

Outsourcing IT complexity

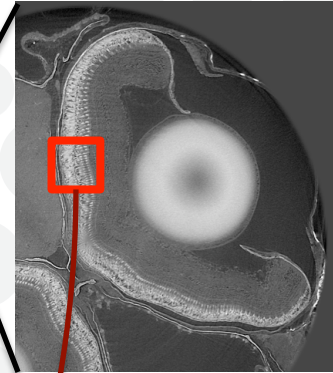
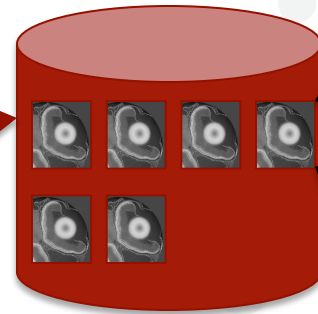
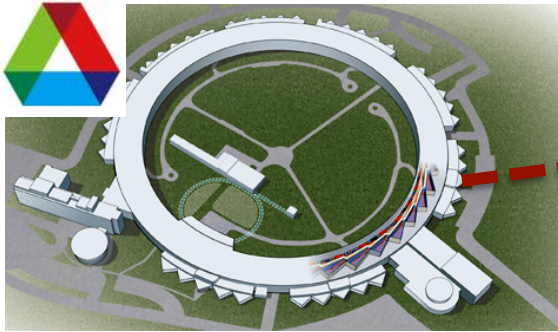
Moving Ultraviz management from the laboratory to the cloud

Ian Foster

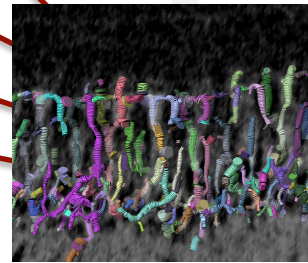
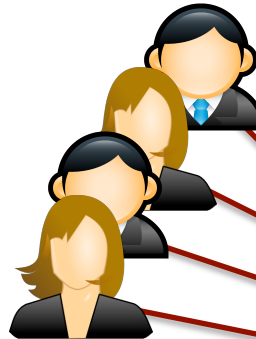
A story of modern science



Keith Cheng: map genotype \rightarrow phenotype for
~3000 zebrafish mutants



**Collect, move, store,
index, analyze, share,
update, millions of files**



Gordon
Kindlmann



IT complexity greater than that of typical enterprise

Sources of complexity in science



- Run experiments
- Collect data
- Manage data
- Move data
- Analyze data
- Run simulations
- Compare experiment with simulation
- Search the literature
- Share results
- Communicate with colleagues
- Publish papers
- Find, configure, install relevant software
- Find, access, analyze relevant data
- Document research
- Order supplies

Outsourcing complexity in business



- Web presence
- Email (hosted Exchange)
- Calendar
- Telephony (hosted VOIP)
- Human resources and payroll
- Accounting
- Customer relationship mgmt

Software
as
a
Service
(SaaS)

Outsourcing complexity in business



- Web presence
- Email (hosted Exchange)
- Calendar
- Telephony (hosted VOIP)
- Human resources and payroll
- Accounting
- Customer relationship mgmt
- Data analytics
- Content distribution
- ...

Software
as
a
Service
(SaaS)

Platform/
Infrastructure
as a Service
(PaaS/IaaS)

Outsourcing complexity in science



- Run experiments
- Collect data
- Manage data
- **Move data**
- Analyze data
- Run simulations
- Compare experiment with simulation
- Search the literature
- Share results
- Communicate with colleagues
- Publish papers
- Find, configure, install relevant software
- Find, access, analyze relevant data
- Document research
- Order supplies

Globus Online = Sci-SaaS

A nuclear physicist sharing data



- Hai Ah Nam, a nuclear physicist from Oak Ridge spoke at GlobusWorld March 2010 about the struggles with moving data
- Transferring 1.6 TB (86 large files) from Oak Ridge to NERSC
- Changed from using scp to GridFTP to reduce transfer from days to hours
- Reduced transferring 137 TB from months to days
- **But, it was not easy...**

Data movement is hard



For many reasons

- SCP is too slow
- No GridFTP at site
- Firewalls
- Space management
- Net (mis)configuration
- Security config, policies
- Other heterogeneities
- Failures, restarts, mirroring, other tasks

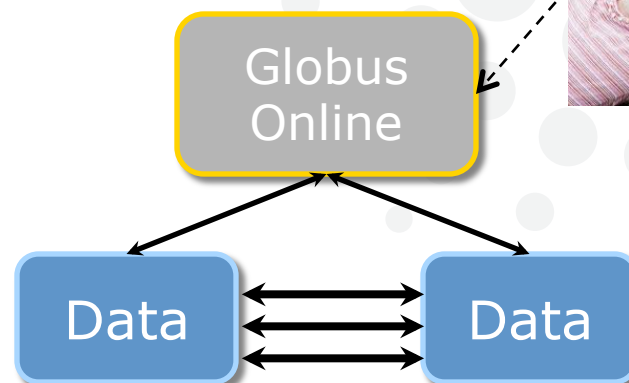
For many people

- Ad-hoc: Non-experts who need to move many files
- Scripted: Experts who want to automate large file movement
- System builders: Don't want to re-engineer solutions to such problems

Globus Online: A “Web 2.0” solution



- Outsource the mundane
 - Manage site configurations, credentials, network configurations, ...
 - Monitor transfers
- Automate the repetitive
 - Retry failed transfers, mirror directories, ...
- Radically simplify interfaces
 - REST
 - CLI 2.0 (including scp)





Software-as-a-Service (SaaS)

Platform-as-a-Service (PaaS)

Infrastructure-as-a-Service (IaaS)



- Service: Built as scale-out web application
 - Hosted on Amazon Web Services
- Client: Minimize software deployment
 - Web 2.0
 - AJAX + REST
 - Notification via email, IM, SMS, Twitter, etc.
 - Enable mash-ups
 - “CLI 2.0”
 - `ssh cli.globus.org ...`
 - Support for heterogeneity in end systems: data transfer and security protocols, etc.

Why SaaS?

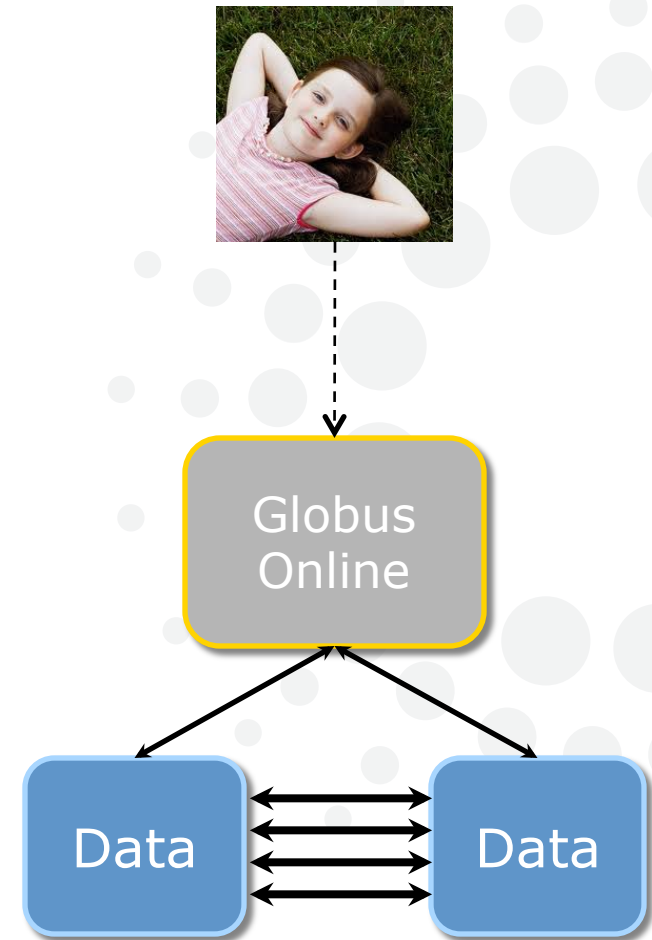


- Requires no user software installation
 - Minimal IT proficiency required
- Efficient software delivery lifecycle
 - Updates developed, tested and deployed quickly
- Consolidated troubleshooting and support
 - An expert group can proactively detect and correct problems
 - Partnering with Argonne/UC IT support group who specialize in support and ops

The Globus.org user can ...



- **Register** with Globus Online
- Update their **profile**
- Define **endpoints** that can then be **activated**
- **Transfer** data between endpoints—creating a task
- Monitor **status** of task(s)
- View **event(s)** for task(s)



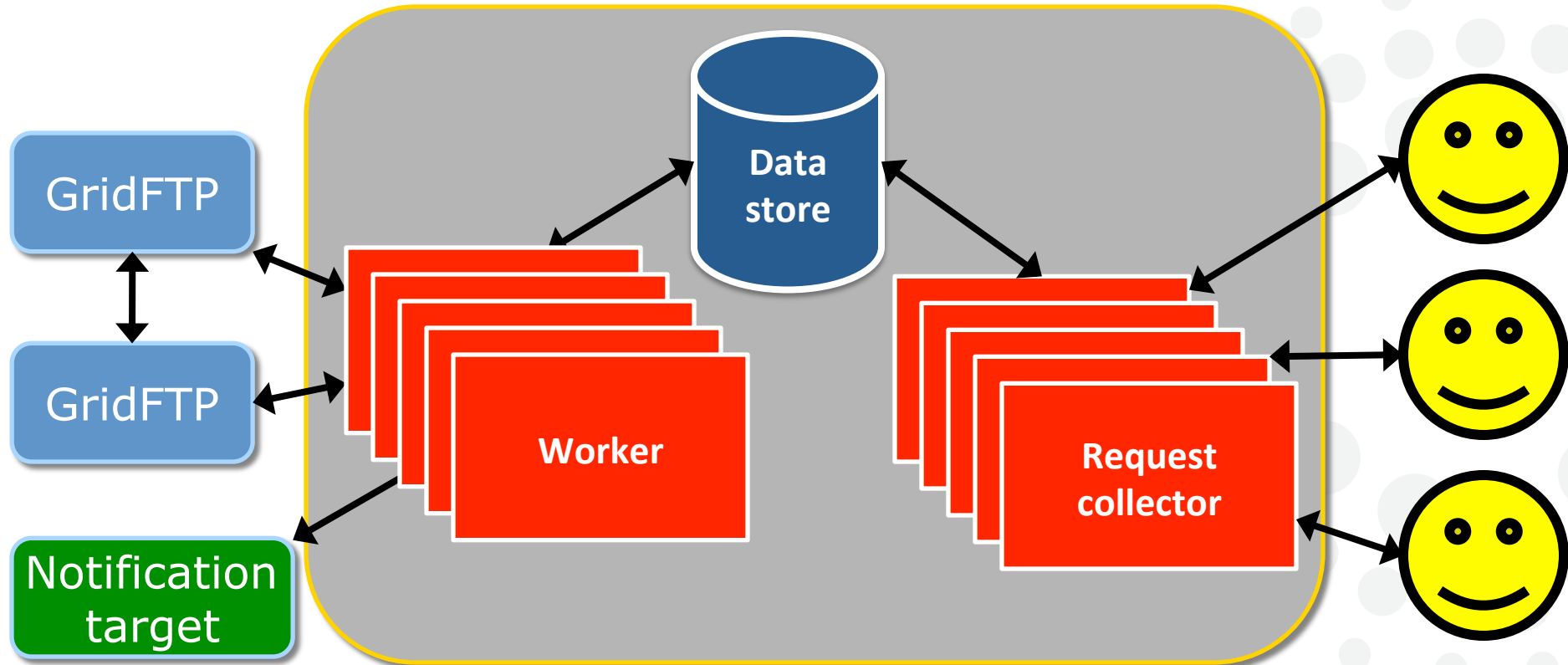


- CLI
 - ssh ME@cli.globusonline.org COMMAND
- REST
 - Same thing, but in HTTP packets
- GUIs
 - As with any Web 2.0 system, a variety of graphical interfaces can be created easily, using Ajax and other technologies

Use as needed to hide password text

```
gsissh <user>@cli.globusonline.org <command> <options> <params>
```

A peek inside Globus Online



- 100s of NERSC users transfer large amounts (>20 GB) of data to/from NERSC with **scp**
- Globus Online **scp** option provides higher performance and reliability ... with simplicity of scp command
- Extra benefits
 - Detached transfer: async, fire and forget
 - Automatic recovery from network, end-system failures
 - Load balancing & fail-over
 - End-to-end verification, sync, ...

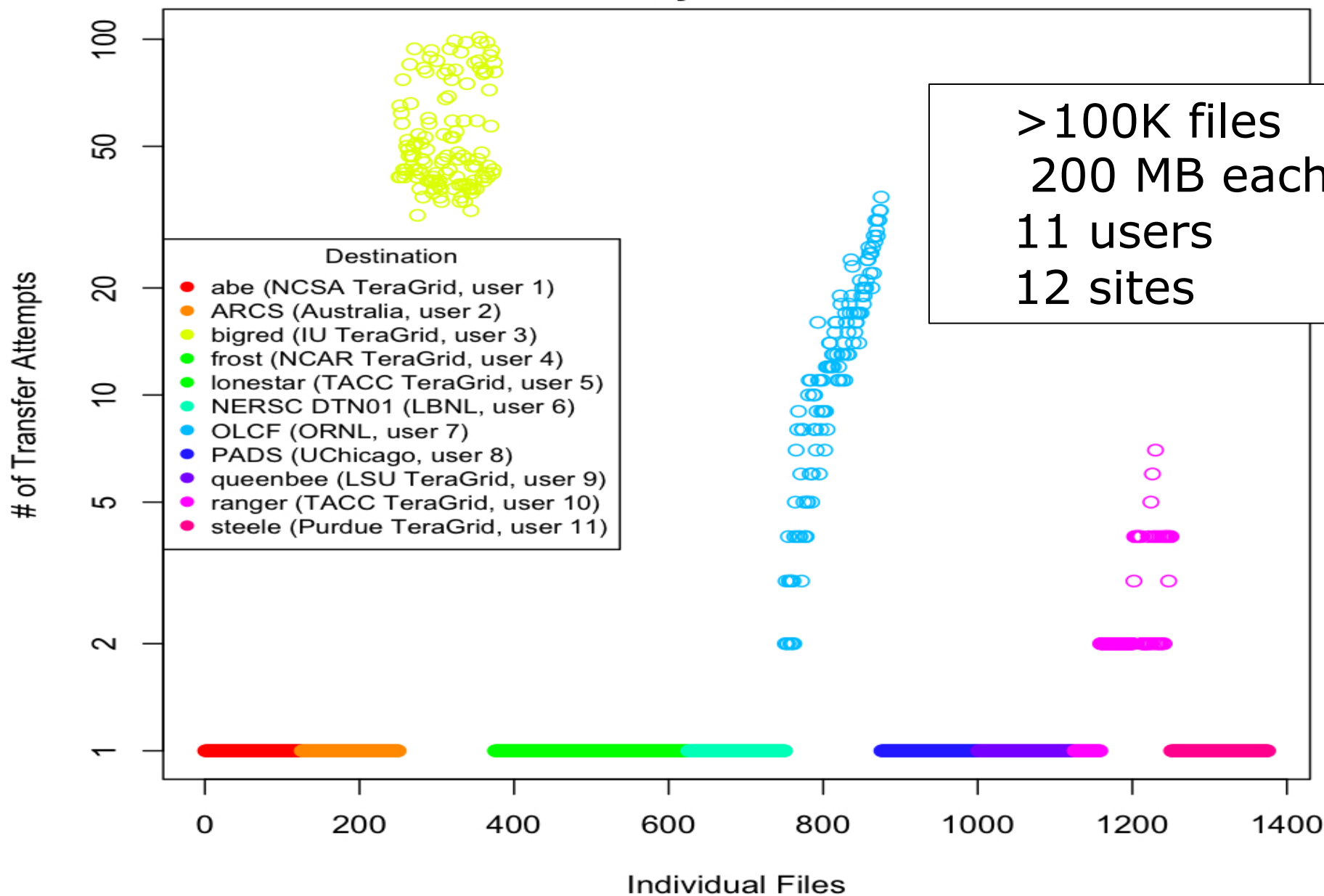
The simple scp command



```
ssh lcc@cli.globusonline.org scp go#ep1:/share/godata/file1.txt
go#ep2:~/myfile.txt
Contacting 'myproxy.tutorial.globusonline.org' ...
Activating 'ep2'
Activating 'ep1'
Task ID: 19029d64-ecec-11df-aa30-1231350018b1
[XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX] 1/1 0.00
mbps
ssh lcc@cli.globusonline.org ls go#ep2/~ /
myfile.txt
```

cancel
details
endpoint-activate
endpoint-add
endpoint-deactivate
endpoint-list
endpoint-modify
endpoint-remove
endpoint-rename
events
ls
profile
scp
status
transfer
versions
wait

CEDPS Data Challenge #3: Attempts (ordered by Destination, Time)
11 users each transferring 125 200M files from ALCF
May 11 2010





- Chan Wilson (GFDL): a sys admin / script developer for a group of climate scientists
- Creating automation scripts for sustaining 80TB / day of simulation output from remote supercomputers to GFDL
- ESnet w/ striped GridFTP endpoints
 - 12 DTN/GridFTP servers attached to GFDL cluster filesystem
- He is counting on Globus Online to help him move that data reliably with required performance



- Condor supports file stage-in/out for each job
- Condor is adding support for Globus.org for staging files to/from compute hosts to off-site locations
- Will be available in version 7.6.0 (~Nov 2010)

- Lightweight transfer agent (firewalls, sites without GridFTP installed)
- Group management
- Higher-level data management capabilities
 - E.g., data publication, replication
- Computation management
 - E.g., Swift computations

Other Sci-SaaS services planned



- VO management
 - Groups, membership, policies (via Grouper)
 - Computation mgmt
 - Run this computation
 - Analyze any files deposited in directory
 - Data management
 - Storage and indexing
 - Archiving and lifecycle
 - Publication and sharing
 - Workflow management
 - Data ingest and analysis pipelines
 - Meta-computations
 - Uncertainty quantification
 - Optimization
- Looking for partners to, e.g.
- Expand set of services
 - Host services in other geographies

Numerous people have contributed to this work, including:

Bryce Allen, Joshua Boverhof, John Bresnahan, Lisa Childers, Paul Dave', Fred Dech, Ian Foster, Dan Gunter, Gopi Kandaswamy, Nick Karonis, Raj Kettimuthu, Jack Kordas, Lee Liming, Mike Link, Stu Martin, JP Navarro, Karl Pickett, Mei Hui Su, Steve Tuecke, Vas Vasiliadis

Funders

- DOE and NSF

See Globus Online in Action at SC10 - Booth 4355

We will be demonstrating the power of Globus Online at SC10, being held in New Orleans, November 16-18, 2010. There are many opportunities available to learn about and use the service.

Globus Online Tutorials

Attend a 30-minute introductory, hands-on tutorial. Learn to use Globus Online and get a cool tee-shirt. Tutorials will be held in Booth 4355 at the following times:

Tuesday, November 16: 11:00am, 1:30pm, and 3:00pm

Wednesday, November 17: 11:00am, 1:30pm, and 3:00pm

Thursday, November 18: 10:30am

We will also host a more advanced tutorial and demonstration at 4:30pm on Tuesday, Nov. 16, and Wednesday, Nov. 17. These sessions are appropriate for users that have attended one of the introductory tutorials above, and want to understand how to use Globus Online more effectively in their specific environment.

[Register for an SC Tutorial](#)