

January 15<sup>th</sup>

# PROTOCOL FOR the RNA-Seq session

## 1. Get data (~ 20 min)

Objective of the project:

Identify the genes differentially expressed in 2 human cell lines

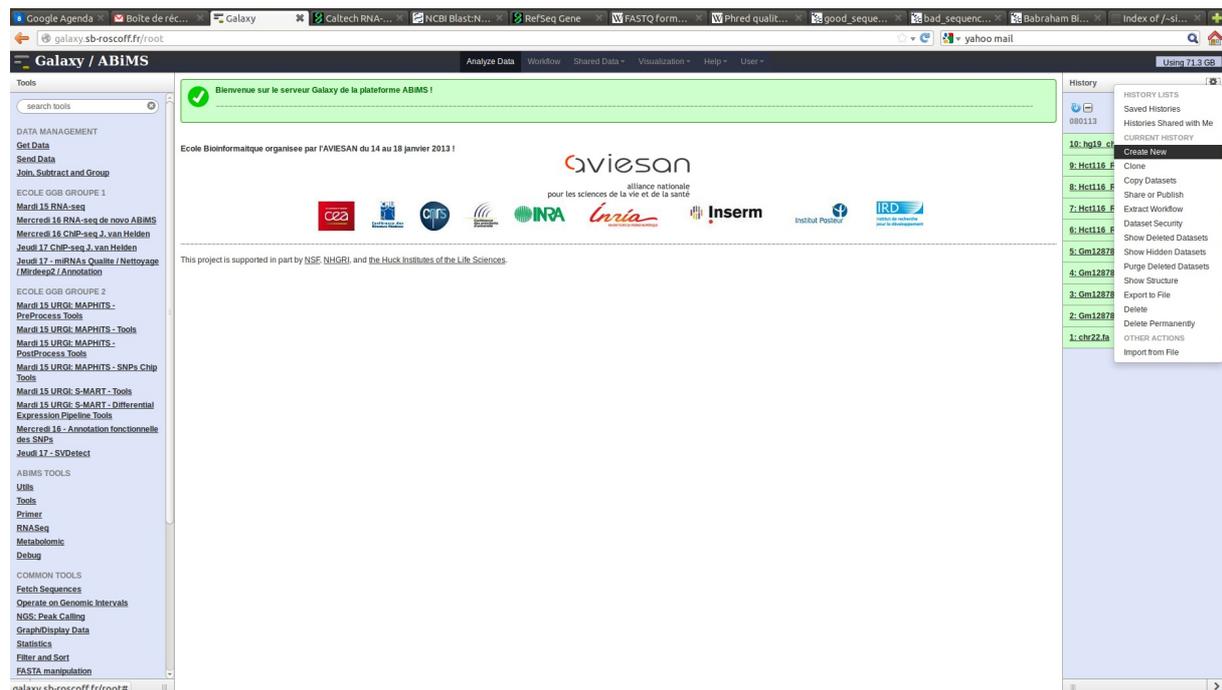
Data retrieved from the ENCODE project:

<http://genome.ucsc.edu/cgi-bin/hgFileUi?db=hg19&g=wgEncodeCaltechRnaSeq>

- 2 human cell lines: - Gm12878 (lymphoblastoid cell line)
- Hct116 (colorectal carcinoma cell line)
- each experiment is performed 2 times leading to a total of 4 samples
- Illumina sequencing, paired-ends 2x75 bp, insert size ~200 bp
- the selected genome region is the chromosome 22

Create a new history: choose a name (for example “RNA-seq”):

History → Create New



The screenshot shows the Galaxy/ABiMS web interface. The main content area displays a welcome message and logos for various institutions including C2A, CNRS, INRA, Inria, Inserm, Institut Pasteur, and IRD. On the right side, there is a 'History' panel with a list of saved histories. The list includes:

- 10: hg19 - cl
- 9: Hct116 - F
- 8: Hct116 - F
- 7: Hct116 - F
- 6: Hct116 - F
- 5: Gm12878
- 4: Gm12878
- 3: Gm12878
- 2: Gm12878
- 1: chr22.fa

The '1: chr22.fa' history is selected, and a context menu is open over it, showing options like 'Clone', 'Copy Datasets', 'Share or Publish', 'Extract Workflow', 'Dataset Security', 'Show Deleted Datasets', 'Purge Deleted Datasets', 'Show Structure', 'Export to File', 'Delete', 'Delete Permanently', and 'Import from File'.

download files :

How to download files?

In shared Data → Data Libraries → RNA Seq Delphine Mardi



Select these files:

reference file for mapping : chr22.fa

Gm12878\_R1\_rep1.fastq (raw data, read1)

Gm12878\_R2\_rep1.fastq (raw data, read2)

Gm12878\_R1\_rep2.fastq

Gm12878\_R2\_rep2.fastq

Hct116\_R1\_rep1.fastq

Hct116\_R2\_rep1.fastq

Hct116\_R1\_rep2.fastq

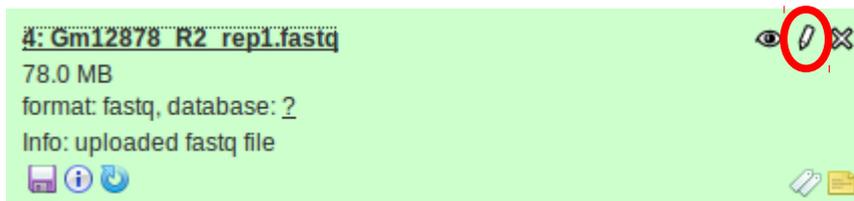
Hct116\_R2\_rep2.fastq

hg19\_chr22.gtf (gene annotation file)

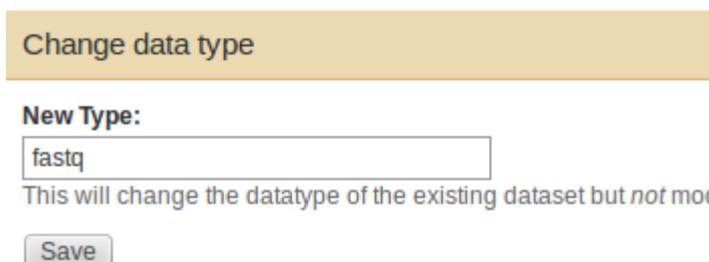
Import these files in your current history.

### **Change fastq format to fastqsanger for each file \*.fastq :**

Click on Edit attributes

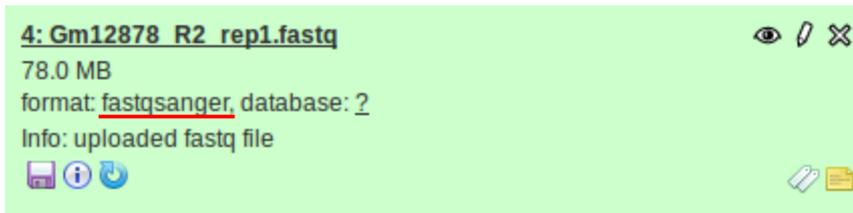


Click on Datatype



Change New Type: fastq to fastqsanger and click on save.

Check if the change has been effective:

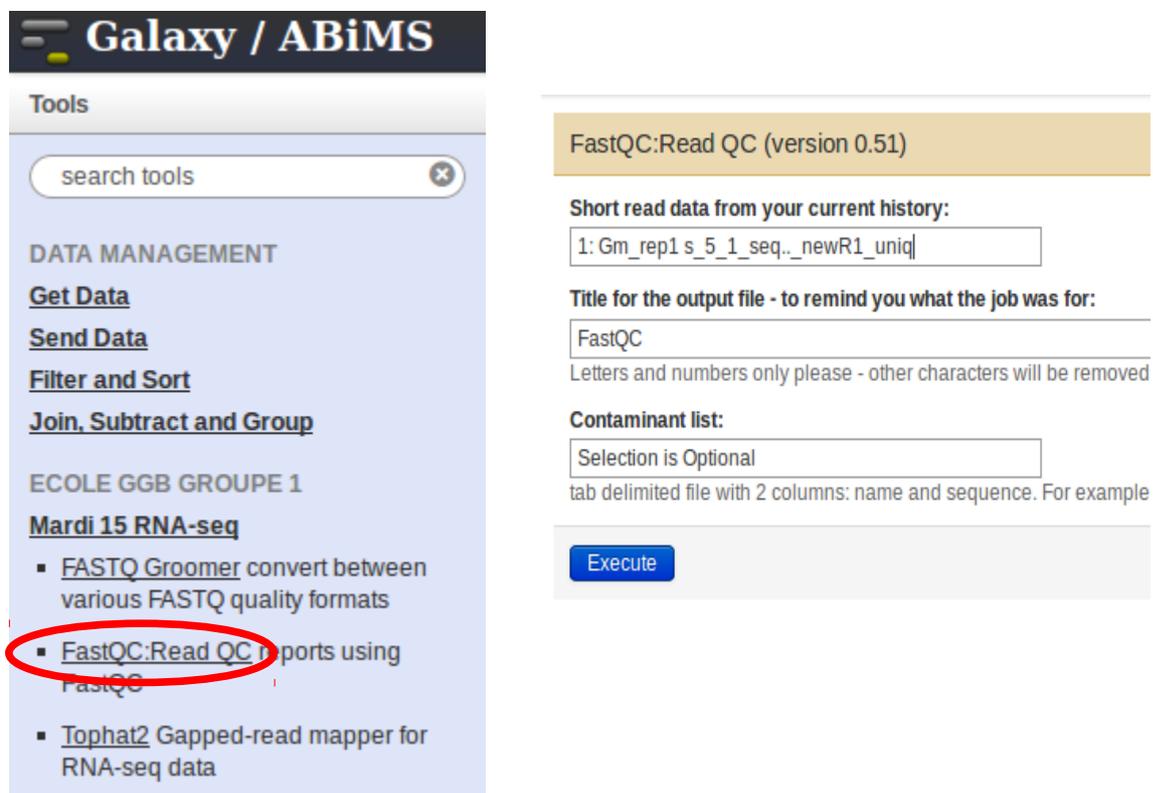


## 2. Data exploration (FastQC) (~ 40 min)

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

FastQC is a quality control tool for high throughput sequence data.

Check the quality of the data contained in the fastq files.



FastQC analysis module details

from <http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc/Help/3%20Analysis%20Modules/>

### Basic Statistics

#### Summary

The Basic Statistics module generates some simple composition statistics for the file analysed.

*Filename:* The original filename of the file which was analysed

*File type:* Says whether the file appeared to contain actual base calls or colorspace data which had to

be converted to base calls

*Encoding:* Says which ASCII encoding of quality values was found in this file.

*Total Sequences:* A count of the total number of sequences processed. There are two values reported, actual and estimated. At the moment these will always be the same. In the future it may be possible to analyse just a subset of sequences and estimate the total number, to speed up the analysis, but since we have found that problematic sequences are not evenly distributed through a file we have disabled this for now.

*Filtered Sequences:* If running in Casava mode sequences flagged to be filtered will be removed from all analyses. The number of such sequences removed will be reported here. The total sequences count above will not include these filtered sequences and will be the number of sequences actually used for the rest of the analysis.

*Sequence Length:* Provides the length of the shortest and longest sequence in the set. If all sequences are the same length only one value is reported.

*%GC:* The overall %GC of all bases in all sequences

#### Warning

Basic Statistics never raises a warning.

#### Failure

Basic Statistics never raises an error.

### **Per Base Sequence Quality**

#### Summary

This view shows an overview of the range of quality values across all bases at each position in the FastQ file.

For each position a BoxWhisker type plot is drawn. The elements of the plot are as follows:

The central red line is the median value

The yellow box represents the inter-quartile range (25-75%)

The upper and lower whiskers represent the 10% and 90% points

The blue line represents the mean quality

The y-axis on the graph shows the quality scores. The higher the score the better the base call. The background of the graph divides the y axis into very good quality calls (green), calls of reasonable quality (orange), and calls of poor quality (red). The quality of calls on most platforms will degrade as the run progresses, so it is common to see base calls falling into the orange area towards the end of a read.

It should be mentioned that there are number of different ways to encode a quality score in a FastQ file. FastQC attempts to automatically determine which encoding method was used, but in some very limited datasets it is possible that it will guess this incorrectly (ironically only when your

data is universally very good!). The title of the graph will describe the encoding FastQC thinks your file used.

### Warning

A warning will be issued if the lower quartile for any base is less than 10, or if the median for any base is less than 25.

### Failure

This module will raise a failure if the lower quartile for any base is less than 5 or if the median for any base is less than 20.

## **Per Sequence Quality Scores**

### Summary

The per sequence quality score report allows you to see if a subset of your sequences have universally low quality values. It is often the case that a subset of sequences will have universally poor quality, often because they are poorly imaged (on the edge of the field of view etc), however these should represent only a small percentage of the total sequences.

If a significant proportion of the sequences in a run have overall low quality then this could indicate some kind of systematic problem - possibly with just part of the run (for example one end of a flowcell).

### Warning

A warning is raised if the most frequently observed mean quality is below 27 - this equates to a 0.2% error rate.

### Failure

An error is raised if the most frequently observed mean quality is below 20 - this equates to a 1% error rate.

## **Per Base Sequence Content**

### Summary

Per Base Sequence Content plots out the proportion of each base position in a file for which each of the four normal DNA bases has been called.

In a random library you would expect that there would be little to no difference between the different bases of a sequence run, so the lines in this plot should run parallel with each other. The relative amount of each base should reflect the overall amount of these bases in your genome, but in any case they should not be hugely imbalanced from each other.

If you see strong biases which change in different bases then this usually indicates an overrepresented sequence which is contaminating your library. A bias which is consistent across all bases either indicates that the original library was sequence biased, or that there was a systematic problem during the sequencing of the library.

### Warning

This module issues a warning if the difference between A and T, or G and C is greater than 10% in any position.

#### Failure

This module will fail if the difference between A and T, or G and C is greater than 20% in any position.

### **Per Base GC Content**

#### Summary

Per Base GC Content plots out the GC content of each base position in a file.

In a random library you would expect that there would be little to no difference between the different bases of a sequence run, so the line in this plot should run horizontally across the graph. The overall GC content should reflect the GC content of the underlying genome.

If you see a GC bias which changes in different bases then this could indicate an overrepresented sequence which is contaminating your library. A bias which is consistent across all bases either indicates that the original library was sequence biased, or that there was a systematic problem during the sequencing of the library.

#### Warning

This module issues a warning if the GC content of any base strays more than 5% from the mean GC content.

#### Failure

This module will fail if the GC content of any base strays more than 10% from the mean GC content.

### **Per Sequence GC Content**

#### Summary

This module measures the GC content across the whole length of each sequence in a file and compares it to a modelled normal distribution of GC content.

In a normal random library you would expect to see a roughly normal distribution of GC content where the central peak corresponds to the overall GC content of the underlying genome. Since we don't know the GC content of the genome the modal GC content is calculated from the observed data and used to build a reference distribution.

An unusually shaped distribution could indicate a contaminated library or some other kinds of biased subset. A normal distribution which is shifted indicates some systematic bias which is independent of base position. If there is a systematic bias which creates a shifted normal distribution then this won't be flagged as an error by the module since it doesn't know what your genome's GC content should be.

#### Warning

A warning is raised if the sum of the deviations from the normal distribution represents more than 15% of the reads.

#### Failure

This module will indicate a failure if the sum of the deviations from the normal distribution represents more than 30% of the reads.

### **Per Base N Content**

#### Summary

If a sequencer is unable to make a base call with sufficient confidence then it will normally substitute an N rather than a conventional base] call

This module plots out the percentage of base calls at each position for which an N was called.

It's not unusual to see a very low proportion of Ns appearing in a sequence, especially nearer the end of a sequence. However, if this proportion rises above a few percent it suggests that the analysis pipeline was unable to interpret the data well enough to make valid base calls.

#### Warning

This module raises a warning if any position shows an N content of >5%.

#### Failure

This module will raise an error if any position shows an N content of >20%.

### **Sequence Length Distribution**

#### Summary

Some high throughput sequencers generate sequence fragments of uniform length, but others can contain reads of wildly varying lengths. Even within uniform length libraries some pipelines will trim sequences to remove poor quality base calls from the end.

This module generates a graph showing the distribution of fragment sizes in the file which was analysed.

In many cases this will produce a simple graph showing a peak only at one size, but for variable length FastQ files this will show the relative amounts of each different size of sequence fragment.

#### Warning

This module will raise a warning if all sequences are not the same length.

#### Failure

This module will raise an error if any of the sequences have zero length.

### **Duplicate Sequences**

#### Summary

In a diverse library most sequences will occur only once in the final set. A low level of duplication may indicate a very high level of coverage of the target sequence, but a high level of duplication is more likely to indicate some kind of enrichment bias (eg PCR over amplification).

This module counts the degree of duplication for every sequence in the set and creates a plot showing the relative number of sequences with different degrees of duplication.

To cut down on the memory requirements for this module only sequences which occur in the first 200,000 sequences in each file are analysed, but this should be enough to get a good impression for the duplication levels in the whole file. Each sequence is tracked to the end of the file to give a representative count of the overall duplication level. To cut down on the amount of information in the final plot any sequences with more than 10 duplicates are placed into the 10 duplicates category - so it's not unusual to see a small rise in this final category. If you see a big rise in this final category then it means you have a large number of sequences with very high levels of duplication.

Because the duplication detection requires an exact sequence match over the whole length of the sequence any reads over 75bp in length are truncated to 50bp for the purposes of this analysis. Even so, longer reads are more likely to contain sequencing errors which will artificially increase the observed diversity and will tend to underrepresent highly duplicated sequences.

#### Warning

This module will issue a warning if non-unique sequences make up more than 20% of the total.

#### Failure

This module will issue an error if non-unique sequences make up more than 50% of the total.

### **Overrepresented Sequences**

#### Summary

A normal high-throughput library will contain a diverse set of sequences, with no individual sequence making up a tiny fraction of the whole. Finding that a single sequence is very overrepresented in the set either means that it is highly biologically significant, or indicates that the library is contaminated, or not as diverse as you expected.

This module lists all of the sequences which make up more than 0.1% of the total. To conserve memory only sequences which appear in the first 200,000 sequences are tracked to the end of the file. It is therefore possible that a sequence which is overrepresented but doesn't appear at the start of the file for some reason could be missed by this module.

For each overrepresented sequence the program will look for matches in a database of common contaminants and will report the best hit it finds. Hits must be at least 20bp in length and have no more than 1 mismatch. Finding a hit doesn't necessarily mean that this is the source of the contamination, but may point you in the right direction. It's also worth pointing out that many adapter sequences are very similar to each other so you may get a hit reported which isn't technically correct, but which has a very similar sequence to the actual match.

Because the duplication detection requires an exact sequence match over the whole length of the sequence any reads over 75bp in length are truncated to 50bp for the purposes of this analysis.

Even so, longer reads are more likely to contain sequencing errors which will artificially increase the observed diversity and will tend to underrepresent highly duplicated sequences.

#### Warning

This module will issue a warning if any sequence is found to represent more than 0.1% of the total.

#### Failure

This module will issue an error if any sequence is found to represent more than 1% of the total.

### **Overrepresented Kmers**

#### Summary

The analysis of overrepresented sequences will spot an increase in any exactly duplicated sequences, but there are a different subset of problems where it will not work.

If you have very long sequences with poor sequence quality then random sequencing errors will dramatically reduce the counts for exactly duplicated sequences.

If you have a partial sequence which is appearing at a variety of places within your sequence then this won't be seen either by the per base content plot or the duplicate sequence analysis.

This module counts the enrichment of every 5-mer within the sequence library. It calculates an expected level at which this k-mer should have been seen based on the base content of the library as a whole and then uses the actual count to calculate an observed/expected ratio for that k-mer. In addition to reporting a list of hits it will draw a graph for the top 6 hits to show the pattern of enrichment of that Kmer across the length of your reads. This will show if you have a general enrichment, or if there is a pattern of bias at different points over your read length.

Any k-mer showing more than a 3 fold overall enrichment or a 5 fold enrichment at any given base position will be reported by this module.

To allow this module to run in a reasonable time only 20% of the whole library is analyzed and the results are extrapolated to the rest of the library.

#### Warning

This module will issue a warning if any k-mer is enriched more than 3 fold overall, or more than 5 fold at any individual position.

#### Failure

This module will issue an error if any k-mer is enriched more than 10 fold at any individual base position.

Compare the FastQC reports of samples Gm12878\_R1\_rep2.fastq and Hct116\_R1\_rep2.fastq.

How can you explain the presence of highly duplicated sequences?

Is it normal?

### 3. Mapping with TopHat2 (~ 45 min)

<http://tophat.cbc.umd.edu/manual.html>

#### 3.1 Slides

#### 3.2 TP

Mapping of the paired end reads on chr22 – this must be performed for all 4 samples.

ECOLE GGB GROUPE 1

**Mardi 15 RNA-seq**

- [FASTQ Groomer](#) convert between various FASTQ quality formats
- [FastQC:Read QC](#) reports using FastQC
- **Tophat2** Gapped-read mapper for RNA-seq data
- [flagstat](#) provides simple stats on BAM files

Use the default parameters except for :

- library type : paired-end ;
- the “mean inner distance between mate pairs” (= 200 bp) ;
- select a reference genome from your history : chr22.fa ;

**Tophat2 (version 0.5)**

**Is this library mate-paired?:**

**RNA-Seq FASTQ file, forward reads:**  
  
Nucleotide-space: Must have Sanger-scaled quality values with ASCII offset 33

**RNA-Seq FASTQ file, reverse reads:**  
  
Nucleotide-space: Must have Sanger-scaled quality values with ASCII offset 33

**Mean Inner Distance between Mate Pairs:**

**Will you select a reference genome from your history or use a built-in index?:**  
  
Built-ins were indexed using default options

**Select the reference genome:**

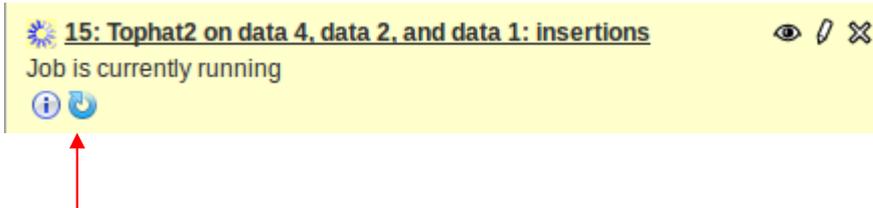
- In the full parameter list :

- maximum number of alignments to be allowed = 1

**Maximum number of alignments to be allowed:**

→ Execute!

To repeat this step 3 times, you can use “run this job again” and just change the input data:



**15: Tophat2 on data 4, data 2, and data 1: insertions** 👁️ ✎️ ✕  
Job is currently running  
📘 ↻

Check the results of Tophat2 with flagstat tool (samtools) on BAM files:



**ECOLE GGB GROUPE 1**  
**Mardi 15 RNA-seq**

- [FASTQ Groomer](#) convert between various FASTQ quality formats
- [FastQC:Read QC](#) reports using FastQC
- [Tophat2](#) Gapped-read mapper for RNA-seq data
- **flagstat** provides simple stats on BAM files

**flagstat (version 1.0.0)**

**BAM File to Convert:**  
9: Gm\_rep1 Tophat2 o..cepted\_hits

**Execute**

## 4. Transcript assembly with Cufflinks (~ 40 min)

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

4.1 Slides

4.2 TP

ECOLE GGB GROUPE 1

**Mardi 15 RNA-seq**

- FASTQ Groomer convert between various FASTQ quality formats
- FastQC:Read QC reports using FastQC
- TopHat2 Gapped-read mapper for RNA-seq data
- flagstat provides simple stats on BAM files
- Cufflinks transcript assembly and FPKM (RPKM) estimates for RNA-Seq data
- Cuffcompare compare assembled transcripts to a reference annotation and track Cufflinks transcripts across multiple experiments

Use as input file the file of mapping results of TopHat2 (in BAM format) named “\*\_accepted\_hits”.

Cufflinks is originally designed for human data, hence use the default parameters.

Execute and repeat the procedure on the 2 mapping results of TopHat2 (rep1 for each cell line for example).

### Cufflinks (version 0.0.5)

SAM or BAM file of aligned RNA-Seq reads:

19: Hct\_rep1 Tophat2 ..cepted\_hits

Max Intron Length:

300000

Min Isoform Fraction:

0.1

Pre mRNA Fraction:

0.15

Perform quartile normalization:

No

Removes top 25% of genes from FPKM denominator to improve accuracy of differential expression calls for low abundance transcripts.

Use Reference Annotation:

No

Perform Bias Correction:

No

Bias detection and correction can significantly improve accuracy of transcript abundance estimates.

Use multi-read correct:

No

Tells Cufflinks to do an initial estimation procedure to more accurately weight reads mapping to multiple locations in the genome.

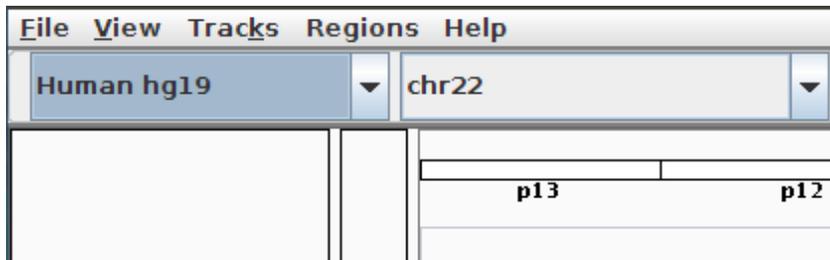
Execute

### 4.3 Visualization with IGV

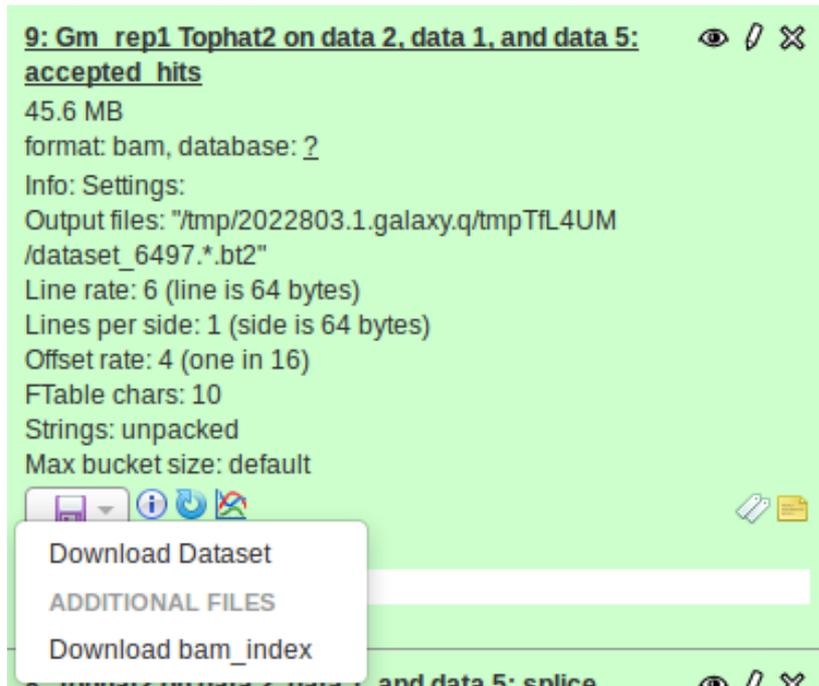
<http://www.broadinstitute.org/igv/>

The **Integrative Genomics Viewer (IGV)** is a high-performance visualization tool for interactive exploration of large integrated genomic datasets.

Open IGV with the shortcut in your desktop, choose the human genome Hg19 (default) and select the chromosome 22:



Download a result of TopHat (bam file) with its index:



- load the bam file in IGV:

file → Load from File

- download the corresponding result from Cufflinks (“assembled transcript” file)

- then load the assembled transcripts file in IGV.

#### 4.4 Cuffcompare launch (comparison with a genome annotation) (~ 20 min)

##### 4.4.1 Slides

##### 4.4.2 TP

This tool enables to compare the assembled transcripts to a reference annotation.

ECOLE GGB GROUPE 1

**Mardi 15 RNA-seq**

- [FASTQ Groomer](#) convert between various FASTQ quality formats
- [FastQC:Read QC](#) reports using FastQC
- [Tophat2](#) Gapped-read mapper for RNA-seq data
- [flagstat](#) provides simple stats on BAM files
- [Cufflinks](#) transcript assembly and FPKM (RPKM) estimates for RNA-Seq data
- **Cuffcompare** compare assembled transcripts to a reference annotation and track Cufflinks transcripts across multiple experiments

Use a reference annotation and use reference file, which is in your history (“hg19\_chr22.gtf”).

Cuffcompare (version 0.0.5)

**GTF file produced by Cufflinks:**  
51: Cufflinks on data..transcripts

**Additional GTF Input Files**  
Add new Additional GTF Input Files

**Use Reference Annotation:**  
Yes

**Reference Annotation:**  
34: hg19\_chr22.gtf  
Make sure your annotation file is in GTF format and that Galaxy knows that your file is GTF--not GFF.

**Ignore reference transcripts that are not overlapped by any transcript in input files:**

**Use Sequence Data:**  
Yes  
Use sequence data for some optional classification functions, including the addition of the p\_id attribute required by Cuffdiff.

**Choose the source for the reference list:**  
History

**Using reference file:**  
5: chr22.fa

Execute

Execute!

The following table shows the code used by Cuffcompare to classify the transcripts issued from Cufflinks in comparison with the reference annotation:

Priority	Code	Description	
1	=	Complete match of intron chain	
2	c	Contained	
3	j	Potentially novel isoform (fragment): at least one splice junction is shared with a reference transcript	
4	e	Single exon transfrag overlapping a reference exon and at least 10 bp of a reference intron, indicating a possible pre-mRNA fragment.	
5	i	A transfrag falling entirely within a reference intron	
6	o	Generic exonic overlap with a reference transcript	
7	p	Possible polymerase run-on fragment (within 2Kbases of a reference transcript)	
8	r	Repeat. Currently determined by looking at the soft-masked reference sequence and applied to transcripts where at least 50% of the bases are lower case	
9	u	Unknown, intergenic transcript	
10	x	Exonic overlap with reference on the opposite strand	
11	s	An intron of the transfrag overlaps a reference intron on the opposite strand (likely due to read mapping errors)	

## ==== Afternoon: Differential Gene Expression ====

### 5. File conversion (Bam to Sam) (~ 10 min)

This step must be performed for all 4 mapping results from TopHat2 in order to launch htseq-count for each sample.

ECOLE GGB GROUPE 1

**Mardi 15 RNA-seq**

- [FASTQ Groomer](#) convert between various FASTQ quality formats
- [FastQC:Read QC](#) reports using FastQC
- [Tophat2](#) Gapped-read mapper for RNA-seq data
- [flagstat](#) provides simple stats on BAM files
- [Cufflinks](#) transcript assembly and FPKM (RPKM) estimates for RNA-Seq data
- [Cuffcompare](#) compare assembled transcripts to a reference annotation and track Cufflinks transcripts across multiple experiments
- [Cuffdiff](#) find significant changes in transcript expression, splicing, and promoter use
- [Paired Read Mate Fixer](#) for paired data
- [htseq-count](#) - Count aligned reads in a BAM file that overlap features in a GFF file
- [DESeq](#) Determines differentially expressed transcripts from read alignments
- **BAM-to-SAM** converts BAM format to SAM format

**BAM-to-SAM (version 1.0.3)**

**BAM File to Convert:**

25: Hct\_rep2 Tophat2 ..cepted\_hits

**Include header in output:**

**Execute**

This step enables to convert a binary format into a readable alignment format.

## 6. Sort file (~ 10 min)

Sort these new SAM files by alphabetical order of the read name (column 1) in order to be accepted by the Htseq-count.

**Filter and Sort**

- Filter data on any column using simple expressions
- **Sort** data in ascending or descending order
- Select lines that match an expression

Sort (version 1.0.1)

Sort Query:  
26: Gm\_rep1 BAM-to-SA..nverted SAM

on column:  
c1

with flavor:  
Alphabetical sort

everything in:  
Ascending order

Column selections  
Add new Column selection

Execute

## 7. Htseq-count (~ 30 min)

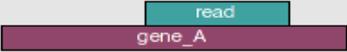
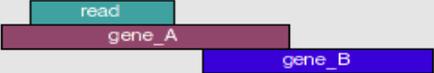
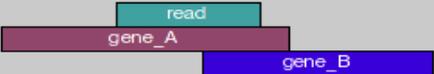
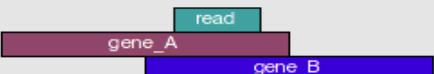
<http://www-huber.embl.de/users/anders/HTSeq/doc/count.html>

This tool counts, for each gene present in the gtf file, the number of reads which maps on it.

**Mardi 15 RNA-seq**

- FASTQ Groomer convert between various FASTQ quality formats
- FastQC:Read QC reports using FastQC
- Tophat2 Gapped-read mapper for RNA-seq data
- flagstat provides simple stats on BAM files
- Cufflinks transcript assembly and FPKM (RPKM) estimates for RNA-Seq data
- Cuffcompare compare assembled transcripts to a reference annotation and track Cufflinks transcripts across multiple experiments
- Cuffdiff find significant changes in transcript expression, splicing, and promoter use
- Paired Read Mate Fixer for paired data
- **htseq-count** Count aligned reads in a BAM file that overlap features in a GFF file

Change mode “Union” by “intersection (non empty)” and stranded “yes” by “no”.

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

## htseq-count (version 0.2.1)

### Aligned SAM/BAM File:

33: Hct\_rep2 Sort on data 29

Paired-End data MUST be sorted by QUERY NA

### GFF File:

34: hg19\_chr22.gtf

### Mode:

Intersection (nonempty)

Mode to handle reads overlapping more than on

### Stranded:

No

Specify whether the data is from a strand-specific

### Minimum alignment quality:

0

Skip all reads with alignment quality lower than t

### Feature type:

exon

Feature type (3rd column in GFF file) to be used.

### ID Attribute:

gene\_id

GFF attribute to be used as feature ID. Several G value for this attribute. The default, suitable for R

### Additional BAM Output:

Write out all SAM alignment records into an outp

Execute

Repeat this step on the 4 sorted sam files.

## 8. Generate Htseq-count result path File (~ 5 min)

In order to use the DESeq soft from the URGI pipeline (presented in the next session), we must construct a file that contains all paths of htseq-count results on the galaxy server.

### Mardi 15 URGI: S-MART - Differential Expression Pipeline Tools

#### DIFFERENTIAL EXPRESSION PIPELINE TOOLS

- [load\\_multiFASTQfiles](#) To load several FASTQ files from different conditions.
- [FASTQ Groomer parallel](#) convert between various FASTQ quality formats for a list of inputs
- [Tophat parallel for Illumina](#) Find splice junctions using RNA-seq data, can have several input RNA-seq data.
- [BAM to SAM](#) converts a list of BAM format files to SAM format
- [CompareOverlapping\\_parallel](#) Shrink or extend the sets of genomic coordinates to get the information between starts of reads and starts of genes.
- [countNumber\\_parallel](#) Calculate the number of reads(annotations) overlapping for each transcript.
- [DESEQ for differential expression analysis](#) Differential expression analysis for reads count data
- [load HTSeqResultFiles](#) To load several HTSeq result files from different conditions.

### load HTSeqResultFiles (version 1.0.0)

#### Condition groups

##### Condition group 1

##### Replicates

##### Replicate 1

##### TABULAR file.:

##### Replicate 2

##### TABULAR file.:

##### Condition group 2

##### Replicates

##### Replicate 1

##### TABULAR file.:

##### Replicate 2

##### TABULAR file.:

## 9. DESeq (~ 1h20)

<http://www-huber.embl.de/users/anders/DESeq/>

DESeq is a tool for computing the differential expression of genes using RNA-Seq data.

Check “If there is a header for your count files, please choose this case”.

**Mardi 15 URGI: S-MART - Differential Expression Pipeline Tools**

**DIFFERENTIAL EXPRESSION PIPELINE TOOLS**

- [load\\_multiFASTQfiles](#) To load several FASTQ files from different conditions.
- [FASTQ Groomer parallel](#) convert between various FASTQ quality formats for a list of inputs
- [Tophat parallel for Illumina](#) Find splice junctions using RNA-seq data, can have several input RNA-seq data.
- [BAM to SAM](#) converts a list of BAM format files to SAM format
- [CompareOverlapping\\_parallel](#) Shrink or extend the sets of genomic coordinates to get the information between starts of reads and starts of genes.
- [countNumber\\_parallel](#) Calculate the number of reads(annotations) overlapping for each transcript.
- **[DESEQ for differential expression analysis](#)** Differential expression analysis for reads count data
- [load HTSeqResultFiles](#) To load several HTSeq result files from different conditions.

**DESEQ for differential expression analysis (version 1.0.0)**

**Input File list:**  
43: HTSeq result files

**If there is a header for your count files, please choose this case.:**

**If your data has not replicates, please choose this case.:**

**Execute**

The DESeq tool performs the data normalization and computes the statistical significance of the differential expression for each gene; it also generates some graphs to visualize statistic analyses of your data. Finally, you obtain a file containing genes that are over-expressed in one condition compared to the other, and a file containing under-expressed genes. These lists allow performing additional functional category analysis, for example GO term.

January 18<sup>th</sup>

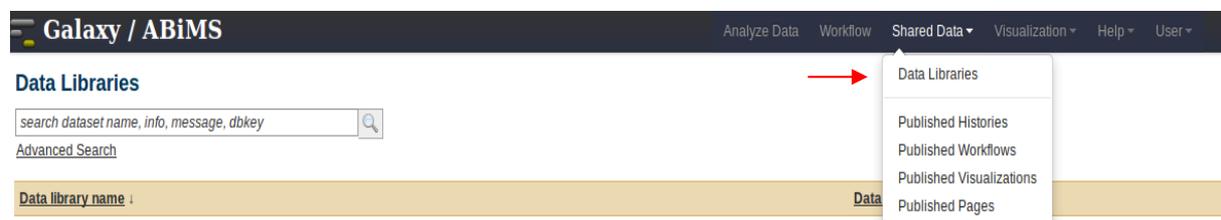
## Exercises of RNA-Seq session

### Data (for the participants without their own data)

Data retrieved from the study of van-Dijk and coworkers (*Nature* 2011, 475:114)

- yeast transcriptome: - cells grown in CSM medium
  - cells grown in lithium-containing CSM medium
- each experiment is performed at least 2 times
- Illumina sequencing\*, single-end 38 nt
- the selected genome regions are the chr 5 and chr 8

In shared Data → data libraries → RNA Seq Friday



reference file for mapping : S.cerevisiae\_atelier\_ok.fa

gene annotation file : S.cerevisiae\_atelier\_ok.gtf

Raw sequencing data\*:

- LiM.Rep1.fastq (lithium -, replicate 1)
- LiM.Rep2.fastq (lithium -, replicate 2)
- LiP.Rep1.fastq (lithium +, replicate 1)
- LiP.Rep2.fastq (lithium +, replicate 2)

\*in fastq Illumina 1.3-1.7 format, need to be converted to fastqsanger format by “FASTQ Groomer”

## Objectives

Analysis of the differential expression of genes in different growth conditions

Hints: I. Main steps:

- quality control,
- mapping,
- counting,
- differential expression

– identification of Go terms

II. Use the protocol of last Tuesday as a reference

III. Try to use the workflow

**Bonus:** since the number of biological duplicates influence strongly the statistical evaluation in the differential expression analysis, it will be interesting for you to check if adding a third duplicate would change the results. A 3<sup>rd</sup> biological replicate of LiP (LiP.Rep3.fastq) is available for this test.

What do you observe?

What can you conclude?