



LABORATOIRE DE BIOINFORMATIQUE POUR LA GÉNOMIQUE ET LA BIODIVERSITÉ

Master de bioinformatique - ingénierie de plate-forme en biologie
UNIVERSITÉ PARIS CITÉ

Rapport d'alternance

Gestion informatique des données de séquençage

31 août 2022

William Amory
sous la responsabilité de Frédérick Gavory



Table des matières

Glossaire	1
1 Introduction	3
1.1 Le LBGB au sein du Genoscope et du CEA	3
1.2 Contexte et missions du LBGB	3
1.3 Présentation du workflow NGS	4
1.4 La technologie MGI	4
2 Objectifs de ma mission	5
3 Matériels et Méthodes	6
3.1 Le cluster de calcul et Slurm	6
3.2 La base de données de référence NGL et la gestion des projets	6
3.3 Le langage de programmation Perl	6
3.4 Logiciels de démultiplexage et génération de fichiers de séquences (bcl2fastq - bcl-convert)	7
3.5 Les pipelines de génération de fichiers de séquences pour les technologies Illumina et Nanopore	7
3.6 Les pipelines de contrôle qualité des lots de séquences pour les technologies Illumina et Nanopore	8
4 Résultats	9
4.1 Etude comparative des logiciels bcl2fastq et bcl-convert	9
4.1.1 Détermination des meilleurs paramètres pour bcl2fastq	9
4.1.2 Comparaison entre bcl2fastq et bcl-convert	10
4.1.3 Préparation de la migration de bcl2fastq vers bcl-convert	11
4.2 Le pipeline de génération de fichiers de séquences pour la technologie MGI	11
5 Discussions et perspectives	18
5.1 perspectives du workflow NGS pour la technologie MGI	18
5.1.1 Améliorations futures du pipeline NGS_RG pour la technologie MGI	18
5.1.2 Développement du pipeline de contrôle qualité pour le technologie MGI	18
5.2 Evaluation d'outils de contrôle qualité	19
Notes	20
Références	21
6 Annexes	21

Glossaire

BGI : *Beijing Genomics Institute*, est une entreprise Chinoise de biotechnologie fondé en 1999.

CEA : Commissariat à l'Énergie Atomique et aux Énergies Alternatives

CNRGH : Centre National de Recherche en Génomique Humaine

CNS : Centre National de Séquençage (Genoscope)

CPU : *Central Processing Unit* (Unité Central de Traitement)

DRF : Direction de la Recherche Fondamentale

ERGA *European Reference Genome Atlas*

IBFJ : Institut de Biologie François Jacob

Illumina : Entreprise Californienne de biotechnologie fondée en 1998, qui réalise : R&D, production et vente d'instruments de séquençage d'ADN à haut débit et très haut débit, ainsi que des logiciels et services d'analyses bio-informatique des données de séquençage.

Jira : Logiciel de gestion de projet, de suivi d'incidents et de bugs développé par l'entreprise Atlassian

LBGB : Laboratoire de Bioinformatique pour la Génomique et la Biodiversité

Lims : *Laboratory Information Management System*

MGI : Filiale du groupe BGI fondée en 2016 dont les missions sont : R&D, production et vente d'instruments de séquençage d'ADN, de réactifs et de produits connexes

NCBI : *National Center for Biotechnology Information*, est un institut national des Etats Unis d'Amérique pour l'information biologique moléculaire. Il développe notamment la base de données de génomes GenBank et la base de données des publications PubMed

NGL : *Next Generation LIMS*

NGS : *Next Generation Sequencing*

Oxford Nanopore : Entreprise Anglaise de biotechnologie fondée en 2005, qui développe et produit des systèmes de séquençage, basé sur les propriétés diélectriques de ces dernières.

PacBio : *Pacific Biosciences of California* est une entreprise Californienne fondée en 2004, qui développe et produit des systèmes de séquençage en temps réel à molécule unique (SMRT) d'ADN

Path : Chemin d'accès à un fichier ou à un répertoire dans le système de fichier

Perl : *Practical Extraction and Report Language*

Ram : *Random Access Memory* (Accès Mémoire Aléatoire, aussi appelé mémoire vive)

Slurm : *Simple Linux Utility for Resource Management* qui est un logiciel open source d'ordonnancement des tâches informatiques

1 Introduction

1.1 Le LBGB au sein du Genoscope et du CEA

Le Genoscope (CNS) a été créé en 1996 pour participer au projet mondial de séquençage du génome humain (*Human Genome Project*) qui a débuté en 1990 et s'est terminé en 2003. Il a notamment participé au séquençage du chromosome 14. Le Genoscope est impliqué dans le développement de programme de génomique en France dans le cadre du projet France génomique. Aujourd'hui les projets phares du Genoscope sont les projets **Tara** (*Pacific*, Océans, *Artic* ...), qui ont pour objectifs l'étude des écosystèmes marins ; Le projet **ERGA**, dont l'objectif est de créer une base de données de références de haute qualité des génomes d'espèces européennes.

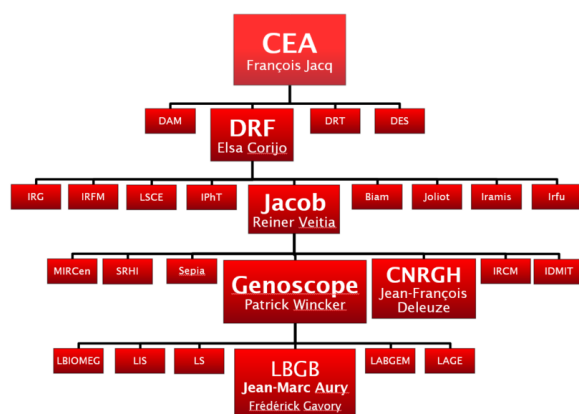


FIGURE 1 – Organigramme situant l'équipe du LBGB au sein du Genoscope et du CEA

Le Laboratoire de Bioinformatique pour la Génomique et la Biodiversité (**LBGB**) dirigé par Jean-Marc Aury, fait partie du Genoscope qui est une composante de l'institut de biologie François Jacob (**IBFJ**) de la direction de la recherche fondamentale (**DRF**) du Commissariat à l'Énergie Atomique et aux Énergies Alternatives (**CEA**), qui a été fondé le 18 octobre 1945 par Charles de Gaulle. L'intégration du Genoscope au CEA a été réalisée en 2007, et en 2017 il devient une composante de l'IBFJ.

1.2 Contexte et missions du LBGB

Les missions qui sont confiées au LBGB sont de réaliser le contrôle qualité des données de séquences issues des différents séquenceurs, d'effectuer l'assemblage¹ des séquences et l'annotation² des génomes, dans l'objectif de mettre à disposition des laboratoires collaborateurs les données avec un premier niveau de valorisation. Le laboratoire est divisé en plusieurs groupes de travail. Le groupe « production » (dont je fais partie), le groupe « assemblage », le groupe « annotation » et le groupe « évaluation des technologies de séquençage ».

Les missions du groupe de « production » sont : de tester des logiciels tiers, ainsi que développer et maintenir des scripts utilisant ces logiciels pour automatiser la prise en charge des données en sortie de séquenceur. Cette prise en charge peut répondre à une demande de la production et des laboratoires du Genoscope et du CNRGH, mais aussi pour des laboratoires extérieurs. L'objectif principal est la mise en place et le main-

tient de pipelines automatisant l'ensemble. Le groupe s'appuie sur un travail de veille et d'évaluation technologique pour chacune de ses missions.

1.3 Présentation du workflow NGS

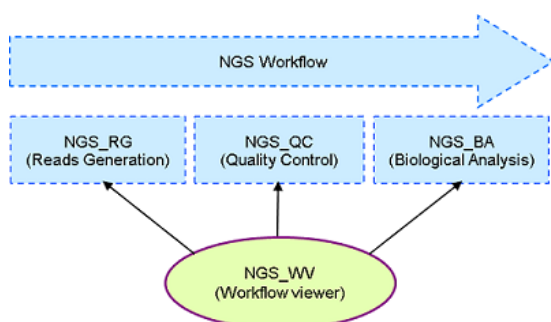


FIGURE 2 – Workflow de génération, de contrôle qualité et d'analyse biologique des fastq

Le workflow NGS est composé de trois pipelines pour les technologies Illumina et Oxford Nanopore. Le premier (*ngs_rg*³), permet la génération des reads⁴ et des fichiers de séquences correspondants aux échantillons. Le second (*ngs_qc*⁵), permet de réaliser leur contrôle qualité. Le dernier (*ngs_ba*⁶), permet de faire les analyses biologiques inter-échantillons (*readset*)⁷.

Ces trois pipelines sont automatisés dans le workflow et permettent de réaliser la distribution des données de séquençage dans des répertoires dédiés, triées par projet, échantillon, runs⁸ et technologie de séquençage. Ils réalisent aussi le nettoyage, l'analyse de ces fichiers et mettent à jour la base de données de référence NGL. Les trois pipelines du workflow NGS sont monitorés par *NGS Workflow Viewer* (NGS_WV), qui est une application web permettant de surveiller l'avancement des pipelines pour les runs pris en charge par le NGS-workflow.

1.4 La technologie MGI

Le genoscope et le CNRGH ont récemment fait l'acquisition de séquenceurs MGI (2 DNBSEQ-G400 et 1 DNBSEQ-T7).

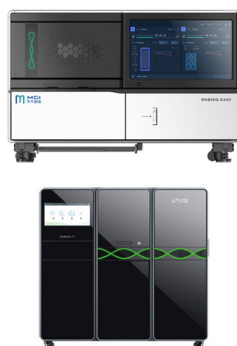


FIGURE 3 – Sequenceurs DNBSEQ-G400 (en haut) et DNBSEQ-T7 (en bas) de MGI
<https://en.mgi-tech.com/products/>

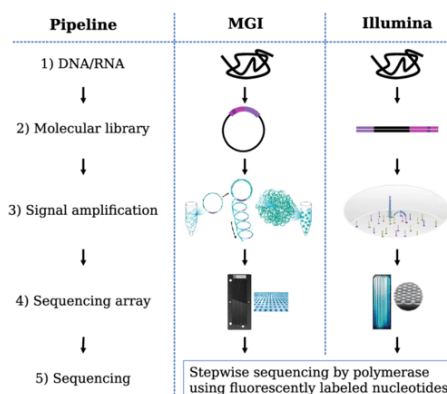


FIGURE 4 – Différences entre Illumina et MGI de technologie NGS

Il s'agit de séquenceurs à haut débit et très haut débit, dont les principales différences entre MGI et Illumina sont dans la création des librairies⁹ et la méthode d'amplification d'ADN. Les librairies sont double brins circulaire pour MGI, alors que pour Illumina elle est double brins linéaire. L'amplification ADN est réalisée en solution et forme des DNB (*DNA-nanoballs*¹⁰), puis déposée sur la Flowcell¹¹ pour MGI, alors que pour Illumina elle est réalisée après immobilisation sur les Flowcell.

Sequencers specifications				
	MGI		Illumina	
	DNBSEQ-G400	DNBSEQ-T7	HiSeq 4000	NovaSeq 6000
Max Number of Flow Cells	2	4	2	2
Max Lane/Flow Cell	4	1	4	4
Run Time	~ 14-37 h	~ 20-30 h	~ 24-84 h	~ 13-44 h
Data ouput/Run	0.27-1.4 Tb	1-6 Tb	0.9-1.8 Tb	1-6 Tb
Max Reads/Run	1.8 billions	5 billions	10 billions	20 billions
Max Read Length	2 × 200 bp	2 × 150 bp	2 × 150 bp	2 × 250 bp

TABLE 1 – Spécification des séquenceurs

2 Objectifs de ma mission

L'objectif principal de ma mission est la mise en place d'un workflow NGS pour les séquenceurs de MGI. Plus précisément il s'agira de créer un pipeline de génération de fichiers de séquences (`ngs_rg_mgi`¹²) et un pour le contrôle qualité de ces fichiers (`ngs_qc_mgi`¹³). Le workflow devra créer et mettre à jour l'état des runs, des *lanes*¹⁴ et de *readset*¹⁵ dans NGL, réaliser le contrôle qualité des fichiers de séquences, au format fastq, obtenus après démultiplexage¹⁶ des runs. Il devra mettre à jour l'avancement du traitement d'un run dans NGL, en y insérant les statistiques obtenues lors du démultiplexage, les résultats des contrôles qualités, etc. Puisque l'objectif est d'obtenir un premier niveau de valorisation des fichiers de séquences, permettant aux autres groupes (« assemblage », « anotation ») de prendre en charge ces fichiers avant de les mettre à disposition des laboratoires collaborateurs.

Je dois également, rechercher et réaliser des évaluations de nouveaux outils pour les différents pipelines des différentes technologies de séquençage. En vue d'un potentiel ajout ou de remplacement d'outils. Il sera donc nécessaire de maintenir les pipelines des différentes technologies de séquençage en conséquence. Par exemple l'évaluation de logiciels de trimming (`Cutadapt`, `Trimmomatic`) en vue d'un remplacement du logiciel `fastx_clean` de

la suite FASTX Toolkit¹⁷ qui est un outil mono-coeur pour un outil multi-coeurs. Ou bien trouver et évaluer un logiciel d'assignation taxonomique plus performant que le logiciel Centrifuge utilisé actuellement.

3 Matériels et Méthodes

3.1 Le cluster de calcul et Slurm

Le Genoscope possède un cluster (*inti*) de calcul de 71 noeuds répartis sur 5 partitions. La partition « normal » est composé de 47 noeuds qui disposent entre 12 et 36 coeurs et entre 96 et 386 Go de Ram. La partition « small » est composée de 8 noeuds dont 4 qui possèdent 8 coeurs et 64 Go de Ram, et 4 autres qui disposent de 16 coeurs et 128 Go de Ram. Cette partition est utilisée pour les processus courts et/ou peu de mémoire. Les partitions « xlarge » et « xxlarge » ont deux noeuds composés de 48 et 2To de Ram, et de 56 coeurs et 6To de Ram respectivement. Ces deux partitions sont utilisées pour les processus demandant plusieurs jours ou semaines de calculs. La partition « production » du cluster *inti* est composé 12 noeuds qui disposent de 16 coeurs et de 257 Go de Ram. L'accès à l'utilisation du cluster et de ses noeuds est réalisé par le logiciel [Slurm](#).

3.2 La base de données de référence NGL et la gestion des projets

Le Genoscope dispose de sa propre base de données de référence NGL. Celle-ci est divisée en plusieurs parties. NGL_BI¹⁸, est la partie de la base de données utilisée par les équipes de bioinformatique. NGL_SEQ¹⁹, est la partie de la base de données utilisée dès la réception des échantillons et jusqu'au séquençage de ces derniers. Il y a également les parties NGL_SUB²⁰, NGL_REAGENT²¹ et NGL_PROJECTS²². La gestion et le suivi du développement informatique sont réalisés par le système de tickets [Jira](#).

3.3 Le langage de programmation Perl

L'écriture du workflow des pipelines pour les séquenceurs MGI est réalisée dans le langage de programmation Perl. L'utilisation de ce langage est rendu nécessaire pour des raisons historiques du laboratoire, puisque de nombreuses librairies et modules qui ont été utilisés dans l'écriture des pipelines sont écrits en Perl.

C'est pour toutes ces raisons qu'il m'a été nécessaire d'apprendre à coder en Perl. j'ai donc commencé par réaliser un programme permettant de faire des analyses statistiques élémentaires sur des fichiers fastq, tel que le taux de GC, la moyenne du score de la qualité, ainsi que plusieurs autres métriques. Le programme est capable de gérer les fichiers fastq issue de séquençage *single end*²³ et *paired end*²⁴. Cela m'a permis de prendre en main les

librairies Perl utilisées pour les différents pipelines déjà en place. Ainsi que de m’habituer à l’environnement de travail, l’utilisation du lancement de job sur les noeuds de calculs et l’utilisation des modules²⁵ pour les différents pipelines.

3.4 Logiciels de démultiplexage et génération de fichiers de séquences (bcl2fastq - bcl-convert)

Ces deux logiciels de *Base Calling* (bcl2fastq et bcl-convert), sont tous deux développés et commercialisés par Illumina. Cette évaluation entre ces deux logiciels est nécessaire pour déterminer les changements qu’il y aura à faire dans les pipelines de génération de fichiers de séquences pour les technologies Illumina, en vue du remplacement de bcl2fastq (qui sera bientôt obsolète) par bcl-convert.

Dans un premier temps, il est nécessaire de déterminer les conditions optimales de bcl2fastq (temps total (*Elapsed time*²⁶), temps CPU (*CPU time*²⁷), pourcentage d’utilisation CPU (*%CPU*²⁸)) en fonction des ressources disponibles sur les noeuds du cluster (*inti*) réservé à la *production*, pour pouvoir comparer les performances des 2 logiciels. Les conditions optimales sont déterminées en fonction des paramètres suivants de bcl2fastq (l’équivalent de bcl-convert est indiqué entre crochets) :

- **r** [bcl-num-decompression-threads] : nombre de *threads*²⁹ accordé pour la décompression et la lecture des *Bases Calls*³⁰
- **p** [bcl-num-conversion-threads] : conversion des *Bases Calls* en fastq
- **w** [bcl-num-compression-threads] : écriture et compression des fichiers fastq

Tous ces tests sont réalisés sur le même noeud de calcul, dans l’objectif de minimiser les biais. La comparaison est effectuée sur le temps total du démultiplexage, ainsi que sur le temps CPU et le pourcentage d’utilisation des CPU.

3.5 Les pipelines de génération de fichiers de séquences pour les technologies Illumina et Nanopore

Les pipelines de générations de fichiers de séquences pour les technologies Illumina et Nanopore réalisent dans un premier temps le démultiplexage permettant la création des fichiers de séquences correspondant aux échantillons et des fichiers de statistiques de ces derniers. Ils créent les runs, les pistes, et les readset dans NGL_BI en y insérant les metriques, graphiques et fichiers permettant leurs évaluations.

Concernant le pipeline de génération de fichiers de séquences pour la technologie MGI, il s’agira de développer un pipeline similaire à celui d’Illumina en prenant en compte que

le démultiplexage est directement réalisé par les séquenceurs. Les métriques, graphiques et fichiers de statistiques sont également différents d'Illumina. Il sera donc nécessaire de trouver comment obtenir les métriques, graphiques et fichiers, ou de les calculer, à partir des données générées lors du séquenceur permettant de les insérer dans NGL_BI

3.6 Les pipelines de contrôle qualité des lots de séquences pour les technologies Illumina et Nanopore

Les pipelines de contrôle qualité des lots de séquences réalisent différentes étapes de contrôle qualité et de nettoyage des lots de séquences. Il réalise le contrôle qualité et l'estimation de duplicat des séquences des fichiers avant et après nettoyage (*trimming*), il retire le *PhiX*³¹ (pour les technologies Illumina), réalise l'assignation taxonomique des séquences, réalise un alignement des séquences si un génome de référence existe, réalise le calcul du pourcentage de séquences qui ont leurs reads *forward* (brin sens) et *reverse* (brin anti-sens) qui se chevauchent et réalise la distribution des fichiers de séquences nettoyés dans leurs répertoires de projet, d'échantillon, de type de technologie et de run.

Concernant le pipeline de contrôle qualité des fichiers de séquences pour la technologie MGI, qui est en cours de développement. Il s'agit de développer un pipeline similaire à celui d'Illumina en prenant en compte qu'avec cette technologie il n'y a pas de *PhiX* à enlever dans les fichiers de séquences.

4 Résultats

4.1 Etude comparative des logiciels bcl2fastq et bcl-convert

4.1.1 Détermination des meilleurs paramètres pour bcl2fastq

Après avoir effectué différentes combinaisons des paramètres, il a été mis en évidence que la variation du paramètre r et w en fixant le paramètre p , n'apportait pas de différences significatives pour le temps total d'exécution, le temps cpu ou le pourcentage d'utilisation cpu, comme on peut l'observer sur la figure 5, pour p fixé à 12. Des résultats similaires ont été obtenus pour p égale à 4, 8 et 16.

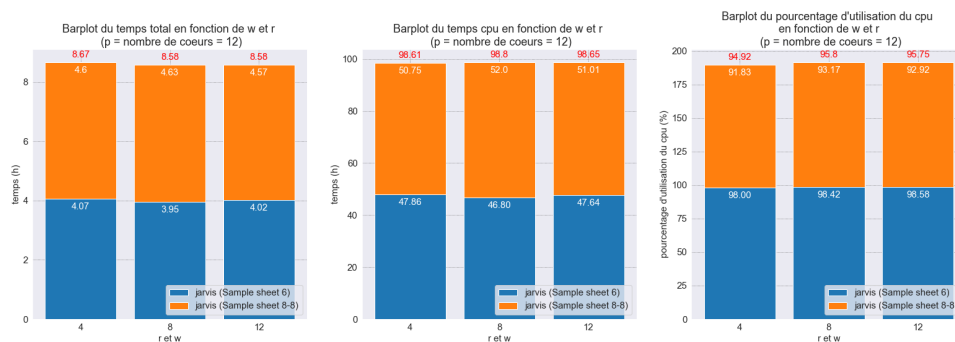


FIGURE 5 – Digrammes en bâtons du temps total d'exécution (à gauche), temps cpu (au milieu) et du pourcentage d'utilisation des cpu (à droite) en fonction des paramètres r et w

Il y a deux *sample sheet*³², car le nombre de bases considérées des *reads index* entre les *lanes* est différent, obligeant à réaliser deux appels différents au logiciel pour générer les fastq et le démultiplexage. Ci-dessous, la figure 6, représente les résultats obtenus en faisant varier p et en fixant les paramètres r et w à 4 (ces deux paramètres sont fixés à 4 pour pouvoir comparer les 4 résultats). On observe que plus on augmente le nombre de coeurs pour p , plus l'exécution est rapide. On observe que le temps cpu augmente bien avec le nombre de coeurs et que le pourcentage d'utilisation des cpu est optimal ($> 90\%$).

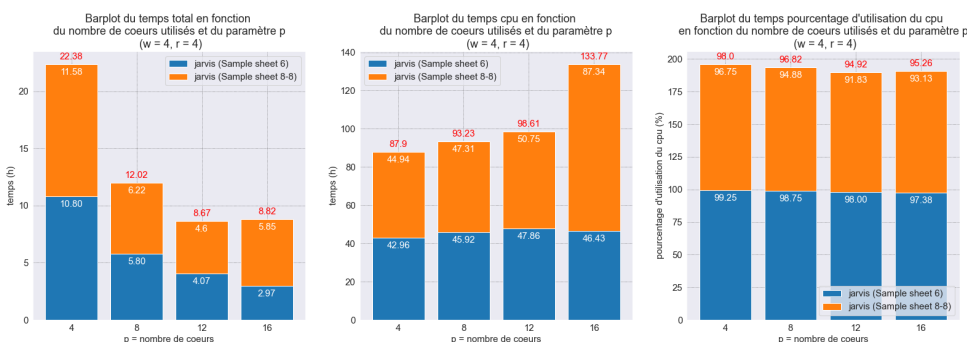


FIGURE 6 – Digrammes en bâtons du temps total d'exécution (à gauche), temps cpu (au milieu) et du pourcentage d'utilisation des cpu (à droite) en fonction du paramètre p

Au vue des résultats obtenus nous avons décidé que les meilleurs paramètres étaient de fixer p à 12, puisque le gain apporté en augmentant à 16 est faible. Néanmoins nous le conserverons pour réaliser la comparaison avec bcl-convert, tout comme p fixé à 8, car il nous permettrait de réaliser deux générations de fastq et de démultiplexage en simultané sur un seul noeud de calcul.

4.1.2 Comparaison entre bcl2fastq et bcl-convert

J'ai donc fait varier les paramètres p , r et w de manière à ce que chacun des paramètres soient égale au nombre de coeurs accordés aux deux logiciels. On observe bien, sur la figure 7, que plus on augmente le nombre de coeurs pour chacun des logiciels (et donc le nombre de *threads* pour p , r et w) plus la génération des fastq et le démultiplexage est rapide. De plus on remarque que bcl-convert permet de réduire le temps d'environ 1/3 par rapport à bcl2fastq.

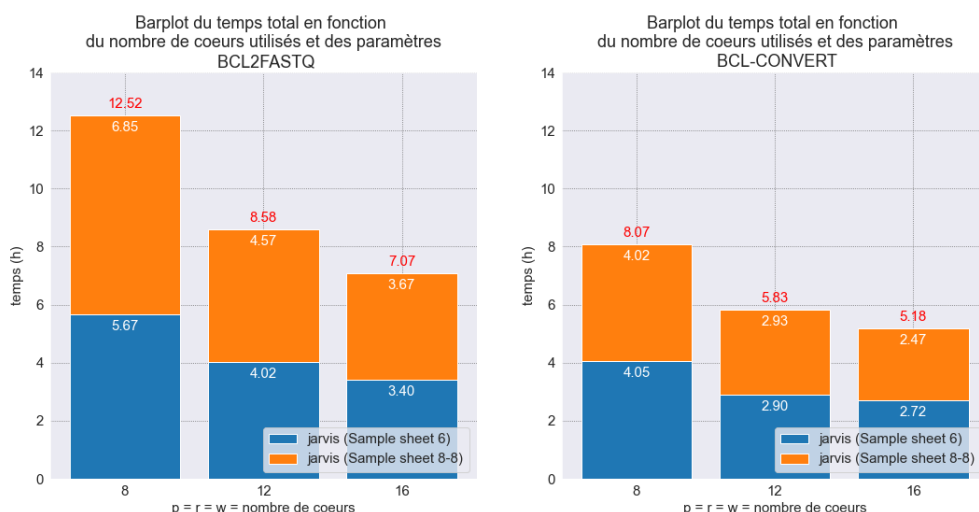


FIGURE 7 – Temps total de génération des fastq pour bcl2fastq et bcl-convert

J'ai également échangé avec le service technique d'Illumina à propos des fichiers de sortie et de l'arborescence de des derniers en utilisant bcl-convert. En effet il s'avère que l'arborescence et les fichiers de sortie sont très différents entre les deux logiciels. Ces échanges avaient pour objectif de savoir s'il on pouvait obtenir une arboresnce similaire à bcl2fastq, pour minimiser l'impact du changement de logiciel sur les pipelines. Le changement de bcl2fstq, qui sera bientôt obsolète, par bcl-convert va nous obliger à réaliser de gros changements dans tous les pipelines qui utilisent ces fichiers de sortie et va demander aussi au laboratoire de séquençage de s'adapter à la nouvelle *sample sheet* de bcl-convert.

4.1.3 Préparation de la migration de bcl2fastq vers bcl-convert

Le logiciel bcl-convert est plus rapide d'environ 1/3 par rapport à bcl2fastq. Sachant également que ce dernier sera bientôt obsolète et que le nombre de coeurs disponibles par noeuds pour la partition « *production* » du cluster de calcul est de 16 coeurs. Nous avons décidé d'attribuer l'intégralité des coeurs d'un noeud de « *production* », c'est à dire 16 coeurs. L'intégralité des changement entre les deux logiciels a été consignés dans un cahier des charges. Il contient, la commande à lancer pour réaliser le *Base Calling*, les modules à charger dans l'environnement, le chemin relatif des fichiers de sorties et leur description, ainsi qu'un exemple d'arborescence des fichiers de sorties. Ce qui permettra au développeur qui se chargera de cette migration de suivre ce cahier des charges et ainsi faciliter cette migration. Dû à la pression actuelle autour de la technologie MGI, c'est un autre développeur qui sera en charge de réaliser cette migration.

4.2 Le pipeline de génération de fichiers de séquences pour la technologie MGI

L'objectif du pipeline NGS_RG_MGI est de générer et distribuer les fichiers de séquences dans le bon répertoire de projet, d'échantillon, de type de séquençage et de run. Tout en créant et mettant à jour les runs, pistes et readsets. Notamment concernant les métriques d'évaluations des ces derniers. Le pipeline est composé de plusieurs grandes étapes.

Création et insertion des métriques du run et des pistes dans NGL

La première étape consiste à créer le run et ses pistes dans la base de données NGL, en y intégrant les métriques permettant d'évaluer le run et les pistes (figure 8). Le nom du run est constitué de la date de séquençage, du nom du séquenceur et de l'identifiant de la flowcell du run.

On y retrouve notamment, le nombre total de reads (Nb Cluster (total)), le nombre de bases totales (Nb Bases (total)) générées par le run, la taille des reads et des index (Nb Cycles). Concernant les pistes on retrouve le nombre total de bases et de reads générés sur la piste, le pourcentage de bases qui ont une qualité supérieur ou égale à Q30, Q20 et Q10. On a également le pourcentage de bases inconnus (%N), ainsi que d'autres métriques qui permettent d'évaluer le run et les pistes. Celles-ci sont détaillées plus précisément en annexes (page 21).

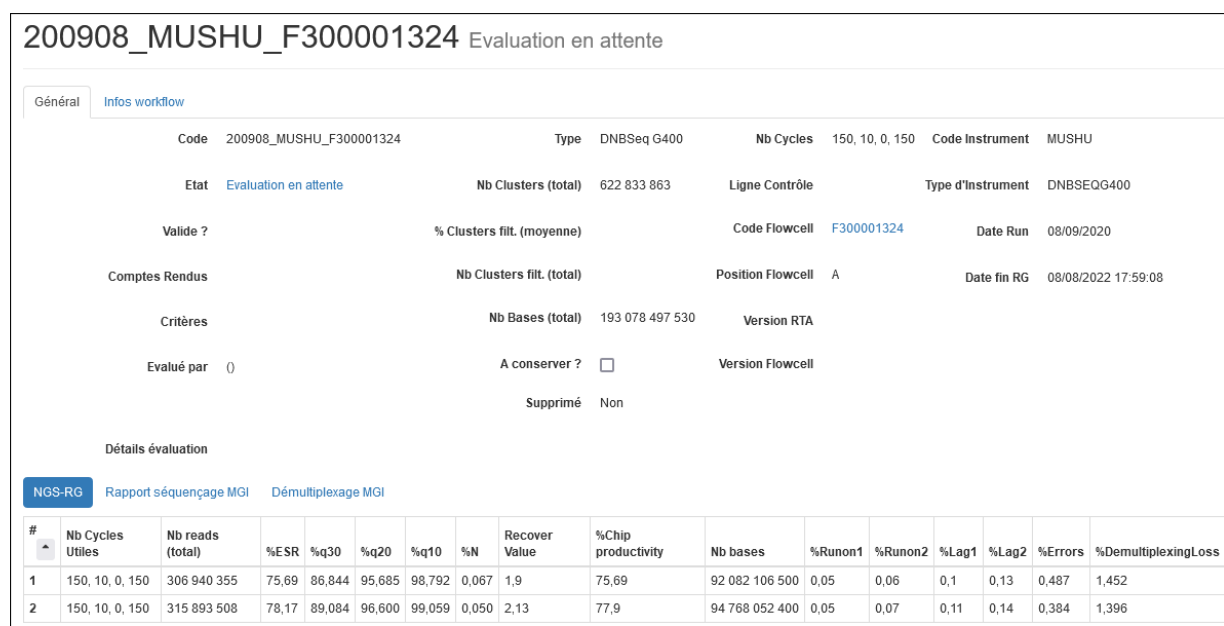


FIGURE 8 – Capture d'écran de la page du run 200908_MUSHU_F300001324 de NGL en cours de génération de fichiers de séquences (étapes d'ajout des métriques d'évaluation du run et des pistes).

Insertion des rapports de séquençage des pistes et de la listes des index dans NGL

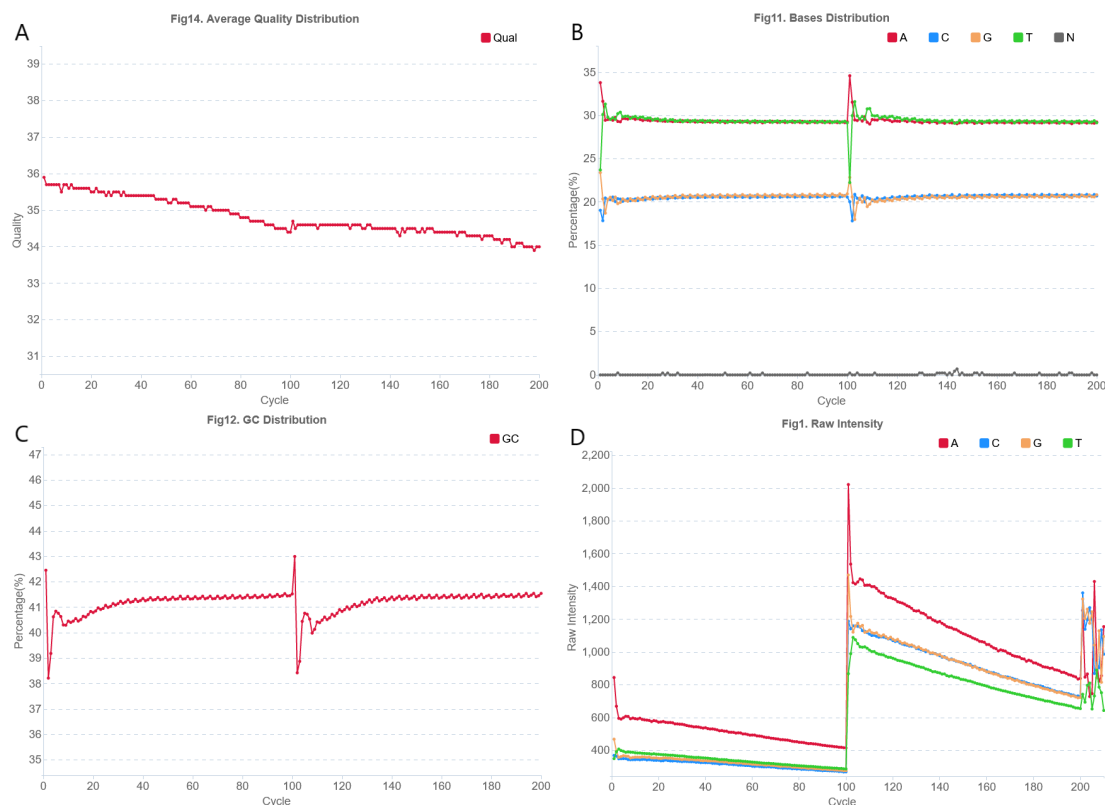


FIGURE 9 – Graphiques des distributions de la qualité moyenne (A), des bases nucléiques (B), du pourcentage de GC (C) et de l'intensité brut (D) au cours des cycles de séquençage

La seconde étape ajoute les rapports de séquençages des pistes que le séquenceur génère en fin de séquençage. Il s'agit de rapport html qui contiennent plusieurs tableaux de métriques et de graphiques permettant d'évaluer les pistes du run. Il y a notamment les graphiques de la distribution de la qualité moyenne en fonction des cycles (figure 9.A), de la distribution des bases nucléiques en fonction des cycles (figure 9.B), de la distribution du pourcentage de Guanine/Cytosine en fonction des cycles (figure 9.C), de la distribution de l'intensité brut au cours des cycles (figure 9.D). Les tableaux et graphiques de ces rapports de séquençage permettent de faciliter l'évaluation du run et de ses pistes.

Toujours dans l'optique de faciliter l'évaluation du run et de ces pistes on ajoute, la liste des index représentées à plus de 0.01% de la pistes, ainsi que les index attendus. Ces index sont triés et affichés par ordre croissant dans NGL (figure 10). Les index attendus sont colorés en vert et les index non-attendus ou inconnus sont colorés en rouge, ce qui permet de vérifier que les index attendus sont bien majoritairement représentés sur les pistes de la flowcell du run.

barcode	count	percent
barcode2	89 597 340	29,190
barcode1	84 106 886	27,402
barcode3	74 172 719	24,165
barcode4	54 607 003	17,791
GATTCGTCCT	206 151	0,067
ATCGGACTAT	181 509	0,059
GATCCGTCCT	156 796	0,051
ATTCCGTCCT	156 103	0,051
CGCAGTAAGT	148 841	0,048
ATCGACCTAT	119 597	0,039
TCAATAGGTT	114 220	0,037
CGGAGTAAGT	99 851	0,033
GGCAGTAAGT	85 114	0,028
ATGGACCTAT	83 324	0,027
ACGGACCTAT	75 840	0,025
CAATAGGTT	71 106	0,023
CGGCATAAGT	70 917	0,023
GATTCCTCCT	59 842	0,019
CGGCAGAAGT	53 283	0,017
barcode29	48 716	0,016
barcode124	37 751	0,012
CGGCGTAAGT	36 893	0,012

FIGURE 10 – Capture d'écran de la page du run 200908_MUSHU_F300001324 de NGL en cours de génération de fichiers de séquences (onglet « Démultiplexage MGI »)

Concaténation des fichiers FASTQ d'un même readset

Ensuite la troisième étape a pour objectif d'obtenir un seul fichier FASTQ par readset. En effet la technologie MGI requiert une homogénéité de dépôt entre les différents index (aussi appelé « barcode ») d'une piste. Ce qui implique qu'il est possible d'avoir plusieurs index associés à un même index, donc qu'un échantillon peut être divisé en plusieurs fractions. Le démultiplexage est directement réalisé par le séquenceur et le réalise à partir des index connus (listes d'index fournis par MGI), on obtient donc un fichier FASTQ par index. Il est impossible de préciser au séquenceur quels index sont associés à un même readset pour le démultiplexage. Cette étape est donc essentielle pour obtenir un seul fichier FASTQ par readset. Si le readset est associé à un seul readset alors on réalise une décompression du fichier FASTQ, alors que s'il est associé à plusieurs index on réalise une décompression et une concaténation des fichiers FASTQ (cf. figure 11).

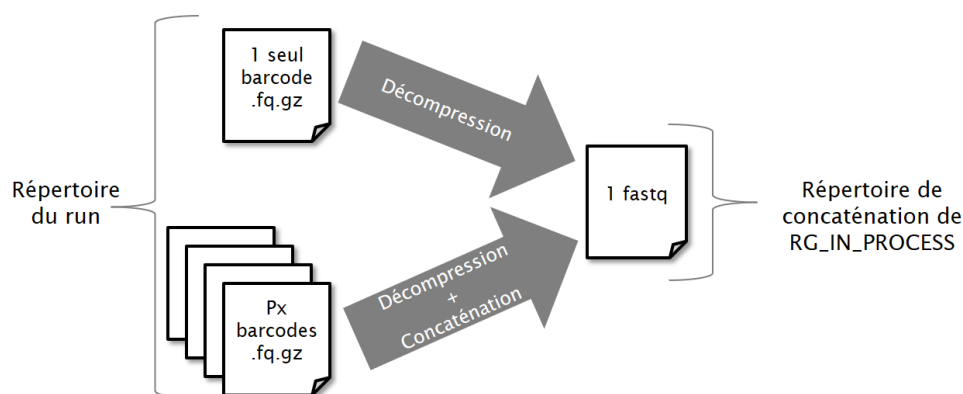


FIGURE 11 – Schéma de l'étape de « concaténation » des fichiers FASTQ d'un readset

Création et insertion des métriques des readset du run dans NGL

La quatrième étape a pour objectif de permettre l'évaluation des readsets, en les créant et en insérant les métriques d'évaluation de ces derniers dans NGL (figure 12). On y retrouve notamment le nombre de bases nucléiques et de reads du readset, ainsi que le pourcentage d'échantillon déposé sur la piste et le pourcentage de séquences valides par rapport au nombre total de séquences de la piste. On y insère également certaines métriques du run dont le readset fait partie, comme le nombre de cycles des reads et des index, la date de run, ect. Toutes ces métriques sont décrites en annexes (page 22)

Le nom du readset est constitué de l'identifiant de projet, de l'identifiant du type de banque utilisée (ADN, ARN ...), de l'identifiant d'échantillon, de l'indice de la piste, de l'identifiant de la flowcell et de l'identifiant du premier barcode.

APY_DA_AEKI_1_F300001324.MGI001 Read generation en cours Mode impression

Général **Avancé** Infos échantillon Infos workflow

Code	APY_DA_AEKI_1_F300001324.MGI001	Nb Séquences utiles	173 704 226	Run / N° Piste	200908_MUSHU_F300001324 / 1
Etat	Read generation en cours	Nb Bases utiles	52 111 267 800	Type de Run	RDNBG400
Valide QC ?	---	Valide BioInfo ?	---	Nb Cycles	150, 10, 0, 150
Comptes Rendus QC		Comptes Rendus BioInfo		Date Run	08/09/2020
Critères QC		Critères BioInfo		Date fin RG	08/08/2022 17:59:08
Évalué par	()	Évalué par	()	Date fin QC	
Détails évaluation					

NGS-RG

Nb reads	% déposé	Nb bases	% séquences valides/piste
173704226	25	52111267800	56,59

FIGURE 12 – Capture d'écran de la page du readset APY_DA_AEKI_1_F300001324.MGI001 de NGL en cours de génération de reads (étapes de création du readset et d'insertion de ces métriques d'évaluation)

On ajoute également la répartition des index au sein d'un readset (figure 13), ce qui permet de vérifier la composition en index du readset et de vérifier l'homogénéité de ces index au sein du readset. Ce nom de readset est unique, ce qui permet de déterminer rapidement et simplement à quel projet, échantillon, ect appartiennent les fichiers séquences de ce readset.

NGS-RG **Répartition des index**

Index	Nb occurrences	% de cet index dans le readset
barcode2	89 597 340	51,580
barcode1	84 106 886	48,420

FIGURE 13 – Capture d'écran de la page du readset APY_DA_AEKI_1_F300001324.MGI001 de NGL en cours de génération de reads (onglet « Répartition des index »)

Au niveaux du run un tableau référençant les readsets et leurs métriques d'évaluation est également ajouté (figure 14).

Readsets (7)

Lanes Q Voir Readsets Evaluer Readsets

N° Piste	Code	Etat	% déposé	% Séquences valides / piste	Nb Séquences valides	Nb Bases	% >= Q30	Score Qualité moyen	Valide QC ?	Valide BioInfo ?
1	APY_DA_AEKI_1_F300001324.MGI001	Read generation en cours	56,59		173 704 226	52 111 267 800	88,52	34,58	---	---
1	CRH_DA_AAAA_1_F300001324.MGI003	Read generation en cours	41,96		128 779 722	38 633 916 600	84,83	33,84	---	---
2	BAY_RA_C_2_F300001324.MGI015	Read generation en cours	23,97		75 704 787	22 711 436 100	89,56	34,82	---	---
2	BAY_RA_C_2_F300001324.MGI013	Read generation en cours	22,95		72 483 387	21 745 016 100	88,11	34,52	---	---
2	BAY_RA_C_2_F300001324.MGI014	Read generation en cours	26,77		84 552 701	25 365 810 300	89,09	34,73	---	---
2	BAY_RA_C_2_F300001324.MGI001	Read generation en cours	0,02		77 446	23 233 800	89,15	34,72	---	---
2	BSW_RA_E_2_F300001324.MGI016	Read generation en cours	24,90		78 665 053	23 599 515 900	89,90	34,90	---	---

FIGURE 14 – Capture d'écran de la page du run 200908_MUSHU_F300001324 de NGL en cours de génération de fichiers de séquences (Tableau des readset du run)

Renommage des fichiers séquences et insertion des méta-données dans NGL

La cinquième étape consiste à renommer les fichiers de séquences des readsets et d'insérer les méta-données de ces derniers dans NGL (figure 16). Le renommage des fichiers est nécessaire pour que chaque fichiers de séquences aient un nom unique et « parlant ».

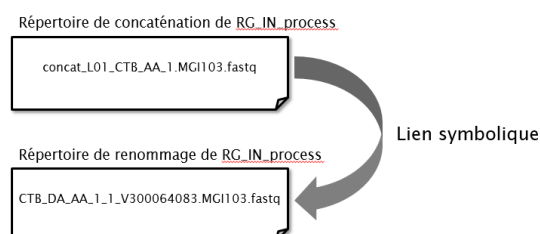


FIGURE 15 – Schéma de l'étape de renommage des fichiers FASTQ d'un readset

Le nom doit permettre d'identifier rapidement et simplement de quel projet, échantillon, flowcell, ect. appartiennent les fichiers. Le renommage des fichiers est effectués en créant un lien symbolique des fichiers obtenus à l'étape de « concaténation » dans un répertoire temporaire (cf figure 15).

Les méta-données, des fichiers de séquences du readset, qui sont insérées dans NGL permet aux utilisateurs de trouver rapidement l'emplacement de ces derniers sur le système de fichier, le type de fichier qui est disponible (brut, nettoyer) et s'ils sont utilisable. On y retrouve donc le chemin vers le repertoire de ces fichiers, leurs noms, leurs types, s'ils sont utilisable, s'il s'agit du read *forward* ou *reverse*, et le type d'encodage de la quality (Pour les séquenceurs MGI l'encodage est en ASCII 33).

APY_DA_AEKL_1_F300001324.MGI001

Read generation en cours

Général
Avancé
Infos échantillon
Infos workflow

SSID netbackup_1659981608

Date de l'archive 08/08/2022 20:08:37

Chemin fichiers utiles /env/cns/proj/projet_APY/AEKI/RunsMGI/200908_MUSHU_F300001324/

Localisation CNS

Envoyé Collaborateur ? ☐

Etat pour la soumission Pas associé à une soumission

Nom du fichier	Type de fichier	Utilisable	Label	Encodage ASCII	Clé codage md5	Nom fichier collaborateur
APY_DA_AEKL_1_1_F300001324.MGI001.fastq	RAW	Oui	READ1	33		
APY_DA_AEKL_1_2_F300001324.MGI001.fastq	RAW	Oui	READ2	33		

FIGURE 16 – Capture d'écran de la page du readset APY_DA_AEKL_1_F300001324.MGI001 de NGL en cours de génération de fichiers de séquences (Onglet « Avancé »)

Distribution des Fichiers séquences et des fichiers de statistiques

La sixième étape est de distribuer des fichiers de séquences « attendus », les fichiers de statistiques du run et les fichiers de séquences « non attendus » dans leurs répertoires dédiés.

Les fichiers attendus sont copiés vers leurs répertoires final si le séquençage a été effectué au genoscope. Sinon, si ce dernier a été effectué au CNRGH, alors les fichiers sont copiés et compressés. Il y a cette différence entre les 2 centres, car le pipeline de contrôle qualité prend en charge de fichiers compressés pour le CNRGH, contrairement à celui utilisé pour le Genoscope.

Les fichiers de statistiques du run sont archivés par pistes et par types (.html, .fq.stat) avant d'être copiés vers leurs répertoires final. Ces fichiers sont conservés dans le cas où une métrique désirés ne fait pas partie de celles insérées dans NGL ou pour tout autres problèmes qui nécessiteraient de récupérer les fichiers de statistiques du run.

Concernant les fichiers « non attendus », il s'agit des fichiers de séquences des index ne faisant pas partie d'un readset. Puisque lors du demultiplexage par les séquenceurs on obtient un fichier FASTQ par index. Ces fichier sont renommés et archivés, avant d'être distribués vers leur répertoire dédiés. Ces fichiers de séquences sont conservés dans l'éventualité d'une mauvaise déclaration d'index par les équipes de séquençages, pour pouvoir récupérer les fichiers fastq appartenant à cet index ou si l'on souhaite étudier les séquences des fichiers « non-attendus ».

Mise à jour de fin de génération de fichiers de séquence dans NGL

L'étape finale du pipeline de génération de fichiers de séquences pour la technologie MGI, est de mettre à jour le run et les readset dans l'état de « fin de génération de reads ». Cela entraine une mise à jour automatique du run à l'état « d'évaluation en attente », ce qui permet d'indiquer aux utilisateurs que le run peut être évalué. Les readset sont aussi automatiquement mis à jour vers l'état « d'attente de contrôle qualité », permettant d'indiquer au pipeline de contrôle qualité qu'il peut effectuer le contrôle qualité des readset de ce run.

5 Discussions et perspectives

5.1 perspectives du workflow NGS pour la technologie MGI

5.1.1 Améliorations futures du pipeline NGS_RG pour la technologie MGI

Le pipeline de génération de fichiers de séquences pour la technologie MGI est similaire au pipeline de la technologie Illumina. Néanmoins il n'est pas possible de comparer ces deux derniers au niveau de leurs performances du fait de leurs différences. En effet le pipeline de génération des fichiers de séquences pour la technologie Illumina, contient les étapes de *Base Calling* et de démultiplexage (Conversion des fichiers *Base Calls* en fichiers FASTQ par échantillon) qui est réalisé par le pipeline NGS_RG_ILLUMINA. À contrario, pour la technologie MGI, cette étape est directement réalisée par le séquenceur. De plus il n'est pas possible de comparer le pipeline de génération de fichiers de séquences avec des pipelines d'autres laboratoire ou outils de génération de fichiers de séquences dû fait de la spécificité du pipeline pour le Genoscope et le CNRGH. En effet l'objectif de celui-ci est de mettre à jour la base de données de référence interne au Genoscope et au CNRGH (NGL), à l'architecture de stockage des fichiers de séquences et aux noms finaux donnés à ces fichiers pour qu'ils soient uniques.

La future amélioration du pipeline NGS_RG_MGI, consistera à la mise en place d'une étape supplémentaire pour les runs qui comporterons de *mids*³³. Cette étape supplémentaire sera donc le démidage, il s'agit d'un second démultiplexage en fonction des *mids* pour la création des readsets et fichiers de séquences.

5.1.2 Développement du pipeline de contrôle qualité pour la technologie MGI

Le pipeline de contrôle qualité des fichiers de séquences pour la technologie MGI qui est en cours de développement, est constitué de deux grande étapes. La première consiste à réaliser un contrôle qualité des fichiers de séquences brut, puis dans un second temps de réaliser un contrôle qualité sur les fichiers de séquences nettoyés ainsi que d'autres traitements sur les fichiers nettoyés. Il prend en charge automatiquement les fichiers dont les readset sont dans l'état « en attente de contrôle qualité » dans NGL.

Avant de réaliser le trimming sur les fichiers de séquences brut, on réalise un échantillonnage de 20000 séquences par fichiers pour réaliser le contrôle qualité des fichiers. Cela permet d'améliorer le temps d'exécution du contrôle qualité, tout en ayant une grande représentativité de la qualité des séquences des fichiers. On réalise donc le contrôle qualité et l'estimation de dupliquas sur les échantillon des fichiers brut.

Ensuite on réalise le trimming des fichiers brut, comme il n'y a pas de *PhiX* à retirer, il s'agit du seul traitement de nettoyage des fichiers brut. On réalise l'échantillonnage de 20000 séquences par fichiers nettoyés pour réaliser leurs contrôle qualité.

Le second contrôle qualité réalisé sur les échantillons des fichiers brut, est composé d'un contrôle de la qualité des séquences, d'une assignation taxonomique des séquences des fichiers, l'estimation des duplicas (séquences retrouvés plusieurs fois dans un fichier de séquences), on réalise aussi un alignement des séquences sur un génome de références si ce derniers existe et est disponible et on calcule le pourcentage de reads qui ont le read *forward* et *reverse* qui se chevauchent.

Une fois le deuxième contrôle qualité effectué, la dernière étape du pipeline est de distribuer les fichiers de séquences nettoyés dans leur répertoire final, de rendre indisponible les fichiers de séquences brut avant de les effacer une fois celle-ci copié sur bande magnétique. Durant toutes les étapes du pipeline, on insert les métriques et graphiques qui permettront de réaliser la validation des readset ou non.

5.2 Evaluation d'outils de contrôle qualité

Les premiers outils à être évalués sont cutadapt et trimmomatic en vue d'un remplacement de fastx_clean de FASTX Toolkit. Ce dernier est un logiciels mono-coeur contrairement à cutadapt et trimmomatic qui sont multi-coeurs. Le temps d'exécution entre ces logiciels sera le critère d'évaluation le plus important, néanmoins on prendra également en compte les différents fichiers de sortie (fichiers de statistiques, fichiers de séquences qui ne passe pas les filtres données ...) pour l'évaluation et le remplacement de fastx_clean.

Un potentiel successeur au logiciels d'assignation taxonomique Centrifuge devra également être effectuer, dans l'optique d'améliorer les pipelines de contrôle qualité. L'objectif est de trouver un logiciels dont les performance son équivalentes ou meilleurs, surtout au niveau du temps d'exécution, mais également au niveau de l'assignation taxonomique des séquences et des fichiers de sortie.

Notes

- ¹Reconstruction d'un génome à partir de fragments de ce dernier
- ²Documenter le plus exhaustivement possible les informations de l'assemblage permettant de prédire la fonction d'une molécule
- ³*Next Generation Sequencing - reads generation*
- ⁴Lecture d'une séquence par un séquenceur d'un fragments d'ADN
- ⁵*Next Generation Sequencing - quality control*
- ⁶*Next Generation Sequencing - biological analysis*
- ⁷Un lot de séquences est une instance de séquences (ou reads) d'un échantillon
- ⁸Séquençage d'un ou plusieurs échantillons sur un séquenceur
- ⁹Collection de fragment d'ADN issue du génome complet d'un organisme ou plusieurs organismes (méta-génomique) et clonés dans un vecteur (le plus souvent dans des plasmides)
- ¹⁰Nanobilles d'ADN générées par la réplication de l'ADN circulaire
- ¹¹Lame d'absorption des fragments d'ADN et cuve réacteur du séquençage
- ¹²*Next Generation Sequencing - reads generation - mgi*
- ¹³*Next Generation Sequencing - quality control - mgi*
- ¹⁴pistes présentes sur la *flowcell*
- ¹⁵Lot de séquences
- ¹⁶Séparation des différents *reads* d'une *lane* en fonction de l'index d'échantillon
- ¹⁷Collection de commandes pour le traitement et l'évaluation de lot de séquences au format FASTA ou FASTQ
- ¹⁸NGL Bioinformatic
- ¹⁹NGL Sequencing
- ²⁰NGL submission (base de données des soumissions de projet (exemple la soumission d'un projet au NCBI))
- ²¹NGL reagent (base de données des réactifs)
- ²²NGL projects (base de données des projets en cours et passé)
- ²³Lecture dans un seul sens des reads par le séquenceur
- ²⁴Lecture dans les deux sens des reads par le séquenceur
- ²⁵Un module contient un ou plusieurs logiciels tiers ou développé par les équipes du genoscope. Il est nécessaire de les charger dans notre environnement de travail pour pouvoir utiliser ces logiciels.
- ²⁶Temps écoulé entre le début du programme et le fin de celui-ci
- ²⁷Temps d'utilisation des cpu par le programme
- ²⁸ $((CPU\ time + \text{temps utilisé par les appels système}) / Elapsed\ time) / \text{nombre de CPU utilisé par le programme}$
- ²⁹Processus : instructions du langage machine d'un processeur.
- ³⁰Fichier d'attribution des bases nucléiques en fonction des pics du chromatogramme lors du séquençage
- ³¹Parties du génome du phage *Lambda* qui sont ajoutés sur les pistes des flowcell avant le séquençage, permettant de contrôler le bon déroulé du séquençage.
- ³²Fichier contenant les informations et instructions pour la génération des fastq et le démultiplexage
- ³³séquence d'une dizaine de nucléotide ajouté en aval du *primer* du read *forward* permettant de réaliser un second démultiplexage

6 Annexes

Description des métriques d'évaluation d'un run et des pistes d'un run MGI dans NGL-BI

Liste et les description des métriques d'évaluation du run et des pistes (cf. figure 8 page 12) :

Nb Cycles Utiles : Nombre de cycles des reads et des index (nombre de cycles pour le read *forward*, nombre de cycles pour le premier index *forward*, nombre de cycles pour le read *reverse*, nombre de cycles pour le second index)

Nb reads (total) : Nombre de reads Total générer par la piste (S'il s'agit d'un run *pair-end* il s'agit du nombre de cluster de reads (read *forward* + read *reverse*))

%ESR : ???

%q30 : Pourcentage de bases qui ont une qualité supérieur ou égale à 30 (pour un encodage de la qualite en ASCII 33)

%q20 : Pourcentage de bases qui ont une qualité supérieur ou égale à 20 (pour un encodage de la qualite en ASCII 33)

%q10 : Pourcentage de bases qui ont une qualité supérieur ou égale à 10 (pour un encodage de la qualite en ASCII 33)

%N : Pourcentage de bases inconnus

Recover value : ???

%Chip productivity : Pourcentage de productivité de la piste (nombre de puits actif de la piste de la flowcell / nombre total de puit de la piste)

Nb bases : Nombre total de bases générés par la piste

%Runon1 : Pourcentage de read *forward* qui ont une incorporation de nucléotide d'avance par rapport au cycles en cours

%Runon2 : Pourcentage de read *reverse* qui ont une incorporation de nucléotide d'avance par rapport au cycles en cours

%Lag1 : Pourcentage de read *forward* qui ont une incorporation de nucléotide de retard par rapport au cycles en cours

%Lag2 : Pourcentage de read *reverse* qui ont une incorporation de nucléotide de retard par rapport au cycles en cours

%Errors : Pourcentage d'erreur d'incorporation de nucléotide

%DemultiplexingLoss : Pourcentage de read écartés lors du démultiplexage

Description des métriques d'un readset d'un run MGI dans NGL-BI

Liste des métriques d'évaluation des readset dans NGL-BI (cf. figure 12 page 15) :

Nb reads : Nombre de reads avant nettoyage des fichiers du readset

%déposé : Pourcentage d'échantillon déposé sur la piste de la flowcell

Nb bases : Nombre de bases avant nettoyage des fichiers séquences du readset

% séquences valides/piste : Pourcentage de séquences de la piste appartenant à ce readset (nombre total de reads du readset / nombre total de reads de la piste)

Liste des métriques d'évaluation des readsets dans le tableau qui référence tous les readsets d'un run (cf. figure 14 page 15) :

%déposé : Pourcentage d'échantillon déposé sur la piste de la flowcell

% séquences valides/piste : Pourcentage de séquences de la piste appartenant à ce readset (nombre total de reads du readset / nombre total de reads de la piste)

Nb Séquences valides : Nombre de reads du readset

Nb Bases : Nombre de bases du readset

% >= Q30 : Pourcentage de bases qui ont une qualité supérieur ou égale à 30 (pour un encodage de la qualité en ASCII 33)

Score Qualité moyen : Moyenne de la qualité des bases du readset

Autres informations à propos d'un run MGI dans NGL-BI

On retrouve également les informations permettant de suivre l'avancement du workflow NGS au niveau de l'onglet « infos workflow » (figure 17).

200908_MUSHU_F300001324 Evaluation en attente		
Général Infos workflow		
Etat	Date	Par
Nouveau	08/08/2022 17:03:38	ngsrg
Séquençage en cours	08/08/2022 17:03:38	ngsrg
Séquençage terminé	08/08/2022 17:03:38	ngsrg
Read generation en attente	08/08/2022 17:03:38	ngsrg
Read generation en cours	08/08/2022 17:03:38	ngsrg
Read generation terminée	08/08/2022 17:59:08	ngsrg
Evaluation en attente	08/08/2022 17:59:08	ngsrg

FIGURE 17 — Capture d'écran de la page du run 200908_MUSHU_F300001324 de NGL en cours de génération de fichiers de séquences (onglet « infos workflow »)

Autres informations à propos d'un readset MGI dans NGL-BI

Il y a deux autres onglets en plus de l'onglet « Général » et « Avancé ». Il s'agit de l'onglet « Infos échantillon » (figure 18) et de l'onglet « Infos workflow » (figure 19). Tout comme pour le run, l'onglet « Infos workflow » permet de suivre l'avancement du workflow NGS pour le readset. Concernant l'onglet « Infos échantillon », référence toutes les informations à propos de l'échantillon du readset. On y retrouve son code, le taxon dont il fait partie et son ID, la catégorie d'échantillon (ADN, ARN ...), la listes des barcodes utilisés et d'autres informations

APY_DA_AEKL_1_F300001324.MGI001 Read generation en cours

Général Avancé **Infos échantillon** Infos workflow

Code d'échantillon	APY_AEKI	% par piste	
Ref. Collaborateur	125SUR0CCKK11	Type processus banque	DA - DNaseq
Taxon Id	408172	Tag	MGI001
Taxon	marine metagenome	Layout Nominal Length (pb)	-1
Type d'échantillon	ADN Métagénomique	Liste tags primaires	MGI001,MGI002
Catégorie d'échantillon	ADN	Orientation brin synthétisé	undef
Code support container	F300001324	Fraction run (%)	
Code container	F300001324_1		

FIGURE 18 – Capture d'écran de la page du readset APY_DA_AEKL_1_F300001324.MGI001 de NGL en cours de génération de reads (onglet « Infos échantillon »)

APY_DA_AEKL_1_F300001324.MGI001 Contrôle qualité en attente

Général Avancé **Infos échantillon** Infos workflow

Etat	Date	Par
Nouveau	08/08/2022 17:36:35	ngsrg
Read generation en cours	08/08/2022 17:36:35	ngsrg
Read generation terminée	08/08/2022 17:59:08	ngsrg
Contrôle qualité en attente	08/08/2022 17:59:08	ngsrg

FIGURE 19 – Capture d'écran de la page du readset APY_DA_AEKL_1_F300001324.MGI001 de NGL en cours de génération de reads (onglet « Infos workflow »)