



# Eco4Cloud on Cisco UCS: Unified Efficiency

*The abstraction of resources provided by Cisco UCS is the perfect starting point for the dynamic/real-time workload consolidation performed by Eco4Cloud. This document highlights the outcomes of the Interoperability Verification Test (IVT) jointly performed at Cisco's labs in San Jose, CA*

Agostino Forestiero<sup>(1,2)</sup>, Raffaele Giordanelli<sup>(1)</sup>, Carlo Mastroianni<sup>(1,2)</sup>, Giuseppe Papuzzo<sup>(1,2)</sup>

(1) Eco4Cloud – [www.eco4cloud.com](http://www.eco4cloud.com)

(2) ICAR-CNR, Institute for High Performance Computing and Networking of the Italian National Research Council, Italy – [www.icar.cnr.it](http://www.icar.cnr.it)

### 1. Executive summary

This whitepaper reports the outcomes of a Proof of Concept test carried out at one of Cisco Interoperability Verification Testing (IVT) labs in San Jose, between March 28<sup>th</sup> and April 2<sup>th</sup>, 2014. The resulting CapEx and OpEx reductions are remarkable, as well as the benefits from the integration of [Eco4Cloud](#) with Cisco's Unified Computing System ([UCS](#)).

[Cisco Unified Computing System](#) (UCS) is a programmable infrastructure component for the data center, built around an architecture natively conceived and optimized for virtualization deployments. UCS allows the system to be integrated into higher-level, data-center-wide management systems as a single, logical entity. Server, network, and I/O resources are provisioned and configured on demand by ways of service profiles.

Once Cisco UCS is wired, system resources become part of a flexible pool that can be quickly used for any workload on demand. The system abstracts the configuration and connectivity of server and I/O resources, allowing administration to be automated.

Cisco UCS offers a single unified system, transcending the traditional boundaries of blade chassis and racks. Cisco UCS brings together server, network, and storage access resources to create a physically distributed but centrally managed system. By abstracting the personalization, configuration, and connectivity of server and I/O resources, these attributes can be programmed automatically.

The abstraction of resources provided by Cisco UCS is the perfect starting point for the workload consolidation provided by Eco4Cloud. Indeed, UCS creates a homogeneous environment tailored to applications migration and dynamic server hibernation and resume.

Eco4Cloud contributes to boost the power efficiency of Cisco UCS, as the following section explains.

[Eco4Cloud](#) is in fact an innovative bio-inspired probabilistic algorithm which consolidates the maximum number of virtual machines on the minimum number of physical servers in a data center, enabling the switch off/hibernation of those freed-up, making them dynamically available as additional capacity for incremental workloads. The Eco4Cloud algorithm was proposed by a research project carried out by the [Institute for High Performance Computing and Networking of the Italian National Research Council](#) (ICAR-CNR) and by the [Politecnico di Torino](#) [1]. Research activities led to an industrial invention filed as an international [PCT patent](#) and owned by CNR and University of Calabria, titled "System for Energy Saving in Company Data Centers" [2].

## Eco4Cloud on Cisco UCS: Unified Efficiency

*“In a virtualized data center, only 20-30% of server capacity is utilized on average.*

*Still, an idle server consumes 65-70% of the power consumed when it is fully utilized.”*

Figure 1 shows an example of workload consolidation with Eco4Cloud on a data center with 100 servers. The figure shows the CPU utilization of the servers, which, before applying Eco4Cloud, is between 20% and 40%. At time-zero, the consolidation algorithm starts: following the VM migrations, some servers increase their respective utilization rate while others are unloaded and eventually switched off. At the end of the experiment, 35 servers have taken the entire load and the remaining 65 servers have been hibernated, leading to energy savings in excess of 50%.

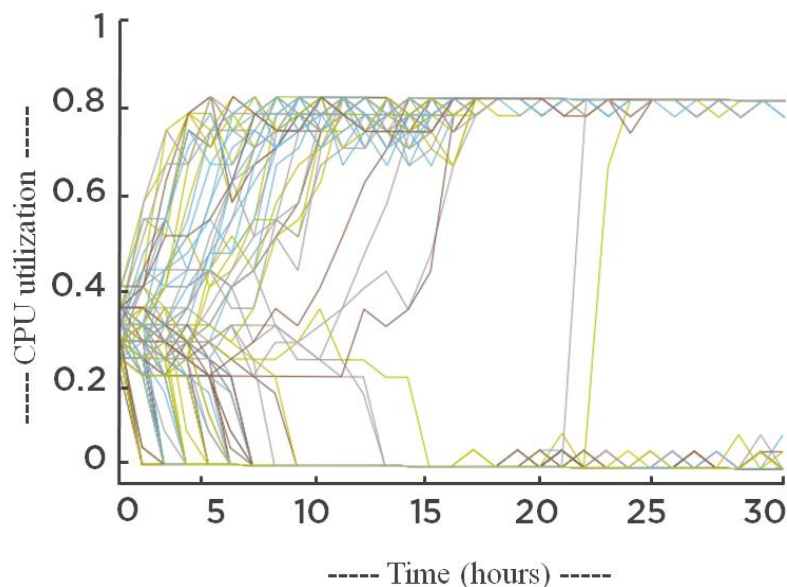


Figure 1. Example of workload consolidation on 100 servers. The workload is consolidated on 35 servers and the other servers are hibernated, leading to great energy savings.

The advantages achievable for a mid/large-scale data center are the following:

### 1) COMPUTATIONAL EFFICIENCY

By using servers at the best of their capacity, efficient consolidation algorithms give the opportunity to execute more tasks (up to 3X) with the same number of physical servers without any impact on the related quality of service, while actually improving it.

### 2) REDUCTION OF DATA CENTER ENERGY BILL

In a virtualized data center, only 20-30% of server capacity is utilized on average. Still, an idle server consumes 65-70% of the power consumed when it is fully utilized. By

*“Eco4Cloud on UCS in mid/large-scale data centers yields tangible advantages in the following areas:*

*Computational efficiency*

*Energy reduction*

*Capacity planning*

*SLA management*

*Scalability”*

consolidating the maximum number of VMs on the minimum number of physical servers, it is possible to reduce the overall energy bill up to 60%. Furthermore, by having less energy consumed by physical servers, power cooling gets reduced as well, contributing to indirect additional energy reductions (OpEx).

### **3) CAPACITY PLANNING**

Thanks to the optimal occupancy of physical resources and to the adaptive optimization of inherently variable workloads, CIOs can improve the budget accuracy and capacity planning for their data centers, hence optimizing their CapEx and OpEx.

### **4) SERVICE LEVEL AGREEMENTS COMPLIANCE (reliability, availability, performance)**

Adaptive and self-organizing solutions allow to proactively/predictively prevent quality of service degradation, for example by instructing the move of VMs from servers that are about to become overloaded. VM migrations are gradual, asynchronous and easy to monitor, whereas classical consolidation algorithms may require the concurrent move of hundreds of VMs.

### **5) SCALABILITY**

Adaptive/self-organized algorithms are much more scalable than classical algorithms, thanks to their capacity of taking decisions on the basis of local and readily/easily available information and to the use of statistical data. Therefore, all the aforementioned benefits improve proportionally with the data center size.

*“Eco4Cloud’s setup only took half an hour to complete and proved quite simple overall”*

### 2. Test Description

Specifically, the Eco4Cloud – Cisco UCS test has been performed on 5 UCS Blade servers running VMWare’s vSphere ESXi 5.5 virtualization platform. The servers are grouped in 1 cluster of 5 servers. The servers have been equipped with the following hardware configurations:

| Cluster  | UCS Blade Servers  | CPU core | RAM (GB) |
|----------|--------------------|----------|----------|
| CiscoE4C | API2-c1-s1-101-130 | 16       | 24       |
| CiscoE4C | API2-c1-s2-101-131 | 16       | 24       |
| CiscoE4C | API2-c1-s3-101-132 | 16       | 24       |
| CiscoE4C | API2-c2-s1-101-134 | 16       | 24       |
| CiscoE4C | API2-c2-s2-101-135 | 16       | 24       |

The Login VSI tool was used to generate the needed number of virtual machines for the synthetic workload exploited in the test.

Eco4Cloud’s setup only took half an hour **to complete** and proved quite simple overall. In fact Eco4Cloud is distributed as a virtual machine, and the VM format is OVA 1.0, hence compatibility is guaranteed with VMware vSphere 4.x and 5.x. Eco4Cloud can be deployed onto vCenter and configured as any other virtual machine running in the datacenter. As soon as the Eco4Cloud virtual machine became active, the Eco4Cloud trial license code and hypervisor connection parameters were entered, and within minutes the virtual appliance started collecting data. The test began with the installation of the Eco4Cloud monitoring tool which allowed to collect the traces of all the physical servers and the VMs.

In a test bed infrastructure, the high level steps which are commonly adopted to evaluate the effectiveness of optimization actions on a data center can be summarized as follows:

1. Identification of typical workloads and resource allocation trends of real data centers
2. Set-up of tools which allow to define repeatable, synthetic loads, resembling the identified workloads on a different scale (the scale depends on the physical resources available in the test bed)
3. Execution and repetition of such workloads under different conditions, e.g. by activating/deactivating optimization strategies
4. Collections of measurements for the key performance indicators

*“The main indicator driving the PoC is the total energy usage of the physical servers during the execution of each test run”*

The reference workload selected for the PoC resembles the typical usage pattern of a Private Cloud data center, identified by analyzing the resource usage data collected in real customer deployments. The adopted workload has been selected to highlight the capabilities of a data center energy optimization platform, because the resource usage trend follows the “working hours” pattern: i.e. it is higher during the day and lower during the night, hence providing many opportunities for optimization during low-peak hours. For our evaluation we considered a two-days timeframe consisting of a typical Friday-Saturday usage pattern. In addition to the optimizations allowed by the fluctuations in the resource usage during a working day, this choice allows to exploit the additional ones that can be obtained in a non-working day, when the resources are typically under-utilized.

A synthetic workload with the above mentioned characteristics has been generated using [LoginVSI](#). The evaluation of the platform is based on the comparison of the performance data (with focus on power consumption) obtained without/with the optimizations performed by Eco4Cloud.

The main indicator driving the PoC is the total energy usage of the physical servers during the execution of each test run. The consumed power is measured through the VMWare APIs which retrieve data from the UCS blade servers.

### 3. Resources utilization before using Eco4Cloud

In this section we report the performance indices obtained by running the synthetic workload without activating the Eco4Cloud consolidation process. Figure 2 reports the CPU utilization of the 5 servers during the monitored period. The workload of the first day represents that of a regular working day, while the second represents a week end day. The percentage of CPU used by each server with respect to the overall CPU capacity of the same server is reported. Similar workload profiles were observed in most PoCs performed by Eco4Cloud. The figure shows that the CPU utilization is far from optimal. Indeed, most servers are heavily under-utilized (utilization  $\approx 35\%$ ) for long intervals. Clearly, this leaves significant room for efficiency improvement, as the next section highlights.

## Eco4Cloud on Cisco UCS: Unified Efficiency

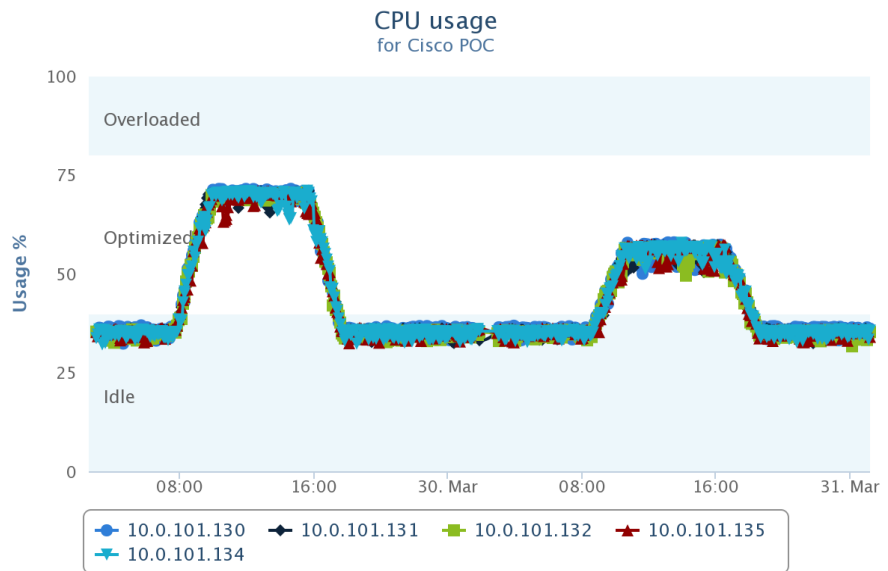


Figure 2 - CPU utilization of the 5 UCS Hosts

Figure 3 shows the RAM allocation (again, in percentage with respect to the total amount of RAM) of the 5 servers. Also, the RAM is under-utilized leading to considerable energy waste.

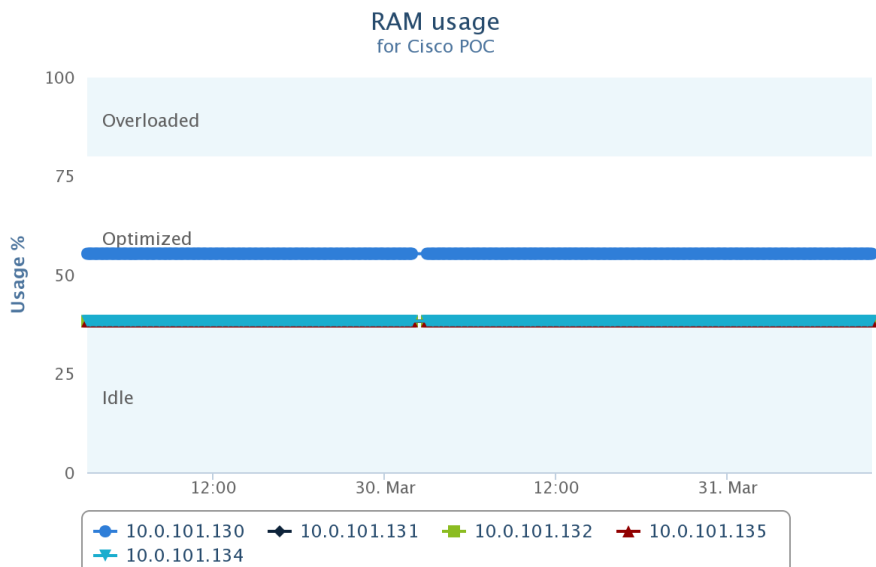


Figure 3 - RAM utilization of the 2 UCS Host

## Eco4Cloud on Cisco UCS: Unified Efficiency

Figure 4 reports the overall power consumption of the data center. No server is ever switched off. The observed variations are due to the corresponding variations of the overall workload. **The overall energy consumption is equal to 44.90 kWh.**

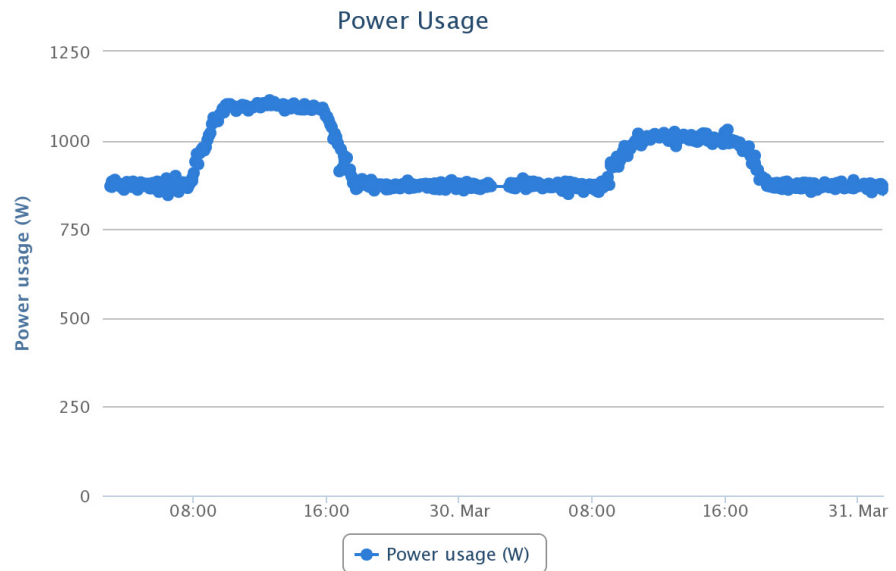


Figure 4 - Power consumption of the data center

For the evaluation of the **quality of service**, VMware suggests to compute two indices:

- the *CPU ready time*;
- the amount of *ballooned memory*.

VMware has set two thresholds for the CPU ready time of a VM: 5% is a warning threshold, while 10% represents an alert. In these tests, the CPU ready time kept constantly below 5% for all the VMs. Moreover, VMware specifies that the presence of ballooned memory should be avoided. In these tests, no event of ballooned memory was detected.



*“Results achieved in this IVT are in line with those of all other Eco4Cloud deployments”*

### 4. Resources utilization when using Eco4Cloud

The synthetic workload was resubmitted to the data center and at the same time Eco4Cloud was activated to consolidate the VMs.

Figures 5 and 6 show, respectively, the CPU and RAM utilization of the servers when using Eco4Cloud. These figures are to be compared to Figures 2 and Figure 3 (CPU and RAM utilization in the real test without Eco4Cloud). Two observations can be done immediately:

- 1) Eco4Cloud consolidates the workload effectively, as expected, and one or two servers are switched off to save energy in periods of under-utilization
- 2) Resources are utilized much more efficiently with respect to the workload observed in Figures 2 and 3 above

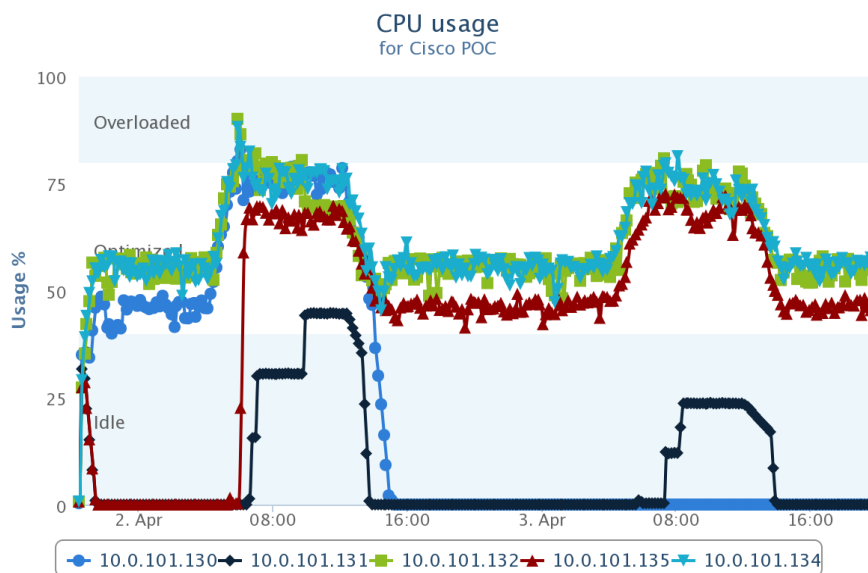


Figure 5 - Performance test with Eco4Cloud: CPU utilization of the 5 servers

Figure 7 shows the consumed power, which must be compared to the previous index computed before using Eco4Cloud (ref. Figure 4). Figure 8 reports the number of migrations per hour executed by Eco4Cloud in order to consolidate the load (low migrations) and prevent overload events (high migrations). It is noticed that during the ascending phases of load dynamics, high migrations are performed in order to balance the load and prevent overload events. During the descending phases instead low migrations are performed in order to off-load and hibernate some servers.

## Eco4Cloud on Cisco UCS: Unified Efficiency

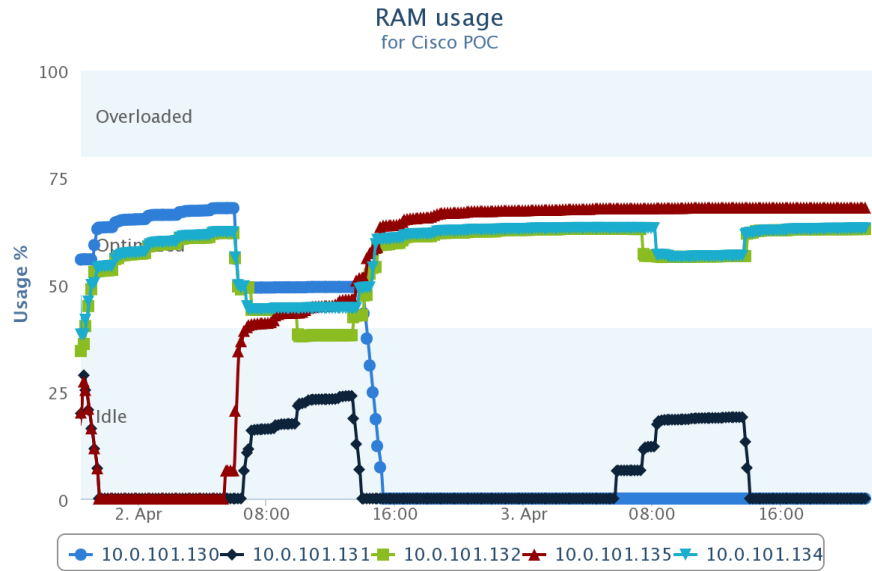


Figure 6 - Performance test with Eco4Cloud: RAM utilization of the 5 servers

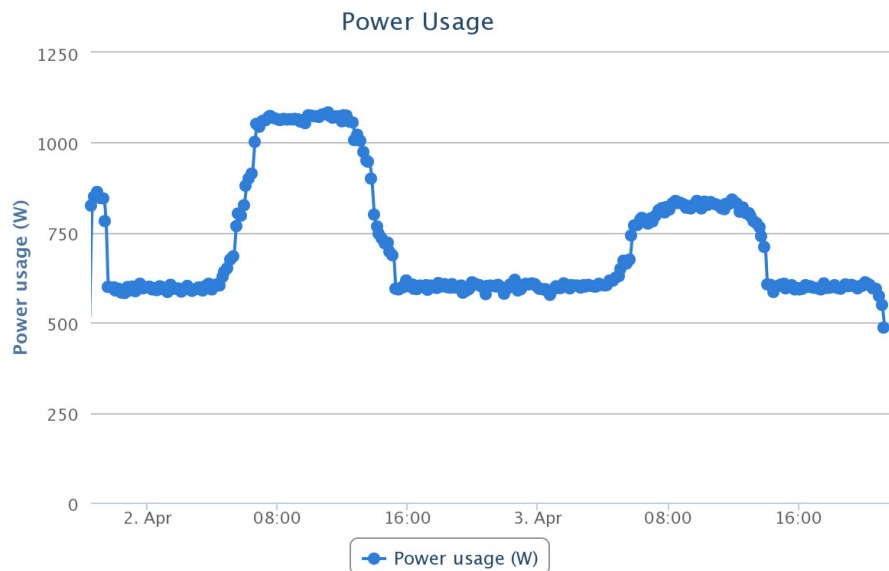


Figure 7 - Performance test with Eco4Cloud: consumed power

The CPU ready time and the amount of ballooned memory are referred to as measures of the quality of service. No event of ballooned memory was detected.

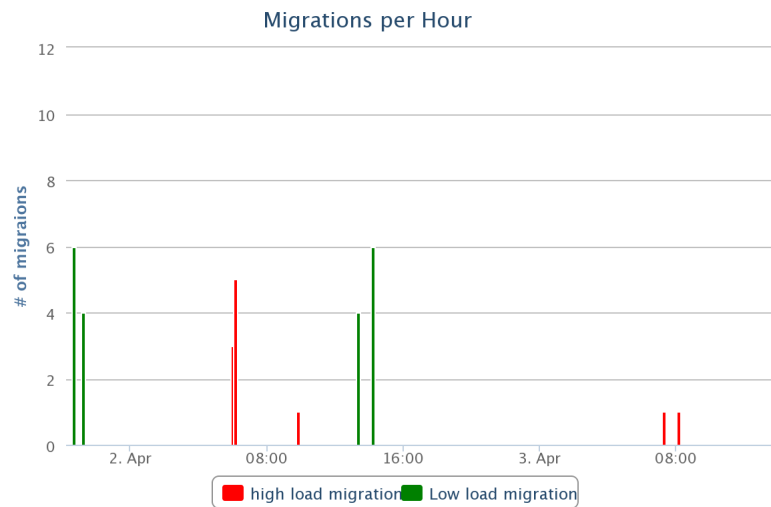


Figure 8 - Performance test with Eco4Cloud: number of migrations per hour

## 5. Test results and conclusions

The results of the test performed on **5 UCS Blade servers** show that Eco4Cloud consolidates the workload on a number of the originally active physical servers that ranges from 3 to 5 in working days, and ranges from 3 to 4 in week-end days. Hence, the results achieved in this IVT are in line with those of all other Eco4Cloud deployments.

## References

[1] C. Mastroianni, M. Meo, and G. Papuzzo, "Probabilistic consolidation of virtual machines in self-organizing cloud data centers," *Cloud Computing, IEEE Transactions on*, vol. 1, no. 2, pp. 215–228, 2013.

[2] International PCT patent: "[System for Energy Saving in Company Data Centers](#)", Italian National Research Council (CNR) and University of Calabria. Info and claims on <http://patentscope.wipo.int>, patent ID: #WO/2013/021407, May 2013.

### For more information

- Eco4Cloud: [www.eco4cloud.com](http://www.eco4cloud.com)
- Cisco UCS: [www.cisco.com/go/ucs](http://www.cisco.com/go/ucs)
- Cisco Marketplace Solutions Catalog: <https://marketplace.cisco.com/catalog/companies/eco4cloud-srl>
- Whitepaper: '[Saving energy in datacenters with Eco4Cloud and Cisco UCS](#)'