

# Deconvolution methods to dissect immunologic profile of the tumor

---

**Petri Pölönen, M.Sc**

University of Eastern Finland

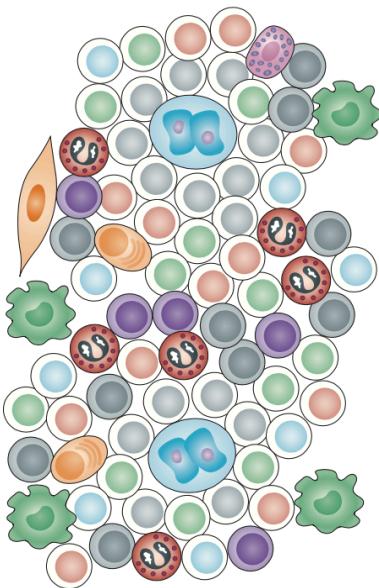
Institute of Biomedicine, school of medicine

Bioinformatics-Systems genomics group (Merja Heinäniemi PhD)



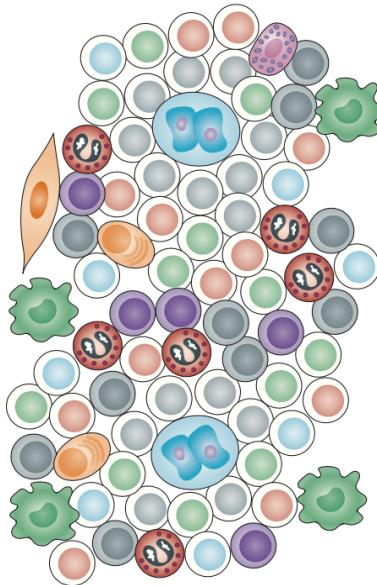
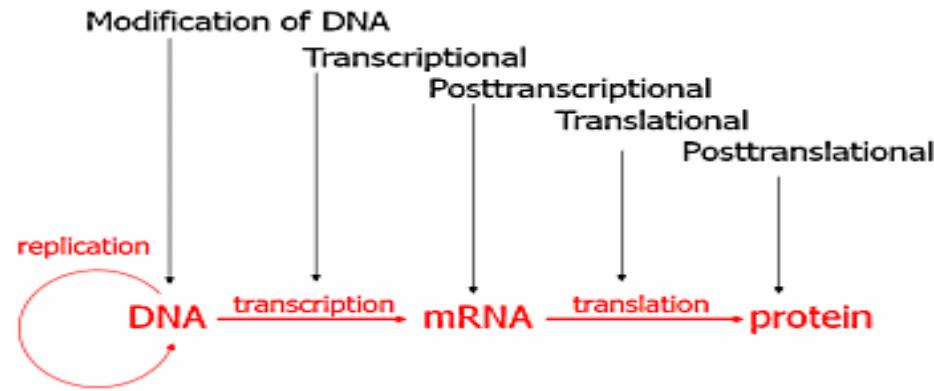
UNIVERSITY OF  
EASTERN FINLAND

“SOURCE separation, or deconvolution, is the problem of estimating individual signal components from their mixtures. This problem arises when source signals are transmitted through a mixing channel and the mixed sensor readings are observed.”



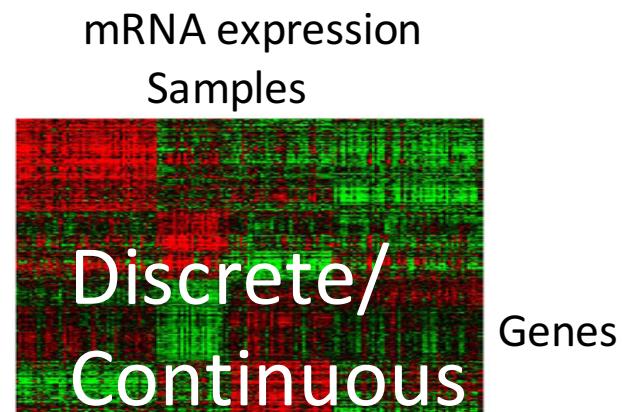
“SOURCE separation, or deconvolution, is the problem of estimating individual signal components from their mixtures. This problem arises when source signals are transmitted through a mixing channel and the mixed sensor readings are observed.”





All cells have same DNA, but they have different mRNA and protein levels

### *RNA-seq microarray*



# Deconvolution methods

## Intuition

- Models assume *linearity*: the expression signature of the mixture is a weighted sum of the expression profile for its constitutive cell-types  
→ approximation that is based on the linearity assumption
- Seek to distinguish features (genes) that closely conform to the linearity assumption, from the rest of the genes.

Mix:

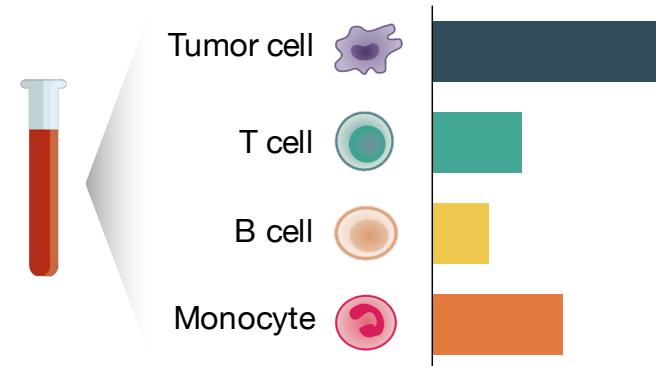
90 % Cells tumor

3 % Cells T-Cells

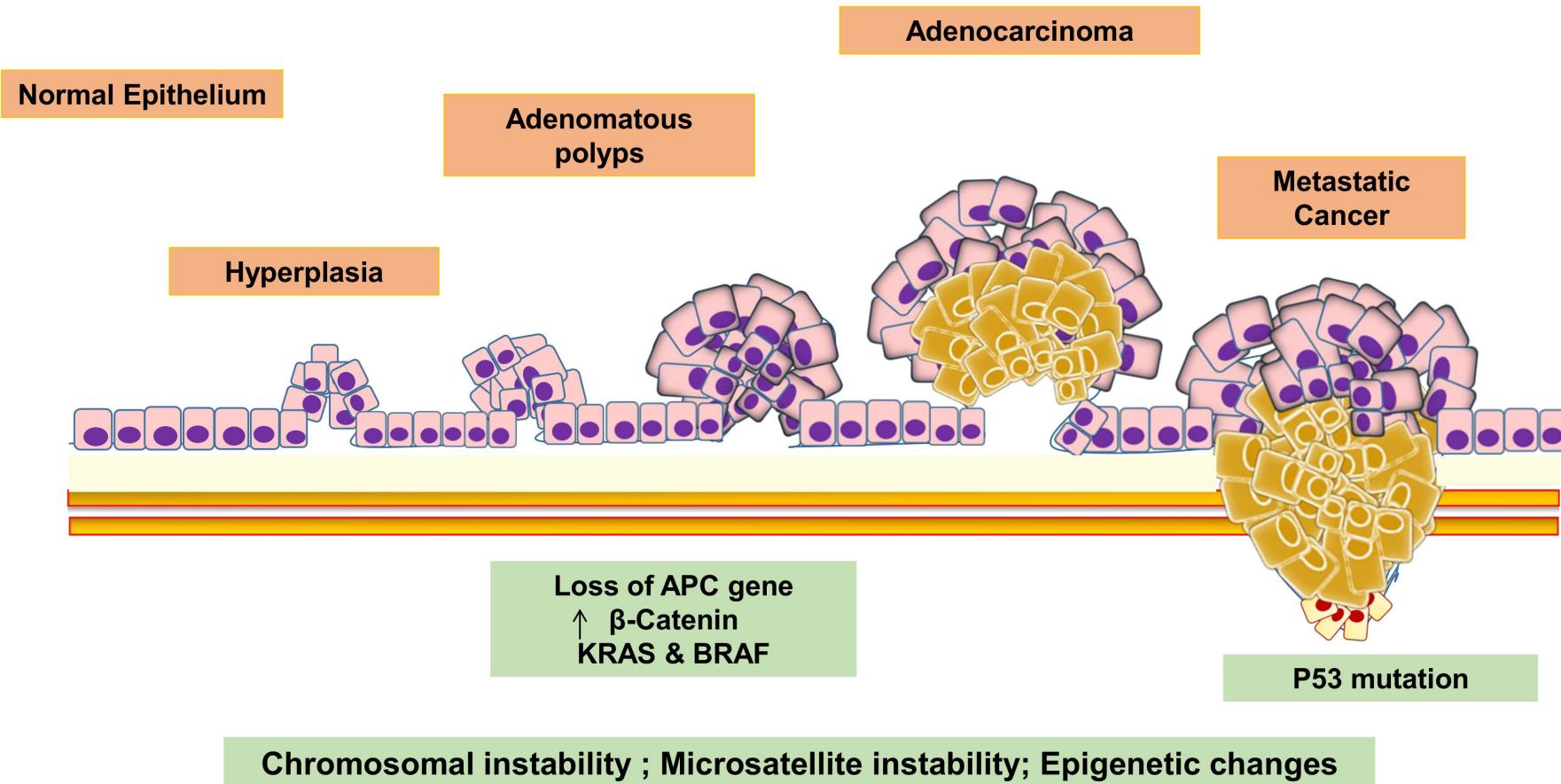
1 % Cells B-Cells

6 % Cells monocyte

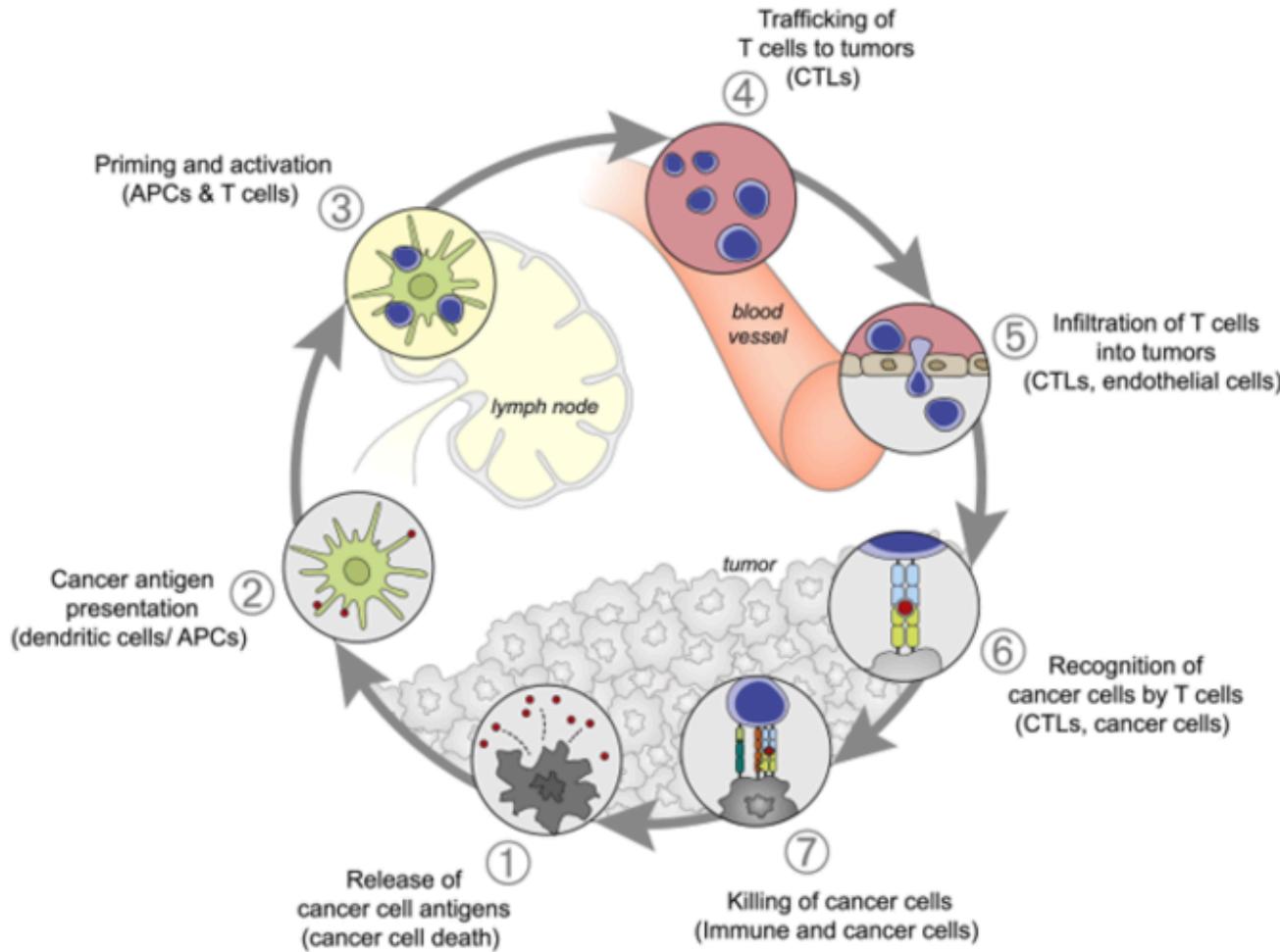
Total 100%



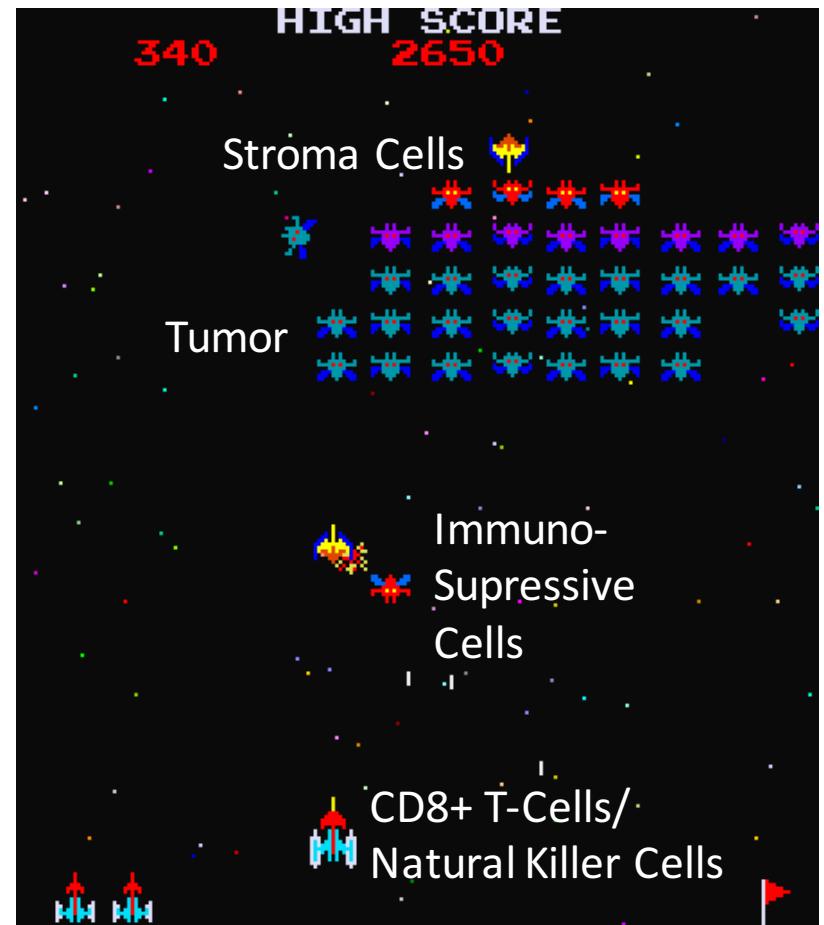
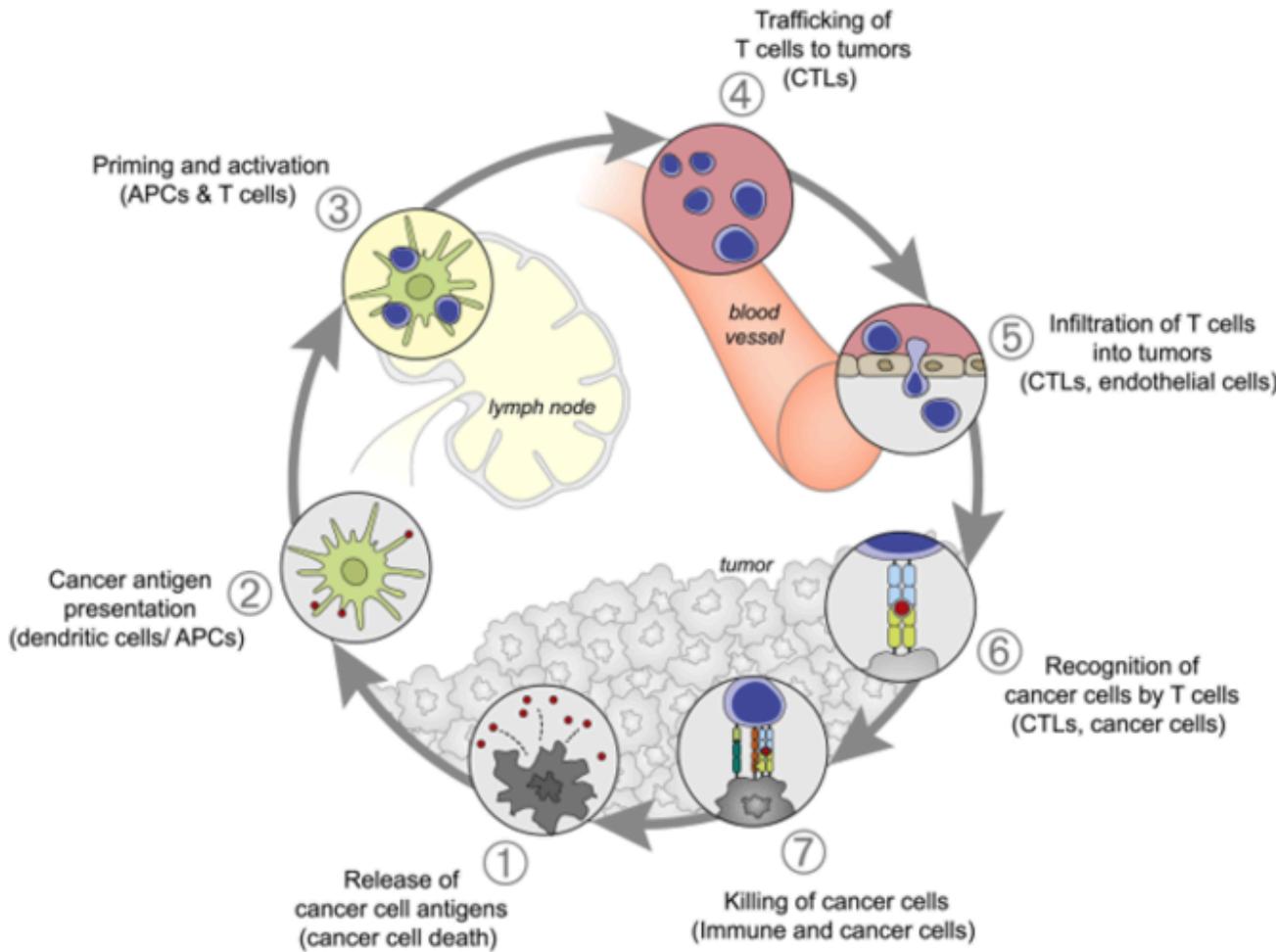
# Cancer development



# Tumor microenvironment matters!



# Tumor microenvironment matters!



# Why is immunology so hot right now?

Cell

Li et al. *Genome Biology* (2016) 17:174  
DOI 10.1186/s13059-016-1028-7

Genome Biology

## Molecular and Genetic Properties of Tumors Associated with Local Immune Cytolytic Activity

Michael S. Rooney,<sup>1,2</sup> Sachet A. Shukla,<sup>1,3</sup> Catherine J. Wu,<sup>1,3,4</sup> Gad Getz,<sup>1,5</sup> and Nir Hacohen<sup>1,4,6,\*</sup>

RESEARCH

Open Access



Comprehensive analyses of tumor immunity: implications for cancer immunotherapy

Bo Li<sup>1,2</sup>, Eric Severson<sup>1,3</sup>, Jean-Christophe Pignon<sup>3</sup>, Haoquan Zhao<sup>1</sup>, Taiwen Li<sup>4</sup>, Jesse Novak<sup>3</sup>, Peng Jiang<sup>1</sup>, Hui Shen<sup>5</sup>, Jon C. Aster<sup>3</sup>, Scott Rodig<sup>3</sup>, Sabina Signoretti<sup>3</sup>, Jun S. Liu<sup>2\*</sup> and X. Shirley Liu<sup>1\*</sup>

Associations to genetic properties (mutations, mutation load, copy number)

## The prognostic landscape of genes and infiltrating immune cells across human cancers

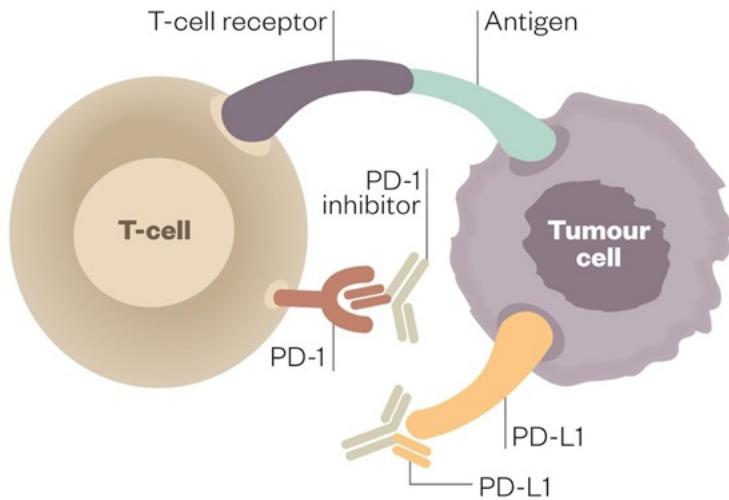
Andrew J Gentles<sup>1,2,12</sup>, Aaron M Newman<sup>3,4,12</sup>, Chih Long Liu<sup>3,4</sup>, Scott V Bratman<sup>3,5,11</sup>, Weiguo Feng<sup>3,5</sup>, Dongkyoon Kim<sup>3</sup>, Viswam S Nair<sup>6</sup>, Yue Xu<sup>7</sup>, Amanda Khuong<sup>7</sup>, Chuong D Hoang<sup>7,11</sup>, Maximilian Diehn<sup>3,5,8</sup>, Robert B West<sup>9</sup>, Sylvia K Plevritis<sup>1,2,13</sup> & Ash A Alizadeh<sup>1,3,4,8,10,13</sup>

938

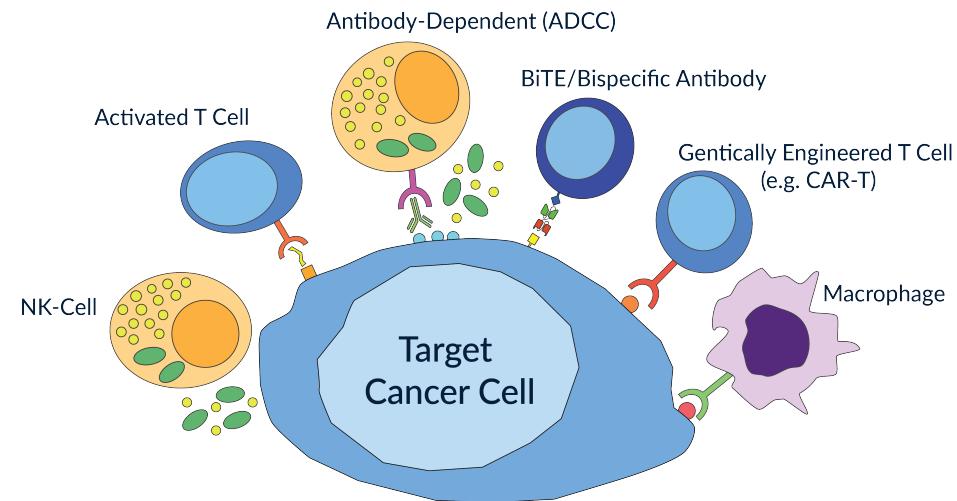
VOLUME 21 | NUMBER 8 | AUGUST 2015 **NATURE MEDICINE**

An increasing proportion of malignant cells, as well as a growing fraction of tumor infiltrating lymphocytes compared to surrounding cells, directly influence tumor growth, metastasis, and clinical outcomes for patients

# Why is immunology so hot right now?



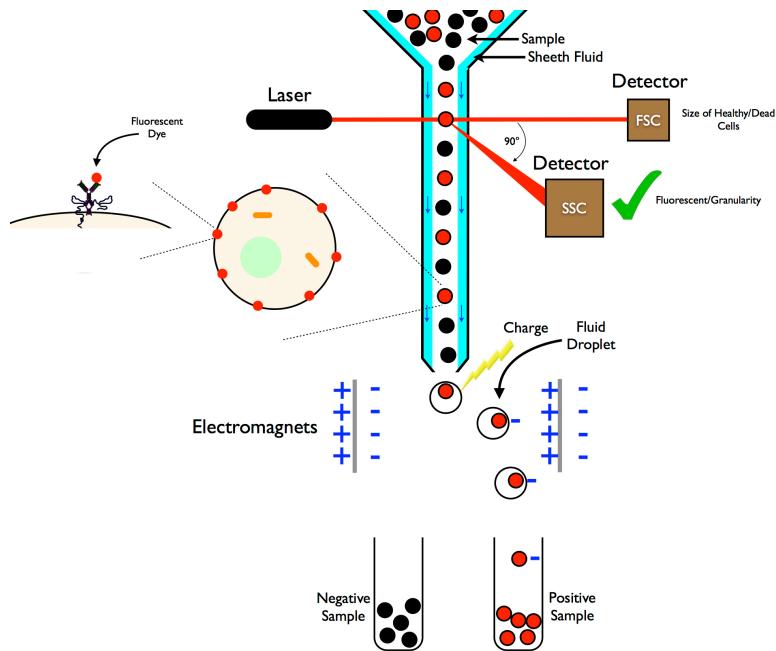
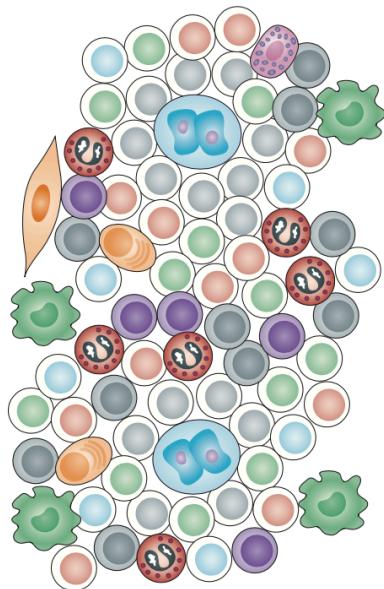
Make T-cells more efficient in killing tumor cells



Make T-cells that are engineered to kill cancer cells

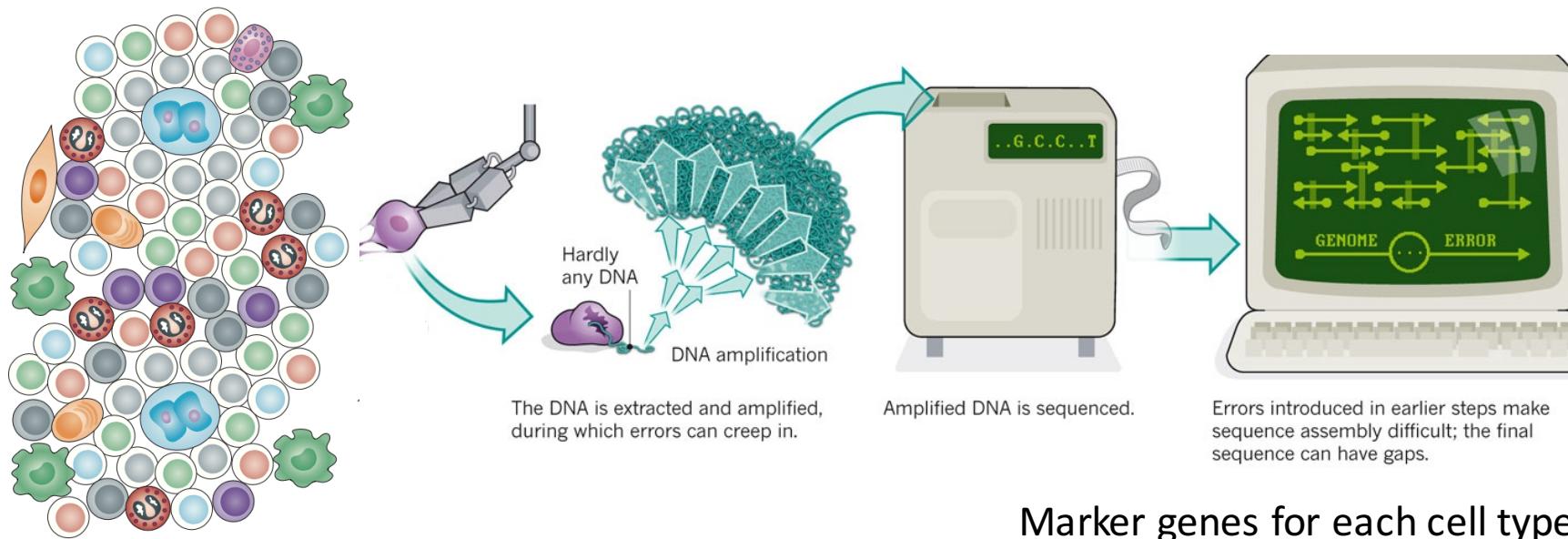
# Profiling immune cells using proteins: Cell sorting

- Cell isolation entails a loss of system perspective and is not currently feasible for all cell subsets



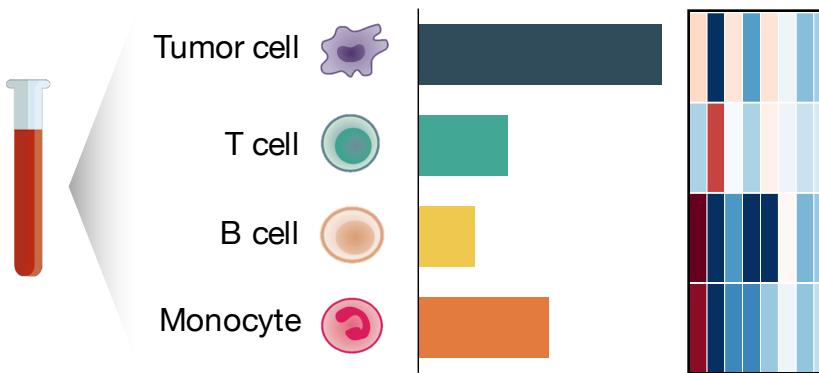
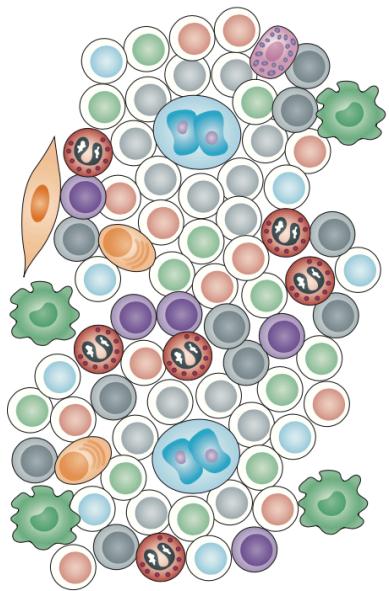
FACS: Depends on surface markers,  
max 8 markers at a time

# Profiling immune cells using single cell mRNA sequencing



Marker genes for each cell type  
Technical artefacts  
Markers not always detected accurately

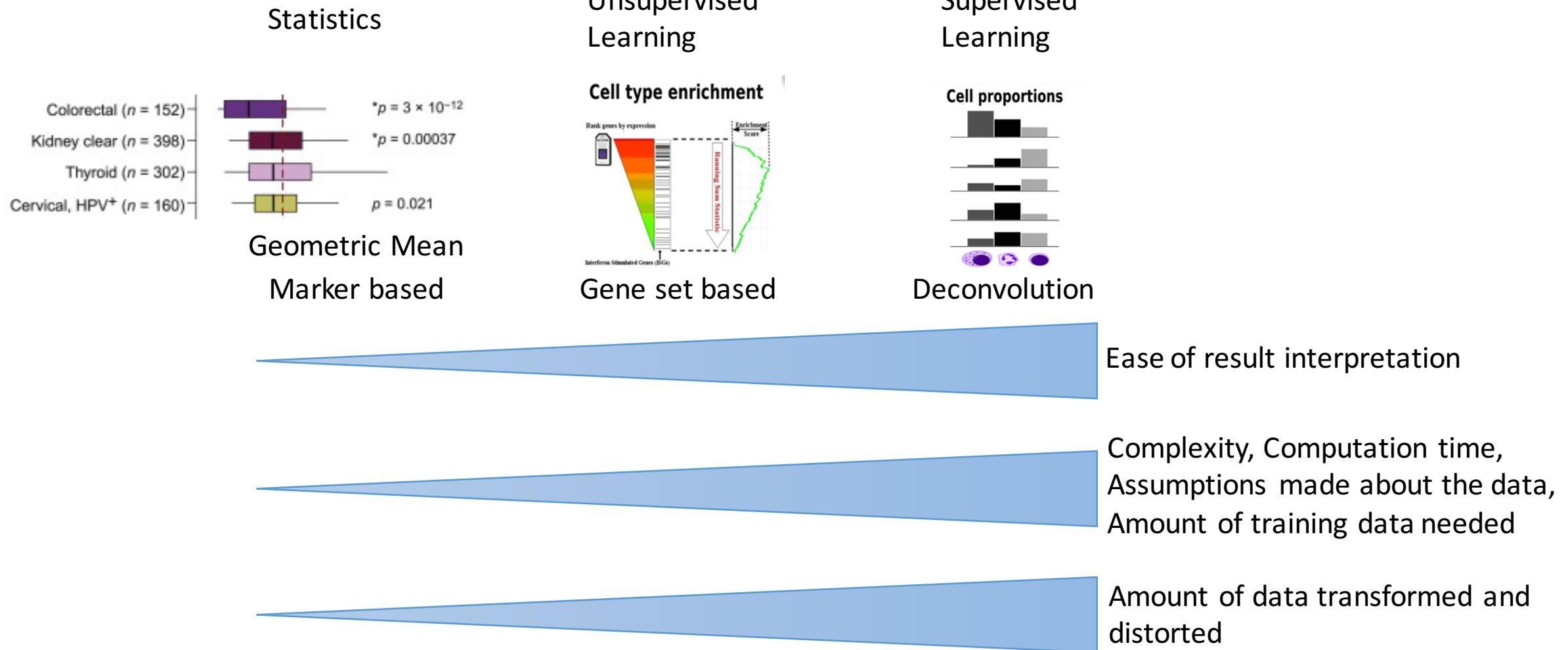
# Profiling immune cells using **bulk mRNA** sequencing, deconvolution to infer fractions



Data exists, cheap, data analysis straight forward

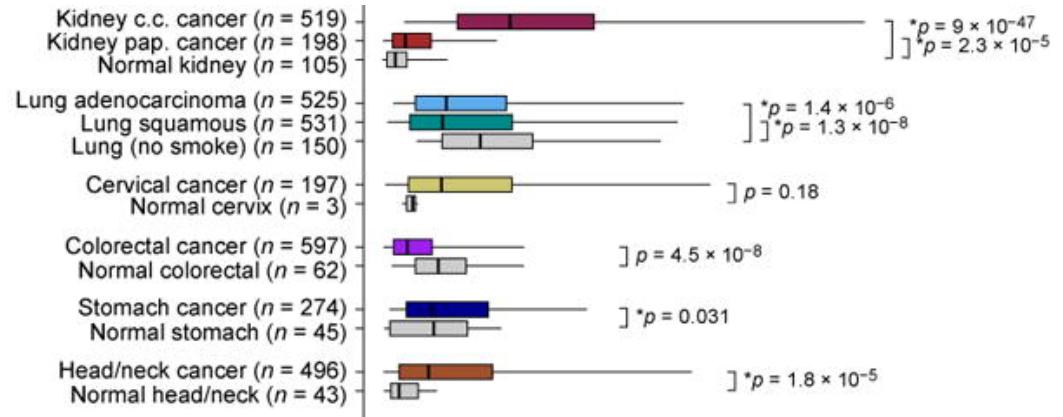
- Computational deconvolution methods can provide a cell-centered system perspective

# Methods



# Marker gene based methods:

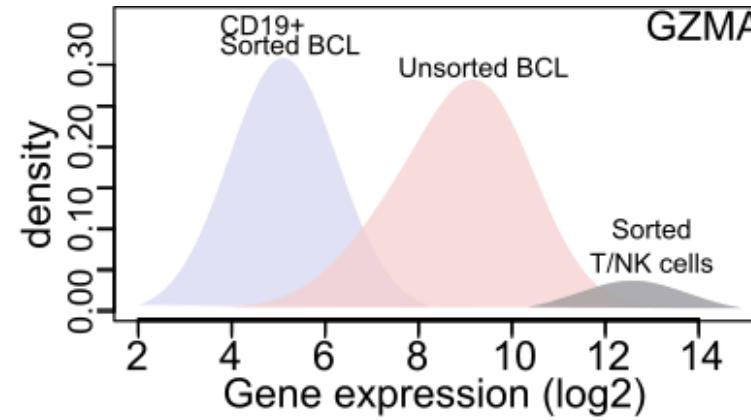
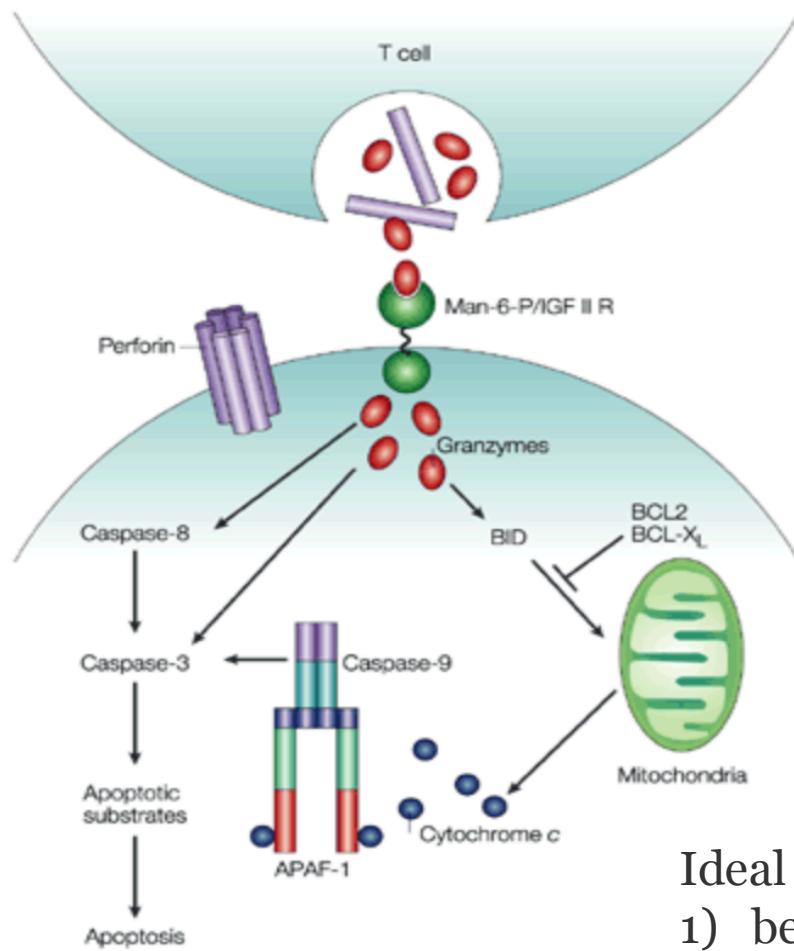
- Geometric mean:  
$$\left( \prod_{i=1}^n x_i \right)^{\frac{1}{n}} = \sqrt[n]{x_1 x_2 \cdots x_n}$$
- Subjective to noise
- Gene selection can be biased, depends on previous knowledge



Cytolytic activity (geometric mean of GZMA-PRF1) in solid tumors

Rooney MS, Shukla SA, Wu CJ, Getz G, Hacohen N. Molecular and genetic properties of tumors associated with local immune cytolytic activity. *Cell*. 2015;160(1–2):48–61.

## Making holes in the cancer cells, or starting programmed cell death of cancer cells



Pölönen & Dufva et. al (unpublished)

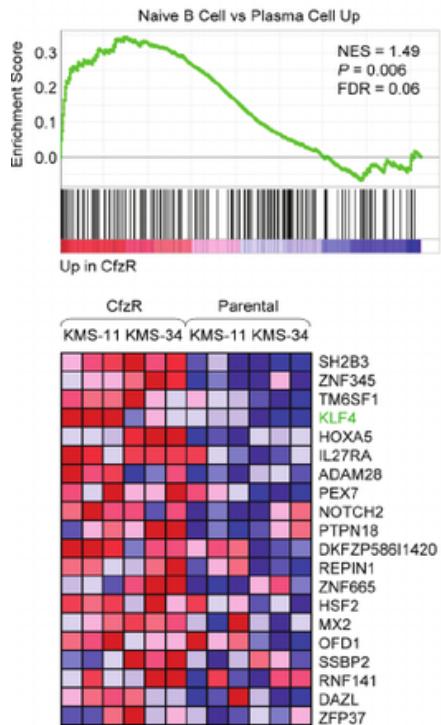
Ideal marker:

- 1) be highly informative of the cell type in which it is expressed
- 2) shows low variance due to spatiotemporal changes in the environment (changes in time, microenvironment or experiment)

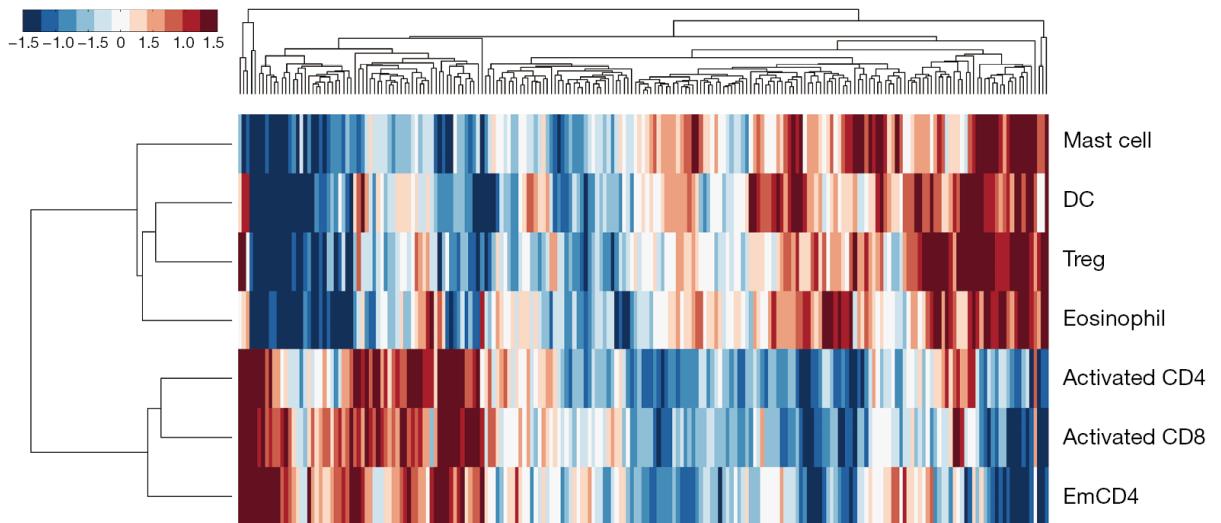
# Gene set based method using unsupervised learning:

For example Bindea et. al. gene sets, derived from immune cell differential expression analysis

- GSEA (two group comparison)



- ssGSEA/GSVA (single sample enrichment of a gene set)



# Deconvolution methods, supervised learning

## Biological data properties:

- *Biological variations*
  - Different developmental stages
  - *Response to environment change*
- *Technical variations*
- Number of parameters is greater than the number of samples ( $p \ll n$ )
- Features are highly correlated
- Some genes are ubiquitously expressed in all cell-types to perform housekeeping functions, whereas other genes exhibit *specificity* for one, or a group of cell-types

## Machine learning challenges:

→ *Noisy data (Training and test data)*

→ *Feature selection*

- Seek to distinguish features (genes) that closely conform to the linearity assumption, from the rest of the (variable) genes.

→ *Over-fitting*

- Model is excessively complex, such as having too many parameters relative to the number of observations

→ *Collinearity*

- Is a phenomenon in which two or more predictor variables in a multiple regression model are highly correlated, meaning that one can be linearly predicted from the others with a substantial degree of accuracy

# Different models used previously

## Regression methods:

- Ordinary Least Squares (OLS) regression (Abbas *et al.*)
- non-negative least squares (NNLS) regression (Qiao et al.)
- quadratic programming (QP) (Gong et. al)
- Non-negative matrix factorization (NMF), coupled with regression (Venet et al., Repsilber et al., Zuckerman et al., Gajoux et al.)
- robust linear regression (RLR) and  $\nu$ -SVR regression (Newman et al.)

## Bayesian methods:

- Bayesian prior and MCMC sampling (Erkkilä et al.)
- Latent Dirichlet Allocation (LDA)

Reference	Method	Loss	Non-negativity	Sum-to-one	Regularizer
Abbas <i>et al.</i> (2009)	Ordinary Least Squares (OLS)	$\mathcal{L}_2$	Imp	Imp	-
Gong <i>et al.</i> (2011)	Quadratic Programming	$\mathcal{L}_2$	Exp	Exp	-
Qiao <i>et al.</i> (2012)	Non-negative Least Squares (NNLS)	$\mathcal{L}_2$	Exp	Imp	-
DCQ- Altboum <i>et al.</i> (2014)	Elastic Net	$\mathcal{L}_2$	Imp	Imp	$\mathcal{L}_1/\mathcal{L}_2$
RLR- Newman <i>et al.</i> (2015)	Robust Linear Regression (RLR)	Huber	Imp	Imp	-
CIBERSORT- Newman <i>et al.</i> (2015)	$\nu$ -SVR	$\epsilon$ -insensitive	Imp	Imp	$\mathcal{L}_2$

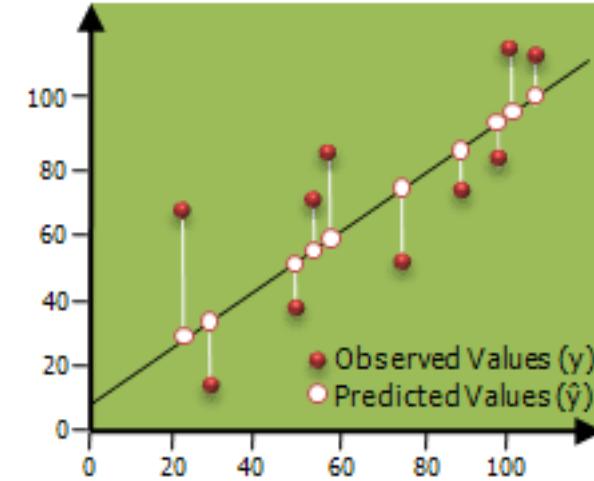
# Formal definition

$\mathcal{L}$  cost of estimation error

m Number of samples

$r_i$  Fitting error

Simple example, Ordinary Least Squares



$\underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{i=1}^m \mathcal{L}(r_i)$  Optimization problem, minimizes the sum of estimation error over all samples

# Formal definition

$\mathcal{L}$  cost of estimation error

$y_i$  expression level of a gene in the mixture

$x_i$  expression level of the same gene in the reference cell type

$w^T$  fraction of each cell type in the mixture, weight vector

$$\underset{f \in \mathcal{F}}{\operatorname{argmin}} \sum_{i=1}^m \mathcal{L}(r_i)$$

Fitting error  
 $r_i = y_i - f(x_i)$

Observed value      f(X) is a function that minimizes the aggregation error over samples,  $f(X_i)$  is a value of estimated value.

$f_w(x) = w^T x$  Linear regression finds a linear function that can predict observed values

$$r_i = y_i - w^T x_i$$

# Formal definition

$\mathcal{L}$  cost of estimation error

$y_i$  expression level of a gene in the mixture

$x_i$  expression level of the same gene in the reference cell type

$w^T$  fraction of each cell type in the mixture, weight vector

$$r_i = y_i - w^T x_i \text{ Fitting error}$$

Objective function

$$\underset{w \in \mathbb{R}^k}{\operatorname{argmin}} \left\{ \underbrace{\sum_{i=1}^m \mathcal{L}(y_i - w^T x_i)}_{\text{Overall loss}} \right\}$$

Finding weights ( $w^T$ )  
for X to minimize the  
prediction error

# Formal definition

$\mathcal{L}$  cost of estimation error

$y_i$  expression level of a gene in the mixture

$x_i$  expression level of the same gene in the reference cell type

$w^T$  fraction of each cell type in the mixture

$\lambda$  parameter controls the relative importance of estimation error  
versus regularization (regularization parameter)

$r_i = y_i - w^T x_i$  Fitting error

Objective function     $\underset{w \in \mathbb{R}^k}{\operatorname{argmin}} \left\{ \sum_{i=1}^m \mathcal{L}(y_i - w^T x_i) \right\}$

loss function  $\mathcal{L}$

The diagram shows the objective function as a minimization problem over the weight vector  $w$ . The objective is the sum of individual loss terms  $\mathcal{L}(y_i - w^T x_i)$  for each sample  $i$  from 1 to  $m$ . A blue bracket under the summation term is labeled "Overall loss".

Function that maps values of one or more variables onto a real number intuitively representing some "cost" associated with the event. An optimization problem seeks to minimize a loss function.

<https://www.youtube.com/watch?v=euhATa4wgzo>

<https://www.youtube.com/watch?v=iSfcRku6euQ>

# Deconvolution methods

## *Choice of Loss (or Cost) Functions*

- *Ordinary Least Squares (OLS)*  $\mathcal{L}_2(r_i) = r_i^2 = (y_i - \mathbf{w}^T \mathbf{x}_i)^2$
- Absolute deviation loss  $\mathcal{L}_1(r_i) = |r_i| = |y_i - \mathbf{w}^T \mathbf{x}_i|$
- Huber's loss function  $\mathcal{L}_{\text{Huber}}^{(M)}(r_i) = \begin{cases} r_i^2, & \text{if } |r_i| \leq M \\ M(2|r_i| - M), & \text{otherwise} \end{cases}$
- *Support Vector Regression* 
$$\begin{aligned} \mathcal{L}_{\epsilon}^{(\epsilon)}(r_i) &= \max(0, |r_i| - \epsilon) \\ &= \begin{cases} 0, & \text{if } |r_i| \leq \epsilon \\ |r_i| - \epsilon, & \text{otherwise.} \end{cases} \end{aligned}$$

# Formal definition

$\mathcal{L}$  cost of estimation error

$y_i$  expression level of a gene in the mixture

$x_i$  expression level of the same gene in the reference cell type

$w^T$  fraction of each cell type in the mixture

$\lambda$  parameter controls the relative importance of estimation error  
versus regularization (regularization parameter)

$r_i = y_i - w^T x_i$  Fitting error

$$\operatorname{argmin}_{f \in \mathcal{F}} \sum_{i=1}^m \mathcal{L}(r_i)$$

$$f_w(x) = w^T x$$

$$r_i = y_i - f(x_i)$$

$$\operatorname{argmin}_{w \in \mathbb{R}^k} \left\{ \underbrace{\sum_{i=1}^m \mathcal{L}(y_i - w^T x_i)}_{\text{Overall loss}} + \underbrace{\lambda \mathcal{R}(w)}_{\text{Regularizer}} \right\}$$

regularizer function  $\mathcal{R}$

Function for selecting the preferred level of model complexity  
so your models are better at predicting (generalizing). Too  
complex model can overfit or be too simple and underfit,  
either way giving poor predictions

<https://www.youtube.com/watch?v=C79kIYkKZ1g>

# Deconvolution methods

## *Choice of Regularizers*

- L1/norm-1 regularizer/Lasso
- L2/norm-2 regularizer/ridge regression
- *L1+L2/elastic net*
- *Group Lasso*

$\mathcal{R}$  regularizer function

$w$  regression coefficient vector

$$\mathcal{R}_1(w) = \|w\|_1 = \sum_{i=1}^k |w_i|$$

$$\mathcal{R}_2(w) = \|w\|_2^2 = \sum_{i=1}^k w_i^2$$

$$\mathcal{R}_{\text{elastic}}(w) = \alpha \mathcal{R}_1(w) + (1 - \alpha) \mathcal{R}_2(w)$$

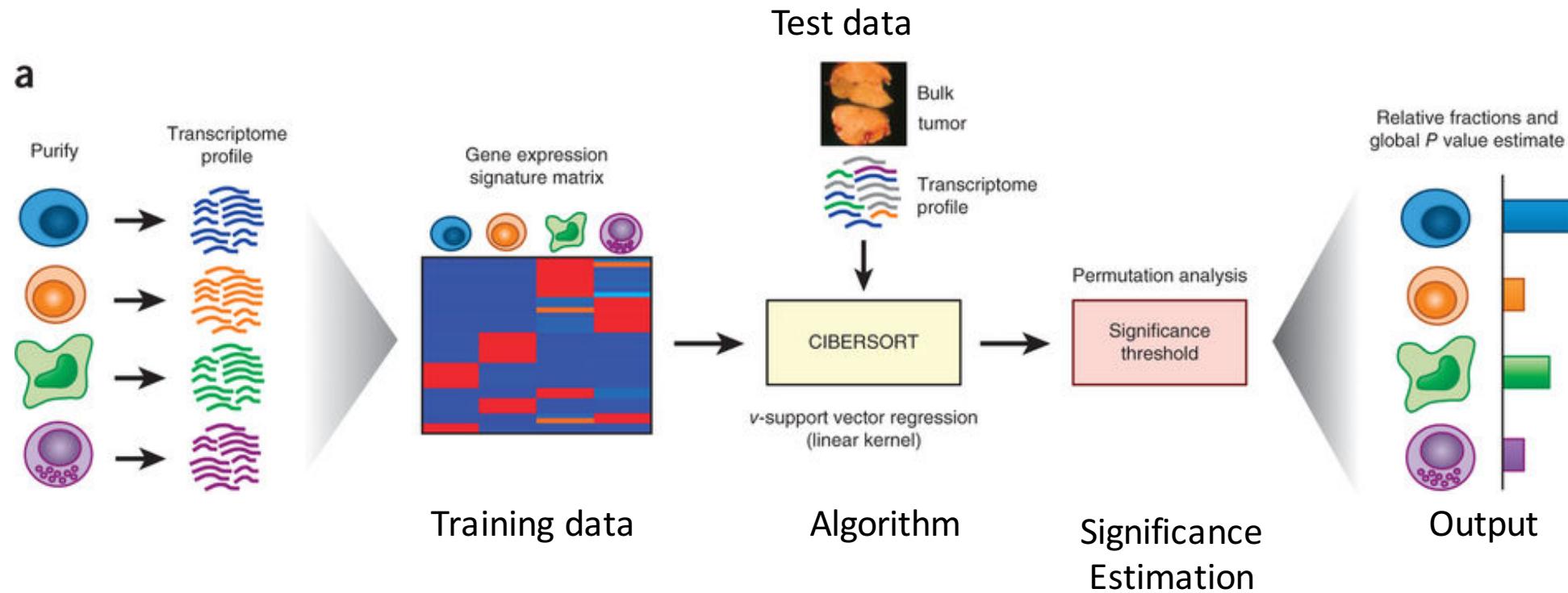
$$\mathcal{R}_{\text{group}} = \sum_{G_i} \mathcal{L}_2(w(G_i))$$

# Objective Functions Used in Practice

- Ordinary Least Squares (OLS)
- Ridge Regression
- Least Absolute Selection and Shrinkage Operator (LASSO) Regression
- Robust Regression
- Support Vector Regression

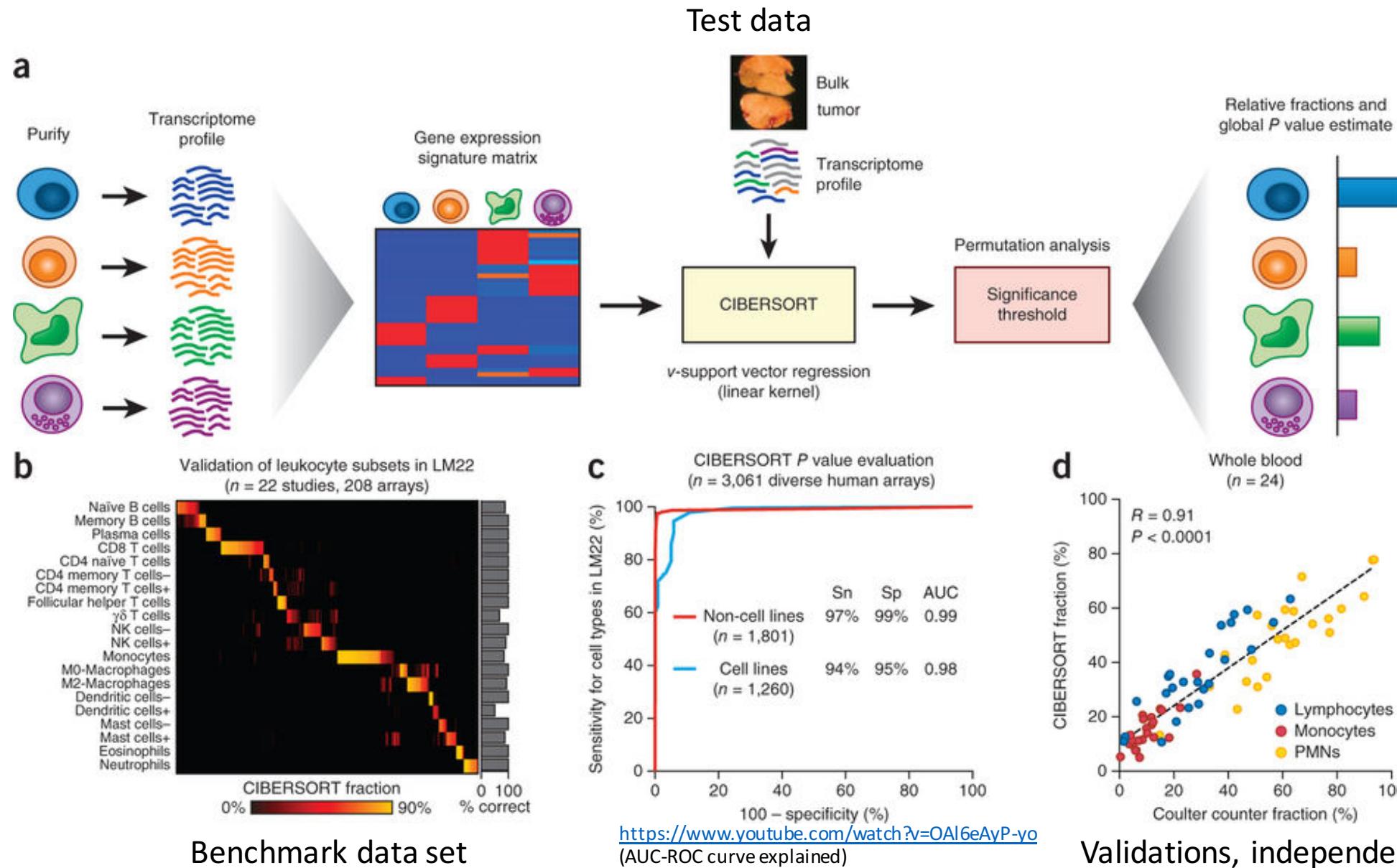
# Deconvolution methods

Example: CIBERSORT



# Deconvolution methods

Example: CIBERSORT



# Cibersort algorithm

- nu-support vector regression (v-SVR)
- Robustness to noise and overfitting owing to both a linear loss function and feature selection of genes from the signature matrix
- Tolerance to collinearity via utilization of the L2-norm penalty function.
- Produces an empirical P-value for the deconvolution using Monte Carlo sampling

# DECONVOLUTION WAR!

2015

A horizontal blue arrow points from left to right, representing time. A vertical blue line segment extends upwards from the start of the arrow, ending at the year '2015' written above it. To the left of this vertical line, the text 'CIBERSORT,  
22 cell type  
deconvolution' is written vertically.  
CIBERSORT,  
22 cell type  
deconvolution

Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, et al. Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods*. 2015;12:453–7.  
2015

Li B, Severson E, Pignon JC, Zhao H, Li T, Novak J, et al. Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. *Genome Biol*.  
2016;17:174.

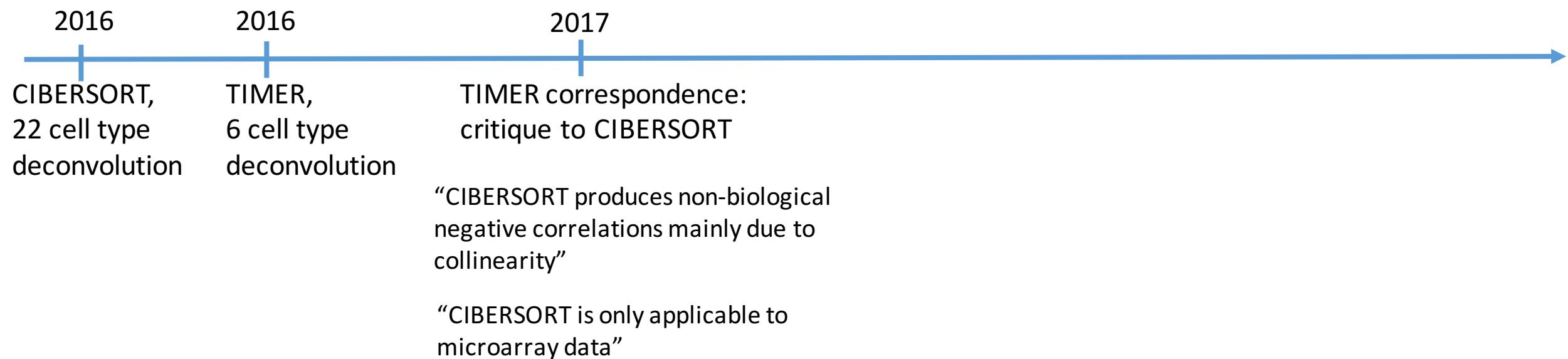
# DECONVOLUTION WAR!



Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, et al. Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods*. 2015;12:453–7.  
2016

Li B, Severson E, Pignon JC, Zhao H, Li T, Novak J, et al. Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. *Genome Biol*.  
2016;17:174.

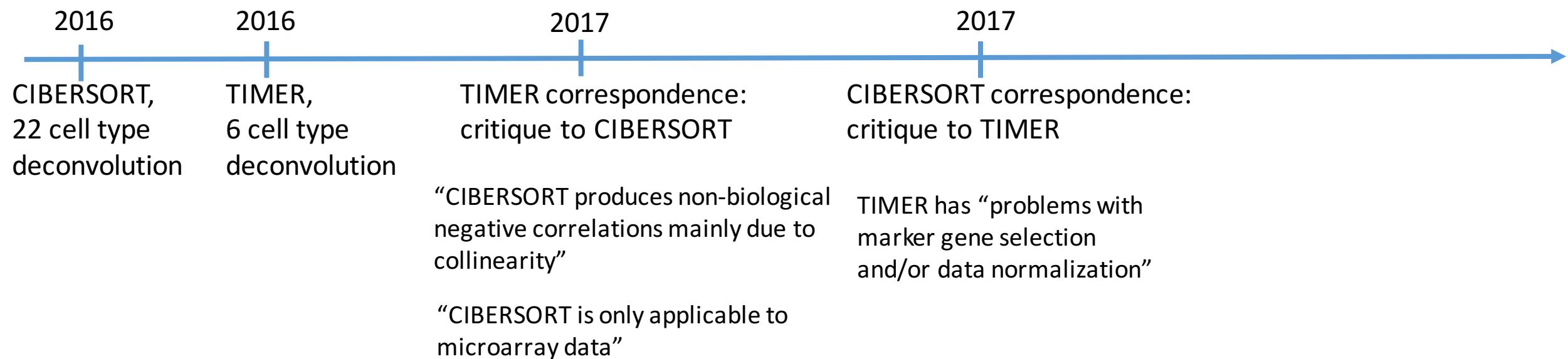
# DECONVOLUTION WAR!



Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, et al. Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods*. 2015;12:453–7.  
2016

Li B, Severson E, Pignon JC, Zhao H, Li T, Novak J, et al. Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. *Genome Biol*.  
2016;17:174.

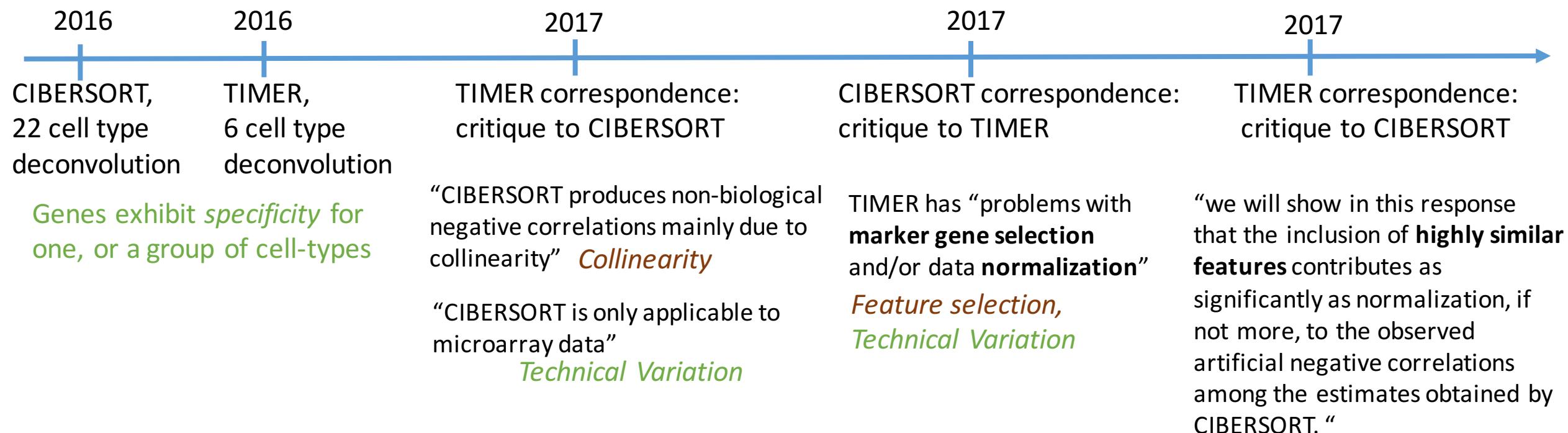
# DECONVOLUTION WAR!



Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, et al. Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods*. 2015;12:453–7.  
2016

Li B, Severson E, Pignon JC, Zhao H, Li T, Novak J, et al. Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. *Genome Biol*.  
2016;17:174.

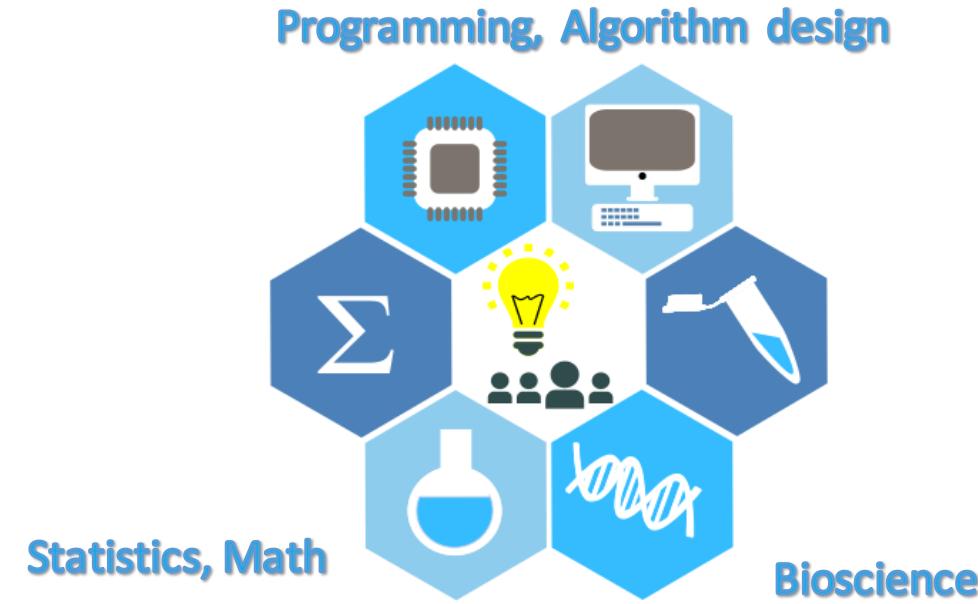
# DECONVOLUTION WAR!



Newman AM, Liu CL, Green MR, Gentles AJ, Feng W, Xu Y, et al. Robust enumeration of cell subsets from tissue expression profiles. *Nat Methods*. 2015;12:453–7.  
2016

Li B, Severson E, Pignon JC, Zhao H, Li T, Novak J, et al. Comprehensive analyses of tumor immunity: implications for cancer immunotherapy. *Genome Biol*. 2016;17:174.

# Bioinformatics is a challenging multidisciplinary field



# For reading and references

- Shen-Orr SS, Gaujoux R. Computational deconvolution: extracting cell type-specific information from heterogeneous samples. *Curr Opin Immunol.* 2013;25:571–8. doi: 10.1016/j.coim.2013.09.015. **REVIEW, biological overview of the problem**
- Mohammadi S, Zuckerman N, Goldsmith A, Grama A. A critical survey of deconvolution methods for separating cell-types in complex tissues. *arXiv.* 2015: 1510.04583 **Technical REVIEW of methods**

# Questions:

Consider you are working with deconvolution to dissect immunologic profile of the tumor:

- Why deconvolution can be useful and what is the main output of the algorithm?
- Which steps are needed to develop a deconvolution algorithm and reason why?  
Hint: Point is not to list all different algorithms that could be used, but to point out the key steps and reasons why they are needed, as referred below.

Training data

Test data

Algorithm

- Cost function
- Regularization
- Feature Selection

Benchmarking/Validation