

# OPTORSIM: A SIMULATION TOOL FOR SCHEDULING AND REPLICA OPTIMISATION IN DATA GRIDS

David Cameron      Ruben Carvajal-Schiaffino      Jamie Ferguson      Paul Millar      Caitriana Nicholson      Kurt Stockinger      Floriano Zini  
 University of Glasgow      ITC-irst      University of Glasgow      University of Glasgow      University of Glasgow      Berkeley Lab      ITC-irst  
 d.cameron@physics.gla.ac.uk      carvajal@itc.it      fergusjk@dcs.gla.ac.uk      p.millar@physics.gla.ac.uk      c.nicholson@physics.gla.ac.uk      Kstockinger@lbl.gov      zini@itc.it

## Introduction

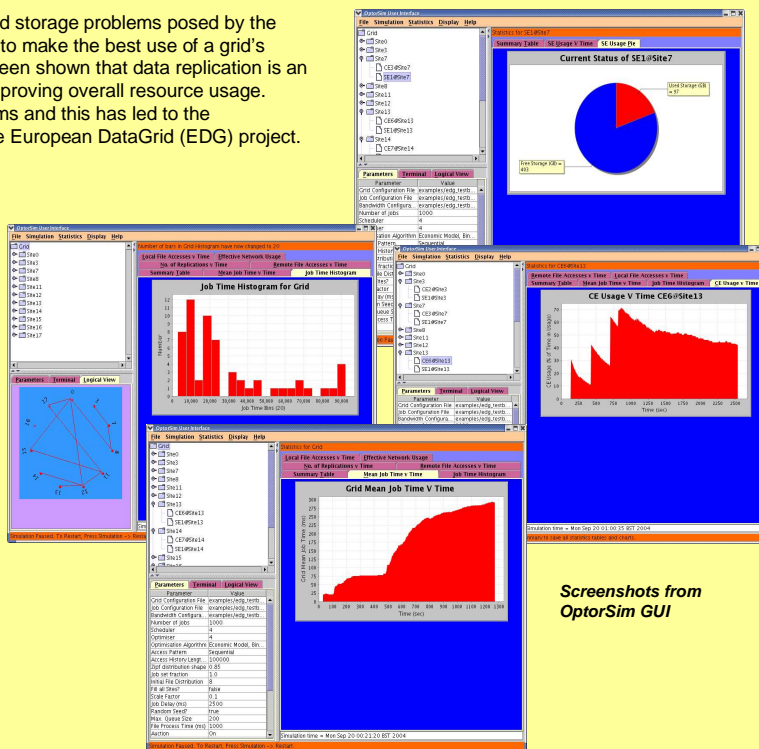
Grid technology is emerging as the solution to the data handling and storage problems posed by the next generation of high energy physics experiments. It is important to make the best use of a grid's resources, whether computational, storage or network, and it has been shown that data replication is an important mechanism for reducing data access times and hence improving overall resource usage. Simulation is a useful way of exploring possible replication algorithms and this has led to the development of the grid simulator OptorSim, originally as part of the European DataGrid (EDG) project.

With OptorSim, it is possible to simulate any grid topology and list of jobs to process by means of a few configuration files. There are several job scheduling and file replication algorithms implemented, and more can easily be added.

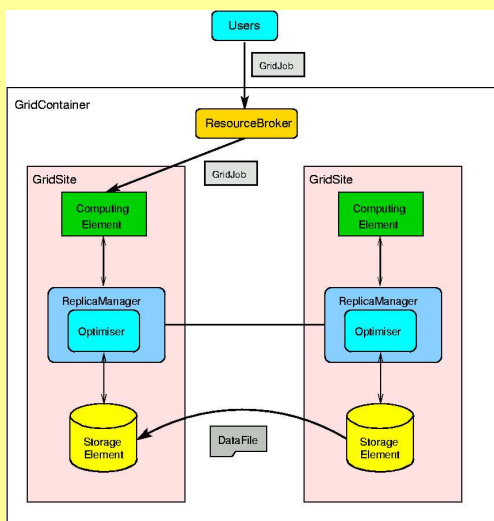
OptorSim can be run from the command line or from a GUI (right). A number of statistics are output:

- Total and individual job times
- CE usage
- Number of replications, local and remote file accesses
- SE usage

The appropriate statistics are output on the level of the grid, individual sites and site components. If the GUI is used, these can also be watched in real time.



Screenshots from OptorSim GUI



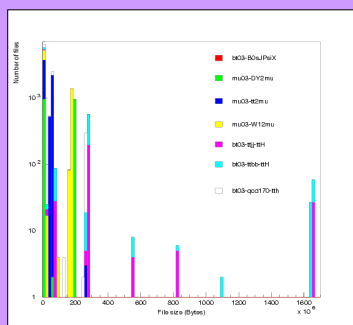
OptorSim architecture

## Architecture

OptorSim's architecture (left) is based on that of the EDG data management components. Computing and storage resources are represented by *Computing Elements* (CEs) and *Storage Elements* (SEs) respectively, which are organised in *Grid Sites*. CEs run jobs by processing data files, which are stored in the SEs. *Users* submit jobs to the grid according to the submission pattern chosen, then a *Resource Broker* (RB) schedules the jobs to Grid Sites. When a job is being processed at a CE, it will go through its list of files to process according to the chosen *access pattern*. If a file is required which is not present on the execution site, it must either be replicated or read remotely. Each site handles its file content with a *Replica Manager* (RM), within which a *Replica Optimiser* (RO) contains the replication algorithm which drives automatic creation and deletion of replicas.

## Inputs

Characterisation of aspects of real grid testbeds such as LCG has been used to get realistic input data, e.g. jobs and files based on the CMS 2004 data challenge (below).



File size distribution from CMS DC04 analysis jobs

## Optimisation Algorithms

There are two different types of optimisation which may be investigated using OptorSim: the scheduling algorithms used by the RB to allocate jobs, and the replication algorithms used by the RM at each site to decide when to replicate a file, which file to replicate and which to delete. The overall aim is to reduce the time it takes jobs to run, and also to make the best use of grid resources. In the short term, an individual user wants their job to finish as quickly as possible, but in the long term the goal is to have the data distributed in such a way as to improve job times for all users, thus giving the greatest throughput of jobs. Currently implemented are:

### Scheduling

- *Random* - schedule to random site
- *Access Cost* - site where time to access all files required by job is shortest
- *Queue Size* - site where job queue is shortest
- *Queue Access Cost* - site where access cost for all jobs in queue is shortest.

### Replication

- *No replication*
- *Least Recently Used (LRU)* - always replicate, delete least recently used file
- *Least Frequently Used (LFU)* - always replicate, delete least frequently used file
- *Economic model (Binomial)* - replicate if economically advantageous, using binomial prediction function for file values
- *Economic model (Zipf)* - replicate if economically advantageous, using Zipf-based prediction function.

