# The CMS CERN Analysis Facility (CAF)

**O.Buchmüller[1], D.Bonacorsi[2], F.Fanzago[3], S.Gowdy[4], P.Kreuzer[5], L.Malgeri[4], R.Mankel[6], S.Metson[7], B.Panzer-Steindel[4], J.Afonso Sanches[8], U.Schwickerath[4], D.Spiga[4, 9], D.Teodoro[8], Rainer Többicke[4]**

[1]Imperial College, [2]Universita & INFN, Bologna, [3]Universita & INFN, Padova, [4] Conseil Europeen Recherche Nucl. (CERN), [5]Rheinisch-Westfaelische Tech. Hoch. (RWTH), [6]Deutsches Elektronen-Synchrotron (DESY), [7]University of Bristol, [8]Universidade do Estado do Rio De Janeiro (UERJ, [9]Universita & INFN, Perugia
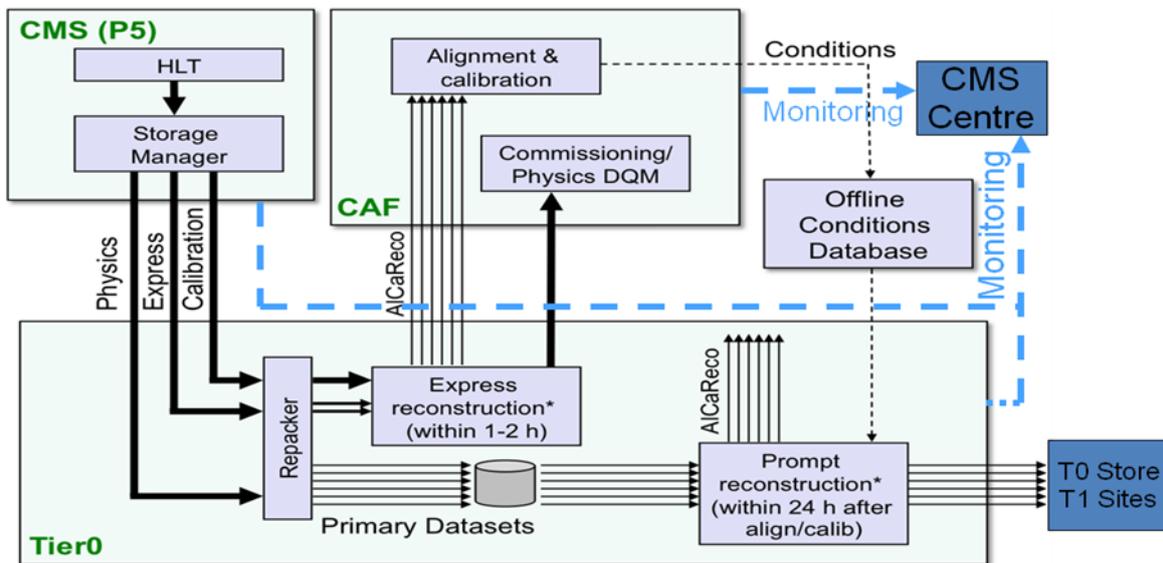
Peter.Kreuzer@cern.ch

**Abstract**. The CMS CERN Analysis Facility (CAF) was primarily designed to host a large variety of latency-critical workflows. These break down into alignment and calibration, detector commissioning and diagnosis, and high-interest physics analysis requiring fast-turnaround. In addition to the low latency requirement on the batch farm, another mandatory condition is the efficient access to the RAW detector data stored at the CERN Tier-0 facility. The CMS CAF also foresees resources for interactive login by a large number of CMS collaborators located at CERN, as an entry point for their day-by-day analysis. These resources will run on a separate partition in order to protect the high-priority use-cases described above. While the CMS CAF represents only a modest fraction of the overall CMS resources on the WLCG GRID, an appropriately sized user-support service needs to be provided. We will describe the building, commissioning and operation of the CMS CAF during the year 2008. The facility was heavily and routinely used by almost 250 users during multiple commissioning and data challenge periods. It reached a CPU capacity of 1.4MSI2K and a disk capacity at the Peta byte scale. In particular, we will focus on the performances in terms of networking, disk access and job efficiency and extrapolate prospects for the upcoming LHC first year data taking. We will also present the experience gained and the limitations observed in operating such a large facility, in which well controlled workflows are combined with more chaotic type analysis by a large number of physicists.

## 1. Introduction

The design of the CMS CERN Analysis Facility (CMS CAF) is based on the assumption of routine data taking mode. Under such regime, while most CMS physics analysis is distributed around the world as described in the CMS computing model [1], the CAF plays a unique role for latency-critical functions at the source of the analysis chain, with no dependence on the GRID infrastructure. Such workflows are based on a highly selected subset of the full production data – so called Calibration and Express Streams (ES), and are made available to restricted analysis groups soon after data taking occurred (~1h). They take precedence over all other activities on the CAF and can be classified as follows in order of priority:

- alignment and calibration
- trigger and detector diagnostics, monitoring and performance analysis
- physics monitoring, analysis of express stream, fast-turnaround high-priority analysis

The overall CMS data flow at CERN including these 3 main CAF activities are shown in Fig.1: the data is sent from the online computing centre near the CMS detector (P5) to the Tier-0 centre at the Meyrin CERN site, where dedicated calibration and express streams are processed and sent to the CAF within a very short latency (~1-2 h). After new calibration and alignment constants have been computed on the CAF, the condition databases are updated and the full event reconstruction can start at the Tier-0 (24 h cycle). Other performance or physics analysis activities on the CAF might be carried out on a much longer timescale and with different priorities, hence requiring careful resource management due to the large variety of workflows.



**Figure 1: CAF Data Flow**

The key hardware components for addressing the low-latency and high-priority requirements are a large dedicated processing farm and a highly accessible disk storage system. The limitation in available resources demands selectiveness in the jobs to be run and in the data to be stored, hence the need of a careful CPU and data management planning and prioritisation. It also involves rigorous user management, by restricting access to persons performing the corresponding priority tasks and by closely monitoring their activities.

In addition, a large number of CMS collaborators physically present at CERN for a longer period of time will need to use the analysis facility as an entry point for their day-by-day work. In this context, the CAF is seen in union with the CMS share of the public interactive and batch services at CERN; it will also support dedicated physics communities in the same spirit as other CMS Tier-2 centres on the GRID [2]. While these extensions of the CMS CAF are currently under planning and implementation, they are not covered in the present paper.
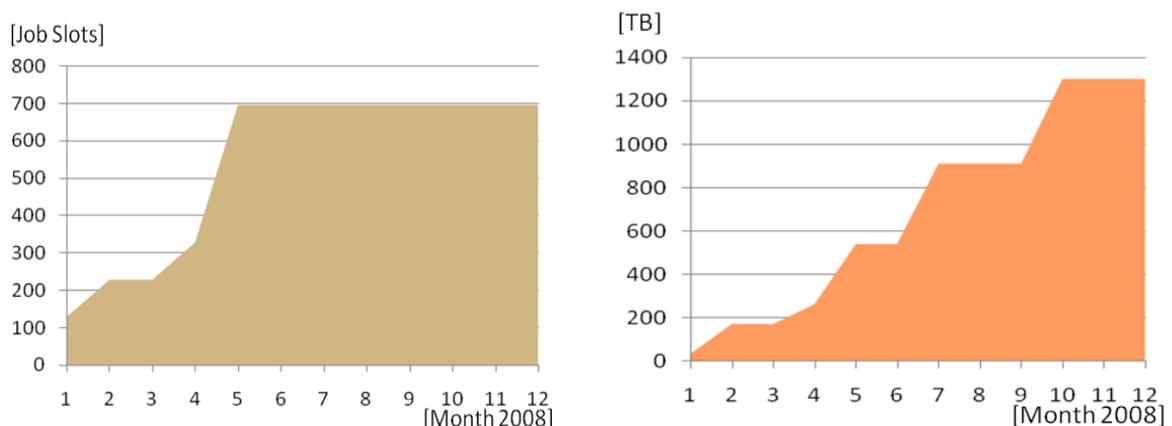
The CAF also plays a major role during the commissioning phase of the detector, when analysis activities are focused on cosmic data, single beam or low luminosity collision data. Under such regime, the increasing complexity as new detectors are brought into operation requires careful and frequent analysis of a large fraction of the data, hence the expectation on the CAF to keep full copies of each reconstructed primary dataset (PD) and occasionally even of the raw data (prior to reconstruction). This requirement differs from the stable operation mode and is not sustainable on the long term in particular once the volume of produced data will increase. In that sense the early commissioning phase provides an excellent stress-test of CAF resources and of the data and workflow management systems, as demonstrated during several run periods in 2008 and as planned for 2009.

In this paper we briefly review the installed CAF resource, followed by CAF performances during commissioning phases and a projection into the expected use cases for the LHC start-up.

## 2. CAF Resources

The technical requirements on the CAF are based on the capability to support high-priority and low-latency workflows, while enforcing controlled and prioritised access policies to hundreds of users. The CAF systems need to be setup such that particular jobs can start with very short latency in parallel to many other concurrent jobs and with highly efficient access to the CMS production data.

A dedicated processing farm has been installed in several steps during 2008, reaching a total of ~100 worker nodes and ~700 cores[1]. Access, scheduling and sharing is handled by the CERN LSF system. The CAF worker nodes can receive 1 job slot per core[2], leading to nearly 700 available job slots as shown on the left of Fig.2. The batch system comprises one main queue ("cmscaf" ~ 630 job slots) available to all CAF users and a special queue ("cmsexpress" ~ 50 reserved job slots) dedicated to very high-priority alignment and calibration workflows.



**Figure 2: CAF CPU and Disk Resource ramp up during 2008**

A large disk storage system based on CASTOR technology [3] has been installed in various steps during 2008, reaching a total of 215 disk servers and a capacity of ~1.3 PB. Highly efficient and asynchronous access to the data coming from the Tier-0 is ensured by the disk-only characteristic of this storage pool, hence avoiding any kind of tape staging latency, that are typical bottle necks encountered for large storage systems such as Tier-1 centres. The plot in Fig.2 also includes a user disk pool (50 TB) hosting the output data from the CAF user analysis. In addition, CAF sub-groups and users are sharing a 2 TB AFS space for POSIX-like storage needed for analysis. Finally, a set of 5 worker nodes (40 cores) are reserved for interactive access to CAF batch and storage resources, in particular for job submission.

The resources described in this section are reserved to CAF users only, who must belong to one of the 3 main groups, according to the main CAF activities: alignment and calibration, detector and trigger commissioning, priority physics analysis. These are further sub-divided into sub-groups, according to specific detector types or analysis channels. There are currently ~30 sub-groups and the goal is to populate each with maximum 10 users, in order to limit the total number of registered users. CAF user registration is done by the sub-group leaders via a single WEB entry-point: this interface is also used to set the LSF fair share fraction of each sub-group and the access permissions to the CASTOR and AFS storage resources. The disk space used on the production storage disk is managed "manually" according to a gentlemen agreement between the main CAF group leaders[3].

---

[1] Some worker nodes contain 8 cores, some 4 cores.

[2] 2 worker nodes were reserved for memory intensive alignment activities (2 cores - 4GB RAM - per job slot).
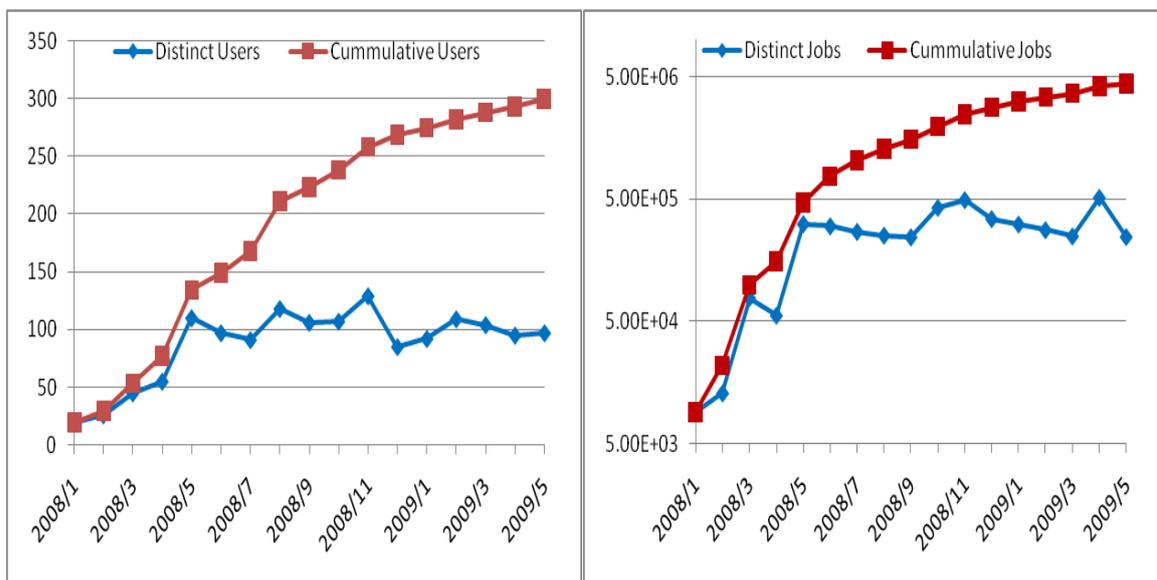
[3] Current CASTOR versions do not offer the possibility of quotas; hence manual space management is required.

### 3. CAF Commissioning and Performances

In parallel to the resource ramp-up achieved throughout 2008, processing activities on the CAF included (i) the integration of key data management and workflow management components, (ii) the commissioning of newly developed alignment, calibration, monitoring and analysis workflows and (iii) the scale testing of the whole CAF machinery, in particular during dedicated "data challenges", based either on simulated physics data or on real cosmic data. So far the role of the CAF has been essentially as described at the end of Sect.1, namely corresponding to the commissioning phase of the experiment, where full copies of the reconstructed CMS data are stored and analysed on the CAF. This mode will continue in 2009 until the LHC startup, at which point CAF activities should progressively evolve towards the design mode, with focus on high-priority activities based on only a subset of the CMS production data. Nevertheless, the achieved data challenges have already provided important feedback in terms of low-latency data management and processing capability of the CAF.

The main CAF data management and workflow management components include the job submission tool "CRAB" [4], the Data Bookkeeping System (DBS) [5] and the Data Transfer system (PhEDEx) [6]. These are standard CMS tools, with only minor adaptations compared to those used for distributed analysis on the GRID.

A good representation of the growing CAF activities is given in Fig.3, where the number of distinct and active users (left) and of jobs (right) is shown from winter 2008 until spring 2009. The CRAFT08 data challenge late Oct.2008 contributed to the highest number of distinct CAF users in a given month (130), while the number of monthly jobs reached 500,000. The cumulative number of active CAF users so far has reached 300 while there were ~450 registered CAF users[4] at the time of editing this paper.
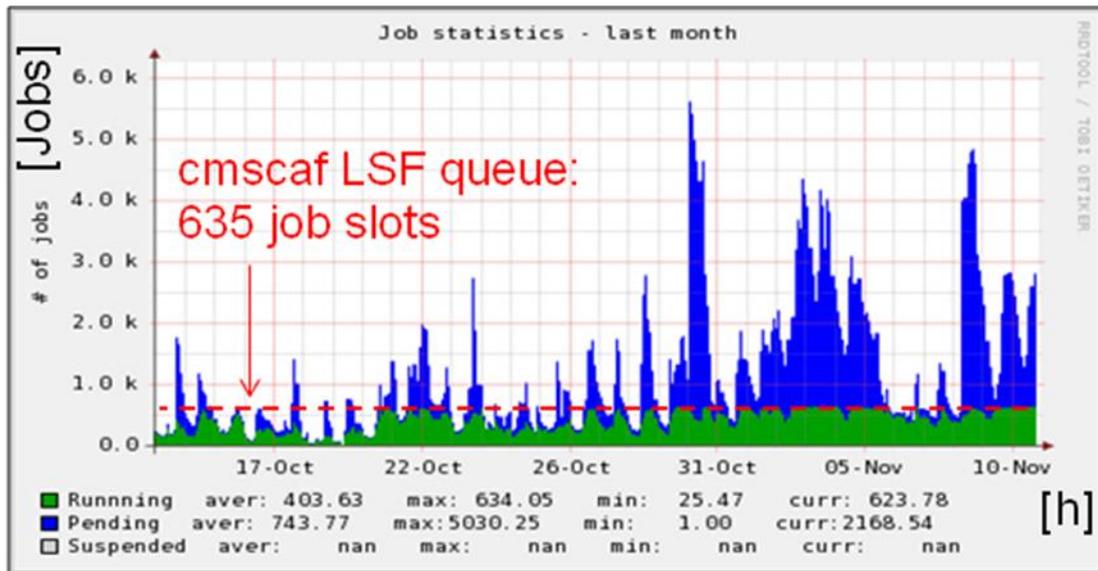


**Figure 3: Number of distinct/cumulative CAF users (left) and number of distinct/cumulative CAF jobs (right)**

The monitoring of user activity and enforcement of access policies has shown to be a challenging task for both Facilities Operations and CAF group leaders. It is an essential ingredient for reaching the design expectation of the CAF; in particular all non-high-priority analysis workflows should be migrated onto the distributed Grid computing fabric in the future.

---

[4] Not all registered CAF users have been active so far.

The "cmscaf" batch queue (~630 job slots) relies on a fair share system, in which ~30 CAF subgroup are allocated equal resources averaged over time. The main CPU management restrictions are a maximum of 100 running jobs per user and a maximum CPU time of 1 week. It has been extensively commissioned, in particular during the CRAFT08 data challenge (see Fig.4), where the queue was saturating 50% of the time and the average job slots usage was of 67%. The utilization of the "express" queue was more modest, however it stays justified as the need to cover very low-latency workflows is seen as crucial for the LHC start-up.
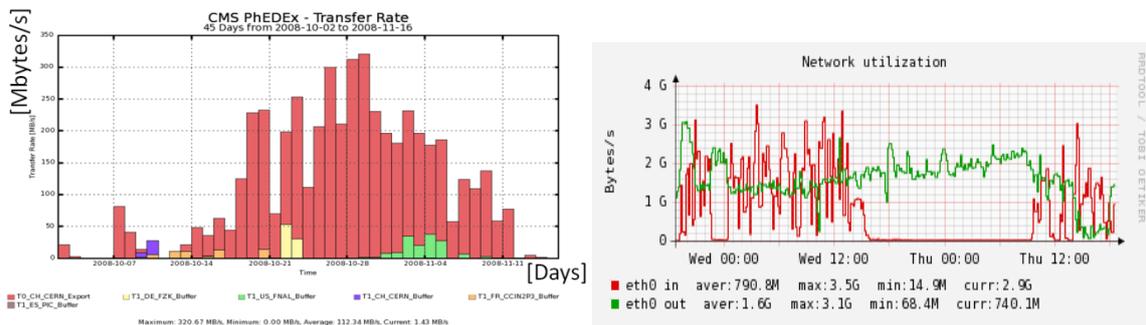


**Figure 4: Number of pending and running CAF jobs (2008)**

No major bottle-neck was encountered with the CAF batch system during the 2008 data challenges. Occasional issues observed were:

- Very large amount of submissions (>500) at a given time by single "heavy" users, sometimes affecting the latency needs of competing workflows
- Low Job CPU efficiency (down to <50%), limiting the overall CPU utilization performance.
- Large job memory consumption, leading to crashes of particularly "heavy" workflows

The first 2 issues can be addressed via a stronger control of the user activity and by more systematically monitoring the job efficiencies, while the latter issue has been addressed by increasing the reserved memory when submitting the job (up to 8 GB / job for some workflows).
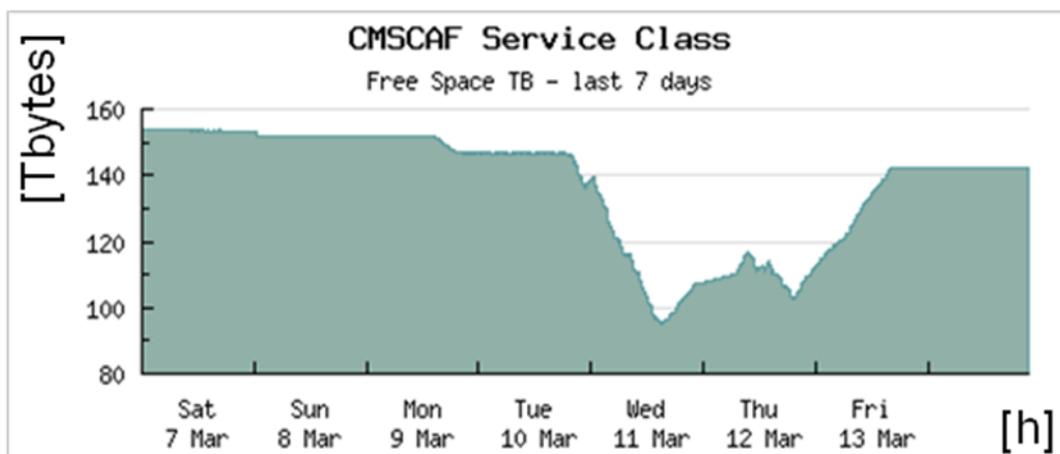
The performance of the CAF production storage pool and related data management has also not encountered any major issue, during the commissioning phase. The Tier-0 to CAF data transfer performances have shown sustained rates above 100 MB/s during several weeks and peak rates near 3 GB/s, as shown in Fig.5. These performances reside well within the design figures.

**Figure 5: Transfer rates from Tier-0 and other Tier-1s into the cmscaf disk pool**

One encountered complication is that the CASTOR namespace where the CAF data reside is common to the public CMS data namespace at CERN, while the access pattern to the former is reserved to CAF users, hence requiring special tuning and monitoring of the CAF data management tools.

Another important aspect is the need for a routine monitoring and management of the free space on the CAF storage pool. While the disk-only characteristic provides highly efficient data accessibility, it has the disadvantage of no automatic garbage collection to tape storage, hence the risk to fill up the disk in the middle of production. To address this issue, a procedure has been established where the CMS computing shift crew receives an alarm once the free space is below a given thresholds and notifies the CAF data managers, including all necessary information for them to promptly trigger data deletion. The result of this procedure is illustrated in Fig.6, where the decrease and increase in free space reflects data transfers into and data deleted from the CAF storage disk.
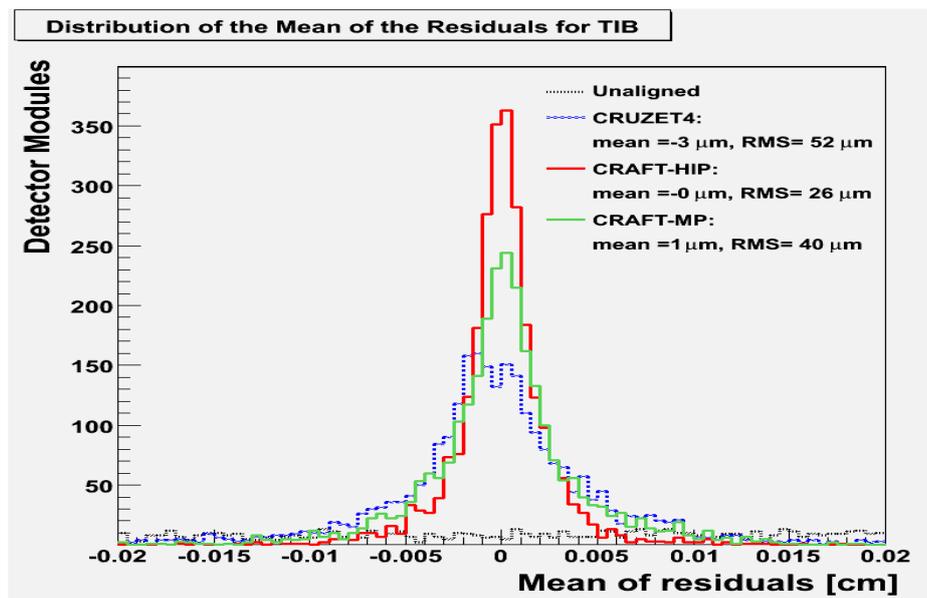


**Figure 6: Free Space on the cmscaf disk pool (March 2009)**

In terms of user input data access, the main encountered difficulties were related to users accessing the wrong storage pool at CERN due to configuration mistake, and very few CASTOR glitches. An efficient ticketing service has been put in place between the CAF users and facilities managers, who in turn solve the problems with the CERN CASTOR team if needed. Finally, as mentioned in Sect.2, a CAF user storage disk has been deployed to host the analysis output data, preventing a mixture of production and non-production data on the production pool and hence further increasing data management difficulties.

Since the purpose of this paper is not to provide a full report about CAF analysis results, we only list for reference important milestones where the CAF computing infrastructure and operations were particularly challenged. The "Computing Software and Analysis challenge of 2008" (CSA08)

consisted of a collection of exercises aimed to test the full scope of data handling and analysis activities needed for LHC data-taking operations. It was based on simulated physics data and several crucial alignment and calibration workflows were tested on the CAF. CMS then turned to cosmic data analysis, first through the "Cosmic Global Run at Zero Tesla" (CRUZET08) followed by a "Cosmic Run at Four Tesla" (CRAFT08). In the latter the CMS experiment ran during 4 weeks continuously, acquiring ~300M cosmic events. Such data taking periods represented an excellent opportunity to validate and improve the CAF analysis workflows. The CAF user activity was largely dominated by calibration, alignment and detector commissioning, while physics analysis workflows were distributed at a more modest level throughout the year. As an example, the result of a CMS tracker alignment analysis in the inner barrel is shown in Fig.7, before and after the alignment procedure.



**Figure 7: Mean Residuals of CMS tracks reconstructed in the Inner Barrel detector, before (dots) and after various alignment steps processed on the CAF**

## 4. Conclusion and Projection into the LHC start-up

The CMS CERN Analysis Facility has been successfully designed, built and commissioned, in particular for the LHC start-up planed for fall 2009. The highly accessible disk storage system and the flexible batch system have met the low latency and high priority requirements, without encountering any major bottle neck.

In view of the LHC start-up, where CAF workflows will increasingly be based on sustained and highly selected input data, the low latency and high priority requirements will become more critical. Various CAF components and aspects need to be strengthened in order to address this evolution, including a routinely job prioritization and user activity monitoring, an improvement in the average job efficiency, a stronger data management policy for old-data deletion, the automation of job submissions to trigger latency-critical calibration workflows as soon as the input data has arrived at the CAF. The remaining running periods without LHC collisions will be exploited by CMS to reach these goals.

## References
[1]    C.Grandi, D.Stickland, L.Taylor et al. "The CMS Computing Model" CERN-LHCC-2004-035/G-083 (2004)
[2]    The CMS Collaboration "CMS Computing Technical Design Report", CERN-LHCC-2005-023,

(2005)

[3]    Lo Presti, G. Barring, O, Earl, A, Garcia Rioja, R.M, Ponce, S, Taurelli, G, Waldron, D, Coelho Dos Santos, M. "CASTOR: A Distributed Storage Resource Facility for High Performance Data Processing at CERN" 24th IEEE Conference on Mass Storage Systems and Technologies, 2007. MSST 2007.

[4]    D.Spiga et al., "The CMS Remote Analysis Builder (CRAB)", Lecture notes in Computing Science, 2007.vol.4873,pp.580-586

[5]    A.Delgado Peris et al. Data location, transfer and bookkeeping in CMS, 18th Hadron Collider Physics Symposium 2007 (HCP 2007) 20-26 May 2007, La Biodola, Isola d'Elba, Italy. Published in Nucl.Phys.Proc.Suppl.177-178:279-280,2008

[6]    T.Barras et al. "Software Agents in Data and Workflow Management", CHEP04 Conference, Interlaken, Switzerland, Published in *Interlaken 2004, Computing in high energy physics and nuclear physics* pp.838-841