# Overview: Data Management & iRODS

#### **Center for Data Intensive Cyber Environments**

University of North Carolina at Chapel Hill

UNC School of Information and Library Science (SILS) Renaissance Computing Institute (RENCI) Institute for Neural Computation (INC) at UC San Diego

diceresearch.org irods.org







₹UCSI







of NORTH CAROLINA

₹UCSI

i-R-O-D-S



### Some questions to ask...

#### How will you share Data Collections?

- Throughout your collaboration and beyond
- Regardless of where the data is

D·I·C·E

- Including diverse data from collaborators, while
- Controlling access in nuanced ways
- Adding/removing data, users, partners, different infrastructure, etc.
- Using the Data Collections with your applications to analyze, visualize, create derived works, etc.

of NORTH CAROLINA

#### Some challenges in Digital Data

Ephemeral; technology obsolescence; exploding size (10 times larger 2006–2011); Babel - proliferating proprietary formats, etc.

**₹UCSE** 



### Toward a Unified Data Space

#### Sharing data across space

- Multiple spaces: geographic, institutional, disciplinary...
- Infrastructure spaces (h.w./s.w.)
- Harness power of cyber technologies
  - Virtual Collections of distributed data
    - Global Name Spaces (data, users, storage, etc.)

₹UCSI

- Beyond single-site repository model (hard copy based)
- Also known as a "Data Grid"



### Toward a Unified Data Space

#### Sharing data across time

- A "memory" for your project
  - Communicating with the future
  - Long-term preservation

Need automated policies that govern a life cycle workflow from ingestion to disposition, access, validate authenticity, access, etc.

# Sharing in space and time require related capabilities, architecture

Trend: people realizing they need to both share and preserve









### Another question to ask...

Does your system design let data be "born" in a comprehensive environment?

From fragmented, ad hoc to intentional

Your decide Policies for managing, when to discard

Cheaper: "A stitch in time saves nine"

- Automation for mushrooming data collections
- Avoids generating more "legacy data" that must be harvested later (difficult, expensive)
- Trend toward annotation at creation
  - Upstream collaboration between data creators and data professionals















### iRODS User Community

- □ iRODS Development Collaborations
  - NARA TPAP Transcontinental Persistent Archive Prototype (NARA funded)
  - NSF SDCI Research in Adaptive Middleware Architecture Systems
  - SHAMAN Sustaining Heritage Access through Multivalent ArchiviNg
  - UK e-Science data grid
- Communities Using DICE Technologies, including Biology, Environment, Psychology, Human Subjects
  - BIRN Biomedical Informatics Research Network (NIH funded)
  - ROADNet Real-time Observatories, Applications, and Data management Network (NSF funded)
  - SEEK Science Environment for Ecological Knowledge (NSF funded)
  - TDLC Temporal Dynamics of Learning Center (NSF funded) Overview
- Physical Sciences Uses
  - CADAC Computational Astrophysics Data Analysis Center (NSF funded)
  - BaBar high energy physics data grid (DOE funded)











## iRODS User Community

- Physical Sciences (continued)
  - NOAO National Optical Astronomy Observatories data grid (NSF funded)
  - NVO National Virtual Observatory (NSF funded)
  - Observatoire de Strasbourg, France, VOSpace Interface
- Persistent Archives and Digital Preservation / Humanities Uses
  - NARA TPAP Transcontinental Persistent Archive Prototype
  - e-Legacy Preserving the Geospatial Data of the State of California
  - DCPC Distributed Custodial Preservation Center (NHPRC funded)
  - DIGARCH UCTV NSF Digital Archiving and Long-Term Preservation (LoC)
  - T-RACES Testbed for Redlining Archives of CA Exclusionary Spaces (IMLS)
- Geosciences Uses
  - OOI Ocean Observatories Initiative (NSF funded)
  - SCEC Southern California Earthquake Center (NSF funded)
- High Performance and Grid Computing
  - NSF TeraGrid
- Plus many international users.
- And growing all the time...





**₹UCSI** 

RCHIVE

#### **Introduction to iRODS Data System**

#### You, Researchers, Students, etc.

Want to easily Find, Access, Use, Move, Share Data, and more... With your Interfaces, your Applications, your Workflows

#### *iRODS Data System – "Middleware"*

A "layer" that "connects the dots" while masking and automating your interactions with diverse infrastructure.

#### The "World of Infrastructure"

Your and other's Storage, Networks, Admin. Domains, Computing Services, Web Services, etc.













#### **iRODS Shows Unified "Virtual Collection"**

User With Client Views & Manages Data

**User Sees Single "Virtual Collection"** 

**My Data** Disk, Tape, Database, Filesystem, etc. **My Data** Disk, Tape, Database, Filesystem, etc.

Partner's Data Remote Disk, Tape, Filesystem, etc.

The iRODS Data System can install in a "layer" over existing or new data, letting you view, manage, and share part or all of diverse data in a unified Collection.







#### Adding Data to iRODS Data System





#### **Preserving Electronic Records with iRODS**



Archivists can use iRODS for preserving Electronic Records, from Appraisal to Access, with Rules enforcing trustworthy respository criteria with audits.

**₹**UCSD

D·I·C·E i·R·O·D·S

THE UNIVERSITY of NORTH CAROLINA at CHAPEL HILL







## More Information about iRODS

- Shared collections assembled from data distributed across different groups, remote storage locations
- 2. Workflow environment executed where data is (server-side on remote storage)
- 3. Policy enforcement engine, with computer actionable Rules applied at remote storage
- **4. Validation environment** for assessment criteria (audit trails)

of NORTH CAROLINA

i-R-O-D-S

**D**·I·C·E

 Consensus building system for establishing collaboration (policies, data formats, semantics, shared collection, etc.)

**₹UCSE** 

### Use Case: Dissertation Collection

- Support for Student Dissertation data in Health Sciences Library (UNC)
  - Organize collection of student simulation data Input files, output files
  - Use iRODS Rules to periodically synchronize student's work area with Collection, registering new files into Collection, and replicating files to tape archive
    - Build templates to describe required metadata for registered files

of NORTH CAROLINA

Use iRODS Rule to verify compliance of metadata for each file with template

**₹UCSD** 



### Use Case: Digital Humanities

- T-RACES: Testbed for the Redlining Archives of California's Exclusionary Spaces
  - A digital humanities collaborative between UNC and UCHRI
  - Building iRODS Data Grid for the digital humanities

**₹UCSE** 

- Provides integrated map, text, and database interfaces
- Extend to redlined cities of North Carolina

of NORTH CAROLINA



### Use Case: NARA Archiving

- NARA Transcontinental Persistent Archive Prototype
  - Federates 7 independent iRODS data grids: Each manages own Storage resources and Metadata Catalog, applies own Policies
  - Use iRODS federation to establish Policies for sharing data between sites.
  - Control operations a remote user can do within your data grid.

of NORTH CAROLINA

Extensible Environment, can federate with additional research and education sites. Each data grid uses different vendor products.

**₹UCSD** 





### iRODS – more details

#### A data grid system - data virtualization

- A distributed file system, based on a client-server architecture.
- Allows users to access files seamlessly across a distributed environment, based upon their attributes rather than just their names or physical locations.
- It replicates, syncs and archives data, connecting heterogeneous resources in a logical and abstracted manner.

#### A distributed workflow system - policy virtualization

- Policies can be coded as functions (micro-services)
- Remote micro-services can be chained
- The chains (workflows) are interpreted at run-time
- Chains can be triggered on an event and condition (Rules)
  - They can also be recursive.
  - Micro-services communicate through parameters, shared contexts, and out-of-band message queues.











## Building a Shared Collection



## Shared Collection Challenges

- Need common naming conventions to identify
  - Collaborators

**D**·I·C·E

- Shared data and their types & methods
- Shared data resources & access policies
- Need discovery metadata
  - Assign attributes to each name space
    - State information (metadata)
- Assign policies between name spaces
  - Access constraints, disposition policy, integrity

Mediate across site and project policies

ORTH CAROLINA

**₹UCS** 



## Discovery: Metadata

#### System Metadata

- User name space
  - □ Address / e-mail / telephone number
  - □ Role (administrator, curator, user)
- File name space
  - □ Creation date / size / location / checksum
  - Owner / access controls
- Storage resource name space
  - Capacity / quotas / Type (archive, disk, fast cache)
- Domain Metadata
  - User-given metadata
    - Key-Value-Unit Triplets, Annotations
    - Relational / XML Metadata
    - Domain-specific Schemas
      - Dublin Core, Darwin Core, FITS, DICOM, ...







### Under the hood - a glimpse



## Policies in iRODS

**Policies:** Express community goals for data access and sharing, management, long-term preservation, uses, etc.

#### Policy Examples

- Run a particular workflow when a "set of files" is ingested into a collection (e.g. make thumbnails of images, post to website).
- Automatically replicate a file added to a collection into 3 geographically distributed sites.
- Automatically extract metadata for a file of a certain type and store in metadata catalog.
- Periodically check integrity of files in a Collection and repair/replace if needed/possible.
- Automatically pick a certain storage location based on user or collection or size or type.
- Let a user access a collection only if using certificate-based login.
- Send a notification when a certain file is ingested.
- etc.



THE UNIVERSITY of NORTH CAROLINA at CHAPEL HILL

₹UCSI





# iRODS Rules

- Implement Policies
- Verify enforcement (audit trails)
- Automate management of exploding data
  - Let you handle petabytes in hundreds of millions of files
- Each Rule defines

i-R-O-D-S

- Event, Condition, Action chains (micro-services, other Rules), Recovery chains
- Rule types

D-I-C-E

- Atomic (immediate), Deferred, Periodic
- Rules are executed by Micro-services
  - Applied where data is (server-side)







#### Micro-services

- Function snippets perform a small, well-defined operation/semantics, e.g.
  - computeChecksum
  - replicateFile
  - integrityCheckGivenCollection
  - zoomImage
  - getSDSSImageCutOut
  - searchPubMed
- Chained to implement iRODS Rules (workflows)
- Invoked by the iRODS Rule Engine
- Recovery micro-services provide roll-back upon failure
- Currently C functions; PHP, Java coming soon
- Can wrap Web-services













# DICE Center

#### Center for Data Intensive Cyber Environments

- University of North Carolina at Chapel Hill (UNC)
  - UNC School of Information and Library Science (SILS)
  - Renaissance Computing Institute (RENCI)
    - Reagan Moore
    - Richard Marciano
    - Arcot Rajasekar
    - Antoine de Torcy, Chien-Yi Hou
- UC San Diego
  - □ Institute for Neural Computation (INC)
    - Mike Wan
    - Wayne Schroeder
    - Sheau-Yen Chen, Lucas Gilbert, Bing Zhu, Paul Tooby
- iRODS development is supported by
  - NSF OCI-0848296 "NARA Transcontinental Persistent Archives Prototype" (2008-2012)

**₹UCSE** 

 NSF SDCI 0721400 "Data Grids for Community Driven Applications" (2007-2010)

THE UNIVERSITY

at CHAPEL HILI

of NORTH CAROLINA



D·I·C·F