# Analysis of next generation sequencing experiments with Galaxy

March 24, 2011

# Previous Hot Topics on Next Generation Sequencing Analysis

- Mapping next generation sequence reads

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

- Analysis of ChIP-seq experiments

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

- RNA-seq: Methods and Applications

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

# Talk Outline

- Introduction to Galaxy

- Data upload

- Format conversion and quality control tools

- Analysis of ChIP-seq experiments with MACS

- Analysis of RNA-seq experiments with Tuxedo tools

- Demo

**BaRC Hot Topics** *Galaxy Next Gen.Seq.*

# What is Galaxy

- A web based platform for analysis of large genomic datasets
- No need of programming experience.
- Integrates many tools within one interface:
  - Easy retrieval of data from UCSC, Biomart and other DBs
  - Powerful text manipulation tools (data preparation)
  - Filter on columns, join, sort, compute etc
  - Format conversion tools (text, tab, bed, GFF …)
  - Integrates tools from other sources. Ex: EMBOSS
  - MSA tools
  - Visualize data in UCSC browser.
    (See Hot topics Dec 09, http://iona.wi.mit.edu/bio/education/hot_topics/galaxy/Galaxy.pdf)

  - **Next Generation Sequencing Toolbox**

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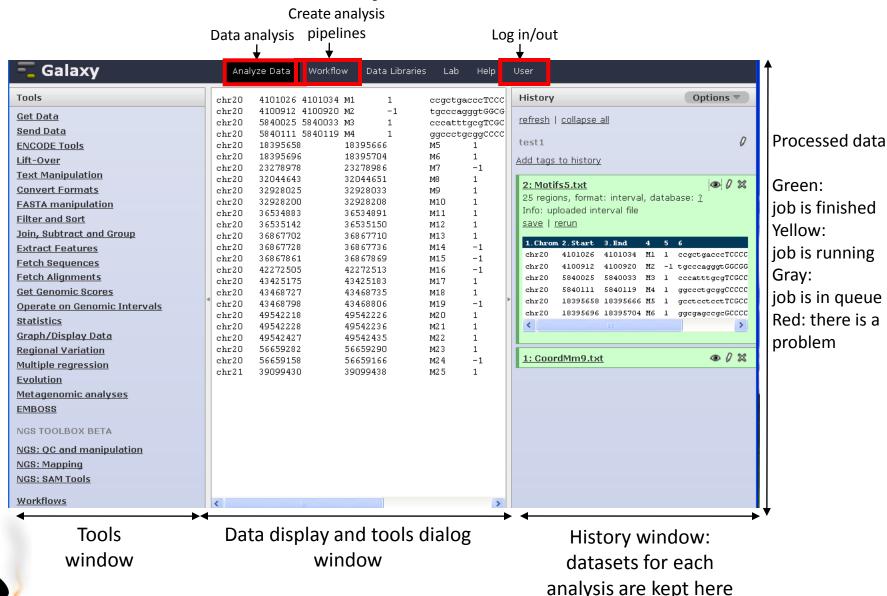
Giardine, B., et al. Galaxy: a platform for interactive large-scale analysis. Genome Research (2005) 15:1451-1455

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

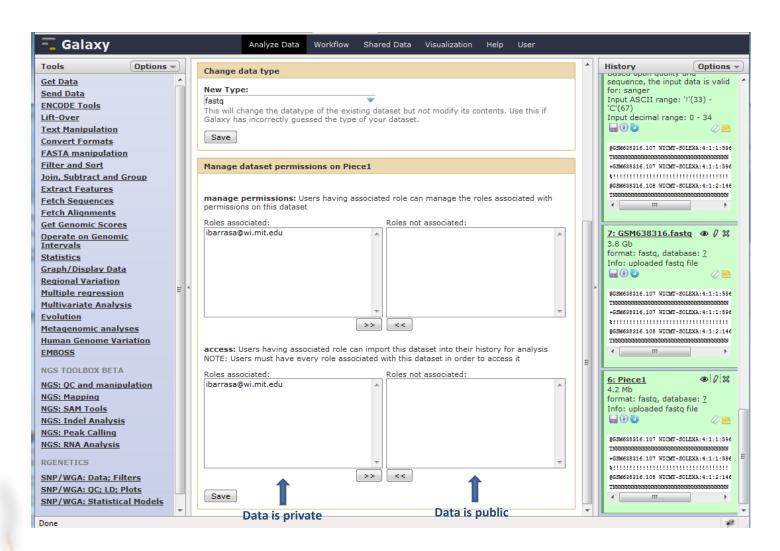Blankenberg D., et al. Manipulation of FASTQ data with Galaxy. Bioinformatics. 2010 Jul 15;26(14):1783-5

*BaRC Hot Topics Galaxy Next Gen.Seq.*

# Galaxy Interface



Create analysis pipelines

Data analysis

Log in/out

Tools window

Data display and tools 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

# Security issues

- Need to register to be able to keep your data and history (log in button).

- Your data has to be public to be able to be visualized at UCSC. By default the data is public.

- You could make your data private, download it and visualize in UCSC or other browser.

# Security issues II

# Talk Outline

- Introduction to Galaxy

- Data upload

- Format conversion and quality control tools

- Analysis of ChIP-seq experiments with MACS

- Analysis of RNA-seq experiments with Tuxedo tools

- Demo

# Data upload I

- For files larger than 2Gb, transfer to the Galaxy server via the file transfer protocol (FTP).
- Log in to tak (**ssh –l userName tak.wi.mit.edu**), and **cd** to the folder that has your files. (See hot topic "introduction to Unix" http://iona.wi.mit.edu/bio/education/hot_topics/unix_2010/slides.pdf)
- Ftp to Galaxy:

  **ftp main.g2.bx.psu.edu**
  *Name (main.g2.bx.psu.edu:ibarrasa):*  **Type your email**
  *Password:* **Type your Galaxy password**

  *230 User ibarrasa@wi.mit.edu logged in*
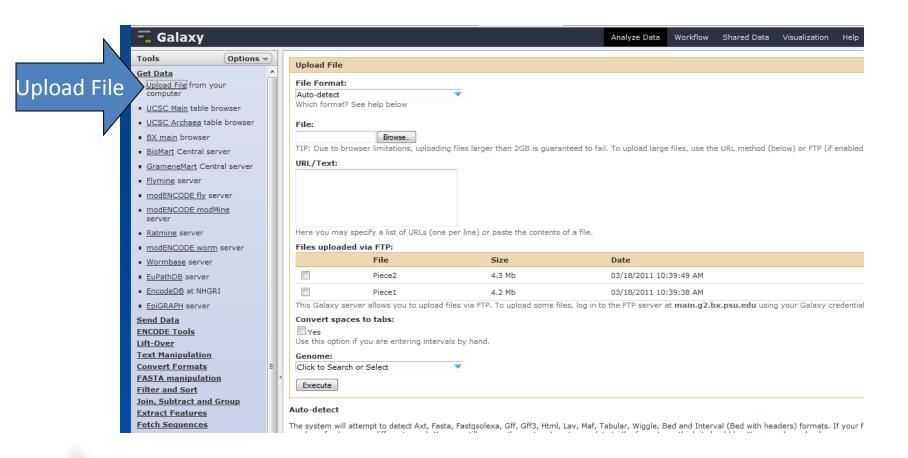  *Remote system type is UNIX.*
  *Using binary mode to transfer files.*
  ftp>

- Upload file

  ftp>  **put FileName**
  ftp> **exit**
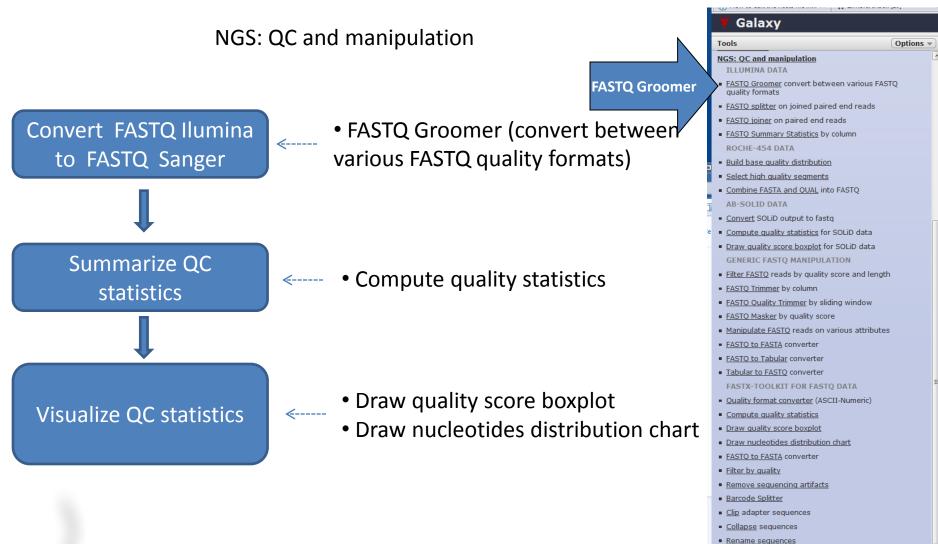
# Data upload II

# Talk Outline

- Introduction to Galaxy

- Data upload

- **Format conversion and quality control tools**

- Analysis of ChIP-seq experiments with MACS

- Analysis of RNA-seq experiments with Tuxedo tools

- Demo

*BaRC Hot Topics Galaxy Next Gen.Seq.*

# Format conversion and quality control tools

NGS: QC and manipulation

**FASTQ Groomer**

**Convert FASTQ Ilumina to FASTQ Sanger** ⟵------ • FASTQ Groomer (convert between various FASTQ quality formats)

**Summarize QC statistics** ⟵------ • Compute quality statistics

**Visualize QC statistics** ⟵------ • Draw quality score boxplot
• Draw nucleotides distribution chart

Note: FastQC is not incorporated in Galaxy but it is installed in tak .

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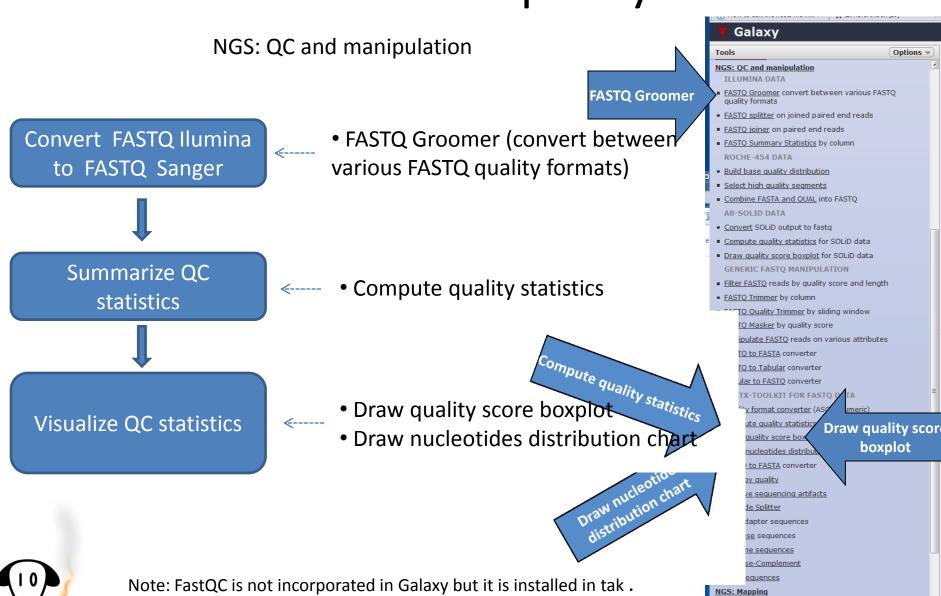
http://en.wikipedia.org/wiki/FASTQ_format

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS...............................
.........................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
............................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII
...............................JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ
!"#$%&'()*+,-./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 have any of the characters: :;<=>?

*BaRC Hot Topics Galaxy Next Gen.Seq.*

# FASTQ Groomer

# Format conversion and quality control tools

NGS: QC and manipulation

Convert FASTQ Ilumina to FASTQ Sanger

• FASTQ Groomer (convert between various FASTQ quality formats)

Summarize QC statistics

• Compute quality statistics

Visualize QC statistics

• Draw quality score boxplot
• Draw nucleotides distribution chart

**FASTQ Groomer**

**Compute quality statistics**

**Draw quality score boxplot**

**Draw nucleotides distribution chart**

**Galaxy**

Tools | Options

**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
- ...TQ Masker by quality score
- ...ipulate FASTQ reads on various attributes
- ...TQ to FASTA converter
- ...TQ to Tabular converter
- ...ular to FASTQ converter
FASTX-TOOLKIT FOR FASTQ DATA
- ...ty format converter (AS... ...meric)
- ...ute quality statistics
- ... quality score box...
- ...nucleotides distribu...
- ... to FASTA converter
- ...y quality
- ...ve sequencing artifacts
- ...de Splitter
- ...dapter sequences
- ...se sequences
- ...ne sequences
- ...se-Complement
- ...equences

**NGS: Mapping**
**NGS: Indel Analysis**
**NGS: RNA Analysis**
Done

Note: FastQC is not incorporated in Galaxy but it is installed in tak .

17

# Quality control visualization tools

**Draw quality score boxplot**

**Draw nucleotides distribution chart**

# How to make a workflow from the history

# Workflow for Quality Control

*BaRC Hot Topics Galaxy Next Gen.Seq.*
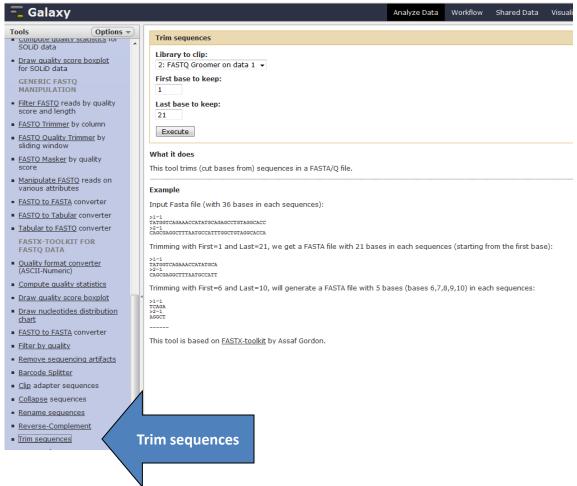
# Remove sequencing artifacts

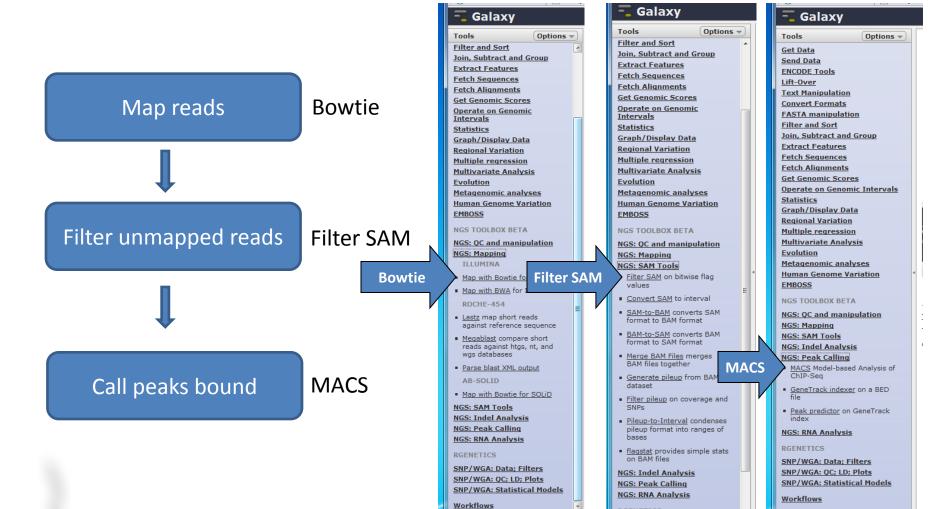# Clip adapter sequences

BaRC Hot Topics *Galaxy Next Gen.Seq.*

# Trim sequences

# Talk Outline

- Introduction to Galaxy

- Data upload

- Format conversion and quality control tools

- **Analysis of ChIP-seq experiments with MACS**

- Analysis of RNA-seq experiments with Tuxedo tools

- Demo

*BaRC Hot Topics Galaxy Next Gen.Seq.*

# Analysis of ChIP-seq experiments



Map reads — Bowtie

Filter unmapped reads — Filter SAM

Call peaks bound — MACS

*BaRC Hot Topics Galaxy Next Gen.Seq.*

# Mapping Reads with Bowtie

# Filtering unmapped reads

*BaRC Hot Topics* *Galaxy Next Gen.Seq.*

# Analysis of ChIP-seq experiments: MACS

*BaRC Hot Topics Galaxy Next Gen.Seq.*

# Workflow for ChIP-seq analysis
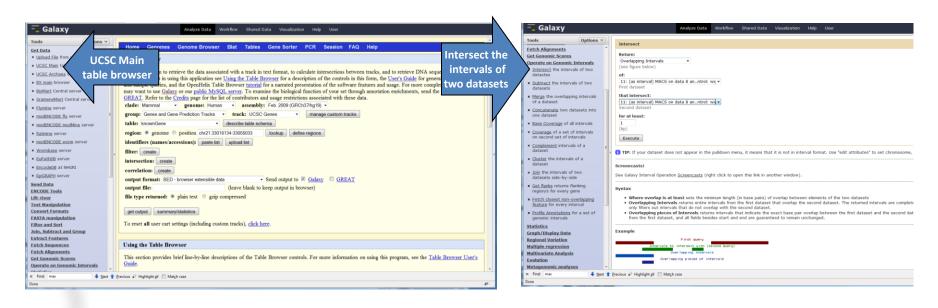
*BaRC Hot Topics Galaxy Next Gen.Seq.*

# MACS output

# Analysis of ChIP-seq experiments: Intersect peaks with promoter regions

1. Download 1Kb regions upstream of genes from UCSC in bed format.
2. Get your bed file with peaks from MACS or other peak finding algorithm.
3. Intersect promoter bed file with peaks bed file.

(See Hot topics Dec 09,

http://iona.wi.mit.edu/bio/education/hot_topics/galaxy/Galaxy.pdf)
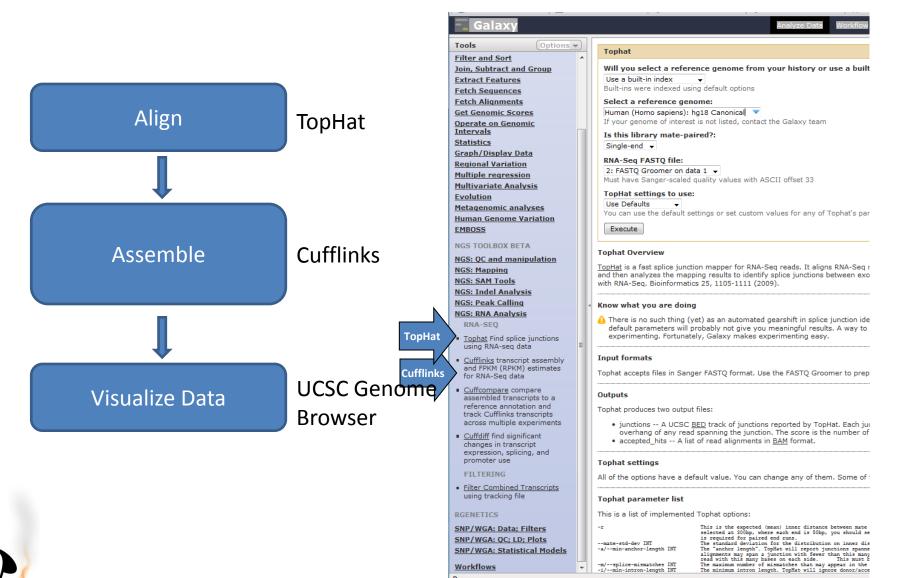
*BaRC Hot Topics Galaxy Next Gen.Seq.*

# Talk Outline

- Introduction to Galaxy

- Data upload

- Format conversion and quality control tools

- Mapping

- Analysis of ChIP-seq experiments with MACs

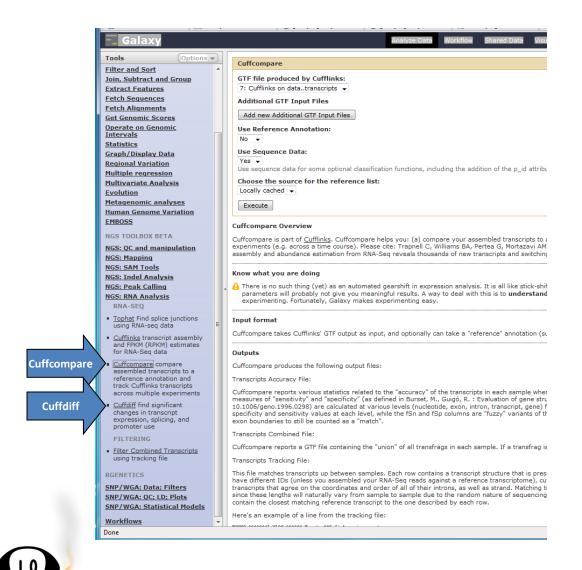- **Analysis of RNA-seq experiments with Tuxedo tools**

- Demo

*BaRC Hot Topics Galaxy Next Gen.Seq.*

# Expression Profiling Workflow
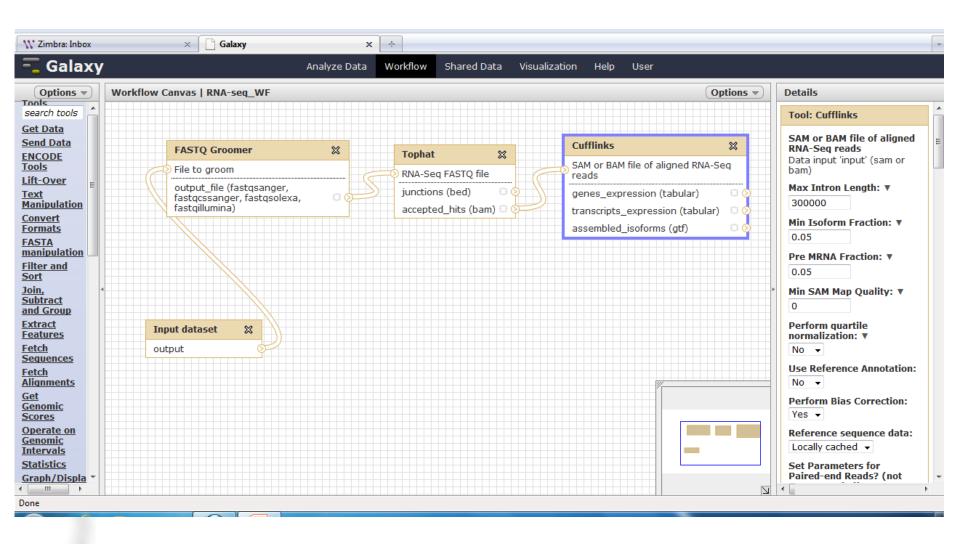
# Other tools for expression profiling



- 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

*BaRC Hot Topics Galaxy Next Gen.Seq.*

# Workflow for RNA-seq analysis

*BaRC Hot Topics Galaxy Next Gen.Seq.*

# Workflow/Demo for ChIP-seq analysis

1. Workflow for quality control
2. Workflow for mapping and running MACS
3. Workflow for RNA-seq

BaRC Hot Topics *Galaxy Next Gen.Seq.*