**EUROPEAN ORGANIZATION FOR NUCLEAR RESEARCH**
**European Laboratory for Particle Physics**

# Online Performance Monitoring of the Third ALICE Data Challenge (ADC III)

**Authors:**

W. Carena, R. Divià, P. Saiz, K. Schossmaier, A. Vascotto, P. Vande Vyvre
for the ALICE Collaboration

# Online Performance Monitoring of the Third ALICE Data Challenge (ADC III)

W. Carena[(1)], R. Divià[(1)], P. Saiz[(2)], K. Schossmaier[(1)],
A. Vascotto[(1)], P. Vande Vyvre[(1)] for the ALICE Collaboration

CERN EP-AID[(1)] and EP-AIP[(2)], CH-1211 Geneva 23, Switzerland

### Abstract

The ALICE data acquisition system has been designed for a maximum bandwidth of 2.5 GB/s for event building and of 1.25 GB/s for mass storage. In order to attain a gradual integration of the overall computing infrastructure, the present hardware components and software prototypes are tested during regular ALICE data challenges. The third one (ADC III) took place from January to March 2001 as a joint effort between the ALICE online/offline team and the CERN IT division. The main goal of this data challenge was to achieve a stable 300 MB/s throughput in the event building network and a 100 MB/s throughput to CASTOR over period of a few days.

Performance monitoring was another goal of this exercise, where a prototype (*dateStat*) was developed to collect and display statistics. In this paper we will introduce this online monitoring system and report on some of the obtained results. It is structured in three parts: (1) An overview will be given on the testbed hardware, the software running on it, and the data flow. (2) The architecture of the monitoring system will be described, which consists of a set of C programs, Perl/gnuplot/CGI scripts, and a MySQL database. It allows to measure individual/aggregate data rates, collected data volumes, and CPU loads. All these values can be visualized on web pages both on a run-by-run and global basis. (3) Various plots will be shown to illustrate the usefulness of this online monitoring system and to document the outcome of the ADC III. Finally, some ideas will be pointed out how to advance *dateStat*.

## 1 Introduction

The ALICE Data Challenges are large-scale high throughput distributed computing exercises, which are jointly conducted by the ALICE experiments [1] and the CERN IT division. Starting in 1998, roughly once a year a data acquisition system is assembled, consisting of the present hardware components and software prototypes, to realize an early integration of the overall computing infrastructure. The idea is

to gradually increase the size and the complexity of each data challenge, since the final ALICE data acquisition system has to provide a maximum bandwidth of 2.5 GB/s for event building and of 1.25 GB/s for mass storage[†].

The third ALICE data challenge (ADC III) has been carried out from January to March 2001. It comprised the CERN Tier0/Tier1 fabric, the mass storage system, the ALICE offline computing, and prototypes of the ALICE data acquisition system. The performance goals of the ADC III were to achieve an aggregate bandwidth of 300 MB/s within the data acquisition system, to reach a stable bandwidth of 100 MB/s over the complete chain for a week, and to store a total of 80 TB in the mass storage system. A comprehensive report [2] on the ADC III is available.

Another functional goal of the ADC III was online performance monitoring, which plays a vital role when executing distributed applications on such large-scale and heterogeneous computing systems. For this reason a prototype (called *dateStat*) has been developed during the data challenge that collects and displays system as well as application specific statistics. In the following we will introduce *dateStat* by giving a brief overview of the testbed (Section 2), by describing the architecture of this monitoring system (Section 3), and by presenting some selected performance results obtained by it during the ADC III (Section 4). Concluding remarks and an outlook how to advance this online monitoring prototype finishes the paper.

## 2  Testbed and Data Flow

The ADC III was running on the LHC testbed, which is a common computing fabric prototype at CERN to evaluate all aspects of the LHC computing. Such a testbed and the repetition of data challenges were recommended by the LHC Computing Review [3]. Our portion of the testbed consisted of the following components, which amounts to approximately 10% of the final ALICE data acquisition setup:

- A farm of 80 standard PCs equipped with dual Pentium III processors at 800 MHz, Intel L440GX+ chipset, and main memory of 512 MB. Only 27 of these machines had a Gigabit Ethernet card (Netgear GA620T) attached, whereas the others had a commodity Fast Ethernet card. On all these machines Red Hat 6.1 Linux was installed running a 2.2.17 kernel.

- An Ethernet network based on 6 switches (Extreme Networks, 3COM) and 2 routers (Cabletron). Trunking between 2 switches (four parallel Gigabit links on copper and fiber media) was used to permit load sharing and to provide an aggregate bandwidth of about 300 MB/s.

- A set of 8 disk servers featured with dual Pentium III processors at 700 MHz, Intel L440GX+ chipset, main memory of 512 MB, and 20 IDE disks that provide 750 GB hardware-mirrored disk space.

---

[†]Throughout this paper kB $\equiv 2^{10}$ Bytes, MB $\equiv 2^{20}$ Bytes, GB $\equiv 2^{30}$ Bytes, TB $\equiv 2^{40}$ Bytes.

- Up to 12 tape drives (STK 9940) in conjunction with standard PCs, where each cartridge has a capacity of 60 GB and a raw bandwidth of 10 MB/s.
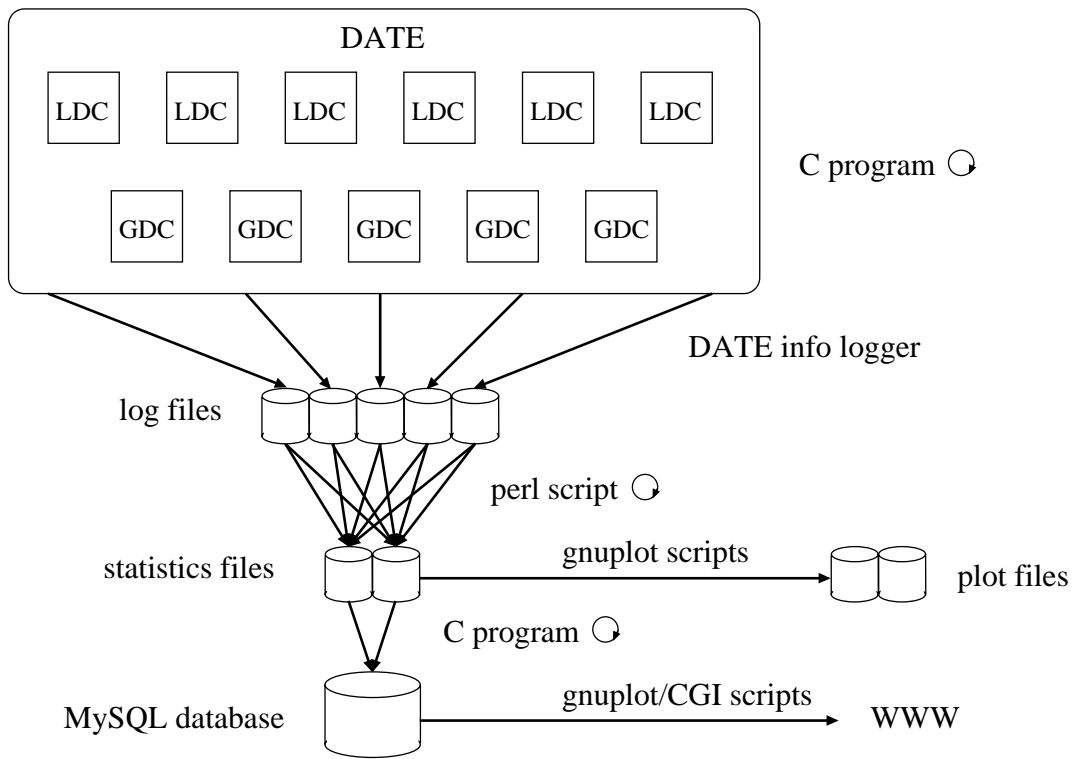
The data flow during the ADC III was generated by the Local Data Concentrators (LDC), which are PCs of the farm that inject simulated physics raw event data of the ALICE TPC, produced in advance by AliRoot [4]. These parallel data streams go over the Fast/Gigabit Ethernet network using the TCP/IP protocol to reach the Global Data Collectors (GDC), which are also PCs of the farm that perform the event building. Both LDCs and GDCs are operated by the ALICE Data Acquisition Test Environment (DATE) system [5] version 3.7. The full events are then formatted with the ROOT I/O library [6] and finally funnelled to the mass storage system CASTOR [7], which manages a huge migrating filesystem. No trigger system was installed and no event filtering was done during this data challenge.

# 3    Performance Monitoring

Online performance monitoring of the overall system was an essential functional improvement of this data challenge. For this purpose the Performance and Exception Monitoring (PEM) software package [8] was planned to be used, however, it could not be made ready for the ADC III. Therefore other existing tools have been applied instead for monitoring purposes:

- A fabric monitoring tool developed by the CERN IT-PDP group displays the current status of the fabric (CPU load, network activities, swapping activities) in form of a pie chart including a history of these parameters. Permanently running agents on fabric components are sending UDP packets to a server which visualizes the performance information by using Tcl/Tk scripts.

- The ROOT I/O formatting program is also measuring the aggregate throughput to the mass storage system. History and histogram information of these measurements are accessible on a web page.

- The CASTOR system itself maintains a set of statistics, for instance the size of the pools on disk and tape, or the used tape volumes. These and many other statistics are also accessible on web pages.

All these tools turned out to be of invaluable help during the whole ADC III, but obviously they were not designed to measure DATE specific parameters such as LDC/GDC throughput (individual and aggregate), data volume (individual and aggregate), or CPU load (user and system). Also, an identification in terms of run number, event number, and fine grained time-stamping was missing. Given this deficiency, a prototype *dateStat* was developed during this data challenge. Its architecture can be seen in Figure 1.
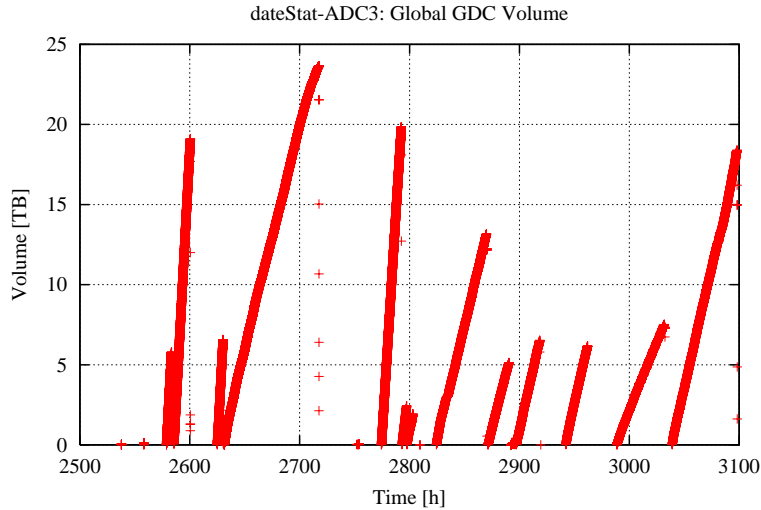
**Figure 1:** Architecture of *dateStat*.

On each LDC and GDC a process is launched at start-of-run and stopped at end-of-run. At regular intervals it makes measurements on DATE (run number, number of events, event data rate, event data volume) and on the system (CPU loads). These measurements are then shipped via the DATE internal message system (info logger) to a single machine which maintains a log file for each running LDC and GDC. These log files grow with about 100 kB/hour by a measurement interval of 30 seconds. A perl script extracts the relevant performance values and calculates aggregate values (sum of appropriate entries with approximately the same time-stamp) for event rate and volume, which are then stored in statistics files.

In an early version of *dateStat* these statistics files were converted to gif files with the help of gnuplot scripts. In the final version these statistics files were inserted at regular intervals in a MySQL database, which made it much easier to retrieve performance data of certain DATE runs and to combine values for global views. Further gnuplot and CGI scripts were feeding a web server with the resulting performance data, hence everybody involved in the ADC III could observe the status and history of the whole system.

Despite the simple design and rather inefficient implementation of *dateStat*, it provided us with many insights about the behavior of the whole data acquisition chain throughout the ADC III. Some performance results can be seen in the following Section 4.

# 4    Performance Results

During the ADC III more than 2200 runs were executed in order to test all the aspects of the data acquisition system on this testbed. Figure 2 shows a series of runs in terms of the aggregate volume. The time axis is scaled in hours with the beginning of year 2001 as origin. Note that this format will be used for all the following time-plots in this section.
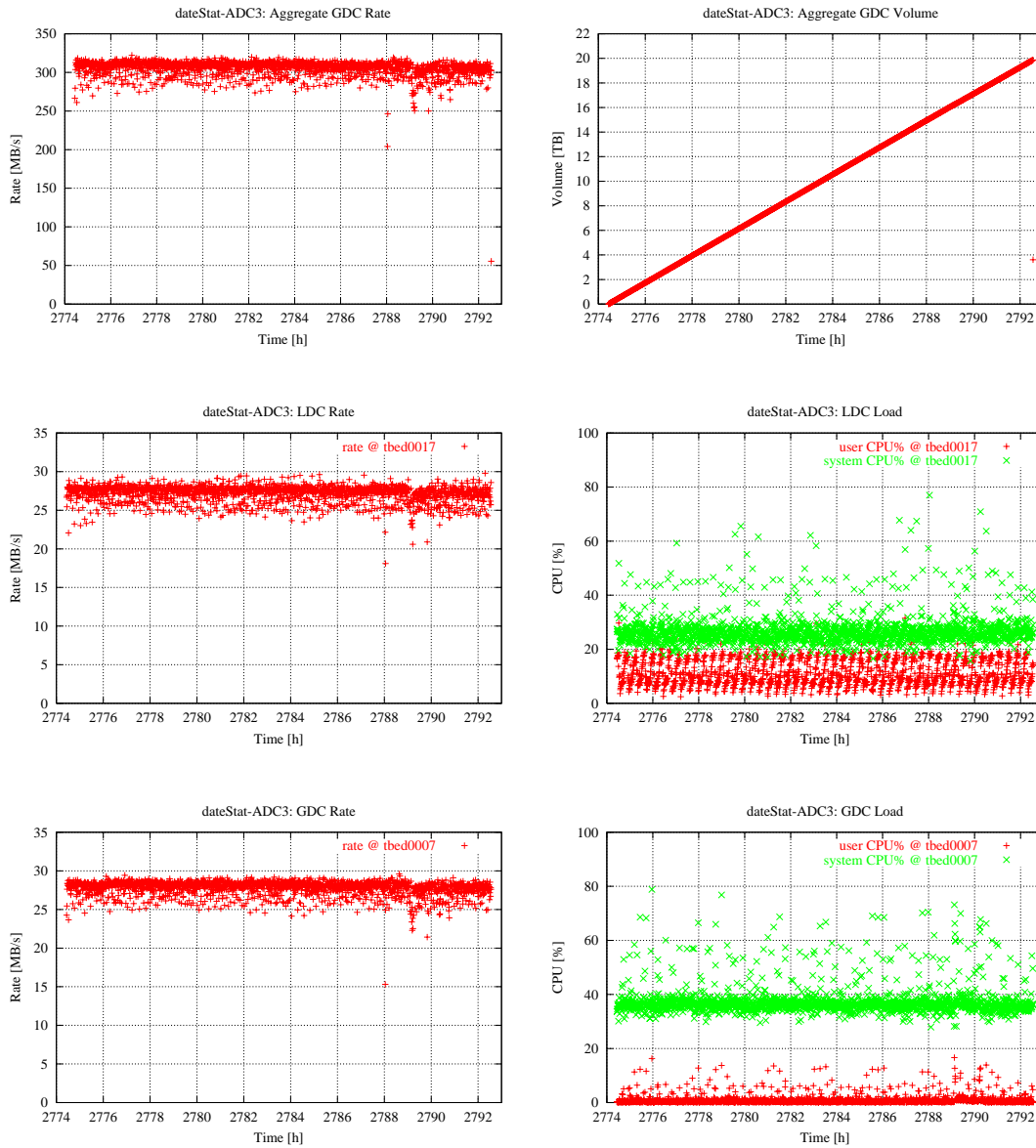


**Figure 2:** Aggregate data volume of DATE runs.

Besides setting different parameters (e.g. assignment of PCs from the farm to run as LDC or GDC, configuration of event size and content) there are several main running conditions of the data acquisition system depending on the destination of the data streams after event building: they are either discarded (DATE standalone), or stored on local disks, or formatted with ROOT I/O and recorded as CASTOR files (complete chain). The following subsections will present performance measurements taken by *dateStat* under different run conditions.
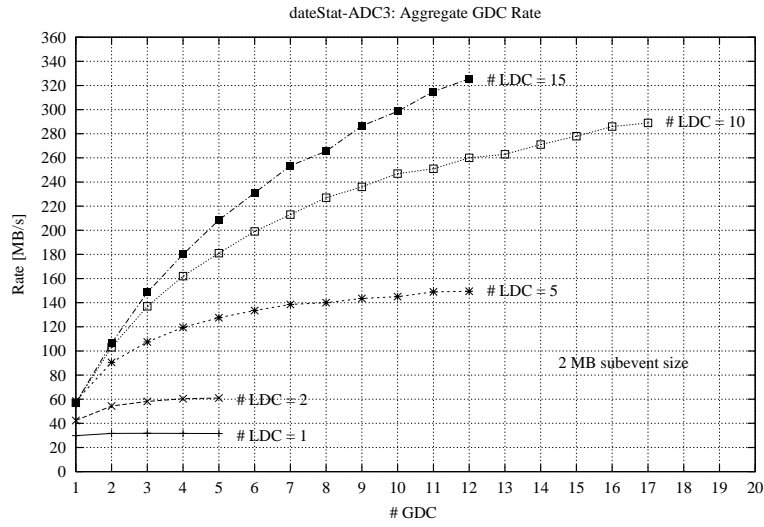
## 4.1    Performance of DATE standalone

Figure 3 shows the performance of a run using 11 LDCs and 11 GDCs, all connected to Gigabit Ethernet, with fixed subevent size and without a recording device. The run lasted 18 hours, had an aggregate rate of 304 MB/s (top left plot), and produced 19.8 TB of full events (top right plot). Some of the LDCs were injecting 420 kB subevents at a speed of 27.1 MB/s (middle left plot) by inducing a 12% user and 27% system CPU load (middle right plot). Each GDC was building events at a rate of 27.7 MB/s (bottom left plot) by consuming 1% user and 37% system CPU (bottom right plot). As a result these plots demonstrate an excellent stability of a high throughput over the switched Ethernet network using commodity equipment.

**Figure 3:** Aggregate/LDC/GDC performance, aggregate volume, and CPU load over time • DATE standalone • 11x11 configuration • 420-440 kB fixed subevent size.

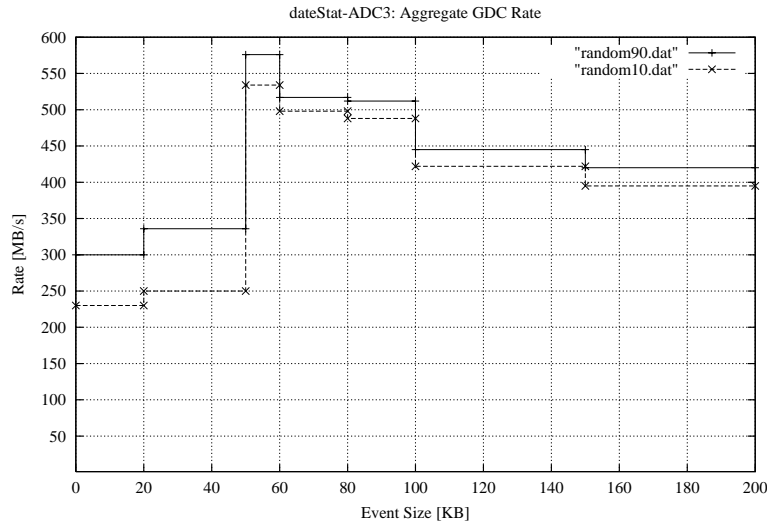## 4.2 Dependence on the number of LDC and GDC

Figure 4 shows the performance of DATE standalone runs with 2 MB fixed subevent size for a varying number of LDCs and GDCs. Two types of saturations are visible: An LDC cannot output more than 30 MB/s (configurations with $|LDC| \ll |GDC|$), and an GDC cannot handle more than 60 MB/s (configurations with $|LDC| \gg |GDC|$). An important outcome of the ADC III was that the scalability of DATE is ensured up to about 30 machines.

**Figure 4:** Aggregate performance over number of LDC/GDC ● DATE standalone ● Gigabit Ethernet ● 2 MB subevent size.

## 4.3 Dependence on the subevent size

Figure 5 shows the performance of DATE standalone runs for a varying subevent data size. The system was composed of 13 LDCs and 13 GDCs, all with Gigabit Ethernet connectivity. This kind of dependence is important to know, since the current estimate of ALICE subevent sizes vary from 20 kB to 440 kB. A maximum performance of 556 MB/s was achieved once during the ADC III by a 1 hour run in the above configuration with subevents of a size between 50 kB to 60 kB.
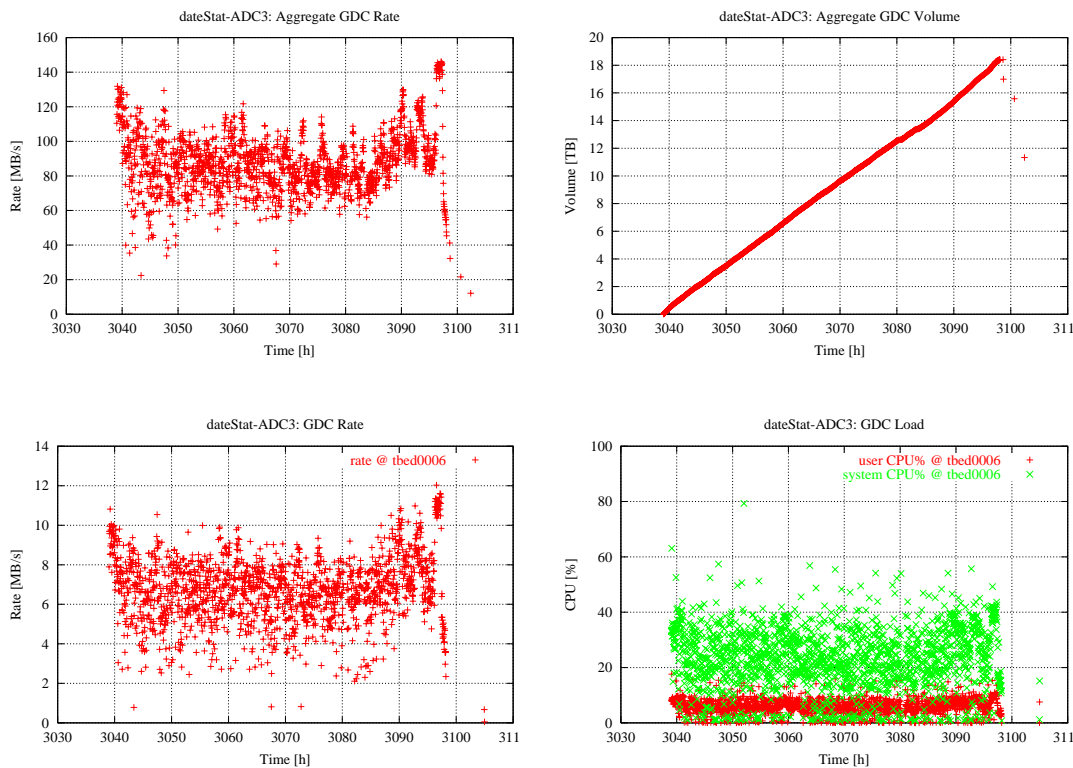


**Figure 5:** Aggregate performance over subevent size ● DATE standalone ● 13x13 configuration ● Gigabit Ethernet.
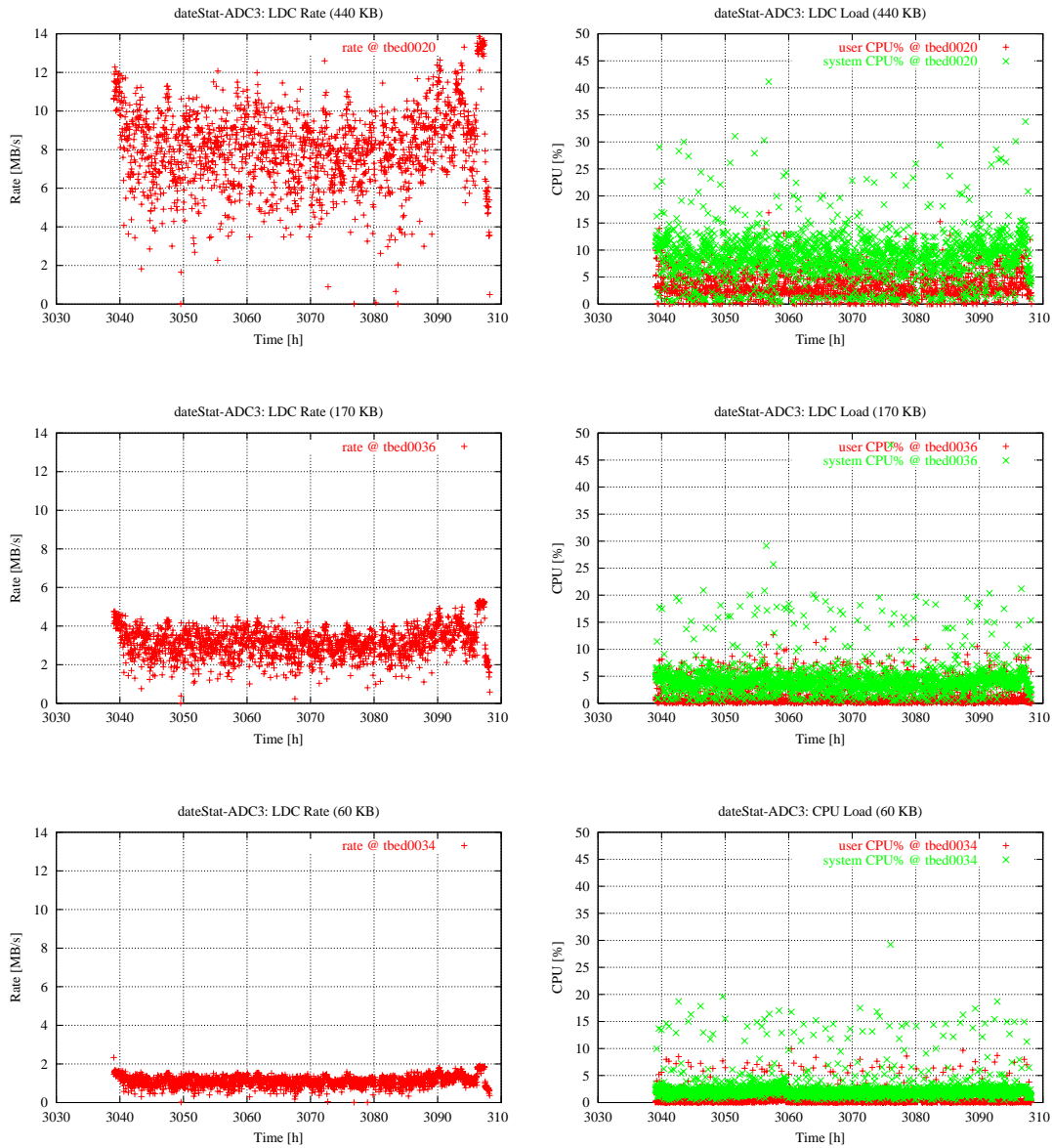
## 4.4 Performance of the complete chain

Most interesting are runs that employ the complete data acquisition chain. In the following example an ALICE-like traffic (a distribution of subevents between 20 kB to 440 kB) was generated by 9 LDCs on Gigabit Ethernet and 11 LDCs on Fast Ethernet, the event building and the ROOT I/O formatting was done on 13 GDCs on Gigabit Ethernet, and eventually the full events were shipped to the CASTOR mass storage system. The latter consisted of 12 tape devices, hence not more than 120 MB/s on average can be achieved. Much effort was dedicated to parameterize and optimize the load balancing of the CASTOR system.

The particular run in Figure 6 lasted 59 hours at an average rate of 87.6 MB/s (upper left plot), where in the first few hours 120 MB/s could be achieved since the data was stored only on disk and no migration to tapes took place. In total $3.6 \cdot 10^6$ events were collected, which amounts to 18.4 TB of data (upper right plot). Each of the 13 GDCs was building and ROOT I/O formatting events at a rate of 6.8 MB/s (lower left plot). This was creating a average of 6% user and 23.4% system CPU load (lower right plot). The traffic generated by the LDCs was proportional to the subevent size, ranging from 0.37 MB/s to 8.09 MB/s, which is depicted in Figure 7.



**Figure 6:** Aggregate/GDC performance, aggregate volume, and CPU load over time • DATE + ROOT I/O + CASTOR • 20x13 configuration • ALICE-like subevents.

**Figure 7:** LDC performance and CPU load over time • DATE + ROOT I/O + CASTOR • 20x13 configuration • ALICE-like subevents.

# 5 Conclusions and Future Work

The ADC III has achieved major performance milestones. The testbed and the software prototypes running on it delivered a stable performance for periods of up to a week. More than 110 TB of data was put into mass storage by an average rate of 86 MB/s, and event building was possible up to 556 MB/s. Numerous runs with different configurations were executed to study for example the influence of the number of LDC/GDC, and the subevent size as well as content.

All these performance measurements were obtained by *dateStat*, which is a simple online monitoring tool developed during the ADC III. It derives performance data at regular intervals by a C program, collects them by the DATE info logger, processes them by Perl scripts, archives them by a MySQL database, and displays them on the web by gnuplot/CGI scripts. Based on this prototype a new monitoring tool called *AFFAIR* ("A Fine Fabric and Application Information Recorder") is being designed and implemented in collaboration with the Rudjer Boskovic Institute and the Mathematics Department of Zagreb University, Croatia. The following main features are foreseen:

- Gathering performance information from all system components in a non-application specific and uniform manner by employing a protocol with low network overhead and by using a round-robin database.

- Processing and permanently storing performance information by using ROOT.

- Visualizing performance information by using a web server to provide a broad spectrum of views on the data acquisition system.

The next ALICE data challenge is planned to take place in the second half of 2002, aiming at higher throughputs (200 MB/s over the complete chain) on a larger setup with better hardware, newer operating system (probably Red Hat 7.2), and improved versions of in-house software (AliRoot, DATE, ROOT I/O, CASTOR, AFFAIR).

# References

[1] ALICE Collaboration, "ALICE - Technical Proposal for A Large Ion Collider Experiment at the CERN LHC", CERN/LHCC 1995-71, December 1995.

[2] J.P. Baud et al, "ALICE Data Challenge III", ALICE internal note, 2001, in preparation.

[3] CERN, "Report of the Steering Group of the LHC Computing Review", CERN/LHCC 2001-004, February 2001.

[4] R. Brun, P. Buncic, F. Carminati, A. Morsch, A. Rademakers, "The ALICE Offline framework, status and perspectives", Proceedings of the CHEP2001 conference, Beijing, China, September 3-7, 2001.

[5] CERN ALICE DAQ Group, "DATE V3.7 User's Guide", ALICE internal note 2000-31, February 2001.

[6] ROOT team, "ROOT User's Guide V3.1", June 2001.

[7] O. Barring, J.P. Baud, J.D. Durand, "CASTOR project status", Proceedings of the CHEP2000 conference, pages 365-369, Padova, Italy, February 7-11, 2000.

[8] CERN IT Division, "The Performance and Exception Monitoring (PEM) project", *http://proj-pem.web.cern.ch/proj-pem*.