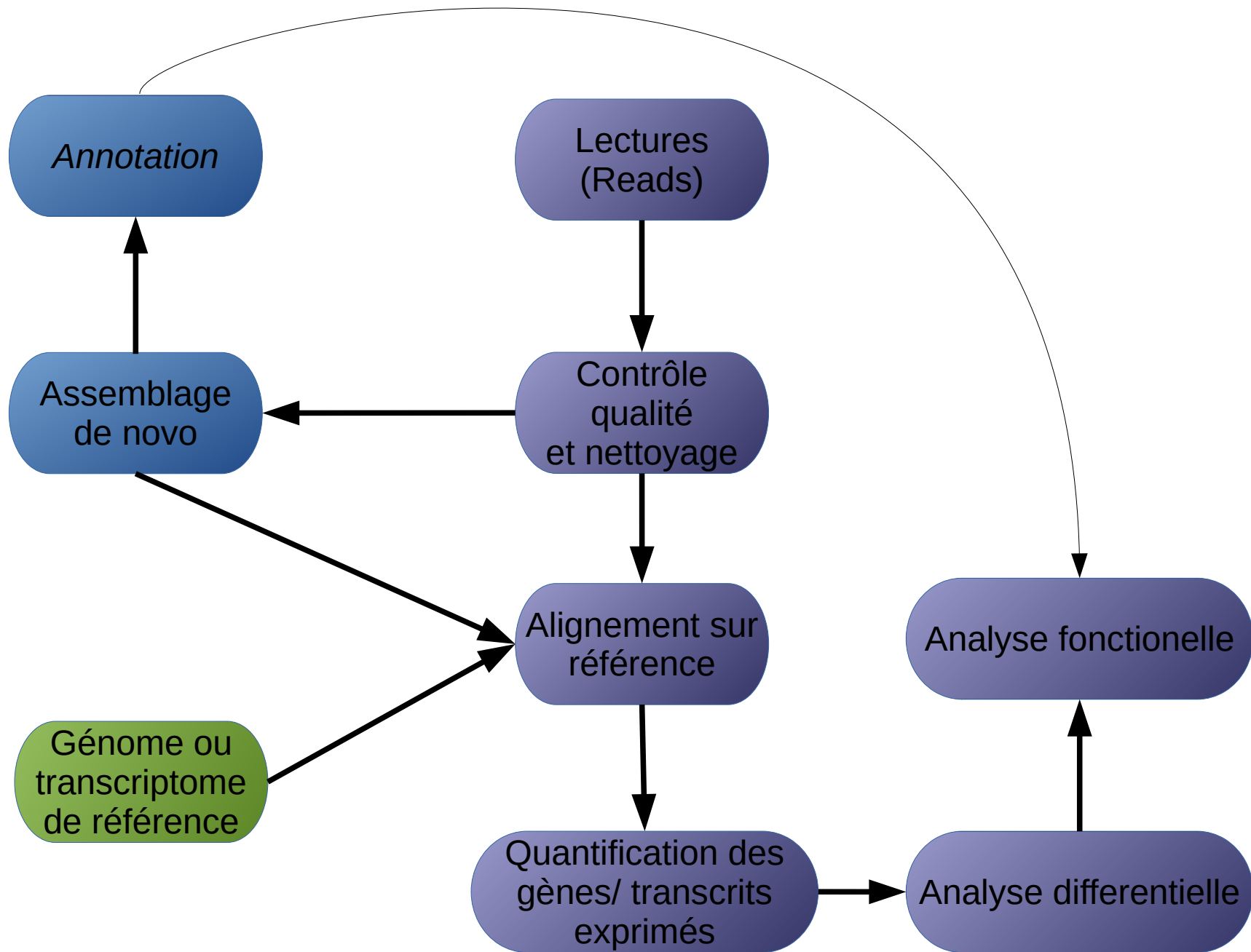




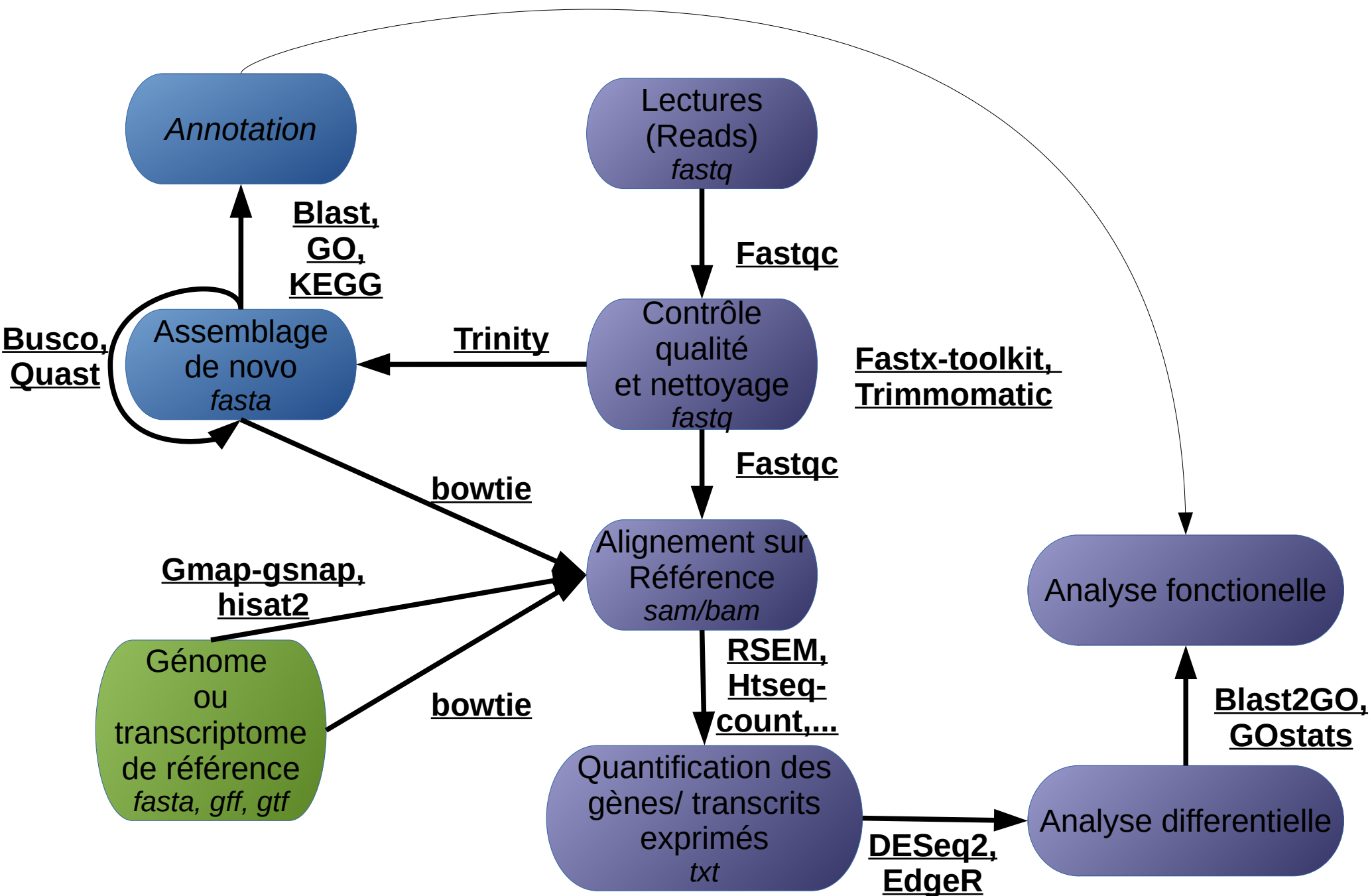
# TRAITEMENT BIOINFORMATIQUE DE DONNÉES RNA-Seq

Anamaria NECSULEA & Adil EL FILALI

# Workflow d'analyse RNA-Seq



# Workflow d'analyse RNA-Seq



## Data preprocessing



**Experimental design**  
19 tools



**Read quality control**  
112 tools



**Adapter trimming**  
108 tools



**Depth of coverage**  
22 tools



**Demultiplexing**  
28 tools



**Base calling**  
33 tools



**Error correction**  
114 tools



**k-mer counting**  
36 tools



**Read clustering**  
24 tools



**Barcode removal**  
4 tools

## Data processing



**Spliced read alignment**  
98 tools



**Read realignment**  
3 tools



**Alignment evaluation**  
33 tools



**Reference-based transcriptome assembly**  
46 tools



**Base quality recalibration**  
3 tools



**Read alignment**  
251 tools



**Indel realignment**  
4 tools



**De novo transcriptome assembly**  
55 tools



**Assembly evaluation**  
10 tools



**Duplicate read removal**  
42 tools

## Quantification



**Read count**  
63 tools



**Novel transcript quantification**  
46 tools



**Batch effect correction**  
13 tools



**Known transcript quantification**  
135 tools



**Alignment-free transcript quantification**  
12 tools



**Normalization**  
174 tools

## Data visualization



**Alternative splicing visualization**  
29 tools



**Transcriptome visualization**  
21 tools



**Circular RNA visualization**  
1 tool



**Read alignment visualization**  
46 tools



**Venn diagram creation**  
29 tools



**Gene visualization**  
21 tools



**Tree visualization**  
12 tools



**Gene expression visualization**  
67 tools



**Gene fusion visualization**  
10 tools



**Genome visualization**  
162 tools



**Heatmap generation**  
36 tools



**Variant visualization**  
44 tools



**Sequence alignment visualization**  
55 tools



**Network visualization**  
210 tools

## Classification/Clustering



**Gene expression clustering**  
174 tools



**Gene expression classification**  
64 tools

## Network analysis



**Gene co-expression prediction**  
74 tools



**Differential co-expression analysis**  
16 tools



**miRNA regulatory modules inference**  
23 tools



**Gene regulatory network inference**  
325 tools

## Enrichment analysis



**RBP motif enrichment**  
2 tools



**Over-representation analysis**  
48 tools

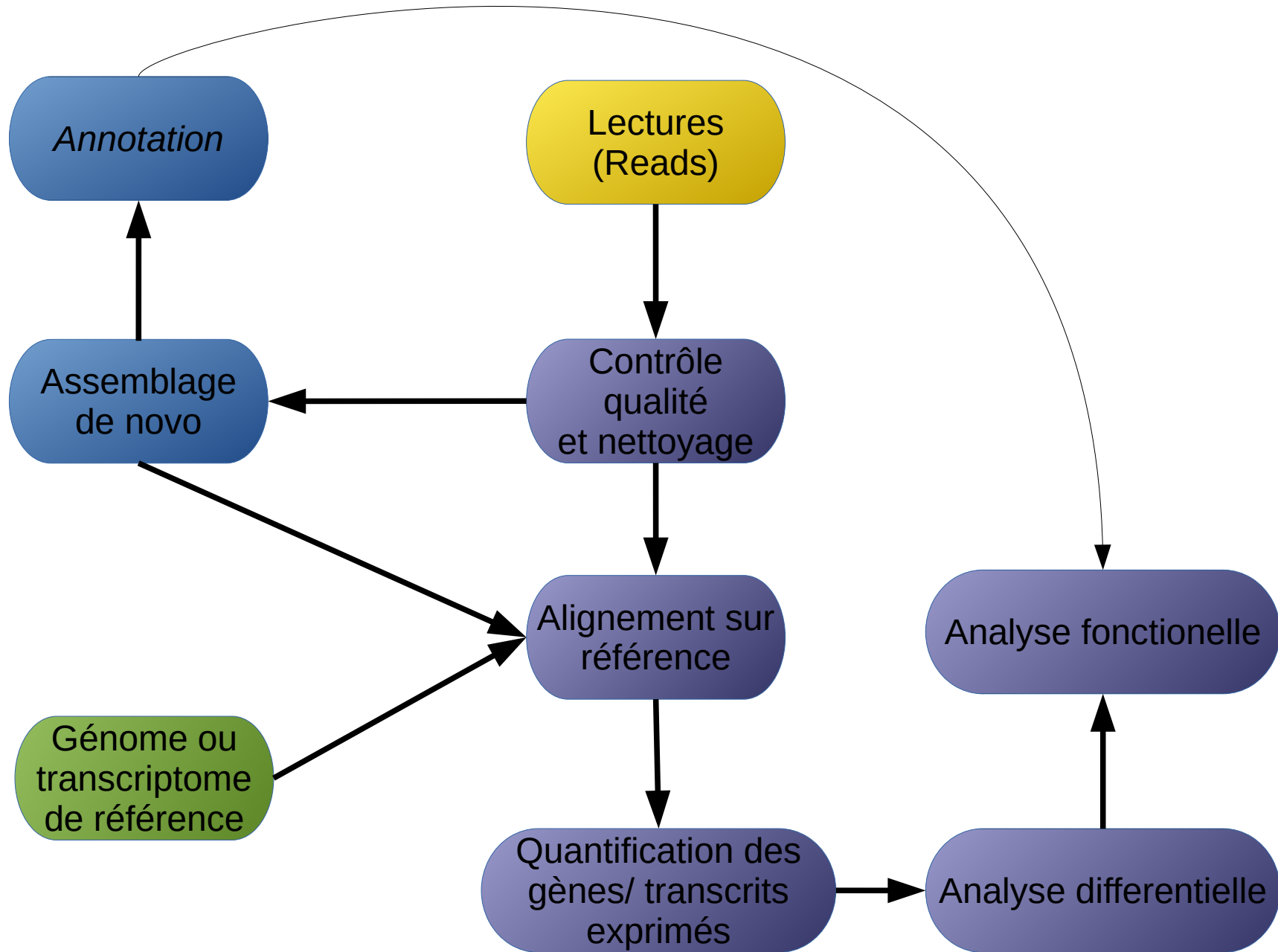


**Gene set enrichment analysis**  
205 tools



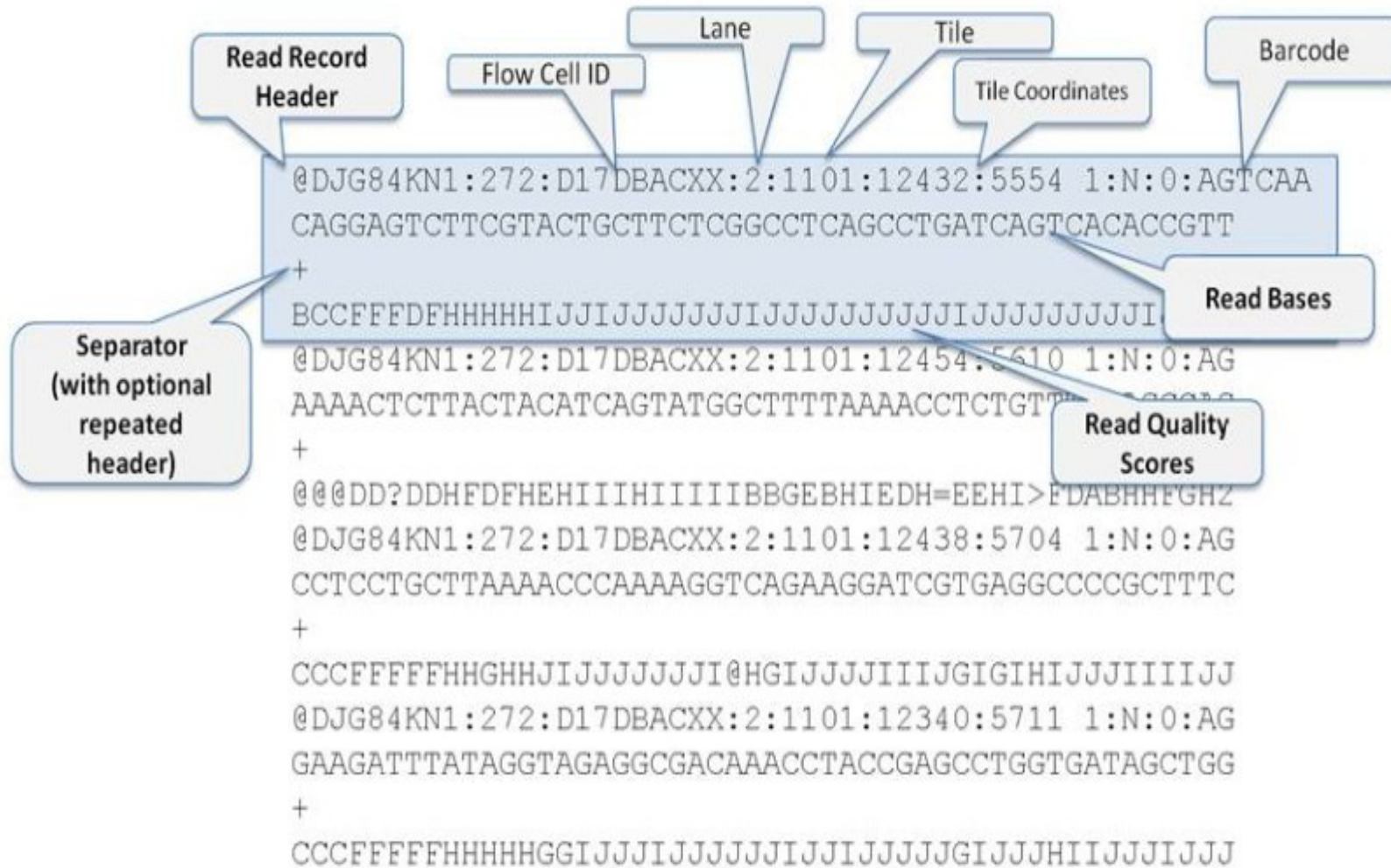
**Topology enrichment analysis**  
70 tools

# Workflow d'analyse RNA-Seq



# Lectures (*fastq*)

2 FASTQ files in case of a paired-end experiment





# Lectures (*fastq*)

Score de qualité = Phred Score

- Pour chaque base, mesure la probabilité que la base assignée soit fausse
- Le score d'une base donné (Q) est donné par l'équation suivante :

$$Q = -10\log_{10}(p)$$

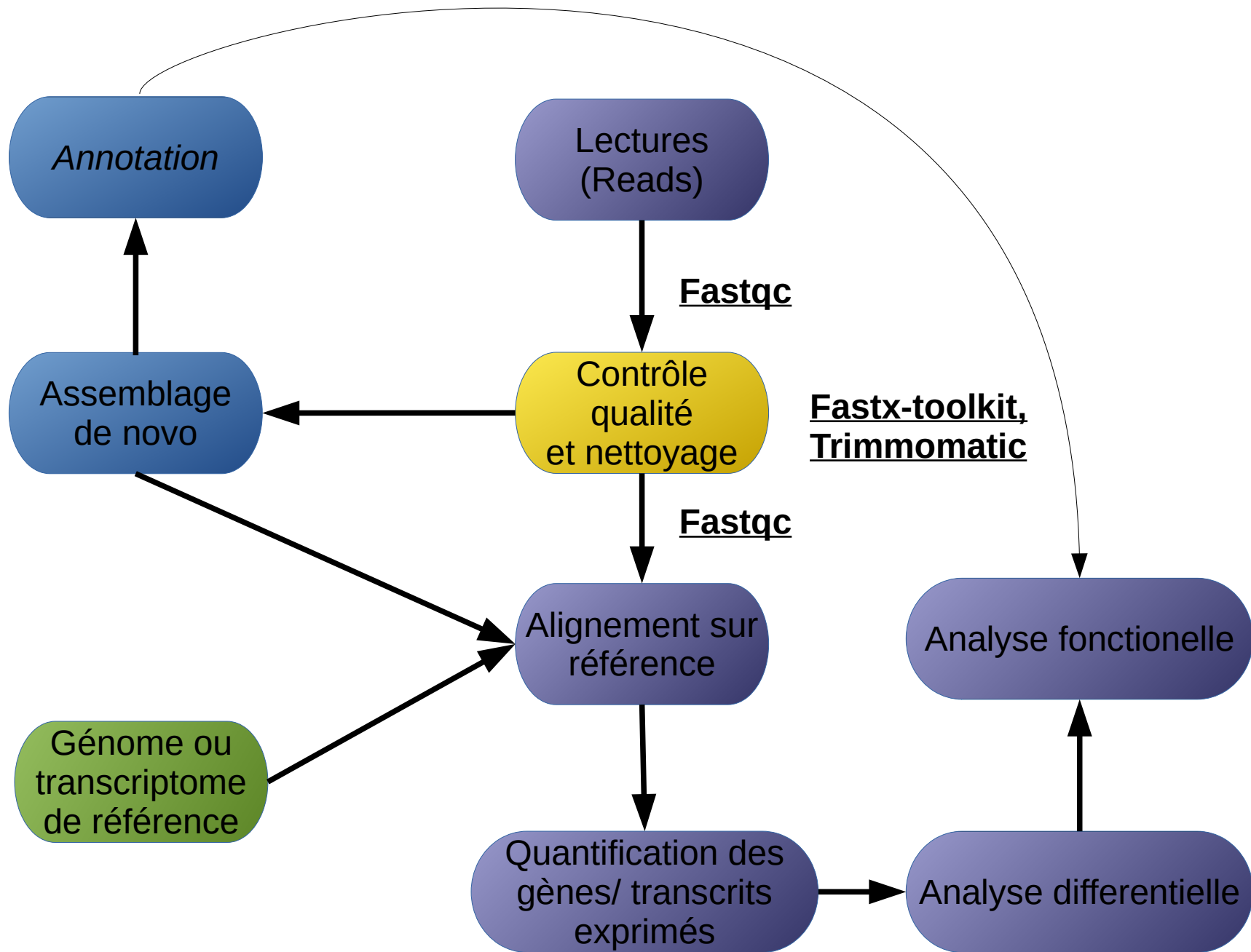
p: Probabilité estimée que la base donnée soit fausse

Donc un haut score indique une plus petite probabilité d'erreur.

Base Quality	P <sub>error</sub> (obs. base)
3	50 %
5	32 %
10	10 %
20	1 %
30	0.1 %
40	0.01 %













# Workflow d'analyse RNA-Seq



<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

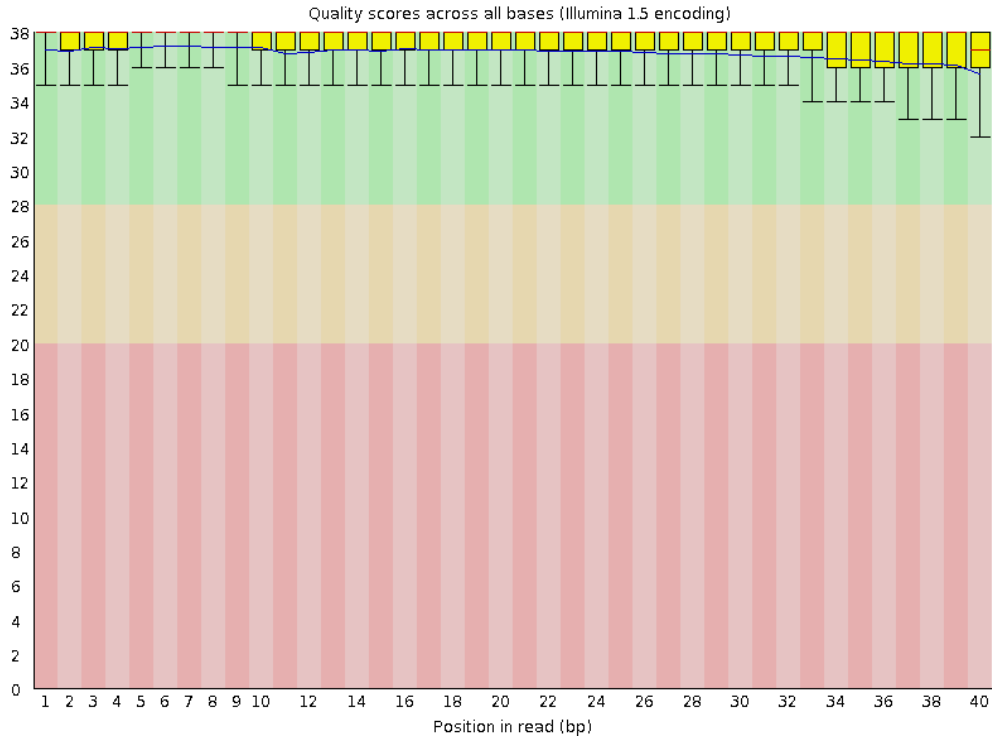
Ce n'est pas un outil de correction mais un outil de visualisation de la qualité des données génomique.

## Summary

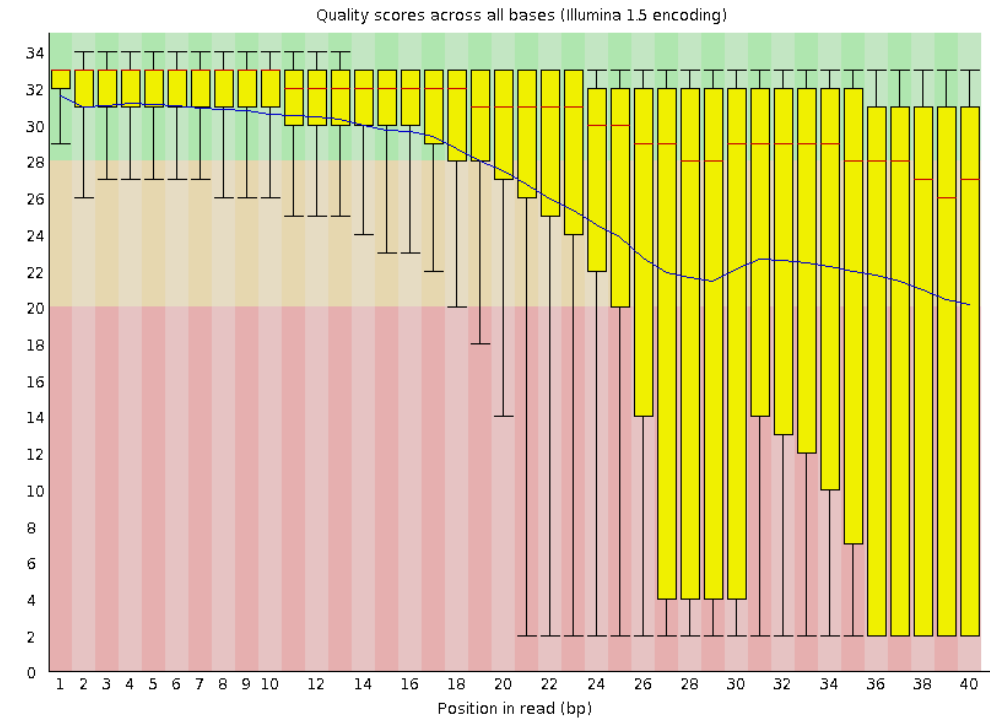
-  [Basic Statistics](#)
-  [Per base sequence quality](#)
-  [Per sequence quality scores](#)
-  [Per base sequence content](#)
-  [Per sequence GC content](#)
-  [Per base N content](#)
-  [Sequence Length Distribution](#)
-  [Sequence Duplication Levels](#)
-  [Overrepresented sequences](#)
-  [Adapter Content](#)

Measure	Value
Filename	454_SRR073599.fastq.gz
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	125000
Sequences flagged as poor quality	0
Sequence length	44-2042
%GC	43

## ✔ Per base sequence quality

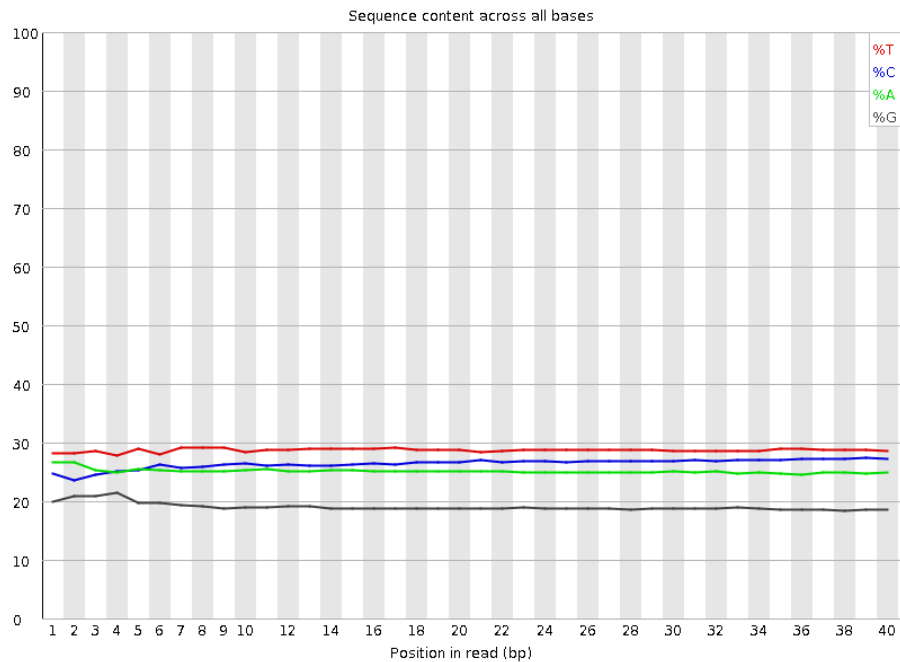


## ✘ Per base sequence quality

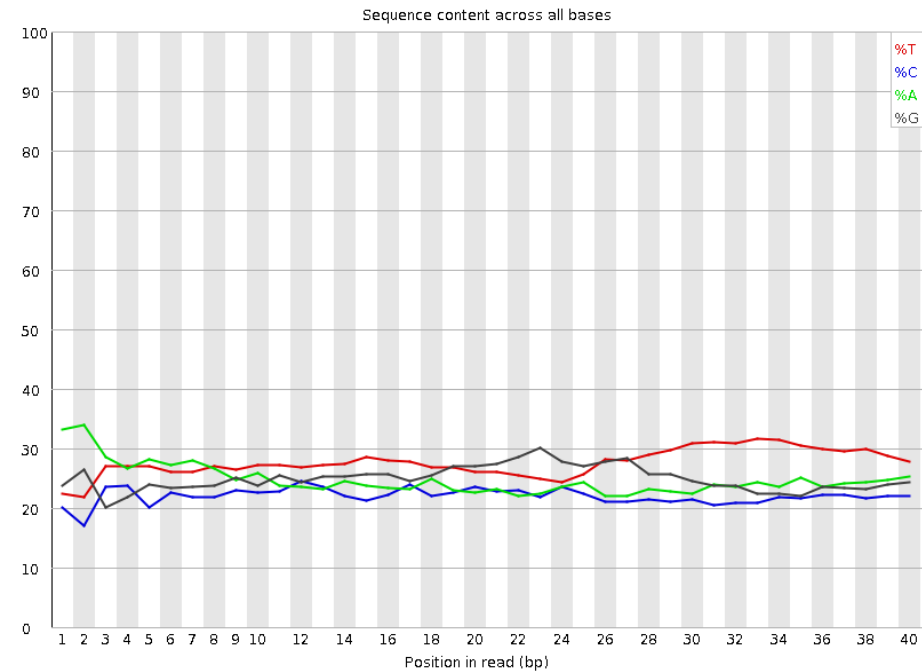


- La ligne rouge centrale correspond a la médiane
- La boxplot jaune represente l'inter-quartile (25-75%)
- Les lignes inférieures et supérieures representent la limite de 10 % et 90 %
- La courbe bleue correspond a la moyenne de la qualité

## ✔ Per base sequence content



## ⚠ Per base sequence content



- Indication de la proportion en base (A/T/C/G) par position
- Dans un jeu de données aléatoire, peu de différence entre bases —> lignes quasi-paralleles
- Dans un jeu de données biaisé, différence dans la proportion des bases



Sur-representation d'une séquence contaminant le jeu de donnée

# Biais : random hexamer priming

- ❖ Fort biais de composition des 13 premières nucléotides en 5'
- spécificité de séquence de la polymérase

Published online 14 April 2010

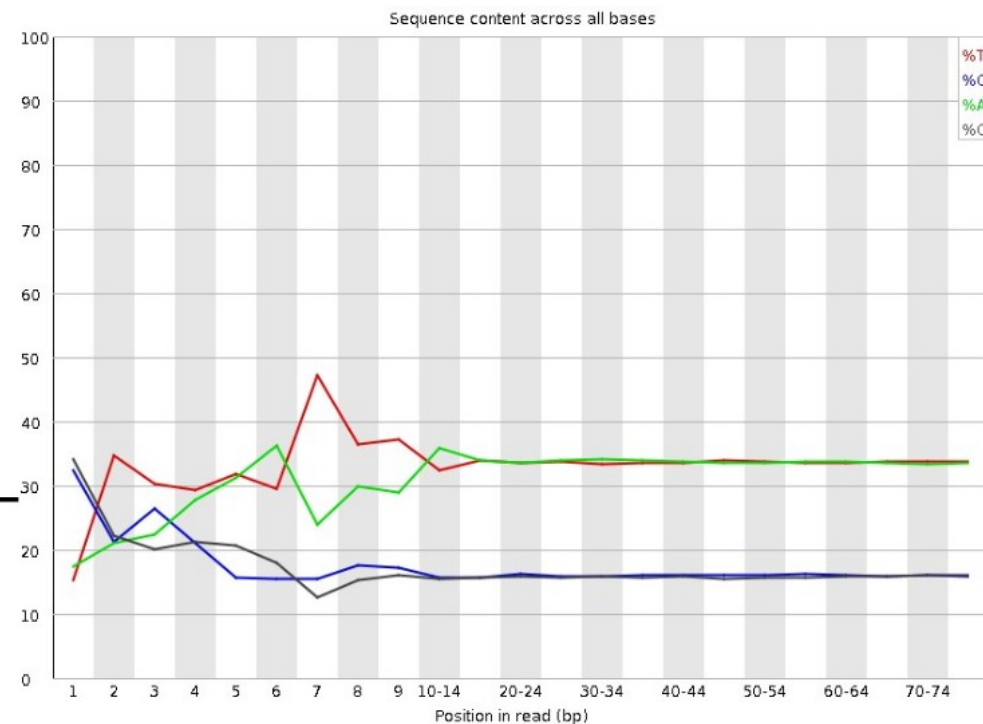
Nucleic Acids Research, 2010, Vol. 38, No. 12 e131  
doi:10.1093/nar/gkq224

## Biases in Illumina transcriptome sequencing caused by random hexamer priming

Kasper D. Hansen<sup>1,\*</sup>, Steven E. Brenner<sup>2</sup> and Sandrine Dudoit<sup>1,3</sup>

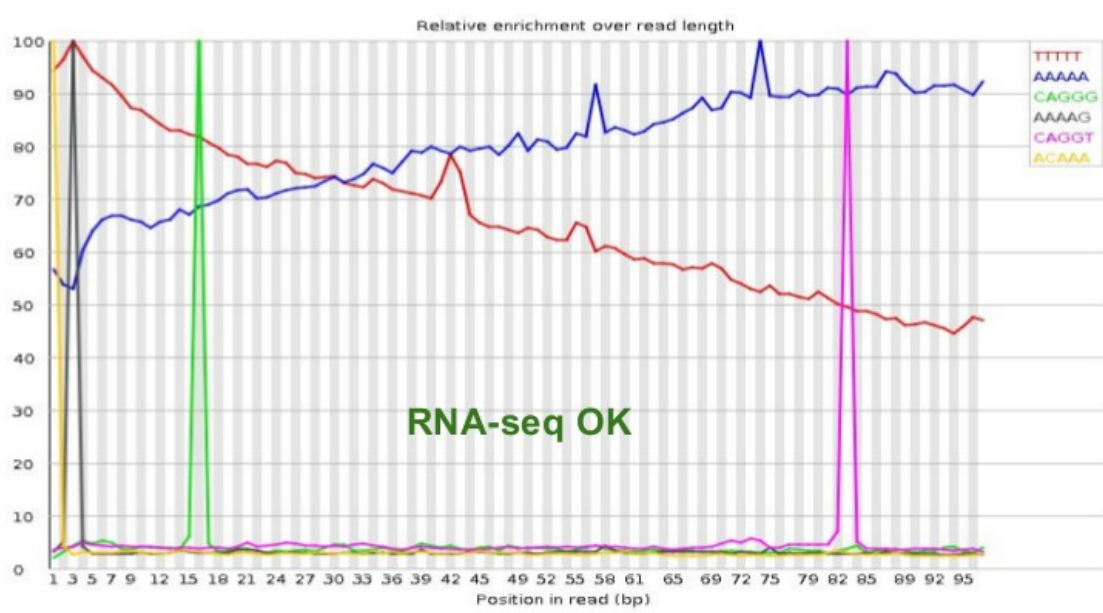
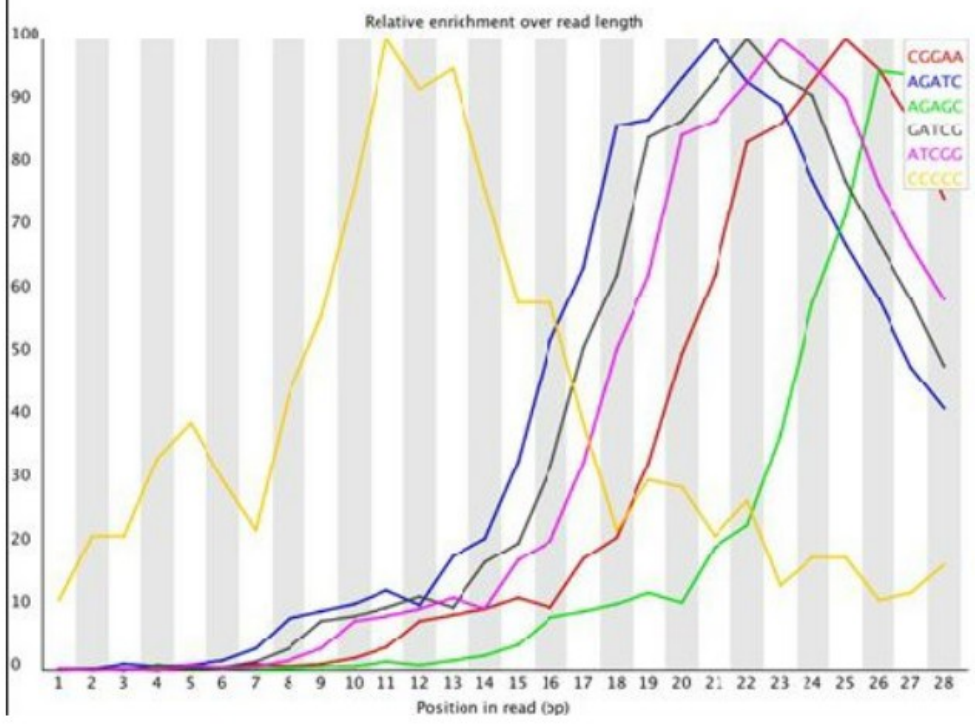
### ABSTRACT

Generation of cDNA using random hexamer priming induces biases in the nucleotide composition at the beginning of transcriptome sequencing reads from the Illumina Genome Analyzer. The bias is independent of organism and laboratory and impacts the uniformity of the reads along the transcriptome. We provide a read count reweighting scheme, based on the nucleotide frequencies of the reads, that mitigates the impact of the bias.



## Surpresentation des k-mers

- Permet de repérer les séquences surreprésentées, donner une bonne impression de toute contamination
- vérifier la présence d'adaptateur ou de séquence poly-A (ou poly-T)



Sequence	Count	Obs/Exp Overall	Obs/Exp Max	Max Obs/Exp Position
TTTTT	47499960	4.84021	7.2762637	3
AAAAA	18101385	4.2297845	5.3006034	74
CAGGG	12486915	2.3769662	49.03375	16
AAAAG	10728075	2.3667703	56.233307	3

# Nettoyage

## But :

Retirer les adaptateurs

Éliminer les reads de faibles qualités

## Software :

BBTools - <https://jgi.doe.gov/data-and-tools/bbtools/>

Bignorm - <https://git.informatik.uni-kiel.de/axw/Bignorm>

Centrifuge - <https://github.com/DaehwanKimLab/centrifuge>

cutadapt - <https://github.com/marcelm/cutadapt>

Falco - <https://github.com/smithlabcode/falco>

fastp - <https://github.com/OpenGene/fastp>

FastQC - <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

khmer - <https://github.com/dib-lab/khmer>

Kraken2 - <https://github.com/DerrickWood/kraken2>

MultiQC - <https://multiqc.info>

rCorrector - <https://github.com/mourisl/Rcorrector>

SortMeRNA - <https://github.com/biocore/sortmerna>

TrimGalore -

[https://www.bioinformatics.babraham.ac.uk/projects/trim\\_galore/](https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)

Trimmomatic - <https://github.com/usadellab/Trimmomatic>

# Nettoyage

**Trimmomatic** est un programme qui effectue une variété de trimming et de filtrage sur les reads (en single-end ou en paired-end), Il requière Java.

<http://www.usadellab.org/cms/?page=trimmomatic>

## Tool: Trimmomatic



Source : Erwan Core, 5ème Ecole de bioinformatique AVIESAN-IFB 2016 <https://www.france-bioinformatique.fr/fr/evenements/EBA2016>

Bolger, A. M. and Lohse, M. and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina sequence data. In *Bioinformatics*, 30 (15), pp. 2114–2120



# Nettoyage

## Cutadapt

- développé en Python
- semi-global alignment algorithm
- détecte et supprime les séquences adaptatrices, les amorces, les queues poly-A et d'autres types de séquences

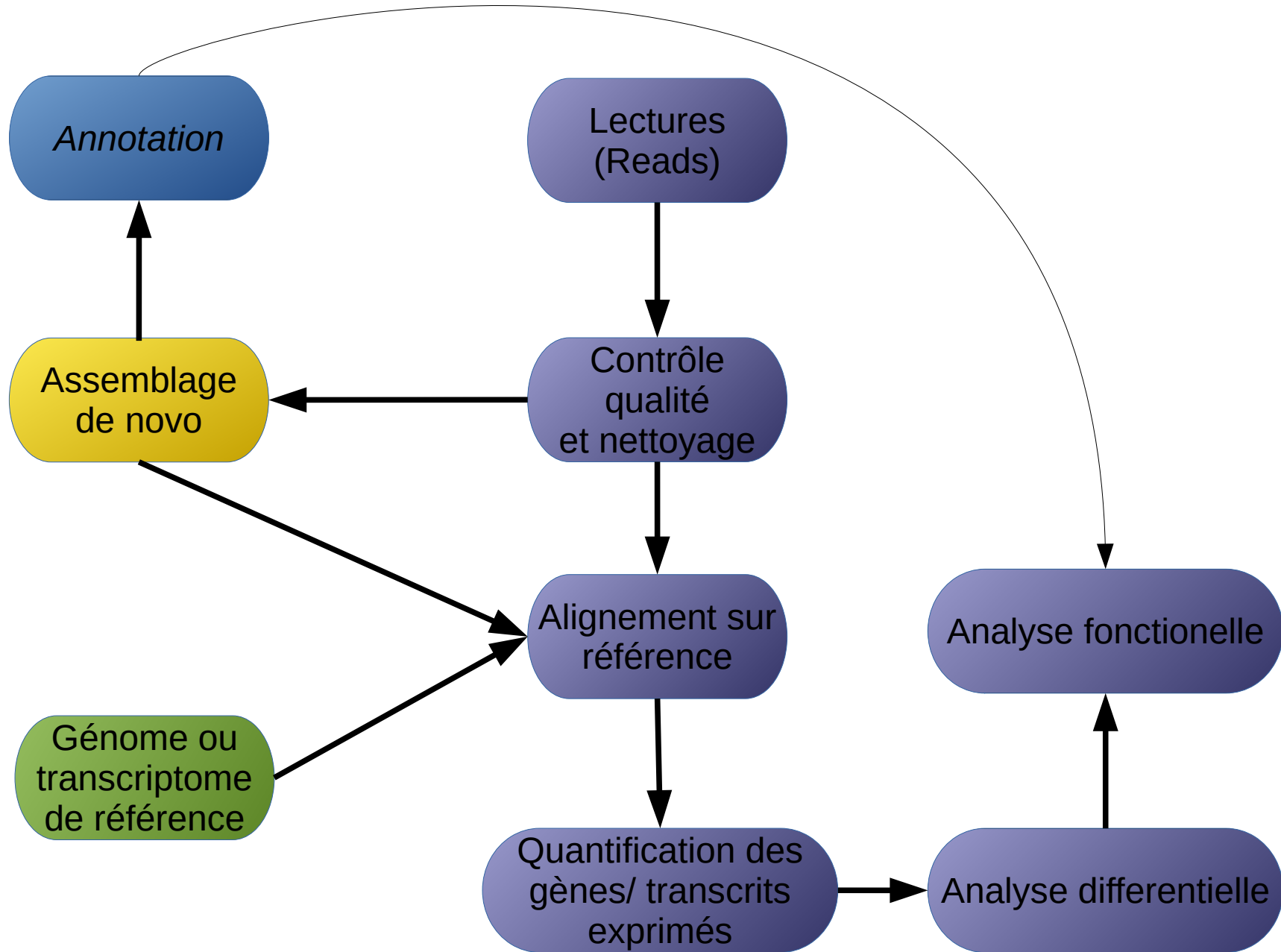
```
$ cutadapt -a ADAPTER-SEQUENCE input.fastq > output.fastq
```

```
$ cutadapt -b TGAGACACGCA -b AGGCACACAGGG input.fastq > output.fastq
```

# Définition d'une bonne qualité pour les reads

- Le nombre de reads produites correspondant au nombre attendu
- Pas de contaminations (K-mer de polyA, polyT, adaptateurs ...)
- Longueur des reads corrects (100-150 pb)
- Bonne qualité
- Bon alignement (re-séquencage avec génome de référence « propre ») : peu de reads non-alignés

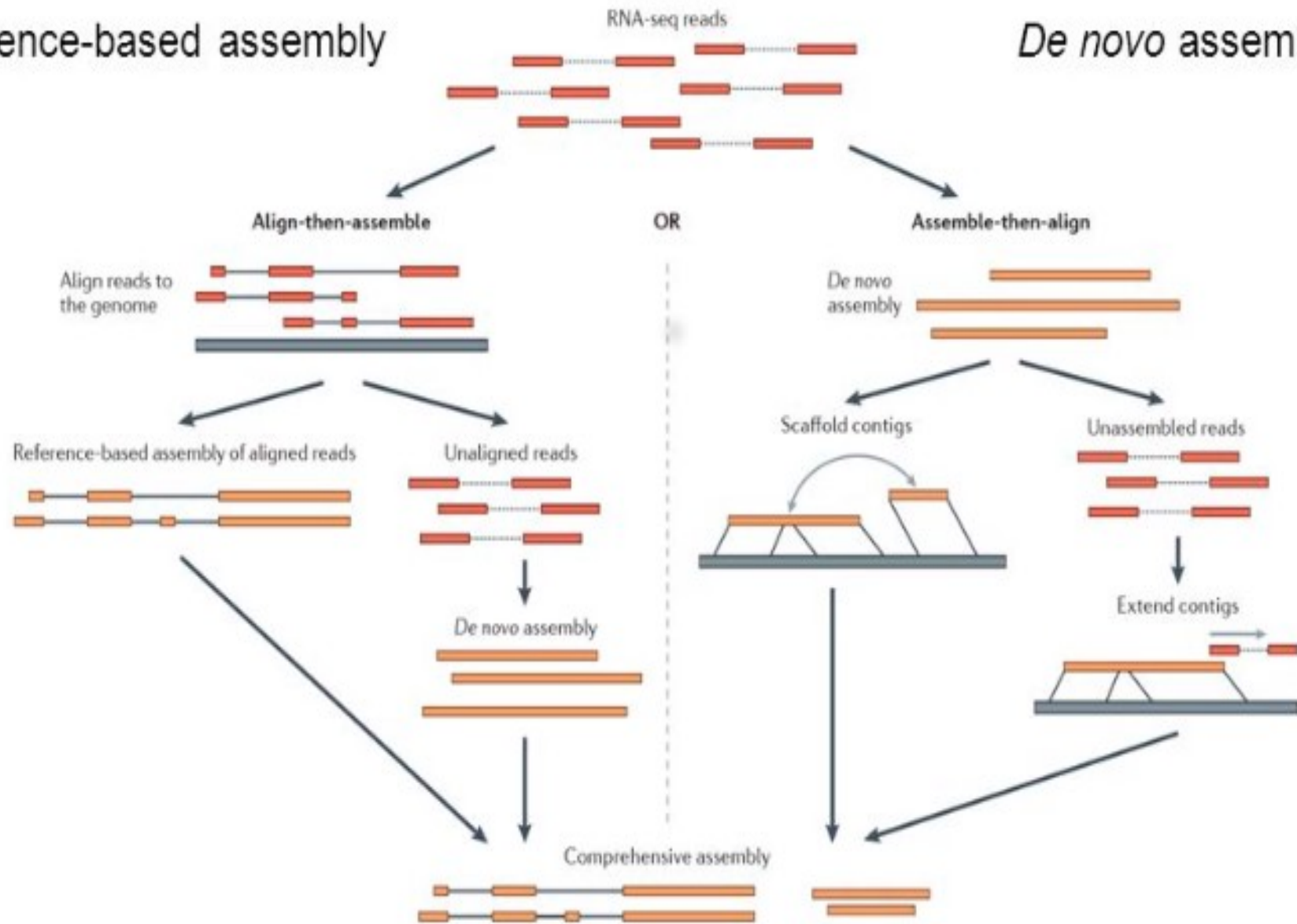
# Workflow d'analyse RNA-Seq



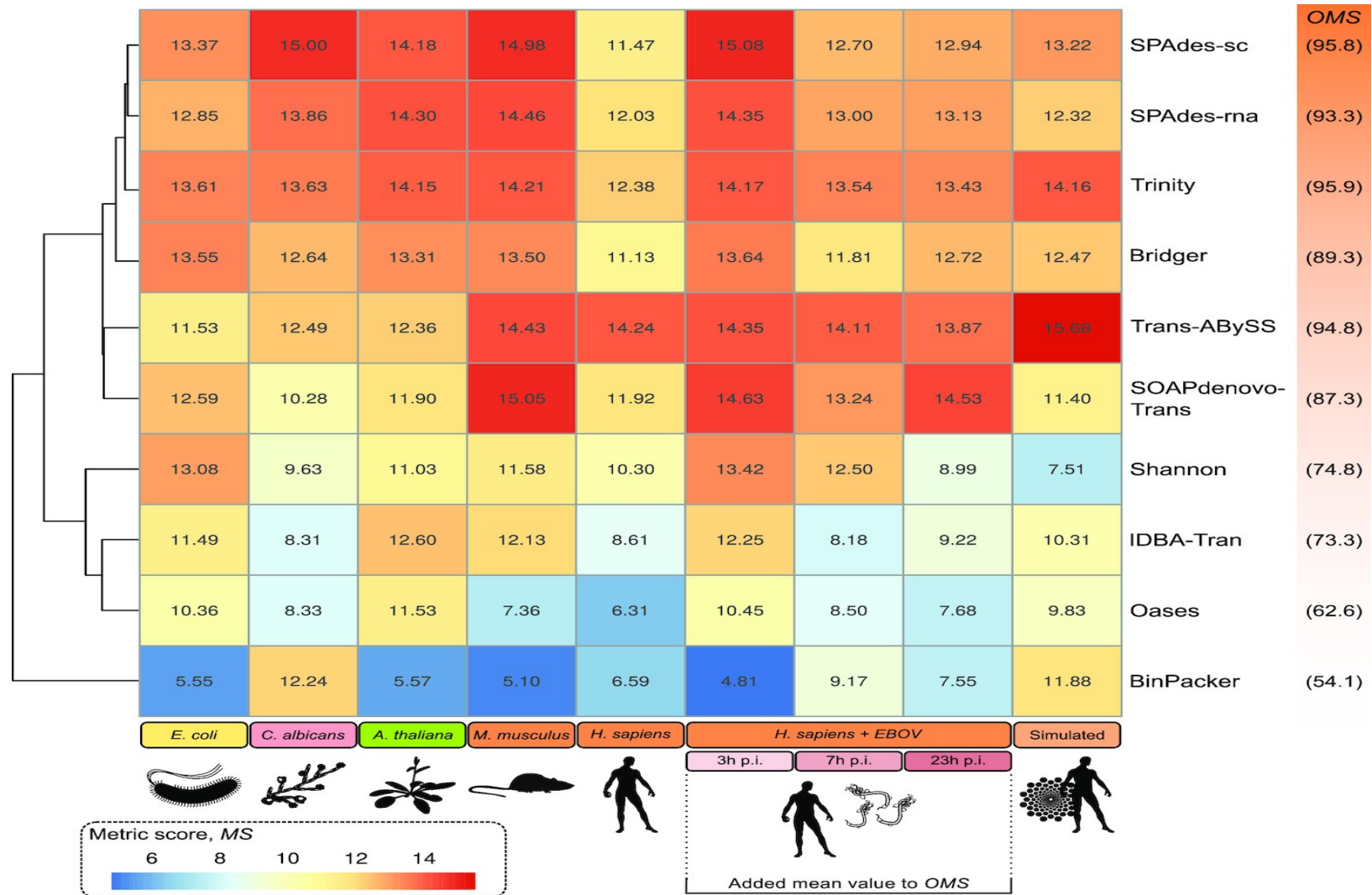
# Assemblage

Reference-based assembly

*De novo* assembly



# Assemblage de-novo



De novo transcriptome assembly: A comprehensive cross-species comparison of short-read RNA-Seq assemblers. Hölzer M, Marz M. Gigascience. 2019 May 1;8(5)

# Assemblage de-novo Trinity

## - Inchworm :

- coupure des reads en k-mer (Jellyfish)
- reconstruction des isoformes a partir des k-mer (gourmand en mémoire vive)
- resultat : ensemble des contigs



## - Chrysalis :

- clusterise les contigs se chevauchant sur une longueur de k-1
- construction d'un graphe de Bruijn pour chaque cluster



## - Butterfly :

- traitement des graphes individuellement
- reconciliation des graphes avec les reads (pairés ou non)
- resultat : gènes et isoformes



# Assemblage de-novo Trinity

## Désavantages :

- basée sur une heuristique
- surestimation des transcrits (DRAP)
- gourmand en mémoire vive

## Avantages :

- intègre toute les étapes bioinformatiques, du traitement des données brutes RNA-seq (trimmomatic, quantification, normalisation, analyse différentielle, annotation)
- choix du k-mer
- nombreux paramètres (sur toutes les étapes de l'assemblage)

# Evaluation de l'assemblage de novo

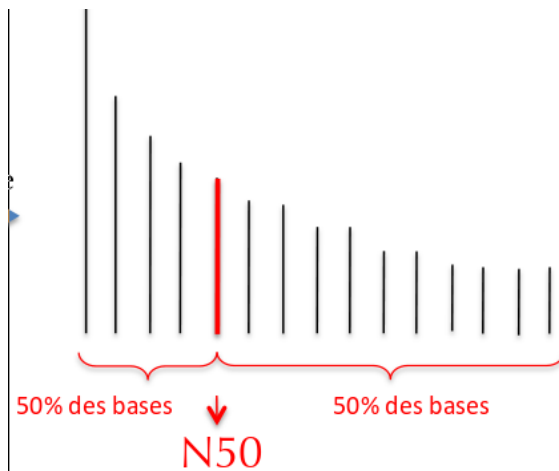
Tool	metrics
HISAT2	Overall mapping rate
rnaQUAST	Transcripts $\geq 1,000$ nt, Misassemblies, Mismatches per transcript, Average alignment length, 95%-assembled isoforms, Duplication ratio
Trinity/Salmon	Ex90N50
Blastx	Full-length transcripts
TransRate	Reference coverage, Mean ORF percentage, Optimal score, Percentage bases uncovered, Number of ambiguous bases
DETONATE	Nucleotide F1, Contig F1, KC score, RSEM-EVAL
BUSCO	Complete BUSCOs, Missing BUSCOs



# Evaluation de l'assemblage de novo : Quast

<http://bioinf.spbau.ru/quast>

- Nombre de contigs
- Nombre des longs contigs (i.e. > 1000 bp)
- Longueur du plus long contigs
- Longueur total de l'assemblage
- N50 : est la taille du scaffold (ou contig) tel que 50% des bases de l'assemblage sont comprises dans des scaffolds de taille supérieures à cette taille
- L50 : Le nombre minimum X tel que X contigs les plus longs couvrent au moins 50% de l'assemblage



All statistics are based on contigs of size  $\geq 500$  bp, unless otherwise noted

## Short report

Basic statistics	L_uniseq_100	R_uniseq_100
# contigs	60 973	48 220
# contigs ( $\geq 0$ bp)	135 904	115 018
# contigs ( $\geq 1000$ bp)	10 517	10 586
Largest contig	5384	10 707
Total length	47 553 100	40 485 523
Total length ( $\geq 0$ bp)	72 636 910	63 616 017
Total length ( $\geq 1000$ bp)	14 543 266	15 807 198
N50	749	827
N75	601	617
L50	21 334	15 476
L75	39 258	29 910

## Misassemblies

### Unaligned

#### Genome statistics

GC (%)	41.04	41.44
# N's	8528	3280
# N's per 100 kbp	17.93	8.1

#### Aligned statistics

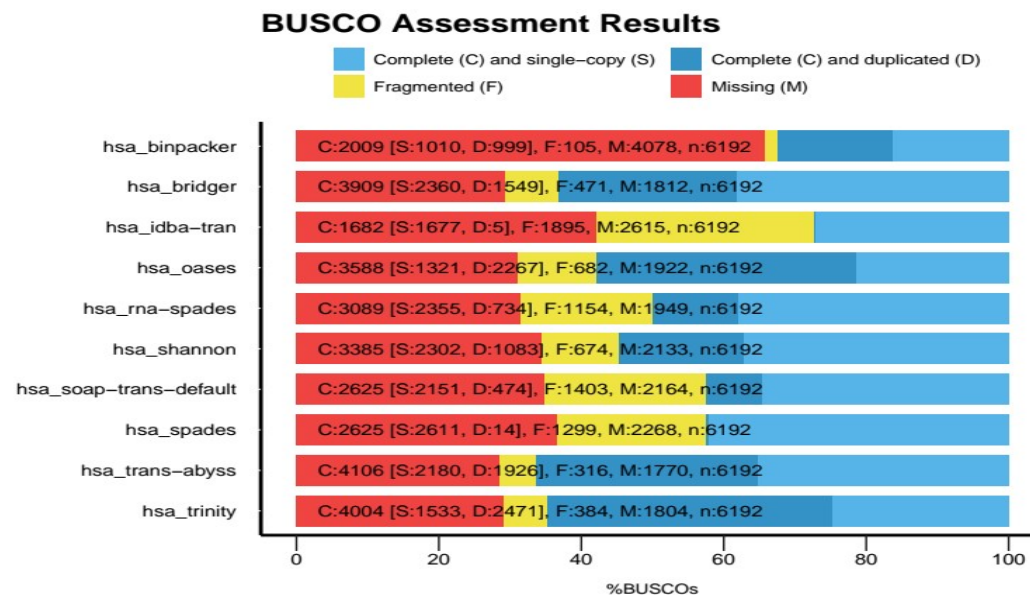
# Evaluation de l'assemblage de novo : BUSCO

BUSCO : Benchmarking Universal Single-Copy Orthologs

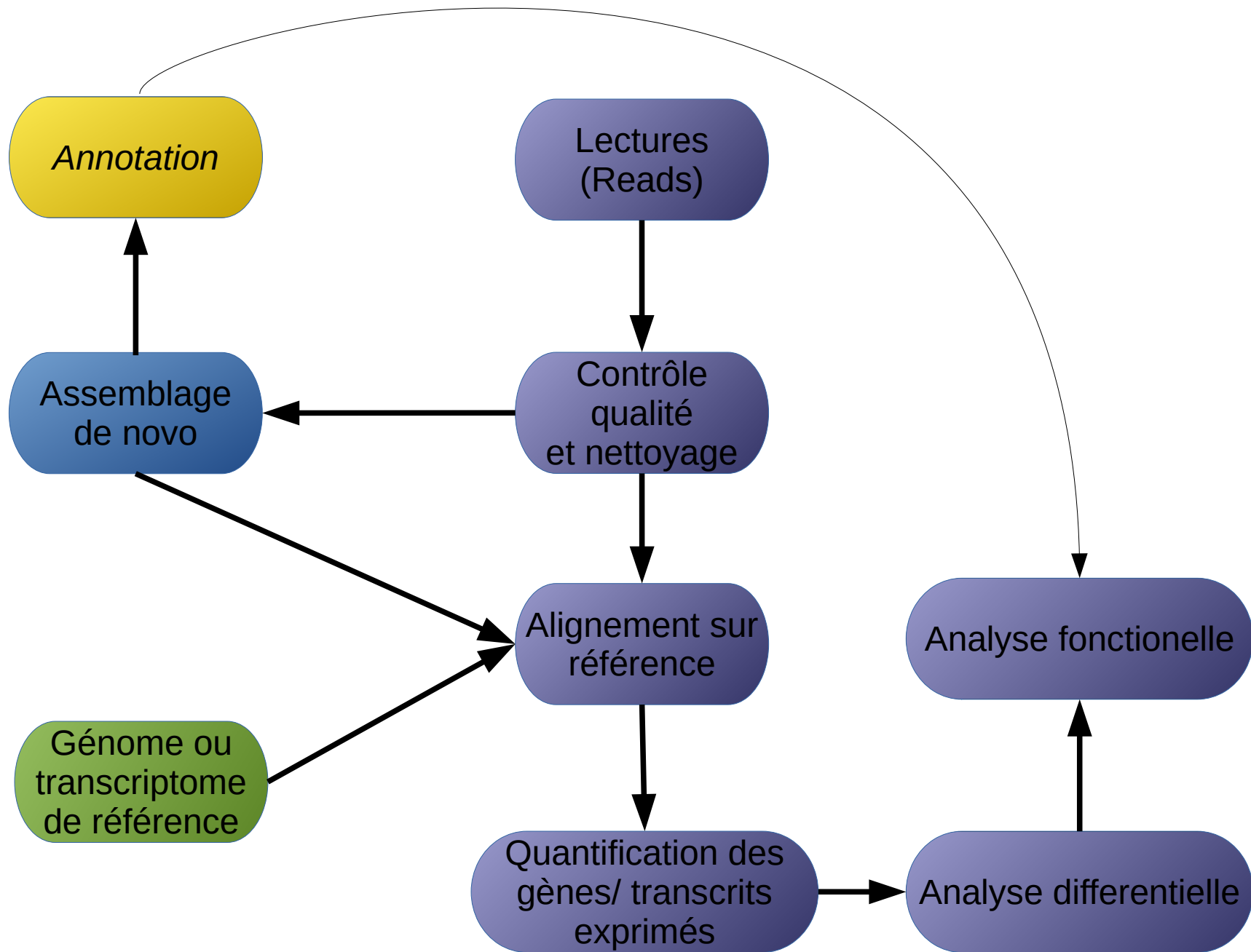
- évaluation de l'assemblage des génomes/transcriptomes en comparant les contigs à des séquences conservées simplescopies et universels

- Le processus général de BUSCO est le suivant:

- tTBLASTn contre des séquences consensus BUSCO.
- Préviation de la structure des gènes en utilisant Augustus avec les profils de blocs BUSCO.
- Assignment : "complete", "duplicated", "fragmented", ou "missing" (si il n y a pas de match).



# Workflow d'analyse RNA-Seq



# Annotation fonctionnelle

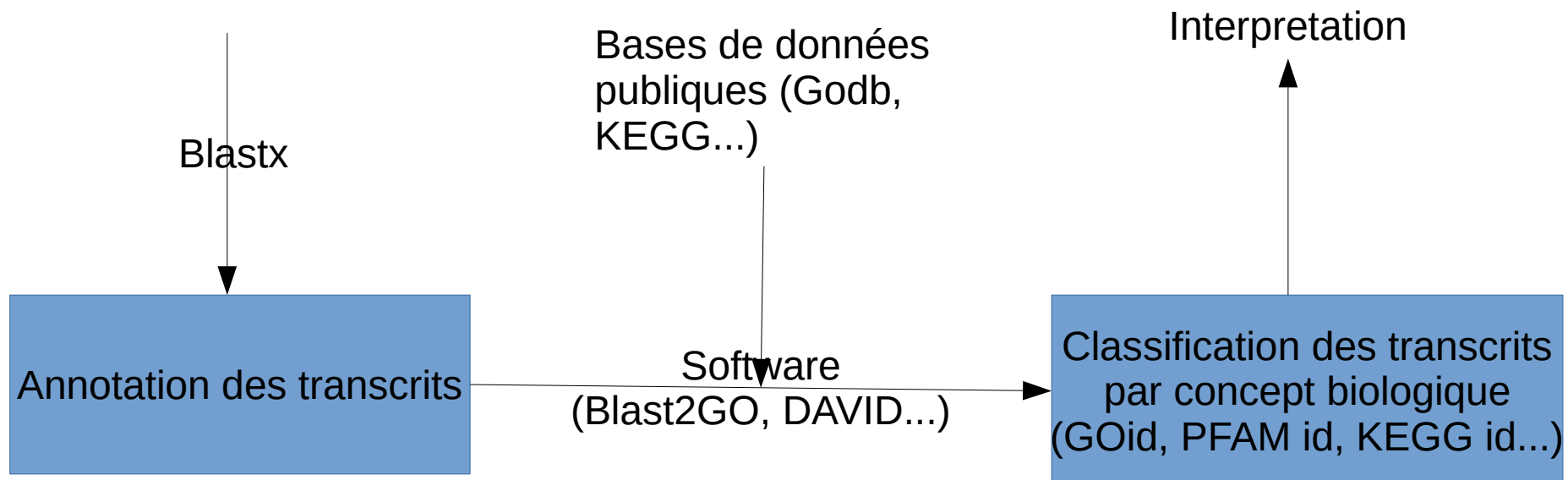
- Annotation fonctionnelle : attribuer une fonction potentielle à une séquence
- Comment : comparaison à des séquences, dont la fonction biologique est connue ou supposée, disponibles dans les bases de données (ex : nr, nt, swissprot...)
- Programmes :
  - blastn : comparaison à une banque nucléotidique (nt)
  - blastx : comparaison à une banque protéique (nr), traduction de la séquence dans les 6 cadre de lectures et comparaison à la banque

# Classification fonctionnelle

Classification fonctionnelle : attribuer à chaque séquence unique une classe de fonction

L'intérêt est de pouvoir cibler certains gènes dans une classe de fonction particulière. Elle vient en complément de l'annotation.

Comment : utilisation de base de données qui relie à un gène donné d'une espèce donnée une (ou des) classes(s) de fonction particulier(s).



# Classification fonctionnelle



Pfam



diamond

PSICQUIC

PANNZER2

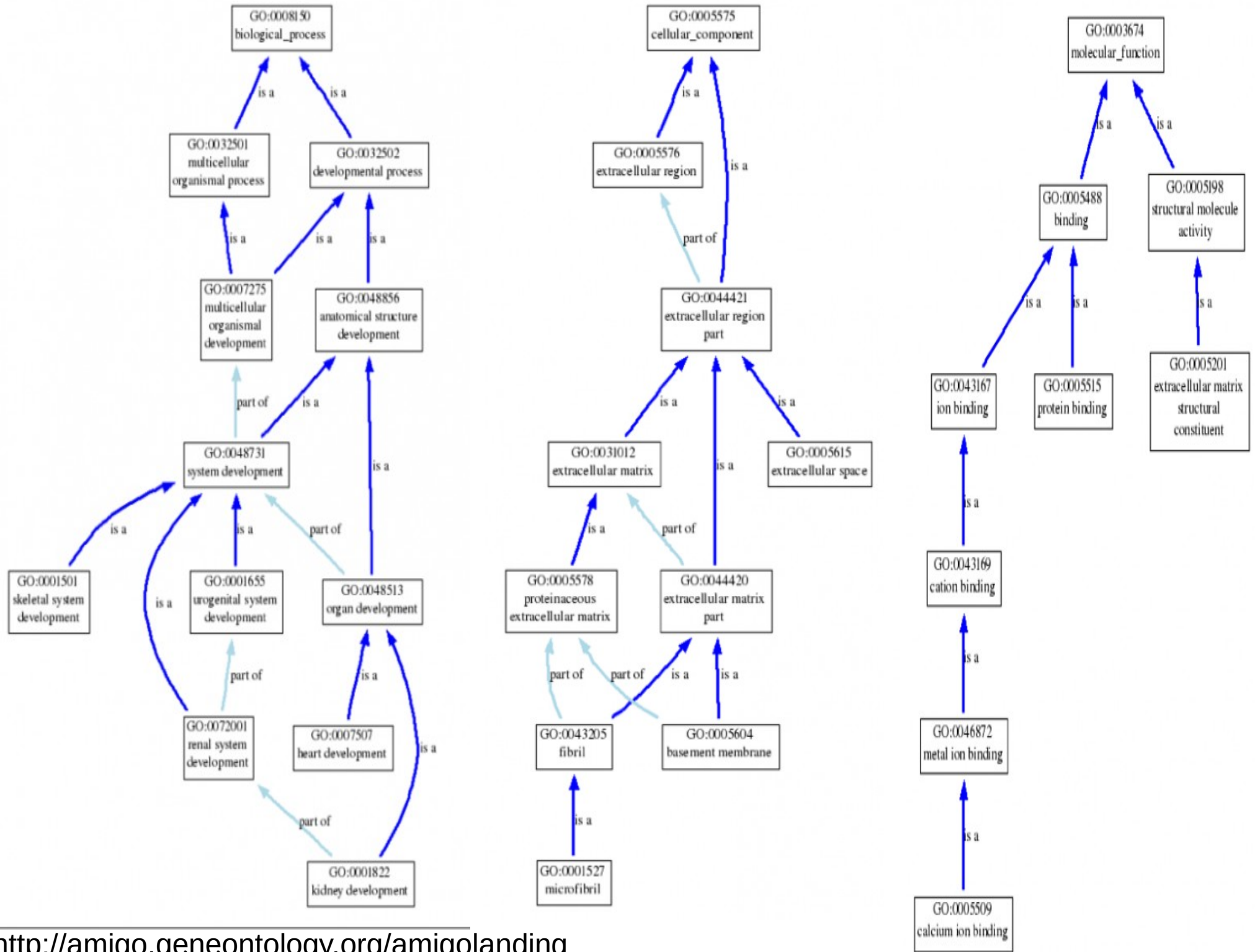


**GENEONTOLOGY**  
Unifying Biology

Le Gene Ontology Consortium est un projet visant à créer une base de données décrivant les organismes de manière commune, permettant aux chercheurs de comparer plus facilement les données entre les espèces

Les ontologies génétiques classent les informations selon trois niveaux différents:

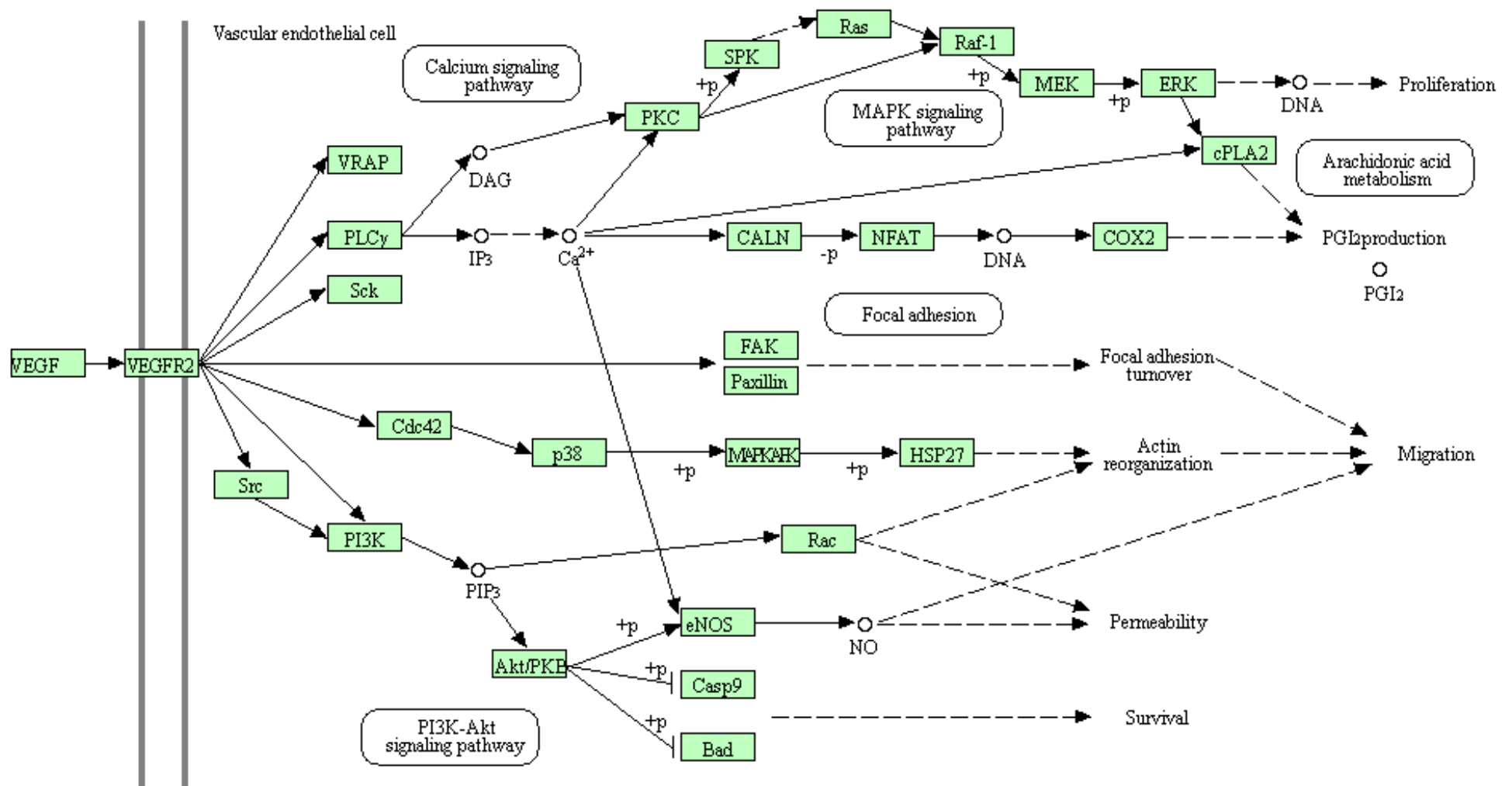
- fonction moléculaire (960900 ids)
- processus biologique (934523 ids)
- composant cellulaire (1019136 ids)



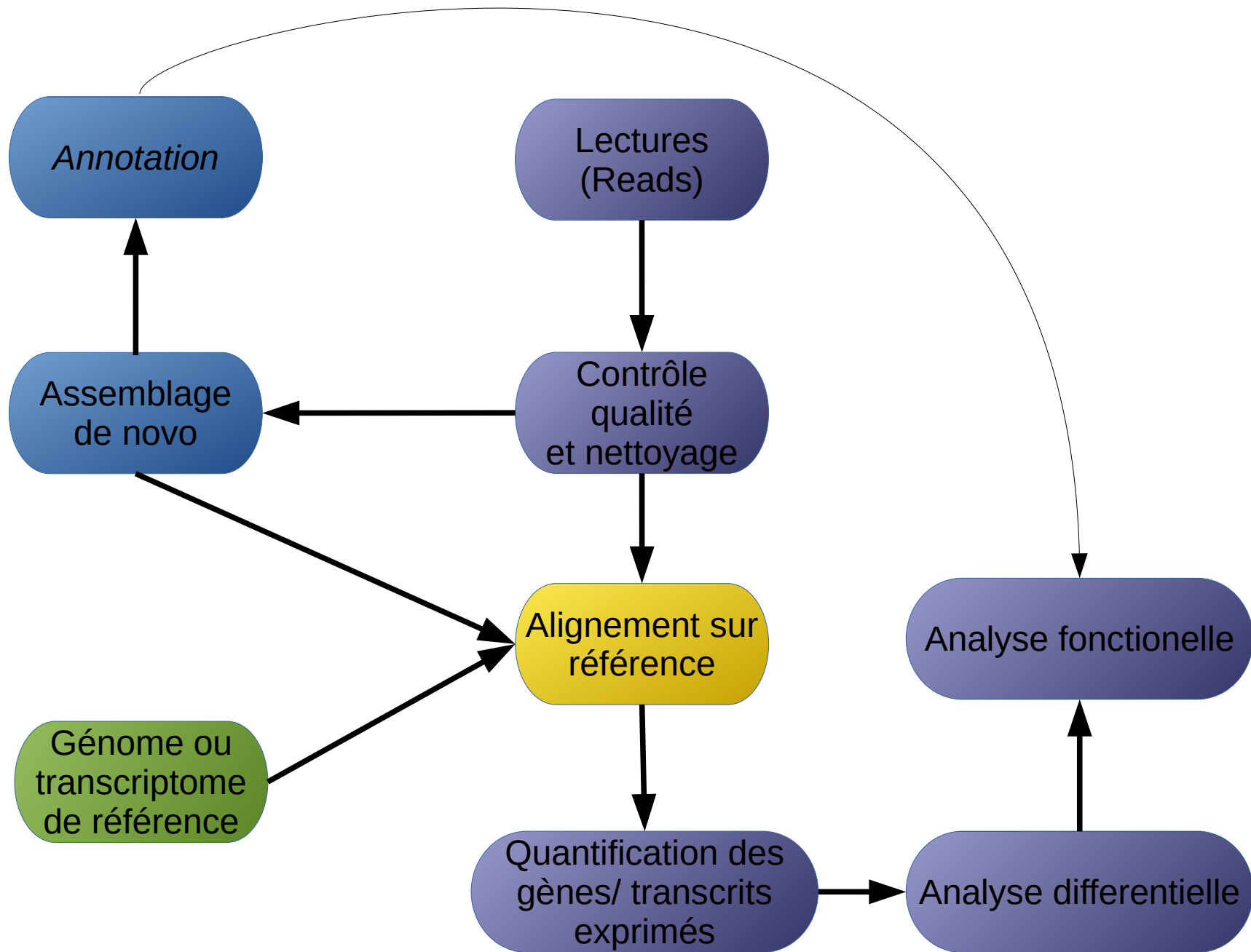


# KEGG

## VEGF SIGNALING PATHWAY



# Workflow d'analyse RNA-Seq



# Définition

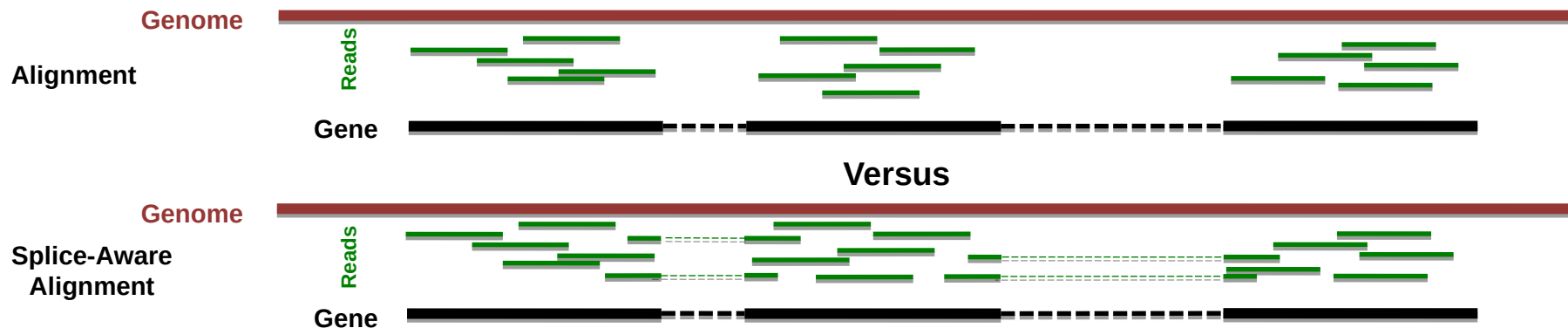
Le mapping est la prédiction du locus dont est originaire la lecture.

- Locus : le résultat est un ensemble de positions génomiques (ex.: chr1:100..150)
- Mapping ARN  $\neq$  Mapping ADN

Aligner les lectures issues du séquençage de dscDNA (transcrits) sur le génome, en tenant compte de l'épissage alternatif

Être capable d'exploiter les listes des jonctions exons-exons connues, mais également d'en détecter de nouvelles

Tout cela dans un temps raisonnable...



# Choix de l'aligneur

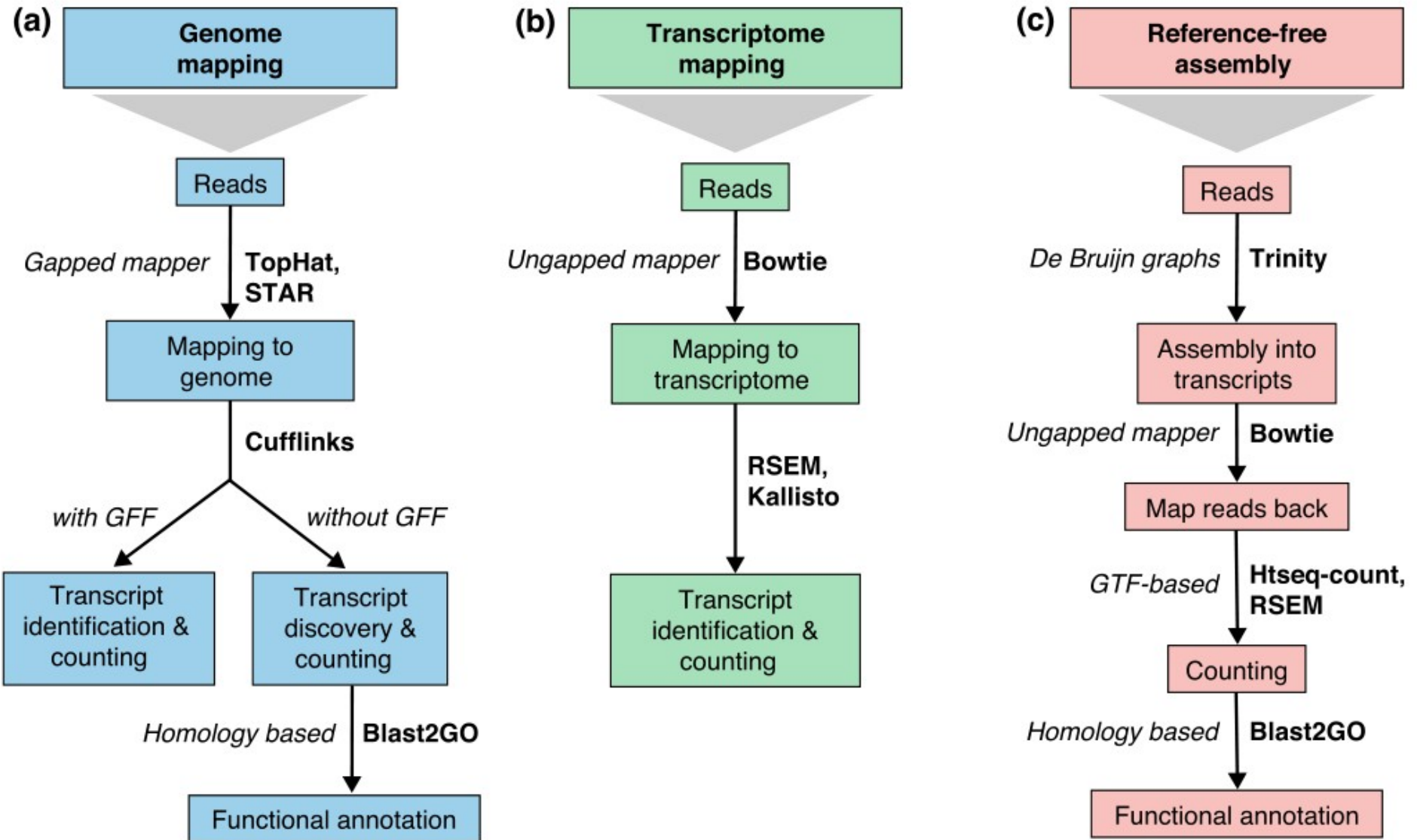
Nous devons aligner les données de séquence sur notre génome d'intérêt

- Si vous alignez les données RNA-Seq sur le génome, choisissez toujours un aligneur considérant l'épissage (sauf si c'est un génome bactérien!)

HiSat2, STAR, MapSplice, SOAPSplice, SpliceMap, GSNAP, HMMSplicer ...

- Il existe d'excellents aligneurs disponibles qui ne sont pas sensibles à l'épissages. Ceci est idéal pour les génomes bactériens

BWA, Novoalign (not free), Bowtie2, SOAPaligner, HiSat2 ...



# Choix de l'aligneur

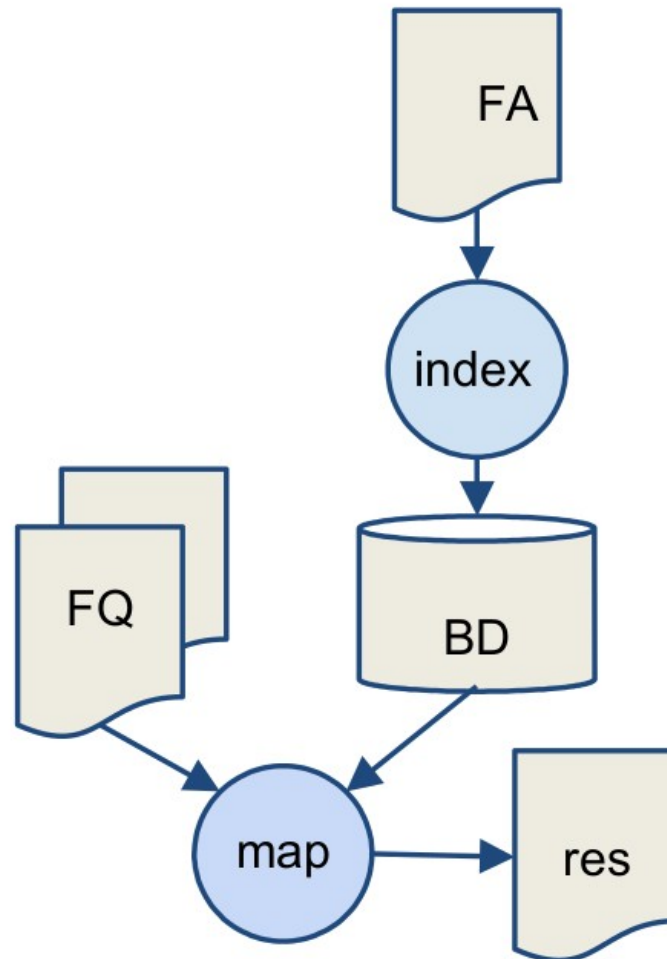
Autres considérations à prendre en compte dans le choix de l'aligneur :

- Comportement de l'aligneur par rapport à l'alignement multiple
- Prise en charge d'alignement en paired-end ou single-end
- Le nombre de mismatches qu'il autorise entre les reads et la référence
- Prise en charge de long-reads ou short-read ou des reads de tailles différentes

...

# Étapes de mapping

- ❖ Indexation du génome une fois pour toutes
- ❖ Mapping des lectures en utilisant l'index



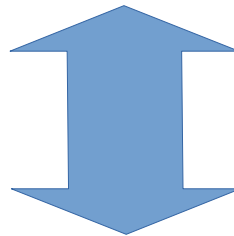
# SAM (Sequence Alignment Map)

```

Coord 12345678901234 5678901234567890123456789012345
ref    AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGCGCCAT
  
```

```

+r001/1      TTAGATAAAGGATA*CTG
+r002        aaaAGATAA*GGATA
+r003        gcctaAGCTAA
+r004                ATAGCT.....TCAGC
-r003                ttagctTAGGC
-r001/2                CAGCGGCAT
  
```



Header section										
@HD VN:1.5 SO:coordinate										
@SQ SN:ref LN:45										
r001	99	ref	7	30	8M2I4M1D3M	=	37	39	TTAGATAAAGGATACTG	*
r002	0	ref	9	30	3S6M1P1I4M	*	0	0	AAAAGATAAGGATA	*
r003	0	ref	9	30	5S6M	*	0	0	GCCTAAGCTAA	* SA:Z:ref,29,-,6H5M,17,0;
r004	0	ref	16	30	6M14N5M	*	0	0	ATAGCTTCAGC	*
r003	2064	ref	29	17	6H5M	*	0	0	TAGGC	* SA:Z:ref,9,+,5S6M,30,1;
r001	147	ref	37	30	9M	=	7	-39	CAGCGGCAT	* NM:i:1

Header section

Alignment section



# SAM (Sequence Alignment Map)

Name	Description
qname	Identifiant du read
flag	Valeur numerique qui encode les détails extra-alignement (voir samtools man page)
rname	Séquence de reference
strand	Brin d'alignement
pos	Coordonnée « start » de l'alignement
qwidth	Longueur du read
mapq	Qualité de l'alignement
cigar	Nomenclature CIGAR qui donne une vision globale de l'alignement
mrnm/rnext	reference chromosome/sequence to which opposite read pair aligned
mpos/pnext	alignment start coordinate to which opposite read pair aligned
isize/tlen	Taille de l'insert (paired-end)
seq	Séquence read qui s'aligne
qual	Score de la qualité
tags	Informations supplementaires sur l'alignement (optionnel) TAG:TYPE:VALUE

# SAM (CIGAR)

Symbole	Description
M	Alignement « matchant »
I	Insertion dans le read
D	délétion dans le read
N	Base illisible dans le read
S	Soft clipping (clipped sequences present in SEQ)
H	Hard clipping (clipped sequences NOT present in SEQ)
P	Padding (silent deletion from padded reference)
=	Séquence match
X	Séquence mismatch

Exemple :

```
RefPos:      1  2  3  4  5  6  7      8  9 10 11 12 13 14 15 16 17 18 19
Reference:   C  C  A  T  A  C  T      G  A  A  C  T  G  A  C  T  A  A  C
Read:                A  C  T  A  G  A  A      T  G  G  C  T
```

POS: 5

CIGAR: 3M1I3M1D5M

# SAM (Sequence Alignment Map)

name	flag	chr	start	mapQ	cigar	nNext	sNext	tlen	seq	qual	tags
r001	99	ref	7	30	8M2I4M1D3M =		37	39	TTAGATAAAGGATACTG	*	
r002	0	ref	9	30	3S6M1P1I4M *		0	0	AAAAGATAAGGATA	*	
r003	0	ref	9	30	5S6M *		0	0	GCCTAAGCTAA	*	SA:Z:ref,29,-,6H5M,17,0;
r004	0	ref	16	30	6M14N5M *		0	0	ATAGCTTCAGC	*	
r003	2064	ref	29	17	6H5M *		0	0	TAGGC	*	SA:Z:ref,9,+,5S6M,30,1;
r001	147	ref	37	30	9M =		7	-39	CAGCGGCAT	*	NM:i:1

#	Decimal	Description of first read
1	1	Read paired
2	2	Read mapped in proper pair
3	4	Read unmapped
4	8	Mate unmapped
5	16	Read reverse strand
6	32	Mate reverse strand
7	64	First in pair
8	128	Second in pair
9	256	Not primary alignment
10	512	Read fails platform/vendor quality checks
11	1024	Read is PCR or optical duplicate
12	2048	Supplementary alignment
<b>Sum</b>	<b>99</b>	

Decimal	Description of second read
1	Read paired
2	Read mapped in proper pair
4	Read unmapped
8	Mate unmapped
16	Read reverse strand
32	Mate reverse strand
64	First in pair
128	Second in pair
256	Not primary alignment
512	Read fails platform/vendor quality checks
1024	Read is PCR or optical duplicate
2048	Supplementary alignment
147	

## Common flags\*

One of the reads is unmapped:  
73, 133, 89, 121, 165, 181, 101, 117,  
153, 185, 69, 137

Both reads are unmapped:  
77, 141

Mapped within the insert size and in  
correct orientation:  
99, 147, 83, 163

Mapped within the insert size but in  
wrong orientation:  
67, 131, 115, 179

Mapped uniquely, but with wrong  
insert size:  
81, 161, 97, 145, 65, 129, 113, 177

\* Collected from [here](#)

<http://www.samformat.info/sam-format-flag>

# BAM (Binary Alignment/Map)

BAM (Binary Alignment/Map) format:

- version compressée et binarisée de SAM
- samtools: reading, writing, and manipulating BAM Files
- la plupart des outils exploitant les fichiers bam, requièrent des bam triés et indexés



Autres formats : CRAM

# Visualisation des alignements via IGV

