

Analyse de données RADseq sous Galaxy : l'exemple de STACKS

Yvan Le Bras, Projet e-Biogenouest, CNRS UMR 6074 IRISA-INRIA, Rennes

Cyril Monjeaud, Projet GenoCloud, CNRS UMR 6074 IRISA-INRIA, Rennes

Avec l'utilisation de ressources réalisées par Julian Catchen, Institute of Ecology and Evolution, University of Oregon

Le but de ces exercices est de familiariser les stagiaires avec l'utilisation des données générées par les séquenceurs de nouvelles générations à partir de RRL (Reduced Representation Libraries) comme les tags associés à des sites de restriction (RAD). Ces librairies sont souvent utilisées dans le cadre du génotypage par séquençage, et peuvent fournir un jeu dense de marqueurs SNP (single nucleotide polymorphism) répartis le long du génome. Les stagiaires acquerront de l'expérience avec un pipeline d'analyse nommé STACKS, créé pour l'analyse de ce type de données. Les données à analyser proviendront d'un organisme sans génome de référence pour générer une cartographie génétique d'une part et identifier de potentielles signatures de sélection. Il est également possible d'utiliser un organisme avec génome de référence.

Les participants vont apprendre à:

1. Préparer les données brutes Illumina RAD pour leur analyse en enlevant les lectures de mauvaise qualité et pour démultiplexer un jeu d'échantillons barcodés.
2. Aligner des séquences RAD contre un génome de référence
3. Utiliser Stacks pour assembler les loci RAD, détecter des SNPs, les génotypes et haplotypes pour chaque individu de deux populations.
4. Calculer des statistiques en génétique des populations

A la fin de cette formation, vous devriez savoir:

5. Manipuler les données brutes Illumina de RAD pour les analyser en utilisant une variété de différents paramètres.
6. Aligner les tags RAD contre un génome de référence pour identifier des signatures potentielles de sélection.
7. Etendre ce qui a été appris vers des problèmes plus complexes, les vôtres.

Nous allons utiliser des jeux de données proposés par Julian dans ces formations plus des jeux de données épinoche de l'article d'Hohenlohe et al. 2010. Pour le nettoyage et l'analyse des données, nous nous reposerons principalement sur Galaxy, le pipeline STACKS et BWA.

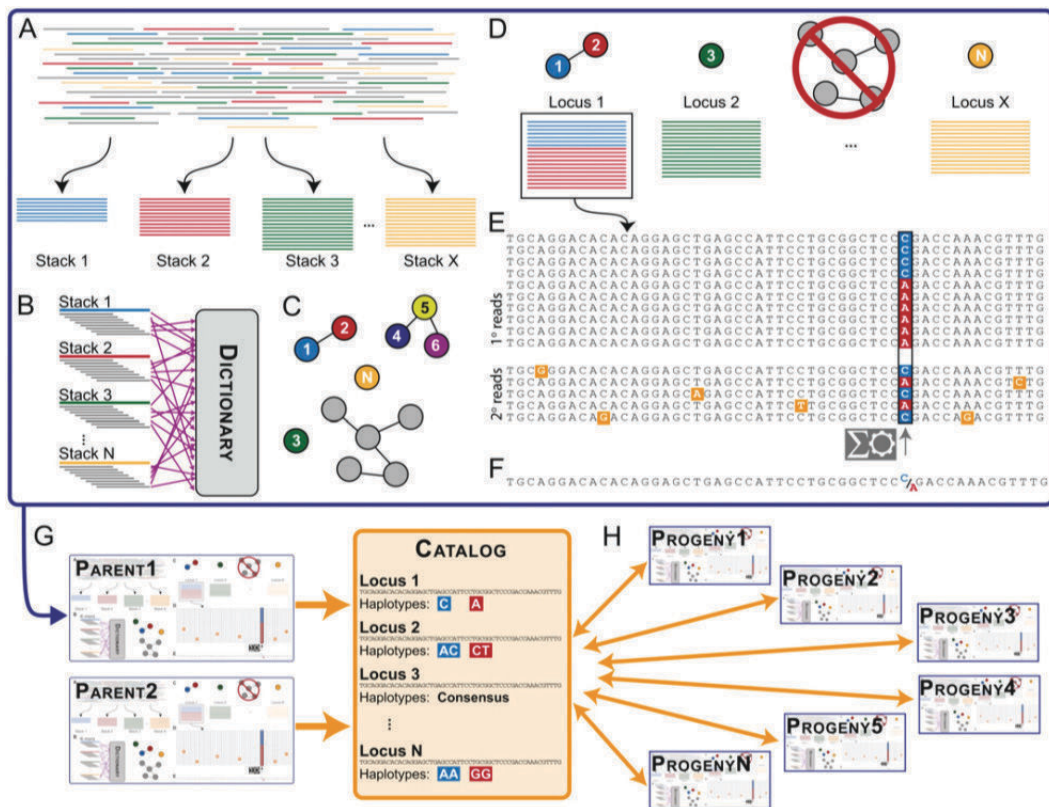
Les jeux de données seront tous produits via des séquenceurs de type Illumina GAI ou HiSeq2000.

Les logiciels sont tous open source

- **BWA** (<http://bio-bwa.sourceforge.net/>) - **BWA** est utilisé pour aligner des séquences contre un génome de référence. Nous l'utiliserons pour aligner les lectures RAD contre le génome de l'épinoche, puis pour analyser ces lectures via le pipeline Stacks. Si nous utilisons BWA ici, plusieurs autres algorithmes et logiciels existent pour effectuer cette tâche (comme Bowtie) et ils peuvent également être utilisés pour cette partie. Nous avons cependant développé un outil spécialement pour STACKS sous Galaxy, basé sur BWA.

- **Stacks** (<http://creskoloab.uoregon.edu/stacks/>) – il s'agit d'un ensemble de programmes open source interconnectés initialement mis en place pour l'assemblage de novo de séquences RAD en loci et cartes génétiques, aujourd'hui étendu pour être utilisé de manière plus flexible dans des études d'organismes présentant ou non un génome de référence. Le pipeline a un wrapper Perl permettant le lancement de l'ensemble des programmes. La modularité de STACKS lui permet d'être appliqué à différents types de scénarios.

A noter que le pipeline STACKS utilisant un génome de référence suit les mêmes étapes. La principale différence s'observe pendant la phase de construction des loci. **Ref_map** construit alors les loci à partir du génome de référence, en utilisant les résultats de mapping des lectures sur ce dernier (réalisé via Bowtie, BWA ou un autre logiciel de mapping). Cette étape se nomme alors **pstacks** et se déroule à la place de **ustacks**. **cstacks** et **sstacks** sont ensuite exécutés.



Il exécutera chaque composant de Stacks. En premier, il exécutera (A-F) **ustacks** sur chacun des échantillons, construisant les loci et détectant les SNPs à chacun d'entre eux. Les séquences présentant des correspondances exactes sont regroupées en piles ("stacks") (A). Un nombre de lectures trop faible dans une pile pouvant provenir d'erreur de séquençage, les piles uniques contenant moins de lectures que le seuil spécifié (*stack depth parameter*) sont désassemblées et les lectures mises de côté. Les lectures finalement conservées dans une pile sont nommées lectures

primaires. Celles mises de côté en raison d'un nombre de lectures trop faible par pile sont nommées lectures secondaires.

A la fin de cette première étape de création de piles, **ustacks** calcule la moyenne de profondeur de couverture et identifie les piles présentant une profondeur supérieure à 2 écarts types au-dessus de la moyenne. Toute les piles dans ce cas de figure (nommées piles bucheronnes), ainsi que celles présentant une séquence proche au nucléotide près, sont exclues car souvent représentées par des éléments répétés (voir cas de figure des 5 piles grises reliées (C et D)).

Ensuite, des sous-ensembles de piles sont produits (C), lorsque des piles sont très proches, un nucléotide de différence. Les piles de chaque sous-ensemble peuvent alors être réunis en un seul locus, comme c'est le cas pour les piles 1 et 2 (D). Ensuite, les lectures qui étaient initialement mises de côté, car présentant une similarité trop faible avec les séquences des piles (polymorphisme sur plus d'un nucléotide) (E), sont comparées à celles constitutives des piles pour identifier si elles peuvent être rattachées à un unique polymorphisme identifié dans les étapes précédentes. Enfin, une séquence consensus est établie (F). En second, (G) **cstacks** sera exécuté pour créer un catalogue de tous les loci à partir des parents du croisement. Finalement, (H) **sstacks** s'exécutera pour évaluer la concordance entre les loci de chaque descendant et le catalogue de loci. Le pipeline identifie ensuite les loci représentant des marqueurs cartographiables (exécution du module genotypes).

En travaillant sur des données NGS, vous noterez probablement que toutes les données générées par le séquenceur ne sont pas de bonne qualité. En général, il faudra enlever les séquences de faibles qualités de vos jeux de données avant de les analyser. En même temps, la stringence de la filtration dépendra de l'application finale. En général, une stringence plus forte est appliquée pour un assemblage de novo comparé à l'utilisation d'alignements contre un génome de référence. Par contre, des données de mauvaise qualité affecteront presque toujours les analyses futures, en produisant des faux positifs, comme par exemple la prédiction de faux SNP.

I. Analyse RAD-seq sous Galaxy : Détection de SNPs

1. L'analyse

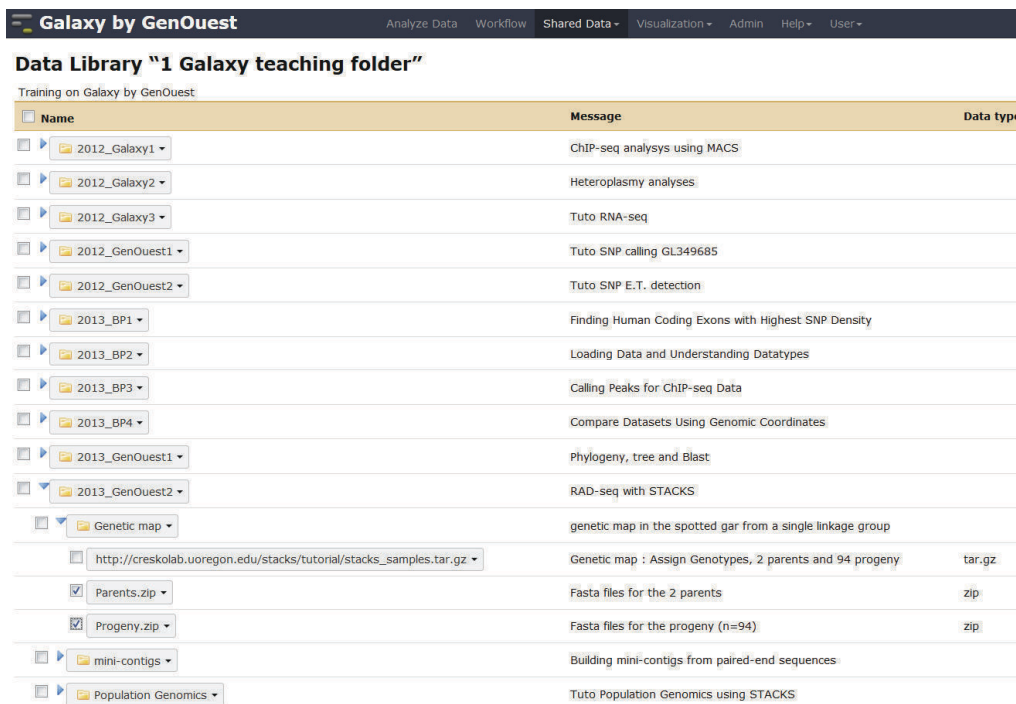
Nous travaillerons sur les données issues de l'archive `pe_sample.tar.gz` disponible sur le site de Stacks (http://creskolab.uoregon.edu/stacks/pe_tutorial/pe_samples.tar.gz). Sous Galaxy, commencez par créer un nouvel historique (ex: "RAD 1 : SNP calling").

Récupération des données brutes dans Shared data/data libraries/1 Galaxy teaching folder/2013_GenOuest2/genetic_map

Vous y trouverez 3 archives :

- une contenant les 95 jeux de données de séquences : `stacks_samples.tar.gz`
- une regroupant les 2 jeux de données de séquences des parents : `Parents.zip`
- une regroupant les 93 jeux de données de séquences des descendants F1 : `Progeny.zip`

Sélectionner les 2 dernières archives `.zip` (**Parents** et **Progeny**), et cliquer sur le bouton "GO" après avoir vérifié que l'action sélectionnée était bien "Import to current history".



Name	Message	Data type
2012_Galaxy1	ChIP-seq analysis using MACS	
2012_Galaxy2	Heteroplasmy analyses	
2012_Galaxy3	Tuto RNA-seq	
2012_GenOuest1	Tuto SNP calling GL349685	
2012_GenOuest2	Tuto SNP E.T. detection	
2013_BP1	Finding Human Coding Exons with Highest SNP Density	
2013_BP2	Loading Data and Understanding Datatypes	
2013_BP3	Calling Peaks for ChIP-seq Data	
2013_BP4	Compare Datasets Using Genomic Coordinates	
2013_GenOuest1	Phylogeny, tree and Blast	
2013_GenOuest2	RAD-seq with STACKS	
Genetic map	genetic map in the spotted gar from a single linkage group	
http://creskolab.uoregon.edu/stacks/tutorial/stacks_samples.tar.gz	Genetic map : Assign Genotypes, 2 parents and 94 progeny	tar.gz
<input checked="" type="checkbox"/> Parents.zip	Fasta files for the 2 parents	zip
<input checked="" type="checkbox"/> Progeny.zip	Fasta files for the progeny (n=94)	zip
mini-contigs	Building mini-contigs from paired-end sequences	
Population Genomics	Tuto Population Genomics using STACKS	

Vous vous retrouvez avec un historique contenant les deux archives au format zip.

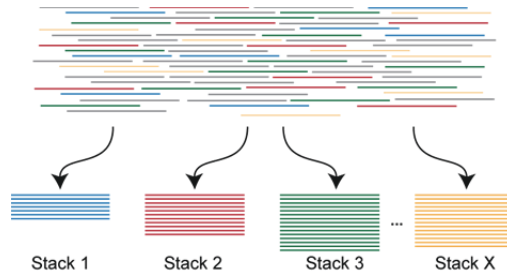
Les données ayant déjà été nettoyées et démultiplexées, il ne sera pas nécessaire d'exécuter **STACKS** : **Process Radtags**.

Nous sommes maintenant prêts à exécuter le pipeline **STACKS denovo**.

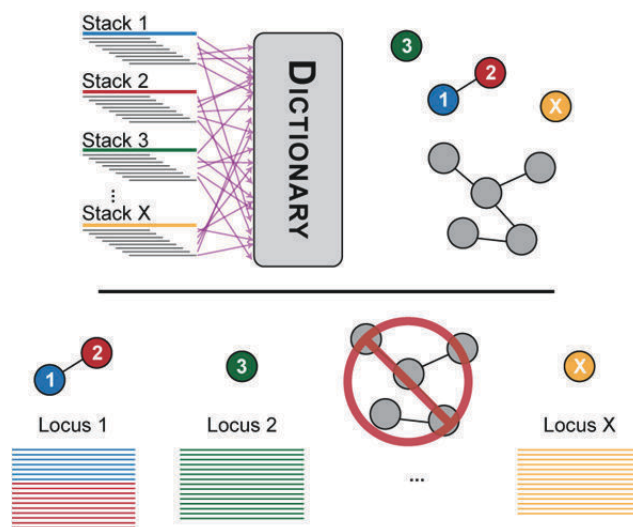
Exécuter **STACKS : de novo map** sur les individus étudiés (par exemple les parents) en prenant comme input l'archive .zip et en spécifiant l'usage "*Genetic map*". Décocher l'option "*create zip outputs*"

Vous pouvez préciser certains paramètres. Ainsi, en sélectionnant le mode Advanced pour le paramètre **Stack assembly options**, vous pouvez préciser:

-la profondeur minimum de la pile (-m), ici "Minimum number of identical". Ce paramètre, passé à **ustacks**, contrôle le nombre de lectures avec une correspondance exacte devant être trouvé pour créer une pile chez un individu. Nous pouvons sélectionner 3, une pile ne sera générée chez un individu que si au moins 3 lectures correspondent exactement. Les lectures alors utilisées pour constituer ces piles sont qualifiées de lectures primaires. Les autres, qui ne correspondent donc pas exactement à des lectures primaires et qui n'étaient pas assez nombreuses pour générer une pile, sont qualifiées de lectures secondaires.



-la distance maximum permise entre des piles pour qu'elles soient fusionnées en un locus potentiel chez un individu. Il s'agit du paramètre -M, ici "Number of mismatches allowed between loci when processing a single individual". Nous précisons ici une distance maximale de 2 nucléotides. En reprenant le schéma ci-dessous, les piles 1 et 2 vont être fusionnées car représentent un seul locus polymorphe (avec 2 allèles différents possédant une différence maximale correspondante au seuil fixé, ici 2), les piles 3 et X vont être associés à deux loci monomorphiques distincts, le gros paquet gris de piles représente un jeu de séquences répétées qui présentent trop d'allèles pour être biologiquement corrects.



-le nombre maximal de différence entre loci du catalogue. Il s'agit du paramètre `-n`, "specify the number of mismatches allowed between loci when building the catalog" ici. Ce paramètre permet notamment de pouvoir créer un locus de type homozygote dans le catalogue alors qu'il est en réalité hétérozygote. Ainsi, si par exemple on fixe la valeur à 2 et que la séquence du locus ne présente pas de polymorphisme au sein de chaque individu mais présente plus de 2 nucléotides de différences entre chaque individu, alors chaque "version" du locus identifié générera un locus distinct des autres. Ce paramètre est notamment pratique pour conserver les loci hétérozygotes entre parents mais homozygote chez chacun d'entre eux. Si la valeur de ce paramètre est trop importante, cstacks risque de considérer des loci distincts en un seul locus. Si cette valeur est trop faible, cstacks générera potentiellement plus de loci qu'il n'y en a réellement. Nous pouvons fixer ce paramètre à 3.

-supprimer ou casser les RAD-tags très répétitifs. Il s'agit de supprimer les piles "bûcheronnes" de l'analyse et briser les piles modérément surdimensionnées. Nous cocherons ici cette option, correspondant au paramètre `-t`.

D'une manière générale, utiliser des options plus stringentes (comme la suppression des RAD-tags très répétitifs) va générer un nombre total de piles plus faible, augmentera le nombre de piles réunis (car manque des piles à n nucléotides près faisant le pont entre deux piles) et diminuera le nombre final de loci du catalogue. (car moins de similarité, plus de bruit de fond, entre les piles des individus).

The screenshot shows the Galaxy web interface. The main panel displays the configuration for the 'Stacks : De novo map (version 1.0.0)' tool. The 'Select your usage' dropdown is set to 'Genetic map'. Under 'Files containing parent sequences', there are three input fields: '1: parents.zip', '2: progeny.zip', and '3: snps_output.zip with STACKS : De novo map on data 1'. The 'Paired-end fasta files' section has a warning: 'be careful, all files must have a paired-end friend'. The 'Use progeny files' dropdown is set to 'No'. The 'Stack assembly options' dropdown is set to 'Default'. The 'SNP Model Options' dropdown is set to 'Default'. The 'Output type' dropdown is set to 'No compression'. An 'Execute' button is visible at the bottom of the configuration panel.

On the right side, the 'History' panel shows a list of jobs. The top job is '3: snps_output.zip with STACKS : De novo map on data 1', which is completed. Below it, there are two jobs related to 'catalog.snps with STACKS : De novo map on data 1', both of which are also completed. The job details for the top job show the format 'zip, database: 2' and the location of the output files.

Ne pas oublier de rafraîchir l'historique une fois le job terminé pour faire apparaître les jeux de données de sorties additionnels!

2. Les fichiers de sortie

Soit le pipeline lancé sur deux individus, "male" et "female". Dans les fichiers générés, nous retrouvons:

a) Un premier jeu de données nommé result.log

Il permet notamment de vérifier le bon déroulement du job.

```
denovo_map.pl started at 2014-02-14 12:37:49
/local/galaxy/stacks-1.09/bin/denovo_map.pl -p /
```

Nous pouvons voir tout d'abord le script lancé (ici `denovo_map.pl`), ainsi que la date et l'heure de départ du job (ici le 14 février 2014 à midi 37, bonne saint Valentin ;)). Ensuite, vient la ligne de commande lancée, `/local/galaxy/stacks-1.09/bin/denovo_map.pl -p...`

S'en suit le listing des étapes effectuées fichier par fichier :

```
Identifying unique stacks: file 1 of 2 [female]
/opt/dependencies/galaxy_stacks/1.1.0/cmonjeau/stacks_dependencies/2ea494032b10/bin/ustacks -t fasta -f /omaha-
beach/galaxy/58/database/tmp/tmpNo7qXs/tmpqVMSAn/female.fa -o /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpqVMSAn -i 1 -p 4 2>&1
Min depth of coverage to create a stack: 3
Max distance allowed between stacks: 2
Max distance allowed to align secondary reads: 4
Max number of stacks allowed per de novo locus: 3
Deleveraging algorithm: disabled
Removal algorithm: disabled
Model type: SNP
Alpha significance level for model: 0.05
Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpqVMSAn/female.fa
Loaded 29517 RAD-Tags; inserted 4922 elements into the RAD-Tags hash map.
0 reads contained uncalled nucleotides that were modified.
Mean coverage depth is 35; Std Dev: 22.8266 Max: 150
Coverage mean: 35; stdev: 22.8266
Deleveraging trigger: 58; Removal trigger: 81
Calculating distance between stacks...
Distance allowed between stacks: 2
Using a k-mer length of 25
Number of kmers per sequence: 51
Minimum number of k-mers to define a match: 1
Merging stacks, maximum allowed distance: 2 nucleotide(s)
715 stacks merged into 450 stacks; deleveraged 0 stacks; removed 1 stacks.
Mean merged coverage depth is 55.8911; Std Dev: 29.1947; Max: 195
Merging remainder rads
4366 remainder sequences left to merge.
Distance allowed between stacks: 4
Using a k-mer length of 15
Number of kmers per sequence: 61
Minimum number of k-mers to define a match: 1
Matched 4346 remainder reads; unable to match 20 remainder reads.
Number of utilized reads: 29497
Writing results
Refetching sequencing IDs from /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpqVMSAn/female.fa... read 29517 sequence IDs.
```

Nous obtenons ainsi, pour le premier des 2 fichiers, le parent "female" :

-la ligne de commande correspondant au lancement de **ustacks**

-les paramètres renseignés lors de la soumission du job (Min depth of coverage,)

-le nombre total de lectures présentant la signature d'un site de restriction, qualifiées de RAD-tags, chargées (ici 29517) ainsi que le nombre d'éléments insérés dans la carte de hachage de RAD-tags (ici 4922).

-le nombre de lectures avec des nucléotides modifiés

-la profondeur de couverture moyenne

-le nombre de stacks fusionnés car proches (ici 715 fusionnés en 450).

-le nombre de lectures restantes concordant avec les 450 stacks préalablement générés. Ici 4346.

Sur les 29517 lectures d'origines "RAD tagguées", seules 20 n'ont pas été utilisées. Il y a donc 29497 lectures finalement utilisées pour générer 450 stacks.

Puis **ustacks** est exécuté sur le parent "male"

Ensuite est exécuté **cstacks**.

```

/opt/dependencies/galaxy_stacks/1.1.0/cmonjeau/stacks_dependencies/2ea494032b10/bin/cstacks -b 1 -o /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S -s /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/male -p 4 2>41
Number of mismatches allowed between stacks: 0
Loci matched based on sequence identity.
Constructing catalog from 2 samples.
Initializing new catalog...
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/female.tags.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/female.snps.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/female.alleles.tsv
Processing sample 2
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/male.tags.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/male.snps.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/male.alleles.tsv
Searching for sequence matches...
  Number of kmers per sequence: 1
  Minimum number of k-mers to define a match: 1
  449 loci in the catalog, 678 kmers in the catalog hash.
Merging matches into catalog...
  0 loci matched more than one catalog locus and were excluded.
Writing catalog...

```

Là encore, un rappel des paramètres utilisés est fait. Trois types de fichiers différents générés par *ustacks* sont utilisés par individu, à savoir *tags.tsv*, *snps.tsv* et *alleles.tsv*. Nous reviendrons sur ces fichiers dans les sections suivantes. Le catalogue de loci est ensuite créé à partir des échantillons parentaux. Ici, le catalogue créé contient 462 loci (voir *catalog.tags*) dont 449 provenant de l'individu de référence (female) et pouvant être partagé avec le second individu (male).

Enfin, *sstacks* est exécuté:

```

/opt/dependencies/galaxy_stacks/1.1.0/cmonjeau/stacks_dependencies/2ea494032b10/bin/sstacks -b 1 -c /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1 -s /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/female -o /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S -p 4 2>41
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.catalog.tags.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.catalog.snps.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.catalog.alleles.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/female.tags.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/female.snps.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/female.alleles.tsv
Searching for sequence matches...
449 stacks compared against the catalog containing 462 loci.
  449 matching loci, 0 contained no verified haplotypes.
  0 loci matched more than one catalog locus and were excluded.
  0 loci contained SNPs unaccounted for in the catalog and were excluded.
  678 total haplotypes examined from matching loci, 678 verified.
Outputting to file /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/female.matches.tsv
/opt/dependencies/galaxy_stacks/1.1.0/cmonjeau/stacks_dependencies/2ea494032b10/bin/sstacks -b 1 -c /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1 -s /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/male -o /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S -p 4 2>41
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.catalog.tags.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.catalog.snps.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.catalog.alleles.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/male.tags.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/male.snps.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/male.alleles.tsv
Searching for sequence matches...
457 stacks compared against the catalog containing 462 loci.
  457 matching loci, 0 contained no verified haplotypes.
  0 loci matched more than one catalog locus and were excluded.
  0 loci contained SNPs unaccounted for in the catalog and were excluded.
  732 total haplotypes examined from matching loci, 732 verified.
Outputting to file /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/male.matches.tsv
/opt/dependencies/galaxy_stacks/1.1.0/cmonjeau/stacks_dependencies/2ea494032b10/bin/genotypes -b 1 -P /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S -r 1 -c -s 2>41
Found 2 input file(s).
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.catalog.tags.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.catalog.snps.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.catalog.alleles.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/female.matches.tsv
  Parsing /omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/male.matches.tsv
Identified parent IDs: 1 2
Populating observed haplotypes for 2 samples, 462 loci.
Performing automated corrections...
0 potential genotypes in 460 markers, 0 populated; 0 corrected, 0 converted to heterozygotes, 0 unsupported homozygotes removed.
Writing 443 loci to genotype file, '/omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.genotypes_1.tsv'
Writing SQL markers file to '/omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.markers.tsv'
Writing SQL genotypes file to '/omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.genotypes_1.txt'
Writing 462 loci to observed haplotype file, '/omaha-beach/galaxy/58/database/tmp/tmpNo7qXs/tmpAChs9S/batch_1.haplotypes_1.tsv'
denovo_map.pl completed at 2014-08-07 12:03:51

```

Pour chaque individu, *sstacks* détermine l'haplotype à chaque locus de chacun des individus du croisement. Il va alors comparer les stacks de chaque individu (female : 449 loci / male : 457 loci) avec le catalogue complet contenant 462 loci. Ces 462 loci sont ceux du fichier *catalog.tags*.

En parcourant ce fichier, identifiez les loci spécifiques à chacun des individus et ceux partagés

Ces 462 loci représentent les 444 de l'individu "female" partagés avec l'individu "male" + 5 loci spécifiques à l'individu "female" + 13 loci spécifique à l'individu "male"). Ici, 678 haplotypes sont trouvés au total pour l'individu "female", 732 pour "male".

b) Les fichiers tags

XXX.tags.tsv file:

Column	Name	Description
1	Sql ID	This field will always be "0", however the MySQL database will assign an ID when it is loaded.
2	Sample ID	Each sample passed through Stacks gets a unique id for that sample.
3	Stack ID	Each stack formed gets an ID.
4	Chromosome	If aligned to a reference genome using pstacks, otherwise it is blank.
5	Basepair	If aligned to ref genome using pstacks.
6	Strand	If aligned to ref genome using pstacks.
7	Sequence Type	Either 'consensus', 'primary' or 'secondary', see the Stacks paper for definitions of these terms.
8	Sequence ID	The individual sequence read that was merged into this stack.
9	Sequence	The raw sequencing read.
10	Deleveraged Flag	If "1", this stack was processed by the deleveraging algorithm and was broken down from a larger stack.
11	Blacklisted Flag	If "1", this stack was still confounded despite processing by the deleveraging algorithm.
12	Lumberja ckstack Flag	If "1", this stack was set aside due to having an extreme depth of coverage.

Notes: For the tags file, each stack will start in the file with a consensus sequence for the entire stack followed by the flags for that stack. Then, each individual read that was merged into that stack will follow. The next stack will start with another consensus sequence.

Exemple du *catalog.tags.tsv* :

#Seq_ID	Sample_ID	Stack_ID	Chromosome	Basepair	Strand	Sequence_Type	Allele	Sequence_ID	Sequence	Deleveraged_Flag	Blacklisted_Flag	Lumberja_ckstack_Flag
0	1	1		0	+	consensus	0_1_1_2_5		TGCAGGTACATCAATCAATCGGACTACATCTGGAACACCTGATCCAAACAACTATGTGTTTTGTCTGCATGG	0	0	0
0	1	2		0	+	consensus	0_1_2_199		TGCAGGAGCTCTTAGACCTCCCTGCTTAGGCCTAGAAGATCATCGAATAATGCCTTGCATTGTGGCCAAA	0	0	0
0	1	3		0	+	consensus	0_1_3_2_16		TGCAGGTTTCATGGTGTCTGTACCACTGTCTTACTGGATTGTGTGGAAAGACTGTAGSATCTCATGCTGGTG	0	0	0
0	1	4		0	+	consensus	0_1_4_2_176		TGCAGGCTAACATGACAGGAGCAGGCAACACACTGTGTGCAATACAGTAGTTCAGTTTGTATTATTCA	0	0	0
0	1	5		0	+	consensus	0_1_5_2_171		TGCAGGCCAGCATGATGABACACTGAAAAGTGAAGGCGAGAGCCGAGAGGTCTCTGCTCCGATCACTG	0	0	0
0	1	6		0	+	consensus	0_1_6_2_267		TGCAGSCCTCTGAAGCTGCTCACTCAAGTATTGATTAGCACATGCAACAATTAATTAAGCTTCTCTCT	0	0	0
0	1	7		0	+	consensus	0_1_7_2_10		TGCAGGAGGTGTGTCTTCAGACAAGACATTTAGACGGAGGTCTGACGCTCTTTAACTAAAGATCCTGTATT	0	0	0
0	1	8		0	+	consensus	0_1_8_2_371		TGCAGGTTATGAATTAAGGAATAGCAAGGCTATAAAATCCAGTAACCTCTTCAATCATATAATAATG	0	0	0
0	1	9		0	+	consensus	0_1_9_2_334		TGCAGGCGAGGAGACCCCTGAGACCATGTTTCCCTGTGABACAGAACTACTCAGGTGGTCTGTTGGCCATA	0	0	0
0	1	10		0	+	consensus	0_1_10_2_14		TGCAGGGAACAATCCAGACAACCTGGAGGGATCAGAGGTTCTGCTGTAGAAGTGCAGACAGGCCAAAGA	0	0	0
0	1	11		0	+	consensus	0_1_11_2_58		TGCAGGCTCTGGGGCTCAAAATCTTTTAAATTTGACGCAAGGTTAGCTCTGCTCTTTAGGAGAAAA	0	0	0
0	1	12		0	+	consensus	0_1_12_2_36		TGCAGGCCAGACTGACAGGACATGTGGGCTTCTGCAACCCCGGTTTTGTACAGACTATCAAAACGCTCAG	0	0	0
0	1	13		0	+	consensus	0_1_13_2_15		TGCAGGGTGTGATATTCCTGCTCTCTGTAAATGABAGATATCAGAGTGTCACTGAGCAGTGGAGTTAG	0	0	0
0	1	14		0	+	consensus	0_1_14_2_13		TGCAGGTCATACATAGTCAATAAATTTATGTTTTAACAACAAATCATGTTTTCCACTAGAAATCAAAATGAA	0	0	0
0	1	15		0	+	consensus	0_1_15_2_17		TGCAGGTAACAGGAGATTATGAGTTTCCCTTTAGTGTGCAACAGCAATGCACCACATGTTTCGGTCTAGATG	0	0	0
0	1	16		0	+	consensus	0_1_16_2_265		TGCAGGACTATTTAAATCTGCAGCCCTTCTACCATGTTTGCATTTTATCAATATCAGAAGGCCACACAGGT	0	0	0
0	1	17		0	+	consensus	0_1_17_2_387		TGCAGGTGATTATACCTTCGCTACTGCTCAGACTGCACACTGTATCAGACACTGCAGCATGCACACTGTCA	0	0	0

Nous pouvons constater que l'identifiant "*Sequence_ID*" est constitué à partir des ID des individus et des piles correspondantes comme suit : "*individu1_stackIDindividu1, individu2_stackIDindividu2*". Il est ainsi directement aisé de savoir quels SNPs sont détectés sur des loci séquencés chez nos 2 individus.

Pour un *sample.tags.tsv* :

*SQ_ID	Sample_ID	Stack_ID	Chromosome	Basepair	Strand	Sequence_Type	Allele	Sequence_ID	Sequence	Deleterevl_Flag	Blacklisted_Flag	Lumbeys_ckptack_Flag
0	1	1				model			TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	GGAGG_2_007_1653_1798_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	GGAGG_2_005_118_125_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	GGAGG_2_009_413_117_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	GGAGG_2_009_429_162_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	GGAGG_2_008_418_269_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CTAGG_2_007_1687_1559_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CTAGG_2_001_1041_1893_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CTAGG_2_008_1118_819_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CTAGG_2_005_195_1323_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CTAGG_2_002_192_332_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CTAGG_2_007_1669_385_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CTAGG_2_010_319_109_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_002_1919_640_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_009_1752_198_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0019_1167_1548_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0018_148_1281_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0018_1063_795_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0021_1543_547_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0028_130_166_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0038_1032_772_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0038_1118_273_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0041_1195_1375_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0048_172_639_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0048_225_405_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0051_874_1617_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0051_1597_494_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0056_1677_815_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0057_178_1779_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0058_13_851_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0062_800_610_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0061_3843_1807_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0069_1473_640_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0072_1026_1898_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0078_870_389_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	0	CAGTC_2_0087_1633_1356_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	1	GGAGG_2_0012_1275_1694_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	1	GGAGG_2_0019_481_1131_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	1	GGAGG_2_0039_606_1603_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	1	GGAGG_2_0062_618_273_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			
0	1	1				primary	1	GGAGG_2_0090_748_1491_179984	TCGAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG			

Primary vs secondary reads :

0	1	1	primary	1	CAGTC_2_0076_1442_1576_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	primary	1	CAGTC_2_0079_1540_762_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	primary	1	CAGTC_2_0080_1139_865_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	primary	1	CAGTC_2_0081_458_1052_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	primary	1	CAGTC_2_0088_1542_1441_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	primary	1	CAGTC_2_0092_629_1368_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	primary	1	CAGTC_2_0093_1732_263_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	primary	1	CAGTC_2_0097_697_1869_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	secondary	1	CAGTC_2_0082_1602_271_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	secondary	1	CAGTC_2_0093_92_1898_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	secondary	1	CAGTC_2_0058_1624_135_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	secondary	1	CAGTC_2_0091_1720_1032_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG
0	1	1	secondary	1	CAGTC_2_0038_1473_926_179984	TGCAAGGTACATCAATCAATCGGACTACACTCTGAACCACCTGATCCACACAAACTATGTGTTTTGTCTGCACGG

Utiliser les fichiers *individu.tags.tsv* pour retrouver le nombre de tags par individu. Vous serez probablement amenés à utiliser l'outil "count"

c) Les fichiers alleles

Column	Name	Description
1	Sql ID	This field will always be "0", however the MySQL database will assign an ID when it is loaded.
2	Sample ID	
3	Stack ID	
4	Haplotype	The haplotype, as constructed from the called SNPs at each locus.
5	Percent	Percentage of reads that have this haplotype
6	Count	Raw number of reads that have this haplotype

Sample1.alleles.tsv à gauche et *catalog.alleles.tsv* à droite

#SQL_ID	Sample_ID	Stack_ID	Haplotype	Percent	Count	#SQL_ID	Sample_ID	Stack_ID	Haplotype	Percent	Count
0	1	1	C	48.0519	37	0	1	1	C	0	0
0	1	1	T	51.9481	40	0	1	1	T	0	0
0	1	7	C	45.283	24	0	1	2	C	0	0
0	1	7	G	54.717	29	0	1	2	G	0	0
0	1	9	A	40.7407	22	0	1	3	A	0	0
0	1	9	G	59.2593	32	0	1	3	G	0	0
0	1	13	C	44.1558	34	0	1	4	A	0	0
0	1	13	T	55.8442	43	0	1	4	G	0	0
0	1	14	G	56.25	36	0	1	5	A	0	0
0	1	14	T	43.75	28	0	1	5	G	0	0
0	1	15	C	48.3871	75	0	1	6	C	0	0
0	1	15	T	51.6129	80	0	1	6	T	0	0
0	1	18	A	46.7742	29	0	1	7	C	0	0
0	1	18	T	53.2258	33	0	1	7	G	0	0
0	1	21	C	59.7403	46	0	1	8	GC	0	0
0	1	21	T	40.2597	31	0	1	8	TA	0	0
0	1	22	G	62.5	35	0	1	9	A	0	0
0	1	22	T	37.5	21	0	1	9	G	0	0
0	1	25	A	47.6923	31	0	1	10	C	0	0
0	1	25	C	52.3077	34	0	1	10	G	0	0
0	1	26	C	66.6667	16	0	1	11	C	0	0

Sample2.alleles.tsv

#SQL_ID	Sample_ID	Stack_ID	Haplotype	Percent	Count
0	2	6	AT	42.3077	11
0	2	6	GC	57.6923	15
0	2	7	C	56.6265	47
0	2	7	T	42.1687	35
0	2	11	C	57.1429	48
0	2	11	T	41.6667	35
0	2	14	C	39.726	29
0	2	14	G	60.274	44
0	2	16	A	72.8571	51
0	2	16	G	27.1429	19
0	2	20	A	59.2593	16
0	2	20	C	40.7407	11
0	2	23	G	47.6923	31
0	2	23	T	52.3077	34
0	2	26	CG	60.5634	43
0	2	26	TA	39.4366	28
0	2	30	C	60	45
0	2	30	T	40	30
0	2	32	A	43.9024	36
0	2	32	G	56.0976	46
0	2	33	A	58.7156	64

Nous pouvons alors constater que si le Stack_ID 5 présente un SNP dans le fichier catalog.alleles.tsv, il n'est retrouvé dans aucun des fichiers samples (1 ou 2)....

C'est parce que le Stack_ID du sample 1 ne correspond pas au Stack_ID du sample 2 qui lui-même ne correspond pas au Stack_ID du catalog...

Pour s'en rendre compte, il suffit de consulter les fichiers *sample**.matches.tsv*

d) Les fichiers matches

Column	Name	Description
1	Sql ID	This field will always be "0", however the MySQL database will assign an ID when it is loaded.
2	Batch ID	
3	Catalog ID	
4	Sample ID	
5	Stack ID	
6	Haplotype	
7	Stack Depth	

Sample1.tsv à gauche et sample2.alleles.tsv à droite

#SQL_ID	Batch_ID	Catalog_ID	Sample_ID	Stack_ID	Haplotype	Stack_depth	#SQL_ID	Batch_ID	Catalog_ID	Sample_ID	Stack_ID	Haplotype	Stack_depth
0	1	1	1	1	C	37	0	1	18	2	1	T	95
0	1	1	1	1	T	40	0	1	22	2	2	T	58
0	1	2	1	2	G	47	0	1	315	2	3	T	81
0	1	3	1	3	G	33	0	1	391	2	4	T	75
0	1	4	1	4	A	53	0	1	1	2	5	T	92
0	1	5	1	5	G	21	0	1	227	2	6	AT	11
0	1	6	1	6	C	73	0	1	227	2	6	GC	15
0	1	7	1	7	C	24	0	1	388	2	7	AC	47
0	1	7	1	7	G	29	0	1	388	2	7	AT	35
0	1	8	1	8	TA	50	0	1	259	2	8	T	48
0	1	9	1	9	A	22	0	1	274	2	9	G	42
0	1	9	1	9	G	32	0	1	7	2	10	G	63
0	1	10	1	10	C	67	0	1	139	2	11	C	48
0	1	11	1	11	C	39	0	1	139	2	11	T	35
0	1	12	1	12	G	49	0	1	44	2	12	AA	32
0	1	13	1	13	C	34	0	1	14	2	13	G	70
0	1	13	1	13	T	43	0	1	10	2	14	C	29
0	1	14	1	14	G	36	0	1	10	2	14	G	44
0	1	14	1	14	T	28	0	1	13	2	15	T	45
0	1	15	1	15	C	75	0	1	3	2	16	A	51
0	1	15	1	15	T	80	0	1	3	2	16	G	19

Le Catalog_ID (= Stack_ID du catalog), reprend le Stack_ID de l'individu de "référence", ici sample 1, mais la numérotation est bien différente de celle du Stack_ID du sample 2.... Ainsi, dans le fichier "catalog.alleles.tsv", le Stack_ID 3 correspond au Stack_ID 16 du sample 2!

e) Les fichiers snps

Column	Name	Description
1	Sql ID	This field will always be "0", however the MySQL database will assign an ID when it is loaded.
2	Sample ID	
3	Stack ID	
4	SNP Column	
5	Likelihood ratio	From the SNP-calling model.
6	Rank_1	Majority nucleotide.
7	Rank_2	Alternative nucleotide.

Notes: If a stack has two SNPs called within it, then there will be two lines in this file listing each one.

sample.snps.tsv à gauche et catalog.snps.tsv à droite

#SQL_ID	Sample_ID	Stack_ID	SNP_Column	Likelihood_ratio	Maj_nt	Alt_nt	#SQL_ID	Sample_ID	Stack_ID	SNP_Column	Likelihood_ratio	Maj_nt	Alt_nt
0	1	1	72	-81.1804	T	C	0	1	1	72	0	T	C
0	1	7	21	-52.261	G	C	0	1	2	15	0	G	C
0	1	9	16	-46.4764	G	A	0	1	3	16	0	A	G
0	1	13	70	-73.6513	T	C	0	1	4	34	0	G	A
0	1	14	58	-60.5197	G	T	0	1	5	50	0	G	A
0	1	15	20	-164.631	T	C	0	1	6	64	0	C	T
0	1	18	28	-63.4613	T	A	0	1	7	21	0	G	C
0	1	21	72	-65.1731	C	T	0	1	8	15	0	G	T
0	1	22	53	-42.6043	G	T	0	1	8	66	0	C	A
0	1	25	69	-67.9754	C	A	0	1	9	16	0	G	A
0	1	26	36	-14.8594	C	T	0	1	10	36	0	G	C
0	1	30	46	-76.5693	C	T	0	1	11	23	0	G	C
0	1	31	38	-77.3799	A	T	0	1	12	43	0	G	A
0	1	33	48	-76.7795	G	T	0	1	13	70	0	T	C
0	1	36	73	-59.2943	G	A	0	1	14	58	0	G	T

Nous avons détecter des SNPs chez nos 2 individus et nous pouvons déterminer lesquels sont situés sur les même loci.

3. Retour à l'analyse

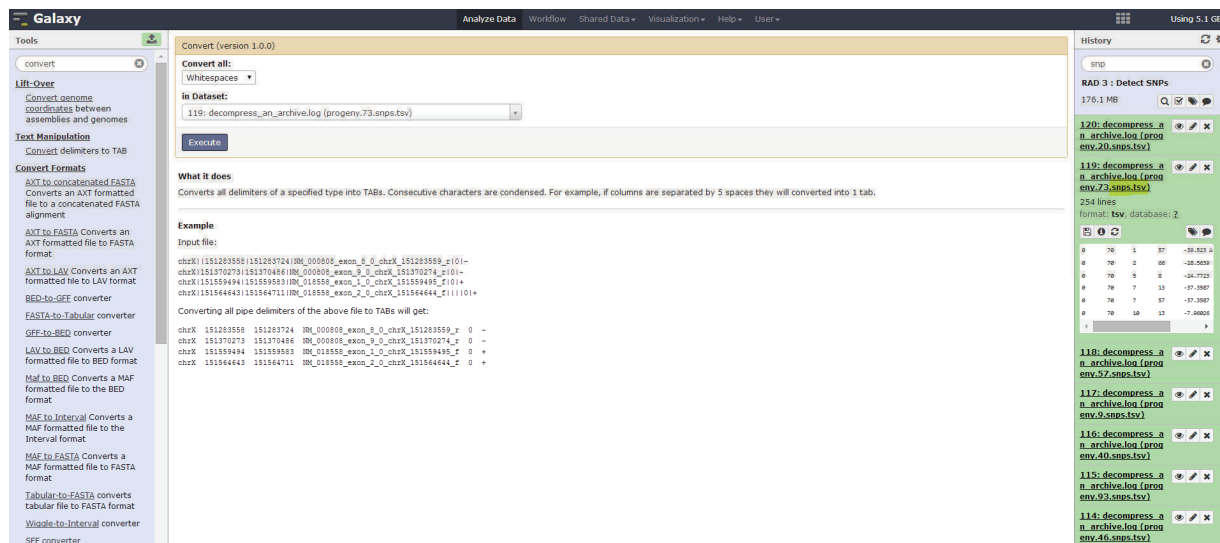
Extraire l'archive *Parents.zip* (outil **Decompress an archive**) et ne pas oublier de rafraîchir l'historique une fois le job terminé. Modifier, si besoin, le type des fichiers extraits en fastq.

Relancer **STACKS denovo** en spécifiant l'usage "*Population*". Au passage, vous pouvez voir que les fichiers nouvellement décompressés sont directement utilisables en entrée de l'outil **STACKS denovo**.

Que pouvez-vous dire des fichiers snps générés par les 2 usages? Vous serez probablement amenés à utiliser l'outil "compare two datasets"

Sélectionnez les snps identifiés chez l'individu femelle présentant un score (likelihood ratio) meilleur que le ¾ des scores. Vous serez probablement amenés à utiliser les outils "summary statistics" et "filter"

Si besoin (normalement pas dans le cas de figure présent), convertir les espaces en tabulations pour un des fichiers *.snps.tsv* via l'outil **Convert**



Il est aussi possible d'utiliser comme solution alternative, la modification directe du datatype à l'aide du petit crayon du dataset concerné.

The screenshot shows the Galaxy web interface. The main panel is titled 'Change data type'. It has a 'New Type' dropdown menu currently set to 'tabular'. Below the dropdown, there is a note: 'This will change the datatype of the existing dataset but not modify its contents. Use this if Galaxy has incorrectly guessed the type of your dataset.' There is a 'Save' button below the note. The left sidebar contains a search bar with 'stacks' and a list of tools under 'STACKS toolsuite'. The right sidebar shows a 'History' panel with a search bar and a list of jobs, including 'RAD 3 : Detect SNPs' and several 'decompress_a_n_archive.log' jobs.

Filtrer le fichier tabulé obtenu par la valeur du *Likelihood ratio*. Plus ce score est important, plus il est significatif :
 A l'aide de l'outil **Summary Statistics**, récupérer les valeurs de likelihood ratio (5^{ème} colonne du fichier *.snps.tsv* tabulé) du premier et troisième quartile (25% et 75%).

The screenshot shows the Galaxy web interface for the 'Summary Statistics (version 1.1.0)' tool. The 'Summary statistics on:' dropdown is set to '121: Convert on data 119'. The 'Column or expression:' field is set to 'c5'. There is an 'Execute' button. Below the configuration area, there are several informational blocks: a warning about input datasets, two tips, and a section for R functions. The 'Syntax' section explains how to reference columns and provides examples. The 'Examples' section shows an input dataset and its corresponding summary statistics.

Input Dataset:

c1	c2	c3	c4	c5	c6
586	chrX	161416	170807	41108_at	16990
73	chrX	505078	532318	35073_at	1700
595	chrX	1361578	1388460	33665_at	1960
74	chrX	1420620	1461919	1185_at	8600

Summary Statistics on column c6 of the above input dataset:

#sum	mean	stdev	0%	25%	50%	75%	100%
29250.000	7312.500	7198.636	1700.000	1895.000	5280.000	10697.500	16990.000

Ici, les valeurs sont toujours négatives, et celle du premier quartile, pour l'individu femelle, est de -32.8659 (-81.1918 avec la version 1.18 de Stacks ;)

Galaxy Analyze Data Workflow

Data Viewer: Summary Statistics on data 121

#sum	mean	stdev	0%	25%	50%	75%	100%
-6837.51	-26.9193	10.0833	-60.7429	-32.8659	-26.3267	-20.5751	-4.60291

Syntax
This tool computes basic summary statistics on a given column, or on a valid expression containing

Nous filtrerons le fichier en récupérant uniquement les SNP pour lesquels le Likelihood ratio est inférieur à cette valeur (ici -32.87).

Galaxy Analyze Data Workflow Shared Data Visualization Help User

Tools

filter

Text Manipulation
Filter on ambiguities in polymorphism datasets

Filter and Sort
Filter data on any column using simple expressions
GFF
Filter GFF data by attribute using simple expressions
Filter GFF data by feature count using simple expressions
Filter GTF data by attribute values list

Fetch Alignments
Filter MAF blocks by Species
Filter MAF blocks by Size
Filter MAF by specified attributes

FASTA manipulation
Filter sequences by length

NGS: QC and manipulation
Filter by quality
Filter FASTQ reads by quality score and length

Workflows
All workflows

Filter (version 1.1.0)

Filter:
121: Convert on data 119
Dataset missing? See TIP below.

With following condition:
c5<=-32.87

Number of header lines to skip:
0

Execute

⚠ Double equal signs, "=", must be used as "equal to" (e.g., c1 == "chr22")

📘 TIP: Attempting to apply a filtering condition may throw exceptions if the data type (e.g., string, integer) in every line of the columns being filtered is not appropriate for the condition (e.g., attempting certain numerical calculations on strings). If an exception is thrown when applying the condition to a line, that line is skipped as invalid for the filter condition. The number of invalid skipped lines is documented in the resulting history

📘 TIP: If your data is not TAB delimited, use Text Manipulation->Convert

Syntax
The filter tool allows you to restrict the dataset using simple conditional statements.
Columns are referenced with **c** and a **number**. For example, **c1** refers to the first column of a tab-delimited file
Make sure that multi-character operators contain no white space (e.g., <= is valid while < = is not valid)
When using 'equal-to' operator **double equal sign** "=" **must be used** (e.g., c1=="chr1")
Non-numerical values must be included in single or double quotes (e.g., c6=="+")
Filtering condition can include logical operators, but **make sure operators are all lower case** (e.g., (c11="chrX" and c11="chrY") or not c6=="+")

Example
c1=="chr1" selects lines in which the first column is chr1
c3-c2<100*c4 selects lines where subtracting column 3 from column 2 is less than the value of column 4 times 100
len(c2-split("\t"))<4 will select lines where the second column has less than four comma separated elements
c2>=1 selects lines in which the value of column 2 is greater than or equal to 1
Numbers should not contain commas - c2<=44,554,350 will not work, but c2<=44554350 will
Some words in the data can be used, but must be single or double quoted (e.g., c3=="exon")

Nous avons bien gardé ~25% des valeurs

Using 5.1 GB

History

search datasets

RAD 3 : Detect SNPs
176.1 MB

126: Filter on data 121

64 lines
format: **tabular**, database: ?

Filtering with $c5 \leq -32.87$, kept 25.20% of 254 valid lines (254 total lines).

1	2	3	4	5	6	7
0	70	1	57	-39.523	A	C
0	70	7	13	-37.3587	T	G
0	70	7	57	-37.3587	A	G
0	70	12	55	-36.931	C	T
0	70	34	57	-34.8981	T	G
0	70	55	40	-35.4551	G	A

125: Summary Statistics on data 121

1 line, 1 comments
format: **tabular**, database: ?

1	2	3	4	5
#sum	mean	stdev	0%	25%
-6837.51	-26.9193	10.0033	-60.7429	-32.86

II. Analyse RAD-seq sous Galaxy : La cartographie génétique

Cet exercice utilise des données générées par les développeurs de STACKS. Ils ont développé une carte génétique pour le Lépisosté tacheté (*Lepisosteus oculatus*) (poisson crocodile parfois surnommé brochet crocodile) et présentent ici les données d'un seul groupe de liaison. La carte génétique provient d'un croisement de type pseudo-test cross F1 entre 2 parents et 93 de leur descendants F1. Sont conservés pour l'exercice, uniquement des marqueurs apparaissant dans un unique groupe de liaison, et les lectures brutes des "stacks" qui ont contribué à ce groupe de liaison. Les fichiers sont déjà nettoyés, il ne sera donc pas nécessaire d'utiliser l'outil "[STACKS : Process radtags](#) Run the STACKS cleaning script" ici.

Les 95 fichiers fasta de départ sont donc déjà démultiplexés et se présentent comme suit :

```
>TTAAT_1_0046_17989_1193_1[67302]
TGCAGGCGAGGAAGTCACAGAGATCCCTGGCCAACACTACTGTAGTTCGAACAGGAACCGAGCTGACAGGGCGCAGA
>TTAAT_1_0092_18487_11460_1[67302]
TGCAGGCGAGGAAGTCACAGAGATCCCTGGCCAACACTACTGTAGTTCGAACAGGAACCGAGCTGACAGGGCGCAGA
>TTAAT_1_0094_8623_4235_1[67302]
TGCAGGCGAGGAAGTCACAGAGATCCCTGGCCAACACTACTGTAGTTCGAACAGGAACCGAGCTGACAGGGCGCAGA
>TTAAT_1_0102_18666_9095_1[67302]
TGCAGGCGAGGAAGTCACAGAGATCCCTGGCCAACACTACTGTAGTTCGAACAGGAACCGAGCTGACAGGGCGCAGA
>TTAAT_1_0114_6838_19507_1[67302]
TGCAGGCGAGGAAGTCACAGAGATCCCTGGCCAACACTACTGTAGTTCGAACAGGAACCGAGCTGACAGGGCGCAGA
>TTAAT_1_0117_2046_2348_1[67302]
TGCAGGCGAGGAAGTCACAGAGATCCCTGGCCAACACTACTGTAGTTCGAACAGGAACCGAGCTGACAGGGCGCAGA
>TTAAT_1_0038_16287_13120_1[22594]
TGCAGGCCTTGTGAAACTGAACAACACAAAAGGTTCTATCAATTAACCGCAGATAATTAGTTGTGTTTCTCCA
>TTAAT_1_0052_11753_10140_1[22594]
TGCAGGCCTTGTGAAACTGAACAACACAAAAGGTTCTATCAATTAACCGCAGATAATTAGTTGTGTTTCTCCA
>TTAAT_1_0054_1743_17715_1[22594]
TGCAGGCCTTGTGAAACTGAACAACACAAAAGGTTCTATCAATTAACCGCAGATAATTAGTTGTGTTTCTCCA
>TTAAT_1_0074_2389_17780_1[22594]
TGCAGGCCTTGTGAAACTGAACAACACAAAAGGTTCTATCAATTAACCGCAGATAATTAGTTGTGTTTCTCCA
>TTAAT_1_0087_18378_6512_1[22594]
```

Pour chaque séquence, on retrouve le barcode (ici TTAAT), et d'autres informations.

Créer un nouvel historique et renommer le (ex : Tuto GenOuest 1 RAD seq : Genetic map).

Comme les données sont les mêmes qu'utilisées dans la partie précédente, vous pouvez les rapatrier en utilisant l'option "*copy datasets*" de l'historique (URL d'origine : http://creskolab.uoregon.edu/stacks/tut_gar.php, les données : http://creskolab.uoregon.edu/stacks/tutorial/stacks_samples.tar.gz). Vous récupérer ainsi 2 archives :

-une regroupant les 2 jeux de données de séquences des parents : Parents.zip

-une regroupant les 93 jeux de données de séquences des descendants F1 : Progeny.zip

Nous sommes maintenant prêts à exécuter le pipeline STACKS denovo.

Sélectionner l'outil "[STACKS : De novo map](#) Run the STACKS denovo_map.pl wrapper"

Vous pouvez préciser certains paramètres. Ainsi, en sélectionnant le mode Advanced pour le paramètre *Stack assembly options*, vous pouvez préciser:

-la profondeur minimum de la pile (-m), ici "Minimum number of identical". Ce paramètre, passé à **ustacks**, contrôle le nombre de lectures correspondant exactement devant être trouvé pour créer une pile chez un individu. Nous pouvons sélectionner 3, une pile ne sera générée chez un individu que si au moins 3 lectures correspondent exactement.

-la profondeur minimum de la pile (-P), ici "Minimum number of identical (progeny)". Ce paramètre, passé à **ustacks**, contrôle le nombre de lectures correspondant exactement devant être trouvé pour créer une pile chez un descendant. Nous pouvons sélectionner 3, une pile ne sera générée chez un descendant que si au moins 3 lectures correspondent exactement.

-la distance maximum permise entre des piles pour qu'elles soient fusionnées en un locus potentiel chez un individu. Il s'agit du paramètre -M, ici "Number of mismatches allowed between loci when processing a single individual". Nous précisons ici une distance maximale de 2 nucléotides.

-le nombre maximal de différence entre tags. Ce paramètre permet notamment de pouvoir créer un locus de type homozygote dans le catalogue alors qu'il est en réalité hétérozygote. Ceci est pratique pour conserver les loci hétérozygotes entre parents mais homozygote chez chacun d'entre eux. Il s'agit du paramètre -n, "specify the number of mismatches allowed between loci when building the catalog" ici. Nous pouvons fixer ce paramètre à 3.

-supprimer ou casser les RAD-tags très répétitifs. Il s'agit de supprimer les piles "bûcheronnes" de l'analyse et briser les piles modérément surdimensionnées. Nous cocherons ici cette option, correspondant au paramètre -t.

On choisira enfin le type de sortie : "compressed by categories"

The screenshot displays the Galaxy web interface for the 'STACKS : De novo map (version 1.0.0)' tool. The configuration panel includes the following settings:

- Select your usage:** Genetic map
- Files containing parent sequences:** 2: Progeny.zip
- Use progeny files:** Yes
- Files containing progeny sequences:** 1: Parents.zip
- Stack assembly options:** Advanced
- Minimum number of identical:** 3
- Minimum number of identical (progeny):** -1
- Number of mismatches allowed between loci when processing a single individual:** 2
- Number of mismatches allowed when aligning secondary reads:** -1
- specify the number of mismatches allowed between loci when building the catalog:** 3
- remove, or break up, highly repetitive RAD-Tags in the ustacks program:**
- disable calling haplotypes from secondary reads:**
- SNP Model Options:** Default
- Output type:** No compression

The 'History' panel on the right shows the workflow history with two files: '2: Progeny.zip' and '1: Parents.zip'.

Dans le cas présent, 12 jeux de données sont générés dans l'historique :

The screenshot displays the Galaxy web interface during a workflow execution. The main window shows the execution log for the tool 'denovo_map.pl'. The log starts with the command: `/local/galaxy/stacks-1.09/bin/denovo_map.pl -p /opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpf2ssdu/male.fa -p /opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpf2ssdu/female.fa`. The log then details the process of identifying unique stacks, merging them, and generating various output files like `result.log`, `catalog.tags.tsv`, `catalog.snps.tsv`, and `matches_output.zip`. The 'Tools' sidebar on the left lists various tools related to NGS and STACKS. The 'History' sidebar on the right shows a list of previous workflow runs.

1. Un premier jeu de données nommé `result.log`, permettant de vérifier le bon déroulement du job.

```
denovo_map.pl started at 2014-02-14 12:37:49
/local/galaxy/stacks-1.09/bin/denovo_map.pl -p /
```

Comme vu lors de la première partie, nous pouvons voir le script lancé (ici `denovo_map.pl`), ainsi que la date et l'heure de départ du job. Ensuite, vient la ligne de commande lancée, `/local/galaxy/stacks-1.09/bin/denovo_map.pl -p..`
S'en suit le listing des étapes effectuées fichier par fichier :

```

Identifying unique stacks; file 1 of 95 [male]
/local/galaxy/stacks-1.09/bin//ustacks -t fasta -f /opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpf2ssdu/male.fa -o /opt/ga
Min depth of coverage to create a stack: 3
Max distance allowed between stacks: 2
Max distance allowed to align secondary reads: 4
Max number of stacks allowed per de novo locus: 3
Deleveraging algorithm: enabled
Removal algorithm: enabled
Model type: SNP
Alpha significance level for model: 0.05
Parsing /opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpf2ssdu/male.fa
Loaded 33021 RAD-Tags; inserted 5129 elements into the RAD-Tags hash map.
0 reads contained uncalled nucleotides that were modified.
Mean coverage depth is 37; Std Dev: 26.6047 Max: 221
Coverage mean: 37; stdev: 26.6047
Deleveraging trigger: 64; Removal trigger: 90
Calculating distance for removing repetitive stacks.
Distance allowed between stacks: 1
Using a k-mer length of 37
Number of kmers per sequence: 39
Miniumum number of k-mers to define a match: 2
Removing repetitive stacks.
Removed 45 stacks.
758 stacks remain for merging.
Calculating distance between stacks...
Distance allowed between stacks: 2
Using a k-mer length of 25
Number of kmers per sequence: 51
Miniumum number of k-mers to define a match: 1
Merging stacks, maximum allowed distance: 2 nucleotide(s)
758 stacks merged into 460 stacks; deleveraged 2 stacks; removed 0 stacks.
Mean merged coverage depth is 62.0022; Std Dev: 32.6096; Max: 221
Merging remainder radtags
4500 remainder sequences left to merge.
Distance allowed between stacks: 4
Using a k-mer length of 15
Number of kmers per sequence: 61
Miniumum number of k-mers to define a match: 1
Matched 3973 remainder reads; unable to match 527 remainder reads.
Number of utilized reads: 32494
Writing results
Refetching sequencing IDs from /opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpf2ssdu/male.fa... read 33021 sequence IDs.

```

Nous obtenons ainsi pour le premier des 95 fichiers, le parent male :

- la ligne de commande correspondant au lancement de **ustacks**
- les paramètres renseignés lors de la soumission du job (Min depth of coverage,)
- le nombre total de lectures, qualifiées de RAD-tags, chargées (ici 33021) ainsi que le nombre d'éléments insérés dans la carte de hachage de RAD-tags (ici 5129).
- le nombre de lectures avec des nucléotides modifiés
- la profondeur de couverture moyenne
- le nombre de stacks répétés supprimés (ici 45)
- le nombre de stacks fusionnés car proches (ici 758 fusionnés en 460).
- le nombre de lectures restantes concordant avec les 460 stacks préalablement générés. Ici 3973.

Sur les 33021 lectures d'origines "RAD tagguées", seules 527 n'ont pas été utilisées. Il y a donc 32494 lectures finalement utilisées pour générer 460 stacks.

Ensuite est exécuté **cstacks**.

```

/local/galaxy/stacks-1.09/bin//cstacks -b 1 -o /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9 -s
Number of mismatches allowed between stacks: 3
Loci matched based on sequence identity.
Constructing catalog from 2 samples.
Initializing new catalog...
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/male.tags.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/male.snps.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/male.alleles.tsv
Processing sample 2
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/female.tags.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/female.snps.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/female.alleles.tsv
Searching for sequence matches...
Distance allowed between stacks: 3
Using a k-mer length of 17
Number of kmers per sequence: 59
Minimum number of k-mers to define a match: 8
426 loci in the catalog, 29323 kmers in the catalog hash.
Merging matches into catalog...
  0 loci matched more than one catalog locus and were excluded.
Writing catalog...

```

Là encore, un rappel des paramètres utilisés est fait. Le catalogue de loci est ensuite créé à partir des échantillons parentaux. Ici, le catalogue créé contient 459 loci (voir *catalog.tags*) dont 426 provenant de l'individu de référence (male) et pouvant être partagé avec le second individu (female).

Enfin, **sstacks** est exécuté:

```

/local/galaxy/stacks-1.09/bin//sstacks -b 1 -c /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/batch_1 -s /
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/batch_1.catalog.tags.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/batch_1.catalog.snps.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/batch_1.catalog.alleles.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/male.tags.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/male.snps.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/male.alleles.tsv
Searching for sequence matches...
426 stacks compared against the catalog containing 459 loci.
426 matching loci, 0 contained no verified haplotypes.
0 loci matched more than one catalog locus and were excluded.
0 loci contained SNPs unaccounted for in the catalog and were excluded.
695 total haplotypes examined from matching loci, 695 verified.
Outputting to file /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/male.matches.tsv

```

Pour chaque individu, **sstacks** détermine l'haplotype à chaque locus de chacun des individus du croisement. Il va alors comparer les stacks de chaque individu avec le catalogue de 459 loci. Ici, 695 haplotypes sont trouvés au total, tous vérifiés.

Enfin, **genotypes** est exécuté. Il récupère les loci contenant les marqueurs identifiés chez les parents, puis y associe les haplotypes des descendants. Si le premier parent présente les haplotypes GA (ex : *aatggtgtGgtccctcgtAc*) et AC (ex : *aatggtgtAgtcctcgtCc*), et le second parent l'haplotype GA (ex : *aatggtgtGgtccctcgtAc*), Stacks déclare un marqueur *ab/aa* pour ce locus. Le programme **genotypes** associe alors GA à *a*, et AC à *b* chez les parents puis scanne les descendants afin d'identifier quels haplotypes sont présents pour chacun d'entre eux et enregistrer les génotypes associés (soit *ab* ou *aa* dans le cas présent).

```

/local/galaxy/stacks-1.09/bin//genotypes -b 1 -P /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9 -t gen -r 1 -c -s 2>#1
Found 95 input file(s).
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/batch_1.catalog.tags.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/batch_1.catalog.snps.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/batch_1.catalog.alleles.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/female.matches.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/male.matches.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/progeny_1.matches.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/progeny_10.matches.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/progeny_11.matches.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/progeny_12.matches.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/progeny_13.matches.tsv
  Parsing /opt/galaxy-dist/database/tmp/tmpIqL3Tn/tmpubNNI9/progeny_14.matches.tsv

```

```

Parsing /opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpubNNI9/progeny_93.tags.tsv
Parsing /opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpubNNI9/progeny_93.snps.tsv
Parsing /opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpubNNI9/progeny_93.alleles.tsv
Parsing /opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpubNNI9/progeny_94.tags.tsv
Parsing /opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpubNNI9/progeny_94.snps.tsv
Parsing /opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpubNNI9/progeny_94.alleles.tsv
42036 potential genotypes in 452 markers, 38874 populated; 1814 corrected, 1216 converted to heterozygotes, 598 unsupported homozygotes removed.
Writing 452 loci to genotype file, '/opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpubNNI9/batch_1.genotypes_1.tsv'
Writing SQL markers file to '/opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpubNNI9/batch_1.markers.tsv'
Writing SQL genotypes file to '/opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpubNNI9/batch_1.genotypes_1.txt'
Writing 459 loci to observed haplotype file, '/opt/galaxy-dist/database/tmp/tmpIq13Tn/tmpubNNI9/batch_1.haplotypes_1.tsv'
denovo_map.pl completed at 2014-02-14 12:38:16

```

Au final, 452 loci, marqueurs, sont conservés pour créer le fichier génotype **batch_1.genotypes_1.tsv**. 459 loci sont répertoriés dans le fichier haplotype observé **batch_1.haplotypes_1.tsv**. denovo_map.pl s'est exécuté en moins d'une minute. Les 42036 génotypes potentiels sont enregistrés dans le fichier **batch_1.genotypes_1.txt**. 452 marqueurs sont enregistrés dans le fichier **batch_1.markers.tsv**.....

Réexécuter **denovo_map** en spécifiant l'option "**compressed all outputs**". Nous obtenons cette fois deux jeux de données en sortie, un fichier de log identique à celui généré lors de la précédente exécution de **denovo_map** et une archive nommée total_output.zip contenant tous les jeux de données de sortie de **denovo_map**. C'est cette archive qui sera utilisée par la suite comme fichier d'entrée de l'outil "**STACKS : genotypes Run the STACKS genotypes program**".

The screenshot shows the Galaxy web interface. The central panel displays the execution of the `denovo_map.pl` tool. The output includes the command used, file paths, and various parameters and results. The right panel shows a history of jobs, with the most recent job being "total_output.zip with STACKS : De novo map on data 2 and data 1". The left panel shows the "Tools" menu with "genotype" selected.

Exécuter l'outil "**STACKS : genotypes Run the STACKS genotypes program**" sur le jeu de données total_output.zip sans spécifier d'option. Cette nouvelle étape permet :

- 1/de spécifier un type particuliers de carte (ex: F1, F2, Doubles Haploid, Back Cross, Cross Pollination)
- 2/d'exporter les informations génotypiques dans différents formats spécifiques de programmes tiers de cartographie (ex: JoinMap, R/qtI). Il faut savoir que le format R/QTl pour un protocole F2 peut servir en entrée d'outils comme MapMaker ou Carthagène
- 3/de spécifier un nombre minimum de descendants et/ou une couverture minimum dans la pile pour considérer un locus

4/d'effectuer des corrections automatiques. Concernant ce dernier point, il est effectivement possible de demander au programme **genotypes** d'effectuer des corrections automatiques pour certaines erreurs comme la vérification des tags homozygotes dans la descendance pour assurer qu'un SNP n'est pas présent. En effet, si le modèle de détection de SNP ne peut pas identifier un site comme hétérozygote ou homozygote, le site est provisoirement marqué comme homozygote pour faciliter la recherche, par **sstacks**, de concordances avec le catalogue de loci. Si un second allèle identifié dans le catalogue (i.e., chez les parents) est présent chez un des descendant à une faible fréquence (<10 des lectures de la pile considérée), le programme **genotypes** corrige le génotype. De même, il supprimera un génotype identifié comme homozygote chez un individu particuliers si le génotype au locus est supporté par moins de 5 lectures. Les génotypes corrigés sont alors notés en majuscule. Il faut savoir que l'utilisation de l'interface web de Stacks permet de modifier manuellement les génotypes. Ceci est intéressant notamment lorsqu'un allèle est sous séquencé chez un individu. Il est alors impossible pour Stacks de le dissocier d'une erreur de séquençage, et ce génotype sera marqué comme homozygote. Si l'utilisateur sait que l'allèle alternatif existe (en observant les génotypes d'autres individus par exemple), il peut alors le corriger manuellement avant de réexécuter **genotypes**.

Reexécuter l'outil "**STACKS : genotypes Run the STACKS genotypes program**" sur le jeu de données total_output.zip en cochant cette fois l'option correction automatique.

Comparer les fichiers batch.1.genotypes.1.tsv générés par les deux jobs.

1. Les fichiers genotypes.tsv

Une ligne par locus, une colonne par individu (aa, ab, AB si correction automatique, bb, bc, ...) avec le génotype observé à chacun des loci.

# Catalog ID	Marker	Cnt	Seg Dist	progeny_1	progeny_10	progeny_11	progeny_12	progeny_13	progeny_14
1	ab/aa	95	1	ab	ab	aa	ab	ab	ab
2	ab/aa	95	1	ab	ab	ab	ab	aa	aa
3	aa/ab	94	1	-	aa	ab	aa	ab	aa
4	ab/ac	94	1	aa	ac	-	ac	ab	bc
5	aa/ab	94	1	ab	aa	ab	aa	ab	aa
6	aa/ab	93	1	-	aa	ab	aa	AB	aa
7	ab/aa	92	1	-	aa	ab	aa	aa	ab
8	ab/aa	95	1	aa	ab	aa	aa	aa	aa
9	aa/ab	90	1	-	aa	AB	aa	-	aa
10	aa/ab	95	1	AB	aa	ab	aa	aa	AB
11	ab/ac	95	1	bc	ab	ac	ab	ac	aa
12	ab/ac	94	1	ac	aa	ac	aa	ac	ab
13	ab/aa	95	1	aa	aa	aa	aa	aa	ab
14	aa/ab	94	1	-	aa	ab	aa	ab	aa
15	aa/ab	95	1	ab	aa	ab	aa	ab	aa
16	ab/aa	84	1	-	ab	ab	ab	ab	aa
17	aa/ab	95	1	ab	ab	aa	ab	aa	aa
18	ab/aa	95	1	ab	ab	ab	ab	ab	aa
19	aa/ab	95	1	ab	aa	ab	aa	ab	aa
20	ab/aa	93	1	AB	aa	ab	ab	aa	aa
21	ab/ac	94	1	ab	bc	aa	bc	bc	ab
22	aa/ab	95	1	aa	ab	aa	ab	aa	AB
23	ab/aa	95	1	aa	aa	ab	aa	ab	ab

2. Les fichiers genotypes.txt

Une ligne par individu et pour chaque individu, à chaque locus du catalogue, le génotype.

# SQL ID	Batch ID	Catalog Locus ID	Sample ID	Genotype
0	1	1	3	ab
0	1	1	4	ab
0	1	1	5	aa
0	1	1	6	ab
0	1	1	7	ab
0	1	1	8	ab
0	1	1	9	ab
0	1	1	10	aa
0	1	1	11	aa
0	1	1	12	aa
0	1	1	13	ab
0	1	1	14	ab
0	1	1	15	ab
0	1	1	16	aa
0	1	1	17	AB
0	1	1	18	aa
0	1	1	19	aa
0	1	1	20	ab
0	1	1	21	ab
0	1	1	22	ab
0	1	1	23	ab
0	1	1	24	ab
0	1	1	25	aa

3. Les fichiers haplotypes.tsv

Cnt	Seg Dist	female	male	progeny_1	progeny_10	progeny_11	progeny_12	progeny_13	progeny_14
95	0.5	A/G	A	A/G	A/G	A	A/G	A/G	A/G
95	0.5	A/T	T	A/T	A/T	A/T	A/T	T	T
94	0.5	T	C/T	-	T	C/T	T	C/T	T
95	0.1	AA/GA	GA/GG	GA	GA/GG	GA	GA/GG	AA/GA	AA/GG
94	0.5	C	A/C	A/C	C	A/C	C	A/C	C
93	0.5	A	A/T	-	A	A/T	A	T	A
92	0.05	A/G	A	-	A	A/G	A	A	A/G
95	0.5	A/T	T	T	A/T	T	T	T	T
90	0.0005	A	A/G	-	A	A	A	G	A
95	0.5	G	C/G	C	G	C/G	G	G	C
95	0.5	AC/GC	AC/GG	GC/GG	AC/GC	AC/GG	AC/GC	AC/GG	AC
94	0.5	AT/CA	AA/CA	AA/CA	CA	AA/CA	CA	AA/CA	AT/CA
95	0.5	A/G	A	A	A	A	A	A	A/G

A la fin de cette partie, nous obtenons des génotypages d'individus pouvant être utilisés dans certains logiciels permettant de créer des cartes génétiques comme Carthagène ou Mapmaker

III. Analyse RAD-seq sous Galaxy : Construction de mini contigs à partir de séquences pairées

Il est possible d'assembler des mini-contigs à partir de lectures pairées issues de RAD-seq. Disposant ainsi de plusieurs centaines de nucléotides génomiques situés à proximité de chaque marqueur, l'ancrage de ces marqueurs à des librairies d'EST, et donc leur connexion à des gènes codant pour des protéines chez d'autres organismes, est facilité.

L'outil "[STACKS : assemble read pairs by locus](#) *Run the STACKS sort_read_pairs.pl and exec_velvet.pl wrappers*" permet de rassembler les lectures pairées associées à chaque stack au sein d'un fichier fasta. Une fois un fichier fasta généré par locus du catalogue, Velvet est utilisé pour assembler les lectures de chaque fichier.

Comme le fichier fasta doit porter dans l'identifiant de la séquence l'identifiant du locus du catalogue, il est possible d'ajouter manuellement des séquences à un locus. Ainsi, en plus des mini-contigs associés automatiquement aux loci par l'outil "[STACKS : assemble read pairs by locus](#) *Run the STACKS sort_read_pairs.pl and exec_velvet.pl wrappers*", nous pouvons ajouter manuellement des séquences d'EST disponibles, ou construites de novo en utilisant des données de RNA-seq. Cette étape peut être réalisée après avoir associé des séquences au catalogue de loci en utilisant Blast ou Bowtie.

Créer un nouvel historique et renommer le (ex : Tuto GenOuest 2 RAD seq : Building mini-contig from paired-end sequences).

Récupération des données brutes dans Shared data/data libraries/1 Galaxy teaching folder/2013_GenOuest2/mini-contig (URL d'origine : http://creskolab.uoregon.edu/stacks/pe_tut.php) . Vous y trouverez

- une archive contenant 4 jeux de données de lectures pairées parentales : pe_samples.zip. Cette archive est constituée de 4 fichiers fastq, 2 par parent (*f0_female.1.fq*, *f0_female.2.fq*, *f0_male.1.fq* et *f0_male.2.fq*).

- un fichier texte "whitelist", reprenant les loci à considérer pour l'assemblage. Nous reviendrons sur ce point par la suite.

Sélectionner les 2 jeux de données, et cliquer sur le bouton "GO" après avoir vérifié que l'action sélectionnée était bien "Import to current history".

Name	Message	Data type	Date uploaded	File size
2012_Galaxy1	ChIP-seq analysys using MACS			
2012_Galaxy2	Heteroplasmy analyses			
2012_Galaxy3	Tuto RNA-seq			
2012_GenOuest1	Tuto SNP calling GL349685			
2012_GenOuest2	Tuto SNP E.T. detection			
2013_BP1	Finding Human Coding Exons with Highest SNP Density			
2013_BP2	Loading Data and Understanding Datatypes			
2013_BP3	Calling Peaks for ChIP-seq Data			
2013_BP4	Compare Datasets Using Genomic Coordinates			
2013_GenOuest1	Phylogeny, tree and Blast			
2013_GenOuest2	RAD-seq with STACKS			
Genetic map	genetic map in the spotted gar from a single linkage group			
mini-contigs	Building mini-contigs from paired-end sequences			
pe_samples.zip	Building mini-contigs from paired-end sequences of 2 F0 parents	zip	2014-02-17	700.6 MB
StacksWhitelist.txt	List of loci to consider to the assembly of mini-contig	txt	2014-02-17	47.3 KB

Vous débutez donc avec les deux jeux de données dans votre historique comme suit

The screenshot shows the Galaxy by GenOuest interface. On the right, the 'History' panel lists two jobs: 'Tuto GenOuest 2 RAD seq : STACKS building mini-contigs from paired-end sequences' (0 bytes) and 'StacksWhitelist.txt' (0 bytes). The main panel displays a welcome message and a list of tools under 'NGS: Building Loci', including 'STACKS : populations', 'Map with BWA for STACKS', 'STACKS : Reference map', 'STACKS : assemble read pairs by locus', 'STACKS : genotypes', 'STACKS : Process radtags', 'STACKS : Prepare population map file for STACKS', and 'STACKS : De novo map'.

Sélectionner l'outil "**STACKS : De novo map** Run the STACKS denovo_map.pl wrapper".

Indiquez que vous souhaitez utiliser la fonction "Genetic map" et que vous n'utilisez pas de fichiers de descendants. ("Use progeny files": No).

STACKS ne gère pas les données pairées. En effet, si d'un côté du fragment le séquençage se fera toujours au niveau du même site (correspondant au site de coupure de l'enzyme), de l'autre, ce n'est pas le cas. Du coup, une belle pile présentant une forte couverture sera déterminée dans le premier sens de séquençage et pas dans l'autre. L'utilisation des lectures pairées entraînant de nombreux biais, en tout cas, présentant un intérêt très limité lors de l'exécution de "**STACKS : De novo map** Run the STACKS denovo_map.pl wrapper", il est conseillé de n'utiliser que les lectures provenant d'un seul sens de séquençage. Toutefois, si vous disposez, comme c'est le cas ici, d'une archive comprenant des lectures pairées (avec toute les paires au complet), et que vous ne voulez pas créer

une archive particulière (avec uniquement les lectures provenant d'un sens de séquençage) pour "[STACKS : De novo map Run the STACKS denovo_map.pl wrapper](#)" et une autre pour "[STACKS : assemble read pairs by locus Run the STACKS sort_read_pairs.pl and exec_velvet.pl wrappers](#)", vous pouvez utiliser l'option "Paired-end fastq files". En cochant cette option lors de l'exécution de "[STACKS : De novo map Run the STACKS denovo_map.pl wrapper](#)", vous indiquez à l'outil de trier les fichiers fastq présents dans l'archive, et de ne s'exécuter que sur les premiers jeux de données pairées, donc les données séquencées en forward. Ainsi, la même archive pourra être utilisée en entrée de l'outil "[STACKS : assemble read pairs by locus Run the STACKS sort_read_pairs.pl and exec_velvet.pl wrappers](#)".

Vous pouvez préciser certains paramètres. Ainsi, en sélectionnant le mode Advanced pour le paramètre **Stack assembly options**, vous pouvez préciser:

-la profondeur minimum de la pile (-m), ici "Minimum number of identical". Ce paramètre, passé à `ustacks`, contrôle le nombre de lectures correspondant exactement devant être trouvé pour créer une pile chez un individu. Nous pouvons sélectionner 3, une pile ne sera générée chez un individu que si au moins 3 lectures correspondent exactement.

-la distance maximum permise entre des piles pour qu'elles soient fusionnées en un locus potentiel chez un individu. Il s'agit du paramètre `-M`, ici "Number of mismatches allowed between loci when processing a single individual". Nous précisons ici une distance maximale de 3 nucléotides.

-supprimer ou casser les RAD-tags très répétitifs. Il s'agit de supprimer les piles "bûcheronnes" de l'analyse et briser les piles modérément surdimensionnées. Nous cocherons ici cette option, correspondant au paramètre `-t`.

On choisira enfin le type de sortie : "compressed all outputs" pour générer une archive contenant tous les résultats du job.

The screenshot displays the Galaxy web interface for the 'STACKS: De novo map' tool. The tool configuration panel includes the following settings:

- Select your usage:** Genetic map
- Files containing parent sequences:** 1: pe_samples.zip
- Use progeny files:** No
- Stack assembly options:** Advanced
- Minimum number of identical:** 3
- Minimum number of identical (progeny):** -1
- Number of mismatches allowed between loci when processing a single individual:** 3
- Number of mismatches allowed when aligning secondary reads:** -1
- specify the number of mismatches allowed between loci when building the catalog:** 0
- remove, or break up, highly repetitive RAD-Tags in the stacks program:**
- disable calling haplotypes from secondary reads:**
- SNP Model Options:** Default
- Output type:** Compressed all outputs

The History panel on the right shows a list of jobs, including 'Tuto GenOuest 2 RAD seq : STACKS building mini-contigs from paired-end sequences' and '4: total_output.zip with STACKS: De novo map on data 1'.

Dans l'archive générée, nous trouvons plusieurs jeux de données:

- f0_male.2.snps.tsv : répertorie tous les snps identifiés chez l'individu considéré, ici 50 353 SNPs.
- f0_male.2.tags.tsv : présente le détail de tous les loci RAD identifiés chez l'individu considéré. Il montre la séquence consensus pour un locus et les piles de lectures constitutives.
- f0_male.2.alleles.tsv : répertorie les haplotypes identifiés à chaque locus.
- f0_male.2.matches.tsv : répertorie les haplotypes vérifiés et se référant aux loci du catalogue. Ici 292 499 haplotypes sont répertoriés. Il s'agit du dernier fichier généré.

Trois fichiers répertorient le même type d'information (loci, SNPs et allèles) mais à l'échelle de tous les échantillons.

- batch_1.catalog.tags.tsv : qui présente les 382 005 loci enregistrés dans le catalogue.
- batch_1.catalog.snps.tsv : qui présente les 72 466 SNPs enregistrés dans le catalogue.
- batch_1.catalog.alleles.tsv : qui présente les 146 753 allèles enregistrés dans le catalogue.
- batch_1.haplotypes_1.tsv : qui présente les 381 735 haplotypes enregistrés dans le catalogue
- batch_1.genotypes_1.tsv : vide car pas de descendants indiqués.
- batch_1.genotypes_1.txt : vide car pas de descendants indiqués.
- batch_1.markers.tsv : vide car pas de descendants indiqués.

En attendant que le job se termine (~50min), vous pouvez débiter la partie génomique des pop.

Sélectionner l'outil "[STACKS : assemble read pairs by locus](#) Run the STACKS sort_read_pairs.pl and exec_velvet.pl wrappers".

Galaxy by GenOuest

Analyze Data Workflow Shared Data Visualization Admin Help User

Using 48%

Tools

stacks

NGS: Building Loci

STACKS : populations Run the STACKS populations program

Map with BWA for STACKS from zip file with fastqsanger files

STACKS : Reference map Run the STACKS ref_map.pl wrapper

STACKS : assemble read pairs by locus Run the STACKS sort_read_pairs.pl and exec_velvet.pl wrappers

STACKS : genotypes Run the STACKS genotypes program

STACKS : Process radtags Run the STACKS cleaning script

STACKS : Prepare population map file for STACKS denovomap and remap

STACKS : De novo map Run the STACKS denovo_map.pl wrapper

Workflows

- All workflows

STACKS : assemble read pairs by locus (version 1.1.1)

Archive from STACKS pipeline regrouping all outputs:
4: total_output.zip with STACKS : De novo map on data 1

Archive with raw reads used to execute previous STACKS pipeline:
1: pe_samples.zip

Whitelist file containing loci that we want to assemble: those that have SNPs:
2: StacksWhitelist.txt

Specify a threshold for the minimum number of reads by locus?:
No

Specify a minimum length for assembled contigs?:
No

Execute

What it does

This program will run each of the Stacks sort_read_pairs.pl and exec_velvet.pl utilities to assemble pair-end reads from STACKS pipeline results

Created by:

Stacks was developed by Julian Catchen with contributions from Angel Amores, Paul Hohenlohe, and Bill Cresko

History

Tuto GenOuest 2 RAD seq : STACKS building mini-contigs from paired-end sequences
0 bytes

4: total_output.zip with STACKS : De novo map on data 1

3: result_log with STACKS : De novo map on data 1

2: StacksWhitelist.txt

1: pe_samples.zip

Indiquez :

- l'archive générée par l'outil "[STACKS : De novo map](#) Run the STACKS denovo_map.pl wrapper". Ici "**4: total_output.zip with STACKS : De novo map on data 1**".

- l'archive regroupant les 4 fichiers de lectures "**1: pe_samples.zip**".

- Si besoin, une liste des loci à assembler. Ici "**2: StacksWhitelist.txt**". Cette liste peut correspondre à l'ensemble des loci du catalogue de loci présentant des SNPs. Pour la créer, vous pouvez par exemple utiliser l'outil "[Cut columns from a table](#)" pour récupérer la troisième colonne du jeu de données "**batch_1.catalog.snps.tsv**" (qui répertorie tous les loci présentant des SNPs dans le catalogue de loci) généré par "[STACKS : De novo map](#) Run the STACKS denovo_map.pl wrapper", puis exécuter l'outil "[remove tab dupes](#) Finds duplicates in tab format files" sur la colonne coupée afin de disposer d'une whitelist sans redondance.

Nous ne spécifierons pas de nombre ou taille minimum de lectures pour assembler un locus.

Galaxy by GenOuest

Analyze Data Workflow Shared Data Visualization Admin Help User

Using 48%

Tools

stacks

NGS: Building Loci

STACKS : populations Run the STACKS populations program

Map with BWA for STACKS from zip file with fastqsanger files

STACKS : Reference map Run the STACKS ref_map.pl wrapper

STACKS : assemble read pairs by locus Run the STACKS sort_read_pairs.pl and exec_velvet.pl wrappers

STACKS : genotypes Run the STACKS genotypes program

STACKS : Process radtags Run the STACKS cleaning script

STACKS : Prepare population map file for STACKS denovomap and remap

STACKS : De novo map Run the STACKS denovo_map.pl wrapper

Workflows

- All workflows

The following job has been successfully added to the queue:
5: collated.fa : STACKS : assemble read pairs by locus on data 1, data 2, and data 4
You can check the status of queued jobs and view the resulting data by refreshing the History pane. When the job has been run the status will change from 'running' to 'finished' if completed successfully or 'error' if problems were encountered.

History

Tuto GenOuest 2 RAD seq : STACKS building mini-contigs from paired-end sequences
0 bytes

5: collated.fa : STACKS : assemble read pairs by locus on data 1, data 2, and data 4

4: total_output.zip with STACKS : De novo map on data 1

3: result_log with STACKS : De novo map on data 1

2: StacksWhitelist.txt

1: pe_samples.zip

IV. Analyse RAD-seq sous Galaxy : La génomique des populations

Nous allons travailler à partir des données de la publication d'*Hohenlohe et al. 2010*.

OPEN ACCESS Freely available online PLOS GENETICS

Population Genomics of Parallel Adaptation in Threespine Stickleback using Sequenced RAD Tags

Paul A. Hohenlohe^{1,2}, Susan Bassham^{1,2}, Paul D. Etter², Nicholas Stiffler², Eric A. Johnson², William A. Cresko^{1*}

1 Center for Ecology and Evolutionary Biology, University of Oregon, Eugene, Oregon, United States of America, 2 Institute of Molecular Biology, University of Oregon, Eugene, Oregon, United States of America, 3 Genomics Core Facility, University of Oregon, Eugene, Oregon, United States of America

Abstract
Next-generation sequencing technology provides novel opportunities for gathering genome-scale sequence data in natural populations, laying the empirical foundation for the evolving field of population genomics. Here we conducted a genome scan of nucleotide diversity and differentiation in natural populations of threespine stickleback (*Gasterosteus aculeatus*). We used Illumina-sequenced RAD tags to identify and type over 45,000 single nucleotide polymorphisms (SNPs) in each of 100 individuals from two oceanic and three freshwater populations. Overall estimates of genetic diversity and differentiation among populations confirm the biogeographic hypothesis that large panmictic oceanic populations have repeatedly given rise to phenotypically divergent freshwater populations. Genomic regions exhibiting signatures of both balancing and divergent selection were remarkably consistent across multiple, independently derived populations, indicating that replicate parallel phenotypic evolution in stickleback may be occurring through extensive, parallel genetic evolution at a genome-wide scale. Some of these genomic regions co-localize with previously identified QTL for stickleback phenotypic variation identified using laboratory mapping crosses. In addition, we have identified several novel regions showing parallel differentiation across independent populations. Annotation of these regions revealed numerous genes that are candidates for stickleback phenotypic evolution and will form the basis of future genetic analyses in this and other organisms. This study represents the first high-density SNP-based genome scan of genetic diversity and differentiation for populations of threespine stickleback in the wild. These data illustrate the complementary nature of laboratory crosses and population genomic scans by confirming the adaptive significance of previously identified genomic regions, elucidating the particular evolutionary and demographic history of such regions in natural populations, and identifying new genomic regions and candidate genes of evolutionary significance.

Citation: Hohenlohe PA, Bassham S, Etter PD, Stiffler N, Johnson EA, et al. (2010) Population Genomics of Parallel Adaptation in Threespine Stickleback using Sequenced RAD Tags. *PLoS Genet* 6(2): e1000862. doi:10.1371/journal.pgen.1000862

Editors: David J. Begun, University of California Davis, United States of America

Received: October 20, 2009; **Accepted:** January 28, 2010; **Published:** February 26, 2010

Copyright: © 2010 Hohenlohe et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Funding: This work was funded by grants from the National Science Foundation (IOS-0642264) and from the National Institutes of Health (1R24GM079486-01A1 and Ruth L. Kirschstein National Research Service Award F32 GM078949). The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing Interests: The authors have declared that no competing interests exist.

* E-mail: wcresko@uoregon.edu

† These authors contributed equally to this work.

Introduction
Population genetics provides a rich and mathematically rigorous framework for understanding evolutionary processes in natural populations. This theory was built over the last hundred years by modeling the processes of selection, genetic drift, mutation and migration in spatially distributed populations [1–6]. The field has concentrated primarily on the dynamics of one or a small number of genetic loci, largely because of methodological limitations. However, genes are not islands, but rather form part of a genomic community, integrated both by physical proximity on chromosomes and by various evolutionary processes [7–10]. With technological advances, such as Next Generation Sequencing (NGS) [11–13], the emerging field of population genomics now allows us to address evolutionary processes at a genomic scale in natural populations [14–20]. Population genetic measures like Wright's *F* statistics [2,21,22], traditionally viewed as point estimates, can now be examined as continuous distributions across a genome [23–29]. As a result, in addition to estimating genome-wide averages for such statistics, we can identify specific genomic regions that exhibit significantly increased or decreased differentiation among populations, indicating regions that have likely been under strong diversifying or stabilizing natural selection [9,30–41]. These signatures of selection can then be used to identify candidate pathways, genes and alleles for targeted functional analyses [42–47].
An excellent opportunity for this type of population genomics approach exists in the threespine stickleback, *Gasterosteus aculeatus* [48–50]. This small fish is distributed holarctically and inhabits a large number of marine, estuarine and freshwater habitats in Asia, Europe and North America. In many regions replicate extant freshwater stickleback populations have been independently derived from oceanic ancestors when stickleback became isolated postglacially in newly created freshwater habitats [49,51]. Population genetic data support this inference, and also indicate that present day oceanic populations can be used as surrogates for stock that gave rise to nearby derived freshwater populations [52–64]. Because of the varied selection regimes in novel habitats,

La génétique des populations est une très vieille discipline riche en théories mathématiques, et utilisant différentes approches statistiques permettant l'inférence de paramètres à partir de données génétiques. Ces statistiques se retrouvent à travers la détermination de la diversité nucléotidique (π), ou de coefficients de différentiation (i.e. F_{st}), ainsi que les mesures de covariances génétique comme le déséquilibre de liaison (D and D'). Cependant, à cause de limitations méthodologiques notamment, la majorité des travaux théoriques, statistiques et empiriques en génétique des populations s'est concentré sur un nombre restreint de loci. Avec l'avènement des NGS, des dizaines ou centaines de milliers de marqueurs génétiques peuvent désormais être examinés sur de nombreux individus, permettant à la discipline nommée génomique des population de devenir une réalité. Une nouvelle activité très excitante en génomique des populations réside dans le fait d'identifier des signatures de sélection dans les populations sauvages. Aujourd'hui, on travaillera sur des données de RAD à partir d'échantillonnage océanique et d'eau douce de populations d'épinoches.

Sous Galaxy, commencez par créer un nouvel historique (ex: "RAD 1 : SNP calling").

Les données seront récupérées via l'outil [EBI SRA ENA SRA](#) de la section Get Data. Les numéros d'accension des 6 jeux de données vont de SRR034310 à SRR034316 (attention, le moteur de recherche est sensible à la casse ;(). Nous chargerons ici le jeux de données SRR034310 en cliquant sur (Fastq files (galaxy))

The screenshot shows the EBI ENA website interface. At the top, there are navigation menus for 'EMBL-EBI', 'ENASRA', and 'European Nucleotide Archive'. A search bar contains 'SRR034310'. Below the search bar, there is a table with the following columns: Submitting Centre, Run Date, Platform, Model, Read Count, and Base Count. The table contains one row of data for SRR034310. Below the table, there is a 'Read Files' section with a table listing study accession, secondary study accession, sample accession, secondary sample accession, experiment accession, run accession, scientific name, instrument model, library layout, fastq files, submitted files, col. tax id, col. scientific name, and reference alignment.

Submitting Centre	Run Date	Platform	Model	Read Count	Base Count
University of Oregon		ILLUMINA	Illumina Genome Analyzer		

Study accession	Secondary study accession	Sample accession	Secondary sample accession	Experiment accession	Run accession	Scientific name	Instrument model	Library layout	Fastq files (ftp)	Fastq files (galaxy)	Submitted files (ftp)	Submitted files (galaxy)	Col. tax id	Col. scientific name	Reference alignment
SRR001747	SRR001747	SAMN00010786	SRR000838	SRR015871	SRR034310	Gasterosteus aculeatus	Illumina Genome Analyzer	SINGLE	File 1	File 1			13700198	Gasterosteus aculeatus Linnaeus, 1758	N

Récupération des informations concernant les barcodes et les populations liées à SRR034310 dans Shared data/data libraries/1 Galaxy teaching folder/2013_GenOuest2/population Genomics/Hohenlohe 2010

Deux possibilités ensuite.

- Si l'archive utilisée n'est pas au format *fastq.gz*, elle doit être décompressée avant utilisation dans **Process_radtags**

The screenshot shows the Galaxy web interface. The main panel displays the 'Decompress an archive (version 1.0.0)' tool. The 'Archive name' field contains '1: EBI SRA: SRR034310 File: ftp://ftp.sra.ebi.ac.uk/vol1/fastq/SRR034/SRR034310/SRR034310.fastq.gz'. The 'Merges all files into one' checkbox is checked. The 'Execute' button is visible. Below the tool, there is a 'Tool documentation' section with text explaining the tool's function and a warning about special characters in filenames. The 'Created and integrated by' section lists Cyril Monjeaud and his affiliation with GenOuest Bio-informatics Core Facility. The history panel on the right shows the tool's output as 'RAD 4 : Population Genomics' with a size of 289.5 MB.

What it does

This program examines raw reads from an Illumina sequencing run and first, checks that the barcode and the RAD outside are intact, and demultiplexes the data. If there are errors in the barcode or the RAD site within a certain allowance process_radtags can correct them. Second, it slides a window down the length of the read and checks the average quality score within the window. If the score drops below 90% probability of being correct (a raw phred score of 10), the read is discarded. This allows for some sequencing errors while eliminating reads where the sequence is degrading as it is being sequenced. By default the sliding window is 15% of the length of the read, but can be specified on the command line (the threshold and window size can be adjusted). The process_radtags program can: handle data that is barcoded, either inline or using an index, or unbarcoded; use combinatorial barcodes; check and correct for a restriction enzyme cutsite for single or double-digested data; filter adapter sequence while allowing for sequencing error in the adapter pattern; process individual files or whole directories of files; directly read gzipped data filter reads based on Illumina's Chastity filter.

- Si l'archive vient de SRA, elle doit probablement être au format *fastq.gz*. Dans ce cas, il suffit de spécifier le format comme indiqué ici :

What it does

This program examines raw reads from an Illumina sequencing run and first, checks that the barcode and the RAD outside are intact, and demultiplexes the data. If there are errors in the barcode or the RAD site within a certain allowance process_radtags can correct them. Second, it slides a window down the length of the read and checks the average quality score within the window. If the score drops below 90% probability of being correct (a raw phred score of 10), the read is discarded. This allows for some sequencing errors while eliminating reads where the sequence is degrading as it is being sequenced. By default the sliding window is 15% of the length of the read, but can be specified on the command line (the threshold and window size can be adjusted). The process_radtags program can: handle data that is barcoded, either inline or using an index, or unbarcoded; use combinatorial barcodes; check and correct for a restriction enzyme cutsite for single or double-digested data; filter adapter sequence while allowing for sequencing error in the adapter pattern; process individual files or whole directories of files; directly read gzipped data filter reads based on Illumina's Chastity filter.

L'outil STACKS : Process Radtags produit un fichier de log. Examinez le et répondez aux questions suivantes:

Combien de lectures brutes étaient présentes ?

Combien ont été retenues ?

Sur celles non retenues, quelles étaient les raisons ?

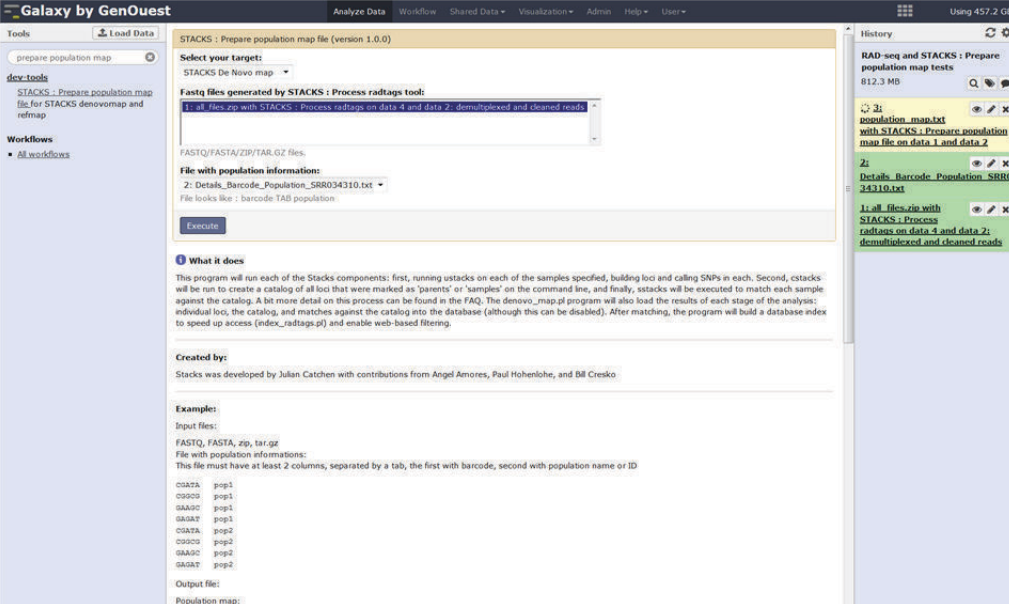
Que vous apprennent ce résultat au sujet de l'analyse de données et du design de barcodes en général?

Vous pourrez effectuer plusieurs exécution de **STACKS : Process Radtags** en jouant sur le score seuil affecté à une fenêtre glissante et en ajoutant les barcodes mentionnés dans le log précédent.

Quel est l'effet de l'augmentation du seuil du score?

Vous pourrez également jouer sur l'ensemble des paramètres de **STACKS : Process Radtags** en spécifiant notamment une mauvaise enzyme de restriction, en faisant varier le score seuil de la fenêtre glissante

Afin de générer un fichier de type population map possédant et mappant les bonnes informations, nous proposons l'utilisation de l'outil **STACKS : Prepare population map file**.



The screenshot shows the Galaxy web interface for the tool "STACKS : Prepare population map file (version 1.0.0)". The interface includes a sidebar with navigation options like "Tools", "dev-tools", and "Workflows". The main panel displays the tool's configuration, including a "Select your target" dropdown, a list of "Fastq files generated by STACKS : Process radtags tool", and a "File with population information" dropdown. The "What it does" section provides a detailed description of the tool's workflow, and the "Created by" section lists the developers. The "Example" section shows input files and a sample of the population map file format.

What it does

This program will run each of the Stacks components: first, running `ustacks` on each of the samples specified, building loci and calling SNPs in each. Second, `cstacks` will be run to create a catalog of all loci that were marked as 'parents' or 'samples' on the command line, and finally, `sstacks` will be executed to match each sample against the catalog. A bit more detail on this process can be found in the FAQ. The `denovo_map.pl` program will also load the results of each stage of the analysis: individual loci, the catalog, and matches against the catalog into the database (although this can be disabled). After matching, the program will build a database index to speed up access (`index_radtags.pl`) and enable web-based filtering.

Created by:

Stacks was developed by Julian Catchen with contributions from Angel Amores, Paul Hohenlohe, and Bill Cresko

Example:

Input files:
FASTQ, FASTA, zip, tar.gz
File with population information:
This file must have at least 2 columns, separated by a tab, the first with barcode, second with population name or ID

03ATA	pop1
03G00	pop1
0AA00	pop1
0AA0T	pop1
03ATA	pop2
03G00	pop2
0AA00	pop2
0AA0T	pop2

Output file:
Population map:

Il prend en entrée les fichiers fastq demultiplexés (organisés ou non dans une archive) ainsi qu'un fichier permettant de lier les barcodes aux populations, ici `Details_Barcode_Population_SRR034310.txt` présentant exactement 2 colonnes, la première colonne contenant les barcodes, la seconde le nom de la population. Pour faire ce fichier, il est possible de partir d'un jeu de donnée contenant plusieurs colonnes et d'utiliser l'outil Galaxy **cut** pour ne garder que les 2 colonnes nécessaires et placées dans le bon ordre (voir ci-dessous). Dans le cas présent, vérifier le format de votre fichier "Details..." pour savoir si cette manipulation est nécessaire.

Galaxy Analyze Data Workflow Shared Data Visualization Admin Help User

Tools

cut

Text Manipulation
Cut columns from a table

Workflows
All workflows

Cut (version 1.0.2)

Cut columns:
c1,c2

Delimited by:
Tab

From:
2: Details_Barcode_Population_SRR034310.txt

Execute

WARNING: This tool breaks column assignments. To re-establish column assignments run the tools and click on the pencil icon in the latest history item.

The output of this tool is always in tabular format (e.g., if your original delimiters are commas, they will be replaced with tabs). For example:

Cutting columns 1 and 3 from:

```
apple,is,good
windows,is,bad
```

will give:

```
apple good
windows bad
```

Galaxy Analyze Data Workflow Shared Data Visualization Admin Help User

Tools

cut

Text Manipulation
Cut columns from a table

Workflows
All workflows

Cut (version 1.0.2)

Cut columns:
c1,c2

Delimited by:
Tab

From:
2: Details_Barcode_Population_SRR034310.txt

Execute

WARNING: This tool breaks column assignments. To re-establish column assignments run the tools and click on the pencil icon in the latest history item.

The output of this tool is always in tabular format (e.g., if your original delimiters are commas, they will be replaced with tabs). For example:

Cutting columns 1 and 3 from:

```
apple,is,good
windows,is,bad
```

will give:

```
apple good
windows bad
```

Nous obtenons finalement un fichier de type Stacks "population map" :

Galaxy Analyze Data Workflow Shared Data Visualization Admin Help User

Data Viewer: population_map.txt with STACKS: Prepare population map file on data 11 and data 13

```
SRR034310.CCCC 0
SRR034310.CCAA 0
SRR034310.CCTT 0
SRR034310.CCGG 0
SRR034310.CACA 0
SRR034310.CAAC 0
SRR034310.CATG 0
SRR034310.CAGT 0
SRR034310.CTCT 1
SRR034310.CTAG 1
SRR034310.CTTC 1
SRR034310.CTGA 1
SRR034310.GGGG 1
SRR034310.GGAA 1
SRR034310.GGTT 1
SRR034310.GGCC 1
```

STACKS: Prepare population map file for STACKS denovomap and reMap

STACKS: De novo map Run the STACKS denovo_map.pl wrapper

Created by:
Stacks was developed by Julian Catchen with contributions from Angel Amores, Paul Hohenlohe, and Bill Cresko

Nous pouvons à présent exécuter **STACKS : de novo map** sur les individus étudiés. Il est possible de charger des fichiers FastQ ou directement l'archive "all_files.zip" générée par l'outil **STACKS : Process Raddtags** comme c'est le cas ici. Il faudra alors préciser un fichier de type STACKS "population map" si l'analyse se fait sur plusieurs populations

Vous pouvez préciser certains paramètres. Ainsi, en sélectionnant le mode Advanced pour le paramètre **Stack assembly options**, vous pouvez préciser:

-la profondeur minimum de la pile (-m), ici "Minimum number of identical". Ce paramètre, passé à **ustacks**, contrôle le nombre de lectures correspondant exactement devant être trouvé pour créer une pile chez un individu. Nous pouvons sélectionner 3, une pile ne sera générée chez un individu que si au moins 3 lectures correspondent exactement.

-la distance maximum permise entre des piles pour qu'elles soient fusionnées en un locus potentiel chez un individu. Il s'agit du paramètre -M, ici "Number of mismatches allowed between loci when processing a single individual". Nous précisons ici une distance maximale de 2 nucléotides.

-le nombre maximal de différence entre tags. Ce paramètre permet notamment de pouvoir créer un locus de type homozygote dans le catalogue alors qu'il est en réalité hétérozygote. Ceci est pratique pour conserver les loci hétérozygotes entre parents mais homozygote chez chacun d'entre eux. Il s'agit du paramètre -n, "specify the number of mismatches allowed between loci when building the catalog" ici. Nous pouvons fixer ce paramètre à 3.

-supprimer ou casser les RAD-tags très répétitifs. Il s'agit de supprimer les piles "bûcheronnes" de l'analyse et briser les piles modérément surdimensionnées. Nous cocherons ici cette option, correspondant au paramètre -t.

On choisira enfin le type de sortie : "compressed all outputs"

The screenshot displays the Galaxy web interface for the 'STACKS : De novo map (version 1.0.0)' tool. The main configuration area includes sections for 'Select your usages', 'Files containing an individual sample from a population', 'Analyzing one or more populations?', 'Specify a population map', 'Stack assembly options', and 'SNP Model Options'. The 'Output types' section is set to 'Compressed all outputs'. The 'History' panel on the right shows a list of generated files, including 'total_output.zip', 'catalog.tags', 'catalog.alleles', 'catalog.snps', 'result.log', 'discard.fastq', 'all_files.zip', 'results.log', and 'Population map'.

Une fois le job terminé, nous pouvons consulter les fichiers "result.log" et "catalog.*" (* couvrant trois types de fichiers déjà mentionnés dans la première partie consacrée à la détection de SNP, à savoir *snps*, *alleles* et *tags*).

The screenshot shows the Galaxy by GenOuest interface. The main window displays the output of the `denovo_map.pl` tool, which includes details about the de novo mapping process, such as the number of stacks identified, the maximum distance allowed between stacks, and the results of the merging process. Below the main output, there are sections for 'Blat', 'Phylostatistics', 'METADATA MANAGEMENT', and 'ISA-Tab tools'. On the right side, a 'History' panel shows a list of previous jobs, including '11: total_output.zip with STACKS : De novo map on data 5 and data 3'.

Afin de spécifier plus d'options et de pouvoir filtrer les résultats, il est possible de ré exécuter le dernier module de **Stacks : de novo map**, à savoir **populations** sur l'archive "total_output.zip" générée dans l'étape précédente.

The screenshot shows the configuration page for the 'populations' tool in Galaxy by GenOuest. The tool is titled 'STACKS : populations (version 1.0.0)'. The main configuration area includes options for 'Archive from STACKS pipeline regrouping all outputs:', 'Batch ID:', 'Specify a population map:', 'Did you want to use file output options?:', 'output results in Variant Call Format (VCF):', 'output results in GenePop Format:', 'output results in Structure Format:', 'output full sequence for each allele, from each sample locus in FASTA format:', 'output genotypes in PHASE/fastPHASE format:', 'output genotypes in Beagle format:', 'output genotypes in PLINK format:', 'output nucleotides that are fixed-within, and variant among populations in Phylip format for phylogenetic tree construction:', 'Include variable sites in the phylip output:', 'write only the first SNP per locus in Genepop and Structure outputs:', 'Did you want to use Kernel-smoothing algorithm options?:', and 'Did you want to use the genomic output option?:'. There is also a section for 'populations advanced options:'. The 'Execute' button is visible at the bottom of the configuration area. On the right side, the 'History' panel shows a list of previous jobs, including '11: total_output.zip with STACKS : De novo map on data 5 and data 3'.

Il existe différentes options avancées. Il est notamment possible de préciser :

File output options :

-VCF

-Genepop

-Structure

-FASTA

-PHASE/fastPHASE

-Beagle

-PLINK

-Phylip

Kernel-smoothing algorithm options :

-le lissage des informations obtenues en moyennant les valeurs dans une fenêtre de taille en paires de base définie. si un génome de référence est utilisé

Genomic output options :

-l'export de chaque position nucléotidique (polymorphe ou non) de tous les individus d'une population dans un fichier. Il faudra alors préciser l'enzyme de restriction utilisée

Population advanced options :

-l'utilisation d'une whitelist. Dans ce cas, un fichier text constitué d'une colonne reprenant les Stack_ID à considérer pour l'analyse sera donné en entrée.

-l'utilisation d'une blacklist. Dans ce cas, un fichier text constitué d'une colonne reprenant les Stack_ID à ne pas considérer pour l'analyse sera donné en entrée.

-le pourcentage minimum d'individus par population pour pouvoir considérer le locus dans l'analyse

-le nombre minimum de populations dans lesquels le locus doit être présent pour le considérer dans l'analyse

-la profondeur minimum d'un "stack" par individu à un locus donné

-une fréquence d'allèle alternatif / minoritaire minimum pour calculer un Fst au locus considéré

-un type de correction à appliquer aux valeurs de Fst

-un seuil de p-valeur pour conserver les valeurs de Fst

-d'effectuer un ré-échantillonnage par bootstrap en précisant la précision et le nombre. Attention, cela n'est applicable qu'à des données initialement mappées contre un génome de référence (via l'utilisation de l'outil **Reference_map**).

Plusieurs types de fichiers sont générés :

1. Le fichier result.log

Il reprend les informations concernant les loci incompatible au sein de chacune des populations puis entre les populations.

NB : Pour consulter la sortie standard dans laquelle l'outil a écrit lors du déroulement du job, il faut cliquer sur le "i" d'un des jeux de données généré par l'outil puis sur "stdout".

Tool: STACKS : populations	
Name:	result.log with STACKS : populations on data 14 and data 19
Created:	Fri Aug 8 13:56:58 2014 (UTC)
Filesize:	6.3 KB
Dbkey:	?
Format:	txt
Galaxy	
Tool	1.0.0
Version:	
Tool	
Standard	stdout
Output:	
Tool	
Standard	stderr
Error:	
Tool Exit	0
Code:	
API ID:	2f4933c7d721d5e8
Full Path:	/omaha-beach/galaxy/58/database/files/000/dataset_224.dat
Job	python /opt/shed_tools/dev-galaxy.genouest.org/repos/monjeau/stacks_toolsum
Command:	beach/galaxy/58/database/files/000/dataset_180.dat --vcf true --genepop false --
Line:	beach/galaxy/58/database/files/000/dataset_225.dat --s /omaha-beach/galaxy/5
	beach/galaxy/58/database/files/000/dataset_228.dat --os=/omaha-beach/galaxy
	beach/galaxy/58/database/tmp

2. Le fichier batch_X.sumstats.tsv

Batch ID

The batch identifier for this data set.

Locus ID

Catalog locus identifier.

Chromosome

If aligned to a reference genome.

Basepair If aligned to a reference genome. This is the alignment of the whole catalog locus. The exact basepair reported is aligned to the location of the RAD site (depending on whether alignment is to the positive or negative strand).

Column The nucleotide site within the catalog locus.

Population ID The ID supplied to the populations program, as written in the population map file.

P Nucleotide The most frequent allele at this position in this population.

Q Nucleotide The alternative allele.

Number of Individuals Number of individuals sampled in this population at this site.

P Frequency of most frequent allele.

Observed Heterozygosity The proportion of individuals that are heterozygotes in this population.

Observed Homozygosity The proportion of individuals that are homozygotes in this population.

Expected Heterozygosity Heterozygosity expected under Hardy-Weinberg equilibrium.

Expected Homozygosity Homozygosity expected under Hardy-Weinberg equilibrium.

pi An estimate of nucleotide diversity.

Smoothed pi A weighted average of p depending on the surrounding 3s of sequence in both directions.

Smoothed pi P-value If bootstrap resampling is enabled, a p-value ranking the significance of p within this population.

FIS The inbreeding coefficient of an individual (I) relative to the subpopulation (S).

Smoothed FIS A weighted average of FIS depending on the surrounding 3s of sequence in both directions.

Smoothed FIS P-value If bootstrap resampling is enabled, a p-value ranking the significance of FIS within this population.

Private allele True (1) or false (0), depending on if this allele is only occurs in this population.

#	Batch ID	Locus ID	Chr	BP	Col	Pop ID	P Nuc	Q Nuc	N	P	Obs Het	Obs Hom	Exp Het	Exp Hom	Pi	Smoothed Pi	Smoothed Pi P-value	Fis	Smoothed Fis	Smoothed Fis P-value	Private
1	1	un	20	19	1	G	T	7	0.928571	0.142857	0.857143	0.132653	0.867347	0.142857	0.0000000000	0	0.0000000000	0.0000000000	0	0	1
1	1	un	20	19	2	G		8	1	0	1	0	1	0	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	4	un	104	7	1	T	G	1	0.5	1	0	0.5	0.5	1	0.0000000000	0	0.0000000000	0.0000000000	0	1	0
1	4	un	104	7	2	T		1	1	0	1	0	1	0	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	4	un	108	11	1	G	A	1	0.5	1	0	0.5	0.5	1	0.0000000000	0	0.0000000000	0.0000000000	0	1	0
1	4	un	108	11	2	A		1	1	0	1	0	1	0	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	4	un	111	14	1	G	A	1	0.5	1	0	0.5	0.5	1	0.0000000000	0	0.0000000000	0.0000000000	0	1	0
1	4	un	111	14	2	A		1	1	0	1	0	1	0	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	4	un	115	18	1	T	C	1	0.5	1	0	0.5	0.5	1	0.0000000000	0	0.0000000000	0.0000000000	0	1	0
1	4	un	115	18	2	C		1	1	0	1	0	1	0	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	4	un	117	20	1	C	A	1	0.5	1	0	0.5	0.5	1	0.0000000000	0	0.0000000000	0.0000000000	0	1	0
1	4	un	117	20	2	A		1	1	0	1	0	1	0	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	15	un	474	25	1	A		6	1	0	1	0	1	0	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	15	un	474	25	2	A	G	7	0.928571	0.142857	0.857143	0.132653	0.867347	0.142857	0.0000000000	0	0.0000000000	0.0000000000	0	1	0
1	17	un	525	12	1	G		6	1	0	1	0	1	0	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	17	un	525	12	2	G	A	5	0.8	0.4	0.6	0.32	0.68	0.355556	0.0000000000	0	-0.1250000000	0.0000000000	0	1	0
1	19	un	608	31	1	T	G	3	0.5	1	0	0.5	0.5	0.6	0.0000000000	0	-0.6666666667	0.0000000000	0	0	0
1	19	un	608	31	2	G	T	3	0.666667	0.333333	0.444444	0.555556	0.533333	0.0000000000	0	-0.2500000000	0.0000000000	0	0	0	
1	20	un	625	16	1	C	A	1	0.5	1	0	0.5	0.5	1	0.0000000000	0	0.0000000000	0.0000000000	0	1	0
1	20	un	625	16	2	A		5	1	0	1	0	1	0	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	20	un	634	25	1	G	T	2	0.75	0.5	0.5	0.375	0.625	0.5	0.0000000000	0	0.0000000000	0.0000000000	0	1	0
1	20	un	634	25	2	G		4	1	0	1	0	1	0	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	20	un	636	27	1	G	A	2	0.75	0.5	0.5	0.375	0.625	0.5	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	20	un	636	27	2	G	A	5	0.6	0.4	0.6	0.48	0.52	0.533333	0.0000000000	0	0.2500000000	0.0000000000	0	0	0
1	24	un	755	18	1	T	C	7	0.571429	0.571429	0.428571	0.489796	0.510204	0.527473	0.0000000000	0	-0.0833333333	0.0000000000	0	0	0
1	24	un	755	18	2	T	C	8	0.625	0.5	0.5	0.46875	0.53125	0.5	0.0000000000	0	0.0000000000	0.0000000000	0	0	0
1	28	un	881	16	1	G		5	1	0	1	0	1	0	0.0000000000	0	0.0000000000	0.0000000000	0	0	0

Accompagné du fichier sumstats_summary.tsv (disposé en 2 parties pour faciliter la visualisation) :

# Batch ID	Locus ID	Pop 1 ID	Pop 2 ID	Chr	BP	Column	Overall Pi	Fst	Fisher's P	Odds Ratio	CI Low	CI High	LOD	Corrected Fst	Smoothed Fst	AMOVA Fst	Corrected AMOVA Fst	Smoothed AMOVA Fst	
1	1	0	1	un	20	19	0.0666667	0.0758293839	0.466667	0.411765	0.0337135	5.02915	0.385351	0.0758293839	0.0000000000	0.0394088670	0.0394088670	0.0000000000	
1	4	0	1	un	104	7	0.5	0.0000000000	0.5	0.333333	0.0166964	6.65479	0.477121	0.0000000000	0.0000000000	0.3333333333	0.3333333333	0.0000000000	
1	4	0	1	un	108	11	0.5	0.0000000000	0.5	3	0.150268	59.8931	0.477121	0.0000000000	0.0000000000	0.3333333333	0.3333333333	0.0000000000	
1	4	0	1	un	111	14	0.5	0.0000000000	0.5	3	0.150268	59.8931	0.477121	0.0000000000	0.0000000000	0.3333333333	0.3333333333	0.0000000000	
1	4	0	1	un	115	18	0.5	0.0000000000	0.5	3	0.150268	59.8931	0.477121	0.0000000000	0.0000000000	0.3333333333	0.3333333333	0.0000000000	
1	4	0	1	un	117	20	0.5	0.0000000000	0.5	3	0.150268	59.8931	0.477121	0.0000000000	0.0000000000	0.3333333333	0.3333333333	0.0000000000	
1	15	0	1	un	474	25	0.0769231	-0.0764331210	1	1.85714	0.149953	22.9989	0.268845	-0.0764331210	0.0000000000	0.0342857143	0.0342857143	0.0000000000	
1	17	0	1	un	525	12	0.17316	0.1675675676	0.194805	4.33333	0.386277	48.6122	0.636822	0.1675675676	0.0000000000	0.1200000000	0.1200000000	0.0000000000	
1	19	0	1	un	608	31	0.5330303	-0.0685714286	0.621212	2	0.194034	20.6149	0.30103	-0.0685714286	0.0000000000	0.0285714286	0.0285714286	0.0000000000	
1	20	0	1	un	625	16	0.166667	0.8695652174	0.166667	11	0.64645	187.176	1.04139	0.8695652174	0.0000000000	0.4545454545	0.4545454545	0.0000000000	
1	20	0	1	un	634	25	0.166667	0.4705882353	0.333333	0.222222	0.0153298	3.22136	0.653213	0.4705882353	0.0000000000	0.1818181818	0.1818181818	0.0000000000	
1	20	0	1	un	636	27	0.494505	-0.0705882353	1	2	0.149615	26.7353	0.30103	-0.0705882353	0.0000000000	0.0200000000	0.0200000000	0.0000000000	
1	24	0	1	un	755	18	0.496552	-0.0308056872	0.333333	1	0.8	0.184946	3.46046	0.09691	-0.0308056872	0.0000000000	0.0029761905	0.0029761905	0.0000000000
1	28	0	1	un	881	16	0.0769231	-0.1818181818	1	1.375	0.110603	17.0939	0.138303	-0.1818181818	0.0000000000	0.0250000000	0.0250000000	0.0000000000	
1	31	0	1	un	977	16	0.198413	-0.0567741935	0.238095	3.71429	0.365843	37.71	0.569875	-0.0567741935	0.0000000000	0.0900000000	0.0900000000	0.0000000000	
1	32	0	1	un	1009	16	0.0666667	-0.0663507109	1	1.875	0.153645	22.8815	0.273001	-0.0663507109	0.0000000000	0.0301724138	0.0301724138	0.0000000000	
1	43	0	1	un	1371	26	0.268421	-0.0700876095	0.242105	3.6	0.336807	38.479	0.556303	-0.0700876095	0.0000000000	0.1176470588	0.1176470588	0.0000000000	
1	44	0	1	un	1391	14	0.253968	0.1108870968	0.284982	0.2	0.0179743	2.2254	0.69897	0.1108870968	0.0000000000	0.0703125000	0.0703125000	0.0000000000	
1	52	0	1	un	1643	10	0.212308	-0.1133069829	0.261538	3.14286	0.306001	32.2795	0.497325	-0.1133069829	0.0000000000	0.0815217391	0.0815217391	0.0000000000	

Pop 2 ID	Chr	BP	Column	Overall Pi	Fst	Fisher's P	Odds Ratio	CI Low	CI High	LOD	Corrected Fst	Smoothed Fst	AMOVA Fst	Corrected AMOVA Fst	Smoothed AMOVA Fst	Smoothed AMOVA Fst P-value	Window SNP Count
1	un	20	19	0.0666667	0.0758293839	0.466667	0.411765	0.0337135	5.02915	0.385351	0.0758293839	0.0000000000	0.0394088670	0.0394088670	0.0000000000	0	0
1	un	104	7	0.5	0.0000000000	0.5	0.333333	0.0166964	6.65479	0.477121	0.0000000000	0.0000000000	0.3333333333	0.3333333333	0.0000000000	0	0
1	un	108	11	0.5	0.0000000000	0.5	3	0.150268	59.8931	0.477121	0.0000000000	0.0000000000	0.3333333333	0.3333333333	0.0000000000	0	0
1	un	111	14	0.5	0.0000000000	0.5	3	0.150268	59.8931	0.477121	0.0000000000	0.0000000000	0.3333333333	0.3333333333	0.0000000000	0	0
1	un	115	18	0.5	0.0000000000	0.5	3	0.150268	59.8931	0.477121	0.0000000000	0.0000000000	0.3333333333	0.3333333333	0.0000000000	0	0
1	un	117	20	0.5	0.0000000000	0.5	3	0.150268	59.8931	0.477121	0.0000000000	0.0000000000	0.3333333333	0.3333333333	0.0000000000	0	0
1	un	474	25	0.0769231	-0.0764331210	1	1.85714	0.149953	22.9989	0.268845	-0.0764331210	0.0000000000	0.0342857143	0.0342857143	0.0000000000	0	0
1	un	525	12	0.17316	0.1675675676	0.194805	4.33333	0.386277	48.6122	0.636822	0.1675675676	0.0000000000	0.1200000000	0.1200000000	0.0000000000	0	0
1	un	608	31	0.5330303	-0.0685714286	0.621212	2	0.194034	20.6149	0.30103	-0.0685714286	0.0000000000	0.0285714286	0.0285714286	0.0000000000	0	0
1	un	625	16	0.166667	0.8695652174	0.166667	11	0.64645	187.176	1.04139	0.8695652174	0.0000000000	0.4545454545	0.4545454545	0.0000000000	0	0
1	un	634	25	0.166667	0.4705882353	0.333333	0.222222	0.0153298	3.22136	0.653213	0.4705882353	0.0000000000	0.1818181818	0.1818181818	0.0000000000	0	0
1	un	636	27	0.494505	-0.0705882353	1	2	0.149615	26.7353	0.30103	-0.0705882353	0.0000000000	0.0200000000	0.0200000000	0.0000000000	0	0
1	un	755	18	0.496552	-0.0308056872	1	0.8	0.184946	3.46046	0.09691	-0.0308056872	0.0000000000	0.0029761905	0.0029761905	0.0000000000	0	0
1	un	881	16	0.0769231	-0.1818181818	1	1.375	0.110603	17.0939	0.138303	-0.1818181818	0.0000000000	0.0250000000	0.0250000000	0.0000000000	0	0
1	un	977	16	0.198413	-0.0567741935	0.238095	3.71429	0.365843	37.71	0.569875	-0.0567741935	0.0000000000	0.0900000000	0.0900000000	0.0000000000	0	0
1	un	1009	16	0.0666667	-0.0663507109	1	1.875	0.153645	22.8815	0.273001	-0.0663507109	0.0000000000	0.0301724138	0.0301724138	0.0000000000	0	0
1	un	1371	26	0.268421	-0.0700876095	0.242105	3.6	0.336807	38.479	0.556303	-0.0700876095	0.0000000000	0.1176470588	0.1176470588	0.0000000000	0	0
1	un	1391	14	0.253968	0.1108870968	0.284982	0.2	0.0179743	2.2254	0.69897	0.1108870968	0.0000000000	0.0703125000	0.0703125000	0.0000000000	0	0
1	un	1643	10	0.212308	-0.1133069829	0.261538	3.14286	0.306001	32.2795	0.497325	-0.1133069829	0.0000000000	0.0815217391	0.0815217391	0.0000000000	0	0

4. Les fichiers de sortie optionnels

Dans la capture d'écran précédente sont par exemple présentés des jeux de données au format :

- VCF, pour être vu par exemple sur des visualisateurs de génome comme IGV ou Trackster (outil embarqué dans Galaxy)

- structure, pour servir de fichier d'entrée du logiciel Structure

5. Utilisation d'un génome de référence

Après avoir récupéré l'archive générique précédemment par **STACKS : Process Radtags**, le fichier à 2 colonnes contenant les informations de population, et avoir téléchargé un génome de référence via l'UCSC (Stickleback gold, sequence), la première étape consiste à aligner les lectures démultiplexées sur le génome. Pour cela, nous utiliserons BWA à travers un outil spécialement réalisé pour STACKS : Map with BWA for STACKS.

The screenshot shows the Galaxy web interface for the tool "Map with BWA for STACKS (version 1.2.3)". The interface includes a left sidebar with "Tools" and "Stacks tool suite" sections. The main panel displays the tool configuration with the following settings:

- Will you select a reference genome from your history or use a built-in index?:** Use one from the history
- Select a reference from history:** 7: UCSC Main on Stickleback: gold (genome)
- Is this library mate-paired?:** Single-end
- FASTQ file:** 8: all_files.zip with STACKS : Process radtags on data 3 and data 1: demultiplexed and cleaned reads
- BWA settings to use:** Commonly Used
- Suppress the header in the output SAM file:** (unchecked)

At the bottom, there is an "Execute" button and a "What it does" section describing the tool: "BWA is a fast light-weighted tool that aligns relatively short sequences (queries) to a sequence database (large), such as the human reference genome. Dürbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform. Bioinformatics, 25, 1754-60."

La spécificité de cet outil est de pouvoir générer une archive contenant les fichiers d'alignement au format SAM portant le même nom que dans le fichier "population map".

Afin de générer un fichier de type population map possédant et mappant les bonnes informations, nous proposons l'utilisation de l'outil [STACKS : Prepare population map file for STACKS denovomap and refmap](#). Il prend en entrée les fichiers fastq demultiplexés (organisés ou non dans une archive) ainsi qu'un fichier permettant de lier les barcodes aux populations, ici `Details_Barcode_Population_SRR034310.txt` présentant exactement 2 colonnes, la première colonne contenant les barcodes, la seconde le nom de la population.

Nous pouvons ensuite exécuter *STACKS : Prepare population map file* en utilisant le fichier à 2 colonnes généré à l'étape précédente pour donner la relation individu / population.

The screenshot shows the Galaxy web interface. On the left is a 'Tools' sidebar with a search bar containing 'stacks' and a list of 'STACKS toolsuite' items. The main panel displays the tool configuration for 'STACKS : Prepare population map file (version 1.0.0)'. The configuration includes a 'Select your target:' dropdown set to 'STACKS De Novo map', a 'Fastq files generated by STACKS : Process radtags tool:' dropdown menu with several options, and a 'File with population information:' dropdown set to '10: Cut on data 2'. An 'Execute' button is visible at the bottom of the configuration area. Below the configuration is an information section titled 'What it does' and a 'Created by:' section.

Tools

stacks

STACKS toolsuite

- [STACKS : populations](#) Run the STACKS populations program
- [STACKS : Process radtags](#) Run the STACKS cleaning script
- [STACKS : Reference map](#) Run the STACKS ref_map.pl wrapper
- [STACKS : assemble read pairs by locus](#) Run the STACKS sort_read_pairs.pl and exec_velvet.pl wrappers
- [STACKS : genotypes](#) Run the STACKS genotypes program
- [Map with BWA for STACKS](#) from zip file with fastqsanger files
- [STACKS : Prepare population map file](#) for STACKS denovomap and refmap
- [STACKS : De novo map](#) Run the STACKS denovo_map.pl wrapper

STACKS : Prepare population map file (version 1.0.0)

Select your target:
STACKS De Novo map

Fastq files generated by STACKS : Process radtags tool:
6: decompress_an_archive.log (SRR034310.fastq)
7: UCSC Main on Stickleback: gold (genome)
8: all_files.zip with STACKS : Process radtags on data 3 and data 1: demultiplexed and cleaned reads
9: Map with BWA for STACKS on data 8 and data 7: mapped reads

FASTQ/FASTA/ZIP/TAR.GZ files.

File with population information:
10: Cut on data 2
File looks like : barcode TAB population

Execute

What it does

This program will run each of the Stacks components: first, running `ustacks` on each of the samples specified, building marked as 'parents' or 'samples' on the command line, and finally, `sstacks` will be executed to match each sample against the database index to speed up access (`index_radtags.pl`) and enable web-based filtering.

Created by:

Stacks was developed by Julian Catchen with contributions from Ángel Amores, Paul Hohenlohe, and Bill Cresko

The screenshot shows the Galaxy web interface for the 'STACKS: Reference map (version 1.0.0)' tool. The interface is divided into a left sidebar and a main content area. The sidebar lists various tools under 'STACKS toolsuite' and 'Workflows'. The main content area contains the tool's configuration options, including a 'Select your usage:' dropdown, a text input for 'Files containing an individual sample from a population:', a 'Specify a population map:' dropdown, a 'Specify the number of mismatches allowed between loci when building the catalog:' input, a 'Minimum depth of coverage:' input, 'SNP Model Options:' dropdown, and an 'Output type:' dropdown. An 'Execute' button is located at the bottom of the configuration area. Below the configuration area, there is a 'What it does' section with a detailed description of the tool's workflow and a 'Created by:' section listing the developers.

Références :

<http://evomics.org/learning/genomics/stacks/> : Cours de Julian sur Evomics. Voir en particuliers les références mentionnées. Je les prends ici pour information :

Core readings for the lecture and workshop

Amores, A., et al. 2011. Genome evolution and meiotic maps by massively parallel DNA sequencing. *Genetics* 188:799-808.

Broman, K. W. 2010. Genetic map construction with R/qtl. Univ. Wisc. Technical Report #214.

Catchen, J. M. et al. 2011. Stacks: building and genotyping loci de novo from short-read sequences. *G3: Genes, Genomes and Genetics* 1; 171-182.

Davey, J. W., et al. 2011. Genome-wide genetic marker discovery and genotyping using next-generation sequencing. *Nature Reviews Genetics* 12:499-510.

Etter, P. D., et al. 2011. SNP Discovery and Genotyping for Evolutionary Genetics using RAD sequencing. in *Molecular Methods in Evolutionary Genetics*, Rockman, M., and Orgonogozo, V., eds. (in press).

Eklom, R., and J. Galindo. 2010. Applications of next generation sequencing in molecular ecology of non-model organisms. *Heredity* 107:1-15.

Hohenlohe, P. A. et al. 2010. Population genomics of parallel adaptation in threespine stickleback using sequenced RAD tags. *PLoS Genetics* 6. 1-23.

NGS population genomics background, concepts and statistical considerations

Broman, K. W., et al. 2003. R/qtl: QTL mapping in experimental crosses. *Bioinformatics* 19:889-890.

Broman, K. W., and S. Sen. 2009. *A Guide to QTL Mapping with R/qtl*. Springer.

Gompert, Z., and C. A. Buerkle. 2011a. A hierarchical Bayesian model for next-generation population genomics. *Genetics* 187:903-917.

Gompert, Z., and C. A. Buerkle. 2011b. Bayesian estimation of genomic clines. *Molecular Ecology* 20:2111-2127.

Lynch, M. 2009. Estimation of allele frequencies from high-coverage genome-sequencing projects. *Genetics* 182:295-301.

Nielsen, R., et al. 2005. Genomic scans for selective sweeps using SNP data. *Genome Research* 15:1566-1575.

Hohenlohe, P. A., et al. 2010. Using population genomics to detect selection in natural populations: Key concepts and methodological considerations. *International Journal of Plant Sciences* 171:1059-1071.

Stapley, J., et al. 2010. Adaptation genomics: the next generation. *Trends in Ecology and Evolution* 25:705-712.

Luikart, G., et al. 2003. The power and promise of population genomics: from genotyping to genome typing. *Nature Reviews Genetics* 4:981-994.

Nielsen, R., et al. 2011. Genotype and SNP calling from next-generation sequencing data. *Nature Reviews Genetics* 12:443-451.

Genetic mapping using RRL and RAD sequencing

Altshuler, D., et al. 2000. An SNP map of the human genome generated by reduced representation shotgun sequencing. *Nature* 407:513-516.

Baxter, S. W., et al. 2011. Linkage mapping and comparative genomics using next-generation RAD sequencing of a non-model organism. *PLoS ONE* 6:e19315.

Chutimanitsakun, Y., et al. 2011. Construction and application for QTL analysis of a Restriction Site Associated DNA (RAD) linkage map in barley. *BMC Genomics* 12: 1-13.

Gore, M. A., et al. 2009. A first-generation haplotype map of maize. *Science* 326:1115-1117.

RAD-seq genotyping methodology

Baird, N. A., et al. 2008. Rapid SNP discovery and genetic mapping using sequenced RAD markers. *PLoS ONE* 3:e3376.

Emerson, K. J., et al. 2010. Resolving postglacial phylogeography using high-throughput sequencing. *Proceedings of the National Academy of Sciences* 107:16196-16200.

Etter, P. D., et al. 2011. Local De Novo Assembly of RAD Paired-End Contigs Using Short Sequencing Reads. *PLoS ONE* 6:e18561

Hohenlohe, P. A., et al. 2011. Next-generation RAD sequencing identifies thousands of SNPs for assessing hybridization between rainbow and westslope cutthroat trout. *Molecular Ecology Resources* 11 Suppl 1:117-122.

Miller, M. R., et al. 2007. Rapid and cost-effective polymorphism identification and genotyping using restriction site associated DNA (RAD) markers. *Genome Research* 17:240-248.

Willing, E. M., et al. 2011. Paired-end RAD-seq for de novo assembly and marker design without available reference. *Bioinformatics* 27:2187-2193.

Other reduced representation library (RRL) methodologies

Andolfatto, P., et al. 2011. Multiplexed shotgun genotyping for rapid and efficient genetic mapping. *Genome Research* 21:610-617.

Elshire, R. J., et al. 2011. A Robust, Simple Genotyping-by-Sequencing (GBS) Approach for High Diversity Species. *PLoS ONE* 6:e19379.

Rigola, D., et al. 2009. High-Throughput Detection of Induced Mutations and Natural Variation Using KeyPoint™ Technology. *PLoS ONE* 4:e4761.

van Orsouw, N. J., et al. 2007. Complexity reduction of polymorphic sequences (CRoPS): a novel approach for large-scale polymorphism discovery in complex genomes. *PLoS ONE* 2:e1172.

van Tassel, C. P., et al. 2008. SNP discovery and allele frequency estimation by deep sequencing of reduced representation libraries. *Nature Methods* 5:247-252.

Useful links

1. Quality scores

1. http://en.wikipedia.org/wiki/FASTQ_format
2. http://en.wikipedia.org/wiki/Phred_quality_score
3. <http://www.phrap.com/phred/>
4. http://www.illumina.com/truseq/quality_101/quality_scores.ilmn

2. Basic Unix, R and PERL commands

1. <http://mally.stanford.edu/~sr/computing/basic-unix.html>
 2. [http://korflab.ucdavis.edu/Unix and Perl/](http://korflab.ucdavis.edu/Unix_and_Perl/)
 3. <http://www.r-project.org/>
 4. <http://cran.r-project.org/doc/manuals/R-intro.html>
 5. <http://manuals.bioinformatics.ucr.edu/home/programming-in-r>
3. Stacks download and tutorials

1. <http://creskolab.uoregon.edu/stacks/>
4. Great site for information on next gen sequencing

1. <http://seqanswers.com/>