

Experience with procuring, deploying and maintaining hardware at remote co-location centre

O Barring¹, E Bonfillou¹, B Clement¹, M Coelho Dos Santos¹, V Dore¹, A Gentit¹, A Grossir¹, W Salter¹, L Valsan¹, A Xafi¹

¹European Organization for Nuclear Research, CERN, CH-1211 Geneva 23. Switzerland

E-mail: olof.barring@cern.ch, eric.bonfillou@cern.ch, benoit.clement@cern.ch,
miguel.coelho.santos@cern.ch, vincent.dore@cern.ch, alain.gentit@cern.ch,
anthony.grossir@cern.ch, wayne.salter@cern.ch, liviu.valsan@cern.ch,
afroditi.xafi@cern.ch

Abstract. In May 2012 CERN signed a contract with the Wigner Data Centre in Budapest for an extension to CERN's central computing facility beyond its current boundaries set by electrical power and cooling available for computing. The centre is operated as a remote co-location site providing rack-space, electrical power and cooling for server, storage and networking equipment acquired by CERN. The contract includes a 'remote-hands' services for physical handling of hardware (rack mounting, cabling, pushing power buttons, ...) and maintenance repairs (swapping disks, memory modules, ...). However, only CERN personnel have network and console access to the equipment for system administration. This report gives an insight to adaptations of hardware architecture, procurement and delivery procedures undertaken enabling remote physical handling of the hardware. We will also describe tools and procedures developed for automating the registration, burn-in testing, acceptance and maintenance of the equipment as well as an independent but important change to the IT assets management (ITAM) developed in parallel as part of the CERN IT Agile Infrastructure project. Finally, we will report on experience from the first large delivery of 400 servers and 80 SAS JBOD expansion units (24 drive bays) to Wigner in March 2013.

1. Introduction

1.1. Background

In 2006, in the beginning of the ramp-up of the CERN [1] IT computing facility as "Tier-0" in the Worldwide LHC Computing Grid (WLCG) [2], an extrapolation of the LHC computing requirements until 2020 revealed that the electrical power required for running equipment would already in 2010 exceed the 2.5MWatts available. A refined estimation, based on more recent hardware the following year, adjusted that prediction to 2012-2013. In the period 2008-2009 the option of building a new data centre on CERN premises was explored but finally abandoned due to the lack of funding. The possibility of deploying computer containers was also investigated but costs for the additional infrastructure (transformers, cooling) were considered too high. In parallel some hosting and co-location facilities in CERN member states had become interested in the problem and submitted their



proposals, which led to the decision in 2011 to launch a competitive tender for remote co-location. Sixteen proposals were received and evaluated. The contract was awarded to the Wigner Data Centre[3] in Budapest and the signature took place in May 2012.

1.2. Wigner Data Centre in Budapest

The data center is a new facility in a refurbished existing building. A total surface of 1100m² is dedicated for computing and split into 4 independent blocks of 275m² each. Three of the four blocks will be used by CERN. Each block has six rows of 21 racks. Two independent and fully redundant power feeds allow for an average equipment load of up to 10kW per rack. The UPS (Uninterruptable Power Supply) and diesel coverage for all IT equipment load and cooling consists of 3 UPS systems per block: one for each of the two IT equipment power feeds and one for cooling. The first block was ready for deployment in early 2013 while the remaining two dedicated to CERN became available in mid-2013 and will allow for a maximum power of 2.7MWatts to be provided to CERN. In-row cooling units are used with N+1 redundancy per row (N+2 per block) and the N+1 chillers function with free cooling technology (below 18°C). The estimated PUE (Power Usage Effectiveness) [4] is 1.5.

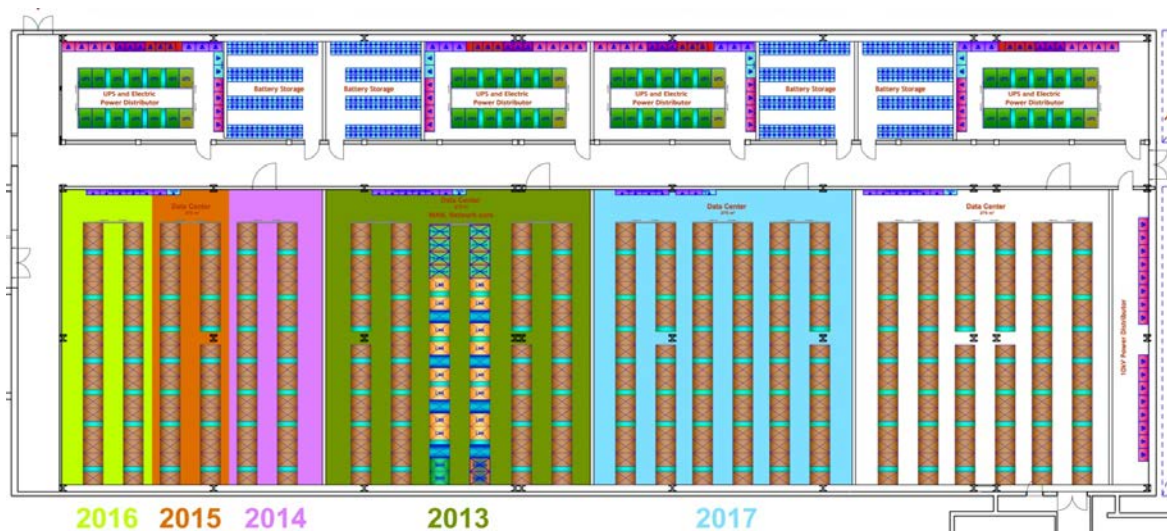


Figure 1: layout of the Wigner data centre machine rooms. The availability of the three blocks for CERN usage is indicated.

2. Preparation

In parallel with the tender for remote co-location, a number of processes and tools related to the hardware procurement, delivery handling and deployment were reviewed:

- Delivery requirements for suppliers
- Hardware handling from reception to rack mounting and cabling
- Stock management of spare parts (related to on-site maintenance below)
- Inventory of computing assets
- Network registration
- Burn-in and performance testing
- Production deployment
- Remote console
- On-site maintenance

In the following subsections the following four processes will be described in detail: delivery requirements, network registration, burn-in and on-site maintenance.

2.1. Delivery requirements

The delivery requirements for the suppliers are partly specified upfront in the tender documents and partly as part of the purchase order. Other than the delivery date (usually 10 weeks for large orders), standard requirements used to be as follows

- full batch must be delivered with an identical hardware configuration. In particular, the manufacturer, model, type, revision number, firmware version, etc. of each hardware component, as reported by the device, must be identical between units as well as within each unit;
- the purchase order includes an attached list with required BIOS settings, e.g. boot list order, PXE (Preboot eXecution Environment) enabled, restore on AC power loss.

Despite this, suppliers failed to fulfil the requirements, in particular the BIOS settings, which caused problems when hardware was powered on. Fixing the wrong settings sometimes required physical access to the machine.

In order to reduce the risk of receiving hardware with wrong settings, a requirement for remote console access to the first produced system while it is still at supplier's premises was added to the tender specification. The console connection is used to configure the remote system's BIOS exactly as desired by CERN. The supplier is expected to clone the settings onto all other systems from the same batch.

A few other additional requirements were recently added for facilitating remote operations:

- BIOS / UEFI, BMC, storage drives, RAID controller and network card(s) firmware files must be provided at delivery time;
- command line tools must be provided at delivery time to flash in an automated fashion these BIOS / UEFI and firmware files;
- a CERN defined "Asset Identifier" must be written into the Field Replaceable Unit (FRU) of the Board Management Controller (BMC) of each system:
 - CERN Contract ID is recorded in the 'Product Asset Tag' (PAT) field;
 - serial number is recorded in the 'Product Serial' (PS) field;
- command line tool allowing information to be written in the FRU fields of the BMC;
- the above tools must run locally under the (64-bit) Red Hat Enterprise Linux (RHEL) and Scientific Linux CERN (SLC) version specified by CERN without requiring recompiling or specific drivers for different kernel versions.

Since the asset number is used to trace the system over its entire lifecycle, the serial number must be different from the serial number of any component inside the enclosure or system unit. Thus, replacing any of the components will not change the asset number used for tracing the piece of equipment.

The BMC FRU settings required above are also checked from the remote console. The Asset Identifier 'PAT-PS' is also required on a barcode sticker at the rear of each system unit, on each enclosure, as well as on the cardboard boxes. Example:



Figure 2: example of the custom barcode required on each system unit and enclosure. The first part before the dash ('-') is the CERN contract identifier and the second part is the vendor serial number. The same information must be burned into the FRU of the BMC of the system.

2.2. Network registration

The network registration used to be a tool assisted manual process executed by the system administrators. The process was error-prone and the tools did not support the registration of machines with multiple physical network interfaces well, which started to become standard with the deployment of 10 gigabit Ethernet. Another issue was that the process relied on the Quattor/Cdb configuration system used by CERN since 2004 but that is now being phased out and replaced by Foreman and Puppet as part of the Agile Infrastructure project launched in 2011.

In order to automate the registration of new machines the following networking changes have been implemented:

- the network switch is configured to provide private (and temporary) IP addresses to unregistered MAC addresses via DHCP;
- the server can use the private IP to PXE boot a live (in-memory) image;
- a new API call was added to the interface to the local area network database (LANDB) at CERN, taking an IP address and a list of MAC addresses as input and returning as response the network switch and port to which the interface is connected to.

The new functionality closely resembles how CERN visitors register WiFi access for laptops and other devices.

Once the node knows the network switches and ports its different network interfaces are connected to it can make a standard call to LANDB to complete the registration of server hostname, MAC addresses, serial number, network service and port. Once the server is properly registered, the network configuration can be reinitialized to acquire the permanent IP(s).

In order for the network registration to run unassisted when the machine is powered up for the first time, the live image must contain a robust method for generating a unique host name. The approach chosen was to build the name based on the FRU information in the BMC (see section 2.1). Unfortunately it turned out that the trivial choice of using the 'PAT-PS' is incompatible with the host name string length limitation in the Windows NETBIOS name, which is 15 characters (a hard limit). There is a possibility to associate a longer name in Active Directory but the NETBIOS name needs to be unique and most of the system tools do not distinguish between NETBIOS names and Active Directory names and hence the risk for encountering recurrent problems when these two entries are not the same. This restriction is in all versions of Microsoft Windows, including Microsoft Windows Server 2012.

The compromise chosen was to concatenate the number part of purchase order stored in the PAT with a locally generated random number. The fixed part for the purchase order number is prefixed with the letter 'P' (Physical node). Example:

Position	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Example	'P'	'0'	'9'	'4'	'7'	'2'	'9'	'6'	'4'	'7'	'5'	'3'	'2'	'7'	'9'
	'P'	Contract CERN doc number								Random decimal number					

Obviously as the random part is generated locally there is a small risk for collisions when there are many servers from the same batch (i.e. purchase order). However, in this case only the first server registering the name will be successful while the second server will get an error and will retry with a different random number.

2.3. Burn-in and performance testing

Once the network registration described in the previous section has completed, the live image launches the burn-in and performance tests. The burn-in phase includes:

- Memory (memtest) and CPU (burnK7 or burnP6, and burn MMX) endurance tests
- Disks endurance tests (badblocks)

All outputs from the tests are logged to a central Splunk [5] server, which provides convenient query interfaces and dashboards for monitoring the progress.

Once the burn-in tests have finished, which can take 1-2 weeks depending on the amount of local memory and disk space, the live image continues with the performance testing. The results from the standard benchmarks HEP-SPEC06 [6] and FIO [7] and/or iohome [8] are uploaded into Splunk, which allows for comparing the values with the expected ones (e.g. from measurements on sample units provided with the tender bid). It also allows for comparing the systems within a batch to detect possible systematic issues. A real example came up when the CPU performance was measured on the storage front-end servers. It turned out that the average performance across the whole batch was significantly lower (~20%) than what had been measured on the sample units benchmarked in the tender evaluation. The problem was investigated and traced to a non-optimal setting of the CPU scaling governor.

2.4. On-site maintenance

The on-site maintenance used to be part of the warranty service of the supply. The vendors were asked to provide 3 years' on-site maintenance. The most common repair target was next business day. As the supplier could be from any CERN member state, it was allowed to sub-contract the on-site maintenance to a company in the local area. Unfortunately this model did not work very well and most suppliers missed the target in more than 50% of the interventions. The quality of the repairs also became an issue when some local service companies used technicians without proper qualifications. The situation was complicated for CERN because the on-site service was part of the supply contract, which cannot easily be adjudicated on a best-value-for-money basis. The maintenance contract was between the supplier and the local service company and therefore not under CERN's control.

Independently of the plans for remote co-location, it was therefore decided to overcome the problems described above by taking better control over the on-site maintenance service. To this end the on-site maintenance repairs were added as part of a tender for renewing the 24/7 operator service contract. The new contract was established in January 2012 and is used for all new deliveries since then.

In order to accommodate for the new on-site maintenance service, the warranty requirements for the supply of servers and storage had to be modified:

- bidder must estimate an initial stock of spare parts and/or standby units sufficiently large to cover expected component failures over a two-month period. For insurance reasons the stock must be owned by CERN and thus part of the supply;
- the bidder must agree to cover all transport charges for failed and replacement spare parts. The nominal frequency is set to monthly but could be more frequent if the estimated stock was insufficient;
- the bidder must agree that the repairs will be performed by CERN or a contractor appointed by CERN.

The service has successfully run for more than 18 months at CERN. It has greatly simplified and improved the quality of the on-site maintenance repairs. The same service was also part of the tender for remote colocation and is therefore now also being established at the Wigner Data Centre.

3. First deployment

The first block for CERN at the Wigner Data Centre was ready on schedule in January 2013. In March 2013 deliveries of in total 400 servers and 80 JBOD arrays (6PB raw disk space in total) were made. The Wigner personnel verified the deliveries and scanned all the CERN asset identifiers on the cardboard boxes and provided CERN with an inventory that could be matched with the detailed configuration sheets the equipment suppliers had provided. Thereafter the boxes were unpacked and the equipment installed in the racks and cabled according to a detailed rack schema agreed between CERN and Wigner. When the servers and JBODs were powered up, the automated network registration and burn-in process (see description in Sections 2.2 and 2.3) started. All logs were

forwarded to a Splunk server at CERN, which allowed for monitoring of the progress and eventual failures on a dashboard. The full process took about 2 weeks to complete with 99% success across the 400 machines and 80 JBODs. The few failures were related to:

- cabling, i.e. wrong port cabled or cable not fully plugged in;
- faulty switch ports: one network switch had to be replaced at Wigner;
- failed components: one mainboard, 9 hard disk drives, 2 memory modules and one power supply

4. Status of remote operation and future plans

The successful execution of the delivery handling, installation and cabling of hardware described in previous sections demonstrated that the hardware management workflow is functioning properly.

The focus has been on streamlining the remote operation. Bi-weekly phone meetings have been sufficient for coordinating unavoidable facility changes and maintenance. CERN personnel have also gone on site to provide detailed training when required, for instance for handling repairs of the deployed hardware. The operations personnel at Wigner have been provided with access to ServiceNow [9], the ticketing system selected by CERN for the implementation of Information Technology Infrastructure Library (ITIL) [10] best practice procedures, notably the Event Management process described in another paper to this conference [11]. Wigner staff is also been given access to the Infor EAM [12] system used by CERN for managing the stock of spare parts. This allows for handling of hardware maintenance repairs with the same workflow at Wigner as locally at CERN.

Based on the positive experience from the first year of operation, the current plan is to deploy 90% of the new capacity acquired in 2014 at Wigner.

5. Conclusions

The use of a remote co-location at the Wigner Data Centre in Budapest is CERN's strategy to grow the Tier-0 computing capacity. The preparation for remote operation of the server and storage equipment necessitated adjustments to processes and workflows some of which have been described in this paper, like the automated network registration and burn-in testing. The changes have also been beneficial for the local operation.

The successful first production deployment of 400 servers and 80 JBOD arrays with 6 petabytes of disk storage has allowed for confidently building up the operation routine between the two sites where the recent establishment of the hardware repair service is the latest example. Based on the positive experience CERN has decided to deploy 90% of its new computing and storage capacity, acquired in 2014, at the Wigner Data Centre.

References

- [1] <http://www.cern.ch>
- [2] Worldwide LHC Computing Grid <http://wlcg.web.cern.ch/>
- [3] Wigner Data Centre <http://www.rmki.kfki.hu/en>
- [4] Power usage effectiveness (PUE) http://en.wikipedia.org/wiki/Power_usage_effectiveness
- [5] <http://www.splunk.com/>
- [6] HEP-SPEC06 (<https://hepix.caspar.it/benchmarks/doku.php?id=bench:howto>) is a computational benchmark derived from the SPEC CPU 2006 industry-benchmark suite for computational performance (<http://www.spec.org>)
- [7] Flexible IO tester (FIO) benchmark suite <http://www.freecode.com/projects/fio>
- [8] IOzone Filesystem Benchmark <http://www.iozone.org/>
- [9] <http://servicenow.com>
- [10] <http://www.itil-officialsite.com>
- [11] R. Alvarez Alonso et al., Migration of the CERN IT Data Center Support System to ServiceNow, CHEP 2013 proceedings
- [12] <http://www.infor.com/solutions/eam/>