



Formation Galaxy

13 Novembre 2014

1-1 Connexion

☐ Connectez-vous sur la plateforme Galaxy SouthGreen à l'adresse suivante : <http://gohelle.cirad.fr/galaxy/>



☐ Utiliser votre adresse email et votre mot de passe pour vous identifier.

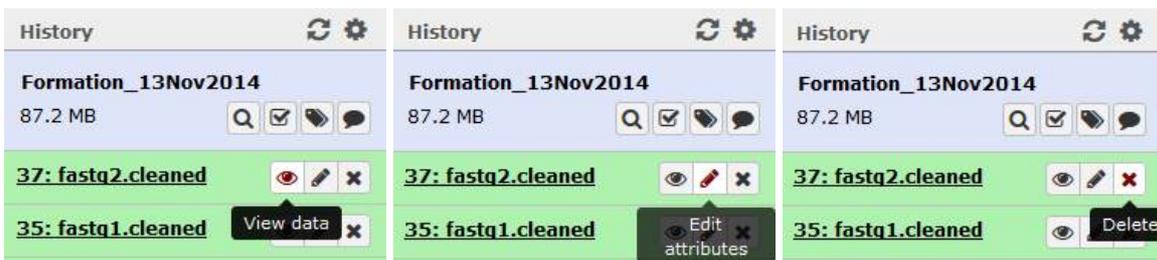
Si vous n'avez pas de compte, utiliser pour aujourd'hui le compte formationN/formationN

1-2 Import de fichier

☐ Dans l'onglet Tools à gauche de l'interface allez dans « Get Data » puis « Upload File ».

Galaxy permet d'importer un fichier via plusieurs moyens :

- Importer un fichier stocké localement sur votre ordinateur en cliquant sur « choisissez un fichier »
- Importer un fichier à partir d'une url en copiant l'adresse dans le cadre « URL/Text »
- Copier/coller directement le contenu du fichier dans le cadre « URL/Text »



Il est ensuite possible de visualiser, supprimer ou éditer les attributs d'un fichier.

Testez les 3 types d'import possibles de données. Vous observerez dans l'onglet « History » à droite de l'interface la progression de l'importation.

Galaxy offre un suivi de l'état de chaque job avec un code couleur :

- Bleu : le job a été soumis
- Jaune : le job est en cour de traitement
- Vert : le job s'est terminé avec succès
- Rouge : le job est en erreur

Import des données de la librairie partagée

Accédez aux données partagées (Shared data => Data libraries) et importez un fichier dans votre historique courant. (répertoire "Formation_Galaxy").

The screenshot shows the Galaxy web interface. At the top, there are navigation tabs: 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Admin', 'Help', and 'User'. The 'Shared Data' tab is selected. Below the navigation, there is a 'Data Libraries' section with a search bar and an 'Advanced Search' link. A dropdown menu is open over the 'Data Libraries' section, showing options: 'Data Libraries', 'Data Libraries Beta', 'Published Histories', 'Published Workflows', 'Published Visualizations', and 'Published Pages'. Below the dropdown, there is a table of data libraries. The table has two columns: 'Data library name' and 'Data lib'. The first row is highlighted in yellow and contains '454_Coffea' and 'Coffea'. Other rows include 'Amalia project', 'Arcad', 'Arcad_Cafe', 'Arcad_Fonio', 'Arcad_Mil', 'Arcad_Monococcum', 'Arcad_Viane', 'Arcad_Yam', 'Arcad_Riz', 'Azucena', 'Banana', 'CIAT_RICE', 'Clotault', 'cocos_chloroplast', 'cocos_fasta_files', 'cocos_fastq_files', 'coffea', 'Coffea_NBS', 'coffea_454', and 'Coffee_Genome_SNP'.

- Cliquez sur la library Formation
- Ouvrez le repertoire Preprocessing and Mapping
- Cochez les fichiers **input1.fastq**, **input2.fastq** et **Solexa_mRNA_primers.txt**
- Ouvrez le repertoire SNP
- Cochez le fichier **reference.fas.txt**
- En bas de page cliquez sur le bouton GO avec l'option "import to current history"

1-3 Exécution de jobs

Nous allons maintenant exécuter plusieurs opérations à partir de ce fichier fastq.

Pour trouver facilement un outil vous pouvez entrer son nom dans la case de recherche « search tools ».

Exécutez les programmes listés ci-dessous :

FastQC : Permet de contrôler la qualité des reads issus d'un séquençage NGS.

Sélectionnez un des fichiers fastq dans « FASTQ reads » puis cliquez sur « Execute ».

Fastq Groomer : Permet de recoder les scores de qualité de séquençage.

Réalisez un recodage de la qualité des fichiers fastq (Illumina 1.5) en qualité Sanger.

Cutadapt : Permet le nettoyage des fichiers fastq bruts (élimination des régions de mauvaises qualité, et des adaptateurs).

Réalisez le nettoyage des fichiers fastq 1 et 2, en renseignant la séquence de 3 adaptateurs contenus dans le fichier "Solexa_mRNA_primers.txt" et en fixant les paramètres suivants: qualité minimale à 20, taille minimale de l'overlap avec un adaptateur égale à 7, taille minimale d'un read égale à 20.

FastQ Interlacer : Permet de dissocier les reads singles et pairés.

Le nettoyage a parfois supprimé certains reads forward et gardé le reverse correspondant et vice-versa. Utilisez cet outil sur chaque fastq nettoyé pour obtenir un fichier fastq ne contenant que les reads en paires.

FastQ de-Interlacer : Permet de séparer les reads forward et reverse.

Utilisez cet outil pour régénérer 2 fichiers fastq gauche et droite.

BWA : Permet le mapping de reads Illumina sur une séquence de référence.

Réalisez le mapping des reads sur le fichier "reference.fas.txt" (3 gènes de Riz).

1-4 Utilisation des historiques

En cliquant sur Option => Saved histories Galaxy ouvre une interface de gestion des historiques.

Saved Histories

search history names and tags

[Advanced Search](#)

<input type="checkbox"/>	Name	Datasets	Tags	Sharing	Size on Disk	Created	Last Updated ↑	Status
<input type="checkbox"/>	mon historique	27	0 Tags		417.8 Kb	Jun 01, 2012	3 days ago	
<input type="checkbox"/>	Unnamed history	17	7	0 Tags	267.9 Kb	Mar 23, 2012	Mar 23, 2012	
<input type="checkbox"/>	Unnamed history	3	2	0 Tags	13.1 Kb	Jan 20, 2012	Jan 20, 2012	current history
<input type="checkbox"/>	ESTtik_pipeline	30	0 Tags		5.6 Mb	Nov 25, 2011	Dec 09, 2011	
<input type="checkbox"/>	Unnamed history	37	2	0 Tags	12.6 Mb	May 10, 2011	Nov 03, 2011	

For 0 selected histories:

Histories that have been deleted for more than a time period specified by the Galaxy administrator(s) may be permanently deleted.

- Donnez un nom à votre historique courant en cliquant sur « Unnamed history ».
- Créez un nouvel historique.

1-5 Création de workflow

Le but ici est de créer un workflow pour reproduire l'enchaînement des briques du paragraphe 1-3.

- Rendez vous dans la rubrique « Workflow » de Galaxy.
- Cliquez sur « Create new workflow » et attribuez lui un nom. Cliquez ensuite sur votre workflow puis sur « edit » pour ouvrir l'interface de création.

En cliquant sur une brique à gauche de l'interface une petite fenêtre apparaît sur le canevas. Cette fenêtre indique les fichiers d'entrée et de sortie de la brique. Pour enchaîner deux briques il suffit de relier le fichier de sortie de la première brique vers le fichier d'entrée de la seconde à l'aide des petites flèches.

The screenshot displays the Galaxy workflow editor. On the left, a 'Tools' sidebar lists various bioinformatics tools. The central 'Workflow Canvas' shows a workflow with the following steps: 'Input dataset' (output) feeds into 'FastQC' (FASTQ reads, Contaminants, report (html)) and 'FASTQ Groomer' (File to groom, output_file (fastq, fastqcssanger, fastqsolexa, fastqillumina)). Another 'Input Dataset' (output) feeds into 'FASTQ Groomer' (File to groom, output_file (fastq, fastqcssanger, fastqsolexa, fastqillumina)). The output of the second 'FASTQ Groomer' feeds into 'Cutadapt' (Fastq file to trim, report (txt), output, rest_output, too_short_output, untrimmed_output). The right sidebar shows the configuration for the 'FASTQ Groomer' tool, including version 1.0.4, file to groom options, and input FASTQ quality scores type set to 'Sanger & Illumina 1.8+'.

A l'aide des briques disponibles à gauche de l'interface, réaliser un workflow qui reproduit l'enchaînement des briques de l'exercice précédent.

This screenshot shows a close-up of the 'FastQC' tool in the workflow canvas. A context menu is open over the tool, displaying the following options: 'Save', 'Run', 'Edit Attributes', 'Auto Re-layout', and 'Close'. The 'FastQC' tool configuration is visible, showing 'FASTQ reads', 'Contaminants', and 'report (html)' as outputs.

- Sauvegarder votre workflow en cliquant sur « Option » puis « Save ».
- Lancez votre workflow en partant du fichier récupéré dans les "shared data" et envoyez les résultats dans un nouvel historique.

Workflow sur fichiers multiples:

Galaxy offre la possibilité de lancer un même workflow sur plusieurs fichiers à la fois, en cliquant sur la petite icone à droite de "Input dataset".

Running workflow "formation_13nov2014" Expand All Collapse

Step 1: Input dataset

Input Dataset

- 4: input1.fastq
- 5: input2.fastq
- 8: Fastq2.sanger
- 9: fastq1.sanger
- 35: fastq1.cleaned
- 37: fastq2.cleaned

type to filter, [enter] to select all

Step 2: FASTQ Groomer (version 1.0.4)

Step 3: Cutadapt (version 0.9.5.a)

Lancer le workflow sur les échantillons RC1, RC2 et RC3 importés depuis la librairie partagée.

1-6 Publication de workflow

Galaxy permet de partager un workflow facilement avec la communauté.

Galaxy Analyze Data Workflow Shared Data Visualization Admin Help User Using 54.3

Your workflows Create new workflow Upload or import workflow

Name	# of Steps
formation_13nov2014	3
imported: ... analysis workflow	7
imported: ...	5
GnpAsso_v ... 3D+Admixture+Tassel_MLM	5
GnpAsso_v ...	5
Workflow c ...	7
ESTik_from_contigs	6
imported: ESTtik	15

Il y a 3 moyens de partager un workflow :

- * Via un lien : crée un lien de partage pour permettre à vos contacts d'importer le workflow.
- * Par publication : rend le workflow public et accessible en en le publiant dans la section «Published Workflows».
- * Par email : partage le workflow avec un autre utilisateur de galaxy.

Partagez votre workflow par email avec l'utilisateur de votre choix.

2 – En roues libres...

Exemples d'analyses possibles :

- Recherche de microsatellites avec définition des primers : Suite d'outils SAT

Galaxy

Tools

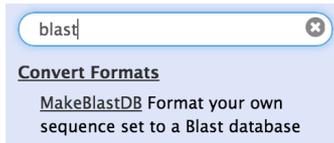
- SNIPlay
- Filter and Sort
- Gene/Protein prediction
- Population Analysis
- SAT**
- SSRIT finds all perfect SSRs in a set of sequences
- Primer3 Batch Pick primers from a set of DNA sequences and a list of targets

- Réaliser un alignement blast d'un fichier contenant plusieurs séquences contre sa propre banque de séquences:

1. Importer sa propre banque de séquences



2. Créer les fichiers index nécessaires pour réaliser une recherche de similarité avec blast avec MakeBlastDB



3. Réaliser la recherche de similarité



- Recherche de SNP (à partir d'un fichier de type bam) :
 - a. Marquer les duplicats techniques **Picard Mark Duplicate reads**
 - b. Réaligner les reads à l'aide des outils de la suite **GATK Realigner Target Creator** et **Indel Realigner**
 - c. Détecter les polymorphismes avec l'outil **GATK Haplotype Caller**

- Expression Différentielle (HT-seq count, EdgeR et DESeq)