

Energy Efficiency Considerations for HPC Procurement Documents: 2014

(revision 1.0)

Energy Efficient High Performance Computing Working Group (EE HPC WG)

Contact: Natalie Bates natalie.jean.bates@gmail.com

Table of Contents

1	Introduction.....	5
2	Measurement Requirements for Power and Energy	7
2.1	System, Platform and Cabinet Level Measurements	9
2.2	Node-Level Measurements	10
2.3	Component-Level Measurement	11
3	Timestamps and Clocks	13
4	Temperature Measurements.....	14
4.1	Cabinet Level Temperature	14
4.2	Node Level Temperature.....	14
4.3	Component Level Temperature.....	15
5	Benchmarks.....	17
6	Cooling.....	18
6.1	Liquid Cooling	18
6.2	Air Cooling.....	18
7	High Level Objectives	20
7.1	Energy Related Total Cost of Ownership (TCO).....	20
7.2	Power Usage Effectiveness (PUE).....	20
7.3	Total Usage Effectiveness (TUE)	20
7.4	Energy Re-Use Effectiveness (ERE).....	21
7.5	Power Distribution	21
8	Usage Cases for Power, Energy and Temperature Management and Control.....	22
8.1	Data Center Infrastructure	22
8.2	System Hardware and Software	23
8.3	Applications, Algorithms, Libraries.....	23
8.4	Schedulers, Middleware, Management	24
A	LIST OF ITEMS TO CONSIDER FOR 2015 VERSION OF THIS DOCUMENT	26
A.1	Liquid Cooling Commissioning	26
A.2	Power API support	26
A.3	Measurement accuracy	26
A.4	Power Bounding requirements	26
A.5	Further guidance with respect to continuous vs. on-demand data collection.....	26
A.6	Sub-component power/energy measurement requirements	26
A.7	Workload metrics like energy to solution and time to solution	26
B	REFERENCES AND LINKS	27
B.1	EE HPC WG	27

B.2 [2013 Document](#)..... 27
B.3 [ANSI C12.20](#) 27
B.4 [EnergyStar](#) 27
B.5 [Firestarter](#) 27
B.6 [ASHRAE Air Cooling](#) 27
B.7 [ASHRAE Liquid Cooling](#) 27
B.8 [TUE](#) 27
B.9 [ERE](#)..... 27

List of Figures

Figure 1: Reported Values vs. Internal Samples.....	8
Figure 2: IT Equipment Environmental Classes.....	19

List of Tables

Table 1: System/Platform/Cabinet Internal Sampling Frequency	10
Table 2: System/Platform/Cabinet External Power/Energy Reported Value Frequency	10
Table 3: Node Internal Sampling Frequency	11
Table 4: Node External Power/Energy Reported Value Frequency	11
Table 5: Component Internal Sampling Frequency	12
Table 6: Component External Power/Energy Reported Frequency.....	12

1 Introduction

This document is written by the Energy Efficient High Performance Computing Working Group (EE HPC WG), whose mission is to encourage the adoption of energy conservation measures and energy efficient design in high performance computing (HPC). It is intended that this document encourage dialogue in the entire community about priorities and specific requirements for HPC system energy efficient features and capabilities. It captures energy efficiency requirements that the Energy Efficient High Performance Computing Working Group (EE HPC WG) recommends as important considerations when writing procurement documents for supercomputer acquisitions. It draws upon recent procurement documents created and used by major supercomputing sites, but also draws upon content experts in energy efficient HPC to modify and supplement the material from these documents. The team that wrote this document has been comprised of members in the user community, but it has been reviewed by members of the vendor community for their feedback.

This document is not a primer or ‘how-to’ on writing procurement documents. The requirements described in this document are intended to be vendor and technology neutral. They are intended to be high level and encourage dialogue, not to set guidelines or define standards. They should encourage innovation and not pick a particular vendor, architecture, technology, product or any other implementation.

The energy efficiency of HPC systems has been improving, but the road ahead still requires improvement. This document sets this year’s vision (2014) for systems to be delivered and accepted in two years (2016). It identifies priorities and sets an immediate bar. It is expected that the priorities will change and the bar will rise over time. This document will be refreshed on a yearly basis and 2013 was the first year. The 2013 version of the document is available at <http://eehpcwg.lbl.gov/sub-groups/equipment-1/procurement-considerations/procurement-considerations-presentations>. This is the 2014 version of the document.

Most of the document describes requirements that could be used to specify system features and capabilities. These requirements are categorized as mandatory, important, or enhancing. In addition to these requirements, the document includes informational content that could be used to set the context for the acquisition, but not be used as a requirement.

Each HPC center has its own unique mission, and priorities may differ greatly between users. The requirements are intended to draw lines in the sand that can be easily re-drawn, not to build isolating fences. Some of these requirements, especially those that are enhancing, may drive up product cost beyond the value of the feature or capability to the user. The authors recognize that there may be trade-offs, but also want to encourage the dialogue that helps to communicate requirements as well as costs. The HPC center has the exclusive responsibility for managing their procurements processes. It is hoped that this document will encourage a consideration of energy efficiency during the execution of those processes.

1. Section 2 describes power and energy measurement requirements. The measurement requirements span from a high level view of the entire system to a low level view of individual components. This section was included in the 2013 version of this document, but it has been substantially enhanced and expanded.

2. Section 3 is new to the 2014 version of this document, and it describes requirements for timestamps and clocks.
3. Section 4 is also new to the 2014 version of this document and it describes requirements for temperature measurements.
4. Section 5 describes requirements for benchmarking power and energy. This section has some updates, but remains mostly unchanged from the 2013 version of the document.
5. Section 6 describes requirements for cooling, both air and liquid. This section covers both the computer system and the data center. Although enhanced and expanded, this section was included in the 2013 version of this document.
6. Section 7 describes requirements for high level objectives, like Total Cost of Ownership and other metrics. Many of these are more specific to the data center than to the computer system. This section remains mostly unchanged from the 2013 version of the document.
7. Section 8 of the document describes usage cases for management and control, but doesn't define requirements. These are suggestive examples that serve to help clarify the requirements set forth in sections above. This section is mostly unchanged from the 2013 version of the document

Conventions

Information: info

Requirements: enhancing
important|
mandatory

2 Measurement Requirements for Power and Energy

- (info)** Power and energy measurement capabilities are necessary to understand the energy demand of the HPC systems in order to properly size support systems and plan for future growth.. These mechanisms may differ in implementation and purpose, and include capabilities for measuring the energy consumption of entire systems, platforms (subsystems), cabinets, nodes and components.
- (info)** This section is primarily focused on measuring the system power and energy, which includes system hardware and software.
- (info)** Section 8 describes usage cases for power and energy management and control.
- (mandatory)** The vendor shall provide the mechanism, interface, hardware, firmware, software, and any other elements that are necessary to capture the individual power and energy measurements.
- (mandatory)** This capability should have no (or minimal and defined) impact on the computation, security, and energy consumption of the equipment. The vendor must describe the impact, preferably in quantitative terms.
- (mandatory)** Scalable tools to extract accumulate and display power, energy and temperature information (accumulated energy and peak, instantaneous as well as average power between any two points in time) should be delivered.
- (mandatory)** The power and energy data must be exportable with at least a comma-separated value (CSV) or a user-accessible application programming interface (API).
- (mandatory)** For power, energy (and discrete current and voltage measurements if available) a detailed description of the measurement capabilities must be provided, including a specified value for measurement precision, accuracy and how data samples are time-stamped. Reference ANSI C12.20. The data can be based on real physical measurements or heuristic event based models.
- (info)** Why hierarchy?

The document is formatted in somewhat of a hierarchical fashion. The purpose of this is to address the various current and anticipated future use cases related to this topic. Component level measurement, for example, is required for fine-grained application power and energy analysis; likewise, component level control could be used to shift power from one component to another based on specific application requirements. Measurement at node level granularity is necessary for understanding the power and energy characteristics of a multi-node application, for example. While cabinet level measurement might have fewer current use cases, cabinet level power capping, as well as node level, are emerging as important requirements in recent procurements. Platform level measurement and

control has many facility inspired use cases and is a critical piece of overall platform management.

(info)

A number of terms are used in this document to describe measurement capabilities. It is important to understand the context in which the terms are used. [Figure 1](#) illustrates these terms. The x-axis of [Figure 1](#) is Time (in generic units). Note, [Figure 1](#) represents a range of possible capabilities that are useful for this discussion, it does not imply that these specific capabilities are a requirement.

- The top horizontal line represents points in time when discrete internal current and voltage measurements are sampled at the device level. These samples are not exposed externally. At each time interval a voltage and current sample is internally measured (v_6, i_6 pair for example).

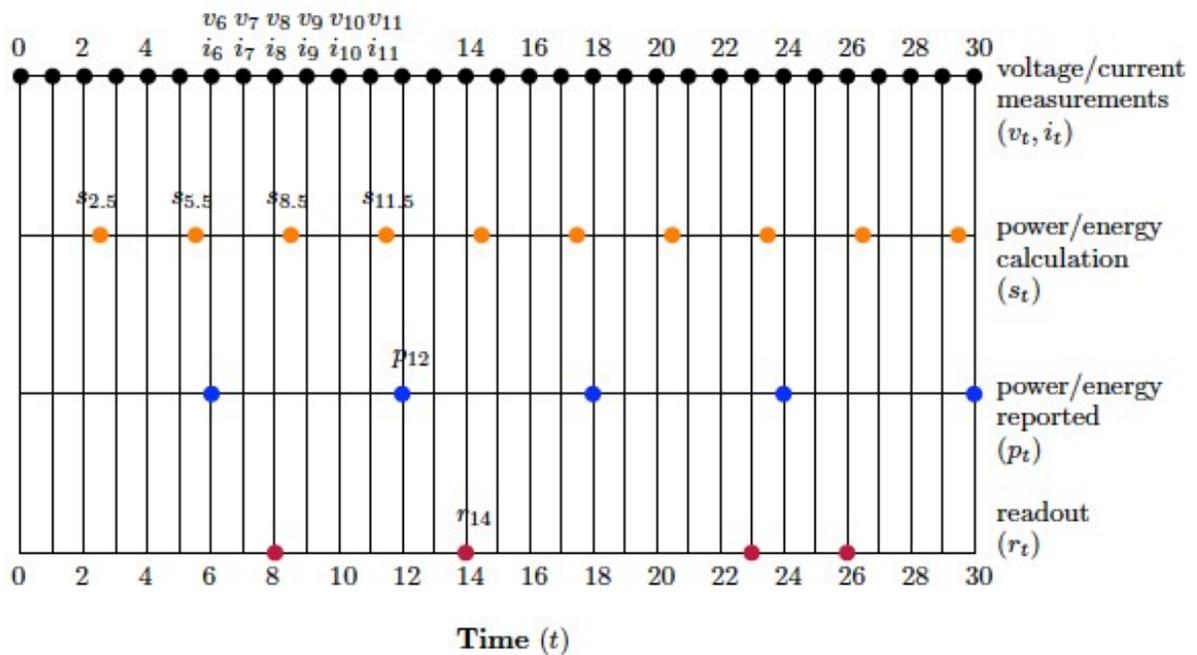


Figure 1: Reported Values vs. Internal Samples

- The second line down represents the points in time when an internal power and/or energy calculation is performed. Again, this is not exposed externally.
- The third line down represents the points in time a reported value is available to be read, externally. Each reported value could represent an average power, an instantaneous power, or an accumulated energy value, depending on the device capabilities. For example, point P_{12} could simply be the power value calculated at $S_{8.5}$ or $S_{11.5}$. P_{12} could also be the average power of points $S_{8.5}$ and $S_{11.5}$, or all of the calculated power samples prior to P_{12} . P_{12} could likewise be an accumulated energy value representing any range of power samples up

to that point in time. The important distinction is the difference between the device's internal sampling capability (frequency of and what the samples represent) and the external reported value capability of the device (again, frequency of and what the values represent).

- Finally, the fourth line down represents when the user actually obtains the reported value readout. It is critical that the timestamp of the reported value represents the time, as accurately as possible, of the measurement. Notice that the actual readout takes place at various time intervals following the availability of the reported value. This emphasizes the importance of time stamping at the time of measurement, not at the time of reading the value.

For example, a measurement device may be capable of producing 100 discrete power samples per second (internally). The power calculation (sample) and availability of the reported value of this same device may be equivalent to the lowest level sampling frequency, but no greater. Both, are typically less than the internal sampling frequency. For example, the same device may have the ability of producing a reported a value at 1 times per second. This reported value could be a power value averaged over 1 second, an accumulated energy value over the past 1 second, or simply a discrete power value for that moment in time.

Generally speaking, the requirements for the frequency of the reported value depend on what the reported value represents. If the reported value is a discrete power value then a higher frequency of reported value is typically desired. If the reported value represents an average power or accumulated energy value, reported frequency is less important than the internal sampling frequency that is used to derive the reported average power or energy value.

2.1 System, Platform and Cabinet Level Measurements

- (info)** The system level may vary by site and architecture, but could be so broad as to include all of the parts of the system that explicitly participate in performing any workload(s). This might include supporting internal and external power and cooling equipment as well as internal and external communication and storage sub-systems.
- (info)** The platform is distinguished from the system so as to differentiate compute from other sub-system equipment (such as external storage) that may be managed distinctly, but together comprise a system.
- (info)** The cabinet (or rack) is the first order discretization of the platform level measurement. The cabinet may be part of the compute, storage or networking platform.
- (mandatory)** Must be able to measure system, platform, and cabinet power and energy.

[Table 1](#) lists the mandatory, important and enhancing requirements for the internal device sampling frequency. The internal samples may be individual current and voltage samples or combined into a discrete power sample (see [Figure 1](#)).

[Table 2](#) lists the mandatory, important and enhancing requirements for the external reported value frequency. This is the data that is exposed externally for consumption (or readout, see [Figure 1](#)). The external reported values can represent a discrete or average power value, or an energy value. The details of the time period represented by the average power and energy values, how power and energy are calculated and time-stamped must be specified. Note that reported rate might differ from readout rate. Readout is when a user chooses to consume the reported value and is limited by the reported rate.

Table 1: System/Platform/Cabinet Internal Sampling Frequency

	Internal Sampling Frequency
Mandatory	≥ 10 per second
Important	≥ 100 per second
Enhancing	≥ 1000 per second

Table 2: System/Platform/Cabinet External Power/Energy Reported Value Frequency

		External Reported Value Frequency
Mandatory	Discrete Power (W)	≥ 1 per second
	Average Power (W)	≥ 1 per second
	Energy (J)	≥ 1 per second
Important	Discrete Power (W)	≥ 10 per second
	Average Power (W)	≥ 1 per second
	Energy (J)	≥ 1 per second
Enhancing	Discrete Power (W)	≥ 100 per second
	Average Power (W)	≥ 1 per second
	Energy (J)	≥ 10 per second

(important) The vendor shall assist in the effort to collect these data in whatever other subsystems are provided (e.g., another vendor’s back-end storage system).

(important) Those elements of the system, platform and cabinet that perform infrastructure-type functions (e.g., cooling and power distribution), must be measured separately with the ability to isolate their contribution to the power and energy measurements.

2.2 Node-Level Measurements

(info) A node level measurement shall consist of the combined measurement of all components that make up a node for the architecture. For example, components may include the CPU, memory and the network interface. If the node contains other components such as spinning or solid state disks they shall also be included in this combined measurement. The utility of the node level measurement is to facilitate measurement of the power and energy characteristics of a single application. The node may be part of the

network or storage equipment, such as network switches, disk shelves and disk controllers.

(important) The ability to measure the power and energy of any and all nodes must be provided.

[Table 3](#) lists the mandatory, important and enhancing requirements for the internal device sampling frequency. The internal samples may be individual current and voltage samples or combined into a discrete power sample (see [Figure 1](#)).

[Table 4](#) lists the mandatory, important and enhancing requirements for the external reported value frequency. This is the data that is exposed externally for consumption (or readout, see [Figure 1](#)). The external reported values can represent a discrete or average power value, or an energy value. The details of the time period represented by the average power and energy values, how power and energy are calculated and time-stamped must be specified. Note that reported rate might differ from readout rate. Readout is when a user chooses to consume the reported value and is limited by the reported rate.

Table 3: Node Internal Sampling Frequency

	Internal Sampling Frequency
Mandatory	≥ 100 per second
Important	≥ 1000 per second
Enhancing	≥ 10000 per second

Table 4: Node External Power/Energy Reported Value Frequency

		External Reported Value Frequency
Mandatory	Discrete Power (W)	≥ 10 per second
	Average Power (W)	≥ 10 per second
	Energy (J)	≥ 1 per second
Important	Discrete Power (W)	≥ 100 per second
	Average Power (W)	≥ 100 per second
	Energy (J)	≥ 10 per second
Enhancing	Discrete Power (W)	≥ 1000 per second
	Average Power (W)	≥ 1000 per second
	Energy (J)	≥ 10 per second

2.3 Component-Level Measurement

(info) Components are the physically discrete units that comprise the node. This level of measurement is important to analyze application energy/performance trade-offs. This level is analogous to performance counters and carries many of the same motivations. Counters are special purpose registers built into CPUs to store the counts of activities and are used for low-level tuning. Components can be any devices that are part of a node for a particular architecture.

(enhancing) The ability to measure the power and energy of each individual component should be provided.

[Table 5](#) lists the mandatory, important and enhancing requirements for the internal device sampling frequency. The internal samples may be individual current and voltage samples or combined into a discrete power sample (see [Figure 1](#)).

[Table 6](#) lists the mandatory, important and enhancing requirements for the external reported value frequency. This is the data that is exposed externally for consumption (or readout, see [Figure 1](#)). The external reported values can represent a discrete or average power value, or an energy value. The details of the time period represented by the average power and energy values, how power and energy are calculated and time-stamped must be specified. Note that reported rate might differ from readout rate. Readout is when a user chooses to consume the reported value and is limited by the reported rate.

Table 5: Component Internal Sampling Frequency

	Internal Sampling Frequency
Mandatory	≥ 1000 per second
Important	≥ 10000 per second
Enhancing	≥ 1000000 per second

Table 6: Component External Power/Energy Reported Frequency

		External Reported Value Frequency
Mandatory	Discrete Power (W)	≥ 100 per second
	Average Power (W)	≥ 10 per second
	Energy (J)	≥ 1 per second
Important	Discrete Power (W)	≥ 1000 per second
	Average Power (W)	≥ 100 per second
	Energy (J)	≥ 10 per second
Enhancing	Discrete Power (W)	≥ 10000 per second
	Average Power (W)	≥ 1000 per second
	Energy (J)	≥ 10 per second

3 Timestamps and Clocks

- (info)** For any post-mortem analysis, measured values need to be associated with a specific time or time frame. Having this time information allows system administrators to recall measured values in the past and correlate them to system events, configuration changes or batch jobs. Similarly, users can correlate energy consumption to application progress in order to improve the application's energy efficiency. All this requires meaningful timestamps to be associated with the measurement values.
- (important)** The vendor shall provide a mechanism to associate a timestamp with each measured value reported by the vendor infrastructure. The timestamp shall indicate the time at which the measurement value is derived and will indicate known accuracy. Any measured value and its associated timestamp shall be provided automatically by the vendor infrastructure.
- (mandatory)** The vendor shall provide a documentation that allows to quantify or bound the "age" of measured values. This shall include involved network latencies and jitter, filter delay, processing times and of course the update rate (see Section 2).
- (info)** Each timestamp is with respect to a reference clock. Possible reference clocks include the compute node clocks that are used in recording application progress and management clocks that are used in recording system events. Timestamps and clocks are never perfect. Therefore, knowing the correctness of a timestamp is important, especially for high update rate.
- (important)** The vendor shall provide information that allows for quantifying the accuracy of timestamps. To that end, the vendor shall describe the applicable factors that have significant impact. They can include:
- On which component and clock, the timestamps are generated.
 - Which clock is used to compute energy from power?
 - The drift of the used clocks.
 - A description of the synchronization mechanisms that are in place between the involved clocks.
 - How the delay between data acquisition and timestamp generation can be quantified.
 - The delay of analog filters or A/D conversion.

4 Temperature Measurements

- (mandatory)** The system must operate safely under all conditions. This includes when thermal emergencies are detected and when thermal sensors are faulty. The system must operate safely even with faulty sensors. Faulty sensors should be identified at the lowest possible level of sensor hierarchy.
- (mandatory)** The data on temperature must be physically accurate, reported in real-time, and provide sufficient detail. The accuracy must be $\pm 0.5^{\circ}\text{C}$ or better. Measurements must be sampled no slower than the quickest thermal response time expected. They must be accurately time-stamped, along with a description of how data samples are time-stamped. Also provided is a detailed explanation of where in the system the temperature is being measured.
- (info)** **There is no consensus yet on which sampling rates are considered sufficient. This bears some more study. For now, whatever sampling rate is used should be stated explicitly.**
- (info)** Generally, referencing existing standard is preferred than creating new ones. The U.S. ENERGY STAR program for computer servers can be used as a reference, although it does not completely capture the requirements in HPC.

4.1 Cabinet Level Temperature

- (important)** The temperature measurements must characterize the range of operating temperatures within the system by type of device (component, node, cabinet, etc.), as well as supply and return temperatures for each coolant. These temperature measurements must include uncertainty bounds, and be reported faster than the shortest thermal response time (e.g., every second).
- (important)** Dew point temperature of the air supplied to cabinets must be measured and reported to the cooling control system in charge of prevention of condensation.
- (info)** Temperature data are more valuable at the platform and cabinet levels than at the system level. Node and component level temperature measurements are also important but for different reasons. These temperatures are monitored to make sure the silicon remains within bounds.

4.2 Node Level Temperature

- (mandatory)** The node level temperature measurements must be representative of temperatures within the node, which must be physically described and justified. Uncertainty in measured temperatures must be stated, and measurements must be delivered faster than the shortest thermal response time of the chosen measurement location. Support must be present for the safety of the node when thermal emergencies are detected, and system

management must be notified in a timely fashion with a detailed account of the incident.

- (mandatory)** **The type of temperature being measured (for example, average or peak) shall be explicitly stated because node-level temperatures vary across the device.**
- (info)** The information about the correlation between temperature and power is more critical in an air-cooled environment.
- (info)** The following list some envisioned use cases of node- and component-level temperature measurement data exposed outside the node/component:
- A better understanding of how the power-consumption behavior of a device is influenced by its surrounding temperature. This also reveals the trade-off between leakage power and temperature. Higher temperature can reduce power on cooling but increase leakage power. But this may still be advantageous if compressor-based cooling can be eliminated.
 - A better modeling of device failure due to thermal effects as well as the development of mechanisms for short-term and long-term failure prediction. Measuring and constraining processor temperature can improve application performance in a faulty environment. However, different applications have different optimal temperatures (A cool way of improving reliability of HPC machines -- SC'13).
 - A better understanding of the thermal distribution within the machine and across machines to optimize the power cost for thermal management.
 - Temperature aware job scheduling: Different applications heat up CPU's to different temperatures. CPU temperature distribution is not homogenous throughout the data center i.e. same workload would heat up different CPUs to different temperatures. An intelligent job scheduler can take CPU-temperature and application-temperature profiles into account while assigning resources. Temperature aware scheduling could be useful in heat re-use as well.
 - A better understanding of influence of temperatures on turbo mode. Different temperatures can result in different maximal frequencies therefore creating an imbalance in computational capability that could have negative impact on parallel applications but also creates a potential for improving the scheduling of load-imbalances.

4.3 Component Level Temperature

- (mandatory)** The temperature data must include specified uncertainty bounds, and be sampled faster than the quickest thermal response time expected. Support must be present for the safety of the component when thermal emergencies are detected, and system management must be notified in a timely fashion with a detailed account of the incident.

(enhancing) The temperature of each individual component should be able to be measured.

5 Benchmarks

- (info)** Since power and energy costs, both operational and capital, are increasingly significant, it is very important to understand the power and energy efficiency requirements of the system. This is best understood when running workloads; either applications or benchmarks. Each site will have to select the workloads to run as part of the procurement and acceptance process. These workloads may differentially exercise or stress various sub-systems; compute (CPU, GPGPU, etc.), I/O, Networks (Internal, facility and WAN). They may focus on applications that are based on integer as well as floating point computations.
- (mandatory)** [Customer] shall specify the set of benchmarks they want. Vendors shall provide the power and energy efficiency requirements, and run times of a set of benchmarks.
- (mandatory)** The types of problem in the benchmarks shall cover compute problems, memory problems, networking problems, idle and sleep system state
- (important)** Benchmarks shall also cover dynamic power consumption. By regularly alternating between high and low power consumption, the measurement accuracy for dynamic power consumption can be verified. This exposes aliasing issues in the measurement. Different rates for alternating the workload shall be tested, with a focus on rates around the measurement sampling rate and other measurement processing rates (see also Figure 1).
- (info)** Suggested examples: HPL (compute problem), Integer-dominant codes (compute problem), Graph500 (memory/networking problem), GUPS, GUPPIE, MySQL and non-MySQL database applications, and systemBurn developed at ORNL and FIRESTARTER developed at TU Dresden (<http://tu-dresden.de/zih/firestarter/>).
- (mandatory)** Customers shall specify the run rules and the measurement quality. Each benchmark must be measurable using the Green500 run rules and attain a Level 2 measurement quality.
- (important)** Customers shall specify the run rules and the measurement quality. Each benchmark must be measurable using the Green500 run rules and attain a Level 3 measurement quality.
- (important)** Vendors shall work with Customers to provide the power and energy efficiency requirements of a set of site-supplied workloads. These workloads will reflect the typical case, not the extremes so that vendors can design around the typical case.
- (important)** Customers may also require application power profiles with power and energy requirements.

6 Cooling

6.1 Liquid Cooling

- (info)** For systems designed to be liquid-cooled, there is an opportunity for large energy savings compared to air-cooled designs. Since liquids have more heat capacity than air, smaller volumes can achieve the same level of cooling and can be transported with minimal energy use. In addition, if heat can be removed through a fluid phase change, heat removal capacity is further increased. By bringing the liquid closer to the heat source, effective cooling can be provided with higher temperature fluids. The higher temperature liquid can be produced without the need for compressor based cooling. Higher return liquid temperatures increase opportunities for recovery of waste heat, thereby increasing system efficiency.
- (info)** [*Customer*] will specify the type of liquid cooling systems contained within the data center. The range of liquid supply temperatures available in the center corresponding to ASHRAE recommended classes (W1-W4) will be provided to the vendor.
- (info)** A traditional data center is cooled using compressor-based cooling (i.e., chillers or CRAC units) and additional heat rejection equipment such as cooling towers or dry coolers. These liquid-cooled systems operate within ASHRAE recommended ranges W1 and W2. Systems designed to operate in these ranges will have limited energy efficiency.
- (important)** For improved energy efficiency and reduced capital expense, many data centers can be operated without compressor based cooling, by using cooling towers or dry coolers combined with water-side economizers. These data centers can operate within the ASHRAE W3 range and accordingly, systems should be requested to operate in this range.
- (enhancing)** In most locations liquid cooling of up to 45° C can be provided using dry coolers. The ASHRAE W4 classification was defined to accommodate this low energy form of cooling. For this type of infrastructure, ASHRAE W4 class should be requested.
- (info)** Parameters like pressure, flow rate and water quality may also be specified by each site in their procurement documents. ASHRAE provides guidance on these parameters, although they are not defined in this guideline.

6.2 Air Cooling

- (Info)** ASHRAE Thermal Guidelines (2011) define environmental classes that allow temperatures up to 40°C and 45°C. These allowable environmental temperature and humidity limits along with the recommended limits are shown in the psychrometric chart below (see Figure 2). Past IT equipment and some IT equipment manufactured today has aligned with operation within the A1 and A2 classes. Some present equipment as well as future

equipment will certainly enable operation within classes A3 and A4 to aid the industry in increasing energy savings. Generally, performance tradeoffs are made to enable operation in A3 or A4 environments. This must be balanced with the potential for energy savings.

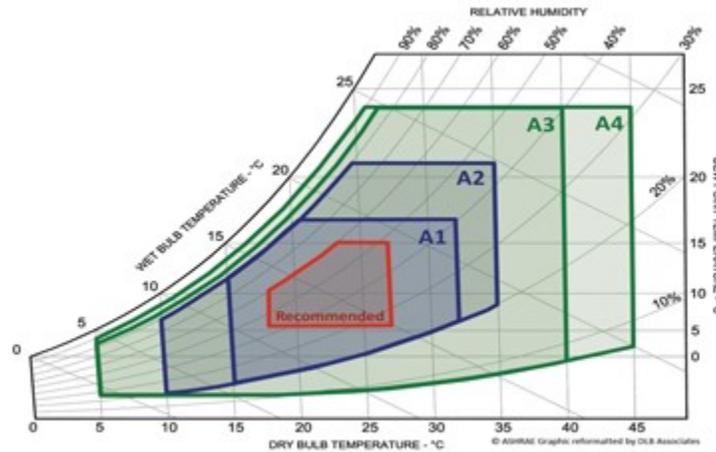


Figure 2: IT Equipment Environmental Classes

- (info)** One must be careful that fan power does not rise to extremes thereby countering potential data center level savings.
- (info)** Please Refer to ASHRAE book – Thermal guidelines for data processing environment 3rd edition. It describes all the trade offs.
- (mandatory)** The system must be able to operate in a Class A1 environment.
- (important)** It is better to operate in a Class A2 environment (important)
- (enhancing)** All other things equal, it is highly desirable to operate in a Class A3 environment.
- (mandatory)** The system must support full performance throughout the allowable range while operating in redundant mode (if applicable). Note: Most redundantly cooled systems can operate with a single cooling component failure.
- (enhancing)** The system must support full performance throughout the allowable range while operating in non-redundant mode (if applicable).

7 High Level Objectives

- (info)** The vendor shall provide [equipment, services and/or resources] that – among other objectives – establish a highly energy efficient solution at justifiable cost. The proposed solutions should demonstrate net benefits under normal production conditions.

7.1 Energy Related Total Cost of Ownership (TCO)

- (enhancing)** It is an objective of [*Customer*] to encourage innovative programs whereby the vendor and/or [*Customer*] are incentivized to reduce the costs for energy and/or power related capital expenditures as well as the operational costs for energy. This may be for the system, data center and/or broader site. By doing this, the vendor would be reducing the energy-related TCO for [*Customer*]. The vendor is encouraged to describe their support for these innovative programs in qualitative as well as quantitative terms.
- (info)** An example of an innovative program for bringing the energy/power element of TCO to the front was used by the Leibniz Supercomputing Center (LRZ). Their procurement was based on TCO whereby the budget covered not just investment and maintenance, but operational costs as well. The intent was to provide a clear incentive for the vendor to deliver a solution that would yield low operational costs and, thereby, lower TCO.

7.2 Power Usage Effectiveness (PUE)

- (info)** It is an objective of [*Customer*] to run a highly energy efficient data center. One measure for data center efficiency is PUE. It is recognized that the metric PUE has limitations. For example, solutions with cooling subsystems that are built into the computing systems will result in a more favorable PUE than those that rely on external cooling, but are not necessarily more energy efficient. In spite of these limitations, PUE is a widely adopted metric that has helped to drive energy efficiency.
- (enhancing)** The US Federal Data Center Consolidation Initiative has set a requirement to achieve an average PUE of 1.4 by 2015. As a result, the vendor is encouraged to qualitatively describe their support for helping [*Customer*] to meet this requirement.

7.3 Total Usage Effectiveness (TUE)

- (info)** TUE [Total Power Usage Effectiveness] (TUE) and IT Power Usage Effectiveness (ITUE)] account for infrastructure elements that are a part of the HPC system (like cooling and power distribution). TUE allows for inter-site comparison and, as such, is an improvement over PUE. iTUE is not only a metric that is necessary for calculating TUE, but stands on its own as a metric for a site to use for improving infrastructure energy efficiency. For more information, see: <https://www.brighttalk.com/webcast/679/96847>

(enhancing) The vendor is encouraged to qualitatively describe their support for measuring iTUE and TUE.

7.4 Energy Re-Use Effectiveness (ERE)

(info) Some sites have the ability to utilize the heat generated by the data center for productive uses, such as heating office space. Energy re-use is not strictly adding to the energy efficiency of either the computing system or the data center, but it can reduce the energy requirements for the surrounding environment. For those sites, it would be an objective of [*Customer*] to achieve an ERE < 1.0.

(enhancing) The vendor is encouraged to qualitatively describe their support for helping [*Customer*] to achieve an ERE < 1.0.

7.5 Power Distribution

(important) The vendor is encouraged to describe energy efficient and innovative solutions that help to a) optimize connection to electrical supply (e.g., electrical grid or on-site generation; and b) optimize electrical distribution within the data center and the HPC equipment. This will consider electrical equipment and conductor sizing; as well as backup and redundancy configurations (e.g., dual power supplies), so as to minimize electrical power conversion losses by considering the entire distribution chain to the processing components within the HPC system.

8 Usage Cases for Power, Energy and Temperature Management and Control

(info) As with the measurement capabilities described above, power and energy management and control capabilities (hardware and software tools and application programming interfaces (APIs)) are necessary to meet the needs of future supercomputing energy and power constraints. It is extremely important that [*Customer*] utilize early capabilities in this area and start defining and developing advanced capabilities and integrating them into a user friendly, production environment.

The vendor shall provide mechanisms to manage and control the power and energy consumption of the system. These mechanisms may differ in implementation and purpose. Below are envisioned usage models for these management capabilities. They are categorized loosely by where the management occurs. It is envisioned that this capability will evolve over time from initial monitoring and reporting capabilities, to management (including activities like 6-sigma continuous improvement), and even to autonomic controls.

These usage models are not requirements for the vendor, but rather suggestive examples that serve to help clarify the requirements for measurement capabilities described in section 4 above. Furthermore, it is recognized that many of these solutions would be provided by a third party, not by the system vendor.

8.1 Data Center Infrastructure

(info) Respond to utility requests or rate structures. For example, cut back usage during high load times, limit power during expensive utility rate times of the day.

“Power capping” may have multiple uses, include one that allows for provisioning the infrastructure for closer to average usage, leading to substantial infrastructure savings compared to those centers which are designed for theoretical peak usage.

Respond to demand requests; including increases in load to accommodate waste heat recovery, renewable energy, etc.

Manage rate of power changes; e.g., avoid spikes. Another example, the large variations of harmonic current produced by computer loads may need to be balanced in the data center as well as the site’s broader infrastructure and even the grid.

Provide an integrated view for the system and the building in a building management system.

8.2 System Hardware and Software

(info)

Reduce power utilization during "design days" so as to enable use of free cooling without backup chillers. Alarm and/or automatic shut-down that responds to environmental temperature excursions that are outside of the facility design envelope by reducing system loads.

Identify higher than normal power draw components needing maintenance and/or replacement. Or, also to identify higher than normal power draw usage from software- perhaps that is "stuck" in an infinite loop-back mode.

Proliferate power scaling and management beyond computation, to memory, communication, I/O and Storage. For example, under and overlocking, OS/hardware control of the total amount of energy consumed

Besides the traditional compiling for performance, the compiler vendor may want to provide the user with mechanisms to compile for energy efficiency. The possible mechanisms may include the following.

- Compiler flags for specifying performance-energy trade-offs or regarding energy efficiency as an optimization goal or a constraint.
- Programming directives for conveying user-level information to the compiler for a better optimization in the context of energy efficiency.
- Program constructs to promote energy as the first-class object so that it can be manipulated directly in source code.
- Compiler-based tools for reporting analyzed results regarding the energy efficiency of applications.

8.3 Applications, Algorithms, Libraries

(info)

Provides programming environment support that leads to enhanced energy efficiency. Some examples are reducing wait states and reducing the power draw in wait states.

Reduce wait-states examples:

- Schedule background I/O activity more efficiently with I/O interface extensions to mark computation and communication dominant phases.
- Use an energy-aware MPI library which is able to use information of wait-states in order to reduce energy consumption.

Reduce the power draw in wait-states examples:

- Attain energy reduction for task-parallel execution of dense and sparse linear algebra operations on multi-core and many-core processors, when idle periods are leveraged by promoting CPU cores to a power-saving C-state.

Scale resources appropriately. Examples are the following:

- Apply the phase detection procedure to parallel electronic structure calculations, performed by a widely used package GAMESS. Distinguishing computation and communication processes have led to several insights as to the role of process-core mapping in the application of dynamic frequency scaling during communications.
- Analyze the energy-saving potential by reducing the voltage and frequency of processes not lying on a critical path, i.e. those with wait-states before global synchronization points.
- Enabling network bandwidth tuning for performance and energy efficiency.

Select appropriate energy-performance trade-off. An example is the following:

- Optimize the power profile of a dense linear algebra algorithm (PLASMA) by focusing on the specific energy requirements of the various factorization algorithms and their stages.

Programming and performance analysis tools.. An example is the following:

- Counters, accumulators, in-band support

Open up control of these policies so that we can turn them on and off. Zero setting if it is detrimental to our applications at scale.

8.4 Schedulers, Middleware, Management

(info)

Putting hardware into the lowest reasonable power state or switching off idle resources (nodes, storage, etc.) when job scheduling cannot allow for full utilization.

Different power states. Careful about how we switch it off. Can't affect reliability. Sleep states is probably the best direction. Response time is much better.

Energy-aware scheduling: Develop mechanism to automatically select processor frequency for which energy to solution is minimized for a specific application.

Demand response – as in the ability to react to electrical grid based incentives – requires enhanced scheduling tools.

Evolving hardware features will likely require enhanced system software and scheduling tools with control at all levels of the hierarchy; from the system down to the components. An example might be a scenario where you have a high priority job, there are available nodes to run the job, but if run at the desired P-state, the system would exceed some notion of a power cap. In this situation, can one dynamically alter the p-state of lower priority jobs to allow them to continue, perhaps at a slower rate, while also accommodating the new, high priority job.

A LIST OF ITEMS TO CONSIDER FOR 2015 VERSION OF THIS DOCUMENT

A.1 Liquid Cooling Commissioning

A.2 Power API support

A.3 Measurement accuracy

A.4 Power Bounding requirements

A.5 Further guidance with respect to continuous vs. on-demand data collection

A.6 Sub-component power/energy measurement requirements

A.7 Workload metrics like energy to solution and time to solution

B REFERENCES AND LINKS

B.1 [EE HPC WG](#)

B.2 [2013 Document](#)

B.3 [ANSI C12.20](#)

B.4 [EnergyStar](#)

B.5 [Firestarter](#)

B.6 [ASHRAE Air Cooling](#)

B.7 [ASHRAE Liquid Cooling](#)

B.8 [TUE](#)

B.9 [ERE](#)