

Plateforme Bioinformatique Midi-Pyrénées




RNA-Seq data analysis

http://genoweb.toulouse.inra.fr/~formation/4_Galaxy_RNAseq/

Delphine Labourette Get-Biopuce / Céline Noirot Bioinfo Genotoul

1

Plateforme Bioinformatique Midi-Pyrénées



Material

http://genoweb.toulouse.inra.fr/~formation/4_Galaxy_RNAseq/

Slides & Exercise leaflet (doc)

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

Data & results files (data)

2

Plateforme Bioinformatique Midi-Pyrénées



Session organisation

- Sequence quality
 - Theory + exercises
- Spliced read mapping
Visualisation
 - Theory + exercises
- expression measurement
 - Theory + exercises
- mRNA calling
 - Theory + exercises

3

Platforme Bioinformatique Midi-Pyrénées

geno toul bioinfo

What you should know

How to connect to Sigenar galaxy workbench?
<http://sigenae-workbench.toulouse.inra.fr/galaxy/>

Platforme Bioinformatique Midi-Pyrénées

geno toul bioinfo

Transcription

Transcription is the process of creating a complementary RNA copy of a sequence of DNA. Transcription is the first step leading to gene expression.

[http://en.wikipedia.org/wiki/Transcription_\(genetics\)](http://en.wikipedia.org/wiki/Transcription_(genetics))

Platforme Bioinformatique Midi-Pyrénées

geno toul bioinfo

Transcription products

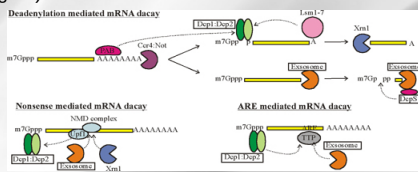
Protein coding gene: transcribed in mRNA
ncRNA : highly abundant and functionally important RNA

- tRNA,
- rRNA,
- snoRNAs,
- microRNAs,
- siRNAs,
- PiRNAs
- lincRNA

http://en.wikipedia.org/wiki/User:Amarchais/RsaOG_RNA

Transcript degradation

After export to the cytoplasm, mRNA is protected from degradation by a 5' cap structure and a 3' poly adenine tail. In the deadenylation dependent mRNA decay pathway, the polyA tail is gradually shortened by exonucleases. This ultimately attracts the degradation machinery that rapidly degrades the mRNA in both in the 5' to 3' direction and in the 3' to 5' direction. Additional mechanisms, including the nonsense mediated decay pathway, bypass the need for deadenylation and can remove the mRNA from the transcriptional pool independently. Interestingly, the same enzymes are responsible for the actual degradation of the mRNA independent of the pathway taken (see figure).



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



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

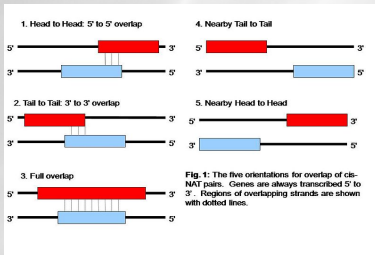


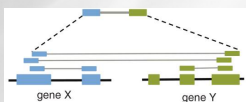
Fig. 1: The five orientations for overlap of cis-NAT pairs. Genes are always transcribed 5' to 3'. Regions of overlapping strands are shown with dotted lines.



Fusion genes

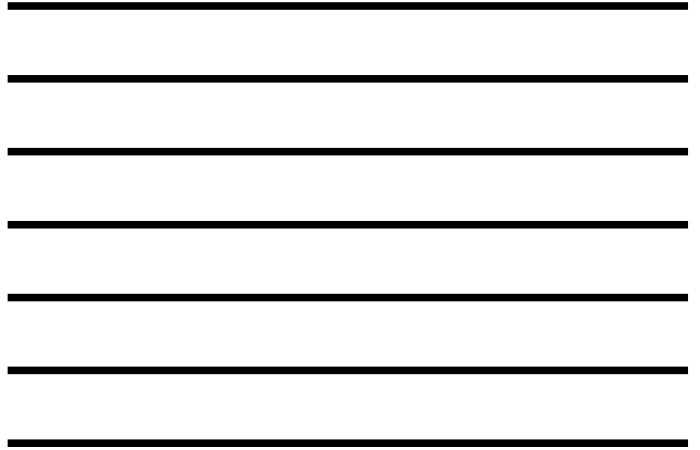
- A fusion gene is a hybrid gene formed from two previously separate genes. It can occur as the result of a translocation, interstitial deletion, or chromosomal inversion. Often, fusion genes are oncogenes.
- They often come from trans-splicing: Trans-splicing is a special form of RNA processing in eukaryotes 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, Murumali A, Kanoaspeetka S, Nicotri D, Honigslav V, Kleiv K, Rie IH, Nuber S, Wolf M, Borresen-Dale AL, Kallioniemi O, Institute for Molecular Medicine Finland (FIMM), Tukohamkatu 8, Helsinki, 00290, Finland, olli.kallioniemi@fimm.fi



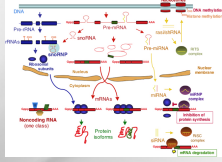
http://en.wikipedia.org/wiki/Fusion_gene

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



Transcriptome variability summary

- 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



13

http://www.nature.com/embj/journal/v25/n5/fig_tab/7601023a_F2.html



UTR length

Lengthening of 3'UTR Increases with Morphological Complexity in Animal Evolution

Cho-Yi Chen^{1,2}, Shih-Tein Chen², Hsueh-Fen Juan^{1,2} and Hsuan-Cheng Huang^{1,2*}
¹Genome and Systems Biology Degree Program, Department of Life Science, Institute of Molecular and Cellular Biology, Graduate Institute of Biomedical Electronics and Bioinformatics, National Taiwan University, Taipei, Taiwan.
²Institute of Biological Chemistry, Academia Sinica, Taipei, Taiwan
³Institute of Biochemical Informatics, Center for System and Synthetic Biology, National Yang-Ming University, Taipei, Taiwan

Associate Editor: Martin Bishop

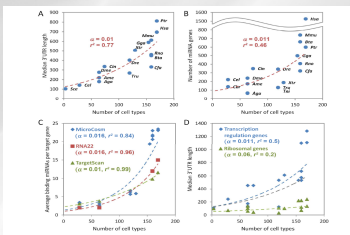
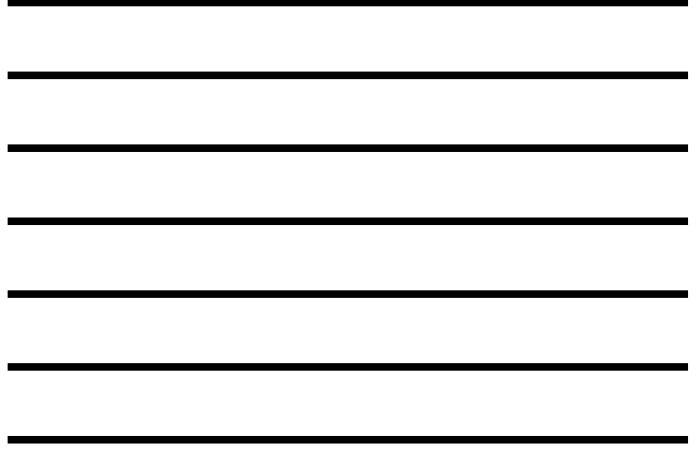


Fig. 1. Empirical correlation between cell type complexity and transcriptome complexity. (A) A strong empirical correlation exists between the median length of 3'UTR and the total number of cell types in vertebrate species, as measured by the number of cell types (1,000 and 2,000). (B) A similar correlation exists between the median length of 5'UTR and the total number of cell types in vertebrate species. The median length of 5'UTR is plotted against the number of cell types. (C) The average length of 3'UTR per gene increases with the number of cell types. (D) The median length of 3'UTR per gene increases with the number of cell types. The data are shown in the main text and Supplementary Figure S1 for species listed. The strength of correlation between 3'UTR length for transcription regulator genes (TRG), immune genes (IG), and housekeeping genes (HKG) and the number of cell types is indicated by the color of the points. The points are color-coded as follows: TRG (red), IG (green), and HKG (blue).

14

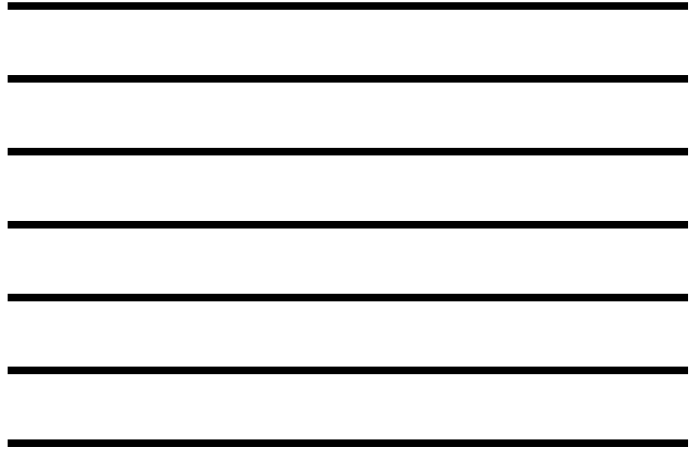


Techniques classification ?

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 novel splice events

15



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

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 Wang¹, Mark Gerstein² & Michael Snyder¹ [About the authors](#) top ¹

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>

16

genotoul bioinfo

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, ...

17

genotoul bioinfo

RNA-Seq platforms comparison

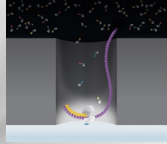
Séquenceurs 2^{ème} génération

Société	Roche				Illumina			Life Technologies				
Plateforme	GS Junior	454	MiSeq	HiSeq 1000	HiSeq 2000	Genome Analyzer IIx	Ion Torrent PGM	SOLID 4	SOLID 5500	SOLID 6600		
Technologie	Titanium	FLX Titanium	FLX+				Chip 314	Chip 316	Chip 318			
Acides nucléiques (matrice)												
Ligation adaptateurs												
Méthode d'amplification	PCR en émulsion		"Bridge PCR"				PCR en émulsion					
Méthode de séquençage	Synthèse (Pyroséquencage)		Synthèse				Ligation					
Durée de séquençage/run	10h	10h-20h	26h	8jrs	8jrs	14jrs	2h	12jrs	8jrs	8jrs		
Capacité (Mb) séquençage/run	50	500-900	1500	100000	200000	950000	>10	>100	>1000	70000	80000	150000
Taille moyenne des reads	400	400-700	150+150	100+100	100+100	150+150	100	>100	>100	50+35	75+35	75+35
Coût (\$) /run	1100	6200	750	10000	20000	11500	500	750	950	8150	6100	10500
Coût machine + annexes (K\$)	110+25	500+30	125	360	690	250	50+20			480+55	350+55	600+55
Exactitude de séquençage (%)	99	99	99,9	99,9	99,9	99,9	99	99,95	99,95	99,95	99,99	99,99

18

Third Generation RNA-Seq

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

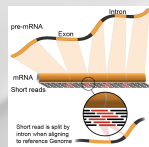


<http://www.genengnews.com/gen-articles/third-generation-sequencing-debuts/3257/>

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



What are we looking for?

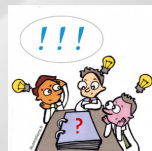
Identify genes

- List new genes

Identify transcripts

- List new alternative splice forms

Quantify these elements → differential expression



genotoul bioinfo

Plateforme Bioinformatique M&P Pyrenees

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 ?

22


genotoul bioinfo

Plateforme Bioinformatique M&P Pyrenees

ENCODE answers

- RNA-Seq is not a mature technology.
- 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 between **0.92 to 0.98**. Experiments with biological correlations that fall below 0.9 should be either be repeated or explained.
- Between **30M and 100M reads** per sample depending on the study.

NB. Guidelines for the information to publish with the data.



Encyclopedia of DNA Elements

<http://encodeproject.org/ENCODE/dataStandards.html>

23

genotoul bioinfo

Plateforme Bioinformatique M&P Pyrenees

Statisticians answers

Efficient experimental design and analysis strategies for the detection of differential expression using RNA-Sequencing

BMC Genomics 2012, **13**:484 doi:10.1186/1471-2164-13-484

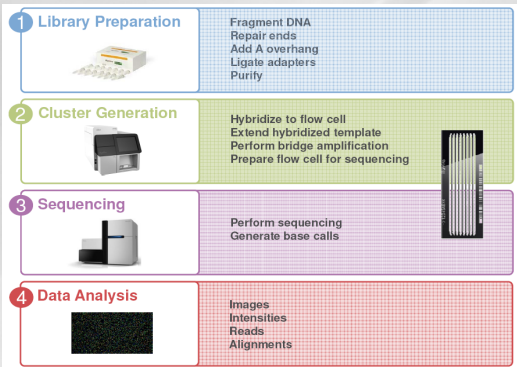
Jose A Robles (jose.robles@csiro.au)

Conclusions

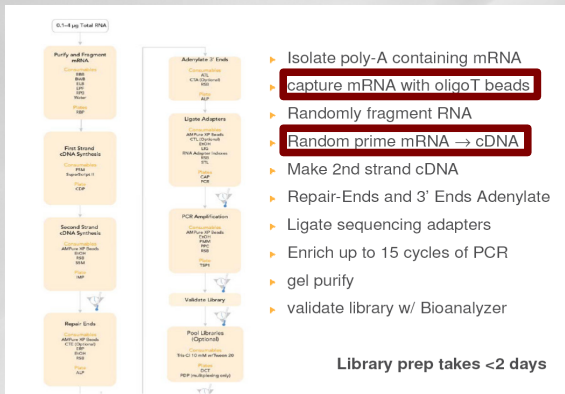
This work quantitatively explores comparisons between contemporary analysis tools and experimental design choices for the detection of differential expression using RNA-Seq. We found that the DESeq algorithm performs more conservatively than edgeR and NBPSeq. With regard to testing of various experimental designs, this work strongly suggests that greater power is gained through the use of biological replicates relative to library (technical) replicates and sequencing depth. Strikingly, sequencing depth could be reduced as low as 15% without substantial impacts on false positive or true positive rates.

24

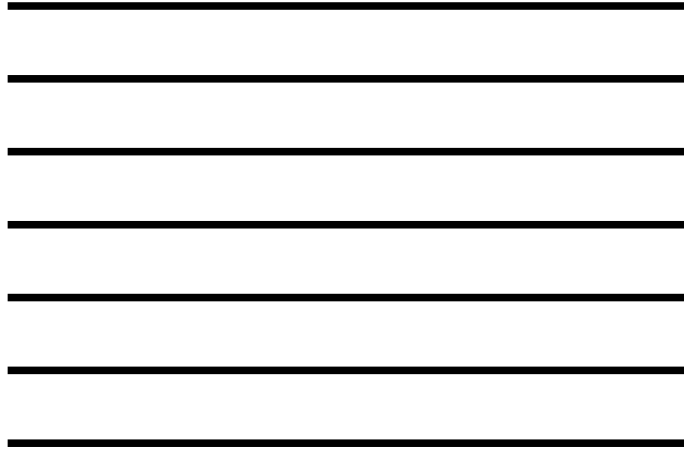
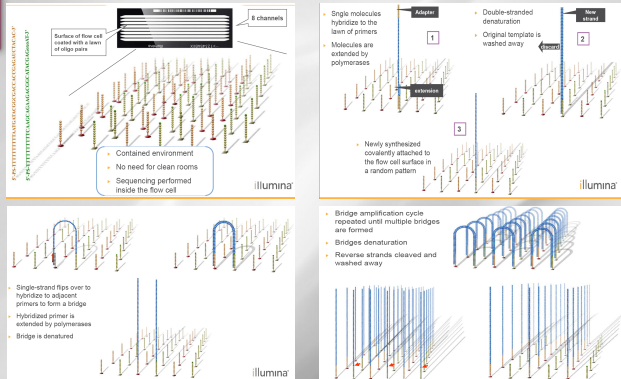
Illumina RNA-Seq protocol



RNA-Seq library preparation



Clusters generation



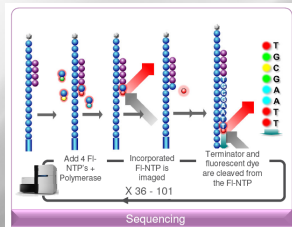
Sequencing

Cluster generation:

- 35 amplification cycles
- 1 cluster → 2 000 identical molecules
- 500 000 clusters / floccell

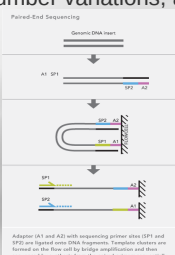
Sequencing:

- Image acquisition:
 - 50 min / cycle
- Ex: 2x100bp → 2x100x50 min



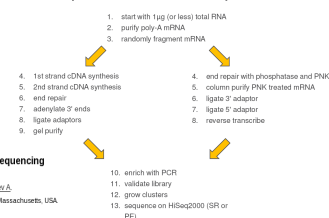
Paired-end sequencing

- 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



Strand specific RNA-Seq protocol

workflow comparison: mRNA-Seq vs directional mRNA-Seq

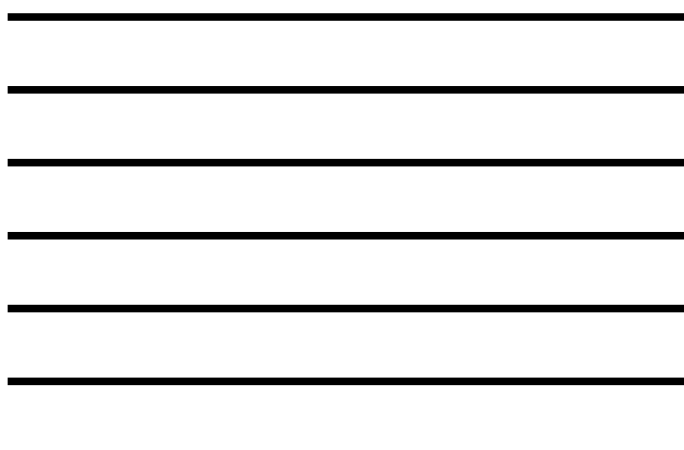
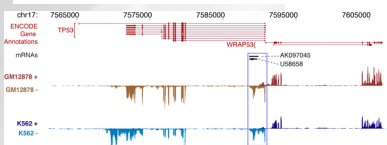


Nat Methods. 2008 Sep 7(9): 709-15. Epub 2008 Aug 15.


Comprehensive comparative analysis of strand-specific RNA sequencing methods.

Levin JZ, Yassour M, Adorciis X, Nusbaum C, Thompson DA, Friedman N, Grinke A, Regev A

Broad Institute of Massachusetts Institute of Technology and Harvard University, Cambridge, Massachusetts, USA
<http://dx.doi.org/10.1038/nmeth1244>



Platforme Bioinformatique M&Pyrénées



Analysis workflow

Data quality control


Spliced mapping

Quantification

Gene and transcript discovery

31

Platforme Bioinformatique M&Pyrénées



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. Cock¹*, Christopher J. Fields², Naohisa Goto³, Michael L. Heuer⁴ and Peter M. Rice⁵

```
@EAS54_6_R1_2_1_413_324
CCCTTCTGTCTTCAGCGTTCTCC
+
!!3!!!!!!!!!!!!!!7!!!!!!88
```

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

32

Platforme Bioinformatique M&Pyrénées



RNAseq specific bias

- Influence of the library preparation
- Random hexamer priming
- Positional bias and sequence specificity bias.
Robert et al. *Genome Biology*, 2011,12:R22
- Transcript length bias
- Some reads map to multiple locations

33

Platforme Bioinformatique M&P Pythons

geno
toul
bioinfo

Hexamer random priming bias

Published online 14 April 2010
Nucleic Acids Research, 2010, Vol. 38, No. 12, e139
doi:10.1093/nar/gkq124

Biases in Illumina transcriptome sequencing caused by random hexamer priming
Kasper D. Hansen¹*, Steven E. Brenner² and Sandrine Dudot^{1,3}

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

- There is a strong distinctive pattern in the nucleotide frequencies of the first 13 positions at the 5'-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

34



Platforme Bioinformatique M&P Pythons

geno
toul
bioinfo

Hexamer random effect

- Orange = reads start sites
- Blue = coverage

35



Platforme Bioinformatique M&P Pythons

geno
toul
bioinfo

Transcript length bias

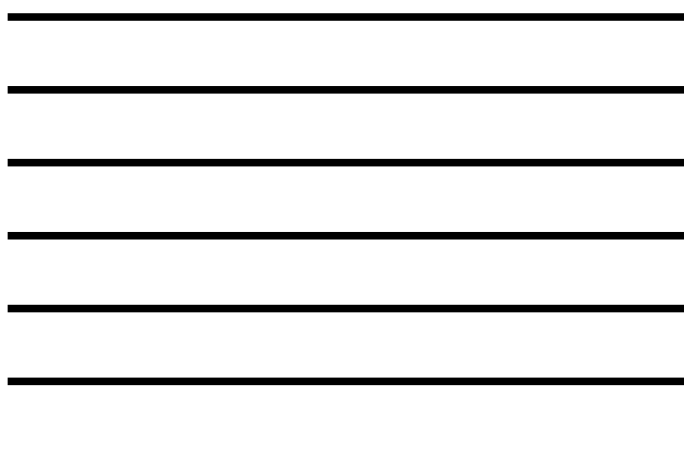
Biol Direct, 2009 Apr 18:4:14
Transcript length bias in RNA-seq data confounds systems biology.
Oshlack A, Wakefield MJ

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 length bias on published data sets. For standard analyses using a 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 Clouston (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^{1,†}, Zhide Fang^{2,†}, Kai Zheng¹, Degui Zhi¹ and Xiangqin Cui^{1,*}

36



Platforme Bioinformatique Multi-Professeurs

geno
toul
bioinfo

Verifying RNA-Seq raw data

FastQC :

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

- Has been developed for genomic data

37

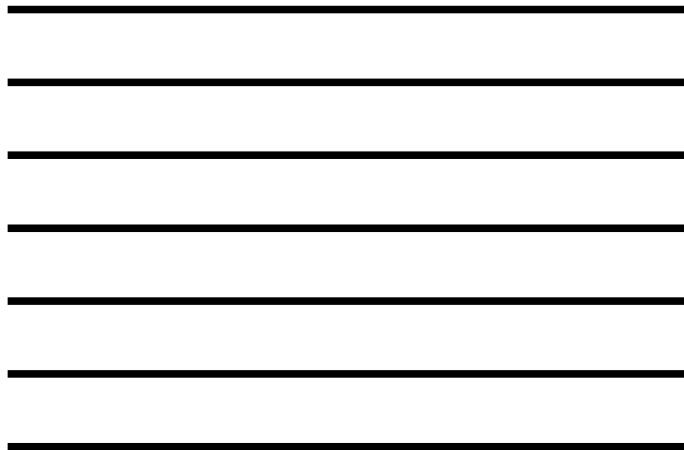


Platforme Bioinformatique Multi-Professeurs

geno
toul
bioinfo

Fourmiz !!!

38



Platforme Bioinformatique Multi-Professeurs

geno
toul
bioinfo

Take home message on quality analysis

Elements to be checked :

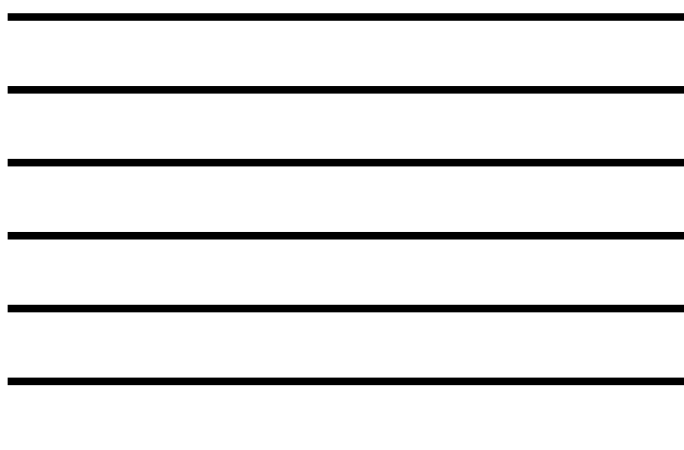
- Random priming effect
- K-mer (polyA, polyT)

Alignment on reference for the second quality check and filtering.


A good run?:

- Expected number of reads produced (2x500millions / flowcell),
- Length of the reads expected (100pb),
- Random selection of the nucleotides and the GC%,
- Good alignment: very few unmapped reads, pairs mapped on opposite strands.

39



Platforme Bioinformatique MSU-Pyrenees



Analysis workflow

Data quality control


Spliced mapping

Quantification

Gene and transcript discovery

40

Platforme Bioinformatique MSU-Pyrenees




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>

41

Platforme Bioinformatique MSU-Pyrenees



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

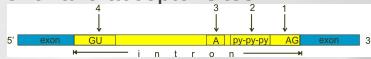
42

Platforme Bioinformatique M&Pyrénées

geno tout bioinfo

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

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

Nucleic Acids Res. 2009 Nov 1;37(21):4364-75.
Analysis of canonical and non-canonical splice sites in mammalian genomes.
Burset M, Saitedzhov JA, Solovnev YV

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

43



Platforme Bioinformatique M&Pyrénées

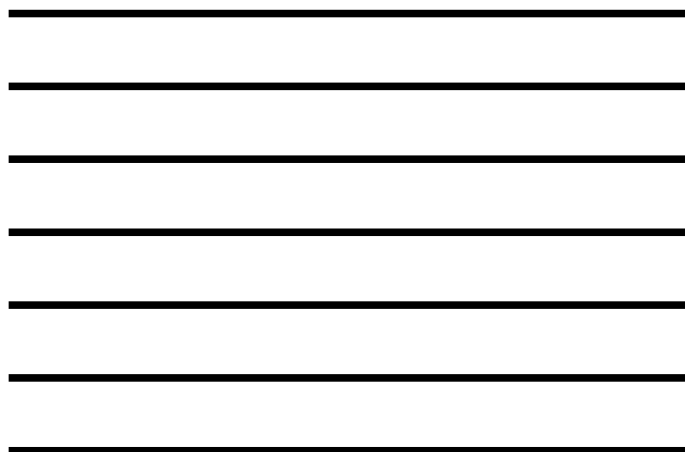
geno tout bioinfo

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 (seed then mapping)

Genome Res. 1998 Sep 8(9):967-74.
A computer program for aligning a cDNA sequence with a genomic DNA sequence.
Flora L, Hanchell G, Zhang Z, Rubin GM, Miller W
Department of Computer Science and Engineering, The Pennsylvania State University, University Park, Pennsylvania 16802 USA.

44



Platforme Bioinformatique M&Pyrénées

geno tout bioinfo

TopHat

BIOINFORMATICS ORIGINAL PAPER Vol. 25 no. 9 2009, pages 1105-1113 doi:10.1093/bioinformatics/btp130

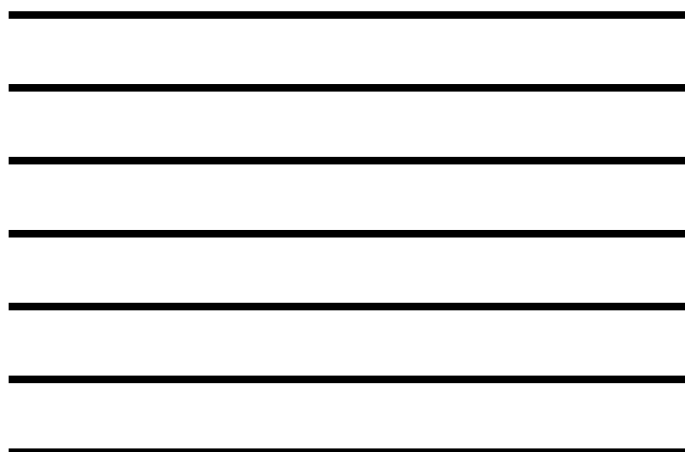
Sequence analysis

TopHat: discovering splice junctions with RNA-Seq
Cole Trapnell^{1,*}, Lior Pachter² and Steven L. Salzberg¹

<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.

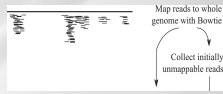
45



TopHat algorithm : first step

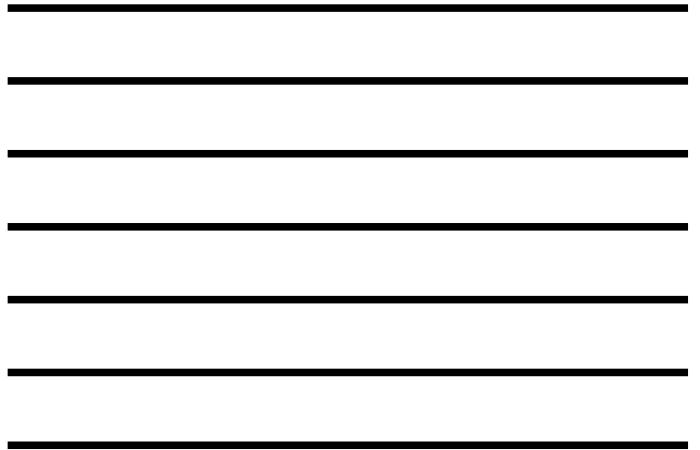
- TopHat finds junctions by mapping reads to the reference:

- 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

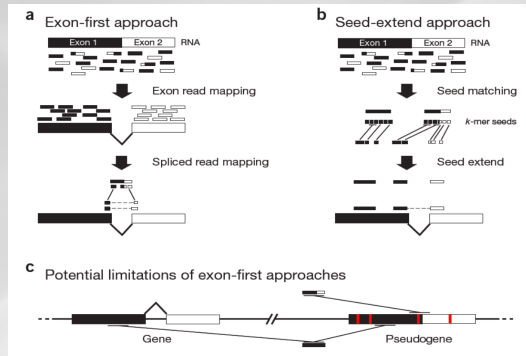


46

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



Exon first approach limitation

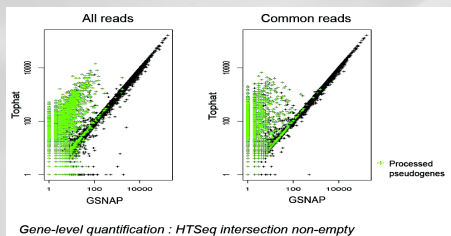


REVIEW
Computational methods for transcriptome annotation and quantification using RNA-seq
Manuel Garber¹, Manfred G Grabher¹, Mitchell Gutman^{2,3} & Cole Trapnell^{1,2}

47



TopHat and pseudogenes



--read-realign-edit-dist
Some of the reads spanning multiple exons may be mapped incorrectly as a contiguous alignment to the genome even though the correct alignment should be a spliced one - this can happen in the presence of processed pseudogenes that are rarely (if at all) transcribed or expressed. This option can direct TopHat to re-align reads for which the edit distance of an alignment obtained in a previous mapping step is above or equal to this option value. If you set this option to 0, TopHat will map every read in all the mapping steps (transcriptome if you provided gene annotations, genome, and finally splice variants detected by TopHat), reporting the best possible alignment found in any of these mapping steps. This may greatly increase the mapping accuracy at the expense of an increase in running time. The default value for this option is set such that TopHat will not try to realign reads already mapped in earlier steps.

48



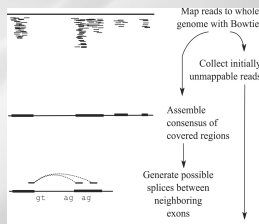
Why does it find small exons?

- In the last tophat versions :

Short read sequencing machines can currently produce reads 100bp or longer but many exons are shorter than this so they would be missed in the initial mapping. **TopHat solves this problem mainly by splitting all input reads into smaller segments which are then mapped independently.** The segment alignments are put back together in a final step of the program to produce the end-to-end read alignments.

Exon assembly process

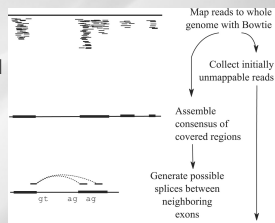
- TopHat then assembles the mapped reads
- Define island: aggregates mapped reads in islands of candidate exons
 - Generate potential donor/acceptor splice sites using neighbouring exons
- Extend islands to cover eventually splice junctions
 - +/- 45 bp from reference on either side of island



Splice junction reference

To map reads to splice junction :

- Enumerate all canonical donor and acceptor sites in islands
 - long (≥ 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
 - $50 \text{ bp} < \text{Intron size} < 500,000 \text{ bp}$



genotoul bioinfo

IUM alignment

- Each possible intron is checked against the IUM

→ seed and extend alignment

left exon gt ag right exon

IUM read

high quality

Map reads to whole genome with Bowtie

Collect initially unmappable reads

Assemble consensus of covered regions

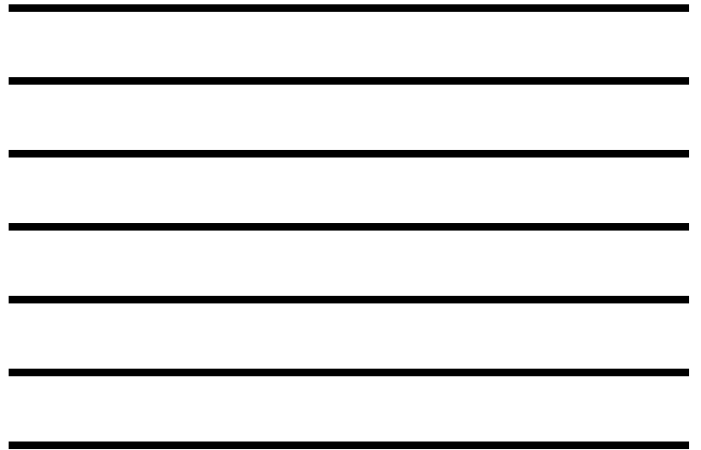
Generate possible splices between neighboring exons

Build seed table index from unmappable read

Map reads to possible splices via seed-and-extend

52

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



genotoul bioinfo

TopHat Inputs

Inputs :

- bowtie2 index of the genome
ftp://ftp.cbcb.umd.edu/pub/data/bowtie_indexes/
<http://bowtie-bio.sourceforge.net/index.shtml>
- file fasta (.fa) of the reference or will be build by bowtie, in the index directory
- File fastq of the reads

! the GTF file and the Bowtie index should have same name of chromosome or contig

Command lines :

```
bowtie2-build <reference.fasta> <index_base>
tophat [options] <index_base> <reads1_1> <reads1_2>
```

53



genotoul bioinfo

TopHat Options

Some useful options (command line) :

- h/--help
- v/--version
- - **bowtie1** (instead of bowtie2)
- o/--output-dir
- r/--mate-inner-dist : no default value
- m/--splice-mismatches : default 0
- i/--min-intron-length : default 50
- l/--max-intron-length : default 500000, prefer 25000 for non human
- max-insertion-length : default 3
- max-deletion-length : default 3
- p/--num-threads

54



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
r-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.
r-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.
r-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.

Platforme Bioinformatique M&Pyrénées

geno toul bioinfo

Bed example

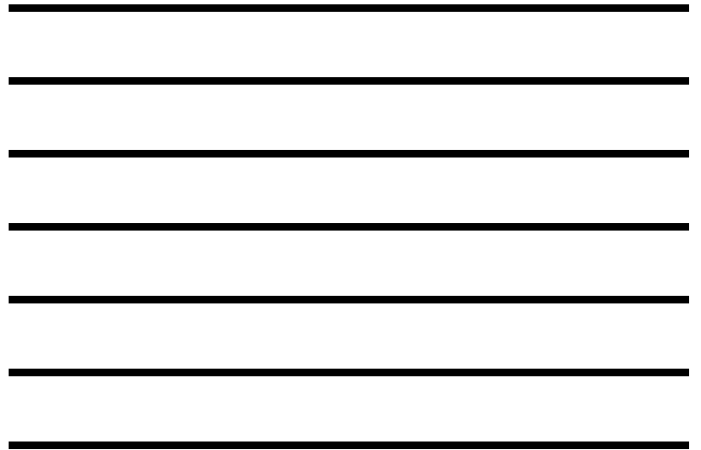
Chrom Start End name score strand drawing RGB Blocks info

```

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 1785 4260 JUNC000000002 1 - 1785 4260 255,0,0,2 26,48 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 12398 21452 JUNC000000012 27 - 12398 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



Platforme Bioinformatique M&Pyrénées

geno toul 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



Platforme Bioinformatique M&Pyrénées

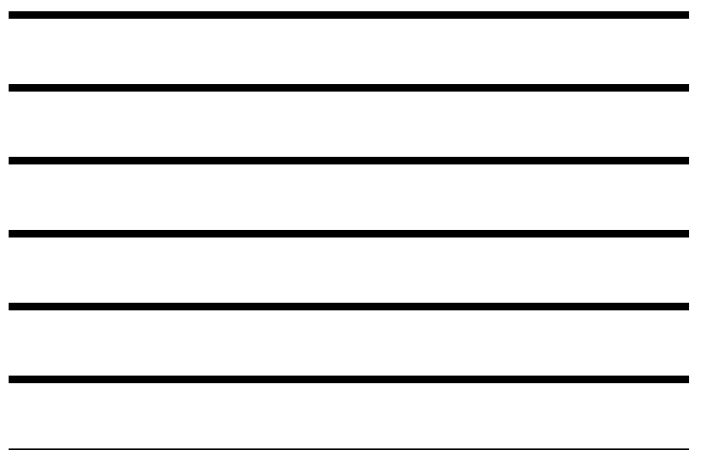
geno toul bioinfo

Trans-splicing with 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

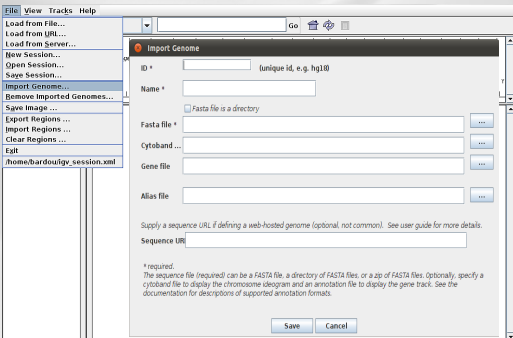
63



genotoul bioinfo

Visualizing alignments on IGV

Import a reference genome



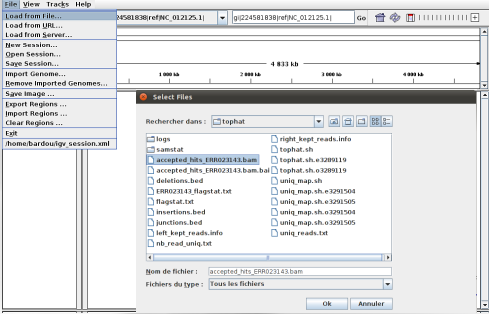
67



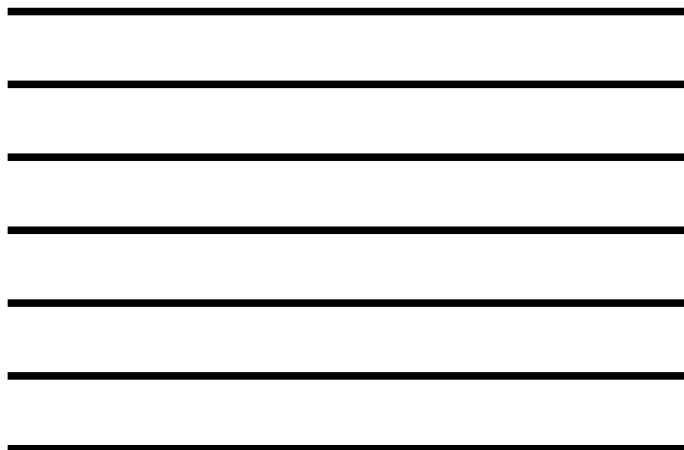
genotoul bioinfo

Visualizing alignments on IGV

– Import your BAM Files



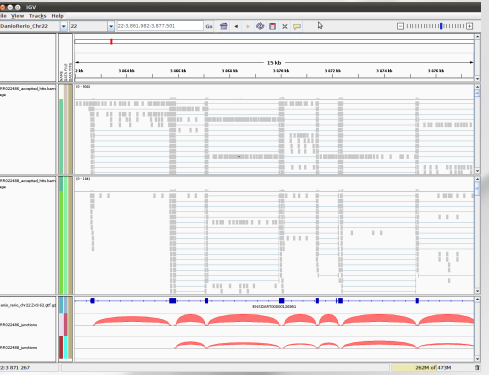
68



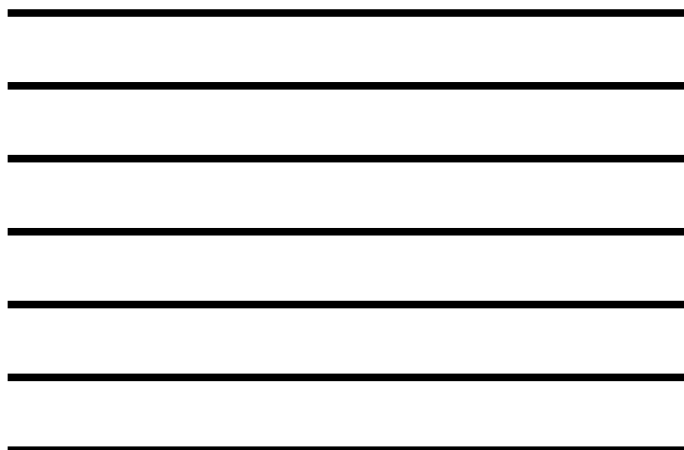
genotoul bioinfo

Visualizing alignments on IGV

– Exemple of bam and bed files visualisation



69



genotoul bioinfo

hands-on : tophat

Tophat location: 8 – Trainings

- RNA-Seq
 - Step 2 : *Alignment and statistics*
 - * Tophat for Illumina Find splice junctions using RNA-seq data

Indexation: * Samtools index

Samtools flagstat

* Tophat for Illumina (version 1.0.0)

Your RNA-Seq FASTQ file (read 1):
3:ERR022488_chr22_read1

Your RNA-Seq FASTQ file (read 2):
4:ERR022488_chr22_read2

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

Number of threads used to align reads:
16

Maximum intron length:
5000

Expected (mean) inner distance between mate pairs:
200

Execute



genotoul bioinfo

Analysis workflow

Data quality control

Spliced mapping

Quantification

Gene and transcript discovery

71



genotoul bioinfo

What do we want to build?

The gene / transcript description file (and corresponding fasta)

```

> protein_coding exon 607700 607947 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "1"
> protein_coding exon 607948 608000 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "2"
> protein_coding exon 608001 608053 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "3"
> protein_coding exon 608054 608107 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "4"
> protein_coding exon 608108 608161 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "5"
> protein_coding exon 608162 608215 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "6"
> protein_coding exon 608216 608269 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "7"
> protein_coding exon 608270 608323 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "8"
> protein_coding exon 608324 608377 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "9"
> protein_coding exon 608378 608431 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "10"
> protein_coding exon 608432 608485 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "11"
> protein_coding exon 608486 608539 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "12"
> protein_coding exon 608540 608593 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "13"
> protein_coding exon 608594 608647 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "14"
> protein_coding exon 608648 608701 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "15"
> protein_coding exon 608702 608755 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "16"
> protein_coding exon 608756 608809 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "17"
> protein_coding exon 608810 608863 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "18"
> protein_coding exon 608864 608917 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "19"
> protein_coding exon 608918 608971 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "20"
> protein_coding exon 608972 609025 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "21"
> protein_coding exon 609026 609079 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "22"
> protein_coding exon 609080 609133 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "23"
> protein_coding exon 609134 609187 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "24"
> protein_coding exon 609188 609241 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "25"
> protein_coding exon 609242 609295 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "26"
> protein_coding exon 609296 609349 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "27"
> protein_coding exon 609350 609403 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "28"
> protein_coding exon 609404 609457 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "29"
> protein_coding exon 609458 609511 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "30"
> protein_coding exon 609512 609565 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "31"
> protein_coding exon 609566 609619 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "32"
> protein_coding exon 609620 609673 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "33"
> protein_coding exon 609674 609727 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "34"
> protein_coding exon 609728 609781 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "35"
> protein_coding exon 609782 609835 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "36"
> protein_coding exon 609836 609889 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "37"
> protein_coding exon 609890 609943 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "38"
> protein_coding exon 609944 609997 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "39"
> protein_coding exon 609998 700000 - - gene_id "ENSG00000007709" transcript_id "ENST00000004420" exon_number "40"
  
```

The count file

chr	start	end	count
1	100	100	100
1	101	101	101
1	102	102	102
1	103	103	103
1	104	104	104
1	105	105	105
1	106	106	106
1	107	107	107
1	108	108	108
1	109	109	109
1	110	110	110
1	111	111	111
1	112	112	112
1	113	113	113
1	114	114	114
1	115	115	115
1	116	116	116
1	117	117	117
1	118	118	118
1	119	119	119
1	120	120	120
1	121	121	121
1	122	122	122
1	123	123	123
1	124	124	124
1	125	125	125
1	126	126	126
1	127	127	127
1	128	128	128
1	129	129	129
1	130	130	130
1	131	131	131
1	132	132	132
1	133	133	133
1	134	134	134
1	135	135	135
1	136	136	136
1	137	137	137
1	138	138	138
1	139	139	139
1	140	140	140
1	141	141	141
1	142	142	142
1	143	143	143
1	144	144	144
1	145	145	145
1	146	146	146
1	147	147	147
1	148	148	148
1	149	149	149
1	150	150	150
1	151	151	151
1	152	152	152
1	153	153	153
1	154	154	154
1	155	155	155
1	156	156	156
1	157	157	157
1	158	158	158
1	159	159	159
1	160	160	160
1	161	161	161
1	162	162	162
1	163	163	163
1	164	164	164
1	165	165	165
1	166	166	166
1	167	167	167
1	168	168	168
1	169	169	169
1	170	170	170
1	171	171	171
1	172	172	172
1	173	173	173
1	174	174	174
1	175	175	175
1	176	176	176
1	177	177	177
1	178	178	178
1	179	179	179
1	180	180	180
1	181	181	181
1	182	182	182
1	183	183	183
1	184	184	184
1	185	185	185
1	186	186	186
1	187	187	187
1	188	188	188
1	189	189	189
1	190	190	190
1	191	191	191
1	192	192	192
1	193	193	193
1	194	194	194
1	195	195	195
1	196	196	196
1	197	197	197
1	198	198	198
1	199	199	199
1	200	200	200

72



If you have the model file

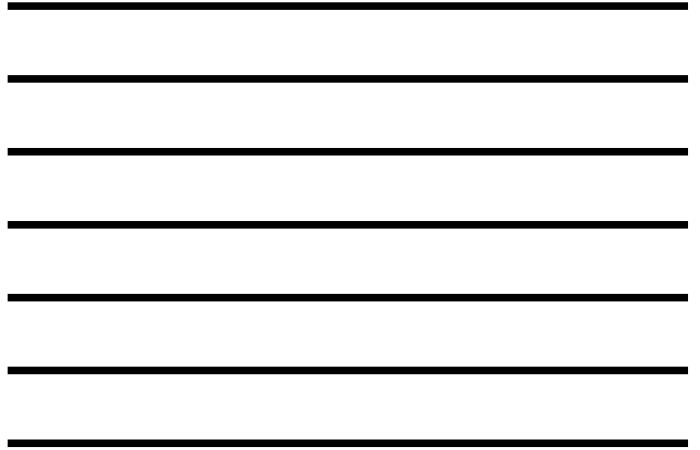
The model is presented in the GTF file (Gene Transfer Format)

Two approaches

- Gene level
- Transcript level

Tools for each approach

- htseq-count
- cufflinks (sigcufflinks)



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



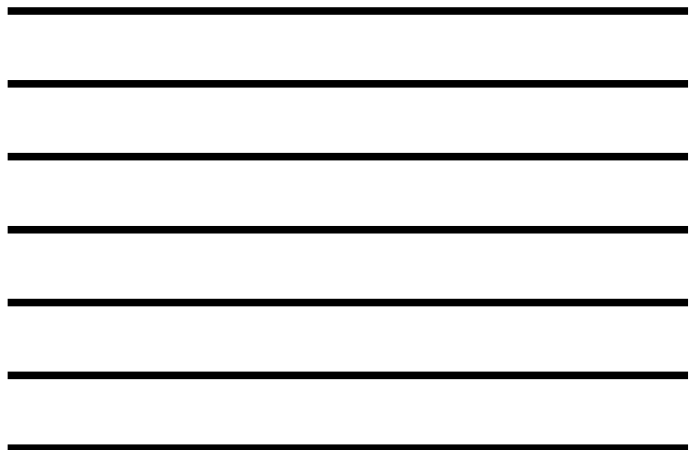
HTSeq-count

- 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:

- 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



Platforme Bioinformatique M&P/Pythons

genotoul bioinfo

HTSeq-count

- Output: a table with counts for each feature and a summary of reads not counted for any feature:

ENSDARG000000095643	967
ENSDARG000000095659	4
ENSDARG000000095667	36
ENSDARG000000095677	98
ENSDARG000000095708	5
no_feature	362748
ambiguous	9937
too_low_aQual	0
not_aligned	0
alignment not unique	239465

- *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.)

76



Platforme Bioinformatique M&P/Pythons

genotoul bioinfo

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, Marjke J van Baren, Steven L Salzberg, Barbara J Wold & Lior Pachter

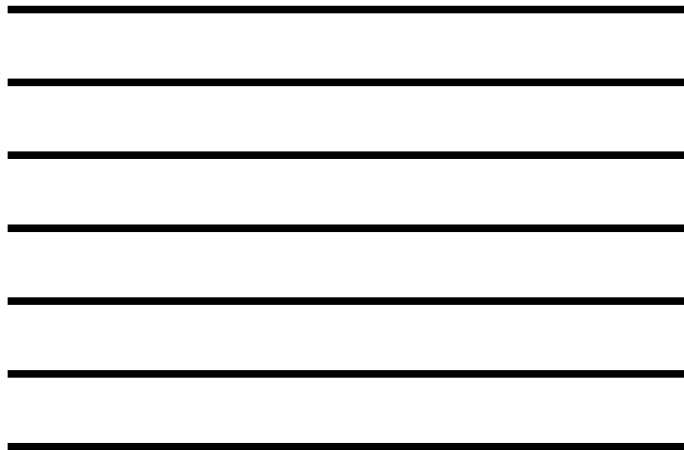
Affiliations | Contributions | Corresponding author

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

77



Platforme Bioinformatique M&P/Pythons

genotoul bioinfo

Cufflinks read attribution

- Violet fragment: from which transcript?

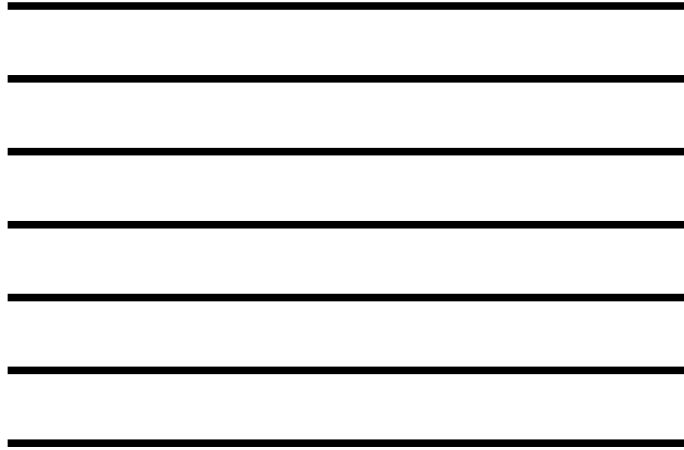
- Use of Fragment length distribution


d Abundance estimation

Transcript coverage and compatibility

Fragment length distribution

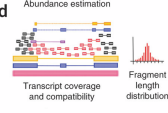
78




Cufflinks expression measurement

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


d Abundance estimation



Transcript coverage and compatibility


Fragment length distribution

e Maximum likelihood abundances



79


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


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

80

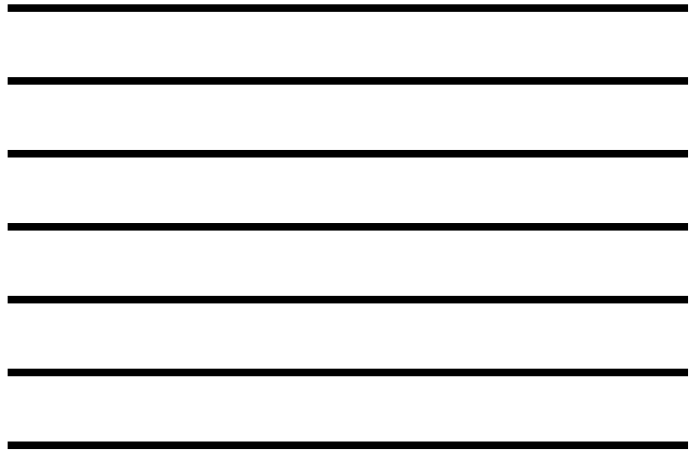

Cufflinks inputs and options

- 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

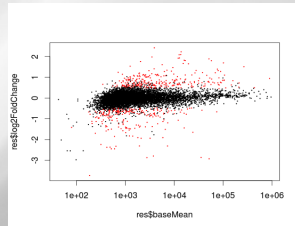
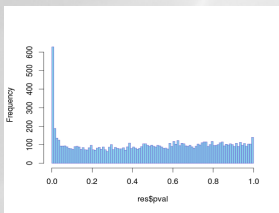
81

- Cufflinks code has been modified by the Sigenae Team of Toulouse in order to obtain raw count of reads: use **sigcufflinks** on **genotoul**
- Run cufflinks, cufflinks outputs + raw_transcripts.tsv:

gene_id	transcript_id	pairs	forward	reverse
CUFF.6	CUFF.6.1	4873	4873	3431
CUFF.6	CUFF.6.2	5222	5222	3769
CUFF.6	ENSDART0000067635	4819	4819	3580



```
> head(res)
  id baseMean baseMeanA baseMeanB foldChange log2FoldChange pval padj
1 mira_c1 3549.2301 3845.3374 3753.1228 1.1218967 0.165939787 0.375560007 0.97718309
2 mira_c2 685.7651 662.2140 709.2163 1.0711284 0.099131556 0.521137250 1.00000000
3 mira_c3 3530.8670 5096.4370 1965.2970 0.3856218 -1.374741648 0.001403322 0.03732238
4 mira_rep_c4 1012.5217 975.4453 1049.5981 1.0760194 -0.105704140 0.795193064 1.00000000
5 mira_rep_c5 1246.1199 12949.4349 12942.6048 0.9994880 -0.000738847 0.95437059 1.00000000
6 mira_rep_c6 4924.7817 5224.1292 4625.4341 0.8853981 -0.175601809 0.290161543 0.92152339
> hist(res$pval, breaks=100, col="skyblue", border="slateblue", main="")
> plotDE <- function(res) { plot(res$baseMean, res$log2FoldChange, log="x", pch="x", cex=3, col = ifelse(res$padj < .1, "red", "black")) }
> plotDE(res)
```



- 1/ Quantify the genes of chromosome 22 using htseq-count and the Ensembl GTF file for both samples.
- 2/ Quantify the genes and transcripts of chromosome 22 using sigcufflinks and the Ensembl GTF file for both samples.
- 3/ In each case merge the files to produce the count tables.



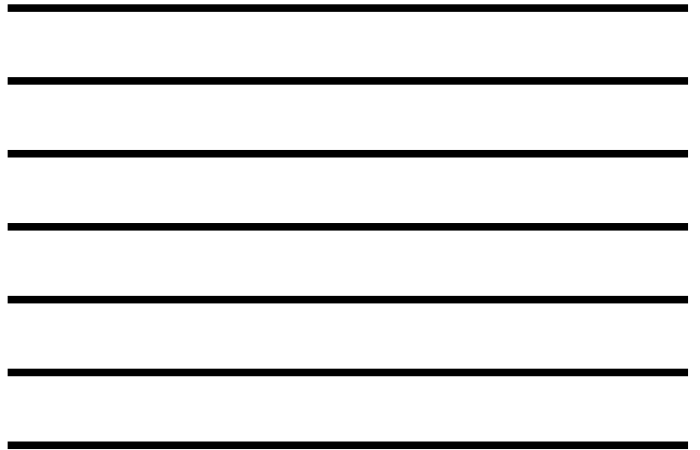
Transcript reconstruction

The different paths :

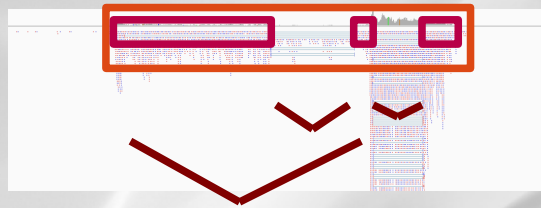
- Finding the gene locations
- Finding the exons
- Finding the junctions :
 - Between pairs junctions
 - Within sequences junction

Defining the model building strategy

- Number of built models
- Intronic reads



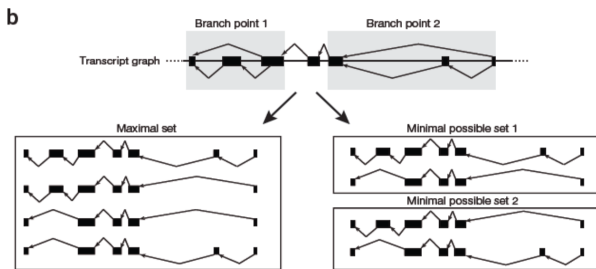
The elements of the model



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




Model building strategies



Computational methods for transcriptome annotation and quantification using RNA-seq

Manuel Garber¹, Manfred G Grabherr¹, Mitchell Guttman^{1,2} & Cole Trapnell^{1,2}




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, 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/>

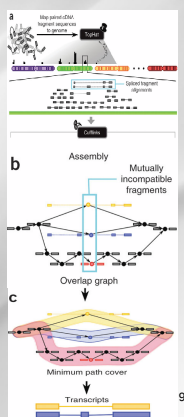
- **assembles transcripts**
- estimates their abundances : based on how many reads support each one
- tests for differential expression in RNA-Seq samples

91



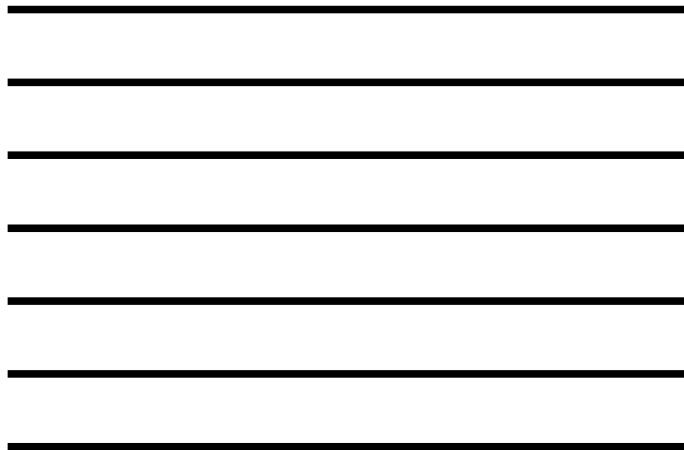

Cufflinks transcript assembly


- 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



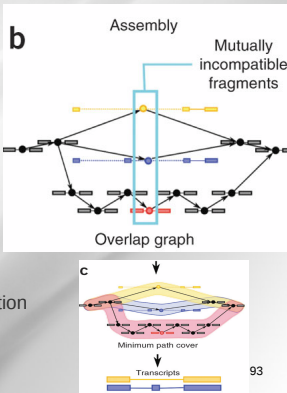
92

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




Cufflinks transcript assembly


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



93

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





Some videos of examples

- Chromosome 3 of the bovine genome, UMD3
- 3 locations
- 3 tracks :
 - Ensembl reference gene
 - Cufflinks model
 - Reads alignment

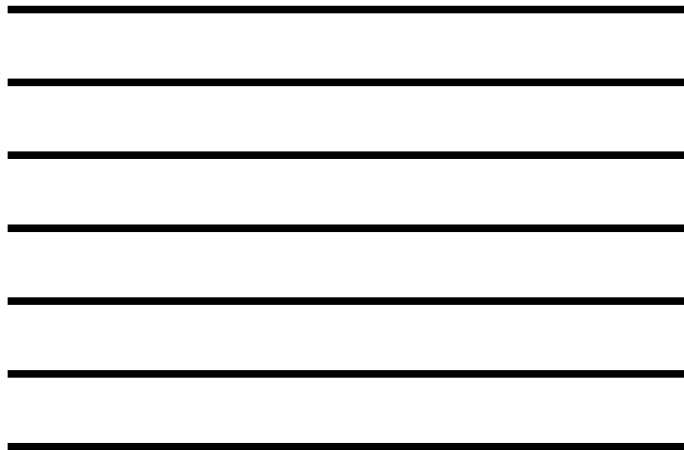
94





Cufflinks inputs and options

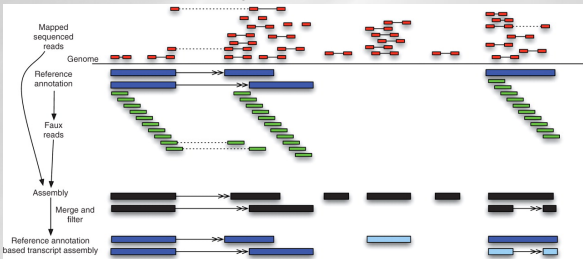
- Command line:
 - `cufflinks [options]* <aligned_reads.(sam/bam)>`
- Some options :
 - h/--help
 - o/--output-dir
 - p/--num-threads
 - G/--GTF <reference_annotation.(gtf/gff)> : estimate isoform expression, no assembly novel transcripts
 - g/--GTF-guide <reference_annotation.(gtf/gff)> : guide RABT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

95




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

96



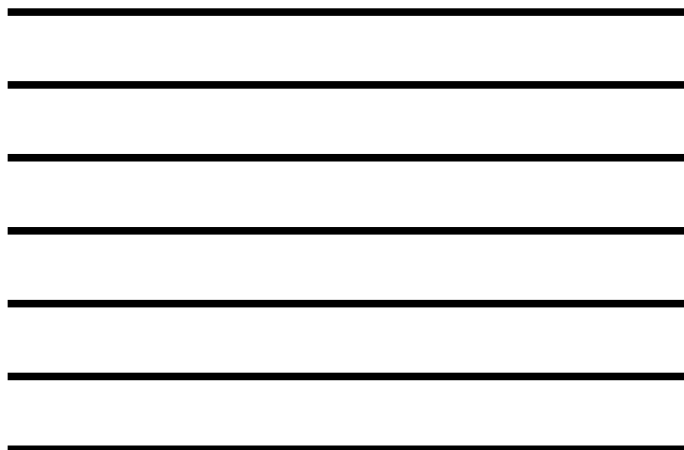
Platforme Bioinformatique Multi-PySnp

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

97



Platforme Bioinformatique Multi-PySnp

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)

GTF format

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.t", FPKM "23.7787563790", fra: "0.443485", conf_lo "0.754478", conf_hi "38.803035", cov "2.840328", full_read_support "yes",
gene_id "CUFF.560", transcript_id "CUFF.560.t", exon_number "1", FPKM "23.7787563790", fra: "0.443485", conf_lo "0.754478", conf_hi "38.803035", cov "2.840328"
```



Platforme Bioinformatique Multi-PySnp

geno toul bioinfo

Cufflinks GTF description

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

99



Platforme Biomédicale Multi-Protéomes

geno toul bioinfo

Cufflinks tracking description

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

Quantification status

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	-	-	22.9743034-9747366	-	-	OK	105.69	77.9404	133.439

- **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	24662.84033	OK	23.7788	8.75448	38.803	
CUFF.560.2	-	-	CUFF.560	-	-	22.9743034-9762309	40208.11967	OK	67.9765	50.3804	85.5727	
CUFF.560.3	-	-	CUFF.560	-	-	22.9743034-9762309	3846.16644	OK	13.9344	-	0.292533	



Platforme Biomédicale Multi-Protéomes

geno toul bioinfo

Cufflink transcript models

Gene XLOC_015511, Chr18

Condition colors

- 1332_ATCACG_1001
- 1338_ATCACG_1001
- 1339_CGATGT_1001
- 1341_TTAGGC_1001
- 1344_CGATGT_1001
- 1345_TTAGGC_1001
- 1349_TGACCA_1002
- 1350_TGACCA_1002
- 1362_ACACATG_1002
- 1393_ACACATG_1002
- 1452_GCCAAAT_1002
- 1454_GCCAAAT_1002
- 1455_ATCACG_1008
- 1455_CACATC_1003
- 1459_ACTTGA_1003
- 1459_CACATC_1003
- 1459_CACATC_1003
- 1474_ACTTGA_1008
- 1474_ACTTGA_1008
- 1474_CACATC_1003
- 1474_TAGCTT_1004
- 1475_CGATCA_1004
- 1475_TAGCTT_1008
- 1475_TAGCTT_1008
- 1477_CTTGTA_1004
- 1479_GGCTAC_1004

101

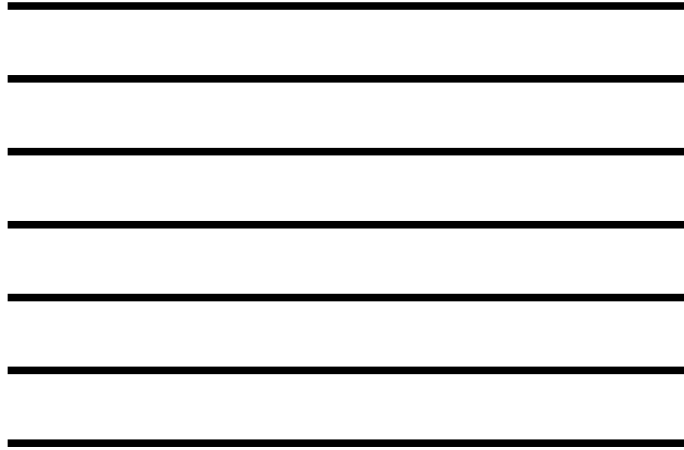


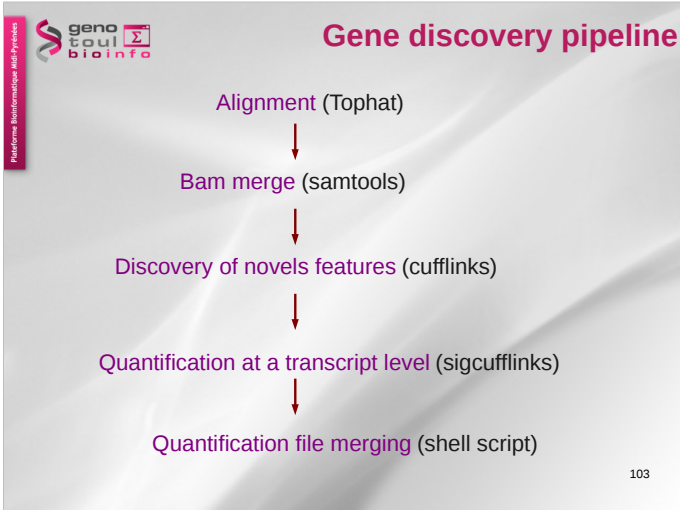
Platforme Biomédicale Multi-Protéomes

geno toul bioinfo

Cufflink transcript models

102





- Quantification strategy**
- First set your gene and transcript model = build a reference GTF file
 - Then use option -G to quantify the same set of elements on all your samples with sigcufflinks
 - Then sort your raw_transcript.tsv files
 - cut the second or third column of the sorted file
 - Paste all the column in the count file
- 104

Hands-on : cufflinks

Merge all bam : Step 5 : RNAseq De Novo

Cufflinks on merge file with -g option (reference annotation as guide) and the Danio gtf file :

Samtools merge (version 1.0.0)

Your first accepted bam file:

Datasets

Dataset 1

Other accepted bam file:

Cufflinks (version 0.2.0)

SAM or BAM file of aligned RNA-seq reads:

Max Intron Length:

Max Isoform Fraction:

Pro MiRNA Fraction:

Perform quantile normalization:

Use Reference Annotation:
 Use reference annotation as guide

Reference Annotation:

Perform Bias Correction:

Set Parameters for Paired-end Reads? (not recommended)

105

Hands-on : file merging

Sigcufflinks with the new gtf file (transcript.gtf of previously step) with -G option

Final count :

Final count file (version 1.0.0)

Select a reference genome (if your genome of interest is not listed, please contact Sigence Team):
Demo: rero dna chromosome 22

Your merged gtf file:
17: Demo_rero_chr22_Zv9-62.gtf

Your first raw transcripts tsv file from sigcufflinks:
34: Sigcufflinks_on_T_scripts_tsv

Datasets

Dataset 1

Other raw transcripts tsv file from sigcufflinks:
30: Sigcufflinks_on_T_scripts_tsv

Remove Dataset 1

Add new Dataset

Execute

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/>
