# Online Computing Hardware Replacement Policies for the LHC Experiments

# 1 – Introduction

**2004 Agreement**

In 2004 an agreement[1] was reached between the LHC experiments and the Scrutiny Group (SG) on a policy for the replacement of online computing hardware. The essential elements of this agreement are as follows:

- disks and file servers are replaced after three years;
- processors are to be replaced typically after four years;
- central network switches are replaced after five years.

The experiments have since based their M&O A budget requests for the replacement of online computing hardware on this policy.

**Request from SG and RRB**

At the meetings of the Resources Review Committees (RRBs) in autumn 2010, the Committee requested a review the policies for online computer replacement in the light of the experience gained operating the experiments with beam. The experiments' budget requests for 2012 and beyond were to be based on a new updated policy.

**Working Group**

The present paper has been prepared by a Working Group with the following composition:

- J. de Groot      Convener
- D. Francis      ATLAS
- B. Jost      LHCb
- F. Meijers      CMS
- B. Panzer-Steindel      CERN IT
- E. Radicioni      TOTEM
- P. Vande Vyvre      ALICE

---

[1] http://committees.web.cern.ch/Committees/LHCRRB/SG/SG_General/online_all_4_reply_sept04.pdf

and:

http://committees.web.cern.ch/Committees/LHCRRB/SG/SG_General/online_sept04.ppt

# 2 - Practice to Date

**LHC Schedule**

The last few years have been characterized by frequent changes in the LHC schedule:

- Startup of the LHC in September 2008, interrupted by a technical malfunction followed by approximately one year of repairs
- Startup of the machine in November 2009 followed by rapid progress towards collisions at injection energy. Decision to run the LHC in 2010 at 2 * 3.5 TeV, half the LHC design energy.
- Following a short winter shutdown, first collisions at 2 * 3.5 TeV are recorded in March 2010. The rest of 2010 saw an exponential increase of the luminosity, eventually surpassing the goal for 2010 of $10^{32}$ cm$^{-2}$s$^{-1}$ by a factor of two. A successful one-month heavy ion run took place at the end of the year as planned.
- In February 2011 the decision was taken to continue LHC running at 2 * 3.5 TeV in 2012 instead of the shutdown initially fore seen.

It is noted that during the extended shutdown periods, the experiment typically continued to take cosmics data for alignment and calibration purposes as well as for tuning of the online and offline software.

Planning of the capacity of the online computing system in such a volatile environment is obviously difficult and so is, a fortiori, planning the replacement of the online computing hardware.

**ALICE**

The full capacity of the ALICE online computer system was reached in 2010 in preparation of the heavy ion run. The online systems installed do not include spare capacity.

ALICE has followed the 2004 model concerning the M&O budget for online computing replacement with the following categories of equipment:

- disks and fileservers 3 years;
- other PCs 4 years;
- central switches (Local and Storage Area Networks): 5 years;
- peripheral switches: 4 years.

This budget has been used to:

- replace some equipment when obsolete, broken after warranty, or not supported by new operating system versions;
- purchase spare parts of custom hardware (DDL and D-RORC);
- purchase the maintenance contracts of the central switch purchased from the IT network frame contract (Force 10);
- purchase the maintenance of a commercial software package (cluster file systems from Quantum or Windows server);

- DB administration paid to IT Dept.

The PCs replaced and still operational have been used as spares or for testing till 2010. Some have been retired at the end of 2010 due to the transition to the new operating system version (SLC5 64 bits) in the DAQ.

The transition to a 64 bit version of Windows could not yet be done in the DCS due to the restrictions imposed by the commercial software (OPC, PVSS). Newer hardware is very often lacking drivers for older operating systems. Some of the old computers are therefore kept running until the transition to 64 bits and the replaced machines are used as spare.

The provisional M&O budget for future years has been adapted in 2010 to the LHC schedule foreseen at that time: a budget decrease the year preceding the shutdown (2011) and an increase at the end of the shutdown (2012).

ALICE has a support contract for the central switch that is included in the CERN/IT blanket purchase agreement with the provider (Force10). A new contractor has been selected to provide this type of equipment to CERN. The large network router will reach 5 years in 2011. It would be unpractical to replace it this year or next year. The situation will be reviewed in 2013: the equipment could either be replaced or be used for another few years if a new warranty agreement is established with the provider. The future M&O model should be flexible enough to accommodate such cases.

**ATLAS**

Due to the reduced operating parameters of the LHC machine, i.e. energy and luminosity, the ATLAS Trigger and Data Acquisition (TDAQ) system has to date not yet been fully deployed. The completion of the TDAQ system to operate at design luminosity and energy will, as foreseen, be funded with CORE construction funds.

ATLAS has followed the 2004 model for online computing replacement. All processors, with few exceptions, are replaced typically after four years. A fraction of processors have been replaced after five years of operations. These nodes provide critical functionality and are, in some cases, unique within the TDAQ system. Their technical specification is for five years of operation. An example of this online computing element is a commercially available dedicated fileserver.

With respect to networking, the initial large central switches will reach 5 years of operation this year. These large central switches are purchased under the CERN blanket purchase agreement with Force10. A recently re-tendering of the blanket purchase resulted in the selection of new supplier. Their replacement has a major impact on TDAQ system availability, therefore, subject to CERN IT purchasing a new maintenance contract with Force10, their replacement could be delayed to take into account the LHC shutdown schedule.

A second type of network switch, so called 'pizza box switches', are replaced after three years of operation, again in line with the agreed 2004 model.

The M&O A budget is also used to purchase:

- the maintenance contract for the network switches as part of the CERN IT blanket purchase with Force 10;
- the software license for the maintenance of commercial software, e.g. Spectrum.

In 2010, in an attempt to take into account the LHC schedule foreseen at that time, i.e. shutdown end 2011, part of the requested M&O A budget was deferred by one year from 2011 to 2012.

**CMS**

M&O budget requests generally use as a starting point the replacement strategy outlined in the 2004 document, taking as period for replacements:

- mass storage: 4 years;
- PCs typically 4 years (HLT 3 years and DAQ controllers and EVB nodes 5 years);
- network for Event Building and HLT: 5 years.

However, the actual budget requests for online computing have been substantially reduced compared to this model, because replacements have been postponed for practical reasons or to adapt to the changing LHC schedule. In particular, the budget request presented to the 2010 RRB meetings assumed an LHC shutdown in 2012 and replacements foreseen for 2011 were delayed.

The M&O material funds have been used for:

- maintenance contracts;
- purchase of spares and spare parts;
- purchase of replacements.

**LHCb**

The LHCb model for replacement of online computing items (without the LAN system, which is subject to an extended guarantee- like system), is based on the acquired experience of recent years. The four main ingredients are:

- buy components as late as possible;
- exchange components only when broken;
- budget for a provision of ~10% of the farm value in M&O Cat. A;
- an assumed decrease of cpu cost at constant performance of 12% per year.

The LHCb online farm is modular, flexible and designed to have a high degree of redundancy, featuring full dynamic allocation and de- allocation of nodes. Individual modules can therefore be swapped or changed without any effect on the data taking. The computing power of the system is designed to be always higher than what is required at any given instance. This allows the possibility to exchange elements in a transparent and backward compatible way.

It is further noted that LHCb, at the end of the 2010 run, was running at the maximum luminosity the experiment can accommodate. The proposed upgrade of the experiment around 2017 will address this issue and will see reading out the whole detector at 40 MHz and analysing each event in a trigger system implemented in software.

**Totem**

TOTEM has a small online system compared to the other experiments. TOTEM does not, at present, have an explicit replacement plan for the online computer hardware.

The online system is composed of the following parts:

- DAQ controller PCs: same order as CMS (2006/2006). These computers are located in the same counting room as the CMS PCs, and it is reasonable to assume that these will follow the replacement strategy of CMS.

- Event builders and mass storage: purchased in 2007. For this part of the equipment, at the time of procurement we considered a lifetime of 4 years.

- Event building network: purchased in 2007, with foreseen lifetime of 5 years.

Given the present LHC schedule, Totem envisages at least a partial replacement of the event building and storage equipment at the end of 2011.

# 3 - Lessons Learned

**Budgeting**

A difficulty encountered by the experiments is that planned budgets are determined by the original installation scenarios of the online systems. To be practical, the physical replacement of the equipment on the other hand should in some cases be determined by the shutdown and running periods of the LHC machine.

Strict annual budgeting without the possibility to carry over a substantial fraction of the budgets from one year to the next complicates the process and leads to an inefficient use of the funds allocated.

**LHC Operating Model**

The 2004 document was written with the LEP operating model in mind. The LHC operating model appears more like 1-2 year shutdown periods and 3-4 year running periods with short winter stops. This necessitates flexibility in the replacement scheme. Replacements of a large amount of equipment can be done more efficiently during shutdown years. This does not undermine the validity of the replacement policy.

**Failure Rate**

Experience has been gained with failure rates of equipment. For the year 2010, CMS PCs had a significant intervention rate (~5%, depending on the model). Power cuts lead to a substantial number of failures. The interventions probably still correspond to the constant failure rate part of the "bathtub" curve and no prediction can be given when we enter the "end of life wear-out" domain. It is indicative that PC manufacturers typically do not offer warranty beyond 3 years.

**Obsolescence of Equipment**

Some experience has been gained with obsolescence of equipment:

- The standard PC I/O bus has evolved from PCI_X to PCIe. ALICE have therefore produced a new generation of custom hardware (D-RORC) for PCIe and purchased PCs supporting both I/O busses.

- The transition from 32 to 64 bits has forced replacement of some PCs not supporting 64 bits.

More cases of obsolescence will appear in the future. The future M&O model should allow addressing them.

**Incompatibility of Equipment**

Although there are cases where equipment can simply be replaced with an up-to-date model, there are situations where a system wide approach is necessary. Replacement of PCs with up-to-date models can have an effect on the core switching network, such as a change from 1 Gbps to 10 Gbps ports. The same holds for replacement of I/O interfaces based on the PCI-X bus.

# 4 - Recommendation

**Replacement Policies**

It is proposed to retain the parameters of the system put in place in 2004 with the following revisions:

- Orient the budget towards a flattened and multi-year budget with the possibility of deferring purchases, leading to savings for some years. This will in particular allow the physical replacement of the equipment to be determined by the shutdown and running periods of the LHC machine and/or other unforeseen events;

- Include a budget line for the maintenance of custom hardware.