

Enabling Inter-repository Access Management between iRODS and Fedora

Bing Zhu,
Uni. of California: San Diego
Richard Marciano
Reagan Moore
University of North Carolina at Chapel Hill

May 18, 2009

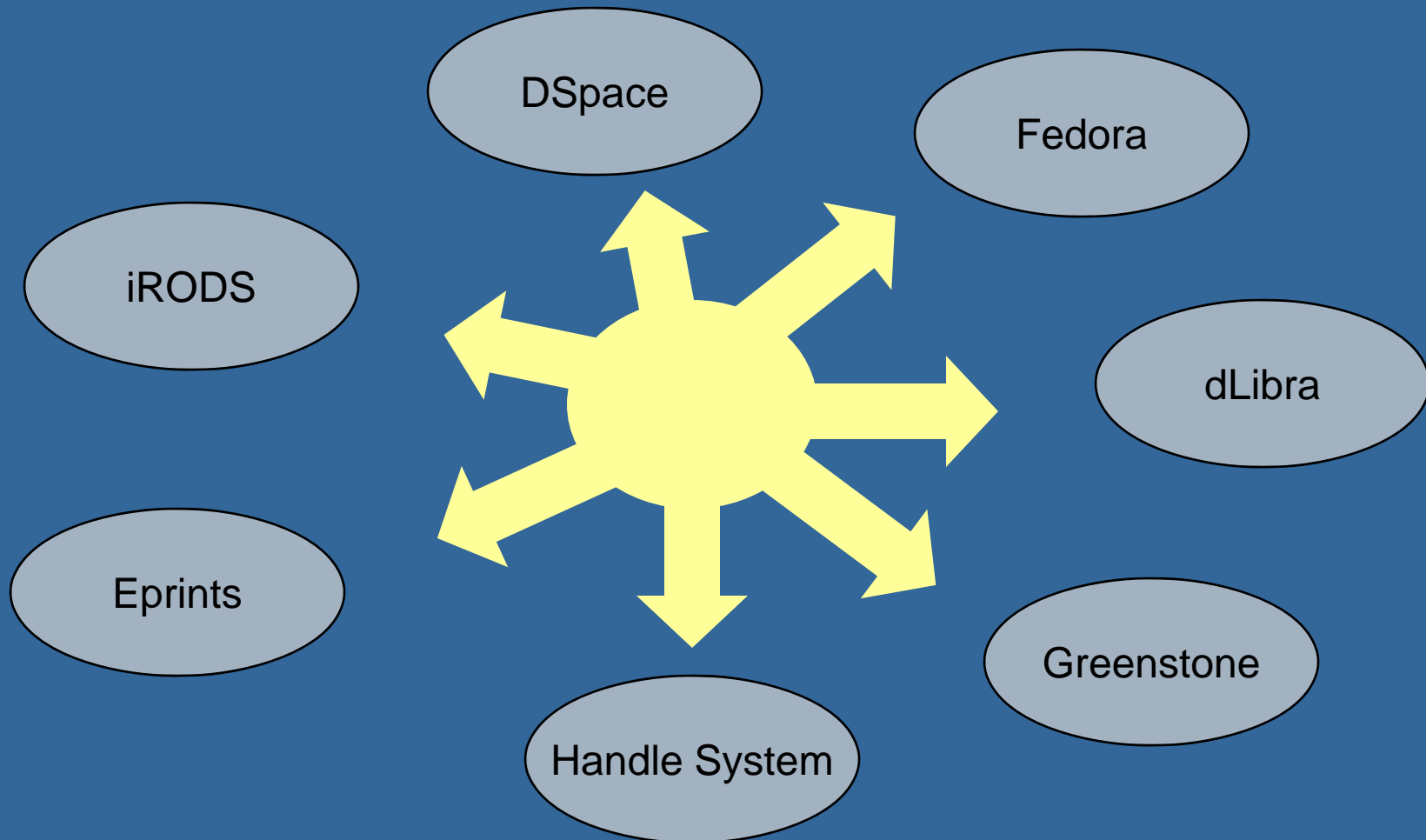
Atlanta



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



An Environment with Heterogeneous Technologies



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Sharing Data Across Repositories

Enabling inter-repository data management allows us to share data by connecting the repositories of:

- Different groups, projects
- Different institutions, locations
- Different disciplines
- Diverse types of data
- Diverse hardware, software infrastructure



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Issues for inter-Repository data management

- Object Model
- Virtual Registration of Digital Objects from One Repository to Another
- Inter-repository Service Management
- Policy Enforcement across Repositories



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



iRODS

iRODS - integrated Rule-Oriented Data System

A middleware providing functions to:

- manage distributed storages
- provide metadata support for digital preservation and search functions
- allow running distributed workflows to enforce system policies and harvest distributed computing power.

iRODS can be used for

- building datagrid
- building digital library
- building digital repositories



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



iRODS (continued)

- **Current version 2.02**
- **Scalable: Managing > 4 PB worldwide**
 - Collections > PB, 100s of millions of files
- **Federation of data grids**
 - Flexible collaborative data sharing. Scaling >1 DB catalog.
- **Transferring data**
 - Parallel transfers ~70% of available bandwidth.
- **Independent evaluations: NASA, etc.**



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



iRODS (Continued)

iRODS Provides:

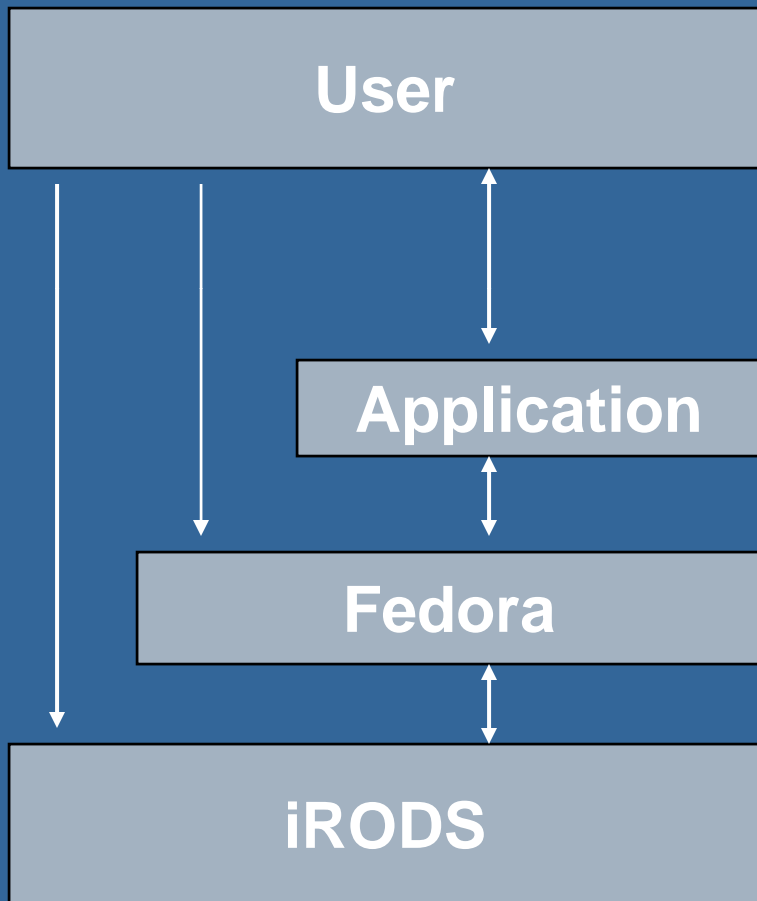
- A uniform global naming space for managing digital datasets
- An efficient data access mechanism using parallel data transfer
- Access to data stored on distributed systems
- Rich interfaces including C/C++ API, Java API, Perl, Web Service, and Python
- Support for remote manipulation of data sets to minimize the amount of data sent over the network
- ...



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Why iRODS + Fedora?



Complex Object Modeling Layer

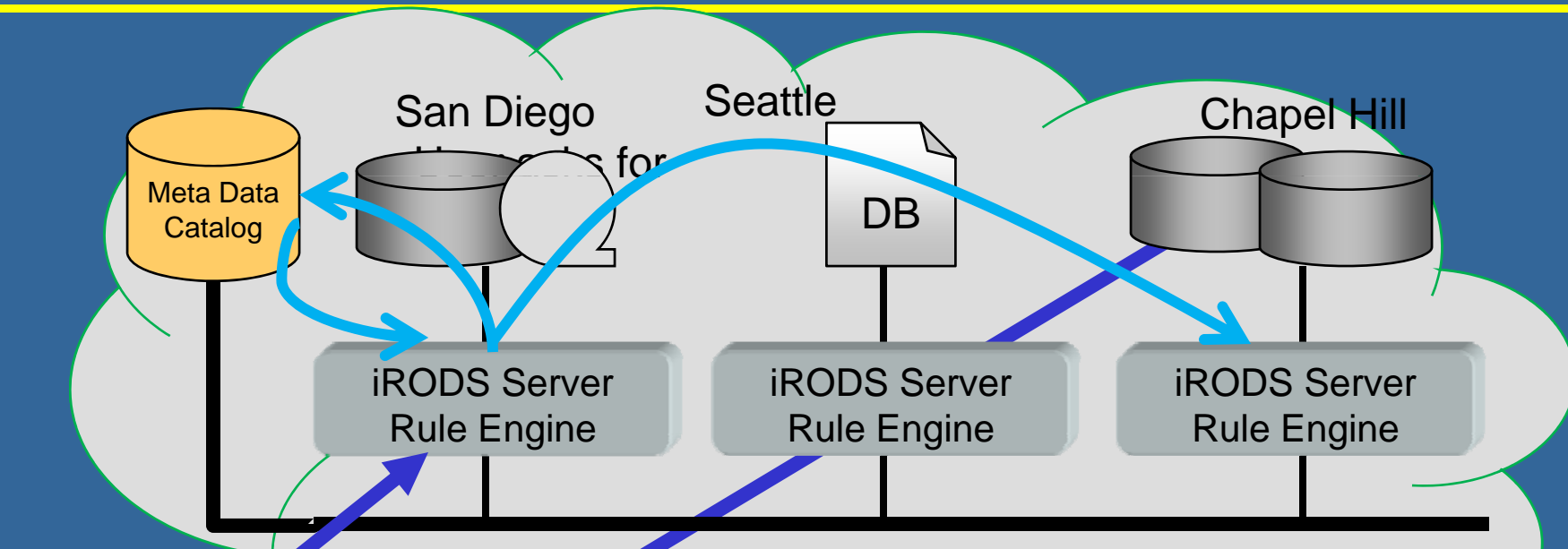
Digital Preservation Layer



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



iRODS for Digital Preservation



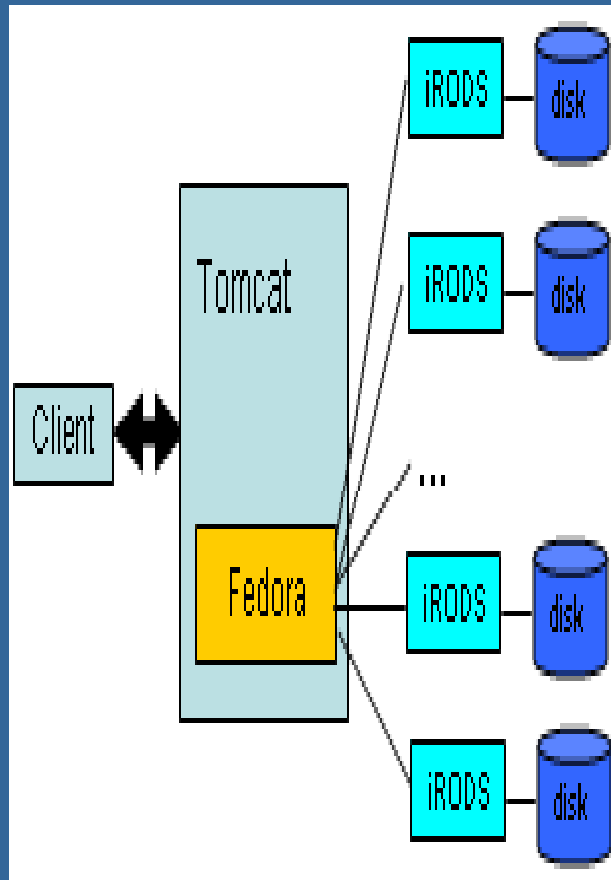
- Data replication service
- Periodic data integrity check
- Distributed storages for disaster recovery
- Metadata support for preservation description information



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



iRODS Storage Module for Fedora



- A replacement for Fedora's local storage module
- A standalone plug-in module (independent of Fedora release)
- iRODS manages both Fedora objects (XML) and data streams for Fedora (for managed content)
- Implemented based on SRB Storage Module for Fedora by DART project

iRODS Storage Module

- Manual Management of Distributed Stores - Admin Selects a Storage Resource for Storing Data Objects
- Auto Management of Distributed Storages - Use a Logical Storage Resource



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Deployment of iRODS Storage Module in Fedora

- Download Jar files from iRODS web site

<https://www.irods.org/index.php/Fedora>

- Copy the jar files into Fedora place under Tomcat

`$CATALINA_HOME/webapps/fedora/WEB-INF/lib`

- **Edit the Fedora config file** `$FEDORA_HOME/server/config/fedora.fcfg`

```
<module role="fedora.server.storage.lowlevel.ILowlevelStorage"
  class="fedorax.server.module.storage.lowlevel.irods.IrodsLowlevelStorageModule">
  <param
    name="file_system" :value="fedorax.server.module.storage.lowlevel.irods.IrodsIFileSystem"/>
  <param name="irods_host" value="irods.sdsc.edu"/>
  ...
```

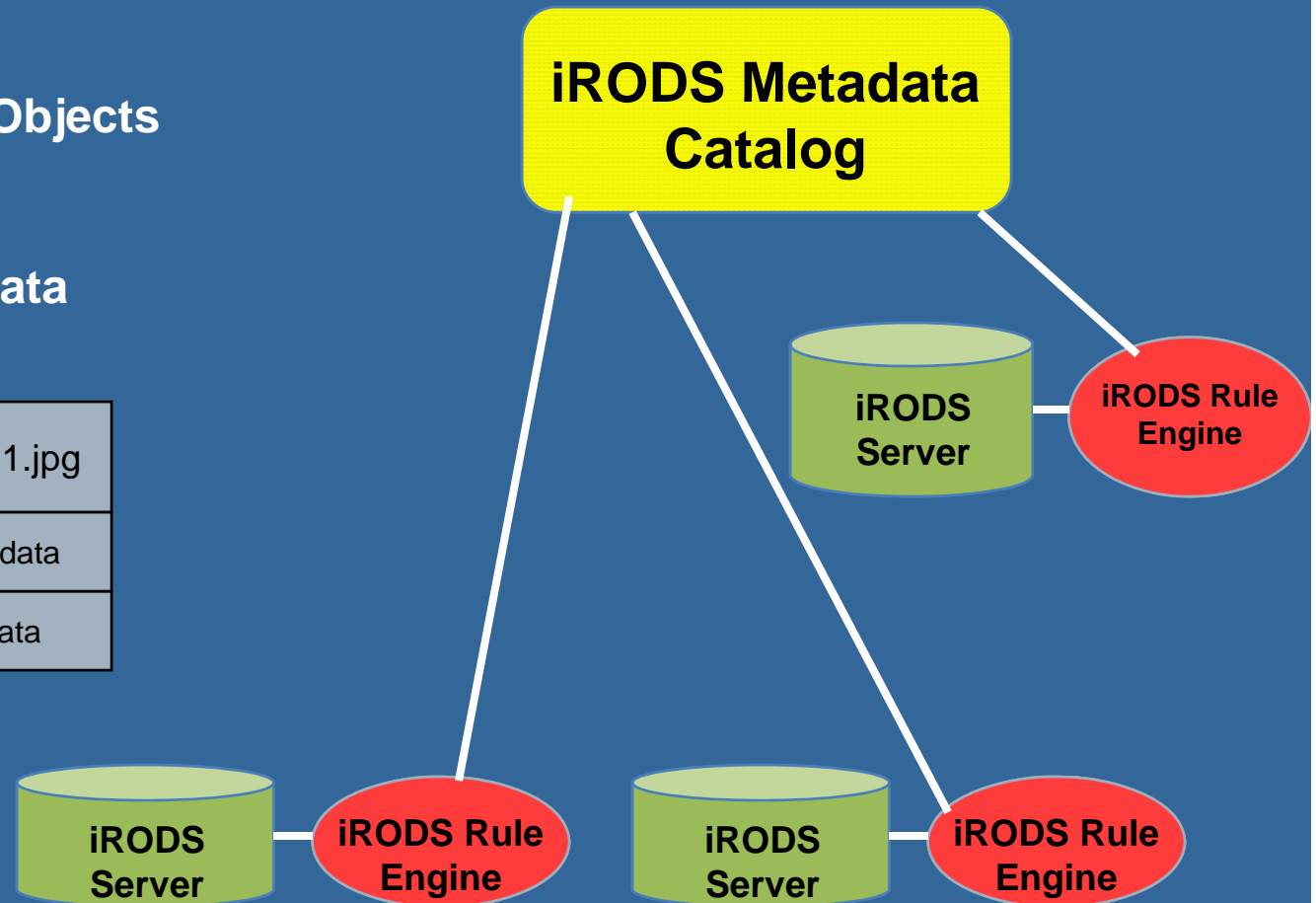
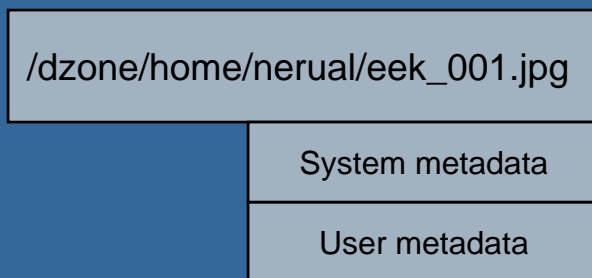


THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



iRODS Digital Object Model

- File Based
- Distributed Digital Objects
- system metadata
- user defined metadata

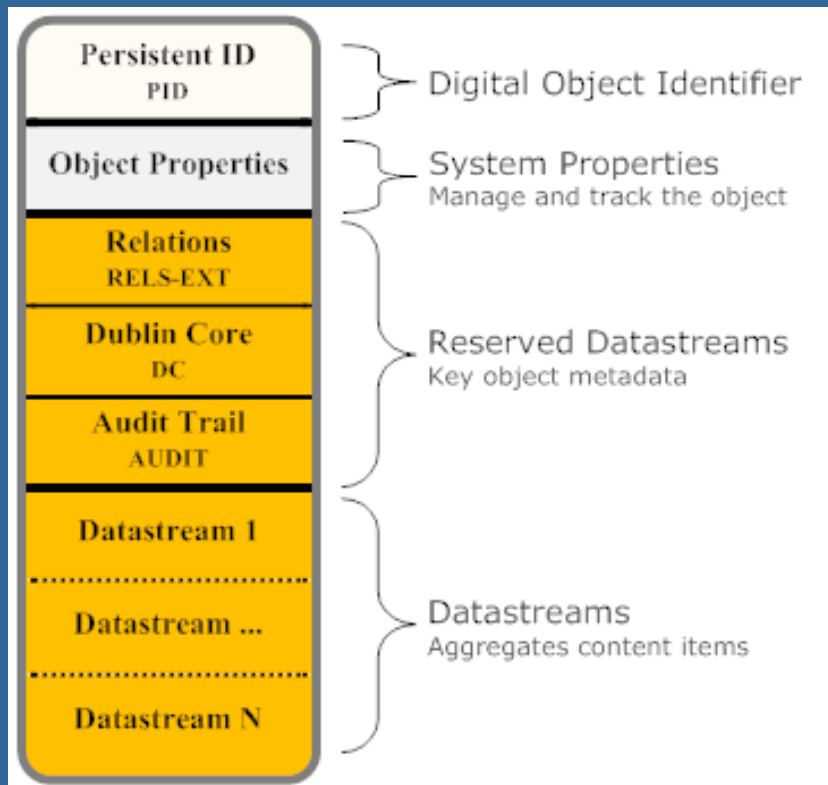


THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Fedora Digital Object Model

A Compound Object Model

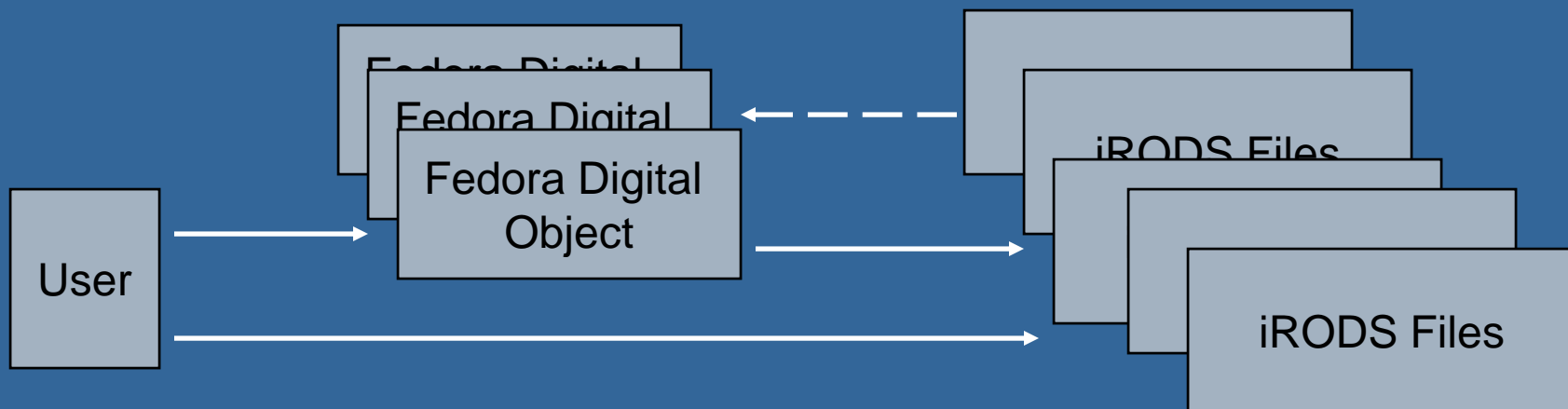


Four Types of control group for Fedora datastream:

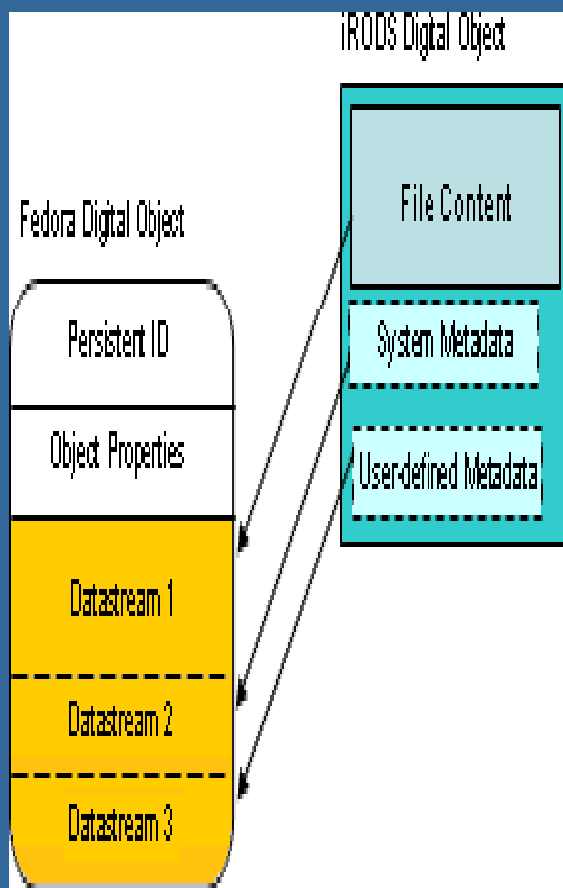
- **Internal XML Metadata**
- **Managed Content**
- **External Reference Content**
- **Redirect**

Referencing iRODS objects as Fedora External Datastreams

- An extension to the Fedora's external content manager
- iRODS files can be virtually linked inside Fedora
- Syntax:
 - `http://irods_hostname:port/irodsfullpath?fs=irods`
 - `future: irods://...`



Registering iRODS Objects into Fedora



- Create a Fedora object. The full path of the iRODS object becomes the label of the Fedora object.
- Create an external reference datastream for the iRODS file.
- iRODS system metadata is registered as an externally referenced datastream in Fedora.
- iRODS user-defined metadata is registered as an externally referenced datastream in Fedora.

Example: Create a Fedora Object

The screenshot shows a window titled "Object - irodsObj:2*" with two tabs: "Properties" (selected) and "Datastreams*". The Properties tab displays the following information:

State	Active ▼
Label	irods://srbrick15.sdsc.edu:7547/pzone/home/testuser/NeSC-Moore-intro.ppt
Created	2009-03-30T18:22:02.744Z
Modified	2009-03-30T18:45:31.238Z
Owner	fedoraAdmin

At the bottom of the dialog, there are five buttons: "View XML", "Export...", "Purge...", "Save Changes...", and "Undo Changes".



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Example: Create a datastream for the iRODS file

Object - irodsObj:2*

Properties | Datastreams*

DC

user-defined-metadata*

system-metadata*

NeSC-Moore-intro.ppt*

New RELS-EXT...

New...

ID	NeSC-Moore-intro.ppt
Control Group	External Reference
State	Active ▼
Versionable	Updates will create new version ▼
Created	2009-03-30T18:45:31.238Z
Label	The PPT file for e-Science workshop
MIME Type	application/ms-powerpoint
Format URI	
Alternate IDs	
Location	http://srbbrick15.sdsc.edu:7547/pzone/home/testuser/NeSC-Moore-intro.ppt?fs=irods
Fedora URL	http://srbbrick7.sdsc.edu:8089/fedora/get/irodsObj:2/NeSC-Moore-intro.ppt
Checksum	DISABLED ▼ none

Export... Purge...

Save Changes... Undo Changes



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Example: a dynamic reference to iRODS system metadata

The screenshot shows the 'Properties' tab for an object named 'irodsObj:2'. The left sidebar lists metadata categories: DC, user-defined-metadata*, system-metadata* (selected), NeSC-Moore-intro.ppt*, New RELS-EXT..., and New... The main area displays the following metadata configuration:

ID	system-metadata
Control Group	External Reference
State	Active
Versionable	Updates will create new version
Created	2009-03-30T18:28:42.481Z
Label	iRODS system metadata
MIME Type	text/xml
Format URI	
Alternate IDs	
Location	http://srbbrick15.sdsc.edu:7547/pzone/home/testuser/NeSC-Moore-intro.ppt?fs=irods&metadata=system
Fedora URL	http://srbbrick7.sdsc.edu:8089/fedora/get/irodsObj:2/system-metadata
Checksum	DISABLED none

At the bottom of the configuration area, there are buttons for 'View', 'Export...', 'Purge...', 'Save Changes...', and 'Undo Changes'.



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Example: a dynamic reference to iRODS user metadata

Object - irodsObj:2*

Properties | Datastreams*

DC

user-defined-metadata*

system-metadata

NeSC-Moore-intro.ppt*

New RELS-EXT...

New...

ID	user-defined-metadata
Control Group	External Reference
State	Active
Versionable	Updates will create new version
Created	2009-03-30T18:28:27.390Z
Label	iRODS user defined metadata
MIME Type	text/xml
Format URI	
Alternate IDs	
Location	http://srbbrick15.sdsc.edu:7547/pzone/home/testuser/NeSC-Moore-intro.ppt?fs=irods&metadata=user
Fedora URL	http://srbbrick7.sdsc.edu:8089/fedora/get/irodsObj:2/user-defined-metadata
Checksum	DISABLED none

View | Export... | Purge...

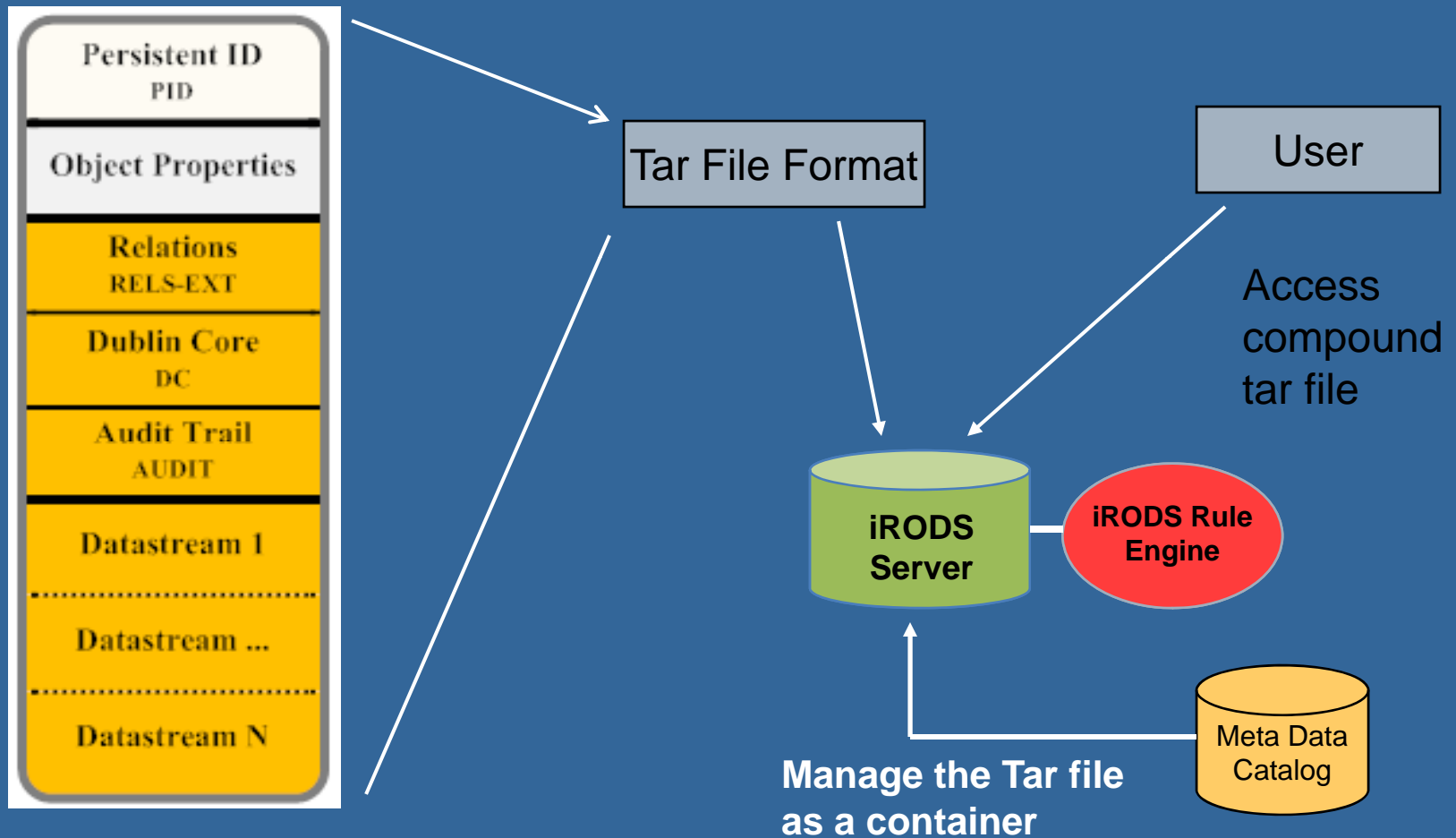
Save Changes... | Undo Changes



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Mapping Fedora Objects in iRODS



iRODS Rules

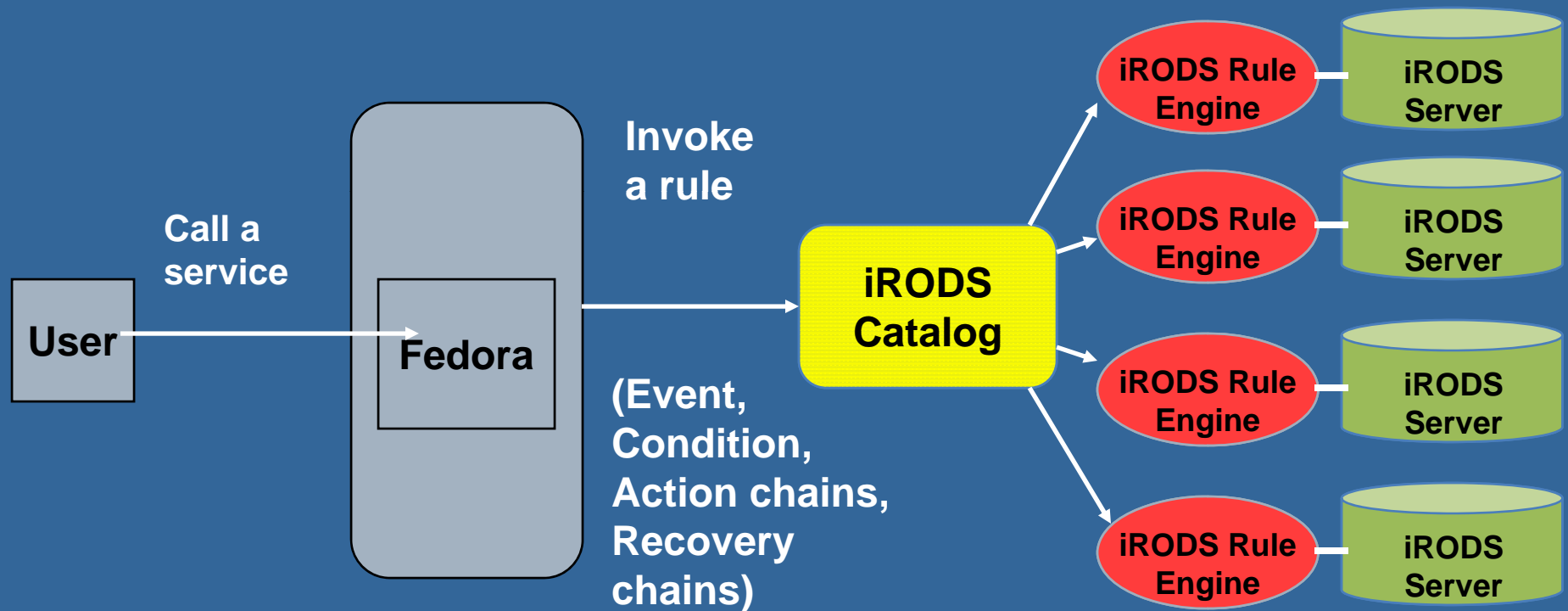
- Implement Policies
- Verify enforcement (audit trails)
- Automate management of exploding data
 - Let you handle petabytes in hundreds of millions of files
- Each Rule defines
 - Event, Condition, Action chains (micro-services, other Rules), Recovery chains
- Rule types
 - Atomic (immediate), Deferred, Periodic
- Rules are executed by iRODS Rule Engine
 - Applied where data is (server-side)



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Invoking iRODS service in Fedora



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



iRODS Policy

- To express community goals for data access and sharing, management, long-term preservation, uses, etc.
- Implemented through iRODS rules
 - Each rule is a chain of micro-services
 - Invoked by the iRODS Rule Engine
 - Currently C functions; PHP, Java coming soon
 - Can wrap Web-services



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Sample iRODS Policies

- 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.



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Policy-level Interoperability

- Can a preservation environment be assembled from two existing repositories with differing management policies?
- Can the policies of the federation be enforced across both repositories, ensuring consistent management of the archives?
- Can policies be migrated between repositories, either by association of the policies with the storage repositories, or through control of repository procedures?
- What fundamental mechanisms are needed within a repository to implement new policies?



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



Embedding Approaches

- ra-DM driving OM. Design policy federation models that are implemented at the storage level through the ra-DM.
- OM driving ra-DM. Design policy federation models in which the workflows within the OM model enforce the policies, but deposit the objects into the ra-DM.
- ra-DM and OM co-driving. Design policy federation models in which policies are enforced by both types of preservation environments.

OM : Objetc Model

ra-DM: rule-aware
distributed model



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL



DICE Center

iRODS:

- <http://www.irods.org>
- <http://www.dice.unc.edu>
- <http://www.diceresearch.org>

Fedora-iRODS Integration:

- <https://www.irods.org/index.php/Fedora>



THE UNIVERSITY
of NORTH CAROLINA
at CHAPEL HILL

