# Introduction to Galaxy

North Carolina State University
Raleigh, NC
July 22, 2013

Dave Clements, Emory University
http://galaxyproject.org/

# Agenda

| | |
|---|---|
| 9:00 | Welcome and Intro |
| 9:30 | Basic Analysis with Galaxy |
| 10:50 | Break |
| 11:10 | Basic Analysis into Reusable Workflows |
| 11:40 | RNA-Seq Example Part I |
| 12:40 | Lunch |
| 1:40 | RNA-Seq Example Part II |
| 2:30 | Galaxy Community |
| 2:40 | Break |
| 3:00 | Sharing, Publishing and Reproducibility |
| 3:20 | Setting up Galaxy on the Amazon Cloud |
| 4:35 | Done |

# Introductions

In 40 seconds or less tell us
- your name
- your affiliation(s)
- something about your research
- something about your goals for today

# Goals

1. Introduce Galaxy
2. Introduce bioinformatics concepts and formats
3. Hands-on experience
   - Load and integrate data
   - Perform bioinformatic analysis with Galaxy
   - Save, share describe and publish your analyses
   - Visualize your results
   - Set up a Galaxy server in the cloud

This workshop will not cover details of how tools are implemented, or new algorithm designs, or which assembler or mapper or … is best for you.

# What is Galaxy?

- **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

- These options result in several **ways to use Galaxy**

http://galaxyproject.org

# Galaxy is available ...

As a free (for everyone) web service

http://usegalaxy.org

However, *a centralized solution cannot scale to meet the analysis needs of the entire world.*

# Galaxy is available ...

- As a free (for everyone) web service

  http://usegalaxy.org

- **As open source software**

  **http://getgalaxy.org**

# As Open Source Software: Local Galaxy Instances

- Galaxy is designed for local installation and customization

    - Easily integrate new tools

    - Run jobs on existing compute clusters

- Requires a computational resource on which to be deployed

**http://getgalaxy.org**

# As Open Source Software: Local Galaxy Instances

- Galaxy is designed for local installation and customization

  - Easily integrate new tools

  - Easy to deploy and manage on nearly any (unix) system

  - Run jobs on existing compute clusters

- Requires a **computational resource** on which to be deployed

**http://getgalaxy.org**

# Got your own cluster?

- Galaxy works with any DRMAA compliant cluster job scheduler (which is most of them).

- Galaxy is just another client to your scheduler.

# Galaxy is available ...

- As a free (for everyone) web service

  **http://usegalaxy.org**

- As open source software

  **http://getgalaxy.org**



- *On the Cloud*

  **http://usegalaxy.org/cloud**

  We are using this right now, and we will demonstrate how to do this later today

  **http://aws.amazon.com/education**

# Galaxy is available ...

- **As a free (for everyone) web service**

- **As open source software**

- **On the Cloud**



- ***With Commercial Support***

  A ready-to-use appliance (BioTeam)

  Cloud-based solutions (Appistry, ABgenomica, AIS)

  Consulting & Customization (Arctix, Deena Bioinformatics)

# Agenda

# Basic Analysis: We have

an assembly of an archaeal organism
gene annotation
TF binding sites

## Which genes have most overlapping TFBs?

**http://cloud1.galaxyproject.org/**
**http://cloud2.galaxyproject.org/**
**http://cloud3.galaxyproject.org/**

**(~ http://usegalaxy.org/galaxy101 )**

# Genes & TFBs: A General Plan

- Get some data
  - Sequence, genes/exons, TFBs
- Mess with it
  - Identify which genes/exons have TFBs
  - Count TFBs per exon
  - Visualize, save, download, ... exons with most TFBs

**http://cloud1.galaxyproject.org/**
**http://cloud2.galaxyproject.org/**
**http://cloud3.galaxyproject.org/**

**(~ http://usegalaxy.org/galaxy101 )**

**Exons**



**TFBs**

Exons

TFBs

Exons

TFBs

Overlap pairings

**Exons**

**TFBS**

**Exons**

**TFBS**

**Overlap pairings**

| | |
|---|---|
| ▬ | I |
| ▬ | I |
| ▬ | 2 |

**Exon overlap counts**

**Exon overlap counts**

**Exons**

**Exon overlap counts**

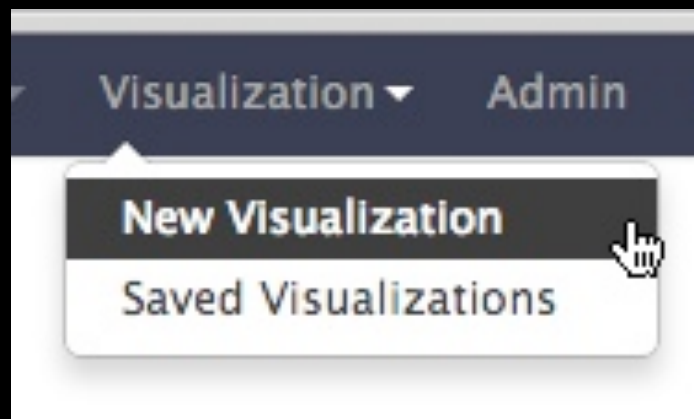**Exons**

**Join on exon name**
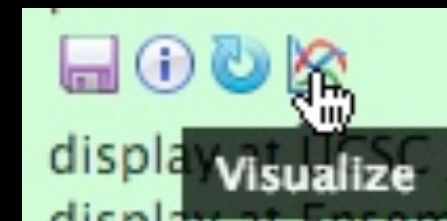
Exon overlap counts

Exons

Join on exon name

Rearrange columns w/ cut

# Visualize results



**or**

# Agenda

9:00   Welcome and Intro

9:30   Basic Analysis with Galaxy

10:50  Break

11:10  Basic Analysis into Reusable Workflows

11:40  RNA-Seq Example Part I

12:40  Lunch

1:40   RNA-Seq Example Part II

2:30   Galaxy Community

2:40   Break

3:00   Sharing, Publishing and Reproducibility

3:20   Setting up Galaxy on the Amazon Cloud

4:35   Done

# Agenda

9:00   Welcome and Intro

9:30   Basic Analysis with Galaxy

10:50  Break

11:10  Basic Analysis into Reusable Workflows

11:40  RNA-Seq Example Part I

12:40  Lunch

1:40   RNA-Seq Example Part II

2:30   Galaxy Community

2:40   Break

3:00   Sharing, Publishing and Reproducibility

3:20   Setting up Galaxy on the Amazon Cloud

4:35   Done

# Some Galaxy Terminology

**Dataset:**

Any input, output or intermediate set of data + metadata

**History:**

A series of inputs, analysis steps, intermediate datasets, and outputs

**Workflow:**

A series of analysis steps
Can be repeated with different data

# Exons and TFBs *History* → Reusable *Workflow?*

- The analysis we just finished was about
  - An archaea
  - Overlap between exons and TFBs
- But, ...
  - there is nothing inherently in the analysis about archaea, exons or TFBs
  - It is a series of steps that sets the score of one set of features to the number of overlaps from another set of features.

# Create a generic *Overlap* Workflow

**Extract Workflow from history**

Create a workflow from this history.
Edit it to make some things clearer.

**Run / test it**

Guided: rerun with same inputs

Did that work?

On your own:

Count # of exons in each TFBS
Did that work?  *Why not?*
Edit workflow: doc assumptions

# Agenda

# RNA-seq Exercise

Shared Data → Published Pages

→ RNA-Seq Analysis Exercise

# RNA-seq Exercise: A Plan

- Get input datasets; hg19, will mostly map to chr19

- Look at quality

- Trim as we see fit.

- Map the reads to the human reference using Tophat

- Run Cufflinks on Tophat output to assemble reads into transcripts

- Visualize it

# RNA-seq Exercise: A Plan

- Get input datasets; hg19, will mostly map to chr19

  - All datasets are FASTQ and from the Body Map 2.0 project

    - Shared Data → Data Libraries → RNA-Seq Example

# What is **FASTQ**?

- Specifies sequence (FASTA) and quality scores (PHRED)

- Text format, 4 lines per entry

```
@SEQ_ID
GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
!''*((((***+))%%%++)(%%%).1***-+*''))**55CCF>>>>>>CCCCCCC65
```

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

```
 SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS
 ...............................IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII
 .........................XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX
 !"#$%&'()*+,-./0123456789:;<=>?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
 |                              |   |        |                                       |          |
 33                             59  64       73                                      104        126

S - Sanger        Phred+33,  93 values  (0, 93) (0 to 60 expected in raw reads)
I - Illumina 1.3 Phred+64,  62 values  (0, 62) (0 to 40 expected in raw reads)
X - Solexa        Solexa+64, 67 values (-5, 62) (-5 to 40 expected in raw reads)
```

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

# RNA-seq Exercise: A Plan

**Look at quality Options 1 & 2:**

1. NGS QC and Manipulation → Compute Quality Statistics

   NGS QC and Manipulation → Draw quality score boxplot

   No control over how it is calculated or presented.

2. NGS QC and Manipulation → FastQ Summary Statistics,

   Graph / Display Data → Boxplot of quality statistics

   Lots of control over what the box plot looks like,
   Statistics in text and graphic formats

# RNA-seq Exercise: A Plan

- Get input datasets; hg19, will mostly map to chr19

- Look at quality: Option 3

  - NGS QC and Manipulation → **FastQC**

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

http://bit.ly/FastQCBoxPlot

# RNA-seq Exercise: A Plan

- Look at quality

- 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

# RNA-seq Exercise: A Plan

- Look at quality

- ~~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.

# RNA-seq Exercise: A Plan

- Look at quality

- 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

+

Option 2

# 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?

  - http://biostars.org/

  - http://seqanswers.com/

  - http://galaxyproject.org/search

# RNA-seq Exercise: A Plan

- Get input datasets; hg19, will mostly map to chr19

- Look at quality

- Trim as we see fit.

- Map the reads to the human reference using 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.*

# Agenda

| | |
|---|---|
| 9:00 | Welcome and Intro |
| 9:30 | Basic Analysis with Galaxy |
| 10:50 | Break |
| 11:10 | Basic Analysis into Reusable Workflows |
| 11:40 | RNA-Seq Example Part I |
| 12:40 | Lunch |
| 1:40 | RNA-Seq Example Part II |
| 2:30 | Galaxy Community |
| 2:40 | Break |
| 3:00 | Sharing, Publishing and Reproducibility |
| 3:20 | Setting up Galaxy on the Amazon Cloud |
| 4:35 | Done |

# Agenda

9:00   Welcome and Intro
9:30   Basic Analysis with Galaxy
10:50  Break
11:10  Basic Analysis into Reusable Workflows
11:40  RNA-Seq Example Part I
12:40  Lunch
1:40   RNA-Seq Example Part II
2:30   Galaxy Community
2:40   Break
3:00   Sharing, Publishing and Reproducibility
3:20   Setting up Galaxy on the Amazon Cloud
4:35   Done

# RNA-seq Exercise: A Plan

- ...

- Map the reads to the human reference using Tophat

- Run Cufflinks on Tophat output to assemble reads into transcripts

  - Tophat does not make any predictions about how the reads it mapped, assemble together into transcripts.

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

# RNA-seq Exercise: A Plan

- ...

- Map the reads to the human reference using Tophat

- Run Cufflinks on Tophat output to assemble reads into transcripts

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

  - Visualize it

# Visualizing Genomics

## Supported external browsers

- UCSC

- Ensembl

- GBrowse

- IGB

- IGV

## Traditional browser strengths:

- Showing what is nearby

- what else is happening here

- highlighting correlations

- integrating many datasets

# Trackster: Galaxy's embedded track browser

# Circster

# Create a visualization in Galaxy



**or**

# Vizualization inside Galaxy

- Levarge visualization to evaluate and refine analyses

- Make the *analyze-visualize-refine* loop seamless and fast

- Enable experimenting with tools and their parameter space

- Support custom genome browsers

# Agenda

9:00   Welcome and Intro
9:30   Basic Analysis with Galaxy
10:50  Break
11:10  Basic Analysis into Reusable Workflows
11:40  RNA-Seq Example Part I
12:40  Lunch
1:40   RNA-Seq Example Part II
2:30   Galaxy Community
2:40   Break
3:00   Sharing, Publishing and Reproducibility
3:20   Setting up Galaxy on the Amazon Cloud
4:35   Done

# Encourage Local Galaxy Instances

- Encourage and support Local Galaxy Instances

  - Support increasingly decentralized model and improve access to existing resources

  - Focus on building infrastructure to enable the community to integrate and share tools, workflows, and best practices

**Galaxy Tool Shed**
http://toolshed.g2.bx.psu.edu

# Encourage **Public** Galaxy Instances
http://bit.ly/gxyServers

## Interested in:

ChIP-chip and ChIP-seq?
✓ Cistrome

Statistical Analysis?
✓ Genomic Hyperbrowser

Protein synthesis?
✓ GWIPS-viz

*de novo* assembly?
✓ CBIIT Galaxy

Reasoning with ontologies?
✓ OPPL Galaxy

Repeats!
✓ RepeatExplorer

Everything?
✓ Andromeda

Plus many more

# Galaxy Resources and Community

Mailing Lists (very active)

Unified Search

Issues Board

Events Calendar, News Feed

Community Wiki

GalaxyAdmins

Screencasts

Tool Shed

Public Installs

CiteULike group

Annual Community Meting

http://wiki.galaxyproject.org

# Galaxy Resources and Community: Mailing Lists
### http://wiki.galaxyproject.org/MailingLists

## Galaxy-Announce

Project announcements, low volume, moderated
Low volume (    42 posts in 2012,  2100+ members)

## Galaxy-User

Questions about using Galaxy and usegalaxy.org
High volume (2900 posts in 2012,  2700+ members)

## Galaxy-Dev

Questions about developing for and deploying Galaxy
High volume (4500 posts in 2012,   900+ members)

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



**Galaxy Web Search**

Google™ Custom Search                    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**

chip-seq

All   Tools   Email   Source code   Shared   Documentation   Abstracts   Requests

About 444 results (0.06 seconds)

Galaxy | Accessible Page | ChIP-seq exercise

*Find*

Everything on …

Tools for …

Email about …

Source code for …

Published Histories, Pages, Workflows, about …

Documentation on …

Papers using Galaxy for …

Related feature requests

# Community can create, vote and comment on issues



## http://bit.ly/gxyissues

# http://wiki.galaxyproject.org

# Events

# News

## Galaxy Wiki

DaveClements  Settings  Logout

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 deadlines that are relevant to 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 Event Calendar, please add it here or send it to ✉ outreach@glaxyproject.org .

### Upcoming Events

Research 2013 Triangle Galaxy Workshop Tour

XSEDE — Extreme Science and Engineering Discovery Environment

| Date | Topic/Event | Venue/Location |
|---|---|---|
| July 18-23 | *Introduction to Galaxy Workshop* — National Institute of Environmental Health Sciences (NIEHS) | |
| | *Introduction to Galaxy Workshop* — University of North Carolina, Chapel Hill | 2013 Research Triangle Workshop Tour, North Carolina, United |
| | *Galaxy Installation Tutorial* **2013 GMOD Summer School** | |
| | *Introduction to Galaxy Workshop* — North Carolina State University | |
| July 19-23 | **ISMB/ECCB, BOSC and MS SIG 2013** — Talks, posters and workshops. Lots of them. | Berlin, Germany |
| July 21-25 | *Experiences in building a Next-Generation Sequencing Analysis Service using Galaxy, Globus Online, and Amazon Web Services* | XSEDE13, San Diego, Calif United States |
| | *A Sustainable National Gateway for Biological Computation* | |
| | *Supporting Genomics and other Biological Research* | |
| September 28 - October 1 | *Galaxy Workshop* | The Genomic Bioinformatics Workshop, Sydney, Australia |
| October 1-3 | *Galaxy* | **Beyond the Genome 2013**, Francisco, California, United |
| October 7-8 | TBD | NGS & Bioinformatics Sum Europe |
| | *Using Galaxy to Provide a NGS Analysis Platform* | |
| October 9-11 | *Galaxy Training Days* | GenoToul bioinformatics facil INRA, Toulouse Auzeville, Fra |
| October 22-26 | *High Throughput Data Analysis and Visualization with Galaxy* | ASHG 2013, Boston, Massach United States |
| November 6-12 | *Computational and Comparative Genomics Course* — Application Deadline: July 15, 2013 | Cold Spring Harbor Laborator York, United States |

## Galaxy Wiki

DaveClements  Settings  Logout  |  Search:

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

#### News Items

New CloudMan Release
SlipStream Appliance: Galaxy Edition
July 2013 Galaxy Update
1000th Galaxy CiteULike Paper
GCC2013 Registration Ends 14 June
June 3, 2013 Galaxy Distribution
June 2013 Galaxy Update
Software Carpentry Boot Camp: Oslo
GCC2013 Early Registration Ends 24 May
Duplicate Accounts on Main

*News Archive*

### News Items

#### New CloudMan Release

**We just released an update to Galaxy CloudMan.** CloudMan offers an easy way to get a personal and completely functional instance of Galaxy in the cloud in just a few minutes, without any manual configuration.

**IMPORTANT - please read**

Any new cluster will automatically start using this version of CloudMan. Existing clusters will be given an option to do an automatic update once the main interface page is refreshed. Note that this upgrade is a major version upgrade and thus the migration is rather complicated. The migration process has been automated but will take a little while to complete. If you have made customizations to your cluster in terms of adding file systems, upgrading the database, or similar, we do not recommend you perform the upgrade. Note that this upgrade comes with (and requires) a new AMI (ami-118bfc78), which will automatically be used when starting an instance via CloudLaunch.

**This update brings a large number of updates and new features, the most prominent ones being:**

- Unification of galaxyTools and galaxyData file systems into a single galaxy filesystem. This change makes it possible to utilize the Galaxy Tool Shed when installing tools into Galaxy.
- Added initial support for Hadoop-type workloads
- Added initial support for cluster federation via HTCondor
- Added a new file system service for an instance's transient storage, allowing it to be used across the cluster over NFS
- Added a service for the Galaxy Reports webapp
- Added optional Loggly based off-site logging support
- Added tags to all resources utilized by CloudMan

For more details on the new features, see the the CHANGELOG and for even more details see, *all 291 commit messages from 7 contributors.*

Enjoy and please let us know what you think,

Enis Afgan

Posted to the *Galaxy News* on 2013-07-08

#### SlipStream Appliance: Galaxy Edition

http://bit.ly/gcc2014

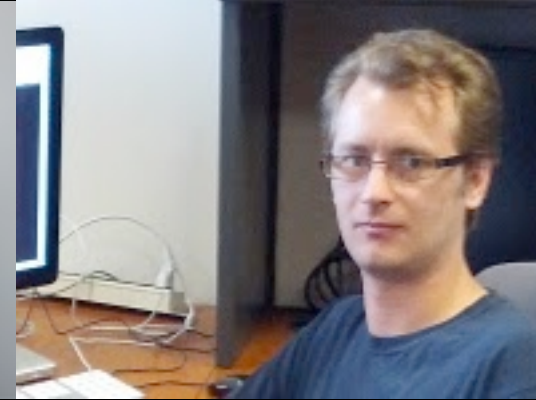# The Galaxy Team



Enis Afgan    Dannon Baker    Dan Blankenberg    Dave Bouvier    Dave Clements
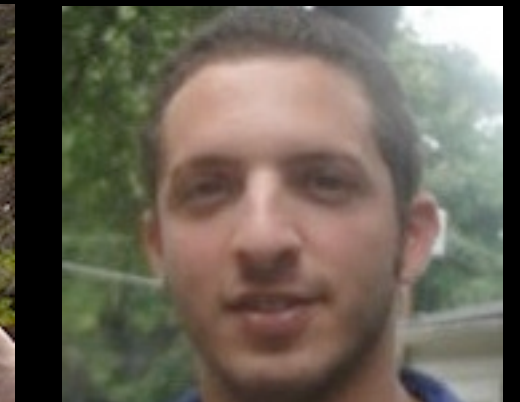
Nate Coraor    Carl Eberhard    Dorine Francheteau    Jeremy Goecks    Sam Guerler

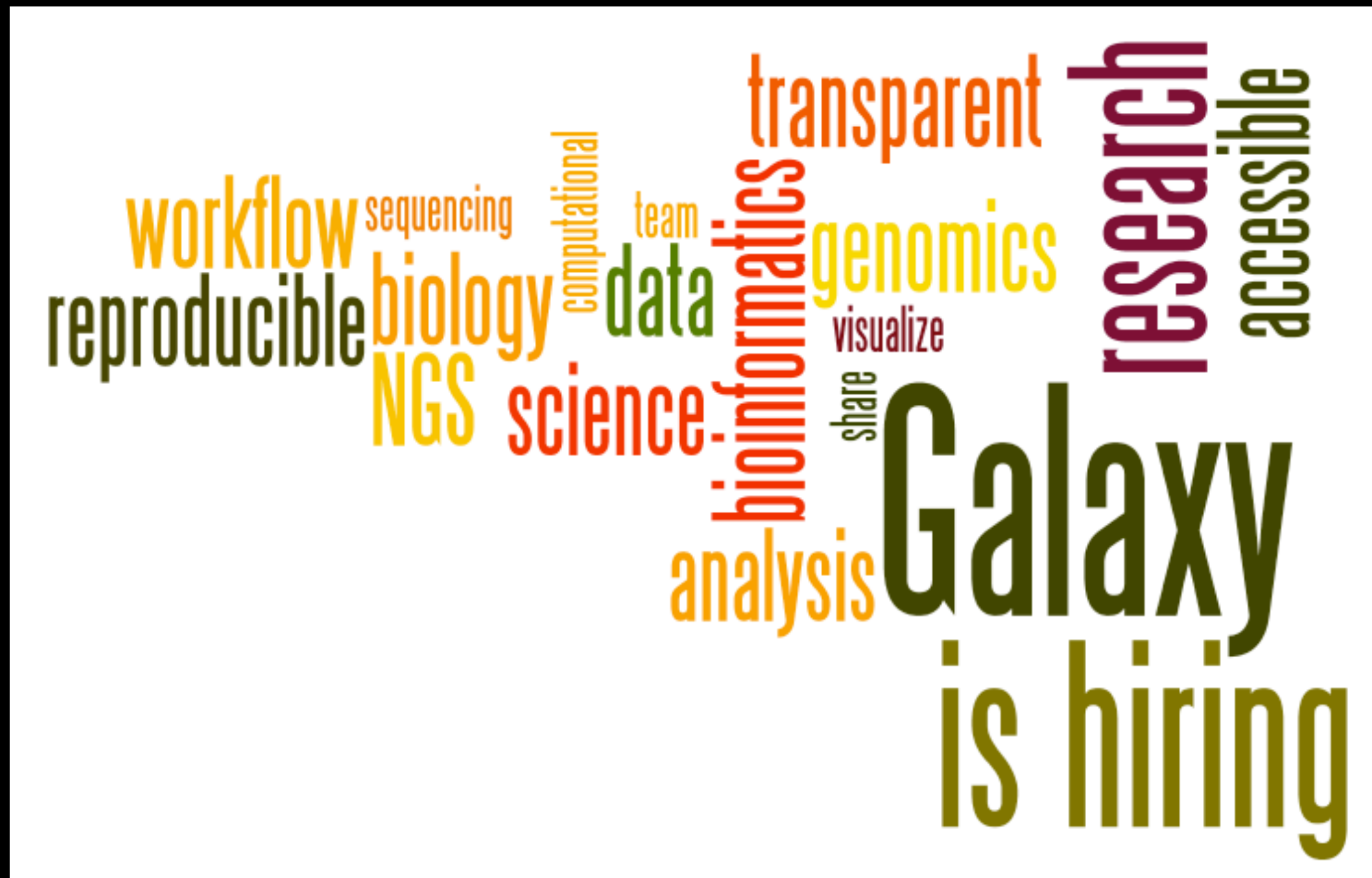Jen Jackson    Greg von Kuster    Ross Lazarus    Anton Nekrutenko    James Taylor

http://wiki.galaxyproject.org/GalaxyTeam

# Galaxy is hiring post-docs and software engineers



Please help.

http://wiki.galaxyproject.org/GalaxyIsHiring

# Agenda

9:00  Welcome and Intro
9:30  Basic Analysis with Galaxy
10:50  Break
11:10  Basic Analysis into Reusable Workflows
11:40  RNA-Seq Example Part I
12:40  Lunch
1:40  RNA-Seq Example Part II
2:30  Galaxy Community
2:40  Break
3:00  Sharing, Publishing and Reproducibility
3:20  Setting up Galaxy on the Amazon Cloud
4:35  Done

# Agenda

# More Galaxy Terminology

**Share:**
  Make something available to someone else

**Publish:**
  Make something available to everyone

**Galaxy Page:**
  Analysis documentation within Galaxy; easy to embed any Galaxy object

# Sharing & Publishing enables Reproducibility

Galaxy aims to push the goal of reproducibility from the bench to the bioinformatics realm

All analysis in Galaxy is recorded without any extra effort from the user.

**Histories, workflows, visualizations** and *pages* can be shared with others or published to the world.

# Sharing & Publishing enables Reproducibility

# Sharing & Publishing enables Reproducibility

# Windshield splatter analysis with the Galaxy metagenomic pipeline: A live supplement

SERGEI KOSAKOVSKY POND[1,2,*], SAMIR WADHAWAN[3,6*], FRANCESCA CHIAROMONTE[4], GURUPRASAD ANANDA[1,3], WEN-YU CHUNG[1,3,7], JAMES TAYLOR[1,5], ANTON NEKRUTENKO[1,3] and THE GALAXY TEAM[1*]

Correspondence should addressed to SKP, JT, or AN.

## How to use this document

This document is a live copy of supplementary materials for the manuscript. It provides access to the **exact** analyses and workflows discussed in the paper, so you can play with them by re-running, changing parameters, or even applying them to your own data. Specifically, we provide the two histories and one workflow found below. You can view these items by clicking on their name to expand them. You can also import these items into your Galaxy workspace and start using them; click on the green plus to import an item. To import workflows you must create a Galaxy account (unless you already have one) – a hassle-free procedure where you are only asked for a username and password.

This is the Galaxy history detailing the comparison of our pipeline to MEGAN:

**Galaxy History | Galaxy vs MEGAN**
Comparison of Galaxy vs. MEGAN pipeline.

This is the Galaxy history showing a generic analysis of metagenomic data. (This corresponds to the "A complete metagenomic pipeline" section of the manuscript and **Figure 3A**):

**Galaxy History | metagenomic analysis**

This is the Galaxy workflow for generic analysis of metagenomic data. (This corresponds to the "A complete metagenomic pipeline" section of the manuscript and **Figure 3B**):

**Galaxy Workflow | metagenomic analysis**
Generic workflow for performing a metagenomic analysis on NGS data.

## Accessing the Data

Windshield Splatter datasets analyzed in this manuscript can be accessed through this Galaxy Library. From

**Author**

aun1

**Related Pages**

All published pages
Published pages by aun1

**Rating**

Community
(6 ratings, 5.0 average)
★★★★★

**Tags**

Community:
paper   galaxy
megan

http://usegalaxy.org/u/aun1/p/windshield-splatter

# Sharing for Galaxy Administrators Too

**Data Libraries**
Make data easy to find

**Genome Builds**
Care about a particular subset of life?

**Galaxy Tool Shed**
Wrapping tools and datatypes

# Agenda

9:00   Welcome and Intro

9:30   Basic Analysis with Galaxy

10:50  Break

11:10  Basic Analysis into Reusable Workflows

11:40  RNA-Seq Example Part I

12:40  Lunch

1:40   RNA-Seq Example Part II

2:30   Galaxy Community

2:40   Break

3:00   Sharing, Publishing and Reproducibility

3:20   Setting up Galaxy on the Amazon Cloud

4:35   Done

# Galaxy CloudMan
## http://usegalaxy.org/cloud

- Start with a **fully configured and populated** (tools and data) Galaxy instance.

- Allows you to scale up and down your compute assets as needed.

- Someone else manages the data center.

- **We are using this today.**



- **You will set up an instance now**

## http://aws.amazon.com/education

# Could do this step by step, but ...
## http://bit.ly/GXYAWSGetStarted

# AWS Credentials

http://bit.ly/ncsuaws

# Instant CloudMan
## http://usegalaxy.org/cloudlaunch



# http://bit.ly/ncsuaws

# Agenda

| | |
|---|---|
| 9:00 | Welcome and Intro |
| 9:30 | Basic Analysis with Galaxy |
| 10:50 | Break |
| 11:10 | Basic Analysis into Reusable Workflows |
| 11:40 | RNA-Seq Example Part I |
| 12:40 | Lunch |
| 1:40 | RNA-Seq Example Part II |
| 2:30 | Galaxy Community |
| 2:40 | Break |
| 3:00 | Sharing, Publishing and Reproducibility |
| 3:20 | Setting up Galaxy on the Amazon Cloud |
| 4:35 | Done, Almost |

# Instant Feedback

**http://bit.ly/20130722Gxy**

# Acknowledgements

Trudy Mackay
You

Tom Randall
Barrie Hayes
The Galaxy Team

NCSU Dept of Genetics

AWS Education Grant

NIH  NSF  Huck Institute
Penn State University   Emory University

http://bit.ly/20130722Gxy

# Thanks

## http://bit.ly/20130722Gxy



### Dave Clements

### Galaxy Project
### Emory University

clements@galaxyproject.org

# RNA-Seq Example: Part III

- Run Cufflinks on Tophat output to assemble reads into transcripts

- Run Cuffdiff on Tophat output to find significant differences in expression.

  - *Imagine pages and pages of discussion on the intricacies and pitfalls of RNA-seq differential expression analysis here.*

http://bit.ly/parsonsrnaseq

# Cuffdiff

- Which Transcript definitions to use?

  - IGenomes

  - Adrenal or Brain from Cufflinks

  - Run Cuffmerge on Adrenal & Brain Cufflinks files

- Depends on what you care about.

  - I care about a timely workshop, so I'll use IGenomes.

# Cuffdiff

- Produces 11 output files, all explained in doc
- We'll focus on gene/transcript differential expression testing files (also care about gene/transcrpt FPKM files)
- Column 7 ("status") can be FAIL, NOTEST, LOWDATA or OK

  - Filter and Sort → Filter

    - c7 == 'OK' or C7 == 'LOWDATA'

- Column 14 ("significant") can be yes or no

  - c14 == 'yes'