# iRODS

*Executive Overview*
*August 12, 2014*

# Agenda

- What is iRODS?
- Who Uses iRODS?
- What Can iRODS Do?
- The Future of iRODS

# WHAT IS iRODS?

# What is iRODS?

iRODS is open source data grid middleware for...

- Data Discovery
- Workflow Automation
- Secure Collaboration
- Data Virtualization

# iRODS is Open Source

iRODS is open source data grid middleware for...
- Data Discovery
- Workflow Automation
- Secure Collaboration
- Data Virtualization

iRODS
— CONSORTIUM —

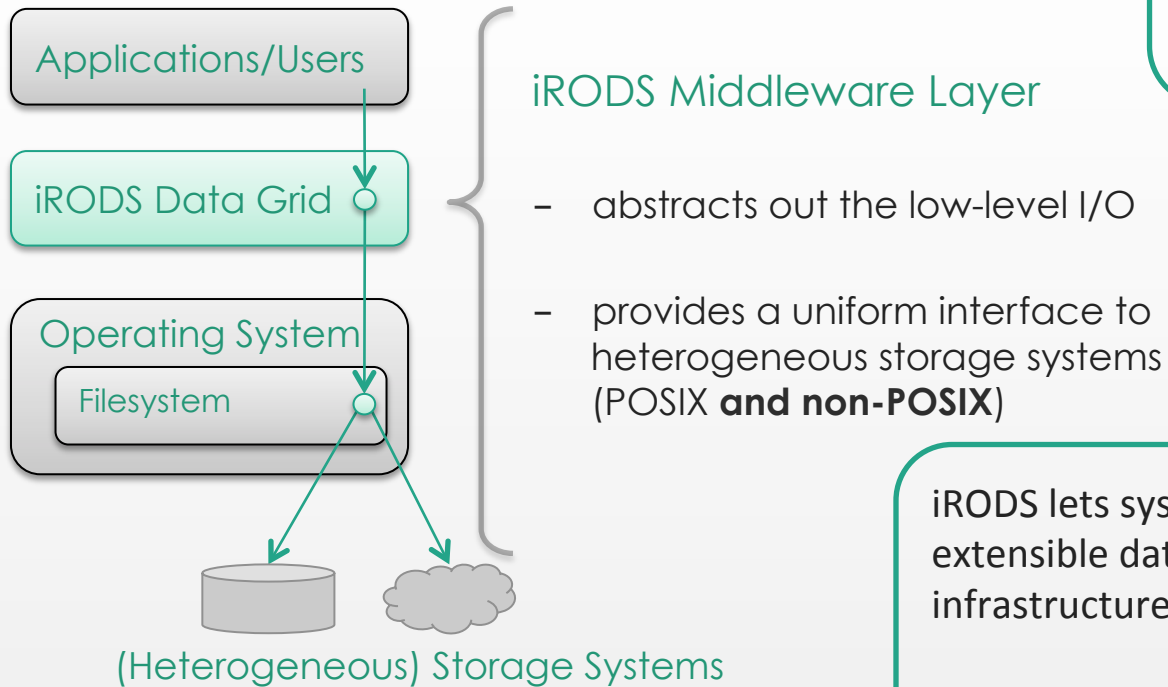The iRODS Consortium exists to ensure the sustainability of iRODS by:
- Ensuring that iRODS source code remains freely available for use and modification.
- Promoting the adaptation of iRODS to a variety of hardware and software platforms.
- Supporting continued development of core iRODS features.
- Facilitating interaction among members of the iRODS developer community.
- Providing a forum for key stakeholders to guide ongoing development of iRODS.

iRODS
CONSORTIUM

# iRODS is Middleware

**iRODS is** open source data grid **middleware** for...
- Data Discovery
- Workflow Automation
- Secure Collaboration
- Data Virtualization

**mid•dle•ware** `midl,we(ə)r

*noun* software that acts as a bridge between an operating system or database and applications, especially on a network

Applications/Users

iRODS Data Grid

Operating System

Filesystem

iRODS Middleware Layer

– abstracts out the low-level I/O

– provides a uniform interface to heterogeneous storage systems (POSIX **and non-POSIX**)

iRODS lets system administrators roll out an extensible data grid **without** changing their infrastructure.

Data is accessed using familiar APIs.
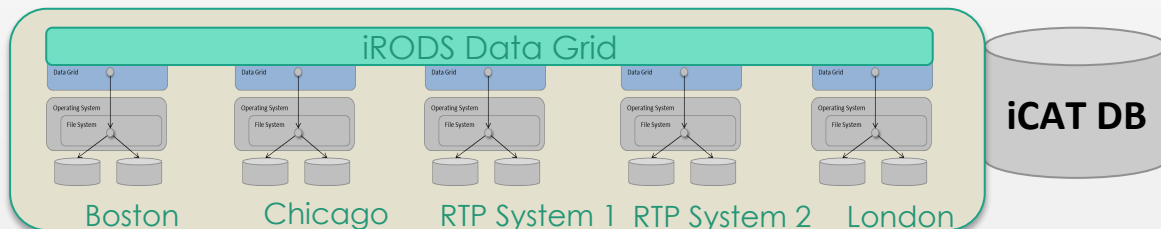
(Heterogeneous) Storage Systems

# Data Discovery

iRODS is open source data grid middleware for...

- **Data Discovery**
- Workflow Automation
- Secure Collaboration
- Data Virtualization

Every iRODS data grid also has a metadata catalog, called the iCAT. The iCAT is used by iRODS to locate data, manage provenance, and to enable automation and access control.

The iCAT also permits user-defined metadata. Altogether, this metadata supports:
- Data discovery based on parameters such as user-defined tags, modification date, outcomes of automation activity.
- Capturing workflows as raw data is processed and used.
- Automation and access control policies.

**Example Metadata:**

*Logical Name (iRODS path):*
*/RDDept/LabX/Flow/Study1*

*Physical Name (Unix path):*
*/London/var1/proj/labx/stuff*

*Lab PI:  Jane Doe*
*Date:  12/1/2010*
*Time:  01:45:12*

*Title:*
*Proliferation optimization studies*

*Data Source:  Flow Cytometer*
*Assay Conditions:  Data captured*
*…*



**iRODS Data Grid**

Boston     Chicago     RTP System 1   RTP System 2   London

**iCAT DB**

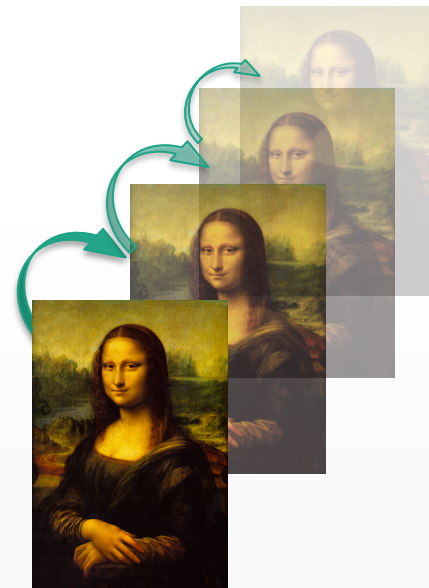# Workflow Automation

iRODS is open source data grid middleware for...
- Data Discovery
- Workflow Automation
- Secure Collaboration
- Data Virtualization

With iRODS, **any agent** can initiate **any action** upon **any trigger.**

This powerful capability allows administrators to automate policies such as:
- Validating checksums every time a new file is placed in a folder.
- Backing up a set of files every second Thursday.
- Archiving data that hasn't been accessed in over 1 month.
- Logging each time a file is replicated or destroyed.
- Permitting a file to be accessed by multiple independently defined user groups.

These operations can be **distributed** to the storage resource or client.
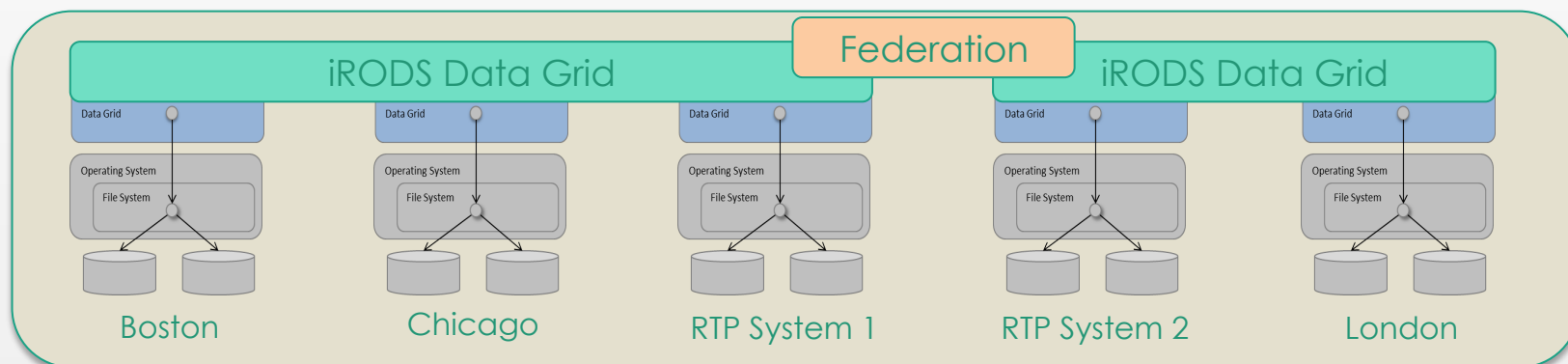
# Secure Collaboration

iRODS is open source data grid middleware for…
- Data Discovery
- Workflow Automation
- Secure Collaboration
- Data Virtualization

iRODS presents centralizes distributed storage systems under a unified namespace.

Administrators can control how the grid is presented to users and implement replication, load-distribution, and archiving policies that are completely transparent to the user.
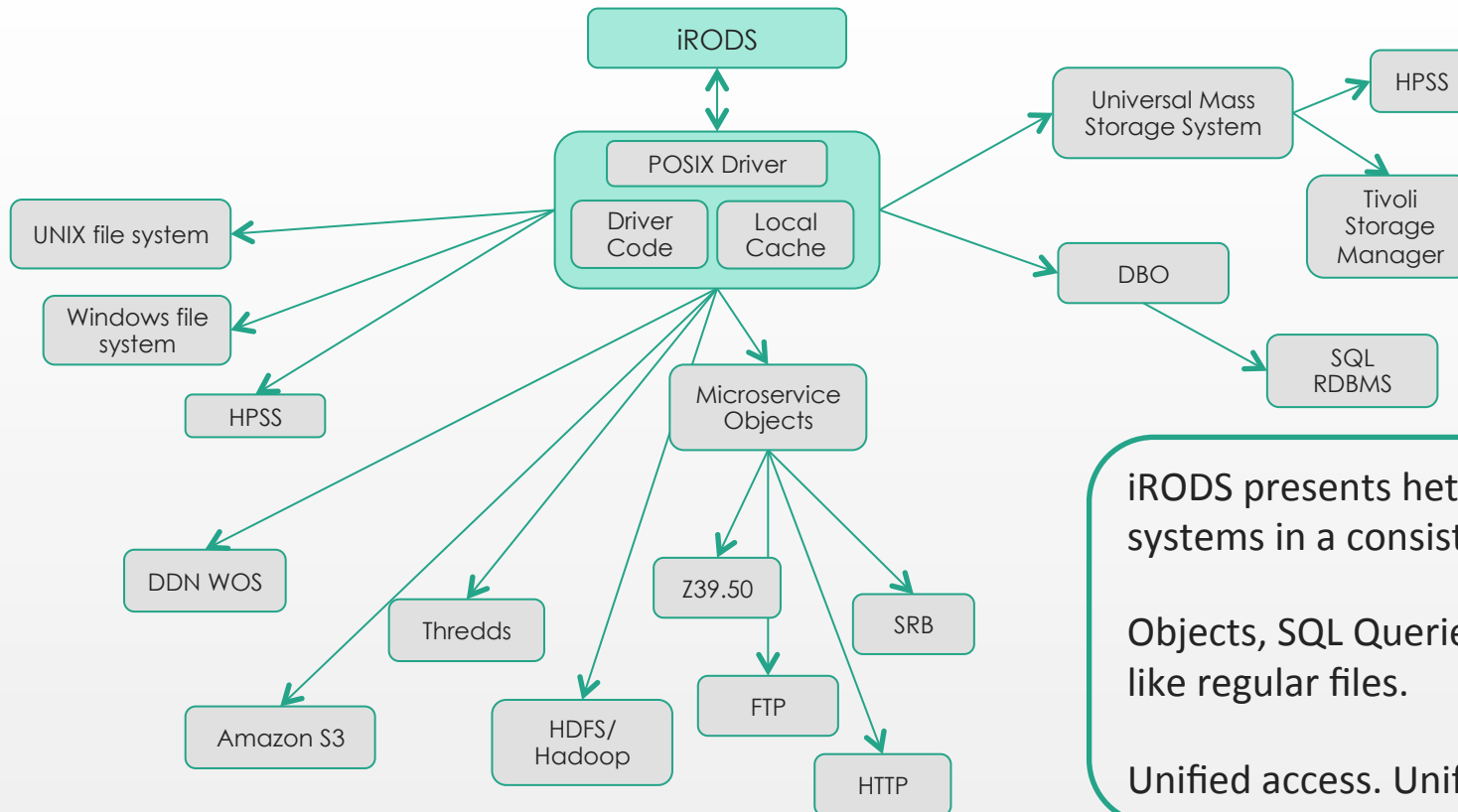
Independent grids can be federated with one another to allow controlled access to remote grids or grids operated by separate workgroups.

# Data Virtualization

iRODS is open source data grid middleware for...
- Data Discovery
- Workflow Automation
- Secure Collaboration
- Data Virtualization



iRODS presents heterogeneous storage systems in a consistent, familiar format.

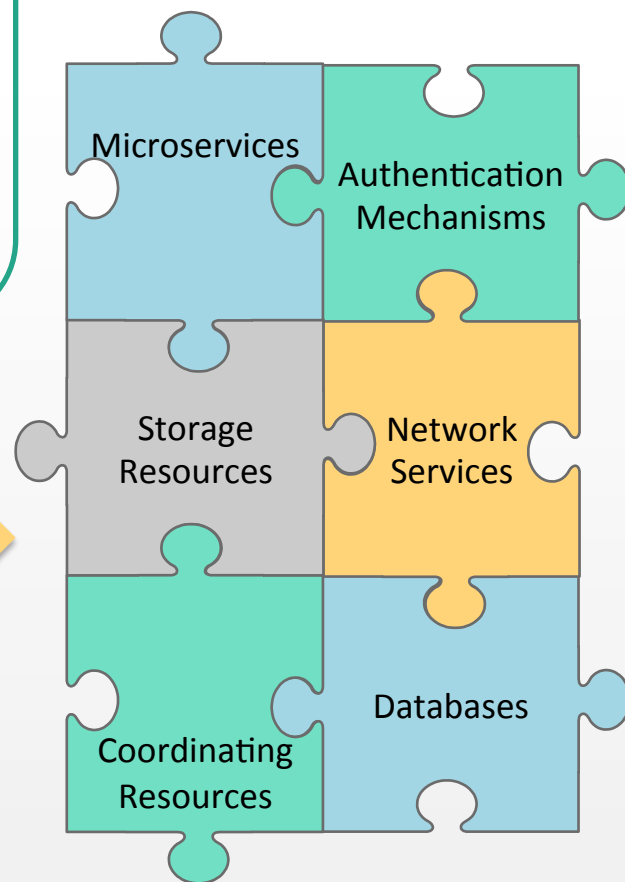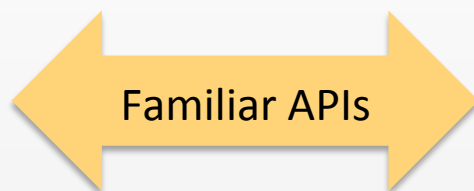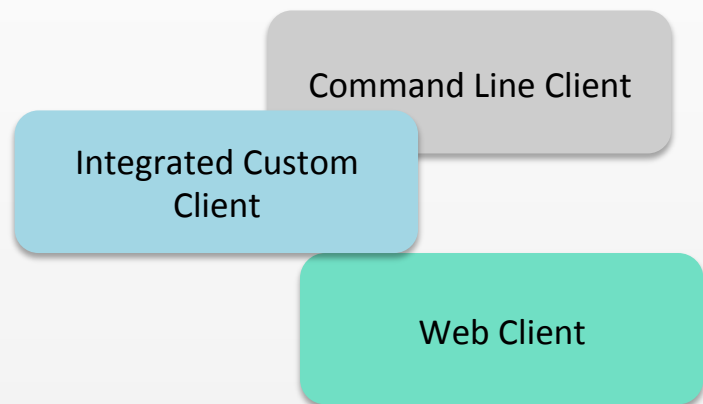Objects, SQL Queries, URLs all accessed like regular files.

Unified access. Unified control.

# iRODS is Extensible

iRODS has a pluggable architecture.

Existing plug-ins support a variety of hardware, communication technologies, database technologies, and storage topologies. Templates are available for new, custom plug-ins.

Command line, web clients, and numerous other clients are available for iRODS. Generic APIs allow developers to build efficient access to iRODS in to their software.

Command Line Client

Integrated Custom Client

Web Client

Familiar APIs

Microservices

Authentication Mechanisms

Storage Resources

Network Services

Coordinating Resources

Databases

# WHO USES iRODS?

# Who Uses iRODS?

- **Federal Users**
  - National Institutes of Health (NIH)
  - National Aeronautics and Space Administration (NASA)
  - National Oceanic and Atmospheric Administration (NOAA)
  - National Optical Astronomy Observatory (NOAO)
  - US Geological Survey (USGS)

- **Storage Vendors and System Integrators**
  - DataDirect Networks
  - EMC
  - Xyratex
  - Distributed Bio
  - Computer Sciences Corporation (CSC)

- **Commercial Users**
  - DOW Chemical
  - Beijing Genomics Institute

- **Research Programs**
  - The iPlant Collaborative
  - Broad Institute
  - International Neuroinformatics Coordinating Facilities (INCF)
  - Wellcome Trust Sanger Institute
  - Computer Center of the French National Institute of Nuclear and Particle Physics (CC-IN2P3)
  - CineGRID

- **Hundreds of academic institutions** worldwide host thousands of users on their iRODS data grids

# iRODS – Proven at Scale

- iPlant: 15,000 users on an iRODS data grid with 100 million files

- IN2P3: over 6 PB of data managed by iRODS

- Sanger Institute: 20+ PB of iRODS data

- NASA Center for Climate Simulations: 300 million metadata attributes

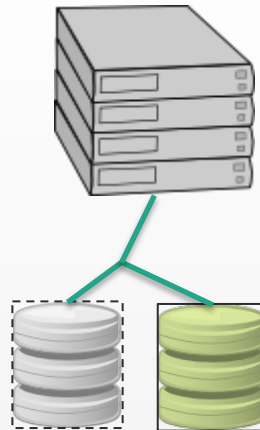- CineGRID: sites distributed across Japan-US-Europe

# WHAT CAN iRODS DO?

# What Can iRODS Do?

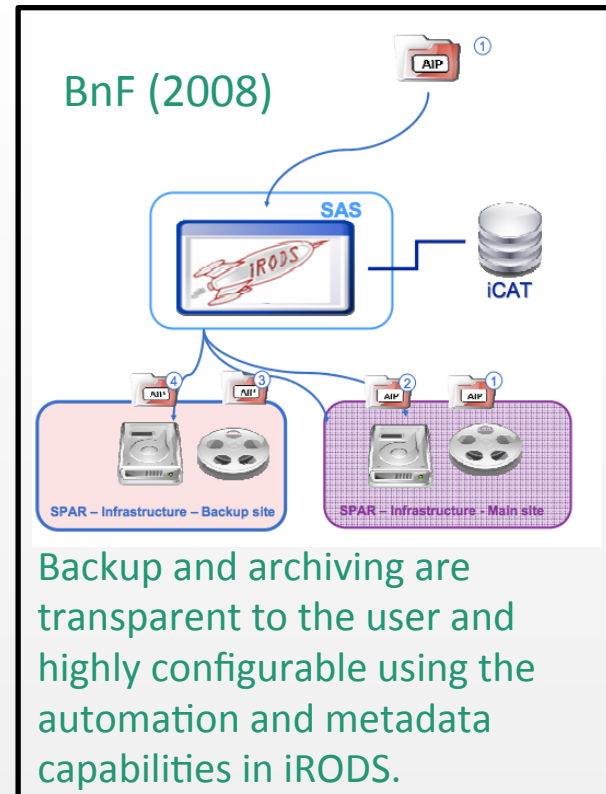For Data Center Managers, iRODS simplifies data grid management.



Data on different storage devices at different locations can be centrally managed.



**In situ** migration to new hardware can be managed by replicating the legacy resource before repurposing or decommissioning it.



BnF (2008)

Backup and archiving are transparent to the user and highly configurable using the automation and metadata capabilities in iRODS.

# What Can iRODS Do?

For Users, iRODS simplifies data discovery, data validation, and data processing.



User-defined and intrinsic metadata make stored data searchable.

Validation and analytical tools can be automated to process incoming data.

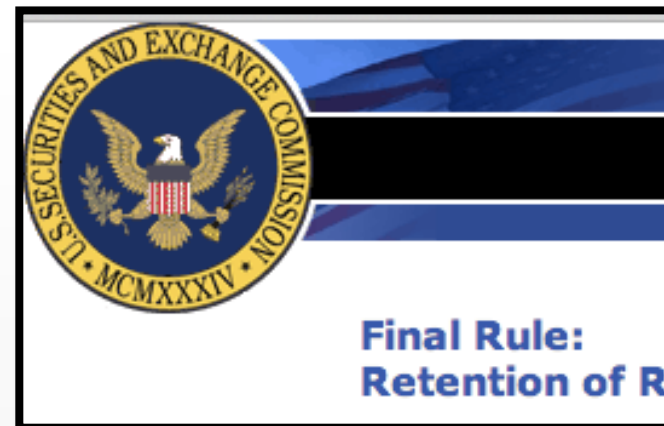The results and process steps can be stored in the iCAT metadata catalog.

# Policy Enforcement and Compliance Verifications

With its metadata catalog and automation capabilities, iRODS presents the infrastructure to enforce mandated data management policies, such as those for records retention and privacy protection.
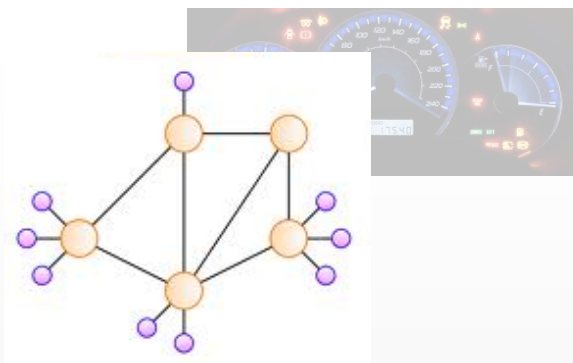
Audit trails generated by iRODS can be used to verify compliance with policy.

# THE FUTURE OF iRODS

# The Future of iRODS: Improving Uptake

- **Plug-In Bundling:** Easier Deployment to Specific Market Segments

- **Registry:** Enables Bundling, Easier Upgrades, Dashboard

- **Simplifying Connection APIs:** Improved Consistency across Programming Languages → More Clients and Plug-Ins

# The Future of iRODS: Full Content Indexing

Efforts underway to:

- Perform full content indexing of incoming data
- Connect iRODS to an external database of indexed search terms

Data discovery based on:

- *Data Content*
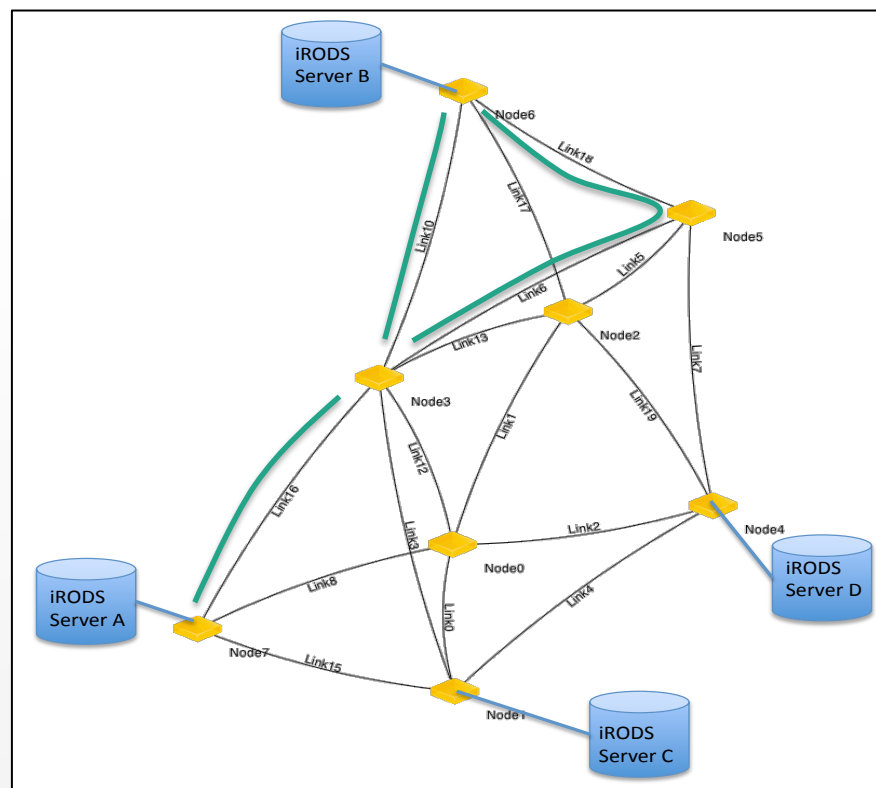- Automatically-generated metadata
- User-generated metadata

# The Future of iRODS: Application-Defined Networking

- iRODS can control a software-defined network (SDN).

- iRODS has been used experimentally with a SDN to maximize bandwidth between storage resources.

- Parallel disjoint network links are dynamically created and torn down in response to network conditions and other parameters.

# BACKUP SLIDES

# Use Case:
## University of North Carolina, Chapel Hill (UNC)

Slides courtesy of Charles Schmitt

# Genomics Primary Physical Infrastructure



Secure Medical Workspace

Sequencing center

Isilon

Storage Media Library

*Primary Processing, Storage, and Archive*

UNC Kure HPC

iRODS Grid

Dell

RENCI Croatan 'Big Data'

DDN, Quantum

RENCI BlueRidge HPC

Genomic Sciences TopSail Hadoop

*Hadoop computing, SAN scratch storage*

*Additional Processing, Secondary Archive*

Storage Media Library

# Example: Unified View of Data



iDROP web client

Spread across:
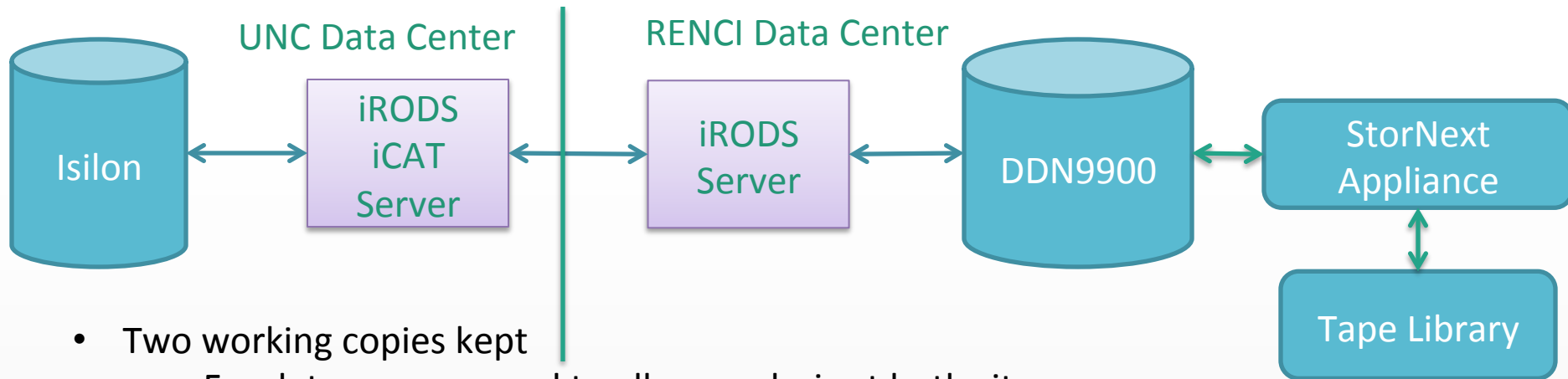1) Disk-storage at UNC, 2) Disk-storage at RENCI, 3) Tape-storage at RENCI
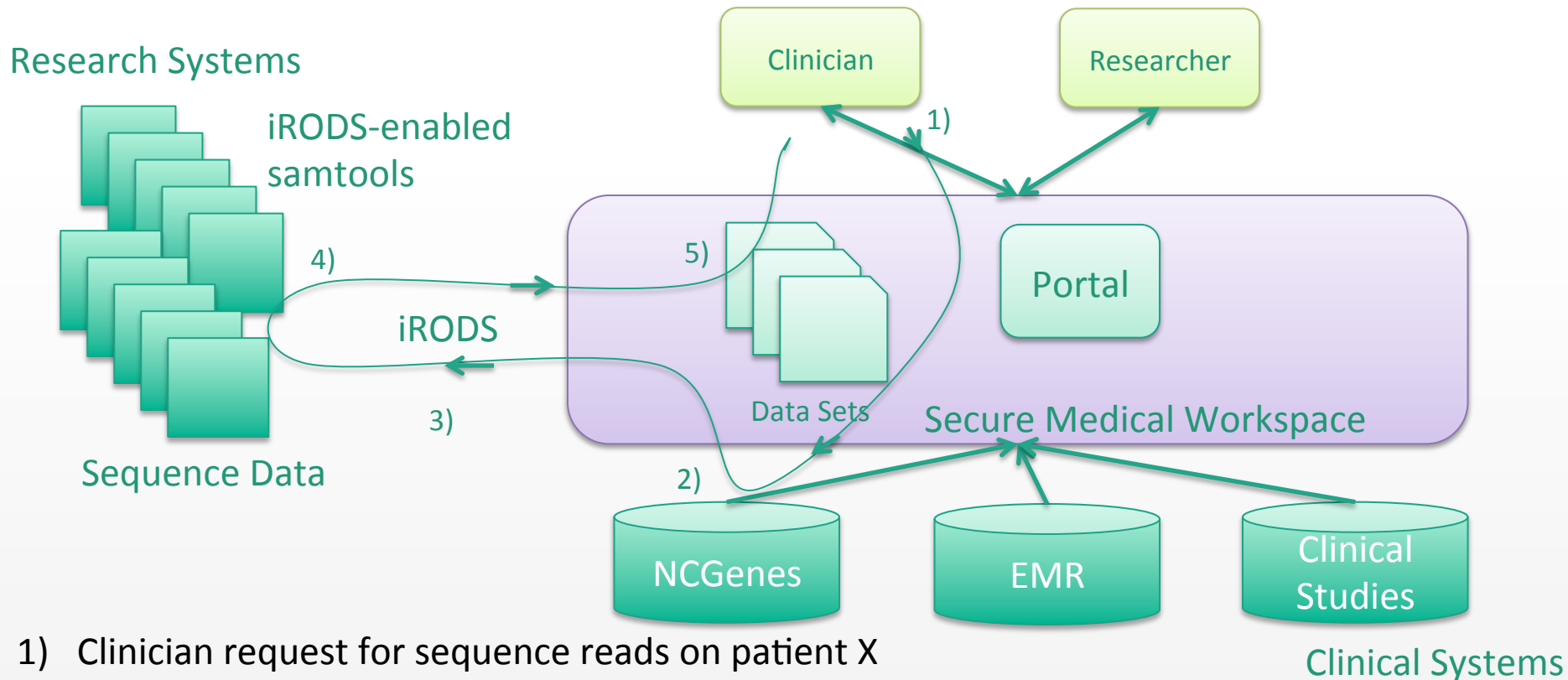
# Example: Data Access Policy

- Challenge
  - Millions of files across different projects, growing daily
  - Hundreds of users across different labs, changing frequently
  - *How to control access*
    - *UNIX ACLs became too unwieldy*
    - *Moving data means reproducing permission and group settings*

- Policy:  access given if user and data belong to the same groups
  - Tag data with group metadata (e.g., Lab X lung tumor study)
  - Access rule: user's group must match data group
    - E.g. (user y member of Lab X lung tumor study)

# Example: Data 'Replication' Policy

**UNC Data Center**

**RENCI Data Center**

Isilon ↔ → iRODS iCAT Server ↔ → iRODS Server ↔ → DDN9900 ↔ → StorNext Appliance ↕ Tape Library

- Two working copies kept
  - For data recovery and to allow analysis at both sites

- 'Copy me' and 'Data copied' metadata control copy process
  - Only on certain files (fastq, 'finished' bam files)

- iRODS rule performs the copy nightly
  - Performs copy, verifies copy successful, resets 'copy me' attribute

- Versioning to allow for re-runs of patient samples

iRODS CONSORTIUM

# Secure Access to Data on the Clinical Side

**Research Systems**

iRODS-enabled samtools

Sequence Data

iRODS

4)

3)

Clinician

Researcher

1)

5)

Portal

Data Sets

Secure Medical Workspace

2)

NCGenes

EMR

Clinical Studies

Clinical Systems

1) Clinician request for sequence reads on patient X
2) Patient id lookup to obtain subject id
3) Subject id lookup in iRODS
4) Data sets packaged in zip file and retrieved
5) Data unzipped and displayed within secure workspace

# Use Case:
## The Sanger Institute (Wellcome Trust Sanger Institute)
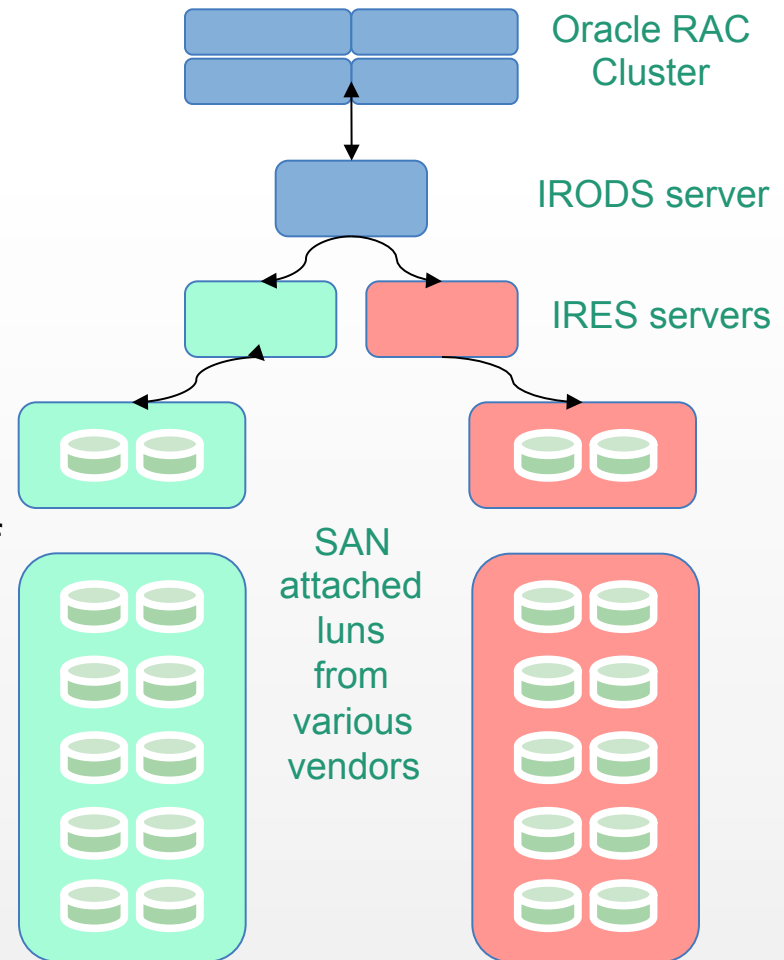
Slides courtesy of Peter Clapham

# iRODS layout

Data lands by preference onto iRES servers in the green data center room

Data is then replicated to red data center room via a resource group rule with checksums added along the way

Both iRES servers are used for r/o access and replication does work either way if bad stuff happens.

Various data and metadata integrity
   Checks are made.

Simple, scalable and reliable

Oracle RAC Cluster

IRODS server

IRES servers

SAN attached luns from various vendors

# Metadata-Rich

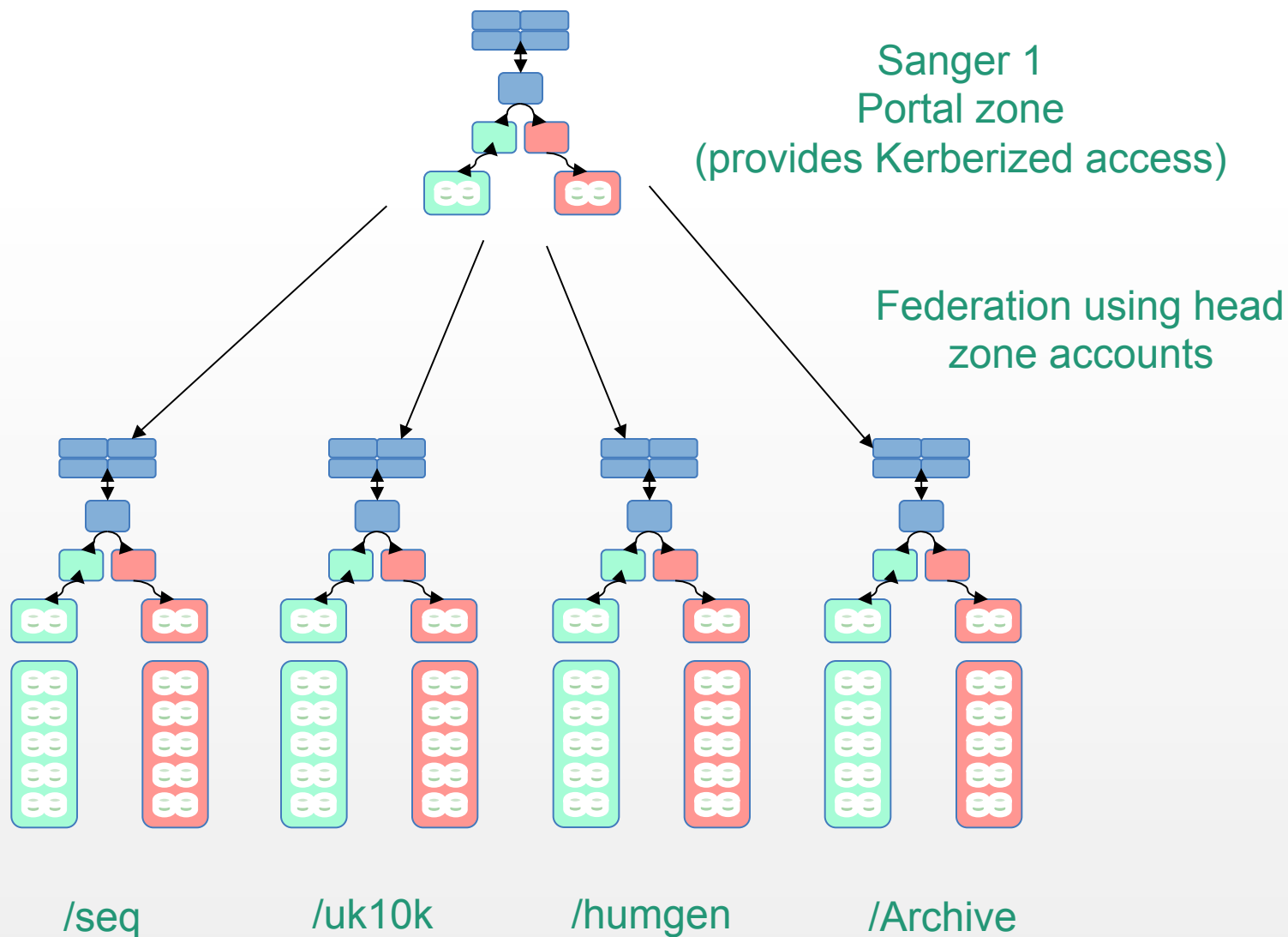Example attribute fields →

Users query and access data largely from
    local compute clusters

Users access iRODS locally via the
    command line interface

attribute: library
attribute: total_reads
attribute: type
attribute: lane
attribute: is_paired_read
attribute: study_accession_number
attribute: library_id
attribute: sample_accession_number
attribute: sample_public_name
attribute: manual_qc
attribute: tag
attribute: sample_common_name
attribute: md5
attribute: tag_index
attribute: study_title
attribute: study_id
attribute: reference
attribute: sample
attribute: target
attribute: sample_id
attribute: id_run
attribute: study
attribute: alignment

# Sanger Zone Arrangement



Sanger 1
Portal zone
(provides Kerberized access)

Federation using head
zone accounts

/seq          /uk10k          /humgen          /Archive

# Baton iRODS "Client"

Thin layer over parts of the iRODS C API

      JSON support

      Connection friendly

      Comprehensive logging

      autoconf build on Linux and OSX

Current state

      Metadata listing

      Metadata queries

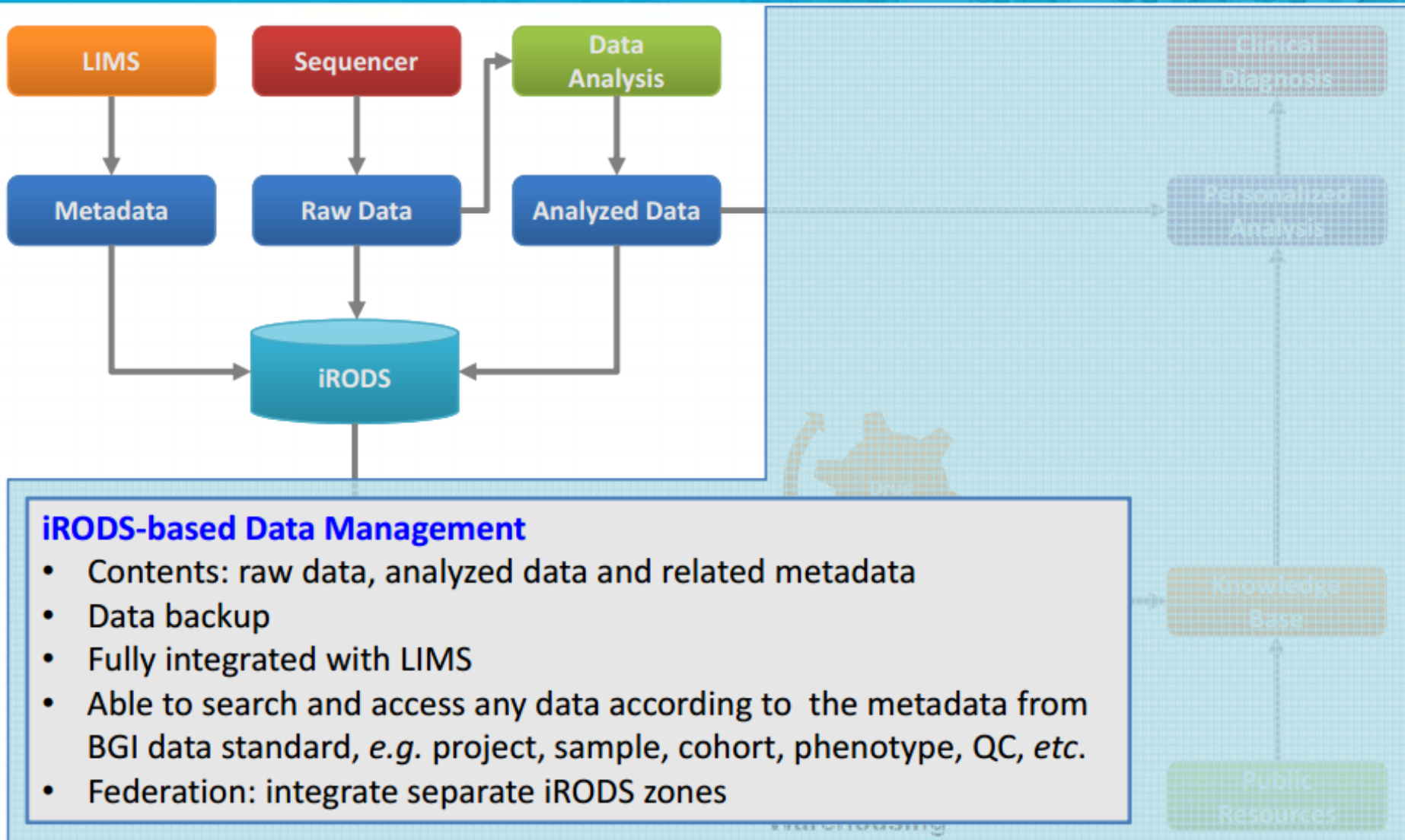      Metadata addition

https://github.com/wtsi-npg/baton.git

# Use Case:
## Beijing Genomics Institute (BGI)

Slides courtesy of Xing Xin

# ● **The world largest genome sequencing center**

- **Started with Human Genome Project in 1999 with only a few sequencers.**

- **Now more than 150 sequencers, 6 TB/day sequencing throughput.**

| MODEL | ABI 3730XL | Roche 454 | ABI SOLiD 4 | Solexa GA IIx | Illumina HiSeq 2000 |
|---|---|---|---|---|---|
| INSTALLATION | 16 | 1 | 27 | 6 | 135 |

**iRODS-based Data Management**

- Contents: raw data, analyzed data and related metadata
- Data backup
- Fully integrated with LIMS
- Able to search and access any data according to the metadata from BGI data standard, *e.g.* project, sample, cohort, phenotype, QC, *etc.*
- Federation: integrate separate iRODS zones

# Acknowledgements

- Presented work funded in part by grants from NIH, NSF, NARA, DHS. Funding also provided by UNC and the iRODS Consortium.

- Teams involved include:
  - *DICE team at UNC and UCSD*
  - *Networking team at RENCI and Duke*
  - *Data sciences team at RENCI*

- *In collaboration with*
  - *UNC Dept of Genetics, Research Computing, Lineberger Comprehensive Cancer Center, NC TraCS Institute, Center for Bioinformatics, Institute for Pharmacogenetics and Personalized Treatment*
  - *UNC HealthCare*

- *Multiple members of the iRODS community*

# Data Policy and iRODS - Definitions

- ## What is (Digital) Data Management?

  - procedures and operations to assure value of digital assets
    - protection: verification, back-ups & replicas, security, access controls, …
    - maintenance: migrating to tape, integrity checking, …
    - control: format conversion, derived data products, …
    - accessibility and usability: discovery, availability, supporting services, …

  - defined by data proprietors or data administrators

  - enables analytics and operations that pull value from the data

- ## What is data policy?
  - statement of data management strategy
  - the ensemble of data management requirements and procedures

- Data policy can be *implemented* and *executed* manually or in automated fashion

# iRODS History

- SRB: initial product begun by DICE, 1997 at the San Diego Supercomputer Center, UCSD and General Atomics

- iRODS: rewrite of SRB by DICE in 2006; current version: iRODS 3.3.1

- Very close interaction with worldwide user communities who drive development

- Enterprise iRODS (e-iRODS): mission critical distribution co-developed by RENCI and DICE in 2012

- iRODS 4.0: merge of the iRODS and e-iRODS codes by iRODS Consortium to form a common core and full deployment of plug-in architecture