

An Introduction to Galaxy

Daniel Blankenberg
The Galaxy Team
<http://UseGalaxy.org>

Overview

What is Galaxy?

What **you** can do in Galaxy

- ✦ analysis interface, tools and datasources
- ✦ data libraries
- ✦ workflows
- ✦ visualization
- ✦ sharing
- ✦ Pages

Galaxy 101 Exercise

The Vision

Galaxy is an **open**, Web-based platform for **accessible, reproducible, and transparent** computational biomedical research

What is Galaxy?

GUI for genomics

- ✦ for complete analyses: analyze, visualize, share, publish

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

Overview

What is Galaxy?

What you can do in Galaxy

- ✦ analysis interface, tools and datasources
- ✦ data libraries
- ✦ workflows
- ✦ visualization
- ✦ sharing
- ✦ Pages

Galaxy 101 Exercise

Galaxy Analysis Workspace

The screenshot displays the Galaxy Analysis Workspace interface. The browser address bar shows <http://main.g2.bx.psu.edu/>. The main navigation bar includes 'Galaxy' and tabs for 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Help', and 'User'. The left sidebar contains a 'Tools' menu with categories like 'Get Data', 'Text Manipulation', 'Filter and Sort', 'Fetch Sequences', 'Operate on Genomic Intervals', 'Statistics', 'Graph/Display Data', 'Regional Variation', 'Multiple regression', 'Multivariate Analysis', 'Evolution', 'Metagenomic analyses', 'EMBOSS', 'NGS TOOLBOX BETA', 'NGS: QC and manipulation', 'NGS: Mapping', 'NGS: SAM Tools', 'NGS: Indel Analysis', 'NGS: Peak Calling', 'RGENETICS', 'SNP/WGA: Data; Filters', 'SNP/WGA: QC; LD; Plots', 'SNP/WGA: Statistical Models', and 'Workflows'. The central workspace shows the 'Map with Bowtie for Illumina' tool configuration. The configuration includes: 'Use a built-in index' selected; 'mm9' selected as the reference genome; 'Paired-end' selected for library mate-pairing; '1: E18 PE.1 Reads' selected for both Forward and Reverse FASTQ files; '1000' entered for maximum insert size; 'FR (for Illumina)' selected for mate orientation; 'Commonly used' selected for Bowtie settings; and the 'Suppress the header in the output SAM file' checkbox checked. An 'Execute' button is at the bottom. Below the configuration is a 'What it does' section describing Bowtie as a short read aligner. The right sidebar shows a 'History' panel with a list of jobs, including '15: Variants from sample E18, consensus different in RefSeq Genes', '14: UCSC mm9 RefSeq_Genes', '13: Variants from sample E18 where consensus base different than ref. base', '10: Variants from sample E18', '9: Generate pileup on data 8', '8: SAM-to-BAM on data 7', '7: Map with Bowtie for Illumina on data 6 and data 5', '6: E18 PE.2 Reads Groomed, Trimmed', '5: E18 PE.1 Reads Groomed, Trimmed', '4: E18 PE.2 Reads Groomed', '3: E18 PE.1 Reads Groomed', '2: E18 PE.2 Reads', and '1: E18 PE.1 Reads'. Each job entry has eye and refresh icons.

Filter and Sort

- Filter data on any column using simple expressions
- Sort data in ascending or descending order
- Select lines that match an expression

GFF FILES

- Extract features from GFF file
- Filter GFF file by attribute using simple expressions
- Filter GFF file by feature count using simple expressions

[Extract Features](#)
[Fetch Sequences](#)
[Fetch Alignments](#)
[Get Genomic Scores](#)
[Operate on Genomic Intervals](#)
[Statistics](#)
[Graph/Display Data](#)
[Regional Variation](#)
[Multiple regression](#)
[Multivariate Analysis](#)
[Evolution](#)
[Metagenomic analyses](#)
[EMBOSS](#)

NGS TOOLBOX BETA

[NGS: QC and manipulation](#)
[NGS: Mapping](#)
[NGS: SAM Tools](#)
[NGS: Indel Analysis](#)
[NGS: Peak Calling](#)

RGENETICS

[SNP/WGA: Data; Filters](#)
[SNP/WGA: QC; LD; Plots](#)
[SNP/WGA: Statistical Models](#)

[Workflows](#)

Galaxy Analysis Workspace

The screenshot displays the Galaxy Analysis Workspace interface. The main panel shows the configuration for the 'Map with Bowtie for Illumina' tool. The configuration includes:

- Reference genome:** mm9
- Library mate-paired?:** paired-end
- Forward FASTQ file:** 1: E18 PE.1 Reads
- Reverse FASTQ file:** 1: E18 PE.1 Reads
- Maximum insert size for valid paired-end alignments (-X):** 1000
- The upstream/downstream mate orientation for valid paired-end alignment against the forward reference strand (--fr/--rf/--ff):** FR (for Illumina)
- Bowtie settings to use:** Commonly used
- Suppress the header in the output SAM file:**

The 'Execute' button is visible at the bottom of the configuration panel. Below the configuration, the 'What it does' section provides a brief description of the Bowtie tool.

The right-hand side of the interface shows the 'History' panel, which lists the workflow steps in a table:

Step ID	Step Name	View	Refresh	Delete
15	Variants from sample E18, consensus different, in RefSeq Genes	👁	🔄	🗑
14	UCSC mm9 RefSeq Genes	👁	🔄	🗑
13	Variants from sample E18 where consensus base different than ref. base	👁	🔄	🗑
10	Variants from sample E18	👁	🔄	🗑
9	Generate pileup on data 8	👁	🔄	🗑
8	SAM-to-BAM on data 7	👁	🔄	🗑
7	Map with Bowtie for Illumina on data 6 and data 5	👁	🔄	🗑
6	E18 PE.2 Reads Groomed, Trimmed	👁	🔄	🗑
5	E18 PE.1 Reads Groomed, Trimmed	👁	🔄	🗑
4	E18 PE.2 Reads Groomed	👁	🔄	🗑
3	E18 PE.1 Reads Groomed	👁	🔄	🗑
2	E18 PE.2 Reads	👁	🔄	🗑
1	E18 PE.1 Reads	👁	🔄	🗑

Filter and Sort

- Filter data on any column using simple expressions
- Sort data in ascending or descending order
- Select lines that match an

Operate on Genomic Intervals

- Intersect the intervals of two queries
- Subtract the intervals of two queries
- Merge the overlapping intervals of a query
- Concatenate two queries into one query
- Base Coverage of all intervals
- Coverage of a set of intervals on second set of intervals
- Complement intervals of a query
- Cluster the intervals of a query
- Join the intervals of two queries side-by-side
- Get flanks returns flanking region/s for every gene
- Fetch closest feature for every interval
- Profile Annotations for a set of genomic intervals

Galaxy Analysis Workspace

The screenshot displays the Galaxy Analysis Workspace interface. The main panel shows a workflow step titled "Bowtie for Illumina". The workflow configuration includes options for selecting a reference genome, indexing, and alignment parameters. The history panel on the right shows a sequence of steps: 1: E18 PE.1 Reads, 2: E18 PE.2 Reads, 3: E18 PE.1 Reads Groomed, 4: E18 PE.2 Reads Groomed, 5: E18 PE.1 Reads Groomed, Trimmed, 6: E18 PE.2 Reads Groomed, Trimmed, 7: Map with Bowtie for Illumina on data 6 and data 5, 8: SAM-to-BAM on data 7, 9: Generate pileup on data 8, 10: Variants from sample E18, 13: Variants from sample E18 where consensus base different than ref. base, 14: UCSC mm9 RefSeq_Genes, and 15: Variants from sample E18, consensus different, in RefSeq Genes.

Filter and Sort

- Filter data on any column using simple expressions
- Sort data in ascending or descending order
- Select lines that match an

Operate on Genomic Intervals

- Intersect the intervals of two queries
- Subtract the intervals of two queries
- Merge the overlapping intervals of a query

NGS: SAM Tools

- Filter SAM on bitwise flag values
- Convert SAM to interval
- SAM-to-BAM converts SAM format to BAM format
- BAM-to-SAM converts BAM format to SAM format
- Merge BAM Files merges BAM files together
- Generate pileup from BAM dataset
- Filter pileup on coverage and SNPs
- Pileup-to-Interval condenses pileup format into ranges of bases

Galaxy Analysis Workspace

The screenshot displays the Galaxy Analysis Workspace interface. The main window shows a workflow titled "SNP Pileup Analysis for Sample E18". The workflow steps are listed in the History panel on the right:

- 1: E18 PE.1 Reads
- 2: E18 PE.2 Reads
- 3: E18 PE.1 Reads Groomed
- 4: E18 PE.2 Reads Groomed
- 5: E18 PE.1 Reads Groomed, Trimmed
- 6: E18 PE.2 Reads Groomed, Trimmed
- 7: Map with Bowtie for Illumina on data 6 and data 5
- 8: SAM-to-BAM on data 7
- 9: Generate pileup on data 8
- 10: Variants from sample E18
- 13: Variants from sample E18 where consensus base different than ref. base
- 14: UCSC mm9 RefSeq_Genes
- 15: Variants from sample E18, consensus different, in RefSeq Genes

The main workspace area shows the configuration for the "Generate pileup" tool. It includes options for selecting a reference genome, indexing, and filtering. The interface is clean and organized, with a navigation bar at the top and a search bar.

Filter and Sort

- Filter data on any column using simple expressions

- Sort data in ascending or descending order

- Select lines that match a query

Operate on Genomes

- Intersect the intervals of two queries
- Subtract the intervals of two queries
- Merge the overlapping intervals of a query

NGS: SAM Tools

- Filter SAM records by values
- Convert SAM records to BAM
- SAM-to-BAM format to BAM
- BAM-to-SAM format to SAM
- Merge BAM files together
- Generate pileup dataset
- Filter pileup on coverage and SNPs
- Pileup-to-Interval condenses pileup format into ranges of bases

Filter pileup

Select dataset:

10: Variants from sample E18

which contains:

Pileup with six columns (simple)

See "Types of pileup datasets" below for examples

Do not consider read bases with quality lower than:

20

No variants with quality below this value will be reported

Do not report positions with coverage lower than:

3

Pileup lines with coverage lower than this value will be skipped

Only report variants?:

Yes

See "Examples 1 and 2" below for explanation

Convert coordinates to intervals?:

No

See "Output format" below for explanation

Print total number of differences?:

No

See "Example 3" below for explanation

Print quality and base string?:

Yes

See "Example 4" below for explanation

Execute

ce

aligner designed to be ultrafast and memory-efficient. It is developed by Ben Langmead, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biology 10:R25.

Filter and Sort

- Filter data on any column using simple expressions

- Sort data in ascending or descending order

- Select lines that match a query

Operate on Genomes

- Intersect the intersection of two queries
- Subtract the intersection of two queries
- Merge the overlap of a query

NGS: SAM Tools

- Filter SAM on values
- Convert SAM to BAM
- SAM-to-BAM format to BAM
- BAM-to-SAM format to SAM
- Merge BAM files together
- Generate pileup dataset
- Filter pileup on coverage and SNPs
- Pileup-to-Interval condenses pileup format into ranges of bases

Filter pileup

Select dataset:

10: Variants from sample E18

which contains:

Pileup with six columns (simple)

See "Types of pileup datasets" below for examples

Do not consider read bases with quality lower than:

20

No variants with quality below this value will be reported

Do not report positions with coverage lower than:

3

Pileup lines with coverage lower than this value will be skipped

Only report variants?:

Yes

See "Examples 1 and 2" below for explanation

Convert coordinates to intervals?:

No

See "Output format" below for explanation

Print total number of differences?:

No

See "Example 3" below for explanation

Print quality and base string?:

Yes

See "Example 4" below for explanation

Execute

aligner designed to be ultrafast and memory-efficient. It is developed by Mark Imbusch and Mark Trapnell. Please cite: Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biol.

History

Options



Variant Analysis for Sample E18

15: Intersect to get Variants from sample E18, consensus different, in RefSeq Genes

14: UCSC mm9 RefSeq Genes

13: Filter to get Variants from sample E18 where consensus base different than ref. base

10: Filter pileup to get Variants from sample E18

9: Generate pileup on data 8

8: SAM-to-BAM on data 7

7: Map with Bowtie for Illumina on data 6 and data 5

6: E18 PE.2 Reads Groomed, Trimmed

5: E18 PE.1 Reads Groomed, Trimmed



This dataset is large and only the first megabyte is shown below.

Show all | Save

chr10	6882036	6882037	A	A	107	0	60	32	0	0	0	0
chr10	14243075	14243076	G	G	107	0	96	0	60	35	0	0
chr10	14243079	14243080	C	C	106	0	106	0	60	35	0	0
chr10	14465082	14465083	T	K	173	176	173	176	60	35	0	0
chr10	14465083	14465084	G	K	144	144	144	144	60	35	0	0
chr10	14465084	14465085	T	T	117	0	117	0	60	38	0	0
chr10	14465085	14465086	G	G	70	0	70	0	60	38	0	0
chr10	14465257	14465258	C	C	79	0	79	0	60	42	0	0
chr10	14465258	14465259	A	A	137	0	137	0	60	46	0	0
chr10	14465263	14465264	A	A	136	0	136	0	60	61	0	0
chr10	14465366	14465367	A	A	101	0	101	0	60	38	0	0
chr10	14465371	14465372	G	G	137	0	137	0	60	50	0	0
chr10	14465410	14465411	G	G	184	0	184	0	60	69	0	0
chr10	14465447	14465448	T	T	186	0	186	0	60	65	0	0
chr10	14465456	14465457	G	G	193	0	193	0	60	70	0	0
chr10	14465465	14465466	T	T	177	0	177	0	60	63	0	0
chr10	14465485	14465486	C	T	129	129	129	129	60	34	0	0
chr10	14465569	14465570	T	T	219	0	219	0	60	84	0	0
chr10	14465581	14465582	G	G	240	0	240	0	60	84	0	0
chr10	14465586	14465587	C	C	248	0	248	0	60	82	0	0
chr10	14465621	14465622	C	C	134	0	134	0	60	49	0	0
chr10	14465658	14465659	C	C	134	0	134	0	60	49	0	0
chr10	14465660	14465661	T	T	153	0	153	0	60	55	0	0
chr10	14465691	14465692	G	G	128	0	128	0	60	42	0	0
chr10	14465778	14465779	C	C	89	0	89	0	60	34	0	0
chr10	14465791	14465792	G	G	104	0	104	0	60	33	0	0
chr10	14465881	14465882	G	G	110	0	110	0	60	41	0	0
chr10	17445088	17445089	A	A	103	0	103	0	60	34	0	0
chr10	17445271	17445272	A	A	55	0	55	0	60	34	0	0
chr10	17731269	17731270	T	T	113	0	113	0	60	42	0	0
chr10	19928287	19928288	G	A	135	135	135	135	60	36	0	0
chr10	19928468	19928469	C	T	132	132	132	132	60	35	0	0
chr10	19928488	19928489	A	A	119	0	119	0	60	44	0	0
chr10	19928494	19928495	C	T	138	138	138	138	60	37	0	0
chr10	19928527	19928528	A	A	134	0	134	0	60	45	0	0
chr10	19928538	19928539	G	G	144	0	144	0	60	52	0	0
chr10	19928543	19928544	A	G	147	147	147	147	60	40	0	0
chr10	19928741	19928742	T	T	80	0	80	0	60	30	0	0
chr10	20799826	20799827	G	G	117	0	117	0	60	37	0	0
chr10	28750217	28750218	C	T	138	138	138	138	60	37	0	0
chr10	28750397	28750398	A	C	154	211	154	211	60	64	0	0
chr10	28750401	28750402	A	A	128	0	128	0	60	47	0	0
chr10	28750423	28750424	C	T	113	113	113	113	60	35	0	0
chr10	28750438	28750439	A	A	95	0	95	0	60	36	0	0
chr10	28750446	28750447	A	G	165	165	165	165	60	46	0	0
chr10	28750487	28750488	A	A	80	0	80	0	60	31	0	0
chr10	28750512	28750513	G	G	220	0	220	0	60	72	0	0
chr10	28750548	28750549	G	C	255	255	255	255	60	97	0	0
chr10	28750574	28750575	T	T	237	0	237	0	60	83	0	0
chr10	28750577	28750578	T	T	234	0	234	0	60	82	0	0
chr10	28750578	28750579	T	T	242	0	242	0	60	76	0	0
chr10	28750593	28750594	G	G	220	0	220	0	60	75	0	0
chr10	28750640	28750641	T	C	165	165	165	165	60	46	0	0
chr10	28750746	28750747	G	A	202	202	202	202	60	58	0	0
chr10	28750766	28750767	A	G	205	205	205	205	60	59	0	0
chr10	28750769	28750770	T	C	175	175	175	175	60	49	0	0

Filter and Sort

- Filter data on any complex or simple expressions
- Sort data in ascending or descending order
- Select lines that match a regular expression
- Operate on Genomes
- Intersect the input queries
- Subtract the input queries
- Merge the output of a query

NGS: SAM To

- Filter SAM values
- Convert SAM to BAM
- SAM-to-BAM format to BAM
- BAM-to-SAM format to SAM
- Merge BAM files together
- Generate pileup dataset

- Filter pileup on coverage and SNPs

- Pileup-to-Interval condenses pileup format into ranges of bases

Analysis

Options

Analysis for Sample E18

Intersect to get Variants from Sample E18, consensus different, 10 Genes

mm9 RefSeq Genes

Intersect to get Variants from Sample E18 where consensus base is different than ref. base

Generate pileup to get Pileup from sample E18

Generate pileup on data 8

Convert BAM to SAM on data 7

Align with Bowtie for Bowtie on data 6 and data 5

6: E18 PE.2 Reads Groomed, Trimmed

5: E18 PE.1 Reads Groomed, Trimmed

aligner designed to be ultrafast and memory-efficient. It is developed by Mark Imamura, Michael Langmead, and Brent E. Tringali. Please cite: Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biol. 2009;10(3):R25.

User Metadata

History

Options ▾

Variant Analysis for Sample E18

Tags:

snp × pileup × bowtie ×

demo × sample:e18 ×

Annotation / Notes:

Perform a variant analysis with default parameters to identify variants in sample E18 that lie in annotated genes.

10: Variants from sample E18

26,742 regions, format: interval, database: mm9

Info:

Tags:

pileup × sample:e18 ×

snps ×

Annotation:

Find variants with coverage ≥ 30 and quality score ≥ 20 .

| display at UCSC [main](#) | view in [GeneTrack](#) | display at Ensembl [Current](#)

1. Chrom	2. Start	3. End	4	5	6
chr10	6882036	6882037	A	A	107
chr10	14243075	14243076	G	G	96
chr10	14243079	14243080	C	C	106
chr10	14465082	14465083	T	K	173
chr10	14465083	14465084	G	K	144
chr10	14465084	14465085	T	T	117

Datasources

Upload file from your computer

- ✦ FTP support for large datasets

UCSC table browser

BioMart

interMine / modMine

EuPathDB server

EncodeDB at NHGRI

EpiGRAPH server

Tool Suites

Text Manipulation

Format Converters

Filtering and Sorting

Join, Subtract, Group

Sequence Tools

Multi-species Alignment Tools

Genomic Interval Operations

Summary Statistics

Graphing / Plotting

Regional Variation

EMBOSS

Evolution / Phylogeny

RNA-seq

ChIP-seq

GATK

Picard

RGenetics

...and more

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

ROCHE-454 DATA

- [Build base quality distribution](#)
- [Select high quality segments](#)
- [Combine FASTA and QUAL](#) into FASTQ

AB-SOLID DATA

- [Convert SOLiD output to fastq](#)
- [Compute quality statistics](#) for SOLiD data
- [Draw quality score boxplot](#) for SOLiD data

GENERIC FASTQ MANIPULATION

- [Filter FASTQ](#) reads by quality score and length
- [FASTQ Trimmer](#) by column
- [FASTQ Quality Trimmer](#) by sliding window

Evolution

Metagenomic analyses

Human Genome Variation

EMBOSS

NGS TOOLBOX BETA

NGS: QC and manipulation

NGS: Mapping

ILLUMINA

- [Map with Bowtie](#) for Illumina
- [Map with BWA](#) for Illumina

ROCHE-454

- [Lastz](#) map short reads against reference sequence
- [Megablast](#) compare short reads against htgs, nt, and wgs databases

- [Parse blast XML output](#)

AB-SOLID

- [Map with Bowtie](#) for SOLiD

NGS: SAM Tools

NGS: Indel Analysis

NGS: Peak Calling

NGS: RNA Analysis

RGNETICS

SNP/WGA: Data; Filters

SNP/WGA: QC; LD; Plots

SNP/WGA: Statistical Models

Workflows

NGS TOOLBOX BETA

NGS: QC and manipulation

NGS: Mapping

NGS: SAM Tools

- [Filter SAM](#) on bitwise flag values
- [Convert SAM](#) to interval
- [SAM-to-BAM](#) converts SAM format to BAM format
- [BAM-to-SAM](#) converts BAM format to SAM format
- [Merge BAM Files](#) merges BAM files together
- [Generate pileup](#) from BAM dataset
- [Filter pileup](#) on coverage and SNPs
- [Pileup-to-Interval](#) condenses pileup format into ranges of bases
- [flagstat](#) provides simple stats on BAM files

NGS: Indel Analysis

NGS: Peak Calling

NGS: RNA Analysis

RGNETICS

SNP/WGA: Data; Filters

SNP/WGA: QC; LD; Plots

SNP/WGA: Statistical Models

Workflows

NGS: SAM Tools

NGS: Indel Analysis

- [Filter Indels](#) for SAM
- [Extract indels](#) from SAM
- [Indel Analysis](#)

NGS: Peak Calling

- [MACS](#) Model-based Analysis of ChIP-Seq
- [GeneTrack indexer](#) on a BED file
- [Peak predictor](#) on GeneTrack index

NGS: RNA Analysis

RNA-SEQ

- [Tophat](#) Find splice junctions using RNA-seq data
- [Cufflinks](#) transcript assembly and FPKM (RPKM) estimates for RNA-Seq data
- [Cuffcompare](#) compare assembled transcripts to a reference annotation and track Cufflinks transcripts across multiple experiments
- [Cuffdiff](#) find significant changes in transcript expression, splicing, and promoter use

FILTERING

- [Filter Combined Transcripts](#) using tracking file

Dozens of tools for different HTS applications packaged with Galaxy

VCF Tools

- Intersect Generate the intersection of two VCF files
- Annotate a VCF file (dbSNP, hapmap)
- Filter a VCF file
- Extract reads from a specified region

NGS: Picard (beta)

QC/METRICS FOR SAM/BAM

- BAM Index Statistics
- Sam/bam Alignment Summary Metrics
- Sam/bam GC Bias Metrics
- Estimate Library Complexity
- Insertion size metrics for PAIRED data
- Sam/bam Hybrid Selection Metrics For (eg exome) targeted data

BAM/SAM CLEANING

- Add or Replace Groups
- Reorder SAM
- Replace Sam Header
- Paired Read Mate Fixer for paired data
- Mark Duplicate reads

FASTQC: FASTQ/SAM/BAM

- Fastqc: Fastqc QC using FastQC from Babraham

NGS: GATK Tools

Alpha

REALIGNMENT

- Realigner Target Creator for use in local realignment
- Indel Realigner – perform local realignment

BASE RECALIBRATION

- Count Covariates on BAM files
- Table Recalibration on BAM files
- Analyze Covariates – perform local realignment

GENOTYPING

- Unified Genotyper SNP and indel caller

Overview

What is Galaxy?

What you can do in Galaxy

- ✦ analysis interface, tools and datasources
- ✦ **data libraries**
- ✦ workflows
- ✦ visualization
- ✦ sharing
- ✦ Pages

Galaxy 101 Exercise

Data Library "Bushman"

Library Actions ▾

These are the data underlying the analyses reported in the paper "Complete Khoisan and Bantu genomes from southern Africa" by S. C. Schuster et al., published in the journal Nature, February 18, 2010. Each data set can be downloaded and/or imported into a Galaxy history. Data will be updated as the project progresses.

Name	Information	Uploaded By	Date	File Size
<input type="checkbox"/> All SNPs in personal genomes ▾	Summary table of SNPs in all individuals	greg@bx.psu.edu	2010-01-28	676.8 Mb
<input type="checkbox"/> Alu insertions in KB1 ▾		greg@bx.psu.edu	2010-02-10	14.9 Kb
<input type="checkbox"/> Alu insertions in NB1 ▾		greg@bx.psu.edu	2010-02-10	6.5 Kb
<input type="checkbox"/> KB1 microsatellites.txt ▾		greg@bx.psu.edu	2010-02-15	3.5 Mb
<input type="checkbox"/> NB1 microsatellites.txt ▾		greg@bx.psu.edu	2010-02-15	828.5 Kb
<input type="checkbox"/> amino acid differences with functional predictions ▾		greg@bx.psu.edu	2010-02-05	1.1 Mb
<input type="checkbox"/> gene copy number of CP2 and other genes in personal genomes ▾		greg@bx.psu.edu	2010-02-15	2.1 Mb
<input type="checkbox"/> indels in ABT ▾		greg@bx.psu.edu	2010-02-03	105.3 Kb
<input type="checkbox"/> indels in KB1 ▾		greg@bx.psu.edu	2010-02-03	14.2 Mb
<input type="checkbox"/> indels in MD6 ▾		greg@bx.psu.edu	2010-02-03	109.8 Kb
<input type="checkbox"/> indels in NB1 ▾		greg@bx.psu.edu	2010-02-03	515.5 Kb
<input type="checkbox"/> indels in TK1 ▾		greg@bx.psu.edu	2010-02-03	123.2 Kb
<input type="checkbox"/> novel SNPs in ABT ▾		greg@bx.psu.edu	2010-02-09	9.4 Mb
<input type="checkbox"/> novel SNPs in KB1 ▾		greg@bx.psu.edu	2010-02-09	16.9 Mb
<input type="checkbox"/> novel SNPs in MD6 ▾		greg@bx.psu.edu	2010-02-09	594.1 Kb
<input type="checkbox"/> novel SNPs in NB1 ▾		greg@bx.psu.edu	2010-02-09	4.1 Mb
<input type="checkbox"/> novel SNPs in TK1 ▾		greg@bx.psu.edu	2010-02-09	722.6 Kb
<input type="checkbox"/> sequenced exon-containing intervals ▾		greg@bx.psu.edu	2010-02-03	3.1 Mb

For selected items:

<http://usegalaxy.org/bushman>

Managing Libraries

Loading Data

- ✦ Upload a single file
- ✦ Import datasets from a Galaxy history
- ✦ Upload a directory of files
- ✦ Directly from Sequencer using Sample Tracking System

Accessing Data

- ✦ Data contents on disk are not copied
- ✦ Dataset security: public, Role-based access control (RBAC)

Annotating Library Data: Library Templates

- ✦ Build user fillable forms
- ✦ Associate at Library, Folder or Dataset level

Overview

What is Galaxy?

What you can do in Galaxy

- ✦ analysis interface, tools and datasources
- ✦ data libraries
- ✦ **workflows**
- ✦ visualization
- ✦ sharing
- ✦ Pages

Galaxy 101 Exercise

Galaxy Workflows

The screenshot displays the Galaxy web interface. At the top, the browser address bar shows `http://main.g2.bx.psu.edu/`. The main navigation bar includes 'Galaxy', 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Help', and 'User'. On the left, a 'Tools' sidebar lists various categories such as 'Get Data', 'Send Data', 'ENCODE Tools', 'Text Manipulation', 'FASTA manipulation', 'Filter and Sort', 'Join, Subtract and Group', 'Extract Features', 'Fetch Sequences', 'Fetch Alignments', 'Get Genomic Scores', 'Operate on Genomic Intervals', 'Statistics', 'Graph/Display Data', 'Regional Variation', 'Multiple regression', 'Multivariate Analysis', 'Evolution', 'Metagenomic analyses', 'EMBOSS', 'NGS TOOLBOX BETA', 'NGS: QC and manipulation', 'NGS: Mapping', 'NGS: SAM Tools', 'NGS: Indel Analysis', 'NGS: Peak Calling', 'RGENETICS', 'SNP/WGA: Data; Filters', 'SNP/WGA: QC; LD; Plots', and 'SNP/WGA: Statistical Models'. The central panel shows a data table with a warning message: 'This dataset is large and only the first megabyte is shown below.' The table contains columns for chromosome (chr10), coordinates, and various metrics. On the right, a 'History' panel is open, showing a list of datasets and workflows. A context menu is visible over the history items, with options like 'Create New', 'Clone', 'Share or Publish', 'Extract Workflow', 'Dataset Security', 'Show Deleted Datasets', 'Show Hidden Datasets', 'Show structure', and 'Delete'. Below the history panel, there are workflow steps: '9: Generate pileup on data 8', '8: SAM-to-BAM on data Z', and '7: Map with Bowtie for Illumina on data 6 and data 5'. The bottom of the history panel shows a table with columns '1. QNAME' and '2. FLAG' and rows of sequencing data.

Galaxy Workflows

The screenshot displays the Galaxy workflow editor interface. On the left, a sidebar lists various tool categories such as 'Get Data', 'Text Manipulation', 'FASTA manipulation', and 'NGS TOOLBOX BETA'. The main workspace is divided into two columns: 'Tool' and 'History items created'. The 'Tool' column contains several tool entries, each with a checkbox to 'Include' it in the workflow. The 'History items created' column shows the resulting workflow steps, numbered 1 through 9. A right-hand panel shows a 'History Lists' menu with options like 'Extract Workflow' and 'Delete'. The bottom right corner shows a preview of a BAM file alignment.

Tool	History items created
Upload File <i>This tool cannot be used in workflows</i>	1: E18 PE.1 Reads <input checked="" type="checkbox"/> Treat as input dataset
Upload File <i>This tool cannot be used in workflows</i>	2: E18 PE.2 Reads <input checked="" type="checkbox"/> Treat as input dataset
FASTQ Groomer <input checked="" type="checkbox"/> Include "FASTQ Groomer" in workflow	3: E18 PE.1 Reads Groomed
FASTQ Groomer <input checked="" type="checkbox"/> Include "FASTQ Groomer" in workflow	4: E18 PE.2 Reads Groomed
FASTQ Trimmer <input checked="" type="checkbox"/> Include "FASTQ Trimmer" in workflow	5: E18 PE.1 Reads Groomed, Trimmed
FASTQ Trimmer <input checked="" type="checkbox"/> Include "FASTQ Trimmer" in workflow	6: E18 PE.2 Reads Groomed, Trimmed
Map with Bowtie for Illumina <input checked="" type="checkbox"/> Include "Map with Bowtie for Illumina" in workflow	7: Map with Bowtie for Illumina on data 6 and data 5
SAM-to-BAM <input checked="" type="checkbox"/> Include "SAM-to-BAM" in workflow	8: SAM-to-BAM on data 7
Generate pileup <input checked="" type="checkbox"/> Include "Generate pileup" in workflow	9: Generate pileup on data 8

Galaxy Workflows

The screenshot displays the Galaxy workflow editor interface. The browser address bar shows the URL: <http://main.g2.bx.psu.edu/workflow/editor?id=a6d94f12f42c1af8>. The page title is "Galaxy" and the navigation menu includes "Analyze Data", "Workflow", "Shared Data", "Visualization", "Help", and "User".

The main workspace is titled "Workflow Canvas | SNP variant detection from paired-end reads". It features a grid background with several workflow steps connected by lines:

- Input dataset** (output) connects to **FASTQ Groomer** (File to groom, output_file: fastqsanger, fastqcssanger, fastqsolexa, fastqillumina).
- FASTQ Groomer** connects to **FASTQ Trimmer** (FASTQ File, output_file).
- FASTQ Trimmer** connects to **Map with Bowtie for Illumina** (Forward FASTQ file, Reverse FASTQ file, output (sam)).
- Map with Bowtie for Illumina** connects to **SAM-to-BAM** (SAM File to Convert, output1 (bam)).
- SAM-to-BAM** connects to **Generate pileup** (Select the BAM file to generate the pileup file for, output1 (tabular)).
- Generate pileup** connects to **Filter pileup** (Select dataset, out_file1 (tabular, interval)).

A sidebar on the left lists various tools, including "Get Data", "Send Data", "ENCODING", "Lift-Over", "Text Manipulation", "Conversion", "FASTA", "Filter", "Join", "Subtract", "Extract", "Fetch", "Fetch A", "Get Gen", "Operate", "Statistic", "Graph/Region", "Multiple", "Multivariate", "Evolution", "Metagenomics", "EMBOSS", "NGS TO", "NGS: QC", "NGS: M", "NGS: SA", "NGS: In", "NGS: Pe", "RGENET", "SNP/WC", "SNP/WC", "SNP/WC", and "Workflo".

At the bottom, a checkbox is checked: "Include 'Generate pileup' in workflow". A small preview window in the bottom right corner shows a visualization of the workflow steps.

Galaxy Workflows

The screenshot displays the Galaxy workflow editor interface. The main window shows a workflow canvas titled "Workflow Canvas | SNP variant detection from paired-end reads". The workflow consists of several steps connected by lines:

- Input dataset** (output) connects to **FASTQ Groomer** (File to groom, output_file: fastqsanger, fastqcssanger, fastqsolexa, fastqillumina).
- FASTQ Groomer** connects to **FASTQ Trimmer** (FASTQ File, output_file).
- FASTQ Trimmer** connects to **Map with Bowtie for Illumina** (Forward FASTQ file, Reverse FASTQ file, output (sam)).
- Map with Bowtie for Illumina** connects to **SAM-to-BAM** (SAM File to Convert, output1 (bam)).
- SAM-to-BAM** connects to **Generate pileup** (Select generate pileup file for output).

The **SAM-to-BAM** tool is highlighted, showing its configuration options:

- Tool: SAM-to-BAM**
- Choose the source for the reference list:**
- SAM File to Convert:** Data input 'input1' (sam)
- Edit Step Actions:**
- Edit Step Attributes:**

At the bottom of the screen, there is a checkbox labeled "Include 'Generate pileup' in workflow" which is checked.

Galaxy Workflows

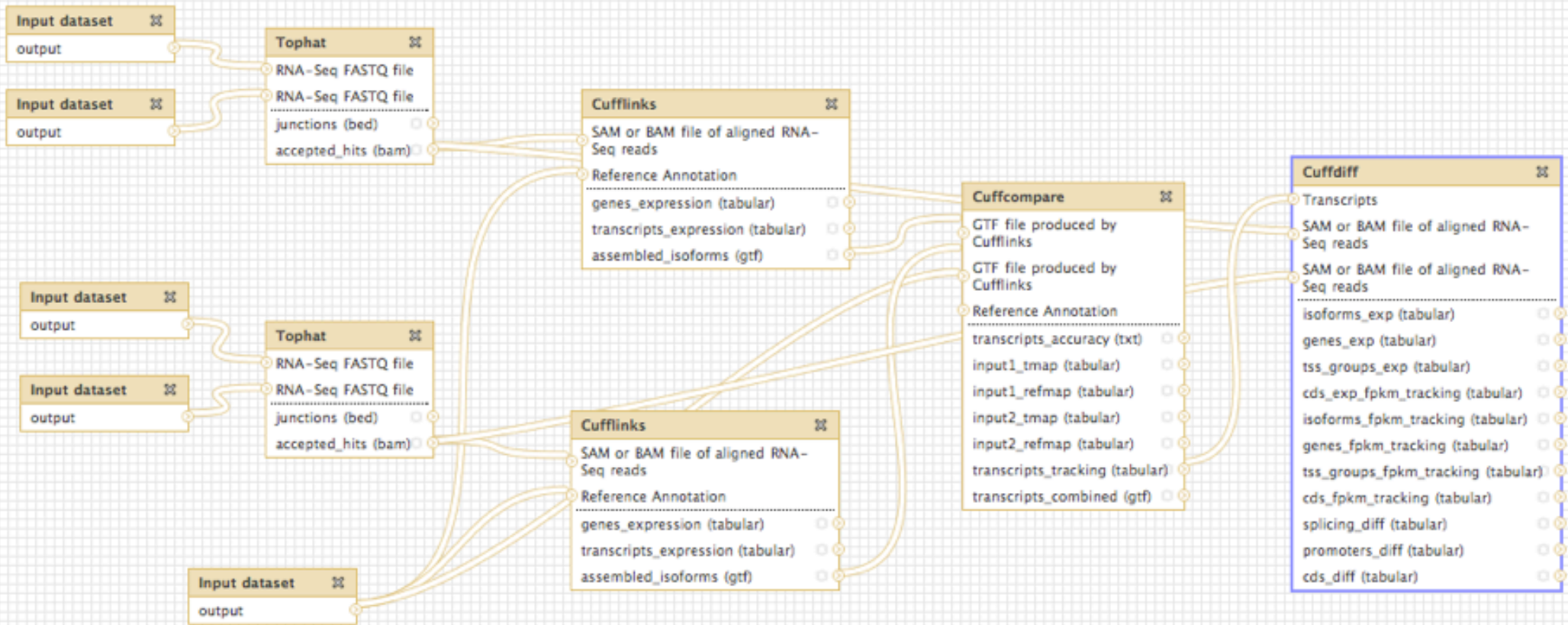
The image displays the Galaxy Workflows interface. In the background, a workflow canvas is visible with steps: 'Input dataset', 'FASTQ Groomer', 'FASTQ Trimmer', and 'FASTQ File'. The 'FASTQ Groomer' step is selected, opening an 'Edit Workflow Attributes' dialog. This dialog contains the following information:

- Name:** SNP identification within annotated genes from NGS PE Data
- Tags:** snp, ngs, pileup, bowtie
- Annotation / Notes:** Identify variants in annotated genes from NGS paired-end data.

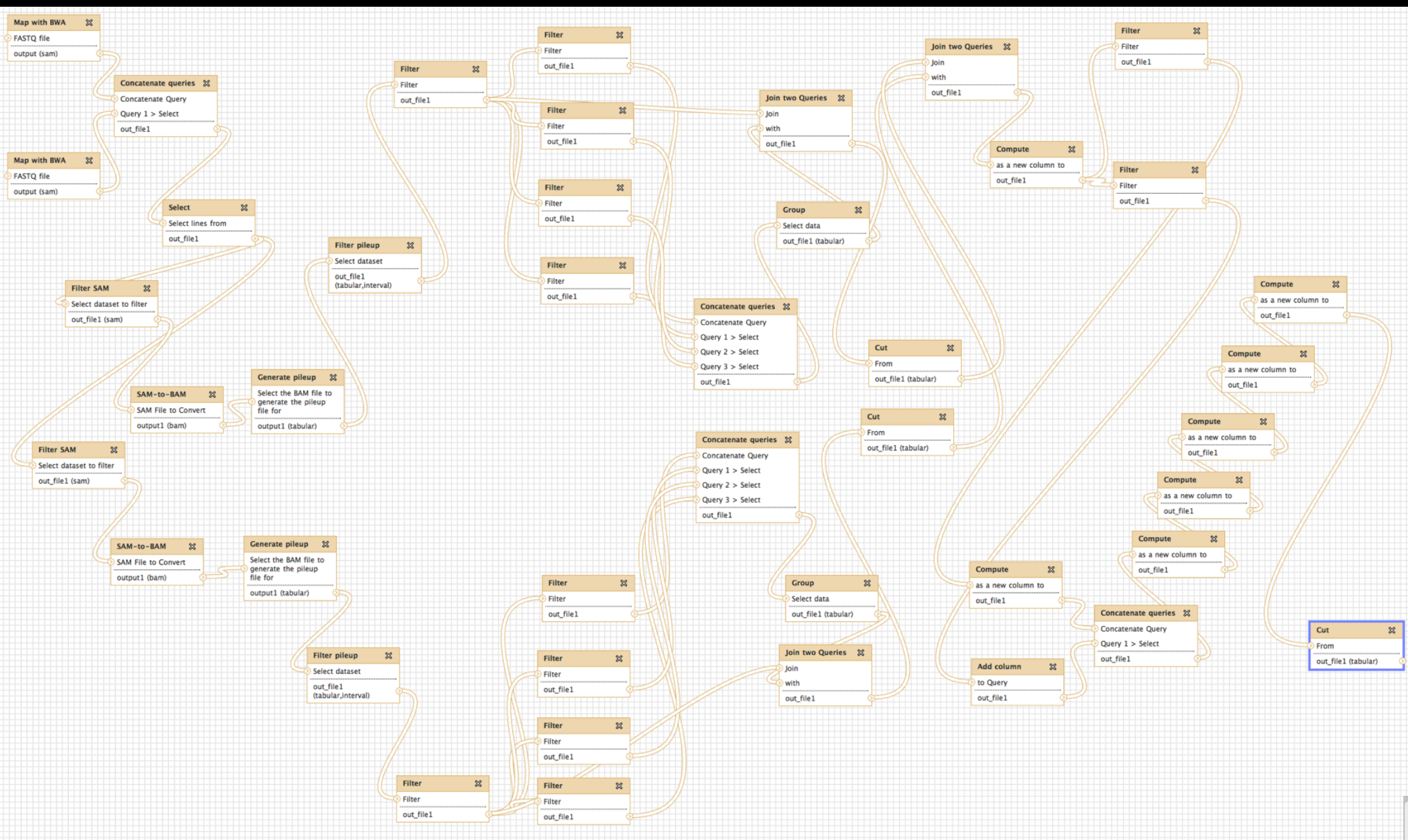
Below the dialog, a 'Tool: SAM-to-BAM' configuration panel is shown with the following settings:

- Choose the source for the reference list:** Locally cached
- SAM File to Convert:** Data input 'input1' (sam)
- Edit Step Actions:** Assign Columns, output1, Create
- Edit Step Attributes:** Annotation / Notes: Convert Bowtie SAM output to BAM format so that pileup can be run.

At the bottom of the interface, a checkbox is checked: 'Include "Generate pileup" in workflow'.



Example: Workflow for differential expression analysis of RNA-seq using Tophat/Cufflinks tools



Example: Diagnosing low-frequency heterosplasmic sites in two tissues from the same individual

Overview

What is Galaxy?

What you can do in Galaxy

- ✦ analysis interface, tools and datasources
- ✦ data libraries
- ✦ workflows
- ✦ **visualization**
- ✦ sharing
- ✦ Pages

Galaxy 101 Exercise

Visualize

Send data results to external genome browsers

Trackster: Galaxy's genome browser

External Genome Browsers

UCSC

Ensembl

GBrowse




IGV

UCSC Genome Browser on Mouse July 2007 (NCBI37)



move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out

position/search chr12:57,795,963-57,815,592 gene jump clear size 17,000 bp. compare

chr12 (qC1) 12qA1.1 qA2 12qA3 qB1 12qB3 12qC1 qC2 12qC3 qD1 qD2 12qD3 12qE 12qF1 qF2

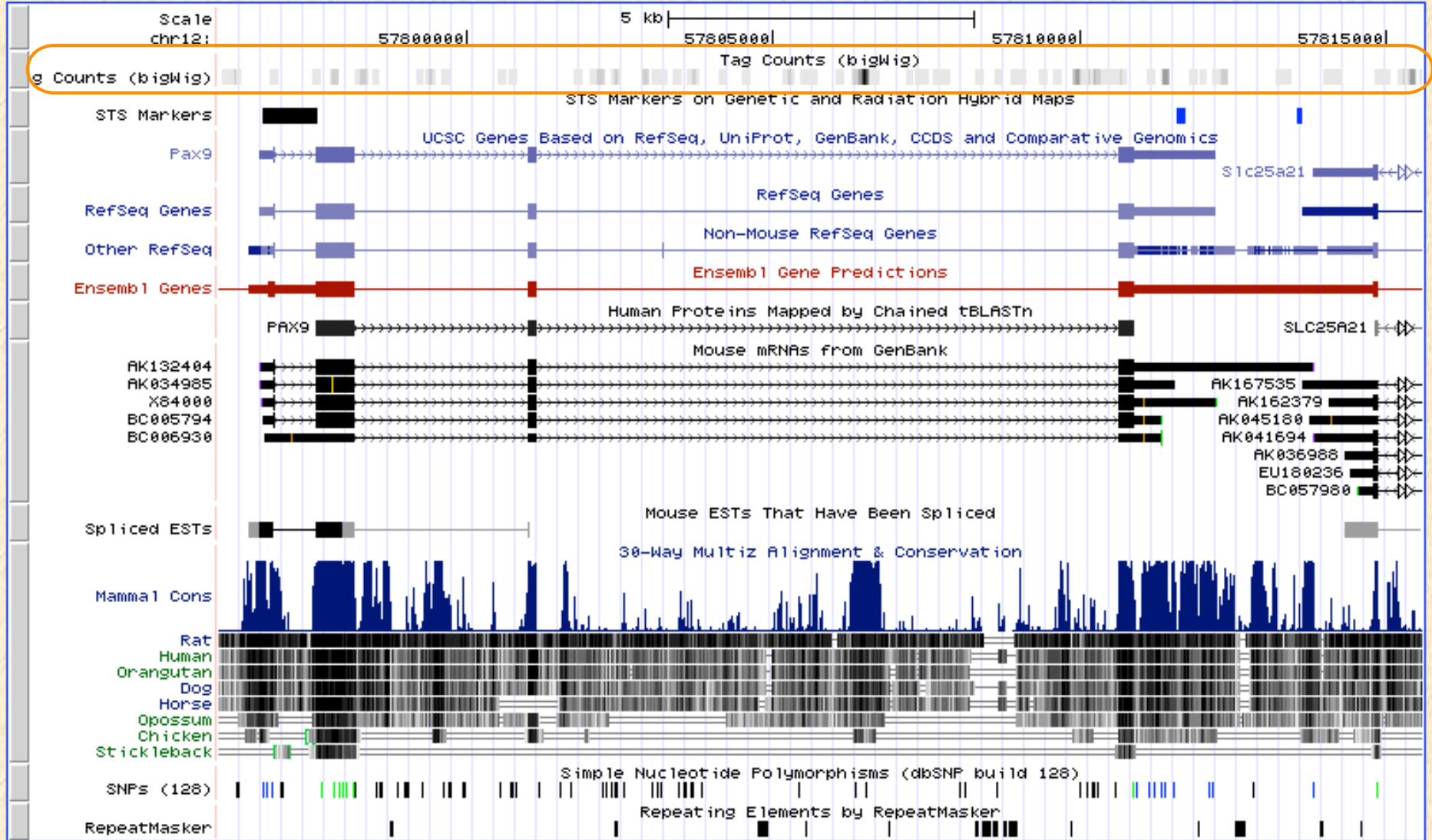
14: Tag Counts (bigWig)   

2.4 Gb, format: bigwig, database: mm9

Info:  

[display at UCSC main](#)

Binary UCSC BigWig file



Integrative Genomics Viewer (IGV)

1: Sample data

1.2 Gb
format: bam, database: mm9
Info: uploaded bam file



display at UCSC [main](#) [test](#)
display at Ensembl [Current](#)
display with IGV [web](#) [local](#)

Binary bam alignments file



The application "IGV 1.5" from "www.broadinstitute.org" is requesting access to your computer.

The digital signature could not be verified.

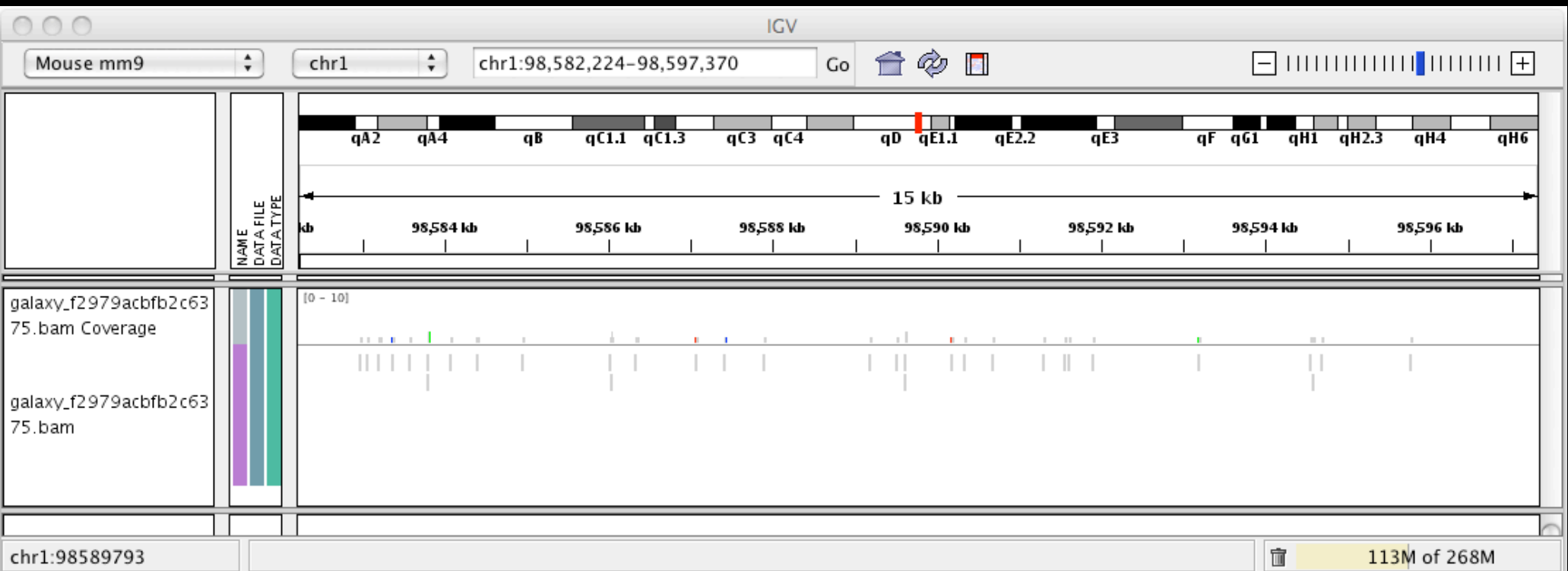
Allow all applications from "www.broadinstitute.org" with this signature



Show Details...

Deny

Allow



Galaxy

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

Genome Browser

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

Trackster

```
graph LR; Galaxy[Galaxy] --> Trackster[Trackster]; GenomeBrowser[Genome Browser] --> Trackster;
```

The diagram illustrates the relationship between Galaxy, Genome Browser, and Trackster. Galaxy and Genome Browser are positioned on the left, each in a light blue rounded rectangle. Two orange curved arrows originate from the right side of these boxes and point towards a larger light blue rounded rectangle on the right labeled Trackster. Galaxy's arrow points from the top right, and Genome Browser's arrow points from the bottom right.

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

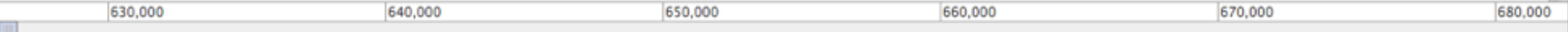
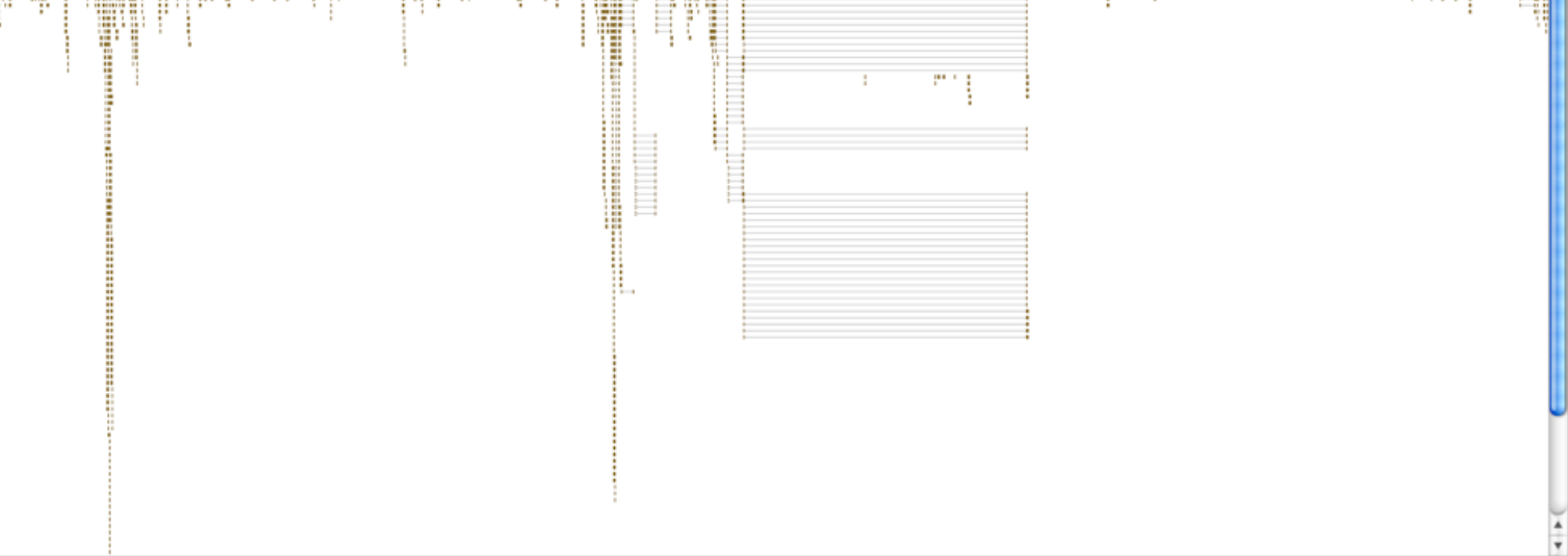
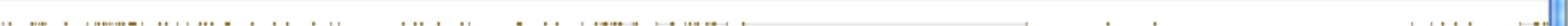
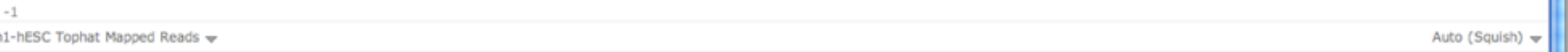
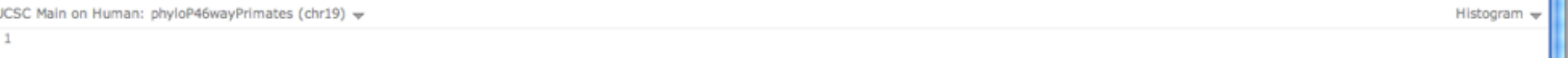
- ✦ custom genomes
- ✦ highly interactive

Published Visualizations | jeremy | GCC2011-1: Viewing and chr19 1,290 - 4,168,475

0 1,000,000 2,000,000 3,000,000 4,000,000



0 1,000,000 2,000,000 3,000,000 4,000,000



Published Visualizations | jeremy | GCC2011-1: Viewing and chr19 663,032 - 663,110

g g c c e g g g e c T C A C C G G C A G G C G C G G G A C G A T C T C C A C G G A G C A G C A G T G G C A G A G T A C C G T C C G G G A T G C G G C G A C

UCSC Main on Human: knownGene (chr19) Auto (Pack)

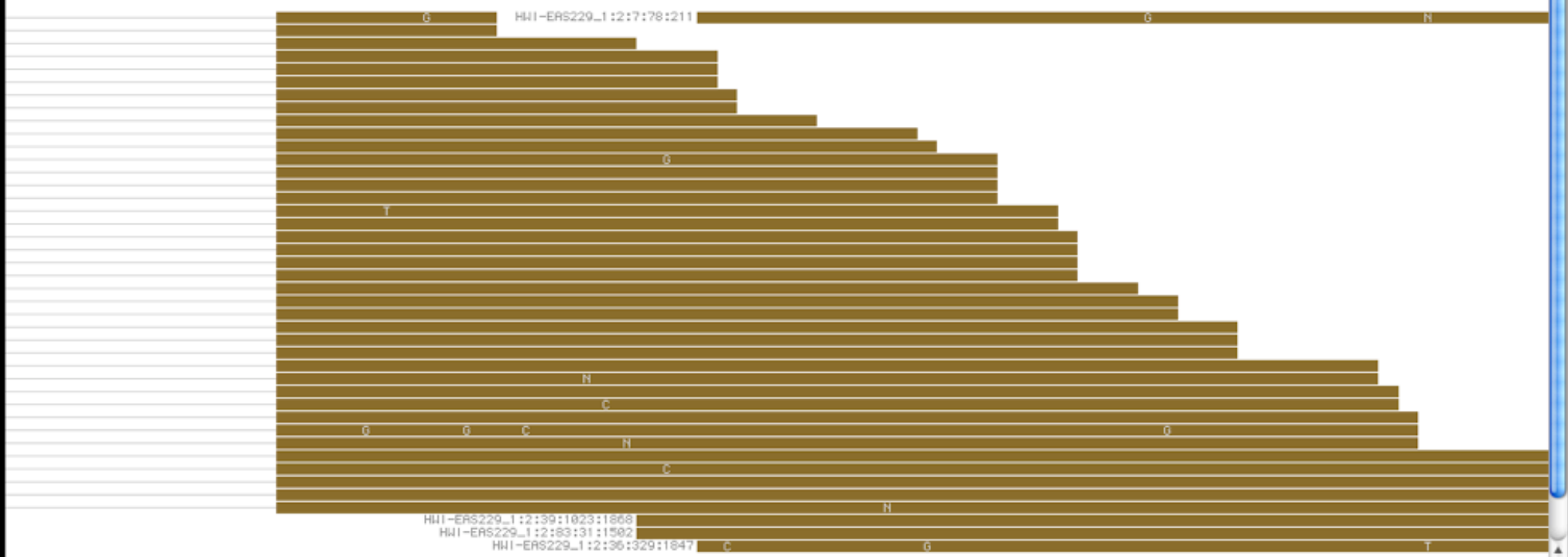
UCSC Main on Human: all_est (chr19) Dense



UCSC Main on Human: phyloP46wayPrimates (chr19) Histogram



h1-hESC Tophat Mapped Reads Auto (Pack)



h1-hESC Cufflinks assembled transcripts Auto (Pack)

g g c c e g g g e c T C A C C G G C A G G C G C G G G A C G A T C T C C A C G G A G C A G C A G T G G C A G A G T A C C G T C C G G G A T G C G G C G A C

Canceled opening the page

But really, why *another* genome browser

From static browsing to **visual analysis**

Visual feedback and experimentation needed for complex tools with many parameters

Leverage Galaxy strengths: a very sound model for abstracting interfaces to analysis tools and already integrates an enormous number

Dynamic Filtering



Integrating Tools and Visualization

Galaxy Analyze Data Workflow Shared Data **Visualization** Admin Help User

GCC3: Running Tools (hg19) chr19 1,523,098 - 1,545,232 1,530,000 1,540,000

UCSC Main on Human: knownGene

h1-hESC Tophat mapped reads

h1-hESC assembled transcripts - region=[all], parameters=[150000, 0.5, 0.05, No]

Cufflinks

Max Intron Length: 150000

Min Isoform Fraction: 0.5

Pre MRNA Fraction: 0.05

Perform quartile normalization: No

Run on complete dataset Run on visible region

CUFF.138.1 CUFF.139.1 CUFF.140.1 CUFF.141.1 CUFF.142.1

h1-hESC assembled transcripts - region=[all], parameters=[150000, 0.5, 0.05, No] ▼

Cufflinks

Max Intron Length

150000

Min Isoform Fraction

0.05

Pre MRNA Fraction

0.05

Perform quartile normalization

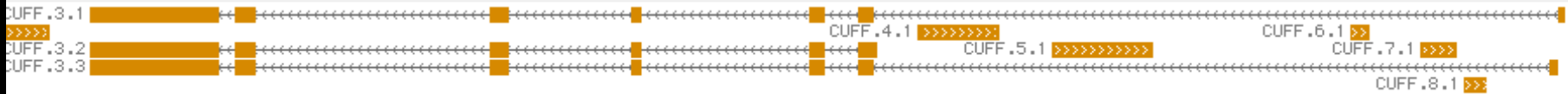
No

Run on complete dataset

Run on visible region



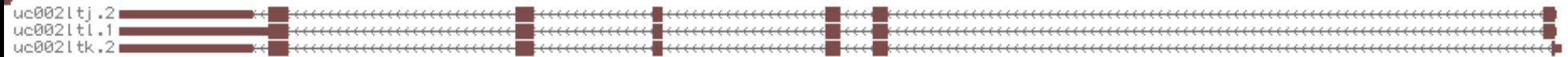
→ Cufflinks - region=[chr19:1523098-1545232], parameters=[150000, 0.05, 0.05, No] ▼



1,530,000

1,540

UCSC Main on Human: knownGene



h1-hESC Tophat mapped reads

h1-hESC assembled transcripts - region=[all], parameters=[150000, 0.5, 0.05, No]

Cufflinks

Max Intron Length
Min Isoform Fraction
Pre MRNA Fraction
Perform quartile normalization



Cufflinks - region=[chr19:1523098-1545232], parameters=[150000, 0.05, 0.05, No]



Cufflinks - region=[chr19:1523098-1545232], parameters=[150000, 0.05, 0.001, No]



Overview

What is Galaxy?

What you can do in Galaxy

- ✦ analysis interface, tools and datasources
- ✦ data libraries
- ✦ workflows
- ✦ visualization
- ✦ **sharing**
- ✦ Pages

Galaxy 101 Exercise

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 [Published Histories](#) 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 Published Histories 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 | Published History | Variant Analysis for Sample E18

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

Galaxy Analyze Data Workflow Shared Data Visualization Help User

Published Histories | jgoecks | Variant Analysis for Sample E18


Galaxy History ' Variant Analysis for Sample E18'

[+ Import history](#)

Annotation: Perform a pileup analysis with default parameters to identify variants in sample E18.

Dataset	Annotation
1: E18 PE.1 Reads	Forward reads from sample E18.
2: E18 PE.2 Reads	Reverse reads from sample E18.
3: E18 PE.1 Reads Groomed	Groom reads to convert quality scores from Solexa 1.0 to Solexa 1.3
4: E18 PE.2 Reads Groomed	Groom reads to convert quality scores from Solexa 1.0 to Solexa 1.3
5: E18 PE.1 Reads Groomed, Trimmed	Trim reads from 3' end to remove low-quality nts.
6: E18 PE.2 Reads Groomed, Trimmed	Trim reads from 3' to remove low-quality nts.
7: Map with Bowtie for Illumina on data 6 and data 5	Map paired-end reads with default parameters.
8: SAM-to-BAM on data 7	Need to convert Bowtie SAM to BAM so that pileup analysis can be performed.
9: Generate pileup on data 8	Pileup analysis with default parameters
10: Filter pileup to get Variants from sample E18	Find variants with coverage ≥ 30 .
13: Filter to get Variants from sample E18 where consensus base different than ref. base	Filter pileup to find variants where the consensus base is different than the reference base.
14: UCSC mm9 RefSeq Genes	UCSC mm9 RefSeq genes.
15: Intersect to get Variants from sample E18, consensus different, in RefSeq Genes	Variants with consensus different that occur in RefSeq genes.

About this History

Author  jgoecks

Related Histories
[All published histories](#)
[Published histories by jgoecks](#)

Rating
 Community (1 rating, 4.0 average) ★ ★ ★ ★ ☆
 Yours ★ ★ ★ ★ ☆

Tags
 Community: snp pileup bowtie demo
sample
 Yours: snp pileup bowtie
demo sample:e18 +

Galaxy | Published Workflow | SNP variant detection from paired-end reads

http://main.g2.bx.psu.edu/u/jgoecks/w/snp-variant-detection-from-paired-end-reads

Galaxy Analyze Data Workflow Shared Data Visualization Help User

Published Workflows | jgoecks | SNP variant detection from paired-end reads

Step 6: FASTQ Trimmer

Trim reads to remove low-quality bases.

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

Map reads using default parameter values.

Will you select a reference genome from your history or use a built-in index?
Use a built-in index

Select a reference genome
/galaxy/data/apiMe13/bowtie_index/apiMe13

Is this library mate-paired?
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 alignments (-X)
1000

The upstream/downstream mate orientation for valid paired-end alignment against 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

Step 8: SAM-to-BAM

Convert Bowtie SAM output to BAM format so that pileup can be run.

Choose the source for the reference list
Locally cached

About this Workflow

Author
jgoecks 

Related Workflows
[All published workflows](#)
[Published workflows by jgoecks](#)

Rating
Community (0 ratings, 0.0 average) ★★★★★
Yours ★★★★★

Tags
Community:
snp bowtie
Yours:
snp bowtie

Published Histories

 search | [Advanced Search](#)

Name	Annotation	Owner	Community Rating ↑	Community Tags	Last Updated
Galaxy vs MEGAN	Comparison of Galaxy vs. MEGAN pipeline.	aun1	★★★★★	metagenomics megan galaxy	Mar 19, 2010
metagenomic analysis		aun1	★★★★★	metagenomics galaxy	Mar 19, 2010
SM_1186088	Datasets correspond to our paper published in Science by Peleg et al. entitled : Altered histone acetylation is associated with age-dependent memory impairment. Experiment layout: This history contains 4 datasets in the form of BED files of uniquely mapped reads produced after chip-seq for histone modifications H4K12ac and H3K9ac in mouse hippocampus of 3 months (young) and 16 months (old) mice after fear conditioning. For detailed information please refer to supplementary materials and methods of the respective work by peleg et al.	fischerlab	★★★★★		Apr 19, 2010
Variant Analysis for Sample E18	Perform a pileup analysis with default parameters to identify variants in sample E18.	jgoecks	★★★★★	snp pileup bowtie demo sample	2 minutes ago
get longest exon		henri	★★★★★	chr22 longest marc exon human workshop	Sep 02, 2010
FASTA to Tabular Test		JJ	★★★★★		Aug 26, 2010
EKLF		yzc109	★★★★★		Aug 24, 2010

Open "http://main.g2.bx.psu.edu/history/list_published?sort=rating&f-tags=All" in a new tab

Sharing Trackster Visualizations

“A picture is worth a 1000 words.”

A fully-interactive visualization is worth many more words

Galaxy | Published Visualization: X

main.g2.bx.psu.edu/u/jeremy/v/gcc2011-1-viewing-and-navigating

Galaxy Analyze Data Workflow **Shared Data** Visualization Admin Help User

Published Visualizations | Jeremy | GCC2011-1: Viewing chr19 1,290 - 4,168,475

0 1,000,000 2,000,000 3,000,000 4,000,000

UCSC Main on Human: knownGene (chr19) Auto (Squish)

UCSC Main on Human: all_est (chr19) Auto (coverage histogram)

11431

UCSC Main on Human: phyloP46wayPrimates (chr19) Histogram

1

-1


h1-hESC Tophat Mapped Reads Auto (coverage histogram)

8732

h1-hESC Cufflinks assembled transcripts Auto (Squish)

0 1,000,000 2,000,000 3,000,000 4,000,000

About this Visualization

Author
jeremy 

Related Visualizations
[All published visualizations](#)
[Published visualizations by Jeremy](#)

Rating
Community (1 rating, 5.0 average) ★★★★★
Yours ★☆☆☆☆

Tags
Community: cool
Yours:

Overview

What is Galaxy?

What you can do in Galaxy

- ✦ analysis interface, tools and datasources
- ✦ data libraries
- ✦ workflows
- ✦ visualization
- ✦ sharing
- ✦ **Pages**

Galaxy 101 Exercise

Galaxy Pages

A web-based, interactive medium for presenting all aspects of an analysis: data, methods, and results

Galaxy Pages

The screenshot shows a web browser window with the URL <http://main.g2.bx.psu.edu/u/jgoecks/p/variant-analysis-for-sample-e18>. The page is titled "Variant Analysis of Embryonic Mouse Brain Tissue" and is authored by Jeremy Goecks, Anton Nekrutenko, James Taylor, and The Galaxy Team. The page content includes a "Results" section, a "Method" section, and a "References" section. The "Results" section describes a variant analysis experiment and lists potential variants. The "Method" section describes the bioinformatics pipeline used. The "References" section lists three scientific papers. The right sidebar contains "About this Page" information, including the author's profile picture, related pages, and a rating system.

Galaxy | Published Page | Variant Analysis for sample E18

Published Pages | jgoecks | Variant Analysis for sample E18

Variant Analysis of Embryonic Mouse Brain Tissue

Jeremy Goecks, Anton Nekrutenko, James Taylor, and The Galaxy Team

Results

To demonstrate how Galaxy can support accessible, reproducible, and transparent NGS re-sequencing studies, we performed a simple variant analysis experiment. This experiment identifies variants from a set of 4,536,964 RNA-seq reads obtained from sequencing a sample of mm9 brain tissue from day 18 of embryonic development.

The initial analysis produced support for 27,742 possible variants. Of these possible variants, there are 5,625 where (a) the consensus base--as determined by the MAQ model--differs from the reference base and (b) read coverage at the base is 30x or greater. Of these potential variants, 2796 occur in known RefSeq Genes. These potential variants are:

[Galaxy Dataset | Intersect to get Variants from sample E18, consensus different, in RefSeq Genes](#)
Variants with consensus different that occur in RefSeq genes.

Method

In the first step of this analysis, the reads were groomed to convert their quality scores from Solexa 1.0 to Solexa 1.3/Fastqsanger. Next, the reads were trimmed from 36bp to 27bp to exclude base pairs with low quality scores; see [1] for this step's rationale and parameter choices. After grooming and trimming, the reads were mapped using the short-read mapper Bowtie [2]. A pileup analysis using SAMtools [3] was then performed and was filtered to identify variants supported by 30+ reads. The complete analysis is contained in this history:

[Galaxy History | Variant Analysis for Sample E18](#)
Perform a pileup analysis with default parameters to identify variants in sample E18.

Here is a workflow for performing this analysis:

[Galaxy Workflow | Variant Identification within annotated genes from NGS PE Data](#)
Identify variants in annotated genes from NGS paired-end data.

References

[1] Han, X. et al. Transcriptome of embryonic and neonatal mouse cortex by high-throughput RNA sequencing. *Proceedings of the National Academy of Sciences* 106, 12741-12746 (2009).

[2] Langmead, B., Trapnell, C., Pop, M. & Salzberg, S.L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* 10, R25 (2009).

[3] Li, H. et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25, 2078-2079 (2009).

About this Page

Author
jgoecks

Related Pages
[All published pages](#)
[Published pages by jgoecks](#)

Rating

Community (0 ratings, 0.0 average) ★★★★★
Yours ★★★★★

Tags
Community: none
Yours:

Galaxy Pages

Galaxy | Published Page | Variant Analysis for sample E18

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

Galaxy Analyze Data Workflow Shared Data Visualization Help User

Published Pages | jgoecks | Variant Analysis for sample E18

The initial analysis produced support for 27,742 possible variants. Of these possible variants, there are 5,625 where (a) the consensus base--as determined by the MAQ model--differs from the reference base and (b) read coverage at the base is 30x or greater. Of these potential variants, 2796 occur in known RefSeq Genes. These potential variants are:

[Galaxy Dataset | Intersect to get Variants from sample E18, consensus different, in RefSeq Genes](#)
Variants with consensus different that occur in RefSeq genes.

Method

In the first step of this analysis, the reads were groomed to convert their quality scores from Solexa 1.0 to Solexa 1.3/Fastqsanger. Next, the reads were trimmed from 36bp to 27bp to exclude base pairs with low quality scores; see [1] for this step's rationale and parameter choices. After grooming and trimming, the reads were mapped using the short-read mapper Bowtie [2]. A pileup analysis using SAMtools [3] was then performed and was filtered to identify variants supported by 30+ reads. The complete analysis is contained in this history:

Galaxy History | Variant Analysis for Sample E18
Perform a pileup analysis with default parameters to identify variants in sample E18.

8: SAM-to-BAM on data 7	Need to convert Bowtie SAM to BAM so that pileup analysis can be performed.
9: Generate pileup on data 8	Pileup analysis with default parameters
10: Filter pileup to get Variants from sample E18	Find variants with coverage ≥ 30 .
13: Filter to get Variants from sample E18 where consensus base different than ref. base	Filter pileup to find variants where the consensus base is different than the reference base.
14: UCSC mm9 RefSeq Genes	UCSC mm9 RefSeq genes.
15: Intersect to get Variants from sample E18, consensus different, in RefSeq Genes	Variants with consensus different that occur in RefSeq genes.

Here is a workflow for performing this analysis:

[Galaxy Workflow | Variant identification within annotated genes from NGS PE Data](#)
Identify variants in annotated genes from NGS paired-end data.

References

About this Page

Author: jgoecks

Related Pages: [All published pages](#), [Published pages by jgoecks](#)

Rating: Community (0 ratings, 0.0 average), Yours

Tags: Community: none, Yours:

Galaxy Pages

Galaxy | Published Page | Variant Analysis for sample E18

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

Galaxy Analyze Data Workflow Shared Data Visualization Help User

Published Pages | jgoecks | Variant Analysis for sample E18

The initial analysis produced support for 27,742 possible variants. Of these possible variants, there are 5,625 where (a) the consensus base--as determined by the MAQ model--differs from the reference base and (b) read coverage at the base is 30x or greater. Of these potential variants, 2796 occur in known RefSeq Genes. These potential variants are:

Galaxy Dataset | Intersect to get Variants from sample E18, consensus different, in RefSeq Genes
Variants with consensus different that occur in RefSeq genes.

Method

In the first step of this analysis, the reads were groomed to convert their quality scores from Solexa 1.0 to Solexa 1.3/Fastqsanger. Next, the reads were trimmed from 36bp to 27bp to exclude base pairs with low quality scores; see [1] for this step's rationale and parameter choices. After grooming and trimming, the reads were mapped using the short-read mapper Bowtie [2]. A pileup analysis using SAMtools [3] was then performed and was filtered to identify variants supported by 30+ reads. The complete analysis is contained in this history:

Galaxy History | Variant Analysis for Sample E18
Perform a pileup analysis with default parameters to identify variants in sample E18.

8: SAM-to-BAM on data 7	Need to convert Bowtie SAM to BAM so that pileup analysis can be performed.
9: Generate pileup on data 8	Pileup analysis with default parameters
10: Filter pileup to get Variants from sample E18	Find variants with coverage >= 30.
13: Filter to get Variants from sample E18 where consensus base different than ref. base	Filter pileup to find variants where the consensus base is different than the reference base.
14: UCSC mm9 RefSeq Genes	UCSC mm9 RefSeq genes.
15: Intersect to get Variants from sample E18, consensus different, in RefSeq Genes	Variants with consensus different that occur in RefSeq genes.

Here is a workflow for performing this analysis:

Galaxy Workflow | Variant identification within annotated genes from NGS PE Data
Identify variants in annotated genes from NGS paired-end data.

References

Open "http://main.g2.bx.psu.edu/history/imp?id=e0b8bd5d661b10c2" in a new tab

About this Page

Author: jgoecks

Related Pages: All published pages, Published pages by jgoecks

Rating: Community (0 ratings, 0.0 average) ★★★★★, Yours ★★★★★

Tags: Community: none, Yours:

Galaxy Pages

The screenshot displays the Galaxy web interface. The browser address bar shows `http://main.g2.bx.psu.edu/`. The main navigation bar includes 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Help', and 'User'. The left sidebar lists various tool categories such as 'Get Data', 'Text Manipulation', 'Filter and Sort', and 'NGS: QC and manipulation'. The central workspace shows the 'Filter pileup' tool configuration. The 'Select dataset' dropdown is set to '9: Generate pileup on data 8'. The 'which contains' dropdown is set to 'Pileup with ten columns (with consensus)'. The 'Do not consider read bases with quality lower than' is set to 20, and 'Do not report positions with coverage lower than' is set to 30. The 'Only report variants?' and 'Convert coordinates to intervals?' options are both set to 'Yes'. The 'Print total number of differences?' and 'Print quality and base string?' options are also set to 'Yes'. An 'Execute' button is visible at the bottom of the tool configuration. The right sidebar shows a 'History' panel with a list of jobs, including '15: Variants from sample E18, consensus different, in RefSeq Genes', '14: UCSC mm9 RefSeq Genes', '13: Variants from sample E18 where consensus base different than ref. base', '10: Variants from sample E18', '9: Generate pileup on data 8', '8: SAM-to-BAM on data Z', '7: Map with Bowtie for Illumina on data 6 and data 5', and '6: E18 PE.2 Reads'. A table of genomic coordinates is visible in the history panel for job 10.

1. Chrom	2. Start	3. End	4	5	6
chr10	6882036	6882037	A	A	107
chr10	14243075	14243076	G	G	96
chr10	14243079	14243080	C	C	106
chr10	14465082	14465083	T	X	173
chr10	14465083	14465084	G	X	144
chr10	14465084	14465085	T	T	117

Open "http://main.g2.bx.psu.edu/tool_runner/rerun?id=1703758" in a new tab

Galaxy Pages

The screenshot shows a web browser window displaying a Galaxy page. The browser's address bar shows the URL: <http://main.g2.bx.psu.edu/u/jgoecks/p/variant-analysis-for-sample-e18>. The page title is "Variant Analysis of Embryonic Mouse Brain Tissue" by Jeremy Goecks, Anton Nekrutenko, James Taylor, and The Galaxy Team. The page is divided into several sections: "Results", "Method", and "References".

Results

To demonstrate how Galaxy can support accessible, reproducible, and transparent NGS re-sequencing studies, we performed a simple variant analysis experiment. This experiment identifies variants from a set of 4,536,964 RNA-seq reads obtained from sequencing a sample of mm9 brain tissue from day 18 of embryonic development.

The initial analysis produced support for 27,742 possible variants. Of these possible variants, there are 5,625 where (a) the consensus base--as determined by the MAQ model--differs from the reference base and (b) read coverage at the base is 30x or greater. Of these potential variants, 2796 occur in known RefSeq Genes. These potential variants are:

[Galaxy Dataset | Intersect to get Variants from sample E18, consensus different, in RefSeq Genes](#)
Variants with consensus different that occur in RefSeq genes.

Method

In the first step of this analysis, the reads were groomed to convert their quality scores from Solexa 1.0 to Solexa 1.3/Fastqsanger. Next, the reads were trimmed from 36bp to 27bp to exclude base pairs with low quality scores; see [1] for this step's rationale and parameter choices. After grooming and trimming, the reads were mapped using the short-read mapper Bowtie [2]. A pileup analysis using SAMtools [3] was then performed and was filtered to identify variants supported by 30+ reads. The complete analysis is contained in this history:

[Galaxy History | Variant Analysis for Sample E18](#)
Perform a pileup analysis with default parameters to identify variants in sample E18.

Here is a workflow for performing this analysis:

[Galaxy Workflow | Variant Identification within annotated genes from NGS PE Data](#)
Identify variants in annotated genes from NGS paired-end data.

References

[1] Han, X. et al. Transcriptome of embryonic and neonatal mouse cortex by high-throughput RNA sequencing. *Proceedings of the National Academy of Sciences* 106, 12741-12746 (2009).

[2] Langmead, B., Trapnell, C., Pop, M. & Salzberg, S.L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* 10, R25 (2009).

[3] Li, H. et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25, 2078 -2079 (2009).

About this Page

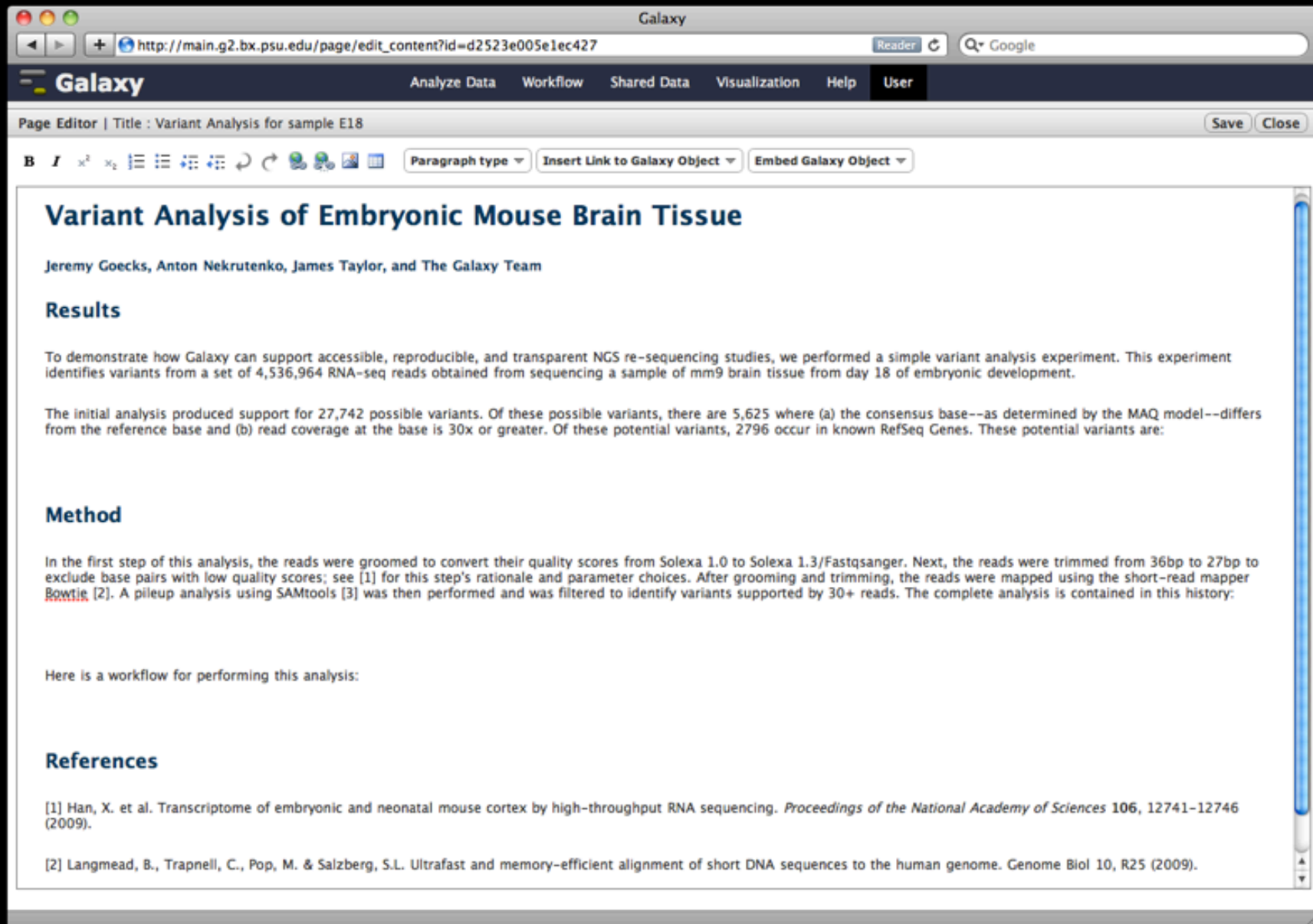
Author
jgoecks

Related Pages
[All published pages](#)
[Published pages by jgoecks](#)

Rating
Community (0 ratings, 0.0 average) ★★★★★
Yours ★★★★★




Tags
Community: none
Yours:

Creating a Page



The screenshot shows a web browser window with the URL http://main.g2.bx.psu.edu/page/edit_content?id=d2523e005e1ec427. The browser title is "Galaxy". The page header includes the Galaxy logo and navigation links: "Analyze Data", "Workflow", "Shared Data", "Visualization", "Help", and "User". The page editor interface shows the title "Variant Analysis of Embryonic Mouse Brain Tissue" and the authors "Jeremy Goecks, Anton Nekrutenko, James Taylor, and The Galaxy Team". The page content is divided into sections: "Results", "Method", and "References".

Page Editor | Title : Variant Analysis for sample E18 Save Close

B *I* \times^2 \times_2    Paragraph type

Variant Analysis of Embryonic Mouse Brain Tissue

Jeremy Goecks, Anton Nekrutenko, James Taylor, and The Galaxy Team

Results

To demonstrate how Galaxy can support accessible, reproducible, and transparent NGS re-sequencing studies, we performed a simple variant analysis experiment. This experiment identifies variants from a set of 4,536,964 RNA-seq reads obtained from sequencing a sample of mm9 brain tissue from day 18 of embryonic development.

The initial analysis produced support for 27,742 possible variants. Of these possible variants, there are 5,625 where (a) the consensus base--as determined by the MAQ model--differs from the reference base and (b) read coverage at the base is 30x or greater. Of these potential variants, 2796 occur in known RefSeq Genes. These potential variants are:

Method

In the first step of this analysis, the reads were groomed to convert their quality scores from Solexa 1.0 to Solexa 1.3/Fastqsanger. Next, the reads were trimmed from 36bp to 27bp to exclude base pairs with low quality scores; see [1] for this step's rationale and parameter choices. After grooming and trimming, the reads were mapped using the short-read mapper [Bowtie](#) [2]. A pileup analysis using [SAMtools](#) [3] was then performed and was filtered to identify variants supported by 30+ reads. The complete analysis is contained in this history:

Here is a workflow for performing this analysis:

References

[1] Han, X. et al. Transcriptome of embryonic and neonatal mouse cortex by high-throughput RNA sequencing. *Proceedings of the National Academy of Sciences* **106**, 12741-12746 (2009).

[2] Langmead, B., Trapnell, C., Pop, M. & Salzberg, S.L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* **10**, R25 (2009).

Creating a Page

The screenshot shows the Galaxy web interface. The browser address bar displays http://main.g2.bx.psu.edu/page/edit_content?id=d2523e005e1ec427. The page title is "Variant Analysis for sample E18". The main content area is titled "Variant Analysis of Embryonic Development" and includes sections for "Results" and "Method". A modal window titled "Embed Histories" is open, showing a search bar and a table of history items.

Name	Tags	Last Updated ↑
<input checked="" type="checkbox"/> Variant Analysis for Sample E18	5 Tags	15 minutes ago
<input type="checkbox"/> Pileup analysis, sample E18	4 Tags	2 days ago
<input type="checkbox"/> Unnamed history	0 Tags	Sep 07, 2010
<input type="checkbox"/> Unnamed history	0 Tags	Dec 17, 2009
<input type="checkbox"/> imported: Hsitory with ~100 items	5 Tags	Dec 10, 2009
<input type="checkbox"/> imported: Galaxy vs MEGAN	0 Tags	Dec 04, 2009
<input type="checkbox"/> imported: Galaxy vs MEGAN	2 Tags	Oct 06, 2009
<input type="checkbox"/> imported: Galaxy vs MEGAN	0 Tags	Oct 06, 2009
<input type="checkbox"/> imported: metagenomic analysis	0 Tags	Sep 30, 2009
<input type="checkbox"/> imported: Galaxy vs MEGAN	0 Tags	Sep 30, 2009

Page: 1 2 | [Show all histories on one page](#)

For 1 selected histories:

Make the selected histories accessible so that they can viewed by everyone.

[Embed](#) [Cancel](#)

Creating a Page

The screenshot shows a web browser window with the URL `http://main.g2.bx.psu.edu/page/edit_content?id=d2523e005e1ec427`. The browser's address bar includes a search engine (Google) and a 'Reader' icon. The page title is 'Variant Analysis for sample E18'. The interface features a dark blue navigation bar with the 'Galaxy' logo and menu items: 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Help', and 'User'. Below the navigation bar is a toolbar with icons for text formatting (bold, italic, subscript, superscript, list, link, unlink, undo, redo) and three dropdown menus: 'Paragraph type', 'Insert Link to Galaxy Object', and 'Embed Galaxy Object'. The main content area contains the following text:

Variant Analysis of Embryonic Mouse Brain Tissue

Jeremy Goecks, Anton Nekrutenko, James Taylor, and The Galaxy Team

Results

To demonstrate how Galaxy can support accessible, reproducible, and transparent NGS re-sequencing studies, we performed a simple variant analysis experiment. This experiment identifies variants from a set of 4,536,964 RNA-seq reads obtained from sequencing a sample of mm9 brain tissue from day 18 of embryonic development.

The initial analysis produced support for 27,742 possible variants. Of these possible variants, there are 5,625 where (a) the consensus base--as determined by the MAQ model--differs from the reference base and (b) read coverage at the base is 30x or greater. Of these potential variants, 2796 occur in known RefSeq Genes. These potential variants are:

Method

In the first step of this analysis, the reads were groomed to convert their quality scores from Solexa 1.0 to Solexa 1.3/Fastqsanger. Next, the reads were trimmed from 36bp to 27bp to exclude base pairs with low quality scores; see [1] for this step's rationale and parameter choices. After grooming and trimming, the reads were mapped using the short-read mapper [Bowtie](#) [2]. A pileup analysis using [SAMtools](#) [3] was then performed and was filtered to identify variants supported by 30+ reads. The complete analysis is contained in this history:

Embedded Galaxy History 'Variant Analysis for Sample E18'

[Do not edit this block; Galaxy will fill it in with the annotated history when it is displayed.]

Here is a workflow for performing this analysis:

References

[1] Han, X. et al. Transcriptome of embryonic and neonatal mouse cortex by high-throughput RNA sequencing. *Proceedings of the National Academy of Sciences* **106**, 12741-12746 (2009).

Open # on this page in a new tab

Creating a Page

The screenshot shows the Galaxy web interface in a browser window. The address bar shows the URL: http://main.g2.bx.psu.edu/page/edit_content?id=d2523e005e1ec427. The Galaxy logo is in the top left, and navigation tabs for 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Help', and 'User' are in the top right. The page title is 'Page Editor | Title : Variant Analysis for sample E18'. Below the title is a toolbar with various editing icons and three dropdown menus: 'Paragraph type', 'Insert Link to Galaxy Object', and 'Embed Galaxy Object'. The main content area contains the following text:

To demonstrate how Galaxy can support accessible, reproducible, and transparent NGS re-sequencing studies, we performed a simple variant analysis experiment. This experiment identifies variants from a set of 4,536,964 RNA-seq reads obtained from sequencing a sample of mm9 brain tissue from day 18 of embryonic development.

The initial analysis produced support for 27,742 possible variants. Of these possible variants, there are 5,625 where (a) the consensus base--as determined by the MAQ model--differs from the reference base and (b) read coverage at the base is 30x or greater. Of these potential variants, 2796 occur in known RefSeq Genes. These potential variants are:

Embedded Galaxy Dataset 'Variants from sample E18, consensus different, in RefSeq Genes'
[Do not edit this block; Galaxy will fill it in with the annotated dataset when it is displayed.]

Method

In the first step of this analysis, the reads were groomed to convert their quality scores from Solexa 1.0 to Solexa 1.3/Fastqsanger. Next, the reads were trimmed from 36bp to 27bp to exclude base pairs with low quality scores; see [1] for this step's rationale and parameter choices. After grooming and trimming, the reads were mapped using the short-read mapper Bowtie [2]. A pileup analysis using SAMtools [3] was then performed and was filtered to identify variants supported by 30+ reads. The complete analysis is contained in this history:

Embedded Galaxy History 'Variant Pileup Analysis for Sample E18'
[Do not edit this block; Galaxy will fill it in with the annotated history when it is displayed.]

Here is a workflow for performing this analysis:

Embedded Galaxy Workflow 'SNP identification within annotated genes from NGS PE Data'
[Do not edit this block; Galaxy will fill it in with the annotated workflow when it is displayed.]

References

[1] Han, X. et al. Transcriptome of embryonic and neonatal mouse cortex by high-throughput RNA sequencing. *Proceedings of the National Academy of Sciences* 106, 12741-12746 (2009).

[2] Langmead, B., Trapnell, C., Pop, M. & Salzberg, S.L. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* 10, R25 (2009).

[3] Li, H. et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25, 2078-2079 (2009).

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

Overview

What is Galaxy?

What you can do in Galaxy

- ✦ analysis interface, tools and datasources
- ✦ data libraries
- ✦ workflows
- ✦ visualization
- ✦ sharing
- ✦ Pages

Galaxy 101 Exercise



EMORY

PENNSTATE.



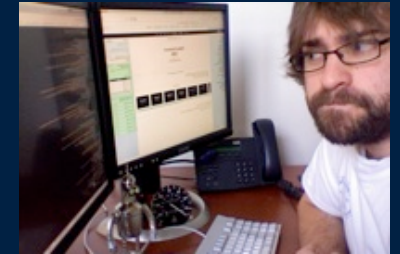
Enis Afgan



Dannon Baker



Dan Blankenberg



Nate Coraor



Dave Clements



Jeremy Goecks



Jennifer Jackson



Greg von Kuster



Kanwei Li



James Taylor



Guru Ananda



Anton Nekrutenko

Supported by the **NHGRI** (HG005542, HG004909, HG005133), **NSF** (DBI-0850103), Penn State University, Emory University, and the Pennsylvania Department of Public Health

Galaxy 101

<http://usegalaxy.org/galaxy101>

A simple question...

- ✦ Which coding exons have highest number of single nucleotide polymorphisms?

Galaxy 101

<http://usegalaxy.org/galaxy101>

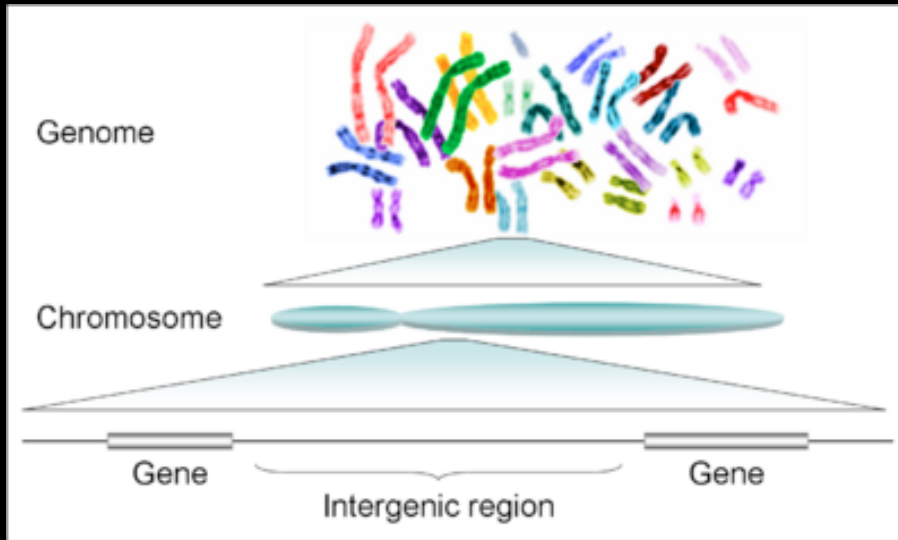
Overview

- ✦ Interactively Analyze Data
- ✦ Create reusable generic Workflow
- ✦ Share analysis Results, History, Workflow

Required Data

- ✦ Genomic Coordinates of coding exons and SNPs

Genomic Coordinates



```
>chr1  
taaccctaaccctaaccctaaccctaaccctaaccctaaccctaacccta  
accctaaccctaaccctaaccctaaccctaaccctaaccctaaccctaac
```

http://library.kiwix.org:4201/A/Human_genome.html

chrom	start	end	name	score	strand
chr1	0	10	first_ten_bases	0	+

see also:

<https://bitbucket.org/galaxy/galaxy-central/wiki/GopsDesc>

https://bitbucket.org/galaxy/galaxy-central/wiki/zero_based_coordinates.pdf

Galaxy 101: Basic Steps

<http://usegalaxy.org/galaxy101>

Get Genomic data from UCSC Table Browser

Determine each SNP that overlaps with a specific coding exon

Calculate count of overlapping SNPs for each exon

Sort and select exons by greatest SNP counts