# Galaxy for Bioinformaticians

# What is Galaxy

- A **Open-Source** platform that allows to run complex bioinformatics analysis with a graphics interface
- More than 350 public Galaxy resources:
  - 120+ public servers, many more non-public
  - Both general-purpose and domain-specific
- "Easy" to install and maintain local instances

#### • Accessibility

• Users without programming experience can easily upload/retrieve data, run complex tools and workflows, and visualize data

#### Reproducibility

 Galaxy captures information so that any user can understand and repeat a complete computational analysis

#### • Transparency

- Users can share or publish their analyses (histories, workflows, visualisations)
- Pages: online Methods for your paper























#### Tools

Protein auerv se	auence(s)	
	No facta or facta dz datasot available	 
(-query)	Ino fasta of fasta.gz dataset available.	 
Subject databas	e/sequences	
Locally installed E	SLAST database	 
Protein BLAST	database	
Select/Unsel	ect all	
,		
Type of BLAST		
⊙ blastp - Traditio	onal BLASTP to compare a protein query to a protein database	
O blastp-short - I	3LASTP optimized for queries shorter than 30 residues	
🔿 blastp-fast - Us	e longer words for seeding, faster but less accurate	
See help text for d	lefault parameter values for each BLAST type. (-task)	
Set expectation	value cutoff	
0.001		
(-evalue)		
Output format		
Tabular (extende	d 25 columns)	
(-outfmt)		 
Advanced Option	15	
Hide Advanced O	ations	

- A tool form contains:
  - $\,\odot\,$  input datasets and parameters
  - $\circ$  help, citations, metadata
  - $\,\odot\,$  an Execute button to start a job, which will add some output datasets to the history
- New tool versions can be installed without removing old ones to ensure reproducibility

#### **T** Galaxy Tool Shed

#### Repositories Groups Help - User -

6532 valid tools on Dec 04, 2018

#### Search

- Search for valid tools
- Search for workflows

#### Valid Galaxy Utilities

- Tools
- Custom datatypes
- Repository dependency definitions
- Tool dependency definitions

#### All Repositories

Browse by category

#### **Available Actions**

Login to create a repository

#### **Repositories by Category**

search repository name, description

Name

e, description Q
Description

Repositories

<u> </u>		
Assembly	Tools for working with assemblies	128
ChIP-seq	Tools for analyzing and manipulating ChIP- seq data.	65
<u>Combinatorial</u> Selections	Tools for combinatorial selection	10
<u>Computational</u> chemistry	Tools for use in computational chemistry	76
<u>Constructive Solid</u> Geometry	Tools for constructing and analyzing 3-dimensional shapes and their properties	12
Convert Formats	Tools for converting data formats	114
	Tools for exporting data to various	-

- Free "app" store: Galaxy Tool Shed
  - Thousands of tools already available
  - Most software can be integrated
  - If a tool is not available, ask the Galaxy community for help!
    - Only a Galaxy admin can install tools

### Data

- Copy/paste from a file
- Upload data from a local computer
- Upload data from internet using URL
- Upload data from online databases: UCSC, BioMart, ENCODE, modENCODE, Flymine etc.
- Import from Shared Data (libraries, histories, pages)
- Upload data from FTP

# Data types

- Tools only accept input datasets with the appropriate datatypes
- When uploading a dataset, its datatype can be either:
- automatically detected
- assigned by user
- Dataset produced by a tool: datatype assigned by the tool

# History

- Location of all analyses
  - collects all datasets produced by tools
  - collects all operations performed on the data

History		C 🕈 🗆	]
search o	latasets	0	)
Galaxy 1	01		
7 shown			
9.07 MB			
7: Comp	<u>are two</u> son data 6	● 🖋 🗙	
and data	1		
5 regions			
format: b	ed, database	e: hg38	
join (GN	U coreutils) 8	3.22	
Foundati	t (C) 2013 F ion, Inc.	ree Software	
License (	GPLv3+: GN	U GPL	
version 3	3 or later <h /apl.html&gt;</h 	ttp://gnu.org	
This is fr	ee software:	you are free	
to chang	e and redistr	ribute it.	
ext		in, to the	
B 0 2	) <u>III</u> ?	<b>&gt;</b>	
display in	IGB View	۵	
display in display wi	IGB <u>View</u> th IGV <u>local</u>	Norman hg38	
display in display wi display at	IGB <u>View</u> th IGV <u>local</u> UCSC <u>main</u>	Numan hg38 test	
display in display wi display at 1.chrom 2 chr22 4	GB View th IGV local UCSC main .start 3.En		
display in display wi display at 1.chrom 2 chr22 4 chr22 1	IGB View           IGB View           th IGV local           UCSC main           .start         3.Em           6256560         4626           5690077         1569	Image: Weight of the second state of the se	
display in display with display at 1.chrom 2 chr22 4 chr22 1 chr22 1	IGB View           IGB View           th IGV local           UCSC main           0.start           6256560           5690077           5528158           1552	Muman hg38           test           d         4.Name           3322         uc003bhh.           0709         uc010gqp.           9139         uc011agd.	
display in display with display at 1.chrom 2 chr22 4 chr22 1 chr22 1 chr22 1	IGB View           IGB View           th IGV local           UCSC main           .start         3.En           6256560         4626           5690077         1569           5528158         1552           5690245         1569	Muman hg38           test           d         4.Name           3322         uc003bhh.           0709         uc010gqp.           9139         uc011agd.           0709         uc062bek.	
display in display wi display wi display at 1.chrom 2 chr22 4 chr22 1 chr22 1 chr22 1 chr22 2	IGB View           IGB View           th IGV local           UCSC main           .start         3.En           6256560         4626           5690077         1569           5528158         1552           5690245         1569           2376182         2237	Muman hg38           test           d         4.Name           3322         uc003bhh.           0709         uc010gqp.           9139         uc011agd.           0709         uc062bek.           6505         uc062cbs.	
display in display wi display at 1.chrom 2 chr22 4 chr22 1 chr22 1 chr22 2	IGB View           IGB View           th IGV local           UCSC main           .start         3.En           6256560         4626           5690077         1569           5528158         1552           5690245         1569           2376182         2237	Auman hg38 test      d     4.Name      3322 uc003bhh.      0709 uc010gqp.      9139 uc011agd.      0709 uc062bek.      6505 uc062cbs.	
display in display wi display wi display at 1.chrom 2 chr22 4 chr22 1 chr22 1 chr22 1 chr22 2 display at chr22 3 chr22 4 chr22 1 chr22 2 chr22 4 chr22 1 chr22 1 chr22 2 chr22 4 chr22 1 chr22 4 chr22 1 chr22 4 chr22 1 chr22 2 chr22 4 chr22 1 chr22 1 chr22 2 chr22 4 chr22 1 chr22 1 chr22 1 chr22 2 chr22 4 chr22 1 chr22 2 chr22 4 chr22 1 chr22 5 chr22 4 chr22 1 chr22 4 chr22 1 chr22 4 chr22 1 chr22 4 chr22 1 chr22 5 chr22 4 chr22 1 chr22 1 chr22 4 chr22 1 chr22 1 chr24 1 chr25 1	IGB View           IGB View           th IGV local           UCSC main           .start         3.En           6256560         4626           5690077         1569           5528158         1552           5690245         1569           2376182         2237		
display in display wi display at     1.chrom 2     chr22 4     chr22 1     chr22 1     chr22 2      chr22 2       chr22 5: Sort or	Idl         ?           IGB View         IGV local           th IGV local         UCSC main           .start         3.En           6256560         4626           5690077         1569           5528158         1552           5690245         1569           2376182         2237           first on         Interval	A .Name      A .Name	
display in display wi display wi display at 1.chrom 2 chr22 4 chr22 1 chr22 1 chr22 2 6: Select data 5 5: Sort or 4: Group	Idl         ?           IGB View         IGV local           IGV local         UCSC main           .start         3.En           6256560         4626           5690077         1569           5528158         1552           5690245         1569           2376182         2237           first on         n           data 4         00 data 3	Muman hg38         test         d       4.Name         3322       uc003bhh.         0709       uc010gqp.         9139       uc011agd.         0709       uc062bek.         6505       uc062cbs.         000       x         (a)       x         (b)       x	
display in display wi display wi display at 1. chrom 2 chr22 4 chr22 1 chr22 1 chr22 1 chr22 2 6: Select data 5 5: Sort on 4: Group	IGB View           IGB View           th IGV local           UCSC main           .start         3.En           6256560         4626           5690077         1569           5528158         1552           5690245         1569           2376182         2237           first on         n           data 4         on data 3		
display in display wi display wi display at 1. chrom 2 chr22 4 chr22 1 chr22 1 chr22 1 chr22 2 6: Select data 5 5: Sort or 4: Group 3: Join or	IIB View         IGB View         th IGV local         UCSC main         .start       3.En         6256560       4626         5690077       1569         5528158       1552         5690245       1569         2376182       2237         first on       1         n data 4       0         on data 2       1		
display in display wi display wi display at 1.chrom 2 chr22 4 chr22 1 chr22 1 chr22 1 chr22 2 6: Select data 5 5: Sort or 4: Group 3: Join or and data	Idl ?         IGB View         th IGV local         UCSC main         .start 3.En         6256560 4626         5690077 1569         5528158 1552         5690245 1569         2376182 2237         first on         n data 4         on data 3         n data 2         1	Human hg38         test         d       4.Name         3322       uc003bhh.         0709       uc010gqp.         9139       uc011agd.         0709       uc062bek.         6505       uc062cbs.         Image: State Stat	

# Workflows

- Workflows are build by:
  - Extracting from a history
  - manually by adding and configuring tools using the canvas
  - Importing an existing shared workflow
- Workflows allow:
  - Create sub-workflows: a workflow inside another workflow
  - Share workflows for publication and with the community

### Workflows



# Example: ATAC Seq

- 1. Upload libraries
- 2. Check Quality
- 3. Map to human genere
- 4. Filter out low quality alignments, unpair reads, and mitochondrial reads
- 5. Remove duplicates
- 6. Visualise insert size (QC)
- 7. Convert BAM to BED
- 8. Call Peaks with MACS2
- 9. Convert BED to BigWig
- 10. Visualise Peaks in genome context



#### But... Who cares? We can write COde



#### But... Who cares? We can write COde

### We still can benefit from...

#### Reproducibility

- For each produced dataset the history tracks:
- name, format, size, creation time, datatype-specific metadata
- tool id, version, inputs, parameters
- standard output (stdout) and error (stderr)
- Automatisation via Workflows
  - Re-run the same analysis on different input data sets
  - Change parameters before re-running a similar analysis
  - Make use of the workflow job scheduling

# The Galaxy API

"In addition to being accessible through a web interface, Galaxy can also be accessed programmatically, through shell scripts and other programs. The web interface is appropriate for things like exploratory analysis, visualisation, construction of workflows, and rerunning workflows on new datasets.

The web interface is less suitable for things like: Connecting a Galaxy instance directly to your sequencer and running workflows whenever data is ready Running a workflow against multiple datasets (which can be done with the web interface, but is tedious) When the analysis involves complex control, such as looping and branching.

The Galaxy API addresses these and other situations by exposing Galaxy internals through an additional interface, known as an Applications Programming Interface, or API."

#### <u>https://galaxyproject.org/develop/api/</u> <u>https://docs.galaxyproject.org/en/latest/api\_doc.html</u>



# we can code



# we can code

# Bindings

- The Galaxy code itself contains JavaScript bindings.
- BioBlend contains a set of Python bindings
- blend4j contains Java bindings largely modeled after BioBlend. JavaDocs
- blend4php contains PHP bindings.
- CLI Bindings: Parsec

# CLI Bindings: Parsec

- <u>https://github.com/galaxy-iuc/parsec/</u>
- parsec is a set of wrappers for BioBlend's API. It builds a set of small, useful utilities for talking to Galaxy servers.
   Each utility is implemented as a subcommand of parsec
- <u>https://parsec.readthedocs.org</u>

# Example

- > parsec init
- > parsec histories get\_histories
- > parsec histories show\_history b381997063dbe2e8
- > parsec workflows get\_workflows
- > parsec workflows show\_workflow fee63819c4abec0e

# Example

- > HISTORY\_ID=\$(parsec histories create\_history | jq .id)
- > parsec histories update\_history --name 'Parsec test' \$HISTORY\_ID
- > parsec tools upload\_file my\_file\_R1 \$HISTORY\_ID
- > parsec tools upload\_file my\_file\_R2 \$HISTORY\_ID
- > # create json input file
- > WORKFLOW\_ID=\$(parsec workflows invoke\_workflow fee63819c4abec0e \
   --inputs inputs.json \
   --history\_id \$HISTORY\_ID \
   | jq .workflow\_id)
- > parsec utils wait\_on\_invocation \$WORKFLOW\_ID

# But... There are many workflow managers

- steeper learning curve
- You need to integrate the tools yourself, and take care off all data typing and compatibility
- Hard to share with other non-techie colleagues