**Worldwide LHC Computing Grid Project
Project Status Report
Resource Review Board – 15th April 2008**

This status report covers the period from September 2007 – February 2008. 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 six month period there have been several significant activities including the preparations for the experiment dress rehearsals and the proposal and execution of the first phase of a Combined Computing Readiness Challenge (CCRC'08). The realistic deployment plan of the SRM v2.2 interface to the Mass Storage Systems was agreed at the end of September, and the deployment and configuration of the required SRM services at all of the Tier 1 sites was achieved according to this plan by the end of December. This was successfully tested in the first phase of CCRC'08 and alleviates one of the major outstanding concerns of the project and the LHCC review last November. In addition, by February, most Tier 2 sites had also already upgraded. As a result of considerable effort by all teams involved, all major problems were resolved during this period. A few issues are still unresolved and will be a priority for the next few months to find and agree fixes or workarounds that can be in place for the May challenge and first data taking. There is still work required to tune the site configurations of the storage systems as the experiments' needs are better defined in the light of experience.

A WLCG Collaboration workshop was held at the time of the CHEP conference in September and was attended by 160 people. This workshop assessed progress in terms of the service and from the experiments' points of view. The need for a complete test of the system at close to the full 2008 data rates with all experiments taking exercising their complete computing models was first realised here and the CCRC'08 was proposed.

In the last quarter of 2007 upgrades for several middleware components were provided, including an FTS version to manage the SRM v2.2 storage, the requested bulk methods in LFC and DPM and a general move to SL4 versions of most of the software.

The Combined Computing Readiness Challenge was designed to bring together all four experiments and to exercise the full computing models from data acquisition through to data analysis at the Tier 2's. It was agreed to be run in two phases in February and in May. The February phase would test components of the system and be limited by the available resources, and the second phase in May will be with a full dress rehearsal at the full 2008 data rates for all experiments with the full 2008 resources in place. Here we report on progress in the first phase and preparation for the second.

The SRM v2.2 mass storage system deployment at Tier 1 sites had been noted as delayed, but was achieved before the start of CCRC'08, and during the challenge showed relatively few problems. In fact in total ~160 sites (Tier 1 + Tier 2) had an SRM v2.2 storage system in production. There is a short list of SRM issues that were highlighted during the challenge, of which only 2 are regarded as high priority. These will be addressed by each of the implementations in the coming months.
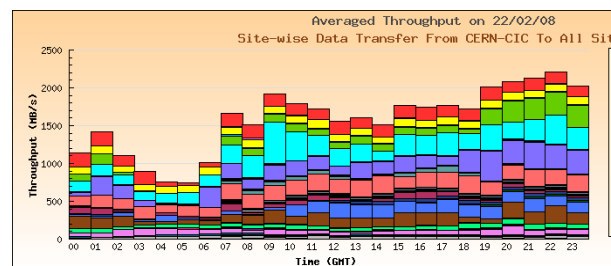


**Figure 1: Data distribution CERN to Tier 1s**

In terms of data transfer, several significant goals were achieved. The total rates transferred out of CERN to the Tier 1 sites were significantly greater than those previously achieved in earlier tes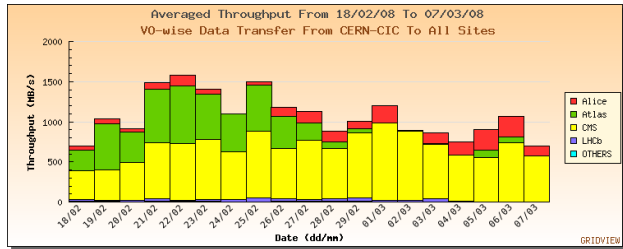ts, and have been sustained over several weeks. All 4 experiments have demonstrated sustained rates in excess of their requirements



**Figure 2: Data transfer by experiment**

for 2008 running. Rates of greater than 2.1 GB/s were achieved in aggregate between all experiments from CERN to all 11 Tier 1 sites. This is shown in the 2 figures above: Figure 1 shows the distribution to sites on 1 day, Figure 2 by experiment over several weeks. As can be seen the experiments have found the testing sufficiently useful that they are continuing.

The performance of the Castor 2 system at CERN had also been of concern, but was demonstrated to perform reliably at rates well in excess of those needed for data taking. CMS in particular were able to demonstrate aggregate rates in and out of Castor of 3-4 GB/s (see Figure 3), and sustained rates to tape of 1.3 GB/s. Unfortunately this level of use with several experiments together was not



**Figure 3: Data rates in/out of Castor2**

demonstrated since ATLAS was later in starting the challenge. In total during the 1 month challenge CMS moved >4.5 PB of data between all participating sites. All of their Tier 1s achieved the targets to receive data from CERN and migrate to tape, and a large fraction of the T1-T1 and T1-T2 targets were also achieved.
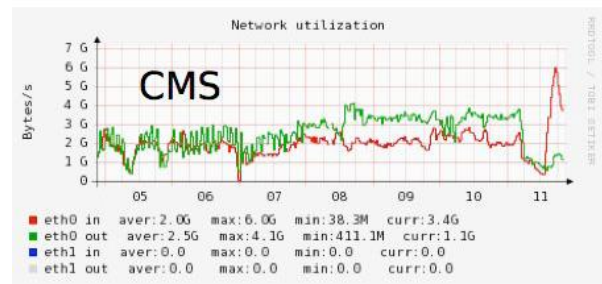
ATLAS started late as the amount of data generated in their Full Dress Rehearsal was rather less than expected. However, using simulated data starting from week 3 rapidly showed the rates mentioned above. They also validated the use of SRM v2 and the Tier 1 storage system setups. They achieved most of their milestones despite the early problems and external dependencies.

ALICE and LHCb also achieved their data rate targets with sustained rates of 80 MB/s and 70 MB/s respectively over several weeks. LHCb tested bulk file deletion with SRM v2. They have tested most of their full computing model, despite the new version of Dirac being available only just before the start of the test.

In summary the February exercise has been a success, with relatively few issues being shown. Some problems of communication – e.g. slow reporting of problems outside of working hours – show that although processes were in place they were not well advertised or used. These points, together with a prioritised list of issues in the storage systems and other middleware services will be addressed for the May challenge. All 4 experiments expressed the desire to keep running at this level from now on. It is important that the full 2008 resources are in place at the Tier 1 sites in time for the May phase so that the complete system can be tested at the full 2008 rates.

## 2. Applications Area

During this period the main activity in the Applications Area has been working towards to the release of the end of the year production versions of Geant4 and ROOT. Particular attention was paid to validation of these releases not only with the standard test suites but also by the experiments themselves since it is likely that the software in this release cycle will be the one used for LHC start-up. The nightly build system was essential to this validation. A new procedure based on the nightly

builds has been put in place to reduce the time needed to deliver validated software releases to experiments.  This build system is being adopted also by Geant4 and LHCb.

The applications area has seen a reduction in staffing with additional reduction anticipated in the next few months due to some staff leaving earlier than anticipated in the staffing plan.  The net result is that some of the activities have been temporarily suspended, mainly in the Physics Validation area. Additional activities will also be affected, and input from the experiments will be requested to prioritize the work and make the best use of the resources.

**Core Libraries.**  The ROOT project has focused on the quality assurance procedure for the new production release 5.18 delivered in January. The QA procedure includes a significant number of tests and validations.  This version of ROOT includes several new packages and consolidation of existing packages.

**Simulation and Validation.**  The 12th Geant4 Workshop took place in September in Hebden Bridge (UK) hosted by the University of Manchester and saw the participation of roughly 80 people; it also included 2 days of user presentations and round tables.

Geant4 version 9.1 was released in December, as planned. It provides a number of fixes and several new features.  Efforts have been undertaken to facilitate the LHC experiments moving to newer Geant4 releases.  Pre-release versions and intermediate development versions were provided to and tested by experiments, providing valuable feedback. Robustness testing was extended with additional, longer testing, enabling the identification and fixing of a number of software issues. Convergence is being sought on using a single recent Geant4 version in production during an agreed period, to enable the concentration of the available effort for the support, maintenance and the provision of fixes.

## 3.  Site Reliability

The results of the site availability metric for CERN and the Tier 1 sites for this period are shown in Table 1.  Full data for each site, including the individual service availability data, is available from the LCG Planning Page.   The site target level was 91% until November, and then 93%.  The project target for the eight best sites was 93% until November and then 95%.  This project target was achieved for all months in the period.  The evolution of the reliabilities for the Tier 1 sites and CERN is shown in Figure 4 and shows a continued general overall improvement.  In particular in the second half of 2007 the stability of most sites had greatly improved.  Unfortunately in February 2008, several sites had a number of problems during CCRC'08 (extended power outages, etc.), although as noted above the overall WLCG service was not affected by these and recovered.

Data on reliability for Tier 2 sites has also been determined regularly and published since October.  However, as yet, not all Tier 2s are publishing this data.  In particular, the US Tier 2 sites rely on Open Science Grid to provide the actual tests that will publish results into SAM.  These tests are not yet in production.  A small group was mandated by the Management Board to assess the equivalence of the OSG proposed tests to those used at EGEE sites.  At the moment only tests for the Compute Elements are defined in OSG and these are not yet running regularly.  A similar situation exists with the Nordic Tier 2 sites.  A similar effort was made to approve the equivalence of tests to run at NDGF, and these are now in production at NDGF – the Tier 1.  However, the Nordic Tier 2 sites are not yet running these tests and publishing results.

| (Colour Schema: Green > Target, Orange > 90% Target, Red < 90% Target) | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Average of the 8 best sites (not always the same 8)** | | | | | | | |
| Jul 93% | Aug 94% | Sept 93% | Oct 93% | Nov 95 | Dec 95 | Jan 95 | Feb 96% |
| **Average of all Tier-0 and Tier-1 sites** | | | | | | | |
| Jul 89% | Aug 88% | Sept 89% | Oct 86% | Nov 92% | Dec 87% | Jan 89% | Feb 84% |

**Detailed Monthly Site Reliability**

| Site | Jul 07 | Aug 07 | Sep 07 | Oct 07 | Nov 07 | Dec 07 | Jan 08 | Feb 08 |
|---|---|---|---|---|---|---|---|---|
| CA-TRIUMF | 97 | 97 | 95 | 91 | 94 | 96 | 97 | 95 |
| CERN | 95 | 99 | 100 | 100 | 98 | 100 | 99 | 97 |
| DE-KIT (FZK) | 75 | 67 | 91 | 76 | 85 | 90 | 94 | 98 |
| ES-PIC | 96 | 94 | 93 | 96 | 95 | 96 | 93 | 99 |
| FR-CCIN2P3 | 94 | 95 | 70 | 90 | 84 | 92 | 95 | 98 |
| IT-INFN-CNAF | 82 | 70 | 80 | 97 | 91 | 96 | 70 | 20 |
| NDGF | n/a | n/a | n/a | 89 | 98 | 100 | 92 | 84 |
| NL-T1(NIKHEF) | 92 | 86 | 92 | 89 | 94 | 50 | 57 | 84 |
| TW-ASGC | 83 | 83 | 93 | 51 | 94 | 99 | 97 | 100 |
| UK-T1-RAL | 98 | 99 | 90 | 95 | 93 | 91 | 92 | 93 |
| US-FNAL-CMS | 92 | 99 | 89 | 75 | 79 | 88 | 93 | 85 |
| US-T1-BNL | 75 | 71 | 91 | 89 | 93 | 44 * | 91 | 63 |
| Target | 91 | 91 | 91 | 91 | 91 | 93 | 93 | 93 |
| **Above Target (+ > 90% Target)** | 7 + 2 | 6 + 2 | 7 + 2 | 5 + 4 | 9 +2 | 6 +4 | 7 +3 | 7 +3 |

(*) The reliability for BNL in Dec 2007 is incorrect because of a mis-configuration of the SAM setup at the site.

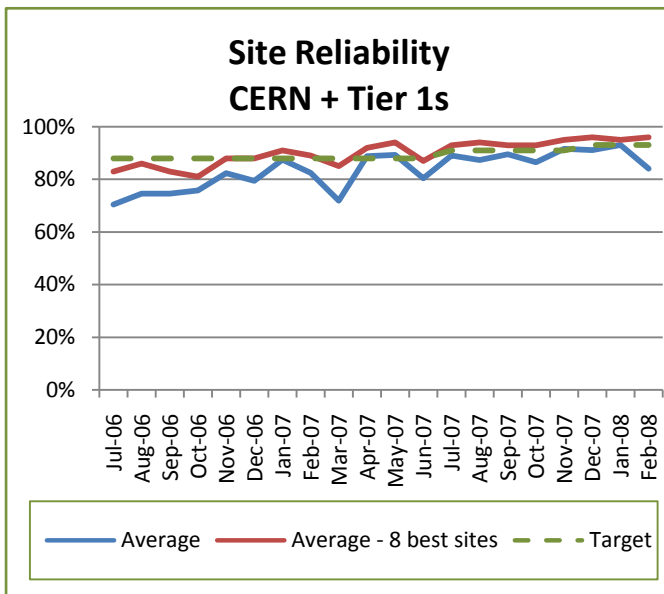Table 1: Site Reliability Summary - July 2007 - Feb 2008

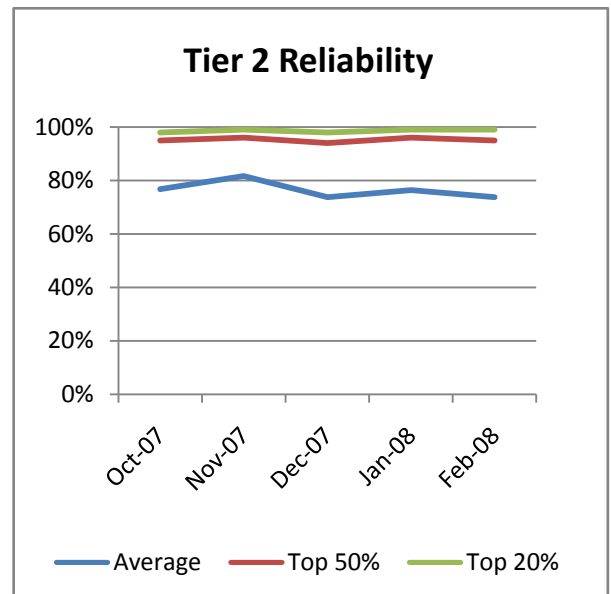Figure 4: Tier 1 Reliability - July 2006 - Feb 2008



Figure 5: Tier 2 Reliability - October 2007 - Feb 2008

Figure 5 shows the situation with reliability of the Tier 2 federations and sites, with 100 sites currently publishing results. Of note is that more than 50% of the Tier 2 sites show very high levels of reliability, which is important as the top 50% of sites provide the majority of resources. It is usually the complexity of the storage systems that causes problems at the Tier 1 sites that affect their reliabilities.

It is important that the Nordic and US Tier 2 sites now rapidly reach the position where they are also publishing the reliability metrics.

## 4. Level-1 Milestones

Most of the previously defined high level milestones have now been achieved. Very few additional milestones are now being added, and there is a general move within the project towards metrics as monitors of the production system performance. However, there are some milestones that are worthy of comment, particularly as several of these are related to ensuring the overall service:

- **24x7 support:** Most Tier 1 sites now have mechanisms defined, tested and in operation for providing support for out-of-hours problem resolution. Three Tier 1s are still to fully finish this milestone, but anticipate this in April, in advance of the May phase of CCRC'08.
- **VOBox SLAs:** Here the progress is still slow. All but 3 Tier 1s have defined an SLA for supporting VOBoxes, but they all anticipate having this in March. Implementations of the SLAs are missing for 6 of the Tier 1s, but these all anticipate this to be achieved in March or April. Acceptance by the experiments can only come once the implementation is done, but in most cases the experiments sign off as the sites define the SLA, with only a subsequent verification that what has been implemented matches what was proposed.
- **Procurement of resources:** see next section.
- **VO-Specific SAM tests:** This milestone was designed to provide VO-specific tests to measure site reliability to be complementary to the standard tests. Although most experiments do use this facility and do run tests, the results are not yet regularly published as the SAM system needed some adaptation to correctly determine the real availability based on this

different set of critical tests.  The adaptations are done, but verification and validation by the experiments, sites, and Management Board are still to be completed.

- **SAM testing for OSG**: This is new and was added in this quarter, and is scheduled to be in place and published by the end of March.
- **Tape efficiency metrics:**  This is a new milestone, which requires the Tier 1 sites to publish a set of metrics that demonstrate that the performance of the MSS systems, particularly tape performance, is adequate.  It is intended that such metrics are published by the Tier 1s for the May phase of CCRC'08.

## 5. Resource procurement

The installation and set up of resources according to the 2008 pledges has proceeded relatively well.  The commitment was to have these resources in place by April and the ramp up from the middle of last year was significant in most cases.   With the second phase of CCRC'08 planned for May it is important that the majority of resources are really in place and available.

In terms of CPU most of the Tier 1 sites will have their full 2008 pledges in place for May 1.  The largest missing contribution is that of the Netherlands which is only expected in November due to problems in the procurement process.  For disk, the 2008 requirement is 23 PB of which 15.5 are expected by May 1, and for tape the requirement is 24 PB of which 15 PB is expected by May.  In the storage area the capacities will catch up later in the year as the need expands.  These levels will be sufficient for the anticipated needs of the May run of CCRC'08.

Several sites reported delays or constraints in their procurement processes that meant the process took longer than anticipated or that equipment was not delivered according to schedule, or was delivered and was not acceptable.  It is vital that in future years, these eventualities are taken into account in the planning and procurement process and allowance for delays, the need to switch vendors, etc. be made from the outset as in those years the resources must be in place for the start of data taking.

## 6. Long term evolution of requirements and pledges

In the previous report it was noted that the long term requirements of the experiments is not fulfilled by the current levels of planned resources for future years.  This situation has not changed.  The Computing Resource Scrutiny Group (C-RSG) has been formed and the mandate agreed, and Domenec Espriu (Spain) has been nominated as the chair.  The group had its first meeting on 20[th] March.  It is intended that the group will provide a level of confidence in the resource requests from the experiments.

In recent years the concerns over power and cooling have become important issues facing many HEP Computer Centres, and several of the WLCG Tier 1s are actively planning or building extensions to their power and cooling infrastructure in order to be able to install the capacity pledged in the MoU.  These concerns are also true for the CERN Computer Centre, and with the current planning for capacity ramp up, the power available in the present centre will only cover the needs of the Tier 0 and CAF until early 2010.  In the present situation the computing capacity that can be provided is limited by the current envelope of 2.5 MW for equipment (plus power for cooling and distribution).  However, of this, 350 kW is "critical" power (backed up by diesel generators).  These critical services include some of the physics database services.  Thus the limit for the remaining physics services is ~2.1 MW.

At the time of planning for the Tier 0 the best industry and technology predictions were that PC power would remain flat at around 100 W per box (a dual processor).  However in the last several

years it became very clear that PC power needs actually scale with the CPU capacity, and some 2 years ago it became clear that the semiconductor industry had no solution.

The load today is around 1.7 MW, with an additional 400 kW anticipated in the next month or so as the 2008 capacity is installed.  Thus this is already reaching the 2.1 MW limit.  Aggressive removal of older equipment will probably allow the installation of the 2009 capacity, as long as the needs for critical power (e.g. physics databases etc.) remain within the 350 kW.  Providing the 2010 capacity will not be feasible with these constraints.

Some indicative early planning shows that:

- The estimated time needed to provide a new or refurbished building to provide 2.5 MW initially and growing to 5 MW ranges from about 27 to more than 40 months;
- External hosting of services is an option that could cover some short term needs, but is expensive (~ 3.6 MCHF/MW/year)
- It is unlikely that the Tier 1 sites could absorb additional Tier 0 capacity as many are in a similar situation regarding power and cooling.

An additional point that must be noted is that the projected increased needs of the experiments in computing capacity after 2009 assume a 30%/year growth.   This is significantly different from the experience over the past 15 years where a 100%/year growth has been typical.

This issue of power and the ability existing Computer Centre infrastructures to provide the computing capacity required for the LHC experiments is of utmost importance, and needs to be addressed with an aggressive and realistic plan.