# Adventures in Scaling Galaxy

https://speakerdeck.com/jxtx     @jxtx / #usegalaxy

…in which I will not talk about the ~~elephant~~ whale in the room…

**Galaxy's motivating questions**

How best can data intensive methods be **accessible** to scientists?

How best to facilitate **transparent communication** of computational analyses?

How best to ensure that analyses are **reproducible**\*?

*\*The state of which is **frighteningly** bad, see doi:10.1038/nrg3305, doi:10.7717/peerj.148*

# Galaxy: accessible analysis system

**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

**An open extensible platform** for sharing tools, datatypes, workflows, ...

Describe analysis tool
behavior abstractly

Analysis environment automatically
and transparently tracks details

Workflow system for complex analysis,
constructed explicitly or automatically

Pervasive sharing, and publication
of documents with integrated analysis

Visualization and visual analytics

Wait… what… a *free* web-service for high throughput sequence data analysis?!

**Galaxy as a Service**

Public web site that anyone can use for free

~1,200 new users, ~20 TB of user data uploaded, and ~180,000 analysis jobs *per month*

Since 2010, disk quotas (250Gb per user) and compute limits (4 concurrent analyses)

**http://usegalaxy.org**

Registered Users versus Jobs Submitted on Galaxy Main

# usegalaxy.org data growth

New Data per Month (TB)

+128 cores for NGS/multicore jobs

Data quotas implemented...

Nate Coraor

# usegalaxy.org frustration growth



Nate Coraor

Scaling plan one: Decentralize!

# Local Galaxy Deployment

Galaxy is designed for local installation and customization... just download and run

Pluggable interfaces to compute resources, easily connect to one or more existing clusters

Ideally, allow users to take advantage of whatever computational resources they already have access to.

# More than 60 known public Galaxy servers

Ballaxy for structure based computational biology,
Cistrome for regulatory sequence analysis,
Genomic Hyperbrowser: statistical integration of genomic data,
GigaGalaxy: integrating workflows published in GigaScience,
Pathogen Portal:comparative analysis of host response to pathogens,
...

# Dozens of large scale private Galaxy instances

**Galaxy / CoSSci**

Analyze Data　Workflow　Shared Data ▾　Visualization ▾　Help ▾　User ▾　　　Using 0 bytes

**Tools**

search tools

**COSSCI TOOLS**

DEf01f Dow Eff

DEf01d Dow Eff

DEf01c Dow Eff

DEf01 Dow Eff

EA Ethnographic Atlas

LRB Lewis R. Binford's forager data

SCCS Standard Cross-Cultural Sample

WNAI Western North American Indians

HPC Tools

Development Tools

**GALAXY TOOLS**

Text Manipulation

Filter and Sort

Join, Subtract and Group

Statistics

Wavelet Analysis

Graph/Display Data

Multiple regression

Multivariate Analysis

# Welcome to CoSSci

Begin an analysis by selecting one of the tools in the left-hand menu (one of EA, LRB, SCCS, or WNAI) and Execute it. You can explore these datasets further by modifying the variables used for the analyses. For more information about using this gateway, including how to select variables for the individual analyses, please visit the Visual Manual for CoSSci.

## How-tos and guides

**How to share histories.** The screencast shows how to share histories and how to access histories shared by other user.



Download the screencast directly

## Recorded presentations

Complex Social Science Gateway: High Performance Computing for Anthropology and the Social Sciences
Society for Applied Anthropology, 74th Annual Meeting, 2014

**History**

Unnamed history

0 bytes

ℹ Your history is empty. Click 'Get Data' on the left pane to start

https://hyperbrowser.uio.no/hb/

**The Genomic HyperBrowser** v1.6 (powered by Galaxy)

Analyze Data    Shared Data    Help    User    Using 0 bytes

## Tools        Options ▾

search tools

**HYPERBROWSER ANALYSIS**

**Statistical analysis of tracks**
- Analyze genomic tracks

**Visual analysis of tracks**
**Specialized analysis of tracks**
**Text-based analysis interface**

**HYPERBROWSER TRACK PROCESSING**

**HyperBrowser track repository**
**Customize tracks**
**Generate tracks**
**Format and convert tracks**
**GTrack tools**

**ARTICLE/DOMAIN-SPECIFIC TOOLS**

**The differential disease regulome**
**MCFDR**
**Monte Carlo null models**
**Transcription factor analysis**
**Gene tools**
**microRNA tools**

**HYPERBROWSER INTERNAL TOOLS**

**Admin of genomes and tracks**
**Development tools**
**Assorted tools**

**STANDARD GALAXY TOOLS**

**Get Data**

---

## The Genomic HyperBrowser

### If you have a *genomic track*, this is the place to analyze it!

To analyze a track, simply:

1. Click **Statistical analysis of tracks: Analyze genomic tracks** in the left-hand menu.
2. Select tracks from your Galaxy history of browse our collection.
   *(To load a track to your history, click **Get data: Upload file**)*
3. Select the analysis you are interested in:
   - any property of a single track
   - any relation between a pair of tracks

For help using the system:

1. Click **The Genomic Hyperbrowser: Help** in the left-hand menu.
2. Or, look through the following screencasts:
   *(further screencasts are available from the help menu)*

# CloudBased Image Analysis & Processing Toolbox

# Image Analysis and Processing *for everyone*.

The Cloud-based Image Analysis and Processing Toolbox project provides access to existing biomedical image processing and analysis tools via remote user-interface using the NeCTAR cloud.

## Use Toolbox

Use project's free server

## Demo

Watch other demos

## Project Blog

Project Blog

https://galaxy.cbio.mskcc.org

# Galaxy / Rätsch Lab Analyze Data

Workflow    Shared Data ▾    Visualization    Cloud ▾    Help ▾    User ▾        Using 0 bytes

## Tools

search tools

**Get Data**

### SEQUENCE ANALYSIS

**Toy Data**

**SVM Toolbox**

**KIRMES**

**Genomic Signals**

### GENE FINDING

**mGene.web (v0.2)**

**mGene.web modules (v0.2)**

**mGene.web modules (v0.4)**

### OQTANS (V0.1)

**Read Mapping**

**Transcript Prediction/Assembly**

**Differential/Quantitative Analysis**

**Enrichment Analysis (v0.1)**

**Read Alignment Filtering (v0.2)**

**GFF Toolkit (v0.1)**
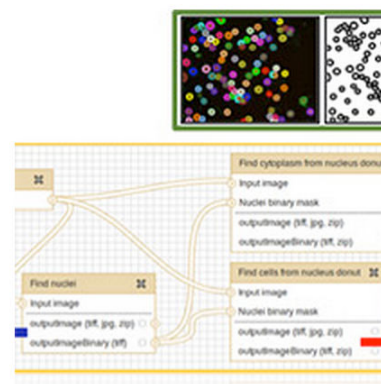
### GENETICS TOOLS

**SInBaD Tools**

**Multiple regression**

### NGS TOOLS

**NGS: QC and manipulation**

**NGS: Mapping**

---

oqtans

ℹ️ **Oqtans Moved to New Data Center**

We moved oqtans instance to our new data center in New Jersey on July 22. Resuming the normal operations.
**Rätsch Lab Galaxy Team**

## Galaxy/Rätsch Lab

# oqtans
## online quantitative transcriptome analysis

This is a customized version of the Galaxy framework, extended with machine learning based tools for sequence and tiling array data analysis. It provides tools developed by members of the Machine Learning in Biology (MLB) Group at cBio@MSKCC in New York City, USA. For problems with any of the non-standard tools, please contact the MLB Galaxy Support Team.

oqtans.org

## History

**Unnamed history**

0 bytes

ℹ️ This history is empty. You can load your own data or get data from an external source

rnaseq.pathogenportal.org

# RNA-Rocket

**Galaxy**

Launch Pad    Project View    Shared Data ▾    How–To    Help ▾    User ▾    Using 0 bytes

View a **list of supported genomes** from EuPathDB, PATRIC, and VectorBase.

Have a question? Contact the Pathogen Portal Team

TRIMMING

DEDUPLICATION

*FASTQ*   *FASTQ*

*BAM*

*FASTQ* → ALIGNMENT & MAPPING → *BAM* → TRANSCRIPT ASSEMBLY → *GTF,FPKM...* → DIFFERENTIAL EXPRESSION ANALYSIS → *LOG RATIOS, P–VALUES...*

READS QUALITY CHECK

MAPPING QUALITY CHECK

**Choose an activity below**

## ⬆ Uploads

From your computer or a URL

From ENA/SRA

## ☑ Quality Control

### Login to get started

Landmark or Region :
ECII_CH01:1..138,201    Search
Examples : AL590443:85000-115000, ECI_CH11:115000..135000.

Data Source
MicrosporidiaDB GBrowse v2.48

Annotate Restriction Sites
Save Snapshot    Load

Scroll/Zoom: « ‹ —

Overview    ECII_CH01
10k  20k  30k  40k  50k  60k  70k  80k

Region

Details    ECII_CH01: 138.2 kbp    50 kbp
30k  40k  50k  60k  70k  80k

**Stream results directly to a BRC**

# CloudMan: a general purpose deployment manager for ANY cloud

CloudMan

amazon
web services™

nectar

Google
Cloud Platform

iPlant
Collaborative™
Atmosphere

Enis Afgan, Dannon Baker

Enis Afgan, Dannon Baker

That was a great plan!

…but users still want one easy to use gateway

Scaling plan two: beg, borrow, steal!

Best place to build this robust entry point is clearly a national supercomputing center

The Texas Advanced Computing Center (TACC) has already built substantial infrastructure in the context of the iPlant project

(Including multi petabyte online storage, cloud infrastructure, collocated with some of the worlds largest HPC machines)

However, the iPlant and TACC cyber-infrastructure was underused; thus we established a collaboration

Since October 2013 Galaxy Main has run from TACC

usegalaxy.org frustration growth

Still not enough!

# Pulsar: Galaxy job runner that can run almost anywhere. No shared filesystem, stages all necessary Galaxy components



Galaxy Server VMs (TACC)

Galaxy Server Processes

Messaging Server

Job control (AMQP)

Data transfer (HTTPS)

Data transfer (HTTPS)

Blacklight (PSC)

Pulsar

Stampede (TACC)

Pulsar

John Chilton

"Big" NGS/Multicore Job

dynamic walltime — Average + std dev

Galaxy dedicated cluster

hit walltime? — No

Yes

Stampede

hit walltime? — Yes → Sorry ಠ⌒ಠ

No → Done! \(• ‿ •)/

Nate Coraor

Result: No waiting for jobs
to run on usegalaxy.org!

(for now…)

# Galaxy

An Ansible role for installing and managing Galaxy servers. Despite the name confusion, Galaxy bears no relation to Ansible Galaxy.

## Requirements

This role has the same dependencies as the `hg` module, namely, Mercurial. In addition, Python virtualenv is required (as is pip, but pip will automatically installed with virtualenv). These can easily be installed via a pre-task in the same play as this role:

```
- hosts: galaxyservers
    pre_tasks:
        - name: Install Mercurial
          apt: pkg={{ item }} state=installed
          sudo: yes
          when: ansible_os_family = 'Debian'
          with_items:
              - mercurial
              - python-virtualenv
        - name: Install Mercurial
          yum: pkg={{ item }} state=installed
          sudo: yes
          when: ansible_os_family = 'RedHat'
          with_items:
              - mercurial
              - python-virtualenv
    roles:
        - galaxy
```

Bringing it all together: automate all the things!
Unified **ansible** playbook for Galaxy main, cloud, and local deployments

# The Galaxy Team



Enis Afgan

Dannon Baker

Dan Blankenberg

Dave Bouvier

Marten Čech

John Chilton

Dave Clements

Nate Coraor

Carl Eberhard

Jeremy Goecks

Sam Guerler

Jen Jackson

Ross Lazarus

Anton Nekrutenko

James Taylor

http://wiki.galaxyproject.org/GalaxyTeam

# Computational Biology, Genomics, and Bioinformatics at Johns Hopkins University



**JOHNS HOPKINS**
WHITING SCHOOL
*of* ENGINEERING

**JOHNS HOPKINS**
SCHOOL *of* MEDICINE

**JOHNS HOPKINS**
BLOOMBERG SCHOOL
*of* PUBLIC HEALTH

**JOHNS HOPKINS**
KRIEGER SCHOOL
*of* ARTS & SCIENCES

## Biomedical Engineering

Joel Bader
Mike Beer
Rachel Karchin
**Steven Salzberg**

## Oncology

Elana Fertig
Luigi Marchionni
Robert Scharpf
Sarah Wheelan

## Biostatistics

Kasper Hansen
Hongkai Ji
Jeff Leek
Ingo Ruczinski
Cristian Tomasetti

## Biology

**James Taylor**

## Computer Science

Alexis Battle
**Ben Langmead**
Suchi Saria

## Medicine

Lilian Florea
**Mihaela Pertea**
Jiang Qian

## Applied Math

Don Geman

http://ccb.jhu.edu

## Tenure-Track Faculty Position in Data Intensive Biology

The Department of **Biology** seeks to hire a tenure-track Assistant Professor who applies data intensive approaches to investigate biological problems in creative and innovative ways… Candidates who apply **computational, quantitative, or data intensive methods in any area of Biology** will be considered…

## Bloomberg Distinguished Professorship in Evolutionary Genomics.

The Johns Hopkins University is searching for an outstanding senior scientist in the area of **Evolutionary Genomics** for an endowed chair as a Bloomberg Distinguished Professor. This position will be held jointly between the Department of **Biology** (Krieger School of Arts and Sciences) and the **Institute for Genetic Medicine** (JHU School of Medicine).

**More Info:** http://www.bio.jhu.edu/Events/Jobs/Default.aspx
**Or contact me:** james@taylorlab.org