


Plateforme Bioinformatique Midi-Pyrénées



**F12d**  
**RNA-Seq data analysis**

Delphine Labourdette Get-Biopuce / Christophe Klopp Bioinfo Genotoul

1

---

---

---

---

---

---

---

---

Plateforme Bioinformatique Midi-Pyrénées



**Session organisation**

<b>Morning (9h30 -12h30) :</b>	<b>Afternoon (14h-17h) :</b>
- Sequence quality	- mRNA calling
- Theory + exercises	- Theory + exercises
- Spliced read mapping	- expression measurement
Visualisation	- Theory + exercises
- Theory + exercises	

2

---

---

---

---


---

---

---

---

Plateforme Bioinformatique Midi-Pyrénées



**Material**

Slides

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

Exercise leaflet

Data, results files and a unix command reminder  
<http://bioinfo.genotoul.fr/index.php?id=119>

3

---

---

---

---

---

---

---

---

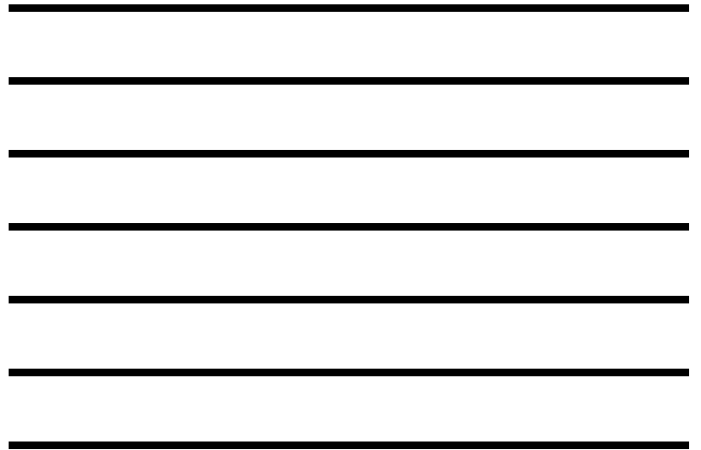
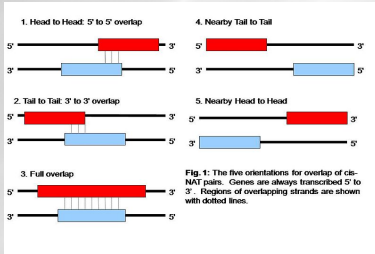




## Cis-natural antisense transcript

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

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



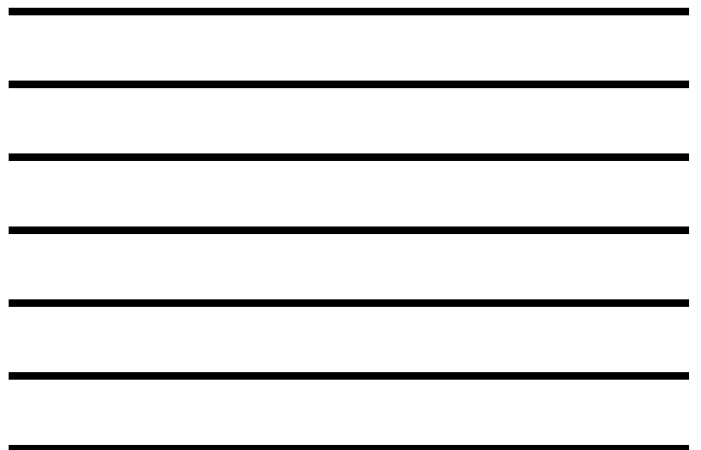
## Fusion genes

- 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.
- 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.

Genome Biol. 2011 Jan 19;12(1):R6. [Epub ahead of print]  
**Identification of fusion genes in breast cancer by paired-end RNA-sequencing.**  
 Edgren H, Murumalau A, Kangasvesika S, Nicorici D, Honaiisto V, Kleivi K, Rye IH, Nuber S, Wolf M, Borresen-Dale AL, Kallioniemi O  
 Institute for Molecular Medicine Finland (FIMM), Tukholmankatu 8, Helsinki, 00290, Finland. olli.kallioniemi@fimm.fi

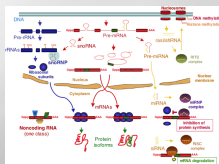


[http://en.wikipedia.org/wiki/Fusion\\_gene](http://en.wikipedia.org/wiki/Fusion_gene)  
<http://en.wikipedia.org/wiki/Trans-splicing>



## Transcriptome variability

- Number of transcripts
  - possible variation factor between transcripts:  $10^6$  or more,
  - expression variation between samples.
- Many types of transcripts
  - mRNA, ncRNA,...
- Isoforms (with non canonical splice sites)
- Intron retention
  - The splicing is not always completed
  - Is a new isoform or a transcription error
- Transcript decay (degradation)
- Allele specific expression



[http://www.nature.com/embojournal/v25/n5/fig\\_tab/7601023a\\_F2.html](http://www.nature.com/embojournal/v25/n5/fig_tab/7601023a_F2.html)



genotoul bioinfo

## How can we study the transcriptome?

Lots of techniques

- EST (Expressed sequence tags)
- PCR (polymerase chain reaction)
- SAGE (serial analysis of gene expression)
- Micro-Arrays
  - Different types: spotting, synthesis
  - Different densities : few thousands up 4M probes / slide

13

---

---

---

---

---

---

---

---

---

---

genotoul bioinfo

## Techniques classification ?

EST	PCR/RT-QPCR	SAGE	MicroArrays
No quantification	quantification		indirect
low throughput	low throughput (up to hundreds)	High (up to thousands)	High throughput (up to millions)
Discovery (Yes)	no		no (except tiling)

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

14

---

---

---

---

---

---

---

---

---

---

genotoul bioinfo

## What is RNA-Seq ?

- use of **high-throughput sequencing technologies** to sequence cDNA in order to get information about a sample's RNA content
- Thanks to the deep coverage and base level resolution provided by next-generation sequencing instruments, RNA-seq provides researchers with efficient ways to measure transcriptome data experimentally

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

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

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.

15

---

---

---

---

---

---

---

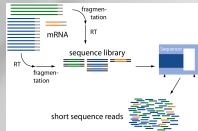
---

---

---

## 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, ...




---

---

---

---

---

---

---

---

---

---

---

---

## RNA-Seq platforms comparison

Platform	454 Roche Titanium	HiSeq2000 Illumina	Solid 3+ Life Technologies
Characteristics	-Titanium chemistry -Pyrosequencing -PCR amplification	- Polymerase-based sequence-by-synthesis -PCR amplification -Multiplexing	-ligation-base-sequencing -PCR amplification
Applications	-De novo sequencing -Small genomes -Transcriptome	-Resequencing -Transcriptome -Epigenomic -Small RNA -Allele specific sequencing	-De novo sequencing -Resequencing -Transcriptome -Epigenomic -Small RNA
Paired end separation	Not used	200bp	200bp
Mb / run	800Mb	600Gb	60Gb
Read length	800 bp	100bp	50bp
Known Biases	- Long homopolymer - makes signal saturation - read duplication	- Rich GC or AT regions: under-representation during amplification - Most error in end of cycle	- read duplication ?

---

---

---

---

---

---

---

---

---

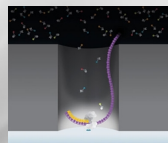
---

---

---

## Third Generation RNA-Seq

- No more amplification
- Single Molecule Sequencing Technology (tSMS)
- Single Molecule Real Time (SMRT) sequencing technology (PacBio RS)




---

---

---

---

---

---

---

---

---

---

---

---

## What are we looking for?

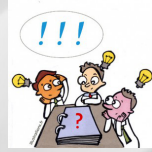
### Identify genes

- List new genes

### Identify transcripts

- List new alternative splice forms

Quantify these elements → differential expression



---

---

---

---

---

---

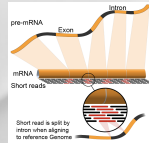
---

---

## Different approaches :

### Alignment to

- De novo
  - No reference genome, no transcriptome available
  - Very expensive computationally
  - Lots of variation in results depending on the software used
- Reference transcriptome
  - Most are incomplete
  - Computationally inexpensive
- Reference genome
  - When available
  - Allow reads to align to unannotated sites
  - Computationally expensive
  - Need a spliced aligner



---

---

---

---

---

---

---

---

## 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 ?

---

---

---

---

---


---

---

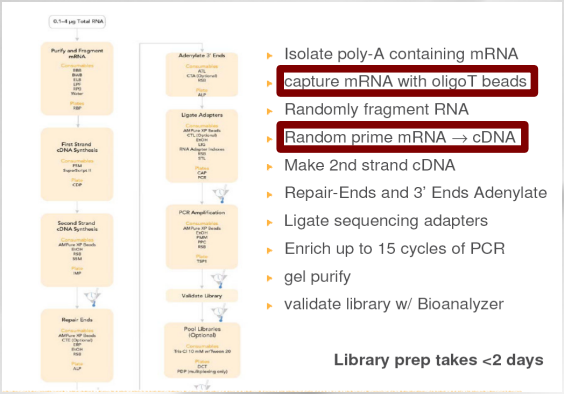
---



Platforme Bioinformatique des Pythéas



## RNA-Seq library preparation




- Isolate poly-A containing mRNA
- capture mRNA with oligoT beads
- Randomly fragment RNA
- Random prime mRNA → cDNA
- Make 2nd strand cDNA
- Repair-Ends and 3' Ends Adenylate
- Ligate sequencing adapters
- Enrich up to 15 cycles of PCR
- gel purify
- validate library w/ Bioanalyzer

Library prep takes <2 days

25

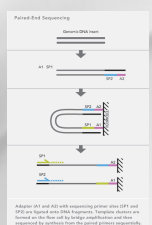


Platforme Bioinformatique des Pythéas

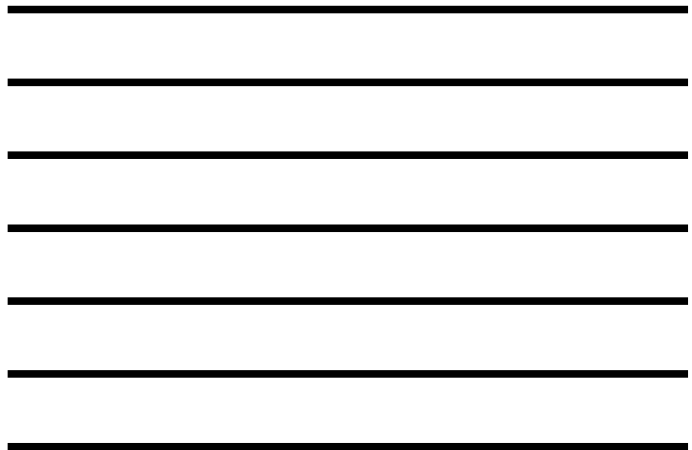


## Paired-end sequencing


- modification to 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



26



Platforme Bioinformatique des Pythéas



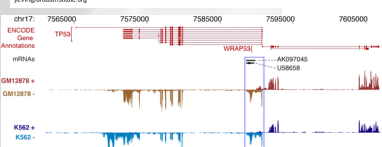
## Strand specific RNA-Seq protocol

workflow comparison:  
mRNA-Seq vs directional mRNA-Seq

- start with 1ug (or less) total RNA
- purify poly-A mRNA
- randomly fragment mRNA
- 1st strand cDNA synthesis
- 2nd strand cDNA synthesis
- end repair
- adenylate 3' ends
- ligate adaptors
- gel purify
- end repair with phosphatase and PNK
- column purify PNK treated mRNA
- ligate 3' adaptor
- ligate 5' adaptor
- reverse transcribe
- enrich with PCR
- validate library
- group clusters
- sequence on HiSeq2000 (SR or PE)

doi:10.1371/journal.pbio.1001046

Comprehensive comparative analysis of strand-specific RNA sequencing methods.  
Levin JZ, Yassour M, Adiconis X, Nusbaum C, Thompson DA, Friedman N, Gertke A, Regev A  
Broad Institute of Massachusetts Institute of Technology and Harvard University, Cambridge, Massachusetts, USA  
jlewin@broadinstitute.org




27

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



Platforme Bioinformatique des Pyrenees



## Analysis workflow

Data quality control

Spliced mapping

Gene and transcript discovery

Quantification

28

---

---

---

---

---

---


---

---

---

---

Platforme Bioinformatique des Pyrenees



## RNAseq specific bias

- Influence of the library preparation
- Random hexamer priming
- Coverage across the transcript may not be random
- Transcript length bias
- Positional bias and sequence specificity bias.

*Robert et al. Genome Biology, 2011,12:R22*

- BWA can force incorrect alignments to meet pairing criteria
- Some reads map to multiple locations

29

---

---

---

---

---

---


---

---

---

---

Platforme Bioinformatique des Pyrenees



## Hexamer random priming bias

Published online 14 April 2010      Nucleic Acids Research, 2010, Vol. 38, No. 12, e111  
doi:10.1093/nar/gkq124

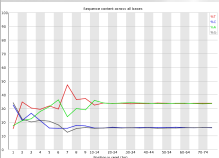
**Biases in Illumina transcriptome sequencing caused by random hexamer priming**

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.

- There is a strong distinctive pattern in the nucleotide frequencies of the first 13 positions at the end of mapped RNA-Seq reads:
  - 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



30

---

---

---

---

---

---


---

---

---

---

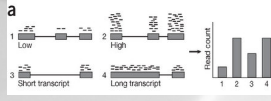
Platforme Bioinformatique M&P Pythons



## Transcript length bias

BioDirect, 2009 Apr 16:4-14  
**Transcript length bias in RNA-seq data confounds systems biology.**  
 Ostlack A, Wakefield M.

**Abstract**  
**Background:** Several recent studies have demonstrated the effectiveness of deep sequencing for transcriptome analysis (RNA-seq) in mammals. As RNA-seq becomes more affordable, whole genome transcriptional profiling is likely to become the platform of choice for species with good genomic sequences. As yet, a rigorous analysis methodology has not been developed and we are still in the stages of exploring the features of the data.  
**Results:** We investigated the effect of transcript length bias in RNA-seq data using three different published data sets. For standard analyses using aggregated tag counts for each gene, the ability to call differentially expressed genes between samples is strongly associated with the length of the transcript.  
**Conclusions:** Transcript length bias for calling differentially expressed genes is a general feature of current protocols for RNA-seq technology. This has implications for the ranking of differentially expressed genes, and in particular may introduce bias in gene set testing for pathway analysis and other multi-gene systems biology analyses.  
**Reviewers:** This article was reviewed by Rohan Williams (nominated by Gavin Hastley), Nicole Cloonan (nominated by Mark Ragan) and James Bullard (nominated by Sandrine Dudot).




– 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  
 Lijun Gao<sup>1,†</sup>, Zhidong Fang<sup>2,†</sup>, Kui Zhang<sup>1</sup>, Dongli Zhu<sup>1</sup> and Xiangjin Cui<sup>1,\*</sup>

31



Platforme Bioinformatique M&P Pythons



## fastq file formats

Published online 16 December 2009  
 Nucleic Acids Research, 2010, Vol. 38, No. 6, 1767–1771  
 doi:10.1093/nar/gkp1137

**SURVEY AND SUMMARY**  
**The Sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants**

Peter J. A. Cock<sup>1,\*</sup>, Christopher J. Fields<sup>2</sup>, Naohisa Goto<sup>3</sup>, Michael L. Heuer<sup>4</sup> and Peter M. Rice<sup>5</sup>

Table 1. The three described FASTQ variants, with columns giving the description, OBF name, ASCII characters, format name used in OBF projects, range of ASCII characters permitted in the quality string (in decimal notation), ASCII encoding offset, type of quality score encoded and the possible range of scores

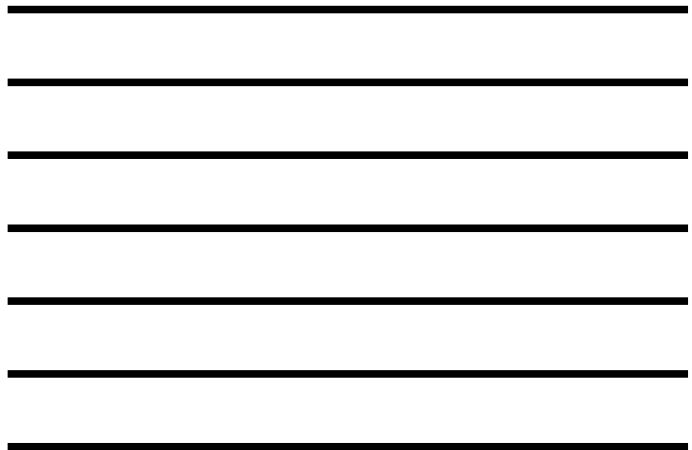
Description, OBF name	ASCII characters		Quality score	
	Range	Offset	Type	Range
Sanger standard fastq-sanger	33–126	33	PHRED	0 to 93
Solexa/Illumina fastq-solexa	59–126	64	Solexa	–5 to 62
Illumina 1.3+ fastq-illumina	64–126	64	PHRED	0 to 62

$$Q_{\text{PHRED}} = -10 \times \log_{10}(P_e)$$


$$Q_{\text{Solexa}} = -10 \times \log_{10}\left(\frac{P_e}{1 - P_e}\right)$$

```
@EAS54_6_R1_2_1_413_324
CCCTTCTGTCTTCAGCGTTTCCTC
+
??3?????????????7?????88
```

32



Platforme Bioinformatique M&P Pythons

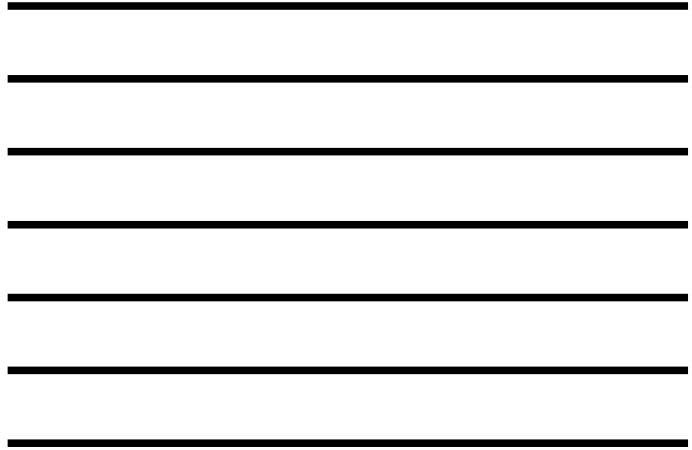


## Verifying RNA-Seq raw data

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

- Import of data from BAM, SAM or FastQ files
- quick overview
- Summary graphs and tables to quickly assess your data
- Export of results to an HTML report
- Offline operation to allow automated generation of reports
- Color code to check quickly the quality

33



**genotoul bioinfo**

## Verifying RNASeq raw data

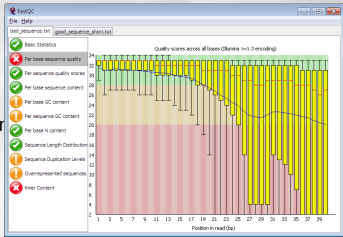
Per base sequence Q :

The higher the score the better the base call:

- very good quality calls (green)
- reasonable quality (orange)
- poor quality (red)

The quality of calls on most platforms will degrade as the run progresses

- ! The lower quartile for any base is less than 10 or the median for any base is less than 25
- ✖ The lower quartile for any base is less than 5 or the median for any base is less than 20



34

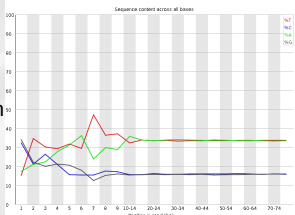


**genotoul bioinfo**

## Hexamer random priming bias

Per base sequence content :

Proportion of each base at each position  
In a random library : no differences, parallel lines



- ! difference between A and T, or G and C is greater than 10% in any position
- ✖ difference between A and T, or G and C is greater than 20% in any position

35

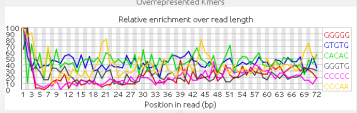


**genotoul bioinfo**

## Kmer content

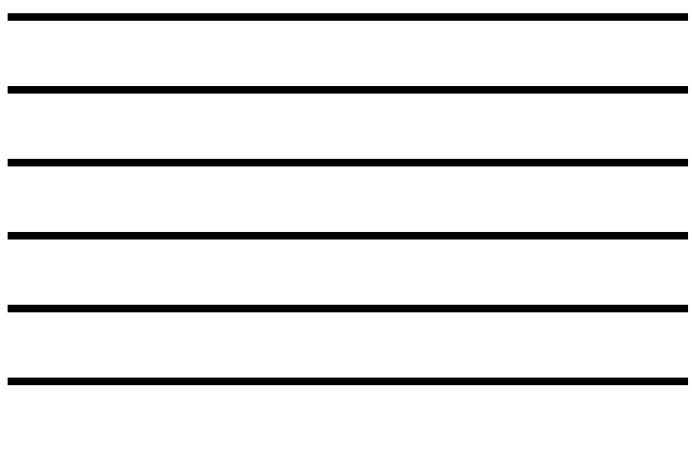
- A kmer is a subsequence of length k
- look for the enrichment of Kmer sequences. This should spot overrepresented sequences. Looking at 5-7mers should give a good impression of any contamination.
- Kmers showing a rise towards the end of the library indicate progressive contamination with adapters.
- Check for adaptor sequence or poly-A sequence

Overrepresented 5-mers



Sequence	Count	Obs/Exp overall	Obs/Exp Max	Max Obs/Exp P.
GGGGG	5120	3.485	15.4641	
GTGTG	17995	3.016	6.9311	
CAACA	17925	3	6.8351	
GGGTG	7905	2.67	10.5762	
CCCCC	3895	2.639	12.4931	


36










**Where to find a reference genome?**

Fasta file


Retrieving the genome file:

- 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>

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

46

---

---

---

---

---

---

---


---

---

---

---

---


**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.

47

---

---

---

---

---

---

---


---

---

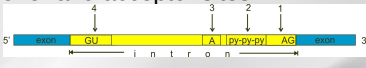
---

---

---


**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

<https://pubmed.ncbi.nlm.nih.gov/128621436/>  
**Analysis of canonical and non-canonical splice sites in mammalian genomes.**  
 Bursel M, Serebriouk S, Salovey V

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

48

---

---

---

---

---

---

---

---

---

---

---

---

genotoul bioinfo

## Spliced alignment

- The recognition of exon/intron junctions can be inferred from the reads that overlap the splicing sites. The resulting spliced reads can produce very short alignments, part of the read will not map contiguously to the reference.
  - therefore this approach requires a dedicated algorithm
- Generation :
  - Sim4
  - Seqanswers : <http://seqanswers.com/wiki/Software/list>
- Idea :
  - Database of potential splice junction sequences (known)
  - splice canonical / non canonical site search (see then mapping)

Genome Res. 1998 Sep 01;8(9):71-74  
A computer program for aligning a cDNA sequence with a genomic DNA sequence.  
Fioravanti A, Harkiss G, Zhanan Z, Rubin GM, Miller W  
Department of Computer Science and Engineering, The Pennsylvania State University, University Park, Pennsylvania 16802 USA.

49



genotoul bioinfo

## Alignment Tools

Tools for splice-mapping:

- Tophat
- Mapsplice:

BIOINFORMATICS ORIGINAL PAPER Vol. 25 no. 9 2008, pages 1105-1111 doi:10.1093/bioinformatics/btn150  
Sequence analysis  
TopHat: discovering splice junctions with RNA-Seq  
Cole Trapnell<sup>1</sup>\*, Lior Pachter<sup>2</sup> and Steven L. Salzberg<sup>1</sup>

Published online 28 August 2010  
Nucleic Acids Research, 2010, Vol. 38, No. 18, e178 doi:10.1093/nar/gkq222

**MapSplice: Accurate mapping of RNA-seq reads for splice junction discovery**  
Kai Wang<sup>1</sup>, Darshan Singh<sup>1</sup>, Zheng Zeng<sup>1</sup>, Stephen J. Coleman<sup>2</sup>, Yan Huang<sup>1</sup>, Glibo L. Sanchez<sup>1</sup>, Xaping He<sup>1</sup>, Piotr Mieczkowski<sup>1</sup>, Sara A. Grimm<sup>2</sup>, Charles M. Perou<sup>1</sup>, James N. MacLeod<sup>1</sup>, Derek Y. Chiang<sup>1</sup>, Jian F. Peiroi<sup>2</sup> and Jinze Liu<sup>1</sup>

<http://www.netlab.uky.edu/p/bioinfo/MapSplice>

50



genotoul bioinfo

## TopHat

BIOINFORMATICS ORIGINAL PAPER Vol. 25 no. 9 2008, pages 1105-1111 doi:10.1093/bioinformatics/btn150  
Sequence analysis  
TopHat: discovering splice junctions with RNA-Seq  
Cole Trapnell<sup>1</sup>\*, Lior Pachter<sup>2</sup> and Steven L. Salzberg<sup>1</sup>

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

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

51



genotoul bioinfo **TopHat**

– TopHat finds junctions by mapping reads to the reference:

- all reads are mapped to the reference genome using Bowtie
- reads not mapped to the genome are set aside as IUM (initially unmapped)
- low complexity reads are discarded
- for each read : allow until 20 alignments

Map reads to whole genome with Bowtie  
Collect initially unmappable reads

52  
Trapnell C et al. Bioinformatics 2009;25:1105-1111



genotoul bioinfo **TopHat**

– TopHat then assembles the mapped reads

– Define island: aggregates mapped reads in islands of candidate exons

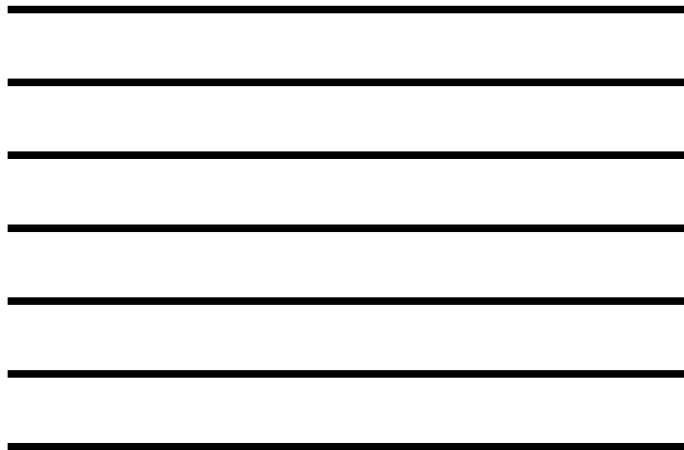
- Generate potential donor/acceptor splice sites using neighbouring exons

– Extend islands to cover eventually splice junctions

- +/- 45 bp from reference on either side of island

Map reads to whole genome with Bowtie  
Collect initially unmappable reads  
Assemble consensus of covered regions  
Generate possible splices between neighboring exons

53  
Trapnell C et al. Bioinformatics 2009;25:1105-1111



genotoul bioinfo **TopHat**

To map reads to splice junction :


- Enumerate all canonical donor and acceptor sites in islands
  - long ( $\geq 75$  bp) reads: "GT-AG", "GC-AG" and "AT-AC" introns
  - Shorter reads: only "GT-AG" introns
- Find all pairings which produce GT-AG introns between islands
  - $70 \text{ bp} < \text{Intron size} < 20,000 \text{ bp}$

Map reads to whole genome with Bowtie  
Collect initially unmappable reads  
Assemble consensus of covered regions  
Generate possible splices between neighboring exons

54  
Trapnell C et al. Bioinformatics 2009;25:1105-1111






**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

58

---

---

---

---

---

---

---


---

---

---

---

---


**TopHat Outputs**

Outputs :
 

- *accepted\_hits.bam* : list of read alignments in SAM format compressed
- *Junctions.bed* : track of junctions, scores : number of alignments spanning the junction
- *Insertions.bed* and *deletions.bed* : tracks of insertions and deletions
- Logs files
- Unmapped reads

59

---

---

---

---

---

---

---


---

---

---

---

---


**TopHat Alignment Strategies**

- Default options :
  - No reference genome
- Alignment with reference genome without discovery:
  - Options -G and --no-novel-juncs

60

---

---

---

---

---

---

---

---

---

---

---

---

genotoul bioinfo

## Sequence alignment and map

- SAM (Sequence Alignment/Map) format:
  - Capture all of the critical information about NGS data in a single indexed and compressed file
  - Sharing : data across and tools
  - Generic alignment format
  - SAMTOOLS: provide various utilities for manipulating alignments in the SAM format: sorting, merging, indexing...
 

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: 19505243

<http://samtools.sourceforge.net/>  
<http://picard.sourceforge.net/explain-flags.html>

61



genotoul bioinfo

## Spliced cigar line

- Extend CIGAR strings

Op	BAM	Description
M	0	alignment match (can be a sequence match or mismatch)
I	1	insertion to the reference
D	2	deletion from the reference
N	3	skipped region from the reference
S	4	soft clipping (clipped sequences present in SEQ)
H	5	hard clipping (clipped sequences NOT present in SEQ)
P	6	padding (silent deletion from padded reference)
=	7	sequence match
X	8	sequence mismatch

- Example: intron de 81 bases

```
ERR022486.8388510 81 22 32099 255 58M81N18M = 27484 -4772
CCTTGGTCTTGC CGAAGTAGATCTCATTGAGAGTGGAG CGGATCTTGTCTCCATTCCCTCCACC
AGGGGTCCGAT @AFADD;GDGAG@GGCBE@GG?GG?GGG?GGGGGGG NM:i:0 XS:A:- NH:i:1
```

62

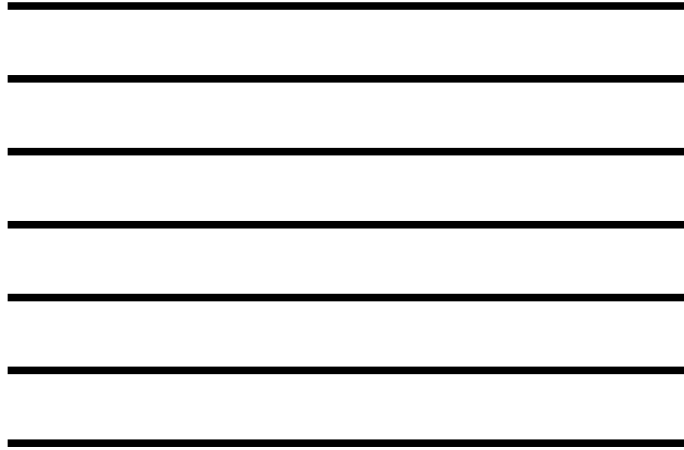


genotoul bioinfo

## Bam & Bed

- BAM (Binary Alignment/Map) format:
  - Compressed binary representation of SAM
  - Greatly reduces storage space requirements to about 27% of original SAM
  - Bamtools: reading, writing, and manipulating BAM files
- Bed (Browser Extensible Data) format:
  - tab-delimited text file that defines a feature track
  - <http://genome.ucsc.edu/FAQ/FAQformat.html#format1>
  - The first three required BED fields are: <chromosome> <start> <end>
  - 9 additional optional BED fields

63



genotoul bioinfo

### Bed example

Chrom Start End name score strand drawing RGB Blocks info

Chrom	Start	End	name	score	strand	drawing	RGB	Blocks info
22	241	1451	JUNC000000001	8	-	241, 1451	255, 0, 0	0, 226
22	1785	4260	JUNC000000002	1	-	1785, 4260	255, 0, 0	0, 2427
22	4285	4485	JUNC000000003	8	-	4285, 4485	255, 0, 0	0, 128
22	4575	4748	JUNC000000004	3	-	4575, 4748	255, 0, 0	0, 107
22	5834	6045	JUNC000000005	1	+	5834, 6045	255, 0, 0	0, 170
22	6143	6776	JUNC000000006	6	-	6143, 6776	255, 0, 0	0, 565
22	6796	7873	JUNC000000007	5	-	6796, 7873	255, 0, 0	0, 226
22	7843	7254	JUNC000000008	6	+	7843, 7254	255, 0, 0	0, 158
22	7220	8877	JUNC000000009	11	-	7220, 8877	255, 0, 0	0, 1595
22	7410	16244	JUNC000000010	2	-	7410, 16244	255, 0, 0	0, 8886
22	7638	7811	JUNC000000011	3	+	7638, 7811	255, 0, 0	0, 136
22	12388	21452	JUNC000000012	27	-	12388, 21452	255, 0, 0	0, 8090
22	16655	27319	JUNC000000013	6	-	16655, 27319	255, 0, 0	0, 10597
22	27711	30684	JUNC000000014	108	-	27711, 30684	255, 0, 0	0, 2901
22	27714	32151	JUNC000000015	303	+	27714, 32151	255, 0, 0	0, 4365
22	30639	32151	JUNC000000016	134	-	30639, 32151	255, 0, 0	0, 1440
22	32085	32388	JUNC000000017	493	-	32085, 32388	255, 0, 0	0, 152
22	32234	33112	JUNC000000018	478	+	32234, 33112	255, 0, 0	0, 686
22	33089	33347	JUNC000000019	292	-	33089, 33347	255, 0, 0	0, 187

64



genotoul bioinfo

### Know your transcriptome!

Ensembl 64  
[www.ensembl.org](http://www.ensembl.org)

65



genotoul bioinfo

### Numbers of exons

66



Platforme Bioinformatique M&Pyrénées

geno toul bioinfo

## Number of exons

number of exons per transcript comparison

— Cufflinks  
— Isoemsl

67

---

---

---

---

---

---

---

---

---

---

---

---

Platforme Bioinformatique M&Pyrénées

geno toul bioinfo

## TopHat technical issues

- Temporary disk space
  - 100 000 000 pair-ends = 0,5Gb 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)

68

---

---

---

---

---

---

---

---

---

---

---

---

Platforme Bioinformatique M&Pyrénées

geno toul bioinfo

## TopHat-Fusion

- an enhanced version of TopHat with the ability to align reads across fusion points
- identify fusions due to chromosomal rearrangements whether inter- or intra-chromosomal
- suggest that reads are at least 50-bp long, where a read is split into two segments (25-bp each)
- Both single and paired-end reads can be used and the output alignments are given in a modified SAM format with a new CIGAR\* operator 'F' to indicate fusion points

69

---

---

---

---

---

---

---

---


---

---

---

---

Platforme Bioinformatique M&P Pythons



**MapSplice**

*MapSplice: Accurate mapping of RNA-seq reads for splice junction discovery*  
 Kai Wang, Darshan Singh, Zheng Zeng, Stephen J. Coleman, Yan Huang, Gleb L. Savich, Xiaping He, Piotr Mieczkowski, Sara A. Grimm, Charles M. Perou, James N. MacLeod, Derek Y. Chiang, Jan F. Prins, Jinze Liu  
 Nucleic Acids Research 2010, doi: 10.1093/nar/gkq622


Features:

- splice junction discovery
- both CPU and memory efficiency.
- detection of small exons.
- discovery of canonical, semi-canonical and non-canonical junctions.
- splice inference based on the alignment quality and diversity of reads mapped to a junction.
- identification of chimeric events (intra-chromosomes and inter-chromosomes, inter-strands)
- support paired-end reads and single-end reads

70



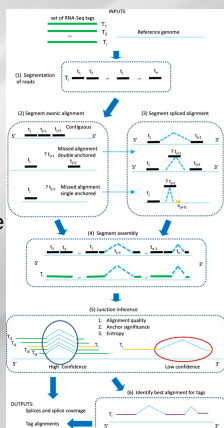
Platforme Bioinformatique M&P Pythons



**MapSplice**

Algorithm in two phases:

- tag alignment (Step 1–Step 4):  
 give a set of candidate alignments for each tag
- splice inference (Step 5–Step 6):  
 determine a splice significance score based on the quality and diversity of alignments that include the splice




INPUTS: set of RNA-Seq tags, Reference genome

OUTPUTS: Exons and splice coverage, Tag alignments



Platforme Bioinformatique M&P Pythons



**MapSplice**

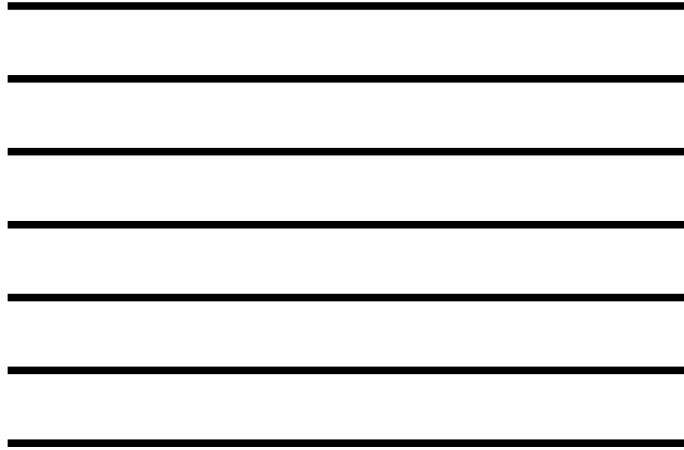
Command line:

- `python bin/mapsplice_segments.py [inputs|options] MapSplice.cfg`
- or
- `python bin/mapsplice_segments.py [inputs|options]`

Options:

- `-Q/--reads-format`: fa or fq
- `-o/--output-dl`
- `-c/--chromosome-files-dir`: directory containing the sequence files of the reference genome (in FASTA format)
- `-u/--reads-file`: for paired-end reads, the order should be: reads1\_end1, reads1\_end2, reads2\_end1, read2\_end2
- `-B/--Bowtieidx`: path and basename of index
- `--paired`
- `-L/--seglen`: Length of read segments in range of [18,25]
- .....

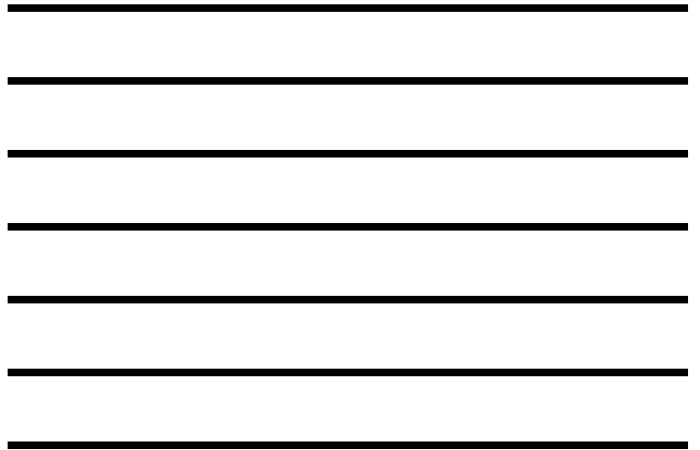
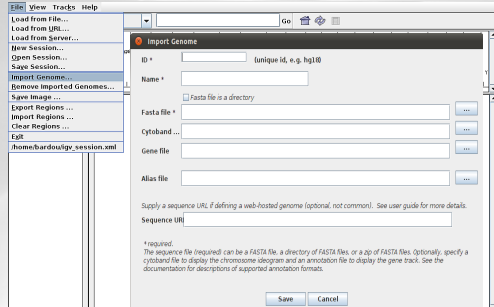
72





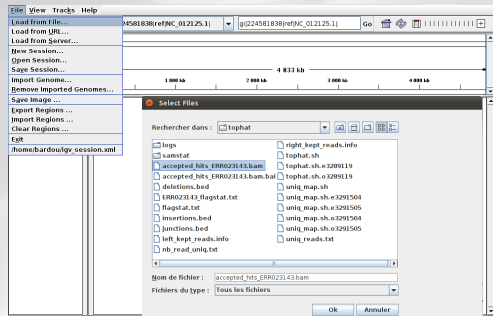
# Visualizing alignments on IGV

## - Import a reference genome



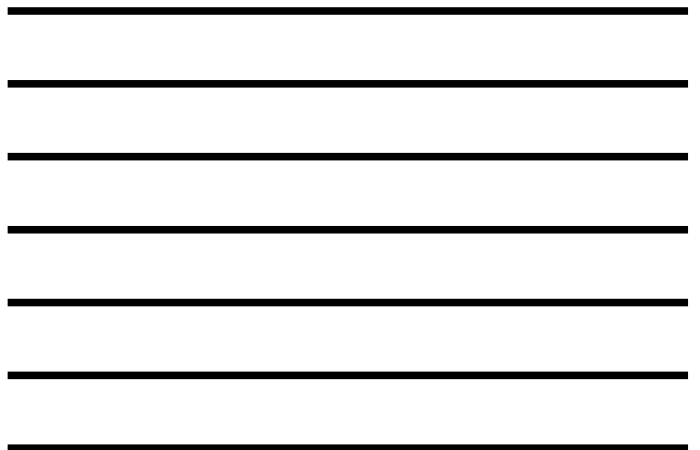
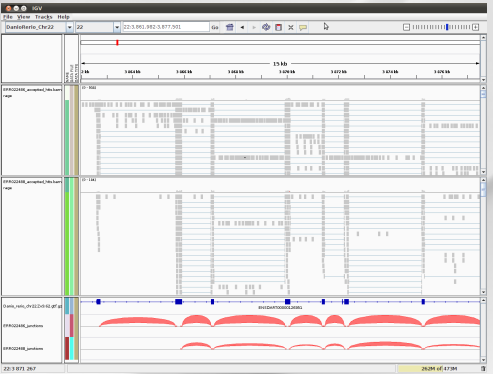
# Visualizing alignments on IGV

## - Import your BAM Files




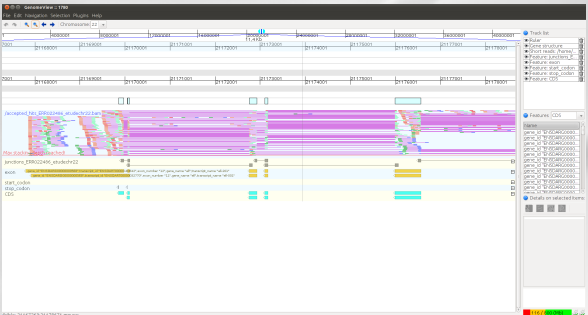
# Visualizing alignments on IGV

## - Exemple de visualisation de bam and bed files






**Exemple on GenomeView**



82

---

---

---

---

---

---

---


---

---

---

---

---


**Exemples et TP : tophat**

Commands guide :

```
bowtie-build sequence.fa index_name
```

```
tophat [options] <index_name> <reads_1,reads_2>
```

Options:

- h/--help
- O/--output-dir
- r/--mate-inner-dist : no default value
- G/--GTF <GTF2.2file>

Indexation:

```
samtools index accepted_hits.bam
```

Samtools flagstat:

```
samtools flagstat accepted_hits_ERR022486_chr22.bam
```

83

---

---

---

---

---

---

---


---

---

---

---

---


**Exemples et TP : tophat**

Example of used commands:

```
bowtie-build Danio_erio.Zv9.62.dna.chromosome.22.fa Danio_erio.Zv9.62_chr22
```

```
qsub -b y tophat --output-dir=tophat_ERR022486 -r 200 -G Danio_erio_chr22.Zv9.62.gtf /work/.../bowtie-index/Danio_erio.Zv9.62_chr22_ERR022486_read1.fastq ERR022486_read2.fastq
```

```
samtools index accepted_hits_ERR022486_chr22.bam
```

```
samtools view accepted_hits_ERR022486.bam | cut -f 1 | sort | uniq -c | cut -c 1-7 | sort -n | uniq -c
```

84

---

---

---

---

---

---

---

---


---

---

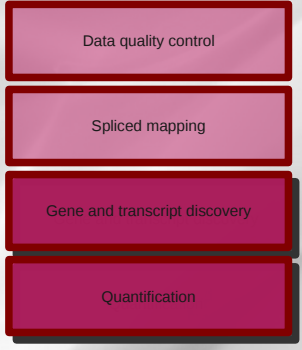
---

---

Platforme Bioinformatique M&P-Praxis



## Analyse workflow



85

---

---

---

---

---

---


---

---

---

---

Platforme Bioinformatique M&P-Praxis



## Transcript reconstruction

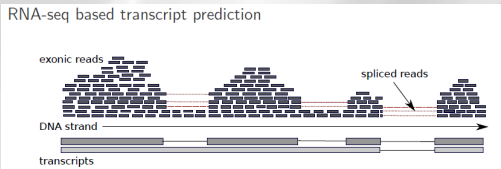
What mapping give us?

- Bam file with splicing localisation and pairing reads

What can we do then?

- Transcript assembly : structure determination
- Novel transcript discovery
- quantification

RNA-seq based transcript prediction



86

<http://www.fmi.tuebingen.mpg.de/raetsch/lectures/RECOMB-mTIM.pdf>

---

---

---

---

---

---


---

---

---

---

Platforme Bioinformatique M&P-Praxis



## Cufflinks in general

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

87

---

---

---

---

---

---

---

---

---

---

## Cufflinks tools

- Cufflinks (1 bam at the time)
  - Assembling : gene and transcript
- Cuffdiff (2 bam together)
  - Differential expression of genes
  - Differential promoter usage, splicing
- Cuffcompare: (n cufflinks results)
  - Comparing assembled transcripts to reference annotation
  - Comparing transcripts across time points
- Cuffmerge : (n cufflinks results)
  - Merge together several cufflinks assemblies
  - Run Cuffcompare

---

---

---

---

---

---

---

---

---

---

---

---

## Cufflinks tools

- Cufflinks (1 bam at the time)
  - Assembling : gene and transcript
- Cuffdiff (2 bam together)
  - Differential expression of genes
  - Differential promoter usage, splicing
- Cuffcompare: (n cufflinks results)
  - Comparing assembled transcripts to reference annotation
  - Comparing transcripts across time points
- Cuffmerge : (n cufflinks results)
  - Merge together several cufflinks assemblies
  - Run Cuffcompare

---

---

---

---

---

---

---

---

---

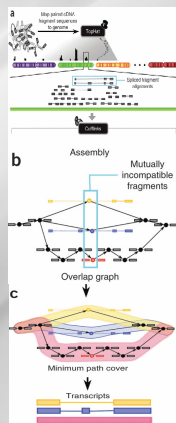
---

---

---

## 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




---

---

---

---

---

---

---

---

---

---

---

---

**genotoul bioinfo** **Cufflinks transcript assembly**

– Transcripts assembly :

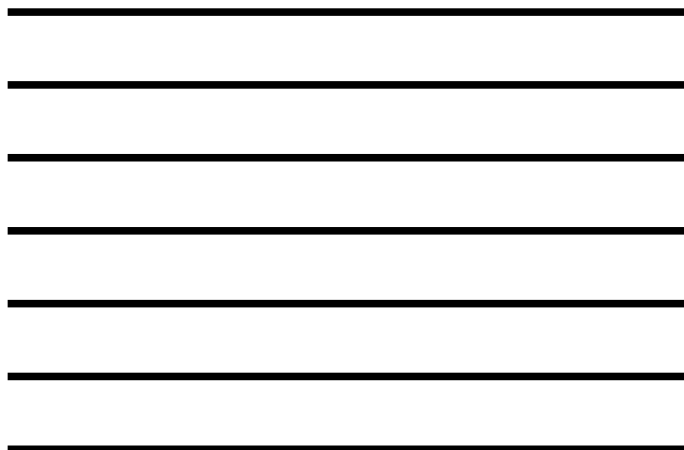
- Identification incompatibles fragments: distinct isoforms
- Compatibles fragments are connected: graphe construction

**b** Assembly  
Mutually incompatible fragments  
Overlap graph

**c** Minimum path cover  
Transcripts

91

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



**genotoul bioinfo** **Cufflinks expression measurement**

– Transcripts abundances:

- statistical model to derive a likelihood function for the abundances
- bayesian inference for isoform with lower abundances, or nb of read very small

→ can estimate the abundances of isoforms

– Fragment length distribution :

- single-end read : no way, use approximate Gaussian distribution
- paired-end reads +assembly : learn the distribution from reads mapped to single isoform genes
- paired-end – no assembly : genomic length of de paired-end reads (without splice) or Gaussian approximation

**d** Abundance estimation  
Transcript coverage and compatibility  
Fragment length distribution

**e** Maximum likelihood abundances

92

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

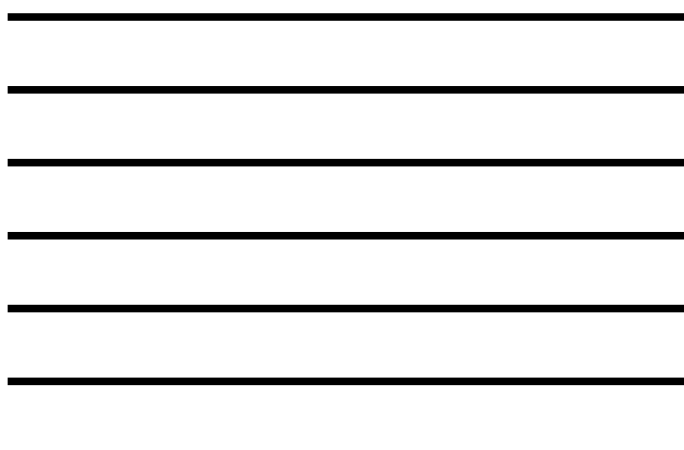


**genotoul bioinfo** **Cufflinks read attribution**

– Violet fragment: from which transcript?

- Use of Fragment length distribution

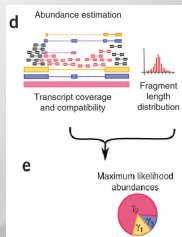
**d** Abundance estimation  
Transcript coverage and compatibility  
Fragment length distribution



genotoul bioinfo

## Cufflinks expression measurement

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



94

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

---

---

---

---

---

---

---

---

---

---

---

---

genotoul bioinfo

## RPKM / FPKM

- Transcript length bias
- **RPKM** : Reads per kilobase of exon per million mapped reads
  - 1kb transcript with 1000 alignments in a sample of 10 million reads (out of which 8 million reads can be mapped) will have:  

$$RPKM = 1000 / (1 * 8) = 125$$
- the transcript length depends on isoform inference
- **FPKM** : for paired-end sequencing
  - A pair of reads constitute one fragment

95

---

---

---

---

---

---

---

---

---

---

---

---

genotoul 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
  - g/--GTF-guide <reference\_annotation.(gtf/gff)> : guide RABT (Reference Annotation Based Transcript) assembly

96

---

---

---

---

---

---

---

---

---

---

---

---

Platforme Bioinformatique Multi-Pratiques

geno toul 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



Platforme Bioinformatique Multi-Pratiques

geno toul bioinfo

## Cufflinks outputs

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

98



Platforme Bioinformatique Multi-Pratiques

geno toul bioinfo

## Cufflinks GTF description

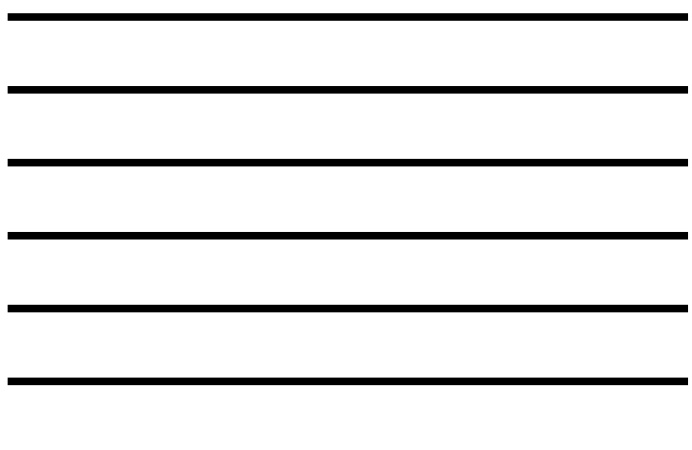
- **transcripts.gtf** (coordinates and abundances): contains assembled isoforms: can be visualized with a genome viewer
  - GTF format + attributes (ids, FPKM, confidence interval bounds, depth or read coverage, all introns and exons covered)

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

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",



**genotoul bioinfo**

## Cufflinks GTF description

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



**genotoul bioinfo**

## Cufflinks tracking description

- **genes.fpkm\_tracking**:
  - contains the estimated gene-level expression values in the generic FPKM Tracking Format

tracking_id	class_code	nearest_ref_id	gene_id	gene_short_name	tss_id	locus	length	coverage	status	FPKM	FPKM_conf_lo	FPKM_conf_hi
CUFF.560	-	-	CUFF.560	-	-	22.9743034-9782309	-	-	OK	106.69	77.9404	133.439

Quantification status

- **isoforms.fpkm\_tracking**: contains the estimated isoform-level expression values in the generic FPKM Tracking Format

tracking_id	class_code	nearest_ref_id	gene_id	gene_short_name	tss_id	locus	length	coverage	status	FPKM	FPKM_conf_lo	FPKM_conf_hi
CUFF.560.1	-	-	CUFF.560	-	-	22.9743034-9747366	24962.84033	OK	23.7788	6.75448	38.803	
CUFF.560.2	-	-	CUFF.560	-	-	22.9743034-9782309	40208.11967	OK	67.9765	50.3804	85.5727	
CUFF.560.3	-	-	CUFF.560	-	-	22.9743034-9782309	3846.166444	OK	13.9344		0.292533	

101



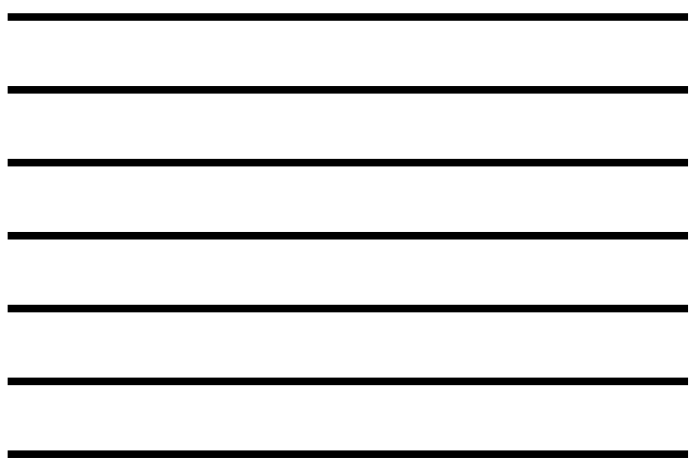
**genotoul bioinfo**


## Using Cufflinks software pieces

```

graph TD
    Tophat --> DiffExp[Identification of differential expression]
    Tophat --> Novel[Discovery of novels features : Cufflinks, Cuffmerge, cuffcompare]
    DiffExp --> DiffExpWith[With discovery : Cufflinks, cuffmerge, cuffdiff]
    DiffExp --> DiffExpWithout[Without discovery : cuffdiff]
  
```

102




**Cuffcompare**

- Help to analyse the transcript assembly
- Compare your assembled transcripts to a reference annotation
- Track Cufflinks transcripts across multiple experiments

- Inputs:

- Cufflinks' GTF output
- optionally can take a "reference" annotation

103

---

---

---

---

---

---

---


---

---

---

---

---


**Cuffcompare command line / options**

- Command line :
  - `cuffcompare [options]* <cuff1.gtf> [cuff2.gtf] ... [cuffN.gtf]`
- Some options:
  - o <outprefix>
  - r : optional "reference" annotation GFF file
  - R: (If -r specified) ignore reference transcripts that are not overlapped by any transcript in one of cuff1.gtf,...,cuffN.gtf

104

---

---

---

---

---

---

---


---

---

---

---

---


**Cuffcompare output**

- Main output:
  - `<outprefix>.tracking`: matches transcripts up between samples

```

TCONS_0001660 XLOC_000895 - u q1:CUFF.560/CUFF.560.11350|0.000000|0.000000|0.000000|0.000000|
TCONS_0001661 XLOC_000895 - q1:CUFF.560/CUFF.560.21100|0.000000|0.000000|0.000000|0.000000|4020 q2:CUFF.520/CUFF.520.11100|0.000000|0.000000|0.000000|0.000000|4065
TCONS_0001662 XLOC_000895 - q1:CUFF.560/CUFF.560.3120|0.000000|0.000000|0.000000|0.000000|3846 q2:CUFF.520/CUFF.520.21170|0.000000|0.000000|0.000000|0.000000|3891
      
```

Cufflinks transfrag id    Cufflinks locus id    **Class code**    Sample 1 :    Sample 2  
 Ref gene id | ref transcript id    <gene\_id>|<transcript\_id>|<FMI>|<FPKM>|<conf\_lo>|<conf\_hi>|<cov>|<len>

105

---

---

---

---

---

---

---

---

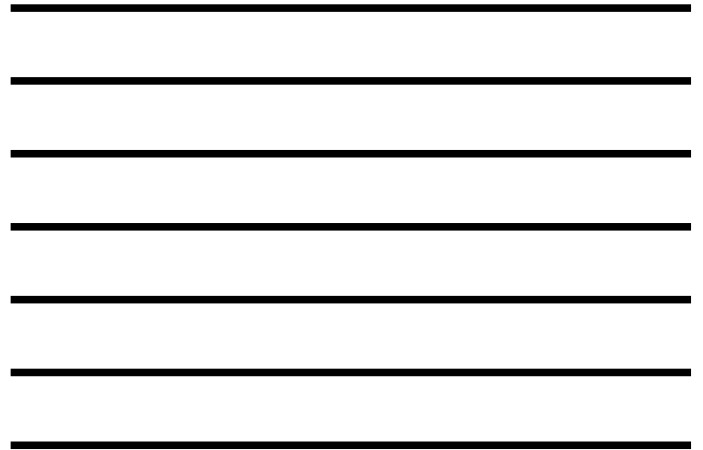
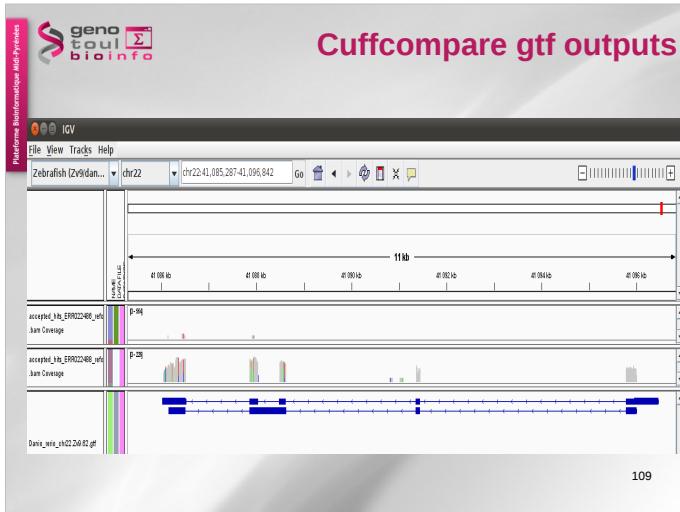
---

---

---

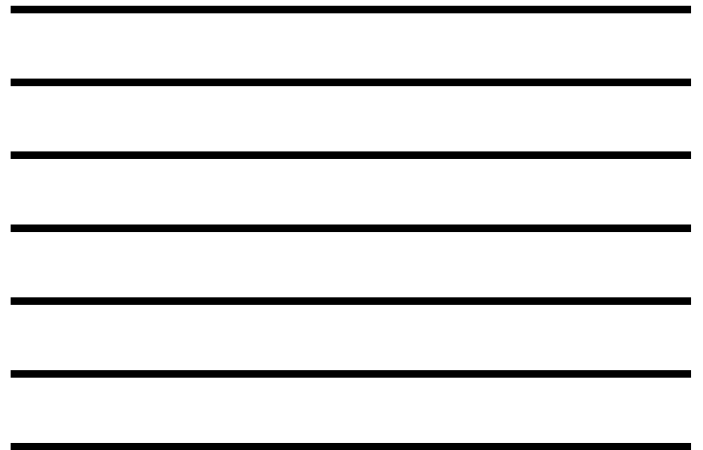
---





## Cuffdiff

- Find significant changes in transcript expression, splicing, and promoter use.
- tracks changes in the relative abundance of transcripts sharing a common transcription start site: see changes in splicing
- tracks changes in the relative abundances of the primary transcripts of each gene: see changes in relative promoter use
- Input :
  - gtf file of transcripts (output of cufflinks or cuffcompare combined.gtf)
  - 2 or more sam alignment files (but analyse 2 by 2)
- Output:
  - **.fpkm\_tracking** : calculates the FPKM of each transcript, primary transcript, and gene in each sample.
  - **exp.diff**: differential expression tests for each pair of samples at gene, isoform... level



## CummeRbund

<http://compbio.mit.edu/cummeRbund/>

- R package that is designed to aid and simplify the task of analyzing Cudiff outputs: explore and visualize the data
- Now included as part of R/Bioconductor
- Takes the output files from cuffdiff and creates a SQLite database of the results describing appropriate relationships between genes, transcripts, transcription start sites, and CDS regions
- numerous plotting functions as well for commonly used visualizations

NATURE PROTOCOLS | PROTOCOL

Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks

Cole Trapnell, Adam Roberts, Loyal Goff, Geo Pertea, Daehwan Kim, David R Kelley, Harold Pimental, Steven L Salzberg, John L Rinn & Lior Pachter

Affiliations | Contributions | Corresponding author

111

Nature Protocols 7, 562–578 (2012) | doi:10.1038/nprot.2012.016  
Published online 01 March 2012



Platforme Bioinformatique M&Pyrénées

geno toul bioinfo

## Sigcufflinks

- Cufflinks code has been modified by the Bioinformatic Platform of Toulouse in order to obtain raw count of reads: use **sigcufflinks** on **snp**
- Run cufflinks, cufflinks outputs + raw\_transcripts.tsv:
 

gene_id	transcript_id	paires	forward	reverse
CUFF.6	CUFF.6.1	4873	4873	3431
CUFF.6	CUFF.6.2	5222	5222	3769
CUFF.6	ENSDART00000067635	4819	4819	3580
- Other scripts are in development : exemple **cuffcompare\_quant\_summary.pl** on **snp**

112

---

---

---

---

---

---

---

---

---

---

---

---

Platforme Bioinformatique M&Pyrénées

geno toul bioinfo

## Cuffcompare\_quant\_summary

- Other scripts are in development : exemple **cuffcompare\_quant\_summary.pl** on **snp**
  - Parse cuffcompare tracking file and produce summary file with quantification for each transcript in each condition.
  - Usage:
 

```
quant_summary.pl [-p <prefix>] [-o <outdir>] [-r transcriptome.gtf] --tracking cuffcmp.tracking lib1name_transcripts.gtf [lib2name_transcripts.gtf ...]
```

    - The gtf files are output from cufflinks.
    - This script also rename the gene and transcript id for each input cufflinks gtf in order to get shared trans and genes id.

113

---

---

---

---

---

---

---

---

---

---

---

---

Platforme Bioinformatique M&Pyrénées

geno toul bioinfo

## Scripture


NATURE BIOTECHNOLOGY | RESEARCH | ARTICLE

**Ab initio reconstruction of cell type-specific transcriptomes in mouse reveals the conserved multi-exonic structure of lincRNAs**

Mitchell Gutman, Manuel Garber, Joshua Z Levin, Julie Donaghey, James Robinson, Xian Adiconis, Lin Fan, Magdalena J Kozlowski, Andreas Grötker, Chad Nusbaum, John L Rinn, Eric S Lander & Aviv Regev

Affiliations | Contributions | Corresponding authors

Nature Biotechnology 28, 503–510 (2010) | doi:10.1038/nbt.1633



- Method to reconstruct the transcriptome using only RNA-seq data and genome sequence (unannotated)
- 6 steps:
  - Use reads uniquely aligned to the genome
  - From the aligned spliced reads: construct a connectivity graph
  - Using all data: segmentation approach to identify significant paths
  - Construct a transcript graph
  - Generate a catalogue of transcripts

114

---

---

---

---

---

---

---

---

---

---

---

---

genotoul bioinfo

## HTSeq-count

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

- Process the output from short read aligners in various formats
- Count how many reads map to each feature (in RNA-Seq, the features are typically genes)
  - counting reads by genes
  - or consider each exon as a feature to check for alternative splicing
- Inputs:
  - file with aligned sequencing reads: bam (or sam) file
  - list of genomic feature; gtf file

115



genotoul bioinfo

## HTSeq-count

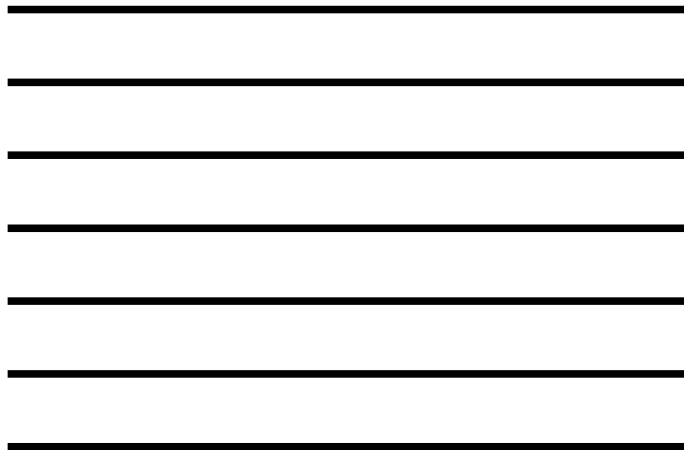
- Command line :
  - `htseq-count [options] <sam_file> <gtf_file>`
  - `samtools view accepted_hits.bam | htseq-count --stranded=no -m intersection-nonempty - file.gtf -q > output.htseq-count.txt &`

Some options:

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

- m <mode> : intersection-strict or intersection-nonempty (default union)
- stranded =<yes, no, or reverse> (default yes)
- t <feature type> : 3rd column in GTF file
- q : quiet
- h : help

116




genotoul bioinfo

## HTSeq-count

- Output: a table with counts for each feature and a summary of reads not counted for any feature:
  - *no\_feature*: reads which couldn't be assigned to any feature
  - *ambiguous*: reads which could have been assigned to more than one feature and hence were not counted for any of these
  - *not\_aligned*: reads in the SAM file without alignment
  - *alignment\_not\_unique*: reads with more than one reported alignment. These reads are recognized from the NH optional SAM field tag. (If the aligner does not set this field, multiply aligned reads will be counted multiple times.)

117




**Exemples et TP : cufflinks**

Commands guide:

Cufflinks command:

```
cufflinks --output-dir=cufflinks_ERR022486
-g Danio_erio_chr22.Zv9.62.gtf
accepted_hits_ERR022486.bam
```

Cuffcompare command:

```
cuffcompare -o cuffcompare_refVSERR022486
-r Danio_erio_chr22.Zv9.62.gtf
cufflinks_ERR022486/transcripts_ERR022486.gtf
```

Cuffdiff command:

```
cuffdiff Danio_erio_chr22.Zv9.62.gtf
tophat/accepted_hits_ERR022486.bam
tophat/accepted_hits_ERR022488.bam
```

118

---

---

---

---

---

---

---


---

---

---

---

---


**Exemples et TP : cufflinks**

Examples of used commands :

Nb of genes in GTF:

```
cat Danio_erio_chr22.Zv9.62.gtf | cut -f9 | cut -d";" -f1 | sort -u | wc -l
```

Nb of transcripts in GTF:

```
cat Danio_erio_chr22.Zv9.62.gtf | cut -f9 | cut -d";" -f2 | sort -u | wc -l
```

Cufflinks:

```
qsub -b y cufflinks --output-dir=/work/.../cufflinks_ERR022486 -g
/work/.../Danio_erio_chr22.Zv9.62.gtf
/work/.../tophat_ERR022486/accepted_hits_ERR022486.bam
```

Nb of genes in cufflinks results:

```
cat genes_ERR022486.fpk_tracking | cut -f1 | sort -u | wc -l
```

119

---

---

---

---

---

---

---


---

---

---

---

---


**Exemples et TP : cufflinks**

Examples of used commands :

Cuffcompare:

```
cuffcompare -o cuffcompare_refVSERR022486 -r
/work/.../Danio_erio_chr22.Zv9.62.gtf
/work/.../cufflinks_ERR022486_refchr22/transcripts_ERR022486_refchr22.gtf
```

Cuffdiff:

```
qsub -b y cuffdiff /work/.../Danio_erio_chr22.Zv9.62.gtf
/work/.../tophat_ERR022486/accepted_hits_ERR022486.bam
/work/.../tophat_ERR022488/accepted_hits_ERR022488.bam
```

Significant differential expression of genes:

```
cat gene_exp.diff | sort -r -k 14 | less -S
```

120

---

---

---

---

---

---

---

---

---

---

---

---

## Exemples et TP : cufflinks

Examples of used commands :

Sigcuffdiff:

```
sigcufflinks --output-dir=/work/.../sigcufflinks_ERR022486 -g  
/.../Danio_rerio_chr22.Zv9.62.gtf  
/.../tophat_ERR022486/accepted_hits_ERR022486.bam
```

Cuffcompare\_quant\_summary:

```
cuffcompare_quant_summary.pl -p quant_86_88 -o  
/work/.../cuffcompare_quantSummary -r /.../Danio_rerio_chr22.Zv9.62.gtf --tracking  
/.../cuffcompare_sigcufflinks/cuffcompare_86_88.tracking  
/.../sigcufflinks_ERR022486/ERR022486_transcripts.gtf  
/.../sigcufflinks_ERR022488/ERR022488_transcripts.gtf
```

---

---

---

---

---

---

---

---

---

---

## Quality for Bioinfo Platform!

Exam :

<http://bioinfo.genotoul.fr/index.php?id=93>

Satisfaction form :

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

---

---

---

---

---

---

---

---

---

---

## Useful links

Seqanswers: <http://seqanswers.com/>

RNAseq blog: <http://rna-seqblog.com/>

Illumina: <http://www.illumina.com/>

---

---

---

---

---

---

---

---

---

---