



## 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**

---

**2 septembre 2022**

William Amory  
sous la responsabilité de Frédérick Gavory



# Table des matières

<b>Glossaire</b>	<b>1</b>
<b>1 Introduction</b>	<b>3</b>
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
<b>2 Objectifs de ma mission</b>	<b>5</b>
<b>3 Matériels et Méthodes</b>	<b>6</b>
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 Librairie et API Perl permettant d'interagir avec la base de données NGL .	7
3.5 Logiciels de démultiplexage et génération de fichiers de séquences (bcl2fastq - bcl-convert) . . . . .	7
3.6 Les pipelines de génération de fichiers de séquences pour les technologies Illumina et Nanopore . . . . .	8
3.7 Les pipelines de contrôle qualité des lots de séquences pour les technologies Illumina et Nanopore . . . . .	8
<b>4 Résultats</b>	<b>9</b>
4.1 Etude comparative des logiciels bcl2fastq et bcl-convert . . . . .	9
4.2 Le pipeline de génération de fichiers de séquences pour la technologie MGI	11
<b>5 Discussions et perspectives</b>	<b>19</b>
5.1 Perspectives du workflow NGS pour la technologie MGI . . . . .	19
5.1.1 Améliorations futures du pipeline NGS_RG pour la technologie MGI	19
5.1.2 Développement du pipeline de contrôle qualité pour le technologie MGI . . . . .	19
5.2 Evaluation d'outils de contrôle qualité . . . . .	22
<b>Notes</b>	<b>23</b>
<b>Références</b>	<b>24</b>
<b>6 Annexes</b>	<b>25</b>

## 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)

**DNB** : *DNA-nanoballs* (Nano « billes » d'ADN générés lors de l'amplification ADN pour les séquenceurs de la technologie MGI)

**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* (bases de données du Genoscope et du CNRGH)

**NGL\_BI** : *NGL Bioinformatic* (base de données des analyses et traitements bio-informatique)

**NGL\_PROJECT** : *NGL projects* (base de données des projets en cours et passé)

**NGL\_REAGENT** : *NGL reagent* (base de données des réactifs)

**NGL\_SEQ** : *NGL Sequencing* (base de données de suivi des échantillons)

**NGL\_SUB** : *NGL submission* (base de données des soumissions de projet ou d'articles (exemple : la soumission d'un projet au NCBI))

**NGS** : *Next Generation Sequencing*

**NGS\_BA** : *Next Generation Sequencing - biological analysis*

**NGS\_QC** *Next Generation Sequencing - quality control*

**NGS\_RG** : *Next Generation Sequencing - reads generation*

**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.

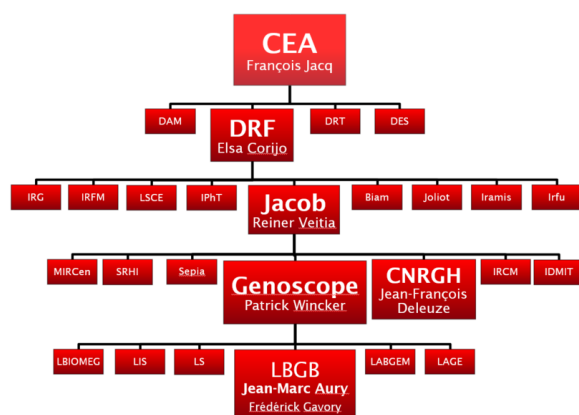


FIGURE 1 – Organigramme situant l'équipe du LBGB au sein du Genoscope et du CEA

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 différentes technologies de séquençage, d'effectuer l'assemblage<sup>1</sup> des séquences et l'annotation<sup>2</sup> des génomes, dans l'objectif de mettre à disposition des laboratoires collaborateurs internes ou externes 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 « d'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

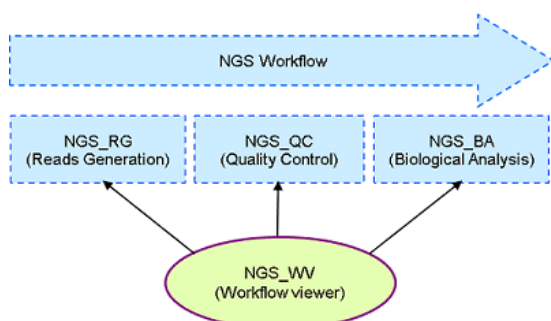


FIGURE 2 – Workflow de génération, de contrôle qualité et d'analyse biologique des fastq

Le workflow NGS est composé de trois pipelines pour les technologies Illumina et Oxford Nanopore. Le premier, NGS\_RG, permet la génération des reads<sup>3</sup> et des fichiers de séquences correspondants aux échantillons. Le second, NGS\_QC, permet de réaliser leur contrôle qualité. Le dernier, NGS\_BA, permet de faire les analyses biologiques inter-échantillons (readset<sup>4</sup>).

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, runs<sup>5</sup> 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 – Séquenceurs DNBSEQ-G400 (en haut) et DNBSEQ-T7 (en bas) de MGI  
<https://en.mgi-tech.com/products/>

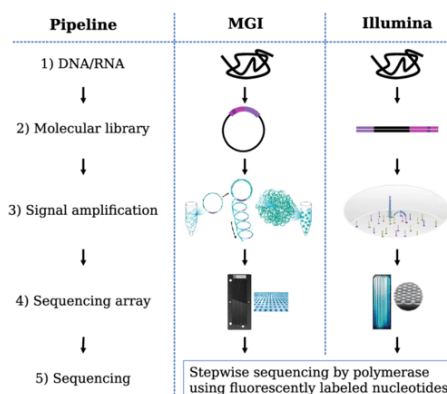


FIGURE 4 – Différences entre Illumina et MGI de technologie NGS

Il s'agit de séquenceurs à haut débit (DNBSEQ-G400) et très haut débit (DNBSEQ-T7), dont les principales différences entre MGI et Illumina sont dans la création des librairies<sup>6</sup> 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*<sup>7</sup>), puis déposée sur la flowcell<sup>8</sup> pour MGI, alors que pour Illumina elle est réalisée après immobilisation sur les flowcells.

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
<b>Data output/Run</b>	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 MGI. Dans un premier temps il s'agira de créer un pipeline de génération de fichiers de séquences (NGS\_RG\_MGI<sup>9</sup>) puis un pour le contrôle qualité de ces fichiers (NGS\_QC\_MGI<sup>10</sup>). Le workflow devra créer et mettre à jour l'état des runs, des pistes ou (*lane*) et de readset<sup>11</sup> dans NGL, réaliser le contrôle qualité des fichiers de séquences, au format FASTQ, obtenus en fin de séquençage. Il devra mettre à jour l'avancement du traitement d'un run dans NGL, en y insérant les métriques et statistiques obtenues lors du démultiplexage<sup>12</sup>, 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 », « annotation ») 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 l'outil fastx\_clean

de l'extension fastxend de la suite FASTX Toolkit<sup>13</sup> 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 composée 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 qui demande peu de mémoire Ram. Les partitions « xlarge » et « xxlarge » ont chacun deux noeuds composés de 48 coeurs et 2To de Ram, 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 composée de 12 noeuds qui disposent de 16 coeurs et de 257 Go de Ram. Les différents pipeline du workflow NGS utilisent cette partition. 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\_BI, est la partie de la base de données utilisée par les équipes de bioinformatique. NGL\_SEQ, 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\_SUB, NGL\_REAGENT et NGL\_PROJECTS. La gestion et le suivi des développements informatiques 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 sera 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 seront utilisés dans le développement 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*<sup>14</sup> et *paired end*<sup>15</sup>. 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 modules<sup>16</sup> pour les différents pipelines.

### 3.4 Librairie et API Perl permettant d'interagir avec la base de données NGL

L'interaction entre les pipelines du workflow NGS et la base de données NGL s'effectue par des fichiers JSON<sup>17</sup>. Cette interaction est possible grâce à une API<sup>18</sup> développé en Perl par l'équipe de « production » du LBGB, elle permet d'ajouter, modifier, récupérer, supprimer des données des fichiers JSON. Une librairie Perl (*DBFactory*) permet d'interagir avec cette API directement depuis une autre librairie ou script Perl, c'est cette dernière qui sera utilisé dans le développement des pipelines du workflow NGS pour la technologie MGI.

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

Ces deux logiciels de génération de fichiers de séquences et de démultiplexage (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 la technologie 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*<sup>19</sup>), temps CPU (*CPU time*<sup>20</sup>), pourcentage d'utilisation CPU (*%CPU*<sup>21</sup>)) en fonction des ressources disponibles sur les noeuds du cluster (*inti*) réservé à la *production*, afin de 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*<sup>22</sup> accordé pour la décompression et la lecture des *Bases Calls*<sup>23</sup>
- **p** [bcl-num-conversion-threads] : nombre de *threads* accordé pour la conversion des *Bases Calls* en fastq
- **w** [bcl-num-compression-threads] : nombre de *threads* accordé l'écriture et la compression des fichiers fastq

J'ai réalisés tous ces tests 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.

### 3.6 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 métriques, graphiques et fichiers permettant leurs évaluations.

Concernant le pipeline de génération de fichiers de séquences pour la technologie MGI, j'ai développé un pipeline dont l'objectif final est le même que 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 par le séquenceur, pour permettre de les insérer dans NGL\_BI

### 3.7 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. Ils réalisent le contrôle qualité et l'estimation des duplicats de séquence des fichiers avant et après nettoyage (*trimming*), ils retirent le *PhiX*<sup>24</sup> (pour les technologies Illumina), réalisent l'assignation taxonomique des séquences, réalisent un alignement des séquences si un génome de référence existe, réalisent le calcul du pourcentage de séquences qui ont leurs reads *forward* (brin sens) et *reverse* (brin anti-sens) qui se chevauchent et réalisent 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, il s'agira de développer un pipeline dont l'objectif est le même que celui d'Illumina en prenant en compte qu'avec cette technologie il n'y auras pas de *PhiX* à enlever dans les fichiers de séquences.

## 4 Résultats

### 4.1 Etude comparative des logiciels bcl2fastq et bcl-convert

#### 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.

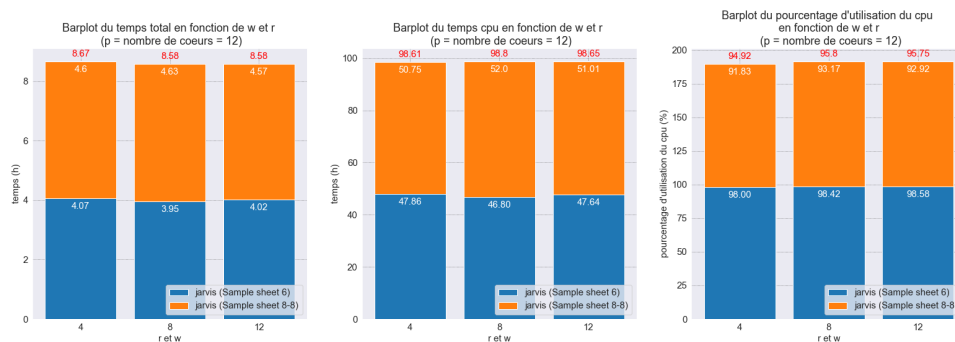


FIGURE 5 – Digrammes en bâtons du temps total d'exécution (à gauche), temps cpu (au milieu) et du pourcentage d'utilisation des cpu (à droite) en fonction des paramètres  $r$  et  $w$

Il y a deux *sample sheet*<sup>25</sup>, car le nombre de bases considérées des index<sup>26</sup> entre les pistes 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 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 ( $> 90\%$ ).

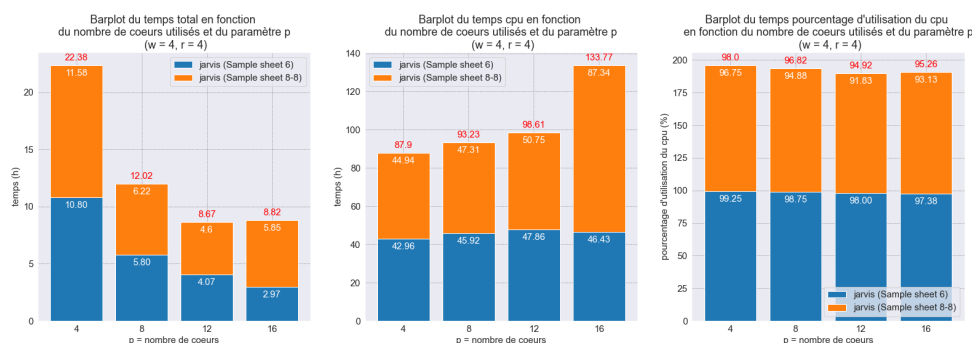


FIGURE 6 – Digrammes en bâtons du temps total d'exécution (à gauche), temps cpu (au milieu) et du pourcentage d'utilisation des cpu (à droite) en fonction du paramètre  $p$

Néanmoins, on remarque qu'il n'y a pas d'amélioration du temps total d'exécution entre  $p$  fixé à 12 et à 16 coeurs. on observe même une augmentation du temps cpu.

Au vue des résultats obtenus j'ai décidé que les meilleurs paramètres étaient de fixer  $p$  à 12, puisque le gain apporté en augmentant à 16 est faible. Néanmoins j'ai décidé de le conserver 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 simultanée sur un seul noeud de calcul de la partition « production » puisqu'ils font 16 coeurs.

### 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, 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.

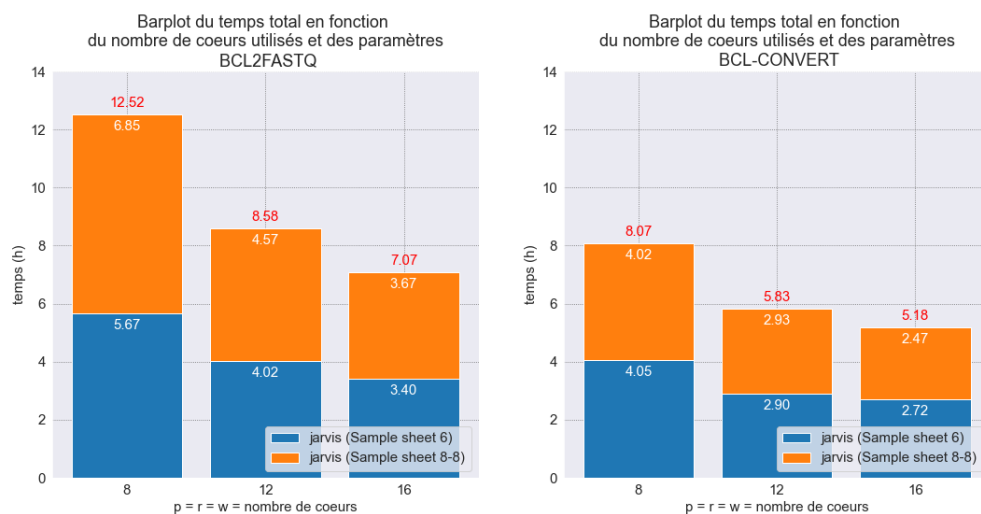


FIGURE 7 – Temps total de génération des FASTQ pour bcl2fastq et bcl-convert

J'ai également échangé avec le service technique d'Illumina à propos des fichiers de sortie et de l'arborescence de ces 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 si l'on pouvait obtenir une arborescence similaire à bcl2fastq, pour minimiser l'impact du changement de logiciel sur les pipelines. Le changement de bcl2fastq, qui sera bientôt obsolète, par bcl-convert va donc nous obliger à réaliser de gros changements dans tous les pipelines qui utilisent ces fichiers de sortie.

## 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 t'attribuer l'intégralité des coeurs d'un noeud de « production », c'est à dire 16 coeurs et de favoriser le temps d'exécution.

J'ai consigné l'intégralité des changements entre les deux logiciels dans un cahier des charges. Il contient, la nouvelle commande à lancer, les modules à charger dans l'environnement, le chemin relatif des fichiers de sorties et leurs descriptions, ainsi qu'un exemple d'arborescence des fichiers de sorties. Ce qui permettra au développeur qui ce 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 de l'équipe, qui sera en charge de réaliser cette migration.

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

Le pipeline NGS\_RG\_MGI que j'ai développé à pour objectif 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 dans NGL\_BI, à l'aide de la librairie Perl permettant d'interagir avec celle-ci (cf. 3.5 page 7). Les différentes étapes du pipeline suivent le schéma (figure 8) que l'on a défini ci-dessous.

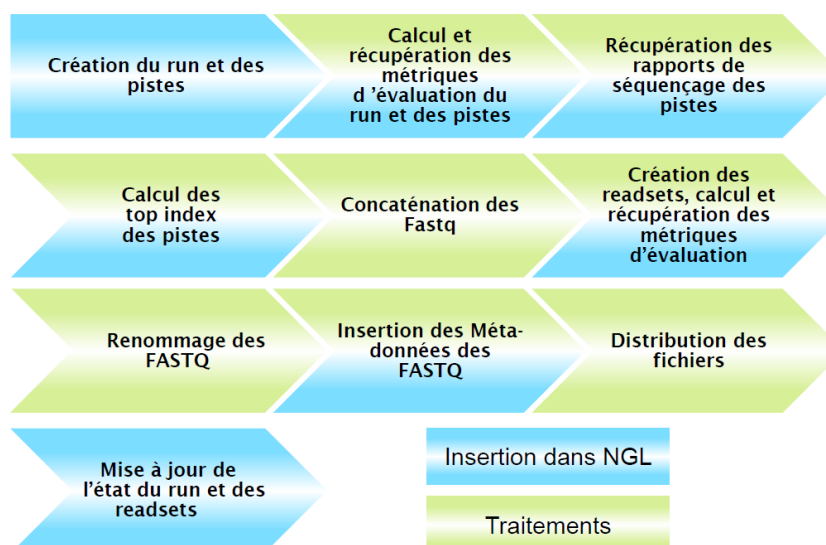


FIGURE 8 – Schéma des différentes étapes du pipeline NGS\_RG\_MGI. Les étapes qui interagissent avec NGL sont en bleu et les étapes demandant un traitement informatique des données de séquençage sont en vert

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

La première étape du pipeline, que je développée, consiste à créer le run et ses pistes dans la base de données NGL, en y insérant les métriques permettant d'évaluer le run et les pistes (figure 9). Le nom du run est constitué de la date de séquençage, du nom du séquenceur et de l'identifiant de la flowcell du run ce qui le rend unique.

Les différentes métriques sont insérées à l'aide des librairies Perl permettant d'interagir avec ngl en postant ses métriques dans le fichier JSON du run de la base de données, ce qui permet l'affichage de ces dernières dans l'interface web de NGL-BI. Toutes ses métriques sont détaillées plus précisément en annexes (page 25).

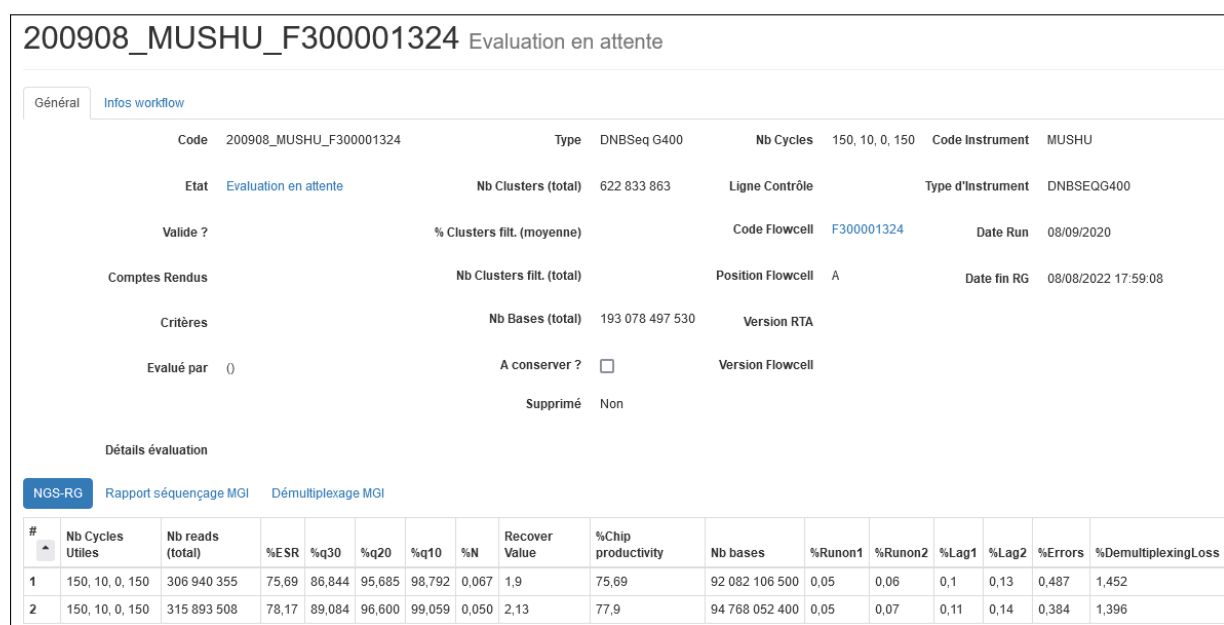


FIGURE 9 – Capture d'écran de la page du run 200908\_MUSHU\_F300001324 de NGL en cours de génération de fichiers de séquences (étapes d'ajout des métriques d'évaluation du run et des pistes).

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

J'ajoute ensuite, en seconde étape, les rapports de séquençage des pistes que le séquenceur génère en fin de séquençage. Il s'agit de rapports 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 10.A), de la distribution des bases nucléiques en fonction des cycles (figure 10.B), de la distribution du pourcentage de Guanine/Cytosine en fonction des cycles (figure 10.C), de la distribution de l'intensité brut au cours des cycles (figure 10.D). Les tableaux et graphiques de ces rapports de séquençage permettent de faciliter l'évaluation du run et de ses pistes.

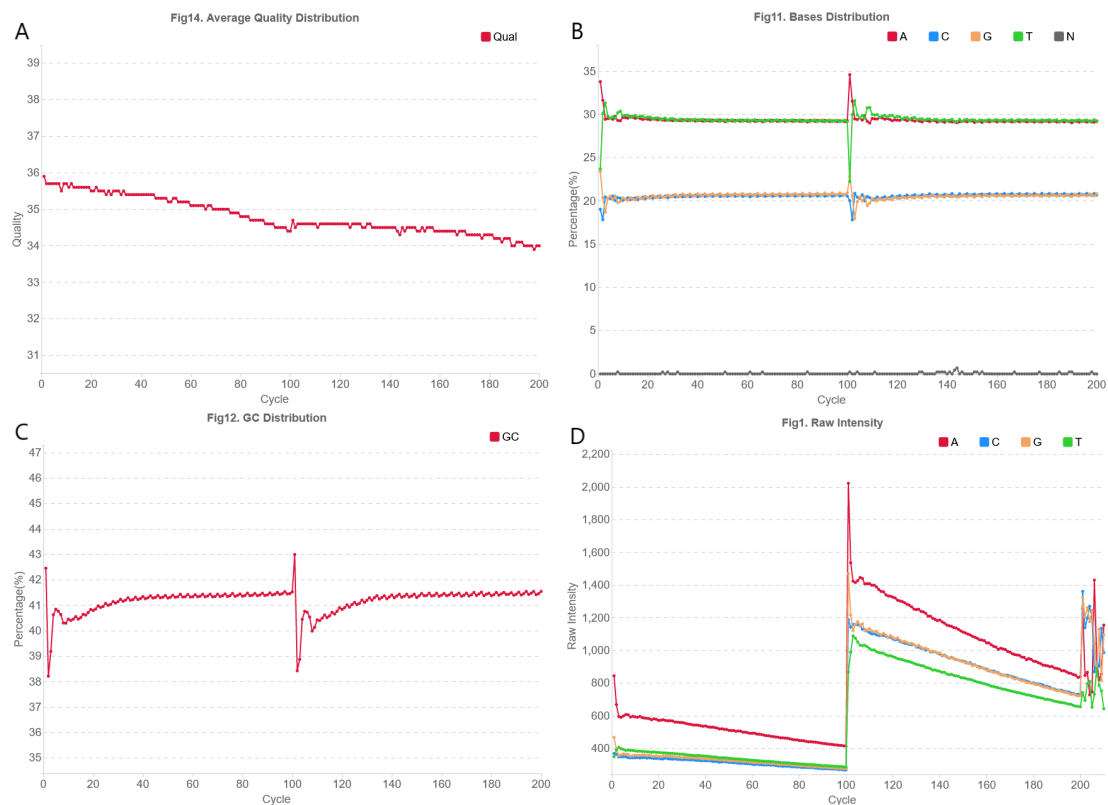


FIGURE 10 – Graphiques des distributions de la qualité moyenne (A), des bases nucléiques (B), du pourcentage de GC (C) et de l'intensité brut (D) au cours des cycles de séquençage

Toujours dans l'optique de faciliter l'évaluation du run et de ces pistes, l'étape suivante du pipeline, a été d'ajouter la liste des index représentées à plus de 0.01% de la piste, ainsi que les index attendus. Ces index sont triés et affichés par ordre décroissant dans NGL (figure 11). 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.

NGS-RG   Rapport séquençage MGI   Démultiplexage MGI

Lane 1

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
GATTCTCCT	59 842	0,019
CGGCAGAAAGT	53 283	0,017
barcode29	48 716	0,016
barcode124	37 751	0,012
CGGCGTAAGT	36 893	0,012

FIGURE 11 – 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 quatrième étape du pipeline que j'ai développé à pour objectif d'obtenir un seul fichier FASTQ par readset. En effet la technologie MGI requiert une homogénéité en composition en base nucléiques (A, T, C, G) au niveau de chaque cycle des index, un déséquilibre étant susceptible d'entraver la récupération du signal pour les bases de ces cycles et donc de fausser leur *base calling* puis le démultiplexage des séquences. Ainsi il est recommandé de « barcoder » les échantillons avec 4 barcodes (index), dont les séquences assurent une composition égale en A, T, C et G à chaque position des barcodes. Dans ce cas, nous obtenons donc plusieurs fichiers par échantillons qu'il faut donc fusionner pour en obtenir qu'un seul par échantillon lors du démultiplexage effectué par le séquenceur.

Si le readset est associé à un seul readset alors on réalise une décompression du fichier FASTQ, à l'inverse si il est associé à plusieurs index on réalise une décompression et une concaténation des fichiers FASTQ (cf. figure 12).



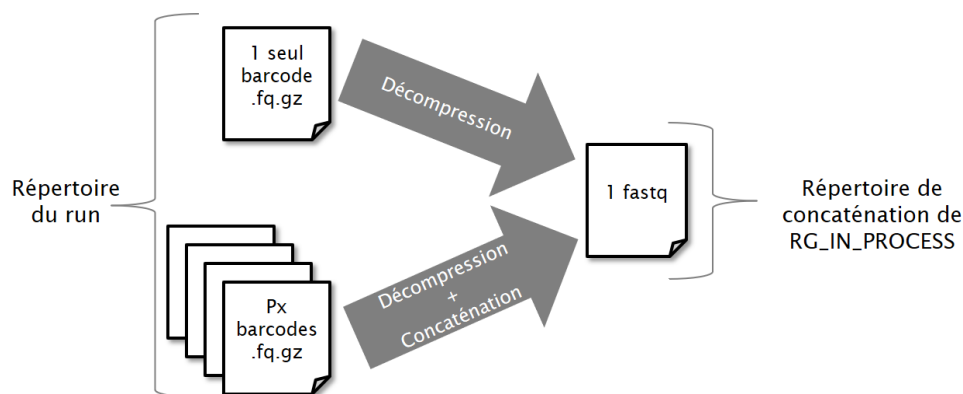


FIGURE 12 – Schéma de l'étape de « concaténation » des fichiers FASTQ d'un readset

### Création et insertion des métriques des readset du run dans NGL

La cinquième étape a pour objectif de permettre l'évaluation des readsets, en les créants et en insérant les métriques d'évaluation de ces derniers, récupérer en parssant des fichiers ou/et en les calculant, dans NGL (figure 13). 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. J'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, etc. Toutes ces métriques sont décrites en annexes (page 26)

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 ce qui le rend unique également.

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 13 – Capture d'écran de la page du readset APY\_DA\_AEKL1\_F300001324.MGI001 de NGL en cours de génération de reads (étapes de création du readset et d'insertion de ces métriques d'évaluation)

L'étape suivante du pipeline est d'ajouter la répartition des index au sein d'un readset (figure 14), ce qui permet de vérifier la composition en index du readset et de vérifier l'homogénéité de ces index au sein du 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 14 – Capture d'écran de la page du readset APY\_DA\_AEKL1\_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 à partir des métriques que j'ai ajoutées dans le fichier JSON du readset. (figure 15).

Readsets (7)										
Lanes										
Rechercher										
Voir Readsets Evaluer Readsets										
Taille (10)										
N° Piste	Code	Etat	% déposé	% Séquences valides / piste	Nb Séquences valides	Nb Bases	% >= Q30	Score Qualité moyen	Valide QC ?	Valide Bioinfo ?
1	APY_DA_AEKL1_F300001324.MGI001	Read generation en cours	56,59		173 704 226	52 111 267 800	88,52	34,58	---	---
1	CRH_DA_AAAA1_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 15 – 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 septième étape, consiste à renommer les fichiers de séquences des readsets et d'insérer les méta-données de ces derniers dans NGL via le fichier JSON du readset (figure 17). Le renommage des fichiers est nécessaire pour que chaque fichiers de séquences aient un nom unique et « parlant ».

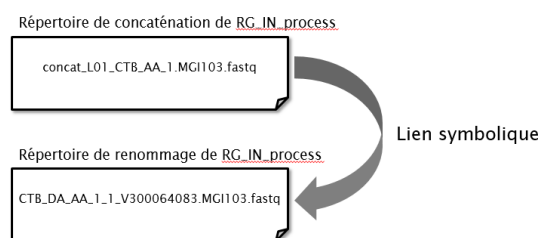


FIGURE 16 – Schéma de l'étape de renommage des fichiers FASTQ d'un readset

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 16).

Les méta-données des fichiers de séquences du readset que j'insère 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 (« raw », « clean ») et s'ils sont « utilisable », c'est à dire si ils représentent le niveau de qualité maximale disponible pour les échantillons en question. On y retrouve donc le chemin vers le répertoire 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 des valeur de la qualité (Pour les séquenceurs MGI l'encodage est en Phred +33).

L'encodage Phred +33, signifie que l'encodage des valeurs de la qualité sont encodées à partir du caractère « ! », qui vaut 33 en base 10 dans la table ASCII<sup>27</sup> pour une valeurs de qualité de 0. Il y a 42 niveaux de valeurs de qualité, 41 étant le niveau maximale de qualité d'une base nucléique. L'encodage de la qualité en Phred +33, est donc encodé à partir du caractère « ! » jusqu'au caractère « J », qui vaut 74 en base 10 dans la table ASCII et représente la valeur maximale de la qualité d'une base.

Un score qualité de 10 indique qu'il y a un probabilité de 90% que le *base call* soit correct (1 risque sur 10 que la base soit incorrect) et un score de 40 indique une probabilité de 99,99% que le *base call* soit correct (1 risque sur 10000 que la base soit incorrect).

## APY\_DA\_AEKI\_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_AEKI_1_1_F300001324.MGI001.fastq	RAW	Oui	READ1	33		
APY_DA_AEKI_1_2_F300001324.MGI001.fastq	RAW	Oui	READ2	33		

FIGURE 17 – Capture d'écran de la page du readset APY\_DA\_AEKI\_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 huitième étape que j'ai développé, permet 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 de séquences « attendus » sont copiés vers leur répertoire final en fonction du centre dans lequel le séquençage a eu lieu (Genoscope, CNRGH) et leurs droits d'accès sont modifiés pour que les utilisateurs aient le droit de lecture, mais n'aient pas les droits d'écriture et d'exécution.

Les fichiers de statistiques du run sont archivés et compressés par pistes et par types (.html, .fq.stat) avant d'être copiés vers leur répertoire final et leurs droits sont changés pour les mêmes raisons. Ces fichiers sont conservés dans le cas où une métrique désirée 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 de séquencer « non attendus », il s'agit des fichiers de séquences des index ne faisant pas partie d'un readset. Puisque lors du démultiplexage par les séquenceurs on obtient un fichier FASTQ par index. Ces fichiers sont renommés, archivés, compressés et leurs droits sont changés pour les mêmes raisons que les fichiers de séquences « attendus », 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 séquençage, 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 readsets dans l'état de « F-RG », correspondant à la fin du pipeline NGS\_RG. De plus, les runs n'étant plus dans un état "IW-RG" ou "IP-RG", correspondant à l'attente de prise en charge par NGS\_RG une fois le séquençage terminé ou en cours de génération de reads, ils ne seront plus éligibles à une prise en charge par NGS\_RG. Cela entraîne donc une mise à jour automatique du run à l'état « IW-V », correspondant à l'attente de validation, ce qui permet d'indiquer aux utilisateurs que le run peut être évalué. Les readsets sont aussi automatiquement mis à jour vers l'état « IW-QC », correspondant à l'attente de prise en charge par le pipeline NGS\_QC, ce qui permet d'indiquer au pipeline de contrôle qualité qu'il peut effectuer le contrôle qualité des readsets 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

Les pipelines de génération de fichiers de séquences est très spécifique à l'environnement de gestion des projets de séquençage du Genoscope et du CNRGH, ainsi qu'à la base de données NGL. Celui pour la technologie MGI est inspiré du pipeline de l'autre technologie de séquençage *short reads* (Illumina) déjà en place, dont la finalité du pipeline est le même, c'est à dire la prise en charge des données de séquençage en fin de séquençage et la mise à disposition des fichiers de séquences et des métriques d'évaluation.

La future amélioration du pipeline NGS\_RG\_MGI, consistera à la mise en place d'une étape supplémentaire pour les runs qui comporterons des *mids*<sup>28</sup>. Un *mid* est une séquence d'une dizaines de nucléotides ajoutés en amont du primer du read *forward* (figure 18). Il s'agit d'un index supplémentaire qui permet lors du séquençage de déposer un nombre plus important d'échantillons différents sur une même piste d'une flowcell.

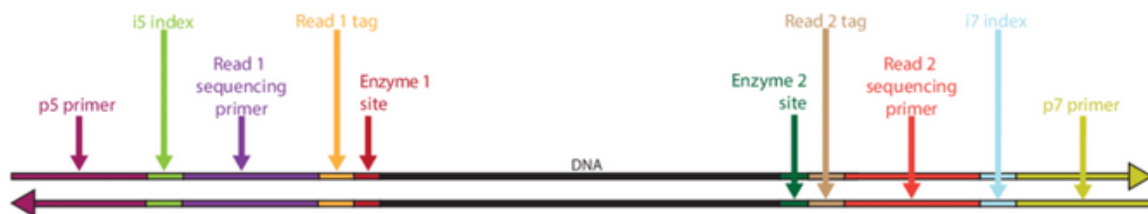


FIGURE 18 – Schéma représentant le vecteur à séquençer qui contient les index (i5 index et i7 index), les mids (Read 1 tag et Read 2 tag), les primers des index (p5 primer et p7 primer) et des reads (Read 1 sequencing primer et Read 2 sequencing primer), ainsi que l'ADN d'intérêt (en noir) et les sites enzymatiques permettant l'insertion de l'ADN d'intérêt dans le vecteur (Enzyme 1 site et Enzyme 2 site).

L'étape supplémentaire sera d'ajouter un second démultiplexage en fonction de ces mids, qu'on appelle le démidage pour la création des readset et des 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, sera constitué de deux grande phases. La première consistera à réaliser un contrôle qualité des fichiers « raw » (fichiers de séquences bruts), puis dans un second temps de réaliser un contrôle qualité sur les fichiers « clean » (fichiers de séquences nettoyés) ainsi que d'autres traitements sur les fichiers « clean ».

Le pipeline devra prendre en charge automatiquement les fichiers dont les readset sont dans l'état « IW-QC » dans NGL, qui représente les readset en attente de contrôle qualité. Celui-ci devra suivre le schéma de traitements et d'analyse (figure 19) déjà en place pour l'autre technologie de séquençage *short reads* (Illumina). L'étape de nettoyage du *PhiX* ne sera pas nécessaire pour la technologie MGI, puisque pour cette dernière il n'est pas utile d'ajouter ces séquences de contrôle.

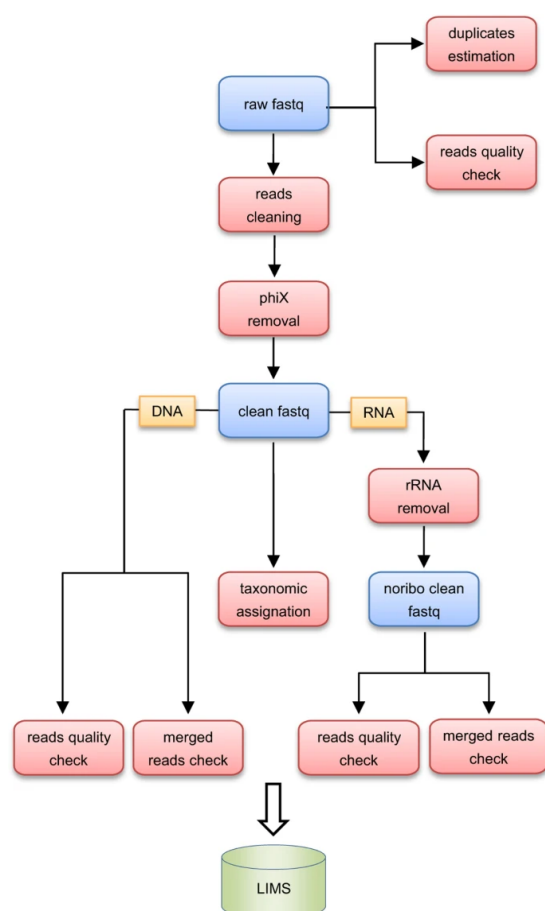


FIGURE 19 – Schéma des étapes du pipeline NGS\_RG-Illumina

Le Trimming consistera à retirer les séquences qui ont une qualité moyenne faible (inférieur à 20), les dernières bases des séquences qui ont une faible qualité (inférieur à 20) seront retirées de la séquence. Si il y a trop de bases inconnus dans la séquence celle-ci sera retirée, de même si dans la séquence on retrouve 3 ou plus de bases inconnus successivement dans la séquence celle-ci sera coupée jusqu'à la fin de la séquence.

Les différentes étapes qui seront à réaliser pour la première phase sont :

- Réaliser le contrôle qualité des séquences des fichiers brut, avec l'outil fastx\_clean

de l'extention [fastxend](#) de la suite [FASTX Toolkit](#) développé par le Genoscope.

- Réaliser l'estimation des duplicats des séquences des fichiers brut, avec l'outil `fastx_estimate_duplicate` de l'extension `fastxend` de la suite [FASTX Toolkit](#) développé par le Genoscope.

Les différentes étapes qui seront à réaliser pour la seconde phase, une fois le Trimming effectué sont :

- Réaliser le contrôle qualité des séquences des fichiers clean, avec l'outil `fastx_clean`.
- Réaliser l'estimation des duplicats de séquences des fichiers clean, avec l'outil `fastx_estimate_duplicate`.
- Réaliser l'assignation taxonomique des séquences des fichiers clean, à l'aide du logiciel [Centrifuge](#).
- Réaliser l'alignement des séquences des fichiers clean sur un génome de référence si celui-ci est disponible, à l'aide de l'outil [BWA](#).
- Réaliser le « *merging* » des séquences des reads *forward* et *reverse* pour les run *pair end* (calcul du pourcentage de reads qui ont le read *forward* et *reverse* qui se chevauchent, calcul de la moyenne et médiane du nombre de bases qui se chevauchent entre les 2 reads, ...), à l'aide de l'outil `fastx_mergepairs` de l'extension `fastxend` de la suite [FASTX Toolkit](#) développé par le Genoscope.

Une fois les étapes de la seconde phase réalisées, il sera nécessaire de réaliser la distribution des fichiers de séquences nettoyés dans leur répertoire final. Les fichiers raw seront alors effacés si ces derniers ont été archivés sur bande magnétique, et rendu indisponible pour les utilisateurs. L'objectif étant de permettre aux utilisateurs d'avoir les fichiers de séquences avec le meilleur niveau de qualité possible.

Durant toutes les étapes du pipeline, on insère les métriques et graphiques obtenus au cours des différentes étapes du pipeline dans `NGL_BI`. Ce qui permettra de réaliser la validation des readset ou non. L'interaction entre le pipeline et la base de données est réalisée, comme pour le pipeline `NGS_RG_MGI`, par la librairie Perl (`DBFactory`) qui permet d'interagir avec `NGL` (cf. [3.5](#) page 7).

L'objectif est d'obtenir un pipeline de contrôle qualité opérationnel le plus rapidement possible, c'est pour cela que les outils utilisés seront les mêmes que ceux utilisés pour le pipeline `NGS_QC_Illumina`. Néanmoins, les outils utilisés évolueront au fil du temps, avec les évaluations d'outils pour les pipelines de contrôle qualité que je réaliserai au cours de l'année à suivre.

## 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 taxinomique des séquences et des fichiers de sortie.



## Notes

<sup>1</sup>Reconstruction d'un génome à partir de fragments de ce dernier

<sup>2</sup>Documenter le plus exhaustivement possible les informations de l'assemblage permettant de prédire la fonction d'un gène, d'une molécule, d'une région de l'ADN, ...

<sup>3</sup>Lecture d'une séquence par un séquenceur d'un fragments d'ADN

<sup>4</sup>Un lot de séquences est une instance de séquences (ou reads) d'un échantillon

<sup>5</sup>Séquençage d'un ou plusieurs échantillons sur un séquenceur

<sup>6</sup>Collection 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)

<sup>7</sup>Nanobilles d'ADN générées par la réplication de l'ADN circulaire

<sup>8</sup>Lame d'absorption des fragments d'ADN et cuve réacteur du séquençage

<sup>9</sup>*Next Generation Sequencing - reads generation - mgi*

<sup>10</sup>*Next Generation Sequencing - quality control - mgi*

<sup>11</sup>Lot de séquences

<sup>12</sup>Séparation des séquences en plusieurs fichiers en fonction de leurs index (séquence d'une dizaines de nucléotides en amont du primer de la séquence)

<sup>13</sup>Collection de commandes pour le traitement et l'évaluation de lot de séquences au format FASTA ou FASTQ

<sup>14</sup>Lecture dans un seul sens des reads par le séquenceur

<sup>15</sup>Lecture dans les deux sens des reads par le séquenceur

<sup>16</sup>Un 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 derniers.

<sup>17</sup>*JavaScript Object Notation* est un format de données textuelles structurées et organisées

<sup>18</sup>*application programming interface* est une interface logicielle qui permet de « connecter » un logiciel ou un service à un autre logiciel ou service afin d'échanger des données et des fonctionnalités

<sup>19</sup>Temps écoulé entre le début du programme et le fin de celui-ci

<sup>20</sup>Temps d'utilisation des cpu par le programme

<sup>21</sup> $((CPU\ time + \text{temps utilisé par les appels système}) / Elapsed\ time) / \text{nombre de CPU utilisé par le programme}$

<sup>22</sup>Processus : instructions du langage machine d'un processeur.

<sup>23</sup>Fichier d'attribution des bases nucléiques en fonction des pics du chromatogramme lors du séquençage

<sup>24</sup>Parties du génome du phage *Lambda* qui sont ajoutés sur les pistes des flowcell avant le séquençage, permettant de contrôler le bon déroulé du séquençage.

<sup>25</sup>Fichier contenant les informations et instructions pour la génération des FASTQ et le démultiplexage

<sup>26</sup>Séquence d'une dizaines de nucléotides en amont du primer de la séquence d'ADN à séquencer, permettant de séparer les séquences de plusieurs échantillons sur une même piste de la flowcell (démultiplexage)

<sup>27</sup>*American Standard Code for Information Interchange* est une norme de l'encodage des caractères

<sup>28</sup>un mid (*molecular identifier*) est une séquence d'une dizaine de nucléotides ajoutés en aval du *primer* du read *forward* permettant de réaliser un second démultiplexage

## Références

- [1] BCL Convert.
- [2] BCL Convert Software Guide v3.7.5 (1000000163594). page 22.
- [3] bcl2fastq2 Conversion Software v2.20 Software Guide (15051736). page 27.
- [4] The Comprehensive Perl Archive Network - [www.cpan.org](http://www.cpan.org).
- [5] perl - The Perl 5 language interpreter - Perldoc Browser.
- [6] BGI goes head-to-head with Illumina. *Nature Biotechnology*, 33(8) :792–792, Aug. 2015.
- [7] S. Anslan, V. Mikryukov, K. Armolaitis, J. Ankuda, D. Lazdina, K. Makovskis, L. Vesterdal, I. K. Schmidt, and L. Tedersoo. Highly comparable metabarcoding results from MGI-Tech and Illumina sequencing platforms. *PeerJ*, 9 :e12254, Sept. 2021.
- [8] C. A. Austin-Tse, V. Jobanputra, D. L. Perry, D. Bick, R. J. Taft, E. Venner, R. A. Gibbs, T. Young, S. Barnett, J. W. Belmont, N. Boczek, S. Chowdhury, K. A. Ellsworth, S. Guha, S. Kulkarni, C. Marcou, L. Meng, D. R. Murdock, A. U. Rehman, E. Spiteri, A. Thomas-Wilson, H. M. Kearney, H. L. Rehm, and Medical Genome Initiative\*. Best practices for the interpretation and reporting of clinical whole genome sequencing. *npj Genomic Medicine*, 7(1) :27, Dec. 2022.
- [9] S. Drmanac, M. Callow, L. Chen, P. Zhou, L. Eckhardt, C. Xu, M. Gong, S. Gablenz, J. Rajagopal, Q. Yang, C. Villarosa, A. Au, K. Davis, A. Jorjorian, J. Wang, A. Chen, X. Zhang, A. Borcharding, X. Wei, M. Zhang, Y. Xie, N. Barua, J. Shafto, Y. Dong, Y. Zheng, L. Wang, L. Zhai, J. Li, S. Liao, W. Zhang, J. Liu, H. Jiang, J. Wang, H. Li, X. Xu, and R. Drmanac. CoolMPS<sup>™</sup> : Advanced massively parallel sequencing using antibodies specific to each natural nucleobase. preprint, Genomics, Feb. 2020.
- [10] S. A. Jeon, J. L. Park, S.-J. Park, J. H. Kim, S.-H. Goh, J.-Y. Han, and S.-Y. Kim. Comparison between MGI and Illumina sequencing platforms for whole genome sequencing. *Genes & Genomics*, 43(7) :713–724, July 2021.
- [11] H.-M. Kim, S. Jeon, O. Chung, J. H. Jun, H.-S. Kim, A. Blazyte, H.-Y. Lee, Y. Yu, Y. S. Cho, D. M. Bolser, and J. Bhak. Comparative analysis of 7 short-read sequencing platforms using the Korean Reference Genome : MGI and Illumina sequencing benchmark for whole-genome sequencing. *GigaScience*, 10(3) :giab014, Mar. 2021.
- [12] J. Patterson, E. J. Carpenter, Z. Zhu, D. An, X. Liang, C. Geng, R. Drmanac, and G. K.-S. Wong. Impact of sequencing depth and technology on de novo RNA-Seq assembly. *BMC Genomics*, 20(1) :604, July 2019.

## 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 9 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** : *Effective spot rate* ( (nombre de reads / nombre total de DNB)  $\times$  100 )

**%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** : Rapport d'intensité entre le read *forward* et *reverse*

**%Chip productivity** : Pourcentage de productivité de la piste (nombre reads qui passent un pré-filtre MGI / nombre total de DNB)

**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 13 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 15 page 16) :

**%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

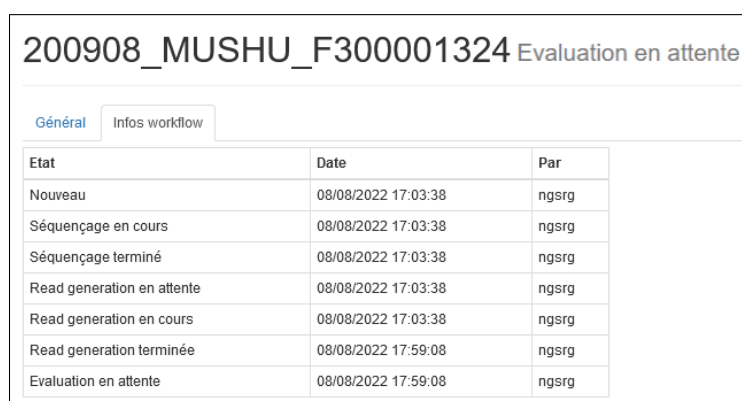
**Nb Bases** : Nombre de bases du readset

**% >= 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 20).



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 20 – 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 21) et de l'onglet « Infos workflow » (figure 22). 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\_AEKI\_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 21 – Capture d'écran de la page du readset APY\_DA\_AEKI\_1\_F300001324.MGI001 de NGL en cours de génération de reads (onglet « Infos échantillon »)

APY\_DA\_AEKI\_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 22 – Capture d'écran de la page du readset APY\_DA\_AEKI\_1\_F300001324.MGI001 de NGL en cours de génération de reads (onglet « Infos workflow »)