# The Worldwide LHC Computing Grid (worldwide LCG)

Jamie Shiers *

*CERN, 1211 Geneva 23, Switzerland*

**Abstract**

The world's largest scientific machine is scheduled to enter production in the second half of 2007 - roughly one year from the time of this conference. In order to exploit the scientific potential of this machine, computational resources way beyond what has been deployed for previous accelerators are required. Given these requirements - and the impracticality of providing such resources in a central place - a distributed solution based on Grid technologies has been developed. This paper describes the overall requirements that come from the Computing Models of the experiments, the state of deployment of the production services, on-going validation of these services as well as the offline infrastructure of the experiments and finally the remaining steps that need to be achieved in the remaining months before the deluge of data arrives.

*Key words:* Grid; Distributed Computing

## 1. Introduction

The European Centre for Nuclear Research - CERN - is situated just outside Geneva, straddling the Franco-Swiss border. It is currently constructing a new particle accelerator, the Large Hadron Collider (LHC). This machine is housed in an existing tunnel, some 27km in circumference, buried around 100m below the ground.

CERN is funded by 20 European member states, but has around 6500 visiting scientists from institutes all over the world, making the LHC a truly global project.

Four major collaborations or "experiments" are scheduled to take data at the LHC - each requiring massive detectors with many sub-components and read-out channels. ATLAS - the largest such detector - is some 44m in length and weighs a total of 7000 tonnes.

The data rates from these detectors are correspondingly massive - of the order of 1PB/s before a multi-level trigger, reducing the data rate to of the order of 100MB/s during proton-proton running and 1.5GB/s during heavy ion running.

In total, some 15PB of 'new data' are created every year, requiring the processing power of approximately 100,000 PCs, with the accelerator expected to operate for some 10 - 15 years.

At the time of writing, we are close to the 50th anniversary of the shipment of the first disk drive, the IBM 305 RAMAC. Launched on September 13th 1956, this drive stored a mere 5MB. If you could afford (to house) them, several billion such drives would be needed to store each year's data generated at the LHC!

## 2. Data Processing Overview

Data processing is performed in a relatively small number of steps, both for real data - i.e. that read out from the detectors as a result of the collision of

---

\* Corresponding author.
  *Email address:* `Jamie.Shiers@cern.ch` (Jamie Shiers).

particles accelerated in the machine - and for simulated data, which corresponds to a simulation of the postulated underlying physical processes, followed by a detailed simulation of the interaction of the final state particles with the detector matter itself.

Each processing step reduces the data volume by approximately one order of magnitude. However, the number of users increases with each such step, as does the degree of randomness of the associated access patterns.

## 3. Physics Goals of the LHC

The LHC is intended as a 'discovery machine', with a few high-level physics goals. It is not the purpose of this paper to describe this in detail, but more to focus on the computational techniques that are deployed to process the data and allow such discoveries to take place.

At the highest level, the physics goals can be stated thus:
– to explore a new energy / distance scale;
– to look for 'the' Higgs boson;
– to look for supersymmetry and / or extra dimensions;
– to find something that was not predicted by the theorists.

## 4. The Worldwide LHC Computing Grid

The Worldwide LHC Computing Grid, or WLCG, is unique in a number of important respects. Firstly, it does not actually exist as such, but is built upon existing Grid infrastructures - primarily the Open-Science Grid (OSG) in the US and the Enabling Grid for EsciencE (EGEE) in Europe and elsewhere. Additional resources are also provided in the Nordic region , which in term run their own middleware stack.

Taking the EGEE Grid as an example, many of the regions that make up this Grid are sub-Grids, managed in their own right with both global and local user communities.

Secondly, and equally importantly, its unique goal is to provide production processing capability for the LHC experiments. It focusses, therefore on the provision of reliable, stable and powerful services, with no R&D component.

Both of these criteria match up well to the 3-point checklist established by Ian Foster - as a service project, a non-trivial level of service is a requirement. By construction, it crosses (numerous) management domains. Furthermore, as it relies on multiple Grid infrastructures, the use of cross-grid standards is mandatory. This is perhaps most evident in the case of storage management, where all sites provide a Storage Resource Manager (SRM) interface to their storage resources (disk and tape). There are 3 main SRM implementations currently in production at the various sites that provide resources to the WLCG (CASTOR, dCache and DPM), with additional implementations either in development (e.g. STORM) or else deployed at a smaller scale (DRM).

Building on the existence of standard storage interfaces, the inter-site reliable file transfer service issue is simplified, although not at the level where multiple interoperable solutions currently exist. The primary solution in production use today being based on the gLite File Transfer Service (FTS), developed as part of the EGEE project, although other solutions, such as the Globus Reliable File Transfer (RFT) system have been evaluated.

## 5. Why a Grid solution?

There has been a great deal of hype concerning Grid computing, with promises of 'free' resources and limitless computing power. Furthermore, comparisons with the development and evolution of the Worldwide Web have become almost mandatory. Whilst there certainly are applications that are able to efficiently exploit unused CPU cycles - SETI@home being perhaps the best known example - this is not the driving force behind the choice of a Grid solution for the purposes of the LHC.

Significantly more pragmatic concerns include the funding model, whereby many institutes in numerous countries are required to provide resources to the overall effort. For these resources to be provided locally, with the added advantage of offering clearly identifiable resources to a local community, is sufficiently easier than a centralised model (which would also create considerable support and even infrastructure problems - such as the availability of the necessary space, together with adequate power and cooling.)

Rather, the Grid offers the promise of seamless access to these distributed resources, hiding the complexity and location of both CPU and storage resources behind consistent interfaces.

There are numerous cases of where the availability of massive resources through standard interfaces has allowed specific problems to be solved in a con-

siderably shorter time than using conventional systems.

Two such examples are UNOSAT: a United Nations initiative to provide the humanitarian community with access to satellite imagary and geographic system services with web based access to the Grid infrastructure using a portal placed in running in cell phones and the ITU: which established a new frequency plan for the introduction of digital broadcasting (band III and IV/V) in Europe, Africa, Arab States and former-USSR States in the 5 week period of the RRC06 Conference using Grid resources.

Many other examples exist in the area of drug discovery, pandemic and disaster response and so forth.

## 6. The WLCG Service Model

The main responsibilities of the different tiers of the WLCG computing model are as follows:
– Tier0 (CERN): safe keeping of RAW data (first copy); first pass reconstruction, **distribution of RAW data and reconstruction output (Event Summary Data or ESD) to Tier1**; reprocessing of data during LHC down-times;
– Tier1: safe keeping of a proportional share of RAW and reconstructed data; large scale **reprocessing** and safe keeping of corresponding output; **distribution of data products to Tier2s** and safe keeping of a share of simulated data produced at these Tier2s;
– Tier2: Handling **analysis** requirements and proportional share of **simulated event** production and reconstruction.

There are variations between the experiments for example LHCb does not plan analysis at Tier2 sites, whereas ATLAS foresee storing two copies of the ESD at Tier1 sites, with a further full copy at Brookhaven National Laboratory in the US. Given the spread in resources that the various ATLAS Tier1 sites will offer, this requires some pairing of sites with approximately balanced resources.

Approximately 10 Tier1 sites have currently been setup, in Europe, Asia-Pacific and North America. Some 100 Tier2 sites also exist, covering all regions of the world.

The resources provided by each tier are approximately equal - that is, the resources provided by the sum of the Tier2s is roughly that of the sum of the Tier1s equal in turn to that provided by the Tier0. Recent updates to the resource allocations suggest that this ratio is evolving in favour of the Tier1 and Tier2 sites.

## 7. Data Flows and Rates

Significant data transfers are required between all levels of the WLCG service hierarchy. The primary data transfers out of the Tier0 are to distribute a fraction of the raw data, where a second copy is stored across the sum of the Tier1 sites, and to distribute the results of the first reconstruction pass. These transfers are predictable - they should occur in pseudo-realtime with data taking, and for the duration of the operating cycle of the accelerator (some 100 days per year). Depending on the resources that a Tier1 provides as a fraction of the total resources requested by an experiment, these data rates vary from some 50MB/s to 200MB/s. Note that these transfers must occur 24 hours a day for roughly 100 days on end - a considerably different problem to a short-term network throughput test. These transfers take place over dedicated 10Gbit/s links from the Tier0 to each Tier1 (some Tier1 sites have additional high throughput connections to other Tier1 sites).

Data transfers from the Tier2 sites into the Tier1 are also predictable, but at much lower rates - some tens of MB/s - and correspond to the result of event and detector simulation that is primarily performed at the Tier2s.

Inter-Tier1 transfers are again scheduled, occuring with reprocessing activities - a responsibility of the Tier1 sites - that takes place approximately twice after the initial processing at the Tier0. The first such reprocessing is performed roughly one month after data taking, with improved detector calibrations. The second reprocessing happens after some months - one year, with both improved calibrations and reconstruction algorithms.

Perhaps the most tricky transfers to handle are the Tier1 to Tier2 transfers, corresponding primarily to the download of analysis data. By definition, these transfers are bursty in nature and are driven with the analysis needs of the end users that are serviced by the Tier2 site in question.

Demonstrating that all of the above transfers can be performed reliably and with sufficient throughput is one of larger challenges that WLCG still faces today. Data export from the Tier0 to the Tier1 sites has been demonstrated up to an aggregate daily rate of 1.6GB/s, but further reliability and service hard-

ening is required in the coming months.

## 8. The WLCG Service

The services that the various sites participating in the WLCG must provide, together with resource requires and availability targets, are defined in a Memorandum of Understanding (MoU) that is signed by each site. Service availability targets are as high as 99% - an aggressive target, particularly in the case of complex, cross-site services.

Currently, we are rather far from these targets, with average availability around the 80% mark. Clearly, much work needs to be done in terms of monitoring, as well as middleware hardening, if these figures are to be improved in time for LHC startup late next year.

Some examples are listed below.

The Host Laboratory shall supply the following services in support of the offline computing systems of all of the LHC Experiments according to their computing models. This includes operation of the Tier0 facility providing:

(i) high bandwidth network connectivity from the experimental area to the offline computing facility (the networking within the experimental area shall be the responsibility of each Experiment);

(ii) recording and permanent storage in a mass storage system of one copy of the raw data maintained throughout the lifetime of the Experiment;

(iii) distribution of an agreed share of the raw data to each Tier1 Centre, in-line with data acquisition;

(iv) first pass calibration and alignment processing, including sufficient buffer storage of the associated calibration samples for up to 24 hours;

(v) event reconstruction according to policies agreed with the Experiments and approved by the C-RRB (in the case of pp data, in-line with the data acquisition);

(vi) storage of the reconstructed data on disk and in a mass storage system;

(vii) distribution of an agreed share of the reconstructed data to each Tier1 Centre;

(viii) services for the storage and distribution of current versions of data that are central to the offline operation of the Experiments, according to policies to be agreed with the Experiments.

## 9. Preparing and hardening the service

In order to be ready to fully exploit the scientific potential of the LHC, a dedicated programme of "Service Challenges" has been established. These challenges are seen as an essential on-going and long-term commitment to achieving the goal of a production quality world-wide Grid at a scale beyond what has previously been achieved in production.

Whilst many of the individual components that make up the overall system are understood or even deployed and tested, much work remains to be done to reach the required level of capacity, reliability and ease-of-use. These problems are compounded not only by the inherently distributed nature of the Grid, but also by the need to get large numbers of institutes and individuals, all with existing, concurrent and sometimes conflicting commitments, to work together on an incredibly aggressive timescale.

The service challenges must be run in an environment that is as realistic as possible, which includes end-to-end testing of all key experiment use-cases over an extended period, demonstrating that the inevitable glitches and longer-term failures can be handled gracefully and recovered from automatically. In addition, as the service level is built up by subsequent challenges, they must be maintained as stable production services on which the experiments test their computing models.

## 10. Status of the WLCG Service

A service providing the basic functionality required by the experiments has been in place since September 2005. Whilst numerous problems have been encountered during this period, many problems have been resolved and the experiments have been able to run large scale productions and validations of their computing models.

Tier0-Tier1 data export tests have demonstrated rates that match the requirements for data taking under full operational conditions. These have been sustained over a two week period, even though not all sites were able to participate fully for the entire time.

More recently, experiment-driven transfers, involving a significantly larger software stack and more realistic file sizes and access patterns, have shown similar results, suggesting that we are on track in this respect.

Two of the large LHC experiments, ATLAS and

CMS, have each exported over 1PB of data in recent months: CMS has managed to export 1PB per month for a 90-day period, with ATLAS exporting 1.25PB in the two month period starting on June 19th.

However, the issues of offering a reliable service distributed over many sites and time zones has not fully been addressed, and much work in the area of monitoring and recovery procedures still remains.

At the time of writing, all experiments are launching large-scale, concurrent production activities, which will test further our readiness for full operation.

## 11. Plans Prior to First Collisions

First collisions in the LHC are currently foreseen for late 2007. These will most likely be at injection energy - 450GeV per proton beam - and not at the design energy of 7 TeV per beam. However, this will allow the machine to be fully debugged as well as providing data for the calibration and alignment of the highly complex LHC detectors.

Although the machine will initially operate at a lower energy, the data rate that must be supported remains the nominal value - the experiments will simply loosen their triggers to compensate for the lower energy collistions.

A number of new or revised sub-services need to be deployed in the remaining months, including new implementations of the SRMs that are in use at the various sites, corresponding to a new version of the specification - SRM 2.2. This also requires that higher level services, such as gLite FTS, are upgraded to conform to this new specification.

New services that are currently foreseen include distributed database services, based either on Web caching or database replication techniques, together with a conditions database serivce layered on the above.

A new release of the gLite middleware is also expected, together with an upgrade to the operation system in use in WLCG, namely migrating from Scientific Linux 3 to 4.

In parallel, the experiments will continue regular data challenges, ramping up in functionality and throughput, culminating in a 'full dress rehearsal' in the summer of 2007.

## 12. Summary

The WLCG is currently offering production services to all of the virtual organisations that will take data at CERN's Large Hadron Collider.

Testing of these services has to date focussed primarily on the Tier0 and Tier1 sites, with tests to prepare the Tier2 sites just beginning. The use of standard interfaces has been particularly successful in the area of storage management, whereas the ability of the LCG to exploit major resources offered by independent grids can be considered a major milestone in the deployment of production grids.

Perhaps one of the main lessons that has been re-learnt in this project is the need to design for failure: in such a highly distributed and complex system, it is the norm that one or more components are not available or not functioning correctly and the services have to be designed to take this into account.

On the other hand, the complexity of building a collaboration at a truly global scale should not be underestimated. We are still a long way from being able to provide such power at the level of basic infrastructure and communication remains one of the biggest problems that has to be constantly addressed. A series of global, regional and local workshops has been shown to solve some of the communications issues but is far from sufficient to solve problems that occur or information that needs to be disseminated on a much shorter timescale.

As a final remark, production, global, (cross-)Grid services are here today. Whether they will cross the chasm from the scientific domain into industry and finally into the home-user commodity market, is a question for a future conference.

The four main LHC experiments are ALICE, AT-LAS, CMS and LHCb.

# References

[1]   The Worldwide LHC Computing Grid (WLCG) -
      http://lcg.web.cern.ch/LCG/.
[2]   The Open Science Grid -
      http://www.opensciencegrid.org/.
[3]   The Enabling Grids for E-sciencE (EGEE) project -
      http://public.eu-egee.org/.
[4]   The Grid middleware development and testbed
      deployment project in the Nordic countries -
      http://www.nordugrid.org/.
[5]   Ian Foster, Argonne National Laboratory and
      University of Chicago, What is the Grid? A Three
      Point Checklist, (2002).
[6]   The Storage Resource Manager Working Group -
      http://sdm.lbl.gov/srm-wg/.
[7]   CASTOR - CERN Advanced STORage manager -
      http://castor.web.cern.ch/castor/.
[8]   dCache - http://www.dcache.org/.
[9]   The LCG Disk Pool Manager (DPM) -
      http://www.grif.fr/spip.php?article16.
[10]  STORM - An SRM Implementation for LHC Analysis
      Farms - http://hst.home.cern.ch/hst/publications
      /storm_chep06.pdf.
[11]  The Disk Resource Manager (or DRM) -
      http://vdt.cs.wisc.edu/components/drm.html.
[12]  Grid Data Management - A Comparison of File
      Transfer Service -
      http://www.gridpp.ac.uk/papers/chep06_stewart.pdf.
[13]  SETI@home - http://setiathome.berkeley.edu/.
[14]  UNOSAT: "Project Gridification: The UNOSAT
      Experience". EGEE User Forum, CERN, 1st March
      2006. Session: "Earth Observation - Archaeology -
      Digital library". Author: P. Mendez.
[15]  ITU: "International Telecommunication Union
      Regional Radio Conference and the EGEE Grid".
      EGEE User Forum, CERN, 1st March 2006. Session:
      "Earth Observation - Archaeology - Digital library".
      Author: A. Manara, M. Cosic, T. Gavrilov, P.N. Hai.