

## **Globus Data Storage Interface (DSI) - Enabling Easy Access to Grid Datasets**

Rajkumar Kettimuthu<sup>1,2</sup>, Michael Link<sup>1,2</sup>, John Bresnahan<sup>1,2,3</sup> and William Allcock<sup>1,2</sup>

<sup>1</sup>Math and Computer Science Division, Argonne National Laboratory, Argonne, IL USA

<sup>2</sup>Computation Institute, The University of Chicago, Chicago, IL USA

<sup>3</sup>Department of Computer Science, The University of Chicago, Chicago, IL USA

In a Grid environment, providing access to data distributed across numerous heterogeneous resources is both an important and challenging task. There are a number of distributed storage systems in use by the scientific and engineering community. These systems have been created in response to specific needs for storing and accessing large datasets. They each focus on a distinct set of requirements and provide distinct services to their clients. For example, some storage systems such as Distributed Parallel Storage System (DPSS) and High Performance Storage System (HPSS) focus on high-performance access to data and utilize parallel data transfer streams and/or striping across multiple servers to improve performance. Other systems like Distributed File System (DFS) focus on supporting high-volume usage and utilize dataset replication and local caching to divide and balance server load. The Storage Resource Broker (SRB) system connects heterogeneous data collections and provides a uniform client interface to these repositories, and also provides metadata for use in identifying and locating data within the storage system. Still other services (HDF5) focus on the structure of the data, and provide services for accessing structured data from a variety of underlying storage systems. Most of these customized storage systems utilize incompatible protocols for accessing data and require the use of their own clients. Applications that require access to data stored in different storage systems must use different methods to retrieve data from each system. It can be challenging to transfer a dataset from one system to another.

GridFTP is a high-performance, secure, reliable data transfer protocol optimized for high-bandwidth wide-area networks. Its been widely used in the Grid environments. To allow GridFTP to be a transfer interface to as many data storage systems as possible, our new GridFTP framework provides a modular pluggable interface to a storage medium. We call that interface as Data Storage Interface (DSI). The DSI presents a modular abstraction layer to a storage system. It consists of several function signatures and a set of semantics. When a new DSI is created, a programmer implements the functions to provide the semantics associated with them. DSIs can be loaded and switched at runtime. When the GridFTP server requires action from the storage system (be it data, meta-data, directory creation, etc) it passes a request to the loaded DSI module. The DSI then services that request and notifies the server when it is finished. An API is provided to the DSI author to assist in implementation. The DSI author is not expected to know the intimate details involved in a GridFTP transfer. Instead this API provides functions for reading and writing data to and from the network. DSIs do exist for the Storage Resource Broker (SRB), the High Performance Storage System (HPSS), and NeST from the Condor team at the University of Wisconsin - Madison. DSIs would confer benefits to both the keepers of large datasets and the users of these datasets. Dataset providers would gain a broader user base, because their data would be available to any client. Dataset users would gain access to a broader range of storage systems and data.