# **Science as a service**
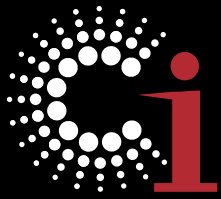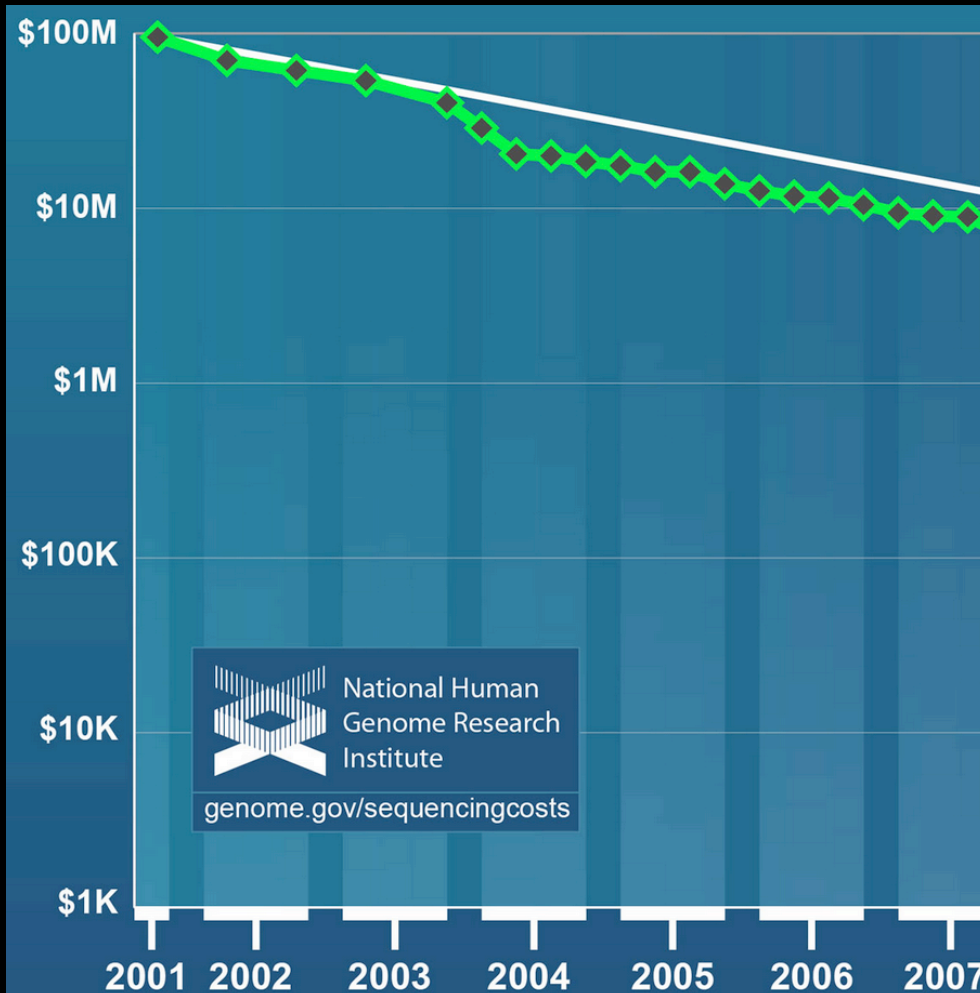## How on-demand computing can accelerate discovery

Ian Foster
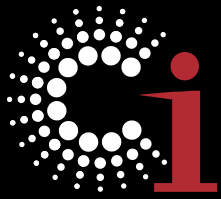
foster@anl.gov

# A time of disruptive change
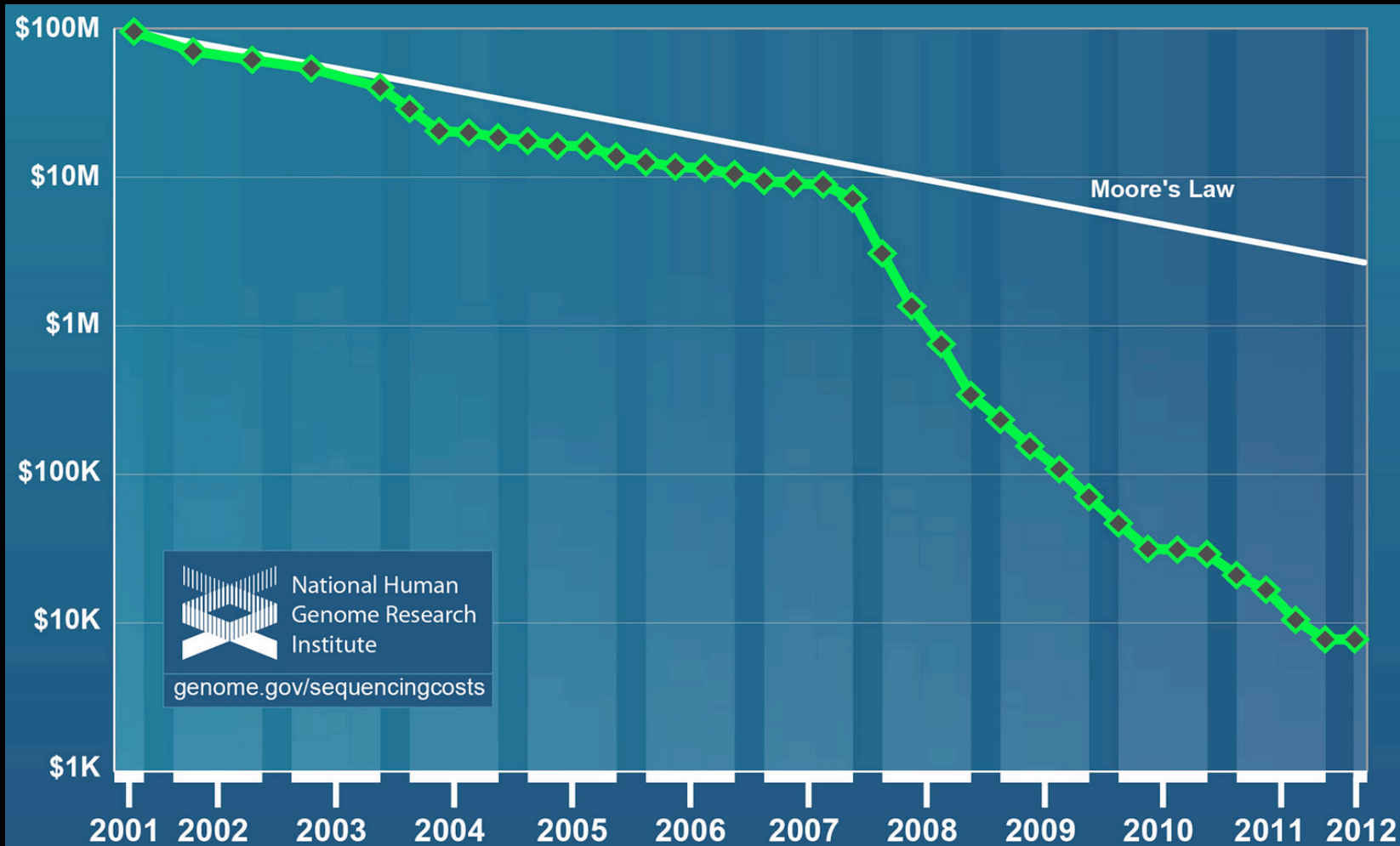## As evidenced by cost per human genome

# But most labs have extremely limited resources



Heidorn: NSF grants in 2007

< $350,000
80% of awards
50% of grant $$

computation**institute**.org

**Automation** is required to apply more sophisticated methods to far more data

**Automation** is required to apply more sophisticated methods to far more data
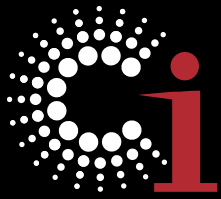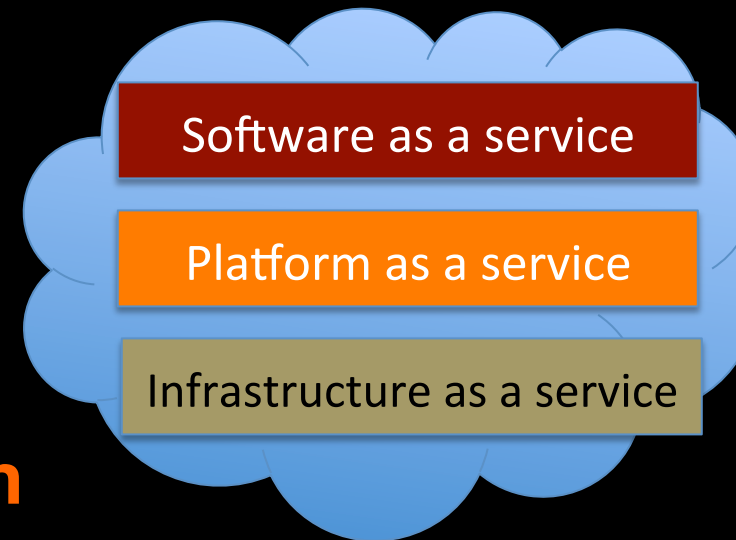


**Outsourcing** is needed to achieve economies of scale in the use of automated methods
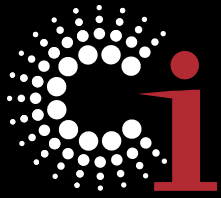
computation*institute*.org

# Building a discovery cloud

- Identify **time-consuming activities** that appear amenable to automation and outsourcing

- Implement as high-quality, low-touch **SaaS solution**

- Leverage **commercial IaaS** for reliability, economies of scale

- Extract common elements as a **research automation platform**

Software as a service

Platform as a service

Infrastructure as a service

Bonus question: Identify methods for delivering Discovery Cloud elements sustainably
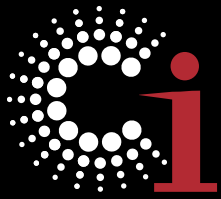
# Where does time go in research?
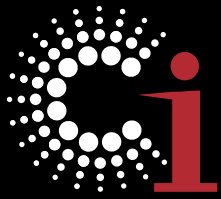
**The FDP Faculty Burden Survey**

42% of the time spent by an average PI on a federally funded research project was reported to be expended on administrative tasks related to that project rather than on research.

## 42%!!

We aspire (initially) to create a great user experience for **research data management**

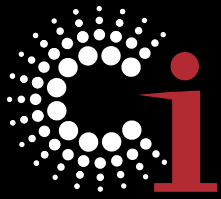What would a "dropbox for science" look like?

- Collect
- Move
- Sync
- Share
- Analyze
- Annotate
- Publish
- Search
- Backup
- Archive

...for BIG DATA

It should be trivial to **Collect, Move, Sync, Share, Analyze, Annotate, Publish, Search, Backup, & Archive** BIG DATA … but in reality it's often very challenging
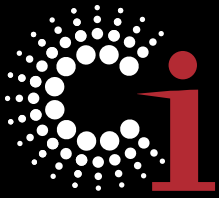
- Collect
- Move
- Sync
- Share
- Analyze
- Annotate
- Publish
- Search
- Backup
- Archive

...for **BIG DATA**
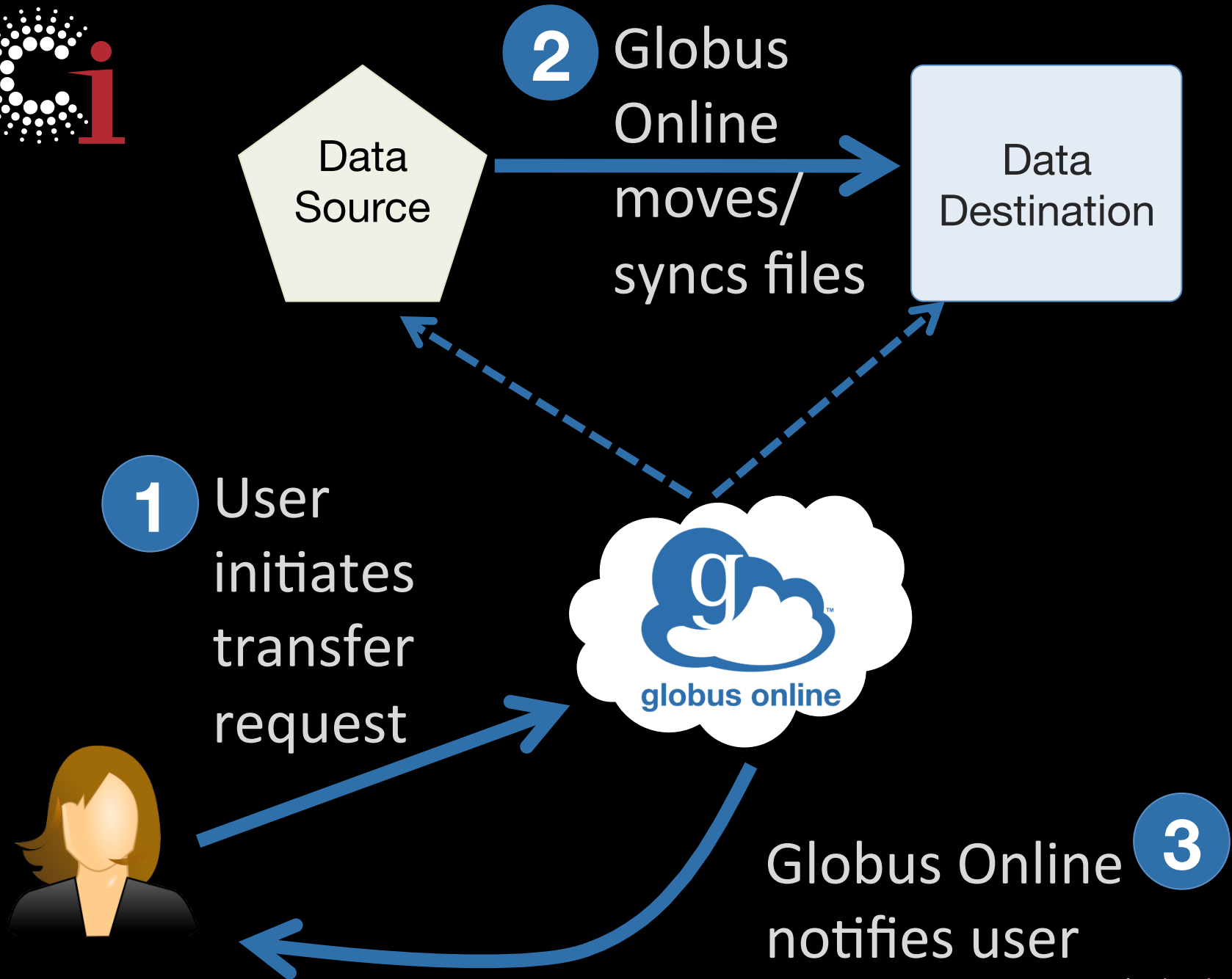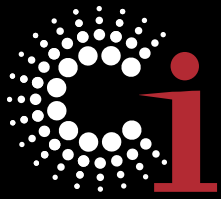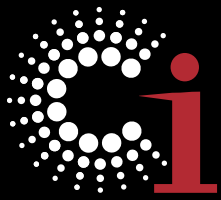
- Collect
- Annotate
- **Move**
- **Sync**
- **Share**

Publish

Search

Backup

Analyze

Archive

globus online

**Capabilities delivered using Software-as-Service (SaaS) model**

computationinstitute.org

**2** Globus Online moves/syncs files

Data Source

Data Destination

**1** User initiates transfer request

globus online

Globus Online notifies user **3**

computationinstitute.org

**2** Globus Online tracks shared files; no need to move files to cloud storage!

Data Source

**1** User A selects file(s) to share; selects user/ group, sets share permissions

globus online

**3** User B logs in to Globus Online and accesses shared file

# Extreme ease of use

- InCommon, Oauth, OpenID, X.509, …
- Credential management
- Group definition and management
- Transfer management and optimization
- Reliability via transfer retries
- Web interface, REST API, command line
- One-click "Globus Connect" install
- 5-minute Globus Connect Multi User install

# Early adoption is encouraging

# Early adoption is encouraging

**10,000 registered users; >100 daily**

**~18 PB moved; ~1B files**

**10x (or better) performance vs. scp**

**99.9% availability**

**Entirely hosted on Amazon**
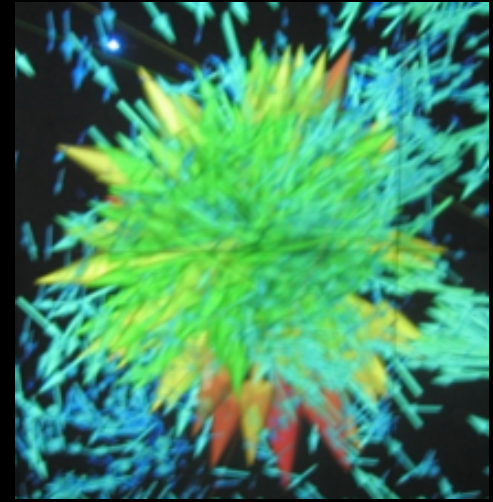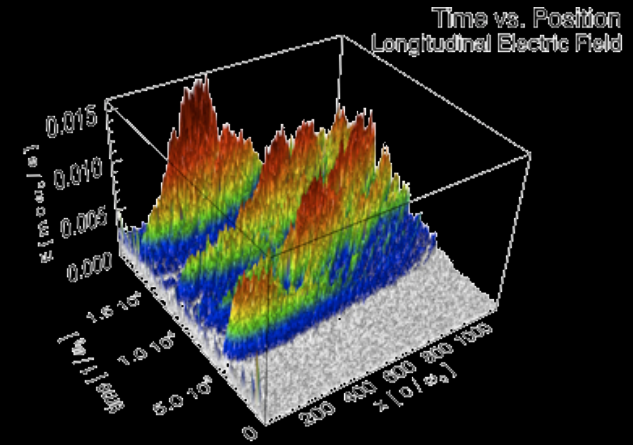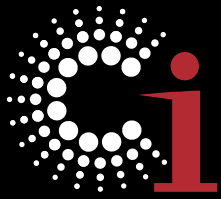
Duration of runs, in seconds, over time.
Red: >10 TB transfer; green: >1 TB transfer.

# K. Heitmann (Argonne) moves 22 TB of **cosmology** data LANL → ANL at 5 Gb/s
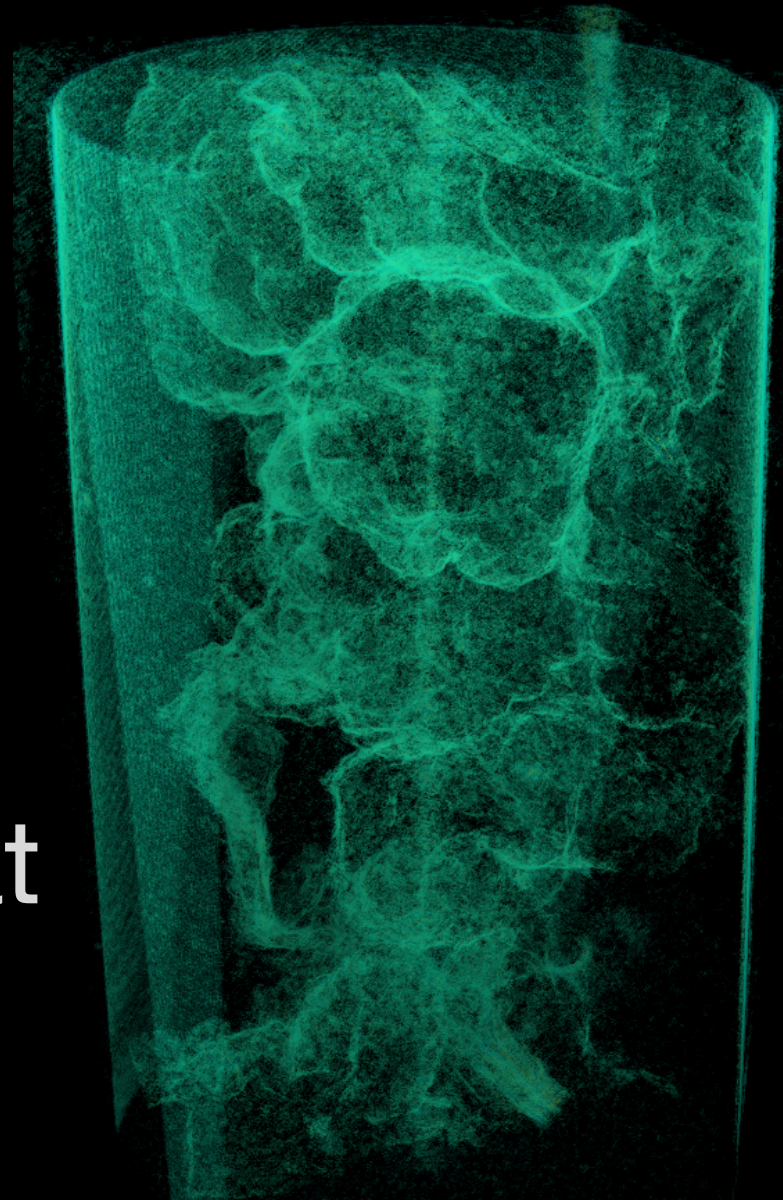
Time vs. Position
Longitudinal Electric Field

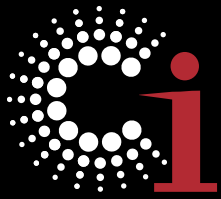# B. Winjum (UCLA) moves 900K-file **plasma physics** datasets UCLA →NERSC

# Dan Kozak (Caltech) replicates 1 PB LIGO **astronomy** data for resilience

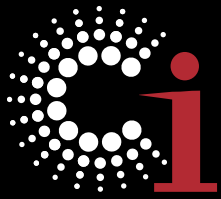Erin Miller (PNNL) collects data at Advanced Photon Source, renders at PNNL, and views at ANL

- Collect
- **Move**
- **Sync**
- **Share**
- Analyze
- Annotate
- Publish
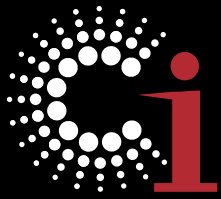- Search
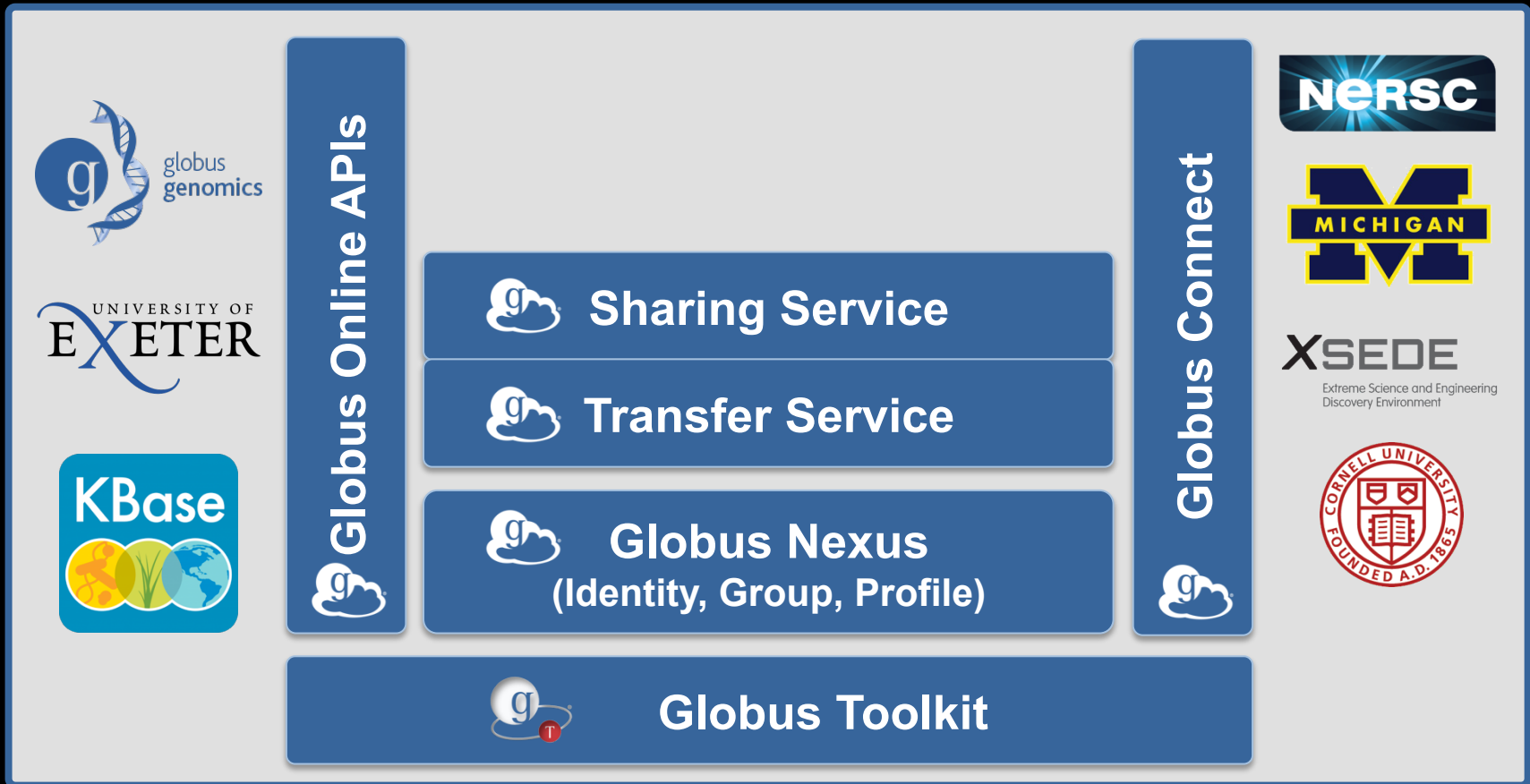- Backup
- Archive

...*for* BIG DATA

- Collect
- Move
- Sync
- Share
- Analyze
- Annotate
- Publish
- Search
- Backup
- Archive

*...for* **BIG DATA**
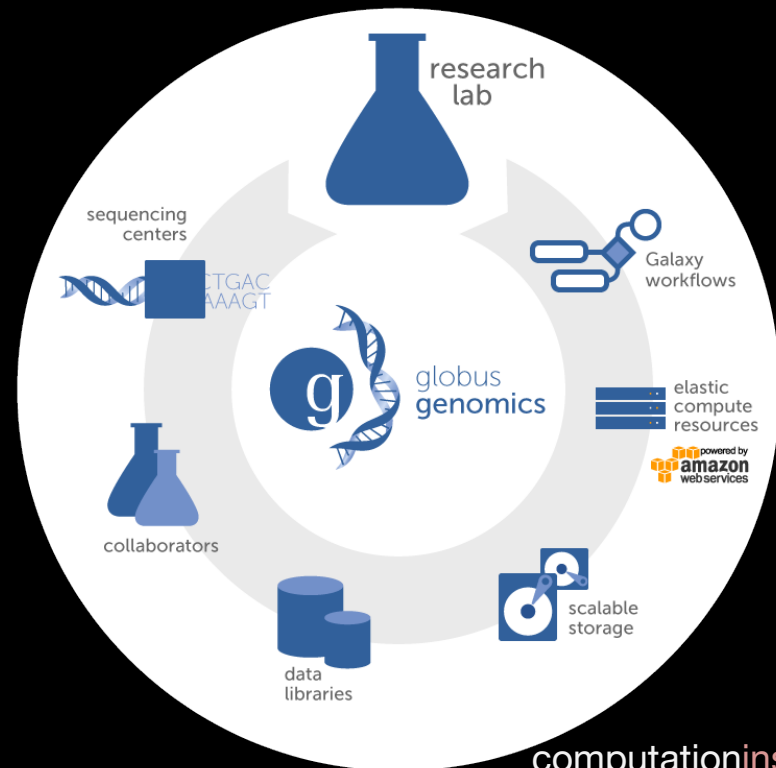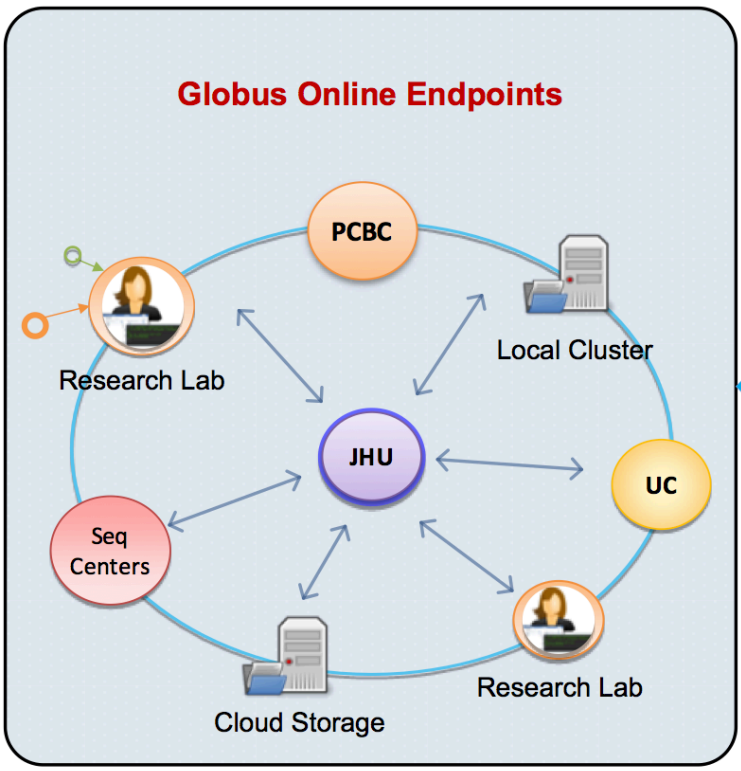
# Globus Online already does a lot

# Data management SaaS (Globus) +
# Next-gen sequence analysis pipelines (Galaxy) +
# Cloud IaaS (Amazon) =

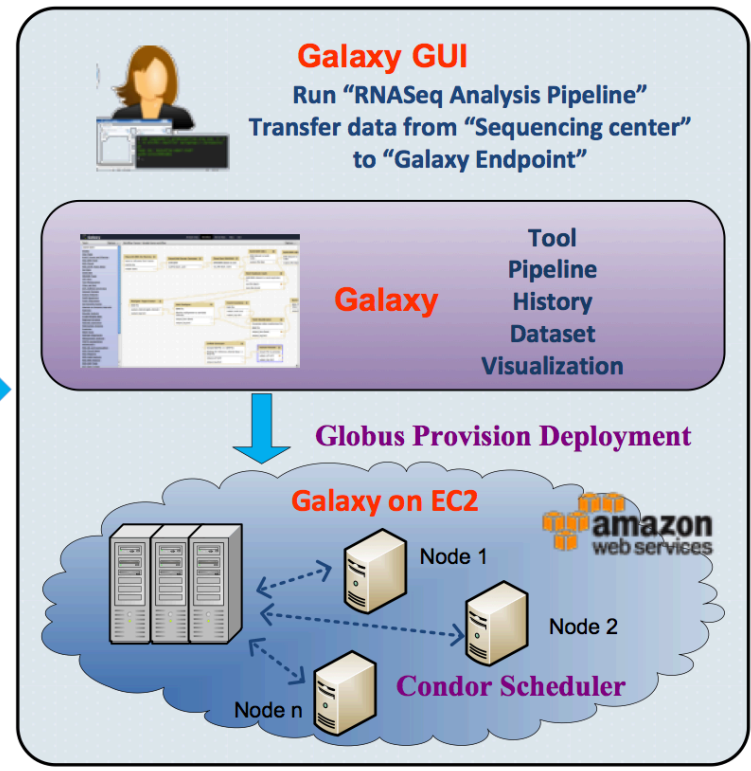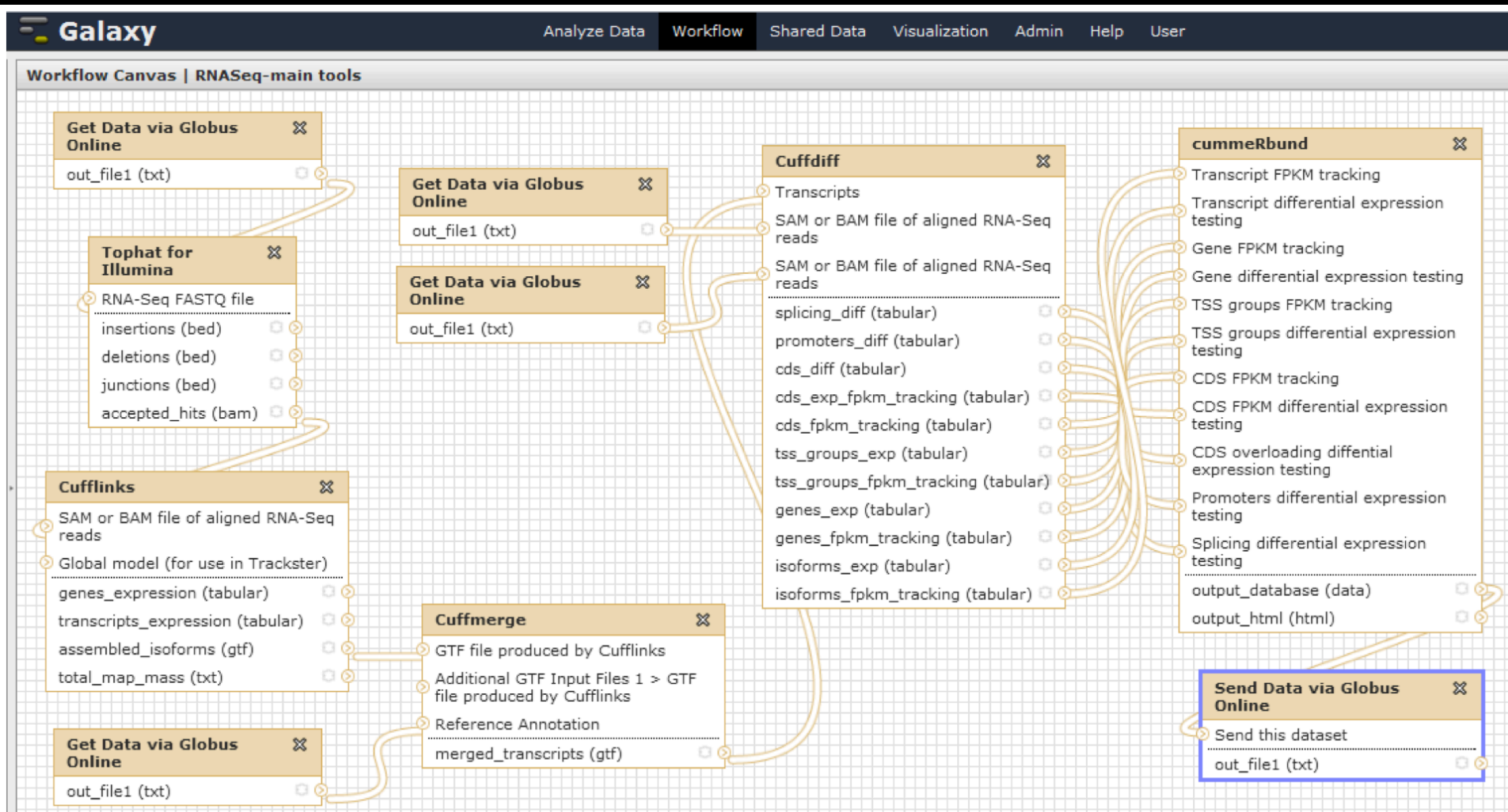**Flexible, scalable, easy-to-use genomics analysis for all biologists**



globus genomics

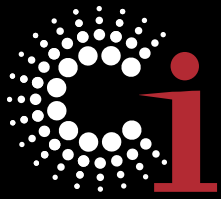Ravi Madduri, Bo Liu, Paul Davé, et al.

computationinstitute.org

# RNA-Seq pipeline

# Amazon pricing for Diffusion Tensor Imaging pipeline



Credit: Kyle Chard

computationinstitute.org

# A platform for integration

# A platform for integration

A platform for integration

# Expanding Globus Online services

- Ingest and publication
  - Imagine a DropBox that not only replicates, but also extracts metadata, catalogs, converts

- Cataloging
  - Virtual views of data based on user-defined and/or automatically extracted metadata

- Computation
  - Associate computational procedures, orchestrate application, catalog results, record provenance

# Looking deeply at how researchers use data

- A single research question often requires the integration of many data elements, that are:
  - In different locations
  - In different formats (Excel, text, CDF, HDF, …)
  - Described in different ways
- Best grouping can vary during investigation
  - Longitudinal, vertical, cross-cutting
- But always needs to be operated on as a unit
  - Share, annotate, process, copy, archive, …

# How do we manage data today?

- Often, a curious mix of ad hoc methods
  - Organize in directories using file and directory naming conventions
  - Capture status in README files, spreadsheets, notebooks
- Time-consuming, complex, error prone

Why can't we manage our data like we manage our pictures and music?

# flickr

Home    You ▾    Organize ▾    Contacts ▾    Groups ▾    Explore ▾                Search ▾

## Hi

## » Your Photostream pro

▾ Recent Uploads  |  Recent Activity

## » Your Contacts

▸ NEW There are new uploads from your contacts.

## » Your Groups                                        Groups activity

▾ Canon DSLR User Group   ( 415,135 items | 6,803 topics )

From Eastanbul    From Bidoll    From rjptn    From DAVÍ    From harold.lloyd

More: photophlow, Nikon D50 Users, We Demand Donuts: April 16 was the 1st
Annual Day of the Donut!, Canon EOS-1Ds Mark III, Nikon DSLR Users, more...

## » Upload Photos & Videos

### Flickr Blog                          Posted 09 Sep 08

**Kitten Tuesday**
It's a very special Kitten Tuesday (back
story here and here). Team Flickr
would like to congratulate Dan and
Charlie Catt on the arrival of...

### Add your photos to a map

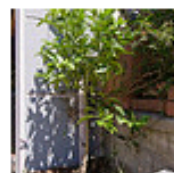Make a note of where you were, and add to the world
map!

### Make, share, and sell books with Blurb                    ✕

It's easy to make, share, and sell
your books with Blurb. Check out
what other people are doing – visit
Blurb's Flickr group.

And even more you can do with your photos:

- Capital One Personalize your credit card NEW
- HP: Prints, Photocubes, Posters and Books

# Introducing the **dataset**

- **Group** data based on use, not location
  - Logical grouping to organize, reorganize, search, and describe usage

- **Tag** with characteristics that reflect content …
  - Capture as much existing information as we can

- …or to reflect current status in investigation
  - Stage of processing, provenance, validation, ..

- **Share** data sets for collaboration
  - Control access to data and metadata

- **Operate** on datasets as units
  - Copy, export, analyze, tag, archive, …

# Builds on catalog as a service

## Approach

- Hosted user-defined catalogs
- Based on tag model
    <subject, name, value>
- Optional schema constraints
- Integrated with other Globus services

## Three REST APIs

**/query/**

- Retrieve subjects

**/tags/**

- Create, delete, retrieve tags

**/tagdef/**

- Create, delete, retrieve tag definitions

Builds on USC Tagfiler project (C. Kesselman et al.)
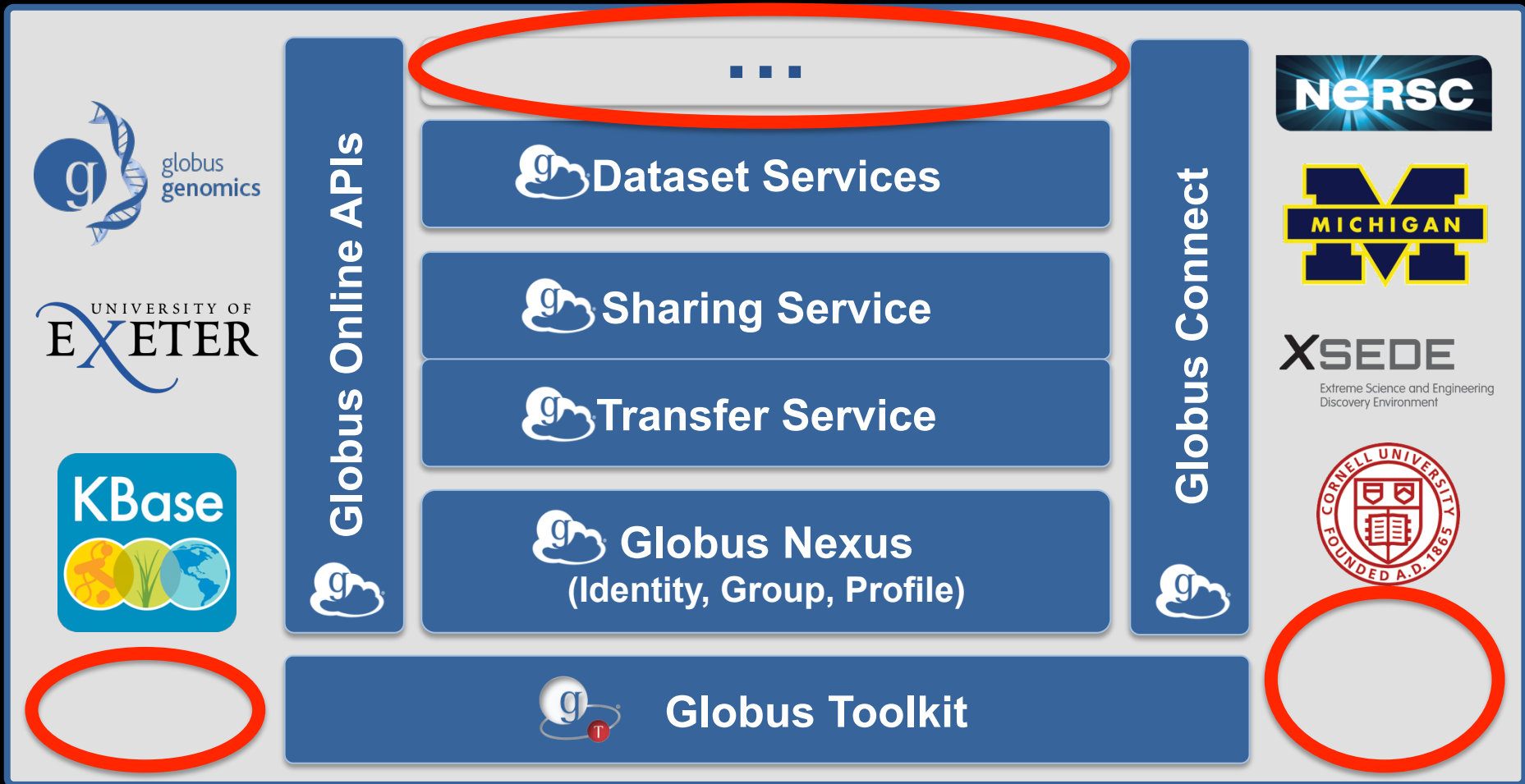
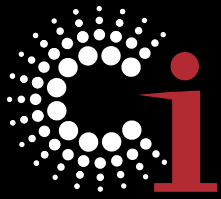# Our vision for a 21st century discovery infrastructure

Provide **more** capability for **more** people at **lower cost** by building a **"Discovery Cloud"**
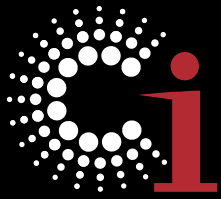
**Delivering "Science as a service"**

# It's a time of great opportunity …
## to develop and apply Science aaS

# Thanks to great colleagues and collaborators

- Steve Tuecke, Rachana Ananthakrishnan, Kyle Chard, Raj Kettimuthu, Ravi Madduri, Tanu Malik, and many others at Argonne & Uchicago

- Carl Kesselman, Karl Czajkowski, Rob Schuler, and others at USC/ISI

- Francesco de Carlo, Chris Jacobsen, and others at Argonne

Thank you to our sponsors!

U.S. DEPARTMENT OF ENERGY

NSF

THE UNIVERSITY OF CHICAGO

NATIONAL INSTITUTES OF HEALTH

Argonne NATIONAL LABORATORY

computationinstitute.org