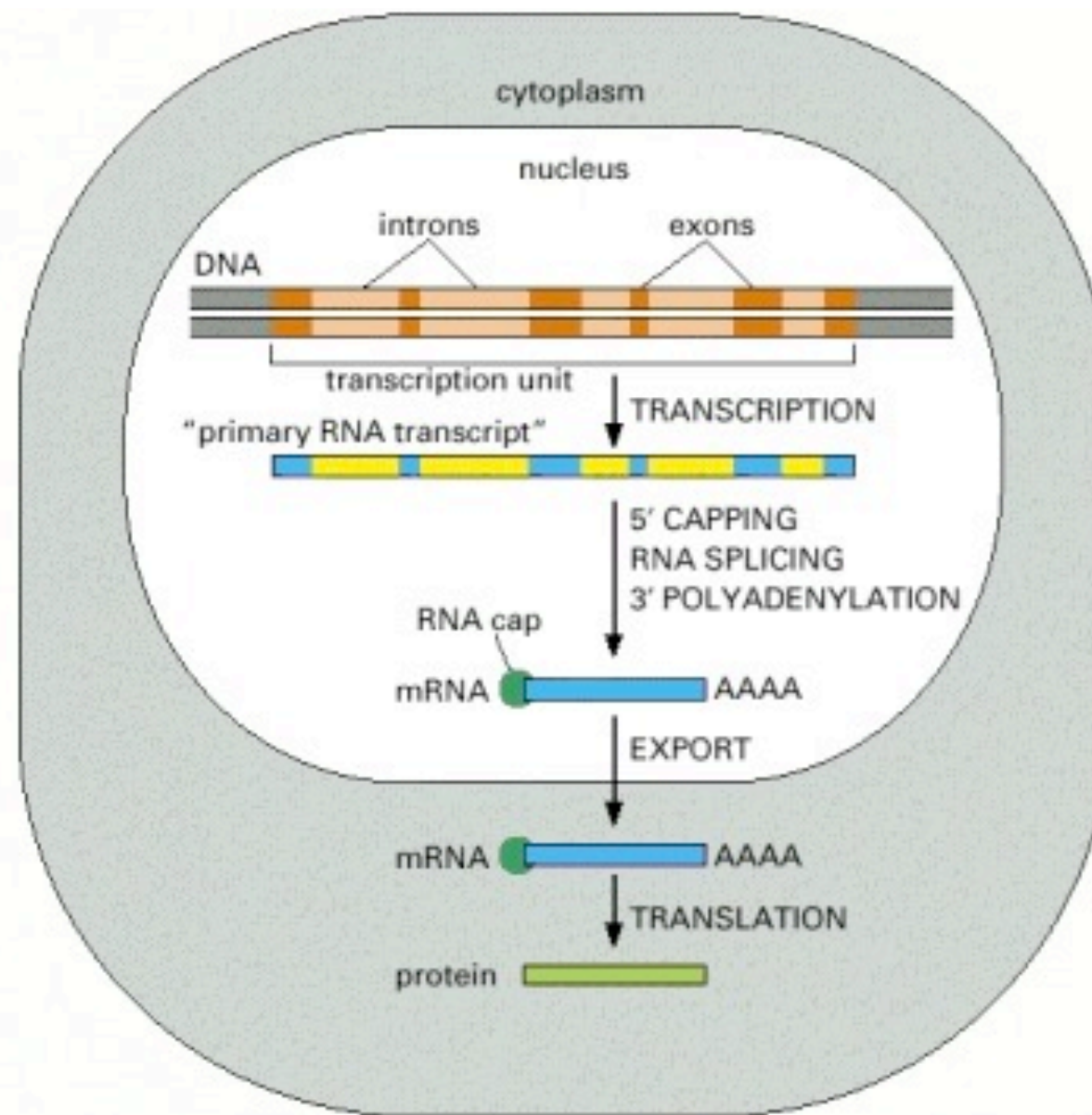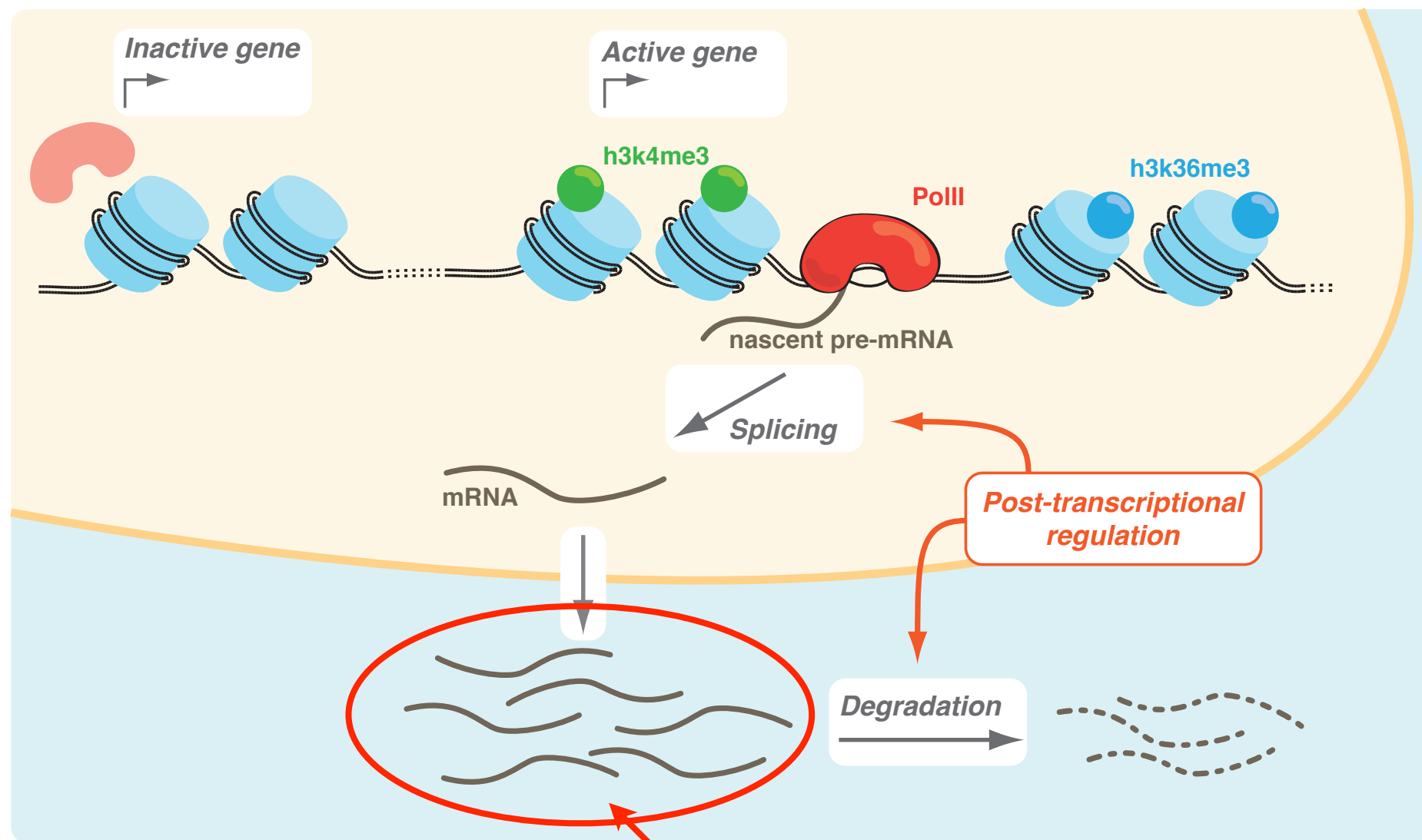# Transcription



Figures from Alberts and co, Mol Biol Cell.

# Transcription



Equilibrium pool of mRNA for each gene

# Simplest model

Transcription rate: $P(t)$ [#mRNA/time]
mRNA pool: $m(t)$ [#mRNA]
Degradation rate: $\gamma$ [1/time]

$$\dot{m}(t) = P(t) - \gamma m(t) \ .$$
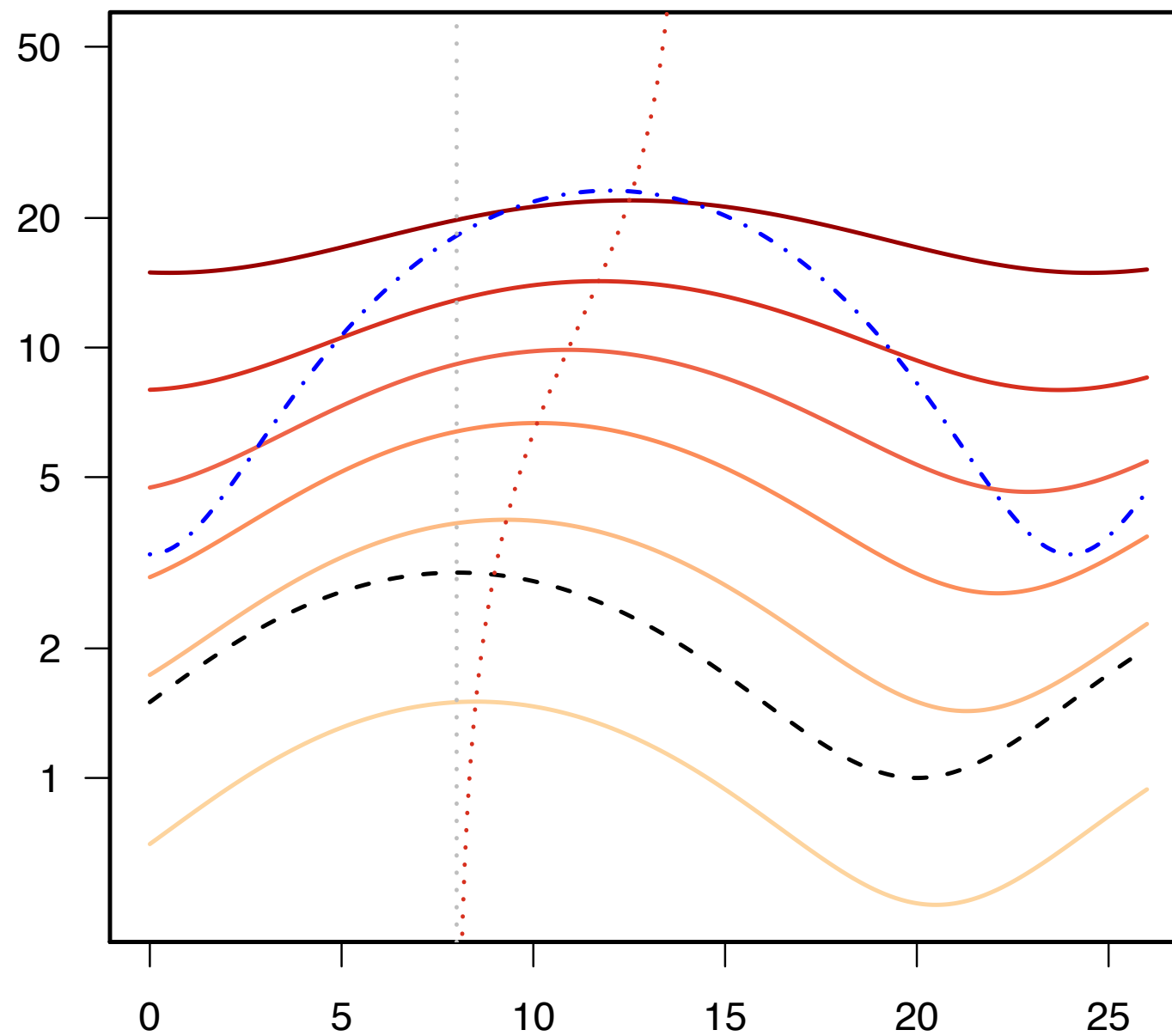
Case 1: $P(t) = P_0$

$$m(t) = \frac{P_0}{\gamma}(1 - e^{-\gamma t}) + e^{-\gamma t} m(0) \ \rightarrow \ \frac{P_0}{\gamma} \ .$$

Case 2: $P(t) = P_0 + \cos(\omega t)$

$$m(t) \ \rightarrow \ \frac{P_0}{\gamma} + \frac{1}{\sqrt{\gamma^2 + \omega^2}} \cos(\omega(t - \tau)) \ .$$

with $\sin \omega\tau = \frac{\omega}{\sqrt{\gamma^2 + \omega^2}}$

# Simplest model

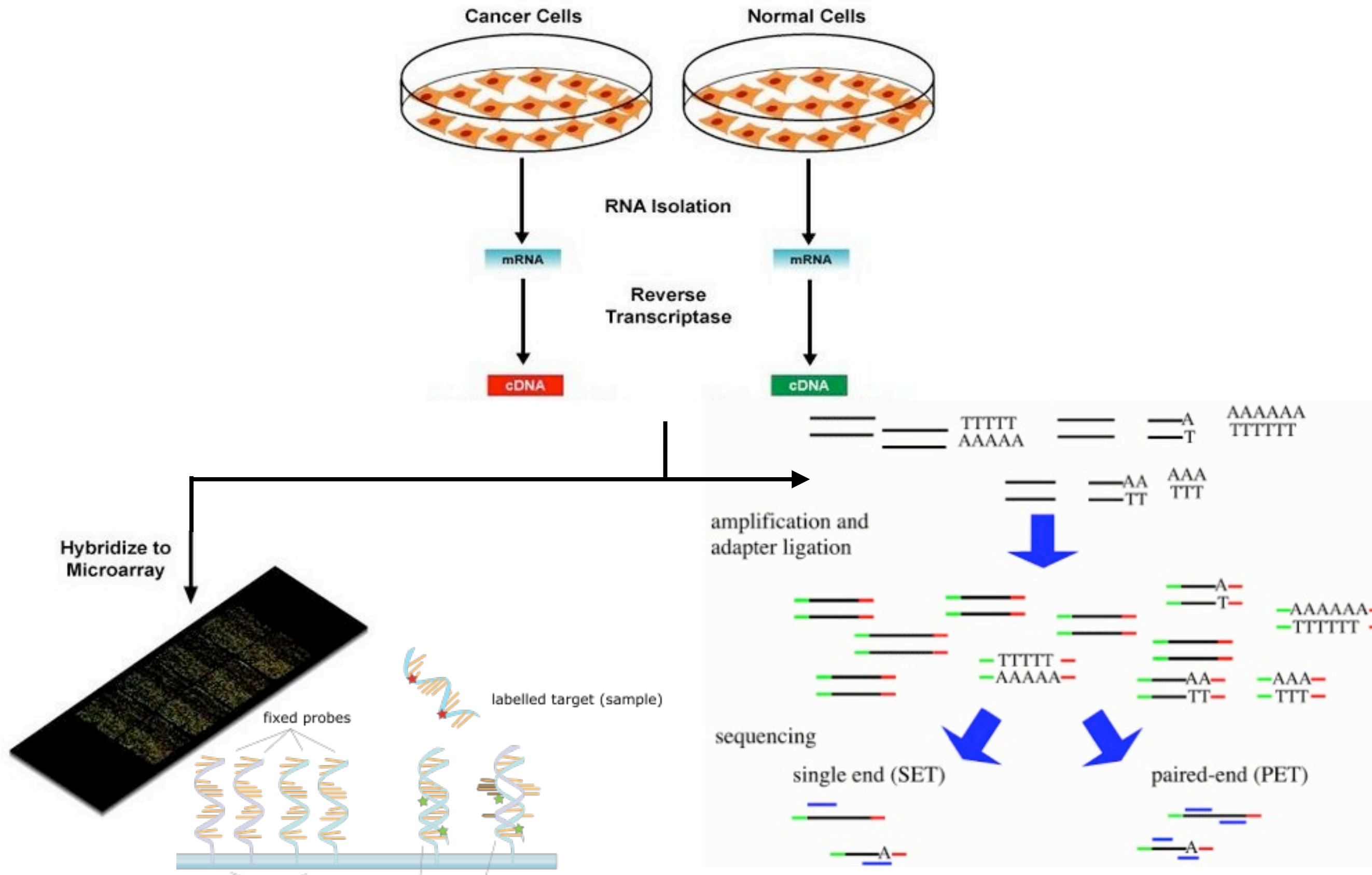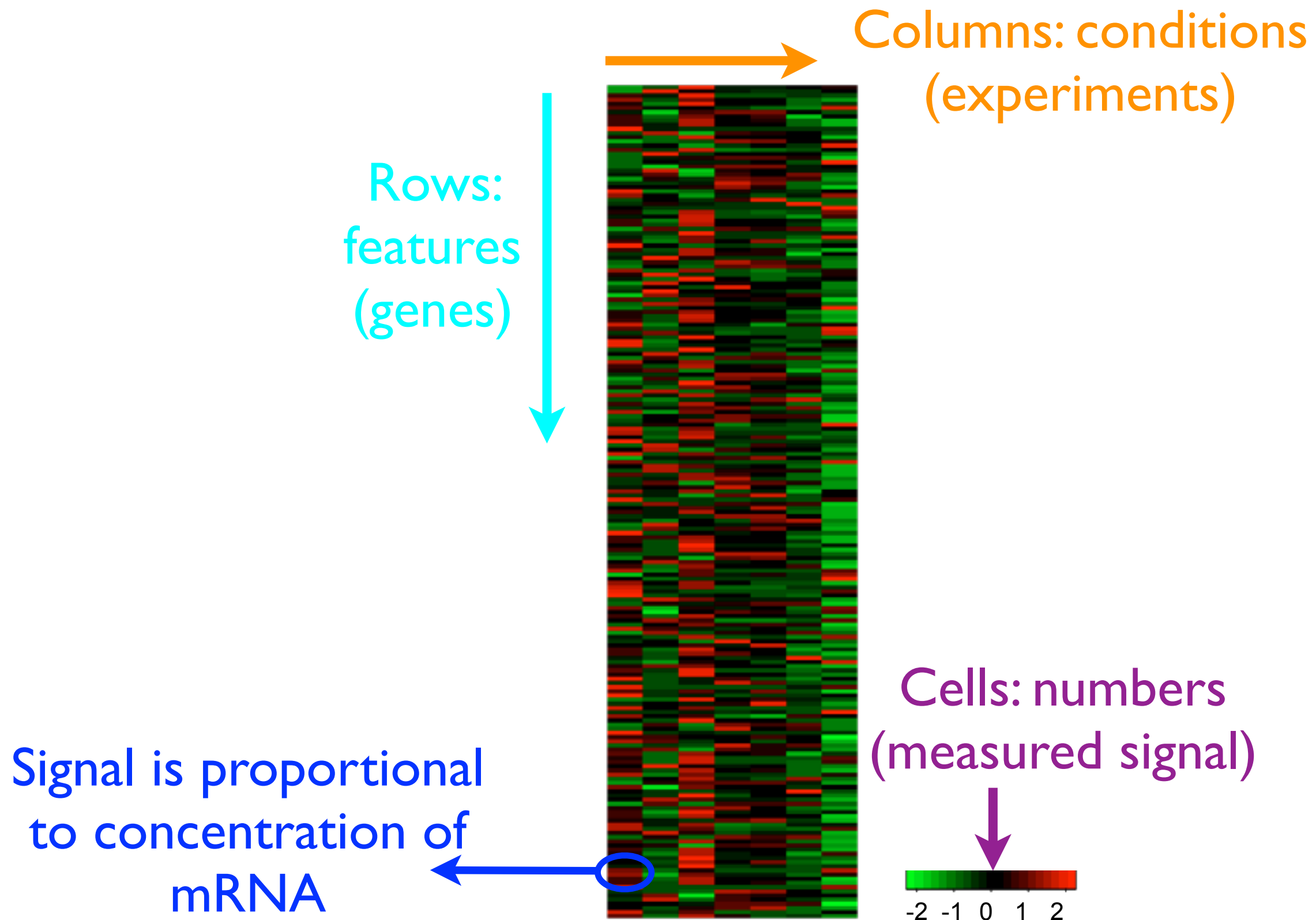# Measuring expression

# Measuring expression

- Goal: for each possible transcript **T**, estimate number of copies of **T** per cell
- We start from a large ($10^6$-$10^7$) population of cells and take a sample of the total mRNA pool
- We then measure a signal which is proportional to the abundance of **T** in the sample
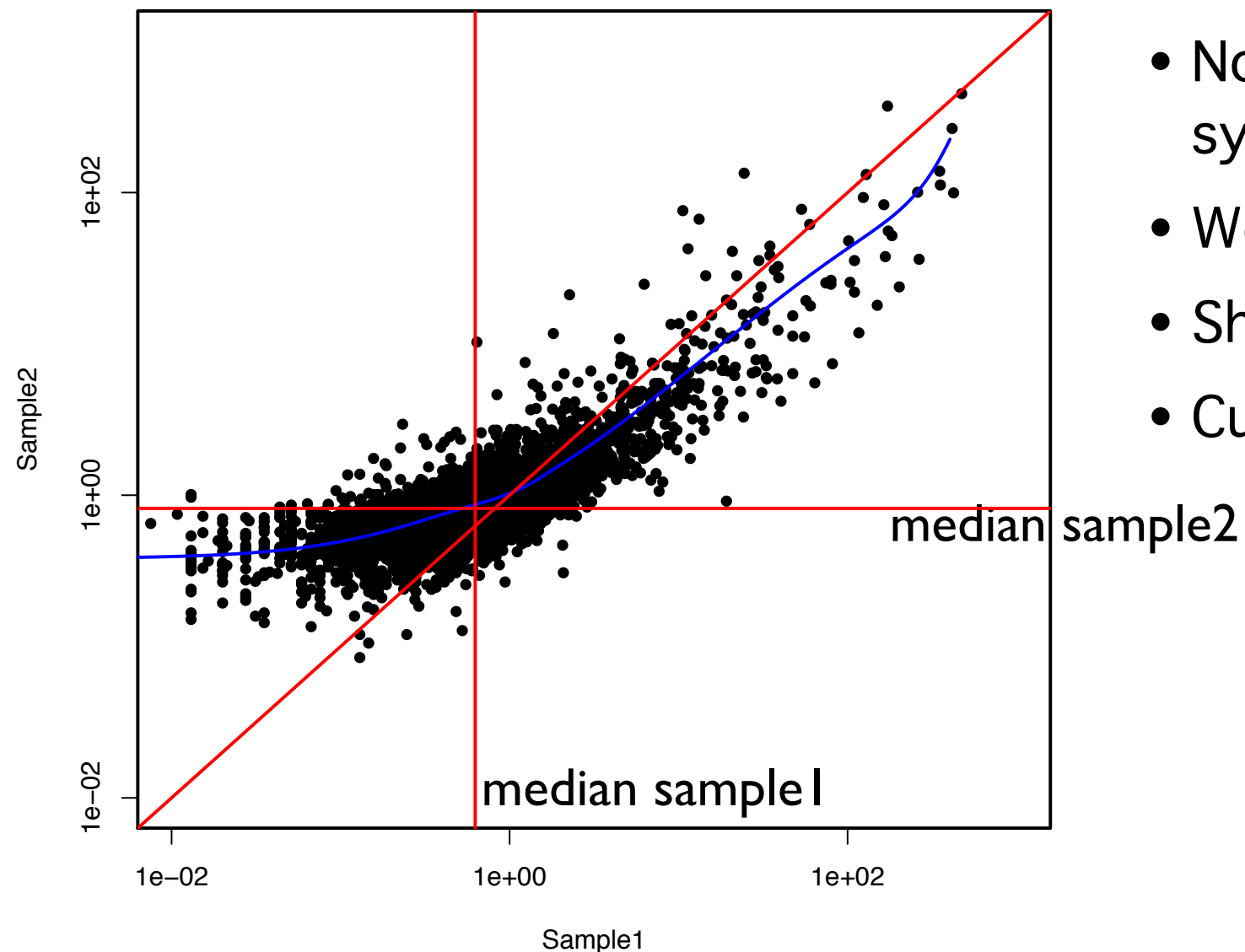- We generally do not know the calibration but we can deduce ratios:

$$\frac{\text{Expr}(g_1)}{\text{Expr}(g_2)} = \frac{N_{\text{copy}}(g_1)/N_{\text{cell}}}{N_{\text{copy}}(g_2)/N_{\text{cell}}}$$

$$= \frac{\alpha\text{Signal}(g_1)/N_{\text{cell}}}{\alpha\text{Signal}(g_2)/N_{\text{cell}}}$$

$$= \frac{\text{Signal}(g_1)}{\text{Signal}(g_2)}$$

# Output



Columns: conditions (experiments)

Rows: features (genes)

Cells: numbers (measured signal)

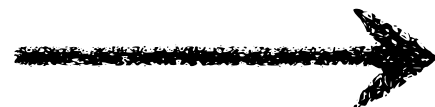Signal is proportional to concentration of mRNA
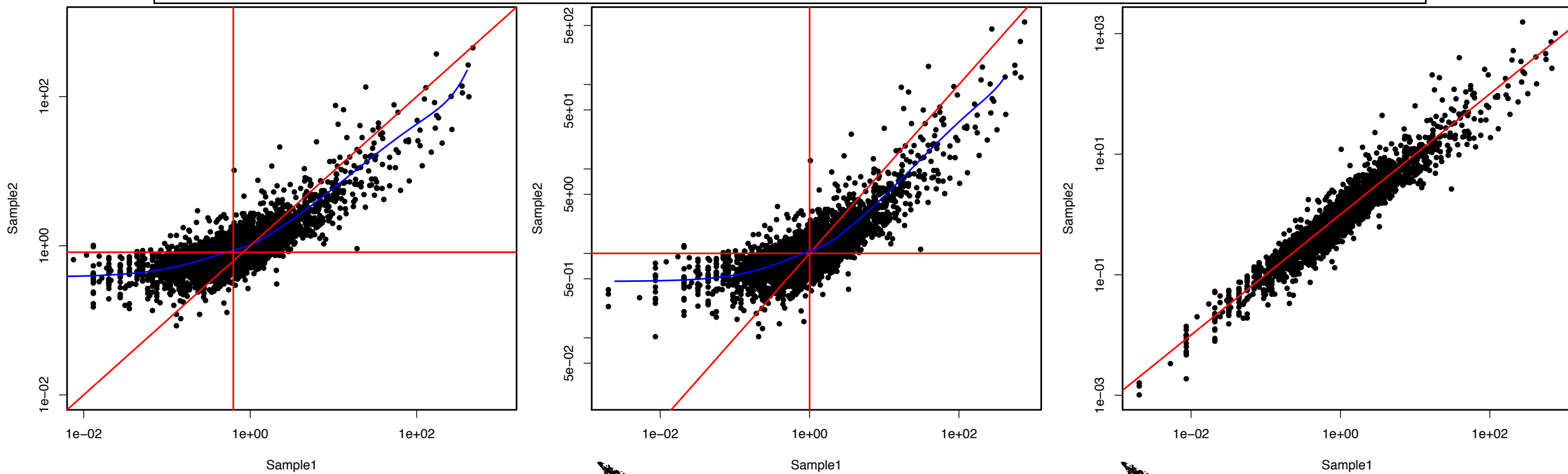
-2 -1 0 1 2

# Normalization

**Average variation between any 2 conditions is 0: Systematic variation MUST BE technical artifact**



- Normalization consists in removing systematic variations
- Work in log-log coordinates
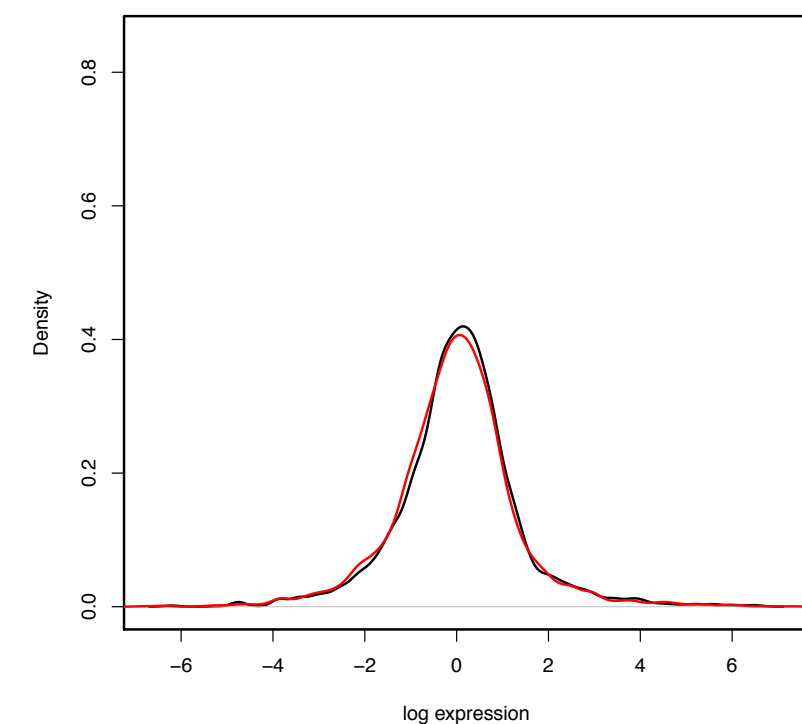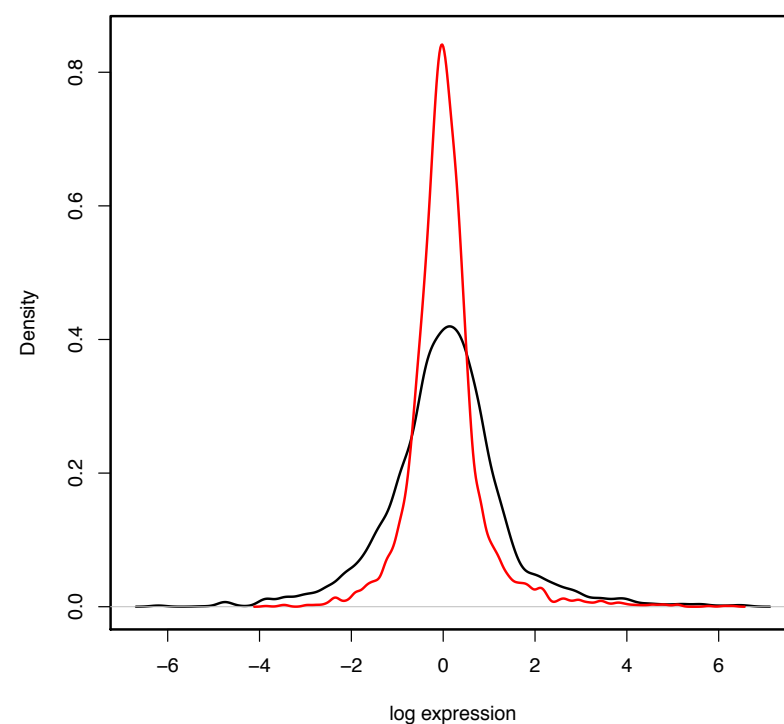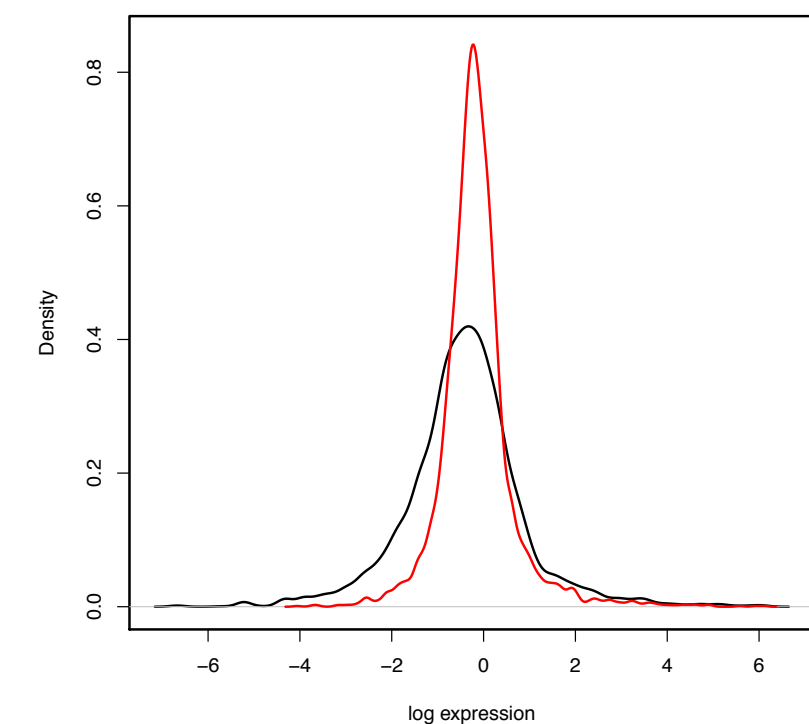- Shift in medians
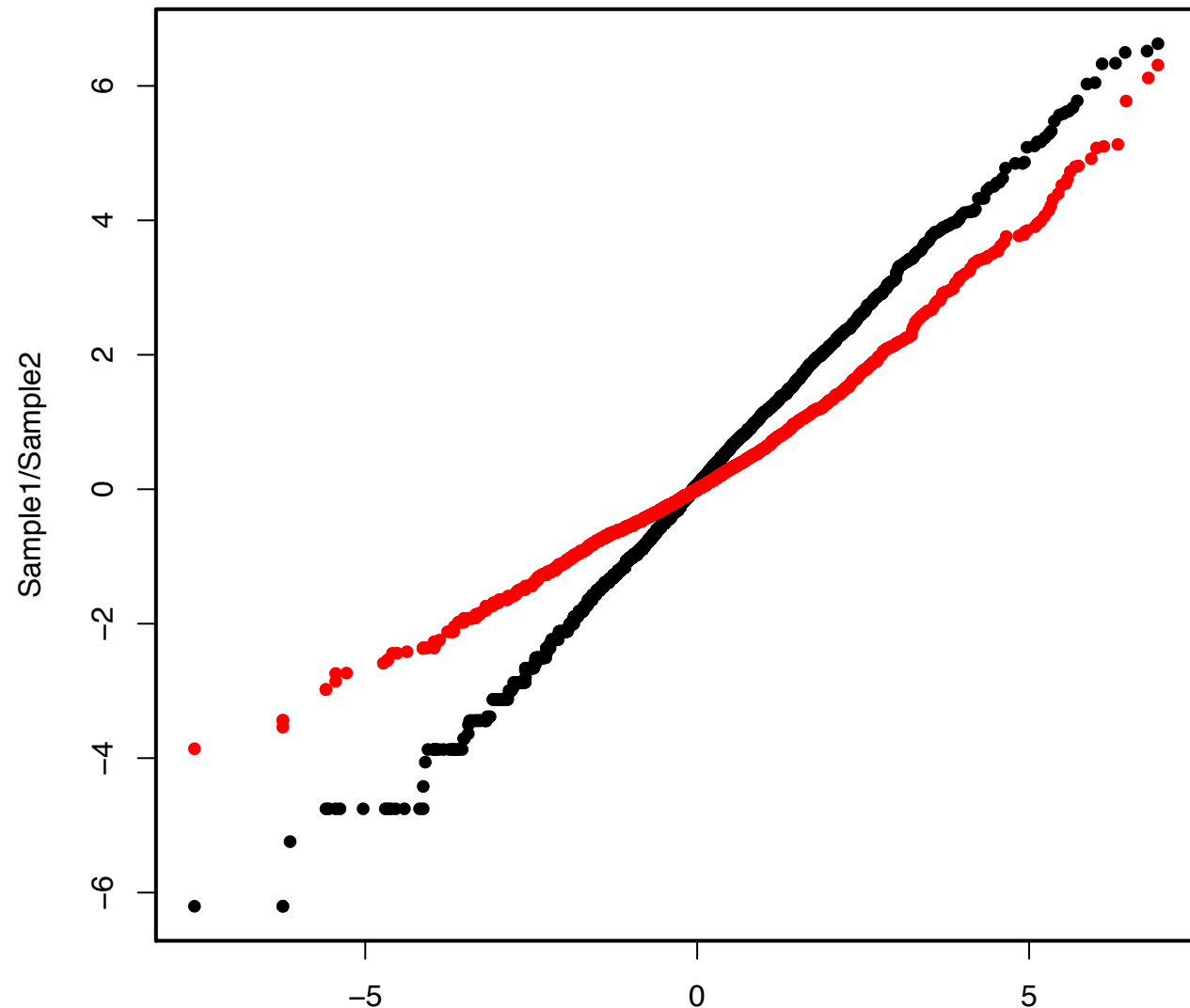- Curved shaped

# Normalization



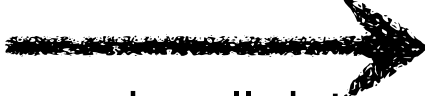sample=sample/median(sample)         sample2=sample2/fit*sample1

# Quantile normalization
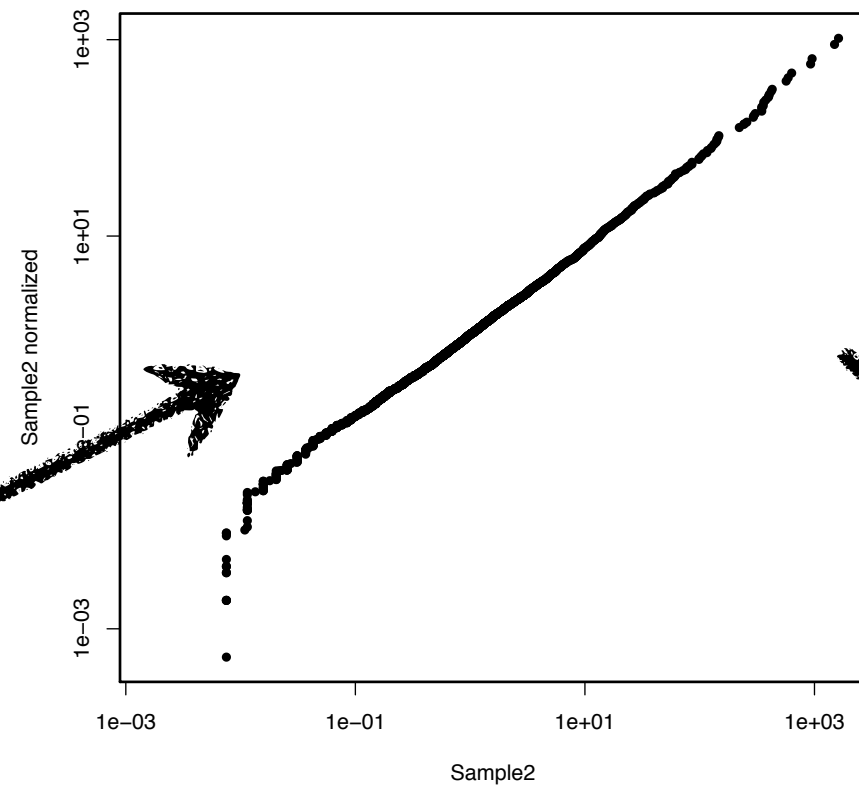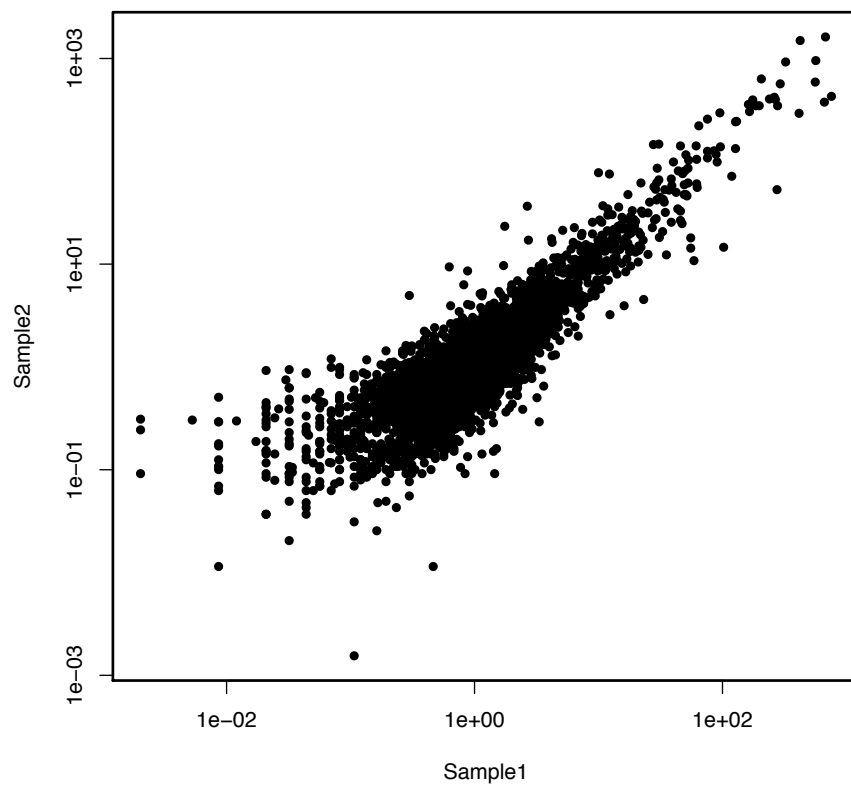


- Substitute ordered values from every sample with ordered values from average (or from specific distribution, e.g. gaussian)

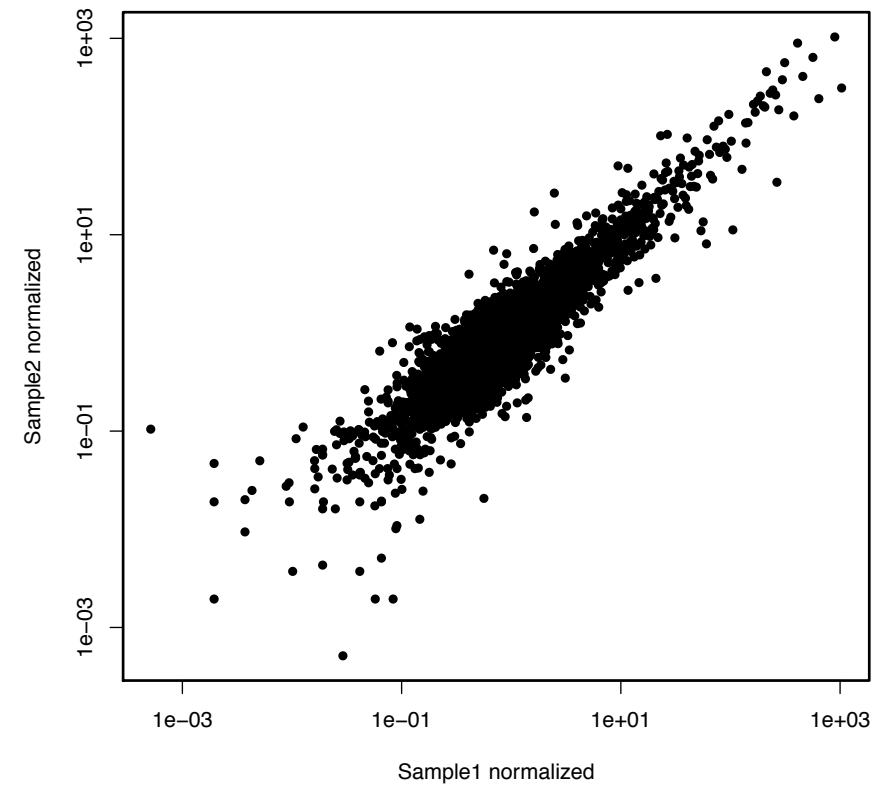# Quantile normalization



before normalization

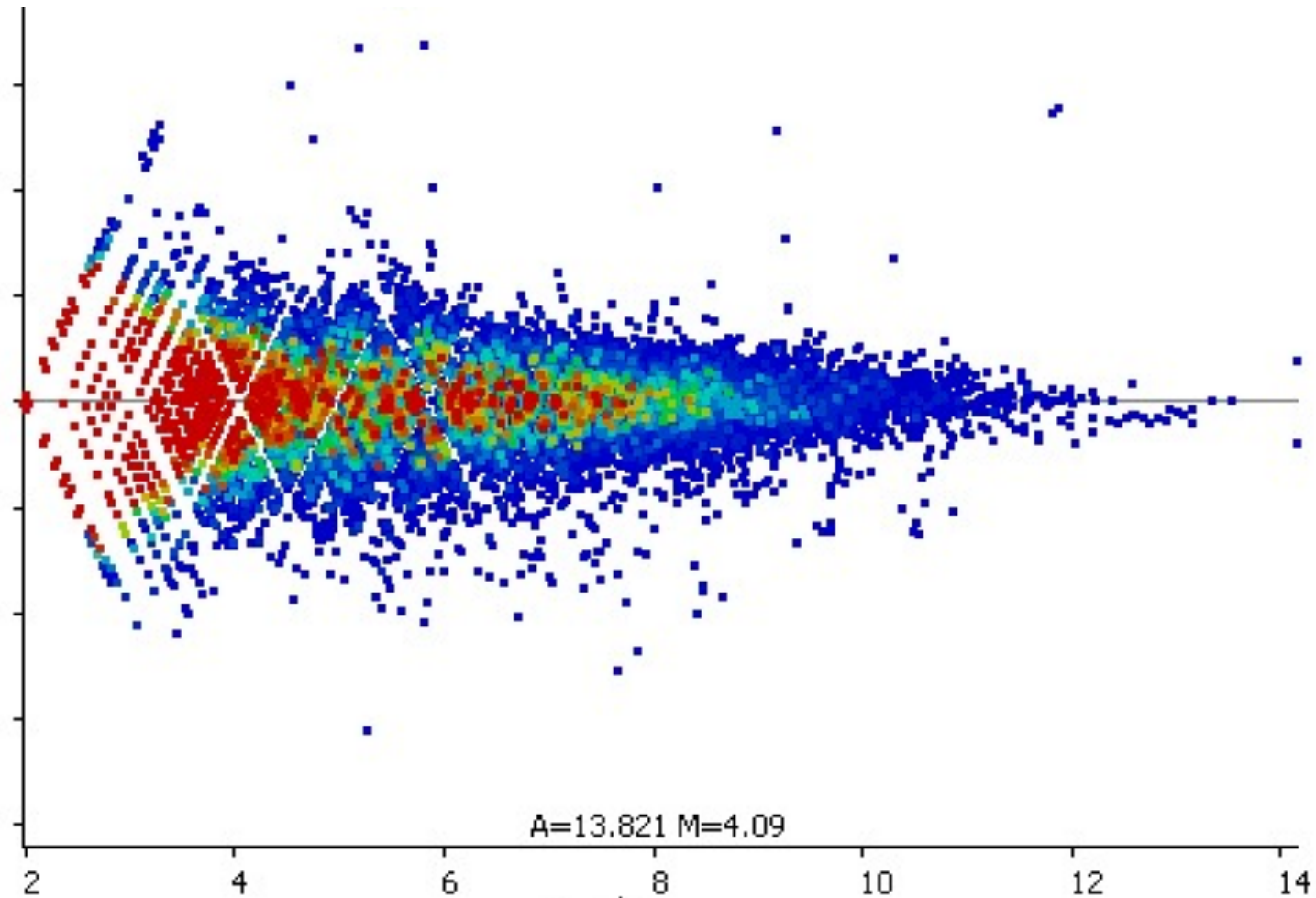after normalization

# MA plot: differential expression

Log2(Sample1/Sample2)



A=13.821 M=4.09

2    4    6    8    10    12    14

Log10(Sample1*Sample2)/2

differential expression

absolute expression

# Clustering



- Same algorithm as UPGMA

- Distance matrix is $1-cor(gene_i, gene_j)$

- Update matrix with distance to average of two groups weighted by size

- Do the same for columns (rotate matrix)



## distance

|   | 1 | 2 |
|---|---|---|
| 2 | 4.3 |   |
| 3 | 2.4 | 4.9 |

## correlation

|   | 1 | 2 |
|---|---|---|
| 2 | 0.99 |   |
| 3 | -0.06 | -0.03 |

# Clustering

# Clustering



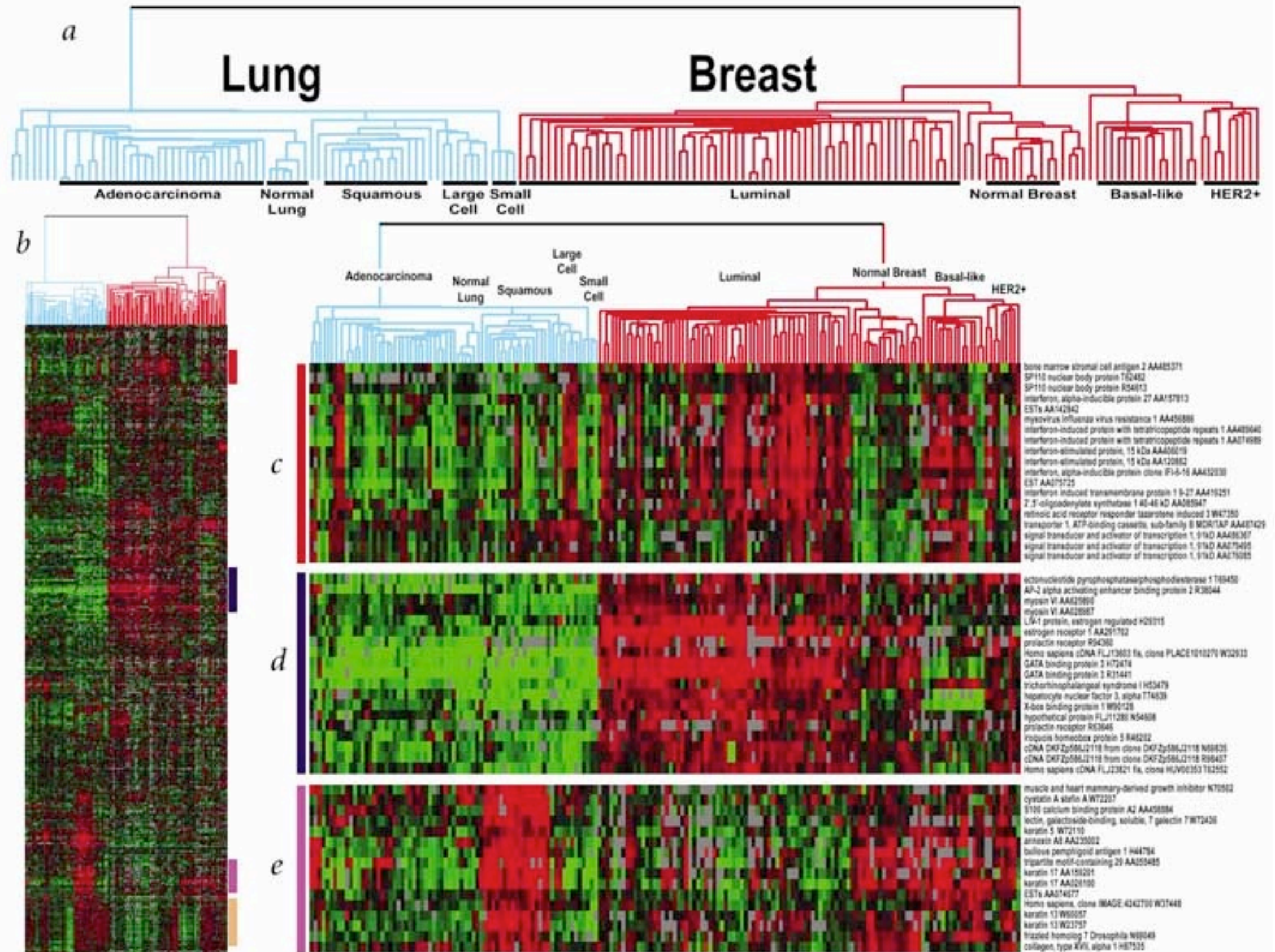- Similar expression patterns across many conditions probably imply a common set of regulators

- Looking for a shared set of functional annotations can help find the regulators

# Linear models

| | treat+WT | treat+KO | no treat+WT | no treat+KO |
|---|---|---|---|---|
| g | $M_{11}$ | $M_{12}$ | $M_{13}$ | $M_{14}$ |

For each gene, make a linear relation between
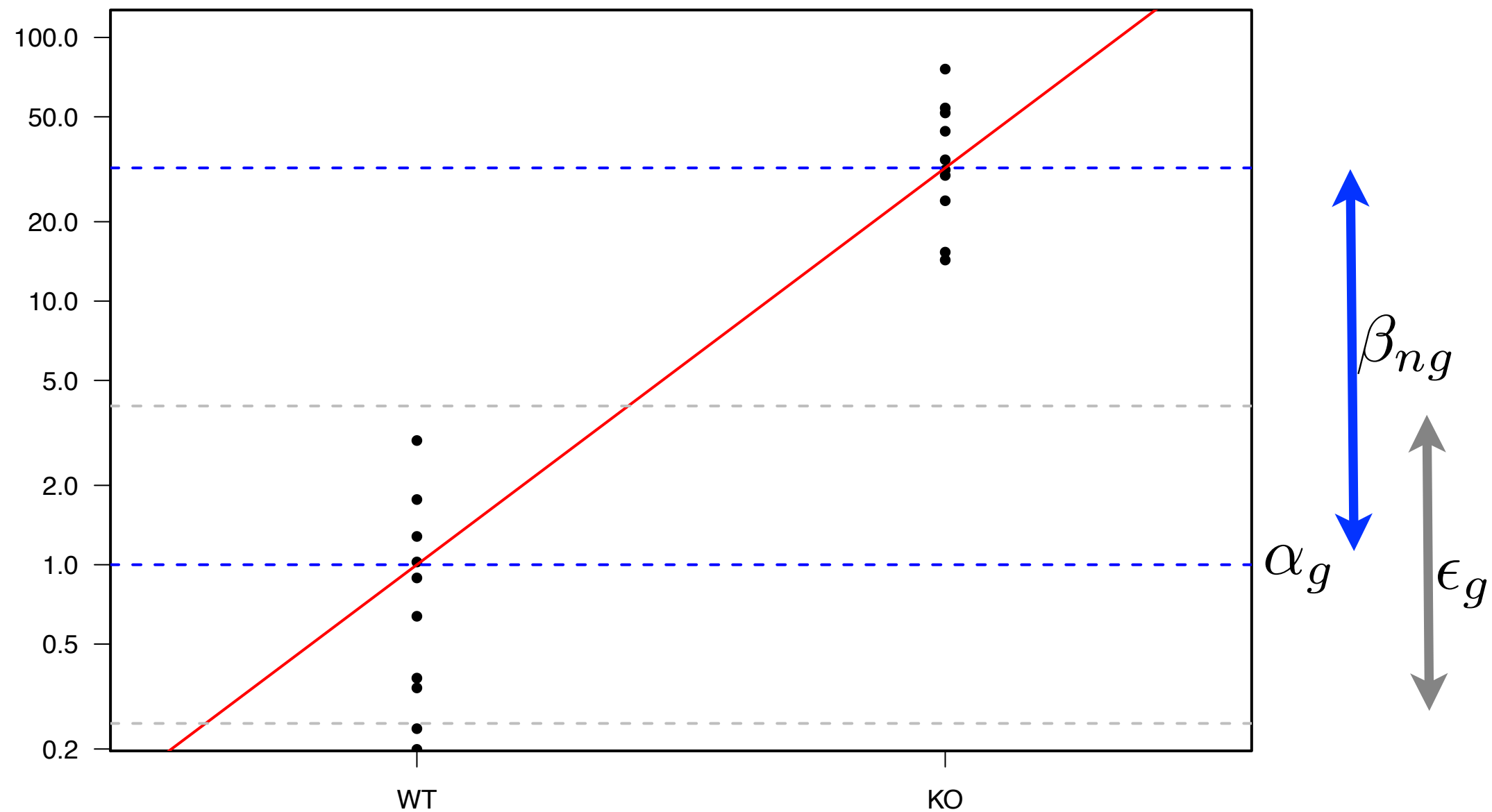effect (expression) and factors (conditions)

$$\log(M_{cg}) \quad = \quad \alpha_g + \sum_n I_{cn}\beta_{ng} + \epsilon_g \ ,$$

Design matrix:

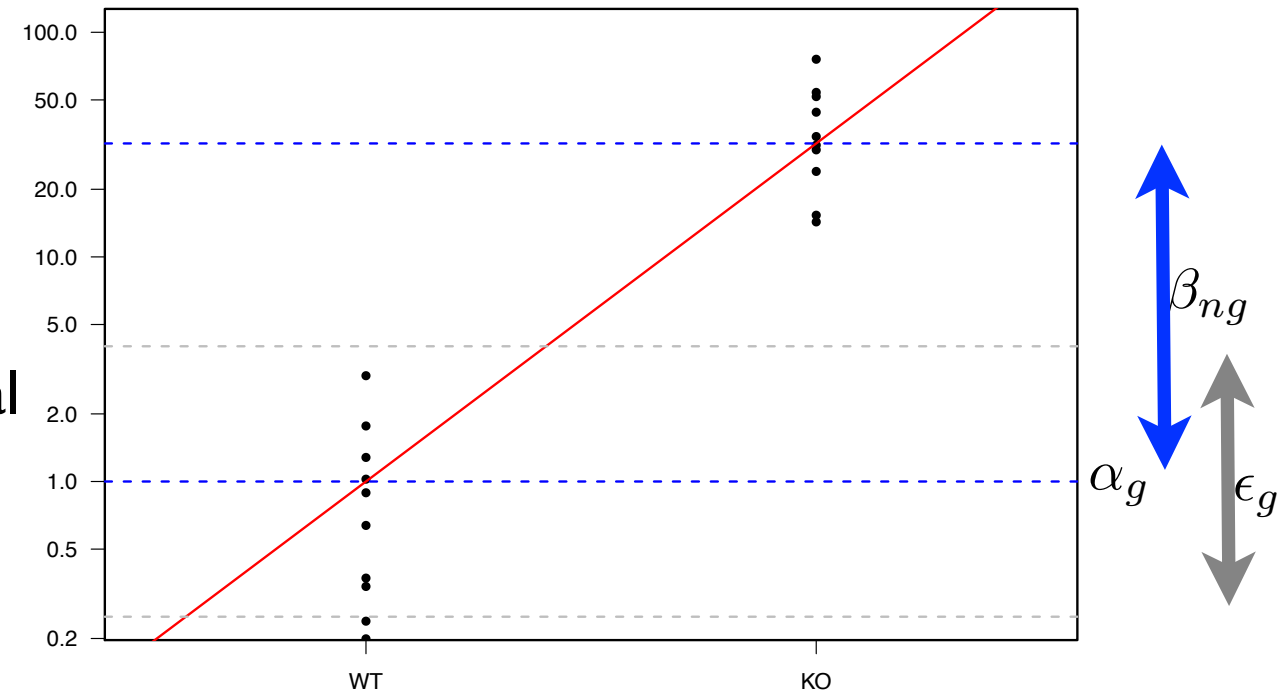| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| treat | 1 | 1 | 0 | 0 |
| KO | 0 | 1 | 0 | 1 |

# Linear models

$$\log(M_{cg}) \; = \; \alpha_g + \sum_n I_{cn}\beta_{ng} + \epsilon_g$$

# Hypothesis Testing



- "Null Hypothesis" H0:

  - average expression of $g$ in WT and KO are equal

- Compute a "statistic":

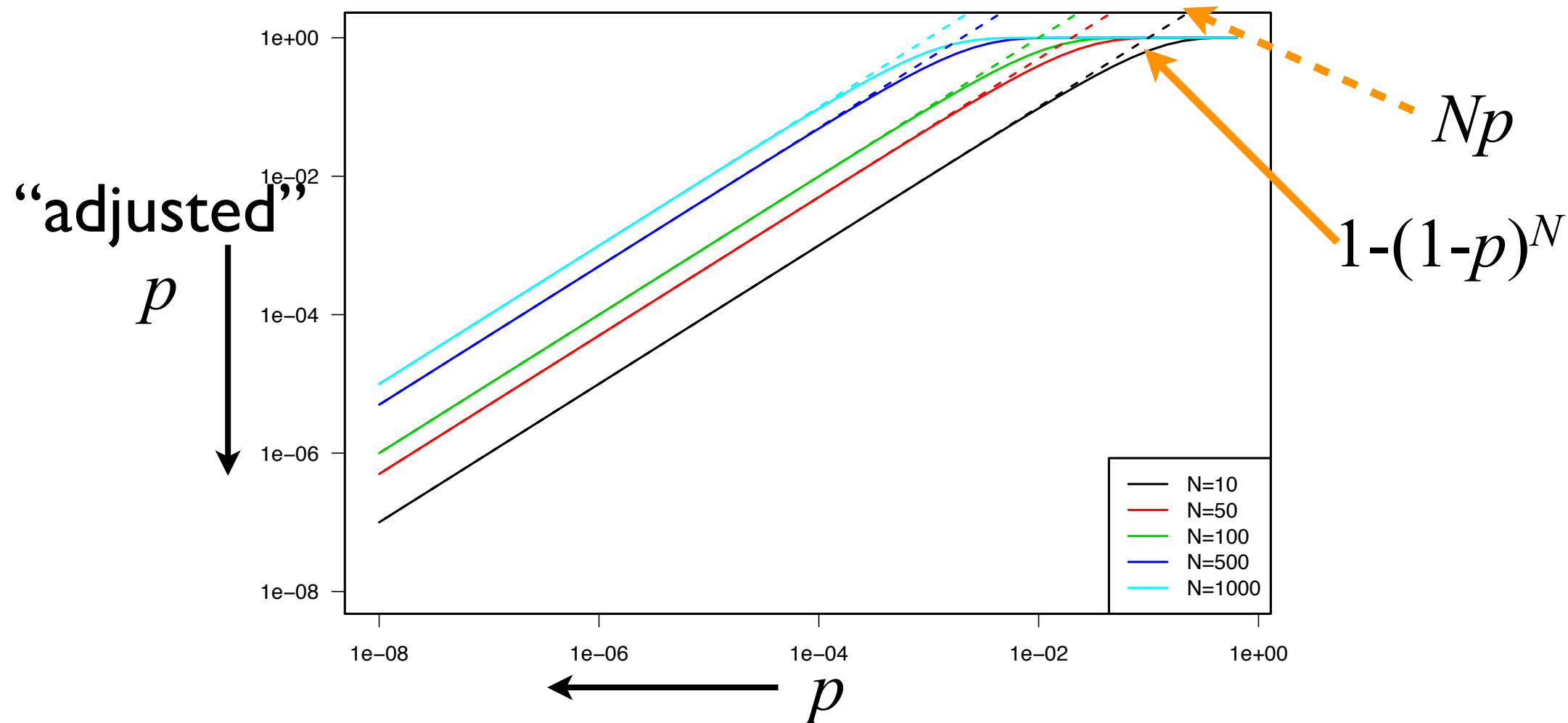$$t(g) = \frac{\beta^2}{\epsilon_{\mathrm{WT}}^2 + \epsilon_{\mathrm{KO}}^2}$$

- Model the data (normal distribution), compute a "p-value":

  - probability of observing a variation larger than $t(g)$ if H0 <u>AND</u> the model are true

$$p(g) = \mathrm{Prob}(t \geq t(g)|H0)$$

# Multiple testing

- 25'000 genes $\Rightarrow$ 25'000 p-values. Suppose all tests are independent:

$$\mathrm{Prob}(\exists n \,:\, t(g_n) \geq t | H0) \;=\; 1 - \mathrm{Prob}(\forall n \,:\, t(g_n) < t | H0)$$

$$=\; 1 - \left(1 - \mathrm{Prob}(t(g) \geq t | H0)\right)^N$$



"adjusted" $p$

$Np$

$1-(1-p)^N$

$p$

# False Discovery Rate

- If we detect 200 differential genes out of 25'000, we may accept a small proportion of false positive

- False Discovery Rate (FDR): proportion of True H0 among set of 200

- Benjamini–Hochberg procedure: find largest $k$ such that $P(g_k)/k < FDR/N$