



# BIO334 Practical Bioinformatics

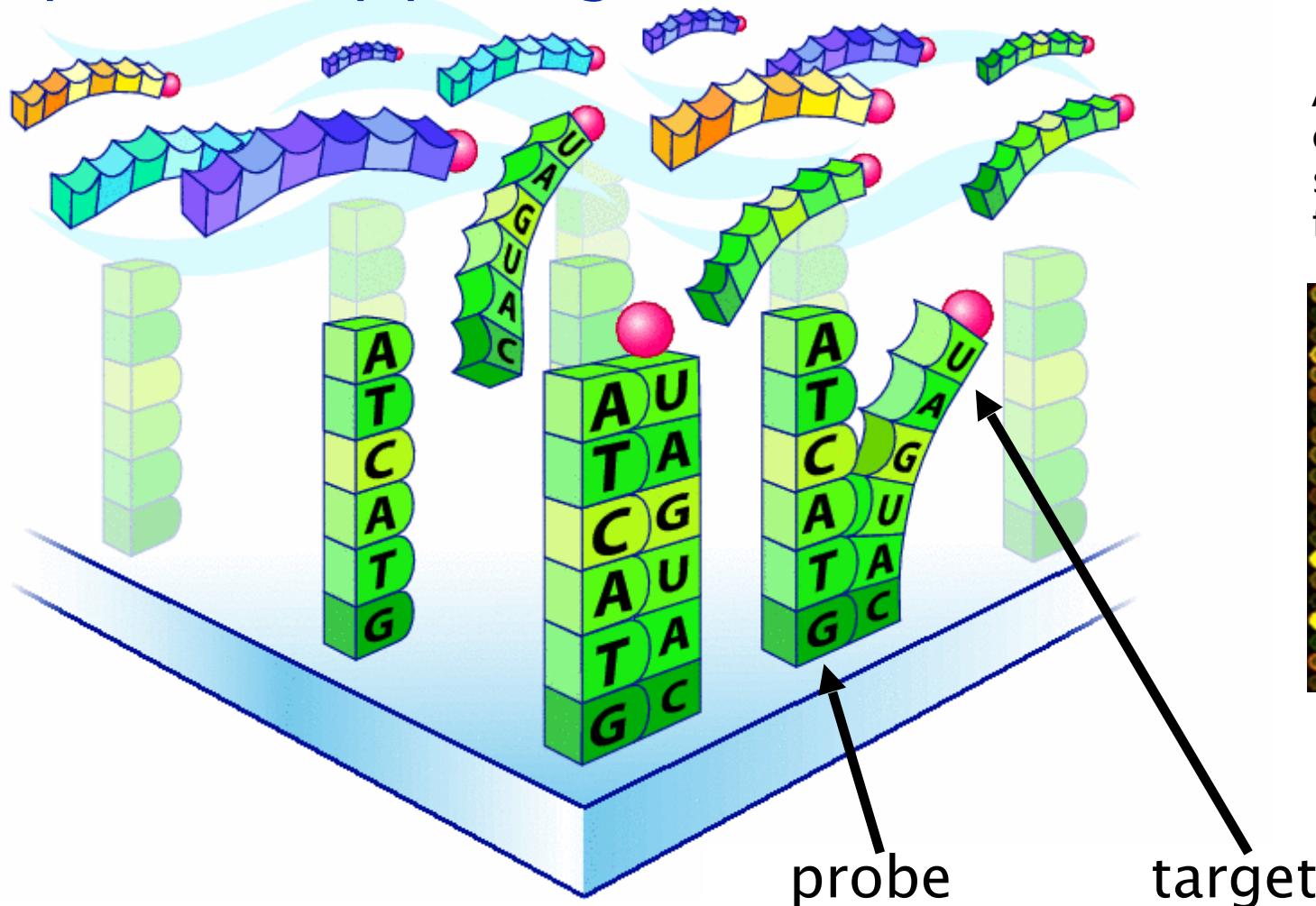
Technologies, Transcriptomic/Proteomic Data  
Analyses (non-exhaustive)

# Technologies in our research area

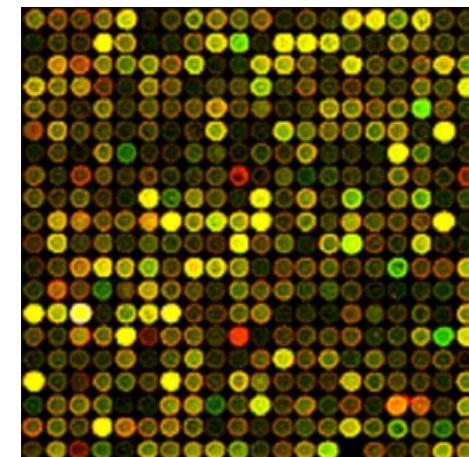
high-throughput sequencing, single cell, cytometry, etc.

“it's just data”

## DNA microarray: parallel northern blots; Nature gives a complementary pairing

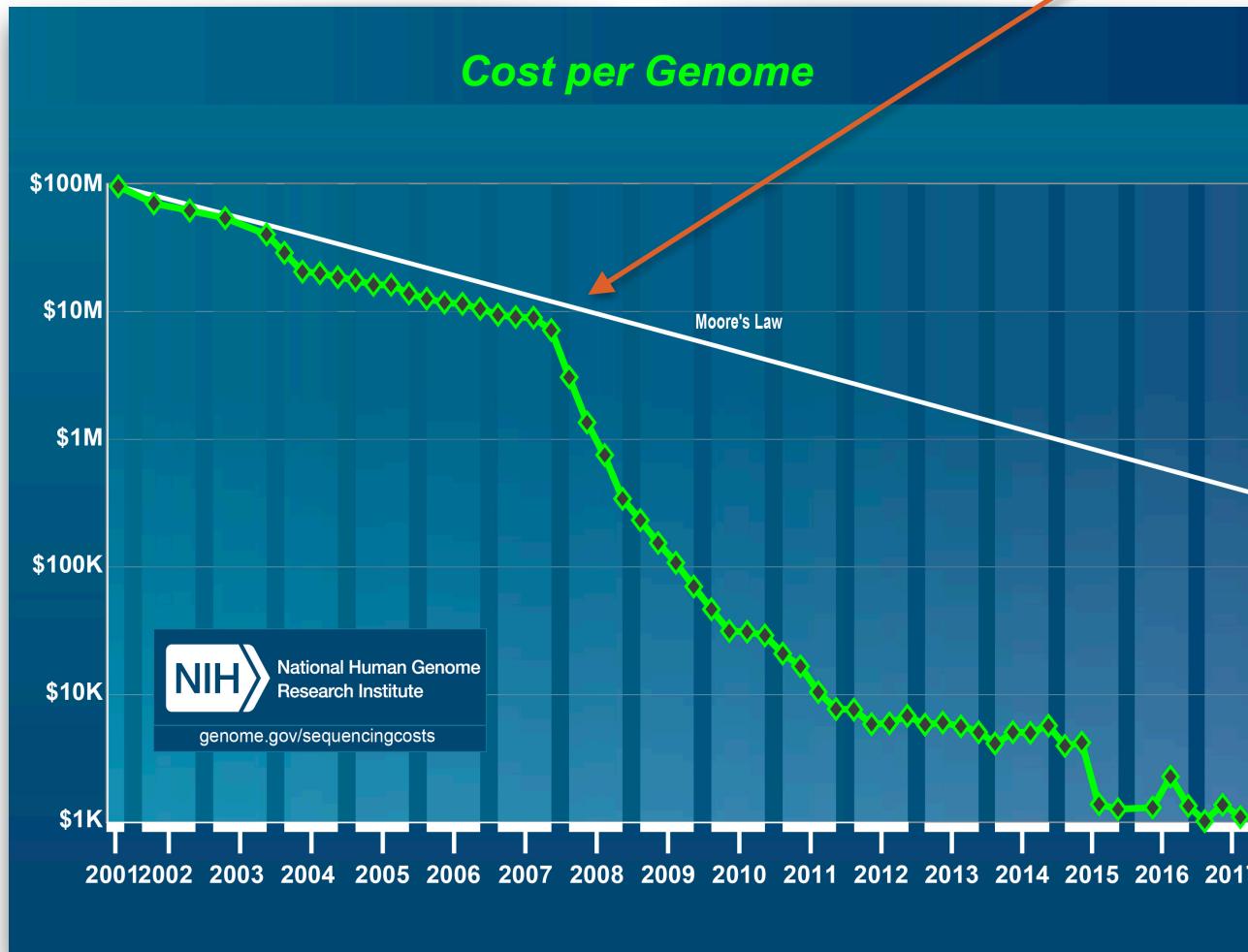


Abundance (of complementary DNA species) measured by fluorescence intensity



## High-throughput sequencing

(Solexa) Illumina

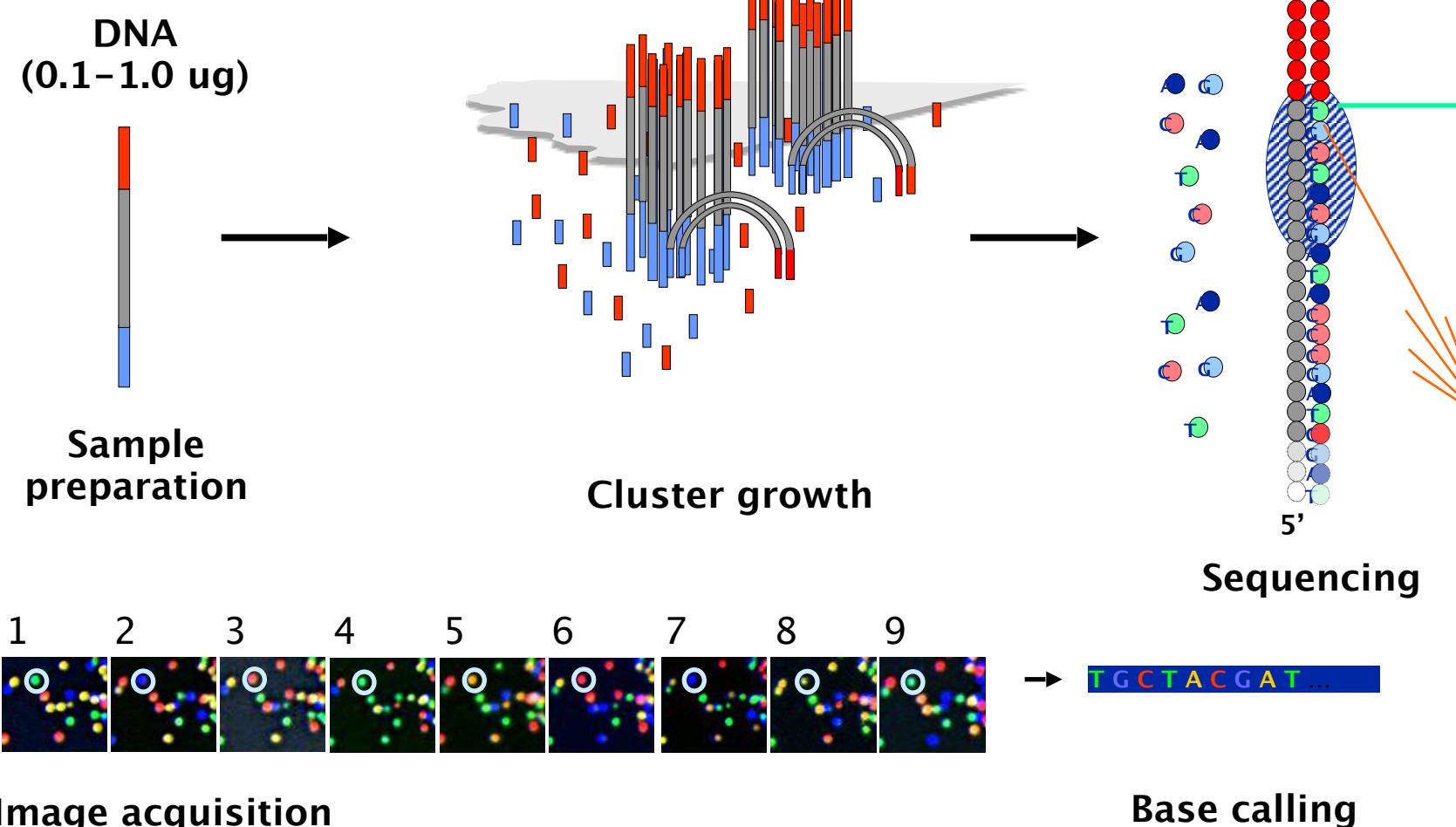


<https://www.statnews.com/2017/01/09/illumina-ushering-in-the-100-genome/>





# Illumina Sequencing Technology



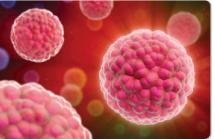


# Applications of high-throughput sequencing

## Common Sequencing Applications

Cancer Research

NGS-based sequencing enables cancer researchers to detect rare somatic variants, tumor subclones, and circulating DNA fragments. [Learn more about sequencing for cancer research.](#)



Microbiology Research

From environmental metagenomics studies to infectious disease surveillance and more, NGS-based sequencing can help researchers gain genetic insight into bacteria and viruses. [Learn more about microbial genomics.](#)



Complex Disease Research

Illumina sequencing is introducing new avenues for understanding immunological, neurological, and other complex disorders on a molecular level. [Learn more about complex disease genomics.](#)



Reproductive and Genetic Health

Illumina sequencing and array technologies deliver fast, accurate information that can guide choices along the reproductive and genetic health journey. [Find reproductive and genetic health solutions.](#)



**ETH zürich**  University of  
Zurich<sup>UZH</sup>

**Functional Genomics Center Zurich**

About us | Working with us | OMICS areas | Education | Research & Publications | FAQ | News & Events

ETH Zurich > UZH > FGCZ

**Services**

Proteomics/Protein analysis services

**Genomics/Transcriptomics services**

Metabolomics/Biophysics services

User Lab Access

Collaboration

FGCZ Policies

Job Offers

FGCZ Terms and Conditions

**Genomics/Transcriptomics services**

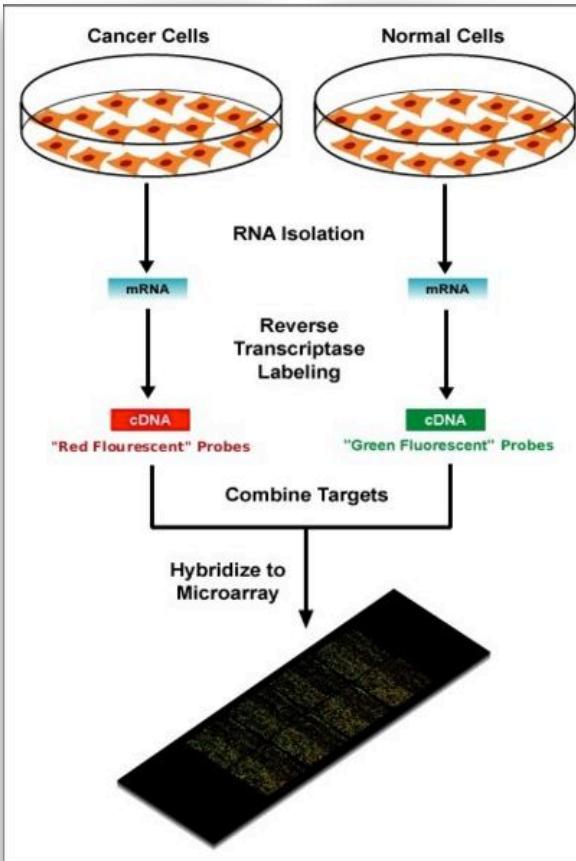
All services in Genomics/Transcriptomics require a project submission via B-Fabric, our project management system.

If you have specific questions about our Genomics/Transcriptomics services please refer to our [FAQ](#) section; alternatively, or in case you would like to request a quote, please do not hesitate to get in touch with our sequencing team at [sequencing@fgcz.ethz.ch](mailto:sequencing@fgcz.ethz.ch)

Application Group	Application	Order via B-Fabric
DNA sequencing	Whole Exome Sequencing	<a href="#">Project</a>
DNA sequencing	Methylation Profiling	<a href="#">Project</a>
DNA sequencing	ChIP-Seq	<a href="#">Project</a>
DNA sequencing	Targeted Sequencing and Metagenomics	<a href="#">Project</a>
DNA sequencing	De novo Genome Assembly	<a href="#">Project</a>
DNA sequencing	Whole Genome Resequencing	<a href="#">Project</a>
RNA sequencing	Transcriptome Profiling	<a href="#">Project</a>
RNA sequencing	Small RNA Profiling	<a href="#">Project</a>
RNA sequencing	De novo Transcriptome Assembly	<a href="#">Project</a>

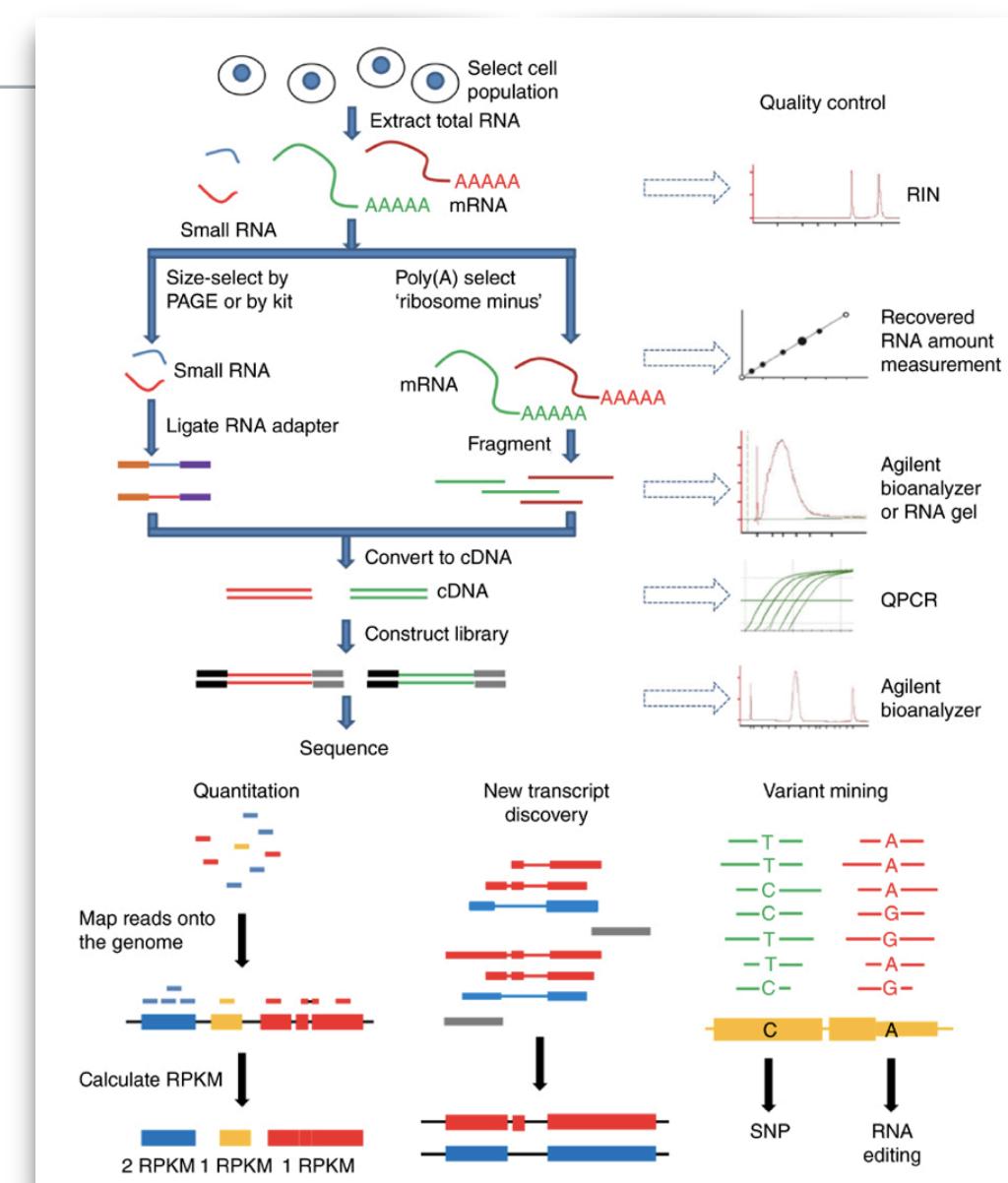


## Abundance by Fluorescence Intensity (DNA microarray)



[http://en.wikipedia.org/wiki/DNA\\_microarray](http://en.wikipedia.org/wiki/DNA_microarray)

## Abundance by Counting (RNA-seq)



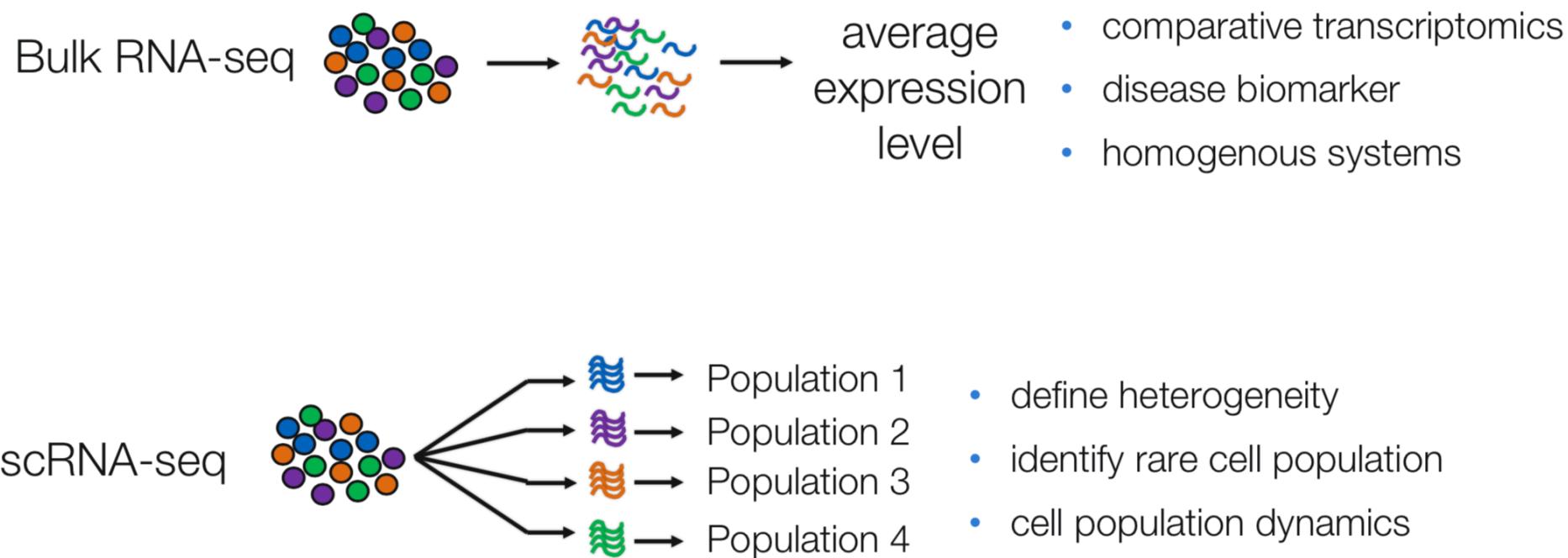
Zeng & Mortazavi, Nature Immunology, 2012



## Gene Expression Profiling: questions of interest

- What genes have changed in expression? (e.g. between disease/normal, affected by treatment)  
*Gene discovery, differential expression*
- Is a specified group of genes all up-regulated in a particular condition?  
*Gene set differential expression*
- Can the expression profile predict outcome?  
*Class prediction, classification*
- Are there tumour sub-types not previously identified? Do my genes group into previously undiscovered pathways?  
*Class discovery, clustering*
- (with single cell gene expression): What changes in cell type composition are observed? What genes have changed in expression in a given subtype of cells?

## Bulk vs Single Cell RNA Sequencing (scRNA-seq)





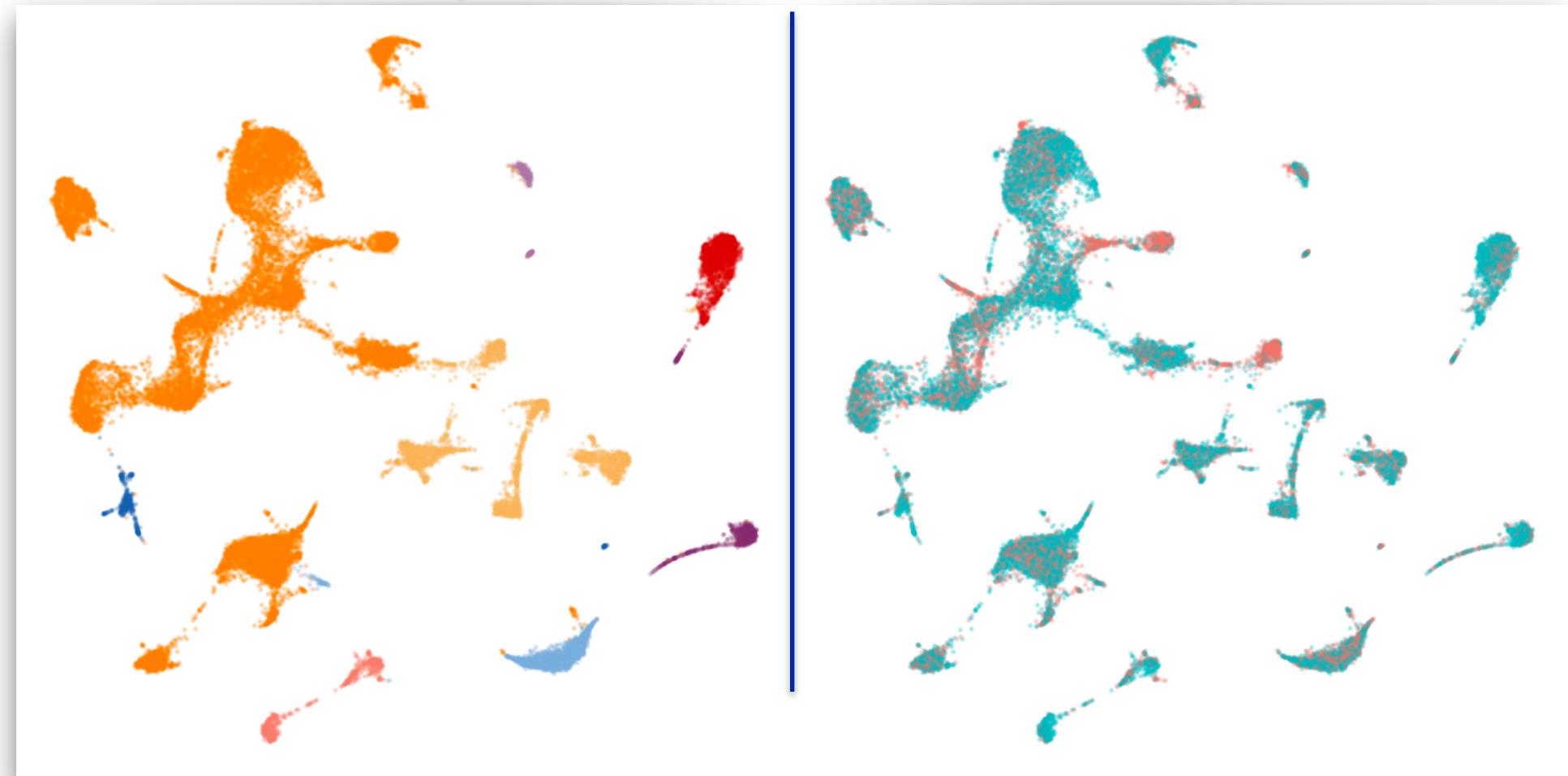
Motivation: Single-cell RNA-seq: finding cell subpopulation-specific changes in state

frontal cortex

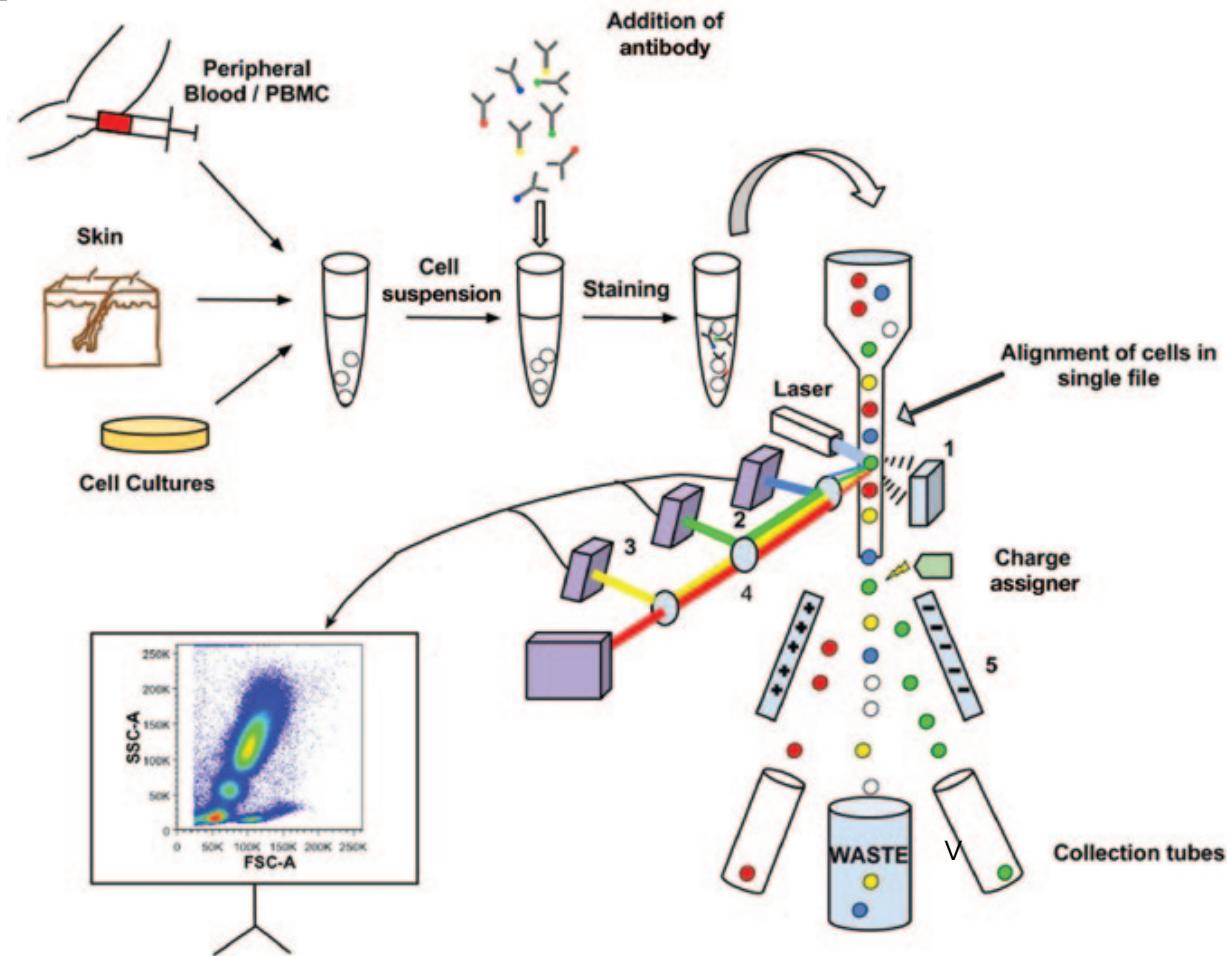
single nuclei RNA-seq  
(10x)

Data from:  
4 mice vehicle treated  
4 mice LPS treated

Each dot is one cell



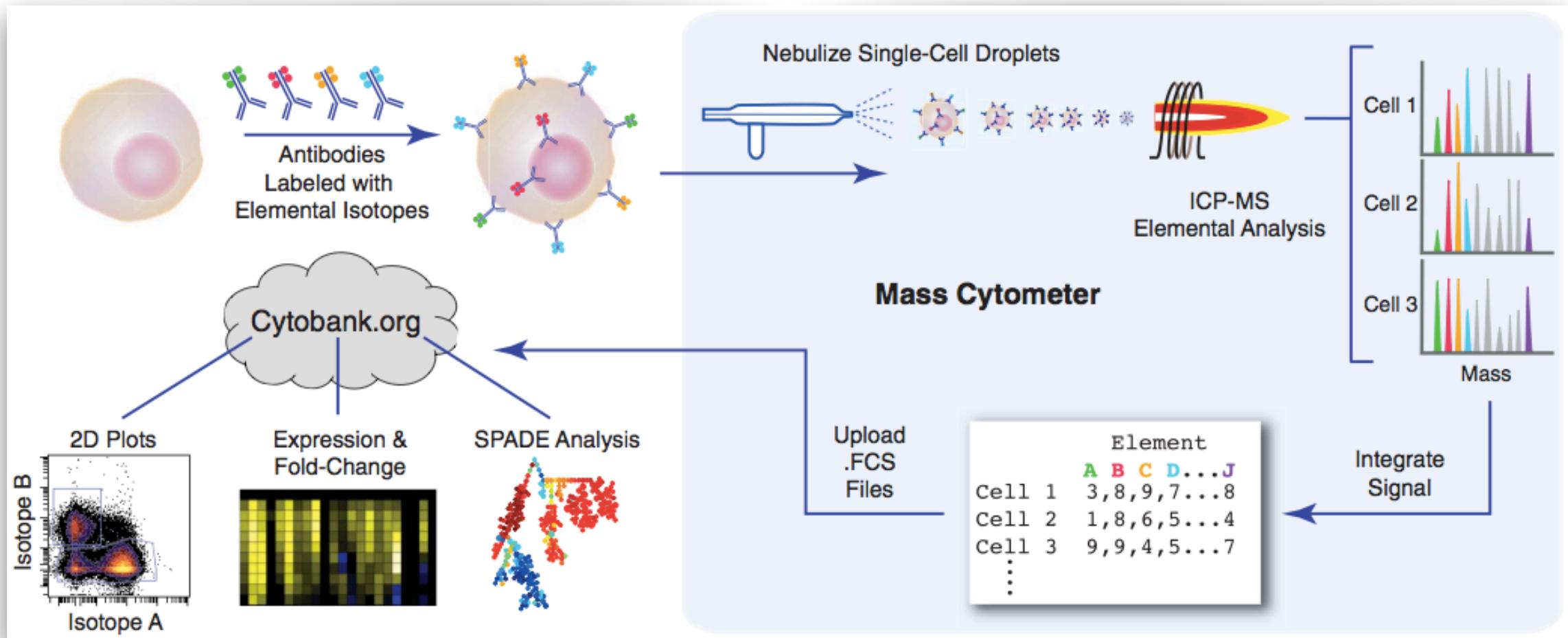
## Flow cytometry



**Figure 1. Schematic representation of a flow cytometer.** For details please see text. (1) Forward-scatter detector, (2) side-scatter detector, (3) fluorescence detector, (4) filters and mirrors, and (5) charged deflection plates.



## Mass cytometry

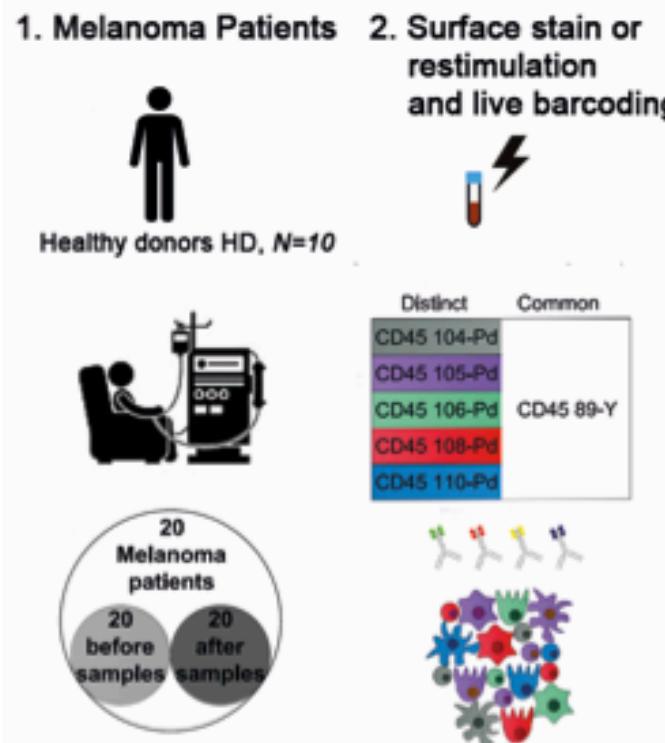


Bendall et al. (2011), Fig. 1A



## Finding molecular biomarkers associated with drug response

### A Workflow



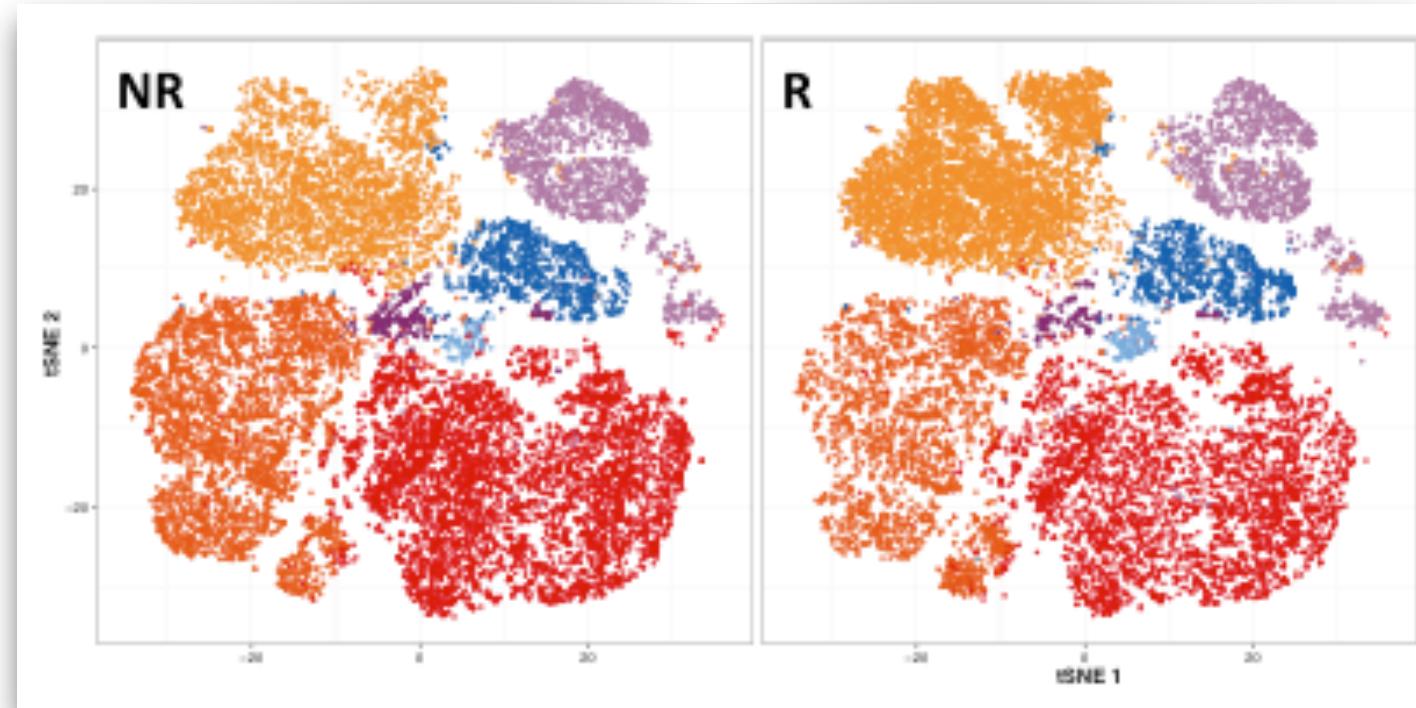
High-dimensional single-cell analysis predicts response to anti-PD-1 immunotherapy

Carsten Krieg<sup>1,6</sup> , Malgorzata Nowicka<sup>2,3</sup>, Silvia Guglietta<sup>4</sup>, Sabrina Schindler<sup>5</sup>, Felix J Hartmann<sup>1</sup> , Lukas M Weber<sup>2,3</sup> , Reinhard Dummer<sup>5</sup>, Mark D Robinson<sup>2,3</sup> , Mitchell P Levesque<sup>5,7</sup>  & Burkhard Becher<sup>1,7</sup> 

## Differential abundance of cell populations

tSNE projection  
(each dot = cell,  
cells from multiple  
patients)

NR: non-responders  
R: responders

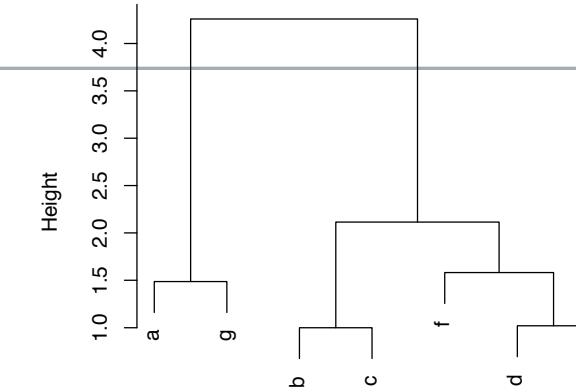


**Under the hood:** Generalized linear mixed model to assess the change in relative abundance of subpopulations.

# Hierarchical Clustering



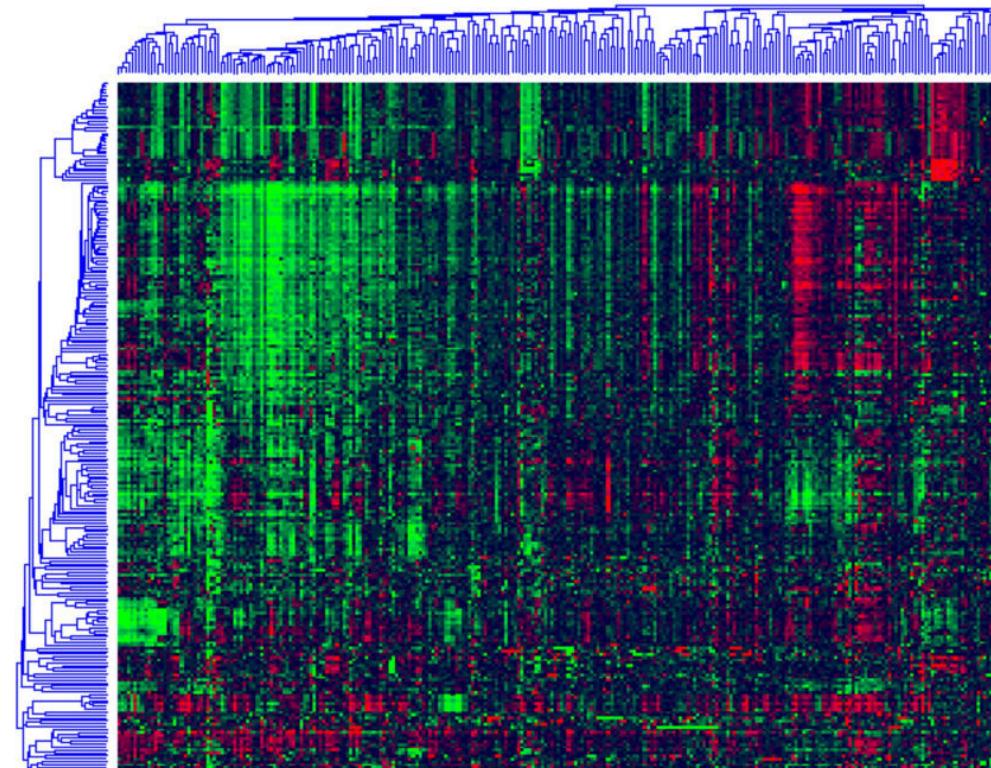
Cluster Dendrogram



Divisive (all features start as 1 cluster, then subsequently split) versus Agglomerative (every feature is it's own cluster, then subsequently merged)

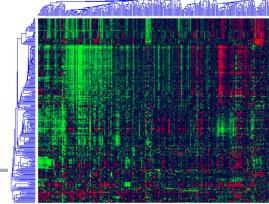
**Metric**: to define how similar any two vectors are.

**Linkage**: determines how clusters are merged into a tree





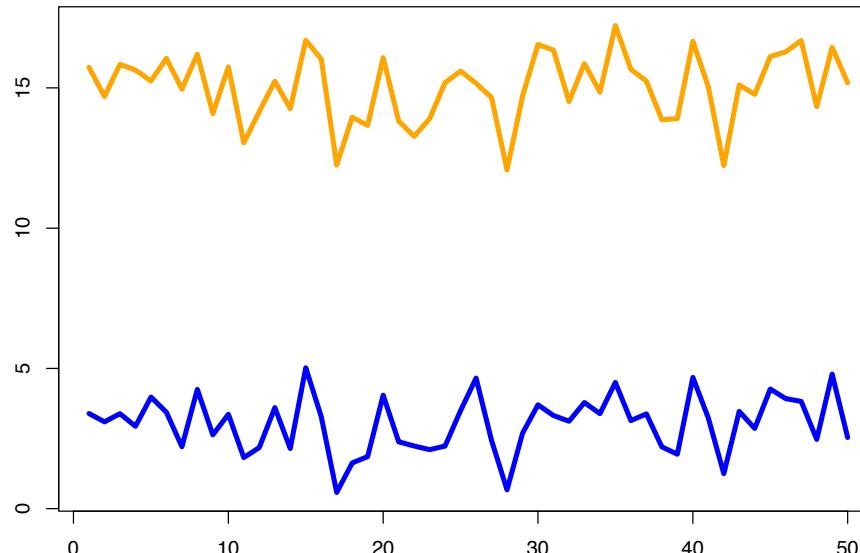
$$d(\mathbf{p}, \mathbf{q}) = d(\mathbf{q}, \mathbf{p}) = \sqrt{(q_1 - p_1)^2 + (q_2 - p_2)^2 + \cdots + (q_n - p_n)^2} = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$$



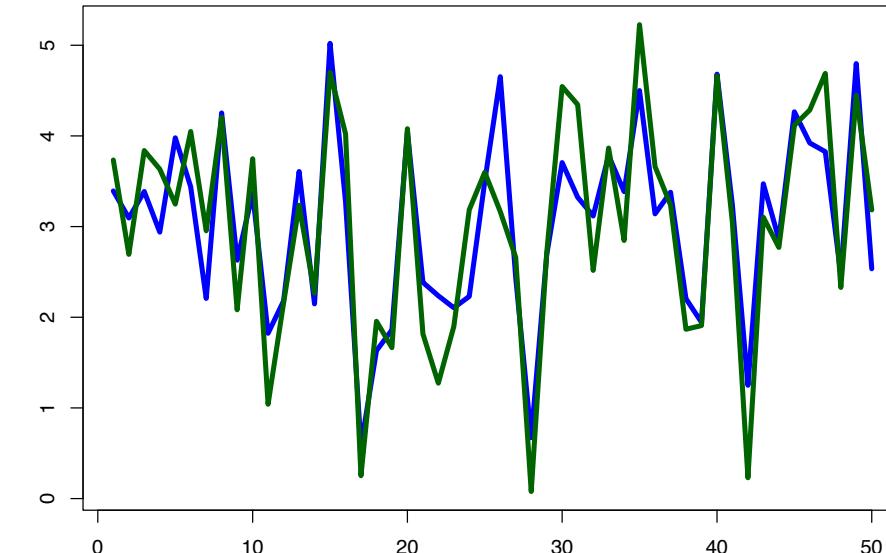
Are these “vectors” similar ?

```
> sqrt(sum((x-(y-12))^2))
[1] 3.926007
> sqrt(sum((x-y)^2))
[1] 84.84028
```

It depends how you define similar.



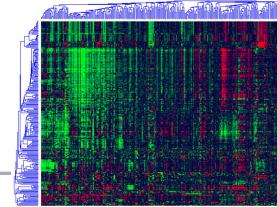
Euclidean distance: 84.84



3.92



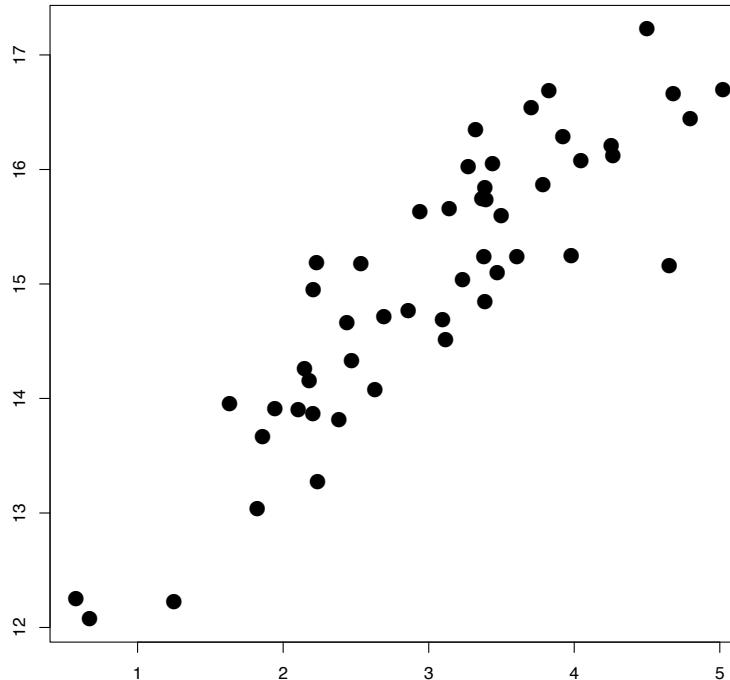
$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1)s_x s_y} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}},$$



Are these “vectors” similar ?

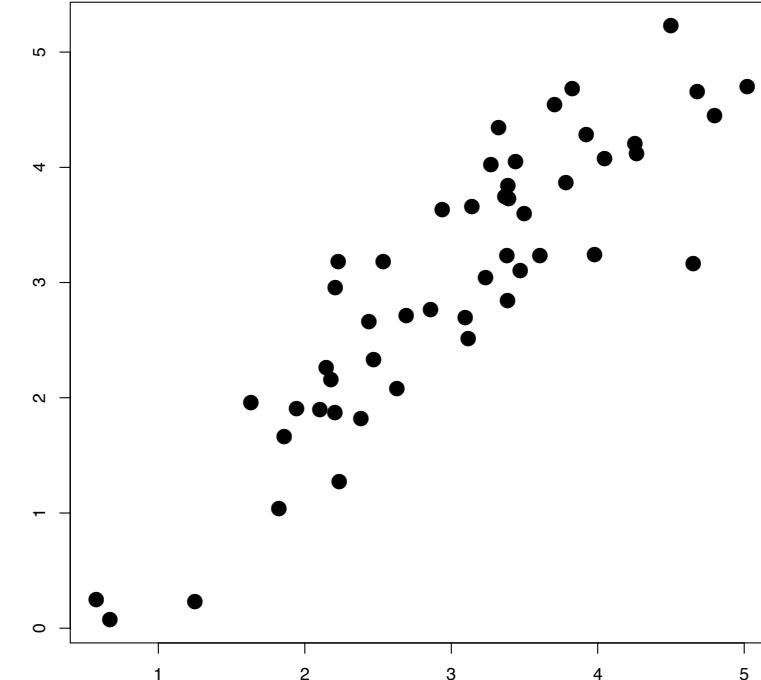
```
> cor(x,y)
[1] 0.8901139
> cor(x,y-12)
[1] 0.8901139
```

It depends how you define similar.

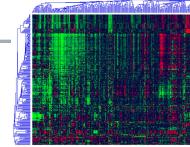


Correlation:

0.89



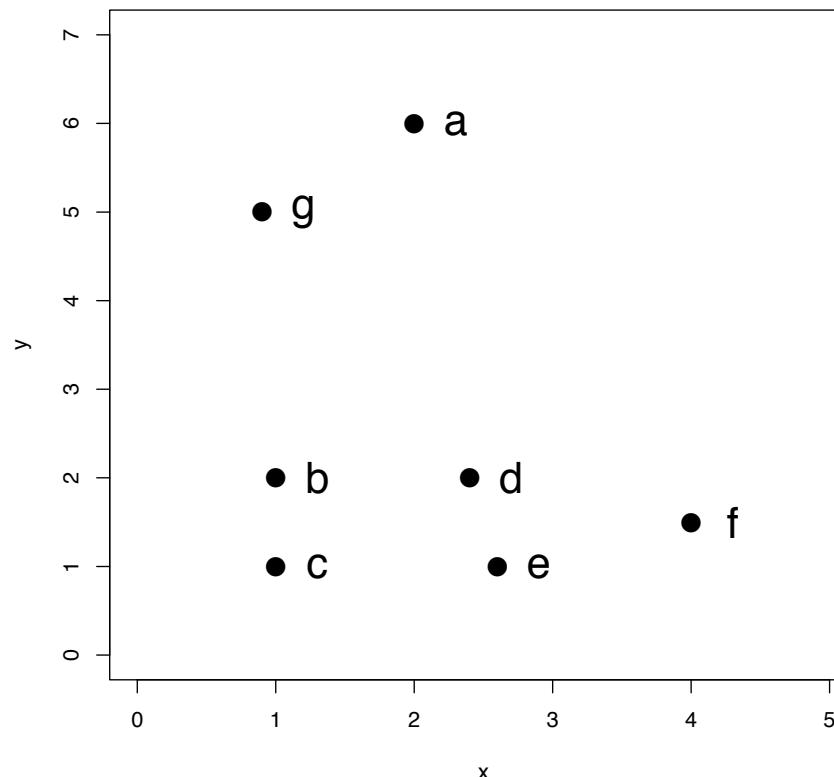
0.89



## Hierarchical (Agglomerative) Clustering

Start with distances.

Linkage: determines how clusters are merged into a tree.

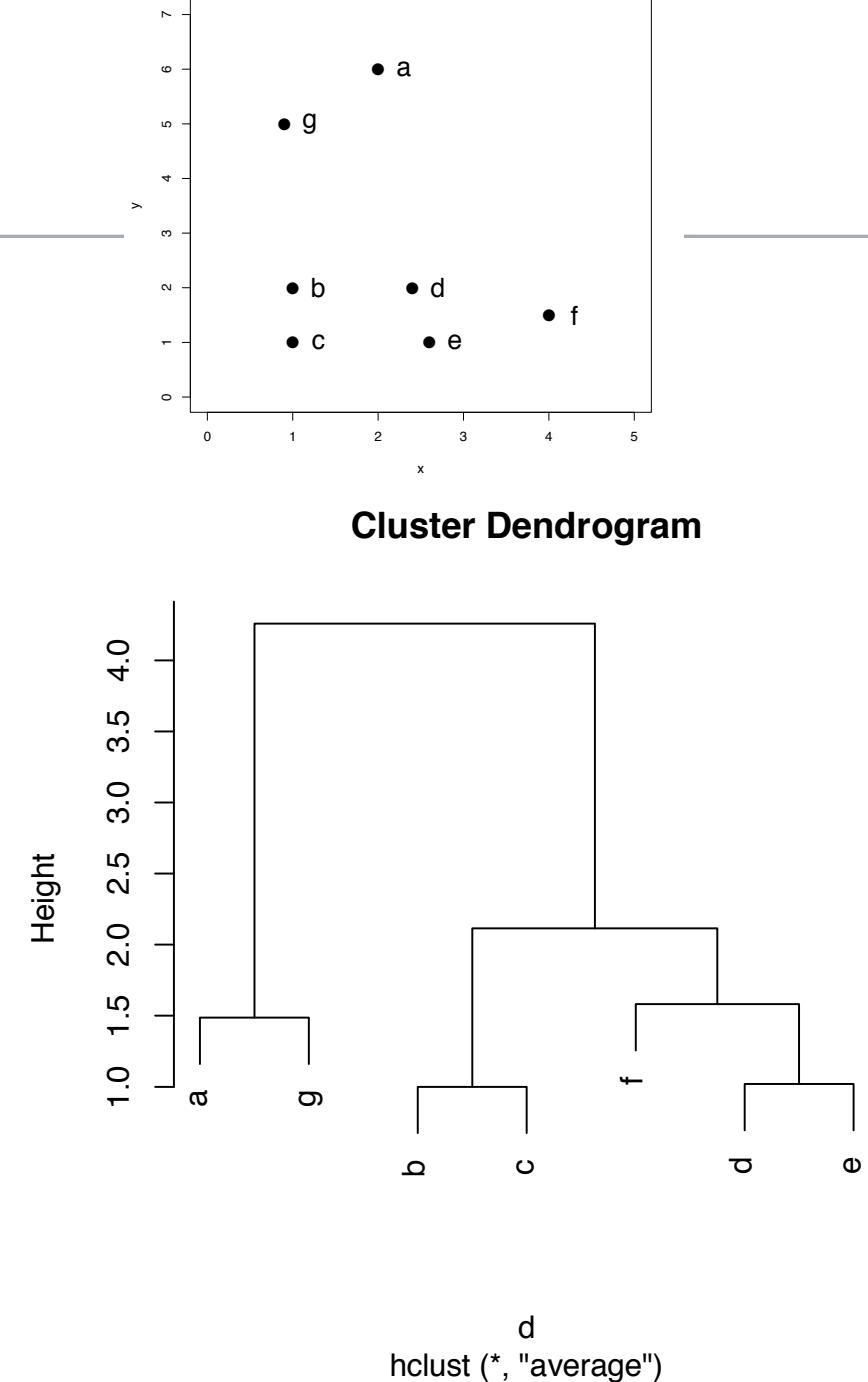
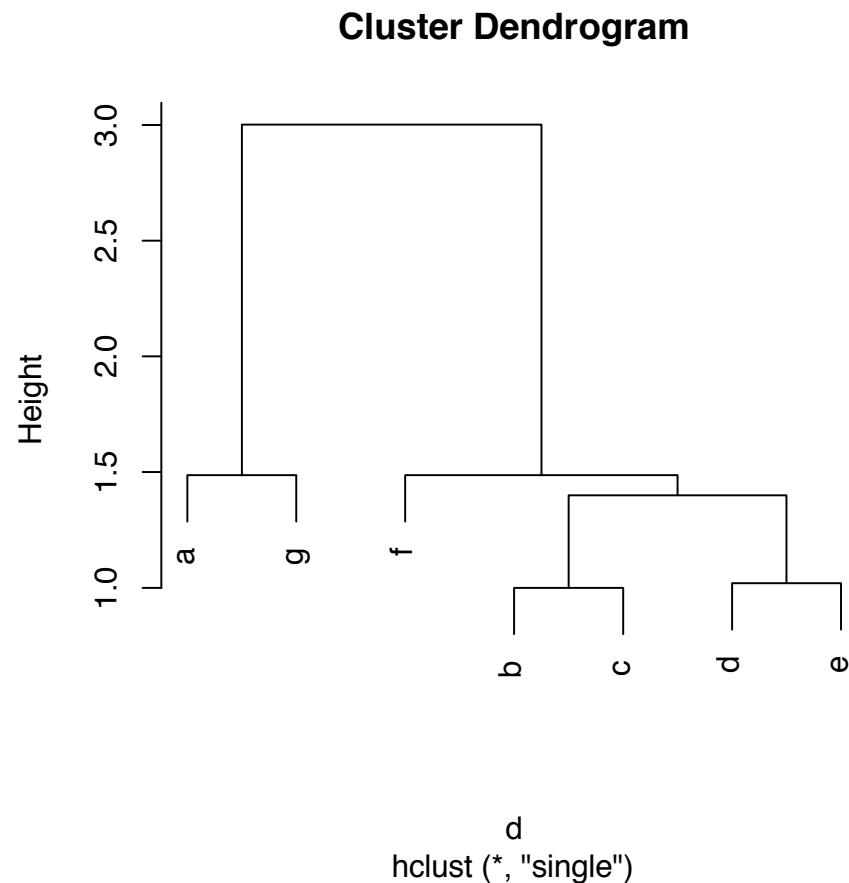


From eyeballing, here is a likely set of merges:

b,c  
d,e  
a,g,  
(d,e),f  
(b,c),((d,e),f)  
ALL



## Different linkages



dimension reduction  
(exploratory data analysis)



“To consult the statistician after an experiment is finished is often merely to ask [them] to conduct a post mortem examination. [They] can perhaps say what the experiment died of.” R. A. Fisher

## Motivation for exploratory data analysis: Case Study

(from Stefano, a former M.Sc. student in my Institute)

He is studying gene expression in fruitfly and is interested in transcriptional responses following “heat shock”.

Basic schematic of experiment:

CTL	t0		t12		
TRT		t4	t12	t24	t72



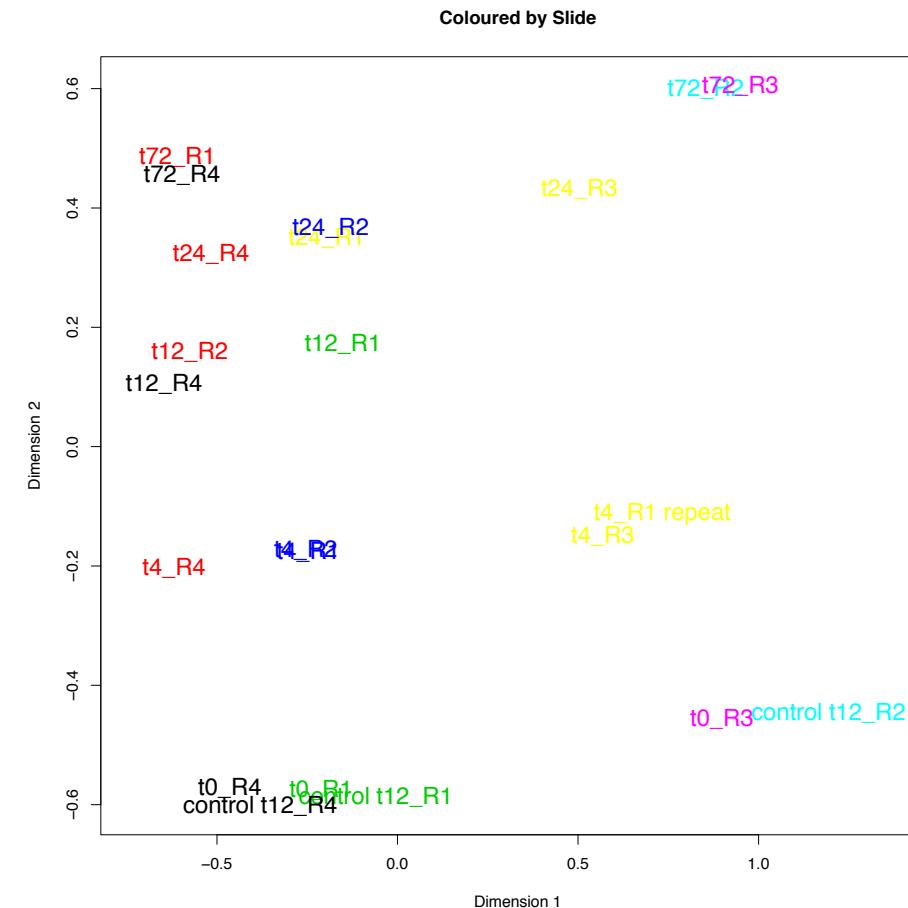
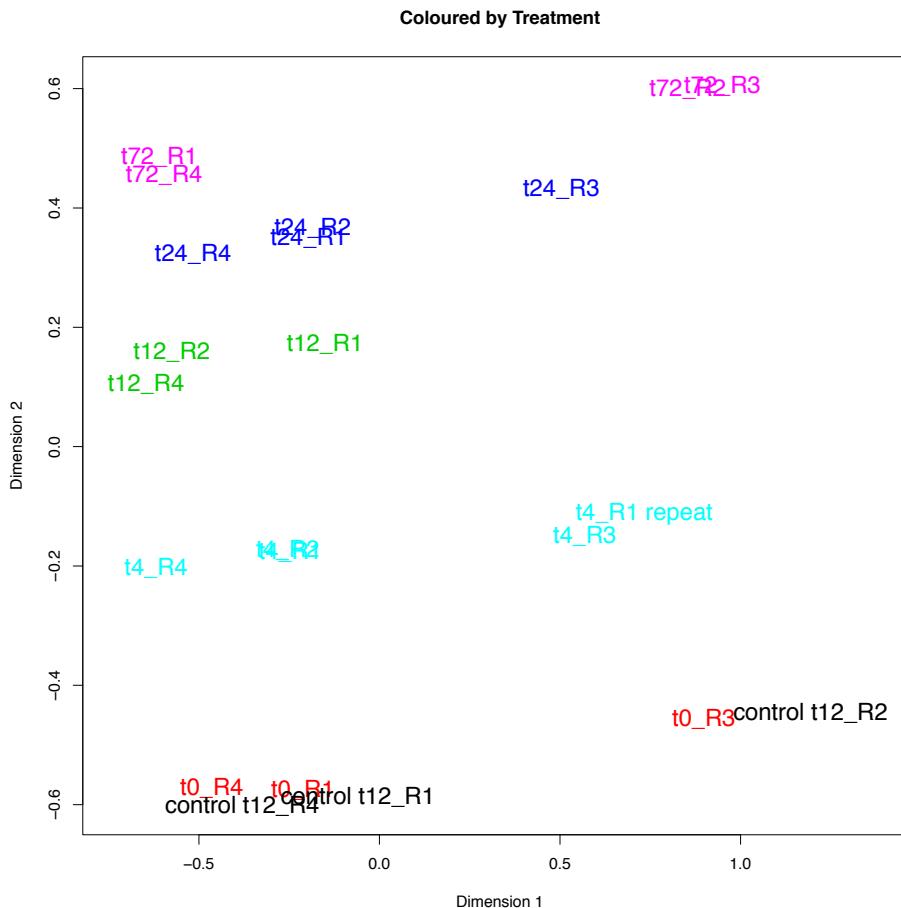
Change to lower  
temperature.

~4 replicates for each condition



```
library(limma)
plotMDS(d) # 'd' is a matrix
"Plot samples on a two-dimensional scatterplot so that
distances on the plot approximate the typical log2 fold
changes between the samples."
```

Take a close look at where the 24 samples are to each other relative to the X- and Y-axes



22 samples x  
~20,000 genes

reduced to 22  
samples x 2  
dimensions



## Magic: Surrogate variable analysis to detect and “remove” batch effects

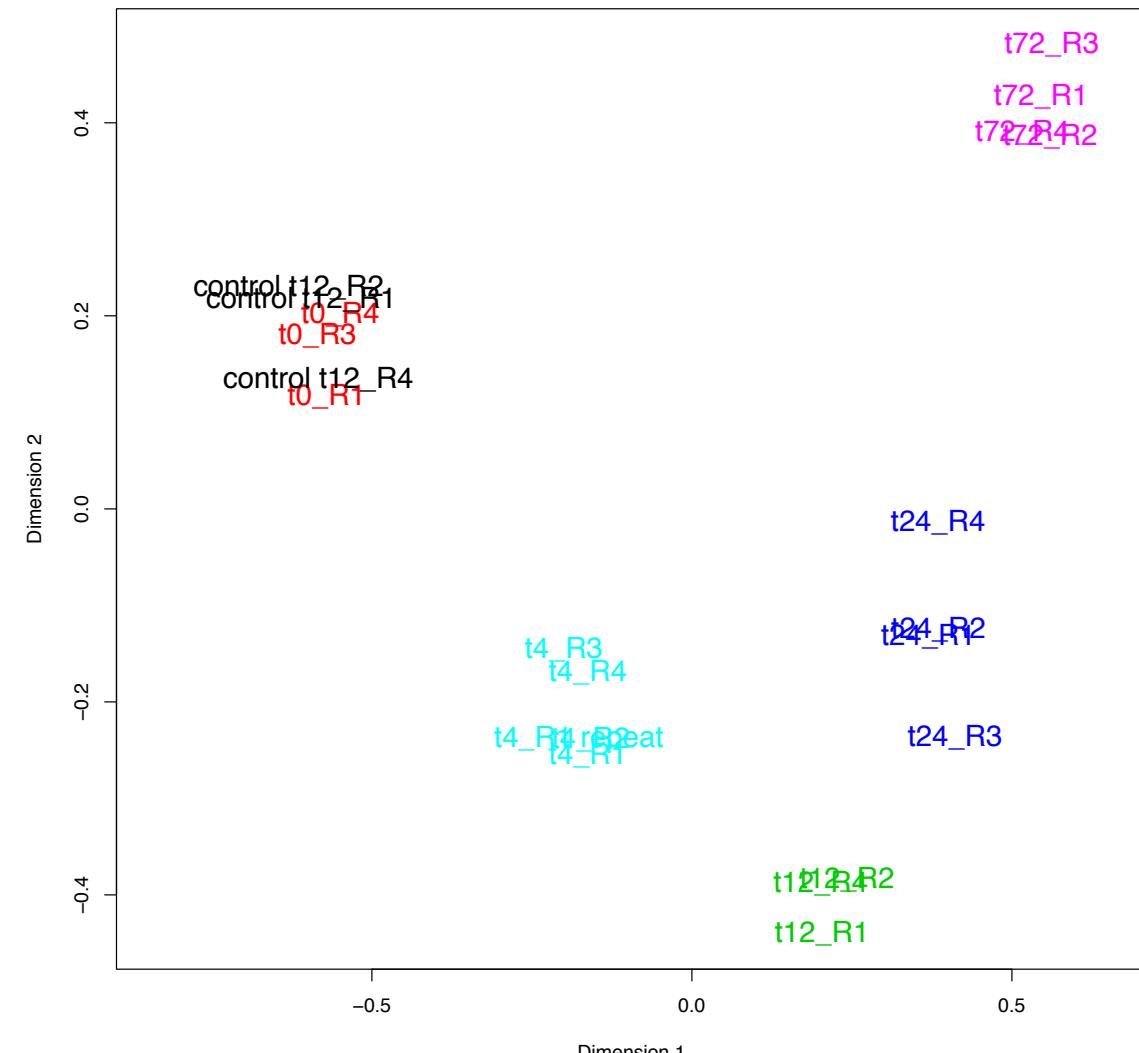
OPEN ACCESS Freely available online

PLOS GENETICS

### Capturing Heterogeneity in Gene Expression Studies by Surrogate Variable Analysis

Jeffrey T. Leek<sup>1</sup>, John D. Storey<sup>1,2\*</sup>

<sup>1</sup> Department of Biostatistics, University of Washington, Seattle, Washington, United States of America, <sup>2</sup> Department of Genome Sciences, University of Washington, Seattle, Washington, United States of America

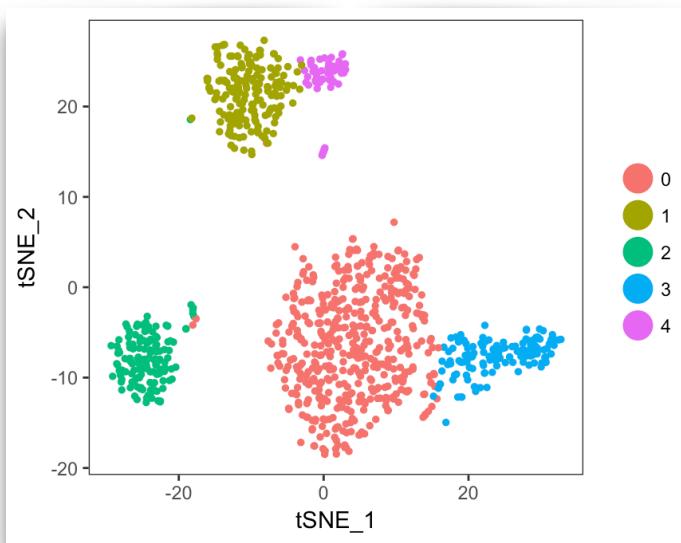


# Dimension reduction: general introduction

- Many types of data come as a matrix of N samples (e.g., cells, patients) x G features (e.g., genes, proteins)
- Each sample is a point in G-dimensional space
- Goal: represent the data in 2-3 dimensions, but preserve **structure** as best as possible (i.e., points that are **close** in G dimensions should be close in 2 or 3 dimensions)

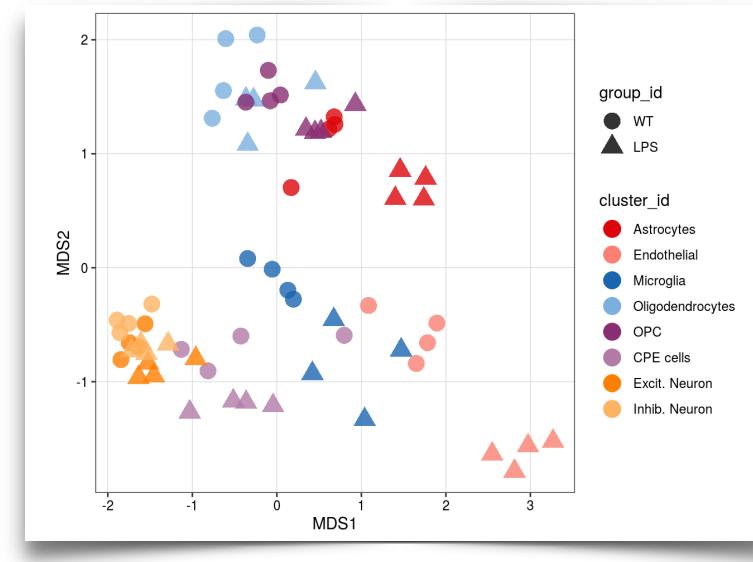
# Dimension reduction is versatile

$K \text{ features} \times N \text{ cells} \rightarrow$   
 $2 \text{ dimensions} \times N \text{ cells}$



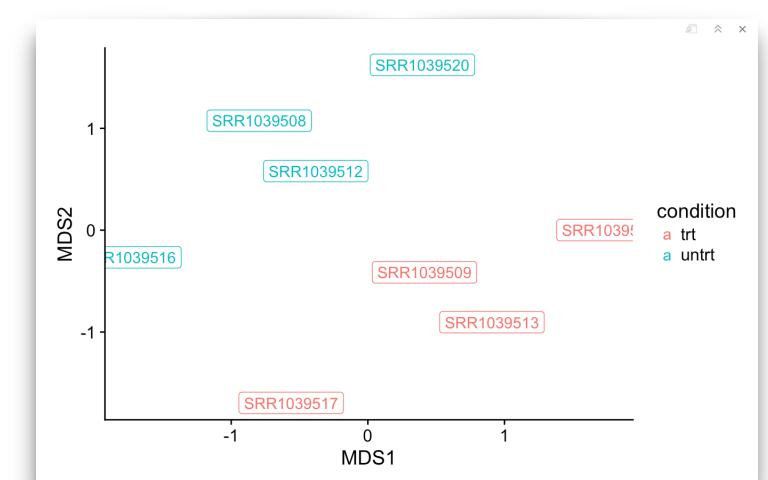
Each point =  
**single cell**  
(10x PBMC)

$N \text{ cells} \times K \text{ features} \rightarrow N \text{ cell}$   
subpopulations  $\times 2 \text{ dimensions}$



Each point =  
**subpopulation from a**  
**single sample** (LPS mouse cortex)

$P \text{ samples} \times K \text{ features} \rightarrow$   
 $P \text{ samples} \times 2 \text{ dimensions}$



Each point =  
**sample** (airway)

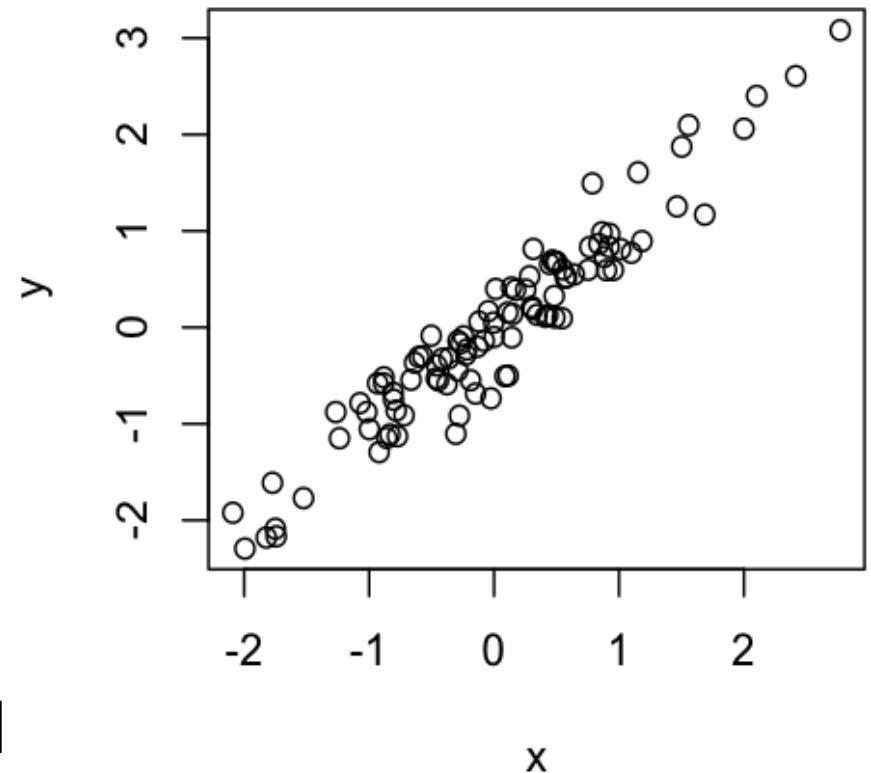
# Introduction to dimension reduction: PCA (principal components analysis)

- Form successive *linear* combinations of the features that are: orthogonal, ordered by variance

$$Y = XA$$

$$Y_{rk} = a_{1k}x_{r1} + a_{2k}x_{r2} + \cdots + a_{pk}x_{rp}$$

- A is the loadings matrix
- Typically, first 2-3 columns ('principal components') of Y are retained for visualisation; often top P PCs are retained for other analyses (e.g., clustering)



Many variations  
(linear/non-  
linear), many  
notions of  
distance, many  
ways to  
“compress”

