# Galaxy
## Web based platform for bioinformatics analysis

June 21, 2012

**Local copy:  https://galaxy.wi.mit.edu/**

**Joint project between BaRC and IT.**

Main site: http://main.g2.bx.psu.edu/

# Talk Outline

- **The Galaxy interface**
- Getting data into Galaxy
- Overview of the tools
- The Next Generation Sequencing tool box:
  - Preprocessing and quality control
  - Analysis of ChIP-seq
  - Analysis of RNA-seq
- Visualizing data on a genome browser and workflows available for analysis

*BaRC Hot Topics Galaxy*

# Galaxy Interface:

A web based platform for analysis of large genomic datasets



✓Type "https://galaxy.wi.mit.edu/" in your browser address.
✓You will be prompted for your name and password (these are the same that you use for your email)

✓No need of programming experience.

✓Integrates many bioinformatics tools within one interface.

✓Keeps track of all the steps performed in an analysis. Even if you delete the datasets, the history keeps the tools used.

LOCAL COPY

✓Faster
✓Customizable
✓250Gb of storage
✓Data is private
✓Jobs are sent to the cluster

# Galaxy Interface: Analyze Data

Data analysis



Tools window

Data display and tool's dialog window

History window: datasets for each analysis are kept here

Processed data

Green:
job is finished
Yellow:
job is running
Gray:
job is in queue
Red: there is a problem

# Galaxy Interface: Workflow

# Galaxy Interface: Shared Data

# Talk Outline

- The Galaxy interface
- **Getting data into Galaxy**
- Overview of the tools
- The Next Generation Sequencing tool box
  - Preprocessing and quality control
  - Analysis of ChIP-seq
  - Analysis of RNA-seq
- Visualizing data on a genome browser and workflows available for analysis

*BaRC Hot Topics Galaxy*

# Getting Data: Upload File

# Getting Data: Uploading Large Files
## Step 1: copy your file to
## /nfs/galaxy/uploads/username@wi.mit.edu
## using a sftp client



CyberDuck

# Getting Data: Uploading Large Files
## Step 2: Select and upload the file within galaxy

# Getting Data from UCSC (local copy)

# Getting Data from UCSC (local copy)

BaRC Hot Topics *Galaxy*

# Data Uploaded

Adapted from OpenHelix tutorial

BaRC Hot Topics *Galaxy*

# View Data

Adapted from OpenHelix tutorial

BaRC Hot Topics *Galaxy*

# Edit Attributes

# History

- All steps are saved.

- Every time we do a new operation a new dataset is created. Data is not overwritten.

- Can share history with other Galaxy users.

- Can create workflow to repeat an analysis.

BaRC Hot Topics *Galaxy*

# History

# History is not removed when datasets are removed

# Talk Outline

- The Galaxy interface
- Getting data into Galaxy
- **Overview of the tools**
- The Next Generation Sequencing tool box
  - Preprocessing and quality control
  - Analysis of ChIP-seq
  - Analysis of RNA-seq
- Visualizing data on a genome browser and workflows available for analysis

BaRC Hot Topics *Galaxy*

19

# Overview of the tools: Lift-Over

# Text Manipulation



**Galaxy / WIBR**

**Tools**          Options ▼

**Get Data**
**Lift-Over**
**Text Manipulation**
- Add column to an existing dataset
- Compute an expression on every row
- Concatenate datasets tail-to-head
- Cut columns from a table
- Merge Columns together
- Convert delimiters to TAB
- Create single interval as a new dataset
- Change Case of selected columns

- Paste two files side by side
- Remove beginning of a file
- Select random lines from a file
- Select first lines from a dataset
- Select last lines from a dataset
- Trim leading or trailing characters
- Line/Word/Character count of a dataset
- Secure Hash / Message Digest on a dataset
- Filter on ambiguities in polymorphism datasets
- Arithmetic Operations on tables

# Filter and Sort:
# Filter data on any column

# Convert Formats:
# GFF-to-BED

# Operate on Genomic Intervals:
# Intersect the intervals of two datasets

# Operate on Genomic Intervals:
# Intersect the intervals of two datasets

BaRC Hot Topics *Galaxy*

# Other tools

Fetch Sequences and
Fetch Alignments

FASTA manipulation

Join, Subtract and Group

*BaRC Hot Topics Galaxy*

# Talk Outline

- The Galaxy interface
- Getting data into Galaxy
- Overview of the tools
- **The Next Generation Sequencing tool box**
  - **Preprocessing and quality control**
  - Analysis of ChIP-seq
  - Analysis of RNA-seq
- Visualizing data on a genome browser and workflows available for analysis

# NGS Tools



**Galaxy / WIBR**

Tools                    Options ▼

**Get Data**
**Lift-Over**
**Text Manipulation**
**Filter and Sort**
**Join, Subtract and Group**
**Convert Formats**
**Extract Features**
**Fetch Sequences**
**Fetch Alignments**
**Get Genomic Scores**
**Operate on Genomic Intervals**
**FASTA manipulation**
**NGS: QC and manipulation**
**NGS: Mapping**
**NGS: RNA Analysis**
**NGS: SAM Tools**
**NGS: Peak Calling**

**Workflows**

Next Generation Sequencing Tools

# NGS: QC and manipulation

## Galaxy / WIBR

**Tools** | **Options**

**NGS: QC and manipulation**

FASTQC: FASTQ/SAM/BAM

- Fastqc: Fastqc QC using FastQC from Babraham

ILLUMINA FASTQ

- 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

GENERIC FASTQ MANIPULATION

## Galaxy / WIBR

**Tools** | **Options**

column

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

## 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
- Barcode Splitter
- Clip adapter sequences
- Collapse sequences
- Rename sequences
- Reverse-Complement
- Trim sequences

29

# Illumina data format

- Fastq format:

/1 or /2 paired-end

```
@ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
GTAGAACTGGTACGGACAAGGGGAATCTGACTGTAG
+ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
hhhhhhhhhhhghhhhhhhehhhedhhhhfhhhhhh
```

@seq identifier

seq

+any description

seq quality values

# Sequence quality values on different FASTQ formats

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

```
  SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS................................
  ...........................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
  ............................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII
  ............................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ
  !"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefgh
  |                         |    |         |                              |
  33                        59   64        73                            104

  S - Sanger        Phred+33,  raw reads typically (0, 40)
  X - Solexa        Solexa+64, raw reads typically (-5, 40)
  I - Illumina 1.3+ Phred+64,  raw reads typically (0, 40)
  J - Illumina 1.5+ Phred+64,  raw reads typically (3, 40)
```

To discriminate between Solexa and Illumina 1.3+ check if your sequences' quality scores have any of the characters ;<=>?

BaRC Hot Topics *Galaxy*

# FASTQ formats and FASTQ Groomer

**ILLUMINA FASTQ**

- FASTQ Groomer convert between various FASTQ quality formats

# NGS: Quality Control



**NGS: QC and manipulation**

FASTQC: FASTQ/SAM/BAM

- Fastqc: Fastqc QC using FastQC from Babraham

http://www.bioinformatics.babraham.ac.uk/projects/fastqc/

# Talk Outline

- The Galaxy interface
- Getting data into Galaxy
- Overview of the tools
- **The Next Generation Sequencing tool box**
  - Preprocessing and quality control
  - **Analysis of ChIP-seq**
  - Analysis of RNA-seq
- Visualizing data on a genome browser and workflows available for analysis

*BaRC Hot Topics Galaxy*

# Analysis of ChIP-seq experiments

**Preprocessing**
Convert to FASTQ Sanger format        FASTQ Groomer

↓

**Step 1**
Map reads        Bowtie

↓

**Step 2**
Call peaks bound        MACS

*BaRC Hot Topics Galaxy*

# Mapping Reads with **Bowtie**

# Mapping Reads with **Bowtie**

*BaRC Hot Topics Galaxy Next Gen.Seq.*

# Analysis of ChIP-seq experiments: **MACS**

# MACS output

# MACS output

# Creating Workflows

# Workflow for ChIP-seq analysis

# Example of downstream analysis: Intersect intervals of two datasets



Genes, promoters, or regions of interest

ChIP-seq peaks

Genes overlapping with peaks

http://jura.wi.mit.edu/bio/education/hot_topics/galaxy/Galaxy.pdf

BaRC Hot Topics *Galaxy*

# Talk Outline

- The Galaxy interface
- Getting data into Galaxy
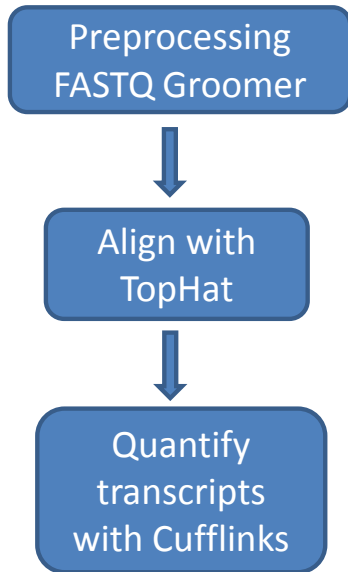- Overview of the tools
- **The Next Generation Sequencing tool box**
  - Preprocessing and quality control
  - Analysis of ChIP-seq
  - **Analysis of RNA-seq**
- Visualizing data on a genome browser and workflows available for analysis
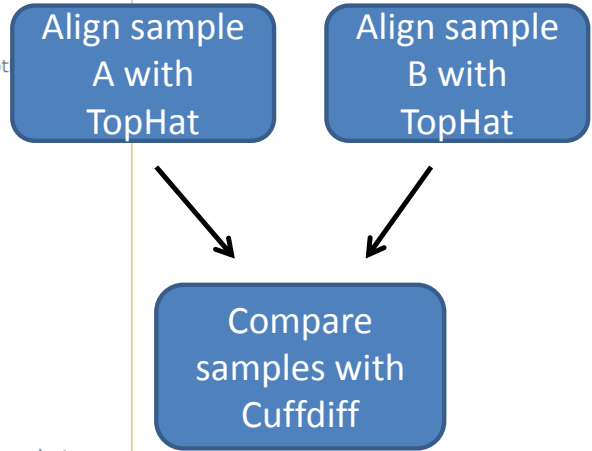
# Expression Profiling Workflow



**Preprocessing FASTQ Groomer**

**Align with TopHat**

**Quantify transcripts with Cufflinks**

*BaRC Hot Topics Galaxy*

# Other tools for expression profiling



**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.

Align sample A with TopHat

Align sample B with TopHat

Compare samples with Cuffdiff

Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat Protoc. 2012 Mar 1;7(3):562-78* . PMID: 22383036

# Talk Outline

- The Galaxy interface
- Getting data into Galaxy
- Overview of the tools
- The Next Generation Sequencing tool box
  – Preprocessing and quality control
  – Analysis of ChIP-seq
  – Analysis of RNA-seq
- Visualizing data on a genome browser and workflows available for analysis

BaRC Hot Topics *Galaxy*

# Visualizing data on UCSC

# Visualizing data on UCSC

BaRC Hot Topics *Galaxy*

# Visualizing data on UCSC: BAM files

# Workflows available inside Whitehead

1. Workflow for ChIP-seq analysis.

2. Workflow for sorting a SAM file.

BaRC Hot Topics *Galaxy*

# Documentation and Tutorials

- OpenHelix tutorials and exercises

  http://www.openhelix.com/cgi/tutorialInfo.cgi?id=82

- Galaxy tutorials

  http://galaxy.psu.edu/screencasts.html

- References

  Galaxy developers: The Center for Comparative Genomics & Bioinformatics, Pennsylvania State University

  Giardine, B., et al. Galaxy: a platform for interactive large-scale analysis. Genome Research (2005) 15:1451-1455

  Taylor, J., et al. Using Galaxy to perform large-scale interactive data analyses. Current Protocols in Bioinformatics (2007) Chapter 10, unit 10.

  Blankenberg D., et al. Manipulation of FASTQ data with Galaxy. Bioinformatics. 2010 Jul 15;26(14):1783-5

**BaRC Hot Topics** *Galaxy*

# Previous Hot Topics

- Previous Hot Topics in Galaxy

  http://jura.wi.mit.edu/bio/education/hot_topics/galaxy/Galaxy.pdf

  http://jura.wi.mit.edu/bio/education/hot_topics/GalaxyNGS/Galaxy_NGS.pdf

- Previous Hot Topics in NGS

  http://jura.wi.mit.edu/bio/education/hot_topics/shortRead_mapping/Mapping_HTseq.pdf

  http://jura.wi.mit.edu/bio/education/hot_topics/ChIPseq/ChIPSeq_HotTopics.pdf

  http://jura.wi.mit.edu/bio/education/hot_topics/RNAseq/RNA_Seq.pdf

*BaRC Hot Topics Galaxy*