

**Energy Efficient High Performance
Computing Power Measurement
Methodology**

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 3 |
| 2 | Checklist for Reporting Power Values | 4 |
| | Quality Level | 4 |
| | Power Measurement Locations | 4 |
| | Measuring Devices | 5 |
| | Workload Requirement | 5 |
| | Level 3 Power Measurement | 5 |
| | Level 2 Power Measurement | 6 |
| | Level 1 Power Measurement | 7 |
| | Idle Power | 7 |
| | Included Subsystems | 8 |
| | Tunable Parameters | 8 |
| | Environmental Factors | 9 |
| 3 | Reporting Power Values (Detailed Information) | 10 |
| 3.1 | Measuring Device Specifications | 10 |
| 3.2 | Measuring Device Terminology | 10 |
| | 3.2.1 Sampling | 10 |
| | 3.2.2 Power-Averaged and Total Energy Measurements | 12 |
| 3.3 | Aspect and Quality Levels | 13 |
| 3.4 | Aspect 1: Granularity, Timespan and Reported Measurements | 14 |
| | 3.4.1 Core Phase | 15 |

1 Introduction

This document recommends a methodological approach for measuring, recording, and reporting the power used by a high performance computer (HPC) system. The document also discusses auxiliary data, such as environmental conditions, related to running a workload.

This document is part of a collaborative effort between the Green500, the Top500, the Green Grid, and the Energy Efficient High Performance Computing Working Group. While it is intended for this methodology to be generally applicable to benchmarking a variety of workloads, the initial focus is on High Performance LINPACK (HPL), the benchmark used by the Top500.

This document defines four aspects of a power measurement and three quality ratings. All four aspects have increasingly stringent requirements for higher quality levels.

The four aspects are as follows:

1. Granularity, time span, and type of raw measurement
2. Machine fraction instrumented
3. Subsystems included in instrumented power
4. Where in the power distribution network the measurements are taken

The quality ratings are as follows:

- Adequate, called Level 1 (L1)
- Moderate, called Level 2 (L2)
- Best, called Level 3 (L3)

The requirements for all four aspects for a given quality level must be satisfied to grant that quality rating for a submission.

2 Checklist for Reporting Power Values

When you are ready to make a submission, go to the Top500 (<http://www.top500.org/>) and Green500 (<http://www.green500.org/>) sites for more information.

This section contains a checklist for all the items of information you need to consider when making a power measurement.

Read through the list and ensure that you can record the needed information when you run your workload.

[] **Quality Level**

Choosing a quality level is the first important decision a submitter must make. Refer to Section 3.3 Aspect and Quality Levels for general information about the three quality levels. Sections 3.4 through 3.8 describe the details of the three quality levels

[] **Power Measurement Locations**

Measurements of power or energy are often made at multiple points in parallel across the computer system. A typical location might be the output of the building transformer. Refer to Section 3.8 Aspect 4: Point where the Electrical Measurements are Taken for more information about power measurement locations.

Note that in some cases, you may have to adjust for power loss. For information about power loss, refer to Section 3.8.1 Adjusting for Power Loss. If you adjust for power loss, how you determined the power losses must be part of the submission.

[] **Measuring Devices**

Specify the measuring device or devices used. A reference to the device specifications is useful.

Refer to Section 3.2 for some terminology about the measuring device specific to the power submissions described in this document. This section describes the difference between power-averaged measurements and total energy measurements.

Refer to Section 3.1 for information about the required measuring device.

If multiple meters are used, describe how the data aggregation and synchronization were performed. One possibility is to have the nodes NTP-synchronized; the power meter's controller is then also NTP-synchronized prior to the run.

[] **Workload Requirement**

The workload must run on all compute nodes of the system. Level 3 measures the power for the entire system. Levels 1 and 2 measure the power for a portion of the system and extrapolate a value for the entire system.

[] **Level 3 Power Measurement**

Level 3 submissions include the average power during the core phase of the run and the average power during the full run.

The core phase is usually considered to be the section of the workload that undergoes parallel execution. The core phase typically does not include the parallel job launch and teardown.

Level 3 measures energy. Power is calculated by dividing the measured energy by the elapsed time. The measured energy is the last measured total energy within the core phase minus the first measured total energy within the core phase.

Refer to Section 3.2.2 Power-Averaged and Total Energy Measurements for information about the distinction between energy and power.

The complete set of total energy readings used to calculate average power (at least 10 during the core computation phase) must be included, along with the execution time for the core phase and the execution time for the full run.

Refer to Section 3.4 Aspect 1: Granularity, Timespan and Reported Measurement for more information about the Level 3 Power Submission.

Refer to Section 3.5 for more information about the format of reported measurements.

For Level 3, all subsystems participating in the workload must be measured. Refer to Section 3.7 Aspect 3: Subsystems Included in Instrumented Power for more information about included subsystems.

With Level 3, the submitter need not be concerned about different types of compute nodes because Level 3 measures the entire system.

[] Level 2 Power Measurement

Level 2 submissions include the average power during the core phase of the run and the average power during the full run.

The complete set of power-averaged measurements used to calculate average power must also be provided. Refer to Section 3.2.2 Power-Averaged and Total Energy Measurements for the definition of a power-averaged measurement and how it differs from a total energy measurement.

For Level 2, the workload run must have a series of equally spaced power-averaged measurements of equal length. These power-averaged measurements must be spaced close enough so that at least 10 measurements are reported during the core phase of the workload. The reported average power for the core phase of the run is the numerical average of the ten (or more) power-averaged measurements collected during the core phase.

Each of the required equally spaced measurements required for L2 must power-average over the entire separating space.

Refer to Section 3.4 Aspect 1: Granularity, Timespan and Reported Measurement for more information about the Level 2 Power Submission.

Refer to Section 3.5 for more information about the format of reported measurements.

For Level 2, all subsystems participating in the workload must be measured or estimated. Level 2 requires that the greater of ? of the compute-node subsystem or 10 kW of power be measured. It is acceptable to exceed this requirement.

The compute-node subsystem is the set of compute nodes. As with Level 1, if the compute-node subsystem contains different types of compute nodes, you must measure at least one member from each of the heterogeneous sets. The contribution from compute nodes not measured must be estimated. Refer to Section 3.7 Aspect 3: Subsystems Included in Instrumented Power for information about heterogeneous sets of compute nodes.

[] **Level 1 Power Measurement**

Level 1 requires at least one power-averaged measurement during the run. Refer to Section 3.2.2 Power-Averaged and Total Energy Measurements for the definition of a power-averaged measurement and how it differs from a total energy measurement.

The total interval covered must be at least 20% of the core phase of the run or one minute, whichever is *longer*.

If the choice is one minute (because it's longer than 20% of the core phase) that minute must reside in the middle 80% of the core phase. If the middle 80% of the core phase is less than one minute, the measurement must include the entire middle 80% and overlap equally on both sides.

If the choice is 20% of the core phase (because this 20% is greater than one minute), this 20% must reside in the middle 80% of the core phase.

Refer to Section 3.4 Aspect 1: Granularity, Timespan and Reported Measurement for more information about the Level 1 power submission.

For Level 1, the only subsystem included in the power measurement is the compute-node subsystem. The compute-node subsystem is the set of compute nodes. Measure the greater of 1/64 of the compute-node subsystem or 1kW of power.

List any other subsystems that contribute to the workload, but do not provide estimated values for their contribution.

For some systems, it may be impossible not to include a power contribution from some subsystems. In this case, list what you are including, but do not subtract an estimated value for the included subsystem.

If the compute node-subsystem contains different types of compute nodes, measure at least one member from each of the heterogeneous sets. The contribution from compute nodes not measured must be estimated. Refer to Section 3.7 Aspect 3: Subsystems Included in Instrumented Power for information about heterogeneous sets of compute nodes.

[] **Idle Power**

Idle power is defined as the power used by the system when it is not running a workload, but it is in a state where it is ready to accept a workload. The idle state is not a sleep or a hibernation state.

An idle measurement need not be linked to a particular workload. The idle measurement need not be made just before or after the workload is run. Think of the idle power measurement as a constant of the system. Think of idle power as a baseline power consumption when no workload is running.

For Levels 2 and 3, there must be at least one idle measurement. An idle measurement is optional for Level 1.

[] **Included Subsystems**

Subsystems include (but are not limited to) computational nodes, any interconnect network the application uses, any head or control nodes, any storage system the application uses, and any internal cooling devices (self-contained liquid cooling systems and fans).

- For Level 1, all subsystems participating in the workload must be listed.

Only the compute-node subsystem must be measured. Not every compute node belonging to the compute node subsystem must be measured, but the contribution from those compute nodes not measured must be estimated. Measure the greater of at least 1/64 of the compute-node system or at least 1kW of power.

- For Level 2, all subsystems participating in the workload must be measured and, if not measured, their contribution must be estimated. The Measured % and the Derived % must sum to the Total %.

In the case of estimated measurements for subsystems other than the compute-node subsystem, the submission must include the relevant manufacturer specifications and formulas used for power estimation.

Measure the greater of at least 10KW of power or 1/3 of the compute-node subsystem.

- For Level 3, all subsystems participating in the workload must be measured.

Include additional subsystems if needed.

Refer to Section 3.7 Aspect 3: Subsystems Included in Instrumented Power for more information about included subsystems.

Refer to Section 3.6 Aspect 2: Machine Fraction Instrumented for information about

[] **Tunable Parameters**

Listing tunable parameters for all levels is optional. Typical tunable values are the CPU frequency, memory settings, and internal network settings. Be conservative, but list any other values you consider important.

A tunable parameter is one that has a default value that you can easily change before running the workload.

If you report tunable parameters, submit both the default value (the value that the data center normally supplies) and the value to which it has been changed.

[] **Environmental Factors**

All levels require information about the cooling system temperature. Reporting other environmental data (such as humidity) is optional.

Submissions require both the in and the out temperature of the cooling system. For air-cooled systems, these are the in and out air temperatures. For liquid-cooled systems, these are the in and out temperatures of the liquid.

Refer to Section 3.9 Environmental Factors for more information.

3 Reporting Power Values (Detailed Information)

Refer to this section for detailed information about the elements of a power submission.

This section describes the information that must be included with a power measurement submission. It also describes some optional information that submitters may decide to include.

The section contains definitions of the terms used to describe the elements of a power submission, some background information, motivation about why the list contains the elements it does, and any other details that may be helpful.

3.1 Measuring Device Specifications

Measuring devices must meet the Level requirements as defined in Sections 3.3 and 3.4. This section lists resources for finding and evaluating meters.

The ANSI specification for revenue-grade meters is ANSI C12.20.

Also, refer to the Power and Temperature Measurement Setup Guide and the list of accepted power measurement devices from the Standard Performance Evaluation Corporation.

http://www.spec.org/power_ssj2008/

http://www.spec.org/power/docs/SPECpower-Device_List.html

3.2 Measuring Device Terminology

Levels 1 and 2 specify power measurements. Level 3 specifies an energy measurement, but reports a power value.

3.2.1 Sampling

For Levels 1 and 2, power measurements must be sampled at least once per second. The actual measurements that constitute a sample may be taken much more frequently than

3 Reporting Power Values (Detailed Information)

once per second.

Sampling in an AC context requires a measurement stage that determines the true power delivered at that point and enters that value into a buffer where it is then used to calculate average power over a longer time. So "sampled once per second" in this context means that the times in the buffer are averaged and recorded once per second. Sampling delivered electrical power in a DC context refers to a single simultaneous measurement of the voltage and the current to determine the delivered power at that point. The sampling rate in this case is how often such a sample is taken and recorded internally within the device.

If the submitter is sampling in a DC context, most likely it will be necessary to adjust for power loss in the AC/DC conversion stage. Refer to Section 3.7 Aspect 3: Subsystems Included in Instrumented Power for details.

3.2.2 Power-Averaged and Total Energy Measurements

The reported power values for Levels 1 and 2 are power-averaged measurements. A power-averaged measurement is one taken by a device that samples the instantaneous power used by a system at some fine time resolution for a given interval. The power-averaged measurement for the interval is the numerical average of all the instantaneous power measurements during that interval and constitutes one reported measurement covering that interval.

Consider Level 1, which requires only one reported power value. This reported power value may consist of several power measurements taken at a frequency of at least once per second and averaged over an interval. That interval must be at least 20

Level 2 also requires that power measurements be taken at least once per second. Level 2 requires that power values be reported for both the core phase of the workload and the total workload. The reported power value for the core phase must be the result of at least 10 power measurements. These 10 power measurements may themselves be power-averaged measurements.

Each of the required equally spaced measurements required for L2 must power-average over the entire separating space. All the values reported by the meter must be used in the calculation.

For example, the meter may sample at one-second intervals and report a value every minute. Assume that the core phase is 600 minutes long. Assume further that the requirement for 10 equally spaced measurements can be satisfied with 10 measurements spaced 50 minutes apart. Using just those 10 measurements does not conform to this specification because all the values reported by the meter during the core phase are not used.

Although those two measurements were equally spaced over 500 minutes, each averaged only over a minute. So some (the majority) of the separating space between the measurements was not included in the average.

For Levels 1 and 2, the units of the reported power values are watts.

Level 3 specifies a total energy measurement that, when divided by the measured time, also reports power. An integrated measurement is a continuing sum of energy measurements. Typically, there are hundreds of measurements per second. Depending on the local frequency standard, there must be at least 120 or 100 measurements per second. The measuring device samples voltage and current many times per second and integrate those samples to determine the next total energy consumed reading.

Level 3 reports an average power value for the core phase, an average power value for the

3 Reporting Power Values (Detailed Information)

Table 3.1: Summary of aspects and quality levels

| Aspect | Level 1 | Level 2 | Level 3 |
|--------------------------------|---|---|--|
| 1a: Granularity | One power sample per second | One power sample per second | Continuously integrated energy |
| 1b: Timing | The longer of one minute or 20% of the run | Equally spaced across the full run | Equally spaced across the full run |
| 1c: Measurements | Core phase average power | <ul style="list-style-type: none"> • 10 average power measurements in the core phase • Full run average power • idle power | <ul style="list-style-type: none"> • 10 energy measurements in the core phase • Full run average power • idle power |
| 2: Machine fraction | The greater of 1/64 of the compute subsystem or 1 kW | The greater of 1/8 of the compute-node subsystem or 10 kW | The whole of all included subsystems |
| 3: Subsystems | Compute-nodes only | All participating subsystems, either measured or estimated | All participating subsystem must be measured |
| 4: Point of measurement | Upstream of power conversion OR Conversion loss modeled with manufacturer data | Upstream of power conversion OR Conversion loss modeled with off-line measurements of a single power supply | Upstream of power conversion OR Conversion loss measured simultaneously |

whole run, at least 10 equally spaced energy values within the core phase, and the elapsed time between the initial and final energy readings in the core phase. The average power value for the core phase is the difference between the initial and final energy readings divided by the elapsed time.

3.3 Aspect and Quality Levels

Table 3 1 summarizes the aspect and quality levels introduced in Section 1 Introduction

3.4 Aspect 1: Granularity, Timespan and Reported Measurements

Aspect 1 has the following three parts. Levels 1, 2, and 3 satisfy this aspect in different ways.

- The granularity of power measurements. This aspect determines the number of measurements per time element.
- The timespan of power measurements. This aspect determines where in the time of the workload’s execution the power measurements are taken.
- The reported measurements. This aspect describes how the power measurements are reported.

For all required measurements, the submission must also include the data used to calculate them. For Level 2 and Level 3 submissions, the supporting data must include at least 10 equally spaced points in the core of the run.

Levels 2 and 3 require a number of equally spaced measurements to be reported for two reasons.

- One is that facility or infrastructure level power measurements are typically taken by a system separate from the system OS and thus cannot be easily synchronized with running the benchmark.
- Secondly, with multiple periodic measurements, more reporting points are included before and after the benchmark run to ensure that a uniform standard of "beginning" and "end" of the power measurement can be applied to all the power measurements on a list.

There is no maximum number of reported points, although one reported measurement per second is probably a reasonable upper limit. The submitter may choose to include more than 10 such points.

The number of reported average power measurements or total energy measurements is deliberately given large latitude. Different computational machines will run long or short benchmark runs, depending on the size of the machine, the memory footprint per node, as well as other factors. Typically the power measurement infrastructure is not directly tied to the computational system’s OS and has its own baseline configuration (say, one averaged measurement every five minutes). These requirements are specified not only to give a rich data set but also to be compatible with typical data center power measurement infrastructure.

3 Reporting Power Values (Detailed Information)

All levels specify that power measurements be performed within the core phase of a workload. Levels 2 and 3 specify that a power measurement for the entire application be reported. Consequently, these levels require measurements during the run but outside of the core phase.

3.4.1 Core Phase

All submissions must include the average power within the core (parallel computation) phase of the run.

Every workload has a core phase where it is maximally exercising the relevant component(s) of the system. More power is often consumed in a workload's core phase than in its startup and shutdown phases. The core phase typically coincides with maximum system power draw.

For example, the core phase of the HPL workload is the portion of the core that actually solves the matrix. It is the numerically intensive solver phase of the calculation. Note that HPL now contains an `HPL_timer()` routine that facilitates power measurements.