# Data Mining and Grid

## Ian Foster
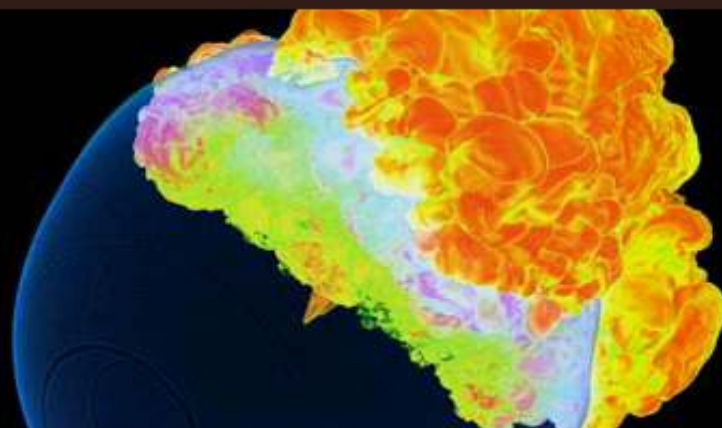
Computation Institute

Argonne National Lab & University of Chicago

http://ianfoster.typepad.com

Data
Mining

Grid

Power
substation

High voltage
transmission lines

Transmission substation

Power pla

Transformer

Power
poles

Transformer

# In the Next 50 Years, We Must …

- Increase energy production by 5, while reducing GHG emissions by 2 or more

- Mitigate and adapt to climate change

- Address increasingly drug resistant diseases

- Provide meaningful livelihoods for 9B people

➔ Innovation

# Innovation
# as a Systems Problem

- Quasi-ubiquitous Internet …

- … connects many potential innovators

  - ◆ Millions of scientists, billions of people

- Who need to leverage

  - ◆ Enormous data of tremendous complexity

  - ◆ Immensely powerful computing

  - ◆ Experimental apparatus of great power

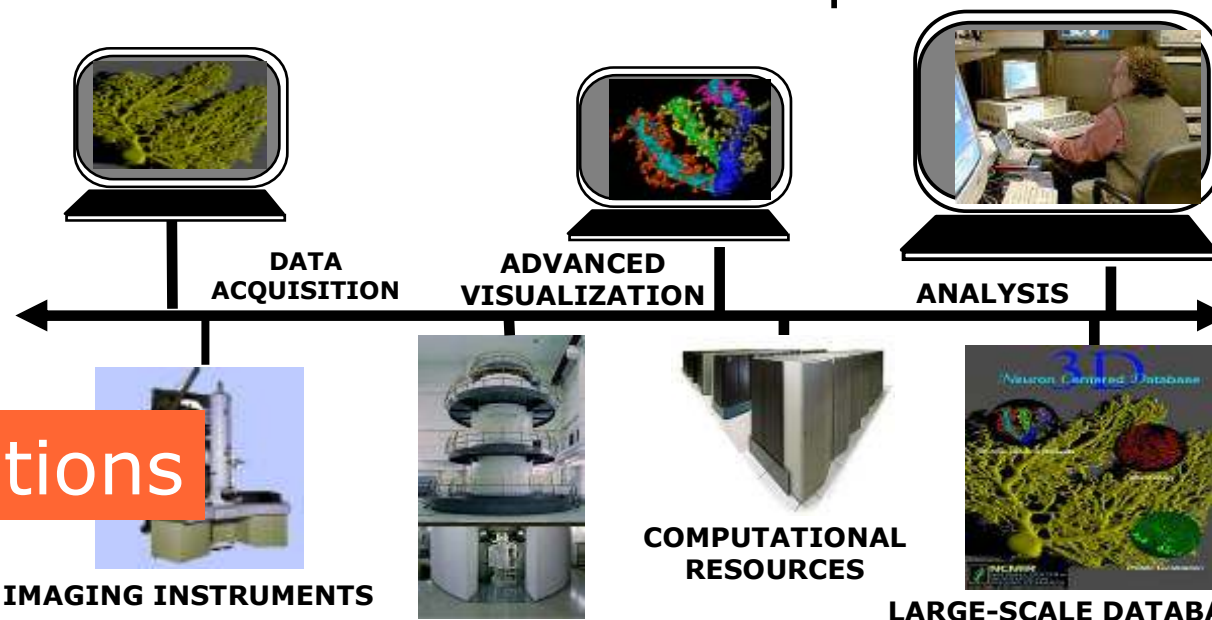➔ We must address problem solving as an distributed, end-to-end, systems problem

# Grid:
# A Unifying Concept & Technology

## Grid enables the **federation** of resources

- Distributed computers, storage, data, people, …
- **Networks** provide connectivity
- **Software** & **standards** provide the "glue"
- **Infrastructure** services facilitate operation

Infrastructure

DATA
ACQUISITION

ADVANCED
VISUALIZATION

ANALYSIS

Applications

Research

IMAGING INSTRUMENTS
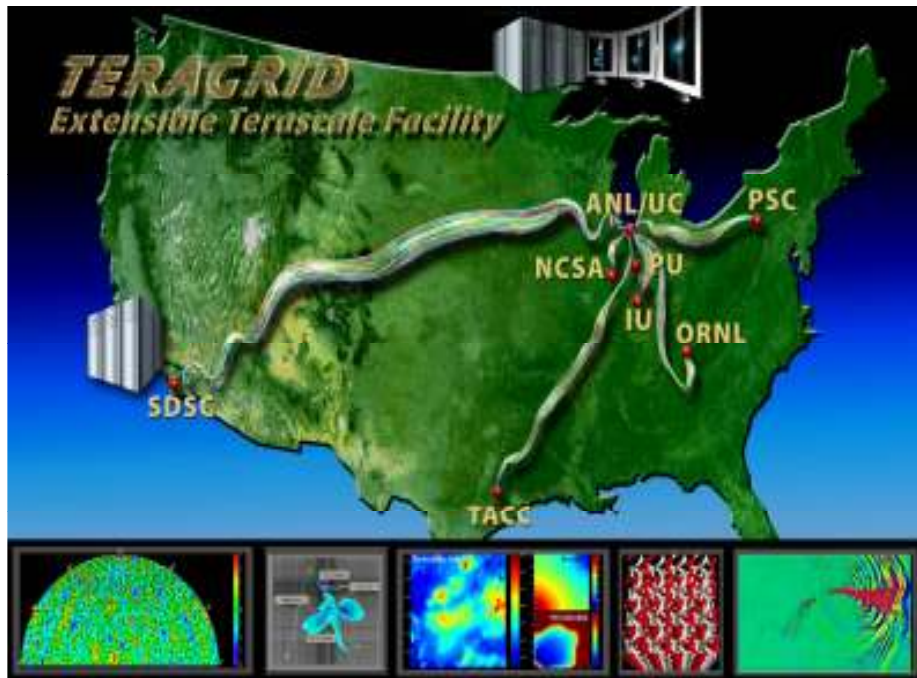
COMPUTATIONAL
RESOURCES

LARGE-SCALE DATABASES

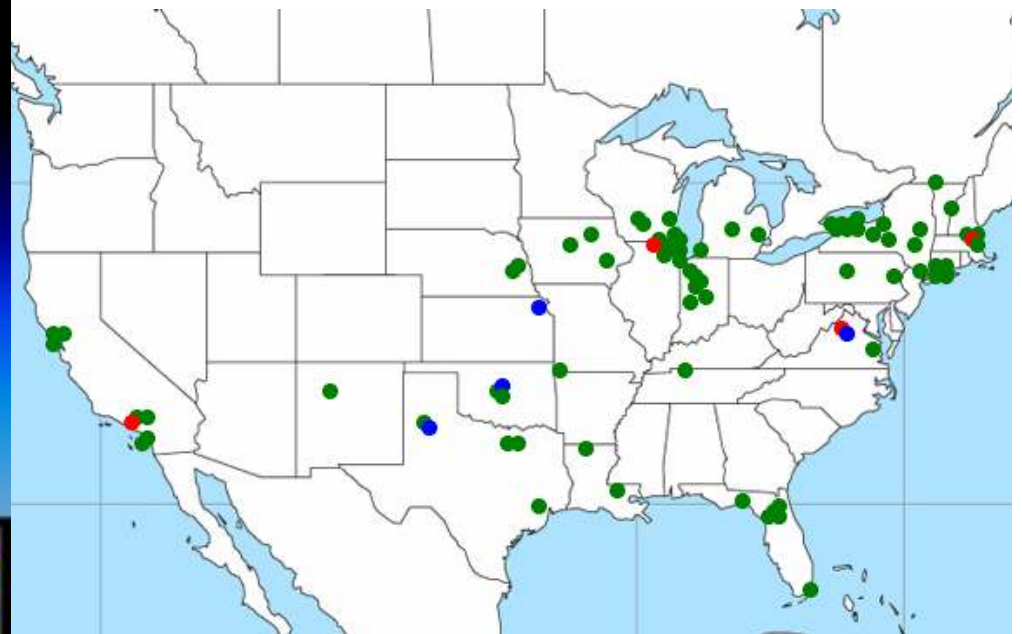Credit: Mark Ellisman

# Grid Infrastructure

- Massive computing and storage
- Service interfaces facilitate access and use
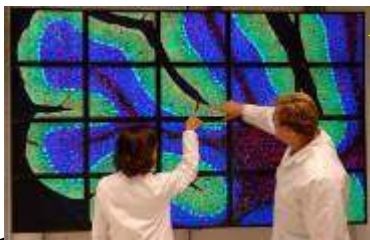


TeraGrid



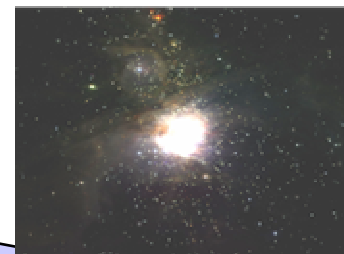Open Science Grid

# Software and Standards

**Domenico Talia**

**Bob Grossman**

Weka 4WS

Angle

Tool

Tool

File Transfer

User Svc

Host Env

Uniform interfaces, security mechanisms, Web service transport, monitoring

Registry

Globus

GRAM

User Svc

Host Env

GridFTP

DAI

Computers

Specialized resource

File system

Database

8

# Globus Downloads Last 24 Hours



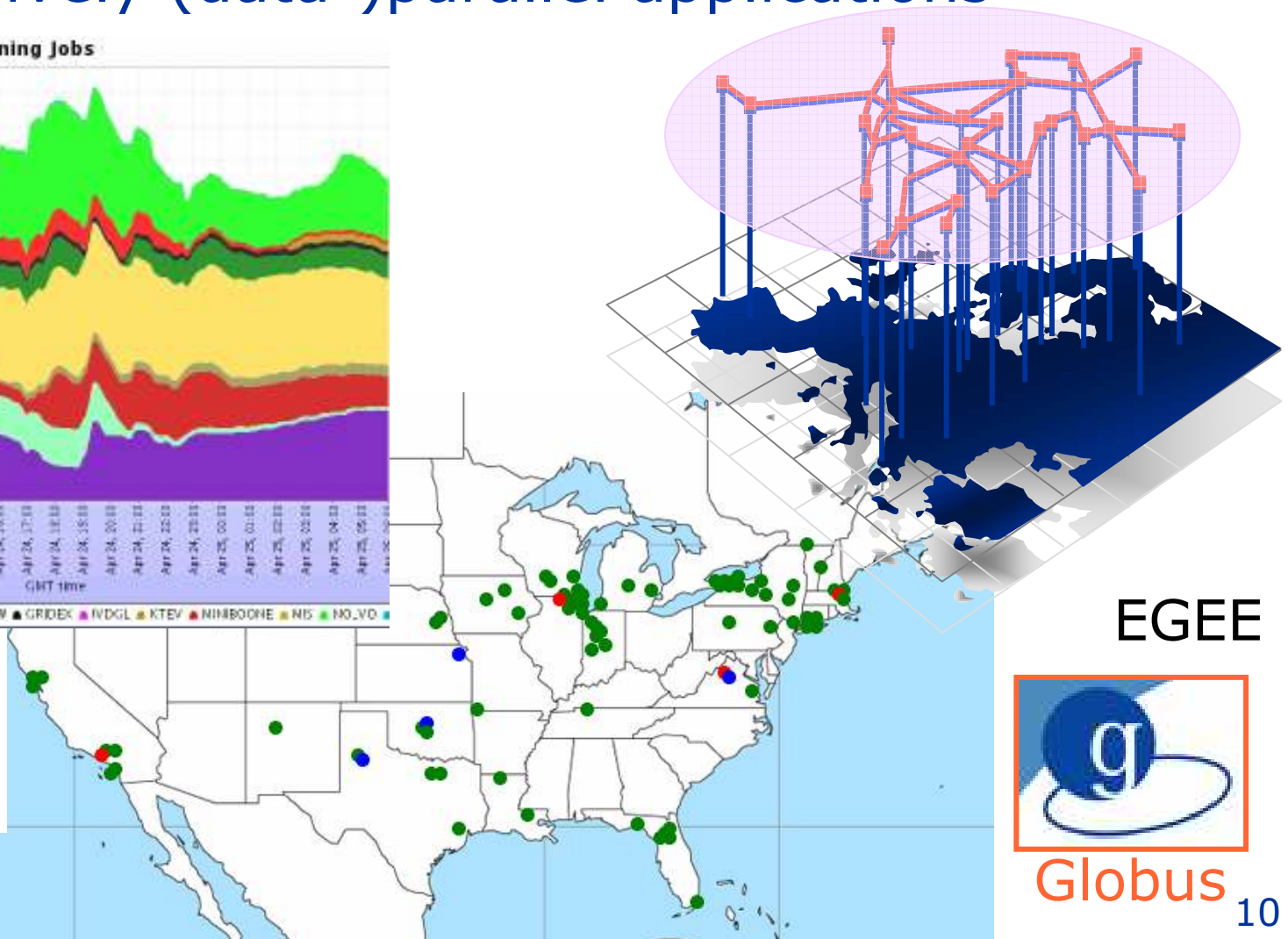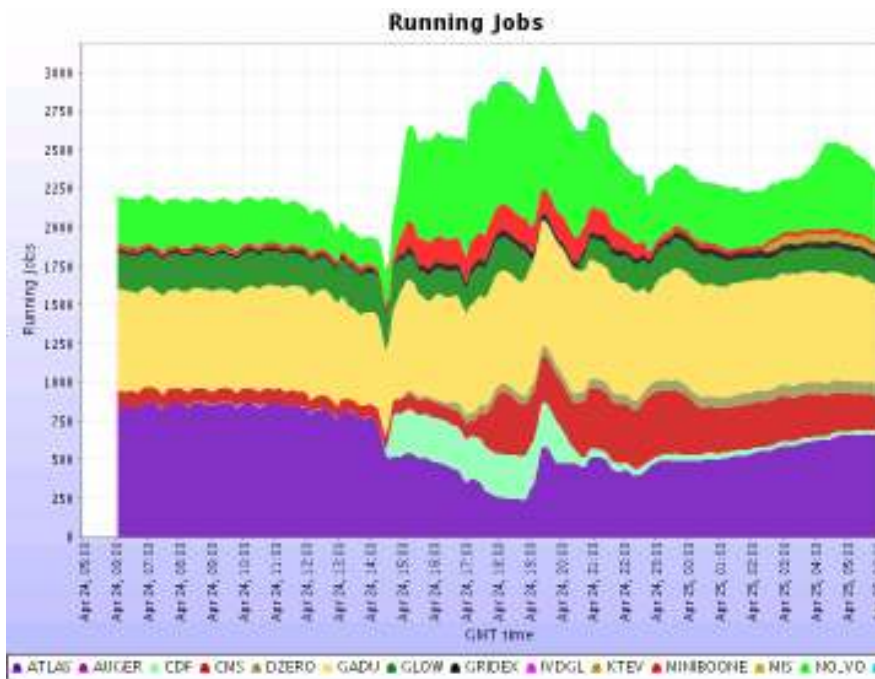**Last month**

# First Generation Grids:
# On-Demand/Batch Computing
## Focus on aggregation of many resources for massively (data-)parallel applications



Running Jobs

ATLAS  AUGER  CDF  CMS  DZERO  GADU  GLOW  GRIDEX  IVDGL  KTEV  MINIBOONE  NIS  NO_VO
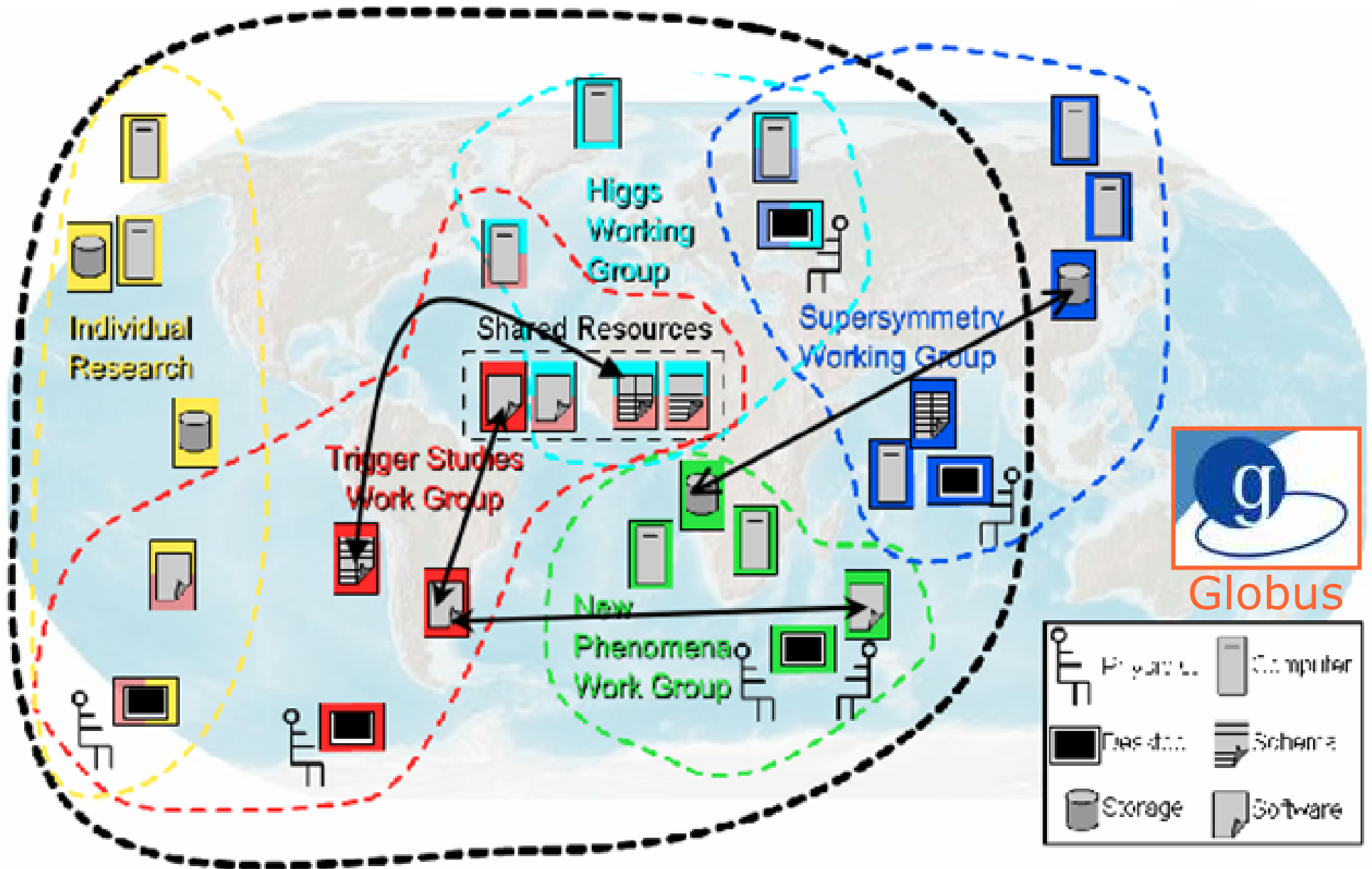
Open Science Grid

EGEE

Globus

# Applications:
# High Energy Physics

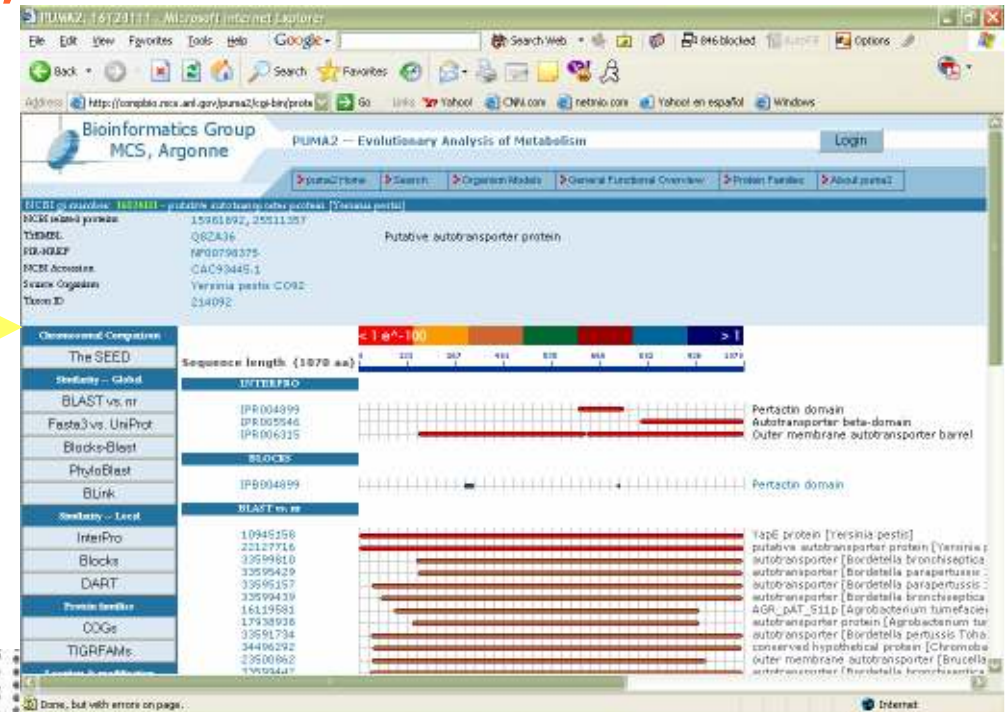# Integrating Data and Computing, on Demand

## Public PUMA Knowledge Base

Information about proteins analyzed against ~2 million gene sequences



## Back Office Analysis on Grid

Millions of BLAST, BLOCKS, etc., on OSG and TeraGrid

Natalia Maltsev et al., http://compbio.mcs.anl.gov/puma2

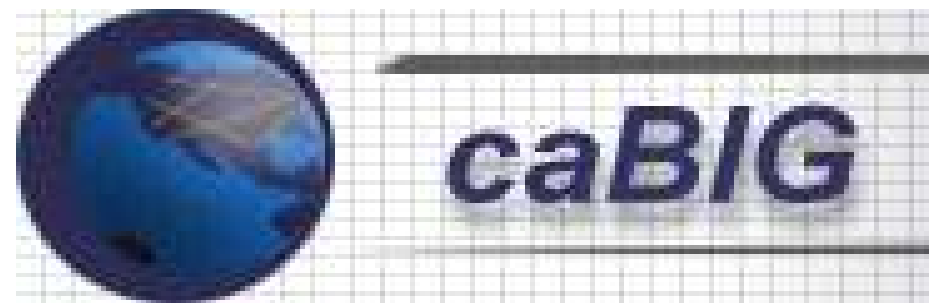# Second Generation Grids: Service-Oriented Science

- Empower many more users by enabling on-demand access to **services**
- Grids become an enabling technology for **service oriented science** (or business)
  - ◆ Grid infrastructures host services
  - ◆ Grid technologies used to build services
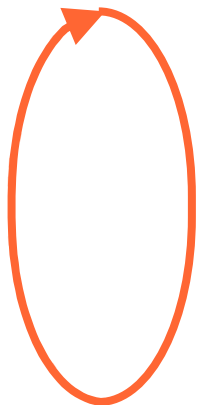
*Science Gateways*

**TeraGrid™**
EMPOWERING DISCOVERY

caBIG

"Service-Oriented Science", *Science*, 2005

# Service-Oriented Science

People **create** services (data or functions) …

which I **discover** (& decide whether to use) …

& **compose** to create a new function …

& then **publish** as a new service.

‼️

→ *I find "someone else" to **host** services, so I don't have to become an expert in operating services & computers!*

→ *I hope that this "someone else" can* **manage** *security, reliability, scalability, …*
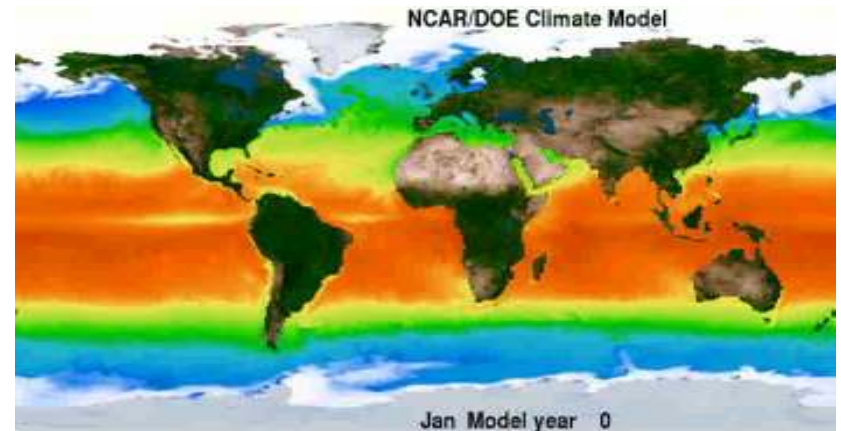
# Earth System Grid

- **On-demand access to climate simulation data**
  - ◆ Multiple archives
  - ◆ Interactive query
  - ◆ Per-collection control
  - ◆ Server-side processing
- **Major scientific impact**
  - ◆ >5000 users
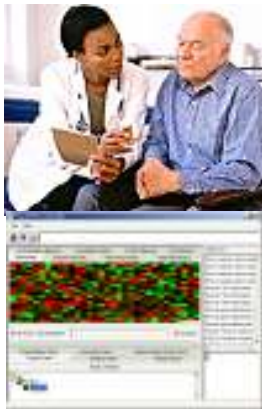  - ◆ >200 TB downloaded
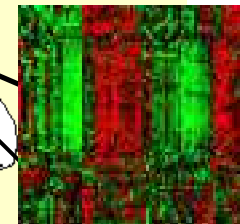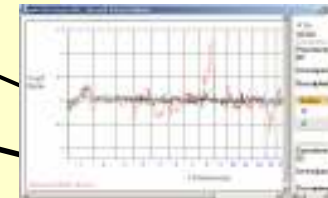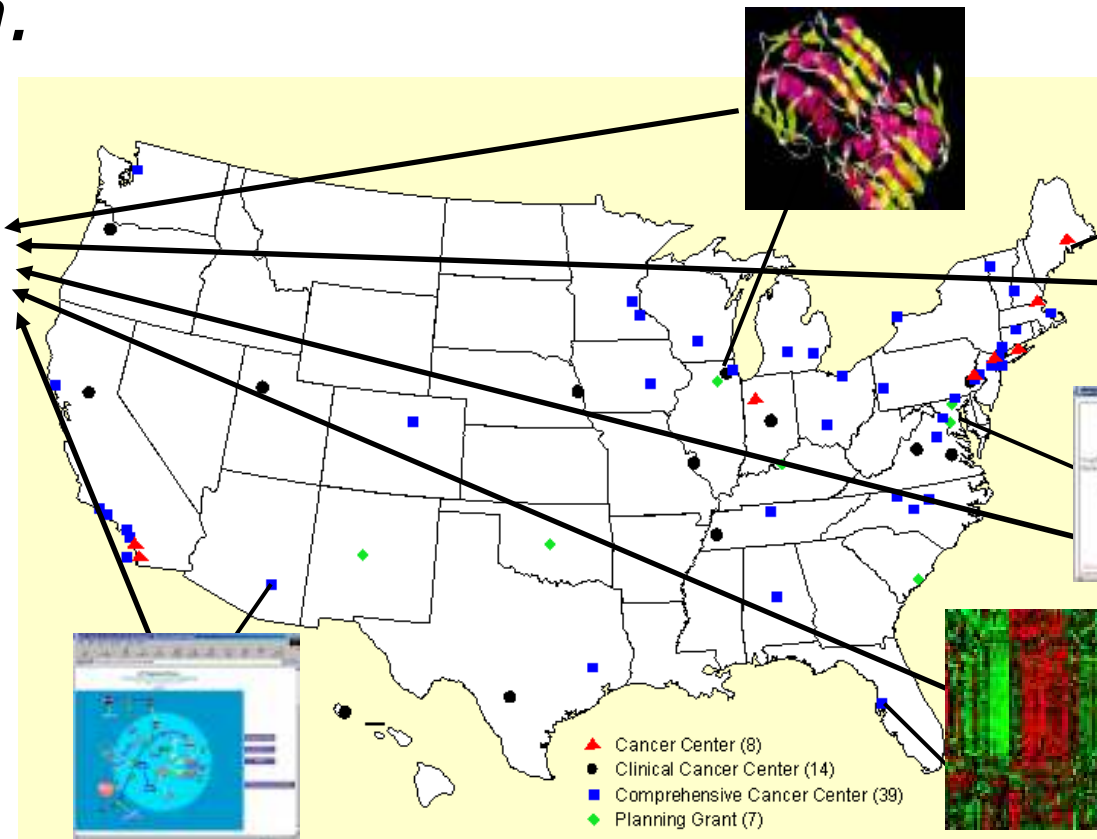  - ◆ >300 scientific papers



Globus

# Cancer Biomedical Informatics Grid (caBIG)

*caBIG: sharing of infrastructure, applications, and data.*



**Data Integration!**

Cancer Center (8)
Clinical Cancer Center (14)
Comprehensive Cancer Center (39)
Planning Grant (7)
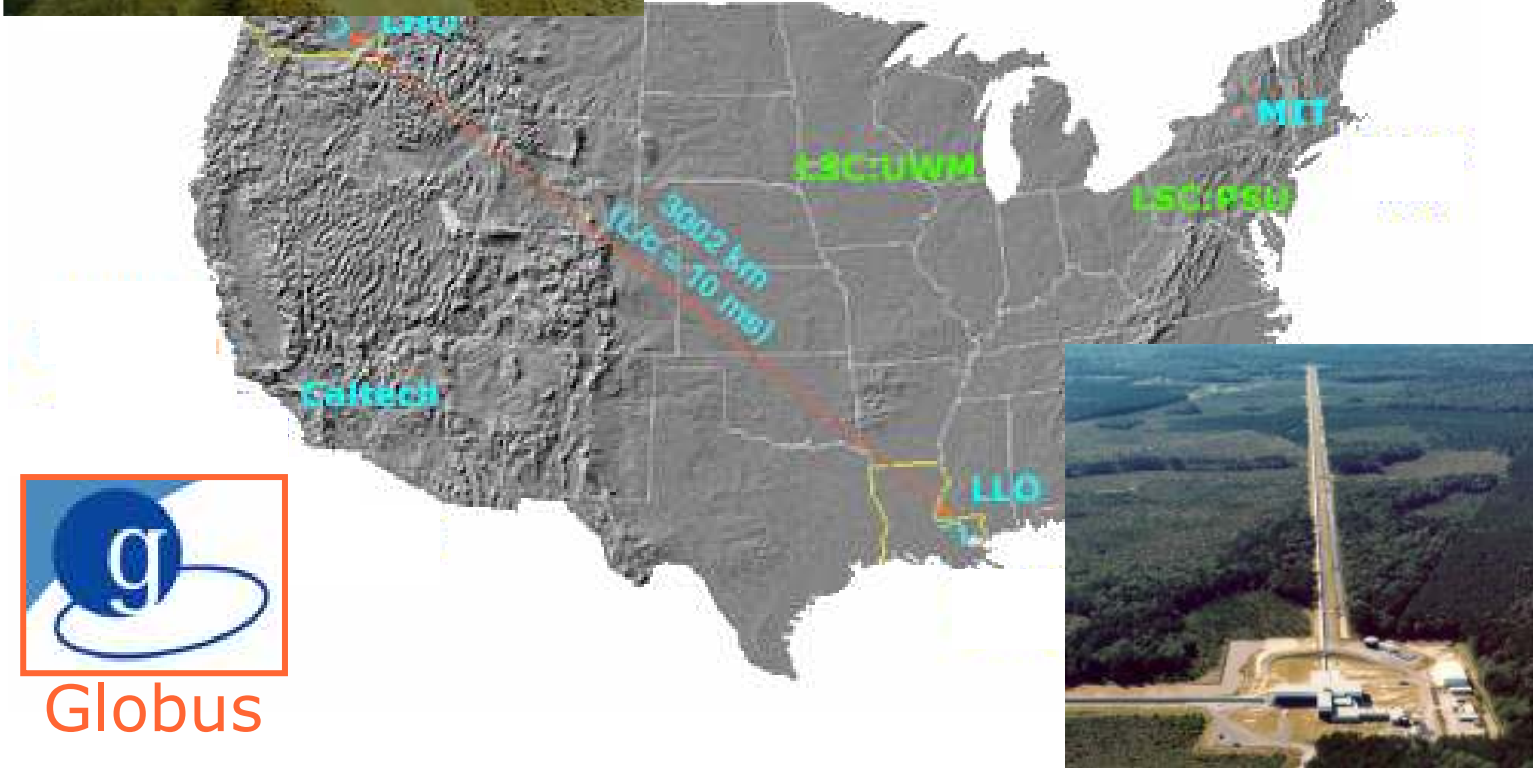
Globus

# caBIG Under the Covers

# LIGO Data Grid

## LIGO Gravitational Wave Observatory

Globus

Replicating >1 Terabyte/day to 8 sites
>150 million replicas so far
MTBF = 1 month   www.globus.org/solutions

The Angle Project

# Social Informatics Data Grid

Global Observation Database (View)

VCR-Style Control Panel

Globus

Animated Text Transcript (Paragraph Representation)

Tag Transcript Editor

Animated Avatar Representation

Animated Graph Panes

Video Displays

Video List

Bennett Berthenthal et al., www.sidgrid.org

# A Few Example Research Themes

- Service discovery, composition, provisioning
  - ◆ SOA, virtualization, cloud computing, …
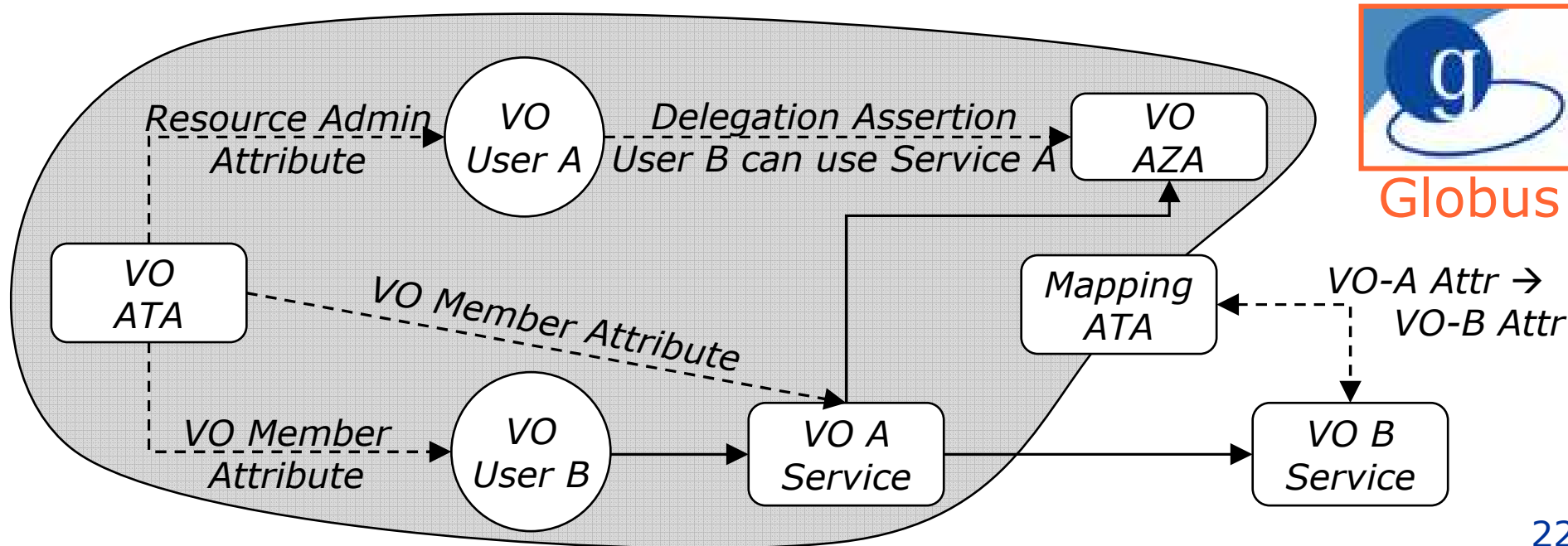- Large-scale (distributed) computation
  - ◆ E.g., Swift, Kepler, Taverna
- Provenance
  - ◆ E.g., "Provenance Challenge"
- "Virtual organizations"
  - ◆ E.g., attribute-based authorization, trust
- Integration of physical systems
  - ◆ Optimization of end-to-end workflows
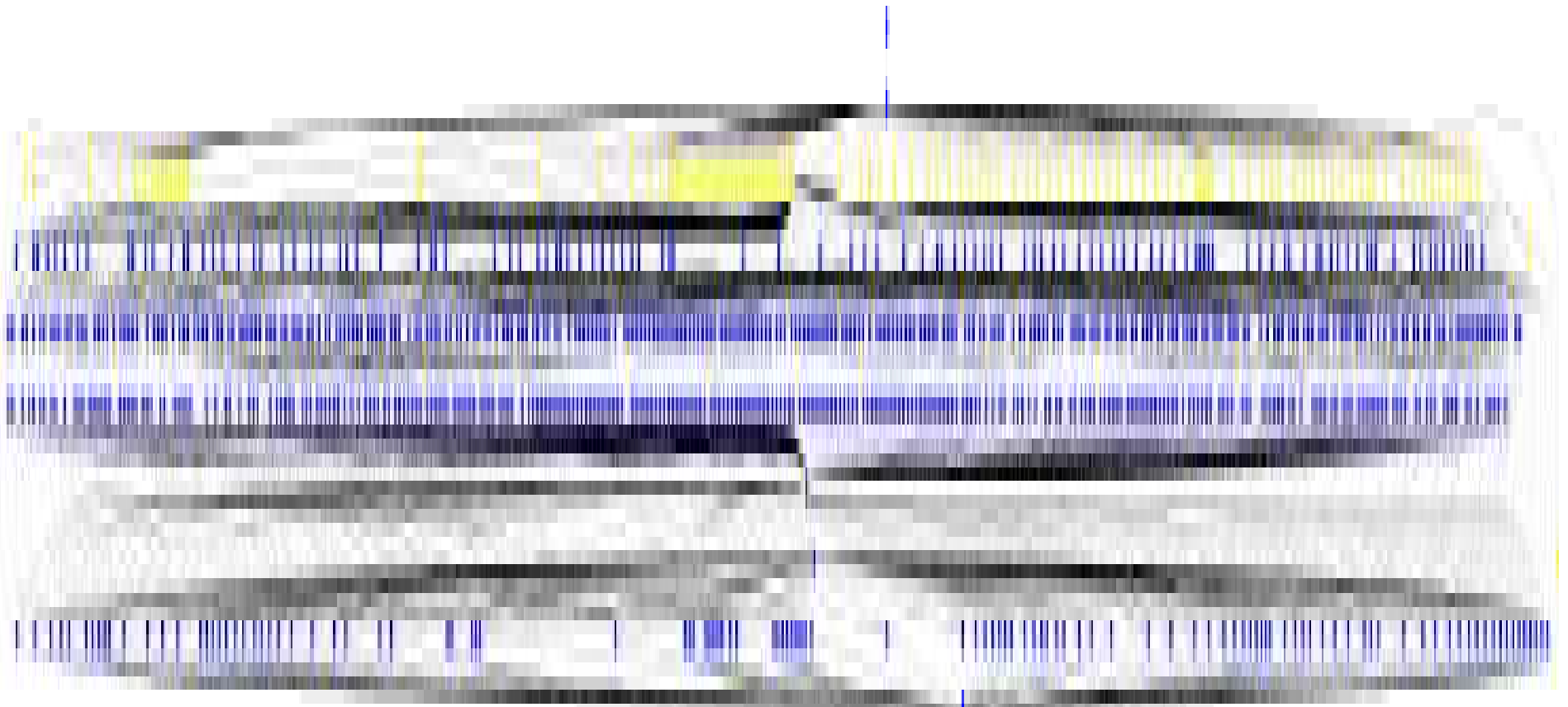
# Security Services for Virtual Organization Policy

- Attribute Authority (ATA)
  - Issue signed attribute assertions (incl. identity, delegation & mapping)
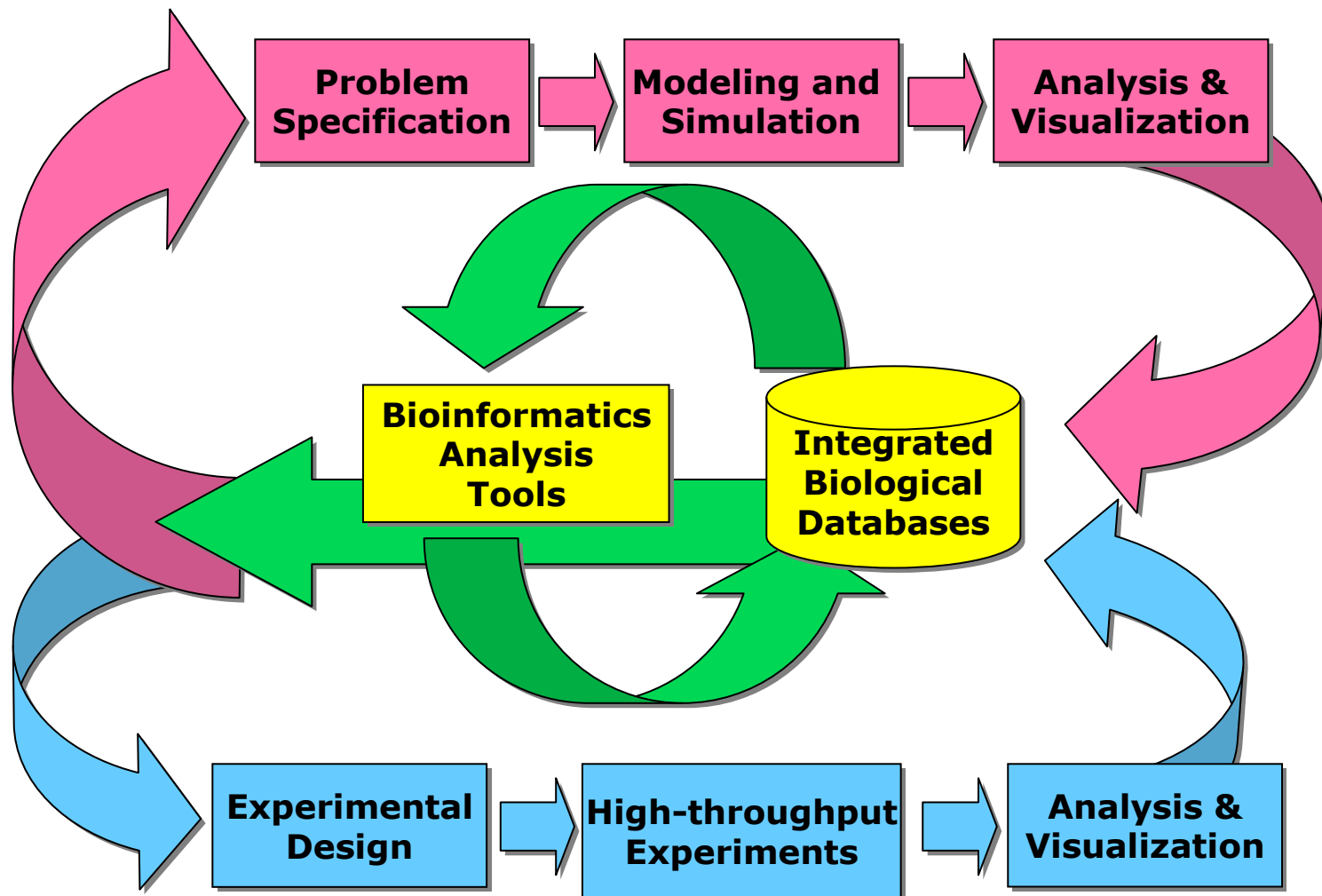- Authorization Authority (AZA)
  - Decisions based on assertions & policy

Globus

# Swift
## (www.ci.uchicago.edu/swift)

# An Integrated View of Modeling, Simulation, Experiment, & Informatics

# Robot Scientist


Biomek 200

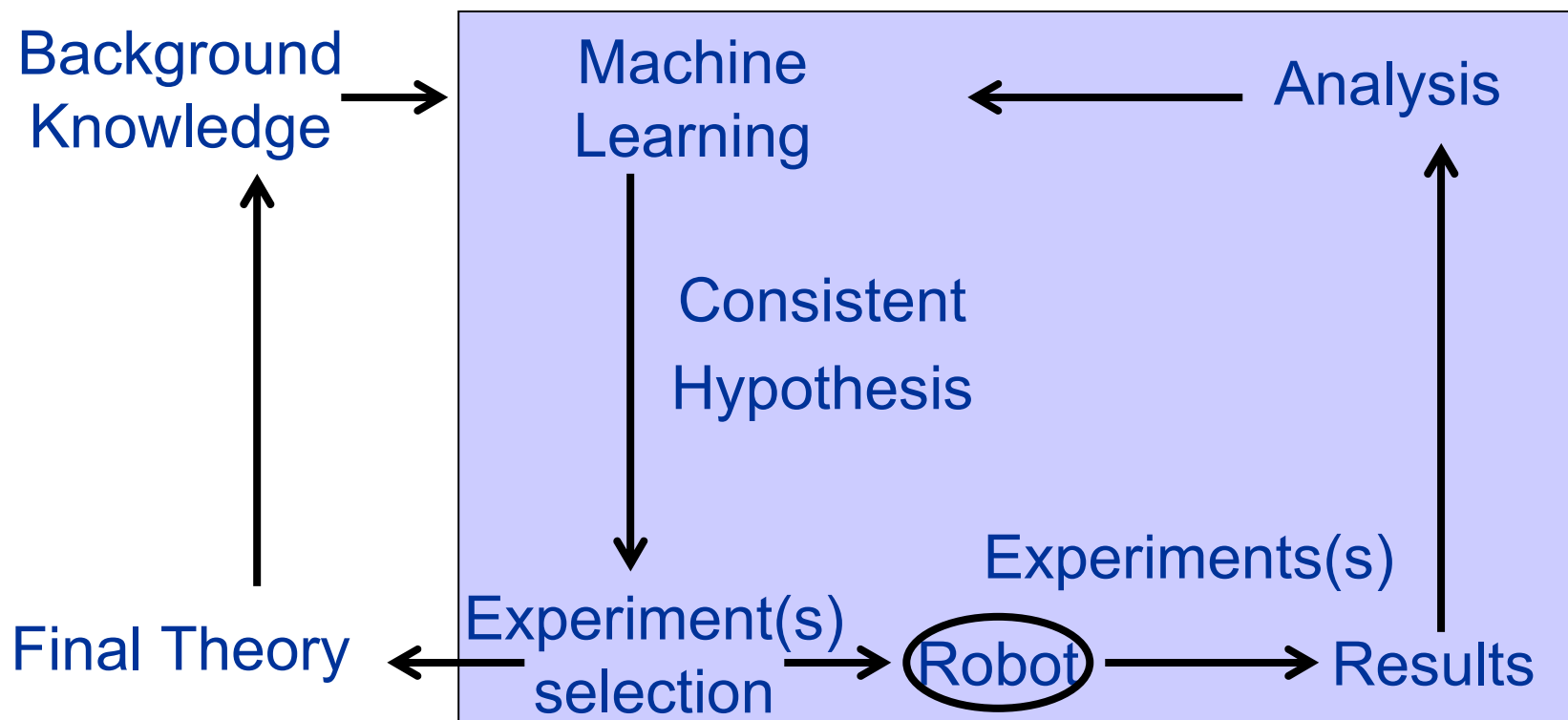"The robot scientist project aims to develop a computer system capable of originating its own experiments, physically doing them, interpreting the results, & then repeating the cycle."

Background Knowledge → Machine Learning ← Analysis

Consistent Hypothesis

Final Theory ← Experiment(s) selection → Robot → Results

Experiments(s)

Stephen Muggleton, Ross King et al., UK

# Team Science meets Data Deluge