

# Cray Hybrid XC30 Installation – Facilities Level Overview

Colin McMurtrie, Ladina Gilly and Tiziano Belotti  
CSCS – Swiss National Supercomputing Centre  
Lugano, Switzerland  
Email: {colin; lgilly; belotti}@cscs.ch

**Abstract**—In this paper we describe, from a facilities point of view, the installation of the 28-cabinet Cray hybrid XC30 system, Piz Daint, at the Swiss National Supercomputing Centre (CSCS). This system was the outcome of a 12 month collaboration between CSCS and Cray and, as a consequence, is the first such system of its type worldwide. The focus of the paper is on the site preparation and integration of the system into CSCS’ state-of-the-art HPC data centre. As with any new system architecture the installation phase brings challenges at all levels. In order to achieve a quick turnaround of the initial bring-up it is essential to ensure that the site design is flexible enough to accommodate unforeseen variances in system environmental requirements. In the paper we detail some of the challenges encountered and the steps taken to ensure a quick and successful installation of the new system.

**Keywords**-XC30; hybrid; HPC data centre; Facilities Management

## I. INTRODUCTION

### A. CSCS Data Centre Overview

The CSCS data centre building comprises three floors. The underground floor is where the electricity and cold water enters the building and is prepared for further distribution. The ground floor (termed the *Installation Deck*) houses all the secondary power distribution units (PDUs) as well as the system-dedicated cooling loops. In order to ensure that this space can adapt to future requirements and to allow ample room for all equipment on this floor it is 5.5m high.

The Machine Room is located on the first floor, directly above the Installation Deck, and provides a contiguous space devoid of pillars. A minimal 3-layered structure of I-beams, raised floor support pedestals and stringers, and raised floor tiles separate the Installation Deck from the Machine Room. Figure 1 shows various features of the Installation Deck including examples of the system-dedicated cooling loops, PDUs and the I-beam structure that supports the Machine Room floor. Note the provision of raised walkways which form part of a support structure for all electrical equipment, including PDUs and the building management equipment as well as the pumps and electrical equipment associated with the dedicated coolings loops. This structure ensures that this equipment (and personnel) is above any water in the event of a widespread water leak. Finally water sensors are placed at strategic locations throughout the installation deck in order to alert facilities staff in the event of water leaks.



Figure 1. Image of the Installation Deck showing the supporting I-beam structure, PDUs and system-dedicated cooling loops.

The building design ensures great permeability between the Installation Deck and the compute floor, making it very easy to bring the desired resources to any location in the Machine Room. It also enables all facilities equipment to be placed in the Installation Deck (e.g. secondary PDUs) thus maximising the floor space available for computer equipment in the Machine Room and shortening cable lengths between PDUs and the computer systems. Furthermore this design allows large installation teams to work comfortably in the Installation Deck in order to prepare the electrical and cooling distribution for a new system, with no impact on operations in the Machine Room. This significantly reduces the risk of contractors accidentally affecting systems and services in operation and allows multiple teams to work in parallel without impeding each other. This latter point is particularly important when installing a new system because the vendor engineers can work freely in the Machine Room, essentially unobstructed by any facilities work taking place on the floor below.

Electricity enters the CSCS building at 16kV and, via step-down transformers and primary PDUs in the underground floor, is run to the PDUs in the Installation Deck at 400V. The entire electrical distribution is monitored by a dedicated electrical network management tool, supplied by Leicom, which is part of the data centre building management infrastructure. The state-of-the-art Leicom system is extremely useful because it allows fine-grained control and monitoring of individual components within the electrical infrastructure, right down to the setting of alarms. It is also possible to remotely manage individual circuit breakers but, for safety reasons, this functionality is not used. The

tool is also extremely useful in its monitoring and diagnostic capabilities; voltages, currents, power and energy etc are monitored and tracked over time. This latter point is particularly useful when assessing the duration and impact of micro-outages, for example, but is indispensable when measuring whole-system energy consumption as in the case of a Top500 and Green500 submission.

More information about the CSCS Data Center can be found online on the CSCS website (see [1]).

#### 1) Lake Water Cooling and Cooling Loop Overview:

The water from lake Lugano that is used to cool the entire capacity of CSCS enters the building at around 7°C, year round. By means of primary heat exchangers it cools the water from the internal primary cooling loops and then returns to the lake; in this way the lake water cooling can be considered as open-loop. There are two internal closed primary cooling loops that operate in different temperature ranges: the low-temperature cooling loop (denoted the TTN loop) is designed for 9°C supply and 17°C return while the medium-temperature cooling loop (denoted the MTN loop) is designed for 21°C supply and 29°C return. For energy efficiency reasons the primary heat exchangers for the TTN and MTN are in series on the lake water cooling loop so that the outlet of the TTN heat exchanger (at 17°C) is fed into the inlet of the MTN heat exchanger. If for some reason the MTN loop needs a colder supply temperature this is achieved by the use of a by-pass, on the lake water loop, which enables controlled mixing of lake water from the inlet to the outlet side of the TTN heat exchanger. This whole process maximises the cooling capacity for the pumping energy required to get the water from the lake to CSCS; however use of the TTN mixing by-pass must be kept to a minimum. All cooling loops for the CSCS computer systems are created by adding heat exchangers and closed loop distribution to either the TTN or MTN loops; in some special cases the system cooling loop can be created by tapping directly into the TTN or MTN primary loops. Finally water returning to the lake is controlled to have a maximum temperature of 25°C, in order to minimise the impact on the lake ecology. Figure 2 shows a high-level schematic of the cooling infrastructure.

2) *Machine Room Ambient Air:* For energy efficiency reasons the entire volume of the CSCS machine room air is not cooled. Rather, systems are required to be room-neutral by the provision of vendor-supplied direct cooling or, when air-cooled, are inserted into one of a number of in-row cooling islands. The machine room environment is, however, monitored and maintained within the ASHRAE 2008 Revised Class 1 & 2 Operating Range [2]. This means that both temperature and humidity levels can fluctuate within quite wide limits.

### B. Facilities Management Overview of Hybrid XC30 System

1) *Cabinet Cooling and Electrical Supply:* All components within the XC30 cabinets are cooled by a stream

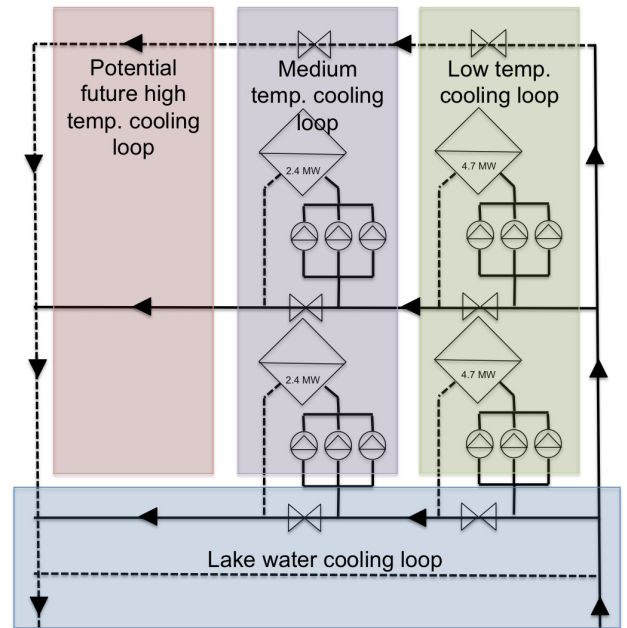


Figure 2. Schematic showing the layout of the low- and medium-temperature primary cooling loops in relation to the lake water cooling loop. Note the provision for a high-temperature cooling loop; this has not yet been implemented but was allowed for in the original data centre design.

of cool air. A key feature of the design is that the air flows horizontally, from left to right, through the cabinets that make up a row. Every cabinet has an air to water heat exchanger (aka radiator) on the downstream side of the cabinet (i.e. on the right when facing the cabinet) and thus each cabinet accepts two 2" water-pipe connections (one supply, one return). In order to control the flow of water within the in-cabinet heat exchanger, each cabinet is also equipped with an electronically controller water valve and control circuitry. Cabinet environmental data is monitored via temperature/pressure transducers in the supply and return water pipework and via sensors place in the airstream, downstream of the cabinet heat exchanger. There are 9 sets of these latter sensors, at the front, middle and rear of the cabinet, located in three vertical locations in the middle of each chassis. The control system of each cabinet acts independently to keep the airstream temperature on the downstream side of its heat exchanger at the temperature set-point by continuously controlling the water valve position and therefore the amount of water flowing through the heat exchanger.

Air is moved down the row via in-row blower cabinets, each containing 6 large fans. A blower cabinet is placed at the start and, optionally, at the end of each row as well as after every second cabinet (aka group) within the row. The cabinet control system operates the fans at 80% under normal conditions. If hot cabinet components (i.e processors) are detected, the fan speed is increased to 100%. In order to

move air down the row in a uniform manner, all fans within a row must operate at the same speed. Each blower cabinet uses 5.5kW of power when its fans are running at 100%. Note that there are no radiators within a blower cabinet.

Given that the blower cabinets do not contain heat exchangers and the fact that the in-cabinet heat exchanger is on the downstream side means that the first cabinet ingests air directly from the ambient machine room environment. If the conditions of the machine room environment warrant it (e.g. high humidity and/or temperature) Cray can supply an optional preconditioner cabinet which can be added upstream of the first blower of each row. The preconditioner cabinet contains its own air-to-water heat exchanger and water control valve and therefore requires its own water supply. Control of the water valve is taken care of by a connection to the cabinet control system of the adjacent full cabinet (i.e. the first cabinet in the row).

Each compute cabinet requires 2 three-phase supplies (either 400/230VAC, 125Amp, 50Hz or 480/277VAC, 100Amp, 60Hz) in a WYE+Neutral+Ground configuration. Hence, in our case where we use 50Hz supply, this means that there are 10 125Amp conductors running to every cabinet.

2) *Blade Layout and Design Features:* Each XC30 compute cabinet contains 3 chassis, each of which has 16 horizontally mounted compute blades. The configuration is such that there are 8 blades (numbered from 0 to 7) on the lefthand side and 8 blades (numbered 8 to 15) on the righthand side of each chassis. As mentioned above, the forced airstream flows from left to right over the blades within each cabinet and as a consequence it gets hotter as it moves over the active blade components (CPUs, GPUs, memory etc); air leaving the righthand blades is therefore hotter than air entering the lefthand blade.

Each blade has an identical layout with 4 nodes and one Aires ASIC per blade. Each node contains one CPU and one GPU and memory DIMMs. There is a design challenge to provide the same degree of cooling to these blade components in an airstream that gets hotter as it moves across the blade. Hence, in the case of the GPUs Cray and Nvidia design engineers employed heat sinks with differing fin counts. All GPUs in the nodes on all lefthand blades therefore have heat sinks with 13 fins, whereas Nodes 2 and 3 of the righthand blades have 20 fins and Nodes 0 and 1 have 30 fins. Figure 3 shows a schematic of the left and righthand blades and identifies 4 distinct node locations, A through D. As it turned out the significance of these 4 distinct node regions became important when it came to understanding the thermal characteristics of the system when running compute-intensive workloads that made heavy use of the GPUs and this is discussed in more detail below.

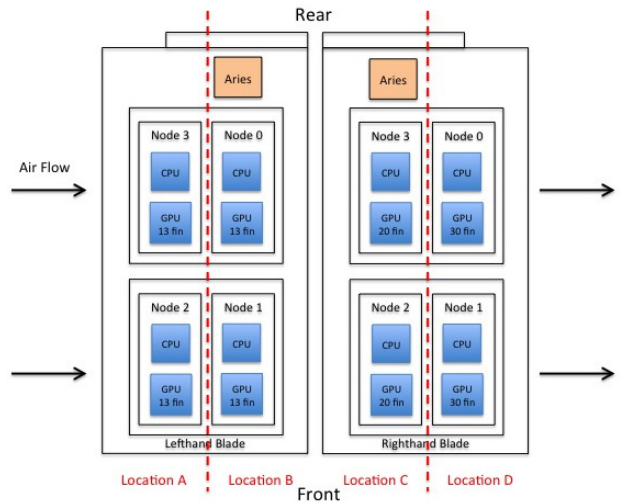


Figure 3. Schematic showing the node layout on left and righthand blades of the hybrid XC30 system. Note the different heat sink fin counts and direction of airflow which gives rise to 4 distinct thermal regions, denoted Locations A to D.

## II. DESIGN OF THE FACILITIES INFRASTRUCTURE FOR THE SYSTEM

### A. Secondary Cooling Loop Design

For the initial 12-cabinet non-hybrid installation in November 2012, which included an additional single-cabinet Test and Development System (TDS), we planned the facilities infrastructure using a preliminary version of the Site Preparation Guide provided by Cray ([3]). This non-hybrid system had been designed to accommodate relatively high inlet air and water temperatures and large  $\Delta T$  across the in-cabinet cooling radiators. In order to benefit from the energy savings this could provide the dedicated cooling loops for this 12-cabinet system were designed in a closed-loop configuration, with a supply temperature of 21°C and a  $\Delta T$  of 12C, resulting in return water at 33°C.

Two separate closed loops were built for this system, one for a row containing 8 cabinets plus the 1 cabinet TDS and one for a row of 10 cabinets but which initially contained only 4 cabinets. The cooling loops for each row were designed for a maximum heat load of 1.2MW and contained two recirculation pumps, configured as an N+1 redundant pair. The control system for these closed loops monitored the pressure drop across groups of cabinets within a row along with the supply water temperature and varied the pump speed (and therefore water flowrate) in order to keep these variables within the specified range. The secondary cooling loop for the row containing the 8+1 cabinets was connected to the medium-temperature (MTN) cooling loop whereas the secondary cooling loop for the row designed for 10 cabinets was connected to the low-temperature (TTN) cooling loop.

Given that the initial installation was designed for a total of 18 cabinets, when it came time to upgrade to the 28 cabinet hybrid system it was only necessary to add a third cooling loop with capacity for 10 cabinets. This secondary loop was also connected to the low-temperature (TTN) cooling loop. At the time the cooling infrastructure for this third row was designed, the hybrid CPU-GPU blades did not exist in any great number and exhaustive testing of multiple cabinets of them had certainly not been conducted, so that Cray and CSCS had to make assumptions about the operating temperatures of these new cabinets. As a consequence Cray's Site Preparation Guide was little changed from the version used for the non-hybrid system and therefore, as before, the planning engineers dimensioned the heat exchanger for 21°C inlet and  $\Delta T$  of 12C. As with the other two rows, a redundant pair of pumps were used and the pump speed was varied by the control system in order to keep the pressure drop across groups of 5 cabinets and supply water temperature within a specified range.

As will be seen below, the assumption that the hybrid CPU-GPU cabinets would have similar operating temperatures to the non-hybrid system turned out to be incorrect and this required last minute changes to the cooling distribution. Thankfully the flexibility of the cooling loop design ensured these changes could be made easily.

### B. Electrical Distribution

The preliminary Site Preparation Guide used in 2012 when planning the 12-cabinet non-hybrid system specified a maximum power requirement of 111kW (113kVA) per cabinet. In order to accommodate future upgrades electrical (and cooling) distribution were planned and built for a maximum power requirement of 115kW per cabinet. This power envelope proved to be more than adequate for the cabinets when equipped with the new hybrid CPU-GPU blades. All compute cabinets are connected to 400V utility power (no UPS) but, based on our experiences with the initial 12-cabinet non-hybrid system, the in-rack power supplies were known to have good ride-through characteristics when subjected to micro outages.

## III. INSTALLATION EXPERIENCE

### A. Pre-installation Data

The project plan agreed between CSCS and Cray involved the very early installation of a 3 cabinet hybrid system at Cray (dubbed *clogin85*) and fully populating the CSCS XC30 TDS (aka *Santis*) with hybrid blades. The *clogin85* system at Cray came online approximately 3 months prior to the start of the full system installation at CSCS, as soon as hybrid blades were available from Cray's hardware manufacturing division. The fully populated *Santis* system came online some weeks later, once the parts had been shipped to CSCS. By the time these systems were operational, planning for the third row cooling had been completed, site

preparation had started and long lead-time items like heat exchangers had been ordered.

As soon as the *clogin85* system was operational, staff from CSCS' Future Systems group gained access to it and started testing its functionality and stability. Cray also had their systems and applications teams look closely at the system, this being the first time the new architecture had been available for complete integrated system testing. The results, on the whole, were very promising and the systems showed promising performance and stability. One thing that became apparent however was that the GPUs in the 4 locations A, B, C and D in the cross-cabinet airstream (see Figure 3) showed different temperature profiles when subjected to a compute-intensive workload like HPL.

Figure 4 shows the time-series temperatures for 16 selected nodes within the middle chassis of the second and third cabinets of the system for one full run of HPL. Blades *c1s3* and *c1s11* in each cabinet (*c1-0* and *c2-0*) were chosen because, as noted above, blades in the XC30 cabinets are placed horizontally and as a consequence these blades are beside each other in Chassis 1, with *c1s11* downstream of *c1s3*; exhaust air leaving the heat sinks of *c1s3* is therefore directed onto the heat sinks *c1s11*. The data points on each plot come straight from the GPU SXM, as reported at semi-regular intervals by Nvidia's *xhpl* code.

As can be seen from the results, GPUs in Location 'B' are consistently hotter than the nodes in the other Locations, for the entirety of the HPL run. Specifically in the case of Node 462 and Node 461 (which are adjacent to each other on blade *c2-0c1s3*, with Node 461 downstream of Node 462) the difference in temperature is almost 15C. Cray gathered data about the relative temperatures of the GPUs in Locations A, B, C and D. Statistical analysis of these data showed that GPUs in Location 'B' were on average 10C hotter than all other GPUs within the system.

Furthermore the CSCS results in Figure 4 show that nodes on blade *c2-0c1s3* are consistently hotter than their corresponding counterpart on blade *c1-0c1s3*. The CSCS team inferred this to mean that the air stream was getting hotter as it moved from left to right down the row; in other words the in-cabinet radiator was not removing enough waste heat from the airstream before it entered the next cabinet.

The fact that cabinet exhaust air temperature was not being kept within the specified range was also supported by results seen on the *Santis* TDS system when this system was fully populated with hybrid blades in early September 2013. When running the compute intensive HPL code from Nvidia it showed extremely high component temperatures with some GPUs thermal capping which results in the GPU reducing its clock frequency to 365MHz, which is half the normal rate. This, in turn, had devastating consequences on the HPL performance. As described above, this cabinet was located at the end of the front row of the full XC30



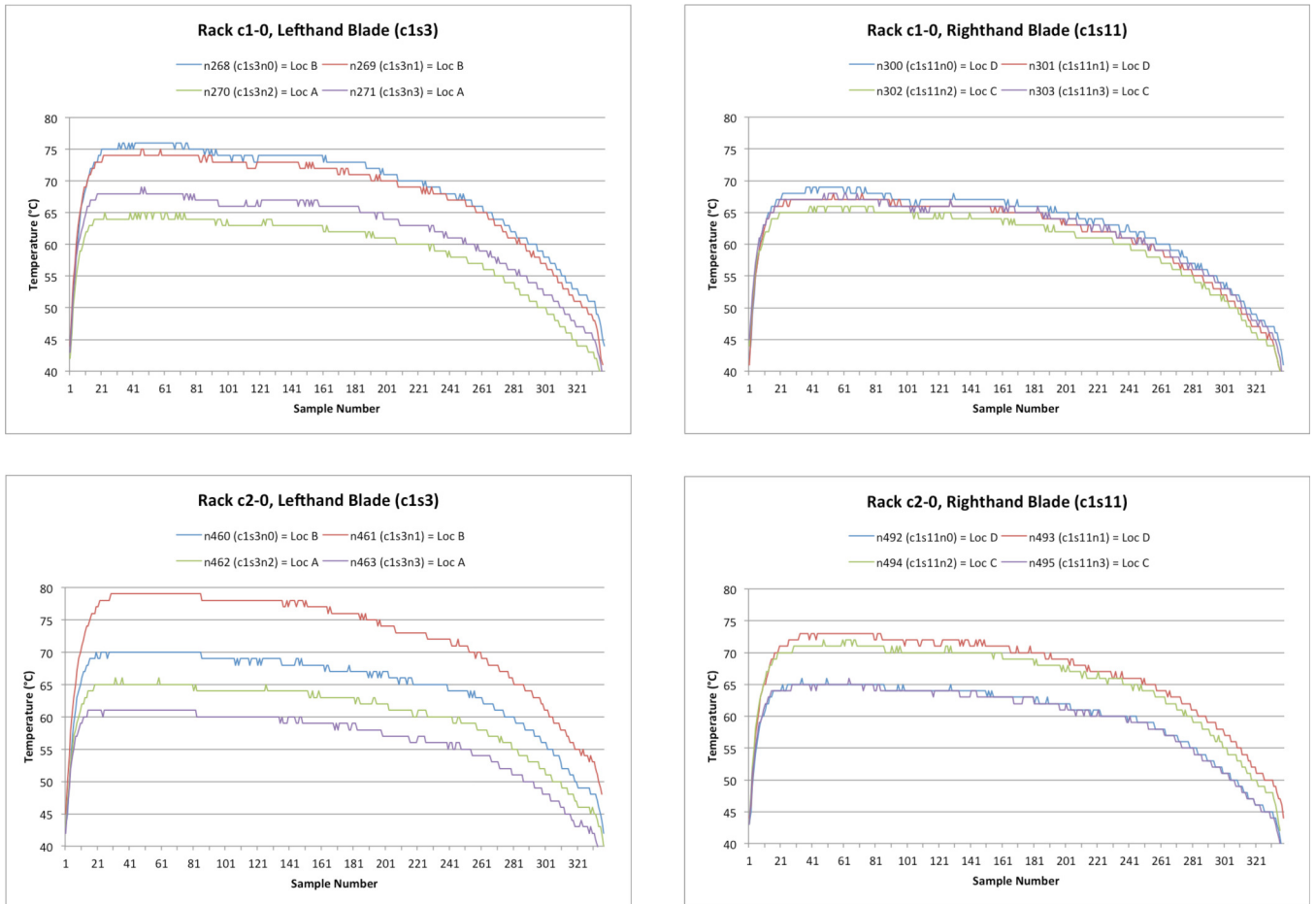


Figure 4. Plots showing the temperature of 16 nodes in the clogin85 system when running HPL. (Results courtesy of Gilles Fourestey, CSCS.)

system and it was discovered that the exhaust air from the other 8 cabinets in this row was around 28°C which, at the time of testing, was up to 9C hotter than the ambient machine room air temperature. Despite there being a 1 meter gap between the end of the 8 cabinet row and the TDS this was not enough to allow adequate mixing of the exhaust airstream and ambient machine room air so that the inlet air temperature was still well above the ambient machine room air temperature (up to 7C). This high inlet air temperature was enough to cause the extremely high component temperatures seen in the HPL runs (recall that the in-cabinet radiator is located on the righthand side of the cabinet so that it is downstream of the blades). This situation was easy to rectify because it was possible to redirect the airflow from the 8 cabinets away from the inlet of the TDS but it raised two important points:

- 1) The main system would require the optional start-of-row preconditioners;
- 2) Why was the exhaust air from the 8 cabinet row not room neutral?

The first point was a consequence of the fact that, as noted above, the machine room air temperature was kept within the ASHRAE 2008 Revised Class 1 & 2 Operating Range which means it can get as high as 27C. The testing on the hybrid TDS had shown that at these temperatures components within the first cabinet of each row would start to overheat.

The second point was discussed with Cray's Site Engineering team and they advised that the cabinet control software would need to be configured in such a way that the last cabinet of each row would cool the air back down to the machine-room ambient temperature. Cray also advised that a Field Change Order (FCO) had been issued which involved a change to the water flow directions within the XC30 cabinet. The FCO was designed to improve the efficiency of the cabinet radiator, thereby providing more cooling capacity for the same inlet water temperature and flow-rate. CSCS advised Cray to make this change to the TDS cabinet but, due to the availability of specialist staff in Europe and the nature of the change it was not possible to do this until the Cray Site Engineering staff were onsite to start the

installation of the main 28 cabinet system in early October 2013. In the meantime, testing continued with the Santis TDS and the clogin85 system at Cray in order to gain more familiarity, in all aspects, with the new system architecture.

### *B. Early On-site Testing and Facilities Changes*

Once Cray Site Engineering staff were at CSCS in early October, the FCO changes were made to the Santis TDS. However additional testing showed that, despite these changes and the cabinet controller setpoint changes, it was still not possible to get the exhaust air temperature to the Cray-recommended 19°C with a coolant water temperature of 16°C. In fact, exhaust air temperature under these conditions was more like 23°C with the cabinet water control valve 100% open. These results were somewhat worrying because the front row of the system was on the MTN cooling loop and hence its water temperature could not be brought below 16°C (and in fact to achieve even this temperature involved using the water bypass on the lake-side cooling loop which was not sustainable in the long term).

By the time these results were available it was mid-October and the main 28 cabinet system was close to initial power up for off-line hardware testing. Thankfully there was an easy fix for the problem with the front-row water temperature, namely to simply move its primary-side supply to the low-temperature (TTN) cooling loop. Thanks to the extreme flexibility of the Installation Deck layout, however, this change involved the construction of a few meters of piping and was easily achieved within a short space of time (little more than 2 days from design to installation); switching the primary-side supply from the middle- to low-temperature cooling loop was then done with the use of valves and hence no service interruption was necessary. Furthermore, no changes were necessary on the secondary side as the were of adequate size and capacity.

For the middle and back rows however the situation was a little more worrisome. The secondary cooling loops for these rows were already connected to the low-temperature (TTN) primary cooling loop. Furthermore, as described above, the heat exchangers on these loops had been sized for different flow-rates and, more importantly, different secondary-side temperatures so that the primary- to secondary-side temperature gradient was much larger in the original design for the 12-cabinet non-hybrid system. The net result was that the heat exchangers were too small to achieve the necessary secondary-side supply temperature for the new hybrid system and, given that there was a long lead-time for delivery of new heat exchangers (6 to 8 weeks), the only viable option in order to stay on schedule with the project and not jeopardise the submission of Top500 results or system acceptance was to remove the existing heat exchangers until new, appropriately sized units could be delivered and installed.



Figure 5. Photograph showing the pipework, prior to lagging, for one of the primary- to secondary-side heat exchanger by-passes, including cross-connect valve.

Work commenced immediately to remove the heat exchangers. New by-pass pipework was designed and manufactured off-site by the cooling distribution contractor in a matter of days and then installed in less than one day. This process was a little more disruptive than that of the front row changes because parts of the secondary-side loops had to be drained. As part of the by-pass design, and as a precaution, a cross-connect valve was installed between the supply and return pipework, the idea being that this valve could be manually set to allow a degree of supply-to-return mixing in order to prevent the secondary-side pumps from “fighting” with those on the primary side, given that now the two loops were directly connected. No other changes were necessary (the original secondary-side pumps were kept) but the building management control system had to be tuned to the new configuration. Figure 5 shows a photograph of one by-pass, prior to lagging. As can be seen the design was compact and well implemented.

### *C. Main System Bring-up*

With the quick turnaround of the cooling loop changes there was little impact on the installation schedule and the system went into off-line diagnostic testing and then online

testing relatively quickly; Cray engineers were able to focus entirely on identifying and eliminating faulty components within the system, as is normal for a new system installation.

Once the system was stabilised from a hardware point of view, aggressive online testing using Nvidia’s latest HPL code began. At this point it became apparent that the Cray recommended 19°C airstream temperature setpoint was too high. At this temperature the CSCS and Nvidia software engineers noticed a large number of GPUs thermal capping which resulted in poor HPL performance. Nvidia advised that the known hotter Location ‘B’ (as described above) in conjunction with normal variation in GPU temperature profile (thanks to the statistical variation in component leakage currents) was producing a situation where component temperature could not be kept within the expected operating range. Hence the decision was taken to lower the cabinet controller airstream setpoint to 16°C. Thanks to the changes made to the cooling loop infrastructure it was easy for the cooling system to accommodate this change, there being ample cooling capacity and flow-rate headroom now that all three rows were supplied by the TTN primary cooling loop. Experiments were also conducted with an airstream setpoint of 14°C but, ironically, this temperature seemed to be too low for the system, with a large number of errors appearing on the High Speed Network (HSN). Given this unexpected behaviour and the fact that 14°C was getting dangerously close to the dew point in the room it was decided to stay with the 16°C airstream setpoint. Note that throughout this process the dew point of the machine room environment was carefully monitored in order to ensure that liquid water did not condense on components within the system.

With the 16°C airstream setpoint the system stabilised considerably and this ultimately lead to the extremely high performance HPL runs with over 80% numerical efficiency (Rmax/Rpeak). For the best performing run the full system electrical power draw was 2.325MW for an average of 83kW per cabinet, well within the design envelope of 115kW per cabinet. The system was also 100% room neutral, thanks to changes made by Cray to the cabinet control system in the last cabinet of each row.

1) *Benefits of Secondary Loop Heat Exchangers:* During this period, given that the secondary cooling loops for the middle and back rows did not contain the usual primary-to-secondary-loop heat exchangers, whereas the front row did, this provided a good opportunity to compare the loop characteristics in both cases.

Figure 6 shows a schematic of the cooling loop layout for the middle row. The cooling loop layout of the other two rows is identical. As can be seen there are 6 temperature sensors in the loop. Obviously, as described above, in the case of the middle and back rows the heat exchanger shown in the schematic had been removed and by-passed so that the primary-side supply was connected directly to the secondary-side supply, and similarly for the returns.

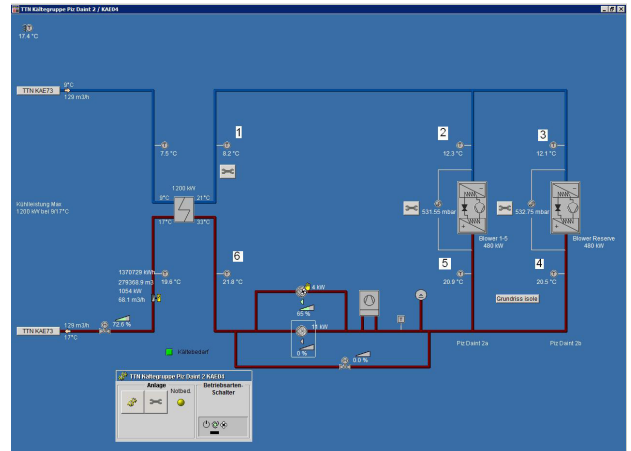


Figure 6. Schematic of the cooling loop layout for middle row showing the position of the various temperature sensors (labeled 1 to 6).

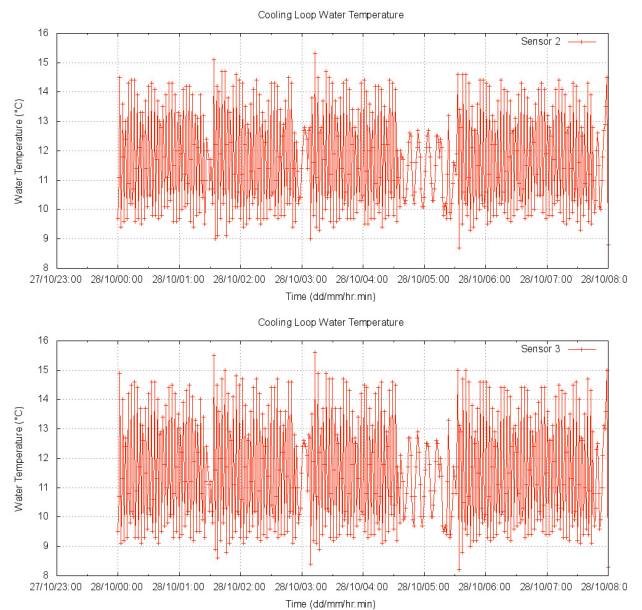


Figure 7. Plots showing the variation in coolant temperature as measured by Sensors 2 and 3 in the cooling loop of the middle row over an 8 hour period.

The secondary loop for the front row remained closed-loop but had sensors in exactly the same locations as shown in Figure 6.

Figure 7 shows the water temperature measured by Sensors 2 and 3 in the cooling loop of the middle row over a period of 8 hours in late October 2013. During this period the system ran a compute-intensive workload (HPL) 4 times and these periods can be recognised within the data. Notably however the temperature fluctuates considerably during the period in a range from 8.5°C to 15.5°C.

Figure 8 shows the water temperature measured by Sen-



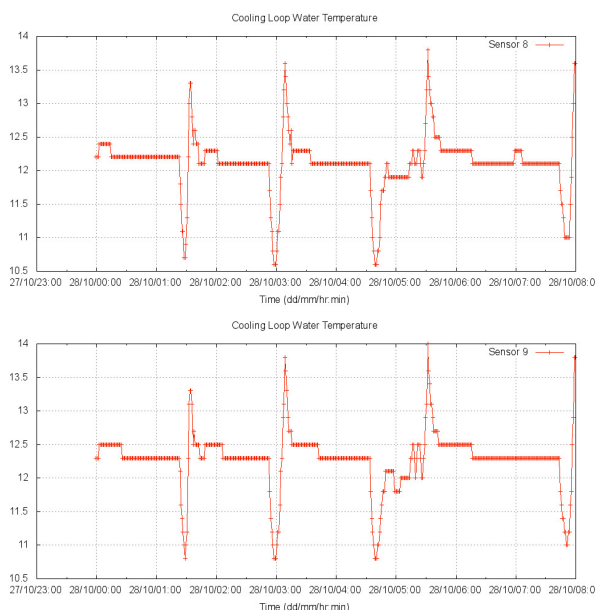


Figure 8. Plots showing the variation in coolant temperature as measured by Sensors 8 and 9 in the cooling loop of the front row over the same 8 hour period as that shown in Figure 7.

sors 8 and 9 in the cooling loop for the front row over the same 8 hour period; these sensors are in the same location on the supply side of this row as Sensors 2 and 3 of the middle row. Four periods where the water temperature varies over a range from 10.5°C to 14°C can be seen. The pattern of behaviour is typical of the step-change response of a damped control system where the control variable undershoots then overshoots the control setpoint, and is the response to the heat load from 4 separate HPL runs (not all runs were successful and hence stopped at different points during the run). However the most notable thing about these plots is that they clearly show the damping behaviour of the primary-to-secondary-loop heat exchanger.

Apparently decoupling the cooling loops by use of a heat exchanger provides the optimal configuration by allowing the control systems of each loop to independently control the water temperature in each loop. In so doing the computer system cabinets receive coolant at a much more stable temperature and this in turn reduces the burden on the cabinet control system which is designed to keep the airstream temperature at the airstream setpoint by controlling the water valve position on the supply side of the in-cabinet radiators. Furthermore there is less fluctuation in pump speed and valve positions on the facilities side and this will serve to lengthen the usable lifespan of these components and reduce energy consumption because the control system is not continuously “hunting” around the control variable setpoint.

There are of course additional benefits to having distinct primary and secondary cooling loops, including the ability to

isolate water leaks and the customisation of water treatment regimes and this all serves to justify the additional expense associated with this configuration.

2) *Monitoring of System Environmental Data:* The cabinet control software must of course monitor all system environmental data on a cabinet by cabinet basis and these data are made available by a script (`envdata`) provided by the Cray Site Engineering team. This script runs on the Software Management Workstation (SMW) and aggregates environmental data for each cabinet (via the `ccsysd` daemon running on the SMW) into a file. The following information is captured by the script each time it is run:

- Fan speed of the 6 blowers in each blower cabinet;
- Cabinet Controller airstream setpoint;
- Average airstream temperature and velocity for the 3 sensors located at the positions at the top, middle and bottom of the cabinet, downstream of the cabinet radiator;
- Airstream humidity
- Cabinet water control valve percentage opening
- Cabinet supply and return water temperature and pressure
- Per cabinet rectifier total power draw (kW)

In order to monitor the time history of these data, CSCS systems staff put an entry in the `crontab` of the SMW to call the script every 2 minutes and pipe the output to a unique file, timestamped at the time it was made. Plotting the time-series data for these various quantities then involved a simple (`awk` and `gnuplot`) script to parse the files and plot the data. This mechanism and parsing script, although simple in their implementation, proved invaluable in the early stages of system commissioning because they provided a simple way to view the large amount of environmental data available on the system.

Examples of problems detected by inspection of the environmental data time history plots include:

- 1) Cabinet controller crashes and reboots which then reset the airstream temperature setpoint;
- 2) Unrecoverable cabinet controller crashes that then stopped the water valve from changing from the last set opening. In some cases this meant that the valve position had to be manually set until the cause of the problem could be traced and rectified;
- 3) Erratic water valve position control which then caused the airstream temperatures to fluctuate widely;
- 4) Unreliable and erroneous cabinet inlet and outlet water temperatures;
- 5) Transposed airstream temperature data for the top and bottom sensors.

Problems of the type described in Item 1) were traced to the fact that the setpoint changes (from 19°C to 16°C) had not been made in a way that was persistent across cabinet reboots, and were easily rectified. The Type 2 problems



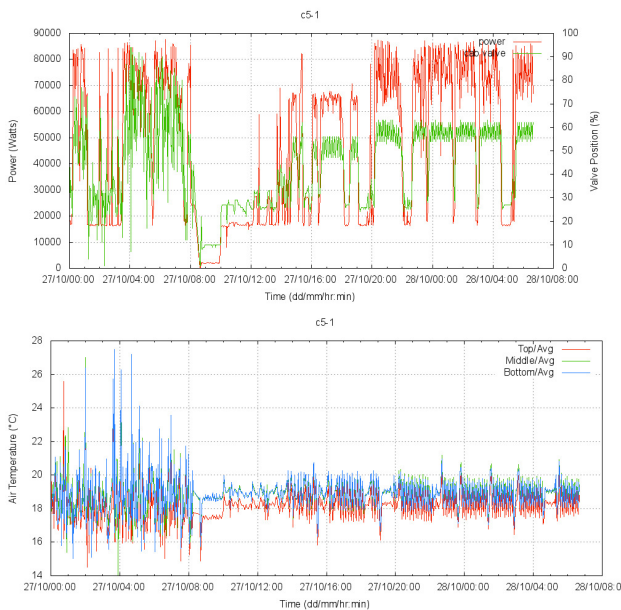


Figure 9. Plot showing cabinet c5-1 water valve position, heat load and corresponding airstream temperatures before and after a system reboot. As can be seen the water valve position is very erratic during the period before 8:30am.

were a little more pernicious and were finally traced to faulty cabinet control boards in some cabinets and, in one case, incorrectly labeled and wired in-cabinet secondary PDUs which prevented a faulty PDU from being correctly identified and replaced.

Figure 9 shows an example of the erratic water valve position for cabinet c5-1. Prior to the cabinet reboot at approximately 8:30am on 27 October the water valve position (green line) is very erratic. After the reboot the control system regains control of the valve position and it tracks the electrical and heat load (red line) in a much better fashion. Also shown are the corresponding airstream average temperatures at the three locations (top, middle and bottom) within the cabinet. As can be seen, when the water valve position was erratic the air temperatures fluctuate wildly and are completely outside of the acceptable operating range (at this time the temperature setpoint for the cabinet was 19°C). Cray finally tracked this behaviour to the use of a new spring-loaded valve actuator, designed to close the cabinet water valves in the event that the cabinet lost power or was powered down. Apparently these actuators caused out-of-range voltages on the Cabinet Control Board (CCB) which in turn caused the erratic water valve positioning. A Field Change Order (FCO) for the replacement of cabinet control system components (including the actuators, CCB, cables and cabinet controller firmware) was issued for this problem.

The Type 4 problems, where the cabinet inlet and outlet water temperatures were unreliable or erroneous was already known by Cray and is due to incorrectly mounted temperature/pressure transducers. Hence, for the time being CSCS rely on the facilities-side sensors for data of this type.

The Type 5 problems were traced to incorrectly labeled output data from the cabinet control software (i.e. the `ccsysd`).

#### D. Post-Install Experience

Installation and acceptance of the system went smoothly, from a facilities point of view, and no further issues were encountered with the facilities infrastructure. Installation of the primary to secondary heat exchangers was done in early 2014, as part of a regular maintenance and took one day. Installation of the preconditions was also completed during a regular maintenance and went smoothly. Some initial teething problems were again experienced with the Cabinet Control Boards however and this was traced to the use of sprung-loaded actuators overloading the control board circuitry. Thankfully the preconditioner valves could be manually set open until new actuators and CCBs could be supplied by Cray and now the system is operating smoothly. On 1 April 2014 the system was officially made part of the full CSCS User Programme and is now the flagship system at the Centre.

#### IV. CONCLUSION

The installation of the 28-cabinet hybrid XC30 system presented a number of challenges from a facilities point of view. Given the leading-edge nature of the machine and the fact that it was the outcome of a joint collaborative design effort between Cray, CSCS and Nvidia meant that the exact facilities requirements of the system were not fully known until very late in the piece. Access to the early test systems and the data they provided proved critically important in allowing all partners to anticipate possible issues, thus ensuring the facility team could prepare for them and react where necessary. Furthermore, once the full system was on-site the flexibility of the CSCS data center design was crucial in obtaining quick turn-around for the necessary last minute changes with virtually no impact on the installation schedule.

When installing leading edge systems such as this there is a need to exercise caution when designing the facilities infrastructure and consider flexibility in conjunction with large safety margins in order to accommodate variances in system requirements at installation time. Moreover data from early systems should be viewed with a critical eye because they may well give indications of divergence from vendor expectations based on data from prototype systems.

#### ACKNOWLEDGMENT

Thanks to the CSCS Facilities Management team and Gilles Fourestey from the CSCS Future Systems group.

Thanks also to Jim Tennesen, Mike Knitter and their colleagues in the Cray Site Engineering team. Finally thanks to the engineering team at Nvidia.

#### REFERENCES

- [1] "CSCS Fact Sheets." [Online]. Available: [http://www.cscs.ch/publications/fact\\_sheets/index.html](http://www.cscs.ch/publications/fact_sheets/index.html)
- [2] "Ashrae Revised Class 1 & 2," 2008. [Online]. Available: <http://www.energystar.gov/>
- [3] C. Site Engineering, "Cray Cascade Liquid Cooled SITE PREPARATION GUIDE," *Cray Internal Document*, 2012.