



LHC Computing Grid Project

Quarterly Status and Progress Report

2006 Q3

Executive Summary

Alberto Aimar
5 November 2006

1 Introduction

This document highlights major achievements and issues of the quarter, also providing an overview of the individual Quarterly Reports that were submitted by sites, experiments and projects.

For further details please refer to the complete document: “Quarterly Status and Progress Report - 2006 Q3” which contains the individual reports. All reports were reviewed by a Review Team (A.Aimar, B.Gibbard, U.Marconi, G.Merino, and L.Robertson). The reviewers then asked for additional information and, in some cases, the reports were modified and re-submitted by the original authors.

2 Level-1 Milestones

SC4-5: Service Challenge 4: Successful completion of service phase

More in detail, the targets to reach were:

- 8 Tier-1s and 20 Tier-2s must have demonstrated availability better than 90% of the levels specified in Annex 3 of the WLCG MoU [adjusted for sites that do not provide a 24 hour service]
- Success rate of standard application test jobs greater than 90% (excluding failures due to the applications environment and non-availability of sites)
- Performance and throughput tests complete: Performance goal for each Tier-1 is the nominal data rate that the site must sustain during LHC operation (see Figure 3): CERN-disk > network > Tier-1-tape. Throughput test goal is to maintain for one week an average throughput of 1.6 GB/s from disk at CERN to tape at the Tier-1 sites. All Tier-1 sites must participate.

Status - This milestone has been completed but the performance targets were not fully met. The reliability of the 8 best sites in each of the three months of SC4 was 83%, 86% and 82% respectively, compared with a target of 90%. The difficulties with the tests and the measurement system during this introductory period have now been resolved. Site reliability metrics were established in April and measured throughout the period of SC4 for CERN and 9 Tier-1s (BNL and NDGF did not participate) using the SAM testing system.

The individual data transfer target for each of the Tier-1s was achieved, except in the case of NDGF, at some time during SC4, and the aggregate target of 1.6 GB/s from CERN to the Tier-1s was met for a short period. However, the target of long-term stable transfers at this rate to tape at the Tier-1 sites was not demonstrated.

DRC-4: 1.6 GB/s data recording demonstration at CERN

Transfer “Data Generator to Disk and to Tape” sustaining 1.6 GB/s for one week using the CASTOR mass storage system.

Status -

In the beginning of October the LCG is running e.g. ATLAS at nominal speed, ALICE at nominal speed and CMS at 25% - all at the same time. The focus now is on direct data challenges together with the experiments.

This milestone was intended to check the concurrent running of all four experiments in one large pool. It is now clear that each of the four experiments will run completely independently with slightly different set-ups. Therefore this milestone will be replaced by experiment-specific milestones in the new planning for 2007.

LCG Quarterly Report – 2006 Q3 – Executive Summary

DBS-1: Full LCG database service in place

Status - This milestone has been completed for all sites that installed the Frontier/Squid framework (CMS Tier-1 sites) and for also for six of the ten Tier-1s concerned with the Oracle/Streams solution.

New milestones, and a closer monitoring of their progress, will be defined for the remaining Tier-1 sites that did not achieve their milestone and that are not yet providing a Distributed Database Service (NDGF, PIC, SARA-NIKHEF, and TRIUMF).

IS-1: Initial LCG Service in operation

Status - Started at end of SC4. It should be capable of handling the full nominal data rate between CERN and Tier-1 sites. The service will be used for extended testing of the computing systems of the four experiments, for simulation and for processing of cosmic-ray data. During the following six months each site will build up to the full throughput needed for LHC operation, which is twice the nominal data rate.

Ramp-up and commissioning plans will be developed during 2006Q4.

3 Service Challenge Progress

Completion of SC4 - The Service Phase of SC4 has been shown to deliver usable, albeit imperfect, services. The meetings and other mechanisms set up (daily and weekly operations meetings, service and experiment coordination meetings) have proven successful in resolving the major problems experienced by the experiments. Almost all of these issues were solved in production prior to the end date foreseen for SC4 (end September 2006), with a couple of remaining issues (LFC and dCache) for which a production release of the corresponding software was made within two weeks of the end of SC4. This provides a solid basis for consolidation.

While there is some way to go in improving reliability, it should be noted that there were also difficulties with the SAM tests and the measurement system during this introductory period that have now been resolved. The job reliability measurement system is under development and is being used within ATLAS and CMS as part of the experiment dashboards for some of their workloads, but there was no overall measure of job reliability during SC4.

As mentioned in the previous section (milestone SC4-5), the performance targets set for SC4 have not been fully met. The individual data transfer targets for each of the Tier-1s was achieved but the target of long-term stable transfers at this rate to tape at the Tier-1 sites was not demonstrated.

A number of additional components need to be added to the LCG service, including production Distributed Database Services, SRM v2.2-based services, as well as support for VOMS groups and roles.

Network Transfers - Stable transfers from Tier0 to tape at all of the participating Tier-1s at full nominal rates have still not been fully demonstrated. Latest estimates of the required data rates lead to an increased nominal rate of 2GB/s, primarily due to a doubling of the ATLAS ESD size (from 500KB to 1MB per event) and an increase in the trigger rate for CMS.

Transfers between Tier-1 and Tier-2 sites also have to be demonstrated as part of stable services - an activity which is expected to take much of 2007 to complete.

Data distribution tests were performed as part of the ATLAS production exercise, starting in June, to all ATLAS Tier-1 sites at the full ATLAS nominal rates. This rate is about 780MB/s out of CERN, with 40% of this data going to tape at the Tier-1 sites (corresponding to the fraction of the data for which that Tier-1 site has custodial responsibility). Concurrent transfers from CMS were also performed during this period at roughly 25% of the nominal rate for CMS (about 150MB/s out of CERN). LHCb transfers have since ramped up to the nominal rates and ALICE has also demonstrated transfers at the rate expected for "pp running", but a consolidated test of all VOs simultaneously still has to be performed.

Workshops and Reviews - A "Service Challenge Technical Day" was held at CERN as scheduled, focusing both on the results and recommendations from the SC4 Service Phase, as well as the issue of monitoring (with a view to delivering MoU quality services). This revealed that a great deal of work has been spent in the area of monitoring (site monitoring, grid monitoring, experiment dashboards, etc) without globally coordination and so with some overlap. This issue has been raised both in the LHCC

LCG Quarterly Report – 2006 Q3 – Executive Summary

Comprehensive Review of LCG and to the MB. A follow up session is foreseen for the WLCG Collaboration Workshop in January 2007.

4 Summary of Progress and Issues

4.1 Procurement

After the more detailed planning for the LHC machine became available, the resource pledges and requirements have been discussed several times. The result has been a better and more detailed understanding of the resources needed and provided in the LCG.

Delays in the delivery of the purchased hardware are announced in some of the quarterly reports; but will be without major consequences on the initial schedule for putting in production these resources. In addition, following the reschedule of the LHC programme for 2007 and 2008, several sites will re-plan their procurement schedule in the next few months. Some sites are reducing their disk expansion purchases, and wait for the availability of lower prices or for the introduction of new technologies.

4.2 Revised Plans

Resource Requirements - Revised estimates of the computing requirements of each of the experiments have been produced, following the announcement of the revised schedule for the accelerator in June. This is possible because of improved information about the operating conditions during 2007/08 and of the latest estimates of program performance and event sizes. In most cases the requirements during 2007 and 2008 are lower than before, offering an opportunity to the funding agencies to fulfill these requirements with the funding that is available.

Data Rates - However, the revised requirements lead to higher nominal data rates between sites, coming primarily from a doubling of the ATLAS ESD size (from 500KB to 1MB per event) and an increase in the CMS trigger rate. On the other hand, the LHCC referees during the comprehensive review of LCG in September recommended that the targets for sustained data rates in 2007 and 2008 be re-assessed taking account of the expected machine efficiency. This has to be discussed as part of the planning activity during the fourth quarter.

Milestones for LCG Initial Service - The quantitative measures in the Level-1 milestone for successful completion of SC4 (SC4-5) cover only reliability and Tier-0/Tier-1 data distribution rates. Additional quantitative targets should now be set for the Initial Service, such as job submission rate, number of simultaneous jobs, catalogue access rates, etc. Milestones should also be defined to follow experiment preparations for the beginning of their data challenges and to measure their success in using the services.

4.3 LCG Middleware Releases and gLite 3.0

gLite 3.0 Deployment - gLite-3.0 is now fully deployed at close to 90% of all EGEE sites and the Tier 1 sites in the US. Subsequent releases (gLite-3.0.1, 3.0.2) have been produced and deployed to address many bugs and address other issues. The release policy for gLite has now evolved, with releases [providing updates for individual components or services. Several of these releases have been produced, together with updates to the YAIM installation manager to support this mode of operation.

gLite Components - Significant effort has been devoted to debugging the gLite WMS to get it into a state suitable for use in the ATLAS and CMS service challenge. A large number of bugs were found and addressed rapidly by the developers. In addition considerable effort was put into understanding the interaction between the distributed components of the WMS and setting appropriate timeouts etc. The result of this work was to demonstrate that a single RB is capable of around 25-30K job submissions per day. This is acceptable for ATLAS and CMS, but needs to improve over the next year. All components needed for job priorities via VOMS roles and groups are included in the WMS and CE now.

The gLite CE is not yet ready for general deployment. No intensive testing and certification has been performed. All available testing effort has been devoted to making the gLite WMS sufficiently stable and scalable in order to use it in the service challenge productions. Similar efforts will now have to be devoted to the gLite CE.

LCG Quarterly Report – 2006 Q3 – Executive Summary

SL4 Porting - There is a delay in building the distribution for SL4 for both 32 and 64-bit machines. A number of issues were found with library incompatibilities in VDT-1.3, which is needed for this version, as well as SSL problems for the 64-bit build. These are being resolved together with the VDT team, and the plan is for a gLite distribution by the end of the year. Meanwhile, it has been shown possible to run the WN code on SL4 (32&64 bit).

4.4 SRM Interface and Implementations

At the Mumbai workshop in February 2006 an initial agreement between the experiments, developers and major sites was reached on the functions to be provided for the initial LCG service and on the way in which these would be implemented within the SRM standard. A final agreement on how to do this was reached at a workshop at FNAL in May, along with an implementation plan. The target was for this new version, SRM v2.2, to be available for the three mass storage systems used in LCG (dCache, DPM, Castor) at the end of October for testing. This should allow for production services based on SRM 2.2 in the first quarter of 2007. A reviewed schedule has to be developed for the deployment of SRM v2.2 in production at sites and for the ramp-up of capacity. This will lead to a series of new Level-1 milestones and these should be translated into individual site targets and milestones.

All production and general user work at CERN was migrated to Castor 2 by June 9th. The CERN Castor 1 services for LHC were removed in July. INFN-CNAF finalized the migration from CASTOR 1 to CASTOR 2 moving all experiments, with CASTOR 1 being phased out at the end of September. RAL has begun migration to Castor 2, and PIC reported little progress on the CASTOR2 installation, but actions have been taken in this quarter for completing the CASTOR 2 migration by early 2007.

4.5 Sites Reliability Monitoring

SAM and GridView - The SAM system is now in full production and used since May 2006. Availability metrics have been published for CERN and nine Tier 1 sites since that time. Some of the tests have shown problems and will need adaptation for certain situations, and the full set of service tests is not yet there. In particular experiments-specific tests are needed in order to provide a realistic verification of the services that are needed by each experiment.

The GridView tool (<http://gridview.cern.ch>) is also in production since several months and has had numerous updates to include display of the SAM metrics and job status statistics as well as the file transfer statistics.

Grid Dashboard - A number of other systems are now in operation to measure service and site performance and reliability (FTS service throughput, resource usage, site availability, etc) and these should be integrated into a "dashboard" showing the status of each site compared with the targets, adding additional factors such as installed capacity vs. current VO requirements, status of new services, etc.

4.6 Accounting and Security

Accounting- Progress has been made on accounting procedures, and an initial prototype of automated storage accounting is being tested. A solution for the user privacy issues for user-level accounting has been proposed but not yet implemented.

Security - A security vulnerability incident, affecting many large HEP sites among others, occurred during the summer. While this was not grid-related (i.e. did not make use of grid mechanisms) it did serve as a test for the incident response procedures and teams. The experience resulted in a number of points that are being followed up.

4.7 Experiments Preparation

Tier-1 - Tier-2 Associations – An activity has been started to integrate data from all four experiments on the expected relationships between Tier-1 and Tier-2 sites, including data rates and storage requirements. This will enable sites to configure and test the necessary FTS services, verify that there is adequate networking performance, and initiate operational relationships with their partner sites.

LCG Quarterly Report – 2006 Q3 – Executive Summary

DAQ-T0-T1 Data Flow - The delay in testing the full data chain, including the DAQ, is of some concern. This is recognised by most of the experiments and it is hoped to begin this testing early in 2007.

24 X 7 Support - The 24x7 support plans are under discussion at Tier-1s and CERN.

VO Boxes Service Level - The VO Boxes needed by the LHC experiments are installed at all required Tier-1 sites and are functioning adequately. Further work is needed between sites and experiments in order to agree on how VO Box configurations and data are stored and backed up (by the sites or by the VOs).

In addition, as agreed at the VO Box Working Group, the services provided by the VO Boxes should be implemented by future services (or improvements) of the general middleware software, making unnecessary the deployment of any specific VO Box. Work in this direction is in progress under the steering of the EGEE Technical Coordination Group (EGEE TCG).

The GDB had agreed that all “class-1” services in the VO boxes should be removed as a priority, but this remains an issue to be resolved.

5 Specific Areas and Projects

5.1 Applications Area

A couple of new complete software configurations (LCG_46 and LCG_47) have been made available during the last quarter in the Applications Area. They include the new releases of the ROOT, CORAL, POOL and COOL packages, which are currently used by the experiments for the various data challenges.

ROOT - The production version of ROOT 5.12.00 was released on July 11th as scheduled with a new version of the User’s Guide. It includes, in addition to several new functionalities and bug fixes, a very interesting set of I/O optimizations for remote access such as read ahead, request grouping, etc. The work in adapting the CINT interpreter to the Reflex data structures is progressing rapidly and getting better estimates of the amount of work involved. It is expected to have a first version for test purposes by the December release of ROOT.

PROOF - For PROOF, the last quarter has seen a lot of work in supporting the ALICE tests on the CAF. The main emphasis has been on implementing the ALICE prompt analysis CAF use cases. A number of new features have been implemented in the area of data file uploading, package management and monitoring. Also a lot of attention has been given to aspects of system reliability, robustness and performance. Quite positive results have been extensively reported during a special ALICE - IT meeting and during the LCG AA internal review.

Generator Services - The Generator services sub-project is being re-structured as a result of a number of concerns expressed by the leading authors of MC generators during the MC4LHC workshop in July at CERN. The project received the strong support from the experiments to continue under its original mandate: to provide services for well-maintained repositories of MC generators on LCG-supported platforms. The new project leader is preparing a new plan taking into account their concerns and this will be presented on a special meeting in which all the stakeholders of the project are invited.

AA Internal Review - The AA internal review took place from September 18th to 20th. This review has been an opportunity for the AA projects to take a close look at the current status and at what still needs to be done in terms of new functionality before the LHC start-up. It has also been very useful to inform the experiments and other projects of the progress that have been made during the last 18 months. The material for the review can be found at <http://lcgapp.cern.ch/project/mgmt/rev200609>. The final review report is not yet available at this time. Once this is done the report will be studied and possibly a number of new milestones will be proposed to cover the recommendations by the review committee.

5.2 Distributed Database Deployment – 3D

During the last quarter the production infrastructure has been largely taken over by the experiments for deployment tests with their applications. All Phase 1 sites are now in direct contact with the experiment teams and performance and stability results have been reported at the last 3D workshop (13-14 Sept at CERN). Based on these results the experiments have confirmed their resource requirements at the Tier 1s

LCG Quarterly Report – 2006 Q3 – Executive Summary

for the next 6 month, when the next scheduled resource review will take place. The replication test with LFC read-only replicas for LHCb has been completed successfully. Streams replication will now move into production between the LHCb database at T0 and CNAF as the first Tier 1 site.

The Phase 2 sites (NDGF, NIKHEF/SARA, PIC and TRIUMF) are all actively working on the commissioning of the requested database setups, but these were not be available as of October 2006. These sites have been asked to provide their plans for completing their deployment.

On the Frontier side with only very few exceptions all installations have been successfully done at CMS Tier 1 and Tier 2 sites. CMS was able to run successfully the first larger scale tests with the new CMS software from 200 client nodes at CERN against the Tier 0 Frontier/Squid installation.

5.3 ARDA

Job Reliability - The results of the analysis of job logs on the EGEE grid have been discussed in several meetings, notably in the EGEE TCG. Since August a daily report is produced showing the success rates for certain CMS job sets. The web-based report supports drill-down functionality to allow the reasons for job failures to be established. The report is being adapted for ATLAS and also for the grid operations team.

Ganga - Ganga public release 4.2 was made available. The user base is increasing, especially in LHCb, after the plenary presentation at the last collaboration meeting. For LHCb the main news is the integration with the LHCb bookkeeping system allowing the users to seamlessly select their datasets and import them into Ganga. For ATLAS, the main news is the integration with the distributed data management system (both to select the input data and to return the output data) and the new gLite Resource Broker.

Software Validation - Within the CMS task force, the job robot system has been used to study performance and reliability of the WMS in view of CSA06. ATLAS production jobs have been submitted via the new WMS. The main result was a substantial improvement in the reliability of the new system and the demonstration that it was satisfying the CMS (and ATLAS) requirements for the near future.

During August, the API-service (bridging users' requests and the grid infrastructure) was upgraded due to instabilities triggered by the increased load. In September, the number of active users ramped up to ~40 as presented in the LHC comprehensive review.

Task Forces Activities – During this quarter the activities connected to the task forces have progressed at a steady pace. The analysis system of ALICE is in use (~40 users in September). The system is the result of the integration of the original ALICE system with contribution of ARDA and EGEE middleware. In parallel, production activities are continuing on LCG resources. The main activity in CMS has been the finalization of the validation of the WMS 3.0.x in view of CSA06. The work during summer (in close contacts with all the players in the field, notably the developers and in collaboration also with the ATLAS production experts) yielded a solid system (now in production). ATLAS is ramping up production activities on LCG. LHCb is in production with DC06.

6 Experiments

6.1 ALICE

The AliRoot software has gained in stability and is approaching the final configuration for processing real data in winter 2007. Technical changes are continuously implemented to improve the modularity of the framework and the robustness of the code. The memory consumption, which was a major concern, has been significantly lowered during data reconstruction. The calibration and alignment framework has been fully implemented and has been successfully exercised in production tests stressing the queries to the Offline Calibration Data Base. The simulation of the raw data flow has been implemented for all detectors and includes: the generation of simulated data in the raw-data format, the reconstruction starting from the raw-data format and the embedding of simulated events into real raw data.

The Physics Data Challenge is in production mode since April 2006. All the ALICE Tier-1 sites, except NGDF, and about 30 Tier-2 sites contribute to the exercise. This Physics Data Challenge is also extremely useful for the training of experts in the sites and the collection of operational experience.

LCG Quarterly Report – 2006 Q3 – Executive Summary

The central ALICE services have reached a near production level. Seven millions of pp events have already been produced as requested by the Physics Board.

The data movement challenge, coordinated with the LCG SC4, involving transfers from the Tier-1 sites to the CERN Tier0 and replication of data from the Tier0 to the Tier-1 site has not yet reached its goal of sustained 300MB/sec transfer rate. So far only a peak rate of 150MB/sec during a short period could be achieved.

The parallel analysis framework based on PROOF and to be deployed on the CAF is in use and tested by several concurrent users. The expected performance could be reached demonstrating the validity of the approach. Many analysis modules integrated in the AliRoot analysis framework are developed within the Physics Working Group and will be tested with the data produced in the PDC either within the distributed environment or in the CAF.

The situation of the computing resources made available to ALICE is evolving slowly with the addition of a few new sites. The deficit of resources has also been slightly reduced in 2007 and to lesser extent in 2008 taking into account the revised LHC start-up scenario. However, in 2009 the situation will become again worrisome. In addition, the resources made available or usable for the PDC by the sites are far from the pledged resources.

6.2 ATLAS

The ATLAS “Tier-0 Functionality Tests” have been performed in June-July and again in September-October, reaching the nominal internal data transfer rates. In addition also “Tier-0 to Tier-1” tests have also been performed, as well as partial “Tier-1 to Tier-2” transfers, during 2006Q2 and Q3. Nominal rates have not been reached at all sites (and NDGF did not take part), several issues were identified and addressed.

The ATLAS Software release 12.0.3 is now in production. Full simulation of 20M events during the next two months is foreseen for Computing System Commissioning (CSC) tests. Tuning of the reconstruction continues in view of release 13, in January 2007.

The performance of the DDM system, when used as part of the distributed production system and to transfer data other than the “Tier-0 to Tier-1” scheduled transfers, is still below expectations in terms of robustness. We are about to review our components and the way we use Grid m/w tools. The “Tier-0 to Tier-1” export data rates are much lower than we thought should be achievable at this stage; the rates absorbed by several sites fluctuate frequently, decreasing the total export rate to an average of 400 MB/s, instead of the target 800 MB/s.

6.3 CMS

CMS computing activity in last three months was mostly carried out under two main streams: (1) SC4 to validate sites and tools, and (2) production of simulated data in preparation of CSA06 (combined Computing, Software and Analysis challenge).

During SC4, the scaling behaviour of the system was tested using automated workloads submitted by JobRobots and ten thousands of jobs were run per day. In a period of 3 months ending in mid-August a total of 3 Petabytes of CMS datasets were transferred between storages systems through the worldwide Grid. After a lot of debugging, integration and commissioning work, data transfers from CERN to the Tier-1 centers regularly achieved rates of 150MB/sec, corresponding to about 25% of the nominal rates. Full rates of up to 450MB/sec were achieved for extended periods, although not regularly. Work on improving throughput and robustness of data transfers is continuing as part of the joint CMS/WLCG Integration Taskforce.

The workflow submitted via the JobRobots were realistic analysis jobs using new CMS Software framework (CMSSW) reading data from local disk storage via the CMS data management catalogs, and this activity allowed to validate and exercise toward CSA06 goals all the CMS Tier-1 sites and 20 Tier-2 sites. The CMS/WLCG Taskforce has also been working on the new gLite WMS toward the goal of 50K jobs submitted per day, which we can not meet with current LCG WMS. In preparation of CSA06, more than 60M events were simulated with the new software and the new computing systems, the simulation was

LCG Quarterly Report – 2006 Q3 – Executive Summary

carried out at all Tier-1 sites and 20 Tier-2 sites. Data have been staged locally at the site, merged into large files and transferred to CERN.

6.4 LHCb

The activities of data distribution and reconstruction have gone smoothly at five Tier-1 sites. There were problems with access data from the storage systems at three others sites. These issues seem solved at two sites. The new version of dCache is expected to solve the remaining issue. The stripping phase was delayed but now is imminent