**Worldwide LHC Computing Grid Project**
**Project Status Report**
**Resource Review Board – 12th April 2011**

This status report covers the period from October 2010 – March 2011. Further details on progress, planning and resources, including accounting and reliability data for CERN and the Tier 1 centres, and detailed quarterly progress reports, can be found in the documents linked to the LCG Planning Page on the web.

## 1. The WLCG Service

During this reporting period, the accelerator provided higher and higher luminosities for the final month of p-p running, as well as a rapid switch to heavy-ions and several weeks of heavy-ion data taking. The technical stop from December to March was also very active for the computing, with all experiments re-processing their full 2010 data samples.

### Tier 0 performance

The p-p run concluded at the end of October, and several weeks of heavy ion data were taken by ALICE, ATLAS, and CMS. In September both CMS and ALICE requested resources at CERN to enable them both to run at significantly higher data rates than had been anticipated or tested for. ALICE requested essentially double their nominal rate (2.5 GB/s in place of 1.25 GB/s) to allow them to
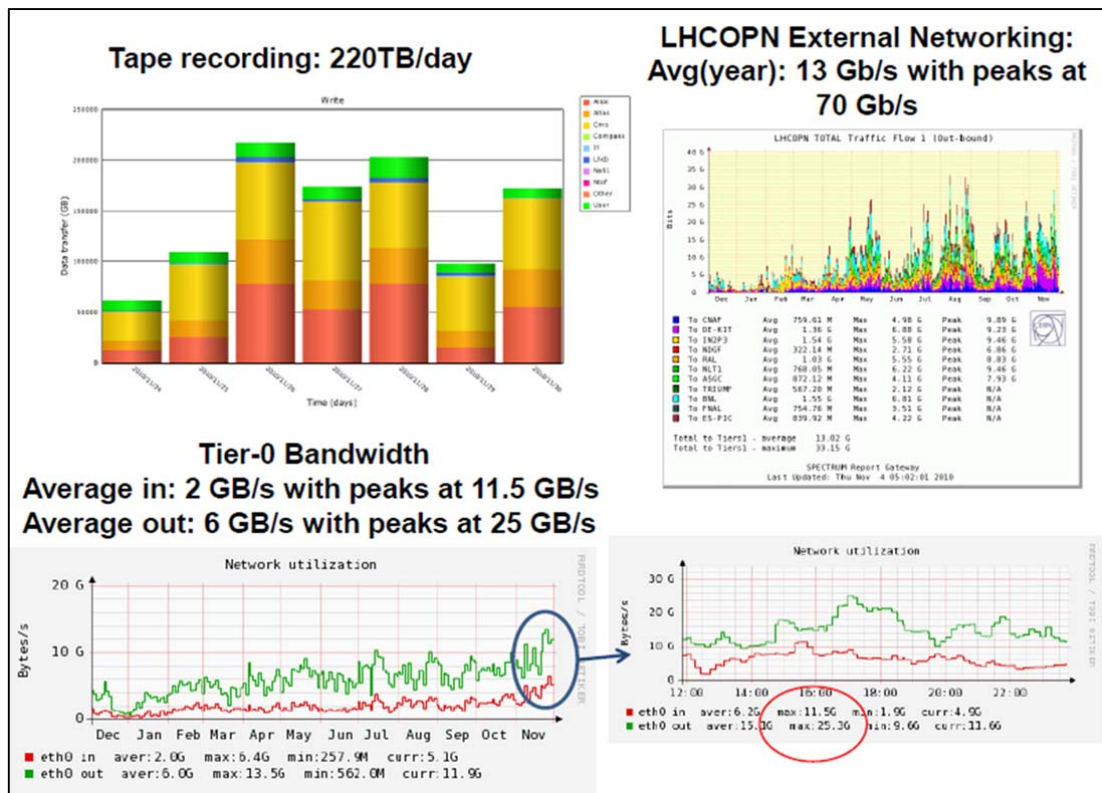


Figure 1: Some Tier 0 performance measures

acquire as much data as possible in the ion run, while CMS planned to take HI data without zero suppression of their detector, which would then be done offline. This would result in a data rate into the Tier 0 of 1.8 GB/s. In addition, because CMS would not be able to export this data before the zero suppression was performed they would be exposed to the same risk of data loss as ALICE during the period before a second copy of the data would be available at the Tier 1s. For these

reasons, CERN added sufficient disk capacity to both ALICE and CMS to allow them to have a full copy of data on disk as well as on tape until the data was copied out to the Tier 1s.  It was able to do this allocation only because the 2011 resource pledges had already been purchased and installed.

For ATLAS the plan was to take HI data at the same rate as for p-p, and to manage the data in exactly the same way as during proton running.

In fact the Tier 0 storage system (Castor) was easily able to manage these very high data rates. Figure 1 shows some of the performance metrics during this time.  The amount of data written to tape reached peaks of over 220 TB/day (a new record), while the bandwidth of data movement into the Tier 0 reached peaks of 11.5 GB/s, and transfer out reached peaks of 25 GB/s.  These rates are far in excess of the original plans, but were managed without problem by the system.  These rates *averaged over the entire year* are 2 GB/s and 6 GB/s respectively.

During p-p running some 2 PB/month were stored into the Tier 0, and in the month of HI running this was 4 PB.  Figure 2 shows the total data written into and read from Castor during 2010, showing that the estimates of ~15 PB/year of data were close to the reality, even in this first year.
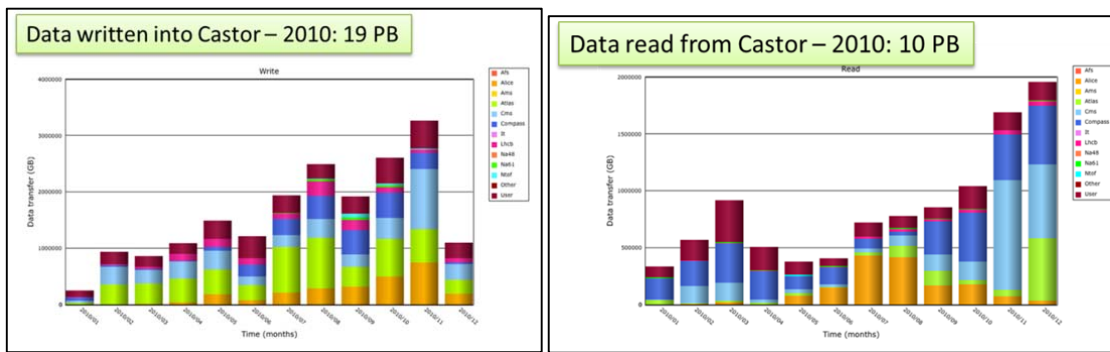


**Figure 2: Totals written and read by Castor in 2010**

### Data transfers

As noted above the main feature of this period has been the heavy ion data taking, Figure 3 shows the transfers of heavy ion data to the Tier 1s – for ALICE in December following the HI run as planned, and for CMS once the zero suppression had been done in the first months of 2011.
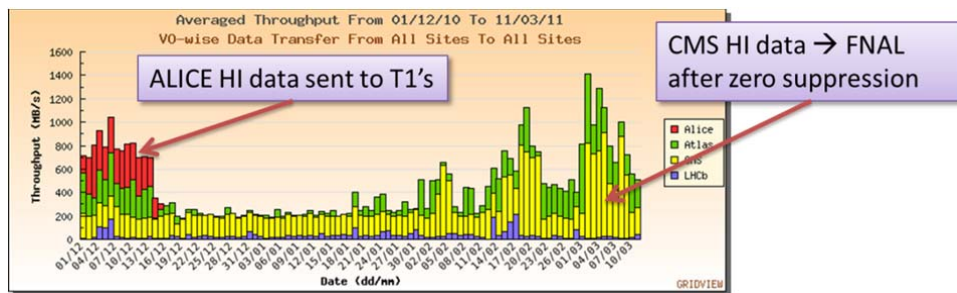


**Figure 3: Data transfers during winter technical stop, showing HI data transfers**

Since ATLAS took data in HI running at a rate close to that during p-p running, their transfers to the Tier 1s were done as the data was being acquired in the usual way.

Traffic on the LHCOPN globally averaged some 13 GB/s during the year with much higher peaks (see also Figure 1).

## Experiment Activities

During the period since the last RRB, the main activities have been the full reprocessing of the entire 2010 p-p data sample for each experiment, as well as the heavy ion data taking and processing. Some of the main features of these activities for each experiment are the following:

For ALICE the raw data has now been fully copied to the Tier 1s following the HI run; they reached a maximum transfer rate of 260 MB/s during the HI run itself. The HI data have been reconstructed once, and the 2nd pass is in progress.

ATLAS has done the full reprocessing of the full 2010 pp dataset and re-distribution of the results - this was completed by the end of 2010. Their HI data sample has been processed and re-processed once.

CMS also completed the reprocessing of the full 2010 pp dataset by the end of 2010. During the HI run CMS reached average transfer rates into Castor of 2 GB/s and average access rates of 3-5 GB/s out to the Tier 0 farm. The zero suppression of HI data is under way, the suppressed raw data is being stored at the Tier 0 and a copy transferred to FNAL. This copy reaches transfer rates of 800 MB/s.

LHCb has also completed the full reprocessing of all the 2010 data by end of 2010. A major MC production campaign is under way and a global disk clean-up campaign is also ongoing in preparation for the 2011 data. Space is at a premium since LHCb can now store fewer copies of the data than planned since the higher pile-up means that event sizes (and thus datasets) are significantly larger.

## Use of resources and workloads

During 2010 the Tier 1 and Tier 2 resources were more and more used, particularly towards the end of the year where 100% usage of the available resources at many sites was common. This is illustrated in the two plots below – showing average usages with respect to the pledges for the Tier 1 and Tier 2 sites.
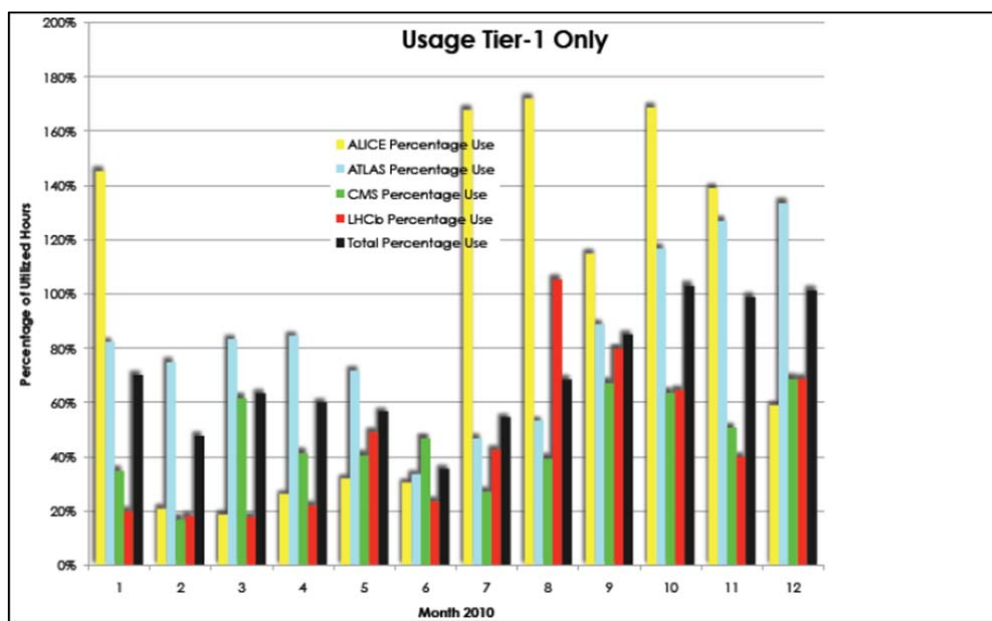


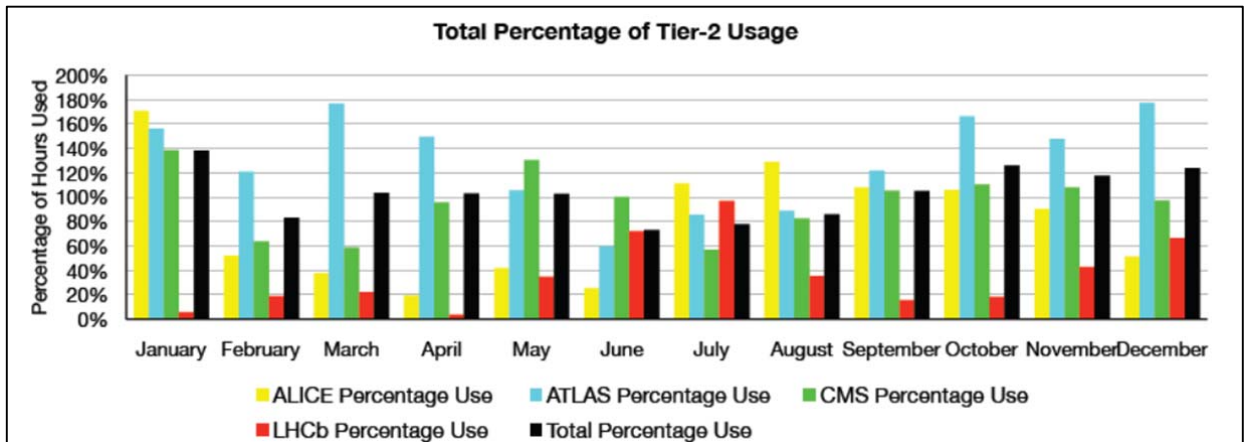Figure 4: Usage of Tier 1 resources in 2010

Figure 5: Usage of Tier 2 sites



Figure 6: Evolution of jobs run per month

The workload in terms of the number of jobs also continues to increase, as shown in Figure 6.

## WLCG Service status

As agreed, significant service interruptions require a documented follow up (Service Incident Report – SIR). The full list for this period (summarised in the Table below) including the full incident reports can be seen as a summary in each Quarterly Report, or consulted on line at https://twiki.cern.ch/twiki/bin/view/LCG/WLCGServiceIncidents. These are followed by the Management Board, with the goal being that lessons are learned and disseminated to other sites.

Table 1: Service incidents requiring a review, Q4 2010 - Q1 2011

| Site | Service Area | Date | Duration | Service | Impact |
|------|-------------|------|----------|---------|--------|
| **Q1 2011** | | | | | |
| IN2P3 | Storage | Mar 19 | 3.5h | SRM | dCache SRM was unusable due to internal overload |
| CERN | Infrastructure | Mar 19 | 12h | Batch system | Job submission became slow, then completely unresponsive |
| IN2P3 | Network | Mar 14 | 40min | Batch system | no connection to other French sites, but no problems observed for jobs |

| CERN | DB | 11-Mar-11 | 5h | CMS offline production db | The database was completely down for ~2 hours and partially not available for 5 hours |
|------|----|-----------|-----|-----------|-----------|
| IN2P3 | Infrastructure | Feb 25-26 | 13h | Batch system | 85% of batch system unavailable, jobs lost |
| IN2P3 | Storage | Feb 13 | 3 h | Storage service | Storage services degraded, no big impact on jobs |
| PIC | Storage | 21-Jan-11 to 08-Feb-11 | 18 days | Storage service | 250TB of ATLAS data partially unavailable |
| KIT | infrastructure | 28-Jan-11 to 02-Feb-11 | 5 days | Batch system, job submission | batch system degraded, reduced # of job slots available |
| CERN | DB | 25-Jan-11 | 8h | FTS, LFC, SAM,VOMS, dashboards | affected services fail, clients may hang |
| IN2P3 | infrastructure | 8-Jul-10 to 7-Jan-11 | 6 months | shared s/w area | jobs fail |
| CNAF-BNL | network | 23-Aug-10 to 20-Jan-11 | months | primary OPN circuit | poor transfer performance; ok when switched to backup |
| **Q4 2010** | | | | | |
| CERN | DB | 18 Dec | 5 days | DB | Service interruption: ATLARC DB following the power cut at CERN CC |
| CERN | infrastructure | 18 Dec | 26 hours for services with weight > 50 | power | Interruption of physics services following power cut |
| CERN | DB | 16 Dec | 2.5h | DB | ATLR database affected (degradation then complete outage) by FC switch replacement |
| CERN | infrastructure | 7 Dec | 7 days | CVS | CMSSW CVS migration problems |
| CERN | DB | Nov/Dec | 8 days | DB | Reboots of Instance 4 of ATLR database |
| KIT | infrastructure | 26 Nov | 1.5h | GGUS | No web access / no ticket update |
| KIT | infrastructure | 16 Nov | 3.5h | GGUS | No web access/ no ticket update |
| IN2P3 | infrastructure | 11 Nov | months | AFS | shared s/w area |
| NL-T1 | DB | 26 Oct | 48h | DB | Inconsistency of data at SARA |
| CERN | infrastructure | 20 Oct | 4.5 h | Batch | Severely degraded response from CERN Batch Service |
| CNAF | storage | 6 Oct | 5 days | CMS storage | CMS storage down (service interruption) due to GPFS bug |
| CERN | infrastructure | 4 Oct | 2.1 h | MyProxy | Outage on myproxy.cern.ch after incorrect certificate renewal |

Out of 23 service incidents, half (11) were infrastructure related. 6 were database problems, 4 were storage-related, and 2 were network issues. The infrastructure problems include major problems such as power or cooling, as well as cases where a site-specific service failed. Some of the database problems should strictly fall into this category also, as could most of the storage-related problems. This illustrates quite strongly that the majority (>~75%) of the problems experienced are not related to the distributed nature of the WLCG at all.

The major power cut at CERN on 18[th] Dec was on the first day of the Christmas closure, and while critical services on reliable power were maintained, for some of the lower priority services it took 24 hours for full restoration, and one of the ATLAS databases had a problem for several days following the cut. A long-standing problem at the Tier 1 in Lyon which affected fraction of LHCb jobs was finally traced to a problem in AFS triggered by the heavy way in which LHCb jobs accessed the software area. This has now been resolved with a new AFS client and some tuning of the worker machines.

The other problem of note here is the ongoing network problem in the OPN between CNAF and BNL. While this has not been truly "resolved", since the fix is to use an alternate route, it has nevertheless been useful in pointing out problems in the management of such 3[rd] party issues (where the actual problem is not in the domain of the sites themselves or their local network providers). As a result of this better procedures have been introduced for tracking and assigning responsibility for follow-up of such cases.

## LHCOPN

As well as the difficulty of the complex network operations model mentioned above, the absence of user-level network monitoring had also been noted in the previous report. However, a prototype LHCOPN dashboard (http://casper.grid.sara.nl) has now been produced, which together with the LHCOPN weather map (https://netstat.cern.ch/monitoring/network-statistics/weathermap/?map=LHCOPN), provide a useful first check on the network operation and performance from a user view point. Snapshots of these tools are shown below.



**Figure 7: Network monitoring tools**

## 2. Site Reliability

The reliabilities for the last 6 months for CERN and the Tier 1 sites are shown in Table 2.

**Table 2: WLCG Tier0/1 Site Reliability – last 6 months**

| | Average of the 8 best sites (not always same 8) | | | | | |
|---|---|---|---|---|---|---|
| Target | Oct-10 | Nov-10 | Dec-10 | Jan-11 | Feb-11 | Mar-11 |
| 98 | 99 | 99 | 100 | 100 | 100 | 100 |
| | | | | | | |
| | Average of ALL Tier 0 and Tier 1 sites | | | | | |
| Target | Oct-10 | Nov-10 | Dec-10 | Jan-11 | Feb-11 | Mar-11 |
| 97 | 96 | 98 | 99 | 99 | 99 | 98 |
| | | | | | | |
| | Detailed Monthly Site Reliability | | | | | |
| Site | Oct-10 | Nov-10 | Dec-10 | Jan-11 | Feb-11 | Mar-11 |
| CA-TRIUMF | 100 | 100 | 100 | 100 | 100 | 100 |
| CERN | 100 | 100 | 100 | 100 | 99 | 100 |
| DE-KIT (FZK) | 99 | 94 | 100 | 98 | 99 | 100 |
| ES-PIC | 97 | 99 | 99 | 100 | 99 | 100 |
| FR-CCIN2P3 | 96 | 98 | 100 | 100 | 97 | 100 |
| IT-INFN-CNAF | 100 | 98 | 99 | 100 | 100 | 100 |
| NDGF | 92 | 99 | 98 | 99 | 100 | 98 |
| NL-T1 | 87 | 93 | 98 | 98 | 98 | 97 |
| TW-ASGC | 83 | 93 | 95 | 100 | 99 | 79 |
| UK-T1-RAL | 100 | 99 | 99 | 99 | 100 | 99 |
| US-FNAL-CMS | 100 | 98 | 100 | 100 | 100 | 100 |
| US-T1-BNL | 100 | 100 | 99 | 100 | 96 | 100 |
| Target | 97 | 97 | 97 | 97 | 97 | 97 |
| Above Target (+ >90% Target) | 8+2 | 9+3 | 11+1 | 12+0 | 11+1 | 11+0 |
| | | | | | | |
| Colours: | Green > Target   Orange > 90% Target   Red < 90% Target | | | | | |

The figures below show the recent evolution of the reliabilities for the Tier 1 and Tier 2 sites. These figures are now quite stable for the Tier 1 sites and the majority of the Tier 2 sites (which provide most of the resources). Some of the smaller Tier 2s could still improve their overall level of reliability quite significantly, and consequently improve their overall usage.
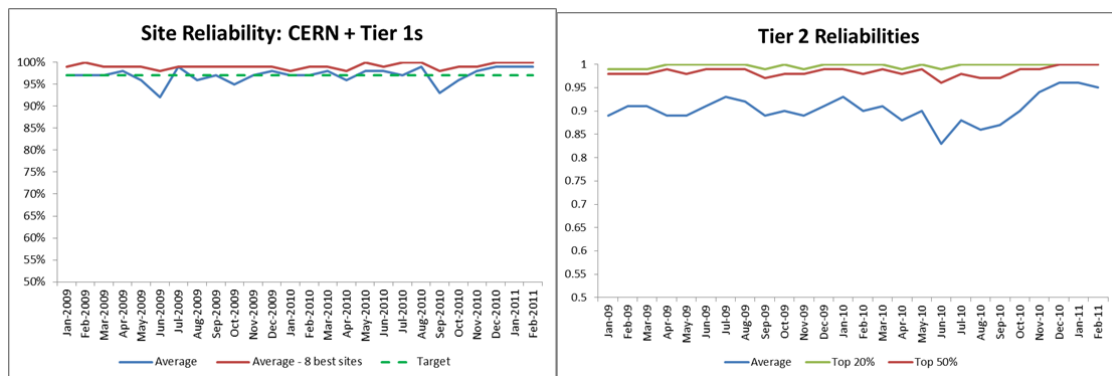


**Figure 8: Site reliability evolution**

This is reflected in the typical site readiness measured by each experiment, and reviewed weekly. Figure 9 gives an example of the site readiness plots for a recent week, for the Tier 1 and Tier 0 sites. This readiness is measured by the experiment-specific tests, while the reliability above is measured by the generic operational tests.
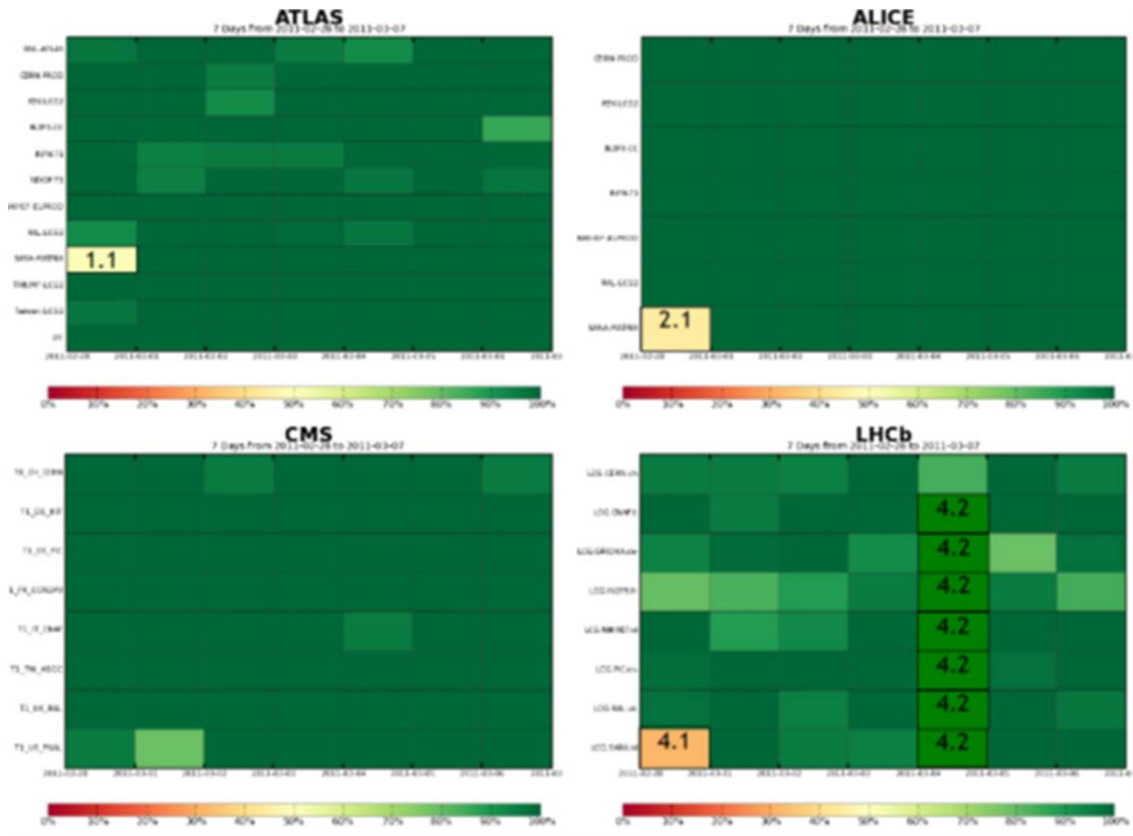


Figure 9: Site readiness for a typical week (March 2011)

All of the availability and reliability reports for all sites can be consulted at: http://lcg.web.cern.ch/LCG/reliability.htm.

## 3. Applications Area

### ROOT

The production release 5.28 at the end of the year went smoothly without any major problems. Large effort has been made to consolidate the existing code and improve its quality by using the coverity tool, better testing infrastructure and additional tests for the GUI subsystem using the event recorder. For the I/O subsystem, major progress has been made to further optimize the reading part of the streaming engine. The new interpreter prototype based on clag/LLVM is clearly progressing: already at this early stage it can replace certain parts of ROOT's current interpreter CINT. Additional improvements have been made in the core mathematical libraries such as an updated interface for numerical integration, minimization and distribution sampling. New algorithms like genetic minimization or kernel density estimation have been added, as well as various bug fixes and improvements have been applied also in Roofit/Roostats packages. Initial investigations have begun

(and prototype) on new GUI technologies based on OpenGL (in particular OpenGL ES) on different platforms.

Significant effort was invested to support PROOF for ALICE during the HI run. The ProofBench suite that allows benchmarking and understanding PROOF performance at any cluster was introduced.

### Persistency Framework

New releases of the persistency projects have been prepared for the two new configurations LCG_59b (for ATLAS, based on ROOT 5.26) and LCG_60 (for LHCb, based on ROOT 5.28), using the same code base for both. The new releases include several bug fixes and enhancements in all three packages, CORAL, POOL, and COOL.

### Simulation

The new public release of Geant4, release 9.4, was made in December as scheduled. Among the features included are: improvements in the Fritiof/FTF hadronic model (better selection of final states at low energies and tuning of the parameters of its Reggeon cascade); a revised choice of modelling in FTFP_BERT to use improved modelling of hyperons and anti-nucleons (adopting the CHIPS model in place of LEP); a new geometrical solid, G4GenericTrap (following an ALICE request); and the first implementation of a new build/installation environment. Validation tests of the new release have been carried on the grid, showing excellent stability. The new release has been chosen by both ATLAS and CMS as the basis of their 2011 simulation productions.

On physics validation, the investigation on the effects of hadronic models transition on the energy resolution has been completed; it has revealed that both resolution and the normalized width do not show the problem, and in particular, the FTFP_BERT physics list is very smooth. "Shower momenta" are now routinely checked with the Simplified-Calorimeter test suite. The results are now available through a web-based application.

The new web service for MC Generators' tuning and validation is now publicly available at http://mcplots.cern.ch.  Further extensions and improvements are planned. Most MC Generator packages have now been ported to MacOSX, and old ones can be added on demand.

### Software Infrastructure

The SPI project has been focusing on consolidation of infrastructure and services. All the web services have been successfully migrated to newer hosts in the CERN computer centre, and in turn integrated into the central service monitoring provided by CERN IT. The AA project websites are still being moved to the central Drupal services of PH/SFT, and an effort to systematically update the documentation of the SPI infrastructure has been started.

There have been two further releases in the "LCG 59" cycle and the first release of the "LCG 60" cycle. The most important change in the "LCG 60" cycle is the move to ROOT 5.28.00 and Boost 1.44. Various other externals were updated as well. This release series is the first one to fully support the Intel icc compiler on Linux. At the same time VC7 was replaced by VC9 on the Windows platform. The Scientific Linux 4 support has now officially been discontinued.

A new prototype project is the usage of CMake as a common build tool for the Applications Area projects, being a potential replacement for CMT and manually managed Makefiles. First attempts have been promising and the next milestone will be to provide the externals of the LCG configuration in form of a CMake environment.

## 4. Level-1 Milestones

A full report on milestones and progress can be found on the WLCG web at http://lcg.web.cern.ch/LCG/milestones.htm. Several of these have been mentioned in sections above.
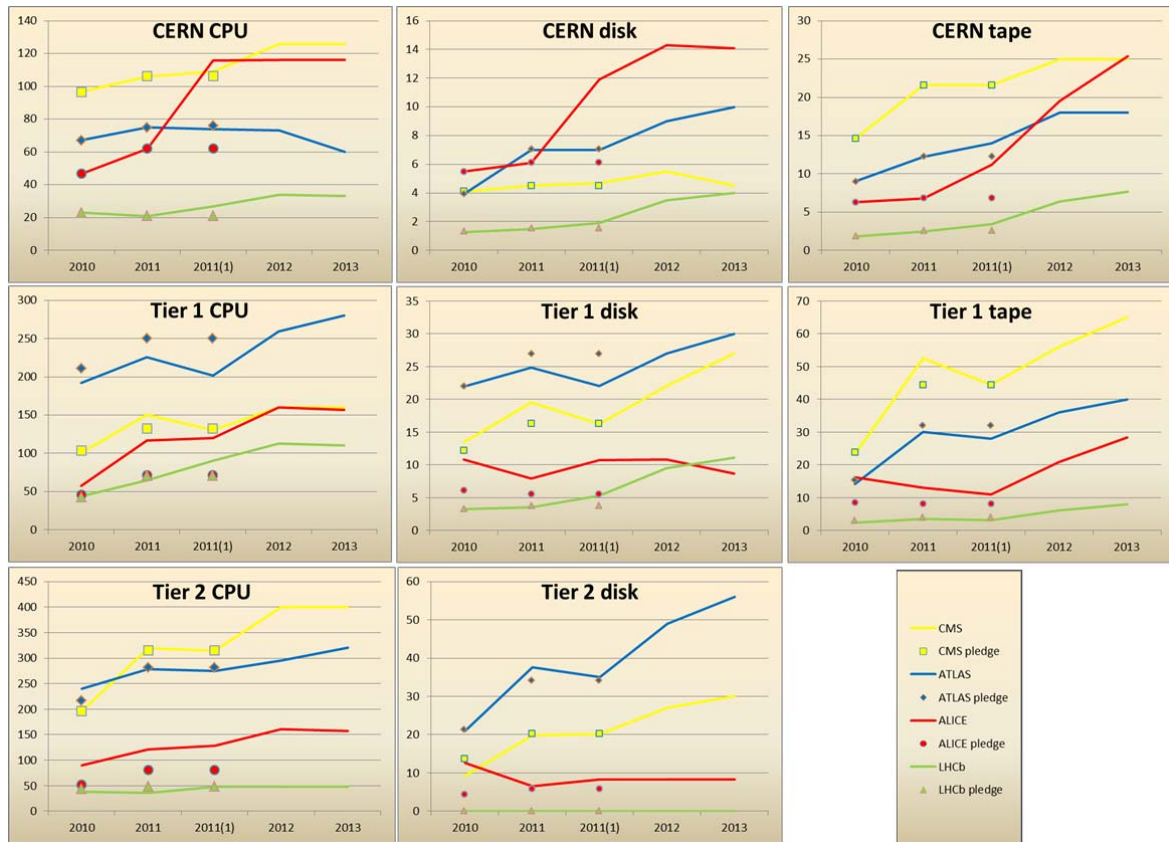
Progress on specific milestones includes:

- Support for multi-user pilot jobs. Deployment continues, with all of the Tier 1 sites now having installed glexec and either SCAS or ARGUS, or GUMS for some US sites. A significant number of Tier 2 sites have now also installed the software and are available to be used by the experiments. As reported previously there has been little reason for the experiments to change their software to use glexec, particularly in the first few months of data taking. In the February MB this issue has again been brought forward, and the agreement is to push the deployment again by using the test probes to check the correct installation at sites and to follow up problems by opening GGUS tickets to the sites. In parallel, the problems reported by ATLAS in using glexec will be followed up. Other experiments are ready to use the facility where available.

- CREAM CE deployment. Cream is now deployed widely and all experiments are satisfied with the performance at the moment. Many sites still run LCG-CEs although this is now no longer necessary. The SAM system still has to complete the update of the testing so that Cream is treated in the same way as the LCG-CE as an alternative CE. This is due to be completed in February. At that point the LCG-CE will no longer be recommended for deployment, and an end-of-life date will be discussed.

- Data Management prototypes. In January the work on the data management prototypes was presented and discussed. Of the original 14 suggested prototypes, around 10 of them are being actively investigated or followed by one or more experiments. It was agreed that the process that had been started in Amsterdam had successfully concluded, and that further work on these tools and services would become part of the work plans of the experiments and software developers. Future work on these will continue to be reported in GDB meetings.

- Automated gathering of installed capacity data. Most sites Tier 1 and Tier 2 now correctly report their overall capacities. However, very few yet correctly publish by VO-share which is what is required in order to correctly report on capacities. This has to be followed up site by site to ensure valid data publication; the Tier 1s have the responsibility to follow up with their Tier 2s.

- Updates to SAM/Nagios to provide more flexible reporting:
    - Validate dashboard applications with Nagios tests (IT/ES); New interface by myEGI is available since mid-January; pre-production service can be used for migration
    - Stop old SAM system as soon as green light from the experiments
    - New ACE availability calculation mechanism:
        - December 2010: Validated the standard availability for OPS - done
        - January: Computation of standard availabilities for LHC experiments (one profile per VO) - done
        - February: Multiple availabilities (different profiles, same algorithm) per VO - done
        - March: Multiple availabilities (different profiles and algorithms: CREAM CE use case) per VO; almost complete.

## 5. Planning and Evolution

### Resource request evolution

The Figure below shows a summary of the evolution of the experiment resource needs for 2011, 2012, and 2013. While it is understood that the 2011 pledged resources are largely already installed based on the requests made in 2010, the experiments have nevertheless re-estimated that need in the light of the 2010 LHC running experience. For 2012 and 2013 the new requirements are based on the updated LHC schedule with a full year of running in 2012 and shutdown in 2013.



**Figure 10: Summary evolution of the experiment resource requests. The 1st 2011 point represents the request for 2011 made in 2010; while the 2nd 2011 point represents the currently understood needs. The solid lines indicate the requirements, the markers represent the pledges.**

The mechanisms of LHC luminosity increase have resulted in much higher pile-up values for all experiments than had been foreseen in the early years of running. This has resulted in much larger event sizes and reconstruction times than anticipated, as well as increasing simulation reconstruction times. In large part this is the reason behind much of the increase in resource needs together with the anticipated data volumes reached in 2012. Increases in 2013 over 2012 are driven by the analysis workloads.

For LHCb, as well as the effects of increased pile-up, the intention is to add an additional 1 kHz of trigger rate for charm physics. This intent is supported by the LHCC (March meeting), but increases the resource needs in 2011 over those planned, with consequent increases in the following years.

The report of the Computing Resources Scrutiny Group, to be presented at the April C-RRB, will analyse these requests in more detail.

## Tier 0

As noted in the previous report, and taking into account the change in the LHC schedule in 2012 and 2013, the need for significant additional new power for the Tier 0 is now foreseen for 2014. Following the solicitation last year of informal bids for co-hosting data centre computing equipment, by the end of November some 28 had been received. A number of visits and discussions to follow up on these are on-going. No decision on further formal steps has yet been made.

## Network Evolution

In June last year a workshop was held to consider how to improve reliability and performance of data access, particularly at Tier 2 and Tier 3 sites. One of the outcomes of that workshop was the realisation that the computing models of the experiments should evolve to make more optimised use of the existing networking infrastructure, but also that it will be equally important to ensure that Tier 2s at least are provisioned sufficiently, and that the overall network architecture will permit data to move between (almost) any set of sites, rather than in the hierarchical manner envisioned so far. The LHCOPN group was tasked with understanding the realistic requirements of the experiments over the next several years and proposing a network architecture to accommodate this. A first proposal has been made, which is currently under discussion in the national and international research and education network organisations. This proposal, named LHCONE (LHC Open Network Environment), builds on the concept of open exchange points. Exchange points will be built in carrier-neutral facilities so that any site can connect with their own fibre or using circuits provided by any telecom provider. LHCONE should enable Tier2s and Tier3s to obtain their data from any Tier1 or other Tier2, and it should shield the general research and education IP infrastructure from LHC traffic.

In addition to soliciting comments on this approach, early limited scale prototypes are being built, and work is needed on the governance and operations model. Appropriate monitoring must also be provided from the outset.

Additional details can be seen at http://www.lhcone.net.