



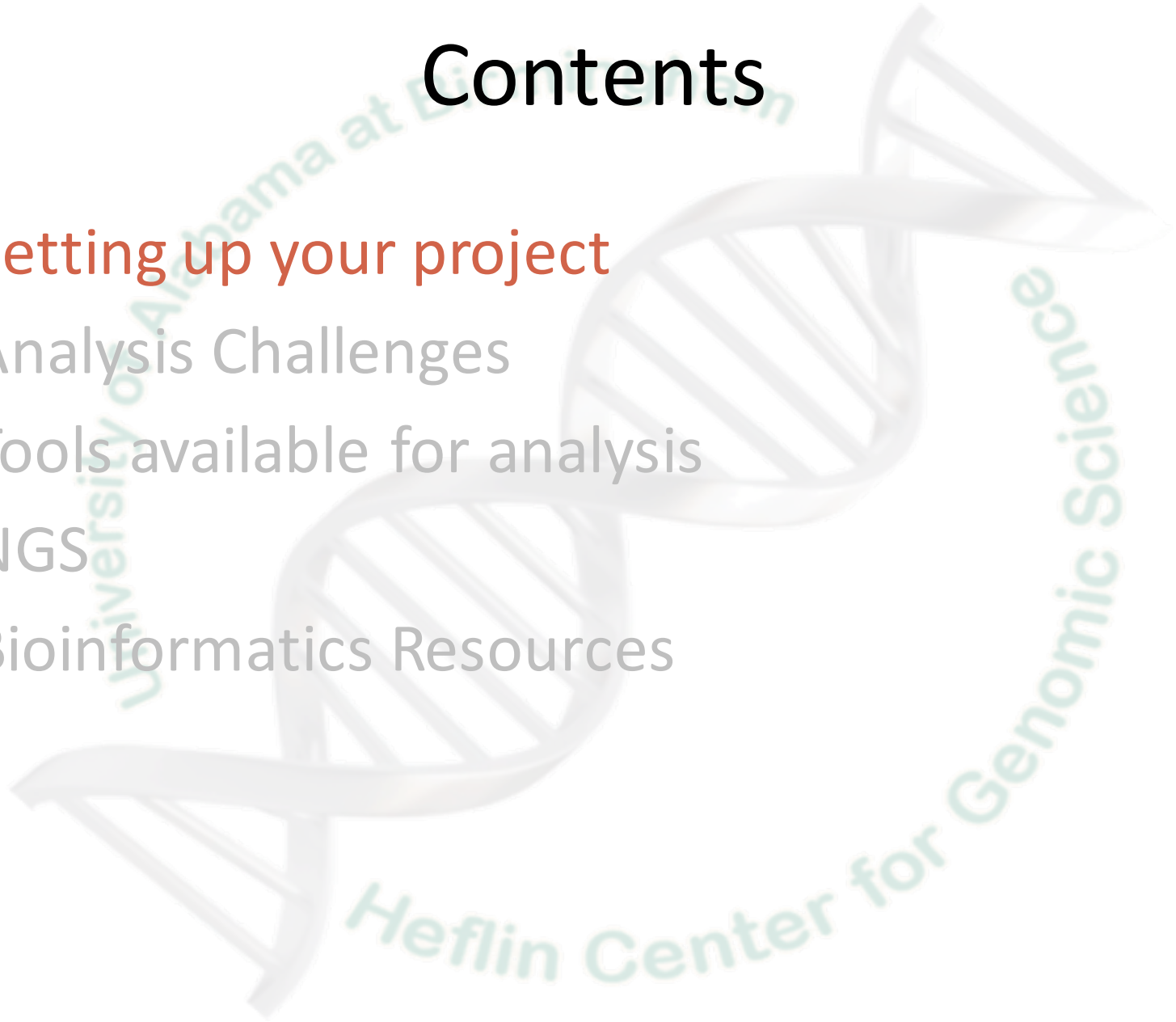
Approaches to Bioinformatics Data Analysis

David Crossman, Ph.D.
UAB Heflin Center for Genomic Science

Immersion Course

Contents

- **Setting up your project**
- Analysis Challenges
- Tools available for analysis
- NGS
- Bioinformatics Resources



Factors to Consider in Genomic Studies

- What is my study design/hypothesis?
- How comprehensive does the data need to be?
- How much work can I afford?
- What is the quality and quantity of my sample
- How fast do I need to have results?
- Will my grant get funded if I don't use the latest technology?

Study Design/Hypothesis

- I have an organism for which there are no off-the-shelf products for gene expression or sequence analysis.
- I need to comprehensively interrogate the entire genome of my model.
- I am studying a rare disease that I hypothesize to be attributed to private/rare/*de novo* mutation.

Next Generation Sequencing is a good choice.

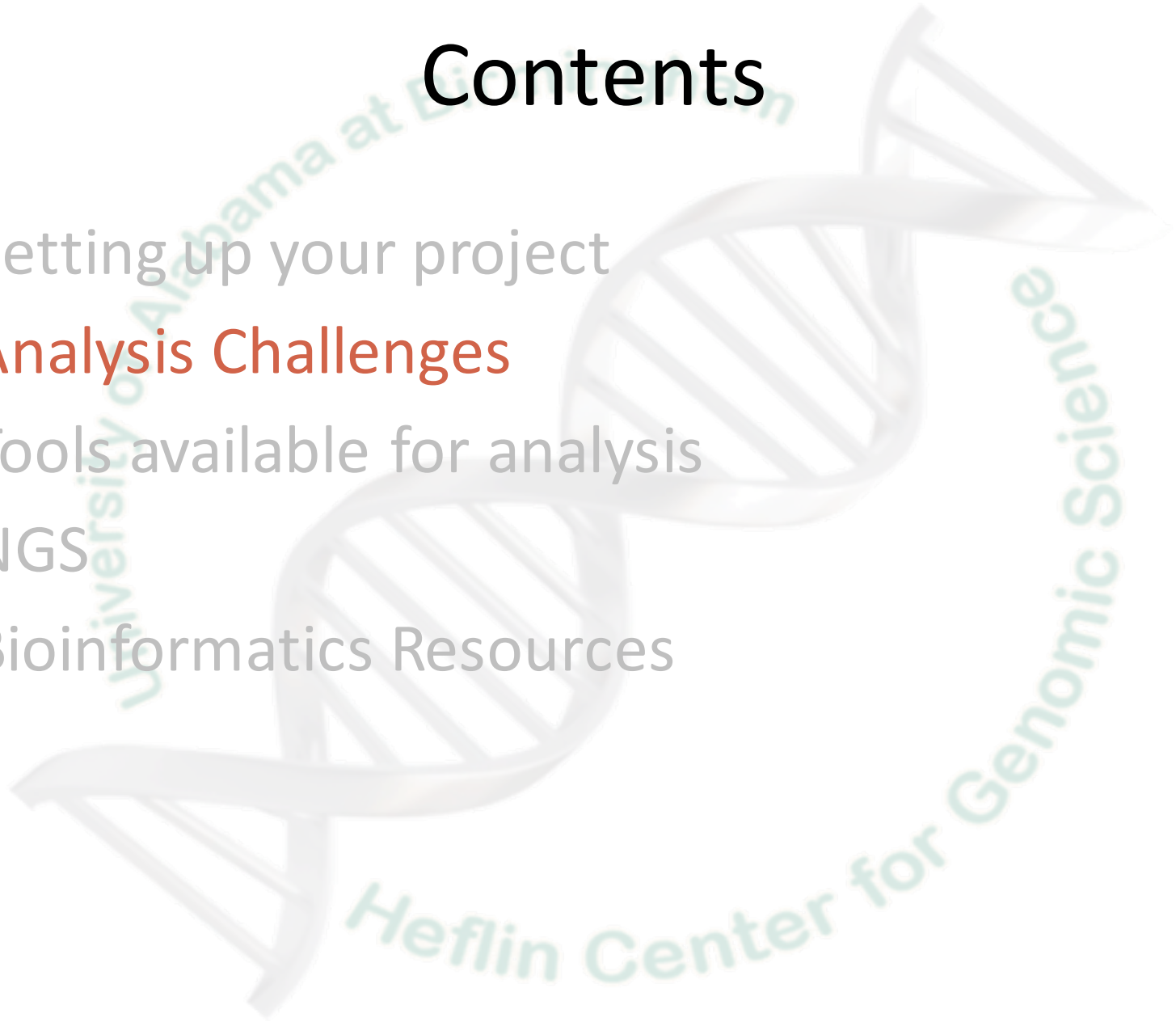
Study Design/Hypothesis

- I have an organism for which there are array-based products available.
- I want to expand my current work to identify new pathways involved in the physiology/organism I am studying.
- I have a large cohort of humans or animals and want to characterize all individuals.

Microarray-based assays are a good choice.

Contents

- Setting up your project
- **Analysis Challenges**
- Tools available for analysis
- NGS
- Bioinformatics Resources



NGS Analysis Challenges

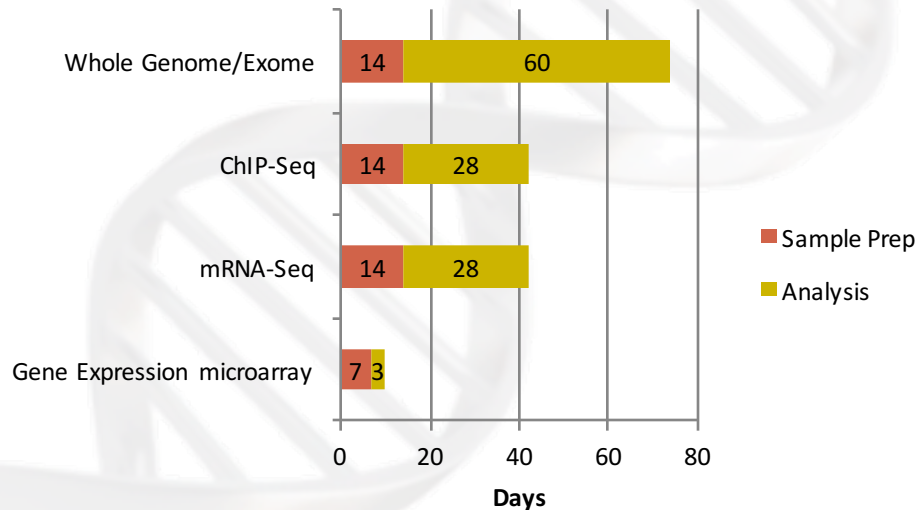
- Advanced technologies require substantial computing resources.
- File sizes:

Per Sample (in GB)	Raw Data	Aligned	Spreadsheet	Total
Gene Expression microarray	0.01	NA	0.002-0.005	0.012-0.015
mRNA-Seq	10-26	10-26	0.01-0.03	20.01-52.03
Exome	10-30	10-30	0.01-0.05	20.01-60.05
Whole Genome	250-500	250-500	0.01-0.1	500.01-1000.1

- Processed NGS files can be several TB in size.
- Software for NGS analysis rapidly evolving.
- Need for realistic understanding of data complexity and timeline for analysis.

Workflow Timeline

- Analysis time varies:



- Additional analyses will extend processing time

Contents

- Setting up your project
- Analysis Challenges
- **Tools available for analysis**
- NGS
- Bioinformatics Resources



Tools available to UAB investigators

- [Galaxy](#) – NGS analysis for those afraid of the “blinking cursor.”
- Command line tools to run on UAB’s Cheaha compute cluster:
 - TopHat
 - Cufflinks
 - Bowtie
 - SAMTools
 - BWA
 - GATK
 - PicardTools
 - MACS2
 - Bismark
 - STAR
 - HISAT
 - Stringtie
 - Ballgown
 - Trinity
 - Freebayes
 - Gemini
 - Circos
 - Homer
 - FastQC
 - SNPEff
 - ANNOVAR
 - muTect
 - MutSig
 - Velvet
 - Abyss
 - vcftools
 - AND MANY MORE!

Contents

- Setting up your project
 - Analysis Challenges
 - Tools available for analysis
 - **NGS**
 - Bioinformatics Resources
- **What is Galaxy**
 - What isn't Galaxy
 - FASTQ anatomy
 - Using Galaxy

University of Alabama at Birmingham

Heflin Center for Genomic Science

What is Galaxy

- GUI for genomics
 - for complete analyses: analyze, visualize, share, publish
- A free (for everyone) web service integrating a wealth of tools, compute resources, terabytes of reference data and permanent storage
- Open source software that makes integrating your own tools and data and customizing for your own site simple

For those afraid of the “blinking cursor!” |

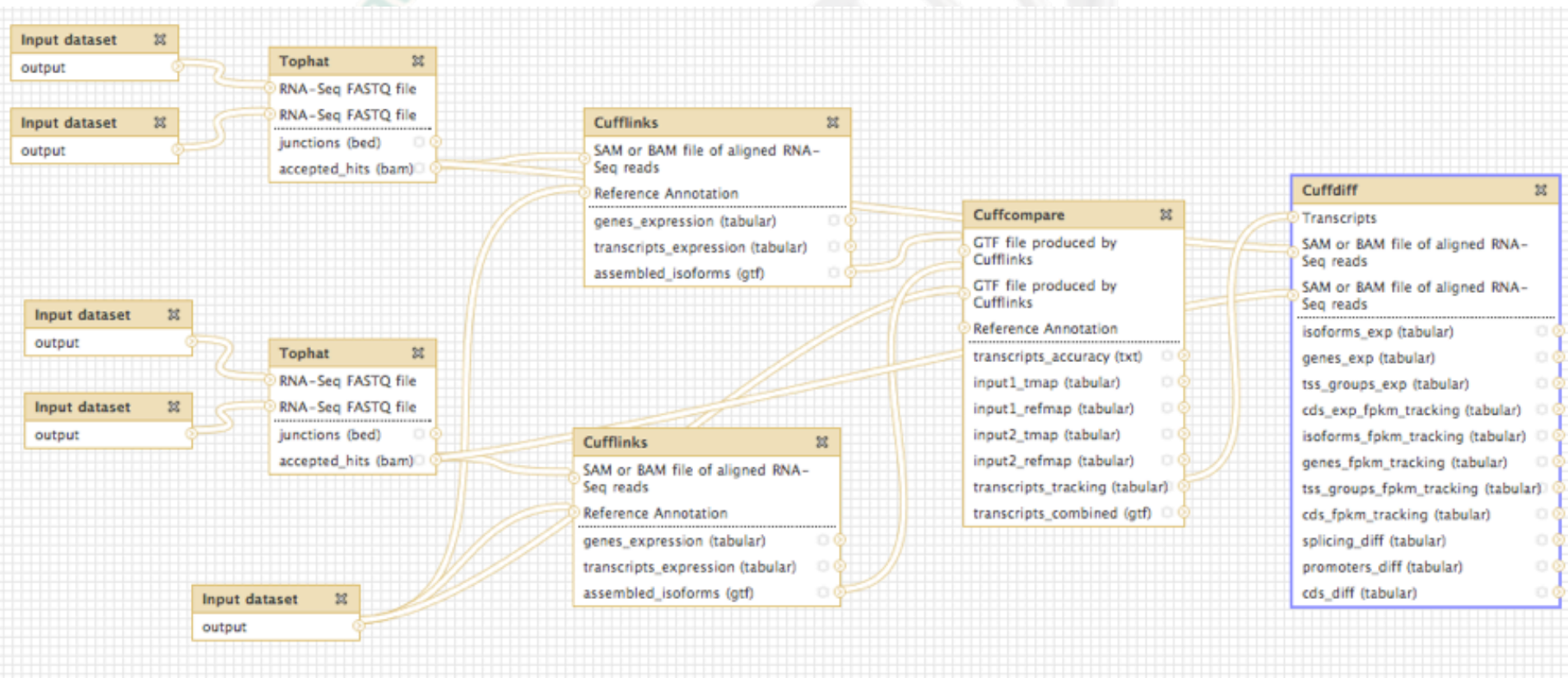
Datasources

- Upload file from your computer
 - FTP support for large datasets
- UCSC table browser
- UCSC Archaea table browser
- BX table browser
- EBI SRA
- BioMart
- Gramene Mart
- Flymine
- modENCODE fly server
- modENCODE modMine
- Ratmine
- YeastMine
- modENCODE worm server
- WormBase
- EuPathDB server
- EncodeDB at NHGRI
- EpiGRAPH server
- GenomeSpace import

Tool Suites

- Generic Tools
 - Text Manipulation
 - Format Converters
 - FASTA Manipulation
 - Filtering and Sorting
 - Join, Subtract, Group
 - Sequence Tools
 - Multi-species Alignment Tools
 - Genomic Interval Operations
 - Summary Statistics
 - Graphing/ Plotting
 - And More!
- NGS
 - QC and manipulation
 - Mapping
 - SAM Tools
 - GATK Tools (beta)
 - Variant Detection
 - Indel Analysis
 - Peak Calling
 - RNA Analysis
 - Picard (beta)
 - BEDTools
 - snpEff

Create Workflows



Sharing and Publishing

Sharing and Publishing History 'Variant Analysis for Sample E18'

Making History Accessible via Link and Publishing It

This history is currently restricted so that only you and the users listed below can access it. You can:

Make History Accessible via Link

Generates a web link that you can share with other people so that they can view and import the history.

Make History Accessible and Publish

Makes the history accessible via link (see above) and publishes the history to Galaxy's [Published Histories](#) section, where it is publicly listed and searchable.

Sharing History with Specific Users

You have not shared this history with any users.

Share with a user

[Back to Histories List](#)

Where you can use and build Galaxy

- Penn State's instance: <https://usegalaxy.org/>
- UAB's instance: <https://www.uab.edu/galaxy>
- Local instance: <http://getgalaxy.org>
 - Galaxy is designed for local installation and customization
 - Just download and run, completely self-contained
 - Easily integrate new tools
 - Easy to deploy and manage on nearly any (Unix) system
 - Run jobs on existing compute clusters
- On the cloud: <http://usegalaxy.org/cloud>
- Tool shed/contributing tools:
<http://toolshed.g2.bx.psu.edu/>

Tool Shed

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

Galaxy Tool Shed Repositories Groups Help User

3598 valid tools on Oct 07, 2015

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

search repository name, description

Name	Description	Repositories
Assembly	Tools for working with assemblies	59
ChIP-seq	Tools for analyzing and manipulating ChIP-seq data.	15
Combinatorial Selections	Tools for combinatorial selection	7
Computational chemistry	Tools for use in computational chemistry	21
Convert Formats	Tools for converting data formats	42
Data Managers	Utilities for Managing Galaxy's built-in data cache	19
Data Source	Tools for retrieving data from external data sources	30
Fasta Manipulation	Tools for manipulating fasta data	67
Fastq Manipulation	Tools for manipulating fastq data	42
Genome-Wide Association Study	Utilities to support Genome-wide association studies	8
Genomic Interval Operations	Tools for operating on genomic intervals	40
Graphics	Tools producing images	40
Imaging	Utilities to support imaging	2
Metabolomics	Tools for use in the study of Metabolomics	20
Metagenomics	Tools enabling the study of metagenomes	44
Micro-array Analysis	Tools for performing micro-array analysis	8
Next Gen Mappers	Tools for the analysis and handling of Next Gen sequencing data	77
Ontology Manipulation	Tools for manipulating ontologies	11
Phylogenetics	Tools for performing phylogenetic analysis	11
Proteomics	Tools enabling the study of proteins	62
RNA	Utilities for RNA	67
SAM	Tools for manipulating alignments in the SAM format	68
Sequence Analysis	Tools for performing Protein and DNA/RNA analysis	334
Statistics	Tools for generating statistics	63
Systems Biology	Systems biology tools	10
Text Manipulation	Tools for manipulating data	56
Tool Dependency Packages	Repositories that contain third-party tool dependency package installation definitions	448
Tool Generators	Tools that make or help make new tools	6
Transcriptomics	Tools for use in the study of Transcriptomics.	28
Variant Analysis	Tools for single nucleotide polymorphism data such as WGA	150
Visualization	Tools for visualizing data	47
Web Services	Tools enabling access to web services	13

Contents

- Setting up your project
 - Analysis Challenges
 - Tools available for analysis
 - **NGS**
 - Bioinformatics Resources
- What is Galaxy
 - **What isn't Galaxy**
 - FASTQ anatomy
 - Using Galaxy
- 

What isn't Galaxy

- Latest version of tools not always available (unless your willing to modify the wrapper for them)
- Not all options for tools are available
 - Examples:
 - TopHat unaligned reads file is not kept
 - Log files not kept
- Your favorite tool isn't there (need to write a wrapper to install it)
- Still buggy (although getting better with each new release!)
 - Example:
 - Job states is complete (by green colored box), but downstream tools can't use it because it didn't completely write all the file.
- Reproducible?

Solution? Blinking Cursor! |

Contents

- Setting up your project
 - Analysis Challenges
 - Tools available for analysis
 - **NGS**
 - Bioinformatics Resources
- What is Galaxy
 - What isn't Galaxy
 - **FASTQ anatomy**
 - Using Galaxy
- 

Heflin Center for Genomic Science

Contents

- Setting up your project
 - Analysis Challenges
 - Tools available for analysis
 - **NGS**
 - Bioinformatics Resources
- What is Galaxy
 - What isn't Galaxy
 - FASTQ anatomy
 - **Using Galaxy**
- 

Heflin Center for Genomic Science

Galaxy Splash Page

<https://www.uab.edu/galaxy>

<https://usegalaxy.org/>

The screenshot shows the Galaxy / UAB splash page. At the top, there is a navigation bar with the following items: **Galaxy / UAB**, **Analyze Data**, **Workflow**, **Shared Data**, **Visualization**, **Admin**, **Help**, **User**, and **Using 2.1 TB**. Below the navigation bar, the page is divided into three main sections: **Tools**, **History**, and a central content area.

Tools (left sidebar):

- search tools
- [Get Data](#)
- [Send Data](#)
- [Demo Tools](#)
- [ENCODE Tools](#)
- [Lift-Over](#)
- [Text Manipulation](#)
- [Filter and Sort](#)
- [Join, Subtract and Group](#)
- [Convert Formats](#)
- [Extract Features](#)
- [Fetch Sequences](#)
- [Get Genomic Scores](#)
- [Operate on Genomic Intervals](#)
- [Statistics](#)
- [Wavelet Analysis](#)
- [Graph/Display Data](#)
- [Regional Variation](#)
- [Multiple regression](#)
- [Multivariate Analysis](#)
- [Evolution](#)
- [Motif Tools](#)
- [Multiple Alignments](#)
- [Metagenomic analyses](#)
- [FASTA manipulation](#)
- [NCBI BLAST+](#)
- NGS TOOLBOX BETA
- [NGS: QC and manipulation](#)
- [NGS: Assembly](#)
- [NGS: Mapping](#)
- [NGS: Indel Analysis](#)
- [NGS: RNA Analysis](#)
- [NGS: SAM Tools](#)
- [NGS: HA GSL Tools](#)
- [NGS: Peak Calling](#)
- [SNP/WGA: Data; Filters](#)
- [SNP/WGA: QC; LD; Plots](#)
- [SNP/WGA: Statistical Models](#)
- [Human Genome Variation](#)
- [SnpEff tools](#)
- [VCF Tools](#)
- [DebugTools](#)

History (right sidebar):

- Unnamed history
- 0 bytes
- Your history is empty. Click 'Get Data' on the left pane to start

Galaxy is back on-line!! (Yellow alert box):

- Galaxy came back online Thursday May 20th.
- However, importing files via `/scratch/imports/galaxy/BLAZERID/` is working **only** for **unzip'ed files**. Importing **gzip'ed files (.gz)** will generate errors.
- See [wiki page](#) for more details.

Welcome to UAB Galaxy! (Green welcome box):

Where all you need is a `BlazerId` and a web browser to run NGS analyses on the UAB Cheaha Cluster!

Local Resources

UAB Galaxy Wiki: [Overview](#), [Data Import](#)
UAB Mailing Lists

- [UAB Galaxy-users](#) (search archive; [subscribe](#)) discuss with other UAB users
- [UAB Galaxy-help](#) ask the UAB admins for help!

UAB Cheaha Computing Cluster

- Cluster Hardware ([wiki](#))
- Request a [cheaha account](#) (needed only for command-line access and bulk data upload)

Internet Resources

[Learn Galaxy](#) - tutorials
[Galaxy Project](#) user mailing list ([searchable archives](#); [subscribe](#); [post](#))
[Galaxy Toolshed](#) plug-ins for additional tools that you can request for installation at UAB Public Galaxy Server at Penn State (PSU): [UseGalaxy.org](#) (more tools, but small disk quotas)

Brought to you by

- UAB IT Research Computing under the Office of the Vice President for Information Technology at UAB
- UAB [CCIS](#) (Center for Clinical and Translational Science under grant UL1 RR025777 from the NIH National Center for Research Resources)
- The [Galaxy Platform](#) is developed by Penn State and Emory University

Live Quickies

- Advanced fastQ manipulation: Galactic quickie # 14
- 454 Mapping: Single End Galactic quickie # 15
- Uploading Data using FTP Galactic quickie # 17
- Managing account histories Galactic quickie # 19

Random Galaxy icons/colors

Colors

Four Galaxy job cards are shown, each with a different background color and icon:

- Grey:** ControlR1FastQC data 4.html (Queued)
- Yellow:** ControlR1FastQC data 4.html (Running)
- Green:** ControlR1FastQC data 4.html (Completed)
- Red:** MF2-3: Cuffmerge on data 42, data 15, and data 46: merged transcripts (Failed)

Queued
Running
Completed
Failed

Download/Save

14: Control Tophat for Illumina on data 3, data 4, and data 2: accepted hits
19.4 Mb
format: bam, database: ?
Info: Settings:
Output files:
"/mnt/galaxyData/tmp/15.1.all.q/t
mpv7F3_v/dataset_12.*.ebwt"
Line rate: 6 (line is 64 bytes)
Lines per side: 1 (side is 64
bytes)
Offset rate: 5 (one in 32)
FTable chars: 10
Strings: unpacked
Max bucket size: d
display in IGB Local Web
Binary bam alignments file

14: Control Tophat for Illumina on data 3, data 4, and data 2: accepted hits
19.4 Mb
format: bam, database: ?
Info: Settings:
Output files:
"/mnt/galaxyData/tmp/15.1.all.q/t
mpv7F3_v/dataset_12.*.ebwt"
Line rate: 6 (line is 64 bytes)
Lines per side: 1 (side is 64
bytes)
Offset rate: 5 (one in 32)
FTable chars: 10
Strings: unpacked
Max bucket size: d
Download Dataset
ADDITIONAL FILES
Download bam_index

Icons

Icons and their functions:

- Display data in browser
- Edit attributes
- Delete
- Edit dataset annotation
- View details
- Run this job again
- View in Trackster
- Edit dataset tags

Edit files in History

Edit Attributes

Name:
Tophat for Illumina on data 3, data *
*

Info:
Settings:
Output files:

Annotation / Notes:
None

Add an annotation or notes to a dataset; annotations are available when a history is viewed.

Database/Build:
Click to Search or Select

Number of comment lines:

Chrom column:
1

Start column:
2

End column:
3

Strand column (click box & select):
 1

Name/Identifier column (click box & select):
 1

Score column for visualization:
1
2
3

*

This will inspect the dataset and attempt to correct the above column values if they are not accurate.

11: Tophat for Illumina on data 3, data 4, and data 2: insertions

11: Control Tophat for Illumina on data 3, data 4, and data 2: insertions

Contents

- Setting up your project
- Analysis Challenges
- Tools available for analysis
- NGS
- **Bioinformatics Resources**



Bioinformatics Resources



Facilities to help



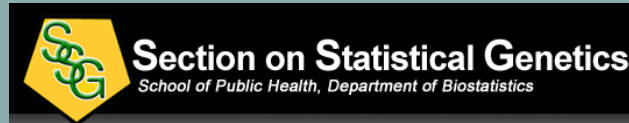
Informatics
Institute

Other
Genomic
Centers

You are
here

Department
of Pathology
Division of
Informatics

CCC Biostatistics
& Bioinformatics
Shared Facility
(BBSF)



Bioinformatics Resources

- **Heflin Center**
 - David Crossman, Ph.D.
 - dkcrossm@uab.edu
 - (205) 996-4045
- **CCTS-BMI**
 - Elliot Lefkowitz, Ph.D.
 - ElliotL@uab.edu
 - (205) 934-1946
- **Section on Statistical Genetics** (School of Public Health)
 - Hemant Tiwari, Ph.D.
 - Htiwari@soph.uab.edu
 - (205) 934-4907
- **Informatics Institute**
 - James Cimino, M.D.
 - ciminoj@uab.edu
 - (205) 996-1958
- **Department of Pathology Division of Informatics**
 - X. Long Zheng, M.D., Ph.D. (Interim)
 - zhengl@uab.edu
 - (205) 975-8161
- **Comprehensive Cancer Center (CCC) Biostatistics and Bioinformatics Shared Facility (BBSF)**
 - Karan Singh, Ph.D.
 - kpsingh@uab.edu
 - (205) 996-6122

Thanks! Questions?

Contact info:

David K. Crossman, Ph.D.

Bioinformatics Director

Heflin Center for Genomic Science

University of Alabama at Birmingham

<http://www.heflingenetics.uab.edu>

dkcrossm@uab.edu