

# Introduction to Galaxy and RNA-Seq workshop

Enis Afgan  
Johns Hopkins University

Slides available at [bit.ly/Gxy-RNA-Seq-J15](http://bit.ly/Gxy-RNA-Seq-J15)

# The Agenda

Galaxy project overview: the core pillars

RNA-Seq exercise

# What is Galaxy?

- A **data analysis and integration** tool
- A **(free for everyone) web service** integrating a wealth of tools, compute resources, terabytes of reference data and permanent storage
- **Open source software** that makes integrating your own tools and data and customizing for your own site simple

# Galaxy project: fundamental questions

When Biology (or any science) becomes dependent on computational methods:

- How can those methods best be made **accessible** to scientists?
- How best to ensure that analyses are **reproducible**?
- How best to facilitate **transparent** communication and reuse of analyses?

Tools



search tools

Get Data

Lift-Over

Text Manipulation

Convert Formats

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

Motif Tools

Multiple Alignments

Metagenomic analyses

Genome Diversity

NGS TOOLBOX BETA

Phenotype Association

NGS: QC and manipulation

NGS: Mapping

NGS: SAM Tools

NGS: GATK Tools (beta)

NGS: Peak Calling

NGS: RNA-seq

Galaxy is an open source, web-based platform for data intensive biomedical research. If you are new to Galaxy [start here](#) or consult our [help resources](#).

Want help?  
Get answers.



**Biostars**  
GALAXY EXPLAINED



## Tweets

[Follow](#)



**Guillaume Lobet** @guillaumelobet 16 Jan  
Great presentation of #usegalaxy by @jxtx at #NGS15 #openscience #reproducibility

Retweeted by Galaxy Project

Expand



**Mark Veugelers** @biotechbase 16 Jan  
#ngs15: JT on future developments for Galaxy, e.g. analysis of 10.000+ samples

Retweeted by Galaxy Project

Expand



**Mark Veugelers** @biotechbase 16 Jan  
#ngs15: James Taylor now speaking on reproducible genomics with Galaxy.  
[How to get reproducibility?](#) [Easy steps](#)

[Tweet to @galaxyproject](#)

PENN STATE



JOHNS HOPKINS  
UNIVERSITY

**TACC**

iPlant  
Collaborative™

The Galaxy Team is a part of the Center for Comparative Genomics and Bioinformatics at Penn State, and the Department of Biology and at Johns Hopkins University.

This instance of Galaxy is utilizing infrastructure generously provided by the iPlant Collaborative at the Texas Advanced Computing Center, with support from the

History



search datasets



Unnamed history

0 bytes



**i** This history is empty. You can [load your own data](#) or [get data from an external source](#)

# Need an analysis? There's a tool for that.


[Get Data](#)

- [Upload File from your computer](#)
- [UCSC Main table browser](#)
- [UCSC Test table browser](#)
- [UCSC Archaea table browser](#)
- [BX main browser](#)
- [EBI SRA ENA SRA](#)
- [Get Microbial Data](#)
- [BioMart Central server](#)
- [BioMart Test server](#)
- [CBI Rice Mart rice mart](#)
- [GrameneMart Central server](#)
- [modENCODE fly server](#)
- [Flymine server](#)
- [Flymine test server](#)
- [modENCODE modMine server](#)
- [Ratmine server](#)
- [YeastMine server](#)
- [metabolicMine server](#)
- [modENCODE worm server](#)
- [WormBase server](#)
- [Wormbase test server](#)
- [EuPathDB server](#)
- [EncodeDB at NHGRI](#)
- [EpiGRAPH server](#)
- [EpiGRAPH test server](#)

[NGS: Mapping](#)

- [Lastz map short reads against reference sequence](#)
- [Lastz paired reads map short paired reads against reference sequence](#)
- [Map with Bowtie for Illumina](#)
- [Map with Bowtie for SOLiD](#)
- [Map with BWA for Illumina](#)
- [Map with BWA for SOLiD](#)
- [Map with BFAST](#)
- [Megablast compare short reads against htgs, nt, and wgs databases](#)
- [Parse blast XML output](#)
- [Map with PerM for SOLiD and Illumina](#)
- [Re-align with SRMA](#)
- [Map with Mosaik](#)

[NGS: RNA Analysis](#)

- RNA-SEQ**
- [Tophat for Illumina Find splice junctions using RNA-seq data](#)
  - [Tophat for SOLiD 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](#)

[NGS: GATK Tools \(beta\)](#)

- ALIGNMENT UTILITIES**
- [Depth of Coverage on BAM files](#)
- 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 – draw plots](#)
- GENOTYPING**
- [Unified Genotyper SNP and indel caller](#)
- ANNOTATION**
- [Variant Annotator](#)
- FILTRATION**
- [Variant Filtration on VCF files](#)
- VARIANT QUALITY SCORE RECALIBRATION**
- [Variant Recalibrator](#)
  - [Apply Variant Recalibration](#)
- VARIANT UTILITIES**
- [Validate Variants](#)
  - [Eval Variants](#)
  - [Combine Variants](#)

# Getting the data

Accessibility

**Get Data**

- Upload File from your computer
- [UCSC Main table browser](#)
- [UCSC Test table browser](#)
- [UCSC Archaea table browser](#)
- [BX main browser](#)
- [EBI SRA ENA SRA](#)
- [Get Microbial Data](#)
- [BioMart Central server](#)
- [BioMart Test server](#)

Analyze Data   Workflow   Shared Data   Visualization   Admin   Help   User   Using 24.4 Gb

Use available data sources  
Upload from local system  
Upload from URL  
Use (shared) Data Library  
Add your own data source

Tools

search tools

**Get Data**

- Upload File from your computer
- [UCSC Main table browser](#)
- [UCSC Test table browser](#)
- [UCSC Archaea table browser](#)
- [BX main browser](#)
- [EBI SRA ENA SRA](#)
- [Get Microbial Data](#)
- [BioMart Central server](#)
- [BioMart Test server](#)
- [CBI BioMart](#)
- [Gramene](#)
- [modENCODE](#)
- [Flymine server](#)
- [Flymine test server](#)
- [modENCODE modMine server](#)
- [Ratmine server](#)
- [YeastMine server](#)
- [metabolicMine server](#)

Home   Genomes   Genome Browser   Blat   Tables   Gene Sorter   PCR   Session   FAQ   Help

## Table Browser

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). To examine the biological function of your set through annotation enrichments, send the data to [GREAT](#). Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data. All tables can be downloaded in their entirety from the [Sequence and Annotation Downloads](#) page.

clade: Mammal   genome: Human   assembly: Feb. 2009 (GRCh37/hg19)

group: Genes and Gene Prediction Tracks   track: UCSC Genes   add custom tracks   track hubs

table: knownGene   describe table schema

region:  genome  ENCODE Pilot regions  position chr21:33,031,597-33,041,!   lookup   define regions

identifiers (names/acccessions): paste list   upload list

filter: create

intersection: create

correlation: create

output format: BED – browser extensible data   Send output to  Galaxy  GREAT

output file: (leave blank to keep output in browser)

file type returned:  plain

get output   summary/statistics

To reset all user cart settings

Output knownGene as BED

Include custom track header:

name= tb\_knownGene  
description= table browser query on knownGene  
visibility= pack  
url=

Create one BED record per:

Whole Gene  
 Upstream by 200 bases  
 Exons plus 0 bases at each end  
 Introns plus 0 bases at each end  
 5' UTR Exons  
 Coding Exons  
 3' UTR Exons  
 Downstream by 200 bases

Note: if a feature is close to the beginning or end of a chromosome and upstream/downstream bases are added, they may be extending past the edge of the chromosome.

Send query to Galaxy

## Concatenate Dataset

Single dataset  Multiple datasets

84: Filter on data 33

## Dataset

1: Dataset

## Select

Single dataset  Multiple datasets

36: Cuffdiff on data 13, data 9, and data 1  
 35: Cuffdiff on data 13, data 9, and data 1  
 34: Cuffdiff on data 13, data 9, and data 1  
 33: Cuffdiff on data 13, data 9, and data 1

This is a batch mode input field. A

+ Insert Dataset

✓ Execute

⚠ WARNING: Be careful not to concatenate datasets being concatenated are in the same

## What it does

Concatenates datasets

## Example

Concatenating Dataset:

```
chrX 151087187 151087355 A 0 -
chrX 151572400 151572481 B 0 +
```

Dataset collection

```
catWrapper.xml x
1 <tool id="cat1" name="Concatenate datasets">
2   <description>tail-to-head</description>
3   <command interpreter="python">
4     catWrapper.py
5     $out_file1
6     $input1
7     #for $q in $queries
8       ${q.input2}
9     #end for
10    </command>
11    <inputs>
12      <param name="input1" type="data" label="Concatenate">
13        <repeat name="queries" title="Dataset">
14          <param name="input2" type="data" label="Select">
15        </repeat>
16      </inputs>
17      <outputs>
18        <data name="out_file1" format="input" metadata_sourc
19      </outputs>
20      <tests>
21        <test>
22          <param name="input1" value="1.bed"/>
23          <param name="input2" value="2.bed"/>
24          <output name="out_file1" file="cat_wrapper_out1.bed"/>
25        </test>
26        <!--TODO: if possible, enhance the underlying test code to handle this test
27          the problem is multiple params with the same name "input2"
28        <test>
29          <param name="input1" value="1.bed"/>
30          <param name="input2" value="2.bed"/>
31          <param name="input2" value="3.bed"/>
32          <output name="out_file1" file="cat_wrapper_out2.bed"/>
33        </test>
34        -->
35      </tests>
36      <help>
37
38 ... class:: warningmark
39
40 **WARNING:** Be careful not to concatenate datasets of different kinds (e.g., sequences with
41 intervals). This tool does not check if the datasets being concatenated are in the same format.
42
43
44 **What it does**
45
46 Concatenates datasets
47
48
49
50 **Example**
51
52 <Concatenating Dataset:>
```

## Adding a tool?

- Automatically generated web UI from a tool wrapper (any tool can be integrated)
- Integrated with other tools

# Accessibility

# Data analysis history

The screenshot displays the Galaxy web interface across four browser tabs, illustrating a complex pipeline for transcript assembly and expression analysis.

- Top Left Panel:** Shows a success message for three Cufflinks jobs (gene expression, transcript expression, and assembled transcripts) and a table of genomic data (QNAME, FLAG, RNAME, POS, MAPQC, CAR, MRNMMPSIZESEQ).
- Top Right Panel:** Displays the "Basic Protocol 4" history, listing eight Cufflinks jobs (gene/transcript expression and assembled transcripts) and a BAM file named "Binary bam alignments file".
- Middle Left Panel:** Shows a list of datasets and histories, including "3: Concatenate datasets on data 1 and data 2", "2: 3.maf", "3: 3.fasta", "11: 3.maf", "10: 123456789012345678901234567890", "9: MAF to Interval on data 4 (hg12)", "8: MAF to Interval on data 4 (panTro1)", "7: MAF to Interval on data 4 (mm5)", "6: MAF to Interval on data 4 (hg17)", "5: New Dataset List", "4: 3.bed", "3: 3.interval", "2: 3.fasta", and "1: 3.bed".
- Middle Right Panel:** Shows a history titled "Galaxy" containing a single history entry: "1: 1.VCF mm5 (hg12)".
- Bottom Left Panel:** Shows a list of datasets and histories, including "15: Filter data 1", "14: Concatenate datasets on data\_10 and data\_7", "13: MAF to Interval on data\_4 (mm5)", "12: MAF to Interval on data\_4 (panTro1)", "11: MAF to Interval on data\_4 (mm5)", "10: MAF to Interval on data\_4 (hg17)", "9: MAF to Interval on data\_4 (mm5) (canFam1)", "8: MAF to Interval on data\_4 (mm5) (panTro1)", "7: Lvcf Qual>60", "6: 3.maf", "5: 3.maf", "4: 3.maf", "3: 3.maf", "2: 3.maf", "1: 3.maf", and a table of genomic data.
- Bottom Right Panel:** Shows a history titled "Galaxy" containing a single history entry: "1: 1.VCF mm5 (hg12)".

# Reproducibility in Genomics

18 *Nat. Genetics* experiments in microarray gene expression

<50% of reproducible

Problems

- missing data (38%)
- missing software, hardware details (50%)
- missing methods, processing details (66%)

Ioannidis, J.P.A. et al. "Repeatability of published microarray gene expression analyses." *Nat Genet* 41, 149-155 (2009)

14 re-sequencing experiments in *Nat. Genetics, Nature, Science*

0% reproducible?

Problems

- missing primary data (50%)
- tools unavailable (50%)
- missing parameter setting, tool versions (100%)

"Devil in the details," *Nature*, vol. 470, 305-306 (2011).

Metadata = Reproducibility

# Automatic metadata

**95: Filter on data 94**

202 lines  
format: **tabular**, database: **hg19**

Filtering with c14=='yes',  
kept 44.40% of 455 valid lines (455 total lines).

**View details** 3 4

test_id	gene_id	gene	locus
AAAS	AAAS	AAAS	chr12:53701239-
ABCC9	ABCC9	ABCC9	chr12:21950323-
ABCD2	ABCD2	ABCD2	chr12:39945021-
ALKBH2	ALKBH2	ALKBH2	chr12:109525992
ANAPC5	ANAPC5	ANAPC5	chr12:121746047

Tool: Filter	
Name:	Filter on data 94
Created:	Sat Jan 17 23:34:56 2015 (UTC)
Filesize:	22.0 KB
Dbkey:	hg19
Format:	tabular
Galaxy Tool	1.1.0
Version:	
Tool Version:	
Tool	
Standard	<a href="#">stdout</a>
Output:	
Tool	
Standard	<a href="#">stderr</a>
Error:	
Tool Exit Code:	0
API ID:	fcc622f6b9afc8d2
Full Path:	/mnt/galaxy/files/001/dataset_1107.dat
Job Command-Line:	python /mnt/galaxy/galaxy-app/tools/stats/filtering.py /mnt/galaxy/files/001/dataset_1107.dat "c14==_sq_ye
Job Runtime (Wall Clock)	2 seconds
Cores Allocated	1
Job Start Time	2015-01-17 23:35:04
Job End Time	2015-01-17 23:35:06
Input Parameter	Value
Filter	94: Fil
With following condition	c14==
Number of header lines to skip	1

# 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  $\geq 30$  and quality score  $\geq 20$ .

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

1.Chrom	2.Start	3.End	4	5	6	7
chr10	6882036	6882037	A	A	107	1
chr10	14243075	14243076	G	G	96	1
chr10	14243079	14243080	C	C	106	1
chr10	14465082	14465083	T	K	173	1
chr10	14465083	14465084	G	K	144	1
chr10	14465084	14465085	T	T	117	1

# Data provenance

**9: Map with BWA on data**

**1**

1,447,842 lines  
format: sam, database: hg19

BWA run on single-end data

Run this job again 2.FLAG 3.RNAME 4.

GA5:3:1:23:402#0	16	chrM
XA:Z:chr17,-22023750,76M,2;		
GA5:3:1:23:532#0	0	chr1
:Z:chrM,+5913,76M,0;		
GA5:3:1:24:280#0	0	chrM
GA5:3:1:24:1454#0	16	chrM

This job was initially run with tool id "bwa\_wrapper", version "1.0.3", which is a derivation of the original tool.

Map with BWA for Illumina (version 1.2.3)

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

Select a reference genome:  
Arabidopsis lyrata: Araly1

An invalid option was selected, please verify

Is this library mate-paired?:  
Single-end

FASTQ file:  
1: pl-c-b-1.fastq

FASTQ with either Sanger-scaled quality values (fastqsanger) or Illumina-style quality values (fastqillumina).

BWA settings to use:  
Commonly Used

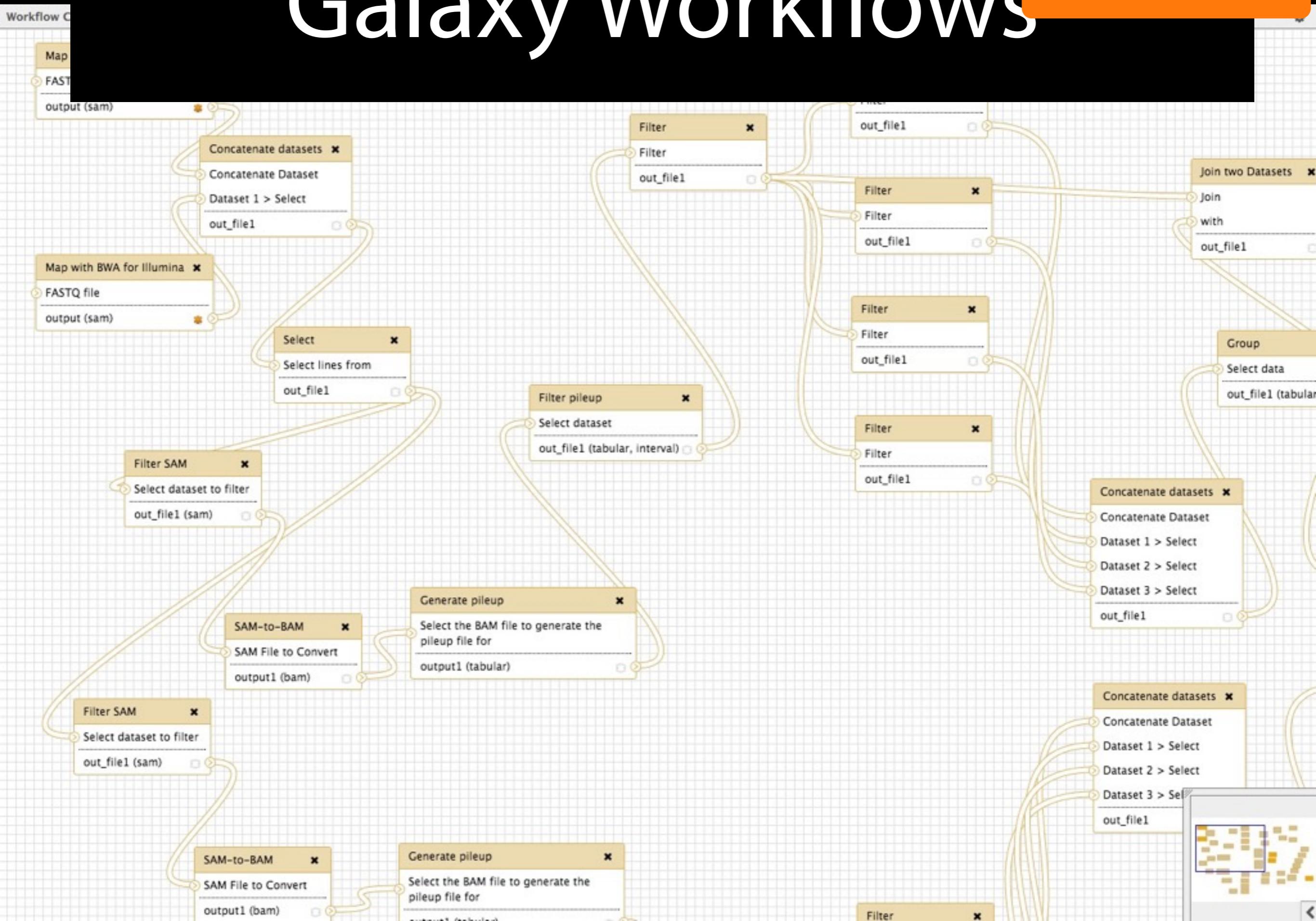
For most mapping needs use Commonly Used settings. If you want full control over the BWA command line options, click here.

Suppress the header in the output SAM file:

BWA produces SAM with several lines of header information

**Execute**

# Galaxy Workflows



# Sharing and Publishing

## Share or Publish History 'Mother child - sample'

### Make History Accessible via Link and Publish 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.

### Share History with Individual Users

You have not shared this history with any users.

[Share with a user](#)

## Published Histories

Q
[Advanced Search](#)

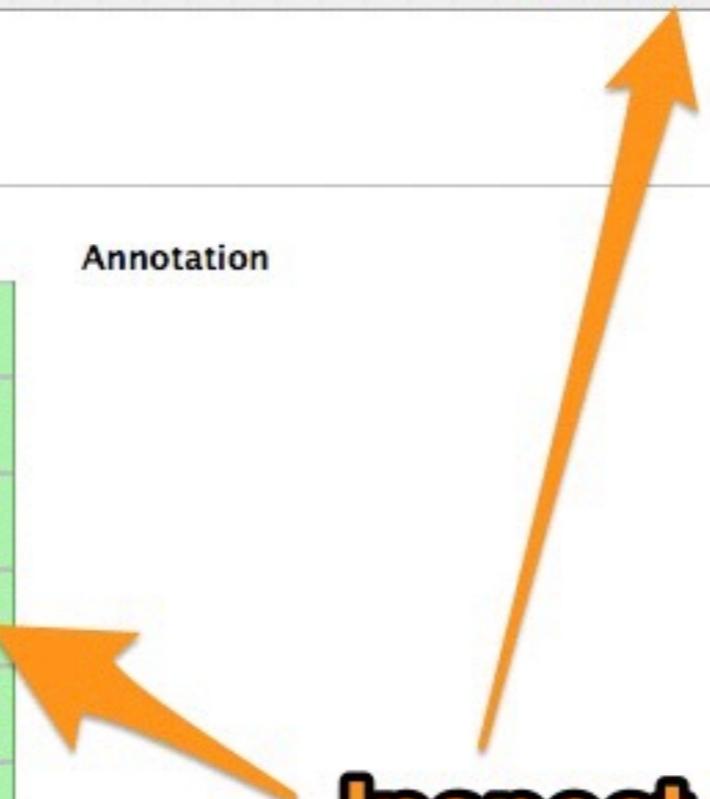
Name	Annotation	Owner	Community Rating	Community Tags	Last Updated
<a href="#">Infravec</a>		dan-lawson			Oct 01, 2014
<a href="#">SM_1186088</a>	Datasets correspond to our paper published in Science by Peleg et al. entitled : Altered histone acetylation is associated with age-dependent memory...	publicdata			Apr 19, 2010
<a href="#">MOL470_Pset3_All</a>		jbgreisman			Nov 18, 2012
<a href="#">SNP Calling</a>		jallen			Feb 06, 2014
<a href="#">RNA-seq shared data</a>		rna-seq-helin-group		illumina rnaseq	Jul 17, 2013
<a href="#">SRR534285</a>		jallen			Feb 04, 2014
<a href="#">ChIP-seq shared data</a>		chip-seq-helin-group		chip illumina	Jun 03, 2013
<a href="#">metagenomic analysis</a>		aun1		nsg megablast galaxy 454 metagenomics	Mar 19, 2010
<a href="#">Galaxy vs MEGAN</a>	Comparison of Galaxy vs. MEGAN pipeline.	aun1		megan galaxy metagenomics	Mar 19, 2010
<a href="#">Variant Analysis for Sample E18</a>	Perform a variant analysis with default parameters to identify variants in sample E18 that lie in annotated genes.			pileup bowtie demo.snp.sample	Sep 24, 2010
<a href="#">RNAseq_DGE_BASIC_Prep</a>		eric-the-red			Aug 26, 2013

## Galaxy History 'week10 RNAseq 241213'

## Dataset

- [1: normal 1.fastq](#)
- [2: normal 2.fastq](#)
- [3: cancerous 1.fastq](#)
- [4: cancerous 2.fastq](#)
- [5: Tophat for Illumina on data 2 and data 1: insertions](#)
- [6: Tophat for Illumina on data 2 and data 1: deletions](#)
- [7: Tophat for Illumina on data 2 and data 1: splice junctions](#)
- [8: Tophat for Illumina on data 2 and data 1: accepted hits](#)
- [9: Tophat for Illumina on data 4 and data 3: insertions](#)
- [10: Tophat for Illumina on data 4 and data 3: deletions](#)
- [11: Tophat for Illumina on data 4 and data 3: splice junctions](#)
- [12: Tophat for Illumina on data 4 and data 3: accepted hits](#)
- [13: Cufflinks on data 8: gene expression](#)
- [14: Cufflinks on data 8: transcript expression](#)
- [15: Cufflinks on data 8: assembled transcripts](#)
- [17: Cufflinks on data 12: gene expression](#)
- [18: Cufflinks on data 12: transcript expression](#)

## Annotation



**Inspect or  
import and  
use**

Author  
arik

## Related Histories

[All published histories](#)  
[Published histories by arik](#)

## Rating

Community  
(0 ratings, 0.0 average)



Yours



## Tags

Community: none

Yours:



[Published Pages](#) | [aun1](#) | Galaxy 101: The first thing you need to try

About this Page

## Author

aun1



## Related Pages

[All published pages](#)

[Published pages by aun1](#)

## Rating

Community

(46 ratings, 4.9 average)



Yours



## Tags

Community:

[tutorial](#)

[snps](#)

[exons](#)

Yours:



## Galaxy 101: The first thing you should try

In this very simple example we will introduce you to bare basics of Galaxy:

Getting data from UCSC

Performing simple data manipulation

Understanding Galaxy's History system

Creating and editing workflows

Applying workflows to your data

You can watch a step-by-step explanation of this entire tutorial [here](#).

## What are we trying to do?

Suppose you get the following question: "Mom (or Dad) ... Which coding exon has the highest number of single nucleotide polymorphisms on chromosome 22?". You think to yourself "Wow! This is a simple question ... I know exactly where the data is (at UCSC) but how do I actually compute this?" The truth is, there is really no straightforward way of answering this question in a time frame comparable to the attention span of a 7-year-old. Well ... actually there is and it is called Galaxy. So let's try it...

## 0. Organizing your windows and setting up Galaxy account

### 0.0. Getting your display sorted out

To get the most of this tutorial open two browser windows. One you already have (it is this page). To open the other, right click [this link](#) and choose "Open in a New Window" (or something similar depending on your operating system and browser):

[Open Link in New Window](#)

[Open Link in New Tab](#)

[Download Linked File](#)

[Download Linked File As...](#)

[Add Link to Bookmarks...](#)

[Copy Link](#)

# Visualization within Galaxy

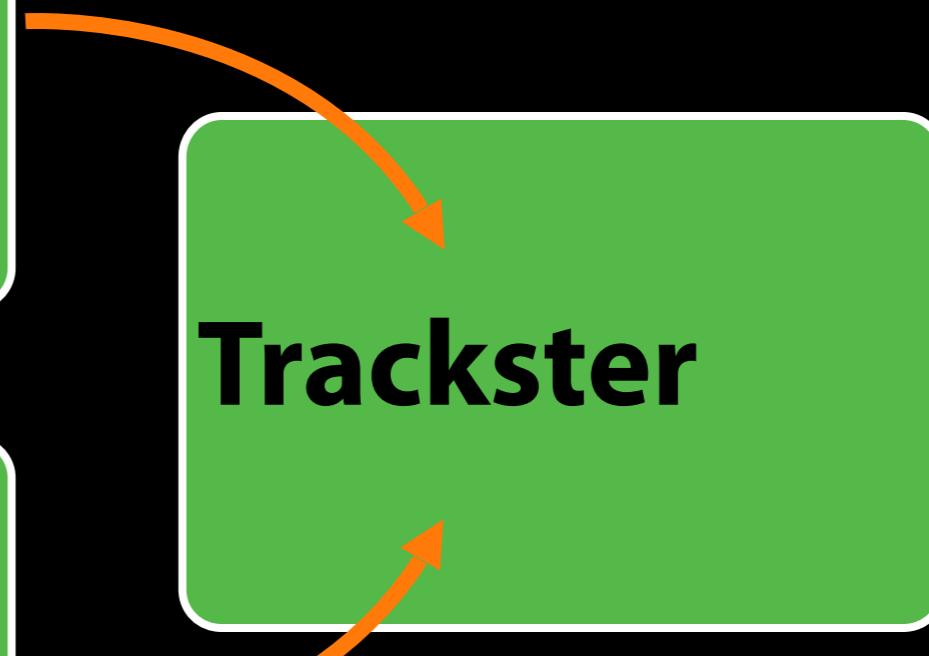
## Galaxy

- Tool integration framework
- Heavy focus on usability
- Sharing, publication framework

## Genome Browser

- Physical depiction of data
- Visually identify correlations
- Find interesting regions, features

## Trackster



[Analyze Data](#) [Workflow](#) [Shared Data](#) [Visualization](#) [Help](#) [User](#)[Published Visualizations](#) | [jeremy](#) | [GCC2011-1: Viewing and](#)

chr19

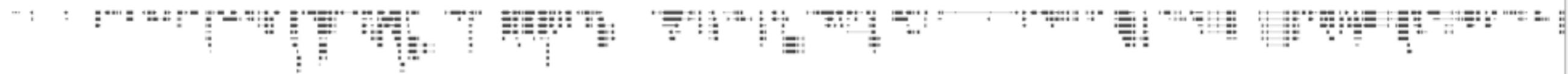
1,290 – 4,168,475



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

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

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

Display a menu

Galaxy

## Analyze Data

## Workflow

## Shared Data

## Visualization

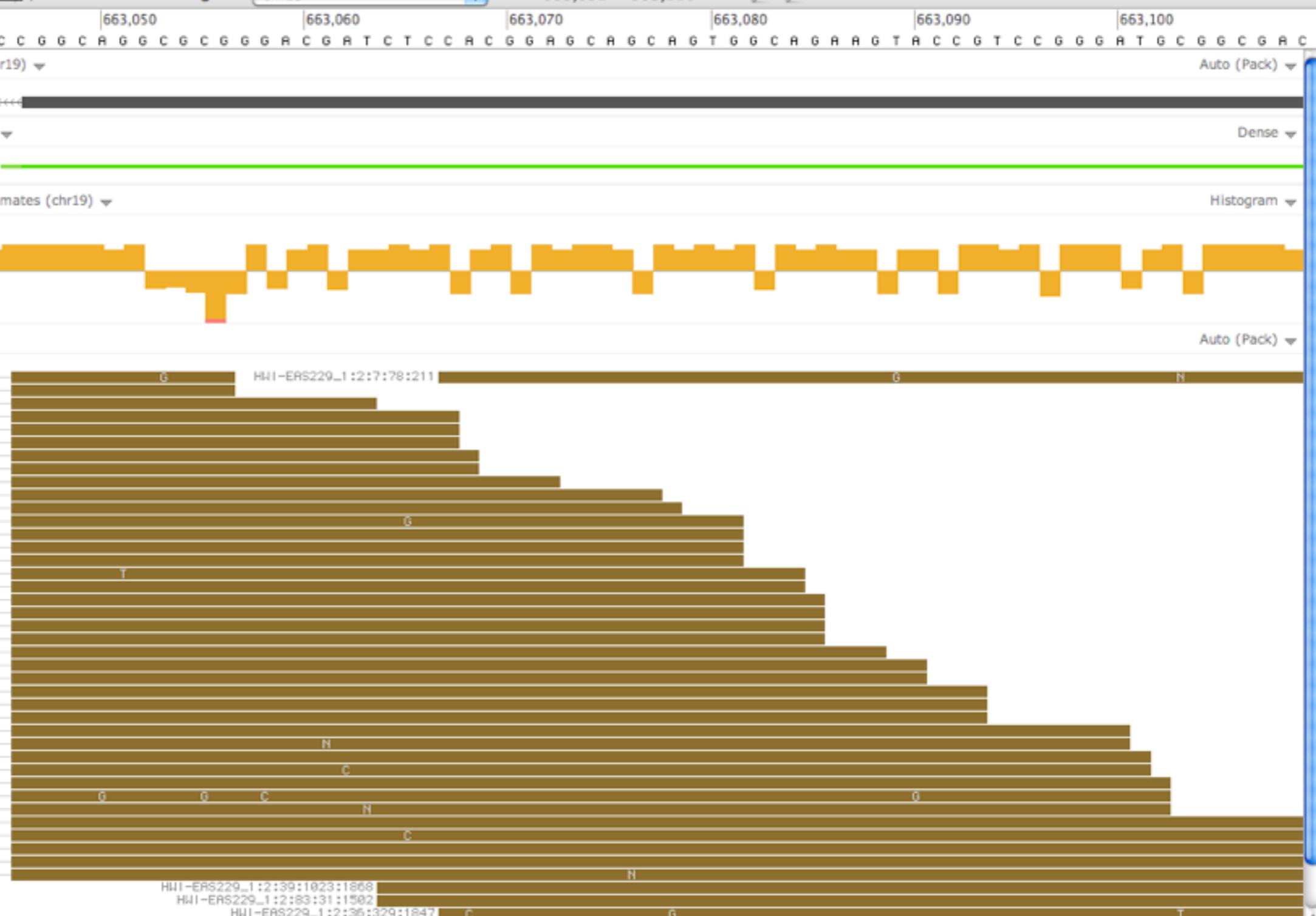
Help

User

Published Visualizations | jeremy | GCC2011-1: Viewing and chr19

chr19

663,032 – 663,110



h1-hESC Cufflinks assembled transcripts

30

# HTS Analysis

- Set parameters
- Run tools / workflow
- Wait...
- Visualize output



Iterate

# Experimentation

- Understand, debug, and “tune” analyses
- Easy with Galaxy
  - Parameter values saved
  - Can rerun tools and analyses
  - All outputs stored

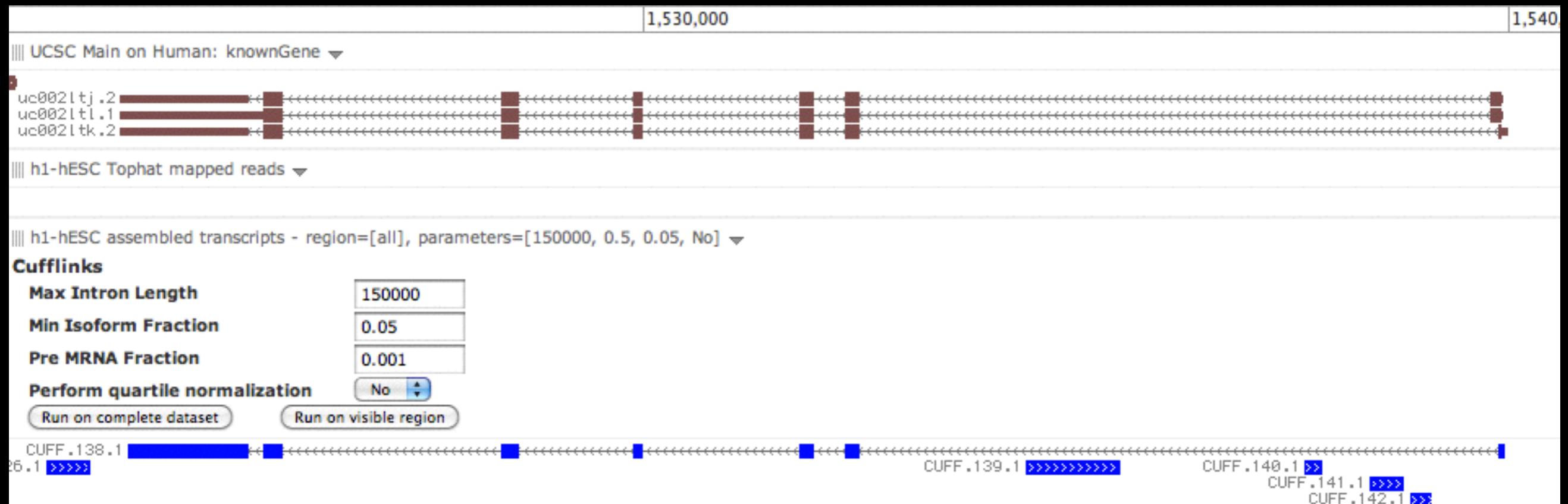
# Visualization Then and Now

- Then
  - Visualization endpoint of analysis
  - Difficult to change analyses based on visualization
- Now
  - Integrate analysis and visualization
  - Make it simple to move between analysis and visualization

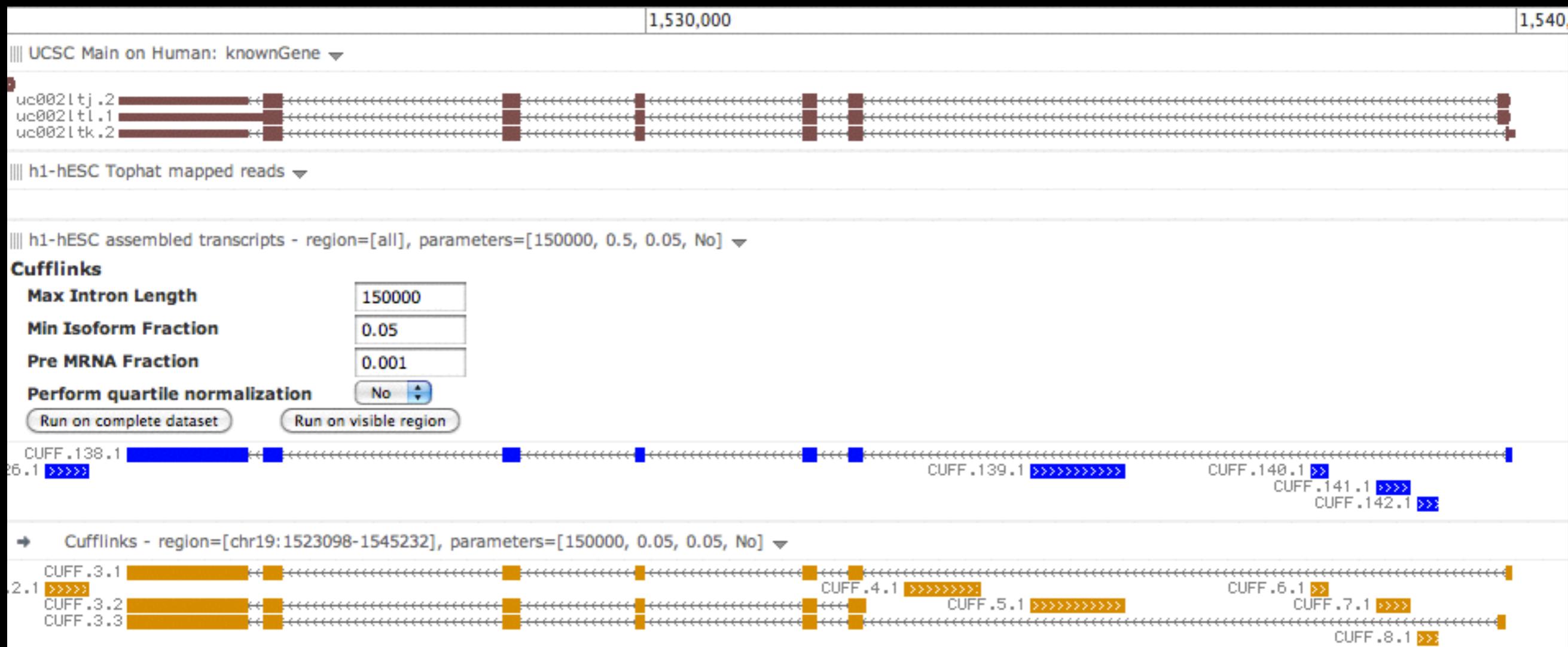
# Dynamic Filtering



# Integrating Tools and Visualization



# Integrating Tools and Visualization



# Three ways to use Galaxy

Public website

Download and Run Locally

Run on the Cloud

# <http://usegalaxy.org> (a.k.a. Main)

- Public web site

- Anybody can use it

- Hundreds of tools

- Persistent

- +500 users/month

- ~200TB of user data

- ~140,000 analysis jobs / month

- Dedicated resources
- Shared resources



# Public Galaxy Servers

<https://wiki.galaxyproject.org/PublicGalaxyServers>

Interested in:

60+ Public Servers

ChIP-chip and ChIP-seq?

✓ Cistrome

Statistical Analysis?

✓ Genomic Hyperbrowser

Sequence and tiling arrays?

✓ Oqtans

Text Mining?

✓ DBCLS Galaxy

Reasoning with ontologies?

✓ GO Galaxy

Internally symmetric protein structures?

✓ SymD

# Local Galaxy Instances

<http://getgalaxy.org>

Galaxy is designed for local installation and customization

- Easily integrate new tools
- Easy to deploy and manage on nearly any (Unix) system

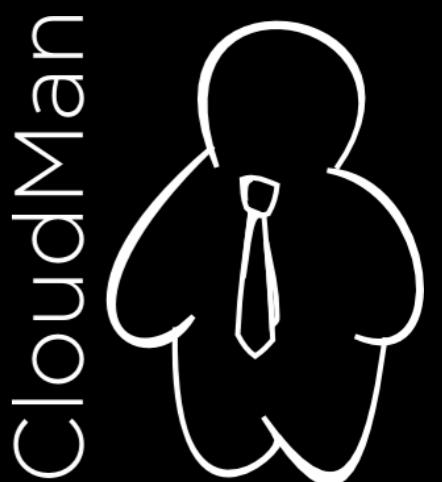
# Got your cluster?

- Move tool execution to other systems
- Galaxy works with any DRMAA compliant cluster job scheduler (which is most of them).
- Galaxy is just another client to your scheduler.



# Need a cluster? Galaxy CloudMan

- Start with a **fully configured and populated** (tools and data) Galaxy instance in the cloud.
- Allows you to scale up and down your compute assets as needed.
- Someone else manages the data center.
- **We are using this today**



# Manage Your Cloud Cluster

 **CloudMan from Galaxy**

[Admin](#) | [Report bugs](#) | [Wiki](#) | [Screencast](#)

## CloudMan Console

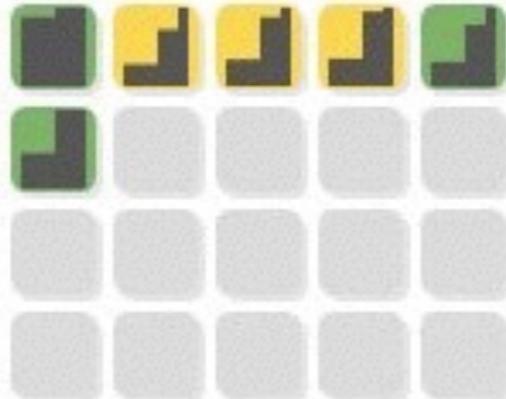


This page allows you to manage this instance cloud cluster and the services provided by it. Once the cluster has initialized, use the controls below to manage your cluster.

**Nodes ▾** **Remove nodes ▾** **Access Galaxy**

Requested: 5

Autoscaling is **off**. Turn on?





# Customizing Galaxy: Tool Shed

- Allow users to share tools, datatypes, workflows, sample data, and automated installation scripts for tool dependencies: **biomedical app store**
- Integration with Galaxy instances to automate tool installation and updates

<http://toolshed.g2.bx.psu.edu>

2857 valid tools on Dec 16, 2014

**Search**

- [Search for valid tools](#)
- [Search for workflows](#)

**Valid Galaxy Utilities**

- [Tools](#)
- [Custom datatypes](#)
- [Repository dependency definitions](#)
- [Tool dependency definitions](#)

**All Repositories**

- [Browse by category](#)

**Available Actions**

- [Login to create a repository](#)

## Repositories by Category

Name	Description	Repositories
<a href="#">Assembly</a>	Tools for working with assemblies	31
<a href="#">ChIP-seq</a>	Tools for analyzing and manipulating ChIP-seq data.	9
<a href="#">Combinatorial Selections</a>	Tools for combinatorial selection	
<a href="#">Computational chemistry</a>	Tools for use in computational chemistry	17
<a href="#">Convert Formats</a>	Tools for converting data formats	34
<a href="#">Data Managers</a>	Utilities for Managing Galaxy's built-in data cache	6
<a href="#">Data Source</a>	Tools for retrieving data from external data sources	18
<a href="#">Fasta Manipulation</a>	Tools for manipulating fasta data	63
<a href="#">Fastq Manipulation</a>	Tools for manipulating fastq data	27
<a href="#">Genome-Wide Association Study</a>	Utilities to support Genome-wide association studies	1
<a href="#">Genomic Interval Operations</a>	Tools for operating on genomic intervals	33
<a href="#">Graphics</a>	Tools producing images	28
<a href="#">Imaging</a>	Utilities to support imaging	1
<a href="#">Metabolomics</a>	Tools for use in the study of Metabolomics	4
<a href="#">Metagenomics</a>	Tools enabling the study of metagenomes	41
<a href="#">Micro-array Analysis</a>	Tools for performing micro-array analysis	7
<a href="#">Next Gen Mappers</a>	Tools for the analysis and handling of Next Gen sequencing data	75
<a href="#">Ontology Manipulation</a>	Tools for manipulating ontologies	8

[Install to Galaxy](#)[Browse repository](#)[Tool Shed Actions](#)

Repository 'bowtie2'

**Revision:**

0:96d2e31a3938 (2014-01-27)

## Dependencies of this repository

**▼ Repository dependencies – installation of these additional repositories is required**[Repository package bowtie2 2 1 0 revision 017a00c265f1 owned by devteam](#)[Repository package samtools 0 1 18 revision 171cd8bc208d owned by devteam](#)**▼ Tool dependencies – repository tools require handling of these dependencies**

Name	Version	Type
samtools	0.1.18	package
bowtie2	2.1.0	package

## Contents of this repository

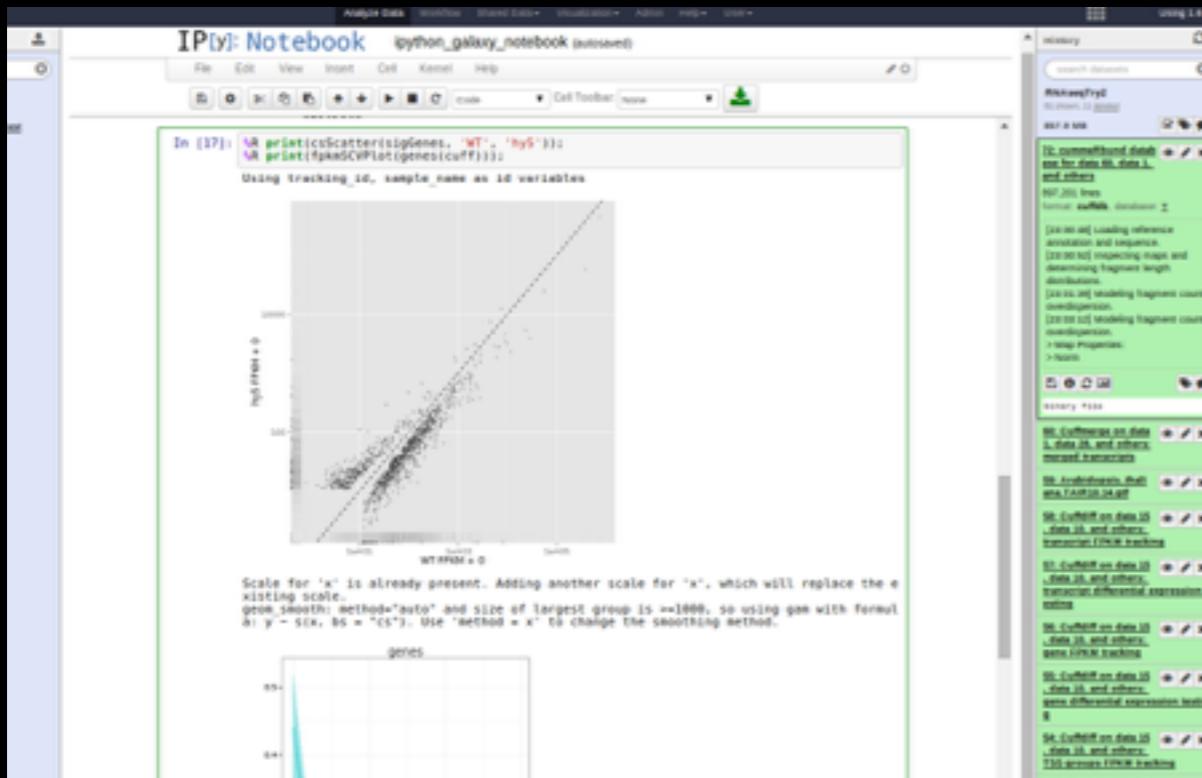
**► Valid tools – click the name to preview the tool and use the pop-up menu to inspect all metadata**

## Automated tool test results

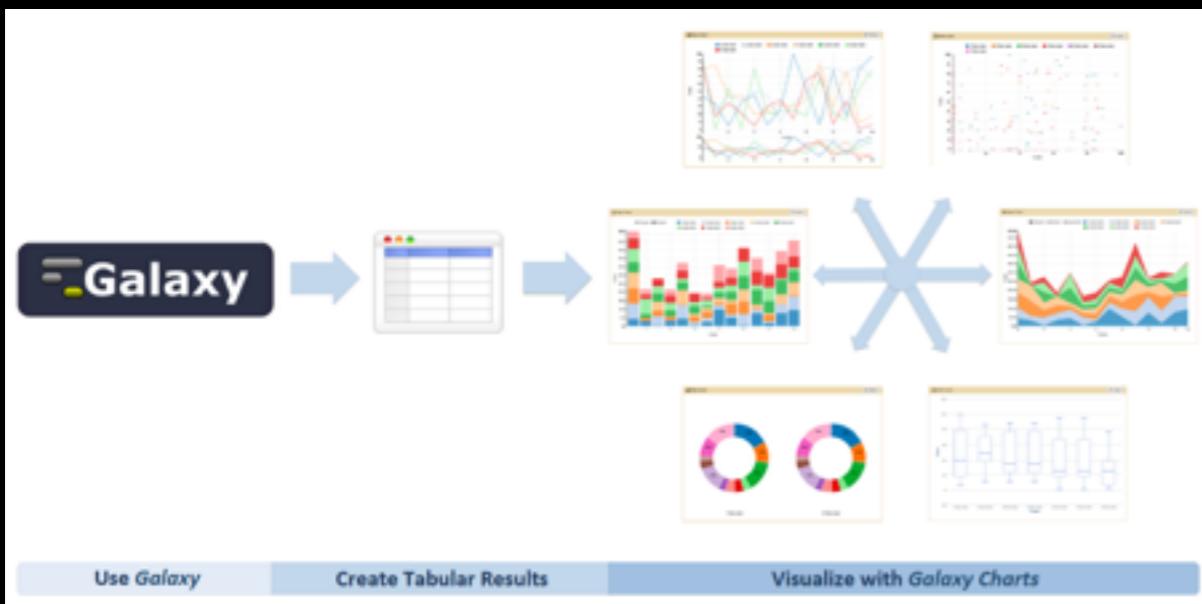
**► Test runs**

# And much, much more...

## Interactive environments



## Embedded charts



## Dataset collections

The figure shows a screenshot of a dataset collection interface. On the left, there are two search boxes labeled 'search datasets' and 'Isolates'. The main area displays a list of datasets:

Index	Dataset Name	Description
8	isolate-1141_2.fastq	
7	isolate-1141_1.fastq	
6	isolate-1140_2.fastq	
5	isolate-1140_1.fastq	
4	isolate-1139_2.fastq	
3	isolate-1139_1.fastq	
2	isolate-1138_2.fastq	
1	isolate-1138_1.fastq	

To the right, there is a detailed view of a dataset named 'isolates-expl' which is described as 'a list of paired datasets'. It contains entries for 'forward' and 'reverse' datasets, each associated with 'isolate-1141', 'isolate-1140', 'isolate-1139', and 'isolate-1138' respectively.

## History structure

The figure shows a screenshot of a history structure interface. At the top, there is an 'Upload File' section with a 'filter' dropdown set to 'vcf files'. Below it, a file named '3.maf' is listed, showing details such as '9 blocks', 'format: maf', 'database: ?', and 'uploaded maf file'. A preview of the MAF file content is shown, including genomic coordinates and variants. To the right, there is a list of processing steps:

- 'Convert.MAF.to.Genomic.Intervals'
- 'Filter' (with a condition '1.vcf Qual>60')
- 'MAF.to.Interval'
- 'Concatenate.datasets' (with a condition 'left-to-right')
- 'Concatenate.datasets.on.data.10.and.data.7'

# The Agenda

Galaxy project overview

RNA-Seq exercise

RNA-Seq in brief

QC

Mapping

Galaxy Community

Differential expression

# Not The Agenda

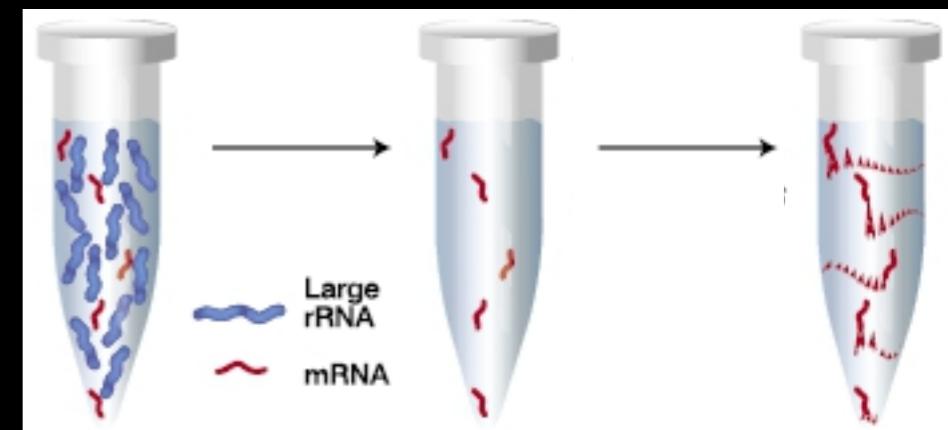
This workshop will *not* cover

- details of how tools are implemented, or
- new algorithm designs, or
- which assembler or mapper or peak caller or ... is best for you.

While this workshop does cover RNA-Seq, we are only using that specific example to learn general principles.

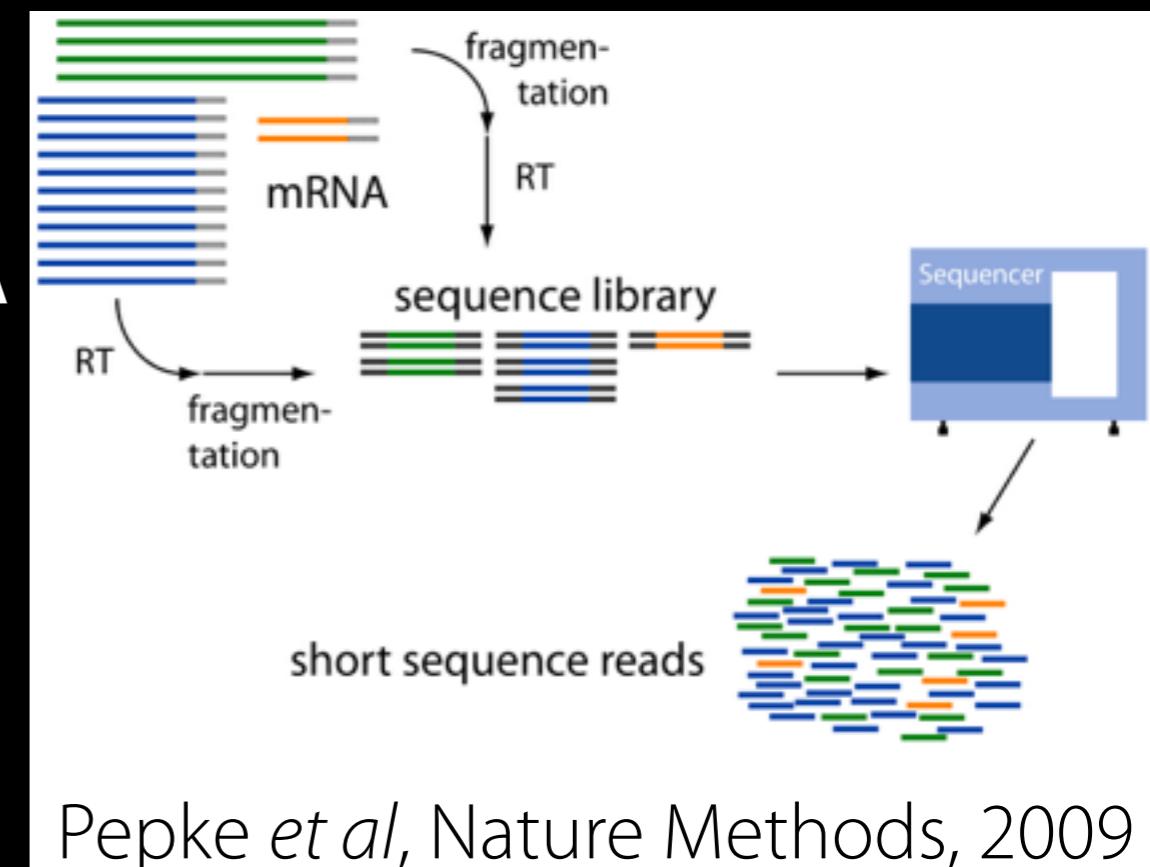
# RNA

- RNA in cells consists of
  - 95% ribosomal rRNA and tRNA
  - other non-coding ncRNA
  - protein coding mRNA
- Sequence is transcribed from genome but
  - Introns spliced out
  - mRNA is polyadenylated (“A”s added to end)



# RNA-Seq

- Deplete rRNA
- or select for polyadenylated RNA
- Fragmentation
- Reverse transcribe to cDNA
- Attach adaptor sequences
- Size selection
- Amplify by PCR
- High-throughput sequencing (eg Illumina)



Pepke *et al*, Nature Methods, 2009

# RNA-Seq output

- Several million “**reads**” per sample
- Reads are RNA sequence starting from random locations within the original mRNA
- May read through into adaptor sequence
- Current typical length ~150 bases  
(really just needs to be long enough to locate in genome)

# RNA-Seq applications

- Sequencing RNA transcripts, applications:
  - *de novo* transcriptome
  - **gene-wise differential expression**
  - novel splice variants
  - differential splicing
  - fusion genes
  - non-coding RNA
  - polyadenylation length
  - post-transcriptional modification
  - ...
- Differential expression analysis most common
  - Commonly called “DGE”

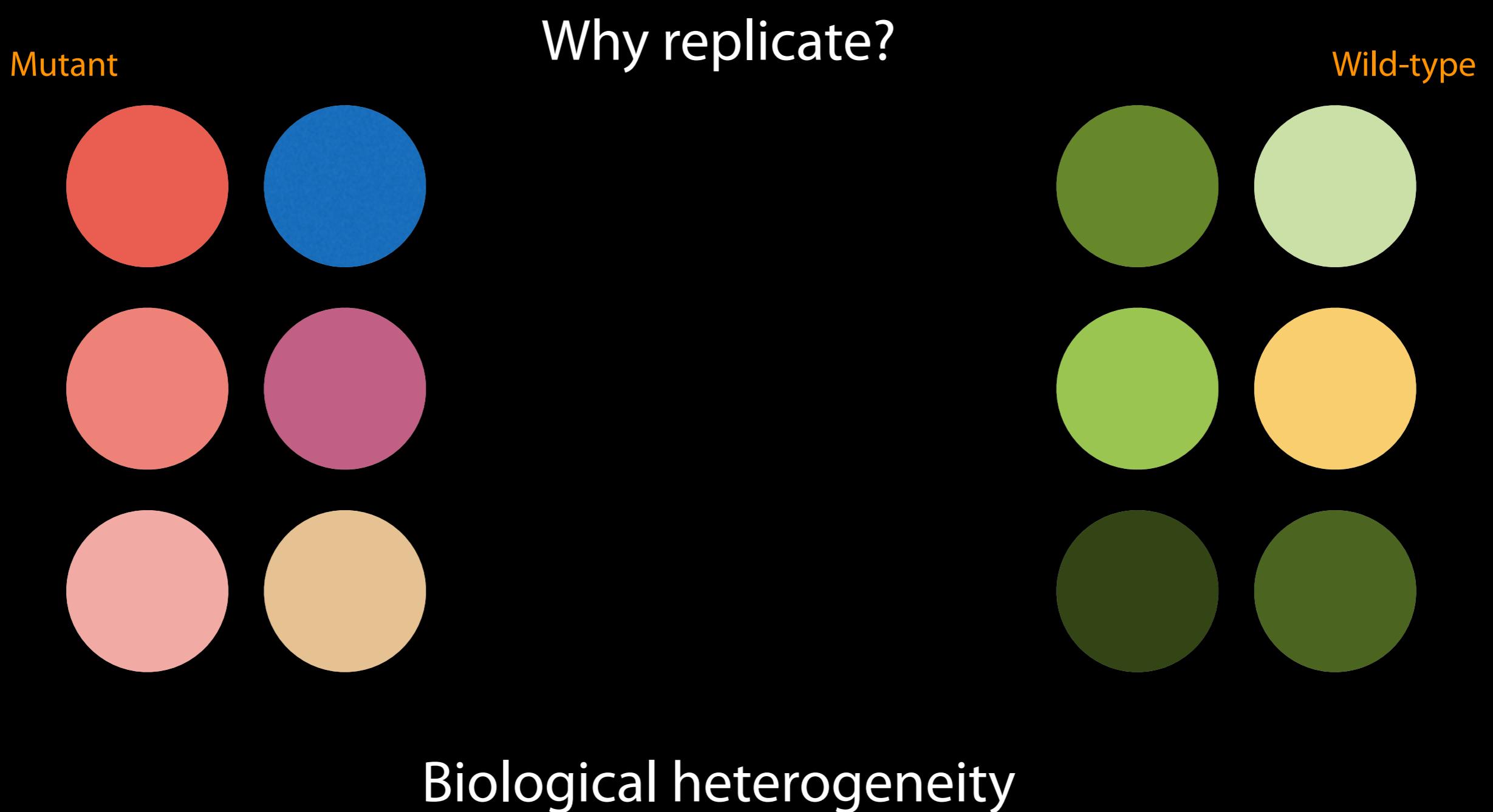
# This workshop

- Analysis of RNA-seq data to detect differential gene expression
- How much are the transcript expression levels changing between conditions
- We'll do a basic **A vs B** analysis

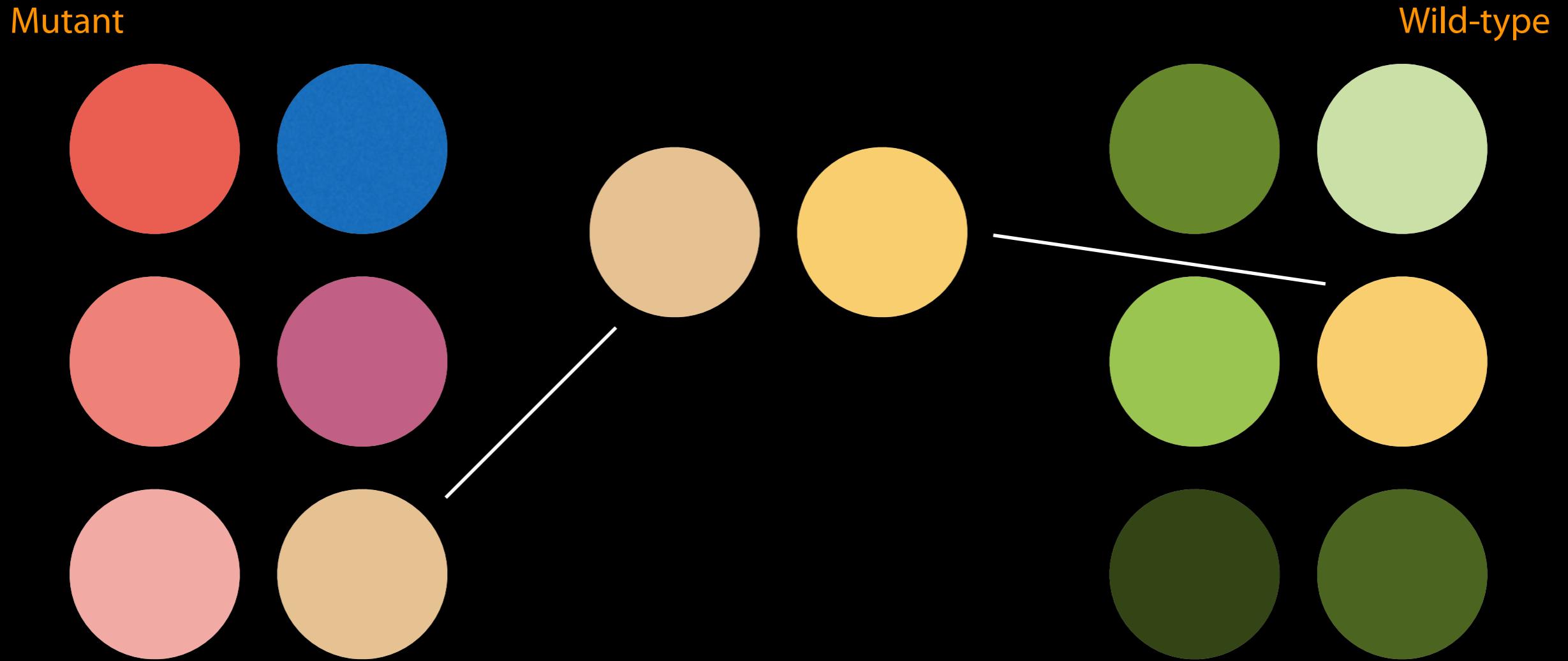
# Example experiment

- 2 conditions of interest
  - eg. mutant versus wild-type, or
  - treated versus untreated, or, ...
- Want to find all genes that are expressed differently between the 2 conditions

# Replicas, replicas, replicas

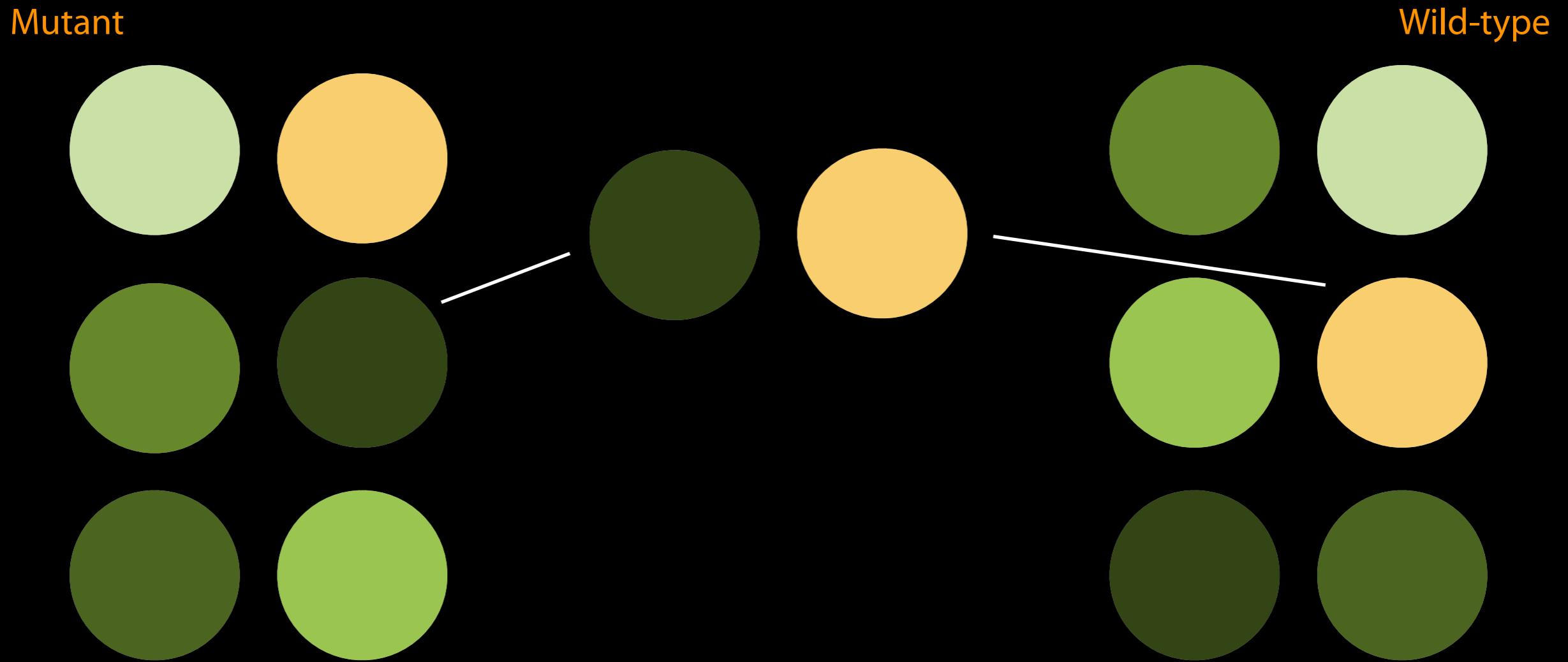


# Unreplicated design



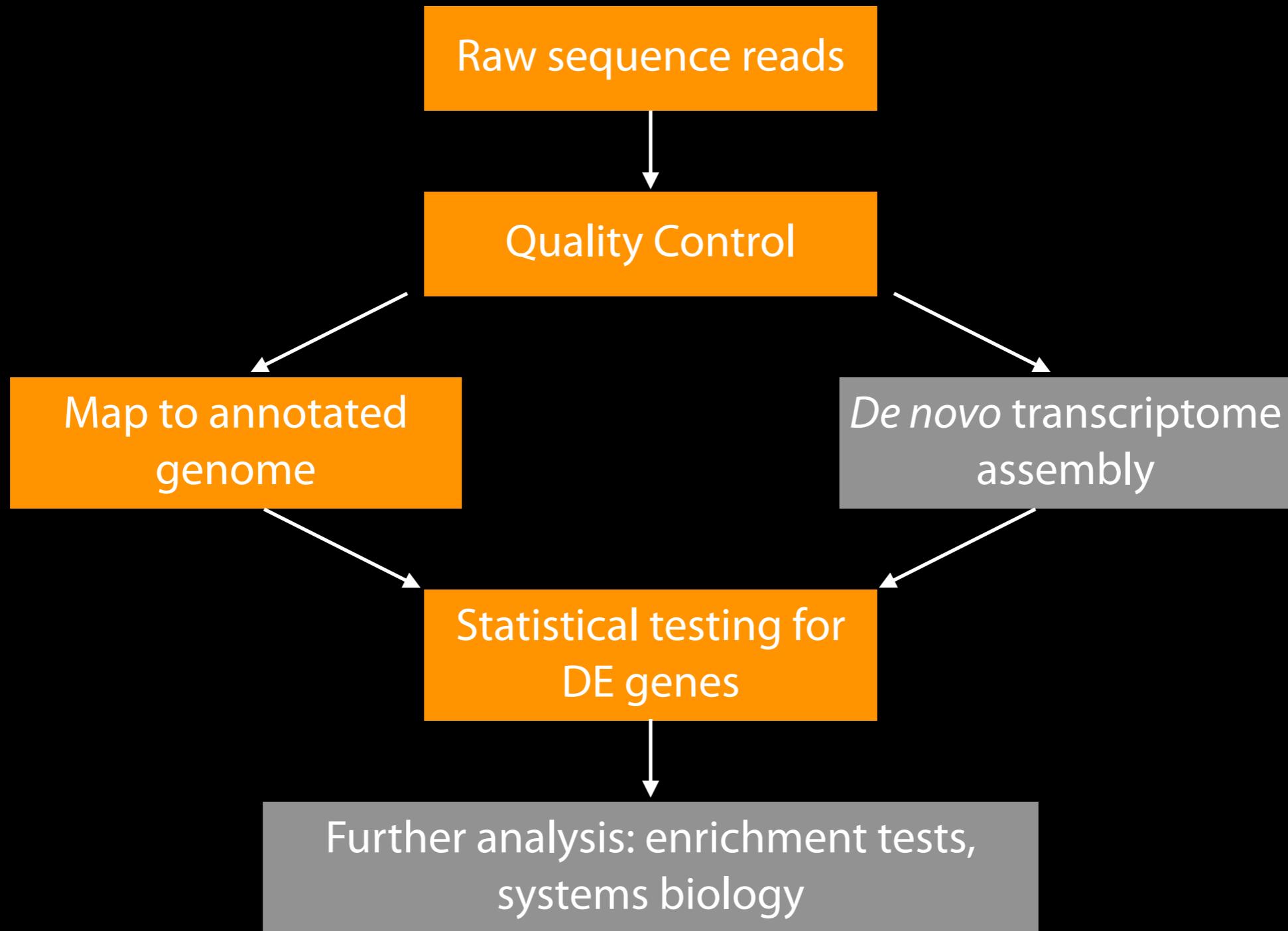
Here, groups differ, but single replicates from each group very similar

# Unreplicated design



Here, groups are similar, but outlying observation from group on right makes it look like there's a big difference in unreplicated experiment

# RNA-Seq DE analysis steps



# NGS Data Quality Control

- FASTQ format
- Examine quality in an RNA-Seq dataset
- Trim/filter as we see fit, hopefully without breaking anything.

Quality Control is not sexy.

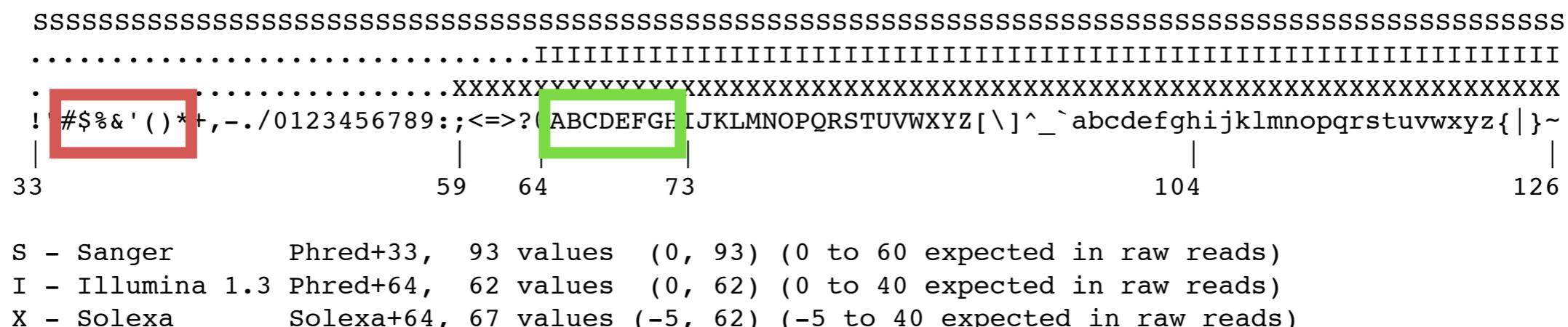
It is vital.

# What is FASTQ?

- Specifies sequence (FASTA) and quality scores (PHRED)
- Text format, 4 lines per entry

```
@SEQ_ID
GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTGTTCAACTCACAGTTT
+
! ' ' * ( ( ( (**+) ) %%++ ) ( %%% ) . 1***-+*' ) ) **55CCF>>>>CCCCCCCC65
```

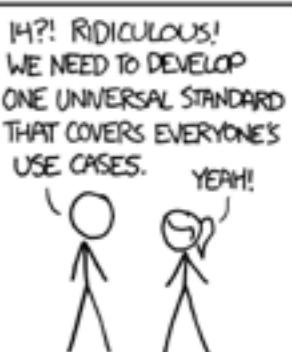
- FASTQ is such a cool standard, there are 3 (or 5) of them!



[http://en.wikipedia.org/wiki/FASTQ\\_format](http://en.wikipedia.org/wiki/FASTQ_format)

HOW STANDARDS PROLIFERATE:  
(SEE A/C CHARACTERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC.)

SITUATION:  
THERE ARE  
14 COMPETING  
STANDARDS.



SOON:  
SITUATION:  
THERE ARE  
15 COMPETING  
STANDARDS.

# Our infrastructure for the day

**<http://cloud1.galaxyproject.org/>**

**<http://cloud2.galaxyproject.org/>**

**<http://cloud3.galaxyproject.org/>**

**<http://cloud4.galaxyproject.org/>**



# NGS Data Quality Exercise

Create new history



(cog) → Create New

Get some data

Shared Data → Data Libraries

→ RNA-Seq Example Data\*

→ Unfiltered Reads

→ Select MeOH\_REPO1\_R1, MeOH\_REPO1\_R2

and then Import to current history



**UCDAVIS** Bioinformatics Core

\* RNA-Seq example datasets from the 2013 UC Davis Bioinformatics Short Course. <http://bit.ly/ucdbsc2013>

# NGS Data Quality: Assessment tools

NGS QC and Manipulation → **FastQC**

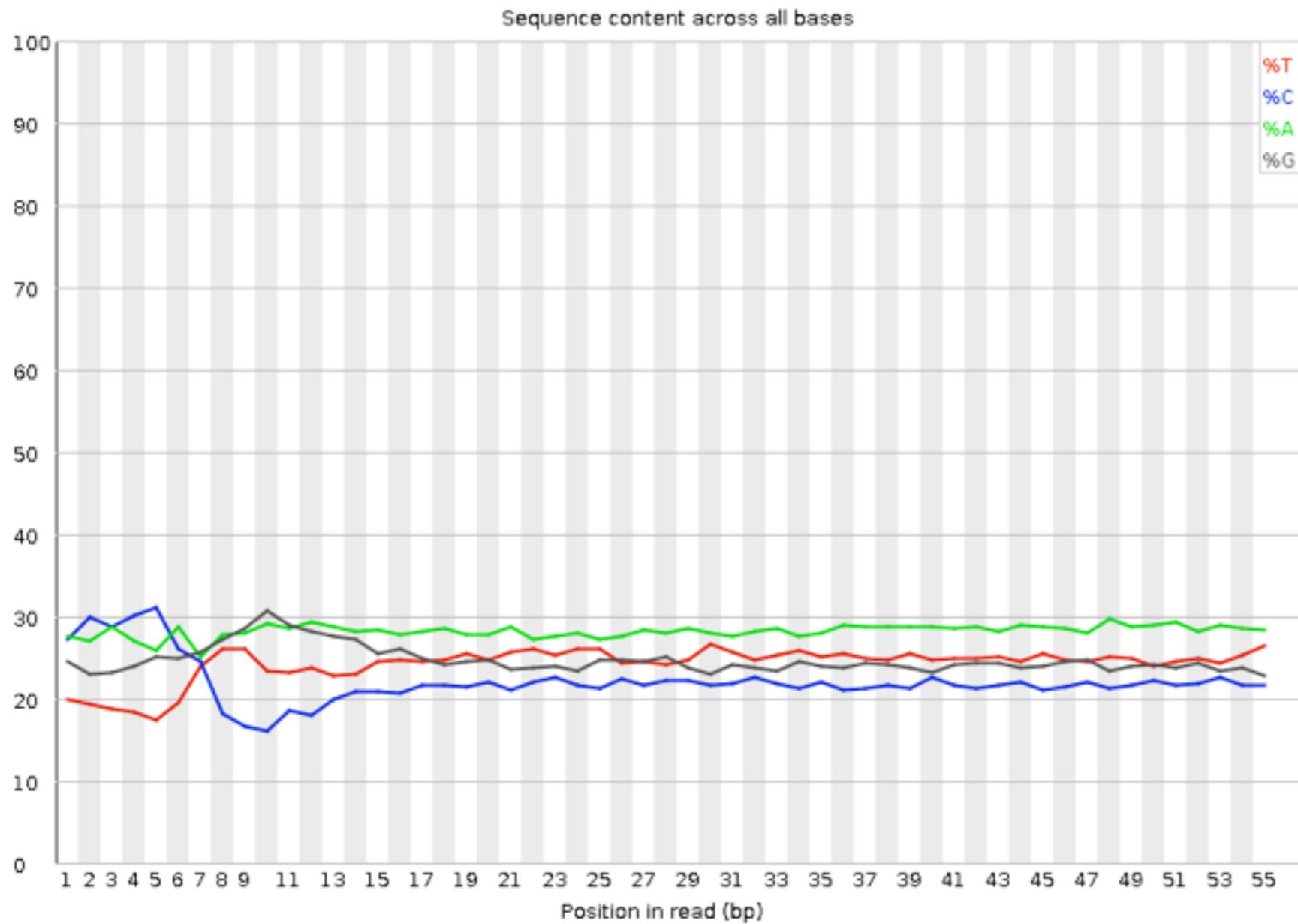
- Gives you a lot a lot of information but little control over how it is calculated or presented.

<http://bit.ly/FastQCBoxPlot>

# NGS Data Quality: Sequence bias at front of reads?



## Per base sequence content

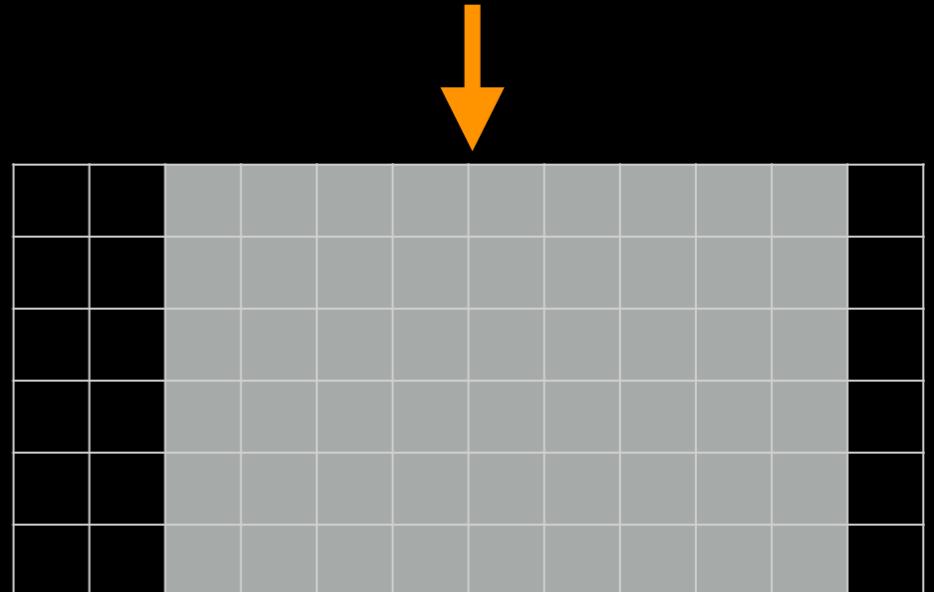
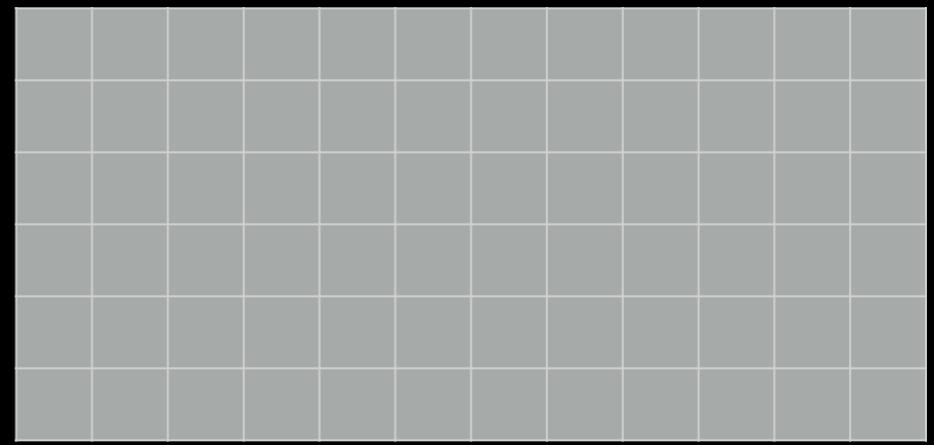


From a sequence specific bias that is caused by use of random hexamers in library preparation.

Hansen, *et al.*, "Biases in Illumina transcriptome sequencing caused by random hexamer priming" *Nucleic Acids Research*, Volume 38, Issue 12 (2010)

# NGS Data Quality: Trim as we see fit

- Trim as we see fit: Option 1
- NGS QC and Manipulation →  
**FASTQ Trimmer by column**
- Trim same number of columns  
from every record
- Can specify different trim for 5'  
and 3' ends



# NGS Data Quality: Base Quality Trimming

- Trim Filter as we see fit: Option 2

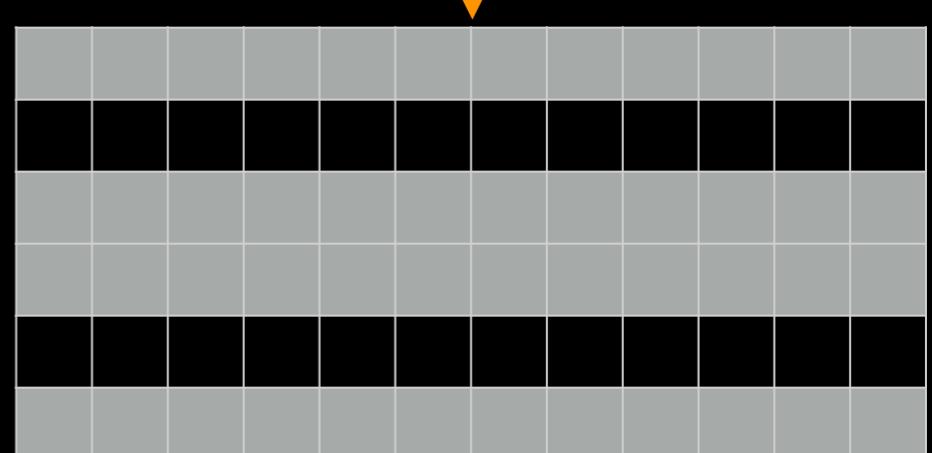
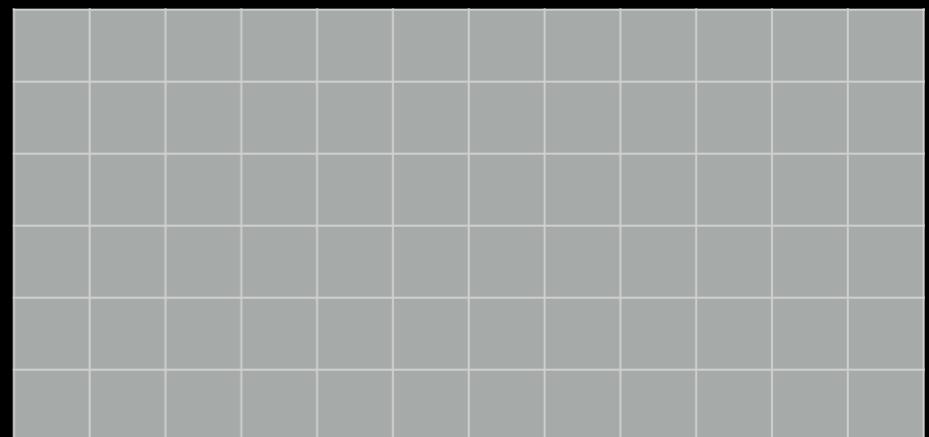
- NGS QC and Manipulation →

**Filter FASTQ reads by quality score and length**

- Keep or discard whole reads

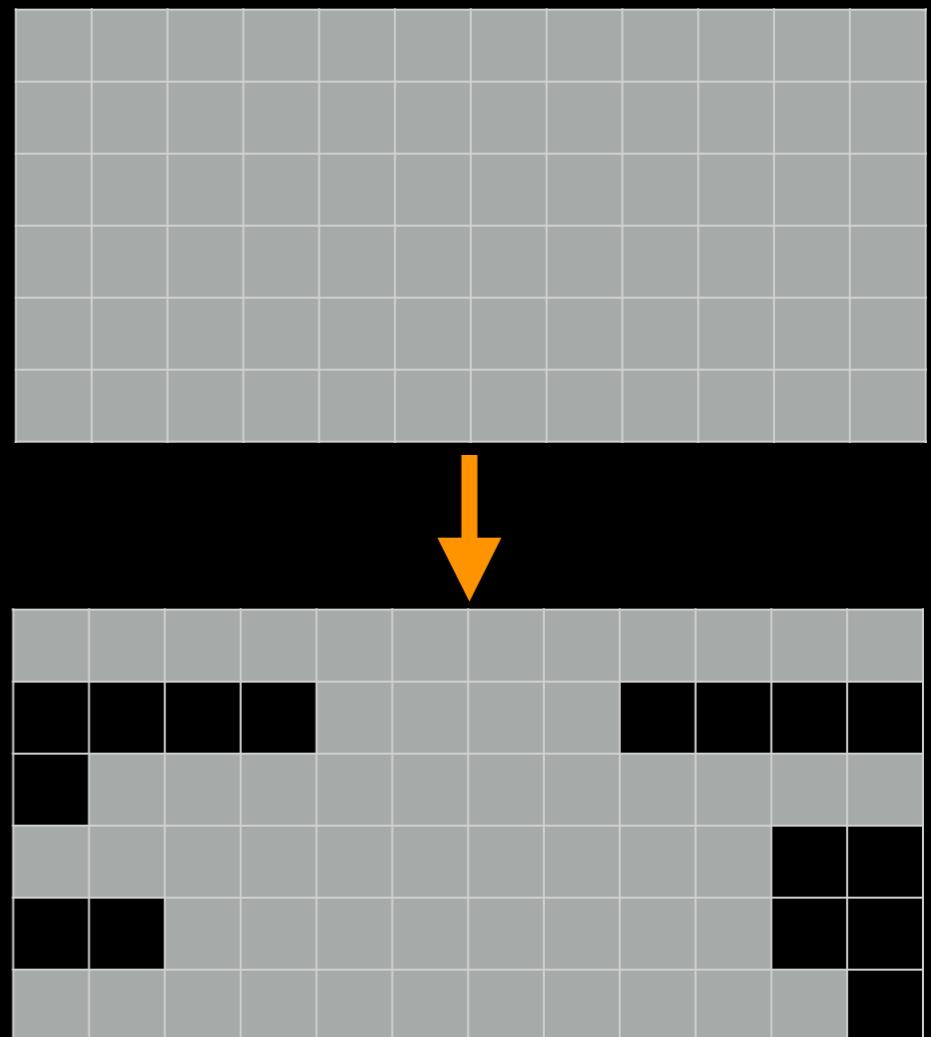
- Can have different thresholds for different regions of the reads.

- Keeps original read length.

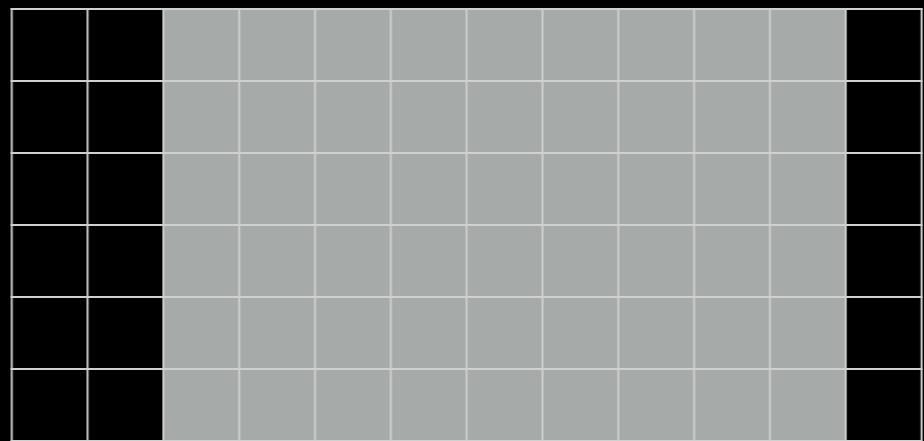
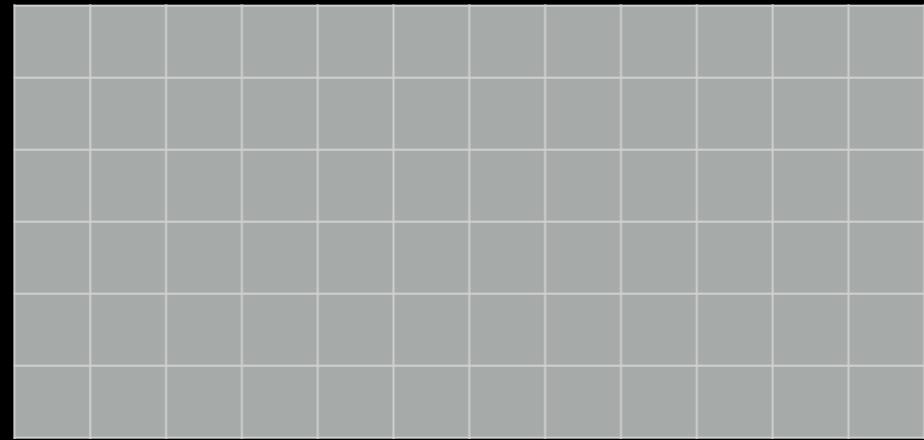


# NGS Data Quality: Base Quality Trimming

- Trim as we see fit: Option 3
- NGS QC and Manipulation →  
**FASTQ Quality Trimmer by  
sliding window**
- Trim from both ends, using  
sliding windows, until you hit a  
high-quality section.
- Produces variable length reads

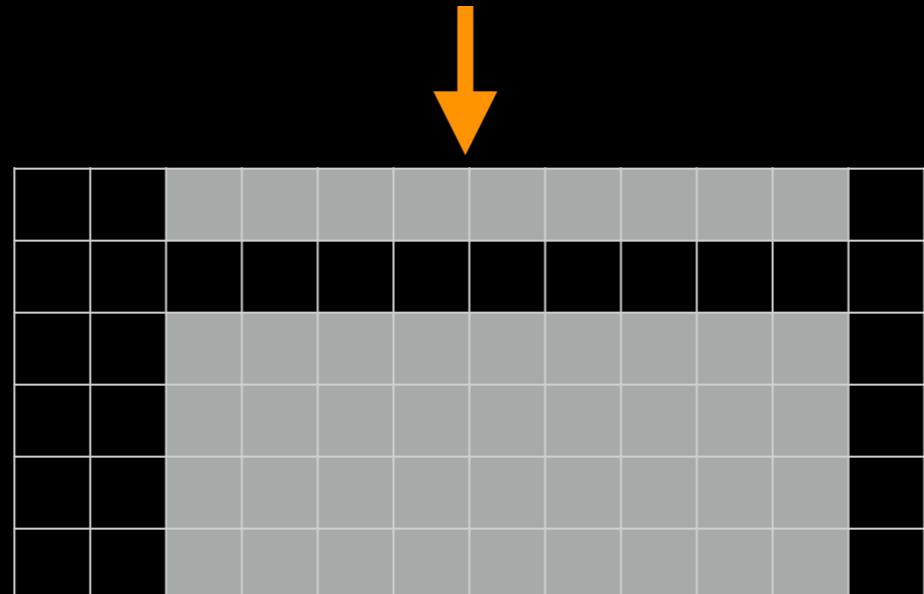


**Options are  
not mutually  
exclusive**



Option 1  
(by column)

+



Option 2  
(by entire row)

# Trim? As we see fit?

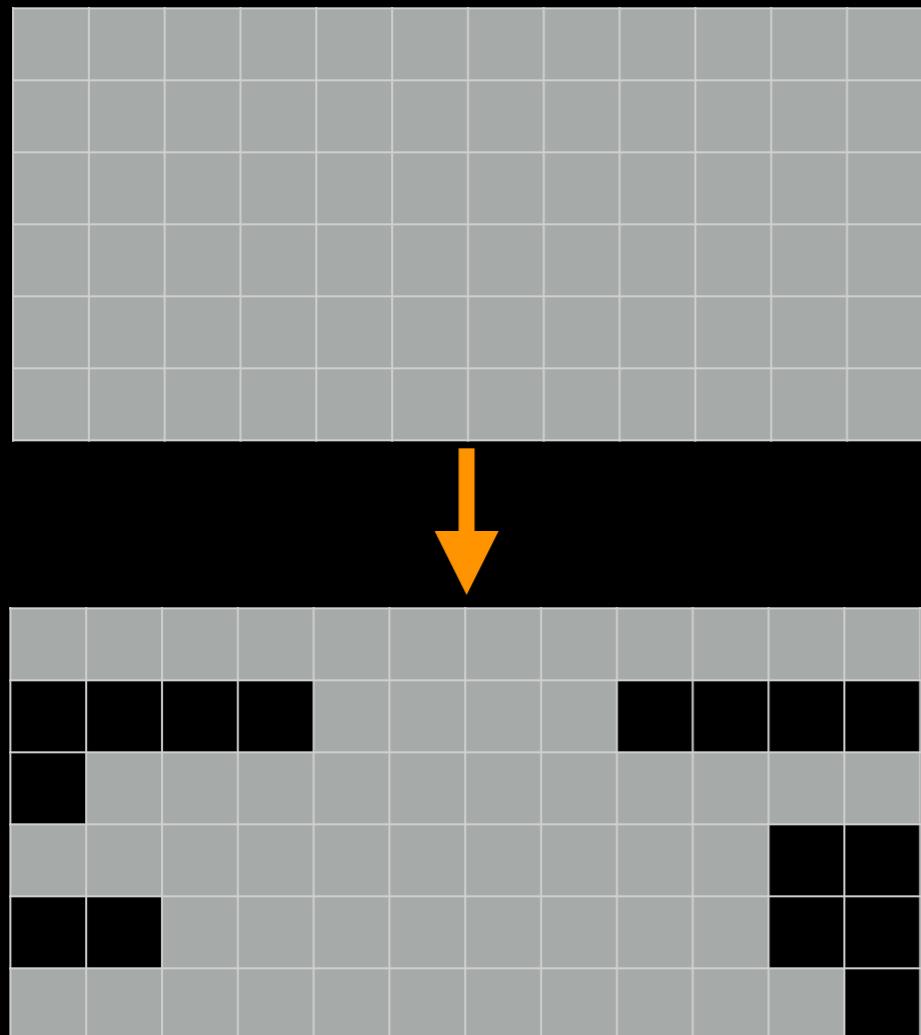
- Introduced 3 options
  - One preserves original read length, two don't
  - One preserves number of reads, two don't
  - Two keep/make every read the same length, one does not
  - One preserves pairings, two don't

# Trim? As we see fit?

- Choice depends on downstream tools
- Find out assumptions & requirements for downstream tools and make appropriate choice(s) now.
- How to do that?
  - Read the tool documentation
    - <http://biostars.org/>
    - <http://seqanswers.com/>
    - <http://galaxyproject.org/search>



# NGS Data Quality: Base Quality Trimming



I really want to use Option 3:

- NGS QC and Manipulation →  
**FASTQ Quality Trimmer by  
sliding window**

but ...

“Mixing paired- and single- end reads together is  
not supported.”

Tophat Manual

“If you are performing RNA-seq analysis, there is no  
need to filter the data to ensure exact pairs before  
running Tophat.”

Jen Jackson

Galaxy User Support Person Extraordinaire

“Dang.”

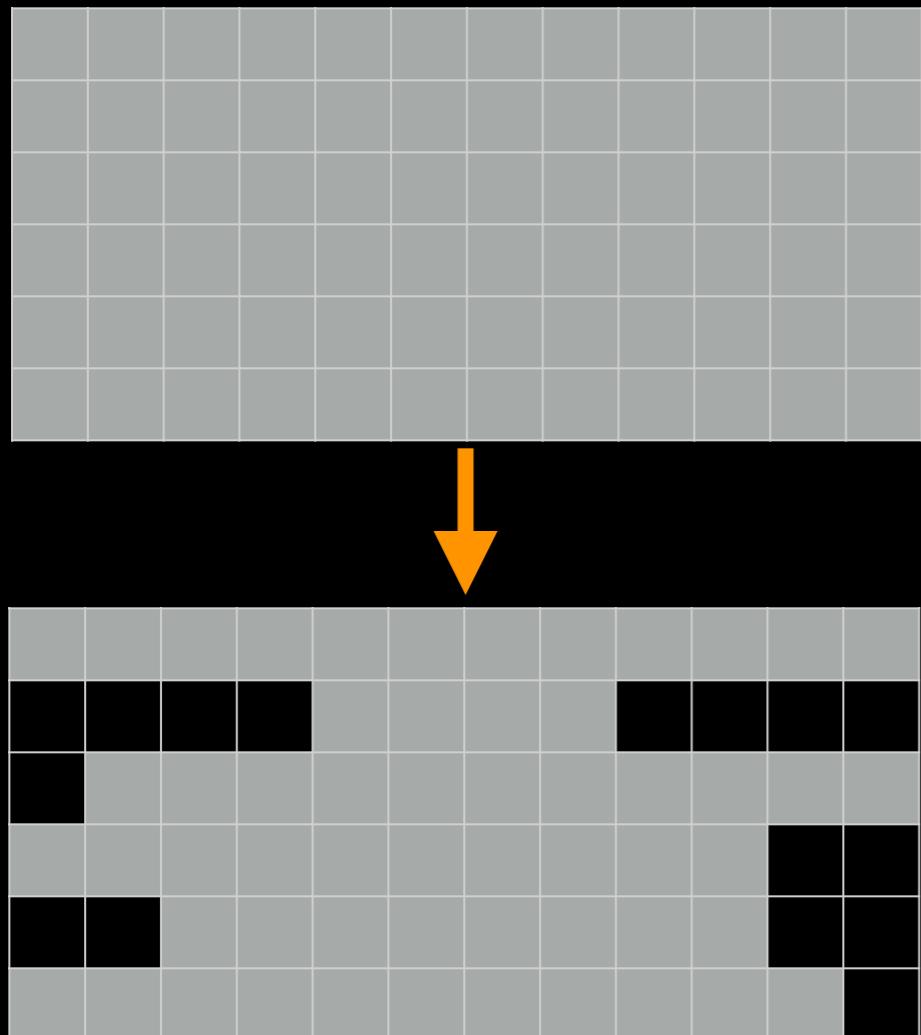
Most of us

Running Tophat on *no-longer-cleanly-paired* data *does map the reads*, but, it no longer keeps track of read pairs in the SAM/BAM file.

# Keeping paired ends paired: Options

- Don't bother.
- Run a workflow that removes any unpaired reads before mapping.
- Run the Picard **Paired Read Mate Fixer** after mapping reads.
- Use sliding windows for QC, but keep empty reads.

# NGS Data Quality: Base Quality Trimming



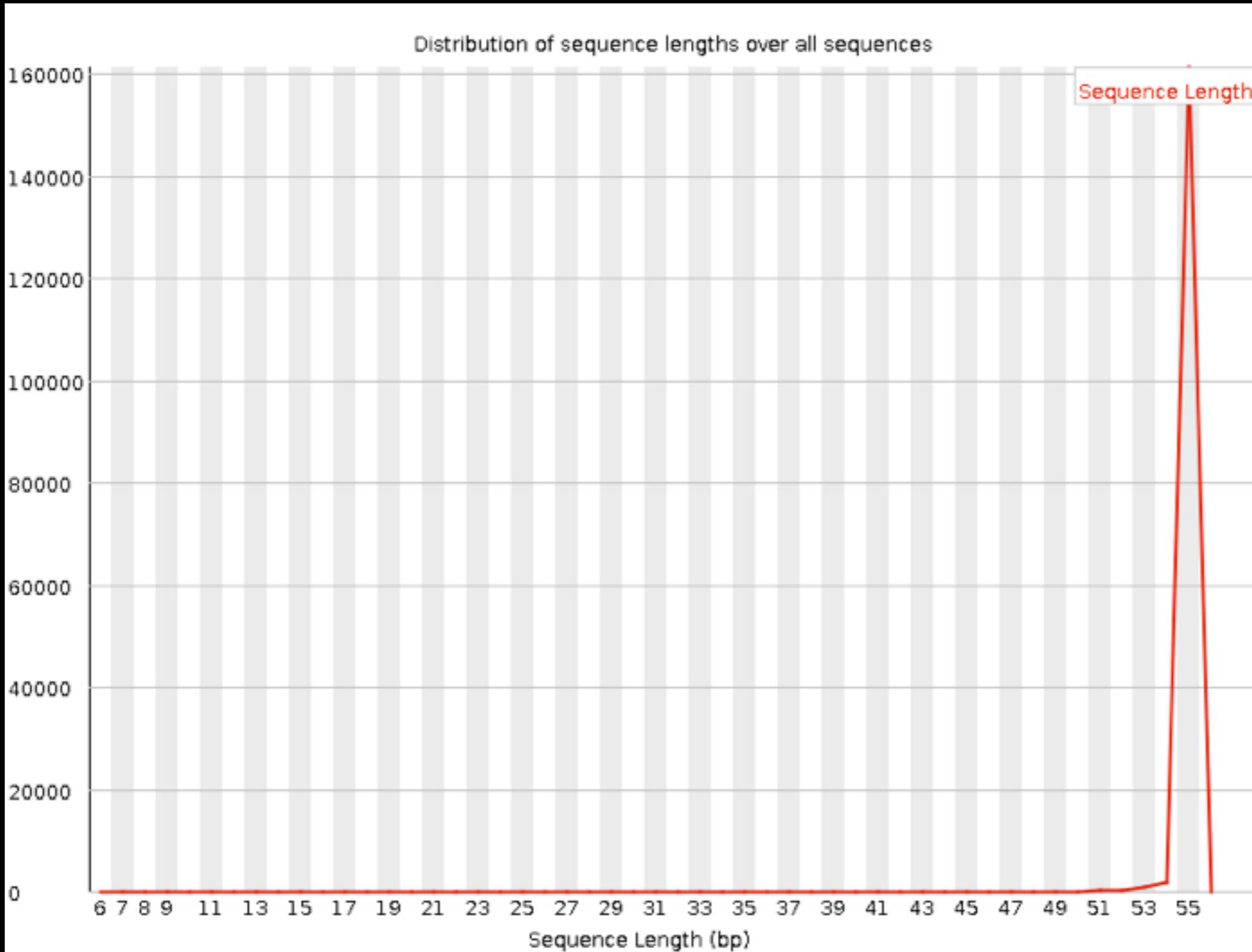
I'll use Option 3 (*but ...*):

- NGS QC and Manipulation → **FASTQ Quality Trimmer by sliding window**  
*Check "Keep reads with zero length"*

Run again:

- NGS QC and Manipulation → **FastQC**  
on trimmed dataset

# NGS Data Quality: Base Quality Trimming



New Problem?

Now some reads are so short they are just noise and can't be meaningfully mapped

Option 2 can fix this (but break pairings).

Or, your mapper may have an option to ignore shorter reads

# NGS Data Quality: Sequencing Artifacts

Repeat this process with MeOH Rep1 R2 (the reverse reads)

... and there's a problem in Overrepresented sequences:

 **Overrepresented sequences**

Sequence	Count	Percentage	Possible Source
CTGTGTATTTGTCAATTTCTTCTCCACGTTCTCTCGGCCCTGTTCCGTAGCCT	590	0.3541692929220167	No Hit
TT	342	0.2052981325073385	No Hit
CGGCCACAAATAAACACAGAAATAGTCAGAATGTCACAGGTCCAGGGCAGAGGA	325	0.19509325457568719	No Hit
CTGCATTATAAAAAGGACAGCCAGATATCAACTGTTACAGAAATGAAATAAGACG	230	0.13806599554587093	No Hit
CGGCCGCAAATAAACACAGAAATAGTCAGAATGTCACAGGTCCAGGGCAGAGGA	199	0.11945710049403614	No Hit
GTCAGCTCAACTTGTAGGCCCAAAAGAAAACAGCGTCTTACTGGGGAGGGATAT	197	0.11825652661972422	No Hit

NGS QC and Manipulation → Remove sequencing artifacts

But this will break pairings.

# NGS Data Quality: Done with 1st Replicate!

**Now, only 3 (or 5) more to go!**

## Workflows:

Create a QC workflow that does all these steps

Or, cheat and import the shared workflow.

Or, really cheat and just import the already trimmed datasets from the shared data library

# NGS Data Quality: Further reading & Resources

FastQC Documentation

Read Quality Assessment & Improvement

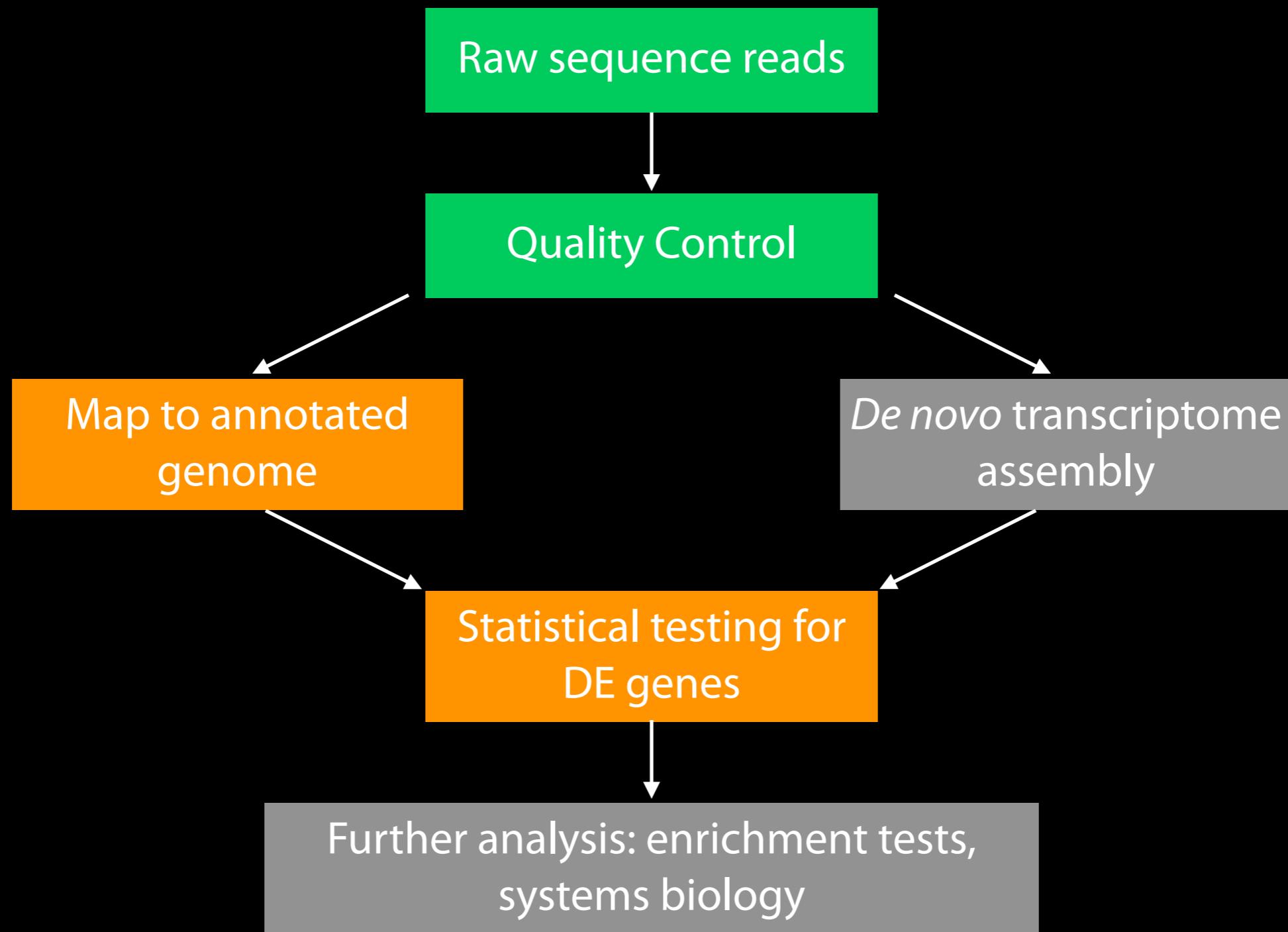
by Joe Fass

From the UC Davis 2013 Bioinformatics Short Course

Manipulation of FASTQ data with Galaxy

by Blankenberg, *et al.*

# RNA-Seq DE analysis steps



# RNA-seq Exercise: Mapping with TopHat

Create a new history

Import all datasets from library:

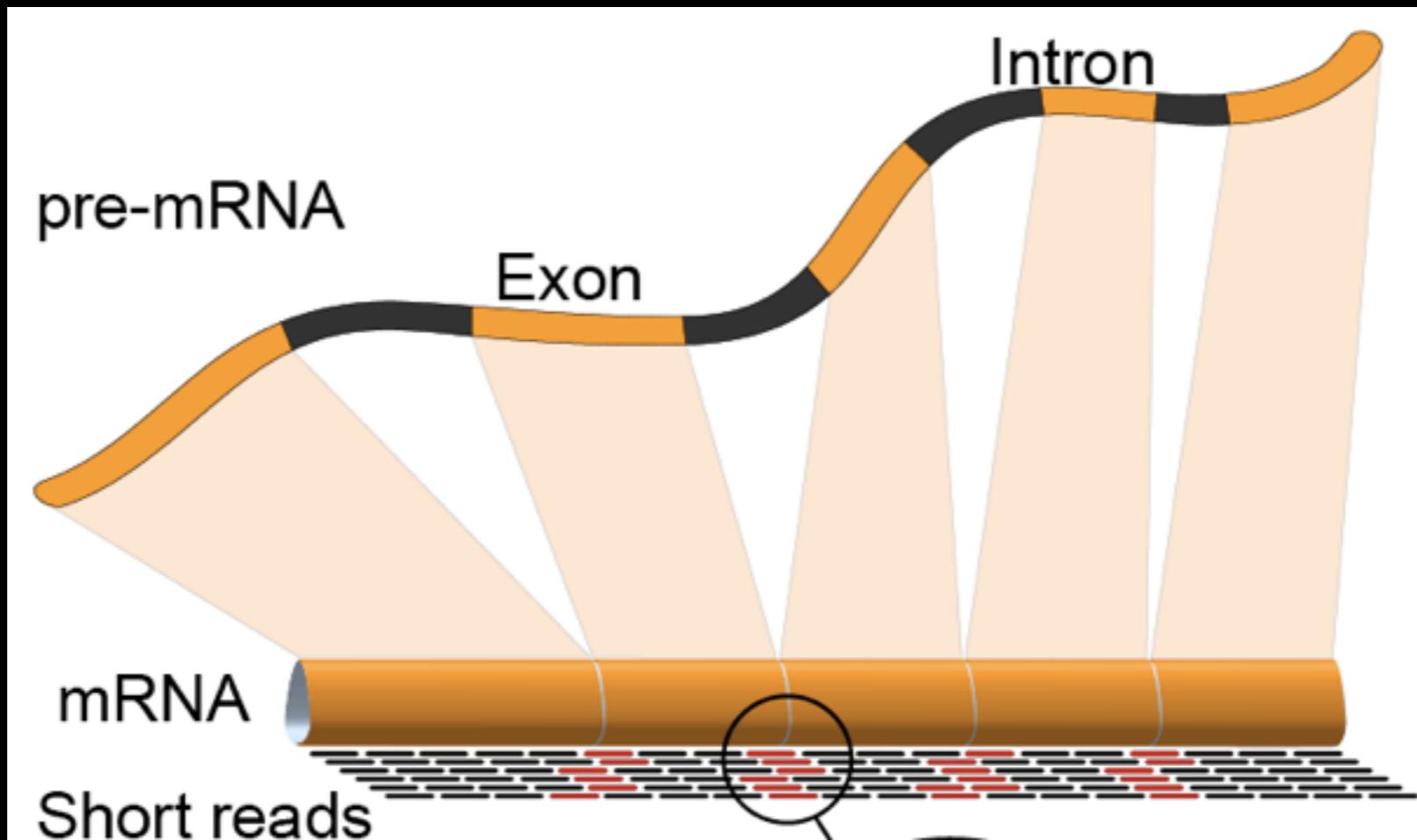
**RNA-Seq Example → Reads, Post-QC, Re-Paired**

All trimmed FASTQ and **genes\_chr12.gtf** (from  
**Reference directory**)

**NGS: RNA Analysis → TopHat2**

# RNA-seq Exercise: Mapping with TopHat

- Tophat looks for best place(s) to map reads, and best places to insert introns



- *Imagine pages and pages of discussion on the intricacies and pitfalls of RNA-seq mapping here.*

# Alignment output files

## SAM

- plain text file, tab separated columns
- "a huge spreadsheet"
- inefficient to read and store

## BAM

- a compressed version of SAM (~80% less storage)
- can be indexed (fast access to subsections)
- needs to be sorted to be useful however

## Standardized format

- readable by most software

# Anatomy of SAM file

```
Read1 113 1 497 37 37M      15 38662 0 CGGGTCTGACC 0;=====9;>>> NM:i:0
Read2 213 1 497 37 37M      15 37662 0 CGGGTCTGACC 0;====45;>>> NM:i:1 XM:i:3
Read3 337 1 497 37 37M      15     38662 0 CGGGTCTGACC ;==9;>>><>; NM:i:0
Read4 615 1 497 37 36MD1 15     447 0 CGGGTCTGACC 0;==5"=69;>> NM:i:0
Read5 844 1 497 37 37M      15    1445 0 CGGGTCTGACC =====9;>>> NM:i:0
```

One line per original read sequence

- Big!
- Where it aligned (if at all)
- How much of it aligned (soft/hard clipping)
- Mapping quality, likelihood correctly aligned
- Any differences to the reference (CIGAR string)
- Lots of other stuff (aligner dependent)
- Does not contain the reference sequence

# TopHat2 basic parameters

- We're using Paired-end data
- Specifying two mate-pair FASTQ input files
- Mapping against hg19 reference genome

# Mapping with TopHat: mean inner distance

Expected distance between paired end reads

- Determined by sample prep
- We'll use **90\*** for **mean inner distance**
- We'll use **50** for **standard deviation**

\* The library was constructed with the typical Illumina TruSeq protocol, which is supposed to have an average insert size of 200 bases. Our reads are 55 bases (R1) plus 55 bases (R2). So, the Inner Distance is estimated to be  $200 - 55 - 55 = 90$

# Mapping with TopHat: Use Existing Annotations?

You can bias TopHat towards known annotations

- Use Own Junctions → Yes
- Use Gene Annotation → Yes
- Gene Model Annotation → `genes_chr12.gtf`
- Only look for supplied junctions → Yes

# Mapping with TopHat: Make it quicker?

## Warning: Here be dragons!

- Allow indel search → No
- Use Coverage Search → No (wee dragons)

TopHat generates its database of possible splice junctions from two sources of evidence. The first and strongest source of evidence for a splice junction is when two segments from the same read (for reads of at least 45bp) are mapped at a certain distance on the same genomic sequence or when an internal segment fails to map - again suggesting that such reads are spanning multiple exons. With this approach, "GT-AG", "GC-AG" and "AT-AC" introns will be found *ab initio*. The second source is pairings of "coverage islands", which are distinct regions of piled up reads in the initial mapping. Neighboring islands are often spliced together in the transcriptome, so TopHat looks for ways to join these with an intron. We only suggest users use this second option (`--coverage-search`) for short reads (< 45bp) and with a small number of reads (<= 10 million). This latter option will only report alignments across "GT-AG" introns

# Mapping with TopHat: Max # of Alignments Allowed

Some reads align to more than one place equally well.

For such reads, how many should TopHat include?

If more than the specified number, TopHat will pick those with the best mapping score.

TopHat break ties randomly.

TopHat assigns equal fractional credit to all  $n$

Instructs TopHat to allow up to this many alignments to the reference for a given read, and choose the alignments based on their alignment scores if there are more than this number. The default is 20 for read mapping. Unless you use --report-secondary-alignments, TopHat will report the alignments with the best alignment score. If there are more alignments with the same score than this number, TopHat will randomly report only this many alignments. In case of using --report-secondary-alignments, TopHat will try to report alignments up to this option value, and TopHat may randomly output some of the alignments with the same score to meet this number.

Mapping with TopHat: Lets do it some more!

NGS: RNA Analysis → TopHat2

for the remaining replicates

Or not.

# RNA-Seq Mapping With TopHat: Resources

[RNA-Seq Concepts, Terminology, and Work Flows](#)  
by Monica Britton

[Aligning PE RNA-Seq Reads to a Genome](#)  
by Monica Britton

both from the [UC Davis 2013 Bioinformatics Short Course](#)

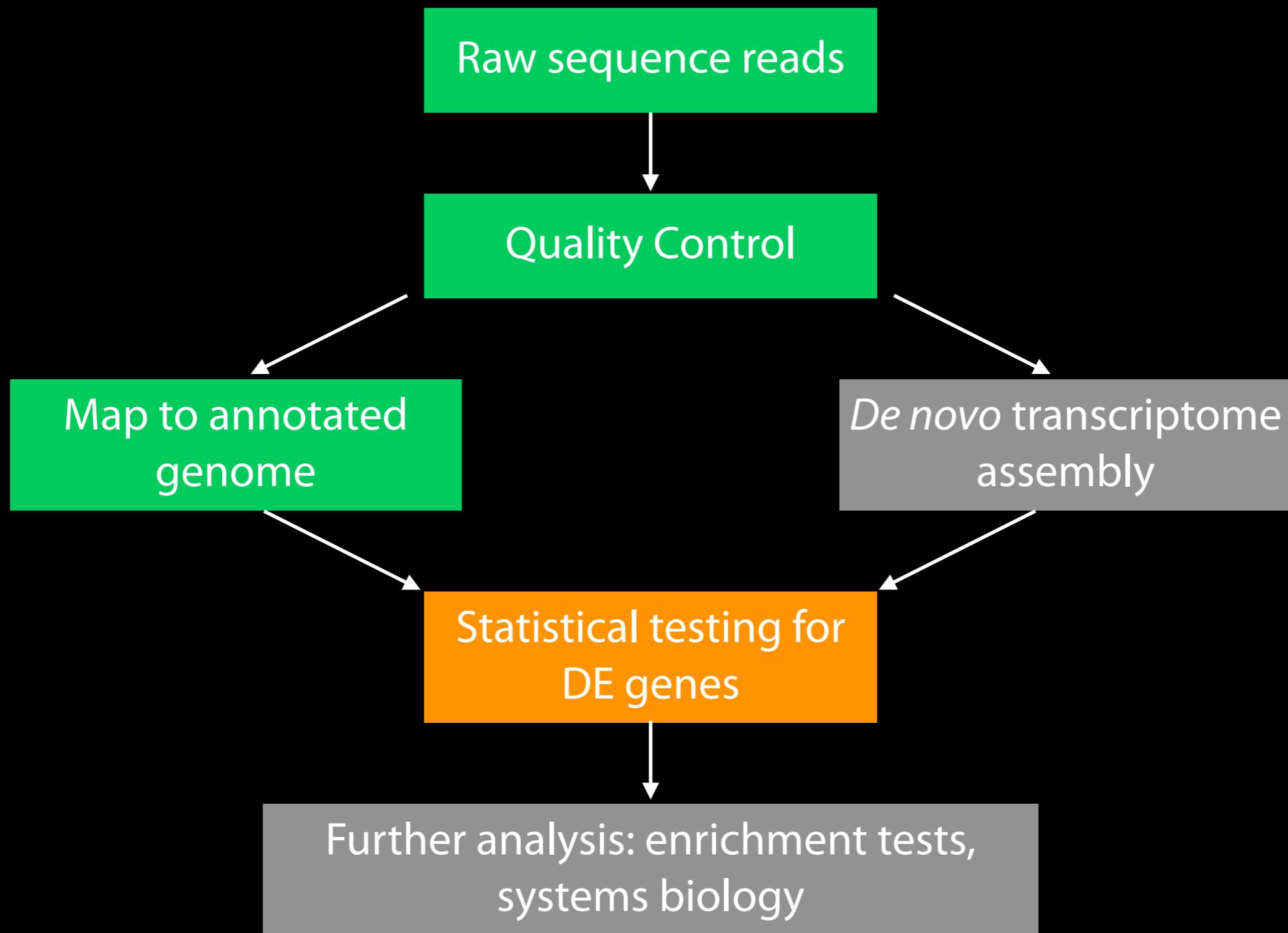
[RNA-Seq Analysis with Galaxy](#)  
by [Jeroen F.J. Laros](#), [Wibowo Arindrarto](#), [Leon Mei](#)

from the [GCC2013 Training Day](#)

[RNA-Seq Analysis with Galaxy](#)  
by [Curtis Hendrickson](#), [David Crossman](#), [Jeremy Goecks](#)

from the [GCC2012 Training Day](#)

# RNA-Seq DE analysis steps



A person is completely buried in a large pile of white, crumpled paper. Only their hands and upper torso are visible above the surface. The person is wearing a white button-down shirt. Their left hand is raised, palm facing forward, while their right hand holds a small, rectangular white card with the word "HELP" printed on it in a simple, black, sans-serif font.

HELP



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

- **Accessible:** Users without programming experience can easily specify parameters and run tools and workflows.
- **Reproducible:** Galaxy captures information so that any user can repeat and understand a complete computational analysis.
- **Transparent:** Users share and publish analyses via the web and create Pages, interactive, web-based documents that describe a complete analysis.

This is the Galaxy Community Wiki. It describes all things Galaxy.



[Early Registration & Abstract Submission](#)

are now open

Galaxy	•	2
Australasia	•	0
Workshop	•	1

**24-25 March,  
Melbourne**

## Use Galaxy

### [Use Galaxy](#)

Galaxy's [public service web site](#) makes analysis tools, genomic data, tutorial demonstrations, persistent workspaces, and publication services available to any scientist. Extensive [user documentation](#) (applicable to any [public](#) or local Galaxy instance) is available on [this wiki](#) and [elsewhere](#).



### [Deploy Galaxy](#)

Galaxy is open source for all organizations. Local Galaxy servers can be set up by [downloading and customizing](#) the Galaxy application.

- Admin
- Cloud
- Galaxy Appliance



## Community & Project

Galaxy has a large and active user community and many ways to [Get Involved](#).

- [Community](#)
- [News](#)
- [Events](#)
- [Support](#)
- [Galaxy Project](#)

## Contribute

- **Users:** Share your histories, workflows, visualizations, data libraries, and [Galaxy Pages](#), enabling others to use and learn from them.
- **Deployers and Developers:** Contribute tool definitions to the [Galaxy Tool Shed](#) (making it easy for others to use those tools on their installations), and code to the core release.
- **Everyone:** [Get Involved!](#)

## Contribute

[Tool Shed](#) • [Share](#)  
[Issues & Requests](#)  
[Teach](#) • [Support](#)



# Events

# News



Events

## Galaxy Event Horizon

Events with Galaxy-related content are listed here.

Also see the [Galaxy Events Google Calendar](#) for a listing of events and discussions in the Galaxy Community. This is also available as an [RSS feed](#) .

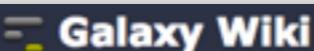
If you know of any event that should be added to this page and/or to the Galaxy project, send it to [outreach@galaxyproject.org](mailto:outreach@galaxyproject.org).

For events prior to this year, see the [Events Archive](#).

## Upcoming Events



Date	Topic/Event	Venue/Location
January 15-16	Accessible and Reproducible Genomics at Scale with Galaxy	Revolutionizing Next-Generation Sequencing Belgium
January 19-20	NGS pipelines with Galaxy	e-Infrastructures for Massively Parallel Computing Sweden
January 22	Baltimore Area Galaxy Meetup	Baltimore, Maryland, United States
January 28	ChIP-Seq Analysis and Visualization using Galaxy and IGB Tutorial	Online
February 9-13	Analyse bioinformatique de séquences sous Galaxy	Montpellier, France
	Accessible and Reproducible Large-Scale	Genome and Transcriptome Analysis, part of Molecular Medicine Tri-



News

## News

Announcements of interest to the Galaxy Community. These can include items from the Galaxy Team or the Galaxy community and can address anything that is of wide interest to the community.

The Galaxy News is also available as an [RSS feed](#) .

See [Add a News Item](#) below for how to get an item on this page, and the [RSS feed](#). Older news items are available in the Galaxy [News Archive](#).

### See also

- [Galaxy News Briefs](#)
- [Galaxy Updates](#)
- [Galaxy on Twitter](#)
- [Events](#)
- [Learn](#)
- [Support](#)
- [About the Galaxy Project](#)

Login | Search:

## News Items

- [January 2014 CloudMan Release](#)
- [GCC2014 Training Day Topics: Vote!](#)
- [January 2014 Galaxy Update](#)
- [2013 Galaxy Day Report](#)
- [Galaxy Community Log Board](#)
- [Galaxy Deployment Catalog](#)
- [Nominate 2014 Training Day Topics](#)
- [December 2013 Galaxy Update](#)
- [Nov 04, 2013 Galaxy Distribution](#)
- [November 2013 Galaxy Update](#)
- [December Bioinformatics Boot Camps](#)
- [GCC2014: Save These Dates!](#)
- [Galaxy Day, 4 décembre à Paris](#)

[News Archive](#)



# Unified Search: <http://galaxyproject.org/search>

Galaxy Web Search

Google™ Custom Search

Search the entire set of Galaxy web sites and mailing lists using Google.

[Run this search at Google.com \(useful for bookmarking\)](#)

Want a [different search?](#)

[Project home](#)

Galaxy Web Search

RNA-Seq

All Tools Email & Biostar Doc Source code Shared Abstracts Requests

About 6,400 results (0.39 seconds)

*Find*

Everything on ...  
Tools for ...  
Email about ...  
Source code for ...  
Published Histories, Pages, Workflows, about ...

About 6,400 results (0.39 seconds)

Related feature requests  
Papers using Galaxy for ...  
Documentation on ...

# Galaxy Resources & Community: Videos

vimeo [Join](#) Log In Create Watch Upload Search

## Galaxy Project PLUS

Joined 1 year ago 9 everywhere

“How to”  
screencasts on  
using and  
deploying  
Galaxy

Talks from  
previous  
meetings.

<http://vimeo.com/galaxyproject>

# Galaxy Resources and Community: Mailing Lists

<http://wiki.galaxyproject.org/MailingLists>

## Galaxy-Announce

Project announcements, low volume, moderated

Low volume ( 47 posts in 2014, 4100+ members)

## Galaxy-Dev

Questions about developing for and deploying Galaxy

High volume (2700 posts in 2014, 1000+ members)

# More interactive help?

<https://biostar.usegalaxy.org/>

LATEST OPEN RNA-SEQ CHIP-SEQ SNP ASSEMBLY FORUM PLANET ALL »

 **Biostars**  
GALAXY EXPLAINED

Enis Afgan · 170 | [Logout](#)

about · faq · rss 

Community Messages Votes My Posts My Tags Following Bookmarks New Post

Live search: start typing... or [Classic search](#)

Limit to: all time <prev · 3,518 results · page 1 of 101 · next> Sort by: update

0 votes	1 answer	11 views	<a href="#">new</a> <a href="#">api</a> <a href="#">student</a> <a href="#">researchproject</a>	written 11 hours ago by <a href="#">jchen015</a> · 0 · updated 7 hours ago by <a href="#">Bjoern Gruening</a> + 1.8k
0 votes	1 answer	10 views	<a href="#">tools</a> <a href="#">galaxy</a>	written 12 hours ago by <a href="#">morteza-sharifinia</a> · 0 · updated 11 hours ago by <a href="#">Bjoern Gruening</a> + 1.8k
0 votes	0 answers	11 views	<a href="#">galaxy</a> <a href="#">rna-seq</a>	written 1 day ago by <a href="#">morteza-sharifinia</a> · 0
0 votes	0 answers	8 views	<a href="#">galaxy</a>	written 1 day ago by <a href="#">nkaplin1</a> · 0
8 votes	8 answers	615 views	<a href="#">biolinux</a> <a href="#">config</a> <a href="#">local</a>	written 7 months ago by <a href="#">alfonso.cervantes</a> · 10 · updated 2 days ago by <a href="#">Suzanne Gomes</a> · 100

**Recent Votes**

- A: Unable to upload (.bam) data file to local Galaxy
- Baltimore Area Galaxy Meetup @ Johns Hopkins, January 22
- January 13, 2015 Galaxy Distribution
- A: Bioblend - toolClient.run\_tool does not work with more than one params ?
- A: Custom toolshed tool: job memory requirements
- A: Custom toolshed tool: job memory requirements
- A: Feedback on using uWSGI

**Recent Locations · All »**

- Australia, 1 hour ago
- United States, 2 hours ago
- Iran, Islamic Republic Of, 3 hours ago
- United States, 5 hours ago
- United States, 5 hours ago

# Community can create, vote and comment on issues

HOME TOUR GOLD BUSINESS CLASS BLOG

Trello

Sign Up

Log In

Want to subscribe, vote or comment on these cards? [Sign up for free](#) or [learn more about Trello](#)

Galaxy: Development Public

Inbox

To add cards, use  
<http://galaxyproject.org/trello>  
↳ 4 votes 2 comments

To request reference genome,  
comment on this card.  
1 vote 5 comments 0 likes

Automatically install tools that  
workflow needs  
3 votes 1 comment

Toolshed installation fails silently  
3 votes 1 comment

Handle cluster job preemption  
2 votes 1 comment

Return code 271 causes traceback  
for PBS torque  
1 vote 2 comments

Visiting Published  
Workflows/History behaves poorly  
1 vote 1 comment

BUG: Job handler error after  
installation of new tool from  
toolshed  
1 vote 1 comment

Taxonomy scripts  
(scripts/taxonomy/\*) should be  
bundled with toolshed tools  
1 vote 1 comment

BUG: Tool shed repository export  
to capsule does not always  
capture all dependencies

Tool Requests

595: Add SAMTools "Sort"  
↳ 4 votes 13 comments

601: SAM-to-BAM tool  
enhancements  
↳ 2 votes 1 comment

biomart data source fails to load  
on Main  
↳ 3 votes

Tools: Add tool to generate  
simulated reads to Main  
↳ 3 votes 1 comment

fastq\_to\_fasta tool on main tool  
shed  
↳ 2 votes 5 comments

default max insert size of Bowtie2  
should be increased  
↳ 2 votes 5 comments

307: A tool to produce a set of  
random intervals.  
↳ 2 votes 2 comments

Wrapper for bigWigToWig  
↳ 2 votes 1 comment

Converter Tool: SAM to BAM  
enhancements  
↳ 2 votes

New Tool: convert IUPAC chars to  
N

Bug Reports

Bug: tools are executed in  
workflows with "negated" boolean  
parameters  
↳ 3 votes 2 comments

Bug: SICER on Main dependency  
issue  
↳ 2 votes 20 comments 3 likes

Profile Annotations bad values  
when "select all"  
↳ 1 vote 5 comments

Filter pileup tool doesn't recognize  
pileup output data  
↳ 1 vote 2 comments

Bug: Odd Fetch Taxonomy tool  
behavior  
↳ 1 vote 1 comment

Strip message after pause jobs  
resumed  
↳ 1 vote 1 comment

The option from\_file="infernal.loc"  
is broken.  
↳ 1 vote 1 comment

At least some confirmation dialogs  
are broken.  
↳ 1 vote

Ideas

697: Workflow job control  
functions  
↳ 10 votes 9 comments

User Metrics and Analytics  
↳ 3 votes 3 comments 1 like

Tuxedo RNA-seq tools: report  
command-line  
↳ 2 votes 3 comments

Tools: Incorporate key Cuffdiff  
output files for Cummerbund  
↳ 2 votes 1 comment 3 likes

Moving objects between Galaxy  
instances, data federation,  
distributed storage, and data  
locality  
↳ 2 votes 1 comment 3 likes

Enhance DataManagers for  
updated UCSC ref genome  
retrieval  
↳ 1 vote 1 comment

Workflow Editor: Provide explicit  
access to implicit datatype  
converter tools  
↳ 1 vote

Workflows and Rerun  
↳ 1 comment 0 likes

Pull Requests / Patches

665: Patch for FASTQ paired-end  
issue  
↳ 17 votes 1 comment

Tools: Bowtie Wrapper Pull  
Requests from Community  
↳ 3 votes 9 comments 1 like

Custom Authentication Providers  
↳ 2 votes 6 comments

add ma seq metrics and  
downsample sam to picard tools  
↳ 2 votes 4 comments 1 like

Please merge patch to bowtie2  
wrapper (adds support for  
mapping fasta files)  
↳ 2 votes 1 comment 1 like

[galaxy-dev] Patch for libxslt and  
libxml packages  
↳ 1 vote 2 comments

Pull Request #627 - Add planemo  
test file to tool shed upload  
blacklist  
↳ 1 vote

Pull Request #625 - Add RNA  
dotplot matrix as datatype.  
↳ 2 votes

Pull Request #624 - Add  
\_getitem\_() to select parameters  
to enable iterations  
↳ 1 vote

Menu

Members



Activity

CE Carl Eberhard added API:  
allow finer control over  
cache/caching of results to  
Ideas and joined. an hour ago

CE Carl Eberhard added UI:  
manually update the js  
libraries to In Progress and  
joined. today at 1:26 pm

G g2robot on Cannot  
download BAM Index when  
dataset is in a collection.  
Submitted by  
@philipmabon  
today at 11:34 am

G g2robot added Cannot  
download BAM Index when  
dataset is in a collection, to  
Inbox. today at 11:34 am

CE Carl Eberhard moved Build,  
UI: restructure mvc/ui.js from  
Testing to In default branch.  
today at 11:30 am

G g2robot added Pull Request  
#639 - Fix for  
LibraryDatasetToolParameter  
to Pull Requests / Patches

<http://bit.ly/gxytrello>



[wiki.galaxyproject.org/Events/GCC2015](http://wiki.galaxyproject.org/Events/GCC2015)



# Galaxy Resources & Community: CiteULike Group



CiteULike   Group: Galaxy   [Search](#)   [Register](#)   [Log in](#)

## Group: Galaxy - library 2020 articles [RSS](#)

[Search](#)   [Copy](#)   [Export](#)   [Sort](#)   [Hide Details](#)

- ✓ [Open pipelines for integrated tumor genome profiles reveal differences between pancreatic cancer tumors and cell lines](#) *Cancer Medicine* (December 2014), pp. n/a-n/a, doi:10.1002/cam4.360  
by [Jeremy Goecks](#), [Bassel F. El-Rayes](#), [Shishir K. Maitel](#), [H. Jean Khoury](#), [James Taylor](#), [Michael R. Rossi](#)  
posted to [methods](#) [shared](#) [usemain](#) by [galaxyproject](#) to the group [Galaxy](#) on 2015-01-09 00:22:41 ★★★★☆  
[Abstract](#)
- ✓ [RNA-Seq reveals a xenobiotic stress response in the soybean aphid, \*Aphis glycines\*, when fed aphid-resistant soybeans](#) *BMC Genomics*, Vol. 15, No. 1. (16 November 2014), 972, doi:10.1186/1471-2164-15-972  
by [Raman Bansal](#), [M. A. R. Mian](#), [Omprakash Mittapalli](#), [Andy P. Michel](#)  
posted to [methods](#) [uselocal](#) by [galaxyproject](#) to the group [Galaxy](#) on 2015-01-09 00:15:52 ★★  
[Abstract](#)
- ✓ [Draft Genome Sequence of the Cellulolytic Bacterium \*Clavibacter\* sp. CF11, a Strain Producing Cold-Active Cellulase](#) *Genome Announcements*, Vol. 3, No. 1. (26 February 2015), e01304-14, doi:10.1128/genomeA.01304-14  
by [Ying Du](#), [Bo Yuan](#), [Yonghui Zeng](#), et al.  
posted to [methods](#) [usemain](#) by [galaxyproject](#) to the group [Galaxy](#) on 2015-01-09 00:03:34 ★★  
[Abstract](#)
- ✓ [Identifying transcriptional cis -regulatory modules in animal genomes](#) *Wiley Interdisciplinary Reviews: Developmental Biology* (December 2014), pp. n/a-n/a, doi:10.1002/wdev.168  
by [Kushal Suryamohan](#), [Marc S. Halfon](#)  
posted to [refpublic](#) by [galaxyproject](#) to the group [Galaxy](#) on 2015-01-09 00:01:02 ★★  
[Abstract](#)
- ✓ [Dynamic shifts in occupancy by TAL1 are guided by GATA factors and drive large-scale reprogramming of gene expression during hematopoiesis](#) *Genome Research*, Vol. 24, No. 12. (01 December 2014), pp. 1945-1962, doi:10.1101/gr.164830.113  
by [Weisheng Wu](#), [Christopher S. Morrissey](#), [Cheryl A. Keller](#), et al.  
posted to [methods](#) by [galaxyproject](#) to the group [Galaxy](#) on 2015-01-08 23:56:32 ★★  
[Abstract](#)
- ✓ [Rampant software errors undermine scientific results](#) *F1000Research* (11 December 2014), doi:10.12688/f1000research.5930.1  
by [David A. W. Soergel](#)  
posted to [workbench](#) by [galaxyproject](#) to the group [Galaxy](#) on 2015-01-08 23:49:39 ★★ [along with 2 people and 1 group](#)  
[Abstract](#)
- ✓ [Reproducibility in Science](#) *Circulation Research*, Vol. 116, No. 1. (02 January 2015), pp. 116-126, doi:10.1161/circresaha.114.303819  
by [C. Glenn Begley](#), [John P. A. Ioannidis](#)  
posted to [reproducibility](#) by [galaxyproject](#) to the group [Galaxy](#) on 2015-01-08 23:21:23 ★★ [along with 1 person](#)  
[Abstract](#)

### Group Tags

All tags in the group Galaxy

Filter:   
[Display as Cloud](#)

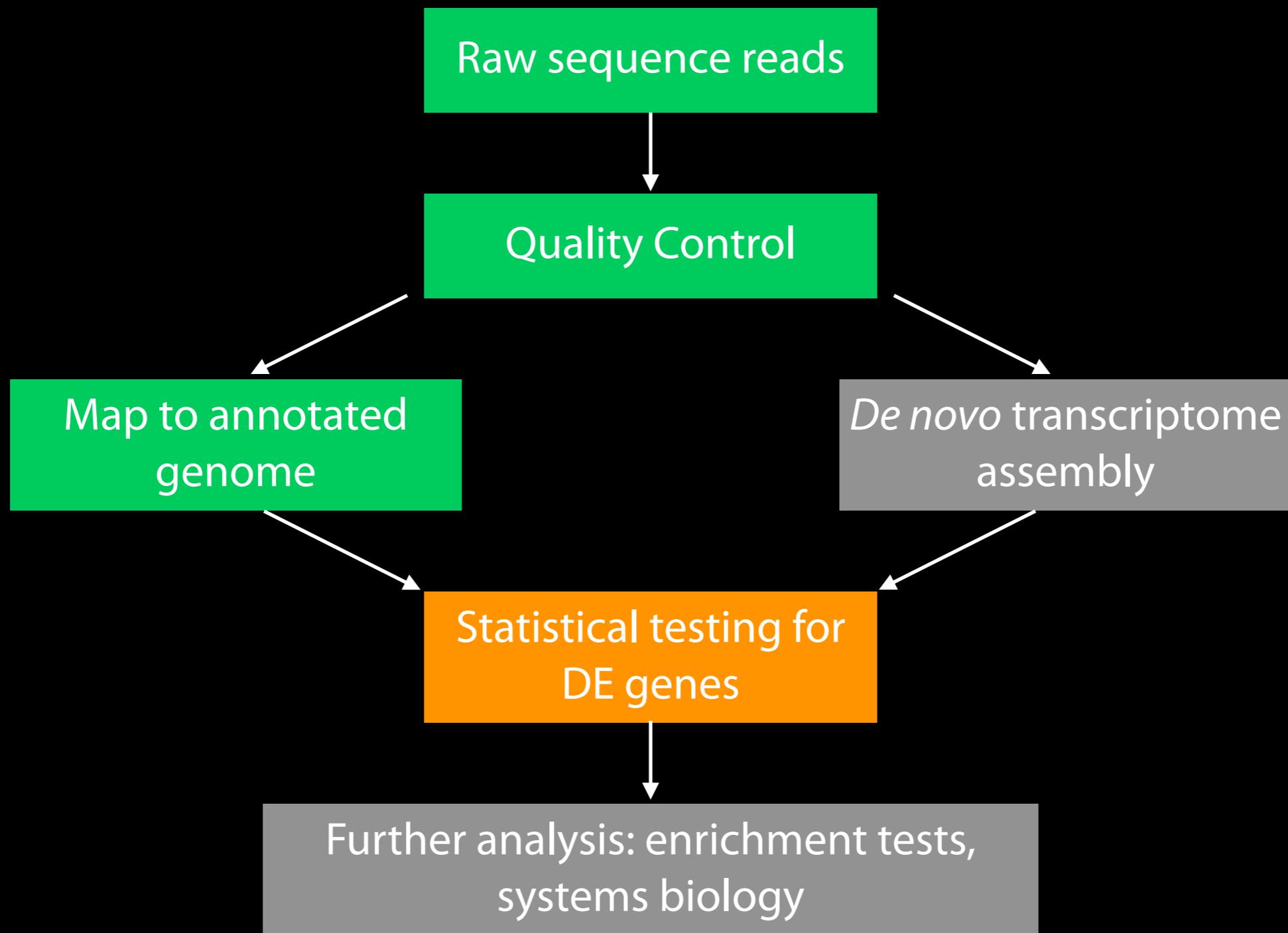
<a href="#">methods</a>	973
<a href="#">workbench</a>	633
<a href="#">usemain</a>	190
<a href="#">tools</a>	146
<a href="#">isgalaxy</a>	115
<a href="#">cloud</a>	77
<a href="#">usepublic</a>	76
<a href="#">uselocal</a>	70
<a href="#">shared</a>	70
<a href="#">other</a>	54
<a href="#">unknown</a>	51
<a href="#">reproducibility</a>	42
<a href="#">howto</a>	41
<a href="#">project</a>	38
<a href="#">refpublic</a>	38
<a href="#">visualization</a>	12
<a href="#">usecloud</a>	3

Over  
2000  
papers

17 tags

<http://bit.ly/gxycul>

# RNA-Seq DE analysis steps



# Cuffdiff

- Part of the Tuxedo RNA-Seq Suite (as are TopHat and Bowtie)
- Widely used and widely installed on Galaxy instances

**NGS: RNA Analysis → Cuffdiff**

# Cuffdiff

- Running with 2 Groups: MeOH and R3G
- Each group has 3 replicates
- genes\_chr12.gtf as the transcript

# Cuffdiff

Produces many output files, all explained in doc

We'll focus on gene differential expression testing

test_id	gene_id	gene	locus	sample_1	sample_2	status	value_1	value_2	log2(fold_change)	test_stat	p_value	q_value	significant
A2M	A2M	A2M	chr12:9217772-9268558	MeOH	R3G	NOTEST	3.32147	3.13694	-0.0824644	0	1	1	no
A2M-AS1	A2M-AS1	A2M-AS1	chr12:9217772-9268558	MeOH	R3G	NOTEST	7.45797	13.9413	0.902515	0	1	1	no
A2ML1	A2ML1	A2ML1	chr12:8975149-9029381	MeOH	R3G	NOTEST	4.83055	7.79884	0.691072	0	1	1	no
A2MP1	A2MP1	A2MP1	chr12:9381128-9386803	MeOH	R3G	NOTEST	2.49656	0	-inf	0	1	1	no
AAAS	AAAS	AAAS	chr12:53701239-53715412	MeOH	R3G	OK	269.035	159.23	-0.756683	-2.22857	0.0005	0.00194017	yes
AACS	AACS	AACS	chr12:125549924-125627871	MeOH	R3G	NOTEST	29.2933	35.0339	0.258178	0	1	1	no
ABCB9	ABCB9	ABCB9	chr12:123405497-123451056	MeOH	R3G	NOTEST	4.68869	1.7732	-1.40283	0	1	1	no
ABCC9	ABCC9	ABCC9	chr12:21950323-22089628	MeOH	R3G	OK	553.247	487.261	-0.18323	-2.02806	0.0004	0.00162143	yes
ABCD2	ABCD2	ABCD2	chr12:39945021-40013843	MeOH	R3G	OK	86.1377	172.795	1.00435	4.3436	5e-05	0.000246739	yes
ACACB	ACACB	ACACB	chr12:109577201-109706030	MeOH	R3G	NOTEST	8.45306	15.5772	0.881885	0	1	1	no
ACAD10	ACAD10	ACAD10	chr12:112123856-112194911	MeOH	R3G	NOTEST	21.8237	27.8326	0.350882	0	1	1	no
ACADS	ACADS	ACADS	chr12:121163570-121177811	MeOH	R3G	NOTEST	38.644	16.1739	-1.25658	0	1	1	no
ACRBP	ACRBP	ACRBP	chr12:6747241-6756580	MeOH	R3G	NOTEST	2.96987	3.26939	0.138621	0	1	1	no
ACSM4	ACSM4	ACSM4	chr12:7456927-7480969	MeOH	R3G	NOTEST	0	0	0	0	0	1	no
ACSS3	ACSS3	ACSS3	chr12:81471808-81649582	MeOH	R3G	NOTEST	0	0	0	0	0	1	no
ACTR6	ACTR6	ACTR6	chr12:100593864-100618202	MeOH	R3G	OK	475.594	421.324	-0.174799	-0.797581	0.1588	0.258406	no
ACVR1B	ACVR1B	ACVR1B	chr12:52345450-52390863	MeOH	R3G	NOTEST	32.5737	38.3075	0.233922	0	1	1	no
ACVRL1	ACVRL1	ACVRL1	chr12:52301201-52317145	MeOH	R3G	NOTEST	1.27713	2.16161	0.759201	0	1	1	no
ADAM1A	ADAM1A	ADAM1A	chr12:112336866-112339706	MeOH	R3G	NOTEST	30.0162	55.2154	0.879331	0	1	1	no
ADAMTS20	ADAMTS20	ADAMTS20	chr12:43748011-43945724	MeOH	R3G	NOTEST	0.453322	0.502067	0.147346	0	1	1	no
ADCY6	ADCY6	ADCY6	chr12:49159974-49182820	MeOH	R3G	NOTEST	9.32722	17.6743	0.922135	0	1	1	no
ADIPOR2	ADIPOR2	ADIPOR2	chr12:1800246-1897845	MeOH	R3G	OK	207.468	179.333	-0.210248	-1.02392	0.09	0.158988	no
AEBP2	AEBP2	AEBP2	chr12:19592607-19675173	MeOH	R3G	OK	143.039	128.293	-0.156957	-0.688267	0.2254	0.344537	no
AGAP2	AGAP2	AGAP2	chr12:58118075-58135944	MeOH	R3G	OK	98.2385	116.302	0.243511	0.935119	0.11475	0.198086	no
AICDA	AICDA	AICDA	chr12:8754761-8765442	MeOH	R3G	NOTEST	78.1514	63.4313	-0.301077	0	1	1	no
AKAP3	AKAP3	AKAP3	chr12:4724675-4754343	MeOH	R3G	NOTEST	6.12385	7.89626	0.366731	0	1	1	no
ALDH1L2	ALDH1L2	ALDH1L2	chr12:105413561-105478341	MeOH	R3G	NOTEST	7.11374	8.11722	0.190377	0	1	1	no
ALDH2	ALDH2	ALDH2	chr12:112204690-112247789	MeOH	R3G	NOTEST	12.8033	8.05635	-0.668321	0	1	1	no
ALG10	ALG10	ALG10	chr12:34175215-34181236	MeOH	R3G	NOTEST	54.8575	59.3459	0.11346	0	1	1	no
ALG10B	ALG10B	ALG10B	chr12:38710556-38723528	MeOH	R3G	NOTEST	43.8157	63.0457	0.524952	0	1	1	no
ALKBH2	ALKBH2	ALKBH2	chr12:109525992-109531293	MeOH	R3G	OK	679.517	297.183	-1.19316	-3.34255	5e-05	0.000246739	yes
ALX1	ALX1	ALX1	chr12:85674035-85695561	MeOH	R3G	NOTEST	0	0	0	0	0	1	no

# Cuffdiff: differentially expressed genes

Column	Contents
test_stat	value of the test statistic used to compute significance of the observed change in FPKM
p_value	Uncorrected P value for test statistic
q_value	FDR-adjusted p-value for the test statistic
status	Was there enough data to run the test?
significant	and, was the gene differentially expressed?

# Cuffdiff

- Column 7 (“status”) can be FAIL, NOTEST, LOWDATA or OK
  - Filter and Sort → **Filter**
    - `c7 == 'OK'`
- Column 14 (“significant”) can be yes or no
  - Filter and Sort → **Filter**
    - `c14 == 'yes'`

Returns the list of genes with

- 1) enough data to make a call, and
- 2) that are called as differentially expressed.

# Cuffdiff: Next Steps

Try running Cuffdiff with different **normalization** and **dispersion estimation** methods.

Try running Cuffdiff with **less replicas**.

Compare the differentially expressed gene lists.  
Which settings have what type of impacts on the results?

# RNA-Seq Differential Expression with Cuffdiff: Resources

RNA-Seq Concepts, Terminology, and Work Flows

by Monica Britton

from the UC Davis 2013 Bioinformatics Short Course

RNA-Seq Analysis with Galaxy

by Jeroen F.J. Laros, Wibowo Arindrarto, Leon Mei

from the GCC2013 Training Day

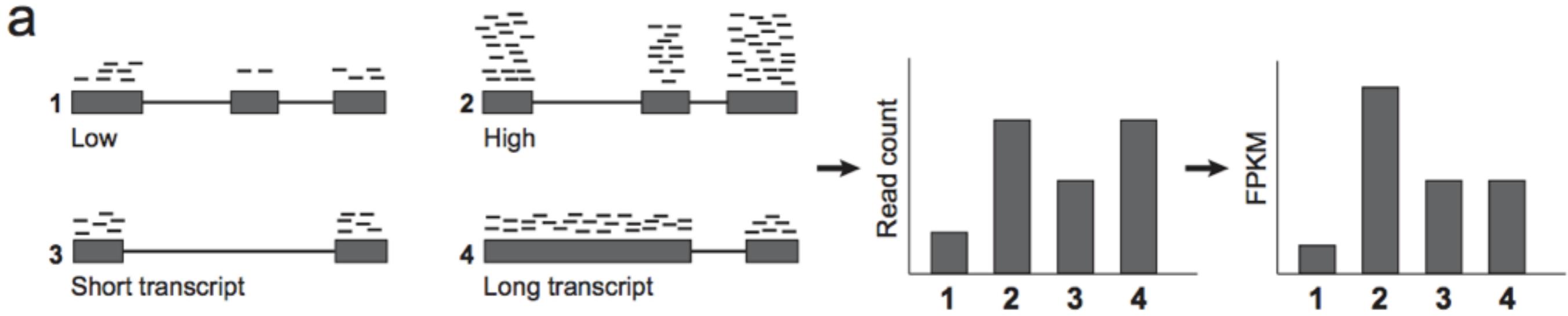
RNA-Seq Analysis with Galaxy and Alternative Tools

by Saskita Hiltemann, Youri Hoogstrate & Leon Mei

from the GCC2014 Training Day

# Cuffdiff?

Cuffdiff uses FPKM/RPKM as a central statistic.  
Total # mapped reads heavily influences FPKM/RPKM.  
Can lead to challenges when you have very highly  
expressed genes in the mix.



Garber et al, Nature Methods, 2011

# Cuffdiff Alternatives

Rapaport, *et al.*, “Comprehensive evaluation of differential gene expression analysis methods for RNA-seq data.”

*Genome Biology* 2013, 14:R95 doi:10.1186/gb-2013-14-9-r95

Reviews 7 packages

Each tool has it's own strengths and weaknesses.

What's a biologist to do?

# Alternatives: What's a biologist to do?

Learn the strengths and weaknesses of the tools you have ready access to. Are they a good match for the questions you are asking?

If not, then research alternatives, identify good options and then work with your bioinformatics/systems people to get access to those tools.

# Cuffdiff Alternatives

**Voom+Limma**

**EdgeR**

**DESeq**

R packages available as Galaxy tools; isolate between features based on expression counts and show a total count of a gene including all its isoforms

# Cuffdiff Alternatives

Take a simple, tab delimited list of features and read counts across different samples.  
First, have to create that list.

`htseq-count`

Is a tool that walks BAM files producing these lists

# Cuffdiff Alternatives

NGS: SAM Tools → htseq-count  
once for each BAM file

Join the 6 HTSeq datasets together on gene name

Cut out the duplicate gene name columns

NGS: RNA Analysis → Differential Count

# Cuffdiff Alternatives: Differential Counts

Output from three tools: Voom, EdgeR, DESeq

Output is a list of genes,  
sorted by adjusted P value,  
with lowest P values listed first

How many genes have an adjusted P value < 0.05  
for each of the tools?

# After DGE RNA-seq?

You get your “gene list”, finished?

- Validate
- Typically expect some false-positives
- Genes not in your list may be differentially expressed

Important to always remember

- Your list of genes is produced with an arbitrary significance threshold!

Next?

- Gene-set enrichment tests
- Novel transcripts, novel splice-variants, ...

# The Galaxy Team



Enis Afgan

Dannon Baker

Dan Blankenberg

Dave Bouvier

Marten Cech

John Chilton



Dave Clements

Nate Coraor

Carl Eberhard

Jeremy Goecks

Sam Guerler



Jen Jackson

Ross Lazarus

Anton Nekrutenko

Nick Stoler

James Taylor

Nitesh Turaga

<http://wiki.galaxyproject.org/GalaxyTeam>

# Thanks



Enis Afgan

Dave Clements

Clare Sloggett

Andrew Sharp

Galaxy Project

Johns Hopkins University

[outreach@galaxyproject.org](mailto:outreach@galaxyproject.org)