



Michigan State University  
<http://binchenlab.org>

# Transcriptomics + drug discovery

**Bin Chen**

Associate Professor

Dept. of Pediatrics and Human Development

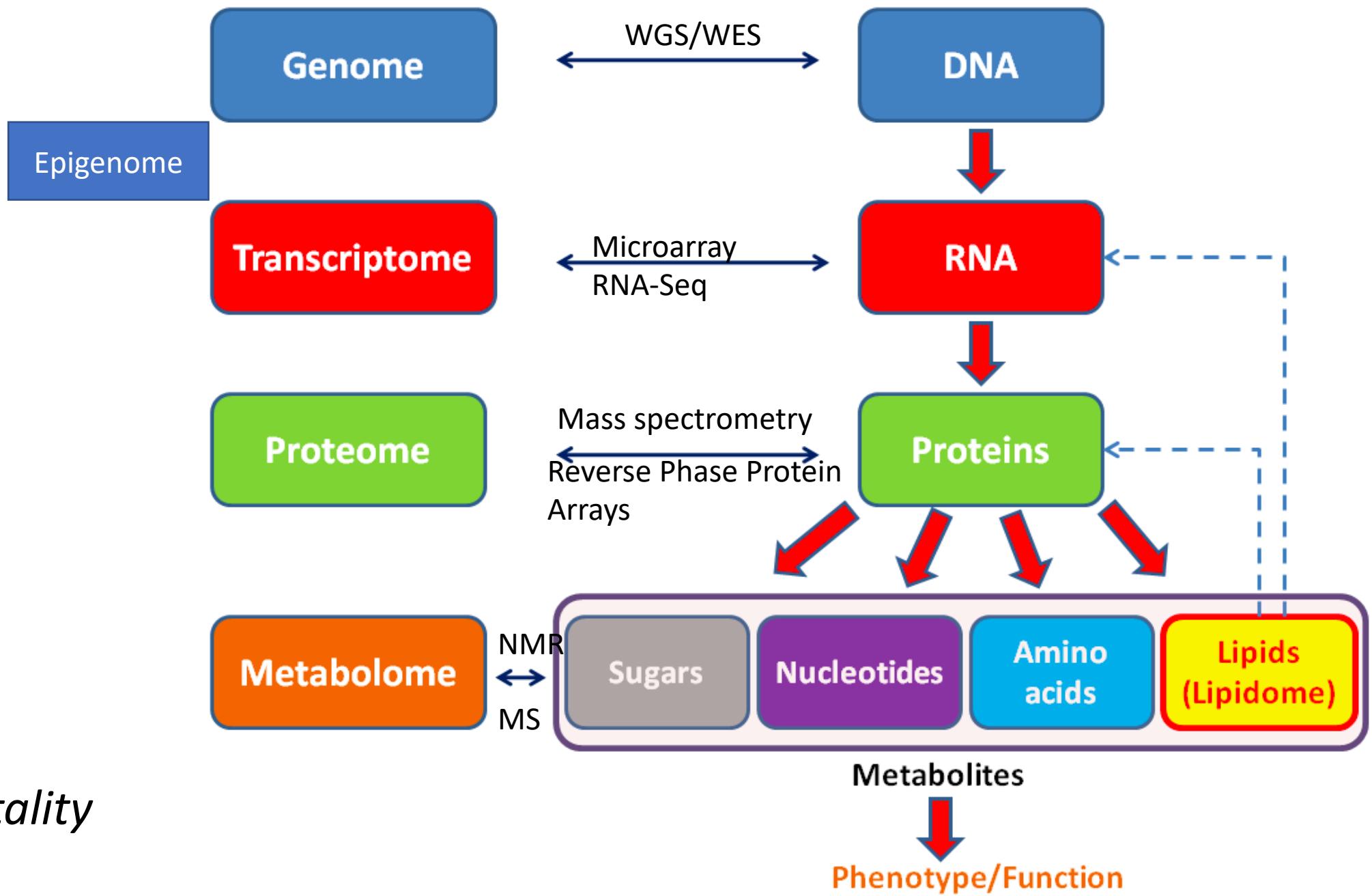
Dept. of Pharmacology and Toxicology

College of Human Medicine

Michigan State University

[Bin.Chen@hc.msu.edu](mailto:Bin.Chen@hc.msu.edu) @DrBinChen

<http://binchenlab.org>



-Ome: *totality*

# Potential questions we can ask using RNASeq

- Identify potential targets/pathways/cellular changes after pharmacological perturbation

Treatment group  
(3 replicates)

Control group  
(3 replicates)

Library preparation

Sequencing  
(done by company)

	T 1	T 2	T 3	C 1	C 2	C 3
Ge ne1						
Ge ne2						
Ge ne3						

Raw counts

Data analysis

Clustering

Differential  
expression analysis

Pathway enrichment

Functional analysis

	T 1	T 2	T 3	C 1	C 2	C 3
Ge ne1						
Ge ne2						
Ge ne3						

TPM/FPKM

```
> count_data <- read.csv(gzfile("~/Downloads/GSE187420_Raw_Read_count.csv.gz"), row.names = 1)
> tpm_data <- read.csv(gzfile("~/Downloads/GSE187420_Normalized TPM_count.csv.gz"), row.names = 1)
> head(count_data)

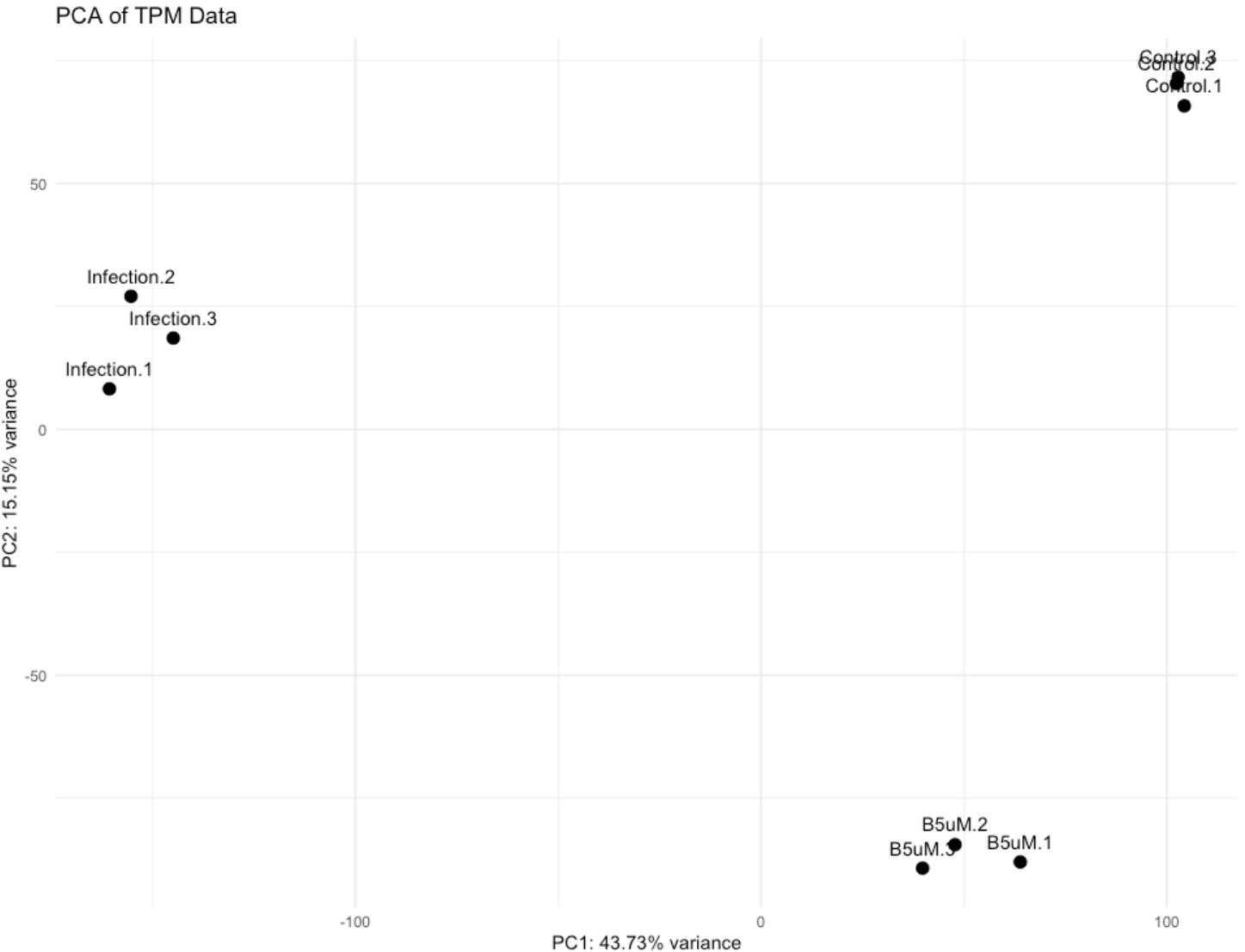
          Control.1 Control.2 Control.3 Infection.1 Infection.2 Infection.3    B5uM.1    B5uM.2    B5uM.3
ENSG00000000003 11.627078 11.384244 11.456868      5.727920      5.807355      5.930737 11.535761 11.179287 11.194757
ENSG00000000005  0.000000  0.000000  0.000000      0.000000  0.00000000 0.00000000 0.00000000 0.00000000 0.00000000
ENSG00000000419 10.817783 10.490851 10.529431      5.930737      6.087463      6.066089 10.303781 9.962896 9.962896
ENSG00000000457  8.926533  8.829786  8.869964      5.891905      6.140779      6.105804 7.676662 7.427439 7.337354
ENSG00000000460  9.336328  9.226364  9.401157      5.479619      5.106013      4.761817 9.288635 9.187006 9.202809
ENSG00000000938  0.000000  0.000000  0.000000      0.000000  1.00000000 1.00000000 0.00000000 0.00000000 0.00000000

> head(tpm_data)

          Control.1 Control.2 Control.3 Infection.1 Infection.2 Infection.3    B5uM.1    B5uM.2    B5uM.3
ENSG00000000003  5.187847  5.050066  5.122673      2.032101      1.9411063     2.2898345 5.171127 4.934988 4.969473
ENSG00000000005  0.000000  0.000000  0.000000      0.000000  0.00000000 0.00000000 0.00000000 0.00000000 0.00000000
ENSG00000000419  5.602291  5.384741  5.415488      3.144046      3.1921942     3.4235782 5.146900 4.915999 4.960234
ENSG00000000457  2.223423  2.195348  2.150560      1.851999      1.9030383     2.1890338 1.189034 1.220330 1.097611
ENSG00000000460  2.933573  2.967169  3.130931      2.077243      1.2927817     1.4646683 3.152183 3.054848 3.140779
ENSG00000000938  0.000000  0.000000  0.000000      0.000000  0.16349870 0.0976108 0.000000 0.000000 0.000000
```

# Sample visualization

```
pca_results <- PCA(t(tpm_data))
```



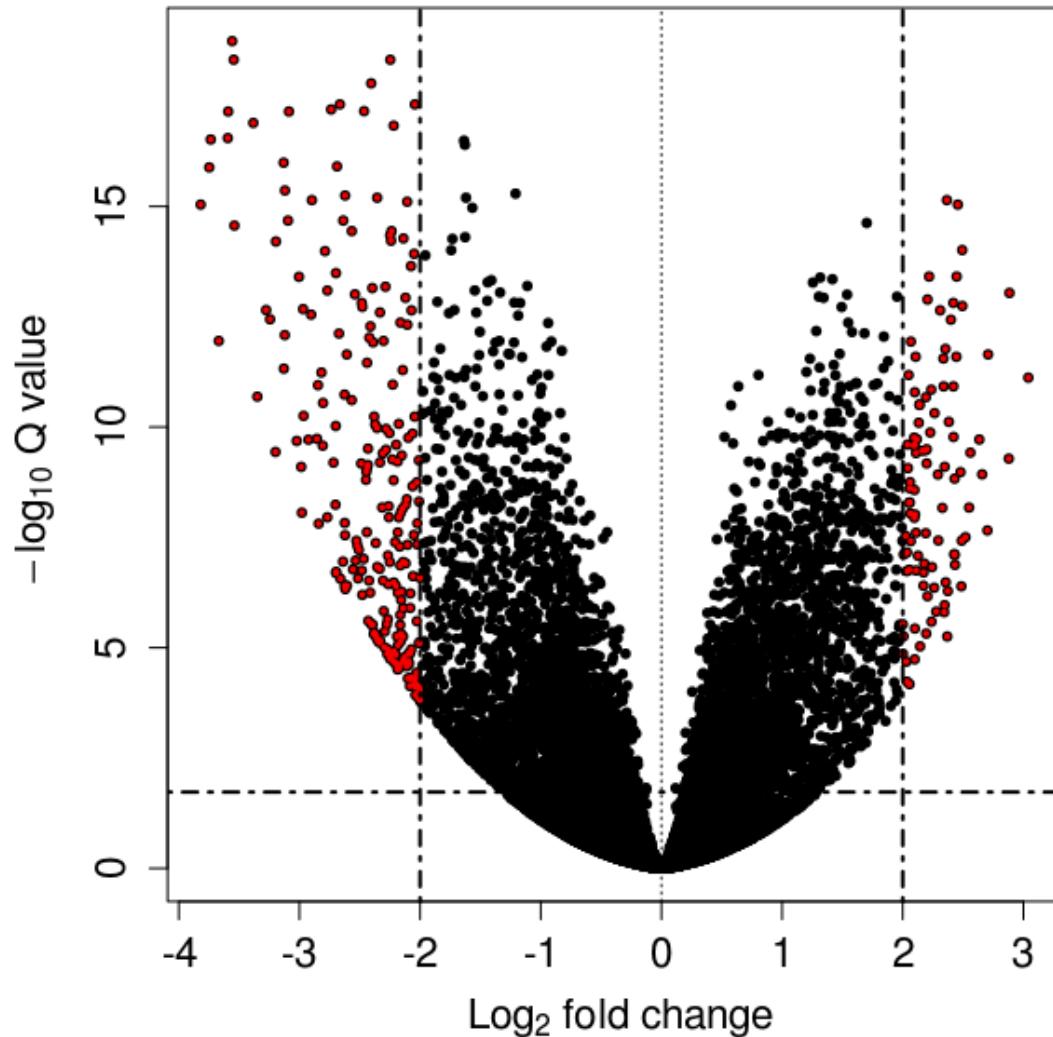
# Gene expression analysis

	T 1	T 2	T 3	C 1	C 2	C 3
Gen e1						
Gen e2						
Gen e3						
Gen e4						

DE analysis  
DESeq2, EdgeR, Limma

The goal is to compute fold change and significance!

Volcano plot



# A naïve approach

	T 1	T 2	T 3	C 1	C 2	C 3
Gene1						
Gene2						
Gene3						
Gene4						

Method 1

$$y = X\beta + \varepsilon$$

$$\beta = (X^T X)^{-1} X^T y$$

Count =  $b_0 + b_1 * \text{treatment}$

fold change, p value

Method 2

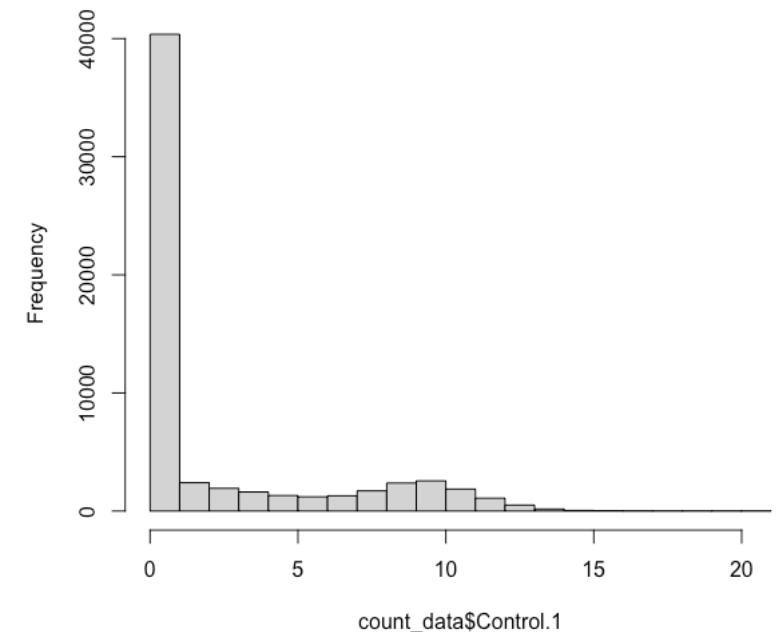
t. test

fold change, p value



Both linear regression and t test require normal distribution of expression count

Histogram of count\_data\$Control.1



# Data normalization

	T1	T2	T3	C1	C2	C3
Gene1	10	20	25	1	1	1
Gene2	24	40	45	1	2	1
Gene3	5	10	8	2	1	2
Gene4	23	20	20	150	160	90

Normalize total read accounts per sample  
Assuming every sample has same transcripts

Normalize each gene across samples to  
adjust for differences in library composition

Regardless of normalization methods, the counts are not normally distributed

# Negative Binomial Distribution

$NB(r, p)$ : # of success before the first  $r$  failure  
given success probability  $p$

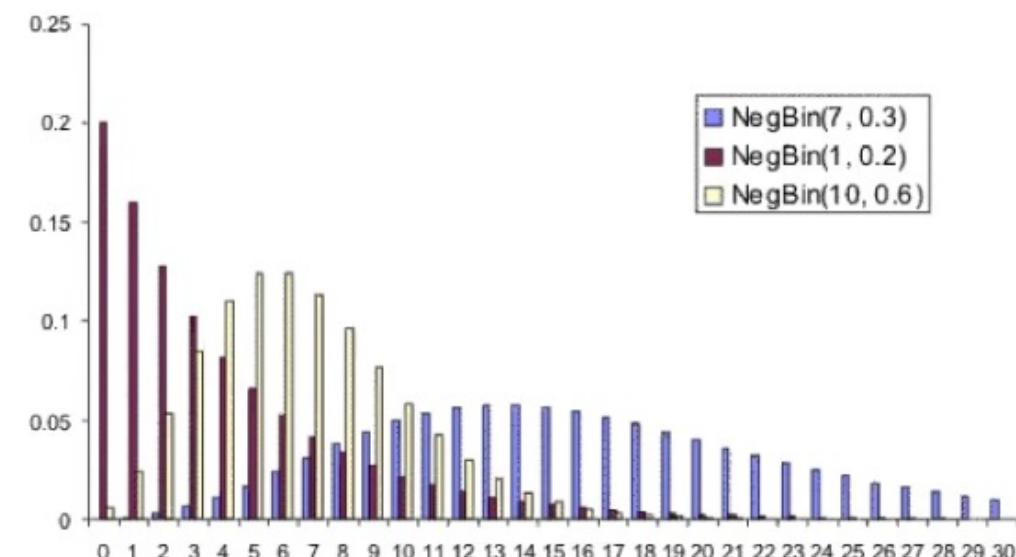
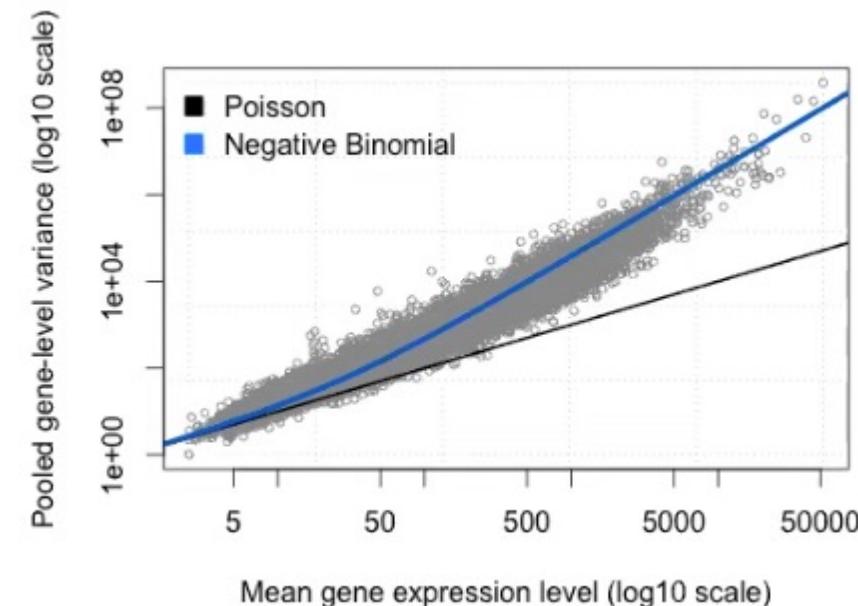
$$K_{ij} \sim NB(s_{ij}q_{ij}, \alpha_i)$$

Raw count for gene  $i$  in sample  $j$

Normalization factor (library size)

dispersion

Expression level



# Generalized linear regression

	T 1	T 2	T 3	C 1	C 2	C 3
Gene1						
Gene2						
Gene3						
Gene4						

Method 1

$$y = X\beta + \varepsilon$$

$$\beta = (X^T X)^{-1} X^T y$$

$$\text{Count} = b_0 + b_1 * \text{treatment} \longrightarrow \text{fold change, p value}$$

$$E(Y|X) = \mu = g^{-1}(X\beta)$$

$$\mu_j = s_j * \exp(X\beta_j)$$

$$\log(\mu_j) = \log(s_j) + X\beta_j$$

Generalized linear model: a unified framework for modeling various types of response variables by specifying a distribution from the exponential family, a linear predictor, and a link function

# DESeq2 approach:

## 1. Negative Binomial Distribution:

Raw counts

$$Y_{ij} \sim \text{NB}(\mu_{ij}, \alpha_i)$$

dispersion

## 2. Mean of Counts:

Link Function:

- The GLM in DESeq2 uses a logarithmic link function to relate the mean count  $\mu_{gi}$  to the linear predictors:

$$\mu_{ij} = s_j \cdot \exp(\eta_{ij}) \longrightarrow \log(\mu_{gi}) = \eta_{gi}$$

- The linear predictor  $\eta_{gi}$  is specified as:

where:

- $\mu_{ij}$  is the mean count for gene  $i$  in sample  $j$ .
- $s_j$  is the size factor for sample  $j$ .
- $\eta_{ij}$  is the linear predictor for gene  $i$  in sample  $j$ .

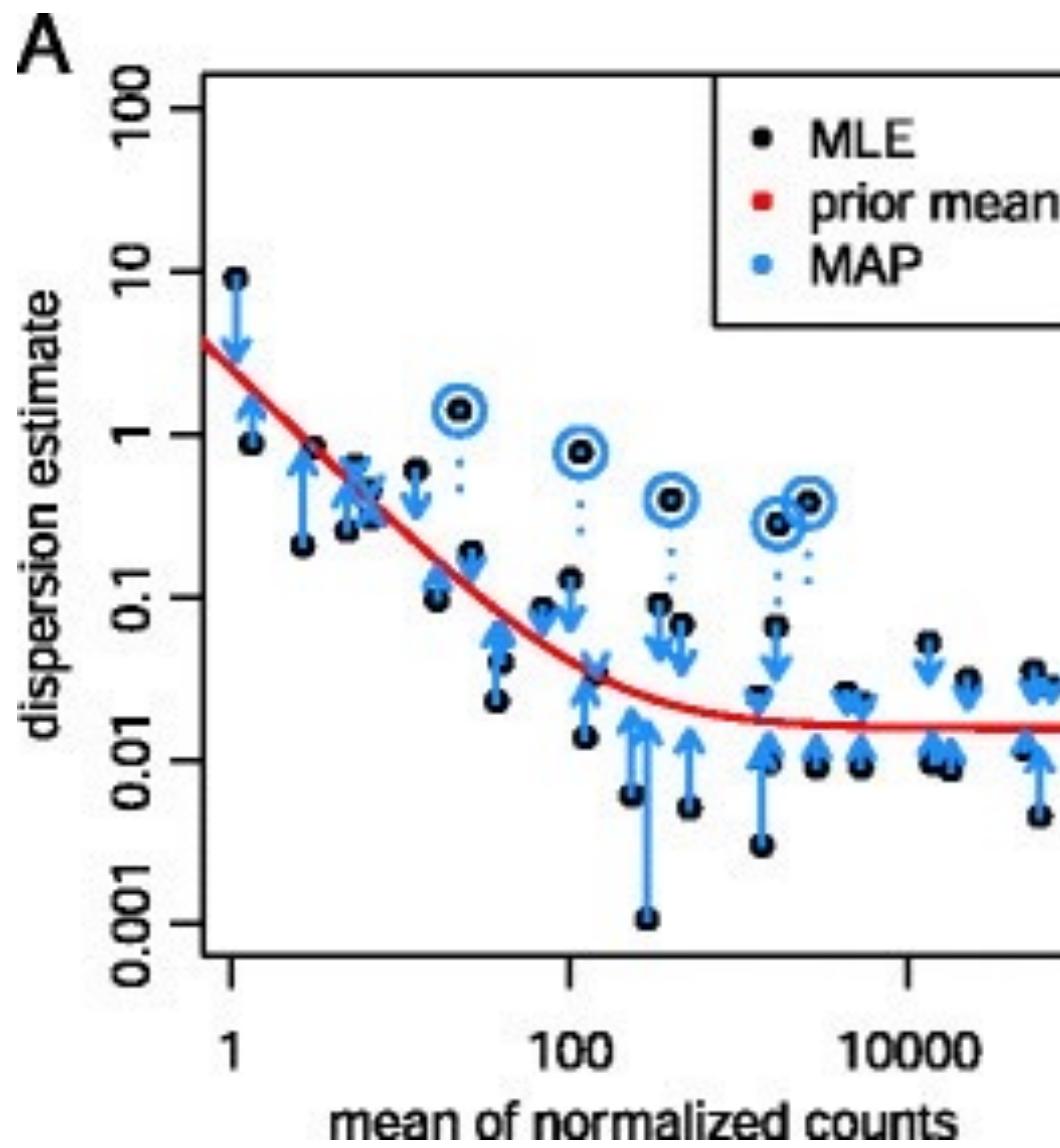
## 3. Variance of Counts:

where:

- $s_i$  is the size factor for sample  $i$ , accounting for differences in sequencing depth.
- $X_i$  is the design matrix row corresponding to sample  $i$ .
- $\beta_g$  is the vector of coefficients for gene  $g$ .

$$\text{Var}(Y_{ij}) = \mu_{ij} + \alpha_i \mu_{ij}^2$$

# DESeq2: Dispersion estimation



# Design Matrix

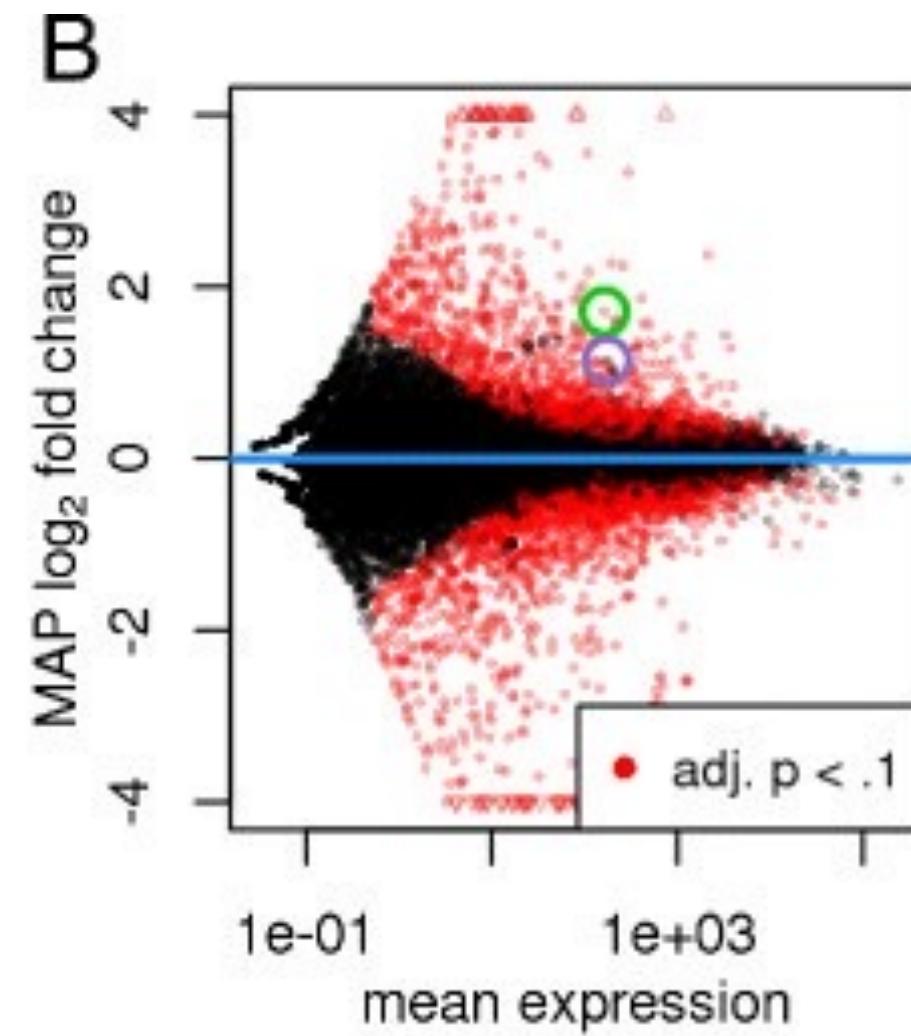
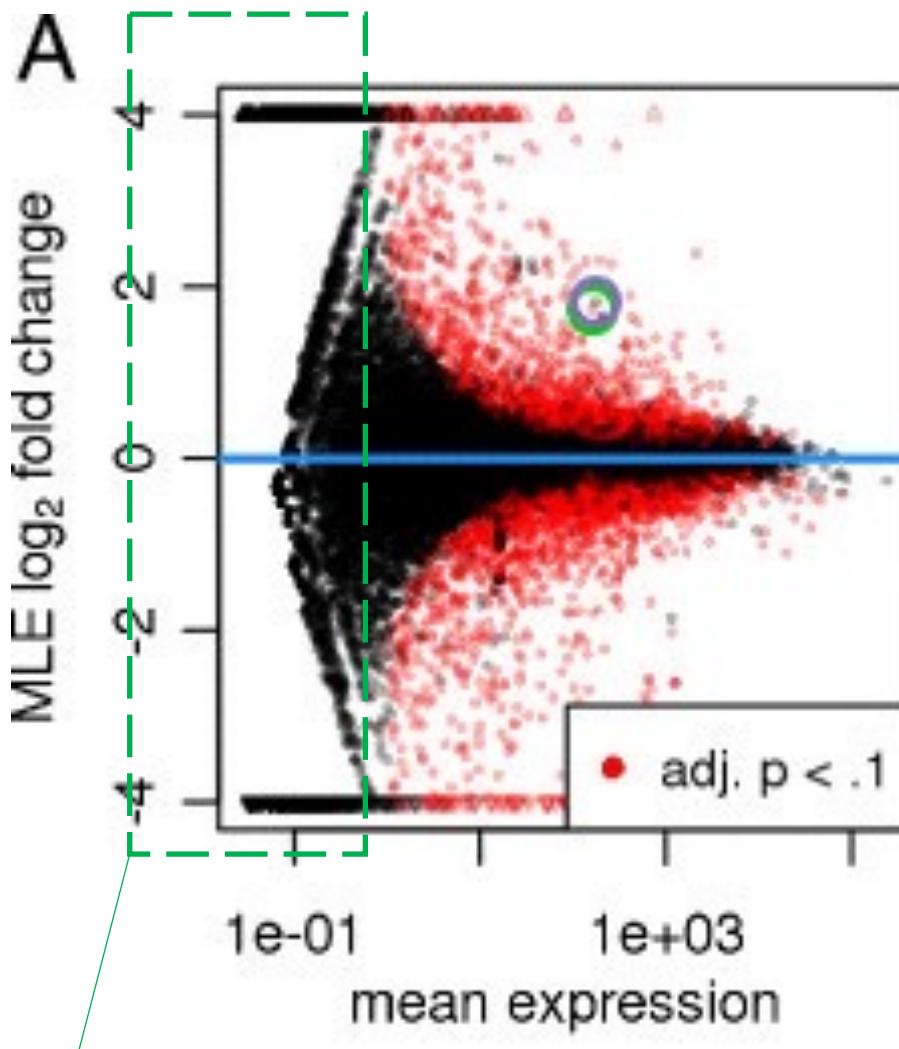
	(Intercept)	conditionInfection	conditionTreatment
Control.1	1	0	0
Control.2	1	0	0
Control.3	1	0	0
Infection.1	1	1	0
Infection.2	1	1	0
Infection.3	1	1	0
B5uM.1	1	0	1
B5uM.2	1	0	1
B5uM.3	1	0	1

attr(,"assign")  
[1] 0 1 1  
attr(,"contrasts")  
attr(,"contrasts")\$condition  
[1] "contr.treatment"

# DESeq2:

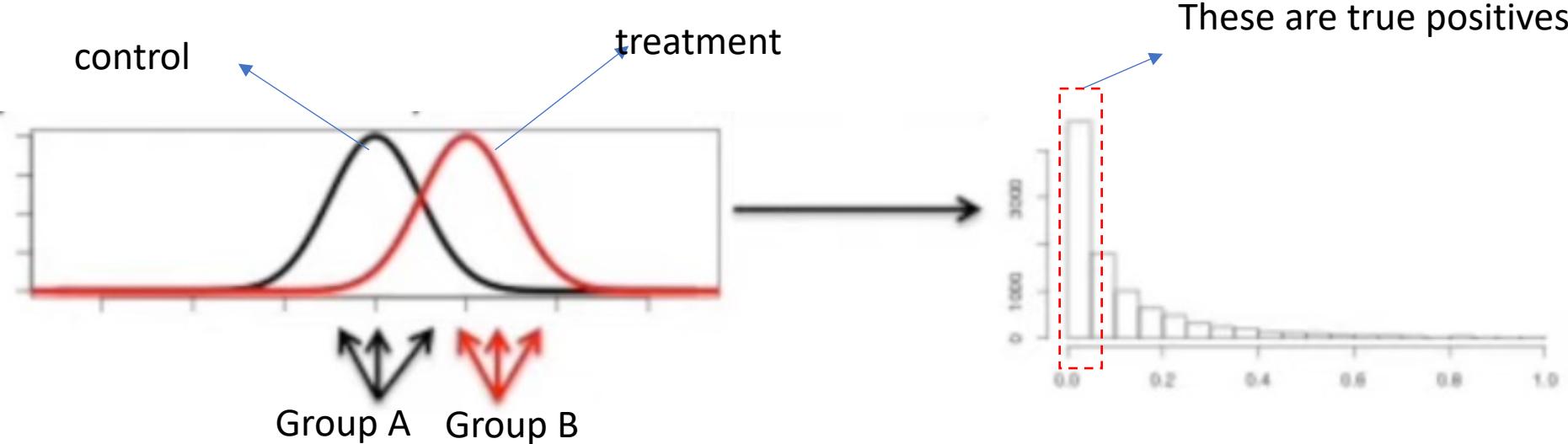
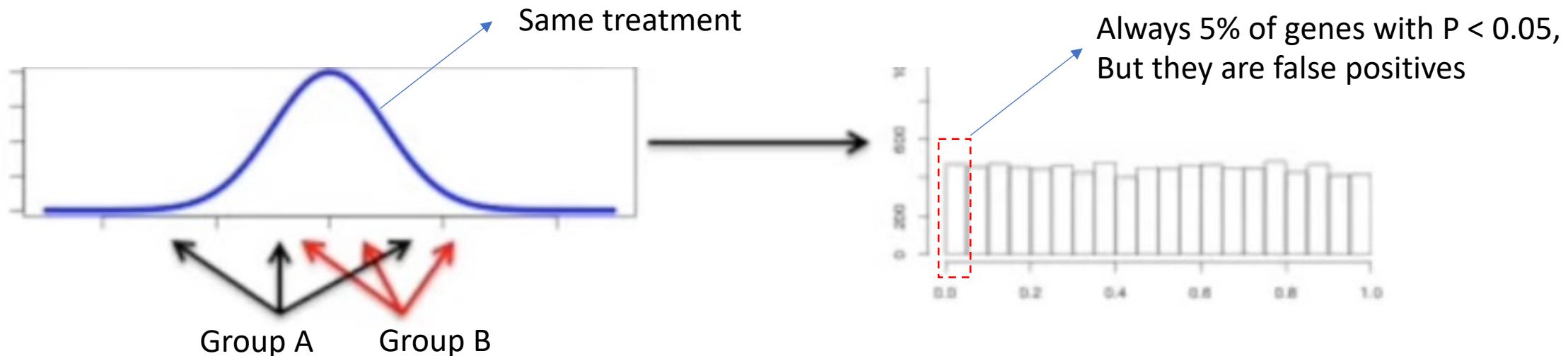
```
dds <- DESeqDataSetFromMatrix(countData = count_data_raw,  
                               colData = sample_info,  
                               design = ~ condition)  
  
# Run the DESeq pipeline  
dds <- DESeq(dds)  
  
# Get the results  
res <- results(dds, contrast=c("condition","Infection","Control"))
```

# DESeq2: shrink fold changes



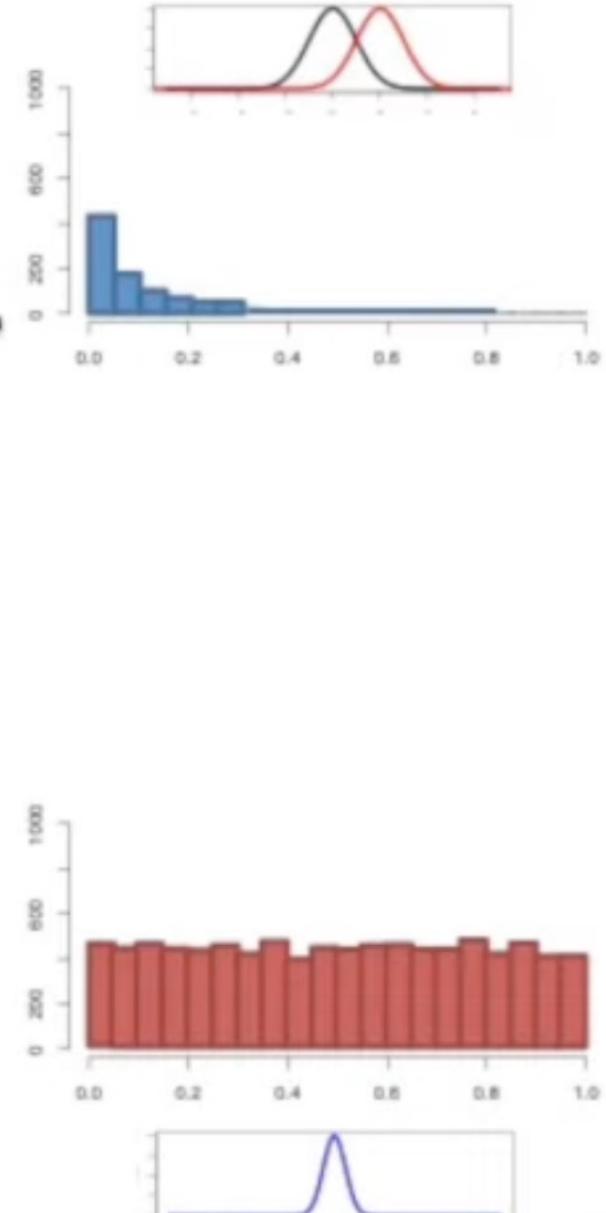
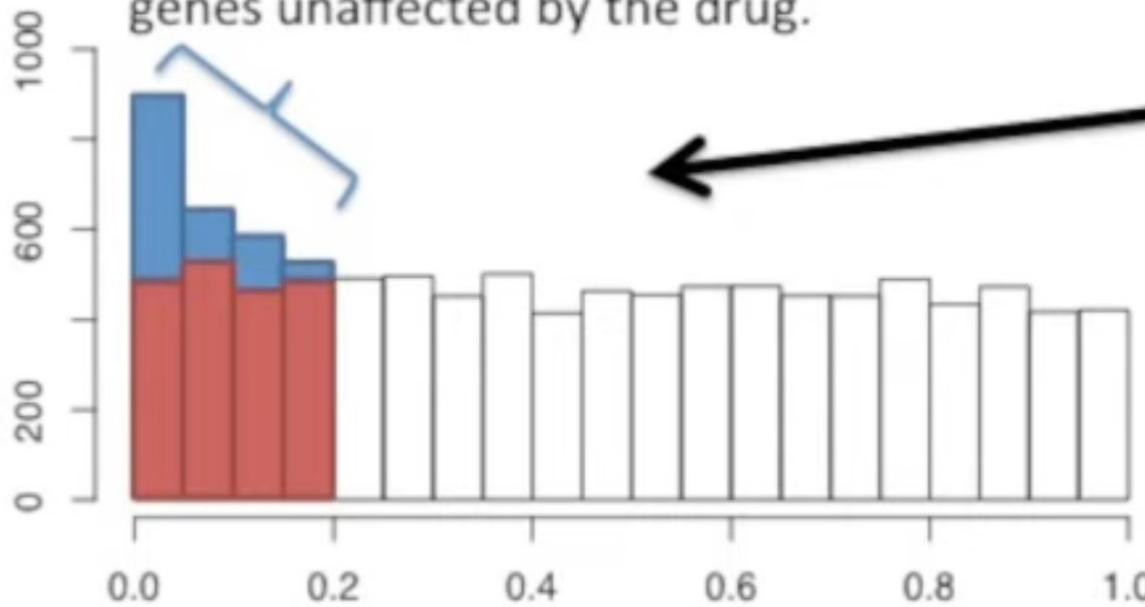
The fold changes of those lowly expressed genes are adjusted

# FDR: why do we need to adjust p value?

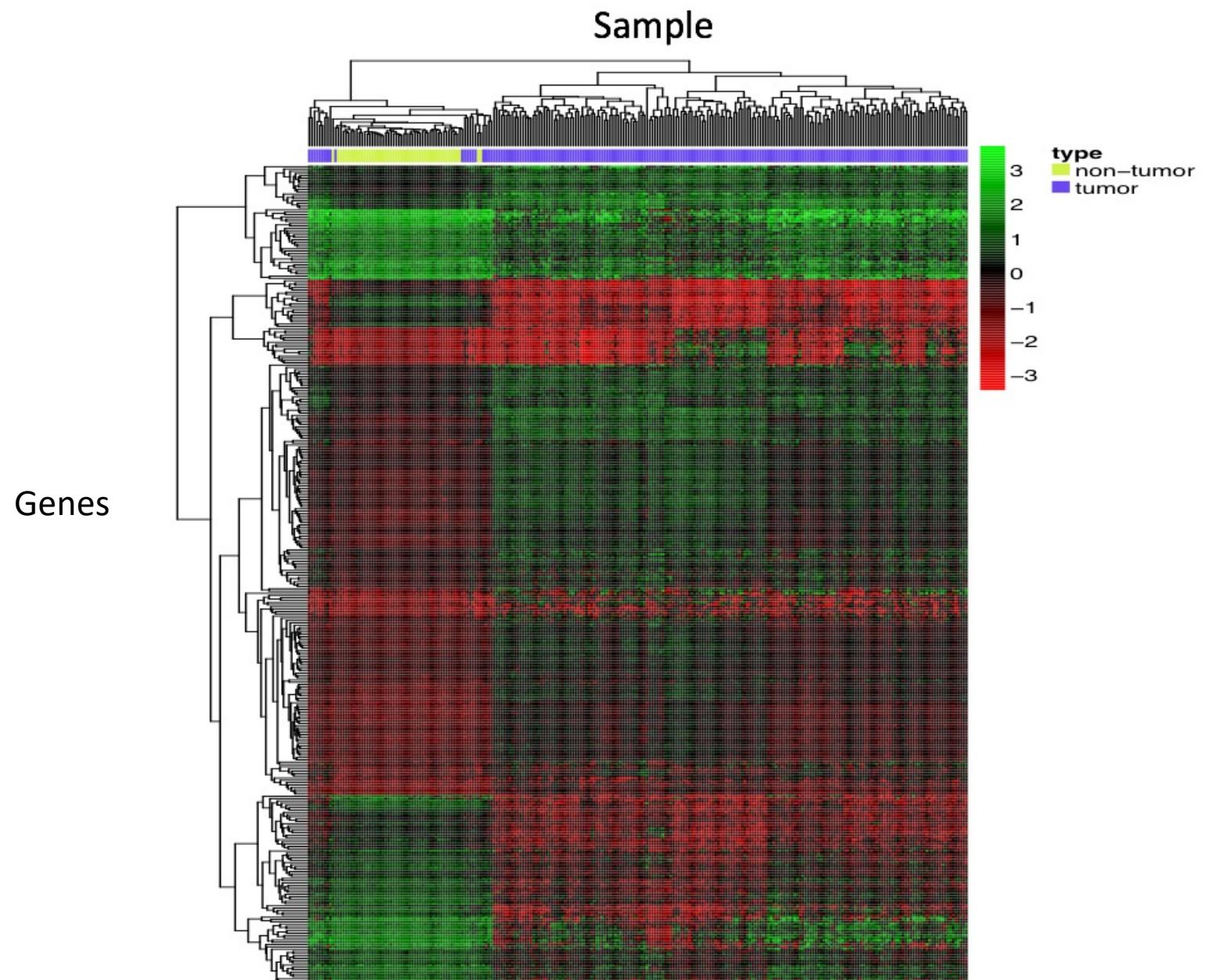


# FDR: how do we correct p value?

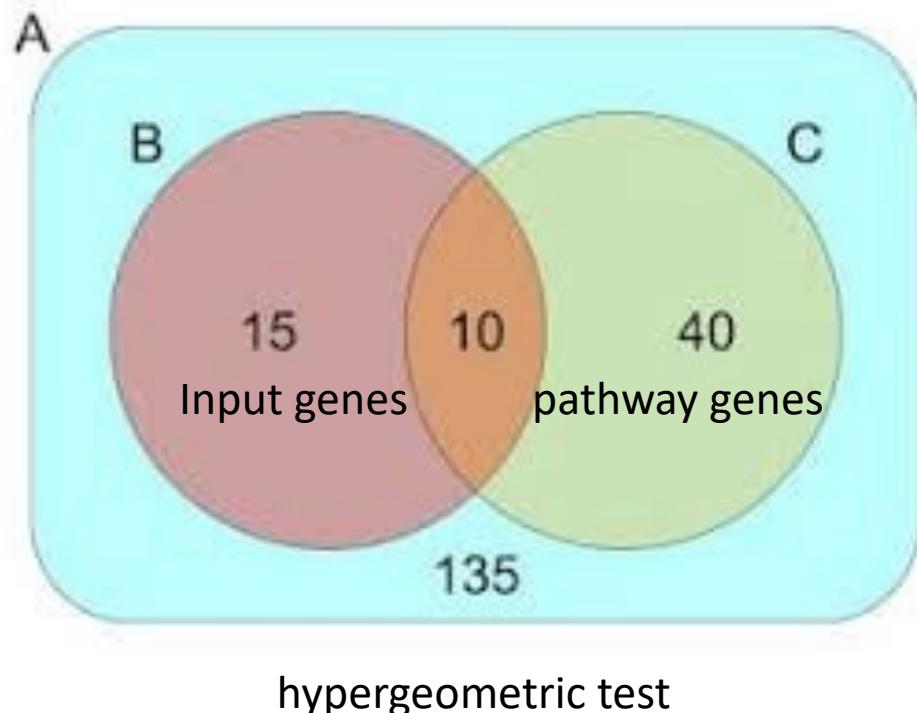
The p-values on the left side are a mixture from genes affected and genes unaffected by the drug.



# Clustering



# Gene Set Enrichment Analysis



$$p_X(k) = \Pr(X = k) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}},$$

**N:** Total number of genes in the background. **K:** Number of genes in the pathway. **n:** Number of genes in your gene list of interest. **k:** Number of genes in the overlap between your gene list and the pathway.

## GO Biological Process 2018

Bar Graph

Table

Clustergram

Appyter

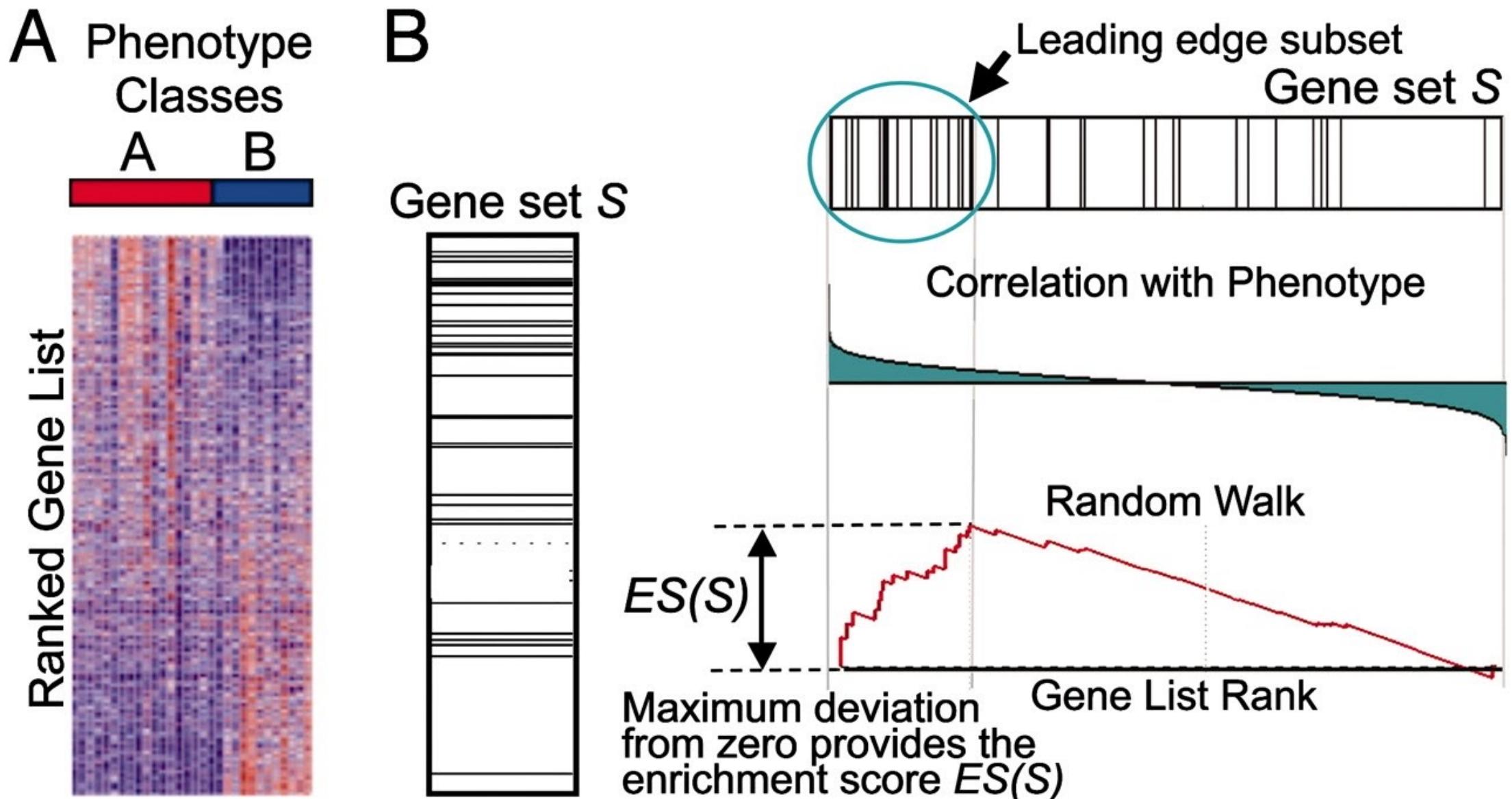


Click the bars to sort. Now sorted by p-value ranking.

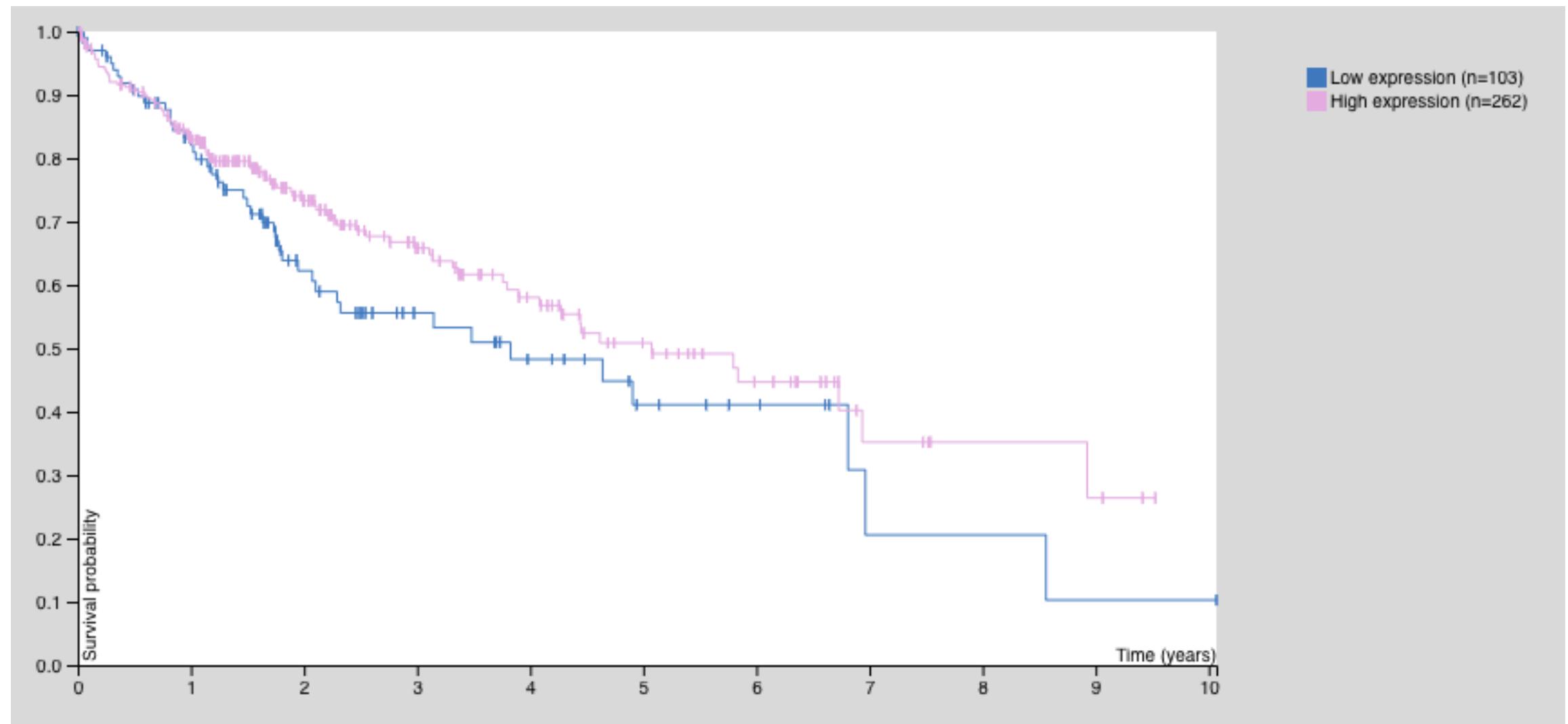
SVG PNG JPG

- mitotic sister chromatid segregation (GO:0000070)
- mitotic spindle organization (GO:0007052)
- microtubule cytoskeleton organization involved in mitosis (GO:1902850)
- mitotic cell cycle phase transition (GO:0044772)
- mitotic metaphase plate congression (GO:0007080)
- metaphase plate congression (GO:0051310)
- regulation of mitotic cell cycle phase transition (GO:1901990)
- cytoskeleton-dependent cytokinesis (GO:0061640)
- mitotic nuclear division (GO:0140014)
- mitotic cytokinesis (GO:0000281)

# Gene Set Enrichment Analysis

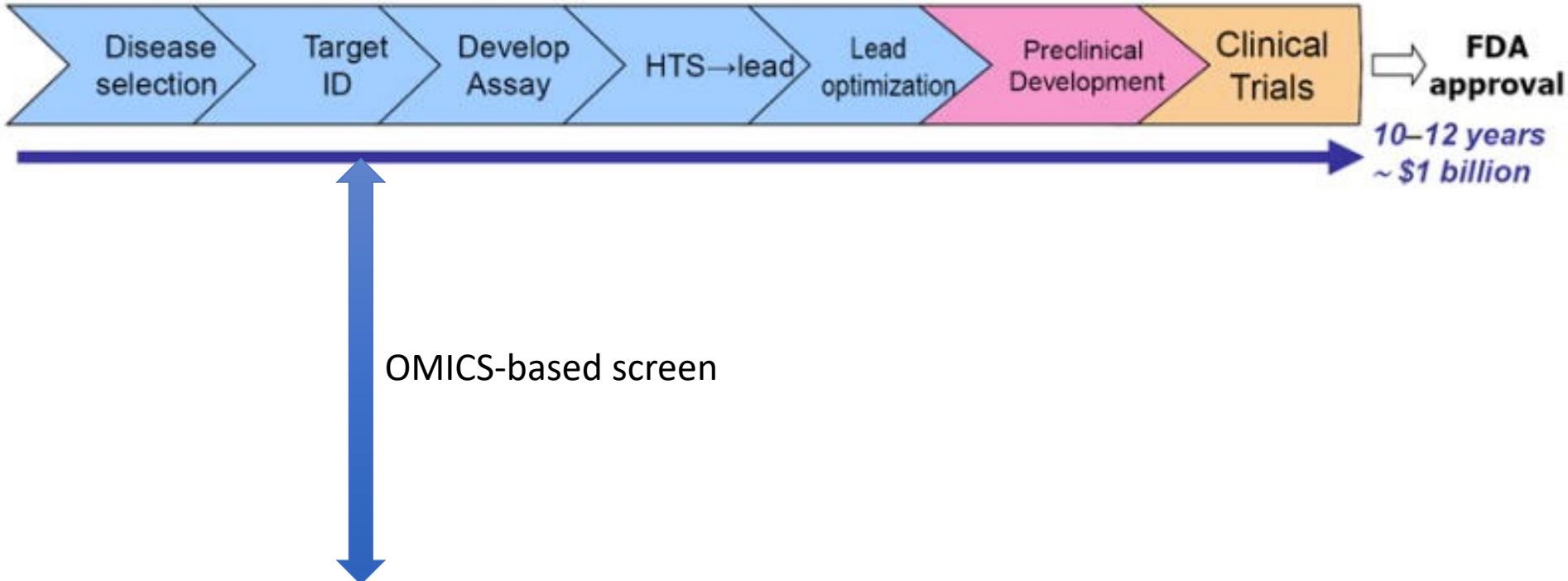


# Survival analysis

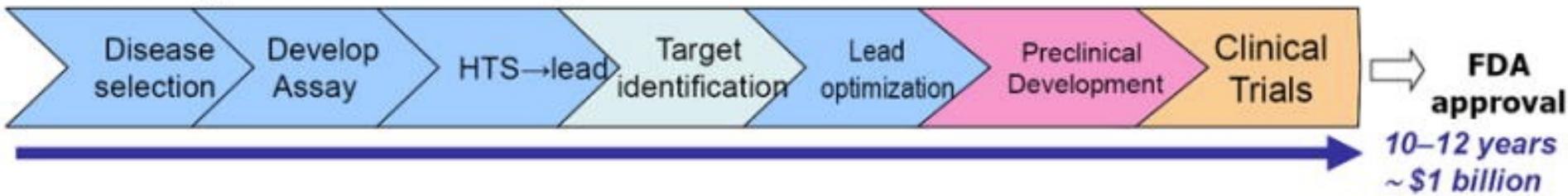


Kaplan–Meier plot

### A. Molecular target screen-based approach:

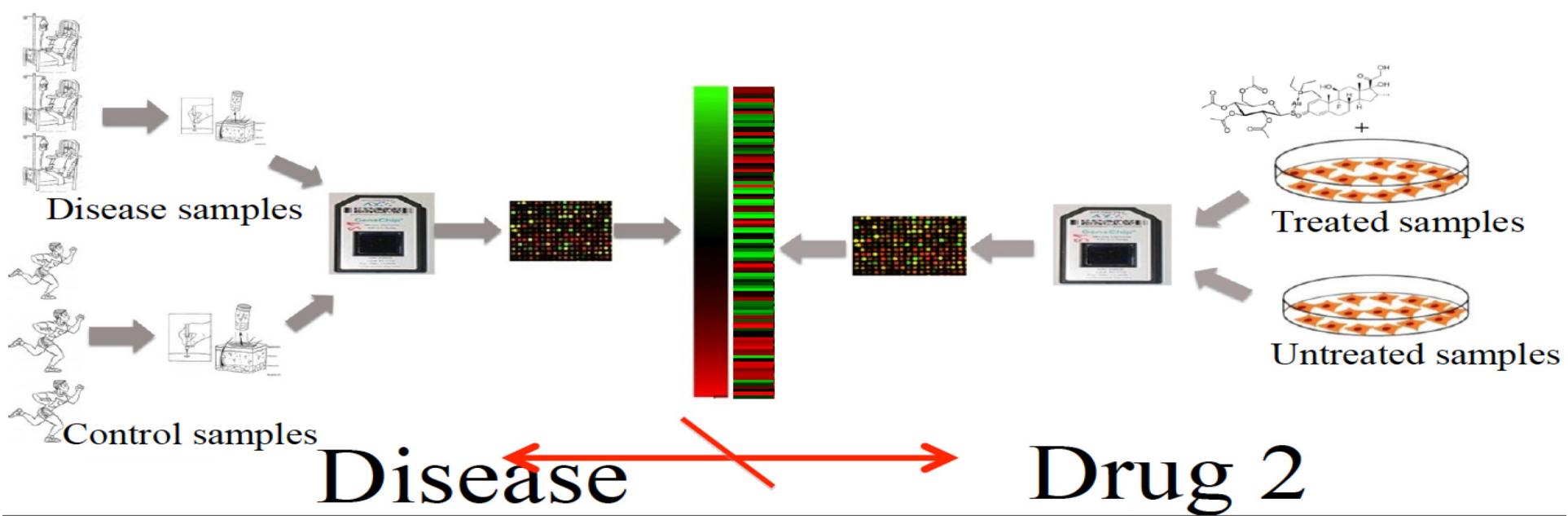
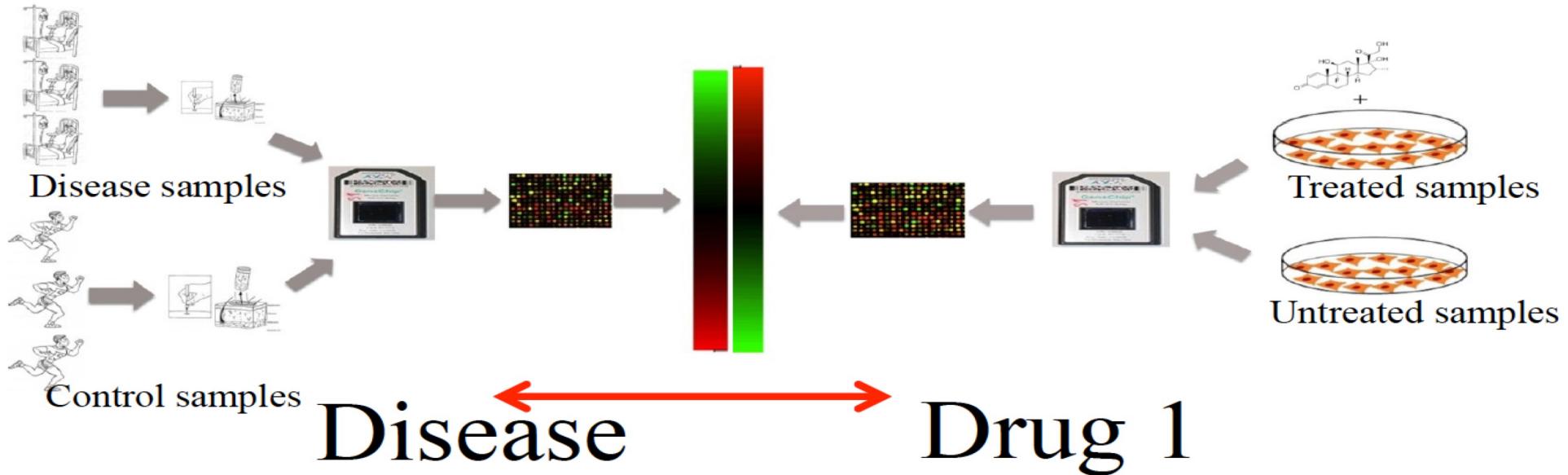


### B. Phenotypic screen-based approach:

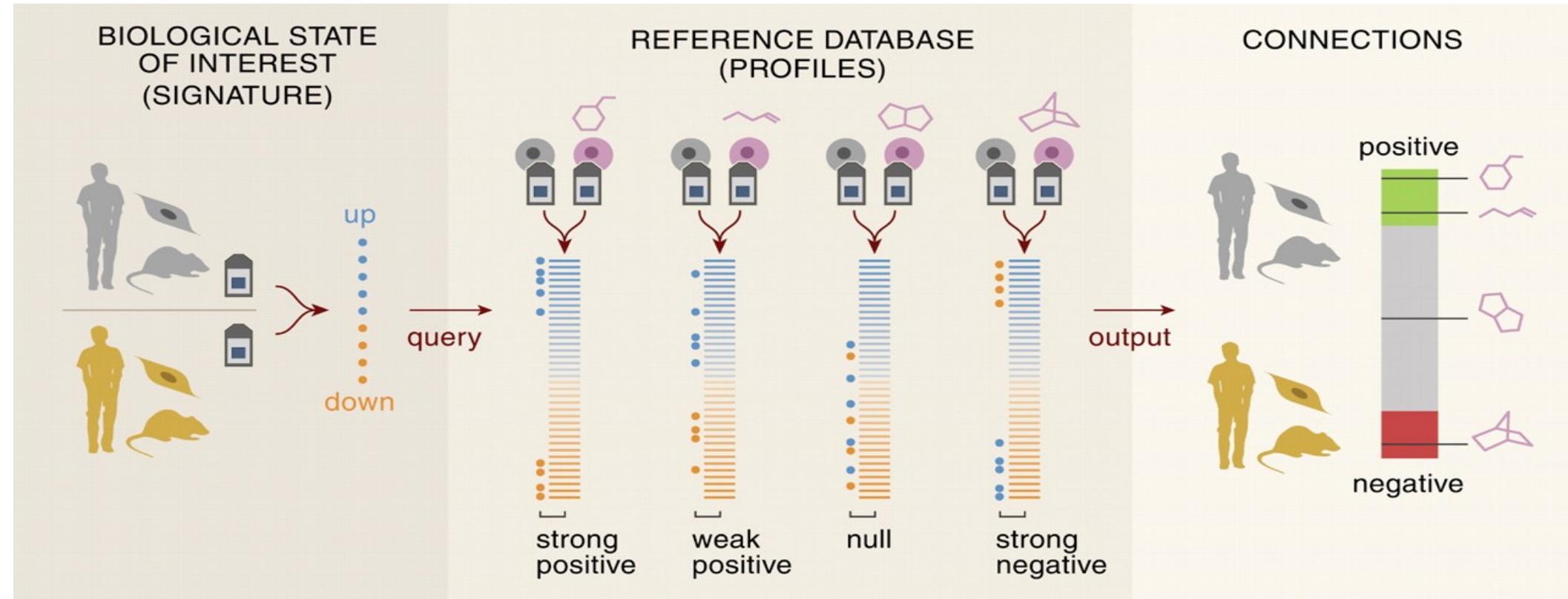


Phenotype = f(omics)

# Systems-based drug discovery approach



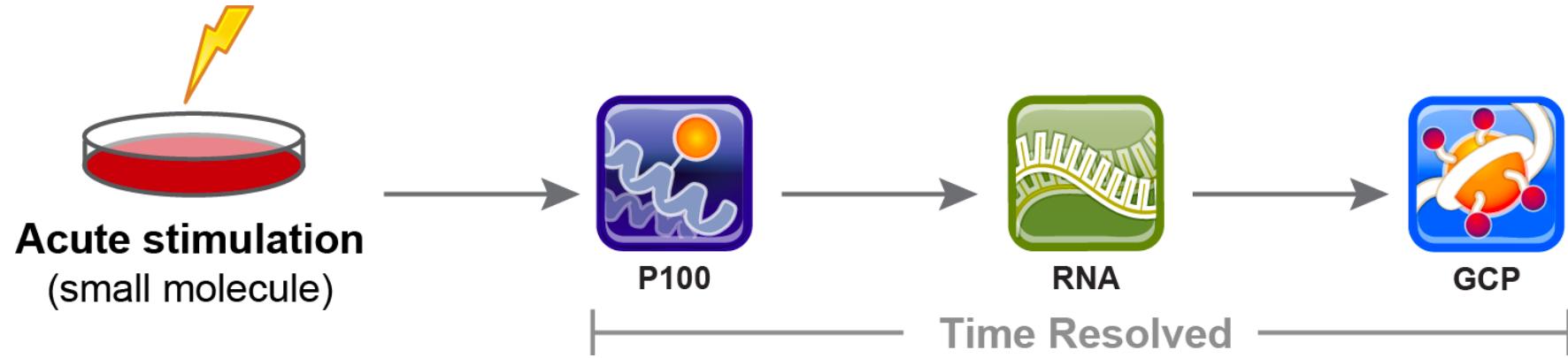
# Connectivity Map 1.0



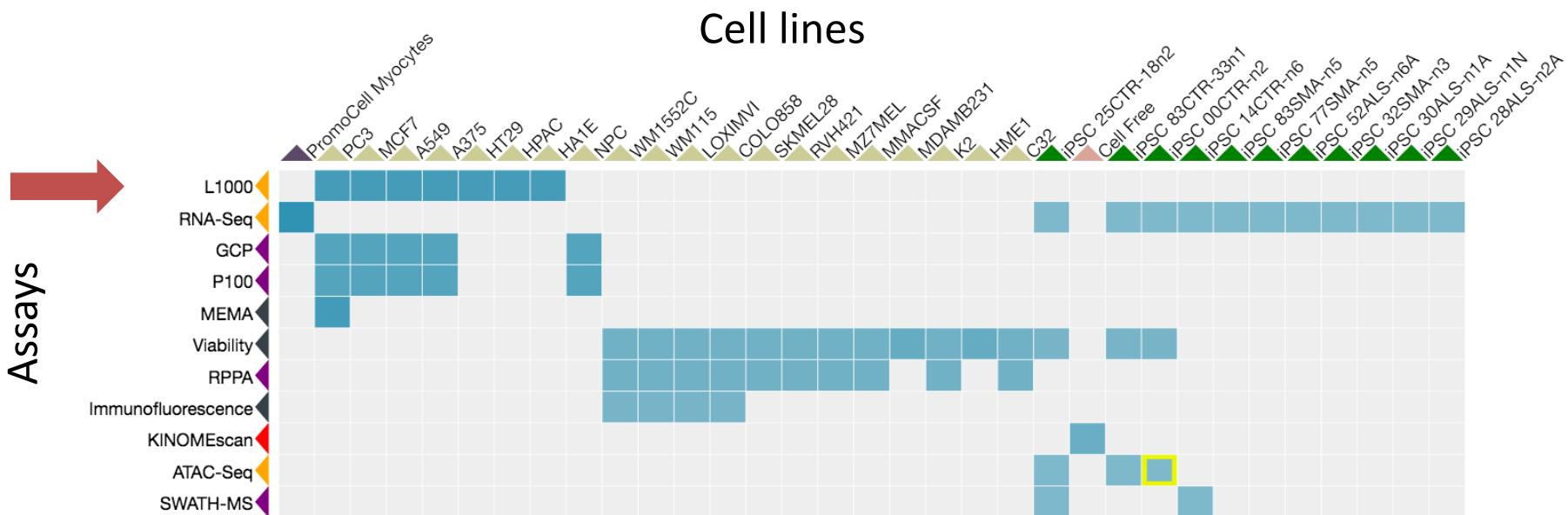
Connectivity Map Score/ reversal score: summary score of enrichment of disease up/down genes on the ranked drug profile

500 gene expression profiles → 6K (CMap) → 1M (LINCS)

# Connectivity Map 2.0 (LINCS)



# Library of Network-Based Cellular Signatures (LINCS)



# 421 Datasets



## 41847 Small Molecules



## 1127 Cells



978 Genes



1469 Proteins  
/155 Peptide  
Probes



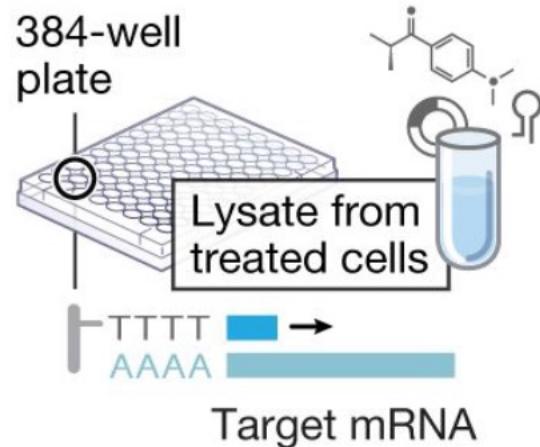
117 Antibodies

# Cost of profiling 1,000,000 samples

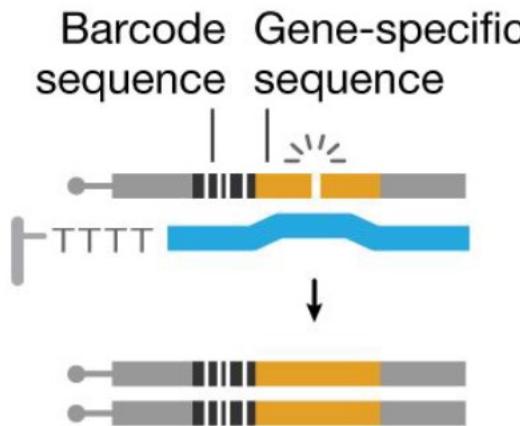
- If using RNAseq
  - $1,000,000 * \$200 \text{ per sample} = \$200\text{M}$
- If using LINCS 1000
  - $1,000,000 * \$2 \text{ per sample} = \$2\text{M}$

# Connectivity Map 2.0 (LINCS L1000)

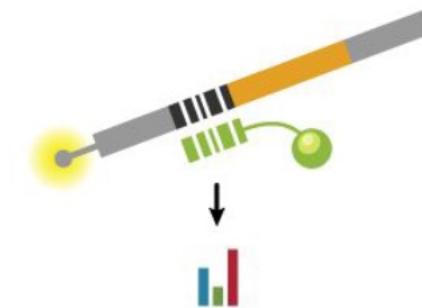
Capture and reverse transcribe mRNA



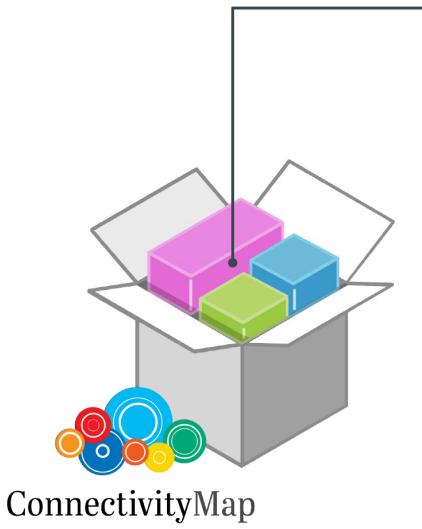
Ligate probes and amplify with biotinylated primers



Hybridize to beads and stain with SAPE



Identify gene and quantify expression



## ConnectivityMap

Subramanian et al, Cell 2017

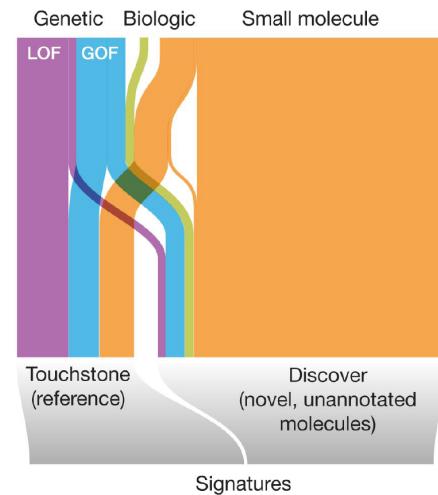
# L1000

---

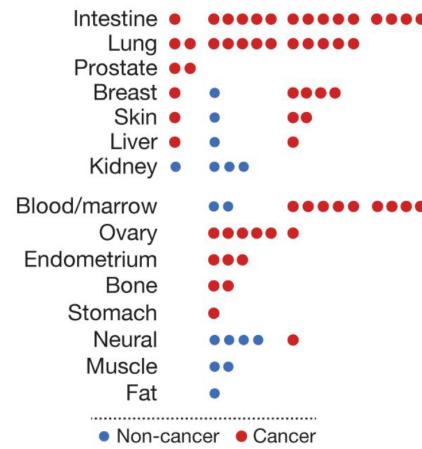
## Pattern-recognition algorithms

**1.3M profiles**

28,000 perturbagens



## 80 Cell lines



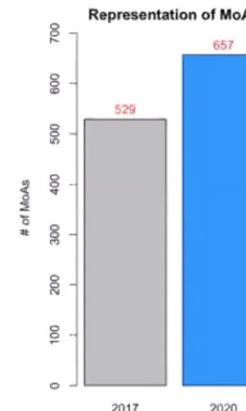
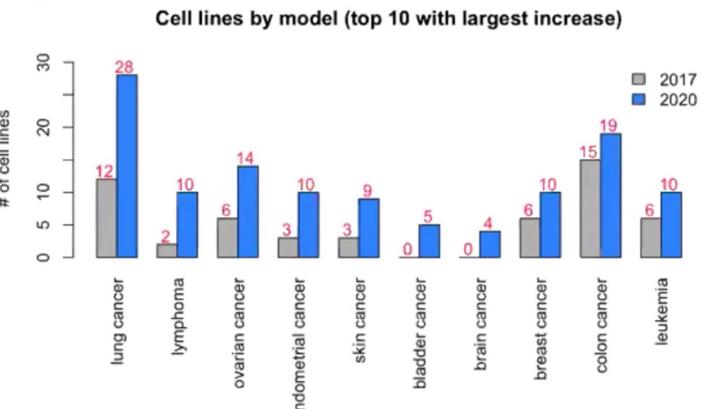
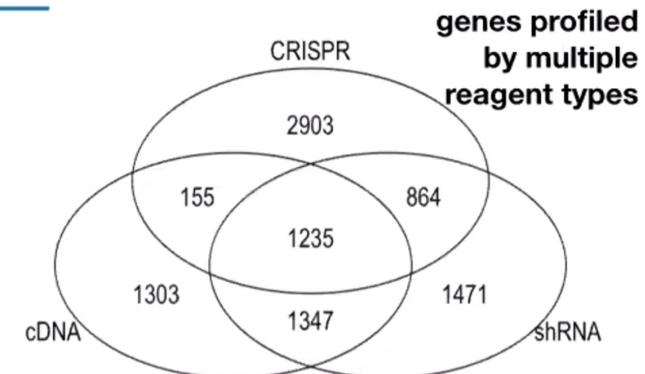
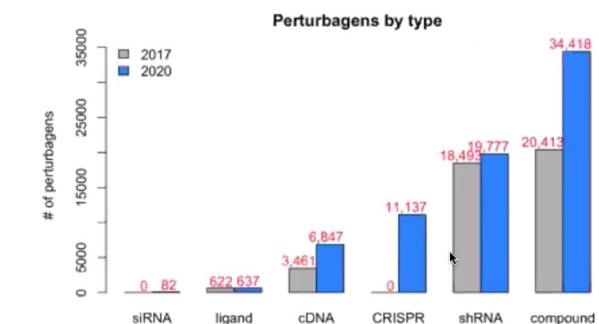
# LINCS L1000 phases

2017: Phase I: 1,319,138 (GSE92742) + Phase II: 354,123 (GSE70138)

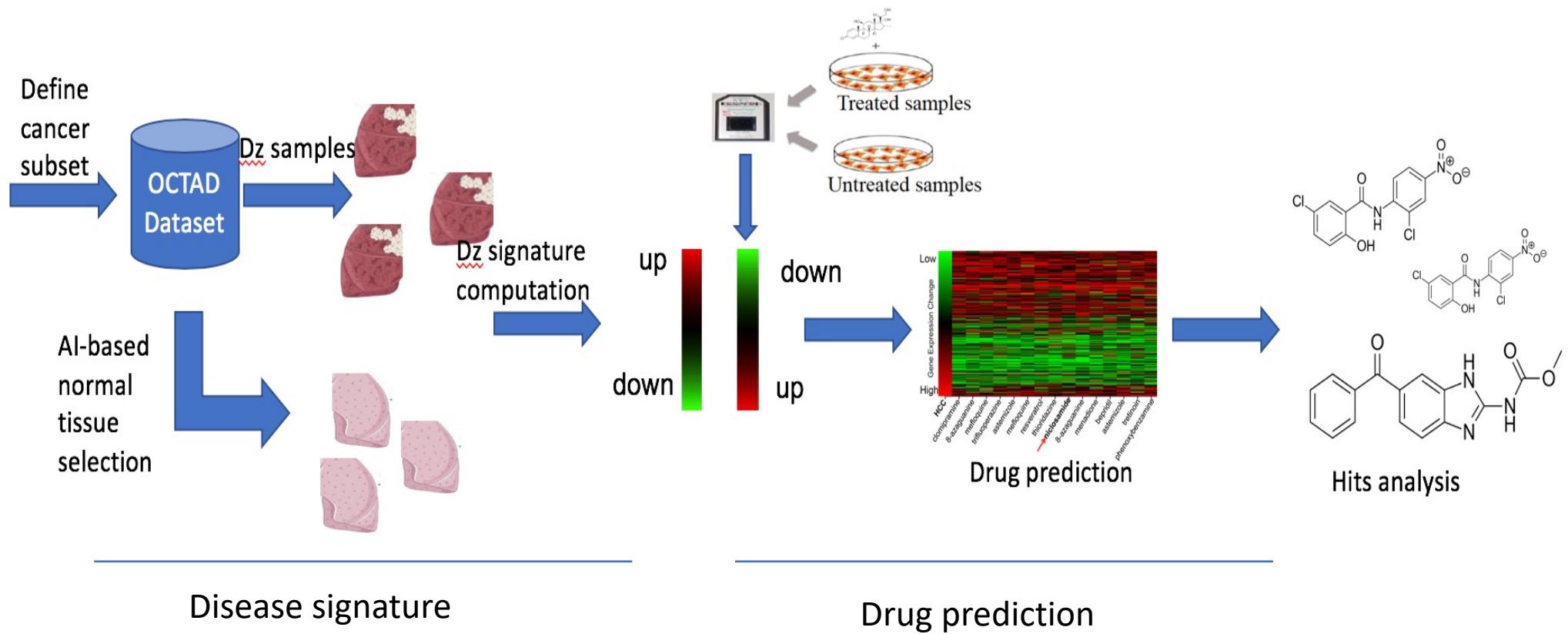
2020: 3M (manuscript unpublished)

## Overview of updated resource expansion upon previous release

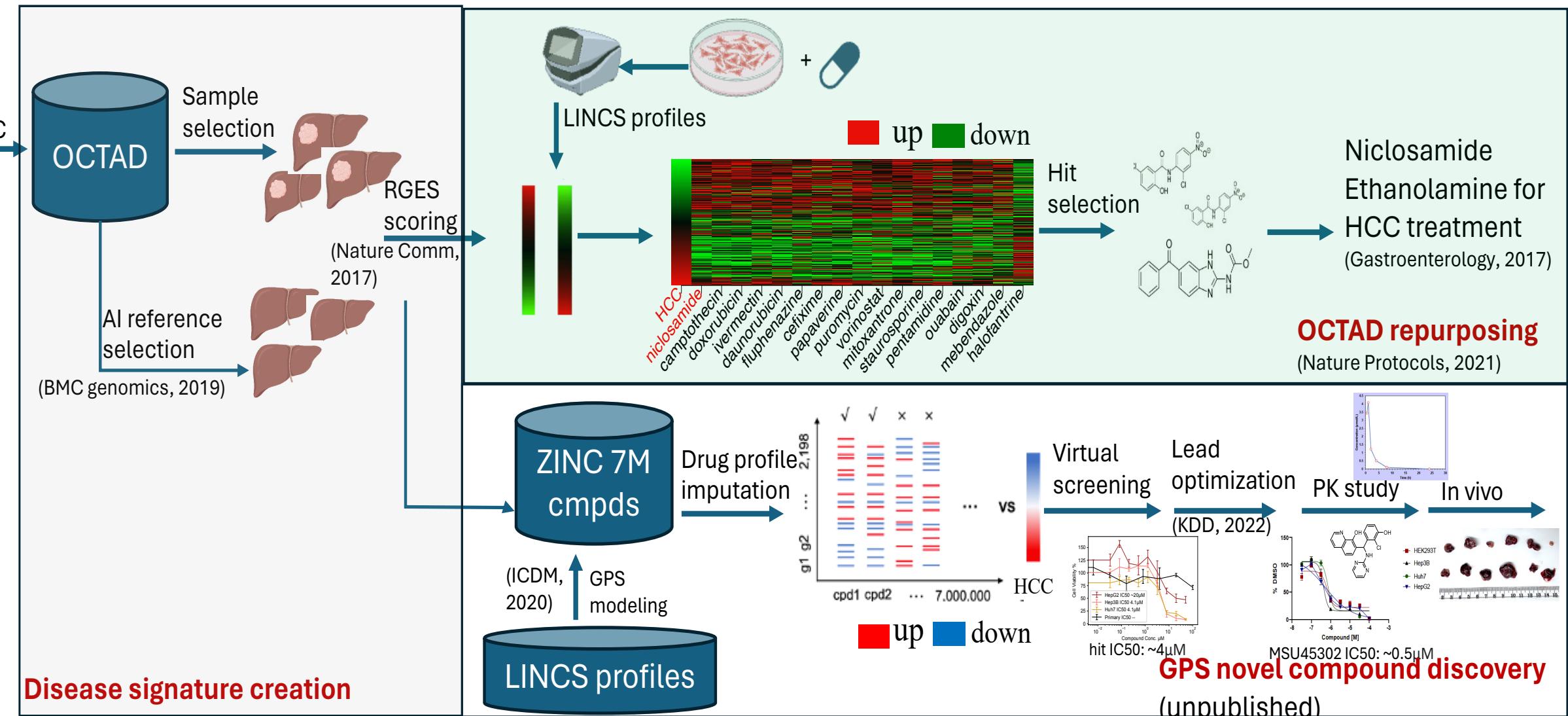
# perturbagens	81,979
# compounds	33,609
# MoAs	657
# genes	9,288
# cell contexts	240 (12 primary)
# profiles	3.02M
# signatures	1.16M



# Transcriptomics based drug repurposing pipeline



# Transcriptomics-based drug discovery platform



# OCTAD Web portal

C

OCTAD--from genomic features to therapeutic candidates in four steps

1. Case      2. Control      3. Signature      4. Drugs      Login

Job management      Case page      Control page      signature page      prediction page

Disease Name: liver hepatocellular carcinoma

Sample selection

Sample list

Mutation	sample Id	Type	Site	Metastatic Site	Cancer
Gain	TCGA-2V-A95S-01	primary	LIVER		liver hepatocellular carcinoma
Loss	TCGA-2Y-A9GS-01	primary	LIVER		liver hepatocellular carcinoma
	TCGA-2Y-A9GT-01	primary	LIVER		liver hepatocellular carcinoma
	TCGA-2Y-A9GU-01	primary	LIVER		liver hepatocellular carcinoma
	TCGA-2Y-A9GV-01	primary	LIVER		liver hepatocellular carcinoma

Showing 1 to 5 of 379 entries      379 rows selected

1. Select case samples      2. Select control samples      3. Create disease signature      4. Predict drugs/ targets

Summary      Save      Previous      Next

#Billy Zeng, #Benjamin S. Glicksberg, #Patrick Newbury, #Evgenii Chekalin, Jing Xing, Ke Liu, Anita Wen, Caven Chow, Bin Chen, OCTAD: an open workplace for virtually screening therapeutics targeting precise cancer patient groups using gene expression features accepted, Nature Protocols

# OCTAD R package

```
case_id=subset(phenoDF,cancer=='liver hepatocellular carcinoma'&sample.type == 'primary', select = c("sample.id"))
```

```
HCC_adjacent=subset(phenoDF,cancer=='liver hepatocellular carcinoma'&sample.type == 'adjacent'&data.source == 'TCGA', select = c("sample.id"))
```

```
res=diffExp(case_id,control_id,source='octad.whole',output=T,n_topGenes=10000,file='octad.counts.and.tpm.h5')
```

```
sRGES=runsRGES(res,max_gene_size=500,permutations=10000)
```

```
head(sRGES)
```