

# Formation sRNAseq Analyse des miRNAs sous Galaxy

# - EXERCICES -

# Babraham Bioinformatics

# cutadapt

BWA

SAMtools

"FastQC is a quality control tool for high throughput sequence data." http://www.bioinformatics.bbsrc.ac.uk/

A tool that removes adapter sequences from DNA sequencing reads.

*cutadapt* removes adapter sequences from high-throughput sequencing data. <u>https://code.google.com/p/cutadapt/</u>

"Burrows-Wheeler Aligner (BWA) is an efficient program that aligns relatively short nucleotide sequences against a long reference sequence such as the human genome." <u>http://biobwa.sourceforge.net</u>

"SAM (Sequence Alignment/Map) format is a generic format for storing large nucleotide sequence alignments." <u>http://samtools.sourceforge.net</u>



## **Objectifs**:

Cette formation a pour objectif de vous aider à traiter les séquences issues de projet de sRNAseq (miRNA). Vous y découvrirez les problématiques spécifiques de l'analyse des petits ARNs non codant, les outils liés et les mettrez en œuvre afin de détecter, annoter, prédire, quantifier, ... les miRNA.

Pour vous connecter à Galaxy sur l'instance de Roscoff:

- Instance de Roscoff : <u>http://galaxy.sb-roscoff.fr/</u>
- Si besoin, instance de Toulouse : <u>http://sigenae-workbench.toulouse.inra.fr</u> (demande de compte sur <u>http://bioinfo.genotoul.fr/index.php?id=74</u>)
   Comptes de formation suivants :
  - Logins : anemone , aster, bleuet, iris, muguet, narcisse, pensee, rose, tulipe, violette, lilas, pervenche, laurier, lavande, lis, capucine, coquelicot, geranium, liseron, arome, chardon.
  - Password : f1o2r3!

L'ensemble des outils utilisés dans ce TP sont regroupés dans la section « miRNA et sRNAseq » du menu « Tools ».



### **Exercice n°1**: Prise en main de votre instance Galaxy

#### 1 - Notion de « dataset » et de « tool »

Un fichier téléchargé dans Galaxy se nomme « dataset » et est visible à droite de votre interface Galaxy, dans la fenêtre « history ».

Nous allons travailler à partir de données Illumina représentant 2 tissus. Ces jeux de données « <u>s2.fastq</u> » et « <u>s1.fastq</u> » peuvent être téléchargés dans Galaxy à partir de ce lien : <u>http://snp.toulouse.inra.fr/~sigenae/Galaxy\_Formation/miRNA/Roscoff/</u>ou directement en récupérant les fichiers dans l'instance Galaxy de Roscoff : shared data -> data libraries -> tp-jeudi-mirna-sarah .

Pour que les fichiers « <u>s2.fastq</u> » et « <u>s1.fastq</u> » soient disponibles dans Galaxy pour vos traitements, veuillez utiliser l'outil « <u>Upload File</u> from your computer » en précisant les URL des fichiers à télécharger. Cet outil télécharge en copiant votre fichier sur le serveur Galaxy. Cette copie diminue donc votre quota sur Galaxy.



Si l'extension de votre fichier est « .fastq », alors ce fichier ne sera pas visible par les outils Galaxy. Pour ce faire, il est nécessaire de modifier l'extension de « fastq » en « fastqsanger » via l'icône « pencil » de votre dataset, menu « change datatype ».

Veuillez renommer les fichiers importés :

Renommer « http://snp.toulouse.inra.fr/~sigenae/Galaxy\_Formation/miRNA/Roscoff/s1.fastq » en « s1 ». Renommer « http://snp.toulouse.inra.fr/~sigenae/Galaxy\_Formation/miRNA/Roscoff/s2.fastq » en « s2 ». Puis changer le datatype de ces deux fichiers, de « fastq » à « fastqsanger ».

#### 2 - Création et partage d'historique et de workflow

#### 2.1 – Archivage automatique de vos traitements dans un historique

Au fur et à mesure que vous faites appel aux différents outils utiles au sein de votre interface depuis le menu « Analyse Data », l'ensemble des étapes sont enregistrées dans un historique qui est automatiquement archivé dans « User / Saved Histories » et que vous pouvez ensuite, si besoin, partager dans « Shared Data / Published Histories ».

Analyze Data Wo	rkflow Sh	ared Data	Visualization	Admin	Help	User	Welcome smamar	n, you are working	in /work	Using 17%
Tools	Options 🔻				-	Logge	d in as smaman@to	ulouse.inra.fr <sub>His</sub> t	tory	Options 🔻
Your user name: smaman Your file path : /work/smar	man/		The following job	has been	success	Jogou	t	ی (ت	-	<i>a</i> = 1
1 - UPLOAD YOUR DATA			120: Add colur	nn on data	111	<u>Saved</u>	Histories	TP :	1 - GALAXY	728.2 Mb
Get Data			You can check th	e status of	foueueo	Saved	Datasets	<u>12</u>	0: Add	● ℓ ×
2 - FILES MANIPULATION	N E		refreshing the <b>Hi</b> change from 'run	story pan	e. When	Public		6 li	nes	
Text Manipulation			problems were e	ncountere	d.			dat	mat: tabular, tabase: ?	·
<ul> <li>Add column to an existin dataset</li> <li>Compute an expression</li> </ul>	ng							Info at I CE	o: Epilog : jol lun. sept. 24 ST 2012	b finished 16:03:22 ≣

Depuis le menu « User » / « Saved Histories », vous avez la possibilité de gérer vos historiques (delete, delete permanently, rename, undelete) en cliquant sur l'intitulé de l'historique.

A chaque connexion à votre instance Galaxy, un nouveau « current history » est automatiquement créé.

Tools Or Linux Username: sigenae File Path : /work/sigenae/gi	alaxy/	Saved Histories search history names and tags	L.		-717		
Get Data Send Data ENCODE Tools		Advanced Search	Datasets	Tags Sharing	Size on Disk	Created	Last Updated ↑
Lift-Over		Unnamed history		0 Tags	0 bytes	less than a minute ago	less than a minute ago
Text Manipulation Filter and Sort		RNA seq statistics	2 4	<u>0 Tags</u>	34.9 Mb	~ 6 hours ago	6 minutes ago
Join, Subtract and Group Convert Formats		Test BWA fichiers Gnome	4	<u>0 Tags</u>	9.0 Mb	Apr 06, 2012	1 day ago
Extract Features		Test region promoters *	5	0 Tags	23.9 Mb	Mar 08, 2012	Apr 06, 2012
Fetch Sequences Fetch Alignments		Unnamed history •	з	0 Tags	60 bytes	Feb 23, 2012	Mar 09, 2012
Get Genomic Scores Operate on Genomic Inter	vals	Unnamed history -	1 1	<u>0 Tags</u>	0 bytes	Mar 07, 2012	Mar 07, 2012
<u>Statistics</u> Wavelet Analysis		Unnamed history	1 1	<u>0 Tags</u>	16.0 Kb	Feb 22, 2012	Mar 05, 2012
Graph/Display Data		For 2 selected histories: Rename	Delete Delete Permar	ently Undelete			
Regional Variation Multiple regression		Histories that have been deleted for more	than a time period specified	by the Galaxy administral	tor(s) may be permane	ently deleted.	

Veuillez nommer votre historique « miRNAs »,

S 200

#### 2.2 – Partage de vos analyses grâce aux workflows



Pour lancer plusieurs jobs en une fois, vous pouvez générer un workflow à partir de votre historique (**History** panel / click **Options** → **Extract Workflow**) ou bien créer directement un workflow depuis une page blanche.

Veuillez générer un workflow depuis votre historique « miRNA » et le nommer « miRNA-Votre/user/login »,

Puis partager le workflow que vous venez de générer avec un autre utilisateur Galaxy présent dans la salle ("Share with a user"). Pour cela, vous avez besoin de connaître son nom d'utilisateur sous Galaxy.

Faire un clone du workflow partagé.



Exercice n°2: Analyse de la qualité et nettoyage (suppression des adaptateurs suppression de la redondance intra « fastq »)

#### 1 - Analyse de la qualité des données : « FastQC »

A partir du jeu de données s2, veuillez utiliser l'outil Galaxy « Fastqc: Fastqc QC using FastQC from Babraham » pour générer les graphiques d'analyse de la qualité.

Fastqc: Fastqc QC (version 0.4)	En cliquant sur l'icône «œ statistiques s'affichent dans la	il », les fenêtre
Short read data from your current history:	centrale :	
38: Suppression des a on data 25 👻	History Options	; 🔻
Title for the output file - to remind you what the job was for:	o 🗆 🖉 🖻	*
FastQC	TEST miRNA 24.1 Gb	
Contaminant list:	44: FastQC data 21.html	
Selection is Optional 🚽	43.4 Kb format: html, database: ?	
tab delimited file with 2 columns: name and sequence. For ex CAAGCAGAAGACGGCATACGA	Info: Epilog : job finished at lun. oct. 15 14:18:24 CEST 2012 	=
Execute	HTML file	

Il vous est possible de récupérer :

• L'ensemble des résultats de FastQC en cliquant sur l'icône « disquette » du dataset.



Remarque : Pour les personnes connectées à l'instance Galaxy de Roscoff (uniquement), l'ouverture des rapports html générés dans le cadre de ce TP s'effectue via un click droit sur l'icône "save" de votre dataset, "enregistrer sous", puis ajouter ".html" au nom du fichier.

Un graphique par un clic droit sur son intitulé, en fin de page :

#### Files created by FastQC

duplication\_levels.png (19.9 KB) fastqc\_data.txt (24.2 KB) fastqc\_report.html (42.1 KB) kmer\_profiles.png (105.8 KB) per\_base\_gc\_content.png (31.4 KB) per\_base\_n\_content.png (7.9 KB) per\_base\_quality.png (9.2 KB)



#### 2 – Nettoyage des données :

#### 2.1 – Suppression des adaptateurs

La suppression des adaptateurs s'effectue grâce au script « cutadapt » disponible dans Galaxy « <u>\* Suppression des adaptateurs</u> ». Veuillez lancer ce traitement sur « s2 ».

Paramètres :

- adaptateur : ATCTCGTATGCCGTCTTCTGCTTG
- taille minimum : 18pb
- taille maximum : 25pb

**Remarques**:

- Attention, la taille minimum des séquences requise pour mirdeep2 est de 18 pb.
- Un astérisque \* devant le nom de l'outil signifie qu'il s'agit d'un module Galaxy rajouté par Sigenae.

Sorties de « cutadapt » :

- fichier « fastq » contenant les séquences nettoyées de leurs adaptateurs
- fichier « log » contenant les informations de traitement

Veuillez renommer le fichier sortant «Suppression des adaptateurs on data 5 » en «cut\_s2\_log » pour le fichier de log et en « cut\_s2\_ sequences» pour le fichier fastq.

#### 2.2 – Statistiques après suppression des adaptateurs

Lancer l'outil « <u>*Production du rapport</u> après élimination							
des	adaptateurs	»,	sur	«cut_s2_log »	et		
« cut_s2_ sequences», pour produire d'un rapport de							
traitement (au format html).							

Veuillez renommer le fichier sortant « Votre rapport apres suppression des adaptateurs « en « Votre rapport apres suppression des adaptateurs sur s2 ».



Pour visualiser ces statistiques, veuillez enregistrer le fichier sur votre PC et ouvrir la page « html » dans votre navigateur.

30004		Sequence n Source: Sigenae/C	umber by "fastq" enotoul miRNA pipeline			Length distribution Source: Sigenae/Genotoul m/RNA pipeline
2500k -				Cutadapt	500k	dataset_1921.dat
admin a scook	ß				300k	
0 1500k					200k —	
500k -					100k	
Ok -			dataset_1921.dat		0k	18 19 20 21 22 23 24 25 26 27 28

Remarque : Pour les personnes connectées à l'instance Galaxy de Roscoff (uniquement), l'ouverture des rapports html générés dans le cadre de ce TP s'effectue via un click droit sur l'icône "save" de votre dataset, "enregistrer sous", puis ajouter ".html" au nom du fichier.

Philippe Bardou - Jérôme Mariette - Christine Gaspin - Olivier Rué - Sarah Maman



Afin de pouvoir relancer ce workflow sur l'ensemble de vos « fastq », veuillez convertir l'historique que vous venez de générer en workflow « miRNA-WF1qualite\_et\_nettoyage\_des\_donnees ».

Pour chaque jeu « fastq », il vous suffira alors de relancer ce WF1.

Relancer votre workflow sur le « fastq » s1 après avoir supprimé les deux boîtes d'upload des vos jeux de données (car ils sont déjà dans votre historique) et avoir vérifier que les outils sont bien paramétrés. Il est préférable de ne pas envoyer les fichiers générés par ce workflow vers un nouvel historique.

	Fastqc: Fastqc QC	3
* Upload local file from 🛛 🗱	Short read data from your current history	
out1 (bam, txt, fastqsanger, csfasta,	Contaminant list	
Jal, bed, gff, gtf, vcf, sam, fasta, O 🔿	html_file (html)	
	* Suppression des adaptateurs	3
	Votre fichier 'lane' Illumina au format fastqsanger	
	output_fq (fastq.cut.fq)	* Production du rapport 🛛 🕱
	output_log (fastq.cut.log)	Votre fichier fq issu de cutadapt
		Votre fichier log issu de cutadapt
		output_html_report (html)
		FASTQ to FASTA 💥

Pour plus de lisibilité, veuillez renommer les datasets suivantes :

Nom proposé par Galaxy	Renommage		
* Suppression des adaptateurs on data 4 (format log)	cut_s1_ log		
* Suppression des adaptateurs on data 4 (format fastq)	cut_s1_ sequences		
* Votre rapport apres suppression des adaptateurs	Votre rapport apres suppression des adaptateurs sur s1		



### Exercice n°3 : Recherche de miRNAs avec mirdeep2

Le WF2 « miRNA-WF2-Mirdeep2 » exécute mirdeep2 à partir des fichiers générés par le l'outil « Suppression des adaptateurs » du WF1.



La première étape de mirdeep2 (« mapper ») permet d'aligner les lectures (fichier « fasta » précédent) sur le génome de référence passé en paramètre (V4\_454Scaffolds\_filter) : outil « \* <u>Process and map reads to the genome ».</u>

La seconde étape de mirdeep2 (« core ») permet d'annoter les miRNAs appartenant à des familles d'ARN connus : outil « <u>\* Annotation des miRNAs</u> ».

Voici vos étapes de traitement :

- Au préalable, il est nécessaire de convertir vos fichiers « fastq » (« cut\_s1\_sequences » et «cut\_s2\_ sequences ») en fichiers « fasta » à l'aide de l'outil « FASTQ to FASTA ». Veuillez renommer les fichiers sortants « [fastq -> fasta] Output File » en « s1\_fasta » et « s2\_fasta ».
- 2. Importation du WF2 intitulé « Galaxy-Workflow-miRNA-WF2-Mirdeep2.ga » sur <u>http://snp.toulouse.inra.fr/~sigenae/Galaxy\_Formation/miRNA/Roscoff/workflows/Galaxy-</u> <u>Workflow-miRNA-WF2-Mirdeep2.ga</u>
- Visualiser (clic sur l'intitulé du workflow → Edit) le workflow importé pour sélectionner la banque V4\_454Scaffolds\_filter pour les 3 outils puis enregistrer ces modifications dans WF2.
- 4. Lancement du workflow WF2 (clic sur l'intitulé du workflow → Run) sur les fichiers entrants « s1\_fasta » et « s2\_fasta », dans l'historique en cours.
- 5. Visualisation des résultats : Veuillez observer les premières lignes du fichier « fasta » généré par l'outil « mapper » et les « pdfs » générés par l'outil d'annotation de mirdeep2.
- 6. Pour plus de lisibilité, veuillez renommer les datasets suivantes :

Nom proposé par Galaxy	Renommage
Fichier fasta généré par l'outil mapper de « Mapper : Process and map reads to the genome. on data 16 and data 18 »	« Mapper on s1 and s2 fasta »
Fichier arf généré par l'outil mapper de « Mapper : Process and map reads to the genome. on data 16 and data 18 »	« Mapper on s1 and s2 arf»
Fichier fasta généré par le convertisseur « * miRDeep2core - bed to fasta file on data 24 ».	« miRDeep2core fasta »





### Exercice n°4 : Recherche des annotations fonctionnelles

Quelques liens (outils / base de données) :

- BWA : <u>http://bio-bwa.sourceforge.net</u>
- BWA man : <u>http://bio-bwa.sourceforge.net/bwa.shtml</u>
- SAMtools : <u>http://samtools.sourceforge.net</u>
- mirBase : <u>ftp://mirbase.org/pub/mirbase/CURRENT/</u> (hairpin.fa.gz )
- Rfam : <u>ftp://ftp.sanger.ac.uk/pub/databases/Rfam/CURRENT/</u> (Rfam.fasta.gz )
- tRNA : <u>http://gtrnadb.ucsc.edu/download.html</u> (eukaryotic-tRNAs.fa.gz )
- rRNA : <u>ftp://ftp.arb-silva.de/current/Exports/</u> ([LS]SURef\_108\_tax\_silva\_trunc.fasta.tgz)

Dans le cadre de ce TP, les banques de référence sont déjà disponibles dans le menu déroulant de l'outil BWA.

L'objectif est de comparer les diagrammes de Venn obtenus après le mapper avec ceux obtenus après le mirdeep2core.

Voici vos étapes de traitement :

 Annoter le « fasta » issu de Mapper (« Mapper on s1 and s2 fasta ») puis le « fasta » issu de miRDeep2core (« miRDeep2core fasta ») avec l'outil « <u>\* Alignement bam et tri</u> sur un fasta ».

Chaque « fasta » doit être annoter 3 fois, une fois contre la banque mirBase (hairpin\_T), puis contre la banque Rfam puis contre la banque tRNA. Chaque alignement doit être filtré selon ces paramètres : Flag 0 ou 16 – 1 mismatch autorisé. Le filtre est nécessaire car il n'est pas possible avec BWA de spécifier qu'on ne veut que les alignements avec un seul mismatch. Les options 0 et 16 permettent de garder les alignements obtenus en forward et en reverse.

2. Pour plus de lisibilité, veuillez renommer les datasets suivantes :

Nom proposé par Galaxy	Renommage
<ul> <li>* Alignement baw, tri et filtre on data 21 » -</li> <li>Fichier au format filter1, sur hairpin_T, issu du</li> <li>« Mapper on s1 and s2 fasta »</li> </ul>	«mapper hairpin_T filter1 » A décliner pour Rfam et tRNA.
<ul> <li>* Alignement baw, tri et filtre on data 21 » -</li> <li>Fichier au format filter2 sur hairpin_T, issu du</li> <li>« Mapper on s1 and s2 fasta »</li> </ul>	«mapper hairpin_T filter2 » A décliner pour Rfam et tRNA.
<ul> <li>* Alignement baw, tri et filtre on data 21 » -</li> <li>Fichier au format filter1, sur hairpin_T, issu du</li> <li>w miRDeep2core fasta »</li> </ul>	«miRDeep2core hairpin_T filter1 » A décliner pour Rfam et tRNA.
<ul> <li>« * Alignement baw, tri et filtre on data 21 » -</li> <li>Fichier au format filter2 sur hairpin_T, issu du</li> <li>« miRDeep2core fasta »</li> </ul>	« miRDeep2core hairpin_T filter2 » A décliner pour Rfam et tRNA.

3. Puis lancer l'outil « <u>\* Comparaison des annotations</u> », à partir des 3 fichiers au format « filter2» générés par l'alignement. L'objectif est de construire un diagramme de Venn.



- 4. Puis lancer l'outil « <u>\* Construction de la matrice d annotations</u> », à partir des 3 fichiers au format « filter1 » générés par l'alignement.
- 5. Comparer les diagrammes de Venn selon que les traitements aient été lancés à partir du « fasta » issu du Mapper ou du « fasta » issu de miRDeep2core.



Cette suite d'outils « miRNA et sRNAseq » est disponible dans le Galaxy Tool Shed (<u>http://toolshed.g2.bx.psu.edu/</u>) avec le mot clés « miRNA » :

- Galaxy Tool Shed	Repositories H	elp <del>+</del> Use		marque poges
2430 valid tools on Jan 11, 2013	Repositories			
Search Search for valid tools	miRNA 🗙 search repository n	ame, descri	ption 🔍	
Search for workflows	<u>Name</u> ↓	<u>Synopsis</u>	Metadata Revisions	Tip Revis
All Repositories   Browse by category  My Repositories and Tools  Repositories I own  My writable repositories  My invalid tools	mirdeep2_and_targetspy	Finding miRNA in NGS data and finding the targets for those miRNA	1:798fe7ba8b5e	1:798fe7b
Available Actions     Create new repository	mirna_mirdeep2	miRNA	3:12ab33cb3511	3:12ab33c