

A SUSTAINABLE NATIONAL GATEWAY FOR BIOLOGICAL COMPUTATION

XSEDE 13

July 24, 2013

Galaxy Team:

james.taylor@emory.edu

anton@bx.psu.edu

nate@bx.psu.edu

PSC Team:

blood@psc.edu

ropelews@psc.edu

josephin@psc.edu

yanovich@psc.edu

rbudden@psc.edu

zhihui@psc.edu

sergiu@psc.edu

Overview

- ▣ Galaxy: The Need, the Framework, the Popularity Bottleneck
- ▣ Distributing Galaxy work and data flows to XSEDE systems: first steps
- ▣ The Future: Galaxy Gateway(s)

643 HiSeqs = 6.5 Pb/year



The Challenge

Biology has rapidly become data intensive, and dependent on computational methods

How can we ensure that these methods are accessible to researchers?

...while also ensuring that scientific results remain reproducible?

Galaxy

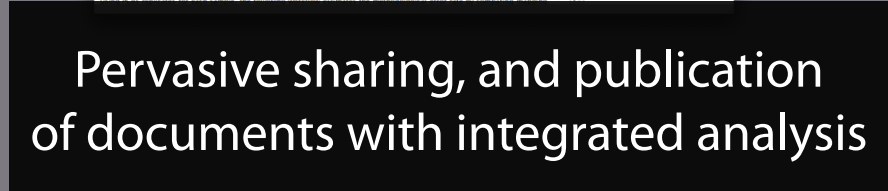
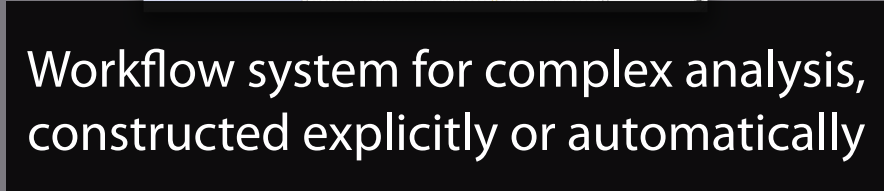
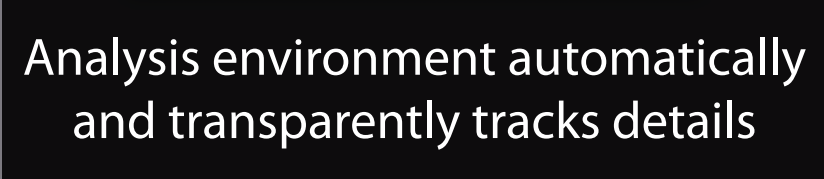
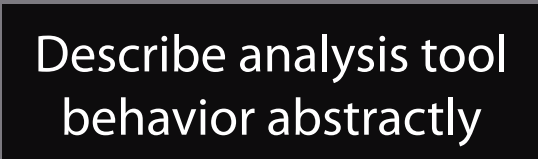
A free (for everyone) web service integrating a wealth of tools, compute resources, terabytes of reference data and permanent storage

Open-source software allowing anyone to freely deploy or extend this platform

A community of users and developers

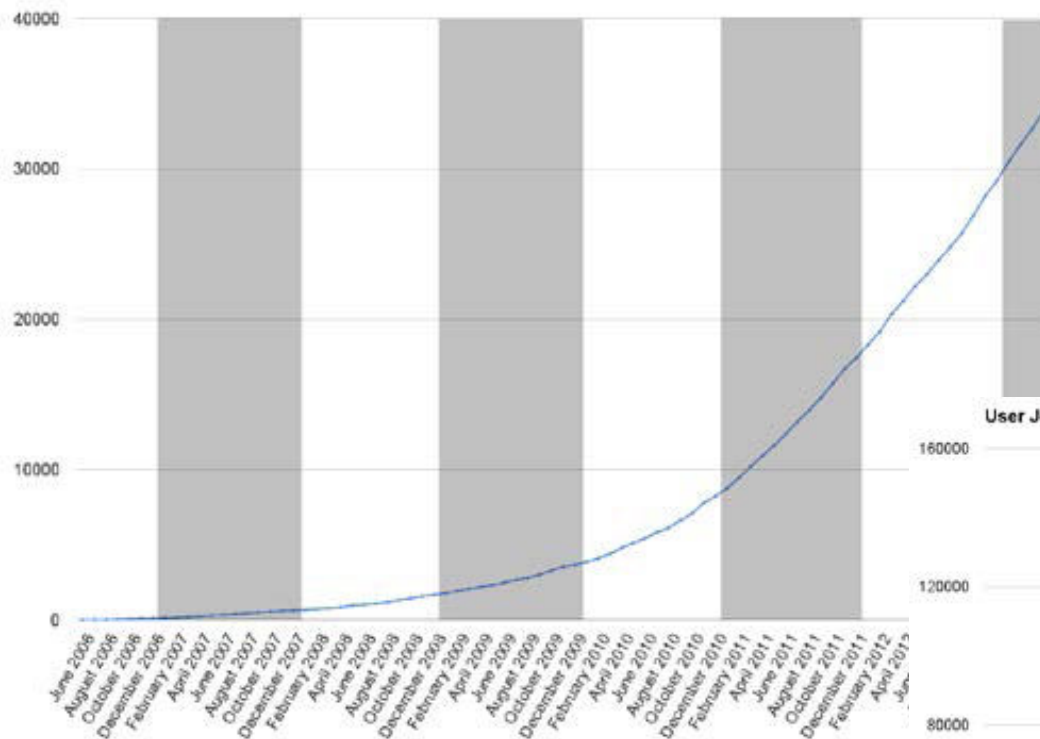
Galaxy integrates existing tools into a uniform framework

- ▣ Mostly command line tools, a declarative XML description of the interface, how to generate a command line
- ▣ Designed to be as easy as possible for tool authors, while still allowing rigorous reasoning
- ▣ Workflows can be constructed from scratch or extracted from existing analysis histories
- ▣ Facilitate reuse, as well as providing precise reproducibility of a complex analysis

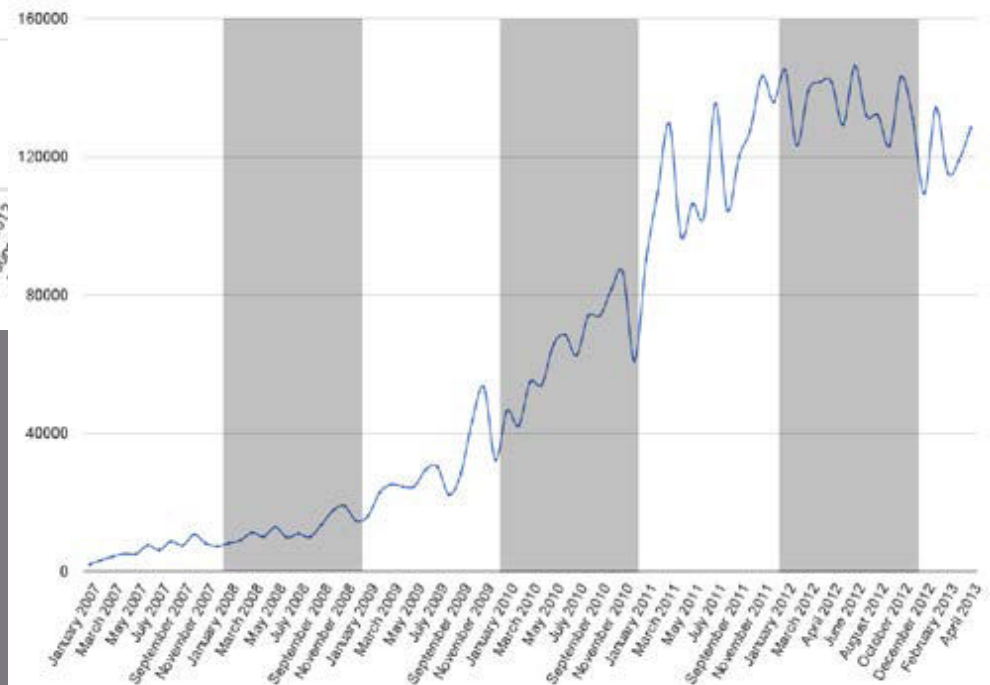


The Popularity Bottleneck

Registered Users on Galaxy Main



User Jobs per month on usegalaxy.org



So, send jobs to XSEDE!

- ▣ Entire Galaxy workflows or component tasks.
- ▣ Especially, tasks that **require HPC**, e.g. de-novo assembly applications *Velvet* (of genome) and *Trinity* (of transcriptome) to PSC *Blacklight* (up to 16 TB of coherent shared memory per process).
- ▣ Should be transparent to the user of usegalaxy.org .

Problems to be solved

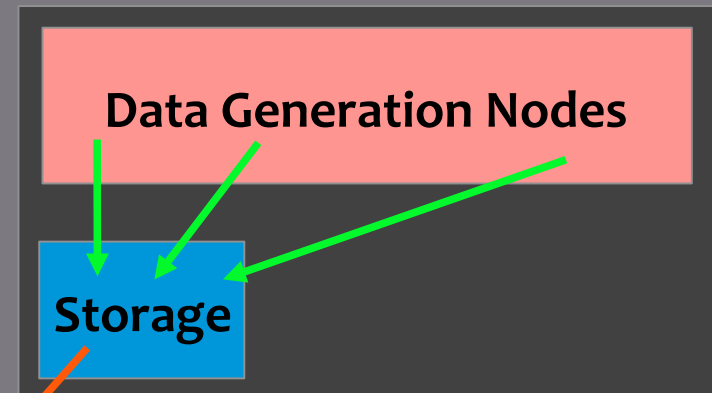
- ▣ **Data Migration:** Galaxy currently relies on a shared filesystem between the instance host and the execution server to store the reference and user data required by the workflow. This is implemented via NFS.
- ▣ **Remote Job Submission:** Galaxy job execution currently requires a direct interface with the resource manager on the execution server.

Initial Galaxy Data Staging to PSC

Transferred 470TB in 21 days from PSU to PSC
(average ~22TB/day; peak 40 TB/day)

rsync used to initially stage and synchronize
subsequent updates

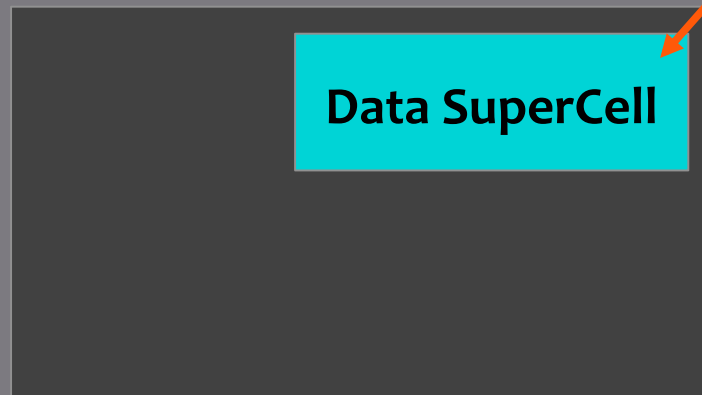
Data copy maintained in PSC in **/arc** file system
available from compute nodes



Penn State

10gigE link

PSC



Data SuperCell

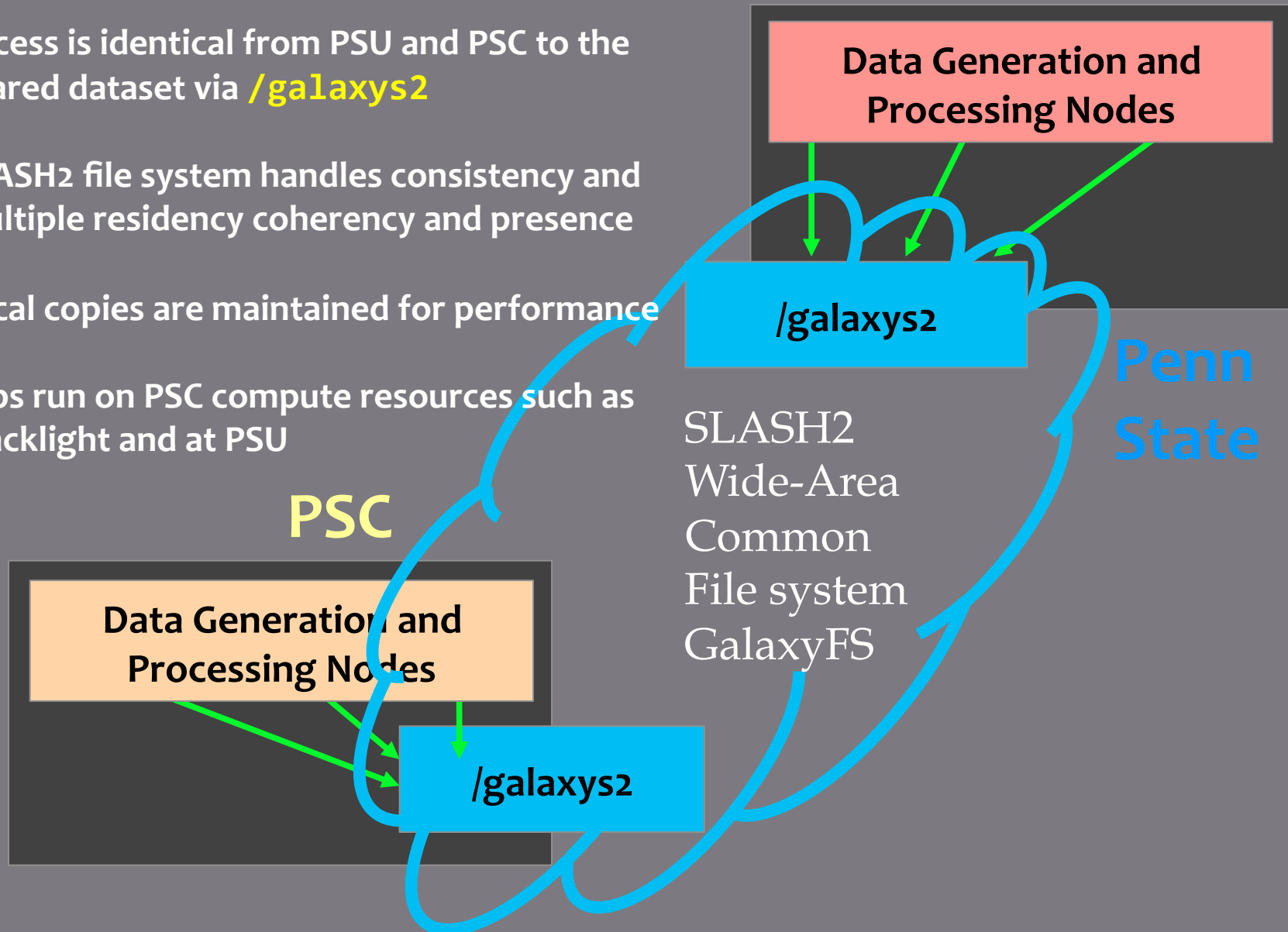
Galaxy Remote Data Architecture

Access is identical from PSU and PSC to the shared dataset via **/galaxys2**

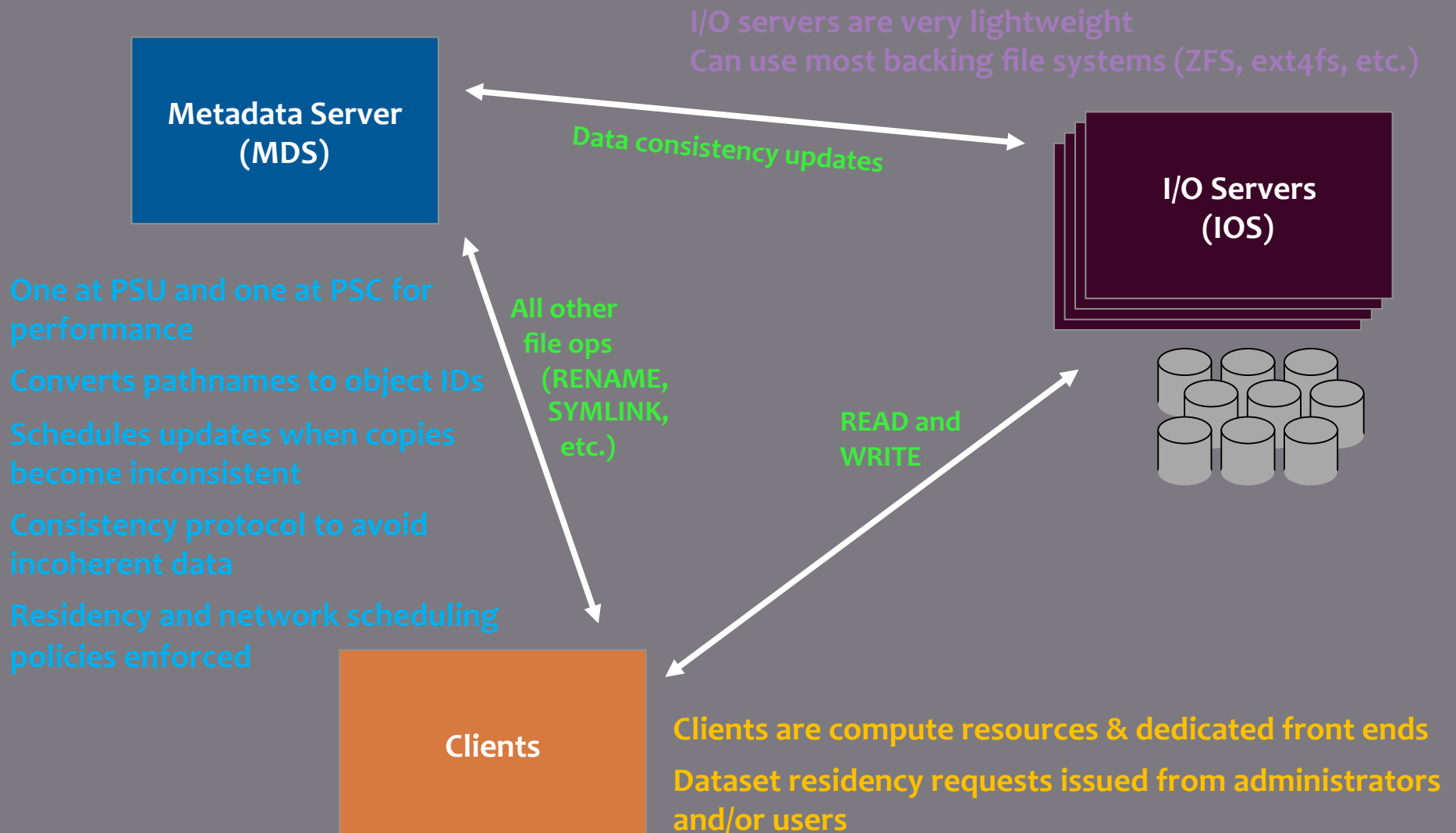
SLASH2 file system handles consistency and multiple residency coherency and presence

Local copies are maintained for performance

Jobs run on PSC compute resources such as Blacklight and at PSU



Underlying SLASH2 Architecture



Submitting to Blacklight

- ▣ Created a new Galaxy job-running plugin for submission via a remote shell program and PSC “*SIMON*” Torque command line: CLI runner.
- ▣ *Velvet* and *Trinity* have been incorporated into the Galaxy web platform using Galaxy’s XML interface.
- ▣ Test jobs have been successfully submitted from Penn State and executed on Blacklight using the data replicated via SLASH2 from Penn State to PSC.

Future Work (1)

- ▣ Integrate this work with the production public Galaxy site, usegalaxy.org
- ▣ Dynamic job submission, allowing the selection of appropriate remote or local resources (cores, memory, walltime, etc.) based on individual job requirements, using an Open Grid Services Architecture Basic Execution Service compatible service, such as *Unicore*.

Future Work (2)

- ▣ Galaxy-controlled data management, to intelligently create replicas as close as possible to the compute resource that will use the data.
- ▣ Authentication with Galaxy instances, using XSEDE or other credentials, e.g., InCommon/CILogon.
- ▣ Additional data transfer capabilities in Galaxy, such as IRODS and Globus Online.

The Vision

Ultimately, we envision that any Galaxy instance (in any lab, not just Galaxy Main) will be able to spawn jobs, access data, and share data on external infrastructure whether this is an XSEDE resource, a cluster of Amazon EC2 machines, a remote storage array, etc.