# Getting prepared for the LHC Run2: the PIC Tier-1 case

**J Flix**[1,2]**, E Acción**[1,3]**, V Acín**[1,3]**, C Acosta**[1,3]**, J Casals**[1,2]**, M Caubet**[1,2]**, R Cruz**[1,3]**, M Delfino**[1,4]**, F López**[1,2]**, A Pacheco**[1,3]**, A Pérez-Calero Yzquierdo**[1,2]**, E Planas**[1,3]**, M Porto**[1,2]**, B Rodríguez**[1,3]**, and A Sedov**[1,3]

[1] Port d'Informació Científica (PIC), Universitat Autònoma de Barcelona, Bellaterra (Barcelona), Spain
[2] Also Centro de Investigaciones Energéticas, Medioambientales y Tecnológicas (CIEMAT), Madrid, Spain
[3] Also at Institut de Física d'Altes Energies, IFAE, Edifici Cn, Universitat Autònoma de Barcelona, Bellaterra (Barcelona), Spain
[4] Also at Universitat Autònoma de Barcelona, Department of Physics, Bellaterra (Barcelona), Spain

E-mail: jflix@pic.es

**Abstract.** The Large Hadron Collider (LHC) experiments will collect unprecedented data volumes in the next Physics run, with high pile-up collisions resulting in events that require a complex processing. Hence, the collaborations have been required to update their Computing Models to optimize the use of the available resources and control the growth of resources, in the midst of widespread funding restrictions, without penalizing any of the Physics objectives. The changes in computing for Run2 represent significant efforts for the collaborations, as well as significant repercussions on how the WLCG sites are built and operated.

This paper focuses on these changes, and how they have been implemented and integrated in the Spanish WLCG Tier-1 centre at Port d'Informació Científica (PIC), which serves the ATLAS, CMS and LHCb experiments. The approach to adapt a multi-VO site to the new requirements, while maintaining top reliability levels for all the experiments, is as well presented. Additionally, a description of work done to reduce the operational and maintenance costs of the Spanish Tier-1 centre, in agreement with the expectations from WLCG, is provided.

## 1. Introduction

The LHC, at the European Laboratory for Particle Physics (CERN, Switzerland), started operating in November 2009. The successful first run (Run1) ended in February 2013, and the accelerator entered into a period of shutdown (LS1) for maintenance and upgrade works. The LHC is expected to start producing physics collisions by June 2015, with a performance at almost the design level. The LHC experiments components were revised and upgraded to prepare for the new period, in which the beams will collide at almost double the energy than previously recorded, and with reduced time between bunch crossings, down to 25 $ns$.

To analyze the unprecedented rate of PetaBytes (PB) of data per year generated by the LHC, a Grid-based computer network infrastructure was built. The largest scientific distributed computing infrastructure in the world adds up the computing resources of more than 170 centres in 34 countries to form the Worldwide LHC Computing Grid (WLCG [1][2]). In Run1, the

infrastructure proved to be a key component for prompt analysis and reconstruction of the LHC data. During LS1, all of the Run1 data was reprocessed and simulations with the new LHC conditions were produced. At the dawn of Run2, the LHC data volume sums up about 300 PB of raw, simulated and processed data, from all of its detectors.

The computing centres are functionally classified in Tiers in WLCG. Eleven of these centres are the so-called Tier-1s, receiving a fraction of the raw data in real time from the Tier-0 at CERN, and in charge of massive data processing, storage and distribution. Spain contributes with a Tier-1 centre: Port d'Informació Científica (PIC), located in the campus of the Universitat Autònoma de Barcelona, near the city of Barcelona. PIC provides services to three of the LHC experiments, ATLAS, CMS and LHCb. It accounts for 5.1% of the total Tier-1 resources of ATLAS and CMS, and 6.5% for LHCb, acting as the reference Tier-1 for the Tier-2 centres in Spain and Portugal, and sites located in Valparaiso (Chile) and Marseille (France).

The next period of LHC running will be producing more data, composed by more complex events, in comparison to Run1. During LS1, the computing models of the experiments underwent a series of revision, in order to cope with the high data volume expected for the second LHC run, while keeping a controlled growth of computing resources elsewhere.

## 2. Computing Upgrades during LS1

The experiments have revised their computing models during the LS1. The goal is to keep a controlled growth of resources for the next years, in a *flat-funding* model [3], i.e. assuming that regions and sites are operated without a substantial budget increase, and benefiting from technology evolution trends in order to provide the needed growth of resources. This had many implications on the way the experiments operate, and the computing models were modified in many aspects: the available resources are exploited in a more flexible manner, and experiment tools have been improved, decreasing the operational efforts while being ready to confront the next LHC run needs. The disk space is better managed: the most popular datasets are kept on disks at the sites, and cleaning mechanisms of non-popular files from disks are present, or being deployed, reducing the disk needs at the sites.

New access protocols have been introduced allowing remote data access. This naturally enables easy failback mechanisms when access to a particular file locally at a site fails. It also allows running activities out of sites, busy in terms of computing, and hence allowing sites to process data without downloading the datasets locally. Many common services used by the experiments have been simplified, with less service instances and more easily manageable deployments and operational procedures. FTS3 is a good example of a service of this type, provided for all of the WLCG community by only a few sites. HLT farms, opportunistic resources, Tier-0 resources, and some other sites, are being exploited regularly by means of Cloud Computing techniques, complementing the resources needs from the experiments.

## 3. PIC Tier-1 compliance with Run2 requirements

PIC is an active and successful participant in the WLCG project since its start, and it has shown its readiness for the LHC data taking periods. It has contributed to prototyping and testing of the Grid middle-ware and services that were being developed, and successfully participated in the Service Challenges carried out by the experiments, testing campaigns aimed to progressively ramp-up the level of load on the infrastructure under conditions as realistic as possible, achieving breakthrough performances. PIC showed good performance results during Run1, and during LS1 it worked to be fully compliant with the Run2 requirements from the experiments, tendered and provisioned resources for data taking re-start, and prepared all of its infrastructure for the resources growth that are expected in the next years. Figure 1 shows the PIC resources growth and usage during Run1, and the expected growth trends for Run2. PIC resources are expected to double during Run2.
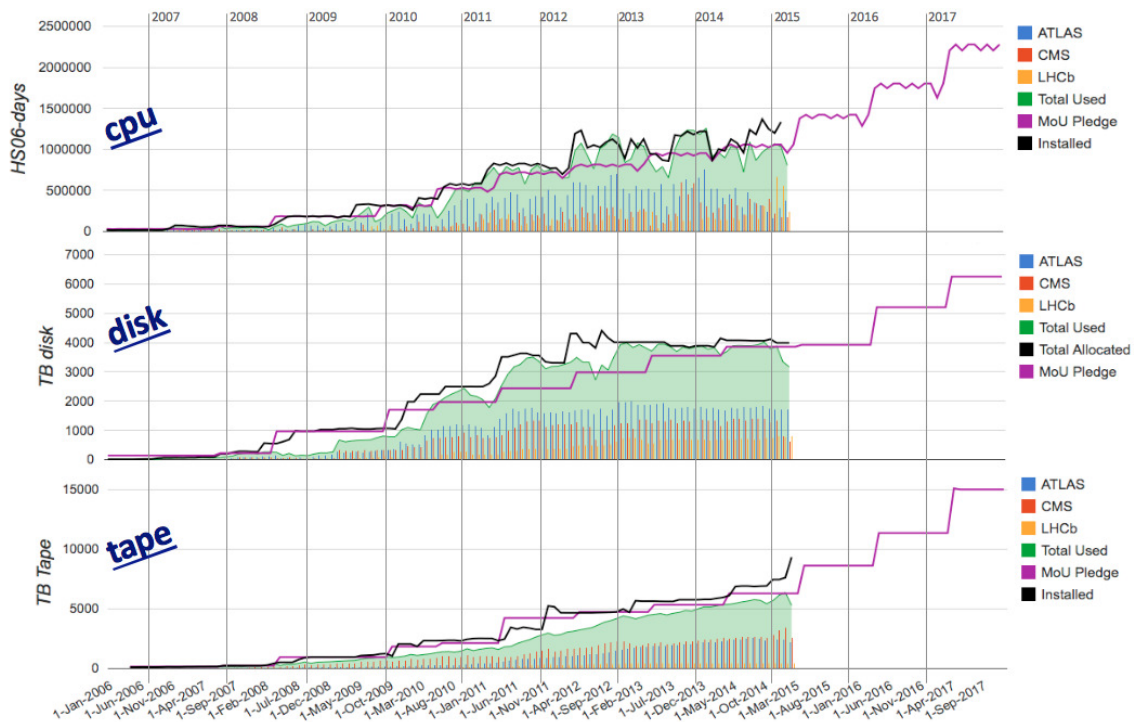
**Figure 1.** CPU (top), disk (center) and tape (bottom) resources installed and used at PIC, since 2008. The expected evolution of the resources growth during LHC Run2 is as well displayed.

As of today, the computing resources installed in PIC comprise around 5500 cores managed by Torque/Maui batch system and scheduler. This corresponds to about 70000 HEP-SPEC06 (HS06, see Ref. [4]). The servers are typically two quad-core x86 CPUs, with at least 2GB RAM per core (recent E5-2640-v3 purchases have 4GB/core). Each of these nodes typically has two 10Gbps Ethernet interfaces which are then aggregated in switches and connected to the storage infrastructure. The main servers consist of Blades (HP) and Dual-Twins (Dell).

The storage service at PIC is managed by dCache [5] and Enstore [6] softwares. The dCache software provides uniform and efficient access to the disk space provided by many file servers, and talks to the Enstore software to interface to magnetic tape storage. As of today, 6 PB of disk space is installed, by means of around 2000 hard disks of 2, 3 and 6 TB, distributed on around 50 servers based on x86 technology, each connected by one or two 10Gbps Ethernet, depending on the hardware. The servers brands comprises DataDirect, SGI, and SuperMicro.

The current tape infrastructure at PIC is provisioned through a Sun StorageTek 8500SL library, providing around 6650 tape slots which are expected to cover the PIC tape needs in the coming years. Enstore manages 12 PB of tape storage, with access to a total of 4.3 million files. The supported technologies are LTO-4, LTO-5 and T10KC, containing 29%, 10% and 61% of the total data respectively, in around 6000 tape cartridges. A total of 24 tape units are installed to read/write the data (12 LTO-4, 4 LTO-5 and 8 T10KC). Aggregated read/write rate has achieved hourly average rates peaking at 3.5GB/s.

### 3.1. Data management upgrades

With better and increased network capabilities among centres, the Tier-1s can naturally become major data servers to the whole Grid. New protocols were deployed in LS1 to allow for remote data access, namely XRootD [7] and standard HTTP/WebDAV.

3

XRootD failback is enabled in all of the PIC computing nodes. This means that if PIC storage is busy or a file of a dataset is missing or not accessible, running jobs can read those files from remote centres. Thus, entire workflows can be executed in PIC reading remote datasets, preventing the need for data pre-placement to PIC. ATLAS and CMS are increasingly using these functionalities. Additionally, ATLAS and CMS disk-only data can be XRootD-accessed from remote centres, a worth of ∼4 PB of data which is accessed by the majority of Tier-1 and Tier-2 sites. LHCb data can be as well HTTP-accessed from remote centres (∼800 TB).

Joining the data federations required substantial R&D and tuning. First of all, it needed a deployment of compatible data management software (i.e. a compatible dCache version). Tape systems need to be protected against uncontrolled accesses, hence disk-pool areas were created and populated. This naturally allowed end user analysis for CMS collaboration to be enabled at the centre. Dedicated experiment monitoring plugins were installed, as well as site 'local' XRootD redirectors (which translate Logical File Names to Physical File Names, and interact with the storage for file discovery).

Exposing the majority of the data to remote centres forced PIC to implement dedicated monitoring and establish protection mechanisms to control/limit the access of the resources, if necessary. This new way of accessing data is being commissioned and it is expected to gain in popularity (hence usage) once the LHC Run2 starts. Since its deployment at PIC, the XRootD exports are at 15%-20% levels of the total data export traffic.

### 3.2. Computing service upgrades

Given the evolution of LHC running conditions at the restart of the data taking in mid 2015, the experiments are developing multi-core applications in order to cope with the analysis of complex events with large pile-up. Many experiment workflows are being migrated to multi-core applications, which will keep the usage of RAM per job during execution at a reasonable level. There are many challenges for sites in this new scenario, and PIC has shown very active in this context co-coordinating the WLCG multi-core deployment task force [8].

The scheduling of both multi-core and single-core jobs at a site should be effective and efficient, as both types of jobs are expected to co-exist from all of the LHC VOs during Run2. Static splitting of resources for both type of jobs should be avoided to maximize site CPU usage. In order to schedule multi-core jobs, n-core slots must be created in the site computing nodes. Effective node *draining* prevents single core jobs taking resources of ending jobs, hence creating the n-core job slot in a node. Using short running jobs while sufficient resources are being reserved to create a multi-core slot might be possible. However, this *backfilling* is not currently practical, as the LHC VOs are not scheduling short jobs with an estimation of the expected duration time being provided at submission time. Draining thus represents wastage at this point, an unavoidable price to be paid. Therefore, once the cost has been paid, batch systems should avoid immediate multi-core slot destruction, to optimize future multi-core jobs allocation and CPU usage in the farms.

Controlled draining and multi-core slot conservation at PIC is achieved with a dynamic partitioning of computing resources. This is achieved by the *mcfloat* tool (as developed by NIKHEF [9]) for Torque/Maui. Controlled ramp up of multi-core resources reduces draining impact on farm utilization. Figure 2 shows a controlled allocation of multi-core jobs in PIC, for ATLAS Tier-1/Tier-2 and CMS Tier-1 jobs. At the end of this ramp-up, half of the slots offered by PIC were used by multi-core jobs. During this ramp up period, the farm occupancy was measured to be at 98%, validating the method for efficient multi-core scheduling in PIC.

Experiment applications are still being adapted to the new multi-core schema, and as of today, some experiments send a fraction of single-core jobs that needs high memory provisioning. PIC has deployed a high-memory queue for ATLAS, which offers 256 slots with 4GB/core, which are used by standard tasks when no ATLAS high-memory jobs are scheduled in PIC.
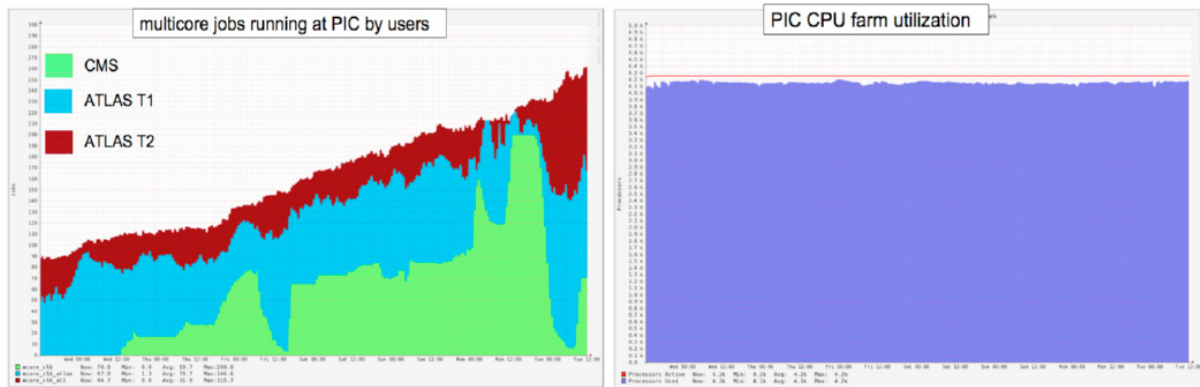
**Figure 2.** Controlled allocation of multi-core jobs in PIC, for ATLAS Tier-1/Tier-2 and CMS Tier-1 jobs, keeping the farm utilization at 98% levels.

*3.3. Network upgrades*
Network access to data allows for a significant cost optimization, as disk is the most expensive resource to be deployed and operated at a site. However, the remote data access scenario has shown to add more load on the network. Therefore, careful planning of network upgrades was required. New switches were acquired, router upgrades were needed, and more powerful firewalls were deployed. There is the tendency in WLCG to consider that network is 'free', but it should be noted that network equipment is typically expensive, and deployment efforts are FTE demanding. These activities, on top of other increasing network demands, might as well result in an increase on site WAN connectivity. PIC is connected to its network provider through a 10Gbps line. This WAN 'last-mile' connection has also a cost, which impacts on the budget of the centre. In PIC, a line of 10Gbps is deployed for LHCOPN [10], 10Gbps for LHCONE, and 2x10Gbps for General IP services. Even if the main 10 Gbps from PIC to our NREN is not yet saturating, WAN bandwidth increase is being drafted with the involved parties.

## 4. Reducing operational and maintenance costs at PIC
In order to pave a secured and efficient road for future site growth, PIC invested efforts in order to reduce maintenance and operational costs at the site. Typical life cycle of equipment at PIC is set to 4 years. This is strictly followed for disk resources, however CPU servers are provisioned with 5 years maintenance support, if possible. Since the beginning of 2013, the PIC computing farm power is adjusted to electricity cost: less CPU power is offered during high cost periods, and vice-versa, without negatively affecting the annual WLCG pledges. The power consumption per unit of performance (Watts/HS06) is more or less constant since a few years, the oldest CPU servers are used to modulate the computing power. The net effect was a reduction of ∼10% in the electricity bill.

At the end of Run1, PIC deployed a RedHat Enterprise Virtualization system (RHEV 3.4.2, KVM-based). The production system is installed over an HPBlade box with seven hypervisors, each of them equipped with 16 cores and 96GB RAM (HP Proliant BL460c) with 2x10GbE ports. The hypervisors are connected via fiber-channel to a NetApp FAS3220 (2TB, Thin Provisioning, with *qcow*2 image formats). Around ∼125 services are run in the system. This reduces the number of physical machines by a factor 10, without any impact on the reliability and performance of the services, and reducing the costs substantially. In order to save licence costs, Ovirt 3.5 is being tested at scale. Four hypervisors, in a similar environment, are used to run up to 60 test services. The new setup is expected to be deployed in production soon.

Constant efforts are made to simplify and improve configuration management and automation. The majority of the services at PIC are managed by Puppet [11]. The implementations are flexible enough to rapidly evolve following changing technologies. All code has been recently migrated to Puppet 3.6. Local repositories for code development projects at PIC have been migrated from SVN to GIT/gitlab, which eases the use and maintenance of the service. Thanks for the large automation of services, PIC is operated with less manpower than the average Tier-1 centre[12].

Cooling is a natural place to save costs for a large computing centre. Given the growth expectations, and coinciding with the LHC shutdown period, PIC improved the energy efficiency of its main computing room. This occurred during fifteen consecutive weeks of work in 2014, without any downtime, interruption and/or negative impact in operations (see [13] for more details). Before the intervention, there was no separation of cold/hot air in the main computing room. Several CRAHs (Computer Room Air Handler) were managing the air through a cold water battery, injecting air at 14° C to get a room temperature of 22-23° C. This system showed a PUE (Power Usage Effectiveness) of 1.8, offering room for improvements.

Three free-cooling units, acting as indirect heat exchangers with outside air and equipped with adiabatic cooling humidifiers, replaced some CRAH units. Direct free-cooling was not considered, as the region were PIC is located has high humidity values and it is within a dusty environment. Figure 3 shows some photographs of the free-cooling units installation. Hot and cold flows are separated in the room, and a ceiling was installed to contain the hot air from the room, as can be seen in Figure 4. Inlet temperature was increased to 20° C, according to ASHRAE recommendations[14].



**Figure 3.** Free-cooling units being installed in PIC, in mid June 2014.

This work was completed in September 2014. A period of one year has been defined to study and adjust the system, in order to reach the maximum energy efficiency possible. Dedicated monitors for the most important climate parameters have been installed to help in this direction. Once tuned, the PUE is expected to be in the range 1.45-1.3. In December 2014, a PUE of 1.3 was measured, even without a fully tuned system, which is really promising. Electricity costs savings in the next 3-3.5 years will amortize the new cooling installation.

As of today, around 200 kW of IT equipment is installed in PIC. The UPS system of the centre was at the end of its lifetime, presenting power losses up to 15%. At the beginning of 2015 a new UPS of 550 kVA, with Insulated-Gate Bipolar Transistor (IGBT) technology that provides efficiency in the range of 97%-99%, was installed, impacting positively on the electricity bill of the centre.

**Figure 4.** PIC main computing room before (left) and after (right) completion of hot/cold air separation and confinement.

Last, but not least, around 70% of CPU resources are installed in a compact module in PIC basement. This module is expected to be upgraded with liquid cooling solutions (immersion) by the end of 2015 or beginning of 2016.

## 5. A reliable, high-capacity Tier-1 service

One of the main characteristics of Tier-1 centres, beyond a very large storage and computing capacity, is to provide these resources through services that need to be extremely reliable. Being closely connected to the detectors' data acquisition, a maximum time for unintended interruption of the services in a Tier-1 is set to 4 hours, and a maximum degradation of Tier-0 to Tier-1 data acceptance of 6 hours [15]. Critical services in a Tier-1 operate in 365x24x7 mode.

Service quality and stability are amongst the cornerstones of the project, therefore they are closely tracked by monitoring two metrics provided by the SAM monitoring framework: site *availability* and *reliability*. These are built from dozens of sensors, for each of the experiments, which hourly probe all of the site Grid services, ensuring peer pressure and guaranteeing that the reliability of WLCG service keeps improving[16].
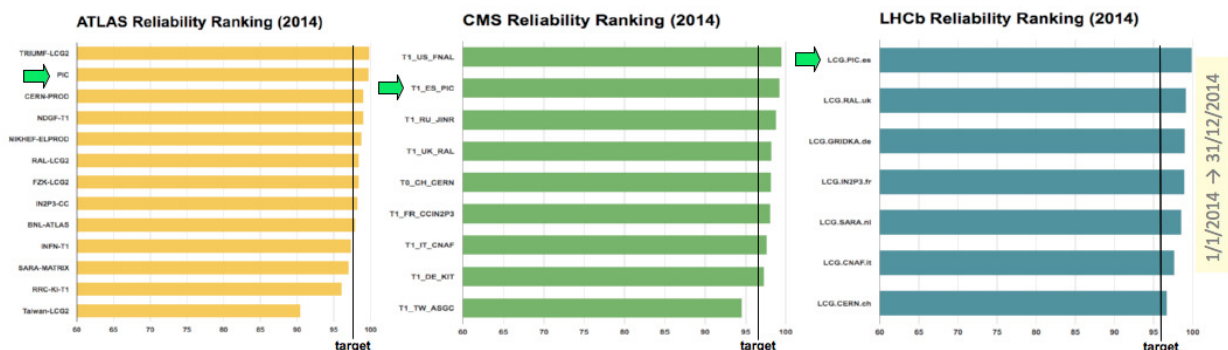


**Figure 5.** Site Reliability ranking plots for Tier-0 and Tier-1s, during 2014. Target for site reliability is set to 97%, according to the WLCG MoU.

Figure 5 shows the reliability ranking results of Tier-0 and Tier-1s for 2014. PIC Tier-1 is at the top of these Reliability rankings (99.9% ATLAS, 99.4% CMS, 99.9% LHCb), and well above the WLCG target, which is set to 97%. These figures were obtained during a year in which many new services were deployed and interventions were made to improve the site efficiency. PIC is one of the smallest Tier-1 centres in WLCG, but supports 3 LHC VOs. PIC has an expert contact person per experiment on site (the liaison), communicating and coordinating priorities with each of the experiments and resolving operational problems. This helps PIC being at top reliability and stability levels.

## 6. Conclusions

LHC LS1 has been a very active period for computing. In a budget-constrained period, experiments were required to modify their computing models and improve the efficiency of resources usage. The resources growth profiles for the next years are compatible with a flat-funding model, in which the growth is achieved by means of technology trends. PIC Tier-1 deployed many new services which were required, and validated them during LS1, demonstrating its readiness for the next LHC data-taking period. Benefiting from the LHC shutdown, PIC improved its computing infrastructure, with new elements that translate in a significant reduction of the electricity costs for the next years. This work gives strength to the centre for the following challenges it will be facing.

## References

[1] *LHC Computing Grid Technical Design Report*, CERN-LHCC-2005-024, 20 June 2005.
[2] *Computing for the Large Hadron Collider*, Ian Bird, Annual Review of Nuclear and Particle Science, Vol. 61: 99-118, November 2011.
[3] Update of the Computing Models of the WLCG and the LHC Experiments, LCG-TDR-002: http://cds.cern.ch/record/1695401/files/LCG-TDR-002.pdf
[4] https://hepix.caspur.it/benchmarks/doku.php.
[5] http://www.dcache.org.
[6] J. Bakken et al., *Enstore Technical Design Document*, http://www-ccf.fnal.gov/enstore/design.html.
[7] http://xrootd.org.
[8] *multi-core job scheduling in the Worldwide LHC Computing Grid*, A Forti, A Pérez-Calero Yzquierdo, et al, in these proceedings.
[9] *Scheduling multi-core workload on shared multipurpose clusters*, J Templon, et al, in these proceedings.
[10] *LHC Optical Private Network*, http://lhcopn.cern.ch.
[11] https://puppetlabs.com/.
[12] *Optimising costs in WLCG operations*, A. Sciabà et al, in these proceedings.
[13] *Free cooling on the Mediterranean shore: Energy efficiency upgrades at PIC*, V Acín, M Delfino, et al, in these proceedings.
[14] See http://www.ashrae.org/File%20Library/doclib/Public/20100901_ASHRAED2468520050330.pdf.
[15] Memorandum of Understanding for Collaboration in the Deployment and Exploitation of the Worldwide LHC Computing Grid, http://wlcg.web.cern.ch/collaboration/mou
[16] See SAM3 probes in http://dashboard.cern.ch/.