

# Energy-Efficient, Thermal-Aware Modeling and Simulation of Data Centers: The CoolEmAll Approach and Evaluation Results

Leandro Cupertino<sup>a</sup>, Georges Da Costa<sup>a</sup>, Ariel Oleksiak<sup>b</sup>, Wojciech Piątek<sup>b</sup>, Jean-Marc Pierson<sup>a</sup>,  
Jaume Salom<sup>c</sup>, Laura Sisó<sup>c</sup>, Patricia Stolf<sup>a</sup>, Hongyang Sun<sup>a</sup>, Thomas Zilio<sup>a</sup>

<sup>a</sup>*Toulouse Institute of Computer Science Research (IRIT), University of Toulouse, France.*

<sup>b</sup>*Poznan Supercomputing and Networking Center (PSNC), Poznan, Poland.*

<sup>c</sup>*Catalonia Institute for Energy Research (IREC), Barcelona, Spain.*

---

## Abstract

This paper describes the CoolEmAll project and its approach for modeling and simulating energy-efficient and thermal-aware data centers. The aim of the project was to address energy-thermal efficiency of data centers by combining the optimization of IT, cooling and workload management. This paper provides a complete data center model considering the workload profiles, the applications profiling, the power model and a cooling model. Different energy efficiency metrics are proposed and various resource management and scheduling policies are presented. The proposed strategies are validated through simulation at different levels of a data center.

*Keywords:* data centers, energy efficiency, metrics, resource management policies, scheduling.

---

## 1. Introduction

IT energy impact is now a major concern from the economical point of view but also from the sustainability one. IT was responsible for around 2% of the global energy consumption making it equal to the demand of aviation industry in 2008 [1]. Focusing on data centers, late 2012 numbers from the European commission [2] shows that European data centers consumed 60TWh during 2012. The same study expects this number to double before 2020.

While this aggregated consumption is high, still nearly a third of organizations (29%) owning data centers did not measure their efficiency in 2012 [3]. Out of this study, for the data centers that measure their Power Usage Effectiveness (PUE) [4], more than a third (34%) have a PUE over or equal to 2, meaning they consume more power on cooling, air movement and infrastructure than on computing itself. The average PUE over all data centers is between 1.8 and 1.89.

Large energy needs and significant CO<sub>2</sub> emissions cause the issues related to cooling, heat transfer, and IT infrastructure location more and more carefully studied during planning and operation of data centers. Even if we take ecological and footprint issues aside, the amount of consumed energy can impose strict limits on data centers. First of all, energy bills may reach millions of euros making computations expensive. Furthermore, available power supply is usually limited so it also may reduce

---

*Email addresses:* fontoura@irit.fr (Leandro Cupertino), dacosta@irit.fr (Georges Da Costa), ariel@man.poznan.pl (Ariel Oleksiak), piatek@man.poznan.pl (Wojciech Piątek), pierson@irit.fr (Jean-Marc Pierson), jsalom@irec.cat (Jaume Salom), lsiso@irec.cat (Laura Sisó), stolf@irit.fr (Patricia Stolf), sun@irit.fr (Hongyang Sun), sun@irit.fr (Hongyang Sun), zilio@irit.fr (Thomas Zilio)

data center development capabilities, especially looking at challenges related to exascale computing breakthrough foreseen within this decade. For these reasons many efforts were undertaken to measure and improve energy efficiency of data centers. Some of those projects focused on data center monitoring and management [5, 6, 7] whereas others on prototypes of low power computing infrastructures [8, 9]. Additionally, vendors offer a wide spectrum of energy efficient solutions for computing and cooling [10, 11].

A variety of possibilities exist at the design level, which have to be simulated in order to be compared and to select the best one. During the lifetime of a data center, smart management can lead to better visibility of the platform behavior and to reduce energy consumption.

In order to optimize the design or configuration of a data center we need a thorough study using appropriate metrics and tools evaluating how much computation or data processing can be done within a given power and energy budget and how it affects temperatures, heat transfers, and airflows within the data center. Therefore, there is a need for simulation tools and models that approach the problem from a perspective of end users and take into account all the factors that are critical to understanding and improving the energy efficiency of data centers, in particular, hardware characteristics, applications, management policies, and cooling. To address these issues the CoolEmAll project [12] aimed at decreasing energy consumption of data centers by allowing data center designers, planners, and administrators to model and analyze energy efficiency of various configurations and solutions. To this end, the project provides models of data center building blocks and tools that apply these models to simulate, visualize and analyze data center energy efficiency.

The structure of the paper is as follows. Section 2 presents relevant related works. Section 3 contains a brief description of the CoolEmAll project. In Section 4 we present the models that are used in the design and management tools. In Section 5 the metrics used to assess the quality of design and management are presented. Section 6 describes smart data center management techniques. In Section 7 we show the results of the simulation experiments and the impact of the proposed models and tools. Section 8 concludes the paper.

## 2. Related Work

Issues related to cooling, heat transfer, IT infrastructure configuration, IT-management, arrangement of IT-infrastructure as well as workload management are gaining more and more interest and importance, as reflected in many ongoing works both in industry and research. There are already software tools available on the market capable to simulate and analyze thermal processes in data centers. Examples of such software include simulation codes along with more than 600 models of servers from Future Facilities [13] with its DC6sigma products, CA tools [14], or the TileFlow [15] application. In most cases these simulation tools are complex and expensive solutions that allow modeling and simulation of heat transfer processes in data centers. To simplify the analysis process Romonet [16] introduced a simulator, which concentrates only on costs analysis using simplified computational and cost models, disclaiming analysis of heat transfer processes using Computational Fluid Dynamics (CFD) simulations. Common problem in case of commercial data center modeling tools is that they use closed limited databases of data center hardware. Although some of providers as Future Facilities have impressive databases, extensions of these databases and use of models across various tools is limited. To cope with this issue Schneider have introduced the GENOME Project that aims at collecting “genes” which are used to build data centers. They contain details of data center components and are publicly available on the Schneider website [17]. Nevertheless, the components are described by static parameters such as “nameplate” power values rather than details that enable simulating and assessing their energy

efficiency in various conditions. Another initiative aiming at collection of designs of data centers is the Open Compute Project [18]. Started by Facebook which published its data center design details, it consists of multiple members describing data centers' designs. However, Open Compute Project blueprints are designed for description of good practices rather than to be applied to simulations.

In addition to industrial solutions significant research effort was performed in the area of energy efficiency modeling and optimization. For example, models of servers' power usage were presented in [19] whereas application of these models to energy-aware scheduling in [20]. Additionally, authors in [21, 22] proposed methodologies of modeling and estimation of power by specific application classes. There were also attempts to use thermodynamic information in scheduling [23]. Nevertheless, the above works are focused on research aspects and optimization rather than providing models to simulate real data centers. In [24], the authors propose a power management solution that coordinates different individual approaches. The solution is validated using simulations based on 180 server traces from nine different real-world enterprises. Second, using a unified architecture as the base, they perform a quantitative sensitivity analysis on the impact of different architectures, implementations, workloads, and system design choices. Shah [25] explores the possibility of globally staggering compute workloads to take advantage of local climatic conditions as a means to reduce cooling energy costs, by performing an in-depth analysis of the environmental and economic burden of managing the thermal infrastructure of a globally connected data center network. SimWare [26] is a data warehouse simulator which compute its energy efficiency by: (a) decoupling the fan power from the computer power by using a fan power model; (b) taking into account the air travel time from the CRAC to the nodes; (c) considering the relationship between nodes by the use of a heat distribution matrix.

### 3. The CoolEmAll Project

CoolEmAll was a European Commission funded project which addresses the complex problem of how to make data centers more energy and resource efficient. CoolEmAll developed a range of tools to enable data center designers, operators, suppliers and researchers to plan and operate facilities more efficiently. The participants in the project included a range of scientific and commercial organizations with expertise in data centers, high performance computing, energy efficient server design, and energy efficient metrics.

The defining characteristic of the CoolEmAll project is that it bridges this traditional gap between IT and facilities approaches to efficiency. The main outcomes of CoolEmAll are based on a holistic rethinking of data center efficiency that is crucially based on the interaction of all the factors involved rather than just one set of technologies. The expected results of the project included a data center monitoring, simulation and visualization software, namely SVD toolkit, designs of energy efficient IT hardware, contribution to existing (and help define new) energy efficiency metrics.

Some commercial suppliers (most notably Data center Infrastructure Management suppliers) and consultants have recently begun to take a more all-encompassing approach to the problem by straddling both IT and facilities equipment. However, few suppliers or researchers up to now have attempted to include the crucial role of workloads and applications. That is beginning to change, and it is likely that projects such as CoolEmAll can advance the state of the art in this area.

As noted in [27], the objective of the CoolEmAll project was to enable designers and operators of a data center to reduce its energy impact by combining the optimization of IT, cooling and workload management. For this purpose CoolEmAll investigated in a holistic approach on how cooling, heat transfer, IT infrastructure, and application-workloads influence overall cooling- and energy-efficiency of data centers, taking into account various aspects that traditionally have been considered separately.

In order to achieve this objective CoolEmAll provided two main outcomes: (i) design of diverse types of computing building blocks well defined by hardware specifications, physical dimensions, and energy efficiency metrics, and (ii) development of simulation, visualization and decision support toolkit (SVD Toolkit) that enables analysis and optimization of IT infrastructures built of these building blocks. Both building blocks and the toolkit take into account four aspects that have a major impact on the actual energy consumption: characteristics of building blocks under variable loads, cooling models, properties of applications, applications workloads, and workload and resource management policies. To simplify selection of right building blocks used to design data centers adjusted to particular needs, data center efficiency building blocks are precisely defined by a set of metrics expressing relations between the energy efficiency and essential factors listed above. In addition to common static approaches, the CoolEmAll approach also enables studies and assessment of dynamic states of data centers based on changing workloads, management policies, cooling methods, environmental conditions and ambient temperature. This enables assessment and optimization of data center energy/cooling efficiency also for low and variable loads rather than just for peak loads as it is usually done today. The main concept of the project is presented in Figure 1.

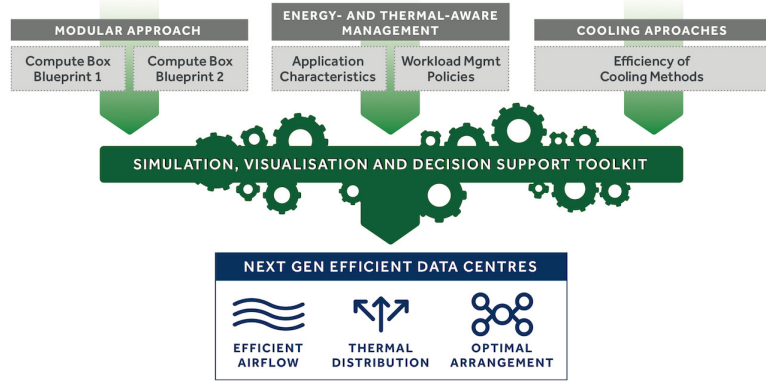


Figure 1: The CoolEmAll concept

## 4. Data Center Modeling

### 4.1. Data Center Efficiency Building Block (DEBB)

One of the main results of the CoolEmAll project is the design of diverse types of data center building blocks on different granularity levels, following a blueprint-specification format called Data center Efficiency Building Block (DEBB).

A DEBB is an abstract description of a piece of hardware and other components reflecting a data center building block on different granularity levels. To illustrate the concept, the DEBB in CoolEmAll was constructed around the RECS (Resource Efficient Computing & Storage) unit [28], a multi-node computer system developed by Christmann [29] with high energy-efficiency and density. The following describes the different granularity levels defined in the DEBB:

1. *Node unit* is the finest granularity of building blocks to be modeled within CoolEmAll. This smallest unit reflects a single CPU module in a RECS.
2. *Node group* is an ensemble of building blocks of level 1, e.g. a complete RECS unit (currently consisting of 18 computing nodes in RECS2.0).

3. *Rack (ComputeBox1)* is a typical rack within an IT service center, including building blocks of level 2 (Node Groups), power supply units and integrated cooling devices.
4. *Room (ComputeBox2)* is an ensemble of building blocks of level 3, placed in a container or compute rooms, with the corresponding CRAC/CRAH (Compute Room Air Conditioner or Air-Handling Unit), chiller, power distribution units, lighting and other auxiliary facilities.

A DEBB on each granularity level is described in the following. More details the on definitions of these components can be found in [30].

- Specification of components and sub-building blocks,
- Outer physical dimensions (black-box description), and optionally arrangements of components and sub-building blocks within particular DEBB (white-box description),
- Power consumption for different load-levels concerning mainly CPU and memory, and optionally IO and storage,
- Thermal profile describing air-flow (including direction and intensity) and temperature on inlets and outlets for different load-levels,
- Metrics describing energy efficiency of a DEBB.

A computing node will be the smallest unit of the modeling process in DEBB. The models established at a lower level, e.g., a Node unit or Node group should provide building blocks to the modeling of larger modules, e.g. full racks or server rooms, for simulations. In this way, DEBBs can improve and facilitate the process of modeling, simulation, and visualization of data centers by delivering predefined models with comprehensive information concerning performance, power consumption, thermal behavior, and shape of data center components.

## 4.2. Workload Characterization

### 4.2.1. Workload Specification

In terms of workload management, workload items are defined as jobs that are submitted by users [31]. In general, workloads may have various shapes and levels of complexity ranging from multiple independent jobs, through large-scale parallel tasks, up to single applications that require single resources. That allows distinguishing several levels of information about incoming jobs. These levels are presented in Figure 2. It is assumed that there is a queue of jobs submitted to the resource manager, and each job consists of one or more tasks. If preceding constraints between tasks are defined, a job may constitute a whole workflow.

The aim of workload profile is to provide information about structure, resource requirements, relationships and time intervals of jobs and tasks that will be scheduled during the workload simulation phase. Having these dependencies established, it is possible to express the impact of each workload item on the system. To this end, each job specified within the workload has to be extended with the particular application characteristic describing its behavior on the hardware. Thus, workload profile contains the references to the corresponding application profiles that are linked during the simulation. To model the application profile in more detail, CoolEmAll follows the DNA approach proposed in [32]. Accordingly, each task can be defined as a sequence of phases that show the impact of the application on the resources that run it. Phases are then periods of time within which the system is stable (cpu load, network, memory, etc.) given a certain threshold. More details concerning application profiles are

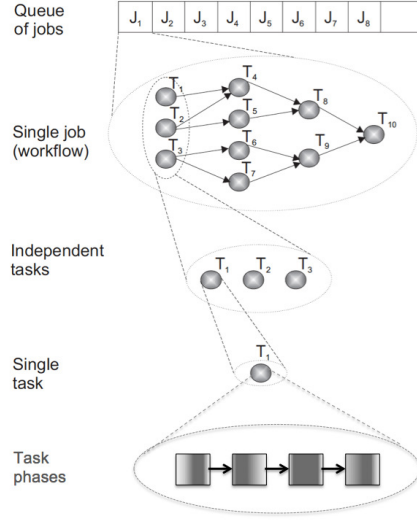


Figure 2: Workload model

provided in the next section. This form of description enables definition of a wide range of workloads: HPC (long jobs, computational-intensive, hard to migrate) or web service (short requests) that are typical for virtualized data center environments.

For the purposes of the workload description within the CoolEmAll project we adopted Standard Workload Format (SWF) [33] that is used for the traces stored in the Parallel Workloads Archive. For now it is one of the main and commonly used formats providing unitary description of both workloads models as well as logs obtained from real systems. In addition to the predefined labels in the header comments, described by Feitelson in [33], we introduce support of a new header label that is used to provide information about types of applications. An example of a workload expressed in SWF is presented in Figure 3.

In general, workload profiles may be taken from real systems or generated synthetically. The main goal of synthetic workloads is to capture the behavior of real observed workloads and to characterize them at the desired level of detail. On the other hand, they are also commonly adopted to evaluate the system performance for the modified or completely theoretical workload models. Usage of synthetic workloads and their comparison to the real ones have been the subject of research for many years. In [34], the authors analyzed both types of workloads in terms of their accuracy and applicability. Today, several synthetic workload models have been proposed [35, 36], which are based on workload logs collected from large scale parallel systems in production use. In a set of experiments depicted in Section 7, we define workloads using arrival rate based on the Poisson process as it has been typically adopted to reflect the task arrivals in supercomputing clusters [35, 36] as well as in web servers [37].

#### 4.2.2. Application Specification

For the purpose of CoolEmAll, applications behavior can be assimilated to its resource consumption. Indeed, CoolEmAll project aims at evaluating the impact of applications from a thermal and energy point of view. Using resources consumption allows evaluating this impact. As applications are usually complex, their resource consumption cannot be assimilated as constant during their lifetime. Applications will be considered as a sequence of phases. One phase will be considered as a duration during which

```

;StartTime: Thu Jan 17 10:00:00 CET 2013
;Application: 1 Linpack
;Application: 2 Abinit
;Application: 3 Mencoder
1 514 -1 8720 1 -1 -1 1 -1 -1 -1 -1 -1 3 -1 -1 -1 -1
2 1057 -1 4513 8 -1 -1 8 -1 -1 -1 -1 -1 2 -1 -1 -1 -1
3 1589 -1 7266 2 -1 -1 2 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
4 2161 -1 6636 4 -1 -1 4 -1 -1 -1 -1 -1 2 -1 -1 -1 -1
5 2742 -1 12845 2 -1 -1 2 -1 -1 -1 -1 -1 3 -1 -1 -1 -1
6 3284 -1 7867 1 -1 -1 1 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
7 3831 -1 8361 4 -1 -1 4 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
8 4409 -1 5644 8 -1 -1 8 -1 -1 -1 -1 -1 3 -1 -1 -1 -1
9 4968 -1 8493 8 -1 -1 8 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
10 5524 -1 3780 4 -1 -1 4 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
11 6103 -1 7568 1 -1 -1 1 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
12 6634 -1 6818 16 -1 -1 16 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
13 7225 -1 5663 8 -1 -1 8 -1 -1 -1 -1 -1 2 -1 -1 -1 -1
14 7810 -1 11140 2 -1 -1 2 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
15 8391 -1 8333 2 -1 -1 2 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
16 8941 -1 4164 8 -1 -1 8 -1 -1 -1 -1 -1 3 -1 -1 -1 -1
17 9527 -1 7568 8 -1 -1 8 -1 -1 -1 -1 -1 2 -1 -1 -1 -1
18 10103 -1 10431 1 -1 -1 1 -1 -1 -1 -1 -1 2 -1 -1 -1 -1
19 10674 -1 6209 4 -1 -1 4 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
20 11279 -1 5760 4 -1 -1 4 -1 -1 -1 -1 -1 2 -1 -1 -1 -1
21 11825 -1 7871 8 -1 -1 8 -1 -1 -1 -1 -1 2 -1 -1 -1 -1
22 12369 -1 5709 8 -1 -1 8 -1 -1 -1 -1 -1 2 -1 -1 -1 -1
23 12912 -1 8616 16 -1 -1 16 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
24 13452 -1 1958 8 -1 -1 8 -1 -1 -1 -1 -1 3 -1 -1 -1 -1
25 14038 -1 7770 4 -1 -1 4 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
26 14573 -1 2370 8 -1 -1 8 -1 -1 -1 -1 -1 3 -1 -1 -1 -1
27 15147 -1 11933 8 -1 -1 8 -1 -1 -1 -1 -1 2 -1 -1 -1 -1
28 15708 -1 7730 2 -1 -1 2 -1 -1 -1 -1 -1 2 -1 -1 -1 -1
29 16238 -1 3576 8 -1 -1 8 -1 -1 -1 -1 -1 3 -1 -1 -1 -1
30 16815 -1 8602 4 -1 -1 4 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
31 17352 -1 8785 4 -1 -1 4 -1 -1 -1 -1 -1 1 -1 -1 -1 -1
32 17936 -1 10676 1 -1 -1 1 -1 -1 -1 -1 -1 2 -1 -1 -1 -1

```

Figure 3: SWF workload definition

Table 1: Resources monitored for profiling applications

Resource name	Content
CPU	Percentage of processor load
network	Percentage of network bandwidth
IO	Percentage of disk bandwidth
memory	Percentage of memory bandwidth

resource consumption can be considered as constant. As the same application will consume different amount of resources depending on the hardware on which it runs, application profiles will encompass the hardware on which it ran. Using a translator it will be possible to take a profile obtained on a particular hardware and to translate it to the probable resource consumption on different hardware. Exact values will be in percentage of maximum available resource. For instance, for CPU this will be the load, and for the network this will be the ratio (in percentage) between the actual bandwidth and the maximum on the platform.

One phase will be characterized by its duration, by the mean resource consumption during this duration and by the reference hardware used to obtain those values. As an example, a simplified XML description of an application could include the section shown in Figure 4 where there are two phases, one of 4 seconds, one of 40 seconds. The first one uses mainly the CPU while not using the memory infrastructure and thus can be labeled as CPU-intensive. The second one loads at the maximum of CPU and memory. Such phase is usually labeled as memory-intensive. The current available resources are shown in Table 1. Once a profile is acquired, it can be displayed as shown in Figure 5.

#### 4.2.3. Power Model

Power measurement is a key feature for the development and maintenance of energy efficient data centers. The use of power models enables the estimation of application's power dissipation, offering a higher granularity for the measurement and leveraging the application scheduling with energy consump-

```

<resourceConsumptionProfile>
  <resourceConsumption>
    <referenceHardware>Intel_i7</referenceHardware>
    <duration>PT4S</duration>
    <behaviour name="cpu">
      <value>99</value>
    </behaviour>
    <behaviour name="memory">
      <value>29</value>
    </behaviour>
  </resourceConsumption>
  <resourceConsumption>
    <referenceHardware>Intel_i7</referenceHardware>
    <duration>PT40S</duration>
    <behaviour name="cpu">
      <value>98</value>
    </behaviour>
    <behaviour name="memory">
      <value>92</value>
    </behaviour>
  </resourceConsumption>
</resourceConsumptionProfile>

```

Figure 4: Example of profile of an application composed of two phases

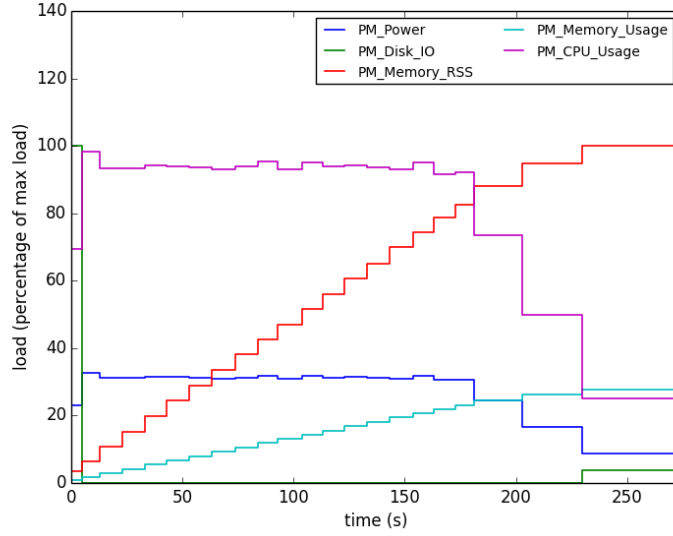


Figure 5: Graphical representation of application profile for a ray tracing benchmark (c-ray)

tion. Besides, even power meters can present some inaccuracy and the use of such models can enhance the measurements.

In CoolEmAll, we target at the RECS compute box, a high density computing system with embedded power meters for each of its computing nodes. In this study, we used a RECS version 2.0 with two types of modules: Intel Core i7-3615QE processor with 16 GB of RAM and Intel Atom N2600 processor with 2 GB of RAM. The compute box is populated with 18 nodes in total – 6 i7 and 12 atom nodes.



The embedded RECS’s power meters have a precision of 1W, which for some usages may not be enough. Even more, when running some experiments, we noticed that the power meter measurement inaccuracy can reach up to 24W. These experiments were executed by stressing one node, while the others remains turned off; in this configuration, the power of the turned off nodes reported by the power meter varies according to the stressed node and reaches from 0 to 9W maximum. So the three most erroneous nodes were stressed and the power of the turned off nodes summed up 24W. This enforces the need for estimating the power even if we dispose of a physical meter. For modeling the power dissipation of a node, we chose one of the nodes which presented less noise in other nodes and included an external power meter to provide higher precision measurements.

A usual way of modeling the power consumption of a node is to create a CPU proportional estimator. As the processor is claimed to be the most power hungry device on a computing system [38], capacitive models are greatly used as follows:

$$P = (cv^2f)u, \quad (1)$$

where the power ( $P$ ) is estimated based on the CPU’s voltage ( $v$ ), frequency ( $f$ ), capacitance ( $c$ ) and usage ( $u$ ). This model does not take into account the type of operation that is executed by the CPU. For instance, the same CPU load can provide different power consumption according to the device it uses [39]. Previous work verified that CPU temperature has a high correlation with the dissipated power [40], even though one cannot decouple the temperature between applications.

In CoolEmAll we use an application level estimator based on performance counters, model specific registers (MSR) and system information. It has been shown that calibrated linear models can provide estimation to generic applications with an accuracy error smaller than 10% [41, 40]. Performance counters (PC) are CPU counters that quantify the number of events done by the processor per core, e.g. cache misses. These counters can be fetched at process level, making the transition to an application level modeling straight forward. MSR provides precise information regarding the processor’s operating frequency and C-states. Although MSRs cannot provide process level information, we can join its information to other process dependent variables such as PC. System information is fetched from the operating system and provides data for networking and memory usage. A complete list of evaluated variables can be founded in Table 2.

A set of synthetic benchmarks was created to simulate a generic application running on this platform. This benchmark set consists of four phases; first we progressively stress the CPU by increasing its usage in 20% steps. This procedure is repeated for three frequencies (1.2, 2.3GHz and Boost) and three CPU intensive benchmarks which stress the control unity, floating point unity and the random number generator. The second phase stressed the memory access, by forcing read/write access to all caches (L1d, L2 and L3) and the RAM. Finally the network is stressed by running Linux’s `iperf` tool and limiting the download/upload to 200, 400 and 1000 Mbit/s. These data were then used to calibrate the capacitive model and to create a linear model using the above mentioned variables. The power profile of this synthetic workload can be found in Figure 6.

The results of the calibration of two models can be seen in Figure 7. One can see that the capacitive model fails when different programs present the same CPU load and lacks the power dissipation due to RAM memory access, presenting an mean absolute error (MAE) of 2.38 W and a correlation factor of 0.6961. The use of performance counters, even as a black box provides a better estimation with a MAE of 0.51 W and a correlation of 0.9831. The results of the black box present a better precision than the embedded power meter, which has 1W precision.

Table 2: Power modeling variables.

perf.cycles	perf.instructions	perf.cache_references
perf.cache_misses	perf.branch_instructions	perf.branch_misses
perf.bus_cycles	perf.idle_cycles_frontend	perf.idle_cycles_frontend
perf.cpu_clock	perf.task_clock	perf.page_faults
perf.context_switches	perf.cpu_migrations	perf.minor_faults
perf.major_faults	perf.alignment_faults	perf.emulation_faults
perf.L1d_loads	perf.L1d_load_misses	perf.L1d_stores
perf.L1d_store_misses	perf.L1d_prefetch_misses	perf.L1i_load_misses
perf.LLC_loads	perf.LLC_load_misses	perf.LLC_stores
perf.LLC_stores_misses	perf.L1d_prefetches	perf.LLC_prefetch_misses
perf.dTLB_loads	perf.dTLB_load_misses	perf.dTLB_stores
perf.dTLB_store_misses	perf.iTLB_loads	perf.iTLB_load_misses
perf.branch_loads	perf.branch_load_misses	perf.node_loads
perf.node_load_misses	perf.node_stores	perf.node_store_misses
perf.node_prefetches	perf.node_prefetch_misses	msr.cpu_freq
msr.cpu_time_in_c0	sys.cpu_use	sys.RRS
sys.received_bytes	sys.sent_bytes	

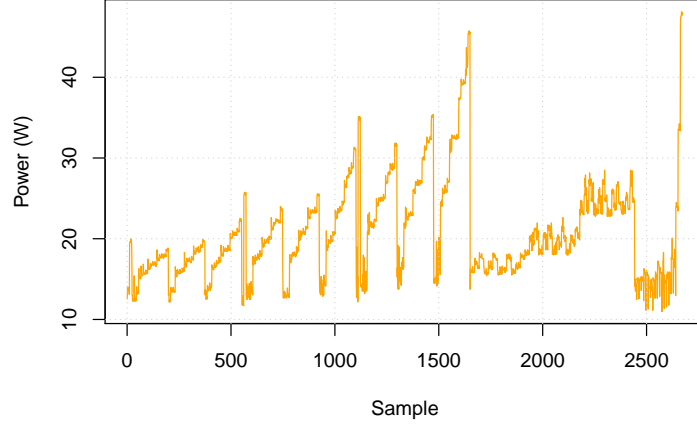


Figure 6: Synthetic workload's power profile.

#### 4.3. Cooling Model

The cooling model defined in CoolEmAll has the objective of calculating the power of the cooling equipment and other electric devices in a data center as a function of IT workload, ambient temperature and room set-up operation temperature. The model is based on a simple data center with a computer room air handler (CRAH), e.g., fan and air-water coil, power distribution unit (PDU) and lighting. All these elements generate thermal load and provide the cooling and power requirements for operating the IT components. Outside the data center, a chiller provides cooling water to the CRAH and dissipates the exhausted heat from the data center to the atmosphere by a dry-cooler (Figure 8 shows details). Other

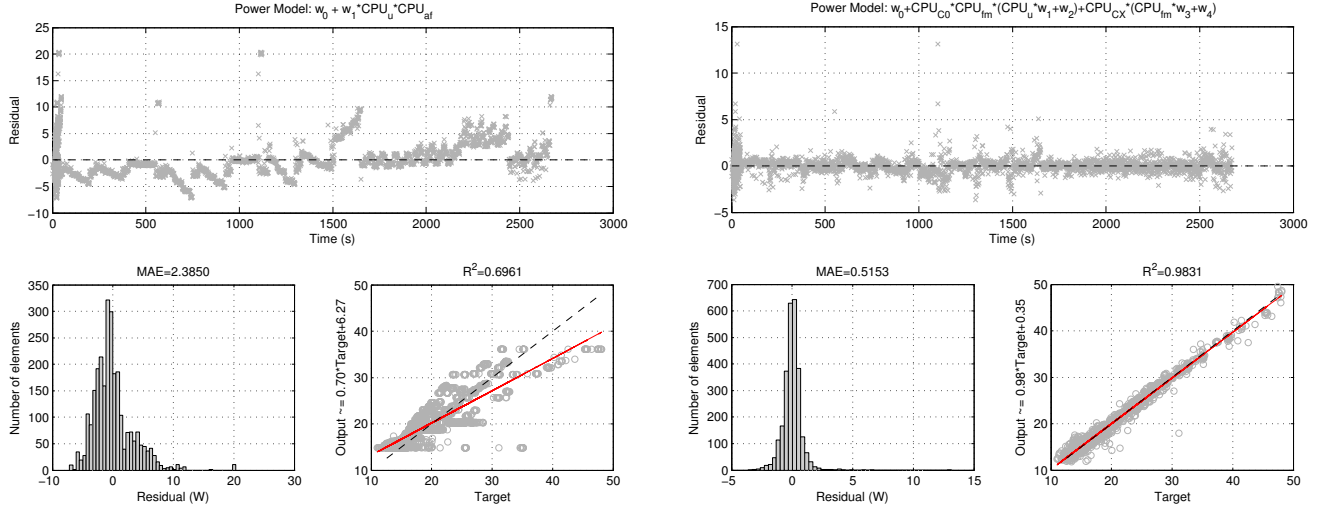


Figure 7: Model comparison for a capacitive and performance counter's linear models.

electric components such as uninterruptible power supply (UPS), back-up generator and transformer are excluded from the present model.

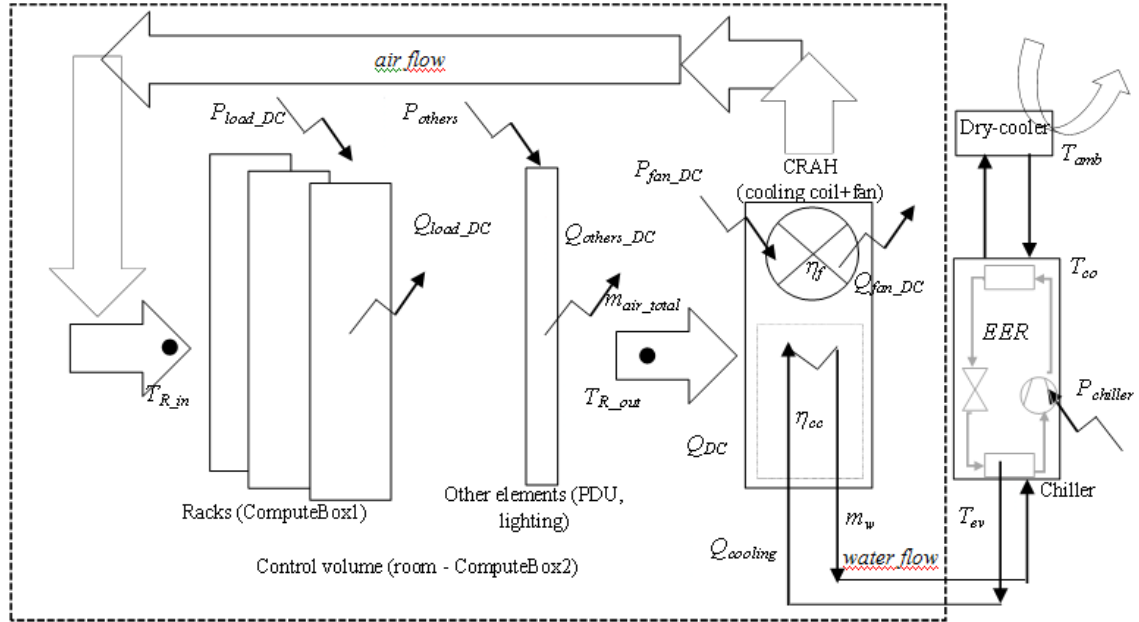


Figure 8: Cooling model of a ComputeBox-based data center

The following describes the cooling model at a single time stamp. In this model,  $Q$  is referred to as the heat dissipation and  $P$  the power consumption. The variables that are varying with time are indicated by the time index  $t$ , otherwise they have constant values in the model. As an overview, the model has been constructed based on basic thermodynamic equations of conservation of mass and

energy. The total power consumption in the data center ( $P_{DC}$ ) is calculated from the knowledge of the IT load in the data center ( $P_{load\_DC}$ ), some boundary conditions such as the inlet air room operation temperature ( $T_{R\_in}$ ), the ambient temperature ( $T_{amb}$ ), the relation between other loads and IT load ( $\alpha$ ), and the performance parameters of chiller, fan, cooling coil of CRAH and coil of dry-cooler.

On the one hand, different thermal loads are defined in the data center, and the total thermal load  $Q_{DC}$  is the sum of them — the heat associated to IT load  $Q_{load\_DC}$ , the heat from other loads such as PDU and lighting  $Q_{others\_DC}$  and the heat from fans distributing air inside the data center room  $Q_{fan\_DC}$  — that are defined as follows:

$$Q_{DC}(t) = Q_{load\_DC}(t) + Q_{others\_DC}(t) + Q_{fan\_DC}(t) \quad (2)$$

On the other hand, the total power consumption of the data center will be sum of powers from these components and the power consumed by the chiller:

$$P_{DC}(t) = P_{load\_DC}(t) + P_{chiller}(t) + P_{fans\_DC}(t) + P_{others}(t) \quad (3)$$

The following shows how to calculate each of these power consumptions.

The heat associated to the IT load  $Q_{load\_DC}$  is assumed to be equal to the power of IT load  $P_{load\_DC}$  [42], which is calculated as the sum of the IT load of each rack  $P_{load\_rack}$ . As it is stated in [42], this assumption is possible since the power transmitted by the information technology equipment through the data lines can be neglected.

$$P_{load\_DC}(t) = Q_{load\_DC}(t) \quad (4)$$

$$P_{load\_DC}(t) = \sum_{j=1}^{N_r} P_{load\_rack}(t) \quad (5)$$

For the loads on lighting and PDU, a factor  $\alpha$  is used to relate their power consumption with IT load. According to [43],  $\alpha$  is estimated to be around 20%, for a typical data center with  $2N$  power and  $N + 1$  cooling equipment, operating at approximately 30% of rated capacity. Therefore, this value for  $\alpha$  can be considered as an example of current energy use in data centers. It is assumed that this power is also transformed into heat inside the data center:

$$P_{others\_DC}(t) = \alpha \cdot P_{load\_DC}(t) \quad (6)$$

$$Q_{others\_DC}(t) = P_{others\_DC}(t) \quad (7)$$

The power consumed by the fan  $P_{fans\_DC}$  is related to the pressure rise over the fan  $\Delta p$  that is equal to the pressure drop provided by CFD calculations, the air density  $\rho$ , and the fan efficiency  $\eta_f$  as follows:

$$P_{fans\_DC}(t) = \frac{\Delta p(t) \cdot m_{air\_total}(t)}{\eta_f \cdot \rho} \quad (8)$$

The heat dissipated is the power consumed that is not transformed in pressure energy as stated below:

$$Q_{fan\_DC} = (1 - \eta_f) \cdot P_{fans\_DC}(t) \quad (9)$$

The total load of a data center  $Q_{DC}(t)$  is determined by the air flow rate  $m_{air\_total}$ , the specific heat  $C_p$ , and the difference between the inlet and outlet air temperatures, i.e.,  $T_{R\_in}$  and  $T_{R\_out}$ , as shown

in the following equation. Note that the air flow also affects the power consumed by the fans  $P_{fan\_DC}$ , and consequently the heat generated by them inside the room  $Q_{fan\_DC}$ .

$$Q_{DC}(t) = m_{air\_total}(t) \cdot C_p \cdot (T_{R\_out}(t) - T_{R\_in}) \quad (10)$$

The cooling demand faced by the chiller  $Q_{cooling}$  includes the thermal load in the data center and the inefficiency  $\eta_{cc}$  in the coil of the CRAH:

$$Q_{cooling}(t) = \frac{Q_{DC}(t)}{\eta_{cc}} \quad (11)$$

To get the power consumption of the chiller, we have to consider a generic performance profile that is function of the condenser temperature ( $T_{co}$ ), the evaporator temperature ( $T_{ev}$ ) and the partial load (PLR). This performance is usually based in certified catalogue data from manufacturer. The following shows directly the relation between the cooling load and the power consumed in the chiller  $P_{chiller}$  by means of energy efficiency ratio (EER):

$$P_{chiller}(t) = \frac{Q_{cooling}(t)}{EER(t)} \quad (12)$$

The partial load ratio (PLR) specifies the relation between the cooling demand in a certain condition and the cooling load in nominal conditions ( $Q_{cooling\_nom}$ ), which corresponds to the operation of the chiller at the chilled water temperature  $T_{ev}$  and condenser water temperature  $T_{co}$ . In addition, PLR is also related to the cooling capacity rated  $Q_{cooling\_rated}$ , which corresponds to load of the chiller in standard condition (full load; temperature of chilled water leaving the chiller at  $7^\circ C$  and temperature of condenser water entering the chiller at  $30^\circ C$ ), in which EER is named  $EER_{rated}$ . The following shows the relations:

$$PLR(t) = \frac{Q_{cooling}(t)}{Q_{cooling\_nom}} \quad (13)$$

$$Q_{cooling\_nom} = Q_{cooling\_rated} \cdot COOL(T_{ev}, T_{co}) \quad (14)$$

$$CoolPR(t) = CoolPR(T_{ev}, T_{co}, PLR(t), EER_{rated}) = \frac{1}{EER(t)} \quad (15)$$

To calculate the chilled water temperature  $T_{ev}$ , it is necessary to know the room inlet air temperature  $T_{R\_in}$  and the minimum temperature difference  $\Delta T_{h-ex}$  on the coil of CRAH between the output of air and the inlet of water:

$$T_{ev} = T_{R\_in} - \Delta T_{h-ex} \quad (16)$$

Common values of  $\Delta T_{h-ex}$  on the commercial coils are between  $5^\circ C$  and  $15^\circ C$ . Since the chiller performance is also affected by  $T_{R\_in}$ , higher operation temperature in the data center room will need less power consumption from the chiller, and hence increase the cooling efficiency.

Figures 9 and 10 show the performance of cooling models. It is presented the relation of power consumption in cooling devices (chiller, fans) with ambient temperature, inlet air room operation temperature and partial load. Also PUE3 (see Section 5) is shown.

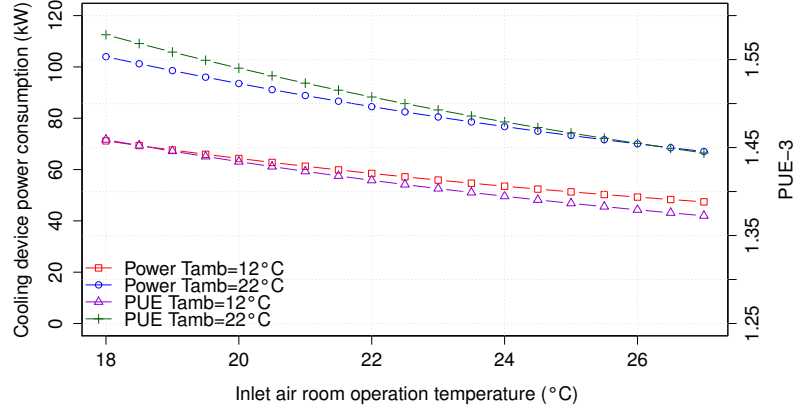


Figure 9: Cooling devices (chiller and fans) power consumption and PUE3 dependence of ambient temperature and inlet air room operation temperature (Maximum IT load 274 kW; Rated chiller cooling capacity (30°C outside air) 250 kW)

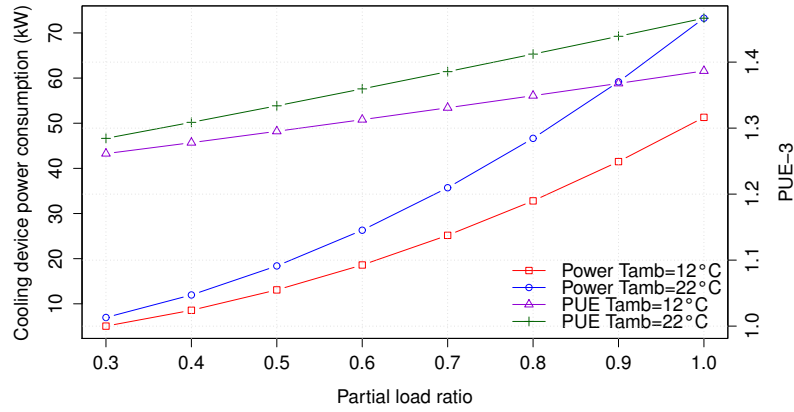


Figure 10: Cooling devices (chiller and fans) power consumption and PUE3 dependence of ambient temperature and partial load (Maximum IT load 274 kW; Rated chiller cooling capacity (30°C outside air) 250 kW)

## 5. Energy Efficiency Metrics

CoolEmAll uses a set of metrics at different levels of analysis defined Section in 4.1. Different metrics have been selected depending on the level at which the experiments are conducted and the purpose of assessment. The following classify the metrics considered:

- *Resource usage metrics* refer to the utilization of a certain resource (CPU, memory, bandwidth, storage capacity, etc), concerning a component (node) or a set of components (node-group, rack).
- *Heat-aware metrics* take temperature as the main indicator for the behavior of a data center.
- *Energy-based metrics* are defined as the consumption of power along a period of time.
- *Impact metrics* that are used to assess the performance of data center in environmental and economic terms.

The complete set of metrics defined in CoolEmAll was described in the public report of the project [44] as well as in some articles [45, 46]. To assess the impact of different strategies used in the simulations conducted in this paper, the following ones are selected: Total energy consumed, Power usage effectiveness (PUE), Productivity, Energy wasted ratio, Carbon emissions, Electricity costs. The following defines these metrics.

Total Energy Consumption (in Wh): This corresponds to the total power consumed by the data center over a certain period of time.

$$E_{DC} = \int_{t_1}^{t_2} P_{DC}(t)dt \quad (17)$$

Productivity: This metric indicates the relation between the useful work ( $W_{DC}$ ) in the data center and the energy required to obtain this useful work during a certain period of time. Useful work [47] identifies the measurable work done by a data center while providing a given service. Useful work is defined on the application level and depending on the application purpose it might be expressed by the number of floating-point operations, number of service invocations, number of transactions, etc.

$$Productivity = W_{DC}/E_{DC} \quad (18)$$

Power Usage Effectiveness (PUE): As defined by The Green Grid [48], this metric (defined as  $PUE_3$ ) is the ratio of the total power consumption in the data center and the power used by the IT equipment. It can be defined at an instantaneous point in time or at the aggregated level over a period of time (in terms of energy).

$$PUE_3 = E_{DC}/E_{IT} \quad (19)$$

In the framework of CoolEmAll and to assess the impact of load management with fans that will stop when they are not used, another level of PUE (referred to as  $PUE_4$ ) is defined, where the consumption of fans in racks is excluded from the IT load. For practical monitoring of this metric, it should be necessary to have separated power meters for fans or a signal to detect its operation mode with an assumption of the fan power consumption. The formula to calculate this metric is expressed as follows:

$$PUE_4 = E_{DC}/(E_{IT} - E_{fans} - E_{PSU}) \quad (20)$$

Energy Wasted Ratio (EWR): This metric assesses how much energy is wasted and is not used for producing useful work.

$$EWR = E_{DC\_not\_useful\_work}/E_{DC} \quad (21)$$

Carbon Emission (in kgCO<sub>2</sub>): This metric converts the total power consumed to CO<sub>2</sub> emissions using carbon emissions factor (CEF), which depends on the country since it is a function of the participation of the different energy sources and technologies (carbon, nuclear, natural gas, wind, hydro, solar, biomass, etc) in the total electricity generation and the efficiency of conversion.

$$Carbon\ Emission = E_{DC} \cdot CEF \quad (22)$$

Electricity Cost (in €): This metric is calculated by multiplying the total energy consumed by the price of electricity.

$$Electricity\ Cost = E_{DC} \cdot Electricity\ Price \quad (23)$$

## 6. Resource Management and Scheduling Policies

In the scope of resource management and scheduling policies, we can usually distinguish three basic components they consist of. These components include scheduling, resource allocation, and resource management. Scheduling is responsible for defining the order of execution for the ready tasks. Resource allocation selects the specific resource(s) for each job to be executed. Finally, resource management means the configuration of the resource states, usually related to their energy efficiency. Quite commonly, these components form the separated phases of various policies as presented in Figure 11. In the following subsections, we will present strategies classified with respect to the convention described above.



Figure 11: Components of workload and resource management policies

### 6.1. Scheduling Algorithms

A scheduling algorithm specifies the order in which tasks are served during the scheduling process (alternatively - it defines the order in which tasks are placed in the queues). The following shows some widely used algorithms, which can be applied to the scheduling of tasks in data centers.

- First Come First Served (FCFS) - a basic scheduling policy in which tasks are served in the order of their arrival in the system. This strategy reduces the waiting time of tasks.
- Last Come First Served (LCFS) - a policy contrary to FCFS, in which the tasks that arrive at the system later are scheduled first.
- Largest Job First (LJF) - tasks are scheduled in order of decreasing size, wherein the size specifies the number of requested processors. The main aim of this strategy is to optimize the utilization of the system.
- Smallest Job First (SJF) - tasks are ordered according to the number of requested processors. This strategy increases the throughput of the system.

The aforementioned scheduling algorithms can be extended with one of the backfilling approaches [49, 50], which exchange the positions of the jobs in the queue based on the availability of the resources and the priorities of the tasks.

- Conservative Backfill - allows a lower priority task to run only if it does not delay any of the higher priority waiting tasks.
- Aggressive Backfill - allows a lower priority task to run if it does not delay the highest priority task.
- Relaxed Backfill - allows a lower priority task to run if it does not delay the highest priority task in the manner that does not exceed a predefined factor.



## 6.2. Resource Allocation Strategies

Resource allocation strategies define the manner in which tasks are assigned to resources. Since tasks are submitted by different users over time, the decision of where to execute each arriving task are usually made in an online manner without knowledge of the future task arrivals. First, we describe three basic allocation strategies that are commonly used to balance the loads of different computing nodes in the system.

- Random - each task is assigned to a randomly chosen node.
- Round-Robin - the tasks are assigned to the nodes in a round robin manner.
- Load Balancing - each task is assigned to a node in order to balance the overall load of the system.

While the previous strategies do not explicitly consider any objective related to the tasks, the following describes three greedy allocation strategies that performance-aware, energy- and thermal-aware, respectively.

- Execution Time Optimization (ExecTimeOpt) - each task is assigned to a node that minimizes its execution/response time.
- Energy Usage Optimization (EnergyOpt) - each task is assigned to a node that minimizes the energy consumed by the task.
- Maximum Temperature Optimization (MaxTempOpt) - each task is assigned to a node that leads to the lowest maximum outlet temperature.

energy consumption. The following describes three thermal-aware allocation strategies that further take cooling into account by considering the temperatures at the server outlets.

The aims of the above three strategies are to minimize the average task response time, the overall energy consumption, and the maximum outlet temperature. For the thermal-aware strategy, the maximum outlet temperature is used as an objective because it has been shown to directly impact the cooling cost of data centers in both homogeneous and heterogeneous environments [51, 52].

Finally, tasks can also be assigned to resources in order to consolidate the workload in a predefined allocation manner. The following describes some consolidation strategies.

- High performance - tasks are assigned to nodes starting from high performance ones.
- Low power - tasks are assigned to nodes starting from low power ones.
- Location-aware - tasks are assigned to nodes with respect to their physical locations.

Depending on the implemented scheduling model (single or multi-level), the presented resource allocation strategies might have different impact on the final allocation of the resources. In case of scheduling at the RECS level, the above strategies are responsible for assigning tasks directly to the nodes with respect to their resource requirements. For scheduling at the room level, a scheduler has to first choose the rack where the task will be assigned, and then the RECS or nodes within which further allocation will be performed. In this case it is possible to mix two or more strategies by applying, for example, the location-aware strategy in order to select a rack, followed by the load balancing strategy to balance the load within the chosen rack, and finally the thermal-aware strategy to minimize the outlet temperature of the chosen RECS.

### 6.3. Resource Management Policies

Resource management policies specify a set of operations performed on the resources during the scheduling process. They usually require supports from the underlying hardware layer and their effectiveness is closely related to the managed IT equipment. The following describes two most popular policies.

- Switching nodes ON/OFF - a node is switched on or off, depending on if it is used or not.
- Dynamic Voltage and Frequency Scaling (DVFS) - the frequency of a processor is scaled up or down, depending on if the processor is used or not.

## 7. Simulations

This section presents simulations and the results obtained by using the Data Center Workload and Resource Management Simulator (DCworms) [53]. Different models of the data center components presented in Section 4 and various resource scheduling policies presented in Section 6 are evaluated.

### 7.1. Simulation Setup

*Resource Description.* In our experiments, five different types of processors are used and their technical specifications are presented in Table 3. All five types of processors were previously profiled in order to obtain their detailed power and performance characteristics. More information can be found in [54]. Moreover, to perform comprehensive studies, different processor configurations are simulated at three different levels, namely the RECS level, the rack level, and the room level. The cooling model adopts that described in Section 4.3.

Table 3: Technical specifications of the simulated processors.

Processor	Max. Frequency	RAM Memory	Number of Cores
Intel Core i7-3615QE	2.3GHz	16GB	4 (8 logical)
Intel Core i7-2715QE	2.1GHz	16GB	4 (8 logical)
Intel Atom D510	1.66GHz	4GB	2 (4 logical)
Intel Atom N2600	1.6GHz	2GB	2 (4 logical)
AMD G-T40N	1GHz	4GB	2 (2 logical)

*Benchmarks and Workloads.* Several types of benchmarks can be used to demonstrate the gains of the proposed system. The three most classical kinds of benchmarks are:

- Micro benchmarks, testing only one particular sub-system like memory accesses;
- Single-host benchmarks, usually used to test a particular host;
- Classical distributed benchmarks from the HPC community like NPB (Nas Parallel Benchmarks).

Benchmarks of the first category were used during development and tests of the monitoring infrastructure and of the application profiling tools. Benchmarks presented in this article include: **fft**, **abinit**, **c-ray**, **lin\_1gb**, **lin\_3gb**, **lin\_tiny**, **tar**, **openssl**. Specifically, **fft** is a tool to compute Fast Fourier Transforms. **abinit** is a scientific tool for electronic simulation at the atomic level. **c-ray** is a raytracing tool. **lin\_1gb**, **lin\_3gb** and **lin\_tiny** are different instances of Linpack (classical High Performance Computing benchmark). **tar** is an archive manipulation tool. Finally, **openssl** is an open-source implementation of cryptographic protocols.

*Power consumption model.* To estimate power consumption of the given processor we followed the model proposed in Section 4.2.3 supported with the gathered application profiles. We replayed the tasks execution, adjusting the frequency level and assumed linear dependency between power processor power drawn and its utilization. Our previous studies [53] show that such an approach presents reliable accuracy, with respect to the data gathered on real hardware, and might be boldly used as a power consumption estimator.

## 7.2. Simulation Results

### 7.2.1. Results for the RECS Level

This subsection shows the results of the simulations performed at the RECS level. Specifically, experiments were conducted to evaluate a system with one single RECS2.0 unit consisting of 18 processors/nodes. The following describes the processor configuration used in the experiment:

- 8 Intel Core i7-2715QE nodes.
- 4 Intel Atom D510 nodes.
- 6 AMD Fusion G-T40N nodes.

The workload consists of 1000 tasks randomly drawn from the benchmarks described in Section 7.1. Tasks arrive according to the Poisson process. The load intensity used in the simulation is proportional to the average arrival rate  $\lambda$  (#jobs/hour), and it is defined as  $\lambda/10$ . Five resource allocation strategies – Random, Round-Robin, ExecTimeOpt, EnergyOpt and MaxTempOpt – were evaluated with FCFS scheduling (which is also used in all subsequent experiments). Besides energy consumption, average response time of the jobs are used as a performance metric, and maximum outlet temperature is used as an indicator for the cooling cost. No resource management technique was applied, so all processors are switched on at all times.

Figure 12 shows the simulation results. Note that only the dynamic energy consumption is shown in the figure, since the nodes are not switched off even when they are idle, so the static part will be identical for all strategies. The simulation results confirm our intuition that ExecTimeOpt provides better average response time, EnergyOpt provides less dynamic energy consumption, and MaxTempOpt provides the lowest maximum outlet temperature. The other two strategies (especially Round-Robin) perform badly for all three metrics, since they are oblivious to the platform and workload characteristics. Moreover, a tradeoff can be observed among the conflicting objectives of performance, energy and temperature (more details concerning such tradeoff can be found in [55]). In particular, the MaxTempOpt strategy reduces the maximum outlet temperature by about 1-1.5°C under light system load. Although the difference in the outlet temperature is small, it can have a strong impact on the cost of cooling, especially when more RECS units are present in the system. The next two subsections study this more general case by applying the ON/OFF resource management policy to save more energy.

### 7.2.2. Results for the Rack Level

This subsection shows the results of the simulations performed at the rack level. Experiments were conducted to evaluate a rack consisting of three RECS2.0 units. The following shows the processor configurations used in the experiments:

- 18 Intel Core i7 nodes: 14 Intel Core i7-3615QE nodes and 4 Intel Core i7-2715QE nodes.
- 18 Intel Atom nodes: 14 Intel Atom N2600 nodes and 4 Intel Atom D510 nodes.

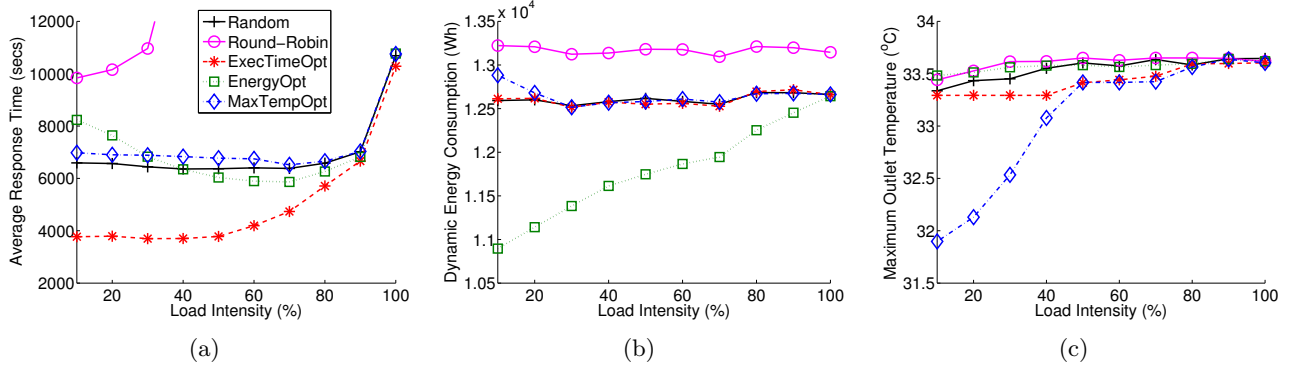


Figure 12: Performance of five resource allocation strategies at the RECS level.

- 18 AMD Fusion G-T40N nodes.

The workload contains 600 openssl tasks with a fixed load intensity. Tasks arrive according to the Poisson process with a submission time range (difference between submission of last and first task) of 2760s. Two types of consolidated resource allocation strategies – High performance and Low power – are evaluated with the ON/OFF resource management policy. The Load Balancing strategy is used as a reference for comparison.

Table 4 and Table 5 present the results according to various energy-efficiency criteria, and Table 6 and Table 7 compare the impact of studied policies on the evaluation criteria. One can see the significant improvement in terms of useful work and productivity for the consolidation on high performance nodes approach. It is obscured, however, by the increase in the scope of energy usage. An improvement on this criterion can be achieved by benefiting from the possibility of switching off unused nodes. Consolidation on high performance resources with additional power management seems to be a good trade-off between energy usage and productivity. On the other hand, Consolidation on low power CPUs can be a good approach to decrease total power usage or to increase capacity. However, it should be noted that this leads to noticeable deterioration of the performance factors.

### 7.2.3. Results for the Room Level

This subsection shows the results of the simulations performed at the room level. Experiments were conducted to evaluate a server room populated with 10 racks. The following shows the configurations of the racks:

- 5 racks equipped with 10 4-unit chassis, each chassis provides a node group containing 4 Intel Core i7 nodes (Intel Core i7-3615QE).
- 5 racks equipped with 40 1-unit chassis, each provides a node group containing 1 AMD Fusion G-T40N node.

The following are the parameters used for the cooling devices:

- Computer Room Air-handling Unit (CRAH): fan efficiency=0.6, cooling coil efficiency=0.95, deltaThEx=10.
- Chiller: max cooling capacity=10000, cooling capacity rate=40000.

Table 4: Energy-efficiency metrics for Consolidation policies at a rack level, part 1

Metrics	Policy				
	Load bal- ancing	Consolidation high perfor- mance	Consolidation high per- formance on/off	Consolidation low power	Consolidation low power on/off
Total processors energy consumption [Wh]	687	853	628	553	380
Total IT energy consumption [Wh]	1671	1838	1613	1538	1365
Total node group fans energy consumption [Wh]	412	412	145	411	246
Total rack energy consumption [Wh]	2394	2586	2021	2240	1852
Total data center fans energy consumption [Wh]	113	113	113	113	113
Total cooling device energy consumption [Wh]	423	423	423	423	423
Total other devices energy consumption [Wh]	48	52	40	45	37
Total energy consumption [Wh]	2978	3174	2597	2821	2425

Table 5: Energy-efficiency metrics for Consolidation policies at a rack level, part 2

Metrics	Policy				
	Load bal- ancing	Consolidation high perfor- mance	Consolidation high per- formance on/off	Consolidation low power	Consolidation low power on/off
Mean rack power [W]	1884	2035	1591	1764	1458
Mean power [W] 2344	2498	2044	2221	1909	
Max rack power [W]	1976	2115	1748	1866	1651
Max power [W]	2438	2579	2205	2325	2106
PUE	1.244	1.227	1.285	1.259	1.309
PUE Level 4	1.782	1.727	1.61	1.835	1.777
Energy waste rate [%]	19.03	13.31	1.22	11.78	0.59
Useful Work [UW units]	556707489	1057942787	1057942787	244375600	244375600
Productivity [UW units/Wh]	232570	409154	523517	109075	131971

- Dry cooler: deltaThDryCooler=10. Dry cooler efficiency=0.02. Details related to the cooling parameters can be found in [54].

A workload containing 6000 openssl tasks is used to drive the simulation. Tasks arrive at the system according to the Poisson process with an average inter-arrival time of 1 second. The same set of resource

Table 6: Comparison of the obtained results with reference to the Load balancing policy, part 1

Metrics	Policy				
	Load bal- ancing	Consolidation high perfor- mance	Consolidation high per- formance on/off	Consolidation low power	Consolidation low power on/off
Total processors energy consumption [%]	0	+24.28	-8.48	-19.38	-44.59
Total IT energy consumption [%]	0	+9.99	-3.47	-7.98	-18.33
Total node group fans energy consumption [%]	0	+0.02	-64.73	-0.02	-40.13
Total rack energy consumption [%]	0	+8.02	-15.58	-6.40	-22.64
Total data center fans energy consumption [%]	0	+0.02	+0.02	-0.02	-0.02
Total cooling device energy consumption [%]	0	+0.02	+0.02	-0.02	-0.02
Total other devices energy consumption [%]	0	+8.02	-15.58	-6.40	-22.64
Total energy consumption [%]	0	+6.58	-12.77	-5.25	-18.57

Table 7: Comparison of the obtained results with reference to the Load balancing policy, part 2

Metrics	Policy				
	Load bal- ancing	Consolidation high perfor- mance	Consolidation high per- formance on/off	Consolidation low power	Consolidation low power on/off
Mean rack power [%]	0	+8.00	-15.60	-6.38	-22.62
Mean power [%]	0	+6.56	-12.79	-5.23	-18.55
Max rack power [%]	0	+6.99	-11.57	-5.59	-16.48
Max power [%]	0	+5.78	-9.56	-4.62	-13.63
PUE [%]	0	-1.37	+3.30	+1.21	+5.23
PUE Level 4 [%]	0	-3.09	-9.65	+2.97	-0.28
Energy waste rate [%]	0	-30.04	-93.58	-38.09	-96.91
Useful Work [%]	0	+90.04	+90.04	-56.10	-56.10
Productivity [%]	0	+75.93	+125.10	-53.10	-43.26

allocation strategies as in the rack-level case are evaluated, again, according to the following criteria: PUE, PUE-level 4, Productivity, Energy waste rate, max IT Power and Total energy used. Table 8 summarizes the results.

According to Table 8, the consolidation policy that favors high performance nodes (Intel i7 in this case) with additional node power management outperforms other strategies with respect to the

Table 8: Assessment of resource allocation policies at room level

Policy	Metrics					
	PUE	PUE level-4	Productivity (rsa1024sign/Wh)	Energy waste rate (%)	Max. IT Power (W)	Total energy (Wh)
Load balancing	1.478	1.983	406269	42.78	22966	30275
HighPerfConsolidation	1.478	1.968	449726	25.9	23424	31649
HighPerfConsolidation + NodePowMan	1.383	1.786	534816	5.639	22027	24909
LowPowerConsolidation	1.479	1.993	391227	29.17	22318	29508
LowPowerConsolidation + NodePowMan	1.365	1.798	435131	2.77	21885	24495

evaluation criteria. Accumulating load on the most efficient (in terms of performance) nodes allows to improve both PUE-related metrics as well as the productivity factor. However, one should note the increase in maximum power consumption and total energy usage for the high performance consolidation, which should be carefully watched in terms of cooling devices capacity. Therefore, the high performance strategies are good for minimization of the total energy consumption and maximization of productivity (useful work per Joule). On the other hand, low power consolidation strategies are better in cases when power usage should be constrained. In this case power capping can be applied to save additional power. As the cooling capacity in a data center is based on the maximum power consumption of the IT infrastructure, power capping can leverage facility’s total cooling capacity (more details can be found in [56]). Besides the evaluation of the specific policies, the presented experiments demonstrate the usefulness as well as the drawbacks of the presented energy-efficiency metrics. For example, values of PUE are reasonably good, but the IT part of PUE includes constant speed fans. These fans are source of inefficiency, that is, their work when nodes are idle is a waste of energy. Thus, PUE Level-4 expresses the actual efficiency in a more adequate way. Applying node-level power management policies affects a ratio between components effectively taking part in computing to other overheads (such as useless fans work, power supply lost) and thereby improves PUE Level-4 to a greater extent than PUE. In this experiment, PUE decreases when node power management techniques are used, because the applied model assumes significant correlation between heat load and power usage of the cooling devices. In other settings, with bigger server room and cooling systems, PUE could even raise after improving efficiency of the IT part. Additional insight is provided by the proposed Energy Waste Rate (EWR), which estimates factor of energy that is wasted in the studied period. It can be easily seen that introducing power management techniques improves its value significantly.

### 7.3. Impact Assessment

The experiments presented in the previous subsection point out the potential of obtaining relevant energy savings when resource management and scheduling policies are applied, especially at the room level. These energy savings can be converted to carbon emissions and electricity cost, as it is shown in Table 9. The carbon emission factor (CEF) used to calculate the carbon emissions is 0.34 kgCO<sub>2</sub>/kWh and the electricity price used to calculate the operation costs is 0.15 €/kWh according to [57].

As we can see, the consolidation policy that favors high performance nodes increases energy consumption by 5%. However, the amount of savings reaches 18% when node power management is applied. The strategy of low performance nodes consolidation provides savings of 3% that increases until 19%

Table 9: Impact assessment of room-level resource management policies.

Policy	Total Energy (Wh)	Carbon Emissions (kgCO <sub>2</sub> )	Electricity Cost (€)	Savings (%)
Load balancing	30275	10.29	4.54	-
HighPerfConsolidation	31649	10.76	4.75	-5%
HighPerfConsolidation + NodePowMan	24909	8.47	3.74	18%
LowPowerConsolidation	29508	10.03	4.43	3%
LowPowerConsolidation + NodePowMan	24495	8.33	3.67	19%

when node power management is included. Furthermore, when extending these strategies to large-scale data centers with size bigger than this model and where the operation runs for 8760 hours per year, the amount of total carbon emissions and operation cost reduced would be substantially worthy.

## 8. Conclusion

In this paper we have presented the approaches and evaluation results of the CoolEmAll project with the aim of making data centers more energy and resource efficient. We have presented workload profiles, application-workloads, power and cooling models used in the approach. In addition, different energy efficiency metrics at different levels of analysis were proposed. Various resource management and scheduling policies, including performance, energy, thermal-aware policies and consolidation policies were presented. Simulations were conducted by using the Data Center Workload and Resource Management Simulator (DCworms) for different levels in a data center, i.e., RECS level, rack level and room level. The experiments validate the specific resource management policies proposed and the energy-efficiency metrics. In future works, CFD simulations will be conducted to validate the simulation results. We will also study resource management allocation considering heat recirculation in the data center and other resource management policies such as DVFS.

## Acknowledgement

This research is funded by the European Commission under contract 288701 through the project CoolEmAll.

## References

- [1] J. G. Koomey, Worldwide electricity used in data centers, *Environmental Research Letters* 3 (3) (2008) 034008. doi:10.1088/1748-9326/3/3/034008.
- [2] EU Joint Research Centre, Harmonised global metrics to measure data centre energy efficiency, <https://ec.europa.eu/jrc/en/news/harmonised-global-metrics-measure-data-centre-energy-efficiency-7014> (November 2012).
- [3] M. Stansberry, J. Kudritzki, Uptime institute 2012 data center industry survey, Tech. rep. (2012).
- [4] C. Belady, A. Rawson, J. Pflueger, T. Cader, Green grid data center power efficiency metrics: PUE and DCiE, Technical report, The Green Grid (2008).
- [5] D. Borgetto, H. Casanova, G. Da Costa, J.-M. Pierson, Energy-Efficient Job Placement on Clusters, Grids, and Clouds, John Wiley & Sons, Inc., 2012, pp. 163–187. doi:10.1002/9781118342015.ch6.
- [6] A. Kipp, L. Schubert, J. Liu, T. Jiang, W. Christmann, M. von dem Berge, Energy consumption optimisation in HPC service centres, in: Proc. of the Second International Conference on Parallel, Distributed, Grid and Cloud Computing for Engineering, Civil-Comp Press, Stirlingshire, UK, 2011, pp. 1–16. doi:10.4203/ccp.95.33.



- [7] A. Berl, E. Gelenbe, M. Di Girolamo, G. Giuliani, H. De Meer, M. Q. Dang, K. Pentikousis, Energy-efficient cloud computing, *Comput. J.* 53 (7) (2010) 1045–1051. doi:10.1093/comjnl/bxp080.
- [8] D. Careglio, G. Da Costa, R. Kat, A. Mendelson, J.-M. Pierson, Y. Sazeides, Hardware leverages for energy reduction in large scale distributed systems, Technical report IRT/RT-2010-2-FR, IRT, University Paul Sabatier, Toulouse, France (May 2010).
- [9] A. Ramirez, European scalable and power efficient HPC platform based on low-power embedded technology, in: International Supercomputing Conference, Leipzig, Germany, 2013.
- [10] EcoCooling Inc., <http://www.ecocooling.org>.
- [11] Colt modular data centre, <http://www.colt.net/uk/en/products-services/data-centre-services/modular-data-centre-en.htm>.
- [12] CoolEmAll project web-site, <http://www.coolmall.eu>.
- [13] Future Facilities, <http://www.futurefacilities.com/>.
- [14] CA technologies, [www.ca.com](http://www.ca.com).
- [15] Innovative Research Inc., TileFlow software, <http://inres.com/products/tileflow/>.
- [16] Romonet, <http://www.romonet.com/overview>.
- [17] Schneider Electric, <http://datacentergenome.com>.
- [18] Facebook, The open compute project foundation, <http://www.opencompute.org>.
- [19] R. Basmadjian, N. Ali, F. Niedermeier, H. de Meer, G. Giuliani, A methodology to predict the power consumption of servers in data centres, in: Proceedings of the 2Nd International Conference on Energy-Efficient Computing and Networking, e-Energy '11, ACM, New York, NY, USA, 2011, pp. 1–10. doi:10.1145/2318716.2318718.
- [20] O. Mammela, M. Majanen, R. Basmadjian, H. D. Meer, A. Giesler, W. Homberg, Energy-aware job scheduler for high-performance computing, *Computer Science - Research and Development* 27 (4) (2012) 265–275. doi:10.1007/s00450-011-0189-6.
- [21] G. Da Costa, H. Hlavacs, K. Hummel, J.-M. Pierson, Modeling the energy consumption of distributed applications, in: I. Ahmad, S. Ranka (Eds.), *Handbook of Energy-Aware and Green Computing*, Chapman & Hall, CRC Press, 2012, Ch. 29, pp. 0–0.  
URL <http://www.crcpress.com/product/isbn/9781466501164>
- [22] M. Witkowski, A. Oleksiak, T. Piontek, J. Weglarz, Practical power consumption estimation for real life HPC applications, *Future Generation Computer Systems* 29 (1) (2013) 208 – 217, including Special section: AIRCC-NetCoM 2009 and Special section: Clouds and Service-Oriented Architectures. doi:10.1016/j.future.2012.06.003.
- [23] T. Mukherjee, A. Banerjee, G. Varsamopoulos, S. K. S. Gupta, Model-driven coordinated management of data centers, *Computer Networks* 54 (16) (2010) 2869–2886. doi:10.1016/j.comnet.2010.08.011.
- [24] R. Raghavendra, P. Ranganathan, V. Talwar, Z. Wang, X. Zhu, No “power” struggles: Coordinated multi-level power management for the data center, *SIGARCH Comput. Archit. News* 36 (1) (2008) 48–59. doi:10.1145/1353534.1346289.
- [25] A. Shah, N. Krishnan, Optimization of global data center thermal management workload for minimal environmental and economic burden, *Components and Packaging Technologies*, *IEEE Transactions on* 31 (1) (2008) 39–45. doi:10.1109/TCAPT.2007.906721.
- [26] S. Yeo, H.-H. Lee, SimWare: A holistic warehouse-scale computer simulator, *Computer* 45 (9) (2012) 48–55. doi:10.1109/MC.2012.251.
- [27] E. Volk, D. Rathgeb, A. Oleksiak, CoolEmAll – optimising cooling efficiency in data centres, *Computer Science - Research and Development* 29 (3-4) (2014) 253–261. doi:10.1007/s00450-013-0246-4.
- [28] Christmann, Description for resource efficient computing system (RECS), <http://shared.christmann.info/download/project-recs.pdf> (2009).
- [29] Christmann, <http://www.christmann.info/>.
- [30] M. vor dem Berge, G. Da Costa, M. Jarus, A. Oleksiak, W. Piatek, E. Volk, Modeling Data Center Building Blocks for Energy-efficiency and Thermal Solutions, Springer, 2013.
- [31] D. G. Feitelson, Workload modeling for computer systems performance evaluation, <http://www.cs.huji.ac.il/~feit/wlmod/>, version 1.0.2 (Sep 2014).
- [32] G. L. T. Chetsa, L. Lefevre, J.-M. Pierson, P. Stolf, G. Da Costa, DNA-inspired scheme for building the energy profile of HPC systems, in: Proceedings of the First International Conference on Energy Efficient Data Centers, E2DC'12, Springer-Verlag, Berlin, Heidelberg, 2012, pp. 141–152. doi:10.1007/978-3-642-33645-4\_13.
- [33] S. Chapin, W. Cirne, D. Feitelson, J. Jones, S. Leutenegger, U. Schwiegelshohn, W. Smith, D. Talby, Benchmarks and standards for the evaluation of parallel job schedulers, in: D. Feitelson, L. Rudolph (Eds.), *Job Scheduling Strategies for Parallel Processing*, Vol. 1659 of Lecture Notes in Computer Science, Springer Berlin Heidelberg, 1999, pp. 67–90, up-to-dated version of the standard workload format (SWF) definition available in <http://www.cs.huji.ac.il/labs/>

- parallel/workload/swf.html. doi:10.1007/3-540-47954-6\_4.
- [34] V. M. Lo, J. Mache, K. J. Windisch, A comparative study of real workload traces and synthetic workload models for parallel job scheduling, in: *Proceedings of the Workshop on Job Scheduling Strategies for Parallel Processing, IPPS/SPDP '98*, Springer-Verlag, London, UK, 1998, pp. 25–46.
  - [35] D. Feitelson, Packing schemes for gang scheduling, in: D. Feitelson, L. Rudolph (Eds.), *Job Scheduling Strategies for Parallel Processing*, Vol. 1162 of *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, 1996, pp. 89–110. doi:10.1007/BFb0022289.
  - [36] A. B. Downey, A parallel workload model and its implications for processor allocation, *Cluster Computing* 1 (1) (1998) 133–145. doi:10.1023/A:1019077214124.
  - [37] B. P. K. Psounis, P. Molinero-Fernandez, F. Papadopoulos., Systems with multiple servers under heavy-tailed workloads, *Performance Evaluation* 62 (2005) 456474.
  - [38] X. Fan, W.-D. Weber, L. A. Barroso, Power provisioning for a warehouse-sized computer, *SIGARCH Comput. Archit. News* 35 (2) (2007) 13–23.
  - [39] L. F. Cupertino, G. Da Costa, J.-M. Pierson, Towards a generic power estimator, *Computer Science - Research and Development* (2014) 1–9doi:10.1007/s00450-014-0264-x.
  - [40] M. Jarus, A. Oleksiak, T. Piontek, J. Węglarz., Runtime power usage estimation of HPC servers for various classes of real-life applications, *Future Generation Computer Systems* 36 (2014) 299 – 310. doi:http://dx.doi.org/10.1016/j.future.2013.07.012.
  - [41] G. Da Costa, H. Hlavacs, Methodology of Measurement for Energy Consumption of Applications, in: *Grid Computing (GRID)*, 2010 11th IEEE/ACM International Conference on, 2010, pp. 290–297. doi:10.1109/GRID.2010.5697987.
  - [42] N. Rasmussen, Calculating total cooling requirements for data centers, White paper 25 rev 3, Schneider Electric's Data Center Science Center, [http://www.apcmedia.com/salestools/NRAN-5TE6HE/NRAN-5TE6HE\\_R3\\_EN.pdf](http://www.apcmedia.com/salestools/NRAN-5TE6HE/NRAN-5TE6HE_R3_EN.pdf).
  - [43] Guidelines for energy-efficient datacenters, White paper, The Green Grid (2007).
  - [44] L. Sisó, R. Fornós, A. Napolitano, J. Salom, D5.1. white paper on energy and heat-aware metrics for computing modules, Tech. rep., CoolEmAll, v1.4 7th FP GA no 288701 (2012).
  - [45] L. Sisó, J. Salom, M. Jarus, A. Oleksiak, T. Zilio, Energy and heat-aware metrics for data centers: Metrics analysis in the framework of coolemall project, in: *Cloud and Green Computing (CGC)*, 2013 Third International Conference on, 2013, pp. 428–434. doi:10.1109/CGC.2013.74.
  - [46] E. Volk, A. Tenschert, M. Gienger, A. Oleksiak, L. Siso, J. Salom, Improving energy efficiency in data centers and federated cloud environments: Comparison of CoolEmAll and Eco2Clouds approaches and metrics, in: *Cloud and Green Computing (CGC)*, 2013 Third International Conference on, 2013, pp. 443–450. doi:10.1109/CGC.2013.76.
  - [47] Harmonizing global metrics for data center energy efficiency, Tech. rep., Global Taskforce Reaches Agreement Regarding Data Center Productivity (March 2014).
  - [48] V. Avelar, D. Azevedo, A. French, PUE<sup>TM</sup>: A comprehensive examination of the metric, White paper 49, The Green Grid (2012).
  - [49] S. Srinivasan, R. Kettimuthu, V. Subramani, P. Sadayappan, Characterization of backfilling strategies for parallel job scheduling, in: *Parallel Processing Workshops, 2002. Proceedings. International Conference on*, 2002, pp. 514–519. doi:10.1109/ICPPW.2002.1039773.
  - [50] J. Ward, William A., C. Mahood, J. West, Scheduling jobs on parallel systems using a relaxed backfill strategy, in: D. Feitelson, L. Rudolph, U. Schwiegelshohn (Eds.), *Job Scheduling Strategies for Parallel Processing*, Vol. 2537 of *Lecture Notes in Computer Science*, Springer Berlin Heidelberg, 2002, pp. 88–102. doi:10.1007/3-540-36180-4\_6.
  - [51] Q. Tang, S. K. S. Gupta, G. Varsamopoulos, Energy-efficient thermal-aware task scheduling for homogeneous high-performance computing data centers: A cyber-physical approach, *Parallel and Distributed Systems, IEEE Transactions on* 19 (11) (2008) 1458–1472. doi:10.1109/TPDS.2008.111.
  - [52] H. Sun, P. Stolf, J.-M. Pierson, G. Da Costa, Energy-efficient and thermal-aware resource management for heterogeneous datacenters, *Sustainable Computing: Informatics and Systems* (2014) doi:10.1016/j.suscom.2014.08.005.
  - [53] K. Kurowski, A. Oleksiak, W. Piatek, T. Piontek, A. Przybyszewski, J. Węglarz, DCworms – a tool for simulation of energy efficiency in distributed computing infrastructures, *Simulation Modelling Practice and Theory* 39 (2013) 135 – 151. doi:http://dx.doi.org/10.1016/j.simpat.2013.08.007.
  - [54] E. Volk, W. Piatek, M. Jarus, G. Da Costa, L. Siso, M. vor dem Berge, Update on definition of the hardware and software models, Deliverable D2.3.1, CoolEmAll (2013).
  - [55] H. Sun, P. Stolf, J.-M. Pierson, G. Da Costa, Multi-objective scheduling for heterogeneous server systems with machine placement, in: *Cluster, Cloud and Grid Computing (CCGrid)*, 2014 14th IEEE/ACM International Symposium on, 2014, pp. 334–343. doi:10.1109/CCGrid.2014.53.
  - [56] A. Oleksiak, W. Piatek, T. Piontek, H. Sun, E. P. P. Montanera, Second set of resource management and scheduling policies, Deliverable D4.6, CoolEmAll (2013).

- [57] R. Kemma, D. Park, Methodology study eco-design of energy-using products meeup, Tech. rep., Final report. VHK. Delft, The Netherlands. [http://ec.europa.eu/enterprise/policies/sustainable-business/ecodesign/methodology/index\\_en.htm](http://ec.europa.eu/enterprise/policies/sustainable-business/ecodesign/methodology/index_en.htm). (2005).