

# STATISTICAL METHODS IN BIOINFORMATICS

## Analysis of RNA sequencing data

Stefan E Seemann

`ses@sund.ku.dk`

Center for non-coding RNA in Technology and Health (<http://rth.dk>)

Dep. of Veterinary and Animal Sciences, SUND,

University of Copenhagen DK

29 April 2025

## Why sequencing?

- Assemble the genome and transcriptome of a species
- Find genomic variation in a population
- Find genomic and transcriptomic associations with diseases and phenotypes
- Find organisms in environmental sample → metagenomics and -transcriptomics
- Identify potential drug targets → personalized medicine
- Tracking of virus variants and mutations → vaccine development

⇒ One of the biggest hammers in the tool box right now – always ask could this experiment be done using sequencing instead?

What is your application?

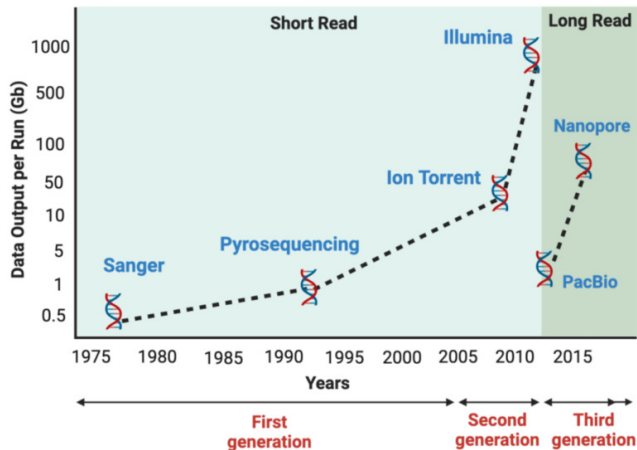
## The beginning

- 1968 — The first 12 bases
- 1973 — 24 bases of the lactose-repressor binding site  
→ two years of work: one base per month
- 1977 — Sanger sequencing and Gilbert sequencing  
→ The Nobel Prize in Chemistry 1980 for Frederick Sanger and Walter Gilbert





# Evolution of sequencing technologies



# Today's programme

## **Alignment and mapping**

0815-0900	<i>Lecture</i>	Alignment methods
0900-1000	<i>Exercise</i>	Dynamical programming of pairwise alignment ( <i>on paper</i> )
1000-1030	<i>Lecture</i>	Read mapping

## **Feature abundance and differential expression**

1030-1100	<i>Lecture</i>	Normalization + Transformation
11-12	<i>Lunch</i>	
1200-1300	<i>Lecture</i>	Unsupervised + Supervised data exploration
1300-1445	<i>Exercise</i>	Differential expression analysis ( <i>in R</i> )
1445-1500		Summary and Discussion

# Alignment methods

Where do we need sequence alignments?

## Where do we need sequence alignments?

- Sequence similarity
- Gene finding by similarity
- Protein structure by similarity
- RNA structure by similarity
- Motif finder
- Genome and transcriptome assembly
- **Gene expression estimation**

# Evolutionary events

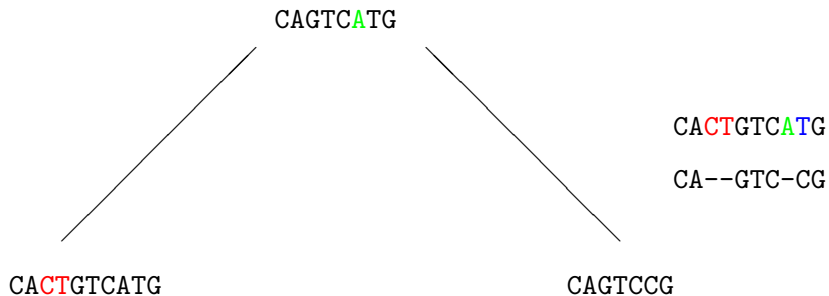
- DNA sequences change in time.
- Find evolutionary related sequences.
- Evolutionary events:

CAGTCATG  $\xrightarrow{\text{INSERTION}}$  CACTGTCATG  $\xrightarrow{\text{DELETION}}$  CACTGTCTG  
 $\xrightarrow{\text{SUBSTITUTION}}$  CACTATCTG

CAGTCATG  $\xrightarrow{\text{DUPLICATION}}$  CAGTGTCATG  $\xrightarrow{\text{TRANSLOCATION}}$  CATCAGTGTG  $\xrightarrow{\text{INVERSION}}$   
CATTGACGTG

# Evolutionary tree

- Finding common ancestors.
- Parsimony principle: Evolution uses minimum number of operations.



- Probabilistic approaches (max. likelihood or sampling).

## Alignments: optimize a score

Score of a given alignment:

C	A	C	T	G	T	C	A	T	G
C	A	-	-	G	T	C	-	C	G

$$S_{\text{tot}} = S \begin{bmatrix} \text{C} \\ \text{C} \end{bmatrix} + S \begin{bmatrix} \text{A} \\ \text{A} \end{bmatrix} + S \begin{bmatrix} \text{C} \\ - \end{bmatrix} + S \begin{bmatrix} \text{T} \\ - \end{bmatrix} + S \begin{bmatrix} \text{G} \\ \text{G} \end{bmatrix} + S \begin{bmatrix} \text{T} \\ \text{T} \end{bmatrix} + S \begin{bmatrix} \text{C} \\ \text{C} \end{bmatrix} + S \begin{bmatrix} \text{A} \\ - \end{bmatrix} + S \begin{bmatrix} \text{T} \\ \text{C} \end{bmatrix} + S \begin{bmatrix} \text{G} \\ \text{G} \end{bmatrix}$$

- Score: substituting a residue in one seq. with a residue in another.
- Find the alignment that have the highest score.
- Try out all alignment combinations? (we deal with this soon)
- So speed of sequence comparisons **matters!**



## Score matrices for DNA

- Identity: 8
- Transition: 2 (eg.  $\{A,G\} \rightarrow \{A,G\}$ ; purine to purine)
- Transversion: -3 (eg.  $A \rightarrow \{C,T\}$ ; purine to pyrimidine).

	A	C	G	T
A				
C				
G				
T				

## What about gaps?

- Gap cost. Cost of indel (Eg.  $d = 10$ ).
- Initiation and elongation.

# Dynamical programming

- Find the alignment between CACTGTCATG and CAGTCTG that has the maximal score?
- What would be a trivial way?

# Dynamical programming

- Find the alignment between CACTGTCATG and CAGTCTG that has the maximal score?
- What would be a trivial way?
- Basic idea: Use sub sequences! → **Dynamic Programming**

## Pairwise global alignments (Needleman–Wunsch)

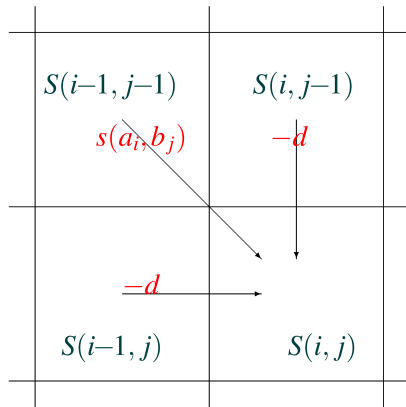
Comparing sequences  $a$  and  $b$ . Given a substitution score  $s(x, y)$  of replacing letter  $x$  with letter  $y$ , the highest scoring alignment can be found by the following recursion:

$$S(i, j) = \max \begin{cases} S(i-1, j-1) + s(a_i, b_j) \\ S(i-1, j) - d \\ S(i, j-1) - d \end{cases}$$

$a_i$  residue at position  $i$  in seq.  $a$

$b_j$  residue at position  $j$  in seq.  $b$

$i = 1, \dots, N; j = 1, \dots, M$



Initialization:  $S(0, 0) = 0$ . Hence:  $S(i, 0) = -id$ ,  $S(0, j) = -jd$ .

Note: the alignment takes time  $O(NM)$ .

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0										
C											
A											
G											
T											
C											
T											
G											

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10										
A	-20										
G	-30										
T	-40										
C	-50										
T	-60										
G	-70										

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2									
G	-30	-12									
T	-40	-22									
C	-50	-32									
T	-60	-42									
G	-70	-52									



## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	16	6	-4	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	4	-6	-16	-26	-36	-46
T	-40	-22	-4	8							
C	-50	-32	-14	4							
T	-60	-42	-24	-6							
G	-70	-52	-34	-16							

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	16	6	-4	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	4	-6	-16	-26	-36	-46
T	-40	-22	-4	8	21	11	12	2	-8	-18	-28
C	-50	-32	-14	4	11	18	13	20	10	0	-10
T	-60	-42	-24	-6	12	8	26	16	17	18	8
G	-70	-52	-34	-16	2	20	10	23	18	14	26

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	16	6	-4	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	4	-6	-16	-26	-36	-46
T	-40	-22	-4	8	21	11	12	2	-8	-18	-28
C	-50	-32	-14	4	11	18	13	20	10	0	-10
T	-60	-42	-24	-6	12	8	26	16	17	18	8
G	-70	-52	-34	-16	2	20	10	23	18	14	26

Back-tracking ...

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	16	6	-4	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	4	-6	-16	-26	-36	-46
T	-40	-22	-4	8	21	11	12	2	-8	-18	-28
C	-50	-32	-14	4	11	18	13	20	10	0	-10
T	-60	-42	-24	-6	12	8	26	16	17	<b>18</b>	8
G	-70	-52	-34	-16	2	20	10	23	18	14	<b>26</b>

Back-tracking ...

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	16	6	-4	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	4	-6	-16	-26	-36	-46
T	-40	-22	-4	8	21	11	12	2	-8	-18	-28
C	-50	-32	-14	4	11	18	13	20	10	0	-10
T	-60	-42	-24	-6	12	8	26	16	17	18	8
G	-70	-52	-34	-16	2	20	10	23	18	14	26

Back-tracking ...

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	16	6	-4	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	4	-6	-16	-26	-36	-46
T	-40	-22	-4	8	21	11	12	2	-8	-18	-28
C	-50	-32	-14	4	11	18	13	<b>20</b>	<b>10</b>	0	-10
T	-60	-42	-24	-6	12	8	26	16	17	<b>18</b>	8
G	-70	-52	-34	-16	2	20	10	23	18	14	<b>26</b>

Back-tracking ...

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	16	6	-4	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	4	-6	-16	-26	-36	-46
T	-40	-22	-4	8	21	11	12	2	-8	-18	-28
C	-50	-32	-14	4	11	18	13	20	10	0	-10
T	-60	-42	-24	-6	12	8	26	16	17	18	8
G	-70	-52	-34	-16	2	20	10	23	18	14	26

Back-tracking ...

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	16	6	-4	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	4	-6	-16	-26	-36	-46
T	-40	-22	-4	8	21	11	12	2	-8	-18	-28
C	-50	-32	-14	4	11	18	13	20	10	0	-10
T	-60	-42	-24	-6	12	8	26	16	17	18	8
G	-70	-52	-34	-16	2	20	10	23	18	14	26

Back-tracking ...



## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	16	6	-4	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	4	-6	-16	-26	-36	-46
T	-40	-22	-4	8	21	11	12	2	-8	-18	-28
C	-50	-32	-14	4	11	18	13	20	10	0	-10
T	-60	-42	-24	-6	12	8	26	16	17	18	8
G	-70	-52	-34	-16	2	20	10	23	18	14	26

Back-tracking ...

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	16	6	-4	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	4	-6	-16	-26	-36	-46
T	-40	-22	-4	8	21	11	12	2	-8	-18	-28
C	-50	-32	-14	4	11	18	13	20	10	0	-10
T	-60	-42	-24	-6	12	8	26	16	17	18	8
G	-70	-52	-34	-16	2	20	10	23	18	14	26

Back-tracking ...

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	8	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	<b>16</b>	<b>6</b>	<b>-4</b>	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	<b>4</b>	-6	-16	-26	-36	-46
T	-40	-22	-4	8	21	11	<b>12</b>	2	-8	-18	-28
C	-50	-32	-14	4	11	18	13	<b>20</b>	<b>10</b>	0	-10
T	-60	-42	-24	-6	12	8	26	16	17	<b>18</b>	8
G	-70	-52	-34	-16	2	20	10	23	18	14	<b>26</b>

Back-tracking ...

## Example of global alignment:

Align the two sequences CACTGTCATG and CAGTCTG

		C	A	C	T	G	T	C	A	T	G
	0	-10	-20	-30	-40	-50	-60	-70	-80	-90	-100
C	-10	<b>8</b>	-2	-12	-22	-32	-42	-52	-62	-72	-82
A	-20	-2	<b>16</b>	<b>6</b>	<b>-4</b>	-14	-24	-34	-44	-54	-64
G	-30	-12	6	13	3	<b>4</b>	-6	-16	-26	-36	-46
T	-40	-22	-4	8	21	11	<b>12</b>	2	-8	-18	-28
C	-50	-32	-14	4	11	18	13	<b>20</b>	<b>10</b>	0	-10
T	-60	-42	-24	-6	12	8	26	16	17	<b>18</b>	8
G	-70	-52	-34	-16	2	20	10	23	18	14	<b>26</b>

Back-tracking ...

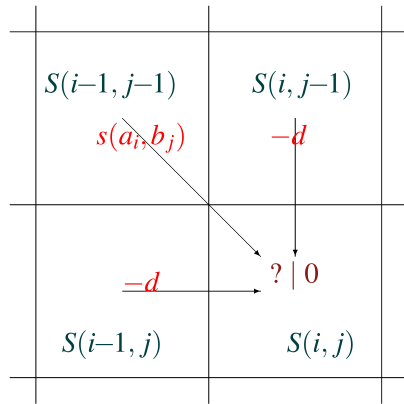
## Pairwise local alignments (Smith–Waterman)

Comparing sequences  $a$  and  $b$ . Given a substitution score  $s(x, y)$  of replacing letter  $x$  with letter  $y$ , the highest scoring alignment can be found by the following recursion:

$$S(i, j) = \max \begin{cases} S(i-1, j-1) + s(a_i, b_j) \\ S(i-1, j) - d \\ S(i, j-1) - d \\ 0 \end{cases}$$

Note only positive numbers!

$i = 1, \dots, N; j = 1, \dots, M$



Initialization:  $S(0, 0) = 0$ . Hence:  $S(i, 0) = ?$ ,  $S(0, j) = ?$ .

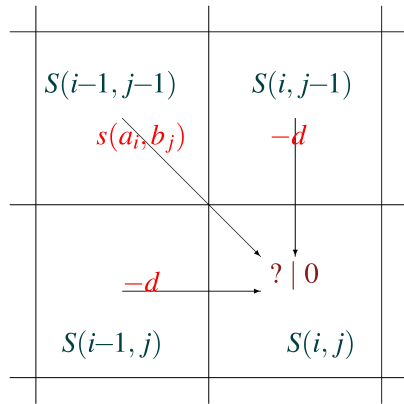
## Pairwise local alignments (Smith–Waterman)

Comparing sequences  $a$  and  $b$ . Given a substitution score  $s(x, y)$  of replacing letter  $x$  with letter  $y$ , the highest scoring alignment can be found by the following recursion:

$$S(i, j) = \max \begin{cases} S(i-1, j-1) + s(a_i, b_j) \\ S(i-1, j) - d \\ S(i, j-1) - d \\ 0 \end{cases}$$

Note only positive numbers!

$i = 1, \dots, N; j = 1, \dots, M$



Initialization:  $S(0, 0) = 0$ . Hence:  $S(i, 0) = 0$ ,  $S(0, j) = 0$ .

## Example of local alignment:

Align the two sequences **AAACTGTTTAAACAG** and **AACAGGGGAAACTG**.

		A	A	A	C	T	G	T	T	T	A	A	C	A	G
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0														
A	0														
C	0														
A	0														
G	0														
G	0														
G	0														
G	0														
A	0														
A	0														
A	0														
C	0														
T	0														
G	0														

## Example of local alignment:

Align the two sequences **AAACTGTTTAAACAG** and **AACAGGGGAAACTG**.

		A	A	A	C	T	G	T	T	T	A	A	C	A	G
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	8	8	8	0	0	2	0	0	0	8	8	0	8	2
A	0	8	16	16	6	0	2	0	0	0	8	16	6	8	10
C	0	0	6	13	24	14	4	4	2	2	0	6	24	14	5
A	0	8	8	14	14	21	16	6	1	0	10	8	14	32	22
G	0	2	10	10	11	11	29	19	9	0	2	12	5	22	40
G	0	2	4	12	7	8	19	26	16	6	2	4	9	12	30
G	0	2	4	6	9	4	16	16	23	13	8	4	1	11	20
G	0	2	4	6	3	6	12	13	13	20	15	10	1	3	19
A	0	8	10	12	3	0	8	9	10	10	28	23	13	9	9
A	0	8	16	18	9	0	2	5	6	7	18	36	26	21	11
A	0	8	16	24	15	6	2	0	2	3	15	26	33	34	24
C	0	0	6	14	32	22	12	4	2	4	5	16	34	30	31
T	0	0	0	4	22	40	30	20	12	10	1	6	24	31	27
G	0	2	2	2	12	30	48	38	28	18	12	3	14	26	39



## Example of local alignment:

Align the two sequences **AAACTGTTTAAACAG** and **AACAGGGGAAACTG**.

		A	A	A	C	T	G	T	T	T	A	A	C	A	G
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	8	8	8	0	0	2	0	0	0	8	8	0	8	2
A	0	8	16	16	6	0	2	0	0	0	8	16	6	8	10
C	0	0	6	13	24	14	4	4	2	2	0	6	24	14	5
A	0	8	8	14	14	21	16	6	1	0	10	8	14	32	22
G	0	2	10	10	11	11	29	19	9	0	2	12	5	22	40
G	0	2	4	12	7	8	19	26	16	6	2	4	9	12	30
G	0	2	4	6	9	4	16	16	23	13	8	4	1	11	20
G	0	2	4	6	3	6	12	13	13	20	15	10	1	3	19
A	0	8	10	12	3	0	8	9	10	10	28	23	13	9	9
A	0	8	16	18	9	0	2	5	6	7	18	36	26	21	11
A	0	8	16	24	15	6	2	0	2	3	15	26	33	34	24
C	0	0	6	14	32	22	12	4	2	4	5	16	34	30	31
T	0	0	0	4	22	40	30	20	12	10	1	6	24	31	27
G	0	2	2	2	12	30	48	38	28	18	12	3	14	26	39

## Example of local alignment:

Align the two sequences **AAACTGTTTAAACAG** and **AACAGGGGAAACTG**.

		A	A	A	C	T	G	T	T	T	A	A	C	A	G
	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	8	8	8	0	0	2	0	0	0	8	8	0	8	2
A	0	8	16	16	6	0	2	0	0	0	8	16	6	8	10
C	0	0	6	13	24	14	4	4	2	2	0	6	24	14	5
A	0	8	8	14	14	21	16	6	1	0	10	8	14	32	22
G	0	2	10	10	11	11	29	19	9	0	2	12	5	22	40
G	0	2	4	12	7	8	19	26	16	6	2	4	9	12	30
G	0	2	4	6	9	4	16	16	23	13	8	4	1	11	20
G	0	2	4	6	3	6	12	13	13	20	15	10	1	3	19
A	0	8	10	12	3	0	8	9	10	10	28	23	13	9	9
A	0	8	16	18	9	0	2	5	6	7	18	36	26	21	11
A	0	8	16	24	15	6	2	0	2	3	15	26	33	34	24
C	0	0	6	14	32	22	12	4	2	4	5	16	34	30	31
T	0	0	0	4	22	40	30	20	12	10	1	6	24	31	27
G	0	2	2	2	12	30	48	38	28	18	12	3	14	26	39

## Time is important

- Dynamic programming: exact in  $O(NM)$ .

When becomes time an issue?

# BLAST (Basic Local Alignment Search Tool)

Less accurate than Smith-Waterman, **BUT** 50 times faster.

Idea: true matches are likely to have short stretches of identity (high score).

- ① List of short words of fixed length that will match the query sequence (word length: 3 for protein; 11 for nucleic acids).
- ② Scan database for these words. Extend matches in both directions in an attempt to find an alignment with a score exceeding  $S$ .

Segment pairs whose scores cannot be improved by extending or trimming are called high scoring pairs (HSPs).

What are the default parameter settings of NCBI blastn and megablast?

Which differences in the alignments do you expect based on their parameters?

## Alignment score statistics

Question: Given a particular scoring system, how many distinct local alignments with score  $\geq S$  can one expect to find by chance from the comparison of two random sequences of lengths  $m$  and  $n$ ?

Or in other words, when can a local alignment be considered statistically significant?

## E-values and P-values

The expected number of local alignments with a score of at least  $S$  is given by the E-value for the score  $S$ :

$$E = Kmne^{-\lambda S}$$

- ① Doubling the length of the query sequence ( $m$ ) or the size of the database ( $n$ ) should double the number of local alignments.
- ② E-value decreases exponentially as score  $S$  increases.

The probability of observing *at least* one alignment with score  $\geq S$

$$p = 1 - e^{-E}$$

⇒ Sequence similarity score  $S$  is *extreme value distributed*



# Summary

- Dynamic programming (DP) saves time in sequence comparisons
- Some assumptions in DP, mention some
- In many applications, heuristics are needed to further speed up the comparison, e.g., use only diagonals in dynamic programming

## Exercise: Dynamical programming of pairwise alignment

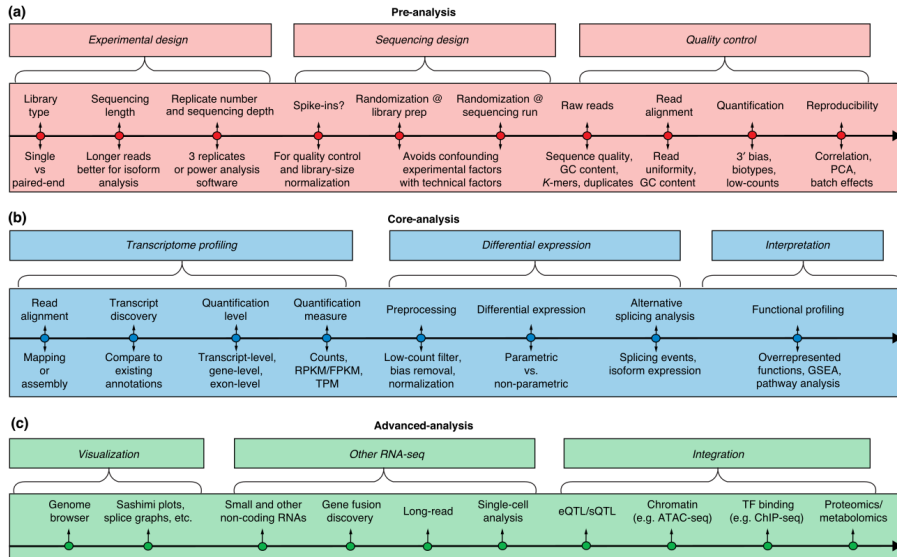
Complete the dynamic programming matrix of a global alignment:

Align the sequences **ACGTG** and **AACGGTG** using a match score of 1, a mismatch of -4 and a gap cost of -10.

		A	C	G	T	G
A						
A						
C						
G						
G						
T						
G						

# Read mapping

# Roadmap for RNAseq analysis



## When becomes time an issue?

- Target is entire genome
- Target is all observed sequences (e.g. RefSeq non-redundant database)
- **Query are millions of reads**

## How many reads do we get from modern Illumina sequencing?

	MiSeq	HiSeq 4000	NovaSeq 6000 S4
Run Time	4-56 hours	2-4 days	36-44 hours
Maximum Output	15 Gb	1500 Gb	2400-3000 Gb
Average Read Output	22 - 25 million	250 - 400 million	2,000 - 2,500 million
Maximum Read Length	2 × 300 bp	2 × 150 bp	2 × 150 bp

# Raw data (Sequencing reads)

The FASTQ format:

```
@ERR459145.1 DHKW5DQ1:219:D0PT7ACXX:2:1101:1590:2149/1
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGATCGGAAGAGCGGTTCAGC
+
@7<DBADDDDBH?DHHI@DH>HHHEGHIIIGGIFFGIBFAAGAFHA'5?B@D
```

- @: begin header
- 2:1101 flowcell lane 2, tile 1101
- x and y coordinates: 1590:2149
- /1 single-end reads;  
/1 and /2 paired-end (mate-paired) reads
- read sequence
- quality encoded ASCII characters

# What is read mapping?

Determine position of a short read on the reference genome or transcriptome.

Reference:	...AA-CGCCTT...	= match
	:-:	: = mismatch
Read:	AGGGGCCTT	- = gap



## Naive mapping

Search for query at each position in reference genome

ACGTTACCGAATCGATCAAAGTCGA

GTTA

$m$  = query length,  $n$  = genome length

## Naive mapping

Search for query at each position in reference genome

ACGTTACCGAATCGATCAAAGTCGA

GTTA

$m$  = query length,  $n$  = genome length

## Naive mapping

Search for query at each position in reference genome

ACGTTACCGAATCGATCAAAGTCGA

GTТА :)

$m = \text{query length}, n = \text{genome length} \rightarrow \text{Time: } O(mn)$

## Naive mapping

- Human Genome (queries) would take far too long:
  - Illumina/Solexa sequencing technology produces **50 – 200 million, 32 – 100 bp short reads**
  - Mapping these reads to a **3.2 billion bp** human genome is a challenge
- Far worse when we allow for Indels and mismatches.

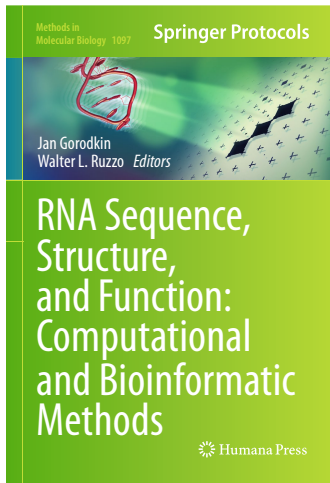
→ Are optimal alignments based on quality scores still feasible?

# Principles of mapping reads

- Most computational time is spent on alignment.
- **Many sequenced reads are redundant**
- We do not need to search the entire genome each time again.

# Book analog

Do not search the entire book, instead search the book **index**.



INDEX	
<b>A</b>	
Ab-initio	407, 408
Abstract shape analysis	102, 215–242
Adenosine	34
Adenosine platform	39
Affinity	40, 364, 493, 495, 498, 499, 501, 504, 505, 509, 513
Alignment	
full	176, 286
gaps	15, 128, 139, 264, 267, 285, 295
seed	111, 116, 172, 176
All-stem plot	286–289
ALPS	446, 449
AMBER	399, 400, 407, 410–412
Ambiguity	
avoidance	12
semantic	97, 100, 101
syntactic	89, 101
Aminoacyl-tRNAs	40
Ancestral correlations	363
Annotations	
annotated	110–112
false	110, 111, 115, 116
pipeline	114, 117, 193, 203
AntiBiotics	40
Antisense	2, 437, 438, 478, 481, 504, 509
Aptamer(s)	40, 227, 264, 396
Argon	167, 187, 189–191, 201–204, 210
ARB	120, 280
Arg-Associated Sequences	252–254, 268, 269
Argonant	458, 500
Assessment	57, 207
Assesses	402, 494
<b>B</b>	
Backbone	283, 297, 299–401, 405, 406, 409, 411, 497–499, 502, 504, 509
Backbone torsion	401, 411
Backtracking	9–11, 79–80, 130, 152
Barrier rise	82, 83, 239, 240
Base pair	
canonical covariance	56, 299, 408
correspondence	402
direct	429
distance	78, 80, 216, 254–255
indirect	429
intermolecular	426, 428–430
intramolecular	421, 425, 428, 429, 432, 483
model non-canonical	18, 271
probability	81, 254, 282, 283, 432
set representation	250
stacking interaction	281, 401, 409
Watson-Crick	49, 166, 171, 180, 187, 338, 380, 381, 384, 385, 390, 405, 406, 410
Base-pairing probability	512
Base pair types	
bifurcated	385
cis Hoogsteen/Hoogsteen	385
cis Hoogsteen/sugar edge	385
cis sugar edge/sugar edge	385
cis Watson-Crick/Hoogsteen	385
cis Watson-Crick/sugar edge	385
cis Watson-Crick/Watson-Crick	385
trans Hoogsteen/Hoogsteen	385
trans Hoogsteen/sugar edge	385
trans sugar edge/sugar edge	385
trans Watson-Crick/Hoogsteen	385
trans Watson-Crick/sugar edge	385
trans Watson-Crick/Watson-Crick	385
Base triple	6, 23, 29, 268, 233, 385, 391
Bchick	167, 187, 189–191, 201, 206, 210
Bellman's GAP	162, 236, 238, 241
Benchmarks	20, 21, 23, 24, 210, 211, 296, 401, 481, 483
Big O	11
BiFold	380
BiFold	482, 483
Bit (unit of information)	12, 89, 95, 171, 176, 178, 179, 182, 184, 189, 190, 249, 264
BLAST	5, 19, 111, 112, 117, 118, 386, 396, 402, 418, 444, 447
Blockbuster	449
Boltzmann sampling	79, 80
Boltzmann weight	80, 218, 220, 221, 230, 230, 235, 236, 423, 424, 426, 483
Boltzmann-weighted energies	218, 423, 424, 426
Boukier ALE	380, 381
Bowtie	448
Breast cancer	496
BWA	448
<b>C</b>	
Carnac	292–294, 297, 298, 307
Carrying capacity	225, 226
CASP	39, 296
Cationic	411

Jan Gorodkin and Walter L. Ruzzo (eds.), RNA Sequence, Structure, and Function: Computational and Bioinformatic Methods, Methods in Molecular Biology, vol. 1097, DOI 10.1007/978-1-62703-709-9, © Springer Science+Business Media New York 2014

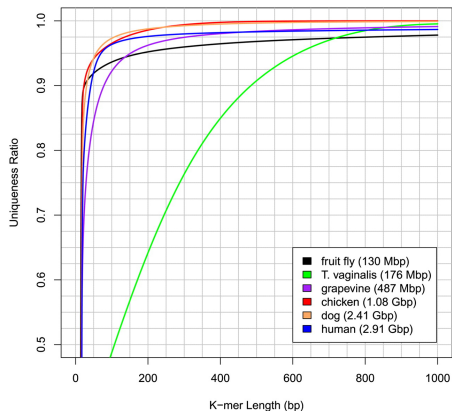
## k-mer

"k-mer" is a substring of length  $k$

For example sequence **GGCGATTTCATCG**:

4-mer GGCG, GCGA, CGAT

3-mer GGC, GCG, CGA, GAT



## Read mapping through indexing

In principle read mapping is to map an exact piece of sequence to the genome.

### Solution

An index is a data structure that improves the speed of data retrieval operations at the cost of additional storage space to maintain the index data structure.

### Pros

Quick search for matches in an entire genome.

### Cons

Index structure of the entire genome takes a lot of memory.



Why might an exact mapping not always be what we want?

Which concerns might you have when mapping genomic sequence?

Which concerns might you have when mapping transcribed sequence?

# Indexing problems

Flexibility and constraints:

- *Errors versus natural variation:*  
trade-off in error threshold.
- *Computational efficiency (time / memory):*  
allowed mismatches / unique mappings.
- Balance between speed, memory and reported mappings.

## Indexing method choice is crucial!

- **Hash-based** (BLAST, Salmon, Kallisto)
- **Suffix arrays** (Salmon, STAR)  
A sorted table of all suffixes (substrings) of a given string
- **Burrows-Wheeler Transform** (BWA, SOAP2, Bowtie2, Hisat)  
A compressed form of suffix arrays

# Indexing is often used for seed matching

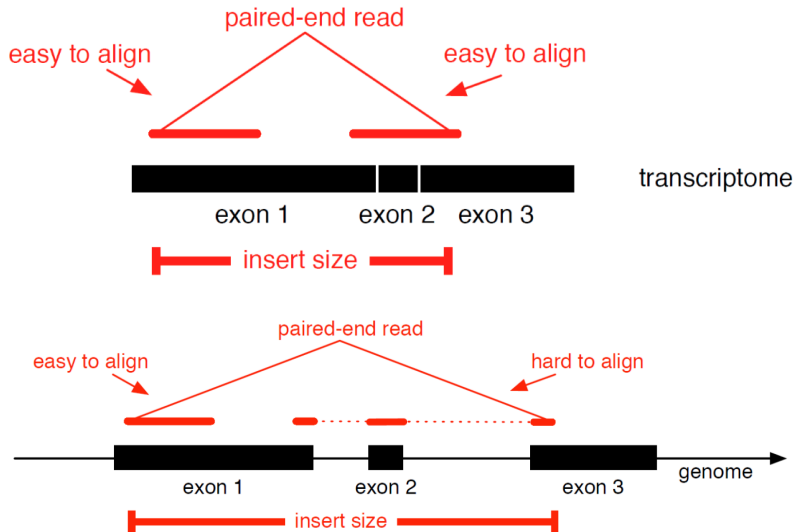
## **”Seed-and-extend“ approach**

- ① find the best possible match of a seed in an index made up from the reference genome
- ② every matched seed is extended on both sides by optimal local alignment

Common software

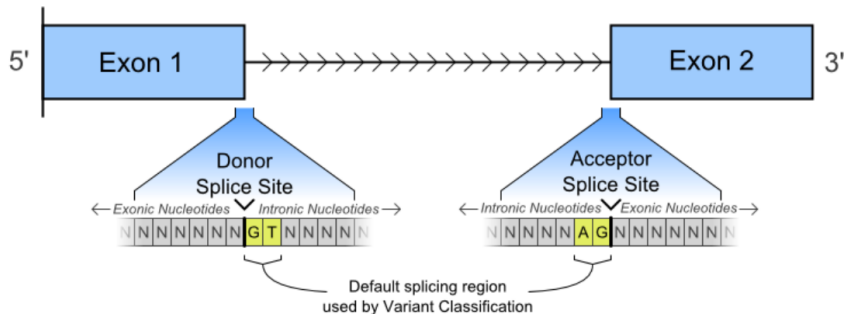
STAR, HISAT2, BLAST

# Transcriptome versus genome mapping



## Splice-aware genome mapping

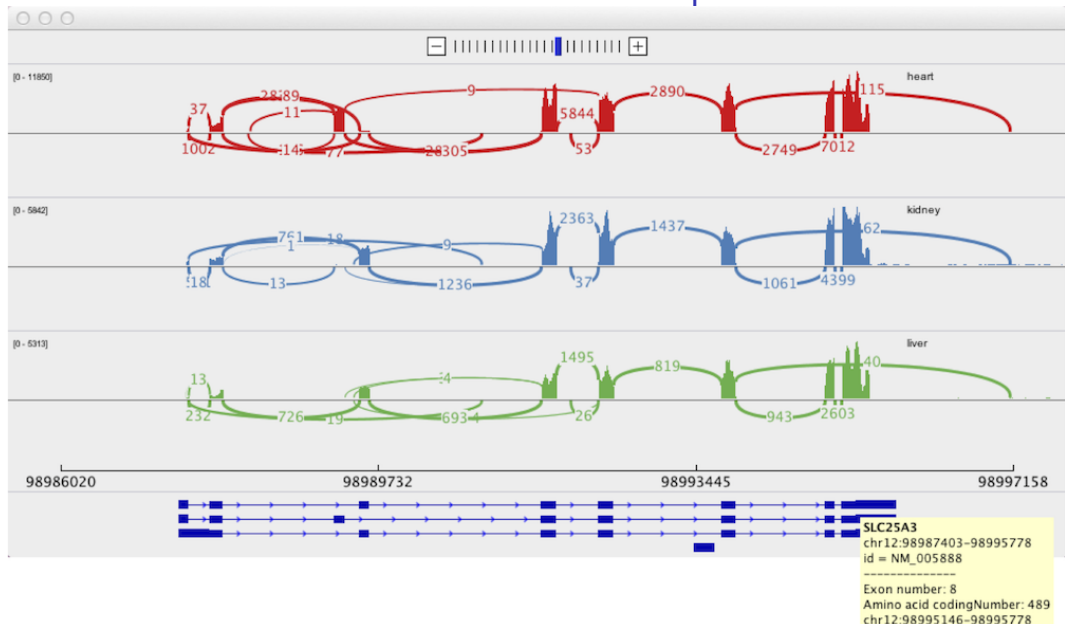
There are many mappers, but they must be able to detect splice junctions to be used in transcriptome assembly and quantification.



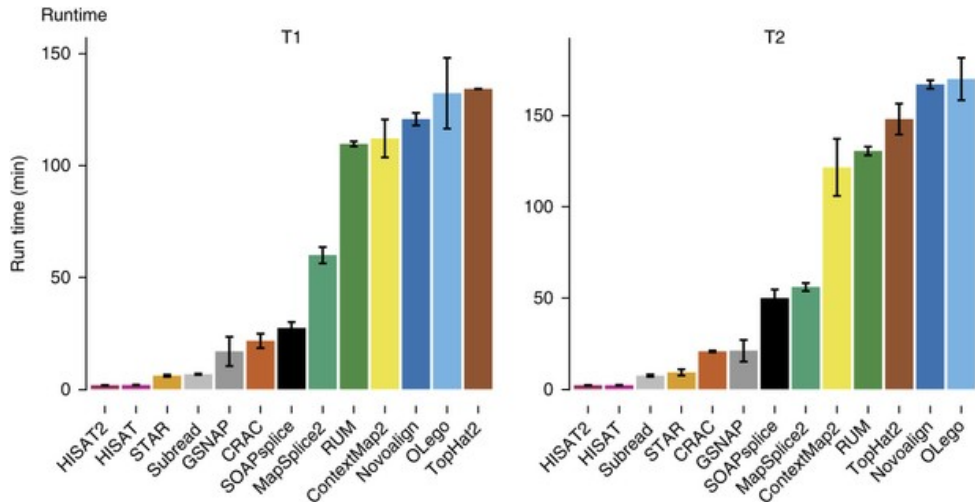
### Common splice-aware alignment software

STAR, HISAT2, BLAT, TopHat (based on Bowtie2), Segemehl

# Sashimi Plot: Visualization of spliced reads

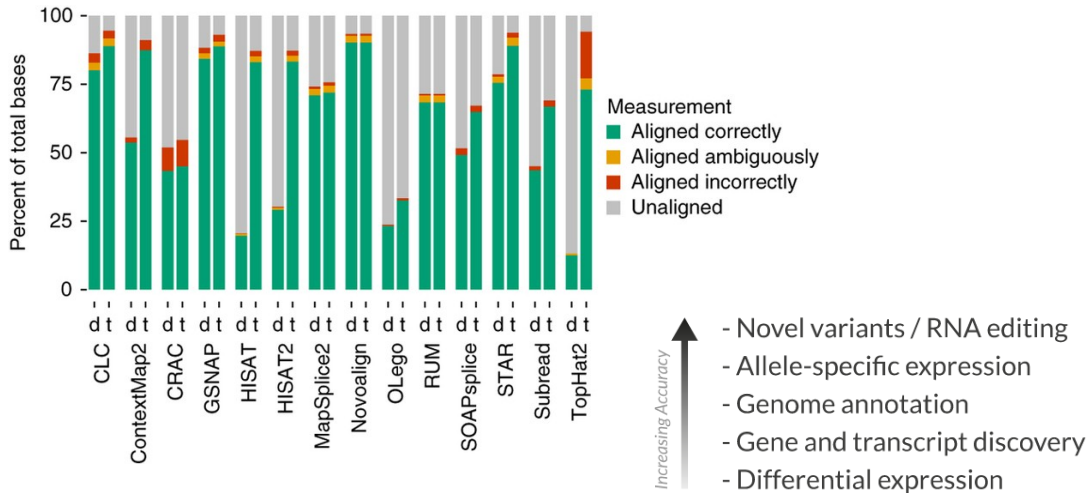


## Aligner's speed





## Aligner's accuracy



## Quantification – estimation of expression

Quantification:

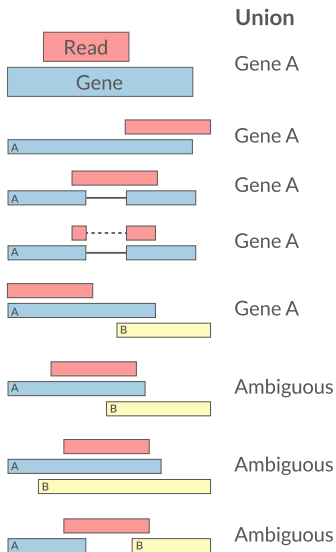
Count reads mapping to a genomic feature (gene, transcript, exon etc)

→ genomic features can be annotated (e.g., NCBI RefSeq genome annotation) or be predicted (e.g., transcriptome assembly)

Assumption:

Number of reads produced from a feature  $\sim$  feature's **relative** abundance in the sample

# Gene-level quantification



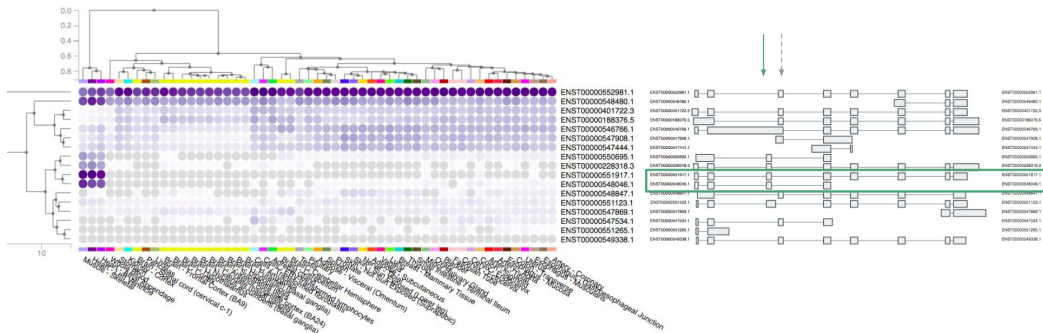
- By default ambiguous reads are not counted
- By default multi-mapping reads (reads aligning with multiple locations) are not counted

Common read counting software

featureCounts, HTSeq

## Transcript-level quantification

Quantify per transcript (not just exon or gene)



Why is it important for transcript-level quantification to consider ambiguous reads?

# Transcript-level quantification

- Statistical methods are used to find the probability that a read originates from a specific transcript.
- ! Beware of annotation quality
- Suggested reading: [CSAMA tutorial](#), Zhang *et al.* “Evaluation and comparison of computational tools for RNA-seq isoform quantification.” BMC Genomics 2017

## Common software

RSEM, Cufflinks, Kallisto, Salmon, Sailfish

## Alignment-free methods ("pseudoaligners")

- ① Focus: Only annotated transcripts (not entire genome!)
- ② Pseudoalign: K-mer composition of reads/transcripts
- ③ Abundance *estimation*
- ④ GC-content, transcript position correction included (e.g. 3' end degradation)

### Pros

dramatic increase in speed; improvements in accuracy for gene-level quantification

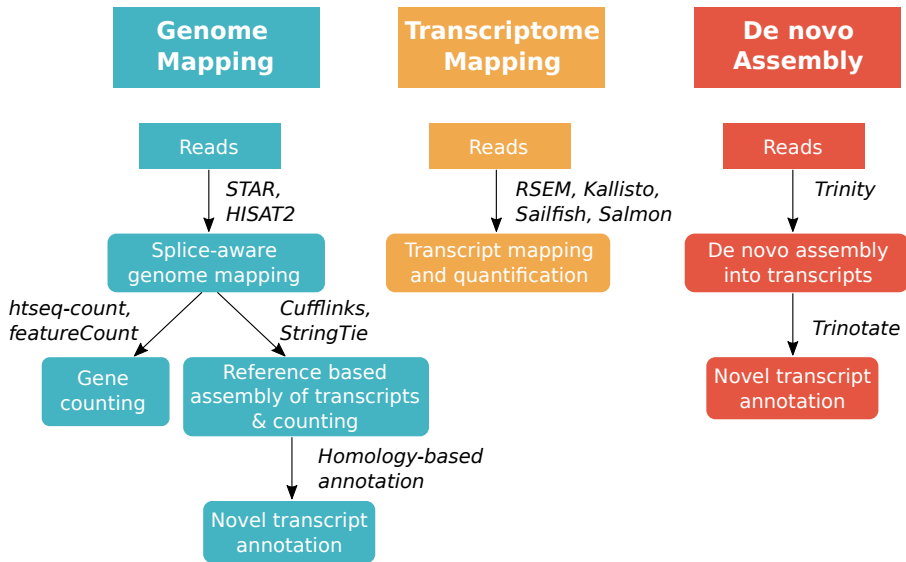
### Cons

absolute reliance on a precise and comprehensive transcript annotation; no information on where each read is mapping

### Common software

Salmon, Kallisto, Sailfish

## Summary: From reads to count matrix





# Normalization + Transformation

## Count matrix

	SRR1039508	SRR1039509	SRR1039512	SRR1039513	SRR1039516	SRR1039517	SRR1039520	SRR1039521
ENSG00000000003	679	448	873	408	1138	1047	770	572
ENSG00000000005	0	0	0	0	0	0	0	0
ENSG000000000419	467	515	621	365	587	799	417	508
ENSG000000000457	260	211	263	164	245	331	233	229
ENSG000000000460	60	55	40	35	78	63	76	60
ENSG000000000938	0	0	2	0	1	0	0	0
ENSG000000000971	3251	3679	6177	4252	6721	11027	5176	7995
ENSG000000001036	1433	1062	1733	881	1424	1439	1359	1109
ENSG000000001084	519	380	595	493	820	714	696	704
ENSG000000001167	394	236	464	175	658	584	360	269
ENSG000000001460	172	168	264	118	241	210	155	177
ENSG000000001461	2112	1867	5137	2657	2735	2751	2467	2905
ENSG000000001497	524	488	638	357	676	806	493	475
ENSG000000001561	71	51	211	156	23	38	134	172
ENSG000000001617	555	394	905	415	727	697	618	599
ENSG000000001626	10	2	9	2	10	6	5	5
ENSG000000001629	1660	1251	2259	1079	2462	2514	1888	1660
ENSG000000001630	59	54	66	23	84	87	31	59
ENSG000000001631	729	692	943	475	1034	1163	731	744
ENSG000000002016	201	161	256	99	268	257	160	137
ENSG000000002079	3	0	3	1	4	0	0	1

## Estimation of gene expression

You are analyzing 2 genes (gene A and B) in two conditions (condition 1 and 2) on the bases of an RNA-seq experiment that resulted in the following number of reads:

	Condition 1	Condition 2
Gene A	1000	3000
Gene B	2000	4000

Are the following statements correct?

- a Both genes A and B are more expressed in condition 2.
- b Gene B is more expressed than gene A.

## Estimation of gene expression

We cannot state any such thing since we do not know

- a **sequencing depth (library size),**  
expression of all other genes within the sample

→ RNA-seq data informs about the relative abundance BUT NOT about the absolute abundance.

- b **gene length,**  
the longer the gene, more reads will be mapped

## Estimation of gene expression

We cannot state any such thing since we do not know

- a **sequencing depth (library size),**  
expression of all other genes within the sample

→ RNA-seq data informs about the relative abundance BUT NOT about the absolute abundance.

- b **gene length,**  
the longer the gene, more reads will be mapped

**Solution:** Control (Normalize) for

- 1 sequencing depth
- 2 compositional bias

## Normalization

$R$  = Reads count for the gene

$G$  = Gene length in kilobases

$T$  = Total number of mapped reads in a sample

**Reads/Fragments Per Kilobase  
transcript per Million mapped reads  
(RPKM/FPKM):**

$$RPKM/FPKM = \frac{\left(\frac{R}{T/10^6}\right)}{G}$$

*RPKM* is used for single end reads

*FPKM* is used for paired end reads.

- First metric used in the old times (e.g. cufflinks).

**Counts Per Million (CPM):**

$$CPM = \frac{R * 10^6}{T}$$

- Sum of all CPMs is constant (1 million).
- Does not consider gene lengths.

## Simple normalization could fail

Genes	Control	Treated
Gene A	10	30
Gene B	30	90
Gene C	5	15
Gene D	1	3
Gene N	1000	240
Total	1046	378

## DESeq2's median of ratios (sample-wise "size factor")

**Assumption: most genes are not differential expressed**

1. For each gene, calculate geometric mean

Genes	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6	Geomean
Gene 1	34	56	23	12	10	30	23
Gene 2	10	6	7	11	12	8	9
Gene n	65	78	67	34	56	23	50

2. For each gene, calculate ratio to geometric mean

Genes	Sample 1	Sample 2	Sample 3	Sample 4	Sample 5	Sample 6
Gene 1	1.5	2.4	1.0	0.5	0.4	1.3
Gene 2	1.1	0.7	0.8	1.3	1.4	0.9
Gene n	1.3	1.6	1.4	0.7	1.1	0.5

3. Take median of these ratios as sample normalization factor ("size factor")

1.3

1.6

1

0.7

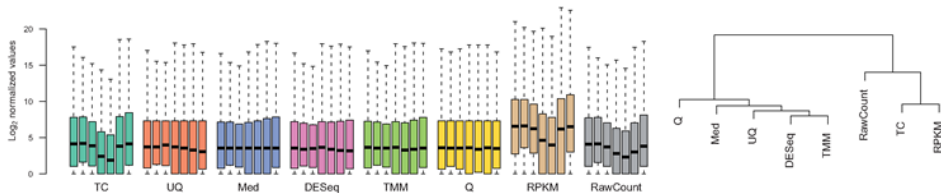
1.1

0.9



## Summary: Normalization

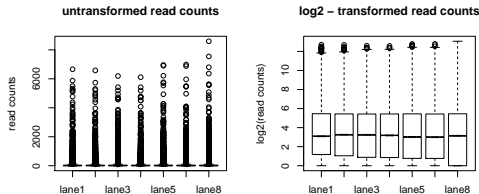
- Use **raw counts** when using DGE packages, e.g., DESeq2 and edgeR, as normalization is done internally
- RPKM/FPKM/CPM are not recommended for DGE analysis
- CPM is usable for visual data exploration (heatmap, abundance comparison, PCA)



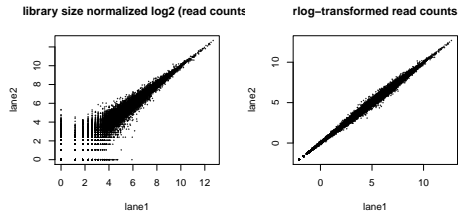
- Median of Ratios (DESeq2) and Trimmed Mean of M-values TMM (edgeR) perform the best for DGE analysis
- other solutions: spike-ins/house-keeping genes

# Transformation of sequencing-depth-normalized read counts

*Log*<sub>2</sub> transformation:



Transformation incl variance shrinkage:



For clustering, heatmaps etc use VST (DESeq2), VROOM (limma) or RLOG (DESeq2)

Challenge your data by different normalization methods → robustness of DGE analysis

## Alternative: Compositional data analysis

DESeq's median of ratios and edgeR's TMM are less suitable in highly **asymmetrical** or **sparse** datasets → unacceptably high false positive DEGs

- ratio transformations capture the relationships between the features in the dataset
- centered log-ratio (clr) transformation:

$$x_{clr} = [\log(x_1/G(x)), \log(x_2/G(x)) \dots \log(x_D/G(x))],$$
$$G(x) = \sqrt[p]{x_1 \cdot x_2 \cdot \dots \cdot x_D}$$

- clr-transformed values are scale-invariant (same ratio with few as well as many read counts)
- to calculate  $G(x)$  we have to delete, replace or **estimate** the 0 count values
- estimate technical variation within each sample using Monte-Carlo instances drawn from the Dirichlet distribution → probability vector prior to clr transformation

Unsupervised + Supervised data exploration

# Explore global and local read count patterns

## ① unsupervised

no *a priori* information is needed

→ to detect technical noise and batch effects

- Dimensional reduction → Principle Component Analysis (PCA)
- Clustering → hierarchical, k-means

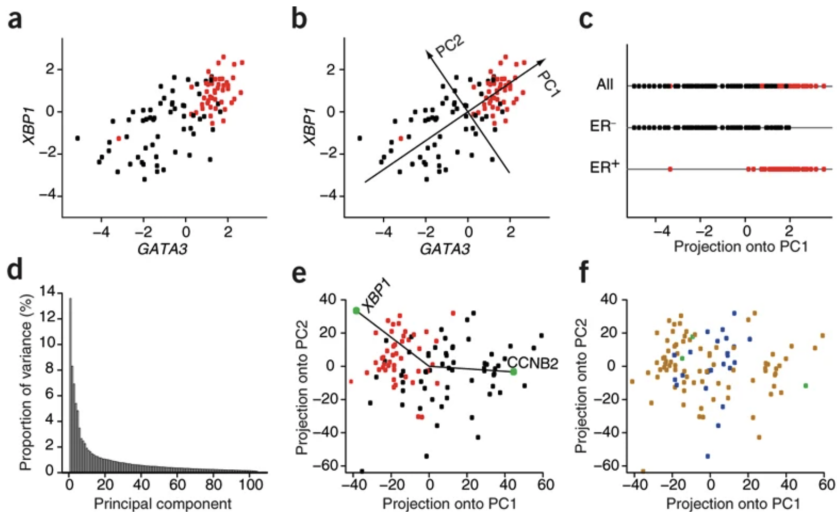
PCA and clustering should be done on normalized and transformed read counts so that high variability of low read counts does not occlude potentially informative data trends.

## ② supervised

usage of known biological labels

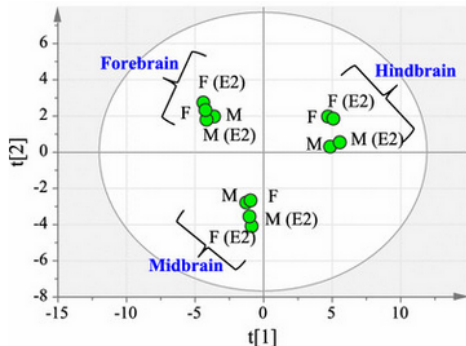
→ differential expression

# Principle component analysis – PCA



# Principle component analysis – PCA

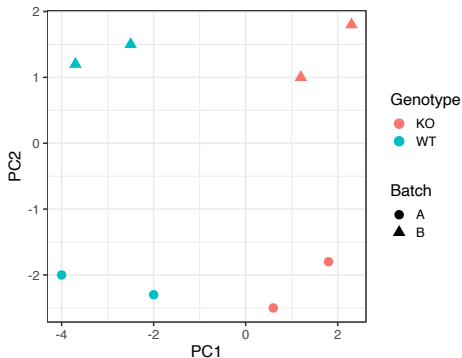
- transforms measurements into new variables that are truly independent
- new variables of most variance are the principal components
- dimensionality reduction



## Applications:

- visualization of your data in lower dimensions (2D, 3D)
- find patterns in numeric data
- identify batch effects or other possible covariates (e.g. male and female) by labeling them

## PCA to detect batch effects



- PC1 separates the *genotype* (group of interest); however
- PC2 separates the *batch* effect (or other covariates)
  - e.g. experiment date, sex, experimenter, different RNA isolation kit (you name it!)
- When batch effect is observed in PCA plots, add it as covariate to your GLM
  - $\sim \text{Batch} + \text{Genotype}$



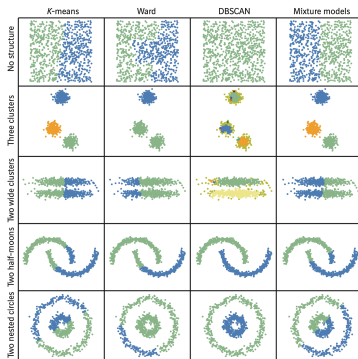
# Clustering

There are several clustering algorithms and more are being developed. Why?

# Clustering

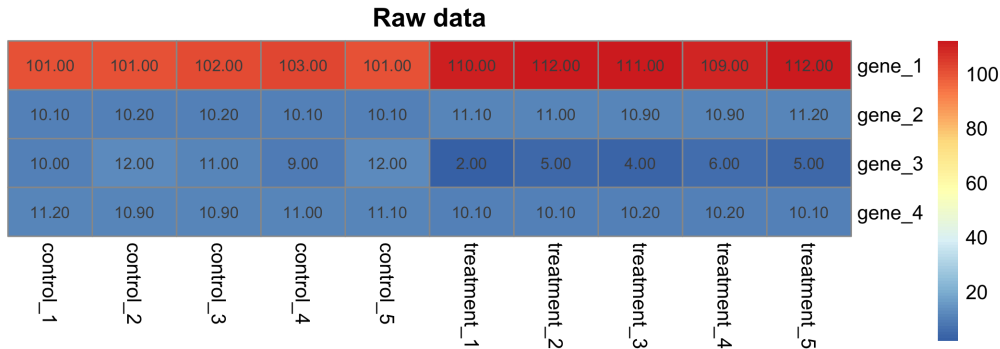
There are several clustering algorithms and more are being developed. Why?

- No clustering algorithm is perfect
- Remember! Always “see” your data and judge if the clusters make sense
- The performance of the clustering algorithm depends on the structure of your data



# Finding patterns in your data

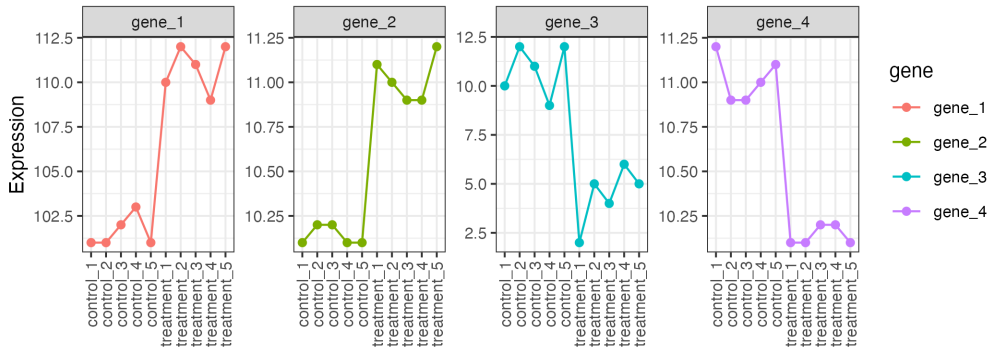
Let's assume that you got these 4 genes differentially expressed in your dataset



Can you identify a pattern?

# Finding patterns in your data

Perhaps this gives you a better hint

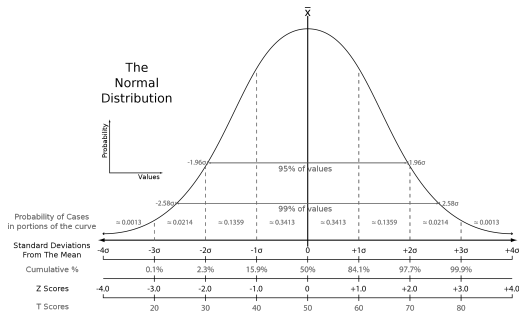


# Z-score

$$Z = \frac{x - \mu}{\sigma}$$

$\mu$  = mean

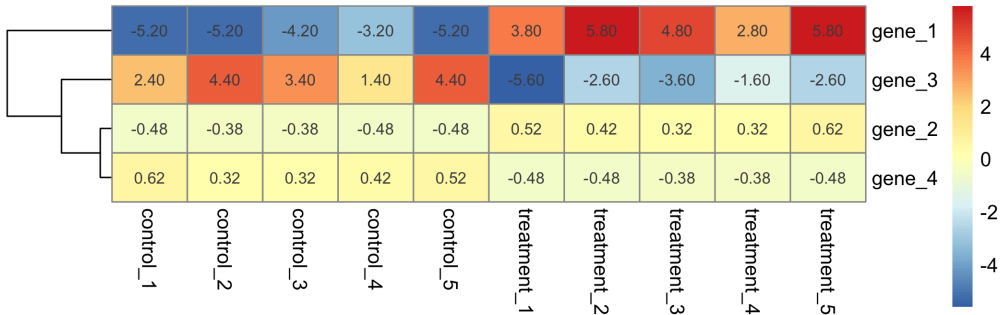
$\sigma$  = standard deviation



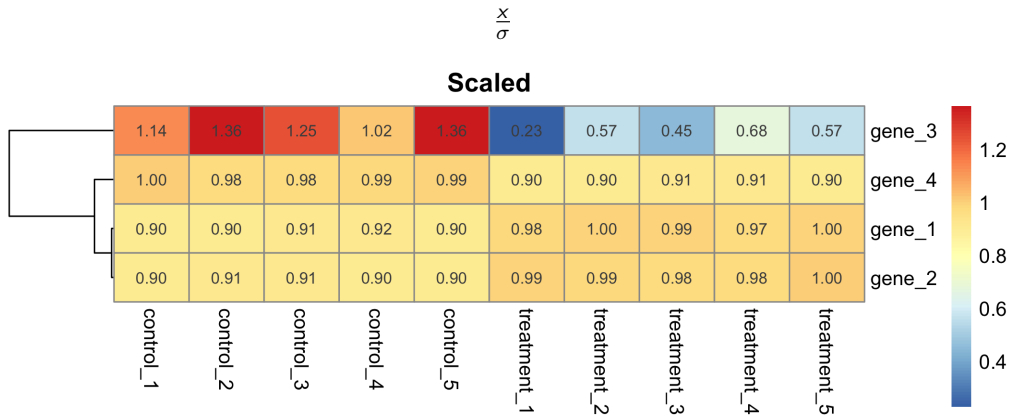
# Centering data

$$x - \mu$$

**Centered**



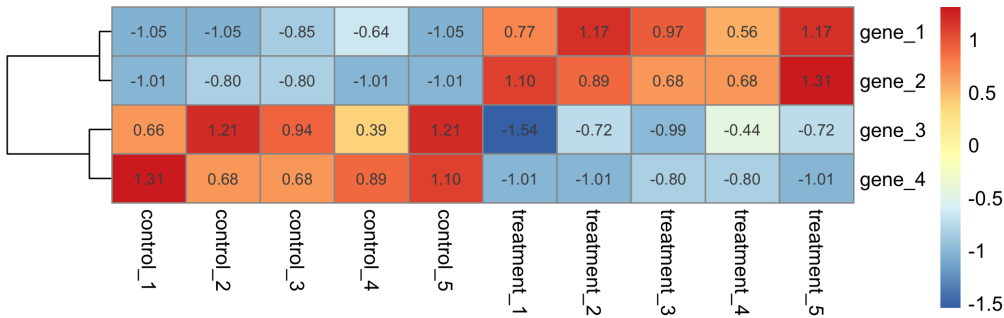
# Scaling data



# Z-score

$$Z = \frac{x - \mu}{\sigma}$$

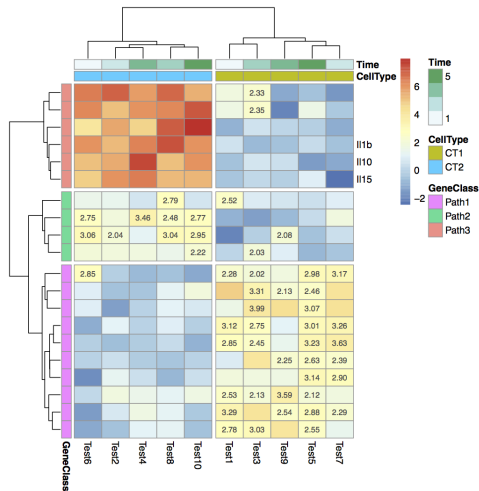
**z-score (centered & scaled)**





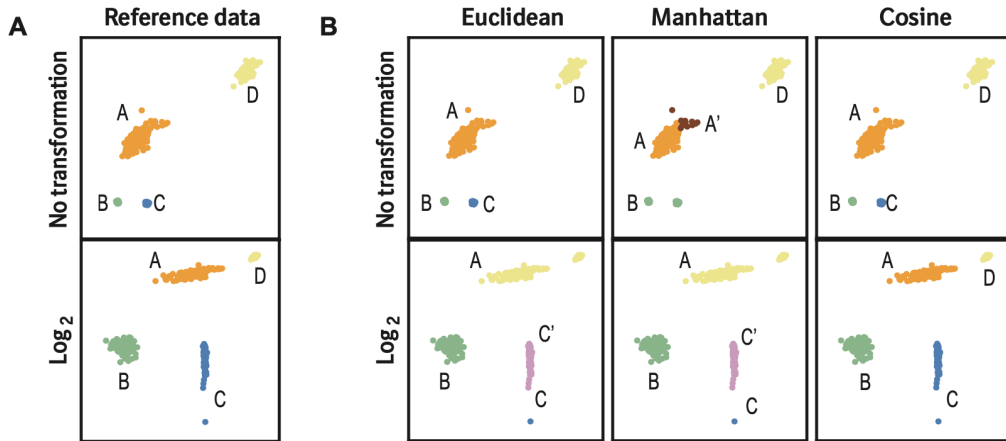
# Hierarchical clustering

- No need for pre-defining the number of clusters
- Results dependent on
  - 1 *distance metric*: euclidean, manhattan, Pearson correlation etc.
  - 2 *clustering method*: Ward, complete, average etc.
- Used with heatmaps, thus allows better visual inspection of data



# Effect of transformations and distance metrics on clustering

- use normalized read counts (rlog, CPM etc.)
- transform your data for better scaling (log2, z-score transformation etc.)



## Differentially expressed genes

Expression level of gene 1

### Control:

<i>Replicate 1</i>	24
<i>Replicate 2</i>	25
<i>Replicate 3</i>	27

### Treated:

<i>Replicate 1</i>	23
<i>Replicate 2</i>	26
<i>Replicate 3</i>	102

Is this a differentially expressed gene?

You might get different answers depending of which software you run.

# Differential gene expression

## **Common scientific question:**

Quantification and statistical inference of systematic changes between conditions.

Principles are the same as for all other significance tests:

- ① Use the replicates (samples of the same conditions) to estimate the within-condition variability (variance) of the expression
- ② Use the expression and variance to test whether the difference between conditions is random or not

The test's statistical power increases with more biological (and technical) replicates!

## Parametric vs. non-parametric methods

It would be nice to not have to assume anything about the expression value distributions but only use rank-order statistics.

However, it is hard to show statistical significance with non-parametric methods if only few replicates are available (less than 8).

## Issues with DGE analysis for RNA-seq

Couldn't we just use a Student's t-test for each gene?

- ① Distribution is not **normal**. Which parametric distribution should I use?
- ② **Variance** across groups may not be **homogeneous** e.g. unequal group size
- ③ The number of replicates is often too small to estimate the variance.
- ④ If we test each gene for DE, we have to account for **multiple testing**!

## Which parametric distribution should I use?

Models for read counts originated from the idea that each read is sampled independently from a pool of reads and hence the number of reads for a given gene follows a ...

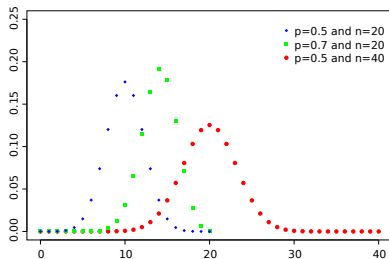
## Binomial distribution

The binomial distribution works when we have a fixed number of trials  $n$ , each with a constant probability of success  $p$ .

The random variable  $X$  is the number of  $k$  successes:

$$p(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

Event: An RNA-seq read "lands" in a given gene (success) or not (failure)



As RNA-seq experiments produce large number of reads ( $n$  is large) the Gaussian distribution can replace the binomial.



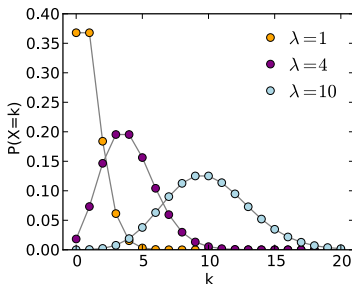
## Poisson distribution

RNA-seq experiments produce large number of reads ( $n$  is large) and probabilities of success are small ( $p$  is small) which can be modelled by the poisson distribution which is an approximation of the binomial.

Instead we know the average number of successes per intervall:

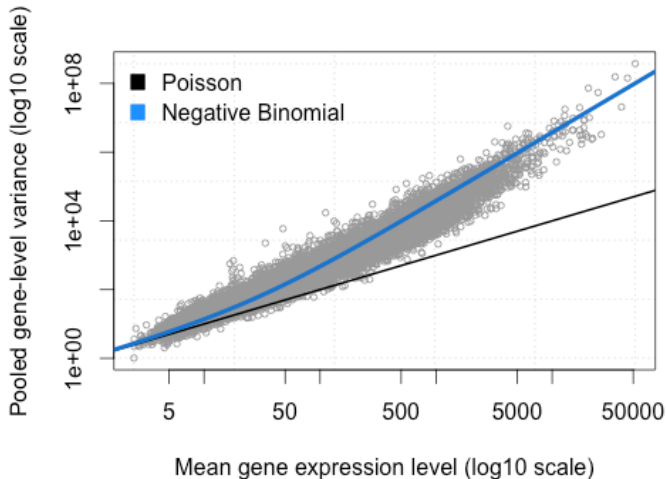
$$\lambda = np$$

For  $X \sim \text{Poisson}(\lambda)$ , both the mean and the variance are equal to  $\lambda$ .



## Poisson versus negative binomial distribution

Many studies have shown that the variance grows faster than the mean in RNAseq data. This is known as **overdispersion**.



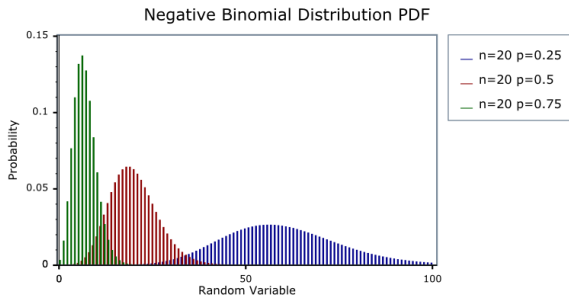
## Negative binomial distribution

The negative binomial distribution works for discrete data over an unbounded positive range whose sample variance exceeds the sample mean.

The random variable  $X$  is the number of trials needed to make  $r$  successes (and  $k$  failures) if the probability of a single success is  $p$ :

$$NB(X = k) = \binom{k + r - 1}{k} p^r (1 - p)^k$$

both the mean and variance can be calculated from  $r$  and  $p$



## Fitting a negative binomial GLM

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

Controls the variance

$$K_{ij} \sim \text{NB}(\text{mean} = \mu_{ij}, \text{dispersion} = \alpha_i)$$

$$\mu_{ij} = s_j q_{ij}$$

Normalization ("size") factor

Normalized count

Design matrix

- Control or Treatment?
- Batch (e.g., flow cell, plate, lab)
- Other co-factors (e.g., gender)

Linear Model:

$$\log_2 q_{ij} = \sum_r x_{jr} \beta_{ir}$$

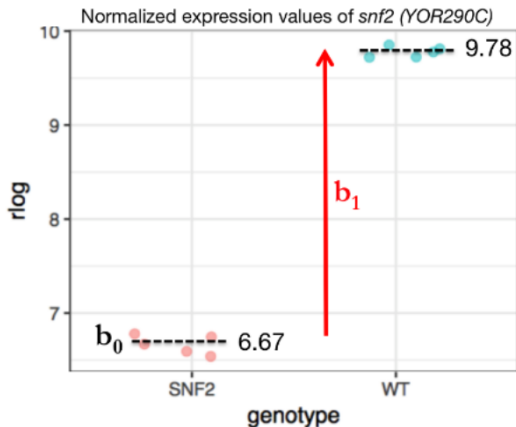
Generalized Linear Model (GLM) coefficients

- One for each Design matrix column (factor)  
= strength of effect  
= **log2 fold change** for each gene

Coefficient

$r$

## Interpretation of the negative binomial GLM



Linear regression model (LM) is evaluated for every gene:  $Y = b_0 + b_1 * x_1 + e$

$Y$  ... describes all read counts for a gene

$b_0$  ... average of baseline group, e.g., control

$x_1$  ... design factor,

e.g., condition (often 0 or 1)

$b_1$  ... coefficient that captures the difference between different conditions

$e$  ... error or uncertainty

→ the closeness of  $b_1$  to zero will be evaluated during statistical testing steps

→ DESeq2 and edgeR use a generalized linear model (GLM)

## Design & contrast matrix

Design matrix (a.k.a. model matrix) has 2 main roles:

- ① defines the form of the model, or structure of the relationship between genes and explanatory variables
- ② is used to store values of the explanatory variable(s)

Contrast matrix is used for:

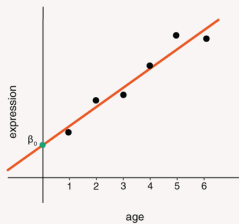
- ① identifying the differences (contrast) between explanatory variables  
e.g. *group*<sub>1</sub> vs *group*<sub>2</sub>

# Basic regression models

**Covariates:** quantitative measurements (e.g. age)

## Regression model

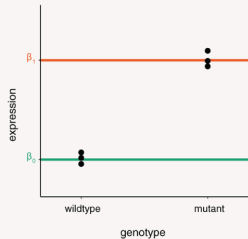
$$\text{expression} = \beta_0 + \beta_1 \text{age}$$



**Factors:** categorical variables (e.g. genotype)

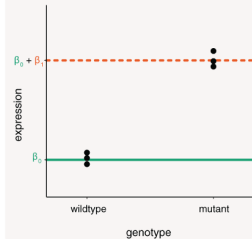
## Means model

$$\text{expression} = \beta_1 \text{wildtype} + \beta_2 \text{mutant}$$



## Mean-reference model

$$\text{expression} = \beta_1 + \beta_2 \text{mutant}$$



## Legend

- Original data points

— Expected gene expression  
(based on model)

- - Expected gene expression  
(of non-reference levels in mean-reference model)

# Design matrix with intercept

x is an indicator variable for **sick** mice:

- $x = 1$  for sick mice
- $x = 0$  otherwise

Model

$$E(y) = 2.95 + 1.62x$$

$$E(y) = 2.95 = 2.95 \quad (\text{for healthy group})$$

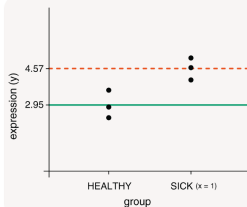
$$E(y) = 2.95 + 1.62 = 4.57 \quad (\text{for sick group})$$

Matrix

```
> model.matrix(~group)
```

$$\begin{array}{c} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{array} \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 1 & 1 \\ 1 & 1 \\ 1 & 1 \end{pmatrix} \begin{array}{l} \text{(Intercept)} \\ \text{groupSICK} \end{array}$$

Plot



$$\text{contrast} = c(0,1)$$



# Design matrix without intercept

x1 is an indicator variable for **healthy** mice:

- $x_1 = 1$  for healthy
- $x_1 = 0$  otherwise

x2 is an indicator variable for **sick** mice:

- $x_2 = 1$  for sick
- $x_2 = 0$  otherwise

Model

$$E(y) = 2.95x_1 + 4.57x_2$$

$$E(y) = 2.95$$

$$= 2.95 \quad (\text{for healthy group})$$

$$E(y) = 4.57$$

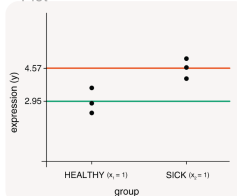
$$= 4.57 \quad (\text{for sick group})$$

Matrix

```
> model.matrix(~0 + group)
```

$$\begin{array}{c} \text{groupHEALTHY} \quad \text{groupSICK} \\ \begin{pmatrix} 1 & 0 \\ 2 & 0 \\ 3 & 0 \\ 4 & 0 \\ 5 & 1 \\ 6 & 1 \end{pmatrix} \end{array}$$

Plot



$$\text{contrast} = c(-1, 1)$$

# Design matrix with intercept

## Model

$$E(y) = 1.03 + 1.09x_1 + 1.97x_2 + 3.87x_3$$

$$E(y) = 1.03 = 1.03 \quad (\text{for control})$$

$$E(y) = 1.03 + 1.09 = 2.12 \quad (\text{for treatment I})$$

$$E(y) = 1.03 + 1.97 = 3.00 \quad (\text{for treatment II})$$

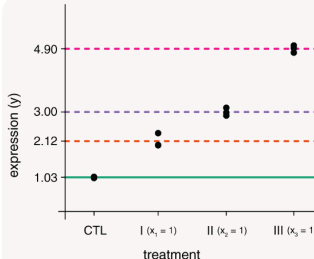
$$E(y) = 1.03 + 3.87 = 4.90 \quad (\text{for treatment III})$$

## Matrix

```
> model.matrix(~treatment)
```

$$\begin{matrix} & \text{(Intercept)} & \text{treatment I} & \text{treatment II} & \text{treatment III} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \\ 10 \\ 11 \\ 12 \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{pmatrix} \end{matrix}$$

## Plot



$$\text{contrast}_{\text{treatment III vs control}} = c(0, 0, 0, 1)$$

# Design matrix without intercept

## Model

$$E(y) = 1.03x_0 + 2.12x_1 + 3.00x_2 + 4.90x_3$$

$$E(y) = 1.03 = 1.03 \quad (\text{for control})$$

$$E(y) = 2.12 = 2.12 \quad (\text{for treatment I})$$

$$E(y) = 3.00 = 3.00 \quad (\text{for treatment II})$$

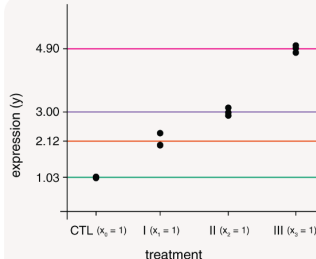
$$E(y) = 4.90 = 4.90 \quad (\text{for treatment III})$$

## Matrix

```
> model.matrix(~0 + treatment)
```

$$\begin{matrix} & \text{treatmentCTL} & \text{treatmentI} & \text{treatmentII} & \text{treatmentIII} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \\ 10 \\ 11 \\ 12 \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix} \end{matrix}$$

## Plot



$$\text{contrast}_{\text{treatmentIII vs control}} = c(-1, 0, 0, 1)$$

# Design matrix for multiple covariates

**group** factor is converted from two factors representing **tissue** samples and **cell types**

Model

$$E(y) = 1.03x_1 + 2.12x_2 + 3.00x_3 + 4.90x_4$$

$$E(y) = 1.03 = 1.03 \quad (\text{for lung B-cells})$$

$$E(y) = 2.12 = 2.12 \quad (\text{for brain B-cells})$$

$$E(y) = 3.00 = 3.00 \quad (\text{for lung T-cells})$$

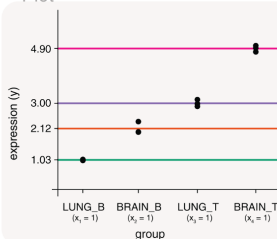
$$E(y) = 4.90 = 4.90 \quad (\text{for brain T-cells})$$

Matrix

```
> model.matrix(~0 + group)
```

$$\begin{matrix} & \text{groupLUNG\_B} & \text{groupBRAIN\_B} & \text{groupLUNG\_T} & \text{groupBRAIN\_T} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \\ 7 \\ 8 \\ 9 \\ 10 \\ 11 \\ 12 \end{matrix} & \begin{pmatrix} 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 1 \end{pmatrix} \end{matrix}$$

Plot



$$\text{contrast}_{\text{BRAIN\_B vs BRAIN\_T}} = c(0, -1, 0, 1)$$

$$\text{contrast}_{\text{LUNG vs BRAIN}} = c(0.5, -0.5, 0.5, -0.5)$$

# Design matrix for cyclic time series

Model

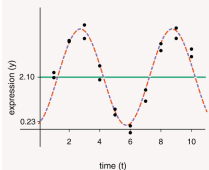
$$E(y) = 2.10 + 0.53\sin(\pi/3 t) + -1.87\cos(\pi/3 t)$$

Matrix

```
> model.matrix(~sinphase + cosphase)
```

	(Intercept)	time	time2
1	1	0.87	0.5
2	1	0.87	0.5
3	1	0.87	-0.5
4	1	0.87	-0.5
5	1	1.2e-16	-1.0
6	1	1.2e-16	-1.0
7	1	-0.87	-0.5
8	1	-0.87	-0.5
9	1	-0.87	0.5
10	1	-0.87	0.5
11	1	-2.4e-16	1.0
12	1	-2.4e-16	1.0
13	1	0.87	0.5
14	1	0.87	0.5
15	1	0.87	-0.5
16	1	0.87	-0.5
17	1	3.7e-16	-1.0
18	1	3.7e-16	-1.0
19	1	-0.87	-0.5
20	1	-0.87	-0.5

Plot



# Design matrix for cyclic time series

Model

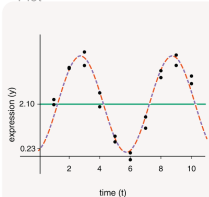
$$E(y) = 2.10 + 0.53\sin(\pi/3 t) + -1.87\cos(\pi/3 t)$$

Matrix

```
> model.matrix(~sinphase + cosphase)
```

	(Intercept)	time	time2
1	1	0.87	0.5
2	1	0.87	0.5
3	1	0.87	-0.5
4	1	0.87	-0.5
5	1	1.2e-16	-1.0
6	1	1.2e-16	-1.0
7	1	-0.87	-0.5
8	1	-0.87	-0.5
9	1	-0.87	0.5
10	1	-0.87	0.5
11	1	-2.4e-16	1.0
12	1	-2.4e-16	1.0
13	1	0.87	0.5
14	1	0.87	0.5
15	1	0.87	-0.5
16	1	0.87	-0.5
17	1	3.7e-16	-1.0
18	1	3.7e-16	-1.0
19	1	-0.87	-0.5
20	1	-0.87	-0.5

Plot



Model

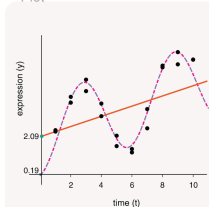
$$E(y) = 2.09 + 0.25t + 0.45\sin(\pi/3 t) + -1.90\cos(\pi/3 t)$$

Matrix

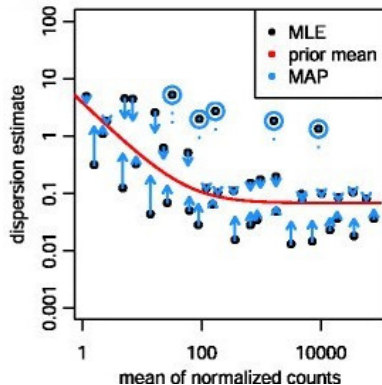
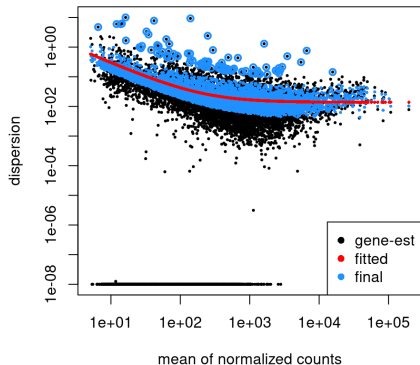
```
> model.matrix(~time+ sinphase + cosphase)
```

	(Intercept)	time	sinphase	cosphase
1	1	1	0.87	0.5
2	1	1	0.87	0.5
3	1	2	0.87	-0.5
4	1	2	0.87	-0.5
5	1	3	1.2e-16	-1.0
6	1	3	1.2e-16	-1.0
7	1	4	-0.87	-0.5
8	1	4	-0.87	-0.5
9	1	5	-0.87	0.5
10	1	5	-0.87	0.5
11	1	6	-2.4e-16	1.0
12	1	6	-2.4e-16	1.0
13	1	7	0.87	0.5
14	1	7	0.87	0.5
15	1	8	0.87	-0.5
16	1	8	0.87	-0.5
17	1	9	3.7e-16	-1.0
18	1	9	3.7e-16	-1.0
19	1	10	-0.87	-0.5
20	1	10	-0.87	-0.5

Plot



## Dispersion estimate



- Not enough replicates to estimate dispersion for individual genes
- Borrow information from genes of similar expression strength among the replicates
- Genes with very high dispersion left as is

## Implementation of DGE testing for RNA-seq

- *Seq. depth normalization*: DESeq2 uses sample-wise size factor, edgeR and Limma-Voom use TMM
- *Assumed distribution*: edgeR and DESeq model the count data using a negative binomial distribution and use their own modified statistical tests based on that. Limma-Voom uses log-normal distribution and  $t$ -test.
- *Dispersion estimate*: edgeR, DESeq2, Limma-Voom (in slightly different ways) "borrow" information across genes to get a better variance estimate.
- *Statistical test to examine if the changes are statistically significant*: DESeq2 provides the Wald test or the likelihood ratio test; edgeR uses quasi-likelihood (QL) F-test or likelihood ratio test
- *Multiple testing issue*: All current packages report false discovery rate FDR (most often Benjamini-Hochberg corrected p values).



## Multiple testing issue

- Assume that you are comparing genes between condition A and B
- You would expect 1 in 20 (5/100) genes to be significant with  $p < 0.05$  level assuming independence of tests

## Multiple testing issue

- Assume that you are comparing genes between condition A and B
- You would expect 1 in 20 (5/100) genes to be significant with  $p < 0.05$  level assuming independence of tests
- Probability of **observing** a type-I error (false positive) in a single test:  
 $\alpha_{single} = 0.05$

## Multiple testing issue

- Assume that you are comparing genes between condition A and B
- You would expect 1 in 20 (5/100) genes to be significant with  $p < 0.05$  level assuming independence of tests
- Probability of **observing** a type-I error (false positive) in a single test:  
 $\alpha_{single} = 0.05$
- Probability of **not observing** a type-I error (false positive) in a single test:  
 $\beta_{single} = 1 - \alpha = 0.95$

## Multiple testing issue

- Assume that you are comparing genes between condition A and B
- You would expect 1 in 20 (5/100) genes to be significant with  $p < 0.05$  level assuming independence of tests
- Probability of **observing** a type-I error (false positive) in a single test:  
 $\alpha_{single} = 0.05$
- Probability of **not observing** a type-I error (false positive) in a single test:  
 $\beta_{single} = 1 - \alpha = 0.95$
- If you run a test with 20 (n) genes:  
 $\beta_{multiple} = (1 - \alpha)^n = 0.95^{20} = 0.36$

## Multiple testing issue

- Assume that you are comparing genes between condition A and B
- You would expect 1 in 20 (5/100) genes to be significant with  $p < 0.05$  level assuming independence of tests
- Probability of **observing** a type-I error (false positive) in a single test:  
 $\alpha_{single} = 0.05$
- Probability of **not observing** a type-I error (false positive) in a single test:  
 $\beta_{single} = 1 - \alpha = 0.95$
- If you run a test with 20 ( $n$ ) genes:  
 $\beta_{multiple} = (1 - \alpha)^n = 0.95^{20} = 0.36$
- Likewise, type-I error for multiple comparisons become:  
 $\alpha_{multiple} = 1 - (1 - \alpha)^n = 0.64$

## Multiple testing issue

- Assume that you are comparing genes between condition A and B
- You would expect 1 in 20 (5/100) genes to be significant with  $p < 0.05$  level assuming independence of tests

- Probability of **observing** a type-I error (false positive) in a single test:

$$\alpha_{single} = 0.05$$

- Probability of **not observing** a type-I error (false positive) in a single test:

$$\beta_{single} = 1 - \alpha = 0.95$$

- If you run a test with 20 ( $n$ ) genes:

$$\beta_{multiple} = (1 - \alpha)^n = 0.95^{20} = 0.36$$

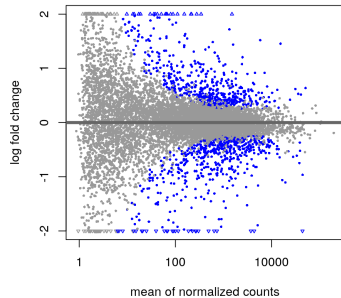
- Likewise, type-I error for multiple comparisons become:

$$\alpha_{multiple} = 1 - (1 - \alpha)^n = 0.64$$

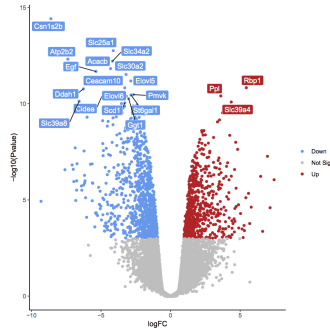
- If the number of tests ( $m$ ) increases, the type-I error rate  $\alpha_{multiple}$  will reach to 1
- This inflation of  $\alpha$  has to be handled by multiple testing correction for p-values
- Most applied method for omics studies is Benjamini-Hochberg method (a.k.a. False Discovery Rate, FDR)

# Visualization of DE analysis

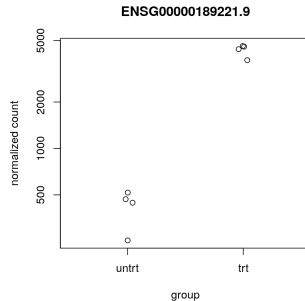
## MA-plot (Bland-Altman)



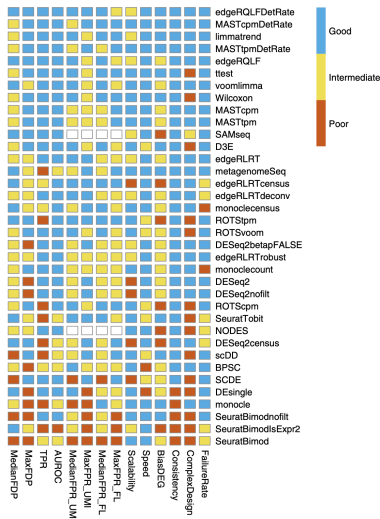
## Volcano plot



## Normalised counts



## An overview of statistical tests for DGE analysis

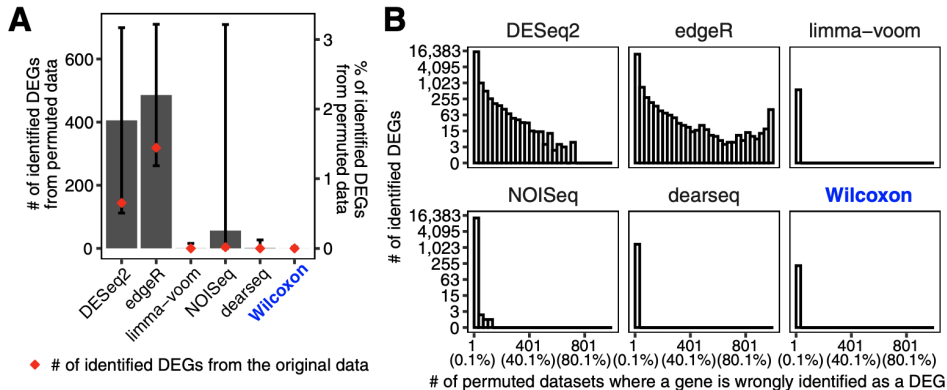


- Classical tests (t-test, Wilcoxon) still perform well
- However, they cannot handle complex designs
- Advanced methods (e.g. DESeq2, edgeR, limma) can handle complex experimental designs
- Choose your method carefully based on your needs
- If you don't know what to do, advanced methods are still the way to go



## Very large sample sizes

- population-level RNA-seq studies
- single cell RNA-seq

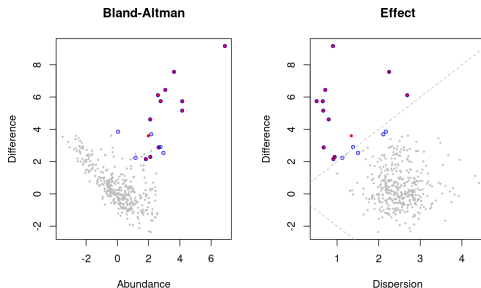


→ Non-parametric approaches, e.g., Wilcoxon rank-sum test, perform best.

## Alternative: Compositional data analysis

ALDex for differential expression analysis:

- ① Add a small prior count to the observed counts for taxa  $j$  across all samples
- ② Draw Monte Carlo samples using the Dirichlet distribution
- ③ Transform the samples using the Centered log ratio (CLR) transform
- ④ Hypothesis testing, e.g., Welch's t-test or Wilcoxon rank test
- ⑤ Report expected values from statistical tests and effect-size estimate



## Summary

- Always challenge your data, think of plausible technical explanation first

*"I'm a scientist and I know what constitutes proof. But the reason I call myself by my childhood name is to remind myself that a scientist must also be absolutely like a child. If [they] see a thing, [they] must say that [they] see it, whether it was what [they] thought [they] were going to see or not. See first, think later, then test. But always see first. Otherwise you will only see what you were expecting. Most scientists forget that."*

– adapted from The Ultimate Hitchhiker's Guide to the Galaxy by Douglas Adams

# PhD course "Bioinformatics analysis of gene expression data (BAGED)"

- ① Do you have your own bulk or single-cell RNA-seq data of a vertebrate?
- ② Do you want to learn how to analyze your data yourself?
- ③ What? 2 to 3 weeks of lectures, tutorials and most importantly student projects
- ④ Hardware and Software? UCloud, (Galaxy), R, Cytoscape
- ⑤ When? BAGED-bulk January 2026; BAGED-single Autumn/Spring 2025/26
- ⑥ Sign up early due to limited number of seats