

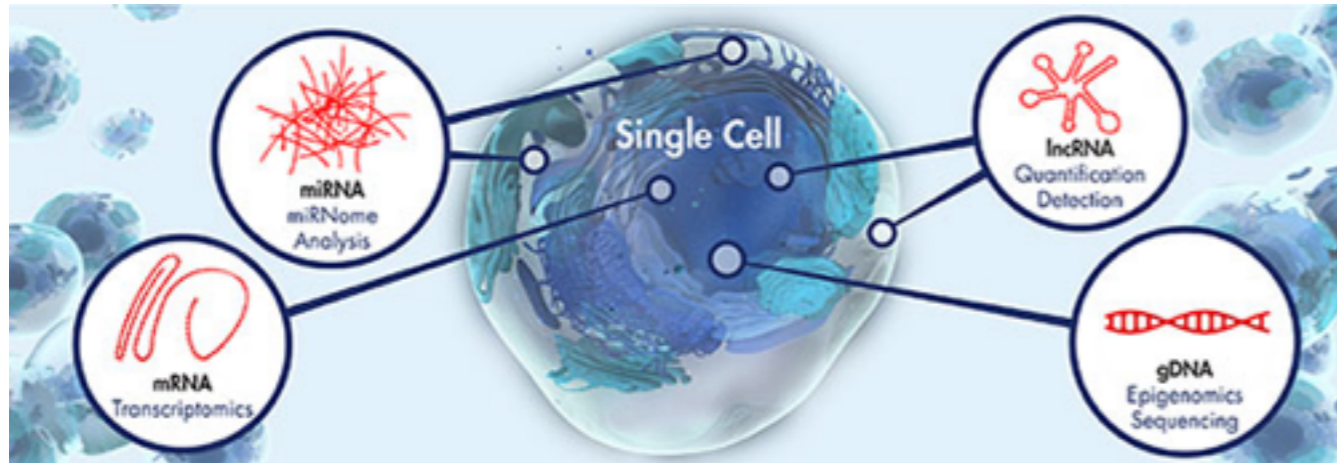


CS262 Winter 2016

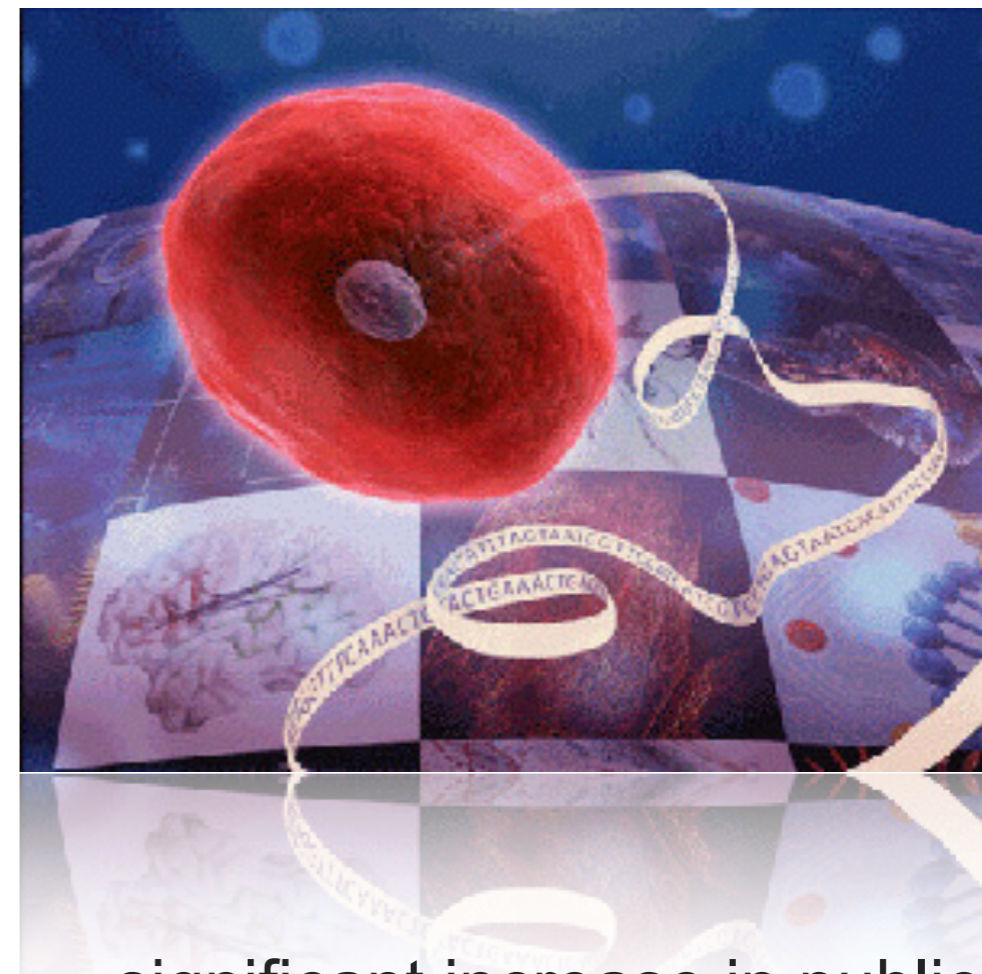
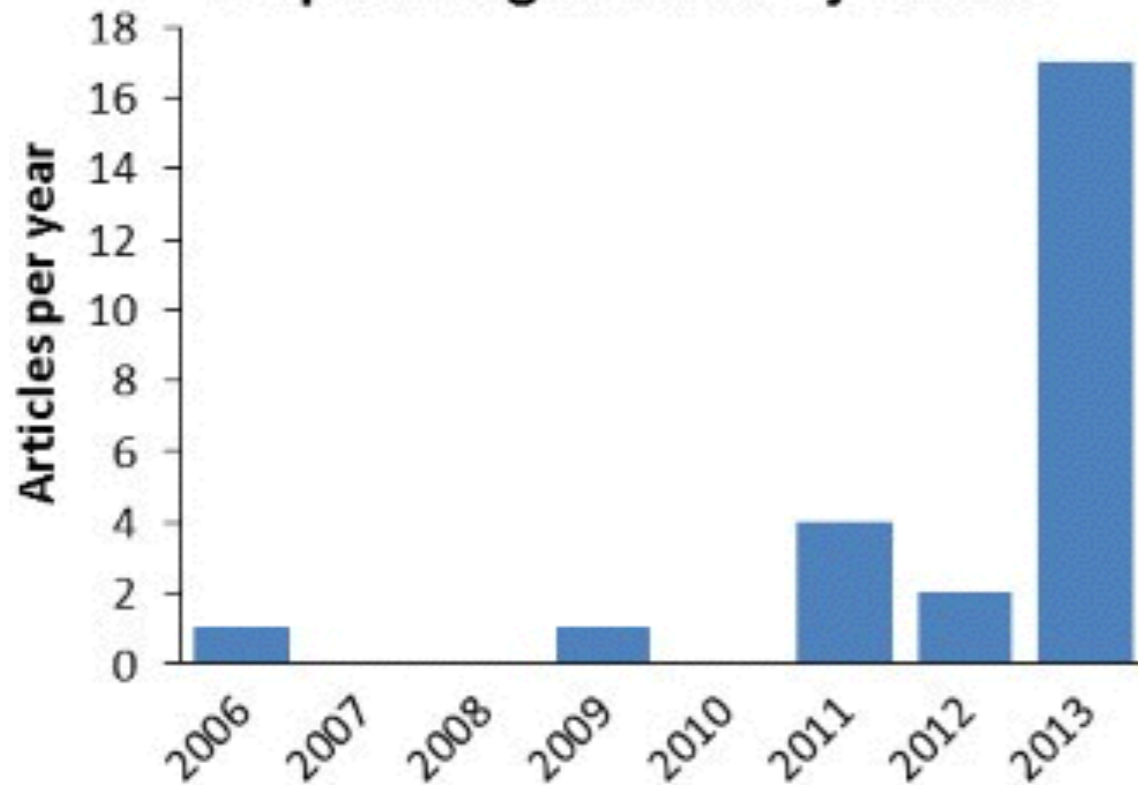
Single Cell Sequencing

# Background

*NATURE METHODS* | METHAGORA



Research articles using single-cell sequencing in Nature journals



significant increase in publications and data in the last two years

# Background



## Single-cell analyses of transcriptional heterogeneity during drug tolerance transition in cancer cells by RNA sequencing

Mei-Chong Wendy Lee<sup>a,1</sup>, Fernando J. Lopez-Diaz<sup>b,1</sup>, Shahid Yar Khan<sup>a,2</sup>, Muhammad Akram Tariq<sup>a,3</sup>, Yelena Dayn<sup>c</sup>, Charles Joseph Vaske<sup>d</sup>, Amie J. Radenbaugh<sup>a</sup>, Hyunsung John Kim<sup>a</sup>, Beverly M. Emerson<sup>b,4</sup>, and Nader Pourmand<sup>a,4</sup>

FOCUS

---

CANCER

Cell

## Single Cell RNA-Sequencing of Pluripotent States Unlocks Modular Transcriptional Variation

Aleksandra A. Kolodziejczyk,<sup>1,2,5</sup> Jong Kyoung Kim,<sup>1,5</sup> Jason C.H. Tsang,<sup>2</sup> Tomislav Ilicic,<sup>1,2</sup> Johan Henriksson,<sup>1</sup> Kedar N. Natarajan,<sup>1,2</sup> Alex C. Tuck,<sup>1,3</sup> Xuefei Gao,<sup>2</sup> Marc Bühler,<sup>3</sup> Pentao Liu,<sup>2</sup> John C. Marioni,<sup>1,2,4,\*</sup> and Sarah A. Teichmann<sup>1,2,\*</sup>

## Combined Single-Cell Functional and Gene Expression Analysis Resolves Heterogeneity within Stem Cell Populations

Nicola K. Wilson,<sup>1,9</sup> David G. Kent,<sup>1,9</sup> Florian Büttner,<sup>2,9</sup> Mona Shehata,<sup>7</sup> Iain C. Macaulay,<sup>1</sup> Manuel Sánchez Castillo,<sup>1</sup> Caroline A. Oedekoven,<sup>1</sup> Evangelia Diamanti,<sup>1</sup> Reiner Schulte-Thyry Voet,<sup>3,6</sup> Carlos Caldas,<sup>7</sup> John Stingl,<sup>7</sup> Anthony R. Green,<sup>1</sup> Fabian J. Theis,<sup>2,8</sup> and

NATURE

## Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq

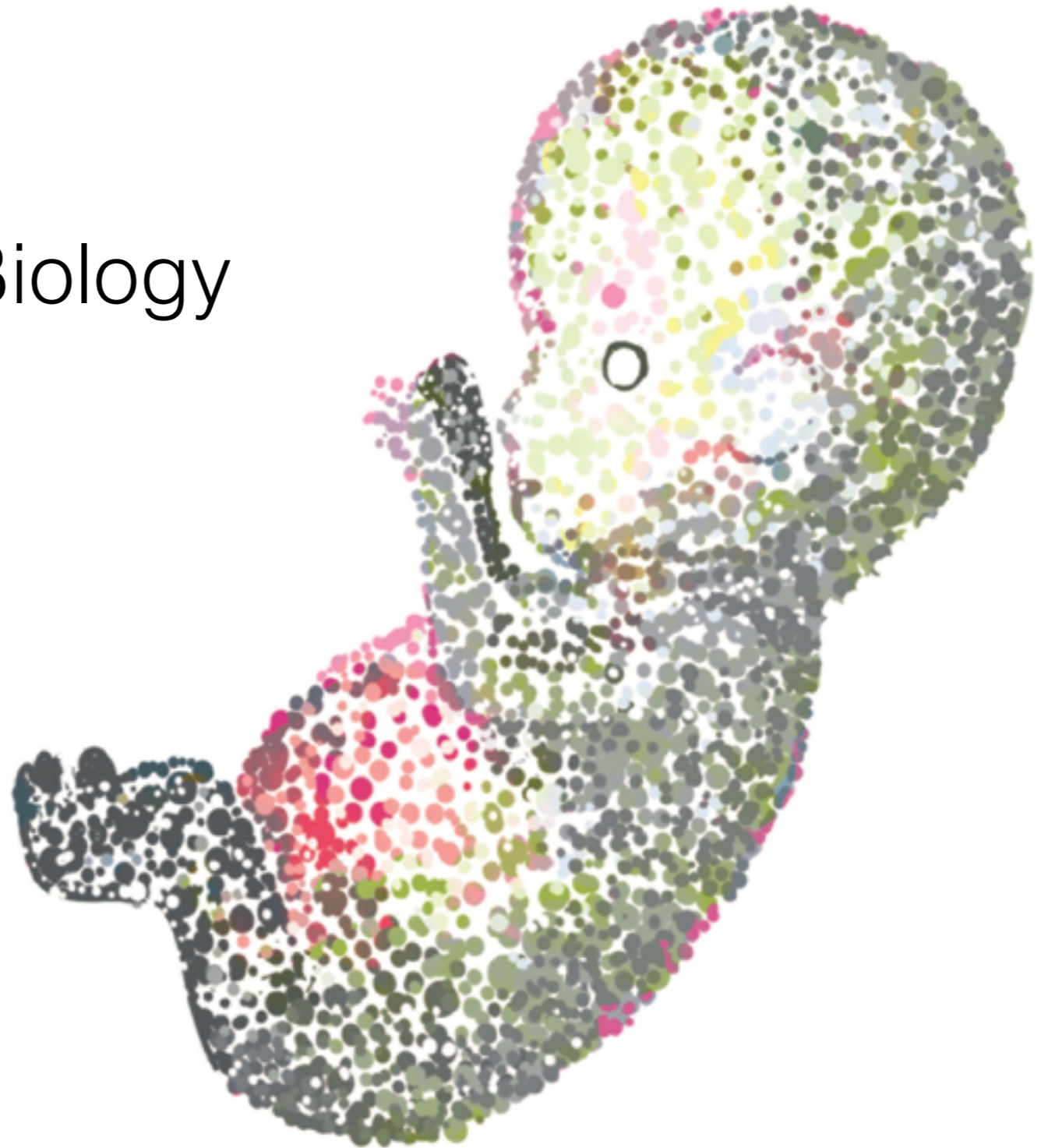
Amit Zeisel,<sup>1\*</sup> Ana B. Muñoz-Manchado,<sup>1\*</sup> Simone Codeluppi,<sup>1</sup> Peter Lönnerberg,<sup>1</sup> Gioele La Manno,<sup>1</sup> Anna Juréus,<sup>1</sup> Sueli Marques,<sup>1</sup> Hermany Munguba,<sup>1</sup> Liqun He,<sup>2</sup> Christer Betsholtz,<sup>2,3</sup> Charlotte Rolny,<sup>4</sup> Gonçalo Castelo-Branco,<sup>1</sup> Jens Hjerling-Leffler,<sup>1†</sup> Sten Linnarsson<sup>1†</sup>

ARTICLE

# Applications



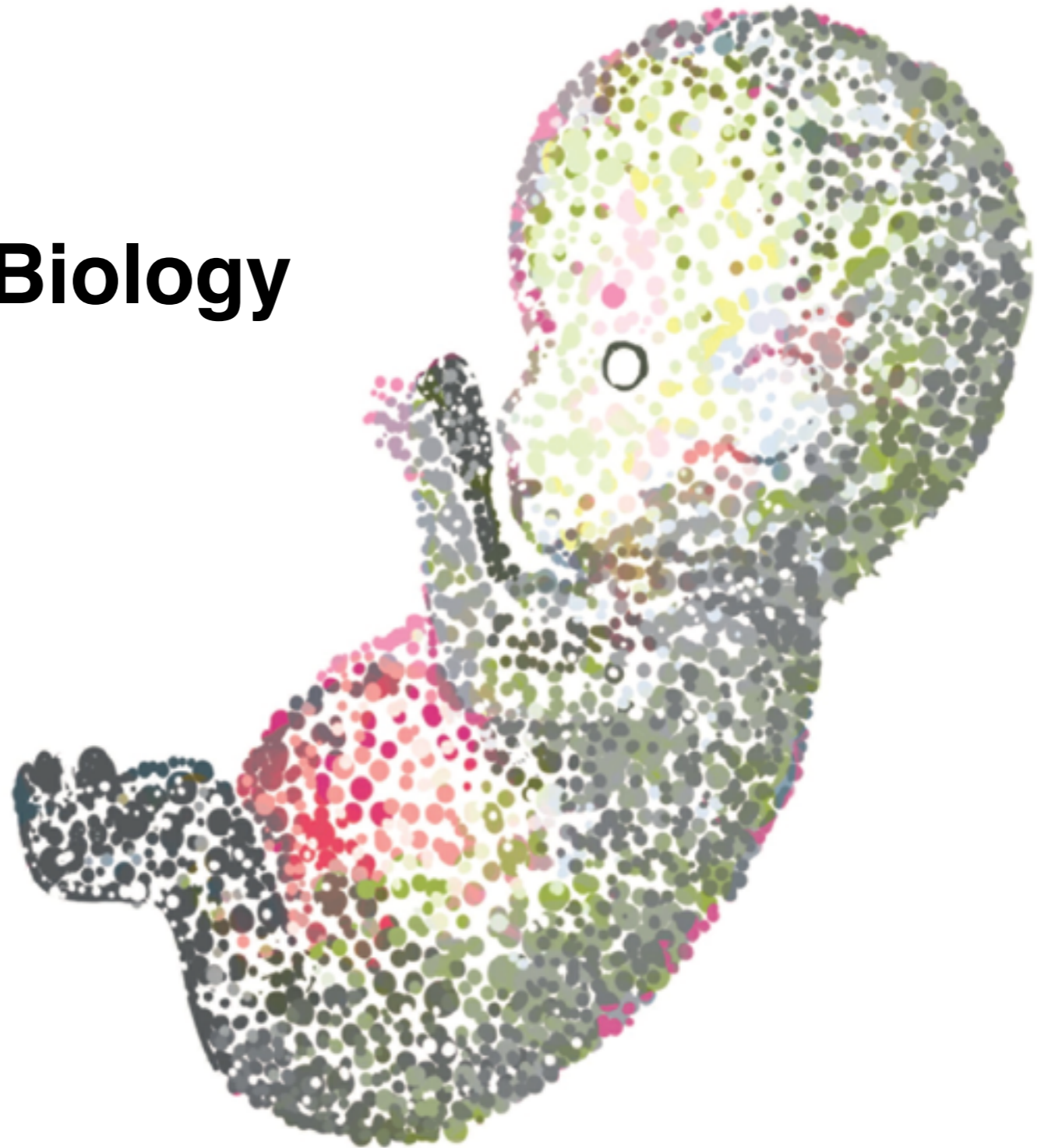
- Developmental Biology
- Cancer Biology
- Microbiology
- Neurology



# Applications

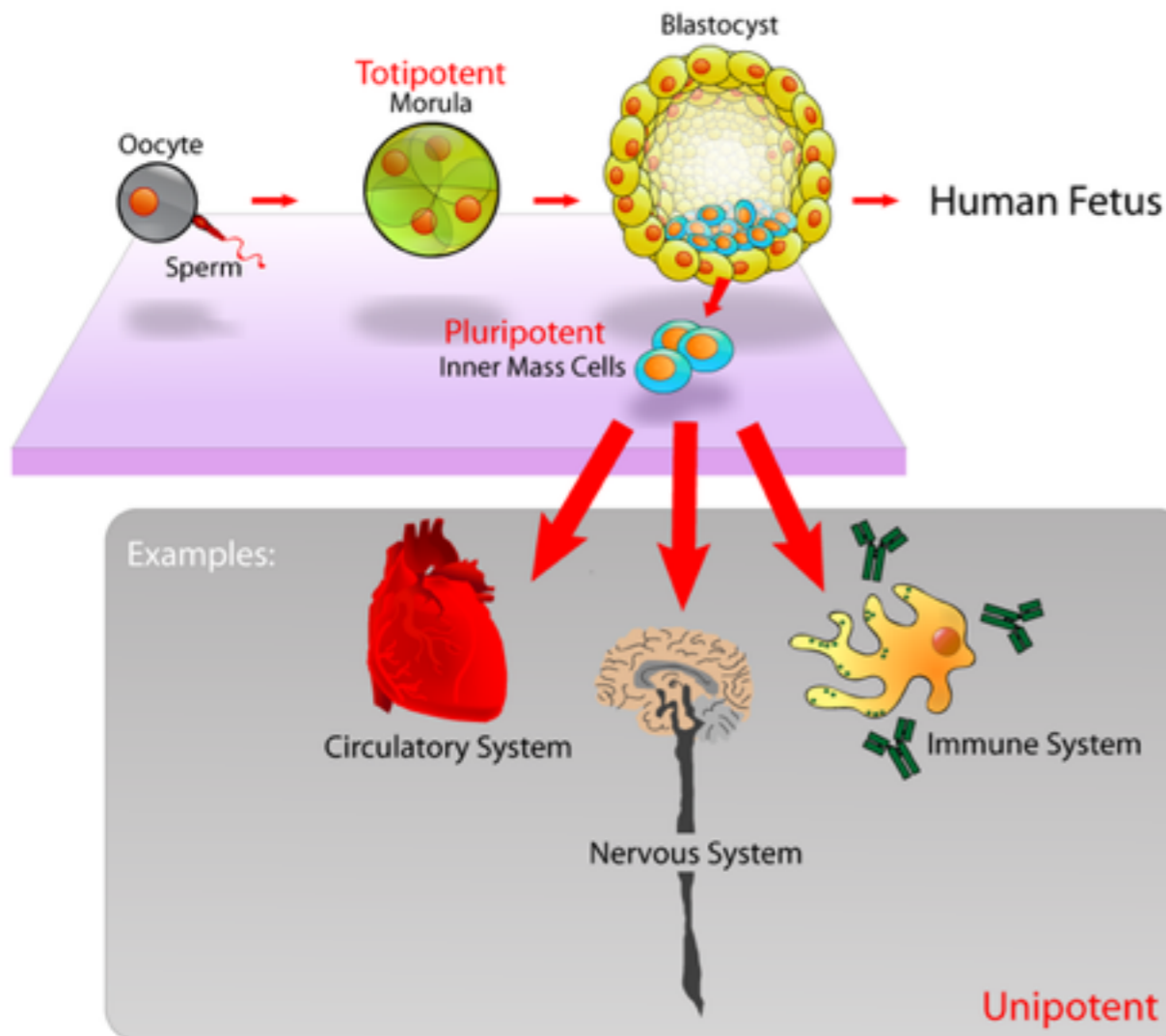


- **Developmental Biology**
- Cancer Biology
- Microbiology
- Neurology

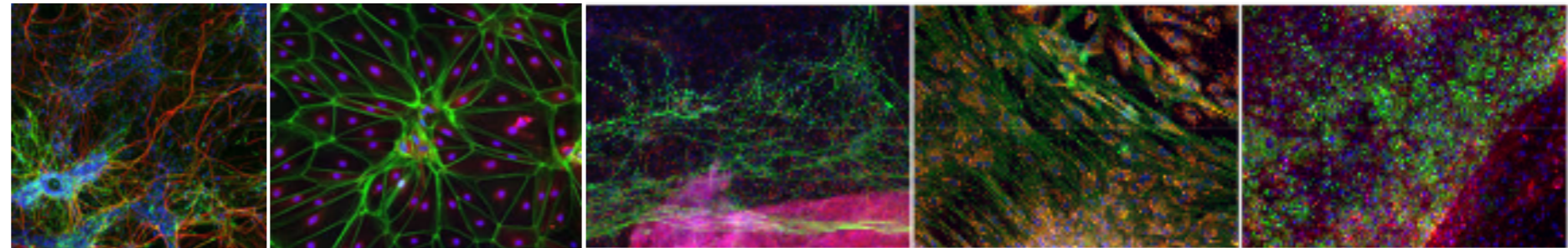
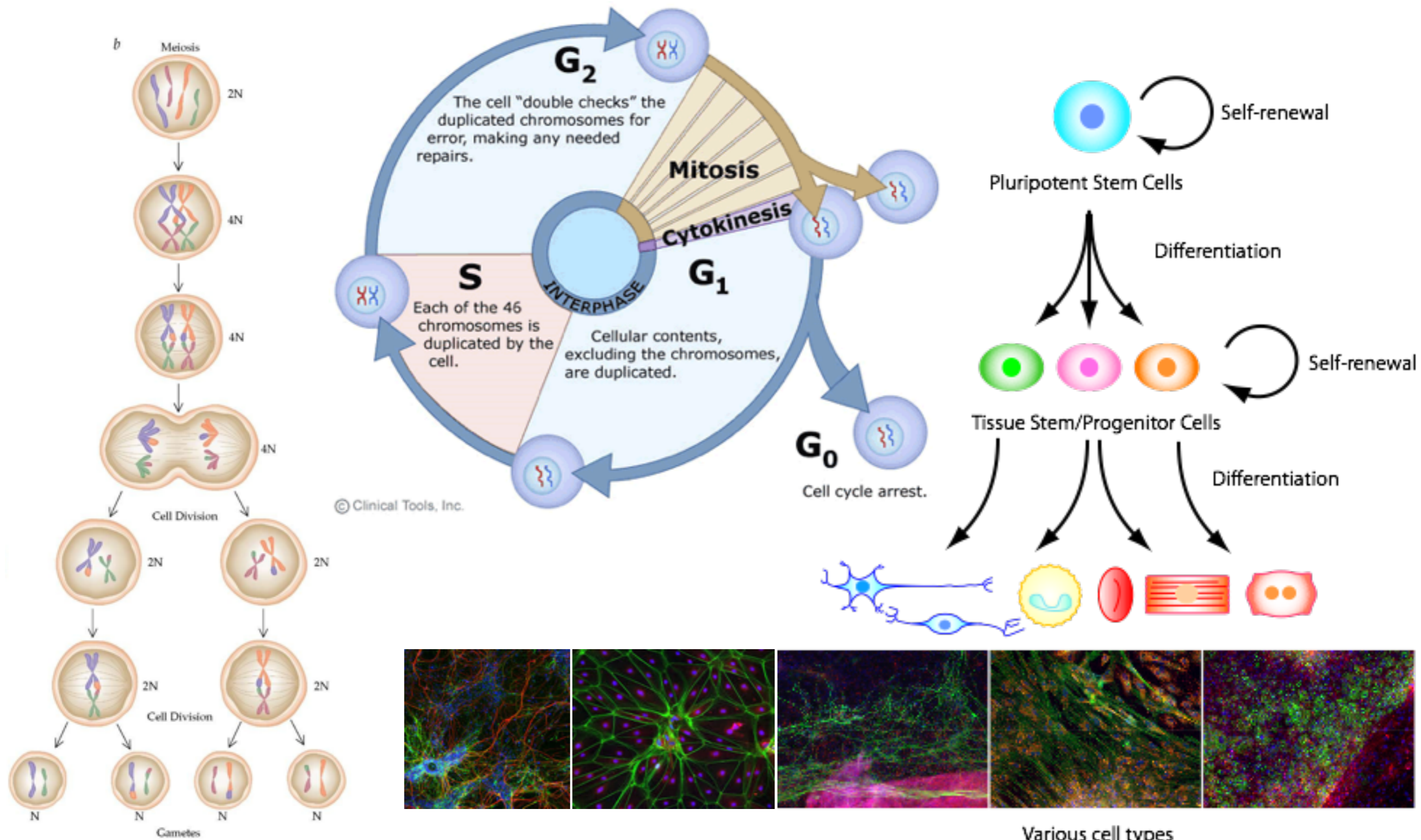


# Developmental Biology

How do animals grow and develop from a single cell?



# Developmental Biology



Various cell types

# Developmental Biology



We need single-cell resolution to:

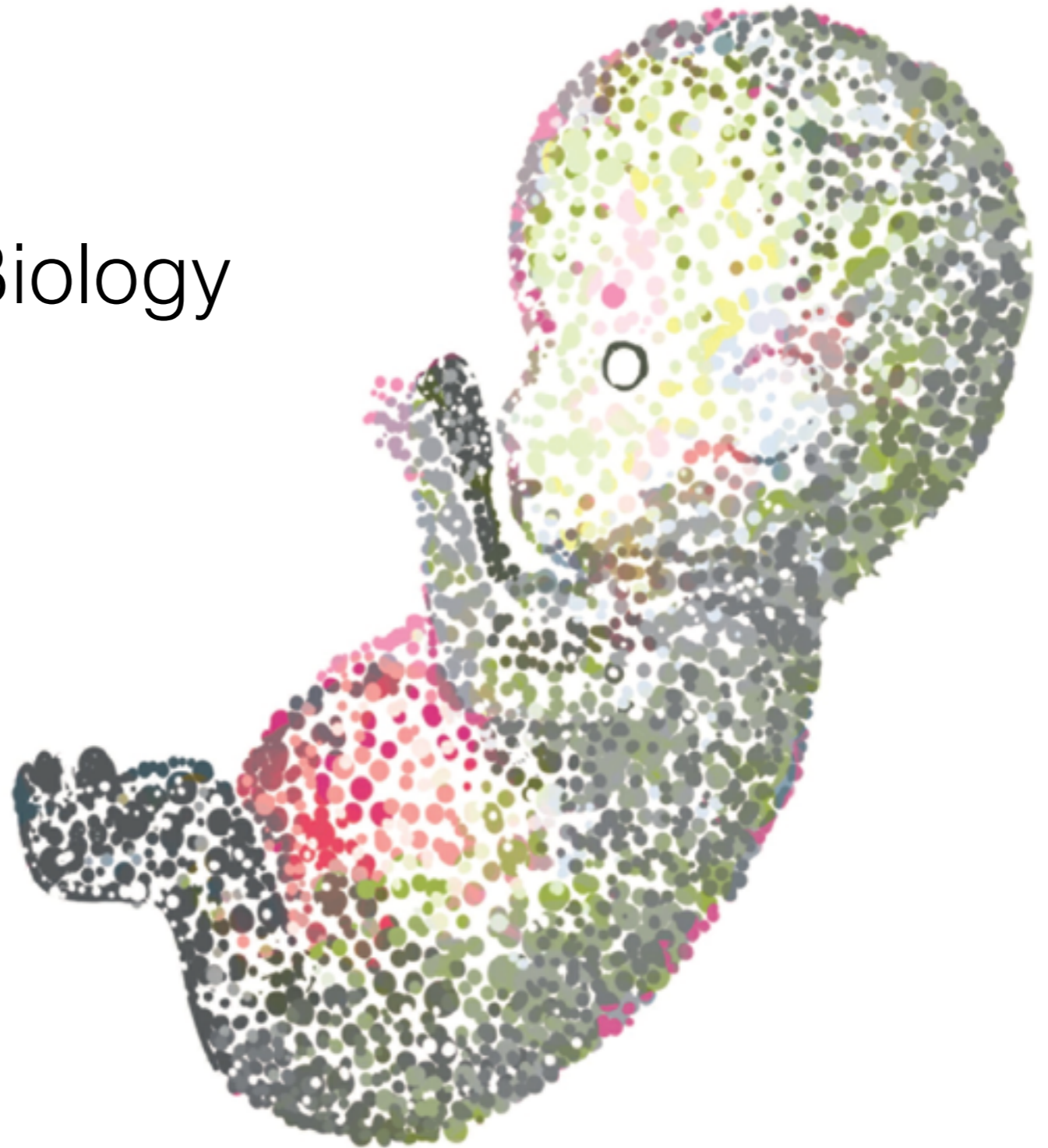
- Discover more complicated mechanisms in cellular development
- Confirm the distinct gene expression signatures across different cell types
- Identify functional differences among the same cell cell type



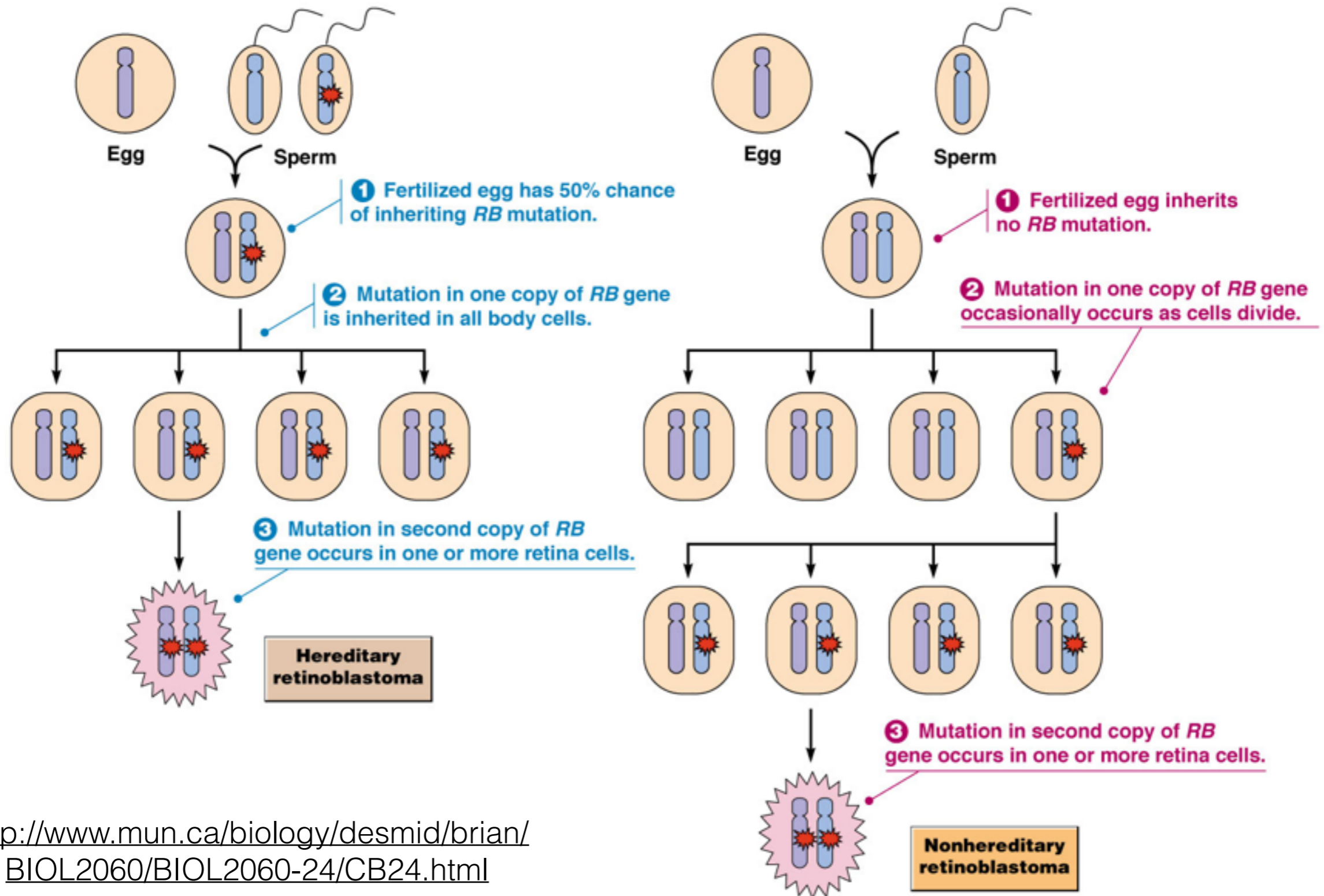
# Applications



- Developmental Biology
- **Cancer Biology**
- Microbiology
- Neurology



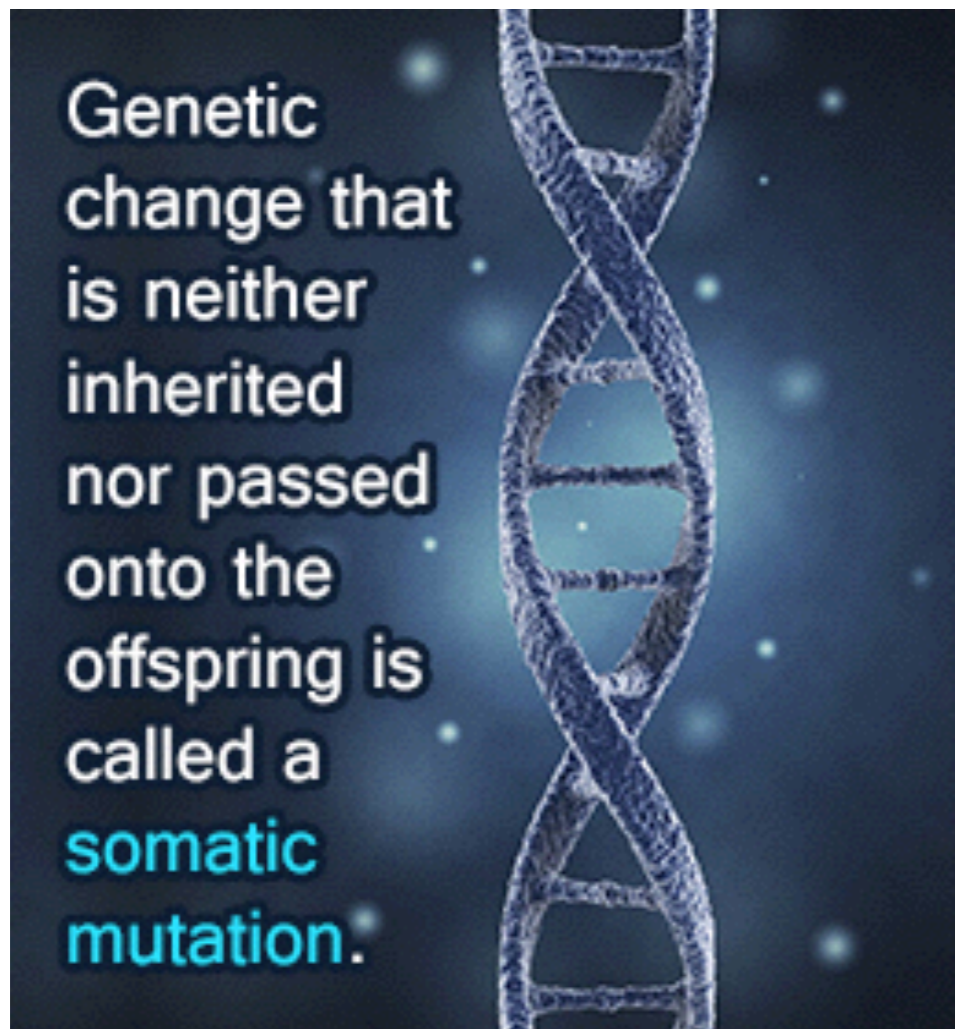
# Cancer Biology



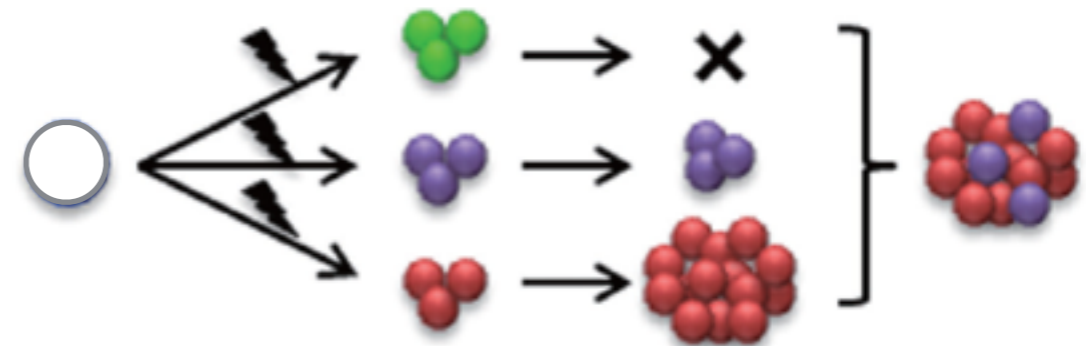
<http://www.mun.ca/biology/desmid/brian/BIOL2060/BIOL2060-24/CB24.html>

# Cancer Biology

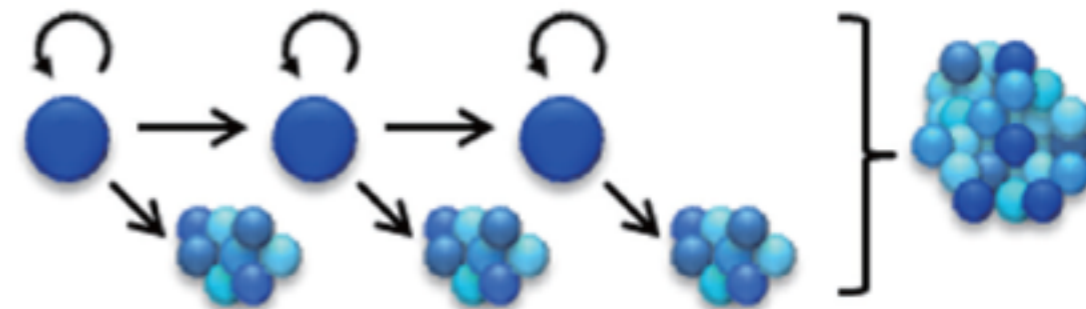
Tumors are composed of genetically and phenotypically **heterogeneous** clones



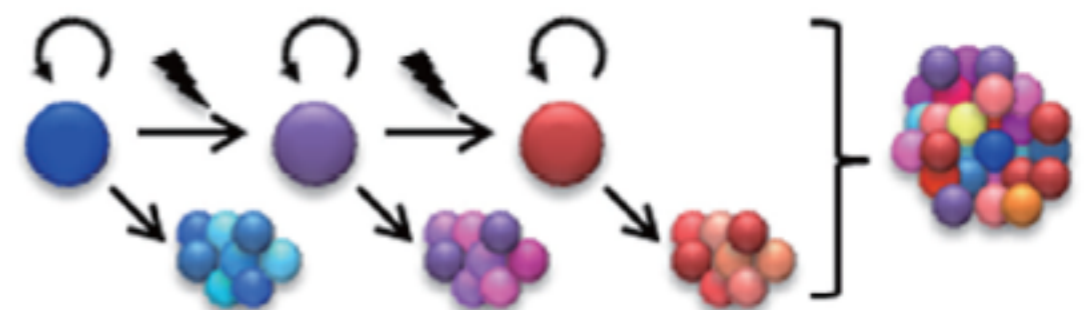
A Stochastic model



B Cancer stem cell model

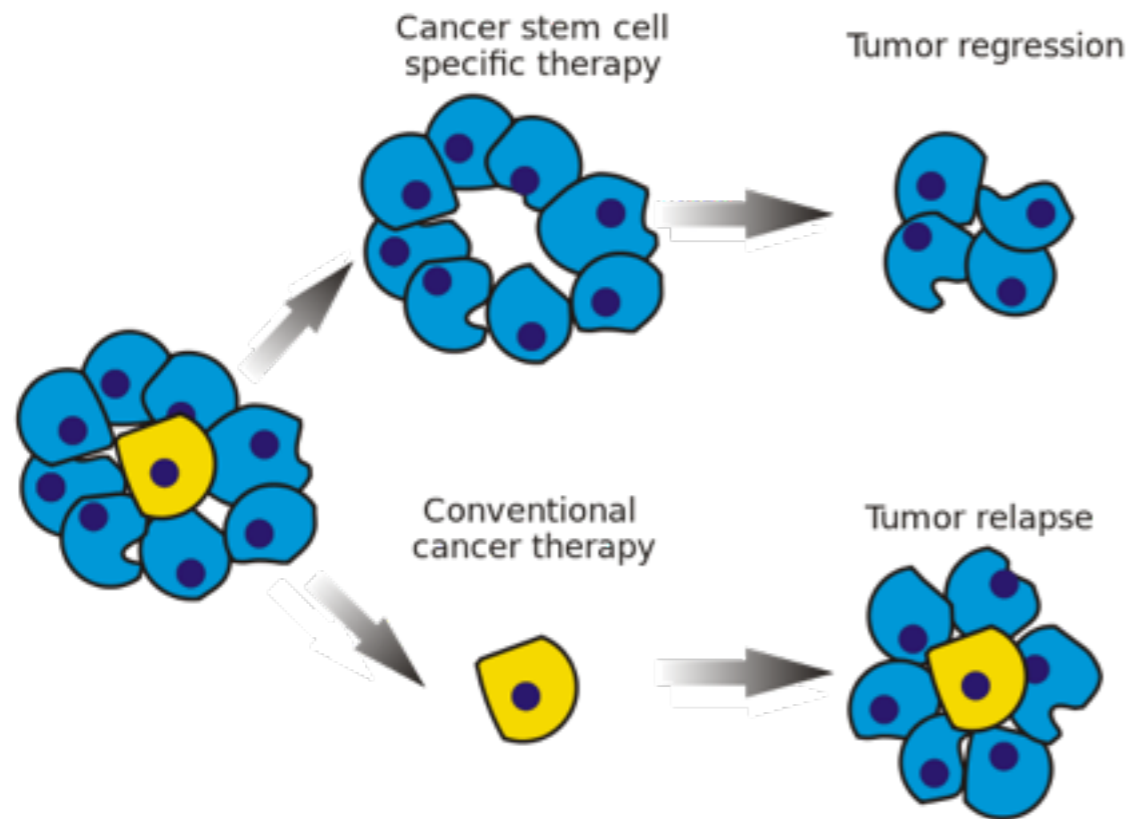


C Combination model

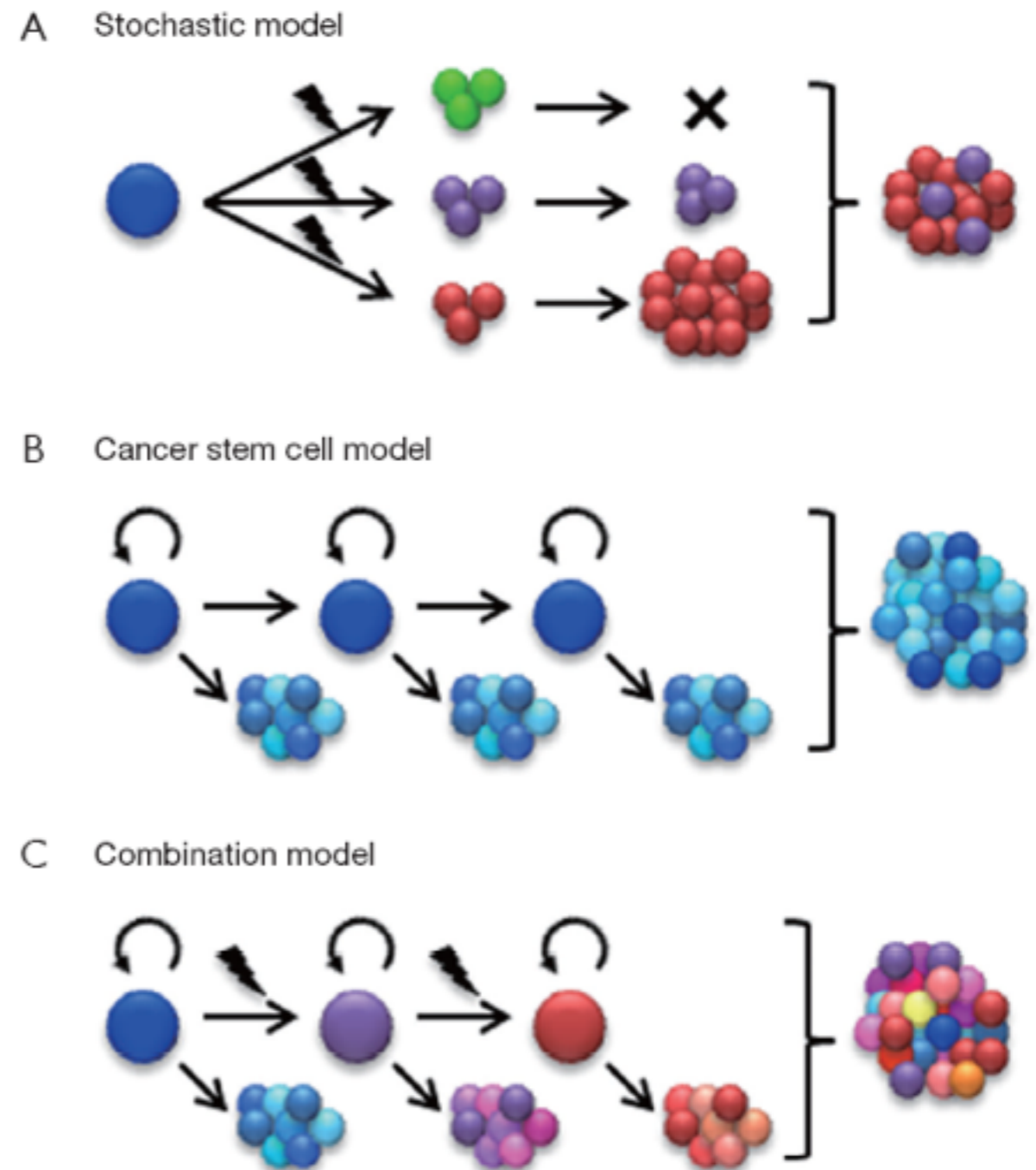


 Major genetic/epigenetic events

# Cancer Biology



Deep (bulk) sequencing can only capture 1% of the cell population (excluding some types such as circulating tumor cells).



⚡ Major genetic/epigenetic events

# Cancer Biology



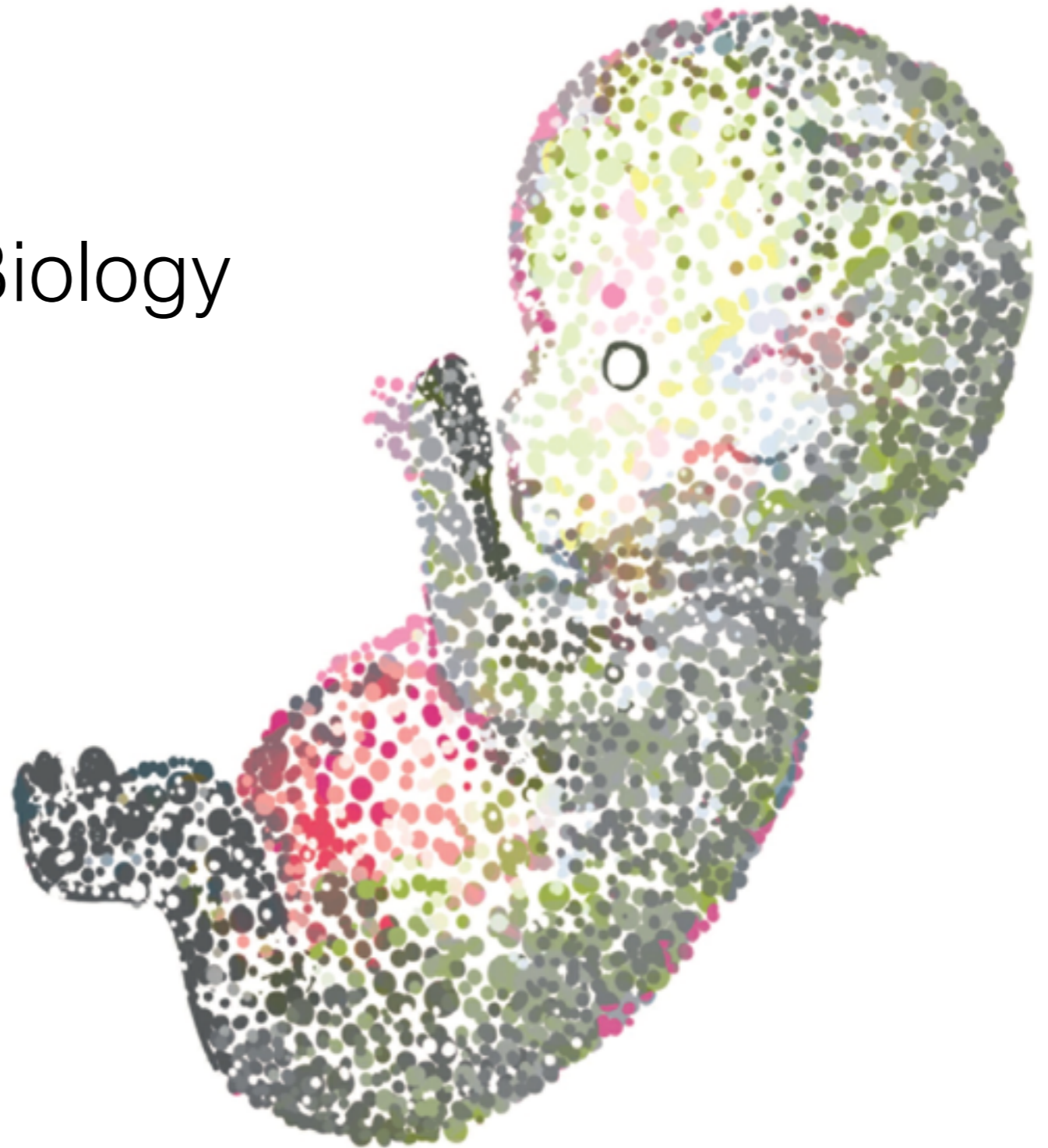
We need single-cell resolution to:

- Find evidence for models of cancer
- Infer timing of mutations and the drivers
- Evaluate effectiveness of targeted therapy

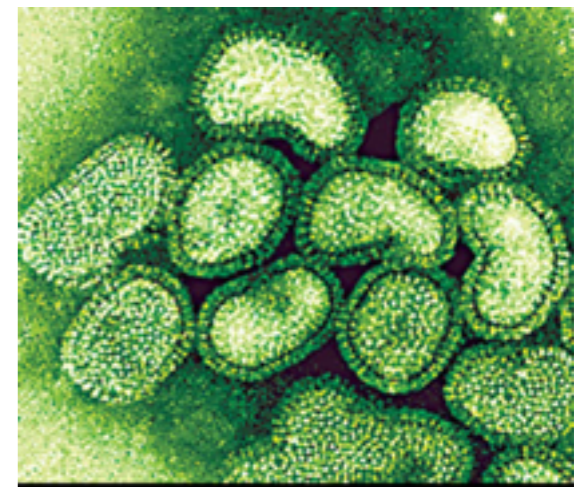
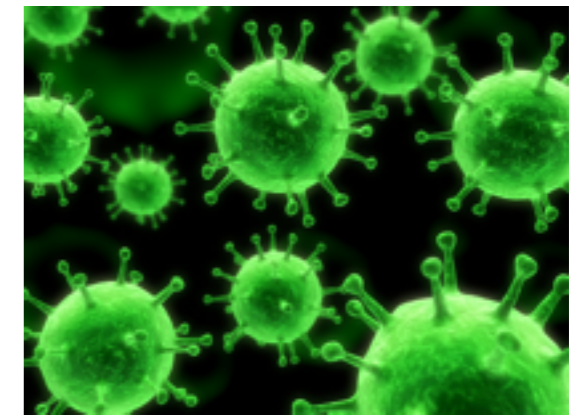
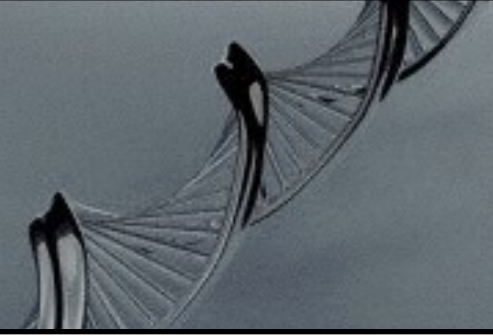
# Applications



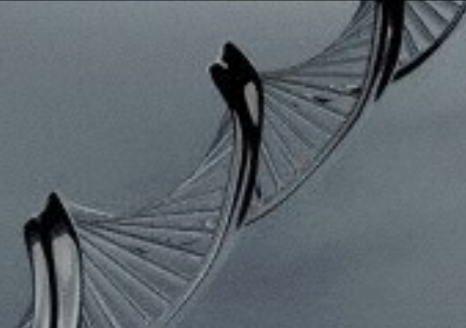
- Developmental Biology
- Cancer Biology
- **Microbiology**
- Neurology



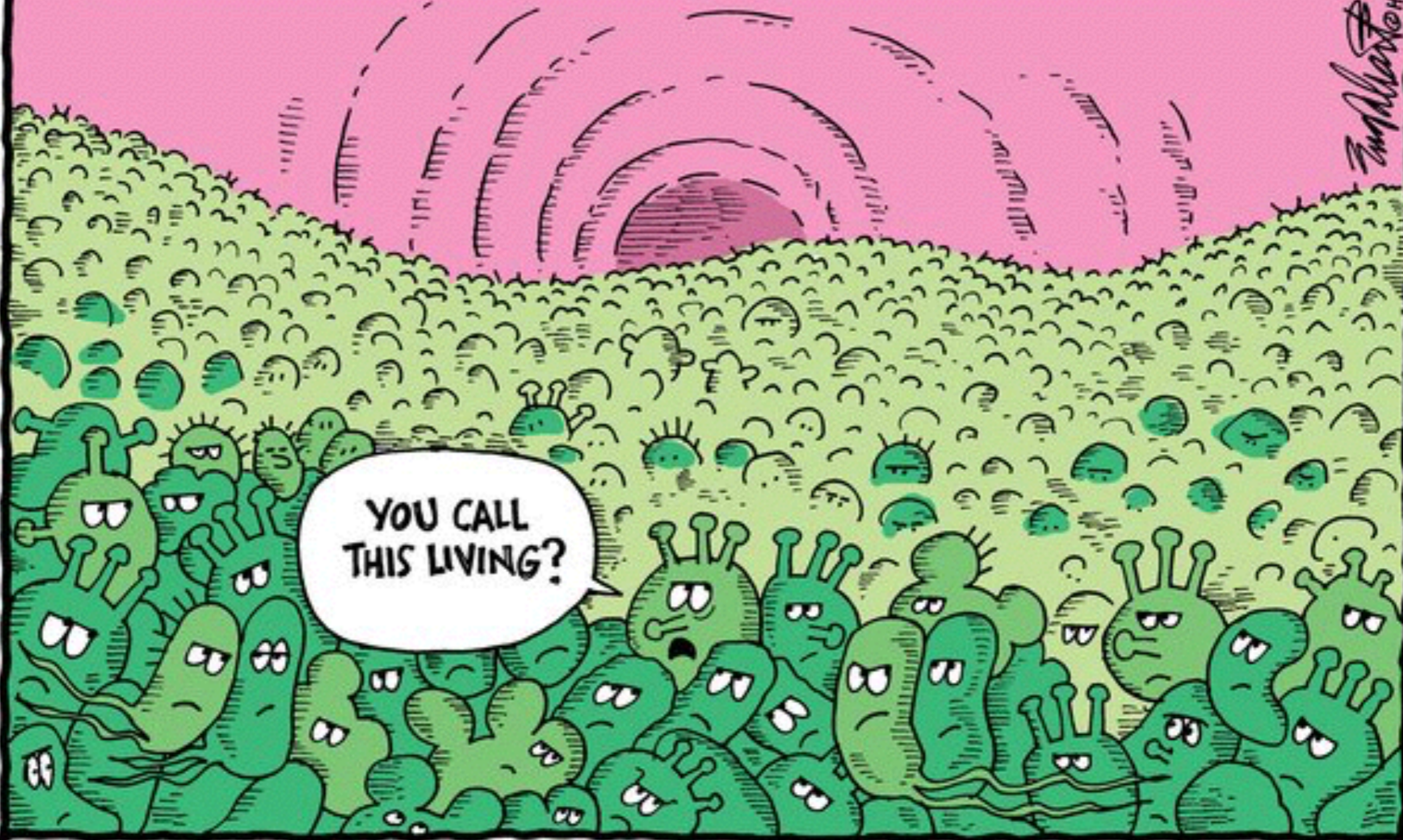
# Microbiology



# Microbiology



THE HUMAN MICROBIOME PROJECT SAYS THE HUMAN BODY HAS 100 TRILLION MICROSCOPIC LIFE FORMS LIVING IN IT.

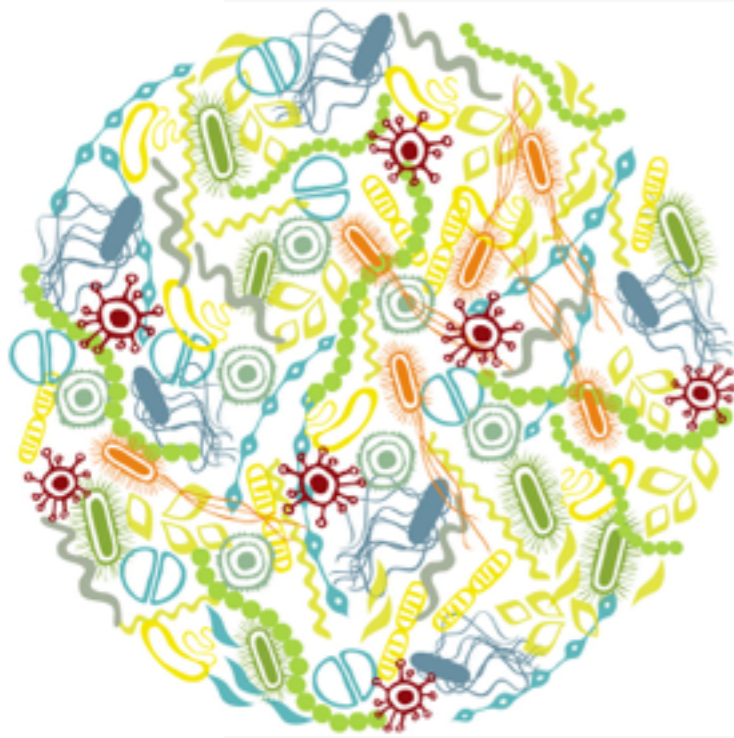
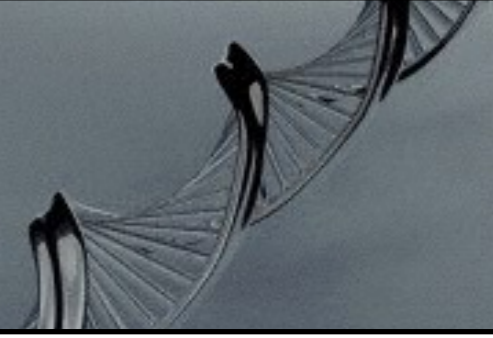


YOU CALL THIS LIVING?

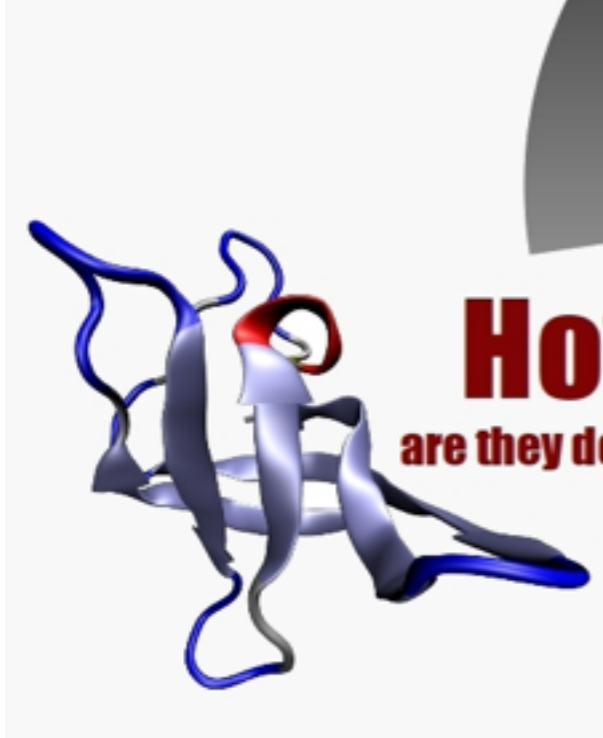
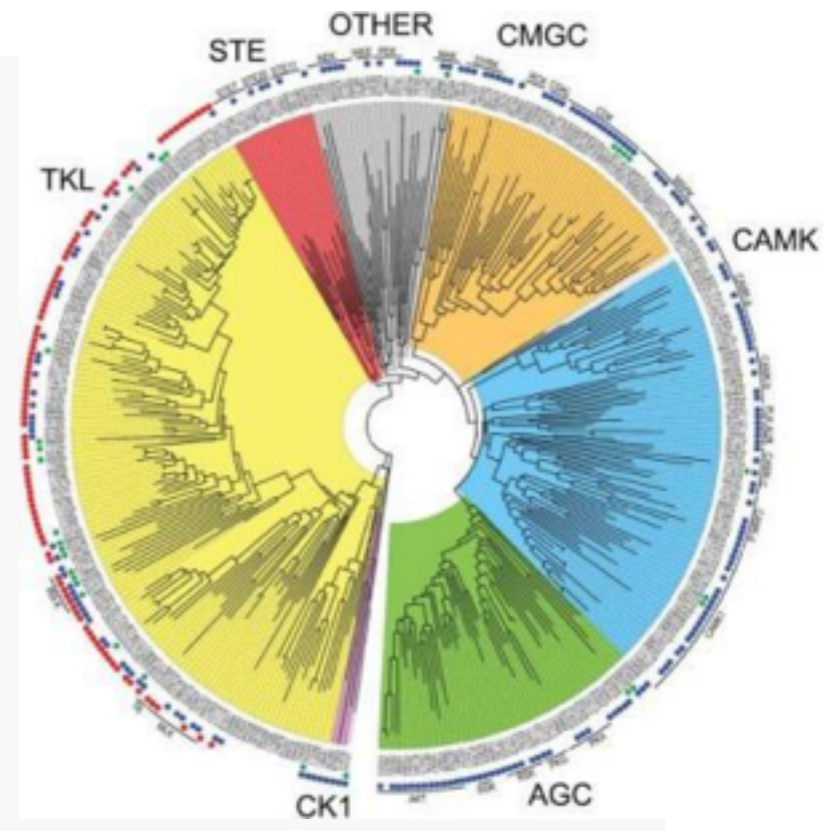
6/15/12  
HARTFORD CONNECTICUT  
GOURANE



# Microbiology



**Who**  
is in there?



**How**  
are they doing it?



**What**  
are they doing?

# Microbiology



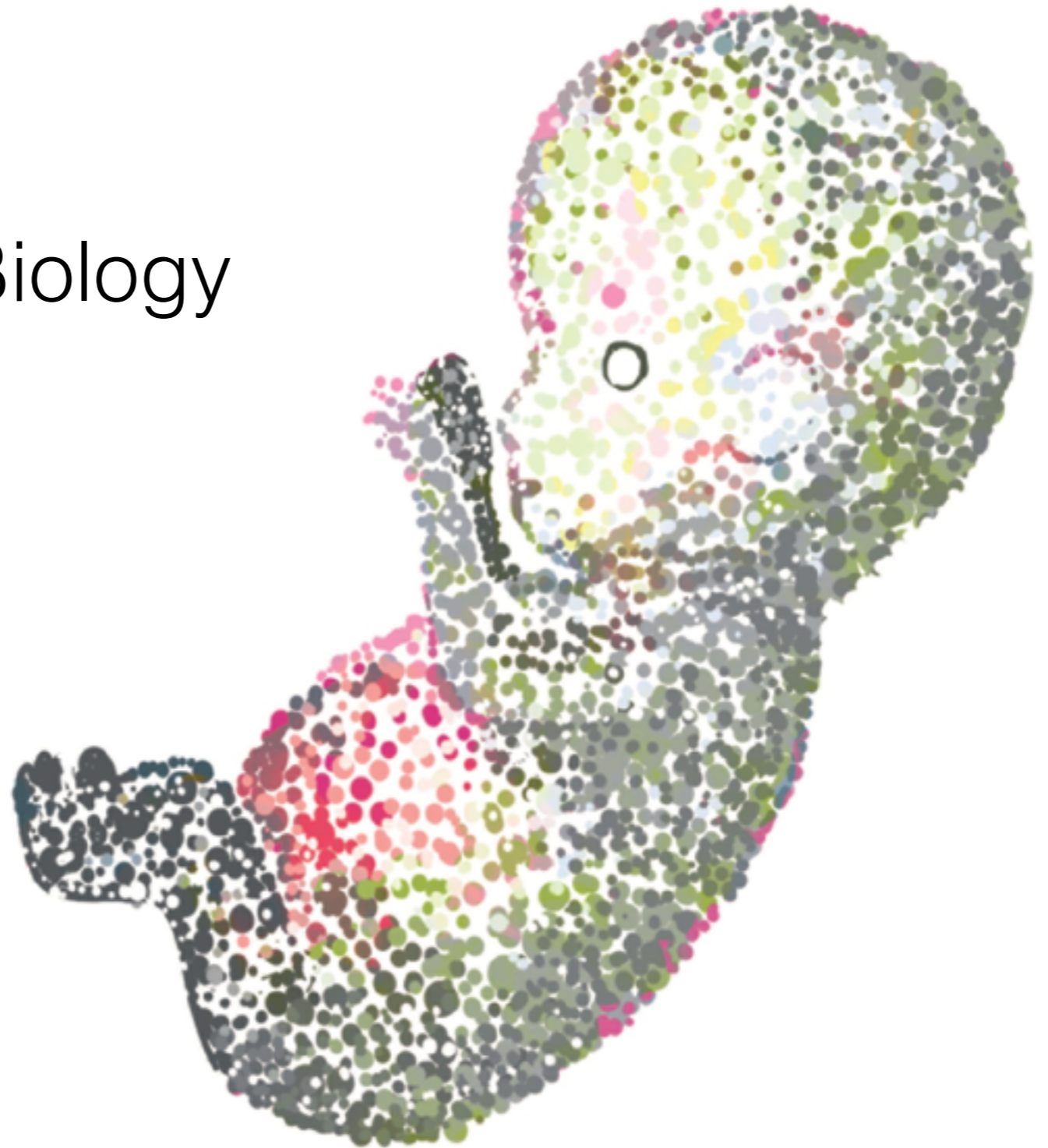
We need single-cell resolution to:

- Discover low-abundance species that are difficult to culture in vitro
- Monitor transcriptional gene activation mechanisms for functional annotation

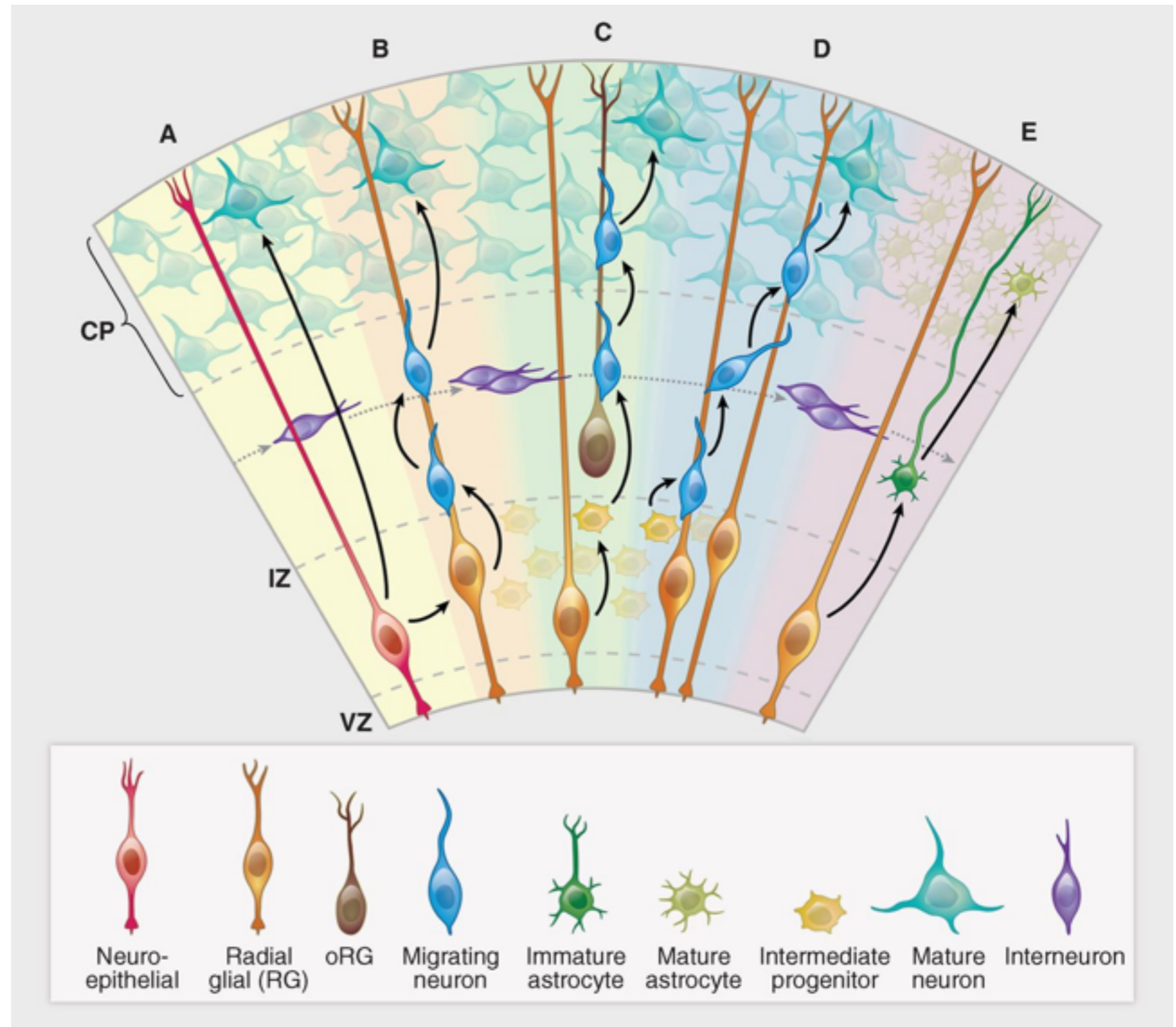
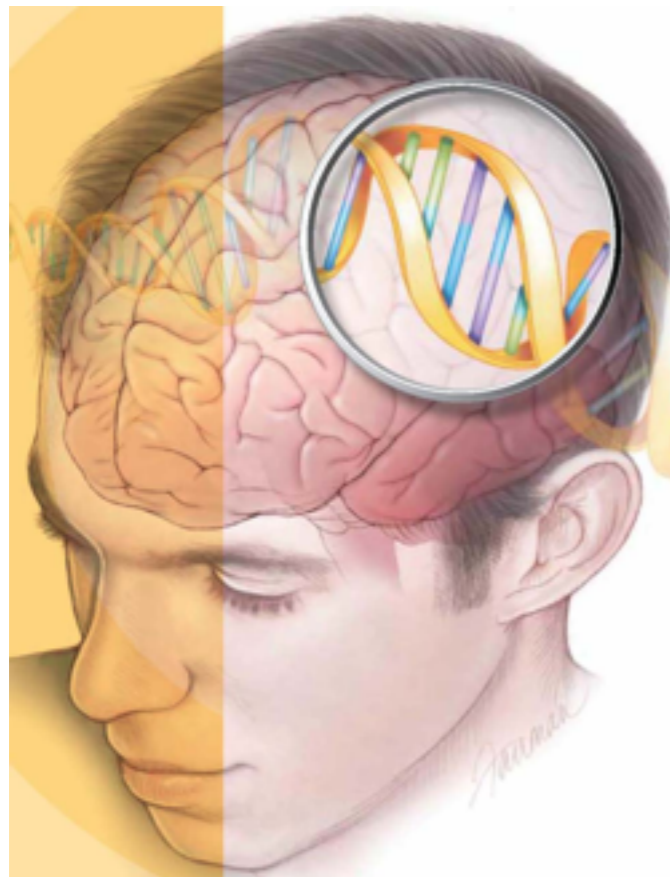
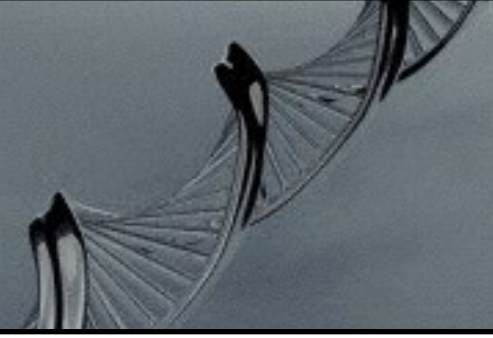
# Applications



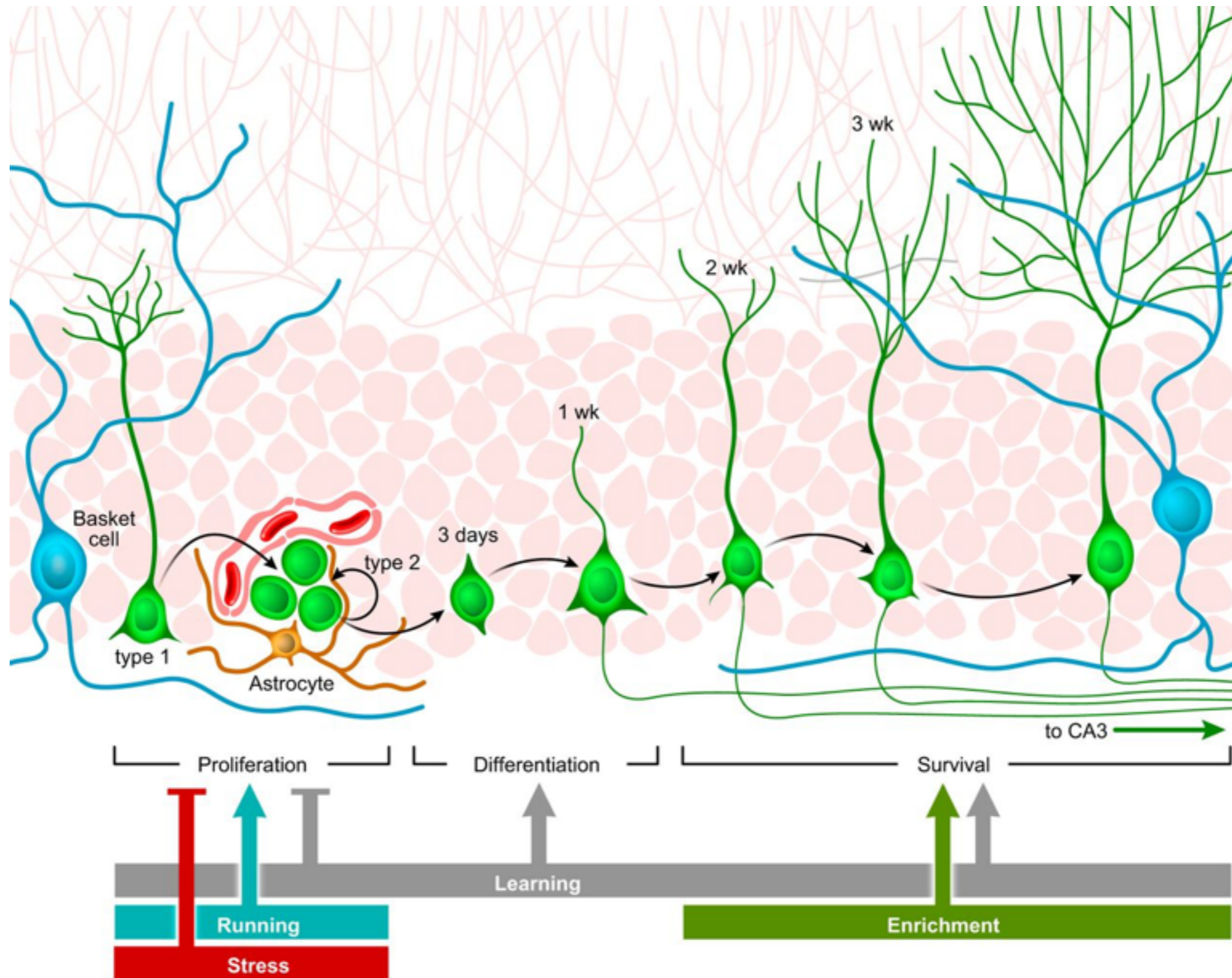
- Developmental Biology
- Cancer Biology
- Microbiology
- **Neurology**



# Neurology



# Neurology



# Neurology



We need single-cell resolution to:

- Study the mosaic genomes of individual neurons and compositions in the brain
- Follow genetic variations during fetal development
- Develop targeted therapy for neurological diseases for specific cell types

# SINGLE CELL

NOW THAT'S IRONIC!

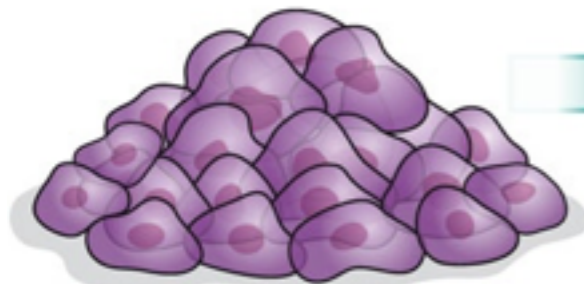
26

# Traditional v.s. Single-cell

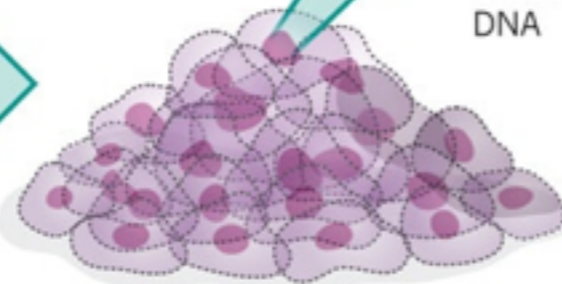
## ONE GENOME FROM MANY

Sequencing the genomes of single cells is similar to sequencing those from multiple cells — but errors are more likely.

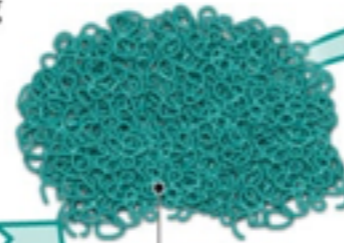
### ► Standard genome sequencing



A sample containing thousands to millions of cells is isolated.



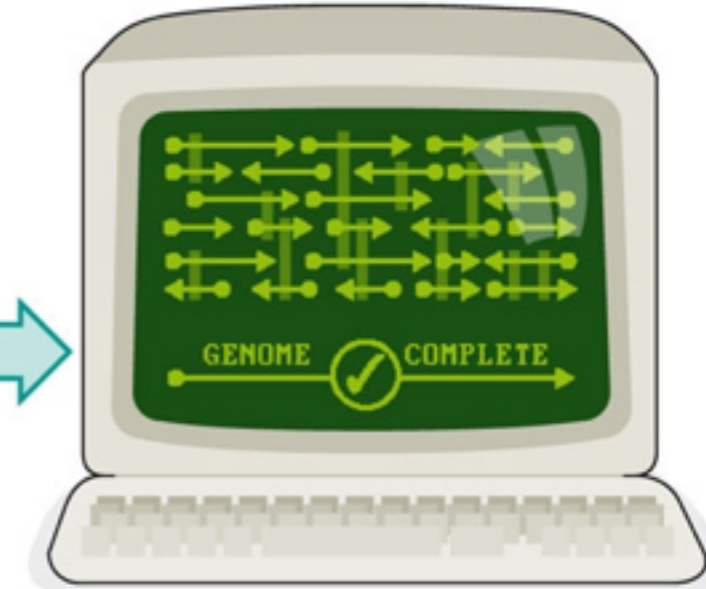
DNA is extracted from all the nuclei.



Loads of DNA

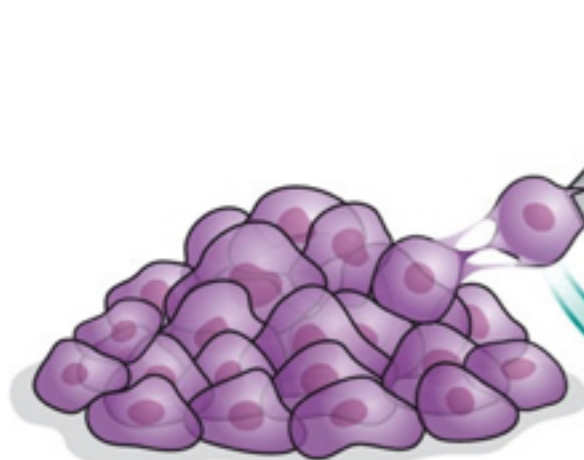


DNA is broken into fragments and then sequenced.

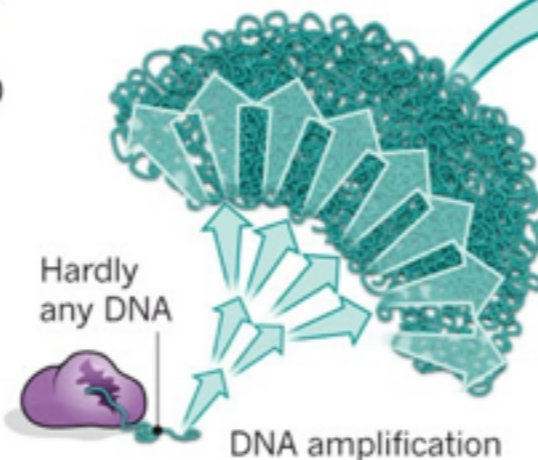


The sequences are assembled to give a common, 'consensus' sequence.

### ► Single-cell sequencing



A single cell is difficult to isolate, but it can be done mechanically or with an automated cell sorter.



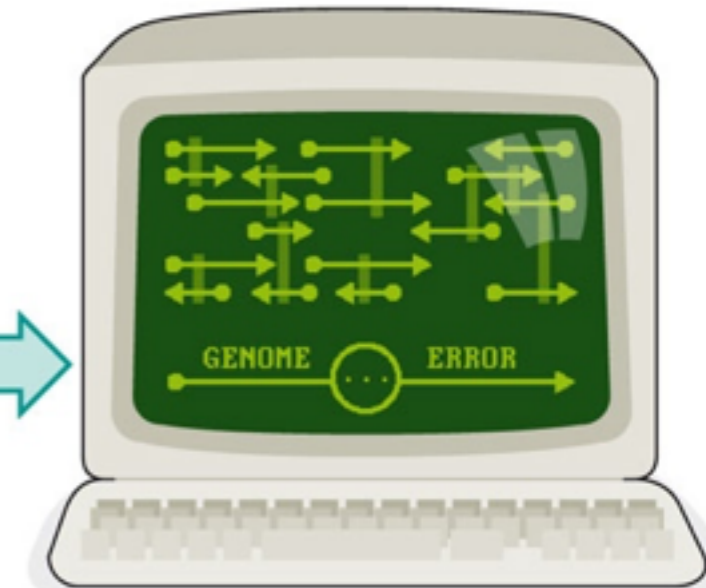
Hardly any DNA

DNA amplification

The DNA is extracted and amplified, during which errors can creep in.



Amplified DNA is sequenced.

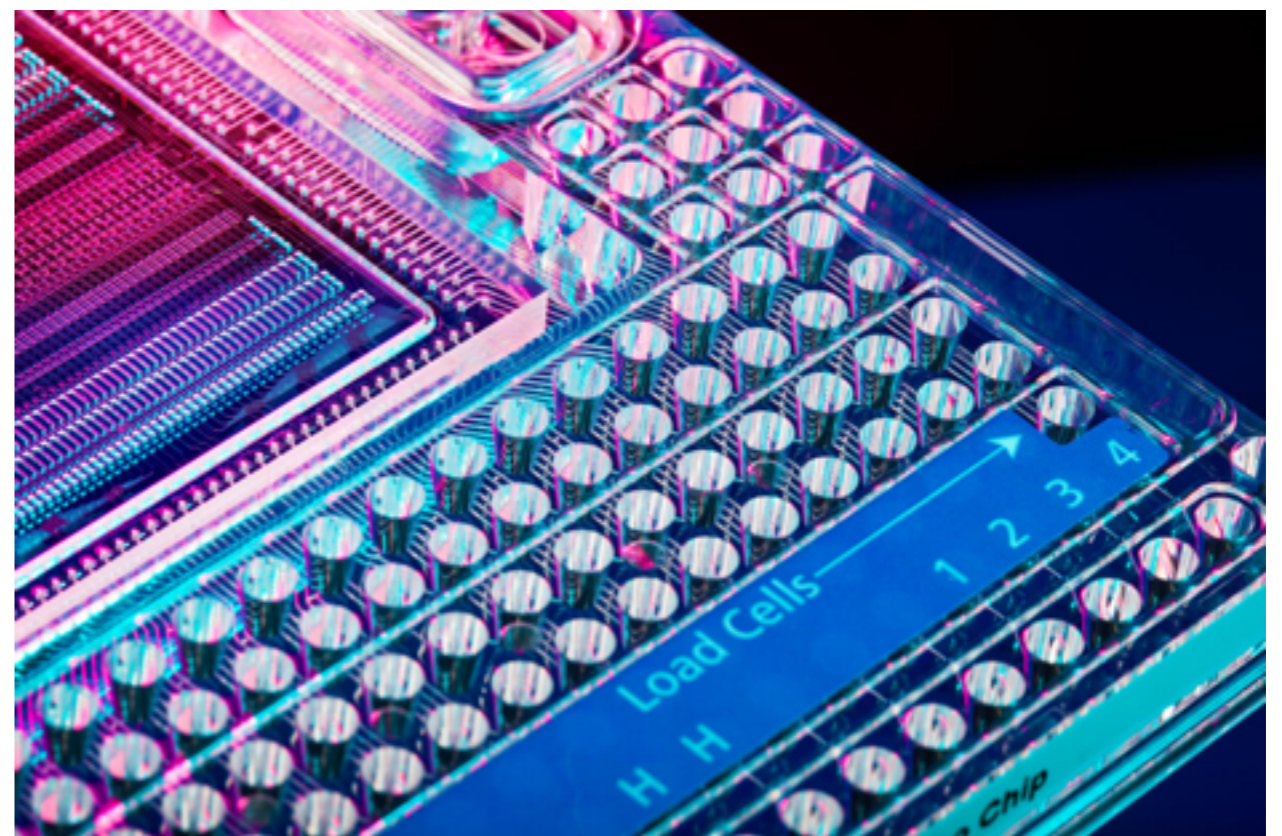
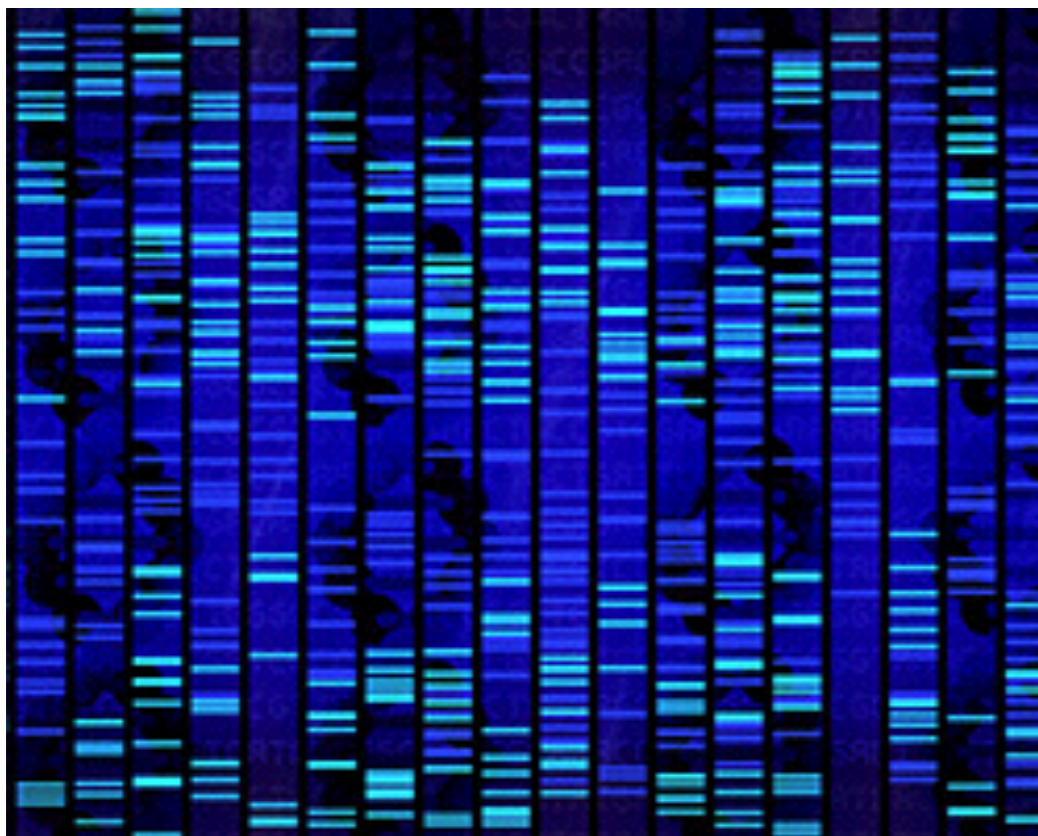
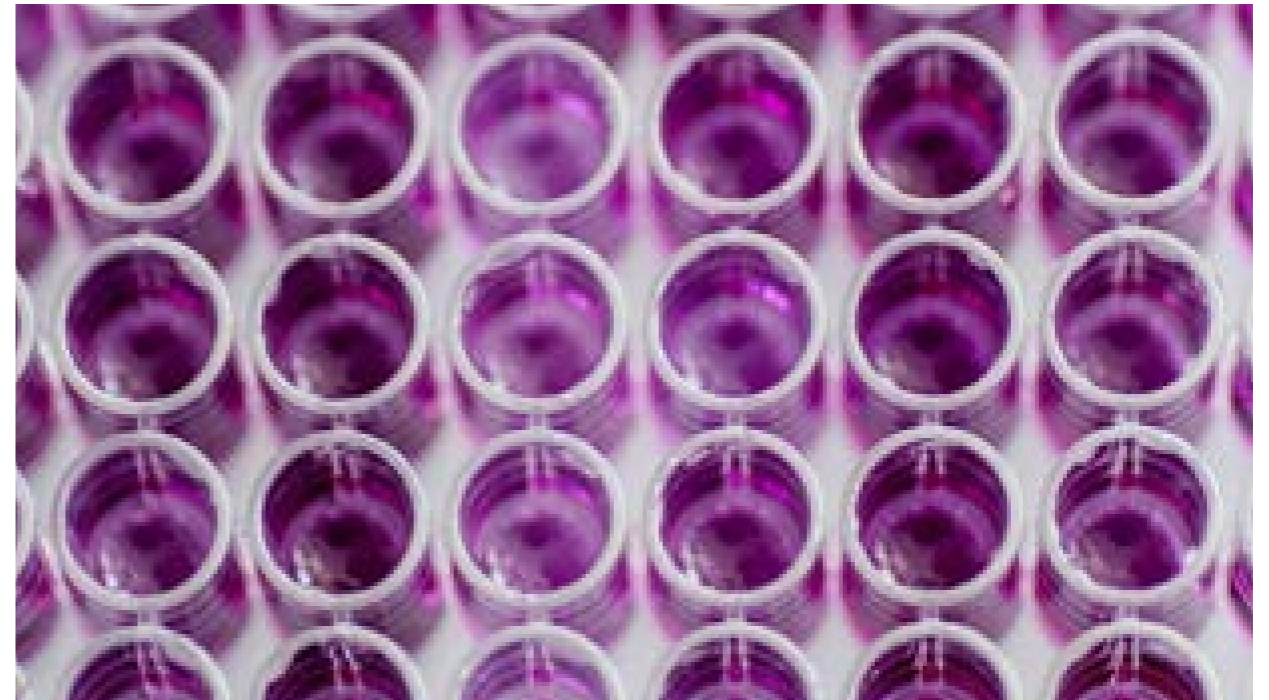


Errors introduced in earlier steps make sequence assembly difficult; the final sequence can have gaps.

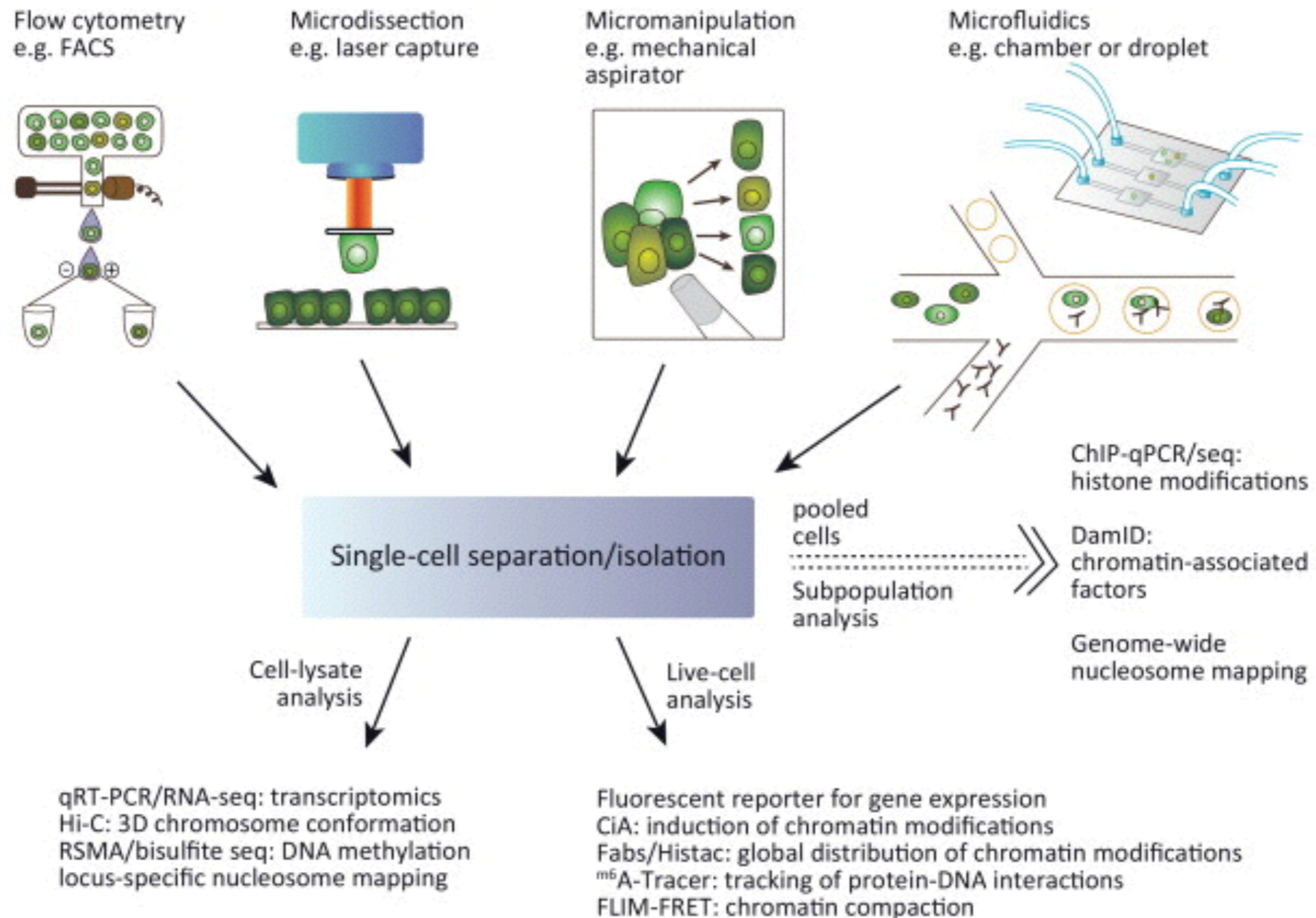


# Single-Cell Technologies

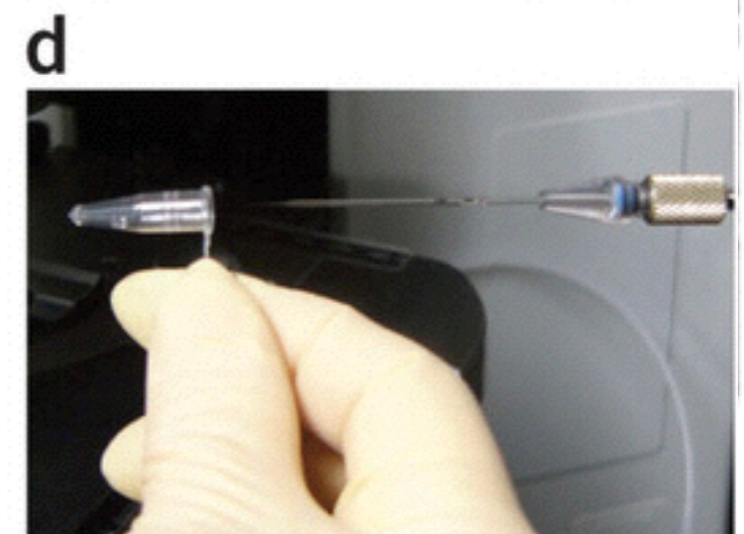
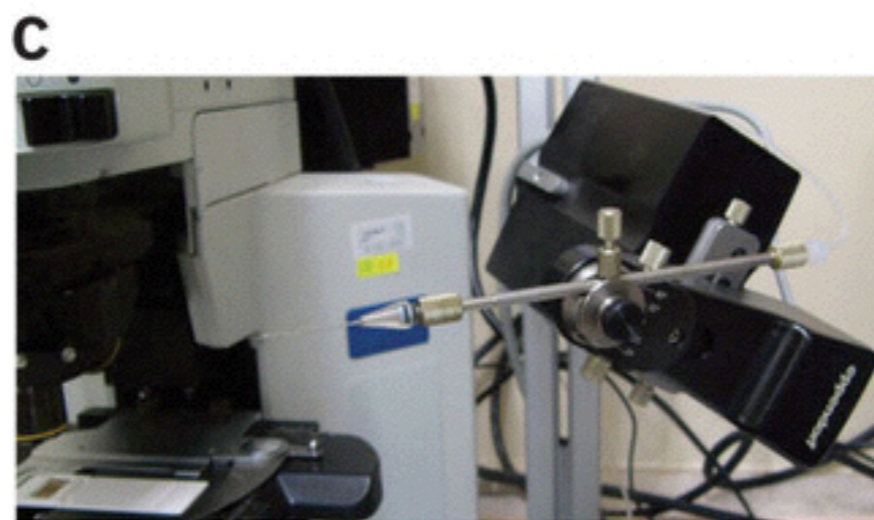
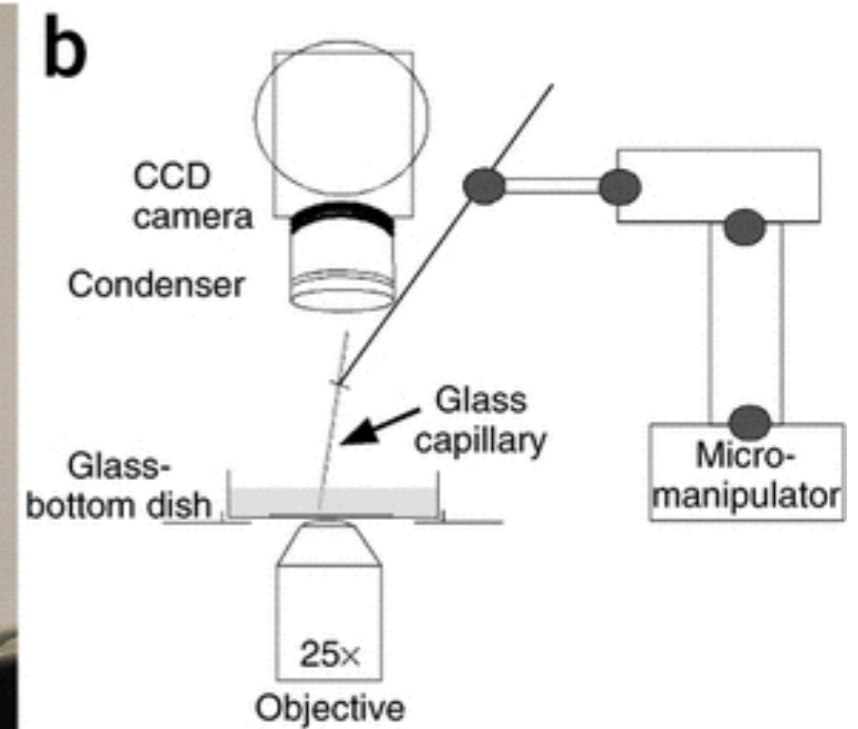
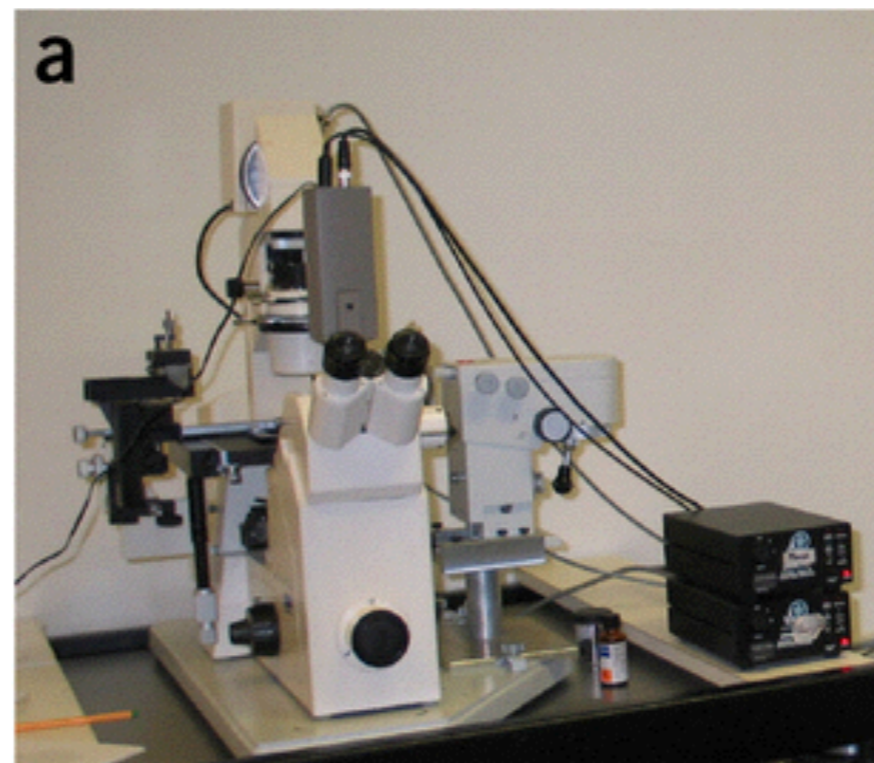
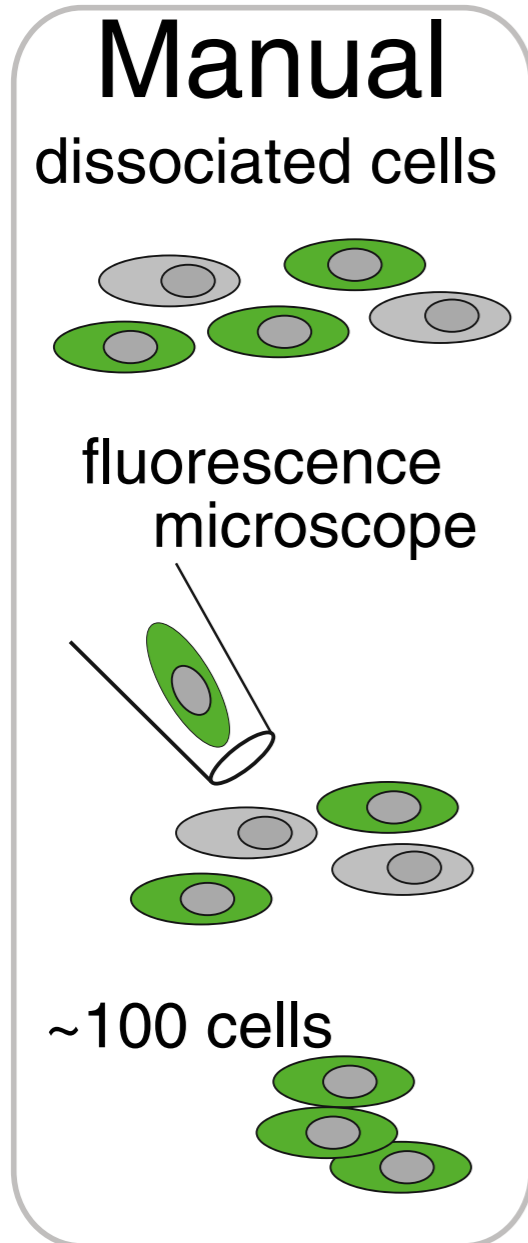
- (i) isolate single cells
- (ii) amplify genome efficiently
- (iii) sequence DNA



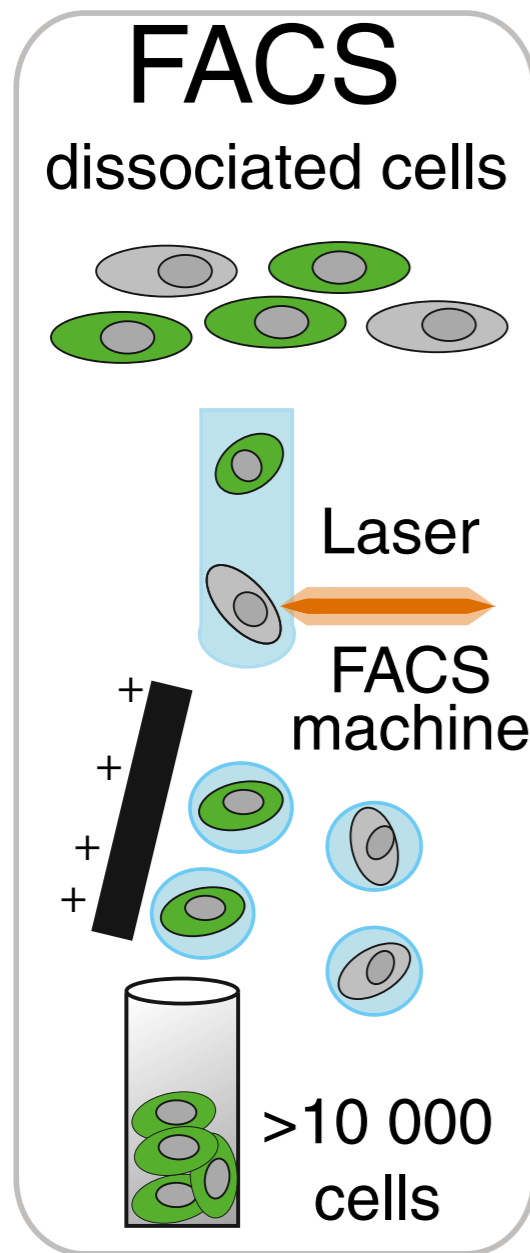
# Single-Cell Technologies



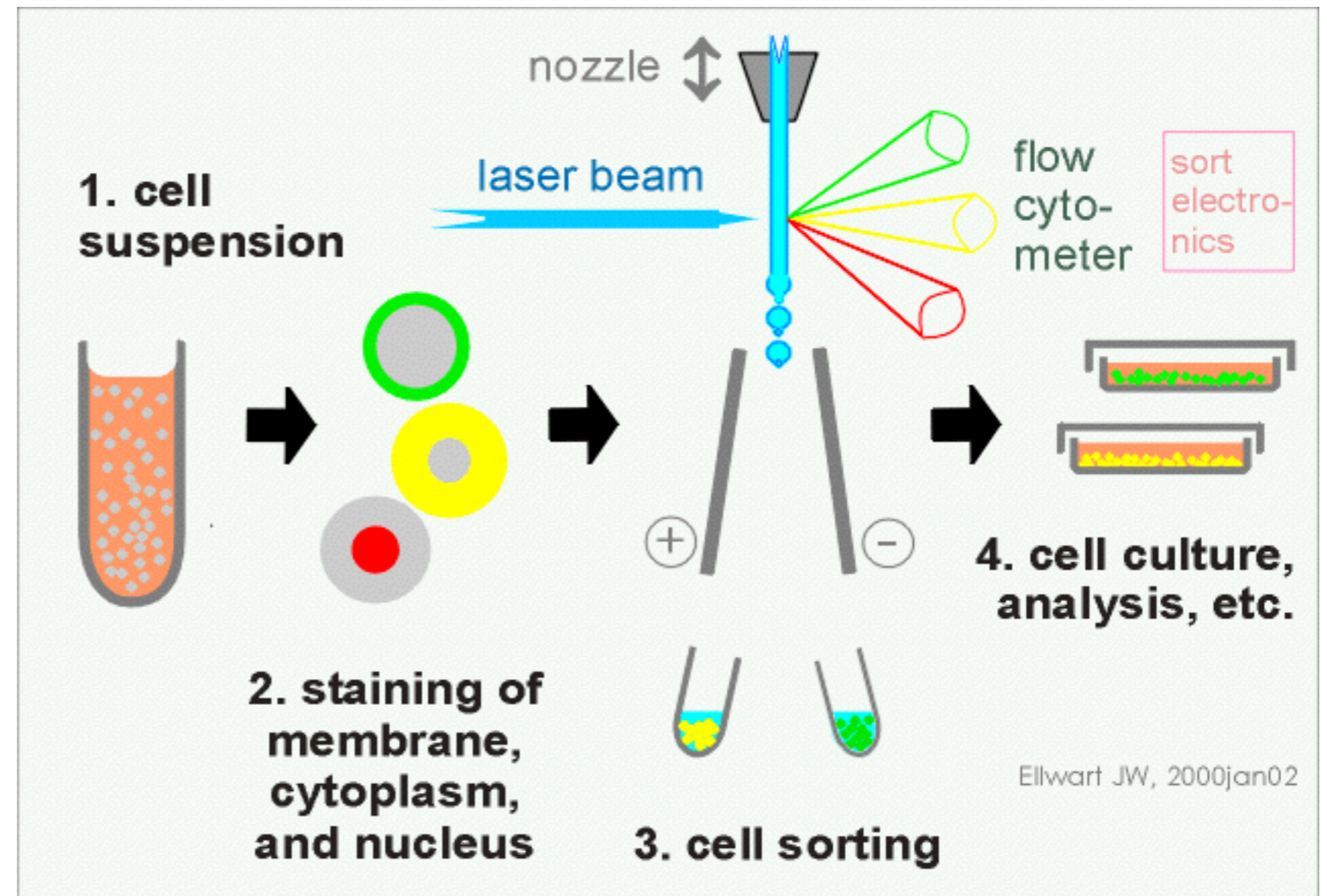
# Cell Sorting



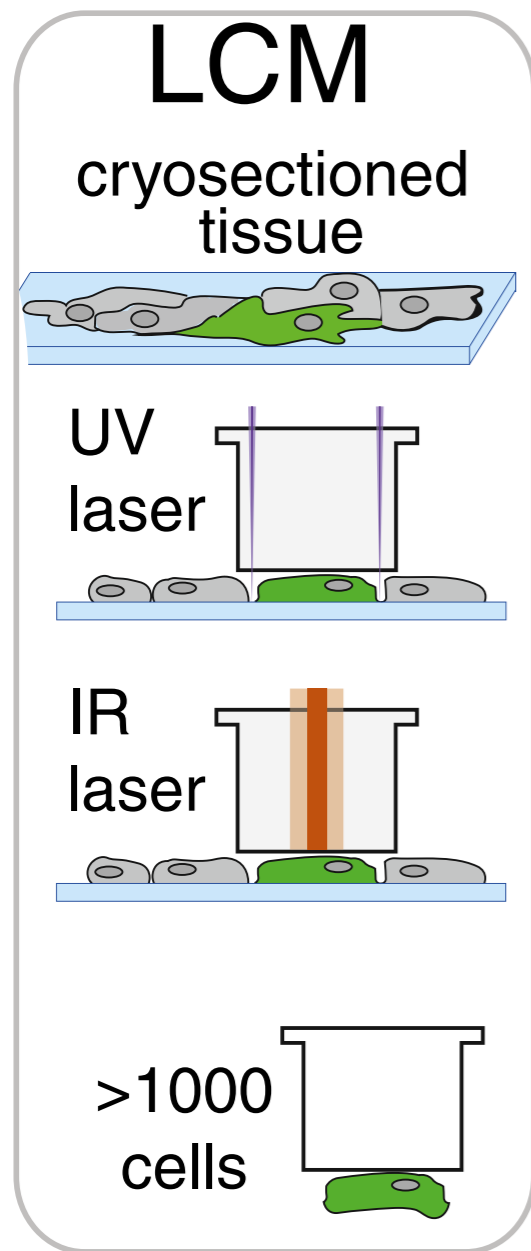
# Cell Sorting



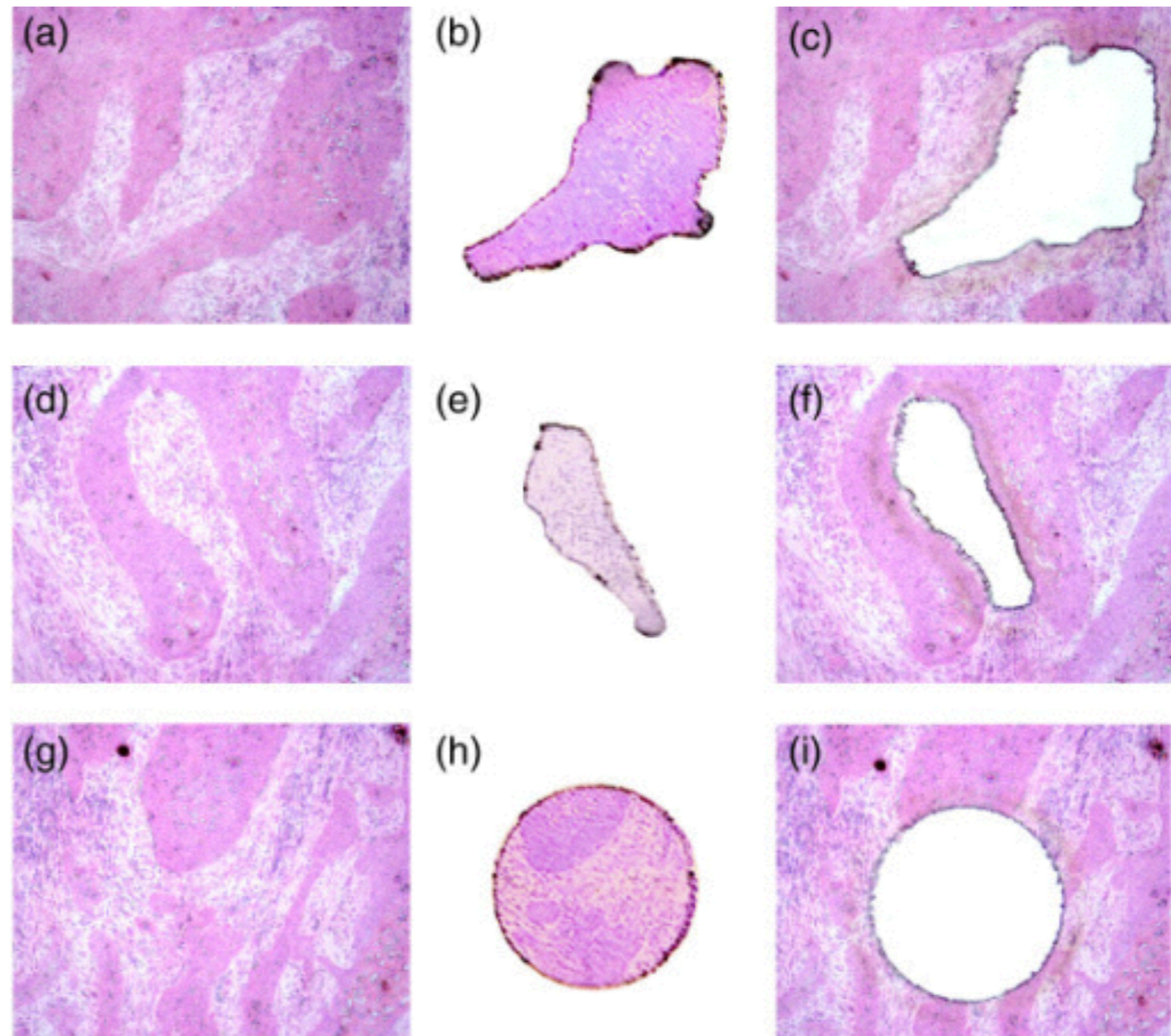
FACS: fluorescence activated cell sorting



# Cell Sorting



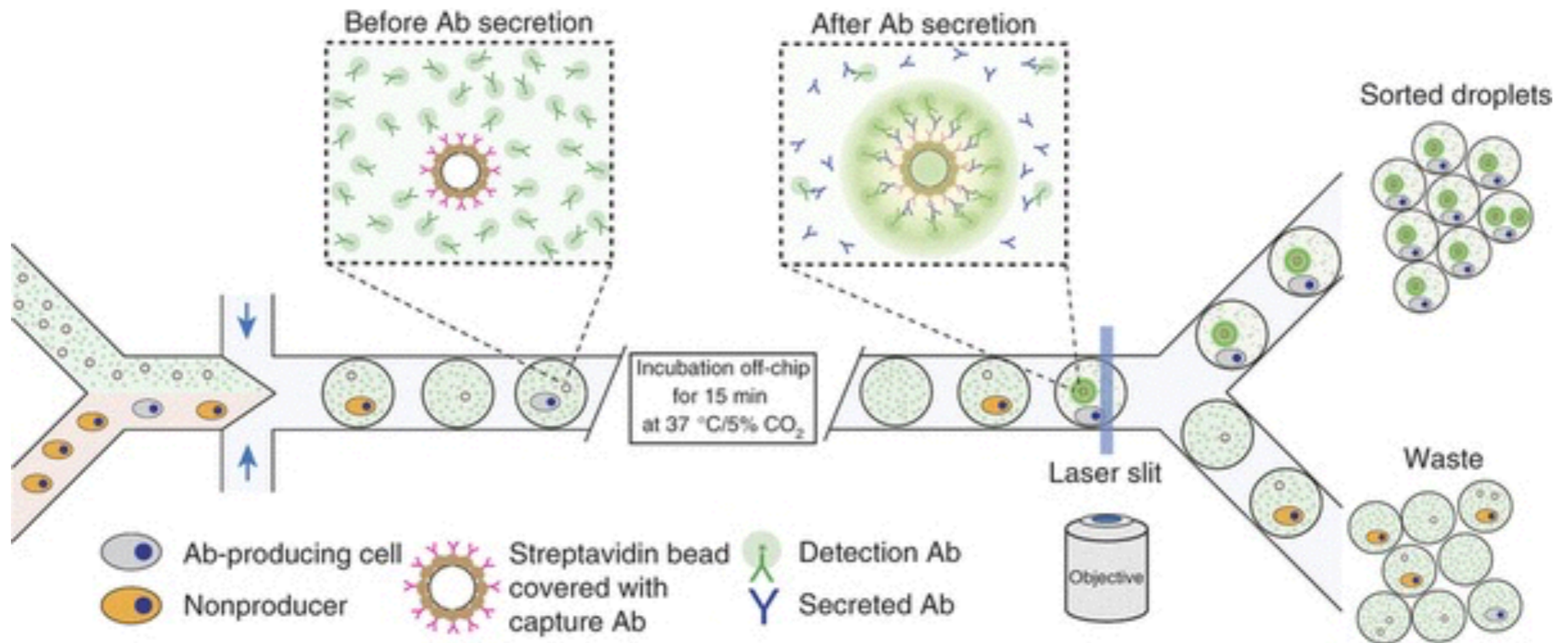
LCM: laser capture microdissection



# Cell Sorting

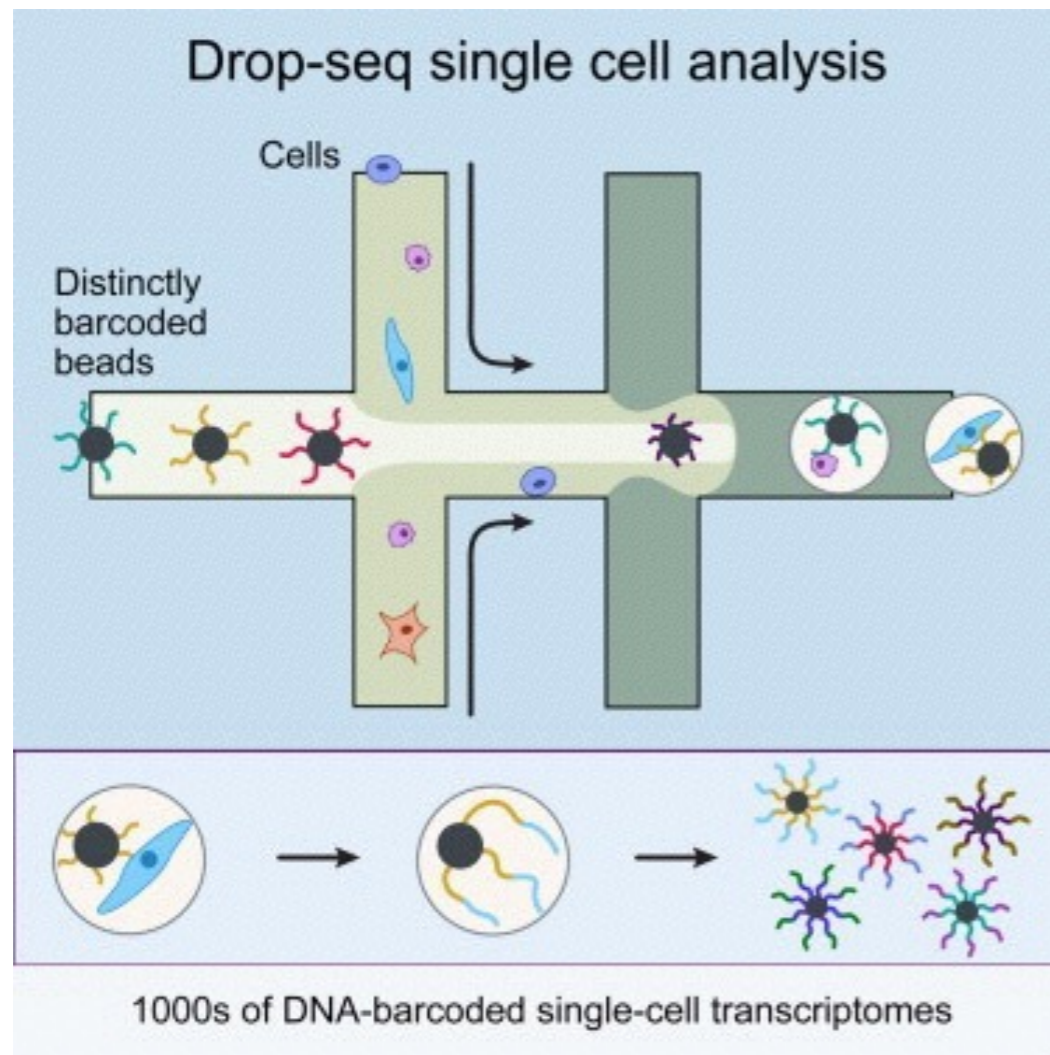


Microfluidics: can isolate rare circulating cells

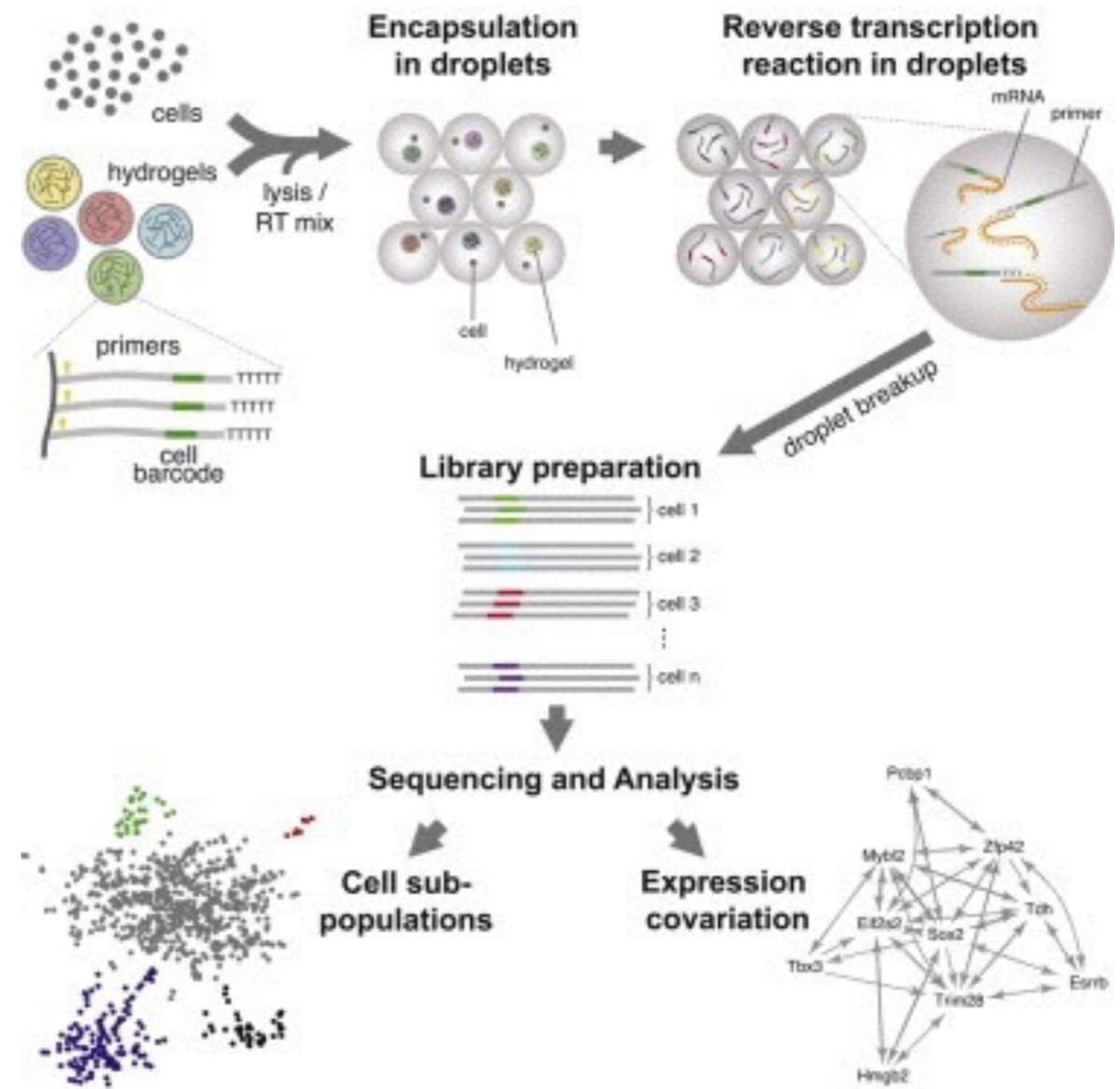


# Cell Sorting

High-throughput (~100,000 cells)



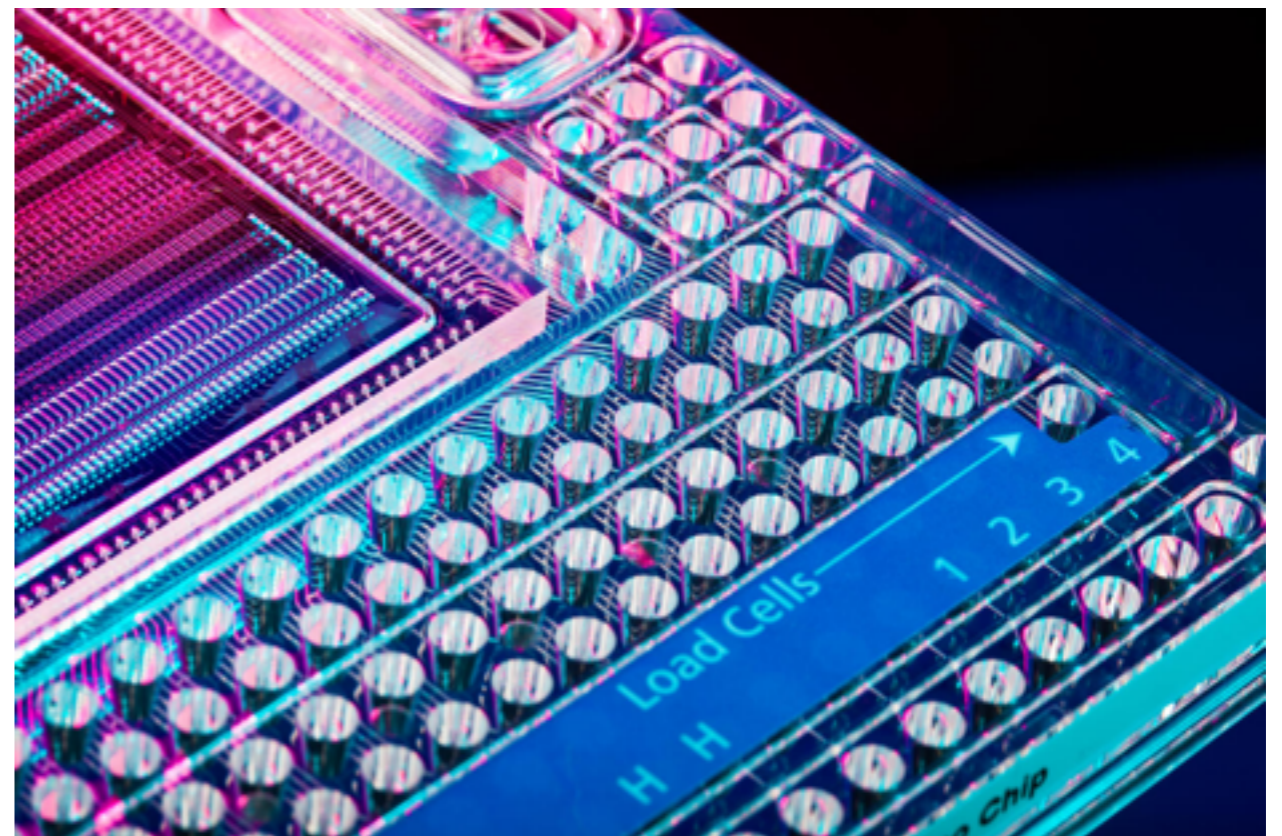
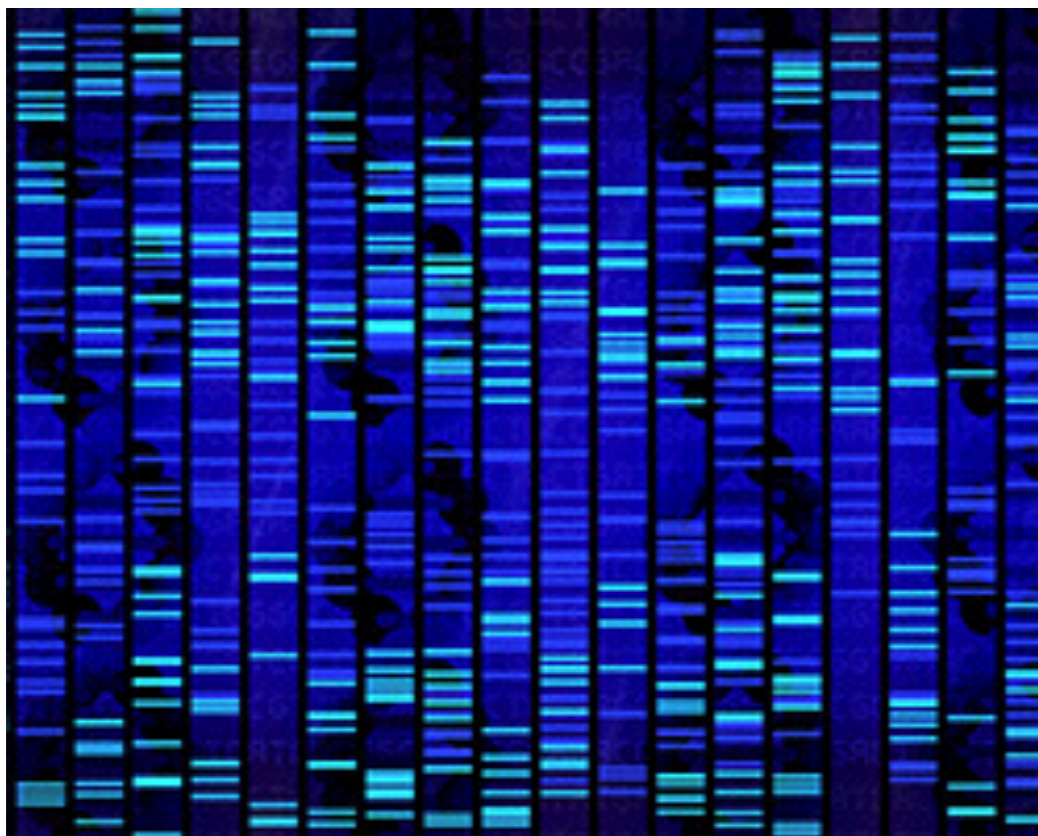
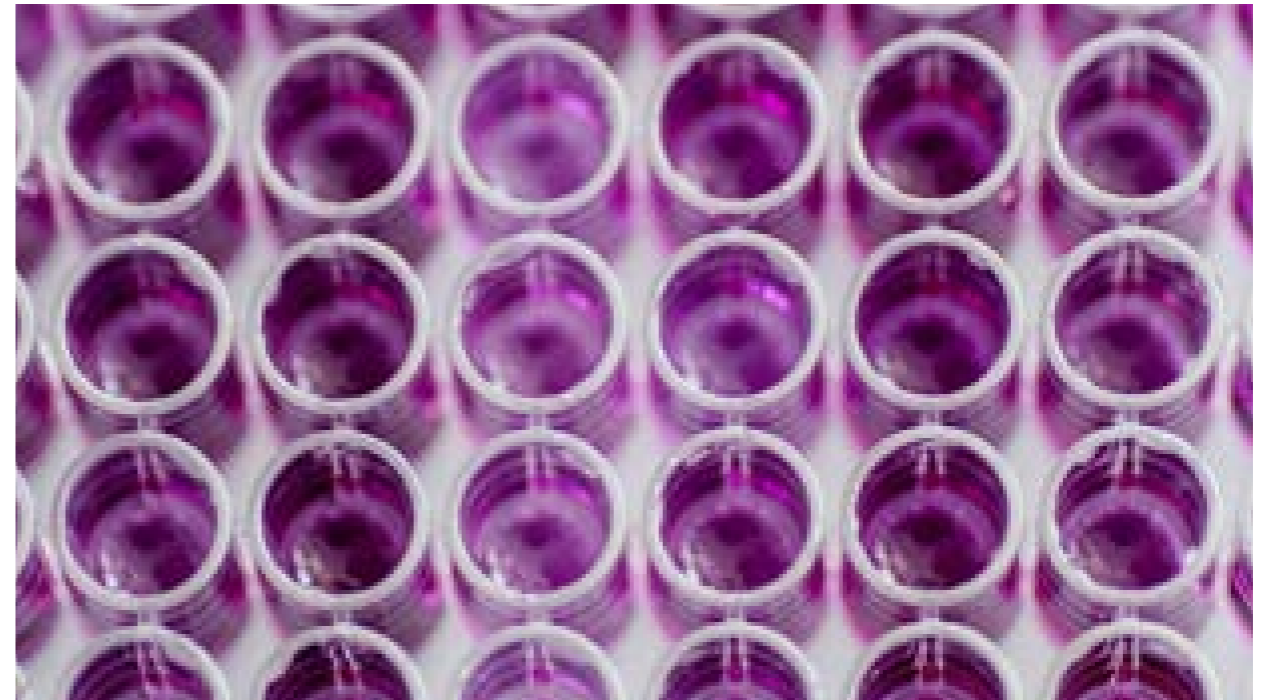
Drop-seq



inDrop

# Single-Cell Technologies

- (i) isolate single cells
- (ii) amplify genome efficiently
- (iii) sequence DNA





# Amplification and Sequencing

## Review: Next Generation Sequencing (NGS)

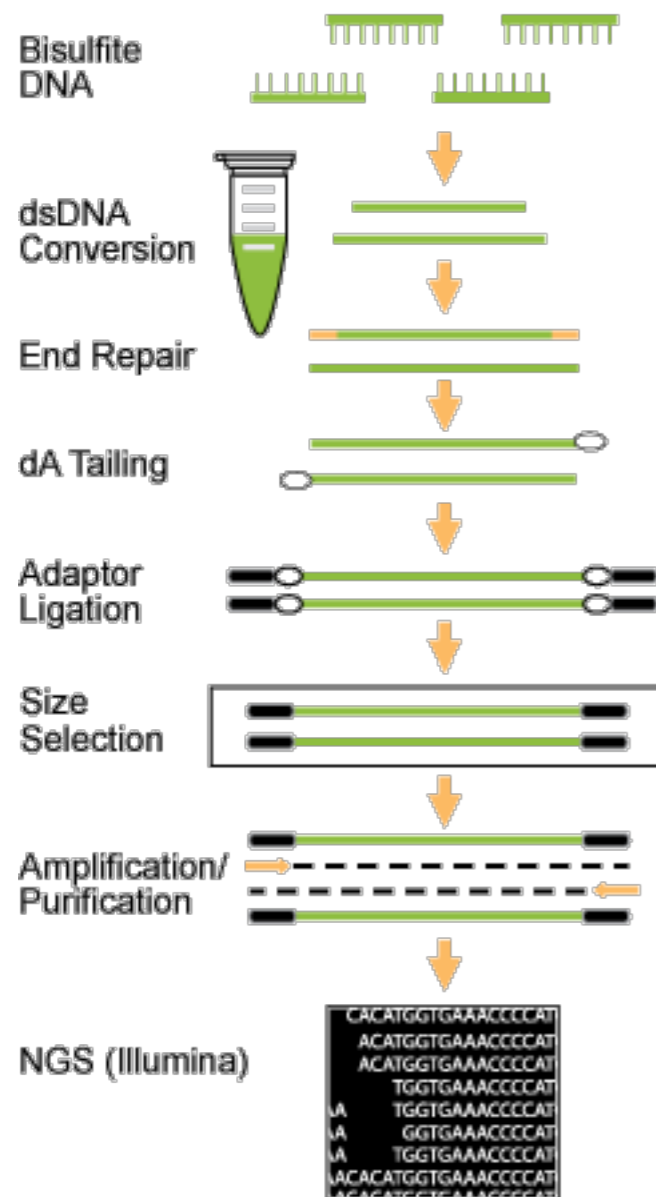


# Amplification and Sequencing

## Review: Next Generation Sequencing (NGS)

<https://www.illumina.com>

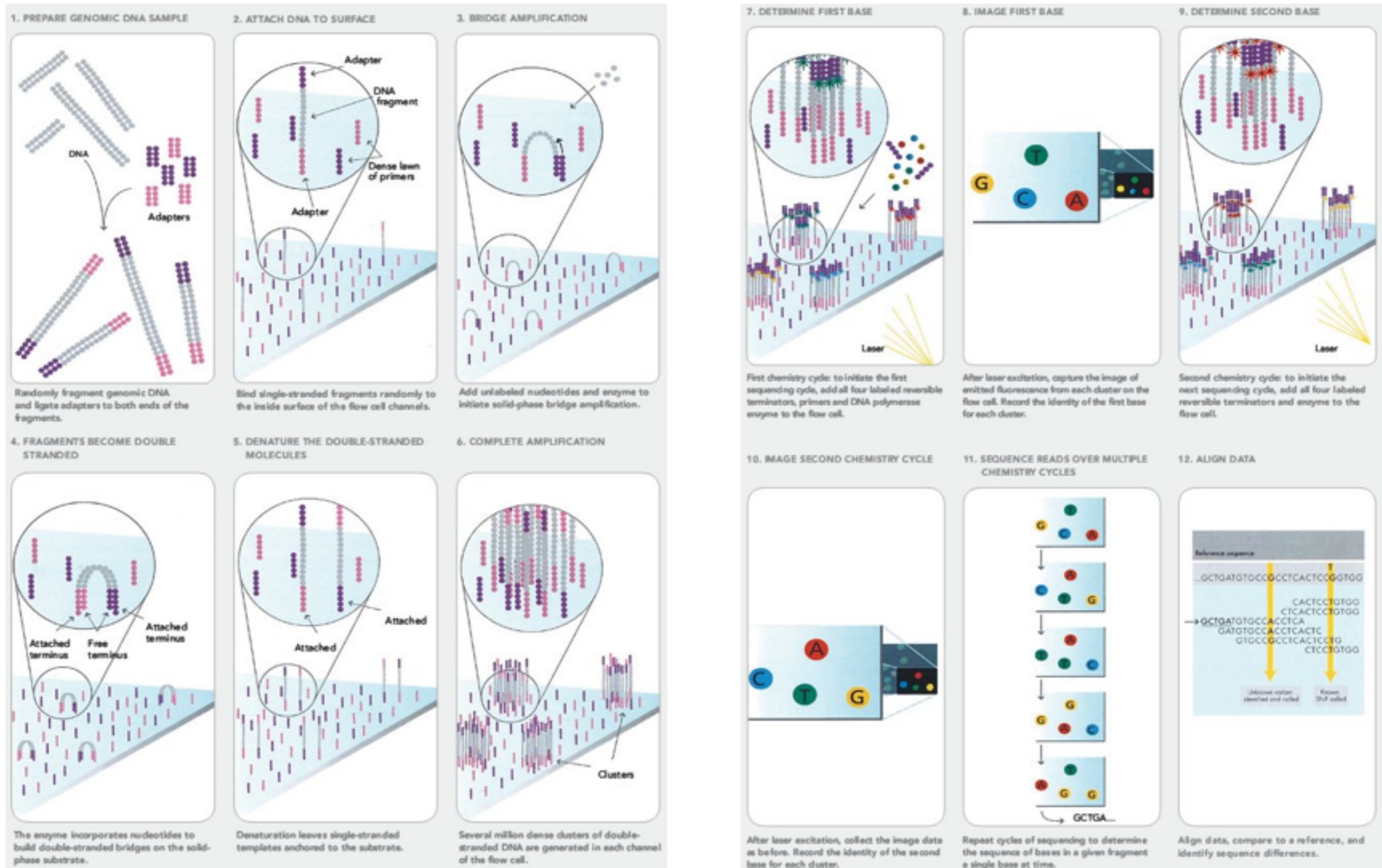
### library preparation



# Amplification and Sequencing

## Review: Next Generation Sequencing (NGS)

<https://www.illumina.com>



# Amplification and Sequencing

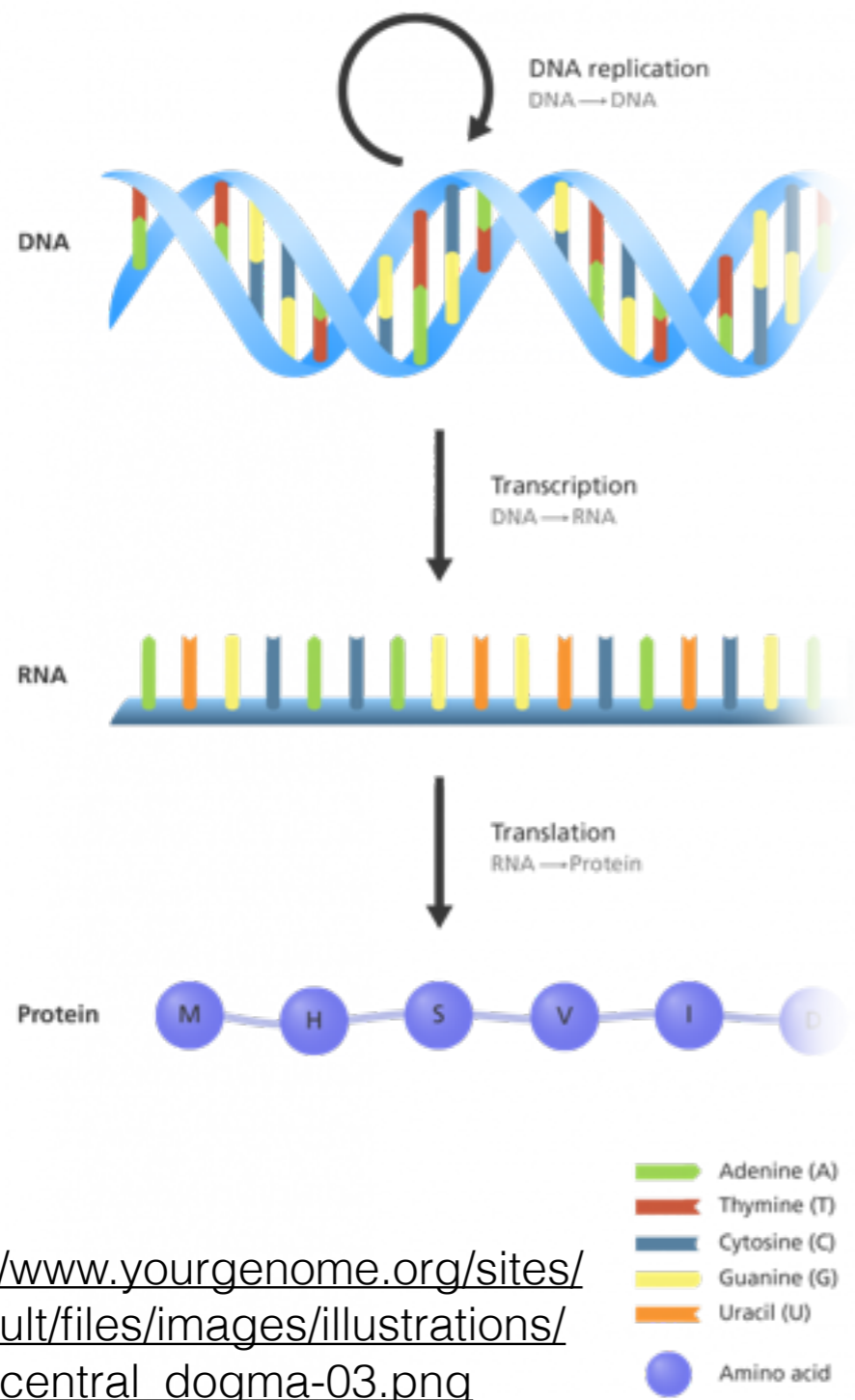


## Review: Next Generation Sequencing (NGS)

<https://www.illumina.com>

# Single-cell Amplification

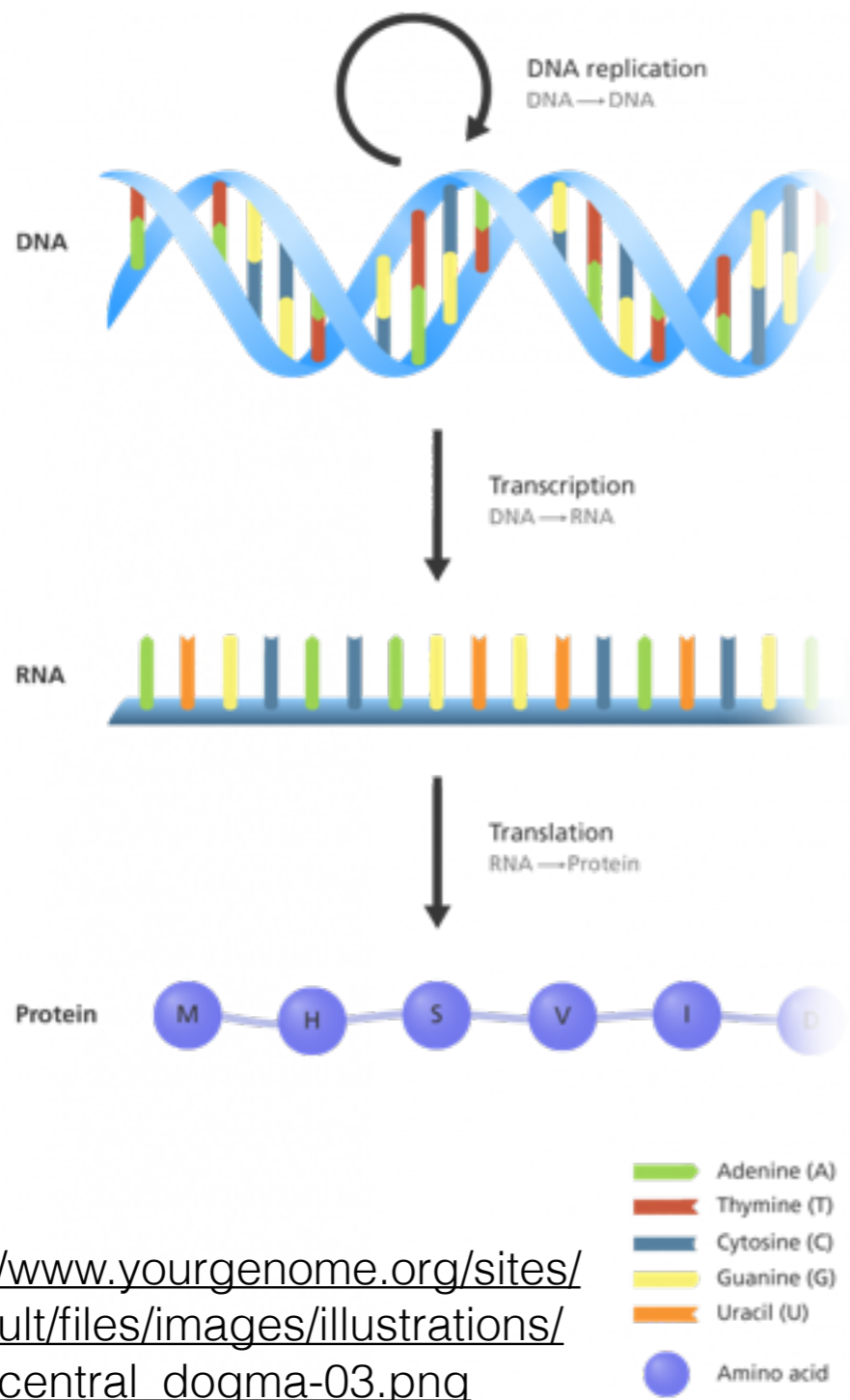
## Review: RNA-Sequencing



[http://www.yourgenome.org/sites/default/files/images/illustrations/central\\_dogma-03.png](http://www.yourgenome.org/sites/default/files/images/illustrations/central_dogma-03.png)

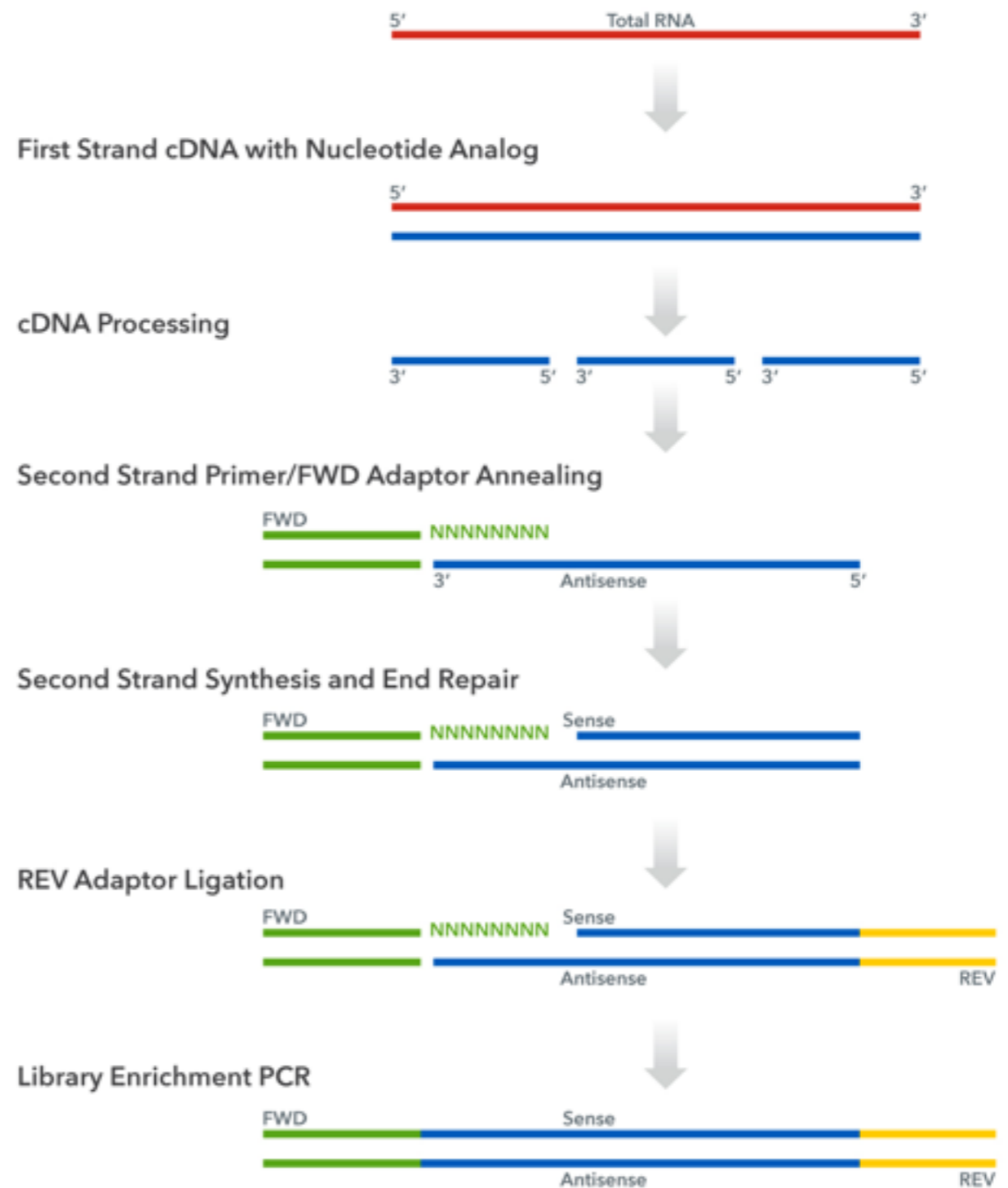
# Single-cell Amplification

## Review: RNA-Sequencing



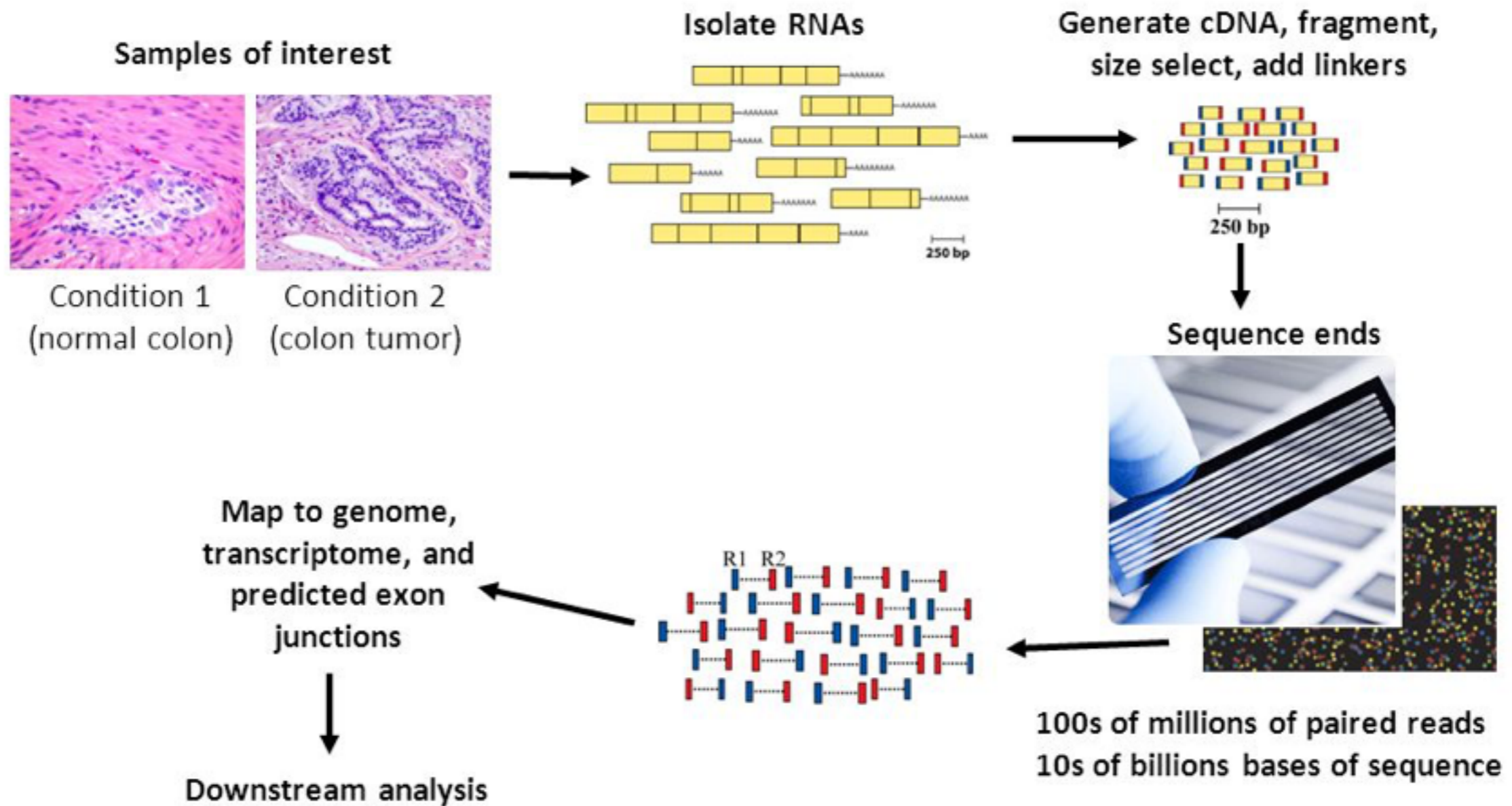
[http://www.yourgenome.org/sites/default/files/images/illustrations/central\\_dogma-03.png](http://www.yourgenome.org/sites/default/files/images/illustrations/central_dogma-03.png)

## library preparation



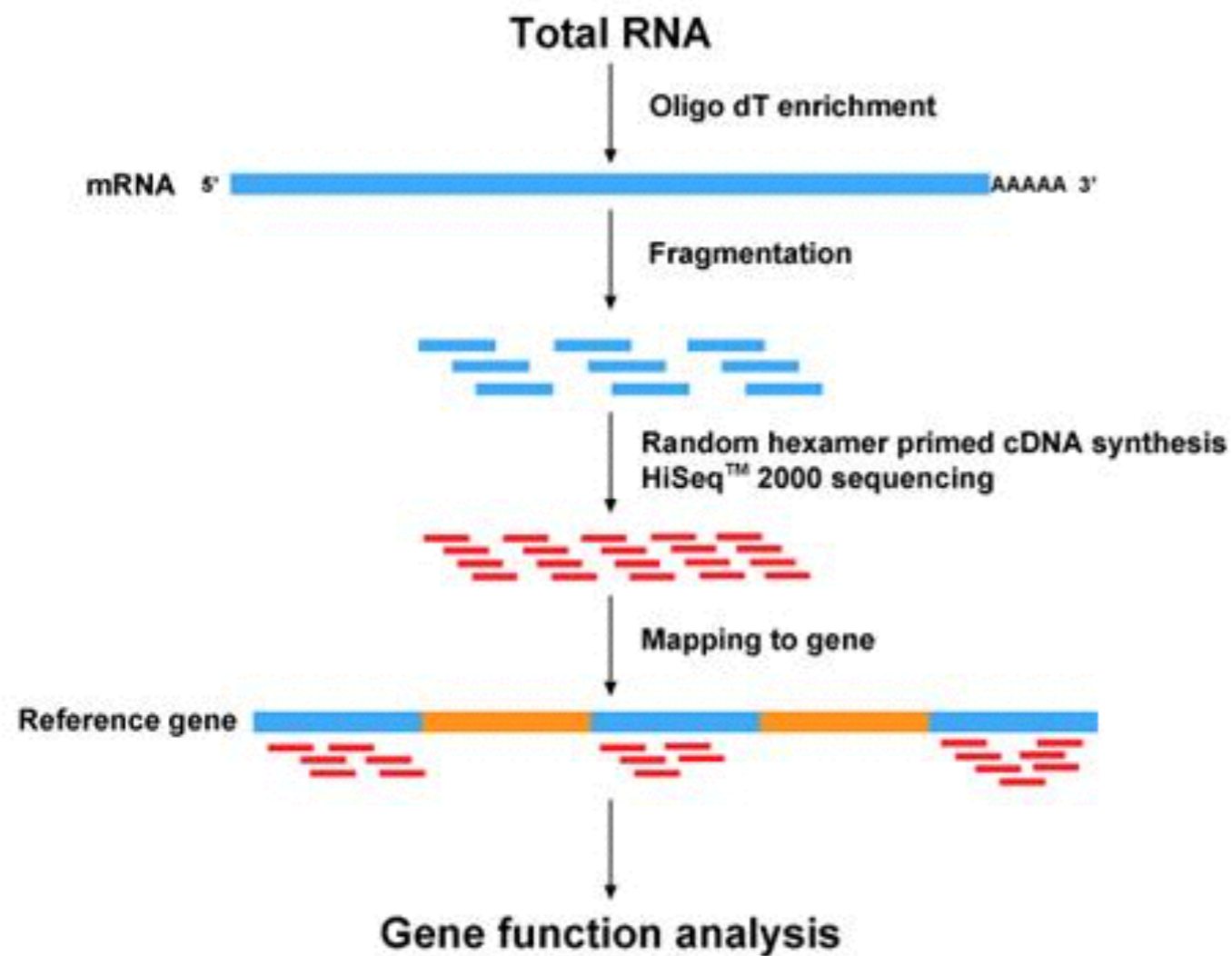
# Amplification and Sequencing

## Review: RNA-Sequencing



# Amplification and Sequencing

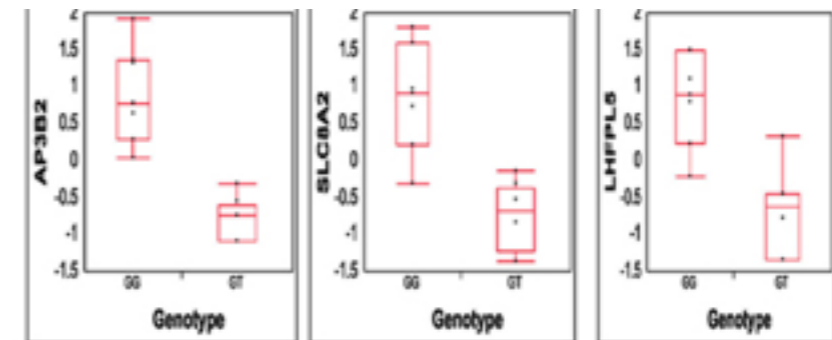
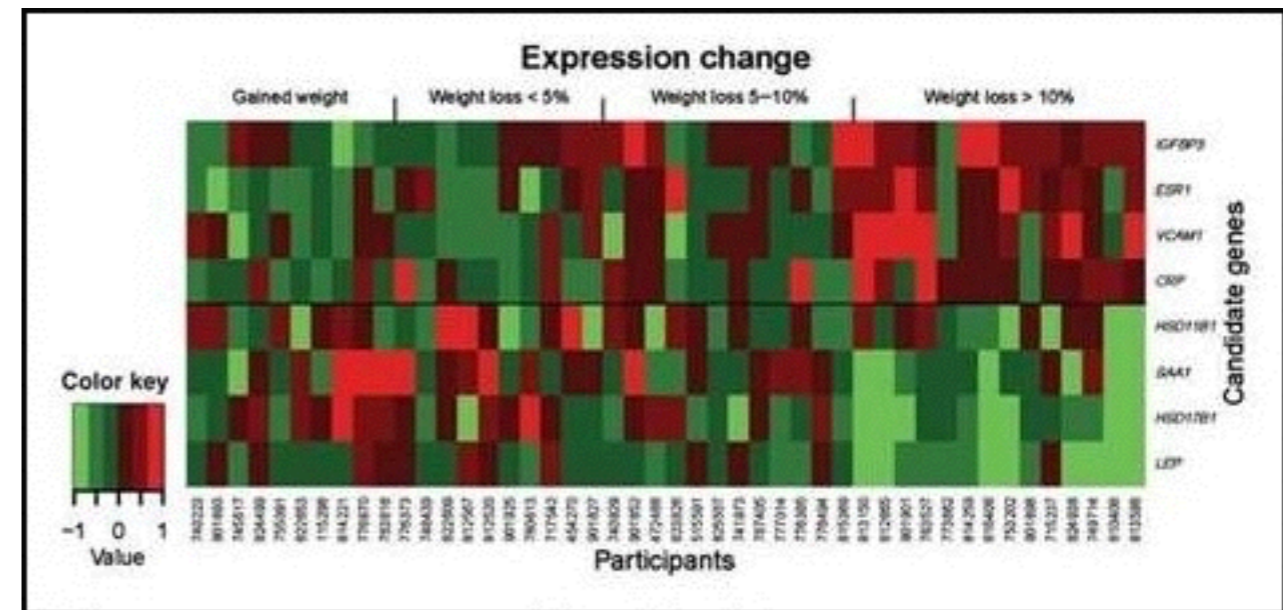
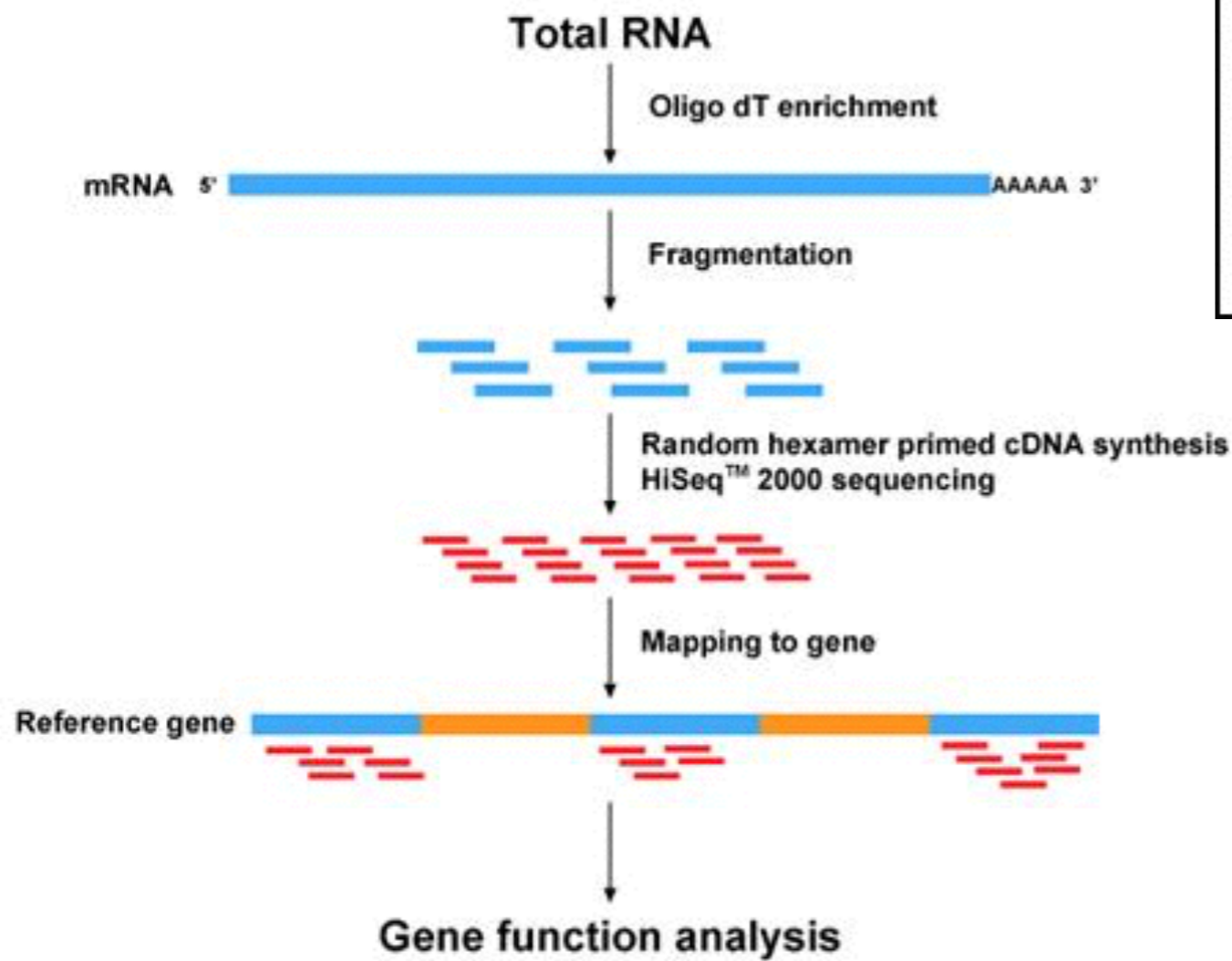
## Review: RNA-Sequencing





# Amplification and Sequencing

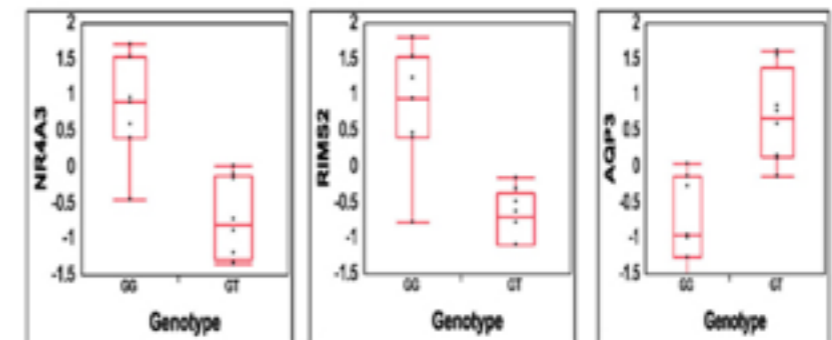
## Review: RNA-Sequencing



AP3B2

SLC8A2

LHPL5



NR4A3

RIMS2

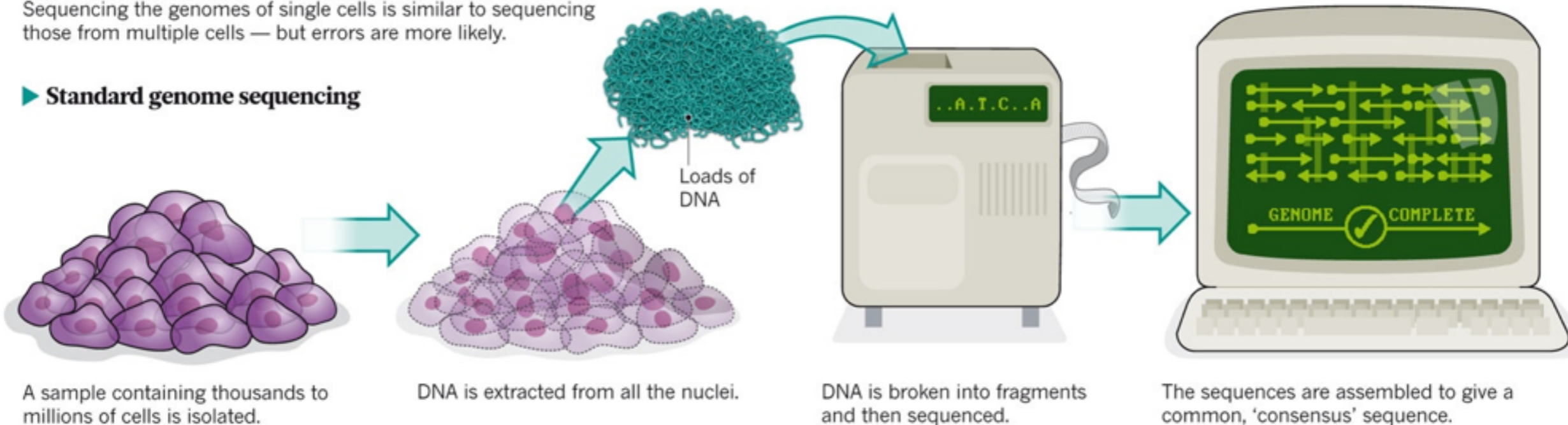
AQP3

# Single-cell Amplification

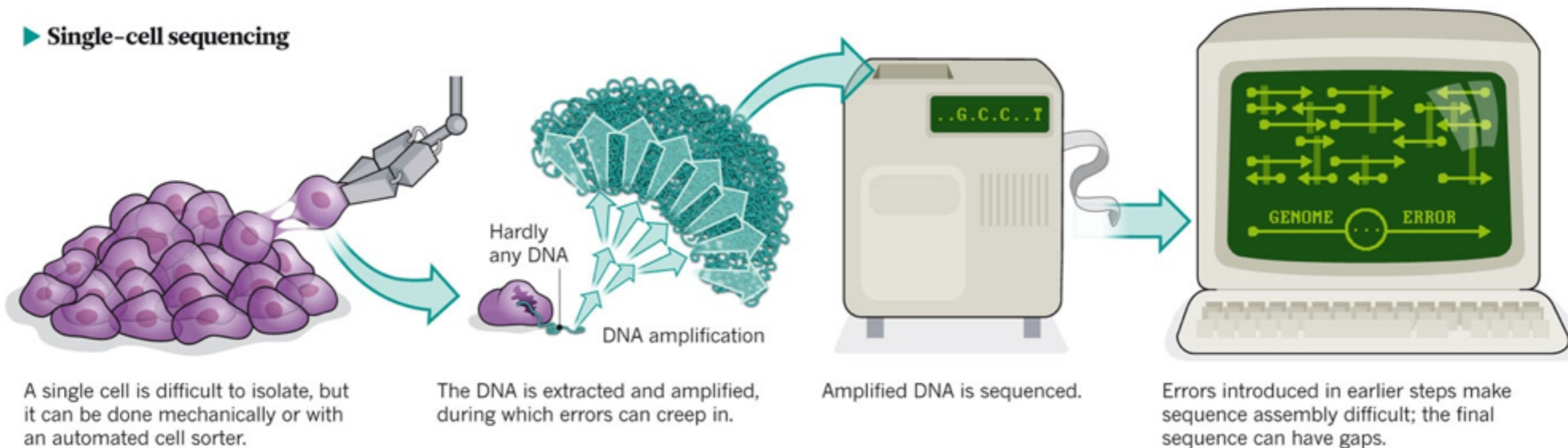
## ONE GENOME FROM MANY

Sequencing the genomes of single cells is similar to sequencing those from multiple cells — but errors are more likely.

### ► Standard genome sequencing

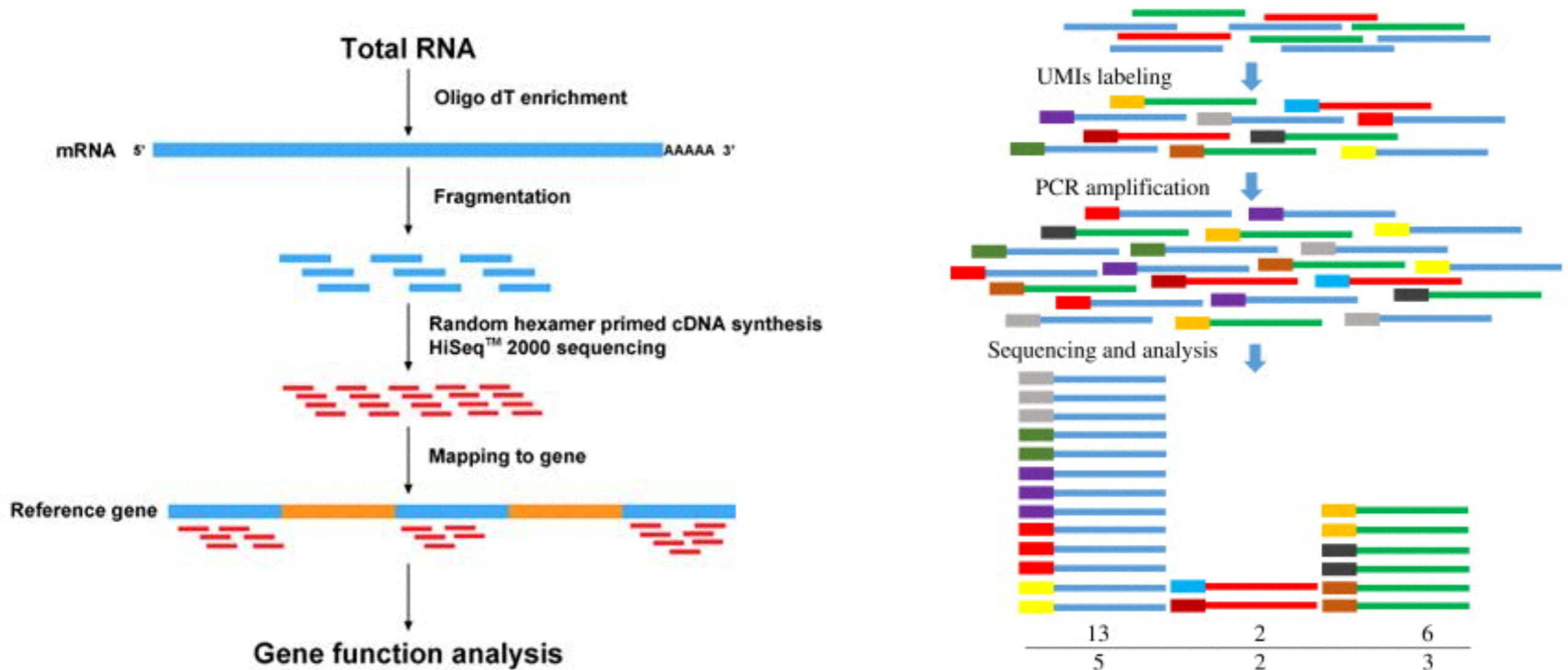


### ► Single-cell sequencing



# Single-cell Amplification

## Digital Expression Matrix: counting unique molecules



# Short Summary



## DNA sequencing:

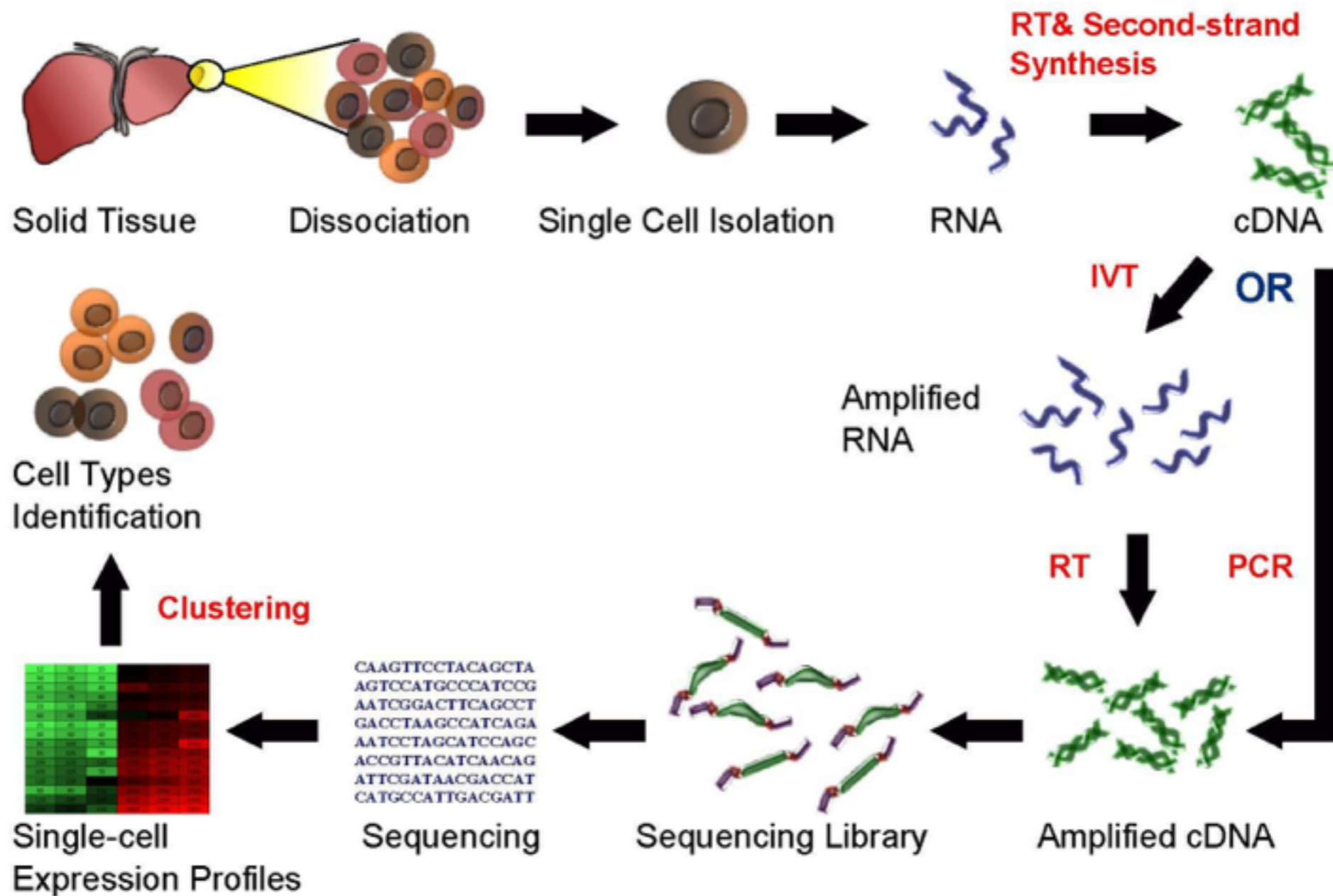
- new amplification methods other than PCR
- statistical methods for SNPs/CNV calling

## RNA sequencing:

- standards created for quality control
- can achieve high sequencing depth
- high cell throughput methods arising

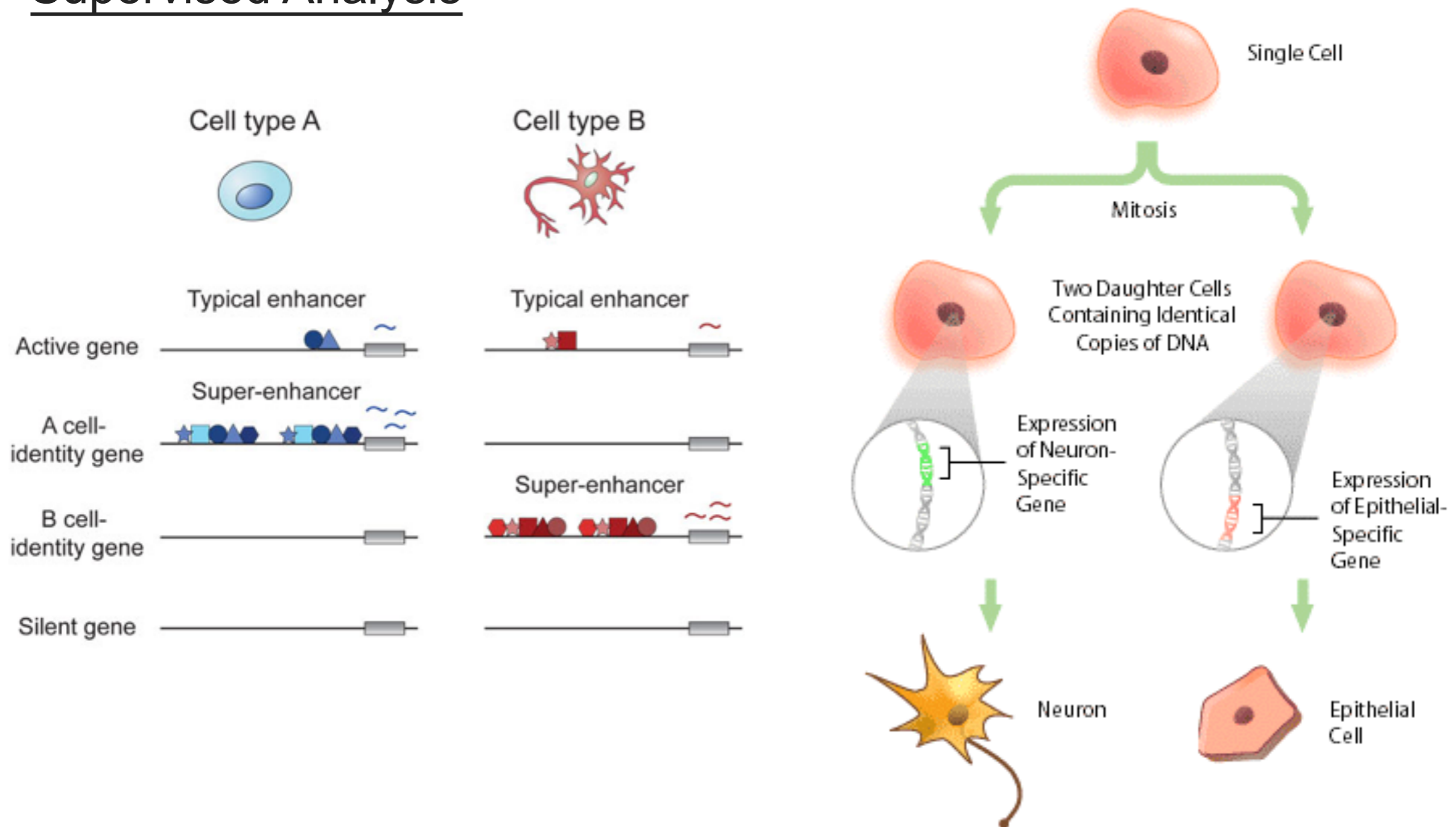
# Short Summary

## Single Cell RNA Sequencing Workflow



# Downstream Analysis

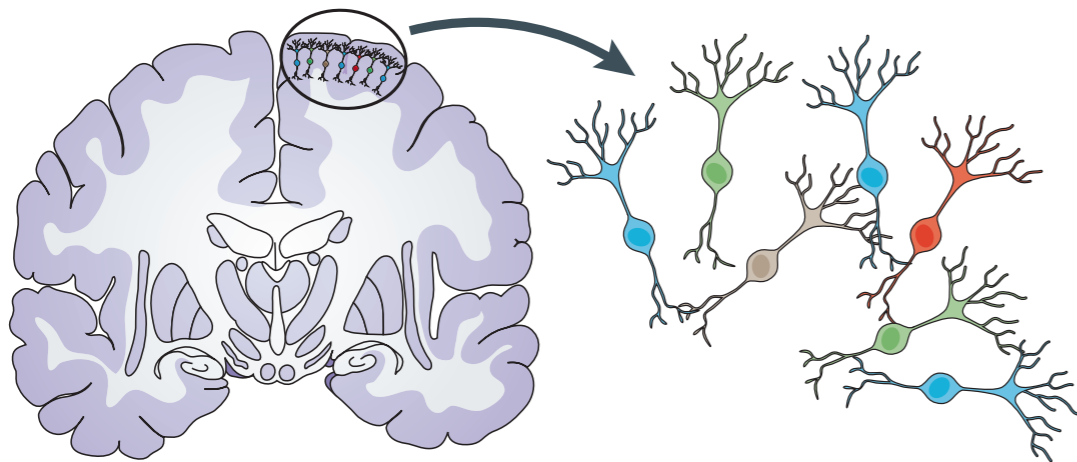
## Supervised Analysis



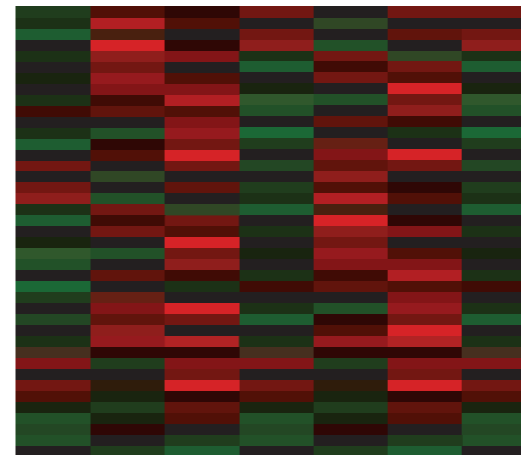
# Cell Population Identification

## Unsupervised Analysis

**a** Obtain an unbiased sample of single cells



**b** Generate single-cell expression profiles



**c** Identify cell types by clustering

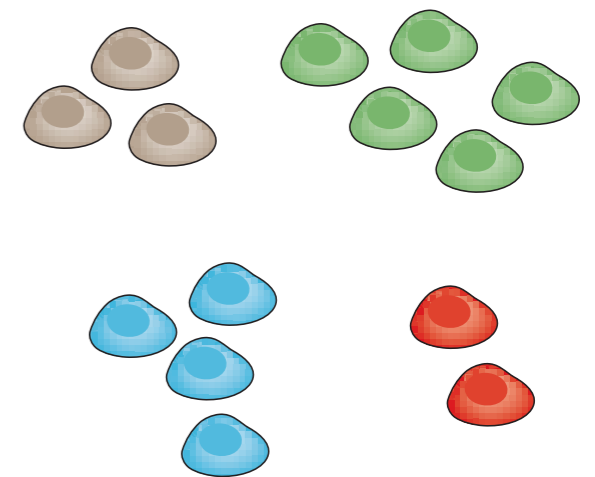


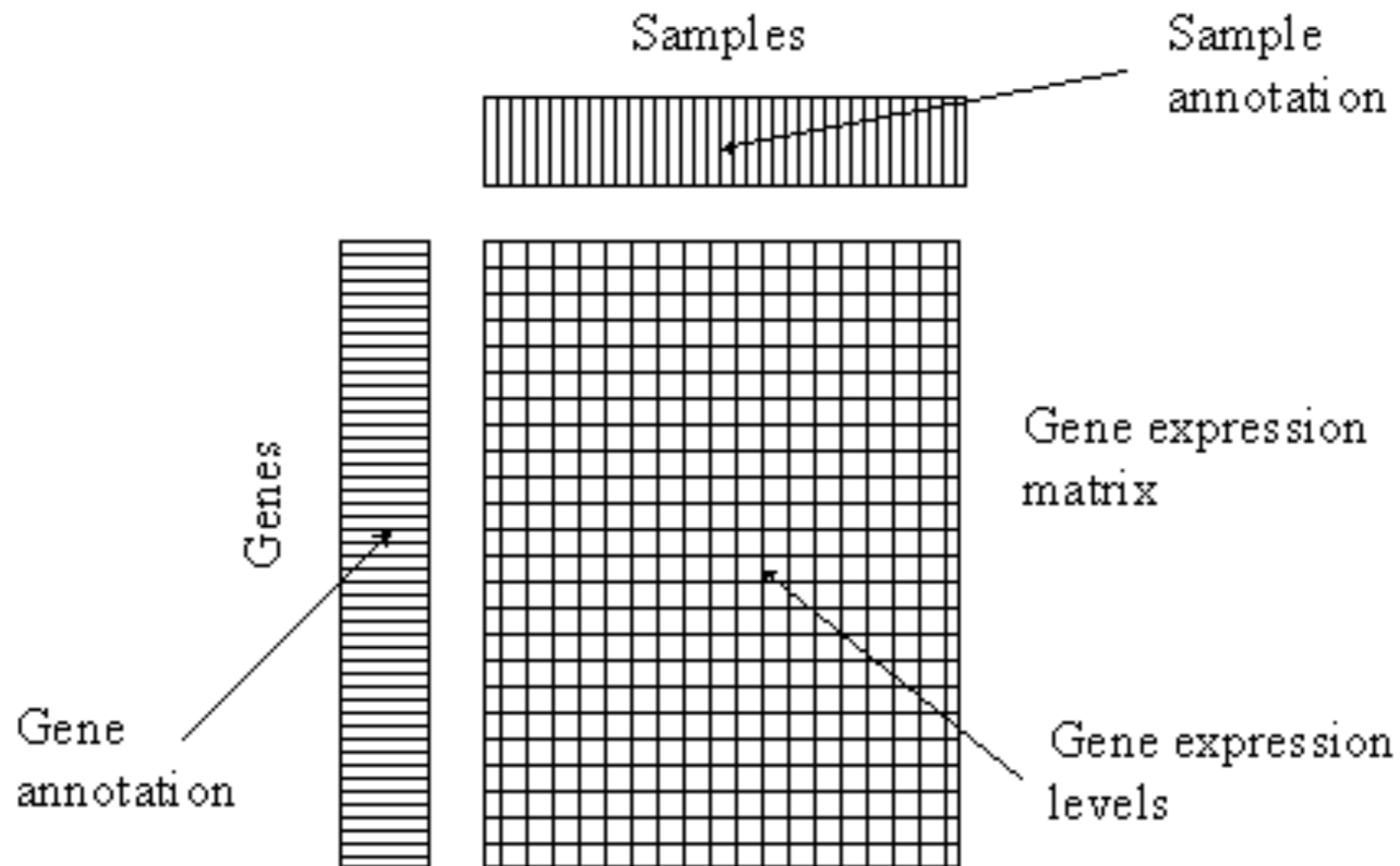
Figure 3 | **Cell-type discovery by unbiased sampling and transcriptome profiling of single cells.**

Shapiro, Ehud, Tamir Biezuner, and Sten Linnarsson. "Single-cell sequencing-based technologies will revolutionize whole-organism science." *Nature Reviews Genetics* 14.9 (2013): 618-630.

# Downstream Analysis

How do cell types differ from each other?

Is there any additional diversity in the same cell type?



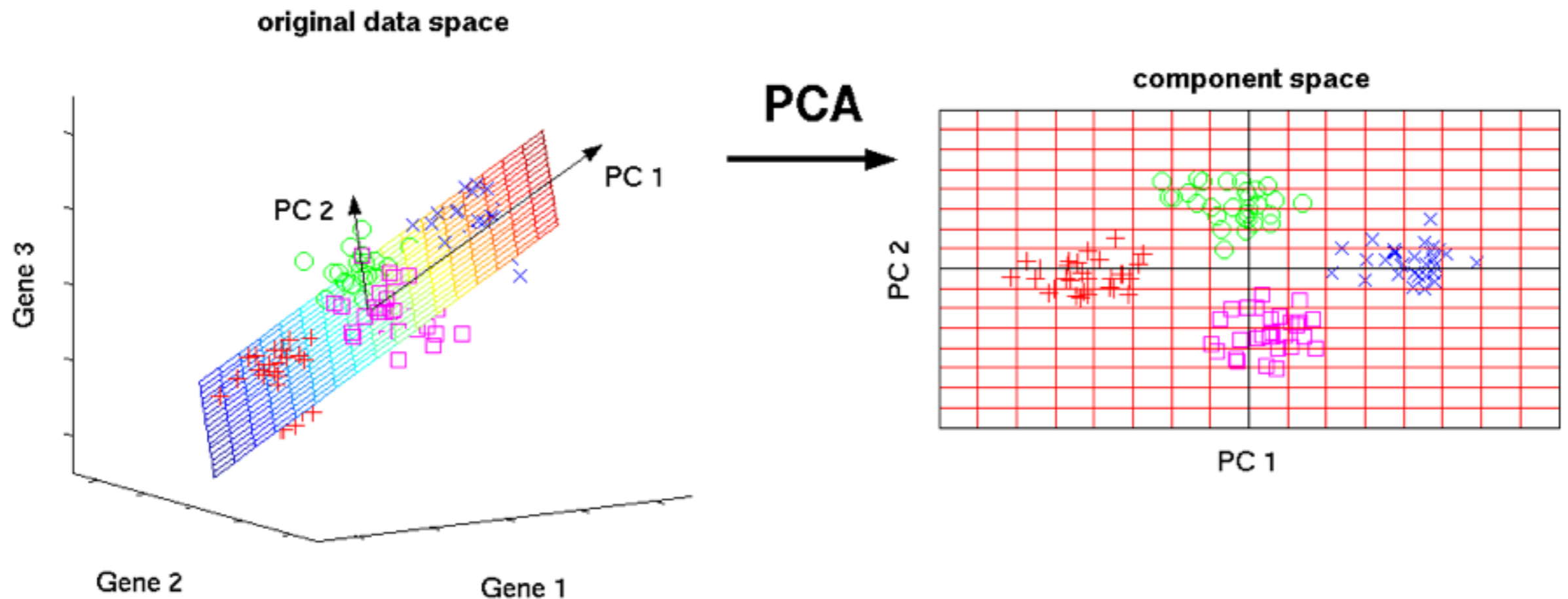


# Dimension Reduction



## Principle Component Analysis (PCA)

e.g., visualizing the samples in a smaller subspace



# PCA



## Probability and Linear Algebra Review

**Variance / Standard Deviation:** measure of the spread of the data

(Calculation: average distance from the mean of the data)

**Covariance:** measure of how much each of the dimensions vary from the mean with respect to each other; measured between 2 dimensions to see if there is a relationship between the 2 dimensions

\* The covariance between one dimension and itself is the variance.

# PCA

## Probability and Linear Algebra Review

E.g. for 3 dimensions, consider random vector  $(x,y,z)$ :

$$C = \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) & \text{cov}(x,z) \\ \text{cov}(y,x) & \text{cov}(y,y) & \text{cov}(y,z) \\ \text{cov}(z,x) & \text{cov}(z,y) & \text{cov}(z,z) \end{bmatrix}$$

Diagonal is the variances of  $x$ ,  $y$  and  $z$

$\text{cov}(x,y) = \text{cov}(y,x)$  hence matrix is symmetrical about the diagonal

$N$ -dimensional data will result in  $n \times n$  covariance matrix

# PCA



## Probability and Linear Algebra Review

- The eigenvalue problem is any problem having the following form:

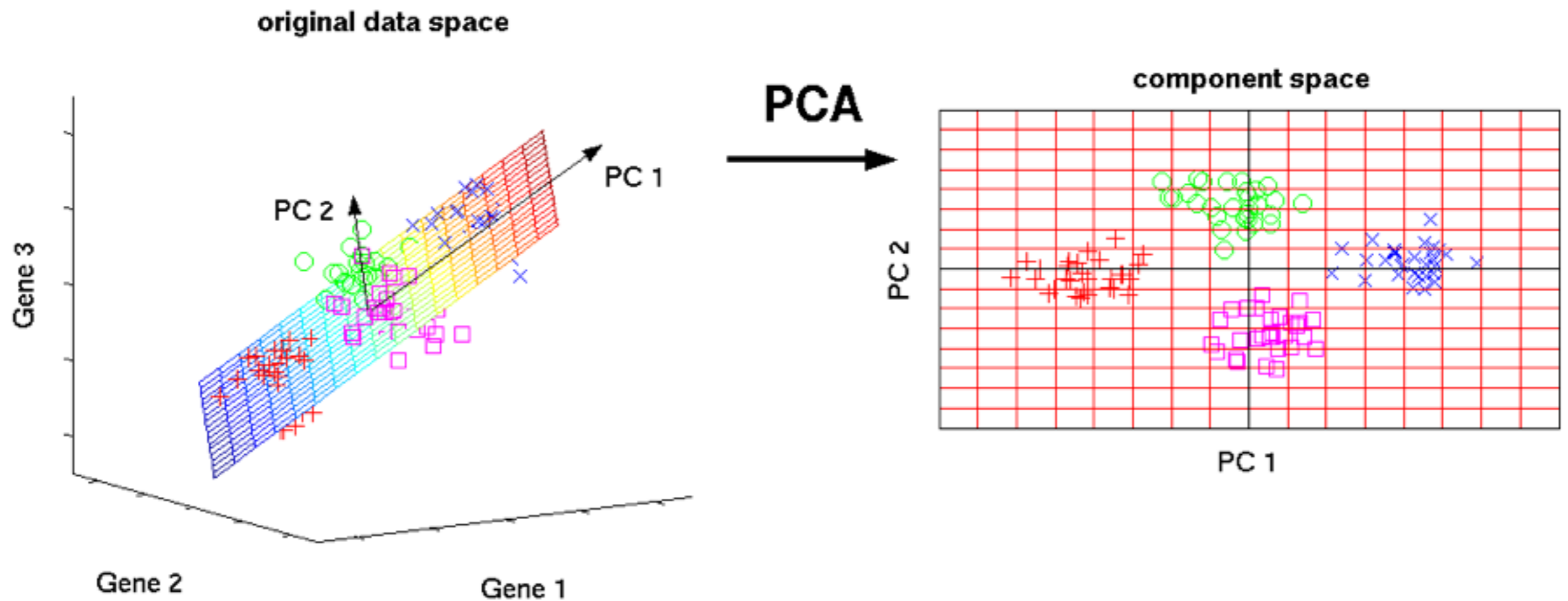
$$A \cdot v = \lambda \cdot v$$

- $A$ :  $n \times n$  matrix
  - $v$ :  $n \times 1$  non-zero vector
  - $\lambda$ : scalar
- Any value of  $\lambda$  for which this equation has a solution is called the eigenvalue of  $A$  and vector  $v$  which corresponds to this value is called the eigenvector of  $A$ .

# Dimension Reduction



## Principle Component Analysis (PCA)



# PCA



Principal component analysis (PCA) converts a set of observations of possibly correlated variables into a set of values of uncorrelated variables called principal components.

The first principal component is the projection of the data into a single dimension that has as high a variance as possible (that is, accounts for as much of the variability in the data as possible); each succeeding component in turn has the highest variance possible under the constraint that it be orthogonal to (uncorrelated with) the preceding components.

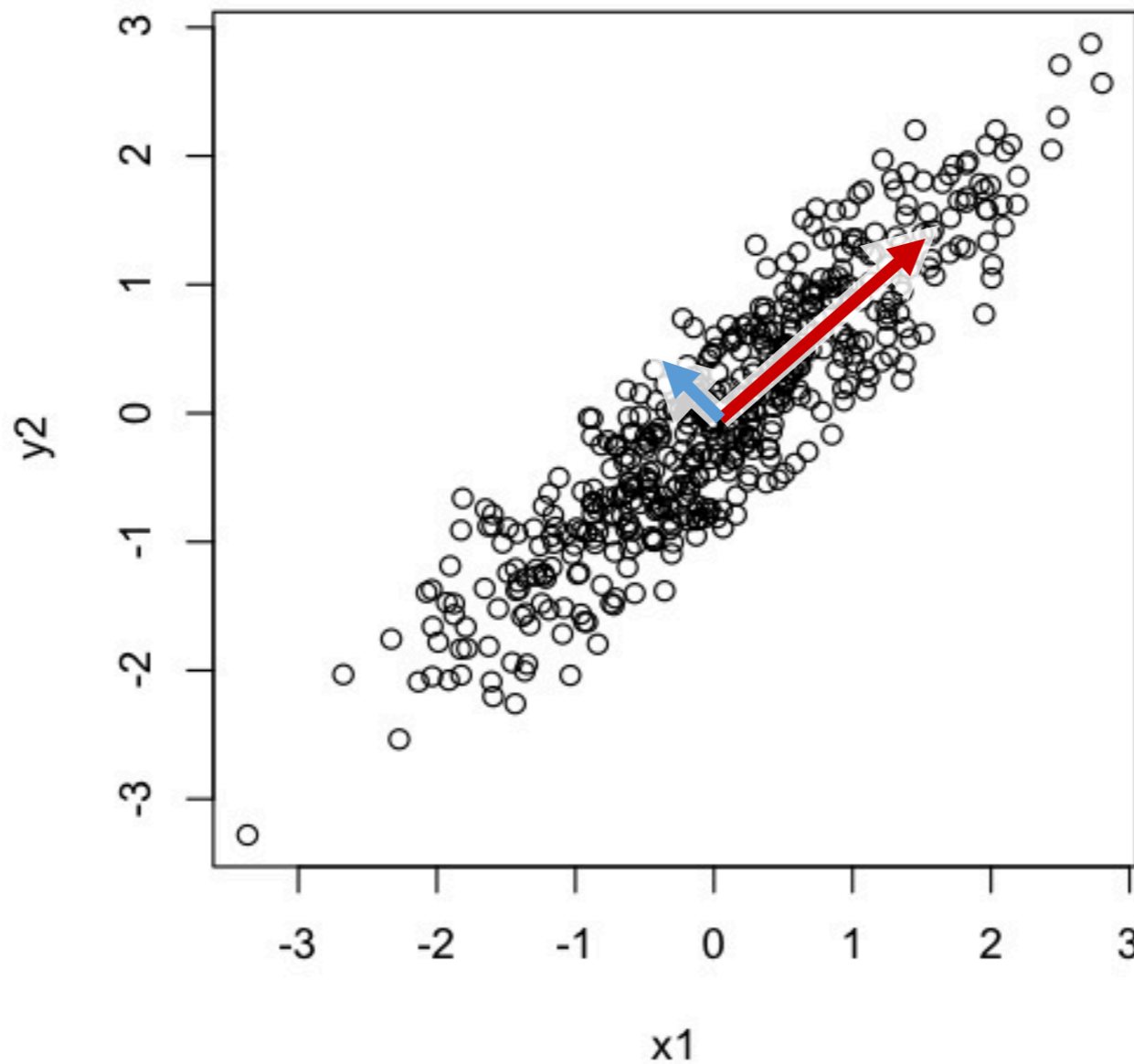
Therefore the PCs provide a view on the structure of the data that best explains its variance.

# PCA



The example data is two-dimensional, but most of the information is contained along a dimension shown here by the **red** vector.

We could thus restrict our analysis to a projection along that vector.



# PCA



## PCA process –STEP 1

- Subtract the mean

from each of the data dimensions. All the  $x$  values have  $\bar{x}$  subtracted and  $y$  values have  $\bar{y}$  subtracted from them. This produces a data set whose mean is zero.

Subtracting the mean makes variance and covariance calculation easier by simplifying their equations. The variance and co-variance values are not affected by the mean value.



# PCA



## PCA process –STEP 1

DATA:

<u>x</u>	<u>y</u>
2.5	2.4
0.5	0.7
2.2	2.9
1.9	2.2
3.1	3.0
2.3	2.7
2	1.6
1	1.1
1.5	1.6
1.1	0.9

ZERO MEAN DATA:

<u>x</u>	<u>y</u>
.69	.49
-1.31	-1.21
.39	.99
.09	.29
1.29	1.09
.49	.79
.19	-.31
-.81	-.81
-.31	-.31
-.71	-1.01

# PCA



## PCA process –STEP 2

- Calculate the covariance matrix

$$\text{cov} = \begin{pmatrix} .616555556 & .615444444 \\ .615444444 & .716555556 \end{pmatrix}$$

- since the non-diagonal elements in this covariance matrix are positive, we should expect that both the x and y variable increase together.

# PCA



## PCA process –STEP 3

- Calculate the eigenvectors and eigenvalues of the covariance matrix

$$\text{eigenvalues} = \begin{pmatrix} .0490833989 \\ 1.28402771 \end{pmatrix}$$

$$\text{eigenvectors} = \begin{pmatrix} -.735178656 & -.677873399 \\ .677873399 & -.735178656 \end{pmatrix}$$

# PCA



## PCA process –STEP 4

- Reduce dimensionality and form *feature vector* the eigenvector with the *highest* eigenvalue is the *principle component* of the data set.
- In our example, the eigenvector with the largest eigenvalue was the one that pointed down the middle of the data.
- Once eigenvectors are found from the covariance matrix, the next step is to order them by eigenvalue, highest to lowest. This gives you the components in order of significance.

# PCA



## PCA process –STEP 4

- Now, if you like, you can decide to *ignore* the components of lesser significance
- You do lose some information, but if the eigenvalues are small, you don't lose much
  - n dimensions in your data
  - calculate n eigenvectors and eigenvalues
  - choose only the first p eigenvectors
  - final data set has only p dimensions.

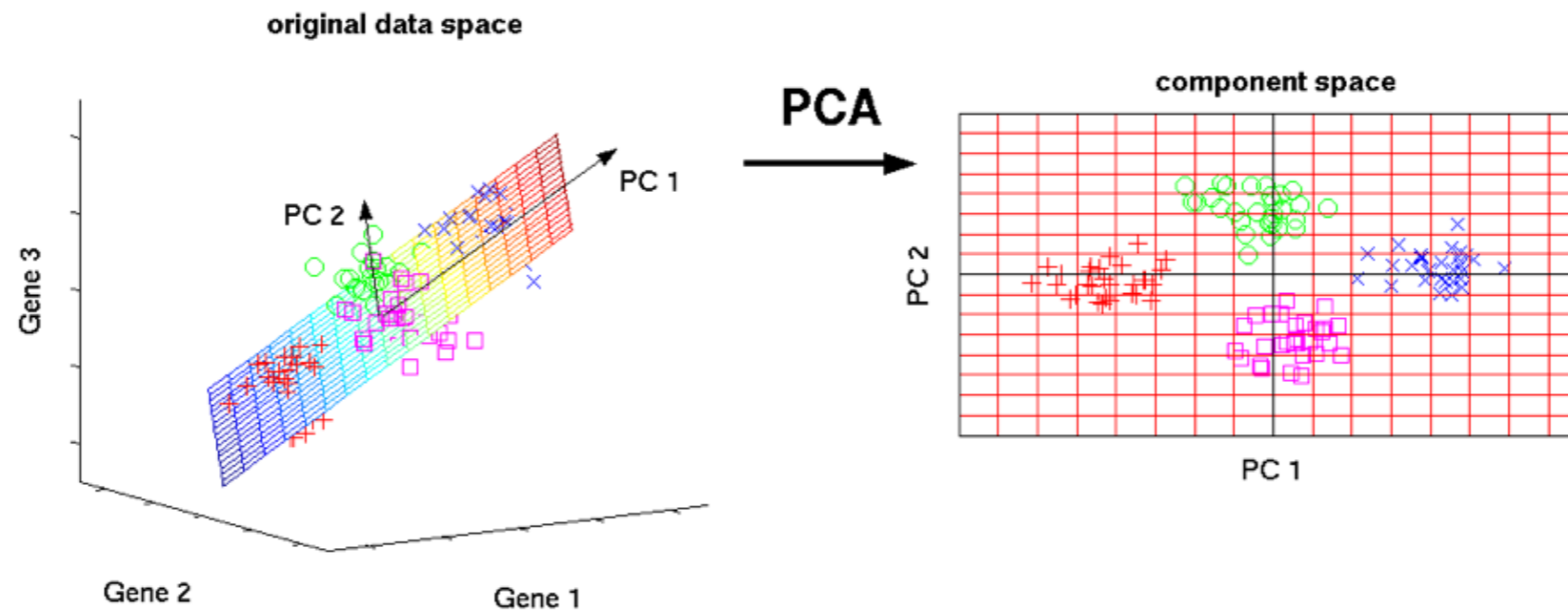
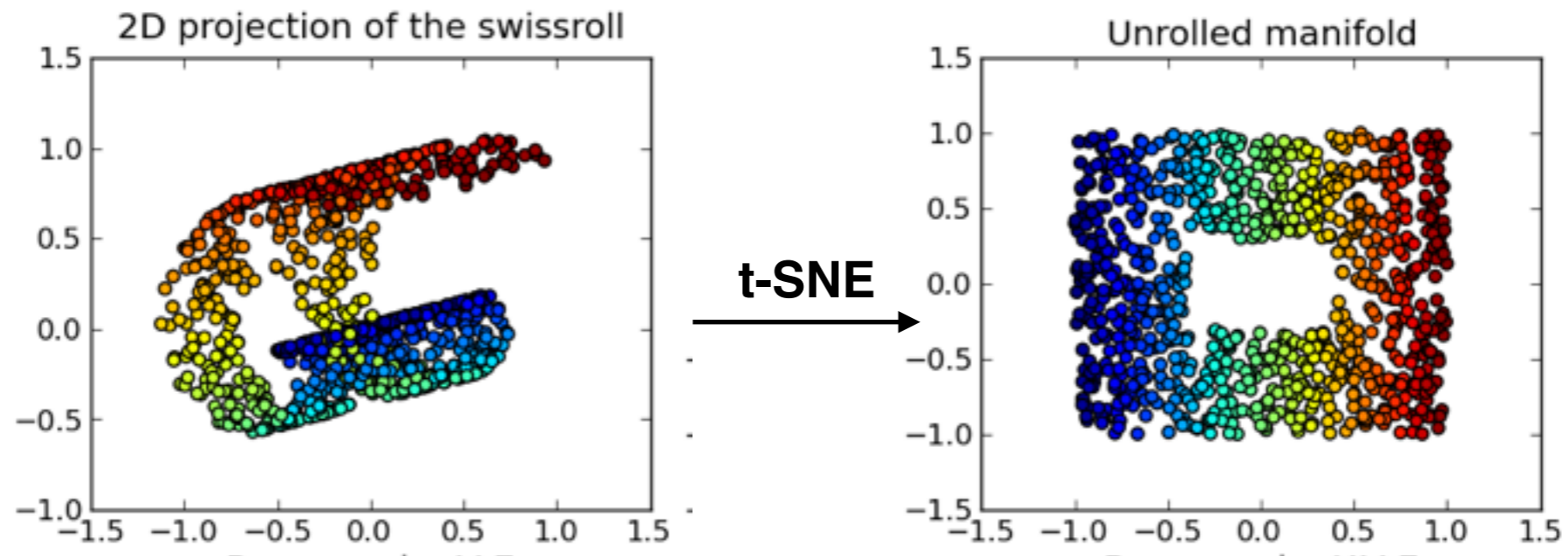
# Dimension Reduction



## **Principle Component Analysis (PCA)**

- linear multivariate statistical analysis
- understand underlying data structures
- identify bias, experimental errors, batch effects
- visualize the samples in a smaller subspace (dimension reduction)
- visualize the relationship between variables (correlation analysis)

# t-SNE



# t-SNE



## Key quantities

high-dimensional  
joint distribution



$$KL(P||Q) = \sum_{i \neq j} p_{ij} \log \frac{p_{ij}}{q_{ij}}$$



Kullback–Leibler divergence  
(to be minimized)



low-dimensional  
joint distribution



# Cluster Analysis



Cluster: a collection of data objects

Similar to the objects in the same cluster (Intraclass similarity)

Dissimilar to the objects in other clusters (Interclass dissimilarity)

Cluster analysis

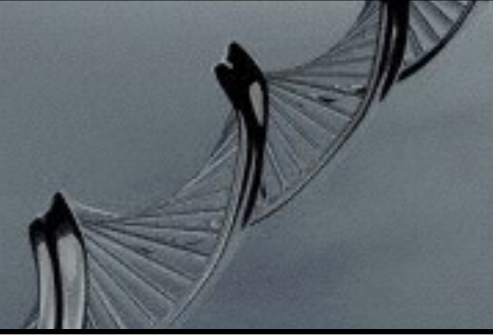
Statistical method for grouping a set of data objects into clusters

A good clustering method produces high quality clusters with high intraclass similarity and low interclass similarity

Clustering is an **unsupervised classification** method

Can be a stand-alone tool or as a preprocessing step for other algorithms

# Cluster Analysis

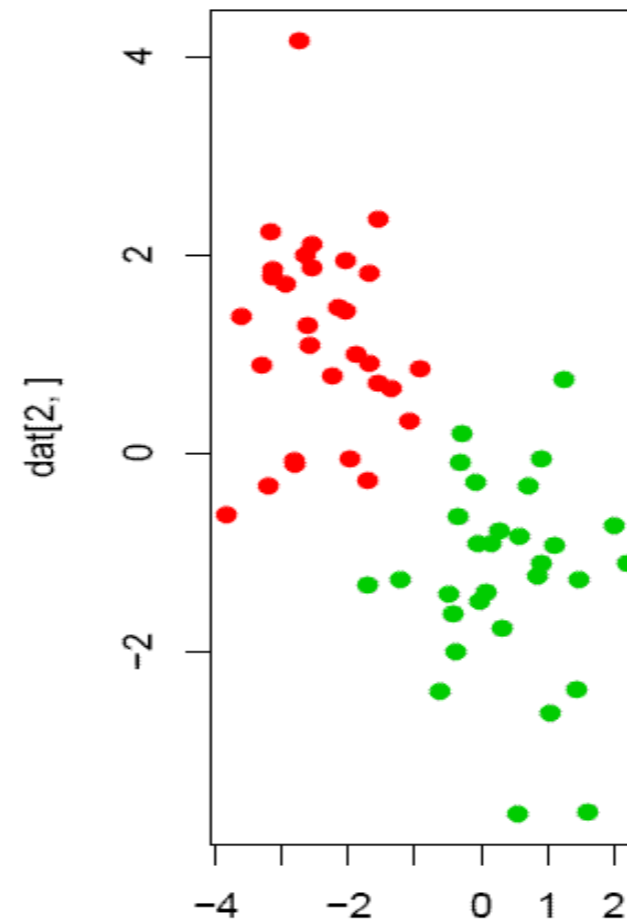
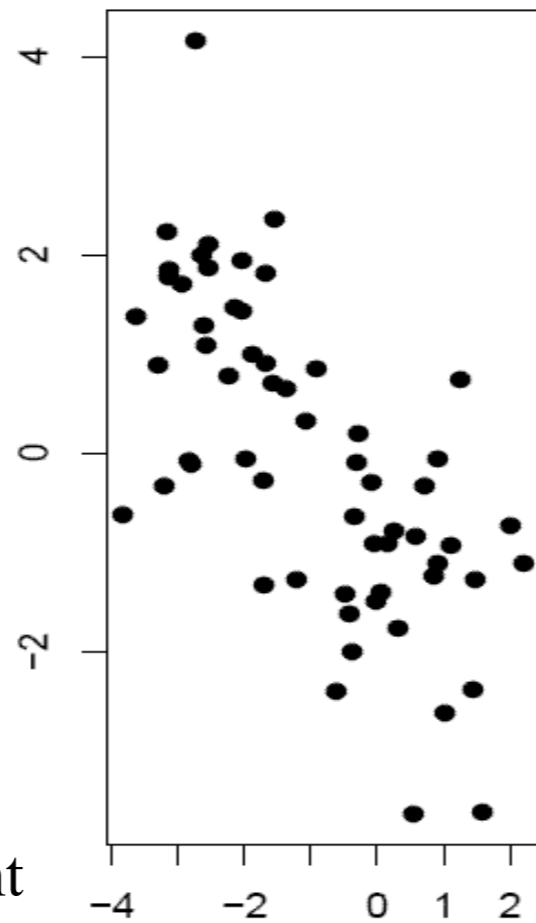


Group objects according to their similarity

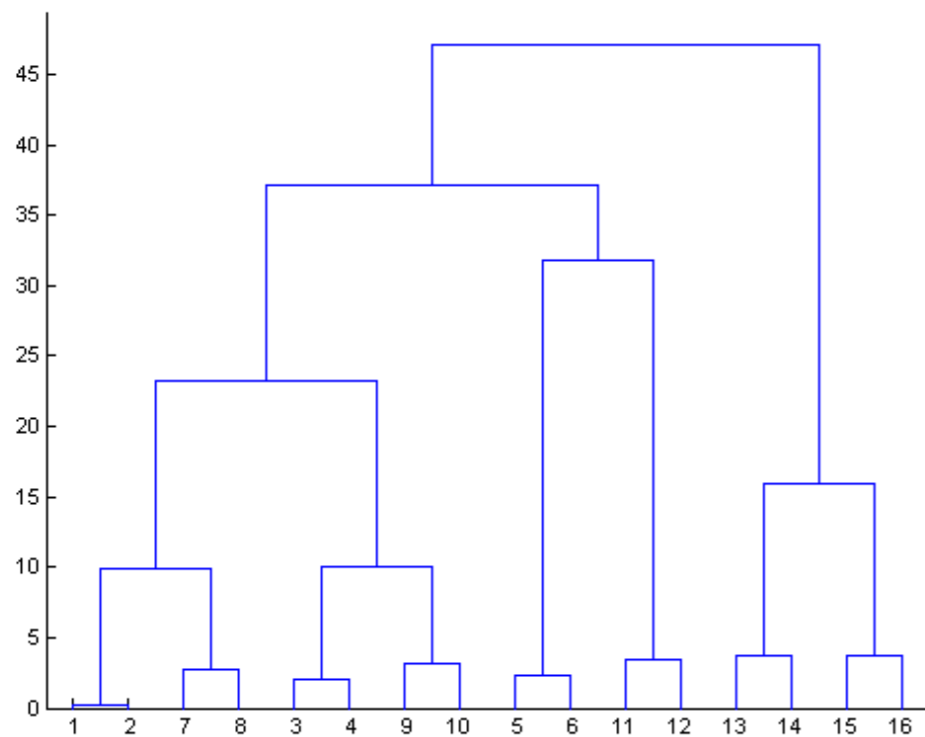
**Cluster:**

a set of objects that are similar to each other and separated from the other objects.

Example: green/red data points were generated from two different normal distributions



# Hierarchical Clustering



- This produces a binary tree or *dendrogram*
- The final cluster is the root and each data item is a leaf
- The height of the bars indicate how close the items are

# Hierarchical Clustering



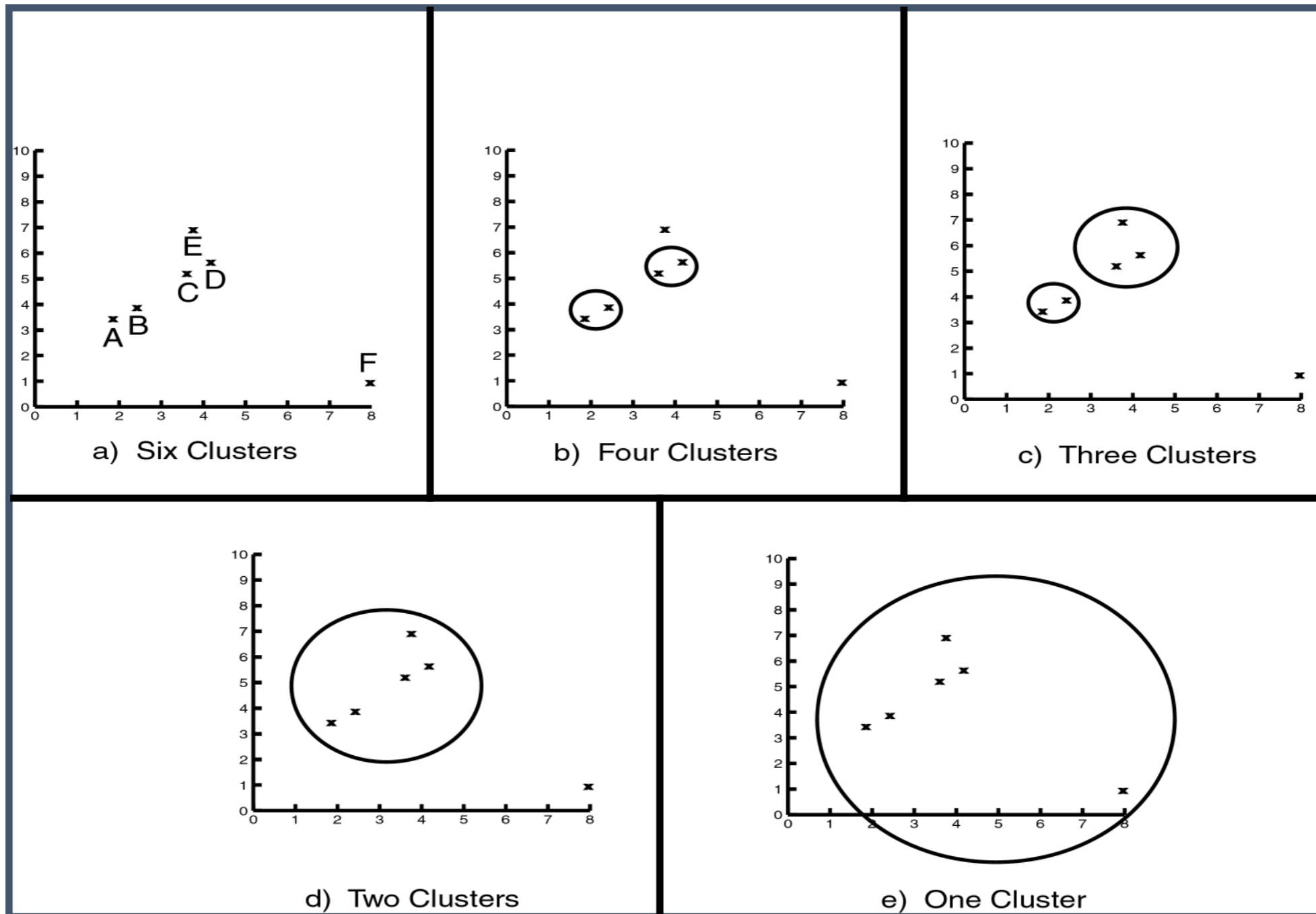
Start with every data point in a separate cluster

Keep merging the most similar pairs of data points/clusters until we have one big cluster left

*This is called a bottom-up or agglomerative method*

# Hierarchical Clustering

## Levels of Clustering



# Hierarchical Clustering



## Linkage in Hierarchical Clustering

We already know about distance measures between data items, but what about between a data item and a cluster or between two clusters?

We just treat a data point as a cluster with a single item, so our only problem is to define a ***linkage*** method between clusters  
As usual, there are lots of choices...

# Hierarchical Clustering



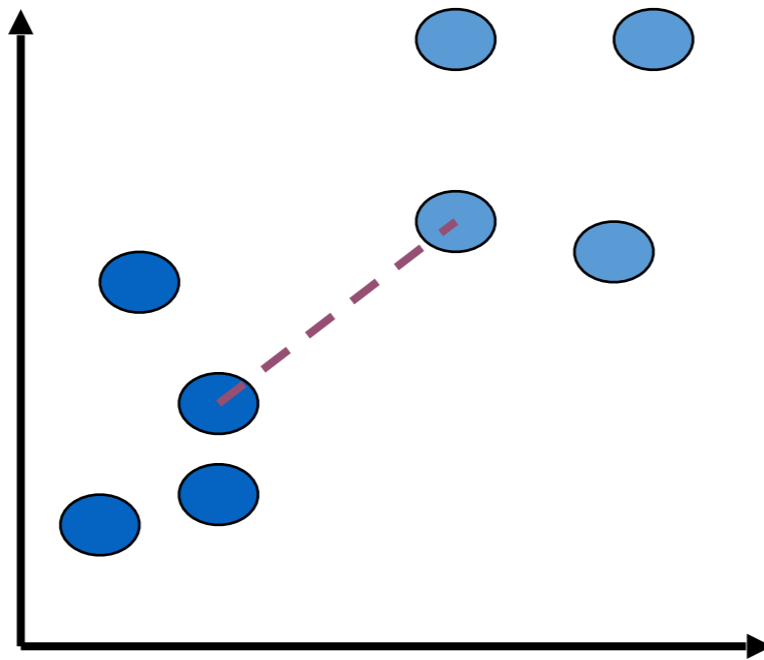
## Average Linkage

- Definition
  - Each cluster  $c_i$  is associated with a mean vector  $\mu_i$  which is the mean of all the data items in the cluster
  - The distance between two clusters  $c_i$  and  $c_j$  is then just  $d(\mu_i, \mu_j)$
- This is somewhat non-standard – this method is usually referred to as centroid linkage and average linkage is defined as the average of all pairwise distances between points in the two clusters

# Hierarchical Clustering

## Single Linkage

- The minimum of all pairwise distances between points in the two clusters
- Tends to produce long, “loose” clusters

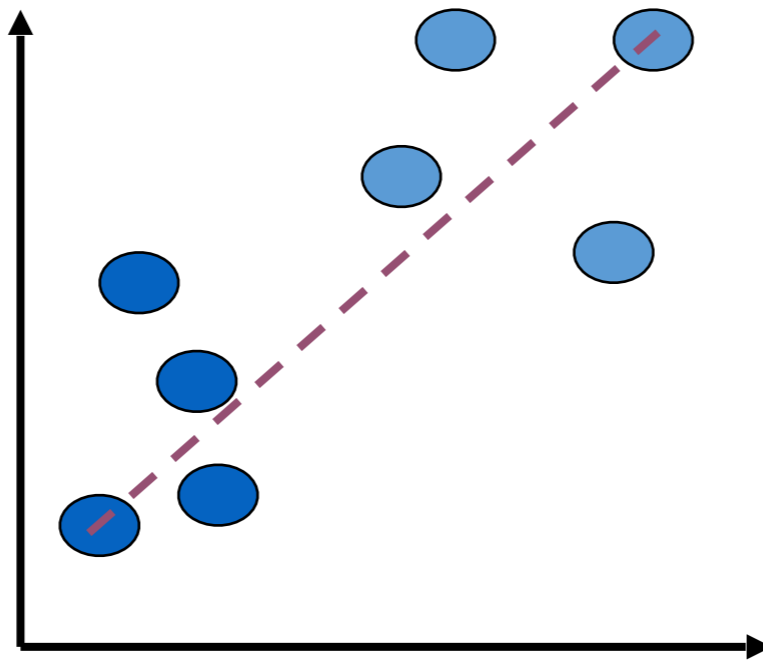




# Hierarchical Clustering

## Complete Linkage

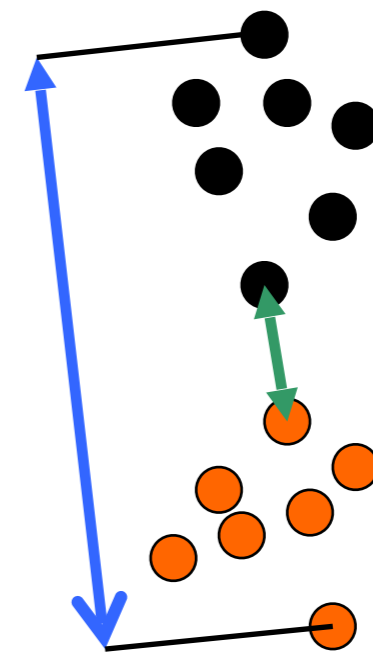
- The maximum of all pairwise distances between points in the two clusters
- Tends to produce very tight clusters



# Hierarchical Clustering

## Distances between clusters (summary)

- Calculation of the distance between two clusters is based on the pairwise distances between members of the clusters.
  - **Complete linkage:** largest distance between points
  - **Average linkage:** average distance between pairs of points
  - **Single linkage:** smallest distance between points
  - **Centroid:** distance between centroids



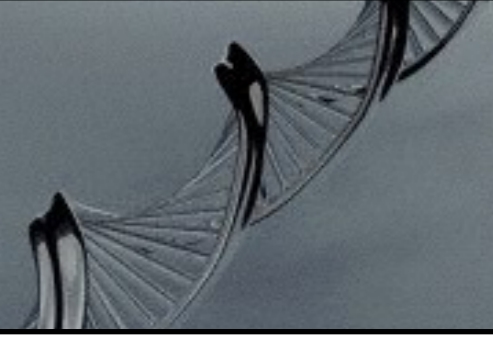
Complete linkage gives preference to compact/spherical clusters. Single linkage can produce long stretched clusters.

# Hierarchical Clustering



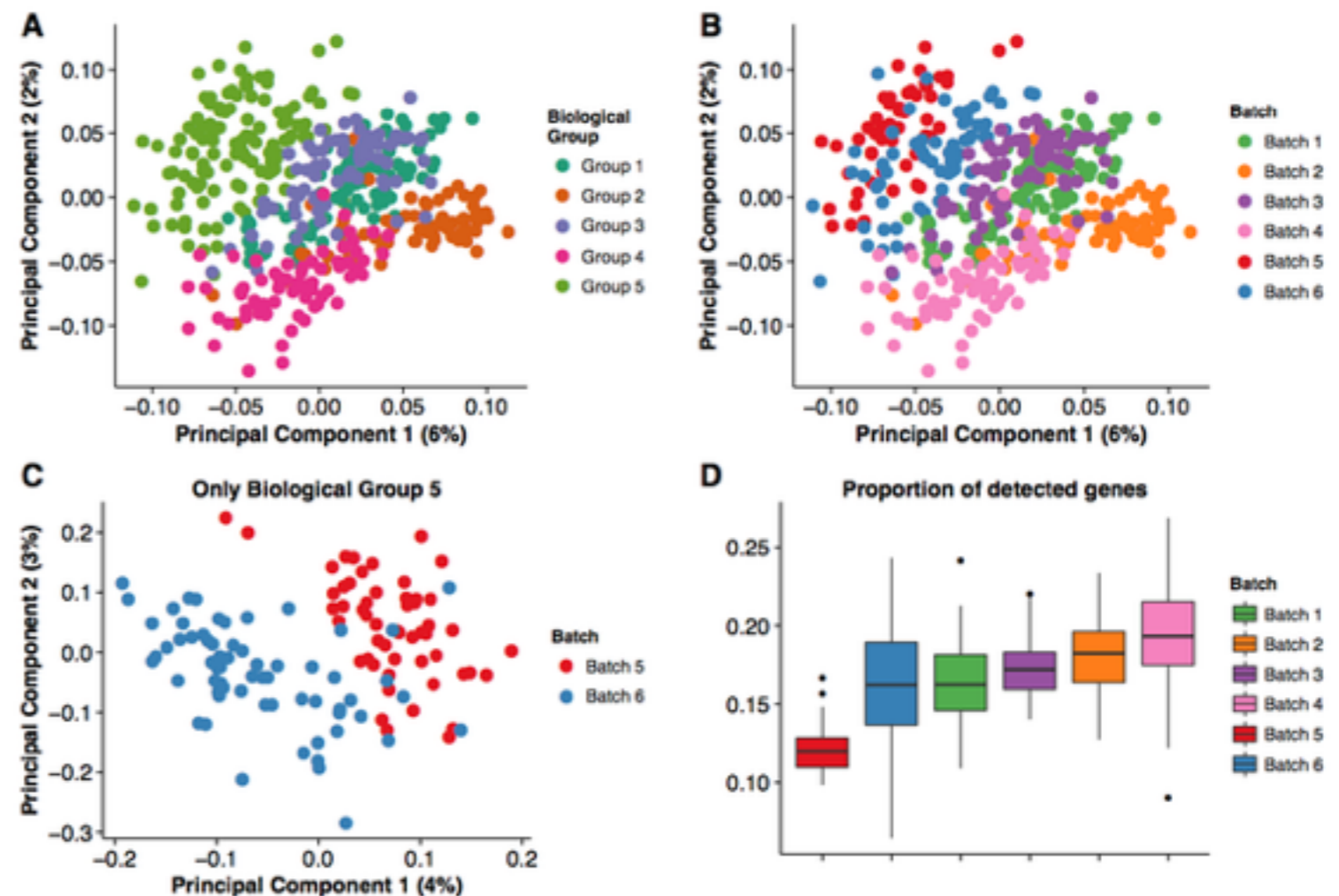
- Major advantage
  - Conceptually very simple
  - Easy to implement → most commonly used technique
- Major weakness of agglomerative clustering methods
  - do not scale well: time complexity of at least  $O(n^2)$ , where  $n$  is the number of total objects
  - can never undo what was done previously → high likelihood of getting stuck in local minima

# Other Challenges



# Batch Effects Occur

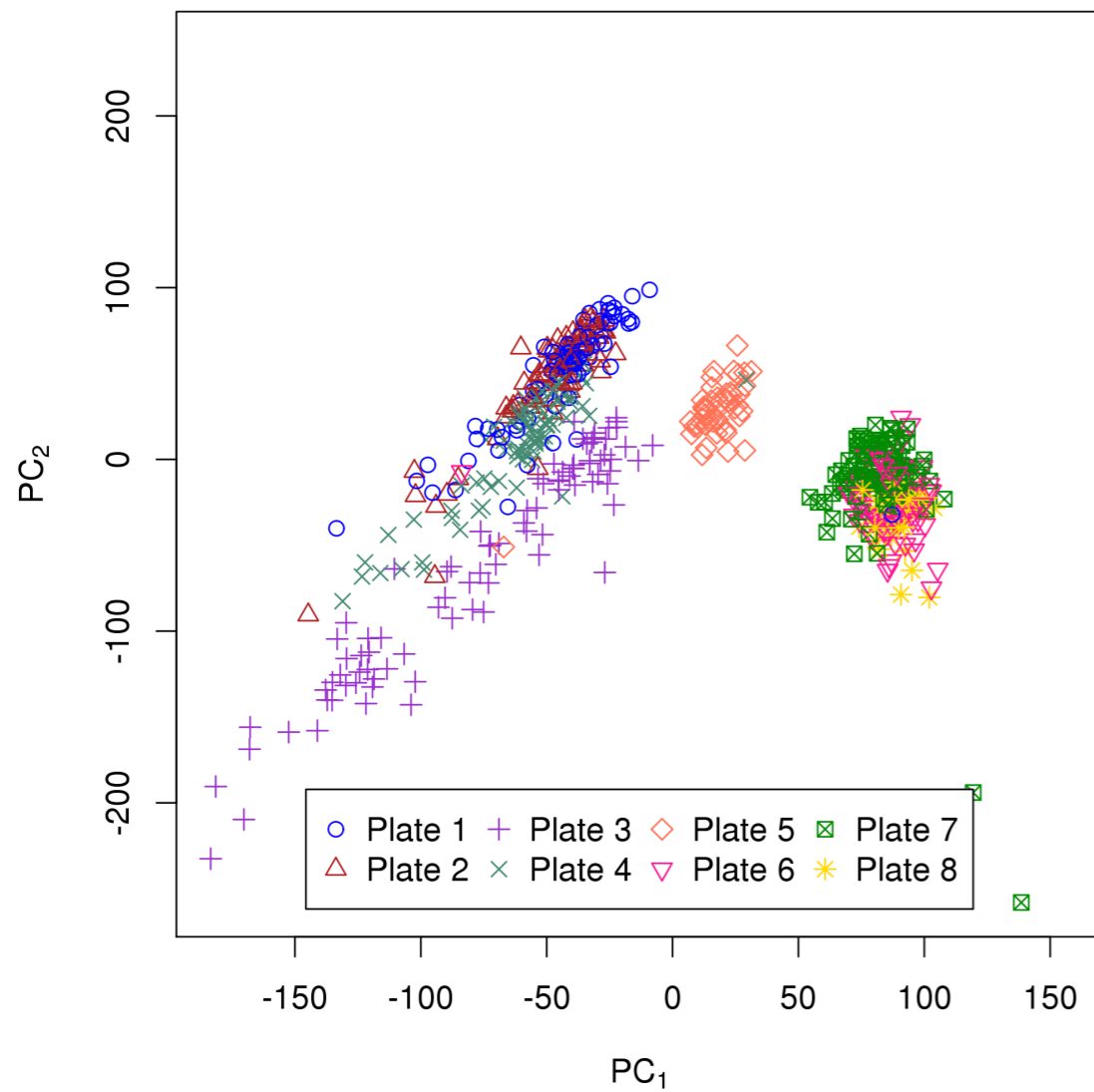
the batch effect represents the systematic technical differences when samples are processed and measured in different batches and which are unrelated to any biological variation recorded



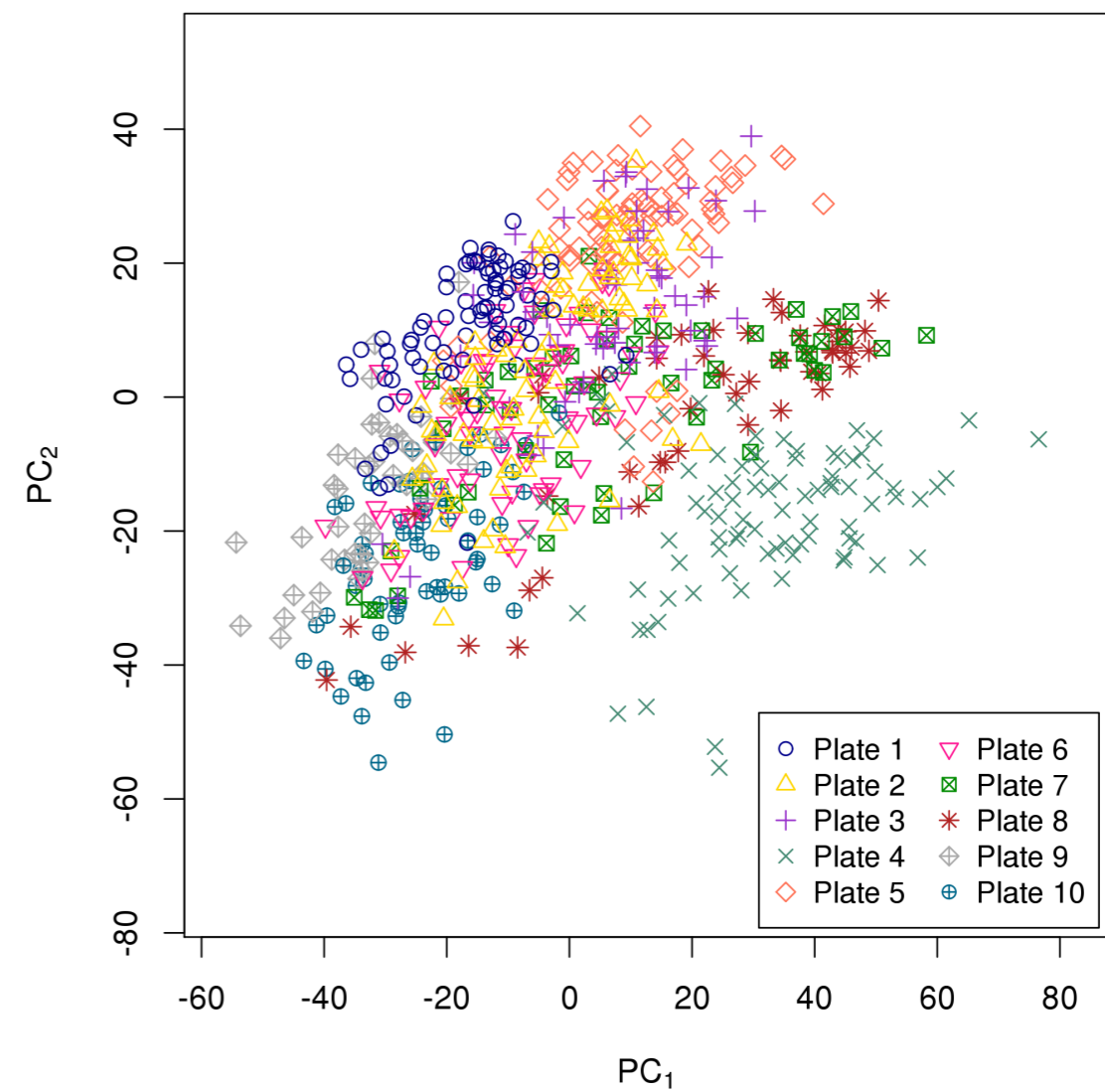
# Batch Effects Occur



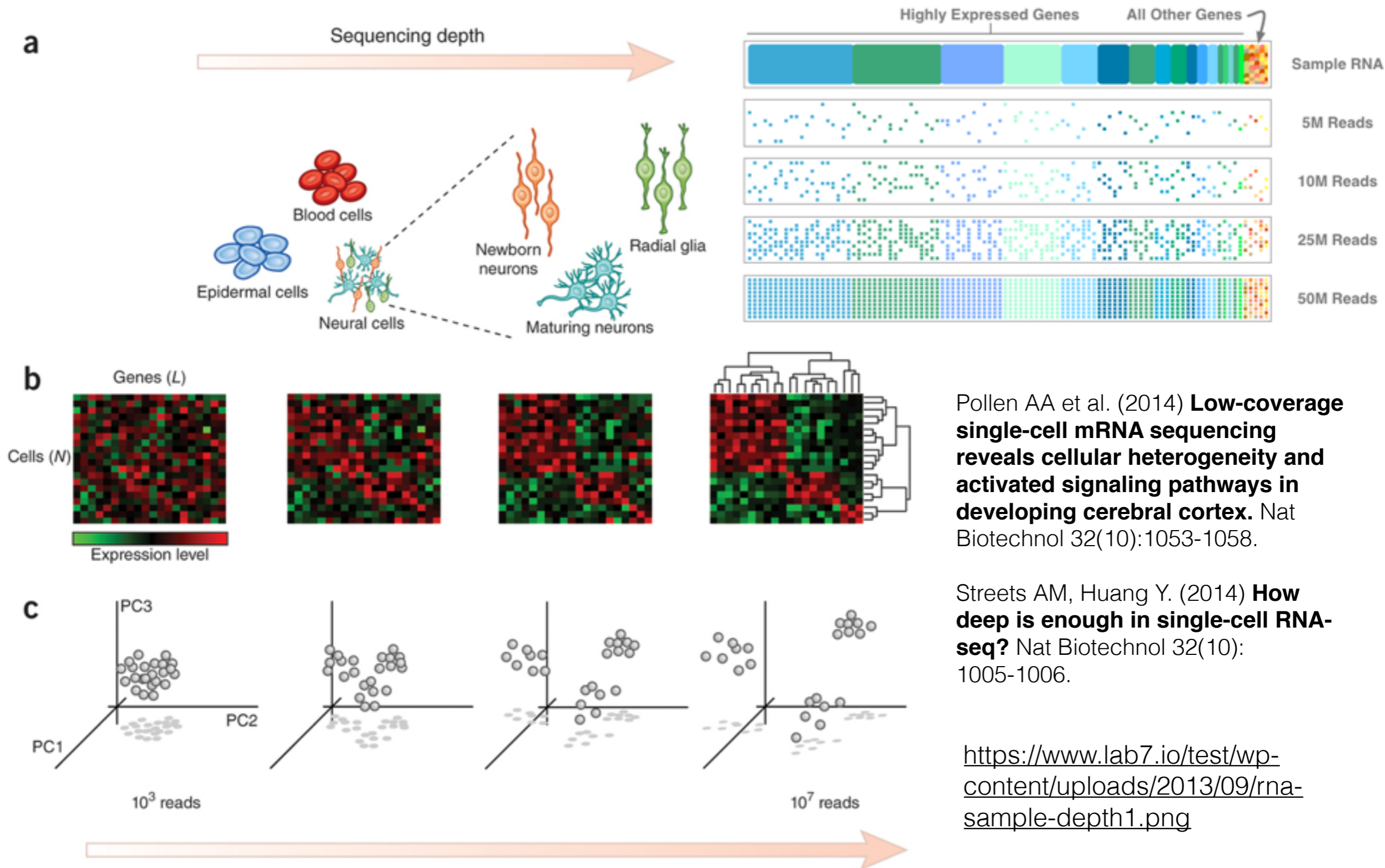
Before batch effect removal



After batch effect removal



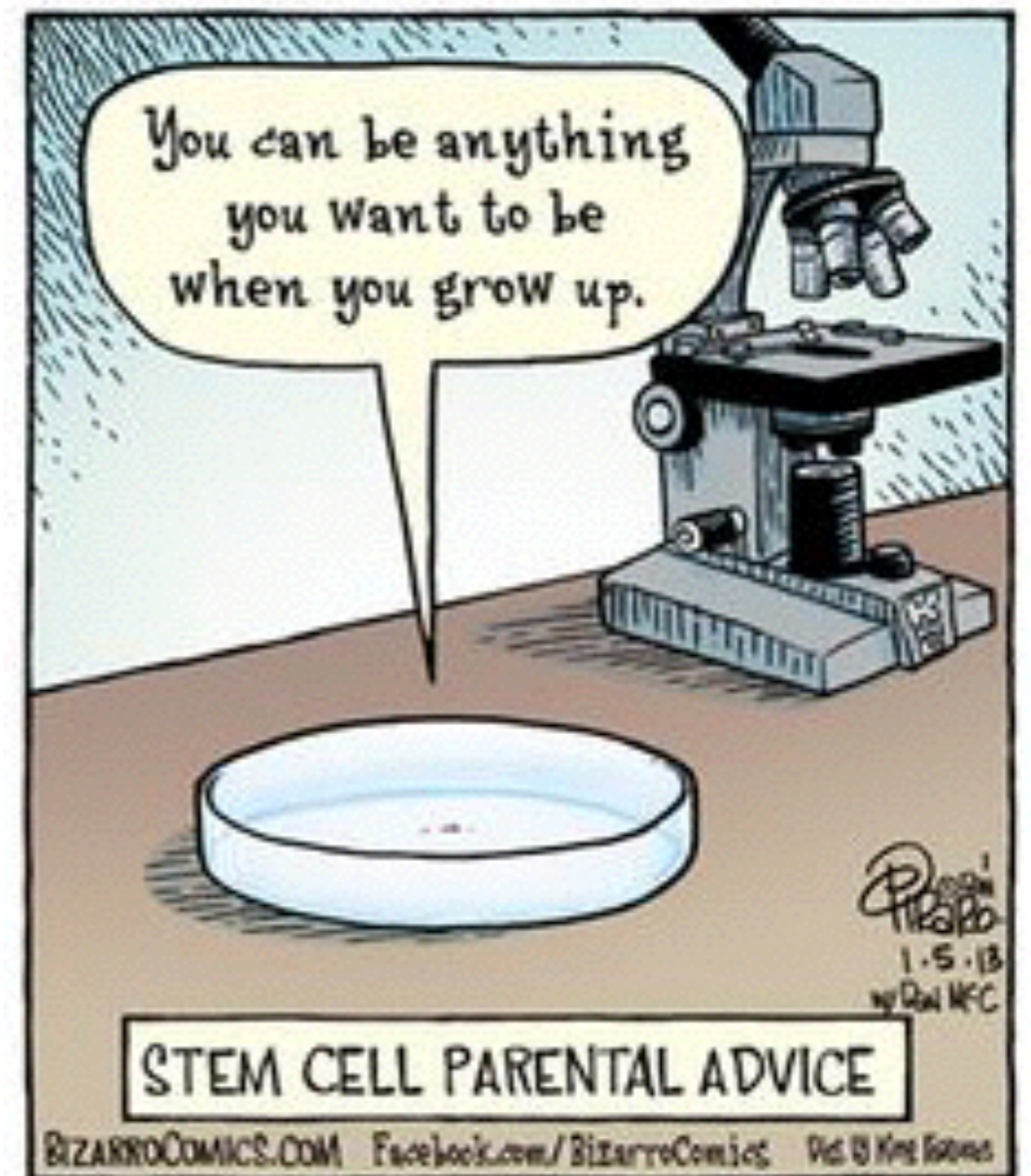
# Sequencing Depth



# Biological Effects



- Cancer: cell lineage
- Metagenomics: cis/trans mechanisms
- Stem Cells: cellular phenotypes
- Immunology: cell type identification
- Neurology: somatic mutations

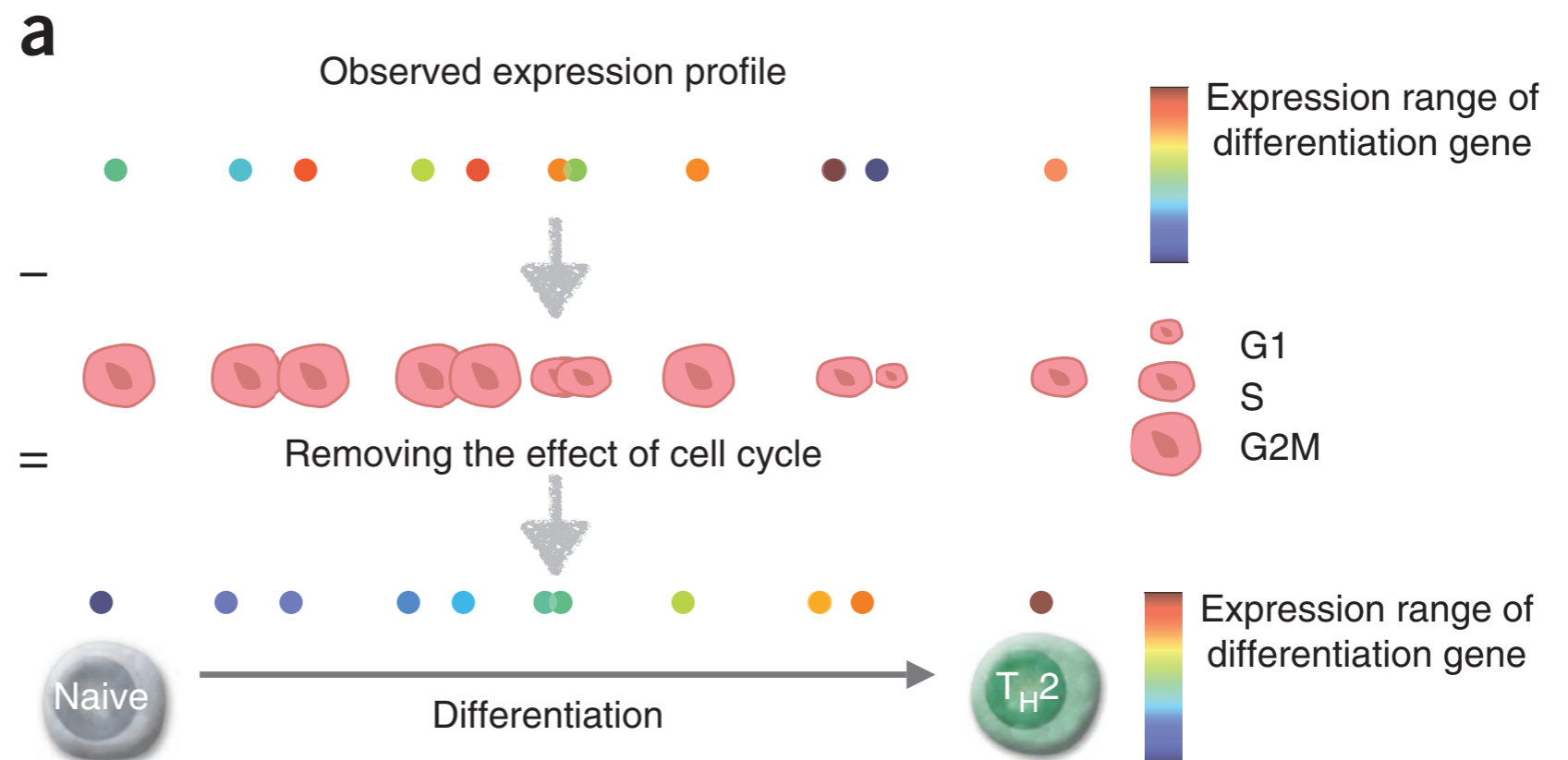
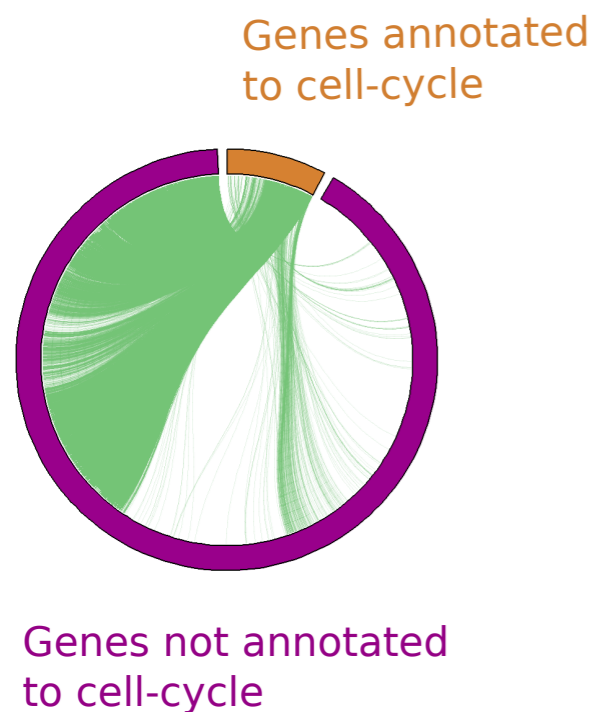




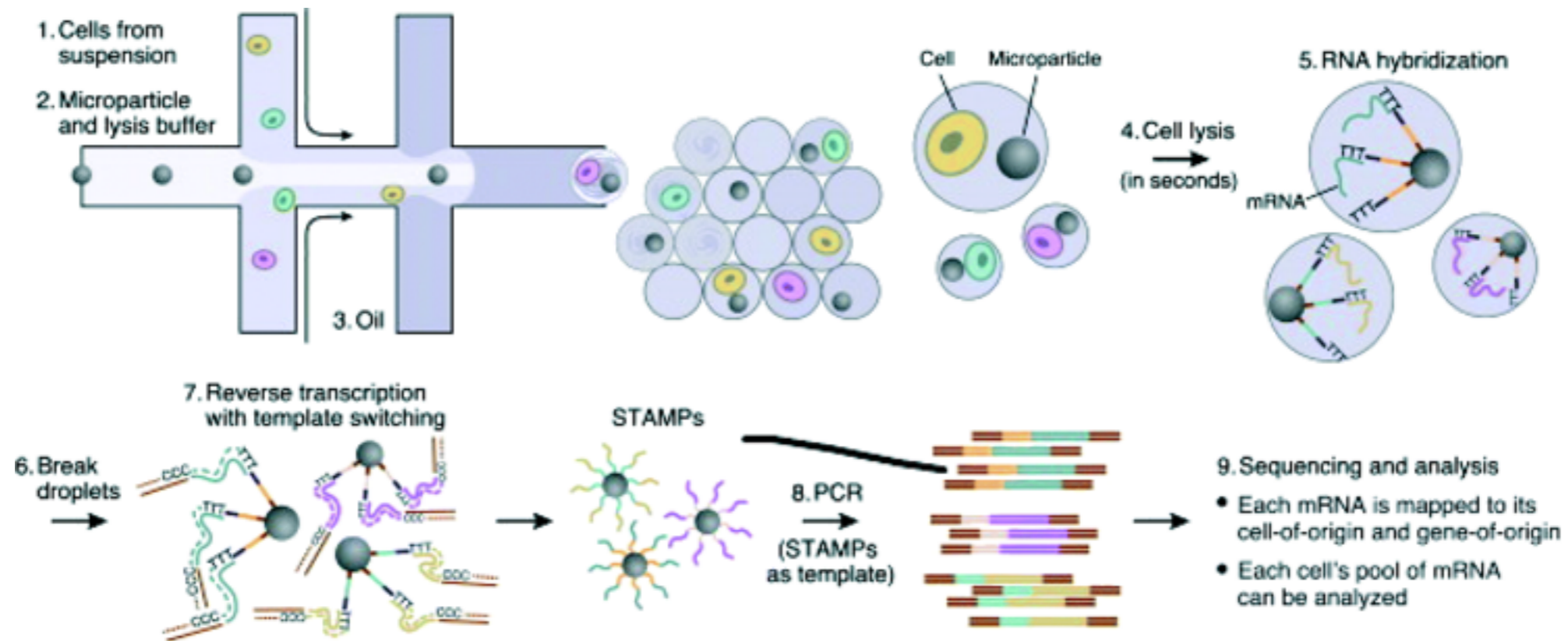
# Example

Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells

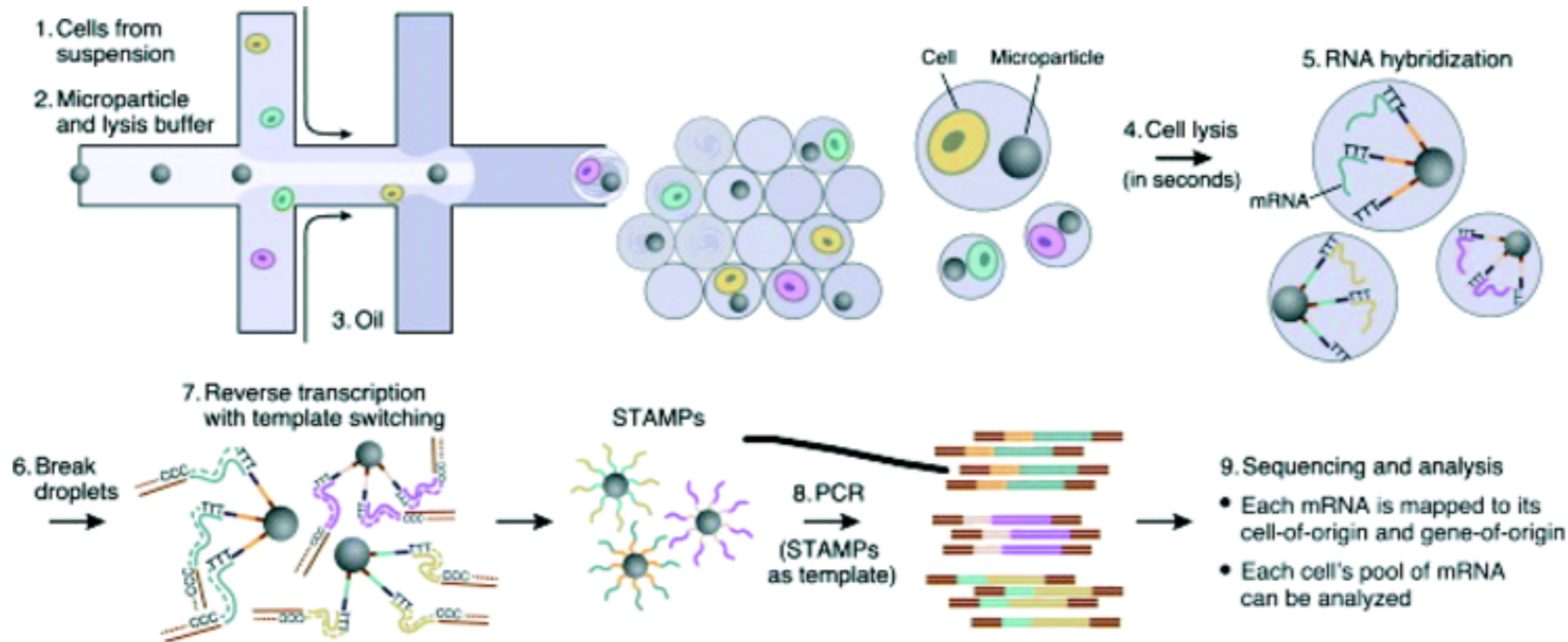
Florian Buettner<sup>1,2,5</sup>, Kedar N Natarajan<sup>2,3,5</sup>, F Paolo Casale<sup>2</sup>, Valentina Proserpio<sup>2,3</sup>, Antonio Scialdone<sup>2,3</sup>, Fabian J Theis<sup>1,4</sup>, Sarah A Teichmann<sup>2,3</sup>, John C Marioni<sup>2,3</sup> & Oliver Stegle<sup>2</sup>



# Full example: DropSeq



# Full example: DropSeq



Cell barcode UMI cDNA (50-bp sequenced)

```

AAATTATGACGATGTGCTTG.....GACTGCAC
CGTTAGATGGCAGGGCCGGG.....CTCATAGT
GACCTACGAGTTAGTTTGTGTA.....GCTCATAA
GTTAAACGTACCCCTAGCTGT.....GATTTTCT
ACGTCACCTTTGTGGGGGT.....ATAAGCTC
TTGCCGTGGTGTATGGAGG.....CCAGCACC
AGTCCATGTCCGGCAGGTTT.....GTTGCCGT
AAATTATGACGAGTTTGTGTA.....AGATGGGG
CCAAAGATGTCTCTAGGCT.....GGGGACGA
GTTAAACGTACCAAGGCTTG.....CAAAGTTC
TTTTTGACCACTCGTAGGG.....TTCCAAGG
ACTGTCCATGCCCTGTGTA.....TGGTACGT
CGTAAACAATAATCCGGTG.....TTAAACCG
.....
.....
.....
    
```

cDNA alignment to genome and group results by cell

```

Cell 1 { TTGCCGTGGTGTGGGGGA.....CGGTGTTA } DDX51
        { TTGCCGTGGTGT TATGGAGG.....CCAGCACC } NOP2
        { TTGCCGTGGTGT TCTCAAGT.....AAAATGGC } ACTB
        { ..... }
Cell 2 { CGTTAGATGGCAAGGGCCGGG.....CTCATAGT } LBR
        { CGTTAGATGGCAACGTTATA.....ACGCGTAC } ODF2
        { CGTTAGATGGCATCGAGATT.....AGCCCTTT } HIF1A
        { ..... }
Cell 3 { AAATTATGACGAGTTTGTGTA.....GGGAATTA } ACTB
        { AAATTATGACGAGTTTGTGTA.....AGATGGGG } RPS15
        { AAATTATGACGAGTTGCTTG.....GACTGCAC }
        { ..... }
Cell 4 { GTTAAACGTACCCCTAGCTGT.....GATTTTCT } GTPBP4
        { GTTAAACGTACCCAGAGAGT.....GTTGCCGT } GAPDH
        { GTTAAACGTACCAAGGCTTG.....CAAAGTTC } ARL1
        { GTTAAACGTACCTCCGGTC.....TCCAGTCG }
        { ..... }
        { ..... }
    
```

Count unique UMIs for each gene in each cell

→

Create digital expression matrix

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

(Hundreds of millions of reads)

(Thousands of cells)