

# THE CCLRC DATA PORTAL

*Glen Drinkwater, Shoaib Sufi*

*CCLRC – Daresbury Laboratory, Daresbury, Warrington, Cheshire, WA4 4AD, UK.*

*E-mail: g.j.drinkwater@dl.ac.uk*

## **Abstract**

*The project aims to provide easy, transparent access to experimental, observational, simulation and visualisation data kept on heterogeneous systems across many sites. Further more it will provide links to other web/grid services, which will allow the user (i.e. scientist) to further use the selected data, e.g. via data mining, simulations or visualisation. The Data Portal aims to work as a broker between the users, the facilities, the data and other services. Users currently have only very limited support in accessing, managing and transferring their data or indeed in identifying new data resources. In a true Grid environment it is essential to ease many of these processes and the aim of the Data Portal is to provide homogenous interaction with user data resourced and automation of access, management and transfer of users data management needs.*

## **1 Introduction**

Currently scientists are forced to manually relate between all the experimental, data, computing and analysis facilities that are available world wide, with little infrastructure support. In the future it is hoped The Grid will provide these functions, enabling the scientists to choose much more easily from a wide range of services, connecting and combining desired services for an optimal working environment. Much of the access to the Grid is envisaged to take place through customisable, community oriented Portals. A range of projects within Council for the Central Laboratory of the Research Councils' (CCLRC) have been chosen to provide the building blocks of an integrated solution for users of experimental, computing and data facilities, showing how technologies can be used to build middleware components that support high level scientific grid applications. Data will play a pivotal role in the success of Grid/e-Science developments. Virtually all envisaged applications will need to be able to draw from and deliver to the distributed heterogeneous

information/data sources with a variety of contents. Hence three major challenges are posed: data accessibility, data transfer and management of personal data. Data accessibility implies the capability to locate information/data without prior knowledge of its physical location or the form in which its contents is described. Data transfer relates to the problem of large data volumes that need to be transferred across the Internet. Management of personal data is concerned with the growing distribution of data produced by scientists within a Grid environment, which required new ways of keeping track and moving data for single scientists and more importantly for research groups.

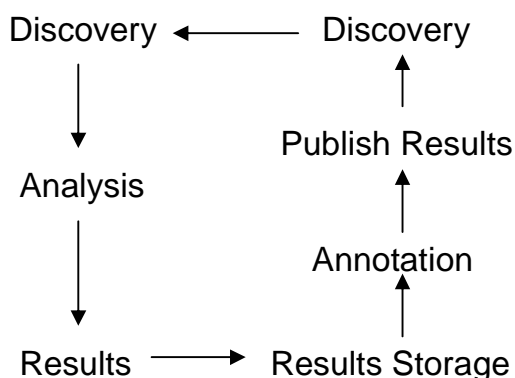
### **1.1. Current Status**

Data Management within CCLRC is currently been used in two projects involving CCLRC, the Environment from the molecular level (e-Minerals) [1] and Simulation of complex materials (e-Materials) projects. There is also a generic CCLRC instance giving access to data from two of our

experimental departments: Synchrotron Radiation (SR at Daresbury) and Neutron Spallation (ISIS at Rutherford (RAL)) as well as data from the British Atmospheric Data Centre (BADC at RAL) and an outside source at the Max Planck Institute for Meteorology in Hamburg, Germany. The latter two are both World Data Centres. The installation is available at <http://dataportal.dl.ac.uk:8080>.

The management of data is achieved using two web service based portals, the Data Portal and the Data Insertion Portal. The Data Portal is for high-level access to multidisciplinary data, a metadata schema that allows the efficient description of data from heterogeneous sources and the Storage Resource Broker [2] (SRB) from the San Diego Super Computing Center (SDSC) to manage the physical location of data both personal as well as archived. The Data Portal is linked to existing data catalogue systems. These catalogues include metadata as well as links to the data itself. The data itself is held in various storage resources from local disks, over databases to multi terabyte tertiary tape systems. The Data Insertion Portal is used to upload data into the system and annotate the data. The annotation of the data creates metadata repositories, which allow the Data Portal to search and locate data.

From a data centric perspective, scientific research can be viewed by the sequence of events shown in Figure 1.



**Figure 1:** Data centric view of scientific workflow.

In general terms, research begins with some data, for example a crystal structure. Some analysis will be performed on this data, which in turn generates more data, i.e. results. These results are then stored and annotated in some fashion after which a subset may or may not be selected for publication or distribution to a wider community.

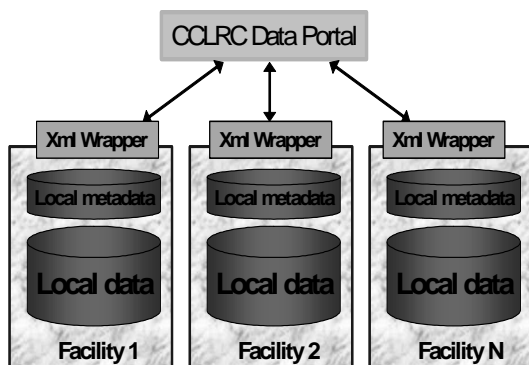
These Portals were designed to work in distributed and heterogeneous environments, our current installations have proven that to be true, integrating a diverse set of system types, operating systems, data resources (databases, Storage Resource Brokers) and sites seamlessly.

## 1.2. Data Portal Architecture

The current version of the Data Portal uses a modular web services model. This is achieved using Apache's Axis implementation of the Simple Object Access Protocol (SOAP) submission to W3C. SOAP is a lightweight protocol for exchange of information in a decentralised, distributed environment. It is a XML based protocol, which defines a framework for representing remote procedure calls and responses.

Using SOAP and web services the Data Portal was decentralised into modules that represent an area of functionality. For example, the Session Manager controls user's state, Authorisation communicates with the MyProxy server to authorise the user to the Data Portal and Query & Reply sends queries to multiple XML Wrappers at each facility. The XML Wrappers basically map the local metadata format into a XML format that the Data Portal understands, allowing seamless integration of multiple facilities metadata information. These services were platform and language independent allowing other services (other portals or clients) to communicate with the Data Portal regardless

of the language that they were written. Each web service transfers an encrypted session id used to obtain state within the core modules of the Data Portal. Any GSI delegation of the user's proxy certificate to another web services is achieved by the delegation of the user's proxy stored in the Session Manager to the web service requiring authentication.



**Figure 2:** Simplified Data Portal Architecture

Vital to this version of the Data Portal is the Lookup module. This is used for the publishing and finding Data Portal web service modules. Essentially this acts as an interface to a Universal Description, Discovery and Integration (UDDI) registry. A module would query the UDDI and receive a Web Services Definition Language (WSDL) file address for the module. This is standard to describe the technical invocation syntax of a web service. A module would use this file to invoke the web service that it needs.

The Data Portal server hosts most of these services but also provides the user interface and manages the interaction with the user and all attached resources. The server provides the user with a web interface to search the existing metadata both on the server itself and the connected data holdings transparently. Incoming requests from the user will be interpreted by the server and a query will be formed and transmitted to the facility's local repositories that are available to the Data Portal. The queries to local repositories are XQueries [3] and are transmitted via web services. The result from

the various local repositories is expected in XML format or the format that the XQuery requested (i.e. HTML). The Data Portal will collate the results and generate the required pages to display the results. The Server is also responsible for the user authentication and session control. In the future the server is also expected to liaise with other data portals as well as other grid/web services.

CCLRC has developed a special multidisciplinary metadata model (The Data Portal project uses an XML implementation) [4] to be able to integrate and make available data from various scientific topics ranging from astronomy to physics. In the following paragraph we will describe how these resources are connected with the server.

Other repositories are expected to have either their own metadata catalogue systems or their own metadata formats describing their data holding or use an extension of the CCLRC Scientific Metadata Model (CSMD). To integrate them each catalogue (facility) will be accompanied by an XML Wrapper, which firstly converts the local metadata catalogue systems into CSMD. The XQuery request from the Data Portal server is executed against the CSMD and results returned to the Data Portal. Currently the metadata catalogues include links to the data location, therefore the data can be transferred to the user (e.g. FTP/HTTP/GridFTP/SRB) or to a third party machine (via e.g. GridFTP).

The core Data Portal system offers the user the possibility to collect all relevant datasets / data files in his personal shopping basket, which can be kept from one session to the next if required. This shopping basket then offers the user a range of functionalities including transfer (using GridFTP, download), delete (from shopping basket), or if available offer other grid/web services to the type of data.

### 1.2.1. Authorisation

The authentication of data within the Data Portal [5] is done by the GSI delegation [6] of the user's proxy certificate to each facility. At each facility sits an Access & Control (ACM) and a XML Wrapper web service. Upon logging on, the user's proxy certificate is delegated to each facility's ACM. The ACM maps the user's distinguished name (DN) to a local user on their system and the access rights are given back to the Data Portal in the form of an XML document. The XML document gives information regarding the read access to the facility, data, metadata respectively and other information, i.e. the users DN, lifetime of the user's access rights, normally 2 hours to match the proxy certificates lifetime. This XML document is known as an Authorisation Token.

```
<?xml version="1.0" encoding="UTF-8"?>
<attributeCertificate>
  <acInfo>
    <version>1.0</version>
    <holder>/C=UK/O=eScience/OU=CLRC/L=DL/CN=glen drinkwater</holder>
    <issuer>EMAILADDRESS=ca-operator@grid-support.ac.uk, CN=CA, OU=Authority, O=eScience, C=UK</issuer>
    <issuerName>ACMEMIN</issuerName>
    <issuerSerialNumber>1</issuerSerialNumber>

    <signatureAlgorithm>SHA1withRSA</signatureAlgorithm>
    <validity>
      <notBefore>2004 0 27 13 35
28</notBefore>
      <notAfter>2004 0 27 14 38 10</notAfter>
    </validity>
    <attributes>
      <DPView>t</DPView>
      <wrapperGroup>t</wrapperGroup>
      <dataAccessGroup>t</dataAccessGroup>
    </attributes>
  </acInfo>
  <signature>e/CUhsww6yhnI9+gGbiTB9o0dcsijIE19PmODHrHca+3/qiRJPusww6yhnI9/+</signature>
</attributeCertificate>
```

**Figure 3:** Format of Authorisation Token

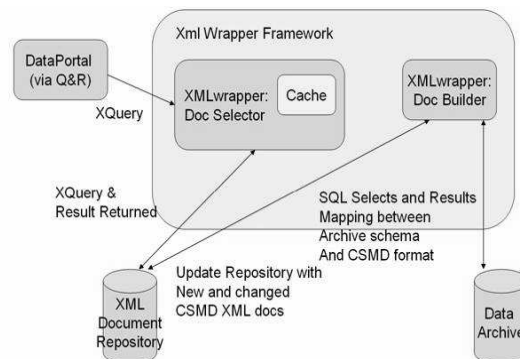
The ACM signs the XML document with the facility's private key and sends the Authorisation Token via web services back to the Data Portal which stores it in a database. When the Data Portal sends a query to the XML Wrapper, it also sends the

Authorisation Token. The XML Wrapper can validate the signature of the Authorisation Token with the facility's public key. Therefore the XML Wrapper can trust the access information regarding the facility given in the Authorisation Token.

### 1.2.2. XML Wrapper

The XML Wrappers are used to convert between the local metadata format of a Data Archive and the CSMD. This allows the integration of metadata repositories which hold information about scientific studies and their associated data in various different formats and present a common interface to them via web services, which allows the Data Portal to seamlessly interact with the different data archives.

It allows the federation of queries to different data archives by placing a wrapper (or adaptor) around the archive which translates the native structure of the data into a form that the Data Portal core modules can understand. It converts the data into a common XML representation and then allows either bringing back the whole data representation or applying XQueries to bring back relevant data. Thus the XML Wrapper allows easier and uniform composable of queries by providing a common interface; once the data format of the Data Portal is known (the CSMD) then XQueries can be written against this to extract relevant metadata.



**Figure 4:** XML Wrapper Architecture

The XML Wrapper has two independent aspects: the building of XML documents into CSMD documents and querying those XML documents.

- 1. XML Wrapper Document Builder:** builds the XML document from the Data Archive format.
- 2. XML Wrapper Document Selector:** processes incoming queries on the built XML documents to retrieve relevant scientific metadata.

The benefits of this type of architecture are that the wrapper allows seamless access to the archive even if the data archive at the facility is down and unavailable. The Document Selector can still process queries and retrieve results using data from the XML Document Repository.

Another benefit of the Wrapper approach is that If Grid Middleware such as OGSA-DAI [7] was used then it would be upon the Data Portal core modules to do the schema discovery and schema mapping from Data Portal queries to that of Data Archive. Thus the Schema mapping would be not the preserve of the Facility but that of the Data Portal this would seriously affect scalability as well as offer an architecture which was still doing an implicit schema conversion or query conversion to extract relevant information from the Data Archive and the rules for such conversion could not be automatic and would have to be done on a per archive basis. Thus schema mapping at or near the facility aids scalability and allows for a much clearer architecture.

## 2. Data Storage

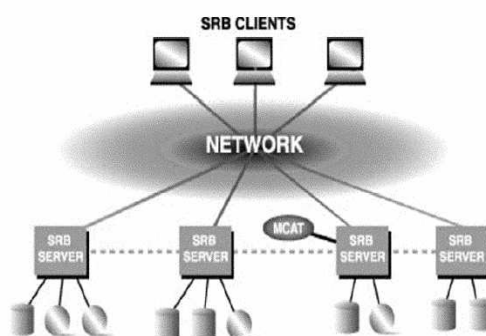
The output of the scientific applications is currently all stored as flat files. These files are often distributed over a number of machines and /or different types of media. There is a data management problem

associated with this scenario, the data is organised physically rather than logically.

In order to access this concern, the Storage Resource Broker has been deployed on the nodes of the projects. The SRB is a tool developed by San Diego Supercomputing, which facilitates data management across a number of distributed heterogeneous file systems.

Essentially the system abstracts the user from the physical location, media and protocols of the underlying storage systems. This allows data to be organised logically into a single virtual file system. In addition, SRB simplifies the sharing of data across a distributed virtual organisation such as the e-Minerals. In particular, it is straightforward to configure the access permissions on individual files to allow other members of the project access. There is also a facility to allow colleagues outside the project access using a ticketing scheme.

The architecture of a SRB domain is shown schematically in Figure 5. Each data storage resource runs a SRB server. The SRB domain needs one master SRB server which is attached to a local database containing the MCAT (Metadata Catalogue) which maps from locations within the virtual file system to physical locations on individual resources.



**Figure 5** SRB Architecture

The Data Portal is able to download or transfer data from a variety of storage systems. Currently the Data Portal can use GSI delegation to access data stored in FTP, HTTP, GASS and GridFTP servers. Recently

the Data Portal was able to access data stored in SRB.

Access to the data via SRB can be through GUIs, web browsers, command line or SRB APIs. This allows the Data Portal to access data sets or files through web services and download them for the user. The access and control is done by SRB. When the data is put into the system, the user specifies which users on the SRB system can have access to the data, similar to a Windows File System, allowing read, write, delete etc access for each file or database blob.

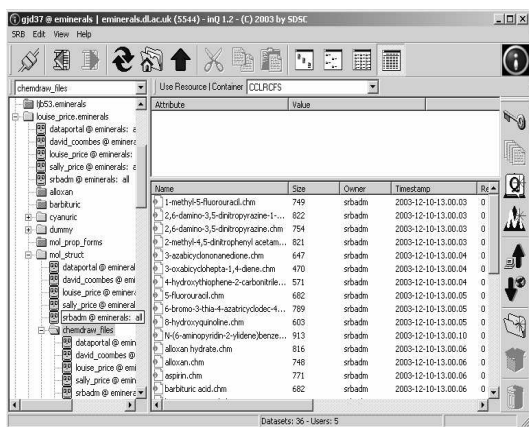


Figure 6: InQ showing a SRB view

The current architecture within CCLRC is that a user would run a remote job on some data and then archive the results for further analysis later. The user would run the job via Globus Gram job and use the SRB command line tool to upload the results into SRB. The user wishing to shard the data to other scientists would annotate this data with the Data Insertion tool and add the link to the data in the metadata.

Scientists can search the metadata via the Data Portal, locate and download the data for further analysis or as input to another job submission.

The uploading and downloading is achieved using a web service interface wrapped around the SRB API Jargon [8]. A user would authenticate themselves to SRB

with a username and password or via GSI delegation with a Globus certificate. Once authenticated the user has access only to their own data and the data that others wish to share with the user.

### 3. Data annotation and publishing

The final stages of the data lifecycle are annotation and publishing. The data files themselves are published by uploading them to the SRB. Although the use of the SRB facilitates sharing and replication of data, the value of the data files would diminish rapidly if they were not annotated with relevant metadata describing the context and method of their generation. In addition, this step is essential if the data files are to be retrievable later in the Data Discovery phase.

The annotation is accomplished by separating the metadata (which is held in a relational database based on the CCLRC Metadata Schema (CSMD)) from the data (which is held in flat files). The metadata contains links to the SRB, which internally maps to the physical file locations. This use of the SRB also maintains referential integrity if the files are moved.

Currently the metadata relating to scientific work is entered using the Data Insertion tool. Within the Insertion tool, individual computation simulations are grouped into studies, which are labelled by various topic hierarchies that allow retrieval via the Data Portal later.

### 4. Future

The Data Portal is will be taking account of new technologies (e.g. Java Portlet API, Java Server Faces, WS-Notification, WS-Resource Framework/Globus Toolkit 4). Further work and research will be undertaken with other projects and the new technologies mentioned above. This will allow new additional web service modules as well as improvements to the current modules.

#### 4.1. Advanced Searching

Using XQuery as the language to search through the XML Documents allows detailed advanced searched to be added to the Data Portal search capabilities. The amount of detail within the CSMD representing the scientific metadata is the only limiting factor in how specific the advanced search can be. This allows the ability for the user to specify their own XQuery or through a GUI to tweak a pre-made advanced search for certain dates, PIs, Institutions etc.

#### 4.2. Shopping Basket sharing

Users will be able to give 'tickets' to other users of the Data Portal allowing restricted access to their Shopping Basket and the information held. This would be useful for PIs giving access to their basket to fellow researchers or post graduates.

#### 4.3. Database Integration

With the increased support for Database structures in the CSMD it is envisaged that the Data Portal will allow access to dataset which are stored in databases. This will not be limited to BLOBs where the database is essentially acting as a file system but also support Named Selects, Queries which have been given a Name by the original data creator such that they return a result and a schema that the result conforms

to (as well as the database system) such that the data can be mined with the appropriate tools. Thus Data Archives which store their data and metadata together or make no distinction will have Wrappers that serve their meta information in XML with enough information to access the data in their databases also in a meaningful way.

Accessing the databases using these named queries maybe an instance where brokering middleware is useful. OGSA-DAI [7] maybe of use in this instance as it would allow authentication of users to get the data based on the Grid certificates which the portal already has when the user logged in.

## 5 References

- [1] Blanshard L, Environment from the Molecular Level e-Science project and its use of CCLRC's Web Services based Data Portal. [http://www.e-science.clrc.ac.uk/documents/staff/kerstin\\_kleese/e-min2003.doc](http://www.e-science.clrc.ac.uk/documents/staff/kerstin_kleese/e-min2003.doc)
- [2] Dr Michael Doherty, SRB in action. [http://www.e-science.clrc.ac.uk/documents/projects/storage\\_resource\\_broker/srb\\_in\\_action.pdf](http://www.e-science.clrc.ac.uk/documents/projects/storage_resource_broker/srb_in_action.pdf)
- [3] W3C, XQuery. <http://www.w3.org/XML/Query>
- [4] Matthews BM, Sufi SA. The CCLRC Scientific Metadata Model - Version 1. The CCLRC Scientific Metadata Model - Version 1 2002. <http://www.dienst.rl.ac.uk/library/2002/tr/dltr-2002001.pdf>
- [5] Ananta Manandhar. Grid Authorisation Framework for the CCLRC DataPortal. [http://www.e-science.clrc.ac.uk/documents/staff/kerstin\\_kleese/Authorisation.pdf](http://www.e-science.clrc.ac.uk/documents/staff/kerstin_kleese/Authorisation.pdf)
- [6] Globus Alliance, Globus and the Grid. <http://www.globus.org/>
- [7] <http://www.ogsadai.org.uk/>, Open Grid Services Architecture Data Access and Integration
- [8] SRB Jargon. <http://www.npaci.edu/DICE/SRB/jargon>