

Encapsulation d'un programme dans l'interface GALAXY

Projet n°7 – Assemblage de lecture en contig et identification de la lecture majoritaire

Objectif : Le But de ce projet est d'intégrer à l'interface galaxy un outil permettant à partir d'un fichier d'alignement BAM de recenser les contigs et d'en déterminer la lecture majoritaire.

Dans un contexte de séquençage d'ARN non codant, l'utilité d'un tel outil serait de déterminer avec le plus de précisions les loci de ces portions d'ARN. En effet la lecture majoritaire dans un contig est la séquence qui a été lue lors du séquençage, c'est donc celle qui représente le mieux et avec le plus de sensibilité la portion du génome qui a été exprimée en tant que ncRNA.

Le traitement des données se fera en deux étapes. Une première qui consistera à aligner les reads sur le génome de référence et de procéder à un nettoyage jusqu'à l'obtention du fichier au format BAM. La deuxième étape sera la recherche des contigs et des loci des lectures majoritaires.

Etape 1 – alignement

Données disponibles :

un fichier de lecture de séquençage au format fasta
http://genoweb.toulouse.inra.fr/~gaspin/cleaned_sequences.fa
un génome de référence au format fasta, compressé en .gz
<http://genoweb.toulouse.inra.fr/~gaspin/genome.fa.gz>

Nous allons utiliser l'aligner Bowtie2, et les fonctions du programme Samtools afin de filtrer et trier les lectures alignées. Les programmes cités ci-dessus seront appelés via le script shell

Le pipeline utilisé pour l'alignement est le suivant :

- Décompression du fasta du génome de référence
- Indexation du génome de référence (Bowtie2-build)
- Alignement des lectures du fastq sur le génome indexé, récupération des reads uniquement alignés (Bowtie2-align, option -k 1)
- Filtrage des reads alignés (Samtools view) sur critère de score d'alignement
- Trie des reads alignés et filtrés (Samtools sort)

En sortie de ce pipeline nous avons un fichier d'alignement au format BAM que nous pouvons utiliser dans la deuxième étape du processus. Le script fasta2BAM.sh appelle les commandes bowtie et samtools nécessaire pour cet alignement, son code est adapté pour être lancé en ligne de commande sur la plateforme genotoul, cependant il peut être facilement adapté en modifiant les chemins d'accès vers les programmes. Notons que nous n'avons pas intégré ce script à galaxy, considérant que notre outil intégré à Galaxy prend en entrée un fichier BAM.

Etape 2 – Generation des loci et identification de la sequence majoritaire

Une fois l'alignement effectué, la tâche consiste à détecter les contigs et déterminer pour chaque contig la lecture majoritaire. Un contig est un ensemble de lecture (read) qui se chevauchent lors de l'alignement et forment ensemble une séquence continue. Nous définissons la lecture majoritaire comme la plus longue séquence commune avec une occurrence maximale entre tous les reads d'un même contig. Ainsi si sur un même contig on retrouve plusieurs séquences dont le nombre d'occurrence est égal et maximal, alors la plus longue sera considérée comme lecture majoritaire. Une fois les positions début et fin des loci des contigs et de leur lecture majoritaire déterminée, nous les stockons sous forme de tableau dans un fichier tsv ou tabular

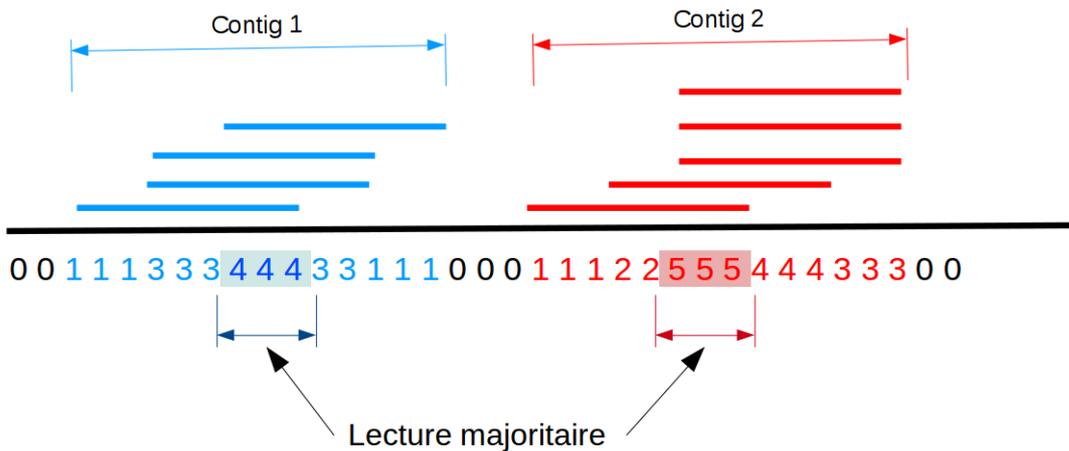


Figure 1 : Explication graphique de la définition du contig et du read.

Pour ce faire, nous décidons de passer par une première étape de pré-traitement du BAM afin d'obtenir la profondeur de séquençage. Utilisation de samtools depth pour obtenir un fichier qui nous donne pour chaque position (paire de base) sur le génome de référence le nombre de reads qui ont mappé sur cette position, chaque ligne correspond à :
 « ID_chromosome position nbr_reads ».

C'est à partir de ce fichier que nous allons pouvoir déterminer les contigs et les lectures majoritaires selon l'algorithme suivant :

Pour chaque ligne:

Initialisation du contig pour la première ligne du fichier.

Si Position précédente +1 == Position actuelle :

Si dans locus :

Si Nombre de read actuel > Meilleur nombre de read :

Définition nouveau meilleur locus.

Si Nombre de read actuel < Meilleur nombre de read :

Fin du locus

Sinon:

Si Nombre de read actuel > Meilleur nombre de read :

Définition nouveau meilleur locus.

Si Nombre de read actuel == Meilleur nombre de read :

Taille du locus +1

Si Taille du locus > Taille du meilleur locus :

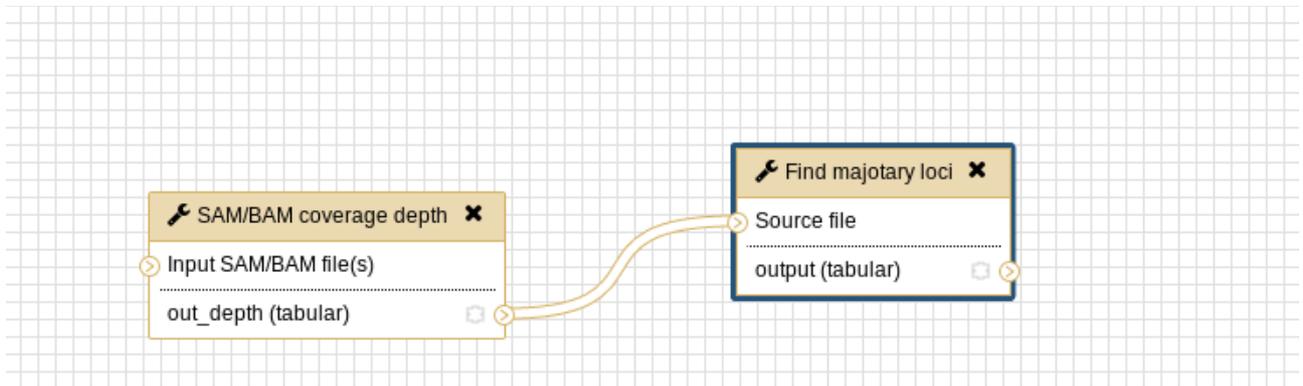
Définition nouveau meilleur locus.

Sinon :

Ecrire dans fichier output le contig

Encapsulation de l'outil :

Pour encapsuler cet outils dans galaxy il faut avoir deux fichier, le premier est notre outil python, le deuxième est notre XML, permettant d'afficher et d'encapsuler dans une iinterface galaxy notre outil. Pour que l'outil fonctionne il faut d'abord faire passer notre fichier BAM par l'outil samtools depth (disponible sur le toolshed ainsi que dans le fichier tools associé à ce Rapport). Les deux outils sont indépendants, mais notre outil est basé sur l'architecture d'un fichier depth. Il faudra l'utiliser comme montré ci dessous :



Discussion :

En théorie cet outil devrait fonctionner sur tous les types de données issues de séquençage à haut débit. Nous pouvons très bien imaginer son application dans un contexte de ChIP-seq pour déterminer avec le plus d'acuité le locus lié par une protéine ou la portion de séquence cible reconnu par un facteur trans-acting. Il devient ensuite trivial de récupérer la séquence avec l'identifiant du chromosome et les positions start et stop sur le genome de reference indexé grâce par exemple à la commande samtool faidx <refgenome.fa> <chr:start-stop>.

Nous supposons cependant que passer par le coverage n'est certainement pas la manière la plus direct d'arriver à nos fins et qu'il aurait été possible de traiter directement le BAM, grâce à des commande incluse dans bedtools ou samtools comme mpileup. Traiter directement le BAM nous aurait certainement permis de récupérer l'orientation du brin (+/-), information que nous perdons en passant par la profondeur.

Cependant nous pouvons adresser comme critique à notre programme le faite de ne pas ponderer les différents candidats pour la lecture principale par leur taille et de les choisir de façon trop naïve. En effet, si 12 reads se chevauchent sur 3 paires de bases, cette séquence va devenir la lecture majoritaire, potentiellement devant une sequence de 20 nucléotides où 11 reads se chevauchent. Pour pallier à ce problème, une solution serait de pondérer les candidats à la lecture majoritaire par la taille de la séquence candidate car selon l'utilisation à venir du résultats, il serait peut être plus interessant d'avoir une lecture majoritaire d'une taille plus grande mais moins représentée plutot qu'une lecture majoritaire de 3 paires de bases.