

# Exploring Power-Performance-Quality trade-offs for Exascale Combustion Simulation

Authors Name/s per 1st Affiliation (Author)  
line 1 (of Affiliation): dept. name of organization  
line 2: name of organization, acronyms acceptable  
line 3: City, Country  
line 4: Email: name@xyz.com

Authors Name/s per 2nd Affiliation (Author)  
line 1 (of Affiliation): dept. name of organization  
line 2: name of organization, acronyms acceptable  
line 3: City, Country  
line 4: Email: name@xyz.com

**Abstract**—The computational demand of high-performance computing (HPC) applications has brought major changes to the HPC system architecture. As a result, it is now possible to run simulations faster and get more accurate results. But behind this, power and energy are becoming critical concerns for HPC systems, e.g. Tianhe-2 has reached speed 33.86 PFLOPS at power of 17.6 MW [1], which cost electric around 30 million per year [2]. U.S. Department of Energy (DOE) has set the goal to achieve exascale performance with a power budget of 20MW [3], this make power efficiency become one of the critical challenges for the exascale research.

Current research efforts have studied power and performance tradeoffs, and how to balance these, e.g., using DVFS to meet power constraints, which significantly impacts performance. However, scientific applications may not tolerate degradation in performance and other tradeoffs need to be explored to meet power budgets, e.g., involving the application in making energy-performance tradeoff decisions.

This research focuses on studying the properties and exploring the performance and power\energy tradeoffs of Adaptive Mesh Refinement (AMR) based simulation applications. Our experimental evaluation provides an empirical evaluation of different application configurations that gives insights into the power-performance tradeoffs space for those AMR applications. The key contribution of this work is a better understanding of the running behavior of these AMR-based applications and the power-performance tradeoffs for these applications under power constraints, which can be used to better schedule power budgets across HPC systems. [4]

**Keywords**-AMR; Energy; Performance; Quality; Power scheduling;

## I. INTRODUCTION

High performance computing (HPC) has been played an important role in the field of computational science. From design perspective, HPC were built to maximize performance while irrespective of power and energy consumption, though their energy consumption has already occupied a large part of cost (i.e., operation cost). However, as we are approaching the exascale era, power is turning from an optimization goal to a critical operation constraint. U.S Department of Energy (DOE) has currently set a bound of 20MW for an exascale system. [5] This strict power constraint poses a hard research challenge with current

hardware and software. Sunway TaihuLight, the top one supercomputer as of June 2016, has a peak performance of 93.0 PetaFLOPS at 15.4 MW, which is 6.04 GigaFLOPS per watt. Sunway's power efficiency has improved three times then Tianhe-2. However, achieving the goal of exascale computing at 20 MW, it requires 50 GigaFLOPS per watt. So that, current HPC system still need at least a 8 times power efficiency improvement towards exascale.

In order to achieve this exascale system power constraint, current research efforts have studied power and performance tradeoffs, and how to balance these. [6]–[11] Many power management strategies have been proposed [12]–[16], but most of the works are tend to choose a performance and then constraint the power consumption under that performance loss. Even if this can meet power constraints, it significantly impacts performance. In fact, most scientific applications may not tolerate degradation in performance, other tradeoffs need to be explored to meet power budgets.

Many research efforts have studied power and performance tradeoffs, and most energy models or strategies are based on runtime (e.g., leveraging MPI slack) for power clamping or power capping techniques, like Dynamic Voltage and Frequency Scaling (DVFS) to constraint the power. However, power and performance are in the two sides of a balance scale, that it is hard to improve one side without scarifying the other one. Therefore, one of the key problems addressed in this research is keeping the power bound (or budget) without losing performance, which is challenging for real world applications targeting exascale. At the same time, it is clear that future HPC system's whole-system power constraint that will be filtered down to job-level power constraint i.e., power budgets need to be managed at application or workflow level. This indicates, as well as we believe that the applications should be involved in making tradeoffs decisions.

This work targeting on AMR-based simulation applications. AMR method considers a hierarchy of grids of differing resolutions ranging from the coarsest to the finest. It can focus computational resources in regions of interest while decrease computing resolution in regions with

less interest. Less resolution means lighter workload and less power consumption. Therefore, this flexible resolution property gives us a potential opportunity to extract power budget for other usage. In order to take best advantage of it, this work first studies the mechanisms and policies to control AMR properties, and then, it characterizes LMC power performance running on different number of cores with different levels of resolution (e.g., levels of refinement in AMR). Finally, the characterization is complemented with power capping techniques (i.e., RAPL). The overarching goal of this work is to understand the tradeoffs between power-performance and quality, and building models taking into account AMR properties for managing power budgets and workflow at scale.

The contributions of this work are summarized below:

- 1) It presents an empirical evaluation of different configurations of application that gives insights into the energy-performance-quality tradeoffs for scientific data-driven workflows.
- 2) It provides a comprehensive study of this LMC simulation performance, quality, and power and energy behavior.
- 3) It presents a proof-of-concept study of potential of power capping and power management to balance power-performance-quality tradeoffs.

The rest of this thesis is organized as follows. Section II summarizes the related work. Section III presents the proof of concepts. Section IV describes the evaluation methodology and presents the results of the experimental evaluation. Section VI concludes the chapter and outlines ongoing and future research.

## II. RELATED WORKS

Energy-efficiency has become a critical concern for HPC applications. There are many approaches have been proposed to obtaining energy savings during HPC application execution. Some of them are to focus on identifying stalls during the execution by measuring architectural parameters from performance counters as proposed in [6]–[8]. In addition to using performance counters, Rountree et al. [12] developed a runtime system called Adagio, by doing the critical path analysis, it can determine which tasks may be slowed down and also suitable opportunities to apply DVFS to minimize the performance loss in the parallel execution. This achieves significant energy saving in scientific application with negligible performance lose. However, this approach appears beneficial when applications have computation or communication imbalances among participating processes, which is typically not the case for a highly efficient parallel application nor suitable for the LMC program targeted in this work. Some approaches in [9], [13] are proposing to determine the communication phases to apply DVFS. Kandalla et al. [17] give algorithms to save energy in the collectives such as MPI\_Alltoall and MPI\_Bcast. Moreover,

Ioannou et al. [14] describe a runtime system for the Intel Single-chip Cloud Computer (SCC) processor to detect repeatable communication phases followed by an application of frequency scaling. Donofrio et al. [10] have studied energy efficiency for extreme-scale science and developed Green Flash, an application-driven design to improve the kernels computational efficiency. Li et al. [15] developed a new power-aware performance prediction model of hybrid MPI/OpenMP programming model and combine memory and disk management techniques to provide performance guarantees for control algorithms. Rodero et al. [16] explored the potential of application-centric aggressive power management of HPC scientific workloads while considering power management mechanisms and controls available at different levels and different subsystems. Gmall et al. [11] explored data-related energy/performance trade-offs for end-to-end co-design simulation workflows on current and ongoing high-end computing systems. Lively et al. [18] explored and investigated energy consumption and execution time of different parallel implementation of scientific applications on multi-cluster systems.

Most of the work described above are from system level and are based on tolerating performance loss and constrain the power consumption under that performance loss. The work presented here discusses the possibility to filter down power constraint to job-level, and then take advantage of applications specific properties, such as resolution, to keep the power budget while maintaining the level of performance.

The cost of provisioning power in data centers is a very large fraction of the total cost of operating a data center. [19], [20] Therefore the ability to cap peak power consumption is a desirable feature in modern data centers. [21] Many power management approaches have been proposed to provide performance guarantees while constraining power budget. Sartori et al. [22] describe a peak power management technique for multi-core systems by choosing the power state for each core that meets the power constraints. Cebrian et al. [23] develop a power balancing strategy that through borrowing power budgets from cores that consume lower power to dynamically adapts the per-core power budgets. While, Gandhi et al. [24] give a power capping strategy to meet the power budget by inserting idle cycles during execution. This approach aims at controlling the average power consumption, but cannot guarantee the peak power. So a number of other approaches are proposed through reconfiguring hardware to meet the power budget. Meng et al. [25] provide a power management strategy through dynamic reconfiguration of cores by cache resizing. Konotorinis et al. [26] propose a table-driven adaptive core reconfiguration technique that configures core resources such as load-store queues and floating point units to meet peak power budget. Most of the power management strategies are using DVFS. Since Intel SandyBridge family processors, Intel provide Running Average Power Limit (RAPL) for

controlling the power constraint on processors and memory. Several studies have begun to evaluate the RAPL power management system. Rountree et al. [27] explore RAPL as a replacement for DVFS in HPC systems by evaluating power consumption for package and memory subsystem. Zhang et al. [28] give a systematic evaluation of RAPL behavior such as stability, setting time, overshoot and, etc. RAPL also been used to study application runtime variability and power optimization for exascale computing in work of Allan et al. [29][21] Sarood et al. [30] use RAPL to set power bounds on across an over-provisioned cluster running homogeneous application processes. This work also takes advantage of RAPL power controlling ability to dynamically manage power consumption.

Measuring the power/energy consumption of software components is the key for energy-aware scheduling, accounting and budgeting. RAPL measurement mechanism is described in [29] and Marcus et al. [31] has investigated the RAPL measurements performance and discussed the practical obstacles that existed in performing these measurements on complex modern CPUs. Vignesh et al. [32] studied the greenness of the in-situ and the post-processing visualization pipelines using RAPL to measure the CPUs and RAMs power consumption and with an average error rate of less than 1%. Venkatesh et al. [33] use RAPL to measure energy consumption in large message-passing applications.

### III. POWER SCHEDULING FOR AMR WORKLOADS

Before going into the experiment details, we has evaluated the relationship among AMR refinement resolution, CPU power level and the energy consumption. As shown in the Figure 10, the execution time will increase as the level of refinement/resolution increases or cap down CPU power. The energy consumption presents the same trend as the execution time, i.e., LMC will consume more energy as the levels of refinement/resolution increase or capping down CPU power. The question here is how to use this characterization to create available power budget. We propose to combine these two factors, adjusting resolution and applying appropriate power capping, to get available power budget for running other tasks (e.g., checkpointing).

#### Configuration

Parameter	Value
Platform	CAPER cluster
Number of cores	64 cores
Resolution	level 3 2 1
Time step	10
Time unit	Second
Power measurement	RAPL meter
Power cap	RAPL

Table I: Configuration for experiment of getting available power budget through power capping and resolution degradation

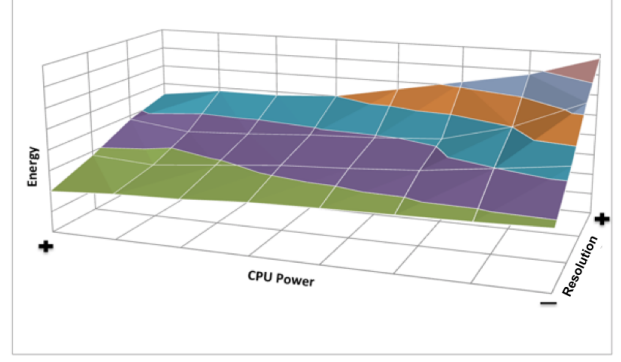


Figure 1: Energy consumption trend for different power caps and refinement levels

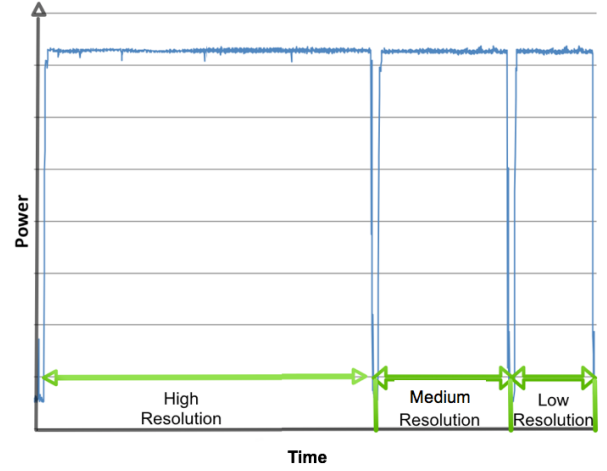


Figure 2: LMC power consumption curve with different (qualitative) resolution levels on 64 cores

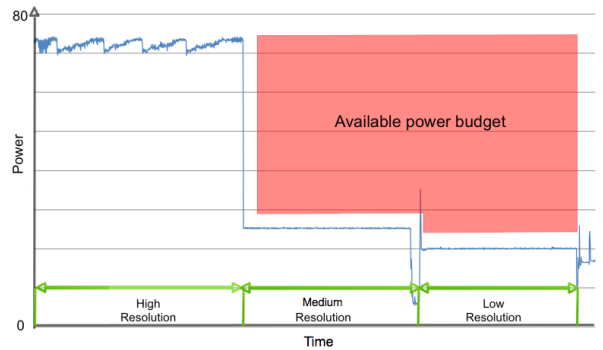


Figure 3: Available power budget from applying resolution degradation and appropriate power capping

To briefly demonstrate this concept, we measure the LMC execution time with different resolution levels on 64 cores, as shown in Figure 2. Obviously, the execution time is decreasing as we degrading the resolution levels. In the

meantime, capping down the CPU power will increase the execution time. Therefore, if applying appropriate power capping to it, then each execution can be the same. In the work, since we define the execution time as the performance criteria, by adjusting resolution and applying appropriate power capping, we can keep the performance while having the available power budget, shown as the red shadow area in Figure 11.

#### IV. EXPERIMENTAL EVALUATION

##### A. Infrastructure overviews

Caliburn Caliburn is the new supercomputer from RDI2, which is ranked #166 worldwide in the Top500 list of June 2016. This system is based on TatTwin SuperServer solution and it has 560 nodes containing 20,160 cores, 140 TB of RAM memory, 218 TB of non-volatile memory and 100 Gbps Omni-Path fabric (OPA) interconnection, which deliver performance of 603 TFLOPS with a peak performance of 677 TFLOPS. Other details of each node's hardware is given in Table II.

H/W Type	H/W Detail
CPU	2 Intel Xeon E5-2695v4
CPU frequency	2.1 GHz
# of Cores	2 18
Cache	2 45 MB
Memory	
Disk	
Network bandwidth	

Table II: Caliburn hardware specification

##### B. Power measurement

1) *RAPL*: From Intel Sandy Bridge family, Intel released Power Gadget API using RAPL to provide power measurements for components within CPU, including current estimated processor power, current processor frequency, base frequency, thermal design power (TDP), current temperature, timestamps and etc. We use it to collect the processor's current estimated power consumption (package power). RAPL measuring power data through reading Machine-Specific Registers (MSRs) which frequency can up to 1 KHz. But to reduce the interference, we set the measuring ratio to 1 Hz. According to [in-situ paper], the overhead of setting RAPL at 1 Hz frequency is increase 0.2 W in average, which is negligible.

2) *SMCIPMITOOL*: The Caliburn supports Intelligent Platform Management Interface (IPMI), which is a set of computer interface specifications that is led by Intel. [wiki [https://en.wikipedia.org/wiki/Intelligent\\_Platform\\_Management\\_Interface](https://en.wikipedia.org/wiki/Intelligent_Platform_Management_Interface)] This IPMI allow administrators to manage the system remotely and monitor platform status as well, such as system power supplies, temperatures, fans and ,etc. The Caliburn vendor Supermicro provides a tool SMCIPMITOOL that allows users to interface

with IPMI devices via a command line interface.[[ftp://ftp.supermicro.com/utility/SMCIPMITool/SMCIPMITool\\_User\\_](ftp://ftp.supermicro.com/utility/SMCIPMITool/SMCIPMITool_User_)]. In the current system set up, the SMCIPMITOOL can provide us system-wider total power measurement of four nodes at frequency of 0.25Hz, one time per four seconds. Due to the limitation of wide nodes' measuring range and relatively low sampling ratio, we will launch our application from the start node, execute relatively long time and repeat the experiment to eliminate extreme value.

##### C. Performance measurement

To quantify the scientific impact of our proposed strategy, in addition to measuring power data, we also using two profiling tools to measure the I/P, memory, networking and CPU utilization. Perf is a profiler tool for Linux, it offers a rich set of commands to collect and analyze performance and trace data. We use 'perf stat' to gather system performance counter statistics, such as number of instructions, CPU cycles, cache-misses and, etc. In the meantime, we also utilize Linux command Sar (System Activity Report), which can report on various system loads, to measure memory/paging, device load and network.

##### D. Application configurations

These four simulation applications are simulating different problems, but they all based on AMR algorithm, which uses a hierarchical grid structure, and fill finer patches on the region of interest. We tune their refinement level via their input configuration, to have high, medium and low resolutions.

##### E. Methodology

###### 1) Finding appropriate power value, using power cap:

In the previous section, we conclude that the application execution time will decrease as level of refinement/resolution decrease or increase as capping down CPU power. In order to keep application execution time the same with high, medium and low resolution, we first record the execution time of high resolution without implement power cap, and then iteratively capping down CPU power by 10 W for running application in medium or low resolution, until their execution time is matching the high resolution's execution time within 5% difference. In this way, we can get the appropriate power cap value for each level of resolution.

2) *Measuring performance*: This experiment is running on 512 cores across 15 nodes. However, due to the limitation of Perf and Sar, they only can profile a single node. Therefore, we only profile the node that launch the experiment. For example, the experiment is running from node 1 to node 15, we only measure the performance of node 1. In addition, the power consumption measurement covers all the 15 nodes. SMCIPMITOOL is running on login node, and can measure every node's power consumption.

\*\*\*\*\*Below is the old contents\*\*\*\*\*

3) *Power measurement methodology*: Our experiment platform CAPER is instrumented with both coarse- and fine-grained power metering at server level. On the one hand, an instrumented Raritan iPDU PX2-4527X2U-K2 provides power measurements at 1 Hz. While, on the other hand, a Yokogawa DL850E ScopeCorder provides voltage and current measurement for all nodes through 1 Ms/s modules, and can sample power data at rate of up to 10 KHz. In addition to those server level measurement, we have RAPL meter to provide power measurement at a up to 20Hz sampling rate in processor level.

#### Use of RAPL power metering

RAPL meter measuring power through reading Machine-Specific Registers (MSRs). Intel released Power Gadget API for using RAPL. This API is a framework that provides very comprehensive information including reading current estimated processor power, current processor frequency, base frequency, thermal design power (TDP), current temperature, timestamps and etc. Also, Intel gives a ready-to-use software-based power monitoring tool called Intel Power Gadget, this has a completed interface for using RAPL. We are using this to measure the processors power usage.

#### RAPL meter sampling frequency considerations

Thought, according to Intels manual, [34] RAPL MSRs would update at rate of once every 1 ms, using very lower frequency (e.g. less than 20 ms) may result in significant overhead and might also increase the power consumption, which would make the data less meaningful. [35] Also, since the instantaneous processor frequency would change very frequently, sampling data may be more useful if you sample often and average the samples overtime. So, Intel recommend sampling frequency of 50 ms or upper. For our experiment, we dont require high resolution power data, so we choose to sample twice per second.

#### Reading power measurement data from PDU

Raritan PDU keep measuring the whole servers power consumption and writing the power measurement data to power log files with timestamp. CAPER has eight nodes, so we extract real-time power measurement data and the timestamp from each nodes power log file after simulation started using Simple Network Management Protocol (SNMP) queries to the PDU from a side script running in an independent node.

4) *AMR resolution adjustment*: LMC simulation program is based on AMR algorithm, which uses a hierarchical grid structure, and fill finer patches on the region of interest. Therefore, the simulations resolution from the computational point of view is mainly directed from the AMR algorithm resolution configuration. There are many parameters controlling AMR solution, but in the work, only those four parameters listed in Table III have been used to tune the simulation resolution.

Parameter	Definition
amr.n_cell	Number of cells in each direction at the coareset level
amr.max_level	Number of levels of refinement above the coarsest level
amr.ref_ratio	Ratio of coarse to fine grid spacing between subsequent levels
amr.regrid_int	How often to regrid

Table III: AMR parameters which are used to tune simulation resolution

Each parameter has different impact on the resolution, which is illustrated with the example use case shown in Figure 4. The parameters of the example case are discussed as follows.

Example case:

- amr.n\_cell = 32 32 32

This would define the domain size (at coarsest level) to have 32 cells in the x-direction, 32 cells in the y-direction and 32 cells in the z-direction. (If it is in the 2D input file, the last number will be ignored). As shown in Figure 4, there are 32 cells in both x and y direction.

- amr.max\_level = 2

This would set a maximum of 2 refinement levels in addition to the coarse level. Within the calculation, the number of refinement level must be  $\leq$  amr.max\_level, but it can be change in time and it is not necessary always be equal to amr.max\_level. Because these additional levels will only be created if the tagging criteria are such that cells are flagged as needing refinement. Figure 4 shows the mesh grids with maximum of 2 refinement levels.

- amr.ref\_ratio = 4 2

Refinement ration means how many individual cells will a cell be divided into. For example in the left-hand side Figure 4, Setting amr.ref\_ratio = 4 2 means dividing cell into 4 cells from levels 0 and 1, and dividing cell into 2 from levels 1 and 2.

- amr.regrid\_int = 2

This would tell the code to regrid every 2 steps. This means level l+1 grids will be created every 2 level l time steps.

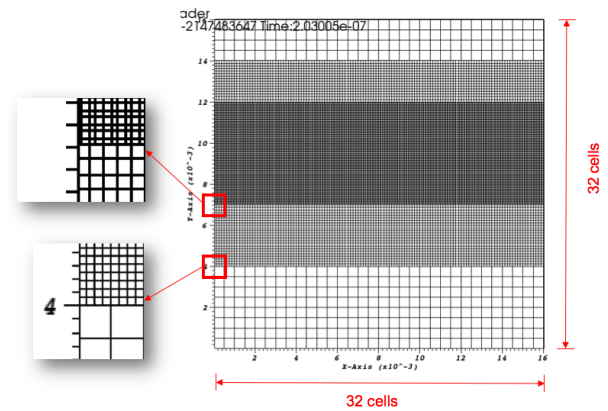


Figure 4: AMR mesh grids outcome of the example configuration

## V. RESULTS

1) *Evaluation of different levels of refinement for different number of cores:* This experiment is first run on CAPER cluster, with a maximum of 128 cores, from level 1 to level 4 on 2, 4, 8, 16, 32, 64 and 128 cores. The execution time is recorded simply with the system tool “time”.

### Configuration

Parameter	Value
Platform	CAPER and DELL cluster
Number of cores	2 to 128 quadratic growth
Level of refinement	level 1 to level 4
Time step	10
Time unit	Second

Table IV: Configuration for experiment of execution time of different levels of refinement under variety number of cores

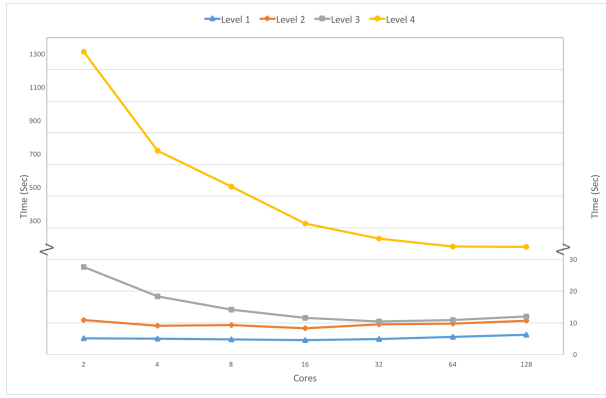


Figure 5: Execution time of running different levels LMC on CAPER cluster through 2 to 128 cores.

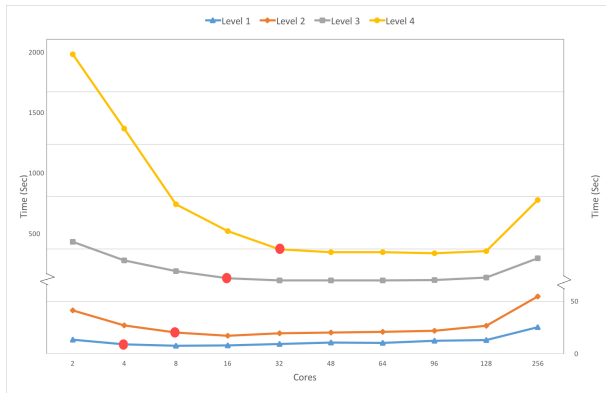


Figure 6: Execution time of running different levels LMC on DELL cluster through 2 to 256 cores.

The results of the first experiment using CAPER cluster are shown in Figure 5. As can be clearly seen from the level 4 line, which is yellow line with solid dot maker, the

more computing resource we use, the faster the execution is. Comparing these four lines, we can also notice that the slope is decreasing from level 4 to level 1, and then the lines go to a relative flat zoom. From this notice, we can see that heavier workload need more computational resources, but provisioning more resources do not bring any performance increase. In Figure 5 it can be seen from level 1 to level 3 lines that their tails are tilting a little bit, but level 4 line doesn't have this trend. In order to show its trend more clearly, we ran this experiment on DELL cluster (256 cores cluster). The result is shown in Figure 6. All levels' curve is like a “U”, that tilt in the head and tile. From level 1 to level 4, the head turning point, which marked by red dot, are respectively 4, 8, 16 and 32 cores. Then they tend to get into a flat zoom, while under the 256 core, the execution time are all increase. From these two experiments we can conclude that, to achieve optimal running time, we need to configure appropriate number of processors to execute certain level of quality for the LMC program and using more processors does guarantee better performance.

2) *Evaluation of energy consumption of different levels of refinement under variety number of cores:* In this experiment, we are aim at exploring the energy consumption of running LMC with different level of refinement under multiple number of cores.

Energy is the product of power and time, therefore, in this experiment, we need to measure the execution time and power respectively. To measure time, we still use Linux command TIME in the script file. And for the CPU power, we are using RAPL power meter to measure it. RAPL power meter is a sub-function of Intel Power Gadget program, which has been mentioned in previous background section. We are sampling power data every 0.5 second.

### Configuration

Parameter	Value
Platform	CAPER cluster
Number of cores	2 128 quadratic growth
Level of refinement	level 1 to level 4
Time step	10
Time unit	Second
Power measurement object	8 nodes processors
RAPL sample rate	2 Hz

Table V: Configuration for experiment of power consumption of different levels of refinement under variety number of cores



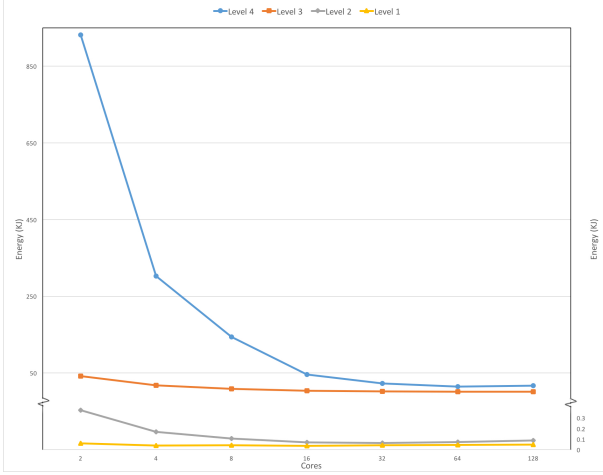


Figure 7: Power consumption of different levels of refinement under different number of cores

Figure 7 shows that there is a huge energy difference on the beginning of this curve. When running on 2 cores, the highest one (level 4) consumes almost 930 KJ while the lowest one (level 1) only take around 0.6 KJ. Since this experiment measure the all 8 node processors power status, the huge energy consumption of level 4 comes from the long execution time. It takes 1,314 seconds to run level 4 job while only 10 seconds to execute level 1. Most energy consumption comes from processors in idle state. But when assign appropriate processor resources to it, the curve reaches the flat zoom, the energy consumption difference is reasonable and stable. From this point of view, an appropriate number of cores for a certain workload is very crucial for energy consumption.

3) *Evaluation of the impact of RAPL power capping on LMC performance* : This experiment is aimed at studying the impact of power capping on LMC power-performance. And it would also give us a reference about the optimal power capping setting for different resolution levels.

We use level 1 to level 4 in LMC as inputs to measure the time and power consumption. Since CAPER CPU power is 95W (Intel Xeon E5-2650 V2), in this experiment, we cap the CPU power to 25W, 35W, 45W, 55W, 65W, 75W, 85W, and 95W and record the package energy consumption (PKG) and execution time respectively. Package energy consumption is the energy used by the CPU chip itself including cores, caches and graphics, etc.

### Configuration

Parameter	Value
Platform	CAPER cluster
Number of cores	64 cores
Level of refinement	level 1 to level 4
Time step	10
Time unit	Second
Power measurement object	8 nodes processors
Power capping	95W to 25W decrease by 10
RAPL sample rate	2 Hz

Table VI: Configuration for experiment of exploring the affection of RAPL power capping on LMC performance

We use Intel *power\_gadget* tool to cap the power. To use it, we assign value to the parameter *MY\_POWER\_LIMIT*, then compile and run it. We use linux command line “time” to measure the execution time. The execution command line is provided as follows: `time mpirun -machinefile hostfile.txt -np 64 .LMC2d.Linux.g++.gfortran.SDC.MPI.ex inputfile amr.max_lev=#level`

For the power measurement, we also use the Intel *power\_gadget* tool. When running it, it will give the real-time power consumption data every 0.5 second.

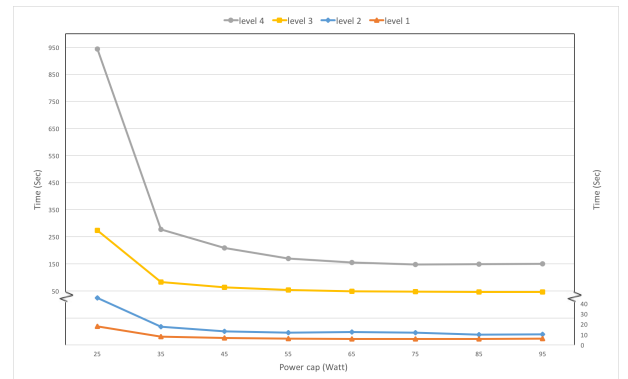


Figure 8: Relationship between different resolution of LMC and their execution time under different power capping levels

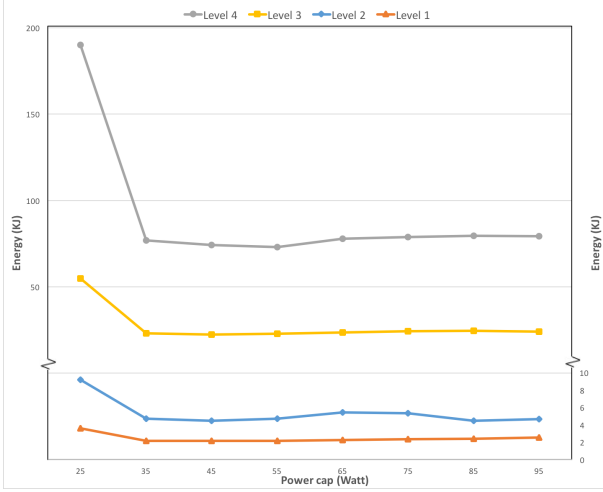


Figure 9: Relationship between different resolution of LMC and the energy consumption under different power capping levels

Figure presents the energy consumption result of executing level 1 to level 4 LMC program on certain number of cores under different power capping configurations. Those curves have the same pattern that sharply decreasing from 25 to 35 and becoming flat between 65 and 95. The lowest energy consumption is appearing when CPU power was capped to 45W to 55W. This can be explaining from execution data that without any CPU power capping, LMC program will boots CPU to 65W to 75W during execution whatever executing on how many cores. This also can be seem from the flat curve on segment 65W to 95W on each chart. Therefore, the advantage of power capping is working when CPU power down below 65W, and optimize between 45W to 55W.

4) *Evaluation of power budget acquisition through power capping and resolution degradation:* The previous experiments have characterized the behavior of LMC. It can be concluded that the execution time will increase as the level of refinement/resolution increases or CPU power is capped down. Also, the energy consumption trend for different power caps and refinement levels is shown in Figure 10. The energy consumption presents the same trend as the execution time, i.e., LMC will consume more energy as the levels of refinement/resolution increase or capping down CPU power. The question here is how to use this characterization to create available power budget. We propose to combine these two factors, adjusting resolution and applying appropriate power capping, to get available power budget for running other tasks (e.g., checkpointing).

#### Configuration

Parameter	Value
Platform	CAPER cluster
Number of cores	64 cores
Resolution	level 3 2 1
Time step	10
Time unit	Second
Power measurement	RAPL meter
Power cap	RAPL

Table VII: Configuration for experiment of getting available power budget through power capping and resolution degradation

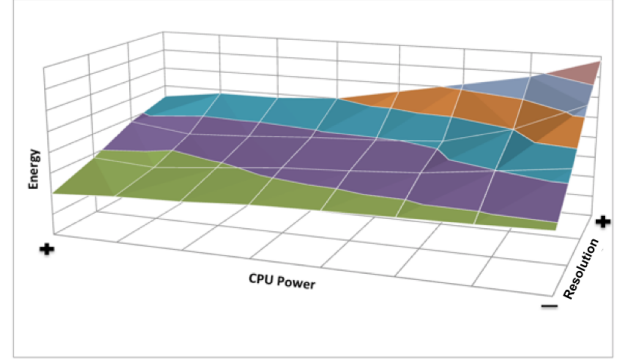


Figure 10: Energy consumption trend for different power caps and refinement levels

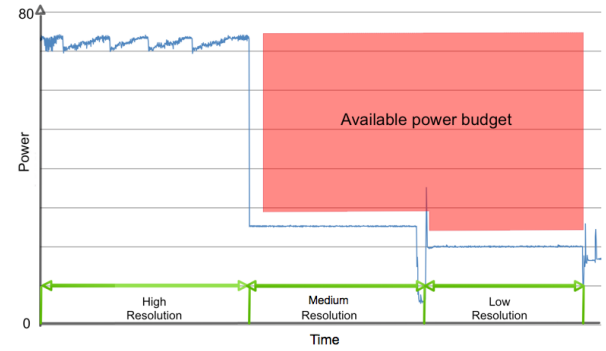


Figure 11: Available power budget from applying resolution degradation and appropriate power capping

In previous Figure ??, we measured the power consumption of running LMC with different resolution. It shows that the execution time decreases when decreasing the resolution. Therefore, we apply appropriate power capping to the medium and low resolution execution to make their execution equal the execution time in highest resolution. As shown in Figure 11, the red region is the available power budget extracted from the total power budget.

5) *Evaluation of power budget management:* We have tested the impact of power capping and level of refinement of LMC on the performance and energy consumption. We



also observed the potential power saving from degrading resolution and implementing power capping. The goal of this experiment is being able to use power budgets opportunistically for running other tasks. In the experiment, we propose to use this power budget to do checkpointing, which will increase the resilience and make the LMC program running more reliably.

In this experiment, we are proposing to degrade LMC resolution and cap its running power, then use this power budget to do checkpointing. However, there are some limitations due to the code characteristics: (i) we can not separate the checkpoint function from simulation program, and (ii) we can not dynamically control LMC resolution.

To address these two issues, we use two LMC executions running alternatively on two set of nodes to emulate one execution dynamically adjusting resolution and doing checkpoint. Each set of nodes contains three nodes for execution. When the first execution start doing checkpoint then the second set start another execution with lower resolution. We focus on system level power consumption.

We implement this using a flag to coordinate these two LMC executions. We've recompiled the first LMC set program to let it write flag when it start to do checkpoint, and the second LMC set keeps reading the flag to start running when the flag turn to be true. As shown in Figure 12, when the blue curve is the first instance (set 1) of LMC finishes its execution, it triggers the second instance (set 2), which is the orange curve. Once the second instance is completed the first one is triggered again and so on. By doing this, we can emulate a single LMC execution running and dynamically change the resolution while giving the available power budget for doing checkpoint.

#### Configuration

Parameter	Value
Platform	CAPER cluster
Number of cores	64 cores
Resolution	level 4 3
Time step	10
Time unit	Second
Power measurement object	3 nodes processors
Power cap	RAPL

Table VIII: Configuration for experiment of exploring the power-performance tradeoffs

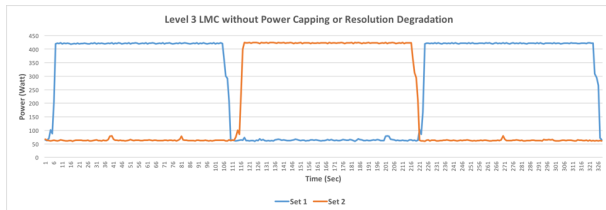


Figure 12: Level 3 LMC power consumption without power capping or resolution degradation

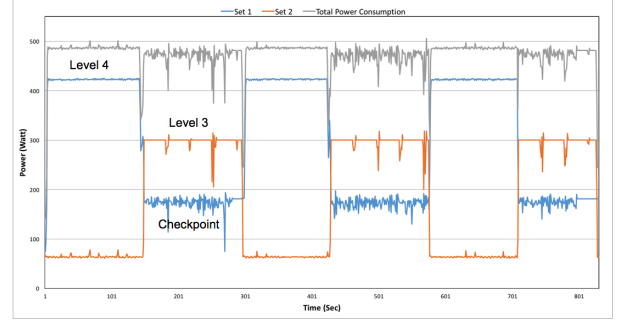


Figure 13: Level 4 LMC power consumption with power capping or resolution degradation

In Figure 13, the blue curve is the set 1 running the high resolution LMC (level 4). Once it starts doing checkpointing, the set 2 (in orange curve) is triggered to run LMC at level 3. Configuring power capping appropriately, the total power consumption, which is shown in gray curve on the top, is kept constant. This means we successfully constraint the power budget of LMC to perform other tasks (i.e., checkpointing).

#### VI. CONCLUSION AND FUTURE WORK

In this work, we have focused on studying the properties and exploring the performance, quality and power/energy tradeoffs of Low-Mach-Number Combustion (LMC) application which is based Adaptive Mesh Refinement (AMR) algorithm. The key contributions of this work are (1) we present an empirical evaluation of different configurations of application that gives insights into the energy-performance-quality tradeoff for this work, (2) we provided a comprehensive study of this LMC simulation performance, power and energy behavior, and (3) we propose a power-performance-quality tradeoff for this application, which can be used to better schedule power budgets across HPC systems.

Our current work investigates how to leverage these insights to implement a runtime to manage power budgets dynamically. Moreover, instead of trading-off with resolution, we also plan to explore the possibility of managing resources dynamically (e.g., cores) to manage power budgets. Future work also includes the management and scheduling of power budgets at whole system scale for assigning power budgets to other applications, not necessarily to the LMC workflow.

#### ACKNOWLEDGMENT

The authors would like to thank... more thanks here

#### REFERENCES

- [1] "Tianhe-2." [Online]. Available: <https://en.wikipedia.org/wiki/tianhe-2>
- [2] "China's supercomputing strategy called out," 2014. [Online]. Available: <https://www.hpcwire.com/2014/07/17/dd/>

- [3] R. Lucas, J. Ang, K. Bergman, S. Borkar, W. Carlson, L. Carrington, G. Chiu, R. Colwell, W. Dally, J. Dongarra *et al.*, "Doe advanced scientific computing advisory subcommittee (ascac) report: Top ten exascale research challenges," USDOE Office of Science (SC)(United States), Tech. Rep., 2014.
- [4] H. David, E. Gorbato, U. R. Hanebutte, R. Khanna, and C. Le, "Rapl: memory power estimation and capping," in *Low-Power Electronics and Design (ISLPED), 2010 ACM/IEEE International Symposium on*. IEEE, 2010, pp. 189–194.
- [5] M. Tolentino and K. W. Cameron, "The optimist, the pessimist, and the global race to exascale in 20 megawatts," *Computer*, vol. 45, no. 1, pp. 0095–97, 2012.
- [6] R. Ge, X. Feng, W.-c. Feng, and K. W. Cameron, "Cpumiser: A performance-directed, run-time system for power-aware clusters," in *Parallel Processing, 2007. ICPP 2007. International Conference on*. IEEE, 2007, pp. 18–18.
- [7] C.-h. Hsu and W.-c. Feng, "A power-aware run-time system for high-performance computing," in *Proceedings of the 2005 ACM/IEEE conference on Supercomputing*. IEEE Computer Society, 2005, p. 1.
- [8] S. Huang and W. Feng, "Energy-efficient cluster computing via accurate workload characterization," in *Proceedings of the 2009 9th IEEE/ACM International Symposium on Cluster Computing and the Grid*. IEEE Computer Society, 2009, pp. 68–75.
- [9] V. W. Freeh and D. K. Lowenthal, "Using multiple energy gears in mpi programs on a power-scalable cluster," in *Proceedings of the tenth ACM SIGPLAN symposium on Principles and practice of parallel programming*. ACM, 2005, pp. 164–173.
- [10] D. Donofrio, L. Oliker, J. Shalf, M. F. Wehner, C. Rowen, J. Krueger, S. Kamil, and M. Mohiyuddin, "Energy-efficient computing for extreme-scale science," *IEEE Computer*, vol. 42, no. 11, pp. 62–71, 2009.
- [11] M. Gamell, I. Roder, M. Parashar, J. C. Bennett, H. Kolla, J. Chen, P.-T. Bremer, A. G. Landge, A. Gyulassy, P. McCormick *et al.*, "Exploring power behaviors and trade-offs of in-situ data analytics," in *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*. ACM, 2013, p. 77.
- [12] B. Rountree, D. K. Lowenthal, B. R. de Supinski, M. Schulz, V. W. Freeh, and T. Bletsch, "Adagio: making dvs practical for complex hpc applications," in *Proceedings of the 23rd international conference on Supercomputing*. ACM, 2009, pp. 460–469.
- [13] M. Y. Lim, V. W. Freeh, and D. K. Lowenthal, "Adaptive, transparent frequency and voltage scaling of communication phases in mpi programs," in *SC 2006 conference, proceedings of the ACM/IEEE*. IEEE, 2006, pp. 14–14.
- [14] N. Ioannou, M. Kauschke, M. Gries, and M. Cintra, "Phase-based application-driven hierarchical power management on the single-chip cloud computer," in *Parallel Architectures and Compilation Techniques (PACT), 2011 International Conference on*. IEEE, 2011, pp. 131–142.
- [15] D. Li, B. R. De Supinski, M. Schulz, K. Cameron, and D. S. Nikolopoulos, "Hybrid mpi/openmp power-aware computing," 2010.
- [16] I. Roder, S. Chandra, M. Parashar, R. Muralidhar, H. Seshadri, and S. Poole, "Investigating the potential of application-centric aggressive power management for hpc workloads," in *High Performance Computing (HiPC), 2010 International Conference on*. IEEE, 2010, pp. 1–10.
- [17] K. Kandalla, E. P. Mancini, S. Sur, and D. K. Panda, "Designing power-aware collective communication algorithms for infiniband clusters," in *Parallel Processing (ICPP), 2010 39th International Conference on*. IEEE, 2010, pp. 218–227.
- [18] C. Lively, X. Wu, V. Taylor, S. Moore, H.-C. Chang, and K. Cameron, "Energy and performance characteristics of different parallel implementations of scientific applications on multicore systems," *International Journal of High Performance Computing Applications*, vol. 25, no. 3, pp. 342–350, 2011.
- [19] S. Pelley, D. Meisner, P. Zandevakili, T. F. Wenisch, and J. Underwood, "Power routing: dynamic power provisioning in the data center," in *ACM Sigplan Notices*, vol. 45, no. 3. ACM, 2010, pp. 231–242.
- [20] J. Hamilton, "Cost of power in large-scale data centers," *Blog entry dated*, vol. 11, p. 28, 2008.
- [21] R. Cochran, C. Hankendi, A. K. Coskun, and S. Reda, "Pack & cap: adaptive dvfs and thread packing under power caps," in *Proceedings of the 44th annual IEEE/ACM international symposium on microarchitecture*. ACM, 2011, pp. 175–185.
- [22] J. Sartori and R. Kumar, "Distributed peak power management for many-core architectures," in *Design, Automation & Test in Europe Conference & Exhibition, 2009. DATE'09*. IEEE, 2009, pp. 1556–1559.
- [23] J. M. Cebrián, J. L. Aragon, and S. Kaxiras, "Power token balancing: Adapting cmps to power constraints for parallel multithreaded workloads," in *Parallel & Distributed Processing Symposium (IPDPS), 2011 IEEE International*. IEEE, 2011, pp. 431–442.
- [24] A. Gandhi, M. Harchol-Balter, R. Das, J. O. Kephart, and C. Lefurgy, "Power capping via forced idleness," 2009.
- [25] K. Meng, R. Joseph, R. P. Dick, and L. Shang, "Multi-optimization power management for chip multiprocessors," in *Proceedings of the 17th international conference on Parallel architectures and compilation techniques*. ACM, 2008, pp. 177–186.
- [26] V. Kontorinis, A. Shayan, D. M. Tullsen, and R. Kumar, "Reducing peak power with a table-driven adaptive processor core," in *Proceedings of the 42nd annual IEEE/ACM international symposium on microarchitecture*. ACM, 2009, pp. 189–200.
- [27] B. Rountree, D. H. Ahn, B. R. De Supinski, D. K. Lowenthal, and M. 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*. IEEE, 2012, pp. 947–953.

- [28] H. Zhang and H. Hoffmann, “A quantitative evaluation of the rapl power control system,” *Feedback Computing*, 2015.
- [29] A. Porterfield, R. Fowler, S. Bhalachandra, B. Rountree, D. Deb, and R. Lewis, “Application runtime variability and power optimization for exascale computers,” in *Proceedings of the 5th International Workshop on Runtime and Operating Systems for Supercomputers*. ACM, 2015, p. 3.
- [30] O. Sarood, A. Langer, L. Kalé, B. Rountree, and B. De Supinski, “Optimizing power allocation to cpu and memory subsystems in overprovisioned hpc systems,” in *Cluster Computing (CLUSTER), 2013 IEEE International Conference on*. IEEE, 2013, pp. 1–8.
- [31] M. Hähnelt, B. Döbel, M. Völp, and H. Härtig, “Measuring energy consumption for short code paths using rapl,” *ACM SIGMETRICS Performance Evaluation Review*, vol. 40, no. 3, pp. 13–17, 2012.
- [32] V. Adhinarayanan, W.-c. Feng, J. Woodring, D. Rogers, and J. Ahrens, “On the greenness of in-situ and post-processing visualization pipelines,” in *Parallel and Distributed Processing Symposium Workshop (IPDPSW), 2015 IEEE International*. IEEE, 2015, pp. 880–887.
- [33] A. Venkatesh, K. Kandalla, and D. K. Panda, “Evaluation of energy characteristics of mpi communication primitives with rapl,” in *Parallel and Distributed Processing Symposium Workshops & PhD Forum (IPDPSW), 2013 IEEE 27th International*. IEEE, 2013, pp. 938–945.
- [34] I. Intel, “and ia-32 architectures software developer’s manual, 2011,” *Intel order Number*, 64.
- [35] “Using the intel power gadget api on mac os x.” [Online]. Available: <https://software.intel.com/en-us/blogs/2012/12/13/using-the-intel-power-gadget-api-on-mac-os-x>