# Proactive and Reactive Thermal Aware Optimization Techniques to Minimize the Environmental Impact of Data Centers

Marina Zapater          José L. Ayala          José M. Moya

Data centers are easily found in every sector of the worldwide economy. They are composed of thousands of servers, serving millions of users globally and 24-7. In the last years, e-Science applications such e-Health or Smart Cities have experienced a significant development. The need to deal efficiently with the computational needs of next-generation applications together with the increasing demand for higher resources in traditional applications has facilitated the rapid proliferation and growing of Data Centers. A drawback to this capacity growth has been the rapid increase of the energy consumption of these facilities. In 2010, data center electricity represented 1.3% of all the electricity use in the world. In year 2012 alone, global data center power demand grep 63% to 38GW. A further rise of 17% to 43GW was estimated in 2013. Moreover, Data Centers are responsible for more than 2% of total carbon dioxide emissions.

The PhD Thesis described here addresses the energy challenge by proposing proactive and reactive thermal and energy-aware optimization techniques that contribute to place Data Centers on a more scalable curve. This work develops energy models and uses the knowledge about the energy demand of the workload to be executed and the computational and cooling resources available at Data Center to optimize energy consumption. Moreover, data centers are considered as a crucial element within their application framework, optimizing not only the energy consumption of the facility, but the global energy consumption of the application.

The main contributors to the energy consumption in a Data Center are the computing power drawn by IT equipment and the cooling power needed to keep the servers within a temperature range that ensures safe operation. Because of the cubic relation of fan power with fan speed, solutions based on over-provisioning cold air into the server usually lead to inefficiencies. On the other hand, higher chip temperatures lead to higher leakage power because of the exponential dependence of leakage on temperature. At the server level, this work focuses on the development of models to describe the leakage-cooling tradeoffs, and proposes strategies to minimize server energy via cooling management.

When scaling to the data center level, a similar behavior in terms of leakage-temperature tradeoffs can be observed. As
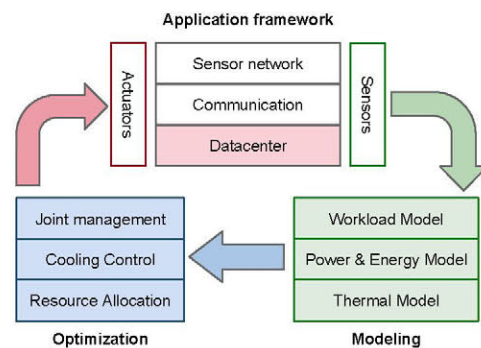


**Figure 1: Overview of the proposed analysis and optimization system**

room temperature raises the efficiency of data room cooling units improves. However, as we increase room temperature, CPU temperature raises and so does leakage power. Moreover, the thermal dynamics of a data room exhibit unbalanced patterns due to both the workload allocation and the heterogeneity of computing equipment. Because data center cooling ensures the safe operation of servers for a worst-case scenario, the maximum CPU temperature is the one limiting the minimum cooling. To leverage energy efficiency at the data center, workload allocation and cooling control strategies need to be applied together to obtain energy savings. At this scope, our goal is the proposal of thermal- and heterogeneity-aware workload management techniques that jointly optimize the allocation of computation and cooling to servers. These strategies need to be backed up by flexible room level models, able to work on runtime, that describe the system from a high level perspective.

Finally, within the framework of next-generation applications, academia has traditionally focused on minimizing the energy consumption of the sensor deployments that support this kind of applications. However, decissions taken at this scope can have a dramatical impact on the energy consumption of lower abstraction levels, i.e. the data center facility. It is important to consider the relationships between all the computational agents involved in the problem, so that they can cooperate to achieve the common goal of reducing energy in the overall system. The techniques proposed at this abstraction level aim to optimize the energy consumption of the overall application by evaluating the energy costs of performing part of the processing in any of the different abstraction layers, from the node to the data center, via workload off-loading techniques.

As summarized in Figure 1, our work proposes a global

solution based on the energy analysis and optimization for next-generation applications from a multi-layer perspective. The main chapters of this work will be detailed next.

# 1. SERVER MODEL AND OPTIMIZATION

Dynamic consumption has historically dominated the power budget. But when technology scales below the 100nm boundary, static consumption becomes much more significant. This issue is intensified by the influence of temperature on leakage current. In a modern enterprise server, CPU is the major contributor to leakage power. Because fan power is a cubic function of fan speed, the sum of leakage power and cooling power describe a convex-like curve, where a workload-dependent optimum point between cooling and leakage power can be found.

Our work in this area has focused on the development of empirical models to estimate the static and dynamic contributions to power consumption in enterprise servers, analyzing the interactions between temperature, leakage, and cooling power. We have provided an analytical model for temperature-dependent CPU leakage, splitting it from CPU dynamic power, and we have proposed a cooling management scheme able to set the optimum fan speed that minimizes the aggregate leakage plus cooling power in the server.

Experimental results on a presently shipping enterprise server demostrate the benefits of including leakage awareness. Our cooling management policy achieves up to 9% energy savings and 30W reduction in peak power when compared to the default cooling control scheme.

Our work in this area can be found in (1).

# 2. DATA CENTER ENERGY REDUCTION

Current thermal data room models are based on static and time costly CFD simulation techniques that often do not match the real dynamics. Our work proposes a methodology to model the temperature of servers and the data room dynamics in a real scenario by using meta-heuristics as a valid alternative to classical analytical techniques. Analytical models require finding the complex relationships between the parameters involved in the modeling to build analytical functions. The modeling of scalable, distributed and highly heterogeneous systems is usually unfeasible for analytical methods. Meta-heuristics are higher-level procedures that make few assumptions about the modeling problem. We propose the usage of Genetic Programming techniques to model and predict server CPU and inlet temperature. To validate the models in a real data center scenario we have developed and deployed a sensor network in the CeSViMa Data Center, of Universidad Politecnica de Madrid, where we gather server and environmental data to train and test the models.

The main goal is to propose resource management techniques that minimize the energy consumption of the facility by assigning tasks to computational resources in the most efficient way and setting the data room cooling appropriately. To this end we propose and solve a Mixed Integer Linear Programming (MILP) problem, consisting on the minimization of the energy used to perform a certain workload in the data center. The proposed approach exploits the heterogeneity of the system from a mixed static/dynamic perspective, and combines the proper selection of cores with the information retrieved during a workload characterization phase.

The static optimization approach performs an off-line configuration of the data center cluster to obtain the most suitable server set-up. Our results show that the best server set is heterogeneous, and outperforms all homogeneous server sets. The dynamic optimization distributes the tasks arriving to the data center between the heterogeneous computing resources, optimizing resource allocation. To validate

results we have collected data on servers of three different architectures: AMD, Intel and SPARC, and compared the optimum allocation to the default allocations of the open-source resource manager SLURM, obtaining from 7.5% to 24% energy savings.

The work developed at the data center scope can be found in the following references: (2) and (3).

# 3. APPLICATION FRAMEWORK

Next generation applications are usually composed of a large number of sensors, connected via wireless through a mobile processing device. Data centers provide the required infrastructure to store, analyze and process all application data. To provide adequate energy management, this heterogeneous distributed computing system is tightly coupled with an energy analysis and optimization system, which continuously adapts the amount of processing that is performed in each layer, and the resources assigned to each task.

The work at this abstraction level leverages the concept of resource management by considering not only the heterogeneity of server architectures in a data center, but the usage of the heterogeneous elements that compose the distributed application framework outside the facility, such as ARM-based smartphones. We use non-optimal lightweight distributed allocation algorithms based on Satisfiability Modulo Theories (SMT) formulas outside the facility, and combine this allocation with MILP-based problems in the data center. A SMT solver is a fast and lightweight tool that checks whether a formula satisfies a condition. Each node runs a solver to check whether a task satisfies the conditions to be executed or needs to be offloaded to the data center.

This way, computation is distributed across all nodes in the network and the data center, in a way that allows us to minimize the energy consumption from 10% to 24% in the overall application framework. This energy savings can be translated into a reduction of almost 70 tons of $CO_2$ anually for the proposed e-Health scenario.

Our work in this area can be found in (4) and (5).

# 4. SUPPORTING PAPER

The supporting paper (5) belongs to a journal publication in *Future Generation Computer Systems* (Elsevier). This publication describes the work at the server, data center and application framework levels, focusing on a case study for e-Health scenarios. Other publications on the dissertation topic are outlined next:

# References

[1] M. Zapater, J. L. Ayala, J. M. Moya, K. Vaidyanathan, K. Gross, and A. K. Coskun, "Leakage and temperature aware server control for improving energy efficiency in data centers," in *DATE'13*, 2013.

[2] J. Pagán, M. Zapater, O. Cubo, P. Arroba, V. Martín, and J. M. Moya, "A Cyber-Physical approach to combined HW-SW monitoring for improving energy efficiency in data centers," in *DCIS*, 2013.

[3] M. Zapater, J. L. Ayala, and J. M. Moya, "Leveraging heterogeneity for energy minimization in data centers," in *CC-GRID'12*, 2012.

[4] M. Zapater, C. Sanchez, J. L. Ayala, J. M. Moya, and J. L. Risco-Martín, "Ubiquitous green computing techniques for high demand applications in smart environments," *Sensors*, vol. 12, no. 8, pp. 10 659–10 677, 2012.

[5] M. Zapater, P. Arroba, J. L. Ayala, J. M. Moya, and K. Olcoz, "A novel energy-driven computing paradigm for e-health scenarios," *Future Generation Computer Systems*, vol. 34, no. 0, pp. 138–154, 2014.