



# Renewable-aware geographical load balancing of web applications for sustainable data centers



Adel Nadjaran Toosi<sup>a,\*</sup>, Chenhao Qu<sup>a</sup>, Marcos Dias de Assunção<sup>b</sup>, Rajkumar Buyya<sup>a</sup>

<sup>a</sup> Cloud Computing and Distributed Computing Laboratory, School of Computing and Information Systems, The University of Melbourne, Australia

<sup>b</sup> Inria, ENS de Lyon, France

## ARTICLE INFO

### Keywords:

Geographical load balancing  
Web applications  
Renewable energy  
Green computing  
Wikipedia  
Auto-scaling  
Brown energy  
Green Energy  
Cost saving  
System prototype

## ABSTRACT

The ever-increasing demand for web applications deployed across multiple data centers results in large electricity costs for service providers and significant impact on the environment. This has motivated service providers to move towards more sustainable data centers powered by renewable or green sources of energy, such as solar or wind. However, efficient utilization of green energy to service web applications is a challenging problem due to intermittency and unpredictability of both application workload and renewable energy availability. One possible solution to reduce cost and increase renewable energy utilization is to exploit the spatio-temporal variations in on-site power and grid power prices by balancing the load among multiple data centers geographically distributed. In this paper, we propose a framework for reactive load balancing of web application requests among Geo-distributed sustainable data centers based on the availability of renewable energy sources on each site. A system prototype is developed, its underlying design and algorithms are described, and experiments are conducted with it using real infrastructure (Grid'5000 in France) and workload traces (real traffic to English Wikipedia). The experimental results demonstrate that our approach can reduce cost and brown energy usage with efficient utilization of green energy and without a priori knowledge of future workload, availability of renewable energy, and grid electricity prices.

## 1. Introduction

Data centers are known to be consuming enormous amount of power leading to high operational cost and high carbon footprint on the environment. According to a report from NRDC<sup>1</sup> (NRDC, 2014), in 2013, US data centers alone consumed 91 billion kilowatt-hours of electricity, equivalent to two-year power consumption of all households in New York city. This is projected to increase to roughly 140 billion kilowatt-hours and is responsible for the emission of nearly 150 million metric tons of carbon dioxide per annum in 2020. These costs and environmental concerns have prompted service providers to reduce their energy consumption and their dependence on power generated from fossil fuels (i.e., Brown energy).

Large companies (e.g., Google,<sup>2</sup> Microsoft<sup>3</sup> and Amazon<sup>4</sup>) are working towards sustainable data centers by using renewable energy

sources (i.e., Green energy) and making direct investments in on-site green power generation. *Photovoltaic solar panels* that directly convert sunlight into electricity and *wind turbines* that capture wind energy and turn it into electricity are among the most popular on-site power sources used by contemporary data centers. For example Amazon Web Services (AWS) is building a wind farm that will be operational by late 2016 and generate 40% of its electrical usage.<sup>5</sup>

Powering data centers entirely with renewable energy sources, unlike brown energy, is challenging due to the intermittent and unpredictable availability of wind and solar energy. For example, photovoltaic (PV) solar energy is only available during the day time and the amount of power produced depends on the weather and geographical location of the data center. To mitigate this variability, besides on-site renewable energy sources, service providers end up using grid power or brown energy as a backup in their data centers.

\* Corresponding author.

E-mail addresses: [anadjaran@unimelb.edu.au](mailto:anadjaran@unimelb.edu.au) (A. Nadjaran Toosi), [cqu@student.unimelb.edu.au](mailto:cqu@student.unimelb.edu.au) (C. Qu), [marcos.dias.de.assuncao@ens-lyon.fr](mailto:marcos.dias.de.assuncao@ens-lyon.fr) (M.D. de Assunção), [rbuyya@unimelb.edu.au](mailto:rbuyya@unimelb.edu.au) (R. Buyya).

<sup>1</sup> Natural Resources Defense Council, <http://www.nrdc.org>.

<sup>2</sup> <http://www.google.com.au/green/energy/>.

<sup>3</sup> <http://www.microsoft.com/environment/renewable.aspx/>.

<sup>4</sup> <http://aws.amazon.com/about-aws/sustainable-energy/>.

<sup>5</sup> <http://www.reuters.com/article/2015/07/14/us-amazon-iberdrola-idUSKCN0PO1PF20150714>.

However, to minimize brown energy usage, they need to obtain the highest possible renewable energy utilization.

Even though data centers can store power generated by renewable sources of energy in batteries for later needs, this approach has many problems (Goiri et al., 2012). For example, 1) batteries lose energy due to internal resistance and self-discharge, 2) battery-related costs can dominate the cost of power systems, and lastly 3) batteries use chemicals that are harmful to the environment. Given aforementioned problems, the best way to take full advantage of the available green energy is to match the energy demand to supply.

There is a large number of studies illustrating the potential of using “geographical load balancing” in reducing brown energy usage and accordingly maximizing renewable energy utilization (Lin et al., 2012; Liu et al., 2011; Adnan et al., 2012). Geographical load balancing (GLB) allows for “follow the renewables” by utilizing resources from geographically distributed data centers (Liu et al., 2011). Additionally, it routes the load to places with lower electricity prices even if renewable power is fully utilized or not available. This eventually leads to significant cost savings.

Among different types of applications, web-applications are highly popular and widely adopted these days. Web applications are ideal for geographical load balancing as they can quickly adapt to changes in the demand and their mostly small sized requests can be easily redirected among multiple data centers. Using multi-tier clustered web server architectures, web applications are able to efficiently allocate resources within and among data centers according to time varying demand and renewable energy availability. In fact, if the Quality of Service (QoS) requirement of web requests in terms of response time can be satisfied, the load balancing algorithm distributes requests among targeted data centers so that overall renewable energy usage is maximized. Please note that web application workload is not deferrable that means every time a request is received, the response should be generated immediately afterwards. Therefore, the load balancing technique must make real-time scheduling of the load which does not delay the current requests.

Fig. 1 depicts an architectural overview of multiple sustainable data centers and a green load balancer. The *power infrastructure* generates and delivers power for the IT equipment and cooling facilities of the data center using grid power and local renewable power supplies. Data centers host resources for the web application and the *load balancer* is

responsible for sharing requests according to the availability of renewable energy sources on each site. In this paper we aim to address the interesting and challenging question posed by this architecture: “Without a priori knowledge of the future demand, dynamic and unpredictable nature of renewable energy sources, and electricity prices, how can the load balancer distribute web application requests among multiple data centers so that the overall renewable energy usage is maximized and the total cost of power consumption is minimized?”

To address this problem, we propose a framework for reactive load balancing of web application requests among multiple geographically distributed data centers based on the availability of renewable energy sources on each site. We develop a prototype, detail its underlying design and algorithms, describe technical aspects of that, and experiment with it using real life infrastructure (i.e., Grid'5000 in France) and workload traces (i.e., real traffic to English Wikipedia). We also model renewable energy power generation using real traces of meteorological data for wind and solar radiation in the location of each data center. Using the prototype, we evaluate the optimization techniques and demonstrate that our simple, yet practical, approach can achieve significant cost savings without advance knowledge of future demands, availability of renewable energy, and electricity prices. The proposed load balancing algorithm is triggered periodically and collects the current available renewable power and electricity price at each data center. Then, based on the present rate of requests at the load balancer, it adjusts the load distribution among data centers. By using this technique, load distribution among data centers is adapted to the dynamic and varying renewable power and electricity prices.

The load balancing technique proposed in this paper is an *online algorithm* that acts without future knowledge of demand, renewable energy availabilities, and electricity prices. Online algorithms have been previously used in the literature to tackle the problem of geographical load balancing (Lin et al., 2011; Adnan et al., 2012). However, since we focus on the implementation aspects of the system, we do not provide the analytical reasoning on the *competitive analysis* of the proposed online algorithm.

One of the unique features of our work is that in contrast to the majority of other studies evaluating their system performance through simulations and analytical reasoning; our experiments are conducted in a real testbed with realistic workload traces. Moreover, results of experiments are generated based on the fine-grained measurements of

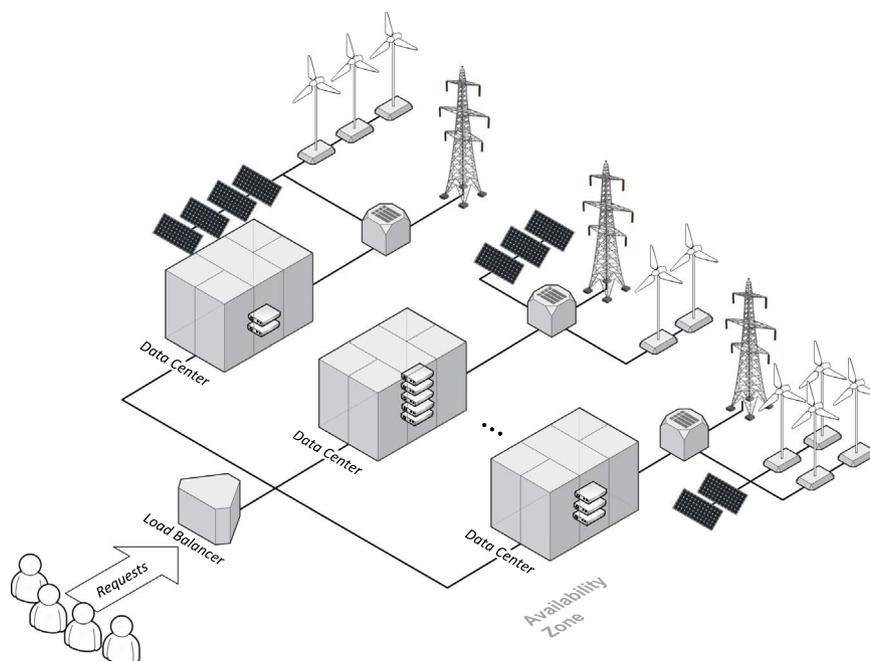


Fig. 1. Architectural overview of sustainable data centers and green load balancing.

power consumption using real-time probes provided by a live monitoring system.

The main contributions of the paper are as follows:

- A reactive load balancing algorithm to distribute a web application load among different data centers in a region to maximize on-site renewable energy utilization at each data center and to minimize overall cost. Our proposed method has a linear order of complexity and does not require any future knowledge of demands, availability of renewable energy, and electricity prices. This removes the need for any prediction component and its simplicity is very appealing in practice.
- Design and implementation of a two-layered load balancing prototype by extending “3-tier architecture” common to web applications.
- Evaluation and validation of the proposed load balancing algorithm using the developed system prototype in a real testbed with realistic workload traces. Results of experiments are generated based on the real-time measurements of actual power consumption. Meteorological data in the location of each data center used to model renewable power generation.

The rest of the paper is organized as follows: Section 2.1 motivates our work. We describe the system architecture in Section 2.2. Section 2.3 presents a detailed discussion on the design and implementation of the system. Section 3 proposes our load balancing and auto-scaling algorithms. The performance evaluation of the system is presented in Section 4. First, we describe the testbed setup, the benchmark algorithm, traces of workload, renewable power and electricity prices in Section 4.1. Then, benchmark policies and the experimental results are discussed in Sections 4.2 and 4.3, respectively. Section 5 discusses related work. Finally, our conclusions are presented in Section 6.

## 2. Load balancing in sustainable data centers

### 2.1. Motivation and background

Internet-scale distributed systems such as web applications use resources provided by geographically dispersed data centers with hundreds of thousands of physical nodes. The enormous and growing energy demands of such data centers have motivated constructing sustainable data centers for both economical and environmental reasons. Renewable sources of energy (e.g. wind and solar) have the potential to play an important role in such sustainable data centers. However, the intermittent and variable nature of renewable energy sources, caused by being heavily dependent on weather conditions, prevents them from being used as primary power supplies for data centers. Essentially, a data center must be operational even when renewable energy is not available. This is possible by feeding the data center with brown energy from the grid in case the local renewable energy is insufficient. Data centers often contract with power companies to pay variable brown electricity prices for their excess usage provided by the utility grid.

Web application providers using resources from sustainable data centers (i.e., those equipped with on-site renewable energy facilities in our definition) in a specific *region* can distribute load by following the renewable energy supply in each site to save cost and reduce carbon emissions. That is, more requests must be redirected to the places with higher availability of renewable energy. The availability of the renewable electricity at each data center might vary due to reasons such as geographical location, weather conditions, workload, and capacity of power generators. Eventually, if there is not enough available renewable power from all data centers to handle the entire load, requests must be routed to places with lower electricity price. Indeed, with widespread adoption of smart grid technologies, spatial and temporal differences in electricity prices, even in a small region, provide an opportunity for the load redirection to save cost.

In order to make this possible, we design and implement a framework for green load balancing of web requests among sustainable data centers hosting the web application and dispersed throughout a region. In this system, it is assumed that the “operational cost” of on-site power generated from renewable sources of energy at each data center is zero. This is a valid assumption as renewable power generation needs one-time installation and very low maintenance cost during the lifespan of the renewable power generators. Decision making regarding the investment on on-site renewable power generation based on its “capital cost” and “return on investment” is not in the scope of this work. Interested readers are referred to Liu et al. (2011) for feasibility of powering internet-scale systems using renewable energy and optimal portfolio of solar and wind energy mix.

Moreover, delay related changes due to geographical routing are not a concern in this work as we only consider data centers located within a certain region. In our terminology, a region is a geographic area (e.g., Europe, US East, Central Asia) with multiple and isolated data centers connected through low latency links. In the performance evaluation section, we show that for a testbed of the same characteristics, response time remains below the acceptable delay requirements of the studied web application (i.e., Wikipedia)<sup>6</sup> and will not significantly be affected due to geographical load balancing within the region.

### 2.2. System architecture

Our approach to green load balancing of web applications extends multi-tier clustered web server architecture by adding an extra layer of load balancing responsible for sharing load at the data-center level as shown in Fig. 2. A standard multi-tier web application often consists of three logical layers (Grozev and Buyya, 2014):

1. Presentation Layer – represents the interface displayed to the end-user, e.g., a web page viewed on a web browser.
2. Business/Domain Layer – implements the core application logic, e.g., core web application deployed on web servers.
3. Data Layer – handles access to the persistent storage, e.g., a database server.

This layered architecture allows for software components of each layer to be deployed within single/multiple separate machines and easily scale out based on the load. For example, the number of *web servers* can increase or decrease dynamically in response to the demand. Web servers of the business/domain layer are often deployed behind a *load balancer*, which redirects the incoming requests among them. In our architecture, we introduce two level of load balancing:

1. *Local Load Balancing* – redirects requests between web servers within a data center, and
2. *Global Load Balancing* – redirects requests among local load balancers, each associated with a data center.

In this paper, we focus on load sharing at the *global load balancer* (Global-LB). Global-LB is the main entry point of the system and all incoming requests at this point are distributed among geographically dispersed data centers based on the proposed policy.

In each data center, there is an *auto-scaler* that dynamically and adaptively adds and removes web servers behind the local load balancer in response to the dynamic workload. Here, we only focus on auto-scaling of the application layer which is a common practice in the real world. Auto-scaling can be done based on statistics gathered by the monitors responsible for fetching corresponding system information such as resource utilization at web server nodes, request rates, etc. According to the obtained information, when it is necessary, the auto-

<sup>6</sup> Web-based encyclopedia project supported by the Wikimedia Foundation.

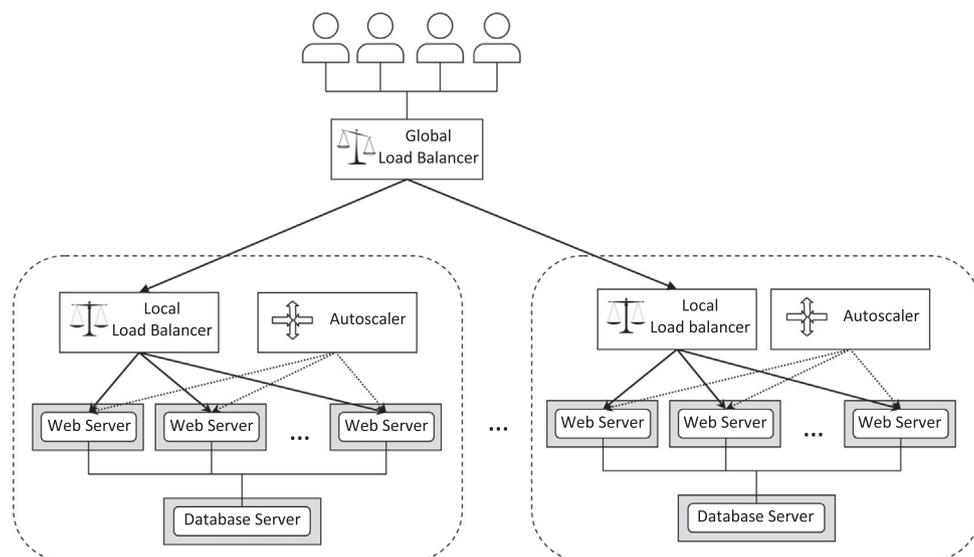


Fig. 2. Overall system architecture and related components.

scaler then makes scaling decisions based on predefined strategies and policies. In Section 2.3.2, we detail our proposed auto-scaling policy.

At the lowest level of the system architecture, Database server(s) provide(s) access to the persistent storage for the web servers. Web servers in the business/domain layer query the database servers of the data layer for required information (e.g., web page content). The data layer can be composed of both transactional relational databases, caches and novel approaches like NoSQL and NewSQL databases.

### 2.3. Design and implementation

Apart from components of multi-tier architecture for web applications, the key components of the proposed system are *global load balancer* and *autoscaler*.

#### 2.3.1. Global load balancer (Global-LB)

The function of Global-LB is to redirect incoming web requests to an appropriate data center site so that the overall renewable energy utilization is maximized. To achieve this goal, we use HAProxy<sup>7</sup> load balancer that distributes requests across local load balancers. The “weighted round robin” is one of the main load balancing algorithms used by HAProxy to determine which server, in the backend, is selected for the next incoming request. A weight parameter can be assigned to each server in the backend to manipulate how frequently the server is selected for the request routing, compared to other servers.

We designed and developed a *controller* to be run besides HAProxy, as shown in Fig. 3. The role of the controller, which is a Java program deployable on the same or separate host as HAProxy, is to assign weights to each server on a regular basis according to the current status of the system. In order to do so, the controller has a *decision making module* calculating weights based on the load balancing policies and information collected by *monitoring modules*. Monitoring modules are responsible for computing the amount of power consumed on each site and the total number of requests redirected to the site in a certain time window. To obtain the number of requests in a time period, the monitoring module queries the statistics measured by HAProxy. To compute the power consumption, the corresponding monitoring module communicates with the power monitoring APIs provided on each data center site. In our prototype, the monitoring module directly fetches the power consumption from every node using probes installed on the node. Having collected the required information, the

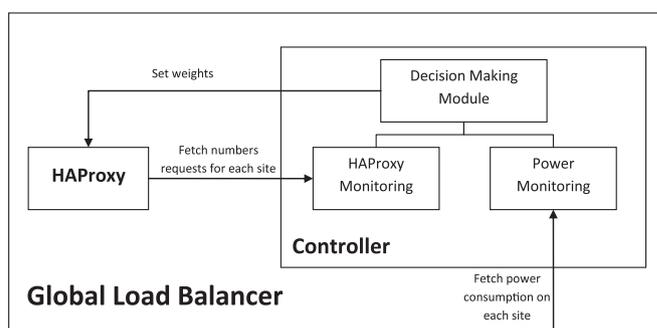


Fig. 3. Global Load Balancer.

decision making module based on the predefined load balancing policy, which will be discussed shortly, sets the weight values relative to each site for HAProxy. The load balancing policy used in our system is presented in Section 3.1.

#### 2.3.2. Auto-scaler

Each *auto-scaler* is packaged in a single JAR file and deployed on the same or a separate host as the local load balancer. The main function of an auto-scaler is to horizontally scale out web servers on demand (to allocate or deallocate web server machines) in the business/domain layer of the web application.

There is a substantial amount of work on designing auto-scaling solutions for multi-tier applications, for example (Urgaonkar et al., 2008; Jiang et al., 2013; Fernandez et al., 2014). Among different approaches, threshold-based auto-scaling methods that work based on performance metrics such as CPU and RAM utilization are among the most widely adopted techniques, e.g., AWS Auto-Scaling Service.<sup>8</sup> More resources are provisioned whenever an upper threshold is exceeded and resources are released whenever a lower threshold is reached. The main scope of this work is not to propose a new methodology of auto-scaling for multi-tier web applications. Any form of auto-scaling method can be plugged into the system without the necessity to modify other components. Nevertheless, in our prototype, we need an efficient auto-scaling mechanism to elastically scale out the web server machines in the virtual cluster once it is required.

Since in our evaluation we use homogeneous server farm<sup>9</sup> for web

<sup>7</sup> HAProxy - The Reliable, High Performance TCP/HTTP Load Balancer, <http://www.haproxy.org/>.

<sup>8</sup> Auto-Scaling - Amazon Web Services, <http://aws.amazon.com/autoscaling/>.

<sup>9</sup> A server farm made up of machines having the same characteristics.

servers in each data center, we employ a simple yet effective auto-scaling method based on profiling data collected from web server machines in each data center. Accordingly, we set an upper threshold for the request rate at or below which web server machines of the specific type can provide responses in a timely manner. Thus, the auto-scaler allocates more web server machines whenever request rate goes higher than the threshold. In [section 3.2](#), we present a detailed discussion on the auto-scaling algorithm.

### 3. Load balancing and auto-scaling policies

In this section, first, we detail our proposed load balancing algorithm used in Global-LB and then the auto-scaling algorithm employed to scale out web servers in each data centers is presented.

#### Algorithm 1. Green Load Balancing (GreenLB) Policy

```

1:  $R \leftarrow 0$ 
2: for all data centers  $d$  in the list do
3:    $c \leftarrow$  Fetch the data center's energy consumption
   in Watt-hour within the time window
4:    $t \leftarrow$  Fetch the number of requests redirected to the
   site within the same time window
5:    $a \leftarrow$  Fetch currently available renewable power at
   the site in Watt
6:    $w \leftarrow$  Compute Watt-hour consumption per request
   ( $c \div t$ )
7:    $r_d \leftarrow$  Compute the request rate (#reqs/hour) data
   center  $d$  can accommodate using renewables ( $a \div w$ )
8:    $R \leftarrow R + r_d$ 
9: end for
10:  $\gamma \leftarrow$  Fetch request rate (#reqs/hour) at Global-LB
11: if  $\gamma < R$ .then
12:   for all data centers  $d$  in the list do
13:     set weight as  $r_d \div R$ 
14:   end for
15: else
16:   Find the data center  $d'$  with the cheapest price of
   brown energy per request.
17:    $L \leftarrow \gamma$ 
18:   for all data centers  $d$  in the list except  $d'$  do
19:     set weight as  $r_d \div \gamma$ 
20:    $L \leftarrow L - r$ 
21:   end for
22:   Set the weight for  $d'$  as  $L \div \gamma$ 
   :
23: end if
24: Update HAProxy weights accordingly

```

#### 3.1. Green load balancing policy

[Algorithm 1](#) shows the pseudo code of the Green Load Balancing (GreenLB) policy used by the Global-LB component of the proposed system. Variable  $R$  is defined to maintain the overall rate of requests that can be handled by the power generated from renewable energy sources in all data centers. Lines 2–9 gather the required information using the monitoring modules and compute request rates at which a data center could accommodate requests only using renewable power within the time window (e.g., for the last 10 min).  $c$  and  $t$  keep track of the energy consumption and total number of requests served by the data center, respectively (Lines 3 and 4).  $a$  is a value for the total amount of renewable power currently available in the site (Line 5). By dividing value of  $c$  to  $t$ , we compute the amount of energy  $w$  consumed at the specific data center to serve each request (Line 6). Accordingly, the maximum rate of requests  $r_d$  specifying the rate at which this data

center can provide service only using renewable power is computed (Line 7). The value of  $R$  is updated based on the calculated rate  $r_d$  iteratively.

When the computation of  $r_d$  is done for each data center, the request rate at which global load balancer received requests in the recent time window,  $\gamma$ , is fetched from HAProxy (Line 10). If  $\gamma < R$ , it means that enough renewable energy is available to handle all incoming requests at the current rate (Line 11) and the weight for each data center is computed proportionally to the availability of renewable energy on each site (Line 13). Otherwise, if there is not enough renewable energy available to handle all requests at this rate (Line 15), each data center will receive requests based on the maximum rate of requests they can serve with renewable energy sources (Line 19) and the one with the lowest price of brown energy consumption per request accommodates the remaining parts of requests (Line 22). Finally, at Line 24, the HAProxy weight parameters are updated according to the computed values.

Please note that our policy never deallocates all web server machines in a specific data center and there is always at least one web server running on each data center. Therefore, if the availability of renewable energy sources at a certain data center is lower than the energy required for running even one web server machine, we set the weight in a way that the data center receives the minimum rate enough for one web server. Details related to this calculation are not shown in [Algorithm 1](#) for the sake of preserving clarity. In addition, we assumed that there are sufficient resources available in each data center to accommodate the whole workload. Therefore, we did not consider the case that the data center with the cheapest brown energy price cannot accommodate the remaining part of the requests. If this is the case, a cap/limit on the maximum rate of requests each data center can handle must be considered by the algorithm.

[Algorithm 1](#) runs a loop on all data centers to fetch required information to calculate values of  $R$  and  $r_d$ . Then according to the condition in Line 11 makes another loop to set weight values for each data center. Therefore, the asymptotic time complexity of algorithm is  $O(n)$ , where  $n$  is the number of data centers.

#### 3.2. Auto-scaling policy

[Algorithm 2](#) shows the details of the auto-scaling policy used in our system. The auto-scaling policy proposed in this paper is a reactive auto-scaling working based on threshold-based rules. By conducting profiling study, we set an upper threshold below which the web server is able to provide responses within an acceptable time frame. The ratio of arrival request rate at local load balancer to the threshold value gives us the total number of required web servers.

Note that, [Algorithm 2](#) is an algorithm without a loop or a recursion and incurs the time complexity of  $O(1)$ . The auto-scaler based on the profiling information sets the request rate at or below which a web server running on a specific machine type can efficiently generate responses for the web requests (Line 1). Then, in Lines 2–4, the total number of required web servers is computed according to the average request rate in the last time window (e.g., last 3 min). The number of required web servers is computed by dividing the average arrival rate of web requests at the local load balancer to the threshold value and taking the ceiling of the result (Line 4). If total number of required web server machines is higher than the number of the currently running web servers, it allocates more web server machines (Line 7); otherwise, if the number is lower, it deallocates excess machines (Line 9). Finally, in Line 13, the local load balancer is informed by the updated list of web servers. Please note that the local load balancing between the web server machines is done in an evenly weighted round robin fashion as the server farm is homogeneous.

**Algorithm 2.** Auto-scaling Policy.

```

1:  $t \leftarrow$  Set the threshold (i.e., appropriate rate for a
   web server) based on the profiling data
2:  $r \leftarrow$  Fetch the request rate in the recent time window
   from the local HAProxy load balancer
3:  $o \leftarrow$  Fetch the number of currently ON web server
   machines
4:  $m \leftarrow \lceil r \div t \rceil$  Compute the total number of web server
   machines required
5:  $n \leftarrow m - o$ 
6: if  $n > 0$  then
7:   Add  $n$  more web servers
8: else if  $n < 0$  then
9:   Remove  $n$  web servers
10: else
11:   No scaling is required
12: end if
13: Inform local load balancer with the new list of web
   servers

```

**4. Performance evaluation**

To evaluate our system, we developed a prototype and performed experimental studies in a real testbed (Grid'5000) using real traces of requests for English Wikipedia pages. In the following, we explain details of the testbed, workload, renewable power traces, and electricity prices. Then, we present results of experiments conducted using this setup. Our aim is to understand the cost and energy consumption performance of the proposed system in realistic settings.

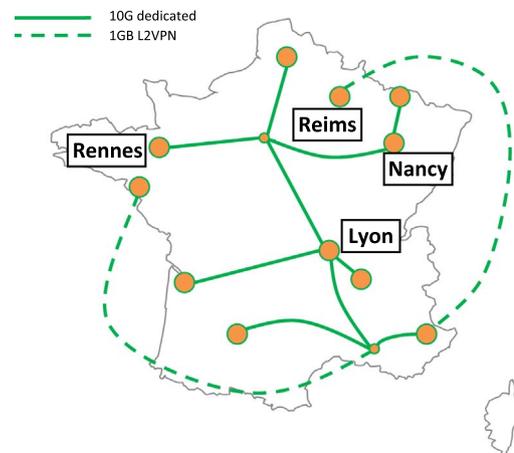
**4.1. Experimental setup****4.1.1. Testbed**

As testbed, we used Grid'5000 (Bolze et al., 2006), a French experimental grid platform.<sup>10</sup> Grid'5000 comprises sites geographically distributed across France. We consider a group of 3 sites equipped with power monitoring APIs (Clouet et al., 2015) in the following locations: *Lyon*, *Rennes* and *Reims* shown in Fig. 4. We set up the experimental testbed by preparing 3 deployable environments as follows:

- **Database (DB):** a mysql database server loaded with the English Wikipedia dataset as of Jan 3rd, 2008 containing roughly 2 million wiki pages and size of 3 GB.
- **Web Server (WS):** an Apache Web Server (version 2.4.10) with the installed Mediawiki application.
- **Local Load Balancer (Local-LB):** HAProxy (version 1.6) load balancer along with the auto-scaler Java program.

All deployable environments run on the *Debian Linux Wheezy* operating system. The characteristics of machines used in each site for hosting web servers are summarized in Table 1. To switch on/off a deployed web server, we used Grid'5000 APIs for accessing the “Wake-on-LAN” interfaces. By switching off, we set the physical machine to the “hard” power off mode (i.e., physical shut down) (Fig. 5).

To replay traces of requests by Wikipedia users, we used wikibench benchmark tool (van Baaren, 2009).<sup>11</sup> Wikibench is a web hosting benchmark allowing the stress-test of systems designed to host web applications. Using *wikijector* software module of wikibench, one can generate traffic by replaying traces of user requests actually made to Wikipedia. For our experimental study, we deploy wikijector on a

**Fig. 4.** Grid'5000 testbed.**Table 1**

The characteristics of machines hosting web servers in different sites.

Site	CPU	Number of cores	Memory
Lyon	AMD Opteron 250 2.4 GHz	2	2 GB
Reims	AMD Opteron 6164 HE 1.7 GHz	2	48 GB
Rennes	AMD Opteron 6164 HE 1.7 GHz	2	48 GB

machine in the *Nancy* site to mimic Wikipedia users sending requests to the Global-LB. Global-LB is deployed on a separate host in the *Nancy* site besides the controller Java program. The architecture of the prototype system deployed on the Grid'5000 testbed is shown in Fig. 5.

**4.1.2. Workload**

We used real traces of requests to the web servers from the Wikimedia Foundation as workload. Our workload contains 5% of all user requests issued to English Wikipedia resources during the period of 19th to 21st of September 2007.<sup>12</sup> Fig. 6 depicts the number of requests per second for the same period.

**4.1.3. Availability of renewable energy**

To capture the availability of solar energy in the location of each data center, we use data traces by SoDa Service<sup>13</sup> with 30 min granularity between 19th and 21st of Sept 2007. The Global Horizontal Irradiance (GHI) in the location of each data center is used to calculate the output for solar photovoltaics (PV) power. We assume that each data center uses power generated by the PV panels of total 4000 m<sup>2</sup> area with tilt angle of 45° degree and PV cell efficiency of 30% (Roughly the highest efficiency reported so far (Green et al., 2015)). We calculate the PV power module output on the tilt surface based on the model in *Solar Radiation on a Tilted Surface*.

We use meteorological data collected from Weather Underground<sup>14</sup> traces to model wind power for the same dates. We presume each data center uses a GE 1.5 MW wind turbine to generate wind power. To estimate the average wind power production, the model proposed by Fripp and Wiser (2008) is employed where the wind speed, the air temperature, and the air pressure measurements in the location of each site are fed into the model.

The summation of power generated from these two sources for

<sup>12</sup> <http://www.wikibench.eu/wiki/2007-09/><sup>13</sup> <http://www.soda-is.com>, The SoDa Service is a broker to a list of services and webservices related to Solar Radiation proposed by several providers in Europe and abroad. The SoDa Service is provided by two mirror sites: one hosted in MINES ParisTech, Sophia Antipolis France, and the other by Transvalor S.A., Mougins.<sup>14</sup> <http://www.wunderground.com/history/><sup>10</sup> The Grid 5000 project, <http://www.grid5000.org/>.<sup>11</sup> Wikibench, the realistic web hosting benchmark, <http://www.wikibench.eu/>.

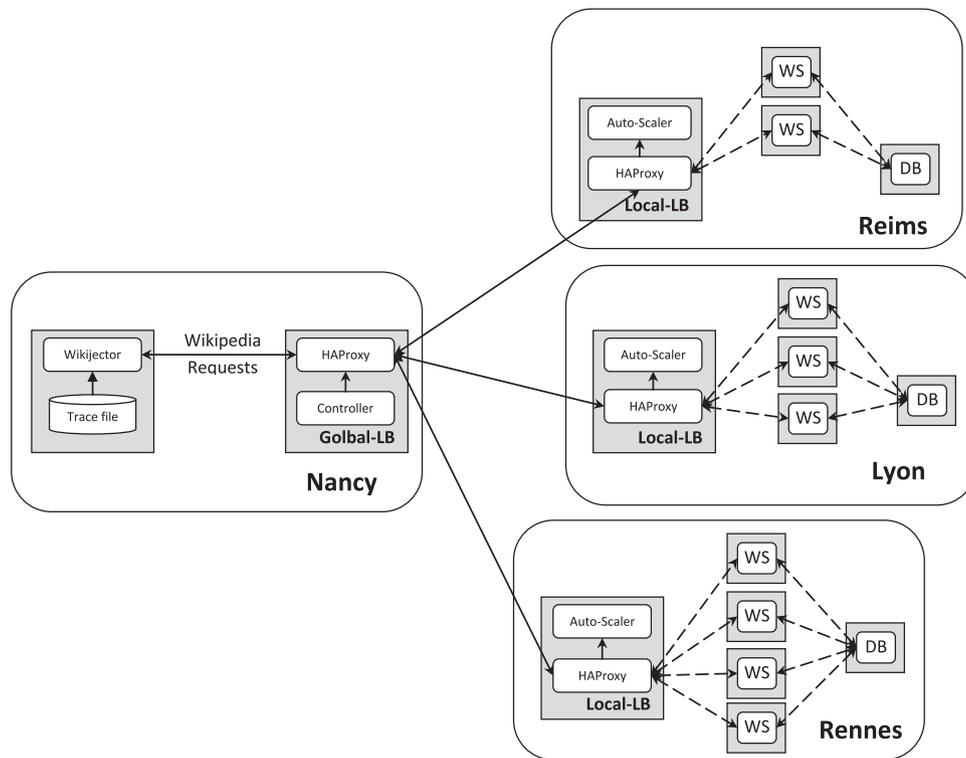


Fig. 5. The architecture of the prototype system.

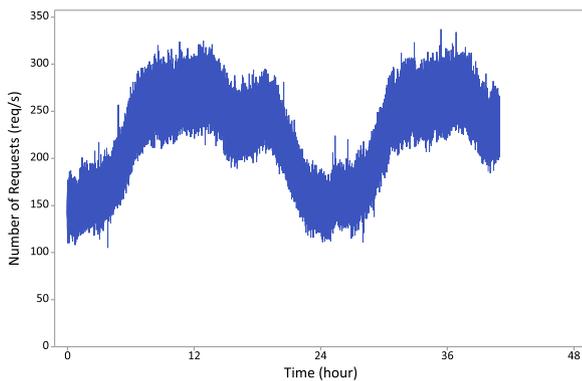


Fig. 6. The English Wikipedia workload for 19th and 20th of September 2007.

every 30 min is computed as available renewable power for every site. Figs. 7 (a), (b) and (c) show respectively the solar, wind, and normalized total power generated from renewable sources on each site for the 2-day period. In our investigation, we scale the normalized value of renewable power availability linearly in a way that the average renewable power generation is roughly enough to serve the average workload. Accordingly, every normalized power value is multiplied by 2128.5 computed in this way.

#### 4.1.4. Electricity prices

In France, the utility grid power is mostly generated by Électricité de France (EDF), French utility company, and is primarily produced from nuclear power sources. Moreover, EDF offers wide range of tariffs which are consistent all over France. In our experiment, we assume on-peak/off-peak scheme, the most common type of variable energy pricing in the market. The electricity prices charged by EDF for on-peak (between 8 am and 8 pm) and off-peak are set to €0.1636/kWh and €0.1150/kWh, respectively. In order to incorporate price variability and effect of brown

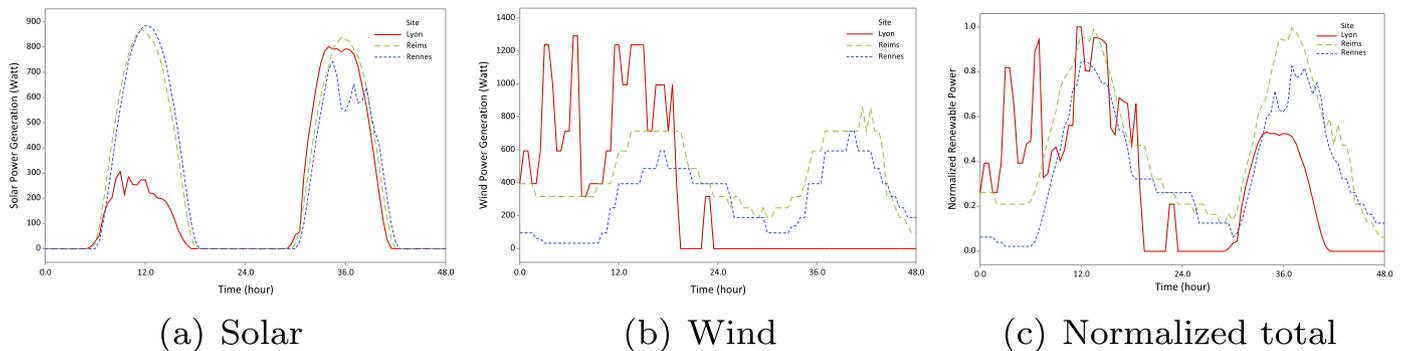


Fig. 7. Renewable Power Generation for two days.

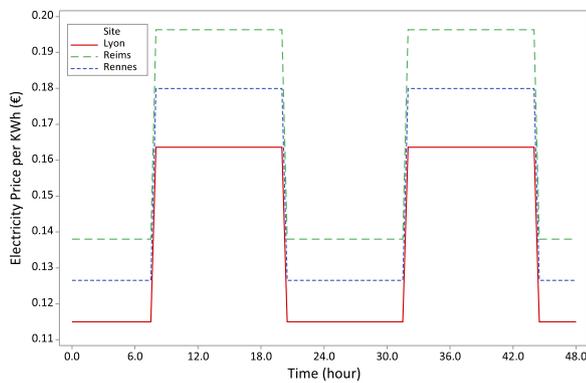


Fig. 8. Electricity prices for two days.

energy into our experiment, we increased the electricity price for each site by a factor based on the closeness of the site to non-renewable plants. Using geographical location of EDF non-renewable plants in France (EDF.), we increased the electricity prices for Reims, Rennes, and Lyon, respectively, by 20%, 10%, and 0%, as it is shown in the Fig. 8.

#### 4.2. Benchmark algorithms

To evaluate the performance of the proposed GreenLB policy, we consider two benchmark algorithms.

##### 4.2.1. Round robin (RR)

The RR policy sets equal weights for Global-LB which results in the even distribution of load among all sites. Since Wikipedia workload and the renewable power generated on every site exhibit a similar diurnal pattern, RR policy that evenly distributes requests among data centers is a competent benchmark policy to evaluate the performance of GreenLB.

##### 4.2.2. Capping

Le et al. (2010) proposed a policy for request distribution across data centers to minimize the overall energy cost. Similar to our method, their optimization method seeks to define the fraction of requests that should be forwarded to different data centers to minimize cost for an Internet service provider. Therefore, we decided to compare the performance of our proposed GreenLB policy with their method which we refer to as “Capping” policy from now on.

Since, their method is different from ours in some aspects and they consider brown energy caps for each data center (not considered in this paper), we modified their method in several ways to adapt it to our settings. First, we set *infinite* brown energy caps for data centers to let their policy solve the optimization problem in the absence of brown energy caps. We also considered zero cost for using green energy because in our settings data centers use free of charge on-site renewable power. Finally, the constraints related to fraction of requests are updated in a way that the minimum percentage of requests forwarded to each center is enough to utilize available renewable power on that site. If there is abundant renewable power (i.e., the renewable power generated by all data centers collectively is more than the power required to handle the entire workload), similar to our policy, we break down the load proportionally to the available renewable power for each data center as no feasible solution can be constructed by their optimization technique. According to the above modifications and assuming that the mixture of green and brown energy for the next time slot is known by the load balancer, their optimization method converts into an LP (Linear problem) which can be solved with an LP solver. Following their method, a solution for the optimization problem is periodically computed on a regular basis of once per hour. Their Auto-Regressive Integrated Moving Average (ARIMA) modeling is also used to predict the request rate for the next hour. Fig. 9 shows the

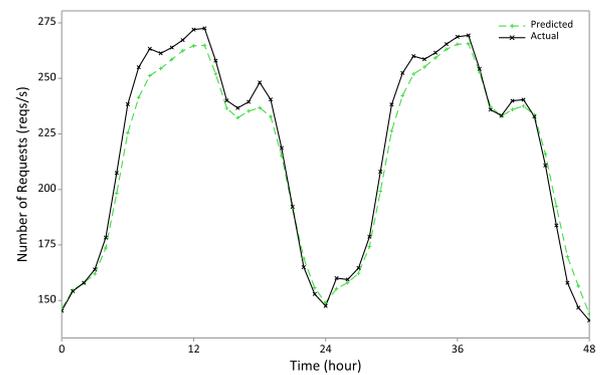


Fig. 9. Hourly actual and predicted request rate used by the Capping policy.

actual and ARIMA-predicted request rates for the two-day Wikipedia workload used in our experiments.

#### 4.3. Experimental results and analysis

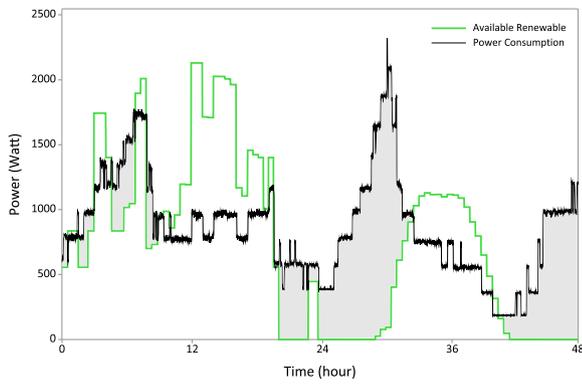
This section presents our experimental results. We run two-day experiments on the pre-configured testbed using the explained traces of the workload, renewable power and electricity prices for GreenLB, Capping, and RR policies. The auto-scaling algorithm on Local-LB are triggered every 2 min for all policies. The value is set based on a preliminary pilot study conducted before the main experiments. The load balancing algorithm on Global-LB is executed every 3 min for the proposed GreenLB policy and every 60 min for the Capping policy. We measure the power consumption via the servers' built-in power monitoring APIs with a granularity of 30 s. Please note that we only collect and report the power consumption for the web servers, since the power consumed by the load balancer and database server machines on each site are largely constant during the experiments.

Fig. 10 shows the power consumption and the brown energy usage for each site when the GreenLB algorithm is used to distribute workload among sites. The GreenLB algorithm successfully follows the available renewable energy on each site as it can be clearly seen in Reims and Rennes data centers. The pattern is different for Lyon because whenever there is insufficient renewable energy available, the algorithm redirects more requests to the Lyon data center which has the cheapest price of electricity per request.

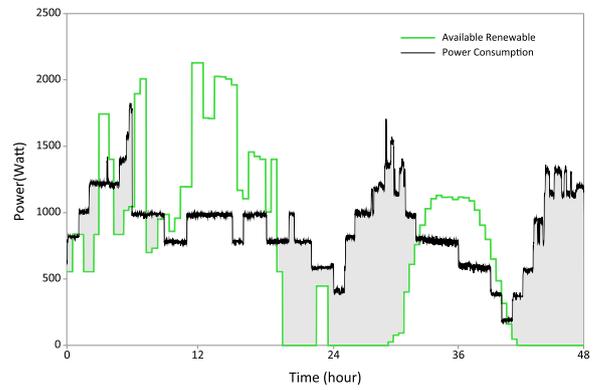
Similarly, Figs. 11 and 12 illustrate that the green and brown power consumptions when Capping and RR algorithms are employed respectively. As shown in Fig. 11, Capping policy demonstrates similar behavior to GreenLB, even though future knowledge regarding requests rate and availability of renewable energy for the next hour are available to this policy. This can be explained by the fact that the Wikipedia workload consists of many small sized requests that can be reactively routed among multiple data centers and in-advance proactive decision making is deemed unnecessary.

A comparison between Figs. 10 and 12 clearly illustrates that substantially more brown energy usage happens in the Reims and Rennes data centers when the RR algorithm is used while less brown energy usage occurs in the Lyon data center. Figs. 13–15 show the number of ON web servers per sites for GreenLB, Capping and RR policies, respectively.

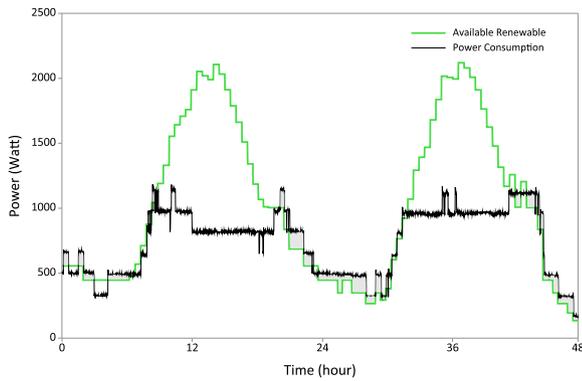
Table 2 summarizes the results shown in the figures. The aggregated total and brown power consumption for all sites shows that even though all algorithms cause similar power consumptions, GreenLB uses 17% less brown energy and saves cost by almost 22% in comparison to RR. The results demonstrate that GreenLB can significantly increase green energy utilization and decrease electricity cost by 8% and 22% compared to Capping and RR policies, respectively, even when sites have similar amount of renewable power production and very competitive price of electricity. The total amount of brown



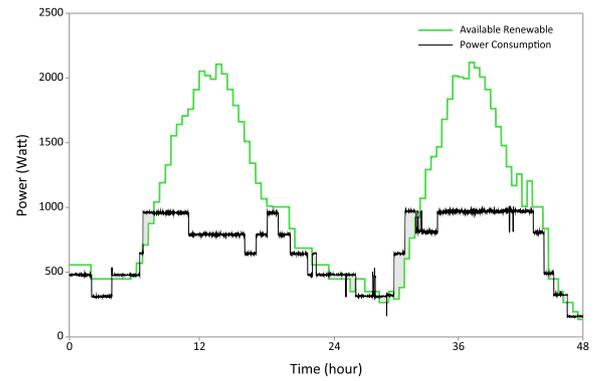
(a) Lyon



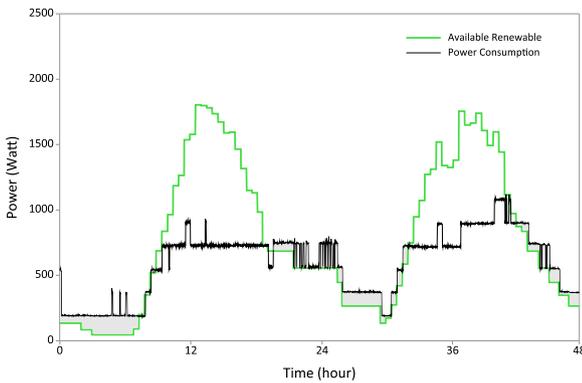
(a) Lyon



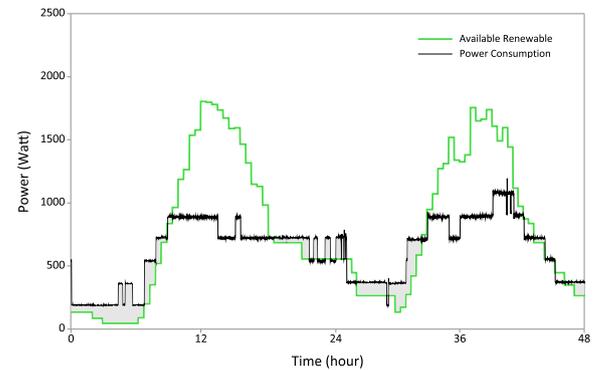
(b) Reims



(b) Reims



(c) Rennes



(c) Rennes

**Fig. 10.** The power consumption for different sites using Green Load Balancing (GreenLB) algorithm. The shaded area shows brown energy usage.

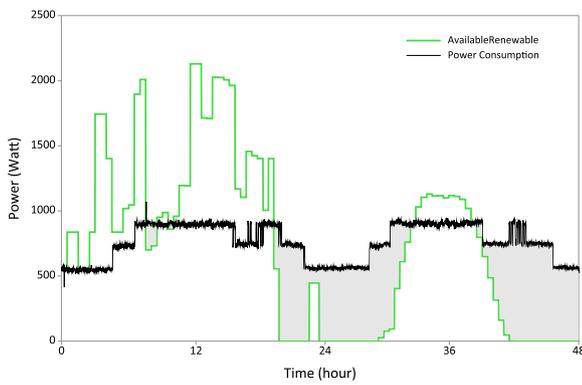
**Fig. 11.** The power consumption for different sites using Capping algorithm. The shaded area shows brown energy usage.

power consumed by GreenLB policy is 2.63 kWh, which is 7% and 17% less than Capping and RR policies, respectively. It is expected that the difference in renewable power utilization and cost saving increases substantially whenever there exists more gaps between the renewable power availability and electricity prices of different sites.

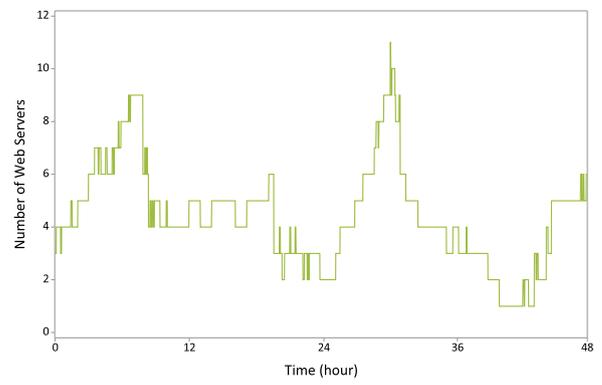
Capping and GreenLB share major similar characteristics and provide competitive solutions for the geographical load balancing problem tackled in this paper. However, the Capping policy works based on linear optimization which is considerably less efficient than GreenLB in terms of time and space complexity. Moreover, errors imposed by future load predication and hourly based decisions performed by Capping policy result in 7% more cost and 8% more

brown power usage compared to our proposed GreenLB algorithm.

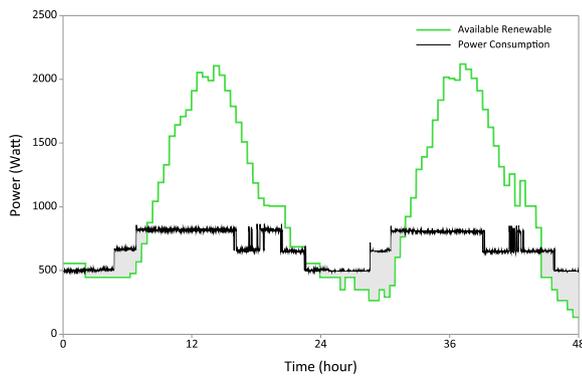
In order to study the impact of our proposed load balancing algorithm on the response time, we measured the real-time response times for all the Wikipedia requests submitted to the system. Fig. 16 shows the CDF of the response time observed by all load balancing algorithms. The graph demonstrates that there is no significant difference in the response time of the algorithms and majority of requests are responded within the acceptable range of hundreds of milliseconds, while RR shows marginally more stable response time. More than 90% of request are responded in about 350 ms and less for the GreenLB algorithm. There are few peaks of high response time upto few seconds for all policies happening by reasons such as Java garbage



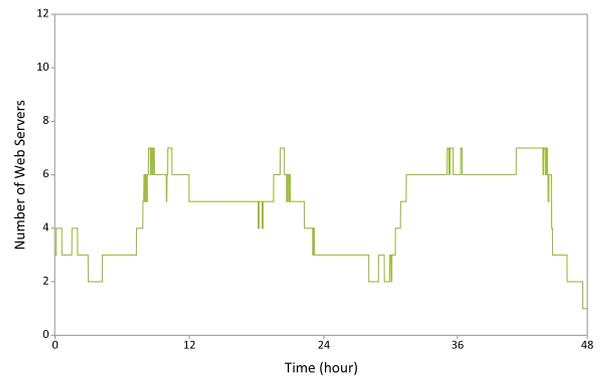
(a) Lyon



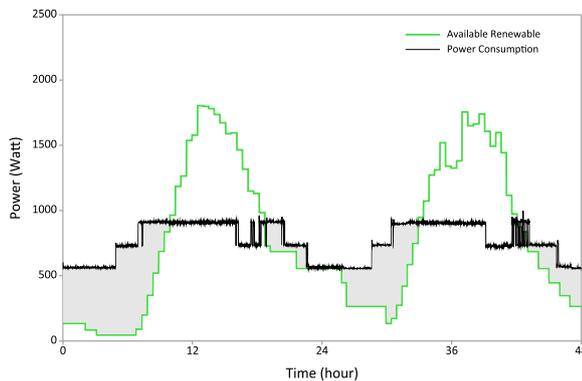
(a) Lyon



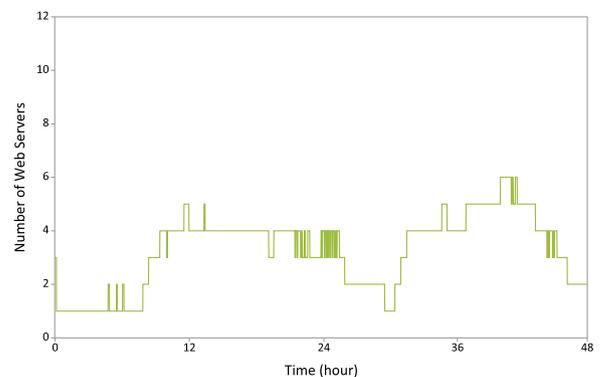
(b) Reims



(b) Reims



(c) Rennes



(c) Rennes

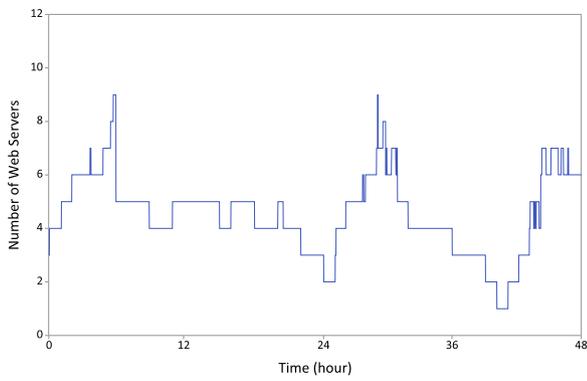
**Fig. 12.** The power consumption for different sites using Round Robin (RR) algorithm. The shaded area shows brown energy usage.

collection for Wikijector module and PHP garbage collections in Apache web servers.

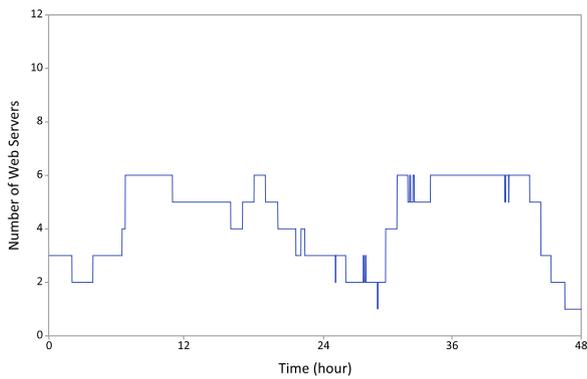
The comprehensive analysis of the experimental results also demonstrates that the impact of the delay associated with switching on and off web servers on the response time of requests is negligible. The main reason that switching on delay is not an issue is that we chose the threshold for scaling out web servers in the auto-scaler algorithm sufficiently below the maximum capacity of the each web server. Therefore, each web server machine has some spare capacity to handle portion of extra load before it becomes fully saturated and affects the response time. This way, our auto-scaler algorithm includes an acceptable level of over-provisioning to avoid Quality of Service

**Fig. 13.** The number of Web Servers for different sites using Green Load Balancing algorithm.

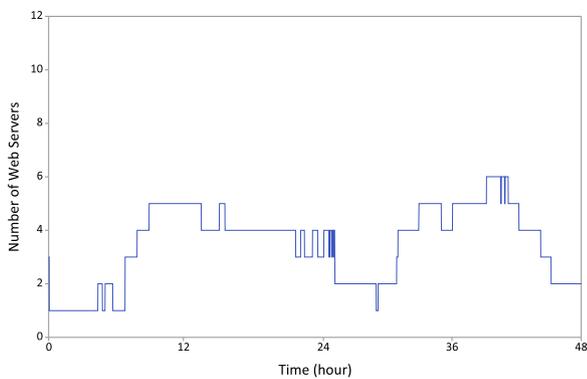
degradation due to switching on (booting) time of web server machines. That is, whenever the incoming request rate at the local load balancer increases, the auto-scaler (at least in case of Wikipedia workload) has enough time to scale out and to add new fully operational web servers. Switching off web server machines also does not affect the response time, as in the first step of the scale-in process, the auto scaler removes the target machine from the list of available machines in the local load balancer. Therefore, this machine will not receive additional requests. Then a signal for switching off is dispatched and the machine only switches off after all remaining requests in the web server queue are responded.



(a) Lyon



(b) Reims

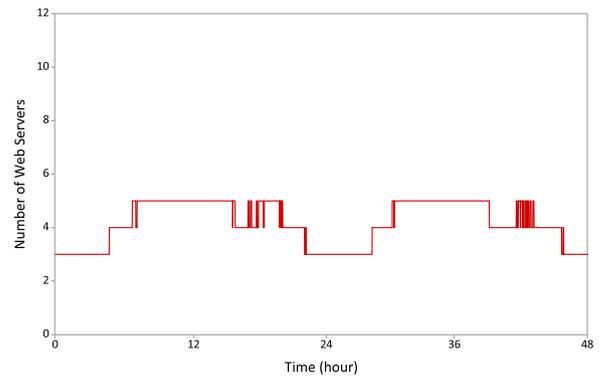


(c) Rennes

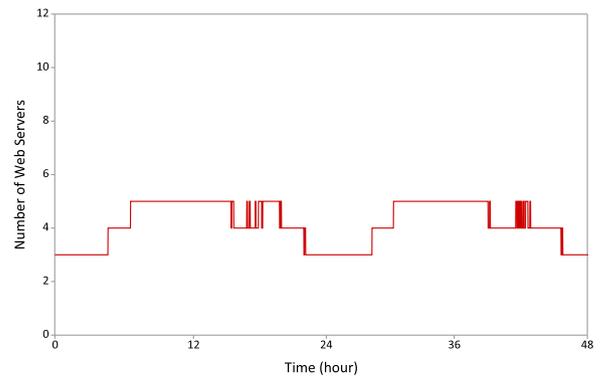
Fig. 14. The number of Web Servers for different sites using Capping algorithm.

5. Related work

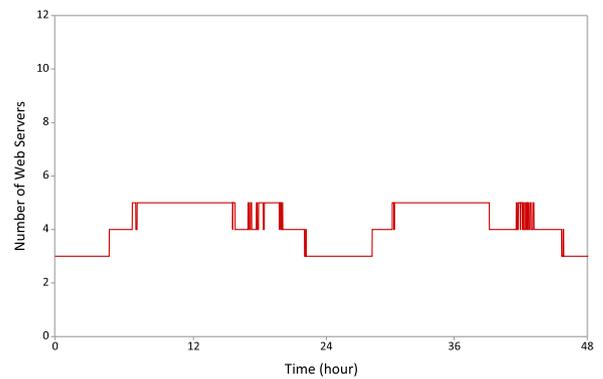
Over the last decade, power management techniques to minimize data centers' costs and environmental impacts have gained considerable attention by both academia and industry. Large data centers such as those used by big companies like Google and Amazon can host thousands of physical servers and require up to tens of megawatts to power them Kong and Liu (2014). As result, service providers are under huge pressure to reduce their energy consumption and its associated costs. This has pushed them towards using more sustainable and green data centers. In a recent study, Shuja et al. (2016a) have provided a survey of



(a) Lyon



(b) Reims



(c) Rennes

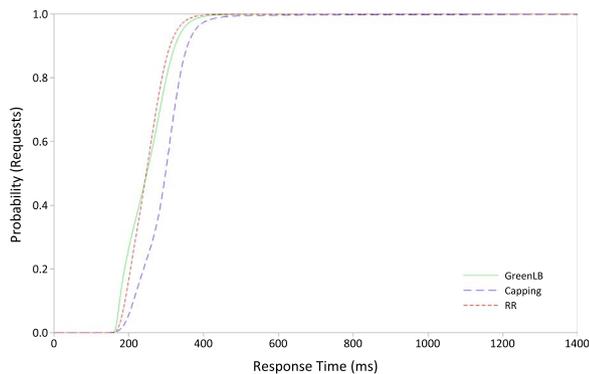
Fig. 15. The number of Web Servers for different sites using Round Robin (RR) algorithm.

enabling techniques and technologies for sustainable and green data centers.

Most of the early research studies on energy efficiency of data centers focus on making green data centers using optimization techniques within a single data center; techniques such as CPU dynamic voltage and frequency scaling (DVFS) (Wu et al., 2014), virtualization and VM consolidation (Beloglazov and Buyya, 2013; Srikantaiah et al., 2008), and workload scheduling (Ghamkhari and Mohsenian-Rad, 2013). An extensive survey and taxonomy of these can be found in Beloglazov et al. (2011). Similarly, Shuja et al. (2016b) conducted a survey of techniques and architectures for designing an energy-efficient data center.

**Table 2**  
Summary of Results.

Site	Metric	RR	Capping	GreenLB
<b>Lyon</b>	Power Consumption (kWh)	36.2	42.9	41.2
	Brown Consumption (kWh)	13.3	19.0	16.9
	Cost (€)	1.71	2.31	2.01
<b>Reims</b>	Power Consumption (kWh)	32.5	32.5	35.4
	Brown Consumption (kWh)	3.1	1.1	1.9
	Cost (€)	0.42	0.15	0.27
<b>Rennes</b>	Power Consumption (kWh)	36.4	29.7	28.3
	Brown Consumption (kWh)	9.3	2.9	2.6
	Cost(€)	1.23	0.39	0.35
<b>Total</b>	Power Consumption (kWh)	<b>105</b>	<b>105</b>	<b>105</b>
	Brown Consumption (kWh)	<b>25.7</b>	<b>23.0</b>	<b>21.4</b>
	Cost(€)	<b>3.36</b>	<b>2.85</b>	<b>2.63</b>



**Fig. 16.** CDF of average response time per second for English Wikipedia requests using different algorithms.

### 5.1. Leveraging green energy

There is a number of studies focused on reducing brown energy or power consumption, monetary costs, and environmental impact using renewable sources of energy for different types of applications, for example, batch processing (Goiri et al., 2013, 2011; Liu et al., 2012), and interactive processing (Stewart and Shen, 2009). Goiri et al. (2013) present Parasol, a green data center prototype. They define a scheduler for planning the workload execution and for selecting the energy sources: solar panels, batteries or grid. The scheduler makes decisions based on workload and energy predictions, battery level, DC characteristics and grid electricity prices. Liu et al. (2012) present an approach to model the energy flows in a data center in order to optimize its operation. They predict renewable energy and IT demand to schedule IT workload and allocates IT resources within a data center according to time varying power supply and cooling efficiency. Authors in Goiri et al. (2011) propose GreenSlot a parallel batch job scheduler for a data center powered by a photovoltaic solar array and the utility grid. Stewart and Shen (2009) also try to maximize green energy use in data center for interactive Internet services. All these studies focus on load or demand shifting to maximize green energy utilization within one data center, while we focus on multiple data center load redirection.

### 5.2. Geographical load balancing (GLB)

A large body of recent literature focuses on reducing energy costs targeting geographically distributed data centers. This group of work mainly devises techniques for workload distribution across geo-distributed data centers in order to achieve performance objectives such as minimizing cost, maximizing renewable energy utilization, and minimizing emission. Rahman et al. (2014) present a comprehensive survey on data center power management using geographic load balancing.

Among different approaches of geographical load balancing, “following the renewables” has gained considerable attentions. This approach requires that dynamic load balancing mechanism be aware of the availability of renewable energy at data centers (Shuja et al., 2016a). One of the early studies on GLB is done by Liu et al. (2011). Using GLB, they propose algorithms to maximize renewable energy utilization and show how dynamic electricity price can affect brown energy usage. They use trace-based numerical simulations to evaluate their algorithms. An extension to this work has been done by Lin et al. (2012), where they propose online algorithms to exploit the potential of geographical diversity of internet-scale services on renewable energy utilization. As part of their research, they show the optimal portfolio of solar and wind energy sources in GLB. Similarly, He et al. (2012) considered the sustainability of data centers by proposing socially-responsible load scheduling for data centers where they consider emission cost as the social cost. Chen et al. (2012) proposed a scheduling algorithm that considers the workload fluctuation, jobs' deadline, variable green energy supply, outside temperature, and data center cooling dynamics. In a recent comprehensive study, Paul et al. (2016) proposed a holistic framework for dynamic load distribution using online algorithms techniques. To minimize cost, they exploit the spatial variation in electricity price and renewable energy generation for a cloud service provider having a large number of data centers collocated with renewable energy sources. Their approach not only maximizes green energy utilization, but also minimizes the number of server switching. They have conducted extensive simulations with real data traces to evaluate their system. Berral et al. (2014) go one step further by proposing a framework that offers suitable locations of data centers to a provider seeking to create a network of sustainable data centers for a follow-the-renewables HPC cloud service.

All these studies consider data centers with on-site free of charge power generations from renewable energy sources. Similar to these, we consider GLB to reduce energy cost and to maximize renewable energy utilization. However, we mainly focus on practical considerations and we evaluate our proposed system in a real environment using realistic traces of workload, renewable energy and electricity prices.

### 5.3. Power capping

Another set of research efforts on dynamic load balancing assume that data center must pay for the power drawn from the off-site renewable energy sources. Le et al. (2010) considered load distribution across data center sites by including limitation on the energy usage from non-renewable sources. Gao et al. (2012) propose a framework, which is compared to the method by Le et al. (2010), for request-routing and traffic engineering considering changes in workload and carbon footprint. They attempt to balance the three-way trade-off between access latency, carbon footprint, and electricity costs. A recent topic of interest that explores capping the brown energy consumption has been studied by several studies such as Le et al. (2010), Abbasi et al. (2014). Policies used in this category propose techniques to abide by carbon caps on brown energy consumption. Le et al. (2010) propose a software framework to distribute requests among different data centers while supporting the capping brown energy consumption. The goal is to abide by the carbon caps for each data center without excessively increasing costs or degrading performance. Similarly, Abbasi et al. (2014), propose online algorithms to tackle the same problem. They remove the need for long-term future prediction via exploiting their proposed online algorithms. All these research efforts assume that data centers have to pay for renewable energy from off-site utility power. Moreover, electricity prices for power generated from off-site renewable sources might be higher than brown power. However, we assume data centers are equipped with on-site renewable power facilities and power generated in this way is free of charge.

#### 5.4. Online application load balancing

Efficient Use of geographically dispersed resources for online interactive application such as web applications has been explored by several authors. Goudarzi and Pedram (2013) investigate the load balancing problem for online service applications using a cloud system comprised of geographically dispersed data centers. Offline and online algorithms are proposed to determine the application placement and migration based on renewable energy generation capacities at different data centers in the cloud system. Similar to our work, Kanizo et al. (2013) use Wikimedia access logs to test their threshold-based load sharing techniques. However, they did not consider renewable energy power generation in their work. Using the same workload, Zhang et al. (2011) proposed GreenWare, a framework to maximize renewable energy utilization of Geo-distributed data centers. Contrary to our work, they use simulation for evaluation purposes and assume that renewable energy is more expensive than brown energy.

#### 6. Conclusions and future work

We proposed a cost and energy efficient load balancing algorithm to distribute web applications requests among multiple data centers geographically distributed in a region. A prototype was implemented and experimental studies in a real testbed are performed using it. A group of 3 data centers from Grid5000 testbed scattered in France and equipped with power monitoring APIs were selected. Meteorological data in the location of each data center used to model solar and wind power generation. Real traces of web requests for English Wikipedia pages were replayed to generate the workload. Our proposed green load balancing algorithm was employed to distribute load among data centers based on the availability of renewable energy and prices of grid brown energy for each site. Finally, total power consumption, brown power usage, and cost of electricity were measured and compared to two benchmark algorithms. Results showed that, even in case similar amount of renewable power production in sites and very competitive price of electricity, our proposed policy is able to reduce the cost by 22% and 8% and brown energy by 17% and 7% in comparison to round robin and policy proposed in Le et al. (2010), respectively. We also demonstrated that the average response time of requests is not affected by our proposed load balancing algorithm.

Studies like the one presented in this article can be used to measure the required amount of energy to handle certain workload. Accordingly, this provides insights regarding the renewable energy facilities required by each data center to handle the load. Future research needs to be done to incorporate our proposed method with support for investment decisions on establishing renewable power generators in a real-world system.

In our proposed system, we only considered stateless web applications in which requests can be responded individually regardless of which user issued them. In future, we focus on the design and implementation of “Sticky load balancing” policies in which after a session between a client and an application server is established, all subsequent requests from this session are redirected to the same server. Moreover, in this paper, we limit our geographical load balancing to data centers dispersed throughout a small region (e.g., France or Europe) where routing requests among different data centers will not significantly affect response time as network delays remains in an acceptable range. We are interested to extend our system for worldwide geographical load balancing in which network proximity of the user is considered in the policy.

The rest of our future work will focus on renewable energy-aware geographical load balancing for other types of workloads/applications such as bag of tasks, scientific workflows, where platforms and tools such as Aneka (Buyya and Barreto, 2015) or workflow engine (Pandey et al., 2011) can be employed. Demand response and capping the brown power consumption to reduce carbon footprint and promote

carbon neutrality are among recent topics of interests which can also be considered as future research targets.

#### Acknowledgments

This work was partially supported by Australian Research Council (ARC) Future Fellowship (grant no. FL170100099) and Discovery Project grants (grant no. DP130101378). Experiments presented in this paper were carried out using the Grid'5000 testbed, supported by a scientific interest group hosted by INRIA and including CNRS, RENATER and several Universities as well as other organizations (see <https://www.grid5000.fr>). Authors would like to thank Adam Wierman from Caltech for his inspiring thoughts on this work and David Margery and Laurent Lefevre from Inria for their technical support of the experiments. They would also like to thank Rodrigo N. Calheiros, the editors of this journal, and three anonymous reviewers for their many helpful comments and suggestions on an earlier version of this article.

#### References

- Abbasi, Z., Pore, M., Gupta, S.K.S., 2014. Online server and workload management for joint optimization of electricity cost and carbon footprint across data centers. In: Proceedings of the 28th IEEE International Parallel and Distributed Processing Symposium, IPDPS '14, pp. 317–326. doi:10.1109/IPDPS.2014.42.
- Adnan, M., Sugihara, R., Gupta, R., 2012. Energy efficient geographical load balancing via dynamic deferral of workload. In: Proceedings of the 5th IEEE International Conference on Cloud Computing, CLOUD '12, pp. 188–195. <http://dx.doi.org/10.1109/CLOUD.2012.45>.
- Beloglazov, A., Buyya, R., 2013. Managing overloaded hosts for dynamic consolidation of virtual machines in cloud data centers under quality of service constraints. IEEE Trans. Parallel Distrib. Syst. 24 (7), 1366–1379. <http://dx.doi.org/10.1109/TPDS.2012.240>.
- Beloglazov, A., Buyya, R., Lee, Y.C., Zomaya, A., 2011. A taxonomy and survey of energy-efficient data centers and cloud computing systems. Adv. Comput. 82 (2), 47–111.
- Berral J.L., Goiri, Nguyen, T.D., Gavald, R., Torres, J., Bianchini, R., 2014. Building green cloud services at low cost. In: Proceedings of the 34th IEEE International Conference on Distributed Computing Systems, ICDCS '14, pp. 449–460. <http://dx.doi.org/10.1109/ICDCS.2014.53>.
- Bolze, R., Cappello, F., Caron, E., Daydé, M., Desprez, F., Jeannot, E., Jégou, Y., Lanteri, S., Leduc, J., Melab, N., et al., 2006. Grid'5000: a large scale and highly reconfigurable experimental grid testbed. Int. J. High Perform. Comput. Appl. 20 (4), 481–494. <http://dx.doi.org/10.1177/1094342006070078>.
- Buyya, R., Barreto, D., 2015. Multi-cloud resource provisioning with aneka: A unified and integrated utilisation of microsoft azure and amazon EC2 instances. In: Proceedings of the 2015 International Conference on Computing and Network Communications, CoCoNet '15, IEEE, pp. 216–229. <http://dx.doi.org/10.1109/CoCoNet.2015.7411190>.
- Chen, C., He, B., Tang, X., 2012. Green-aware workload scheduling in geographically distributed data centers. In: Proceedings of the 4th IEEE International Conference on Cloud Computing Technology and Science, CloudCom '12, pp. 82–89. <http://dx.doi.org/10.1109/CloudCom.2012.6427545>.
- Clouet, F., Delamare, S., Gelas, J.-P., Lefevre, L., Nussbaum, L., Parisot, C., Pouilloux, L., Rossignaux, F., 2015. A unified monitoring framework for energy consumption and network traffic. In: International Conference on Testbeds and Research Infrastructures for the Development of Networks & Communities, TridentCom '15, pp. 10.
- EDF, Fossil-fired energy URL:[http://energie.edf.com/fichiers/fckeditor/Commun/En\\_Direct\\_Centrales/collection\\_nos\\_energies/edf\\_thermique\\_bd\\_va.pdf](http://energie.edf.com/fichiers/fckeditor/Commun/En_Direct_Centrales/collection_nos_energies/edf_thermique_bd_va.pdf).
- Fernandez H., Pierre G., Kielmann, T., 2014. Autoscaling web applications in heterogeneous cloud infrastructures. In: Proceedings of IEEE International Conference on Cloud Engineering, IC2E '14, pp. 195–204. <http://dx.doi.org/10.1109/IC2E.2014.25>.
- Fripp, M., Wiser, R.H., 2008. Effects of temporal wind patterns on the value of wind-generated electricity in California and the northwest. IEEE Trans. Power Syst. 23 (2), 477–485. <http://dx.doi.org/10.1109/TPWRS.2008.919427>.
- Gao, P.X., Curtis, A.R., Wong, B., Keshav, S., 2012. It's not easy being green. SIGCOMM Comput. Commun. Rev. 42 (4), 211–222. <http://dx.doi.org/10.1145/2377677.2377719>.
- Ghamkhar, M., Mohsenian-Rad, H., 2013. Energy and performance management of green data centers: a profit maximization approach. IEEE Trans. Smart Grid 4 (2), 1017–1025. <http://dx.doi.org/10.1109/TSG.2013.2237929>.
- Goiri, I., Le, K., Haque, M.E., Beauchea, R., Nguyen, T.D., Guitart, J., Torres, J., Bianchini, R., 2011. Greenslot: Scheduling energy consumption in green datacenters. In: Proceedings of the 2011 International Conference for High Performance Computing, Networking, Storage and Analysis, SC '11, ACM, pp. 20:1–20:11. <http://dx.doi.org/10.1145/2063384.2063411>.
- Goiri, I., Le, K., Nguyen, T.D., Guitart, J., Torres, J., Bianchini, R., 2012. Greenhadoop: Leveraging green energy in data-processing frameworks. In: Proceedings of the 7th

- ACM European Conference on Computer Systems, EuroSys '12, ACM, pp. 57–70. <http://dx.doi.org/10.1145/2168836.2168843>.
- Goiri, I., Katsak, W., Le, K., Nguyen, T.D., Bianchini, R., 2013. Parasol and greenswitch: managing datacenters powered by renewable energy. SIGARCH Comput. Archit. News 41 (1), 51–64. <http://dx.doi.org/10.1145/2490301.2451123>.
- Goudarzi, H., Pedram, M., 2013. Geographical load balancing for online service applications in distributed datacenters. In: Proceedings of the 6th IEEE International Conference on Cloud Computing, Cloud '13, pp. 351–358. <http://dx.doi.org/10.1109/CLOUD.2013.77>.
- Green, M.A., Emery, K., Hishikawa, Y., Warta, W., Dunlop, E.D., 2015. Solar cell efficiency tables (version 45). Prog. Photovolt.: Res. Appl. 23 (1), 1–9. <http://dx.doi.org/10.1002/ppp.2573>.
- Grozev, N., Buyya, R., 2014. Multi-cloud provisioning and load distribution for three-tier applications. ACM Trans. Auton. Adapt. Syst. (TAAS) 9 (3), 13:1–13:21. <http://dx.doi.org/10.1145/2662112>.
- He, J., Deng, X., Wu, D., Wen, Y., Wu, D., 2012. Socially-responsible load scheduling algorithms for sustainable data centers over smart grid. In: Proceedings of the Third IEEE International Conference on Smart Grid Communications, SmartGridComm '12, pp. 406–411. <http://dx.doi.org/10.1109/SmartGridComm.2012.6486018>.
- Jiang, J., Lu, J., Zhang, G., Long, G., 2013. Optimal cloud resource auto-scaling for web applications. In: Proceedings of the 13th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing, CCGrid '13, pp. 58–65. <http://dx.doi.org/10.1109/CCGrid.2013.73>.
- Kanizo, Y., Raz, D., Zlotnik, A., 2013. Efficient use of geographically spread cloud resources. In: Proceedings of the 13th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing, CCGrid '13, pp. 450–457. <http://dx.doi.org/10.1109/CCGrid.2013.18>.
- Kong, F., Liu, X., 2014. A survey on green-energy-aware power management for datacenters. ACM Comput. Surv. 47 (2), 30:1–30:38. <http://dx.doi.org/10.1145/2642708>.
- Le, K., Bilgir, O., Bianchini, R., Martonosi, M., Nguyen, T.D., 2010. Managing the cost, energy consumption, and carbon footprint of internet services. In: Proceedings of the ACM International Conference on Measurement and Modeling of Computer Systems, SIGMETRICS '10, ACM, pp. 357–358. <http://dx.doi.org/10.1145/1811039.1811085>.
- Le, K., Bianchini, R., Nguyen, T.D., Bilgir, O., Martonosi, M., 2010. Capping the brown energy consumption of internet services at low cost. In: Proceedings of the International Green Computing Conference, IGCC '10, pp. 3–14. <http://dx.doi.org/10.1109/GREENCOMP.2010.5598305>.
- Lin M., Wierman, A., Andrew, L.L.H., Thereska, E., 2011. Dynamic right-sizing for power-proportional data centers. In: Proceedings of the INFOCOM 2011, pp. 1098–1106. <http://dx.doi.org/10.1109/INFOCOM.2011.5934885>.
- Lin, M., Liu, Z., Wierman, A., Andrew, L., Online algorithms for geographical load balancing. In: Proceedings of the International Green Computing Conference, IGCC '12, pp. 1–10. <http://dx.doi.org/10.1109/IGCC.2012.6322266>.
- Liu, Z., Lin, M., Wierman, A., Low, S.H., Andrew, L.L., 2011. Greening geographical load balancing. In: Proceedings of the ACM SIGMETRICS Joint International Conference on Measurement and Modeling of Computer Systems, SIGMETRICS '11, ACM, pp. 233–244. <http://dx.doi.org/10.1145/1993744.1993767>.
- Liu, Z., Lin, M., Wierman, A., Low, S.H., Andrew, L.L., 2011. Geographical load balancing with renewables. SIGMETRICS Perform. Eval. Rev. 39 (3), 62–66. <http://dx.doi.org/10.1145/2160803.2160862>.
- Liu, Z., Chen, Y., Bash, C., Wierman, A., Gmach, D., Wang, Z., Marwah, M., Hyser, C., 2012. Renewable and cooling aware workload management for sustainable data centers. In: Proceedings of the 12th ACM SIGMETRICS/PERFORMANCE Joint International Conference on Measurement and Modeling of Computer Systems, SIGMETRICS '12, ACM, pp. 175–186. <http://dx.doi.org/10.1145/2254756.2254779>.
- NRDC, 2014. Anthesis, Scaling up energy efficiency across the data center industry: Evaluating key drivers and barriers, Tech. rep., Natural Resources Defense Council.
- Pandey, S., Karunamoorthy, D., Buyya, R., 2011. Workflow engine for clouds. In: Buyya, R., Broberg, J., Goscinski, A. (Eds.), Cloud Computing: Principles and Paradigms. John Wiley & Sons, Inc., 321–344. <http://dx.doi.org/10.1002/9780470940105.ch12>.
- Paul, D., Zhong, W.-D., Bose, S.K., 2016. Energy efficiency aware load distribution and electricity cost volatility control for cloud service providers. J. Netw. Comput. Appl. 59, 185–197. <http://dx.doi.org/10.1016/j.jnca.2015.08.012>.
- Rahman, A., Liu, X., Kong, F., 2014. A survey on geographic load balancing based data center power management in the smart grid environment. IEEE Commun. Surv. Tutor. 16 (1), 214–233. <http://dx.doi.org/10.1109/SURV.2013.070813.00183>.
- Shuja, J., Gani, A., Shamshirband, S., Ahmad, R.W., Bilal, K., 2016a. Sustainable cloud data centers: a survey of enabling techniques and technologies. Renew. Sustain. Energy Rev. 62, 195–214. <http://dx.doi.org/10.1016/j.rser.2016.04.034>.
- Shuja, J., Bilal, K., Madani, S.A., Othman, M., Ranjan, R., Balaji, P., Khan, S.U., 2016b. Survey of techniques and architectures for designing energy-efficient data centers. IEEE Syst. J. 10 (2), 507–519. <http://dx.doi.org/10.1109/JSYST.2014.2315823>.
- Solar Radiation on a Tilted Surface, <http://pvducation.org/pvcdrom/properties-of-sunlight/solar-radiation-on-tilted-surface>.
- S. Srikantaiah, A. Kansal, F. Zhao, 2008. Energy aware consolidation for cloud computing. In: Proceedings of the USENIX 2008 Conference on Power Aware Computing and Systems, Vol. 10 of HotPower '08, San Diego, California.
- Stewart, C., Shen, K., 2009. Some joules are more precious than others: Managing renewable energy in the datacenter. In: Proceedings of the Workshop on Power Aware Computing and Systems (HotPower) 22nd ACM Symposium on Operating Systems Principles, SOSP '09, pp. 15–19.
- Urgaonkar, B., Shenoy, P., Chandra, A., Goyal, P., Wood, T., 2008. Agile dynamic provisioning of multi-tier internet applications. ACM Trans. Auton. Adapt. Syst. (TAAS) 3 (1), 1:1–1:39. <http://dx.doi.org/10.1145/1342171.1342172>.
- van Baaren, E.-J. 2009. Wikibench: A distributed, wikipedia based web application benchmark, Master's thesis, VU University Amsterdam.
- Wu, C.-M., Chang, R.-S., Chan, H.-Y., 2014. A green energy-efficient scheduling algorithm using the dvfs technique for cloud datacenters. Future Gener. Comput. Syst. 37, 141–147. <http://dx.doi.org/10.1016/j.future.2013.06.009>.
- Zhang Y., Wang, Y., Wang, X., 2011. Greenware: Greening cloud-scale data centers to maximize the use of renewable energy. In: F. Kon, A.-M. Kermarrec (Eds.), Middleware 2011, Vol. 7049 of Lecture Notes in Computer Science, Springer Berlin Heidelberg, pp. 143–164. [http://dx.doi.org/10.1007/978-3-642-25821-3\\_8](http://dx.doi.org/10.1007/978-3-642-25821-3_8).