



# REPORT ON REQUIREMENTS FOR INTEGRATION AND INTEROPERABILITY WITH DATAGRID

## WP5.2

---

Document Filename:	<b>CG5.2-D5.2.4-v1.0-CYF030-DataGridCollaboration.doc</b>
Work package:	<b>WP5.2</b>
Partner(s):	<b>CYFRONET</b>
Lead Partner:	<b>CYFRONET</b>
Config ID:	<b>CG5.2-D5.2.4-v1.0-CYF030-DataGridCollaboration</b>
Document classification:	<b>CONFIDENTIAL</b>

---

Abstract: This document details the areas of cooperation between CrossGrid and DataGrid, including sharing of software modules, operational procedures, testbed collaboration and other contributions.



### Delivery Slip

	Name	Partner	Date	Signature
<b>From</b>				
<b>Verified by</b>				
<b>Approved by</b>				

### Document Log

Version	Date	Summary of changes	Author
0.5	22/1/2003	First draft	Piotr Nowakowski
1.0	26/1/2003	Refinements and updates	Piotr Nowakowski, Robert Pająk

---

## CONTENTS

<b>1. INTRODUCTION .....</b>	<b>4</b>
<b>2. ABBREVIATIONS.....</b>	<b>5</b>
<b>3. REFERENCES .....</b>	<b>6</b>
<b>4. THE DATAGRID COLLABORATION FRAMEWORK.....</b>	<b>7</b>
<b>5. DATAGRID SOFTWARE MODULES ADOPTED BY CROSSGRID .....</b>	<b>8</b>
5.1. EDG REPLICA MANAGER AND REPLICA CATALOG (EDG WORK PACKAGE 2).....	8
5.2. EDG RESOURCE BROKER AND JOB SUBMISSION SERVICE (EDG WORK PACKAGE 1) .....	9
<b>6. OTHER DATAGRID CONTRIBUTIONS TO CROSSGRID.....</b>	<b>10</b>
6.1. STANDARD OPERATING PROCEDURES.....	10
6.2. PARTICIPATION IN THE GLUE FORUM.....	10
<b>7. CROSSGRID CONTRIBUTIONS TO EDG DEVELOPMENT .....</b>	<b>11</b>
7.1. STORAGE OPTIMIZATION.....	11
7.2. PARTICIPATION IN THE GRIDSTART FORUM .....	11
<b>8. TESTBED COLLABORATION .....</b>	<b>12</b>
8.1. CROSSGRID TESTBED STATUS .....	12
8.2. CROSSGRID TESTBED DEVELOPMENT AND ITS LINKS TO EDG .....	12
8.3. TESTS AND DEMONSTRATIONS.....	13
<b>9. FUTURE ISSUES AND CONCLUSION .....</b>	<b>15</b>

---

## 1. INTRODUCTION

The aim of CrossGrid is to extend Grid functionality to a new area of applications (mainly interactive tasks), basing on current Grid middleware and infrastructure; most notably those developed within the European DataGrid. DataGrid (hereafter called EDG) predates CrossGrid (CG) by over a year, so it is obviously at a more advanced stage than CrossGrid and has already yielded valuable results, which can be incorporated into other research projects. Therefore it should come as no surprise that the co-operation between EDG and CG is heavily biased towards adoption of solutions and mechanisms developed within EDG by CG programmers, Technical Architecture Team and project management. However, there are selected aspects of interaction between the two projects where EDG can and does benefit from close contacts with CG (mainly the optimization of data access and testbed sharing).

This document is an attempt to summarize the current state of cooperation between EDG and CG by pointing out which software modules and mechanisms are shared by both projects and in what way the software developed as part of EDG affects the architecture and functionality of CG.

---

## 2. ABBREVIATIONS

CA	Certification Authority
CE	Computing Element
CERN	Centre Européenne pour la Recherche Nucléaire
CG	CrossGrid
EDG	European DataGrid
GBJ	Grid Batch Job
GIJ	Grid Interactive Job
GLUE	Grid Laboratory Unified Environment
IST	Information Society Technologies Programme
JDL	Job Description Language
LDAP	Lightweight Directory Access Protocol
LHC	Large Hadron Collider
OGSA	Open Grid Services Architecture
OGSI	Open Grid Services Infrastructure
RAS	Roaming Access Server
RB	Resource Broker
RC	Replica Catalog
RM	Replica Manager
RPM	RPM Package Manager
SE	Storage Element
VO	Virtual Organization
VOMS	Virtual Organization Membership System

---

### 3. REFERENCES

- [FRAME] The CrossGrid - DataGrid Collaboration Framework; [http://www.eu-crossgrid.org/cooperation\\_dg.htm](http://www.eu-crossgrid.org/cooperation_dg.htm)
- [INIT] CrossGrid D4.3: Initial CrossGrid testbed (Confidential)
- [PORT] CrossGrid WP3.1 (Roaming Access and Portals) Design Document (Confidential)
- [SOP] CrossGrid Standard Operating Procedures; <http://kinga.cyf-kr.edu.pl/~tat/docs/Deliverables/D5.2.3/CG5.2-D5.2.3-v1.0-CYF020-StandardOperatingProcedures.doc>
- [VALID] CrossGrid D4.2: validation of testbed architecture; <http://www.eu-crossgrid.org/Deliverables/M6pdf/CG4.4-D4.2-v1.1-LIP011-ValidationOfTestbedArchitecture.pdf>

#### 4. THE DATAGRID COLLABORATION FRAMEWORK

The framework of collaboration between CrossGrid and DataGrid is governed by a special contract signed upon the creation of the CrossGrid Project [FRAME]. In it, the management teams of both projects agree to support one another through access to software and computing resources. The main areas of collaboration include:

- reuse of middleware developed within both projects,
- joint design and implementation of new, interactive-oriented Grid services,
- establishing joint teams for evaluation of OGSA, participation in the GLUE effort and evaluation of software engineering tools,
- elaboration of common programming and naming conventions
- elaboration of joint proposals for the GGF
- participation of the chairmen of the EDG Architecture Board and Technical Forum and the CG Architecture Team in common meetings.

The issues presented above will be discussed in the following sections of this document. A special section will be devoted to testbed sharing and common initiatives regarding EDG and CG infrastructures. Other points will be addressed as necessary.

## 5. DATAGRID SOFTWARE MODULES ADOPTED BY CROSSGRID

The authors have established which DataGrid components are being relied upon by CrossGrid software development tasks, based on the CrossGrid Design Documents.

Software development in CrossGrid is organized in three separate Work Packages: WP1 (CrossGrid application development), WP2 (Application development support) and WP3 (new Grid services and tools). Relatively speaking, the highest reliance on DataGrid components can be found in WP3, which deals with Grid middleware and the corresponding software modules (i.e. job scheduling, portal access, job definition language, performance metrics and benchmarks).

Figure 5.1 shows the general component diagram of the CrossGrid architecture. Elements inherited from DataGrid are marked in purple. The following section lists DataGrid software components and explains where they fit in the CrossGrid architecture.

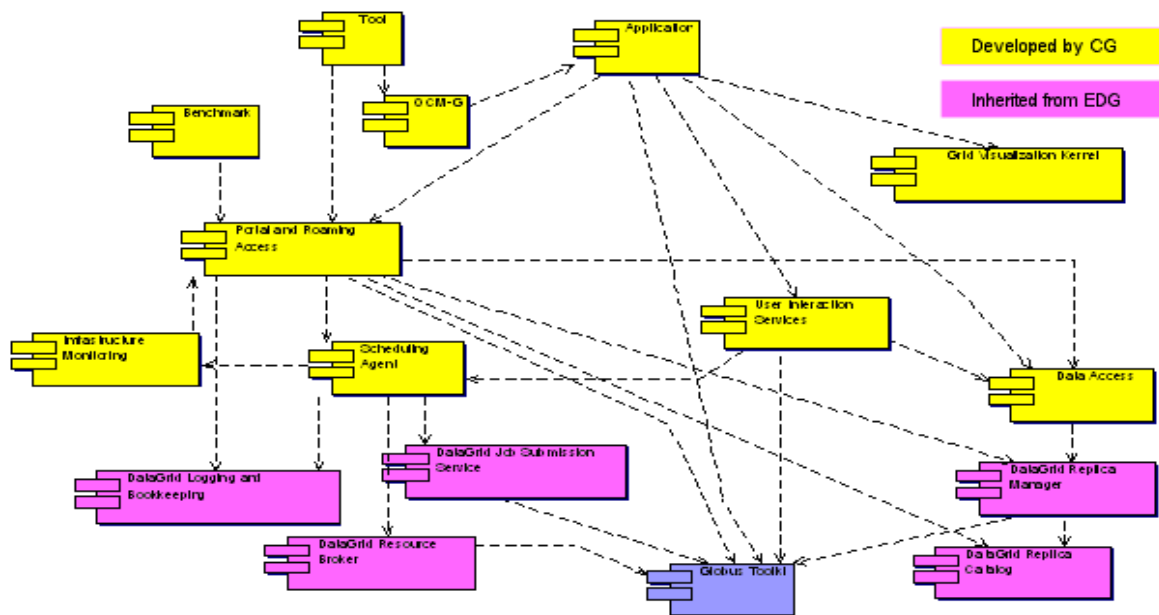


Figure 5-1: CrossGrid Architecture as implemented in the first set of prototypes

### 5.1. EDG REPLICAMANAGER AND REPLICACATALOG (EDG WORK PACKAGE 2)

The DataGrid Replicamanager and Replicacatalog will be used by CrossGrid Task 3.1 (Portals and Roaming Access) for managing and accessing files. The schema of interaction is as follows:

- Upon receiving a file request, the Roaming Access Server (RAS) queries the Replicacatalog as to the location of the best replica.
- Using this information RAS issues a download command from a storage element using the replicamanager software.

---

Likewise, when a file is being uploaded, the RAS is asked to provide a locale for storing the data, which is then copied via GridFTP to the appropriate Storage Element (SE), managed by a Replica Catalog. Finally, the corresponding entry in the LDAP database is updated.

## **5.2. EDG RESOURCE BROKER AND JOB SUBMISSION SERVICE (EDG WORK PACKAGE 1)**

The Resource Broker (RB) is one of the key elements of EDG and as such it has been adapted by CG for early use, prior to being extended by providing additional functionality (related to support for interactive applications and the “person in a loop” concept native to CrossGrid architecture). The EDG RB will initially be incorporated into CrossGrid Task 3.1 (Portals and Roaming Access), where it will underlie the job submission service provided by CG portals.

The role of the Resource Broker is to match requests with the appropriate resources on the Grid and to submit jobs for processing on these resources. In DataGrid, this is done via a dedicated Job Description Language (JDL). JDL files constitute definitions of individual jobs and the resources needed for their completion. CrossGrid Task 3.1 will inherit the JDL concept and the associated job submission service from DataGrid, then extend it to cover new aspects of Grid processing (interactivity and control).

The current version of JDL and the job submission service as implemented by DataGrid enables the user to submit Grid Batch Jobs (GBJ). These are accompanied by properly written JDL files and submitted to the Resource Broker for processing. The user may monitor the current status of the job, but there is no control over the execution process. Once completed, the job reports its results back to the RB, which then makes them visible to the user.

In the near future, DataGrid aims to extend the functionality of its job submission service and RB to cover a new type of Grid jobs - namely, Grid Interactive Jobs (GIJs). A GIJ is a job, which has been declared as interactive through the use of an appropriate flag in the JDL file. During the whole session a bidirectional channel is open between the user and the job running on a remote machine. This channel is used for interactive communication: the user submits input to the job and the job sends back output and/or error messages. Details on how the GIJ mechanism can be adopted by CrossGrid are to be found in [PORT].

## 6. OTHER DATAGRID CONTRIBUTIONS TO CROSSGRID

Besides software modules, DataGrid also provides CrossGrid with organizational support. Testbed collaboration will be covered separately (see Section 8); this section lists contributions that are not directly related to sharing of software modules or computing facilities.

### 6.1. STANDARD OPERATING PROCEDURES

The CrossGrid Standard Operating Procedures document is based on an analogous document developed and implemented in DataGrid. The main areas of convergence include:

- Incremental release preparation procedures
- Software naming conventions
- Automated documentation tools
- Issue reporting and handling
- Document templates

While borrowing some concepts and ideas from EDG, CrossGrid extends and adapts the Standard Operating Procedures to comply with local requirements and conditions. See [SOP] for further details.

### 6.2. PARTICIPATION IN THE GLUE FORUM

Like Gridstart, GLUE is a high-level initiative aimed at fostering cooperation between various Grid projects. However, unlike Gridstart, which focuses on dissemination and awareness initiatives, GLUE is strictly technical. Its aim is to create a uniform, ready-to-use Grid suite of middleware tools which can be adopted for existing and new Grid undertaking. This involves providing recommendations to the Globus development team (in fact Globus developers are actively participating in the Glue effort). Both EDG and CG are active members of the GLUE research forum and mailing lists and have participated in developing common GLUE schemas for Computing Elements (CEs) and Storage Elements (SEs), based on existing EDG/CG architectures.

## 7. CROSSGRID CONTRIBUTIONS TO EDG DEVELOPMENT

Being a younger project, CrossGrid offers little in the way of direct software contributions to DataGrid. Nevertheless, EDG does benefit from contacts with CG.

### 7.1. STORAGE OPTIMIZATION

There is active co-operation within CrossGrid Task 3.4, dealing with optimization of access to secondary and tertiary storage. This task is currently working on extending the DataGrid Replica Manager with functions, that evaluate access times to given storage systems/data sources. DataGrid will then be able to integrate this functionality in an improved version of the Replica Manager.

### 7.2. PARTICIPATION IN THE GRIDSTART FORUM

Both projects actively participate in Gridstart events, as can be demonstrated by the Gridstart presentation at the IST2002 event in Copenhagen (November 2002), where both projects shared a common stand and several CrossGrid testbed sites took part in the WorldGrid demo conducted by DataGrid personnel (see Section 8 for details on testbed collaboration).

---

## 8. TESTBED COLLABORATION

The initial testbed is simultaneously a production facility to test the first CrossGrid application prototypes, a test facility to gain experience with the existing EDG and Globus middleware and the starting point for the deployment of the first CrossGrid testbed that will include middleware developed by the project. Hence the initial testbed has been deployed and maintained in the context of the tasks 4.1 (International Testbed Organization) and 4.4 (Verification and Quality Control).

The middleware for the initial CrossGrid testbed prototype is based on the Globus and DataGrid distributions. This ensures compatibility with other sites running Globus and EDG middleware thus extending the geographic coverage of the Grid in Europe and at the same time providing a basis for the development and test of CrossGrid middleware and applications.

### 8.1. CROSSGRID TESTBED STATUS

The first CrossGrid testbed activities started at the project very beginning, some partners were already involved in grid testbed deployment and research activities, namely through the involvement in the DataGrid testbed and in the CERN LHC computing. These partners brought valuable know how and expertise to the CrossGrid project providing a seed for the initial CrossGrid testbed.

Since the CrossGrid middleware is still being developed it was decided that the initial testbed middleware would have to be based entirely in DataGrid (EDG) and Globus middleware distributions. Delays in the release of a stable EDG testbed 1 middleware also delayed the deployment of many CrossGrid sites. The efforts to establish an integrated CrossGrid testbed started with the release of EDG 1.2, however several problems were found in the first installed sites. The EDG release 1.2.2 with improved stability has allowed more sites to join in spite of some serious middleware limitations. Currently EDG 1.2.2 and 1.2.3 are deployed and version 1.4.3 is being tested at several sites. It is expected that it will overcome many of the major limitations of the previous versions and will allow the interconnection of both CrossGrid and DataGrid testbeds [INIT].

### 8.2. CROSSGRID TESTBED DEVELOPMENT AND ITS LINKS TO EDG

A CG VO server has been installed at LIP and is being used to support CrossGrid activities. Close collaboration with the DataGrid certification authorities task force has been established aiming to coordinate CA activities and mutual cross acceptance of CA's used by the two projects. This was a major step towards the interoperation of both testbeds. DataGrid now accepts all CrossGrid CA's with the exception of the Cyprus CA that has been recently deployed.

The deployed CA's mentioned above are issuing X.509 authentication certificates to the CrossGrid users and systems involved in the current testbed activities. One of the most frequent problems found in the setup of the CrossGrid sites until recently was the installation of the certificates and CRL distribution points for the new CrossGrid CA's. The problem appeared since the new CA's were not being distributed nor installed within the standard EDG middleware. Therefore new RPMs had to be built to solve this problem until DataGrid could recognize the CA's. Currently all CA certificates and distribution URLs are installed properly. The CA's are generally behaving correctly regarding both the certificate issuance and the CRL issuance, which is essential for stable testbed operation. The

---

CrossGrid CVS repository has been established at Karlsruhe and will play an important role as the main distribution point for all the CrossGrid RPMs including the CA RPMs [INIT, VALID].

Two statistics modules have been developed in Perl and Shell to obtain information about the usage of CE and RB systems and publish it in a web server. The RB statistics module collects information from the Logging and Bookkeeping database while the CE statistics module collects information from the Globus gatekeeper log files. Information is collected twice a day. The published information contains counters with the total values and the variations since the last data retrieval. Collaboration with DataGrid on this subject has started with the manifestation of interest by DataGrid of using the CrossGrid statistics modules.

The Mapcenter tool has been enhanced with the capacity of probing of new services, and the addition of links to web pages containing CE and RB statistics developed at LIP. Collaboration with DataGrid on the improvement of Mapcenter has started with the implementation of some corrections and changes.

The EDG VO management tools have been modified to support the CrossGrid testbed needs and to correct several problems and insufficiencies. Some new scripts have been written to ease and automate the VO management process.

VOMS is the new EDG software to replace the currently used LDAP-based VO servers. VOMS is a more sophisticated system that uses a relational database to store the users' rights and roles within a Virtual Organization. VOMS is being tested at LIP.

The LCFG installation profiles provided by DataGrid have been reviewed, cleaned and enhanced. These changes have been reported back so that they can be used in a future EDG release. Several CrossGrid partners have submitted LCFG bug reports and enhancements to DataGrid. An extensive manual about site deployment with LCFG has been written by Demokritos covering all the site deployment steps. The new EDG distribution of LCFGng has been deployed and tested in a small cluster of three machines at LIP using RedHat 7.2. A report of the problems found has been submitted to DataGrid.

DataGrid has adopt Pinger and Iperf as monitoring tools, for CrossGrid this choice has the advantage of both tools being available has RPMs in the EDG middleware releases. The EDG Pinger and Iperf have been modified to fetch configuration information from a central server and publish the monitoring results into the MDS information tree. A central configuration web server will be required to provide the necessary configuration information [VALID].

### **8.3. TESTS AND DEMONSTRATIONS**

Basic tests covering the Globus and EDG middleware functionalities have been performed using the testbed. These tests cover the job submission through Globus, job submission through the EDG RB, file transfer with GSI ftp, file replication with GDMP, file replication with the Replica Manager, the VO server and the MDS information system.

Three CrossGrid sites (FZK - Germany, IFIC - Spain and LIP - Portugal) have participated at the IST 2002 demonstration event showing the interoperation of the CrossGrid sites with other testbeds. During the demonstration several jobs have been successfully executed in the three participating CrossGrid sites [INIT].

In addition, CrossGrid has been testing new EDG middleware and reporting problems to DataGrid.

## 9. FUTURE ISSUES AND CONCLUSION

The most pressing issue to be considered in the future is a potential switch to an OGSI-compliant architecture as implemented in Globus Toolkit 3.0 (currently in preparation). While EDG is currently in no position to perform such migration (the amount of work put into already functional EDG prototypes precludes a major paradigm shift at this point in the development process), CrossGrid is in a better position to explore the potential advantages of OGSI/OGSA.

Several members of the CrossGrid Architecture Team are conducting research into potential applicability of Web Services in selected CrossGrid tasks. This knowledge may at a later stage prove useful to DataGrid development teams, possibly as a follow-up project to EDG itself. There is collaboration on specifying Web Service interfaces so that CG software can easily interact with EDG WP2 software and cooperation on code development, verification and testing of EDG WP2 replication software (in particular Reptor).