

*Electrical Grid and Supercomputing Centers: An Investigative Analysis of Emerging Opportunities and Challenges**

Natalie Bates · Girish Ghatikar
Ghaleb Abdulla · Gregory A. Koenig
Sridutt Bhalachandra
Mehdi Sheikhalishahi · Tapasya Patki
Barry Rountree · Stephen Poole

Introduction

Supercomputing centers (SCs) with petascale¹ systems for high-performance computing (HPC) can have an outsized impact on their electricity service providers (ESPs), with peak power demands exceeding 20 MW and instantaneous power fluctuations of up to 8 MW. As the HPC community moves towards exascale computing², we anticipate that a growing number of facilities will be reaching or exceeding these service levels, with a significant potential effect on electrical grid reliability. In this paper we seek to understand how these anticipated usage patterns can be integrated safely into the power grid with minimal cost and disruption in order to manage this risk.

Being a “good citizen” on the electrical grid has several historical precedents. In the past, electrically-intensive industries such as aluminum smelters have received preferential pricing in return for predictable loads and flexibility in reducing power during periods of high consumption. A mutual understanding of concerns between SCs and ESPs can produce a symbiotic relationship that goes beyond the current producer–consumer paradigm, paving the way for possible integration of SCs with the grid. HPC-grid integration in the context of this study refers to the dynamic interaction and value between the demand-side resources (SCs) and the supply-side

resources (ESPs), as well as the relationship between the electricity grid and its markets.

The Energy Efficient HPC Working Group (EE HPC WG) investigates opportunities for large supercomputing sites to integrate more closely with their ESPs. We seek to understand the willingness of SCs to cooperate with their ESPs, their expectations from their ESPs, and the feasible measures that SCs could employ to help their ESPs. To achieve our objectives we developed a questionnaire and distributed it to the Top100 SCs in the United States.

This paper leverages prior work on datacenter and grid integration opportunities done by Lawrence Berkeley National Laboratory’s (LBNL) Demand Response Research Center. This prior work describes the challenges and opportunities for

DOI 10.1007/s00287-014-0850-0
© Springer-Verlag Berlin Heidelberg 2014

Natalie Bates
Energy Efficient HPC Working Group,
Washington, USA
E-Mail: natalie.jean.bates@gmail.com

Girish Ghatikar
Lawrence Berkeley National Laboratory,
California, USA

Ghaleb Abdulla · Barry Rountree
Lawrence Livermore National Laboratory,
California, USA

Gregory A. Koenig · Stephen Poole
Oak Ridge National Laboratory,
Tennessee, USA

Sridutt Bhalachandra
University of North Carolina,
North Carolina, USA

Mehdi Sheikhalishahi
University of Calabria,
Calabria, Italy

Tapasya Patki
University of Arizona,
Arizona, USA

* This work was partially performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344 as well as the U.S. Department of Energy’s Office of Energy Efficiency and Renewable Energy (EERE) by Lawrence Berkeley National Laboratory under Contract No. DE-AC02-05CH11231.

¹ Petascale computing refers to computing systems capable of at least 1015 operations or floating point instructions per second (FLOPS).

² Exascale computing refers to computing systems capable of at least 1018 operations or FLOPS.

Abstract

Some of the largest supercomputing centers (SCs) in the United States are developing new relationships with their electricity service providers (ESPs). These relationships, similar to other commercial and industrial partnerships, are driven by a mutual interest to reduce energy costs and improve electrical grid reliability. While SCs are concerned about the quality, cost, environmental impact, and availability of electricity, ESPs are concerned about electrical grid reliability, particularly in terms of energy consumption, peak power demands, and power fluctuations. The power demand for SCs can be 20 MW or more – the theoretical peak power requirements are greater than 45 MW – and recurring intra-hour variability can exceed 8 MW. As a result of this, ESPs may request large SCs to engage in demand response and grid integration.

This paper evaluates today's relationships, potential partnerships, and possible integration between SCs and their ESPs. The paper uses feedback from a questionnaire submitted to supercomputing centers on the Top100 List in the United States to describe opportunities for overcoming the challenges of HPC-grid integration.

datacenters and ESPs to interact with each other and how this integration can advance new market opportunities [21, 24]. This integration model describes programs that are used by some of the ESPs to encourage particular responses by their customers and methods used to balance the electrical grid supply and demand. This is referred to as *demand response (DR)*.

Eleven sites responded to the aforementioned questionnaire. Based on these responses, we noted a few primary observations:

- Only 20 % of SCs currently communicate with their ESPs about DR issues.
- SC managers believe that the candidate solutions most likely to be effective in responding to ESP requests involve coarse-grained power management techniques, job scheduling techniques, and shutting down computing resources.

- A stronger relationship, including DR capabilities, between SCs and ESPs, can lead to both energy savings and cost savings over time, and in some cases such capabilities might become a requirement for large SCs located in energy-challenged locations.

One of the most straightforward ways that SCs can begin the process of engaging in integration is by participating in efforts to develop software infrastructure to manage their electricity requirements in a tightly coupled manner with their ESPs, facilitating both energy efficiency and grid reliability. This will provide for extensive funding and cost analysis and help the community base future requirements for SCs and ESPs on facts and a proven set of measurements.

Our analysis in this paper focuses on SCs in the United States. However, the findings can be extended to and may relate to SCs in other countries with similar practices. Electrical grid infrastructure and market design are highly dependent on governmental regulations that vary across geographies. We restricted the initial analysis to the understanding of electricity markets in the United States. Future work can extend the analysis to electricity markets in Europe and other countries.

The paper is organized as follows. Sections “Electricity Service Providers” and “Supercomputing Centers” of this paper describe in greater detail the model for integrating the electrical grid and SCs. Section “Prior Work” reviews prior work in SC strategies. Section “Survey Results” provides the results of the questionnaire. Section “Opportunities, Solutions, and Barriers” discusses the several opportunities, solutions, and barriers that have been highlighted by the survey results. We offer our conclusions in Section “Conclusions and Next Steps” along with our plan for future work. The Appendix summarizes survey questions.

Electricity Service Providers

An ESP seeks to supply efficient and reliable generation, transmission, and distribution of electricity. *Market-based programs* employed by ESPs and consumers' participation are key to managing these electricity supply goals. While the goals describe ESPs overarching the objective for electricity supply, the *programs* describe the market products that the ESPs can offer to their consumers to achieve those goals. Such electricity market goals and demand-

side programs have been well studied for non-SC customer sectors [35].

Electricity Market Goals and Programs

Although critical to ESPs, the goals are generally not visible to the electricity consumer because they operate within the supply-side of the electric grid (for example, generation). These programs are the means by which customers get to engage in the electricity markets. The following is a summarized list and brief definitions of the key goals.

- *Transmission congestion*: The goal is to resolve congestion that occurs when there is not enough transmission capability to support all requests for transmission services. Transmission system operators must re-dispatch generation or, in the limit, deny some of these requests to prevent transmission lines from becoming overloaded.
- *Distribution congestion*: The goal is to resolve congestion that occurs when the distribution control system is overloaded. It generally results in deliveries that are held up or delayed.
- *Frequency response*: The goal is to keep grid frequency constant and in balance. Generators are typically used for frequency response, but any appliance that operates to a duty cycle (such as air conditioners and heat pumps) may be used to provide a constant and reliable grid balancing service by timing their duty cycles in response to system load.
- *Peak and reserve capacity*: The overall generation and extra capacity for supply during the peak or unforeseen high demand days.
- *Renewable integration*: The goal is to manage the variable uncertain generation nature of many renewable resources.

For efficient management of these goals, ESP programs encourage customer-side responses to manage demand for electricity at different time scales. Such market-based programs can be day-ahead or day-of. Day-ahead programs refer to timescales of notification and responses from customers that are determined based on advanced forecasting and capacity planning (for example, day-ahead, hourly wholesale electricity prices). The programs that are day-of are the ones when the notification and responses support

same-day capacity planning and/or emergency response.

An example of an ESP program that encourages energy efficiency would be to provide home consumers rebates and financial incentives to replace single pane windows with double pane windows. An example that illustrates programs that help with day-ahead or day-of demand management would be to provide credits and financial incentives to reduce load during high demand periods (such as hot summer afternoons when air conditioners are heavily utilized). The following is a summarized list and brief definitions of these key programs.

- *Energy efficiency*: Programs offered to reduce overall electricity consumption, thus eliminate the need for electricity generation.
- *Peak load reduction*: Programs used to *shed* load during peak times. Here the load reduced during the peak is either not used at a later time, or the load is *shifted* to, typically, non-peak hours.
- *Dynamic pricing*: Time varying pricing programs used to enable changes in the electricity consumption. The two types of pricing are peak and real-time. Peak pricing is pre-scheduled; however, the consumer does not know if a certain day will be a peak or a non-peak day until day-ahead or day-of. Real-time pricing is not pre-scheduled; prices can be set day-ahead or day-of, reflecting the real-time electricity system prices.
- *Regulation (up or down)*: Programs used to dispatch the portion of electricity generation reserves that are needed to manage changing demand at all times. Raising supply is *up* regulation and lowering supply is *down* regulation. There are many types of reserves (for example, operating reserves, ancillary services), distinguished by who manages them and what they are used for.

The following example illustrates the potential relevance of these programs to SCs. The generation capacity requirements and the timescales of customers' response vary across the United States for ESPs and system operators. One example is the New England independent system operator (ISO-NE) reserve capacity planning, which relies heavily on a day-ahead market program. This provides an opportunity for demand-side resources – such as SCs with local generation sources or flexible loads – to participate in the ISO-NE electricity markets. It also

makes the ISO-NE particularly less sensitive to major changes in electricity demand, which, as is discussed further in the questionnaire section, is an emerging characteristic of some of the largest SCs.

Supercomputing Centers

In November 2004, the Blue Gene/L system at Lawrence Livermore National Laboratory became the fastest computer in the Top 500 [1], displacing the NEC Earth Simulator, the previous champion. This change marked the transition from supercomputing gains based on ever-higher-performance components to systems that comprised far larger numbers of slow but energy-efficient components. However, total system power consumption continued to rise, and we are now poised to begin a second transition to “power-limited computing” and “power-aware computing”. The new model has been exemplified by the US Department of Energy issuing guidance that the first DOE exascale machine should not exceed 20 MW; effectively a $1000\times$ performance improvement with only a $3\times$ increase in power.

However, the problem is not as simple as provisioning 20 MW. Ultimately, SCs optimize for performance per dollar, not performance per watt, and flexibility in power consumption can be expected to result in lower overall prices. The use of green technologies such as wind and solar energy may also lead to cheaper but less predictable sources of power. In addition, as described in Section “Electricity Service Providers”, ESPs may request a change in timing and/or magnitude of demand by SCs. To adapt to this new landscape, SCs may employ one or more strategies to control their electricity demand.

The EE HPC WG took as their starting point a model developed by LBNL’s Demand Response Research Center. This model describes strategies that datacenters might employ for utility programs to manage their electricity and power requirements to lower costs and benefit from utility incentives. The EE HPC WG adopted this model to reflect the supercomputing environment focus (as opposed to the datacenter focus described by LBNL’s Demand Response Research Center).

It is important to highlight the differences between SCs and datacenters. Unlike datacenters, SCs are more performance oriented, have significantly higher system utilization, and use little or no

virtualization. Additionally, supercomputing applications are distinguished by their lack of geographical portability due to security concerns, data size, and machine-specific optimizations. We also note that SCs tend to be more energy efficient than datacenters. Power usage effectiveness (PUE) is a good measure for energy efficiency. PUE is the ratio of the total energy supplied for the facility to the amount of energy that actually reaches the IT infrastructure. A PUE of 1.0 is ideal. In our survey, none of the SCs exceeded a PUE of 1.53, while the average PUE for a datacenter falls in the range of 1.91 and 2.9 [34].

Strategies

Below we describe some of the strategies that SCs may use to adapt to the changing landscape.

- *Fine-grained power management* refers to the ability to control SC system power and energy with tools that offer high-resolution control and can target specific low level sub-systems. A typical example is CPU voltage and frequency scaling.
- *Coarse-grained power management* also refers to the ability to control SC system power and energy, but contrasts with fine-grained power management in that the resolution is low and it is generally done at a more aggregated level. A typical example is power capping.
- *Load migration* refers to temporarily shifting computing loads from an SC system in one site to a system in another location that has a stable power supply. This strategy can also be used in response to changes in electricity prices.
- *Job scheduling* refers to the ability to control SC system power by understanding the power profile of applications and queueing the applications based on those profiles.
- *Back-up scheduling* refers to deferring data storage processes to off-peak periods.
- *Shutdown* refers to a graceful shutdown of idle SC equipment. It usually applies when there is redundancy.
- *Lighting control* allows for datacenter lights to be shutdown completely.
- *Thermal management* is widening temperature set-point ranges and humidity levels for short periods.

These strategies can be used temporarily to modify loads in response to a request from an ESP. Addi-

tionally, some of these strategies could eventually be used at all times to improve overall energy efficiency if the SC sees no operational issues. Two examples may help to clarify this distinction. Temporary load migration is an example of a strategy that is well suited to responding to an ESP request, but is not likely to improve energy efficiency (lowering aggregate energy use). Fine-grained power management, on the other hand, can be used at all times and is more likely to be used for improving overall energy efficiency, unless the strategy is specifically used in response to an ESP's request.

Implementations

SC system power management has a very broad range of implementations and warrants greater exploration. For example, the coarse-grained and fine-grained strategies described above can be implemented at many levels of the system hierarchy – from node-level to site-level. We discuss these implementation approaches below.

- *Node level:* Controlling power ultimately requires control of individual components. Historically, this control has been accomplished through dynamic voltage/frequency scaling (DVFS), which allows the processor to use a lower voltage at the cost of a slower clock frequency. Newer technologies such as Intel's running average power limit leverage DVFS to guarantee that a user-specified processor power bound will, on average, not be exceeded over the duration of a short time window. DVFS can also be found on accelerator components such as NVIDIA's Kepler GPGPU. Other efforts reduce DRAM power by optimizing reads and writes, thus allowing the memory to spend more time in a lower-power state. Several processor configuration options have indirect but significant effects on power consumption. For example, the choice of the number of cores to use, whether or not to enable hyperthreading, and the use of "turbo" modes will change the power/performance curve.
 - *Job level:* Each of the node-level controls requires a tradeoff between power and performance. SC resources are typically oversubscribed, so degrading performance to save power and energy ultimately results in less science getting done. However, at the job level, load imbalance provides opportunities to slow nodes that are off the critical path of execution without slowing the overall job execution time.
- Traditionally, load rebalancing strategies have focused on moving bytes around the job allocation. With power control, we can now re-balance power as well as work. In American English, commas, semi-/colons, periods, and question and exclamation marks are located within quotation marks only when a complete thought or name is cited, such as a title or full quotation. When quotation marks are used, instead of a bold or italic typeface, to highlight a word or phrase, punctuation should appear outside of the quotation marks. A parenthetical phrase or statement at the end of a sentence is punctuated outside of the closing parenthesis (like this). (A parenthetical sentence is punctuated within the parentheses.)
- *System level:* While most SCs use time and space partitioning (where a node only runs a single job at a time), there are still shared resources that must be managed across jobs. Periodic checkpointing saves sufficient job states to a file system shared across jobs, so that a job may be restarted from a recent point in case a fault occurs. Because these checkpoints involve much more data motion than normal execution, power spikes can be observed at the node level (particularly DRAM), network, and file system. These checkpoints may need to be coordinated across large jobs to prevent unnecessary performance degradation.
 - *Scheduler level:* Up through the system level, power control is evaluated using the execution time of individual jobs. The scheduler optimizes for overall throughput rather than individual job performance. At this point, scheduling is a two-dimensional problem: jobs request a certain number of nodes for a certain duration. As power-limited computing becomes more common, schedulers will add power bounds to this mix: a job will be allowed nodes, time, and a certain number of watts (the responsibility for not exceeding the job power bound rests with the system software, not the user or application). The scheduler not only determines when jobs in the queue begin execution, but also what happens when a job exits the system. Depending on the priorities of already-running jobs and the priorities of jobs in the queue, the best solution in terms of throughput may be to idle the recently-freed nodes and redistribute the freed power to running jobs.
 - *Site level:* At the level of the machine room (or multiple machine rooms), decisions must be made as

to how much power should be allocated for cooling versus computation, which requires understanding how temperature interacts with performance. A higher intake air temperature uses less cooling power but results in higher static processor power and may limit opportunities for “turbo” mode in processors where it is available. As cooling power varies with outside air temperature, a single machine room temperature set point may not be the optimal solution in terms of overall performance.

Prior Work

This paper pulls together several diverse research domains. In this section, we provide an overview of prior work in these areas.

Power Management

Processor power management can be divided into two distinct eras. First, with the introduction of DVFS, users were able to change the CPU clock speed of their processors, lowering both voltage and, in most cases, energy: the workload used less power and ran longer, but the quadratic relation of power to frequency biased the results towards overall energy savings. Early work included several modeling efforts focused on the effects of CPU and memory-boundedness on delay and energy in MPI programs [5, 22, 26, 30, 42]. This work led to the CPUMiser [22] and Jitter runtime systems, which were designed to maximize energy saving consistent with a user-specified delay [28]. Treating energy savings as an optimization problem led to a linear programming solution [40]. The follow-on Adagio runtime system slowed only computation that could be proven to be off the critical path, leading to significant energy savings with only negligible slowdown [38]. These techniques were also applied to non-MPI datacenter workloads [18].

Other power saving approaches were attempted that did not use DVFS, but most were not deemed relevant to the supercomputing environment. A notable exception is dynamic concurrency throttling, where energy savings are realized by varying the number of threads at runtime [8–10, 37].

The research landscape changed considerably with the introduction of Intel’s Sandy Bridge processor. Turbo mode allowed higher clock frequencies to be reached so long as fewer cores were in use, making for a non-trivial power- performance tradeoff calculation. The running average power limit (RAPL)

technology provided an onboard power model that allowed the processor to both estimate power and, using rapid dithering of CPU clock frequencies, enforce a user-specified power bound across a short time window [11, 27]. For the first time, users were able to ask questions about performance under power bounds. This new capability arrived concurrently with the Department of Energy guideline that exascale machines would be subject to power (as opposed to energy) bounds.

Initial work showed that while processor performance at a fixed frequency was reproducible across processors, execution in turbo mode or under a power bound revealed significant performance variation [39]. Further work demonstrated a 2× performance improvement between conservative and optimal processor configurations while executing under a power bound [36].

Thermal Management

Thermal management is a key driver for improving the energy efficiency of datacenters as well as SCs. There are many strategies for thermal management that can improve energy efficiency, such as free cooling and proper airflow. This paper discusses two thermal management strategies that have an opportunity for grid integration. The first strategy is controlling the inlet temperature to the computing equipment, raising it as high as possible without causing reliability-induced hardware failures. The second strategy is using thermally aware job scheduling.

In 2011, the American Society of Heating, Refrigeration and Air Conditioning (ASHRAE) datacenter Technical Committee TC9.9 published guidelines that expanded the environmental range for datacenters and SCs [4]. The environmental range includes factors such as temperature, humidity and dew point, and the allowable rate of change. This expansion allows for maintaining high reliability while achieving gains in energy efficiency. These guidelines continue to be updated and the range continues to expand as the industry collects more historical data showing tradeoffs between reliability and environmental factors.

It is implicit in the ASHRAE guidelines that an SC might be able to increase temperature as a response to a request from an ESP. The guideline defines both recommended and allowable environmental ranges. It also specifies a maximum rate

of change, which is most stringent for tape drives. For SCs, the difference between the maximum recommended and allowable dry bulb temperature is a minimum of 9° F. The rate of change for tape drives is 9° F per hour (36° F for solid state computing systems). Therefore, assuming that SCs normally operate within the recommended range and that they are willing to operate on occasion in the allowable range (or beyond), it is theoretically possible to stay within ASHRAE thermal guidelines and use temperature excursion as a grid-integration strategy.

ASHRAE has also published a guideline on liquid cooling environmental ranges. At this point, however, the guidelines do not document the rate of change for liquid temperature. Although it is not explored in this paper, it may be possible to use increases in liquid cooling temperature as a grid-integration strategy as well.

Ghatikar et al. [24] describe field studies on using thermal management as a grid-integration strategy. They demonstrate increasing facility HVAC temperature set points in order to decrease HVAC power demand in two different field locations. There was only a small electricity demand decrease demonstrated.

Runtime cooling strategies are mostly job-placement centric. These techniques either aim to place incoming computationally intensive jobs in a thermally-aware manner on servers with lower temperatures or attempt to migrate or load-balance jobs from high-temperature servers to servers with lower temperatures.

Kaushik et al. [29] proposed T^* , a system that is aware of server thermal profiles and reliability as well as data semantics (computation job rates, job sizes, etc.). This system saves cooling energy costs by using thermally-aware job placements without trading off performance.

Sarood et al. [41] designed a runtime system that does temperature-aware load balancing in datacenters using DVFS and task migration. They also discussed how hotspots could be avoided in datacenters, and showed that cooling costs can be reduced by up to 48 % with temperature-aware load balancing.

Job Scheduling

The problem of scheduling jobs has been extensively studied. Most resource managers implement the first come first served (FCFS) policy as a simple but fair

strategy for scheduling jobs. However, FCFS suffers from low system utilization. A common optimization is *backfilling* [17, 31, 33]. Backfilling improves system utilization by executing jobs with small resource requests out of order on idle nodes.

Fan et al. [16] discussed power-aware job scheduling in the datacenter domain. They discussed a power monitoring system that could use power capping (based on a power estimation method such as RAPL or direct power sensing) and a power throttling mechanism. Such a system works well when there is a set of jobs with loose service level guarantees or low priority that can be forced to reduce consumption when the datacenter approaches the power cap value. Etinski et al. [12–15] explored scheduling under a power budget in supercomputing and analyzed bounded slowdown of jobs. In their series of papers, they introduced three policies. Their first policy looks at current system utilization and uses DVFS during job launch time to meet a power bound. Their second policy meets a bounded slowdown condition without exceeding a job-level power budget. Their third policy improves upon the former by analyzing job wait times and adding a reservation condition.

There are many use cases in a grid computing environment that require Quality of Service guarantees in terms of guaranteed response time, including time-critical tasks that must meet a deadline. Foster et al. [19, 20] proposed *advance reservations* to achieve time guarantees. Advance reservation is a guarantee for the availability of a certain amount of resources to users and applications at specific times in the future. The advance reservation feature requires scheduling systems to support reservation capabilities in addition to backfilling-based batch scheduling. Modern resource management systems such as Sun Grid Engine, PBS, OpenPBS, Torque, SLURM, Maui, and Moab support advance reservation capabilities.

Load Migration

Chiu et al. [7] discussed an electrical grid balancing problem that was experienced in the Pacific Northwest. In order to match electricity supply and balance the electrical grid, they proposed low-cost geographic load migration. They also suggested that a symbiotic relationship between datacenters and electrical grid operators that leads to mutual cost benefits could work well. Ganti et al. [21] looked at two applied cases for distributed datacenters. The

results show that load migration is possible in both homogenous and heterogeneous systems. Their migration strategies were based on a manual process and can benefit from automation.

Datacenter Participation in Smart Grid Programs

Aikema et al. [2] explored the potential for HPC centers to adapt to dynamic electrical prices, to variation in carbon intensity within an electrical grid, and to availability of local renewables. Their simulations demonstrated that 10–50 % of electricity costs could potentially be saved. They also concluded that adapting to the variation in the electrical grid carbon intensity was difficult and that adapting to local renewables could result in significantly higher cost savings.

Power-aware resource management without degrading utilization has been proposed as a DR strategy to reduce electricity costs [44, 45]. The novelty of the proposed job scheduling mechanism is its ability to take the variation in electricity price (dynamic pricing) into consideration as a means to make better decisions about job start times. Experiments on an IBM Blue Gene/P and a cluster system, as well as a case study on Argonne's 48-rack IBM Blue Gene/Q system have demonstrated the effectiveness of this scheduling approach. Preliminary results show a 23 % reduction in the cost of electricity for HPC systems.

Chen et al. [6] studied the potential of datacenter participation in the demand side regulation services. They proposed a dynamic control policy that modulates the datacenter power consumption in response to independent service operator (ISO) requests by leveraging server power capping techniques and various server power states. Results show that datacenters can decrease their energy costs by around 50 % by providing regulation service reserves, without a major deterioration in the quality of service.

Liu et al. [32] introduced a way to reduce cost and environmental impacts using a holistic approach that integrates energy and cooling supply control with IT workload planning to improve the overall attainability of datacenter operations. The results demonstrated a reduction of the recurring power costs and the use of non-renewable energy by as much as 60 % compared to existing techniques, while still meeting service level agreements.

Aikema et al. [3] also analyzed a number of different potential advanced power markets for datacenters to participate in and showed energy cost reductions of up to 12 % with only a small impact on the quality of service provided to users. Ghamkhari et al. [23] built an analytical profit model to show that datacenters can noticeably increase their profit by participating in voluntary load reduction to offer ancillary services and help the grid achieve better service quality and reliability.

Survey Results

We used a questionnaire to understand the current experiences of interaction between SCs and their ESPs. We restricted the analysis to sites in the United States because the results of the survey and practices of DR are highly correlated and driven by energy policies in the country [43].

Nineteen Top100 List sized sites in the United States were targeted for the questionnaire. Eleven sites responded – Oak Ridge National Laboratory (ORNL), Lawrence Livermore National Laboratory, Argonne National Laboratory (ANL), Los Alamos National Laboratory (LANL), Lawrence Berkeley National Laboratory (LBNL), Wright Patterson Air Force Base, the National Oceanic Atmospheric Administration (NOAA), the National Center for Supercomputing Applications (NSCA), San Diego Supercomputing Center (SDSC), Purdue University, and Intel Corporation. The questionnaire was sent to a sample that was not randomly selected. It was sent to those sites where it was relatively easy to identify an individual based on membership within the EE HPC WG. The sample is more representative of Top50 sized sites (one Top50 sized site was not in the sample and 60 % (9/15) of the sample responded). Only 4 additional sites were sampled from the Top51-Top100 List and, of those, 2 responded (Intel and NOAA).

The total power load as well as the intra-hour fluctuation of these sites varied significantly (Fig. 1). The total power load includes all computing systems plus ancillary systems such as power delivery and cooling components. There were four sites with a total power load greater than 10 MW, two sites with approximately 5 MW total power load, and five sites with less than 2 MW of the total power load. For those with a total power load greater than 10 MW, the intra-hour fluctuation (maximum variability) varied from less than 3 MW to 8 MW. One



Table 1

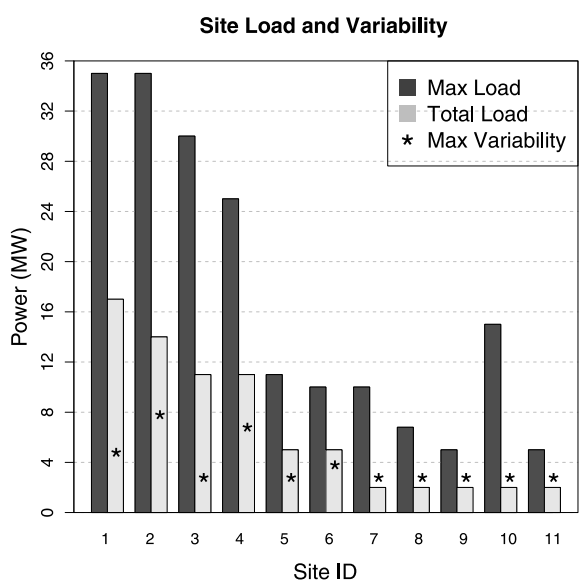


Fig. 1 Site load and variability

of the 5 MW sites said that they experienced 4 MW variability. We chose less than 3 MW intra-hour variability as the bottom of the scale because we assumed that the ESPs would not be affected by 3 MW (or less) fluctuations. The rest of the sites all reported less than 3 MW intra-hour fluctuation. Most of the intra-hour variability was due to preventative maintenance.

For every respondent, the theoretical peak energy or maximum load is approximately twice the total energy, which is indicative of expected future growth in power and energy requirements for SCs. Some of the design parameters that may affect theoretical peak limits are the customer switchgear, transformer, and chiller water capacities. In some cases, there are also limits based on regional ESP capacity constraints.

We asked whether the SCs had talked to their ESPs about the programs and methods used to balance the grid supply and the demand for electricity (see Table 1). About half of them had had some discussion about it, but it was mostly been limited to programs (e. g., peak shed, dynamic pricing) and not methods (e. g., regulation, frequency response, congestion).

Approximately half of the respondents are not currently interested in shedding load during peak demand. LANL reports that the “technical feasibility” and “business case has yet to be developed.” There is slightly more interest in shifting than

Discussions with ESPs

Discussions with ESPs	%Yes
<i>Demand-side programs</i>	
Shedding load during peak demand	54
Responding to pricing incentive programs	45
Shifting load during peak demand	36
<i>Supply-side programs</i>	
Enabling use of renewables	36
Congestion, regulation, frequency response	18
Contributing to electrical grid storage	10

shedding load. SDSC reports that “automatic load shedding is being explored/deployed today” for the entire campus, not just the SC.

Responding to pricing incentive programs is also not considered currently interesting by approximately half of the respondents, although the reasons for this low interest may be organizational. Several open-ended comments revealed that pricing is fixed and/or done by another organization at the site level and is outside of their immediate control.

Only 20 % of the respondents had had discussions with their ESPs about congestion, regulation, and frequency response. LANL is one of the two who have had discussions and who commented that they are “learning about the process” and that it is “outside of [their] visibility or control”.

There were many more respondents who had had discussions with their ESPs about enabling the use of renewables; 36 % had already had discussions and more than half were interested in further and/or future discussions. SDSC already has a site-wide program: “the campus has a large fuel cell (2.5+ MW) and works with the utility with renewables.” Other responses suggest that the interest is at the site level and not unique to the SC.

An open-ended question was posed as to whether or not there was information either requested from the SCs by their ESPs or, conversely, requested from the ESPs by the SCs. In both cases, well over 75 % of the respondents answered “no”. LLNL and LANL were the exceptions. LLNL is “responding to requests for additional data on an hourly, weekly, and monthly basis.” They are also working to develop an automated capability to share data with their ESPs, which would provide automated additional detailed forecasting and ultimately real time data. LANL has also been requested to pro-



Table 2

HPC strategies for responding to electricity provider requests (listed from highest to lowest interest + impact)	% Interested	% High Impact	% Medium Impact
Coarse grained power management	64	46	27
Facility shutdown	36	64	10
Job scheduling	36	27	18
Load migration	10	36	18
Re-scheduling back-ups	45	0	10
Fine-grained power management	27	0	36
Temperature control beyond ASHRAE limits	27	0	18
Turn off lighting	18	0	0
Use back-up resources (e. g., generators)	0	10	27

vide average “power projections, hour by hour, for at least a day in advance.” Additionally, LANL has asked their ESP for more information on “sensitivity of power distribution grid to rapid transients (random daily step changes of 10 MW up or down within a single AC cycle).”

Given the low levels of current engagement between the ESPs and the SCs, it is not surprising that none of the SCs are currently using any power management strategies to respond to grid requests by their ESPs. SDSC’s supercomputer center is not an exception, but they did respond that their entire “campus is leveraging parallel electrical distribution to trigger diesel generators and other back-up resources to respond to grid and non-grid requests.”

It was suggested by ORNL that some of the power management strategies are of questionable business value even for energy efficiency, let alone grid integration. For example, ORNL comments that “these assets have very clear depreciation schedules, and the modest cost savings in terms of electricity consumption due to some of these methods may not (or frequently will not) outweigh the capital investment cost in the computer. That is, if a site spent \$ 100M for a computer that will remain in production for 60 months, then the apparent benefit of power capping, etc., could easily be outweighed by lost productivity of the consumable resource.”

Similarly, another comment by ORNL suggested that the rapid deployment of hardware features, like P-states, may outpace the need for strategies like power-aware job scheduling.

We tried to evaluate whether power management strategies will be considered relevant and effective for grid integration at some point in the future. Two questions were asked: is there interest in using the strategies and what impact did they think that the strategies would have? When combining interest and impact, the results showed that power capping, shutdown, and job scheduling were both potentially interesting and of high impact (see Table 2).

Load migration, back-up scheduling, fine-grained power management, and thermal management were of medium interest and impact. Lighting control and back-up resources were of low interest and impact.

Temperature control and lighting management are utilized as strategies but are considered to be of medium to low interest and impact for responding to requests from ESPs. The infrastructure energy efficiency of the responding supercomputer sites is high, as reflected in the reported PUE (Fig. 2). Two sites reported a PUE below 1.25; the majority were between 1.25 and 1.5 and the highest was 1.53. Approximately half of the respondents said that they used temperature control and lighting management as strategies, but not for grid requests. Temperature control and lighting management are well documented and well understood strategies for improving energy efficiency, so it is not surprising that sites with PUEs below 1.5 are using them.

NOAA comments that their “lights automatically shut off 24 × 7 when there is no motion in the data-center.” There is a value in lighting control for energy efficiency purposes, as demonstrated by its having

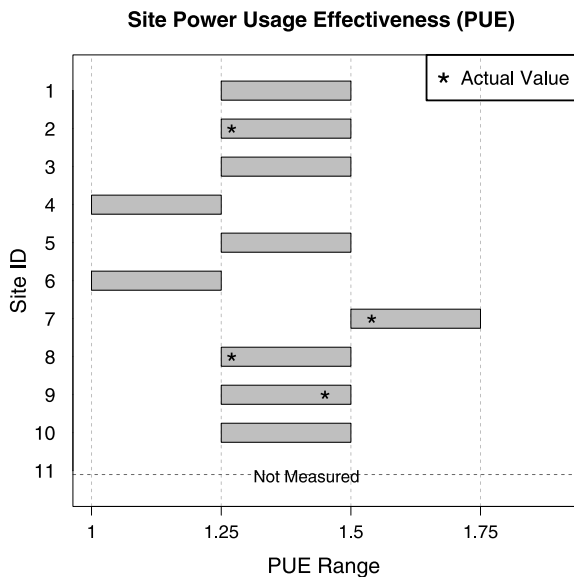


Fig. 2 Site power usage effectiveness

been fully implemented. NOAA also comments that the impact of further lighting control “is so small compared to the HPC demand load that” they would “be surprised if the utility is interested.”

LLNL reports that they “took 3 years to raise the temperature in their center by 18° F. It was done in conjunction with a failure rate analysis of the systems, as well as a measurement of the electrical savings prior to moving to the next set point.” LLNL is currently operating in the ASHRAE recommended range, but expresses concerns with increasing temperature as a grid-integration response. The concerns include hardware failures, tape storage read/write errors, and compromising dew point requirements where liquid and air-cooling are co-located.

Distinguishing interest from impact sheds further insight; some strategies are considered high impact, but not interesting enough to consider deployment. Facility shutdown is rated as having a high impact, but only considered interesting by 36 % of the respondents. NOAA commented that,

“We’ve had too many HPC instability and equipment failures to utilize this as a strategy.” This divide is even more apparent with load migration. It is rated as having a high impact by 36 % of the respondents, but is only interesting to 10 %.

Opportunities, Solutions, and Barriers

The responses to the questionnaire presented in Sect. “Survey Results” represent a variety of de-

sires and experience regarding interactions between SCs and ESPs. For example, the responses from the two SCs with the largest power draws, LLNL and ORNL, diverge in several areas. This divergence is perhaps primarily due to characteristics of their respective ESPs. In contrast, SDSC stands out as a leader in integrating with their ESP on a site-wide level. To that end, the responses from SDSC may exemplify some of the opportunities available to other SCs that are willing to pursue this degree of integration.

The responses to the questionnaire also suggest that some ESPs request that their SC customers develop capabilities for informing the provider of expected periods of exceptional power consumption and for responding to requests from the provider to consume less power for specified periods of time. Upon initial consideration, this idea might seem to run counter to the primary mission objective of most SCs of delivering as many uninterrupted computational cycles as possible to their users. In some extreme cases, SCs may not have a choice in the matter as the size and energy requirements of supercomputers increase; indeed, some ESPs may require large centers to develop a DR capability. However, a direct business case may exist to encourage SCs to develop this negotiation capability on their own. For example, if ESPs were to offer electricity at a significantly reduced rate on the condition that the SC customer develop DR capabilities, the long-term cost savings to the center could make undertaking such a project worthwhile.

Perhaps one of the most straightforward ways that SCs can begin the process of developing a DR capability is by enhancing existing system software used for managing computing resources within the center. Indeed, the questionnaire responses from Sect. “Survey Results”, as well as the literature review presented in Sect. “Prior Work”, both strongly support the idea that the greatest opportunities for SCs to develop integration capabilities are related to system software. Specifically, and presented in approximate order of decreasing interest and expected impact to the questionnaire respondents, system software in this context consists of coarse-grained power management (such as uniform processor power capping across the cluster), job scheduling, load migration, rescheduling backups, and fine-grained power management (such as dynamic, per-processor power capping).

Coarse-grained power capping may be one of the most straightforward methods of power management. In the simplest form, this technique may entail human intervention to adjust computing resources so that they operate at a reduced capacity or to entirely shut down some of the computing capacity of an SC. By attenuating resources, the SC manager can ensure that power consumption stays below some defined level. This defined level may be a pre-arranged power cap negotiated between the SC and the ESP and maintained on an ongoing basis, or, perhaps more likely, a power draw level that is requested by the ESP to handle unanticipated loads somewhere else in the ESP's system. Note that the savings in power may not need to come entirely from attenuating computing resources. Rather, reducing power consumption in computing resources is likely to result in a corresponding reduction in thermal load within the SC, which may allow significant power savings in the cooling system as well.

The coarse-grained power capping technique described above assumes that the SC environment has some amount of instrumentation and metering that allows for the collection of power telemetry data. This telemetry is necessary for the SC facility manager in order to understand how the power supplied by the ESP is distributed to resources within the center. Further, this telemetry is likely to be important to automated solutions for power management, such as the job scheduling techniques described below. In light of the fact that many system integrators such as Cray and IBM are now delivering supercomputing systems that include telemetry capabilities, the assumption that this information is available seems acceptable. According to the responses to the questionnaire presented in the previous section, SC facility managers perceive this accounting data as distinct from per-user or per-job accounting data that is typically collected and indicate that this data should be retained for electricity provisioning planning purposes.

Techniques that involve job scheduling may offer more automated approaches to power management. Due to the unique role that the job scheduler and resource manager play within an SC, these techniques may involve adjusting either the workflow of jobs within the center or characteristics of the computational resources within the center.

On the one hand, the job scheduler has knowledge of and control over the upcoming workflow

within the SC simply by examining and manipulating the job queue. One easily-accessible technique is for a human operator to use capabilities such as advanced reservations to reserve pre-arranged blocks of time in which jobs with high power loads will run. These blocks of time could be negotiated with the ESP on an ongoing basis or could be in response to on-demand requests made by the ESP. Even more automatic techniques are possible if the job scheduler is given enough information about the workflow to make intelligent decisions about job scheduling. For example, jobs may be submitted with various metadata that enable the job scheduler to understand the characteristics of each job such as priority, the relative importance of a job compared to other jobs, and urgency, the rate at which the value of a job decreases as time elapses. These characteristics are not only important to a job scheduler for ensuring efficient utilization of an SC's resources under traditional circumstances, but they are also a vital piece of successfully implementing a DR capability for at least two reasons. First, they provide a set of metrics by which the SC can estimate the cost in terms of the "lost opportunity" of responding to an ESP's request to run with attenuated resources. Second, they allow the SC to prioritize jobs in the queued workflow in order to understand how to best utilize computational resources. This capability is important under normal circumstances, but becomes even more essential in a DR scenario.

At a lower level, schedulers and runtime systems can exercise fine-grained, dynamic DR capability. For example, the job scheduler knows which nodes within a supercomputer are occupied with running jobs or are expected to become occupied in the near future. To that end, the job scheduler can use its control over the resource management process to place idle nodes into a sleep state in which they draw significantly reduced power. This strategy is especially effective in supercomputing environments containing at least some resources that are used at irregular intervals, allowing opportunities to utilize sleep states effectively during periods when the resources are idle.

In environments where all computing resources are heavily utilized, fine-grained power scheduling will be directed by the runtime system. For example, in the presence of load imbalance within a job, traditional applications may rely on periodically moving data around the allocated nodes to ensure all pro-

processors are performing a roughly equal amount of work. This load-balancing process is both time and energy-intensive. By relocating power instead of data, processors with lighter loads can surrender power and run more slowly, allowing more heavily-loaded processors to use additional power to run faster. Combining both techniques should lead to improved execution time as well as more efficient power utilization.

Even more interesting scenarios are possible in cases where the job scheduler combines its knowledge of the upcoming queued workflow with its knowledge and control over the computational resources within the SC. These scenarios are most appropriate when the supercomputing scenario contains a pervasively heterogeneous mix of computational resources. For example, many contemporary SCs contain several different types of compute nodes with various types of processors and accelerator cards. In some circumstances, the job scheduler may be able to choose which resource to use for running a given job among several candidate resources. The tradeoff here is not only in terms of the time necessary to complete the job (that is, different resources could potentially complete the job in very different amounts of time) but also in terms of the energy consumed in completing the job (that is, different resources could potentially consume very different amounts of energy in completing the job). Further, other resources such as memory access patterns, disk access patterns, and network use affect the energy signature of a job and may be observed by the scheduler. By maintaining a database of job-to-resource mappings that record the time and energy taken for each job, the scheduler can, over time, improve its ability to decide which jobs have the highest affinity to each type of resource. Using this knowledge to optimize an SC's workflow in terms of job throughput or energy consumption is admittedly complex, but the potential rewards are likely to be compelling both to the day-to-day operation of the center and to DR capabilities.

Opportunities may also exist for SCs to cooperate with each other in scenarios in which computational loads are migrated from one site to another where energy costs are less expensive. This scenario is challenging for both technical and business reasons. Technical challenges include issues such as user authentication and authorization (i. e., a user may be authorized to use resources at one site but not at another site) and data movement (i. e., it

may be infeasible to migrate large datasets from one site to another site). To some extent, some of these technical challenges may be mitigated by the use of advanced reservation capabilities in the scheduling systems at each site, allowing resources to be simultaneously reserved while large datasets are properly staged. Business challenges include the notion that an SC currently has little incentive to migrate jobs to another "competing" center. Indeed, the questionnaire results reflect low interest in load migration strategies. It seems likely that in order to be a feasible scenario, the structure of payment and rewards to an SC to cooperate with other centers would need to be structured differently than they are currently.

In a very broad sense, DR techniques such as job scheduling, power capping, and load migration can be considered to be coarse-grained approaches because they involve considering "big picture" views of the workload and computational resources in an SC. According to the questionnaire results presented in the previous section, facilities managers view these approaches as the most likely candidates for creating effective DR capabilities.

Finally, this section has focused heavily on the opportunities available to SCs that come from developing DR capabilities. This notion is primarily due to the fact that the questionnaire presented in Sect. "Survey Results" was distributed to SCs in the United States, not to ESPs. That said, opportunities do exist for ESPs that develop DR capabilities. At one level, the negotiation process itself requires integration in terms of the communication and messaging protocols that are necessary. To that end, opportunities exist for adapting and extending existing standards currently used within the industry, thus creating new use cases and capabilities for ESPs. At a higher level, ESPs will most likely need to improve their ability to determine in near real time the important places within the electrical grid where demands exceed supply. Determining this is likely to be a complex optimization problem. While this section focuses on solving these problems to the end of developing a DR strategy in conjunction with SCs, these capabilities are likely applicable to a broad range of customers.

Conclusions and Next Steps

This paper explores the possibility of a new relationship between ESPs and SCs with increased communication and engagement from both parties.

Because SCs have an increasingly large and fluctuating power demand, they challenge their providers to supply a reliable source of electricity. ESPs are interested in partnering with customers, like SCs, to create a more dynamic and resilient grid by obtaining predictable demand forecasts and engaging in programs like DR.

We focused our attention on the largest SCs in the United States. The two SCs with the largest electricity demand, ORNL and LLNL, have had very different experiences. ORNL's experience is that its electricity demand and fluctuations are not significant factors for their ESP. LLNL's experience is opposite to that of ORNL. Because of large swings in power usage, the LLNL SC was approached by their ESP with a request for daily predictable demand forecasts. That request began an ongoing relationship.

The LANL SC's experience is similar to that of LLNL. SDSC has an even tighter relationship with their ESP, but this relationship involves the entire campus and not just the SC.

As previous research with datacenters has shown [25], SCs can serve as resources to the grid. To enable this, automation technologies and data communication standards, which can link the SCs with the electric grid as well as on-site power management strategies will play a key role to ease adoption and lower the participation costs. Power capping, shutdown, and job scheduling are identified as the most interesting management strategies with the highest leverage for responding to requests from ESPs.

Nonetheless, the business case for the grid integration of SCs remains to be demonstrated. SCs have concerns that deploying these strategies might have an adverse impact on their primary mission. One of the key enablers for SCs to participate in electricity markets (for example, DR, electricity prices) is having markets that value their participation. In other areas like commercial buildings and select industrial facilities, benefits to both ESPs and customers are well documented. However, as the electrical grid and new dynamic loads such as SCs evolve, the markets need mechanisms to identify and provide value of participation (for example, cost, energy, carbon).

We are planning to pursue several areas in our future work. We are planning a similar survey for Europe to explore whether there is a more compelling business case in other geographies. We

expect the business value of such grid integration to be enhanced where the price of electricity is expensive, or where the supply is constrained, or varies dynamically.

We plan on following up with the ESPs that support these US-based SCs. We note that this work's focus was from the perspective of the SC, and we are interested in hearing from the ESPs about what makes a customer more or less interesting or challenging with respect to grid integration.

Finally, we want to better understand the specific information that could be exchanged between SCs and ESPs.

Additional Authors

Lou Ahlen, Purdue University
Anna Maria Bailey, LLNL
Susan Coghlan, Argonne NL
Bob Conroy, OSISOFT
Ayse K. Coskun, Boston University
Thomas E. Durbin, NCSA
Lucio Grandinetti, UNICAL
Ted Kubaska, IEEE
Jim Rogers, ORNL
Dale Sartor, LBNL
Darren Smith, NOAA

Appendix

For the purposes of this paper, this appendix contains a summary of the questionnaire. The questionnaire is divided into the following three sections:

- Facility energy. The total facility energy and the total HPC load should be the same number that you use when calculating PUE, as defined by the Green Grid Whitepaper #49.
- Management and control. Please answer whether or not you employ any of the strategies described below for managing and controlling total facility energy in response to a request from your electrical/utility provider. You may use some of these same strategies for improving energy efficiency. Answer "yes" only when the strategy is used at least in part for grid response. Answer "yes" only when the strategy is used at least in part for grid response. Answer "no" if the strategy is only used for improving energy efficiency.
- Electrical/utility provider information. Answers to these questions help us understand the nature

of any relationship you might have between your HPC facility and your site's electric utility/provider. Please answer "yes" if you have had any communication about the following programs and methods with your site's electric/utility provider. For each program and/or method for which there has been communication, please describe the nature of that communication in the comments.

Facility Energy

1. What is your total facility energy?
2. What is your total HPC load?
3. What is your facility PUE?
4. What is your facility's theoretical peak energy, as the infrastructure is currently fit up.
5. What is the maximum variation in total facility energy that is likely to re-occur?
6. How often does this variation occur?
7. If there is any regular pattern to this variation, please describe the circumstances. Include the reason for the variation, the magnitude and duration if possible. For example, "There is a 5MW drop every two weeks for a 6 h period during preventative maintenance periods."

Management and Control

8. COARSE-GRAINED POWER MANAGEMENT: manage power for the HPC system or subsystem (could include storage, networking as well as compute sub-systems). Example: power capping.
9. FINE-GRAINED POWER MANAGEMENT: intelligent built-in power management. Examples: voltage and frequency governors, hibernation.
10. LOAD MIGRATION: shift computing loads to a different electrical grid.
11. JOB SCHEDULING: job shifting or queuing (scheduling) has historically been used as a strategy for managing CPU utilization, but could also be used to manage the energy utilization of IT equipment.
12. BACK-UP SCHEDULING: defer data storage processes to off-peak periods
13. SHUTDOWN: graceful shutdown of idle HPC equipment loads. Usually applies when there is redundancy
14. LIGHTING CONTROL: with advance warning, datacenter lights could be shutdown completely.

15. TEMPERATURE ADJUSTMENT: widen acceptable (ASHRAE thermal conditions) temperature set point ranges and humidity levels for short periods.
16. BACK-UP RESOURCES: using generators and other electrical storage devices.
17. Are there any other strategies that you use to manage and control your total facility energy in response to a request from your energy/utility provider? Please describe.
18. Please evaluate as high, medium or low the MW impact of each of these strategies as a response to a grid request.
 - Power capping
 - Load migrations
 - Temperature adjustments - clock speeds
 - Lighting control
 - Job scheduling
 - Back-up scheduling
 - Idle management
 - Shutdown
 - Back-up resources

Electrical/Utility Provider Information

19. PEAK SHEDDING: utility provider arrangements used to reduce peak load, where the reduced load is not shifted to another time.
20. PEAK SHIFTING: utility provider arrangements where the load during peak times is moved, typically to non-peak hours.
21. DYNAMIC PRICING: Time-varying pricing arrangements used to increase, shed, or shift electricity consumption. There are two types of pricing, peak and real time. Peak pricing is pre-scheduled; however, the consumer does not know if a certain day will be a peak or a non-peak day until day-ahead or day-of. Real time pricing is not pre-scheduled; prices can be set day-ahead or day-of.
22. GRID SCALE STORAGE: methods used to store electricity on a large scale. Pumped-storage hydroelectricity is the largest-capacity form of grid energy storage.
23. RENEWABLES: variability in the electric power generation from renewable resources and the methods used to respond to that variability.
24. FREQUENCY RESPONSE: methods used to keep grid frequency constant and in balance. Generators are typically used for frequency

response, but any appliance that operates to a duty cycle (such as air conditioners and heat pumps) could be used to provide a constant and reliable grid balancing service by timing their duty cycles in response to system load.

25. **REGULATION** (up or down): methods used to maintain that portion of electricity generation reserves that is needed to balance generation and demand at all times. Raising supply is up-regulation and lowering supply is down-regulation. There are many types of reserves (e. g., operating, congestion), distinguished by who controls them and what they are used for.
26. **CONGESTION**: methods used to resolve congestion that occurs when there is not enough transmission capability to support all requests for transmission services. Transmission system operators must re-dispatch generation or, in the limit, deny some of these requests to prevent transmission lines from becoming overloaded. Or, methods used to resolve congestion that occurs when the distribution control system is overloaded. It generally results in deliveries that are held up or delayed.
27. Is there any information you would like from your provider that you are not getting? If yes, please describe what you would like to know.
28. Is your provider asking for information from you that you are not able to provide? If yes, please describe what they are asking for.
29. Do you experience any power quality issues at your HPC facility? If yes, please describe it.

References

1. Top500 Supercomputer Sites, <http://www.top500.org/lists/2013/11>.
2. Aikema D, Simmonds R (2011) Electrical Cost Savings and Clean Energy Usage Potential for HPC Workloads. In: 2011 IEEE International Symposium on Sustainable Systems and Technology (ISSST), pp 1–6
3. Aikema D, Simmonds R, Zareipour H (2012) Data Centres in the Ancillary Services Market. In: 2012 International Green Computing Conference (IGCC), IEEE, pp 1–10
4. ASHRAE (2012) Thermal Guidelines for Data Processing Environments, Special Publication, third edition. American Society of Heating, Refrigerating and Air-Conditioning Engineers, Inc., Atlanta, GA
5. Cameron KW, Feng X, Ge R (2005) Performance-Constrained Distributed DVS Scheduling for Scientific Applications on Power-Aware Clusters. In: Supercomputing, Seattle, Washington
6. Chen H, Caramanis MC, Coskun AK (2014) The Data Center as a Grid Load Stabilizer. In: 19th Asia and South Pacific Design Automation Conference (ASP-DAC), IEEE, pp 105–112
7. Chiu D, Stewart C, McManus B (2012) Electric Grid Balancing Through Lowcost Workload Migration. SIGMETRICS Perform Eval Rev 40(3):48–52
8. Curtis-Maury M, Jevic FB, Antonopoulos C, Nikolopoulos DS (2008) Prediction-Based Power-Performance Adaptation of Multithreaded Scientific Codes. IEEE Trans Parallel Distrib Syst 19(10):1396–1410
9. Curtis-Maury M, Dzierwa J, Antonopoulos CD, Nikolopoulos DS (2006) Online Power-Performance Adaptation of Multithreaded Programs Using Hardware Event-Based Prediction. In: International Conference on Supercomputing, New York, NY, USA, ACM
10. Curtis-Maury M, Shah A, Jevic FB, Nikolopoulos DS, de Supinski BR, Schulz M (2008) Prediction Models for Multi-Dimensional Power-Performance Optimization on Many Cores. In: International Conference on Parallel Architectures and Compilation techniques, New York, NY, USA, ACM
11. David H, Gorbatov E, Hanebutte UR, Khanna R, Le C (2010) RAPL: Memory Power Estimation and Capping. In: Proceedings of the 16th ACM/IEEE International Symposium on Low Power Electronics and Design, ISLPED '10, New York, NY, USA, ACM, pp 189–194
12. Etinski M, Corbalan J, Labarta J, Valero M (2012) Parallel Job Scheduling for Power Constrained HPC Systems. Parallel Comput 38(12):615–630
13. Etinski M, Corbalan J, Labarta J, Valero M (2010) Optimizing Job Performance Under a Given Power Constraint in HPC Centers. In: Green Computing Conference, pp 257–267
14. Etinski M, Corbalan J, Labarta J, Valero M (2010) Utilization Driven Power-Aware Parallel Job Scheduling. Comput Sci R&D 25(3-4):207–216
15. Etinski M, Corbalan J, Labarta J, Valero M (2011) Linear Programming Based Parallel Job Scheduling for Power Constrained Systems. In: International Conference on High Performance Computing and Simulation, pp 72–80
16. Fan X, Weber W-D, Barroso LA (2007) Power Provisioning for a Warehouse-Sized Computer. In: The 34th ACM International Symposium on Computer Architecture
17. Feitelson DG, Schwiegelshohn U, Rudolph L (2004) Parallel Job Scheduling – a Status Report. In: Lecture Notes in Computer Science, Springer-Verlag, pp 1–16
18. Femal ME, Freeh VW (2005) Safe Overprovisioning: Using Power Limits to Increase Aggregate Throughput. In: International Conference on Power-Aware Computer Systems
19. Foster I (2001) The Anatomy of the Grid: Enabling Scalable Virtual Organizations. In: Proceedings of the First IEEE/ACM International Symposium on Cluster Computing and the Grid, pp 6–7
20. Foster I, Kesselman C, Lee C, Lindell B, Nahrstedt K, and Roy A (1999) A Distributed Resource Management Architecture That Supports Advance Reservations and Co-Allocation. In: Seventh International Workshop on Quality of Service, IWQoS '99, pp 27–36
21. Ganti V, Ghatikar G (2012) Smart Grid as a Driver for Energy-Intensive Industries: A Data Center Case Study. In: Grid-Interop 2012
22. Ge R, Feng X, Feng W, Cameron KW (2007) CPU Miser: A Performance-Directed, Run-Time System for Power-Aware Clusters. In: International Conference on Parallel Processing, Xi'An, China
23. Ghamkhari M, Mohsenian-Rad H (2012) Data Centers to Offer Ancillary Services. In: IEEE Third International Conference on Smart Grid Communications (Smart-GridComm), pp 436–441, IEEE
24. Ghatikar G, Ganti V, Matson N, Piette MA (2012) Demand Response Opportunities and Enabling Technologies for Data Centers: Findings From Field Studies. In: PG&E/SDG&E/CEC/LBNL
25. Ghatikar G, Riess D, Piette MA (2014) Analysis of Open Automated Demand Response Deployments in California and Guidelines to Transition to Industry Standards. Technical Report LBNL-6560E, Lawrence Berkeley National Laboratory, 1 Cyclotron Rd, Berkeley, CA 94720
26. Hsu C-H, Feng W-C (2005) A Power-Aware Run-Time System for High-Performance Computing. In: Supercomputing
27. Intel (2011) Intel-64 and IA-32 Architectures Software Developer's Manual, Volumes 3A and 3B: System Programming Guide
28. Kappiah N, Freeh VW, Lowenthal DK, Pan F (2005) Exploiting Slack Time in Power-Aware, High-Performance Programs. In: Supercomputing
29. Kaushik RT, Nahrstedt K (2012) T*: A Data-Centric Cooling Energy Costs Reduction Approach for Big Data Analytics Cloud, SC'12, Los Alamitos, CA, USA, IEEE Computer Society Press, Article no. 52
30. Li J, Martínez JF (2006) Dynamic Power-Performance Adaptation of Parallel Computation on Chip Multiprocessors. In: 12th International Symposium on High-Performance Computer Architecture, Austin, Texas
31. LiKa DA (1995) The ANL/IBM SP Scheduling System. In: Job Scheduling Strategies for Parallel Processing, Springer-Verlag, pp 295–303
32. Liu Z, Chen Y, Bash C, Wierman A, Gmach D, Wang Z, Marwah M, Hyser C (2012) Renewable and Cooling Aware Workload Management for Sustainable Data Centers. ACM SIGMETRICS Perform Eval Rev 40:175–186
33. Mu'alem AW, Feitelson DG (2001) Utilization, Predictability, Workloads, and User Runtime Estimates in Scheduling the IBM SP2 with Backfilling. IEEE Trans Parallel Distrib Syst 12(6):329–543
34. Niccolai J (2013) New Data Center Survey Shows Mediocre Results for Energy Efficiency. IT World IDG News Service, www.itworld.com/article/2709163/hardware/new-data-center-survey-shows-mediocre-results-for-energy-efficiency.html

35. Palensky P, Dietrich D (2011) Demand Side Management: Demand Response, Intelligent Energy Systems, and Smart Loads. *IEEE T Ind Inform* 7(3):381–388
36. Patki T, Lowenthal DK, Rountree B, Schulz M, de Supinski BR (2013) Exploring Hardware Overprovisioning in Power-Constrained, High Performance Computing. In: *International Conference on Supercomputing*, pp 173–182
37. Porterfield AK, Olivier SL, Bhalachandra S, and Prins JF (2013) Power Measurement and Concurrency Throttling for Energy Reduction in OpenMP Programs. In: *27th International Parallel and Distributed Processing Symposium Workshops & PhD Forum (IPDPSW)*, IEEE, pp 884–891
38. Rountree B, Lowenthal D, de Supinski BR, Schulz M, Freeh V, Bletch T (2009) Adagio: Making DVS Practical for Complex HPC Applications. In: *International Conference on Supercomputing*
39. Rountree B, Ahn DH, de Supinski BR, Lowenthal DK, Schulz M (2012) Beyond DVFS: A First Look at Performance Under a Hardware-Enforced Power Bound. In: *IPDPS Workshops*, IEEE Computer Society, pp 947–953
40. Rountree B, Lowenthal DK, Funk S, Freeh VW, de Supinski B, Schulz M (2007) Bounding Energy Consumption in Large-Scale MPI Programs. In: *Supercomputing*
41. Sarood O, Kalé LV (2011) A “Cool” Load Balancer for Parallel Applications. In: *Proceedings of the 2011 ACM/IEEE Conference on Supercomputing*, Seattle, WA
42. Springer R, Lowenthal DK, Rountree B, Freeh VW (2006) Minimizing Execution Time in MPI Programs on an Energy-Constrained, Power-Scalable Cluster. In: *ACM Symposium on Principles and Practice of Parallel Programming*
43. Torriti J, Hassan MG, Leach M (2010) Demand Response Experience in Europe: Policies, Programmes and Implementation. *Energy* 35(4):1575–1583
44. Yang X, Zhou Z, Wallace S, Lan Z, Tang W, Coghlan S, Papka ME (2013) Integrating Dynamic Pricing of Electricity Into Energy Aware Scheduling for HPC Systems. In: *Proceedings of SC13: International Conference for High Performance Computing, Networking, Storage and Analysis, SC '13*, New York, NY, USA, ACM
45. Zhou Z, Lan Z, Tang W, Desai N (2013) Reducing Energy Costs for IBM Blue Gene/P via Power-Aware Job Scheduling. *IPDPS Workshop*, IEEE Computer Society