

INRA
SCIENCE & IMPACT

Formation RNA-Seq



Formation RNAseq Bioinfo/Biostat

Formateurs



❖ Bioinformatique :

- Sarah Maman
- Céline Noirot
- Olivier Rué
- Matthias Zytnicki

❖ Biostatistique :

- Ignaccio Gonzales
- Annick Moisan
- Nathalie Villa-Vialaneix

Organisation



- ❖ Formations de 3,5 jours
- ❖ Présentations seront entre-coupées de TP
- ❖ Horaires
 - Matin : 9h30 - 12h30
 - Après-midi : 14h00 – 17h

Tour de table



- ❖ Vous ? Votre labo ? Vos données ?
- ❖ Avez vous déjà traité des NGS ?
- ❖ Avez vous déjà utilisé Galaxy et R ?
- ❖ Vos attentes ?

Qu'allez vous apprendre ...

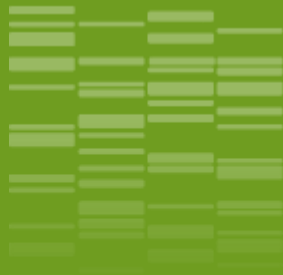


❖ A l'issue de 3 jours de formation **vous saurez** :

- Utiliser un environnement **Galaxy** pour **quantifier** et **découvrir de nouveaux transcrits**,
- Utiliser **Rstudio** pour détecter le gène DE correspondant à un **plan d'expérience simple**.

❖ Vous **ne saurez pas** :

- Étudier un transcriptome **sans génome** de référence,
- (Traiter bioinformatiquement de **nombreuses librairies**.)
- Détecter les gènes DE correspondant à un plan d'expérience **complexe**.



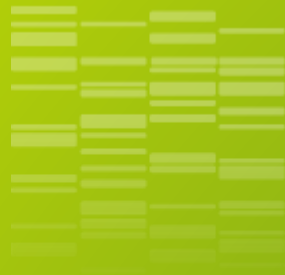
TRAITEMENT BIOINFORMATIQUE DE DONNÉES RNA-Seq



Plan



- ❖ **Mardi - Matin (9h30 - 12h30)**
 - Introduction au RNAseq (Biologie & Protocole)
 - Présentation Galaxy et prise en main
- ❖ **Mardi - Après-midi(14h – 17h)**
 - Vérification de la qualité
 - Algorithmes d'alignement
- ❖ **Mercredi - Matin (9h30 - 12h30)**
 - Les formats de fichier d'alignement
 - Visualisation
 - Assemblage de transcrits
- ❖ **Mercredi - Après-midi (14h00 - 17h)**
 - Quantification de gènes
 - Quantification de transcrits
 - Créer vos workflow RNAseq



_01

Rappels biologiques

Un peu de vocabulaire

- ❖ **Transcriptome** : Ensemble des transcrits d'un organisme
- ❖ **RNAseq de novo** : Etude du transcriptome sans génome de référence.
- ❖ **Read** : Lecture
- ❖ **Fragment** : Paire de lecture

Rappels biologiques

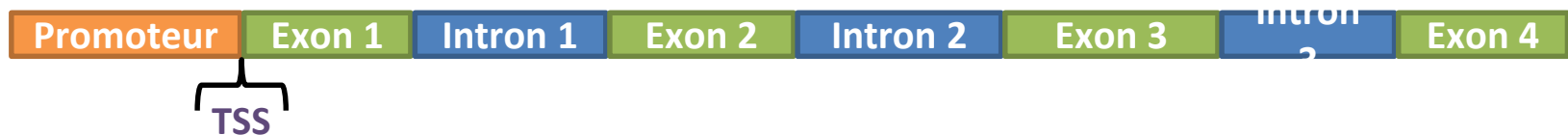


Qu'est-ce qu'un gène ?

Rappels biologiques

Qu'est-ce qu'un gène ?

- o **Gène** : unité fonctionnelle de l'ADN qui contient les instructions nécessaires à la création d'un produit fonctionnel



- o **Promoteur** : zone de fixation des ribosomes
- o **TSS** : site de départ de transcription
- o **Exon** : région codante de l'ARNm inclus dans le transcrit
- o **Intron** : région non codante

Rappels biologiques

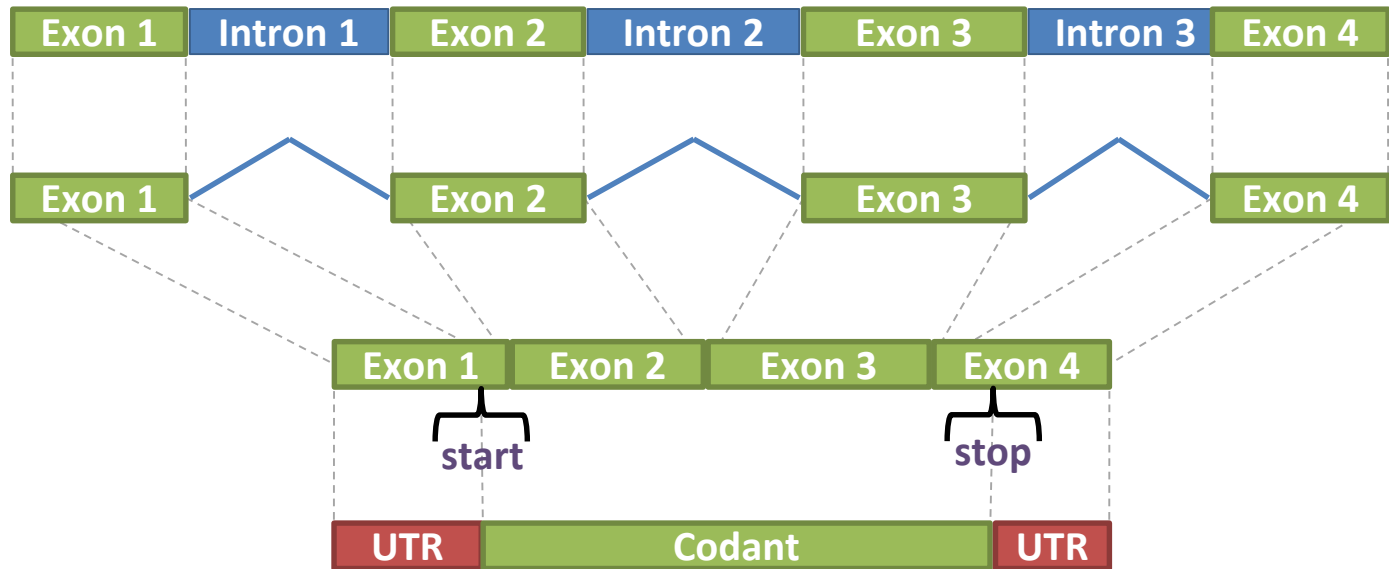


Qu'est-ce qu'un transcrit?

Rappels biologiques

Qu'est-ce qu'un transcrit ?

- o **Epissage** : Excision des introns avant traduction



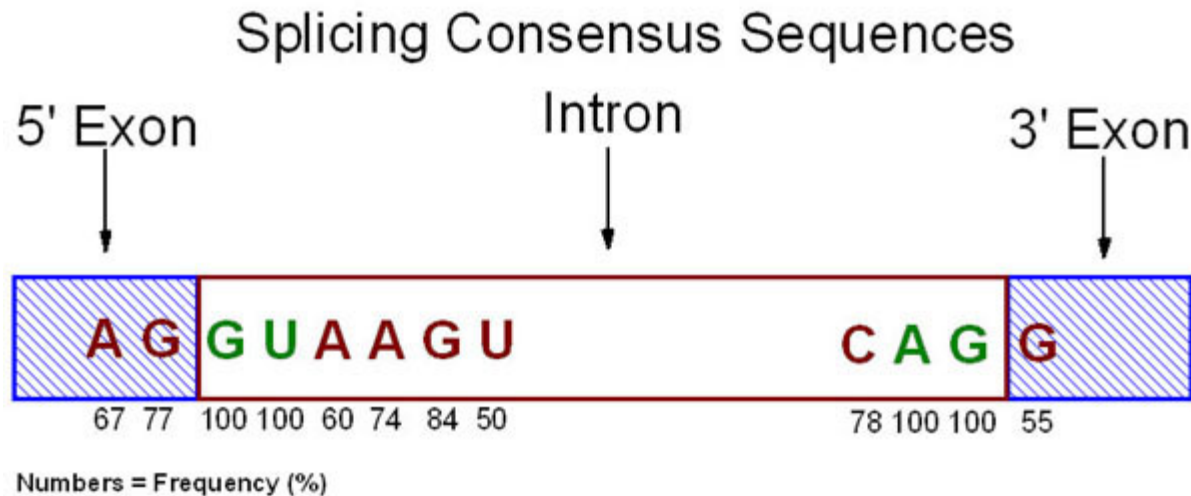
- o **Transcrit** : portion d'ADN transcrite en molécule d'ARN
- o **UTR** : région transcrite mais pas traduite

Rappels biologiques

Qu'est-ce qu'un site d'épissage?

o Site d'épissage canonique :

- plus de 99% de **GT** et **AG** comme sites **donneurs** et **accepteurs**

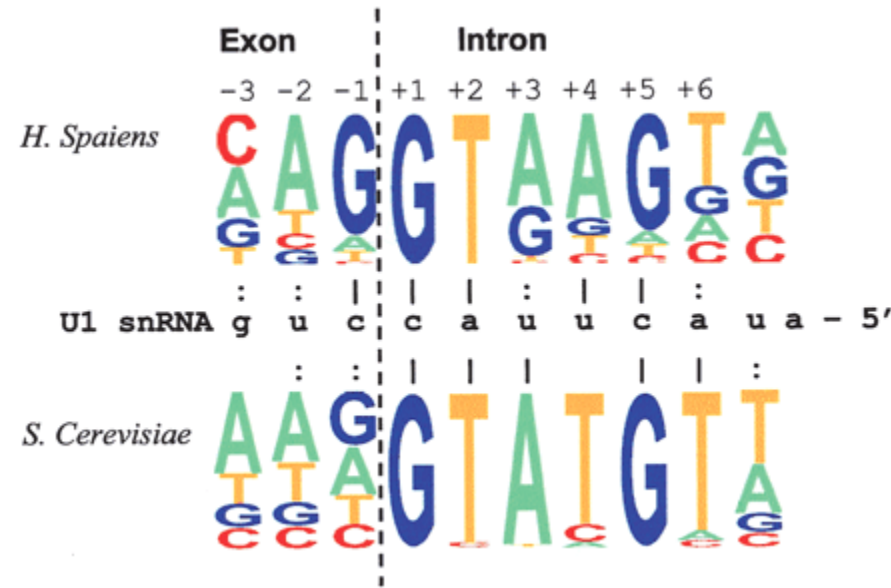


<http://themedicalbiochemistrypage.org/rna.php>

Rappels biologiques

Qu'est-ce qu'un site d'épissage?

- o **Site d'épissage non-canonique :**
 - GC-AG ou AT-AC comme sites **donneurs** et **accepteurs**
- o **Mammifère :**
 - 0.69% GC-AG
 - 0.05% AT-AC
- o **Autre exemple :**



<http://rnajournal.cshlp.org/content/10/5/828.full>

Rappels biologiques

Epissage alternatif et isoformes

- o Excision d'exon



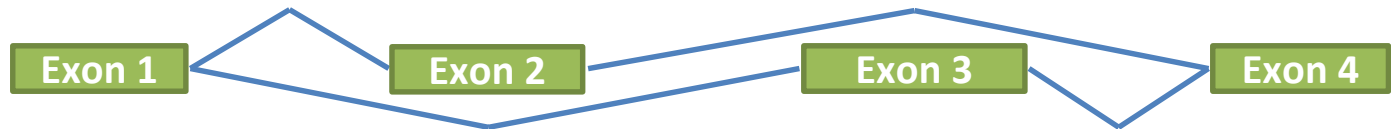
- o Rétention d'intron



- o TSS alternatif



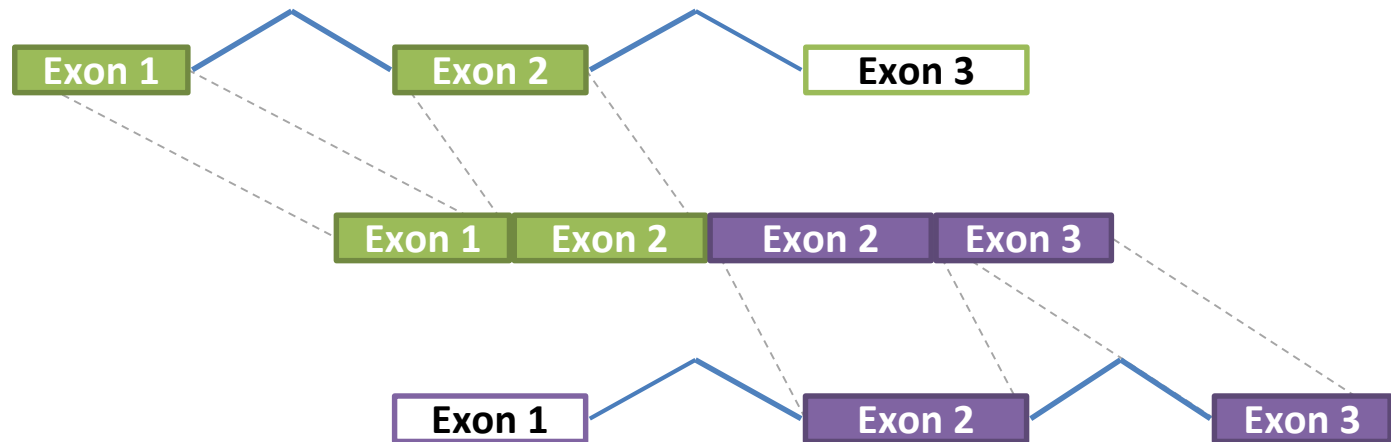
- o Exons exclusifs



Rappels biologiques

Et plus encore ?

- o Fusion de gènes ou Trans-épissage

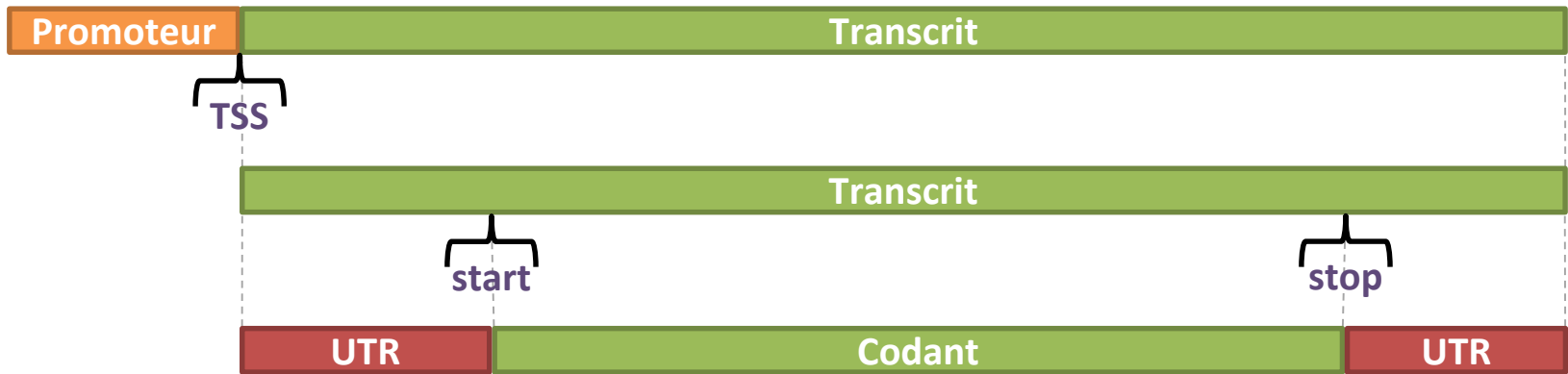


- o Chimère biologique

Rappels biologiques

Gène procaryote / gène eucaryote

- o Pas d'intron chez les procaryotes





Le RNA-Seq

Modes d'étude du transcriptome

- ❖ EST
- ❖ rt-PCRq
- ❖ puce d'expression
- ❖ tiling array

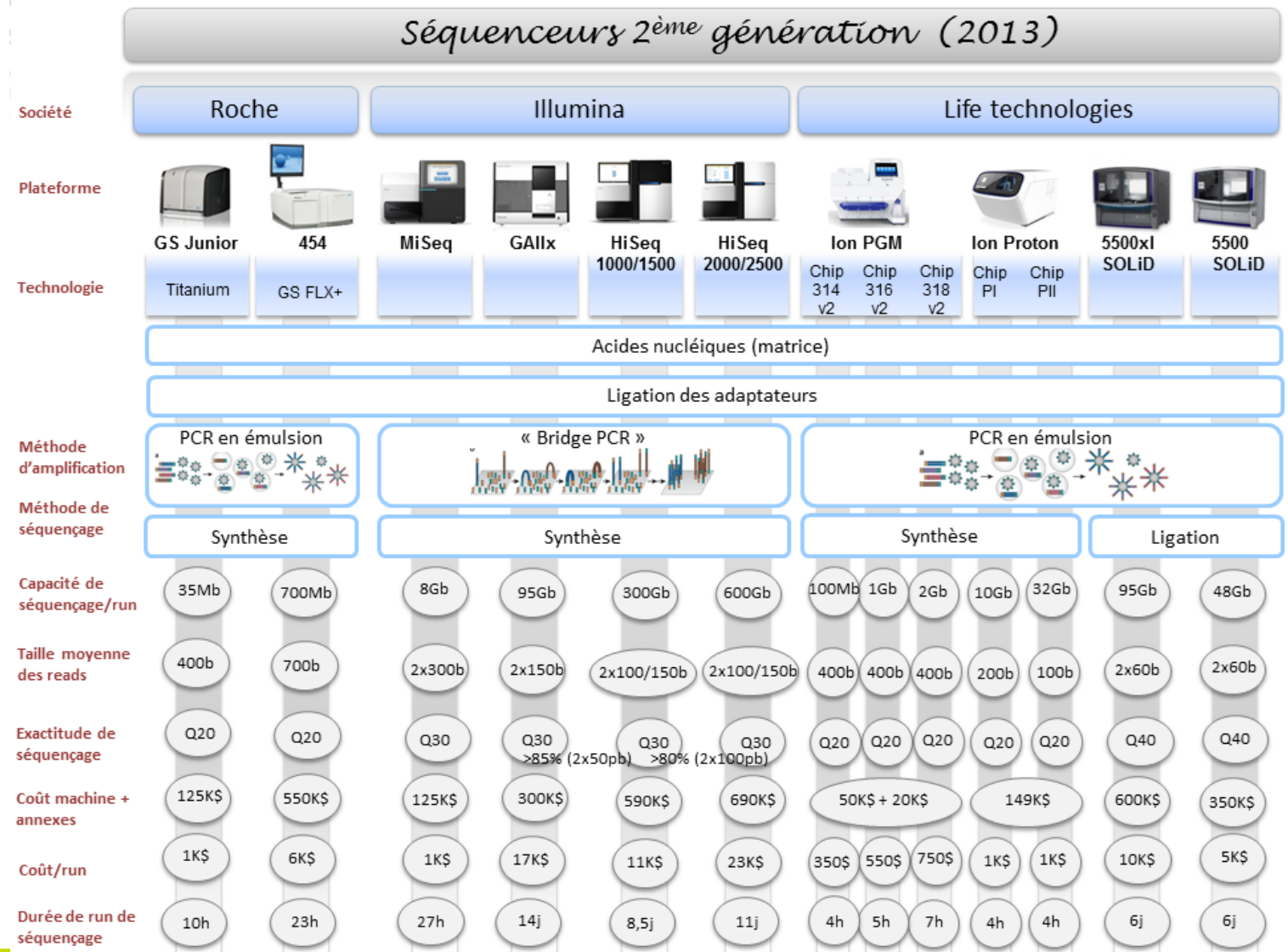
- ❖ RNA-Seq

Quelles sont les principales différences ?

Modes d'étude du transcriptome

- ❖ Pas besoin d'avoir de connaissance sur la séquence
- ❖ Spécificité de ce que l'on mesure
- ❖ Augmente l'échelle de mesure
- ❖ Quantification directe
- ❖ Très bonne reproductibilité
- ❖ Différents niveau d'étude : gènes, transcrits, spécificité allélique, variant de structure
- ❖ Découverte de nouveaux : transcrits, isoformes, (ncRNA), structures (fusion...)
- ❖ Détection possible of SNPs, ...

Les séquenceurs

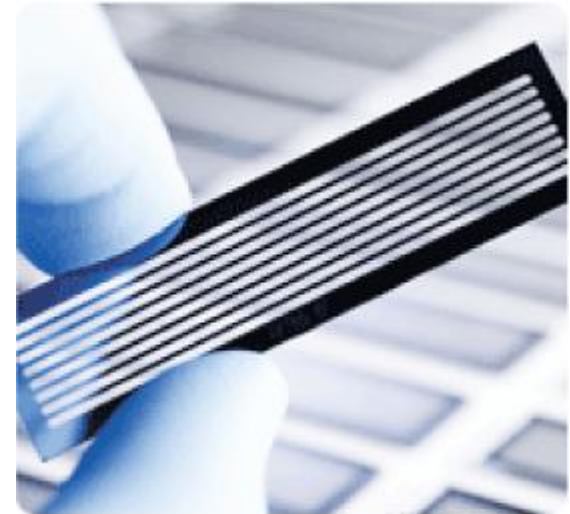


Les séquenceurs

Séquenceurs 2 ^{ème} génération (2013)													
Société	Roche		Illumina				Life technologies						
Plateforme													
Technologie	Titanium	GS FLX+			1000/1500	2000/2500	Chip 314 v2	Chip 316 v2	Chip 318 v2	Chip PI	Chip PII	SOLiD	SOLiD
Génome humain	✗	✗	✗	✓	✓	✓	✗	✗	✗	✗	✗	✓	✓
Exome	✗	✗	✓	✓	✓	✓	✗	✗	✗	✓	✓	✓	✓
Petit génome (Bactéries, levures)	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Régions ciblées	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Transcriptome	✗	✗	✗	✓	✓	✓	✗	✗	✗	✓	✓	✓	✓
Chip-Seq	✗	✗	✓	✓	✓	✓	✗	✗	✓	✓	✓	✓	✓
Métagénomique	✗	✓	✓	✓	✓	✓	✗	✓	✓	✓	✓	✓	✓

Illumina sequencing vocabulary

- ❖ **Flowcell : 1 plaque (en général 1 run)**
- ❖ Lane : ligne de séquençage
- ❖ 1 Flowcell : 8 Lane
- ❖ 1 flowcell Hiseq 2500 : 2 Milliard de reads single ou 4 Milliard de reads paired.
- ❖ Hiseq 2000 / Hiseq 2500 : séquençage possible de 2 flowcells en parallèle.



Le protocole RNAseq



Préparation des Echantillons biologiques pour le RNAseq

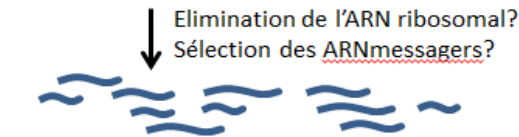
1. ARN messager ou ARN total



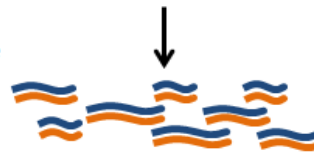
2. Elimination de l'ADN contaminant



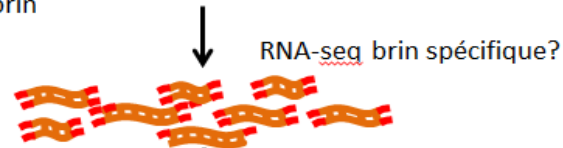
3. Fragmentation de l'ARN



4. Retro-transcription de l'ARN en cDNA, hybride d'ADN/ARN



5. Synthèse du second brin d'ADN et ligation d'adaptateurs



6. Sélection des fragments par la taille



7. Séquençage des extrémités et production de « reads »



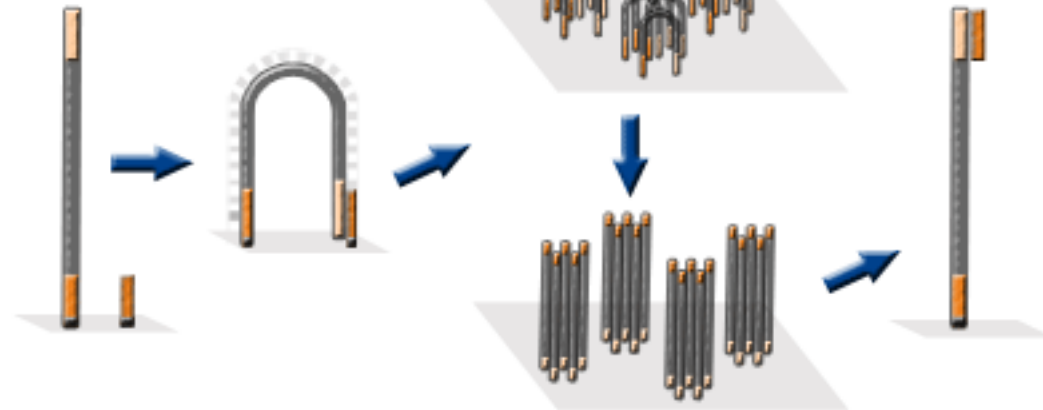
Séquençage illumina

1. Attach DNA to flow cell

2. Perform bridge amplification

3. Generate clusters

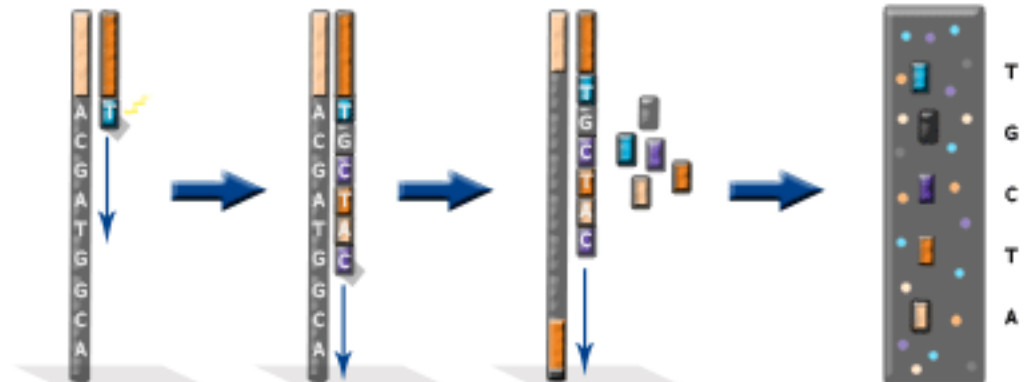
4. Anneal sequencing primer



5. Extend first base, read, and deblock

6. Repeat step above to extend strand

7. Generate base calls



Quels choix quand on fait du RNA-Seq ?



- ❖ **Déplétion / enrichissement**
- ❖ **Paired-end / single-end**
- ❖ **Séquençage en tenant compte du sens du brin**
- ❖ **Nombre de séquence / de réplicats**
- ❖ **Multiplexage**

Déplétion / Enrichissement

❖ Déplétion :

- Suppression des rRNA

❖ Enrichissement polyA :

- Pas de transcrits sans queue PolyA ou partiellement dégradés

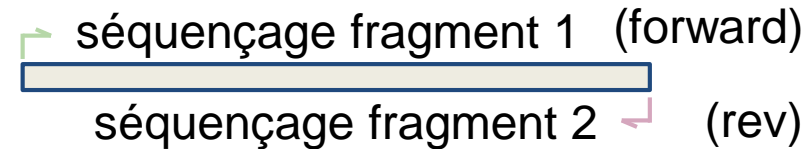
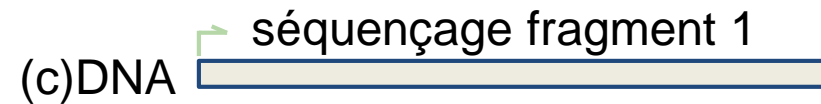
❖ Résultats semblables d'après :

Comparison of RNA-Seq by poly (A) capture, ribosomal RNA depletion, and DNA microarray for expression profiling, BMC Genomics , 2014

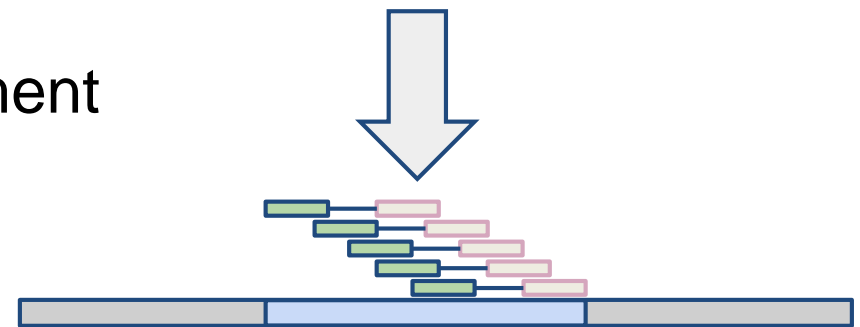
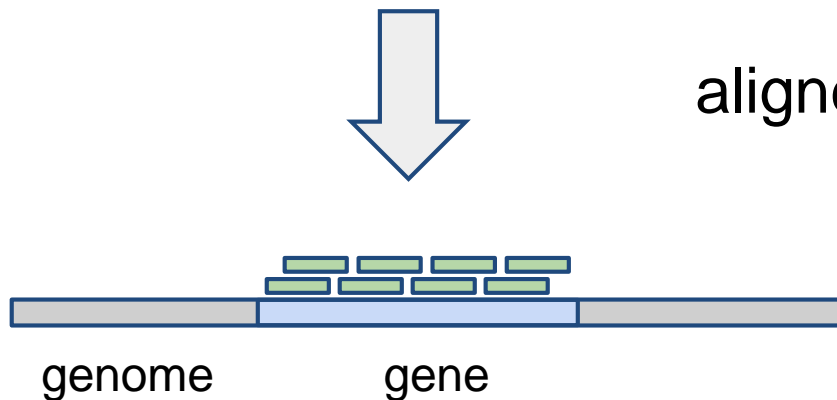
Single-end vs Paired-end

Single-end

Paired-end



alignement



- ❖ La taille des cDNA détermine la taille d'insert (p. ex. 200-500 pb).
- ❖ Les fragments sont habituellement en Forward-Reverse.

Paired-end



- ❖ Protocole différent (Adaptateurs spécifiques)
- ❖ Améliore le mapping
- ❖ Aide à la détection de variant alternatif
- ❖ Plus généralement aide à la détection de : variation structurale de génome (insertion/délétion), CNV, réarrangement génomique



L'intérêt des librairies brin spécifique

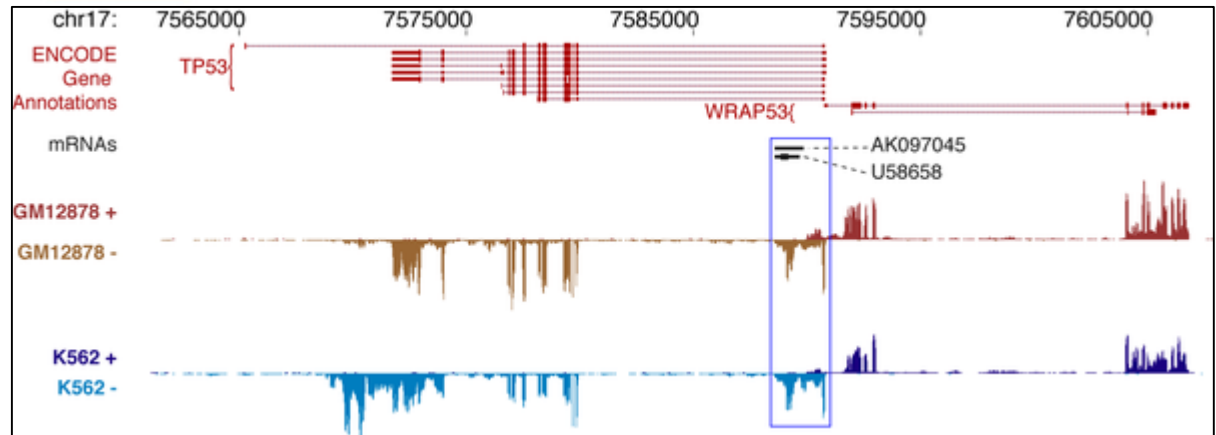
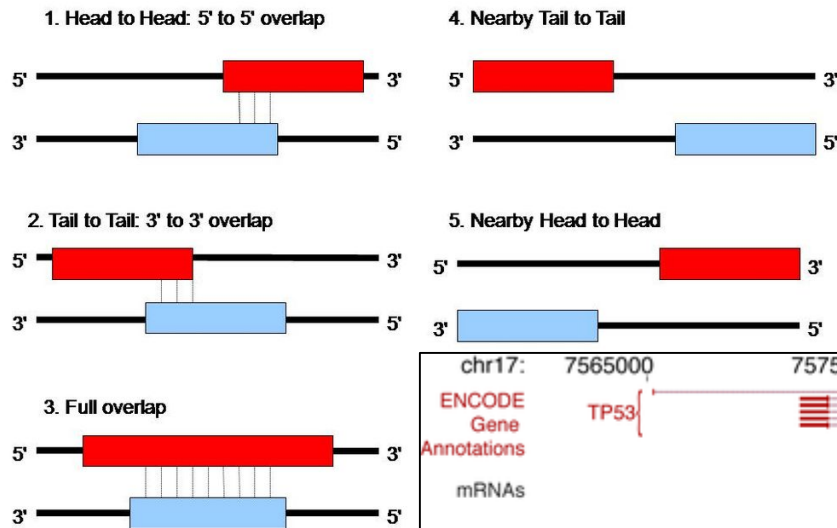
Nat Methods. 2010 Sep;7(9):709-15. Epub 2010 Aug 15.

Comprehensive comparative analysis of strand-specific RNA sequencing methods.

Levin JZ, Yassour M, Adiconis X, Nusbaum C, Thompson DA, Friedman N, Gnirke A, Regev A.

Broad Institute of Massachusetts Institute of Technology and Harvard University, Cambridge, Massachusetts, USA.
jlevin@broadinstitute.org

Abstract



Profondeur / Répétitions ?

- ❖ **Equilibre profondeur / nombre de répétitions :**
 - directives du consortium ENCODE en 2011
 - **plus de deux répétitions biologique**
- ❖ **Chez l'humain 100M de lectures** sont suffisantes pour détecter **90 % des transcrits** de **81 % des gènes** du **transcriptome humain**.
- ❖ **20M de lectures** (75bp) permettent de détecter des **transcrits exprimés à un niveau moyen ou faible** chez le **poulet**.
- ❖ **10 M de lectures** permettent que **90% des transcrits** (humain, zebrafish) soient **couverts par 10 lectures en moyenne**.

(Plus d'informations : Toung et al. 2011 ; Wang et al. 2011 ; Hart et al. 2013)

Profondeur / Répétitions ?

❖ Pourquoi augmenter le nombre de répétitions biologiques ?

- Généraliser les résultats à la population
- Estimer avec plus de précision la variation de chaque transcrit individuellement (*Hart et al. 2013*)
- Améliorer la détection des transcrits différentiels et le contrôle du taux de faux positifs : **VRAI à partir 3** ([Zhang et al. 2014](#), *Sonenson et al. 2013*, *Robles et al 2012*)

Profondeur / Répétitions ?

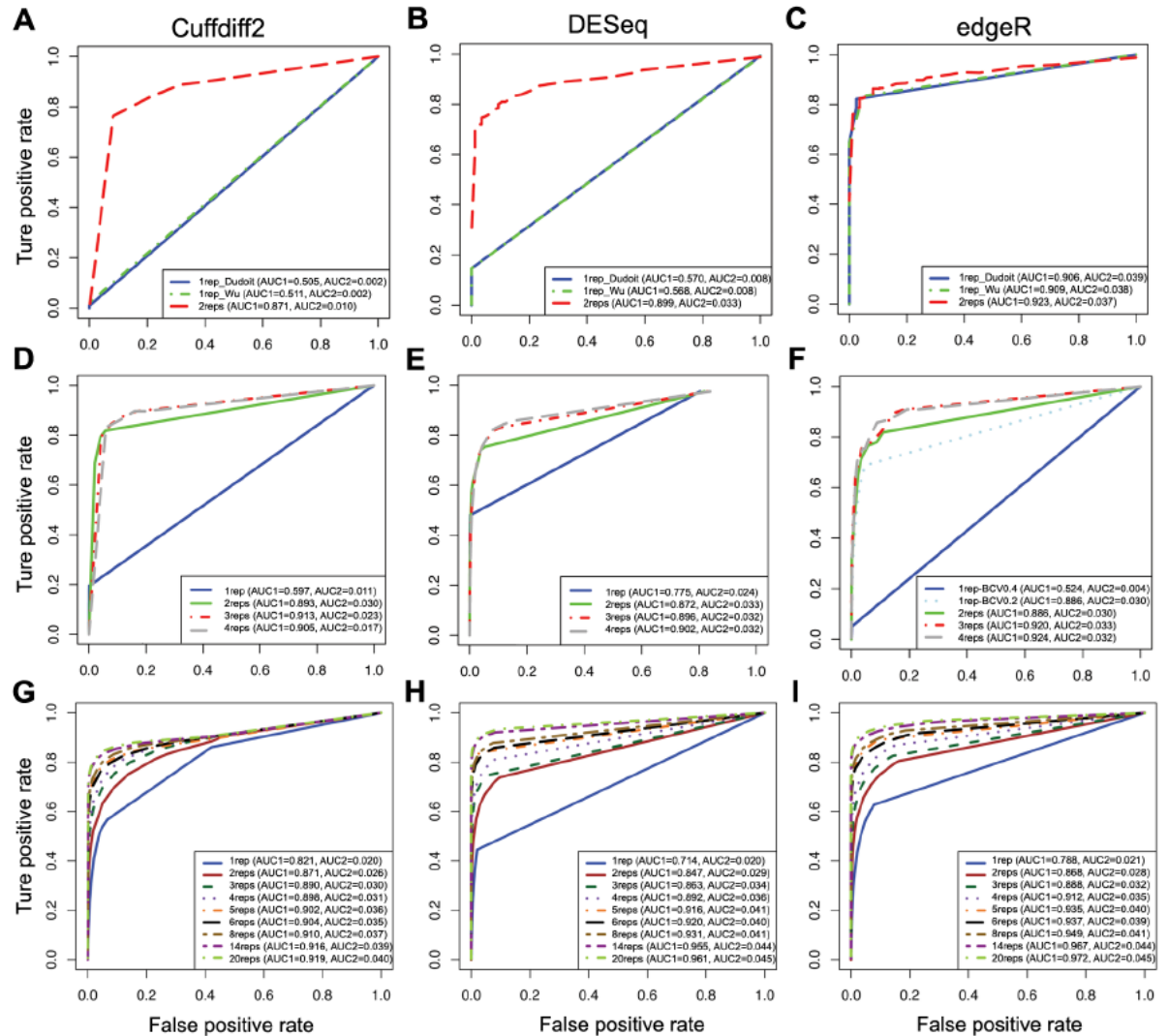
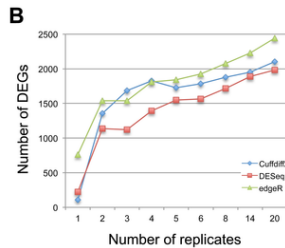
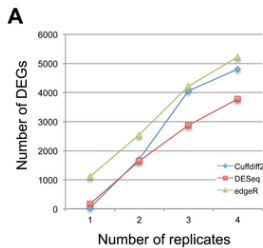


L'effet du nombre de réplicats sur le taux de vrai positifs et de faux positifs

K_N

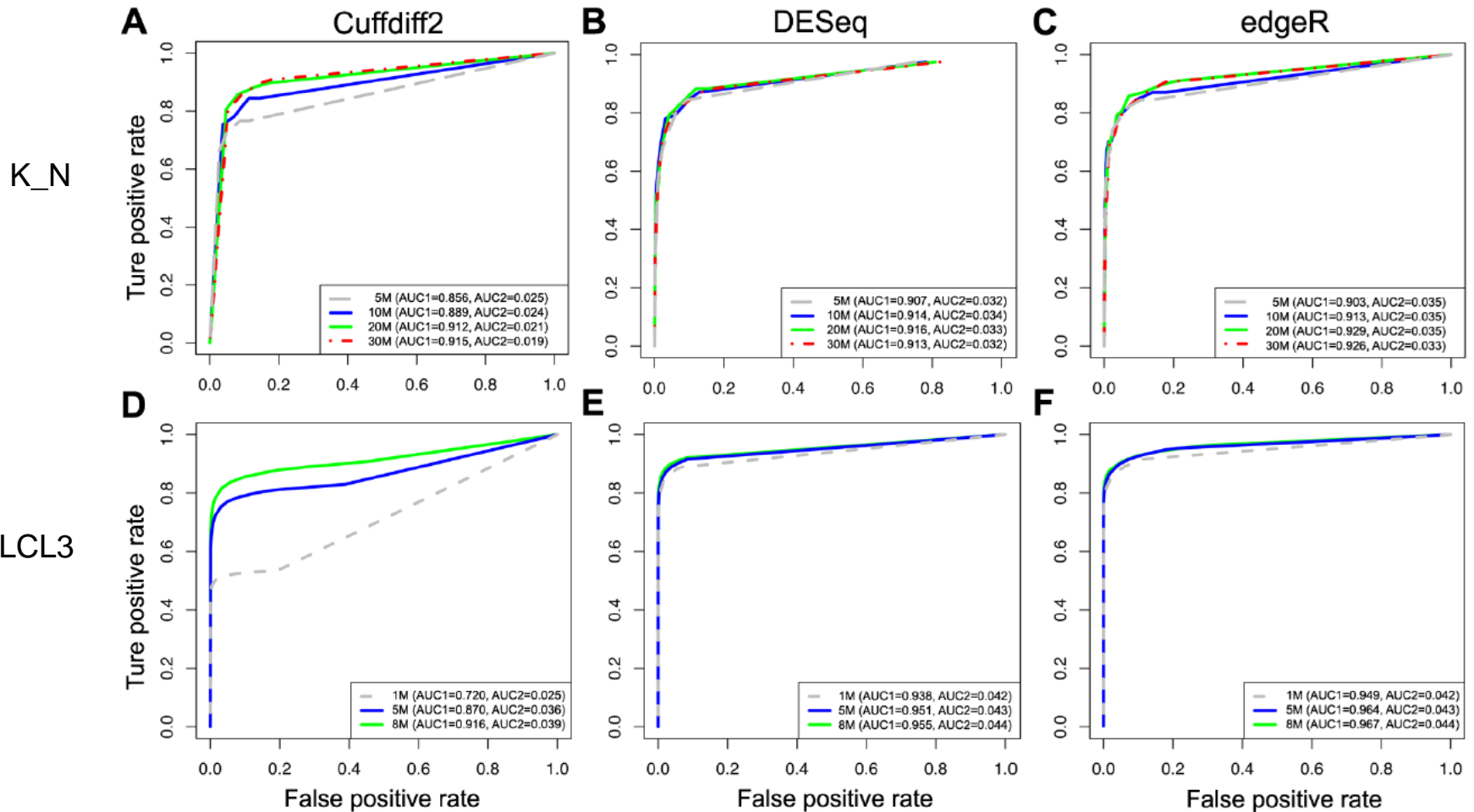
K_N

LCL2



Profondeur / Répétitions ?

L'effet de la profondeur.



Profondeur / Répétitions ?

- ❖ **Quel choix ? Plus de profondeur *ou plus de répétition* ?**

- ❖ **Ça dépend !** (Haas et al. 2012, Liu Y. et al 2013)
 - Détection de transcrits différentiels :
(+) répétitions biologiques

 - Construction/annotation transcriptome :
(+) profondeur & (+) conditions

 - Recherche de variants :
(+) répétitions biologiques & (+) profondeur

A quelles questions biologiques PEUT répondre le RNA-seq ?

- ❖ L'**analyse d'expression différentielle** (différence d'expression) au niveau du transcriptome
- ❖ L'étude de l'**épissage alternatif** (isoformes) et recherche de **nouveaux transcrits**
 - amélioration des annotations structurales existantes
- ❖ La recherche d'**allèles spécifiques** et la **quantification** de leur **expression**
- ❖ La construction d'un **transcriptome *de novo*** (organismes non modèles)

Stratégie d'analyse en fonction des données disponibles

❖ De novo :

- Pas de génome/transcriptome de référence
- Outils en évolution permanente
- Ressources (cpu/disque) +++

❖ Transcriptome de reference

- Dépendant de la qualité de l'annotation structurale
- Peu couteux

❖ Génome de référence

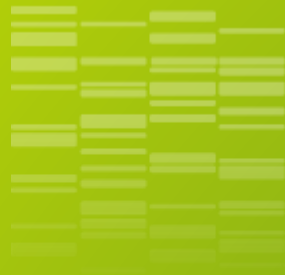
- Permet une approche combinée :
 - sur **transcriptome**
 - recherche de **nouveaux transcrits**
- Ressources ++
- Alignement épissé

Pipeline d'analyse RNA-Seq : avec référence

- ❖ **Contrôle qualité**
- ❖ **Pre-nettoyage** des lectures
 - suppression des adaptateurs de séquençage
 - (suppression des adaptateurs de multiplexage)
- ❖ **Nettoyage** des lectures
 - tronquer les extrémités de mauvaise qualité des lectures
- ❖ **Alignement des lectures sur la référence**
 - gènes ou génome complet
- ❖ **Comptage** des gènes / transcrits



Galaxy : un outil de traitement de données

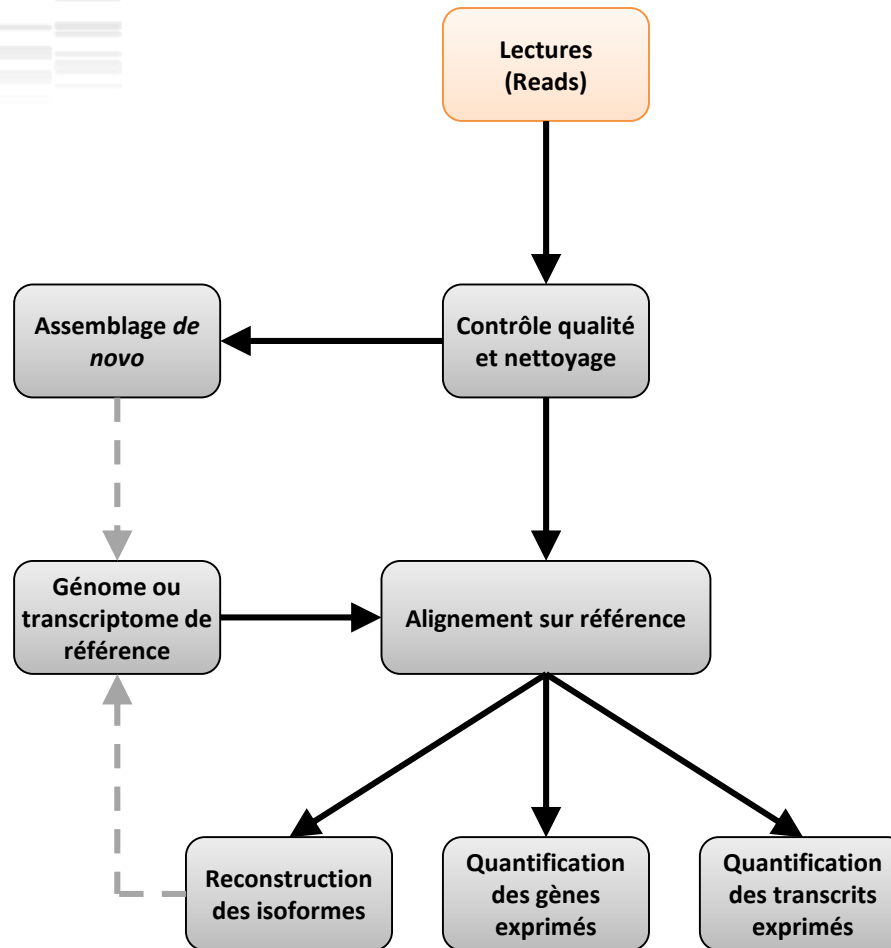


_02

DONNÉES BRUTES

Obtenir des séquences de qualités

Workflow d'analyse RNA-Seq



Plan



- ❖ **Le format fastq**
- ❖ **Les biais connus**
- ❖ **Vérification de la qualité avec FastQC**
- ❖ **Nettoyage des lectures avec Sickle**

Format Fastq

- o 1 séquence = 4 lignes dans le fichier

```
@SEQ_ID
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*(((((***+))%%&&+))(%&&&).1***-+*'' )**55CCF>>>>>CCCCCCC65
```

- o 1 ère ligne = identifiant de la séquence

```
@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG
```

EAS139	the unique instrument name
136	the run id
FC706VJ	the flowcell id
2	flowcell lane
2104	tile number within the flowcell lane
15343	'x'-coordinate of the cluster within the tile
197393	'y'-coordinate of the cluster within the tile
1	the member of a pair, 1 or 2 (<i>paired-end or mate-pair reads only</i>)
Y	Y if the read fails filter (read is bad), N otherwise
18	0 when none of the control bits are on, otherwise it is an even number
ATCACG	index sequence

Format Fastq



- o 4ème ligne = Qualité

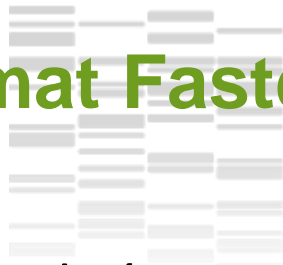
```
@SEQ_ID
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*((( (***) )#####) (#####) .1***-+*'')**55CCF>>>>>CCCCCCC65
```

- o Appelée aussi Phred quality score (Sanger format)

$$Q_{\text{sanger}} = -10 \log_{10} p$$

Probabilité qu'une base soit incorrecte

Format Fastq



o Qualité Encodée en ASCII

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS.....
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....
.....IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII.....
.....JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ.....
..LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL.....
!"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMN...
|                                     |                                     |
33                                     59 64 73                                     104
0.....26...31.....40
-5...0.....9.....40
0.....9.....40
3.....9.....40
0.2.....26...31.....41
```

S - Sanger Phred+33, raw reads typically (0, 40)
X - Solexa Solexa+64, raw reads typically (-5, 40)
I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)
J - Illumina 1.5+ Phred+64, raw reads typically (3, 40)
with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
(Note: See discussion above).

L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)



Objectif du contrôle qualité

- ❖ Vérifier que les séquences sont **conformes au niveau de prestation attendu (taille, nombre, qualité,...)**
- ❖ Vérifier que les séquences peuvent **répondre au questions biologiques** posées :
 - **Biais techniques**
 - **Biais biologiques**
- ❖ Aider au paramètres pour le nettoyage des données

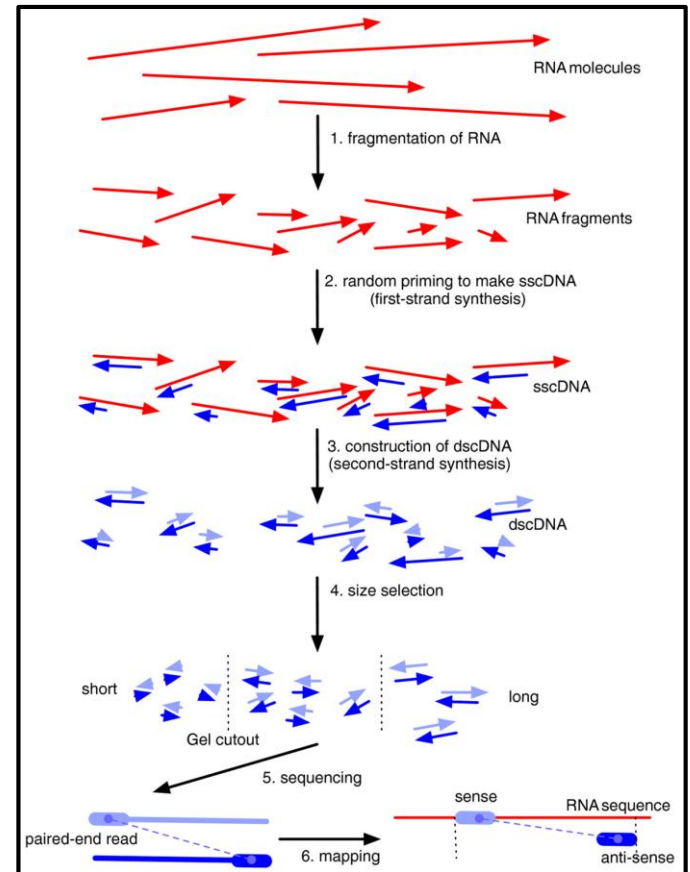


Biais spécifiques au RNA-Seq

- ❖ Influence du mode de préparation de la banque
 - amplification hexamérique aléatoire (**Random hexamer priming**)
- ❖ Influence du séquençage
 - biais de position, de composition en séquence (contenu en GC)
 - influence de la longueur des transcrits
- ❖ « Mapabilité » du génome/transcriptome

Préparation de la banque

- ❖ Extraction ARN total
- ❖ Déplétion (queue polyA)
- ❖ Fragmentation, reverse transcription avec des hexamères aléatoires -> dscDNA
- ❖ Séquençage



Roberts et al. Genome Biology 2011, 12:R22

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

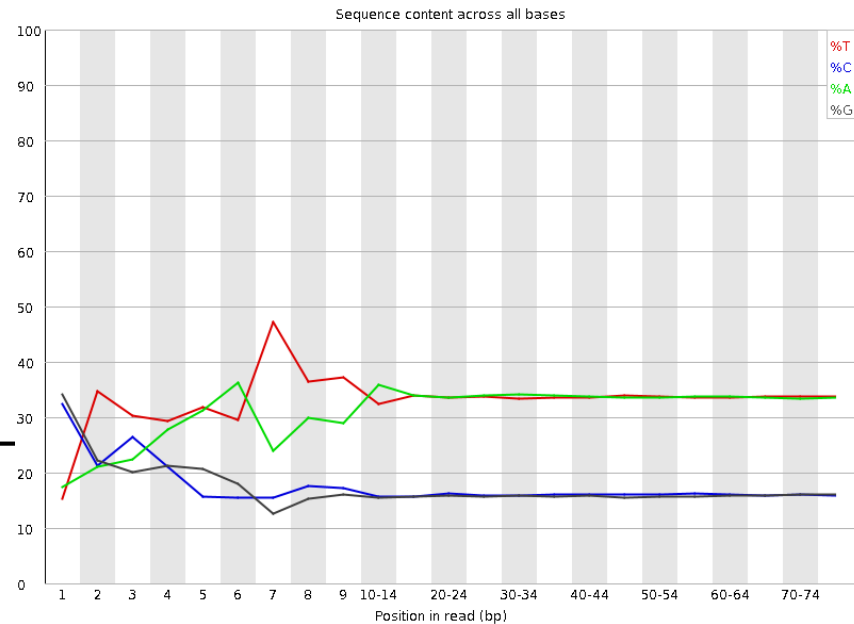
Nucleic Acids Research, 2010, Vol. 38, No. 12 e131
doi:10.1093/nar/gkq224

Biases in Illumina transcriptome sequencing caused by random hexamer priming

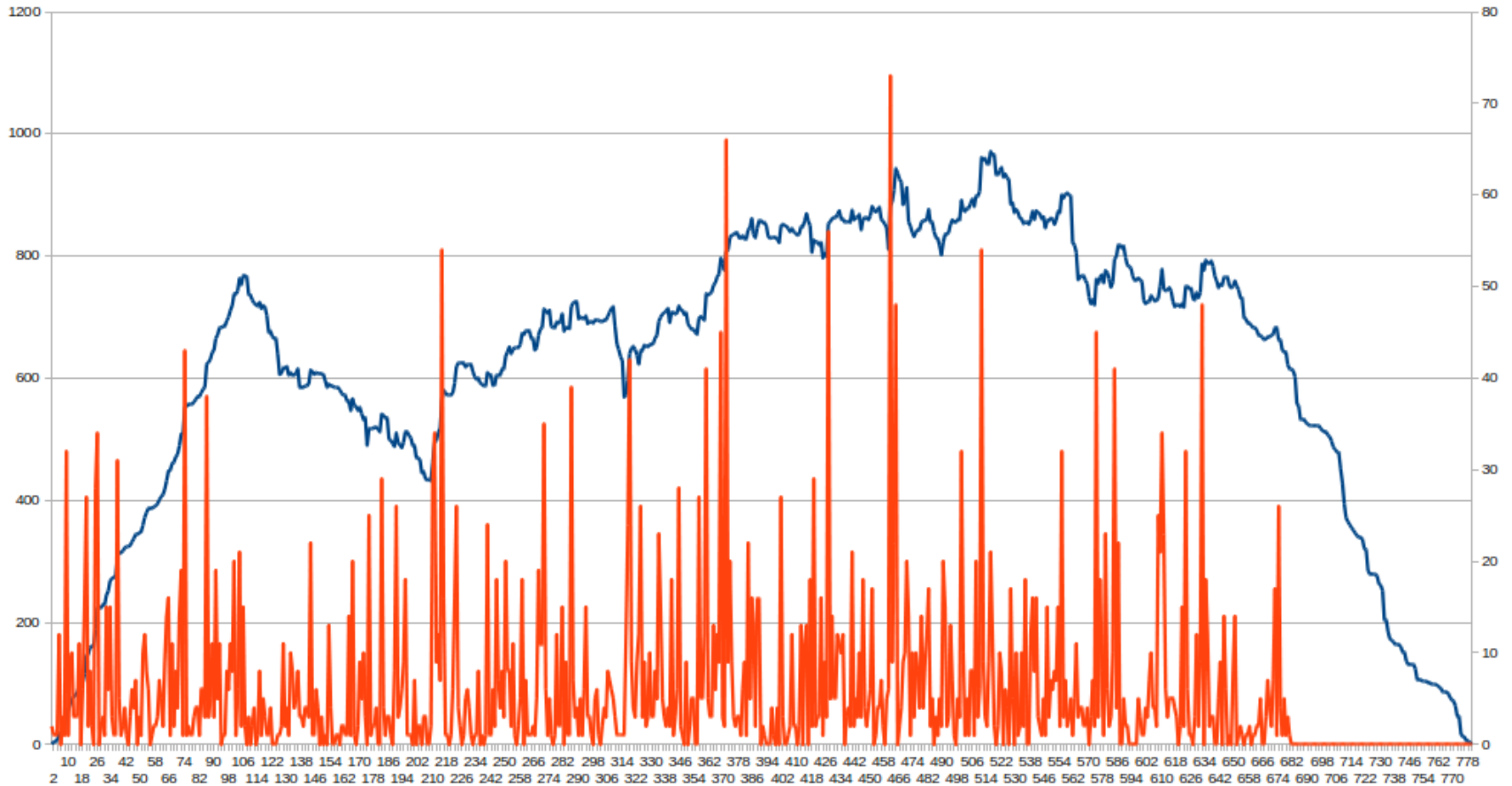
Kasper D. Hansen^{1,*}, Steven E. Brenner² and Sandrine Dudoit^{1,3}

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.



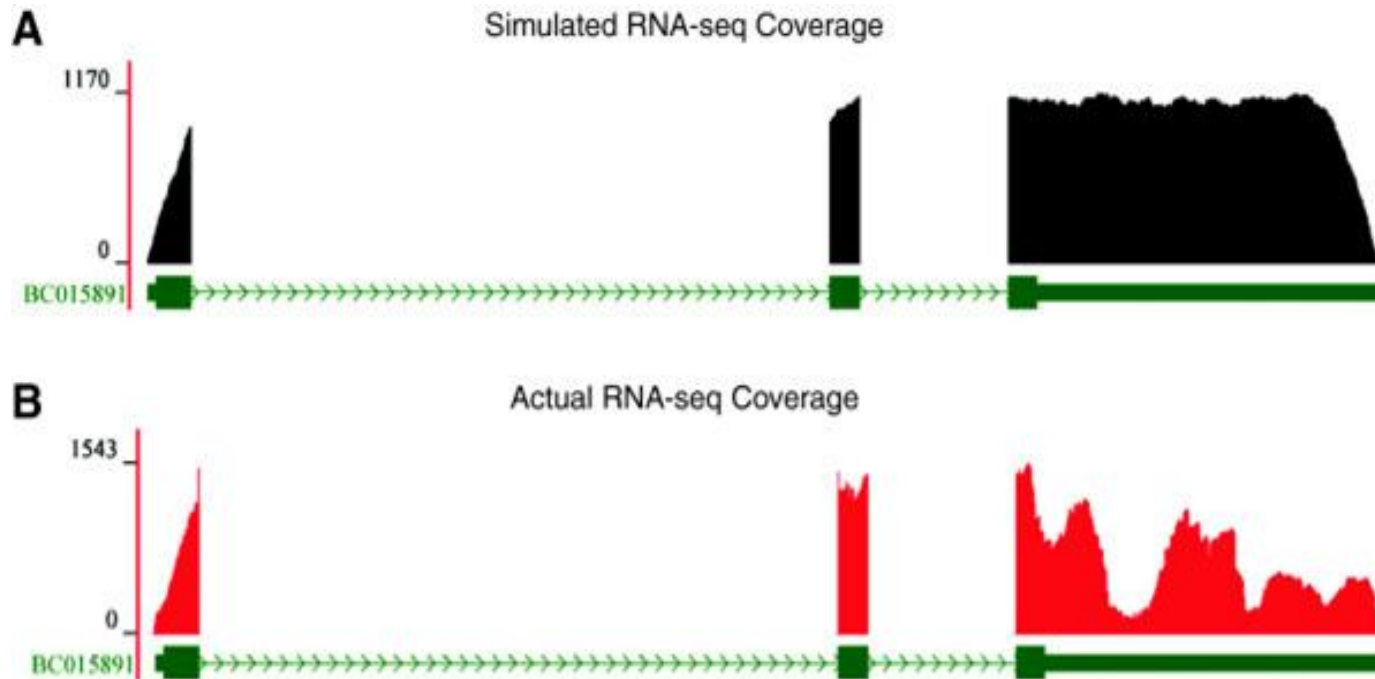
Biais : *random hexamer priming*



Orange = reads start sites

Blue = coverage

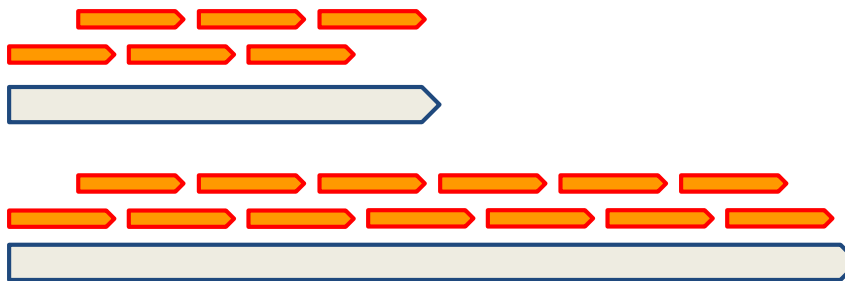
Biais : *Préparation de la librairie*



IVT-seq reveals extreme bias in RNA sequencing, Lahens et al 2014
<http://genomebiology.com/2014/15/6/R86>

Biais : longueur des transcrits

- o La capacité, en utilisant des **comptages** obtenus par **RNA-Seq**, à observer un transcrit comme étant **différentiellement exprimé** est **directement reliée** à sa **longueur**.
- o Pour un **même gène** ayant **deux isoformes**, l'une faisant la moitié de l'autre, exprimé en **même abondance dans deux conditions différentes** :
 - L'isoforme la plus courte sera deux fois moins « comptée » que la plus longue

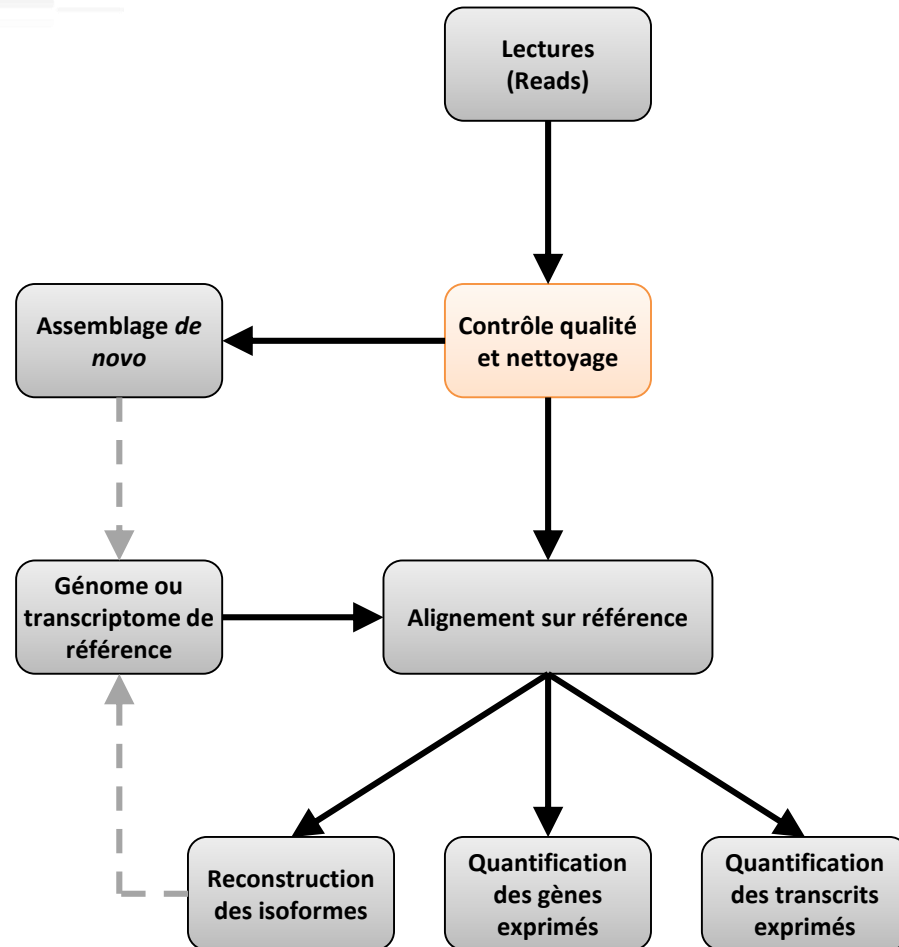




Biais : « mappabilité »

- o Les étapes bioinformatiques peuvent être **influencées** par :
 - La **qualité** de la **référence**
 - ✓ **assemblage**
 - ✓ **finition**
 - La **composition** de la **séquence**
 - ✓ **zones répétées**
 - La **qualité** de l'**annotation**

Workflow d'analyse RNA-Seq



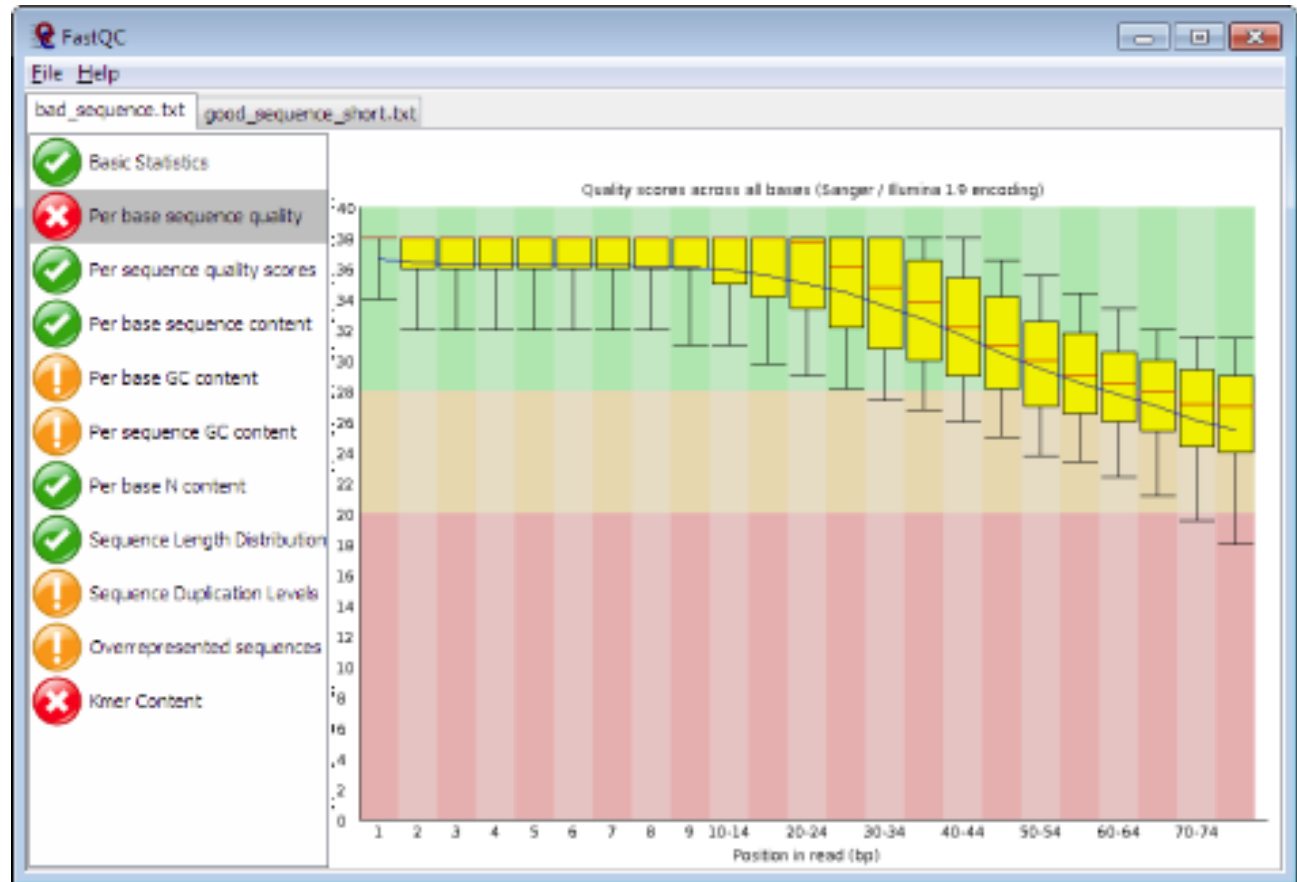
Contrôle qualité

Objectifs :

- ❖ Vérifier que les séquences sont **conformes au niveau de prestation attendu (taille, nombre, qualité,...)**
- ❖ Vérifier que les séquences peuvent **répondre au questions biologiques** posées :
 - **Biais techniques**
 - **Biais biologiques**
- ❖ Aider au paramètres pour le nettoyage des données

Contrôle qualité avec FastQC

❖ orienté DNA-Seq

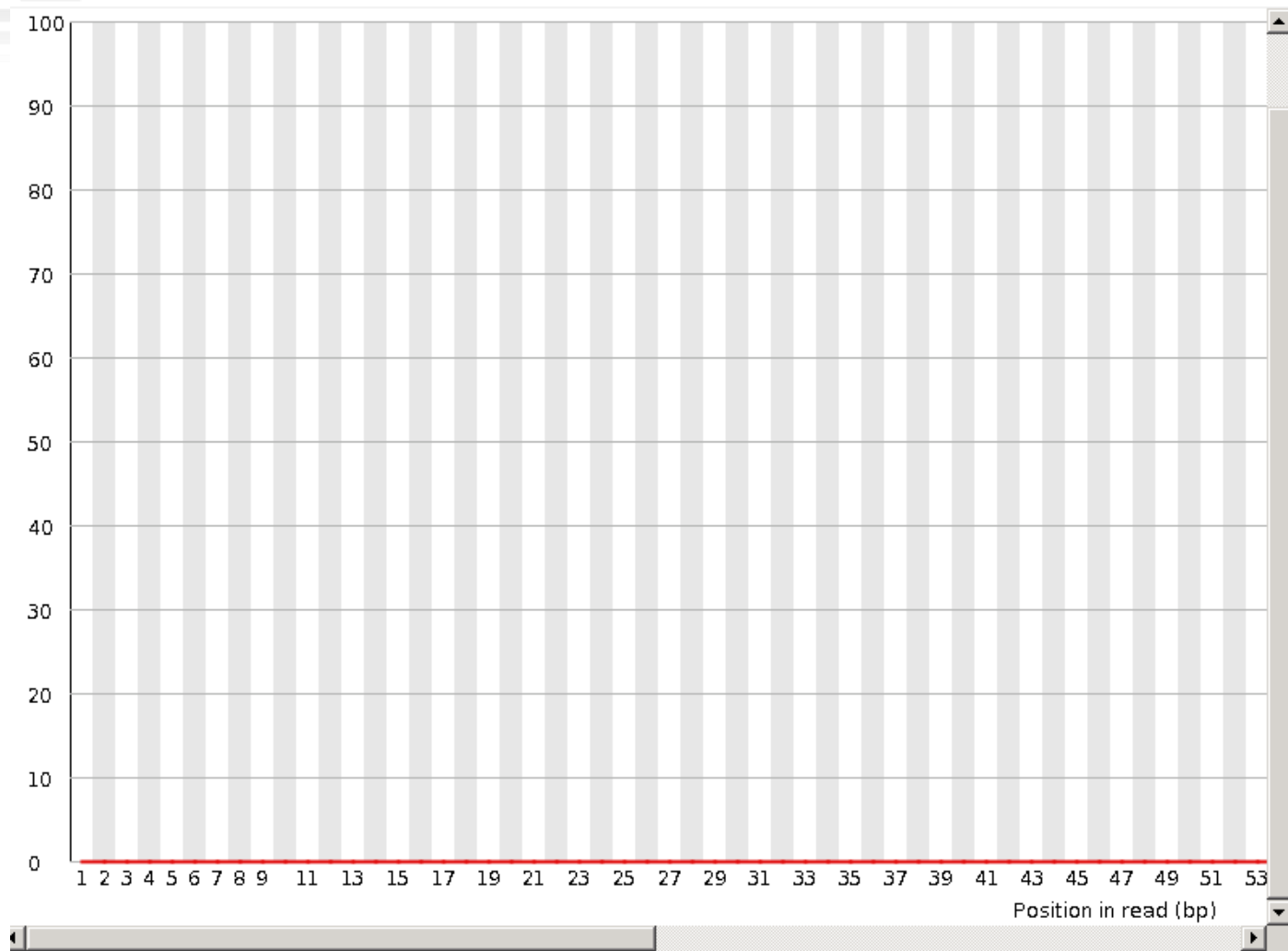


<http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc>

Contrôle qualité

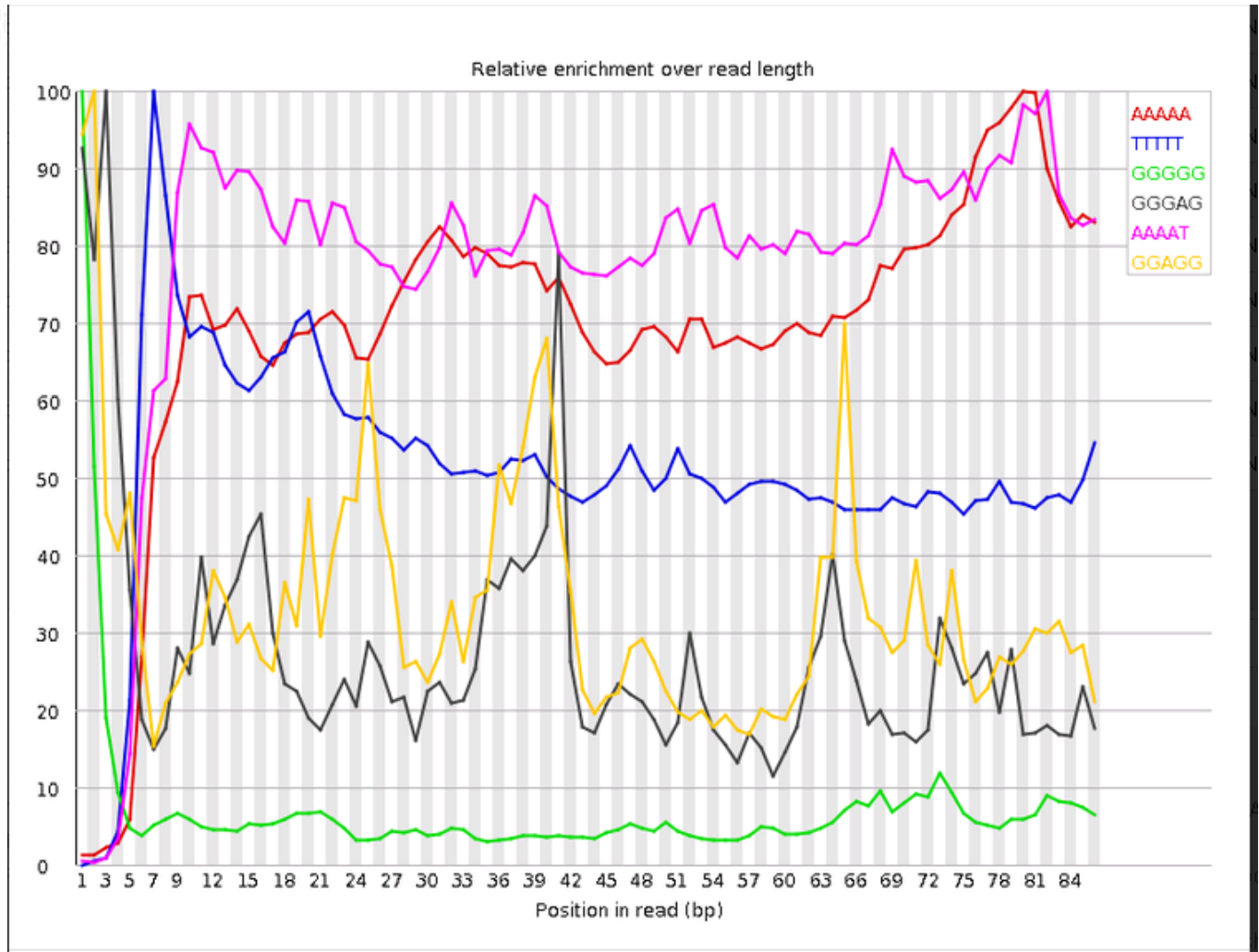


N content



Contrôle qualité

Kmer content



Contrôle qualité

Over-expressed sequences

! Overrepresented sequences

Sequence	Count	Percentage	Possible Source
TT	64433	0.20451959970239228	No Hit

[Back to summary](#)

❖ Détecter les adaptateurs



Nettoyage des données

❖ Nettoyage « optionnel »

❖ L'alignement permettra de supprimer les lectures

- De mauvaise qualité
- D'adaptateurs
- Contaminantes

❖ Les outils :

- Cutadapt : Nettoyage des adaptateurs & Tags
- Prinseq : Nettoyage des lectures de mauvaise qualité
- Sickle : Nettoyage des lectures de mauvaise qualité

Nettoyage des données

❖ Principe de Sickle :

- Traite les paires ensemble
- Fenêtre glissante : 10% de la taille des reads
- Calcul de la qualité moyenne des lectures

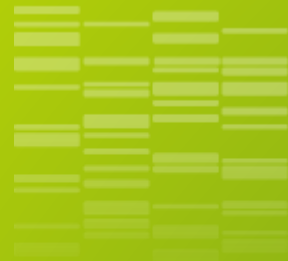
exemple : Longueur = 23

A	C	T	T	G	A	T	C	A	T	G	C	A	T	C	G	A	T	C	G	T	A	G
30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	30	25	20	18	18	10

Travaux pratiques

Présentation des objectifs

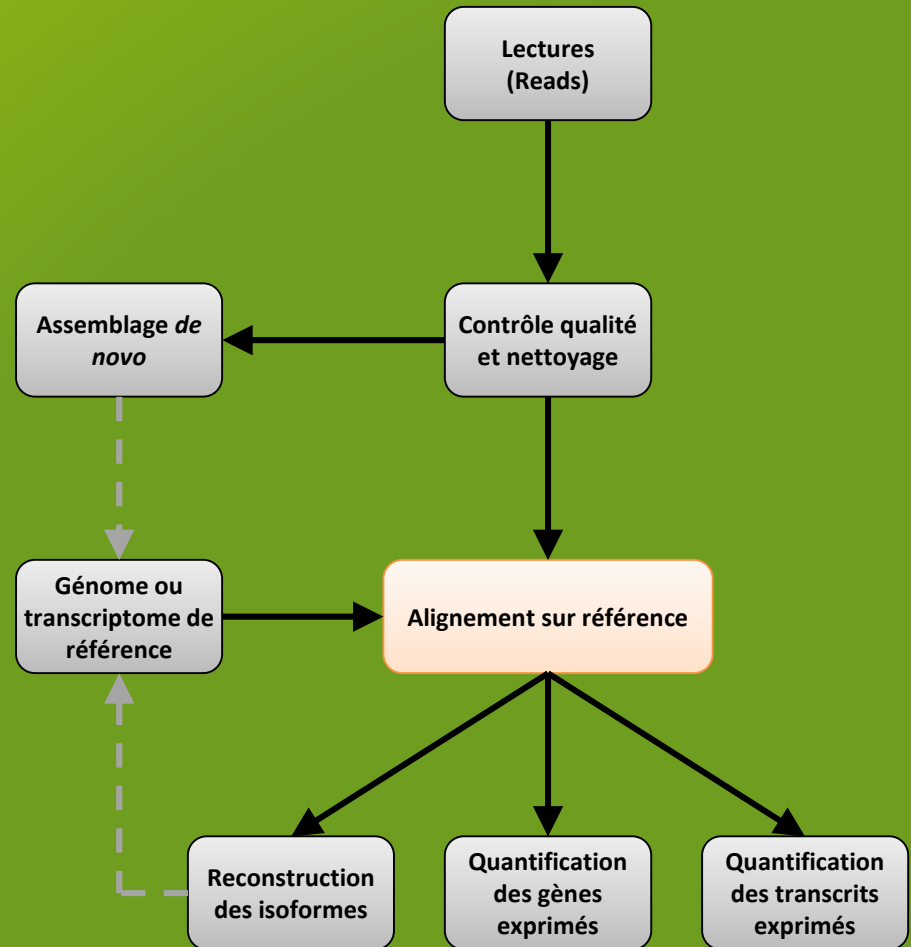
- ❖ **Aborder les différentes étapes indispensables au traitement bioinformatique de données RNA-Seq à travers un exemple issu de données réelles :**
- ❖ **Séquençage de la tomate :**
 - Wt : wild type, PAIRED
 - Mt : mutant type , PAIRED



03

MAPPING

et Visualisation





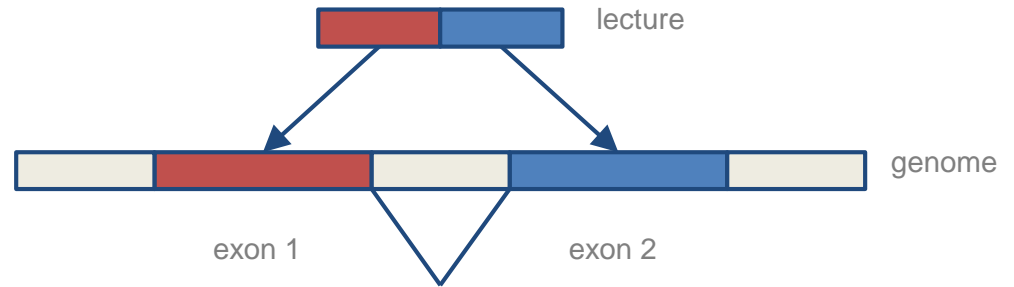
Alignement épissé

Objectifs :

- ❖ **Aligner** les **lectures** issues du séquençage de **dscDNA** (transcrits) sur le **génom**e, 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...**

Introduction

Définition



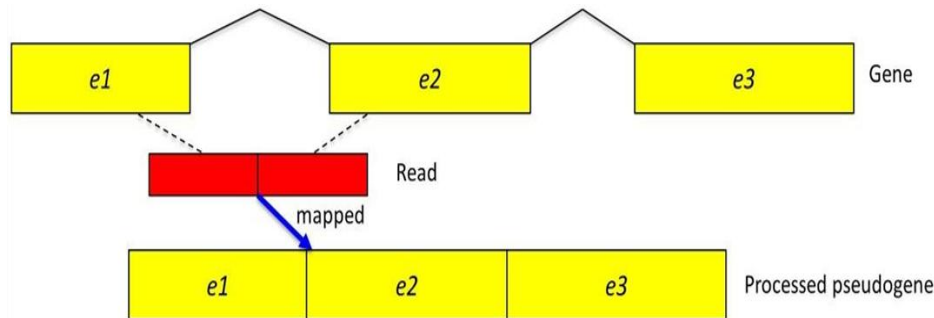
Le *mapping* est la *prédiction* du *locus* dont est originaire la lecture.

- **Prédiction :** chaque outil propose un/plusieurs locus. Ils peuvent ne pas être les bons.
- **Locus :** le résultat est un ensemble de positions génomiques (ex.: chr1:100..150)
- Mapping ARN \neq Mapping ADN
- Mapping \neq Alignement

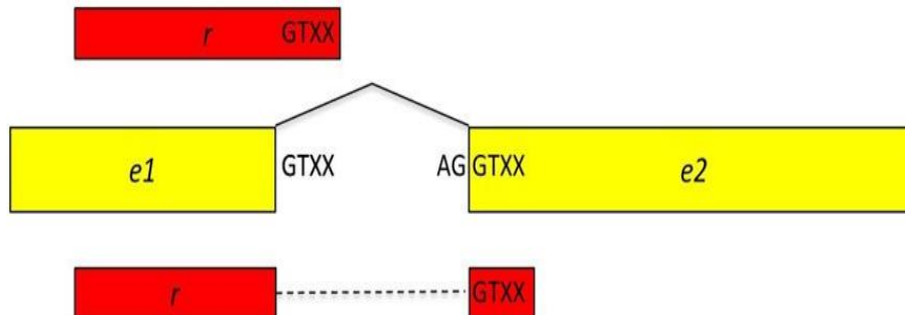
Les outils de mapping font de mauvais alignements (sauf aux jonctions).

Cas difficiles

- Beaucoup de différences : erreurs de séquençage ou locus muté
- Séquence répétée
- Lecture sur 3+ exons
- (Variante :) Gène ou pseudo-gène ?



- Fin de la lecture sur un exon propre



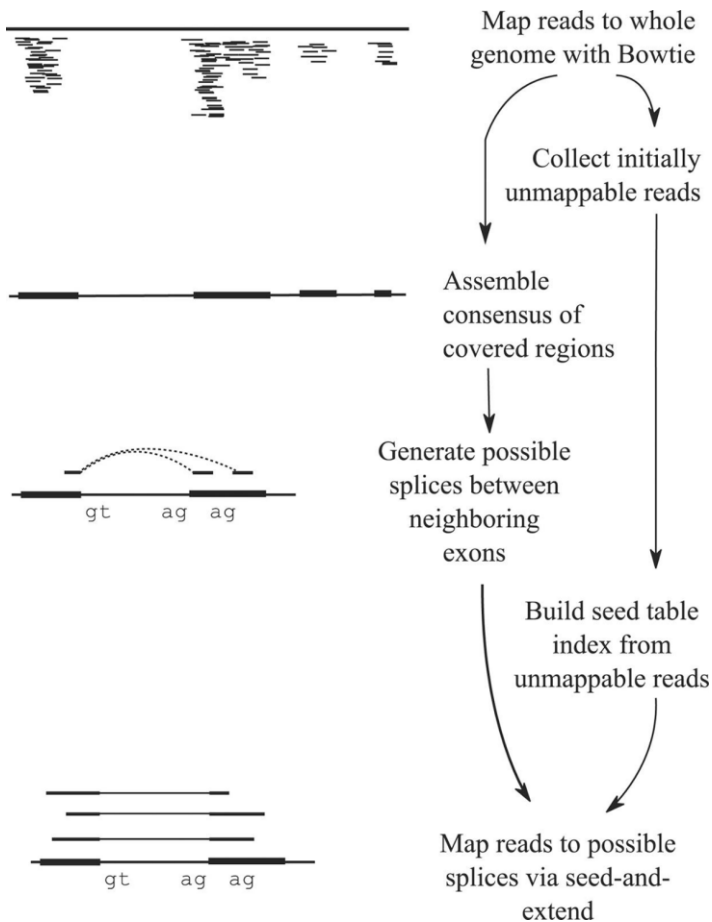
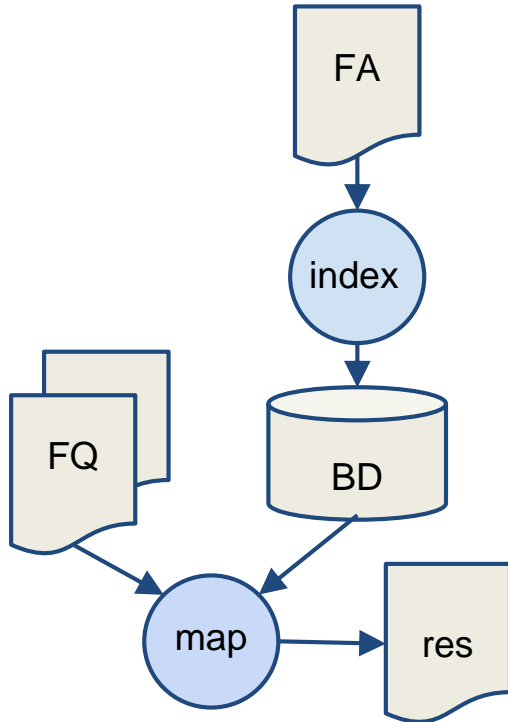
(Kim et al, Genome Biology, 2013)

- Lecture sur une jonction non-connue d'un gène peu exprimé

Tophat1

Étapes de mapping :

- ❖ *Indexation du génome 1 fois pour toutes*
- ❖ *Mapping des lectures utilise l'index*

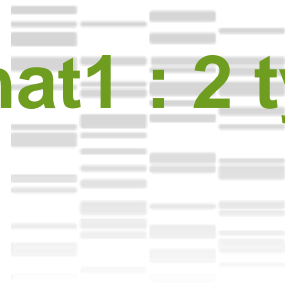


Utilise Bowtie pour mapper les lectures sur le génome.

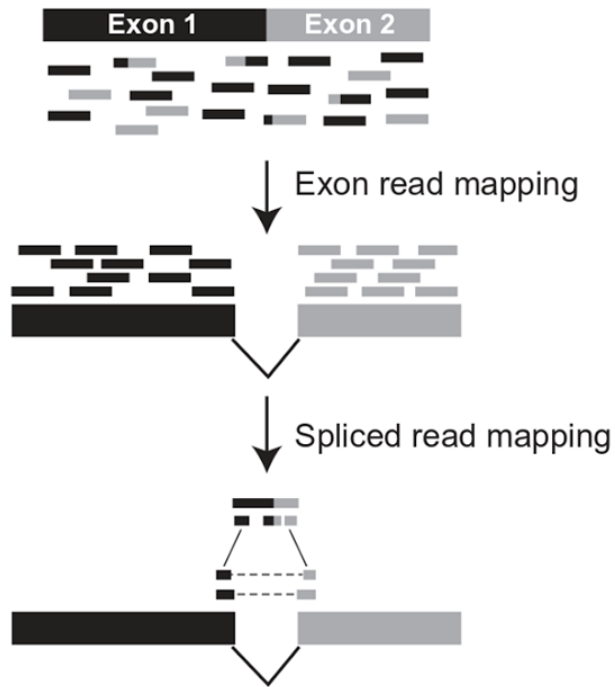
Problème pour les pseudo-gènes !

(Trapnell et al, Bioinformatics, 2009)

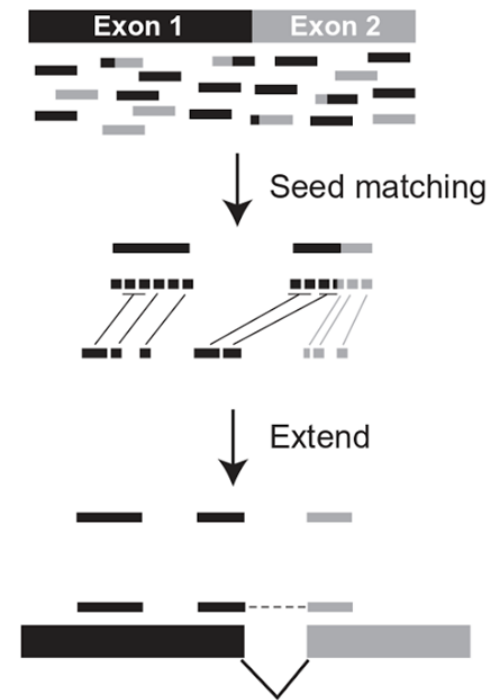
Tophat1 : 2 types d'algo



Exon-First Approach

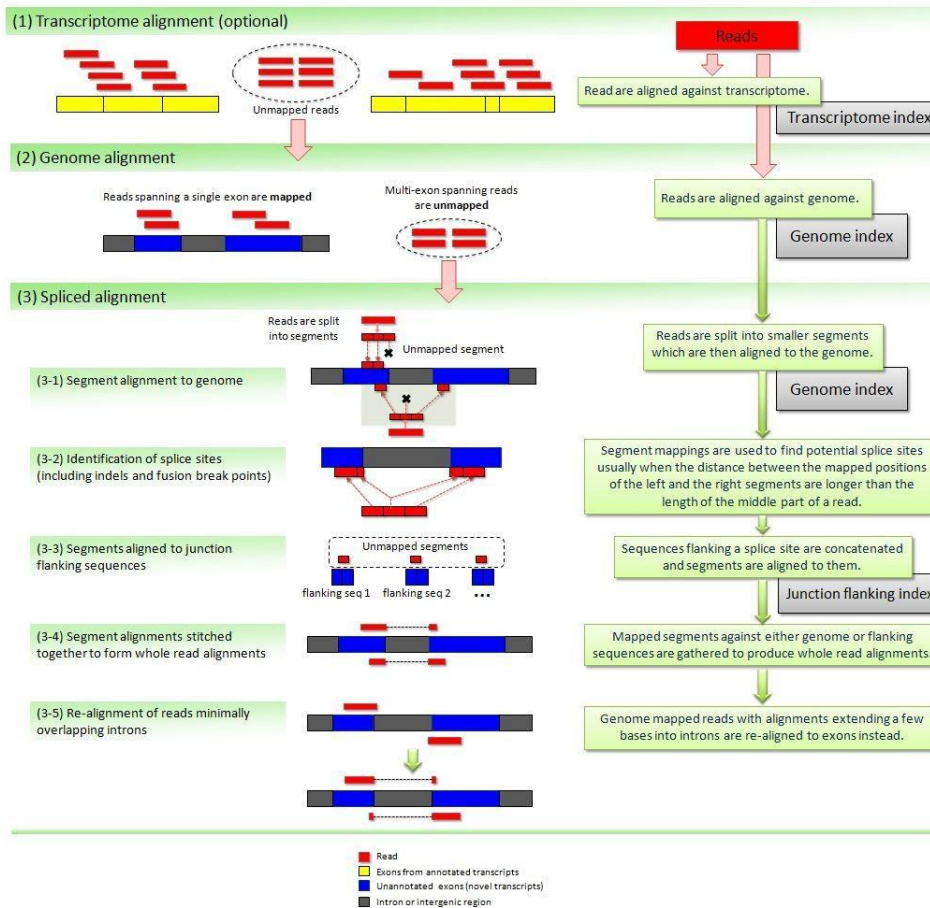


Seed-Extend Approach



(Garber et al, Nature Methods, 2011)

Tophat2

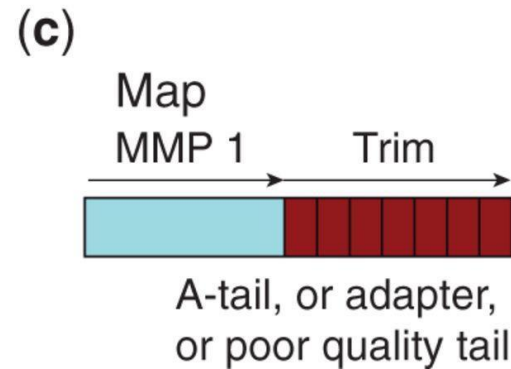
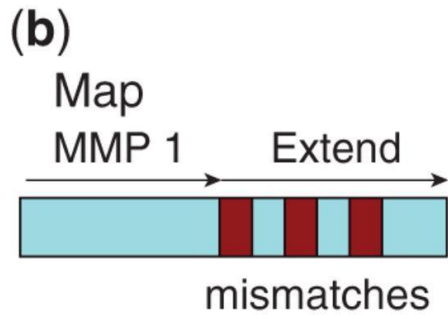
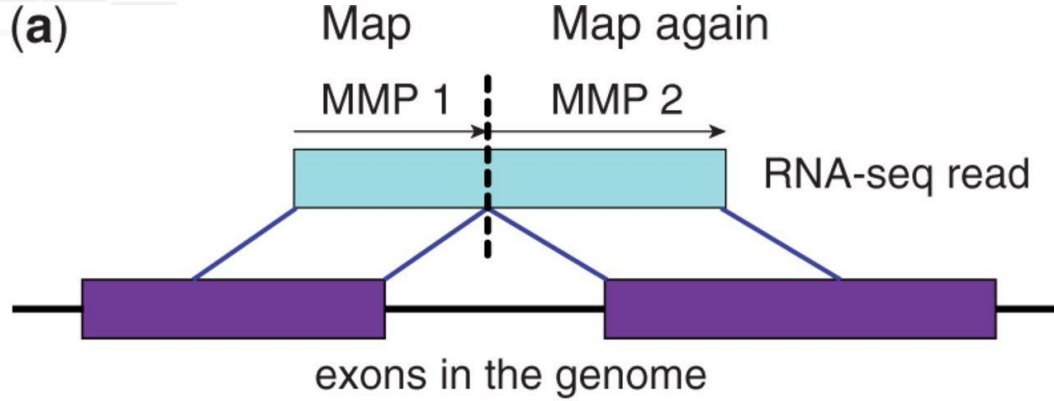
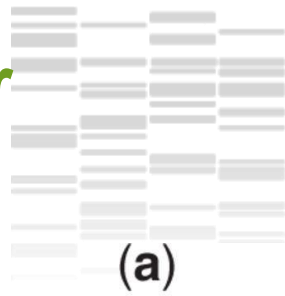


Tophat2 est constitué de beaucoup d'étape pour résoudre chaque cas difficile.

Chaque étape contient des heuristiques dont les paramètres sont à fixer.

(Kim et al, Genome Biology, 2013)

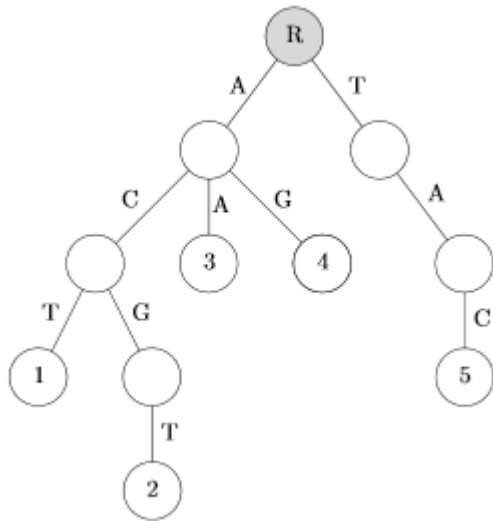
Star



(Dobin et al, Bioinformatics, 2011)

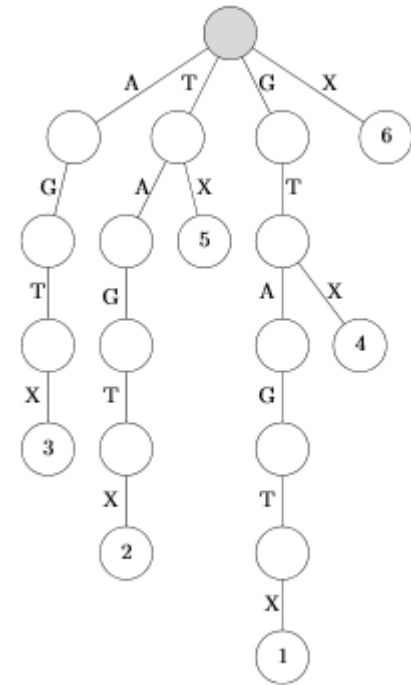
(RNA-)STAR

Aligneurs index BWT
(BWA, Bowtie, SOAP)



ACT, ACGT, AA, AG, et TAC.

STAR



GTAGT.

Outils existants

- ❖ **Tophat2 (le plus utilisé, le plus suivi)**
- ❖ **Star (runner-up)**
- ❖ **Crac (français !)**
- ❖ **GSNAP**
- ❖ **RUM**
- ❖ **MapSplice**
- ❖ **Gem**
- ❖ **...**

Outils existants

La plupart des outils

- ❖ utilise des sites de jonctions donnés par l'utilisateur pour “s'aider”
- ❖ suppose des sites canoniques GT-AG

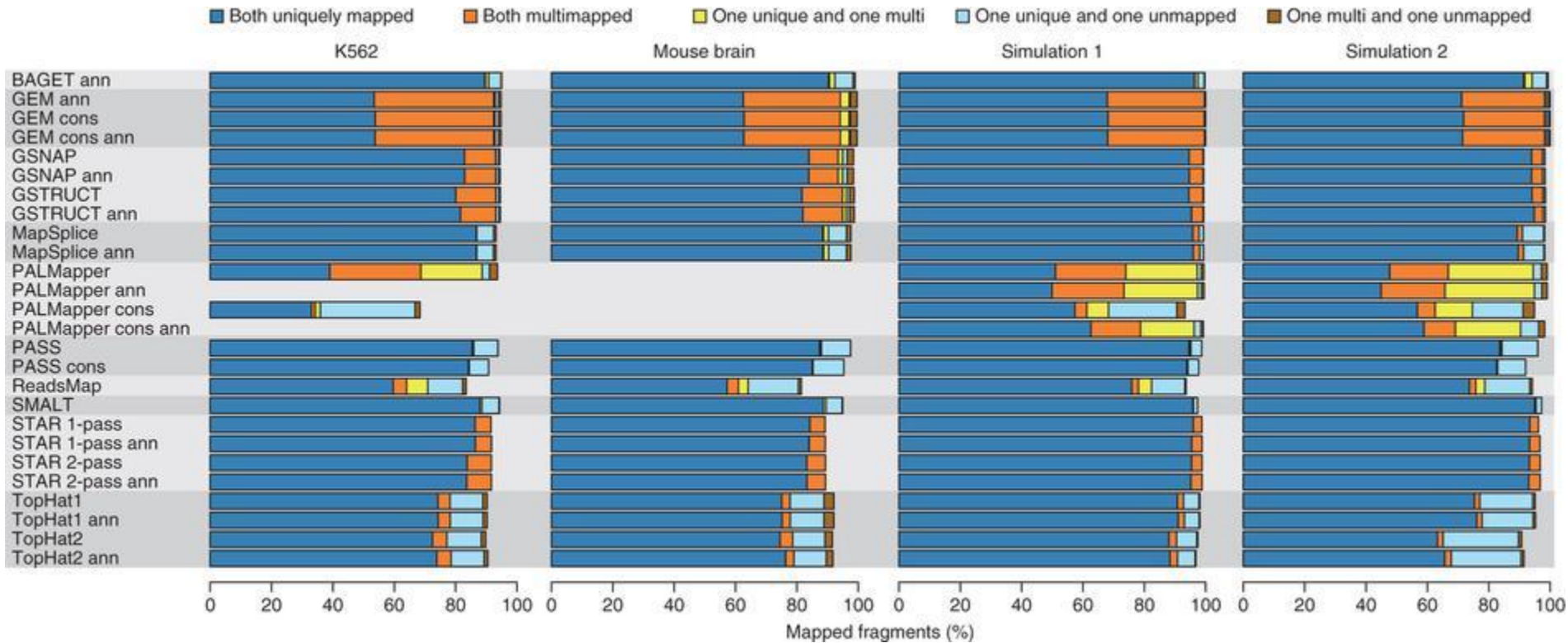
Comment évaluer un outil ?

- ❖ Sensibilité (mappe le plus de lectures)
- ❖ Spécificité (ne se trompe pas)
- ❖ ... sur les lectures et sur les jonctions
- ❖ Temps
- ❖ Mémoire

En général, les critères sont contradictoires.

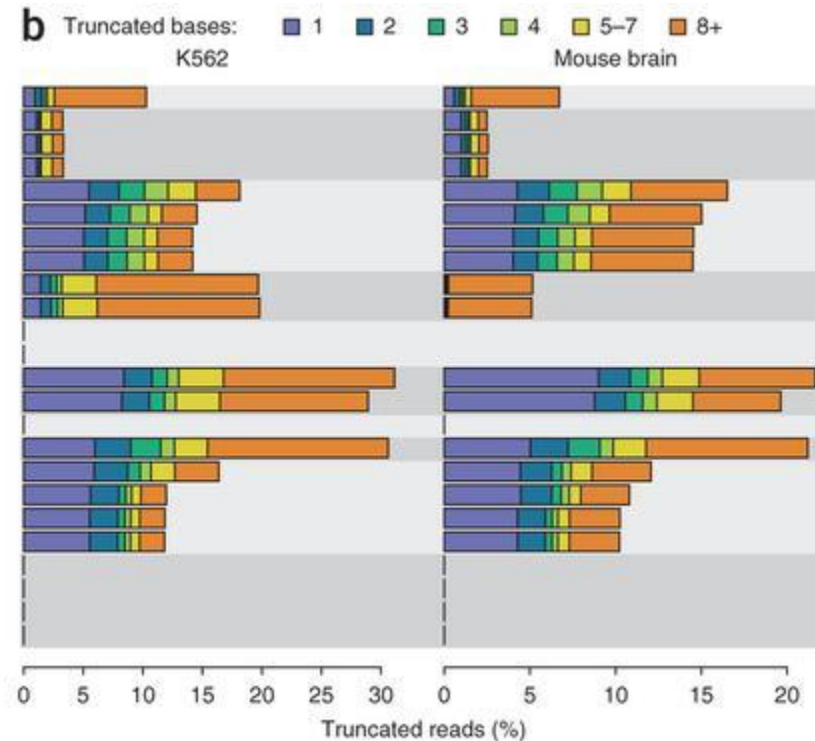
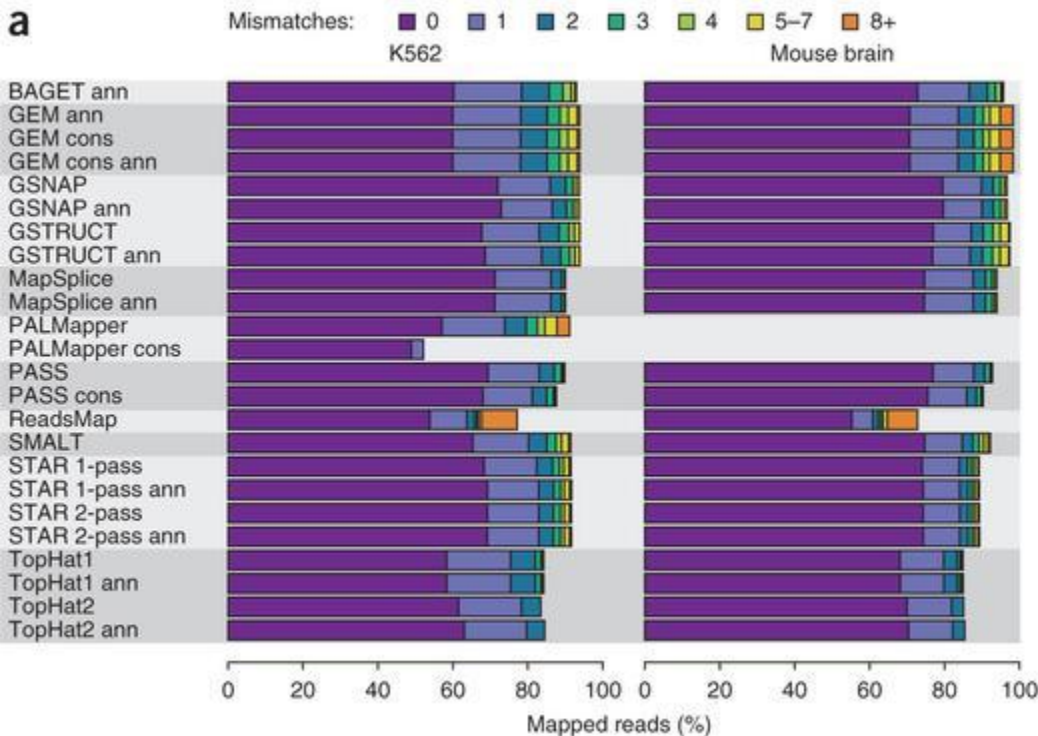
RGASP 3

The RNA-seq Genome Annotation Assessment Project
(Engström et al., Nature Methods, 2013)



RGASP 3

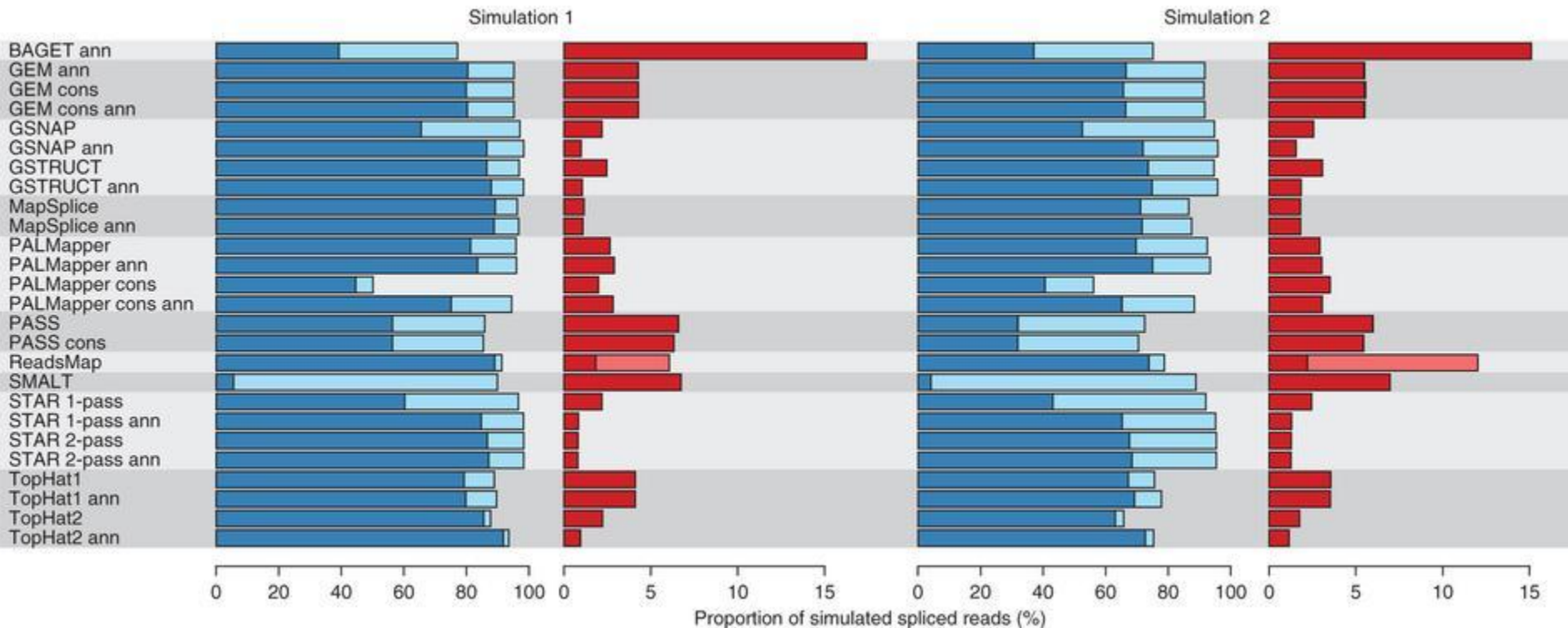
The RNA-seq Genome Annotation Assessment Project
(Engström et al., Nature Methods, 2013)



RGASP 3

The RNA-seq Genome Annotation Assessment Project
(Engström et al., Nature Methods, 2013)

■ Perfectly mapped
 ■ Part correctly mapped
 ■ Mapped, no base correct
 ■ No base correctly mapped but intersecting correct location



RGASP 3

The RNA-seq Genome Annotation Assessment Project
(Engström et al., Nature Methods, 2013)

Les phrases clés

« Mapping properties are largely dependent on software algorithms even when the genome and transcriptome are virtually identical »

« Exon detection results based on K562 data were similar for GEM, GSNAP, GSTRUCT, MapSplice, STAR and TopHat »

RGASP 3

The RNA-seq Genome Annotation Assessment Project
(Engström et al., Nature Methods, 2013)

STAR

+

-

-

-

VS

de lectures alignées

de lectures correctement alignées

sensibilité aux variations

sensibilité aux annotations

TopHat

-

+

+

+

Alignement : données initiales

- ❖ Lectures (brutes / nettoyées ?)
- ❖ **Génome de référence éventuellement annoté :**
 - Séquence nucléique (fasta)
 - Annotation structurale (GTF)
- ❖ **Où trouver un génome et un transcriptome de référence ?**
 - Ensembl
 - NCBI
- ❖ **Exo : trouver votre génome préféré et son annotation.**

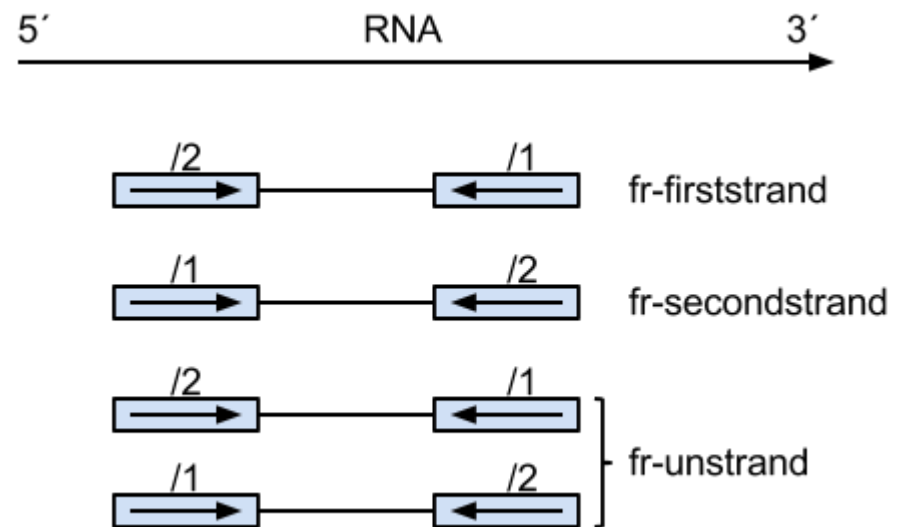
Tophat

❖ En entrée :

- lectures (.fastq)
- index bowtie2 de la référence
- annotation structurale du génome (.gtf) [optionnel]
- *jonction* (.bed) [optionnel]
- *insertions / délétions* (.bed) [optionnel]

❖ Les options :

- si paired, type de librairie.



Attention aux paramètres par défaut !

❖ Dans le manuel :

Please Note TopHat has a number of parameters and options, and their default values are tuned for processing mammalian RNA-Seq reads.

If you would like to use TopHat for another class of organism, we recommend setting some of the parameters with more strict, conservative values than their defaults.

Usually, setting the maximum intron size to 4 or 5 Kb is sufficient to discover most junctions while keeping the number of false positives low.

<http://tophat.cbcb.umd.edu/manual.shtml>

Format GTF : Gene Transfert Format

- ❖ **Dérivé** du format généraliste GFF (General Feature Format)
- ❖ Contient l'**annotation structurale** du **génom**e (gène, transcrits)

```
<seqname> <source> <feature> <start> <end> <score> <strand> <frame> [attributes] [comments]
```

```
3R protein_coding exon 380 509 . + . gene_id "FBgn0037213"; transcript_id "FBtr0078961";  
exon_number "1"; gene_name "CG12581"; transcript_name "CG12581-RB";
```

❖ Le champ attribut doit :

- Commencer par le *gene_id* : identifiant **unique** du gène
- Être suivi par *transcript_id* : identifiant **unique** du transcrit prédit

- ❖ Les identifiants du chromosome (**Fasta** et **1^{ère} colonne du GTF**) doivent être les **mêmes**

<http://genome.ucsc.edu/FAQ/FAQformat.html#format4>

TP : Assemblage avec Tophat

Tophat for Illumina (version 1.0.0)

Your RNA-Seq FASTQ file (read 1):
3: ERR022488_read1 ▼

Your RNA-Seq FASTQ file (read 2):
4: ERR022488_read2 ▼

Select a reference genome:
Danio rerio Zv9 62 chr 22 ▼

Number of threads used to align reads:
8

Maximum intron length:
5000

Expected (mean) inner distance between mate pairs:
200

Your RNA-seq FASTQ file are zipped:
 Yes
Please check this option if your files are zipped.

GTF file available:
Yes ▼
Do you have a gtf file available ?

Your GTF file:
5: http://genoweb.toulouse.inra.fr/~formation/OLD_4_Galaxy_RNAseq/data/reference/Danio_rerio.gtf ▼

Library type:
fr-unstranded ▼

Execute

❖ **Lancer Tophat maintenant ! Avant de continuer la présentation.**

<http://tophat.cbcb.umd.edu/manual.shtml>

Alignement : Format SAM/BAM

- ❖ Le partage des données est un problème majeur dans le projets “1000 génomes”
- ❖ Capturez toute l'information critique sur les données de NGS dans un seul fichier indexé et comprimé
- ❖ Alignement format générique
- ❖ Prise en charge reads de taille variable (454 - Solexa - Solid ... PacBio)
- ❖ Flexible dans le style , de taille compacte , efficace en accès aléatoire

Website :

<http://samtools.sourceforge.net>

Paper :

Li H.*, Handsaker B.*, Wysoker A., Fennell T., Ruan J., Homer N., Marth G., Abecasis G., Durbin R. and 1000 Genome Project Data Processing Subgroup (2009) The Sequence alignment/map (SAM) format and SAMtools. *Bioinformatics*, 25, 2078-9. [PMID: 19505943]

Alignement : Format SAM

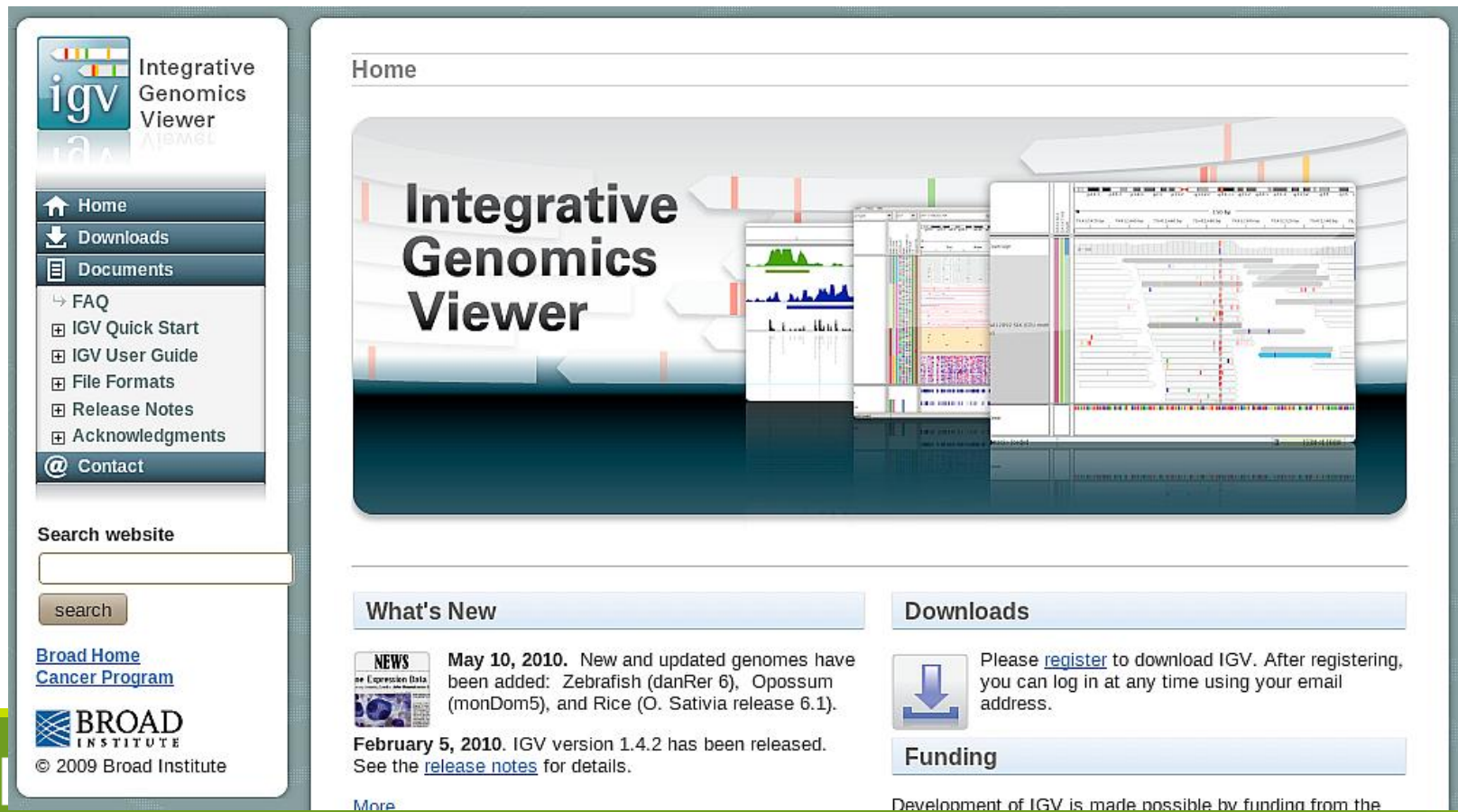


Quelles informations doivent être stockées dans un fichier d'alignement SAM ?

http://genoweb.toulouse.inra.fr/~formation/2_Galaxy_SG_S-SNP/.SAMformat/sam.html

Visualisation des alignements avec IGV

- ❖ IGV : Integrative Genomics Viewer
- ❖ Website : <http://www.broadinstitute.org/igv>



Integrative Genomics Viewer

Home

Home
Downloads
Documents
FAQ
IGV Quick Start
IGV User Guide
File Formats
Release Notes
Acknowledgments
Contact

Search website

search

[Broad Home Cancer Program](#)

BROAD INSTITUTE
© 2009 Broad Institute

What's New

NEWS May 10, 2010. New and updated genomes have been added: Zebrafish (danRer 6), Opossum (monDom5), and Rice (O. Sativa release 6.1).

February 5, 2010. IGV version 1.4.2 has been released. See the [release notes](#) for details.

Downloads

Please [register](#) to download IGV. After registering, you can log in at any time using your email address.

Funding

Development of IGV is made possible by funding from the

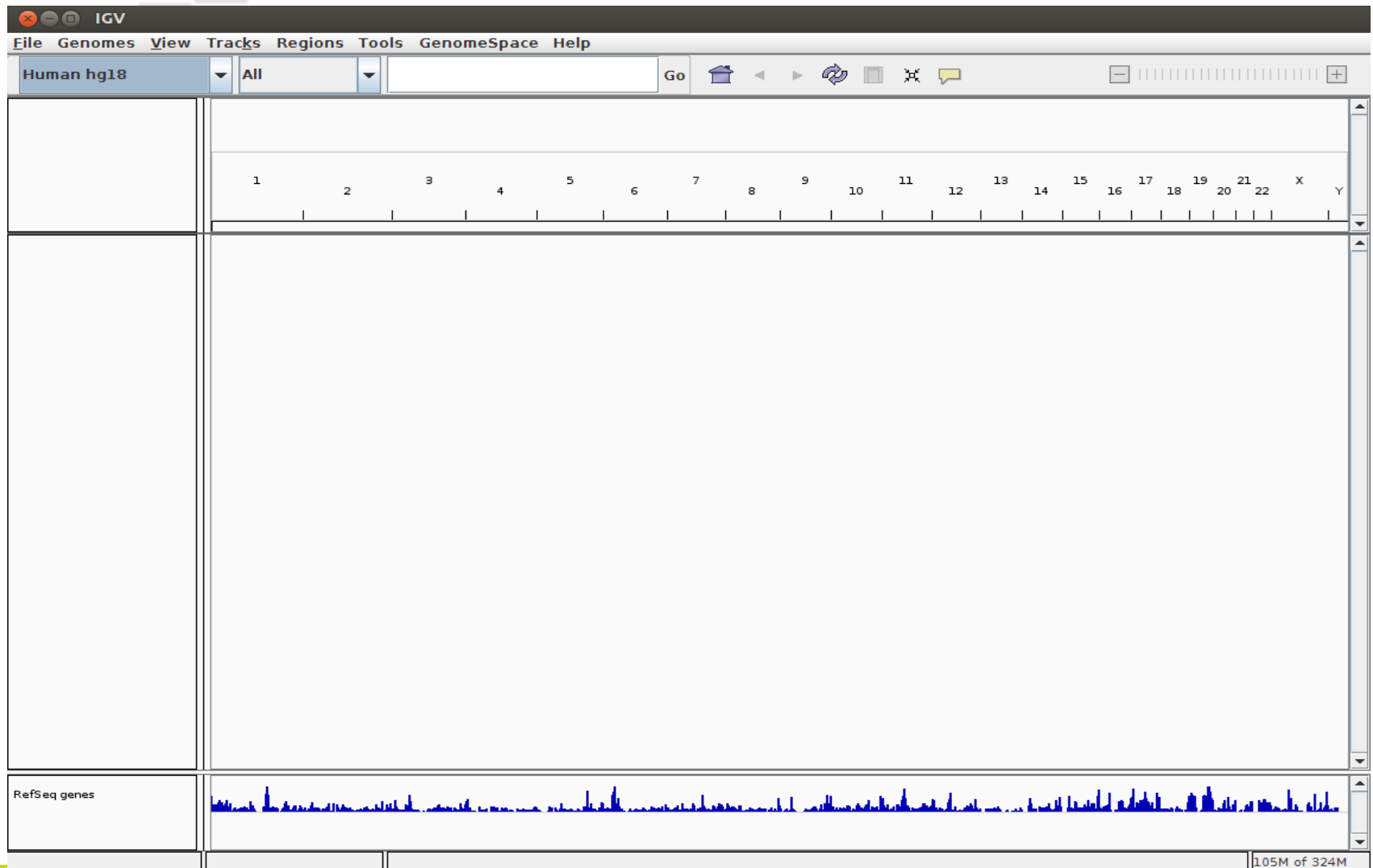
Visualisation des alignements avec IGV

- ❖ High-performance visualization tool
- ❖ Interactive exploration of large, integrated datasets
- ❖ Supports a wide variety of data types
- ❖ Documentations
- ❖ Developed at the Broad Institute of MIT and Harvard

File Formats

- [File Extension Identifies Format](#)
- [Recommended File Formats](#)
- [BAM](#)
- [BED](#)
- [CBS](#)
- [CN](#)
- [Cytoband](#)
- [FASTA](#)
- [GCT](#)
- [genePred](#)
- [GFF](#)
- [GISTIC](#)
- [HDF5](#)
- [IGV](#)
- [LOH](#)
- [Birdsuite Files](#)
- [MUT](#)
- [RES](#)
- [SAM](#)
- [Sample Information](#)
- [SEG](#)
- [SNP](#)
- [TAB](#)
- [TDF](#)
- [Track Line](#)
- [Type Line](#)
- [WIG](#)

Visualisation des alignements avec IGV



IGV : Chargement de la référence

The screenshot shows the IGV application window with the 'File' menu open. The 'Load Genome from File...' option is highlighted. A 'Load Genome' dialog box is open, displaying a file browser interface. The dialog box shows a list of files and folders, including 'bin', 'boot', 'cdrom', 'dev', 'etc', 'home', 'lib', 'lib64', 'lost+found', 'media', 'mnt', 'proc', 'root', 'run', 'sbin', 'selinux', 'srv', 'sys', 'tmp', 'usr', 'var', 'initrd.img', 'initrd.img.old', 'vmlinuz', and 'vmlinuz.old'. The 'Nom du fichier' field is empty, and the 'Fichiers de type' dropdown is set to 'Tous les fichiers'. The 'Ouvrir' and 'Annuler' buttons are visible at the bottom of the dialog box. A text box on the right side of the dialog box contains the text: 'Select a fasta file, the index .fai must exists in the same directory'. The background shows the IGV interface with a track labeled 'RefSeq genes' and a genomic track.

Select a fasta file, the index .fai must exists in the same directory

IGV : Chargement de l'annotation

The screenshot shows the IGV interface with the 'File' menu open. The 'Load from File...' option is highlighted, and a red box is drawn around the 'File' menu and the 'chr1' dropdown. The main display area is mostly empty, with a large text overlay in the center that reads 'Charger le fichier GTF, pour avoir la piste d'annotation'. A blue arrow points from this text towards the 'Load from File...' menu item. The interface includes a top menu bar with 'File', 'Genomes', 'View', 'Tracks', 'Regions', 'Tools', 'GenomeSpace', and 'Help'. Below the menu is a toolbar with various icons for navigation and zooming. The main display area shows a genomic track with a scale from 90 mb to 98 mb and a 10 mb zoomed-in view. The bottom track shows 'RefSeq genes' with labels for PKN2, GBP4, LRRC8D, BARHL2, HFM1, BRDT, GF1, MIF2, BCAR3, ABCD3, ALG14, PTBP2, DPVD, and MIR2682. The status bar at the bottom indicates '2 tracks loaded', 'chr1:88 899 520', and '106M of 480M'.

Charger le fichier GTF, pour avoir la piste d'annotation

IGV : Chargement des alignements

The screenshot displays the IGV interface with a file selection dialog box open. The dialog box is titled "Rechercher dans : CORRECTION" and lists the following files:

- bam.intervals
- empty.vcf
- empty.vcf.idx
- ERR000017.bam
- ERR000017.bam.bai
- ERR000017.fastq
- ERR000017.sai
- ERR000017.sam
- ERR000017_rmdup.bam
- ERR000017_rmdup.bam.bai
- ERR000017_rmdup_realign.bai
- ERR000017_rmdup_realign.bam
- ERR000017_rmdup_realign_re...
- ERR000017_rmdup_realign_re...
- ERR003037.bam
- ERR003037.bam.bai
- ERR003037.fastq
- ERR003037.sai
- ERR003037.sam
- ERR003037_rmdup.bam
- ERR003037_rmdup.bam.bai
- ERR003037_rmdup_realign.bai
- ERR003037_rmdup_realign.ba...
- ERR003037_rmdup_realign_re...

The file "ERR003037.bam" is highlighted with a red box. Below the file list, the "Nom de fichier" field contains "ERR000017.bam" "ERR003037.bam" and the "Fichiers du type" dropdown is set to "Tous les fichiers". The "Ok" and "Annuler" buttons are visible at the bottom of the dialog.

Select a bam file, the index .bai must exists in the same directory

The background shows the IGV interface with the following tracks visible:

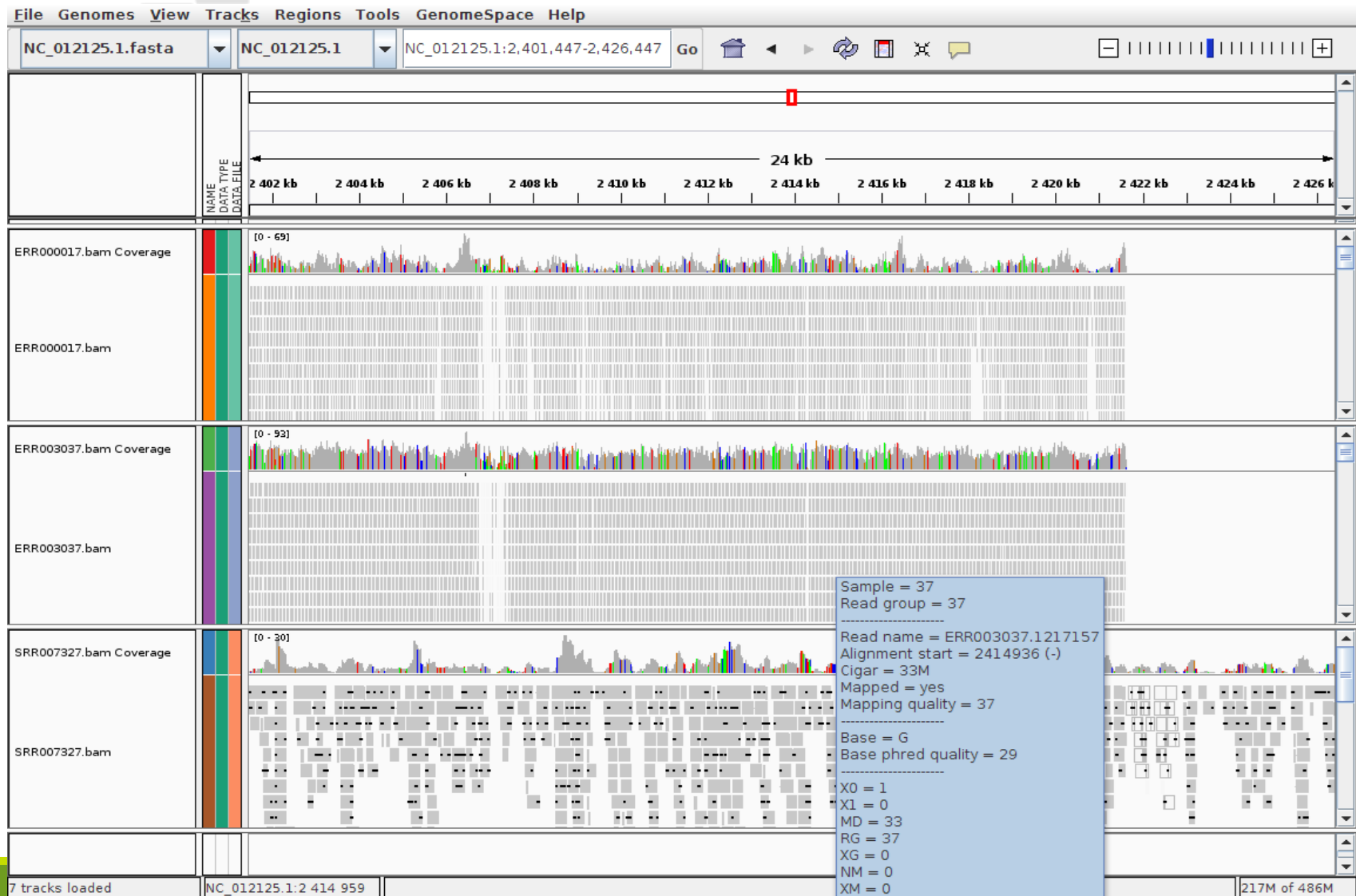
- File menu (highlighted in red)
- Genomes menu
- View menu
- Tracks menu
- Regions menu
- Tools menu
- GenomeSpace menu
- Help menu
- Load from File... (highlighted in red)
- Load from URI
- Load from Service (highlighted in red)
- Load from DAS...
- New Session...
- Open Session...
- Save Session...
- Save Image ...
- Exit
- Rechercher dans : CORRECTION
- ERR000017.bam
- ERR000017.bam.bai
- ERR000017.fastq
- ERR000017.sai
- ERR000017.sam
- ERR000017_rmdup.bam
- ERR000017_rmdup.bam.bai
- ERR000017_rmdup_realign.bai
- ERR000017_rmdup_realign.bam
- ERR000017_rmdup_realign_re...
- ERR003037.bam
- ERR003037.bam.bai
- ERR003037.fastq
- ERR003037.sai
- ERR003037.sam
- ERR003037_rmdup.bam
- ERR003037_rmdup.bam.bai
- ERR003037_rmdup_realign.bai
- ERR003037_rmdup_realign.ba...
- ERR003037_rmdup_realign_re...
- Nom de fichier : "ERR000017.bam" "ERR003037.bam"
- Fichiers du type : Tous les fichiers
- Ok
- Annuler
- RefSeq genes
- PKN2
- GBP4
- LRRC8D
- BARHL2
- HFM1
- BRDT
- GFI1
- MTF2
- BCAR3
- ABCD3
- ALG14
- PTBP2
- DPVD
- MIR2682
- 2 tracks loaded
- chr1:88 899 520
- 106M of 480M

IGV : Chargement des alignements

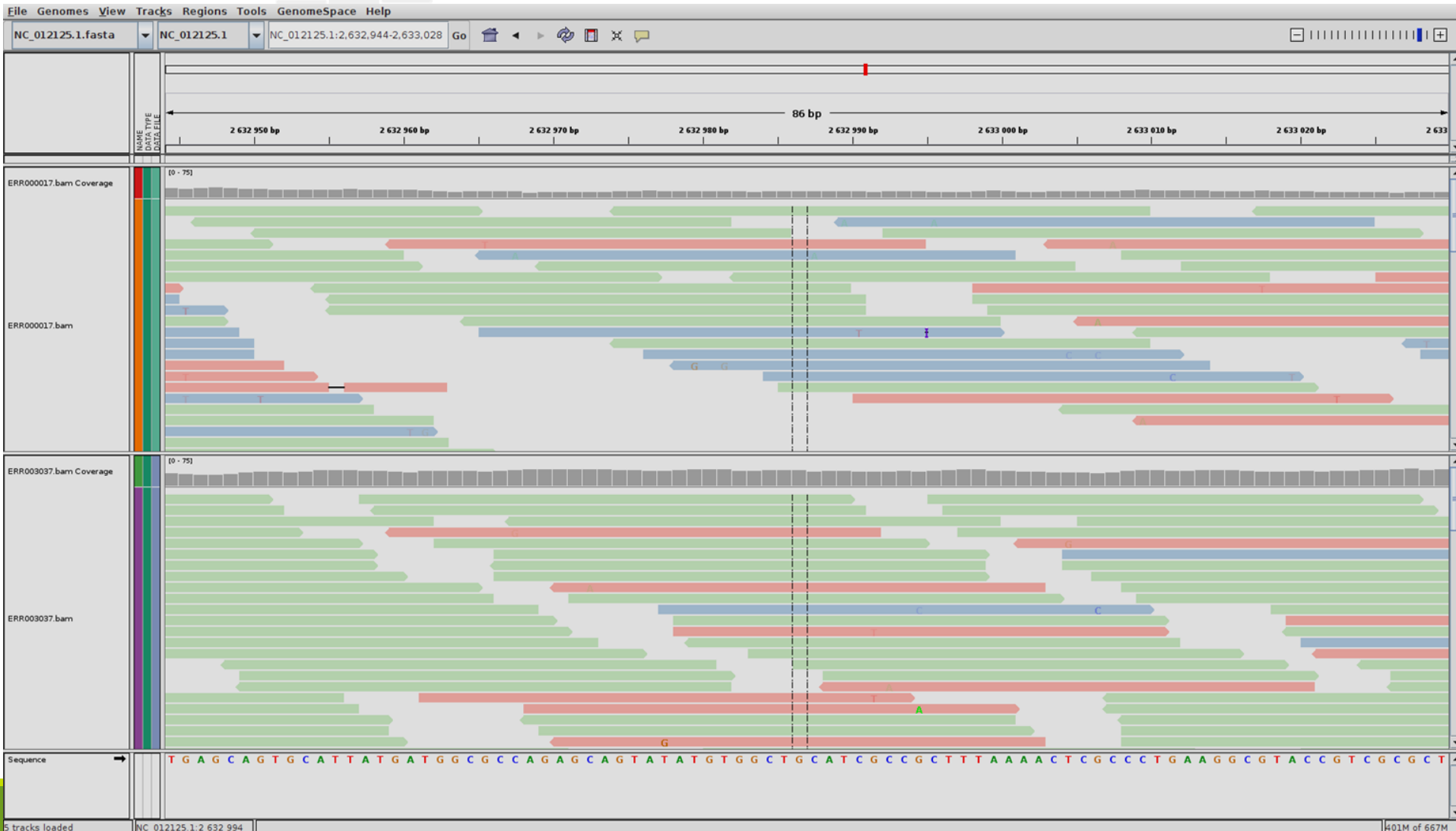
The screenshot displays the IGV interface with the following components:

- Menu Bar:** File Genomes View Tracks Regions Tools GenomeSpace Help
- File Path:** NC_012125.1.fasta | NC_012125.1 | NC_012125.1 | Go
- Genomic Scale:** A horizontal axis showing a 4,822 kb region, with markers at 1,000 kb, 2,000 kb, 3,000 kb, and 4,000 kb.
- Track 1:** ERR000017.bam Coverage (range [0 - 69]) and ERR000017.bam. The track content is empty with the text "Zoom in to see alignments."
- Track 2:** ERR003037.bam Coverage (range [0 - 93]) and ERR003037.bam. The track content is empty with the text "Zoom in to see alignments."
- Track 3:** SRR007327.bam Coverage (range [0 - 30]) and SRR007327.bam. The track content is empty with the text "Zoom in to see alignments."
- Bottom Bar:** 7 tracks loaded | NC_012125.1:26 069 | 200M of 486M

IGV : Chargement des alignements

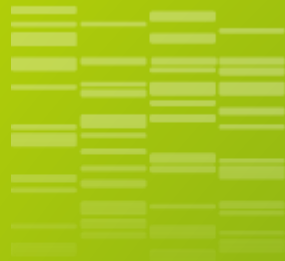


IGV : Chargement des alignements



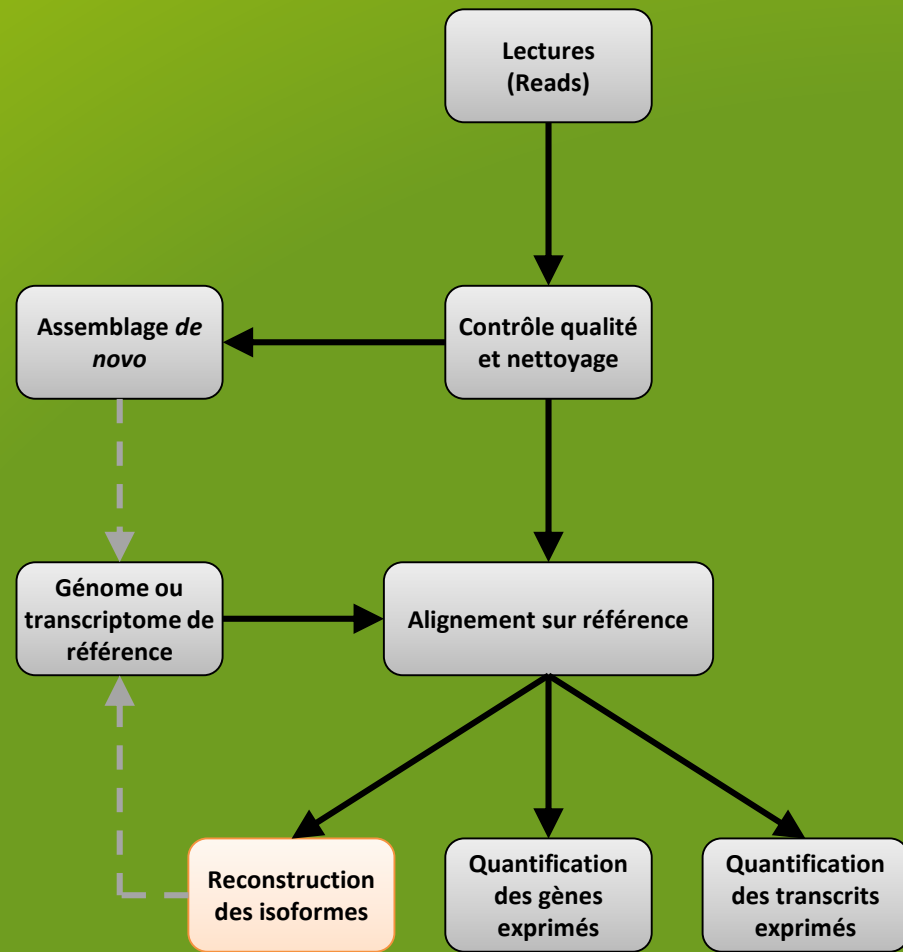


TP : Visualisation avec IGV



_04

Reconstruction de transcrit



Cufflinks

❖ Pipeline / suite logiciel de traitement RNA-Seq :

- assemble les transcrits (cufflinks)
- quantifie l'abondance des transcrits (cufflinks)
- compare les annotations des transcrits (cuffcompare)
- analyse l'expression différentielle des transcrits (cuffdiff)

nature
biotechnology

LETTERS

Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation

Cole Trapnell¹⁻³, Brian A Williams⁴, Geo Pertea², Ali Mortazavi⁴, Gordon Kwan⁴, Marijke J van Baren⁵, Steven L Salzberg^{1,2}, Barbara J Wold⁴ & Lior Pachter^{3,6,7}

<http://cufflinks.cbcb.umd.edu/>

The element of the model



Gene :

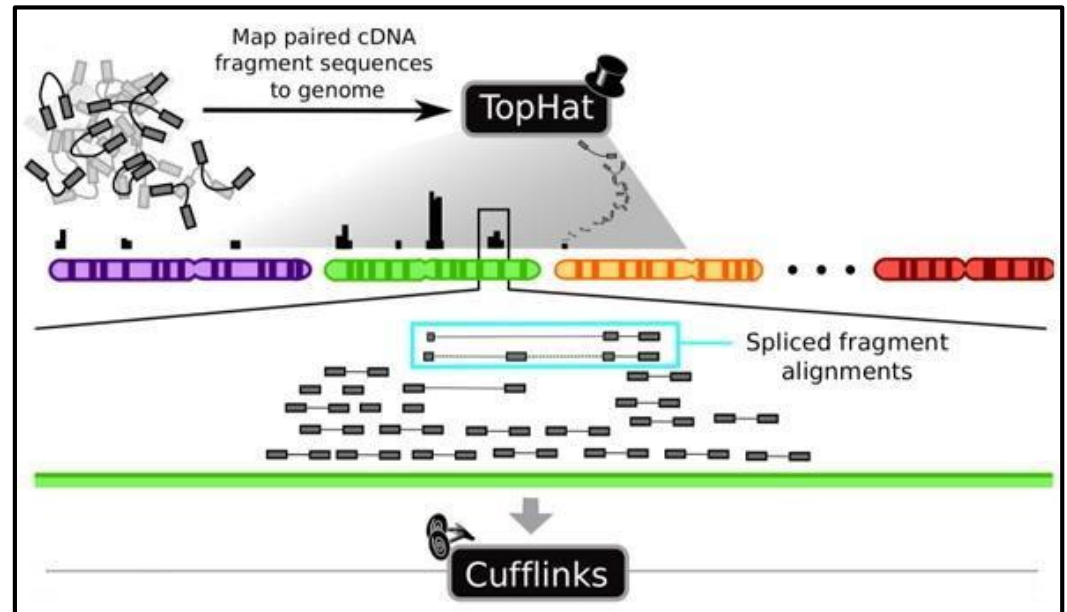
Exons :

Jonctions (dans les paires & les Reads)

Cufflinks

Reconstruction de transcrits

- ❖ Fragments divisés en *loci* non chevauchants
- ❖ Chaque *locus* est assemblé indépendamment



Trapnell et al. Nat Biotechnol. 2010

Cufflinks

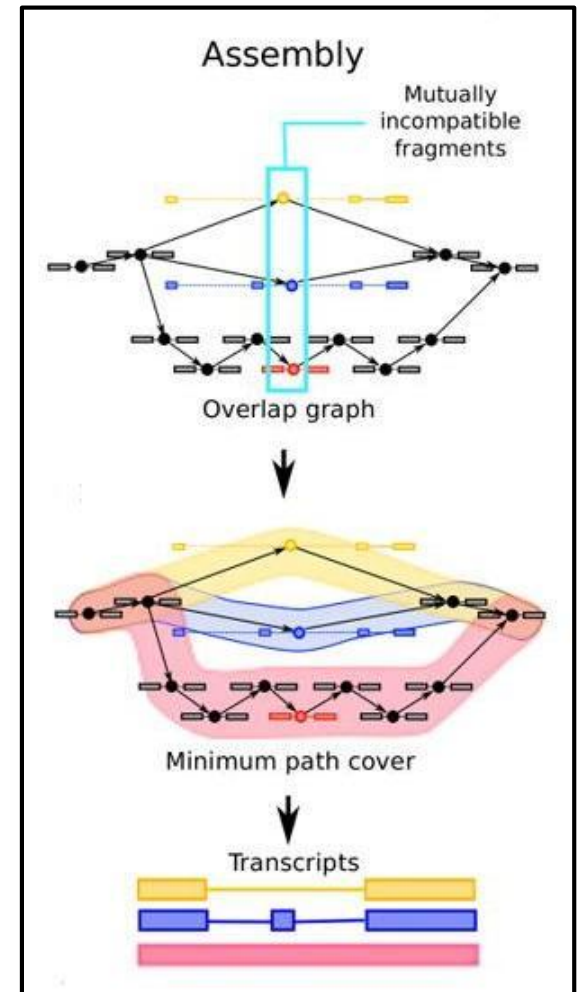
Reconstruction de transcrits

❖ Les différents chemins :

- trouver les **positions des gènes**
- trouver les **exons**
- trouver les **jonctions** :
 - entre les paires
 - dans les séquences

❖ Stratégie de construction du modèle :

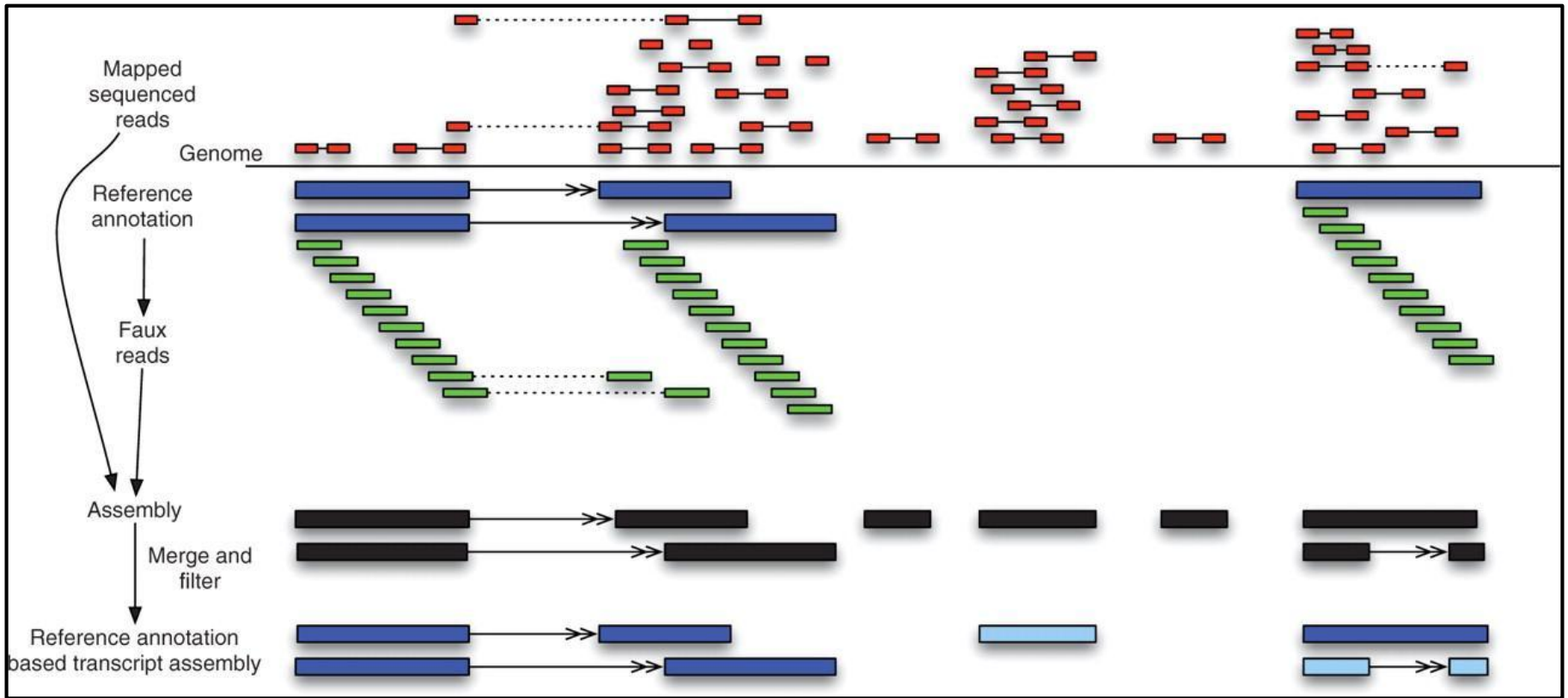
- trouver le **nombre minimum de modèles qui expliquent les lectures** :
 - **minimum de chemins**, théorème de Dilworth
 - **nb de lectures incompatibles**
= **nb minimum de transcrits nécessaires**
 - **1 chemin = 1 isoforme**



Trapnell et al. Nat Biotechnol. 2010

Cufflinks

Reference Annotation Based Transcripts Assembly



Roberts et al. Bioinformatics 2011

Cufflinks

- ❖ Reference fasta (génomome)
- ❖ Référence gtf (transcriptome)
- ❖ 1 bam par échantillon
- ❖ Quelles sont les stratégies possibles pour identifier le **maximum** de transcrits ?

Cufflinks

Reconstruction des transcrits

❖ En entrée :

- lectures (.sam/.bam)
- Use guide transcript assembly : annotations (.gtf)

❖ En sortie :

- **transcrits (.gtf)** :
 - positionnement et quantification des isoformes
- **gènes (.fpkm_tracking)** :
 - F/RPKM des gènes
- **isoformes (.fpkm_tracking)** :
 - F/RPKM des isoformes

Cufflinks - Cuffcompare

Class code de cuffcompare

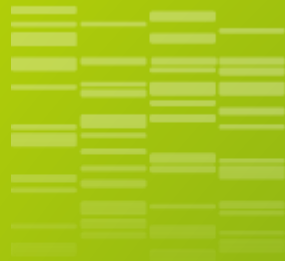
=	identité	
c	inclus	
j	nouvel isoforme	
e	exon	
i	intron	
o	chevauchant	
p	polymerase run-on	
r	répétition	
u	autre	
x	exon antisens	
s	intron antisens	

http://cufflinks.cbc.umd.edu/manual.html#class_codes

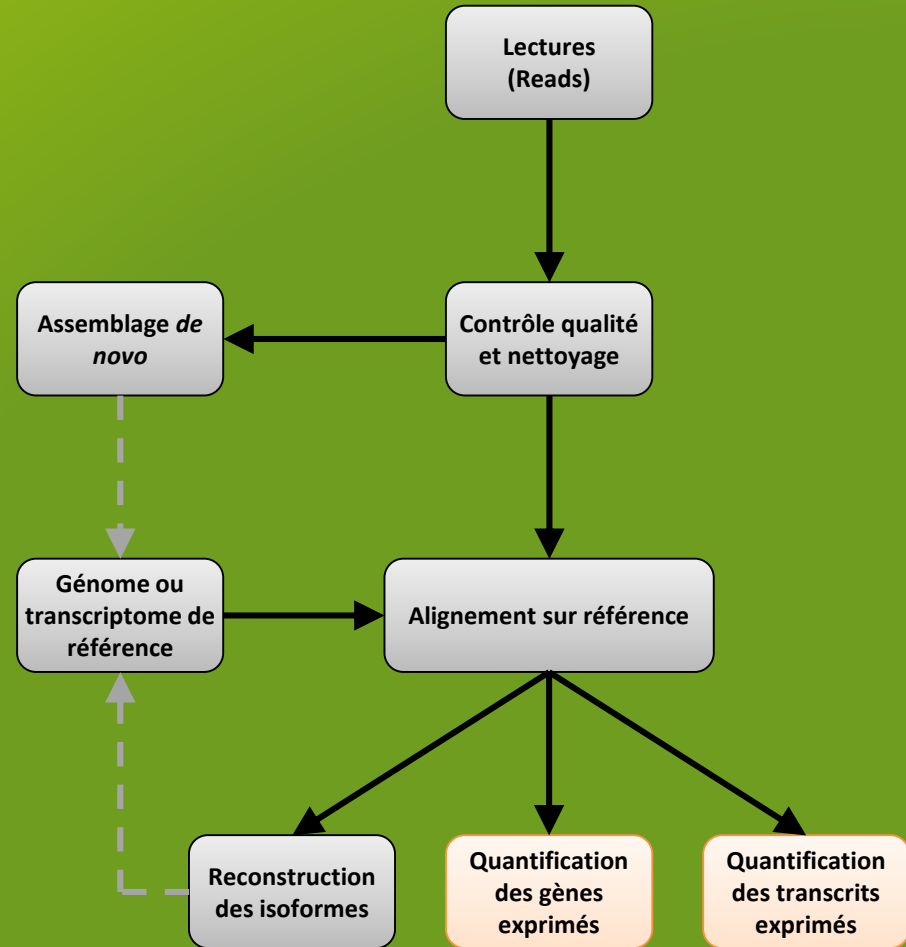
TP - Découverte de transcrit



- ❖ Fusionner les alignements (samtools merge)
- ❖ Supprimer les duplicats (samtools rmdup)
- ❖ Détecter les nouveaux transcripts (cufflinks)
- ❖ Ouvrir le nouveau transcriptome dans IGV



05 Quantification



Quantification

Que cherche-t-on à compter ?

❖ Quel *feature* compter ?

- gènes
- exons
- transcrits

chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	aggregate_gene	7529	9484	.	.	.	gene_id "FBgn0031208"
chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	exonic_part	7529	8116	.	+	.	transcripts "FBtr0300689+FBtr0300690"; exon_c_part_number
chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	exonic_part	8193	8589	.	.	+	transcripts "FBtr0300689+FBtr0300690"; exon_c_part_number
chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	exonic_part	8590	8667	.	.	+	transcripts "FBtr0300689"; exon_c_part_number "003"; gene
chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	exonic_part	8668	9484	.	.	+	transcripts "FBtr0300689+FBtr0300690"; exon_c_part_number
chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	aggregate_gene	9836	21372	.	.	.	gene_id "FBgn002121"
chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	c_part_number "001"; gene_id "FBgn002121"	exonic_part	9836	11344	.	.	transcripts "FBtr0078170+FBtr0078171+FBtr0078169+FBtr0078172"
chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	exonic_part	11410	11518	.	.	.	transcripts "FBtr0078170+FBtr0078171+FBtr0078169+FBtr0078172"
chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	c_part_number "002"; gene_id "FBgn002121"	exonic_part	11779	12221	.	.	transcripts "FBtr0078170+FBtr0078171+FBtr0078169+FBtr0078172"
chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	c_part_number "003"; gene_id "FBgn002121"	exonic_part	12286	12928	.	.	transcripts "FBtr0078170+FBtr0078171+FBtr0078169+FBtr0078172"
chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	c_part_number "004"; gene_id "FBgn002121"	exonic_part	13520	13625	.	.	transcripts "FBtr0078170+FBtr0078171+FBtr0078169+FBtr0078172"
chr2L	Drosophila_melanogaster.BDGP5.25.62.gtf.gz	exonic_part	13683	14874	.	.	.	transcripts "FBtr0078170+FBtr0078171+FBtr0078169+FBtr0078172"

❖ Comptage brut sur les gènes ou les exons :

- htseq-count

❖ Comptage brut sur les gènes ou les exons ou les transcrits:

- featureCount

❖ Estimation de l'abondance des transcrits reconstruits :

- Cufflinks

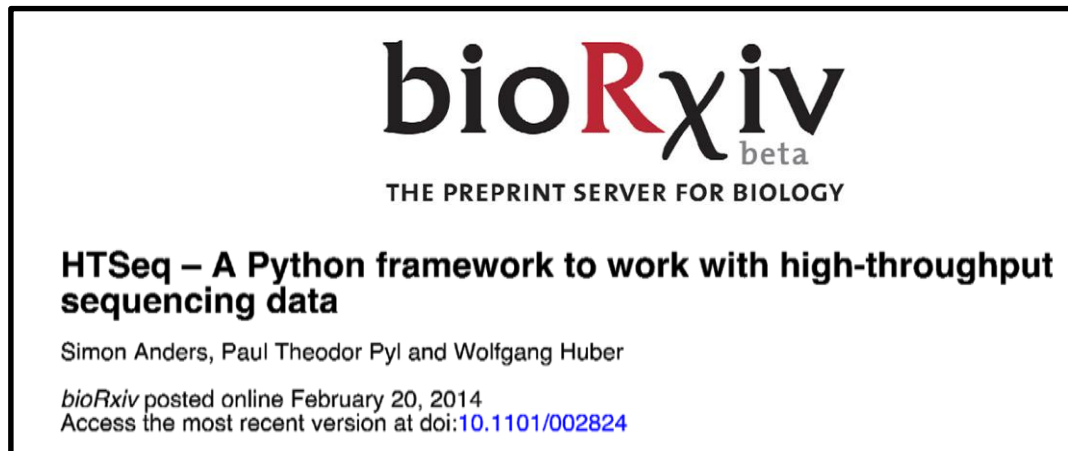
❖ Dépend des données disponibles

gene_id	untreated1	untreated2	untreated3	untreated4	treated1
FBgn0000003	0	0	0	0	1
FBgn0000008	92	161	76	70	140
FBgn0000014	5	1	0	0	4
FBgn0000015	0	2	1	2	1
FBgn0000017	4664	8714	3564	3150	6205
FBgn0000018	583	761	245	310	722
FBgn0000022	0	1	0	0	0
FBgn0000024	10	11	3	3	10
FBgn0000028	0	1	0	0	1
FBgn0000032	1446	1713	615	672	1698

Comptage des gènes et exons

HTSeq-count

- ❖ **Comptage des lectures** s'alignant sur une *feature* donnée :
 - gène
 - exon
- ❖ **Utilise** les fichiers d'**alignement** (SAM/BAM) et une **annotation**



<http://www-huber.embl.de/users/anders/HTSeq/doc/count.html#count>

HTSeq-count



	union	intersection_strict	intersection_nonempty
	gene_A	gene_A	gene_A
	gene_A	no_feature	gene_A
	gene_A	no_feature	gene_A
	gene_A	gene_A	gene_A
	gene_A	gene_A	gene_A
	ambiguous	gene_A	gene_A
	ambiguous	ambiguous	ambiguous

```
Usage: htseq-count [options] alignment_file gff_file

This script takes an alignment file in SAM/BAM format and a feature file in GFF format and calculates for each feature the number of reads mapping to it. See http://www-huber.embl.de/users/anders/HTSeq/doc/count.html for details.

Options:
-h, --help                show this help message and exit
-f SAMTYPE, --format=SAMTYPE
                           type of <alignment_file> data, either 'sam' or 'bam'
                           (default: sam)
-r ORDER, --order=ORDER
                           'pos' or 'name'. Sorting order of <alignment_file>
                           (default: name). Paired-end sequencing data must be
                           sorted either by position or by read name, and the
                           sorting order must be specified. Ignored for single-
                           end data.
-s STRANDED, --stranded=STRANDED
                           whether the data is from a strand-specific assay.
                           Specify 'yes', 'no', or 'reverse' (default: yes).
                           'reverse' means 'yes' with reversed strand
                           interpretation
-a MINAQUAL, --minaqual=MINAQUAL
                           skip all reads with alignment quality lower than the
                           given minimum value (default: 10)
-t FEATURETYPE, --type=FEATURETYPE
                           feature type (3rd column in GFF file) to be used, all
                           features of other type are ignored (default, suitable
                           for Ensembl GTF files: exon)
-i IDATTR, --idattr=IDATTR
                           GFF attribute to be used as feature ID (default,
                           suitable for Ensembl GTF files: gene_id)
-m MODE, --mode=MODE
                           mode to handle reads overlapping more than one feature
                           (choices: union, intersection-strict, intersection-
                           nonempty; default: union)
-o SAMOUT, --samout=SAMOUT
                           write out all SAM alignment records into an output SAM
                           file called SAMOUT, annotating each line with its
                           feature assignment (as an optional field with tag
                           'XF')
-q, --quiet                suppress progress report
```

HTSeq-count

fichiers de sortie

- ❖ Une **table de comptage** pour chaque *feature* ainsi qu'un **résumé**
 - **__no_feature** : lectures non assignées
 - **__ambiguous** : lectures assignables à plus d'un feature, non comptées
 - **__too_low_aQual** : lectures filtrées sur la qualité d'alignement (-a)
 - **__not_aligned** : lectures non alignées du fichier d'entrée
 - **__alignment_not_unique** : lectures avec alignement multiple (BAM)

FBgn0264694	246	
FBgn0264706	0	
FBgn0264712	251	
__no_feature	65531	
__ambiguous	11617	
__too_low_aQual	0	
__not_aligned	0	
__alignment_not_unique		1967339

featureCounts

- ❖ htseq-count ++.
- ❖ Niveau exon, gène, transcrit.
- ❖ 1 read peut être attribué à plusieurs Feature.
- ❖ Reads avec alignement multiples peuvent être pris en compte.
- ❖ Brin-spécifique très bien géré.
- ❖ 2 Notions :
 - *feature* (e.g. exon)
 - *meta-feature* : agrégation de feature (e.g. gene)

featureCounts options

Feature Counts (version 1.0.0)

Your annotation file (gtf file):

39: Cufflinks on merged: assembled transcripts

Give the name of the annotation file. The program assumes that the provided annotation file is in GTF format. Use -F option to specify other annotation formats.

First SAM/BAM file:

29: {WT_rep1_1_Ch6.fastq}-Tophat_mapped.bam

Give the names of input read files that include the read mapping results. Format of input files is automatically determined (SAM or BAM). Paired-end reads will be automatically read each other. Multiple files can be provided at the same time.

Add another BAM/SAM datasets

Add another BAM/SAM dataset 1

Other SAM/BAM files:

32: {MT_rep1_1_Ch6.fastq}-Tophat_mapped.bam

Remove Add another BAM/SAM dataset 1

Add new Add another BAM/SAM dataset

Specify feature type:

exon

Only rows which have the matched feature type in the provided GTF annotation file will be included for read counting. 'exon' by default

Specify the attribute type used to group features (eg. exons) into meta-features (eg. genes), when GTF annotation is provided:

gene_id

Reads will be allowed to be assigned to more than one matched meta-feature:

Yes

Indicate if strand-specific read counting should be performed:

unstranded

Multi-mapping reads/fragments will be counted:

Yes

Only primary alignments will be counted:

Yes

Minimum number of overlapped bases required to assign a read to a feature:

30

Negative values are permitted, indicating a gap being allowed between a read and a feature.

Optional paired-end parameters:

Paired-end reads

featureCounts : options

Multi-mapping reads/fragments will be counted:

Yes ▾

Only primary alignments will be counted:

Yes ▾

Minimum number of overlapped bases required to assign a read to a feature:

30

Negative values are permitted, indicating a gap being allowed between a read and a feature.

Optional paired-end parameters:

Paired-end reads ▾

Fragments (or templates) will be counted instead of reads. The two reads from the same fragment must be adjacent to each other in the provided SAM/BAM file:

Fragments NOT counted instead of reads ▾

Paired-end distance will be checked when assigning fragments to meta-features or features:

Paired-end distance will NOT be checked. ▾

Minimum fragment/template length:

50

Minimum fragment/template length, 50 by default.

Maximum fragment/template length:

600

Maximum fragment/template length, 600 by default.

If specified, only fragments that have both ends successfully aligned will be considered for summarization:

Not only fragments with both ends successfully aligned ▾

If specified, the chimeric fragments (those fragments that have their two ends aligned to different chromosomes) will NOT be included for summarization:

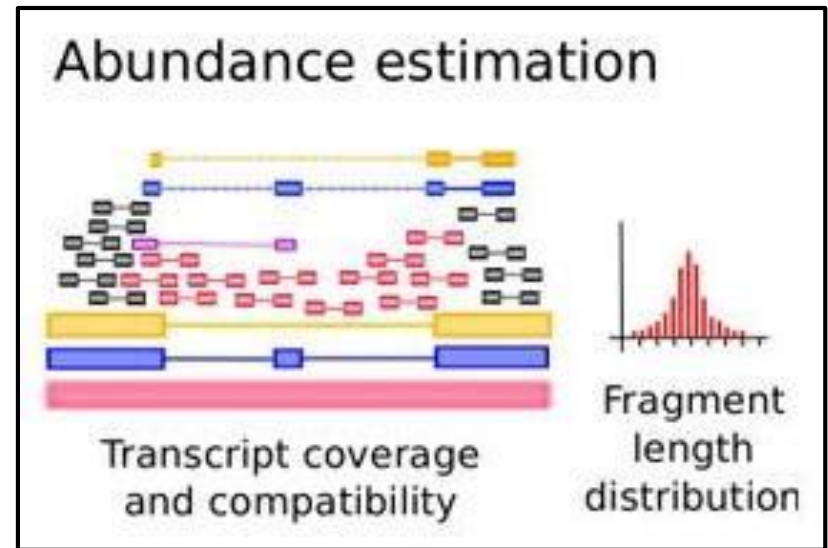
The chimeric fragments will NOT be included ▾

Execute

Cufflinks

Principes

- **Assignment des lectures à un transcrit**
- **Estimation de l'abondance de chaque transcrit** mesurée en :
 - RPKM (*single reads*)
 - FPKM (*paired-end reads*)



Trapnell et al. Nat Biotechnol. 2010

Cufflinks

RPKM / FPKM

❖ Permet de corriger les **biais de longueur** des transcrits

❖ RPKM :

○ **Reads Per Kilobase** of exon per **Million** fragments mapped :

R = Nombre de read mappés

N = Nombre total de read de la librairie

L = taille des exons du gène en bp

$$\text{RPKM} = \frac{10^9 \times R}{N \times L}$$

❖ FPKM :

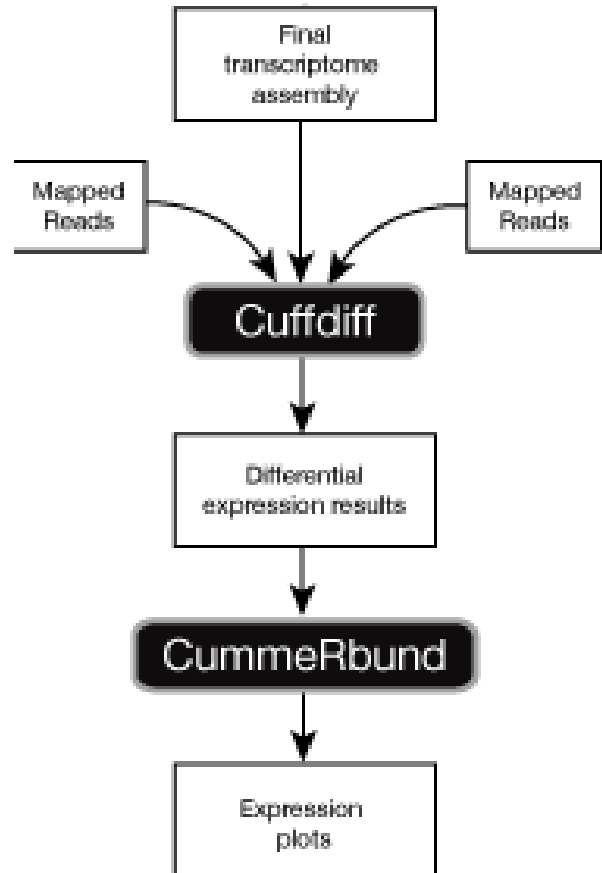
○ **Fragments Per Kilobase** of exon per **Million** fragments mapped

○ **1 paire de lecture = 1 fragment**

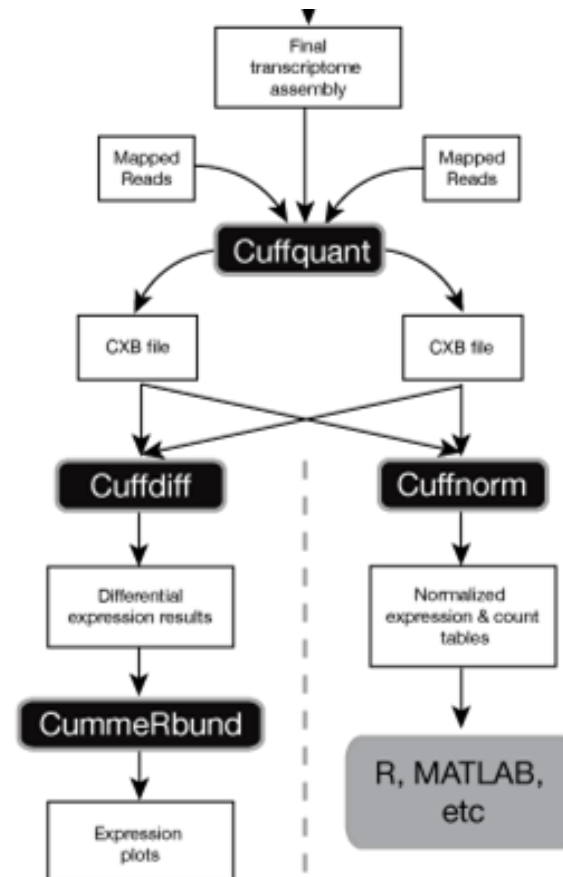
Mortazavi et al. Nature Methods 2008

Cufflinks - Estimation de l'abondance

<2.2.0



>=2.2.0



Cufflinks

❖ En **entrée** :

- transcriptome de référence
- bam d'un échantillon
- mode : assemble ou quantifie ?

❖ En **sortie** :

- **assembled transcripts (.gtf)** :
 - positionnement et quantification des isoformes
- **gene expression (.fpkm_tracking)** :
 - F/RPKM des gènes
- **transcript expression (.fpkm_tracking)** :
 - F/RPKM des isoformes

Cufflinks outputs

Description du format GTF

❖ transcripts.gtf :

- coordonnées et abondance des isoformes
- annotations (.gtf)

❖ Score :

- le plus abondant = 1000
- moins abondant : ratio = minor FPKM / major FPKM

3R	<u>Cufflinks</u>	transcript	30207	41011	901	-	.
3R	<u>Cufflinks</u>	<u>exon</u>	30207	31507	901	-	.
3R	<u>Cufflinks</u>	<u>exon</u>	40845	41011	901	-	.

```
gene_id "FBgn0037217"; transcript_id "FBtr0111206"; FPKM "3.7298141180"; frac "0.003967"; co  
gene_id "FBgn0037217"; transcript_id "FBtr0111206"; exon_number "1"; FPKM "3.7298141180"; f  
gene_id "FBgn0037217"; transcript_id "FBtr0111206"; exon_number "2"; FPKM "3.7298141180"; f
```

```
"; conf_lo "0.000000"; conf_hi "7.855077"; cov "0.330100";  
80"; frac "0.003967"; conf_lo "0.000000"; conf_hi "7.855077"; cov "0.330100";  
80"; frac "0.003967"; conf_lo "0.000000"; conf_hi "7.855077"; cov "0.330100";
```

Cufflinks outputs

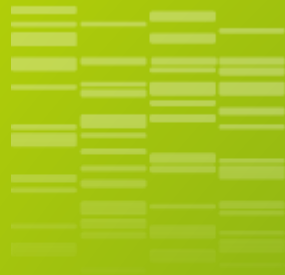
- ❖ Isoform estimation
- ❖ Gene estimation
- ❖ Format :

tracking_id	class_code	nearest_ref_id	gene_id	gene_short_name	tss_id	locus
CUFF.1.1	-	-	CUFF.1	-	-	2L:5374-5556
CUFF.2.1	-	-	CUFF.2	-	-	2L:5382-5555
CUFF.3.1	-	-	CUFF.3	-	-	2L:5707-5856

length	coverage	FPKM	FPKM_conf_lo	FPKM_conf_hi	FPKM_status
182	351.21	12.2989	8.54337	16.0544	OK
173	121.069	4.23967	1.69117	6.78817	OK
149	620.192	21.7182	13.0174	30.419	OK



TP : Quantification

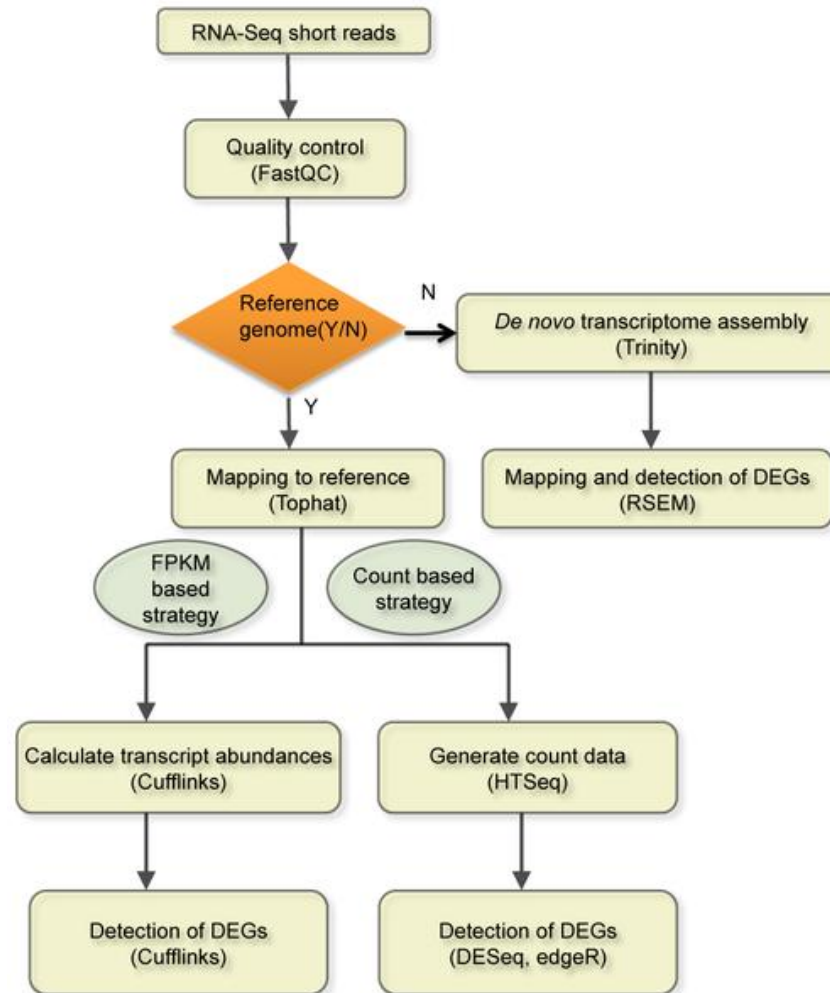


06 Conclusion

Conclusion générale

- ❖ Workflow galaxy à construire
- ❖ Choix des outils dépendent des données disponibles et de la question biologique
- ❖ Tous les outils sont dispo sur Migale et Galaxy
- ❖ Et maintenant en avant pour les stats !

Figure 1. The workflow of differential expression analysis for RNA-Seq data.



Zhang ZH, Jhaveri DJ, Marshall VM, Bauer DC, Edson J, et al. (2014) A Comparative Study of Techniques for Differential Expression Analysis on RNA-Seq Data. PLoS ONE 9(8): e103207. doi:10.1371/journal.pone.0103207
<http://127.0.0.1:8081/plosone/article?id=info:doi/10.1371/journal.pone.0103207>

Liens utiles

- ❖ **Seqanswer** : <http://seqanswers.com/>
- ❖ **Biostar** : <https://www.biostars.org/>
- ❖ **RNA-Seq blog** : <http://rna-seqblog.com/>

Remerciements

- ❖ Le groupe de travail « **Planification d'expériences et RNA-seq** »
du **PEPI IBIS**

Satisfaction form :

<http://bioinfo.genotoul.fr/index.php?id=79>