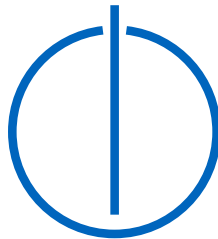# FAKULTÄT FÜR INFORMATIK

## DER TECHNISCHEN UNIVERSITÄT MÜNCHEN

Dissertation in Informatik

# Assessing the Energy Efficiency of High Performance Computing (HPC) Data Centers

Torsten Wilde

**TECHNISCHE UNIVERSITÄT MÜNCHEN**

Fakultät für Informatik
Lehrstuhl für Rechnertechnik und Rechnerorganisation

# Assessing the Energy Efficiency of High Performance Computing (HPC) Data Centers

Torsten Wilde

Vollständiger Abdruck der von der Fakultät für Informatik der Technischen Universität München zur Erlangung des akademischen Grades eines

Doktors der Naturwissenschaften (Dr. rer. nat.)

genehmigten Dissertation.

Vorsitzender: Prof. Dr. Thomas Huckle
Prüfer der Dissertation:

1. Prof. Dr. Arndt Bode

2. Prof. Dr. Dieter Kranzlmüller
   *Ludwig-Maximilians-Universität München*

Die Dissertation wurde am 13.11.2017 bei der Technischen Universität München eingereicht und durch die Fakultät für Informatik am 19.02.2018 angenommen.

Ich versichere, dass ich diese Dissertation selbständig verfasst und nur die angegebenen Quellen und Hilfsmittel verwendet habe.

München, den 8. April 2018                                  Torsten Wilde

# Acknowledgments

This dissertation came to fruition through the support of numerous people, who I wish to thank.

First and foremost, my sincere gratitude goes to my supervisor Professor Arndt Bode for his guidance, support, comments, and remarks.

I would also like to express my eternal gratitude to my wife, Jeanette, for her unwavering support and all the hours spent proofreading all my research papers, including this work, helping me to get my thoughts, ideas, and results written down in a coherent and understandable way.

My gratitude goes also to Professor Kranzlmüller for making space in his very busy schedule to be my second supervisor and for providing valuable comments and suggestions.

Furthermore, I would like to express my thankfulness to all my friends and colleagues at LRZ for their support and interesting discussions. In particular to Axel Auweter, who, as my room mate, had to survive many in depth discussions and brain storming sessions, and to Hayk Shoukourian for providing many interesting research ideas and valuable remarks.

I would also like to thank my parents and sister for supporting me in my decision to live for 11 years in the USA, without which this dissertation might not have happened, and for their encouragements over the last 4 years. Thank You!

Last but not least, I would like to thank my funding agencies that enabled me to work on the research presented in this dissertation. The work presented here was partially funded by:

And finally, I would like to acknowledge the Dagstuhl Seminar *Power-Bounded HPC Performance Optimization (15342)* from 2015 that inspired part of the presented work in this dissertation.

# Abstract

Data centers are faced with an ever increasing need for computational power. Over the last decade this need resulted in a substantial increase in data center power consumption world wide. At this point in time the energy consumption is becoming a main constraint for data center operation and expansion. Therefore, a major research area is the improvement of the energy efficiency of all data center systems including the cooling infrastructure, power distribution, and IT systems. In order to approach data center energy efficiency in a structured and unified way, a common view of the domain is needed. This common view should help to: identify gaps in current efforts; identify additional areas of improvement; and guide the development of needed tools for data collection, monitoring, analysis, and control.

This dissertation will contribute to the advancement of data center energy efficiency research in three important areas.

It will introduce a common frame of reference for the data center energy efficiency research domain by defining the 4 Pillar Framework which consists of the 4 pillars: Building Infrastructure; IT System Hardware; IT System Software; and Applications.

It will show that by using the newly defined 4 Pillar Framework important gaps can be identified and from this a new metric is introduced, Data Center Energy Efficiency (DCEE), that enables the calculation of the energy efficiency of a data center for a specific workload mix.

It will further show the existence of node power variability in homogeneous high performance computing systems and will discuss the usage of this IT hardware property to save energy by defining three techniques: Node Power Aware Scheduling; Node Power Aware System Partitioning; and Node Ranking Based on Power Variation.

The presented work will help to guide required work across traditional boundaries (like the connection between data centre building and infrastructure management systems and IT management) and will enable power and energy modelling comprising the whole data center. This can allow the characterization and evaluation of the impact of singular changes anywhere in the data center on the complete facility energy efficiency.

# Zusammenfassung

Rechenzentren stehen vor einem immer größeren Bedarf an Rechenleistung. Im Laufe des letzten Jahrzehnts hat diese Notwendigkeit zu einem erheblichen Anstieg des Stromverbrauchs der Rechenzentren weltweit geführt. Dieser immer noch ansteigende Energieverbrauch ist ein kritischer Faktor, der für den Betrieb und die Erweiterung des Rechenzentrums von großer Wichtigkeit ist. Ein großes Forschungsgebiet ist daher die Verbesserung der Energieeffizienz von Rechenzentrumssystemen, einschließlich der Kühlinfrastruktur, der Stromverteilung und der IT-Systeme. Um die Energieeffizienz des Rechenzentrums auf strukturierte und einheitliche Weise zu erreichen, ist eine gemeinsame Sicht auf die Domäne erforderlich. Diese gemeinsame Sicht sollte dazu beitragen, die Lücken bei den derzeitigen Anstrengungen zu erkennen; weitere Ver-besserungsbereiche zu identifizieren; und die Entwicklung von benötigten Werkzeuge für die Datenerfassung, Überwachung, Analyse und Kontrolle voranzutreiben.

Die Arbeiten, die in dieser Dissertation vorgestellt werden, tragen zur Verbesserung der Energieeffizienzforschung für Rechenzentren in drei wichtigen Bereichen bei.

Es wird ein ganzheitlicher Bezugsrahmen für den Forschungsbereich Energieeffizienz im Rechenzentrum einfgeührt. Dieser wird mittels einer 4-Säulen Grundstruktur definiert. Die vier Säulen sind: Gebäudeinfrastruktur; IT-System-Hardware; IT-Systemsoftware; und Anwendungen.

Durch die Anwendung der neu definierten 4-Säulen-Grundstruktur wird eine wichtige Lücke in der Bewertung der Energieeffizienz von Rechenzentren identifiziert und daraus eine neue Metrik mit der Bezeichnung Data Center Energy Efficiency (DCEE) eingeführt, die die Berechnung der Energieeffizienz eines kompletten Rechenzentrums für eine bestimmte Arbeitslast (definierte Anzahl von Anwendungen) ermöglicht.

Ferner wird das Vorhandensein von Knotenleistungsschwankungen in homogenen Höchstleistungscomputersystemen aufgezeigt. Diese Eigenschaft wird analysiert, inwieweit sie zur Energieeinsparung verwendet werden kann. Es werden drei Techniken definiert: Node Power Aware Scheduling; Node Power Aware System Partitionierung; und Node Ranking basierend auf der Leistungsvariation.

Die hier vorgestellte Arbeit wird dazu beitragen, traditionelle Grenzen zu verbinden (wie die Verbindung zwischen Rechenzentrumsgebäudeautomatisierung und Infrastrukturmanagementsystemen mit IT-Management) und die Energiemodellierung des gesamten Rechenzentrums zu ermöglichen. Dadurch kann der Einfluss singulärer Änderungen im Rechenzentrum auf die gesamte Energieeffizienz der Anlage charakterisiert und bewertet werden.

# Contents

# 1 Introduction

> The first step to knowledge is the
> confession of ignorance.
>
> — Gerald M. Weinberg

Energy consumption of data centers has increased substantially over the past decade and is becoming a major cost factor for data center operators.

Koomey [1], [2] estimated the power consumption for data centers world wide to increase from 8.1GW in 2000 to 17.4GW in 2005 and to somewhere between 23.2GW and 45.4GW in 2010.



Figure 1.1: Trend of Energy Consumption of Data Centers in Germany.

According to a survey done by DataCenter Dynamics [3] there was a world wide slow-down in data center power consumption from 19% in the year 2011-2012 to just over 7% in 2012-2013. At the same time the power consumption in Europe increased by 6% from 12.7GW to 13.5 GW. The same survey estimated the 2014 world wide data center power consumption to be 38.84GW. The Power Usage Effectiveness (PUE) [4] was between 1.81 and 2.0 for 2013. Meaning that for each amount of IT power an additional 81% to 100% was needed for data center overheads (for details on PUE see subsection 3.3.1 Existing metrics for assessing energy efficiency). The major part of this overhead can be attributed

to power transmission and conversion losses and data center cooling. A different survey done by the Uptime Institute [5] estimates the average PUE for 2013 at 1.65 for data centers world wide. Taking the average of both surveys (1.78), 17GW of the total data center power consumption was spent on the infrastructure needed to run the IT systems.

A latest study by the Fraunhofer Institute and Bordestep Institute [6] showed that the energy consumption from data centers in Germany increased from 10.5TWh(a) in 2010 to 12TWh(a) in 2015. It estimates a further increase in 2020 and 2025 which will raise the total energy consumption to 16.4TWh(a). An increase of 56% in 15 years.

This world wide development is driving an increased interest to save power and energy. PUE is still the main focus and multiple industry groups have taken on the challenge, for example, American Society of Heating, Refrigeration, and Air Conditioning Engineers (ASHRAE) guidelines [7] and the GreenGrid [8].

Since data centers are becoming more and more important for the society, governments provide guidelines (e.g. European Commission Code of Conduct for Energy Efficient Data Centers [9]) and funding agencies offer grants related to power and energy efficiency improvements in data centers.



Figure 1.2: Trend of Energy Consumption for LRZ from 2000 till 2016, and prediction for 2017.

Figure 1.2 and Figure 1.3 show the energy consumption and associated costs at the data center of the Leibniz Supercompting Centre (LRZ). LRZ is an institute of the Bavarian Academy of Sciences and Humanities (Bayerische Akademie der Wissenschaften, BADW) and functions as the IT service provider for Munich's universities and a growing number of scientific institutions in the greater Munich area and in the state of Bavaria. LRZ operates SuperMUC, a leadership class supercomputer with 241,000 x86 cores and a peak performance of over 6 PFlop/s, as well as a number of general purpose and specialized clusters including a 100% warm water-cooled Intel Xeon Phi Cluster (CoolMUC-3).

As can be seen from the Figures, the operating costs of the LRZ data center increased from less than 500 k€ in 2000 to 6.5 Mio.€ in 2015, an increase by a factor of 13. Clearly,

Figure 1.3: Trend of Energy Costs for LRZ from 2000 till 2016, and prediction for 2017.

this is not a growth rate that is sustainable in the future. Also the biggest consumer is the HPC system which consumes 2/3 of the overall data center energy.

As outlined above, improving the energy efficiency of data centers is a very important topic.



Figure 1.4: Generic continuous improvement process.

This dissertation will approach this topic in a generic way enabling a continuous improvement process [10] (Figure 1.4) and applies the results using the LRZ data center as an

example.

The first part of this dissertation will define a formal approach to data center energy efficiency which will be the foundation for improving the energy efficiency in a data center (*System* - Figure 1.4).

The second part of this dissertation discusses work related to measuring power and energy information and will propose a new metric closing a gap in data center energy efficiency assessment (*Measure* and *Analyze* - Figure 1.4).

The third part of this dissertation discusses the existence of node power variability in homogeneous HPC systems and investigates ways in which this property could be used to improve the energy efficiency of data centers (*Analyze* and *Improve* - Figure 1.4).

The final part of this dissertation summarizes the accomplishments. It shows the impact of the thesis and concludes with a look at future work enabled by the presented results.

**Research questions this dissertation will answer are:**

**Question:** What are the parts of an HPC data center that are important for energy efficiency?

**Thesis:** There exists a unified way to approach data center energy efficiency that is applicable to all data centers.

**Question:** How can power and energy related data be measured, collected, and evaluated?

**Thesis:** It should be possible to have one Key Performance Indicator that measures the energy efficiency of a data center.

**Question:** Is there any energy saving potential not realized related to the HPC system hardware?

**Thesis:** IT hardware manufacturing tolerances will influence the energy consumption of large scale homogeneous IT systems.

# 2 A Wholistic View of The Data Center

> If you can't measure it, you can't
> improve it.
>
> ——— Attributed to Peter Drucker

The author would like to add to this: "If you don't know an area of improvement you can't measure it."

In order to improve data center energy efficiency the first step is to identify all influencing areas and to understand the interaction between the areas. This will allow one to collect the right data for all further work. This chapter will introduce "The 4 Pillar Framework for energy efficient HPC data centers" (partly based on the author's paper [11]) which shows the major areas influencing a data center's energy efficiency.

## 2.1 Challenges

Improving data center Energy Efficiency means to find a scientific solution to operate a data center with specific IT systems, for varying system loads, and with the minimum of energy.

Energy consumption is one of the challenges for exa-scale computing in light of the 20MW challenge[1] [12]. The increasing energy consumption of current and expected systems made it of interest to facility managers as well. Research in energy efficiency has taken on a new emphasis. There are many publications related to efforts to improve the energy efficiency ranging from specific application tuning [13], scheduling improvements (for example, energy aware scheduling [14] and temperature aware scheduling [15]), application co-design [16], new energy efficient HPC hardware architectures [17], new system and data center cooling technologies [18], data center level energy and power management [19], and power and energy management of federated data centers [20], [21].

To approach the improvement of data center energy efficiency in a scientific way, one will need to:

- identify all sources of energy consumption

- determine the impact of each source on the overall energy consumption

- determine the dependencies between the sources

- optimize locally (each source) as well as globally (over multiple sources)

---

[1] in effect at the time of this publication

HPC data centers are unique entities. Each data center is different in setup and operation. Defining a process for energy efficiency improvement that is data center specific will not help the community overall. Therefore, the best solution would be provided by a generic and systematic approach together with measurable metrics that define the energy efficiency for specific aspects of a data center.

Improving the data center energy efficiency in any meaningful way requires data. The first step for any data center serious about improving their energy efficiency would be to identify all important areas influencing its energy consumption, then to define a method of measurement for those areas, and to classify and correlate current research work and analyze available information and solutions related to the area of improvement, and lastly apply this work.

Currently, no generic framework exists that allows data centers to approach energy efficiency as a well defined process considering a complex system rather than focusing on single metrics, such as Power Usage Effectiveness (PUE) or Flops/W (see subsection 3.3.1 Existing metrics for assessing energy efficiency for details on PUE and Flops/W). The identification of the right stakeholders is another critical aspect of the energy efficiency improvement process to ensure that any improvements will have the right impact.

A basic view of the data center based on primary input and outputs is shown in Figure 2.1.

The only input on the operational site is *electrical power* used to power the data center. This power is converted into *heat* [22] and needs to be removed from the data center (operational output). From an operational perspective, a data center is a warehouse size space heater.

Since data centers exist to provide IT services it has users. Users submit *work* (for example, in HPC, directly by submitting scientific applications via batch scheduling systems, or for web services, like Google, indirectly via web queries). The *results* produced are returned to the user.



Figure 2.1: High level input/output of a data center.

## 2.2 The *4 Pillar Framework* for energy efficient HPC data centers

The *4 Pillar Framework* for energy efficient HPC data centers was developed to provide a fundamental structure that allows for a common view of the energy efficiency domain for data centers. It provides a frame of reference that helps with understanding the different aspects of a data center related to energy efficiency improvements. Figure 2.2 shows the basic *4 Pillar Framework*.



Figure 2.2: The *4 Pillar Framework* for energy efficient HPC data centers.

The 4 Pillars, which are explained in more details in their own sections, are:

- *Pillar 1* - Building Infrastructure representing everything in the data center required to run the HPC system, for example, the cooling infrastructure and electrical infrastructure (see chapter 2.2.3)

- *Pillar 2* - HPC System Hardware representing, in the context of the paper specifically, the hardware components of the complete HPC system, for example, the CPU, internal cooling system, memory, etc. (see chapter 2.2.4)

- *Pillar 3* - HPC System Software representing the system software stack of the HPC system (see chapter 2.2.5)

- *Pillar 4* - HPC Applications (chapter 2.2.6) representing both a Workload and the Workload-mix run on the HPC system

The 4 pillars are encompassed within in a data center (chapter 2.2.2) which is affected by external influences and constraints (chapter 2.2.1).

Figure 2.3 shows the *4 Pillar Framework* in relation to the high level input and output of a data center. Again, there are the two inputs (Work and Electrical Power) and the two outputs (Results and Heat).



Figure 2.3: The *4 Pillar Framework* showing the high level input/output of a data center and effected pillars.

The main objective (work) of an HPC data center is to enable advanced scientific discovery using scientific applications, therefore, applications (work) are submitted and the results from running those applications are returned. At LRZ, the User submits his application (job) via the SLURM [23] or LoadLeveler [24] batch scheduling system. The scheduling system allocates required resources and starts the job. At job end any results are returned to the user. The user is mainly interested in the underlying architecture of the HPC system in order to optimize the performance of their applications.

Traditionally, this junction of work and power has separated the IT side of a data center from the infrastructure side. But in the context of energy efficiency, one can no longer maintain this strong separation (see chapter 3 Measuring Energy Efficiency). For example, future scheduling systems might use the characteristics of the data center cooling infrastructure power consumption to make energy optimized scheduling decisions.

### 2.2.1 External influences and constraints which affect a data center

The main "External Influences and Constraints" are: weather and climate, geographical features (rivers, lakes), power contract, and heat re-use opportunities.

Weather and climate effects the efficiency of exchanging the generated heat with the environment (figure 2.3). Geographical features might favour a specific cooling technology. The power contract is an agreement with an external entity and determines possible optimization opportunities. If one would like to re-use part of the generated heat than possible users in the area are of importance.

Figure 2.4 shows the LRZ data center power profile from January 2014. The black line shows the whole data center power consumption. As can be seen, the peaks correlate very well with the green area on the bottom of the graph which shows the outside wet bulb temperature. Interesting enough this is not true for all peaks. The first and last peak of the wet bulb temperature has no correlated peak in the data center power consumption, similarly the first peak of the data center power consumption is apparently not related to the outside conditions.



Figure 2.4: LRZ power profile and wet bulb temperature from January 2014.

Determining how exactly the outside conditions influence the data center and what the best data center operational policy for different operation points are can help a data center to save power and energy.

Geographical features play an important roll during data center site selection in relation to the use of on-site renewable energy and to available energy efficient cooling technologies. For example, Google selects sites where the data center heat can be exchanged with the environment via water (Douglas County, Georgia USA, uses city waste water; Hamina, Finland, uses sea water ([25], [26])).

The contract with the power provider has a strong impact on possible energy or power optimization opportunities and solutions. For example, if one pays for a fixed and narrow power band, one might be very concerned with staying in the power band and getting as close to the upper band as possible since one has already paid for it. On the other hand, if one pays for energy at the end of the year one might find every possible way to save energy. At LRZ the power contract provides boundary conditions (see chapter 4.1). The energy budget is estimated 2 years in advance due to funding requirements. LRZ has a 100% renewable energy contract. The use of renewable energy does not make a data center more efficient but does increase its sustainability. Using on-site renewable energy depends on the location of a data center and can be added during the lifetime of a data center but needs to be part of initial planing and the data center infrastructure needs to be prepared before hand.

Heat re-use can help to recuperate operating costs but can also be a source of additional constraints mainly related to required water temperatures.

### 2.2.2 Data Center

The "Data Center" itself is part of the *4 Pillar Framework* since data centers have specific operational policies, mandates, and other administrative guidelines that will influence possible energy or power optimization opportunities and solutions. For example, if the data center goal is to generate Gordon Bell prise winners, than saving energy might not be a high level optimization goal. On the other hand, if the energy budget for a year is fixed, optimizing the energy-to-solution (EtS) of applications might be a high level target.

Data center specific constraints are, for example, Service Level Agreements (which might require a redundant infrastructure), maximum power dissipation for data center cooling technologies (which might put a limit on the data center power consumption independent of the physical power feed), and power feed capacity (the hard limit of the power draw, at LRZ this is 10MW for example).

For completeness, it needs to be said that the data center can include some power consumers like on-site office buildings. At LRZ those have a relatively fixed power consumption and the optimization possibilities are very low in comparison to the 4 Pillars. As an example, Figure 2.5 shows the power consumption of the LRZ office buildings for the 3rd quarter of 2016.
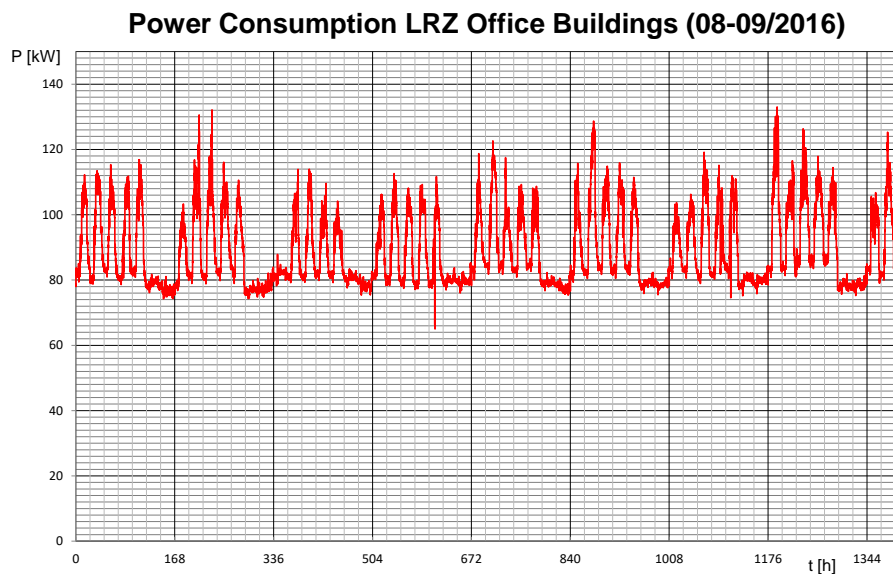


Figure 2.5: Power consumption of the LRZ office buildings from 3rd quarter 2016.

The power consumption is relatively constant with 80kW on weekends and varies slightly during working hours between 100kW and 120kW with peak loads of 130kW. This represents 2-3% of the total LRZ data center power consumption.

### 2.2.3 *Pillar 1 -* **Building Infrastructure**

*Pillar 1*, data center "Building Infrastructure", represents everything that is required to provide power to the IT systems and to remove the generated heat. The complexity depends on each data center design. As an example, Figure 2.6 shows an overview of the LRZ cooling infrastructure and Figure 2.7 shows an overview of the LRZ power distribution infrastructure. As can be seen, it is a complex cooling infrastructure and nearly all current cooling technologies are in use (air, cold water, hot-water, well water, and evaporative cooling). The electrical infrastructure is divided into different circuits providing power to different parts of the data center and power of a different quality is used. EV1 is the power quality that comes directly from the power supply line. EV2 is power filtered (mainly using dynamic flywheel UPS systems). EV3 is battery backed (static UPS systems) and EV4 is diesel generator backed.



Figure 2.6: LRZ cooling infrastructure in a nut-shell.

*Pillar 1* has a high impact on data center energy efficiency. Since the building infrastructure is designed in advance of the actual construction of the data center and is an integral part of the data center it is not easy to upgrade later on. Inefficiencies can have an impact over a very long time (>20 years) and any physical upgrade is cost intensive. Automated controls can be improved over the lifetime depending on available data, analyses tools, and the automation and control systems being used. Additional sensors can be added later but can be costly and time intensive to install and integrate. So far, analysis tools that can help to identify and analyse inefficiencies are missing.

Figure 2.8 shows the power used for removing a specific heat load for one cooling tower of the hot-water cooling circuit at LRZ. What can be seen is that there are 3 distinctly different operating points. To choose the right one can improve the energy efficiency of the cooling loop substantially.

Figure 2.7: LRZ electrical infrastructure in a nut-shell.



Figure 2.8: Infrastructure optimization opportunities.

The best known metric used in *Pillar 1* is, currently, PUE (for details see subsection 3.3.1). PUE was developed specifically to assess the infrastructure overhead. According to an Uptime Institute data center survey [5], the introduction of PUE in 2006 led to data centers reducing their overhead from an average of 2.5 (for each 1W into the IT systems, 1.5W went into the infrastructure showing a 150% overhead) in 2007 to an average of 1.65 in 2013. Recently, the shortcomings of PUE are becoming more of an issue and in response new metrics have been developed. For example, ASHRAE replaced PUE with electrical loss component (ELC) and mechanical load component (MLC) in its latest draft for Standard 90.4P "Energy Standard for Data Centers" [27] and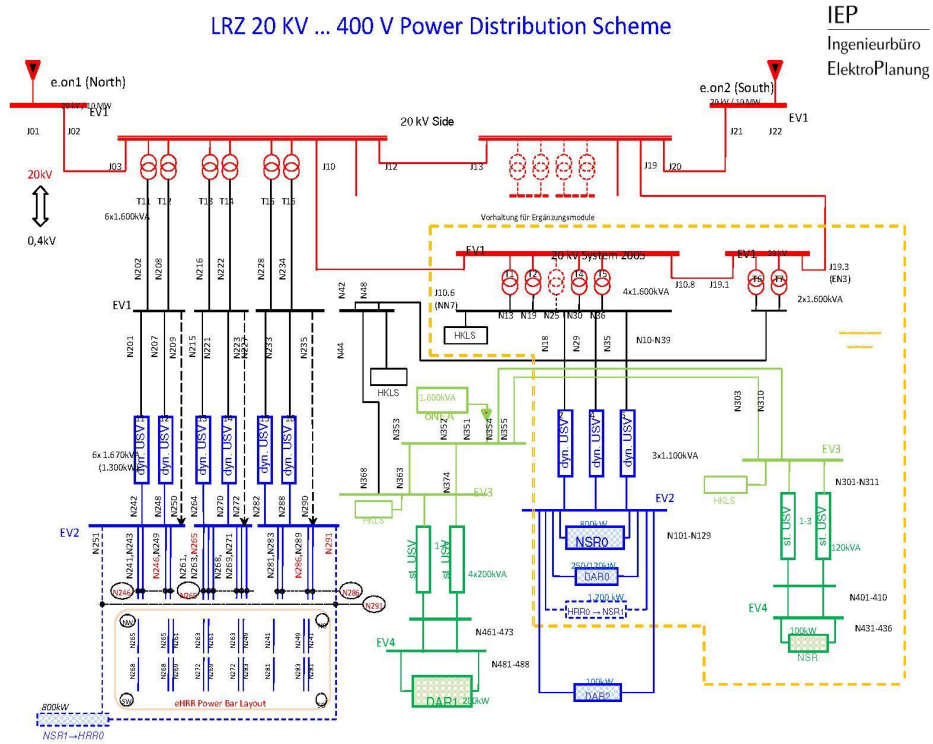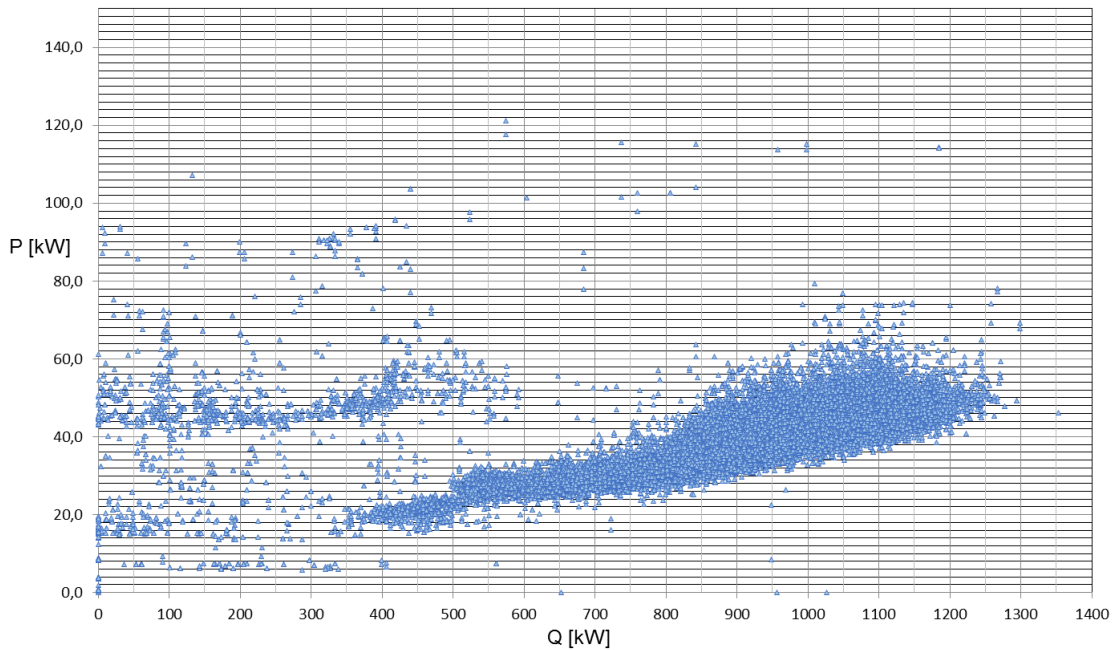 the Energy Efficient High Performance Computing Working Group (EEHPCWG) introduced IT-power Usage Effectiveness (ITUE) (a PUE like metric for the IT system) and Total Usage Effectiveness (TUE) (a combination of ITUE and PUE) [28]. This dissertation will introduce the Data Center Energy Efficiency (DCEE) metric which combines workload energy efficiency with IT system and data center cooling infrastructure overheads (for details see section 3.3).

### 2.2.4 *Pillar 2 -* HPC System Hardware

*Pillar 2* includes all of the IT hardware in the data center. For an HPC data center (like LRZ) the main focus is the HPC system since its behavior dominates the data center power consumption (see Figure 2.4). SuperMUC Phase1 (second line from top - pink) is the biggest power consumer and it's power profile is reflected in the complete data center power profile (first line from top - black).

Energy efficiency in *Pillar 2* can be improved by using more efficient cooling technologies such as chiller-less direct liquid cooling (also referred to as high temperature direct liquid cooling (HT-DLC), or warm or hot water cooling) which is used by SuperMUC. Figure 2.9 shows the inside of a node from SuperMUC Phase1 and Figure 2.10 shows the inside of a node of Phase2. In both node designs the CPUs, memory, and network chips are water cooled. In both nodes the only active air cooled part is the power supply. In Phase1 it sits beside the node and, therefore, does not draw air over the water cooled node. In Phase2 this changed. Here the power supplies are behind the node drawing air over the water cooled node. In praxis this node design is less optimal but Phase1 did not use standard racks (not as deep) whereas Phase2 does. The effectiveness of the system cooling technology is included in the ITUE metric. The DWPE and DCEE metrics (introduced by this work, see section 3.3 for more details) capture the efficiency of the complete system cooling technology.

Since the upgrade cycle for most HPC systems is 5 years it is important to procure the right system. This means getting the most energy efficient hardware (which often is the newest technology) for the application mix at the data center. The right choice of accelerators, CPU architecture, memory per node, and network which all support the application requirements is very important.

Since SuperMUC was procured in two phases, the improved hardware energy efficiency of newer technology can be seen. SuperMUC Phase1 ranks number 210 (with 846.42 MFlops/W) and Phase2 ranks number 92 (with 1,900.03 MFlops/W) in the Nov 2015 Green500 list [29].

There are some techniques that can be used during the runtime of the HPC system to improve the data center energy efficiency. One technique is Dynamic Voltage and Fre-
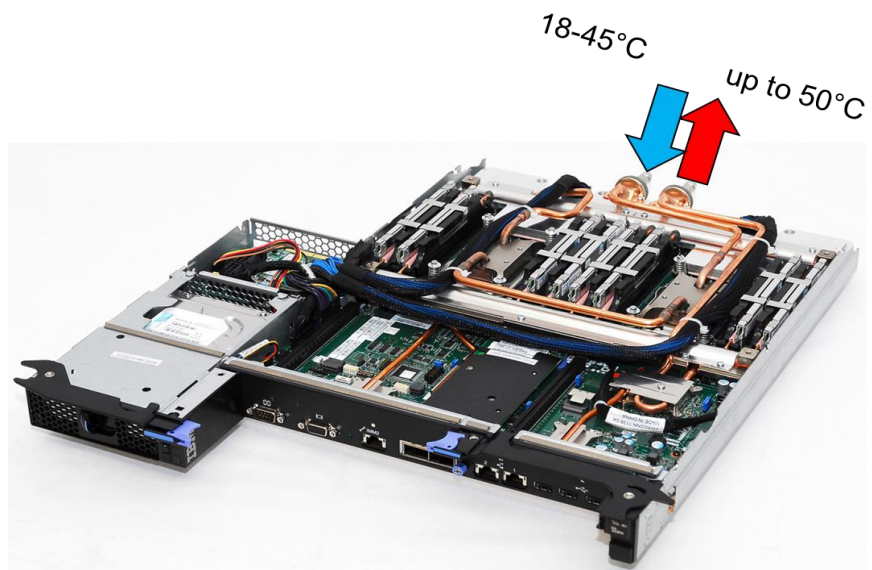
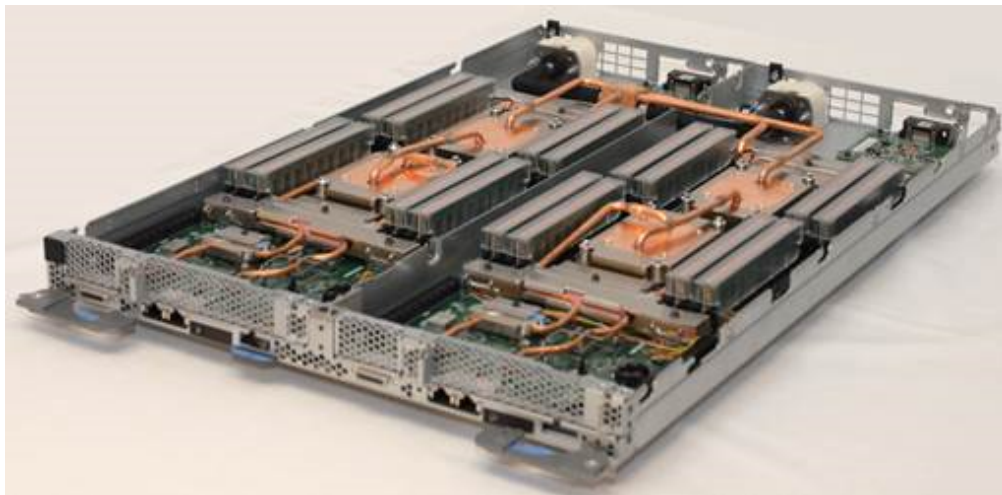Figure 2.9: SuperMUC Phase1 IBM iDATAPLEX dx360 M4 node picture.



Figure 2.10: SuperMUC Phase2 Lenovo NeXtScale nx360M5 WCT node picture (picture by Torsten Bloth, Lenovo).

quency Scaling (DVFS) which is widely discussed in the literature (for example, [30], [31], and [32]). Another possibility is to take advantage of hardware properties, for example, the node power variation of homogenous HPC systems (for more details see chapter Saving energy by taking advantage of node power variability in homogenous HPC systems). And yet another technique is to switch off unused hardware which is mainly used in virtual environments together with workload consolidation [33].

There is also the idea of hardware over-provisioning in power-constrained High Performance Computing ([34], [35]) which is not a technique to save energy but is mainly proposed because of specific contracts with the power provider mainly in the USA where some national labs pay for a specific *n* amount power and not the consumed energy.

### 2.2.5 *Pillar 3* **- HPC System Software**

*Pillar 3* represents the complete HPC system software stack. The software stack includes: the scheduling system (at LRZ, for example, SLURM [23] and LoadLeveler [24]); the operating system (at LRZ, SUSE SLES Linux); and all tools and libraries that can be used by the user and/or user application (such as profiling tools and optimized libraries).

Linux comes with different power governors and, depending on an applications behavior, one governor might be more energy efficient than another ([36], [37]).

The scheduling system is an important component that influences the systems energy and power efficiency. Current work is looking into energy aware scheduling ([38], [39]) which is used at LRZ ([14]) and thermal aware scheduling ([40], [41]).

In the future, the scheduling system might become data center infrastructure aware since the energy efficient use of resources does depend on the data center infrastructure. For example, if the infrastructure is very inefficient (high overhead) because of outside conditions or possible weather predictions, scheduling less power hungry jobs might be beneficial. Or if the infrastructure is running efficiently the scheduler should run power hungry jobs. In this way a new "data center aware" scheduler could take advantage of system and data center infrastructure properties.

As the examples above show, *Pillar 3* is one area where a data center has many knobs that can be adjusted during the lifetime of the HPC system to improve the overall data center energy efficiency. Unfortunately, the implemented energy efficiency improvements might not transfer easily to the next generation HPC system.

### 2.2.6 *Pillar 4* **- HPC Applications**

*Pillar 4* represents the applications running on the HPC system. Scientific discovery is the main reason for the existence of HPC data centers such as LRZ. Figure 2.11 shows the scientific domains of all applications run on SuperMUC in 2014. As can be seen, LRZ provides compute resources for nearly all major scientific domains [42].

The number of different applications running on SuperMUC is estimated to be over 240. This variety makes energy saving in *Pillar 4* a challenge since each application needs to be analysed and optimized. Luckily, optimizing for performance is also optimizing for energy efficiency [43]. A side effect is that performance optimized code usually consumes more power since the IT hardware is better utilized. Therefore, power capping might be counter productive to performance optimization.

Figure 2.11: SuperMUC applications by science area (2014)

Even though application improvements can lead to substantial savings when starting from less optimized code, optimizing an already well optimized code might not provide enough benefits. Also, one has to keep in mind that the lifetime of a scientific code might span multiple hardware generations of HPC systems and it will not be possible to optimize for fundamentally different hardware architectures multiple times. That said, adapting code parts (solvers, libraries) closer to the underlaying architecture might be beneficial and feasible over multiple HPC system generations. But here a strong tie between application developers, users, and the HPC data center is essential. One example of such activity is the LRZ extreme scaling workshop [44]. That said, *Pillar 4* is the Pillar over which most multi-science HPC data centers, such as LRZ, have the least influence. For example, LRZ provides some application support but the development and optimization of applications is the responsibility of the application owner.

In general, one can envision three possible goals in *Pillar 4*:

- Solve the same problem with less energy (reduce Energy-to-Solution (EtS)), this can include code optimization or support for different hardware architectures

- Increase the level of details solved or simulated in the application but stay in the same energy envelope as before, for example, as described for SeisSol in [45]

- Provide application support for power consumption management depending on data center needs

## 2.3 Usage

The generic *4 Pillar Framework* can be used to generate specific instances which identify areas of improvements and/or provide a drill down for different pillars. These can be used to focus on specific aspects such as different IT systems in the data center, a complete interaction chain, showing project coverage, or classifying information according to the different pillars.

Figure 2.12 uses the *4 Pillar Framework* to show the multi-layer complexity of improving the energy efficiency of data centers.



Figure 2.12: The *4 Pillar Framework* for energy efficient HPC data centers showing the complexity of the domain.

As can be seen, each pillar provides additional details. For example, *Pillar 1* shows possible different cooling technologies that might be used (air cooling, water cooling with different classifications according to ASHRAE [7]), different cooling tower options, and optional seasonal support technologies. Using the *4 Pillar Framework* in this way shows in a graphical representation the high complexity of the data center energy efficiency domain. Figure 2.13 shows the use of the *4 Pillar Framework* to present the work areas of a research project. Here, the focus areas of the SIMOPEK project [46] are presented.

At LRZ the final goal is to develop a global optimization strategy that encompasses all aspects of the data center. The SIMOPEK project is one step towards that goal. SIMOPEK, which has received funding from the German Federal Ministry for Education and Research under grant number 01IH13007A, developed a model and simulation of the chiller-less cooling infrastructure at LRZ in order to assess how much energy can be saved via multi-criteria optimization. The simulation took the behavior of the IT systems into account (CooLMUC, SuperMUC). In order to create and simulate a model, relevant information (sensor data) needs to be collected. This is done via PowerDAM (for details see section 3.2 Wholistic Power and Energy Data Collection - PowerDAM). PowerDAM V1.0 was developed to assess the energy efficiency of applications (calculating EtS) by collecting power

Figure 2.13: The *4 Pillar Framework* showing the coverage of the SIMOPEK project.

data from the HPC system power distribution units (*Pillar 2*) and data from the scheduling system (*Pillar 3*). In SIMOPEK, PowerDAM was extended (PowerDAM V2.0) to collect information from *Pillar 1*. Another aspect of the project was to specify and assess (via the simulation) data center specific adsorption chiller designs [47].

The *4 Pillar Framework* is also used in subsection 3.3.1 Existing metrics for assessing energy efficiency to classify Key Performance Indicators (KPI's) related to measuring the energy efficiency of data centers (Figure 3.12).

## 2.4  Related Work

As discussed in my previous work [11], the following existing works share some ideas with the *4 Pillar Framework* but leave out significant parts (directly quoted from the author's paper *The 4 Pillar Framework for energy efficient HPC data centers* [11]).

> For example, "Energy Efficiency in Data Centers: A new Policy Frontier" [48] is referring to parts of *Pillar 1 and 2* only and, thus, it is solely focusing at the data center energy efficiency from an operational point of view. Another example, the "DPPE: Holistic Framework for Data Center Energy Efficiency" [49], references *Pillar 1, 2, and 3* to define the data center energy flow. This flow is then used to find areas of improvement and the corresponding division (operating unit) in the data center. It considers each 'pillar' as a separate improvement area with their own KPI (GEC - Green Energy Coefficient, PUE - Power Usage Effectiveness, ITEE - IT Equipment Energy Efficiency, ITEU - IT Equipment Utilization) that can be measured. This approach also has a strong operational focus.
>
> Both examples can be seen as one specific implementation of the *4 Pillar Framework*. They show that the framework can be the foundation for creating energy

efficiency models for specific data centers and for modeling data center energy flow chains. This is possible because the *4 Pillar Framework* acknowledges that different data centers have different goals and requirements, that applications play an important part for the energy efficiency of HPC data centers, and that the addition of cross pillar interactions allow for more fine-tuned energy efficiency related decisions.

## 2.5 Summary

A wholistic view is required to improve a data center's energy efficiency. The *4 Pillar Framework* provides that view. It defines a generic frame of reference that allows data centers to approach energy efficiency as a well defined process. This process needs to be adapted to data center specific requirements and constraints. Using this framework helps to systematically identify possible opportunities to improve a data centers energy efficiency and leads to a better understanding of the complexity of the task. By using the *4 Pillar Framework*, gaps in existing work can be identified and addressed (see section 3.3 A New Metric to Measure Data Center Energy Efficiency - Data Center Energy Efficiency (DCEE)).

# 3 Measuring Energy Efficiency

> In God we trust, all others must bring
> data.
> _____
> — American Statistician W. Edwards
> Deming

This chapter will discuss the challenge of collecting and using power and energy data to determine data center energy efficiency. It will introduce the developed tool Power-DAM which is used to collected data from all Pillars of the data center (partly based on the author's publications [50], [51], and [47]). This chapter will also discuss metrics that are used for data center energy efficiency assessments and will introduce two new metrics, namely Data center Workload Power Efficiency (DWPE) (partly based on the author's publication "DWPE, a new data center energy-efficiency metric bridging the gap between infrastructure and workload" [52]) and Data Center Energy Efficiency (DCEE).

## 3.1 What is Energy for a Data Center?

Electrical energy is the product of power consumption over time. For a data center this is the time integral of the total power going into the data center.

$$\text{Energy} = \text{Power} * \text{Time}$$

According to [53], Energy efficiency is:

- The ratio of the energy delivered (or work done) by a machine to the energy needed (or work required) in operating the machine. The efficiency of any machine is always less than one due to forces such as friction that use up energy unproductively.

- The ratio of the effective or useful output to the total input in any system.

For a data center, energy efficiency is a way of managing and restraining the growth of its energy consumption. A data center can be more energy efficient if it delivers more services for the same energy input or the same services for less energy input. For example, when an application uses less energy to solve the same problem as another application it is considered to be more energy efficient. Similarly, a cooling loop with a higher coefficient of performance (COP) is more energy efficient because it uses less electrical power to deliver the same amount of thermal cooling energy.

$$\text{Energy Efficiency} = \frac{\text{Work done}}{\text{Total energy consumed}}$$

For High Performance Computing, *WORK* can be defined in general as executing scientific applications. For each application, energy efficiency can be defined (see Figure 3.1 ) on the IT node level (e.g. how energy efficient a specific technology is (CPU, GPU, etc.)), on the IT system level (e.g. how energy efficient is the complete system including internal infrastructure overhead), and on the data center level (e.g. how energy efficient was the application when data center infrastructure overhead is included; this translates directly into how much an application costs). The author's journal paper "Analysis of the Efficiency Characteristics of the First High-Temperature Direct Liquid Cooled Petascale Supercomputer and Its Cooling Infrastructure" [54] shows this in measurements for the LRZ SuperMUC system.

Figure 2.3 shows that power in a data center is mainly used to run the data center infrastructure (*Pillar 1*) and the IT systems (Pillar 2) which are used to perform the work submitted by the users. The distribution (and associated effectiveness) of power throughout the data center depends strongly on the required redundancy levels and power type (AC vs. DC). Figure 3.1 provides an overview of a typical HPC data enter power distribution taking the LRZ data center as example. Power comes in from the utility connection (*Data Center Power*) and is distributed via additional power equipment (such as backup systems) to the outlets where the IT systems are connected (*Wall*). In parallel, power is provided to the data center cooling infrastructure controlled by a building automation system (*Infrastructure Automation*).

Energy efficiency can be defined at all levels. The IT systems consume a specific amount of energy to run a user application. Here energy efficiency improvement means to run the same application with less energy. Also, the cooling infrastructure consumes energy depending on the power needed to remove the heat generated by the IT systems. Energy consumption is very closely related to the Coefficient of Performance (COP) of the different cooling technologies used in the data center (see Figure 2.6 for an overview of the LRZ cooling infrastructure).

$$\mathrm{COP} = \frac{\mathrm{Q}_{\text{Cold energy delivered}}}{\mathrm{P}_{\text{Total power consumed}}}$$

The less power one needs for generating cold (removing heat) the more efficient the cooling infrastructure. The COP varies with time since part of it is the exchange of heat with the environment which depends on outside properties such as weather and climate.

For the power distribution infrastructure, power effectiveness can be defined. Its not efficiency since the power distribution is only used to transfer power. The effectiveness is defined as:

$$\text{Power effectiveness} = \frac{\mathrm{P}_{\text{delivered}}}{\mathrm{P}_{\text{supplied}}}$$

The best effectiveness is 100%; meaning that no power losses occur over the power distribution.

For the complete data center, a well know effectiveness measurement is PUE (for detailed explanation see subsection 3.3.1 Existing metrics for assessing energy efficiency).

Currently, there are two main concerns for HPC data centers. One is the increased and varying power consumption of flagship supercomputers and how this will effect the data

center infrastructure [55]. The other is the complete energy consumption of the data center. This a major concern related to the sustainability and operating cost of data centers. For example, energy consumption translates directly into costs for LRZ (see Figure 1.3), therefore, saving energy translates directly into lower operating costs.

Power and energy measurements in *Pillar 1* have a strong seasonal (cooling infrastructure) and load (power conversion and distribution) dependency, therefore, for data center characteristics a yearly average is used. Whereas measurements in *Pillar 2* mainly show the energy efficiency of the IT system for all characteristics (compute, I/O, Memory, Network) of an application. Here fine grain node measurements are used for application profiling [56] and coarser grain measurements are used for system characterization [57].

## 3.2 Wholistic Power and Energy Data Collection - PowerDAM

The first step for a wholelistic approach is the availability of measurement data from all Pillars. Since no such tool existed, LRZ started the development of PowerDAM (Power Data Aggregation Monitor) in 2012.

Many monitoring systems already exist in a data center that collect different power and energy sensor information. Figure 3.1 shows a generic overview of the most common monitoring and data collection areas related to the electrical power distribution and heat flow (Figure 2.3).



Figure 3.1: Generic overview of Data Center monitoring and sensor systems.

The areas in Figure 3.1 are directly related to the *4 Pillar Framework*. Areas 1 and 2 provide information related to *Pillar 1*. Area 3 provides information concerning Pillar 2. Area 4 is related to *Pillar 3* and *Pillar 4* since performance counters are provided by the OS and are mainly used in application performance analysis. Currently, each area has its own specialised tools which PowerDAM does not replace.

The main idea behind PowerDAM is shown in Figure 3.2. Instead of replacing the func-

Figure 3.2: LRZ data consolidation overview.

tionality of already existing tools, PowerDAM is used to centralize all sensor data and measurement information related to power and energy from all 4 Pillars. This is done via plug-ins in PowerDAM and remote agents for each monitored system. Having a central data repository allows for generic data processing and analysis (for example, the calculation of an application's Energy-to-Solution) and the calculation of important KPIs such as PUE, ERE, COP, etc. Some of the major KPIs related to data center power and energy efficiency are discussed in subsection 3.3.1 Existing metrics for assessing energy efficiency.

At LRZ, data from the building infrastructure, some HPC systems, and used re-use technology are collected. Figure 3.3 shows the detailed systems for each pillar monitored via PowerDAM at LRZ.

The first PowerDAM version (PowerDAM V1.0) was developed to calculate EtS for the CoolMUC PRACE 1IP-WP9 prototype [58] collecting data from *Pillar 2* and *3*. It collected power data from the CoolMU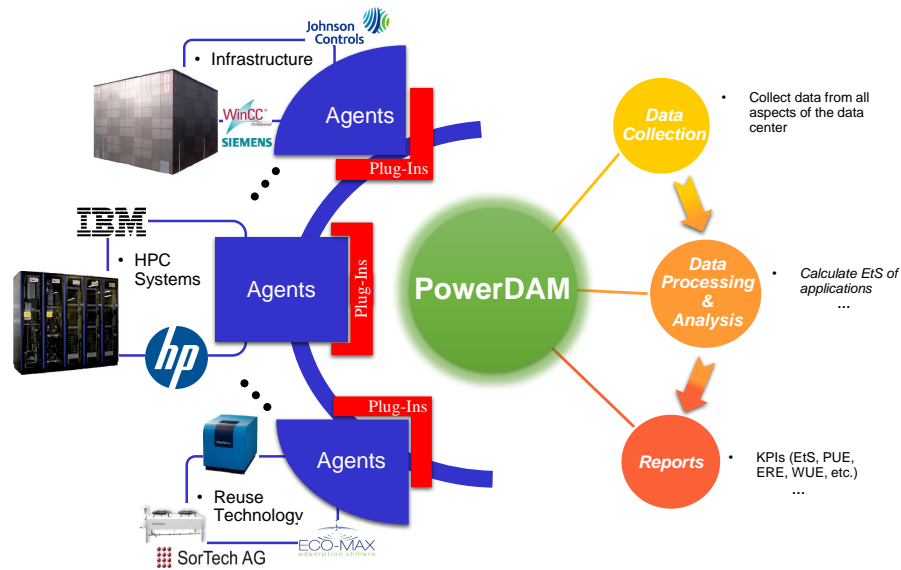C cooling infrastructure and the management software Clustware, job data from SLURM [23], and some additional information such as CPU temperature via the IPMI interface. The EtS calculation included the system cooling share based on the proportion of node power and overall system node power for each timestep and the network consumption share based on the fraction of number of job nodes and overall active system nodes for each timestep.

PowerDAM V2.X was developed during the SIMOPEK project [46]. It extended the previous capabilities to include data from *Pillar 1*. The new version is able to collect data from the LRZ building infrastructure automation system, Johnson Controls, and the power infrastructure monitoring system WinCC. In *Pillar 2* additional HPC systems were added including SuperMUC Phase1 and Phase2, CoolMUC2, and the iCinga cluster monitoring tool used by IBM for monitoring the adsorption chillers connected to CoolMUC2. And in *Pillar 3* the connection to the LoadLeveler scheduling system was added.

Figure 3.3: PowerDAM V1 and V2 interfaced LRZ systems overview.

## 3.2.1 PowerDAM Design

One challenge with the analysis of data from different systems is that the access and data formats are all different. Also, changes to the data structure of monitored systems might require a substantial code re-write if the data access and format is not abstracted. Therefore, PowerDAM uses its own unified data format standardizing sensor data access for data processing and analysis tools. Figure 3.4 shows this.



Figure 3.4: PowerDAM high level design overview.

The sensors from a monitored system are transcoded into a tree-like structure. The Root

Resource is the system name (like jci for the Johnson Controls Building Automation System). The depth of the Resource tree depends on the availability of exploitable layers. For jci, the circuit (position 4-8 in the jci sensor name schema Figure 3.8), like KLT72, and the device in each circuit (position 9-10 in the jci sensor name schema) provide a natural tree abstraction. Sensors of a specific SensorType can be attached to each resource. This can be physical sensors, like temperature and electrical power, but also virtual sensors, sensors that combine other sensors via a mathematical formula [59].

PowerDAM internal sensors can be either accessed directly via their sensor name or through the Resource tree structure. The internal API (communication between remote Agents and PowerDAM) uses this format in addition to value and timestamp information. If a sensor name changes in the source system its old information is still accessible in PowerDAM. To support the collection of only a subset of system sensors and to allow for sensor removal/addition/re-naming, PowerDAM uses a xml file that is generated during system 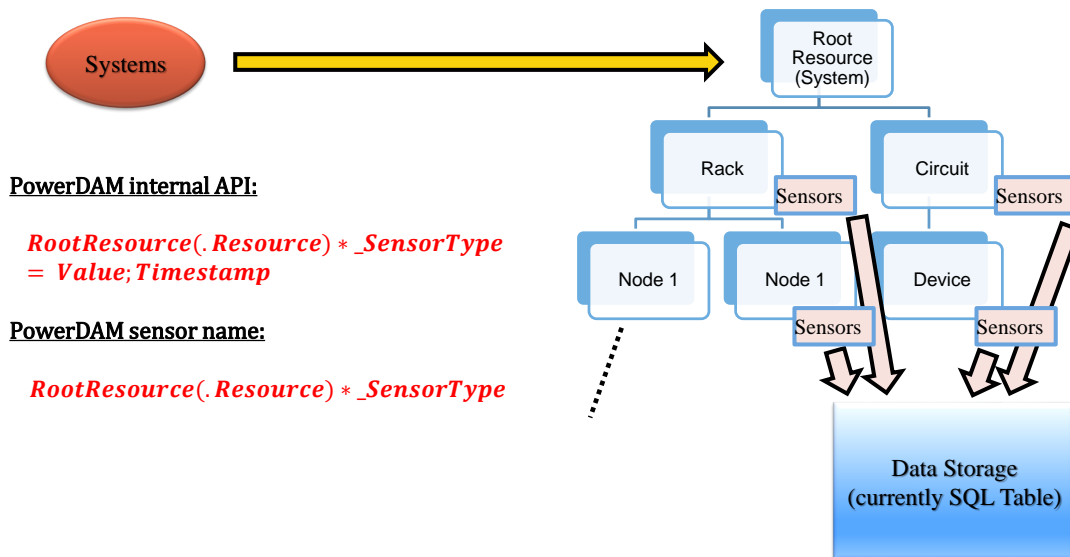initialization which is the first step required to add a new system to PowerDAM. Figure 3.6 shows the first part of the xml file.

```xml
1   <?xml version="1.0" encoding="utf-8"?>
2   <SensorList version="1.1">
3     <head>
4       <title>JCI Sensor List</title>
5       <dateCreated>2015-03-19 17:34:33.777677</dateCreated>
6       <dateModified>2015-06-23 18:10:37.196000</dateModified>
7     </head>
8     <sensors orgRecordCount="2297">
9       <sensor active="1" collect="0" id="2329">
10        <datatypeName>Analog</datatypeName>
11        <orgName>.REGKLT71VT__ANST01</orgName>
12        <powerDamName>jci.REGKLT71.VT01_ST__AN_gear</powerDamName>
13        <location>R.EG.</location>
14        <powerDamConversionMultiplier>1.0</powerDamConversionMultiplier>
15      </sensor>
16      <sensor active="1" collect="0" id="1251">
17        <datatypeName>Analog</datatypeName>
18        <orgName>H3ORLT01VRAB__ST01</orgName>
19        <powerDamName>jci.H3ORLT01.VR01_ST_AB_binary</powerDamName>
20        <location>H.3O.AB</location>
21        <powerDamConversionMultiplier>1.0</powerDamConversionMultiplier>
22      </sensor>
23      <sensor active="1" collect="0" id="1324">
24        <datatypeName>Analog</datatypeName>
25        <orgName>H3ORLT01VRAB__ST02</orgName>
26        <powerDamName>jci.H3ORLT01.VR02_ST_AB_binary</powerDamName>
27        <location>H.3O.AB</location>
28        <powerDamConversionMultiplier>1.0</powerDamConversionMultiplier>
29      </sensor>
```

Figure 3.5: PowerDAM sensor XML file syntax.

The main structure for the SensorList consists of the head area (storing file specific information) and the sensors area (where the sensor information is stored).

Information encoded in the xml file includes:

- SensorList version - this allows for future changes in the xml file syntax

- dateCreated - shows the date time when the xml file was generated

- dateModified - shows the date time when the sensor information in the xml file was last updated which helps to answer, very quickly, the question of whether the file was updated after sensors in the monitored system changed

- orgRecordCount - shows the number of sensors found for the monitored system which is very helpful to debug sensor name transcoding issues

- active - if set to "1" than the sensor exists in the system, if set to "0" the sensor no longer exists

- collect - if set to "1" than the sensor data should be sent to PowerDAM, if set to "0" no data needs to be collected

- id - a programmer definable entry that can be used to simplify access to the sensor data (here the jci data base internal sensor id is used which speeds up sensor data access and helps with sensor data related debugging)

- orgName - the original sensor name in the source system

- powerDamName - the sensor name transcoded into the PowerDAM sensor schema

- powerDamConversionMultipier - since sensor types in PowerDAM have a fixed unit (like W) other sensor values need to be converted (for a sensor value in kW this would be 1000)

```
1   # -*- coding: utf-8 -*-
2
3   [Collect]
4
5   # sets the group names for which all sensors should be collected
6   GROUP_NAMES=R2OKLT72, R2OSS_71, R2OKLT71, R3OKLT71, R3OKLT72, R2OKLT13, R2OKLT14, R3OKLT11, R2OKLT12,R3OKLT12,R1OKLT13,R1OKLT14,R4OKLT
7
8   # sets specific sensors that should be collected, needs to be the complete name shown in the JCI db
9   #DEVICE_NAMES=R4OKLT11RK__T_MW01, @JCSQL:NAE054-01:NAE054-01/N2 Trunk 1.NAE054-MIG118.REGSS_60EZ__PEMW04.Present Value, \
10  #              @JCSQL:NAE054-01:NAE054-01/N2 Trunk 1.NAE054-MIG122.REGSS_60EZ__PEMW04.Present Value, REGSS_60EZ__PEMW04
11  DEVICE_NAMES=R4OKLT11RK__F_MW01, R4OKLT11RK__T_MW01, @JCSQL:NCE052-01:NCE052-01/Programming.REGKLT70.REGKLT60AN__VPMW02.Present Value
12
13  # sets sensor names that should be ignored (no longer active in JCI), needs to be the complete name shown in the JCI db
14  INVALID_NAMES=R4OKLT17RK__LTMW01,R4OKLT14AN__VPSW01,R4OKLT13AN__VPSW01,R4OKLT12AN__VPSW01,R4OKLT11AN__VPSW01
```

Figure 3.6: PowerDAM sensor activation config file.

The activation of sensors for collection is currently done via an additional system specific config file (Figure 3.6). This is done to support the future plan of allowing a PowerDAM GUI to activate/deactivate sensors for collection.

Three sensor activation options are available. Adding an entry to the GROUP_NAMES activates all sensors that belong to the group. Individual sensor names that should be activated can be added to the DEVICE_NAMES list. Sensors that should be explicitly disabled for collection can be added to the INVALID_NAMES list.

Since PowerDAM provides a plug-in infrastructure for reports, different exporting options of the sensor data can be implemented (e.g. graphical reports and different file formats). For example, during the SIMOPEK project a MYNTS [60] specific sensor xml file needed to be generated in order to feed sensor data into MYNTS. Figure 3.7 shows part of the xml file.

Since PowerDAM is able to mark data as invalid (in some cases) it can provide this information to MYNTS as well.

```
 1 <root>
     <resource name="jci.KLT16.F01_MW_T">
       <sensorType id="1210" key="temperature" unit="°C">
         <sensorEntry value="12.1347761154" timestamp="2015-04-03 00:00:17" quality="valid"/>
         <sensorEntry value="12.3123874664" timestamp="2015-04-03 00:01:17" quality="valid"/>
         <sensorEntry value="12.3123874664" timestamp="2015-04-03 00:02:18" quality="valid"/>
         <sensorEntry value="12.3123874664" timestamp="2015-04-03 00:03:19" quality="valid"/>
         <sensorEntry value="12.193980217" timestamp="2015-04-03 00:04:19" quality="valid"/>
         <sensorEntry value="12.3123874664" timestamp="2015-04-03 00:05:20" quality="valid"/>
10       <sensorEntry value="12.3123874664" timestamp="2015-04-03 00:06:21" quality="valid"/>
         <sensorEntry value="12.3123874664" timestamp="2015-04-03 00:07:22" quality="valid"/>
         <sensorEntry value="12.3123874664" timestamp="2015-04-03 00:08:22" quality="valid"/>
         <sensorEntry value="12.3123874664" timestamp="2015-04-03 00:09:23" quality="valid"/>
         <sensorEntry value="12.3123874664" timestamp="2015-04-03 00:10:24" quality="valid"/>
         <sensorEntry value="12.3716554642" timestamp="2015-04-03 00:11:24" quality="valid"/>
         <sensorEntry value="12.3716554642" timestamp="2015-04-03 00:12:25" quality="valid"/>
         <sensorEntry value="12.3716554642" timestamp="2015-04-03 00:13:26" quality="valid"/>
         <sensorEntry value="12.2531833649" timestamp="2015-04-03 00:14:27" quality="valid"/>
         <sensorEntry value="12.3716554642" timestamp="2015-04-03 00:15:28" quality="valid"/>
20       <sensorEntry value="12.3716554642" timestamp="2015-04-03 00:16:28" quality="valid"/>
         <sensorEntry value="0.0" timestamp="2015-04-03 00:17:29" quality="invalid"/>
         <sensorEntry value="12.3716554642" timestamp="2015-04-03 00:18:29" quality="valid"/>
```

Figure 3.7: PowerDAM sensor file generated for the SIMOPEK project.

### 3.2.2 PowerDAM Lessons Learned

During the current usage of PowerDAM at LRZ some challenges were encountered. These are most likely not limited to the LRZ setup.

The first challenge was related to sensor name handling. Figure 3.8 shows the sensor name schema for JCI.

At first glance the schema seemed well defined. Letters 4-8 identify the circuit and letters 1-3 identify the location. The first PowerDAM mapping was based on the circuit name. Unfortunately, it turned out that a very small number of sensors used the circuit+location as a unique identifier. For example, R3OKLT72 (Location: R3O, Circuit KLT72) and R2OKLT72 (Location: R2O, Circuit KLT72) have the same Circuit name (KLT72) but are physically not the same circuit. Also, there were other naming inconsistencies that required special name handling making the final solution more complicated than it should have been.

For example, the sensor name: IUGSS_41EZ__PEMW04 follows the naming schema and maps to: the PowerDAM sensor name jci.IUGSS41.EZ04_MW__PE_power.

The sensor name: @JCSQL:NAE054-01:NAE054-01/N2 Trunk1.NAE054-MIG136.IUGSS_41EZ__PEMW04.Present Value on the other hand does not adhere to JCI internal naming schema and, therefore, maps to: jci.IUGSS41.EZ04_MW__PE_power as well.

The second JCI related challenge was related to the replacement of the physical JCI server with a more powerful virtual machine. Here the database was not copied from the old machine to the new one but a fresh database was initialized. Accessing the new data base with the old sensor xml file worked fine but generated some strange values for certain sensors. An investigation showed all internal sensor ids had changed and a re-initialization of the sensor xml file was required. An additional challenge was that the old data base was copied under a different database name to the new server requiring different database access information than the new active database. This was not foreseen in the initial design of PowerDAM. The latest version of PowerDAM supports multiple data

Figure 3.8: LRZ JCI sensor name schema.

base access configurations and sensor xml files for one system defined using timestamp ranges; this is only required for backfilling sensor data and not for accessing sensor data already stored in PowerDAM.

WinCC provided other challenges. For example, because of a WinCC system update, the exported sensor data for that timeframe was done in two lines instead of the standard one line for each timestamp. Figure 3.9 shows a short segment of the exported data file.

Zeit;MS Nord MWh;MS Süd MWh;Trafo 1 MWh

21.03.2015 05:30:00;;;

21.03.2015 05:30:00;60793,527;57764,809;26884,48

21.03.2015 05:35:00;;;

21.03.2015 05:35:00;60793,703;57764,969;26884,488

21.03.2015 05:40:00;;;

21.03.2015 05:40:00;60793,879;57765,133;26884,488

Zeit;MS Nord kW;MS Süd kW;kW Trafo 1

21.03.2015 05:30:00;2102,07;-1979,7;

21.03.2015 05:30:00;;;54,8

21.03.2015 05:35:00;2098,35;1975,45;

21.03.2015 05:35:00;;;54,3

21.03.2015 05:40:00;2098,17;1973,19;54,5

21.03.2015 05:40:00;;;

Figure 3.9: PowerDAM WinCC double line bug.

For the energy file (left segment) an additional empty line was inserted during export. For the power file (right segment) one sensor provided data on its own line but was part of the complete line later on. PowerDAM can now deal with sensor data for one timestamp

distributed over multiple lines.

Summarizing the experiences of PowerDAM at LRZ the following recommendations can be made:

- The data center operations and monitoring groups (as well as the software engineering team) should be involved in defining a global data center sensor naming schemas (valid for all monitoring systems) as well as a process on how to add new sensors later on. This will simplify the automatic processing of data by other data center tools.

- When considering names and descriptions, stick to the English alphabet everywhere in the data center monitoring and control systems. This creates less worries when processing the data later on. Otherwise, be aware of possible UTF-8 related issues.

- Avoid data center monitoring and automation tools that require proprietary data access tools. If this is not possible, any extra tools and installations needed to access the data should be included in the procurement.

- Collected sensor data should be treated with an appropriate amount of skepticism. In most automation systems only a limited amount of sensors are used for controls. If one starts to collect non-critical (not used by the control system) sensor data and uses those for system analysis sometimes results are generated that are not physically possible. For example, during the early LRZ data center power and energy analysis the data center generated power. As it turned out, some of the internal power measurement equipment didn't work correctly. The LRZ experience has shown that the correct explanation for strange behavior is in most cases invalid sensor data.

### 3.2.3 PowerDAM Future Work

During the PowerDAM usage it became clear that adding new systems for collection is cumbersome and that the pulling mechanism is not scaling well with the increase of the number of monitored systems. Therefore, a new publish-subscribe communication model between the collection agents and PowerDAM was designed. Figure 3.10 shows the use-cases for the publish-subscribe communication model.

The main idea is to separate agents and PowerDAM completely, meaning that new monitoring agents do not require any PowerDAM source code changes. An agent first registers with PowerDAM. Upon successful registration sensor data will be published by the agent. The agent can de-register or update their own registration. PowerDAM can request data backfilling from the agent specifying the timeframe and sensors involved. The agent and PowerDAM can use *Out of Bound* communication to inform each other of important information such as when new sensors were added to the system monitored by the agent or if the PowerDAM user changes the collection sensor selection or sensor collection interval.

Each Use-Case can be further defined using activity diagrams. Figure 3.11 shows the high level activity diagram for the agent registration use case.

The starting point is the start of the agent. If the agent hasn't already registered it assembles and sends the registration message to PowerDAM. After receiving a positive acknowledgment from PowerDAM it sends it's system sensor information. After a successful processing of the sensor information by PowerDAM it saves it's registration information and

Figure 3.10: PowerDAM Publish-Subscribe communication use-cases.

Figure 3.11: PowerDAM Publish-Subscribe agent registration.

writes a registration successful message in it's logfile. In all other cases an error message will be written and the agent terminates without registering with PowerDAM.

## 3.3 A New Metric to Measure Data Center Energy Efficiency - Data Center Energy Efficiency (DCEE)

This chapter is partly based on the author's publication "DWPE, a new data center energy-efficiency metric bridging the gap between infrastructure and workload" [52].

As mentioned in chapter 3.1, data center energy efficiency can be defined as the whole data center energy consumption related to the run of an application. LRZ is interested in answering the following questions related to energy efficiency:

1. What is the cost for an application run?

2. How energy efficient would a specific HPC system technology be at LRZ?

Current metrics lack the complete data center coverage required to answer either question. For example, the Green500 list (November 2013) results for 2 different HPC systems (SuperMUC Phase1 at LRZ and PLX at CINECA) are shown in Table 3.1.

Both systems have a very similar Green500 ranking. This might imply that the answer to the above questions will be the same for both systems. But is this actually true?

|              | SuperMUC | PLX |
|--------------|----------|-----|
| Rank (Nov 2013) | 95    | 98  |
| WPE (MFLOPS/W)  | 908   | 892 |

Table 3.1: Green500 list rank for four different HPC systems and DWPE for LRZ

### 3.3.1 Existing metrics for assessing energy efficiency

Figure 3.12 shows the coverage areas for a selection of important power/energy related efficiency metrics using the *4 Pillar Framework*. A solid olive green box indicates primary power or energy related metrics. Metrics in a doted gray box are not directly related to energy efficiency but measure additional factors one needs to keep in mind when considering a data center's overall power and energy efficiency. An explanation of each metric follows.



Figure 3.12: Currently popular metrics related to data center energy efficiency.

Power Usage Effectiveness (PUE) [4] is the ratio of the total power consumed by a data center to the total power consumed by the IT equipment inside the data center. It shows the ratio of power coming into the data center (*Pillar 1*) from the energy provider to the power going into the IT equipment (*Pillar 2*). PUE is defined as an annual average since it is prone to external influences (e.g. temperature, humidity, etc.).

$$\text{PUE} = \frac{\text{Total power consumed by the data center}}{\text{Total power consumed by the IT equipment}}$$

Energy Reuse Effectiveness (ERE) [61] is adding an energy re-use component to PUE. Opposite to PUE, ERE can go below 1 and can be 0 in the best case (all energy consumed by the data center is re-used).

$$\text{ERE} = (1 - \frac{\text{Reuse Energy}}{\text{Total Energy consumed by the data center}}) * \text{PUE}$$

IT-power Usage Effectiveness (ITUE)[28] is similar to PUE but for the HPC system covering Pillar 2. It is defined as the ratio of the power coming into the system from *Pillar 1* and the power consumed by HPC IT. The metric was introduced by the Energy Efficient HPC Working Group (EEHPCWG) [62]. It is similar to IT Equipment Utilization (ITEU) defined by the Japan's Green IT Promotion Council [63].

$$\text{ITUE} = \frac{\text{Total power consumed by the HPC system}}{\text{Total power consumed by HPC IT}}$$

Total-power Usage Effectiveness (TUE) [28] as defined by the EEHPCWG connects *Pillar 1* and *Pillar 2* by combining ITUE and PUE.

$$\text{TUE} = \text{ITUE} * \text{PUE}$$

Green Energy Coefficient (GEC) as defined by the Global Metrics Harmonization Task Force [64] quantifies the portion of a facility's energy that comes from green sources. It connects *Pillar 1* with utility providers.

$$\text{GEC} = \frac{\text{Green energy used by the data center}}{\text{Total data center source energy}}$$

IT Equipment Energy Efficiency (ITEE) as defined by the Japanese Green IT Promotion Council [63] tries to capture the capability/capacity of the IT equipment in relation to its power consumption and, therefore, is a metric for *Pillar 2*. ITEE is not a measurement but is taken from manufacture provided specifications. If one can have more capability/capacity with the same amount of power the equipment is considered to be more efficient.

$$\text{ITEE} = \frac{\text{Total server capacity}}{\text{Rated power of IT equipment}} + \frac{\text{Total storage capacity}}{\text{Rated power of IT equipment}} + \frac{\text{Total network equipment capacity}}{\text{Rated power of IT equipment}}$$

The Data center Performance Per Energy (DPPE) metric was defined by the Japanese Green IT Promotion Council [63] to provide a more wholistic data center energy efficiency metric. It combines ITEU, ITEE, PUE, and GEC.

$$\text{DPPE} = \frac{\text{ITEU} * \text{ITEE}}{\text{PUE} * (1 - \text{GEC})}$$

The Performance per Power (Watt) metric is a measurement of the energy efficiency of the compute hardware used in an HPC system for a specific workload. The more Performance per Watt a system can deliver the more energy efficient it is. The metric covers

Pillar 2 (in most cases only the IT part), Pillar 3, and Pillar 4. A well known example of this metric is the FLOPS/W metric used by the Green500 List [29].

$$\text{Performance/W} = \frac{\text{average achieved performance}}{\text{average HPC IT power used}}$$

Energy to Solution (EtS) is another metric for measuring the energy efficiency of an HPC system for a specific workload. It is complementary to the Performance per Power metric. It shows the energy consumed by an application for solving a specific problem but doesn't require performance details. One can find different definitions in the literature. Minartz et al. [65] doesn't really define EtS but hints to:

$$\text{EtS}_{\text{application}} = \int_{i=startTime}^{endTime} \sum \text{node power}_{\text{i}}$$

Auweter et al. [14] considers only the compute node power consumption and defines EtS as:

$$\text{EtS}_{\text{application}} = \text{average power per node} * \#\text{nodes} * \text{runtime}_{\text{application}}$$

Shoukourian et al. [51] takes the power supply efficiency, system network power consumption, and system cooling power consumption into consideration defining EtS as:

$$EtS(J, S) = \sum_{i=startIteration}^{endIteration} \Delta t_i \cdot P_i(J, S)$$

where $P_i(J, S)$ includes: the average node power consumption (measured at system PDU outlets) times the number of nodes used by the application; a fraction of the average system cooling power consumption; and a fraction of the system networking power consumption depending on the number of nodes used by the application vs the overall number of nodes of the system.

Data Center Performance Efficiency (DCPE) was defined by the GreenGrid [66] to be the ultimate energy efficiency metric for a data center.

$$\text{DCPE} = \frac{\text{total useful work performed}}{\text{total power consumed by the facility}}$$

Total Cost of Ownership (TCO) is a metric that indicates what the expenses are for running a system throughout its lifetime. It includes investment and running costs as well as personal costs. This is not a metric that shows energy efficiency but it has a strong relationship with energy efficiency. The more energy efficient a data center is the lower TCO becomes if all other costs stay the same.

### 3.3.2 Metric shortcomings

As Figure 3.12 shows, no metric covers the complete data center with the exception of DCPE.

PUE only measures the power distribution and cooling infrastructure effectiveness of a data center. It does not reflect the power efficiency of the IT equipment operated in a data

center. In the worst case, a good PUE could mean that the data center can waste IT power very effectively.

ITUE, similar to PUE, does not indicate if the power used by HPC IT is utilized efficiently, i.e. it doesn't indicate if a system is energy efficient for a particular workload mix.

Since TUE combines PUE and ITUE it doesn't consider the IT equipment energy efficiency and is, therefore, not a good measure for HPC data center energy efficiency.

GEC is not an indicator of data center energy efficiency since the power source does not determine the energy efficiency of the data center.

ITEE doesn't indicate how energy efficient the IT components are for the specific work they need to do. The available capacity does not influence how well an HPC system is suited for running different workloads. Therefore, ITEE is no indicator of a HPC data center energy efficiency.

DPEE is a combination of metrics that do not consider the Workload-mix energy efficiency and, therefore, is not an indicator for HPC data center energy efficiency.

By itself, Performance/W is not an indicator of how energy efficient a data center is because it does not include either the power losses in the power distribution infrastructure or the power of the data center cooling infrastructure needed for the operation of the IT equipment. Also there is no formal definition requiring the inclusion of more than the power consumption of the compute node IT components.

EtS is similar to the Performance/W metric in the sense that it doesn't include the additional power needed to operate the IT equipment in a data center and, depending on its definition, it might not include system cooling power, system power losses, or networking power consumption. EtS has the potential for a complete data center metric but no formal definition exists [67]. The work presented in chapter 3.3.4 and 3.3.5 on the new metrics which extends the Performance/W metric can also be applied to the EtS metric.

DCPE is the only metric that covers the whole data center. Unfortunately, no one has been able to define what the term "total useful work performed" really means, or how it can be measured. ITEE is one possible definition but it is not applicable for HPC. An additional challenge for DCPE is that one would need to collect and analyze data from all parts of the *4 Pillar Framework*.

Despite the lack of a comprehensive metric, PUE and FLOPS/W are often used to indicate wether a data center is energy efficient. Nevertheless, these two metrics are a good starting point but a new metric is needed that can extend the current Performance/W metric to cover the complete data center.

### 3.3.3 Metric proposal

Data center Workload Power Efficiency (DWPE) is intended to be an energy efficiency metric for one specific workload covering the complete data center. It can be seen as one instance of DCPE. DWPE makes the connection between workload power efficiency of the HPC system and the data center infrastructure by combining the Performance/W metric of the IT system with the data center overheads specific to the IT system. DWPE will help data centers track their energy efficiency over time including:

- assessing when an HPC system needs to be replaced

- assessing the energy efficiency of new machines to determine wether it will be a good fit for the needs of the data center

- assessing the impact on the data center energy efficiency if the workload-mix changes

To help with the definition of DWPE one additional metric is introduced, namely Workload Power Efficiency (WPE) which is defined in the section 3.3.4.

Figure 3.13 shows the areas the new metrics are covering. WPE extends the Performance/W metric to include the complete HPC system (*Pillar 2*). DWPE combines WPE with the data center overhead incurred by running the HPC system (system PUE or sPUE) and, therefore, shows the energy efficiency for one workload for the complete data center. Multiple DWPE's can be combined to show the energy efficiency for a particular workload mix in a data center. From this, the metric Data Center Energy Efficiency (DCEE) can be derived. It is the first real workload-mix energy efficiency metric for a data center and is especially useful during the procurement of new systems.



Figure 3.13: Proposed metrics WPE and DWPE.

### 3.3.4 Workload Power Efficiency (WPE)

Workload Power Efficiency (WPE) extends the Performance/W as well as the EtS metric to include all aspects of *Pillar 2* (figure 3.13). This can be done by using the IT Performance/W result and combining it with ITUE. Another option would be to directly measure the complete HPC system power consumption, including the HPC system infrastructure and HPC system IT, during the benchmarking process.

$$\text{WPE} = \frac{\text{average achieved performance}}{\text{average HPC system power used}}$$

Table 3.2 highlights why WPE is an important metric when talking about HPC system power efficiency. The table shows the HPL core phase result for SuperMUC Phase1 submitted to the Top500 List. As can be seen, the WPE measurement shows a lower efficiency rating then the Green500 measurement would indicate. This also highlights that the Green500 list is not a measure of system efficiency but a measure for IT component efficiency for running HPL.

|  | SuperMUC Phase1 Flops/W |
| --- | --- |
| Green500 (IT part of system only) | 908 |
| WPE | 855 |

Table 3.2: Green500 vs WPE result for SuperMUC Phase1 (HPL core phase only)



Figure 3.14: SuperMUC Phase1 Green500 HPL run power consumption graph.

Figure 3.14 shows the power consumption of SuperMUC Phase1 during the Green500 run for different power measurement levels. As expected, the IT only power (Power (PDU, kW)) shows the lowest power consumption resulting in the highest Green500 score (908 Flops/W). The power consumption that includes the SuperMUC Phase1 cooling infrastructure (Power (Machine Room, kW)) is slightly higher, resulting in an decreased Green500 score (855 Flops/W). The complete power consumption which includes the data center infrastructure power consumption (Power (Infrastructure, kW)) is around 20% higher

than the Machine Room power consumption. To calculate the resulting Green500 score, the new metric Data center Power Efficiency (DWPE) is needed.

### 3.3.5 Data center Power Efficiency (DWPE)

Data Center Workload Power Efficiency (DWPE) can be used to calculate the HPC system efficiency for a specific workload in a given data center.

It is tempting to just divide WPE by PUE (equation 3.1). PUE, however, depends on multiple factors like cooling power used in the data center and the consumption of the IT equipment. In reality, the PUE of a data center will change if the current HPC system is replaced with a different system. This is especially true if the cooling technology changes. Each Watt of electrical power going into a HPC system is converted into heat and needs to be removed by the data center cooling system. Current systems mainly use three different cooling technologies:

1. air cooling - where heat is removed using air as a transfer medium

2. chilled water cooling - where heat is removed directly (e.g on-chip cooling) or indirectly (e.g. rear door heat exchanger) via cold water which is generated using chillers

3. chiller-less direct water cooling - where heat is removed directly via water which is cooled without the use of chillers (also called free cooling)

Each of these different cooling technologies has a different cost (e.g. how much additional power is needed to remove 1W of heat from the system/data center).

Due to the principle of energy conservation, the electrical power supplied to the IT system is converted into heat. Therefore, one can substitute IT equipment power ($P_{IT}$) in PUE (equation 3.1) with IT Heat Quantity ($Q_{IT}$) (equation 3.2). Additionally, the power used by each cooling technology can be replaced (equation 3.3). The resulting equation 3.4 shows clearly that the PUE of a data center is strongly dependent on the cooling technologies used in the IT equipment and the efficiency of the data center cooling infrastructures.

$$\text{PUE} = \frac{\sum\limits_{i=1}^{n} P_i}{P_{IT}} \tag{3.1}$$

Where *i* is: IT power consumption, electrical power distribution and conversion losses (PDCL), and cooling infrastructure power consumption. Other power consumers (like attached offices, lights, etc.) are not considered since they present a very small fraction of the HPC data centers and can be easily added as a fixed factor if required similar to the power delivery and conversion losses $\frac{P_{PDCL}}{P_{IT}}$ (equation 3.4).

$$P_{IT} = Q_{IT} = Q_{hotwater} + Q_{coldwater} + Q_{air} \tag{3.2}$$

$$\text{COP} = \frac{Q}{P_{used}} \qquad => \qquad P_{used} = \frac{Q}{\text{COP}} \tag{3.3}$$

$$\text{PUE} = 1 + \frac{P_{\text{PDCL}}}{P_{\text{IT}}} + \frac{1}{Q_{\text{IT}}} * \frac{Q_{\text{hotwater}}}{\text{COP}_{\text{hotwater}}} + \frac{1}{Q_{\text{IT}}} * \frac{Q_{\text{coldwater}}}{\text{COP}_{\text{coldwater}}} + \frac{1}{Q_{\text{IT}}} * \frac{Q_{\text{air}}}{\text{COP}_{\text{air}}} \qquad (3.4)$$

Using equation 3.4, a data center and IT system specific PUE (sPUE equation 3.5) can be defined as:

$$\text{sPUE} = 1 + \text{Overhead}_{\text{PDCL}} + \sum_{k=1}^{n}(w_k * \frac{1}{\text{COP}_{\text{k}}}) \qquad (3.5)$$

$$with$$

$$w_k = \frac{Q_k}{Q_{IT}} \qquad (3.6)$$

$$1 = \sum_{k=1}^{n} w_k \qquad (3.7)$$

$$Overhead_k = \frac{w_k}{COP_k} \qquad (3.8)$$

The HPC system heat removed by different cooling technologies is represented by equation 3.6. Where $w_k$ is the fraction of heat removed via each heat removal technology $k$, and the sum of all $w_k$ is 1 (equaling 100% heat removal, 3.7). Overhead$_{PDCL}$ is the additional power needed by the data center to provide 1W of IT power. The data center overhead incurred by cooling the system for one specific cooling technology $k$ can be represented by Overhead$_k$ which is the power needed to remove 1W of heat via the heat removal technology $k$ (equation 3.8).

Finally, DWPE can be defined (equation 3.9).

$$\text{DWPE} = \frac{\text{WPE}}{\text{sPUE}} \qquad (3.9)$$

This metric is very useful for comparing the energy efficiency of different HPC systems and cooling solutions for running one particular workload.

**DWPE and EtS**

Calculating the EtS$_{datacenter}$ which includes the data center infrastructure overhead can be done using the same principle as the DWPE calculation. Opposite to the Performance per watt metric (DWPE), one can not use average PUE or average COP's since they are averaged over a year but can change substantially over time, with load, outside conditions, cooling temperatures, etc. Therefore, to use DWPE a data center needs to be able to measure or calculate the COP for the time interval of an application run. Then equation 3.10 can be used to calculate an application EtS$_{datacenter}$ over the whole data center.

Another difference between EtS$_{datacenter}$ and DWPE is that EtS$_{datacenter}$ will increase from less comprehensive measurements like EtS$_{system}$ or EtS$_{computeonly}$ whereas DWPE decreased from WPE. Which makes sense since energy consumption increases the more systems are included resulting in a decreased performance per watt. Equation 3.10 shows how EtS$_{datacenter}$ can be calculated assuming EtS$_{system}$ was measured.

$$\text{EtS}_{\text{datacenter}} = \text{EtS}_{\text{system}} * (1 + \text{Overhead}_{\text{PDCL}} + \sum_{k=1}^{n} \text{Overhead}_{\text{k}_{\text{average for time interval}}}) \quad (3.10)$$

PowerDAM can be used to calculate the COP's and, therefore, the different cooling surcharges for each monitoring timestep as shown by the author in [68].

Figure 3.15 shows the COP of LRZ for the chiller-less (also called warm water or hot water) and chiller-supported (cold water) cooling infrastructure for the months of May to July for 2016 (straight lines indicate times where no valid data was collected from the monitoring systems).



Figure 3.15: COP of the chiller-less (Warm Water) and chiller-supported (Cold Water) cooling infrastructure at LRZ (May to July 2016).

As can be seen, the COP of the chiller-less cooling infrastructure is more variable than the COP traditional chiller-supported cooling infrastructure. It fluctuates between 10 and 20. Also the chiller-less cooling infrastructure is, in the worst case, two times more efficient and, in the best case, 4 to 5 times more efficient than the chiller-supported cooling infrastructure.

### 3.3.6 DWPE Usage Example

Returning to the question related to table 3.1: Have SuperMUC Phase1 and PLX the same efficiency at LRZ?

The cooling distribution for the two systems is shown in table 3.3. The heat generated by SuperMUC is removed by the following cooling technologies: 10% is removed via air cooling, 18% via cold water cooling, and 72% via hot water cooling. PLX is 100% air cooled.

|                        | SuperMUC | PLX  |
|------------------------|----------|------|
| removed via air        | 10%      | 100% |
| removed via cold water | 18%      | 0%   |
| removed via hot water  | 72%      | 0%   |

Table 3.3: Cooling system distribution for SuperMUC and PLX

The 2013 cooling overheads for the LRZ data center are shown in table 3.4. The electrical distribution and conversions overhead is 7.5% (0.075). The overhead for removing 1W via W4 (Ashrae Water Categories [7] - chiller-less high temperature water cooling) is 5%, via W2 (chiller supported cold water cooling) is 40%, and via air is 50%.

|                           | LRZ   |
|---------------------------|-------|
| electrical overhead       | 0.075 |
| air cooling overhead      | 0.500 |
| cold water cooling overhead | 0.400 |
| hot water cooling overhead  | 0.050 |

Table 3.4: LRZ electrical and cooling overheads from 2013

Table 3.5 shows the DWPE results for SuperMUC Phase1 and PLX for the LRZ data center. Due to the lack of WPE data for PLX, the benchmark measurements were treated as WPE results.

DWPE for SuperMUC Phase1 was calculated as follows:

$$\text{sPUE}_{\text{SuperMUC Phase1}} = 1 + 0.075 + (0.1 * 0.5) + (0.18 * 0.4) + (0.72 * 0.05) = 1.233$$

$$\text{DWPE}_{\text{SuperMUC Phase1}} = \frac{908}{1.233} = 736$$

The results (Table 3.5) show that SuperMUC Phase1 is more efficient for running High Performance Linpack at LRZ than the PLX system. The "real world" efficiency difference would be 170 MFLOPS/W at LRZ whereas the Green500 values show only a difference of 16 MFLOPS/W.

To summarize, DWPE is a better real world energy efficiency indicator then WPE and the HPC system cooling technology has a strong impact on the "real world" efficiency of a HPC system.

|  | SuperMUC Phase1 | PLX |
|---|---|---|
| Green500 (MFLOPS/W) | 908 | 892 |
| DWPE (MFLOPS/W) | 736 | 566 |

Table 3.5: MFLOPS/W and DWPE for SuperMUC Phase1 and PLX for the LRZ data center

### 3.3.7 Data Center Energy Efficiency (DCEE)

When using DWPE, how can a data center determine which system would be best for its workload-mix? This question can be answered by defining a new metric called Data Center Energy Efficiency (DCEE) which combines multiple DWPE results. Currently, the main use for DCEE is for the procurement of a new system. By selecting a new system according to DCEE, a data center can show that it procured the most energy efficient system for its requirements at that time.

To be able to describe a workload-mix using multiple DWPEs which can include different instances of the Performance/W metric (such as FLOPS/W, MTEPS/W, or Iterations(It)/W, etc.) , one needs to find a way to remove the unit dependencies. For DCEE this is done by considering how the system under evaluation differs in it's Performance/W from the best value. The calculated performance difference for one DWPE is then weighted according to its importance in the data center workload-mix.

DCEE is defined as the sum of the weighted DWPE factors; where each factor is calculated by dividing the measured DWPE for each workload by the best DWPE for each workload. Because the best DWPE for each workload will change over time it is necessary to add a date to DCEE.

$$DCEE_{date} = \sum_{i=1}^{n} w_i * \frac{DWPE_i}{DWPE_{Nr.1}} \tag{3.11}$$

Where $w_i$ is the particular weight for a representative workload $DWPE_i$ which is part of the data center workload mix. The sum of all weights $w_i$ equals 1:

$$1 = \sum_{i=1}^{n} w_i \tag{3.12}$$

A DCEE of 1 would indicate that the data center is running the most energy efficient system for all workloads in its workload mix. A lower DCEE indicates a less efficient system for the data center and its work load mix.

Figure 3.16 depicts a generic representation of the DCEE evaluation process.

The process starts with an official HPC system ranking (right side of the picture). From the official ranking for a specific workload (Green500 for example), a new ranking is calculated using WPE adjustments, if required. For real world energy efficiency evaluation, all HPC Performance/W ranking lists should include, besides the compute hardware ranking, a WPE ranking. This WPE ranking of the systems is than combined with DWPE to generate a data center specific HPC system ranking for the specific workload. This is done with as many workload rankings as necessary to represent the expected data center workload-mix. Finally, DCEE is calculated for each HPC system.

Figure 3.16: DCEE generic evaluation process.

As a concrete example [1], the best supercomputer for the LRZ data center is determined considering the three different workload mixes as shown in Figure 3.17.



Figure 3.17: Example evaluation workload mix.

The following systems are considered:

- SuperMUC Phase1: an iDataPlex DX360M4, Xeon E5-2680 8C system and one of the PRACE Tier0 systems, installed at LRZ

- EURORA: a PRACE 2IP prototype, based on Eurotech Aurora HPC 10-20, Xeon E5-2687W 8C, with two NVIDIA K20 accelerators per node, installed at CINECA

---

[1] The shown examples are *not* an evaluation of the efficiency of the used applications. They are used to drive the point home that a system that is efficient for one workload is not necessarily the best for another application. This was also shown in the PRACE [69] 1st Implementation Phase Work Package 9 (1IP-WP9) deliverable 9.3.3 [58].

- Fermi: a BlueGene Q and one of the PRACE Tier0 systems, installed at CINECA

- PLX: an iDataPlex DX360M3, Xeon E5645 6C system, with one NVIDIA M2070 per node, installed at CINECA

First, one needs the cooling distribution for the HPC systems under consideration. The cooling distribution for the four systems are shown in table 3.6. The information for Fermi was taken from [70].

|                        | SuperMUC Phase1 | Eurora | Fermi | PLX  |
|------------------------|-----------------|--------|-------|------|
| removed via air        | 10%             | 0%     | 9%    | 100% |
| removed via cold water | 18%             | 0%     | 91%   | 0%   |
| removed via hot water  | 72%             | 100%   | 0%    | 0%   |

Table 3.6: Cooling system distribution for SuperMUC, Eurora, Fermi (BlueGene Q), and PLX

The cooling distribution depends strongly on the temperature difference between the machine room temperature and the cooling medium since the IT racks and internal cooling pipes are not insulated. For example, if the inlet temperature of the hot water cooling circuit is increased more heat will radiated into the air. This effect is even more visible if air is drawn over the HT-DLC components.



Figure 3.18: CoolMUC-2 heat transfer to direct liquid cooling.

Figure 3.18 shows, for different water inlet temperatures, how much of the electrical power consumed by the CooLMUC-2 system is transferred into the hot water cooling loop. Since the CooLMUC-2 node design has the power supplies behind the nodes, the power supply fans draw air over the node. At an inlet temperature of 40°C, 80% of the consumed power goes into the water. At 55 °C this drops below 60%.

Next, one needs the rankings for each of the workloads in the workload-mix. Table 3.7 shows the $\mathrm{DWPE}_{HPL_{2013-11}}$ as FLOPS/W values for each of the systems running HPL

at LRZ (the Green500 list results were assumed to be WPE results). By using the LRZ overhead information from table 3.4 the DWPE for each system is calculated. For example, DWPE for SuperMUC Phase1 is calculated using:

$$\text{DWPE} = \frac{\text{WPE}_{\text{HPL}_{2013-11}}}{1 + \text{Overhead}_{\text{electric}} + \text{Overhead}_{\text{air}} + \text{Overhead}_{\text{coldWater}} + \text{Overhead}_{\text{hotWater}}} \tag{3.13}$$

$$\text{DWPE}_{\text{HPL}_{2013-11}} = \frac{908}{1 + 0.075 + 0.10 * 0.500 + 0.18 * 0.400 + 0.72 * 0.050} = 736 \tag{3.14}$$

As can be seen, Eurora would be the most energy efficient system for running HPL at LRZ.

| | $\text{WPE}_{HPL_{2013-11}}$ | $\text{DWPE}_{LRZ_{HPL_{2013-11}}}$ | rank |
|---|---|---|---|
| SuperMUC Phase1 | 908 | 736 | 3 |
| Eurora | 3209 | 2852 | 1 |
| Fermi | 2176 | 1466 | 2 |
| PLX | 892 | 566 | 4 |

Table 3.7: DWPE for the Green500 workload for different system for the LRZ data center

Since there is no official ranking for Quantum Espresso, dataset:Ta2O5-2x1xz-552 with 20 iterations was used to benchmark each system:

1. on SuperMUC Phase1: 16 nodes, 64 MPI tasks, 4 OpenMP threads per MPI rank

2. on Eurora: 5 nodes, 10 K20 GPU, 10 MPI tasks, 8 OpenMP threads per MPI rank

3. on Fermi: 64 nodes, 256 MPI Task, 8 OpenMP threads per MPI rank

4. on PLX: 5 nodes, 10 M2070 GPU, 10 MPI tasks, 6 OpenMP threads per MPI rank

Table 3.8 shows the $\text{DWPE}_{QuantumEspresso_{2014-02}}$ as Iterations/W values for Quantum Espresso for the four systems. Here SuperMUC would be the most energy efficient system at LRZ.

| | $\text{WPE}_{QuantumEspresso_{2014-02}}$ | $\text{DWPE}_{LRZ_{QuantumEspresso_{2014-02}}}$ | rank |
|---|---|---|---|
| SuperMUC Phase1 | 35.3784E-06 | 28.6929E-06 | 1 |
| Eurora | 10.9583E-06 | 9.74074E-06 | 2 |
| Fermi | 3.11953E-06 | 2.10211E-06 | 4 |
| PLX | 6.55271E-06 | 4.16045E-06 | 3 |

Table 3.8: DWPE for the Quantum Espresso workload for different system for the LRZ data center

Now that we have the DWPE's for workloads in our workload-mix (Figure 3.17), we can calculate DCEE for the different workload distributions (equation 3.11).

For Workload Mix 1 running on SuperMUC Phase1 the following calculation needs to be solved:

$$DCEE_{LRZ.workloadmix1.SuperMUC} = 0,5*\frac{718MFLOPS/W}{2852MFLOPS/W}+0,5*\frac{27.9781E-06It/W}{27.9781E-06It/W} = 0,626$$

(3.15)

The DCEE for LRZ with different systems and the three different workload distributions are shown in table 3.9.

|  | Workload Mix 1 | Workload Mix 2 | Workload Mix 3 |
|---|---|---|---|
| SuperMUC Phase1 | 0.626 | 0.925 | 0.327 |
| Eurora | 0.674 | 0.413 | 0.935 |
| Fermi | 0.295 | 0.119 | 0.470 |
| PLX | 0.174 | 0.154 | 0.194 |

Table 3.9: DCEE for LRZ with different systems and different workloads

According to DCEE, for Workload Mix 1, Eurora would be the most energy efficient machine for LRZ. If Workload Mix 2 resembles the real distribution then SuperMUC Phase1 is the best choice. For a distribution like Workload Mix 3 the Eurora system would be the best choice again.

Since DCEE can combine different Performance/W metrics one should be conscious about the representative workload selection. If one has a significant higher value than others, the Performance/W selection could create a very strong bias towards that specific benchmark. Here the weight should be adjusted according to data center preferences.

### 3.3.8 Summary

This chapter discussed the term of Energy Efficiency and what it can mean for an HPC data center. It presented the PowerDAM tool which is used to collect power and energy data from *Pillar 1, 2, and 3* of the *4 Pillar Framework* providing the data required for a complete data center energy efficiency analysis. And finally, two new metrics DWPE and DCEE were introduced that can be used to calculate the energy efficiency from single applications or from a workload mix for a specific HPC system for a specific data center. DWPE is useful for application costs calculations whereas DCEE is more useful in a procurement context.

A white paper by CGG [71], using DWPE, showed that their free air cooled data center (PUE 1.05) consumed more power for their workload mix than their immersion cooling based data center (PUE 1.26).

# 4 Improving Energy Efficiency

> Discovery consists of seeing what
> everybody has seen and thinking what
> nobody has thought.
>
> _____
> Albert Szent-Györgyi

This chapter applies the work detailed in chapters 1, 2, and 3.

## 4.1 Power vs. Energy - Why power optimization is not energy optimization

Energy and power are used interchangeably when talking about efficiency improvement in HPC since power and energy are closely related. Energy is a product of consumed power over time. Even though Energy and Power are related, the techniques used and the efforts require to improve energy efficiency or to control power consumption are different. Currently there are two major research areas. One focuses on energy efficiency improvement of the HPC system and data center (European HPC research projects and research focus of LRZ) whereas the other tries to optimize already spent money by getting the power consumption of an HPC system/data center as close as possible to an already paid for power bound (US DOE (Department of Energy) funded projects, mainly investigated by Lawrence Livermore National Laboratory (LLNL) and Oak Ridge National Laboratory (ORNL)).

Power is a physical constraint in a sense that a data center is designed for a specific power load (for example, the LRZ infrastructure is designed for 10MW of power). Therefore, power impacts the operation of the data center. The main challenges are:

- power spikes - a short time excessive power consumption impacting the power delivery system and affecting power contract savings opportunities

- brown out - disturbances on one phase of the power delivery system impacting data center systems

- heat load - the consumed power directly impacts the heat load on the data center cooling infrastructure

- temporary disturbance in data center cooling capacity - if some part of the cooling infrastructure fails the power consumption of the IT systems might need to be adjusted if the available cooling capacity is insufficient

- high power slopes - very short high power changes due to application workload (power ramp-up and/or ramp-down)

- trapped capacity - power that is allocated to the HPC system but not used in normal operation [72]

- stranded capacity - power capacity for which there is insufficient distribution in the data center [72]

Power is mainly related to *Pillar 1* and *Pillar 2*; but it is also important for *Pillar 4* since it can be used to characterize applications in terms of min, max, and average power consumption.

For example, one technique that is actively advocated to address trapped capacity (especially important if a DC has a utility contract where power is paid) is called "over provisioning under a power bound" [73]. In order to use all available power more compute hardware is installed; resulting in a system where not all compute hardware can be run with maximum power draw [74].

Energy, on the other hand, is for most data centers a budget constraint. Depending on individual power contracts this could be a dynamic resource meaning any energy efficiency improvement will save money, or a fixed resource where any energy efficiency improvement would allow more IT hardware to operate. Since energy consumption can only increase over the year there is no energy peak but a finite amount available for the year. Energy can be used to associate costs to the run of an application via Energy-To-Solution (EtS) using DWPE. The main challenges are:

- a finite supply - since the monetary budget is fixed, there is a limit to the energy that can be used for a whole year but energy is a dynamic resource

- a fixed supply - if power was paid than energy becomes a fixed resource

- the best way to associate costs to applications - shared data center cooling infrastructure makes it non trivial

- the low frequency measuring interval of the data center power infrastructure

For example, at LRZ energy is a dynamic cost factor. LRZ pays for consumed energy over one year (8760h). Energy is measured using the average power over 15min intervals. The consumed energy is paid by the end of the year but needs to be budgeted 2 years in advance.

### 4.1.1 Utility Provider Contract

The power contract with the utility provider provides a major optimization constraint. Depending on the individual contract, the main optimization focus is either power usage or energy consumption, or a mix of both. LRZ's power contract is a typical example of an energy contract in Germany. The yearly payment consists of two parts: a connection fee (charged by the owner of the physical power lines) which is a fixed charge per kWh; and the energy consumption charge which depends on final (yearly) energy consumption (kWh).

Since the major cost factor is energy consumption (kWh), LRZ's main goal is to improve its data center overall energy efficiency. A secondary goal is to limit power peaks to 10%

of data center average power consumption in a billing cycle (15min at LRZ) because if the power consumption characteristic is good for the utility provider, 50% of the connection fee (around 266 k€ for 2015) will be refunded. To save 50% of the connection fee the yearly energy consumption needs to be more than the energy in the integral of the maximum power peak during the complete time interval (8760h in the year) over 80% of the time interval (80% of 8760h = 7000h).

To explain the connection fee saving opportunity, Figure 4.1 shows a 80h plot (operational data from February 2013, and for demonstration purpose equaling one year) without a power peak.
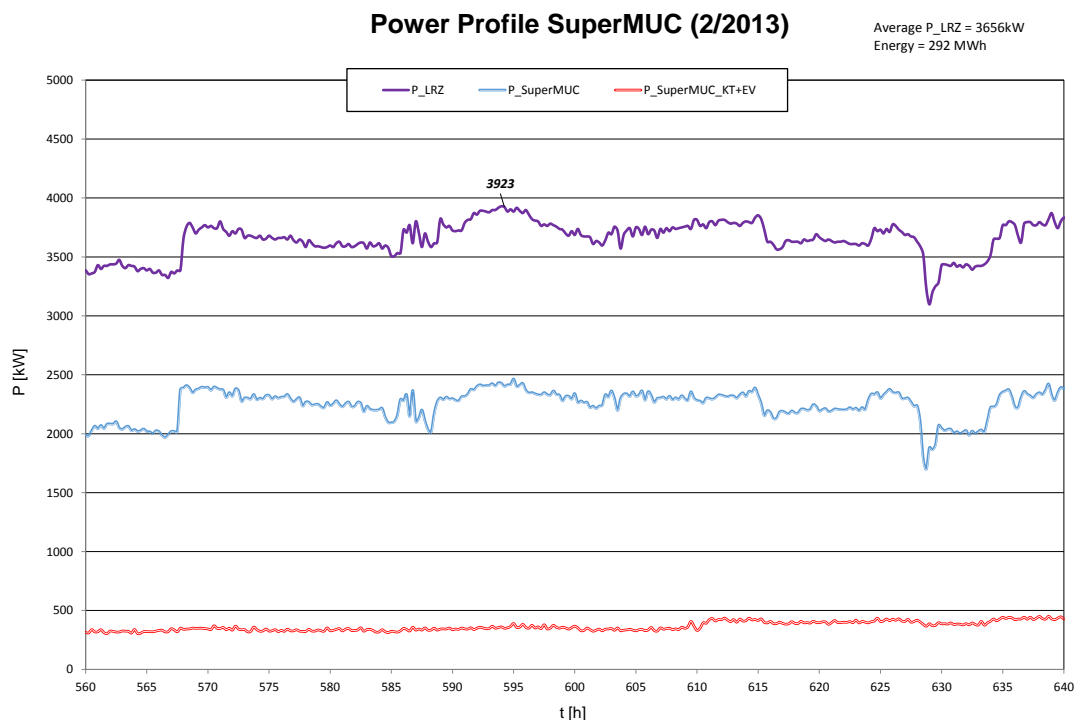


Figure 4.1: LRZ power profile for 80h in 2013 (without power peak).

50% of the connection fee can be saved if the energy consumption for that 80h time frame at least covers the area defined by a line through the highest power peak and ending at 80% of the measured timeframe (80% of 80h = 64h). Figure 4.2 shows the power profile with the energy area.

The average power consumption over the 80h time frame was 3656kW resulting in an energy consumption of 292MWh. The maximum power peak was 3923kW. The area inside the artificial energy box is: 3923kW * 64h = 251MWh. Since the overall energy consumption was higher than the energy described by the energy box, 50% of the connection fee would be refunded.

An additional example is shown in Figure 4.3. The same power profile is used but extended to 85h (representing a year for this example). It shows a power peak of 4250kW (added to the data as an example).

In this case, the area inside the box is: 4711W * 68h = 320MWh (Figure 4.4). The average
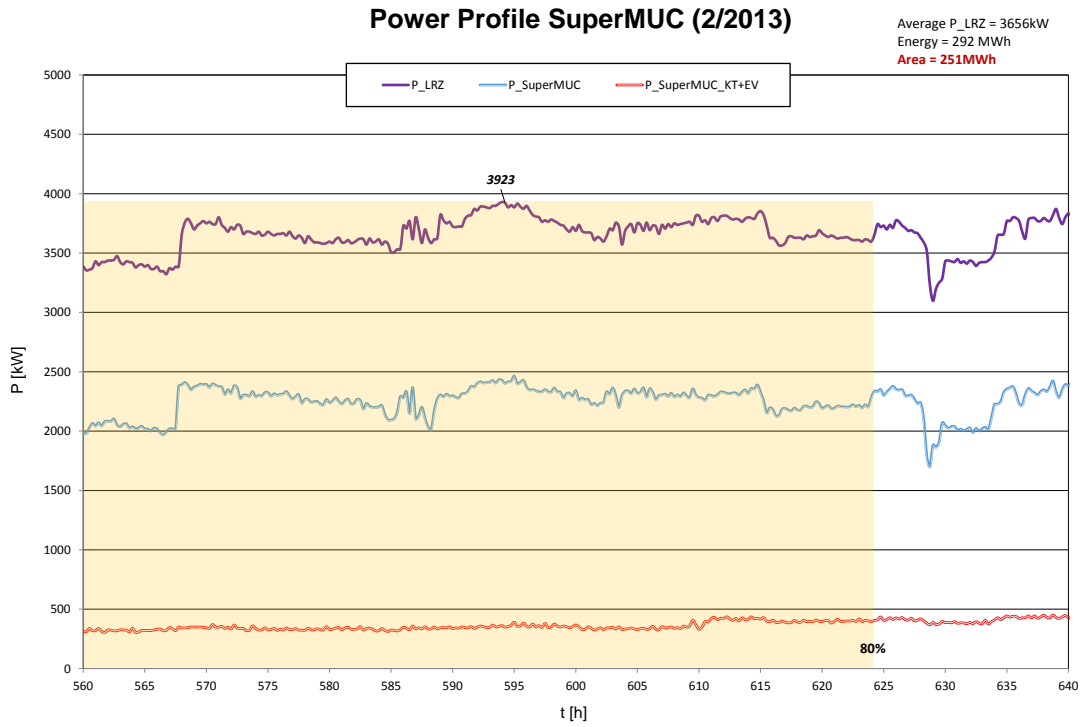
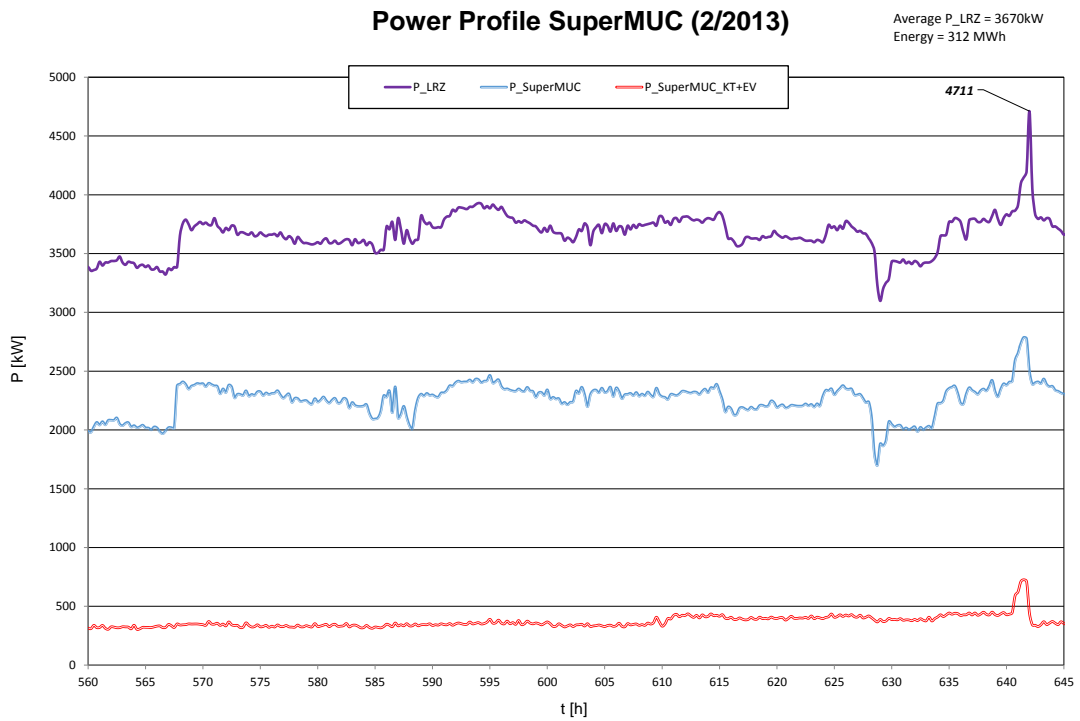Figure 4.2: LRZ power profile for 80h in 2013 (without power peak) with contract energy area.



Figure 4.3: LRZ power profile for 85h in 2013 (with artificial power peak).

power consumption of the time frame was 3670kW. This results in an energy consumption of 312MWh which is not larger than the area of the energy box. Therefore, no connection fee would be refunded.
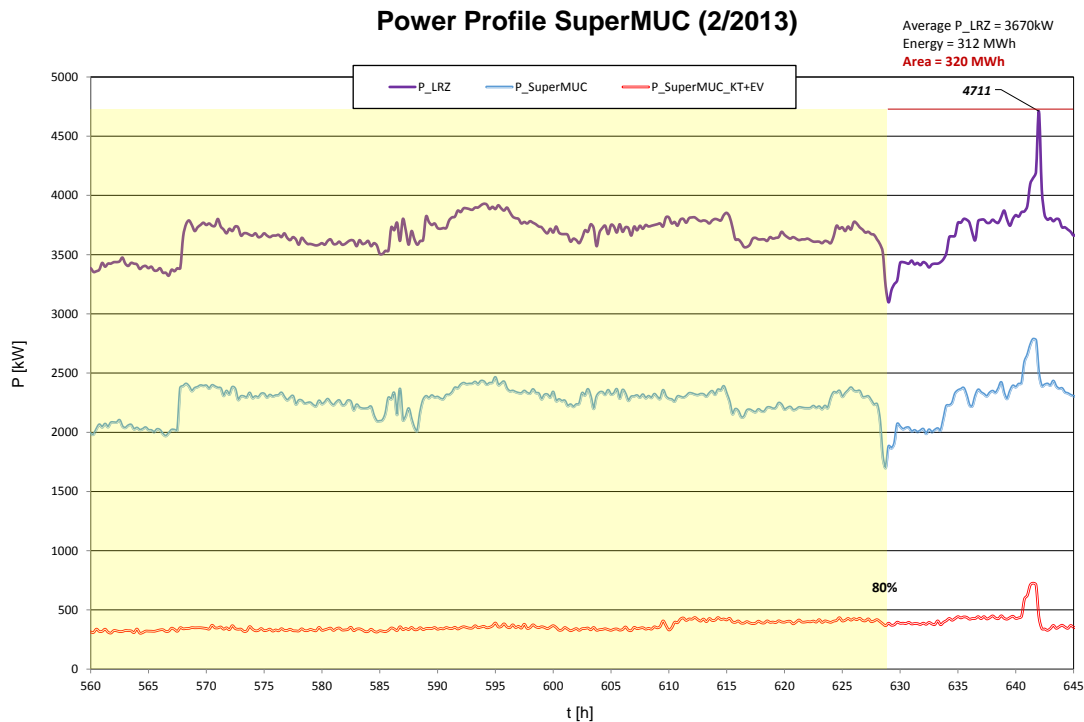


Figure 4.4: LRZ power profile for 85h in 2013 (with artificial power peak) with contract energy area.

In general it can be said that for the LRZ power contract a power peak of less than 25% of the yearly average power consumption would be allowed in order to get the 50% refund of the connection fee. For LRZ the energy billing interval is 15 min. Inside that interval the maximum allowed peak power would depend on the duration of the peak power consumption. This could be, for example, a power peak of 50% (assuming a constant average DC power draw) for 7.5 min of the billing interval if during the remaining 7.5 min only the average power is consumed. Since the average power of the 15min interval is used for billing purposes, a power peak of only 25% above average is recorded.

Saving on the connection fee is a possible optimization criteria if the HPC system is in stable operation where power peaks can be better controlled and avoided. At LRZ this is the case if the HPC system is older than 2 years. In the future the power consumption of the HPC system should be monitored in real time which would allow the development of automatic techniques which can limit the systems peak power consumption.

## 4.2 Identifying the most opportune pillars for a data center

The *4 Pillar Framework* can help data centers to identify areas that can provide the most benefit for overall energy efficiency improvements.

In general, savings in Pillar 4 (Applications) can be substantial since every Watt saved here will save more energy when looking at the complete data center [75]. Unfortunately, the effort required can be quite substantial and increases with the amount of applications a data center is running. For SuperMUC, this is estimated to be over 240 different applications. Also, with the replacement of the super computer every 5 years application codes might need to be optimized again. Since multi purpose, multi science data centers, such as LRZ, don't have direct control over the application developers, focusing on the other 3 Pillars might prove more fruitful.

Pillar 3 (HPC System Software) depends on Pillar 2. One area of interest is dynamic voltage and frequency scaling (DVFS) which allows the influencing of a systems power and energy consumption. A comprehensive overview of available energy efficiency techniques can be found in [76]. Another area seeing a lot of attention is the resource management system (also called scheduler) since it allows data centers to influence how applications are running. For example, LRZ uses Energy Aware Scheduling for SuperMUC [14].

Pillar 2 (HPC System Hardware) is under direct control of the data center since the data center writes the system request for proposal (RFP) document. Part of the requirements at LRZ is the required performance for key applications. The Energy Efficient HPC Working Group is working on RFP guidelines concerning energy efficient super computing [77].

Pillar 1 (Building Infrastructure) has a big impact on the data center energy efficiency since any additional overhead here will apply to all IT power consumption. Also, *Pillar 1* has the longest update cycle of any data center pillar. For example, the LRZ infrastructure is designed to last for at least 20 years. A switch to chiller-less direct liquid cooling reduces the cooling system power consumption overhead from 50% (indirect cold water mechanical chiller supported cooling) to 5%. Any savings that can be realized in Pillar 1 are complementary to any application optimization.

## 4.3 Saving energy by taking advantage of node power variability in homogenous HPC systems

This section discusses the existence of node power variation in homogenous HPC systems and its possible use to save energy. This chapter is partially based on the author's paper "Taking Advantage of Node Power Variation in Homogenous HPC Systems to Save Energy" [78] and journal paper "Analysis of the Efficiency Characteristics of the First High-Temperature Direct Liquid Cooled Petascale Supercomputer and Its Cooling Infrastructure" [54].

Node power variability was first quantified at LRZ for the CooLMUC system through the use of PowerDAM. CooLMUC (Figure 4.5) was built by MEGWARE and is the first AMD based high temperature direct liquid cooled (HT-DLC) HPC cluster (inlet temperature ASHRAE W4-W5) with 178 nodes (8 nodes interactive, 166 nodes batch, and 4 nodes reserved for internal use). A single node contains two AMD Opteron 6128HE CPUs (MagnyCours) with 8 cores each and 12MB L3 cache. In their standard setting, the CPUs run at 2GHz clock frequency. Each node is equipped with 16GB RAM arranged in eight 2GB DDR3 modules. The main interconnect network is InfiniBand QDR using a fat tree topology. In addition, each node has two Gbit Ethernet ports for IPMI and a service network which is used to boot the diskless nodes and to provide the root filesystem over NFS.
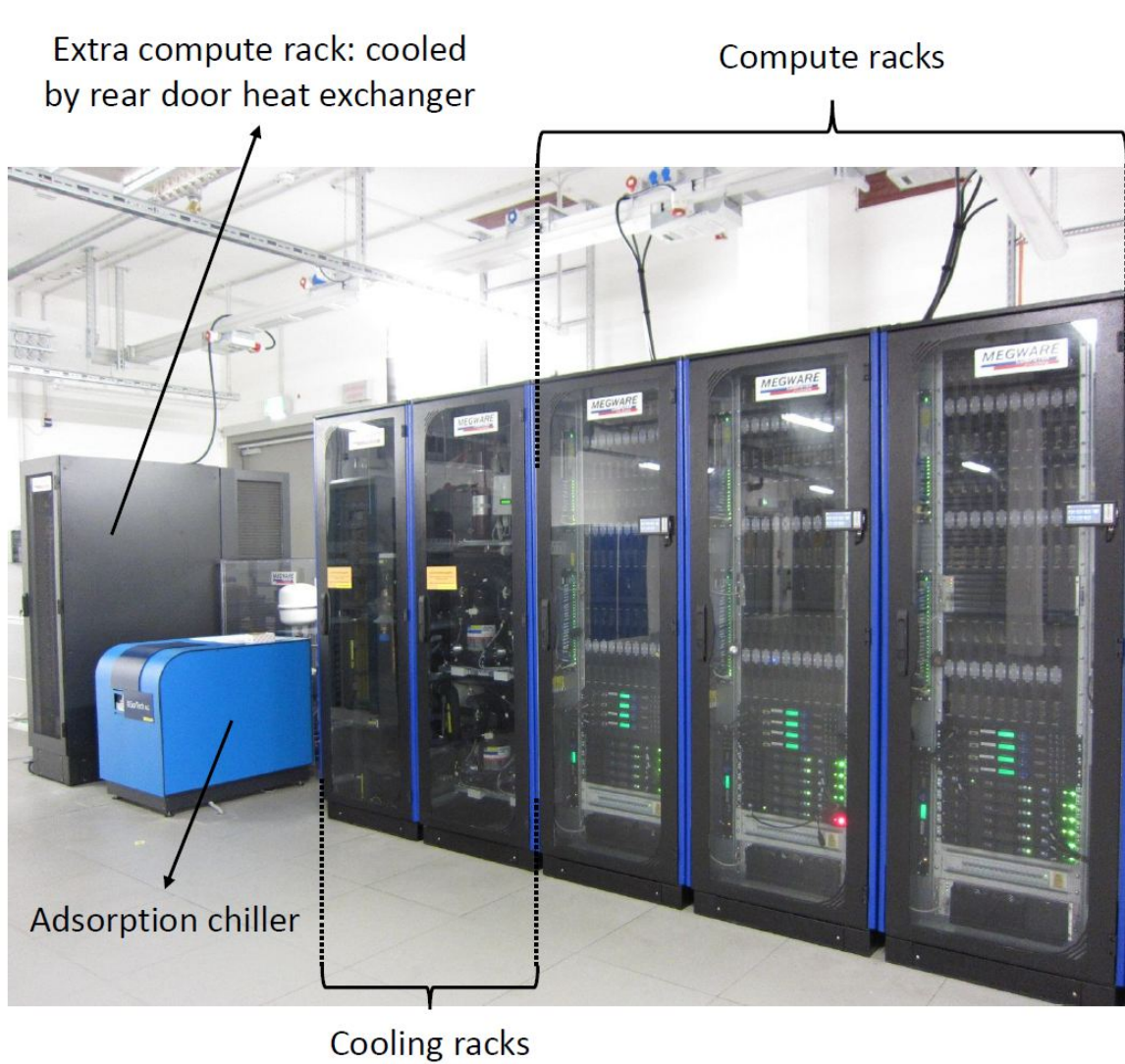
Figure 4.5: CooLMUC HPC cluster prototype at LRZ.

The cluster is completely room neutral; meaning that there is no requirement for computer room air conditioning (CRAC) units. Power measurements on CooLMUC are based on smart PDUs which report 1-minute average power values per node. Sufficiently long benchmark times (> 40min) are used to minimize the error of the 1-minute readouts.

Figure 4.6 shows the node power histogram for the CooLMUC system when running single node MPrime benchmarks on the complete system.

As can be seen, the power histogram for CooLMUC shows a non standard Gaussian distribution where 80% of the values are in an interval of one standard deviation $\sigma$. The average power consumption is 251W. The standard deviation is 4W (1.6% from average).



Figure 4.6: CooLMUC-1 node power histogram when running single node MPrime at 2.0GHz (AMD MagnyCore).

This discovery led to the following questions:

- Is this is a system hardware property that exists on all newer homogeneous HPC systems?

- If so, can it be used to save energy?

To answer these questions the SuperMUC system was analyzed. PowerDAM was used to collect measurements from the SuperMUC-Phase1 system. In addition, the node power variability for SuperMUC-Phase2 was measured during acceptance testing.

SuperMUC-Phase1 (Figure 4.7), which was Nr.4 on Top500 List (Jul 2012), was built by IBM based on iDataPlex technology with a peak performance of 3 PetaFLOPS. It is a Gauss Center for Supercomputing (GCS) system made available to PRACE users. SuperMUC's

Phase1 thin node islands have 147.456 processor cores in 9216 compute nodes. Each node has two Intel Sandy Bridge-EP Xeon E5-2680 8C processors, 32GB memory, and is direct liquid cooled using ASHRAE W4 water. The interconnect is Infiniband FDR10, a fat tree inside one island, and a Pruned Tree (4:1 blocking factor) between islands. SuperMUC Phase1 provides multiple levels of power measurements. For this analysis the "IBM Active Energy Manager" was used which collects power and energy consumption data at the power supply of each node.



Figure 4.7: The SuperMUC system[1]

SuperMUC-Phase2 (Figure 4.7), which was Nr.21 on Top500 List (Jun 2015), was built by IBM/Lenovo based on Lenovo NeXtScale nx360M5 WCT technology with a peak performance of 3 PetaFLOPS. It is a Gauss Center for Supercomputing (GCS) system made available to PRACE users. SuperMUC's Phase2 thin node islands have 86.016 processor cores in 3072 compute nodes. Each node has two Intel Haswell Xeon Processor E5-2697 v3 processors, 64GB memory, and is direct liquid cooled using ASHRAE W4 water. The interconnect is Infiniband FDR14, a fat tree inside an island, and a Pruned Tree (4:1 blocking factor) between islands. SuperMUC Phase2 provides multiple levels of power measurements. For this analysis the "IBM/Lenovo Active Energy Manager" was used which collects power and energy consumption data at the power supply of each node.

Figure 4.8 shows the node power histogram of island5 (512 nodes) of SuperMUC Phase1. 512 nodes are representative enough for the complete system characterization according to [79].

As can be seen, the distribution is also a non standard Gaussian distribution. The average node power consumption is 210W. 77% of the nodes are in an interval of one standard deviation $\sigma$. The standard deviation of 5.3W (2.5% from average) is slightly higher than

---

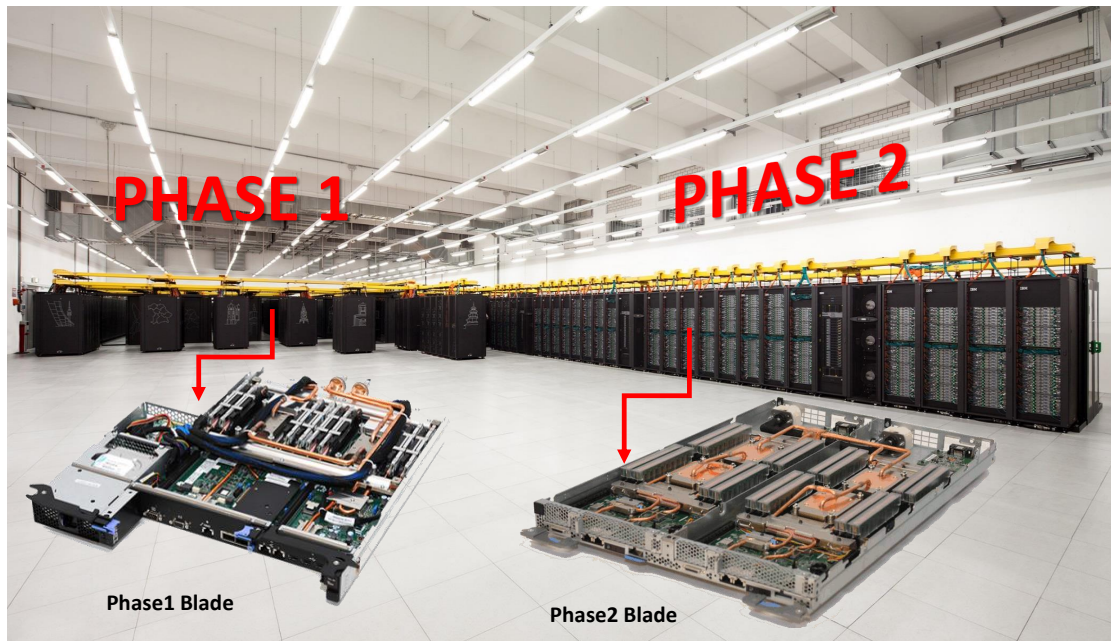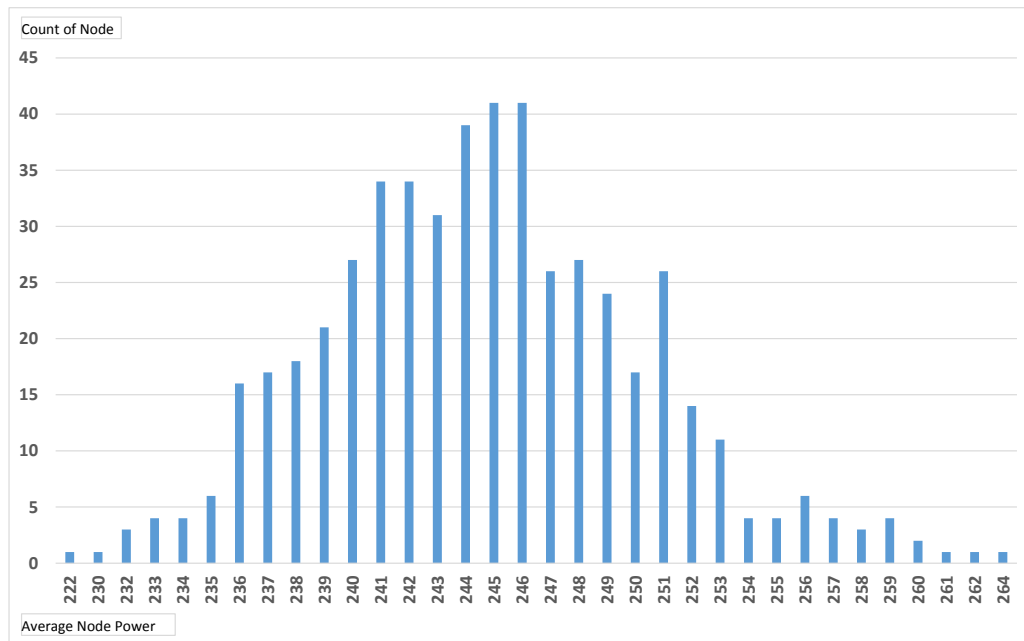[1]Blade pictures curtesy of Torsten Bloth, Lenovo, Germany

Figure 4.8: SuperMUC-Phase1 (island5 - 512 nodes) node power histogram running single node Firestarter at 2.3GHz (Turbo-Off, Intel Sandy Bridge).

the CooLMUC standard deviation.

Since in both distributions the major number of nodes are in the interval of one standard deviation $\sigma$, outlining nodes could indicate compounding issues in addition to manufacturing tolerances. For example, not well seated CPU heat sinks can cause higher CPU temperatures leading to higher node power consumption due to the increase of CPU leakage currents.

As a side note, the probability of receiving a good node when replacing a bad one is very high.

During the SuperMUC Phase2 node power characterization an interesting anomaly was detected [54]. Figure 4.9 shows a double peak in the distribution of DC (Direct Current) node power when running single node HPL (Turbo ON) on the whole system.

This double peak distribution was not seen by Hackenberg et al. [80] and Huang et al.[81] since both investigated only a single Haswell node. Also, Inadomi et al. [82] did not see it since Haswell was not part of the investigated system architectures. Running HPL with Turbo ON triggers Intel's power capping. According to Intel, all CPUs should be capped at 145 W; meaning that all nodes should consume approximately the same amount of power. But, as can be seen, this is not the case for SuperMUC Phase2. There are two distinct power distributions. Those could be traced back to the Motherboard Voltage Regulators [83] which are produced by two different manufactures. Since Intel uses, in the Haswell processor, for the first time, power and/or temperature for TDP capping, and the power measurements are provided by the Motherboard Voltage Regulator, the accuracy of the measurements provided by the external voltage regulator become important. For SuperMUC Phase2 the power reported by one manufacture is clearly lower than 145 W. There is an average difference of 32 W per node for the two motherboard types translating
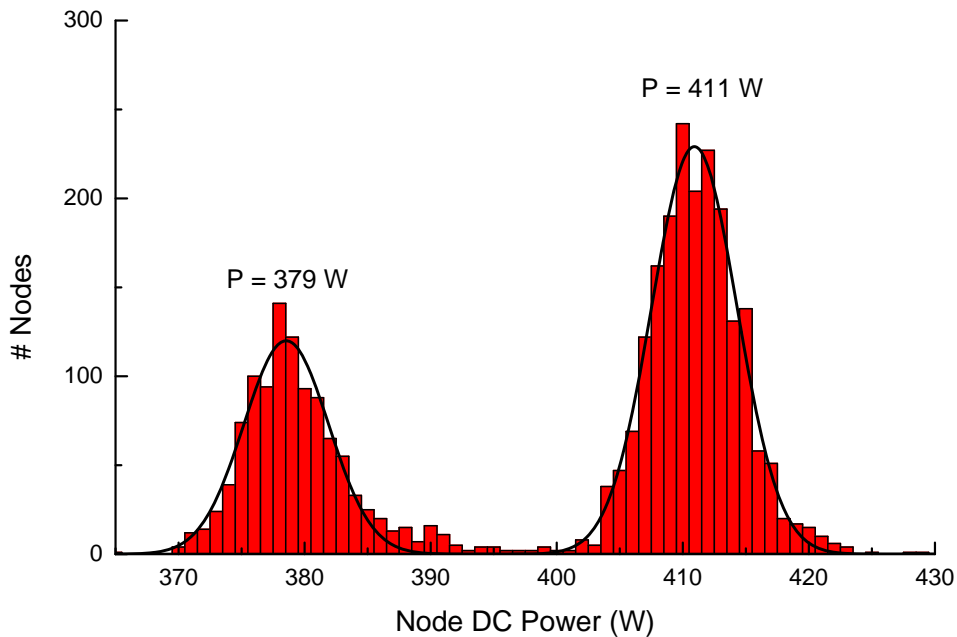
Figure 4.9: SuperMUC-Phase2 3072 nodes power histogram running single node HPL (Turbo-ON, Intel Haswell).

to a power difference of 16 W (slightly more than the 10% of TDP) per CPU.

The existence of node power variation has been shown as well by Hackenberg et al. [84]. Their analysis used the SPEC MPI benchmark to quantify power consumption variations of HPC systems. It showed a node power variation of 7% when idle and 5% under maximum load for 16 double nodes of an AMD Opteron cluster. Davis et al. [85] looked into variability of large-scale cluster power models. They stated that inter-node variations in power consumption is one reason that single node power models, when scaled to a large-cluster, show high errors. Inadomi et al. [82] analysed the impact of the power variability in a power constraint computing environment. Here the power variably translates into performance variability which could reach up to 64% performance variation across HPC application ranks. By implementing a variation-aware power budgeting framework, controlling the power budget for each compute node directly, the authors showed a substantial improvement in performance under a power bound since the direct control minimized the performance variation of the system nodes. However, no paper considered the use of node power variation for energy savings.

Our first proposed question, "Is this is a system hardware property that exists on all newer homogeneous HPC systems?" can then be answered with a *YES*.

Figure 4.10 shows how *Node Power Variation* fits into the *4 Pillar Framework* when looking at improving the energy efficiency of applications. Node Power Variation is a HPC system property and could potentially be used to improve Energy Aware Scheduling on SuperMUC [14]. Still, unless the application is fully optimize for the HPC system hardware application performance improvement will provide, in most cases, the greatest energy saving.

Since node power variability is a hardware property of the current and future HPC sys-
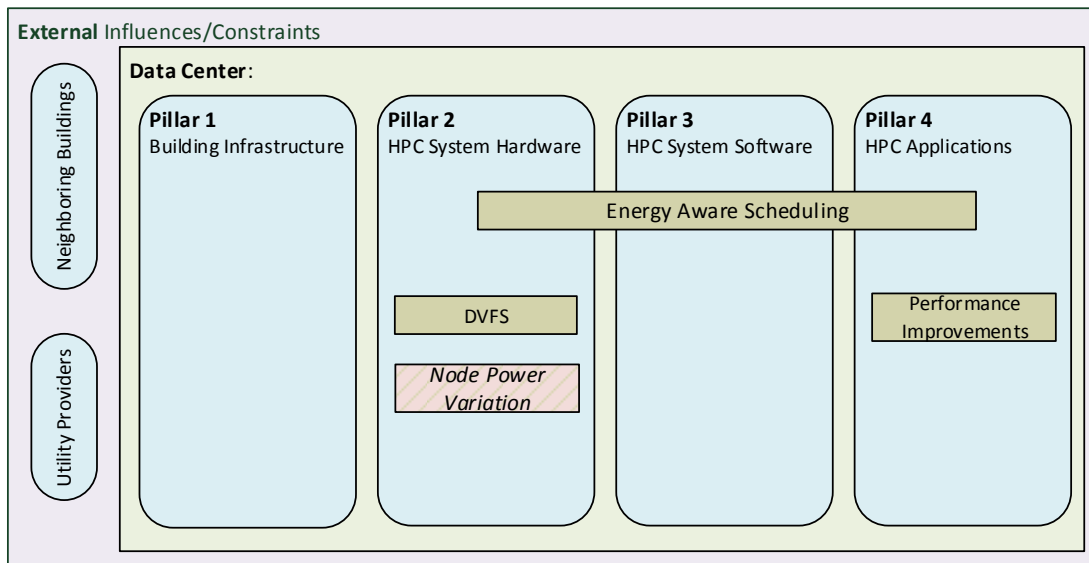
Figure 4.10: Node Power Variation in the 4Pillar Framework.

tem, node power variation lends itself to at least 3 possible energy savings opportunities.

1. Use this variation to schedule jobs on different nodes according to their power profile. For example, low power jobs on the more power consuming nodes and high power jobs on the less power consuming nodes.

2. Move the nodes with the higher power consumption into scheduling queues that show less usage.

3. Rank all nodes inside a system scheduling partition according to their power consumption and modify the scheduling system to use the higher ranked nodes before the lower ones.

### 4.3.1 Node power aware scheduling

Using hardware properties of the IT system to make better scheduling decisions in order to save energy is possible. Banerjee et al. [41] and Wang et al. [40] have shown that thermal and hot spot aware scheduling can save energy in distributed and cloud computing environments. Since node power is also a hardware property, similar scheduling techniques might be useful. Before investing into the development of these new scheduling techniques, however, the possible savings potential needs to be evaluated.

Table 4.1 shows the power and energy saving potential for running the High Performance Conjugate Gradient Benchmark (HPCG) [86] on the 10 worst vs. the 10 best nodes of island5 (node power distribution Figure 4.8). It can be seen that the application performance is not affected by the selection of nodes. Therefore, power and energy optimization can be used interchangeably. There is a difference of 26.9 W between running on the 10 worst nodes and running on the 10 best nodes. This represents an increase of 14.4% in power consumption when running on the worst nodes. This result seems promising.

|  | EtS (kWh) | average node power (W) | TtS (s) |
|---|---|---|---|
| 10 worst nodes | 0.658 | 213.4 | 1110 |
| 10 best nodes | 0.575 | 186.5 | 1110 |

Table 4.1: HPCG at 2.3GHz on 10 best and worst nodes of island5 of SuperMUC (Intel Sandy Bridge)

However, most HPC systems show an average system utilization of over 90%. Table 4.2 shows the energy savings for running HPCG on the the 256 best vs. the 256 worst compute nodes. Here a lower savings (4.3%) can be realized. This is due to the relatively narrow node power distribution where 77% of the nodes are in the interval of one standard deviation σ.

|  | EtS (kWh) | average power (W) | TtS (s) |
|---|---|---|---|
| 256 worst nodes | 15.55 | 197.7 | 1106 |
| 256 best nodes | 15.03 | 189.6 | 1115 |

Table 4.2: HPCG at 2.3GHz on 256 best and worst nodes of island5 of SuperMUC (Intel Sandy Bridge)

Another implication of the high node utilization in HPC is that even if one application can run on the most power efficient nodes another will be running on the worst. Therefore, any possible gain depends on the difference of power consumption of the running applications in both configurations. Assuming two applications $A$ and $B$, the possible energy savings depends on the difference between two configurations: running $A$ on the worst and $B$ on the best nodes; or running $A$ on the best and $B$ on the worst nodes. Figure 4.11 shows this graphically.
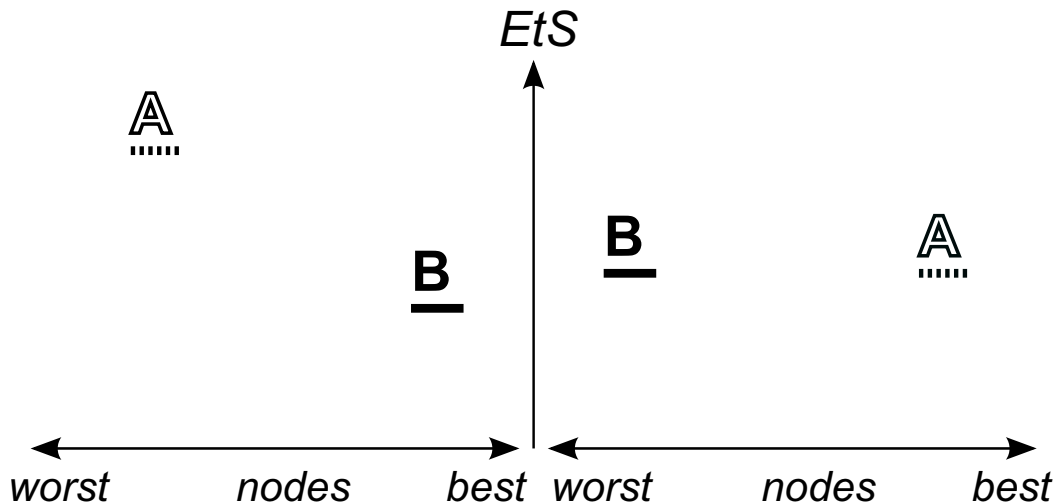


Figure 4.11: Node power aware scheduling energy savings.

The savings potential can be defined via the following formula:

$$\text{Saving}_{\text{max.theoretical}} = |\text{Gain}_{\text{A}} - \text{Gain}_{\text{B}}| = |(P_{\text{A.on.Worst}} - P_{\text{A.on.Best}}) - (P_{\text{B.on.Worst}} - P_{\text{B.on.Best}})|$$

By using this formula and, for example, the measurement data for HPCG (Table 4.2) as application $A$ and Epoch [87] (Table 4.3) as application $B$, an expected savings of 1.7W per node, equaling 870.4 W for the 512 nodes of island5, can be calculated. Both applications reflect real HPC data center workloads.

| | EtS (kWh) | average power (W) | TtS (s) |
|---|---|---|---|
| 256 worst nodes | 9.20 | 164.2 | 788 |
| 256 best nodes | 8.83 | 157.8 | 787 |

Table 4.3: Epoch at 2.3GHz on 256 best and worst nodes of island5 of SuperMUC (Intel Sandy Bridge)

The conclusions for node power aware scheduling that can be drawn from the presented examples are as following:

- To co-schedule applications to save energy, one application needs to be less sensitive to the used node set than the other.

- A combination of CPU bound applications (high CPU frequency and high power consumption) running on the best nodes and memory/io bound applications (low CPU frequency and low power consumption) running on the worst might show a slightly better savings potential.

- By using the best case savings from the examples, the expected savings is overestimated when compared to reality where the node allocation would be random.

In summary, the possible savings of node power aware scheduling for HPC systems is currently not substantial enough to account for the efforts required to enable the scheduling system to allow for this selection and to increase the scheduling complexity by another dimension.

*Nevertheless*, node power aware scheduling can provide *substantial benefits* for *distributed* and *cloud computing* where the average utilization rate for the system is between 10% to 50% [88], [89], [90], [91].

### 4.3.2 Node switch-off, node power aware system partitioning, and node ranking based on power variation

Node switch-off is one of the techniques that has the potential to save energy [92], [93]. Reducing the usage of the nodes with the highest power consumption by switching those off first could improve the savings potential. Unfortunately, with large scale HPC systems, node switch-off has some pitfalls. Since most large scale HPC systems use a disk-less node design, the nodes need to be re-initialized again after each switch off. Besides delaying the availability of the node, the re-initialization has been shown to lead to software stack

problems (driver initialization) which does not guarantee a 100% successful re-boot. Additionally, the InfiniBand auto-routing feature which is triggered after each node switch-off and re-boot has been shown to lead to cluster availability issues. SuperMUC Phase1 was the first very large scale InfiniBand FDR10 installation. The first tryout of the node switch-off feature disabled the complete HPC system since the route recalculation took to long and triggered network timeouts resulting in crashed applications. Therefore, node switch-off is not used on SuperMUC Phase1. The route re-calculation problem was later fixed but re-boot software stack issues still exist.

Node switch-off is not the only potential usage of node power variability. Figure 4.12 shows the runtime for every node in the two scheduling queues for the CooLMUC HPC cluster system for 2014. To calculate the node runtime, the individual job information (start and end time, node list) was used. For each job where the node was part of the node list the job runtime was calculated and summed up. The Batch queue (top line) has the highest node usage (around 300 days out of the year) whereas the Interactive queue (bottom line) shows the least usage (less than 100 days during the year). Additionally, in both queues not all nodes are used for the same amount of time over the year. Table 4.4 shows the average, minimum, and maximum utilization for the nodes of each partition for the year 2014.
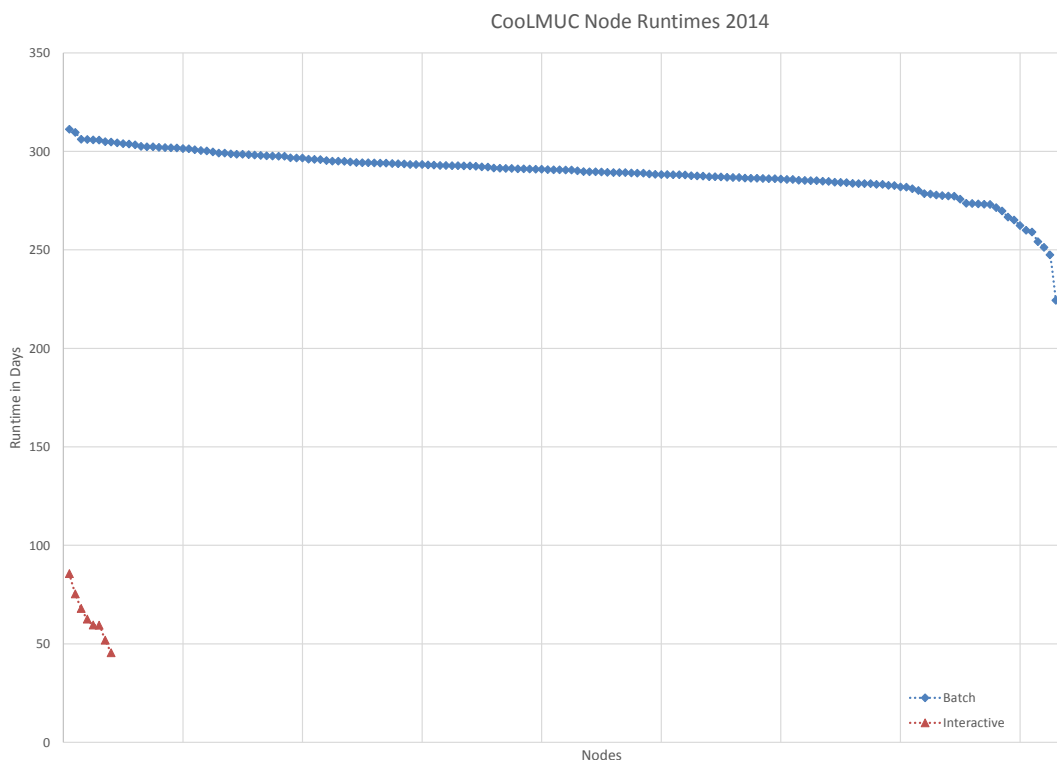


Figure 4.12: LRZ CooLMUC-1 scheduling queues node runtime for 2014.

For example, the node with the highest runtime in the Batch queue ran for 311.19 days whereas the one with the least usage ran for only 224.43 days.

Two important observations can be made from this data:

| runtime per node | interactive | batch |
|:---:|:---:|:---:|
| min | 45.49 days | 224.43 days |
| max | 85.62 days | 311.19 days |
| average | 63.47 days | 289.07 days |

Table 4.4: CooLMUC partition usage for 2014

1. A significant difference between partitions, depending on their usage pattern, exists. For example, the interactive partition is mainly used during normal working hours (8h a day equals 83 days/year per node) whereas the batch partition can be used 24/7 with a maximum of 365 days per node.

2. The utilization varies between the nodes even in the same partition.

By using the observed system properties, it is possible to define at least three techniques that use the average node power distribution to save energy. Firstly, the worst nodes should be moved to the partition with the least usage (ideal would be spare nodes). Secondly, the scheduler should prioritize the nodes according to the power consumption for each partition. Preference should always be given to the nodes with the lowest power consumption. And finally, because nodes are never utilized 100% in a year one should switch-off nodes especially for partitions that do not have a 24/7 usage pattern (if possible).

### 4.3.3 Large Scale HPC System Scheduling Queues Analysis

One important research question is: Are the CooLMUC scheduling queues representative of large scale HPC systems?

This section will analyse the SuperMUC-Phase1 scheduling queues and will draw conclusions related to the usefulness of node power aware system partitioning and node ranking based on power variation for large scale HPC systems.

SuperMUC-Phase1 has 8 official queues. From those, 4 are available for all users while others are only available after special permission. Figure 4.13 (source LRZ website [94]) shows the queue definitions of the 4 public queues:

The *test* and *micro* queue are limited to jobs with node requirements of 1 to 32 nodes. *test* is runtime limited to 30min whereas all other queues have a wall clock limit of 2 days (48h). *general#* is restricted to jobs requiring between 33 and 512 nodes. Those are scheduled on one island. *large#* is reserved for large scale jobs needing more than 1 island but less than 4 islands. The full machine can only be used during the special scaling workshops hosted at LRZ.

There are 4 other job queues. *tmp1*, *special*, and *preempt* are internal test queues. *tmp2* is a one island queue that can be reserved by customer request and is not shared during the runtime of the reservation.

Figure 4.14 shows the runtime of all the nodes associated at one point in time with the *test* queue sorted by most used nodes to least used ones.

The orange triangles show the complete node runtime for 2014. The blue circles show the runtime of the node in this particular queue. This is important since nodes can be part

| Class Name | Purpose | Max. Jobsize in Islands§ | min.-.max. Nodes | Wall Clock Limit | Run limit per user |
|---|---|---|---|---|---|
| | | | | | |
| test | Test and interactive use | 1 | 1 - 32 | 30 min | 1 |
| micro | Small jobs, pre- and postprocessing runs (internally restricted to run only on some specific islands) | 1 | 1 - 32 | 48 h | ~ 8 |
| general# | Medium-sized production runs fitting into a single island | 1§ | 33 - 512 | 48 h | ~ 8 |
| large# | Large Jobs, spanning more than one Island. | 4§ | 513 - 2048 | 48 h | ~ 8 |

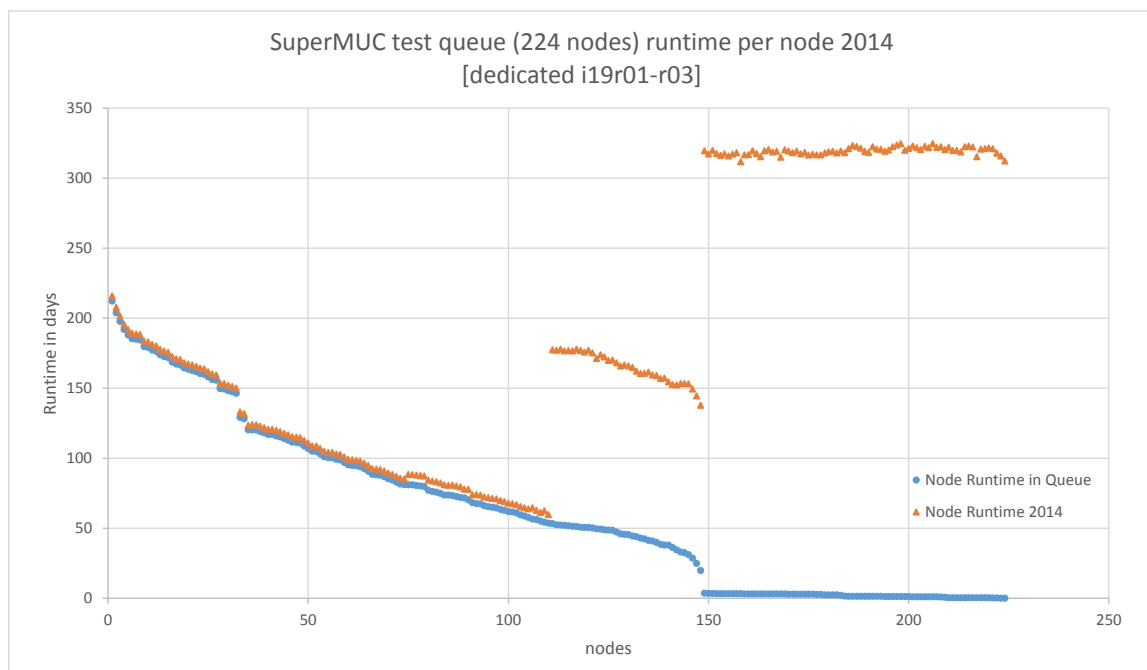Figure 4.13: LRZ SuperMUC Phase1 scheduling queues definition.



Figure 4.14: LRZ SuperMUC Phase1 Nodes in *test* queue runtime for 2014.

of multiple queues. What can be seen is that for the nodes where the blue and orange line match up the runtime shows a significant variation (50 days to over 200 days). Some nodes were part of another queue as well (between 100 and 150) but still show a distinct runtime variation. The nodes between 150 and 224 were part of the *test* queue for a very short time. For those, the runtime does not vary much.

Figure 4.15 shows the runtime of nodes in the *micro* queue. Besides a small number of nodes, all nodes show a high runtime (over 300 days). Still, the line drops from 330 days to 300 days.
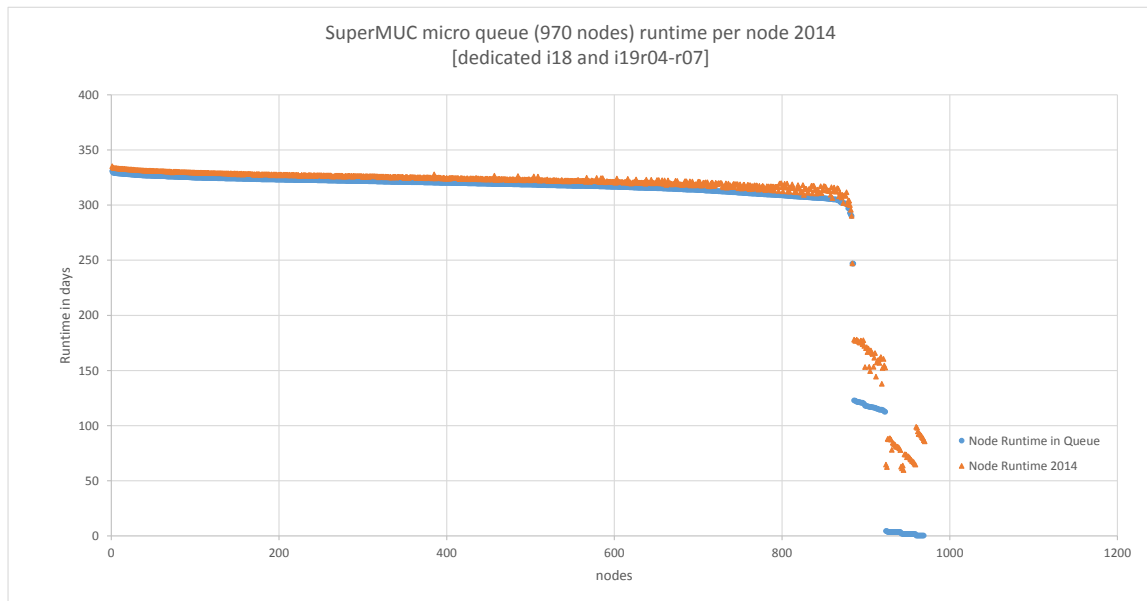


Figure 4.15: LRZ SuperMUC Phase1 Nodes in *micro* Queue runtime for 2014.

The next Figure 4.16 shows the runtime for nodes associated with the *general*<sup>#</sup> queue. This is the queue where jobs with less than 33 nodes are not allowed. Here the main bulk of the nodes run from 250 to 180 days. There are a small number of nodes with less runtime, but for most nodes the complete runtime is in a band between 280 and 310 days.

Figure 4.17 shows the runtime for nodes associated with the *large*<sup>#</sup> queue. Since only jobs requiring more than one island are allowed, the node runtime varies from above 100 days to less than 10 days. As can be seen, all nodes were part of other queues as well and the complete node runtime for the majority of nodes is in a band between 310 and 270 days. Nevertheless, the complete yearly usage drops to 100 days for some nodes.

Figure 4.18 shows the node runtime for the *tmp2* queue. This queue is now dedicated to island17 (512nodes) but other islands were used for a very short time as well. The queue is reserved on user request for whole island runs. Node usage varies from 275 to 250 days. A small number of nodes were used less. An interesting observation is that because of the first numbered node first allocation policy the usage is highest for the #1 node (275 days) and the least for the #512 node (150 days). The spare nodes (#513-#516) show less than 100 days of usage.

Figure 4.19 shows the complete SuperMUC-Phase1 node runtime sorted per island and node. Each of the 18 islands (starting with island-2 and ending with island-19) is clearly
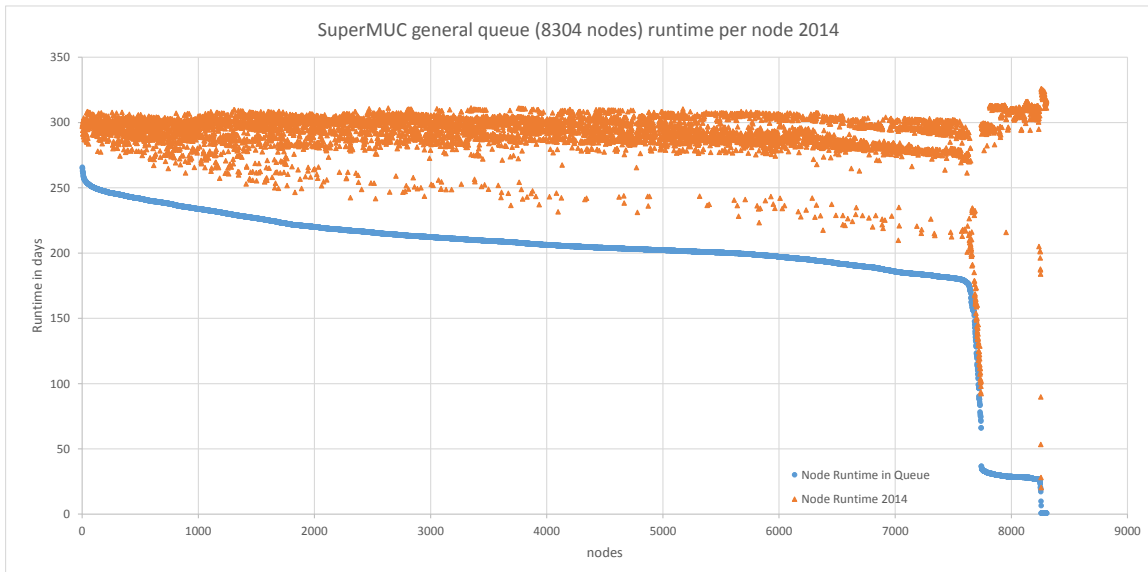
Figure 4.16: LRZ SuperMUC Phase1 Nodes in *general*[#] Queue runtime for 2014.
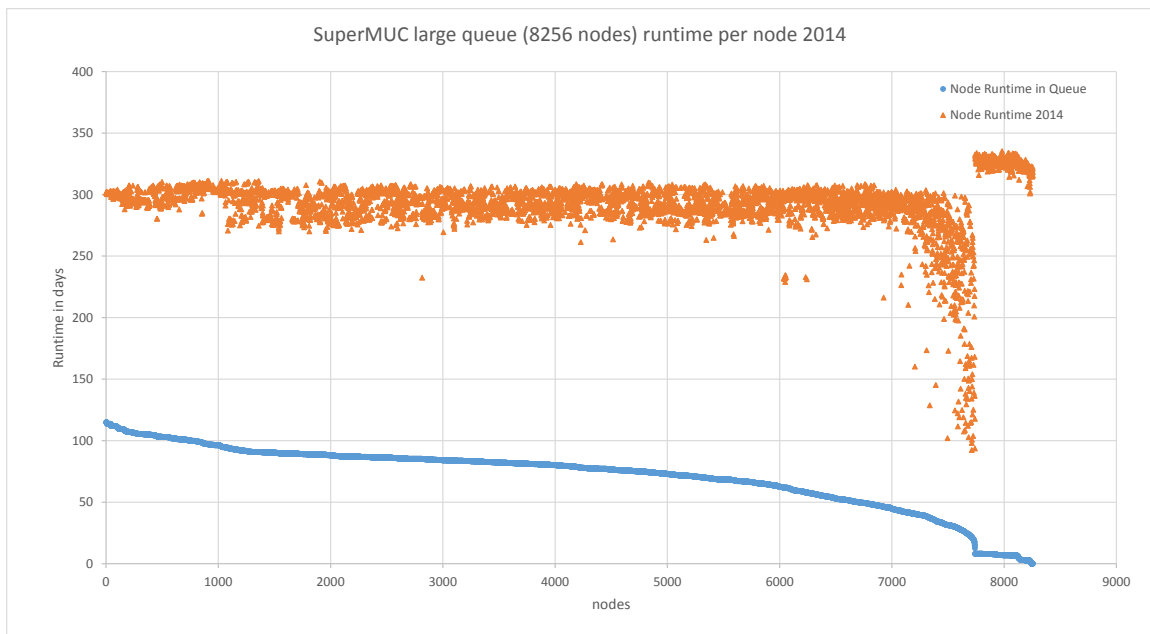


Figure 4.17: LRZ SuperMUC Phase1 Nodes in *large*[#] Queue runtime for 2014.

Figure 4.18: LRZ SuperMUC Phase1 Nodes in *tmp2* Queue runtime for 2014.

visible since the spare nodes in each island show the least amount of usage.



Figure 4.19: LRZ SuperMUC Phase1 Nodes runtime sorted per island for 2014.

The following observations can be made from the figure:

1. There is a visible node runtime variation.

2. The last nodes in each island (#500-#516) are used much less.

3. Some islands show a visible higher runtime for nodes than others.

4. The nodes at the beginning of the last island (#19) are used less than all other active nodes.

5. Two islands (#8 and #15) show a two line usage pattern.

6. The first nodes of most islands are not the most used nodes.

7. The 4 spare nodes in each island are clearly visible with the lowest runtime during the year.

To explain the observations, Figure 4.21 shows the same information as Figure 4.19 with the addition of the runtime for the different scheduling queues for each node.



Figure 4.20: LRZ SuperMUC Phase1 Node runtime per queue sorted per island for 2014.

The "first-node-first" scheduling policy can be clearly seen for the *large*# queue (third line from top). The usage for the first two thirds of the nodes in each island is relatively similar. After that there is a definit drop in runtime.

This drop-off can be exploited using back filing of small job sizes (*general*#). Here the runtime is opposite to the *large*# queue: the last nodes of each island are used most. Both

main queues combined show a relatively constant runtime for most nodes inside an island. Still, the last nodes of each island show a drop in yearly runtime. Even though the nodes of island #17 are mainly reserved for whole island runs they still show distinct runtime differences indicating that some of the applications scheduled do not use all available 512 nodes (black dots - second line for island #17). There is some backfilling with the *general*# queue but much less than for island #2-#16.

The *micro* queue (yellow - island #18 and #19) achieves the highest overall utilization of nodes since single node jobs are the norm.

The node in the *test* queue (green - beginning nodes of island #19) show the least runtime during the year since they are exclusively reserved for LRZ personnel to do software development or software testing. Some nodes were removed from the queue during the year and added to the *micro* queue.

The other queues are not relevant since they were used for a very short time on special occasions (training workshop, LRZ large scaling workshop, system testing after maintenance, etc).

The following additional observations can be made from Figure 4.21:

1. The node runtime variation comes from the "first-node-first" scheduling policy and the usage pattern of different queues combined with backfilling.

2. The last nodes of each island are used less because of the "first-node-first" scheduling policy and the usage pattern of different queues combined with backfilling.

3. Since nodes are shared between queues the combined runtime can vary for the nodes in each island.

4. Since the first 128 nodes of island #19 are exclusively reserved for the LRZ internal *test* queue they are utilized less. Also the nodes are used mainly during working hours.

5. For two islands (#8 and #15) the scheduling for the *general*# queue seems to be more scattered than for the other islands. Together with backfilling, this generates two visible lines for the complete runtime plot.

6. Interestingly, the first nodes of the islands are not always the most used nodes in the *large*# queue (third line from top) even though the queue requires jobs sizes of more than 1 island. But since the job size doesn't have to be an exact multiple of islands it could lead, in combination with backfilling, to higher usage of middle nodes.

7. The 4 spare nodes of each island are rarely used since they will be taken back out of any production queue if a failed node was replaced.

Figure 4.21 shows the node runtime for all SuperMUC-Phase1 nodes for 2014 sorted by runtime. The runtime of nodes dedicated to the LRZ *test* queue is the lowest compared to the nodes dedicated to the *temp2* (whole island run) queue and the node shared between all other *batch* queues.

Overall, the node utilization of SuperMUC-Phase1 is high with an average node runtime of 292.5 days. With the scheduled maintenance and outages of 23 days during 2014, this

Figure 4.21: LRZ SuperMUC Phase1 Node runtime for 2014.

leads to an overall node runtime of 315.5 days during the year. This corresponds to an average node usage of 86.5%.

The scheduling queue analysis of SuperMUC-Phase1 showed that:

- The observed CooLMUC node runtime variation (Figure 4.12) is not unique. The large scale HPC system SuperMUC-Phase1 shows a similar behavior.

- For nodes in multiple queues, the runtime in each queue can not be looked at individually. The overall runtime is important and might look much different from the single queue usage pattern.

"Node Power Aware System Partitioning" can be used for large scale HPC systems. Having multiple queues with intermixed nodes is beneficial for node usage but reduces the possible benefits of "Node Power Aware System Partitioning". The worst nodes should be moved into a dedicated partition with low yearly runtime (here the LRZ internal partition *test* and all SuperMUC-Phase1 spare nodes).

"Node Ranking Based on Power Variation" can be useful since there is a very distinct runtime difference of nodes (Figure 4.21). This technique can be used regardless of scheduling policies.

The high average node runtime of close to 90% limits the savings potential of the two techniques in HPC but the impact will increase with lower system utilization rates (standard data centers).

### 4.3.4 Savings analysis[2]

This section will quantify possible savings related to node power aware partitioning and node power ranking using the CooLMUC HPC system (174 nodes [8 nodes interactive, 166 nodes batch] are considered). First, the theoretically best and worst system distribution scenarios are compared to get an idea about the maximum possible savings potential. In reality, a system without node power aware partitioning and node ranking based on power variation will have a distribution between the best and worst case scenarios. Therefore, the real-world savings using the actual CooLMUC setup are discussed. PowerDAM [50] was used to collect and analyze power and energy data.

At first, the measured node power distribution (Figure 4.6) was normalized using the average power consumption of all nodes during the MPrime benchmark (251.14W). Given the power consumption of an application on an average node, the normalized power distribution can be used to derive the power consumption of the application when running on good nodes or bad nodes respectively. After that the following system and node statistics were collected:

1. System:
   - average power per node when running jobs on interactive partition
   - runtime of all jobs on the interactive partition
   - average power per node when running job on batch partition
   - runtime of all jobs on the batch partition

2. Nodes:
   - runtime for each node (annual utilization time of each node)
   - average power for each node (average of all power measurements when the node was used by a job during the year)

Table 4.5: CooLMUC partitions statistic for 2014

|                          | interactive | batch      | system    |
|--------------------------|-------------|------------|-----------|
| average node power (W)   | 156.20      | 217.58     |           |
| runtime (h)              | 12186.22    | 1151654.97 |           |
| energy consumption (kWh) | 1903.49     | 250577.09  | 252480.58 |

Table 4.5 shows the CooLMUC system partitions power and runtime statistics for 2014. The system energy consumption is the sum of the energy consumption of the interactive and batch partitions. The energy consumption of one partition is the *average yearly node power when running jobs on partition X* multiplied by *partition X runtime*.

Table 4.6 shows the data for node power aware system partitioning without taking node ranking into consideration. In the worst case, the best nodes are moved into the interactive partition. In the best case, the worst nodes are moved into the interactive partition. Using this technique a maximum theoretical savings of 663.28 kWh per year can be achieved.

---

[2]This chapter is nearly completely quoted from the author's paper "Taking Advantage of Node Power Variation in Homogenous HPC Systems to Save Energy" [78]

Table 4.6: Best possible savings for CooLMUC using node power aware system partitioning for 2014

|  | interactive | batch | system |
|---|---|---|---|
| average node power (worst in interactive) | 162.26 W | 217.17 W |  |
| energy consumption (worst in interactive) | 1977.39 kWh | 250108.26 kWh | 252085.65 kWh |
| average node power (best in interactive) | 152.08 W | 217.86 W |  |
| energy consumption (best in interactive) | 1853.27 kWh | 250895.65 kWh | 252748.93 kWh |
| possible max. savings |  |  | 663.28 kWh |

Table 4.7 shows the results for power aware system partitioning and node ranking based on power variation. In the worst case the best nodes are moved into the interactive partition and in each partition the nodes are arranged so that the longer the runtime the worse its power consumption. In the best case the worst nodes are moved into the interactive partition and in each partition the nodes are arranged so that the longer the runtime the better its power consumption.

Table 4.7: Best possible savings for CooLMUC using node power aware system partitioning and node ranking based on power variation for 2014

|  | interactive | batch | system |
|---|---|---|---|
| energy consumption (worst in interactive) | 1962.81 kWh | 250059.26 kWh | 252022.07 kWh |
| energy consumption (best in interactive) | 1841.05 kWh | 251112.44 kWh | 252953.49 kWh |
| possible max. savings |  |  | 931.42 kWh |

Using node ranking based on power variation in each partition saves nearly 50% more energy than power aware system partitioning alone; increasing the possible theoretical energy savings to 931.42 kWh.

In reality, CooLMUC consumed 252480.58 kWh in 2014 (table 4.5 - system energy consumption). This leads to a possible savings of 251.77 kWh for 2014 (252480.58 kWh minus 252022.07 kWh (taken from table 4.7 - energy consumption (worst in interactive))).

Since each 1kWh IT power saved will also save cooling costs, the HPC system internal

cooling overhead (1.23 for CooLMUC [95]) and sPUE (as defined in equation 3.5) need to be considered.

Table 4.8: LRZ electrical and cooling overhead

|  | LRZ |
| --- | --- |
| PDCL overhead | 0.075 |
| Air cooling overhead | 0.500 |
| W1 cooling overhead | 0.400 |
| W4 cooling overhead | 0.050 |

Using the LRZ data center overhead information from Table 4.8, and the knowledge that CooLMUC is 100% cooled using W4, sPUE can be determined:

$$\text{sPUE}_{\text{CooLMUC}} = 1 + 0.075 + 1 * 0.05 = 1.125$$

Using the system internal cooling overhead and sPUE, the final yearly (for 2014) energy savings for CooLMUC would be:

$$\text{Savings}_{\text{CooLMUC}} = 251.77 * 1.23 * 1.125 = 348.39\text{kWh}$$

If one, in the most simplistic way, scales the CooLMUC (43 kW average system power) result to SuperMUC (2.4 MW average power consumption) then the possible yearly savings would be: 56 times 348.39 kWh = 19509.84 kWh.

### 4.3.5 Summary

This chapter 4 Improving Energy Efficiency discussed the difference between power and energy optimization and quantified the existence of node power variation in current homogeneous HPC systems. This system property has an impact on other HPC activities related to power and energy management. It impacts the energy and power prediction for applications (and by extension the HPC system) [96] and should be part of power aware scheduling [97]. It has also been shown by the author that this property is frequency dependent [98].

This variation can be used to save energy. Three techniques were proposed and evaluated, namely:

1. Node Power Aware Scheduling - which, unfortunately, is not beneficial for current HPC systems.

2. Node Power Aware System Partitioning - should be used (depending on interconnect requirements and system size) but might have low impact on systems where all nodes are in multiple queues.

3. Node Ranking Based on Power Variation - should be used since it requires very low effort for most scheduling systems. Unfortunately, for HPC systems, the savings potential is relatively low since a node utilization of over 90% limits the effectiveness of this technique.

Taking advantage of node power variation can save energy. The proposed techniques can be used in conjunction with any other energy saving effort and are not limited to HPC data centers. From the three techniques, "Node Power Aware System Partitioning" and "Node Ranking Based on Power Variation" are the most practical since they can provide energy savings over the lifetime of the IT system without requiring to much effort.

Even though the use of node power variation doesn't provide impressive energy savings for HPC data centers, it can potentially save a lot of energy for data centers with low average node utilization. Liu [88] reports an average utilization of AMAZON cloud of 15% for one week. Industry surveys approximate the average utilization rate of most server clusters at 15% [89], [90]. A white paper written by the Natural Resources Defense Council in 2014 [91] states that current hyper-scale cloud providers (which accounted for only 4% of the overall data centers power consumption in 2011) can realize average utilization rates of 40%. And lastly, research done by Google indicates that typical server clusters (which accounted for 95% of the overall data centers power consumption in 2011) have an average utilization anywhere from 10 to 50 percent [91].

Seeing that the average utilization of standard server clusters (which account for 95% of the overall data centers) is between 10-50%, taking advantage of node power variation can save a substantial amount of energy. For an average cluster utilization of 50% picking the best nodes (chapter 4.3.4) would save 4.3% off energy for the CoolMUC cluster.

# 5 Summary

The work presented here answered the following research questions and proposed thesis:

**Question:** What are the parts of an HPC data center that are important for energy efficiency?

**Thesis:** There exists a unified way to approach data center energy efficiency that is applicable to all data centers.

***Answer:*** The developed *4 Pillar Framework* (section 2.2 The *4 Pillar Framework* for energy efficient HPC data centers) provides for the first time a common view of energy relevant components of a data center which helps to formalize energy efficiency research. The 4 Pillars are:

1. Data Center Infrastructure

2. IT System Hardware

3. IT System Software

4. Applications

The 4 Pillars are influenced by outside conditions (utility contract, climate conditions, possible waste heat re-use, etc.) and data center policies (SLA's, political requirements, internal constraints).

The *4 Pillar Framework* can be used to: classify current research; show areas of activities; identify gaps; and allow for a common base from which to approach data center energy efficiency improvements.

**Question:** How can power and energy related data be measured, collected, and evaluated?

**Thesis:** It should be possible to have one Key Performance Indicator that measures the energy efficiency of a data center.

***Answer:*** The developed tool *PowerDAM* (section 3.2 Wholistic Power and Energy Data Collection - PowerDAM) allows for the unified collection of power and energy relevant data from different data center monitoring and control systems. This unified data archive enables the calculation of KPI's. DCEE (section 3.3 A New Metric to Measure Data Center Energy Efficiency - Data Center Energy Efficiency (DCEE)) is a metric to calculate the energy efficiency of the complete data center for a specific workload or workload mix. It is the first metric bridging the gap between IT system and data center infrastructure.

**Question:** Is there any energy saving potential not realized related to the HPC system hardware?

**Thesis:** IT hardware manufacturing tolerances will influence the energy consumption of large scale homogeneous IT systems.

**Answer:** The presented work showed and quantified the existence of node power variability in homogenous large scale cluster systems (section 4.3 Saving energy by taking advantage of node power variability in homogenous HPC systems). The three proposed techniques: *Node Power Aware Scheduling*; *Node Power Aware System Partitioning*; and *Node Ranking Based on Power Variation* can be used to take advantage of the hardware manufacturing tolerances to save energy (subsection 4.3.1 Node power aware scheduling and subsection 4.3.2 Node switch-off, node power aware system partitioning, and node ranking based on power variation). Unfortunately, this doesn't work so well for supercomputers since the node usage is very high during the year (around 90%). Fortunately, these techniques work very well for standard IT systems where the node runtime throughout the year is between 10% and 50%.

# 6 Conclusion and Future Work

The systematic approach to energy efficiency facilitated by the *4 Pillar Framework* led to:

- The identification of the current KPI gap between data center infrastructure and IT systems leading to the definition of *DWPE* and *DCEE*.

- The development of PowerDAM since data from all data center pillars and external influences need to be collected for any wholistic data center power and energy optimization.

- The discovery of node power variability in homogeneous large scale cluster systems through the use of PowerDAM.

Next steps are to define the major influencing factors for the data center energy efficiency so that the available control knobs can be defined. The first step towards that goal is the prediction of the data center Coefficient of Performance (COP). The author's paper "Using Machine Learning for Data Center Cooling Infrastructure Efficiency Prediction" [68] discusses the machine learning approach to COP prediction used at LRZ. From this work 6 major factors could be identified that define the COP of the chiller-less cooling infrastructure at LRZ.

These factors are:

1. IT Power - the power consumption of the IT system (this can be predicted based on running jobs [96] [98] and can be influenced via job scheduling)

2. $\Delta$T - temperature difference between the HT-DLC inlet and outlet (this can be controlled via the building automation system)

3. HT-DLC inlet temperature - influences how efficient the cooling infrastructure can remove the IT heat (this can be controlled via the building automation system)

4. Number of Cooling Towers - since the chiller-less cooling circuit at LRZ has 4 cooling towers, the number of active towers is important (one could trade off fan power of the cooling towers with additional or less cooling towers)

5. Wetbulb Temperature - this effects the efficiency of hybrid and wet cooling towers. (cannot be controlled directly but infrastructure presets, like inlet temperature, could be adjusted by the building automation system)

6. Flow Rate - the flow rate inside the chiller-less cooling circuit (this can be controlled via the building automation system)

Some factors influence each other. For example, changing the ΔT will affect the flow rate. For a higher ΔT a lower flow rate is required whereas for a smaller ΔT a higher flow rate is needed.

Figure 6.1 shows the impact of the HT-DLC inlet water temperature on the power consumption of the LRZ cooling infrastructure (here specifically the power consumption and cold generation from one cooling tower circuit KLT14).
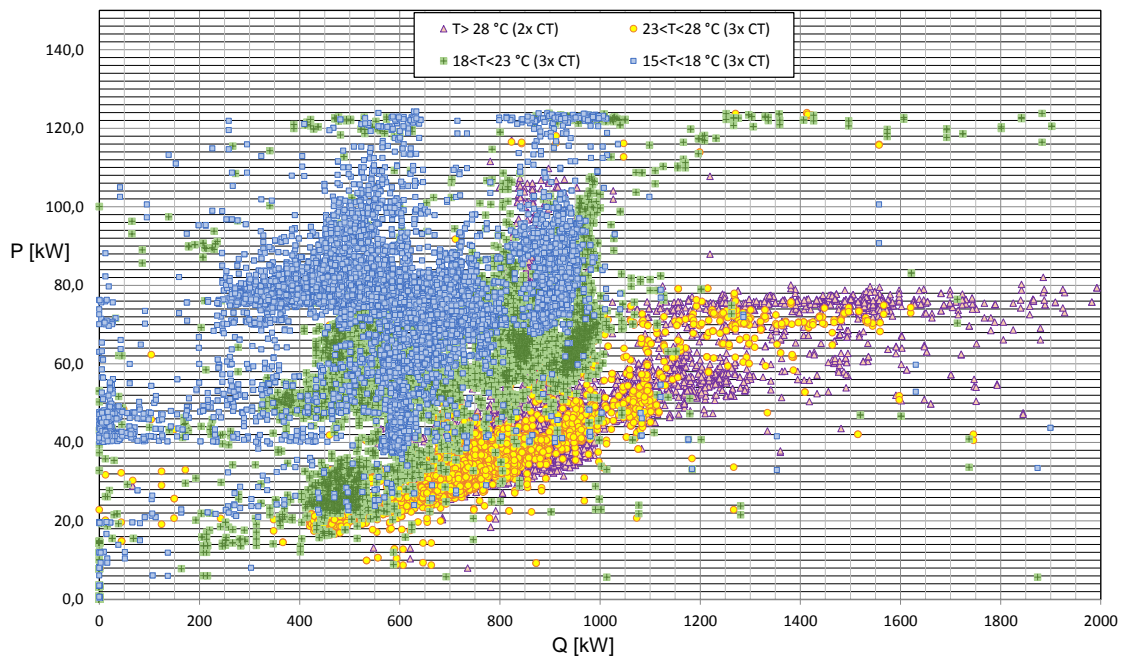
## KLT14 P(Q) (15.3.-15.4.2015)



Figure 6.1: Cooling loop KLT14 power consumption vs. generated cold.

Some observations that can be made are:

- The higher the HT-DLC inlet temperature the lower the power consumption (what is not visible here is that the higher inlet temperature leads to higher IT power consumption, as shown by the author in [99]).

- Running less cooling towers in parallel is more energy efficient.

- With 80kW of electrical power, cold energy ranging from 150kW till 2000kW can be generated.

- The maximum possible power consumption of the cooling tower circuit is slightly above 120kW.

Because of the complexity of the interdependency of the major factors (*IT Power*, *ΔT*, *HT-DLC inlet temperature*, *Number of Cooling Towers*, *Wetbulb Temperature*, and *Flow Rate*), more research is needed to quantify this interdependency, to identify and quantify opposing

relationships and to propose robust multi-parameter energy optimization strategies for the cooling infrastructure that can react automatically to the change of one or more of the factors.

This information needs to be combined with IT system hardware and cooling properties. For example, the journal paper by the author "Analysis of the Efficiency Characteristics of the First High-Temperature Direct Liquid Cooled Petascale Supercomputer and Its Cooling Infrastructure" [54] shows for the first time how the energy consumption of an IT system, using HT-DLC, is influenced by the cooling water temperature.

This closer connection between *Pillar 1* and *Pillar 2* needs to be further extended with application power and energy consumption information. An example is also found in the paper by the author [54] showing that application power and energy consumption is affected by Dynamic Voltage and Frequency Scaling (DVFS), power capping, node power variability, and applications.

And finally, all information needs to be integrated into the batch scheduling system. This will allow future scheduling systems to have the new capability to make decisions based on the power and energy policies at individual data centers leading finally to true data center energy aware scheduling.

# Bibliography

[1] Jonathan G Koomey. Worldwide electricity used in data centers. http://iopscience.iop.org/1748-9326/3/3/034008/pdf/1748-9326_3_3_034008.pdf, 2008.

[2] Jonathan Koomey. Growth in data center electricity use 2005 to 2010. http://www.twosides.us/content/rspdf_218.pdf, 2011.

[3] DCD Intelligence. Is the industry getting better at using power? *Data Center Dynamics FOCUS 33, January/February 2014*, 33:16 − 17, 2014.

[4] Dan Azevedo, Disney Alan French, and Emerson Network Power. PUE: A COMPREHENSIVE EXAMINATION OF THE METRIC. 2012.

[5] Matt Stansberry. 2013 uptime institute annual data center industry survey report and full results. http://www.data-central.org/resource/collection/BC649AE0-4223-4EDE-92C7-29A659EF0900/uptime-institute-2013-data-center-survey.pdf, 2014.

[6] L. Stobbe, M. Proske, H. Zedel, R. Hintemann, J. Clausen, and S. Beucker. Entwicklung des IKT-bedingten Strombedarfs in Deutschland - Studie von Fraunhofer IZM und Borderstep im Auftrag des Bundesministeriums für Wirtschaft und Energie, 2015.

[7] ASHRAE TC 9.9.2011. 2011 thermal guidelines for liquid cooled data processing environments. *White paper*, 2011.

[8] The Green Grid. http://www.thegreengrid.org/.

[9] EUROPEAN COMMISSION; DIRECTORATE-GENERAL; JOINT RESEARCH CENTRE; Institute for Energy and Transport - Renewable Energy Unit. European Code of Conduct on Data Centre Energy Efficiency. http://iet.jrc.ec.europa.eu/energyefficiency/sites/energyefficiency/files/introductory_guide_v2_0_2.pdf.

[10] Karen J. Fryer, Jiju Antony, and Alex Douglas. Critical success factors of continuous improvement in the public sector: A literature review and some key findings. *Total Quality Management*, 19(5):497–517, 2007.

[11] Torsten Wilde, Axel Auweter, and Hayk Shoukourian. The 4 Pillar Framework for energy efficient HPC data centers. http://dx.doi.org/10.1007/s00450-013-0244-6, 2014.

[12] P. Kogge. Exascale computing study: Technology challenges in achieving exascale systems. Univ. of Notre Dame, CSE Dept. Tech. Report TR-2008-13, Sept. 28, 2008.

[13] J. W. Choi, D. Bedard, R. Fowler, and R. Vuduc. A roofline model of energy. In *2013 IEEE 27th International Symposium on Parallel and Distributed Processing*, pages 661–672, May 2013.

[14] Axel Auweter, Arndt Bode, Matthias Brehm, Luigi Brochard, Nicolay Hammer, Herbert Huber, Raj Panda, Francois Thomas, and Torsten Wilde. A Case Study of Energy Aware Scheduling on SuperMUC. In *Proceedings of the 29th International Conference on Supercomputing - Volume 8488*, ISC 2014, pages 394–409, New York, NY, USA, 2014. Springer-Verlag New York, Inc.

[15] Q. Tang, S. K. S. Gupta, and G. Varsamopoulos. Energy-efficient thermal-aware task scheduling for homogeneous high-performance computing data centers: A cyber-physical approach. *IEEE Transactions on Parallel and Distributed Systems*, 19(11):1458–1472, Nov 2008.

[16] DEEP. http://www.deep-project.eu/deep-project/EN/Home/home_node.html.

[17] B. Videau (CNRS), D. Brayford (LRZ), M. Allalen (LRZ), P. Lanucara (CINECA), N. Sanna (CINECA), F. Mantovani (BSC), R. Halver (JSC), D. Broemmel (JSC), JH. Meinke (JSC), K. Pouget (CNRS), J.-F. Mehaut (CNRS), L. Genovese (CEA), C. Gomez (BSC), and A. Rico (BSC). MONT-BLANC Deliverables D4.4, D4.5, and D4.6 Report from Woek Package 4 "Exascale Applications". http://montblanc-project.eu/sites/default/files/MB_D4.4-D4.5-D4.6-FINAL.pdf, 2015.

[18] Torsten Wilde, Detlef Labrenz, Michael Ott, Axel Auweter, Ingmar Meijer, Patrick Ruch, Markus Hilger, Steffen Kühnert, and Herbert Huber. CooLMUC-2: A Supercomputing Cluster with Heat Recovery for Adsorption Cooling. *accepted: SEMI-THERM 33, 13th-17th March 2017, San Jose, USA*, 2017.

[19] Jonathan Eastep, Steve Sylvester, Christopher Cantalupo, Federico Ardanaz, Brad Geltz, Asma Al-Rawi, Fuat Keceli, and Kelly Livingston. Global Extensible Open Power Manager: A Vehicle for HPC Community Collaboration Toward Co-Designed Energy Management Solutions. In *7th International Workshop in Performance Modeling, Benchmarking and Simulation of High Performance Computer Systems (PMBS16)*, 2016.

[20] COOL-EM-ALL. http://tricoryne.man.poznan.pl/web/guest.

[21] ECO2CLOUDS. http://eco2clouds.com/.

[22] P. van der Meer, A. van Staveren, and A. van Roermund. *Low-Power Deep Sub-Micron CMOS Logic: Sub-threshold Current Reduction*. The Springer International Series in Engineering and Computer Science. Springer US, 2012.

[23] SchedMD. Simple Linux Utility for Resource Management Workload Scheduler (SLURM). http://www.schedmd.com/slurmdocs/slurm.html.

[24] IBM. Tivoli workload scheduler LoadLeveler. `http://www.ibm.com/systems/software/loadleveler/`.

[25] Google Datacenters Information. `https://www.google.com/about/datacenters/efficiency/internal/#water-and-cooling`.

[26] Joe Kava. 10 things we know to be true - a behind the scenes look at google data center financial and operational efficiencies. `https://www.youtube.com/watch?v=pghac0H0x_U`, 2013.

[27] BSR/ASHRAE. Standard 90.4p - 3rd isc public review draft energy standard for data centers. `https://osr.ashrae.org/sitepages/showdoc2.aspx/ListName/Public/Review/Draft/Standards/ItemID/1427/IsAttachment/N/BSR_ASHRAEStd90.4P3rdISCPPRDraft.pdf`, 2016.

[28] Michael K Patterson, Stephen W Poole, Chung-Hsing Hsu, Don Maxwell, William Tschudi, Henry Coles, David J Martinez, and Natalie Bates. TUE, a New Energy-Efficiency Metric Applied at ORNLs Jaguar. In *Supercomputing*, pages 372–382. Springer, 2013.

[29] Green500. `http://www.Green500.org/`.

[30] Rong Ge, Xizhou Feng, and Kirk W. Cameron. Performance-constrained distributed dvs scheduling for scientific applications on power-aware clusters. In *Proceedings of the 2005 ACM/IEEE conference on Supercomputing*, SC '05, pages 34–, Washington, DC, USA, 2005. IEEE Computer Society.

[31] Vincent W Freeh, David K Lowenthal, Feng Pan, Nandini Kappiah, Robert Springer, Barry L Rountree, and Mark E Femal. Analyzing the energy-time trade-off in high-performance computing applications. *Parallel and Distributed Systems, IEEE Transactions on*, 18(6):835–848, 2007.

[32] Chung-hsing Hsu and Wu-chun Feng. A power-aware run-time system for high-performance computing. In *Proceedings of the 2005 ACM/IEEE conference on Supercomputing*, SC '05, pages 1–, Washington, DC, USA, 2005. IEEE Computer Society.

[33] Anton Beloglazov, Jemal Abawajy, and Rajkumar Buyya. Energy-aware resource allocation heuristics for efficient management of data centers for Cloud computing. *Future Generation Computer Systems*, 28(5):755 – 768, 2012. Special Section: Energy efficiency in large-scale distributed systems.

[34] Barry Rountree, Dong H. Ahn, Bronis R. de Supinski, David K. Lowenthal, and Martin Schulz. Beyond dvfs: A first look at performance under a hardware-enforced power bound. *2013 IEEE International Symposium on Parallel And Distributed Processing, Workshops and Phd Forum*, 0:947–953, 2012.

[35] Tapasya Patki, David K. Lowenthal, Barry Rountree, Martin Schulz, and Bronis R. de Supinski. Exploring hardware overprovisioning in power-constrained, high performance computing. In *Proceedings of the 27th International ACM Conference on International Conference on Supercomputing*, ICS '13, pages 173–182, New York, NY, USA, 2013. ACM.

[36] Jenifer Hopper. Reduce Linux power consumption, Part 3: Tuning results. `http://www.ibm.com/developerworks/library/l-cpufreq-3/l-cpufreq-3-pdf.pdf`.

[37] Venkatesh Pallipadi and Alexey Starikovskiy. The ondemand governor. In *Proceedings of the Linux Symposium*, volume 2, pages 215–230. sn, 2006.

[38] Luigi Brochard, Raj Panda, Don Desota, Francois Thomas, and Rob Bell Jr. Optimizing performance and energy of high performance computing applications. *Parallel Computing: From Multicores and GPU's to Petascale*, 19:455, 2010.

[39] R.H. Bell, L. Brochard, D.R. DeSota, R.D. Panda, and F. Thomas. Energy-aware job scheduling for cluster environments, December 17 2013. US Patent 8,612,984.

[40] Lizhe Wang, Samee U Khan, and Jai Dayal. Thermal aware workload placement with task-temperature profiles in a data center. *The Journal of Supercomputing*, 61(3):780–803, 2012.

[41] Ayan Banerjee, Tridib Mukherjee, Georgios Varsamopoulos, and Sandeep K. S. Gupta. Cooling-aware and thermal-aware workload placement for green hpc data centers. In *Proceedings of the International Conference on Green Computing*, GREENCOMP '10, pages 245–256, Washington, DC, USA, 2010. IEEE Computer Society.

[42] Leibniz Supercomputing Centre of the Bavarian Academy of Sciences and Humanities. Leibniz Supercomputing Centre yearly report 2014. `http://www.lrz.de/wir/berichte/JB/JBer2014.pdf`, 2014.

[43] Jee Whan Choi and Richard Vuduc. A roofline model of energy. In *Technical report no. GT-CSE-12-01*. Georgia Institute of Technology, School of Computational Science and Engineering, Atlanta, GA, USA, 2012.

[44] Momme Allalen, Christoph Bernau, Arndt Bode, David Brayford, Matthias Brehm, Nicolay Hammer, Herbert Huber, Ferdinand Jamitzky, Anupam Karmakar, Carmen Navarrete, and Helmut Satzger. Extreme scaling workshop at lrz. 01 2013.

[45] Alexander Breuer, Alexander Heinecke, Leonhard Rannabauer, and Michael Bader. *High Performance Computing: 30th International Conference, ISC High Performance 2015, Frankfurt, Germany, July 12-16, 2015, Proceedings*, chapter High-Order ADER-DG Minimizes Energy- and Time-to-Solution of SeisSol, pages 340–357. Springer International Publishing, Cham, 2015.

[46] SIMOPEK. `http://www.simopek.de`.

[47] Torsten Wilde, Tanja Clees, Hayk Shoukourian, Nils Hornung, Michael Schnell, Inna Torgovitskaia, Eric Lluch Alvarez, Detlef Labrenz, and Horst Schwichtenberg. Increasing data center energy efficiency via simulation and optimization of cooling circuits - A practical approach. In *Energy Informatics - 4th D-A-CH Conference, EI 2015, Karlsruhe, Germany, November 12-13, 2015, Proceedings*, pages 208–221, 2015.

[48] Joe Loper and Sara Parr. Alliance to save energy, energy efficiency in data centers: A new policy frontier. http://www.ase.org/resources/energy-efficiency-data-centers-new-policy-frontier, Jan 2007.

[49] Japan National Body/Green IT Promotion Council. DPPE: Holistic framework for data center energy efficiency - kpis for infrastructure, it equipment, operation (and renewable energy). http://home.jeita.or.jp/greenit-pc/topics/release/pdf/dppe_-e_20120824.pdf, Aug 2012.

[50] Hayk Shoukourian, Torsten Wilde, Axel Auweter, and Arndt Bode. *Monitoring Power Data: A first step towards a unified energy efficiency evaluation toolset for HPC data centers*. Elsevier, 2013.

[51] Hayk Shoukourian, Torsten Wilde, Axel Auweter, Arndt Bode, and Petra Piochacz. Towards a unified energy efficiency evaluation toolset: an approach and its implementation at Leibniz Supercomputing Centre (LRZ). In *ICT4S 2013: Proceedings of the First International Conference on Information and Communication Technologies for Sustainability*, pages 276–282, 2013.

[52] T. Wilde, A. Auweter, M.K. Patterson, H. Shoukourian, H. Huber, A. Bode, D. Labrenz, and C. Cavazzoni. DWPE, a new data center energy-efficiency metric bridging the gap between infrastructure and workload. In *High Performance Computing Simulation (HPCS), 2014 International Conference on*, pages 893–901, July 2014.

[53] dictionary.com. http://www.dictionary.com/browse/efficiency/.

[54] Hayk Shoukourian, Torsten Wilde, Herbert Huber, and Arndt Bode. Analysis of the Efficiency Characteristics of the First High-Temperature Direct Liquid Cooled Petascale Supercomputer and Its Cooling Infrastructure. *Journal of Parallel and Distributed Computing*, 2017.

[55] Natalie Bates, Girish Ghatikar, Ghaleb Abdulla, Gregory A. Koenig, Sridutt Bhalachandra, Mehdi Sheikhalishahi, Tapasya Patki, Barry Rountree, and Stephen Poole. Electrical grid and supercomputing centers: An investigative analysis of emerging opportunities and challenges. *Informatik-Spektrum*, 38(2):111–127, 2015.

[56] Michael Knobloch, Maciej Foszczynski, Willi Homberg, Dirk Pleiter, and Hans Böttiger. Mapping fine-grained power measurements to hpc application runtime characteristics on ibm power7. *Computer Science - Research and Development*, 29(3):211–219, Aug 2014.

[57] Hayk Shoukourian, Torsten Wilde, Herbert Huber, and Arndt Bode. Analysis of the efficiency characteristics of the first high-temperature direct liquid cooled petascale supercomputer and its cooling infrastructure. *Journal of Parallel and Distributed Computing*, 107:87 – 100, 2017.

[58] L. Johnsson(KTH), G. Netzer(KTH), E. Boyer(CINES), S. Graf(Juelich), W. Homberg(Juelich), G. Koutsou(CaSToRC), J. Jakic(IPB), R. Januszewski(PSNC), N. Puzovic(BSC), T. Roeblitz(SIGMA/UiO), O.-W. Saastad(SIGMA/UiO), B. Schembera(HLRS), G. Schwarz(Juelich), H. Shoukourian(LRZ), V. Strumpen(JKU),

S. Thiell(CEA), G.-C. deVerdiere(CEA), and T. Wilde(LRZ). PRACE 1IP-WP9 Deliverable D9.3.3 - Report on prototypes evaluation. http://www.prace-ri.eu/IMG/pdf/d9.3.3.pdf, 2013.

[59] Hayk Shoukourian. *Adviser for Energy Consumption Management: Green Energy Conservation*. PhD thesis, München, Technische Universität München (TUM), 2015.

[60] K. Cassirer, T. Clees, B. Klaassen, I. Nikitin, and L. Nikitina. *MYNTS User's Manual, Release 3.3*. Fraunhofer SCAI, Sankt Augustin, 2015. www.scai.fraunhofer.de/mynts.

[61] BILL Tschudi, OTTO Vangeet, J Cooley, and D Azevedo. Ere: A metric for measuring the benefit of reuse energy from a data center. *White Paper*, 29, 2010.

[62] Energy Efficiency High Performance Computing Working Group (EEHPCWG. https://eehpcwg.llnl.gov/.

[63] Japan National Body/Green IT Promotion Council. DPPE: Holistic framework for data center energy efficiency - kpis for infrastructure, it equipment, operation (and renewable energy). http://home.jeita.or.jp/greenit-pc/topics/release/pdf/dppe_-e_20120824.pdf, Aug 2012.

[64] Global Metrics Harmonization Task Force. Global taskforce reaches agreement on measurement protocols for gec. *ERF, and CUE–Continues Discussion of Additional Energy Efficiency Metrics*, 2012.

[65] Timo Minartz, Julian M Kunkel, and Thomas Ludwig. Simulation of power consumption of energy efficient cluster hardware. *Computer Science-Research and Development*, 25(3-4):165–175, 2010.

[66] Green Grid Metrics. Describing datacenter power efficiency. *Green Grid Technical Committee White Paper*, 2007.

[67] Arndt Bode. Keynote europar 2013 - energy to solution: A new mission for parallel computing. 2013.

[68] H. Shoukourian, T. Wilde, D. Labrenz, and A. Bode. Using machine learning for data center cooling infrastructure efficiency prediction. In *2017 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pages 954–963, May 2017.

[69] Partnership for Advanced Computing in Europe. http://www.prace-ri.eu/.

[70] James Milano and Pemela Lembke. *IBM System Blue Gene Solution: Blue Gene/Q Hardware Overview and Installation Planning*. IBM Redbook, 2013.

[71] Gracia, Michael and Barragy, Ted and Blanc, Jean-Yves. Experiences with Oil Immersion Cooling in a Processing Datacenter 2017. http://www.cscs.ch/fileadmin/events/Experiences_with_Oil_Immersion_Cooling_in_a_Processing_Datacenter_2017_Garcia.pdf.

[72] Barclay (DoD), Craig. Maximizing Mission Capability By Mitigating Trapped Power Capacity 2015. http://hpm.ornl.gov/Archives/HPM15/documents/HPM2015_Barclay.pdf.

[73] Barry Rountree, Dong H Ahn, Bronis R De Supinski, David K Lowenthal, and Martin Schulz. Beyond dvfs: A first look at performance under a hardware-enforced power bound. In *Parallel and Distributed Processing Symposium Workshops & PhD Forum (IPDPSW), 2012 IEEE 26th International*, pages 947–953. IEEE, 2012.

[74] Tapasya Patki, David K Lowenthal, Barry Rountree, Martin Schulz, and Bronis R De Supinski. Exploring hardware overprovisioning in power-constrained, high performance computing. In *Proceedings of the 27th international ACM conference on International conference on supercomputing*, pages 173–182. ACM, 2013.

[75] Andrey Semin. Keynote Ena-HPC 2011 - HPC systems energy efficiency optimization thru hardware-software co-design on Intel technologies. www.ena-hpc.org/2011/talks/semin-slides.pdf, 2011.

[76] Giorgio Luigi Valentini, Walter Lassonde, Samee Ullah Khan, Nasro Min-Allah, Sajjad A. Madani, Juan Li, Limin Zhang, Lizhe Wang, Nasir Ghani, Joanna Kolodziej, Hongxiang Li, Albert Y. Zomaya, Cheng-Zhong Xu, Pavan Balaji, Abhinav Vishnu, Fredric Pinel, Johnatan E. Pecero, Dzmitry Kliazovich, and Pascal Bouvry. An overview of energy efficiency techniques in cluster computing systems. *Cluster Computing*, 16(1):3–15, 2013.

[77] Energy Efficient HPC Working Group. Energy Efficiency Considerations for HPC Procurement Documents: 2014. https://eehpcwg.llnl.gov/documents/compsys/aa_procurement_2014.pdf, 2014.

[78] Torsten Wilde, Axel Auweter, Hayk Shoukourian, and Arndt Bode. Taking Advantage of Node Power Variation in Homogenous HPC Systems to Save Energy. In *High Performance Computing - 30th International Conference, ISC High Performance 2015, Frankfurt, Germany, July 12-16, 2015, Proceedings*, pages 376–393, 2015.

[79] Thomas Scogland, Jonathan Azose, David Rohr, Suzanne Rivoire, Natalie Bates, Daniel Hackenberg, and Torsten Wilde. Node variability in large-scale power measurements: Perspectives from the green500, top500 and eehpcwg. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, SC '15, pages 74:1–74:11, New York, NY, USA, 2015. ACM.

[80] D. Hackenberg, R. Schöne, T. Ilsche, D. Molka, J. Schuchart, and R. Geyer. An Energy Efficiency Feature Survey of the Intel Haswell Processor. In *International Parallel and Distributed Processing Symposium Workshop (IPDPSW)*, pages 896–904, May 2015.

[81] Song Huang, Michael Lang, Scott Pakin, and Song Fu. Measurement and characterization of haswell power and energy consumption. In *Proceedings of the 3rd International Workshop on Energy Efficient Supercomputing*, E2SC '15, pages 7:1–7:10, New York, NY, USA, 2015. ACM.

[82] Yuichi Inadomi, Tapasya Patki, Koji Inoue, Mutsumi Aoyagi, Barry Rountree, Martin Schulz, David Lowenthal, Yasutaka Wada, Keiichiro Fukazawa, Masatsugu Ueda, et al. Analyzing and mitigating the impact of manufacturing variability in power-constrained supercomputing. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, page 78. ACM, 2015.

[83] E. A. Burton, G. Schrom, F. Paillet, J. Douglas, W. J. Lambert, K. Radhakrishnan, and M. J. Hill. FIVR - fully integrated voltage regulators on 4th generation Intel Core SoCs. In *Applied Power Electronics Conference and Exposition (APEC)*, pages 432–439. IEEE, March 2014.

[84] Daniel Hackenberg, Robert Schöne, Daniel Molka, Matthias S. Müller, and Andreas Knüpfer. Quantifying power consumption variations of hpc systems using spec mpi benchmarks. *Computer Science - Research and Development*, 25(3-4):155–163, 2010.

[85] John D Davis, Suzanne Rivoire, Moises Goldszmidt, and Ehsan K Ardestani. Accounting for variability in large-scale cluster power models. *Exascale Evaluation and Research Techniques*, 2011.

[86] Jack Dongarra and Michael A Heroux. Toward a new metric for ranking high performance computing systems. *Sandia Report, SAND2013-4744*, 312, 2013.

[87] T. Arber and et al. EPOCH: Extendable PIC Open Collaboration. http://ccpforge.cse.rl.ac.uk/gf/project/epoch/, 2014.

[88] Huan Liu. A measurement study of server utilization in public clouds. In *Dependable, Autonomic and Secure Computing (DASC), 2011 IEEE Ninth International Conference on*, pages 435–442, Dec 2011.

[89] Mark Aggar (Microsoft). The IT Energy Efficiency Imperative. *White paper*, 2011.

[90] Ravi A. Giri (Staff Engineer, Intel IT) and Anand Vanchi (Solutions Architect, Intel Data Center Group). Increasing Data Center Efficiency with Server Power Measurements. *IT@Intel White Paper*, page 7, 2011.

[91] Josh Whitney and Pierre Delforge. Scaling up energy efficiency across the Data Center Industry: evaluating Key Drivers and Barriers. *Data Center Efficiency Assessment*, 2014.

[92] F. Alvarruiz, C. de Alfonso, M. Caballer, and V. Hern'ndez. An energy manager for high performance computer clusters. In *2012 IEEE 10th International Symposium on Parallel and Distributed Processing with Applications*, pages 231–238, July 2012.

[93] Vikas Ashok Patil and Vipin Chaudhary. Rack aware scheduling in HPC data centers: an energy conservation strategy. *Cluster Computing*, 16(3):559–573, 2013.

[94] LRZ Loadleveler Documentation. https://www.lrz.de/services/compute/supermuc/loadleveler/.

[95] L. Johnsson(KTH), G. Netzer(KTH), E. Boyer(CINES), G. Koutsou(CaSToRC), R. Januszewski(PSNC), P. Carpenter(BSC), O.-W. Saastad(SIGMA/UiO), and

T. Wilde(LRZ). PRACE 1IP-WP9 Deliverable D9.3.4 - Final Report on Prototypes Evaluation. `http://www.prace-ri.eu/IMG/pdf/d9.3.4_1ip.pdf`, 2013.

[96] Hayk Shoukourian, Torsten Wilde, Axel Auweter, and Arndt Bode. Predicting the Energy and Power Consumption of Strong and Weak Scaling HPC Applications. *Supercomputing Frontiers And Innovations*, 1(2):20–41, 2014.

[97] Hayk Shoukourian, Torsten Wilde, Axel Auweter, and Arndt Bode. Power Variation Aware Configuration Adviser for Scalable HPC Schedulers. In *Proceedings of the 13 International Conference on High Performance Computing & Simulation, HPCS*, pages 71–79. IEEE, July 2015.

[98] Hayk Shoukourian, Torsten Wilde, Axel Auweter, Arndt Bode, and Daniele Tafani. Predicting energy consumption relevant indicators of strong scaling hpc applications for different compute resource configurations. In *Proceedings of the Symposium on High Performance Computing*, HPC '15, pages 115–126, San Diego, CA, USA, 2015. Society for Computer Simulation International.

[99] M. Ott, T. Wilde, and H. Ruber. ROI and TCO analysis of the first production level installation of adsorption chillers in a data center. In *2017 16th IEEE Intersociety Conference on Thermal and Thermomechanical Phenomena in Electronic Systems (ITherm)*, pages 981–986, May 2017.

# Definitions

Definitions used in this dissertation:

**Building Infrastructure** encompasses all energy transportation systems inside a data center needed to operate the IT equipment, in context of this dissertation this is mainly related to data center cooling and power conversion and power distribution.

**Continuous improvement** The Institute of Quality Assurance who defined continuous improvement as a gradual never-ending change which is: "focused on increasing the effectiveness and/or efficiency of an organisation to fulfil its policy and objectives. It is not limited to quality initiatives. Improvement in business strategy, business results, customer, employee and supplier relationships can be subject to continual improvement. Put simply, it means 'getting better all the time'."

**Compute-subsystem** That part of an HPC system that performs computational work, e.g. CPU, GPU, memory etc.

**Energy-to-Solution (EtS)** The energy used for running a specific application, this should include networking, cooling, and data center overheads but it is in most cases limited to the compute-subsystem.

**Efficiency** The ratio of the useful work performed by a machine, or in a process, to the total energy expended.

**Effectiveness** The degree to which something is successful in producing a desired result.

**Framework** in context of this dissertation is used in the more general meaning of frame of reference or foundation and not in the context of computer science.

**HPC Hardware** is encompassing all physical parts and components of an HPC system.

**HT-DLC** stands for High Temperature Direct Liquid Cooling which is a HPC system cooling technology not requiring the use of mechanical chillers.

**IT equipment (IT)** All data center information technology hardware, e.g. switches, router, computers, storage etc.

**IT system, HPC system, or supercomputer** One specific type of IT equipment.

**Workload** A specific user application running on the HPC system.

**Workload-mix** A combination of multiple workloads running on the HPC system.

**Pull Communication Model** Pull communication or pulling is a style of network communication where an active request for information is send to a receiver. The reverse is known as push communication or pushing, where information is automatically transmitted without a specific request.