# Achieving energy efficiency in data centers with a performance-guaranteed power aware routing

CrossMark

E. Baccour [a,b,*], S. Foufou [a], R. Hamila [c], Z. Tari [d]

[a] CSE department, College of Engineering, Qatar University, Doha, Qatar
[b] LE2i Lab, University of Burgundy Dijon, France
[c] EE department, College of Engineering, Qatar University
[d] School of Science, RMIT University

## ABSTRACT

Nowadays, data centers are designed to offer the highest performance in case of high traffic load and peak utilisation of the network. However, in a realistic data center environment, the peak capacity of the network is rarely reached and the average utilisation of devices varies between 5% and 25% which results into a huge loss of energy since most of the time links and servers are idle or under-utilized. The high impact of this wasted power on environmental effects, energy needs and electricity costs raised the concerns to seek for an efficient solution to make data centers more power effective while keeping the desired quality of service. In this paper, we propose a power-aware routing algorithm that saves a considerable amount of energy with a negligible trade-off on the performance of the network and a guaranteed reliability of the system. The key idea is to keep active only the vital and critical nodes participating in the communication traffic and ensuring the reliability while the unneeded devices are turned-off. Vital nodes between clusters (parts of the network) are calculated only once during the initialization of the system and consequently used with a constant time complexity. Besides its short computation time, our routing algorithm guarantees over 50% of energy saving by maintaining the minimum number of needed devices and over 20% when adding backup routes. This power efficiency is accompanied by a guaranteed performance and reliability against failures.

© 2017 Elsevier B.V. All rights reserved.

## 1. Introduction

Scaling from server rooms that support small companies, to enterprise infrastructures that control the medium sized organizations, to servers farms that host cloud computing services offered by Google, Facebook, Microsoft and Amazon, data centers present the backbone of all digital contents, new technologies, e-commerce, social media, video streaming and others [1]. Thus, the dependence on data centers to provide all our daily activities made it one of the fastest consumers of energy in the world [2]. While most of the research attention focuses on how to make data centers scalable and efficient to offer peak performance and rapid services to the users, many concerns are raised about the huge amount of energy consumed by this gigantic infrastructure. Statistics in [3] and [4] show that the power consumed by all data centers in the world in 2011 accounts for 1.1%-1.5% from the world power consumption and it will increase to 8% by 2020. As an ac-

tion, many efforts have been made to design a green data center network from different aspects: Some companies such as Google and Facebook are saving the energy by using renewable energy [5,6]. Others are focusing on the computing part and how to make voltage or frequency adjusted depending on the CPU workload [7,8]. In addition, the virtualization of network resources and the consolidation in physical machines seem to be promising for conserving the energy [9].

It is noted also that a typical data center operates only at 5% to 25% of its maximum capacity depending on the duration where the utilisation of the network fluctuates according to the incoming loads [10,31]. If the load is low, the servers are kept idle and they may still consume up to 70% of their peak power which causes a great waste of energy [11]. This situation is worsened by the traditional routing algorithms that do not take into consideration the non-proportionality between the traffic load and the energy consumption. Based on these observations, we can deduce that when the network is under-utilized, the traffic can be satisfied by only a subset of devices and the idle ones can be powered off to save the energy.

* Corresponding author.
*E-mail address:* ebaccour@qu.edu.qa (E. Baccour).

In this context, many proposals are trying to make the power consumption proportional to the traffic workload. The idea is to calculate first the best routes in terms of demand satisfaction [10], bandwidth allocation [12] or throughput [13]. Then, the nodes impacted in these routes are kept active and the others are disabled. However, even though these schemes are solving the problem of non-proportionality of the network and conserving a large amount of energy, they still suffer from long computation time and complexity of routes calculation which can affect the traffic delivery delay and consequently the network performance. Based on these challenges, the aim of this paper is to propose a power aware routing algorithm that contributes in saving the energy while offering a short computation time. The proposed algorithm suggests not to calculate the best routing paths in real time when receiving the traffic pattern. Instead, vital nodes between any two clusters communicating in the network are calculated using different methods (betweenness, closeness, degree, etc). These nodes are pre-calculated once when conceiving the network and used directly with a constant computation complexity. At a given time $t$, when receiving the traffic matrix, only the vital nodes are kept active. Regarding the fault tolerance of the system, we applied the DFS (depth-first search) algorithm to identify nodes causing failures then we calculated vital nodes in the backup routes to ensure the robustness against faulty conditions. The main contributions of this paper can be summarized as follows:

1. The formulation of the energy saving problem and proposing a minimization solution.
2. The calculation of vital nodes participating in the best routing paths between different clusters of the network and the ones participating in the backup routes.
3. The description of the power-aware routing algorithm aiming at maximizing power saving and minimizing the trade-off in terms of network performance, computation complexity and reliability.
4. The implementation and evaluation of the performance of the power-efficient algorithm in the context of the Flatnet [14] and PTNet [15] data center networks under various conditions.

The rest of the paper is organized as follows: Section 2 gives an overview about the most promising existing works in the area of energy efficiency in data centers. In Section 3, we provide the formulation of the proposed energy-aware optimization technique. In the Section 4, we describe the calculation of the vital nodes and the use of the DFS algorithm to maintain the robustness of the system. Section 5 introduces the power aware routing algorithm. A detailed experimental evaluation is provided in Section 6 then a further discussion and conclusions are drawn in Sections 7 and 8 respectively.

## 2. Related works

Large-scale data centers (Google, Amazon, Microsoft) and even small data centers (companies, universities, government establishment) suffer from energy consumption issues. Therefore, considerable efforts and investigations to achieve a green data center have been conducted in both industry and academia. In fact, by minimizing the energy consumption, the network cost can be reduced as well as the impact on the environment. A short review about the most promising approaches in this field is summarized:

### 2.1. Usage of renewable sources of energy

This approach exploits the green resources such as water, wind, solar energy, pumps and heat to reduce the budget related to energy. For example, authors in [16] proposed to make all data center devices almost powered entirely by the renewable energy. To realize a testbed, they added temperature and humidity sensors, solar powered systems and wind power systems to a data center. Using this method, they obtained a successful results and witnessed the launch of several international collaborations including the US Greenlight and the New York State Energy Research. Another proposal to introduce the green energy in the data center is the net-zero networks [17]. This new concept consists of producing an amount of energy per day that is equal to the same amount of energy consumed by the network. However, the renewable energy sources used to power the network are limited by several factors such as the weather conditions and the location of the data center, in addition to the high cost of the infrastructure to generate the power and deliver it.

### 2.2. Energy-aware hardware

This method focuses on making the energy consumed by switches and servers proportional to the traffic load. To achieve this, some specific techniques are applied such as vary-on/vary-off (VOVO). Workloads in VOVO are concentrated in a subset of servers instead of distributing it across all servers which ensures more energy efficiency. Another technique is the Dynamic Voltage (DVFS) [18] where the CPU power is adjusted depending on the arriving load. The idea is based on the fact that the power in the chip is proportional to $V^2.f$, where $V$ is the voltage level and $f$ is the frequency. This ensures that the power consumption of a chip is proportional to the workload. However, these approaches optimize only the energy consumed by the CPU while the remaining network devices are untouched and are still working on their normal energy level.

### 2.3. Energy-aware architectures

This relates to the design of a data center that conserves energy thanks to its architecture. NovaCube [19] and flattened butterfly topology [20] are designed to interconnect servers directly without an intermediate switches (switchless networks). These topologies save energy consumed by switches, racks and associated cooling machines. Another example is the Nano data centers [21], where services and content are stored in the home gateways instead of being stored in data centers. To access to the content in the gateways, a P2P infrastructure is used. This approach can only be applied to a small sized network and needs more investigation to be feasible in a real data center network. Wireless data center topologies [22,23] are also an attempt to minimise the energy consumed. In fact, by relying on transceivers with a minor energy consumption compared to switch interfaces [24], communications are delivered without wasting a huge amount of energy.

### 2.4. Virtualization

The virtualization technology is based on creating multiple virtual machine (VM) instances on a single physical server [9,25]. VMs can act like a real machine with an operating system. Thus, by powering only one device while multiple machines are processing, the amount of hardware in use is reduced, the resource utilisation is improved and the energy consumed is optimized. Recently, virtualization tools are available to test and use by some vendors such as VMware [26]. GreenCloud [27] is one of the optimisations used for migration and placement of virtual machines in data centers. In this project, the authors studied the workloads of the applications in cloud networks. Depending on the workloads, they proposed algorithms for mixing and mapping virtual machines to the existing resources while considering the QoS (Quality of Services) requirements and energy efficiency. The GreenCloud approach consists of two phases. The first one deals with the collection of new