

Using Galaxy for High-throughput Sequencing (HTS) Analysis and Visualization

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

Overview

High-throughput Sequencing (HTS) Data

Using Galaxy to Analyze HTS Data

- ✦ Prepare, quality control and manipulate reads
- ✦ Read Mapping
- ✦ SNP & INDEL analysis
- ✦ Binding sites analysis and peak calling
- ✦ Transcriptome analysis

Galaxy exercises: ChIP-seq and RNA-seq

HTS Data

From the Sequencer:

- ✦ reads and quality scores (FASTQ)

In the Analysis Pipeline / Workflow:

- ✦ alignments against reference genome (SAM, BAM)
- ✦ annotations (GFF, BED)
- ✦ genome Assemblies (FASTA)
- ✦ quantitative tracks, e.g. conservation (WIG)

Getting Your Data into Galaxy

Cannot upload any file larger than 2GB via Web browser

- ✦ Galaxy does not currently support compressed files

Use FTP client, e.g. FileZilla: <http://filezilla-project.org/>

Overview

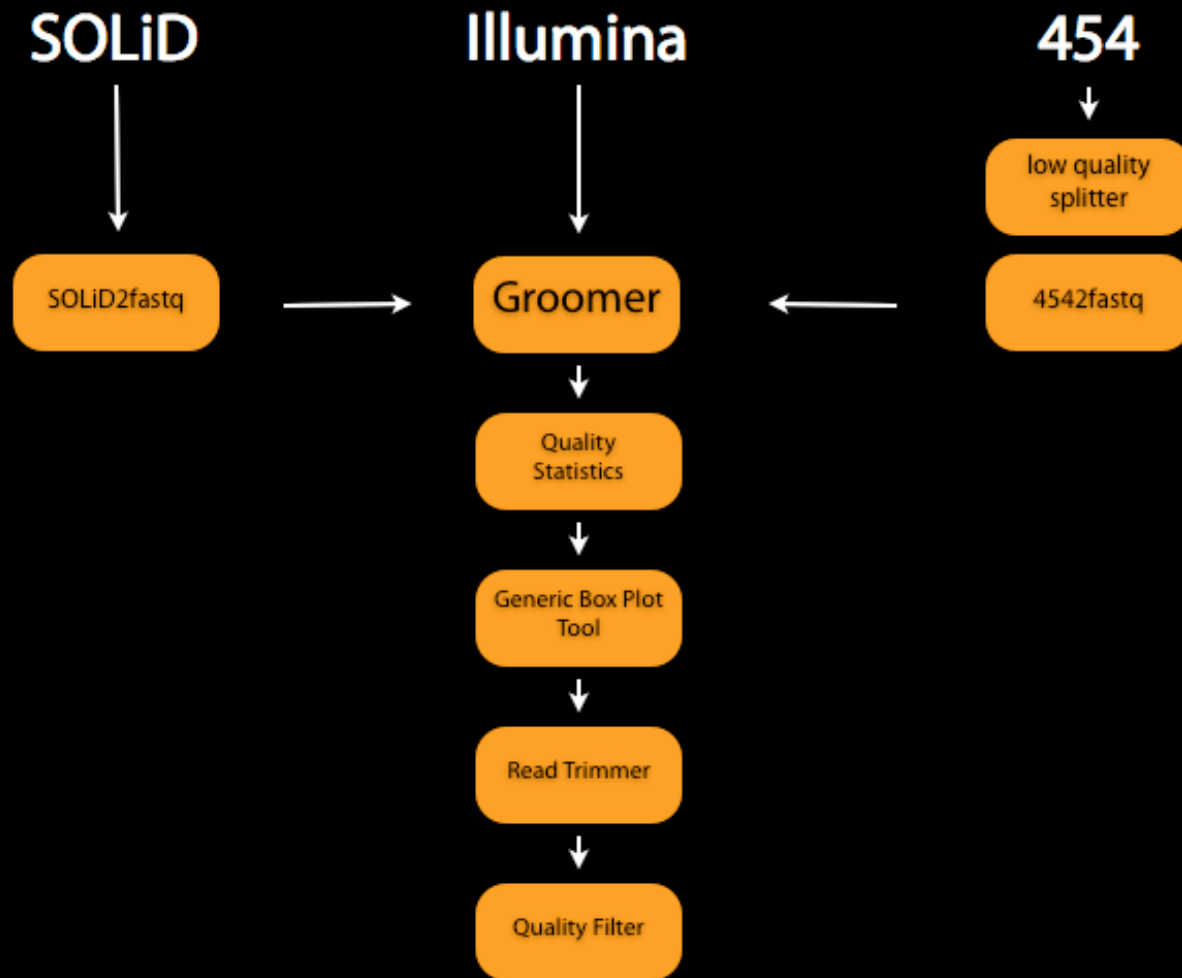
High-throughput Sequencing (HTS) Data

Using Galaxy to Analyze HTS Data

- ✦ Prepare, quality control and manipulate reads
- ✦ Read Mapping
- ✦ SNP & INDEL analysis
- ✦ Binding sites analysis and peak calling
- ✦ Transcriptome analysis

Galaxy exercises: ChIP-seq and RNA-seq

Prepare and Quality Check



Combining Sequences and Qualities

The screenshot displays the Galaxy web interface. At the top, the navigation bar includes 'Galaxy', 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Admin', 'Help', and 'User'. On the left, a 'Tools' sidebar lists various bioinformatics tools, with 'Combine FASTA and QUAL into FASTQ' highlighted under the 'AB-SOLID DATA' section. The main workspace shows the configuration for the 'Combine FASTA and QUAL' tool. The 'FASTA File' is set to '1: 454.fasta', the 'Quality Score File' is '2: 454.qual', and the 'Force Quality Score encoding' is 'ASCII'. An 'Execute' button is visible. Below the configuration, a 'What it does' section explains that the tool joins a FASTA file to a Quality Score file to create a single FASTQ block for each read. The 'History' panel on the right shows two jobs: '2: 454.qual' (52 lines, format: qual454) and '1: 454.fasta' (18 sequences, format: fasta). The bottom of the screen shows a preview of the FASTQ output, including sequence headers like '@EYKX4VC01B65GS' and sequence data.

Tools Options ▾

- 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 read score and length
- FASTQ Trimmer
- FASTQ Quality Trimming sliding window
- FASTQ Masker by

Combine FASTA and QUAL

FASTA File:
1: 454.fasta ▾

Quality Score File:
2: 454.qual ▾

Force Quality Score encoding:
ASCII ▾

Execute

What it does

This tool joins a FASTA file to a Quality Score file, creating a single FASTQ block for each read.

Specifying a set of quality scores is optional; when not provided, the output will be fastqsanger or fastqcssanger (when a csfasta is provided) with each quality score being the maximal allowed value (93).

Use this tool, for example, to convert 454-type output to FASTQ.

History Options ▾

Combine QUAL and Sequence

2: 454.qual

52 lines
format: qual454, database: ?
Info: uploaded qual454 file

1: 454.fasta

18 sequences
format: fasta, database: ?
Info: uploaded fasta file

```
@EYKX4VC01B65GS length=54 xy=0784_1754 region=1 run=R_2007_11_07_16_15_57_
CCGGTATCCGGGTGCCGTGATGAGCGCCACCGGAACGAATTCGACTATGCCGAA
+
B8C:===ABC%==@6=<<=====B8=B9E<@6==B;B9<=====A8=C:
@EYKX4VC01BNCSP length=187 xy=0558_3831 region=1 run=R_2007_11_07_16_15_57_
CTTACCGGTCACCACCGTGCCTTCAGGATTGATCGCCAGATCGGTCCGTGCGTCAGGCGGGGTGACATCGCCACACCGGTACTCACTGGCTCTGGTCTCCGGCGGCATCGGAG
+
<D; :F=F=:<E=<E=<E=<E<?4<E=B8E<<=<=<F?>;<99E<;=E=9:6=9=:C;LE7*84====;=HA-<E==;F==;====<E<=<E=<=<E=<E=HA-D>;F>====F>=E
@EYKX4VC01CD9FT length=115 xy=0865_1719 region=1 run=R_2007_11_07_16_15_57_
GGGGGCTTTGGCTGTCGTCGGCACCTCGCAAGAGCTACAGCAGGCGCGGCTGGCGATCATCGGCGGCACGCGGCCCTATATGTCGCGGGAACACACCACCCGCACCCAACGGC
+
D91*#<HB.E<E<====<=B8F==E<====E<====<====F====F;=E<====F====D;=<====E<D:A7====C:E<C:====E<=D>'====F?)B9=<<<
@EYKX4VC01B8FW0 length=95 xy=0799_0514 region=1 run=R_2007_11_07_16_15_57_
TAAATTTCAAGGAATGCAAATCAGGGTCTGTGTTTAGACTTCGGCTTTAGAGACCTGAATACGTCAAAAACATAACTTTCATGATATCTTGCAGT
+
=IC0D='<B8C9A7===JC2===F?*====<=F?)<=<D;<D;=F?*====C:==A7;====<LE8-"=6=<1=AB<=<====A7;=<;<=
@EYKX4VC01BCGYW length=115 xy=0434_3926 region=1 run=R_2007_11_07_16_15_57_
GGCCAGCCGGGACAGCGTTGTTGGGCTGCATGGCGACGAGCTAAAAGTCCGCATCACCGCCCGCGGTTGATGGGCAGGCTAATGCCCATCTGTTAAAAACTTTCTCGCCAAAC
+
=';0<=F=JD2=6=86<E<9E=IC/7:=9<=F;=<====<LE7)=;<;/=:5=C9:IB3"4<1E=E=6<:JC17=F>;<D<;<J1====<F>;LE8-",HA--25==2E>(9
@EYKX4VC01AZXC6 length=116 xy=0292_0280 region=1 run=R_2007_11_07_16_15_57_
GGGGGGCTTTGGCTGTCGTCGGCACCTCGCAAGAGCTACAGCAGGCGCGGCTGGCGATCATCGGCGGCACGCGGCCCTATATGTCGCGGGAACACACCACCCGCACCCAACGGC
+
```


Grooming --> Sanger

Tools 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 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

FASTQ Groomer

File to groom: 3: Combine FASTA and.. and data 2

Input FASTQ quality scores type: Sanger

What it does

This tool offers several conversions options relating to the FASTQ format. When using *Basic* options, the output will be *sanger* formatted (Sanger). When converting, if a quality score falls outside of the target range, the minimum or maximum quality score will be used. When converting between Solexa and the other formats, quality scores are scaled using the equations found in Cock PJ, Fields CJ, Goto N, Heuer ML, et al. (2009) quality scores, and the Solexa/Illumina FASTQ variants. Nucleotide bases are lost or gained; if gained, the base 'G' is used as the adapter base if there is no adapter present in the color space sequence. Any masked or ambiguous nucleotides in base space will be converted to 'N's when determining color space encoding.

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

```
@EYKX4VC01B65GS length=54 xy=0784_1
CCGGTATCCGGGTGCCGTGATGAGCGCCACCGGAA
+
BBC:---ABC<4---@6--<<-----BB=B9E<@6
@EYKX4VC01BNCSP length=187 xy=0558_
CTTACCGGTCACCACCGTGCCCTTCAGGATTGATCG
```

History

Combine QUAL and Sequence

3: Combine FASTA and QUAL on data 1 and data 2

18 sequences
format: fastqsanger, database: ?
Info: Combined 18 of 18 sequences with quality scores (100.00%).

```
@EYKX4VC01B65GS length=54 xy=0784_1
CCGGTATCCGGGTGCCGTGATGAGCGCCACCGGAA
+
BBC:---ABC<4---@6--<<-----BB=B9E<@6
@EYKX4VC01BNCSP length=187 xy=0558_
CTTACCGGTCACCACCGTGCCCTTCAGGATTGATCG
```

2: 454.qual

52 lines
format: qual454, database: ?
Info: uploaded qual454 file

```
>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 25 28 27 36
```

Quality Score Comparison

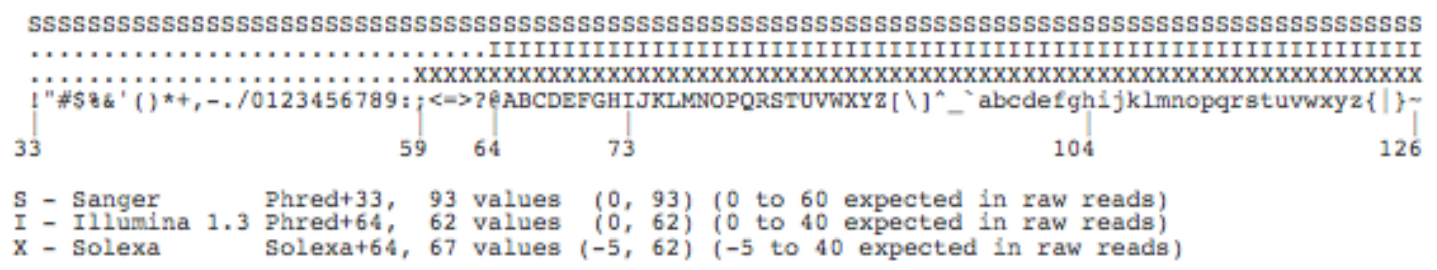


Diagram adapted from http://en.wikipedia.org/wiki/FASTQ_format

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

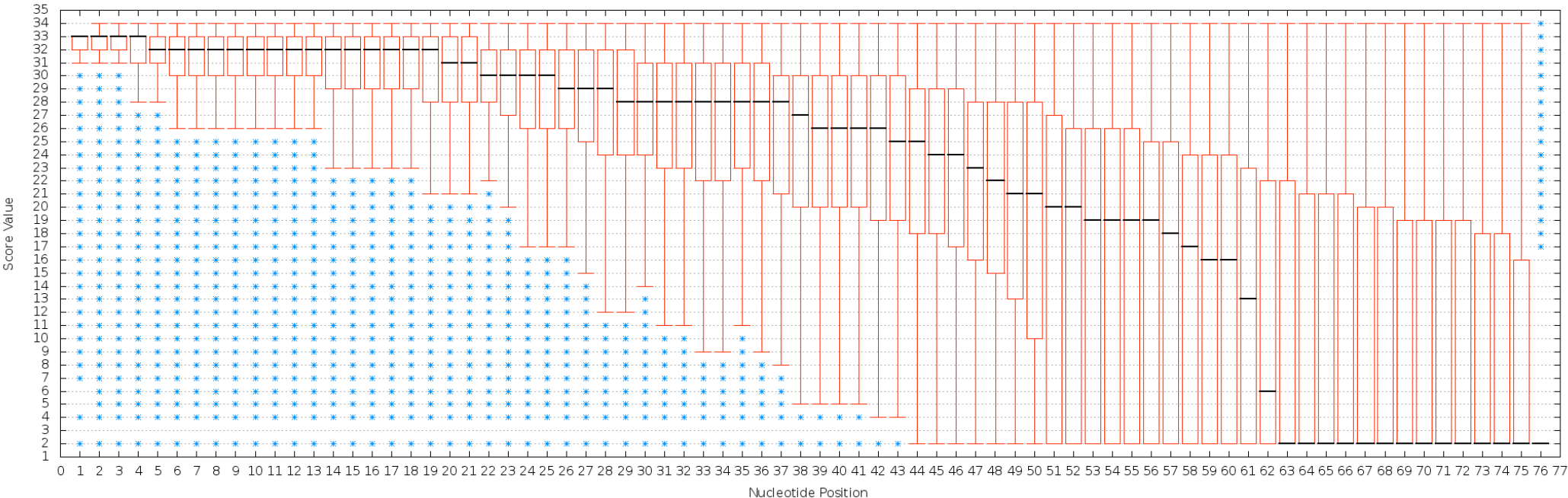
Quality Statistics and Box Plot Tool

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

Box plot in Galaxy

Quartiles 
Medians 
Outliers 



FastQC

The screenshot shows the Galaxy web interface with a FastQC report. The browser address bar shows 'main.g2.bx.psu.edu'. The Galaxy navigation bar includes 'Analyze Data', 'Workflow', 'Shared Data', 'Visualization', 'Admin', 'Help', and 'User'. The left sidebar shows 'Tools' with 'fastqc' selected under 'NGS: QC and manipulation'. The main content area displays the 'FastQC Report' for 'dataset_1750787.dat' dated 'Mon 20 Jun 2011'. The 'Summary' section lists 12 metrics with status icons: Basic Statistics (green check), Per base sequence quality (red X), Per sequence quality scores (green check), Per base sequence content (green check), Per base GC content (yellow exclamation), Per sequence GC content (yellow exclamation), Per base N content (green check), Sequence Length Distribution (green check), Sequence Duplication Levels (red X), Overrepresented sequences (yellow exclamation), and Kmer Content (red X). Below this is the 'Basic Statistics' section with a table:

Measure	Value
Filename	dataset_1750787.dat
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9

The right sidebar shows a 'History' panel with a list of jobs, including '120: FastQC.html' (11.6 Kb, format: html, database: hg18) and several 'Cuffdiff' jobs.

Read Trimming

Tools

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 Tabular converter
 - Tabular to FASTQ converter
- FASTX-TOOLKIT FOR FASTQ DATA
- Quality format converter (ASCII-Numeric)
 - Compute quality statistics
 - Draw quality score boxplot
 - Draw nucleotides distribution chart
 - FASTQ to FASTA converter
 - Filter by quality
 - Remove sequencing artifacts

FASTQ Trimmer

FASTQ File:

2: imported: GM12878..ple Dataset

Define Base Offsets as:

Absolute Values

Use Absolute for fixed length reads (Illumina, SOLiD)
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 offsets to be trimmed. When using the percent-based method, offsets are relative to the read length.

For example, if you have a read of length 36:

```
@Some FASTQ Sanger Read  
CAATATGTNCTCACTGATAAGTGGATATNAGCNCCA  
+  
-@@.0;B-?78>CBA@>7@7BBCA4-48%<;%<B@
```

And you set absolute offsets of 2 and 0:

FASTQ Quality Trimmer

FASTQ File:

7: FASTQ Trimmer on data 2

Keep reads with zero length:

Trim ends:

5' and 3'

Window size:

1

Step Size:

1

Maximum number of bases to exclude from the window during aggregation:

0

Aggregate action for window:

min score

Trim until aggregate score is:

>=

Quality Score:

0.0

Execute

Filter FASTQ

FASTQ File:

7: FASTQ Trimmer on data 2

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

Minimum Size:

0

Maximum Size:

0

A maximum size less than 1 indicates no limit.

Minimum Quality:

0.0

Maximum Quality:

0.0

A maximum quality less than 1 indicates no limit.

Maximum number of bases allowed outside of quality range:

0

This is paired end data:

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:

Absolute Values

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

Offset from 5' end:

0

Values start at 0, increasing from the left

Offset from 3' end:

0

Values start at 0, increasing from the right

Aggregate read score for specified range:

min score

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

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

Match Reads 1

Match Reads by:

Sequence Content

Sequence Match Type:

Regular Expression

Match by:

N

Remove Match Reads 1

Add new Match Reads

Manipulate Reads

Add new Manipulate Reads

Execute

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

Match Reads 1

Match Reads by:

Sequence Content

Sequence Match Type:

Regular Expression

Match by:

N

Remove Match Reads 1

Add new Match Reads

Manipulate Reads

Manipulate Reads 1

Manipulate Reads on:

Miscellaneous Actions

Miscellaneous Manipulation Type:

Remove Read

Remove Manipulate Reads 1

Add new Manipulate Reads

Execute

Overview

High-throughput Sequencing (HTS) Data

Using Galaxy to Analyze HTS Data

- ✦ Prepare, quality control and manipulate reads
- ✦ **Read Mapping**
- ✦ SNP & INDEL analysis
- ✦ Binding sites analysis and peak calling
- ✦ Transcriptome analysis

Galaxy exercises: ChIP-seq and RNA-seq

Mapping HTS Data

Collection of interchangeable mappers

- ✦ accept fastq format, produce SAM/BAM

Mappers for

- ✦ DNA
- ✦ RNA
- ✦ Local realignment

Mappers

DNA

- ✦ short reads: Bowtie, BWA, BFAST, PerM
- ✦ longer reads: LASTZ

Metagenomics

- ✦ Megablast

RNA / gapped-reads mapper

- ✦ Tophat

Commonly Used/Default Parameters

Lastz

Align sequencing reads in:

Against reference sequences that are:

Using reference genome:

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

Output format:

Lastz settings to use:

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 95% identity
Roche-454 90% identity
Roche-454 85% identity
Roche-454 75% identity
Illumina 95% identity
Illumina 85% identity

reference name?:

by this identity (%):

Do not report matches above this identity (%):

Do not report matches that cover less than this percentage of each read:

Convert lowercase bases to uppercase:

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

Output format:

SAM

Lastz settings to use:

Full Parameter List

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

Full Parameter List

Which strand to search?:

Both

Select seeding settings:

Seed hits require a 19 bp word with matches in

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

Do you want to modify the reference name?:

No

Full Parameter List

Do you want to modify the reference name?:

No

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 fastq and other method specific formats.

Overview

High-throughput Sequencing (HTS) Data

Using Galaxy to Analyze HTS Data

- ✦ Prepare, quality control and manipulate reads
- ✦ Read Mapping
- ✦ **SNP & INDEL analysis**
- ✦ Binding sites analysis and peak calling
- ✦ Transcriptome analysis

Galaxy exercises: ChIP-seq and RNA-seq

SNPs & INDELS

SNPs from Pileup

- ✦ Generate
- ✦ Filter

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

The screenshot shows the Galaxy web interface with the 'Indel Analysis' tool selected. The tool configuration includes a dropdown menu for the input file (54: BAM-to-SAM on dat..nverted SAM) and a frequency threshold input field set to 0.015. Below the configuration is a table of reference sequences and their corresponding indel statistics.

Indel Analysis

Select sam file to analyze:
54: BAM-to-SAM on dat..nverted SAM

Frequency threshold:
0.015

Cutoff

Execute

What it does

Given an input sam file, this tool provides analysis of the indels. It filters out matches that do not meet the frequency threshold. The way this frequency of occurrence is calculated is different for deletions and insertions. The CIGAR string's "M" can indicate an exact match or a mismatch. For SAM containing the following bits of information (assuming the reference "ACTGCTCGAT"):

CHROM	POS	CIGAR	SEQ
ref	3	2M1I3M	TACTTC
ref	1	2M1D3M	ACGCT
ref	4	4M2I3M	GTTCAAGAT
ref	2	2M2D3M	CTCCG
ref	1	3M1D4M	AACCTGG
ref	6	3M1I2M	TTCAAT
ref	5	3M1I3M	CTCTGTT
ref	7	4M	CTAT
ref	5	5M	CGCTA
ref	3	2M1D2M	TGCC

The following totals would be calculated (this is an intermediate step and not output):

POS	BASE	NUMREADS	DELPROPCALC	DELPROP	INSPROPSTARTCALC	INSSTARTPROP	INSPROPENDCALC	INSENDPROP
1	A	2	2/2	1.00	---	---	---	---
2	A	1	1/3	0.33	---	---	---	---
	C	2	2/3	0.67	---	---	---	---
3	C	1	1/5	0.20	---	---	---	---
	T	3	3/5	0.60	---	---	---	---
	-	1	1/5	0.20	---	---	---	---
4	A	1	1/6	0.17	---	---	---	---

GATK Tools

Local re-alignment

Base re-calibration

Genotyping

Alpha status

- ✦ please try, report bugs
- ✦ available on test server:
<http://test.g2.bx.psu.edu/>

NGS: GATK Tools

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

Unified Genotyper

Inputs

- ✦ BAM files

Lots of possible parameters

Output

- ✦ VCF file(s)

Unified Genotyper

Choose the source for the reference list:

Sample BAM files

Sample BAM file 1

BAM file:

Using reference genome:

dbSNP reference ordered data (ROD):

Binding for reference-ordered datas

The minimum phred-scaled confidence threshold at which variants not at 'trigger' track sites should be called:

The minimum phred-scaled confidence threshold at which variants not at 'trigger' track sites should be emitted (and filtered if less than the calling threshold):

Basic or Advanced GATK options:

Basic or Advanced Analysis options:

Overview

High-throughput Sequencing (HTS) Data

Using Galaxy to Analyze HTS Data

- ✦ Prepare, quality control and manipulate reads
- ✦ Read Mapping
- ✦ SNP & INDEL analysis
- ✦ **Binding sites analysis and peak calling**
- ✦ Transcriptome analysis

Galaxy exercises: ChIP-seq and RNA-seq

Peak Calling / ChIP-seq analysis

Punctate binding

- ✦ transcription factors

Diffuse binding

- ✦ histone modifications
- ✦ PolII

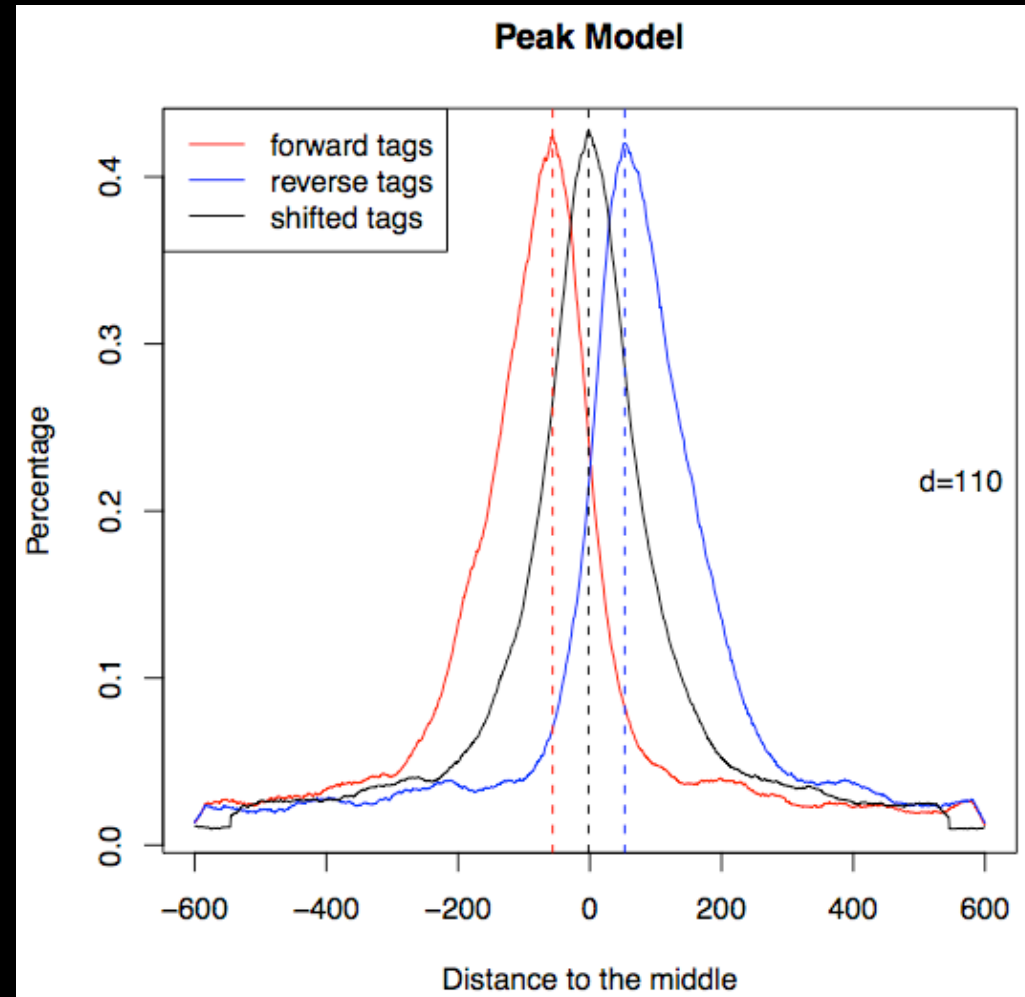
Punctate Binding --> MACS

Inputs

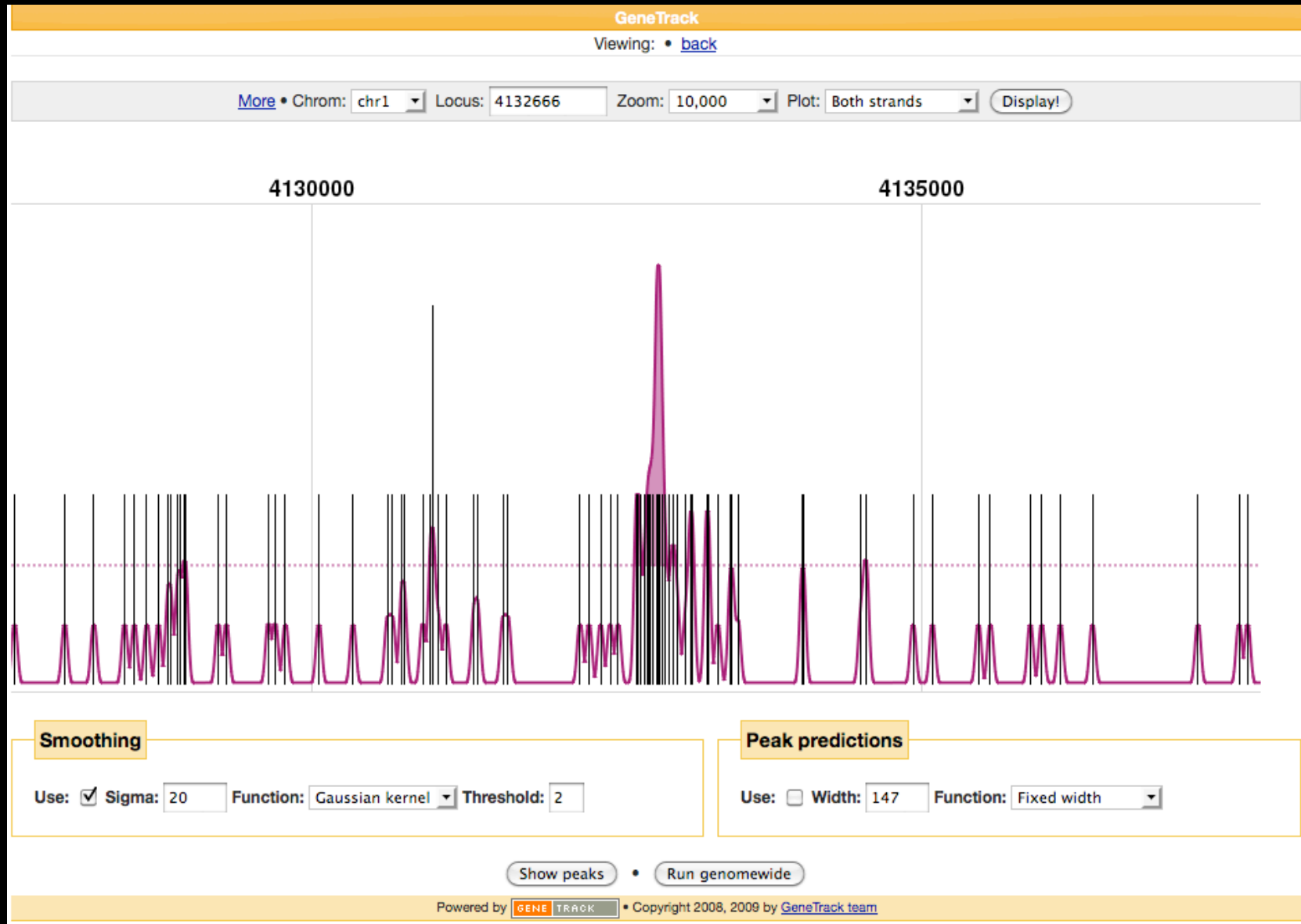
- ✦ Enriched Tag file
- ✦ Control / Input file (optional)

Outputs

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

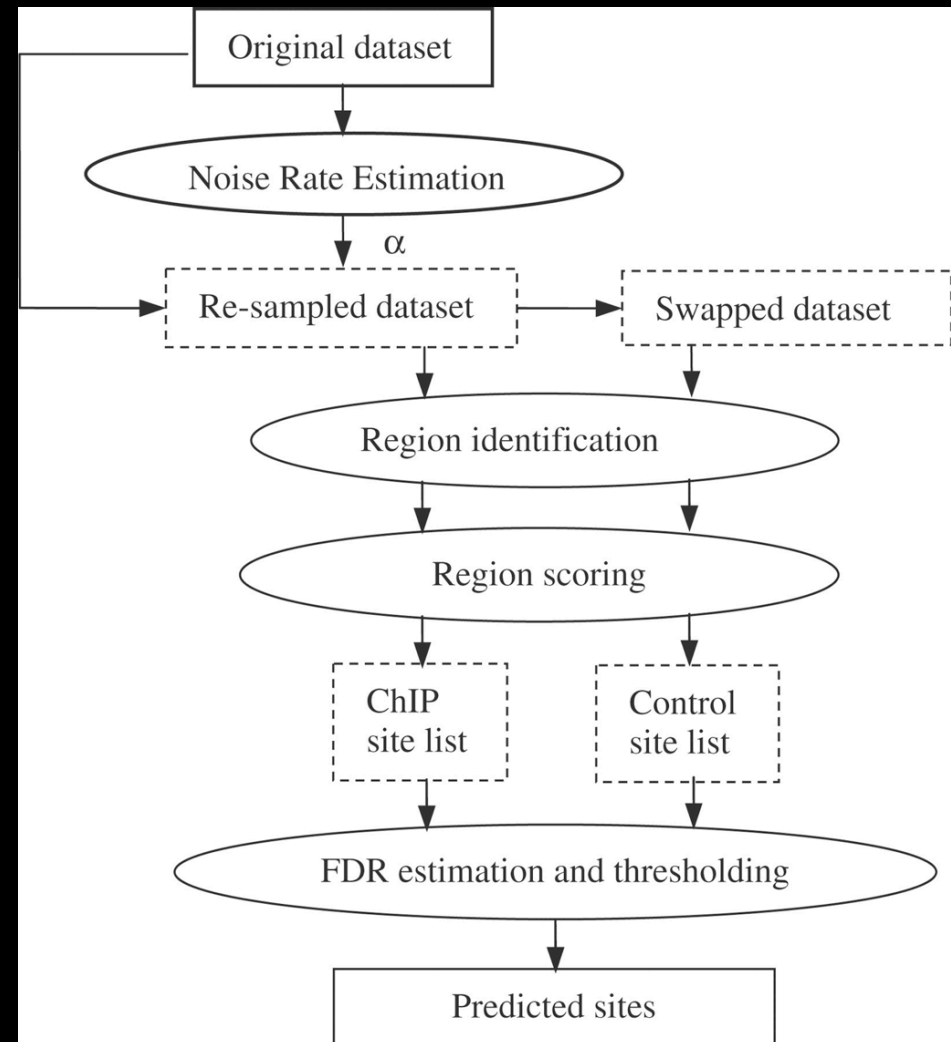
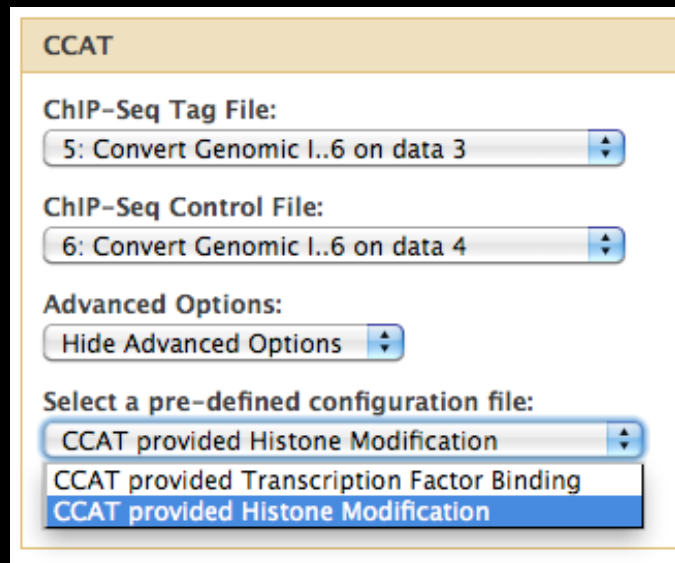


MACS --> GeneTrack

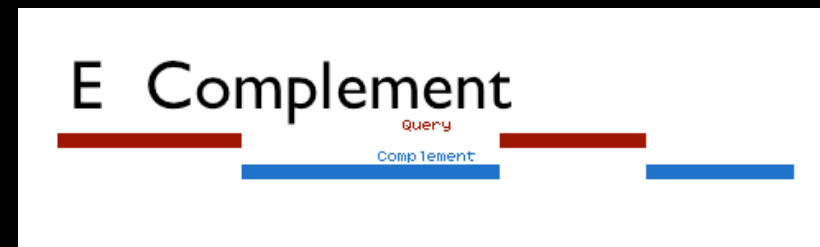
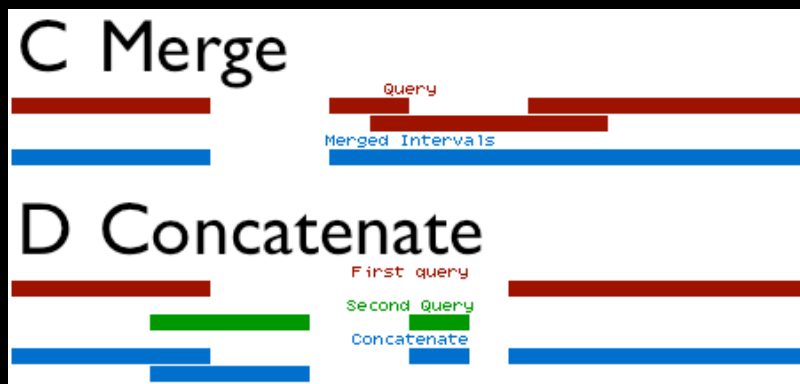
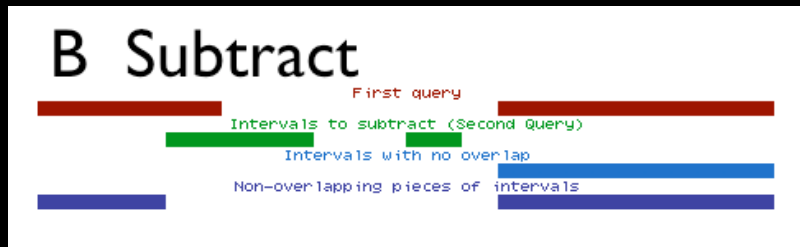
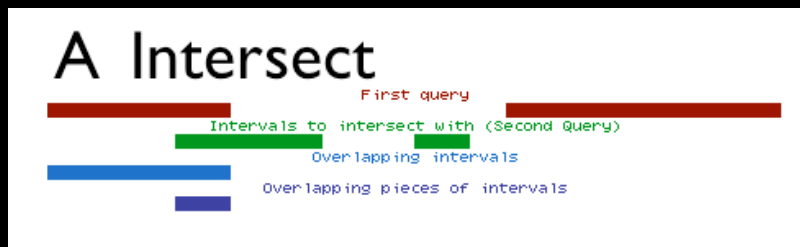


Diffuse Binding

CCAT (Control-based ChIP-seq Analysis Tool)



I have Peaks, now what?



Compare to other annotations using interval operations

Secondary Analysis

A simple goal: determine number of peaks that overlap a) **coding exons**, b) **5-UTRs**, c) **3-UTRs**, d) **introns** and d) **other** regions

Get Data

Import Peak Call data

Retrieve Gene location data from external data resource

Extract exon and intron data from Gene Data (**Gene BED To Exon/Intron/Codon BED expander** x4)

Create an Identifier column for each exon type (**Add column** x4)

Create a single file containing the 4 types (**Concatenate**)

Complement the exon/intron intervals

Force complemented file to match format of Gene BED expander output (**convert to BED6**)

Create an Identifier column for the 'other' type (**Add column**)

Concatenate the exons/introns and other files

Determine which Peaks overlap the region types (**Join**)

Calculate counts for each region type (**Group**)

Secondary Analysis

Galaxy Analyze Data Workflow Shared Data Admin Help User

Tools Options ▾

- [Get Data](#)
- [Send Data](#)
- [ENCODE Tools](#)
- [Lift-Over](#)
- [Text Manipulation](#)
- [Filter and Sort](#)
- [Join, Subtract and Group](#)
 - [Join two Queries](#) side by side on a specified field
 - [Compare two Queries](#) to find common or distinct rows
 - [Subtract Whole Query](#) from another query
 - [Group](#) data by a column and perform aggregate operation on other columns.
 - [Column Join](#)
- [Convert Formats](#)
- [Extract Features](#)
- [Fetch Sequences](#)
- [Fetch Alignments](#)
- [Get Genomic Scores](#)
- [Operate on Genomic Intervals](#)
- [Statistics](#)
- [Wavelet Analysis](#)
- [Graph/Display Data](#)
- [Regional Variation](#)
- [Multiple regression](#)
- [Multivariate Analysis](#)
- [Evolution](#)

3 UTR	803	
5 UTR	574	
coding exons		2743
introns	13746	
other	12499	

History Options ▾

[2: MACS peak calls \(broadPeak\)](#)

21,728 regions, format: interval, database: mm9

Info:

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

1.Chrom	2.Start	3.End	4	5	6	7	8	9
chr1	4132666	4133002	.	0	.	16.04	14.366	0.
chr1	4322446	4323079	.	0	.	27.07	26.185	0.
chr1	4336241	4336651	.	0	.	23.06	18.736	0.
chr1	4406740	4407268	.	0	.	16.20	23.794	0.
chr1	4506655	4507162	.	0	.	20.30	21.868	0.
chr1	4758431	4758873	.	0	.	24.01	30.691	0.

[1: UCSC Main on Mouse: refGene \(genome\)](#)

28,108 regions, format: bed, database: mm9

Info: UCSC Main on Mouse: refGene (genome)

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

1.Chrom	2.Start	3.End	4.Name	5	6..
chr1	134212701	134230065	NM_028778	0	+
chr1	134212701	134230065	NM_001195025	0	+
chr1	33510655	33726603	NM_008922	0	-
chr1	58714963	58752833	NM_175370	0	-
chr1	25124320	25886552	NM_175642	0	-
160945,328960,353082,363947,364951,389516,393...					

Annotation Profiler

One click to determine base coverage of the interval (or set of intervals) by a set of features (tables) available from UCSC galGal3, mm8, panTro2, rn4, canFam2, hg18, hg19, mm9, rheMac2

Profile Annotations

Choose Intervals:
34: UCSC Main on Mous..na (genome) ▾

Keep Region/Table Pairs with 0 Coverage:
Discard ▾

Output per Region/Summary:
Per Region ▾

Choose Tables to Use:

- [+] Comparative Genomics
- [+] Genes and Gene Prediction Tracks
- [+] Mapping and Sequencing Tracks
- [+] Phenotype and Allele
- [+] Expression and Regulation
- [+] mRNA and EST Tracks
- [-] Variation and Repeats
 - Microsatellite
 - Simple Repeats
 - SNPs (128)
- [+] Uncategorized Tables

Selecting no tables will result in using all tables.

Overview

High-throughput Sequencing (HTS) Data

Using Galaxy to Analyze HTS Data

- ✦ Prepare, quality control and manipulate reads
- ✦ Read Mapping
- ✦ SNP & INDEL analysis
- ✦ Binding sites analysis and peak calling
- ✦ **Transcriptome analysis**

Galaxy exercises: ChIP-seq and RNA-seq

Transcriptome Analysis (with a reference genome)

TopHat

Cufflinks/compare/diff

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

1. Trapnell, C., Pachter, L. and Salzberg, S.L. TopHat: discovering splice junctions with RNA-Seq. *Bioinformatics* 25, 1105-1111 (2009).
2. Trapnell et al. Transcript assembly and abundance estimation from RNA-Seq reveals thousands of new transcripts and switching among isoforms. *Nature Biotechnology* doi:10.1038/nbt.1621

TopHat

Map RNA (FASTQ) to a reference Genome

- ✦ gapped mapper

Outputs

- ✦ BAM file of accepted hits
- ✦ BED file of splice junctions

TopHat

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

Built-ins were indexed using default options

Select a reference genome:
Human (Homo sapiens): hg18 Canonical

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

Is this library mate-paired?:

RNA-Seq FASTQ file:
1: imported: h1-hESC..ple Dataset

Must have Sanger-scaled quality values with ASCII offset 33

TopHat settings to use:

You can use the default settings or set custom values for any of Tophat's parameters.

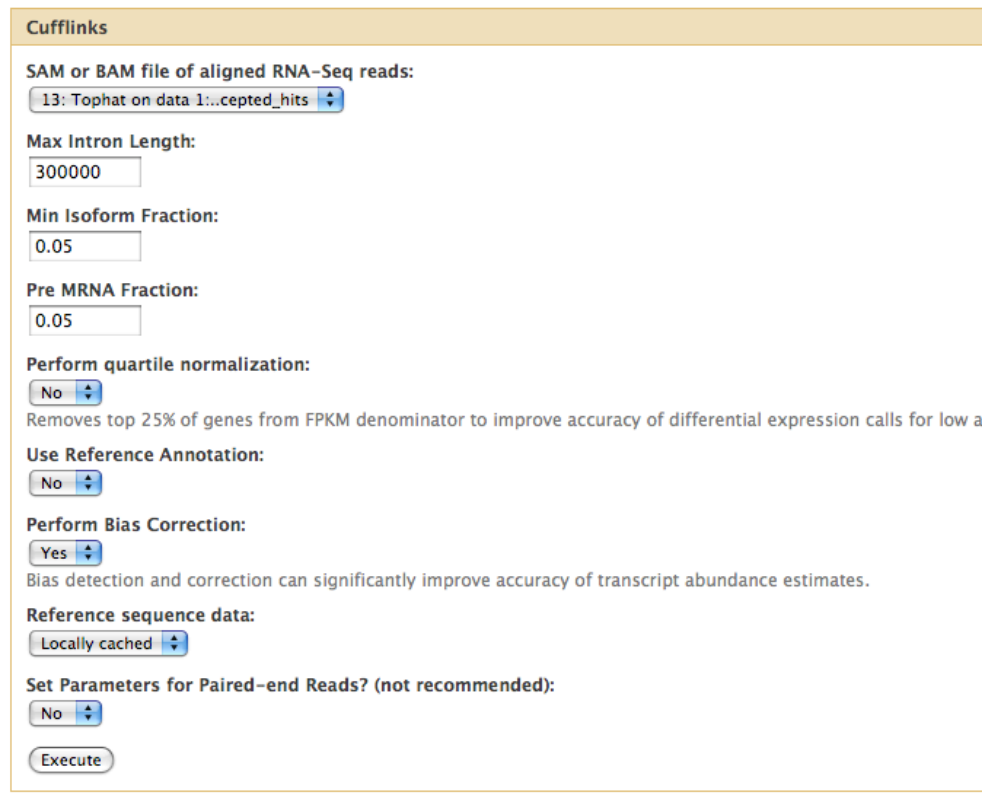
Cufflinks

Goal: transcript assembly and quantitation

Input: aligned RNA-Seq reads, usually from TopHat

Outputs

- ✦ assembled transcripts (GTF)
- ✦ genes' and transcripts' coordinates, expression levels



Cufflinks

SAM or BAM file of aligned RNA-Seq reads:
13: Tophat on data 1...cepted_hits

Max Intron Length:
300000

Min Isoform Fraction:
0.05

Pre MRNA Fraction:
0.05

Perform quartile normalization:
No
Removes top 25% of genes from FPKM denominator to improve accuracy of differential expression calls for low a

Use Reference Annotation:
No

Perform Bias Correction:
Yes
Bias detection and correction can significantly improve accuracy of transcript abundance estimates.

Reference sequence data:
Locally cached

Set Parameters for Paired-end Reads? (not recommended):
No

Execute

Cuffcompare

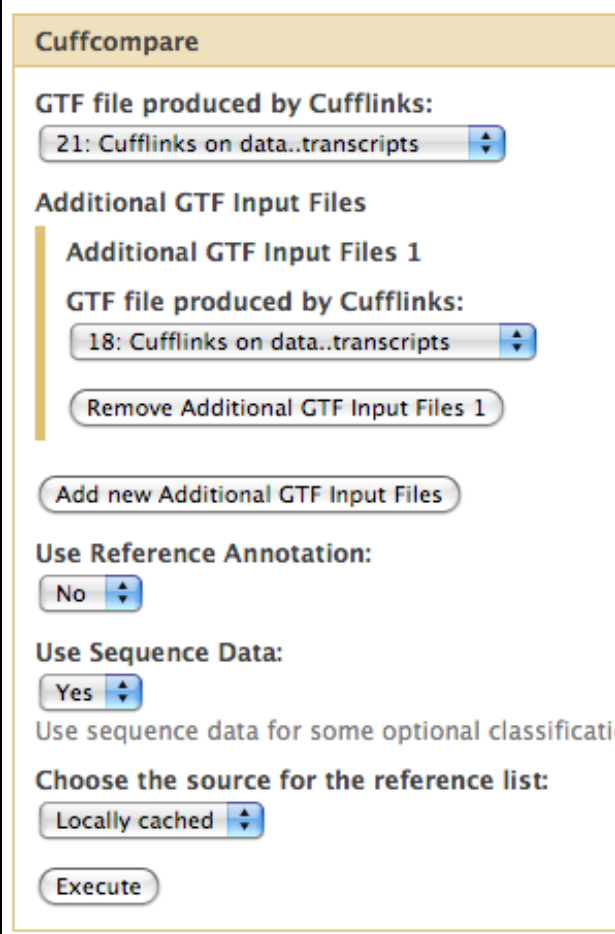
Goals

- ✦ generate complete list of transcripts for a set of transcripts
- ✦ compare assembled transcripts to a reference annotation

Inputs: assembled transcripts from Cufflinks

Outputs:

- ✦ Transcripts Combined File
- ✦ Transcripts Accuracy File
- ✦ Transcripts Tracking Files



The screenshot shows the Cuffcompare web interface with the following elements:

- Cuffcompare** (Title)
- GTF file produced by Cufflinks:** 21: Cufflinks on data..transcripts
- Additional GTF Input Files** (Section Header)
- Additional GTF Input Files 1** (Section Header)
- GTF file produced by Cufflinks:** 18: Cufflinks on data..transcripts
- Remove Additional GTF Input Files 1** (Button)
- Add new Additional GTF Input Files** (Button)
- Use Reference Annotation:** No
- Use Sequence Data:** Yes
- Use sequence data for some optional classification** (Text)
- Choose the source for the reference list:** Locally cached
- Execute** (Button)

Cuffdiff

Goals

- ✦ differential expression testing
- ✦ transcript quantitation

Inputs

- ✦ Combined set of transcripts
- ✦ mapped reads from 2+ samples

Outputs

- ✦ differential expression tests for transcripts, genes, splicing, promoters, CDS
- ✦ quantitation values for most elements

Cuffdiff

Transcripts:

A transcript GTF file produced by cufflinks, cuffcompare, or other source.

Perform replicate analysis:

Perform cuffdiff with replicates in each group.

SAM or BAM file of aligned RNA-Seq reads:

SAM or BAM file of aligned RNA-Seq reads:

False Discovery Rate:

The allowed false discovery rate.

Min Alignment Count:

The minimum number of alignments in a locus for needed to conduct significance testing or

Perform quartile normalization:

Removes top 25% of genes from FPKM denominator to improve accuracy of differential expr

Perform Bias Correction:

Bias detection and correction can significantly improve accuracy of transcript abundance est

Reference sequence data:

Set Parameters for Paired-end Reads? (not recommended):

Next Steps

Filtering

- ✦ for differentially expressed elements
- ✦ combined transcripts (e.g. for those differentially expressed between samples)

Extract transcript sequences and profile sequences for function

Filter Combined Transcripts

Cufflinks assembled transcripts:
130: Cuffcompare on da..transcripts

Cuffcompare tracking file:
130: Cuffcompare on da..transcripts

Sample Number:
1

Execute

Filter

Filter:
130: Cuffcompare on da..transcripts

Dataset missing? See TIP below.

With following condition:
c14=='yes'

Double equal signs, ==, must be used as

Execute

Extract Genomic DNA

Fetch sequences for intervals in:
130: Cuffcompare on da..transcripts

Interpret features when possible:
Yes

Only meaningful for GFF, GTF datasets.

Source for Genomic Data:
Locally cached

Output data type:
FASTA

Execute

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

Transcripts shown: 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

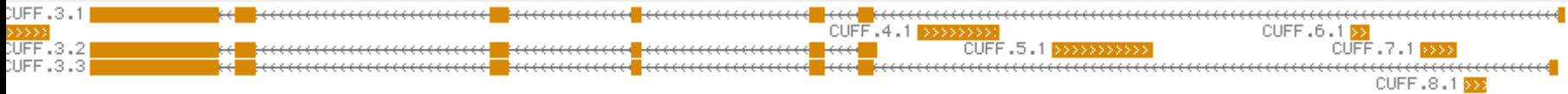
No ▼

Run on complete dataset

Run on visible region



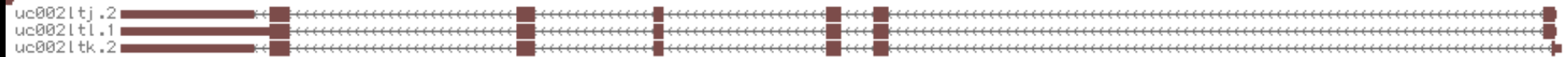
→ Cufflinks - region=[chr19:1523098-1545232], parameters=[150000, 0.05, 0.05, No] ▼



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]



Working to add GATK Unified Genotyper (and **more!**) to Trackster as well

Working with HTS Tools

Often challenging

- ✦ many parameters
- ✦ time intensive
- ✦ evaluating results difficult

Good options

- ✦ filter early, filter often: easier to understand fewer results
- ✦ experimentation: can rerun tools, workflows
- ✦ visualization: use tools in Trackster when possible

Overview

High-throughput Sequencing (HTS) Data

Using Galaxy to Analyze HTS Data

- ✦ Prepare, quality control and manipulate reads
- ✦ Read Mapping
- ✦ SNP & INDEL analysis
- ✦ Binding sites analysis and peak calling
- ✦ Transcriptome analysis

Galaxy exercises: ChIP-seq and RNA-seq



EMORY

PENNSYLVANIA STATE UNIVERSITY



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

Using Galaxy

Use public Galaxy server: UseGalaxy.org

Download Galaxy source: GetGalaxy.org

Galaxy Wiki: GalaxyProject.org

Screencasts: GalaxyCast.org

Public Mailing Lists

- ✦ galaxy-bugs@bx.psu.edu
- ✦ galaxy-user@bx.psu.edu
- ✦ galaxy-dev@bx.psu.edu

ChIP-seq and RNA-seq exercises

Chip-seq

- ✦ <http://usegalaxy.org/u/james/p/exercise-chip-seq>

RNA-seq

- ✦ <http://usegalaxy.org/u/jeremy/p/galaxy-rna-seq-analysis-exercise>
 - start Tophat mapping first (second section), then look at QC (first section)
- ✦ Add various outputs to a Trackster visualization and play with filtering and reruning tools

Variant Detection

