

- Galaxy -

Formation à l'analyse de données RNA-seq

- EXERCICES -

Liens utiles

Données publiques :



The European Nucleotide Archive (ENA) provides a comprehensive record of the world's nucleotide sequencing information, covering raw sequencing data, sequence assembly information and functional annotation.




<http://www.ebi.ac.uk/ena/>







The Ensembl project produces genome databases for vertebrates and other eukaryotic species, and makes this information freely available online.

<http://www.ensembl.org/index.html>

Logiciels utilisés :

	<p>FastQC aims to provide a simple way to do some quality control checks on raw sequence data coming from high throughput sequencing pipelines. It provides a modular set of analyses which you can use to give a quick impression of whether your data has any problems of which you should be aware before doing any further analysis.</p> <p>http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc/</p>
	<p>Cutadapt Cutadapt finds and removes adapter sequences, primers, poly-A tails and other types of unwanted sequence from your high-throughput sequencing reads. Cutadapt helps to trim reads by finding the adapter or primer sequences in an error-tolerant way. It can also modify and filter reads in various ways. Adapter sequences can contain IUPAC wildcard characters. Also, paired-end reads and even colorspace data is supported. If you want, you can also just demultiplex your input data, without removing adapter sequences at all.</p>
<p>STAR</p>	<p>STAR is a Spliced Transcripts Alignment to a Reference.</p> <p>https://github.com/alexdobin/STAR</p>
	<p>Cufflinks assembles transcripts, estimates their abundances, and tests for differential expression and regulation in RNA-Seq samples. It accepts aligned RNA-Seq reads and assembles the alignments into a parsimonious set of transcripts. Cufflinks then estimates the relative abundances of these transcripts based on how many reads support each one, taking into account biases in library preparation protocols.</p> <p>http://cufflinks.cbc.umd.edu/</p>
<p>SAMtools</p>	<p>SAM (Sequence Alignment/Map) format is a generic format for storing large nucleotide sequence alignments. http://samtools.sourceforge.net/</p>



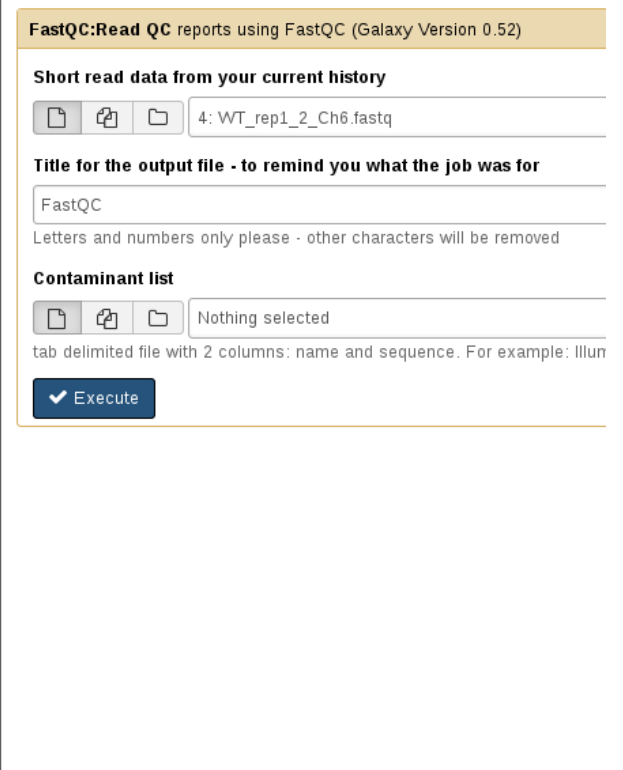


	RSEM: accurate quantification of gene and isoform expression from RNA-Seq data
	The Integrative Genomics Viewer (IGV) is a high-performance visualization tool for interactive exploration of large, integrated datasets. It supports a wide variety of data types including sequence alignments, microarrays, and genomic annotations. http://www.broadinstitute.org/igv/
	Bioconductor provides tools for the analysis and comprehension of high-throughput genomic data. Bioconductor uses the R statistical programming language, and is open source and open development. http://bioconductor.org/
	R is a free software environment for statistical computing and graphics. R provides a wide variety of statistical (linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering, ...) and graphical techniques, and is highly extensible. http://www.r-project.org/

Objectifs:

Cette formation a pour objectif de vous aider à traiter les séquences issues des SGS (Second Generation Sequencing), en particulier les plates-formes Illumina HiSeq. Vous y découvrirez les nouveaux formats de séquences, les biais connus et mettrez en œuvre des logiciels d'alignement épissé sur génome de référence, la recherche de nouveaux gènes, de nouveaux transcrits et la quantification de l'expression de ces gènes et transcrits.


Exercice n°1: Analyse de la qualité des données avec FastQC

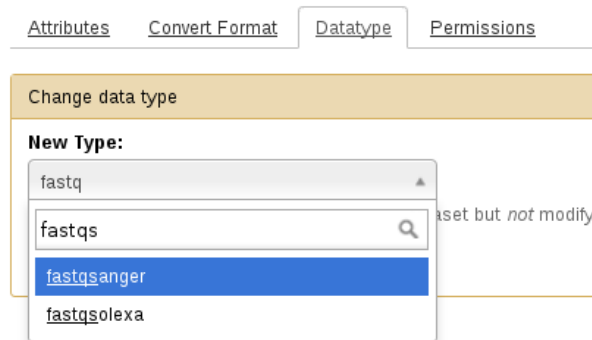
<p>Recherchez l'outil « FastQC:Read QC reports using FastQC » à l'aide du champs “search” de Galaxy:</p>	
<p>Puis lancez cet outil, 4 fois, pour chacun des fichiers fastq. Renommez les fichiers obtenus :</p> <p>FastQC MT_rep1_1_Ch6.fastq FastQC MT_rep1_2_Ch6.fastq FastQC WT_rep1_1_Ch6.fastq FastQC WT_rep1_2_Ch6.fastq</p> 	

- 1) Quelle est la longueur des lectures ? Est-ce la même que celle que vous avez obtenue à l'exercice précédent ?
- 2) La qualité du séquençage vous paraît-elle correcte ?
- 3) Regardez les résultats concernant les biais décrits lors du cours : lesquels retrouve-t-on ?

Exercice n°2: Nettoyage des données avec cutadapt et Sickle

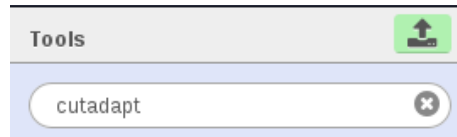
1) Supprimez les adaptateurs à l'aide de cutadapt (outil "Remove adaptators with cutadapt").

L'outil cutadapt ne reconnaît en entrée les fichiers fastq que si ces derniers sont au format "fastqsanger". Pour modifier le datatype de fastq à fastqsanger, veuillez cliquer sur l'icône  pour accéder au menu "datatype".

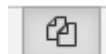


Spécifiez que la taille minimale des reads trimmées doit être de 36 bp.

Adaptateurs : -a AGATCGGAAGAGCACACGTCTGAACTCCAGTCAC -A AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGTAGATCTCGGTGGTCGCCGTATCATT



Il est possible de traiter plusieurs fichiers fastq en entrée de cutadapt afin de ne pas avoir à relancer l'outil 4 fois. Il s'agit du mode "batch". Pour cela, il suffit de cliquer sur l'icône "Multiple datasets"




L'outil cutadapt sera lancé une première fois pour les reads R1 pour enlever l'adaptateur AGATCGGAAGAGCACACGTCTGAACTCCAGTCAC puis cutadapt sera lancé une seconde fois en R2 pour supprimer l'adaptateur AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGTAGATCTCGGTGGTCGCCGTATCATT De plus, spécifiez que la taille minimale des reads trimmées doit être de 36 bp.



Remove adaptators with cutadapt (Galaxy Version 1.0.0)

Your Illumina 'lane' file (fastqsanger format)

 This is a batch mode input field. A separate job will be triggered for each dataset.

Complete your(s) primer(s)

Choose the way you add your primer

Primer used

Minimum size (in pb)

Maximum size (in pb)

Quality base

Chaque job génère deux fichiers sortants pour un fichier fastq entrant : un fichier fastq nettoyé {MT_rep1_1_Ch6.fastq}-cutadapt.fastq, et un fichier de log {MT_rep1_1_Ch6.fastq}-cutadapt.log

Pensez aussi à aller voir les log du logiciel.

2) Testez un autre type de nettoyage avec l'outil "Sickle Window Adaptive Trimming of FastQ data", enlevez les reads trimmées lorsqu'elles sont plus courtes que 36 pbs, pour WT et pour MT.

Aidez-vous du manuel pour bien paramétrer Sickle: <https://github.com/najoshi/sickle>



Tools sickle

Generic FASTQ manipulation

[Sickle Window Adaptive Trimming of FastQ data](#)

Sickle Window Adaptive Trimming of FastQ data (Galaxy Version 1.0.0)

Single-End or Paired-End reads?
Paired-End

Paired-End Forward Strand FastQsanger Reads
3: WT_rep1_1_Ch6.fastq

Paired-End Reverse Strand FastQsanger Reads
4: WT_rep1_2_Ch6.fastq

Quality Threshold
20

Length Threshold
36

Disable 5'-end trimming
Yes No

Discard any sequence with any number of Ns
Yes No

Execute

Pensez à renommer vos datasets sortantes.

- 3) A quelle valeur de qualité sickle trimme-t-il par défaut ?
- 4) Relancez fastQC sur les résultats.
- 5) Quel est le nombre de séquences après nettoyage ?

Exercice n°3: alignement/visualisation

1) Alignement avec STAR

L'outil "Map with STAR 2.4.0i with GTF and reference" indexe la référence, ou bien utilise les index déjà disponibles sur le serveur Genotoul, dans /bank/STARdb/, puis réalise les alignements épissés.

- Recherchez le manuel de « rnaSTAR » sur internet.
- Quelle version de STAR est utilisée par défaut ?
- Chargez le fichier fasta de référence à l'aide de l'outil "Upload File from your computer" dans la section "Get data" de Galaxy. Voici l'URL d'accès au fichier :

RNAseq Alignement

[Map with STAR 2 \(STAR_2.4.0i\)](#)
aligns RNA-seq single and paired reads

[Map with STAR 2.4.0i with GTF and reference](#)

"http://genoweb.toulouse.inra.fr/~formation/4_Galaxy_RNAseq/2018/data/reference/ITAG2.3_genomic_Ch6.fasta"

- Renommez la dataset sortante en "ITAG2.3_genomic_Ch6.fasta".
- Lancez "Map with STAR 2.4.0i with GTF and reference" en indiquant les paramètres suivants:

Map with STAR 2.4.0i with GTF and reference (Galaxy Version 1.0.0) Options

Paired or single reads
Paired reads

First input fastq gzipped file (read1.fastq.gz)
5: http://genoweb.toulouse.inra.fr/~formation/4_Galaxy_RNAseq/2018/data/reads/WT_rep1_1_Ch6.fastq

Second input fastq gzipped file (read2.fastq.gz)
6: http://genoweb.toulouse.inra.fr/~formation/4_Galaxy_RNAseq/2018/data/reads/WT_rep1_2_Ch6.fastq

compressed fastq file
Not compressed
fastq files are compressed or not

Threads number
8

alignIntronMin
20

alignIntronMax
1000000

outFilterMismatchNmax
10

Genotoul reference genome or your own fasta file
Your own fasta file

Your own reference genome
2: http://genoweb.toulouse.inra.fr/~formation/4_Galaxy_RNAseq/2015/data/reference/ITAG2.3_genomic_Ch6....

Your own GTF file
1: http://genoweb.toulouse.inra.fr/~formation/4_Galaxy_RNAseq/2015/data/reference/ITAG_pre2.3_gene_mo...

Do you want to perform Cufflinks quantification after ?
Yes
Cufflinks needs strand information. This option add STAR options: --outSAMstrandField intronMotif --outFilterIntronMotifs RemoveNoncanonical --outFilterType BySJout

Transcriptome alignment
Yes
STAR will outputs alignments translated into transcript coordinates (for RSEM, eXpress, etc.), works if indexation has been performed with gtf.

- f) Quels sont les fichiers de sortie ? Les renommer (WT_...)
- g) Combien de reads sont alignés de façon unique et de façon multiple ? (voir Log.final.out)

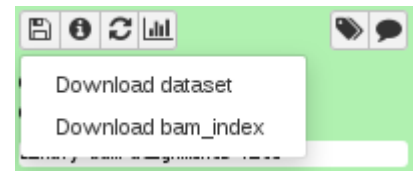


2) Alignement avec HISAT2

Lancez “**HISAT2** A fast and sensitive alignment program (Galaxy Version 2.0.5.2)” sur une paire de fastq. **HISAT2 remplace Tophat.**

3) Visualisation des résultats de STAR

Téléchargez sur votre ordinateur les fichiers de résultats de STAR (genome bam et SJ.out.tab) et le fichier d'indexation (bai)



- a) Utilisez IGV pour visualiser les résultats sur votre poste de travail.
- b) Lancez IGV depuis « download » du site web de la formation (en bas de la page): <http://www.broadinstitute.org/software/igv/download>
 - a) Chargez les annotations (fichier gtf mis à disposition dans <http://genoweb.toulouse.inra.fr/~formation/LigneCmd/RNAseq/2017/data/reference/>)
 - b) Chargez les .bam et les .bai
 - c) Explorez l'interface, en utilisant le clic-droit (pour visualiser toutes les isoformes, les paires de lectures associées....)
 - d) Regardez les régions suivantes :
 - ~ SL2.40ch06:34,298,666-34,306,292
 - ~ SL2.40ch06:34,209,900-34,260,000
 - ~ SL2.40ch06:2,786,806-2,807,064
 - ~ SL2.40ch06:38,479,173-38,483,269
 - ~ SL2.40ch06:10,694,176-10,704,838
 - ~ Solyc06g009140.2.1
 - ~ SL2.40ch06:7,973,823-7,977,708

NB. Pour des jeux de données de tailles plus conséquentes, pensez à trier et indexer vos fichiers GTF avec les outils igvtools sort et igvtools index.

Exercice 4 : Recherche de nouveaux transcrits et comparaison des gtf

Cufflinks transcript assembly and FPKM (RPKM) estimates for RNA-Seq data (Galaxy Version 0.0.7)

SAM or BAM file of aligned RNA-Seq reads

Max Intron Length

Min Isoform Fraction

Pre mRNA Fraction

Perform quartile normalization

Removes top 25% of genes from FPKM denominator to improve accuracy of differential expression calls

Use Reference Annotation

Reference Annotation

Gene annotation dataset in GTF or GFF3 format.

- 1) Dans Cufflinks, que signifie RABT ?
- 2) Quelle version de cufflinks est disponible sur genotoul ? Et sur internet ?
- 3) Lancez “Cufflinks transcript assembly and FPKM (RPKM) estimates for RNA-Seq data “ sur chaque condition avec les options suivantes :
 - -g pour faire un assemblage RABT
 - max-intron-length : 5000

Max Intron Length

Use Reference Annotation

Reference Annotation

Gene annotation dataset in GTF or GFF3 format.

4) Utilisez cuffmerge pour fusionner les 2 nouveaux gtf.

Cuffmerge merge together several Cufflinks assemblies (Galaxy Version 2.2.1.0)

GTF file(s) produced by Cufflinks

- 76: MT_Cufflinks on data 5 and data 71: assembled transcripts
- 68: WT_Cufflinks on data 5 and data 36: assembled transcripts
- 5: ITAG_pre2.3_gene_models_Ch6.gtf

- Combien de transcrits obtenez-vous ? Comparer ce résultat au comptage de l'exercice 6.
- L'outil "Cuffcompare compare assembled transcripts to a reference annotation and track Cufflinks transcripts across multiple experiments " permet d'obtenir une comparaison entre deux fichiers d'annotation.

Cuffcompare compare assembled transcripts to a reference annotation and track Cufflinks transcripts across multiple experiments (Galaxy Version 2.2.1.0) Options

GTF file(s) produced by Cufflinks

- 85: Cuffcompare on data 34, data 5, and others: combined transcripts
- 78: Cuffmerge on data 76 and data 68: merged transcripts
- 76: MT_Cufflinks on data 5 and data 71: assembled transcripts
- 68: WT_Cufflinks on data 5 and data 36: assembled transcripts
- 5: ITAG_pre2.3_gene_models_Ch6.gtf

Additional GTF Inputs (Lists)

+ Insert Additional GTF Inputs (Lists)

Use Reference Annotation

Yes

Reference Annotation

5: ITAG_pre2.3_gene_models_Ch6.gtf
Requires an annotation file in GFF3 or GTF format.

Ignore reference transcripts that are not overlapped by any input transfrags

Yes No
consider only the reference transcripts that overlap any of the input transfrags (Sn correction)

Ignore input transcripts that are not overlapped by any reference transcripts

Yes No
consider only the input transcripts that overlap any of the reference transcripts (Sp correction). Warning: this will discard all 'novel' loci!

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

34: ITAG2.3_genomic_Ch6.fasta



- Extrayez du fichier tmap les lignes dont la troisième colonne n'est pas '=' et allez voir pour chaque type de transfrag un exemple issu du nouveau gtf (merged.gtf) dans IGV, puis retournez voir les zones citées dans l'exercice 5.

Filter data on any column using simple expressions (Galaxy Version 1.1.0)

Filter

Dataset missing? See TIP below.

With following condition

Double equal signs, ==, must be used as shown above. To filter for an arbitrary string, use the Select tool.

Number of header lines to skip

Execute

Dans IGV, pensez à étendre la piste du gtf (view « extended »).

Exercice n°5 : mesure d'expression au niveau gènes/transcripts

1) Manipulation du GTF, se familiariser avec sa référence :

À partir du fichier ITAG_pre2.3_gene_models_Ch6.gtf, comptez combien il y a de transcripts ?

Count occurrences of each record (Galaxy Version 1.0.2)

from dataset

Dataset missing? See TIP below

Count occurrences of values in column(s)

Select/Unselect all

Multi-select list - hold the appropriate key while clicking to select multiple

Delimited by

How should the results be sorted?

1	2
3	gene_id "Solyc06g005000.2.1"; transcript_id "Solyc06g005000.2.1";
1	gene_id "Solyc06g005010.1.1"; transcript_id "Solyc06g005010.1.1";
6	gene_id "Solyc06g005020.1.1"; transcript_id "Solyc06g005020.1.1";
1	gene_id "Solyc06g005030.1.1"; transcript_id "Solyc06g005030.1.1";
1	gene_id "Solyc06g005040.1.1"; transcript_id "Solyc06g005040.1.1";
3	gene_id "Solyc06g005050.2.1"; transcript_id "Solyc06g005050.2.1";
2	gene_id "Solyc06g005060.2.1"; transcript_id "Solyc06g005060.2.1";
2	gene_id "Solyc06g005070.1.1"; transcript_id "Solyc06g005070.1.1";
24	gene_id "Solyc06g005080.2.1"; transcript_id "Solyc06g005080.2.1";
2	gene_id "Solyc06g005090.2.1"; transcript_id "Solyc06g005090.2.1";
8	gene_id "Solyc06g005100.2.1"; transcript_id "Solyc06g005100.2.1";
4	gene_id "Solyc06g005110.2.1"; transcript_id "Solyc06g005110.2.1";
9	gene_id "Solyc06g005120.1.1"; transcript_id "Solyc06g005120.1.1";

2) Quantification brute avec FeatureCount

- a) Lancez feature count avec le GTF de référence sur l'ensemble des échantillons (utilisez les paramètres définis dans le tutoriel)
- b) Testez également le comptage au niveau des gènes avec cet utilitaire.

Exercice 6 : Sauvegarder votre chaîne de traitements en convertissant votre historique en workflow

Créez un workflow à partir des traitements bioinformatiques précédemment réalisés :

Etape 1 : Depuis le menu « Options » de votre historique en cours, choisir « Extract Workflow ».

Un workflow est ainsi généré automatiquement et disponible depuis le menu « Workflow ».

Etape 2 : Exportez ensuite ce workflow en cliquant sur la flèche noire à côté de l'intitulé du workflow et choisir « Download or Export ».

Etape 3 : Download to File : Veuillez cliquer sur « Download workflow to file so that it can be saved or imported into another Galaxy server.»

Puis enregistrez ce fichier sur votre PC. Il vous sera ensuite possible de le ré-importer dans votre instance Galaxy.

Création de workflow :

- A partir de rien : Menu « Workflow » puis « Create a new workflow »
- A partir d'un historique : « History panel » Options → « Extract workflow »



Comme pour les historiques, il est possible de partager des workflows.