

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 :



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/>



STAR is a Spliced Transcripts Alignment to a Reference.

<https://github.com/alexdobin/STAR>



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/>

Objectifs:

Cette formation a pour but de vous aider à traiter les séquences issues des SGS (Seconde 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.

Pré-requis: savoir utiliser un environnement Unix.



Pour réaliser l'ensemble de ces exercices, connectez-vous sur votre **compte « genotoul »** en utilisant « putty » depuis un poste windows ou la commande ssh depuis un poste linux.

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

Pour les traitements « lourds » utilisez le cluster avec la commande « **qlogin** » ou « **qrsh** ».

Sur genotoul, créer, dans votre répertoire work, un répertoire de travail : **tp_rnaseq**.

Exercice n°1: Data Quality

- Récupérer les données re-formatées pour l'étude du chromosome 6 de la Tomate depuis la page web de la formation: <http://bioinfo.genotoul.fr/index.php?id=119>.
MT_rep1_1_Ch6.fastq.gz
MT_rep1_2_Ch6.fastq.gz
WT_rep1_1_Ch6.fastq.gz
WT_rep1_2_Ch6.fastq.gz



Vous pouvez télécharger les fichiers fastq directement sur votre compte « genotoul » en utilisant la commande « wget » depuis genotoul (en copiant l'adresse du lien et coller), penser à vous placer dans le répertoire correspondant sur genotoul.

Analyse de la qualité des données (en se connectant sur un noeud):

- Utilisation de FastQC sans visualisation directe java (afin d'éviter la surcharge du serveur !)
- Lancer FastQC pour analyser la qualité des lectures sans visualisation directe, mais en passant par la création de fichiers résultats :
/usr/local/bioinfo/src/fastqc/current/fastqc MT_rep1_1_Ch6.fastq.gz
vous pouvez utiliser l'option -nogroup pour ne pas moyenner dans le module base_sequence_quality
- Analyser les résultats (depuis genotoul) :
 - Aller dans le répertoire xxx_fastqc et ouvrir avec firefox la page html générée :
firefox fastqc_report.html &
 - OU visualiser les images contenues dans le dossier Images :
display per_base_quality.png &

OU

- Sinon pour une utilisation avec visualisation java sur votre poste local, il faut installer FastQC sur votre poste
- Télécharger et installer FastQC :
http://www.bioinformatics.bbsrc.ac.uk/projects/fastqc/fastqc_v0.10.0.zip
instructions disponibles ici :
- <http://www.bioinformatics.bbsrc.ac.uk/projects/download.html#fastqc>
- Charger les fastq (il est possible de charger plusieurs fichiers à la fois)
- Vous pouvez sauvegarder un rapport pour chaque analyse

- Quelle est la longueur des lectures ?
- Quelle est la qualité du séquençage ?
- Regarder les résultats concernant les biais décrits lors du cours, lesquels retrouve-t-on ?

Exercice n°2: alignement/visualisation

Quelques liens:

- Tophat: <http://tophat.cbcb.umd.edu/>
- Samtools: <http://samtools.sourceforge.net/>
- Bowtie: <http://bowtie-bio.sourceforge.net/index.shtml>
- FTP download de Ensembl: <http://www.ensembl.org/info/data/ftp/index.html>
- STAR: <https://github.com/alexdobin/STAR>

Aujourd'hui nous allons nous focaliser sur l'alignement sans transcriptome de référence avec les paramètres de base. Pour lancer l'alignement il vous faut une référence.

- Créer un répertoire bowtie2-index.
 - Depuis la page de la formation sur le web, récupérer la séquence du chromosome 6 (ITAG2.3_genomic_Ch6.fasta) sur genotoul.
 - Générer l'index bowtie2 des séquences fasta (suivre les indications suivantes)
syntaxe : bowtie2-build [options]* <reference_in> <bt2_index_base>
- a) Saisir bowtie2-build sans paramètres pour obtenir l'aide.
 - b) Lancer la commande sur le fichier fasta précédemment téléchargé en spécifiant comme nom d'index 'bowtie2-index/tomato_chr6'.
 - c) Lister le contenu du répertoire bowtie2-index. A quoi correspondent les fichiers *.bt2 ?



Sur le serveur genotoul les génomes sont déjà indexés pour vous dans /bank/bowtie2db/. Vous pouvez directement les utiliser pour réaliser l'alignement.

- Réaliser les alignements épissés (suivre les indications suivantes):
syntaxe : tophat [options] <bowtie_index> <reads1.1[,reads2.1,...]>\[reads1.2[,reads2.2,...]]

a) Saisir tophat pour obtenir l'aide du logiciel d'alignement.

b) Quelle version de tophat est utilisé ? (tophat -v)
Quelle est la dernière version de tophat disponible sur internet ?
Quelle est la version la plus récente disponible sur genotoul ? (lister le répertoire /usr/local/bioinfo/src/tophat/)

b) Lancer 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 WT et MT contre la l'index nommé 'tomato_chr6', nommer respectivement le répertoire de sortie aln_tophat_wt et aln_tophat_mt



Rappel :

*Pour lancer une commande sur le cluster en réservant 4 CPU utiliser la commande :
qsub -N job_name -pe parallel_smp 4 -b Y 'ma commande'*

*Pour vérifier l'avancement des calculs utiliser la commande :
qstat -u nom_utilisateur*

c) Vérifier que votre job tourne sur le cluster et est lancé sur 4 CPU (qstat)

En option, si vous souhaitez pendant que les calculs tournent, réaliser un alignement STAR suivre les indications suivantes, sinon passer cette section grisée.

Voir le manuel : <https://github.com/alexdobin/STAR/blob/master/doc/STARmanual.pdf>

Penser à vérifier la version du logiciel. (STAR --version)

Créer et se déplacer dans un repertoire nommé star.

Indexer la référence :

mkdir star-index

STAR --runMode genomeGenerate --genomeDir star-index --genomeFastaFiles
../ITAG2.3_genomic_Ch6.fasta

nb. Si vous avez le GTF de référence il est recommandé de l'utiliser pour le TP nous ne l'utiliserons pas.

Aligner le WT contre l'index créé :

STAR --runThreadN 4 --genomeDir star-index --readFilesIn
../WT_rep1_1_Ch6.fastq.gz ../WT_rep1_2_Ch6.fastq.gz --readFilesCommand zcat
--alignIntronMax 5000 --outFileNamePrefix ./aln_star_wt --outSAMtype BAM
SortedByCoordinate

Important si votre librairie n'est pas brin spécifique et que vous souhaitez faire de la

découverte de nouveaux transcrits avec cufflinks, il faut ajouter l'option – outSAMstrandField intronMotif

- Visualiser le contenu du fichier align_summary.txt dans chacun des répertoires de sortie.
- Utiliser « samtools flagstat » sur les fichiers accepted_hits.bam pour obtenir le nombre de read alignés.
Syntaxe : samtools flagstat <in.bam>
- Quelle sont les différences entre le fichier align_summary.txt et ces résultats ?
- Indexer le fichier bam avec samtools (samtools index) pour pouvoir ensuite le visualiser avec IGV sur votre ordinateur.
- Télécharger sur votre ordinateur les fichiers de résultats de tophat (bam et junctions.bed) et le fichier d'indexation (bai)
- En local, renommer ces fichiers WT.bam, WT.bed, WT.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>
- Chargez les annotations (fichier gtf mis à disposition dans
<http://genoweb.toulouse.inra.fr/~formation/LigneCmd/RNAseq/data/reference/>)
- Chargez les .bam, .bed
- Explorez l'interface, utiliser le clic-droit (pour visualiser toutes les isoformes, les paires de lectures associées....)
- Regardez les régions montrées dans le cours ainsi que les régions suivantes :
SL2.40ch06:2,786,806-2,807,064
SL2.40ch06:22,595-27,402
SL2.40ch06:38,480,842-38,484,938
SL2.40ch06:10,694,176-10,704,838
SL2.40ch06:9,839,693-9,862,815
Solyc06g009140.2.1

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

Manipulation du GTF, se familiariser avec sa référence :

- Depuis la page de la formation sur le web, récupérer le gtf ne contenant que le chromosome 6 : ITAG_pre2.3_gene_models_Ch6.gtf

- A partir de ce fichier :
 - Combien y a-t-il de gènes ?(utiliser cut sur colonne 9, cut selon « ; », sort -u et wc)?
 - Combien y a-t-il de transcrits ?

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

- Saisir featureCounts sans options pour obtenir l'aide et la version.
- Quelle est la version disponible sur genotoul et la dernière version disponible sur internet ?
- lancer featureCount (sur un nœud du cluster) sur les deux bam en sachant que l'on compte les reads :
 - s'alignant sur les exons,
 - à regrouper par gène,
 - un read peut être assigné plusieurs fois,
 - la librairie n'est pas brin spécifique,
 - on ne veut compter que les alignements primaires et les reads mappés de façon unique,
 - les reads ayant un alignement avec une qualité minimum de 20 (attention la qualité STAR n'est pas standard)
 - le nombre minimum de bases chevauchantes doit être > 10,
 - on ne compte que les fragments (les paires),
 - la distance entre deux reads doit être en 60 et 600 bp
 - on ne compte que les fragments dont chacun des reads est correctement alignés.

Exercice 4 : Recherche de nouveaux transcrits :

- Créer un répertoire 'cufflinks' pour l'analyse par cufflinks de l'ensemble du jeu de données.
- Fusionner les alignements obtenu par échantillons dans un seul fichier.
Syntaxe : samtools merge merge.bam fichier1.bam fichier2.bam ..
- Que signifie RABT ? A quoi sert l'option -g du cufflinks?
- Quelle version de cufflinks est disponible sur genotoul ? Et sur internet ?
- Lancer cufflinks en utilisant le fichier bam fusionné (afin d'obtenir un gtf complet correspondant à nos échantillons) avec les options suivante :
 - -g pour faire un assemblage RABT
 - library-type : fr-unstranded
 - max-intron-length : 5000
 - si vous souhaitez paralléliser utiliser l'option -p
- Combien de transcrits obtenez vous ? Comparer ce résultat au comptage de l'exercice 3.

- L'outil cuffcompare permet d'obtenir une comparaison entre deux fichiers d'annotation.
Syntaxe : cuffcompare -r reference.gft cufflink1.gtf cufflink2.gtf ...

Extrayez du fichier tmap, les lignes dont la troisième colonne n'est pas '=' et allez voir pour chaque type de transfrag un exemple dans IGV, puis retournez voir les zones citées dans l'exercice 2.

- Lancer a nouveau featureCounts avec ce nouveau transcriptome de référence.

Exercice 5 : Un pas vers les statistiques (en option)

Toutes les informations sur l'étape de biostatistique sont disponibles en bas de la page suivante : <http://bioinfo.genotoul.fr/index.php?id=119>

- constituer la matrice attendue par le script R :

```
#gene_id  mt  wt
Solyc06g005000.2.1    240  72
Solyc06g005010.1.1     0   0
Solyc06g005020.1.1     1   0
Solyc06g005030.1.1     0   0
Solyc06g005040.1.1     0   0
Solyc06g005050.2.1    20   5
```

- Utiliser le script /usr/local/bioinfo/Scripts/bin/Normalization.R sur la matrice de comptage issue de featureCount.

```
/usr/local/bioinfo/src/R/R-3.2.2/bin/Rscript /usr/local/bioinfo/Scripts/bin/Normalization.R -f
tomato_count.R.txt -o norm
```

Copier en local le répertoire de sortie pour visualiser les images.

- Utiliser le script d'expression différentielle /usr/local/bioinfo/Scripts/bin/DEG.R sur un des résultats de la normalisation

```
/usr/local/bioinfo/src/R/R-3.2.2/bin/Rscript /usr/local/bioinfo/Scripts/bin/DEG.R -f
tomato_count.R.txt --norm norm/RLE_info.txt --pool1 mt --pool2 wt -o DEG
```

Copier en local le répertoire de sortie pour visualiser les images.