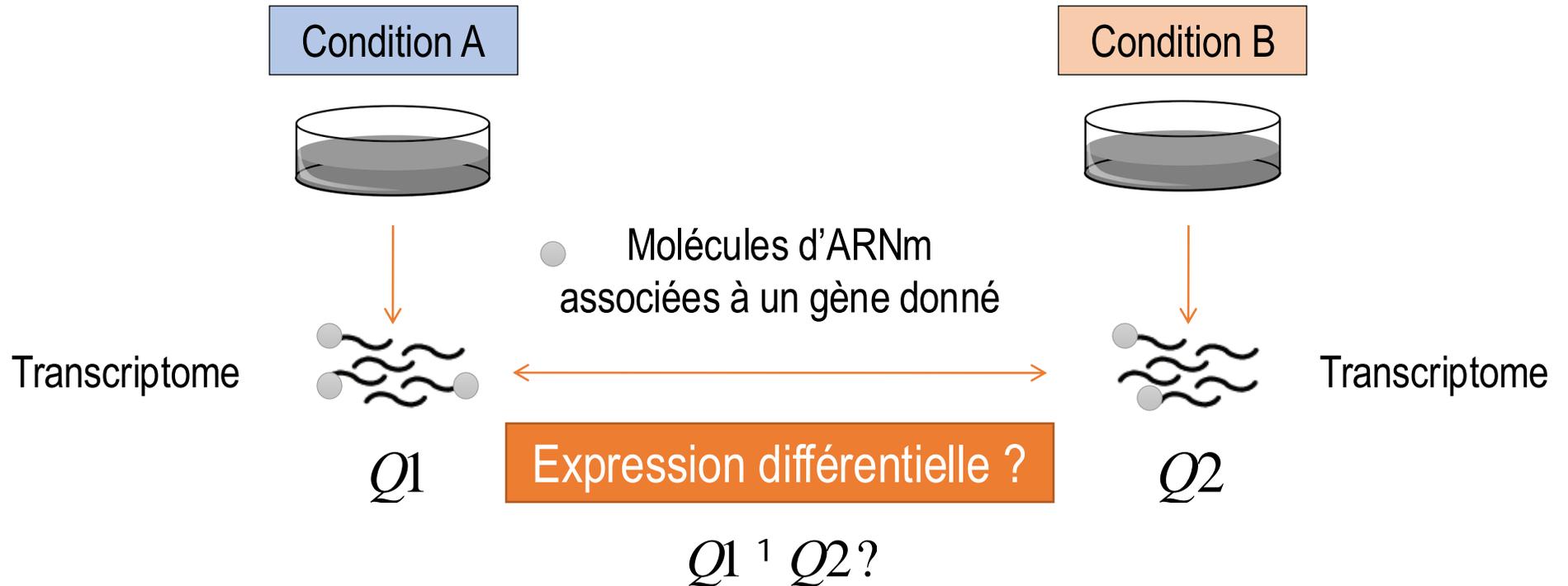




# L'analyse différentielle, taille d'effet, dispersion et force d'expression

Gaëlle LELANDAIS

# Quelle est la question ? (analyse différentielle)

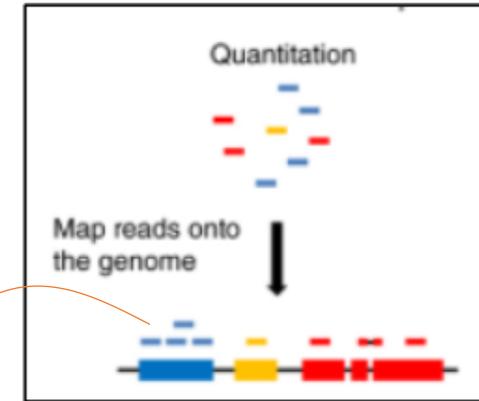


Gène différentiellement exprimé

- Tous les gènes sont étudiés successivement (répétition de la question plusieurs dizaines de milliers de fois)

# Les données analysées

➤ Table de comptages (*counts*) :



	Condition A			Condition B		
	R1	R2	R3	R1	R2	R3
A1BG	4	6	2	7	6	7
A1CF	41	33	42	32	42	32
A2M	1	3	1	4	3	7
A2ML1	3	2	2	6	7	3
A2MP1	3	2	2	1	1	0
A3GALT2	1	4	4	3	2	1
A4GALT	420	344	291	327	360	371
A4GNT	1	1	2	1	3	3
AA06	0	0	0	0	0	0
AAAS	2452	2192	1977	2054	2134	2100
AACS	3234	2804	2609	1678	1670	1742
AACSP1	1544	1369	1300	1926	2015	1963

# Problématique de normalisation (1/2)

- Le nombre de sequences (*reads*) depend de la longueur des transcrits ...

A1BG	4
A1CF	41
A2M	1
A2ML1	3
A2MP1	3
A3GALT2	1
A4GALT	420
A4GNT	1
AA06	0
AAAS	2452
AACS	3234
AACSP1	1544



8 *reads*  
associés



4 *reads*  
associés

# Problématique de normalisation (2/2)

- Le nombre de sequences (*reads*) dépend de la somme totale\* (par colonne) ...

A1BG	4	7
A1CF	41	32
A2M	1	4
A2ML1	3	6
A2MP1	3	1
A3GALT2	1	3
A4GALT	420	327
A4GNT	1	1
AA06	0	0
AAAS	2452	2054
AACS	3234	1678
AACSP1	1544	1926

Somme 1

Somme 2

Biais systématique ?

(si par exemple Somme 1 = 40 000 et  
Somme 2 = 30 000 )

Table de comptage normalisée

\* La somme des comptages (*counts*) par expérience est souvent nommé *library size* dans les articles.

# Quantifier l'expression différentielle

## ➤ Ratio

$$R = \frac{Q_A}{Q_B} = \frac{\text{Valeur}(s) \text{ de comptage Condition A}}{\text{Valeur}(s) \text{ de comptage Condition B}}$$

## ➤ Fold Change

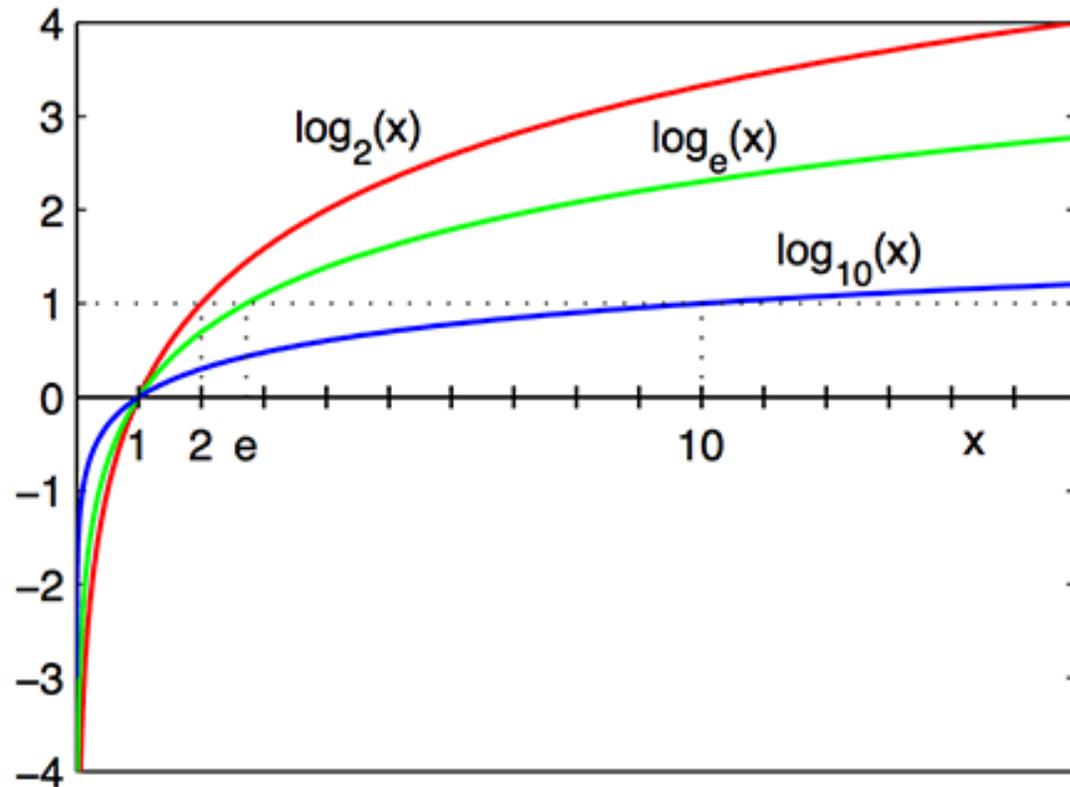
$$FC = \begin{cases} R & \text{si } R \geq 1 \\ \frac{-1}{R} & \text{si } R < 1 \end{cases}$$

## ➤ Log fold change

$$\log FC = \log_2(R) = \log_2(Q_A) - \log_2(Q_B)$$

# Interprétation du logFC (taille d'effet)

- Les valeurs du logFC peuvent être positives et négatives. L'utilisation de la base 2 permet de traduire un doublement par une unité de variation (+/-).



<https://fr.wikipedia.org/wiki/Logarithme>

$$\log FC_g > 0 \hat{=} Q1_g > Q2_g$$

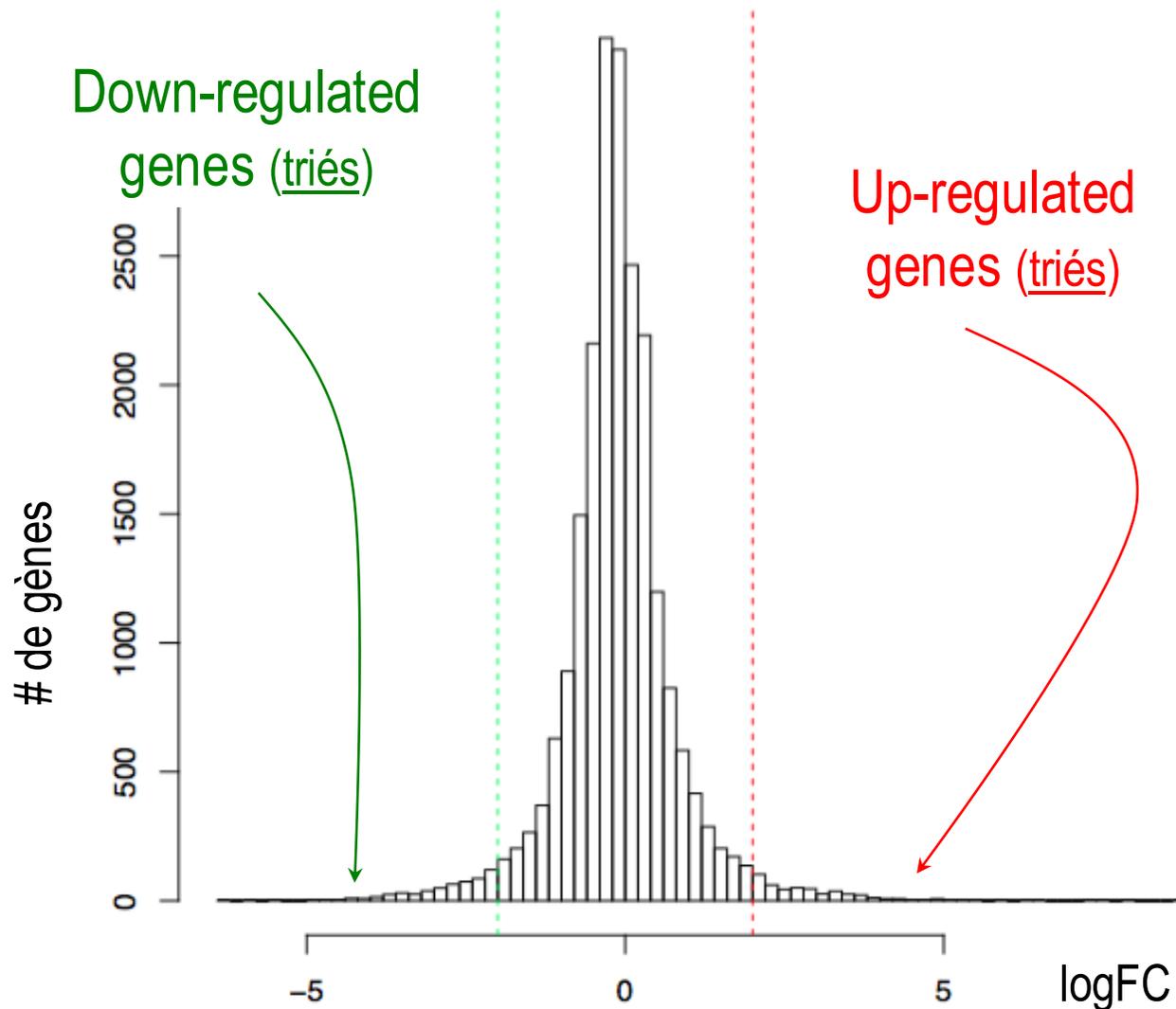
↳ "Up-regulated gene"

$$\log FC_g < 0 \hat{=} Q1_g < Q2_g$$

↳ "Down-regulated gene"

Notion de taille d'effet

# Utiliser la taille d'effet



Source des données : Yang et al., Mol Cell Biol (2016)

Liste de gènes, candidats pour être différentiellement exprimés

# Tenir compte de la reproductibilité

(dispersion)

- Dans la pratique, les expériences sont répétées. La taille d'effet mesurée est fondée sur plusieurs observations.

$$\text{Gène 1: } \log FC_{G1} = \frac{1}{3}(0.5 + 2.5 + 6) = 3$$

$$++++ \text{ Gène 2: } \log FC_{G2} = \frac{1}{3}(3 + 2.8 + 3.2) = 3$$

Meilleure cohérence entre les observations

Notion de dispersion

# Quantifier la dispersion\*

➤ Variance (estimateur)

$$Var(X) = \frac{1}{n-1} \sum (x_i - m)^2$$

➤ Ecart type

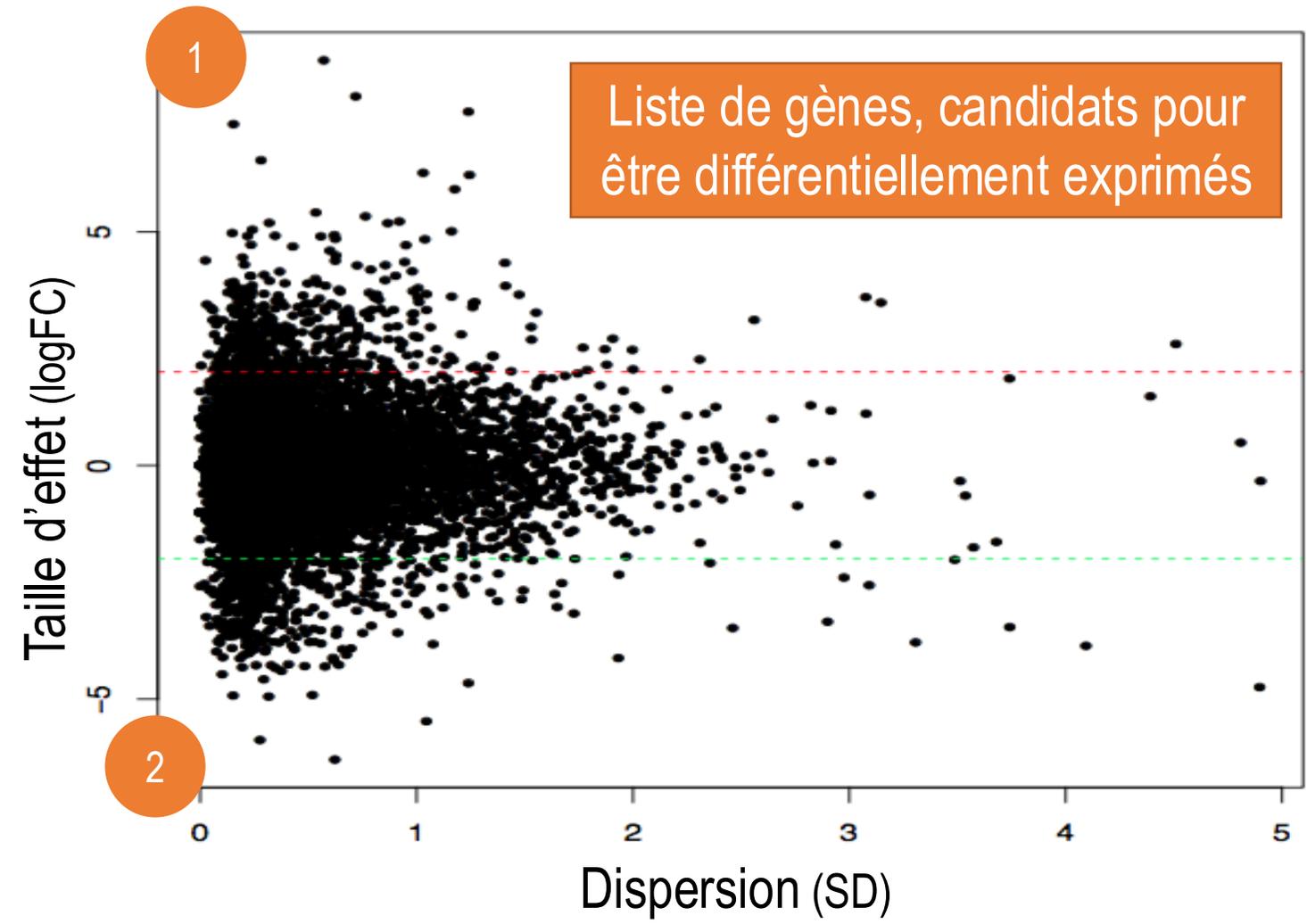
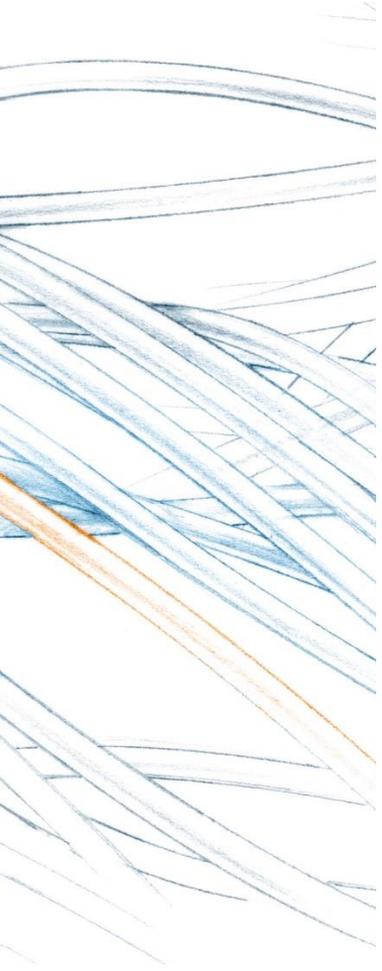
$$SD = \sqrt{Var(X)}$$

➤ Exemple :

	logFC	Variance	SD
Gène 1	3	7.75	2.78
→ Gène 2	3	0.04	0.20

\* Mais aussi le coefficient de variation biologique, l'erreur standard, etc.

# Combiner taille d'effet et dispersion

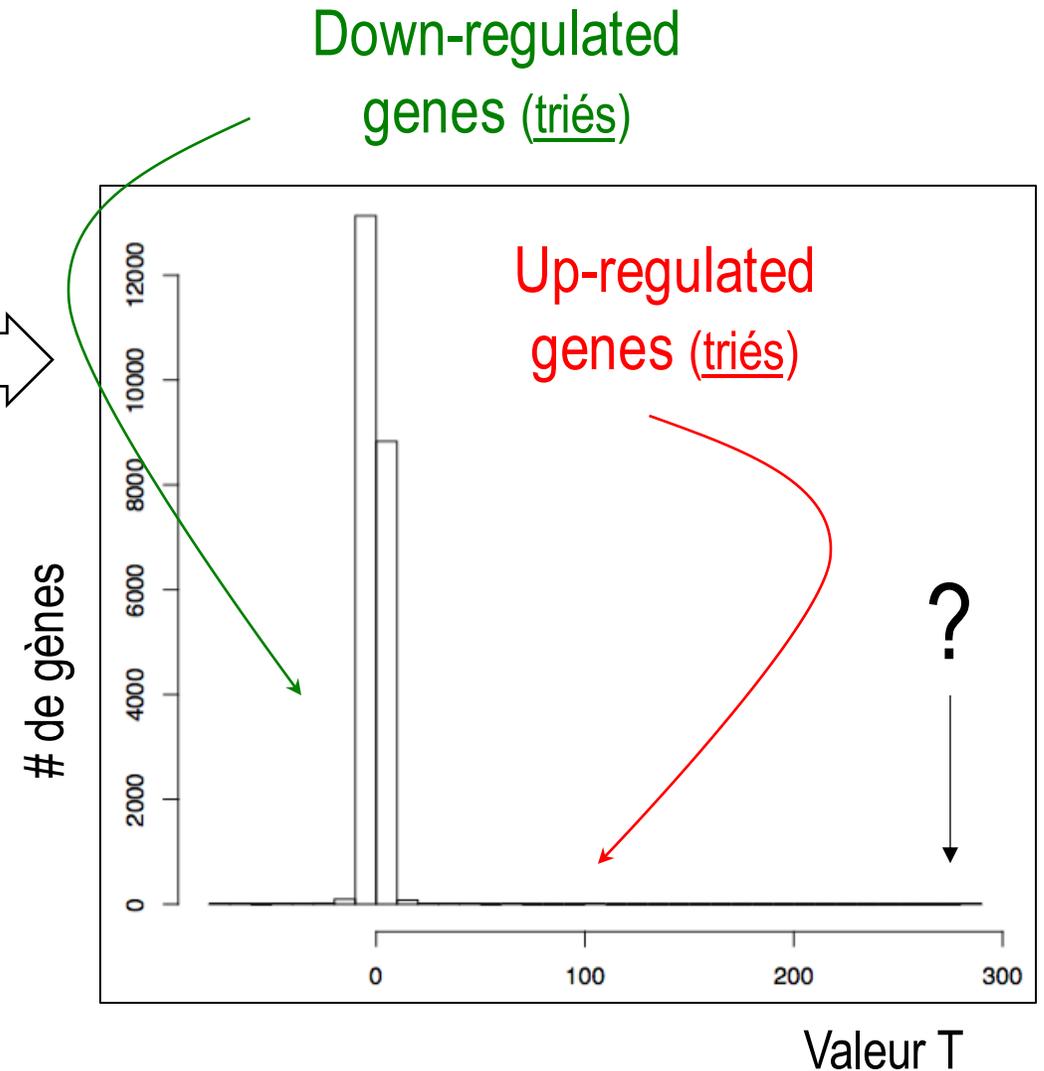
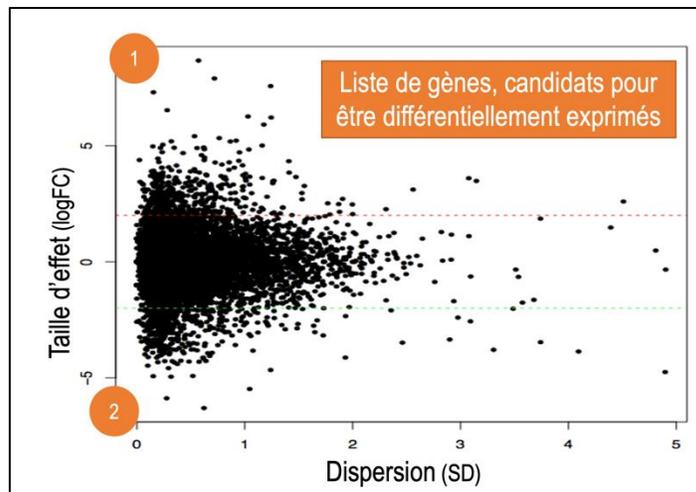


Source des données : Yang et al., Mol Cell Biol (2016)

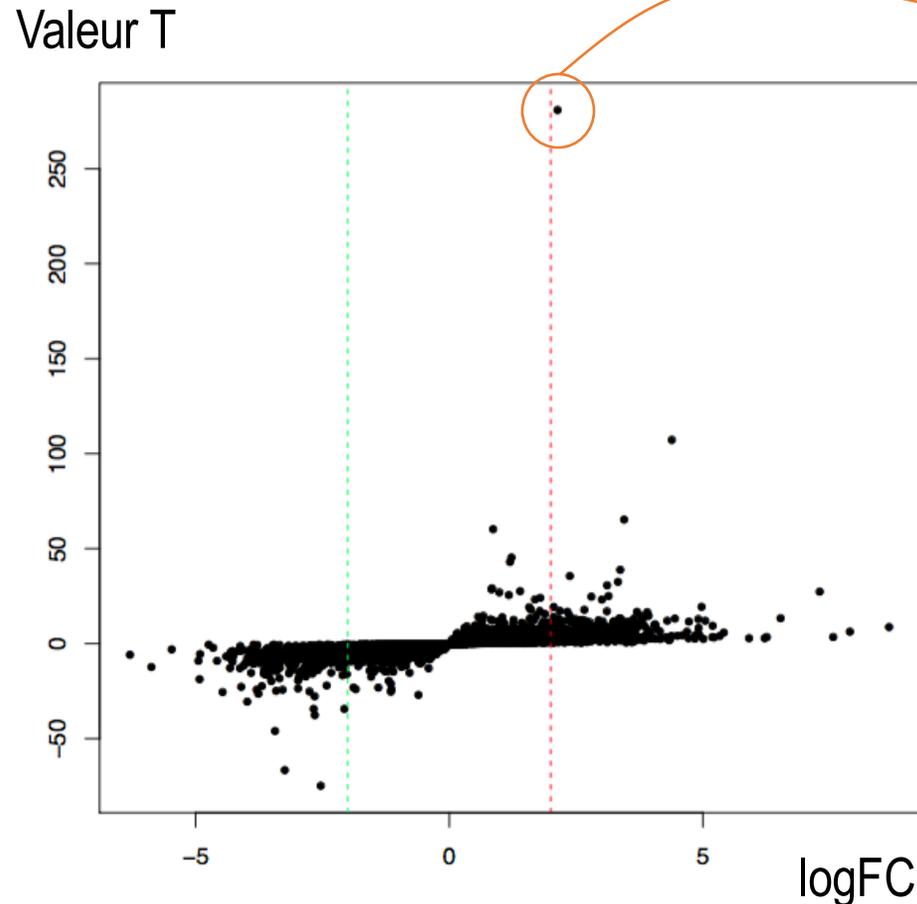
# Statistique de Student

- Taille d'effet, dispersion, nombre d'observations :

$$t_g = \frac{\log FC_g}{\frac{SD_g}{\sqrt{n}}}$$



# ... malheureusement non !



Importance trop grande donnée au paramètre de dispersion ( $n = 3$ )

- Estimer correctement la variabilité associée aux observations de comptage d'un gène est un **défi statistique**.

A1BG	4	6	2	7	6	7
A1CF	41	33	42	32	42	32
A2M	1	3	1	4	3	7
A2ML1	3	2	2	6	7	3
A2MP1	3	2	2	1	1	0
A3GALT2	1	4	4	3	2	1
A4GALT	420	344	291	327	360	371
A4GNT	1	1	2	1	3	3
AA06	0	0	0	0	0	0
AAAS	2452	2192	1977	2054	2134	2100
AACS	3234	2804	2609	1678	1670	1742
AACSP1	1544	1369	1300	1926	2015	1963

# Tenir compte du nombre de séquences alignées (force d'expression)

➤ Une même taille d'effet peut être obtenues avec des données de comptages différentes.

Gene 1:

$$\log FC_{G_1} = \log 2(8000 / 1000) = 3$$

Gene 2:

$$\log FC_{G_2} = \log 2(8 / 1) = 3$$

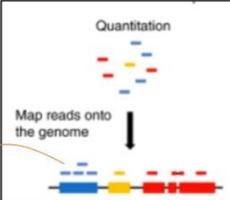
Variations aléatoires ?

La notion de “force d'expression” est intéressante à prendre en compte pour établir une liste de gène candidats, avec une meilleure fiabilité.

\* Cette information est nommée *expression strength* dans les articles statistiques.

## Les données analysées

➤ Table de comptages (*counts*) :



	Condition A			Condition B		
	R1	R2	R3	R1	R2	R3
A1BG	4	6	2	7	6	7
A1CF	41	33	42	32	42	32
A2M	1	3	1	4	3	7
A2ML1	3	2	2	6	7	3
A2MP1	3	2	2	1	1	0
A3GALT2	1	4	4	3	2	1
A4GALT	420	344	291	327	360	371
A4GNT	1	1	2	1	3	3
AA06	0	0	0	0	0	0
AAAS	2452	2192	1977	2054	2134	2100
AACS	3234	2804	2609	1678	1670	1742
AACSP1	1544	1369	1300	1926	2015	1963

# Bilan

- L'analyse différentielle est un problème difficile qui nécessite l'utilisation de méthodologies d'analyses adaptées.

DESeq et DESeq2 →

edgeR →

LIMMA →

**TABLE 1.** RNA-seq differential gene expression tools and statistical tests

Name	Assumed distribution	Normalization	Description	Version	Citations <sup>d</sup>	Reference
<i>t</i> -test	Normal	DESeq <sup>a</sup>	Two-sample <i>t</i> -test for equal variances	–	–	–
log <i>t</i> -test	Log-normal	DESeq <sup>a</sup>	Log-ratio <i>t</i> -test	–	–	–
Mann-Whitney	None	DESeq <sup>a</sup>	Mann-Whitney test	–	–	Mann and Whitney (1947)
Permutation	None	DESeq <sup>a</sup>	Permutation test	–	–	Efron and Tibshirani (1993a)
<i>Bootstrap</i>	Normal	DESeq <sup>a</sup>	Bootstrap test	–	–	Efron and Tibshirani (1993a)
<i>baySeq</i> <sup>c</sup>	Negative binomial	Internal	Empirical Bayesian estimate of posterior likelihood	2.2.0	159	Hardcastle and Kelly (2010)
<i>Cuffdiff</i>	Negative binomial	Internal	Unknown	2.1.1	918	Trapnell et al. (2012)
<i>DEGseq</i> <sup>c</sup>	Binomial	None	Random sampling model using Fisher's exact test and the likelihood ratio test	1.22.0	325	Wang et al. (2010)
<i>DESeq</i> <sup>c</sup>	Negative binomial	DESeq <sup>a</sup>	Shrinkage variance	1.20.0	1889	Anders and Huber (2010)
<i>DESeq2</i> <sup>c</sup>	Negative binomial	DESeq <sup>a</sup>	Shrinkage variance with variance based and Cook's distance pre-filtering	1.8.2	197	Love et al. (2014)
<i>EBSeq</i> <sup>c</sup>	Negative binomial	DESeq <sup>a</sup> (median)	Empirical Bayesian estimate of posterior likelihood	1.8.0	80	Leng et al. (2013)
<i>edgeR</i> <sup>c</sup>	Negative binomial	TMM <sup>b</sup>	Empirical Bayes estimation and either an exact test analogous to Fisher's exact test but adapted to over-dispersed data or a generalized linear model	3.10.5	1483	Robinson et al. (2010)
<i>Limma</i> <sup>c</sup>	Log-normal	TMM <sup>b</sup>	Generalized linear model	3.24.15	97	Law et al. (2014)
<i>NOISeq</i> <sup>c</sup>	None	RPKM	Nonparametric test based on signal-to-noise ratio	2.14.0	177	Tarazona et al. (2011)
<i>PoissonSeq</i> <sup>c</sup>	Poisson log-linear model	Internal	Score statistic	1.1.2	37	Li et al. (2012)
<i>SAMSeq</i> <sup>c</sup>	None	Internal	Mann-Whitney test with Poisson resampling	2.0	54	Li and Tibshirani (2013)

<sup>a</sup>See Anders and Huber (2010).  
<sup>b</sup>See Robinson and Oshlack (2010).  
<sup>c</sup>R (v3.2.2) and bioconductor (v3.1).  
<sup>d</sup>As reported by PubMed Central articles that reference the listed reference (December 21, 2015).

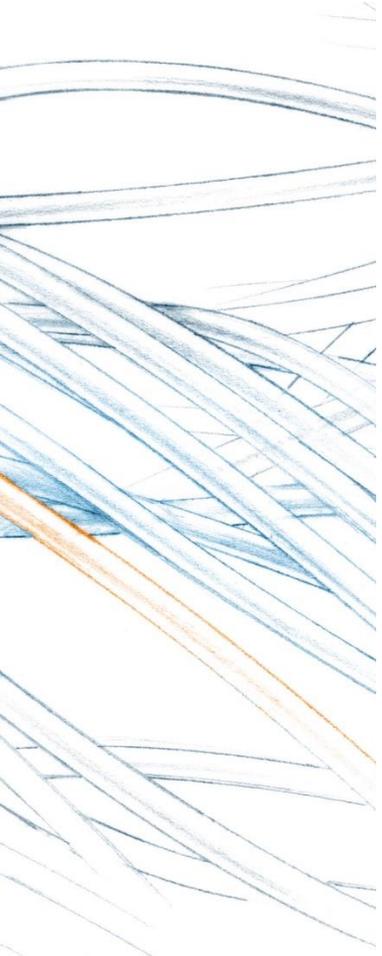
Taille d'effet

Dispersion

Force expression

# Ressource(s) complémentaire(s)

- Commentaires enregistrés du cours (16 minutes) :
  - Lien Youtube : <https://youtu.be/eztUy86l4ds>





Sauf mention contraire, ce contenu est mis à disposition selon les termes de la licence Creative Commons Attribution - Partage dans les mêmes conditions 4.0 International (CC BY-SA 4.0)

Gaëlle LELANDAIS

Version du document : 11/02/2025