# Integrating Galaxy with Globus Online: Lessons Learned from the CVRG Project

Bo Liu

Computation Institute, University of Chicago
Argonne National Laboratory

Computation Institute

CardioVascular Research Grid

Argonne NATIONAL LABORATORY

THE UNIVERSITY OF CHICAGO

www.ci.anl.gov
www.ci.uchicago.edu

# Agenda

- Introduction
- **New Galaxy Tools**
  - Globus Online
  - CRData
  - Picard/GATK via Condor
  - CummeRbund
  - Miso
- System Implementation
  - Amazon Cloud
  - Globus Provision
- CVRG Use Case
- Integration of Galaxy and Globus
- Conclusions & Future Plan

# Introduction

- ## Computation Institute

  - A joint institute of Argonne National Laboratory and the University of Chicago

  - Interdisciplinary approaches to challenging systems problems

  - Development and application of advanced computational methods

  - Co-hosts of the GCC 2012

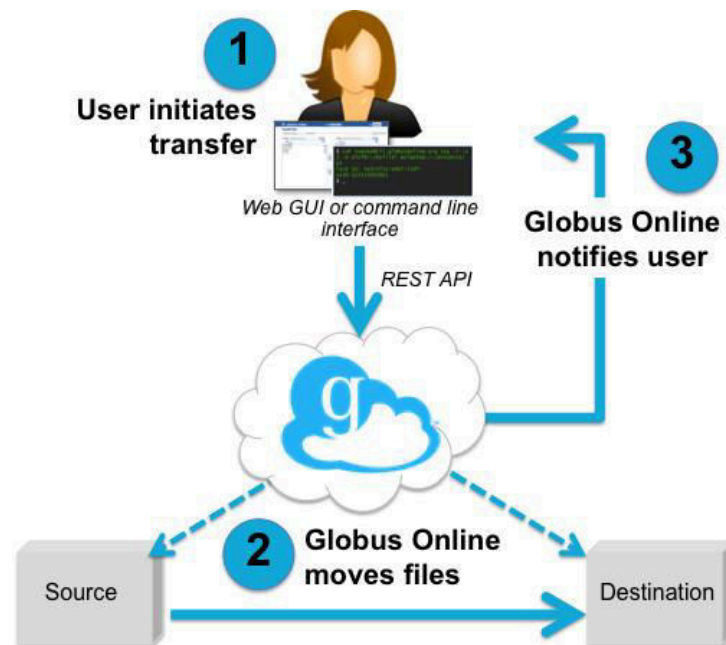- ## CVRG project

  
  CardioVascular Research Grid

  - Funded by NHLBI

  - Aim: create an infrastructure for sharing cardiovascular data and data analysis tools

  - Collaborator: Institute for Computational Medicine at the Johns Hopkins University

  http://cvrgrid.org/

# Introduction

- **Galaxy**
  - Web-based platform, ease of use and distribution
  - Various sequencing analysis tools
  - Easy data/workflow/history sharing
- **Challenges**
  - Distributed data at different locations
    - Transfer and sharing data with other researchers
  - Inefficient ways of data movement
    - Browser Upload: 2GB limit
    - FTP: not stable enough
  - Distributed computing for high performance
  - Requirements for new tools
    - CRData, execute R script
    - CummeRbund , analyze Cufflinks RNA-Seq output
    - Miso, compute Psi values for genes/events
  - Automatic deployment of Galaxy on Cloud

# Globus Transfer

- ## Globus Online - No IT required
  - Software as a Service (SaaS)
  - Consolidated support & troubleshooting
  - Works with existing GridFTP servers
  - Globus Connect solves "last mile problem"

- ## Reliable file transfer
  - Easy "fire-and-forget" transfers
  - Automatic fault recovery
  - High performance
  - Across multiple security domains

- ## > 5500 registered users, > 6PB moved

- ## Recommended and used by DOE Facilities, NSF Supercomputing centers, many campuses



**1** User initiates transfer — Web GUI or command line interface — REST API

**3** Globus Online notifies user

**2** Globus Online moves files

Source → Destination

https://www.globusonline.org/



Reliable, high-performance, secure file transfer.
**Move files fast. No IT required.**

WATCH A VIDEO — Globus Online in a nutshell

GET STARTED — Sign up and get moving

6,048,216,038 MB TRANSFERRED

CVRG | Argonne NATIONAL LABORATORY | THE UNIVERSITY OF CHICAGO | www.ci.anl.gov www.ci.uchicago.edu

# Integration of Galaxy and Globus Online

- Transfer large quantities of data in and out of Galaxy fast and reliably
- 3 Globus Online tools
  - GO transfer
    - transferring data between different endpoints
  - Get data via Globus Online
    - getting data from GO endpoint to Galaxy server
  - Send data via Globus Online
    - sending data from Galaxy to GO endpoint
- Set "Deadline" for transfer
- Email notification
- Transfer monitoring

# CRData

- R
  - A software for statistical computing and graphics

- CRData
  - Originally created as part of crdata.org
  - Wrap 35 tools
    - heatmap_plot_demo.R
    - sequenceDifferentialExperssion.R
    - affyClassify.R
    - ......
  - Execute a set of BioConductor R scripts
  - Used for analyzing ECG Data in CVRG

- Complements the functionality of executing CRData R scripts in Galaxy



CRData
- scotttest.R
- arith.R
- Heston.R
- hostInfo.R
- Test_Script.R
- heatmap_plot_demo.R
- Demo_of_RSBML_and_RgraphViz
- sequenceDifferentialExperssion.R
- ggplot_demo.R
- Demo_search_for_a_motif_in_a
- Hmisc_xYplot_Demo.R
- CherryAdler_mutual_repression_-course_simulation.R
- PICS_plot_demo.R
- Demo_of_conditional_plotting.R
- Demo_sequence_logo_visualizat
- Demo_upload_a_file_from_a_rer
- Demo_stochastic_kinetic_simulat
- Demo_3D_plot_with_user_select
- SakaSmith_bistable_morphogen

# CRData

# Picard/GATK via Condor

- Wrap all the Picard and GATK tools to run through Condor job scheduler
- Leverage local cluster or cloud based scalable computational resources for parallelizing the tools
- Condor
  - Parallel execution
  - Automatically job scheduling
  - Condor pool: multiple worker nodes, easily add or remove
- Condor Runner for Galaxy
  - The executable along with command line options get passed to a condor runner
  - Applications are run on a worker node instead of the galaxy node
  - Enable faster and more efficient execution of Galaxy jobs

**NGS: GATK Tools via Condor**

ALIGNMENT UTILITIES
- Depth of Coverage (via Condor) on BAM files
- Print Reads (via Condor) from BAM files

REALIGNMENT
- Realigner Target Creator (via Condor) for use in local realign
- Indel R
BASE
- Count
- Table R
- Analyze
GENO
- Unified
ANNOT
- Variant
FILTR
- Variant
- Select
VARIA
- Variant
- Apply
VARIA
- Validate
- Eval Va
- Combir

**NGS: Picard via Condor**

CONVERSION
- FASTQ to BAM (via Condor) creates an unaligned BAM file
- SAM to FASTQ (via Condor) creates a FASTQ file

QC/METRICS FOR SAM/BAM
- BAM Index Statistics (via Condor)
- SAM/BAM Alignment Summary Metrics (via Condor)
- SAM/BAM GC Bias Metrics (via Condor)
- Estimate Library Complexity (via Condor)
- Insertion size metrics (via Condor) for PAIRED data
- SAM/BAM Hybrid Selection Metrics (via Condor) for targeted resequencing data

BAM/SAM CLEANING
- Add or Replace Groups (via Condor)
- Reorder SAM/BAM (via Condor)
- Replace SAM/BAM Header (via Condor)
- Paired Read Mate Fixer (via Condor) for paired data
- Mark Duplicate reads (via Condor)
- Picard SAM Format Converter (via Condor)
- Build BAM Index (via Condor)
- Sort Sam (via Condor)
- Clean Sam (via Condor)

# CummeRbund

- An R package designed to aid and simplify the task of analyzing Cufflinks RNA-Seq output

- Input: a backend database file from the history, or build a new database using cuffdiff output

- Plot type: Density, Boxplot, Scatter, Volcano, Heatmap, Cluster, Expression Plot, Expression Bar Plot

# CummeRbund

# Miso

- ## MISO (Mixture-of-Isoforms)
  - a probabilistic framework
  - Quantitate the expression level of spliced genes from RNA-Seq data
  - Identify regulated isoforms or exons across samples
- ## Wrap 5 Miso tools
  - ### Index GFF
    - o   create an index for the GFF annotations
  - ### Compute Psi values for genes
  - ### Compute Psi values for all events
    - o  compute Psi values for a given GFF annotation of either whole mRNA isoforms or isoforms produced by single alternative splicing events
  - ### Summarize Samples
    - o  summarize MISO output and obtain confidence intervals for Psi values
  - ### Compare Samples
    - o  compute comparison statistics between the two given samples

# Miso

www.ci.anl.gov
www.ci.uchicago.edu

# System Implementation

- ## Amazon EC2/EBS
  - scalable computing and storage capabilities
- ## Deploying Galaxy on Cloud
  - on-demand resource configuration
  - usage-based pricing
- ## Globus Provision
  - a tool for deploying fully-configured Globus systems
  - Automate the process of deploying Galaxy on EC2
  - Provision a EC2 cluster with Galaxy, Globus Online, Condor, GridFTP server and a set of users in <15mins
  - Reusable Chef recipes to provision production galaxy instances on demand

# CVRG Use Case



**Galaxy GUI**
Run "RNASeq Analysis Pipeline"
Transfer data from "Sequencing center"
to "Galaxy Endpoint"

**Globus Online**

**Galaxy**

- RNASeq analysis
  - Tophat
  - Cufflink
  - CummeRbund
  - Miso
  - CRData
  - ……

**GO Endpoint**

FTP, SCP, others

FTP, SCP, HTTP

Seq Centers

FedEx
Home Delivery

~~s Online~~ ~~ides a~~

- High-perfo~~
- Fault-tolerant
- Secure~~

file transfer Service between all data-endpoints

FTP, SCP

Research Lab

SCP

JHU

PCBC

UC

Cloud Storage

Local Cluster

**Globus Provision Deployment**

**Condor Scheduler**

amazon web services

- Analytical tools are automatically run on scalable compute resources

**Galaxy on EC2**

Data Management

Data Analysis

# RNASeq Analysis Pipeline

# Miso Pipeline

# Towards "research IT as a service"



**Research Data Management as a Service**

| Globus Transfer | Globus Storage | Globus Collaborate | Globus Catalog | ... SaaS |

Globus Integrate platform — PaaS

(1) Collect Data
(2) Move to Storage Store
(3) Ingest Processing
(4) Move to Community Store
(5) Publish in Registry
(6) Validate
(7) Backup
(8) Mirror
(9) Search, Browse Analyze, Update, Annotate

# Integrate Galaxy and Globus Online

## Globus Services

- Transfer
- User (Authentication)
- Group
- Collaborate
- Storage
- Compute

## Globus

Integrate

## GO-Galaxy

- GO Transfer tools
- Log in to Galaxy via GO account
- Use GO's Group management
- Share history with user/group at GO
- Use the data at GO storage
- Add clusters to condor pool for distributed computing

# Conclusions & Future Plan

- ## Galaxy + Globus Online
  - Transfer large-scale datasets in and out of Galaxy in a secure, efficient and fast way
- ## New Galaxy tools
  - CRData: processing R scripts
  - CummeRbund, Miso: RNA analysis
  - Condor: distributed computing capabilities
- ## Globus Provision
  - Automatically deploy Galaxy on EC2 with user-specific configuration
- ## Future Plan
  - Thorough integration of Galaxy and Globus
  - Performance optimization
- ## Galaxy community Toolshed contributions

# Thank you for your attention

Bo Liu:  boliu@uchicago.edu
Ravi Madduri:  madduri@mcs.anl.gov

CVRG | Argonne NATIONAL LABORATORY | THE UNIVERSITY OF CHICAGO | www.ci.anl.gov
www.ci.uchicago.edu