



Session de Formation

« annotation fonctionnelle & analyse du polymorphisme »

Vendredi 16 Octobre 2015

LIPM

Programme

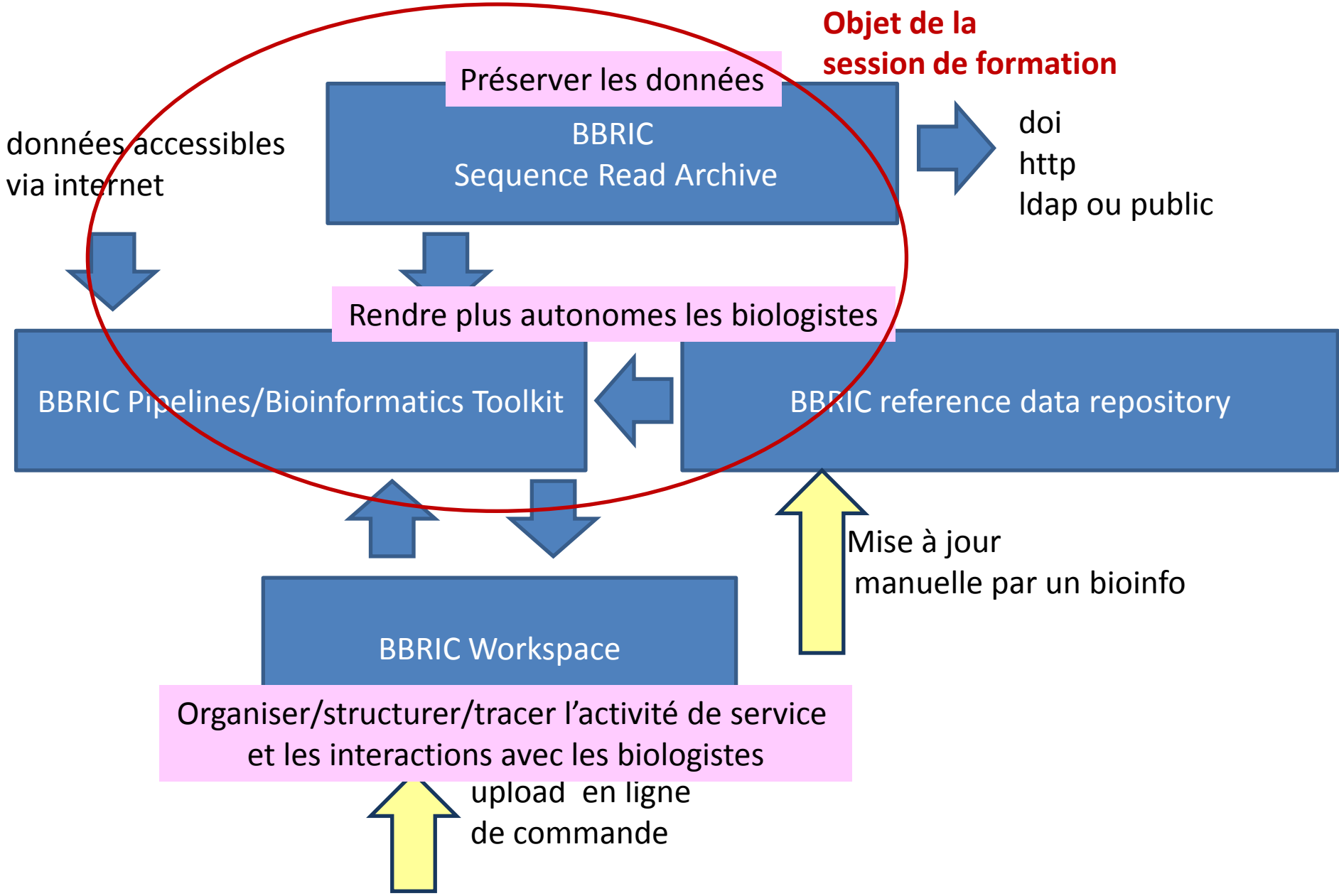
MODULE	RESPONSABLE	PLAGE
Introduction Galaxy - Archive - Reference	S. Carrere	9h-10h
Analyse du polymorphisme - Détection de SNPs à partir d'un génome de référence	L. Cottret	10h-11h30
PAUSE	10h45-11h	
Analyse du polymorphisme - Prédiction de l'effet	S. Carrere	11h30-12h
Analyse du polymorphisme - Construction d'une matrice de fréquences alléliques		12h-12h30
PAUSE	12h30-13h30	
RNASeq - Mesure de l'expression & Contrôle Qualité	E. Sallet	13h30-15h15
PAUSE	15h15-15h30	
Génomique comparative - Analyse de familles de protéines orthologues	L. Cottret	15h30-16h30
Métabolisme - Identification de voies métaboliques	L. Cottret	16h30-17h30

Modules	Expert
Détection des SNPs à partir de données de re-séquençage (<i>GATK</i>)	Fabrice Legeai (IGEPP)
Détection des SNPs (<i>samtools/VarScan</i>) pour analyse de leurs effets	Ludovic Legrand (LIPM)
Détection des SNPs (<i>samtools/VarScan</i>) pour calcul des fréquences alléliques	Ludovic Legrand (LIPM)
Mesure de l'expression a partir de données RNAseq	Ludovic Legrand (LIPM)
Indicateurs qualité d'une expérience de RNAseq	Joseph Tran (IJPB) Adeline Simon (BIOGER)
Analyse comparative de protéomes d'espèces apparentées (<i>orthoMCL</i>)	Martial Briand (IRHS) Corinne Rancurel (ISA) Sébastien Carrere & Ludovic Cottret (LIPM)

Présentation de l'architecture bioinformatique et du portail BBRIC

Sébastien Carrere

Architecture bioinformatique cible



https://bbric.toulouse.inra.fr



BBRIC Network Bioinformatics Hub

Home

Archive

Analysis Protocols

Reference

Workspace

Tutorial

Deja-vu/Bioinfo

Login



Welcome to the BBRIC network bioinformatics hub

You will find here access to:



ARCHIVE: Deposit, share and retrieve raw sequence files on the BBRIC Archive network



ANALYSIS PROTOCOLS: Access to the BBRIC analysis protocols (assembly, annotation, variant calling, etc...) to run your own analyses



REFERENCE DATA: Access to the BBRIC Reference datasets (assembly, annotation) to use for your own analyses



WORKSPACE: Retrieve the analysis results shared with you by your bioinformatics collaborators




TUTORIALS: Visualize screencasts to help you using the BBRIC resources



DEJA-VU/BIOINFO: Access to a list of bioinformatics tools tested by the BBRIC network collaborators

If you have any question feel free to contact one of us or send an email to bbric.contact@toulouse.inra.fr

If you used our computational facilities to publish results, we only ask you to send us the reference of the paper to bbric.contact@toulouse.inra.fr

 Follow @CATI_BBRIC

<https://bbric.toulouse.inra.fr>

En termes d'analyse de données

- ➔ Il n'y a pas UNE façon de faire les choses
- ➔ Pour répondre à une question, les données et les programmes changent au fil du temps
- ➔ On n'a pas forcément *a priori* connaissance de là où se trouve l'outil qui va permettre de répondre à une question
- ➔ On interroge le système BBRIC à partir d'un problème (ex: Bactérie/Assemblage)

Protocoles d'analyses BBRIC



Portail Bioinformatique de la communauté BBRIC

Home

Analysis Protocols

Archive

Login



login

Quick search:

go

Include BBRIC Archive

Query BBRIC protocols

Species:

Insects (9)
Plants (8)
Metazoa (6)
Bacteria (6)
Fungi (6)

Applications:

Gene expression data processing (4)
Transcriptome assembly (4)
Gene and gene component prediction (3)
File reformatting (2)
Sequence assembly (1)
Phylogenomics (1)
Protein site detection (1)

Launch



login

Quick search: go

Include BBRIC Archive

title	description	url	keywords	categories	domains	authors	submitter
Small genome assembly	Small genome assembly from illumina paired-end reads	link	assembly illumina genome	Sequence assembly	Bacteria	lipm.bioinfo@toulou	Ludovic.Legrand...


Le système propose la liste des programmes potentiellement pertinents avec un lien hypertexte qui renvoie vers le formulaire permettant de faire l'analyse



Authentification simplifiée

   **BBRIC Network Bioinformatics Hub**

Home Archive Analysis Protocols Workspace Tutorial Deja-vu/Bioinfo Login

 **Welcome to the BBRIC network bioinformatics**

 **BBRIC Archive Portal - Collection**

Sequences Login

 **Protein family analyses** 

OrthoMCL-Companion

Home

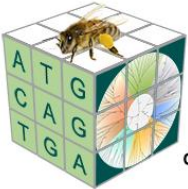
BBRIC Archive Portal






search in sequence

Analyze Data Workflow Shared Data Visualization Admin Help User

 **Bioinformatique
Biodiversité
Représentation
& Intégration
des Connaissances**

 **INRA**
SCIENCE & IMPACT
Service d'authentification de l'INRA

Vous souhaitez accéder à un service qui nécessite une authentification

Votre identifiant LDAP:

Mot de passe:

Prévenez-moi avant d'accéder à d'autres services.

En cas de difficultés, supprimez les cookies de votre navigateur.
Si le problème persiste adressez-vous à assistance-websso@inra.fr

Pour des raisons de sécurité, veuillez vous déconnecter et fermer votre navigateur lorsque vous avez fini d'accéder aux services authentifiés.

Copyright © 2005-2007 JA-SIG. All rights reserved.
Powered by [JA-SIG Central Authentication Service 3.3.1](#)

Pour accéder au service Inra - BBRIC Pipelines sélectionnez ou cherchez l'établissement auquel vous appartenez.

Se souvenir de mon choix définitivement et contourner cette étape à partir de maintenant.

Présentation de l'Archive Séquences BBRIC

Sébastien Carrere



Archive BBRIC

Objectifs

- Conserver sur le **long terme** les données « brutes » de séquence et les informations associées à la génération des données (métadonnées)
- Faciliter la **soumission** des séquences aux **banques publiques** lors des publications
- Fournir un **minimum d'information** permettant l'analyse (automatique) des données



BBRIC Archive : 182 espèces ; 2518 échantillons (une centaine de publiés!) ; 69 projets ; 7.3 Tb compressées (oct 2015)



BBRIC Archive Portal - Collection

	Sequences	Account	Admin	Logout
			Statistics	Users

Statistics by species

Species	Data size (Go)	Entry count	Molecule	Projects
Acidovorax citrulli	0.71	1	genomic_DNA	AIP Bioressources
Aeschynomene evenia	21.83	2	polyA_RNA	SESAM2
Alopecurus myosuroides	92.43	26	total_RNA	ALOMY
Arabidopsis thaliana	1235.26	216	genomic_DNA,polyA_RNA	PoolSeq48,QUANTIREX,Xanthomix,MYBLYON,ARATH-DamSeq-RRS1,RESURRECTION,CBМК,MYBSeq
Bacillus sp.	4.97	7	genomic_DNA	RAINS-Project
Brassica oleracea	214.20	22	other,total_RNA	metaSEED,Brassichax
Burkholderia caryophylli	0.78	1	genomic_DNA	AIP Bioressources
Clavibacter michiganensis sbsp insidiosus	0.52	1	genomic_DNA	AIP Bioressources
Clavibacter michiganensis sbsp michiganensis	0.39	1	genomic_DNA	AIP Bioressources
Cryphonectria parasitica	16.16	18	genomic_DNA	Cryphomic
Curtobacterium flaccumfaciens pv flaccumfaciens	0.34	1	genomic_DNA	AIP Bioressources
Dickeya dianthicola	0.52	1	genomic_DNA	AIP Bioressources
Diplocarpon rosae	5.79	1	genomic_DNA	ROADMOVIE
Ditylenchus dipsaci	7.10	1	total_RNA	Nematargets
Erwinia amylovora	0.56	1	genomic_DNA	AIP Bioressources
Erwinia sp.	1.86	2	genomic_DNA	RAINS-Project
Lolium sp.	197.82	30	total_RNA	LOLSS
Medicago sativa	386.64	4	genomic_DNA	SEQLuz
Medicago truncatula	519.18	126	other,total_RNA,polyA_RNA	DIAGNOGENE,MIRMED,Symbimics,NFYAKIN,ANR-09-BLAN-0241 MycSignalling
Meloidogyne arenaria	6.25	1	total_RNA	Nematargets
Meloidogyne incognita	22.08	9	total_RNA	Nematargets
Meloidogyne javanica	6.11	1	total_RNA	Nematargets
Microbacteriaceae sp.	0.65	1	genomic_DNA	RAINS-Project

jerome.gouzy@toulouse.inra.fr
jgouzy@inra.fr

Interface Archive (1/2)

ludovic.legrand@toulouse.inra.fr
llegrand@inra.fr



BBRIC Archive Portal - Sequence Collection

Recherche rapide

Navigation menu with tabs: Sequences, Account, Admin, Logout. Sub-menu items: Submit, Search, Browse, Manage.

Quick search Envoyer (use * for partial search, ex: 'phospho*' will match with for 'phosphorylase', 'phosphokinase')

Filter(s) Project Molecule

Filtrage et tri des données

4 results

Contributor	Date	Title	Species	Molecule	Project	More Info
CATI BBRIC	20140314	Sb-2404-G1	Sinorhizobium bbric	genomic_DNA	BBRIC	+
CATI BBRIC	20140314	Sb-2404-33-R1	Sinorhizobium bbric	total_RNA	BBRIC	+
CATI BBRIC	20140314	Sb-2404-21-R1	Sinorhizobium bbric	total_RNA	BBRIC	+
CATI BBRIC	20140314	Sb-2404-36-R1	Sinorhizobium bbric	total_RNA	BBRIC	+

All files are public All files are shared or belong to you Some files are shared All files are private

Liste des données accessibles

Détails sur les données



Interface Archive (2/2)

Visualisation des métadonnées

Téléchargement des métadonnées

Détails

Contributor	Date	Title	Species	Molecule	Project	Metadata
CATI BBRIC	20140314	Sb-2404-G1	Sinorhizobium bbric	genomic_DNA	BBRIC	
		S.bbric-250K-100x.2.fastq.gz	4.25 Mo	jerome.gouzy@toulouse.inra.fr	public	
		S.bbric-250K-100x.1.fastq.gz	4.24 Mo	jerome.gouzy@toulouse.inra.fr	public	

■ All files are public
 ■ All files are shared or belong to you
 ■ Some files are shared
 ■ All files are private

Contributor	CATI BBRIC
Date	20140314
Date update	20140314
Institution	INRA
Project	BBRIC
Extract protocol	unknown
Molecule	genomic_DNA
Species	Sinorhizobium bbric
Strain	2404
Title	Sb-2404-G1
Read length	100 (x2)
Name	Sb-2404-G1
Format	fastq
Repeat	none
Instrument model	Illumina HiSeq 2000
Sequencing center	Genotoul PLAGÉ
Average insert	300
Standard deviation	10
Type	pe
Library construction protocol	art
Library strategy	OTHER
Submitter	jerome.gouzy@toulouse.inra.fr

Téléchargement des données

Résumé des métadonnées



Soumission de séquences à l'Archive (1/4)

Choix du type de librairie

The screenshot shows the 'Sequences' submission page. At the top right, there is a 'Submit' button. Below it, the section is titled 'Select data type'. A 'Data Type:' dropdown menu is open, showing the following options: '(unoriented) single end', '(unoriented) single end', '(unoriented) paired end' (highlighted in red), 'oriented single end', 'oriented paired end', and 'mate pair'. A 'continue' button is visible to the left of the dropdown. An arrow points from the text 'Choix du type de librairie' to the dropdown menu.

Personnes impliquées dans la production de séquence

The screenshot shows the 'Specify Meta Data' page. At the top, there are navigation tabs: 'Sequences', 'Account', 'Admin', and 'Logout'. Below these are buttons for 'Submit', 'Search', 'Browse', and 'Manage'. The main section is titled 'Specify Meta Data'. It includes a table with columns 'File name', 'File size', 'Percent uploaded', and 'Server Data', with a message 'No files have been selected.' Below this, there is a welcome message for 'ludovic.legrand@toulouse.inra.fr' and a dropdown menu to select a group to share the data, currently set to 'ludovic.legrand@toulouse.inra.fr'. The 'Contributors' section is highlighted in red and contains a table with the following data:

*Contributor	
CATI BBRIC	
Add other Contributor	

Below the contributors table, there are sections for 'Sample', 'Sequencing', and 'Funding', each with an orange asterisk icon. At the bottom left, there is an 'Upload files and metadata' button. An arrow points from the text 'Personnes impliquées dans la production de séquence' to the 'Contributor' table.

Soumission (2/4): Description de l'échantillon biologique



Contributors



Sample



*Title Sb-2404-G1

Summary

Source name

Informations générales

Organisms

Add other Organism

Organism

*Species Sin

Sinorhizobium meliloti

Strain

Sinorhizobium bbric

Genotype

TaxID

Métadonnées concernant l'organisme

Characteristics

Add other Characteristic

Characteristics

Add other Characteristic

Molecule

genomic DNA
total RNA
polyA RNA
cytoplasmic RNA
nuclear RNA
genomic DNA
protein
other

Growth protocol

Treatment protocol

*Extract protocol

Description

Type of molecule that was extracted from the biological material. Include one of the following: total RNA, polyA RNA, cytoplasmic RNA, nuclear RNA, genomic DNA, protein, or other.

Protocoles appliqués à l'échantillon

Sequencing



Soumission (3a/4): Description du protocole de séquençage



Sample

Sequencing

Protocol

Library strategy: RNA-Seq

*Library construction protocol: RNA-Seq, RNA-Seq (size fractionation), RNA-Seq (CAGE), RNA-Seq (RACE), CTS, ChIP-Seq, MNase-Seq, MBD-Seq, MRE-Seq, Bisulfite-Seq, Bisulfite-Seq (reduced representation), MeDIP-Seq, DNase-Hypersensitivity, OTHER

Library construction kit

Library sequencing kit

Data processing

Basecalling protocol

Data processing steps

Genome build

Processed data files format and content

Library

Average insert

Standard deviation

Protocoles et informations concernant le séquençage

Informations provenant du centre de séquençage

Platform

*Sequencing center: Genotoul PLAGE

*Instrument model: Illumina HiSeq 2000

Collection: MiSeq

Repeat: Illumina

Format: Illumina Genome Analyzer, Illumina GAlx, 454 GS FLX Titanium, Illumina Genome Analyzer Ix, 454 GS FLX, Illumina HiSeq 2500, 454 GS FLX

Add other Files

Files

*Name: 454 GS FLX

Orientation

File

*File name: Select a file from your local directories

Plateforme de séquençage

Include one of the following models: Illumina Genome Analyzer, Illumina Genome Analyzer II, Illumina Genome Analyzer Ix, Illumina HiSeq 2000, Illumina HiSeq 1000, Illumina MiSeq, AB SOLiD System, AB SOLiD System 2.0, AB SOLiD System 3.0, AB SOLiD 4 Systemn AB SOLiD 4hq System, AB SOLiD PI System, AB SOLiD 5500xl SOLiD System, AB SOLiD 5500 SOLiD System, 454 GS, 454 GS 20, 454 GS FLX, 454 GS Junior, 454 GS FLX Titanium, Helicos HeliScope, PacBio RS, Complete Genomics, Ion Torrent PGM.



Soumission (3b/4): Upload des données

Collection

Repeat

Format

Files

*Name

Orientation

File

*File name

S.bbric-250K-100x.1.fastq.gz : 4.45 Mo

*Read length

File

*File name

S.bbric-250K-100x.2.fastq.gz : 4.45 Mo

*Read length

Nom du jeu de données

Taille des lectures

Sélection des fichiers

- directement à partir de son disque dur
- liste de fichiers après un transfert par ftp



Soumission (4/4) Financement et soumission

	Sequences	Account	Admin	Logout
	Submit	Search	Browse	Manage

Specify Meta Data

File name File size Percent uploaded Server Data

No files have been selected.

Welcome ludovic.legrand@toulouse.inra.fr

Please select a group to share the data

Contributors



Sample



Sequencing



Funding



Add other Source

Source

*Project

Institution

Upload files and metadata

Informations sur le financement



Contrôles d'intégrité et droits d'accès

- Validation des métadonnées et des données
 - contrôle de la taille des lectures
 - contrôle d'intégrité sur les paires pour les données paired-end et mate-pairs
 - ➔ Envoi d'un email (échec ou succès)
- Possibilité d'associer les données à des groupes d'utilisateurs



Best-practice @ LIPM

- → *Mail fournisseur: données disponibles*
 1. **Vous** transférez à votre **bioinformaticien**
 2. **Votre bioinformaticien** charge les fichiers dans votre espace
 3. **Vous** saisissez 1 première fiche
 4. **Le bioinformaticien** valide avec vous
 5. **Vous** utilisez cette fiche comme *template* pour vos 300 échantillons suivants

Présentation des Références BBRIC

Sébastien Carrere

Reference BBRIC

Objectifs

- **Centraliser** les données de **référence** qui sont **utiles** pour vos analyses
 - Assemblages (génomomes, transcriptomes)
 - Annotations structurales
 - Annotations fonctionnelles
 - Extractions de séquences (genes/mRNA/CDS/prot/ncRNA)
- Fournir un **minimum d'information** permettant de tracer leur origine

[References](#)[Login](#)Quick search (use * for partial search, ex: Xantho* will match all data produced on Xanthomonas)

Filter(s) Type Species

8 results

Species	Title	Type	Date	Accession	XRef	More Info
<i>Helianthus annuus</i>	HaT13l	Annotation Sequence:Proteins Sequence:mRNA	20150806		www.heliagene.org/HaT13l	+
<i>Medicago truncatula</i>	Mt20120830-Symbimics	Annotation Sequence:CDS Sequence:Genes Sequence:Genome Sequence:Proteins Sequence:mRNA Sequence:ncRNA	20150810		doi:10.1111/tpj.12442 www.ncbi.nlm.nih.gov/sra/SRP028599	+
<i>Rosa chinensis</i>	ROSCH-Transcriptome	Annotation Sequence:Proteins Sequence:mRNA	20150810		doi:10.1186/1471-2164-13-638 bbirc-archive.toulouse.inra.fr:sequence	+
<i>Lolium sp.</i>	LOLSS-Transcriptome	Annotation Sequence:Proteins Sequence:mRNA	20150806		doi:10.1007/s11103-015-0292-3 bbirc-archive.toulouse.inra.fr:sequence	+
<i>Helianthus annuus</i>	Ha412v1r1-bronze	Annotation Sequence:CDS Sequence:Genes Sequence:Genome Sequence:Proteins Sequence:mRNA Sequence:ncRNA	20150806		sunflowergenome.org:HA412.v1.1.bronze.20141015.fasta.gz www.heliagene.org:about_annotation.html	+
<i>Medicago truncatula</i>	JCVI-Mt4.0	Annotation Sequence:CDS Sequence:Genes Sequence:Genome Sequence:Proteins Sequence:mRNA Sequence:ncRNA	20150806		ftp.jcvi.org:Mt4.0v1_genes_20130731_1800.gff3 ftp.jcvi.org:JCVI.Medtr.v4.20130313.fasta	+
<i>Alopecurus myosuroides</i>	ALOMY-Transcriptome	Annotation Sequence:Proteins Sequence:mRNA	20150806		doi:10.1186/s12864-015-1804-x	+
<i>Arabidopsis thaliana</i>	TAIR10	Annotation Sequence:CDS Sequence:Genes Sequence:Genome Sequence:Proteins Sequence:mRNA Sequence:ncRNA	20150812		www.arabidopsis.org:README_TAIR10.txt	+

 All files are public All files are shared or belong to you Some files are shared All files are private



Visualisation des métadonnées

Téléchargement des métadonnées

Species	Title	Type	Date	Accession	XRef	Metadata
<i>Medicago truncatula</i>	JCVI-Mt4.0	Annotation Sequence:CDS Sequence:Genes Sequence:Genome Sequence:Proteins Sequence:mRNA Sequence:ncRNA	20150806		ftp.jcvi.org:Mt4.0v1_genes_20130731_1800.gff3 ftp.jcvi.org:JCVI.Medtr.v4.20130313.fasta	
	JCVI-Mt4.0.gff3	116.95 Mo	lipm_bioinfo, sebastien.carrere@toulouse.inra.fr, public			
	JCVI-Mt4.0_cds.fasta	96.52 Mo	lipm_bioinfo, sebastien.carrere@toulouse.inra.fr, public			
	JCVI-Mt4.0_ncrna.fasta	6.03 Ko	lipm_bioinfo, sebastien.carrere@toulouse.inra.fr, public			
	JCVI-Mt4.0_genome.fasta	399.33 Mo	lipm_bioinfo, sebastien.carrere@toulouse.inra.fr, public			
	JCVI-Mt4.0_mrna.fasta	119.32 Mo	lipm_bioinfo, sebastien.carrere@toulouse.inra.fr, public			
	JCVI-Mt4.0_embi	639.40 Mo	lipm_bioinfo, sebastien.carrere@toulouse.inra.fr, public			
	JCVI-Mt4.0_gene.fasta	169.22 Mo	lipm_bioinfo, sebastien.carrere@toulouse.inra.fr, public			
	JCVI-Mt4.0_prot.fasta	40.09 Mo	lipm_bioinfo, sebastien.carrere@toulouse.inra.fr, public			

Téléchargement direct des données

Liens vers source des données:

- Données brutes ayant été utilisées (Archive)
- Publication
- Origine des données (TAIR/NCBI)

Présentation de l'environnement web « Galaxy »

Sébastien Carrere

- ❖ Portail web
- ❖ Accès simplifié à de nombreux outils bioinformatiques
- ❖ Moins puissant que la ligne de commande, mais suffisant pour beaucoup d'analyses
- ❖ Projet international très actif, grande communauté
- ❖ Aide : tutoriaux, vidéos d'exemples
 - <http://wiki.galaxyproject.org/Learn>

❖ Beaucoup de fonctionnalités

- Nombreuses formations dédiées à Galaxy

❖ Aujourd'hui

- Utilisation de workflows conçus par BBRIC
- Lien avec l'architecture BBRIC
- Traitement de données issues de l'Archive

Galaxy / BBRIC

Analyze Data Workflow Shared Data Visualization Help User Using 33.5 MB

Tools

search tools

- Get Data
- BBRIC protocols
- QC tools
- BBRIC tools
- Tools
- NCBI Blast+
- Workflows
 - All workflows

Bioinformatique Biodiversité Représentation & Intégration des Connaissances

Compatibilité des navigateurs

Nous recommandons l'utilisation de **Firefox** et de **Chrome** pour le site Galaxy.
En ce qui concerne Internet Explorer, si le site fonctionne plutôt bien avec des versions récente (IE11), il reste des problèmes qui peuvent être gênants.

Transférer des données vers Galaxy

Depuis votre PC, utiliser l'outil **Upload File** dans la section **Get Data**
Depuis BBRIC Archive, utiliser l'outil **BBRIC Archive** dans la section **Get Data**
Depuis un PC du **centre INRA de Toulouse**, vous pouvez utiliser un client FTP (ex: Filezilla) avec les informations suivantes

serveur 147.100.164.25
protocole FTP
login votre email affiché dans le menu 'User' en haut a droi
mot de passe API key disponible dans le menu 'User' en haut a droite

History

search datasets

test_wf
5 shown, 7 hidden
31.2 MB

- 12: Genotypes
- 6: Catalog
- 4: Snp positions
- 2: progeny.zip
- 1: parents.zip

Galaxy is an open, web-based platform for data intensive biomedical research. The [Galaxy team](#) is a part of [BX at Penn State](#), and the [Biology and Mathematics and Computer Science](#) departments at [Emory University](#). The [Galaxy Project](#) is supported in part by [NHGRI](#), [NSF](#), [The Huck Institutes of the Life Sciences](#), [The Institute for CyberScience at Penn State](#), and [Emory University](#).

Lancement des outils
Visualisation des résultats
...

Liste d'outils

Historique qui contient des Datasets

Analyze Data

Workflow

Shared Data ▾

Visualization ▾

Help ▾

User ▾

❖ Analyze Data

- Lancement d'outil, visualisation de résultats

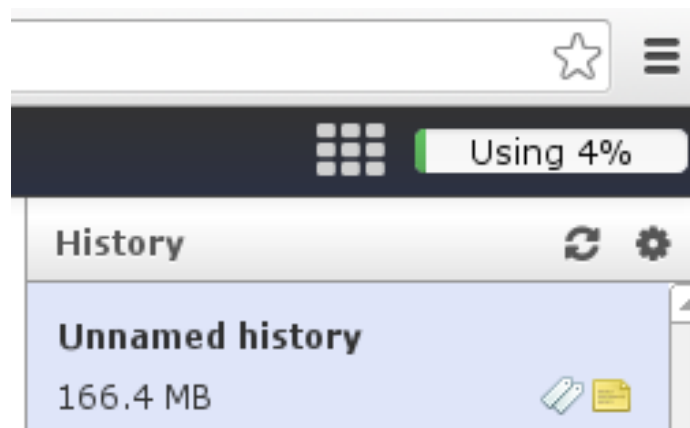
❖ Workflow

- Conception et **utilisation** de workflows

❖ Shared Data

- Accès à des données mises à disposition par les administrateurs

- ❖ Sources de données multiples
 - Upload depuis le portail web
 - **Shared Data**
 - **Archive BBRIC**
 - FTP
- ❖ Volume total limité par utilisateur
 - Quota défini par l'administrateur



Chaque dataset :

- Résultat d'un outil
- Donnée ajoutée (upload, archive)

Using 33.5 MB

History

search datasets

test_wf
5 shown, 7 [hidden](#)
31.2 MB

12: Genotypes

6: Catalog

4: Snp positions

2: progeny.zip

1: parents.zip

Aperçu du contenu

Modification

-nom

-type

Suppression

Les datasets s'empilent
au fur et à mesure

Gris = outil pas encore lancé

Jaune = l'outil s'exécute

Vert = terminé avec succès

Rouge = terminé avec une erreur

Détails d'un dataset : clic sur le titre

Télécharger le fichier

Infos sur l'origine

Relancer l'outil
avec les mêmes paramètres

Visualisation

The screenshot shows the Galaxy History interface. At the top, there's a search bar for datasets. Below it, the dataset 'test_wf' is listed with 5 shown and 7 hidden items, and a size of 31.2 MB. A green panel highlights the '12: Genotypes' dataset, showing 362 lines and 1 comment, in a tabular format. Below this, a table displays the data. At the bottom, another green panel highlights the '6: Catalog' dataset.

1	2	3	4	5			
#	Catalog	ID	Cnt	Seg	Dist	female	male
1		12	1		A	A/G	
2		12	1		C/T	T	
3		12	1		T	C/T	
4		12	1		G	A/G	
5		12	1		G	A/G	

Statistiques/aperçu

Nom de l'historique

Using 33.5 MB

History

search datasets

test_wf
5 shown, 7 [hidden](#)
31.2 MB

12: Genotypes

6: Catalog

4: Snp positions

2: progeny.zip

1: parents.zip

Menu de l'historique

The image shows a screenshot of the Galaxy web interface's History panel. A context menu is open, listing various actions. Annotations with arrows point to specific items in the menu:

- Liste de nos historiques**: Points to the **HISTORY LISTS** section header.
- Historiques partagés par d'autres**: Points to the **Shared Histories** section, which includes "Saved Histories" and "Histories Shared with Me".
- Nouvel historique**: Points to the **Create New** option under the **CURRENT HISTORY** section.
- Partage de l'historique actuel**: Points to the **Share or Publish** option.
- Création de workflow**: Points to the **Extract Workflow** option.
- Suppression de l'historique actuel**: Points to the **Delete Permanently** option.

The menu items are as follows:

- HISTORY LISTS**
- Saved Histories
- Histories Shared with Me
- CURRENT HISTORY**
- Create New
- Copy History
- Copy Datasets
- Share or Publish
- Extract Workflow
- Dataset Security
- Resume Paused Jobs
- Collapse Expanded Datasets
- Unhide Hidden Datasets
- Delete Hidden Datasets
- Purge Deleted Datasets
- Show Structure
- Export Citations
- Export to File
- Delete
- Delete Permanently
- OTHER ACTIONS**
- Import from File

Saved Histories

[Advanced Search](#)

<input type="checkbox"/>	Name	Datasets	Tags	Sharing	Size on Disk	Created	Last Updated ↑	Status
<input type="checkbox"/>	test_wf	5	0 Tags		31.2 MB	Mar 27, 2015	~4 days ago	current history
<input type="checkbox"/>	Unnamed history	2	0 Tags		4.2 MB	Apr 17, 2014	Apr 17, 2014	
<input type="checkbox"/>	Unn		0 Tags		0 bytes	Mar 06, 2014	Mar 06, 2014	

For 0 s

- Switch
- View
- Share or Publish
- Copy
- Rename
- Delete
- Delete Permanently

Histories that have not been updated for more than a time period specified by the Galaxy administrator(s) may be permanently deleted.

Lancement d'un outil

1. Choix d'un outil

The screenshot shows the Galaxy web interface. On the left is a 'Tools' sidebar with a search bar and a list of tool categories: Get Data, BBRIC protocols, QC tools, BBRIC tools, Tools, and NCBI Blast+. Under NCBI Blast+, the tool 'NCBI BLAST+ blastn' is selected. The main panel displays the configuration for 'NCBI BLAST+ blastn' (Galaxy Tool Version 0.1.01). The configuration includes:

- Nucleotide query sequence(s):** A dropdown menu showing '1: Xbbric genomic sequence'.
- Subject database/sequences:** A dropdown menu showing 'FASTA file from your history (see warning note below)'.
- Nucleotide FASTA file to use as database:** A dropdown menu showing '1: Xbbric genomic sequence'.
- Type of BLAST:** Radio buttons for 'megablast' (selected), 'blastn', 'blastn-short', and 'dc-megablast'.
- Set expectation value cutoff:** A text input field containing '0.001'.
- Output format:** A dropdown menu showing 'Tabular (extended 25 columns)'.
- Advanced Options:** A dropdown menu showing 'Hide Advanced Options'.

At the bottom of the configuration panel is an 'Execute' button. Below the configuration panel, a note states: 'Note. Database searches may take a substantial amount of time. For large input datasets it is advisable to allow overnight processing.' Below the note is the text 'What it does'.

2. Réglages des options

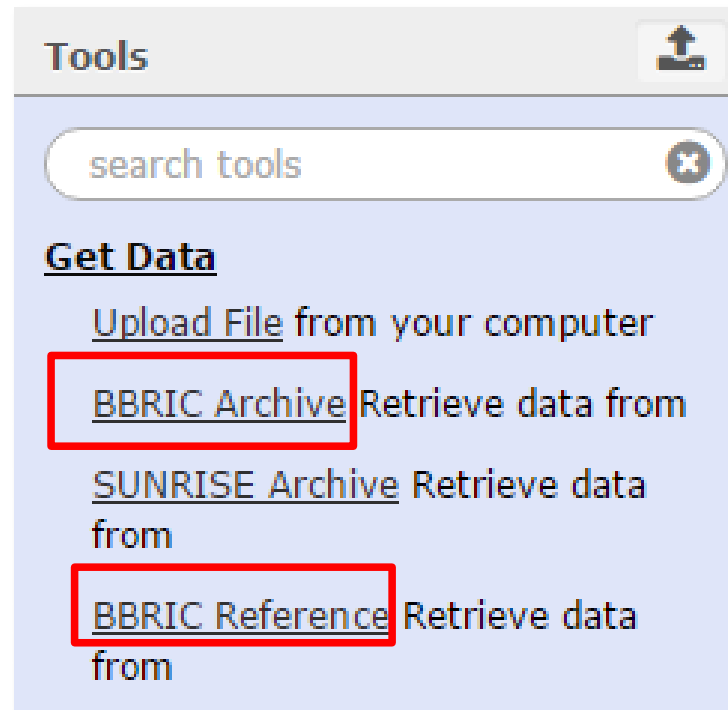
3. Lancement du calcul (créé un nouveau dataset)

❖ Comment créer les premiers datasets ?

- **Depuis l'archive BBRIC**
- **Depuis les Shared data**
- Upload d'un fichier
- ...

Depuis l'archive : outil « BBRIC Archive »

Depuis les références: outil « BBRIC Reference »



Depuis l'archive : choix d'un jeu de données

Analyze Data Workflow Shared Data Visualization Help User





BBRIC Archive Portal - Sequence Collection

Submit Sequences Search Account Browse Logout

Quick search Valider (use * for partial search, ex:'phospho*' will match with for 'phosphorylase', 'phosphokinase')

Filter(s) Project Molecule

100 results

Contributor	Date	Title	Species	Molecule	Project	More Info
Laurent, Noel Ralf, Koebnik	20130923	Xcs-CFBP4642_PE_genomic_seq	Xanthomonas cassavae	genomic_DNA	Xanthomix	
Laurent, Noel Lionel, Gagnevin	20130923	Xaa-CFBP6369_PE_genomic_seq	Xanthomonas axonopodis pv. allii	genomic_DNA	Xanthomix	
Laurent, Noel Marie-Agnes, Jacques	20130923	Xaf-CFBP3836_PE_genomic_seq	Xanthomonas alfalfae sbsp alfalfae	genomic_DNA	Xanthomix	
Laurent, Noel Marie-Agnes, Jacques	20130923	Xag-CFBP2526_PE_genomic_seq	Xanthomonas axonopodis pv. glycines	genomic_DNA	Xanthomix	

Depuis l'archive : import dans galaxy

BBRIC Archive Portal - Sequence Collection

Sequences Account Logout
Submit Search Browse Manage

Quick search Valider (use * for partial search, ex:'phospho*' will match with for 'phosphorylase', 'phosphokinase')

Contributor	Date	Title	Species	Molecule	Project	Metadata
Laurent, Noel Ralf, Koebnik	20130923	Xcs-CFBP4642_PE_genomic_seq	Xanthomonas cassavae	genomic_DNA	Xanthomix	
Xcs-CFBP4642-G1.NG-5862_CFBP4642_lib8477.pe.1.fastq.gz						3.77 Go <input checked="" type="checkbox"/>
Xcs-CFBP4642-G1.NG-5862_CFBP4642_lib8477.pe.2.fastq.gz						3.81 Go <input checked="" type="checkbox"/>
Xcs-CFBP4642-G1.NG-5862_CFBP4642_lib8477second.pe.1.fastq.gz						1.42 Go <input type="checkbox"/>
Xcs-CFBP4642-G1.NG-5862_CFBP4642_lib8477second.pe.2.fastq.gz						1.45 Go <input type="checkbox"/>

All files are public All files are shared or belong to you Some files are shared All files are private

anthony.bretaudeau@rennes.inra.fr
abretaudeau@inra.fr

Bioinformatique
Biodiversité
Représentation
Intégration
des
Connaissances

Histo
SE:
test_
5 show
31.2 M
12: Ge
6: Cat
4: Snp
2: pro
1: par

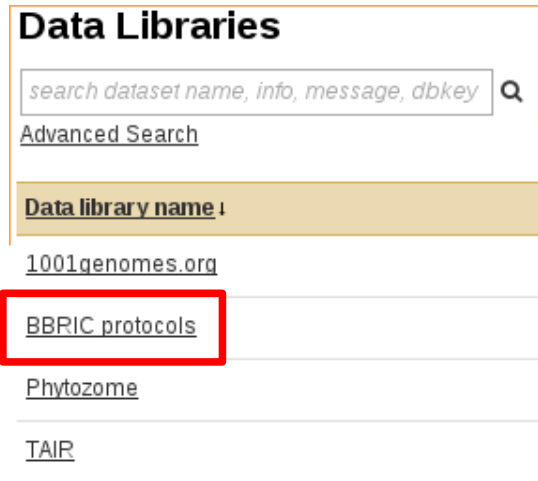
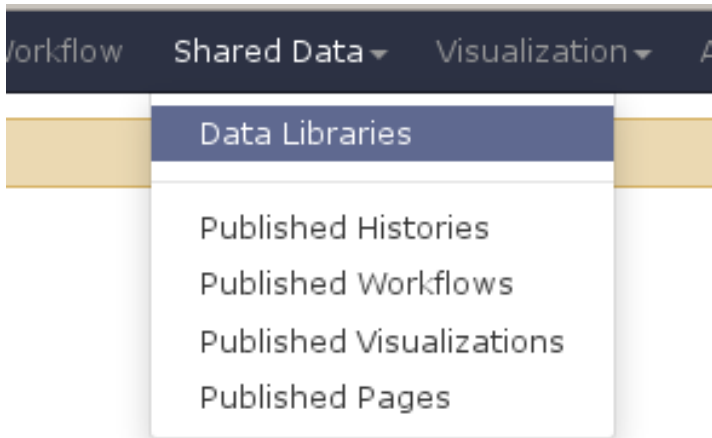
Depuis l'archive : nouveau dataset

The screenshot shows the 'History' panel in Galaxy. At the top, there is a search bar labeled 'search datasets'. Below it, a dataset named 'test_wf' is listed with a size of 7.6 GB and 7 shown items (7 hidden). Below 'test_wf', two specific files are listed:

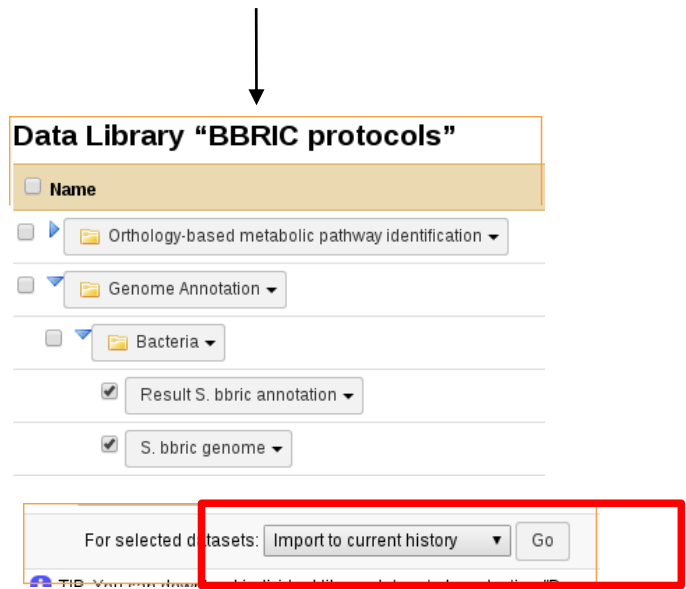
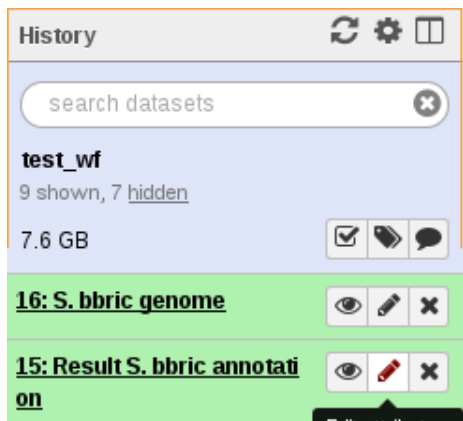
- 14: Xcs-CFBP4642-G1.NG-5**
862_CFBP4642_lib8477.p
e.2.fastq.gz
3.8 GB
format: **fastq.gz**, database: ?
- 13: Xcs-CFBP4642-G1.NG-5**
862_CFBP4642_lib8477.p
e.1.fastq.gz

Each file entry includes icons for viewing, editing, and deleting. The 'test_wf' entry has icons for selection, tagging, and commenting. The file entries also have icons for saving, information, and refreshing. A text input field at the bottom of the file entry shows 'Compressed fastq file'.

Depuis les data libraries



... BBRIC is expected to fail. To instead leave file



Upload de fichier



Download data directly from web or upload files from your disk

Name	Size	Type	Genome	Settings	Status
<input type="checkbox"/> Arabidopsis_thaliana_DNA.fasta	0.1 MB	Auto-detect	unspecified (?)		OK

You added 1 file(s) to the queue. Add more files or click 'Start' to proceed.

Choose local file Choose FTP file Paste/Fetch data **Start** Pause Reset Close

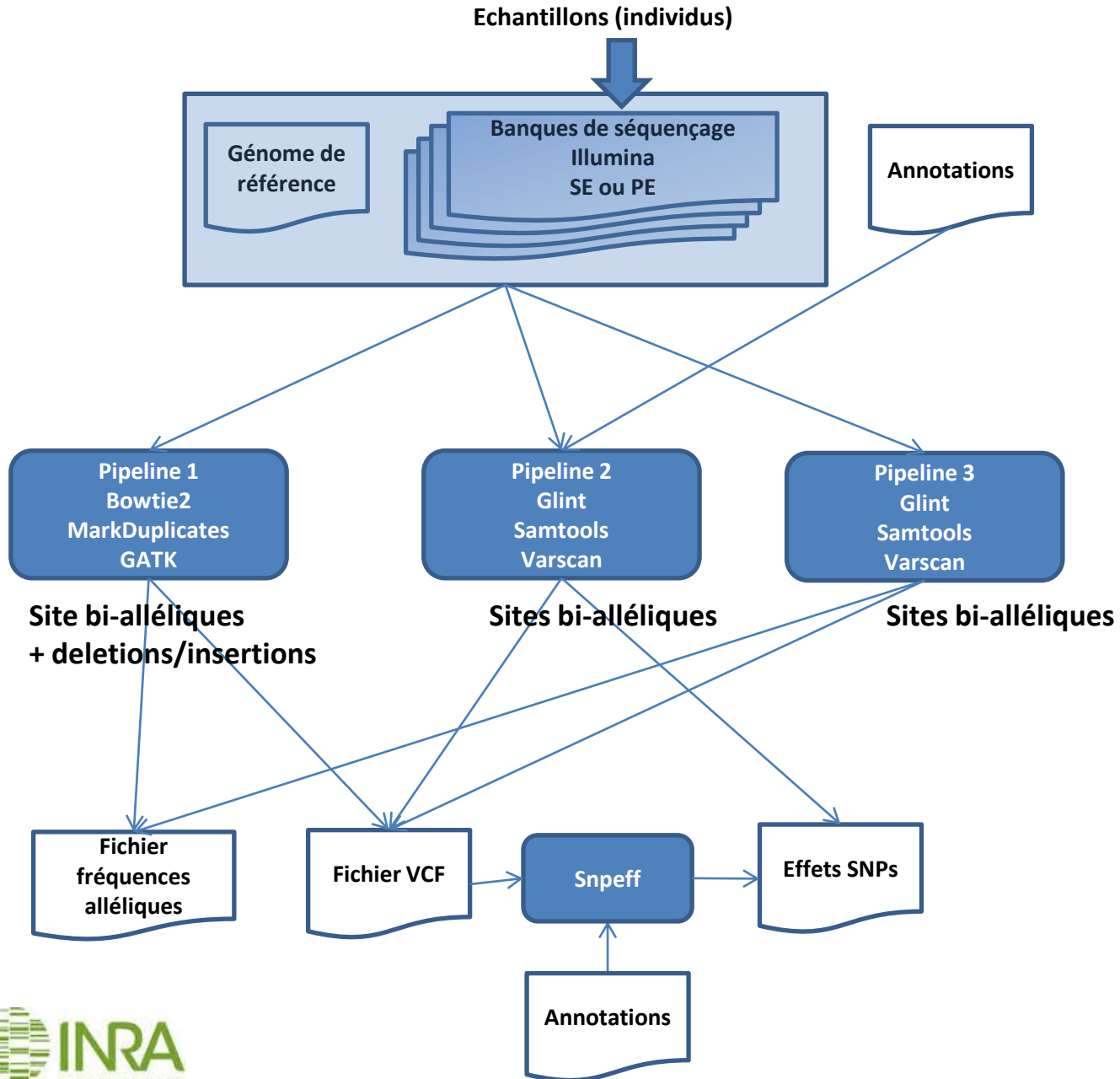
Déposer les fichiers

Download data directly from web or upload files from your disk

Name	Size	Type	Genome	Settings	Status
<input type="checkbox"/> Arabidopsis_thaliana_DNA.fasta	0.1 MB	Auto-detect	unspecified (?)		100%

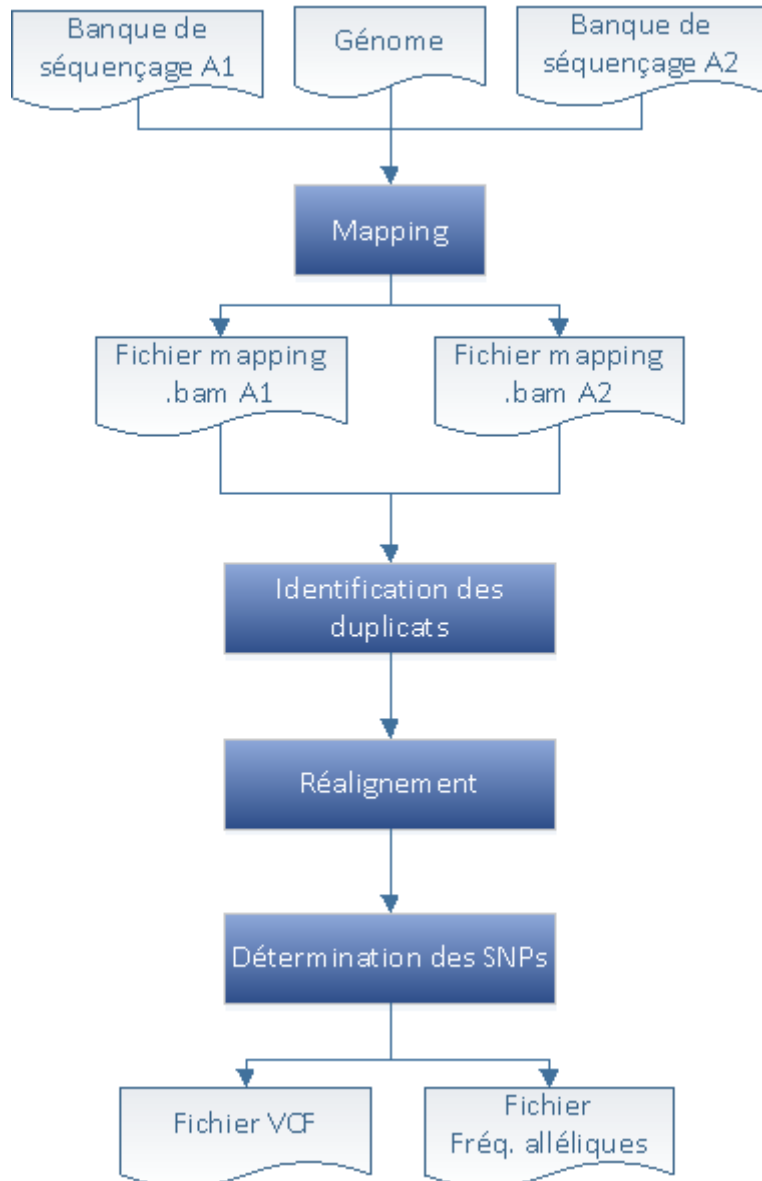


Session Polymorphisme



Détection des SNPs à partir de données de re-séquencage (GATK)

Expert : Fabrice Legeai (INRA Rennes)



Données entrées :

- Banques de séquençage de type Illumina PE ou SE
- Génome

Objectifs du pipeline :

Identifier les SNP sur un ou plusieurs individus par rapport à un génome de référence.

Limitation du pipeline :

Assez long

Assez peu documenté (voir doc GATK)

Ne fonctionne que sur des données Illumina, Single-end ou Pair-ends

Définitions :

Un **SNP** (Single Nucleotide Polymorphism) est un type de polymorphisme de l'ADN constitué par la variation d'une seule paire de base du génome entre individus d'une même espèce.

Individus 1 :CATAT**G**CGCATA.....

Individus 2 :CATAT**A**CGCATA.....

Un **allèle** représente l'une des formes que peut prendre le polymorphisme (G ou A dans le cas précédent).

Un polymorphisme est dit **hétérozygote** chez un individu ou dans une population lorsque plusieurs allèles y sont présents.

Un polymorphisme est dit **homozygote** chez un individu ou dans une population lorsque un unique allèle y est présent.

Un **génotype** est un ensemble des allèles identifiés chez un individu

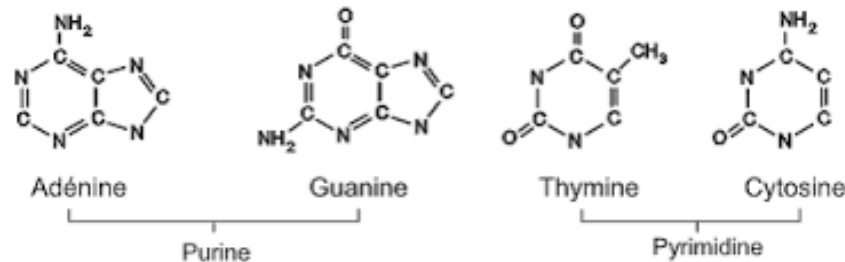
Un **haplotype** est un ensemble d'allèles situés sur un même chromosome. Un individu diploïde comporte 2 haplotypes.

Fréquence :

- Un SNP tous les 100 à 1000 nucléotides. (Génome humain de l'ordre 10^7).
- Distribution au hasard
- Moyenne 10 SNP par gènes mais certains peuvent en avoir 0.

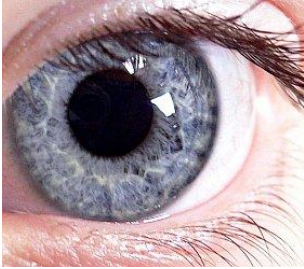
Les différents types de SNP :

- Les **Transitions** : changement d'une purine par une purine ($A \leftrightarrow G$) ou d'une pyrimidine par une pyrimidine ($C \leftrightarrow T$).
- Les **Transversions** : changement d'une purine par une pyrimidine et vice-versa (A ou $G \leftrightarrow C$ ou T).



- Les **Insertions** : ajout d'un nucléotide dans la séquence.
- Les **Délétions** : suppression d'un nucléotide dans la séquence.

Lien avec d'un phénotype



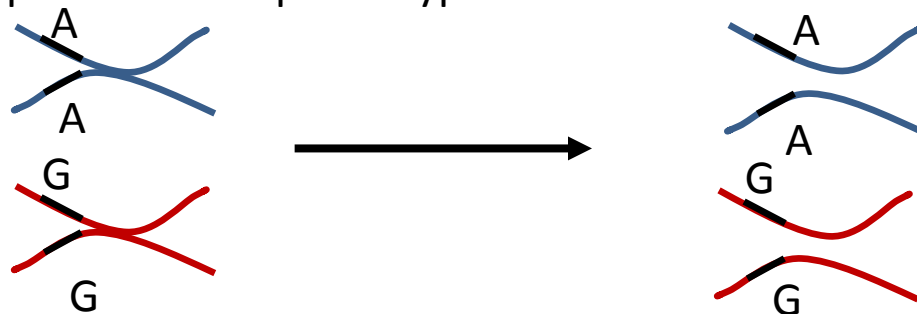
Origine des yeux bleus ?

Analyse des SNPs sur plus de 3000 individus -> le polymorphisme de type SNP sur le gène HERC2 qui régule le gène OCA2 est responsable de la pigmentation des yeux.

A global view of the OCA2-HERC2 region and pigmentation. Donnelly MP & al . Hum Genet. 2012 May;131(5):683-96. doi: 10.1007/s00439-011-1110-x. Epub 2011 Nov 8

Des **méthodes statistiques** dites d'**association** ou de **scan génomique**, permettent d'associer des régions chromosomiques à des phénotypes.

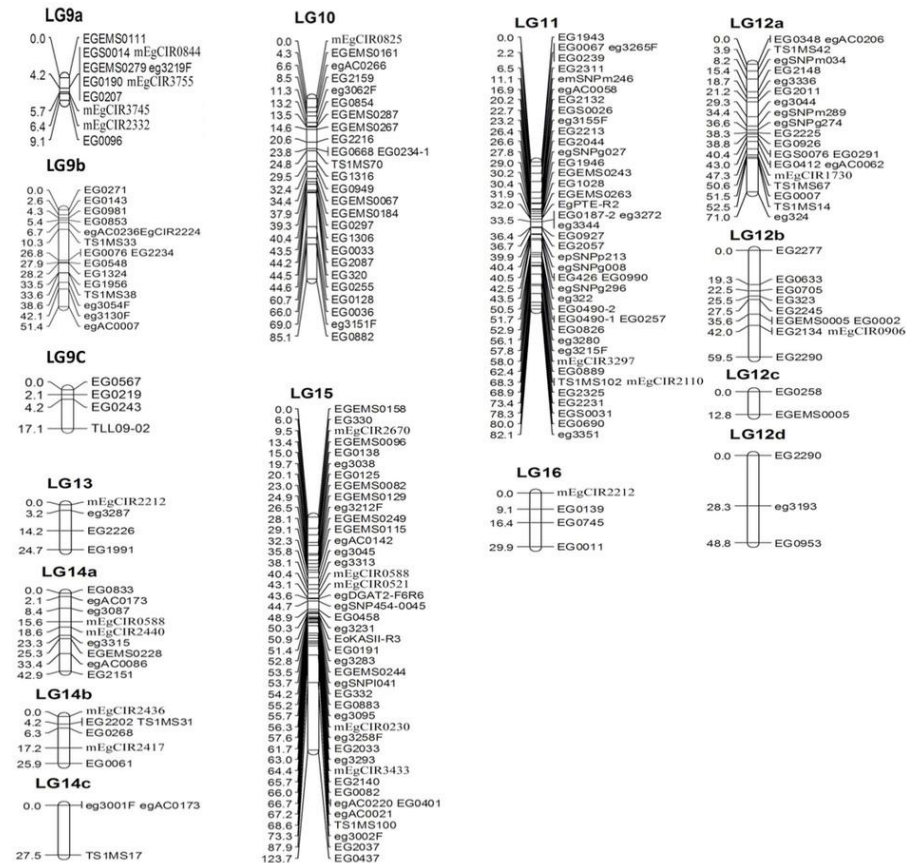
En effet, les allèles localisés à proximité de gènes responsables de phénotypes ségrégent avec les allèles responsables du phénotype.



On cherche alors à identifier des haplotypes caractéristiques dans des populations portant un caractère phénotypique particulier.

Cartes génétiques

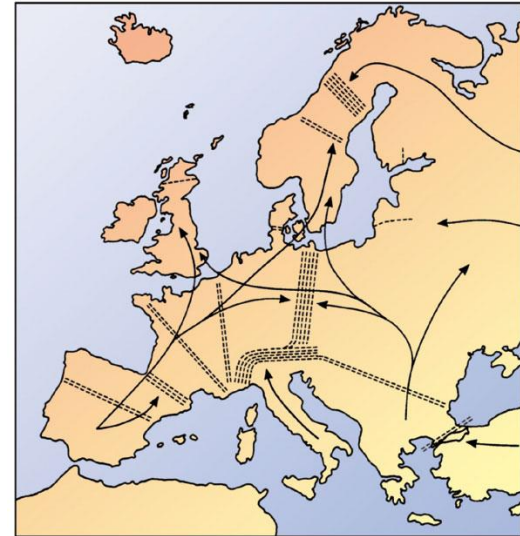
En étudiant la ségrégation des allèles (génotype) des SNPs chez des individus apparentés, des méthodes statistiques permettent de connaître la distance entre ces marqueurs et ainsi d'élaborer des cartes génétiques



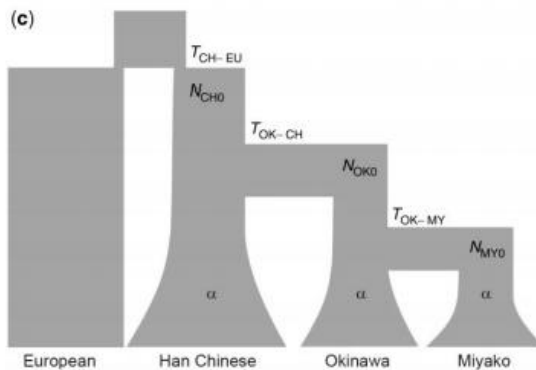
Lee et al.. A consensus linkage map of oil palm and a major QTL for stem height. Sci Rep.

En génomique des populations : Analyse des SNPs pour comprendre le déplacement et l'évolution des populations.

Le polymorphisme ou plutôt **les fréquences des polymorphismes dans les populations** sont des **marqueurs des croisements** (transferts de gènes, des échanges entre les populations), et peuvent être interprétés pour illustrer des traits d'histoire de vie ou des épisodes démographiques.



The genetic legacy of the Quaternary ice ages
Godfrey Hewitt
Nature 405, 907-913(22 June 2000)



Genome-wide SNP analysis reveals population structure and demographic history of the ryukyu islanders in the southern part of the Japanese archipelago. Sato & al. Mol Biol Evol. 2014

Mapping :

- Positionner les lectures de séquençage sur le génome de référence.

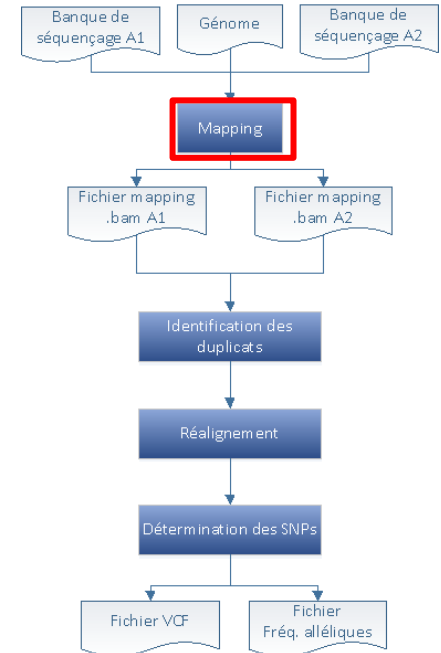
Paramétrage :

- Alignement global End to End

Ref: ---ACTGTCGACTGCGATCTCGACATCGGTGC---
 Read: GACTGGGCGATCTCGACTTCG



Ref: ---ACTGTCGACTG -- CGATCTCGACATCGGTGC---
 Read: GACTGGGCGATCTCGACTTCG



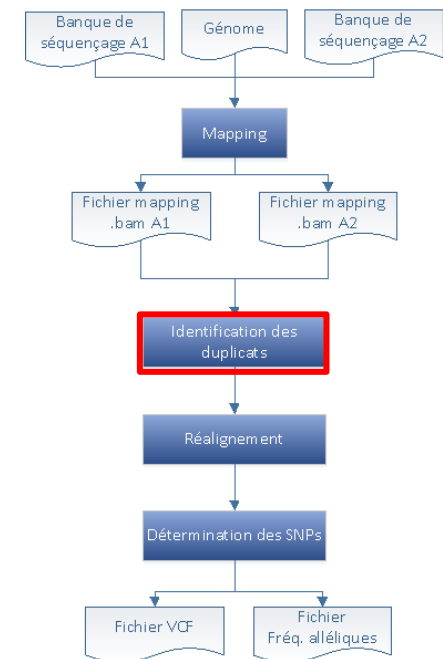
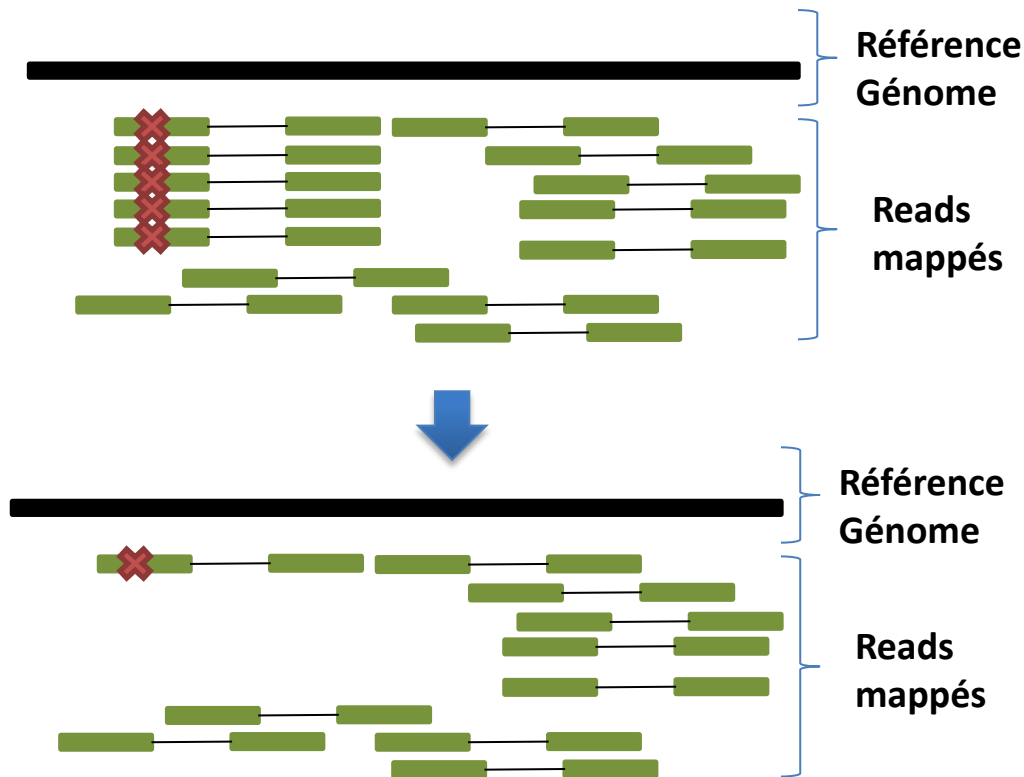
- Conserve le meilleur alignement
- Si les fragments d'ADN ont été séquencés par paires (Paired-Ends), et que les alignements des 2 lectures sont discordants, c'est à dire que la distance entre les séquences de la paire ou les orientations sont incohérentes, on élimine les alignements



 **Tool : Bowtie2**

-> Réduire le biais lié à la PCR d'amplification réalisée pendant la construction de la banque.

Enlever les duplicats techniques (même couple de séquences lu deux fois) pour éviter la surestimation des profondeurs

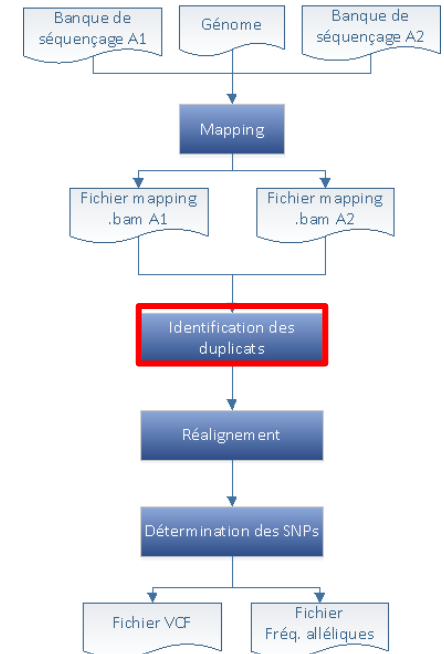


↳ Tool : Mark duplicates (picards Tools)

Réf: TAGGCCGCCGCGCTATGCCGTGCA
 R1: CGC*CGCGCTATGCCGTGC
 R2: CCGCGCTATGCCGTGC
 R3: GCC*GCGCTATGCCGTGC
 R4: GCC****GCGCTATGCCGTGC



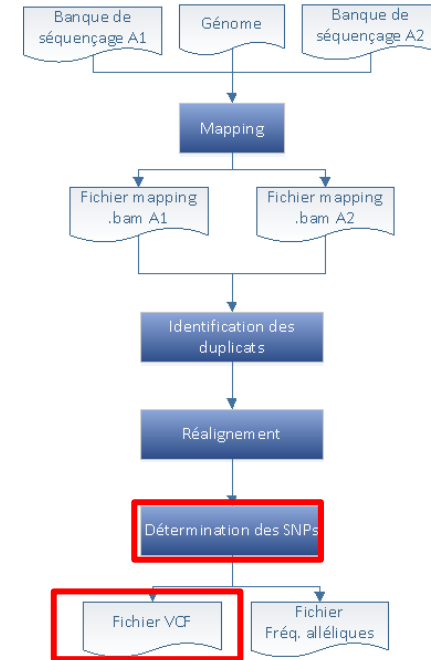
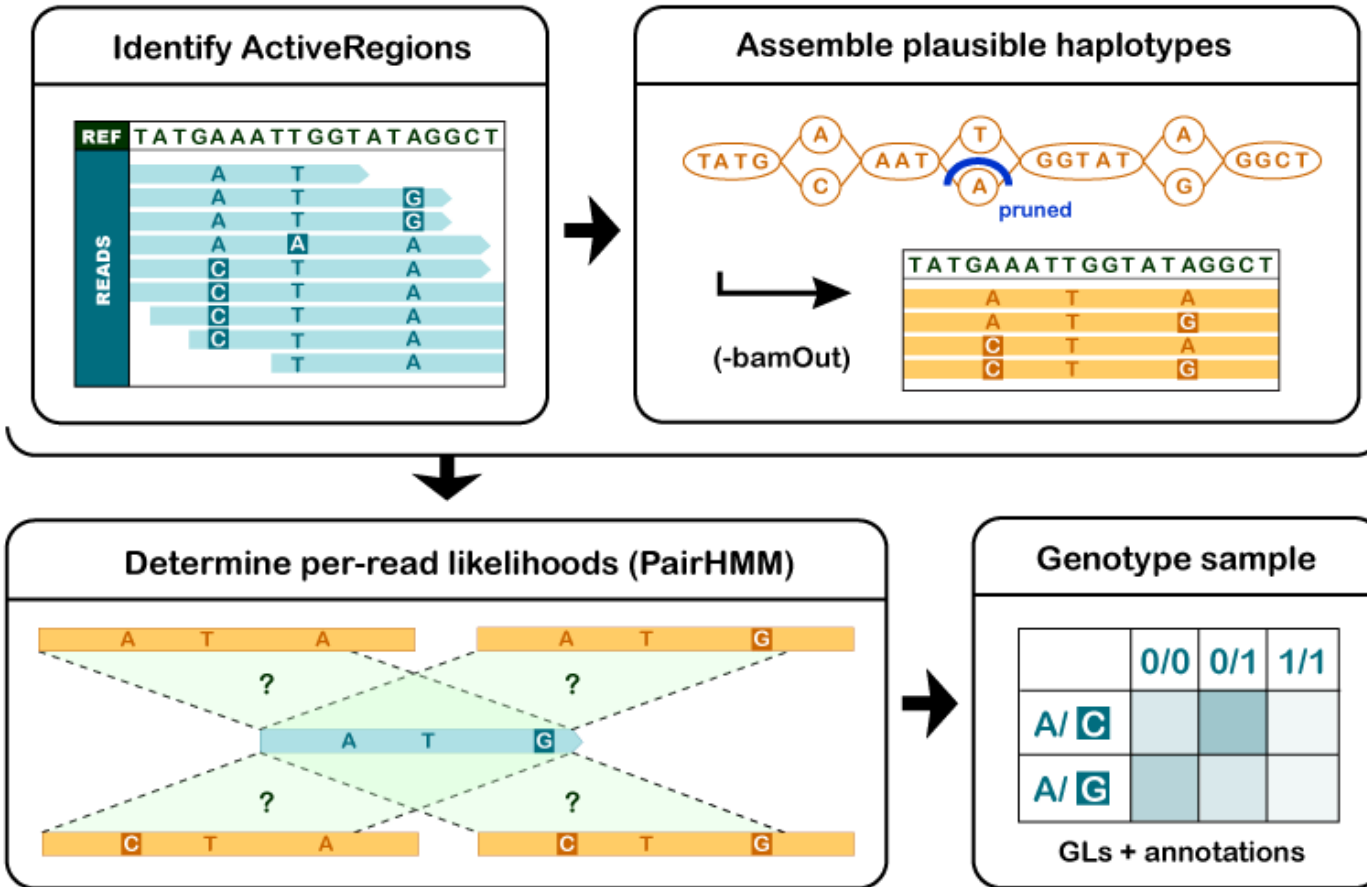
Réf: TAGGCCGCCGCGCTATGCCGTGCA
 R1: CGCC*GCGCTATGCCGTGC
 R2: CC*GCGCTATGCCGTGC
 R3: GCC*GCGCTATGCCGTGC
 R4: GCC*GCGCTATGCCGTGC



Réalignement pour l'identification des insertions /délétions sur les zones de faible complexité.



Tool : RealignerTargetCreator /IndelRealigner (GATK)



<https://www.broadinstitute.org/gatk/guide/article?id=4148>



Résultats : Fichier VCF (Variant Call Format)

Fichier Format VCF :

Format de fichier permettant de stocker la liste des variants

```
##fileformat=VCFv4.1
##fileDate=20090805
##source=myImputationProgramV3.1
##reference=file:///seq/references/1000GenomesPilot-NCBI36.fasta
##contig=<ID=20,length=62435964,assembly=B36,md5=f126cdf8a6e0c7f379d618ff66beb2da,species="Homo sap
##phasing=partial
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FILTER=<ID=q10,Description="Quality below 10">
##FILTER=<ID=s50,Description="Less than 50% of samples have data">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
```

Meta-données = description du contenu des données

Titres des colonnes

#CHROM	POS	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA00001	NA00002	NA00003
20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51	1/1:43:5:
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0.017	GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3:5:65,3	0/0:41:3
20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10;AF=0.333,0.667;AA=T;DB	GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2	2/2:35:4
20	1230237	.	T	.	47	PASS	NS=3;DP=13;AA=T	GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51	0/0:61:2
20	1234567	microsat1	GTC	G,GTCT	50	PASS	NS=3;DP=9;AA=G	GT:GQ:DP	0/1:35:4	0/2:17:2	1/1:40:3

Ecole bioinfo Roscoff 18-22/11/2013

Informations sur chaque site

Génotypes: format + n échantillons

#chr	pos	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA00001	NA00002	NA00004
20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51	1/1:43:5 , .
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0,017	GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3::5:65,3	0/01:41:3
20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10,AF=0.333,0.667;AA=T;DB	GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2	2/2:35:4
20	1230237	.	T	.	47	PASS	NS=3,DP=13;AA=T	GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51	0/0:61:2

- Chromosome
- Position
- Nom
- Allèle sur la référence
- Allèle(s) alternatifs
- **QUAL** Qualité
- **FILTER**
- **INFO**

Informations sur le génotype
(site x échantillon)

- Critères de sélection ou de filtrage :
- Les informations générales du site (lignes)
 - Les échantillons (colonnes)
 - Les génotypes (sites x échantillons)

Exemple	REF	ALT
Substitution	G	A
Insertion	C	CT
DéletionCC	C	
2 Génotypes	T	A,G

Score de qualité phred	Probabilité d'une identification incorrecte	Précision de l'identification d'un base
10	1 pour 10	90 %
20	1 pour 100	99 %
30	1 pour 1000	99.9 %
40	1 pour 10000	99.99 %
50	1 pour 100000	99.999 %



Informations sur chaque site

Génotypes: format + n échantillons

#chr	pos	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA00001	NA00002	NA00004
20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51	1/1:43:5 ,.
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0,017	GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3::5:65,3	0/01:41:3
20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10,AF=0.333,0.667;AA=T;DB	GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2	2/2:35:4
20	1230237	.	T	.	47	PASS	NS=3,DP=13;AA=T	GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51	0/0:61:2

INFO : Liste de champs + valeurs

- champ1=xx:champ2=yy:champ3=zz
- Exemple ici : NS, DP, AF, AA, DB, H2
- Signification : cf méta données

```
##INFO=<ID=NS,Number=1,Type=Integer,Description="Number of Samples With Data">
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total Depth">
##INFO=<ID=AF,Number=A,Type=Float,Description="Allele Frequency">
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">
##INFO=<ID=DB,Number=0,Type=Flag,Description="dbSNP membership, build 129">
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
```

Les champs ont des formes différentes

- DP (Total depth) contient 1 entier (number=1, type=Integer)
- AF (Allele Frequency) contient une liste de valeurs décimales car il peut y avoir plus d'1 allèle alternatif (number=Array, type=float)
- H2 (Hapmap2 membership) est présent ou absent (number=0, type=flag)

Les sites n'ont pas tous les mêmes champs

Informations sur chaque site

Génotypes: format + n échantillons

#chr	pos	ID	REF	ALT	QUAL	FILTER	INFO	FORMAT	NA00001	NA00002	NA00004
20	14370	rs6054257	G	A	29	PASS	NS=3;DP=14;AF=0.5;DB;H2	GT:GQ:DP:HQ	0 0:48:1:51,51	1 0:48:8:51,51	1/1:43:5 ,.
20	17330	.	T	A	3	q10	NS=3;DP=11;AF=0,017	GT:GQ:DP:HQ	0 0:49:3:58,50	0 1:3::5:65,3	0/01:41:3
20	1110696	rs6040355	A	G,T	67	PASS	NS=2;DP=10,AF=0.333,0.667;AA=T;DB	GT:GQ:DP:HQ	1 2:21:6:23,27	2 1:2:0:18,2	2/2:35:4
20	1230237	.	T	.	47	PASS	NS=3,DP=13;AA=T	GT:GQ:DP:HQ	0 0:54:7:56,60	0 0:48:4:51,51	0/0:61:2

Format des informations des colonnes génotype de tous les échantillons

- champ1:champ2:champ3
- Exemple ici : GT, GQ, DP, HQ
 - Signification : cf méta données

Valeurs des génotypes ; site x échantillon

```
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=HQ,Number=2,Type=Integer,Description="Haplotype Quality">
```

Les informations ont des formes différentes

- DP (Read depth) contient 1 entier (number=1, type=Integer)
- GT (Génotype) contient un 1 mot (number=1, type=string)
- HQ (Haplotype Quality) contient 2 entiers (number=2, type=Integer)

Les sites n'ont pas tous les mêmes champs

- Comment traiter les données manquantes lors du filtrage ? => option des outils

GT :
 Unphased : allèle1 / allèle 2
 Phased : allèle1 | allèle 2

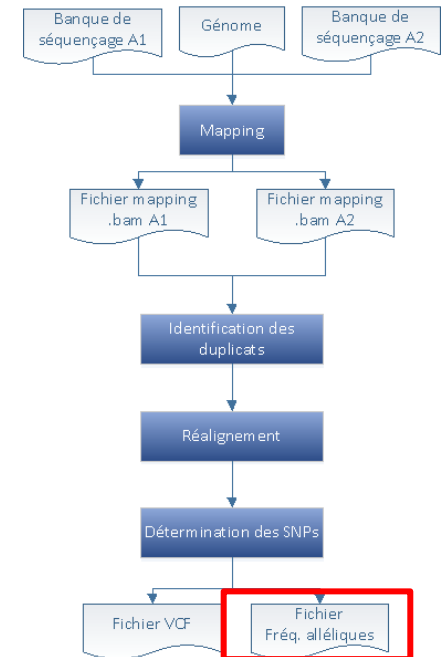
Allèle :
 0=REF, 1=ALT1, 2=ALT2 ...



CHROM	POS	REF	ALT	#ALLELES	ind1 #REF	ind1 #TOT	ind2 #REF	ind2 #TOT
GL349685	5347	G	A	1	0	2	0	0
GL349685	5353	G	T	1	0	2	1	1
GL349685	5395	G	T	1	0	2	3	3
GL349685	5397	T	C	1	0	2	3	3
GL349685	5458	T	G	1	0	2	0	8
GL349685	5513	T	TAATCGCC	1	0	2	9	9
GL349685	5522	T	A	1	0	2	9	9
GL349685	5644	GAC	G	1	0	0	0	0
GL349685	6425	T	C	1	0	2	4	4
GL349685	6431	G	GGT	1	0	2	4	4
GL349685	6661	C	T	1	0	10	7	7
GL349685	6672	C	T	1	2	11	7	7

ind2 #REF ind2 #TOT

9



- CHROM : scaffold ou chromosome
- POS : position sur la référence
- REF : la séquence sur la référence
- ALT : la séquence alternative
- #ALLELES : Nombre d'allèles différents de la référence
- ind1 #REF : Nombre de reads correspondant à la référence pour échantillon 1
- ind1 #TOT : Nombre de reads total sur ce site pour échantillon 1
- ind2 #REF : Nombre de reads correspondant à la référence pour échantillon 2
- ind2 #TOT : Nombre de reads total sur ce site pour échantillon 2

Banques Single-End

Variant detection and allele frequency (version 1.0)

reference genome: ①

input files

input files 1

Target source: ②

Reads group: ③

input: ④

Add new input files

Execute

Banques Paire-End

Variant detection and allele frequency (version 1.0)

reference genome: ①

input files

input files 1

Target source: ②

Reads group: ③

input 1: ④

input 2: ④

minimal insert size: ⑤

maximal insert size: ⑤

Add new input files

Execute

index

- ① Génome de référence
- ② Choix du type de librairie
- ③ Nom de l'échantillon
- ④ Choix de la librairie
- ⑤ taille du fragment séquencé (en cas de séquençage pair-end)

POLYMORPHISM

[Genetic variant annotation and effect prediction \(snpEff\)](#)

[SNP detection and effect prediction](#)

[SNP detection without reference](#)

[SNP detection for building an allelic frequency matrix](#)

[Variant detection and allele frequency GATK based](#)

Je veux identifier des SNPs et des indels à partir de données de séquençage et d'un génome de référence

Je vais sur

L'outil permet

L'outil ne permet pas

Paramètres clés

Pièges

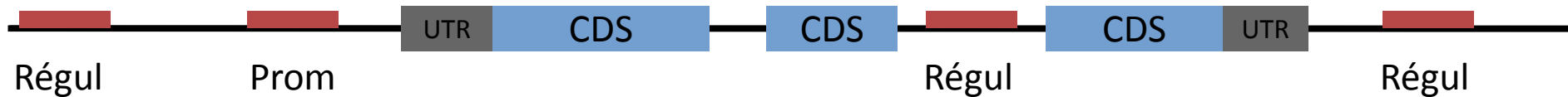
Formats et fichiers

Détection des SNPs (samtools/VarScan) pour analyse de leurs effets

Responsable et Intervenant principal : Sébastien Carrere
Expert : Ludovic Legrand

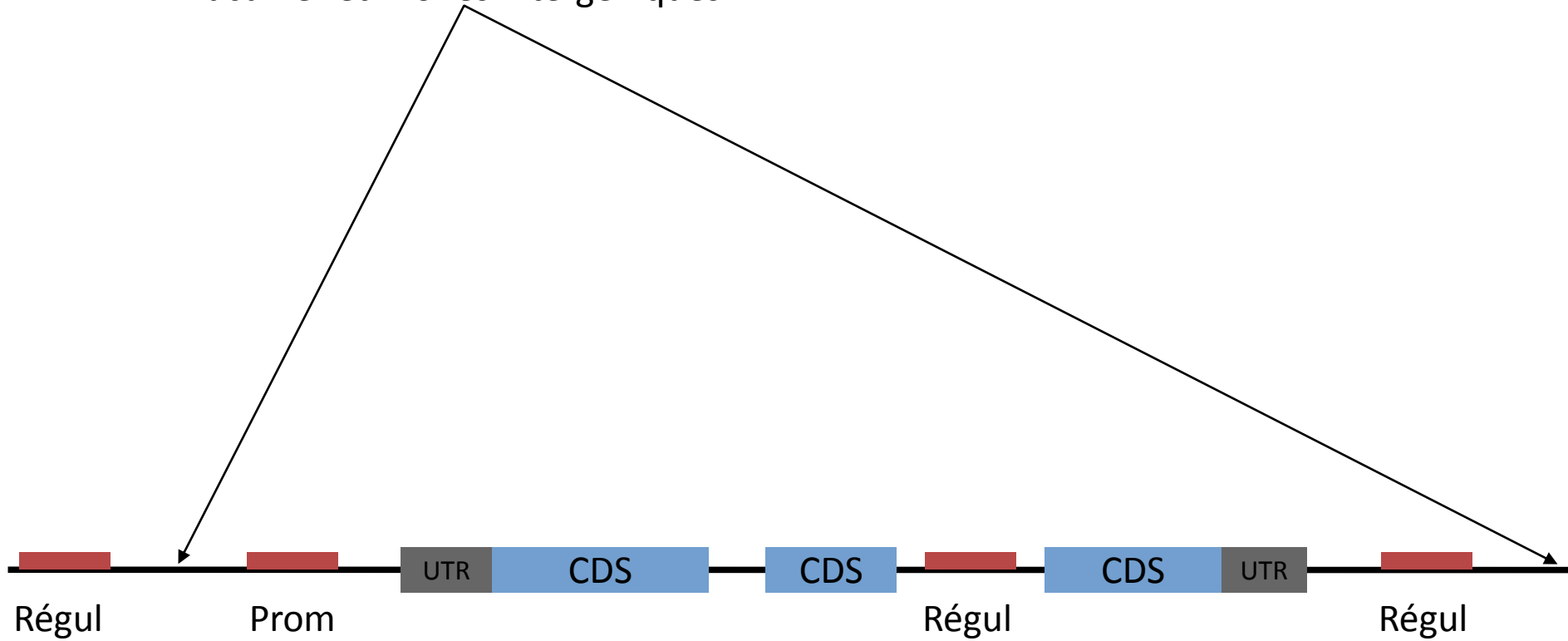
Impact des SNPs

- Impact de chaque SNP variable selon sa position et son environnement :
 - Aucun effet
 - Modification de la protéine produite
 - Effet sur la régulation de l'expression (« eSNP » = expression SNP)
 - Sites de fixation de facteurs de transcription
 - Sites d'épissage
 - Dégradation des mRNA



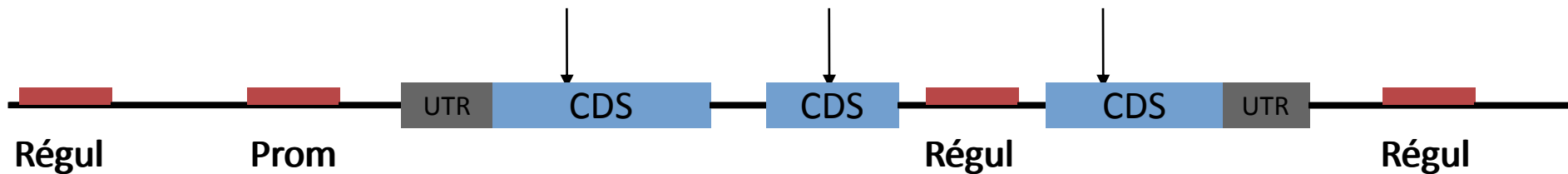
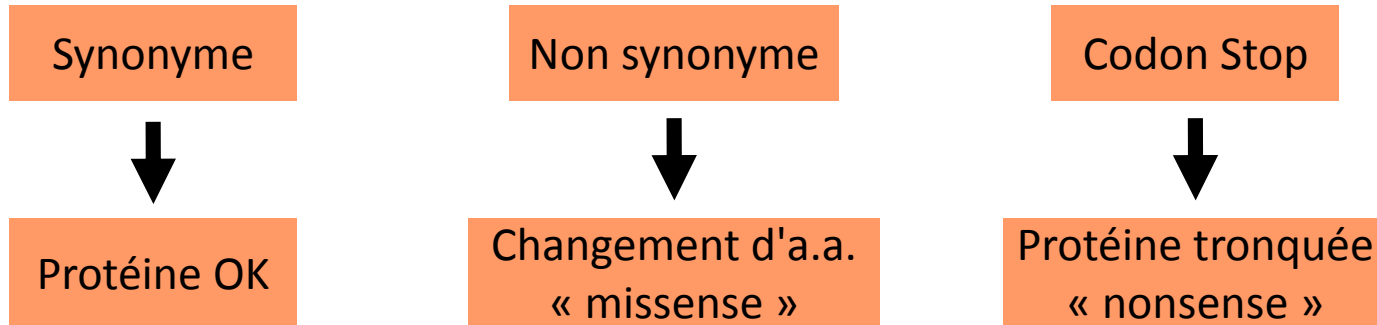
Impact des SNPs

- Impact de chaque SNP variable selon sa position et son environnement :
 - Aucun effet : zones intergéniques



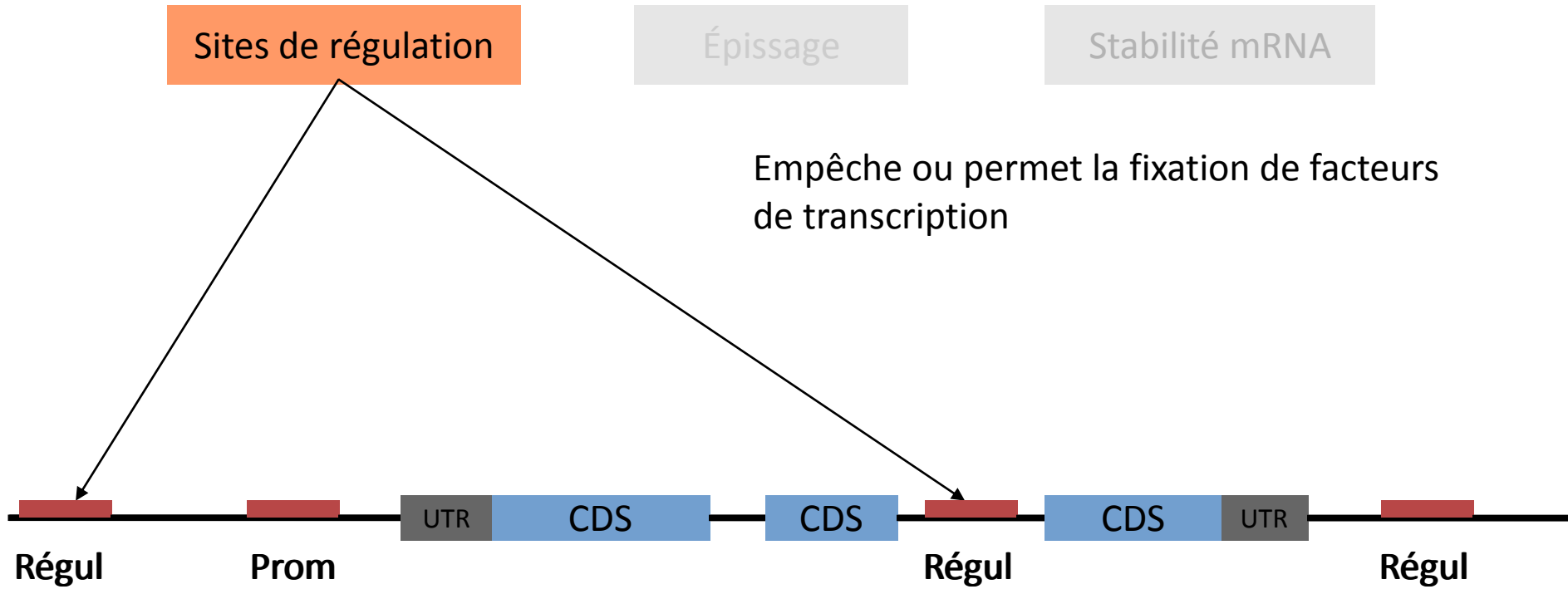
Impact des SNPs

- Impact de chaque SNP variable selon sa position et son environnement :
 - Modification de la protéine produite (nonsense, missense)



Impact des SNPs

- Impact de chaque SNP variable selon sa position et son environnement :
 - Effet sur la régulation de l'expression (« eSNP » = expression SNP)



Impact des SNPs

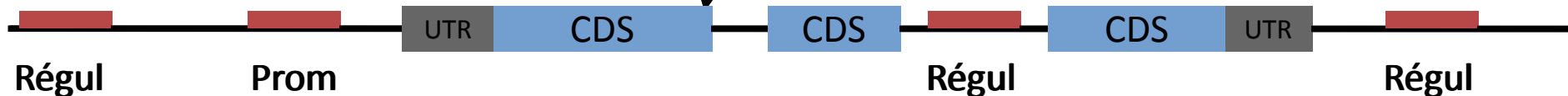
- Impact de chaque SNP variable selon sa position et son environnement :
 - Effet sur la régulation de l'expression (« eSNP » = expression SNP)

Sites de régulation

Épissage

Stabilité mRNA

Épissage = sites donneurs et accepteurs
Supprime ou crée ces sites
+ Effets sur sites de régulation



=> Impact sur les isoformes exprimées

Impact des SNPs

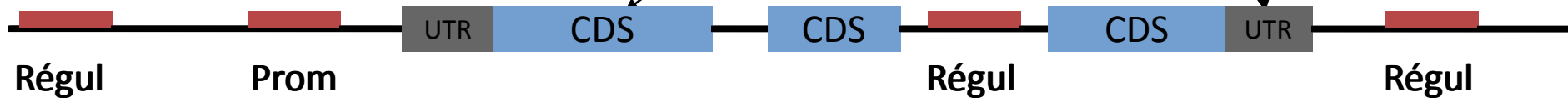
- Impact de chaque SNP variable selon sa position et son environnement :
 - Effet sur la régulation de l'expression (« eSNP » = expression SNP)

Sites de régulation

Épissage

Stabilité mRNA

Durée de vie des mRNA
Liaison à des RNA-binding proteins

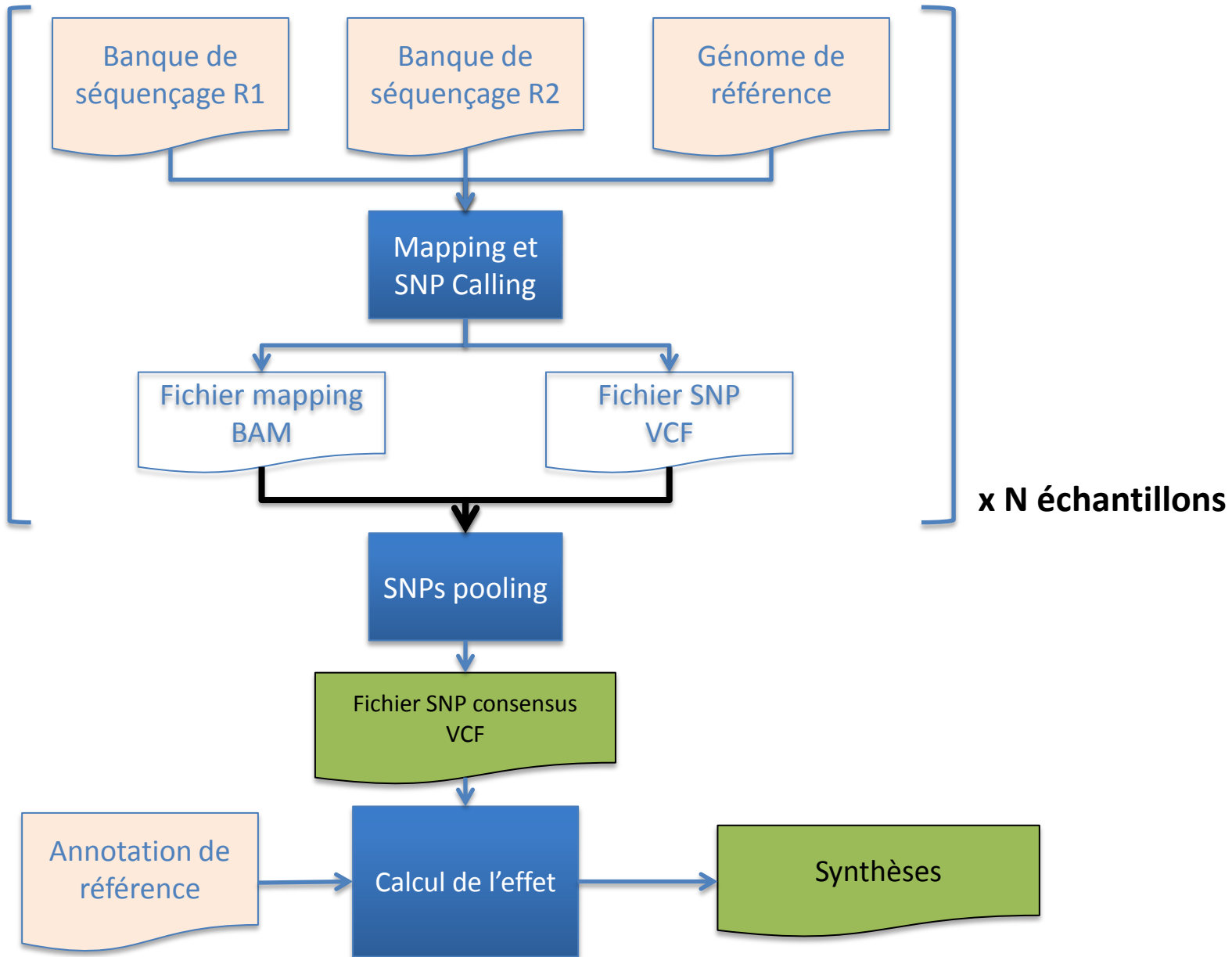


=> Impact sur la durée de vie des mRNA => expression

Le workflow

- Entrées:
 - Génome de référence connu et annotation structurale disponible
 - Re-séquençage du génome de plusieurs individus : Illumina pairé (ou single)
- Etapes:
 - Mapping
 - Calling de SNPs pour chaque échantillon
 - Génotypage de tous les individus pour tous les SNPs
 - Prédiction des effets induits par ces SNPs
- But :
 - Trouver des SNPs pouvant expliquer les différences de phénotypes observées
- Limites :
 - Importance de l'annotation de référence
 - Effets cachés des SNPs : SNP intergénique ou synonyme peut avoir un fort impact

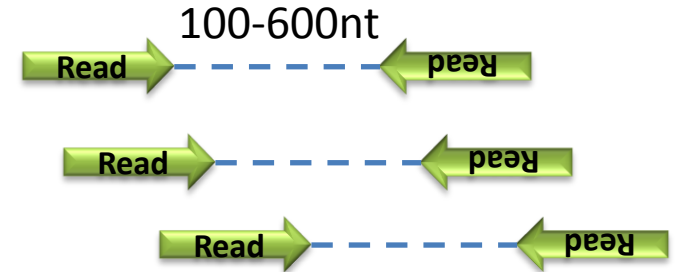
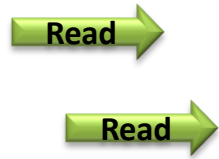
Fonctionnement du pipeline



Fonctionnement du pipeline: mapping et SNP calling

Mapping

Par échantillon



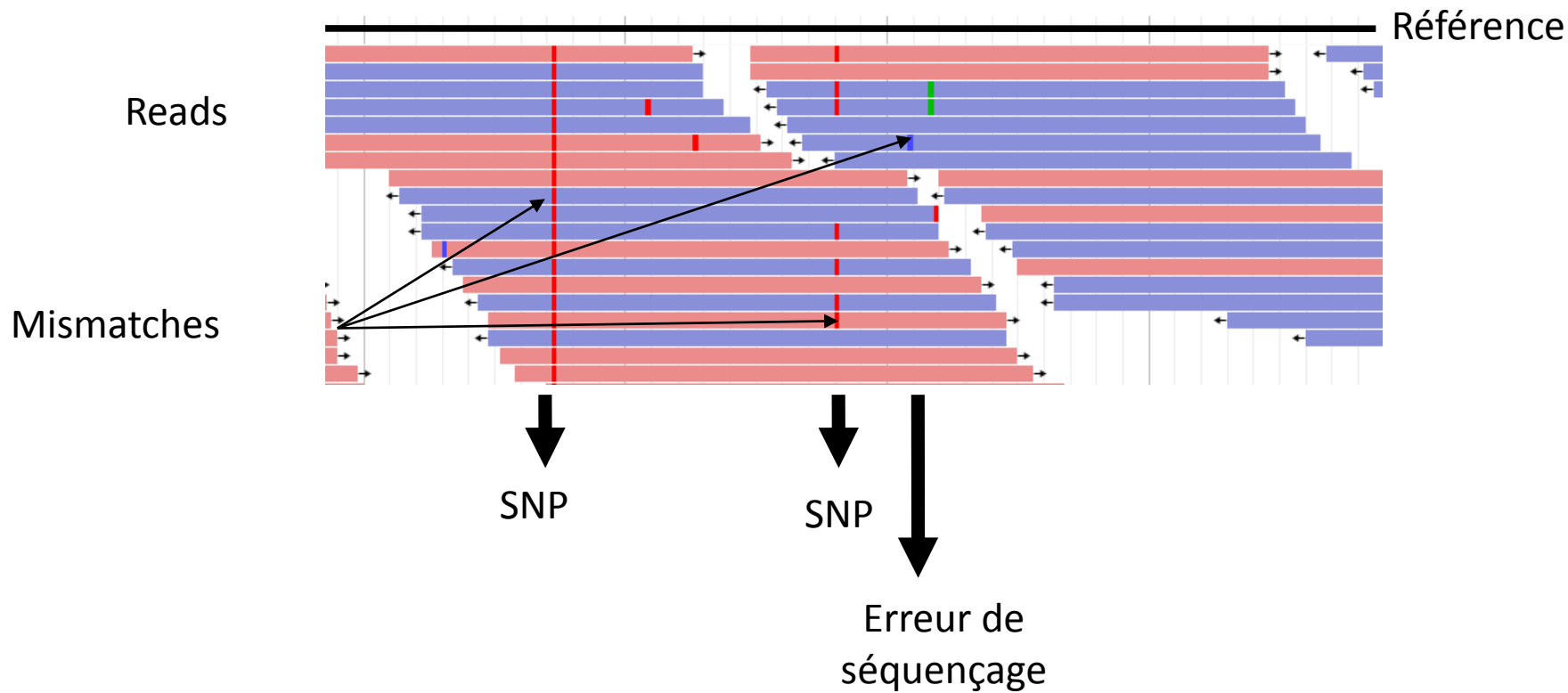
Librairie Single-end

Librairie Paired-end

Paramètres par défaut de glint: `-lmin=80nt -mmis=4`



SNP calling



Fonctionnement du pipeline: mapping et SNP calling

SNP calling

Par échantillon

Fichier mapping
BAM

samtools mpileup

VarScan mpileup2snp

Fichier SNP
VCF

option -B: désactive le calcul du BAQ (Base Alignment Quality)
=> Trop strict, perte de SNPs vrais positifs

options:

Minimum position coverage

Profondeur minimum en reads à une position donnée pour faire le calling
valeur par défaut du pipeline: 20

Minimum variant coverage

Nombre minimum de reads supportant un variant pour faire le calling
valeur par défaut du pipeline: 10

Minimum variant frequency

Fréquence allélique minimum du variant pour le conserver
valeur par défaut du pipeline: 0,2

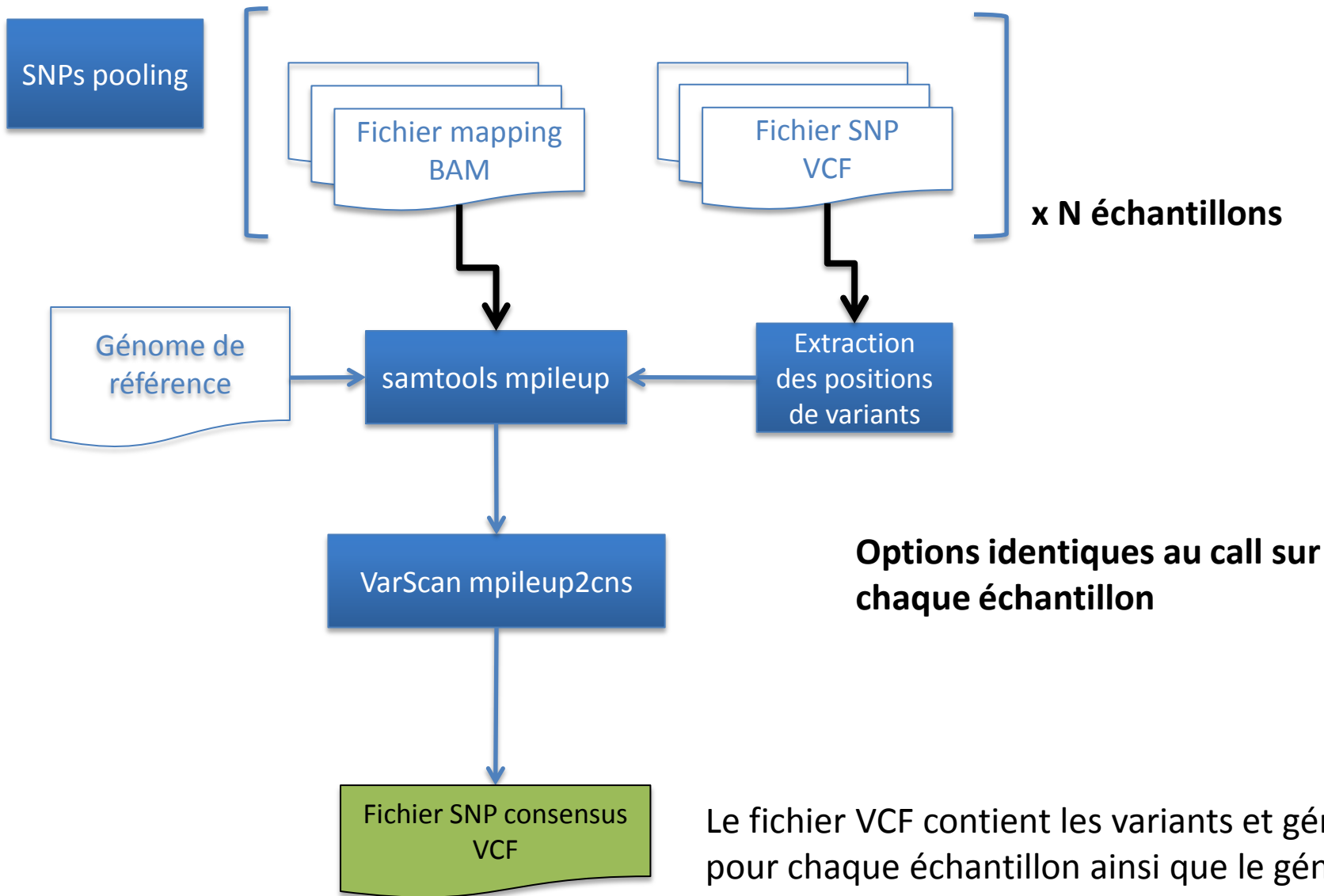
Minimum homozygous frequency

Fréquence allélique minimum pour attribuer un génotype homozygote
valeur par défaut du pipeline: 0,75

<http://samtools.sourceforge.net/mpileup.shtml>

<http://varscan.sourceforge.net/using-varscan.html>

Fonctionnement du pipeline: SNP pooling



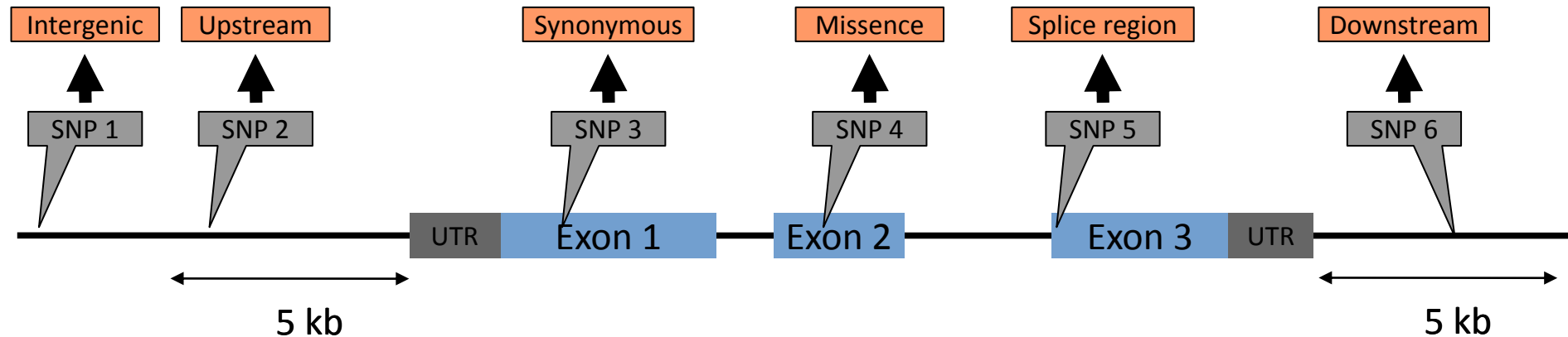
Options identiques au call sur chaque échantillon

Le fichier VCF contient les variants et génotypes pour chaque échantillon ainsi que le génotype consensus pour l'ensemble des échantillons.

SNPEff

- Prédiction des effets des SNPs avec SNPEff
 - Se base sur une annotation de référence
 - Résultat
 - VCF annoté
 - Rapport HTML
 - Tableau des effets par gène

SNPEff déduit un effet probable en fonction de la position du SNP par rapport à l'annotation



Galaxy : import des données

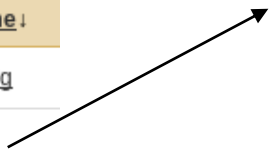
Data library name :

1001genomes.org

[BBRIC protocols](#)

[Phytozome](#)

[TAIR](#)



<input checked="" type="checkbox"/>	X.bbric	
<input checked="" type="checkbox"/>	Xbbric_A	
<input checked="" type="checkbox"/>	Xbbric_A.pe.1.250k.fastq.gz	Individu 1
<input checked="" type="checkbox"/>	Xbbric_A.pe.2.250k.fastq.gz	
<input checked="" type="checkbox"/>	Xbbric_B	
<input checked="" type="checkbox"/>	Xbbric_B.pe.1.250k.fastq.gz	Individu 2
<input checked="" type="checkbox"/>	Xbbric_B.pe.2.250k.fastq.gz	
<input checked="" type="checkbox"/>	Xbbric_C	
<input checked="" type="checkbox"/>	Xbbric_C.pe.1.250k.fastq.gz	Individu 3
<input checked="" type="checkbox"/>	Xbbric_C.pe.2.250k.fastq.gz	
<input checked="" type="checkbox"/>	Xbbric genomic sequence	Génome (fasta)
<input checked="" type="checkbox"/>	Xbbric structural annotation	Annotation (gff)
<input type="checkbox"/>	Xiphinema	

For selected datasets:

Galaxy : Formulaire

POLYMORPHISM

[Genetic variant annotation and effect prediction \(snpEff\)](#)

[SNP detection and effect prediction](#)

[SNP detection without reference](#)

[SNP detection for building an allelic frequency matrix](#)

[Variant detection and allele frequency GATK based](#)

SNP detection and effect prediction (version 0.1)

Paired-end library:




Maximum distance between paired reads (nt):

300

Paired-end libraries

Paired-end library 1

Read file 1:  

1: Xbbtric_A.pe.1.250k.fastq.gz

fastq,fastq.gz

Read file 2:  

2: Xbbtric_A.pe.2.250k.fastq.gz

fastq,fastq.gz

Sample name:

A

Name without space like: genotype1, genotype_2...



Remove Paired-end library 1

Paired-end library 2

Read file 1:  

3: Xbbtric_B.pe.1.250k.fastq.gz

fastq,fastq.gz

Read file 2:  

4: Xbbtric_B.pe.2.250k.fastq.gz

fastq,fastq.gz

Sample name:

B

Name without space like: genotype1, genotype_2...

Remove Paired-end library 2

Paired-end library 3

Read file 1:  

5: Xbbtric_C.pe.1.250k.fastq.gz

fastq,fastq.gz

Read file 2:  

6: Xbbtric_C.pe.2.250k.fastq.gz

fastq,fastq.gz

Sample name:

C

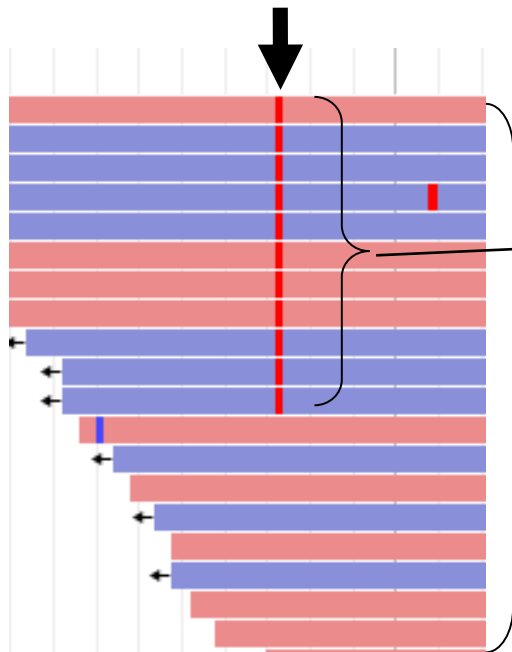
Name without space like: genotype1, genotype_2...

Remove Paired-end library 3

Add new Paired-end library

3 échantillons

Galaxy : Formulaire



Reference genome file (fasta):

7: Xbbtric genomic sequence

Genome annotation file (GFF3):

8: Xbbtric structural annotation

Minimum hit length:

80

Maximum number of mismatches:

4

Minimum position coverage:

20

Minimum variant coverage:

10

Minimum variant frequency:

0.2

Minimum homozygous frequency:

0.75

Mapping (échantillons = 76~100bp)

SNP calling

Galaxy : Formulaire

Effect results filter:

Select All

Unselect All

- Hide DOWNSTREAM/UPSTREAM effect
- Hide INTERGENIC effect
- Hide INTRON effect
- Hide 3/5 PRIME UTR effect

} SNP non annotés dans le VCF
(a priori les moins intéressants)

Upper case nucleotides:



use masked nucleotides (lower case) for mapping

Execute

Galaxy : Résultats

snp_effects
11 shown
113.5 MB

- 11: SNP effect by gene**
- 10: SNP effect report**
- 9: SNP matrix file**

Number of effects by impact

Type (alphabetical order)	Count	Percent
HIGH	5	0.012%
LOW	2,818	6.755%
MODERATE	747	1.791%
MODIFIER	38,148	91.443%

Number of effects by functional class

Type (alphabetical order)	Count	Percent
MISSENSE	747	20.924%
NONSENSE	5	0.14%
SILENT	2,818	78.936%

Missense / Silent ratio: 0.2651

Number of effects by type and region

Génome compact

Type	Count	Percent	Region	Count	Percent
downstream_gene_variant	18,220	43.674%	DOWNSTREAM	18,220	43.674%
intergenic_region	479	1.148%	INTERGENIC	479	1.148%
missense_variant	747	1.791%	EXON	3,567	8.55%
splice_region_variant+stop_retained_variant	3	0.007%	NONE	8	0.019%
stop_gained	5	0.012%	SPLICE_SITE_REGION	3	0.007%
synonymous_variant	2,815	6.748%	UPSTREAM	19,441	46.601%
transcript	8	0.019%			
upstream_gene_variant	19,441	46.601%			

Galaxy : Résultats

snp_effects
11 shown
113.5 MB

11: SNP effect by gene

10: SNP effect report

9: SNP matrix file

The following table is formatted as tab separated values.

#GeneName	GeneId	variants_effect_downstream_gene_variant	variants_effect_missense_variant
GNL	gene:Xbbirc.3991	6	1
V	gene:Xbbirc.3969	24	0
Xbbirc-rRNA-16s_rRNA-4561300-4562834	gene:Xbbirc-rRNA-16s_rRNA-4561300-4562834	66	0
Xbbirc-rRNA-23s_rRNA-4557919-4560799	gene:Xbbirc-rRNA-23s_rRNA-4557919-4560799	38	0
Xbbirc-rRNA-5s_rRNA-4557675-4557789	gene:Xbbirc-rRNA-5s_rRNA-4557675-4557789	4	0
Xbbirc.3908	gene:Xbbirc.3908	88	3
Xbbirc.3909	gene:Xbbirc.3909	90	8
Xbbirc.3911	gene:Xbbirc.3911	93	1
Xbbirc.3912	gene:Xbbirc.3912	80	9
Xbbirc.3913	gene:Xbbirc.3913	50	6
Xbbirc.3915	gene:Xbbirc.3915	35	1
Xbbirc.3916	gene:Xbbirc.3916	35	0
Xbbirc.3919	gene:Xbbirc.3919	1	1
Xbbirc.3921	gene:Xbbirc.3921	16	1
Xbbirc.3922	gene:Xbbirc.3922	26	0
Xbbirc.3923	gene:Xbbirc.3923	3	0
Xbbirc.3925	gene:Xbbirc.3925	11	0
Xbbirc.3926	gene:Xbbirc.3926	66	1

(...)

Nom/id des gènes

Nombres de SNP trouvés
dans chaque catégorie

Galaxy : Résultats

snp_effects
11 shown
113.5 MB

11: SNP effect by gene

10: SNP effect report

9: SNP matrix file

Fichier VCF : standard, VCFTools

#CHROM	POS	ID	REF	ALT	QUAL	FILTER
Xbbric	29309	.	G	A	.	PASS

Chromosome

Position

SNP ID (. = indéfini)

Allèle de référence

Allèle alternatif

Qualité (. = indéfinie)

Qualité suffisante ou non (indéfini ici)

Galaxy : Résultats

snp_effects
11 shown
113.5 MB

11: SNP effect by gene

10: SNP effect report

9: SNP matrix file

Fichier VCF : standard, VCFTools

```
INFO
ADP=49;WT=1;HET=0;HOM=2;NC=0;ANN=A|stop_gained|HIGH|Xbbri.3932|gene:Xbbri.3932|transcript|mRNA:Xbbri.3932|
|Coding|1/1|c.23G>A|p.Trp8*|23/339|23/339|8/112||WARNING_TRANSCRIPT_NO_START_CODON
```

Liste de paires de « clé=valeur » séparées par des « ; »

ADP : Profondeur moyenne par échantillon pour les bases ayant un score Phred = 20

WT : nombre d'échantillons wild-type

HET : nombre d'échantillons hétérozygotes

HOM : nombre d'échantillons homozygotes (=ont le snp x 2)

NC : échantillons « not called »

ANN : annotation ajoutée par SNPEff (ici apparition d'un codon stop)

Galaxy : Résultats

snp_effects
11 shown
113.5 MB

11: SNP effect by gene

10: SNP effect report

9: SNP matrix file

Fichier VCF : standard, VCFTools

FORMAT
GT:GQ:SDP:DP:RD:AD:FREQ:PVAL:RBQ:ABQ:RDF:RDR:ADF:ADR

A	B	C
1/1:255:49:49:0:49:100%:3.925E-29:0:38:0:0:31:18	1/1:255:47:47:0:47:100%:6.1512E-28:0:36:0:0:33:14	0/0:100:53:53:53:0:0%:1E0:69:0:23:30:0:0

Infos par échantillons selon codage défini dans la colonne FORMAT :

GT : genotype (1 = allele alt)

GQ : genotype quality

SDP : Raw Read Depth

DP : Quality Read Depth of bases with Phred score ≥ 20

RD : Depth of reference-supporting bases

AD : Depth of variant-supporting bases

FREQ : Variant allele frequency

PVAL : P-value from Fisher's Exact Test (not computed here : default value)

RBQ : Average quality of reference-supporting bases

ABQ : Average quality of variant-supporting bases

RDF/RDR : Depth of reference-supporting bases on forward/reverse strand

ADF/ADR : Depth of variant-supporting bases on forward/reverse strand

Je veux identifier des SNPs et évaluer leurs effets sur l'annotation d'un génome

Je vais sur

L'outil permet

L'outil ne permet pas

Paramètres clés

Pièges

Formats et fichiers



Détection des SNPs et calcul des fréquences alléliques pour les positions bi-alléliques

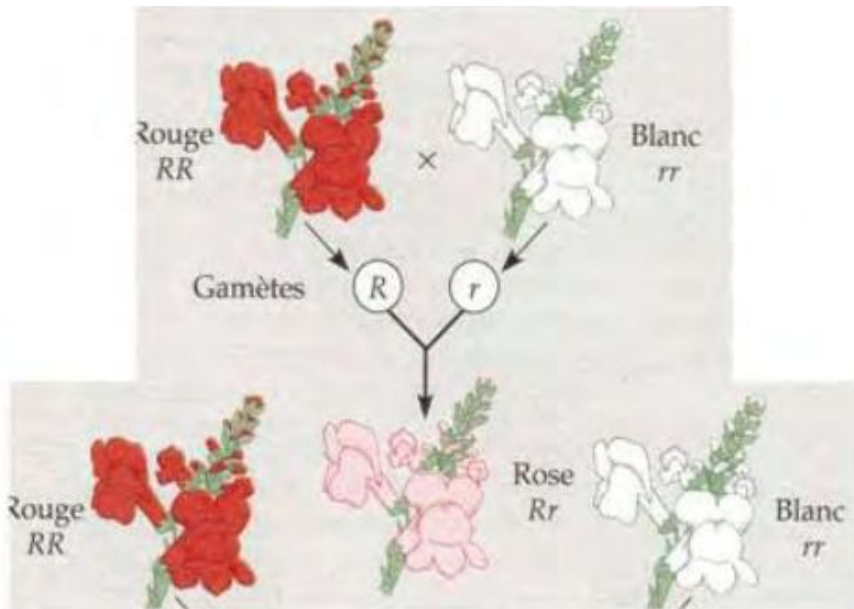
Responsable et intervenant principal: Sébastien Carrere
Expert: Ludovic Legrand

Plan

- Une matrice de comptage/fréquences alléliques pour quoi faire?
 - Contexte et objet d'étude nécessitant ce type de matrice
- Le pipeline pour la détection des SNPs et le calcul des fréquences alléliques (positions bi-alléliques)
 - Objectif et périmètre du pipeline
 - Fonctionnement général du pipeline
 - Mapping et calling de variants (SNPs)
 - Génération de la matrice de SNPs poolés
 - Filtre des sites bi-alléliques et génération de la matrice de comptages/fréquences alléliques
 - Fichiers de sortie:
 - VCF
 - Fichier de comptage allélique
 - Interface galaxy du pipeline
- Perspectives: exploitation de la matrice de fréquences alléliques

Qu'est-ce qu'une fréquence allélique?

- Un **allèle** est un variant, une des versions d'un même gène/locus (R ou r).
- Pour un organisme diploïde, une paire d'allèles, hérités de chacun des parents, représente le **génotype** d'un gène spécifique (RR, Rr ou rr).



Estimation des fréquences alléliques dans le polymorphisme floral chez les Gueules-de-loup (*Antirrhinum majus*)

Echantillon de **400** plantes d'une population **diploïde**:

165 Rouges; 190 Roses; 45 Blanches

$$P = \text{fréquence de l'allèle R dans l'échantillon} \\ = (2 \times 165 + 190) / (400 \times 2) = 0.65$$

$$Q = \text{fréquence de l'allèle r dans l'échantillon} \\ = (190 + 2 \times 45) / (400 \times 2) = 0.35$$

$$\text{Vérification: } p + q = 0.65 + 0.35 = 1.00$$

Qu'est-ce qu'une fréquence allélique?

Fréquence Allélique

C'est une mesure de la variabilité génétique au sein d'une population. Elle quantifie à un locus donné la probabilité d'observer un allèle (variant) donné par rapport à l'ensemble des allèles connus à ce locus dans une population.

Si on considère un gène A possédant de multiples allèles A_1, \dots, A_n en un locus d'une population diploïde,

$$p(i) = (2n(A_i A_i) + \sum_{i \neq j} n(A_i A_j)) / 2N$$

Où $p(i)$ représente la fréquence de l'allèle A_i en ce locus.

Fréquence Génotypique

Une mesure de la diversité génétique au sein d'une population. Elle quantifie à un locus donné la probabilité d'observer un génotype spécifique par rapport à l'ensemble des génotypes recensés à ce même locus dans la population.

$$F(AA) = n(AA)/N$$

Où $F(AA)$ représente la fréquence des individus ayant le génotype AA, $n(AA)$ est le nombre d'individus ayant ce génotype, et N le nombre total d'individus.

Les fréquences alléliques pour quoi faire?

- **Génétique/Génétique des populations**

Traite de l'impact des différentes forces évolutives (sélection, mutations, migration) sur la répartition de la diversité génétique entre les populations et dans les populations.

- **Les fréquences alléliques permettent de pouvoir:**

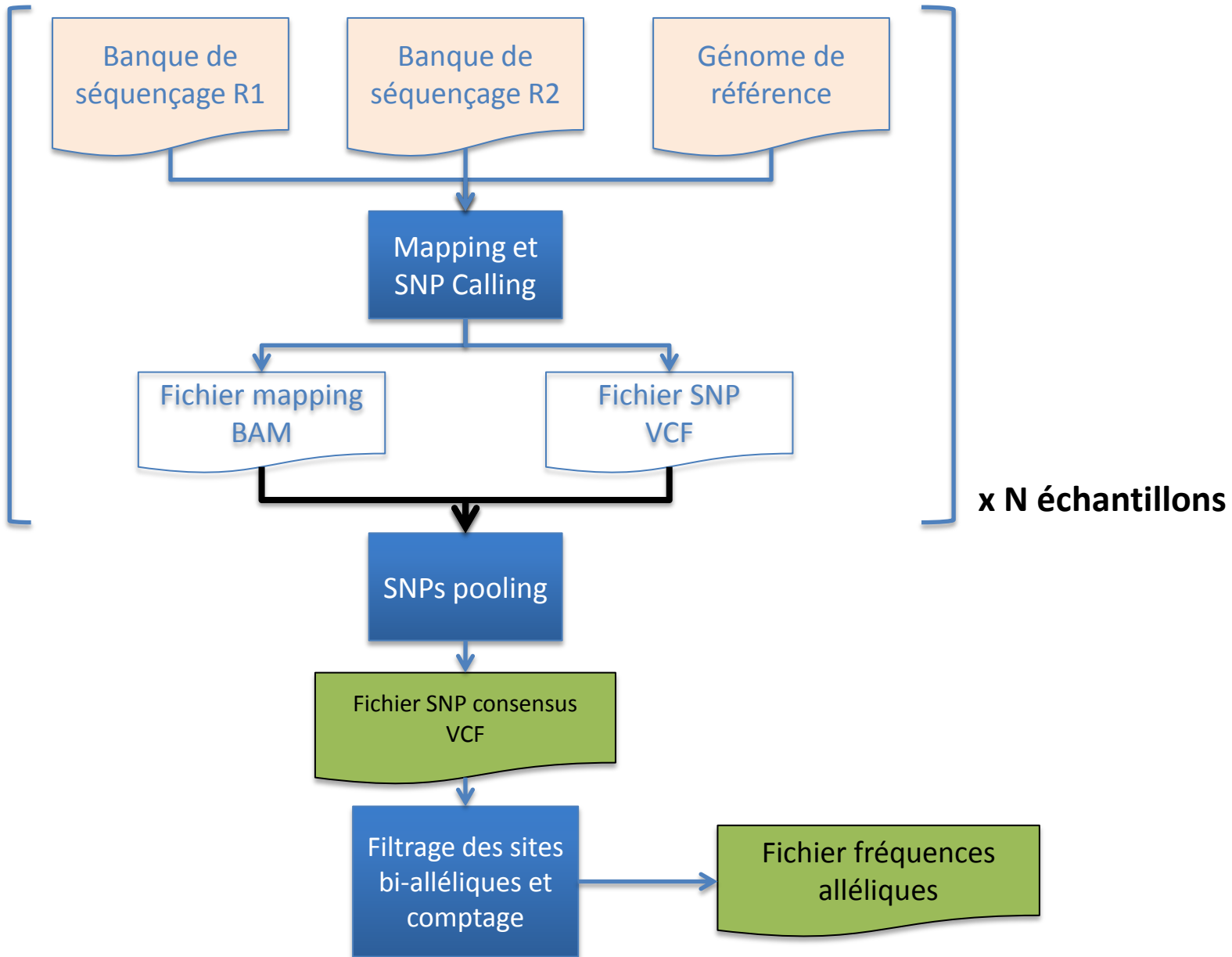
- Quantifier et typer des **variations génétiques** présentes dans les populations (MAF, Minor Allele Frequency): génique \Leftrightarrow genotype
- Chercher des **associations** entre variants et des traits phénotypiques (GWAS): genotype \Leftrightarrow phenotype

- Les **changements dans la fréquence des allèles** dans le temps peuvent indiquer la présence d'une **dérive génétique** ou bien que de **nouvelles mutations** ont été introduites dans la population.

Objectif et périmètre du pipeline

- **Objectif:** Détecter les variants de type SNP afin de calculer sur un ensemble d'échantillons les fréquences alléliques associées aux SNPs uniquement aux positions bi-alléliques.
- **Périmètre:** principalement pour réaliser des études d'association (GWAS) en génétique/génétique des populations et/ou caractériser des allèles peu fréquents dans une population (Pool-seq).
- **Limitations:** SNPs bi-alléliques uniquement, exclusion des sites multi-alléliques et des InDels

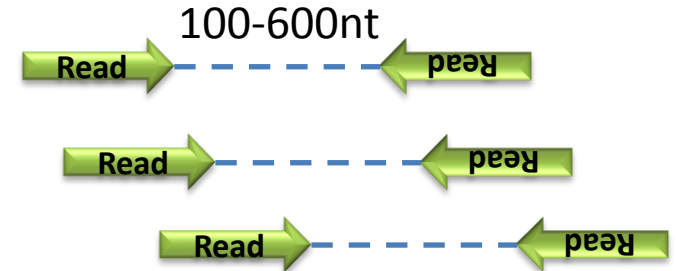
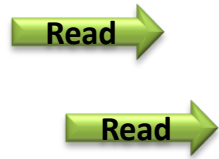
Fonctionnement du pipeline



Fonctionnement du pipeline: mapping et SNP calling

Mapping

Par échantillon



Librairie Single-end

Librairie Paired-end

Paramètres par défaut de glint: `-lmin=50nt -mmis=4`



Fonctionnement du pipeline: mapping et SNP calling

SNP calling

Par échantillon

Fichier mapping
BAM

samtools mpileup

VarScan mpileup2snp

Fichier SNP
VCF

option -B: désactive le calcul du BAQ (Base Alignment Quality)
=> Trop strict, perte de SNPs vrais positifs

options:

Minimum position coverage

Profondeur minimum en reads à une position donnée pour faire le calling
valeur par défaut du pipeline: 10

Minimum variant coverage

Nombre minimum de reads supportant un variant pour faire le calling
valeur par défaut du pipeline: 2

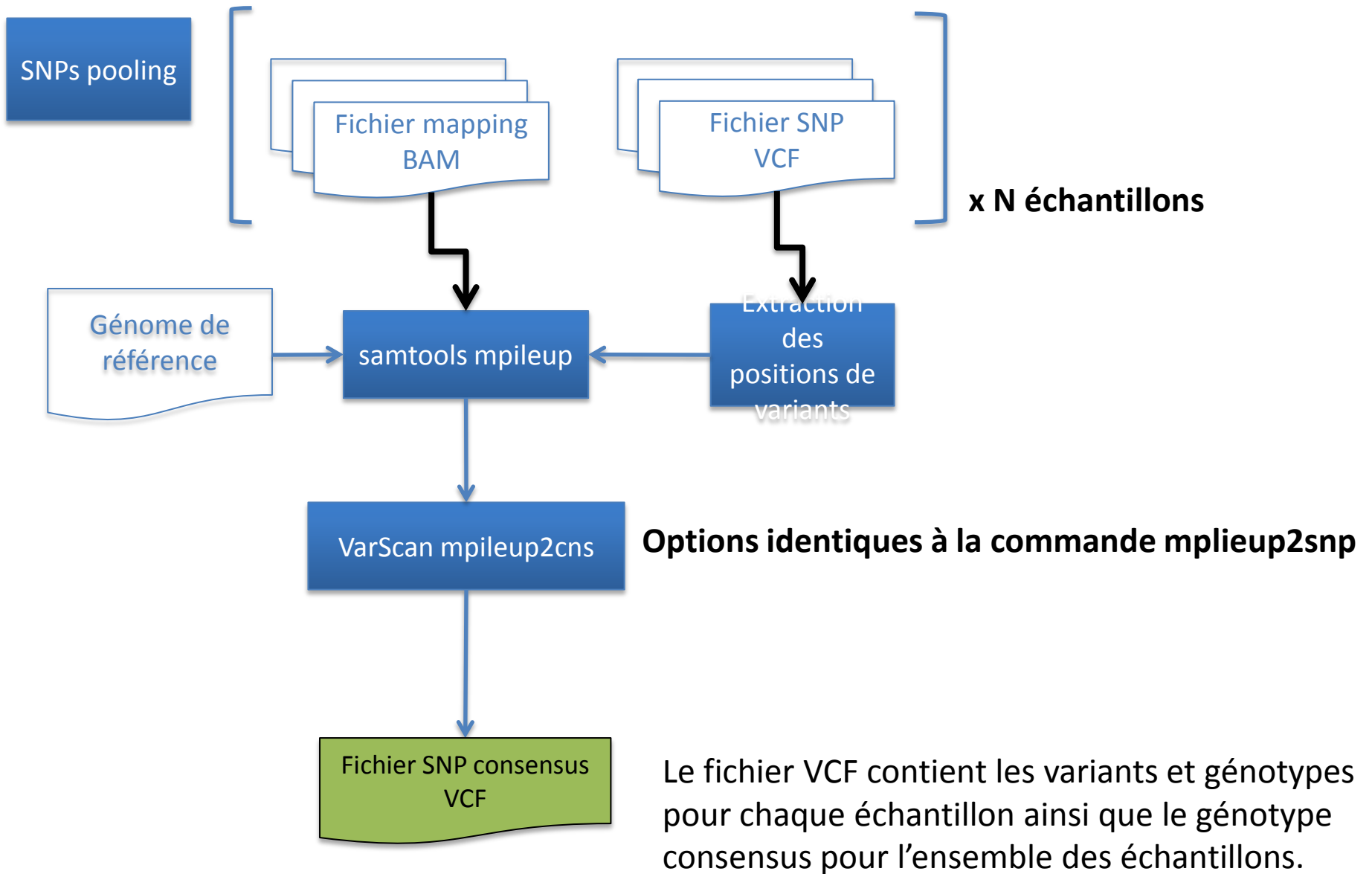
Minimum variant frequency

Fréquence allélique minimum du variant pour le conserver
valeur par défaut du pipeline: 0,01

<http://samtools.sourceforge.net/mpileup.shtml>

<http://varscan.sourceforge.net/using-varscan.html>

Fonctionnement du pipeline: SNP pooling



Fonctionnement du pipeline: Filtrage et comptage

Filtrage des sites bi-alléliques et comptage

Fichier SNP consensus VCF

Site bi-allélique

Site multi-alléliques

Reference TAGGCCGCCCGCGCTATGCCGTGCA

GACTGCGATCTCGACATCG

Reads

C
C
C
C
T
T
T
T
T
T
C
T
C
C

C
G
G
C
T
T
T
T
T
T
C
T
G
C

comptage



Fichier fréquences alléliques

Fichier fréquences alléliques

Allelic count

Tabular file with Allelic count by sample.

CHROM	POS	REF	ALT	#ALLELES	Spl1 #REF	Spl1 #TOT	Spl2 #REF	Spl2 #TOT
ID1	475	G	A	1	NA	NA	0	13
ID1	781	A	G	1	0	10	15	0

NA: position not in sample (missing data)

CHROM: chromosome or scaffold ID

POS: variant position

REF: base on reference sequence

ALT: alternative base

#ALLELE: number of allele

Spl1 #REF: number of reads with same base than reference for sample Spl1

Spl1 #TOT: number of reads for sample Spl1

Spl2 #REF: number of reads with same base than reference for sample Spl2

Spl2 #TOT: number of reads for sample Spl2

Fichier SNP consensus VCF

Cf. la fiche de description du format VCF.

Interface Galaxy du pipeline

Data Input

POLYMORPHISM

[Genetic variant annotation and effect prediction \(snpEff\)](#)

[SNP detection and effect prediction](#)

[SNP detection without reference](#)

[SNP detection Identifying variant using GATK](#)

[SNP detection for building an allelic frequency matrix](#)

Banques Single-End

Single-end libraries

Single-end library 1

Read file:



fastq,fastq.gz

③

Sample name:

④

Name without space like: genotype1, genotype_2...



Add new Single-end library

Banques Paired-End

①



Paired-end library:



Maximum distance between paired reads (nt):

②

Paired-end libraries

Paired-end library 1

Read file 1:

6: Xbbric_C.pe.2.250k.fastq.gz

③

fastq,fastq.gz

Read file 2:

6: Xbbric_C.pe.2.250k.fastq.gz

③

fastq,fastq.gz

Sample name:

④

Name without space like: genotype1, genotype_2...



Add new Paired-end library

Génome de référence

Reference genome file (fasta):



⑤

- ① Choix banque paired-end (single-end par défaut)
- ② Taille de l'insert si paired-end
- ③ Sélection des fichiers de séquençage
- ④ Nom de l'échantillon
- ⑤ Génome de référence

Interface Galaxy du pipeline

Mapping and SNP
calling/SNP pooling

Glint

Minimum hit length:

Couverture minimum le long de la read

Maximum number of mismatches:

Nombre maximum de substitutions autorisées

Minimum position coverage:

Profondeur minimum en reads à une position donnée pour faire le calling

VarScan mpileup2snp/mpileup2cns

Minimum variant coverage:

Nombre minimum de reads supportant un variant pour faire le calling

Minimum variant frequency:

Fréquence allélique minimum du variant

Execute



Perspectives: post-traitements

- GWAS avec PLINK/gPLINK/Haploview
 - PLINK: toolset pour l'analyse de données genotype/phenotype
 - <http://pngu.mgh.harvard.edu/~purcell/plink/>
 - gPLINK: interface graphique pour la gestion de commandes PLINK
 - <http://pngu.mgh.harvard.edu/~purcell/plink/gplink.shtml>
 - Haploview: analyse d'haplotypes
 - <http://www.broadinstitute.org/scientific-community/science/programs/medical-and-population-genetics/haploview/haploview>
- GWAS avec Efficient Mixed-Model Association eXpedited (EMMAX)
 - Transformation/conversions de format de génotypes
 - <http://pngu.mgh.harvard.edu/~purcell/plink/>
 - GWAS
 - <http://genetics.cs.ucla.edu/emmax/>

Je veux identifier des SNPs et calculer des fréquences alléliques

Je vais sur

L'outil permet

L'outil ne permet pas

Paramètres clés

Pièges

Formats et fichiers



MESURE DE L'EXPRESSION A PARTIR DE DONNÉES RNASEQ

Responsable et intervenant principal: Erika Sallet
Expert: Ludovic Legrand



Objectif et périmètre

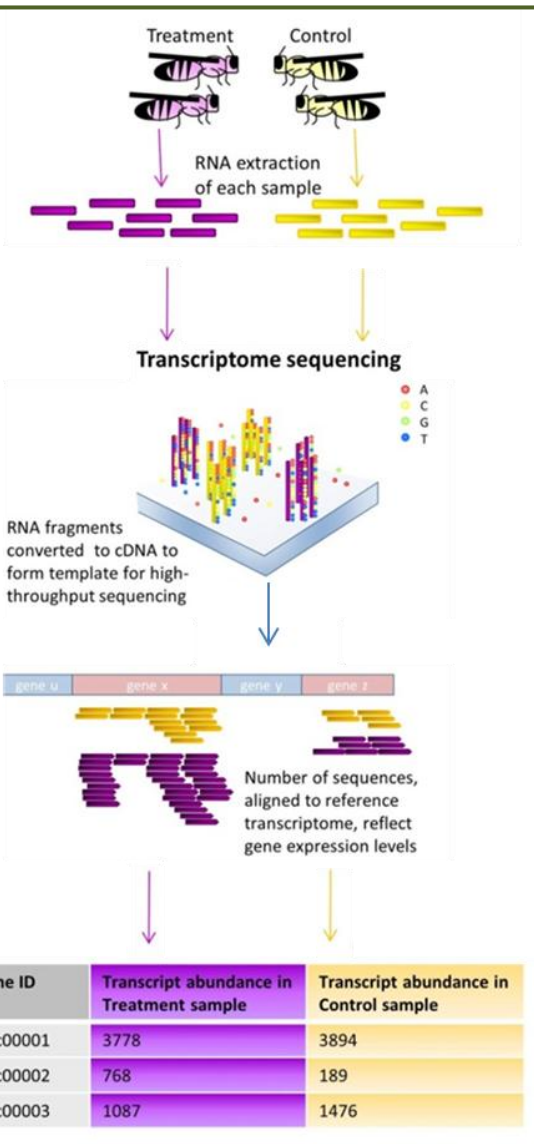
- Objectif : Mesure de l'expression

- Périmètre du pipeline :

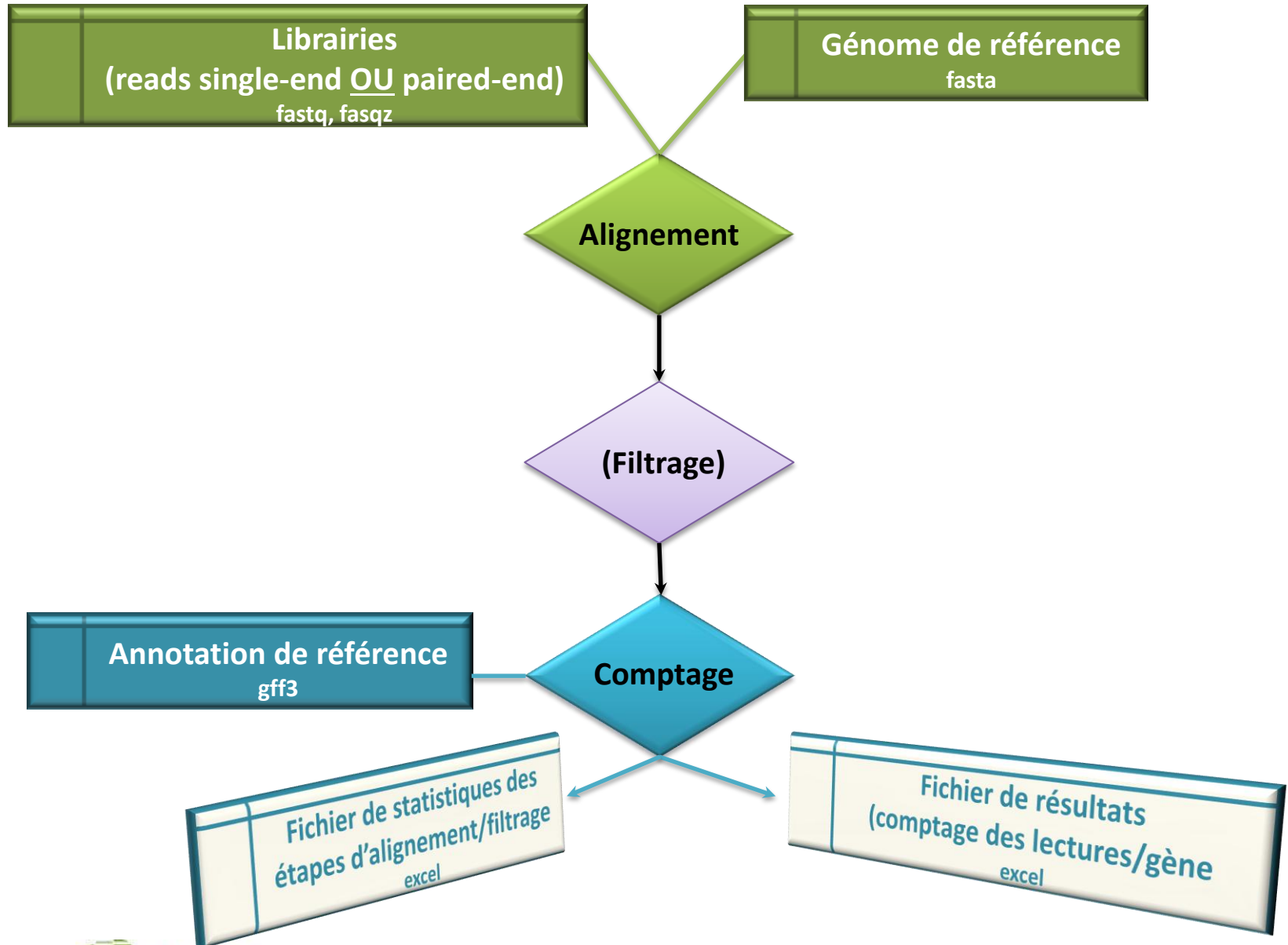
- Comptage des lectures sur des gènes définis dans le fichier d'annotation (pas de découverte de nouveaux transcrits)
- Pro / eucaryotes

- Fichiers requis :

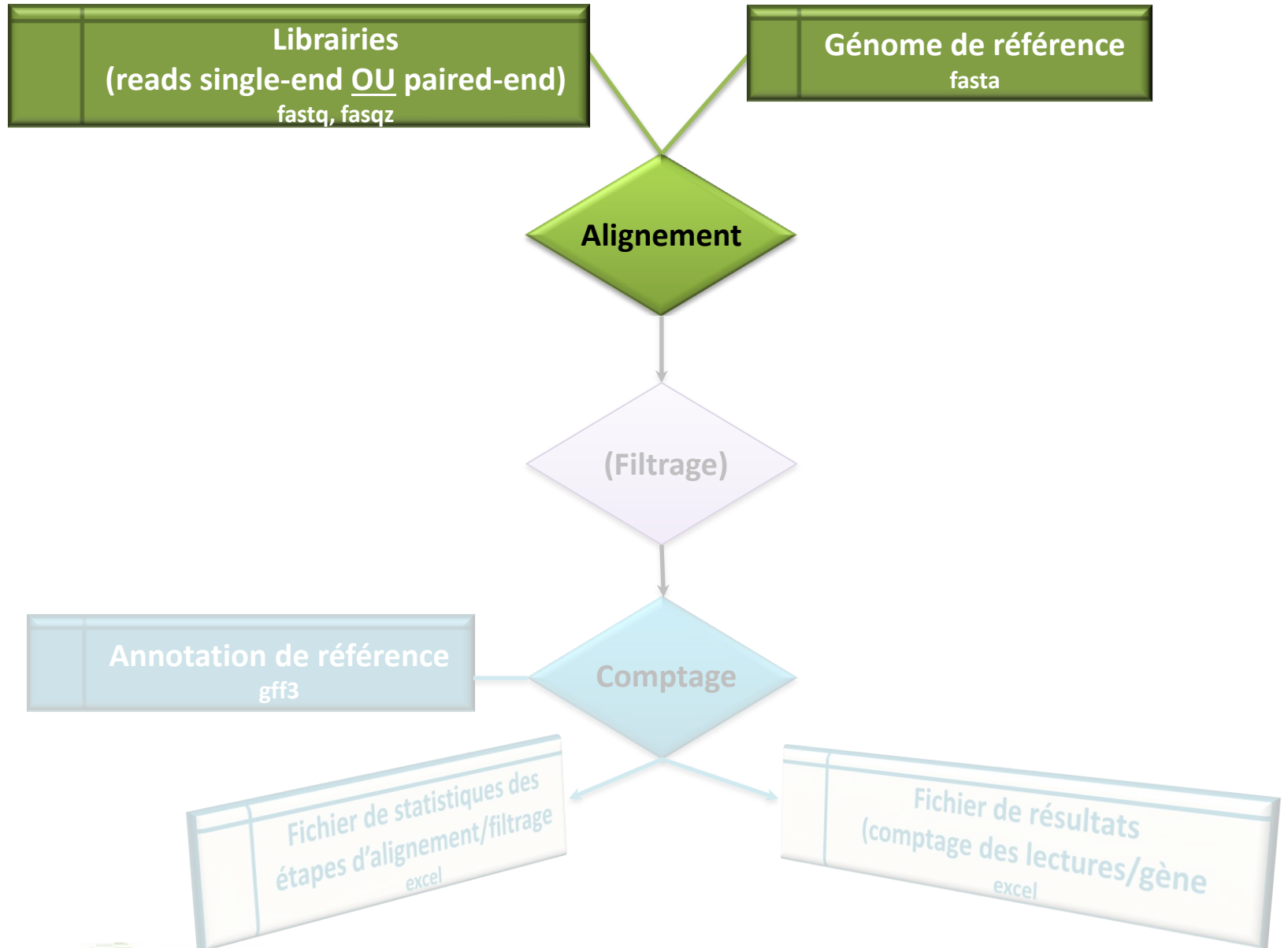
- Bibliothèques RNASeq paired-end ou single-end (fastq ou fastqz)
- Génome ou transcriptome (fasta)
- Annotation (GFF3)



Fonctionnement

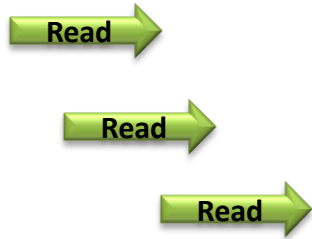


Alignement

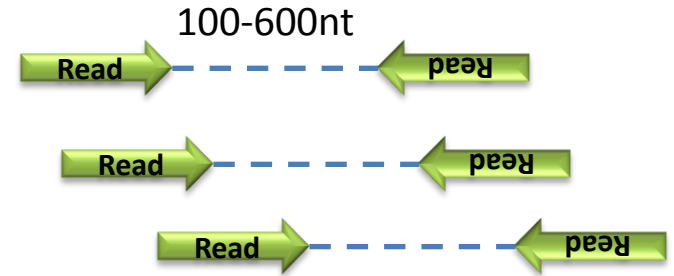


Alignement

Contig/Chromosome



Librairie Single-end



Librairie Paired-end

Paramètres influant sur l'alignement :

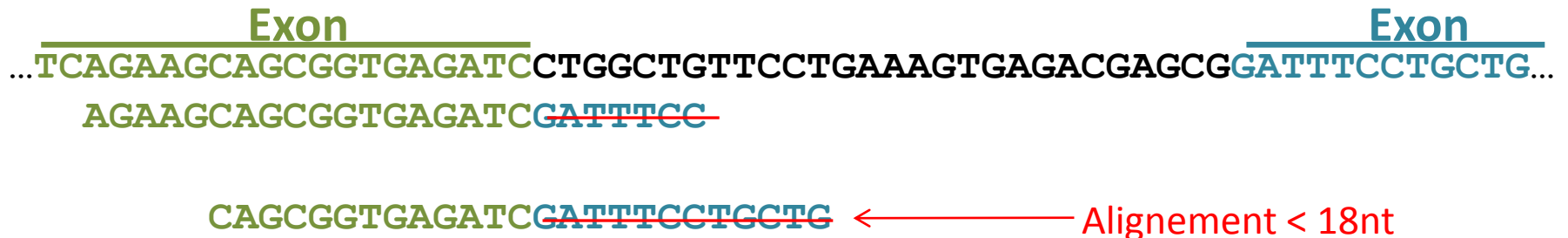
- longueur minimale du hit
- nombre maximal de mismatches
- distance maximale entre les paires
- petit ARN non codant (small RNA)

Alignement : influence des paramètres

Longueur minimale du hit

- compromis entre bruit et information
- perte des exons plus petits que ce paramètre
- valeur par défaut : 18 nt

Exemple read: 25nt lmin: 18nt



Alignement : influence des paramètres

Nombre maximal de mismatches

- erreurs de séquençage (~1%)
- variabilité avec la référence
- hétérozygotie
- valeur par défaut : 0

Exemple: lecture de 100nt paired-end

- 5 mismatches autorisés
 - prise en compte des variations alléliques et des erreurs
 - perte de spécificité d'alignement sur les lectures compensée par l'alignement de la paire

Alignement : influence des paramètres

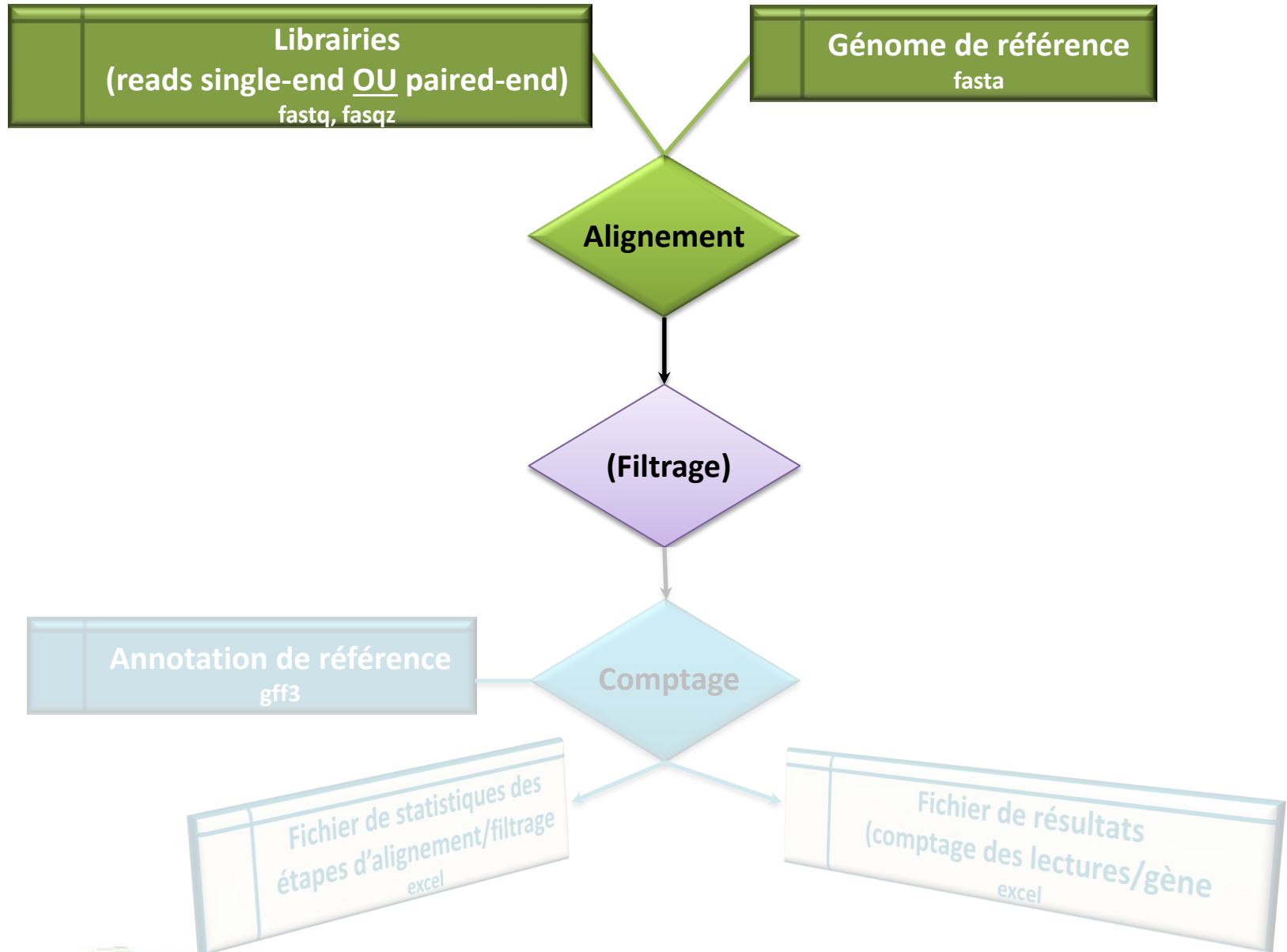
Distance maximale entre les paires

- cohérence de l'alignement de 2 reads d'une paire
- généralement entre 100 et 600 nt

Petit ARN non codant (small RNA)

- lecture entre 20 et 50 nt mais hits entre 18 et 30 nt
- 0 mismatch autorisé
 - pas d'informations complémentaires avec la paire
 - limite les alignements dus au hasard
- adaptation des seuils pour les petits alignements

Filtrage



Filtrage

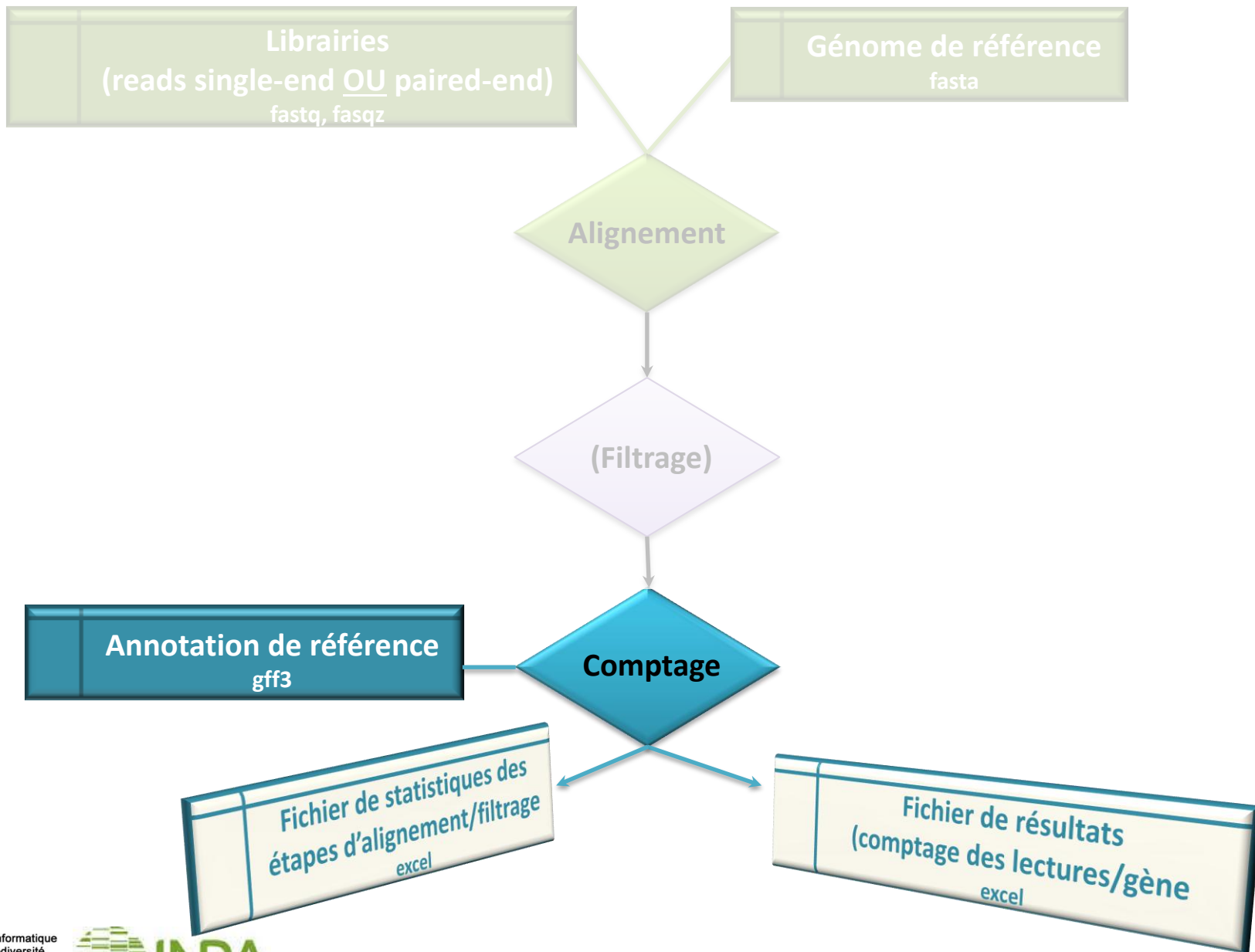
Filtrage des lectures non spécifiques (ambiguës)

- alignement identique de score maximum sur plusieurs positions
- duplication de gènes, transposons, gènes conservés

⇒ on préfère perdre de l'information en écartant les hits ambigus

⇒ garder les hits ambigus, c'est additionner le niveau d'expression de N objets biologiques très probablement **régulés différemment**

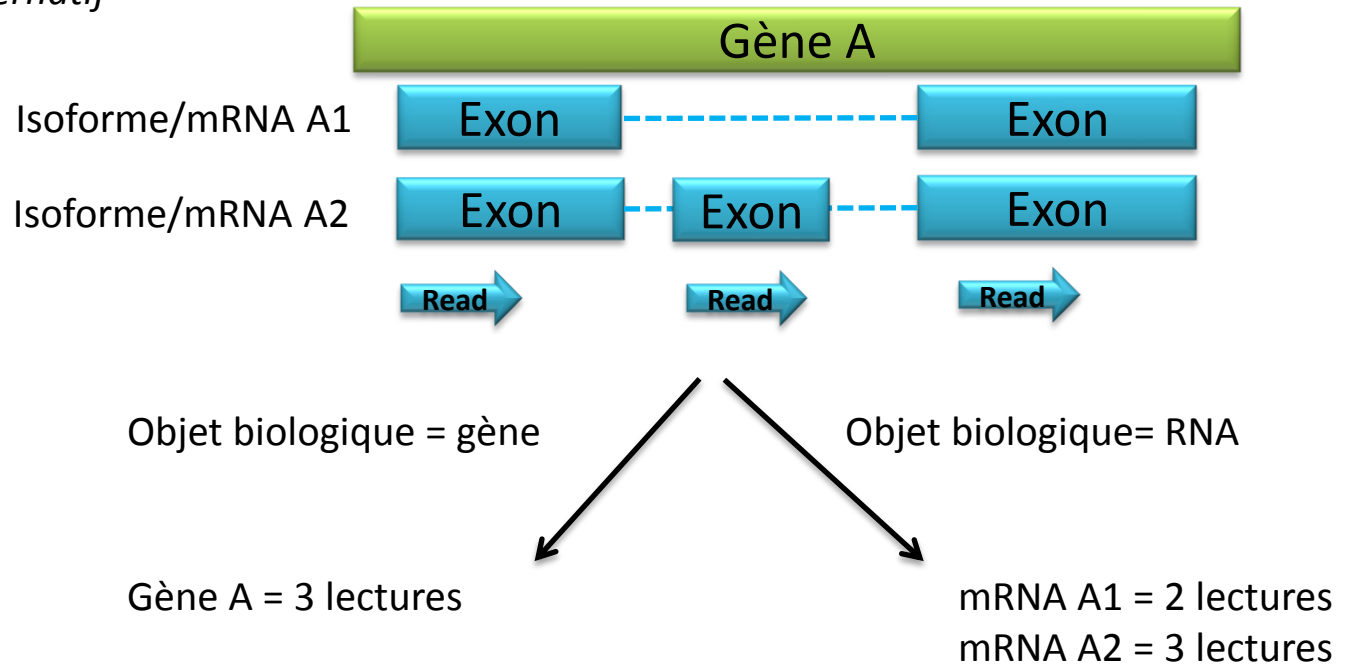
Comptage



Comptage : choix de l'objet biologique

Comptage du nombre de lectures par objet biologique

Exemple d'un gène A
avec épissage alternatif



Comptage : choix de l'objet biologique

Comptage du nombre de lectures par gène

- moins sensible à la qualité de l'annotation
- ne tient pas compte de l'épissage alternatif

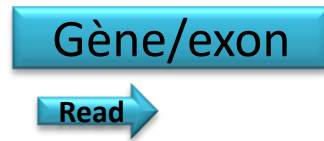
Comptage du nombre de lectures par RNA

- sensible à l'annotation des exons
- tient compte de l'épissage alternatif

Comptage : influence des paramètres

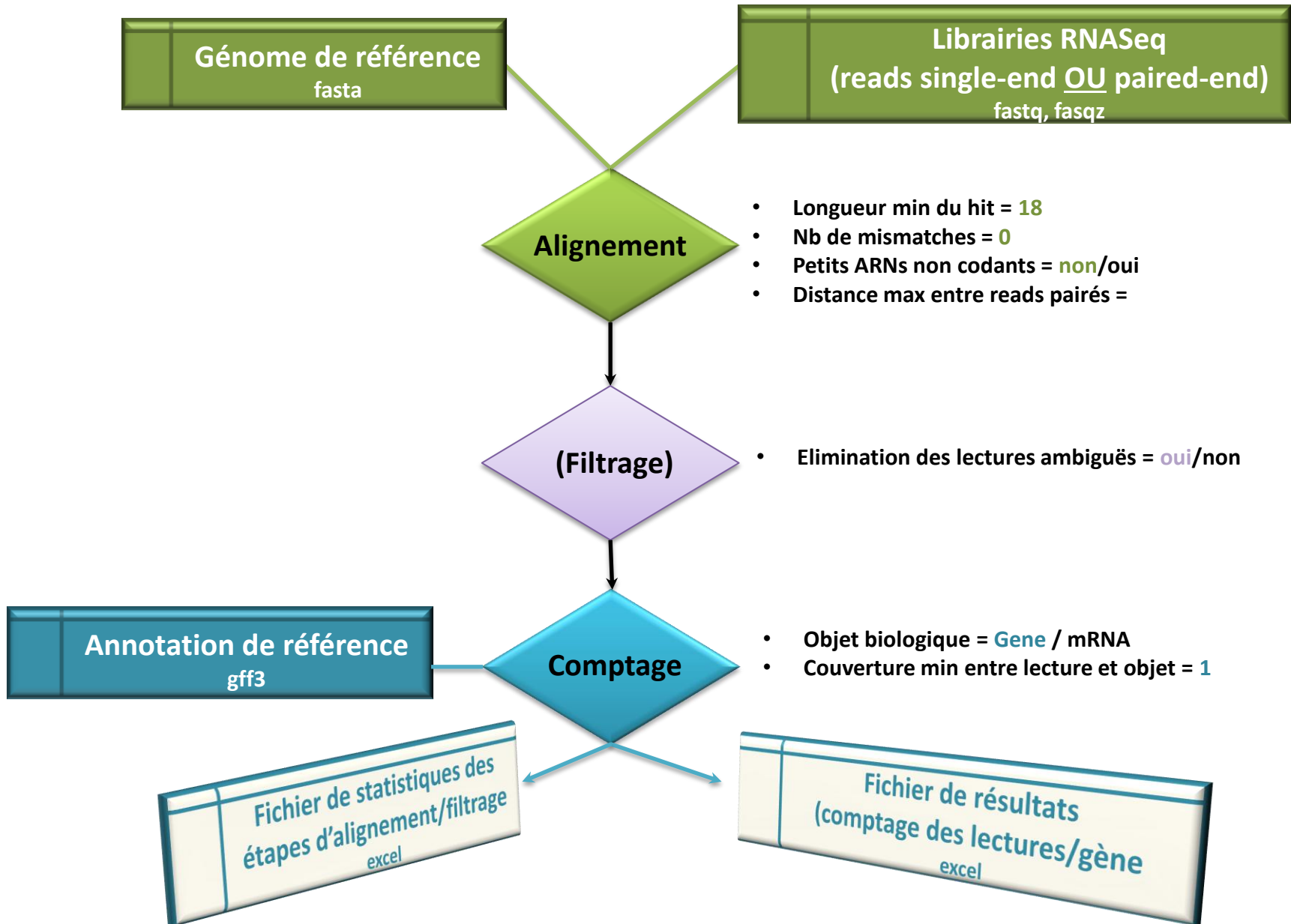
Couverture min de la lecture sur l'objet biologique

- Valeur par défaut = 1 (100%)
 - suppose une annotation de bonne qualité



- diminuer la couverture
 - autoriser les hits partiels sur les objets
 - annotation de faible ou moyenne qualité





Interface Galaxy

Expression measure (version 1)

Library type:

Paired-end

Paired-end libraries

Paired-end library 1

Read file 1:

11: S.bbric-RbmSmall-GGK36.ope.2.fastq.gz
fastq,fastq.gz

Read file 2:

11: S.bbric-RbmSmall-GGK36.ope.2.fastq.gz
fastq,fastq.gz

Add new Paired-end library

Maximal distance between paired reads (nt):

- Distance max entre reads pairés

Reference genome file (fasta):

19: Genomic sequence from Annotation on bacterial genome on data 8, data 9, and others

Génome de référence
fasta

Genome annotation file (GFF3):

12: Annotation on bacterial genome on data 8, data 9, and others

Annotation de référence
gff3

Expression reported for:

Gene

- Objet = Gene / RNA

Select biological objects (gff3 type) for which the expression is reported

Small inserts analysis:

- Petits ARNs non codants = non/oui

Change mapping parameters

Minimal hit length:

18

- Longueur min du hit = 18

Maximum number of mismatches:

0

- Nb de mismatches = 0

Advanced parameters:

Use no specific mapped reads:

- Elimination des lectures ambiguës = oui/non

Considering all mapped read/pairs (default: considering unambiguously mapped reads/pairs only)

Minimum overlap:

1.0

- Couverture min entre lecture et objet = 1 (100%)

Minimum overlap required as a fraction of the mapped read (intersecBed parameter -f)

Execute

Librairies
(reads single-end OU paired-end)
fastq, fastq

Génome de référence
fasta
Annotation de référence
gff3

Fichier de statistiques

Nombre d'alignements spécifiques sur génome

Nombre d'alignements sur génome

Nombre de lectures brutes

Nombre de lectures alignées sur le génome

Nombre d'alignements sur objet biologique

alignements sur objet biologique / alignements spécifiques sur génome

lib	specific_hits	mapping_hits	raw_reads/pairs_count	mapped_reads/pairs_count	feature_overlap_ping_hits	feature_overlapping_hits/specific_hits_percent
S.bbric_Rbm Long_GGK2 1.ope	1178	1181	9697	1179	1153	97.88
S.bbric_Rbm Small_GGK3 6.ope	24405	24411	64272	24407	20674	84.71

Contig/Chromosome



1181 hits totaux
1179 lectures qui s'alignent

Fichier de statistiques

Nombre d'alignements spécifiques sur génome

Nombre d'alignements sur génome

Nombre de lectures brutes

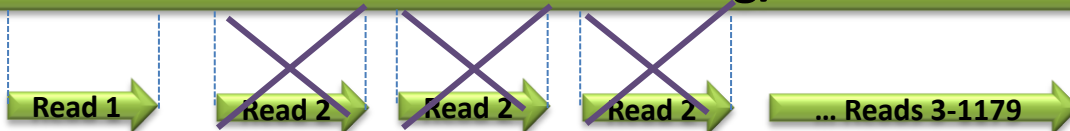
Nombre de lectures alignées sur le génome

Nombre d'alignements sur objet biologique

alignements sur objet biologique / alignements spécifiques sur génome

lib	specific_hits	mapping_hits	raw_reads/pairs_count	mapped_reads/pairs_count	feature_overlap_ping_hits	feature_overlapping_hits/specific_hits_percent
S.bbric_Rbm Long_GGK2 1.ope	1178	1181	9697	1179	1153	97.88
S.bbric_Rbm Small_GGK3 6.ope	24405	24411	64272	24407	20674	84.71

Contig/Chromosome



1178 hits spécifiques =
1178 lectures avec alignement spécifique

Fichier de résultats

Gene/RNA ID

Contig Id

Positionnement du gène/RNA sur le contig

Début

Fin

Brin

Taille du Gène/RNA

Type (Gène/RNA)

Comptage brut (Lectures/Objet)

Comptage normalisé (RPKM)

id	0_seqid	1_start	2_end	3_strand	4_length	5_type	6_Note	S.bbric_RbmLong_GGK21.ope-count	S.bbric_RbmLong_GGK21.ope-rpkm	S.bbric_RbmSmall_GGK36.ope-count	S.bbric_RbmSmall_GGK36.ope-rpkm
SBBRIC1.1	SBBRIC1	384	685	-	302	gene		5	14054.58	2	260.41
SBBRIC1.10	SBBRIC1	9162	11163	+	2002	gene		14	5936.34	107	2101.63
SBBRIC1.100	SBBRIC1	90298	90594	+	297	gene		0	0.00	9	1191.58

Fichier de résultats

Comptage

- nombre de lectures alignées par objet
- utilisé pour les analyses d'expression différentielle après normalisation

[Nat Methods](#), 2008 Jul;5(7):621-8. doi: 10.1038/nmeth.1226. Epub 2008 May 30.

Mapping and quantifying mammalian transcriptomes by RNA-Seq.

[Mortazavi A](#)¹, [Williams BA](#), [McCue K](#), [Schaeffer L](#), [Wold B](#).

Normalisation RPKM

- **R**ead **P**er **K**ilobase per **M**illion mapped reads

$$\text{RPKM (X)} = \frac{\text{Nb de lectures/gène (comptage brut)}}{\text{million lectures alignées} \times \text{taille objet biologique(kb)}}$$

Normalisation entre librairies

Normalisation entre objets

Limites : normalisation RPKM

BRIEFINGS IN BIOINFORMATICS, VOL 14, NO 6, 671–683
Advance Access published on 17 September 2012

doi:10.1093/bib/bbs046

A comprehensive evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis

Marie-Agnès Dillies*, Andrea Rau*, Julie Aubert*, Christelle Hennequet-Antier*, Marine Jeanmougin*, Nicolas Servant*, Céline Keime*, Guillemette Marot, David Castel, Jordi Estelle, Gregory Guernec, Bernd Jagla, Luc Jouneau, Denis Laloë, Caroline Le Gall, Brigitte Schaeffer, Stéphane Le Crom*, Mickaël Guedj*, Florence Jaffrézic* and on behalf of The French StatOmique Consortium

Key points

- Normalization of RNA-seq data in the context of differential analysis is essential in order to account for the presence of systematic variation between samples as well as differences in library composition.
- The Total Count and RPKM normalization methods, both of which are still widely in use, are ineffective and should be definitively abandoned in the context of differential analysis.
- Only the DESeq and TMM normalization methods are robust to the presence of different library sizes and widely different library compositions, both of which are typical of real RNA-seq data.

Table 3: Summary of comparison results for the seven normalization methods under consideration

Method	Distribution	Intra-Variance	Housekeeping	Clustering	False-positive rate
TC	–	+	+	–	–
UQ	++	++	+	++	–
Med	++	++	–	++	–
DESeq	++	++	++	++	++
TMM	++	++	++	++	++
Q	++	–	+	++	–
RPKM	–	+	+	–	–

A '–' indicates that the method provided unsatisfactory results for the given criterion, while a '+' and '++' indicate satisfactory and very satisfactory results for the given criterion.

Aller plus loin...

- Normalisation et analyses statistiques pour identifier les gènes différentiellement exprimés avec R
 - DESeq ou DESeq2 (bioconductor)
 - TMM (edgeR)

Je veux mesurer l'expression de gènes connus à partir de données de séquençage

Je vais sur

L'outil permet

L'outil ne permet pas

Paramètres clés

Pièges

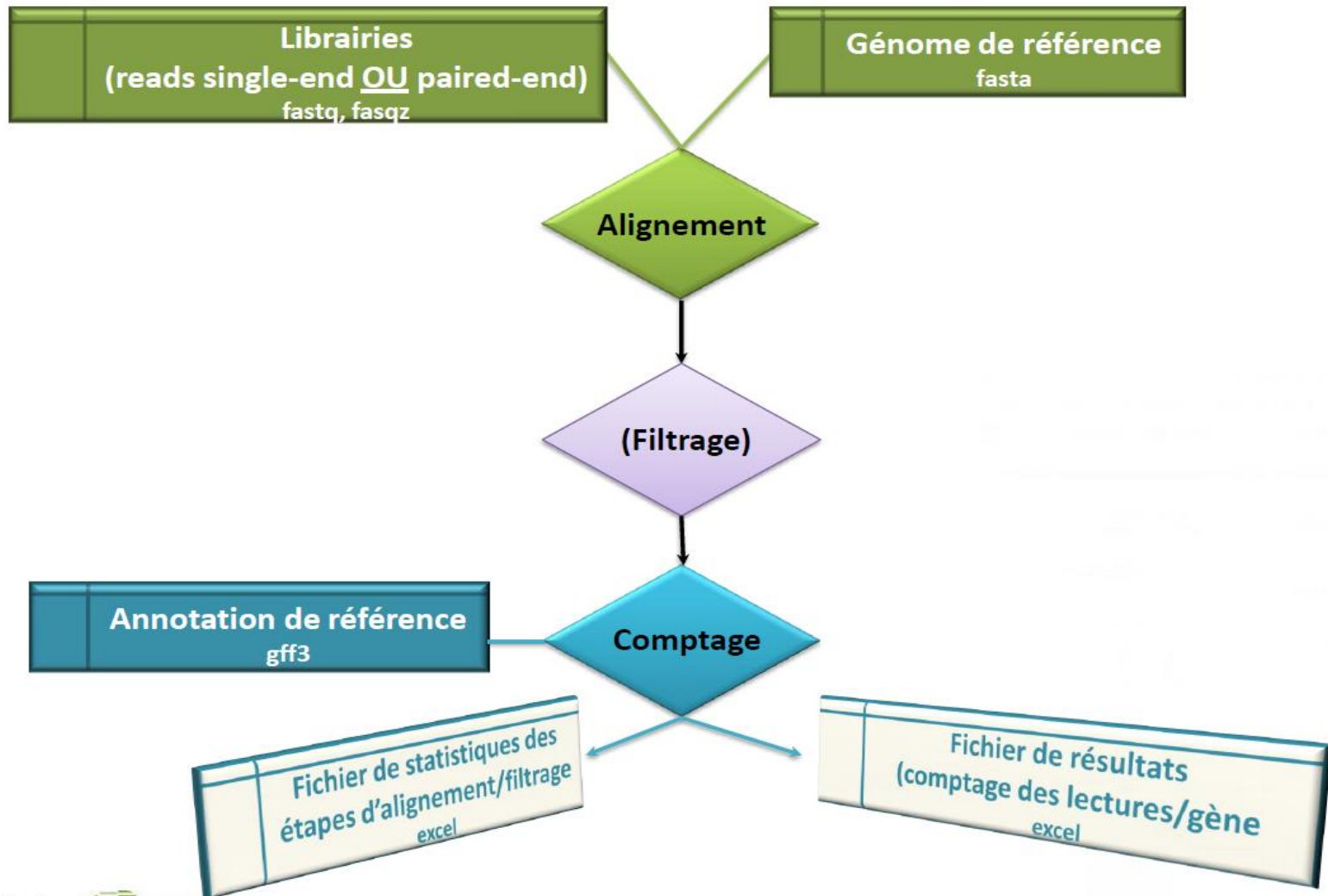
Formats et fichiers

CONTRÔLE QUALITÉ DE MESURES D'EXPRESSION

Responsable et intervenant principal: Erika Sallet

Expert: Adeline Simon et Joseph Tran

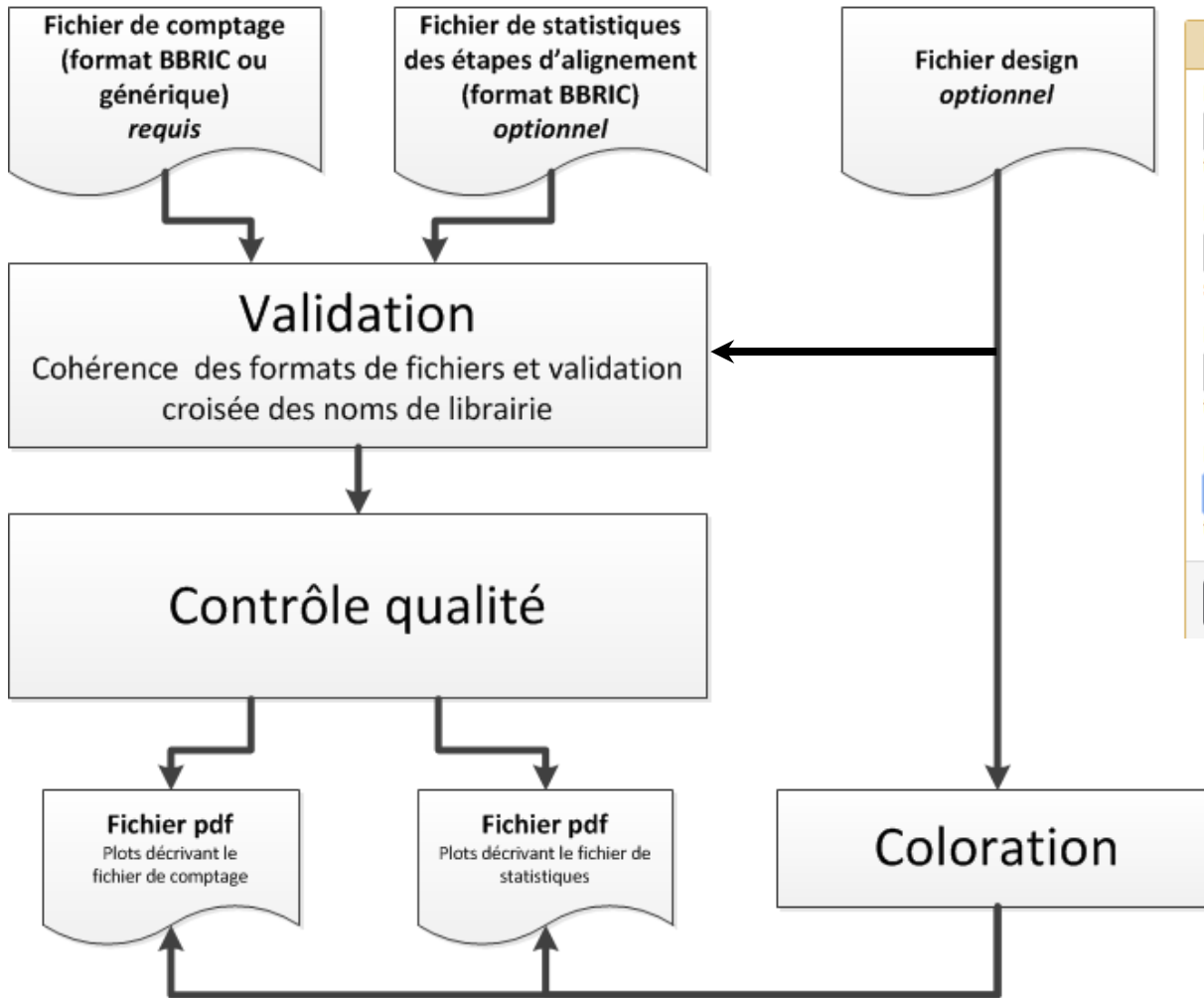
Rappel : pipeline Mesures de l'expression



Objectif et périmètre du pipeline

- Objectif : identifier facilement des problèmes de qualité dans les résultats issus d'une mesure d'expression
- Périmètre : sortie du pipeline BBRIC ou pipelines externes équivalents

Description du Workflow



RNAseq libraries (version 0.1)

RNASeq count expression file :

2: DataTest_Count_expression_BBRIC.txt ▾
tabular file, see below(Input files) for format description

RNASeq count expression file format:

BBRIC ▾
see below(Input files) for description

RNASeq mapping statistics file :

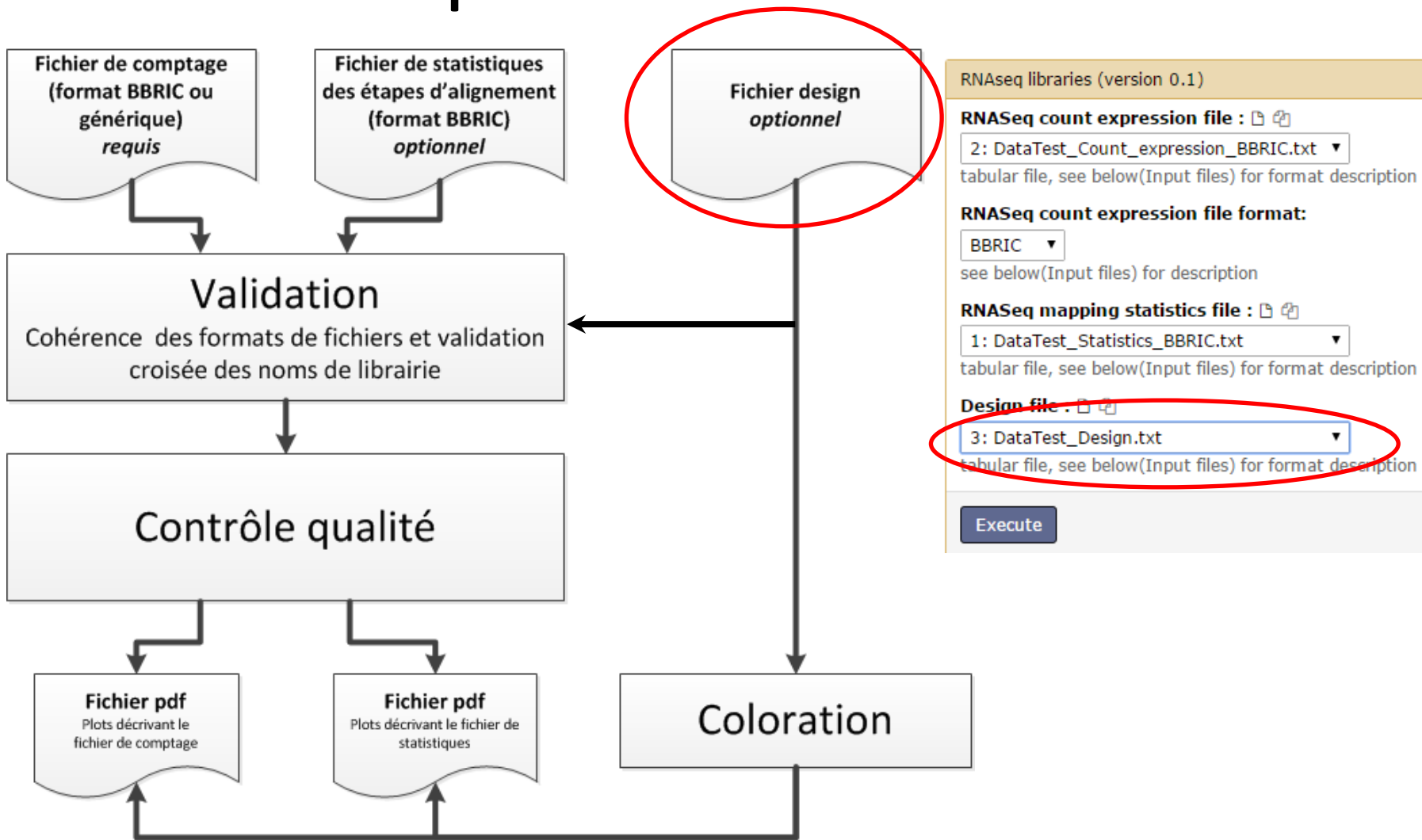
1: DataTest_Statistics_BBRIC.txt ▾
tabular file, see below(Input files) for format description

Design file :

3: DataTest_Design.txt ▾
tabular file, see below(Input files) for format description

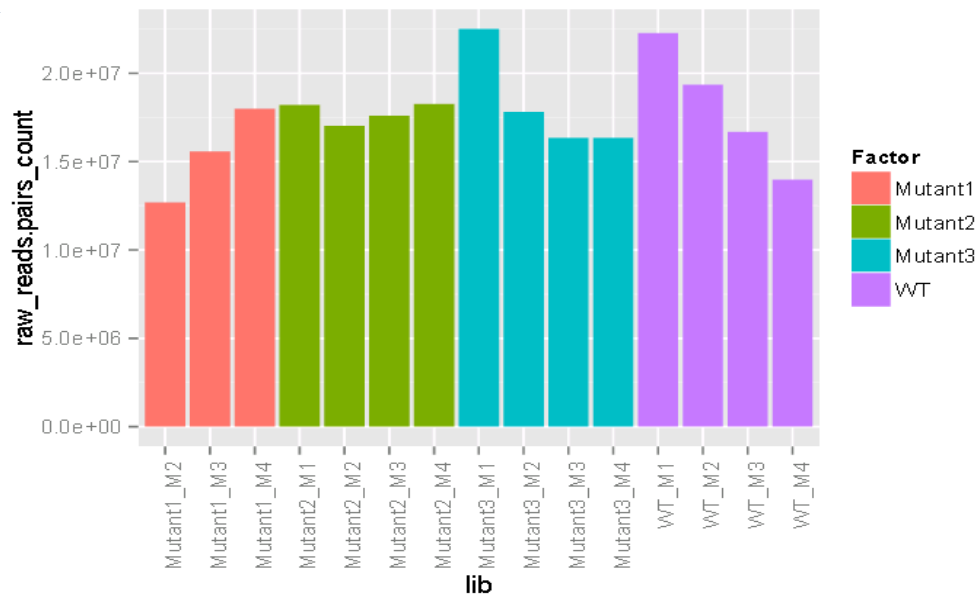
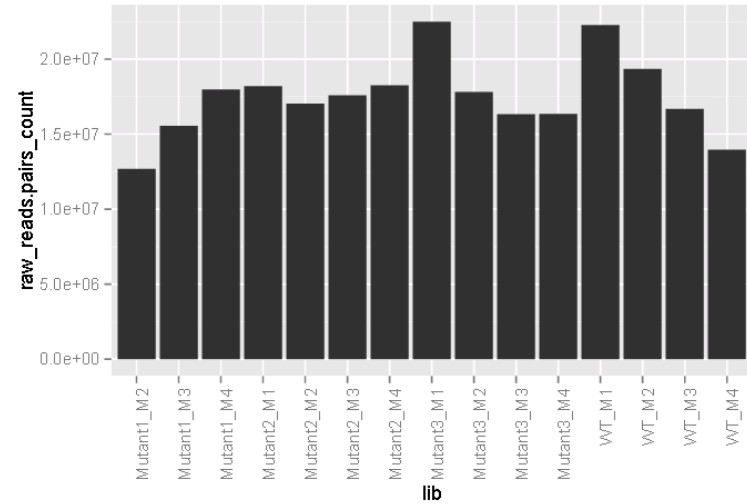
Execute

Description du Workflow

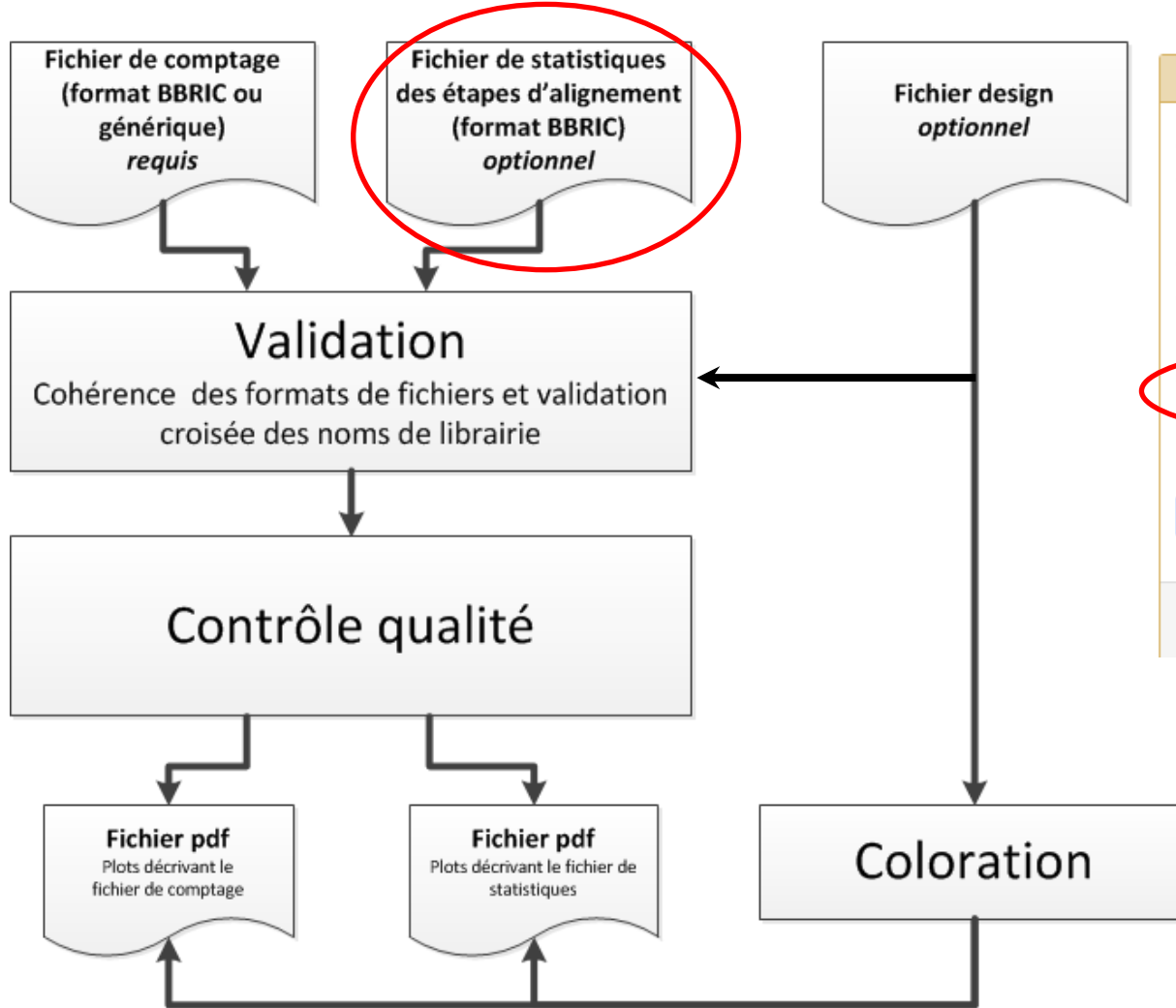


Format du fichier design

lib	Strain
Mutant1_M2	Mutant1
Mutant1_M3	Mutant1
Mutant1_M4	Mutant1
Mutant2_M1	Mutant2
Mutant2_M2	Mutant2
Mutant2_M3	Mutant2
Mutant2_M4	Mutant2
Mutant3_M3	Mutant3
Mutant3_M4	Mutant3
Mutant3_M1	Mutant3
Mutant3_M2	Mutant3
WT_M1	WT
WT_M2	WT
WT_M3	WT
WT_M4	WT



Description du Workflow



RNAseq libraries (version 0.1)

RNASeq count expression file :

2: DataTest_Count_expression_BBRIC.txt ▾
tabular file, see below(Input files) for format description

RNASeq count expression file format:

BBRIC ▾
see below(Input files) for description

RNASeq mapping statistics file :

1: DataTest_Statistics_BBRIC.txt ▾
tabular file, see below(Input files) for format description

Design file :

3: DataTest_Design.txt ▾
tabular file, see below(Input files) for format description

Execute

Fichier de comptage

Nombre d'alignements spécifiques sur génome

Nombre d'alignements sur génome

Nombre de lectures brutes

Nombre de lectures alignées sur le génome

Nombre d'alignements sur objet biologique

alignements sur objet biologique / alignements spécifiques sur génome

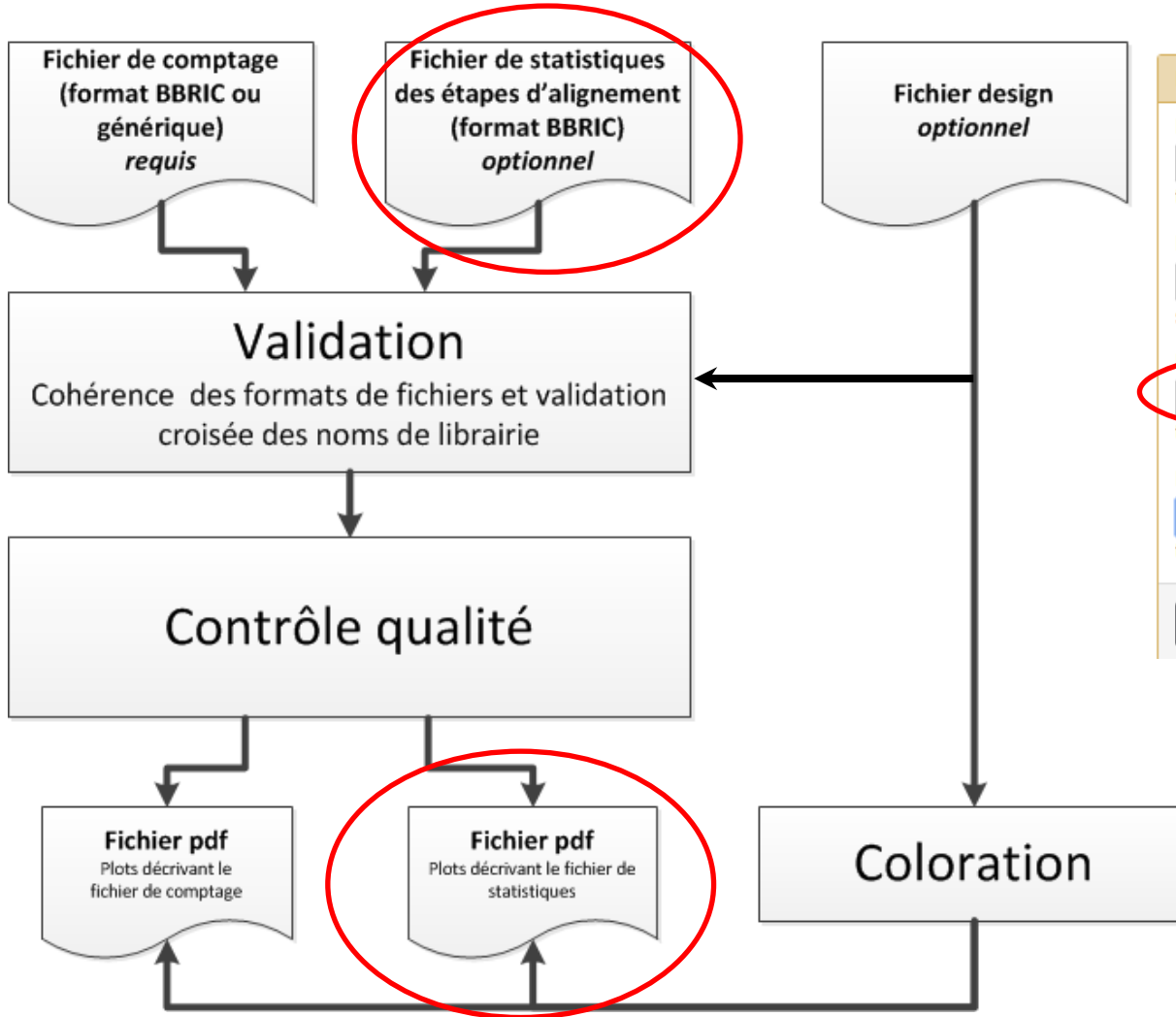
lib	specific_hits	mapping_hits	raw_reads/pairs_count	mapped_reads/pairs_count	feature_overlapping_hits	feature_overlapping_hits/specific_hits_percent
S.bbric_Rb mLong_GGK 21.ope	1178	1181	9697	1179	1153	97.88
S.bbric_Rb mSmall_GG K36.ope	24405	24411	64272	24407	20674	84.71

Contig/Chromosome



1181 hits totaux
1179 lectures qui s'alignent

Description du Workflow



RNAseq libraries (version 0.1)

RNAseq count expression file :
tabular file, see below(Input files) for format description

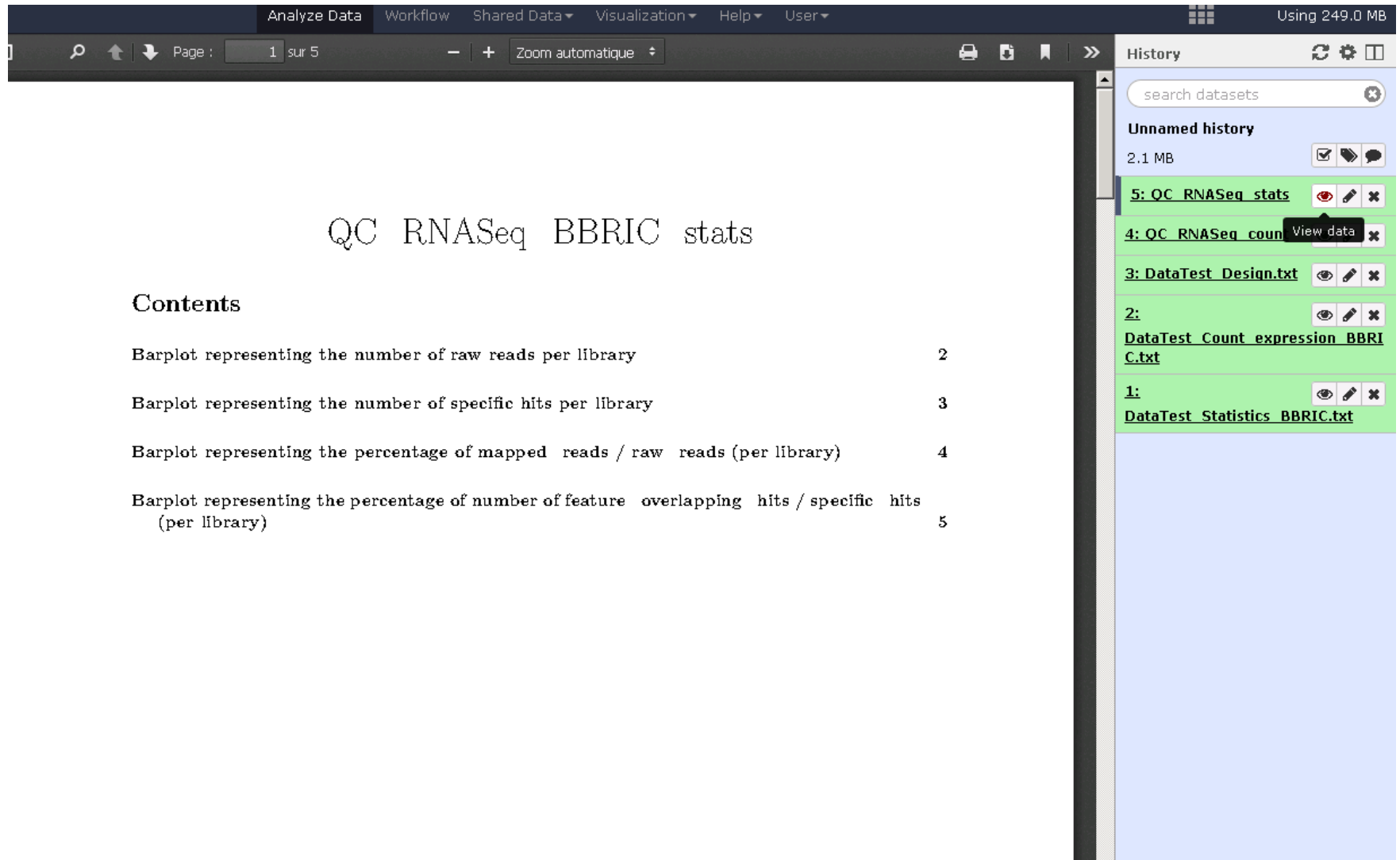
RNAseq count expression file format:

see below(Input files) for description

RNAseq mapping statistics file :
tabular file, see below(Input files) for format description

Design file :
tabular file, see below(Input files) for format description

Plots à partir du fichier de comptage



The screenshot shows a software interface with a dark top bar containing menu items: 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Help', and 'User'. The top right corner indicates 'Using 249.0 MB'. Below the menu bar, there is a navigation area with 'Page : 1 sur 5' and a 'Zoom automatique' button. The main content area displays a table of contents for a document titled 'QC RNASeq BBRIC stats'. To the right, a 'History' panel lists several datasets, with the most recent one being '5: QC RNASeq_stats'.

QC RNASeq BBRIC stats

Contents

Barplot representing the number of raw reads per library	2
Barplot representing the number of specific hits per library	3
Barplot representing the percentage of mapped reads / raw reads (per library)	4
Barplot representing the percentage of number of feature overlapping hits / specific hits (per library)	5

History

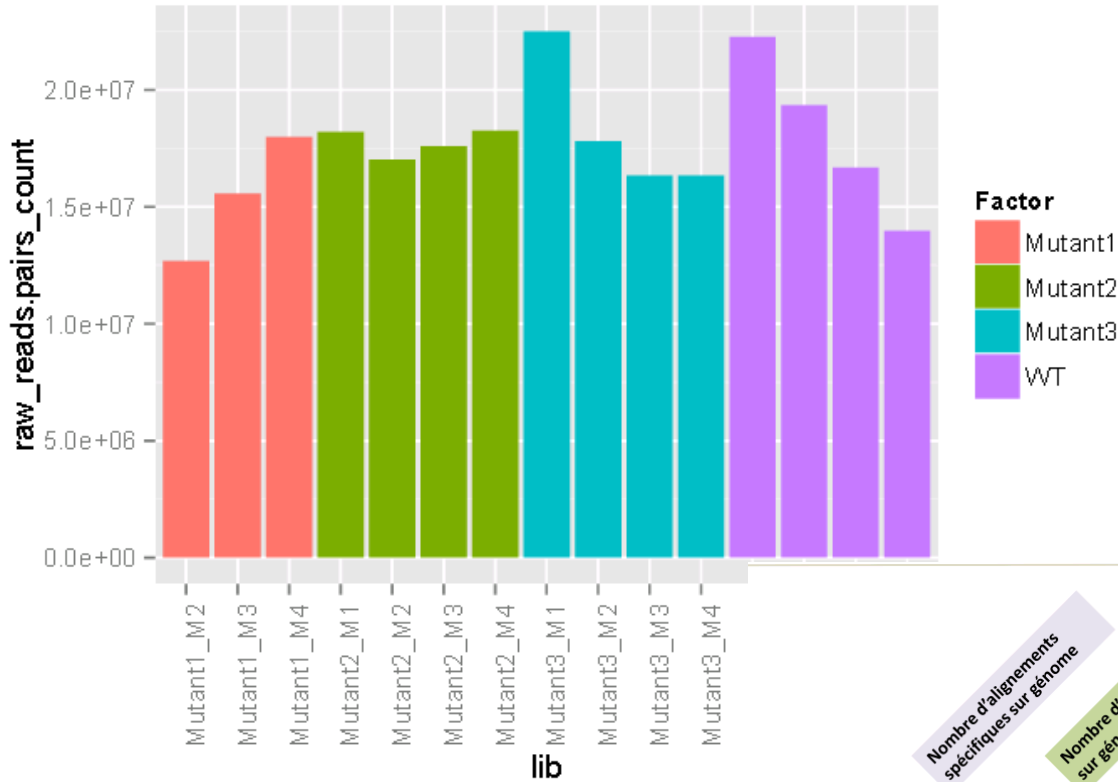
search datasets

Unnamed history
2.1 MB

- 5: QC RNASeq_stats
- 4: QC RNASeq_count View data
- 3: DataTest_Design.txt
- 2: DataTest_Count_expression_BBRIC.txt
- 1: DataTest_Statistics_BBRIC.txt

Nombre de lectures brutes par librairie

Barplot representing the number of raw reads per library



Objectif :

Vérifie que le nombre de reads est équivalent entre les échantillons pour repérer les problèmes expérimentaux

Question :

Le prestataire a-t-il bien fourni le nombre de reads « garanti » pour tous les échantillons ?

lib	specific_hits	mapping_hits	raw_reads/pairs_count	mapped_reads/pairs_count	feature_overlapping_hits	feature_overlapping_hits /specific_hits_percent
S.bbric_Rb mLong_GGK 21.ope	1178	1181	9697	1179	1153	97.88
S.bbric_Rb mSmall_GG K36.ope	24405	24411	64272	24407	20674	84.71

Nombre d'alignements spécifiques sur génome

Nombre d'alignements sur génome

Nombre de lectures brutes

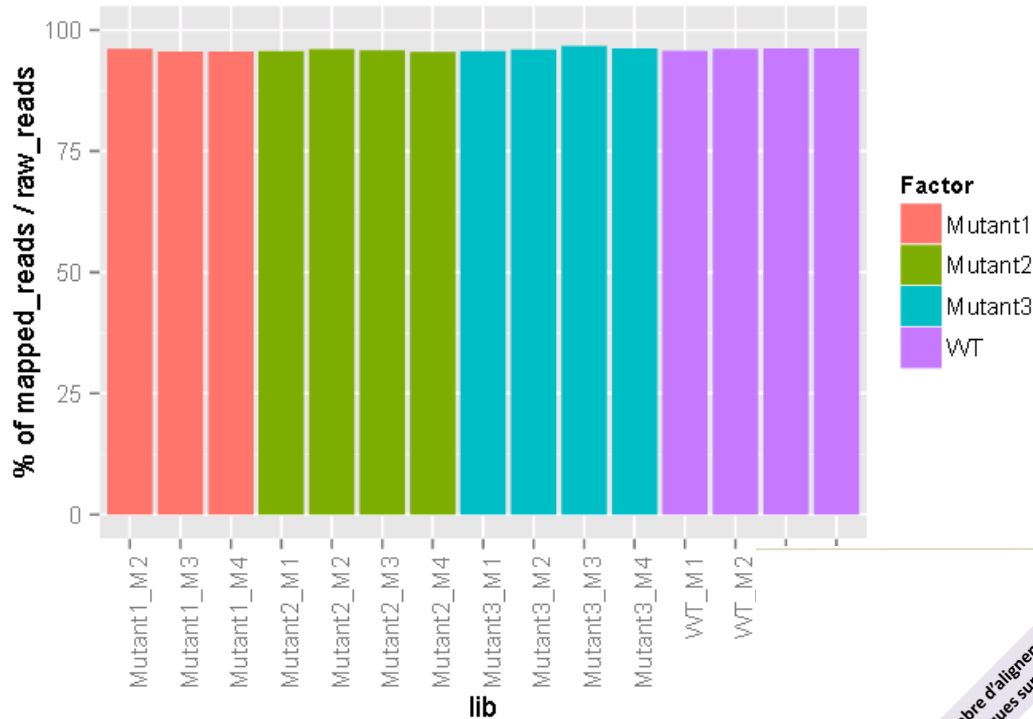
Nombre de lectures alignées sur le génome

Nombre d'alignements sur objet biologique

alignements sur objet biologique / alignements spécifiques sur génome

Proportion de lectures mappées

Barplot representing the percentage of mapped_reads / raw_reads (per library)



Objectif :

Vérifie que la qualité du mapping est la même pour toutes les librairies.

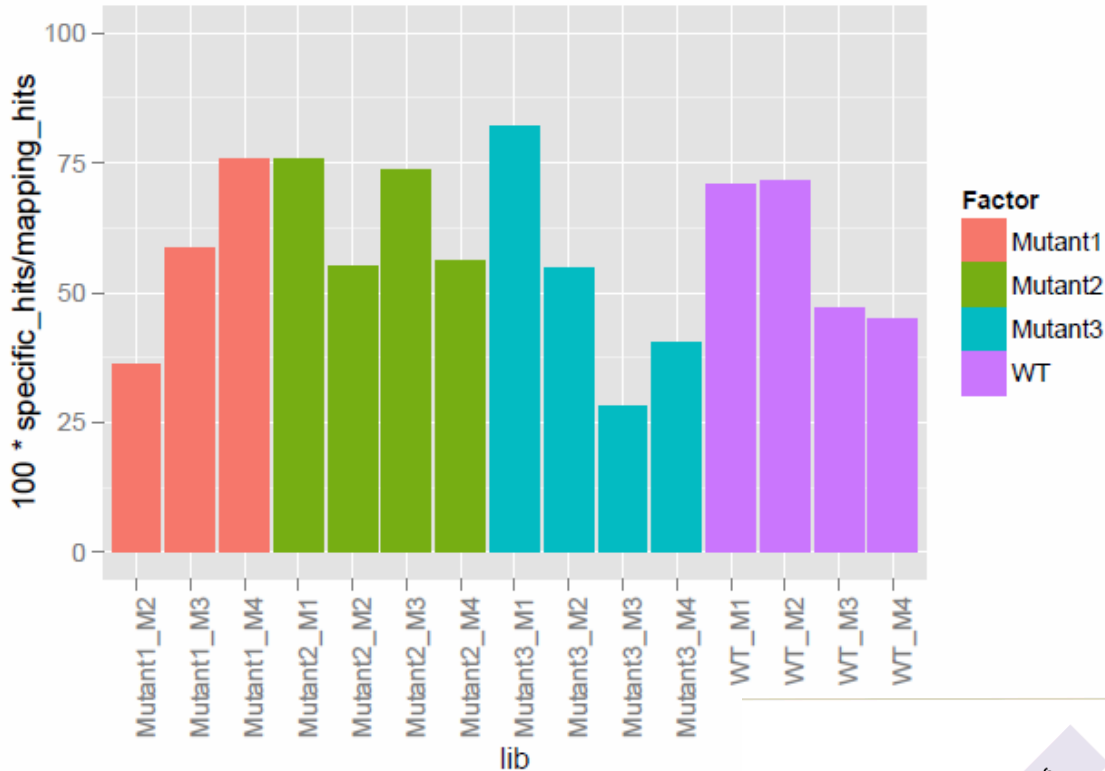
Si le pourcentage est trop bas, peut être la preuve d'une contamination

Questions :

- Un échantillon est-il contaminé par un autre organisme ?
- Les paramètres de mapping sont-ils corrects pour tous les échantillons ?

lib	specific_hits	mapping_hits	Nombre d'alignements sur génome		Nombre de lectures alignées sur le génome	
			raw_reads/pairs_count	mapped_reads/pairs_count	feature_overlapping_hits	feature_overlapping_hits /specific_hits_percent
S.bbric_Rb mLong_GGK 21.ope	1178	1181	9697	1179	1153	97.88
S.bbric_Rb mSmall_GG K36.ope	24405	24411	64272	24407	20674	84.71

Proportion de hits avec alignement spécifique



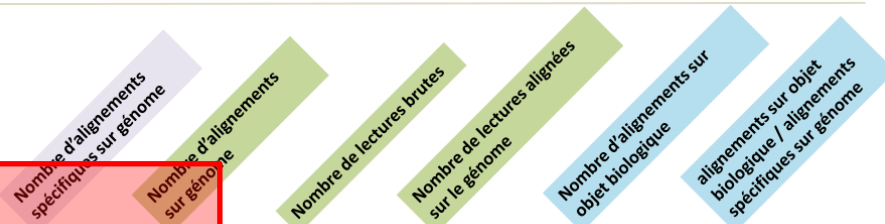
Objectif :

Déterminer la proportion de hits spécifiques (qui ne s'alignent qu'à un seul endroit dans le génome) par rapport à l'ensemble des hits mappant sur le génome

Questions :

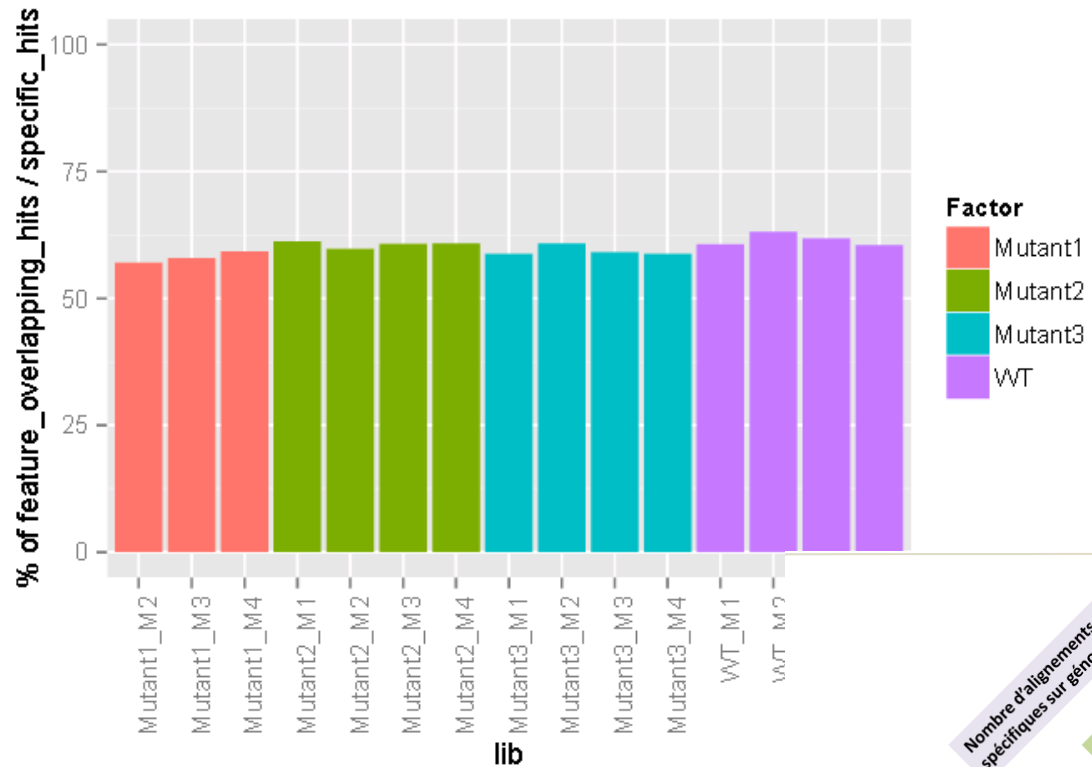
- Les paramètres de mapping permettent-ils d'assurer des hits spécifiques ?
- Un échantillon est-il contaminé (ARN ribosomique par exemple) ?

lib	specific_hits	mapping_hits	raw_reads/pairs_count	mapped_reads/pairs_count	feature_overlap/ing_hits	feature_overlapping_hits/specific_hits_percent
S.bbric_Rb mLong_GGK 21.ope	1178	1181	9697	1179	1153	97.88
S.bbric_Rb mSmall_GG K36.ope	24405	24411	64272	24407	20674	84.71



Proportion de hits recouvrant une feature

Barplot representing the percentage of `feature_overlapping_hits / specific_hits` (per library)

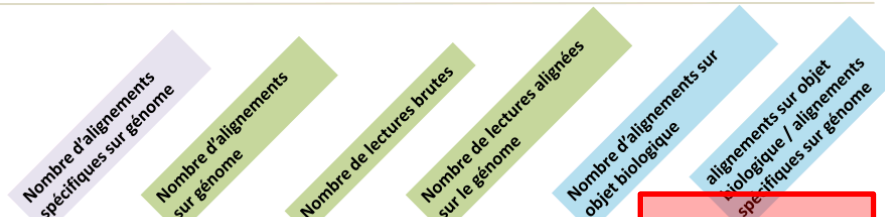


Objectif :

Vérifie que la qualité du mapping est la même pour toutes les librairies

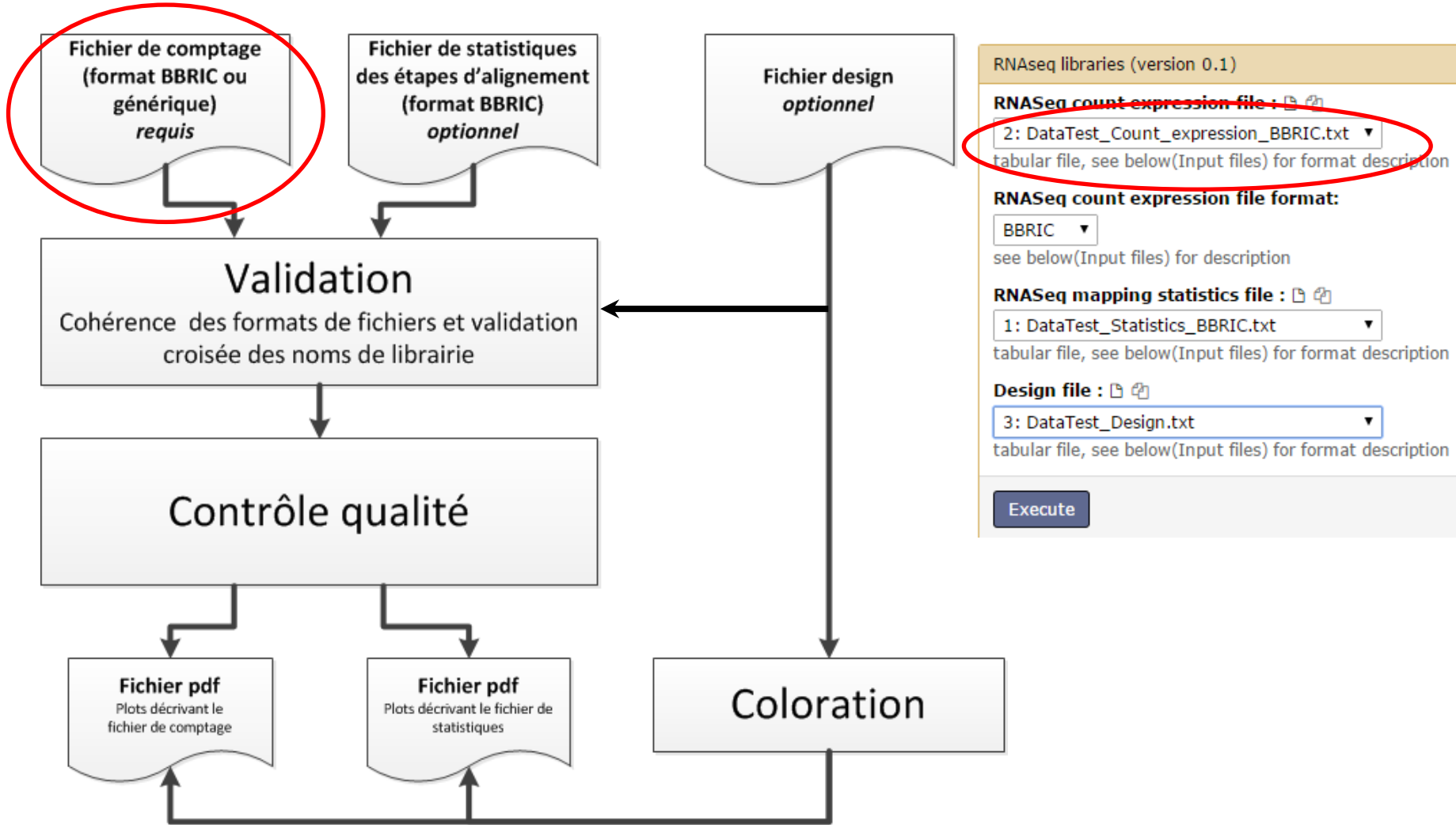
Questions :

- L'annotation structurale est-elle valide ?
- Les fichiers pair-end sont-ils utilisés dans la bonne orientation ?



lib	specific_hits	mapping_hits	raw_reads/pairs_count	mapped_reads/pairs_count	feature_overlapping_hits	feature_overlapping_hits / specific_hits_percent
S.bbric_Rb mLong_GGK 21.ope	1178	1181	9697	1179	1153	97.88
S.bbric_Rb mSmall_GG K36.ope	24405	24411	64272	24407	20674	84.71

Description du Workflow



Fichier de comptage (format BBRIC)

Gene/RNA ID

Contig Id

Positionnement du gène/RNA sur le contig

Début

Fin

Brin

Taille du Gène/RNA

Type (Gène/RNA)

Comptage brut (Lectures/Objet)

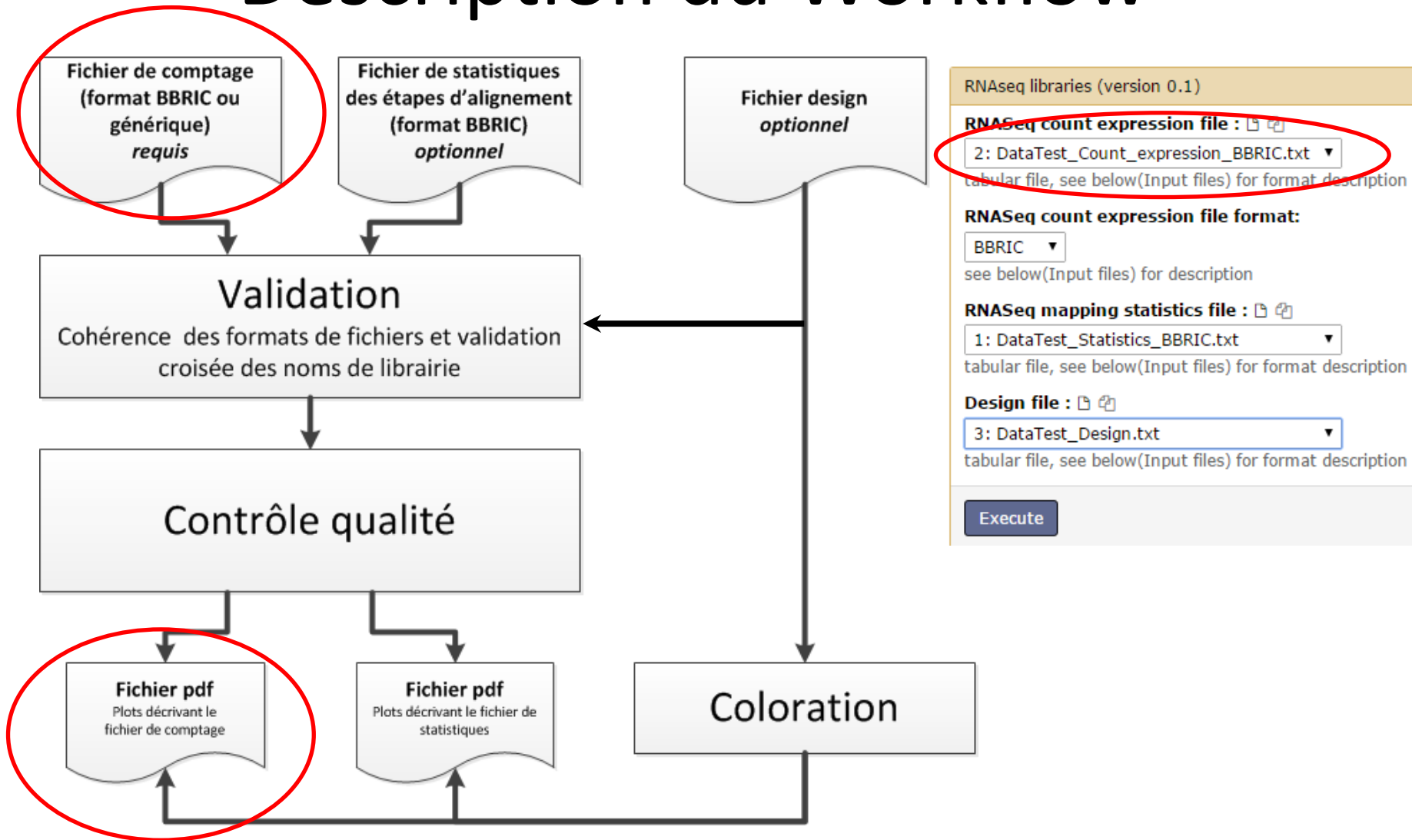
Comptage normalisé (RPKM)

id	0_seqid	1_start	2_end	3_strand	4_length	5_type	6_Note	S.bbric_RbmLong_GGK21.op e-count	S.bbric_RbmLong_GGK21.op e-rpkm	S.bbric_RbmSmall_GGK36.op e-count	S.bbric_RbmSmall_GGK36.op e-rpkm
SBBRIC1.1	SBBRIC1	384	685	-	302	gene		5	14054.58	2	260.41
SBBRIC1.10	SBBRIC1	9162	11163	+	2002	gene		14	5936.34	107	2101.63
SBBRIC1.100	SBBRIC1	90298	90594	+	297	gene		0	0.00	9	1191.58

Fichier de comptage (format générique)

id	Lib1	Lib2
gene1	0	0
gene2	314	497
gene3	10991	11553

Description du Workflow



Résultats de comptage

The screenshot shows a web browser window displaying the Galaxy interface. The main content area shows a workflow titled "QC_RNASeq_generic_count" with a table of contents. The table of contents lists the following items:

Item	Page Number
Cluster dendrogram representing the hierarchical clustering of libraries (log values)	2
Barplot representing for each library the number of features (genes / RNA) that have at least 10 reads mapped	3
Box-plot (log values)	4
Summary of values	4

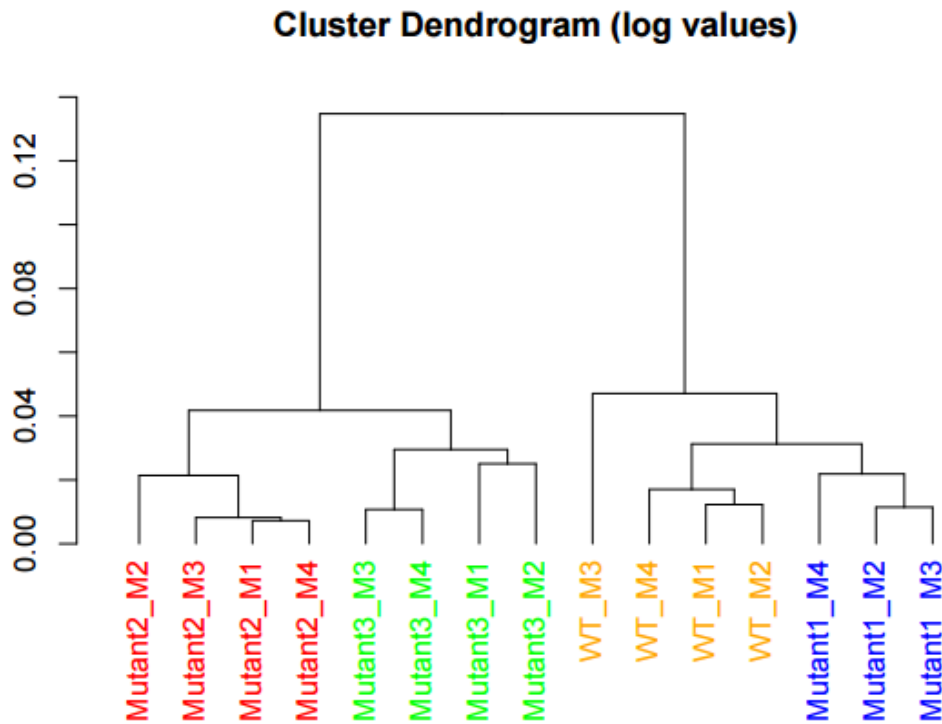
The right-hand side of the interface shows a "History" panel with a list of workflow steps, including:

- NP detection on data 47, data 49, and data 48
- 50: VCF of SNP detection on data 47, data 49, and data 48
- 49: Xbbriic genomic sequence.fasta
- 48: xbbriic2.fasta.gz
- 47: xbbriic1.fasta.gz
- 46: QC_RNASeq_stats
- 45: QC_RNASeq_count
- 44: QC_RNASeq_stats
- 43: QC_RNASeq_count
- 42: DataTest_Statistics_BBRIC.txt
- 41: DataTest_Design.txt
- 40: DataTest_Count_expression_generic.txt
- 39: DataTest_Count_expression_BBRIC.txt
- 38: Output gene data file
- 37: Output tabular file
- 36: Output SBML file
- 35: RalstoEcoliOrthologs.tab
- 34: iJO1366.xml
- 33: out.fasta

The browser address bar shows the URL: <https://bbriic-pipelines.toulouse.inra.fr/galaxy-dev/>. The page title is "Galaxy / BBRIC DEV". The status bar at the bottom indicates "Using 331.8 MB".

Clustering hiérarchique des bibliothèques

Cluster dendrogram representing the hierarchical clustering of libraries (log values)



Objectif:

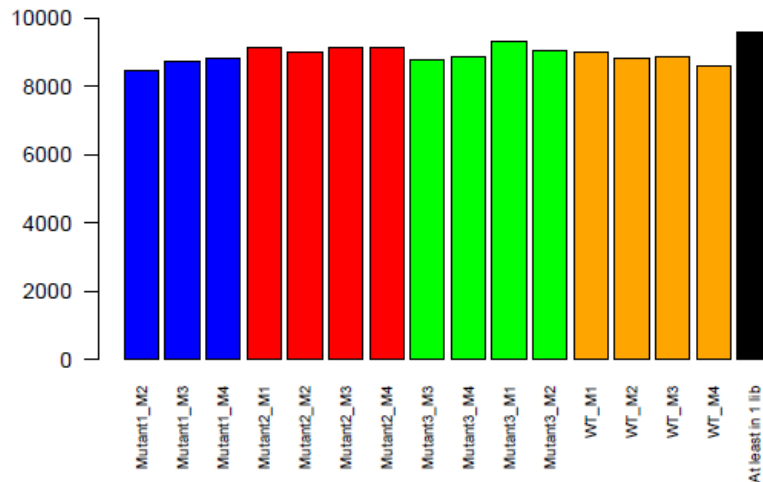
- Vérifier que les duplicats sont bien dans le même groupe
- Classification des conditions entre elles

Question:

Y-a-t-il une inversion des échantillons ?

Barplot representing the number of features mapped

Barplot representing for each library the number of features (genes / RNA) that have at least 10 reads mapped



```
##          lib nb features with at least 10 reads nb features total
## 1      Mutant1_M2                8471          10345
## 2      Mutant1_M3                8730          10345
## 3      Mutant1_M4                8803          10345
## 4      Mutant2_M1                9112          10345
## 5      Mutant2_M2                9001          10345
## 6      Mutant2_M3                9128          10345
## 7      Mutant2_M4                9121          10345
## 8      Mutant3_M3                8796          10345
## 9      Mutant3_M4                8872          10345
## 10     Mutant3_M1                9300          10345
## 11     Mutant3_M2                9024          10345
## 12           WT_M1                9006          10345
## 13           WT_M2                8815          10345
## 14           WT_M3                8857          10345
## 15           WT_M4                8605          10345
## 16 At least in 1 lib            9563          10345
```

Transcrits prédits sans reads : $10345 - 9563 = 782$

Objectifs :

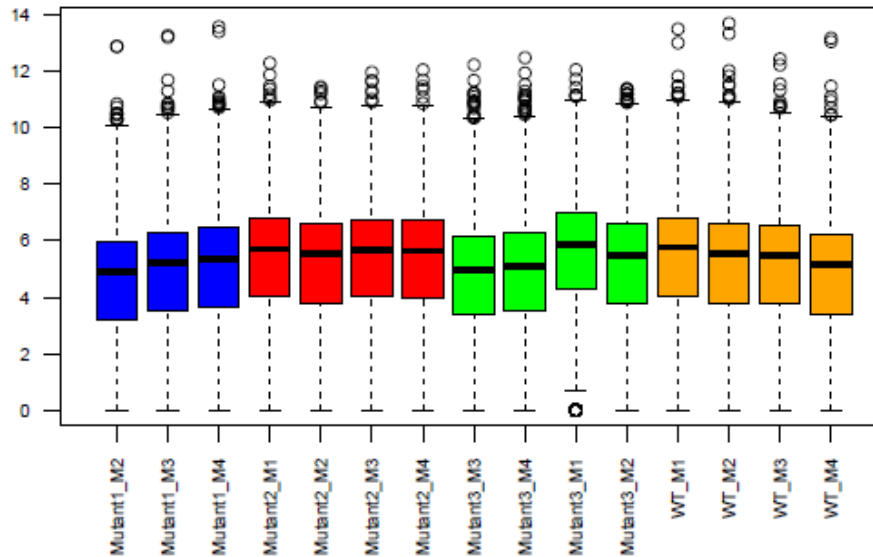
- Permet de connaître le nombre de transcrits exploitables par librairie
- Permet de valider la prédiction des transcrits
- Permet de repérer un problème expérimental si l'histogramme est déséquilibré

Question :

Y-a-t-il des soucis avec l'annotation structurale?

Box plot du nombre de reads par feature

BoxPlot of log(count)



Summary of values

```
## Mutant1_M2      Mutant1_M3      Mutant1_M4
## Min.   :    0.0   Min.   :    0.0   Min.   :    0.0
## 1st Qu.:   23.0   1st Qu.:   32.0   1st Qu.:   37.0
## Median :  134.0   Median :  189.0   Median :  211.0
## Mean   :  529.2   Mean   :  751.1   Mean   :  936.7
## 3rd Qu.:  386.0   3rd Qu.:  548.0   3rd Qu.:  640.0
## Max.   :390252.0   Max.   :573790.0   Max.   :780274.0
## Mutant2_M1      Mutant2_M2      Mutant2_M3
## Min.   :    0.0   Min.   :    0.0   Min.   :    0.0
## 1st Qu.:   54.0   1st Qu.:   43.0   1st Qu.:   54.0
## Median :  297.0   Median :  254.0   Median :  289.0
## Mean   :  982.1   Mean   :  838.7   Mean   :  936.1
## 3rd Qu.:  884.0   3rd Qu.:  744.0   3rd Qu.:  843.0
## Max.   :216767.0   Max.   :93211.0   Max.   :156509.0
```

Objectifs :

- Vérifier que la distribution des valeurs de comptage est homogène à travers les librairies

- Permet de vérifier que le regroupement par la classification hiérarchique n'est pas dû à une ressemblance de la distribution des valeurs

Question :

- La distribution des valeurs de comptage nous permet-elle de lancer une analyse statistique efficace ?

Je veux contrôler la qualité de mes banques RNAseq et de mes mesures d'expression

Je vais sur

L'outil permet

L'outil ne permet pas

Paramètres clés

Pièges

Formats et fichiers



ANALYSE COMPARATIVE DE PROTEOMES D'ESPECES APPARENTEES (ORTHOMCL)

Experts : Martial Briand, Sébastien Carrère, Ludovic Cottret, Corinne Rancurel

C'est déterminer ce qui est **spécifique** à chaque espèce et ce qui est **commun** à chaque espèce

La comparaison génomique sous-entend que les génomes à comparer ont un ancêtre commun :

C'est la recherche de parenté entre les gènes de différents génomes

Un génome n'est pas statique :
il évolue au cours du temps avec des additions et des pertes de gènes

⇒ Il y a une évolution des gènes ancestraux avec des différences qualitatives (changement de la nature des gènes) et des différences quantitatives (variations du nombre de gènes)

C'est comparer les gènes (protéines) homologues entre eux au niveau de leur séquence

→ Identifier les gènes orthologues et paralogues

la motivation première

C'est inférer des connaissances sur une séquence à partir des connaissances attachées à une autre

Quelques applications :

- constituer des familles de gènes
- identifier les gènes uniques de chaque génome
- aider à l'annotation fonctionnelle
- étudier l'expansion/ réduction de familles de gènes au cours du temps
- identifier les orthologues de gènes décrits chez les organismes modèles (ex. Arabidopsis, *C. elegans*, Drosophile) dans nos organismes d'intérêt (ex. Légumineuses, Meloidogyne, insectes ravageurs)

Qu'est-ce qu'une relation d'orthologie

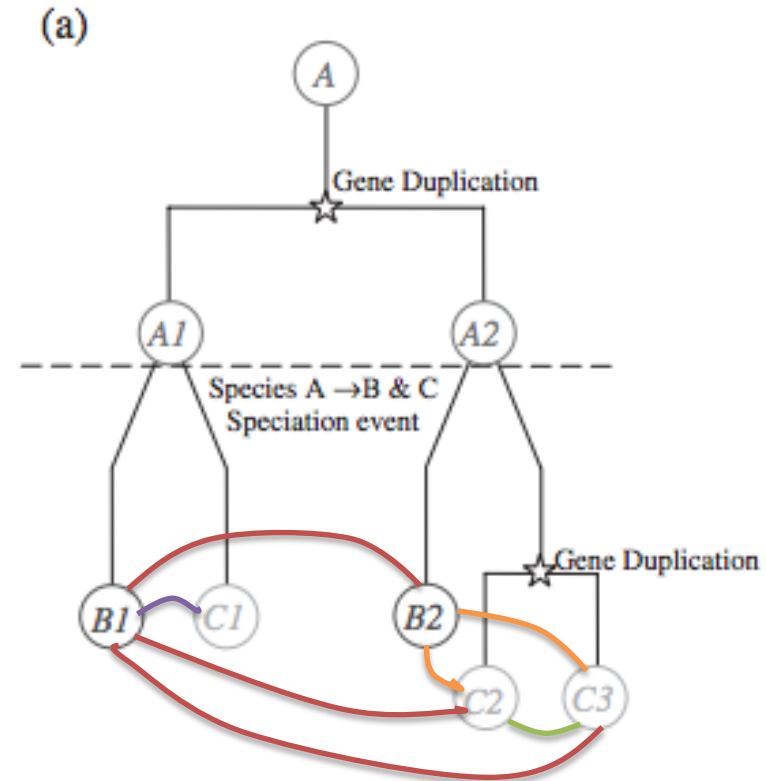
Définitions :

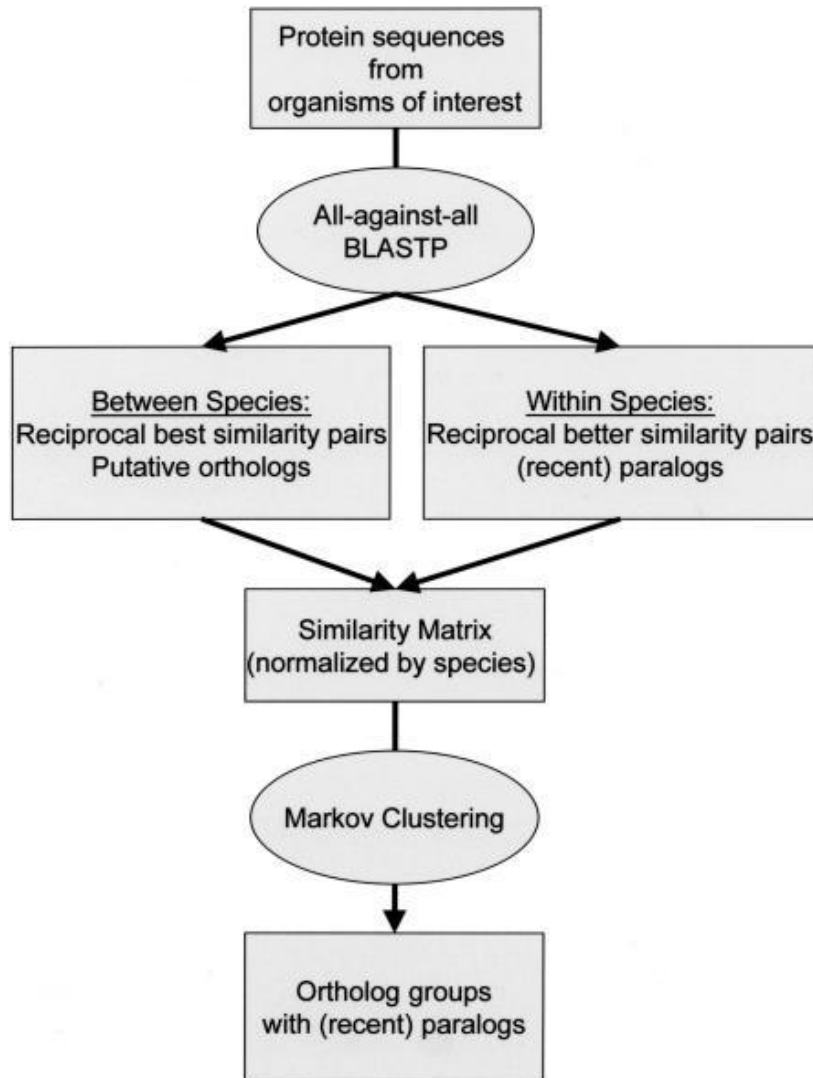
Paralogues : gènes qui ont co-évolué suite à un événement de duplication dans un génome

Orthologues : gènes qui ont un même ancêtre commun et qui ont co-évolué suite à un événement de spéciation

In-paralogues : paralogues qui ont été dupliqués suite à un événement de spéciation, ils sont alors **co-orthologues** aux gènes de l'autre espèce

Out-paralogues : paralogues qui ont été dupliqués avant un événement de spéciation (ils n'appartiennent pas nécessairement à la même espèce)

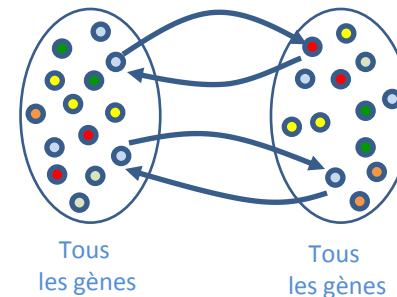




OrthoMCL est un logiciel qui construit des clusters d'orthologues (incluant les in-paralogues récents) à partir de fichiers multifasta contenant des CDS.

La méthode utilisée est basée sur :

la méthode du « blast hit réciproque » (BHR)
recherche des meilleures similitudes réciproques



+

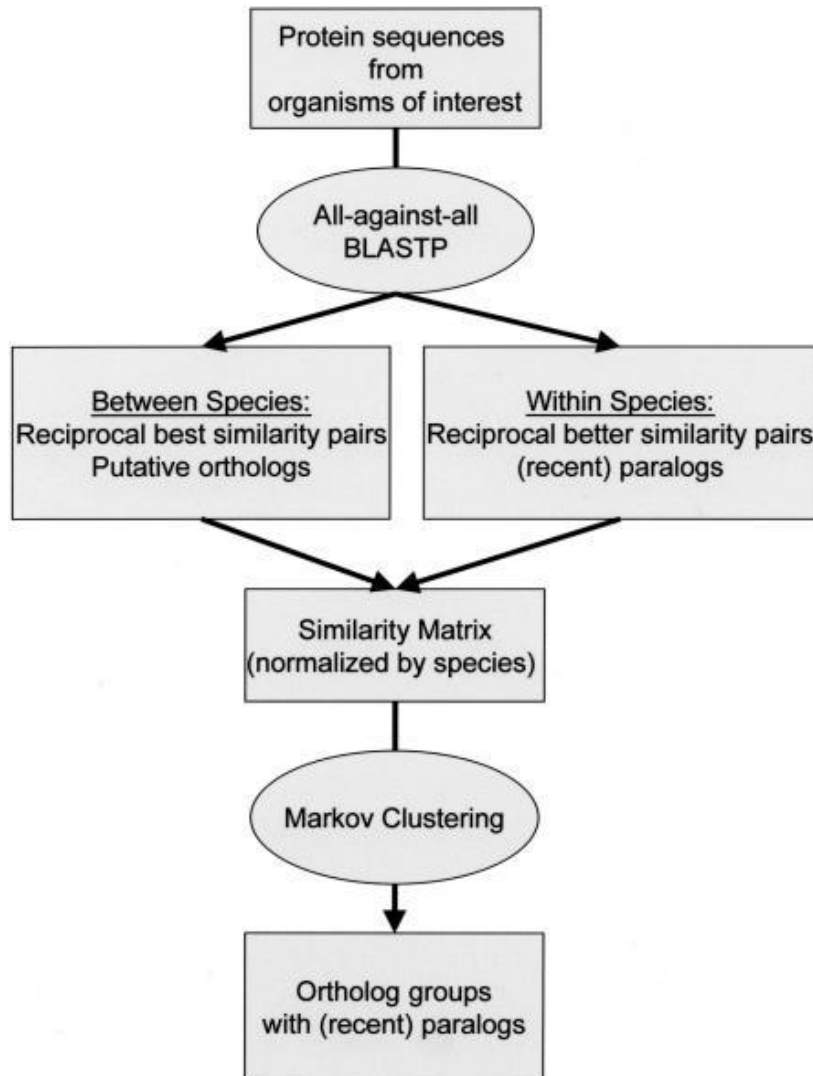
l'algorithme MCL (Markov Cluster Algorithm)
pour trouver des regroupements d'éléments similaires

Approche = création du graphe de similarité puis partitionnement (clustering)

Li et al.,

2003





pv_cutoff=<Float>

1e-5 (DEFAULT)

P-Value or E-Value Cutoff in BLAST search and/or ortholog clustering
=> Plus la e-value est faible, plus l'alignement est significatif. On fixe un seuil ("*cutoff*") au dessus duquel on ne veut pas conserver les résultats.

pi_cutoff=<Int>

0 (DEFAULT)

Percent Identity Cutoff <0-100> in ortholog clustering
=> Utiliser pour éliminer les hits avec de petites valeurs de pourcentage d'identité (pi). Le pi est la moyenne des pi de tous les hsp pour cette paire (query-subject).

pmatch_cutoff=<Int>

0 (DEFAULT)

Percent Match Cutoff <0-100> in ortholog clustering
=> Pour former une paire, le nombre d'acide aminés de la protéine la plus courte doit être inclus d'au moins **pmatch_cutoff** % de la longueur de l'autre séquence.

inflation=<Float>

1.5 (DEFAULT)

Markov Inflation Index, used in MCL algorithm,
Increasing this index increases cluster tightness, and the number of clusters
=> Une valeur de 1,5 donne un équilibre entre sensibilité et sélectivité.

Li et al.,

2003



Extrait de quelques groupes d'un fichier orthomcl.out

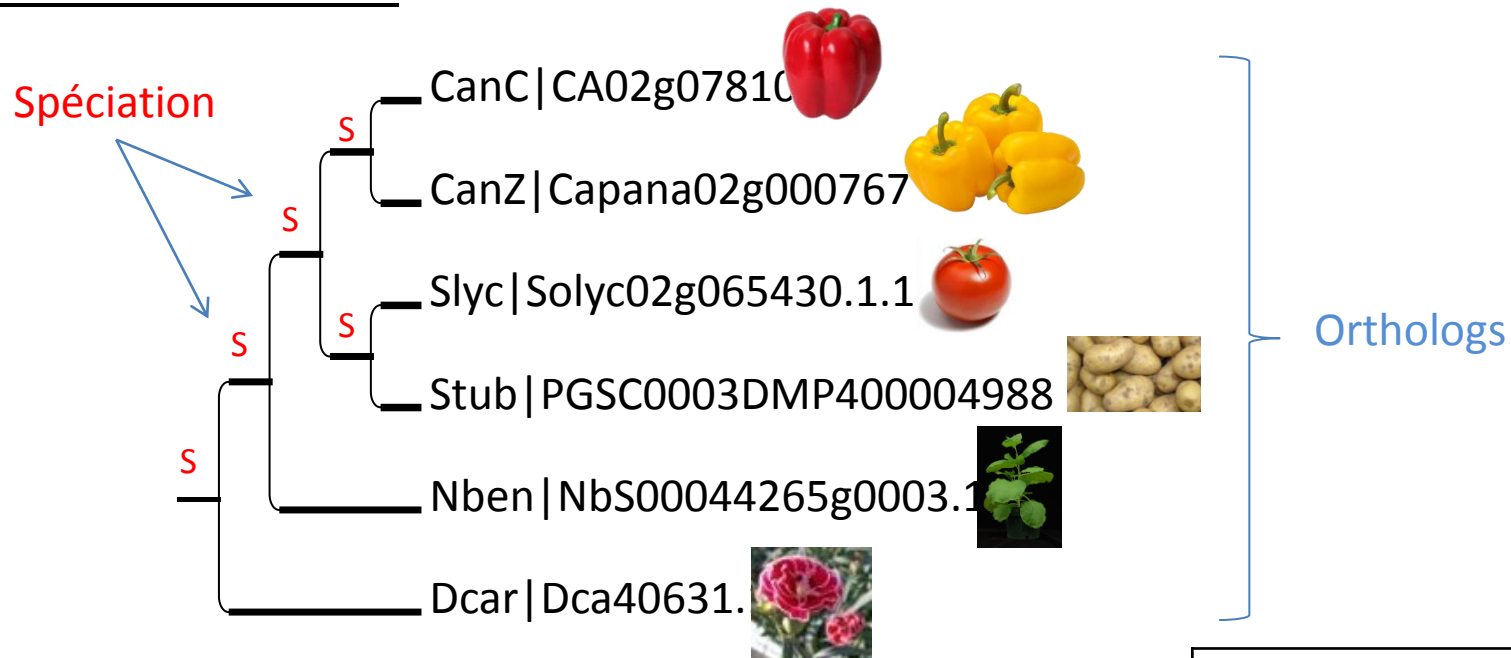
GRP22264: Mdom|MDP0000120086 Mdom|MDP0000255483 Mdom|MDP0000287109
Mdom|MDP0000388642 Mdom|MDP0000554052 Mdom|MDP0000298540
Mdom|MDP0000570301

GRP26560: Nben|NbS00044265g0003.1 Dcar|Dca40631.1 CanC|CA02g07810
CanZ|Capana02g000767 Slyc|Solyc02g065430.1.1 Stub|PGSC0003DMP400004988

GRP27024: Bdis|Bradi4g23995.1 Bdis|Bradi4g24001.1 Osa|BGIOGA001239-PA
Ljap|chr3.CM0786.230.r2.d OsaJ|LOC_Os11g15620.1

GRP26560: Nben|NbS00044265g0003.1 Dcar|Dca40631.1 CanC|CA02g07810
 CanZ|Capana02g000767 Slyc|Solyc02g065430.1.1
 Stub|PGSC0003DMP400004988

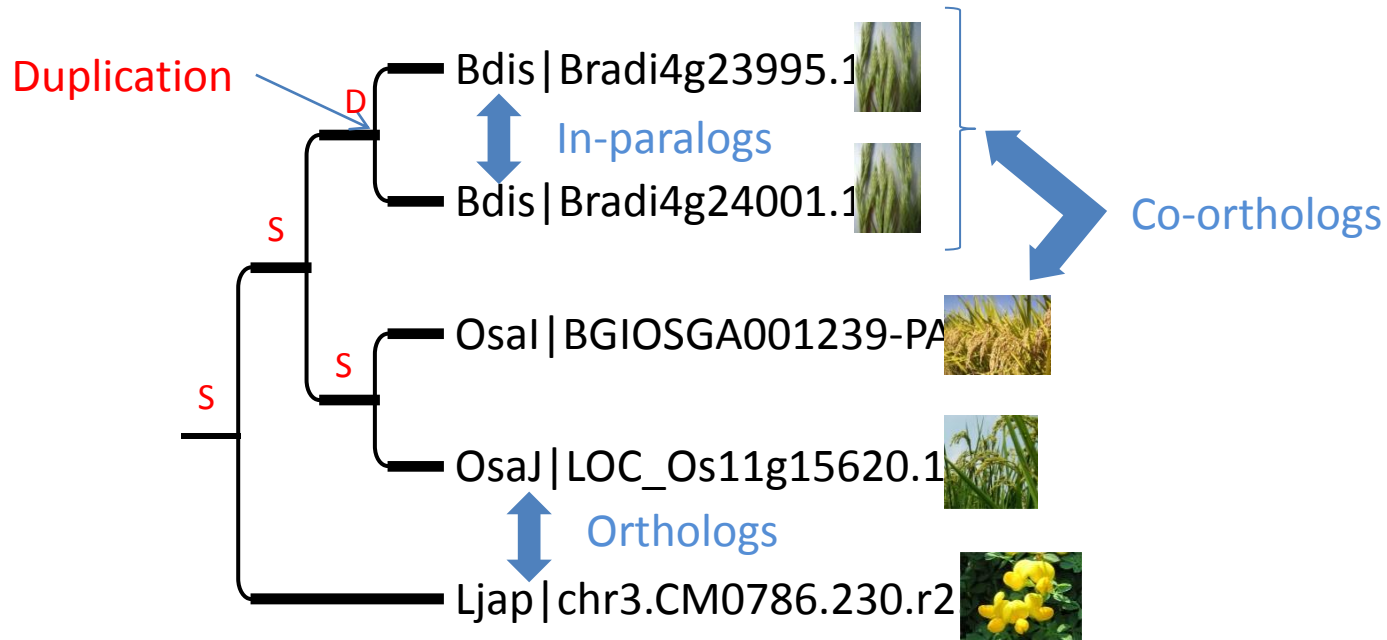
Traduction en arbre:



Species	Abbrev
<i>Capsicum annuum C</i>	CanC
<i>Capsicum annuum Z</i>	CanZ
<i>Dianthus caryophyllus</i>	Dcar
<i>Nicotiana benthamiana</i>	Nben
<i>Solanum lycopersicum</i>	Slyc
<i>Solanum tuberosum</i>	Stub

GRP27024: Bdis | Bradi4g23995.1 Bdis | Bradi4g24001.1 Osa | BGIOGA001239-PA
Ljap | chr3.CM0786.230.r2 OsaJ | LOC_Os11g15620.1

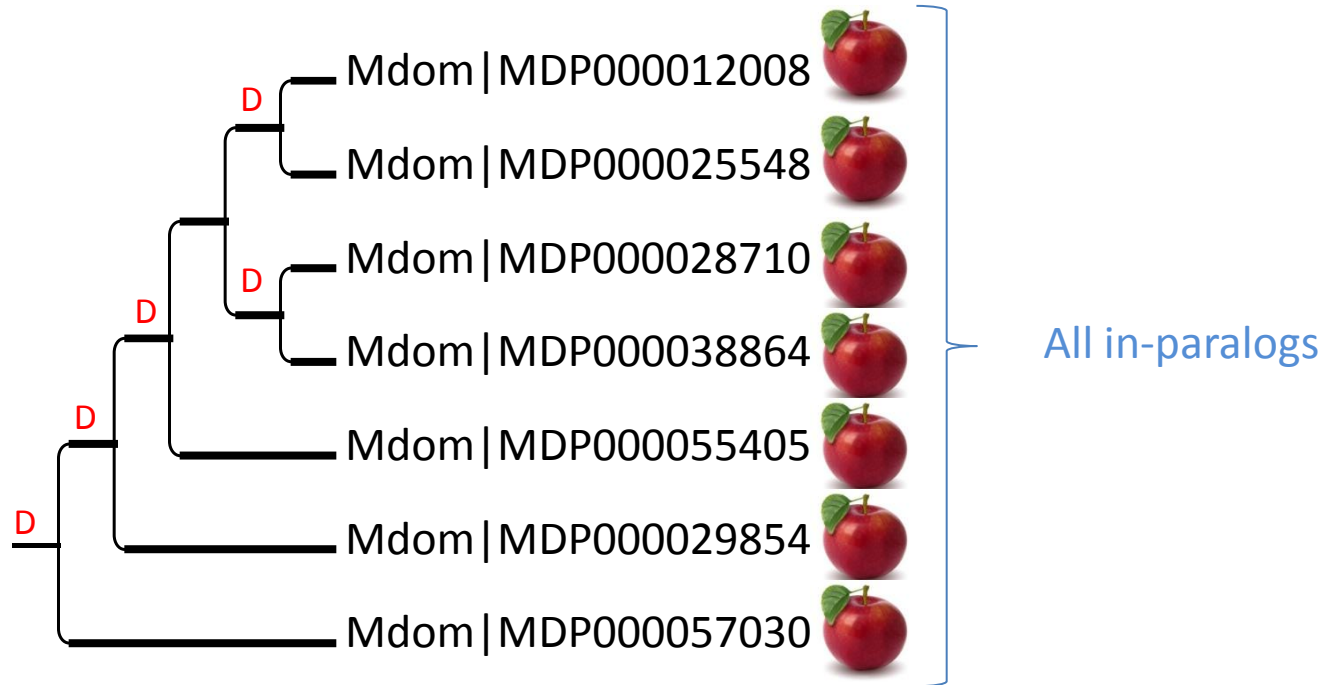
Traduction en arbre:



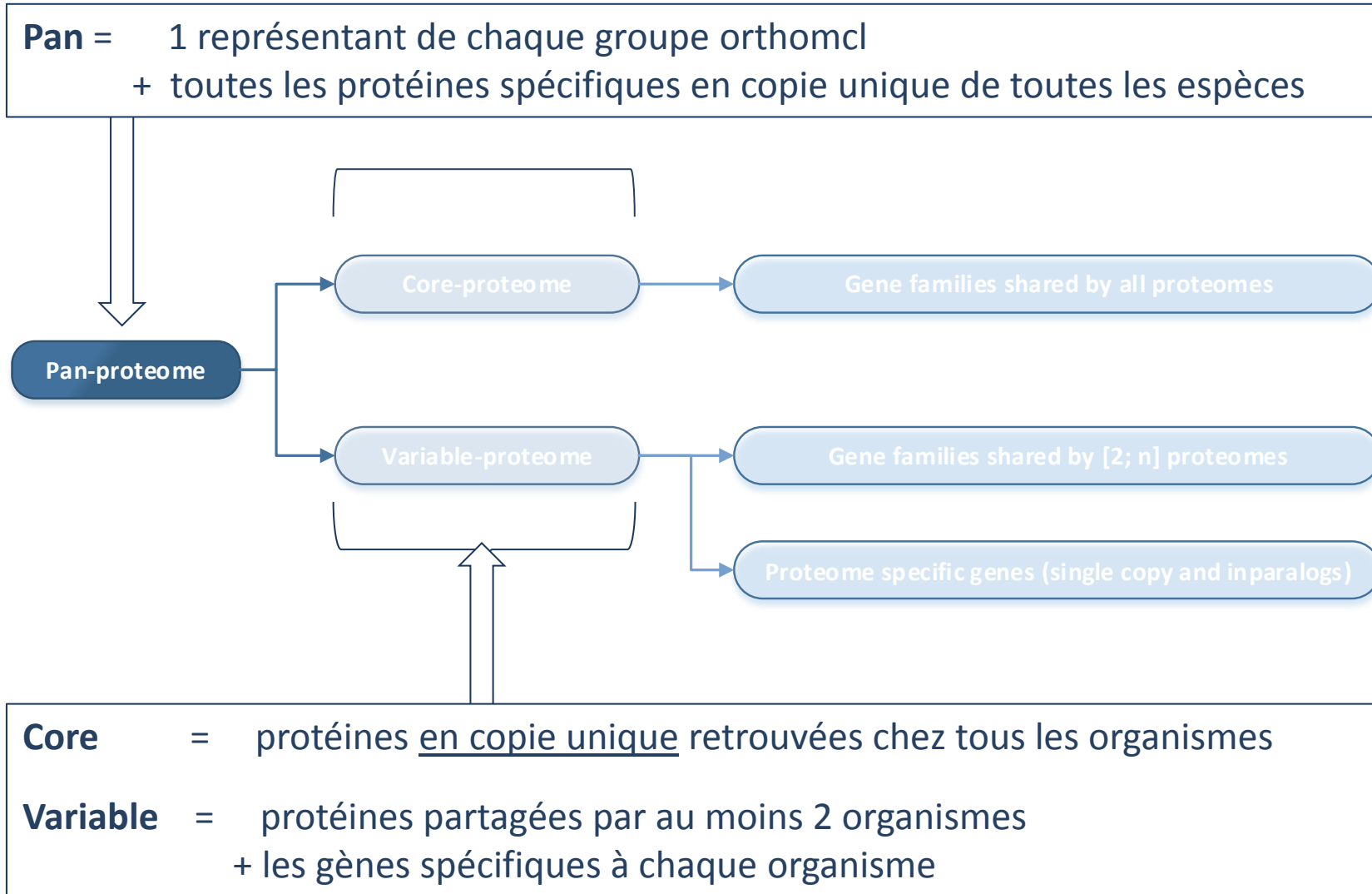
Species	Abbrev
<i>Brachypodium distachyon</i>	Bdis
<i>Lotus japonicus</i>	Ljap
<i>Oryza sativa I</i>	Osa
<i>Oryza sativa J</i>	OsaJ

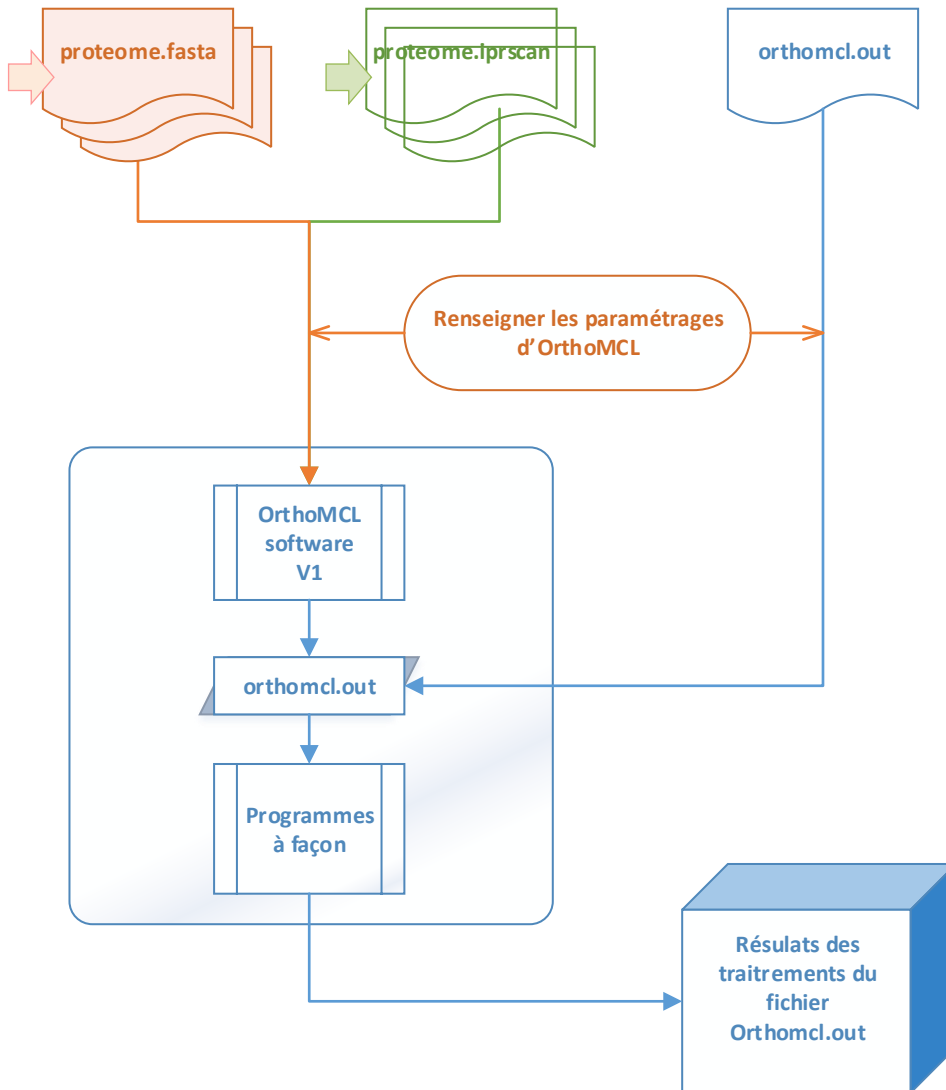
GRP22264: Mdom | MDP0000120086 Mdom | MDP0000255483 Mdom | MDP0000287109
 Mdom | MDP0000388642 Mdom | MDP0000554052 Mdom | MDP0000298540
 Mdom | MDP0000570301

Traduction en arbre:



Species	Abbrev
<i>Malus domestica</i>	Mdom





➔ **Données minimales :**
les fichiers de séquences des protéomes

➔ **Valeur ajoutée :**
Possibilité de fournir les fichiers d'annotation au format InterProScan

L'outil accepte en entrée un fichier orthomcl.out version : 1 ou 2

Nomenclature obligatoire :

filename . extension

extension

.fa .faa .fas .fasta

.iprscan .ipr

filename = **code court et concis**

Il sert d'identifiant de votre protéome

proteome.fasta

proteome.iprscan

orthomcl.out

Obligatoirement: les mêmes identifiants (IDs) des séquences doivent être utilisés dans les 3 types de fichier pouvant être fournis en entrée

CODE = strm5

strm5.fasta

strm5.iprscan

```
1 >Sma1_0001
2 MDAWSRSLERLEAEFPDPDVHTLWKLPLQADLRVDSLVLVYAPNAFVDOQVR
3 ELYLARIRELLAHFAGFSDVLEI GSRPRPVEAQNPASTPSAHVSSEPO
4 VPFAGLNDNHYTFANFVEGSRNLGLAAFAQAKPKGDRAHNPLLLYGGT
5 GLGKTHLMFAAGNAMRQANPGAKVLVLRSEOFFSAMIRALQKTDHQFKR
6 FQOQVDALLDDIOFFAAGKDRTOQEFFFHTFNALFDGKQOILLTCDRYPRE
7 VEGLEEARLKSRALWGLSVAIQEPDFETRAI VLAKAPERGAEPDDVAFL
8 IAKKMRSNRVDELEGALNLTARANPFTGRAITTEFAQETLRDLLRAQQOAI
9 SIPNIQKTVADYVGLQIKDLLSKRRTRSLARPRQVAMALTKEHTEHSLPE
10 IGDAFAGRDRHTTVLHACRQIRTLMEADGKLRWDKILRKLSE
11 >Sma1_0002
12 MRFTLQREAFKPLAQVNVVRRQTLPLVLANFLVQVQNGQLSLTGTDL
13 VEMYSRIA VEDAQDGETTIPARKLFEIIRALPDGSRITVYSGDKITVQA
14 GRSRFTLATLPSNDPFSVDEVEATERVAI GEATLKELIERTAFAMAQQDV
15 RYLLNGLLFDLRGDALRT VATDGHRLALCF
16 TELQRLLLESQDREIELEVGRSHVRVKKDDI
17 GADREKVNRESLRASLQRAAILSNEKYR
18 EAQEEIEADTTVSDLAIGFNWNYLLDALS
19 ESSSEKSRHVVMPLRL
20 >Sma1_0003
21 MQIRRLALHQLRRFNAVELSPQGLNLLTI
22 SFRGRVDRGLVROGQEALEIFVWEQORA
23 LDGEDVAQLGNLCAALAVVTFEPGSHALY
24 FLSLWRRYSRALKORNALLKOGGSPRLMD
25 ERLQRTVALAELAPQLGIQAMELSPGW
26 GYTSVGPFRADWSVDFHNI PRGDALSRGQ
27 EWPVIALDDLAASELDRTHQARVLERLLGG
28 ARFHVEHAQI VAVP
29 >Sma1_0004
30 MSDFONTPNANNMYDANSTTALFGLFAVR
```

1	Sma1_0002	1C24AD072989A2A8	366	HMPFfam	PF02768	DNA_pol3_beta_3	246	365	2.90000000000000062E-30	T
2	Sma1_0002	1C24AD072989A2A8	366	HMPFfam	PF00712	DNA_pol3_beta_1	1	118	3.5000000000000003E-35	T
3	Sma1_0002	1C24AD072989A2A8	366	HMPFfam	PF02767	DNA_pol3_beta_2	130	243	1.099999999999999873E-34	T
4	Sma1_0002	1C24AD072989A2A8	366	Gene3D	G3DSA:3.10.150.10	DNA_polIII_beta_1	1	121	11.200000000000000	C
5	Sma1_0002	1C24AD072989A2A8	366	Gene3D	G3DSA:3.10.150.10	DNA_polIII_beta_124	248	81.600000000000000	C	
6	Sma1_0002	1C24AD072989A2A8	366	Gene3D	G3DSA:3.10.150.10	DNA_polIII_beta_249	366	21.199999999999999	C	
7	Sma1_0002	1C24AD072989A2A8	366	HMMTigr	TIGR00663	dnan	1	366	1.5999883531364446E-107	T
8	Sma1_0002	1C24AD072989A2A8	366	superfamily	SSF55979	SSF55979	1	116	51.499993387561E	C
9	Sma1_0002	1C24AD072989A2A8	366	superfamily	SSF55979	SSF55979	125	244	11.3999989204987	C
10	Sma1_0002	1C24AD072989A2A8	366	superfamily	SSF55979	SSF55979	245	366	81.5000302480906	C
11	Sma1_0002	1C24AD072989A2A8	366	HMMSmart	SM00480	POL3Bc	17	362	0.0	T
12	Sma1_0005	F069F515A31EA231	281	HMPFfam	PF02517	Abi_175	273	1.99999999999999946E-17	T	
13	Sma1_0006	228F1D6CEF7F9751	268	HMPFfam	PF01435	Peptidase_M48	33	243	4.5999999999999999E-41	T
14	Sma1_0006	228F1D6CEF7F9751	268	HMPPanther	PTHR22726	PTHR22726	20	258	2.100026783402	C

Préconisation

utilisation des mêmes fichiers fasta pour lancer InterProScan et OrthoMCL

Texte brut v Largeur des tabulations

```
1357.ORTHOMCL1356(13 genes,13 taxa):
1358.ORTHOMCL1357(13 genes,13 taxa):
1359.ORTHOMCL1358(13 genes,13 taxa):
1360.ORTHOMCL1359(13 genes,13 taxa):
1361.ORTHOMCL1360(13 genes,13 taxa):
1362.ORTHOMCL1361(13 genes,13 taxa):
1363.ORTHOMCL1362(13 genes,13 taxa):
1364.ORTHOMCL1363(13 genes,13 taxa):
1365.ORTHOMCL1364(13 genes,13 taxa):
1366.ORTHOMCL1365(12 genes,1 taxa):
1367.ORTHOMCL1366(12 genes,1 taxa):
1368.ORTHOMCL1367(12 genes,8 taxa):
1369.ORTHOMCL1368(12 genes,10 taxa):
1370.ORTHOMCL1369(12 genes,9 taxa):
1371.ORTHOMCL1370(12 genes,11 taxa):
1372.ORTHOMCL1371(12 genes,12 taxa):
1373.ORTHOMCL1372(12 genes,12 taxa):
1374.ORTHOMCL1373(12 genes,12 taxa):
1375.ORTHOMCL1374(12 genes,12 taxa):
1376.ORTHOMCL1375(12 genes,12 taxa):
1377.ORTHOMCL1376(12 genes,12 taxa):
1378.ORTHOMCL1377(12 genes,12 taxa):
```

PD0005(xy ft)	PX0_03480(xanor)	Sma1_0004(st rm5)	XAC0004(xanac)	XALC_0004(xanal)	XAUC_15250(xanfu)	XCC0004(xanccp)	XCV0004(xanc5)	XFF4834R_chr00040(XFF4834R)	XGA_3527(xanga)	X
PD0003(xy ft)	PX0_03481(xanor)	Sma1_0003(st rm5)	XAC0003(xanac)	XALC_0003(xanal)	XAUC_15240(xanfu)	XCC0003(xanccp)	XCV0003(xanc5)	XFF4834R_chr00030(XFF4834R)	XGA_3526(xanga)	X
PD0002(xy ft)	PX0_03482(xanor)	Sma1_0002(st rm5)	XAC0002(xanac)	XALC_0002(xanal)	XAUC_11630(xanfu)	XCC0002(xanccp)	XCV0002(xanc5)	XFF4834R_chr00020(XFF4834R)	XGA_3525(xanga)	X
PD0001(xy ft)	PX0_03483(xanor)	Sma1_0001(st rm5)	XAC0001(xanac)	XALC_0001(xanal)	XAUC_11620(xanfu)	XCC0001(xanccp)	XCV0001(xanc5)	XFF4834R_chr00010(XFF4834R)	XGA_3524(xanga)	X

https://bbric-pipelines.toulouse.inra.fr/orthomcl-companion/web/index.html

Protein family analyses
OrthoMCL-Companion

login

Home

Perform OrthoMCL analyses, browse, visualise and share the results!

Nowadays, comparative genomics is a classical tool for biologists to address various scientific questions such as pathogeny determinism, bacteria host specificity, duplication and gene family expansion in a species. The large number of available finished or draft genome sequences associated with the full automatization of structural annotation pipelines leads to an era of high throughput comparative genomics. Tools like **OrthoMCL** provide an easy way to compare tens of species at a time and build homologous gene families. However, output results, classically a matrix with groups as lines and genes as columns, remains too raw to be interpreted by an end user without programming skills. In order to help users in the post process of these data, we developed a web user interface that automatically extract frequently asked datasets such as core and pan proteomes, specific proteins, abundance and presence/absence matrix (aka. phylogenetic profiles). Providing optional InterPro annotation will also lead to build tables and charts representing terms (eg : PFAM, GO) count in each previously described datasets. Our tool uses standard formats as input and outputs to ease post-process analyses.

Test Dataset

[3 proteomes of about 400 proteins each annotated with InterPro 51.0](#)

Results

[Test dataset analyses results \(OrthoMCL with high coverage parameter - pvmatch cut off = 80\)](#)

Analyse raw proteomes

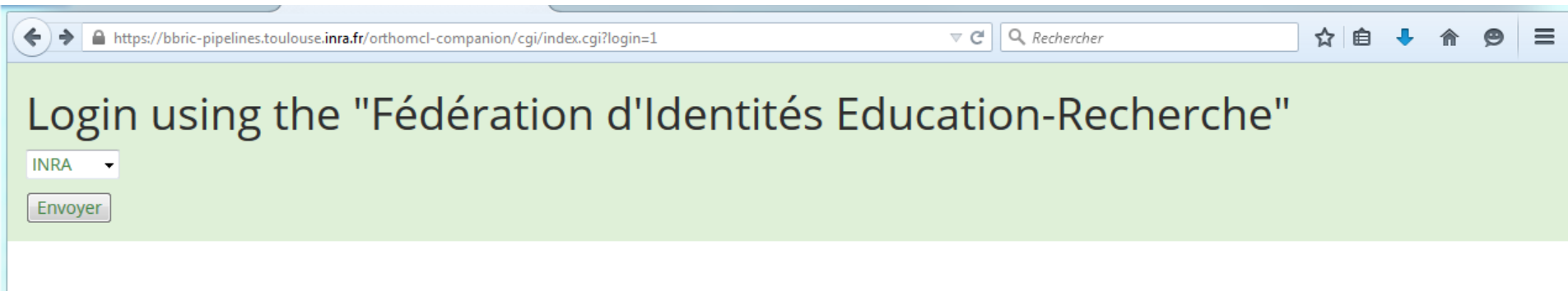
Analyse annotated proteomes

Browse

Share

Delete

Fast launch



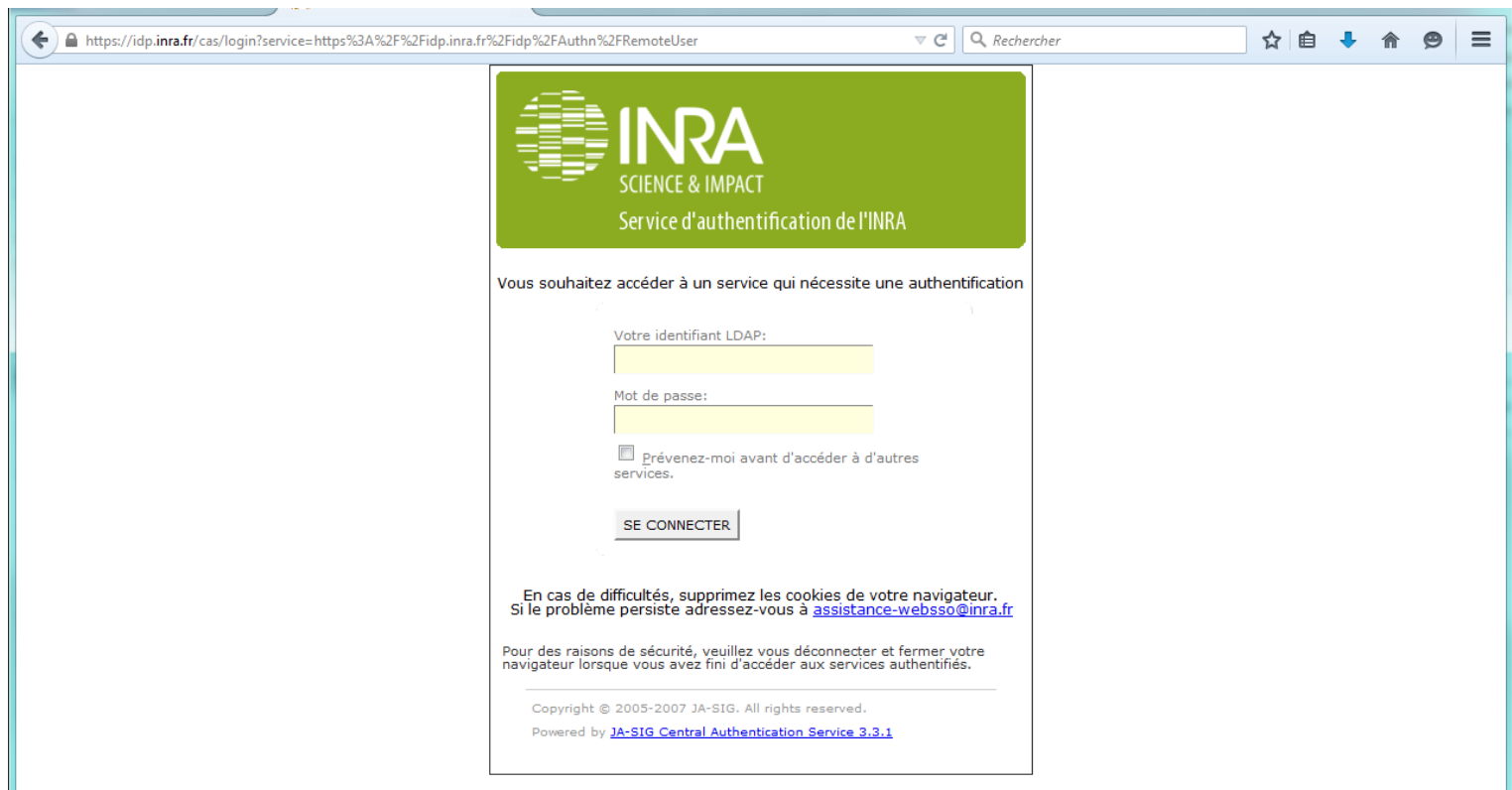
https://bbric-pipelines.toulouse.inra.fr/orthomcl-companion/cgi/index.cgi?login=1

Rechercher

Login using the "Fédération d'Identités Education-Recherche"


INRA

Envoyer



https://idp.inra.fr/cas/login?service=https%3A%2F%2Fidp.inra.fr%2Fidp%2FAuthn%2FRemoteUser

Rechercher



INRA
SCIENCE & IMPACT
Service d'authentification de l'INRA

Vous souhaitez accéder à un service qui nécessite une authentification

Votre identifiant LDAP:

Mot de passe:

Prévenez-moi avant d'accéder à d'autres services.

SE CONNECTER

En cas de difficultés, supprimez les cookies de votre navigateur.
Si le problème persiste adressez-vous à assistance-websso@inra.fr

Pour des raisons de sécurité, veuillez vous déconnecter et fermer votre navigateur lorsque vous avez fini d'accéder aux services authentifiés.

Copyright © 2005-2007 JA-SIG. All rights reserved.
Powered by [JA-SIG Central Authentication Service 3.3.1](#)

<https://bbric-pipelines.toulouse.inra.fr/orthomcl-companion/web/index.html>

Protein family analyses

OrthoMCL-Companion

Home New analysis

Perform OrthoMCL analyses, browse, visualise and share the results!

Nowadays, comparative genomics is a classical tool for biologists to address various scientific questions such as pathogeny determinism, bacteria host specificity, duplication and gene family expansion in a species. The large number of available finished or draft genome sequences associated with the full automatism of structural annotation pipelines leads to an era of high throughput comparative genomics. Tools like **OrthoMCL** provide an easy way to compare tens of species at a time and build homologous gene families. However, output results, classically a matrix with groups as lines and genes as columns, remains too raw to be interpreted by an end user without programming skills. In order to help users in the post process of these data, we developed a web user interface that automatically extract frequently asked datasets such as core and pan proteomes, specific proteins, abundance and presence/absence matrix (aka. phylogenetic profiles). Providing optional InterPro annotation will also lead to build tables and charts representing terms (eg : PFAM, GO) count in each previously described datasets. Our tool uses standard formats as input and outputs to ease post-process analyses.

Test Dataset
[3 proteomes of about 400 proteins each annotated with InterPro 51.0](#)

Results
[Test dataset analyses results \(OrthoMCL with high coverage parameter - pmatch cut off = 80\)](#)

Analyse raw proteomes

Analyse annotated proteomes

Browse

Share

Delete

Fast launch



Protein family analyses

OrthoMCL-Companion

[Home](#) [New analysis](#) [List of analyses](#)

[Launch analysis](#) [Reset](#)

Input data zip file

Select a zip file:

File name: No file defined

Analysis description

Title: *

Description:

General parameters

Add IprScan results?

OrthoMcl

OrthoMcl result file

Select an OrthoMCL file:

File name: No file defined

OR

Parameters

pv_cutoff: *

pi_cutoff: *

pmatch_cutoff: *

inflation: *

Proteomes

Proteome 1

Proteome description

Code: *

Proteome description:

Is reference?:

Proteome file

Select a fasta file: *

File name: No file defined

IprScan file

Select an IprScan file: *

File name: No file defined

[Add Proteome](#)

[Launch analysis](#) [Reset](#)

Home New analysis List of analyses

Launch analysis Reset

Input data zip file

Select a zip file:

File name: No file defined

Analysis description

Title: *

Description:

General parameters

Add IprScan results?

OrthoMcl

OrthoMcl result file

Select an OrthoMCL file:

File name: No file defined

OR Parameters

pv_cutoff: *

pi_cutoff: *

pmatch_cutoff: *

inflation: *

Proteomes

Proteome 1

Proteome description

Code: *

Proteome description:

Is reference?:

Proteome file

Select a fasta file: *

File name:

IprScan file

Select an IprScan file: *

File name: No file defined

Add Proteome

Launch analysis Reset

Home New analysis List of analyses

Launch analysis Reset

Input data zip file

Select a zip file:

File name: No file defined

Analysis description

Title: *

Description:

General parameters

Add IprScan results?

OrthoMcl

OrthoMcl result file

Select an OrthoMCL file:

File name: No file defined

OR Parameters

pv_cutoff: *

pi_cutoff: *

pmatch_cutoff: *

inflation: *

Proteomes

Proteome 1

Proteome description

Code: *

Proteome description:

Is reference?:

Proteome file

Select a fasta file: *

File name:

IprScan file

Select an IprScan file: *

File name: strm5.iprscan

Add Proteome

Launch analysis Reset

Home New analysis List of analyses

Launch analysis Reset

Input data zip file

Select a zip file:

File name: No file defined

Analysis description

Title: *

Description:

General parameters

Add IprScan results?

OrthoMcl

OrthoMcl result file

Select an OrthoMCL file:

File name: No file defined

OR Parameters

pv_cutoff: *

pi_cutoff: *

pmatch_cutoff: *

inflation: *

Proteomes

Proteome 1

Proteome description

Code: *

Proteome description:

Is reference?:



Proteome file

Select a fasta file: *

File name:

IprScan file

Select an IprScan file: *

File name:

Add Proteome

Launch analysis Reset



[Home](#) | [New analysis](#) | [List of analyses](#)

[Launch analysis](#) | [Reset](#)

Input data zip file

Select a zip file:

File name: No file defined

Analysis description

Title:

Description:

General parameters

Add IprScan results ?

OrthoMcl

OrthoMcl result file

Select an OrthoMCL file:

File name: orthomcl.out

OR

Parameters

pv_cutoff: *

pi_cutoff: *

pmatch_cutoff: *

inflation: *

Proteomes

- Proteome 1
- Proteome 2
- Proteome 3
- Proteome 4

Proteome description

Code: *

Proteome description:

Is reference ?

Proteome file

Select a fasta file: *

File name: xancb.fasta

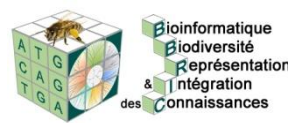
IprScan file

Select an IprScan file: *

File name: xancb.iprscan

[Add Proteome](#)

[Launch analysis](#) | [Reset](#)





Home | **New analysis** | List of analyses

Launch analysis | Reset

Input data zip file

Select a zip file:

File name: No file defined

Analysis description

Title: *

Description:


General parameters

Add IprScan results ?

OrthoMcl

OrthoMcl result file

Select an OrthoMCL file:

File name: orthomcl.out 

OR Parameters

pv_cutoff: *

pi_cutoff: *

pmatch_cutoff: *

inflation: *

Proteomes

- Proteome 1
- Proteome 2
- Proteome 3
- Proteome 4

Proteome 4

Proteome description

Code: *

Proteome description:

Is reference ?

Add Proteome

Proteome file

Select a fasta file: *

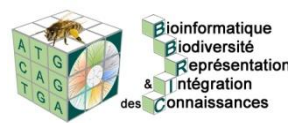
File name: xancb.fasta

IprScan file

Select an IprScan file: *

File name: xancb.iprscan

Launch analysis | Reset





Home | **New analysis** | List of analyses

Launch analysis | Reset

Input data zip file

Select a zip file:

File name: No file defined

Analysis description

Title: *

Description:

General parameters

Add IprScan results?

OrthoMcl

OrthoMcl result file

Select an OrthoMCL file:

File name: **orthomcl.out**

OR

Parameters

pv_cutoff: *

pi_cutoff: *

pmatch_cutoff: *

inflation: *

OBLIGATOIRE

Proteomes

Proteome 1

Proteome 2

Proteome 3

Proteome 4

Proteome description

Code: *

Proteome description:

Is reference?:

Proteome file

Select a fasta file: *

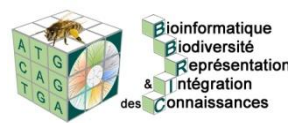
File name: xancb.fasta

IprScan file

Select an IprScan file: *

File name: xancb.iprscan

Launch analysis | Reset





Home | **New analysis** | List of analyses

Launch analysis | Reset

Input data zip file

Select a zip file:

File name: No file defined

Analysis description

Title: *

Description:

General parameters

Add IprScan results ?

OrthoMcl

OrthoMcl result file

Select an OrthoMCL file:

File name: orthomcl.out

OR Parameters

pv_cutoff: *

pi_cutoff: *

pmatch_cutoff: *

inflation: *

Proteomes

- Proteome 1
- Proteome 2
- Proteome 3
- Proteome 4

Proteome 4

Proteome description

Code: *

Proteome description:

Proteome file

Select a fasta file: *

File name: xancb.fasta

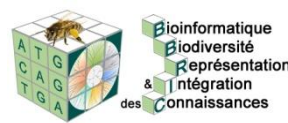
IprScan file

Select an IprScan file: *

File name: xancb.iprscan

Add Proteome

Launch analysis | Reset



Success



analysis launched
Please regularly press the reload button in the panel 'List Of Analyses' to follow the progress of your analyses

OK



Protein family analyses

OrthoMCL-Companion

Home New analysis **List of analyses**

Reload

status	public	actions	url	title	date	log
	<input type="checkbox"/>		zip	90Krx5FDVvk	20151009 13:40:38	INFO - Orthomcl version 1.4 (blast+) CMD - (cd /mnt/scrat



Home New analysis **List of analyses**

Reload

status	public	actions	url	title	date	log
	<input type="checkbox"/>		zip	90Krx5FDVvk	20151009 13:40:38	Finished : 20151009 15:41:22



Protein family analyses

OrthoMCL-Companion

Home | **New analysis** | List of analyses | 90Krx5FDVk

Global parameters

Contact Ludovic.Cottret@toulouse.inra.fr
Date 20151009 15:40:44

OrthoMcl parameters

Outfile [orthomcl_output/orthomcl.out](#)
Version 1.4 (blast+)
Parameters --pi_cutoff=0 --pv_cutoff=0.00001 --pmatch_cutoff=80 --inflation=1.5

Proteomes

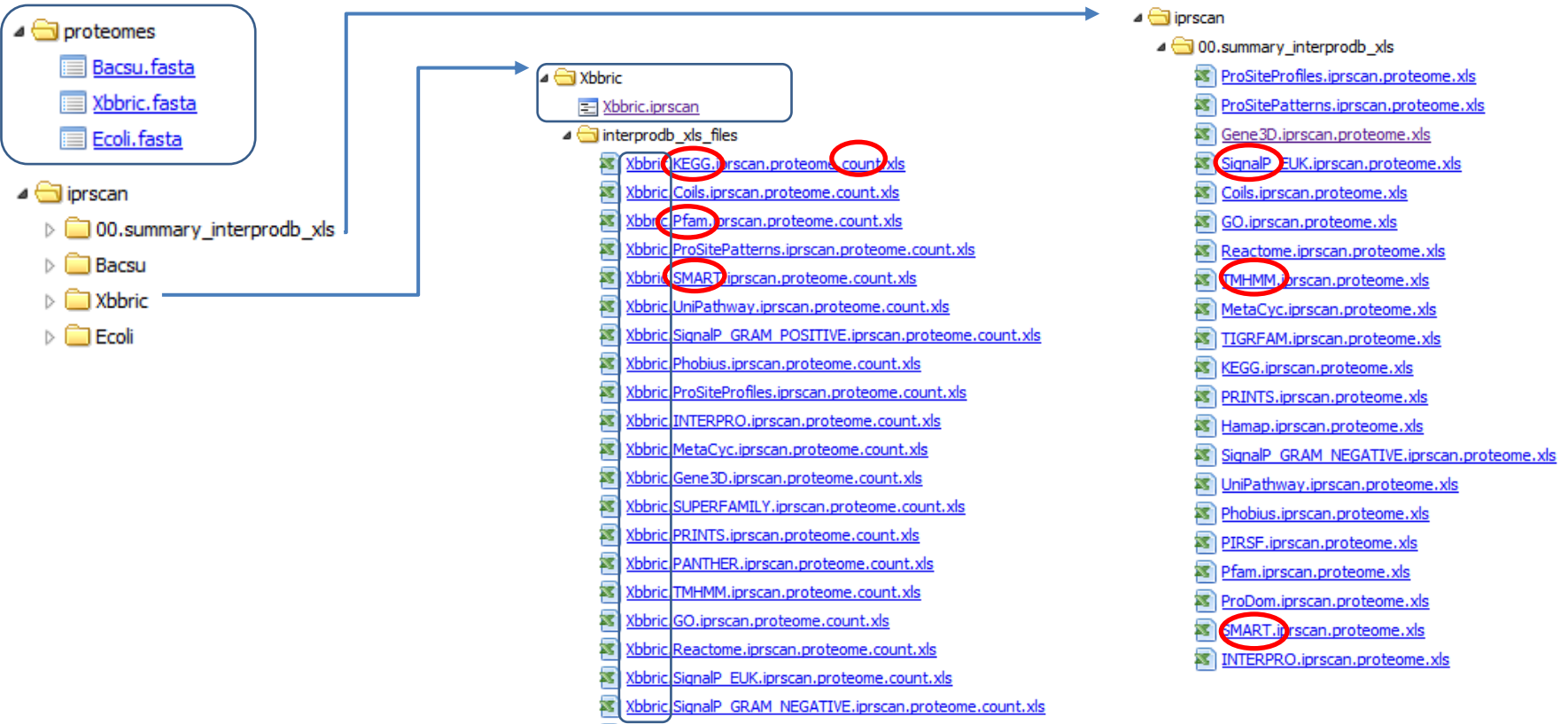
code	reference	title
Xbbric	false	X. bbric str. mini
Ecoli	false	E. coli str. mini
Bacsu	false	B. subtilis str. mini

Results

Status	Name	Description
	input_data	Input data
	orthomcl_output	OrthoMCL Group description and annotation
	summary	Analysis summary - Global counts
	specific_proteins	Extraction of Specific protein sequences and annotation
	orthologous_proteins	Extraction of "orthologous" protein sequences and annotation
	core_proteome	Core proteome pseudomolecule
	venn_diagrams	Lists of all groups with a protein of this proteome + unique protein by proteome
	pan_proteome	Directory with fasta multi-sequence of pan-proteome
	statistics_matrices	Descriptive Statistics: number of OrthoMCL group - Number of proteins by OrthoMCL group - Number of taxa by OrthoMCL gro...

Status	Name	Description
	input_data	Input data
	summary	Analysis summary - Global counts
	specific_proteins	Extraction of Specific protein sequences and annotation
	orthologous_proteins	Extraction of "orthologous" protein sequences and annotation
	core_proteome	Core proteome pseudomolecule
	venn_diagrams	Lists of all groups with a protein of this proteome + unique protein by proteome
	pan_proteome	Directory with fasta multi-sequence of pan-proteome
	statistics_matrices	Descriptive Statistics: number of OrthOMCL group - Number of proteins by OrthOMCL group - Number of taxa by OrthOMCL group; Phylogenetics matrices : Presence/Absence matrix - ...

input_data



Status	Name	Description
	input_data	Input data
	summary	Analysis summary - Global counts
	specific_proteins	Extraction of Specific protein sequences and annotation
	orthologous_proteins	Extraction of "orthologous" protein sequences and annotation
	core_proteome	Core proteome pseudomolecule
	venn_diagrams	Lists of all groups with a protein of this proteome + unique protein by proteome
	pan_proteome	Directory with fasta multi-sequence of pan-proteome
	statistics_matrices	Descriptive Statistics: number of OrthOMCL group - Number of proteins by OrthOMCL group - Number of taxa by OrthOMCL group; Phylogenetics matrices : Presence/Absence matrix - ...

input_data

proteomes

- Bacsu.fasta
- Xbbric.fasta
- Ecoli.fasta

iprscan

- 00.summary_interprodb.xls
- Bacsu
- Xbbric
- Ecoli

BG13696	Bacsu	PF07733	Bacterial DNA polymerase III alpha subunit
BG13696	Bacsu	PF02811	PHP domain
BG13696	Bacsu	PF14579	Helix-hairpin-helix motif
b0532	Ecoli	PF13953	PapC C-terminal domain
b4210	Ecoli	PF00892	EamA-like transporter family
b4293	Ecoli	PF04542	Sigma-70 region 2
b4293	Ecoli	PF08281	Sigma-70, region 4
b0604	Ecoli	PF13098	Thioredoxin-like domain
Xbbric.3982	Xbbric	PF05598	Transposase domain (DUF772)
Xbbric.4079	Xbbric	PF00535	Glycosyl transferase family 2
Xbbric.3926	Xbbric	PF13659	Methyltransferase domain
Xbbric.4259	Xbbric	PF01551	Peptidase family M23

#TERM	DESCRIPTION	HITS	RATIO
PF00106	short chain dehydrogenase	6	0.0159574468085106
PF01609	Transposase DDE domain	5	0.0132978723404255
PF00005	ABC transporter	5	0.0132978723404255
PF00126	Bacterial regulatory helix-turn-helix protein, lysR family	4	0.0106382978723404
PF00072	Response regulator receiver domain	4	0.0106382978723404
PF04542	Sigma-70 region 2	4	0.0106382978723404
PF13276	HTH-like domain	4	0.0106382978723404
PF00535	Glycosyl transferase family 2	4	0.0106382978723404

Status	Name	Description
▶	input_data	Input data
▶	summary	Analysis summary - Global counts
▶	specific_proteins	Extraction of Specific protein sequences and annotation
▶	orthologous_proteins	Extraction of "orthologous" protein sequences and annotation
▶	core_proteome	Core proteome pseudomolecule
▶	venn_diagrams	Lists of all groups with a protein of this proteome + unique protein by proteome
▶	pan_proteome	Directory with fasta multi-sequence of pan-proteome
▶	statistics_matrices	Descriptive Statistics: number of OrthOMCL group - Number of proteins by OrthOMCL group - Number of taxa by OrthOMCL group; Phylogenetics matrices : Presence/Absence matrix - ...

summary

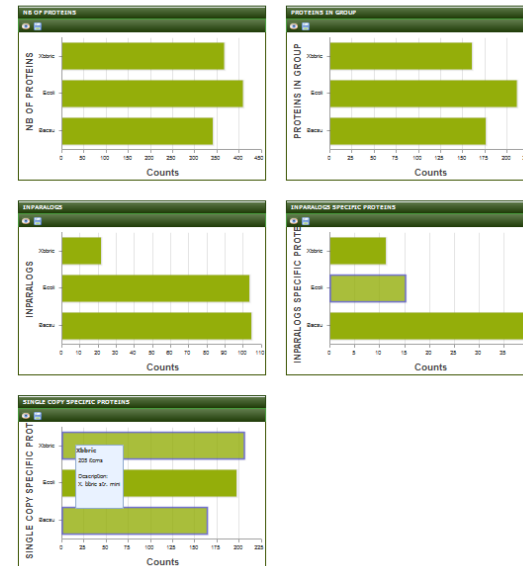
summary

- summary.xls
- summary.count.html
- summary.percentage.html

#CODE	TITLE	NB OF PROTEINS	PROTEINS IN GROUP	INPARALOGS	INPARALOGS SPECIFIC PROTEINS	SINGLE COPY SPECIFIC PROTEINS
Bacsu	B. subtilis str. mini	338	175	104	40	163
Ecoli	E. coli str. mini	406	210	103	15	196
Xbbri	X. bbric str. mini	364	159	21	11	205

Percentage of every kind of datasets for all species

Count of every kind of datasets for all species



Status	Name	Description
	input_data	Input data
	summary	Analysis summary - Global counts
	specific_proteins	Extraction of Specific protein sequences and annotation
	orthologous_proteins	Extraction of "orthologous" protein sequences and annotation
	core_proteome	Core proteome pseudomolecule
	venn_diagrams	Lists of all groups with a protein of this proteome + unique protein by proteome
	pan_proteome	Directory with fasta multi-sequence of pan-proteome
	statistics_matrices	Descriptive Statistics: number of OrthOMCL group - Number of proteins by OrthOMCL group - Number of taxa by OrthOMCL group; Phylogenetics matrices : Presence/Absence matrix - ...

specific-proteins

specific_proteins

- 00.summary_interprodb_xls
- Bacsu
- Xbbric
- Ecoli
- specific_proteins
 - 00.summary_interprodb_xls
 - KEGG.iprscan.specific.singlecopy.xls
 - TIGRFAM.iprscan.specific.inparalogs.xls
 - INTERPRO.iprscan.specific.inparalogs.xls
 - SMART.iprscan.specific.singlecopy.xls
 - ProSiteProfiles.iprscan.specific.singlecopy.xls
 - SignalP_GRAM_POSITIVE.iprscan.specific.inparalogs.xls
 - Gene3D.iprscan.specific.singlecopy.xls
 - SignalPEUK.iprscan.specific.singlecopy.xls
 - Reactome.iprscan.specific.inparalogs.xls
 - SignalP_GRAM_NEGATIVE.iprscan.specific.inparalogs.xls
 - TMHMM.iprscan.specific.inparalogs.xls
 - GO.iprscan.specific.singlecopy.xls

Bacsu

- Bacsu.inparalogs.quality.xls
- Bacsu.iprscan.specific.inparalogs.count.html
- Bacsu.iprscan.specific.singlecopy.count.html
- Bacsu.singlecopy.quality.xls
- fasta_files
 - Bacsu.inparalogs.fasta
 - Bacsu.singlecopy.fasta
- interprodb_xls_files
 - Bacsu.INTERPRO.iprscan.specific.singlecopy.count.xls
 - Bacsu.SUPERFAMILY.iprscan.specific.inparalogs.count.xls
 - Bacsu.PIRSF.iprscan.specific.singlecopy.count.xls
 - Bacsu.GO.iprscan.specific.singlecopy.count.xls
 - Bacsu.Reactome.iprscan.specific.inparalogs.count.xls
 - Bacsu.ProSiteProfiles.iprscan.specific.inparalogs.count.xls
 - Bacsu.PIRSF.iprscan.specific.inparalogs.count.xls
 - Bacsu.GO.iprscan.specific.inparalogs.count.xls
 - Bacsu.Pfam.iprscan.specific.singlecopy.count.xls
 - Bacsu.TMHMM.iprscan.specific.singlecopy.count.xls

#seqid	groupid	length	%of undetermined AA (X or *)
BG12077	ORTHOMCL157	308	0
BG11300	ORTHOMCL159	254	0
BG12932	ORTHOMCL70	513	0
BG12005	ORTHOMCL159	253	0
BG11346	ORTHOMCL168	641	0
BG11867	ORTHOMCL74	358	0
BG13442	ORTHOMCL158	363	0
BG11392	ORTHOMCL70	544	0

Bacsu.inparalogs.quality

Protéines d'un groupe orthomcl mono-espèce

#seqid	groupid	length	%of undetermined AA (X or *)
BG12780	null	73	0
BG11314	null	184	0
BG10712	null	181	0
BG13114	null	154	0
BG13882	null	80	0
BG10239	null	106	0
BG12369	null	438	0
BG12620	null	157	0

Bacsu.singlecopy.quality

Protéines spécifiques à une espèce et en une seule copie

Status	Name	Description
	input_data	Input data
	summary	Analysis summary - Global counts
	specific_proteins	Extraction of Specific protein sequences and annotation
	orthologous_proteins	Extraction of "orthologous" protein sequences and annotation
	core_proteome	Core proteome pseudomolecule
	venn_diagrams	Lists of all groups with a protein of this proteome + unique protein by proteome
	pan_proteome	Directory with fasta multi-sequence of pan-proteome
	statistics_matrices	Descriptive Statistics: number of OrthOMCL group - Number of proteins by OrthOMCL group - Number of taxa by OrthOMCL group; Phylogenetics matrices : Presence/Absence matrix - ...

Orthologous_proteins

orthologous_proteins

- 00.summary_interprodb_xls
- Bacsu
- Xbbric
- Ecoli

orthologous_proteins

- 00.summary_interprodb_xls
 - SUPERFAMILY.iprscan.specific.orthologous.xls
 - SignalP_GRAM_POSITIVE.iprscan.specific.orthologous.xls
 - ProSiteProfiles.iprscan.specific.orthologous.xls
 - PIRSF.iprscan.specific.orthologous.xls
 - SMART.iprscan.specific.orthologous.xls
 - TIGRFAM.iprscan.specific.orthologous.xls
 - MetaCyc.iprscan.specific.orthologous.xls
 - UniPathway.iprscan.specific.orthologous.xls
 - INTERPRO.iprscan.specific.orthologous.xls
 - PRINTS.iprscan.specific.orthologous.xls
 - Phobius.iprscan.specific.orthologous.xls
 - Hamap.iprscan.specific.orthologous.xls
 - GO.iprscan.specific.orthologous.xls
 - ProDom.iprscan.specific.orthologous.xls
 - KEGG.iprscan.specific.orthologous.xls

Ecoli

- Ecoli.orthologous.quality.xls
- Ecoli.iprscan.specific.orthologous.count.html
- fasta_files
 - Ecoli.orthologous.fasta
- interprodb_xls_files
 - Ecoli.MetaCyc.iprscan.specific.orthologous.count.xls
 - Ecoli.SMART.iprscan.specific.orthologous.count.xls
 - Ecoli.SUPERFAMILY.iprscan.specific.orthologous.count.xls
 - Ecoli.KEGG.iprscan.specific.orthologous.count.xls
 - Ecoli.PANTHER.iprscan.specific.orthologous.count.xls
 - Ecoli.SignalP_EUK.iprscan.specific.orthologous.count.xls
 - Ecoli.INTERPRO.iprscan.specific.orthologous.count.xls
 - Ecoli.PIRSF.iprscan.specific.orthologous.count.xls
 - Ecoli.ProDom.iprscan.specific.orthologous.count.xls
 - Ecoli.PRINTS.iprscan.specific.orthologous.count.xls
 - Ecoli.Phobius.iprscan.specific.orthologous.count.xls
 - Ecoli.UniPathway.iprscan.specific.orthologous.count.xls
 - Ecoli.Hamap.iprscan.specific.orthologous.count.xls
 - Ecoli.TMHMM.iprscan.specific.orthologous.count.xls

Protéines provenant de groupes orthomcl possédant une et une seule séquence par espèce et où toutes les espèces sont présentes.



Status	Name	Description
	input_data	Input data
	summary	Analysis summary - Global counts
	specific_proteins	Extraction of Specific protein sequences and annotation
	orthologous_proteins	Extraction of "orthologous" protein sequences and annotation
	core_proteome	Core proteome pseudomolecule
	venn_diagrams	Lists of all groups with a protein of this proteome + unique protein by proteome
	pan_proteome	Directory with fasta multi-sequence of pan-proteome
	statistics_matrices	Descriptive Statistics: number of OrthOMCL group - Number of proteins by OrthOMCL group - Number of taxa by OrthOMCL group; Phylogenetics matrices : Presence/Absence matrix - ...

Core-genome

- core_proteome
 - [aligned.core.proteome.fasta](#)

Protéines orthologues présentes en une seule copie dans toutes les espèces



Alignements



Concaténation des alignements

Status	Name	Description
	input_data	Input data
	summary	Analysis summary - Global counts
	specific_proteins	Extraction of Specific protein sequences and annotation
	orthologous_proteins	Extraction of "orthologous" protein sequences and annotation
	core_proteome	Core proteome pseudomolecule
	venn_diagrams	Lists of all groups with a protein of this proteome + unique protein by proteome
	pan_proteome	Directory with fasta multi-sequence of pan-proteome
	statistics_matrices	Descriptive Statistics: number of OrthOMCL group - Number of proteins by OrthOMCL group - Number of taxa by OrthOMCL group; Phylogenetics matrices : Presence/Absence matrix - ...

venn_diagrams

- venn_diagrams
 - venn.html
 - lists
 - Xbbri.list.txt
 - Ecoli.list.txt
 - Bacsu.list.txt

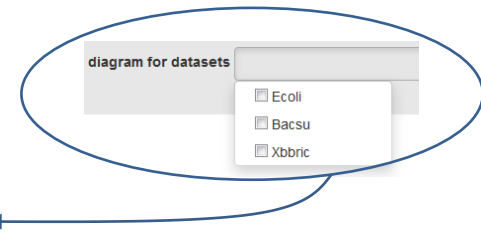
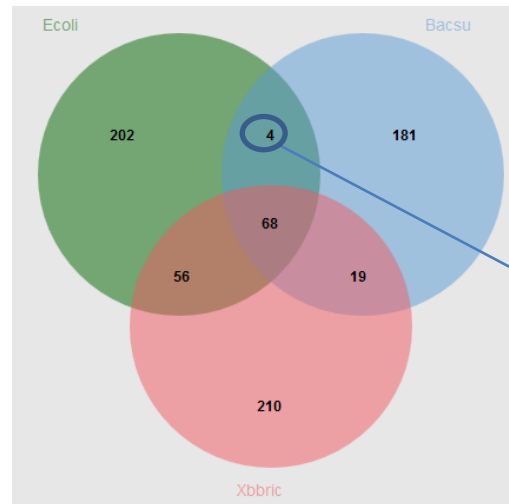
Orthomcl Groups content analysis

Generate Venn diagrams that show shared and specific families proteins of selected proteomes.

On these Venn diagrams :

- intersection between proteomes contain number of orthomcl groups shared by these proteomes.
- number of specific elements of a proteome is the sum of :
 - * specific inparalogs groups of this proteome
 - * specific singlecopy protein of this proteome.

Sub directory lists contain, for each proteome, the list of each group that contain a protein of this proteome and specific singlecopy protein of this proteome.



Click on a venn diagram figure to display the linked elements:

Common elements in Ecoli Bacsu :

- ORTHOMCL161
- ORTHOMCL153
- ORTHOMCL148
- ORTHOMCL160

Status	Name	Description
	input_data	Input data
	summary	Analysis summary - Global counts
	specific_proteins	Extraction of Specific protein sequences and annotation
	orthologous_proteins	Extraction of "orthologous" protein sequences and annotation
	core_proteome	Core proteome pseudomolecule
	venn_diagrams	Lists of all groups with a protein of this proteome + unique protein by proteome
	pan_proteome	Directory with fasta multi-sequence of pan-proteome
	statistics_matrices	Descriptive Statistics: number of OrthoMCL group - Number of proteins by OrthoMCL group - Number of taxa by OrthoMCL group; Phylogenetics matrices : Presence/Absence matrix - ...

statistics-matrices

statistics_matrices

- Number Of Proteins by OrthoMCL group.classes.html
- Number Of Taxa by OrthoMCL group.classes.html
- MatrixAbundance.txt
- MatrixPresenceAbsence.txt

xls_files

- Number Of Taxa by OrthoMCL group.classes.count.xls
- Number Of Proteins by OrthoMCL group.classes.count.xls
- Number Of Taxa by OrthoMCL group.xls
- Number Of Proteins by OrthoMCL group.xls

GROUP	Bacsu	Ecoli	Xbbric
ORTHOMCL0	0	1	1
ORTHOMCL1	1	1	1
ORTHOMCL10	1	1	1
ORTHOMCL100	0	1	1
ORTHOMCL101	0	1	0
ORTHOMCL102	0	1	1
ORTHOMCL103	0	1	0
ORTHOMCL104	0	1	1
ORTHOMCL105	0	1	1
ORTHOMCL106	0	1	1
ORTHOMCL107	0	1	1
ORTHOMCL108	0	1	0
ORTHOMCL109	0	1	1
ORTHOMCL11	1	1	1
ORTHOMCL110	0	1	1
ORTHOMCL111	0	1	0
ORTHOMCL112	0	1	1
ORTHOMCL113	0	1	1
ORTHOMCL114	0	1	1
ORTHOMCL115	0	1	1
ORTHOMCL116	0	1	1
ORTHOMCL117	0	1	1
ORTHOMCL118	0	1	1
ORTHOMCL119	0	1	1
ORTHOMCL12	1	1	1

GROUP	Bacsu	Ecoli	Xbbric
ORTHOMCL0	0	19	2
ORTHOMCL1	11	2	1
ORTHOMCL10	1	2	2
ORTHOMCL100	0	1	1
ORTHOMCL101	0	2	0
ORTHOMCL102	0	1	1
ORTHOMCL103	0	2	0
ORTHOMCL104	0	1	1
ORTHOMCL105	0	1	1
ORTHOMCL106	0	1	1
ORTHOMCL107	0	1	1
ORTHOMCL108	0	2	0
ORTHOMCL109	0	1	1
ORTHOMCL11	3	1	1
ORTHOMCL110	0	1	1
ORTHOMCL111	0	2	0
ORTHOMCL112	0	1	1
ORTHOMCL113	0	1	1
ORTHOMCL114	0	1	1
ORTHOMCL115	0	1	1
ORTHOMCL116	0	1	1
ORTHOMCL117	0	1	1
ORTHOMCL118	0	1	1
ORTHOMCL119	0	1	1
ORTHOMCL12	3	1	1

Je veux comparer des protéomes complets

Je vais sur

L'outil permet

L'outil ne permet pas

Paramètres clés

Pièges

Formats et fichiers