

# Grid Data Management Pilot (GDMP): A Tool for Wide Area Replication

Asad Samar<sup>1</sup>, Heinz Stockinger<sup>2,3</sup>

1) California Institute of Technology, Pasadena, USA

2) CERN, European Organization for Nuclear Research, Geneva, Switzerland

3) Institute for Computer Science and Business Informatics, University of Vienna, Austria

Asamar@cacr.caltech.edu, Heinz.Stockinger@cern.ch

## Abstract

The CMS experiment at CERN, the European Organization for Nuclear Research, is currently setting up a Data Grid which will provide access to several Terabytes of data, distributed and replicated. In the real production environment data will be produced in different countries on both sides of the Atlantic by the end of year 2000. The stringent requirements of data consistency, security and high-speed transfer of huge amounts of data, imposed by the physics community need to be satisfied by an asynchronous replication mechanism. A pilot project called the Grid Data Management Pilot (GDMP) has been initiated which is responsible for asynchronously replicating large object-oriented data stores over the wide-area network to globally distributed sites. We present the design, architecture, functionality and performance results of our first working prototype. Different replication policies and protocols are supported that range from strictly synchronous to rather relaxed asynchronous models in terms of data consistency.

We believe that this first prototype can be regarded as a pioneer step towards a Data Grid and as a prototype for replication management within other Data Grid approaches like DataGrid, GriPhyN and PPDG [1, 2, 3].

**keywords:** Grid computing, data management, replication

## 1 Introduction

Next generation High Energy Physics applications are characterised by large amounts (several Petabytes) of mostly read-only data that are distributed and replicated around the globe. The visualisation and analysis of physics data will bring the physics community a step further in investigating the ultimate constituencies of matter, the big bang and hence our earth. Within the Compact Muon Solenoid (CMS) experiment at CERN, the European Organization for Nuclear Research, we are setting up a Data Grid infrastructure required to fulfil the needs of the physics community. Note that

this real world production system has a smaller scale than the Data Grid project [4, 1] that starts officially in the beginning of year 2001.

Recently, Grids [5] have become very popular in the distributed and parallel computing communities. Whereas in the past we have been speaking about meta and cluster computing, the trend has turned to Grid computing where the network between different nodes spans several countries and even continents. Existing Grid technology can be categorised into two major fields, the traditional computational Grids and data intensive Grids, called Data Grid.

Within the Compact Muon Solenoid (CMS) experiment at CERN, first data production tests are currently being done on the Grid. We expect several Terabytes of data to be produced and stored this year which will be a preparation for the Petabyte-scale data taking in the years 2005 and onwards. Physics data in CMS are stored in object-oriented databases. The DBMS of choice is Objectivity/DB [6]. Existing commercial database management systems provide replication features but they fail to satisfy the stringent requirements of consistency, security and high-speed transfers of huge amounts of data, imposed by the physics community. An asynchronous replication mechanism that supports different levels of consistency, a uniform security policy and an efficient data transfer is necessary.

The need for such a production system as well as the requirement of a prototype for assessing existing Grid technologies and evaluating new strategies in this upcoming field, were the triggers for the Grid Data Management Pilot (GDMP) project [7]. GDMP can be regarded as the first prototype of a Data Grid that is used in a production environment. It is currently being used in the High Energy Physics (HEP) community, but the design is flexible enough to be applied to other data intensive applications as well.

GDMP is responsible for asynchronously replicating large object-oriented data stores over the wide-area network to globally distributed sites. We present the design, architecture, functionality and performance results of our first working prototype. The middle-ware

of choice is the Globus [8] toolkit that provides promising functionality. We present test results which prove the ability of the Globus toolkit to be used as an underlying technology for a world-wide Data Grid. The required data management functionality includes high-speed file transfers, secure access to remote files, selection and synchronisation of replicas and managing the meta information. The whole system is expected to be flexible enough to incorporate site specific policies. Our first prototype is being used by the CMS experiment for the management of simulated physics data over the wide area on both sides of the Atlantic. We will be one of the pioneers to use the Data Grid functionality in a running production system. Although there is recently much effort going on in the Grid community to deal with the management of large amounts of data, the GDMP project can be viewed as an evaluator of different strategies, a test for the capabilities of middle-ware tools and a provider of basic Grid functionalities.

Different replication policies and protocols are supported that range from very stringent synchronous to rather relaxed asynchronous models in terms of data consistency [9, 10].

## 2 Architecture

The GDMP software consists of several modules that closely work together but are easily replaceable. In this section we describe the modules and the software architecture of GDMP. The core modules are Control Communication, Request Manager, Security, Database Manager and the Data Mover. An application which is visible as a command-line tool uses one or several of these modules.

### 2.1 Control Communication Module

This module takes care of the control communication between the clients and the servers, and uses the Globus IO library as the middle-ware and builds high level functionality on top. It takes care of the intricacies related to socket communication over the wide area between nodes with heterogeneous architectures. The functionality includes starting and stopping the server, connecting and disconnecting the client to and from the server and sending and receiving messages at both the client and server ends. This module provides services to the modules in the layers above.

### 2.2 Data Mover Module

The main purpose of GDMP is to move files over the wide area in an automatic, efficient and fault tolerant way. This is the module which actually transfers files physically from one location to another one. It uses the NC-FTP client libraries to open a connection

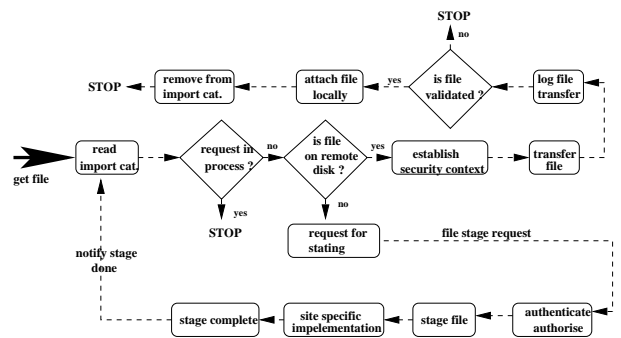


Figure 1. The actual file transfer

with the WU-FTP server on the server end. Since we use the GSI-NCFTP client and the GSI-WUFTP server, we have the same security mechanism (Globus Security Infrastructure) for both the Control Communication and the Data Mover. The functionality includes client authentication and authorisation before the transfer starts, transferring files to and from the server, validating the transferred files using the file size and checksum information, resuming the file transfer from the latest checkpoint after a network failure, using a progress meter to output the progress during a file transfer and finally logging all file transfers that have been completed. Our intention is to replace the use of NC-FTP with the upcoming Globus FTP libraries. This will enhance the performance of the Data Mover Module because of new features like partial file transfers, parallel file streaming and better fault recovery behaviours expected in the Globus FTP. Figure 1 gives the complete picture of the file transfer.

### 2.3 Security Module

Security is one of the main concerns in a Grid. Allowing people from outside one's domain to use the resources is a big issue for most organisations. Sensitivity of data and unauthorised use of network bandwidth to transfer huge files are the main security issues we are dealing with. This is done by the Security Module. It is based on the Globus Security Infrastructure which is an implementation of the Generic Security Service (GSS) API. It uses the Public Key Infrastructure as the underlying mechanism. The Security Module provides methods to acquire credentials, initiate context establishment on the client side and accept context establishment requests on the server side (context is established to authenticate the client), encrypting and decrypting messages and client authorisation. The server authenticates and authorises any client before servicing its request, hence, the software protects a site from any un-wanted file transfers and blocks any un-authorised requests.

## 2.4 Request Manager Module

Every client-server system has to have a way to generate requests on the client side and interpret these requests on the server side. The Request Manager Module does exactly that. It contains several request generator/handler pairs and more can be added for customised use. This module is based on the Globus DC library which provides methods to convert data between different formats supported by variable machine architectures. The requests are generated on the client side by filling a buffer with the appropriate function handler, which is to be called on the server end, and any arguments required by that function. On the server side this buffer is unfolded and the data types are converted according to the local architecture. Finally, the respective function is called with the given arguments. The Request Manager Module basically mimics a limited Remote Procedure Call (RPC) functionality with the advantage of being light-weight and extensible.

## 2.5 Database Manager Module

This is the module which interacts with the actual Database Management System (DBMS). The module relies on the APIs provided by the particular data storage system being used. In our case the DBMS of choice is the Objectivity/DB, hence we use Objectivity's APIs to implement this module. To use GDMP in a different system, this is the only module which has to be swapped by the one which can interact with the specific DBMS. The functionality includes retrieving the database catalogue, containing information about the files currently present in the database, and attaching files to the DBMS once they arrive on the client side and are validated.

## 2.6 The GDMP Applications

GDMP includes some applications which are the customers of the services provided by the above outlined modules. These applications include various clients and one server.

### 2.6.1 The GDMP server

The GDMP server is a daemon constantly running on sites which produce data or want to export their data to other sites. We expect a large number of clients connecting to the GDMP server from all over the globe, transferring huge files which might take days to complete the transfer. Under such conditions the server is required to be very robust, extremely fault tolerant and able to cope with multiple clients simultaneously.

The server itself uses the communication module for receiving requests from application clients. Since thread creation is rather time consuming, the server

uses a thread pool with a certain amount of threads which can be adapted. Each time a client is connecting, a thread is allocated to a single client. This is again a performance aspect, since the client and the server can communicate over one connection channel as long as the client is connected to the server. When the client disconnects, i.e. the application program terminates, the socket connection is closed and the thread is returned to the thread-pool.

For each client one thread is used. Thus, the number of threads in the server corresponds to the number of concurrent connections to the server. If no free thread is available, the client's request is put into a waiting queue. As soon as one thread terminates, the first request in the queue is served. The number of elements in the queue can be set on compilation time of the server.

### 2.6.2 Client Application Programs

An application program sends a request through the Request Manager to any other module that serves the user request. For instance, a user starts the tool `gdmprreplicate_file_get` which internally formulates a request that is generated by the Request Manager. This module uses the Globus Data Conversion library for transforming a request into a byte string for internal socket communication in the communication module. The communication module then communicates the required information to the server, which uses the Data Mover module to transfer the file from site A to site B. Once the file is transferred, it has to be integrated into the local Objectivity federation. Since this is a very database specific function, the DB-Manager takes care of vendor specific features for integrating files. Thus, the DB-Manager is very specific to the data storage requirements of the end-user and currently tightly coupled to Objectivity. However, the DB-Manager is extensible to any database management system or other storage systems.

Since all the access to data has to be done in a secure environment, a client has to be authorised and authenticated before he can request a service from the server. A client has to have a Globus proxy running in order to start the authentication process which is handled by the security module that is based on GSI, the Global Security Infrastructure. We use the *single login* procedure which is available through Globus, i.e. once a client has successfully got the proxy on one machine, he can send requests to any server without any further password entering (provided the local client is authorised to access the server).

Figure 2 shows the current architecture and all the modules of the software. On top of the architecture we have the Globus application, which can be multi-threaded. All the software modules are written in C++ and run on Solaris 2.6, 7 and Linux RedHat 6.1.

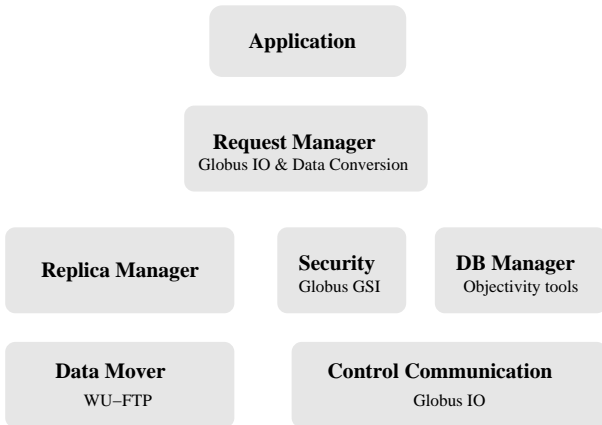


Figure 2. The GDMP Architecture.

The clients provide command-line tools to the users, offering different functionalities. These basically act as the user interface to GDMP. Each client is customised to perform a specific task by communicating with the remote server(s).

### 3 Replication Policies and the Data Model

The GDMP software tool performs automatic, asynchronous replication of Objectivity database files in a Data Grid environment. Currently, the software is restricted to replicate Objectivity/DB files only, but future extensions will allow it to replicate files of any data type. The restriction to replicate only Objectivity files is due to the use of native Objectivity federation catalogue to handle files in GDMP. Once the Objectivity file catalogue is replaced by the announced Globus Replica Catalogue [8], a more flexible replication model can be supported.

#### 3.1 Interfacing with the Data Production Software

In principle, a site, where Objectivity files are produced, has to trigger the GDMP software which notifies all the "subscriber" sites in the Grid about the new files. It is the responsibility of the data production software to trigger GDMP only after the Objectivity files have been completely written and are ready to be transported. In detail, the data production software at the local site uses the command-line tools provided by the GDMP software.

The "subscriber" (destination) sites receive a list of all the new files available at the source site and can determine themselves when to start the actual data transfer. The data transfer is done with a WU-FTP server and an NC-FTP client. Since the usage of different machines in a Grid is a big security issue, a user

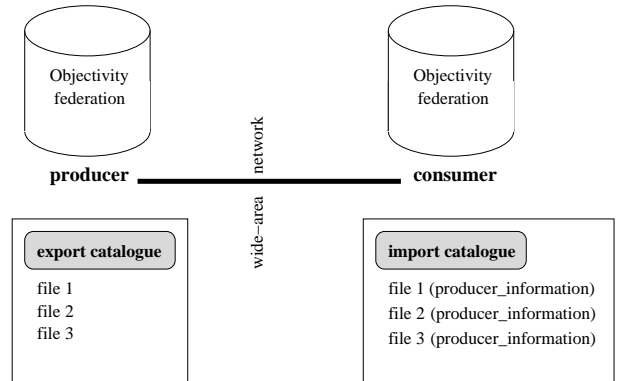


Figure 3. The role of the export and import catalogues

has to be authenticated and authorised before contacting any remote site. The security is based on the GSI security toolkit available from Globus.

#### 3.2 File Catalogues and a Subscription Model

We now elaborate on how GDMP manages the transfer of files and their integration into the Objectivity catalogue. We illustrate the data flow by a producer-consumer example. The producer is the site where one or several Objectivity files are originally written, and the consumer is the site that wants to replicate these files locally. Once the producer has finished writing a set of files (or just a single file), it publishes this information by creating and sending an *export catalogue* to all the "subscribed" consumers. The export catalogue has a listing of all the newly generated files and their related information. The consumer creates an *import catalogue* where it lists all the files that are published by the producer and have not yet been transferred to the consumer site. Figure 3 illustrates this model graphically.

The currently implemented replication policy is an asynchronous replication mechanism with a subscription model. A producer can choose at what time new files are written into the export catalogue and thus made publicly available for consumers. Hence, the producer can decide the degree of data consistency by delaying the publication of new files. In the example above we only have one consumer for demonstration purpose. In reality, the number of consumers can be large and depends on the number of sites in the Grid. The subscription model means that the subscribed consumers get informed immediately when new files are published in the export catalogue. Each consumer that wants to be notified about changes in the producer's export catalogue, subscribes to a producer. Depending on the degree of interest in a producer's data, consumers might want to subscribe to

only some of the producers in the Grid.

Since the data exchanged has to be done in a controlled and secure way, a consumer first has to be “registered Grid user” at the producer site, i.e. the user has to be added to the grid-mapfiles of the producer site. These files contain all the users who are allowed to talk to GDMP servers running on a site. Thus a producer has total control over who subscribes to and transfers files from its site. Once this is done, a consumer is allowed to subscribe to the producer site. The producer then adds the new consumer and its related information in a file called *host-list*.

Once a producer has decided to publish the new entries made in the export catalogue (the tool `gdmp_publish_catalogue` is used), the producer sends the whole export catalogue to all the subscribed consumers. At the consumer site, the GDMP software creates the corresponding entries in the local import catalogue of the consumer. The export catalogue can be regarded as an intermediate buffer that contains a list of newly created files. In more detail, the current Objectivity federation<sup>1</sup> catalogue and the old catalogue (the one which was available since the latest publication of the catalogue) are compared and the files which are new are inserted into the export catalogue. The export catalogue only contains newly created files and does not propagate information about deleting of files. This is regarded as a HEP-specific feature where files are not deleted but versions of old files are created.

The *export catalogue* contains the necessary information about the new files (host-name, port, file-name with full directory path). The consumer can then decide when to start the file transfer from the producer to the consumer site with the tool `gdmp_replicate_file_get`. The tool reads the import catalogue and starts FTP sessions to get the necessary file. Once a file has safely arrived and is integrated into the local Objectivity federation, the file entry is deleted from the import catalogue. Section 4 describes what happens in case of broken connections and network failures.

Every transfer of a file, either successfully or not successfully done, is logged in a file called `replicate.log`. Figure 4 shows the control flow for a catalogue update.

### 3.3 Partial Replication: Filtering Files

GDMP allows a partial-replication model where not all the available files in a federation are replicated. In other words, one or several filter criteria can be applied to the import and/or export catalogue in order to sieve out certain files. This allows for a partial replication model where the producer as well as the consumer can limit the amount of files to be replicated.

<sup>1</sup>A federation is the highest granularity of storing data in Objectivity. A federation consists of several database files which are managed by a federation catalogue.

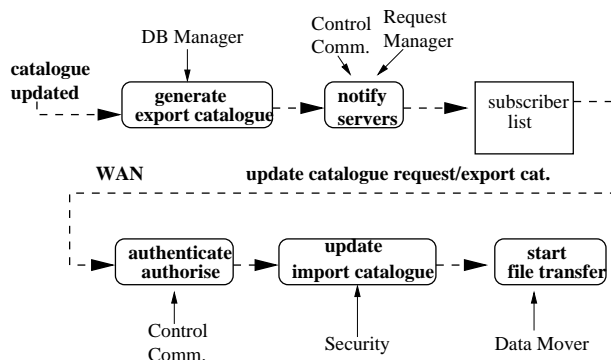


Figure 4. The control flow for a catalogue update.

## 4 Fault tolerance and failure recovery

In the current version of GDMP, each site is itself responsible for getting the latest information from any other site. This is also expressed by the subscription system, where a site has to explicitly subscribe to another site in order to get the file catalogue. Furthermore, when a site recognises a local error which has caused the broken connection, this site has to request from the peer site the required information. A site which publishes information to a subscribed site does not re-send information nor logs that a site could not receive the information. A site can retry to send the message again to the destination site within a particular time frame which can be set by a timeout parameter. If the re-sending fails again, the sending sites stops trying to contact the sites and hands over the responsibility to the destination site to recover from the broken connection.

To sum up, the policy is the following. Each site has to be aware of its state (connection okay or broken). Then it has to search for the origin of the broken connection. If it detects that the error is on the own site, it has to recover otherwise the peer site is responsible for failure recovery.

A site can be unavailable for several hours or even days. Meanwhile several producers can have created and published files, and the entries in the export catalogues may already have been overwritten. Recall that a producer deletes the entries from the export catalogue once the catalogue has been successfully published to at least one consumer. A consumer can recover from the site failure by issuing the command `gdmp_get_catalogue`. Once a producer site has published its catalogue, the catalogue is available to be transferred to any consumer’s site. GDMP at the consumer site then compares the consumer and producer catalogue and creates the necessary information in the consumer’s import catalogue. Since the implementation of WU-FTP has a “resume transfer” feature, not the entire file has to be re-transferred but only the part

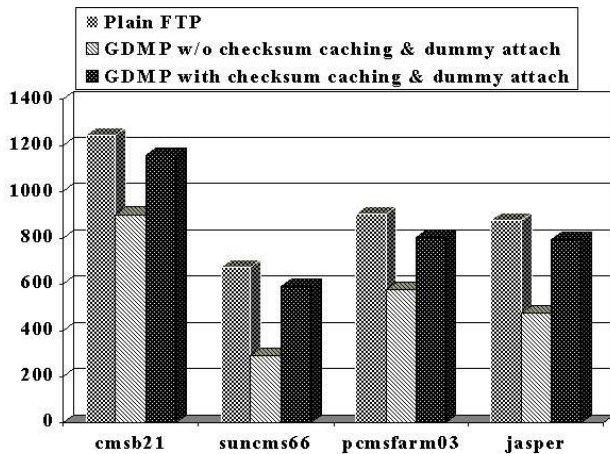


Figure 5. Experimental Results.

of the file that is still missing since the last check point in the file. This allows for an optimal utilisation of the bandwidth in case of network errors.

## 5 Experimental Results

We have used 4 different test machines at both sides of the Atlantic at CERN, Pisa (Italy) and Caltech for our experiments as well for the production use. suncms66 and cmsb21 have installed Solaris 2.6 whereas the other machines run Linux RedHat 6.1. Our performance tests do not emphasise the network performance but the performance relative to the NC-FTP implementation on Linux and Solaris, see Figure 5.

We did a suite of different tests in order to measure the raw aggregate throughput over the network link. We measured the latency for the whole replication process of a number of files. The latency includes the calculation of the CRC check sum, which increases with the file size, attaching a DB file to the federation and deleting the entry from the import catalogue. The graph proves that with checksum caching and smart dummy attach methods, we get almost the same aggregate throughput with GDMP as with plain NC-FTP. Hence, GDMP provides much more automated and reliable replication with tolerable overhead.

## 6 Conclusion and Future Work

We have been developing a data replication tool that allows for secure and fast data transfers over the wide-area network in a Data Grid environment. With our production ready software we have proven that Globus can be used as a middle-ware toolkit in a Data Grid. This has been a pioneer step in the direction of a Data Grid and to the best of our knowledge first software approach where a wide-area replication tool based on

Globus is used in a production system.

The current architecture is restricted to Objectivity files but the system is kept flexible and extensible to include the announced Globus Replication Manager. Once this is integrated, we can provide a replication mechanism for any kind of files and the replica catalogue is managed by the Globus toolkit.

## References

- [1] DataGrid Project: <http://www.cern.ch/grid/>
- [2] GriPhyN Project, <http://griphyn.org>
- [3] Particle Physics Data Grid, <http://www.ppdg.net>
- [4] Wolfgang Hoschek, Javier Jaen-Martinez, Asad Samar, Heinz Stockinger, Kurt Stockinger, Data Management in an International Data Grid Project, *1st IEEE, ACM International Workshop on Grid Computing (Grid'2000)*, Bangalore, India, Dec. 2000.
- [5] Ian Foster and Carl Kesselman (editors), *The Grid: Blueprint for a New Computing Infrastructure*, Morgan Kaufmann Publishers, USA, 1999.
- [6] Objectivity Inc., <http://www.objectivity.com>
- [7] Mehnaz Hafeez, Asad Samar, Heinz Stockinger, A DataGrid Prototype for Distributed Data Production in CMS, *VII International Workshop on Advanced Computing and Analysis Techniques in Physics Research (ACAT2000)*, October 2000.
- [8] The Globus Project: <http://www.globus.org>
- [9] Yuri Breitbart and Henry Korth, Replication and Consistency: Being Lazy Helps Sometimes, *Proc. 16 ACM Sigact/Sigmod Symposium on the Principles of Database Systems*, Tucson, AZ, 1997.
- [10] Jim Gray, Pat Holland, Patrick O'Neil, Dennis Shasha, The Dangers of Replication and a Solution, *SIGMOD*, 1996.