*Title:*   The ASCI Grid:  Initial Deployment

*Author(s):*   Randal Rheinheimer, Los Alamos National Laboratory

Steven L. Humphries,
Hugh P. Bivens,
Judy I. Beiriger,
Sandia National Laboratories, New Mexico

# The ASCI Computational Grid:  Initial Deployment

Randal Rheinheimer
*Los Alamos National Laboratory, P.O. Box 1663, Los Alamos, NM 87545*
*randal@lanl.gov*


Judy I. Beiriger, Hugh P. Bivens, Steven L. Humphreys
*Sandia National Laboratories[1], P. O. Box 5800, Albuquerque, NM 87185-1137*
*{jibeiri, hpbiven, slhumph}@sandia.gov*

## 1.  Overview

The Accelerated Strategic Computing Initiative (ASCI) computational grid consists of a handful of very large SMPs, some of which have dedicated visualization and IO nodes contained within them.  There are also standalone visualization clusters of various architectures and three High Performance Storage Systems (HPSS) within the computing network.  The systems are geographically widely distributed, but are connected by four stripes of OC-12 bandwidth.   The user community is small by grid standards, with only a few analysts accounting for a large percentage of computing cycles and storage bandwidth.

The goal of the Distributed Resource Management (DRM) project in this context is to simplify access to the diverse computing, storage, network, and visualization resources and to provide superior monitoring and job control mechanisms. To this point, our efforts have focused on implementing the grid infrastructure necessary to allow a user to submit, monitor, and control jobs in a secure manner.  The final link in the initial deployment is the user interface itself.  The next six months, as the system is introduced to users more accustomed to individualized scripts than to unified grids, will be a telling time for this particular grid computing environment.

## 2.  Architecture description

One key requirement for the ASCI grid architecture is that local control by disparate batch scheduling systems (LSF, DPCS, PBS, and NQS) be preserved.  The Globus MetaComputing Toolkit[1] was selected to provide the core interface to those batch scheduling systems, for job control and for job monitoring. The DRM project provides grid services to users and to higher level applications. Numerous extensions to the Globus GRAM scripts were made to mesh better with local implementations and to account for the Globus services installation on front end machines rather than on the compute resources themselves, particularly in the area of information gathering for machine monitoring and resource brokering purposes.  Additionally, a major effort was made to ensure that Kerberos authentication mechanisms functioned properly with the entire Toolkit.

The next layer of grid services were developed by DRM and consist of the Production Wizard (a GUI), Workflow Manager, Resource Broker, Information Service, and Monitoring Services.  These will be discussed in more detail in the next section.
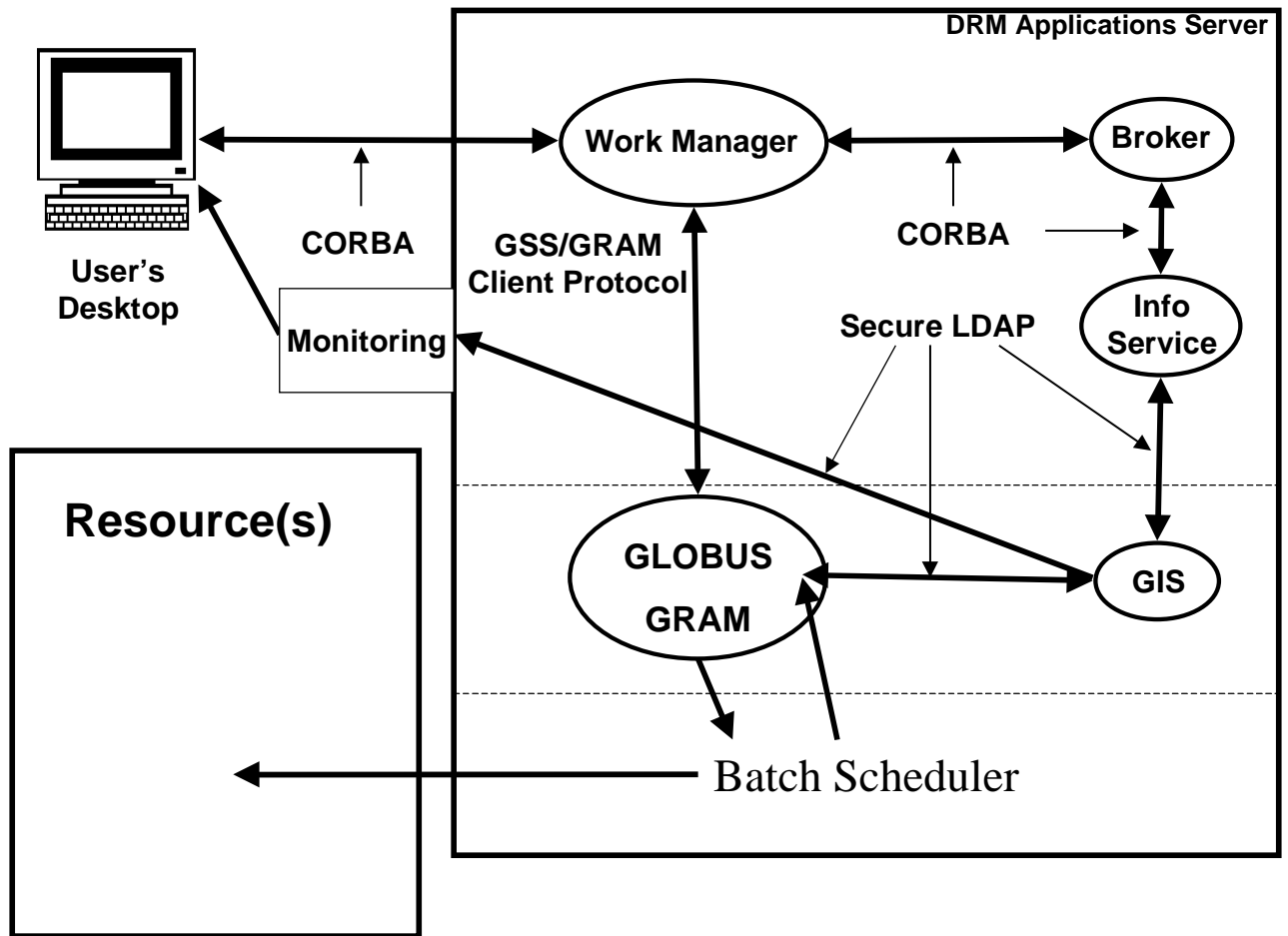
For data movement, DRM uses a custom file transfer client and daemon that run in batch mode.  The daemon, developed by Sandia National Laboratories, New Mexico, is very similar to the HPSS file movement protocol with the mechanisms for communicating with HPSS movers removed, so that a standard parallel file transfer protocol (PFTP) client can communicate both with a storage system and with any machine which has this daemon installed.  Recently, a collaboration has begun which aims to bring the Globus GridFTP project and the HPSS data movement protocol together, and it is anticipated that the future data transfer mechanism will be GridFTP, both for storage and machine-to-machine.

Another element of a robust, secure, grid-computing environment is being addressed by the ASCI DisCom² program in the form of a distributed Credential Agent for a Kerberos/DCE environment.  In an environment in which jobs may run for days or weeks, and workflows for months, it is necessary for the grid infrastructure to provide a mechanism for users to

---

create and designate credentials for use by various processes without circumventing established security regulations and practices.  In the short term, the DRM Workflow Manager will handle the task of refreshing and forwarding user credentials on a limited basis until DRM can become a client of a more flexible Credential Agent.



**Figure 1 :  Architecture, including DRM Services, Globus, and local batch scheduler**

### 3.  Implementation

   The DRM Production Wizard (PW), Workflow Manager, Information Service, and Resource Broker are all CORBA-based.  The Workflow Manager and PW are Java tools; their means of communication is an XML workflow definition vocabulary, the Grid Access Language for HPC Environments (GALE)[2].  A workflow specification, expressed in GALE, is the output of the Production Wizard and the input required by the Workflow Manager.  The intent is for any application to be able to create workflow specifications for the Workflow Manager, but a necessary first deployment step is to have a dedicated tool for that purpose.

   Phase One of the Production Wizard, then, allows an end user to generate a GALE workflow specification by choosing functional computational steps such as job submission, file transfer, or resource queries in graphical mode and filling in pre-defined data fields: source and sink for a data transfer, for instance.  Existing GALE specifications may be imported, altered, and saved, and serial dependencies defined graphically.  Standard error and output may be viewed via the PW, if desired.

   Given a workflow specification, the Workflow Manager sets up serial dependencies, executes resource queries to the Resource Broker, and has as its output a globusrun command with a well-formed resource contact string and Globus Resource Specification Language (RSL) argument, or a file transfer command.  Current job dependencies are strictly

serial; future extensions to the Workflow Manager and to GALE will allow better control for automatic computation restart and parallel dependencies. The Workflow Manager also manages and forwards appropriate user credentials and handles things like reconnecting a new instance of a PW to an existing workflow. Data transfer is handled within the Workflow Manager, as well, with both parallel and serial data transfer capabilities automatically chosen based on data size and source and sink capabilities.

The Resource Broker and Information Service provide information about the "best" resources to use to the Workflow Manager, given a minimal set of user requirements. The information resides in a Lightweight Directory Access Protocol (LDAP)-based Grid Information Server (GIS) which is populated primarily by Globus' original push model, in which job data is culled from local batch schedulers and published to the GIS every thirty seconds, and system configuration data is obtained and published every four hours. The information published is being enhanced to include data like queue run and processor limits, user service ratios, predicted (by the underlying batch scheduler) job start time, and machine (not front end) reliability statistics. The purpose of these enhancements is primarily to enable the Monitoring Services, discussed below, to provide a common user view across our platforms. Some of this data might eventually be used by the Broker in the service of sophisticated brokering algorithms, particularly one based on historical queue wait time for similar jobs, with similarity defined by fairshare service ratio, CPUs and runtime requested, etc. The current brokering algorithms are fairly simple. To provide the best resource for a job submission, for instance, the Broker might use the Information Server to discover where the required software (user defined) is installed. It then returns the computational resource with the least jobs currently queued in the default queue of that resource. Because the demand for brokering services from the ASCI user community is low, however, little attention is currently being given to their enhancement, and development resources have been directed elsewhere.

Monitoring Services, as noted above, are dependent on the GIS; these services consist of a number of HTML and CGI scripts which query the GIS whenever a monitoring page is accessed. The main computational resource information page, which requires proper authentication for access, displays all users' running and pending jobs in graphical format, including the number of processors used. Also displayed are the total number of processors in use for each compute resource, and a graphical node-by-node representation of the usage of Sandia National Laboratories', New Mexico, Cplant resource. Monitoring services have been identified as a useful bridge between current user practices and fully enabled grid computing. Work is ongoing to upgrade these services in a manner that will help to integrate the grid view into the user perspective.

## 4.  Supported Grid Services

Security is one of the first considerations in the deployment of any ASCI software. Kerberos is the primary authentication mechanism, implemented via the Generalized Security Framework (GSF)[3], a security framework developed at Sandia National Laboratories. This mechanism enables secure user-to-user and process-to-process communication throughout the ASCI grid, except for access to data in the GIS. That security is implemented via a Netscape Directory Service (NDS) plug-in that provides a Simple Authentication and Security Layer (SASL) mechanism for authentication using GSSAPI over Kerberos. Access to data is authorized using standard LDAP Access Control Instructions (ACIs).

Within this secure grid infrastructure, ASCI grid services provide job and system status monitoring, job submission and control capabilities with elementary computational resource brokering, and automated data transfer capabilities.

## 5.  Project Status and Future Plans

The ASCI grid infrastructure has been deployed at three national laboratories, encompassing a computational grid of some 19 teraOPS and storage on the order of petabytes. Another 30 teraOPS resource will soon be included. The monitoring features are in production use, and the job submission and control features are due for production use by November 2001.

Near-future security issues include better credential management with more precise delegation and coordination of compute system and Globus authorization mechanisms. Data transfer mechanisms will no doubt improve with the recent establishment of collaborations between HPSS, ASCI data management personnel, and Globus GridFTP. Parallel job control, automated restart mechanisms based on signal reading from job output files, and better error control and reporting are current development tasks within the Workflow Manager. An important step in advancing user acceptance of grid concepts will be enhanced monitoring of local and remote job and systems status; an important key to that is continued enhancement of the Globus GRAM interfaces to publish relevant information to the GIS. More sophisticated brokering systems are available for development, should demand increase. Research and development efforts into co-scheduling are underway, and some early research into network Quality of State (QoS) could be revived, should the need arise.

Development of the Production Wizard will necessarily continue both in reaction to user demand and in anticipation of it, and collaborations with other tool developers to grid-enable their utilities will be an important part of future work.

November 2001 marks a turning point in the way that laboratory analysts and scientists can access the computational resources that are available in the ASCI computing environment. No longer need they be tied to their local resources, or spend time mastering the peculiarities of each new system in order to perform their work. Grid services will provide ubiquitous access to all of the ASCI compute resources. Future grid services will include access to non-compute resources and coordination of multiple resources in support of advanced workflows (conditional, concurrent, and interactive workflow mechanisms). The challenge today is to get the analysts and scientists to use the new grid services and further provide them with tools that will allow them to stop viewing ASCI resources as individual chunks of hardware and start perceiving them as a collection of compute, storage, and visualization services.

## 6. References

[1] Foster, I., and C. Kesselman, Globus: A Metacomputing Infrastructure Toolkit, *International Journal of Supercomputer Applications*, 1997

[2] http://vir.sandia.gov/~hpbiven

[3] Detry, R., Kleban, S., Moore, P., and Berg R., The Generalized Security Framework, Presented at CSCoRE 2000, http://www.ccs.bnl.gov, Brookhaven National Laboratory, NY.