


Plateforme Bioinformatique Midi-Pyrénées



# RNA-Seq data analysis

Cédric Cabau Siganae / Céline Noirot Bioinfo Genotoul

1

---

---

---

---


---

---

---

---

Plateforme Bioinformatique Midi-Pyrénées



## Material

<http://bioinfo.genotoul.fr/index.php?id=119>

Slides & Exercise leaflet (doc)

- pdf : one per page
- pdf : three per page with comment lines

Data & results files (data)

<http://genoweb.toulouse.inra.fr/~formation/LigneCmd/RNAseq/data/>

2

---

---

---

---

---

---

---

---

Plateforme Bioinformatique Midi-Pyrénées



## Session organisation

<p><b>Afternoon (14h-17h) :</b></p> <ul style="list-style-type: none"> <li>- Sequence quality           <ul style="list-style-type: none"> <li>• Theory + exercises</li> </ul> </li> <li>- Spliced read mapping           <ul style="list-style-type: none"> <li>• Theory + exercises</li> </ul> </li> </ul>	<p><b>Morning (9h00 -12h30) :</b></p> <ul style="list-style-type: none"> <li>- Visualisation           <ul style="list-style-type: none"> <li>• Exercises</li> </ul> </li> <li>- expression measurement           <ul style="list-style-type: none"> <li>• Theory + exercises</li> </ul> </li> </ul> <p><b>Afternoon (14h-17h) :</b></p> <ul style="list-style-type: none"> <li>- mRNA calling           <ul style="list-style-type: none"> <li>• Theory + exercises</li> </ul> </li> <li>- some statistics ...</li> </ul>
--	--

3

---

---

---

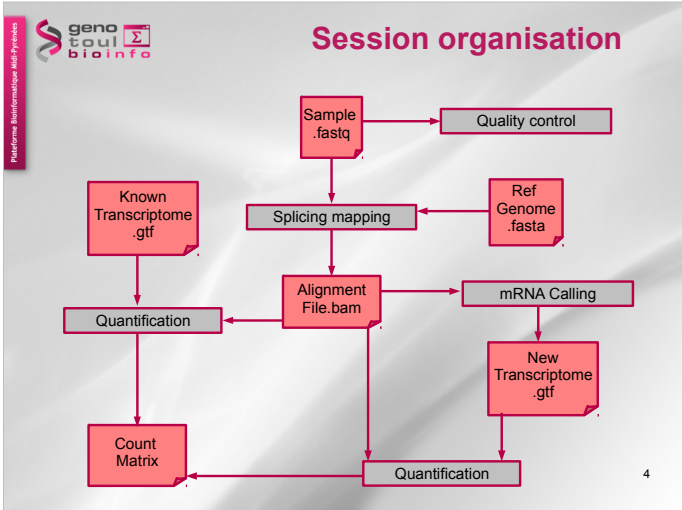
---

---

---

---

---




---

---

---

---

---

---

---

---

---

---

---

---

**What you should know**

How to connect to [genotoul.toulouse.inra.fr](http://genotoul.toulouse.inra.fr)  
 How to use unix commands?

wget URL  
 qllogin -pe parallel\_smp 4

5

---

---

---

---

---

---

---

---

---

---

---

---

- Summary - Sequence quality**
- Context, vocabulary, transcriptome variability ...
  - Methods to analyse transcriptoms
  - What is RNAseq ?
  - High throughput sequencers
  - Illumina protocol, paired-end library, directional library
  - Known biases
  - How to check quality ?
- 6

---

---

---

---

---

---

---

---

---

---

---

---



Platforme Bioinformatique M&P Pythons

geno toul bioinfo

## Vocabulary

**Gene** : functional units of DNA that contain the instructions for generating a functional product.

**Exon** : coding region of mRNA included in the transcript  
**Intron** : non coding region  
**TSS** : Transcription Start Site  $\neq$  1<sup>st</sup> amino acid  
**Transcript** : stretch of DNA transcribed into an RNA molecule

10

---

---

---

---

---

---

---

---

---

---

---

---

Platforme Bioinformatique M&P Pythons

geno toul bioinfo

## Alternative splicing

**Alternative splicing** (or differential splicing)

- the exons are reconnected in multiple ways during RNA splicing.
- different mRNAs translated into different protein isoforms
- a single gene may code for multiple proteins.

**Intron Retention**

**Post-transcriptional modification** (eukaryotic cells) eg: the conversion of precursor messenger RNA into mature mRNA (mRNA), editing...

[http://en.wikipedia.org/wiki/Alternative\\_splicing](http://en.wikipedia.org/wiki/Alternative_splicing)

11

---

---

---

---

---

---

---

---

---

---

---

---

Platforme Bioinformatique M&P Pythons

geno toul bioinfo

## Transcript degradation

- mRNA export to the cytoplasm,
- protected from degradation by a 5' cap structure and a 3' polyA tail.
- the polyA tail is gradually shortened by exonucleases
- the degradation machinery rapidly degrades the mRNA in both in directions.
- others mechanisms, bypass the need for deadenylation and can remove the mRNA from the transcriptional pool independently.

<http://www.eb.tuebingen.mpg.de/research-groups/remco-sprangers>

12

---

---

---

---

---

---

---

---

---

---

---

---

genotoul bioinfo

## Cis-natural antisense transcript

- Natural antisense transcripts (NATs) are a group of RNAs encoded within a cell that have transcript complementarity to other RNA transcripts.

Fig. 1: The five orientations for overlap of cis-NAT pairs. Genes are always transcribed 5' to 3'. Regions of overlapping strands are shown with dotted lines.

[http://en.wikipedia.org/wiki/Cis-natural\\_antisense\\_transcript](http://en.wikipedia.org/wiki/Cis-natural_antisense_transcript)

13



genotoul bioinfo

## Fusion genes

- A fusion gene is a hybrid gene formed from two previously separate genes. It can occur as the result of a translocation, interstitial deletion, or chromosomal inversion. Often, fusion genes are oncogenes.

Genome Biol. 2011 Jan 19;12(1):R6. [Epub ahead of print]

**Identification of fusion genes in breast cancer by paired-end RNA-sequencing.**  
 Edoren H, Mucumagai A, Kanagaspassia S, Nicotri D, Honjojo V, Klehi K, Rye H, Nuber S, Wolf M, Borresen-Dale AL, Kallioniemi O. Institute for Molecular Medicine Finland (FIMM), Tukholmankatu 8, Helsinki, 00290, Finland. olli.kallioniemi@fimm.fi

[http://en.wikipedia.org/wiki/Fusion\\_gene](http://en.wikipedia.org/wiki/Fusion_gene)- They often come from trans-splicing : Trans-splicing is a special form of RNA processing in eukaryotes where exons from two different primary RNA transcripts are joined end to end and ligated.

<http://en.wikipedia.org/wiki/Trans-splicing>

14



genotoul bioinfo

## Transcriptome variability

- Many types of transcripts (mRNA, ncRNA ...)
- Many isoform (non canonical splice sites, intron retention ...)
- Number of transcripts
  - possible variation factor between transcripts:  $10^6$  or more,
  - expression variation between samples.
- Allele specific expression

15



## How can we study the transcriptome?

### Techniques classification

EST	PCR/RT-QPCR	SAGE	MicroArrays
No quantification	Quantification	Quantification	Indirect quantification
Low throughput	Low throughput (up to hundreds)	Low throughput (up to thousands)	High throughput (up to millions)
Discovery (Yes)	No	No	Discovery (Yes)

- Need transcript sequence partially known
- Difficulties in discovering novels splice events

---

---

---

---

---

---

---

---

---

---

---

---

## What is RNA-Seq ?

- use of **high-throughput sequencing technologies** to sequence cDNA in order to get information about a sample's RNA content
- the deep coverage and base level resolution => measure transcriptome data experimentally

Nature Reviews Genetics 10, 57-63 (January 2009) | doi:10.1038/nrg2484

ARTICLE SERIES: Applications of next-generation sequencing

INNOVATION

**RNA-Seq: a revolutionary tool for transcriptomics**

Zhong Wang<sup>1</sup>, Mark Gerstein<sup>2</sup> & Michael Snyder<sup>1</sup> [About the authors](#)

top

RNA-Seq is a recently developed approach to transcriptome profiling that uses deep-sequencing technologies. Studies using this method have already altered our view of the extent and complexity of eukaryotic transcriptomes. RNA-Seq also provides a far more precise measurement of levels of transcripts and their isoforms than other methods. This article describes the RNA-Seq approach, the challenges associated with its application, and the advances made so far in characterizing several eukaryote transcriptomes.

<http://en.wikipedia.org/wiki/RNA-Seq>

---

---

---

---

---

---

---

---

---

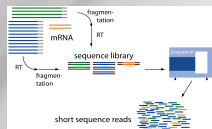
---

---

---

## What is different with RNA-Seq ?

- No prior knowledge of sequence needed
- Specificity of what is measured
- Increased dynamic range of measure, more sensitive detection
- Direct quantification
- Good reproducibility
- Different levels : genes, transcripts, allele specificity, structure variations
- New feature discovery: transcripts, isoforms, ncRNA, structures (fusion...)
- Possible detection of SNPs, ...




---

---

---

---

---

---

---

---

---

---

---

---



## Usual questions on RNA-Seq !

- How many replicates ?
  - Technical or/and biological replicates ?
- How many reads for each sample?
- How many conditions for a full transcriptome ?
- How long should my reads be ?
- Single-end or paired-end ?

---

---

---

---

---

---

---

---

---

---

---

---

## Depth VS Replicates

- Encode (2011) : <https://www.encodeproject.org/data-standards/>
  - Experiments should be performed with **two or more biological replicates**, unless there is a compelling reason why this is impractical or wasteful
  - A typical  $R^2$  (Pearson) correlation of gene expression (RPKM) between two biological replicates, for RNAs that are detected in both samples using RPKM or read counts, should be [0.92 - 0.98] Experiments with biological correlations < 0.9 should either be repeated or explained.
- Between **30M and 100M reads** per sample depending on the study.
- On Human 100M reads are enough to detect 90% of transcript from 81% of genes.
- Zhang et al. 2014 : From 3 replicates improve DE detection and control false positive rate.

---

---

---

---

---

---

---

---

---

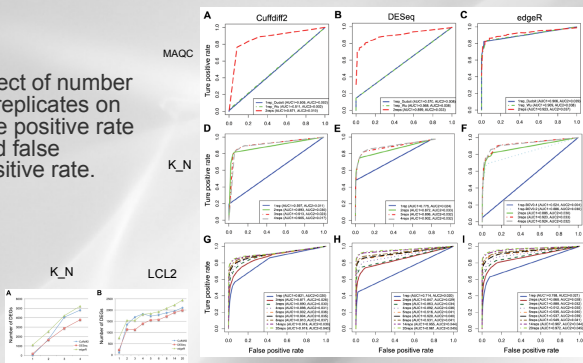
---

---

---

## Depth VS Replicates

Effect of number of replicates on true positive rate and false positive rate.




---

---

---

---

---

---

---

---

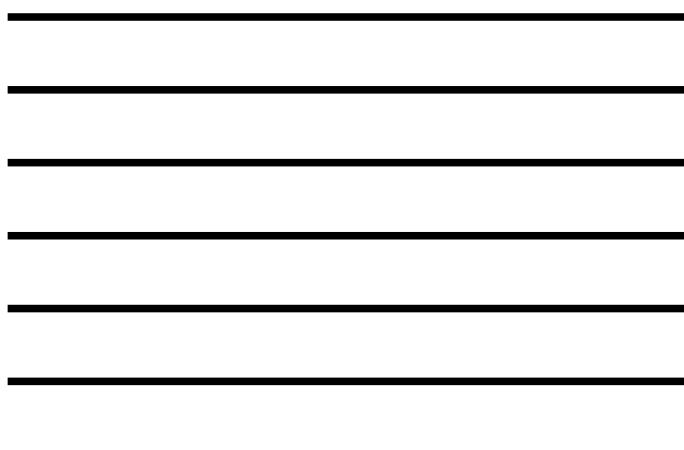
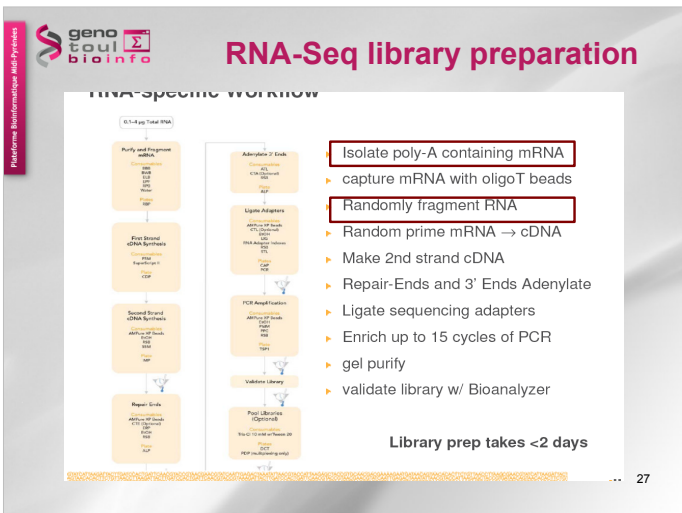
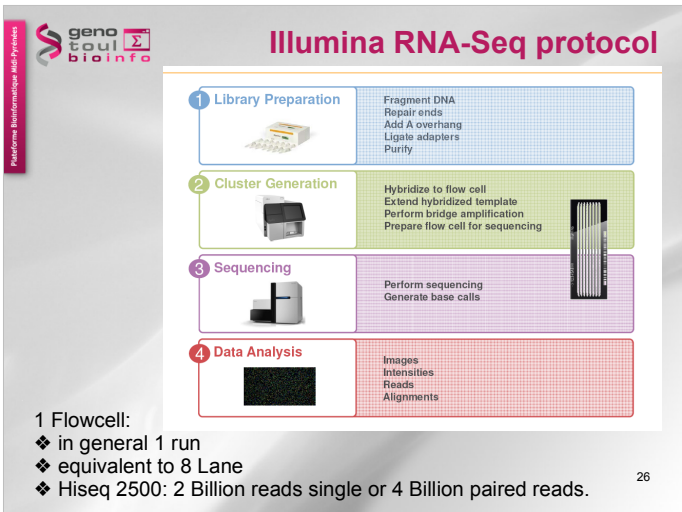
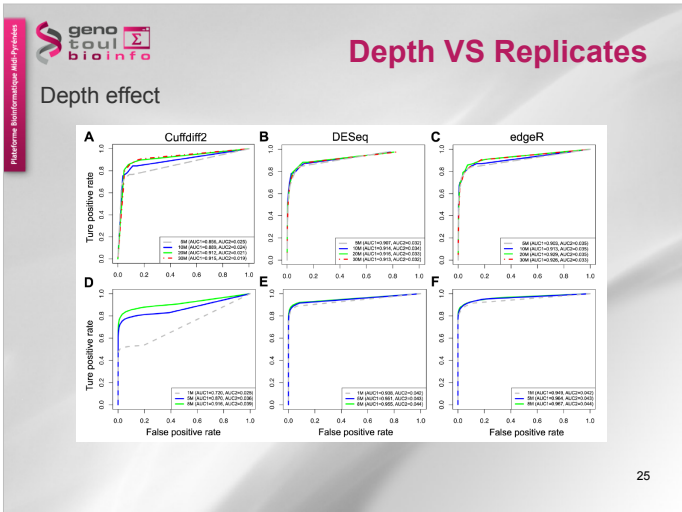
---

---

---

---





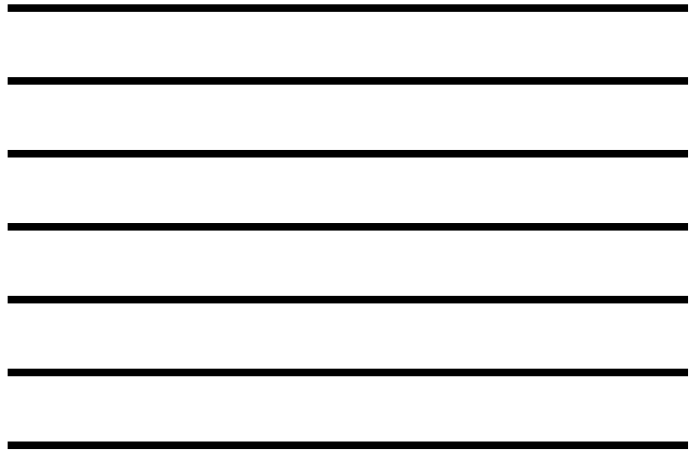
Platforme Bioinformatique MSB-PyBiochem

geno **to**ul **Σ** **bio**info

## Clusters generation / Sequencing

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

28



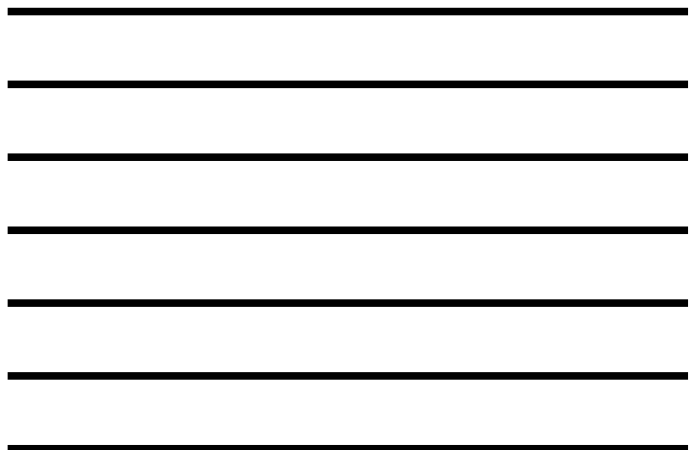
Platforme Bioinformatique MSB-PyBiochem

geno **to**ul **Σ** **bio**info

## Paired-end sequencing

- Modification of the standard single-read DNA library preparation facilitates reading both ends of each fragment
- Improvement of mapping
- Help to detect structural variations in the genome like insertions or deletions, copy number variations, and genome rearrangements

29



Platforme Bioinformatique MSB-PyBiochem

geno **to**ul **Σ** **bio**info

## Strand specific RNA-Seq protocol

workflow comparison: mRNA-Seq vs directional mRNA-Seq

1. start with total RNA (or total mRNA)
2. purify poly-A mRNA
3. randomly fragment mRNA
4. 1st strand cDNA synthesis
5. 2nd strand cDNA synthesis
6. end repair
7. adenylate 3' ends
8. ligate adaptors
9. gel purify
4. end repair with phosphatase and PINK
5. column purify PINK treated mRNA
6. ligate 5' adaptor
7. ligate 3' adaptor
8. reverse transcriptase
10. enrich with PCR
11. validate library
12. grow clusters
13. sequence on HiSeq2000 (SR or PE)

Next Methods, 2009 Sep 7;9:109-115. Epub 2009 Aug 15.

**Comprehensive comparative analysis of strand-specific RNA sequencing methods.**

Leung J, Szechar M, Adkins S, Nishizumi C, Thompson DA, Friedman N, Golub A, Brenner S  
 Broad Institute of Massachusetts Institute of Technology and Harvard University, Cambridge, Massachusetts, USA.  
 jleung@broadinstitute.org

Abstract

30

<http://www.plosbiology.org/article/info:doi/10.1371/journal.pbio.1001046>





Platforme Bioinformatique M&P Pythons

genotoul bioinfo

## Hexamer random priming bias

Published online 14 April 2010  
Nucleic Acids Res. 2010, Vol. 38, No. 12, e139  
doi:10.1093/nar/gkq254

**Biases in Illumina transcriptome sequencing caused by random hexamer priming**  
Kasper D. Hansen<sup>1</sup>\*, Steven E. Brenner<sup>2</sup> and Sandrine Dudot<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.

- A strong distinctive pattern in the nucleotide frequencies of the first 13 positions at the 5'-end :
  - sequence specificity of the polymerase
  - due to the end repair performed
- Reads beginning with a hexamer over-represented in the hexamer distribution at the beginning relative to the end are down-weighted

34



Platforme Bioinformatique M&P Pythons

genotoul bioinfo

## Hexamer random effect

- Orange = reads start sites
- Blue = coverage

35



Platforme Bioinformatique M&P Pythons

genotoul bioinfo

## Transcript length bias

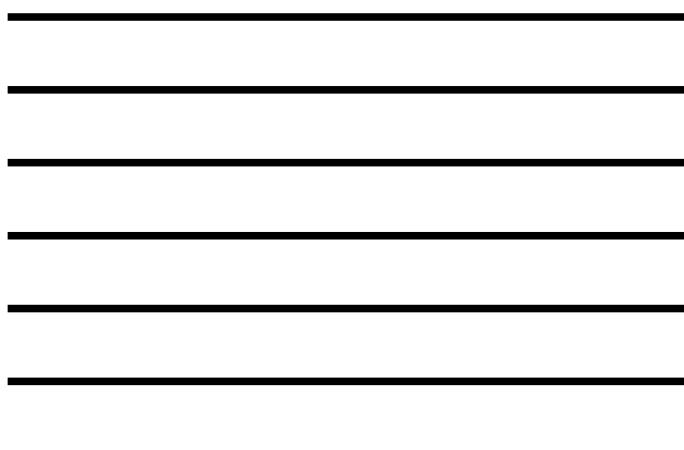
Biol Direct. 2009 Apr 18;4:14.  
**Transcript length bias in RNA-seq data confounds systems biology.**  
Oshlack A, Wakefield MJ

**Abstract**  
**Background:** Several recent studies have demonstrated transcriptome analysis (RNA-seq) in mammals. genome transcriptional profiling is likely to become genomic sequences. As yet, a rigorous analysis is still in the stages of exploring the features of the **Results:** We investigated the effect of transcript published data sets. For standard analyses using to call differentially expressed genes between using transcript. **Conclusions:** Transcript length bias for calling differentially expressed genes using RNA-seq technology. This expressed genes, and in particular may introduce other multi-gene systems biology analyses. **Reviews:** This article was reviewed by Robert C. O'Connor (nominated by Mark Ragan) and James B.

*- the differential expression of longer transcripts is more likely to be identified than that of shorter transcripts*

BIOINFORMATICS ORIGINAL PAPER  
Gene expression  
Length bias correction for RNA-seq data in gene set analyses  
Liyen Gao<sup>1,2</sup>, Zhide Fang<sup>2,3</sup>, Kai Zhang<sup>1</sup>, Degui Zhi<sup>1</sup> and Xiangqin Gu<sup>1,4</sup>

36



genotoul bioinfo

## Verifying RNA-Seq raw data

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

- Has been developed for genomic data

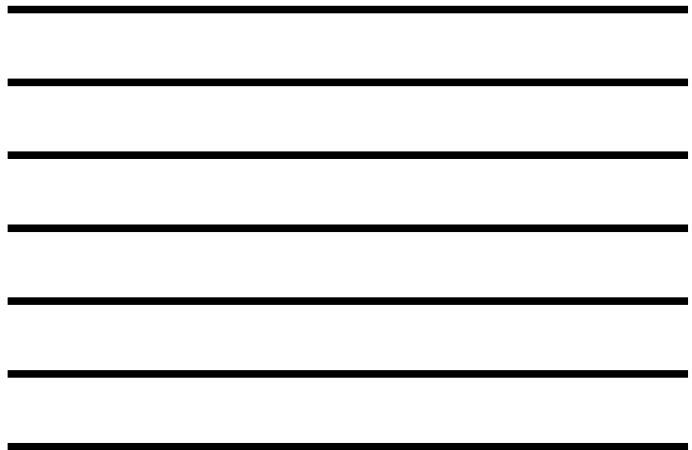
37



genotoul bioinfo

## FastQC graphics

38



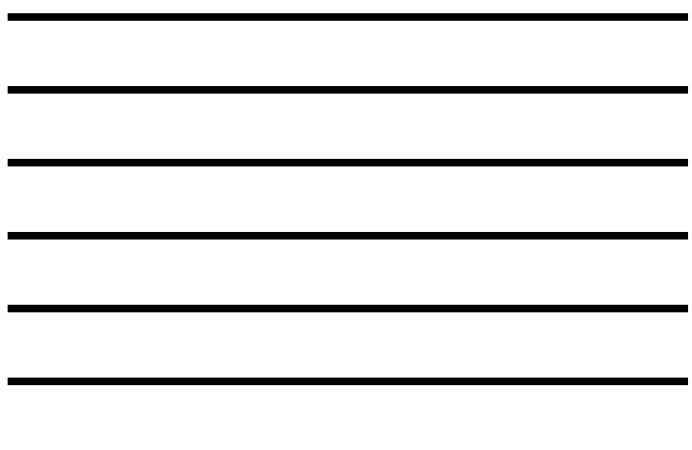
genotoul bioinfo

## FastQC graphics : Kmer content

- A kmer is a subsequence of length k
- Should spot overrepresented sequences, give a good impression of any contamination.
- Kmers showing a rise towards the end of the library indicate progressive contamination with adaptors.
- Check for adaptor sequence or poly-A sequence

Sequence	Count	Obs/Exp Overall	Obs/Exp Max	Max Obs/Exp Position
TTTTT	2749960	8.8623	7.3762037	3
AAAAA	18310185	0.2207845	0.3009034	74
CAAGG	12488915	2.3789662	49.53375	16
AAAAA	10728075	2.3667303	56.22267	3

39



## Take home message on quality analysis

Elements to be checked :

- Random priming effect
- K-mer (polyA, polyT)

Alignment on reference for the second quality check and filtering.

A good run?:

- Expected number of reads produced (2x 2 billions / flowcell),
- Length of the reads expected (150pb),
- Random selection of the nucleotides and the GC%,
- Good alignment: very few unmapped reads, pairs mapped on opposite strands.

---

---

---

---

---

---

---

---

## Hands-on : data quality

Connection genotoul : `ssh -X nom@genotoul`

To connect to the processing node : `qlogin`

Training accounts : *anemone aster*  
*bleuet iris*  
*muguet narcisse*  
*pensee rose*  
*tulipe violette*

FastQC location : `/usr/local/bioinfo/src/FastQC/current/fastqc`

---

---

---

---

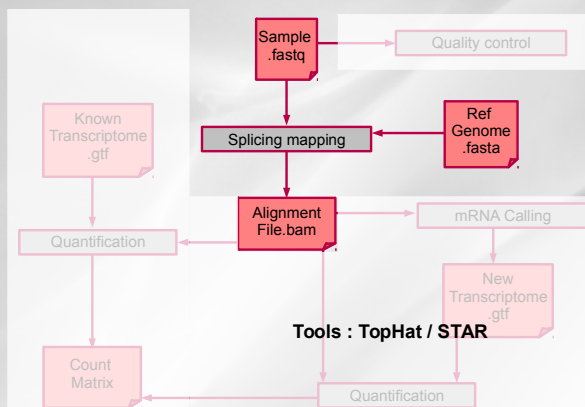
---

---

---

---

## Analysis workflow



---

---

---

---

---

---

---

---

### Aim -

#### Spliced read mapping & Visualisation

- Discover the true location (origin) of each read with respect to the reference
- Obviously features of the reference (repetitive regions, assembly errors, missing information) will render this objective impossible for a subset of the reads
- Because sequencing library was constructed from transcribed RNA, account for reads that may be split by potentially thousands of bases of intronic sequence
- Take advantage of intron/exon boundary annotations and be able to split reads across exons from no additional information (de novo spliced alignment)
- Do it in/with reasonable time/resources

---

---

---

---

---

---

---

---

### Summary -

#### Spliced read mapping & Visualisation

- Reference genome & Reference transcriptome files formats
- What is a spliced aligner ?
- Tophat principle and usage
- BAM & Bed files formats
- STAR usage
- Visualisation with IGV

---

---

---

---

---

---

---

---

### Where to find a reference genome?

Retrieving the genome file (fasta):

- The Genome Reference Consortium

<http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/>

- ! NCBI chromosome naming with « | » not well supported by mapping software

- Prefer EMBL

<http://www.ensembl.org/info/data/ftp/index.html>

---

---

---

---

---

---

---

---

**Reference transcriptome file**


What is a GTF file ?

- derived from GFF (General Feature Format, for description of genes and other features)
- Gene Transfer Format: <http://genome.ucsc.edu/FAQ/FAQformat.html#format4>

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

The [attribute] list must begin with:

- gene\_id value : unique identifier for the genomic source of the sequence.
- transcript\_id value : unique identifier for the predicted transcript.

 The chromosome name should be the same in the gtf file and fasta file

46

---

---

---

---

---

---

---

---

---


---

---

---

**Splice sites**

- Canonical splice site: which accounts for more than 99% of splicing  
GT and AG for donor and acceptor sites



[http://en.wikipedia.org/wiki/RNA\\_splicing](http://en.wikipedia.org/wiki/RNA_splicing)

- Non-canonical site: GC-AG splice site pairs, AT-AC pairs

Nucleic Acids Res. 2000 Nov 12;28(21):4364-75.  
Analysis of canonical and non-canonical splice sites in mammalian genomes.  
Bursat M, Soteldtsou IA, Solovnev VV.

- Trans-splicing: splicing that joins two exons that are not within the same RNA transcript

47

---

---

---

---

---

---

---

---

---

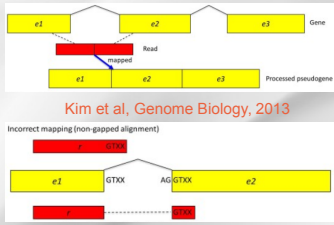
---

---

---

**Hard case**

- Lot of variations (sequencing errors, mutations)
- Repeats
- Reads spanning 3+ exons
- Gene or pseudogene



Kim et al, Genome Biology, 2013

- Small end "anchor"
- Unknown junction inside poorly rarely expressed gene

48

---

---

---

---

---

---

---

---

---

---

---

---



genotoul bioinfo

## Alignment Tools

Tools for splice-mapping:

- Tophat:
 

**BIOINFORMATICS ORIGINAL PAPER** doi:10.1093/bioinformatics/btt111

Sequence analysis

**TopHat: discovering splice junctions with RNA-Seq**

Cole Trapnell<sup>1</sup>\*, Lior Pachter<sup>2</sup> and Steven L. Salzberg<sup>3</sup>

*Genome Biol.* 2013 Apr 25;14(4):R36. doi: 10.1186/gb-2013-14-4-R36.

**TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions.**

Kim D. Pertea<sup>1,2</sup>, Trapnell C., Pimentel H., Kelley B., Salzberg S.L.
- STAR:
 

**STAR: ultrafast universal RNA-seq aligner**

Alexander Dobin<sup>1</sup>, Carrie A. Davis<sup>1</sup>, Felix Schlesinger<sup>1</sup>, Jorg Drenkow<sup>1</sup>, Chris Zaleski<sup>1</sup>, Sonali Jha<sup>1</sup>, Philippe Batut<sup>1</sup>, Mark Chaisson<sup>1</sup> and Thomas R. Gingeras<sup>1</sup>

<sup>1</sup>Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA

<sup>2</sup>Paetec Biosciences, Menlo Park, California, USA

Associate Editor: Dr. Ivanc Bircik

49

---

---

---

---

---

---

---

---

---

---

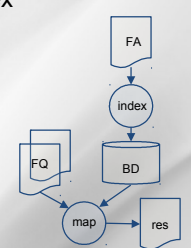
---

---

genotoul bioinfo

## Mapping steps

- Indexing reference (once only)
- Mapping reads using index



```

graph TD
  FA[FA] --> index((index))
  index --> BD[BD]
  FQ[FQ] --> map((map))
  BD --> map
  map --> res[res]
  
```

50

---

---

---

---

---

---

---

---

---

---

---

---

genotoul bioinfo

## TopHat

<http://ccb.jhu.edu/software/tophat>

- Aligns RNA-Seq reads to a reference genome with Bowtie2
- splice junction mapper for reads without knowledges
- identify splice junctions between exons

51

---

---

---

---

---

---

---

---

---

---

---

---



## Special note on the website

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.

---

---

---

---

---

---

---

---

## More topHat options

Your own junctions :

- G/--GTF <GTF2.2file>
- j/--raw-juncs <.juncs file>
- no-novel-juncs (ignored without -G/-j)

Your own insertions/deletions:

- insertions/--deletions <.juncs file>
- no-novel-indels

---

---

---

---

---

---

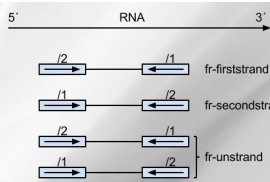
---

---

## Library types

--library-type  
TopHat will treat the reads as strand specific. Every read alignment will have an XS attribute tag. Consider supplying library type options below to select the correct RNA-seq protocol.

Library Type	Examples	Description
fr-unstranded	Standard Illumina	Reads from the left-most end of the fragment (in transcript coordinates) map to the transcript strand, and the right-most end maps to the opposite strand.
fr-firststrand	dUTP, NSR, NNSR	Same as above except we enforce the rule that the right-most end of the fragment (in transcript coordinates) is the first sequenced (or only sequenced for single-end reads). Equivalently, it is assumed that only the strand generated during first strand synthesis is sequenced.
fr-secondstrand	Ligation, Standard SOLID	Same as above except we enforce the rule that the left-most end of the fragment (in transcript coordinates) is the first sequenced (or only sequenced for single-end reads). Equivalently, it is assumed that only the strand generated during second strand synthesis is sequenced.




---

---

---

---

---

---

---

---



genotoul bioinfo

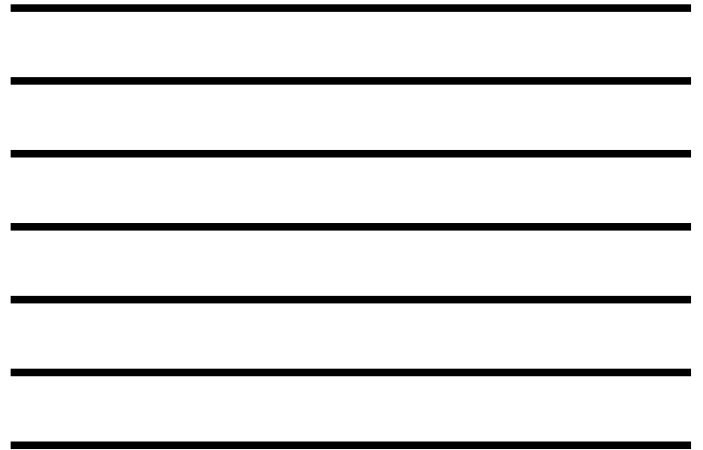
## Bed example

Chrom Start End name score strand drawing RGB Blocks info

```

track name=junctions ERR022486_etudechr22.bed description="TopHat junctions"
22 241 1451 JUNC000000001 8 - 241 1451 255,0,0,2 67,66 0,1144
22 1705 4260 JUNC000000002 1 - 1705 4260 255,0,0,2 28,48 0,2427
22 4285 4485 JUNC000000003 8 - 4285 4485 255,0,0,2 55,72 0,128
22 4575 4748 JUNC000000004 3 - 4575 4748 255,0,0,2 32,66 0,107
22 5834 6045 JUNC000000005 1 + 5834 6045 255,0,0,2 35,41 0,1170
22 6143 6776 JUNC000000006 6 - 6143 6776 255,0,0,2 61,68 0,565
22 6796 7873 JUNC000000007 5 - 6796 7873 255,0,0,2 71,51 0,226
22 7843 7254 JUNC000000008 6 + 7843 7254 255,0,0,2 66,01 0,158
22 7220 8877 JUNC000000009 11 - 7220 8877 255,0,0,2 64,62 0,1595
22 7410 16244 JUNC000000010 2 - 7410 16244 255,0,0,2 48,28 0,8886
22 7638 7811 JUNC000000011 3 + 7638 7811 255,0,0,2 58,37 0,136
22 12388 21452 JUNC000000012 27 - 12388 21452 255,0,0,2 78,72 0,8990
22 16655 27319 JUNC000000013 6 - 16655 27319 255,0,0,2 26,67 0,10597
22 27711 30684 JUNC000000014 108 - 27711 30684 255,0,0,2 74,72 0,2901
22 27714 32151 JUNC000000015 303 - 27714 32151 255,0,0,2 71,72 0,4365
22 30639 32151 JUNC000000016 134 - 30639 32151 255,0,0,2 68,72 0,1440
22 32085 32388 JUNC000000017 493 - 32085 32388 255,0,0,2 71,71 0,152
22 32234 33112 JUNC000000018 478 + 32234 33112 255,0,0,2 69,72 0,886
22 33089 33347 JUNC000000019 292 - 33089 33347 255,0,0,2 68,71 0,187
  
```

61



genotoul bioinfo

## TopHat technical issues

- Temporary disk space
  - 100 000 000 pair-ends = 0,5 To of temporary disk space
- Number of cpus
  - 100 000 000 pair-ends = 5-7 cpu days on the local cluster
- New platform cluster:
  - 34 cluster nodes with 4\*12 cores and 384 GB of ram per node: 1632 cores
  - 1 hypermem node (32 cores and 1024 GB of ram)
  - A scratch file system (157 To available, 6 Gbps bandwidth)

62



genotoul bioinfo

## An other aligner : STAR

OXFORD JOURNALS Bioinformatics (Oxford, England)

Bioinformatics, 2013 Jan; 29(1): 15-21. PMID: PMC3530905  
 Published online 2012 Oct 25. doi: 10.1093/bioinformatics/bts605

**STAR: ultrafast universal RNA-seq aligner**

Alexander Dobin,<sup>1,\*</sup> Carrie A. Davis,<sup>1</sup> Felix Schlesinger,<sup>1</sup> Jora Drenkow,<sup>1</sup> Chris Zaleski,<sup>1</sup> Sonali Jha,<sup>1</sup> Philippe Batut,<sup>1</sup> Mark Chaisson,<sup>2</sup> and Thomas R. Gingeras<sup>1</sup>

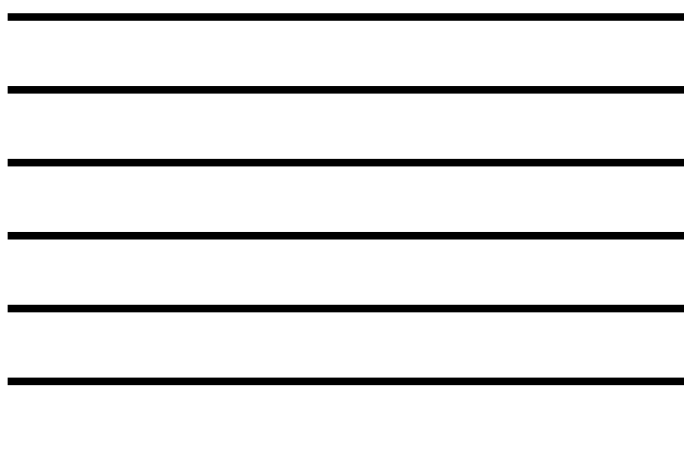
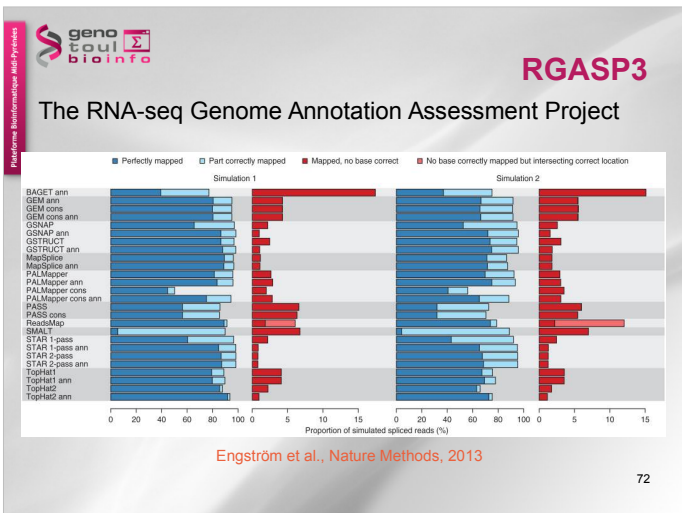
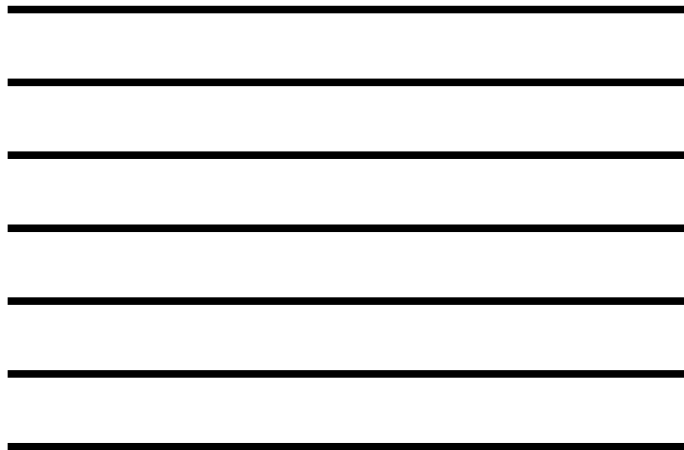
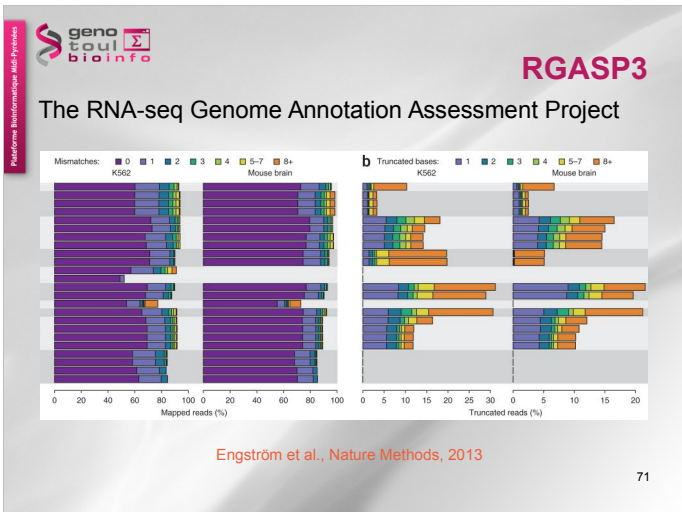
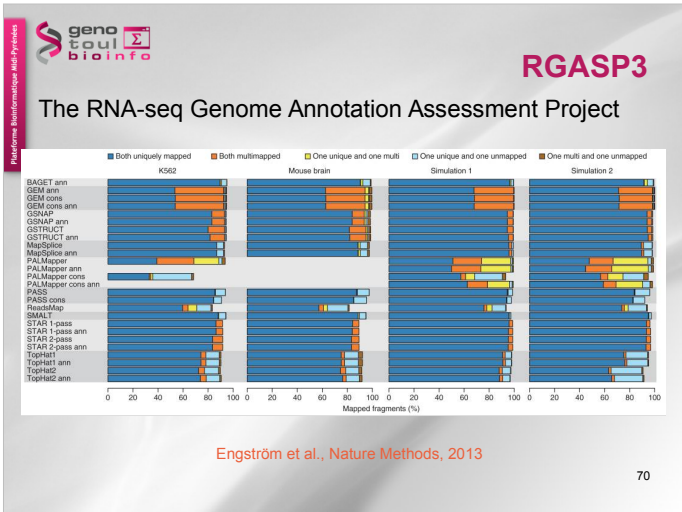
- Spliced Transcripts Alignment to a Reference
- Outperforms other aligners by more than a factor of 50 in mapping speed

63











## TopHat vs STAR

The RNA-seq Genome Annotation Assessment Project

STAR	vs	TopHat2
+	# lectures alignées	-
-	# lectures correctement alignées	+
-	Sensibilité aux variations	+
-	Sensibilité aux annotations	+

---

---

---

---

---

---

---

---

---

---

---

---

## Visualizing alignments on IGV



<http://www.broadinstitute.org/igv/home>

NATURE BIOTECHNOLOGY | OPINION AND COMMENT | CORRESPONDENCE

### Integrative genomics viewer

James T Robinson, Helga Thorvaldsdóttir, Wendy Winckler, Mitchell Guttman, Eric S Lander, Gad Getz & Jill P Mesirov

Affiliations | Corresponding authors

Nature Biotechnology 29, 24–26 (2011) | doi:10.1038/nbt.1754  
 Published online 10 January 2011

---

---

---

---

---

---

---

---

---

---

---

---

## hands-on : tophat

Example of used commands:

```
bowtie2-build ITAG2.3_genomic_Ch6.fasta index-bowtie2/tomato_chr6
```

```
qsub -N tophat_wt -pe parallel_smp 4 -b Y 'tophat2 -o aln_tophat_wt --max-intron-length 5000 --mate-inner-dist 200 bowtie2-index/tomato_chr6_WT_rep1_1_Ch6.fastq.gz WT_rep1_2_Ch6.fastq.gz'
```

```
samtools index file.bam
```

```
samtools view file.bam | cut -f 1 | sort | uniq -c | cut -c 1-7 | sort -n | uniq -c
```

Or

```
samtools flagstat file.bam
```

---

---

---

---

---

---

---

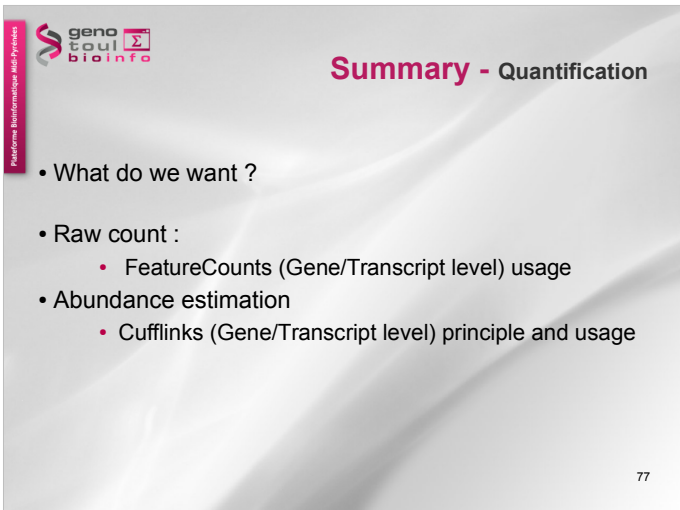
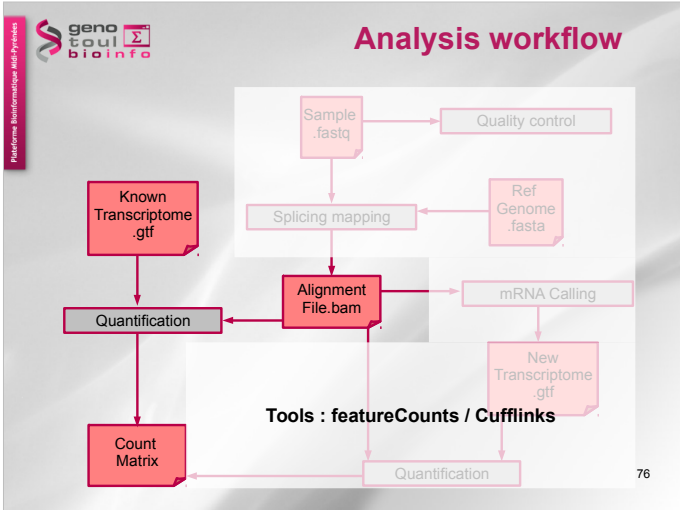
---

---

---

---

---



**What do we want to build?**

The gene / transcript description file (and corresponding fasta)

```

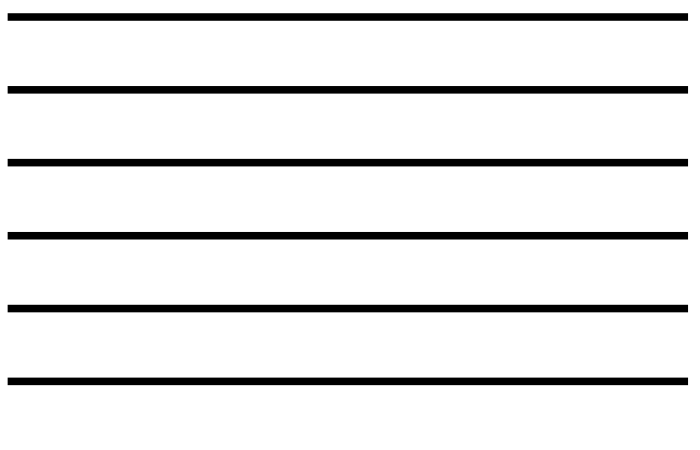
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "1"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "2"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "3"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "4"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "5"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "6"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "7"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "8"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "9"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "10"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "11"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "12"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "13"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "14"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "15"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "16"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "17"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "18"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "19"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "20"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "21"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "22"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "23"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "24"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "25"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "26"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "27"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "28"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "29"
> protein_coding exon 607702 607947 - - gene_id "ENSG00000277030" transcript_id "ENSG00000244202" exon_number "30"
  
```

The count file

```

raw.names SRR11377 SRR11378 SRR11379 SRR11380 SRR11381 SRR11382 SRR11383 SRR11384 SRR11385 SRR11386 SRR11387
1 raw_c1 1851 1851 4888 4888 1025 1025 2885 2885 6878 6878 6184 6184
2 raw_c2 393 818 828 828 854 854 393 789 441 1515 1515 152
3 raw_c3 1261 1261 2583 2583 1245 1245 2900 1246 4652 12612 1261 4258
4 raw_t1_1 897 789 1074 1100 1101 937 1051 898 1289 1331
5 raw_t1_2 2385 4737 6457 5312 4562 2399 7610 5134 8163 4758
6 raw_t1_3 289 148 927 107 106 247 889 522 168 264
7 raw_t1_4 824 1893 1874 1889 2000 891 1387 1351 1287 1352
8 raw_t1_5 782 1881 2738 2839 1075 720 2305 2512 2188 2643
9 raw_t1_6 111 117 684 886 895 346 639 581 1245 1099
10 raw_t1_7 289 148 927 107 106 247 889 522 168 264
11 raw_t1_8 111 117 684 886 895 346 639 581 1245 1099
12 raw_t1_9 538 513 944 1156 1175 515 1020 933 1447 1444
13 raw_t1_10 148 148 289 289 318 308 1134 2845 275 6144 1623
14 raw_t1_11 443 1849 1738 8576 9558 3954 8432 4339 8373 4758
15 raw_t1_12 148 148 289 289 318 308 1134 2845 275 6144 1623
16 raw_t1_13 443 1849 1738 8576 9558 3954 8432 4339 8373 4758
17 raw_t1_14 148 148 289 289 318 308 1134 2845 275 6144 1623
18 raw_t1_15 148 148 289 289 318 308 1134 2845 275 6144 1623
19 raw_t1_16 443 1849 1738 8576 9558 3954 8432 4339 8373 4758
20 raw_t1_17 148 148 289 289 318 308 1134 2845 275 6144 1623
21 raw_t1_18 443 1849 1738 8576 9558 3954 8432 4339 8373 4758
22 raw_t1_19 148 148 289 289 318 308 1134 2845 275 6144 1623
23 raw_t1_20 148 148 289 289 318 308 1134 2845 275 6144 1623
24 raw_t1_21 148 148 289 289 318 308 1134 2845 275 6144 1623
25 raw_t1_22 148 148 289 289 318 308 1134 2845 275 6144 1623
  
```

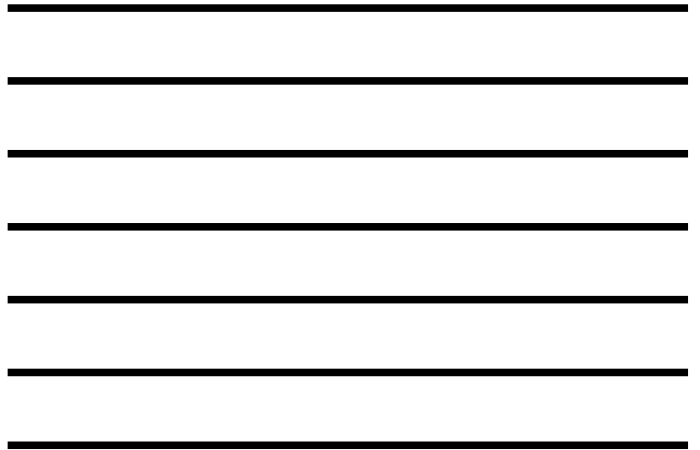
78





## featureCounts

- Q The minimum mapping quality score a read must satisfy in order to be counted. For paired-end reads, at least one end should satisfy this criteria. 0 by default.
- primary If specified, only primary alignments will be counted.
- minReadOverlap Specify the minimum number of overlapped bases required to assign a read to a feature. 1 by default.
- p If specified, fragments (or templates) will be counted instead of reads.
- P If specified, paired-end distance will be checked when assigning
- d Minimum fragment/template length, 50 by default.
- D Maximum fragment/template length, 600 by default.
- B If specified, only fragments that have both ends successfully aligned will be considered for summarization.



## Cufflinks in general

NATURE BIOTECHNOLOGY | RESEARCH | LETTER

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

Colin Trapnell, Brian A Williams, Geo Pertea, Ali Mortazavi, Gordon Kwan, Marjke J van Baren, Steven L Salzberg, Barbara J Wold & Lior Pachter

Affiliations | Contributions | Corresponding author

Nature Biotechnology 28, 511–515 (2010) | doi:10.1038/nbt.1621  
 Received 02 February 2010 | Accepted 22 March 2010 | Published online 02 May 2010

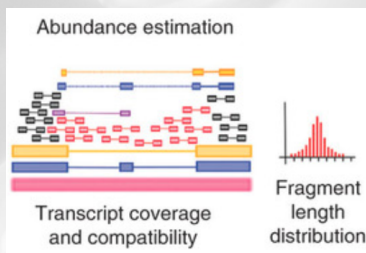
<http://cufflinks.cbc.umd.edu/>

- *assembles transcripts*
- **estimates their abundances** : based on how many reads support each one
- tests for differential expression in RNA-Seq samples



## Cufflinks read attribution

- Violet fragment: from which transcript?
  - Use of Fragment length distribution



Platforme Bioinformatique M&Pyrénées

geno toul bioinfo

## Cufflinks expression measurement

- Fragments attribution
- Isoforms abundances estimation:
  - RPKM for single reads
  - FPKM for paired-end reads

85

Trapnell C et al. Nature Biotechnology 2010;28:511-515

---

---

---

---

---

---

---

---

---

---

Platforme Bioinformatique M&Pyrénées

geno toul bioinfo

## RPKM / FPKM

- Transcript length bias
- **RPKM** : Reads per kilobase of exon per million mapped reads
  - R = Number of mapped reads
  - N = Total mapped reads
  - L = Exon gene length in bp

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

If my gene length (L) is : 200pb  
 Number of reads mapped (C) : 400  
 Total mapped reads (sum for all genes) (N) : 10<sup>8</sup>  
 RPKM = (10<sup>9</sup> \* 400) / (10<sup>8</sup> \* 200) = 20

- the transcript length depends on isoform inference
- **FPKM** : for paired-end sequencing
  - A pair of reads constitute one fragment

86

---

---

---

---

---

---

---

---

---

---

Platforme Bioinformatique M&Pyrénées

geno toul bioinfo

## Cufflinks inputs and options

- Command line:
  - `cufflinks [options]* <aligned_reads.(sam/bam)>`
- Some options :
  - h/--help
  - o/--output-dir
  - p/--num-threads
  - G/--GTF <reference\_annotation.(gtf/gff)> : estimate isoform expression, no assembly novel transcripts

87

---

---

---

---

---

---

---

---

---

---



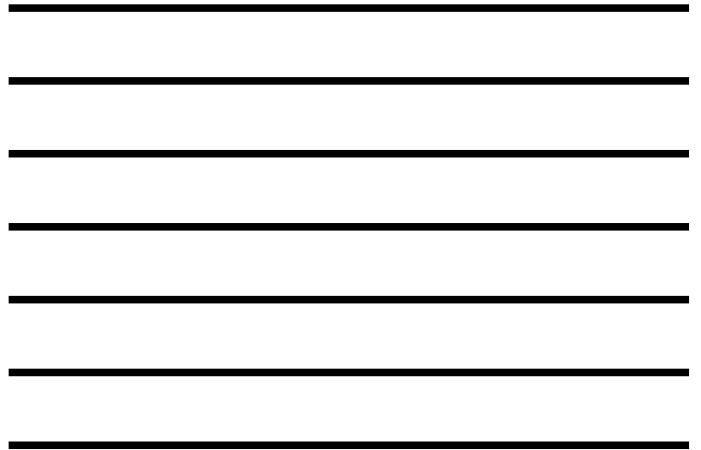
Platforme Bioinformatique M&P/Pythons

geno  
toul  
bioinfo

## Transcript reconstruction

Gene location ———  
 Exon location ———  
 Junctions :  
 - Between read pair junction   
 - Within read junction

91



Platforme Bioinformatique M&P/Pythons

geno  
toul  
bioinfo

## Model building strategies

b

Transcript graph ... Branch point 1 ... Branch point 2 ...

Maximal set

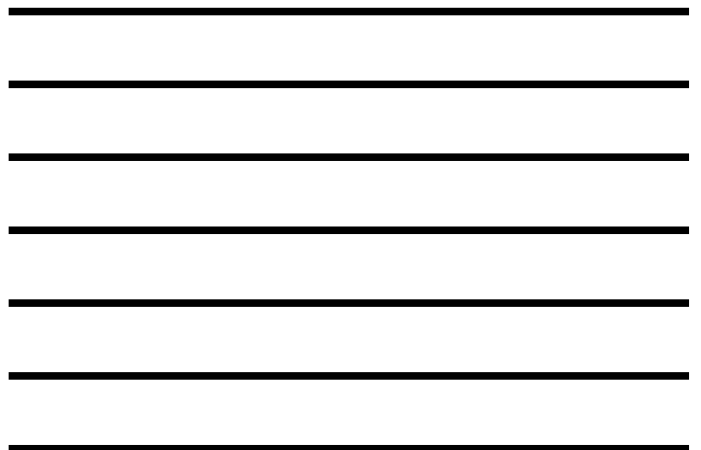
Minimal possible set 1

Minimal possible set 2

REVIEW

Computational methods for transcriptome annotation and quantification using RNA-seq  
 Manned Garber<sup>1</sup>, Manfred G Grabherr<sup>1</sup>, Mitchell Gutman<sup>1,2</sup> & Cole Trapnell<sup>1,2</sup>

92



Platforme Bioinformatique M&P/Pythons

geno  
toul  
bioinfo

## Cufflinks

NATURE BIOTECHNOLOGY | RESEARCH | LETTER

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

Cole Trapnell, Brian A Williams, Geo Pertea, Ali Mortazavi, Gordon Kwan, Marijke J van Baren, Steven L Salzberg, Barbara J Wold & Lior Pachter

Affiliations | Contributions | Corresponding author

Nature Biotechnology 28, 511–515 (2010) | doi:10.1038/nbt.1621  
 Received 02 February 2010 | Accepted 22 March 2010 | Published online 02 May 2010

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

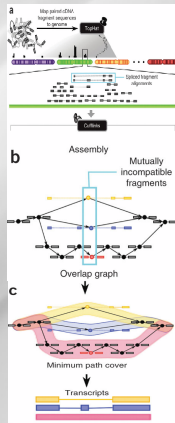
- **assembles transcripts**
- estimates their abundances : based on how many reads support each one
- tests for differential expression in RNA-Seq samples

93

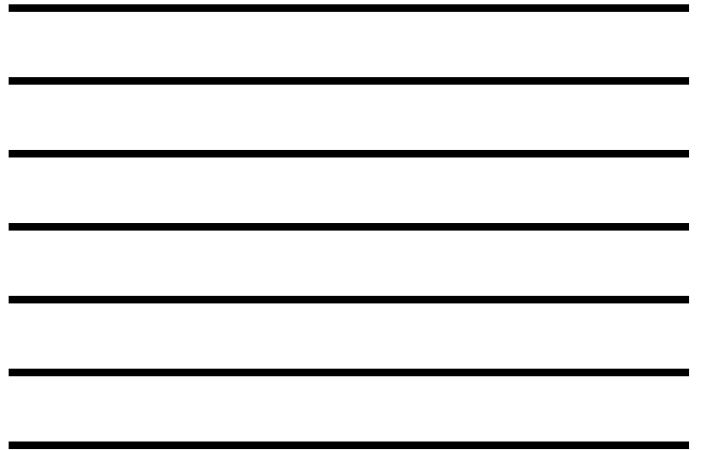


## Cufflinks transcript assembly

- Transcripts assembly :
  - Fragments are divided into non-overlapping loci
  - each locus is assembled independently :
- Cufflinks assembler
  - find the mini nb of transcripts that explain the reads
  - find a minimum path cover ( Dilworth's theorem ) :
    - nb incompatible read = mini nb of transcripts needed
    - each path = set of mutually compatible fragments overlapping each other

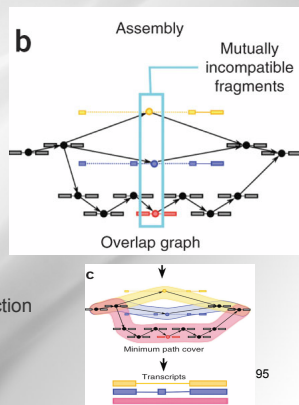


94



## Cufflinks transcript assembly

- Transcripts assembly :
  - Identification incompatibles fragments: distinct isoforms
  - Compatibles fragments are connected: graphe construction



95



## Cufflinks inputs and options

- Command line:
  - `cufflinks [options]* <aligned_reads.(sam/bam)>`
- Some options :
  - h/--help
  - o/--output-dir
  - p/--num-threads
  - G/--GTF <reference\_annotation.(gtf/gff)> : estimate isoform expression, no assembly novel transcripts
  - g/--GTF-guide <reference\_annotation.(gtf/gff)> : guide RABT (Reference Annotation Based Transcript) assembly

96





genotoul bioinfo

## Cufflinks RABT assembly option

– Some options :

**-g/--GTF-guide <reference\_annotation.(gtf/gff)>** : guide RABT assembly

Roberts A et al. Bioinformatics 2011;27:2325-2329 97



genotoul bioinfo

## Cufflinks outputs

- **transcripts.gtf** : contains assembled isoforms (coordinates and abundances)
- **genes.fpk\_tracking**: contains the genes FPKM
- **isoforms.fpk\_tracking**: contains the isoforms FPKM

98



genotoul bioinfo

## Cufflinks GTF description

- **transcripts.gtf** (coordinates and abundances): contains assembled isoforms: can be visualized with a genome viewer
  - GTF format + attributes (ids, FPKM, confidence interval bounds, depth or read coverage, all introns and exons covered)

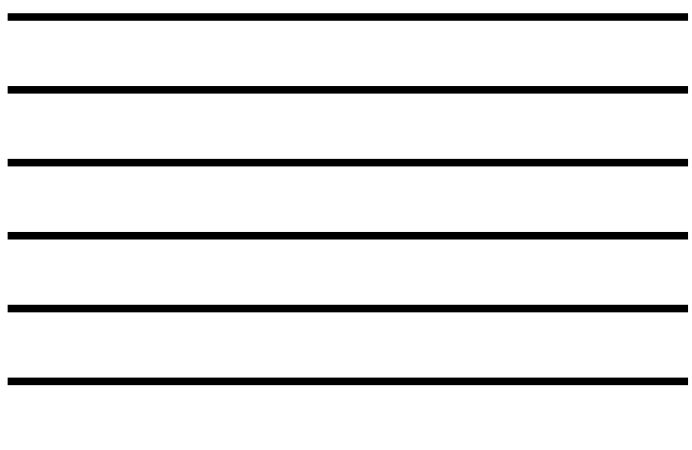
Chr	Source	Feature	Start	End	strand	Frame
22	Cufflinks	transcript	9743035	9747366	349-	.
22	Cufflinks	exon	9743035	9745254	349-	.

Attributes

Score: Most abundant isoform = 1000  
Minor : ratio=minor Fpkm/major FPKM

Whether or not all introns and exons were fully covered by Reads (with -g)

gene\_id "CUFF.560", transcript\_id "CUFF.560.1", FPKM "23.7787563790", frac "0.143485", conf\_lo "8.754478", conf\_hi "38.803035", cov "2.840328", full\_read\_support "yes", gene\_id "CUFF.560", transcript\_id "CUFF.560.1", exon\_number "1", FPKM "23.7787563790", frac "0.143485", conf\_lo "8.754478", conf\_hi "38.803035", cov "2.840328"



Platforme Bioinformatique M&P Pyrenees

geno  
toul  
bioinfo

## Cufflinks GTF description

- **transcripts.gtf** (coordinates and abundances): contains assembled isoforms: can be visualized with a genome viewer
- Exemple VISUALISATION IGV

100



Platforme Bioinformatique M&P Pyrenees

geno  
toul  
bioinfo

## Cuffcompare

- Comparison of transcriptoms files
- Command:  
cuffcompare -r <reference\_mrna.gtf> -o <outprefix> <input1.gtf> ...
- Outputs:
  - Overall summary statistics: <outprefix>.stats  
The Sn and Sp columns show specificity and sensitivity values at each level, while the fSn and fSp columns are "fuzzy" variants of these same accuracy calculations, allowing for a very small variation in exon boundaries to still be counted as a "match".
  - The "union" of all transfrags in all assemblies: <outprefix>.combined.gtf
  - Transfrags matching to each reference transcript: <cuff\_in>.refmap
  - Best reference transcript for each transfrag: <cuff\_in>.tmap
  - Tracking transfrags through multiple samples: <outprefix>.tracking

101



Platforme Bioinformatique M&P Pyrenees

geno  
toul  
bioinfo

## Cuffcompare

### Class code de cuffcompare

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

[http://cufflinks.cbc.umd.edu/manual.html#class\\_codes](http://cufflinks.cbc.umd.edu/manual.html#class_codes) 102



Platforme Bioinformatique M&Pyrénées

geno  
toul  
bioinfo

## Gene discovery pipeline

```

graph TD
    A[Alignment (Tophat)] --> B[Bam merge (samtools)]
    B --> C[Discovery of novels features (cufflinks)]
    C --> D[Quantification file with (featureCounts)]
  
```

103

---

---

---

---

---

---

---

---

Platforme Bioinformatique M&Pyrénées

geno  
toul  
bioinfo

## Hands-on : cufflinks

Commands :

Merge all bam :

```
samtools merge merge_all.bam file1.bam file2.bam
```

Cufflinks command:

```
cufflinks -p 4 --output-dir=cufflinks
-g reference_transcript.gtf
merge_all.bam
```

Cuffcompare command :

```
cuffcompare -f reference_transcript.gtf -o compare cufflink_transcripts.gtf
```

104

---

---

---

---

---

---

---

---

Platforme Bioinformatique M&Pyrénées

geno  
toul  
bioinfo

## Analysis workflow

```

graph TD
    A[Sample .fastq] --> B[Quality control]
    A --> C[Splicing mapping]
    D[Ref Genome .fasta] --> C
    C --> E[Alignment File .bam]
    E --> F[mRNA Calling]
    G[Known Transcriptome .gtf] --> F
    F --> H[New Transcriptome .gtf]
    H --> I[Quantification]
    E --> I
    J[Known Transcriptome .gtf] --> K[Quantification]
    K --> L[Count Matrix]
    I --> L
  
```

105

---

---

---


---

---

---

---

---


**Differential expression**

- Biostatistics Genotoul Platform
- Training :
  - <http://perso.math.univ-toulouse.fr/biostat/category/formation/>
  - Tutotial of RNAseq analysis [www.nathalievilla.org/teaching/naseq.html](http://www.nathalievilla.org/teaching/naseq.html)
- R scripts available on Genotoul cluster
  - See <http://bioinfo.genotoul.fr/index.php?id=119>

106

---

---

---

---

---

---

---


---

---

---

---

---


**Differential expression**

**But :** trouver les *gènes significativement* différentiellement exprimés entre 2 conditions.

**Méthode:**

- Normalisation
- Estimation de l'expression
- Test

**Outils :**

- DESeq, **EdgeR**, DESeq2, etc. (en R)
- **CuffDiff** (suite Tuxedo)

107

---

---

---

---

---

---

---


---

---

---

---

---


**Differential expression Normalization**

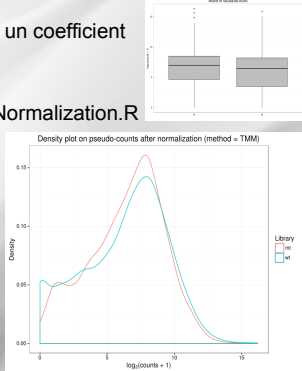
**Problème:** Le nombre de lectures est différent d'un réplicat à l'autre.

**Idée:** Appliquer à chaque échantillon un coefficient multiplicatif.

**Command :**

```
Rscript /usr/local/bioinfo/Scripts/bin/Normalization.R
-f tomato_count.R.txt
-o norm
```

Script that perform edgeR normalizations (RLE, TMM upperquartile) and provide graphs




---

---

---

---

---

---

---

---

---

---

---

---

## Differential expression Test

**But :** Comparer les distributions d'expression dans 2 conditions.

**Résultats:** p-value et q-value (p-value avec correction de tests multiples)

**Command :**

```
Rscript /usr/local/bioinfo/Scripts/bin/DEG.R
-f tomato_count.R.txt
--norm norm/RLE_info.txt
--pool1 mt
--pool2 wt
-o DEG
```

---

---

---

---

---

---

---

---

---

---

---

---

## Differential expression Test

```
Options:
-f CHARACTER, --file=CHARACTER
  tabulated raw count matrix with library name as header
  (e.g. #gene_id lib1 lib2...)
-n CHARACTER, --norm=CHARACTER
  file with normalized factors with library name
  (e.g. sample.name lib.size norm.factors). This file is obtained
  with script 'Normalization.R'
-o CHARACTER, --out=CHARACTER
  folder path where results are stored
--pool1=CHARACTER
  library name in pool 1 separated by ','
--pool2=CHARACTER
  library name in pool 2 separated by ','
--filter=CHARACTER
  if TRUE low expressed genes are removed [default=TRUE]
--alpha=CHARACTER
  significance level of the tests (i.e. acceptable rate of
  false-positive in the list of differentially expressed genes)
  [default=0.05]
--correct=CHARACTER
  method used to adjust p-values for multiple testing
  ('BH', 'BY' or 'fdr') [default=BH]
--Mplots=CHARACTER
  if TRUE all MA plots are saved [default=FALSE]
```

---

---

---

---

---

---

---

---

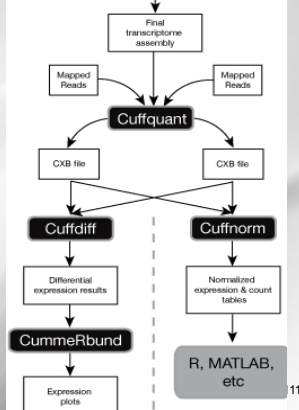
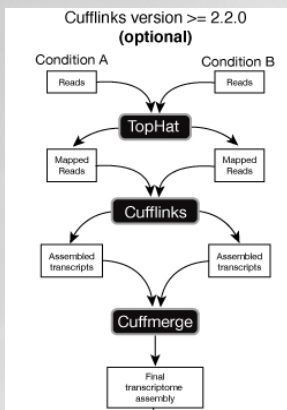
---

---

---

---

## Cufflinks >= 2.2.0




---

---

---

---

---

---

---

---

---

---

---

---

## Quality for Bioinfo Platform!

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

---

---

---

---

---

---

---

---

## Useful links

Seqanswer: <http://seqanswers.com/>  
RNAseq blog: <http://rna-seqblog.com/>  
Illumina: <http://www.illumina.com/>

---

---

---

---

---

---

---

---