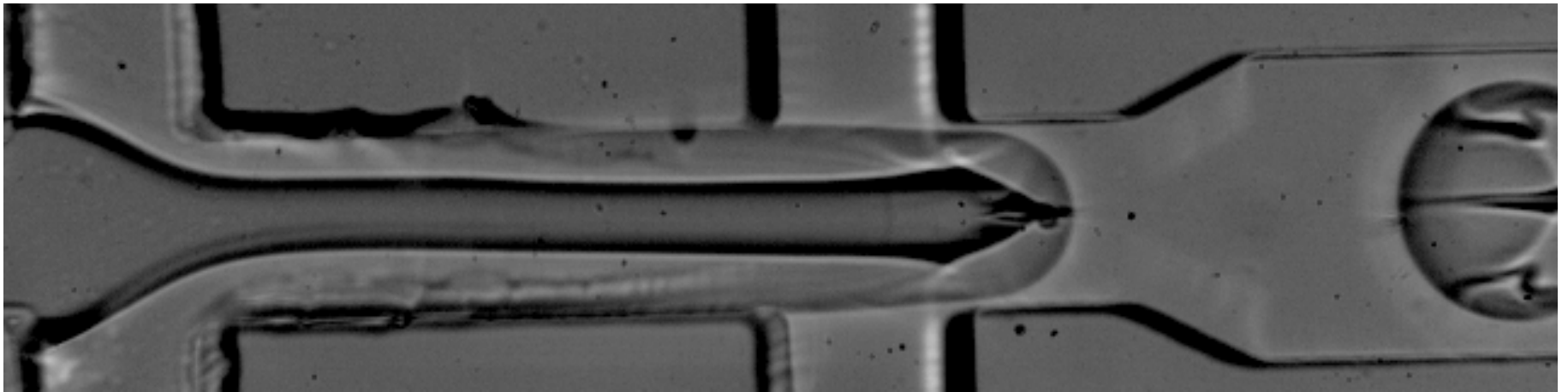


# Emulsion-based / Droplets





# Drop-seq – Overview



# Emulsion-based / Droplets

- **Advantages:**

- Very scalable → thousands of cells per experiment
- Smaller volumes → higher detection & better reproducibility
- Smaller volumes → cheaper reagent cost
- Sequencing cost become bottleneck → often shallow sequencing

- **Disadvantages:**

- High cell input required (DropSeq) though low cell capture
- Variable quality of beads → can increase cost
- Need to be familiar with microfluidics (unless opt for commercial option like 10X)
- Droplet-based assays can have leaky RNA (unlike plate → compartmentalization)
- Capture less transcripts than plate-based (lower resolution)
- Only 3' end readout

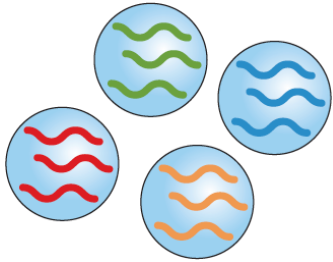
- **Some pointers:**

- Before library generation wash off any medium (inhibits library generation)
- Adding PBS & BSA (0.05-0.01%) can help protect the cells
- Filter all reagents with micron strainer before loading on microfluidic
- Some purchased devices come with hydrophobic coating  
→ Can deteriorate (2 months at best) → recoating works

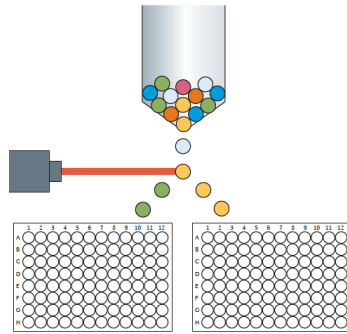
# Selecting scRNAseq protocol

## Cell Isolation

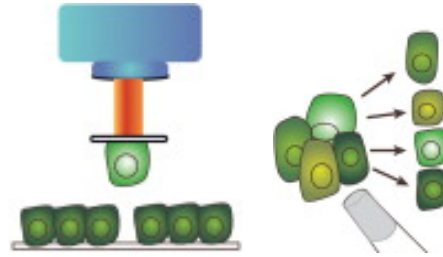
Trapping cells in droplets



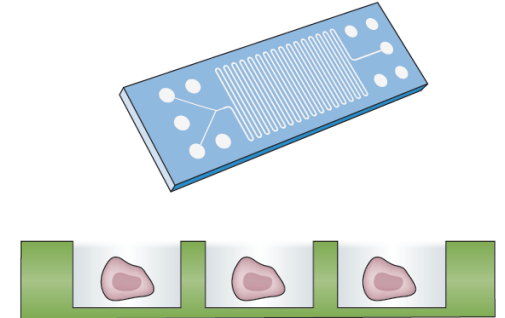
FACS / cell sorter



Microdissection & micromanipulation

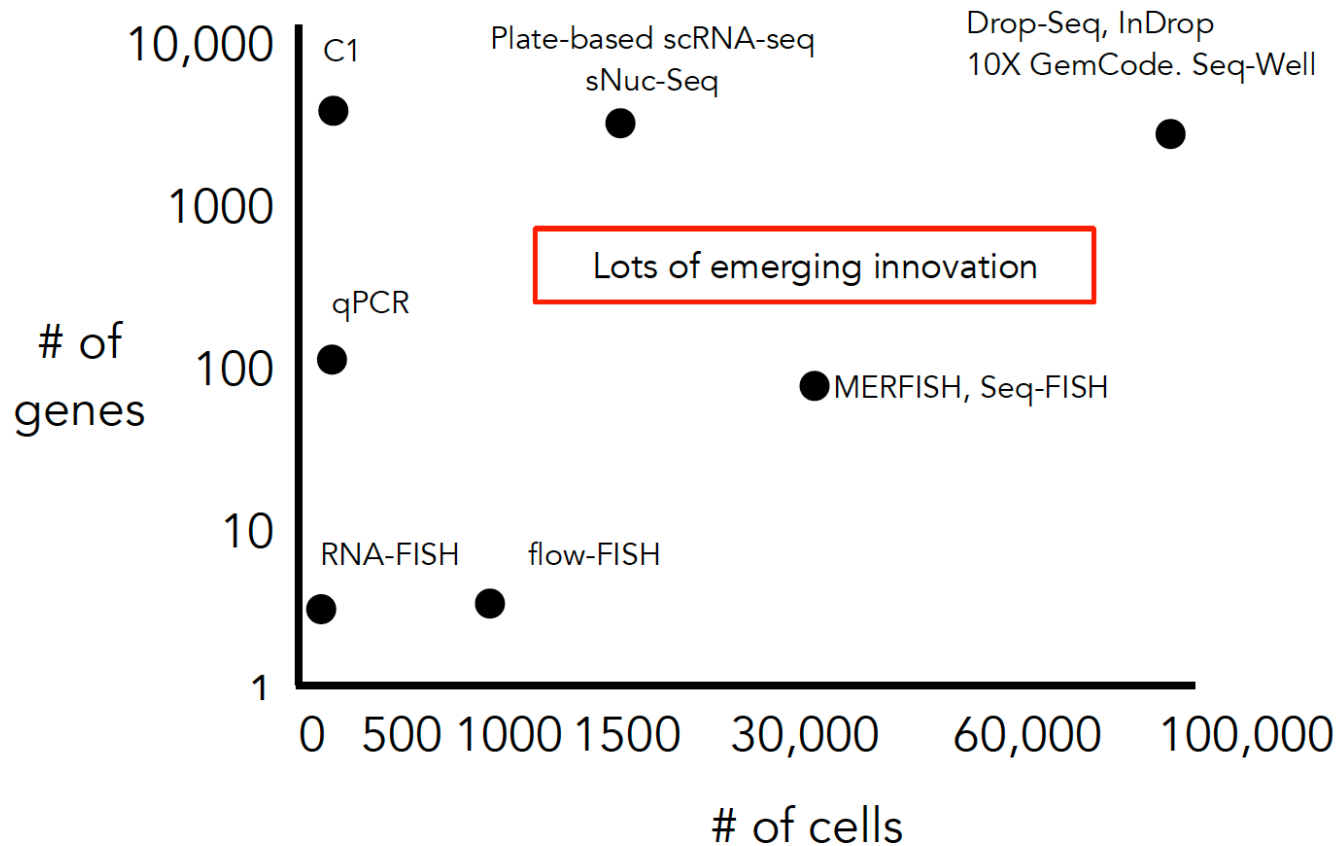


Microfluidics & microwells



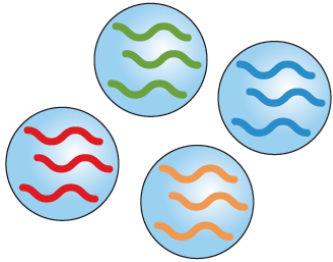
## Applications Sequencing Amplification method

### Tradeoff between scale & resolution

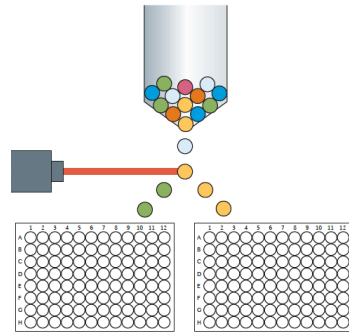


## Cell Isolation

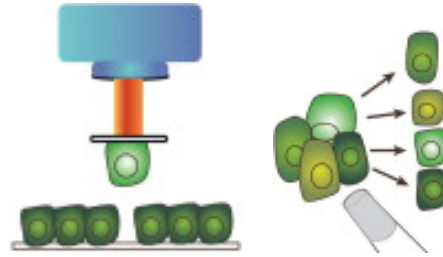
Trapping cells in droplets



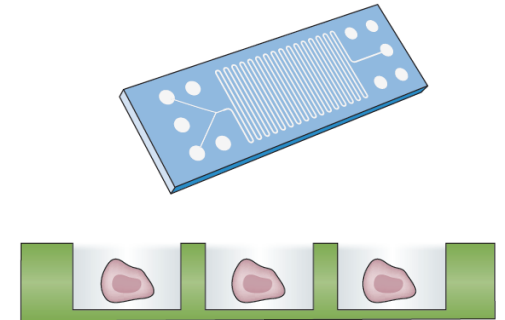
FACS / cell sorter



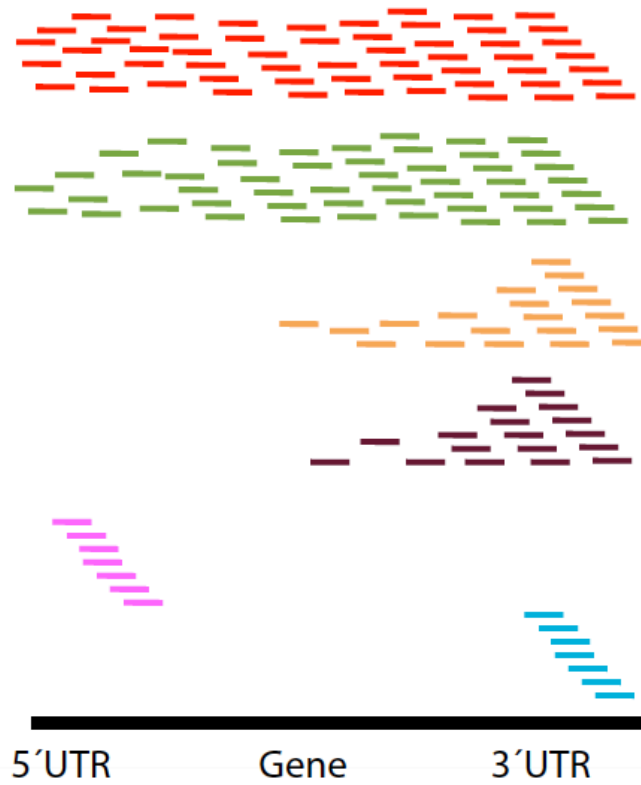
Microdissection & micromanipulation



Microfluidics & microwells



## Applications Sequencing Amplification method



SmartSeq2  
(Picelli et al. Nature Methods 2014)

SmartSeq – SMARTer kit  
(Ramsköld et al. Nature Biotech 2012)

Quartz-seq  
(Sasagawa et al. Genome Biology 2013)

Tang et al.  
(Nature methods 2009)

STRT  
(Islam et al. Genome Res 2011)

CEL-Seq  
(Hashimshony et al. Cell Reports 2012)

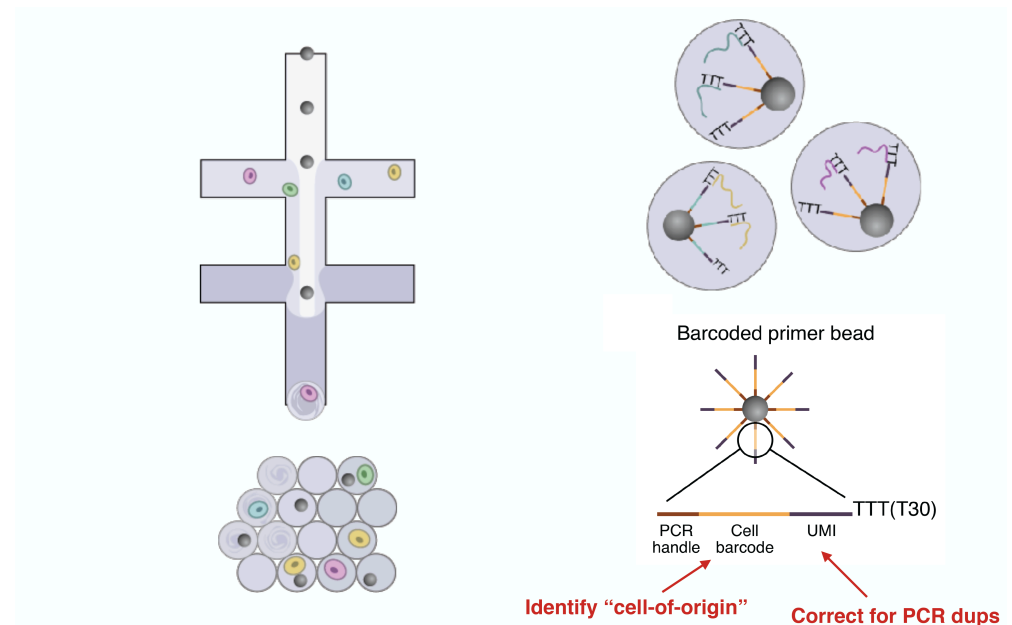
# Unique molecular identifies (UMIs) and cellular barcodes

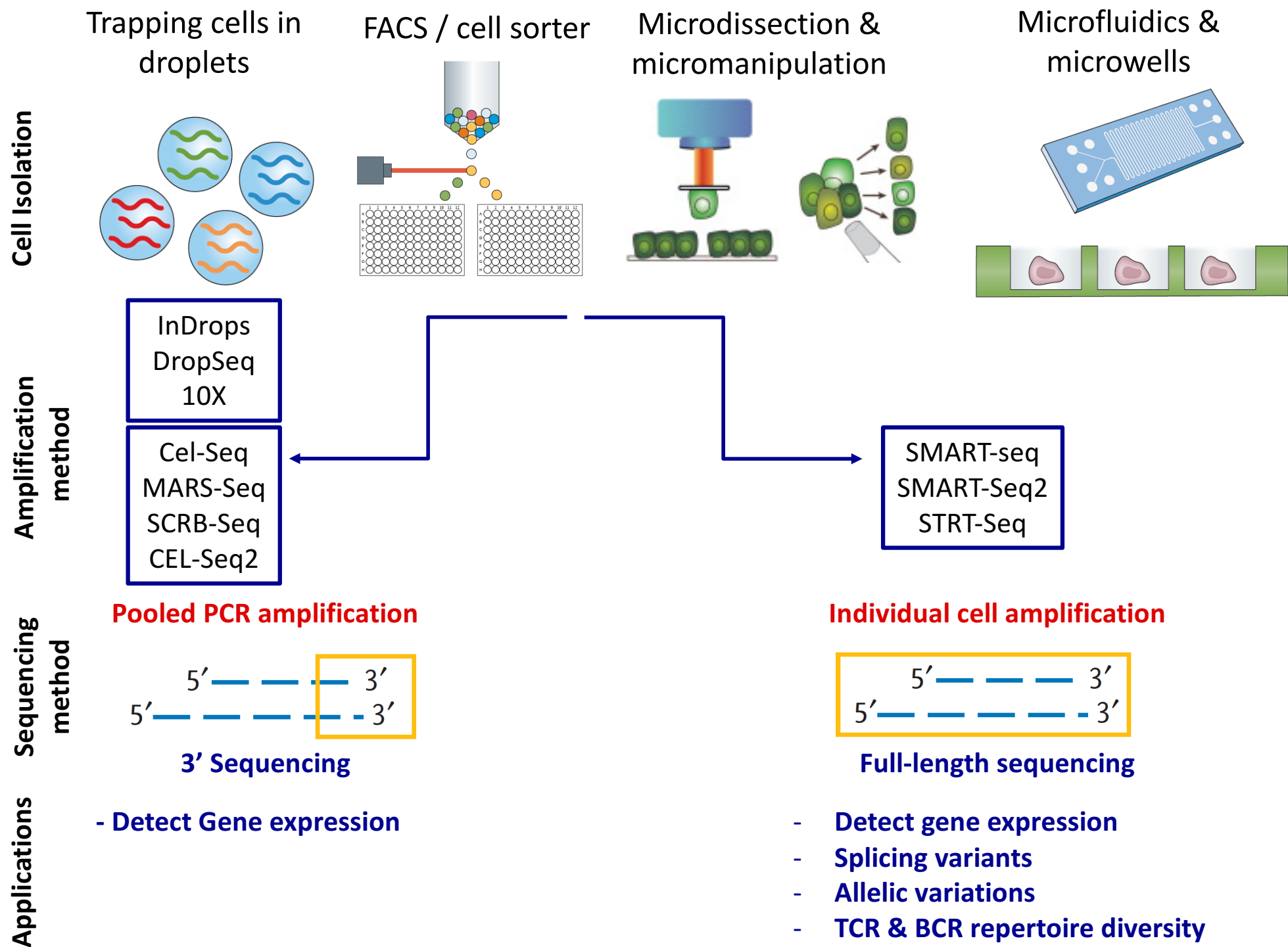
- **Cellular barcodes**

- Introduced at RT step with one unique sequence per cell
- Enables pooling many libraries into one tube for subsequent step (reduces cost & technical errors)

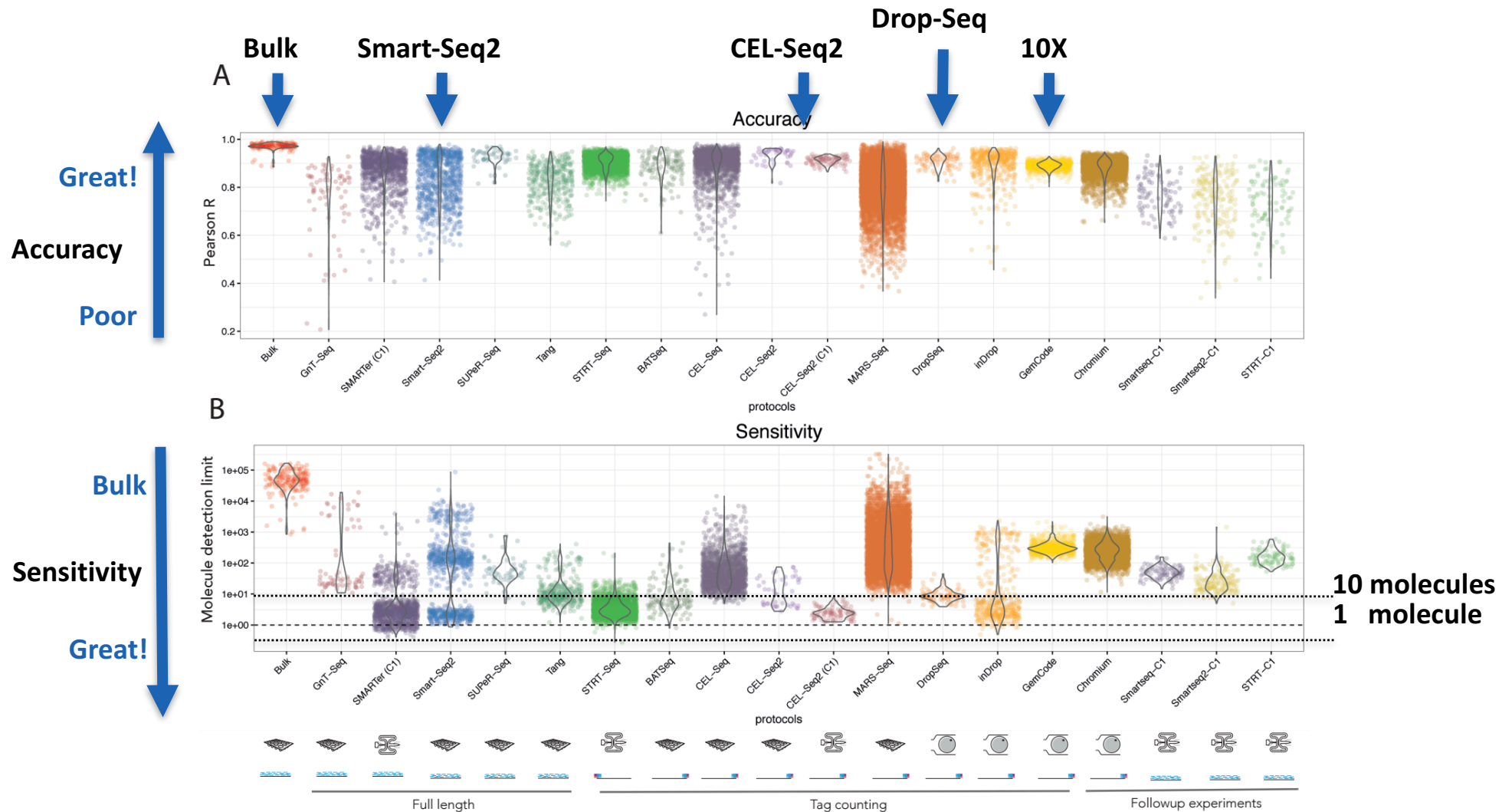
- **UMIs**

- Introduce random sequences at the beginning of each sequence
- Reduces effect of amplification bias by removing PCR duplicate





# Sensitivity and Specificity

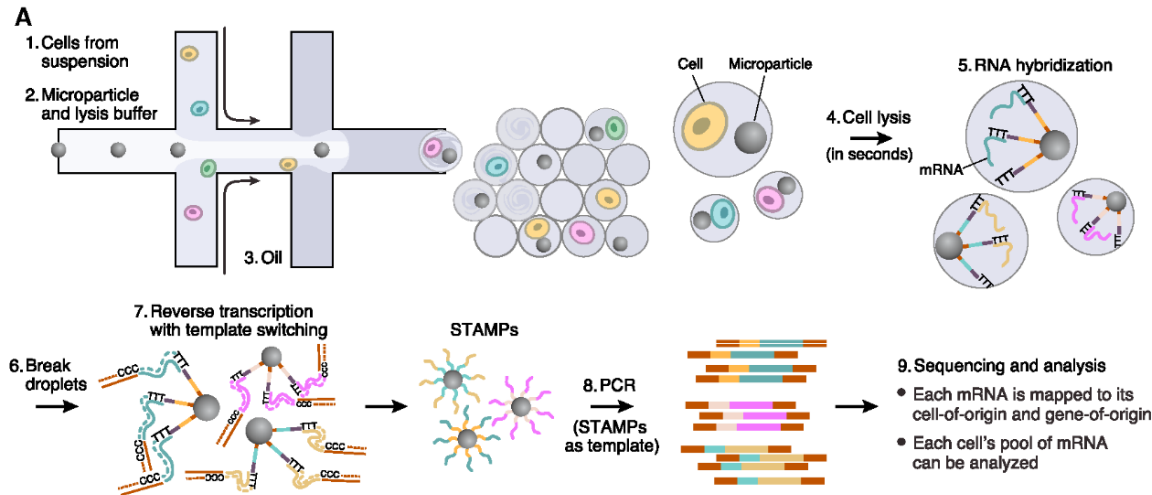


- All better than bulk
- Many between 1 and 10 molecule detection
- Sensitivity dependent on sequencing depth → can sequence more!
- Sensitivity = critical when studying lowly expressed genes



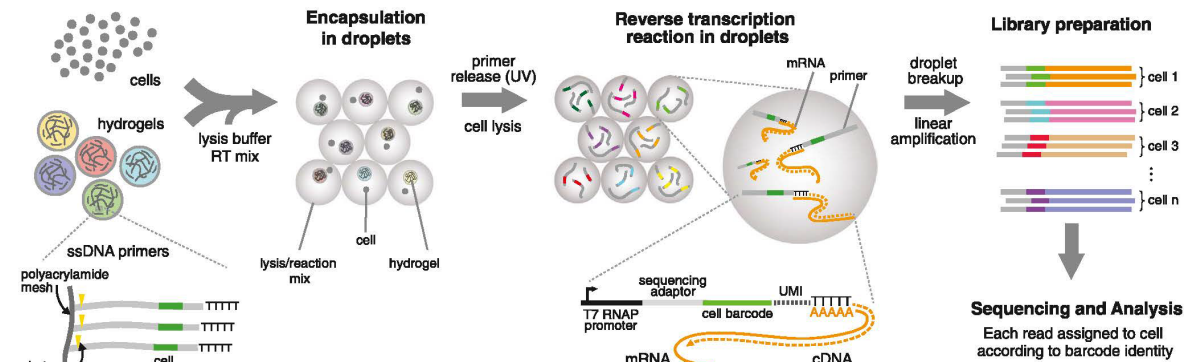
# Scalability – Massively parallel scRNAseq approaches

## DropSeq



Macasko et al. Cell 2015

## InDrops



Klein et al. Cell 2015

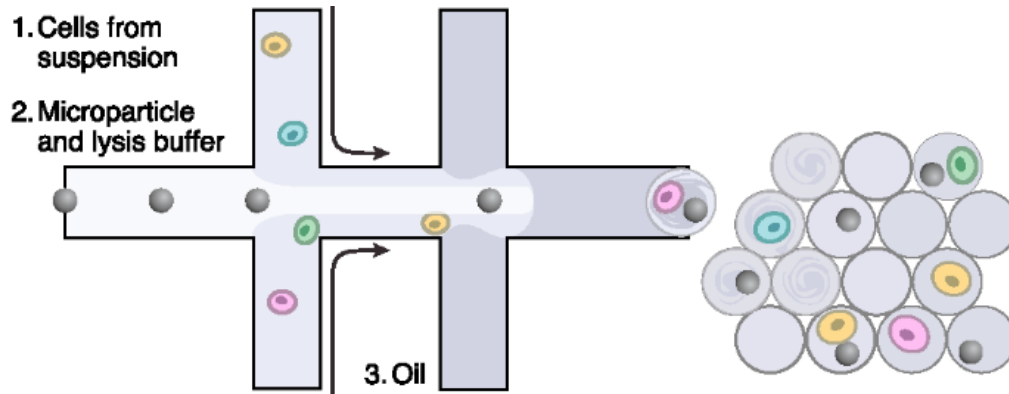
- Cell lysed in the drop & hybridize to primers attached to beads
- STAMP: single cell transcriptomes attached to microparticles
- Droplets are broken & RT/template switching occurs on pool

- Lysis and reverse transcription occurs in the droplet
- Samples are frozen after RT as RNA:DNA in gel

Adapted from Boswell S.  
<https://iccb.med.harvard.edu/single-cell-core>

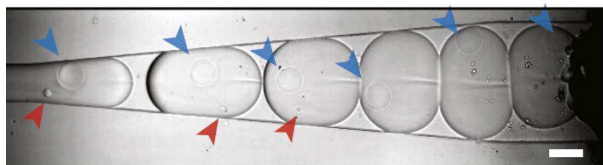
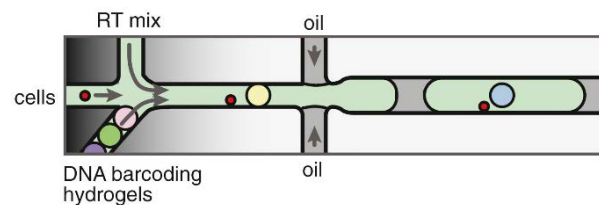
# DropSeq vs. InDrops

## DropSeq



- 1/10 droplets contain microparticle
- 1/10 droplets contain cell
- 1/20 droplets contain both cell and microparticle

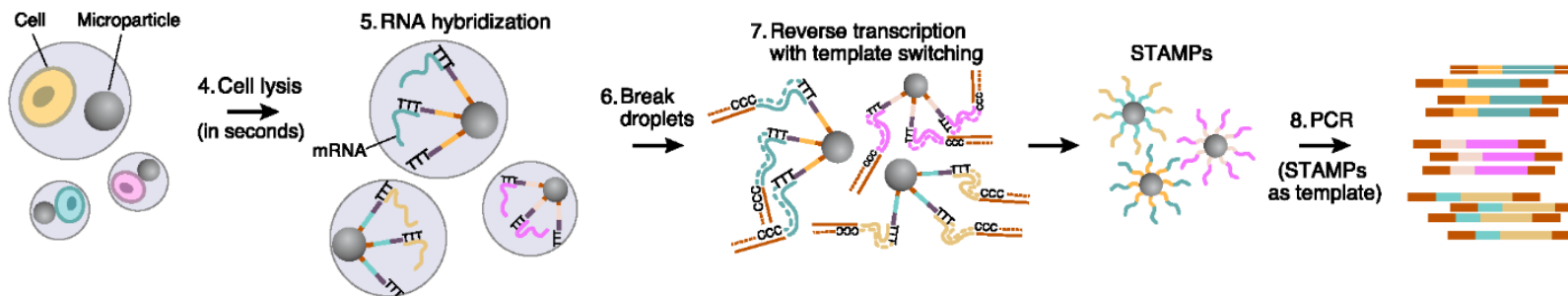
## InDrops



- Match speed of bead injection with speed of droplet generation
- Nearly every droplet loaded with one barcode

# DropSeq vs. InDrops

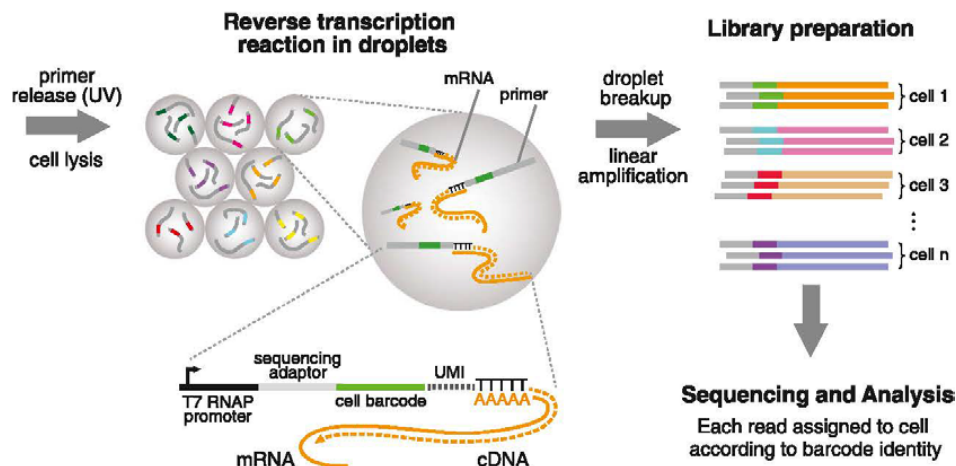
## DropSeq



→ Smart-Seq: RT/template-switching to tagmentation

→ Immediate lysis

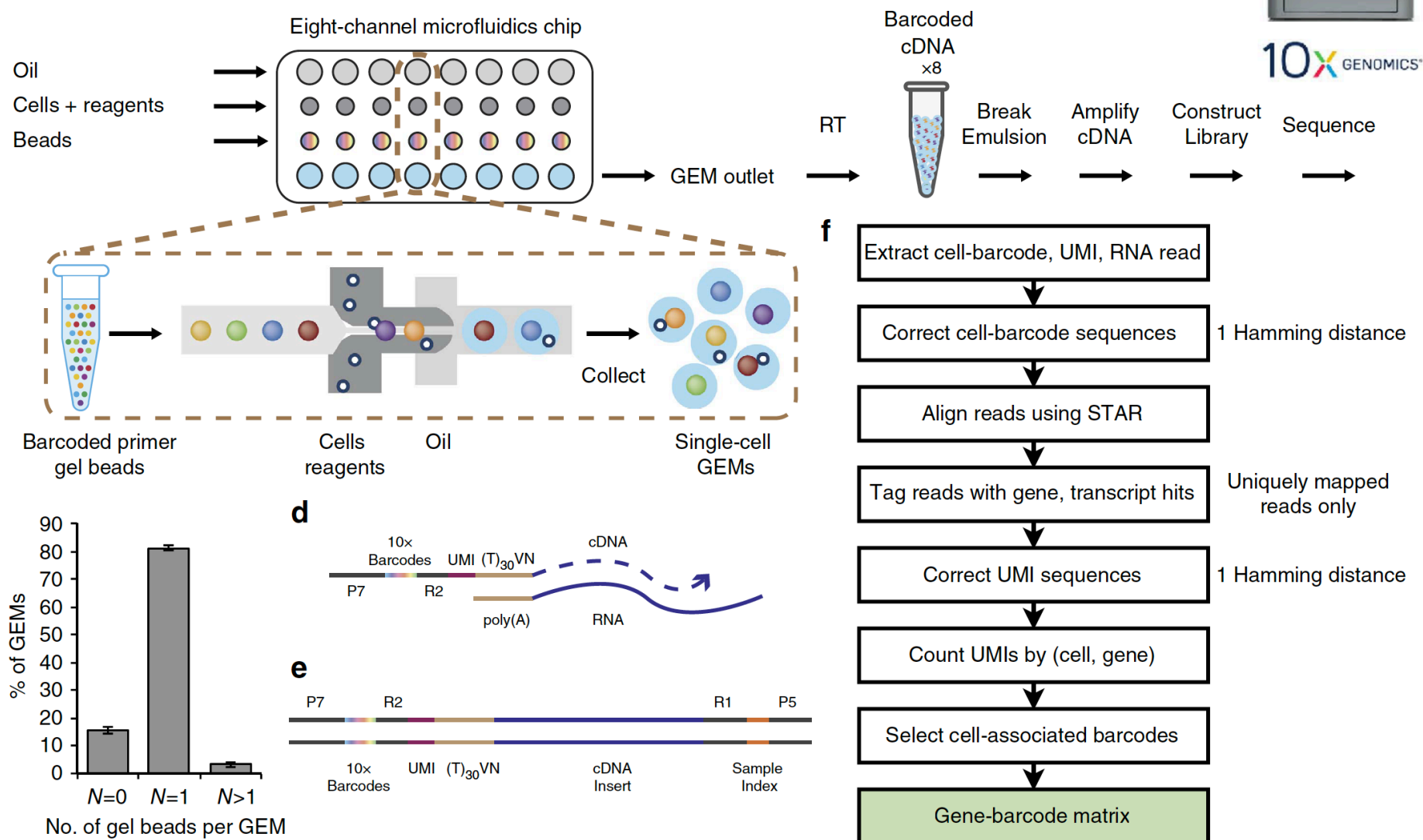
## InDrops



→ CEL-Seq: RT/second strand synthesis to IVT and RNA fragmenting

→ Gentle lysis that may not be completed until sample collection finished

# 10X Genomics 3' mRNA sequencing



A higher throughput “plug & play” version of InDrop

# inDrops, DropSeq, 10X Genomics

## 3' mRNA sequencing

	Capture Efficiency*	Doublet Rate	Number of samples at once	Library prep
InDrops	50-90%	3%	1	CEL-Seq
10X	50-60%	3%	8	CEL-Seq
DropSeq	5-10%	10%	1	Smart-Seq

*\* Capture efficiency is of the cells that reach the device*

- InDrops and 10X are very similar technologies
- InDrops & DropSeq → more labor intensive but customizable & cheaper; need some expertise in handling microfluidics
- 10X → more scalable (8 samples in parallel), “plug & play”, comes with standardized pipeline, but much more expensive (upfront cost \$25k)
- DropSeq requires 100,000 cells as input vs. 7,000 cells for 10X
- Number of transcripts detected varies between approaches (also tissue dependent)
- Cost per library varies greatly!

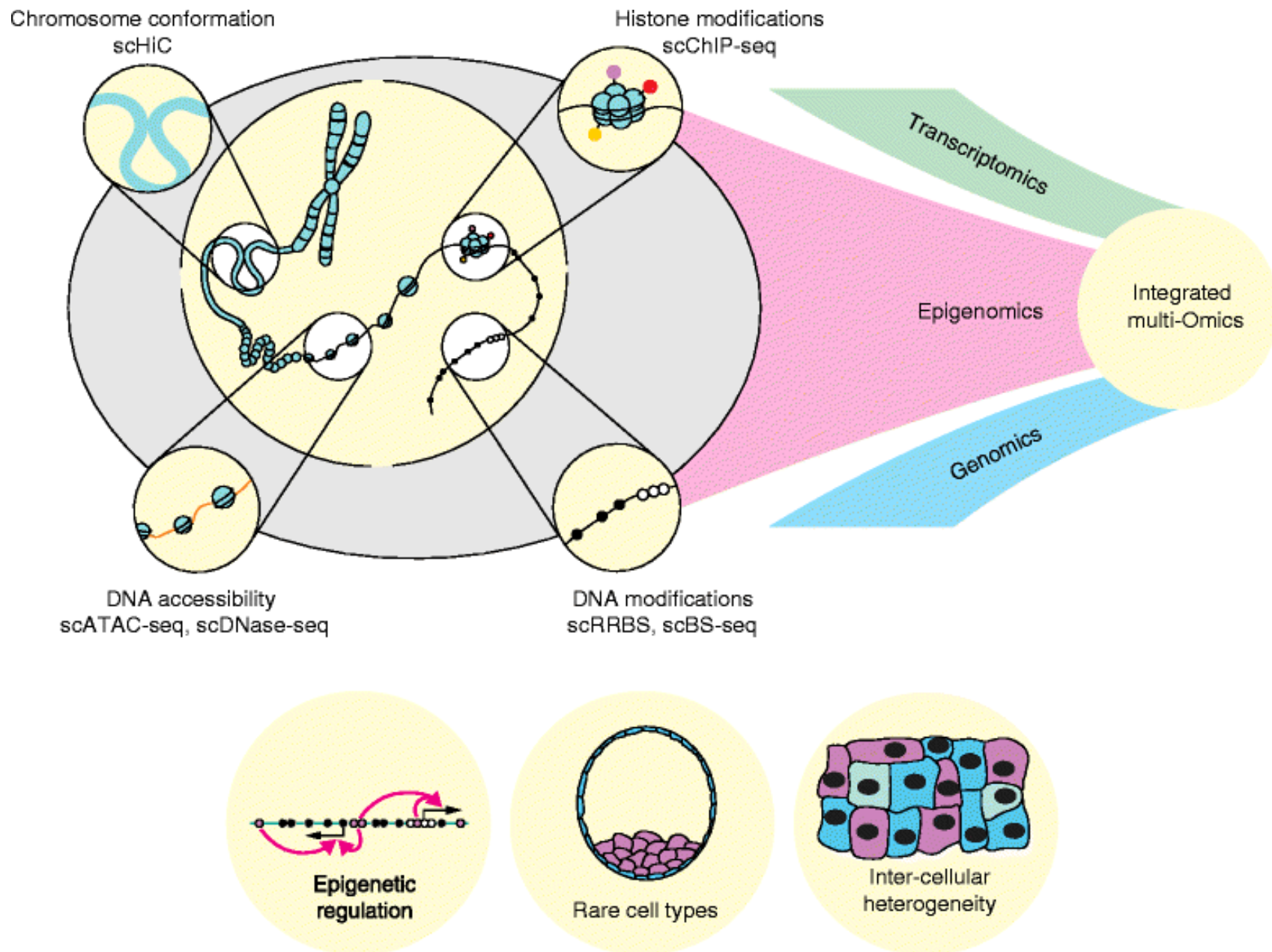
# Your biological question will dictate which method(s) to pursue

- **Different scRNAseq have pros and cons**
- **Needing scalability**
  - Do you know which cell type you want to study?
  - Looking to generate cell census?
  - Are you trying to map very rare cell subsets for which you do not know markers?
  - Dissecting tissue (healthy/disease) ecosystem?
  - Mapping response to treatment (pre vs. post), not knowing which cells would be affected?
- **Needing higher sensitivity and/or full-length transcripts**
  - Predicting binding specificity of TCR receptors?
  - Interested in studying a particular population, potentially rare?
  - Want to map at higher resolution the transcriptome of signaling components/less abundant transcripts to dissect particular biology / pathway?
  - Interested in mapping allelic expression, x-chromosome inactivation, or spliced isoforms?
- **Interested into lineage tracing?** Consider combining DNA/RNA seq and ATACseq
- **Trying to decipher interacting partners** → spatial omics
- Consider **combining different approaches** in your study design!

## **Other single cell readouts & multi-omics approaches**



# Single cell epigenomics



# Technological evolving landscape – stay tuned!

## Split & pool barcoding (not relying on microfluidics)

- SPLiT-Seq (Rosenberg et al. 2017)
- SCI-Seq (Cao et al. 2017)

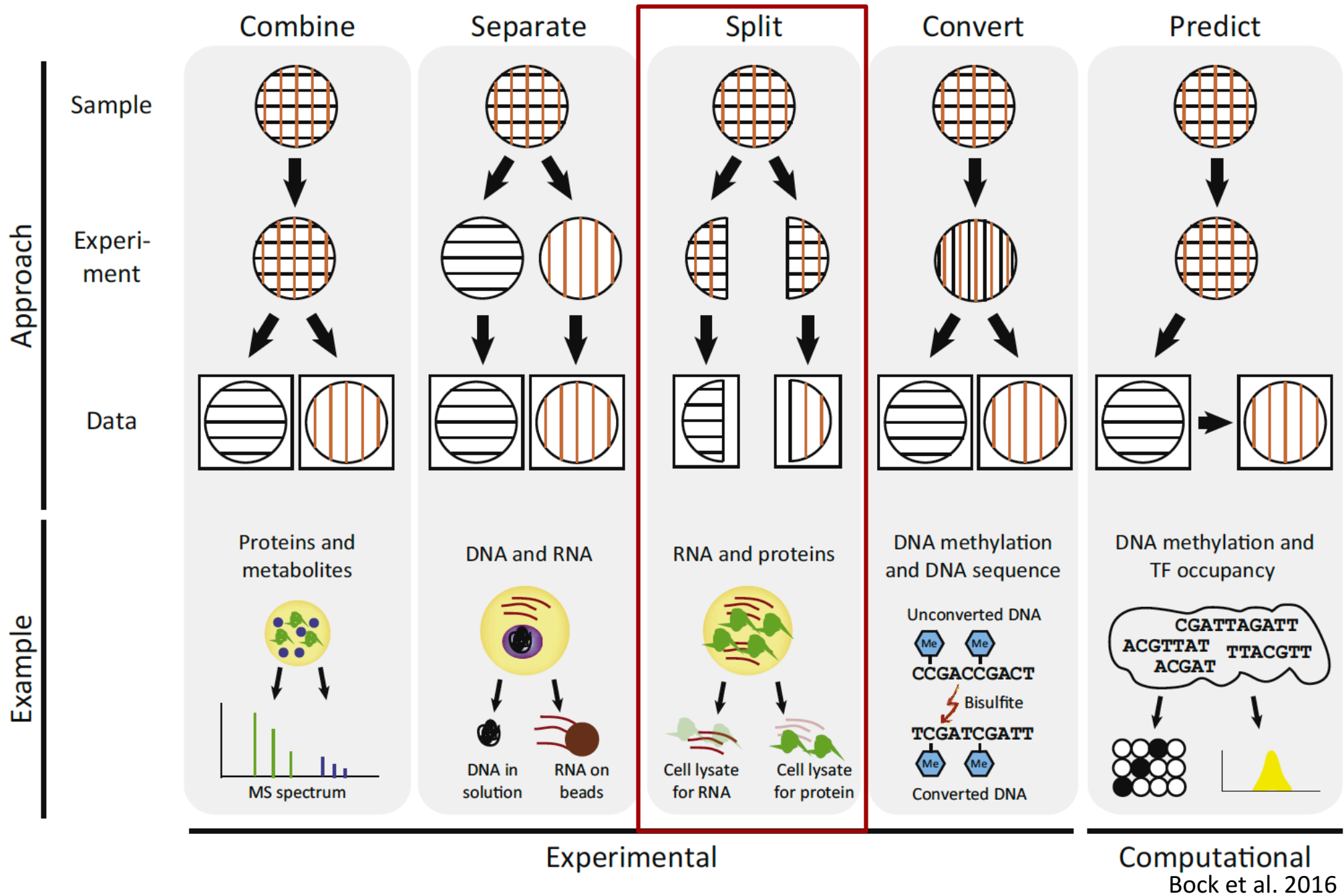
## Spatial ‘Omics’

- Multiplex FISH (Seq-FISH, MERFISH)
- *In situ* RNA-seq (e.g. FISSEQ)

## Multi-omics

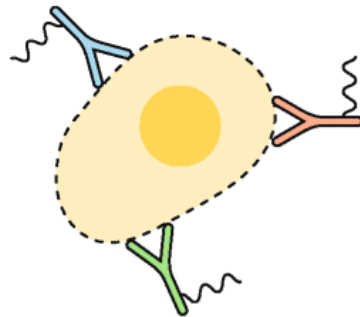
- DNA + RNA (G + T)
- RNA + protein (T + P)
- Epigenome + RNA

# Multi-omics strategies

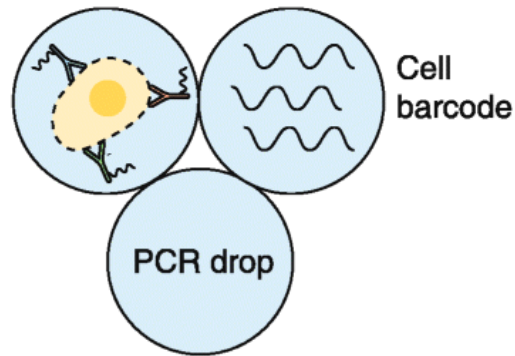


# RNA-Seq & cell-surface proteomics in a drop!

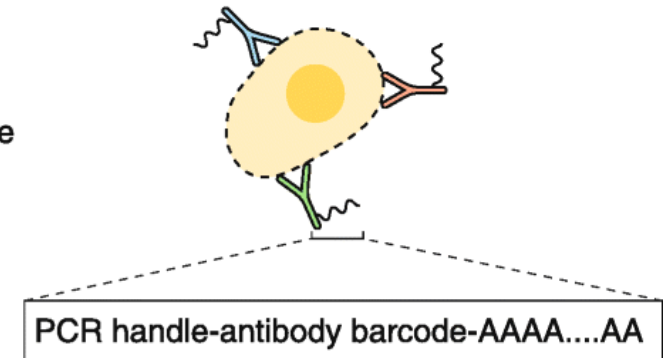
**ABCD** Ullal et al. 2014



**Abseq** Payam et al. 2017



**CITE-Seq** Stoeckius et al. 2017



Cell surface proteins	✓	✓	✓
Intracellular proteins	✓	✗	✗
Multiplexed	✗	✓	✓
RNA-seq	✗	✗	✓
Sequencing technology	Nanostring	Illumina	Illumina

**Analysis – I have generated some  
scRNAseq data ... what are the next steps**

# Starting with the gene expression matrix

## “Raw data”

```

AAATTATGACGATGTGCTTG.....GACTGCAC
CGTTAGATGGCAGGGCCGGG.....CTCATAGT
AAATTATGACGAAGTTTGTA.....GCTCATAA
GTTAAACGTACCCTAGCTGT.....GATTTTCT
TTGCCGTGGAGTGTGGGGGT.....ATAAGCTC
TTGCCGTGGGTGTATGGAGG.....CCAGCACC
GTTAAACGTACCGCAGGTTT.....GTTGGCGT
AAATTATGACGAAGTTTGTA.....AGATGGGG
CGTTAGATGGCATCTAGGCT.....GGGGACGA
GTTAAACGTACCAAGGCTTG.....CAAAGTTC
TTGCCGTGGAGTCGTGAGGG.....TTCCAAGG
CGTTAGATGGCACCTGTGTA.....TGGTACGT
GTTAAACGTACCATCCGGTG.....TTAAACCG

```

.....  
 .....  
 .....  
 (Hundreds of millions of reads)

## “Processed” data

	Cell: 1	2	...	N
GENE 1	1	2		14
GENE 2	4	27		8
GENE 3	0	0		1
⋮	⋮	⋮		⋮
GENE M	6	2		0

- Sequences derived from different scRNAseq assays are complex and vary
- Different pipelines are needed to address different sequence formats
- Common steps include:
  - Aligning
  - QC
  - Read counting

Credit: Karthik Shekhar

# Starting with the gene expression matrix

## “Raw data”

```
AAATTATGACGATGTGCTTG.....GACTGCAC
CGTTAGATGGCAGGCCGGG.....CTCATAGT
AAATTATGACGAGTTTGTA.....GCTCATAA
GTTAAACGTACCTAGCTGT.....GATTTTCT
TTGCCGTGGAGTGTGGGGGT.....ATAAGCTC
TTGCCGTGGTGTATGGAGG.....CCAGCACC
GTTAAACGTACGCAGGTTT.....GTTGGCGT
AAATTATGACGAGTTTGTA.....AGATGGGG
CGTTAGATGGCATCTAGGCT.....GGGGACGA
GTTAAACGTACCAAGGCTTG.....CAAAGTTC
TTGCCGTGGAGTCGTGAGGG.....TTCCAAGG
CGTTAGATGGCACCTGTGTA.....TGGTACGT
GTTAAACGTACCATCCGGTG.....TTAAACCG
.....
.....
.....
```

(Hundreds of millions of reads)

## Qualifications

- Full length vs. 3' vs 5'
- Poly A vs. Random priming
- Strand-specific vs non-specific
- UMI vs. non-UMI

## “Processed” data

	Cell: 1	2	...	N
<i>GENE 1</i>	1	2		14
<i>GENE 2</i>	4	27		8
<i>GENE 3</i>	0	0		1
⋮	⋮	⋮		⋮
⋮	⋮	⋮		⋮
<i>GENE M</i>	6	2		0



# Starting with the gene expression matrix

## “Raw data”

```
AAATTATGACGATGTGCTTG.....GACTGCAC
CGTTAGATGGCAGGCCGGG.....CTCATAGT
AAATTATGACGAGTTTGTA.....GCTCATAA
GTTAAACGTACCCTAGCTGT.....GATTTTCT
TTGCCGTGGAGTGTGGGGGT.....ATAAGCTC
TTGCCGTGGTGTATGGAGG.....CCAGCACC
GTTAAACGTACGCAGGTTT.....GTTGGCGT
AAATTATGACGAGTTTGTA.....AGATGGGG
CGTTAGATGGCATCTAGGCT.....GGGGACGA
GTTAAACGTACCAAGGCTTG.....CAAAGTTC
TTGCCGTGGAGTCGTGAGGG.....TTCCAAGG
CGTTAGATGGCACCTGTGTA.....TGGTACGT
GTTAAACGTACCATCCGGTG.....TTAAACCG
.....
.....
.....
```

(Hundreds of millions of reads)

## Qualifications

- Full length vs. 3' vs 5'
- Poly A vs. Random priming
- Strand-specific vs non-specific
- UMI vs. non-UMI

- Quality filtering
- Cell barcode stratification
- Alignment
- Multimapping reads/intronic reads
- Quantification / UMI collapse

## “Processed” data

	Cell: 1	2	...	N
<i>GENE 1</i>	1	2		14
<i>GENE 2</i>	4	27		8
<i>GENE 3</i>	0	0		1
⋮	⋮	⋮		⋮
⋮	⋮	⋮		⋮
<i>GENE M</i>	6	2		0

# Starting with the gene expression matrix

## “Raw data”

```
AAATTATGACGATGTGCTTG.....GACTGCAC
CGTTAGATGGCAGGCCGGG.....CTCATAGT
AAATTATGACGAGTTTGTA.....GCTCATAA
GTTAAACGTACCCTAGCTGT.....GATTTTCT
TTGCCGTGGAGTGTGGGGGT.....ATAAGCTC
TTGCCGTGGTGTATGGAGG.....CCAGCACC
GTTAAACGTACGCAGGTTT.....GTTGGCGT
AAATTATGACGAGTTTGTA.....AGATGGGG
CGTTAGATGGCATCTAGGCT.....GGGGACGA
GTTAAACGTACCAAGGCTTG.....CAAAGTTC
TTGCCGTGGAGTCGTGAGGG.....TTCCAAGG
CGTTAGATGGCACCTGTGTA.....TGGTACGT
GTTAAACGTACCATCCGGTG.....TTAAACCG
.....
.....
.....
```

(Hundreds of millions of reads)

## Qualifications

- Full length vs. 3' vs 5'
- Poly A vs. Random priming
- Strand-specific vs non-specific
- UMI vs. non-UMI

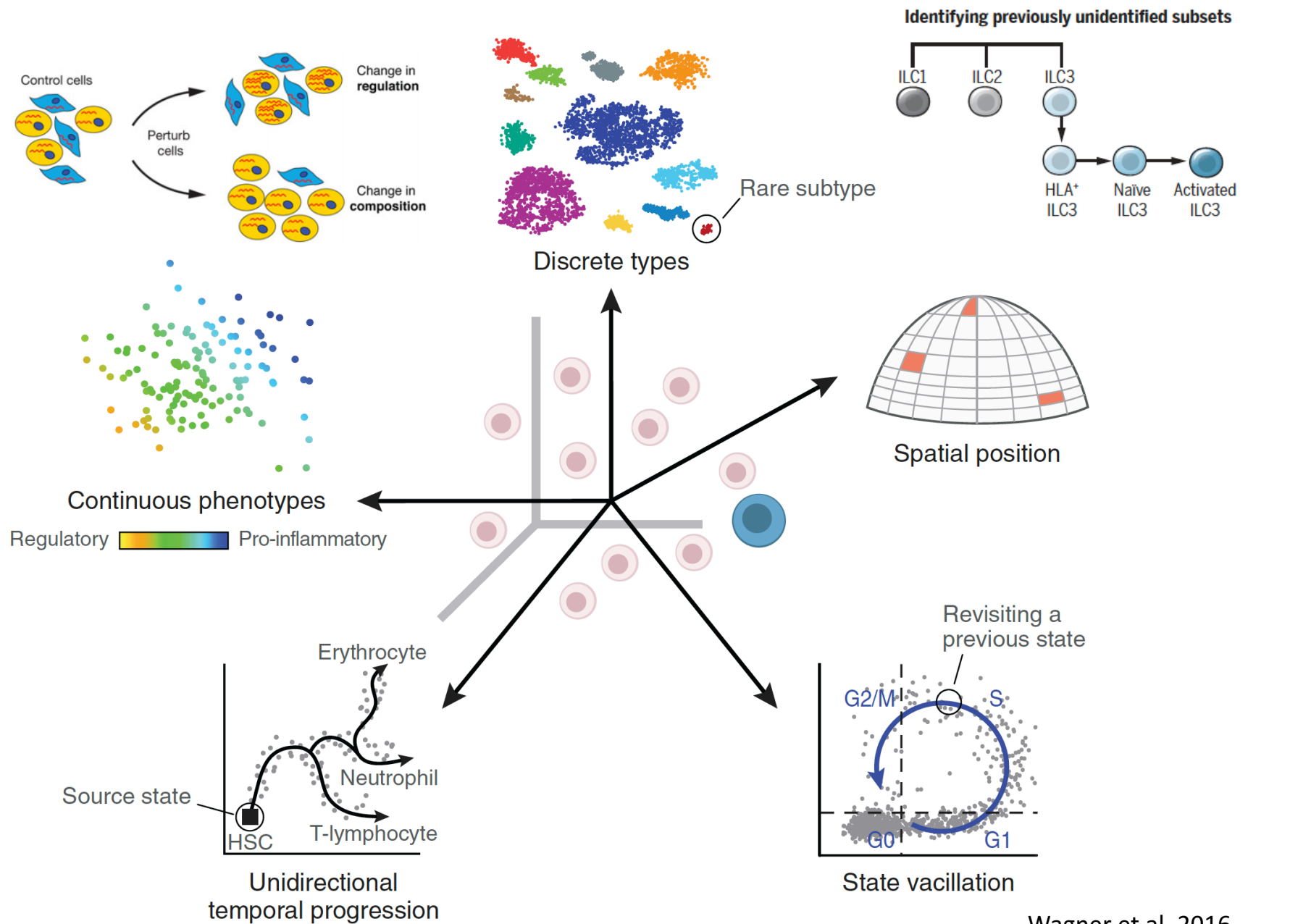
- Quality filtering
- Cell barcode stratification
- Alignment
- Multimapping reads/intronic reads
- Quantification / UMI collapse

## “Processed” data

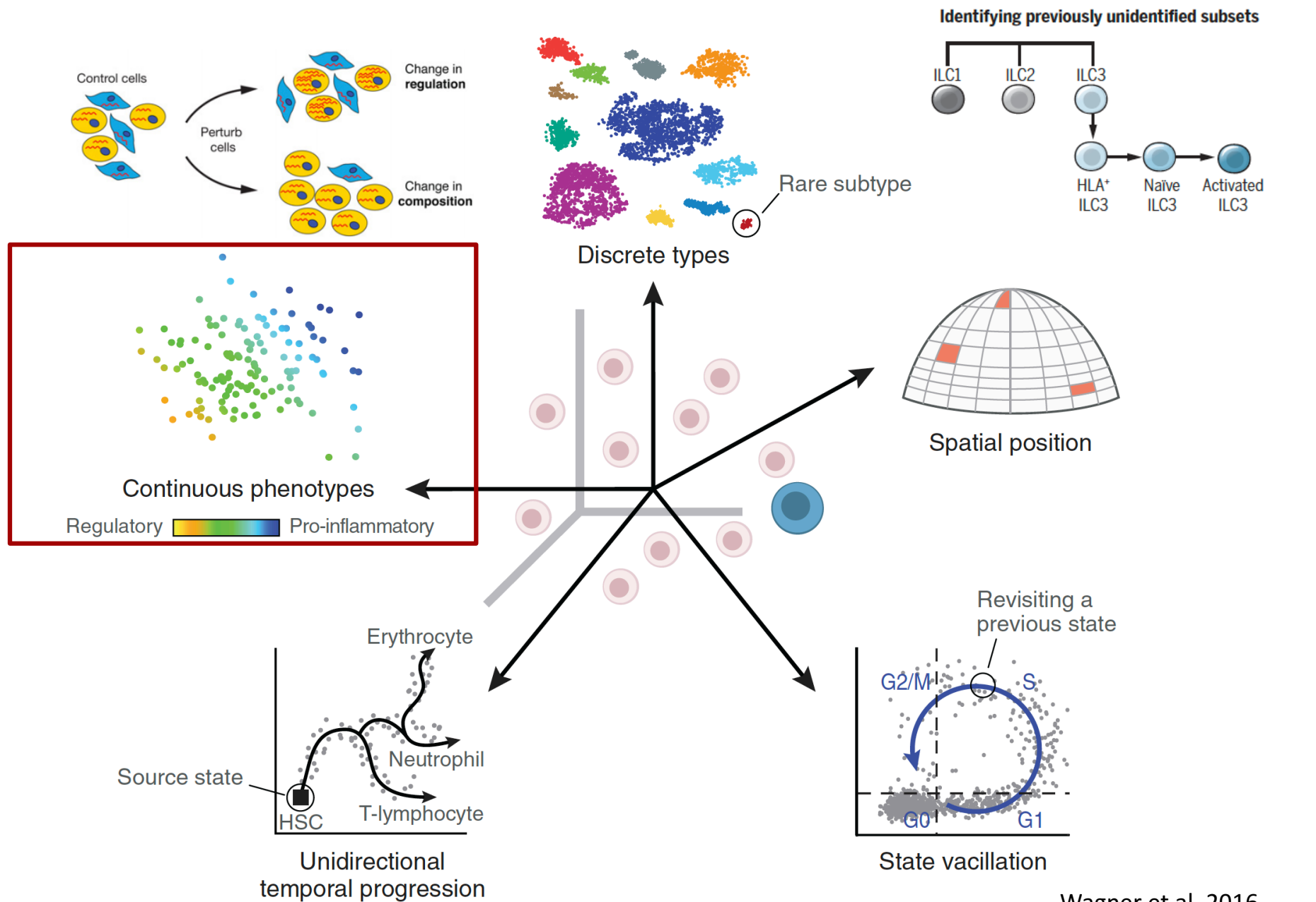
	Cell: 1	2	...	N
<i>GENE 1</i>	1	2		14
<i>GENE 2</i>	4	27		8
<i>GENE 3</i>	0	0		1
⋮	⋮	⋮		⋮
⋮	⋮	⋮		⋮
<i>GENE M</i>	6	2		0

Once I have my gene expression matrix, what's next?

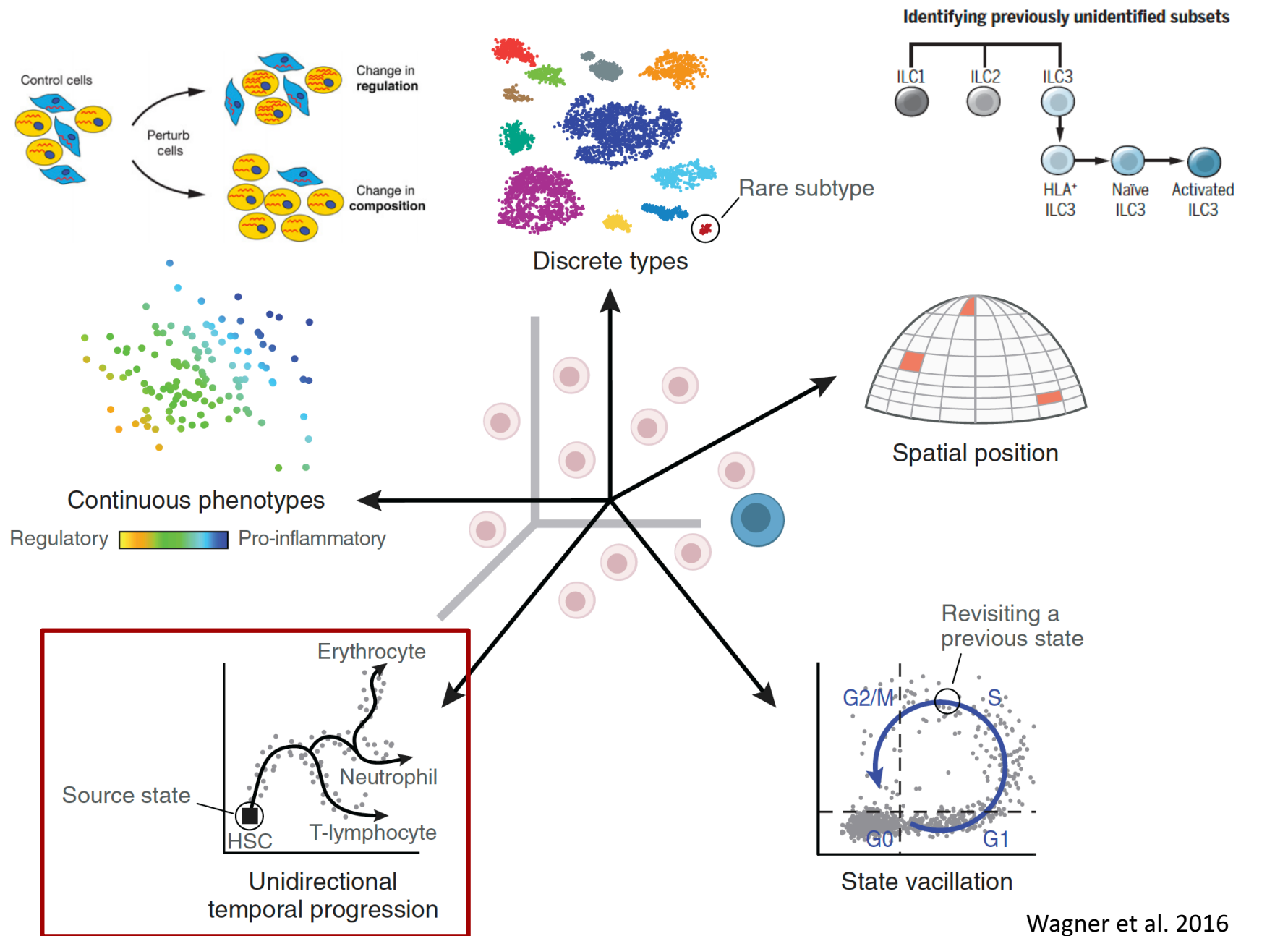
# Inference – from data to biology



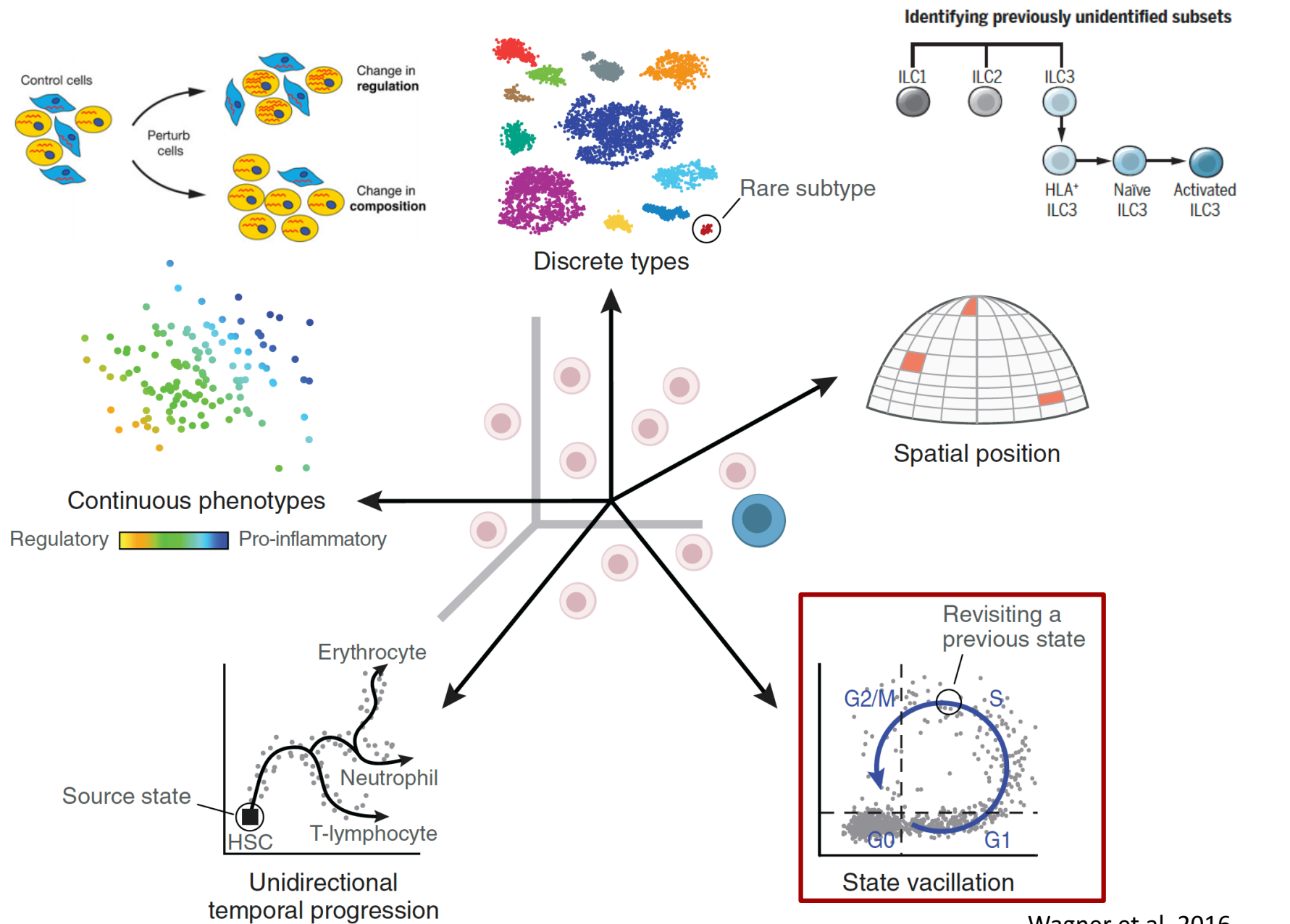
# Inference – from data to biology



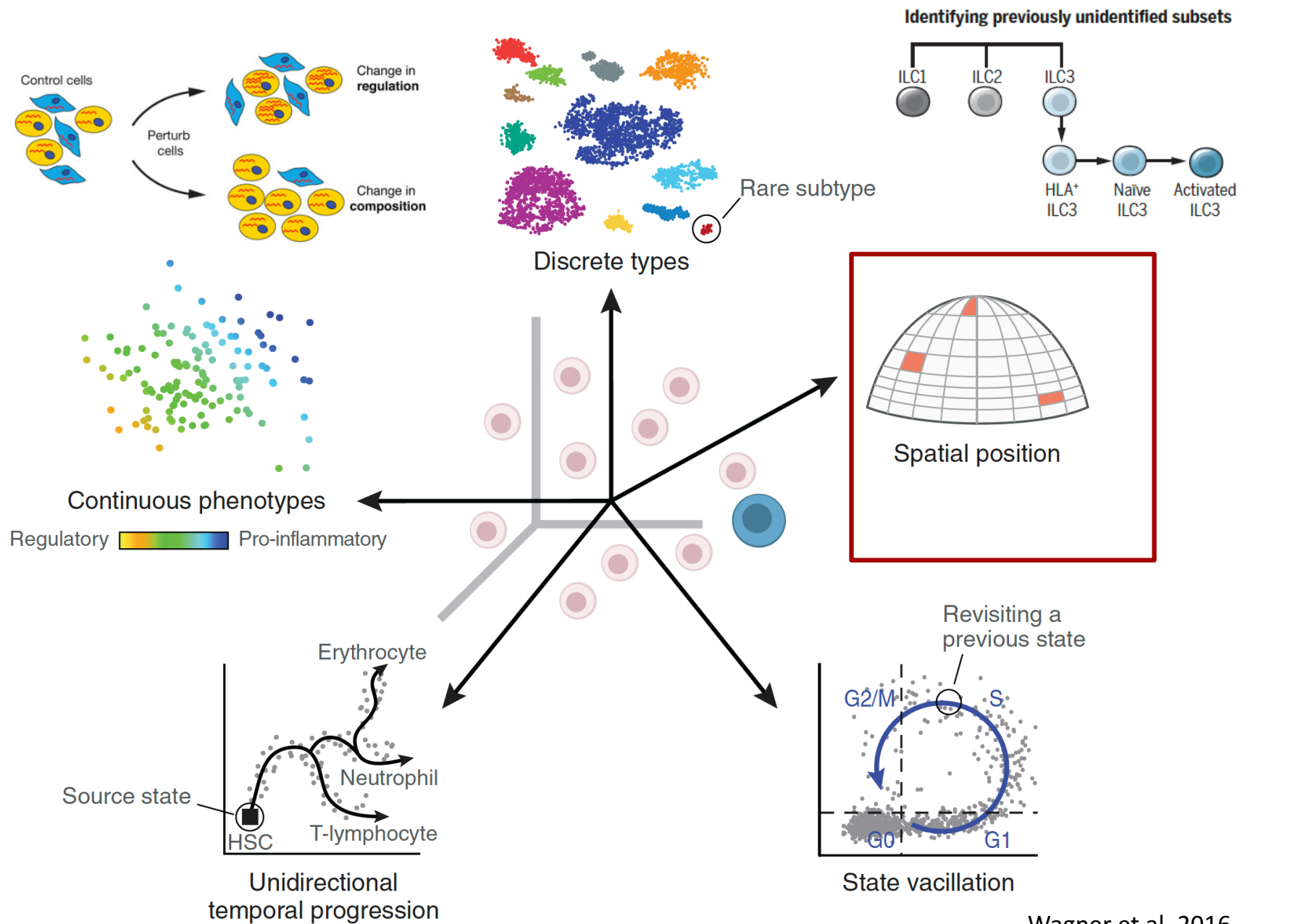
# Inference – from data to biology



# Inference – from data to biology

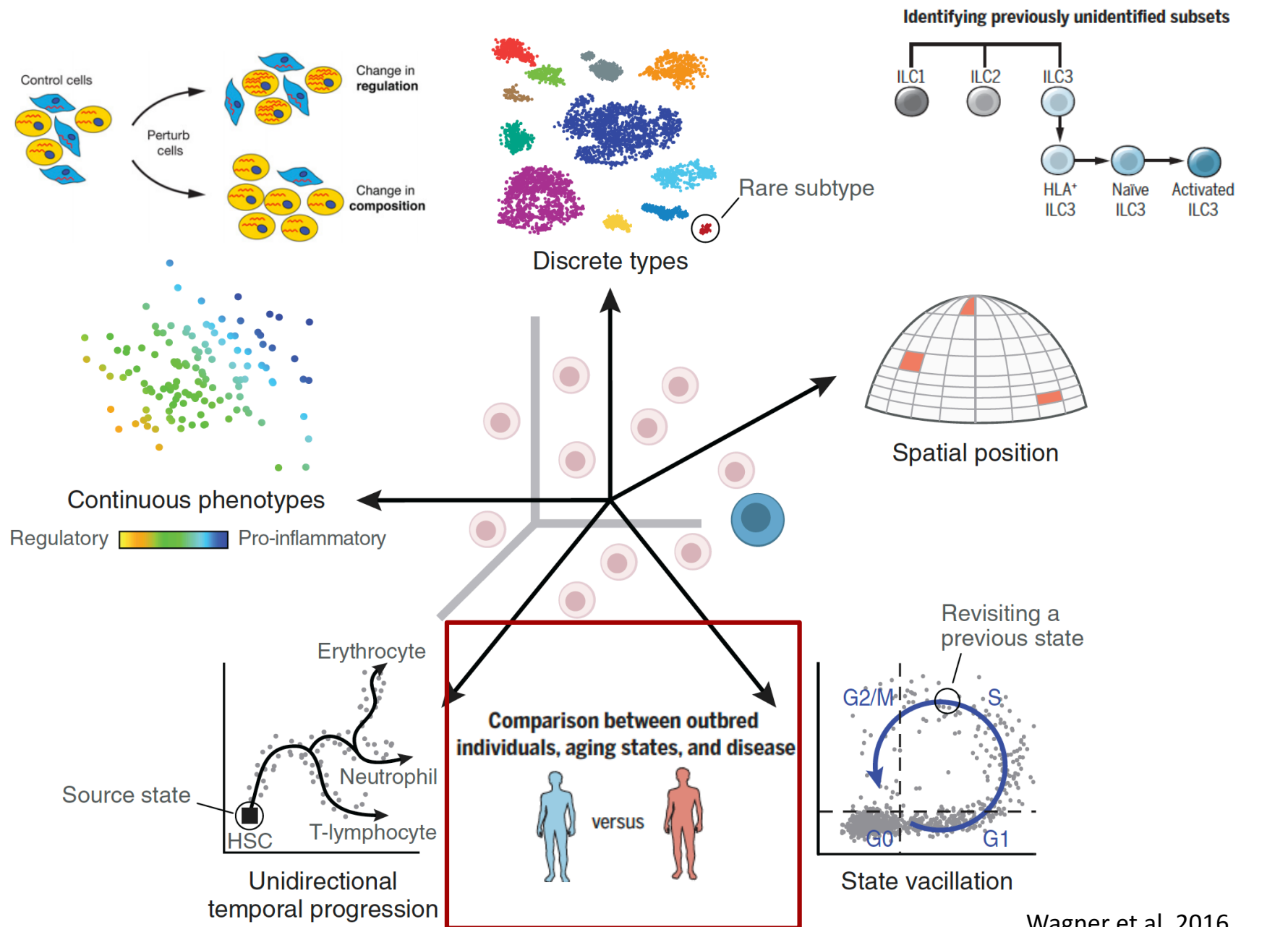


# Inference – from data to biology



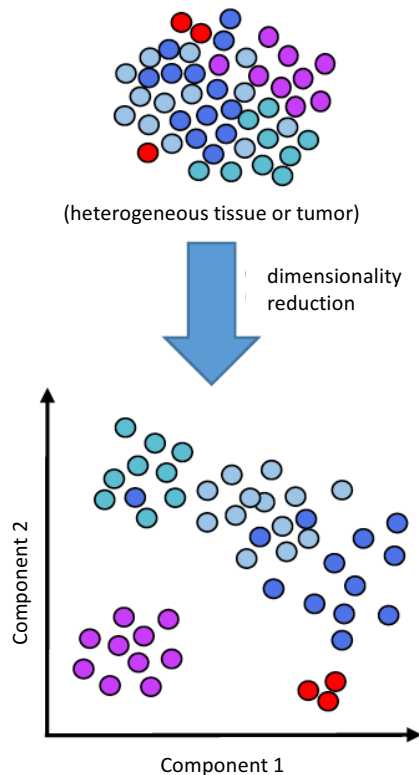


# Inference – from data to biology



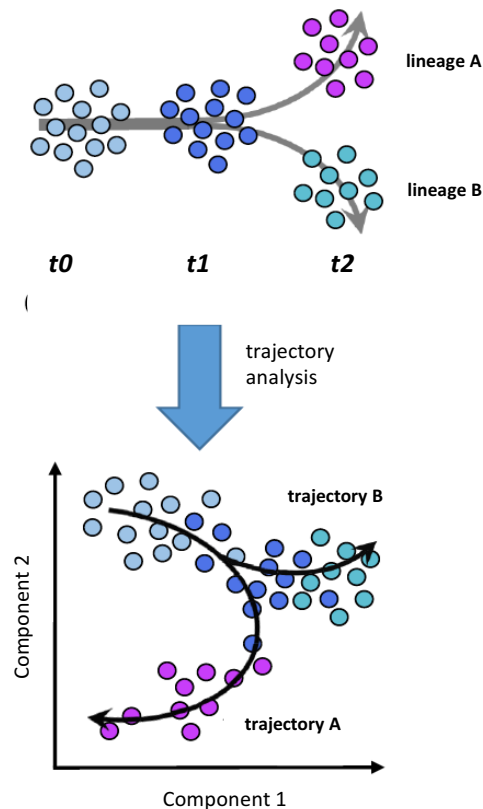
# The type of biological questions will dictate analyses to be undertaken

## (1) Deconvolution of heterogeneous population



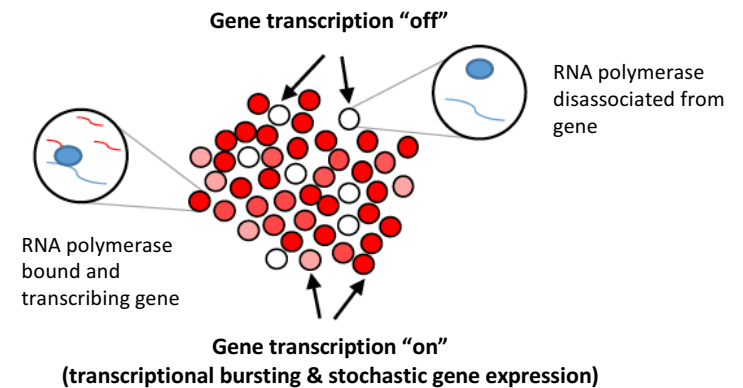
- Cell atlas
- Diseased vs. healthy
- Pre- vs. post-therapy

## (2) Trajectory analysis

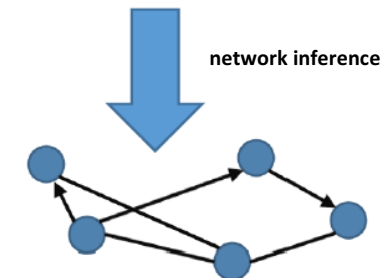
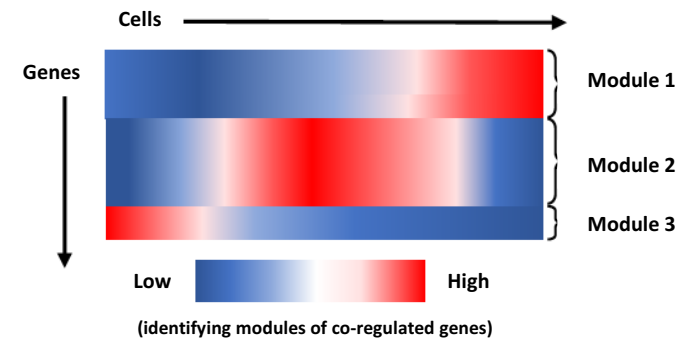


- Cell state transition:
  - cell differentiation
  - response to stimulus
- Development

## (3) Dissecting transcription mechanics

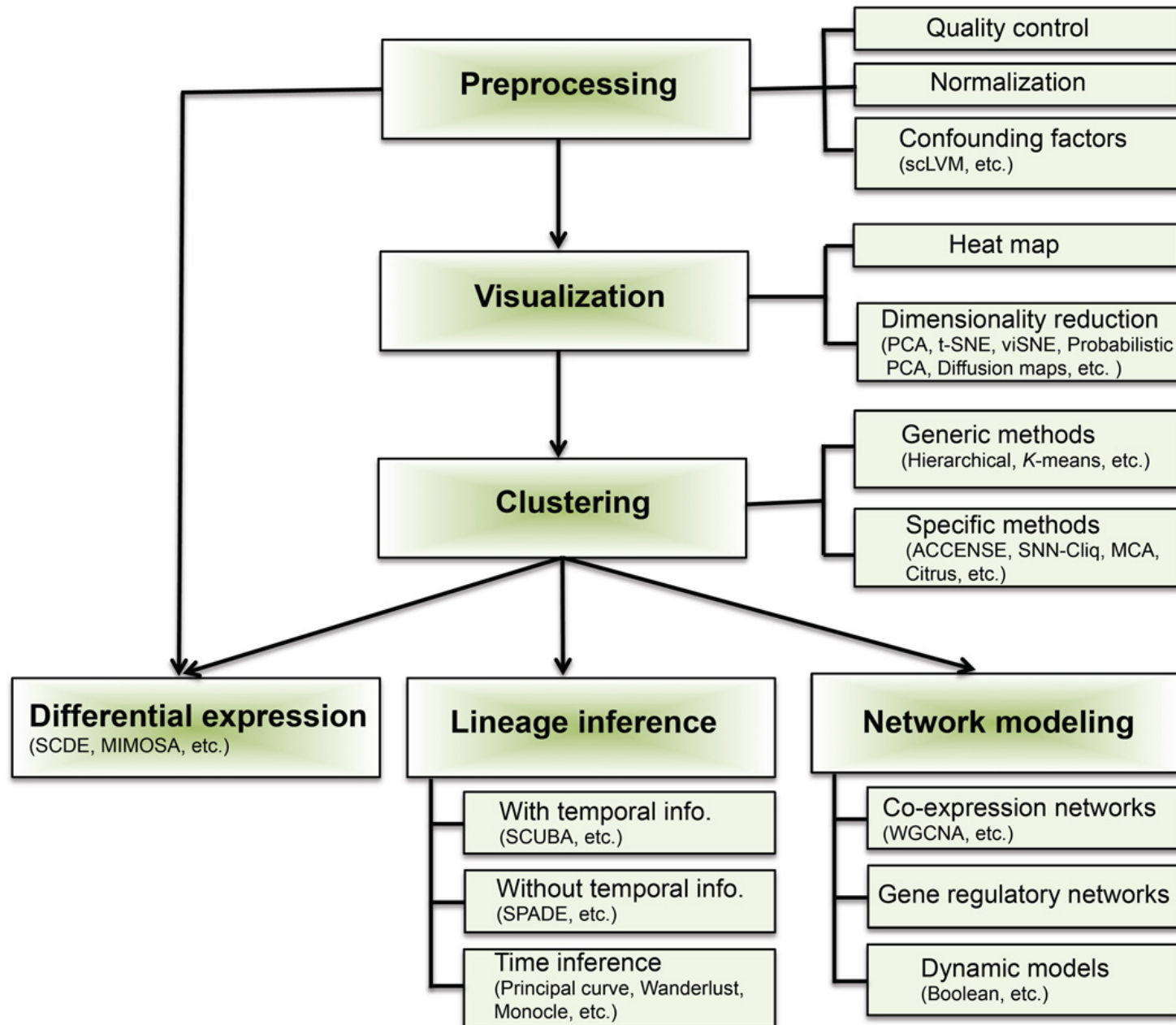


## (4) Network Inference



(inference of gene regulatory networks/subnetworks)

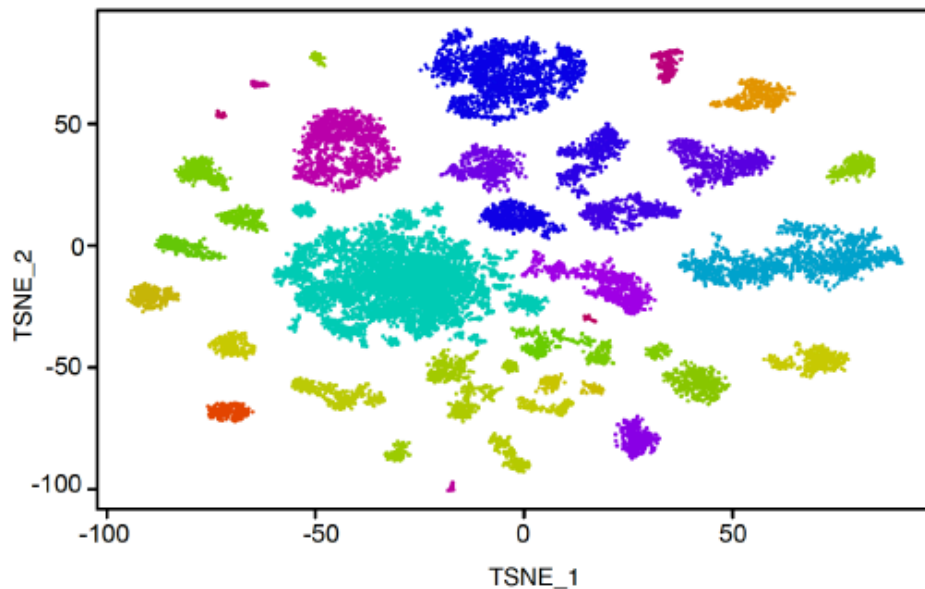
# Typical chart for scRNAseq analysis



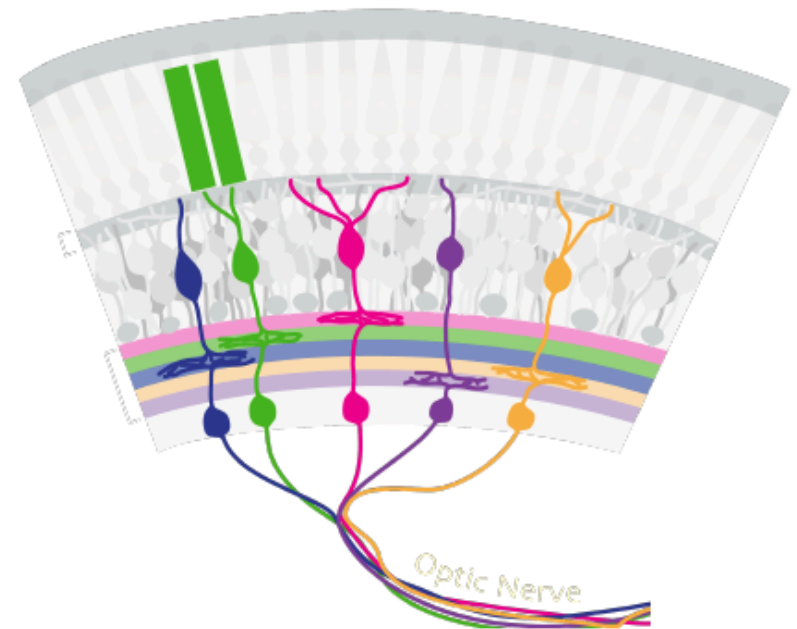
# Unsupervised clustering for cell type discovery

Full Retina Drop-seq (45000 cells, 39 types)

Macosko et al., *Cell*, 2015



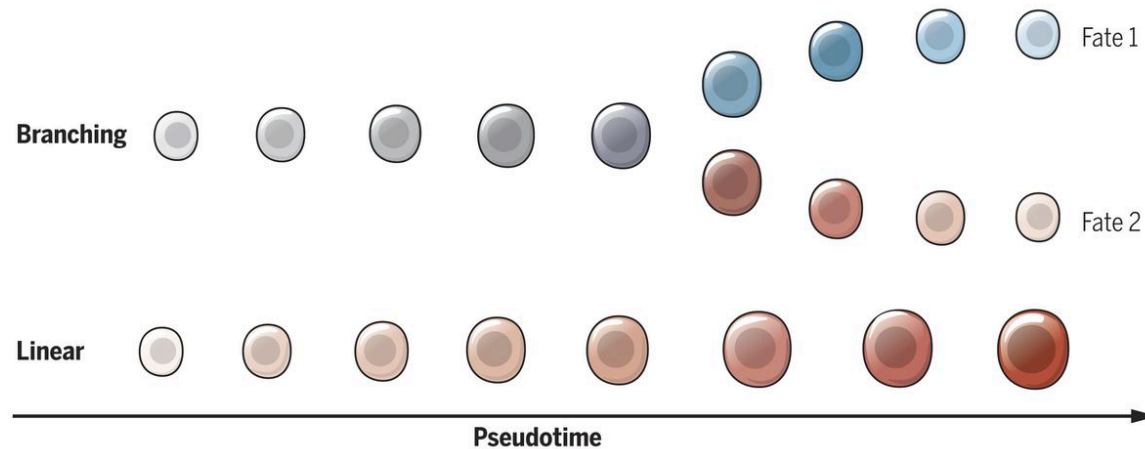
Proof of  
Principle



# Single cell trajectory analyses

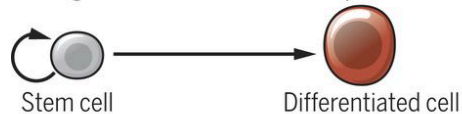
## A Development and differentiation of lymphocytes are studied with time series

Pseudotime measures the progress of cells through a differentiation process

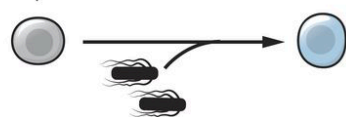


## B Examples of biological processes

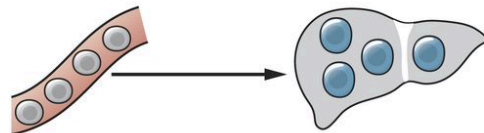
### 1 Progression of stem cell development



### 2 Response of naive immune cells to infection

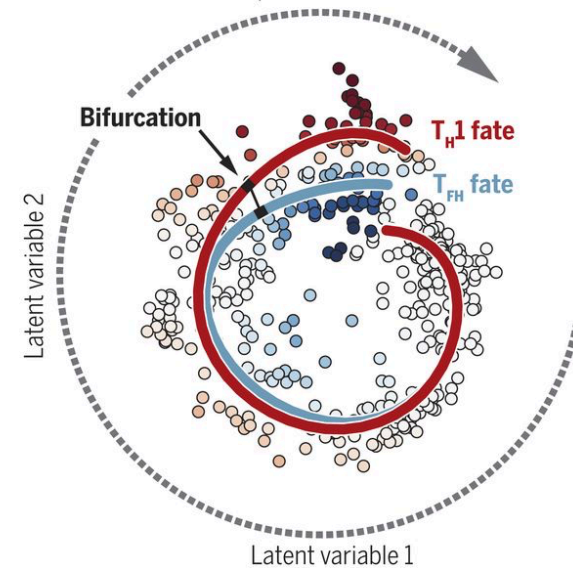


### 3 Adaptation of circulating immune cells to the tissues where they ultimately reside



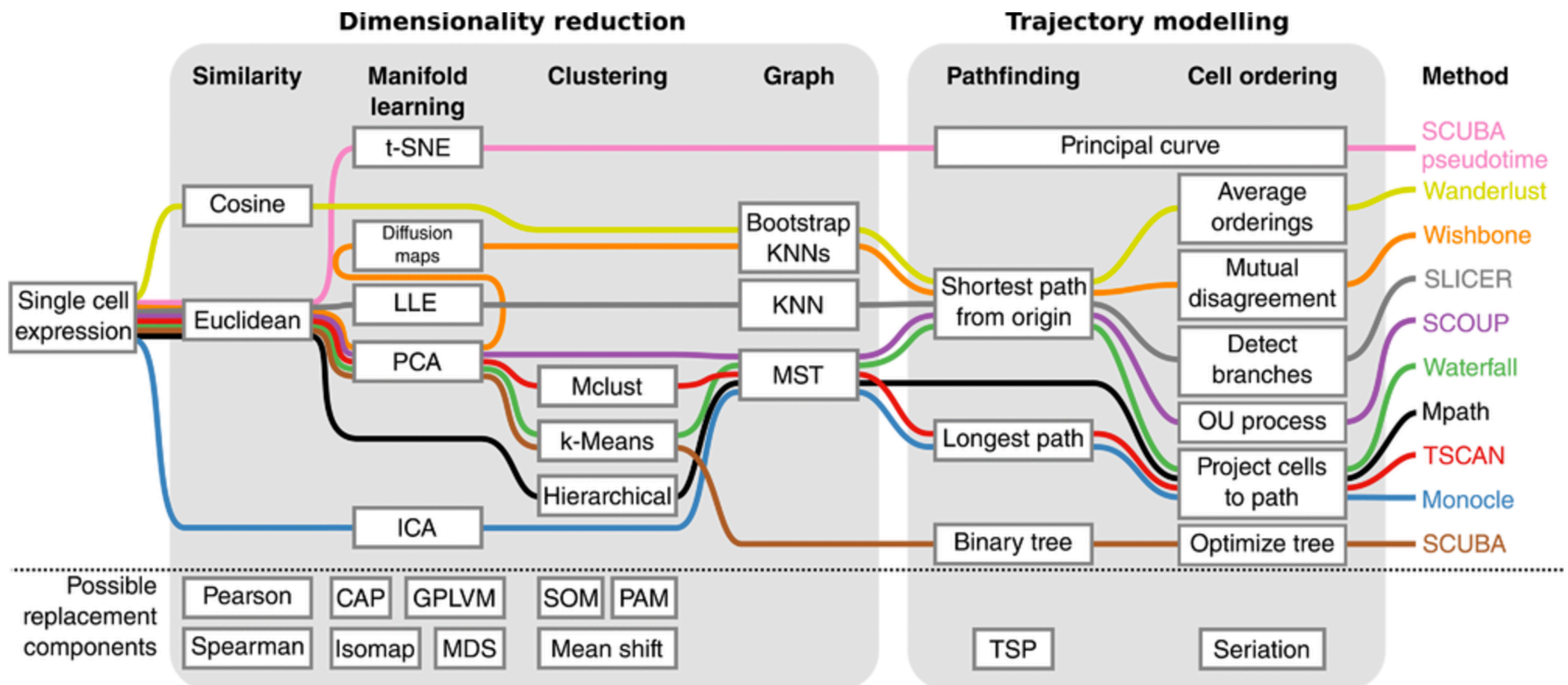
## C Bifurcating pseudotime trajectory

Inferred from scRNA-seq data






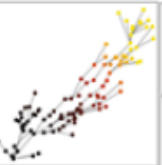
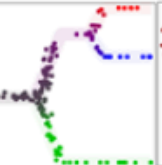
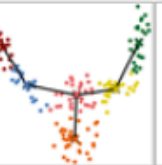
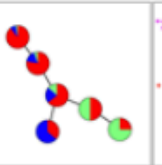
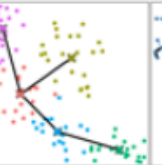
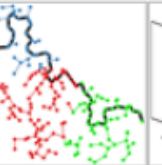
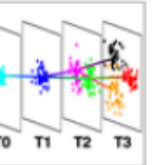
# Single cell trajectory analyses

Simplified representation of dataset





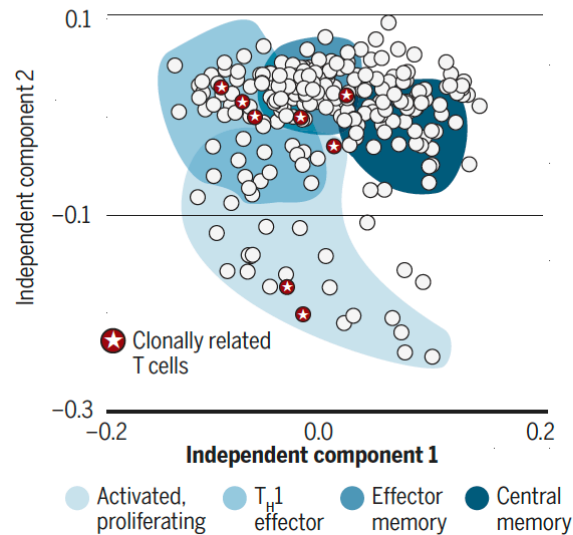
# Single cell trajectory analyses

Method	SCUBA pseudotime	Wanderlust	Wishbone	SLICER	SCOUP	Waterfall	Mpath	TSCAN	Monocle	SCUBA
Visual abstract										
Structure	Linear	Linear	Single bifurcation	Branching	Branching	Linear	Branching	Linear	Branching	Branching
Robustness strategy	Principal curves	Ensemble, starting cell	Ensemble, starting cell	Starting cell	Starting population	Clustering of cells	Clustering of cells using external labelling	Clustering of cells	Differential expression	Simple model
Extra input requirements	None	Starting cell	Starting cell	Starting cell	Starting population	None	Time points	None	Time points	Time points
Unbiased	+	±	±	±	±	+	-	+	-	-
Scalability w.r.t. cells	-	-	±	±	-	±	+	+	-	±
Scalability w.r.t. genes	+	+	+	+	-	+	±	±	±	+
Code and documentation	-	±	+	±	+	±	+	+	+	±
Parameter ease-of-use	+	+	+	+	-	±	-	+	+	+

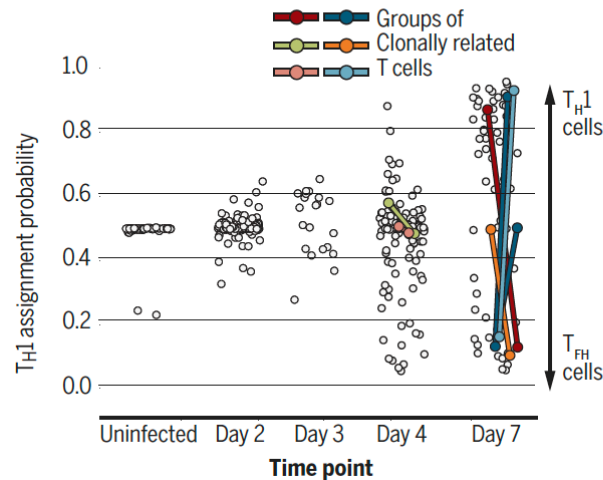


# Revealing T clone distributions between transcriptional state by analyzing TCR (requires full-length or custom primers)

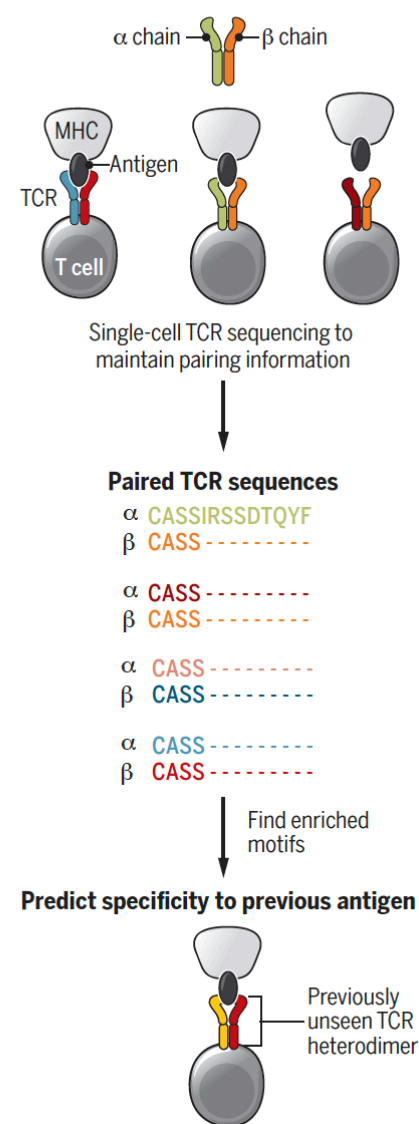
**A** TCR sequences assembled from scRNA-seq reads during *Salmonella* infection in mice



**B** TCR analysis during the immune response to malaria



**C** Prediction of binding specificity of TCR receptors



# Examples of additional analyses

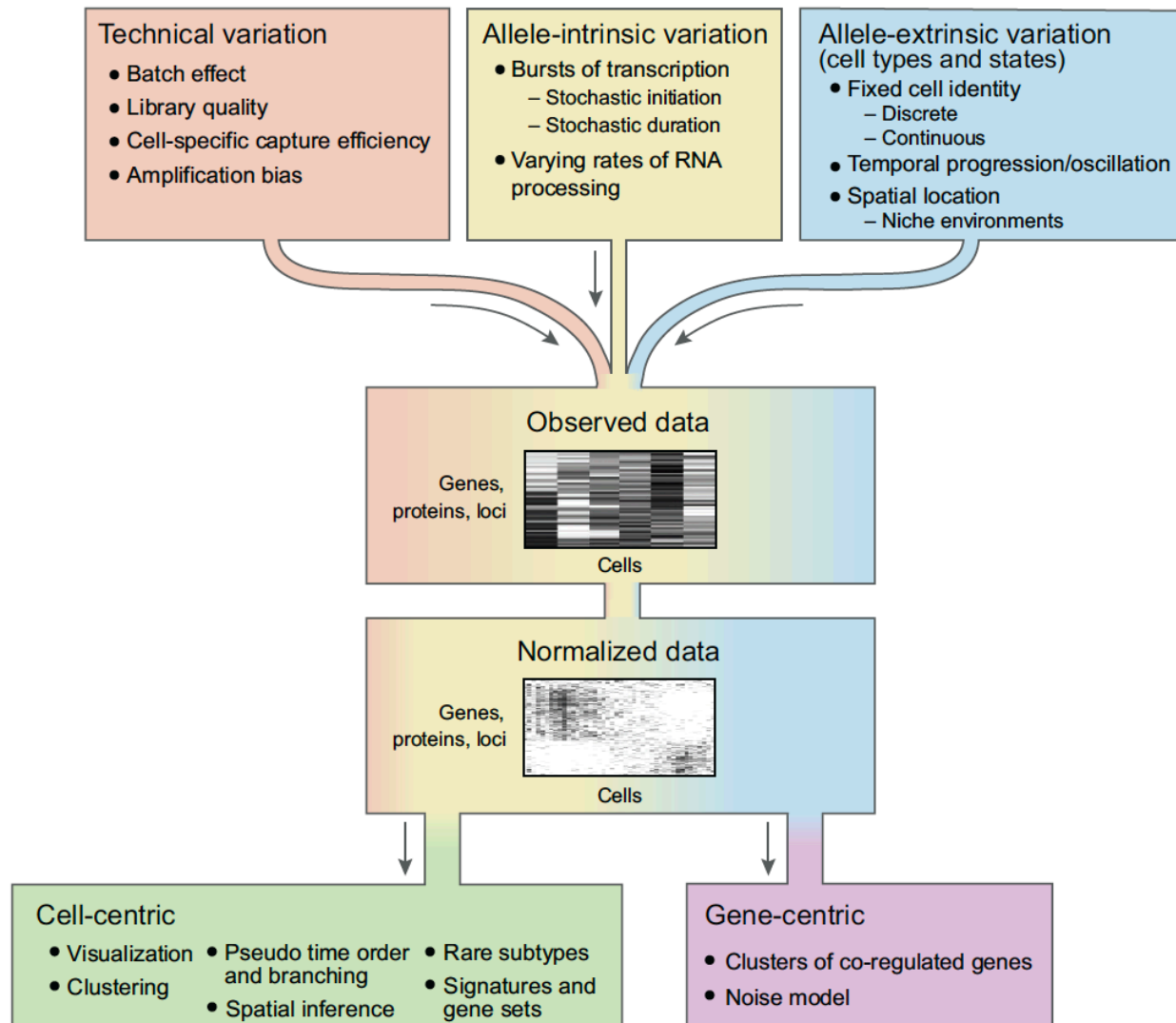
- Pathway and Geneset OverDispersion Analysis (PAGODA; Fan et al. Nat. Methods 2016)
- Alternative splicing
- Allelic expression
- Copy-number variation
- N.B. : alternative splicing and allelic expression require full length methods
  - Can draw conclusions with certainty only for highly expressed genes with good coverage
  - Take into consideration the drop-out rate → a unique splice form/allele in a single cell may be the results of detection issue

# List of references for methods & tutorial

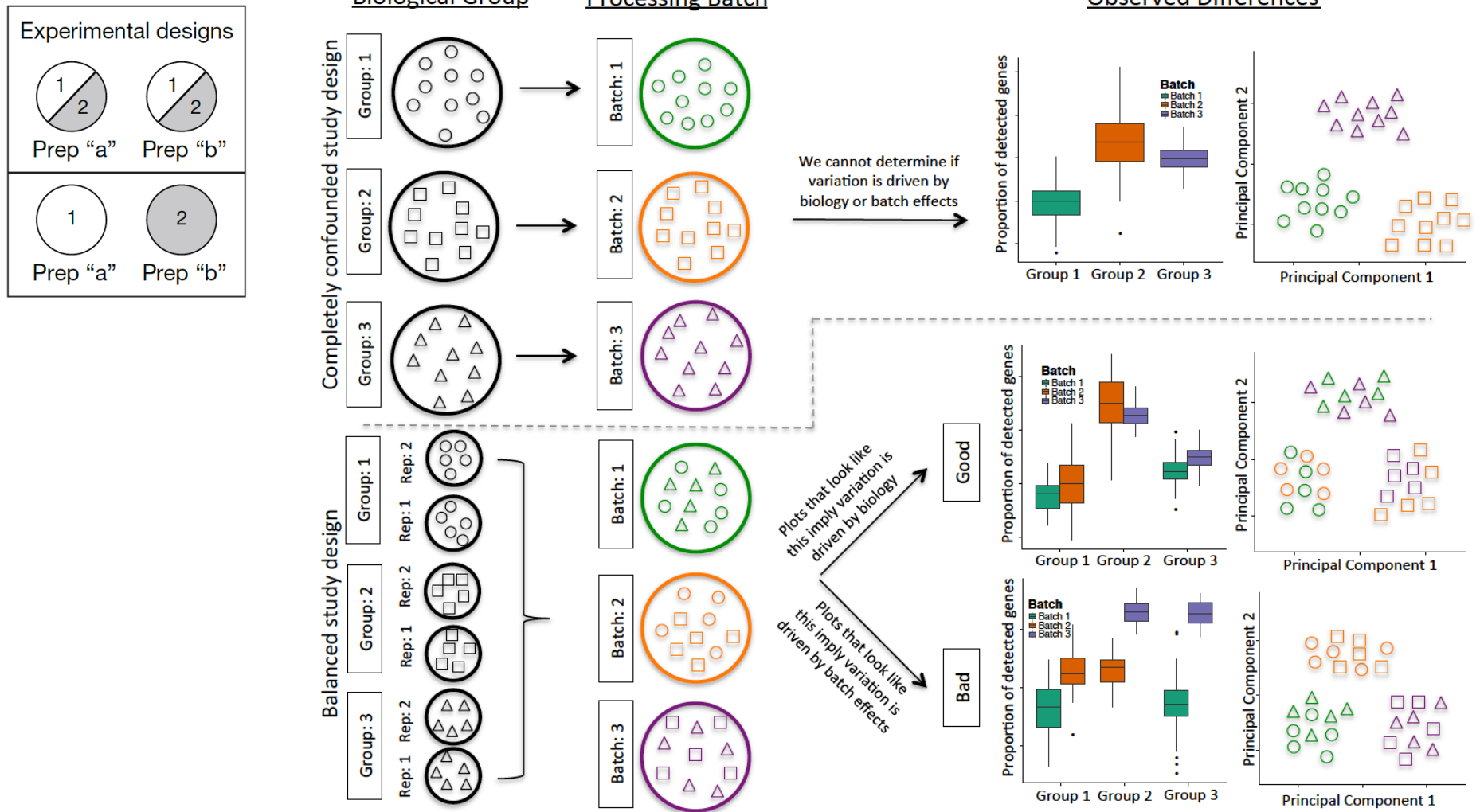
- Thank you to Sean Davis for the “Awesome single cell” compilation of software packages (and the people developing these methods) for analysis, including RNA-seq, ATAC-seq, etc.
  - <https://github.com/seandavi/awesome-single-cell>
- Examples of tutorials to get started:
  - [Seurat \(v2.0\) - Guided Clustering Tutorial:](http://satijalab.org/seurat/pbmc3k_tutorial.html)  
[http://satijalab.org/seurat/pbmc3k\\_tutorial.html](http://satijalab.org/seurat/pbmc3k_tutorial.html)
  - Sanger, [Hemberg Lab scRNA-seq course materials:](http://hemberg-lab.github.io/scRNA.seq.course)  
<http://hemberg-lab.github.io/scRNA.seq.course>
  - Harvard Single Cell Workshop (hosted by Peter Kharchenko ):  
<http://hms-dbmi.github.io/scw/>

# **Technical challenges in scRNAseq**

# Biological and technical factors driving gene expression readout

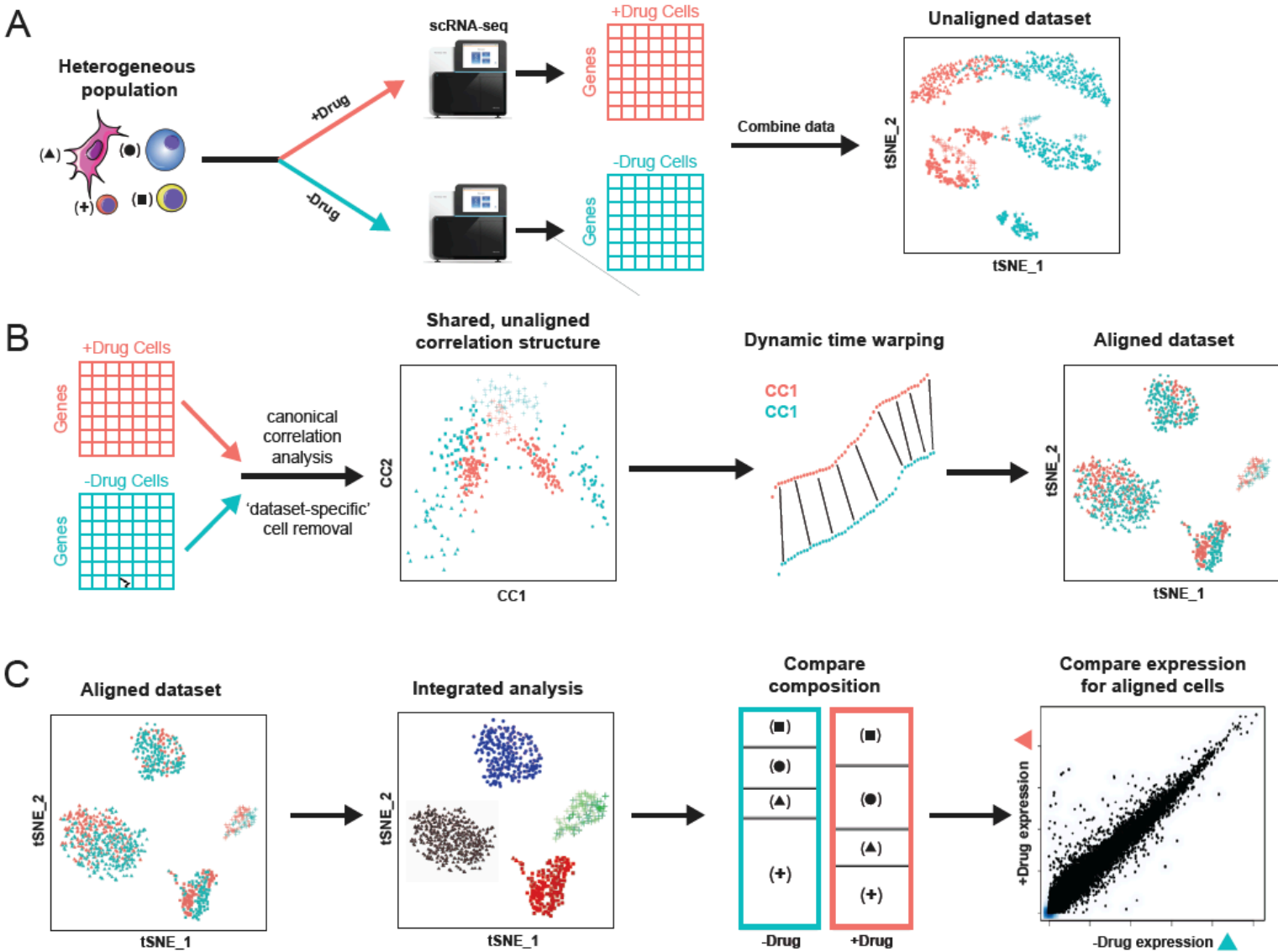


# Technical confounders in scRNAseq: Batch effect



# Integrated analysis of single cell transcriptomic data across conditions, technologies, and species

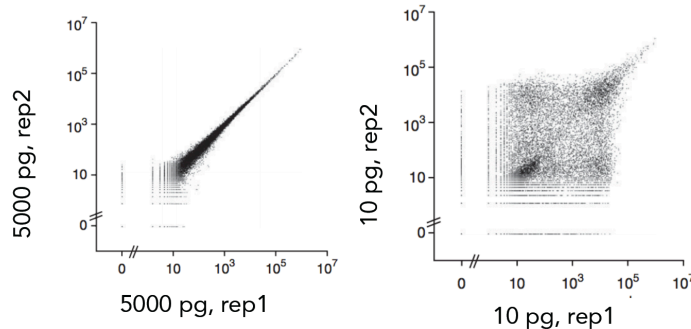
Andrew Butler<sup>1,2</sup> and Rahul Satija<sup>1,2,#</sup>



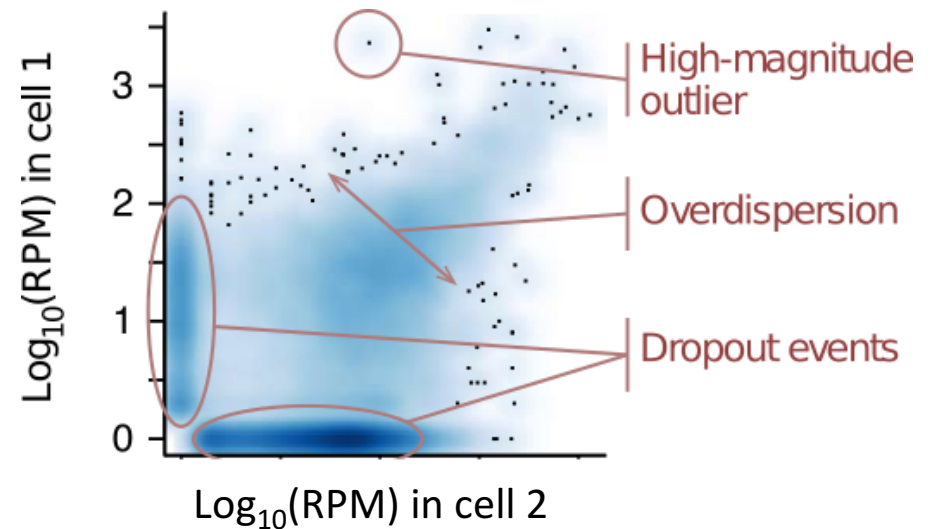


# Technical confounders in scRNAseq: Dropouts

(inefficient mRNA capture → sparse data / ~10% non-zero values)



- **Zero inflation**
  - Drop-out event during reverse-transcription
  - Genes with more expression have less zeros
  - Complexity varies
- **Transcription stochasticity**
  - Transcription bursting
  - Coordinated transcription of multigene networks
  - Over-dispersed counts
- **Higher Resolution**
  - More sources of signal



## BRIEF COMMUNICATIONS

### Bayesian approach to single-cell differential expression analysis



Peter V Kharchenko<sup>1-3</sup>, Lev Silberstein<sup>3-5</sup> &  
David T Scadden<sup>3-5</sup>

© 2014 Nature America, Inc.

# Technical confounders in scRNAseq: Dropouts

(inefficient mRNA capture → sparse data / ~10% non-zero values)

**Solution: imputing missing data!**

bioRxiv preprint first posted online Feb. 25, 2017; doi: <http://dx.doi.org/10.1101/111591>. The copyright holder for this preprint (which was not peer-reviewed) is the author/funder. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

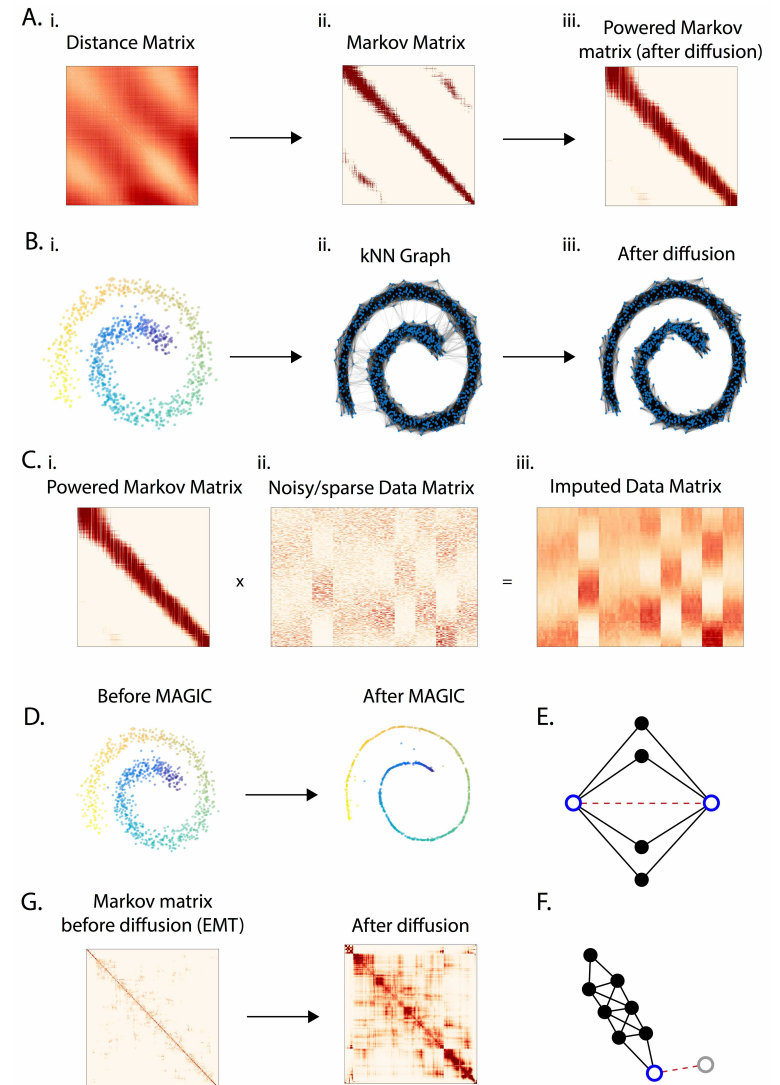
## MAGIC: A diffusion-based imputation method reveals gene-gene interactions in single-cell RNA-sequencing data

David van Dijk<sup>1</sup>, Jozas Nainys<sup>2,4</sup>, Roshan Sharma<sup>1,3</sup>, Pooja Kaithail<sup>1,4</sup>, Ambrose J. Carr<sup>1,4</sup>, Kevin R. Moon<sup>5,6</sup>, Linas Mazutis<sup>1,2</sup>, Guy Wolf<sup>5</sup>, Smita Krishnaswamy<sup>6\*</sup>, Dana Pe'er<sup>1\*</sup>

**MAGIC** = *Markov Affinity-based Graph*

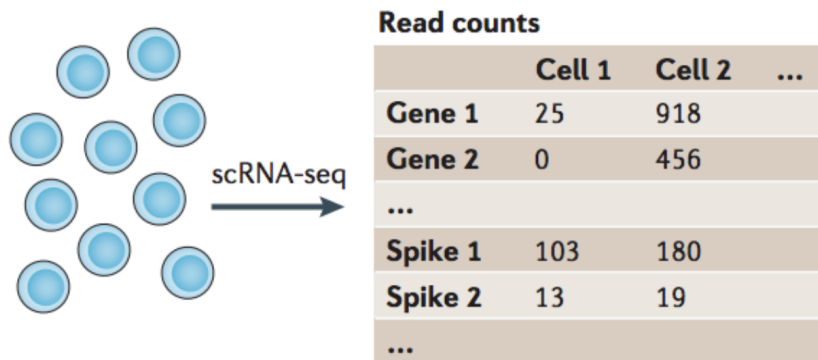
*Imputation of Cells*

→ *Method for imputing missing values & restoring structure in the data*

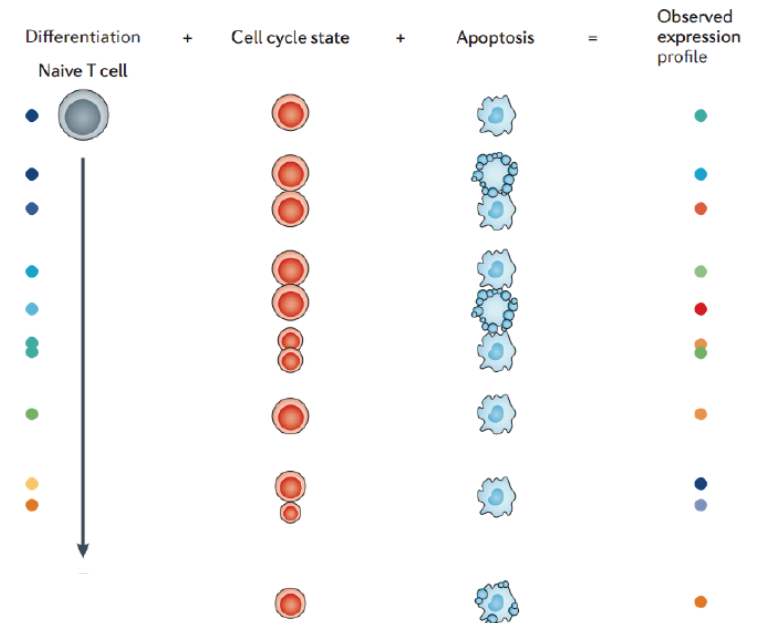


# Other technical confounders in scRNAseq

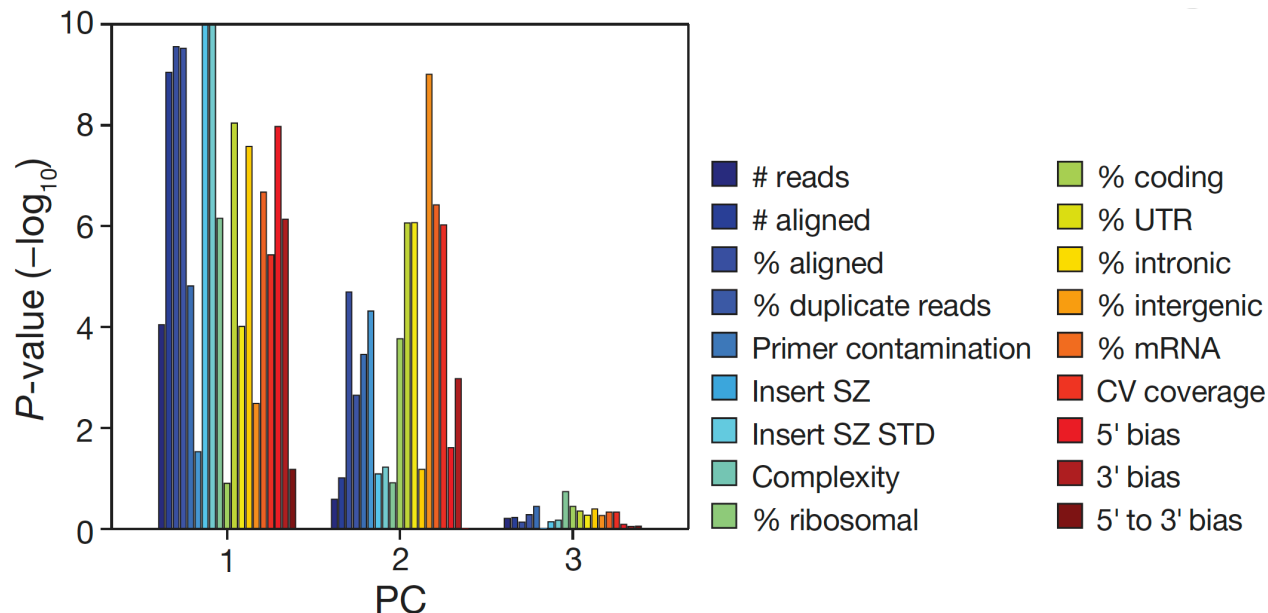
## 1- Variation in cell size and quality



## 2- Observed gene expression is a convolution



## 3- Variation dominated by “technical factors”

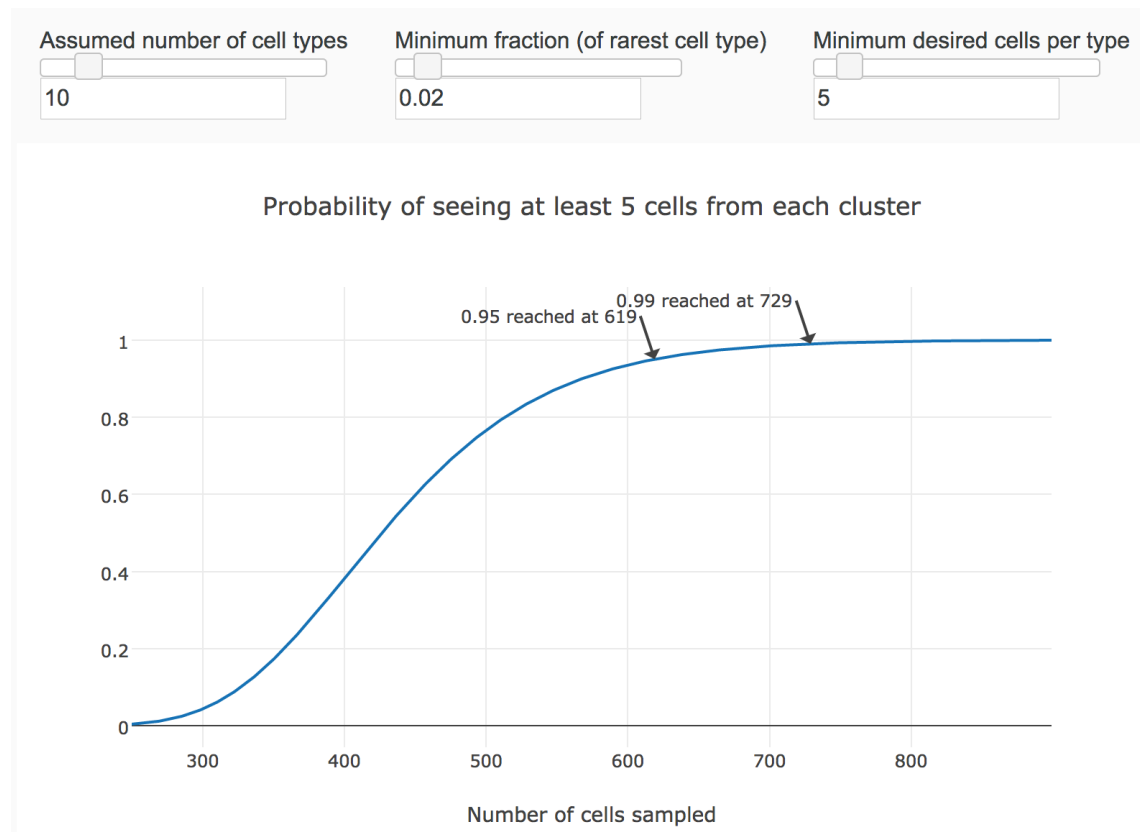


Buettner et al. 2015  
Wagner et al. 2017

# **Experimental design & common questions**

# How many cells should I be profiling?

- Can change depending on the variability of the biology and the expectation of finding rare populations.
- Satija lab online tool – [satijalab.org/howmanycells](https://satijalab.org/howmanycells)



# Cell number & Read depth

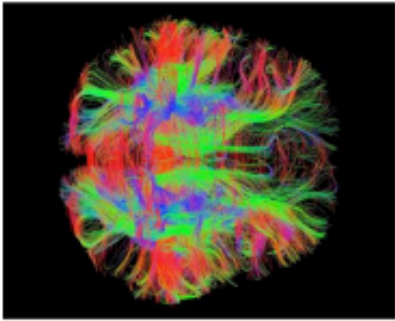
- For initial pilot study → aim for around 25-30 cells from each type
  - Sample with minor cell types < 5% will require sequencing at least 400 cells
  - Cell preselection/enrichment may be necessary, but unbiased cell selection is always preferred
- To study gene expression only, sequencing depth doesn't always have to be deep (depends on questions)
  - Multiplexing hundreds of samples on one sequencing lane is common
  - Cell clustering & cell-type identification benefits from large number of cells and doesn't always require high sequencing depth (~100,000 reads per cell)
  - Gene detection starts saturating from 1 million reads per cell
  - Transcription factor detection (regulatory networks) require higher read depth and most sensitive protocols

# Applications

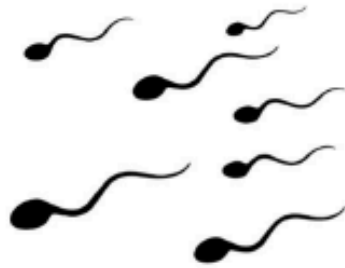


# Applications – Cancer biology

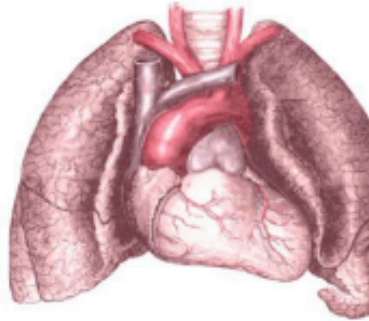
Neurobiology



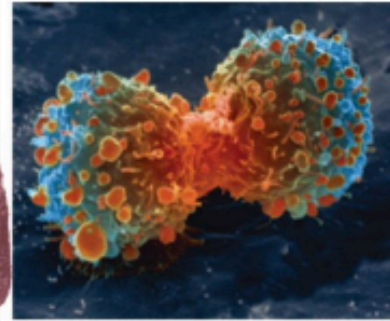
Germline Transmission



Organogenesis



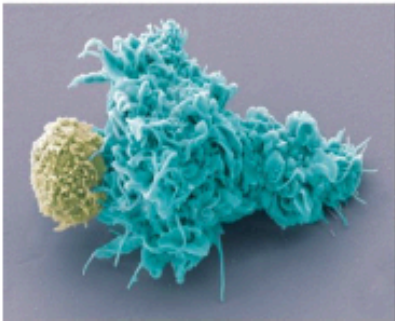
Cancer biology



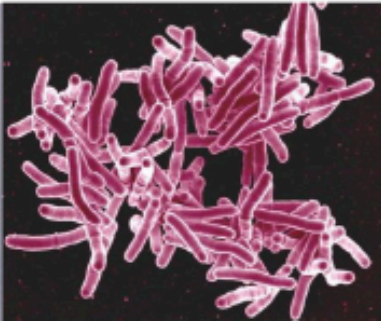
Clinical diagnostics



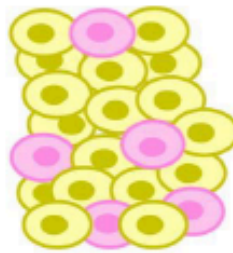
Immunology



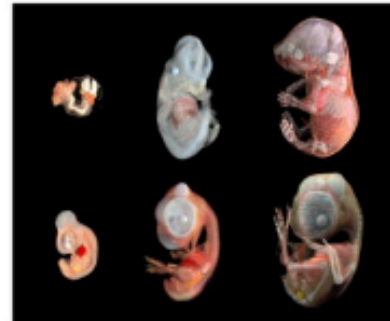
Microbiology



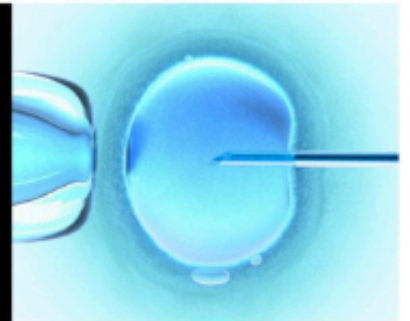
Tissue Mosaicism



Embryology

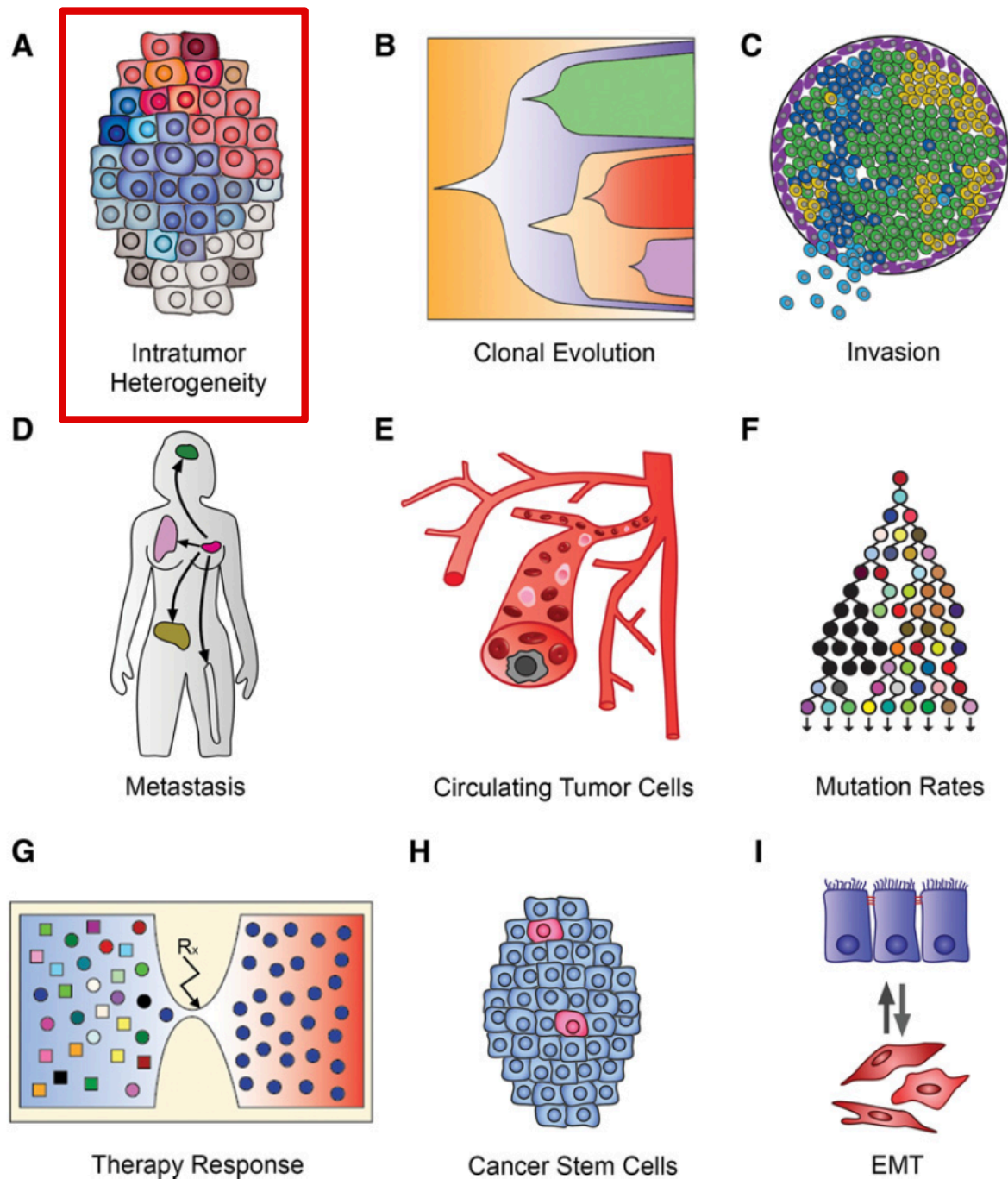


Prenatal-genetic diagnosis



Wang et al., 2015

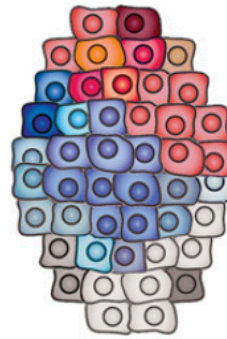
## A- Resolving intratumoral heterogeneity & dissecting microenvironment



**A- Resolving intratumoral heterogeneity  
& dissecting microenvironment**

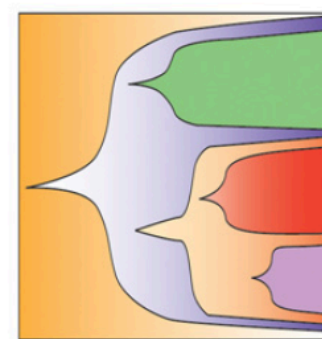
**B- Investigating clonal evolution in  
primary tumors**

**A**



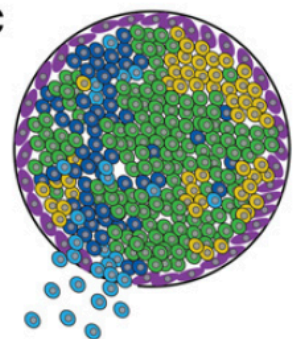
Intratumor  
Heterogeneity

**B**



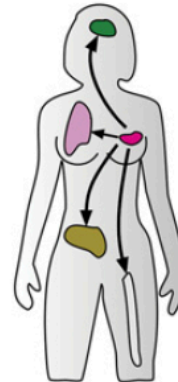
Clonal Evolution

**C**



Invasion

**D**



Metastasis

**E**



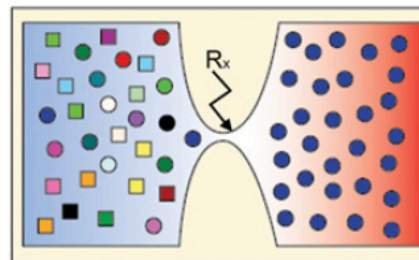
Circulating Tumor Cells

**F**



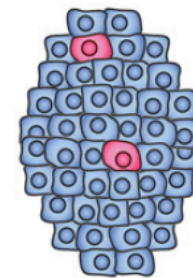
Mutation Rates

**G**



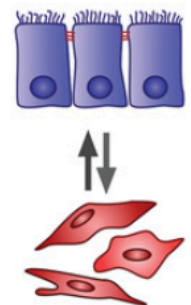
Therapy Response

**H**



Cancer Stem Cells

**I**



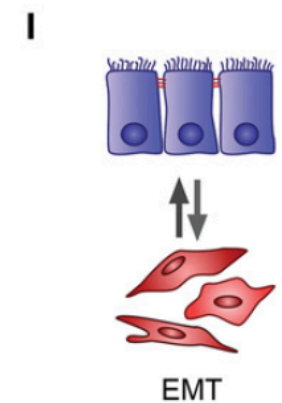
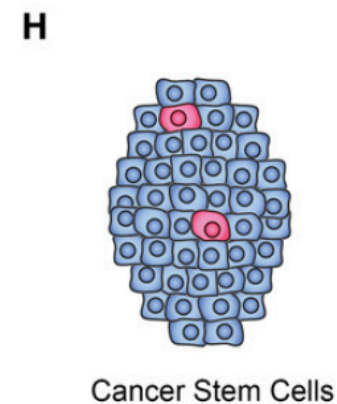
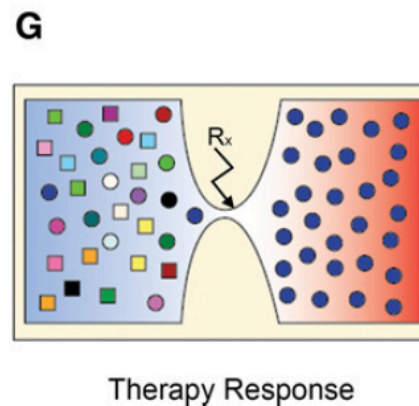
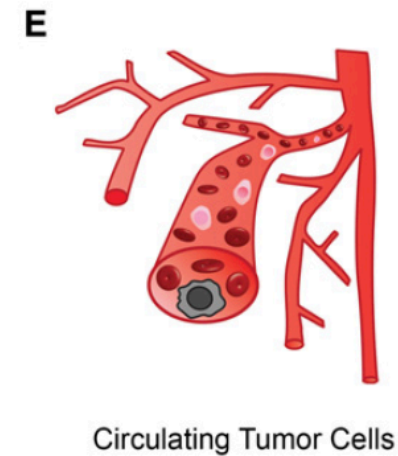
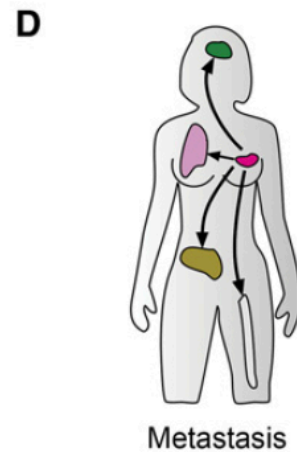
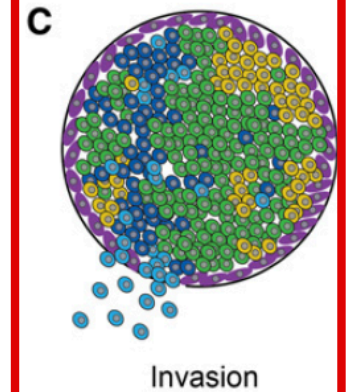
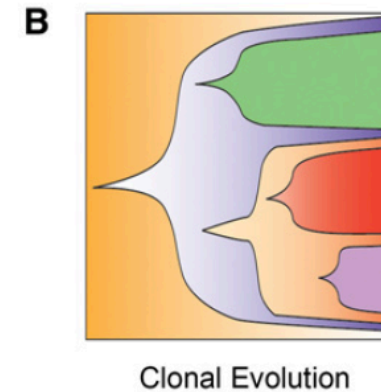
EMT



**A- Resolving intratumoral heterogeneity  
& dissecting microenvironment**

**B- Investigating clonal evolution in  
primary tumors**

**C- Studying invasion in early stage cancers**

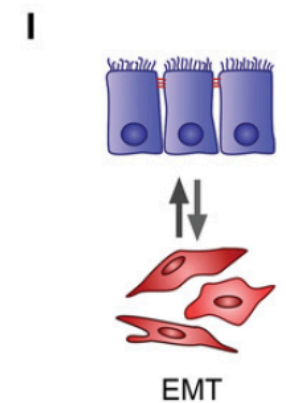
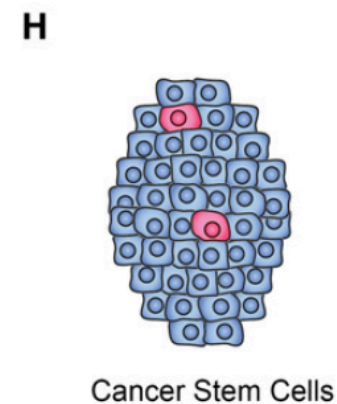
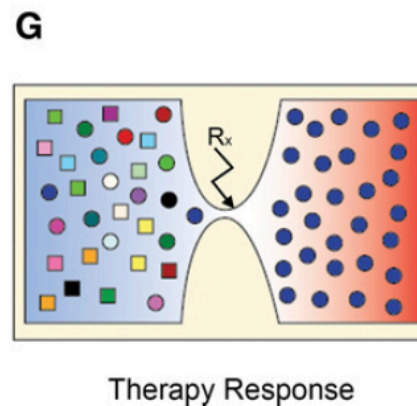
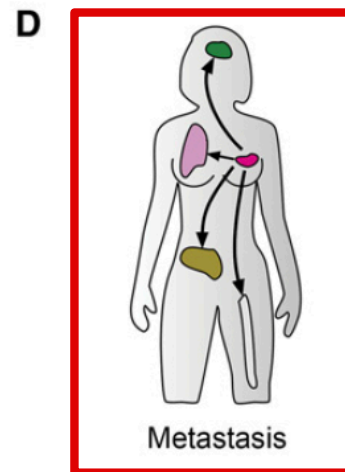
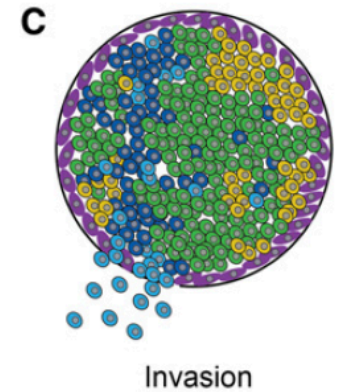
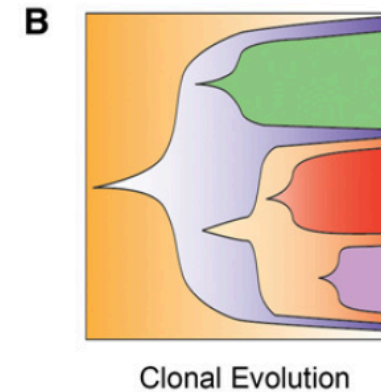


**A- Resolving intratumoral heterogeneity  
& dissecting microenvironment**

**B- Investigating clonal evolution in  
primary tumors**

**C- Studying invasion in early stage cancers**

**D- Tracking metastatic dissemination**



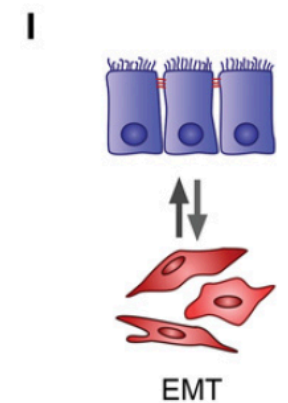
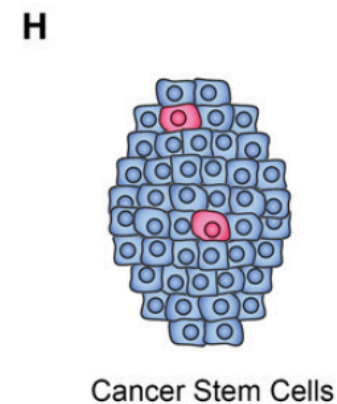
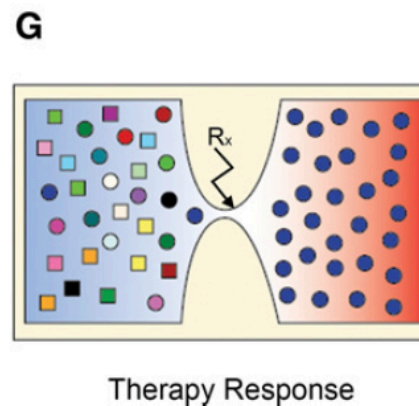
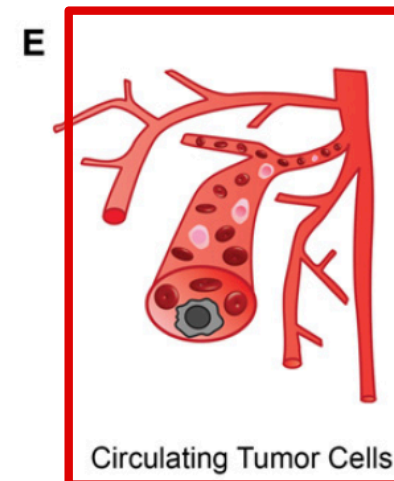
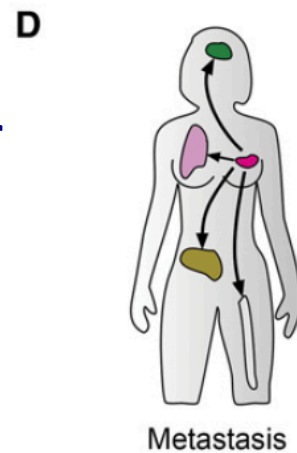
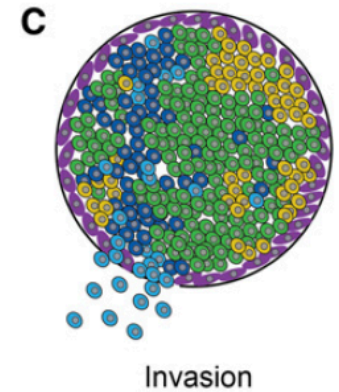
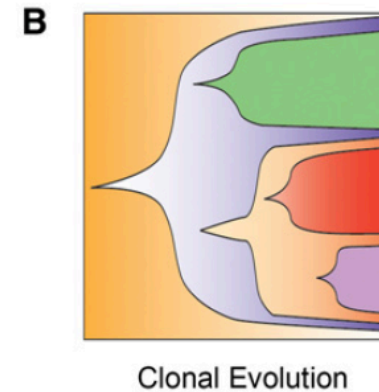
**A- Resolving intratumoral heterogeneity  
& dissecting microenvironment**

**B- Investigating clonal evolution in  
primary tumors**

**C- Studying invasion in early stage cancers**

**D- Tracking metastatic dissemination**

**E- Genomic profiling of circulating tumor  
cells**





**A- Resolving intratumoral heterogeneity & dissecting microenvironment**

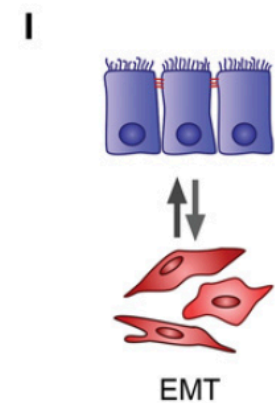
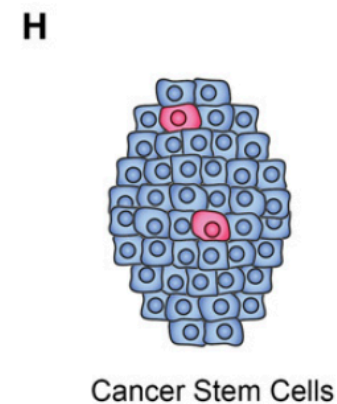
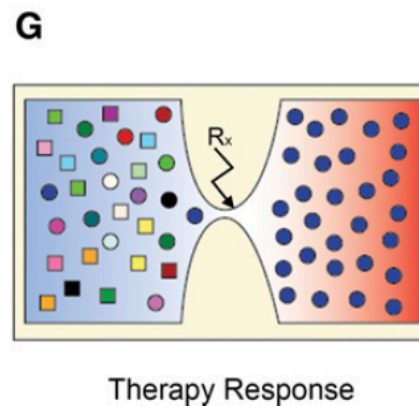
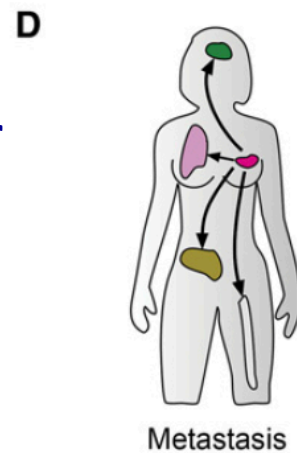
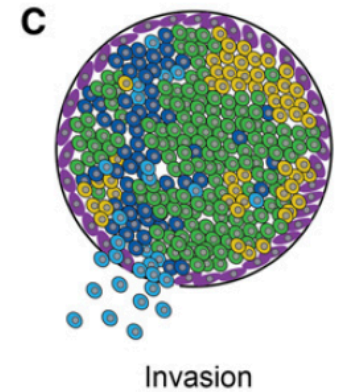
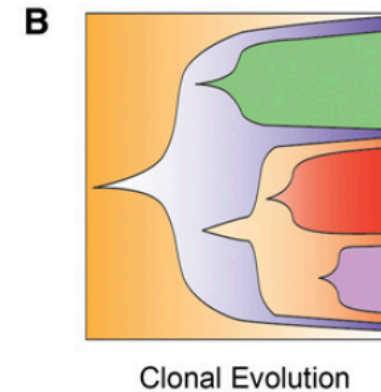
**B- Investigating clonal evolution in primary tumors**

**C- Studying invasion in early stage cancers**

**D- Tracking metastatic dissemination**

**E- Genomic profiling of circulating tumor cells**

**F- Studying mutation rare and mutated phenotypes**





**A- Resolving intratumoral heterogeneity & dissecting microenvironment**

**B- Investigating clonal evolution in primary tumors**

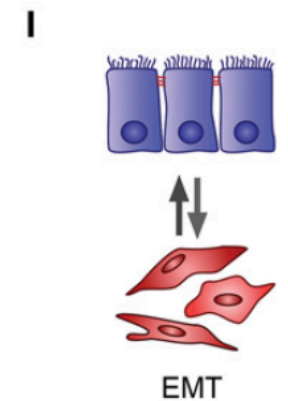
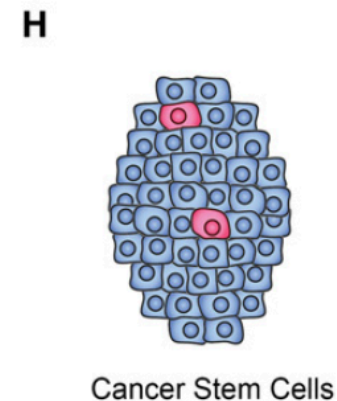
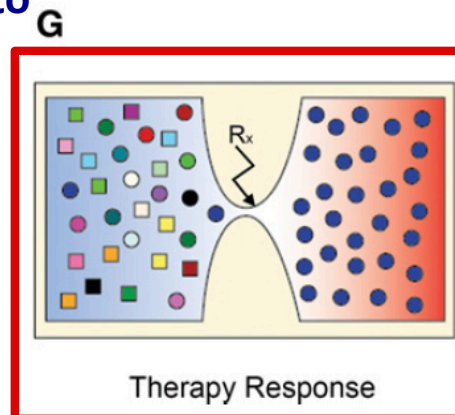
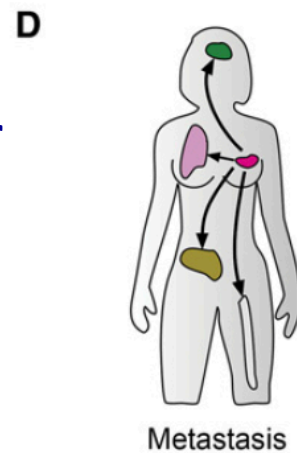
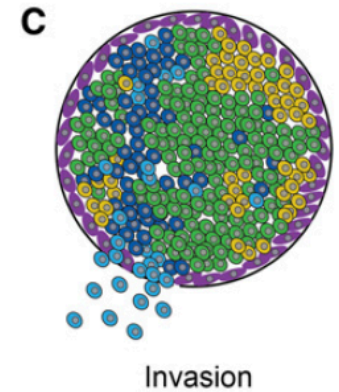
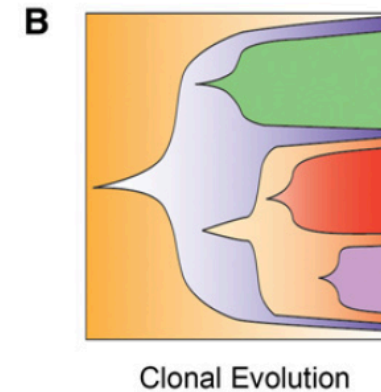
**C- Studying invasion in early stage cancers**

**D- Tracking metastatic dissemination**

**E- Genomic profiling of circulating tumor cells**

**F- Studying mutation rare and mutated phenotypes**

**G- Understanding resistance evolution to therapy**



**A- Resolving intratumoral heterogeneity & dissecting microenvironment**

**B- Investigating clonal evolution in primary tumors**

**C- Studying invasion in early stage cancers**

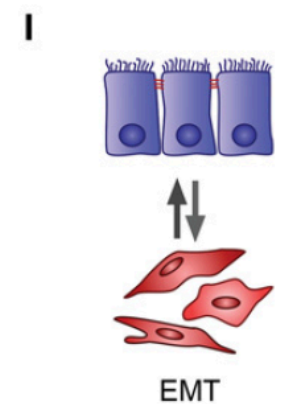
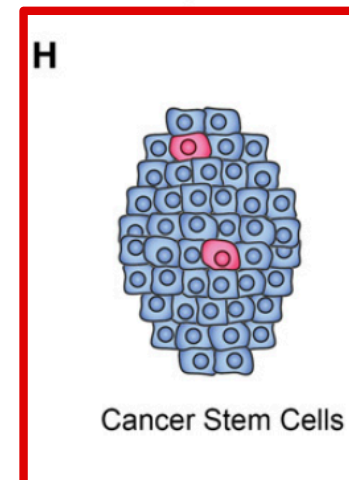
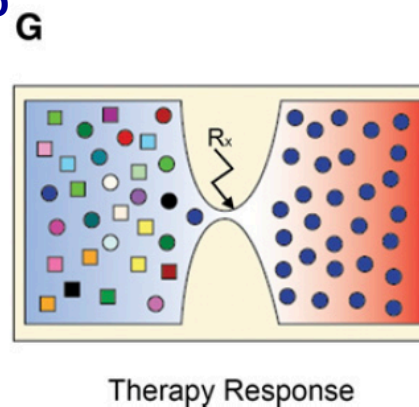
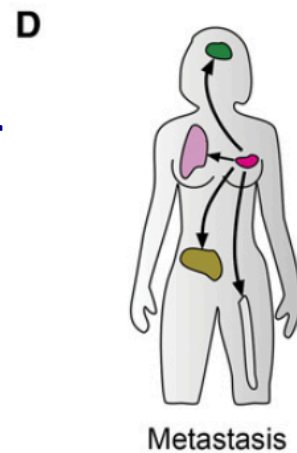
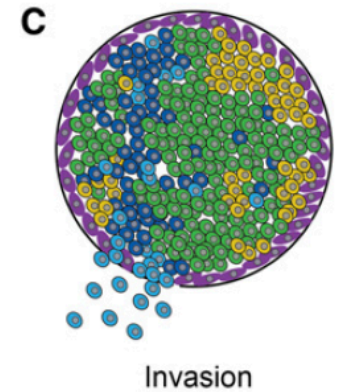
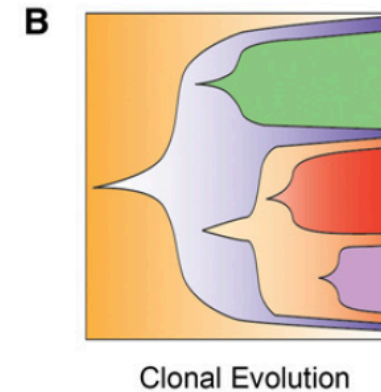
**D- Tracking metastatic dissemination**

**E- Genomic profiling of circulating tumor cells**

**F- Studying mutation rare and mutated phenotypes**

**G- Understanding resistance evolution to therapy**

**H- Understanding cancer stem cell & cell hierarchies**



**A- Resolving intratumoral heterogeneity & dissecting microenvironment**

**B- Investigating clonal evolution in primary tumors**

**C- Studying invasion in early stage cancers**

**D- Tracking metastatic dissemination**

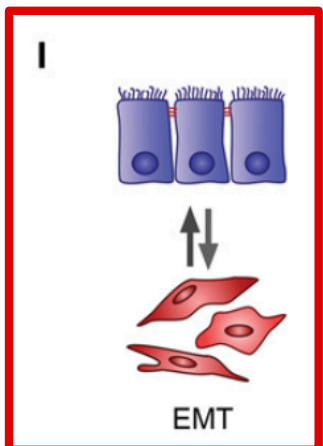
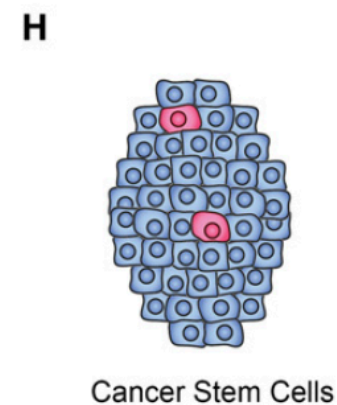
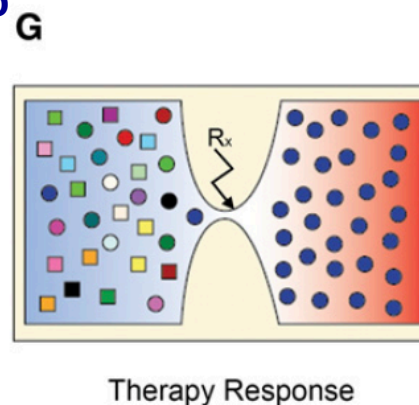
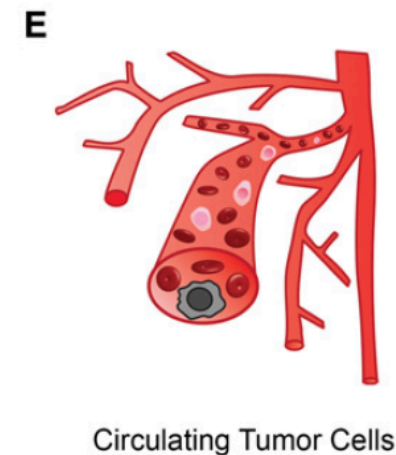
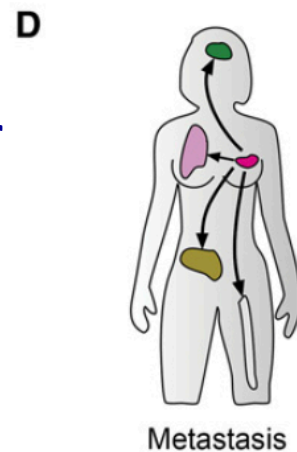
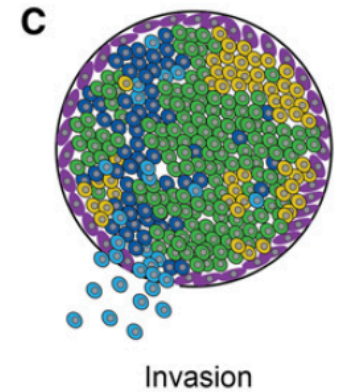
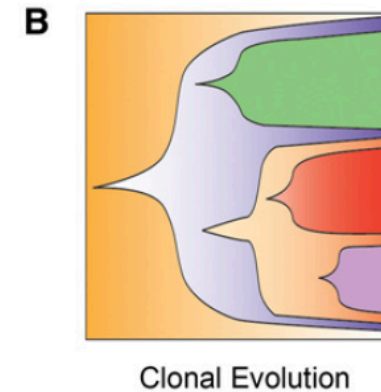
**E- Genomic profiling of circulating tumor cells**

**F- Studying mutation rare and mutated phenotypes**

**G- Understanding resistance evolution to therapy**

**H- Understanding cancer stem cell & cell hierarchies**

**I- Studying cell plasticity and epithelial-to-mesenchymal transition**



# **Strategies for census and validation**

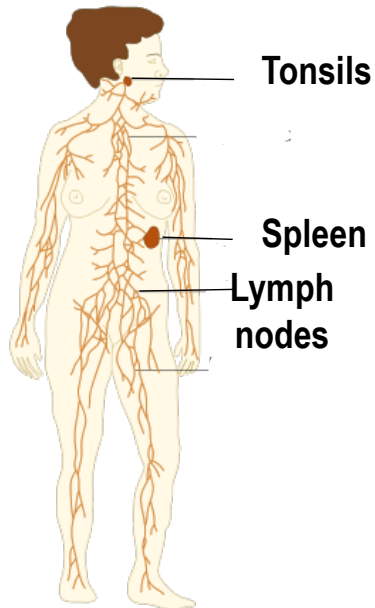




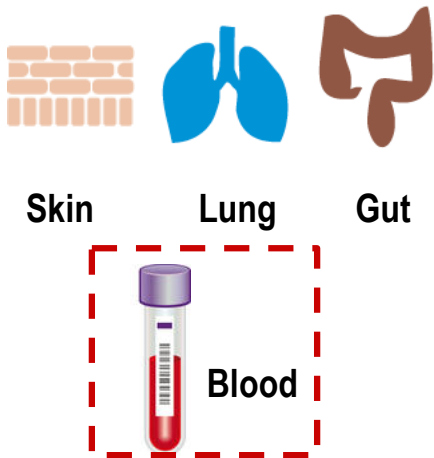
# Phase I: Generating unbiased DC map

Healthy tissue to be profiled

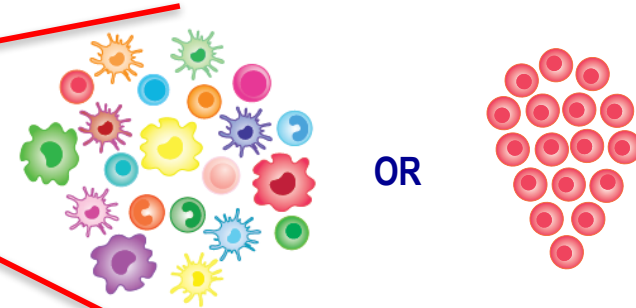
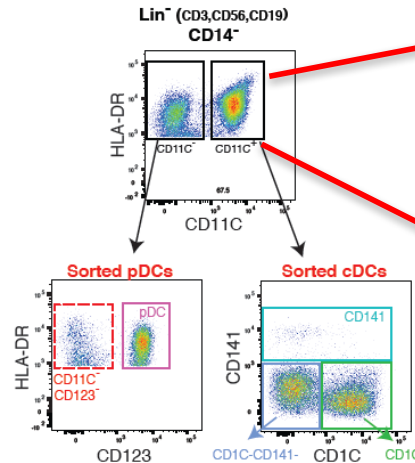
Lymphoid Organs



Non-Lymphoid "Barrier" Organs



Sample dissociation, enrichment

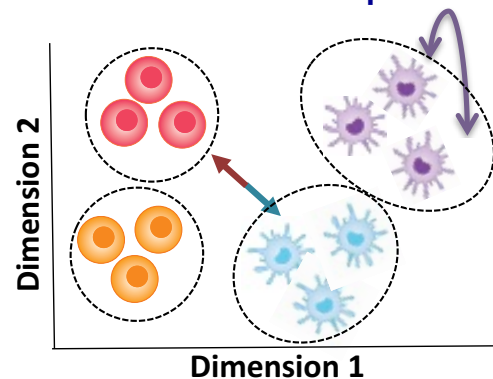


Single cell RNA-Sequencing → Plate-based (e.g. **SS2**, Cel-Seq2, SCRB-Seq)  
→ Droplet-based (e.g. 10X, DropSeq, InDrop)



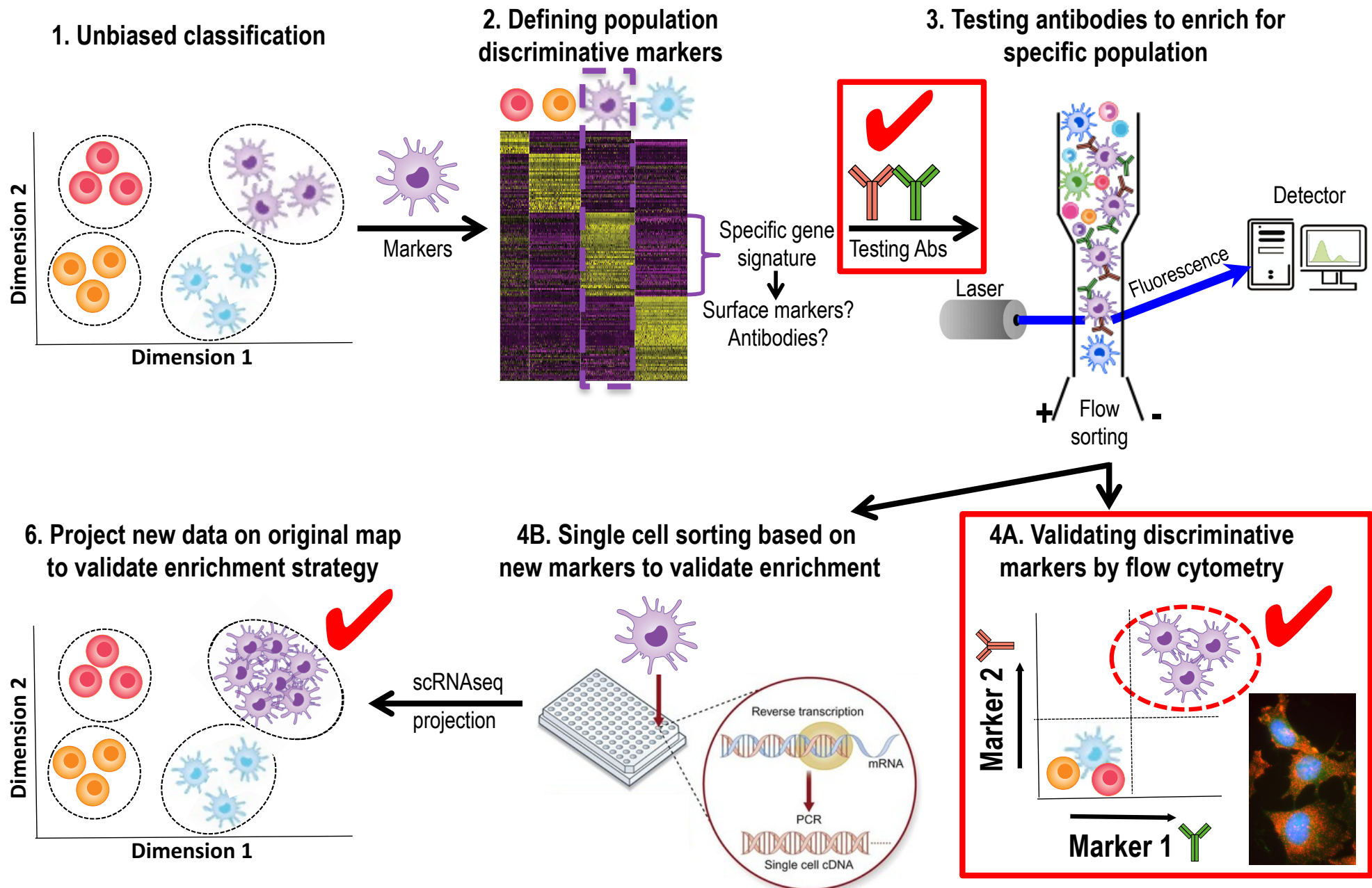
Expression profile clustering

Cell subsets map



- ✓ Cell type identification
- ✓ Deconvolution of population structure
- ✓ Identification of markers
- ✓ Variability of transcription
- ✓ Regulatory network inference

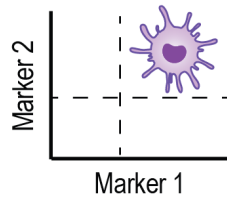
# Phase II: Enriching for new predicted cell populations – developing & validating reagents and isolations strategies



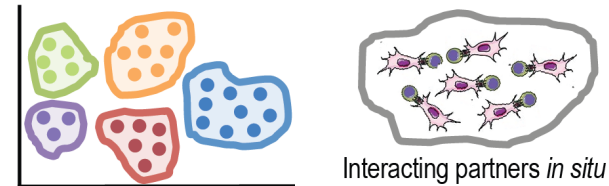
# Phase III: Functionally defining uniqueness of predicted new cell population in health and disease

## A- Functional Study & Characterization

1. Validating new markers



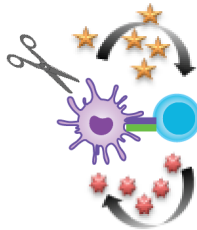
2. CyTOF, FACS, secretion, functional analysis



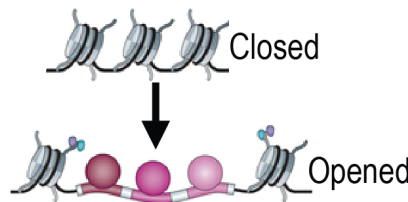
Static & dynamic  
characterization



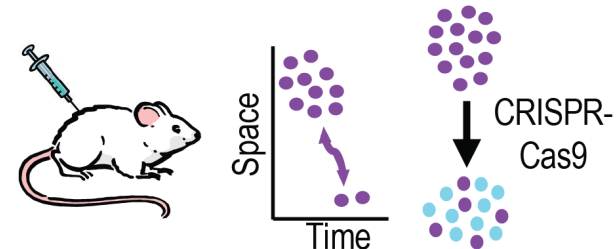
5. Model cell-cell  
interaction *in vitro*



4. Regulatory  
Landscape



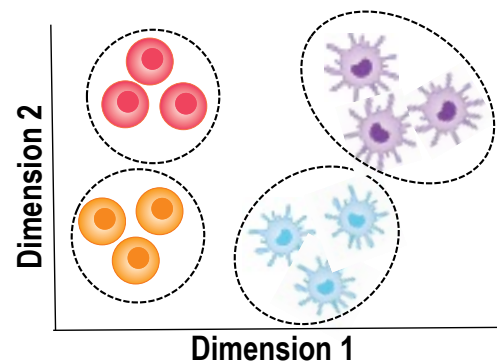
3. Humanized mice studies: live imaging  
over time, loss-of-function experiments



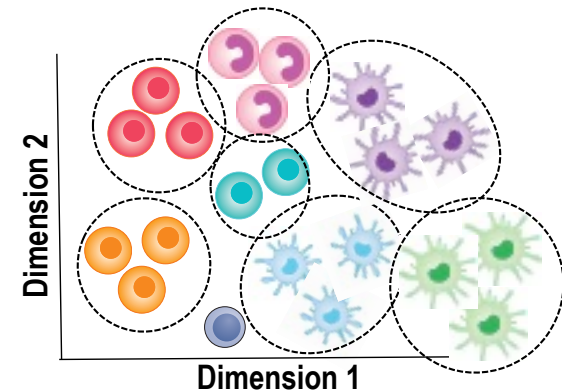
---

## B- Mapping Disease

Healthy State

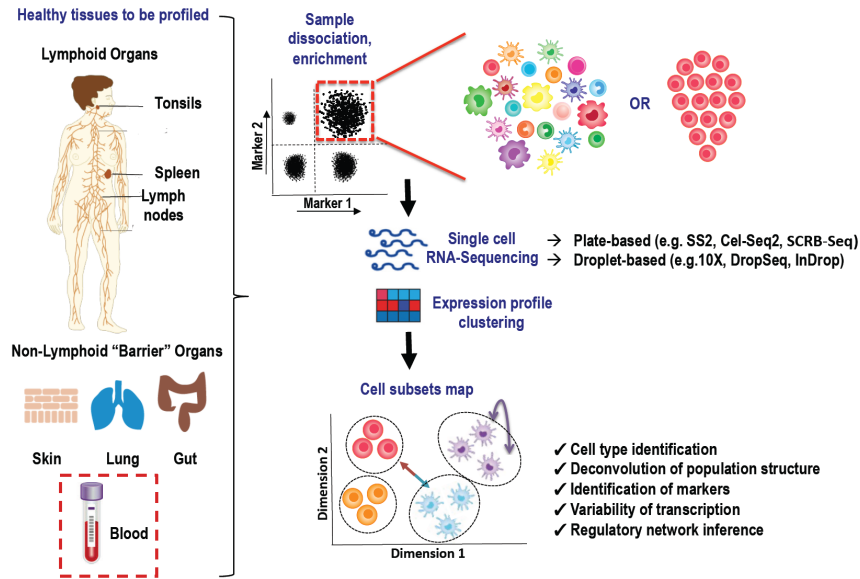


Disease State

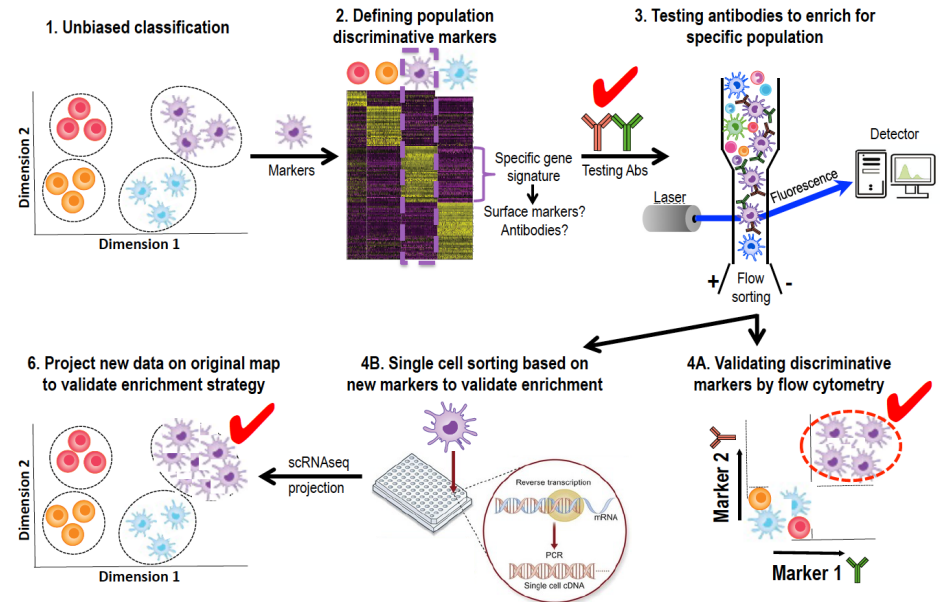




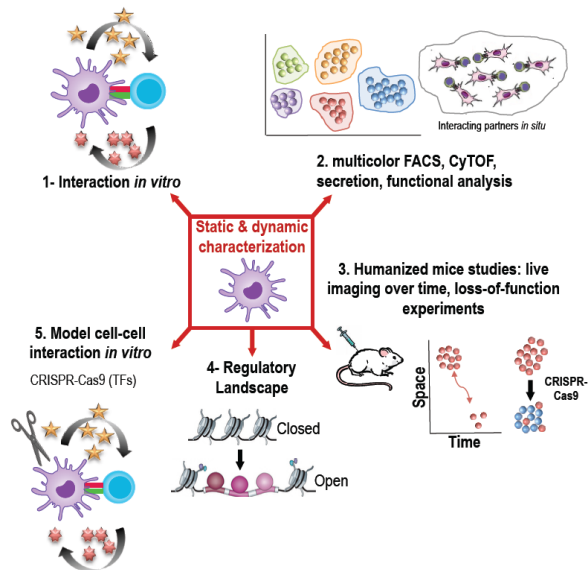
## 1. unbiased DC map



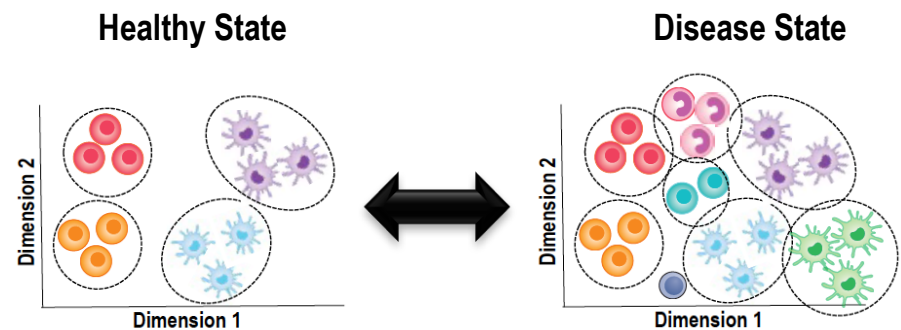
## 2. Identifying & validating new markers and gating strategies



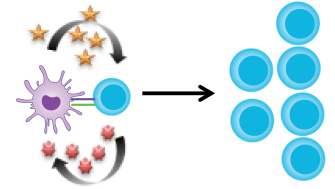
## 3. Functional characterization



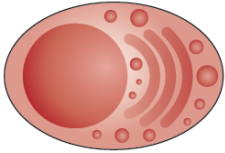
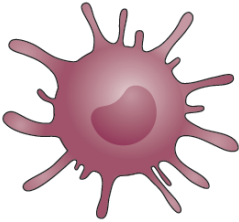
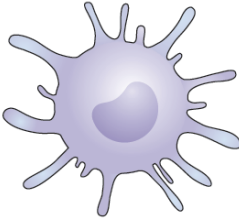
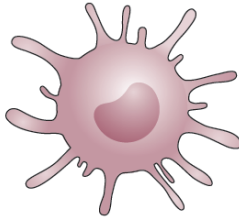



## 4. Mapping & studying disorders



# Dendritic Cells (DCs) & Monocytes



- DCs  $\approx$  1-3% & monocytes  $\approx$  10-25% in blood
- DCs function in pathogen sensing, antigen presentation, T cell activation
- Monocytes role in phagocytosis, cytokine production, macrophage source
- Involved in several auto-immune diseases & cancers; therapeutic target
- Several subtypes have been defined:

pDC	CD141 <sup>+</sup>	CD1c <sup>+</sup>	CD1c <sup>-</sup> CD141 <sup>-</sup>	CD14 <sup>+</sup> CD16 <sup>lo</sup>	CD14 <sup>+</sup> CD16 <sup>+</sup>	CD16 <sup>+</sup> CD14 <sup>lo</sup>
						
15-20% of DC	3-5% of DC	19-25% of DC	50-70% of DC	75-80% of mono	2-5% of mono	10-15% of mono
Interferon production	Antigen presentation to CD8 <sup>+</sup> T cells	Inflammation: Ag presentation to CD4 <sup>+</sup> T cells				

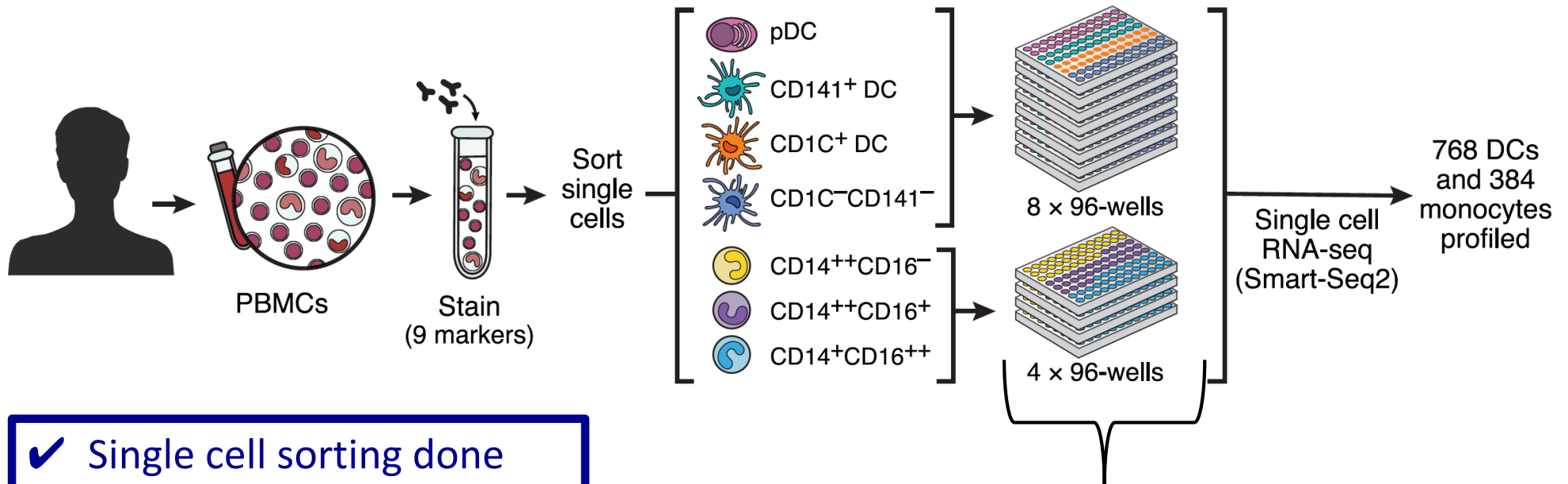
## Answering key questions to discover & characterize all blood dendritic cell (DC) & monocyte subsets

- 1) How many subsets can be found in blood?
- 2) Do they have the expected markers?
- 3) Can we identify better markers?
- 4) Is there heterogeneity within the major subsets?
- 5) Are there previously uncharacterized subsets?
- 6) Can these subsets be used to map cells in human disease?

## Answering key questions to discover & characterize all blood dendritic cell (DC) & monocyte subsets

- 1) How many subsets can be found in blood?
- 2) Do they have the expected markers?
- 3) Can we identify better markers?
- 4) Is there heterogeneity within the major subsets?
- 5) **Are there previously uncharacterized subsets?**
- 6) Can these subsets be used to map cells in human disease?

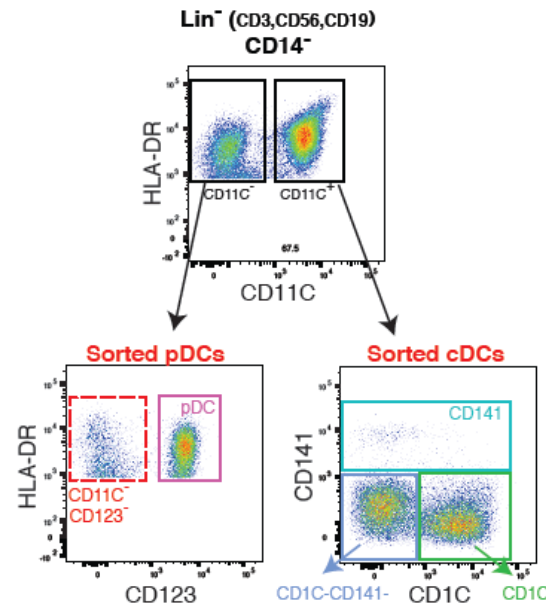
# How should we discover DC subsets?



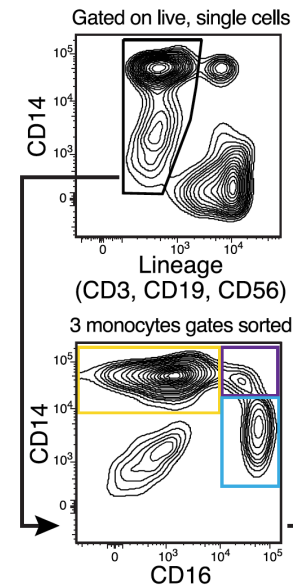
- ✓ Single cell sorting done from a constant source
- ✓ Cell sorting with optimized panel of markers
- ✓ Deep Sequencing (1-2M reads/cells)

## Adaptive sampling strategy:

### DCs:

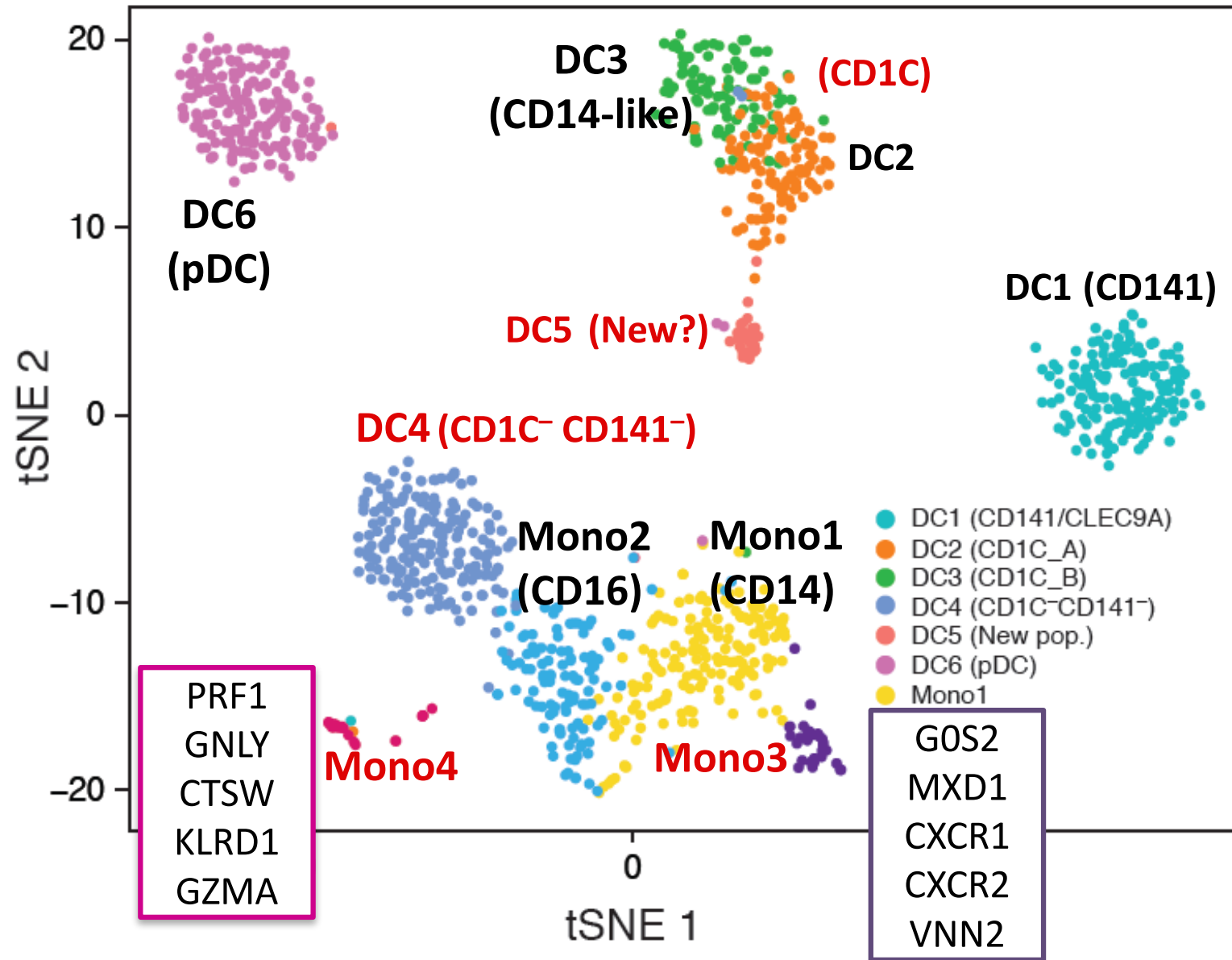


### Monocytes:



**How many subsets can be  
found in blood?**

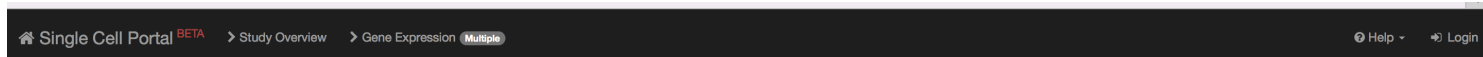
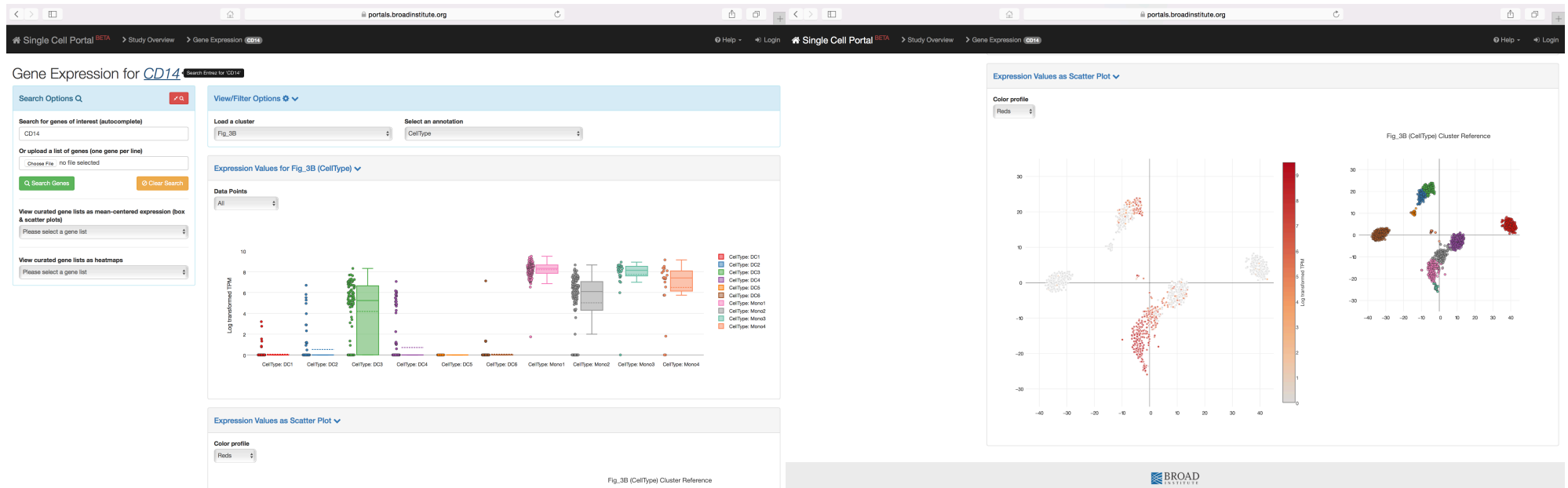
# Observed 6 DC & 4 Monocyte clusters in blood





# Data uploaded to single cell portal

[https://portals.broadinstitute.org/single\\_cell/study/atlas-of-human-blood-dendritic-cells-and-monocytes](https://portals.broadinstitute.org/single_cell/study/atlas-of-human-blood-dendritic-cells-and-monocytes)

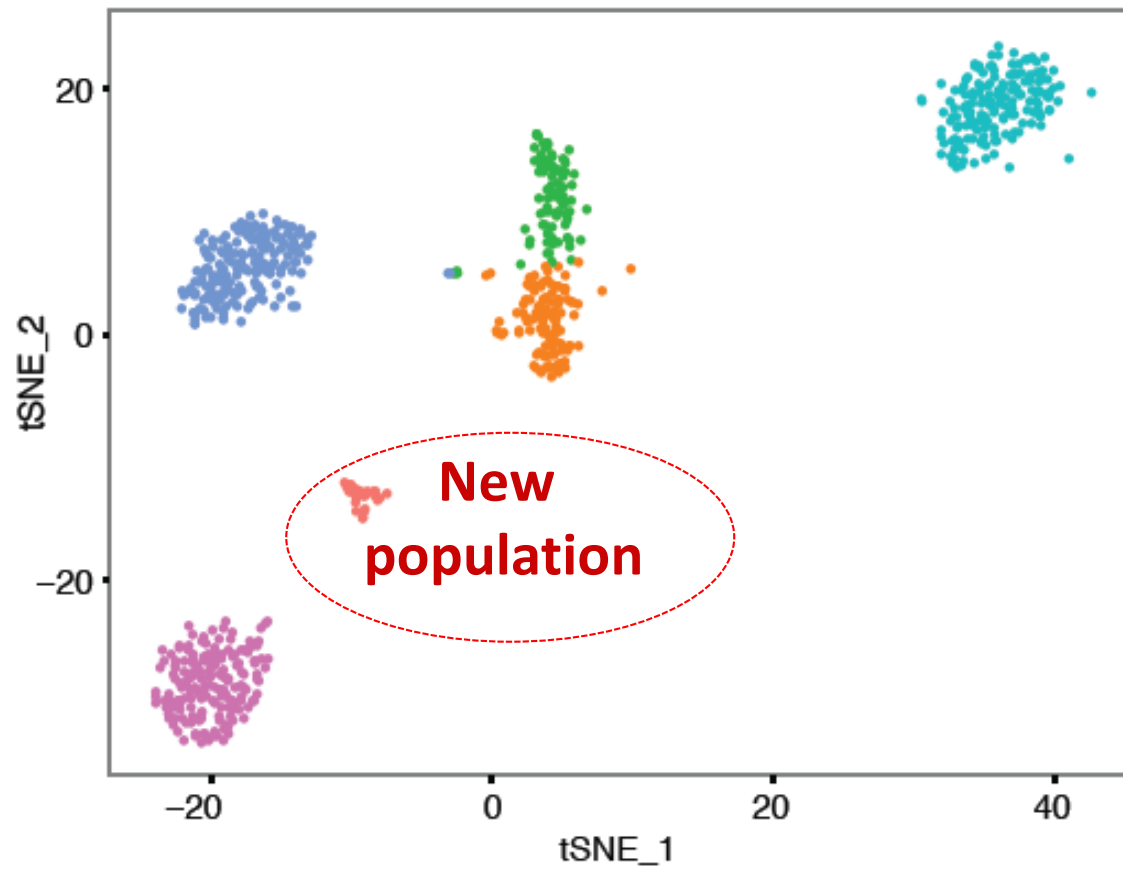


Gene Expression for *PRF1*, *GNLY*, *CTSW*, *FGFBP2*, *IL2RB*... [7 more](#)



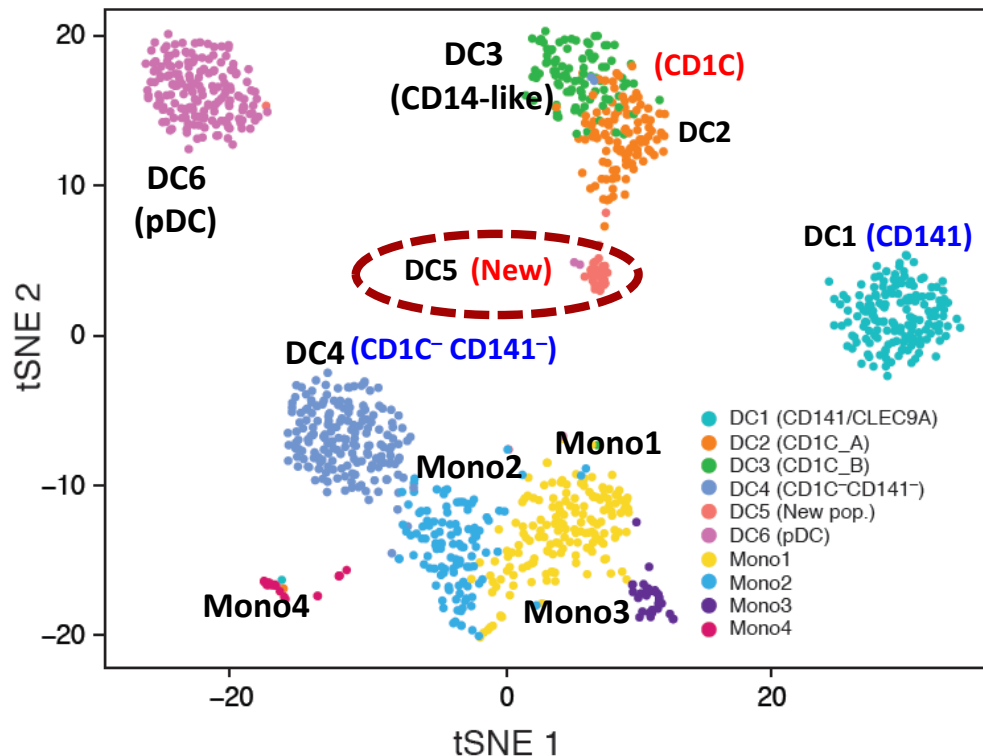
Tim Tickle

# What is the uncharacterized DC subset?



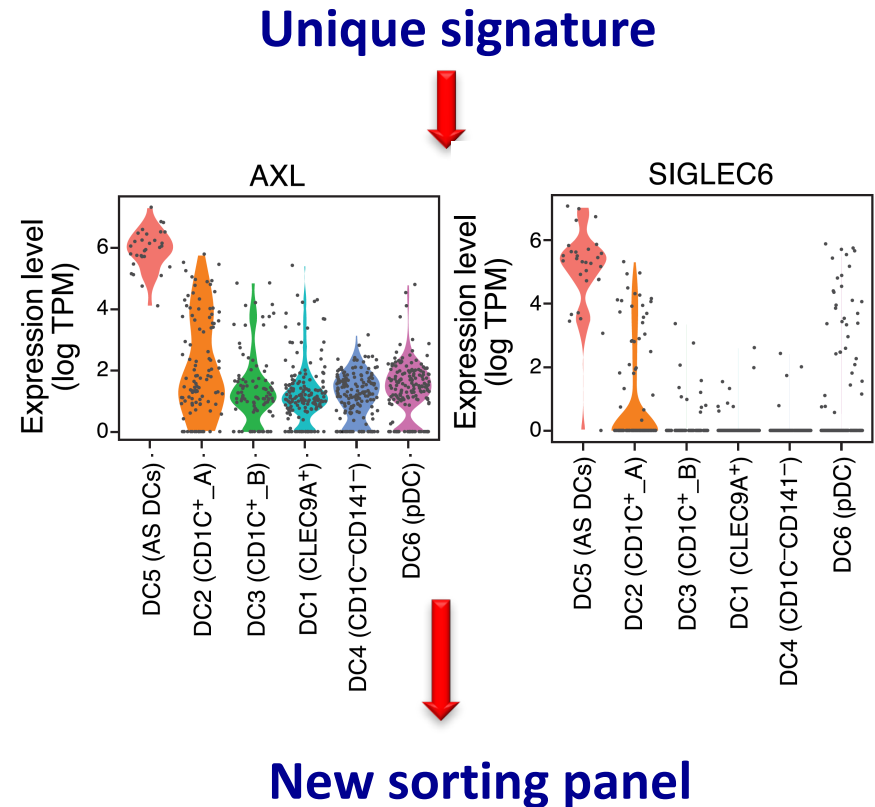
# Automatic multi-dimensional classifier predicts the presence of rare new DC subset

## DC5 Challenge: rare ( $\approx 0.06\%$ of PBMCs)

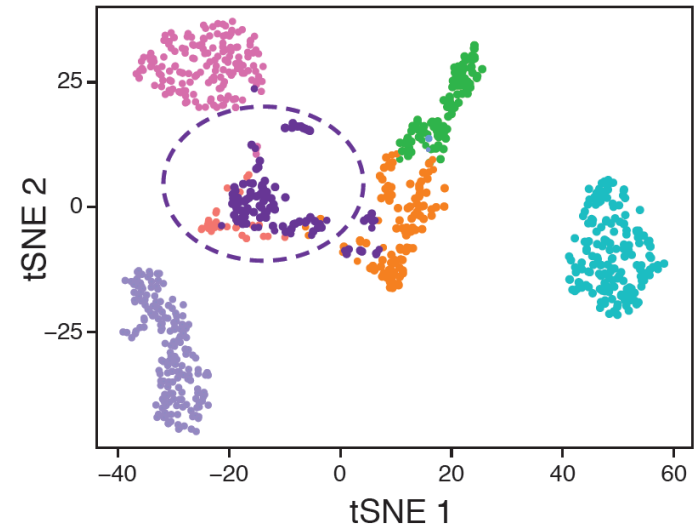
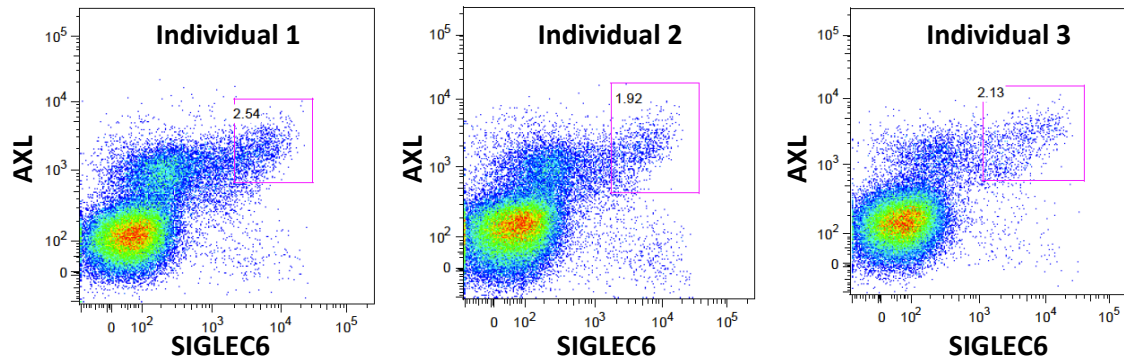


## Approach:

- (1) Find markers specific for population
- (2) Develop new sorting panel
- (3) Profile cells from additional individuals



# Validation of DC5 population existence by flow cytometry & scRNAseq of prospectively isolated cells



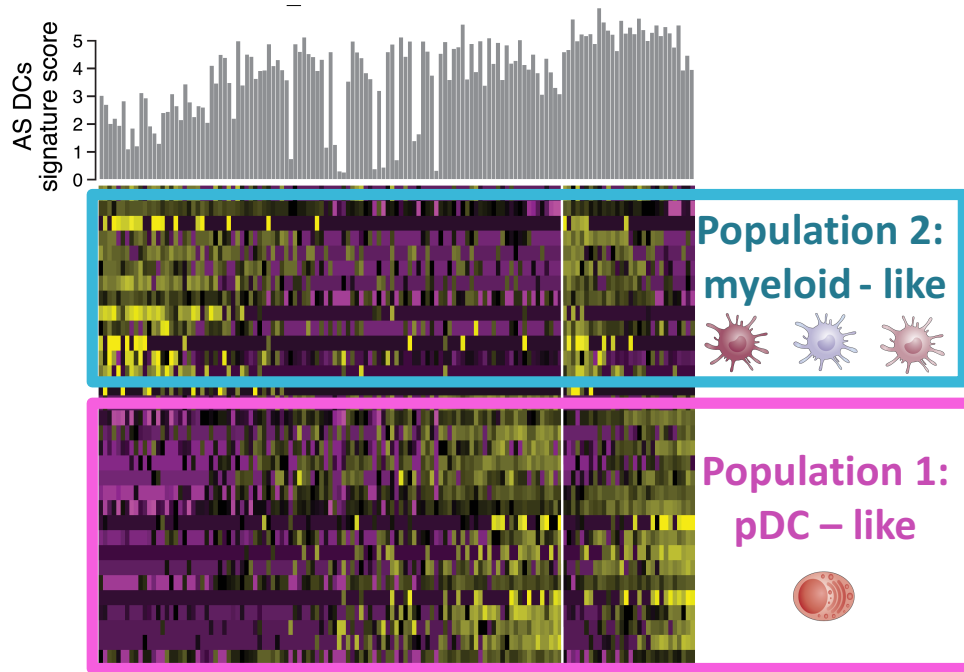
→ New DC population observed in **ALL 10 patients**  
→ Represents **1.9-3.2% of the DC / 0.04-0.064% of PBMCs**

- **What do they look like?**
  - Transcriptionally – what's distinct and common
  - Morphology
- **What are its communication capabilities with other cells**
  - Receptors, secreted factors
  - Co-culture with other cell types
- **Who are the direct interacting partners**
  - *In-situ* co-localization staining
- **Where are these cells in DC gating strategy?**

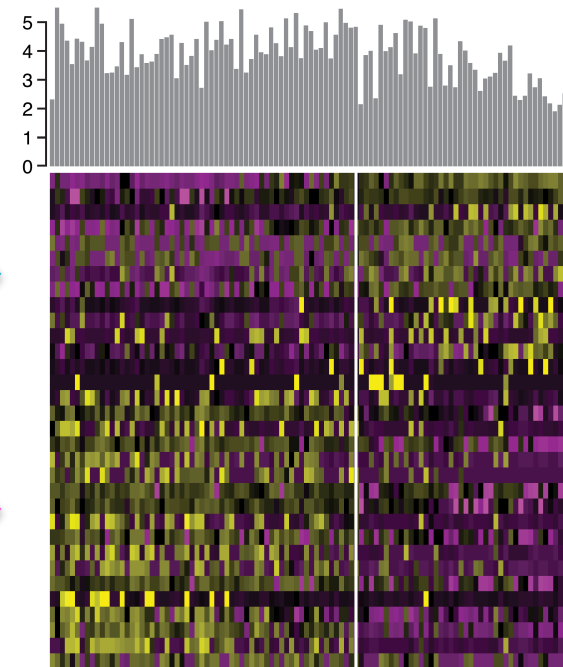
# DC5 new population falls along continuum with 2 clear extremes

## Successful enrichment of both subsets & validation in 10 healthy individuals

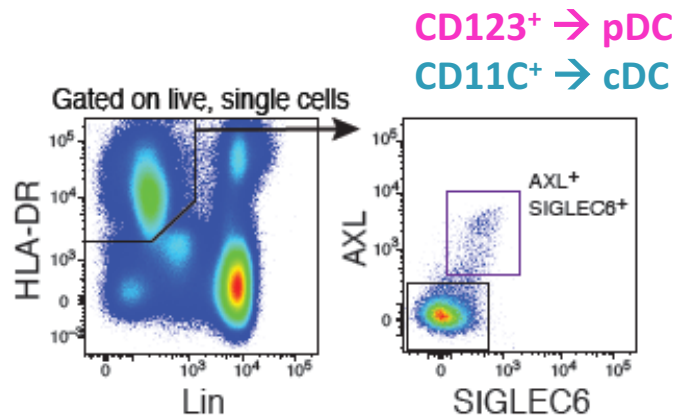
DC5 falls across 2 extremes



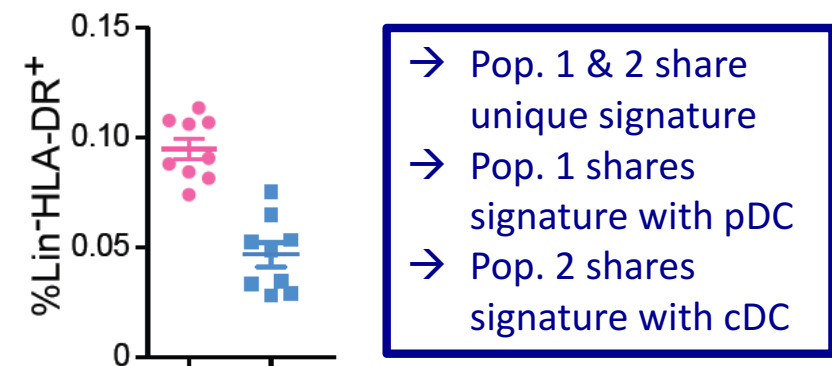
B- Validation of 2 putative subsets



A- NEW panel to enrich for both putative populations



C- Validation across 10 individuals



Villani et al. *Science* 2017 Apr 21;356(6335).

# Concluding thoughts

- Single cell genomics methods are becoming an essential tool for dissecting biology at an unprecedented resolution
- Single cell multi-omics will empower new definition of cell types/states and tissue
- Being able to track live cells over time will be truly transformative
- Scale will continue growing and price will come down  
→empowering translational efforts!
- New analyses techniques and framework are needed to handle such large dataset

# A Word of Caution

**“Tempering some of the enthusiasm are myriad challenges inherent to the process, from the isolation of cells, to amplification of their genomes or transcriptomes, to making sense of the data. Cost is also a consideration leaving good reason to carefully select situations that justify going to the single-cell level.”**

## **Bottom Line:**

*Single cell transcriptomics is not the solution to answering every biological question!*



# Acknowledgements



**Rahul Satija**

**Hacohen Group**

**Siranush Sarkizova**

Weibo Li

All lab members

**Christophe Benoist**

David Puyraimond-Zemmour

**Philip De Jager**

Alina Von Korff

Laura Glick



**Nir Hacohen**

**Muzlifah Haniffa**

Gary Reynolds

James Fletcher

Laura Jardine

Andrew Filby

**Ragon Institute:**

Marcus Altfeld

Morgane Griesbeck

Ryan Park

Michael Waring

Adam Chicoine



**Aviv Regev**

**Regev Group**

Karthik Shekhar

**Orit Rozenblatt-Rosen**

**Dana-Farber Cancer Institute**

**Andrew A. Lane**

Suzan Lazo-Kallanian

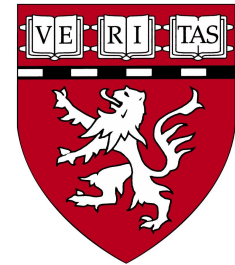
**Olink Proteomics**

Ida Grundberg

Emil Nilsson

Questions: [cvillani@broadinstitute.org](mailto:cvillani@broadinstitute.org)

# Harvard Medical School MGH Single Cell Genomics Research Program



**Villani Lab: postdoc positions available**  
**Contact: [cvillani@broadinstitute.org](mailto:cvillani@broadinstitute.org)**