Analyze Data

Workflow

Shared Data

Visualization





### GALAXY pour vos traitements bioinformatiques http://sigenae-workbench.toulouse.inra.fr

### « sig-learning » pour votre auto-formation continue en ligne http://sig-learning.toulouse.inra.fr



29/01/2013



# Vos traitements bioinformatiques avec GALAXY







Plateforme

Vos données

Historique

Workflow

Bioinfo

Vous



Présentation de la plateforme Galaxy.

Comment récupérer vos données (privées et publiques) ?

Notions d'outils, d'historique et de workflow.

Lancement de traitements bioinformatiques.

Guide pour les utilisateurs Galaxy.

### **Galaxy Project**

Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences Jeremy Goecks<sup>1</sup>, Anton Nekrutenko<sup>2\*</sup>, Jemes Taylor<sup>1\*</sup> and The Galaxy Team



Equipe "Galaxy project" :

Le Center for Comparative Genomics and Bioinformatics - Penn State,
Des départements "Biology" et "Mathematics and Computer Science" de l'Université d'Emory.





**Plateforme** 

Vos données

Historique

Bioinfo

Vous





Anton Nekrutenko Penn State



Nate Coraor Penn State



James Taylor Emory

### Une « Galaxy » parmi tant d'autres





Serveur public (https://main.g2.bx.psu.edu/ ): •Gratuit

•Quota limité : pour se familier à l'outil sur des petits jeux de donneés.

Données non protégées



### Nombreuses autres instances : •Curie (Nebula) •URGI



Plateforme

Vos données

Historique

Workflow

Bioinfo

Vous



# Une communauté nationnale et internationnale très active : Listes de diffusion (US, FR) Wiki

•Twitter

•"Galaxy tour de France"

L'instance locale Sigenae de Galaxy :

- •Maintenue par Sigenae.
- •Intégration des outils et scripts "locaux".
- ightarrow Présentation des particuliarités de l'instance Sigenae.



	Galaxy « la bioinformatique pour tous »	
	<ul> <li>Galaxy est :</li> <li>• Open source ».</li> <li>• Développé et maintenu par une communauté active.</li> <li>• Une plateforme proposant un ensemble d'outils bioin</li> <li>• Accessible : http://sigenae-workbench.toulouse.inra</li> <li>• Une "constellation" d'outils (analyser, manipuler, visu</li> </ul>	oformatiques. <b>1.fr/</b> aliser)
Distoformo	Les biologistes peuvent :	
Platelorme	•Lancer des traitements sans Linux.	
Vos données	•Dupliquer des traitements.	
Historique	Partager des analyses complètes.      The second de manière très intuitive l	
Workflow		
Bioinfo	Les bioinformaticiens peuvent : •Faire ajouter des outils.	Not Refer Lands, proc. weakly (S. vert, verse)         EVENUE           Edge         State         State           Edge         State         State           Edge         State         State           Edge         State         State
Vous	Partager des outils (Tool Shed).	
	Partager des worktriows.     Setting     Aussienterstein     Aussienterstein	NY TRANSC Upon Silon. And Pathon and All Solitan All S

	Contexte d'utilisation dans un laboratoire
	<ul> <li>✓ Complémentaire au « cahier de laboratoire », avec archivage des données de séquençage</li> <li>→ Retrouver les données, les outils, les références pour la publication</li> </ul>
Plateforme	✓Manipuler les informations contenues dans un fichier, de façon simple et rapide.
Vos données Historique	<ul> <li>✓ Autres fonctionnalités intéressantes :</li> <li>✓ "mapping" des séquences,</li> <li>✓ analyse des régions de variation ("indel", substitution)</li> </ul>
Workflow Bioinfo	✓ Construction de worflow résumant l'ensemble des fonctionnalités utilisées.
Vous	✓ Intégration de nos propres outils (outils très utiles et fréquemment utilisés)
	→ Galaxy devient VOTRE BOITE A OUTILS.



Vos données sont protégées. Vos jobs sont envoyés sur le cluster. Inutile de savoir programmer De nombreux outils bioinformatiques sont intégrés dans Galaxy.











Analyze Data Workflow Sha	ed Data Visualization Admin Help Wolcome.cmama	Using 13%
User Options	* Upload local file from filesystem path (version 1.0.0)	History Options -
<ul> <li>Your file path : /work/smaman/</li> <li>1 - UPLOAD YOUR DATA</li> <li>Get Data</li> <li>2 - NUES MANIPULATION</li> <li>Text Manipulation</li> <li>Filter and Sort</li> <li>Join, Subtract and Group</li> <li>Convert Formats</li> </ul>	File Name: phiX174_rea( File type: Fastq ▼ Path to file: /work/smaman/phiX174_reads.fastqsanger Execute	<ul> <li>Unnamed history</li> <li>O bytes</li> <li>Your history is empty. Click 'Get Data' on the left pane to start</li> </ul>
FASTA manipulation FASTA manipulation SAM/BAM manipulation : Picard (beta) SAM/BAM manipulation : SAM Tools		▶
4 - MAPPING <u>BWA - Bowtie</u> 5 - INDEL ET SNP <u>Indel Analysis</u> <u>RNA-Seq</u> <u>GATK Tools (beta)</u>	•	

A

1



Analyze Data Workflow	Shared Data	Visualization Admin	Help	no cmaman	vou aro working in	Using 13%
Analyze Data       Workflow         User       Option         Your User Hame: Simamary Hour file path : /work/smamar       Option         Your User Hame: Simamary Hour file path : /work/smamar       Option         1 - UPLOAD YOUR DATA       Option         Get Data       Option         Your User Hame: Simamary Hour file path : /work/smamary       Option         1 - UPLOAD YOUR DATA       Option         Get Data       Option         Text Manipulation       Filter and Sort	Shared Data	Visualization Admin Upload local file from filesys le Name: hiX174_read le type: astq v ath to file:	Help Wolcor		History Galaxy sensibilisation 2 - BWA and FastQO 14: phiX174 reads.fast 1.0 Mb format: fastqsanger () () ()	Using 13% (work / cmaman Options On - TP12.1 Mb tgsanger , database: <u>?</u>
Join, Subtract and Group Convert Formats 3 - SEQUENCES MANIPULATION FASTA manipulation FASTQ manipulation SAM/BAM manipulation : Pic (beta) SAM/BAM manipulation : SA Tools	E Card	vork/smaman/phiX174_reads ∃xecute	fastqsanger		0080917-and-080922:5: GATGTTATTCTTCATTTGG + IIIIIIIIIIII 0080917-and-080922:5: GTTTCTTCTGCGTCAGTAAG	1:185:82 GGTAAAACCTCTTAT HIIQFIBA/IOII4I 1:1366:223 BAACGTCAGTGTTTCC
4 - MAPPING <u>BWA - Bowtie</u> 5 - INDEL ET SNP <u>Indel Analysis</u> <u>RNA-Seq</u> <u>GATK Tools (beta)</u>	*			Ŧ		

ß

1

















### Interface divisée en 4 parties :

- 1 Liste des outils disponibles.
- 2 Visualisation de l'outil utilisé, historique ou workflow en construction.
- 3 Historique ou workflow détaillé.
- 4 Menu .





### **Principaux onglets**



•ANALYSE DATA	: Page d'accueil	de Galaxy.
---------------	------------------	------------

- •WORFLOW : Liste des workflows .
- •SHARED DATA : Liste des datasets, historiques et workflows partagés.
- •VISUALIZATION : Outil de visualisation de vos fichiers résultats.
- •USER : Accès à vos historiques et datasets sauvegardés.



No workflows have been shared with you.







### 2 méthodes de téléchargement de vos données privées





### Téléchargement de données publiques



### Données UCSC, Ensembl, NG6, BIOMART :

Plate

Vos d

Hist

Wo

Bi

1

	<b>^</b>	Genomes	Genome Browser	Tools	Mirrors	Downloads	My Data	About Us	He
	Table Brow	wser							
	clade: Ma	mmal 🔻	genome: Huma	an 👻	assem	bly: Feb. 2009 (	GRCh37/hg19)	•	
	group: Ge track hubs table: know	nes and Gene	Prediction Tracks	track:     describe tab	UCSC Gene ole schema	es 🗸	add custor	n tracks	
eforme	region: 💿	genome 💿	ENCODE Pilot re	egions 💿 po	sition chr2	21:33031597-3304	1570 looku	ıp define re	gions
	identifiers	(names/acc	essions): paste	list upload	list				
onnées	filter: crea	ate							
	intersectio	on: create							
orique	correlation	n: create							
	output for	mat: BED - b	rowser extensible dat	а	- Sen	d output to 🗵	Galaxy 🔲 🤆	<u>GREAT</u>	
rkflow	output file	:		(leave blank	to keep o	utput in browse	er)		
	file type re	eturned: 💿	plain text 💿 gzip	compresse	b				
pinfo	Galaxy		Analyze Data Workflow	Shared Data Admin H	elp User <u>course</u>	Welcome smaman, you are	wa		
	Tools Options	EMBL-EBI	Research Training Industry	About Us Help		Find Help   Feedba	ck		
ous	Serveur :galaxy Get Data	ENA	EBI Home » ENA Home >					•	
	<ul> <li><u>Upload local file from</u> <u>filesystem path</u> Upload data to history without copying on</li> </ul>	ENAHome	European Nucleotide Archive				hio	• ma	rt
	<ul> <li><u>Upload File</u> from your computer</li> </ul>	= Search & Browse = Submit & Update = About ENA	The European Nucleotide Archive (ENA) prov nucleotide sequencing information, covering information and functional annotation more	vides a comprehensive record of raw sequencing data, sequence <u>e</u>	assembly EN			•••••	Iι
	EBI SRA ENA SRA	= Contact	Access to ENA data is provided though the bi download and through the API.	rowser, through search tools, larg	e scale file Europe	an Nucleotide Archive	New	Count R	esults
	<ul> <li><u>UCSC Main</u> table browser</li> <li><u>UCSC Test</u> table browser</li> </ul>	ANNOUNCEMENTS	Text search						
	UCSC Archaea table browser     RX main browser	released 7 Mar 2012	Enter search query, for exam	mple: BN000065	Search		Dataset		
	<u>Get Microbial Data</u>	CRAM toolkit 0.7 has been released. More information					Dataset		
	<u>BioMart</u> Central server	with download, installation and upage instructions are							





### **Gestion de vos historiques**



	History	Options 🔻	
	🕑 🗖 TP FastQC	⊘ 📄 54.0 Mb	
	8: FastQC data 5.htm	<u>1</u> @ / X	
	6: GM.fastqsanger	• / ×	
	5: h1.fastqsanger	• / ×	
Plateforme	4: FastQC data 18.html	• / ×	
os données	3: FASTQ Summary Statistics on data 18	• / %	ш
Historique	2: FASTO Summary Statistics on data 18 76 lines, 1 comments	• / %	
Workflow	format: tabular, databa Info: 99115 fastq read processed.	ase: <u>?</u> s were	
Bioinfo	sequence characters, t data is valid for: sange Input ASCII range: '#'(	ides and the input tr (35) -	
Vous	Input decimal range: 2 Epilog : job finished at	- 34 ven mai	
00	11 10:36:43 CEST 201	2	
	1 2 3 4 5 ‡column count min max su 1 99115 2 33 33	6 um mean 194703 32.2	
	2 99115 2 34 35	156652 31.8	
	2 00115 2 24 21	145060 21 7	

•Conserver toutes les étapes de vos analyses .

### •Partager vos analyses.

•A chaque run d'un outil, une nouvelle dataset est créée. Les données ne sont pas écrasées.

• Répéter, autant de fois que nécessaire, une analyse.

For 0 selected hist	ories: Re	name	Delete	Delet	te Permanentl	y
indexation v	1		<u>0 Taqs</u>		46 bytes	Ju 20
TP FastQC 🔻	12	16	<u>0 Taqs</u>		54.0 Mb	Ma 20
TP:NGS- Polymorphisme	8	2	<u>0 Taqs</u>	<u>Shared</u>	6.6 Gb	Ар 20
FastQC 🔻	6		<u>0 Taqs</u>	<u>Shared</u>	17.4 Mb	Ар 20
SwanPorc 🔻	18		<u>0 Taqs</u>	<u>Shared</u>	0 bytes	Ju 20

### Comment lancer un job sans ligne de commande ?







### 2 – Choisir un outil dans « Tools » :

		•	
	or		
<b>N14</b>	UT.		

**Vos données** 

Workflow

Bioinfo

Vous



#### NGS: Mapping

- Lastz map short reads against reference sequence
- Lastz paired reads map short paired reads against reference sequence
- Map with Bowtie for Illumina
- Map with Bowtie for SOLiD
- Map with BWA for Illumina

Map with BWA for Illumina (version 1.2.2)

Will you select a reference genome from your Use one from the history -

• 1 ×

Select a reference from history: 11: phiX174\_genome.fa 🔻

Is this library mate-paired?: Single-end 👻

FASTQ file:

14: phiX174\_reads.fastqsanger 🝷 FASTQ with either Sanger-scaled guality values (

3 – Lancer le job en cliquant sur « Executer ». L'Execution du job en cours est visible dans votre historique. Fini les lignes de commande !

😂 15: Map with BWA for 👁 🖉 💥 Illumina on data 14 and data 11: mapped reads Job is waiting to run 1







1	Créer un workflow	
	Depuis une page blanche, vous pouvez concevoir un workflow Aide : les résultats produits sont typés, il n'est donc pas pos une dataset sur un mauvais tool !	v. sible de brancher
Plateforme		
Vos données	* Upload local file from & filesystem path	
Historique	out1 (bam, txt, fastqsanger, csfasta, qual, bed, gff, gtf, vcf, sam, fasta, pdf, xsq) Map with BWA for	Illumina 💥
Workflow	PASTQ file	
Bioinfo	output (sam)	
Vous		
020		

### Exporter votre historique en workflow.





**Plateforme** 

**Historique** 

Workflow

Bioinfo

Vous

Depuis votre fenêtre « History », vous pouvez extraire un workflow.

Workflow name Workflow constructed from history 'IGV bai Create Workflow Check all Uncheck all History items created Tool \* Upload local file from filesystem path 1: ERR000017.bam ► Include "\* Upload local file from filesystem path" in workflow \* Upload local file from filesystem path ► 8: ERR000017.sorte Include "\* Upload local file from filesystem path" in workflow \* BAM sorted to BAI for IGV 11: \* BAM sorted to ► Include "\* BAM sorted to BAI for IGV" in workflow

![](_page_31_Picture_1.jpeg)

Options 1

### Cliquer sur le menu « Workflow » pour lister vos workflows :

![](_page_31_Picture_3.jpeg)

Bioinfo

Plateforme

**Vos données** 

**Historique** 

Workflow

Vous

![](_page_31_Picture_6.jpeg)

Vous pouvez ensuite, depuis le menu « Options », soit :
•Editer votre workflow pour le commenter et/ou le modifier.
•Run workflow pour lancer simultanément vos jobs.

![](_page_31_Picture_8.jpeg)

![](_page_32_Figure_0.jpeg)

5 – Editer, partager et lancer vos traitements à volonté (run de votre workflow).

![](_page_33_Figure_0.jpeg)

![](_page_33_Figure_1.jpeg)

### Les principaux outils Galaxy

![](_page_34_Picture_1.jpeg)

### **GET DATA** :

Télécharger vos données privées. Télécharger des données publiques : UCSC, Ensembl, Biomart ...

#### **FILES MANIPULATION :**

Manipulation de fichiers texte ou autres : Couper, coller, comparer, soustraire, merger, concatener, selectionner, filtrer, trier, convertir, grouper ...

### **SEQUENCE MANIPULATION :**

Travailler sur des fichiers FASTA et FastQ, analyse qualité FasQC, Picard tools et samtools.

### **MAPPING** : BWA, Bowtie, indexation de génome.

Autres : Recherche d'indel, SNP, RNAseq (TopHat, Cufflinks).

## Autres outils en cours de tests et d'ajout .. Selon vos besoins.

#### Plateforme

Vos données

Historique

Workflow

Bioinfo

Vous

![](_page_34_Picture_17.jpeg)

Get Data 2 - FILES MANIPULATION

1 - UPLOAD YOUR DATA

Text Manipulation Filter and Sort Join, Subtract and Group Convert Formats

3 - SEQUENCES MANIPULATION

FASTA manipulation

FASTQ manipulation SAM/BAM manipulation : Picard (beta) SAM/BAM manipulation : SAM Tools

4 - MAPPING

BWA - Bowtie

5 - INDEL ET SNP

**Indel Analysis** 

![](_page_35_Picture_0.jpeg)

 Blankenberg D, Von Kuster G, Coraor N, Ananda G, Lazarus R, Mangan M, Nekrutenko A, Taylor J. "Galaxy: a web-based genome analysis tool for experimentalists". Current Protocols in Molecular

 Giardine B, Riemer C, Hardison RC, Burhans R, Elnitski L, Shah P, Zhang Y, Blankenberg D, Albert I, "Galaxy: a platform for interactive large-scale genome analysis." *Genome Research*. 2005 Oct; 15(1)

Biology. 2010 Jan; Chapter 19:Unit 19.10.1-21.

![](_page_36_Picture_1.jpeg)

Ces outils sont nombreux et constituent une bonne alternative à la ligne de commande.

Les traitements sont automatiquement lancés sur GENOTOUL (qsub).

Voici les principaux outils « non bioinfo » proposés :

- •Join (des fichiers lourds), Substract and Group
- Text Manipulation
- •Filter and sort
- Convert Formats

Select first (version 1.0.0)
Select first:
10
ines
from:
4: UCSC Main on Humane (genome) 🔻
Execute

#### What it does

This tool outputs specified number of lines from the **beginning** of a dataset

#### Example

Selecting 2 lines from this:

chr7 56632 56652 D17003 CTCF R6 310 + chr7 56736 56756 D17003 CTCF R7 354 + chr7 56761 56781 D17003 CTCF R4 220 + chr7 56772 56792 D17003 CTCF R7 372 + chr7 56775 56795 D17003 CTCF R4 207 +

will produce:

chr7 56632 56652 D17003\_CTCF\_R6 310 + chr7 56736 56756 D17003\_CTCF\_R7 354 +

![](_page_36_Picture_17.jpeg)

Plateforme

Vos données

**Historique** 

Workflow

**Bioinfo** 

Vous

![](_page_36_Picture_24.jpeg)

![](_page_37_Picture_0.jpeg)

1ê HES

http://www.genomenewsnetwork.org/

**Plateforme** 

**Vos données** 

**Historique** 

**Workflow** 

**Bioinfo** 

CFTR →

![](_page_37_Figure_1.jpeg)

![](_page_37_Picture_2.jpeg)

### Mapper un FASTQ sur une référence avec BWA.

![](_page_37_Figure_4.jpeg)

- Lastz map short reads against reference sequence
- Lastz paired reads map short paired reads against reference sequence
- Map with Bowtie for Illumina
- Map with Bowtle for SOLID
- Map with BWA for Illumina

![](_page_37_Figure_10.jpeg)

![](_page_37_Figure_11.jpeg)

### Visualiser la qualité des données avec FASTQC Report.

Vous

![](_page_37_Picture_14.jpeg)

Visualiser un génome avec UCSC.

![](_page_37_Figure_16.jpeg)

	En résumé
	De nombreux outils disponibles : •Outils de traitement de fichiers •BWA, FastQC, SAM Tools, Picard Tools
Plateforme	Façilité d'ajout de nouveaux scripts / outils selon vos besoins.
Vos données	Par exemple : •GATK,
Historique	•Mirdeep2 •Cutadapt
Workflow	<ul> <li>Indexation de génomes</li> <li>Autros N'hésitez pas à en faire la demande l</li> </ul>
Bioinfo	$\rightarrow$ Mise à jour du menu avec l'ajout d'outils.
Vous	

![](_page_39_Figure_0.jpeg)

![](_page_39_Figure_1.jpeg)

### FAQ et formation en ligne

**Bioinfo** 

Vous

![](_page_40_Picture_1.jpeg)

### Une FAQ et le lien vers « sig-learning » sont disponibles depuis la page d'accueil.

![](_page_40_Picture_3.jpeg)

•Demandes via des tickets Redmine ou mail à sigenae-support@listes.inra.fr

### **En conclusion ...**

![](_page_41_Picture_1.jpeg)

### GALAXY

✓ Simplicité d'utilisation (sans Linux).

✓ Partage de vos datasets, historiques et workflows.

✓ Présentation schématique de vos traitements grâce aux workflows.

✓ Possibilité d'ajout de nouveaux outils, selon vos besoins.

![](_page_41_Picture_7.jpeg)

![](_page_42_Picture_0.jpeg)

# Votre auto-formation continue en ligne avec « sig-learning »

![](_page_42_Picture_2.jpeg)

![](_page_43_Figure_0.jpeg)

![](_page_44_Picture_1.jpeg)

Il vous est possible de vous inscrire directement en ligne à une formation : « Trainings » « Trainings management » puis « Subscribe to training » :

![](_page_44_Picture_3.jpeg)

Plateforme

Vos formations L'inscription s'effectue via une recherche de la formation par mots clés. Voici donc la liste des formations

(disponibles au 01/2013):

![](_page_44_Picture_8.jpeg)

![](_page_44_Picture_9.jpeg)

![](_page_45_Picture_1.jpeg)

Outre une introduction et un carrousel permettant d'accéder aux principaux chapitres de la formation, la page d'accueil de la formation donne accès :

![](_page_45_Picture_3.jpeg)

![](_page_45_Picture_4.jpeg)

**TRAINING PLAN** : Parcours pédagogique avec les supports en ligne.

FORUM	: support de communication entre stagiaires / formateurs.

**TESTS** : Tests et exercices.

LINKS : Liens utiles.

![](_page_46_Picture_1.jpeg)

1 – Demande à compte sur la plateforme BIOINFO GENOTOUL : http://bioinfo.genotoul.fr/index.php?id=81 Ou : bioinfo.genotoul.fr puis «menu « Help », puis « Create an account ». Vous recevrez un login et mot de passe LDAP Genotoul. 2 – Puis utilisez ce login et mot de passe LDAP Genotoul lorsque vous souhaitez accéder à : formations Instance Sigenae de Galaxy : http://sigenae-workbench.toulouse.inra.fr/ .« Sig-learning » : http://sig-learning.toulouse.inra.fr/ 3 - Pour demander une augmentation de votre quota utilisateur sur Galaxy, veuillez vous adresser à : sigenae-support@listes.inra.fr

Vos