

Rapport de Stage

CHARRIAT Florian

29 mars 2018

Matériels et Methodes

1. Matériels

2. Methodes

2.1 Assemblage

Les données de séquençage étant des données short read, l'outil ABySS (Assembly By Short Sequences) [4] à été utilisé pour réaliser les assemblages. L'assemblage avec cet outil se réalise en deux étapes. Une première étape consiste à générer tous les k-mères de longueur donnée à partir des reads. Cette ensemble de k-mères est ensuite utilisée pour construire les contigs initiaux. Puis la deuxième étape utilise les informations de pair end pour étendre les contigs et résoudre les ambiguïtés de chevauchement des contigs.

Pour obtenir le meilleur assemblage possible, chaque souche est assemblée 8 fois avec une taille de k-mère différents. Les tailles de k-mère utilisées vont de 20 à 90 avec un pas de 10. Le meilleur assemblage de chaque souche est ensuite sélectionné après visualisation des données de qualité.

Une fois la sélection réalisée, le script formatFastaName est utilisé. Pour chaque assemblage sélectionné, le script élimine, dans un premier temps, les scaffolds de taille inférieure à 500 pb qui risquent de poser problème pour l'annotation. Puis les scaffolds seront numérotés en fonction de leur longueur afin d'homogénéiser la nomenclature des scaffolds. Le plus grand scaffold sera alors nommé scaffold_1.

2.2 Éléments répétés

Une fois l'assemblage réalisé, il est nécessaire de masquer les éléments répétés des assemblages. En effet, les éléments répétés peuvent poser problèmes lors d'annotation automatique. Ils peuvent être confondus avec des gènes codant pour des protéines. Mais peuvent aussi perturber la structure des modèles de gène, en s'insérant dans les introns par exemple. L'outil repeatMasker est donc utilisé pour masquer les éléments répétés des assemblages sélectionnés. Cet outil compare les séquences provenant d'une base de données d'éléments répétés avec ceux provenant de l'assemblage. Il va ensuite masquer ces éléments trouvés en les remplaçant par des N. L'outil repeatMasker permet aussi d'obtenir une annotation détaillée des répétitions présentes dans l'assemblage. La base de données utilisée pour repeatMasker est Repbase auquel a été ajouté les éléments répétés découverts par l'équipe à l'aide de l'outil RepeatModeler.

Une fois les éléments masqués, les scaffolds ne comportant que des éléments répétés sont éliminés des assemblages. En effet, les outils utilisés pour l'annotation automatique n'arrivent pas à annoter les assemblages comportant des scaffolds uniquement composés de N.

2.3 Annotation automatique

M.ozzae étant un eukaryote, l'outil Braker a été utilisée pour l'annotation automatique. Cet outil combine les avantages des outils GeneMark-ET et AUGUSTUS.

Tout d'abord, GeneMark-ET génère des structures géniques. Pour cela l'outil utilise un ensemble de paramètres du hidden semi-Markov model (HSMM) défini initialement pour prédire des régions codantes dans l'assemblage.

Une fois ces gènes prédites, un sous ensemble d'exon et intron est utilisé pour re-estimer les paramètres du HSMM. L'intégration des régions codantes dans le set d'entraînement nécessite que l'exon prédit possède au moins un site d'épissage. Les sites d'épissage sont ceux prédits indépendamment par les deux méthodes, l'ab initio et l'alignement des données RNA-Seq. Les étapes de prédiction des régions codantes et de la ré-estimation des paramètres du HSMM sont réalisées de manière itérative tant que les régions codantes prédites varient.

Deuxièmement, AUGUSTUS utilise les gènes prédits par GeneMarker-ET pour l'entraînement, puis intègre les informations de mapping des données RNA-seq dans les prédictions de gènes finaux.