

Concepts in Distributed Data Management or History of the DICE Group

Reagan W. Moore¹, Arcot Rajasekar¹, Michael Wan³, Wayne Schroeder², Antoine de Torcy¹, Sheau-Yen Chen², Mike Conway¹, Hao Xu¹

¹University of North Carolina at Chapel Hill,

²University of California, San Diego

³San Diego, California

The Data Intensive Computing Environments group has been developing data grid technology for twenty years. Two generations of technology were created, the Storage Resource Broker - SRB (1994-2006) and the integrated Rule Oriented Data System - iRODS (2004-2014). When developed, both products represented pioneering technology in distributed data management and were widely applied by communities interested in sharing data, publishing data, and preserving data. Applications included national data grids, national digital libraries, national archives, and international collaborations. The success of the software was strongly driven by basic concepts that still represent the state-of-the-art for data management systems. These foundational concepts are built upon evolutionary concepts in virtualization and the abstractions needed to manage data, information, and knowledge. The concepts included policy-based data management, collection life cycle, and federation. The development, evolution, and application of these concepts in distributed data management systems is reviewed in this paper.

I. Introduction:

The implementation of two successful data grid software systems – the Storage Resource Broker (SRB) [2] and the integrated Rule Oriented Data System (iRODS) [13], represents an example of a software development life cycle. User requirements from the academic Science community drove the implementation of data grid technology, and the evolution data grids from data management to information and knowledge management. These systems pioneered significant conceptual ideas and technologies in large-scale data management and indeed added multiple terms to the ever changing vocabulary of the field. The current emergence of Big Data as a full-fledged field can be traced to some of the concepts implemented by these two systems – concepts such as data-intensive computing, infrastructure independence, virtualization and policy-based data management. The two software systems were developed by the Data Intensive Computing Environments group (DICE), which was started in 1994 at the San Diego Supercomputer Center (SDSC) to pursue the goal of implementing software systems that would enable collaborative research through large-scale sharing of multi-disciplinary data files. In the following we track the history of this development.

The selection of the initial software development goal was based on observations of research requirements in computational plasma physics, observations of technology management requirements within the San Diego Supercomputer Center, results from a prior collaboration on an Alternative Architecture study for the Earth

Observing System [1], and research in high-performance networking within the CASA Gigabit Network project [2]. For example, in computational plasma physics, the analysis of the stability of toroidal plasma physics configurations was being done at institutions on the East and West coasts of the United States in the 1980s. A collaboration environment was needed to enable researchers to compare stability analyses and independently verify results. This required the ability to share input files, as well as output results, across institutional boundaries.

Within the San Diego Supercomputer center, which started in 1986, technology was replaced every three years to track and take advantage of the emergence of cheaper and higher performance systems. In particular, by 1994, the third version of an archival storage system had been implemented, using a third generation of tape technology. A mechanism was needed to simplify migration of the archived data between old and new systems.

The Earth Observing System analysis proposed that data products should be organized as a collection, and that relational database technology should be used to manage the system state information. Data replication was proposed between two centers, with data streaming to support processing of the contents.

In the CASA Gigabit Network, theoretical predictions were made of the maximal achievable performance of a distributed, heterogeneous computational environment. The concept of superlinear speedup through the federation of heterogeneous computing resources was analyzed, and a practical demonstration was made that showed a speedup of a factor of 3.3 across two supercomputers. This indicated that management of heterogeneous resources was important for optimizing performance across distributed systems.

The combination of these prior research efforts pointed to the need for researchers to be able to provide a context for interpreting shared data, while managing technology evolution. These requirements for a distributed data management system were the seeds for the development of the first data grid software - the Storage Resource Broker data grid system. Its essential capabilities included:

- Management of data from multiple institutions as a shareable collection through virtualization mechanisms. This was implemented by managing universal name spaces for files, collections, users, and storage systems independently of the physical storage systems where the objects were stored, and independently of the administrative domains at each institution. Authentication and authorization on the universal user name space was implemented as third-party services, mapping logically named resources into physically located storage systems.
- Organization of data files as a collection – independently of the physical characteristics of the data file. That is, a collection provides a virtual “grouping” of files that might be stored on distributed resources of various types, created and owned by multiple users and groups but having some

common properties that warrant bundling them into the same virtual group. Not all objects in the collection need to be files, but can also be dynamic relational queries, sensor streams or self-aggregated/described objects such as tar files or HDF files.

- Association of descriptive metadata with objects in a collection to provide a context for interpreting the data and capturing domain-centric and systems-centric structured information.
- Management of system state information in a relational database. System metadata were associated with files, collections, users, and storage systems. This enabled rapid queries on a much richer set of attributes than normally provided by file systems. The abstraction of a common set of attributes masked the differences between the types of resources being used in the physical layer and provided a uniform system information management layer.
- Management of the properties of the collection, independently of the properties of the storage system in which the files were stored. This was a key goal based on the virtualization of the data collection instead of the virtualization of the storage systems.
- Implementation of a single sign-on authentication system. The files that were shared were owned by the data grid. Users authenticated to the data grid, and in turn, the data grid authenticated itself to the remote storage system. The files were stored under an account that represented the data grid. This meant that the data grid had to both authenticate users, and authorize actions on resources and data. Access controls were managed by the data grid independently of the administrative domain – again providing a common service across the distributed environment.
- An architecture based on a peer-to-peer server environment. Users could connect to any server and the data grid would redirect the request to the correct location for the desired file operation. This meant that users could request a file without knowing where the file was located, without knowing the local name of the file (physical path name), without having an account of the remote storage system, and without knowing the network access protocol required by the storage system. The data grid managed the mapping from the logical file name to the physical path name, managed information about the file location, translated the request by the user client to the protocol required by the remote storage location, and initiated operations on behalf of the user.
- Fault-tolerant semantics. The intent was to build a system that tolerated failure, by redirecting data storage to locations that could provide the space. This was implemented through the concept of storage resource groups. Writing to a resource group succeeded when a file was written to at least one member of the group. Thus some of the storage systems could be off-line, or down for maintenance, and the success of the operation could still be ensured. Another type of fault tolerance was achieved through replication. Since the data grid provided a mapping from the logical name to the physical

address location, it was easy to extend this mapping to multiple physical addresses – hence providing management of synchronized copies of a data object distributed across multiple resources. If access to one copy was unavailable, the system automatically provided access to its replica.

II. Storage Resource Broker:

The development of the Storage Resource Broker was funded initially by DARPA through the “Massive Data Analysis Systems” project [3]. The effort to build software to manage distributed data was viewed as a sufficiently risky objective to warrant DARPA funding. When the approach was presented at a meeting with the tape storage vendor Storage Tek, the response was that they were used to leading edge projects, but the DICE group was halfway down the cliff. The initial development integrated multiple types of technology:

- Use of relational database technology to manage the system state information. As part of the EOSDIS alternative architecture study (1994), a centralized architecture was proposed in which all data were managed by a relational database. The SRB data grid was designed to store system state information in a relational database, while maintaining links to files on distributed storage systems. At that time, holding and accessing hierarchical path information in relational systems was considered to be a performance bottleneck. We chose to do this in order to achieve scalability, since the file systems at that time dealt with less than 2 million files.
- Virtualization of data collections versus virtualization of storage. The SRB focused on managing the properties of the data collection, instead of managing the storage systems. This made it possible to implement operations directed at data manipulation in addition to data storage. Vendors were beginning to implement storage virtualization but considered data/collection virtualization to be too risky.
- Support for heterogeneous storage systems. In order to manage interactions with multiple types of storage system protocols, the SRB software was designed to map from a standard protocol that was based on extensions to Posix I/O, to the protocol used by a specific type of storage system such as the IBM High Performance Storage System, the UniTree storage system, the Network File System, and the Cray File System etc. The protocol conversion was implemented as a modular and extensible software driver. The data grid tracked all operations performed through the middleware, and updated persistent state variables consistently.
- Extended support for data manipulation operations. The SRB data grid implemented operations for replication, versioning, synchronizing, auditing, aggregation in containers, staging of files, archiving of files, checksum creation, metadata extraction, and metadata loading. Since the additional operations were initiated through both Unix utilities and web browsers, a key property of the data grid was the decoupling of access mechanisms from the data management middleware.

- Support for multiple types of client interfaces. A second layer of virtualization was needed to manage mapping from the protocol used by client software, to the standard I/O protocol supported within the data grid. For the SRB, the clients that were supported included web browsers, Java load library, C I/O library, and Fortran I/O library. The protocol used by the client did not have to match the protocol required by the storage system. In effect, the SRB implemented brokering technology between clients and storage.
- Support for multiple authentication environments. Since the data grid managed information that were spread across multiple administrative domains, it needed to deal with the different types of authentication that were supported by these system administrations. To perform authorization across users as well as files, multiple types of authentication systems were supported, including Unix passwords, Kerberos, and Grid Security Infrastructure through the Generic Security Service API. For each type of authentication environment, the associated information was stored in the metadata catalog as attributes on the user account name. The authentication mechanism used to authenticate a person to the data grid did not have to be the same as the authentication mechanism used to authenticate data grid access to a remote storage system. Hence, the system also worked as an authentication broker.
- Schema indirection. Each user community had different definitions for the descriptive metadata that they associated with files and collections. Schema indirection was used to store a triplet consisting of the attribute name, the attribute value, and an attribute unit or comment. This allowed each community to use the data grid as generic infrastructure and implement their domain specific descriptive metadata. Association of name spaces to form an entity set (e. g. Dublin Core, FITS metadata, DICOM metadata, etc.) was also possible.
- Extensible generic infrastructure. Since multiple types of applications built upon the SRB data grid, new features were implemented through appropriate forms of virtualization. This ensured that the system would remain compatible with prior versions, and that extensions to the software could build upon multiple versions of storage technology. The highly extensible architecture ensured long-term sustainability of the software through continued application to additional science and engineering domains.

The SRB can be viewed as an interoperability mechanism that enabled use of multiple types of storage technology, multiple types of authentication systems, and multiple types of access clients. The interoperability enabled by the SRB software is shown in Figure 1. The SRB data grid was implemented as multiple software servers that may reside on different computers or may be co-located on a single computer. Each software server ran as a user-level application on the computer. The servers communicated over a network using a protocol written specifically for

the Storage Resource Broker. External clients accessed the data grid over a network. Each access was authenticated, and each operation was authorized by the data grid. One of the servers managed interactions with a metadata catalog, which in turn composed the SQL needed to interact with a relational database that stores the catalog. The SRB had drivers for interacting with multiple types of storage systems (tape archives, file systems, objects in databases, object ring buffers) and multiple databases (DB2, Oracle, Sybase, Postgres, MySQL, and Informix). Any of the listed clients (C library, Java, Unix shell command, C++ library, web browser, Kepler workflow actor, Python load library, Perl load library, Dspace digital library, GridFTP transport tool) could discover, retrieve, or load files within the distributed environment using a uniform API or the SRB communication protocol.

The development of the SRB was funded by 22 projects that represented collaborations with groups sharing data, groups managing large-scale distributed data, groups organizing digital libraries, and groups building preservation environments. The very wide range of applications ensured that generic infrastructure was developed, with appropriate virtualization mechanisms used to support the domain features of each application.

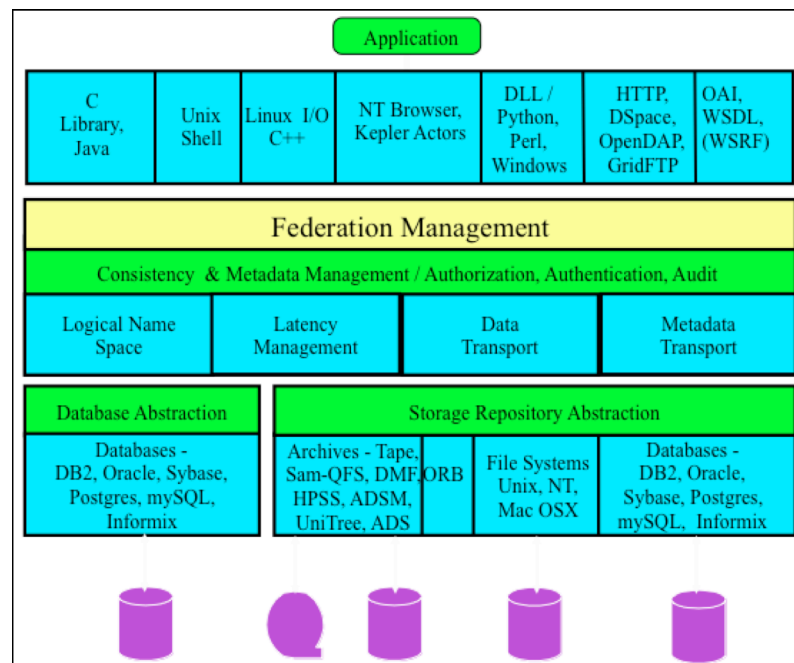


Figure 1. Storage Resource Broker Data Grid Components

III. Data Management Concepts:

Within each collaboration, data management concepts were developed to represent how generic infrastructure could be used to support all types of data management applications. The concepts are useful in that they help define standard semantics for discussing data management. In many cases, the DICE group had to invent terms, or extend the meaning of terms in order to describe what was being done. Eventually, many terms gained broader acceptance within the academic world. Each example of a concept is illustrated within the context of the collaboration project that supported the development of the associated generic infrastructure. We describe the various concepts and their timeline during the SRB development.

Logical File Name and Collection (1996): In the SRB data grid, we needed a term that differentiated the name space used to organize distributed data from the names used within the physical file system. We used the term “logical file name” to denote the identifier for a file as managed by the data grid. The “logical file name” could be organized into “logical collections”, making it possible to associate files that were stored on different storage systems within the same logical collection.

Data Grid (1998): A data grid is the software infrastructure that organizes distributed data into a shareable collection. A paper describing the Storage Resource Broker data grid was presented at the CASCON conference in 1998 [4]. This paper subsequently won an award as one of the top fourteen CASCON First Decade High Impact Papers.

Middleware definition (1998): At an NSF middleware workshop, the question of “What is middleware?” was discussed [5]. The answer based on the SRB data grid was:

“Middleware is the software system that manages distributed state information.”

This definition was extended to include support for services across the distributed environment, but the relationship of middleware to network infrastructure was not codified. Data grid middleware manages distributed state information about file location and membership in collections. Networks also manage distributed state information within their routing tables. The resolution of this dichotomy was recently achieved within the iRODS data grid software, with the integration of policy-based data management with policy-based network routing. See the concept Software Defined Networks.

Persistent Archive (2000): In the Transcontinental Persistent Archive Prototype, a project funded by the National Archives and Records Administration, the DICE group needed a term to describe the preservation of an archive [6]. Note that the word archive (from the computer science discipline) is used to denote the infrastructure that is used to preserve records. In the preservation community, the word “archives” is used to denote the records that are being preserved. A “persistent archive” provides a way to archive a collection independently of the preservation environment, and then retrieve the archives for instantiation of the preservation environment on new technology, overcoming technology obsolescence.

Preservation through Interoperability Mechanisms (2000): There is an equivalence between access to heterogeneous resources across space versus access to heterogeneous resources over time. At the point in time when records are migrated to new technology, data grid middleware can provide the interoperability mechanisms that enable access to both the old and the new technology [7]. Thus preservation infrastructure needs to provide the virtualization mechanisms that abstract preservation properties from the current choice of storage technology. In a sense, application of interoperability across spatial resources was taken to the next

level by providing interoperability across time. The SRB provided a convenient mechanism for performing the temporal jumps in a seamless manner. What resulted is an “organic system” that enabled migration of data objects across time overcoming technology obsolescence through codification of infrastructure independence.

Persistent Objects (2003): Preservation communities previously considered two basic approaches for long term preservation: 1) Emulation, in which the supporting software infrastructure was emulated to ensure that the record could be parsed using the original application; 2) Transformative migration, in which the format of the record was transformed to the format that could be parsed by modern display applications. Persistent objects is a third approach, in which the record is preserved in an unaltered form, while the preservation environment virtualizes I/O operations, enabling access to the record by modern access protocols. This viewpoint considers that the purpose of the preservation environment is to provide an interface between an original record and the ever-changing data management technology.

Consider Figure 2. Data grid technology implements persistent objects [8] by mapping from the actions requested by the display application to the protocol of the storage system where the record is located. In the iRODS data grid, this concept was extended to include the ability to write a display application in a rule language, ensuring independence from the original operating system that was used to support the display application. In both cases, the original record was not changed. Instead the preservation environment was modified to support interactions with new technologies.

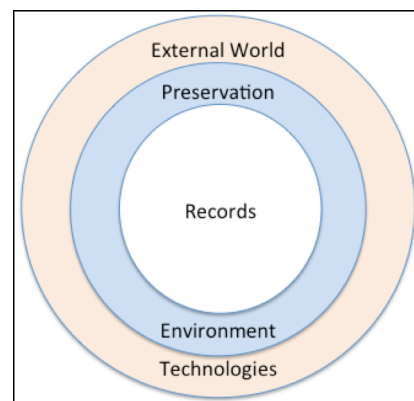


Figure 2. Managing Technology Evolution – Persistent Objects

Policy-based Data Management (2006): One of the applications of the Storage Resource Broker was in the United Kingdom eScience Data Grid. The SRB ensured consistency by encoding within the software middleware explicit management constraints. The constraints were applied by each of the distributed servers, ensuring that the properties of the system were appropriately conserved. However, within the UK data grid, incommensurate management constraints were needed. An archive collection was desired in which no changes to records was allowed, not even by the data grid administrator. Also, a publication collection was desired in which the data grid administrator could replace bad files. Finally, a research collection was needed in which a researcher could replace files at will. Three different management policies were needed within the same data grid.

In the iRODS policy-based data management system, we identified each location in the software middleware where consistency constraints were imposed, and replaced the control software with a policy-enforcement point. On execution of the policy-enforcement point, the system would retrieve the appropriate rule from a rule base, and then execute the associated procedure. The rule controlled the procedure using state information stored in the data grid metadata catalog. Thus the rule could retrieve the name of the collection, and then enforce the appropriate deletion policy. This enables virtualization of policy management, providing both administrators and users with a declarative way to define and control actions that happen at the data storage level. Hence, one can view iRODS as defining a new generation of servers that is completely configurable and capable of enforcing user-centric actions.

Preservation as Communication with the Future (2008): The projects sponsored by the National Archives and Records Administration focused on development of an understanding of the principals behind data preservation. The traditional preservation objectives are authenticity, integrity, chain of custody, and original arrangement. These objectives are all aspects of a higher level goal, that of enabling communication with the future. The traditional representation information defined by the Open Archival Information System model provides a context for correctly interpreting a record by a future knowledge community through creation of preservation metadata. In the future, the knowledge community will have enough information from the associated representation information to correctly interpret a record. This viewpoint needed to be augmented with a characterization of the representation information that describes the preservation environment. Within policy-based data management systems, the environment representation information is characterized by the policies and procedures that are used to manage the records along with the associated system state information. It is then possible for an archivist in the future to verify communication from the past, and validate that the preservation objects have been appropriately preserved [9].

If preservation is communication with the future, then policy-based systems enable verification of the validity of communication from the past.

IV. Integrated Rule Oriented Data System

In 2006, the Storage Resource Broker development was deprecated, in favor of developing an Open Source version of data grid technology. At the same time, a decision was made to go beyond data and information virtualization, to also support knowledge virtualization. The basic approach was to turn policies into computer actionable rules, turn procedures into computer executable workflows, and use policy enforcement points to decide when policies should be applied.

The architecture of the policy-based data management systems was similar to the SRB, as shown in Figure 3. Multiple peer-to-peer servers managed interactions with remote storage locations, and a central metadata catalog stored state information in a relational database. The integrated Rule-Oriented Data System (iRODS) also implemented servers to manage message passing, and to manage a queue of

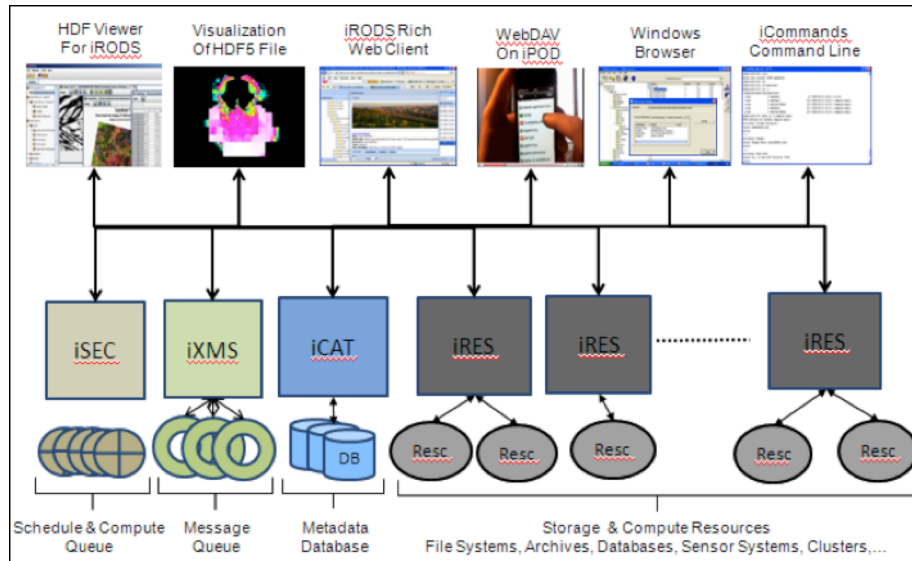


Figure 3. Policy-based Data Management Architecture

outstanding rule requests [10].

A comparison of policy-based systems with distributed data management systems shows how the concepts related to data management have been evolving. Figure 4 illustrates the central concepts behind traditional file systems, and also behind the Storage Resource Broker. External events interact with the data management system through a well defined protocol. The data management system uses state information to control the execution of operations on the stored files, and the state information is appropriately updated. The file system (i-nodes, v-nodes, etc.) environment in some sense is synonymous with the state information that is managed about the files. A key component of a file system is the consistent update of the state information after every operation that is performed upon the files. The SRB answered the challenge of self-consistent update of state information in a distributed environment, across heterogeneous storage systems, across multiple administrative domains.

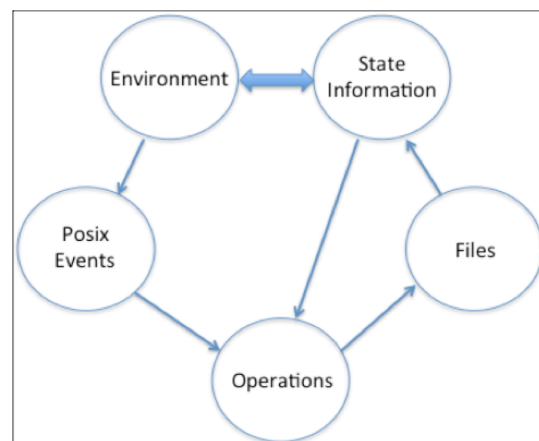


Figure 4. File System Characterization

In policy-based data management systems, operations are replaced by policies that control updates through procedures, and files are replaced by objects that may include workflows, active or realized objects, and databases, as well as files. Figure 5 lists the characteristics of policy-based data management, representing the evolution from traditional file-based systems to information and knowledge based systems. As before, the data management environment is synonymous with the consistent management of state information.

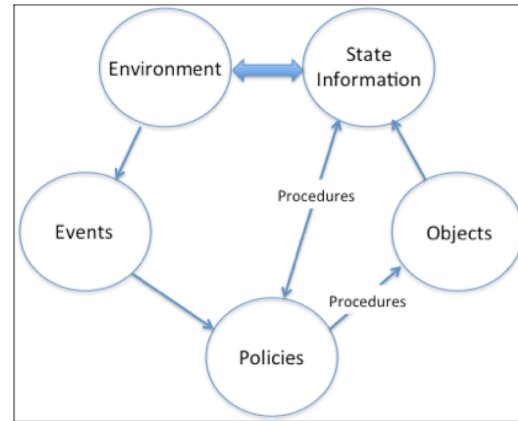


Figure 5. Policy-based System Characterization

However, in the policy-based system, the environment is governed by the set of policies that are implemented as computer actionable rules. Thus a description of the environment must include not only the state information, but also the policies and procedures that are being enforced. Similar to the SRB, the development of iRODS also required several new concepts, which we describe along with a timeline.

Computer Actionable Knowledge (2012): A major goal of data grid technology has been the extension of data management systems to also support information management and knowledge management through computer actionable forms. The Storage Resource Broker augmented data management with information management, by associating state information as metadata attributes on an appropriate name space. The types of information that were managed included provenance information, descriptive information, representation information, and system administrative information.

Policy-based data management systems augment information management with knowledge management. The knowledge required to execute a protocol, or manipulate a file, or access a remote repository is encapsulated in procedures, known as micro-services. In a sense, a file (or object) is not viewed in isolation, but along with all policies and procedures that governs its usage and existence. The application of knowledge requires the dynamic execution of procedures. The result of the execution is stored as system state information, and is assigned as metadata on objects within a name space. In effect, the reification of a knowledge procedure is turned into administrative information that is stored as metadata in a relational database. One can view the metadata as inherent properties (labels) on the objects that codify the derived knowledge obtained through application of procedures.

This approach to knowledge management through computer actionable forms can be quantified as follows:

- Data consists of bits (zeros and ones)

- Information consists of labels applied to data
- Knowledge evaluates relationships between labels
- Wisdom imposes relationships between relationships.

Within the iRODS data grid, data are managed as files in a file system, or objects in an object store. Information is managed as metadata in a relational database. Knowledge is applied as computer actionable rules through a rule engine. Wisdom (within the confines of the user-configurable iRODS system) is applied through policy enforcement points which determine when and where the knowledge procedures should be executed.

Note that the concept of relationships has been extended to include:

- Semantic or logical relationships
- Spatial or structural relationships
- Temporal or procedural relationships
- Functional or algorithmic relationships
- Systemic or epistemological relationships

Thus a procedure is the application of a functional relationship to a digital object to generate either information about the digital object, or a new digital object [11].

The differentiation between information and knowledge is complex. In order to assign a label to a digital object, a knowledge relationship between existing labels needs to be evaluated. However each existing label required the prior application of knowledge relationships. Information generation is an infinite recursion on the application of knowledge procedures. Each knowledge procedure evaluates relationships between labels that were previously generated. The recursive nature is closed by reducing the information labels to a well known set that are interpreted the same way by the entire user community. The simplest way to separate information and knowledge is to view information as the reification of knowledge. Information is a static property, while knowledge is the active evaluation of a relationship.

The first attempt to characterize information and knowledge was expressed as a matrix, with the goal of differentiating between ingestion, management, and access services for digital objects [12]. This characterization focused on services that were used to manipulate data, information and knowledge, within the context of a data grid. Figure 6 shows the components of the characterization, with the data grid represented by the matrix that links together the individual components related to the types of service.

This characterization is realized in the iRODS policy-based data management system. The services to manipulate data are the operations supported upon digital objects. The storage systems for data are accessed through storage drivers. The services to manipulate information are the operations supported upon metadata attributes. The information repository is the metadata catalog, stored in a relational database. The knowledge relationships between concepts are implemented as micro-services that are controlled by computer actionable rules. The knowledge repository is implemented as a rule base. The knowledge-based grid can be viewed spatially, as a shared distributed service or temporally, as a persistent archive.

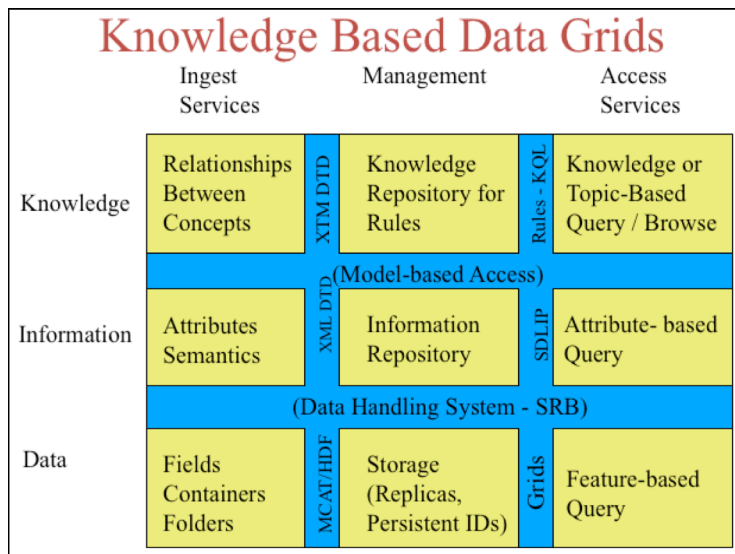


Figure 6. Knowledge-based Grids

The access services remain an area of active development, and are further discussed in the feature-based indexing concept.

Knowledge Virtualization (2010): The iRODS data grid provides virtualization of data, information, and knowledge. Figure 7 shows a simple architecture view of the interoperability mechanisms. An access interface virtualizes access by mapping from the access protocol to the iRODS interaction protocol. Each interaction is trapped at policy enforcement points where a rule base is consulted to determine which policy to execute. The policies control the execution of procedures that are composed by chaining together basic functions, called micro-services. This requires that the middleware manage exchange of structured information between the chained micro-services.

The micro-services perform operations such as I/O manipulation, metadata extraction, and domain-specific

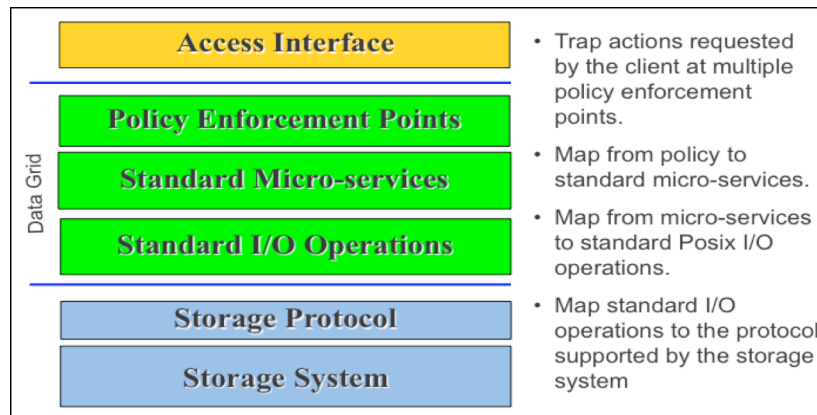


Figure 7. iRODS data grid virtualization mechanisms

operations. Each micro-service invokes standard Posix based I/O operations. The data grid middleware then translates between the standard I/O and the protocol required by the remote storage location. Thus the micro-services are operating system independent. The same micro-services run on Windows, Unix, and Mac computers, enabling the migration of policies and procedures across operating systems. The ability to manage application of knowledge procedures, independently of the choice of storage environment, can be viewed as a form of knowledge encapsulation.

Policies as Intellectual Property (2013): A major goal of the development of policy-based data grid middleware has been the conversion of management policies into computer actionable rules that control computer executable procedures. This enabled multiple communities, shown below, to apply the technology. The users of the software span multiple science and engineering disciplines, and include national data grids, national libraries, and international projects:

Archives	Taiwan National Archive, Chronopolis
Astrophysics	Auger supernova search
Atmospheric science	NASA Langley Atmospheric Sciences Center
Biology	Phylogenetics at CC IN2P3
Climate	NOAA National Climatic Data Center
Cognitive Science	Temporal Dynamics of Learning Center
Computer Science	GENI experimental network
Cosmic Ray	AMS experiment on the International Space Station
Dark Matter Physics	Edelweiss II
Earth Science	NASA Center for Climate Simulations
Ecology	CEED Caveat Emptor Ecological Data
Engineering	CIBER-U
High Energy Physics	BaBar / Stanford Linear Accelerator
Hydrology	Institute for the Environment, UNC-CH; Hydroshare
Institutional Repositories	Carolina Digital Repository
Genomics	Broad Institute, Wellcome Trust Sanger Institute, NGS
Libraries	French National Library, Texas Digital Libraries
Medicine	Sick Kids Hospital
Neuroscience	International Neuroinformatics Coordinating Facility
Neutrino Physics	T2K and dChooz neutrino experiments
Oceanography	Ocean Observatories Initiative
Optical Astronomy	National Optical Astronomy Observatory
Particle Physics	Indra multi-detector collaboration at IN2P3
Plant genetics	the iPlant Collaborative
Quantum Chromodynamics	IN2P3
Radio Astronomy	Cyber Square Kilometer Array, TREND, BAOradio
Seismology	Southern California Earthquake Center
Social Science	Odum, TerraPop

Each community implemented different choices for semantics, policies, and procedures. A generalization of the observed usage patterns is to identify the intellectual properties of each community with the policies and procedures that

they implemented. The underlying data grid middleware was generic infrastructure that provided the mechanisms needed to virtualize interactions with data, information, and knowledge. The policies and procedures encapsulated the knowledge that was needed to apply the middleware within each domain.

This means that intellectual property can be captured and applied within generic data management infrastructure to cater to the specific needs of each domain. This idea is extended in Figure 8, which describes a general approach towards quantifying intellectual property.

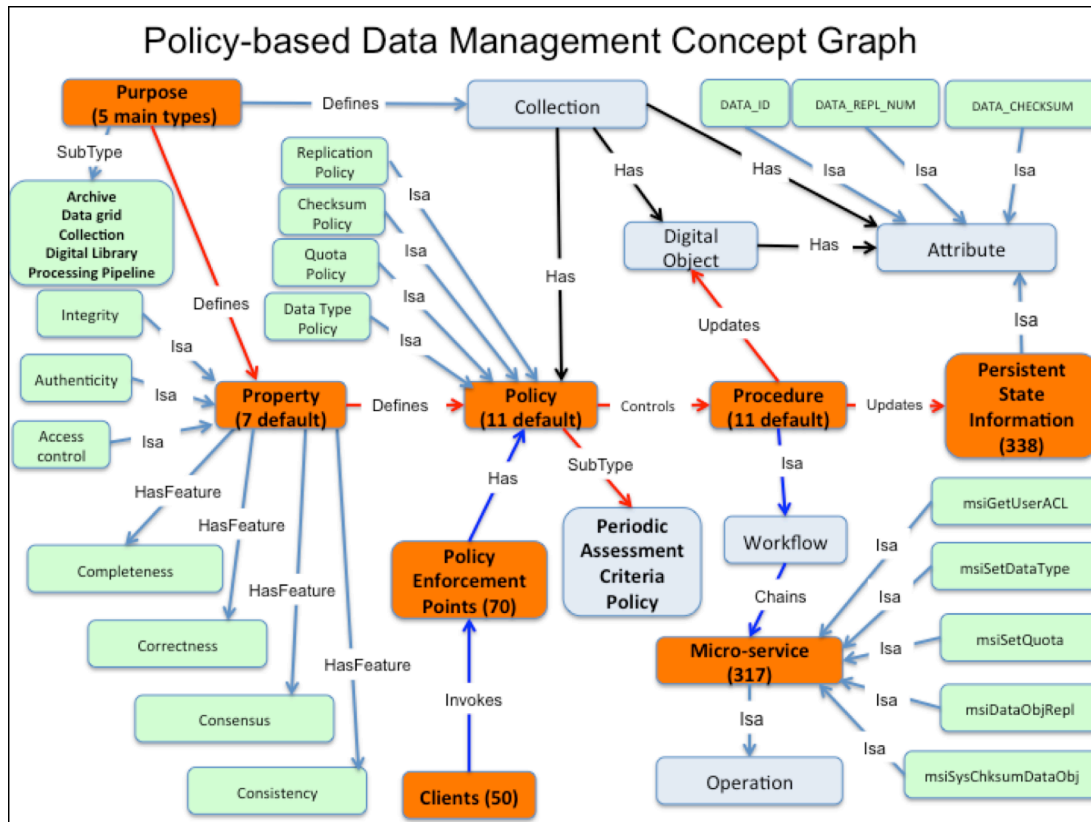


Figure 8. Conceptualizing intellectual property as policies and procedures.

Each domain is characterized by:

- **Purpose** driving the formation of a data collection. The purpose represents a consensus of the persons collaborating on a data management project.
- **Properties** that will be maintained for the data collection. The properties are dependent upon the driving purpose. If the intent is preservation, then properties related to authenticity, chain of custody, integrity, and original arrangement are desired. If the intent is formation of a local project repository, then properties related to file format and access controls may be desired. The properties comprise assertions made about the collection by the developers of the collection. Other domain centric elements (such as requisite metadata etc.) can also be defined as part of these properties.

- **Policies** that enforce the desired properties. The policies control when and where management procedures are executed. Multiple policies may be needed for each desired property. In general, policies are needed to generate the desired property. Policies are also needed to validate whether the desired property has been maintained over time. Since the distributed environment is subject to multiple forms of risk (network outage, storage system maintenance, operator error, policy change), assessment criteria are needed that can be checked to verify compliance with the desired collection properties. Policies are turned into computer actionable rules. Example domain centric policies include enforcing authority (e. g. HIPAA policies), integrity checks, data cleansing, metadata extraction, etc..
- **Procedures** codify policies and apply the operations needed to generate a desired property. Examples include procedures to create a replica, extract metadata, set access controls, manage a quota, check a retention period, apply disposition, etc. Procedures are turned into computer executable workflows.
- **Persistent state information** is generated each time a procedure is run. The persistent state is stored as metadata attributes on one of the name spaces managed by the data grid. The state information can be queried for compliance at a point in time. To verify compliance over time, the system parses audit trails. Persistent state information in turn codify the properties of a collection.

A viable policy-based data management system must be sufficiently sophisticated to handle a wide variety of data management applications. The iRODS data grid provides 317 micro-services that can be used to compose procedures, and manages 338 persistent state information attributes. In practice, each domain implements a small number of policies. Out of the box, the iRODS data grid source provides 11 default policies for enforcing data sharing properties. Communities typically add another 5 policies on the average to control desired features. However, the range of policies that are required to support a fully customized data grid may be very large.

Each policy and procedure set encapsulates the domain knowledge needed to manage a specific domain application.

Federation through Interoperability Mechanisms: Within the DataNet Federation Consortium [14], the iRODS data grid is being used to create national data cyberinfrastructure through the federation of existing data repositories. In the process, interoperability mechanisms have been implemented that enable three basic functions:

1. Micro-services that retrieve data from a remote repository using the protocol of the remote repository. This is a traditional approach similar to brokering, in which data are retrieved for analysis at the local computer. The data are moved to the processing engine.

2. Middleware servers that enable application of the desired operations at the remote repository. In this case, the operations are moved to the data.
3. Policies that control where the operations are performed.

Using these three mechanisms, the DataNet Federation Consortium has been able to support interoperability with web services, sensor networks, union catalogs, data repositories, workflow environments, databases, message buses, and systems that communicate over the internet.

The expectation is that these three interoperability mechanisms are sufficient to federate all existing data management applications. The DataNet Federation Consortium currently federates systems across national projects in oceanography, cognitive science, plant biology, engineering, hydrology, and social science.

Quantifying the Broadening of Impact: A notable requirement for National Science Foundation funding is the demonstration that the research results will impact a broad user community. A mechanism has been needed to quantify the impact. One way to do this has been the observation that the set of policies and procedures used to manage a collection evolve over time to represent the current requirements of each broader user community. It is possible to quantify impact by tracking the policy evolution. This can be represented through a collection life cycle:

- Project collection – usually the team members have complete tacit knowledge about the acceptable semantics, data formats, and analysis procedures used with the team data sets. The data sets are organized in a project collection with minimal metadata. The data sharing is limited to the group, and is mostly through shared and mounted file spaces.
- Shared collection – when data products are shared with other groups and institutions, the tacit knowledge must be made explicit. Policies are needed to govern the application of semantic terms, and the transformation of data to required data formats. Policies are also needed to enforce authentication, access controls and data distribution. Policies for data manipulation may also be needed.
- Published collection – when the results are formally published, policies are needed to enforce domain standards for semantics and data formats. Policies are also needed to generate persistent identifiers, to validate integrity, and to track provenance.
- Processing pipeline – when the data sets are used in an analysis service, procedures are needed that support the manipulation and transformation of the data.
- Preserved reference collection – when the results are archived for use by future researchers, a sufficient context is needed that enables a person in the future to interpret the data. This is typically encapsulated in representation information. At the same time, the policies and procedures also need to be preserved so a future archivist can verify that the collection was managed correctly.

The broadening of user impact can be quantified through the evolution of the policies and procedures that are used to manage a data collection.

V. Future Data Management Infrastructure and Conclusion

The current generation of data grid middleware is still evolving. New opportunities to apply policies to control the data management environment are emerging. We consider three specific extensions, the inclusion of policies within storage controllers, the integration of policy-based data management with policy-based networks, and the extension of a knowledge grid into a wisdom grid.

Feature-Based Indexing: A major challenge in constructing a collection is the assignment of appropriate descriptive metadata. This is a laborious task, which potentially is non-scalable. A major question is whether the act of description can be turned into the application of a knowledge procedure, that is automatically applied by a policy-based system. Normally descriptive metadata are used to provide a context for the contents of a file. An alternative approach is to use descriptive metadata to define features present within a file. If the desired features can be extracted by a knowledge procedure, then the generation of descriptive metadata can be automated.

This approach is being explored in collaboration with storage vendors. The Data Direct Networks storage controllers now support virtual machine environments that can be used to run the iRODS data grid. When a file is written to the storage system, the data grid can apply feature extraction procedures automatically, and index the stored data by the features present within each record. Hence, one can construct a domain-centric data-grid appliance that can perform automated data management including automated data description.

Software Defined Networks (2013): Policy-based systems are also appearing within networks that are based on the OpenFlow router. Routing decisions can be controlled by policies that are used to manage path-selection within the router. A demonstration of the use policy-based data grids to control policy-based routing was given at the Supercomputing '13 conference [15]. The iRODS data grid managed information about the location of files, their access controls, and the availability of replicas. Within the iRODS data grid, a parallel data transfer was set up, with subsets of the file sent in parallel over the network. The iRODS data grid communicated with the OpenFlow router to select a disjoint network path for each of the parallel data transfer channels.

The idea here is that a traditional data grid views the network as a black box (and vice versa, the network is opaque with respect to the applications at the end-points of the communication pipeline). If the data grid is able to export some of its policies to be implemented by the network (through the OpenFlow router) and also is able to get feedback from these routers about network topology, congestion and

statistics, the two can work together to mutual advantage and improve performance. Having this exchange of information can be used in multiple ways to improve data grid operations.

One way to exchange information is through the integration of control policies between data grids and networks. Since both systems are managing distributed state information, it is reasonable to think about formally moving data grid middleware into network routers. It will then be possible to access data by name (or metadata attribute) instead of an IP address, enforce access controls within the network, cache data within the network, and debug data transfers by single-stepping through the data grid procedures (currently supported in iRODS).

The approach would rely upon the data grid to provide a context for the files through their organization in collections. A file would be referenced by its membership in a collection, with the data grid controlling the access (authentication and authorization). The data grid would negotiate with the network for selection of the replica to use as the starting point, and the network path to use for data delivery. In the long term, data grid middleware should disappear as separate infrastructure, and be subsumed within the network. The upshot of this would be collection-oriented addressing of objects instead of name-oriented or ip-oriented addressing for data ingestion, movement and access.

Wisdom: Current virtualization mechanisms focus on data, information, and knowledge. Future data management systems will also need to support virtualization of wisdom. If we can think of wisdom as the evaluation of relationships between relationships, then we can build a computer actionable form of wisdom. Within the iRODS data grid, wisdom is captured as hard-coded policy-enforcement points. To make application of wisdom a dynamic process, the system will need to implement mechanisms that enable wisdom-based decisions to be selected as systemic processes that apply to all interactions. This will require processing information about each access session, information about the collections, and information about the user community to infer which set of knowledge procedures should be applied.

Acknowledgements:

The original members of the DICE group included Reagan Moore, who led the group; Michael Wan, who was the chief architect and designed the communication and management protocols for the peer-to-peer server architecture; Arcot Rajasekar, who proposed and implemented key innovations related to virtualization, management of state information, and policy-based data management; Wayne Schroeder, who implemented the security environment, unix-style utilities, and the testing environment; Lucas Gilbert, who implemented a Java I/O library; Sheau-Yen Chen, who was the grid administrator for the group through multiple evolutions of the grid; Bing Zhu, who ported the system to Windows; and Chaitan Baru, Richard Marciano, Ilkay Altintas, Bertram Ludaescher, and Amarnath Gupta who applied the technology. The DICE group maintained a core set of

developers for twenty years, while adding expertise including Mike Conway, who developed advanced Java library interfaces; Hao Xu, who optimized and extended the iRODS distributed rule engine, and Antoine de Torcy, who developed iRODS micro-services for application domains.

The SRB and iRODS technologies were developed and applied across more than 30 funded projects, and multiple funding agencies. These included:

EarthCube Layered Architecture	NSF	4/1/12-3/31/13
DFC Supplement for Extensible Hardware	NSF	9/1/11-8/31/15
DFC Supplement for Interoperability	NSF	9/1/11-8/31/15
DataNet Federation Consortium	NSF	9/01/11-8/31/16
SDCI Data Improvement	NSF	10/1/10 – 9/30/13
National Climatic Data Center	NOAA	10/1/09 – 9/1/10
Subcontract: Temporal Dynamics of Learning Center	NSF	1/1/10 – 12/31/10
NARA Transcontinental Persistent Archive Prototype	NSF	9/15/9-9/30/10
Subcontract: Temporal Dynamics of Learning Center	NSF	3/1/9-12/31/9
Transcontinental Persistent Archive Prototype	NSF	9/15/8-8/31/13
Petascale Cyberfacility for Seismic Community	NSF	4/1/8-3/30/10
Data Grids for Community Driven Applications	NSF	10/1/7 9/30/10
Joint Virtual Network Centric Warfare	DOD	11/1/6-10/30/7
Petascale Cyberfacility for Seismic	NSF	10/1/6-9/30/9
LLNL Scientific Data Management	LLNL	3/1/5 – 12/31/08
NARA Persistent Archives	NSF	10/1/4-6/30/08
Constraint-based Knowledge	NSF	10/1/4-9/30/6
NDIIPP California Digital Library	LC	2/1/4-1/31/7
NASA Information Power Grid	NASA	10/1/3-9/30/4
NARA Persistent Archive	NSF	6/1/2-5/31/5
National Science Digital Library	NSF	10/1/2-9/30/6
Particle Physics Data Grid	DOE	8/15/1-8/14/4
SCEC Community Modeling	NSF	10/1/1-9/30/6
Terascale Visualization	DOE	9/1/98-8/31/02
Grid Physics Network	NSF	7/1/00-6/30/05
Persistent Archive	NARA	9/99-8/00
Digital Library Initiative UCSB	NSF	9/1/99-8/31/04
Digital Library Initiative Stanford	NSF	9/1/99-8/31/04
Information Power Grid	NASA	10/1/98-9/30/99
Persistent Archive	NARA	9/98-8/99
NPACI data management	NSF	10/1/97-9/30/99
DOE ASCI	DOE	10/1/97-9/30/99
Distributed Object Computation Testbed	DARPA/USPTO	8/1/96-12/31/99
Massive Data Analysis Systems	DARPA	9/1/95-8/31/96

References:

1. Frank Davis, William Farrell, Jim Gray, Roberto Mechoso, Reagan Moore, Stephanie Sides, Michael Stonebraker, "EOSDIS Alternative Architecture", Submitted to HAIS, 1994.
2. Reagan Moore, "Distributing Applications Across Wide Area Networks", General Atomics Technical Report GA-A20074, May 1990.

3. Reagan Moore, Chaitanya Baru, Richard Frost, Richard Marciano, Arcot Rajasekar and Michael Wan}, "MDAS - A Massive Data Analysis System", In Proceedings of Interface97 Symposium, 1997.
4. Chaitanya Baru, Reagan Moore, Arcot Rajasekar, Michael Wan, "The SDSC Storage Resource Broker," Proc. CASCAN'98 Conference, Nov.30-Dec.3, 1998, Toronto, Canada, p. 5.
5. B. Aiken, J. Strassner, B. Carpenter, I. Foster, C. Lynch, J. Mambretti, R. Moore, B. Teigelbaum, "Network Policy and Services: A Report of a Workshop on Middleware", The Internet Society, 2000.
6. Reagan Moore, Chaitan Baru, Arcot Rajasekar, Bertram Ludaescher, Richard Marciano, Michael Wan, Wayne Schroeder, Arcot Rajasekar, "Collection-based Persistent Digital Archives – Part 2", D-Lib Magazine, April 2000.
7. Reagan Moore, Chaitan Baru, Arcot Rajasekar, Bertram Ludaescher, Richard Marciano, Michael Wan, Wayne Schroeder, Amarnath Gupta, "Collection-Based Persistent Digital Archives – Part 1", D-Lib Magazine, April, 2000.
8. Reagan Moore, "The San Diego Project: Persistent Objects," Archivi & Computer, Automazione E Beni Culturali, l'Archivio Storico Comunale di San Miniato, Pisa, Italy, February, 2003.
9. Reagan Moore, "Towards a Theory of Digital Preservation", IJDL Volume 3, Issue 1, pp. 63-75, August 2008.
10. Arcot Rajasekar, Michael Wan, Reagan Moore, Wayne Schroeder, "A Prototype Rule-based Distributed Data Management System", HPDC workshop on "Next Generation Distributed Data Management", May 2006, Paris, France.
11. Reagan Moore, "Automating Data Curation Processes", NSF Curating for Quality Workshop, September 11, 2012, Arlington, Virginia.
12. Reagan Moore, "Knowledge-based Grids," Proceedings of the 18th IEEE Symposium on Mass Storage Systems and Ninth Goddard Conference on Mass Storage Systems and Technologies, San Diego, April 2001.
13. Arcot Rajasekar, Michael Wan, Reagan Moore, Wayne Schroeder, Sheau-Yen Chen, Lucas Gilbert, Chien-Yi Hou, Christopher Lee, Richard Marciano, Paul Tooby, Antoine de Torcy, Bing Zhu, "iRODS Primer: Integrated Rule-Oriented Data System", Morgan & Claypool, 2010.
14. R. Moore., A. Rajasekar, "Reproducible Research within the DataNet Federation Consortium", International Environmental Modeling and Software Society 7th International Congress on Environmental Modeling and Software, San Diego, California, June 2014, <http://www.iemss.org/society/index.php/iemss-2014-proceedings>.
15. S. Huang, H. Xu, Y. Xin, "A Framework for Integration of Rule-oriented Data Management Policies with Network Policies", the 3rd GENI Research and Educational Experiment Workshop [GREE 2014], Atlanta, GA, March 2014.

Appendix A

Multiple versions of the SRB and iRODS software were developed:

SRB releases

SRB 3.5	Dec 3, 2007	Bind variables, bulk replication, transfer restart
SRB 3.4	Oct 31, 2005	Master/slave MCAT, HDF5 integration
SRB 3.3	Feb 18, 2005	ACL inheritance, bulk move, GT3 GSI
SRB 3.2	July 2 2004	Client initiated connections, Database access
SRB 3.1	April 19, 2004	Synchronization, trash can, checksums
SRB 3.0	Oct 1, 2003	Federation
SRB 2.0	Feb 18, 2003	Parallel I/O, bulk load, metadata access control
SRB 1.1.8	Dec 15, 2000	Encrypted passwords, large file size
SRB 1.1.7	May 2000	GSI authentication
SRB 1.1.6	Nov 1999	Stream support, Oracle support
SRB 1.1.4	May 1999	Containers
SRB 1.1.3	Feb 1999	Recursive replication
SRB 1.1.2	Dec 1998	Monitoring daemon
SRB 1.1	Mar 1998	Query support
SRB 1.0	Jan 1998	Unix commands

iRODS releases

iRODS 3.3.1	Feb 24, 2014	SHA2 hash, Rule looping, WSO extensions
iRODS 3.3	July 17, 2013	NetCDF support, HDFS, PAM authentication
iRODS 3.2	Oct 3, 2012	WSO objects, direct access resources
iRODS 3.1	March 16, 2012	Tickets, locks, group-admin updates
iRODS 3.0	Sept. 30, 2011	New rule language, soft links
iRODS 2.5	Feb 24, 2011	Database resources, Fortran I/O library
iRODS 2.4	July 23, 2010	Bulk upload, monitoring,
iRODS 2.3	March 12, 2010	Extensible iCAT, quotas, group-admin
iRODS 2.2	Oct 1, 2009	HPSS driver, S3 driver, compound resource
iRODS 2.1	July 10, 2009	mySQL driver, Kerberos, policy enforcement
iRODS 2.0	Dec 1, 2008	federation, master/slave catalog, bundling
iRODS 1.1	June 27, 2008	GSI, mounted structured files, HDF5, Jargon
iRODS 1.0	Jan 23, 2008	Oracle driver, FUSE interface, rule language
iRODS 0.9	June 1, 2007	replication, metadata, trash, integrity checking
iRODS 0.5	Dec 20, 2006	policy enforcement points, rule engine