

Formation à l'analyse de données RNA-seq Galaxy

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 :



Galaxy is an open, web-based platform for data intensive biomedical research. Whether on this free public server or [your own instance](#), you can perform, reproduce, and share complete analyses.



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



TopHat is a fast splice junction mapper for RNA-Seq reads. It aligns RNA-Seq reads to mammalian-sized genomes using the ultra high-throughput short read aligner [Bowtie](#), and then analyzes the mapping results to identify splice junctions between exons. <http://tophat.cbcb.umd.edu/>



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



SAM (Sequence Alignment/Map) format is a generic format for storing large nucleotide sequence alignments. <http://samtools.sourceforge.net/>



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/>



Objetctifs:

Cette formation a pour but de vous aider à traiter les séquences issues des SGS (Seconde Generation Sequencing) en particulier les plates-formes Illumina (GAIIx, 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.

Pré-requis: savoir utiliser un environnement Galaxy.



Pour réaliser l'ensemble de ces exercices, connectez-vous avec votre utilisateur genotoul <http://sigenae-workbench.toulouse.inra.fr/> depuis un navigateur.

Les données que nous utiliserons sont accessible à cette adresse : http://genoweb.toulouse.inra.fr/~formation/4_Galaxy_RNAseq/

Vous pouvez également utiliser un des comptes formation : anemone aster bleuet iris muguet narcissse pensee rose tulipe violette

Exercice n°1: Data Management

Quelques liens:

- EMBL-ENA (European Nucleotide Archive) : <http://www.ebi.ac.uk/ena/>

Étude des données publiques disponibles à EMBL ENA:

- données correspondantes aux runs ERR022486 et ERR022488
- sujet de l'étude :
ERR022486 : RNA from Zebrafish 1day, embryo (31 504 560 reads)
ERR022488 : RNA from Zebrafish 3 days, embryo (24 920 613 reads)
- type de séquenceur utilisé : Illumina Genome Analyzer II
- type de librairie de séquençage (protocole) utilisé :
PAIRED

Abstract: **Paired-end sequence** data has been generated using **polyA selected RNA** from a range of zebrafish tissues and developmental stages using the **Illumina Genome Analyzer**. These data have been used to improve the gene annotation of the zebrafish genome. Study description: Zebrafish total RNA was extracted from embryonic and adult tissue, then polyA selected. After fragmentation and reverse transcription Illumina sequencing libraries were prepared. Paired-end sequence runs were performed with 36, 37, 54 and **76 base reads** on the Illumina Genome Analyzer.



- Récupérez les données re-formatées pour l'étude du chromosome 22 : utilisez par le « Shared data » / « Published histories », importer l'historique « Data RNAseq » de dlabourdette.

Name	Annotation	Owner	Community Rating	Community Tags	Last Updated ↓
Data RNAseq		dlabourdette	★★★★★		3 minutes ago
TEST VARIANT EF		rodriguez	★★★★★		Oct 15, 2013
workflow sRNAseq 2		pervenche	★★★★★		Oct 14, 2013
SNP DATA		orue	★★★★★		Oct 11, 2013
Human Genome		smaman	★★★★★		Jun 04, 2013
Coding exons & SNPs		smaman	★★★★★		Jun 04, 2013

Galaxy History ' Data RNAseq'

Dataset	Annotation
1: ERR022486_read1	None
2: ERR022486_read2	None
3: ERR022488_read1	None
4: ERR022488_read2	None
8: Danio rerio_chr22.Zv9.62.gtf	

About this History

Author
dlabourdette

Related Histories
All published histories
Published histories by dlabourdette

Rating
Community (0 ratings, 0.0 average) ★★★★★
Yours ★★★★★

Tags
Community: none
Yours:

- Visualisez le contenu de chacun des fichiers pour vérifier leur chargement.

```

1: /ERR022488_chr22_read1.fastq
2: http://genoweb/~formation/galaxy_rnaseq/data/fastq/ERR022486_chr22_read2.fastq
    
```




- Lancez tophat (sur 4 CPU) en paired-end avec une taille d'insert de 200bp et une taille maximale d'intron de 5000bp pour les jeux de données ERR022486, ERR022488 contre la référence nommée « Danio rerio Zv9 62 chr 22 »

*** Tophat for Illumina (version 1.0.0)**

Your RNA-Seq FASTQ file (read 1):

Your RNA-Seq FASTQ file (read 2):

Select a reference genome:

Number of threads used to align reads:

Maximum intron length:

Expected (mean) inner distance between mate pairs:

Vous obtenez en résultat 3 fichiers (pouvant être renommés) :

- Fichier de jonction (bed)
- Fichier d'alignement (bam)
- Fichier des reads non alignées (unmapped.bam)
- Utilisez « samtools flagstat » sur le fichier bam de chaque alignement, pour obtenir un résumé des statistiques d'alignement.
- Indexez le fichier bam avec samtools (samtools index) pour pouvoir ensuite le visualiser avec IGV sur votre ordinateur.
- Téléchargez sur votre ordinateur les fichiers de résultats de tophat (bam et bed) et le fichier d'indexation (bai)
- Renommez ces fichiers ERR022488.bam, ERR022488.bed, ERR022488.bam.bai

Visualisation des résultats :

- Utilisez IGV pour visualiser les résultats sur votre poste de travail.
- Lancez IGV depuis « download » du site web de la formation (en bas de la page):
<http://www.broadinstitute.org/software/igv/download>
 Le génome zebrafish est déjà intégré dans IGV
- Vous pouvez également charger les annotations correspondant au chromosome 22 en tant que fichier (disponible à http://genoweb.toulouse.inra.fr/~formation/4_Galaxy_RNAseq/data/reference/Danio_rerio_chr22.Zv9.62.gtf.gz).
- Chargez les .bam, .bed
- Explorez l'interface, utiliser le clic-droit (pour visualiser toutes les isoformes, les paires de lectures associées....)

**Pour info :**

Lors d'une analyse future, il se pourra que les noms des chromosomes soient différents entre ceux intégrés dans IGV et ceux de votre gtf

Dans ce cas, il faut utiliser un fichier de correspondance ; par exemple :

Zv9_alias.tab à mettre dans le répertoire :

/home/...../igv/genomes

<http://www.broadinstitute.org/software/igv/LoadData/#aliasfile>

Ce fichier devra alors contenir les correspondances entre les noms utilisés par IGV et les noms de votre GTF :

```
1   chr1
2   chr2
3   chr3
4   chr4
5   chr5...
```

....





Exercice n°3: mesure d'expression brute au niveau gène/transcripts :

Manipulation du GTF, se familiariser avec sa référence :

- A partir du fichier référence du chromosome 22 :
Danio_rerio_chr22.Zv9.62.gtf
- Combien y a-t-il de gènes (voir outil d'analyse de GTF)?
- Combien y a-t-il de transcrits ?

*** Ensemble GTF statistics (version 1.0.0)**

Your GTF file on which you would count genes:

17: Danio_rerio_chr22.Zv9.62.gtf

Execute

This tool count uniq line or column.

Your GTF contains :

Number of genes : 1276

Number of transcripts : 2133

Quantification au niveau gènes à l'aide du gtf de référence et Htseq-count :

- triez les alignement à l'aide de samtools sort selon les read name (pré-requis pour htseq-count sur des données paired-ends) (pensez à cocher yes)

*** Samtools sort (version 1.0.0)**

Your accepted hits bam file:

12: {ERR022488_read1}-Tophat.bam

Sort by read names rather than by chromosomal coordinates:

Yes

Use this option if you want to sort by read names

Execute

This tool sort alignments by leftmost coordinates. File out.prefix.bam will be created.

Command : samtools sort -n in.bam out.prefix

OPTION: -n Sort by read names rather than by chromosomal coordinates



- lancez htseq count sachant que les données ne sont pas strand-spécifique et que l'on veut les intersections de gène non vide, on utilisera comme attribut le gene_id (par défaut c'est le transcript_id).

*** htseq (version 1.0.0)**

Your accepted hits bam file:

Your gtf or gff file:

Use this option if you want to skip all reads with alignment quality lower than the given minimum value (default: 0):

Use this option to feature type (3rd column in GFF file) to be used, all features of other type are ignored:

GFF attribute to be used as feature ID (default,suitable for Ensembl GTF files: gene_id):

gene_id

Select whether the data is from a strand-specific assay. Specify 'yes', 'no', or 'reverse' (default: yes). 'reverse' means 'yes' with reversed strand interpretation:

Select mode to handle reads overlapping more than one feature(choices: union, intersection-strict, intersection-nonempty; default: union):

- afin d'avoir un fichier de comptage complet, utilisez l'outil « Merge Htseq count» qui concatène les résultats de chaque colonne pour chaque ligne.

*** Merge Htseq count (version 1.0.0)**

Your first htseq count file:

First htseq count file name:

Datasets

Dataset 1

Other htseq count files:

htseq count file name:



Quantification au niveau transcrits à l'aide du gtf de référence et sigcufflinks :

- lancez sigcufflinks avec le GTF de référence sur chaque échantillon séparément (sans recherche de nouveaux transcrits).

*** Sigcufflinks (version 1.0.0)**

Your accepted hits bam file:

Your gtf file:

G or g ?:

- Afin d'avoir un fichier de comptage complet, utilisez l'outil « Merge sigcufflinks » qui concatène les résultats de chaque colonne pour chaque ligne.

*** Merge sigcufflinks (version 1.0.0)**

Select a reference genome (if your genome of interest is not listed, please contact Siganae Team):

Your annotation file (defined feature to count):

Your first raw transcripts tsv file from sigcufflinks:

Your first raw transcripts tsv file name:

Datasets

Dataset 1

Other raw transcripts tsv file from sigcufflinks:

Other raw transcripts tsv file name:

Pour info :

Quelques outils pour l'analyse d'expression différentielle:

Il existe toute une batterie de package Bioconductor disponible pour l'analyse différentielle de données RNA-seq. Ces outils sont en plein développement et ne sont pas encore mature (difficulté dans le choix de la méthode à appliquer, choix de normalisation, évolution rapide des versions avec de fort changements d'une version à l'autre souvent...)

1. DESeq: <http://bioconductor.org/packages/release/bioc/html/DESeq.html>
2. EdgeR <http://www.bioconductor.org/packages/2.11/bioc/html/edgeR.html>



Exercice 4 : Recherche de nouveaux transcrits :

- Créez un fichier bam contenant tous les alignements de tous les échantillons (avec samtools merge)

*** Samtools merge (version 1.0.0)**

Your first accepted bam file:

Datasets

Dataset 1

Other accepted bam file:

- lancez cufflinks avec le GTF de référence, avec le fichier bam complet (afin d'obtenir un gtf complet correspondant à nos échantillons) en prenant l'option recherche de nouveaux transcrits (use ref as guide).

Cufflinks (version 0.0.5)

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 for low abundance transcripts.

Use Reference Annotation:

Reference Annotation:

Gene annotation dataset in GTF or GFF3 format.

Perform Bias Correction:

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

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

Mean Inner Distance between Mate Pairs:

Standard Deviation for Inner Distance between Mate Pairs:



- lancez sigcufflinks avec le GTF issu de l'étape précédente (assembled transcripts), chaque échantillon séparément (sans rechercher de nouveaux transcrits : option quantitate against ref).

*** Sigcufflinks (version 1.0.0)**

Your accepted hits bam file:

Your gtf file:

G or g ?:

- lancez les comptages bruts comme précédemment.

En fin de formation:

Veillez, en fin de TP, nettoyer votre compte de formation ("Delete permanently" de l'ensemble des "histories" créés).

Pour répondre à vos questions à propos de Galaxy:

* Mail : sigenae-support@listes.inra.fr

* Une FAQ et un manuel utilisateur sont disponibles depuis la page d'accueil de l'instance Sigenae de Galaxy.

* Les formations de la plateforme BIOINFO GENOTOUL sont disponibles sur <http://sig-learning.toulouse.inra.fr>