An Introduction to = Galaxy

Daniel Blankenberg The Galaxy Team | Penn State <u>http://GalaxyProject.org</u>

Goals

What is Galaxy

Mission, Terminology, and Capabilities Using Galaxy Play along with a basic Analysis How to get involved with Galaxy The community and how to participate Adding Tools Making your code Accessible and Reproducible

Bonus: Brief Intro to Sequence Analysis in Galaxy

Goals

What is Galaxy

Mission, Terminology, and Capabilities Using Galaxy Play along with a basic Analysis How to get involved with Galaxy The community and how to participate Adding Tools Making your code Accessible and Reproducible

Bonus: Brief Intro to Sequence Analysis in Galaxy

Overview

What is Galaxy?

Analysis Interface, Tools, and Datasources

Workflows

Visualizations

Sharing, Publishing, and Galaxy Pages

ToolShed

Galaxy Project Mission

Galaxy is a data integration and analysis platform emphasizing accessibility, reproducibility, and transparency

Accessible: Users without programming experience can easily specify parameters and run tools and workflows.

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

Transparent: Users share and publish analyses via the web and create Pages, interactive, web-based documents that describe a complete analysis.

What is Galaxy?

A data analysis and integration tool

A free (for everyone) web service integrating a wealth of tools, compute resources, terabytes of reference data and permanent storage

Open source software that makes integrating your own tools and data, and customizing for your own site simple

There are several ways to use Galaxy

Using Galaxy - 4 ways

- Public Main Galaxy web instance: *usegalaxy.org*
- Local instance: *getgalaxy.org*
- Cloud instance: *usegalaxy.org/cloud*
- Other Public Galaxy web instances hosted by various groups:

wiki.galaxyproject.org/PublicGalaxyServers



Galaxy as a Genomics WorkBench

Dataset:

Any input, output or intermediate set of data + metadata. A record of a specific data or analysis step.

History:

A series of inputs, analysis steps, intermediate datasets, and outputs. A record of a group of data and analysis steps.

Tool:

An operation within Galaxy that acts upon dataset(s) as an analysis step. May be developed by Galaxy team or a 3rd party program that has been "wrapped" for Galaxy.

Workflow:

A series of analysis steps executed as a unit.

More Galaxy Terminology

Share:

Make something available to someone else

Publish:

Make something available to everyone

Galaxy Page:

Analysis documentation within Galaxy; easy to embed and link to any Galaxy object (histories, datasets, workflows, visualization) or external resource (video, graphics, publications).

Visualize:

External resources. Trackster. Galaxy Charts & Visualizations.

Overview

What is Galaxy?

Analysis Interface, Tools, and Datasources

Workflows

Visualizations

Sharing, Publishing, and Galaxy Pages

ToolShed

Galaxy Analysis Workspace

- Galaxy	Analyze Data Workflow Shared Data - Visualization - Cloud - Admin Help - User -	Usi	ng 10.0 TB
Tools	Map with BWA for Illumina (version 1.2.3)	History	C 🕈
search tools	Will you select a reference genome from your history or use a built-in index?: Use a built-in index	Galaxy 101 NGS Variant 313.4 MB	Q 📎 🗩
Send Data	Select a reference genome:	21: Filter on data 20	• / ×
Lift-Over Text Manipulation	Human (Homo sapiens) (hg19 with mtDNA replaced with rCRS): Homo_sapiens_nuHg19_mtrCRS 🔹	20: Filter on data 19	• # ×
Convert Formats FASTA manipulation	Is this library mate-paired?: Paired-end	<u>19: Variant Annotator on</u> <u>data 17</u>	• # ×
Filter and Sort Join, Subtract and Group Extract Features	Forward FASTQ file: 1: raw_child-ds-1.fq FASTO with either Sanger-scaled quality values (fastosanger) or Illumina-scaled quality values (fastoillumina)	<u>18: FreeBayes on data 15</u> (variants)	• / ×
Fetch Sequences	Reverse FASTQ file:	<u>17: Naive Variant Caller</u> on data 15	• / ×
Get Genomic Scores	FASTQ with either Sanger-scaled quality values (fastqsanger) or Illumina-scaled quality values (fastqillumina)	16: child-mother Merge BAM Files.log	● / ×
Statistics	BWA settings to use:	15: child-mother.bam	• / ×
<u>Graph/Display Data</u> Regional Variation	For most mapping needs use Commonly Used settings. If you want full control use Full Parameter List	14: Add or Replace	• / ×
Multiple regression	Suppress the header in the output SAM file:	Groups on data 12: bam with read groups replaced	
<u>Multivariate Analysis</u> Evolution	BWA produces SAM with several lines of header information	13: Add or Replace Groups on data 11: bam	• / ×
Motif Tools Multiple Alignments	Execute	with read groups replaced	
Metagenomic analyses	What it does	12: SAM-to-BAM on data 10: converted BAM	• / ×
Genome Diversity	BWA is a fast light-weighted tool that aligns relatively short sequences (queries) to a sequence database (large), such as the human reference genome. It is developed by Heng Li at the Sanger Insitute. Li H. and Durbin R. (2009) Fast and accurate short read alignment with Burrows-Wheeler transform.	11: SAM-to-BAM on data 9: converted BAM	• / ×
Phenotype Association NGS: QC and manipulation	Bioinformatics, 25, 1754-60.	10: Filter SAM on data 8	• / ×
<	Know what you are doing		>

Galaxy Analysis Workspace

Using 10.0 TB

2 4

Q > 9

• / ×

• / ×

ator on 💿 🖋 🗙

data 15 💿 🖋 🗙

Caller 💿 🖋 🗙

Merge 💿 🖋 🗙

• / ×

• / ×

• / ×

on data 💿 🖋 🗙

on data 💿 🖋 🗙

data 8 💿 🖋 🗙

>

e 2: bam replaced

ce 11: bam replaced

			8: A job that will surely	• / %		
- Galaxy		Analyze Data Workflow Shared	fail			
Tools	Load Data	Map with BWA for Illumina (version 1.2.3)				History
search tools	0	Will you select a reference genome from your history or use	<u>7: top 5 exons</u>	• 0 %		Galaxy 101 NGS Variant
Get Data		Use a built-in index +		~ ~ ~		
Lift-Over		Select a reference genome:	5 Select first on data 5	@ (/ X		21: Filter on data 20
Text Manipulation		Human (Homo sapiens) (ng19 with mtDNA replaced with rCKs		A		20: Filter on data 19
Convert Formats		Is this library mate-paired?:	5: Sort on data 4	• (%		19: Variant Annotator on
FASTA manipulation	ĺ	Faired-end +	15,310 lines			data 17
Filter and Sort	roup	Forward FASTQ file:	format: tabular, database: ho	19		18: FreeBayes on data 15
Extract Features	<u>toup</u>	FASTQ with either Sanger-scaled quality values (fastqsanger) or		20 🕞		<u>(variants)</u>
Fetch Sequences		Reverse FASTQ file:				17: Naive Variant Caller
Fetch Alignments		2: raw_child-ds-2.fq \$				On data 15
Get Genomic Scores		FASTQ with either Sanger-scaled quality values (fastqsanger) or				16: child-mother Merge BAM Files.log
Operate on Genomic	<u>: Intervals</u>	BWA settings to use:	uc003qqn.2_cds_0_0_chr6_1570	099238_f 1		15: child-mother ham
Graph/Display Data		Commonly Used \$	uc003qqo.2_cds_0_0_chr6_1570	099238_f 1		13. child Hiothensum
Regional Variation		For most mapping needs use Commonly Used settings. If you v	uc003000 2 cds 0 0 chr6 1570	099238 £ 1		14: Add or Replace Groups on data 12: bam
Multiple regression		Suppress the header in the output SAM file:				with read groups replace
Multivariate Analysis	<u>s</u>	BWA produces SAM with several lines of header information	uc003hqu.2_cds_4_0_chr4_885:	34937_1 1		13: Add or Replace
Motif Tools			uc001vqv.2_cds_1_0_chr13_110	0434389_r 8		Groups on data 11: bam with read groups replace
Multiple Alignments		Execute	uc001vsb.1 cds 0 0 chr13 112	2721973 f 8		12: SAM-to-PAM on data
Metagenomic analys	ies	What it does	6	224 1		10: converted BAM
Genome Diversity		RWA is a fast light-weighted tool that aligns relatively short seg			ince genome. It is	11: SAM-to-BAM on data
NGS TOOLBOX BETA		developed by Heng Li at the Sanger Insitute. Li H. and Durbin R.			isform.	9: converted BAM
Phenotype Associati	on	Bioinformatics, 25, 1754-60.				10: Filter SAM on data 8
NGS: QC and manipu	llation					
<		Know what you are doing				

Tool Suites

Text Manipulation Format Converters Filtering and Sorting Join, Subtract, Group Sequence Tools Multi-species Alignment Tools Genomic Interval Operations Statistics Graphing / Plotting RNA Structure Prediction Regional Variation EMBOSS Evolution / Phylogeny RNA-seq ChIP-seq GATK Picard Metagenomics Motif Tools ...and more

Datasources

Upload file from your computer

- Import via URL
- FTP support for large datasets

UCSC table browser

EBI SRA

BioMart

interMine / modMine

GenomeSpace

add more!

Overview

What is Galaxy?

Analysis Interface, Tools, and Datasources

Workflows

Visualizations

Sharing, Publishing, and Galaxy Pages

ToolShed

Create Workflow Automatically

Extract Workflow from History Create a workflow from a History that

you created interactively.



unning workflow "metagenomic analysis"	Expand All	Collapse
eneric workflow for performing a metagenomic analysis on NGS data.		
<u>Step 1: Input dataset</u> 454 Reads		0
reads D 1: 454 reads type to filter		
<u>Step 2: Input dataset</u> 454 Quality Dataset		
qualities 2: 454 qualities type to filter		
itep 3: Select high quality segments (version 1.0.0) Here we select segments of reads with contiguous high quality bases above threshold phred score of 20		
<u>Step 4: FASTA-to-Tabular</u> (version 1.1.0) Convert to tabular format so that column for additional metadata can be added		
Add colur Step 14: Find lowest diagnostic rank (version 1.0.1)		
Get reads specific to ranks below Kingdom level itep 6: Ta Convert b Step 15: Summarize taxonomy (version 1.0.0) Tabulate list of taxonomic groups contained in reads from dataset 14		
Step 7: M Step 16: Draw phylogeny (version 1.0.0) Step 8: M		



Run workflow

Workflow Editor



Overview

What is Galaxy?

Analysis Interface, Tools, and Datasources

Workflows

Visualizations

Sharing, Publishing, and Galaxy Pages

ToolShed

Visualizing Genomics

Supported external browsers

- UCSC
- Ensembl
- GBrowse
- IGB
- IGV

Traditional browser strengths:

- Showing what is nearby
- what else is happening here
- highlighting correlations
- integrating many datasets



Integrative Genomics Viewer (IGV)



000							IG	V										
Mouse mm9	•	chr1	•	chr1:98,5	582,224-9	8,597,3	370	Go	Ê	Ø [1			Ξ]			+
		qA	A2 qA4	qB	qC1.1	qC1.3	qC3 q	C4	qD	qE1.1	qE2.2	qE3	qF	q61	qH1	qH2.3	qH4	qH6
	NAME DATA FILE DATA TYPE	db 	98,584 kb 	1	98,586 kb 	Î	98,588 	kb	15 9	kb — 8,590 kb 	I	98,592 kb 	Î	98,594 	t kb	1	98,596 kb 	
galaxy_f2979acbfb2c63 75.bam Coverage galaxy_f2979acbfb2c63 75.bam		[0 - 10] 		Ì	i i I I I	1	Î.Î		I II I II		1	II I	l		11 1		1	
chr1:98589793															1	11	<mark>3M</mark> of 268	м

Galaxy

- tool integration framework
- heavy focus on usability
- sharing, publication framework

Trackster

Genome Browser

- physical depiction of data
- visually identify correlations
- + find interesting regions, features

Trackster

View your data from within Galaxy

- No data transfers to external site
- Use it locally, even without internet access

Supports common filetypes

+ BAM, BED, GFF/GTF, WIG

Unique features

- custom genomes
- highly interactive

Trackster: Galaxy's embedded track browser

- Gal	axy / Sable	Analyze Data	Workflow Sh	ared Data-	Visualization -	Admin	Help- U	ser –	Using 59	3.1 MB
RNA-Seq	Example (hg19)	chr	19	;	3,221,594 - 3,3	01,240	₽₽		849	88
0,000 IIIGenomes U	3,230,000 CSC hg19, chr19 gene annotation	3,240,000	3,250,000	3,26	50,000	3,270,0	000	3,280,000	3,290,000	3,3
			1				1			
Brain: asser	nbled transcripts from Cufflinks	-				_			·	_
Brain: splice	e junctions from TopHat									
1. A.			1				8 8	· · · · · · · · · · · · · · · · · · ·		
Adrenal ass	embled transcripts from Cufflinks									-
Adrenal spli	ce junctions from TopHat									
Adrenal acc	epted hits from TopHat									
									1	
Brain: accep	oted_hits from TopHat									
0.000	3,230,000	3,240,000	3,250,000	3.26	50,000	3,270,0	000	3,280,000	3,290,000	<

Circster



Create a visualization in Galaxy



or

28: Brain: transcript 211 lines format: gt Info: cuffli cufflinks - 300000 -F p 4 P 4	assemble s from Cu f, database nks v2.0.2 qno-up 0.100000	d @ fflinks a: hg19 odate-check b -j 0.1500	0 X (-1 00 -
displa Visu display at	lalize mair	<u>n</u> Surrent	
1.Segname	2.Source	3.Feature	4.St
chr19	Cufflinks	transcript	3348:
chr19	Cufflinks	exon	3348:
chr19	Cufflinks	transcript	3349:
chr19	Cufflinks	exon	3349:
chr19	Cufflinks	transcript	3351
chr19	Cufflinks	exon	3351
1)+

Visualization inside Galaxy

- Leverages visualization to evaluate and refine analyses
- Exposes basic analyses in visualization to make it more informative
- Makes that analyze-visualize-refine loop seamless and fast
- Enables learning tools and exploring their parameter space
- Supports custom genome browsers, without a predefined reference genome





Overview

What is Galaxy?

Analysis Interface, Tools, and Datasources

Workflows

Visualizations

Sharing, Publishing, and Galaxy Pages

ToolShed

Sharing & Publishing enables Reproducibility

Reproducibility: Everybody talks about it, but ...

Galaxy aims to push the goal of reproducibility from the bench to the bioinformatics realm

All analysis in Galaxy is recorded without any extra effort from the user.

Histories, workflows, visualizations and *pages* can be shared with others or published to the world.

Sharing and Publishing

Sharing and Publishing History 'Variant Analysis for Sample E18'

Making History Accessible via Link and Publishing It

This history is currently restricted so that only you and the users listed below can access it. You can:

Make History Accessible via Link

Generates a web link that you can share with other people so that they can view and import the history.

Make History Accessible and Publish

Makes the history accessible via link (see above) and publishes the history to Galaxy's <u>Published Histories</u> section, where it is publicly listed and searchable.

Sharing History with Specific Users

You have not shared this history with any users.

Share with a user

Back to Histories List

Sharing and Publishing

Sharing and Publishing History 'Variant Analysis for Sample E18'

Making History Accessible via Link and Publishing It

This history accessible via link and published.

Anyone can view and import this history by visiting the following URL:

http://main.g2.bx.psu.edu/u/jgoecks/h/variant-analysis-for-sample-e18 🥖

This history is publicly listed and searchable in Galaxy's Published Histories section.

You can:

Unpublish History

Removes history from Galaxy's Published Histories section so that it is not publicly listed or searchable.

Disable Access to History via Link and Unpublish

Disables history's link so that it is not accessible and removes history from Galaxy's <u>Published Histories</u> section so that it is not publicly listed or searchable.

Sharing History with Specific Users

You have not shared this history with any users.

Share with a user

Back to Histories List

Galaxy F Galaxy F Galaxy F Galaxy F Galaxy F	ublished I for-sample	History Variant Analysis for Sample E18	
- Galaxy Analyze Data	Workflow	Shared Data Visualization Help User	
Published Histories jgoecks Variant Analysis for Sample E18			About this History
Galaxy History ' Variant Analysis for Sample E18' Annotation: Perform a pileup analysis with default parameters to identify var Dataset	riants in sa	mple E18.	Author jgoecks
<u>1: E18 PE.1 Reads</u> 2: E18 PE.2 Reads	@ @	Forward reads from sample E18. Reverse reads from sample E18.	Related Histories All published histories Published histories by igoecks
3: E18 PE.1 Reads Groomed 4: E18 PE.2 Reads Groomed	@ @	Groom reads to convert quality scores from Solexa 1.0 to Solexa 1.3 Groom reads to convert quality scores from Solexa 1.0	Rating Community (1 rating, 4.0 average)
5: E18 PE.1 Reads Groomed, Trimmed	Ф	to Solexa 1.3 Trim reads from 3° end to remove low-quality nts.	Yours Tags
6: E18 PE.2 Reads Groomed, Trimmed 7: Map with Bowtie for Illumina on data 6 and data 5	40 40	Trim reads from 3' to remove low-quality nts. Map paired-end reads with default parameters.	snp pileup bowtie demo sample
8: SAM-to-BAM on data 7 9: Generate pileup on data 8	@ @	Need to convert Bowtie SAM to BAM so that pileup analysis can be performed. Pileup analysis with default parameters	Yours: snp x pileup x bowtie x demo x sample:e18 x 4
10: Filter pileup to get Variants from sample E18	۲	Find variants with coverage >= 30.	
13: Filter to get Variants from sample E18 where consensus base different than ref. base	æ	Filter pileup to find variants where the consensus base is different than the reference base.	
14: UCSC mm9 RefSeq Genes	40	UCSC mm9 RefSeq genes.	
15: Intersect to get Variants from sample E18, consensus different, in RefSeq Genes	æ	Variants with consensus different that occur in RefSeq genes.	

000	Galaxy Published Workflo	w SNP variant detec	tion from paired	-end reads		
+ http://main.g2.bx.psu.e	du/u/jgoecks/w/snp-variant-detecti	ion-from-paired-end-	reads	¢ Q+ G	oogle)
- Galaxy	Analyze Data Workflo	w Shared Data Visu	alization Help	User		
Published Workflows jgoecks SNP variant det	ection from paired-end reads				About this Workflow	
Step 6: FASTQ Trimmer FASTQ File Output dataset 'output_file' from step 4 Define Base Offsets as Absolute Values Offset from 5' end 0 Offset from 3' end 9 Keep reads with zero length False Step 7: Map with Bowtie for Illumina Will you select a reference genome from your 1 Use a built-in index Select a reference genome /galaxy/data/apiMel3/bowtie_index/apiMel3 Is this library mate-paired?	history or use a built-in index?	Trim reads to remov	t low-quality bases.		Author jgoecks Related Workflows All published workflows by is Rating Community (0 ratings, 0.0 average) Yours Tags Community: snp bowtie Yours: snp bowtie x	toecks *****
Paired-end Forward FASTQ file Output dataset 'output_file' from step 6 Reverse FASTQ file Output dataset 'output_file' from step 5 Maximum insert size for valid paired-end alig 1000 The upstream/downstream mate orientation f the forward reference strand (fr/rf/ff) FR (for Illumina) Bowtie settings to use Commonly used Suppress the header in the output SAM file True	nments (-X) or valid paired-end alignment against					
Step 8: SAM-to-BAM Choose the source for the reference list Locally cached		Convert Bowtie SAM can be run.	output to BAM format	t so that pileup		

00		Galax	y Published H	Histories			
< - + 🙆 http	://main.g2.bx.psu.edu/history/list_published				1000	¢ (Q+ Ge	oogie
Galaxy	Analyze Data	Workflow	Shared Data	Visualization	Help	User	
ublished His	tories						
search	Advanced Search						
Name	Annotation		Owner	Communi Rating †	ĽΧ	Community Tags	Last Updated
<u>Galaxy vs MEGAN</u>	Comparison of Galaxy vs. MEGAN pipeline.		aunl	***	k de	metagenomics megan galaxy	Mar 19, 2010
<u>metagenomic</u> analysis			aunl	****	t de	metagenomics galaxy	Mar 19, 2010
<u>SM 1186088</u>	Datasets correspond to our paper published in Scie Peleg et al. entitled : Altered histone acetylation is associated with age-dependent memory impairmen Experiment layout: This history contains 4 datasets form of BED files of uniquely mapped reads produc chip-seq for histone modifications H4K12ac and H mouse hippocampus of 3 months (young) and 16 m (old) mice after fear conditioning. For detailed infor please refer to supplementary materials and metho respective work by peleg et al.	ence by at. in the ed after 3K9ac in nonths rmation ds of the	fischerlab	****	k skr		Apr 19, 2010
Variant Analysis for Sample E18	Perform a pileup analysis with default parameters t variants in sample E18.	o identify	jgoecks	****	r it	snp pileup bowtie demo sample	2 minutes ago
get longest exon			henri	***	ksk	chr22 longest marc exon human workshop	Sep 02, 2010
FASTA to Tabular Test			n	***	in .		Aug 26, 2010
EKLE			yzc109	***	init -		Aug 24, 2010

The power of Galaxy publishing

Galaxy's publishing features facilitate access and reproducibility without any extra leg work

One click grants access to the *actual analysis* you performed to generate your original results

- Not just data access: the full pipeline
- Annotate each step
- Anyone can import your work and immediately reproduce or build on it

Galaxy Pages: A web-based, interactive medium for presenting all aspects of an analysis: data, methods, and results
Sharing and Publishing Your Work



Histories, workflows, visualizations and *pages* can be shared with others or published to the world.

http://usegalaxy.org/u/aun1/p/windshield-splatter

🗧 Galaxy

⊕ 🗗

Using

Published Pages | aun1 | Windshield Splatter

Windshield splatter analysis with the Galaxy metagenomic pipeline: A live supplement

SERGEI KOSAKOVSKY POND^{1,2,*}, SAMIR WADHAWAN^{3,6*}, FRANCESCA CHIAROMONTE⁴, GURUPRASAD ANANDA^{1,3}, WEN-YU CHUNG^{1,3,7}, JAMES TAYLOR^{1,5}, ANTON NEKRUTENKO^{1,3} and THE GALAXY TEAM^{1*}

Correspondence should addressed to SKP, JT, or AN.

How to use this document

This document is a live copy of supplementary materials for <u>the manuscript</u>. It provides access to the **exact** analyses and workflows discussed in the paper, so you can play with them by re-running, changing parameters, or even applying them to your own data. Specifically, we provide the two histories and one workflow found below. You can view these items by clicking on their name to expand them. You can also import these items into your Galaxy workspace and start using them; click on the green plus to import an item. To import workflows you must <u>create a Galaxy account</u> (unless you already have one) – a hassle-free procedure where you are only asked for a username and password.

This is the Galaxy history detailing the comparison of our pipeline to MEGAN:



This is the Galaxy history showing a generic analysis of metagenomic data. (This corresponds to the "A complete metagenomic pipeline" section of the manuscript and Figure 3A):



Galaxy History | metagenomic analysis



Galaxy Workflow | metagenomic analysis
 Generic workflow for performing a metagenomic analysis on NGS data.
 Generic workflow for performing a metagenomic analysis on NGS data.
 Generic workflow for performing a metagenomic analysis on NGS data.
 Generic workflow for performing a metagenomic analysis on NGS data.
 Generic workflow for performing a metagenomic analysis on NGS data.
 Generic workflow for performing a metagenomic analysis on NGS data.
 Generic workflow for performing a metagenomic analysis on NGS data.
 Generic workflow for performing a metagenomic analysis on NGS data.

Accessing the Data

Windshield Splatter datasets analyzed in this manuscript can be accessed through this Galaxy Library. From

http://usegalaxy.org/u/aun1/p/windshield-splatter





aun1

Related Pages

About this Page

All published pages Published pages by aun1

Rating

Community (6 ratings, 5.0 average)





>

Overview

What is Galaxy?

Analysis Interface, Tools, and Datasources

Workflows

Visualizations

Sharing, Publishing, and Galaxy Pages

ToolShed



Enables sharing of Galaxy Utilities: tools proprietary datatypes exported Galaxy workflows

Automatically install tools and tool suites, and their dependencies, into a Galaxy instance

Galaxy Utilities can be created and shared by any member of the community

https://wiki.galaxyproject.org/ToolShed



private Galaxy installations

Galaxy Tool Shed

Repositories Help- User

2464 valid tools on May 27, 2014

Search

- Search for valid tools
- Search for workflows

Valid Galaxy Utilities

- <u>Tools</u>
- Custom datatypes
- <u>Repository dependency definitions</u>
- Tool dependency definitions
- All Repositories
- Browse by category

Available Actions

٢

Login to create a repository

search	repository	name, de	scription

Categories

Name	Description	Repositories
Assembly	Tools for working with assemblies	31
ChIP-seq	Tools for analyzing and manipulating ChIP-seq data.	7
Computational chemistry	Tools for use in computational chemistry	18
Convert Formats	Tools for converting data formats	31
Data Managers	Utilities for Managing Galaxy's built-in data cache	4
Data Source	Tools for retrieving data from external data sources	16
Fasta Manipulation	Tools for manipulating fasta data	51
Fastq Manipulation	Tools for manipulating fastq data	21
Genome-Wide Association Study	Utilities to support Genome-wide association studies	1
Genomic Interval Operations	Tools for operating on genomic intervals	37
Graphics	Tools producing images	22
Imaging	Utilities to support imaging	
Metabolomics	Tools for use in the study of Metabolomics	3
Metagenomics	Tools enabling the study of metagenomes	25
Micro-array Analysis	Tools for performing micro-array analysis	7
Next Gen Mappers	Tools for the analysis and handling of Next Gen sequencing data	58
Ontology Manipulation	Tools for manipulating ontologies	7
Phylogenetics	Tools for performing phylogenetic analysis	6
Proteomics	Tools enabling the study of proteins	50

toolshed.g2.bx.psu.edu

Installing from the ToolShed

Requires Admin access on Galaxy Server Point and Click to Install Minimal additional configuration needed Dependency directory

Accessible Galaxy tool sheds	
Galaxy main tool shed 👻	
Galaxy test tool shed -	
Browse valid repositories	-
Search for valid tools	
Search for workflows	

Search repositories for valid tools

Valid tools are those that properly load in Galaxy. Enter any combination of the following tool attributes to find repositories that contain valid tools matching the search criteria.

Comma-separated strings may be entered in each field to expand search criteria. Each field must contain the same number of comma-separated strings if these types of search strings are entered in more than one field.

Tool id:

Tool name: bwa

Tool version:

Exact matches only:

Θ

Check the box to match text exactly (text case doesn't matter as all strings are forced to lower case).

Search repositories

Repositories with matching tools

tool id:

tool name: bwa tool version: exact matches only: False Synopsis Revision **Repository** name Owner bwa base 35f793d775ea bwa_base 👻 greg Install to Galaxy bwa_color 17da959633b0 greg Galaxy BWA short read aligner 3365c42b076d devteam bwa_mappers -Galaxy BWA short read aligner 95d5133dd241 devteam bwa_mappers -Demonstrates a tool_dependencies.xml complex_repository_dependency_on_bwa_059 file that defines a complex repository 96a15cf82c50 greg dependency on bwa_059 data_manager_bwa_index_builder blankenberg Builds bwa indexes fe6508204acc For 0 selected items: Install to Galaxy

The core Galaxy development team does not maintain the contents of many Galaxy tool shed repositories. Some repository tools may include code that produces malicious behavior, so be aware of what you are installing.

If you discover a repository that causes problems after installation, contact <u>Galaxy support</u>, sending all necessary information, and appropriate action will be taken.

Contact the repository owner for general questions or concerns.

Confirm dependency installation

These dependencies can be automatically handled with the installed repository, providing significant benefits, and Galaxy includes various features to manage them.

Handle repository dependencies?

☑ Un-check to skip automatic installation of these additional repositories required by this repository.

Repository dependencies – requires bwa color

Handle tool dependencies?

Un-check to skip automatic handling of these tool dependencies.

Tool dependencies – repository tools require handling of these dependencies

Choose the tool panel section to contain the installed tools (optional)

Shed tool configuration file:

/home/g2test/var/shed_tool_conf.xml ‡

Your Galaxy instance is configured with 1 shed-related tool configuration file, so repositories will be installed using it's tool_path setting.

Add new tool panel section:

Add a new tool panel section to contain the installed tools (optional).

Select existing tool panel section:

- Get Data
- BEDtools
- SNP Eff
- Send Data
- ENCODE Tools

Monitor installing tool shed repositories					
Name [Description	Owner	Revision	Status	
bwa base b	bwa_base	greg	35f793d775ea	New	
bwa_color b	bwa_color	greg	17da959633b0	Setting tool versions	

Installed tool shed repositories

Update tool shed status

search repository name Q

Advanced Search

<u>Name</u> ↓	Description	Owner	Revision	Installation Status	Tool shed
✓ allele_counts_1 +	Parse the output of the Naive Variant Detector tool, count alleles	nick	db6f217dc45a	Installed	testtoolshed
allele_counts_1	Parse the output of the Naive Variant Detector tool, count alleles	nick	ecalea054d0d	Installed	testtoolshed
🕗 bwa_base 👻	bwa_base	greg	35f793d775ea	Installed	testtoolshed
🕗 bwa_color 🚽	bwa_color	greg	17da959633b0	Installed	testtoolshed

Goals

What is Galaxy

Mission, Terminology, and Capabilities Using Galaxy Play along with a basic Analysis How to get involved with Galaxy The community and how to participate Adding Tools Making your code Accessible and Reproducible

Bonus: Brief Intro to Sequence Analysis in Galaxy

Basic Analysis

Which 5 coding exons have most overlapping SNPs? Use Human, hg19, Chromosome 22

gcc2015-X.dblankenberg.org
(where X is 1 through 3)

(~ http://usegalaxy.org/galaxy101)

Exons & Repeats: A General Plan

- Get some data
 - Get Data → UCSC Table Browser
- Identify which exons have SNPs
- Count SNPs per exon
- Visualize, save, download, ... exons with most SNPs

(~ http://usegalaxy.org/galaxy101)





SNPs

(Identify which exons have SNPs)









Operate on Genomic Intervals \rightarrow Join (Identify which exons have SNPs)







Overlap pairings





Join, Subtract, and Group → Group (Count SNPs per exon)

Find Top 5

Sort by counts

Filter and Sort -> Sort

Select top 5

Text Manipulation -> Select First





Exon overlap counts

Exons

We've answered our question, but we can do better. Incorporate the overlap count with rest of Exon information







Join on exon name

Join, Subtract, and Group \rightarrow Join

(Incorporate the overlap count with rest of Exon information)







Text Manipulation \rightarrow Cut

(Incorporate the overlap count with rest of Exon information)

Visualize within a Genome Browser

Click display at UCSC link Zoom, scroll, examine

Do it again on Repeats

Extract Workflow

Edit Workflow to include input labels, rename outputs

Get Genome-wide repeat and exon information

Run Workflow

Drink Coffee

Goals

What is Galaxy

Mission, Terminology, and Capabilities Using Galaxy Play along with a basic Analysis How to get involved with Galaxy The community and how to participate Adding Tools Making your code Accessible and Reproducible

Bonus: Brief Intro to Sequence Analysis in Galaxy

The Galaxy Team







Dannon Baker



Dan Blankenberg



Dave Bouvier



Marten Čech

John Chilton



Dave Clements



Nate Coraor



Carl Eberhard



Jeremy Goecks



Sam Guerler



Jen Jackson



Ross Lazarus



Nick Stoler







Nitesh Turaga

Anton Nekrutenko



James Taylor





The Extended Galaxy Team & Contributors



Peter Cock (TJHI)



Greg Von Kuster (PSU)



Eric Rasche (CPT)





Yousef Kowsar (VLSCI)

Future You?



Nicola Soranzo (TGAC)



Nuwan Goonasekera (VeRSI)



Björn Grüning (Uni Freiburg)

Questions about Using Galaxy



https://biostar.usegalaxy.org/

How to Help

Report your bugs and request Features galaxyproject.org/trello galaxy-bugs@lists.galaxyproject.org Contribute (check trello) https://github.com/galaxyproject write and publish tools Edit wiki.galaxyproject.org

Chat

IRC Channel

Server: irc.freenode.net Channel #galaxyproject Twitter: @galaxyproject #usegalaxy

Contributing Code

All code is add via Pull Requests

Info for contributors

https://github.com/galaxyproject/galaxy/blob/dev/ CONTRIBUTING.md

Contributing

This document briefly describes how to contribute to the core galaxy project - also checkout our 2013 Galaxy Community Conference presentation on the topic (video, presentation). For information on contributing more broadly to the Galaxy ecosystem and a deeper discussion of some of these points - please see the Develop section of the Galaxy Wiki.

Before you Begin

If you have an idea for a feature to add or an approach for a bugfix - it is best to communicate with Galaxy developers early. The most common venue for this is the Galaxy Trello board. Browse through existing cards here and if one seems related comment on it. If no existing cards seem appropriate, a new issue can be opened using this form. Galaxy developers are also generally available via IRC and on the development mailing list.

.

The Intergalactic Utilities Commission (IUC)

Community Group tasked with Developing Tool Best Practices Policing the ToolShed

http://galaxy-iuc-standards.readthedocs.org/en/latest/best_practices.html



Goals

What is Galaxy

Mission, Terminology, and Capabilities Using Galaxy Play along with a basic Analysis How to get involved with Galaxy The community and how to participate Adding Tools Making your code Accessible and Reproducible

Bonus: Brief Intro to Sequence Analysis in Galaxy

Creating Your Own Tools

The Problem

You have written a Python script to analyze genomic data and you want to share it with command-line averse colleagues

The Galaxy Solution

Solution: Integrate the script as a new Tool into your own Galaxy server

Steps:

- Obtain and install Galaxy source code (GetGalaxy.org)
- Write an XML file describing the inputs and outputs and how to execute the script
- Instruct Galaxy to load the tool

Adding your Own

Write or download a command-line executable

Determine number and kind of

- Input and Output Datasets
- Input Parameters

Construct a descriptive tool configuration XML file

Write a wrapper script, only if required

A Basic Tool





Compute GC content

</tool>

Source file:

1: Uploaded FASTA File

Execute)

This tool computes GC content from a FASTA file.

\$





cluster.xml Cluster <tool id="gops cluster 1" name="Cluster"> <description>[[Cluster]] the intervals of a query</description> Cluster intervals of: 6: UCSC Main on Human: knownGene -3 <command interpreter="python2.4"> gops cluster.py \$input1 \$output -1 \$input1 chromCol,\$input1 startC max distance between 1 5 -d \$distance -m \$minregions -o \$returntype intervals: (bp) 6 </command> 7 min number of <inputs> 8 <param format="interval" name="input1" type="data"> intervals per cluster: 9 <label>Cluster intervals of</label> Return type: Merge clusters into single intervals -10 </param> 11 <param name="distance" size="5" type="integer" value="1" help="(bp</pre> Execute 12 <label>max distance between intervals</label> 13 </param> 14 <param name="minregions" size="5" type="integer" value="2"> TIP: If your query does not appear in the pulldown menu -> it is not in 15 <label>min number of intervals per cluster</label> interval format. Use "edit attributes" to set chromosome, start, end, and 16 </param> strand columns 17 <param name="returntype" type="select" label="Return type"> 18 <option value="1">Merge clusters into single intervals</option> 19 Screencasts! <option value="2">Find cluster intervals; preserve comments and 20 <option value="3">Find cluster intervals; output grouped by clus See Galaxy Interval Operation Screencasts (right click to open this link in 21 <option value="4">Find the smallest interval in each cluster</op</pre> another window). 22 <option value="5">Find the largest interval in each cluster</opt</pre> 23 </param> 24 </inputs> Syntax 25 <help> 26 Maximum distance is greatest distance in base pairs allowed between .. class:: infomark 27 intervals that will be considered "clustered". Negative values for 28 distance are allowed, and are useful for clustering intervals that overlap. **TIP:** If your query does not appear in the pulldown menu -> it is n · Minimum intervals per cluster allow a threshold to be set on the 29 30 minimum number of intervals to be considered a cluster. Any area with 31 less than this minimum will not be included in the ouput. ---- Merge clusters into single intervals outputs intervals that span the **Screencasts!** entire cluster. Find cluster intervals; preserve comments and order filters out non-cluster intervals while maintaining the original ordering and See Galaxy Interval Operation Screencasts (right click to open this) comments in the file. Find cluster intervals; output grouped by clusters filters out .. Screencasts: http://www.bx.psu.edu/cgi-bin/trac.cgi/wiki/GopsDesc

38

39

40

41

42

Syntax

Line: 87 Column: 8 D XML

Maximum distance is greatest distance in base pairs allowed betw

Minimum intervals per cluster allow a threshold to be set on the

45 - **Merge clusters into single intervals** outputs intervals that span 46 - **Find cluster intervals; preserve comments and order** filters out

Find cluster intervals: output grouped by clusters filters out n

OF Tabs: 2 + −

4 1

 Find cluster intervals; output grouped by clusters filters out non-cluster intervals, but outputs the cluster intervals so that they are grouped together. Comments and original ordering in the file are lost.


External Display Application

<datatype extension="bam" type="galaxy.datatypes.binary:Bam"
 mimetype="application/octet-stream" display_in_upload="true">
 <display file="ucsc/bam.xml" />
</datatype>



BAM at UCSC



Goals

What is Galaxy

Mission, Terminology, and Capabilities Using Galaxy Play along with a basic Analysis How to get involved with Galaxy The community and how to participate Adding Tools Making your code Accessible and Reproducible

Bonus: Brief Intro to Sequence Analysis in Galaxy

Using Galaxy to Analyze NGS Data

- Prepare, quality control and manipulate reads
- Read Mapping
- ChIP-Seq: Binding sites analysis and peak calling

Prepare and Quality Check



Blankenberg D, Gordon A, Von Kuster G, Coraor N, Taylor J, Nekrutenko A; Galaxy Team. Manipulation of FASTQ data with Galaxy. Bioinformatics. 2010 Jul 15;26(14):1783-5.

What is **FASTQ**?

Specifies sequence (FASTA) and quality scores (PHRED)

• Text format, 4 lines per entry



• FASTQ is such a cool standard, there are 3 (or 5) of them!

SSSSSSSSSSSSSSSS	SSSSSSSSSSS	SSS	SSSSSSS	SSSSS	SSSS	SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS	SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS	SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS
• • • • • • • • • • • • • • •	• • • • • • • • • • • •		IIIII	IIIII	IIII	IIIIIIIIIIIIIIIIIIII	IIIIIIIIIIIIIIIII	IIIIIIIIIIIIIII
		XXX	XXXXXXXX	XXXXX	XXXX	******	******	******
! "#\$%&'()*+,/(0123456789:;	;<=>	>?@ABCDI	EFGHI	JKLMI	NOPQRSTUVWXYZ[\]^	`abcdefghijklmno	pqrstuvwxyz{ }~
33	59)	64	73			104	126
S - Sanger	Phred+33,	93	values	(0,	93)	(0 to 60 expected	in raw reads)	
I - Illumina 1.3 X - Solexa	Phred+64, Solexa+64,	62 67	values values	(0, (-5,	62) 62)	(0 to 40 expected (-5 to 40 expected	in raw reads) d in raw reads)	

http://en.wikipedia.org/wiki/FASTQ_format

Combining Sequences and Qualities

	Galaxy			Analyze Data	Workflow	Shared Data	Visualization	Admin H	lelp User					
Tools		Options 🔻	Combine FAST								History	0	Options 🔻	0
 <u>FAS</u> end 	TQ splitter on joine I reads	ed paired	FASTA File:	A and QUAL							Combine QUA	L and Seque	nce	r
= <u>FAS</u> rea	TQ joiner on paired	l end	1: 454.fasta	•							2: 454.qual		• / %	
<u>FAS</u> colu	TQ Summary Statis	<u>tics</u> by	Quality Score 2: 454.qual	File:							format: qual45	4, database: qual454 file	2	
RO	CHE-454 DATA		Force Quality	Score encoding:							60		Ø 🖻	
 Buil Sele 	ld base quality distr	ibution ments	Execute								>EYKX4VC01B65G 33 23 34 25 28	S length=54 > 28 28 32 23	xy=0784_1 34 27 4	
E Cor	nbine FASTA and Q	UAL into	What it does								>EYKX4VC01BNCS 27 35 26 25 37 22 9 23 19 28	P length=187 28 37 28 25 28 28 28 26	xy=0558_ 28 27 36 28 39 32	
AB	-SOLID DATA		₄ This tool joins a	FASTA file to a Qua	ality Score file	e, creating a sing	le FASTQ block fo	or each read.		1	26 27 37 29 28	26 28 36 28	26 24 38	
• <u>Cor</u>	<u>ivert</u> SOLiD output t	to fastq	Specifying a set of csfasta is provide	of quality scores is ed) with each quali	optional; whe	en not provided, g the maximal a	, the output will be llowed value (93).	e fastqsanger (or fastqcssanger (w	vhen a			● / ×	
SOL	iD data		Use this tool, for	example, to conve	ert 454-type	output to FASTQ	l.				1: 454.fasta 18 sequences format: fasta	database: 7	• / 23	
SOL	W quality score box ID data @EYK	x4VC01B65GS	length=54 xy=078	4 1754 region=	1 run=R 20	07 11 07 16 3	15 57				Info: uploaded	fasta file	$D \square$	
GEI	NERIC FASTC CCGG NIPULATION +	TATCCGGGTGCC	GTGATGAGCGCCACCO	GAACGAATTCGACT	ATGCCGAA						>EYKX4VC01B65G	S length=54 ;	xy=0784_1	
Filtersco	er FASTQ rea @EYK re and length CTTA	X4VC01BNCSP CCGGTCACCACC	length=187 xy=05 GTGCCTTCAGGATTGA	58_3831 region	=1 run=R_2 GGTGCGTCAG	007_11_07_16 GCGGGGGTGACAT	15_57_ CGCCCACCACGGTA	CTCACTGGCTG	GCTCTGGTTCCCGGG	CGGCATCGGAG	CCGGTATCCGGGTG >EYKX4VC01BNCS	CCGTGATGAGCGG P length=187	CCACCGGAA	
= <u>FAS</u>	TQ Trimmer <d;:< td=""><td>F=F=:=<e<=e< X4VC01CD9FT</e<=e< </td><td>==<e<?4<=e=8e<<= length=115 xy=08</e<?4<=e=8e<<= </td><td><<=F><;<99E<;=) 65 1719 region</td><th>E=9:6=9=;C =1 run=R 2</th><td>:;LE7*84==== 007 11 07 16</td><td>;=HA-<e==;f==; 15 57</e==;f==; </td><td>===<=; E<<<e< td=""><td>E=<<==E<e=ha-d=;< td=""><td>; F>===F>=E</td><td>GTTACCGGTCACCA GGTGACATCGCCCA</td><td>CCGTGCCTTCAG</td><td>GATTGATCG</td><td></td></e=ha-d=;<></td></e<></td></d;:<>	F=F=:= <e<=e< X4VC01CD9FT</e<=e< 	== <e<?4<=e=8e<<= length=115 xy=08</e<?4<=e=8e<<= 	<<=F><;<99E<;=) 65 1719 region	E=9:6=9=;C =1 run=R 2	:;LE7*84==== 007 11 07 16	;=HA- <e==;f==; 15 57</e==;f==; 	===<=; E<< <e< td=""><td>E=<<==E<e=ha-d=;< td=""><td>; F>===F>=E</td><td>GTTACCGGTCACCA GGTGACATCGCCCA</td><td>CCGTGCCTTCAG</td><td>GATTGATCG</td><td></td></e=ha-d=;<></td></e<>	E=<<==E <e=ha-d=;< td=""><td>; F>===F>=E</td><td>GTTACCGGTCACCA GGTGACATCGCCCA</td><td>CCGTGCCTTCAG</td><td>GATTGATCG</td><td></td></e=ha-d=;<>	; F>===F>=E	GTTACCGGTCACCA GGTGACATCGCCCA	CCGTGCCTTCAG	GATTGATCG	
= <u>FAS</u> slid	ing window +	GCTTTGGCCTGT	CGTCCGGCACCTCGC	AGAGCTACAGCAGG	CGCGGCTGGC	GATCATCGGCGG	CACGCCGGCCTATA	TGTCGCCGGAA	CACACCACCCGCACC	CCAACGCG		TATICCCCTCGG) 4 >	
= FAS	TO Masker b @EYK TAAA	X4VC01B8FW0 TTTCAAGGAATG	length=95 xy=079 CAAATCAGGGTCGTGT	9_0514 region= GTTTAGACTTCGGC	1 run=R_20 TTTAGAGACC	07_11_07_16_3 TGAATACGTCAA	AAACATAACTTCAT	GATATCTTGCA	г. Абт	189=				
	+ =IC0 @EYK GGCC	D=' <b8c9a7== X4VC01BCGYW AGCCGGGACAGC</b8c9a7== 	=JC2===F?*===== length=115 xy=04 GTTGTTGGGCTGCATC	=F?)==<=D; <d;= 34_3926 region GCGACGAGCTAAAA</d;= 	F?*=<===C: =1 run=R_2 GTCGCCATCA	==A7;==== <le0 007_11_07_16 .ccgccccgccgg</le0 	8-"=6=<1=A8<=< _15_57_ TTGATGGGCAGGCT	==<===A7=;;	<= TGGTAAAAACTTTCTC	CGCCAAAC				
	=';0 @EYK GGGG	<=F=JD2=6=86 X4VC01AZXC6 GCGTTTGGCCTG	<e<9e=ic 7:="9<=E<br">length=116 xy=02 TCGTCCGGCACCTCGC</e<9e=ic>	'=;=<<==== <le7) 92_0280 region AAGAGCTACAGCAG</le7) 	=;=<;/=:5= =1 run=R_2 GCGCGGCTGG	C9:IB3"4<1E=1 007_11_07_16 CGATCATCGGCGG	E=6<:JC17=F>;; _15_57_ GCACGCCGGCCTATA	D<=;JCl==<=	F>:LE8-",HA-=29	5==2E>(9 CCCAACGCG				

Grooming --> Sanger

- Galaxy	Analyze Data Workflow Shared Data	a Visualization Admin Help User		
Tools Options				History Options 🔻
 NGS TOOLBOX BETA NGS: QC and manipulation ILLUMINA DATA FASTQ Groomer convert between various FASTQ quality formats FASTQ splitter on joined paired end reads FASTQ splitter on paired end reads FASTQ joiner on paired end reads FASTQ Summary Statistics by column Build base quality distribution Select high quality segments Combine FASTA and QUAL into 	A and and data 2 ality scores type: ger veral conversions options relating to the FA options, the output will be <i>sanger</i> formatter f a quality score falls outside of the target haximum).	4: FASTQ Groomer on data 3 18 sequences format: fastqsanger, database: ? Info: Groomed 18 sanger reads into sanger reads. Based upon quality and sequence, the input data is valid for: sanger Input ASCII range: '!'(33) - 'L'(76) Input decimal range: 0 - 43 Image: Imag	Space able value (i.e.	Combine QUAL and Sequence 3: Combine FASTA and OND 3: Combine FASTA and OND OUAL on data 1 and data 2 18 sequences format: fastqsanger, database: ? Info: Combined 18 of 18 sequences with quality scores (100.00%). Combined 18 of 18 sequences with quality scores (100.0
FASTQ When converting b AB-SOLID DATA Generat SOLiD output to fastq Convert SOLiD output to fastq When converting b Compute quality statistics for SOLiD data When converting b Draw quality score boxplot for SOLiD data When converting b GENERIC FASTQ Quality Score Comparing	etween Solexa and the other formats, qual d in <u>Cock PJ, Fields CJ, Goto N, Heuer ML, I</u> the Solexa/Illumina FASTQ variants. Nucle etween color space (csSanger) and base/se if gained, the base 'G' is used as the adapt. ent in the color space sequence. Any maske ning color space encoding.	@EYKX4VC01BNCSP length=187 xy=0558_ CTTACCGGTCACCACCGTGCCTTCAGGATTGATCG O Image: Comparison of the state o	scales using ces with adapter bases space if there be converted to	2: 454.qual ● Ø X 52 lines format: qual454, database: ? Info: uploaded qual454 file ● Im ● Ø E >EYKX4VC01B65GS length=54 xy=0784_1 33 23 34 25 28 28 28 32 23 34 27 4 >EYKX4VC01BNCSP length=187 xy=0558_ 27 35 26 25 37 28 37 28 37 28 25 28 27 36
SSSSSSSSSSSSSSSSSSSSSSSSSSS I"#\$%&'()*+,/01234567 33 S - Sanger Phred+3 I - Illumina 1.3 Phred+6 X - Solexa Solexa+	SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS 	SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS	SSSSSSSSSSSS IIIIIIIII xxxxxxxxxxx stuvwxyz{ }~ 126	

Diagram adapted from http://en.wikipedia.org/wiki/FASTQ_format

NGS TOOLBOX BETA

NGS: QC and manipulation

ILLUMINA DATA

- FASTQ Groomer convert between various FASTQ quality formats
- FASTQ splitter on joined paired end reads
- FASTQ joiner on paired end reads
- FASTQ Summary Statistics by column

Score Value

9

2

Quality Statistics and Box Plot Tool

Graph/Display Data

- Histogram of a numeric column
- Scatterplot of two numeric columns
- Plotting tool for multiple series and graph types
- Boxplot of quality statistics



- - **-*-* - **--*-* *--- **-- **-

**--* -*-



Nucleotide Position

FastQC



Read Trimming

3

🖥 Galaxy	Analyze Data Workflow Shared Data Visualization Admin Help User
ools Options	
GENERIC FASTQ MANIPULATION Filter FASTQ reads by quality score and length FASTQ Trimmer by column FASTQ Quality Trimmer by sliding window FASTQ Masker by quality score Manipulate FASTQ reads on various attributes FASTQ to FASTA converter FASTQ to FASTA converter FASTQ to Tabular converter Tabular to FASTQ converter FASTX-TOOLKIT FOR FASTQ DATA Quality format converter (ASCII- Numeric) Compute quality statistics Draw quality score boxplot	FASTQ Trimmer FASTQ File: 2: imported: GM12878ple Dataset : Define Base Offsets as: Absolute Values Use Absolute Values Use Absolute for fixed length reads (Illumina, SOLip Use Percentage for variable length reads (Roche/45: Offset from 5' end: 0 Values start at 0, increasing from the left Offset from 3' end: 16 Values start at 0, increasing from the right Keep reads with zero length: Execute This tool allows you to trim the ends of reads. You can specify either absolute or percent-based offs You can specify either absolute or percent-based offs 0 0 0 0 0 0 0 0 0 0 0
 <u>Draw nucleotides distribution</u> <u>chart</u> <u>FASTQ to FASTA</u> converter <u>Filter by quality</u> 	For example, if you have a read of length 36: <pre></pre>
Remove sequencing artifacts	And you get absolute offects of 2 and 0: Quality Score:
	0.0 Execute

Filter FASTQ

FASTQ File:

7: FASTQ Trimmer on data 2	\$	è
----------------------------	----	---

Requires groomed data: if your data does not appear here try using the FASTQ groomer.

Minimum Size:

1.1	-	
	•	

Maximum Size:

0		
υ		

A maximum size less than 1 indicates no limit.

Minimum Quality:

1	n		n	í.	
4	Υ	•	Y		

Maximum Quality:

-	~
0	U

A maximum quality less than 1 indicates no limit.

Maximum	number	of	bases	allowed	outside	of	quality	range
Maximum	numper	01	Dases	anoweu	outside	01	quanty	range

0	
υ	

This is paired end data:

-	_	_	_	-	
c					١.
ε.					
					E.
L.				_	ε.

Quality Filter on a Range of Bases

Add new Quality Filter on a Range of Bases

Execute

Quality Filter on a Range of Bases

Quality Filter on a Range of Bases 1

Define Base Offsets as:



Use Absolute for fixed length reads (Illumina, SOLID) Use Percentage for variable length reads (Roche/454)

+

Offset from 5' end:

\sim		
0		

0

Values start at 0, increasing from the left

Offset from 3' end:

Values start at 0, increasing from the right

Aggregate read score for specified range:



Keep read when aggregate score is:



Quality Score:

0.0

Remove Quality Filter on a Range of Bases 1

Add new Quality Filter on a Range of Bases

Execute

Manipulate FASTQ

Manipulate FASTQ

FASTQ File:

7: FASTQ Trimmer on data 2

Requires groomed data: if your data does not appear here try using the FASTQ groomer.

\$

Match Reads

Add new Match Reads

Manipulate Reads

Add new Manipulate Reads

Execute

Mani	pula	te F	ASTO

FASTQ File:

7: FAST	Q Trimme	er on o	da	ta 2		÷		
equires	groomed	data:	if	your	data	does	not	

appear here try using the FASTQ groomer.

Match Reads

	Match Reads
	Match Reads
t 🛟	Sequence C
ype:	Sequence M
on 🛟	Regular Exp
	Match by:
	N
ads 1	Remove Ma
ds	Add new Mate
	lanipulate Re
Reads	Add new Man
	Execute
e Reads	Match by: N Remove Ma Add new Mate Ianipulate Re Add new Man Execute

STQFI	le:	
: FASTO	2 Trimmer on data 2	v da ta a la c
quires near he	groomed data: if your	data does
tch Re	ads	g groomer
Match	Reads 1	
Match	Reads 1	
Match	Reads by:	
Seque	ance Content	
Sequer	nce Match Type:	
Regu	ar Expression 💲	
Match	hv	
Match	wy.	
N Remo	ove Match Reads 1	
N Remo	v Match Reads 1	
N Remo Add nev	v Match Reads 1 v Match Reads	
N Remo Add nev anipula Manipu	v Match Reads 1 v Match Reads te Reads ulate Reads 1	
N Remo Add nev anipula Manipu	v Match Reads 1 v Match Reads te Reads ulate Reads 1 ulate Reads on:	
N Remo Add nev Anipula Manipu Manipu Misce	v Match Reads 1 v Match Reads te Reads ulate Reads 1 ulate Reads on: Ilaneous Actions \$	
N Remo Add nev Anipula Manipu Manipu Miscel	v Match Reads 1 v Match Reads te Reads ulate Reads 1 ulate Reads on: Ilaneous Actions \$	1 Туре:
N Remo Add nev Anipula Manipu Manipu Misce Miscel Remo	ove Match Reads 1 v Match Reads ute Reads ulate Reads 1 ulate Reads on: Ilaneous Actions \$ laneous Manipulation ve Read \$	n Type:

Execute

Using Galaxy to Analyze NGS Data

- Prepare, quality control and manipulate reads
- Read Mapping
- ChIP-Seq: Binding sites analysis and peak calling

Mapping NGS Data

Collection of interchangeable mappersaccept fastq format, produce SAM/BAM

Mappers for

- + DNA
- + RNA

Mappers

DNA

- short reads: Bowtie, BWA, BFAST, PerM
- Ionger reads: LASTZ

Metagenomics

Megablast

RNA / gapped-reads mapper

Tophat

Commonly Used/Default Parameters

Lastz
Align sequencing reads in:
Against reference sequences that are:
locally cached 🗘
Using reference genome:
Aedes aegypti: AaegL1 🔍
If your genome of interest is not listed, contact the Galaxy team
Output format:
SAM \$
Lastz settings to use:
Commonly used
For most mapping needs use Commonly used settings. If you want full control use Full List
Select mapping mode:
Roche-454 98% identity 🛟
Roche-454 98% identity
Roche-454 95% identity
Roche-454 85% identity
Roche-454 75% identity v this identity (%):
Illumina 95% identity
Do not report matches above this identity (%):
100
Do not report matches that cover less than this percentage of each read
0
Convert lowercase bases to uppercase:
Yes 🗘
Execute

Lastz

Align sequencing reads in: 53: FASTQ to FASTA on data 7

Against reference sequences that are:

+

locally cached

Using reference genome:

Aedes aegypti: AaegL1

If your genome of interest is not listed, contact the Galaxy team

0	u	t	p	u	t	f	0	r	m	a	t:	
_	_	_		_	_	-	_	۰.		_		

SAM

Lastz settings to use:

Full Parameter List 📫

use Commonly used settings. If you want full control use Full List Commonly used

+

Full Parameter List Both

Select seeding settings:

Seed hits require a 19 bp word with matches i

allows you set word size and number of mismatches

Select transition settings:

Allow one transition in each seed hit +

affects the number of allowed transition substitutions

Perform gap-free extension of seed hits to HSPs (high scoring segment pairs)?:

No ‡

Perform chaining of HSPs?:

No ‡

Gap opening penalty:

400

Gap extension penalty:

30

X-drop threshold:

910

Y-drop threshold:

9370

Set the threshold for HSPs (ungapped extensions scoring lower are discarded):

3000

Set the threshold for gapped alignments (gapped extensions scoring lower are discarded):

3000

Involve entropy when filtering HSPs?:

No ‡

No ‡

Do you want to modify the reference name?:

Full Parameter List

No 1
Do not report matches below this identity (%):
0
Do not report matches above this identity (%):
100
Do not report matches that cover less than this percentage of each read
0
Convert lowercase bases to uppercase:
Yes ÷
Execute

What it does

LASTZ is a high performance pairwise sequence aligner derived from BLASTZ. It is written by Bob Harris in Webb Miller's laboratory at Penn State University. Special scoring sets were derived to improve runtime performance and quality. This Galaxy version of LASTZ is geared towards aligning short (Illumina/Solexa, AB/SOLiD) and medium (Roche/454) reads against a reference sequence. There is excellent, extensive documentation on LASTZ available here.

Input formats

LASTZ accepts reference and reads in FASTA format. However, because Galaxy supports implicit format conversion the tool will recognize fastg and other method specific formats.

Using Galaxy to Analyze NGS Data

- Prepare, quality control and manipulate reads
- Read Mapping
- ChIP-Seq: Binding sites analysis and peak calling

Peak Calling / ChIP-seq analysis

Punctate binding

transcription factors

Diffuse binding

- histone modifications
- + Polli

Punctate Binding --> MACS

Inputs

- Enriched Tag file
- Control / Input file (optional)

Outputs

- Called Peaks
- Negative Peaks (when control provided)
- Shifted Tag counts (wig, convert to bigWig for visualization)



Zhang et al. Model-based Analysis of ChIP-Seq (MACS). Genome Biol (2008) vol. 9 (9) pp. R137

Diffuse Binding

00

SICER (version 0.0.1)

ChIP-Seq Tag File:

1: Mapped Tags

ChIP-Seq Control File:

2: Mapped Control 💲

Fix off-by-one errors in output files:

\$

 \checkmark

SICER creates non-standard output files, this option will fix these coordinates

Redundancy Threshold:

1

The number of copies of identical reads allowed in a library

Window size:

200

Resolution of SICER algorithm. For histone modifications, one can use 200 bp

Fragment size:

150

for determination of the amount of shift from the beginning of a read to the center of the DNA fragment represented by the read. FRAGMENT_SIZE=150 means the shift is 75.

Effective genome fraction:

0.74

Effective Genome as fraction of the genome size. It depends on read length.

Gap size:

600

Needs to be multiples of window size. Namely if the window size is 200, the gap size should be 0, 200, 400, 600, ...

Statistic threshold value:

0.01

FDR (with control) or E-value (without control)

Execute

Zang C, Schones DE, Zeng C, Cui K, Zhao K, Peng W. A clustering approach for identification of enriched domains from histone modification ChIP-Seq data. Bioinformatics. 2009 Aug 1;25(15):1952-8. Epub 2009 Jun 8.

SICER



Hands on NGS: ChIP-Seq

This time we will be using Data from a Shared Data Library

gcc2015-X.dblankenberg.org
 (where X is 1 through 3)