
Energy-based cost model of containers provisioning on clouds

Aline S. Moreira, Charles C. Miers,
Guilherme P. Koslovski and Maurício A. Pillon*

Graduate Program in Applied Computing,
Santa Catarina State University,
Joinville, Santa Catarina, Brazil
Email: aline.moreira@edu.udesc.br
Email: charles.miers@udesc.br
Email: guilherme.koslovski@udesc.br
Email: mauricio.pillon@udesc.br

*Corresponding author

Nelson M. Gonzalez

IBM Watson Research Centre,
Yorktown Heights, New York, USA
Email: nelson@ibm.com

Abstract: Cloud computing revolutionised the development and execution of distributed applications by providing on-demand access to virtual resources. Containerisation simplifies management and support of the cloud infrastructure and applications. Clouds typically are consumed in a pay-as-you-go pricing model. However, when applied to containerised environments, such traditional models do not consider resource utilisation values, leading to inaccurate estimates. Moreover, these models do not consider energy consumption, a dominant component of the data centre's total cost of ownership. This paper proposes Energy Price Cloud Containers (EPCC), a cost model based on energy consumption that accounts for containers' effective resource utilisation. We compare EPCC with AWS Fargate to highlight the benefits of using an energy-based pricing model. Thus, by comparing the cost of an application running using AWS Fargate with the estimated cost of that application in nome, it is possible to identify the benefits of using an energy-based pricing model. The weekly costs estimated when running computational resources at nome vary between US\$ 2.31 and US\$ 10.59. In contrast, when estimating the same amount of resources on AWS Fargate, the costs vary between US\$ 2.71 and US\$ 29.94. Nome resulted in a cost reduction of up to 35%.

Keywords: pricing model; containers; cloud computing; energy consumption.

Reference to this paper should be made as follows: Moreira, A.S., Miers, C.C., Koslovski, G.P., Pillon, M.A. and Gonzalez, N.M. (2022) 'Energy-based cost model of containers provisioning on clouds', *Int. J. Grid and Utility Computing*, Vol. 13, No. 6, pp.607–623.

Biographical notes: Aline S. Moreira received her Master's degree from Santa Catarina State University (UDESC) and Bachelor's degree from the Educational Society for Santa Catarina (UNISOCIESC).

Charles C. Miers is Professor at Santa Catarina State University (UDESC) in Joinville/SC and Member of the LabP2D (Laboratory of Parallel and Distributed Processing) of UDESC which has a private OpenStack cloud. He received his Doctorate degree in Computer Engineering from University of São Paulo (USP), the Master's degree in Computer Science from the Federal University of Santa Catarina (UFSC) and Bachelor's degree in Data Processing from the Santa Catarina State University. He was Security Consultant at LockNet Security Solutions, in the Software Area, from 1999 to 2003, having worked on projects of national companies (private and public) and multinationals.

Guilherme P. Koslovski is a Professor at Santa Catarina State University (UDESC) in Joinville/SC, and Member of the LabP2D (Laboratory of Parallel and Distributed Processing) of UDESC which has a private OpenStack cloud. He received his Doctorate degree from The École Normale Supérieure de Lyon at Lyon/France, Master's degree from Federal University of Santa Maria (UFSM) and Bachelor's degree in Computer Science from the UFSM.

Maurício A. Pillon is Professor at Santa Catarina State University (UDESC) in Joinville/SC and Member of the LabP2D (Laboratory of Parallel and Distributed Processing) of UDESC which has a private OpenStack cloud. He received his Doctorate degree from the Institut National Polytechnique of Grenoble at France, Master's degree from Pontifical Catholic University of Rio Grande do Sul (PUC-RS) and Bachelor's degree in Informatics from the Regional University of the Northwest of the State of Rio Grande do Sul. He also received his Post-doctorate degree from the Federal University of Rio Grande do Sul (UFRGS), conducted in the Research Group – Parallel and Distributed Processing Group (GPPD).

Nelson M. Gonzalez received his BSc, MSc and PhD degrees in Electrical/Computing Engineering from the University of Sao Paulo, Escola Politecnica (Poli-USP), Brazil in 2011, 2013 and 2015, respectively, and Post-doctoral degree from IBM Thomas J. Watson Research Centre. Currently, he works at the IBM Thomas J. Watson Research Centre in Yorktown Heights, New York, USA.

1 Introduction

Virtualisation technologies are in constant evolution. Containerisation, based on OS-level virtualisation (Sharma et al., 2016; Souppaya et al., 2017), provides scalability and flexibility to develop and deploy applications while simplifying management and adapting to customers' needs. Cloud services leverage virtualisation technologies at scale. Business models are based on dynamic on-demand resource allocation, pay-as-you-go pricing and resource consumption (Bindu et al., 2018). Several cloud providers offer elastic virtualisation services, e.g., Amazon AWS, Microsoft Azure, Google Cloud, Rackspace, Heroku (Begum and Khan, 2011). An increasing number of organisations rely on cloud computing as their core pool of resources to serve their own customers. Cloud-related expenses represent a considerable part of the IT organisation budget (Columbus, 2017). A precise specification of resources to be allocated and provisioned is crucial to improve performance and costs while addressing each service's complexities (Bittencourt et al., 2018; Wu et al., 2019).

IaaS providers adopt various pricing models based on monthly contracts, licenses and SLA policies. For example, Amazon, Microsoft Azure and Google Cloud Platform offer pay-as-you-go services charging a predetermined price for each resource, per task or by the number of calls of a service (Wu et al., 2019). From the cloud provider perspective, DC energy consumption stands out among all operational costs (Danilak, 2017). This directly impacts the costs seen by customers in the cloud service catalogue. Virtualisation contributes to efficient energy management through resource consolidation and isolation techniques, improving DC resource utilisation (Comerford, 2015). Compared to traditional hypervisor-based virtualisation, containerised infrastructures demonstrate even higher efficiency by sharing the available physical resources and the OS kernel.

Virtualisation technologies and application settings influence the utilisation level of computational resources such as CPU, memory, networking, disks, etc. For instance, bandwidth and latency requirements may come from different network settings used by an application (i.e., NAT, bridge, host-only (Danilak, 2017; Mentz et al., 2020)) as well as from the properties of the workload. Consequently, these aspects

also have a direct impact on energy consumption. Cloud customers typically face a trade-off between reducing resource consumption (i.e., cost) and improving performance. From the customer perspective, reducing the energy consumption is not a priority since there is no transparency about how energy resources are consumed or any financial reward (Souppaya et al., 2017; Hinz et al., 2018; Kurpicz et al., 2018).

We propose nome brings energy efficiency into the equation by defining an adaptable pricing model that accounts for energy consumption as a resource utilisation function, not just fixed time-based and allocation-based quotas. This provides the means to establish a clear financial incentive to save energy since it directly ties energy consumption to service costs. Our assessment of the impact of container utilisation on energy consumption reveals that each resource has a different impact on the final estimates. Energy consumption rapidly increases proportionally to the percentage of CPU utilisation. In contrast, memory and network observe a milder increase. Storage resources typically display stable energy consumption, even for I/O-intensive applications. This fine-grained analysis led to the stratification of containers' energy consumption. Nome formalises these components in a precise and extensible pricing model.

The contributions of this work are threefold: (i) an energy-aware cost model to containers; (ii) stratification of container energy consumption and (iii) a comparison of energy consumption among servers, virtual machines and container environments.

We compare estimated energy costs from AWS Fargate with nome cost model. AWS Fargate follows a pay-as-you-go charging model that does not consider the resource utilisation levels of leased containers. Since the AWS Fargate invoice comprises all business costs, we consider that AWS Fargate energy costs represent 5%, 10% and 15% of the total business costs (Hinz et al., 2018).

Nome estimated weekly values vary between US\$ 2.31 and US\$ 10.59, while AWS Fargate US\$ 2.71 and US\$ 29.94, a reduction of up to 35.39%. The AWS Fargate account was only cheaper than our nome model when we estimated that the cost of energy in the AWS Fargate model represents 5% of the total value of the contracted service. However, when the container is allocated and remains for a long period without physical resource demand, the server

usage is set to idle. The entire allocation period is counted with the total contract for AWS Fargate, even the server usage idle period. Thus, a container application is hampered by the AWS Fargate cost model if the allocation resource is overestimated.

This paper is organised as follows. Section 2 addresses background and prior work on cost models for virtualised resources. Section 3 describes and applies the method to provide a breakdown of energy consumption based on container application resource utilisation. Section 4 presents nomenclature. Section 5 presents a case study to exercise EPCC using the AWS Fargate pricing table. Section 6 concludes and proposes future work.

2 Cost models for virtualised resources

The consolidation of computational resources directly influences the DC energy reduction. Organisations have been migrating their services to the cloud computing service model since its emergence. Organisations still might maintain their DC to ensure strict control of management policies or because of legacy software. Although the concentration of physical resources on large DC reduces global energy consumption (Bawden, 2016), they still are an energy bottleneck.

In 2015, the energy consumption reached 416 TWh with a forecast to double every four years (Bawden, 2016). In the USA, the DC consumed more than 90 TWh (Danilak, 2017). The required energy supply system causes an economic and environmental impact. Thus, green clouds, renewable energy supplies and cost models have become interesting topics for scientific research (Hu et al., 2013; Pandikumar et al., 2012; Jain et al., 2013; Zhang et al., 2018; Garg and Buyya, 2012; Sharma et al., 2017).

Energy cost represents about 50% of the DC operating expenses (Guitart, 2017; Comerford, 2015). Providers constantly look for new methods, equipment and practices to improve their DC energy efficiency. In contrast, cloud tenants typically are not aware of the energy consumed by their applications, nor are they rewarded for improving it.

2.1 DC servers virtualisation

Virtualisation allows servers to be consolidated, therefore improving energy usage and simplifying DC management (Kominos et al., 2017). Although it exists since the mid-1960s, virtualisation has seen several improvements specially after its adoption by cloud infrastructures to transparently provide resources such as CPU processing, storage and virtual environments. This allows IaaS providers to grant access to remote, configurable, and shared sets of computing resources. IaaS providers can efficiently afford and make resources available to their tenants while minimising management, communication and energy costs (Mell et al., 2011; Hammadi and Mhamdi, 2014).

There are two primary technologies used to implement virtualisation: (i) *hypervisor-based*, which requires each VM to load a complete OS (leading to higher overhead) and (ii) *container-based*, a lightweight approach that shares the host OS kernel among all instances (Souppaya et al., 2017).

Containers provide an abstraction at the application layer by packaging the code and the necessary dependencies for its operation (Da Silva et al., 2018). Figure 1 presents some of these virtualisation environments adopted by organisations, including bare metal (a), VM Type-1 (b), container on bare metal (c) and container atop VM (d). Using containers does not restrict the use of VM; in fact, it allows the encapsulation of applications in several layers (Estrada et al., 2014). Docker is one of the most prominent containerisation technologies. Docker started in 2016, reaching US\$761 million in market value and this market could grow more than 35 times up to 2020 which would represent US\$27 billion (Kim et al., 2018). Gartner (Moore, 2020) expects that up to 15% of enterprise applications will run in a container environment by 2024, up from less than 5% in 2020, hampered by application backlog, technical debt and budget constraints.

The transition from internally hosted services to cloud computing can result in 4–5x improvement in resource utilisation efficiency (Wu et al., 2019) – up to 6x with container-based virtualisation. The adoption of virtualisation is, therefore, a solid technological trend – but what about energy? How to factor in economic aspects such as supply and demand while encouraging tenants to use DC resources consciously?

2.2 Cost models for VM

The main pricing models used by cloud providers consider value, supply and market. The value aspect focuses on customer demand, while the supply aspect focuses on cost. The market-based model seeks a balance between supply and demand. Cloud providers offer a vast range of services comprising infrastructure models and ways to allocate, use, and reserve resources. This evolution was naturally followed by creating cost models suitable to an extensive range of tenants with different profiles (Wu et al., 2019). However, despite constant improvements to products and services, current cost models still have gaps.

Figure 2 presents the evolution of AWS cost models. The X-axis represents the number of sales for each service model, q . The Y-axis represents the service pricing per unit, p . Over the years, AWS has designed new service and cost models focusing on the market and operating costs to fill provisioning gaps. We identify gaps in cloud computing in two distinct moments: (i) in its early days (2009), when AWS concentrated its resources on centralised DC infrastructures and (ii) in current days (2019), when AWS offers geographically distributed and fully decentralised infrastructures. Supply and demand was the foundation of the core services, providing discounts up to 90% compared to the same server in the AWS standard offer model. In 2015, AWS started offering a specific model for executing fault-tolerant workloads. In the same year, AWS products portfolio included two other pricing models for the spot service: *spot blocks*, in which instances have a specified duration; and *spot fleets*, sets of spot instances that run based on specified criteria. This flexibility allowed tenants to specify their resource needs with precision, reducing the resources' underutilisation (Wu et al., 2019).

Figure 1 Common approaches for providing computational resources to an application

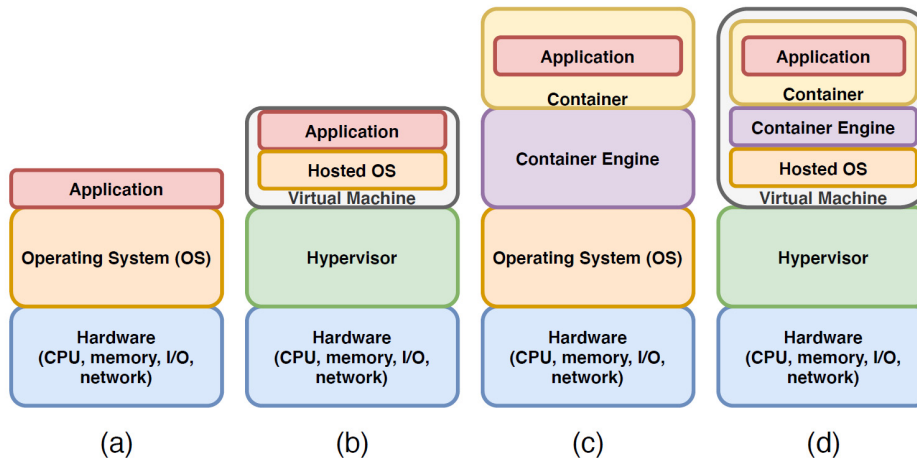


Figure 2 reveals that AWS had notable success in reducing under provisioning – white areas in the diagrams. Competition from other providers such as Microsoft and Google led to creating of new and more efficient offerings. The existence of low-cost instances opened a new market niche for customers with relaxed requirements. From the providers’ viewpoint, this increased the DC utilisation rate by using of idle resources.

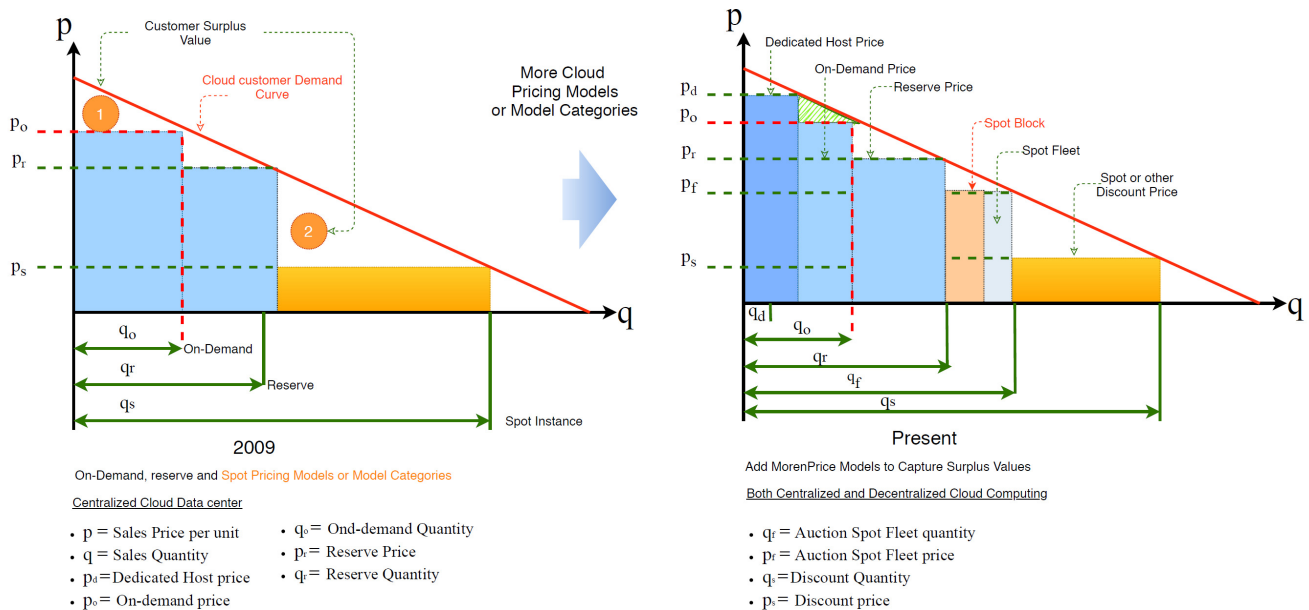
Although the evolution in service offerings is evident, there still are gaps to be addressed. This evolution is guided by providers’ interests, policies and perceptions of cost and service package models. The information from tenants usually is disregarded, e.g., the optimisation of their applications for lower energy consumption or politically correct actions, such as environmental preservation. Rewarding tenants for being aware of their energy

consumption can be a strong motivator for optimising applications. This paper highlights the need to fill the price gap per reservation, described in Figure 2 by the green hatched area in the *Present* years diagram.

2.3 Related work

There are several researches in the specialised literature regarding energy consumption in computational environments. These studies examine different environments (e.g., bare metal servers, VM, containers), manipulate various resources (e.g., CPU, network, memory, storage) in addition to having several purposes (e.g., monitoring tool, predictions). Table 1 summarises the related work. In short, the proposals have different characteristics, but none of them meet all the resources addressed by nome.

Figure 2 AWS costs models evolution



Source: Adapted from Wu et al. (2019).

Table 1 Models and tools from the perspective of energy consumption of VM in DC or IaaS cloud providers

<i>Author</i>	<i>Objective</i>	<i>Environment</i>	<i>Resources</i>	<i>Objective</i>
Hinz et al. (2018)	Analysis of energy costs in VM and hypervisors	IaaS	CPU, network traffic	Pricing
Kurpicz et al. (2018)	Proportional energy cost model in environments based on cloud providers	IaaS	CPU	Pricing
Bhattacharya et al. (2013)	Energy consumption measurement tool	DC	CPU, network, memory and storage	Monitoring
Zakarya and Gillam (2018)	Relates the energy efficiency of workloads to CPU consumption	Cloud providers	CPU and memory	Monitoring
Brondolin et al. (2018)	DEEP-mon hardware energy monitoring of containers	OS	CPU, network and memory	Monitoring
Leitner et al. (2016)	Container deployment cost model	IaaS and PaaS	CPU, memory, and storage	Pricing
nome	Container cost model based on energy consumption	IaaS	CPU, network, memory and storage	Pricing

Hinz (2018) introduced pricing models for IaaS clouds focusing on VM and guided by energy consumption. PSVE bases its pricing model on identifying the VM individualised energy consumption through the hypervisor, accounting the resources individually and collectively. Collective costs are prorated by VM according to the number of CPUs allocated in the period (Hinz, 2018). This benefits the customer's lowest cost and the providers' PTC and goes to green IT actions. However, PSVE does not consider the relevance of containers in the computational clouds and restricts the proposed model only to VM. EPAVE (Kurpicz et al., 2018) is a pricing model for accounting VM dynamic consumption and the proportional static cost of the cloud infrastructure, assigning general energy costs per VM. EPAVE and PSVE differ in the energy costs division policy and do not consider the network traffic energy consumption. A set of techniques to measure energy consumption was described by Bhattacharya et al. (2013) and (Zakarya and Gillam, 2018). Regarding the workload impact on energy consumption, Zakarya and Gillam (2018) claimed that the impact on energy efficiency varies according to CPU models.

Finally, two works stand out associated with virtualisation-based on OS and containerised applications: (i) the DEEP-mon tool which monitors the CPU, network and memory resources at the level of OS (containers) to measure energy consumption (Brondolin et al., 2018); and (ii) a pricing model for deploying cloud applications based on micro-services (Leitner et al., 2016). A research carried out in 2019 shows that more than 34% of companies that invest over US\$100 thousand in technologies are adhering to the container environments atop bare metal servers, migrating the applications of VM for containers (DIAMANTI, 2019). These companies aim to increase their performance, reduce infrastructure complexity, improve flexibility and reduce costs. Containers are efficient, allowing operators to do more with less and manage critical applications with high availability. The usage of containers by companies ranges from modernising legacy applications to analysing big data (DIAMANTI, 2019). Such research highlights the growing impact of containers for organisations, attesting the scientific

gap in the containers pricing models study for cloud providers guided by energy consumption (one of the motivations of this work).

Cloud providers have been improving their service offered over the years, filling the gaps of underutilised processes and seek to maximise the energy efficiency in their DC (Wu et al., 2019). However, on the client side, the providers' efforts fall short in the cost models offered. The cost models offered by providers do not provide economic benefits to their tenants to develop energy-conscious applications. A better cost model would benefit not only the customer's final cost but also environmental preservation. This action reflects the reduction of PTC of providers making the offer convenient for the provider and their tenants. The DC must offer cost models that consider the energy consumption in their pricing, as proposed by PSVE (Hinz et al., 2018). However, this model only includes hypervisor-based environments, excluding containers in bare metal and therefore having lower support for containers per physical server in VM environments (DIAMANTI, 2019). Moreover, providers offer container-based services priced by resource allocation, not utilisation. This pricing and provisioning gap requires an appropriate cost model for containers focused on energy consumption so that tenants are encouraged to develop energy efficient applications. State of art in energy-aware consumption reinforces the importance of a pricing model for the IaaS context and energy monitoring methods for containerised applications. The focus of this work is to identify and address this gap regarding the energy-aware containerisation-based pricing model.

3 Energy consumption analysis

This section presents an analysis of the energy consumption of computational resources. The objective is to establish a baseline to quantify and compare each resource's energy consumption impact (CPU, memory, storage and network) in bare metal, VM and container-based setups.

3.1 Experimental protocol

We elaborated the experimental environment using the Grid5000 platform (GRID5000, 2020) Grenoble site (<https://www.grid5000.fr/w/Grenoble:Hardware>). The selected cluster, *Yeti*, comprises four Dell PowerEdge R940 computing nodes with Intel Xeon Gold 6130 (Skylake, 2.10GHz, 4 CPUs/node, 16 cores/CPU), 768Gb RAM and an Intel Ethernet Controller X710 for 10GBE SFP+. A set of non-intrusive wattmeters monitors *Yeti*'s energy consumption. The software comprises three layers: (i) Debian 10 GNU/Linux bare-metal installation; (ii) KVM hypervisor 3.1.0 and (iii) Docker 19.03.

We used four benchmarks to analyse energy consumption. Each benchmark applies a different workload to a specific resource. *StressNG* applies a configurable workload to the server's CPU (King, 2019). *Stream* performs vector calculations to saturate the server's memory and CPU. The user defines the amount of memory consumed, CPU load and a number of processes to be allocated (McCalpin et al., 1995). *Fio* (flexible I/O) executes asynchronous read/write storage operations (Axboe, 2017). Finally, *iperf3* generates network traffic to stress bandwidth usage, allowing parallel communication flows between server and client (Mortimer, 2018). These benchmarks allow us to observe each resource's behaviour and its correlation to energy consumption (Li et al., 2017).

We quantify the energy consumption of each physical resource for each active container. Our baseline is the bare

metal setup (see Figure 3(a)). This baseline is compared to three other scenarios: (i) VM, (ii) Docker container and (iii) Docker container inside VM, as shown in Figures 3(b) to 3(d). Each benchmark was executed with the same workload 10 times.

3.2 CPU analysis

Figure 4(a) illustrates the energy consumption behaviour according to the increase in CPU usage. The *X*-axis represents the number of threads (CPUs) assigned to the benchmark. The *Y*-axis represents the energy consumed (in joules) during the experiment. The idle consumption is around 175 J (1 J standard deviation) for all platforms. For each platform, we increased the number of benchmark threads from 1 to 128. Running one thread doubles energy consumption compared to idle. For 128 threads, energy consumption increases by more than 4x. Also, VM and container atop VM platforms consume more energy than bare metal and containers-only. Container atop VM consumption is slightly higher than VM only, except when the server is saturated. The energy consumption gap among platforms widens as CPU usage increases, mainly due to the management performed by the KVM hypervisor. Using only 1 CPU, the difference in consumption between environments is less than 2%. For 8 CPUs, the difference is higher than 9%. This reveals that CPU usage is the essential component for a cost model; two customers may pay the same price for a service, but one consumes 3x more energy than the other.

Figure 3 Experimental environment

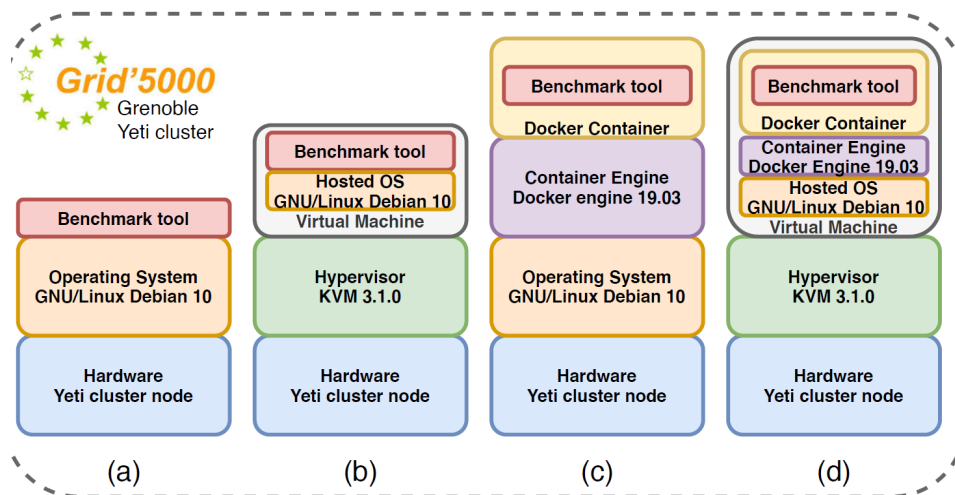
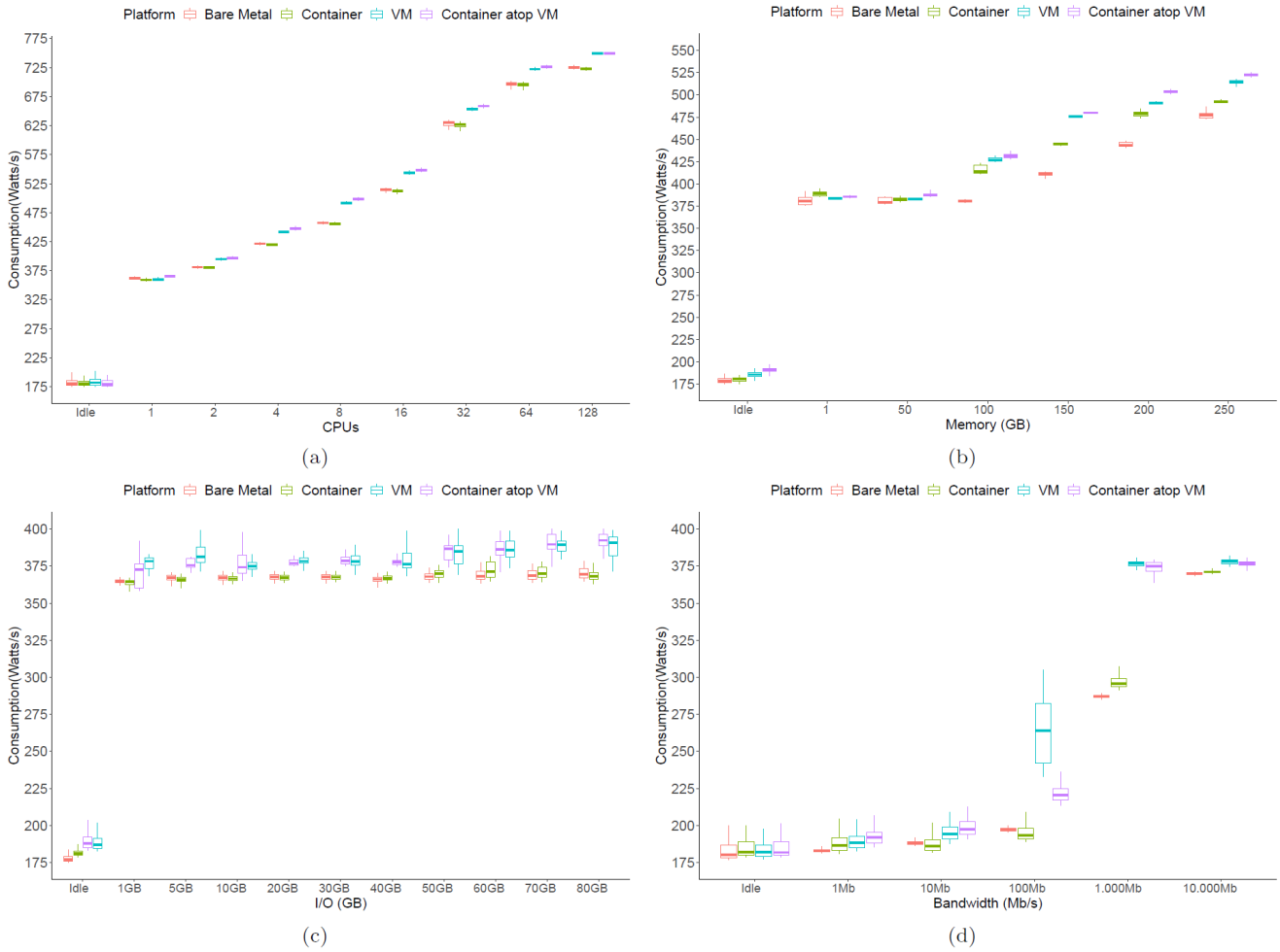


Figure 4 Physical resource energy consumption

3.3 Memory analysis

Figure 4(b) shows the results for the experiments with memory utilisation. X -axis represents memory utilisation in GB. Y -axis indicates energy consumption in joules. The experiments with STREAM (McCalpin et al., 1995) use seven memory workloads processed by only one thread. We compare the baseline (idle system) to six memory utilisation values: 1 GB, 50 GB, 100 GB, 150 GB, 200 GB and 250 GB. The energy consumption curve is visually similar to CPU but with a narrower range, up to 525 J. The experiments with 1 GB show an energy consumption increase between 375 J and 400 J. Increasing memory utilisation to 50 GB, however, does not drastically change energy consumption. Also, subsequent experiments increase energy consumption by 7% on average compared to the previous one. This reveals that, although memory utilisation does influence energy consumption, its impact is not as pronounced as CPU.

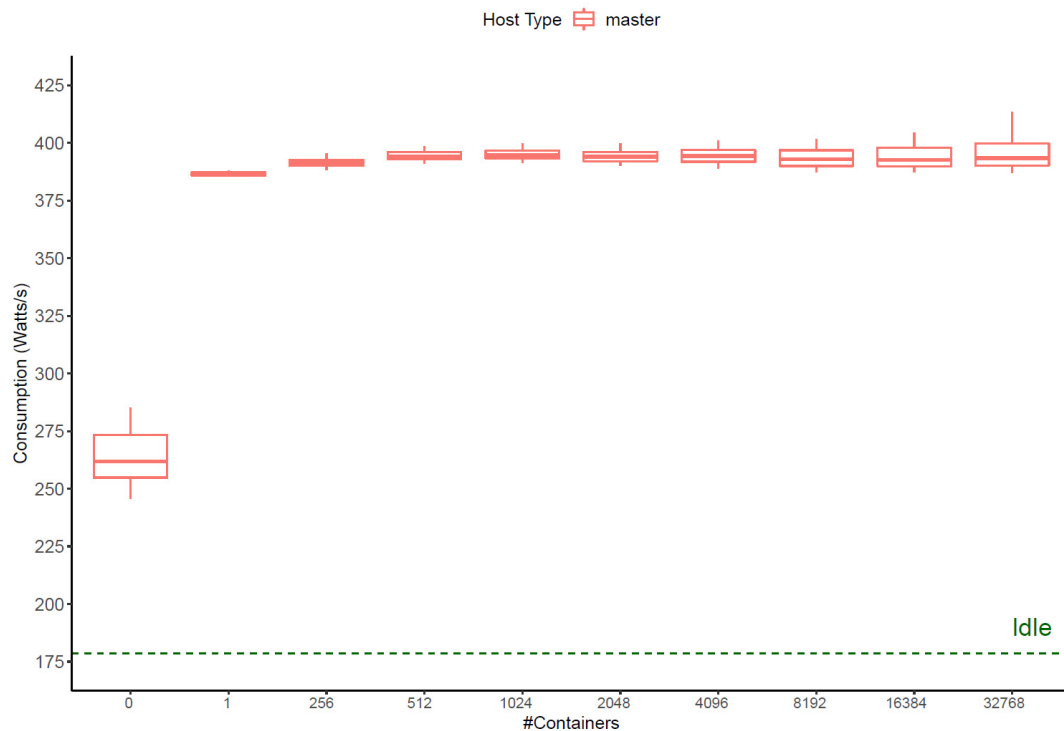
Memory management represents an important resource demand for all platforms. The energy consumption gap compared to bare metal is more accentuated than the one for

CPU resources. Not even the container setup can compete with bare metal. In the worst case, when comparing idle and 250 GB memory setup used for the container atop VM scenario, the energy consumption increases by 214%.

3.4 Storage analysis

Consumption on SSD drives tends to be much lower, but many DC still have hard drives. Thus, the storage analysis focuses on energy consumption from hard disks. Flexible I/O (Axboe, 2017) is the storage benchmark selected for the evaluation. The storage demand varies from 1 GB to 80 GB. In all cases, the benchmark uses only one thread.

Figure 4(c) depicts the behaviour of energy consumption with storage. X -axis represents the storage demand (in GB). Y -axis shows energy consumption in joules. Energy consumption from storage remains stable from 1 to 80 GB for all scenarios. The variation coefficient is around 8.4%. This reveals that increasing the storage demand does not impact energy consumption. Compared to idle, the energy consumption increases 2x.

Figure 5 Kubernetes orchestrator's energy consumption

3.5 Network analysis

The last resource analysed is the network. Network management in virtualised environments behave differently in terms of energy consumption depending on network type (e.g., host, bridge, overlay, macvlan) (Morabito, 2015). In our experiments, all scenarios were set up using bridges. We executed iperf3 (Mortimer, 2018) to create network workloads while evaluating energy consumption. We used two nodes: a client and a server. Energy consumption was measured on the server-side. We executed iperf3 with several different workload configurations: 1 MB, 10 MB, 100 MB, 1 GB and 10 GB.

Figure 4(d) shows the energy consumption observed through the execution of the experiments. *X*-axis represents the bandwidth settings applied in the benchmark. *Y*-axis represents the energy consumed (in joules) by the server. We observe a smooth increase in the beginning up until 10 Mbps – up to 28% compared to idle. For bandwidths higher than 100 Mbps, we noted a spike in energy consumption, up to 117% compared to idle. Network utilisation is, therefore, an important component of energy consumption, depending on its utilisation level.

3.6 Container orchestration analysis

So far, this section has analysed the influence of physical resources in energy consumption. However, large cloud

environments comprise thousands of containers that must be properly managed. Our experiments used a Kubernetes (2020) deployment comprising a single dedicated master. Even though Kubernetes is a remote service, the objective of this section is to see how customer demands impact the energy consumption behaviour of the Kubernetes master.

Figure 5 illustrates this behaviour. *X*-axis represents the number of active containers. The *Y*-axis represents the energy consumption in joules. Activating the orchestrator corresponds to an increase of 90 J (compared to idle). Adding the first container increases energy consumption to about 390 J. Consumption remains fairly constant up to 32,768 containers.

4 Containers cost model based on energy consumption

Energy costs in the composition of IaaS cloud expenses are undeniable. One of the main contributions of this work is the proposal and description of a pricing model for IaaS providers that considers containers' energy consumption. Tenants are encouraged to optimise their applications when priced based on energy consumption.

Table 2 Notation related to the container's energy consumption in a period of time

Notation	Description
P_w	Price for watts.
W_{idle}	Watts consumption of idle container server.
C_{idle}	Energy costs to maintain active idle server.
T_u	Server's available resources usage rate.
T_c	Server's available resources complement rate.
$ Cont_u $	Tenants containers in a server.
G_p	Providers expenses to maintain available and not allocated resources.
$ Cont_s $	Containers per server.
P_{cpu}	CPU resource price.
W_{cpu}	CPU resources consumed in <i>watts</i> by container.
C_{cpu}	Energy consumption related to physical CPU usage.
$U_{cpu}(t, c)$	CPU usage of a specific container <i>c</i> .
W_{mem}	Watts consumed by the container memory resource.
C_{mem}	Energy cost related to memory usage.
(t, c)	Given instant <i>t</i> of a specific container <i>c</i> .
$U_{mem}(t, c)$	Container's memory usage.
fvc	Variation factor of consumption by containers allocated on the server.
C_{net}	Energy costs related to network.
$U_{net}(t, c)$	Network usage of a container.
W_{net}	Watts consumption by network container resources.
C_{disk}	Cost related to input and output storage operations.
P_{disk}	Storage price.
Q_{disk}	Amount of storage used.
W_{orc}	Watts consumed by the container orchestrator.
C_{orc}	Container orchestrator's related energy cost.
$Cost$	Energy cost of a container over a period of time.
C_{total}	Energy cost related to the resources usage by containers of a given customer.

The proposed pricing model follows the notation described in Table 2. Initially, P_w represents the price charged by the cloud provider's electrical distribution company, which tax the provider by watts consumed and according to the power's supplement region. W_{idle} describes the energy consumption in watts for a given server in idle status (i.e., without hosting containers). C_{idle} is the server's minimum energy cost, or in other words, is given by minimal energy that the provider expends to keep it active and available for the execution of containers. Each server has a different C_{idle} resulted from hardware-specific characteristics (technologies used, power supplies, CPUs, GPUs, drives, storage, among others). For this reason, the cost of C_{idle} changes depending on the

platform on which the containers run. Thus, each specific server has a fixed cost in idle and must distribute it through the proportional apportionment between clients and containers allocated on the server. Therefore, equation (1) defines the cost to keep the server active.

$$C_{idle} = P_w \times W_{idle} \quad (1)$$

For the proportional share of C_{idle} among the hosted containers, the split of W_{idle} is required. This apportionment consists of 5 elements:

- 1) T_u denoting the computing resources (the percentage of the resources allocated by clients that run containers on a server);

- 2) T_c , a fee for complementing the server's computational resources, applied to exempt the customer from incurring a higher cost than effectively allocated and used. This rate equalises the resource allocation percentage. When a single client allocates a server, and the containers do not use the server as a whole, the consumption of W_{idle} should not be considered 100% for pricing this client. The provider is also responsible for the minimum consumption and assumes a percentage, which we call T_c . Therefore, if the allocated tenants' rate T_u on the platform is less than T_c , the cost's apportionment is carried out between tenants and the provider.
- 3) G_p , a management cost attributed to the provider due to its responsibility to scale and manage the server computational resources, since keeping a specific machine available is considered a fixed cost;
- 4) $|Cont_u|$ is the element that counts and assigns the number of the allocated containers by each client on a server;
- 5) $|Cont_s|$ counts and identifies by unit the allocated containers on the server.

In summary, for a given tenant, the fraction of watts consumption in idle mode (W_{idle}) is given by equation (2).

$$W_{idle} = \begin{cases} |Cont_u| \times \frac{G_p \times (100 + (T_u - T_c))}{|Cont_s|} & \text{if } T_u < T_c \\ |Cont_u| \times \frac{G_p}{|Cont_s|} & \text{else} \end{cases} \quad (2)$$

The equation (2) verifies if the utilisation rate T_u is less than the complementary rate T_c . So the apportionment of this cost occurs between: the number of containers per tenants ($|Cont_u|$); the number of containers per server ($|Cont_s|$); and between the cloud provider (G_p). When T_u is greater than T_c , W_{idle} is apportioned only between $|Cont_u|$ and $|Cont_s|$.

From the identification of the server's cost in idle mode, it is necessary to understand the costs of each resource used by containers (CPU, memory, network, storage, and orchestration).

The CPU price is given by the P_{cpu} element, representing the price paid to the energy distribution company for the server's CPU consumption. This element consists of W_{cpu} , which identifies the number of watts consumed by CPU, and is composed of the elements P_w and C_{idle} , the watts price and the idle server's cost, respectively, as given by equation (3).

$$P_{cpu} = (W_{cpu} \times P_w) - C_{idle} \quad (3)$$

Equation (3) enables the formulation of the total CPU cost (C_{cpu}) at a given time interval (T_1, T_2). The total CPU cost

is composed of the element $U_{cpu}(t, c)$, which is equivalent to the CPU usage at a given time t of a specific container (c), and by the allocated containers sum per server. So the total CPU cost is expressed in equation (4).

$$C_{cpu} = \sum_{t=T_1}^{T_2} \sum_{c=1}^{|Cont_u|} U_{cpu}(t, c) \times P_{cpu} \quad (4)$$

Following, W_{mem} reflects the memory consumption's energy, identifying the number of watts consumed by the memory resource. Also, $U_{mem}(t, c)$ is equivalent to the memory usage at a given time t of a specific container (c). Hence, it is possible to identify the energy cost of memory represented by the element C_{mem} , addressed in equation (5), which depicts the memory pricing in the container's execution.

$$C_{mem} = \sum_{t=T_1}^{T_2} \sum_{c=1}^{|Cont_u|} U_{mem}(t, c) \times ((W_{mem} \times P_w) - C_{idle}) \quad (5)$$

In turn, the W_{net} accounts for the consumption in watts an application has when using the network resource, highlighting that this resource has little significant variation in the number of active containers on the machine, maintaining a variation only to the network bandwidth used. Therefore, to calculate the element C_{net} it is necessary to analyse two elements in addition to W_{net} : (i) the $U_{net}(t, c)$ and (ii) P_{net} . The element $U_{net}(t, c)$ identifies the network usage at a given time t for a specific container (c), given by the size of the bandwidth used in the container's execution in a given period, through the sum of the time interval, as well as the same-server containers' sum. P_{net} is composed of fv_c , which refers to a price variation factor depending on the same-server existing containers' consumption. Therefore, P_{net} is detailed in equation (6) while C_{net} is expressed by equation (7).

$$P_{net} = ((W_{rede} \times P_w) - C_i) \times fv_c \quad (6)$$

$$C_{net} = \sum_{t=T_1}^{T_2} \sum_{c=1}^{|Cont_u|} U_{net}(t, c) \times P_{net} \quad (7)$$

The parameter fv_c represents the variation in energy consumption (i.e., even if not very significant, this work understands that it is relevant to assign it to the client), according to the number of containers running on the server. Thus, if the server has only one client running its containers, fv_c is 1. If there is more than one client running containers on the server, then the containers' number per client is divided by the total number of containers allocated on the machine and then added 1.

$$fv_c = \begin{cases} 1 & \text{if } |Cont_u| = |Cont_s| \\ \frac{|Cont_u|}{|Cont_s|} + 1 & \text{else} \end{cases} \quad (8)$$

Regarding the storage resource, the experimental analyses (see Section 3) demonstrated that the resource is constant in terms of energy consumption when there is only one container running in the server, regardless of the size of the disc I/O. However, when running more than one container per server, the storage's energy consumption is variable. In this sense, the storage cost C_{disk} (described in equation (9)) is composed by U_{disk} (i.e., the storage's size used, and for W_{disk} representing the consumed watts by the storage resource).

$$C_{disk} = \begin{cases} \sum_{t=T_1}^{T_2} \sum_{c=1}^{|Cont_u|} U_{disk}(t,c) \times (W_{disk} \times P_w - C_{idle}) & \text{if } |Cont_s| > 1 \\ \sum_{t=T_1}^{T_2} U_{disk}(t) \times (W_{disk} \times P_w - C_{idle}) & \text{else} \end{cases} \quad (9)$$

If the container number is higher than 1, then the storage cost is obtained by double sums. The inner sum considers all containers' storage usage requested by the customer, while the outer sum realises the sum of the general storage usage in a given interval. Otherwise, the cost is obtained only by the sum of the time interval accounting for the server allocated container's storage usage. In addition to the energy consumption of computational resources in active containers, there is also the container management's energy consumption. Container orchestrating technologies realises the management of these containers (i.e., Kubernetes, Swarm, Mesos, among others). Therefore, orchestration expenses must be accounted.

The container orchestration analysis (see Section 3) reveals the master node is not influenced energetically by the number of managed containers. For this reason, it is necessary that the energy sharing of this management takes place between all the orchestrated containers. The orchestration's energy pricing involves the utilisation and compensation fees, i.e., T_u and T_c , respectively. Also, the W_{orc} represents the watts consumed by the orchestrator, defining the cost of container orchestration (C_{orc}). Thus, C_{orc} is defined in equation (10).

$$C_{orc} = \begin{cases} |Cont_u| \times \frac{((W_{orc} \times P_w) - C_{idle}) \times \left(\frac{100 + (T_u - T_c)}{100} \right)}{|Cont_s|} & \text{if } T_u < T_c \\ |Cont_u| \times \frac{[(W_{orc} \times P_w) - C_i]}{|Cont_s|} & \text{else} \end{cases} \quad (10)$$

Equation (10) determines if the utilisation fee T_u is less than the complement fee T_c , then the apportionment of this cost occurs among the number of containers per tenants $|Cont_u|$ and enter the containers number per server $|Cont_s|$.

However, if T_u is greater than T_c , the apportionment of this cost occurs only between $|Cont_u|$ and $|Cont_s|$. Thus, if the cloud provider allocates a customer into a machine with no other user, the customer will not pay the full price of the active server. Therefore, equation (11) presents the container cost, where C_{total} is the cost of using resources per container for a given customer.

$$C_{total} = C_{idle} + C_{cpu} + C_{mem} + C_{rede} + C_{disk} + C_{orc} \quad (11)$$

The proposed cost model is simplified and stratified according to the elements that make up an OS virtualisation environment in IaaS cloud providers. Consequently, the calibration of the model components depends on the analysis of the behaviour of each resource.

5 AWS case study

AWS Fargate (2020) is a serverless compute service, allowing container execution without the need to provision servers from customer side. We selected this service due to its pricing model based on on-demand resource allocation. Tenants pay for the number of CPUs allocated, and the amount of memory, storage, and communication (data size) used (Vohra, 2018) over time. Tenants can adapt resource demands second by second. Unfortunately, in practice, container management optimisation is a complex task. Furthermore, tenants are responsible for monitoring resource utilisation and scheduling.

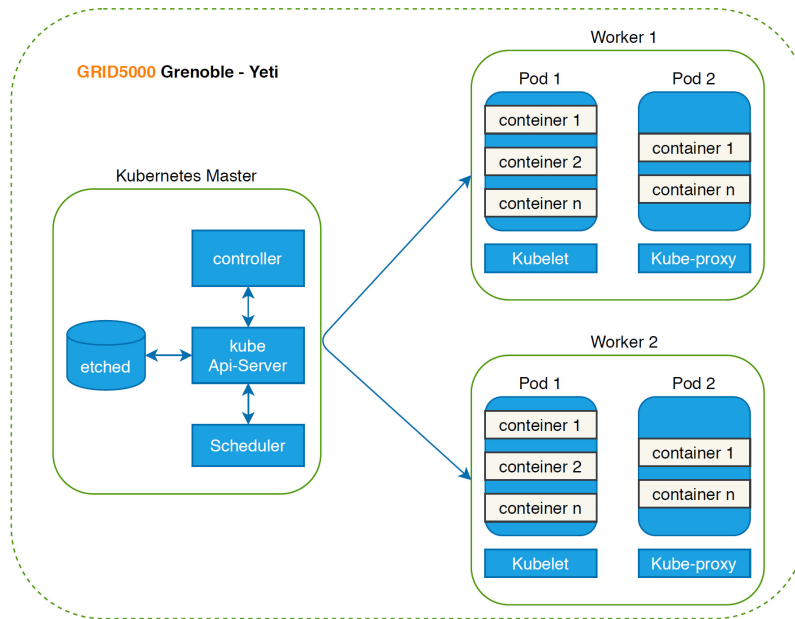
5.1 Energy consumption with Fargate settings

Section 3 depicts the energy consumption analysis of physical resources where the dedicated server has its resources saturated. Since AWS does not allow low-level access to physical resources, we built a similar platform using our Yeti cluster. We used Yeti to reproduce AWS Fargate's set of supported configurations (<https://aws.amazon.com/fargate/pricing/>). We used Docker as the containerisation platform and Kubernetes for management since Fargate also uses both technologies. The experimental environment consists of three servers. One server is reserved for the Kubernetes master orchestrator. The others servers run the Docker containers through the Kubernetes workers (see Figure 6).

5.1.1 AWS Fargate service offers

Table 3 presents the AWS Fargate offers analysed. We evaluate energy consumption for each resource at maximum utilisation. Our experiments consist of two benchmarks, StressNG and Stream, executing 21 different vCPU and memory consumption scenarios.

Figure 6 Containers management



Source: Adapted from Kubernetes (2020)

Table 3 AWS Fargate service offers

ID	CPU	Memory values
#1	0.25 vCPU	0.5 GB, 1 GB and 2GB
#2	0.5 vCPU	Min. 1GB and Max. 4GB, in 1 GB increments
#3	1 vCPU	Min. 2GB and Max. 8GB, in 1 GB increments
#4	2 vCPU	Min. 4GB and Max. 16GB, in 1GB increments
#5	4 vCPU	Min. 8GB and Max. 30GB, in 1GB increments

Source: Adapted from Amazon Web Services (AWS) (2020a).

VCPU energy consumption is obtained for each ID by keeping memory utilisation at idle while stressing CPU. Memory energy consumption is obtained with maximum CPU and memory utilisation.

Figure 7 Experiment results (a) CPU and memory mixed energy consumption (b) storage energy consumption

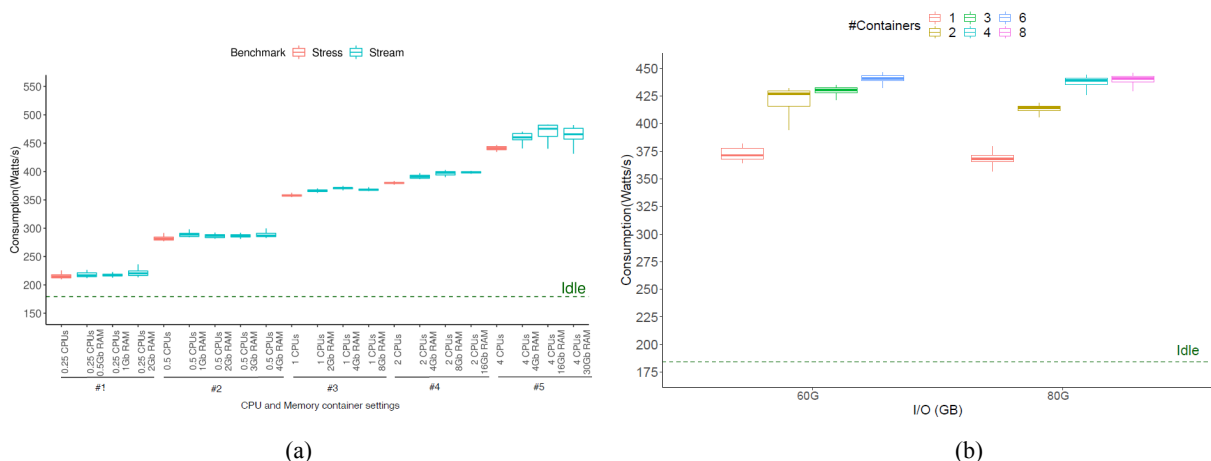


Figure 7(a) depicts energy impact from resources (CPU and memory) based on Fargate service offers. X-axis represents each scenario evaluated. Y-axis represents energy consumption in joules. The dashed green line represents the minimum energy consumption of the server with no workload; i.e., when there is a resource allocation but its applications do not require resources. All scenarios show a gap between the idle line and their energy consumption values. This gap reaches up to 300 joules. However, Fargate does not take this into account when charging customers.

5.1.2 Additional AWS Fargate services

AWS considers storage and data transfers as additional services. For the storage service (Amazon EBS Volumes), tenants have five types of volumes and pay per GB-month of provisioned storage, as described in Table 4. We selected the cold HDD (sc1) volume because it is similar to Yeti's hardware.

Table 4 Additional AWS offers

Storage services	Data transfers OUT
General Purpose SSD (gp2) Volumes	Up to 1 GB / Month
Provisioned IOPS SSD (io2) Volumes	Next 9.999 TB / Month
Provisioned IOPS SSD (io1) Volumes	Next 40 TB / Month
Throughput Optimised HDD (st1) Volumes	Next 100 TB / Month
Cold HDD (sc1) Volumes	Greater than 150 TB / Month

Source: Adapted from Amazon Web Services (AWS) (2020c, 2020b)

Storage devices consume energy to remain in service, regardless of I/O activity. Storage transactions (writes and reads, HDD on our experiments) demand additional energy. Amazon EBS volume service disregards storage transactions and, consequently, the energy consumption variation. Figure 7(b) depicts storage energy consumption considering disk transaction activities. The *X*-axis represents the amount of available HDD space to perform HDD-intensive transactions. The *Y*-axis represents the energy consumed in joules. We used the flexible I/O benchmark (Axboe, 2017) to generate HDD-intensive transactions. We varied the number of containers (from one to eight) and HDD space size (60 GB and 80 GB). The results show that the energy consumption behaviour is similar for both space sizes. However, the number of containers has a clear energy impact. For instance, eight containers consume about 70 J more than only one container. Compared to idle (dashed green line, without disk-intensive transactions), the gap reaches 260 J.

Regarding data transfers, Amazon EC2 offers free data transfers into Amazon (inbound). Outbound data transfers are charged according to Table 4. Amazon's data transfer service considers the energy consumption of network equipment to be fixed. However, network consumption may be high in one month and demand no resources in the next month.

We used iperf3 (Mortimer, 2018) to generate network load. Figure 8 presents client and server network flow behaviour. The *X*-axis represents the number of containers and bandwidth. The *Y*-axis represents the energy consumption in joules. We gradually increase network load from 1 MB up to server saturation at 10 GB. Then we oversaturate the server network using loads of 25 GB, 50 GB and 100 GB to verify its behaviour under high traffic. The results show a gradual increase in energy consumption up to 500 MB, with a maximum difference of 50 J when compared to idle server – an increase of 42%. Values higher than 1 GB behaved differently, showing non-linear growth in energy consumption and leading the energy consumption gap to a maximum of 270 J compared to the idle – an increase of 157%. If the application aims to reduce energy consumption, it is important to reduce communication between containers. It is possible to characterise the energy impact of the relevant storage resource by varying the number of containers that use storage. For the network resource, the energy impact is evident when the resource is saturated.

5.2 nome vs. AWS Fargate

This section compares nome to AWS Fargate model. The first scenario comprises two configurations with two and four CPUs.

Table 5 reveals that Fargate only considers the amount of CPU and memory provisioned, not their actual usage (percentages). The prices, therefore, are the same for 50% and 100% CPU resource utilisation. Nevertheless, one CPU (corresponding 50% of CPU provisioned demand) represents 5.57% less energy consumption than the 100% considered by AWS Fargate. The customer's bill must consider this difference. Since containers leverage a flexible platform, the cost model should also be flexible. A flexible cost model motivates customers to make their applications more efficient, encouraging energy savings and green IT practices.

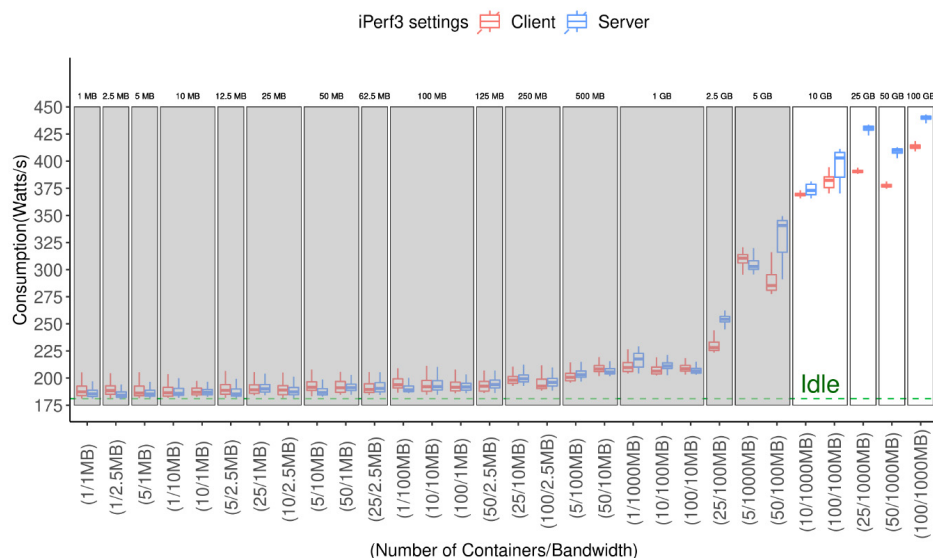
Figure 8 Network energy consumption

Table 5 AWS Fargate price

CPU's	Usage (%)	Price (US\$/hour)
2	50%	0.08
	100%	0.08
4	50%	0.16
	100%	0.16

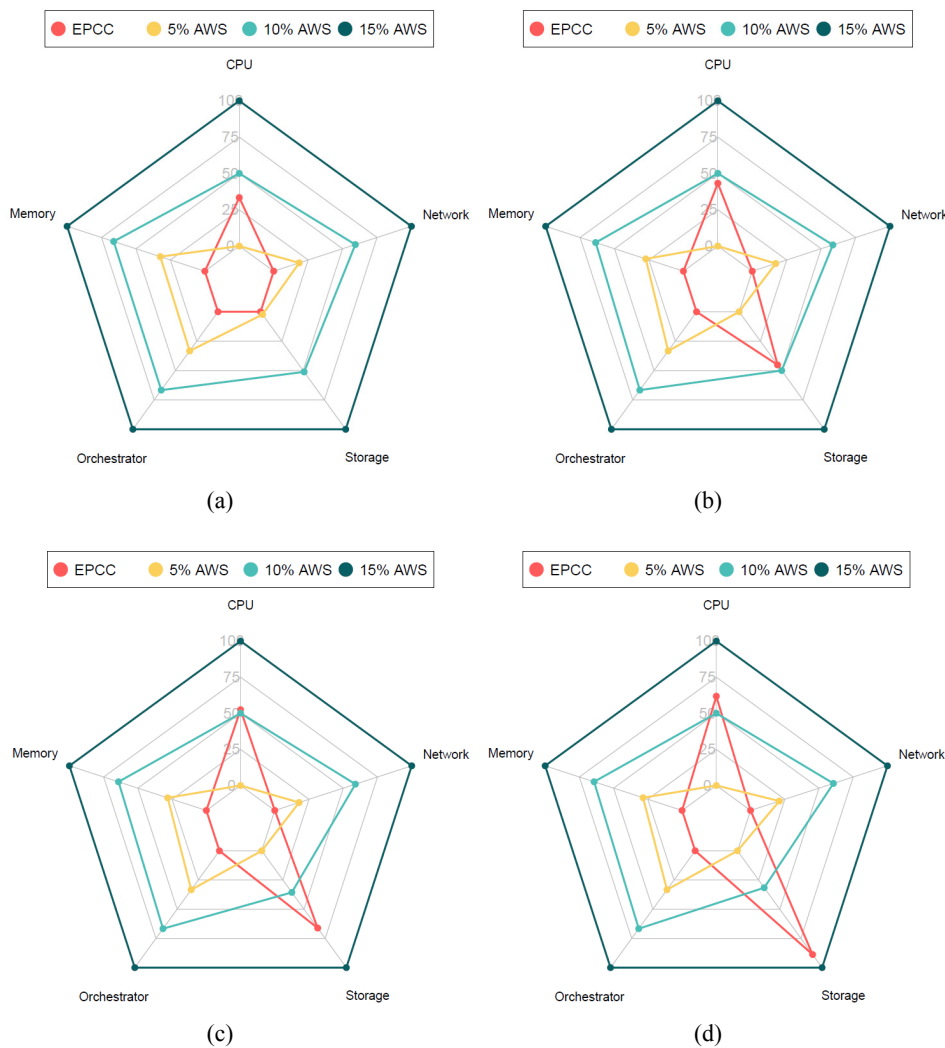
5.2.1 Estimated price of individual cost components

EPCC considers five components' energy consumption to calculate the cost: CPU, memory, storage, network, and container orchestration. In contrast, since AWS has a private cost model, the proportion of the energy cost in the final customer's billing is unknown to us, and we cannot identify these publicly available data. The literature suggests that energy costs reach up to 50% of total DC costs (Comerford, 2015). Operating costs include air conditioning and other energy-intensive equipment. In PSVE (Hinz et al., 2018), the

authors estimate that energy consumption can be between 5 and 15% of IaaS instance price. Our work follows this estimate. Also, we adopt the energy price of US\$ 0.0773 per kWh (This is the value from the electrical distributor in the northern Virginia region in the United States, where the AWS DC is located).

Figure 9 presents estimated prices. As none considers components' usage, the results show the estimated price for four different scenarios (25%, 50%, 75% and 100%). Each radio graph has a proportional axis, in which the maximum value is 100%, and the minimum is 0%. Therefore, the highest prices are on edge and the lowest one in the centre. Each axis represents the estimated price, with energy cost to 5%, 10%, 15% and none, for each cost model component (CPU, memory, network, storage, and orchestrator). We represented the scenarios using different colours using solid lines linked by dashes, forming a polygonal surface. The analysis of combined components is realised by calculating the polygon's surface area. The smaller the surface area, the lower the price.

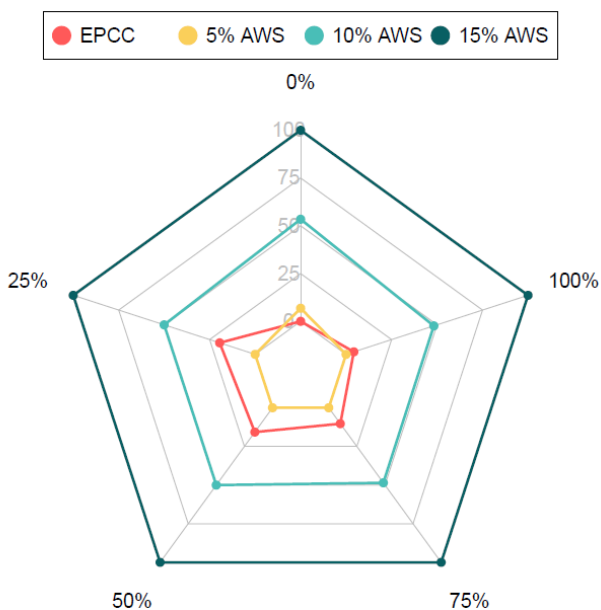
Figure 9 EPCC vs. AWS Fargate: estimated price of individual cost components according to usage (a) Components use by 25% (b) Components use by 50% (c) Components use by 75% and (d) Components use by 100%



In all scenarios (Figure 9(a) to 9(d)), the AWS 15% estimate (dark green) is the worst case, as expected, concerning price (highest). Although AWS 10% (light green) is better than nome when CPU and storage usages are 75% and 100%, the surface area of nome remains smaller than the AWS 10% estimate's surface area. Finally, comparing the AWS 5% (orange estimate) and nome, nome has at least three of five dots in the centre (the best result), but it has an estimate price higher than AWS 5% on two scenarios (75% and 100%). As expected, nome is best when component use is low. More importantly, the AWS pricing model cannot reflect lower utilisation of energy, which corroborates to our main point the importance of aggregating financial benefit to engender energy savings.

The total price considers energy cost to keep an active idle server in addition to the individual component prices. Figure 10 shows a radio graph including five different scenarios based on components usage; e.g., 0% shows idle server and for 100% all components are saturated. For all scenarios, nome leads to total prices lower than AWS 15% and 10%. The AWS 5% estimated has the lowest value for the total price, except when the server is idle. Finally, the price estimates from nome adapt to energy consumption according to components usage. Thus, the tenants are encouraged to improve energy consumption, optimise resource utilisation and apply green IT concepts.

Figure 10 Total container energy price – C_{total}



6 Considerations and future works

Energy is one of the main components of data centres' total cost and physical infrastructure in cloud computing. Since the emergence of cloud computing, cost models have evolved while service offers expanded. Among other characteristics, resource optimisation is one of the key characteristics of the

cloud service model. Cloud tenants are aware of where resources are allocated because it directly impacts their bills. However, energy still is missing in this equation, especially since its impact is not perceived by tenants who are not encouraged to optimise its utilisation.

This work proposed nome, a cost model that rewards energy awareness by accounting for the usage of allocated resources. nome combines the energy component with the allocation component to determine where and how to contract services and use resources. Our analysis of individual resources and their energy consumption reveals that an application with resource-intensive demands can consume 3-4x more energy. Finally, a resource-intensive application can count up to US\$ 0.049 /hour more than an idle one.

Considering the limitations identified during the realisation of this research, it was observed the possibility of considering in a future work a detailed analysis of the Kubernetes energy consumption. Since in the nome proposal, the Kubernetes master was deployed in an isolated structure (i.e., the only service running on the VM). Thus, this specific energy consumption of each computational resource (i.e., network, memory, storage and CPU) was not measured.

Another point to be measured in future work is the apportionment of the energy consumption of different virtualisation services hosted on the same host node. For example, when running a VM and a container in parallel on the same host node, then how should the energy cost of computational resources be attributed to the different environments?

Future work will also comprise expanding the research to evaluate the impact of resource utilisation on energy consumption using other cloud providers and services, such as OpenStack and its services: Nova, Magnum and Ironic.

Acknowledgements

This study was supported by CAPES, FAPESC, UDESC and LabP2D. 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 organisations.

References

- Amazon Web Services (AWS) (2020a) *AWS Fargate Pricing*. Available online at: <https://aws.amazon.com/fargate/pricing/>
- Amazon Web Services (AWS) (2020b) *Data Transfer*. Available online at: https://aws.amazon.com/ec2/pricing/on-demand/#Data_Transfer
- Amazon Web Services (AWS) (2020c) *Amazon EBS volumes*. Available online at: <https://aws.amazon.com/fargate/pricing/>
- Axboe, J. (2017) *I. fio -flexible i/o tester rev. 3.23a*.
- Bawden, T. (2016) 'Global warming: data centres to consume three times as much energy in next decade, experts warn', *The Independent*, Vol. 23.

- Begum, S. and Khan, M.K. (2011) 'Potential of cloud computing architecture', *Proceedings of the IEEE International Conference on Information and Communication Technologies (ICICT)*, IEEE, pp.1–5.
- Bhattacharya, A.A., Culler, D., Kansal, A., Govindan, S. and Sankar, S. (2013) 'The need for speed and stability in data center power capping', *Sustainable Computing: Informatics and Systems*, Vol. 3, No. 3, pp.183–193.
- Bindu, G.B.H., Ramani, K. and Bindu, C.S. (2018) 'Energy aware multi objective genetic algorithm for task scheduling in cloud computing', *International Journal of Internet Protocol Technology*, Vol. 11, No. 4, pp.242–249.
- Bittencourt, L.F., Goldman, A., Madeira, E.R.M., Da Fonseca, N.L.S. and Sakellariou, R. (2018) 'Scheduling in distributed systems: a cloud computing perspective', *Computer Science Review*, Vol. 30, pp.31–54.
- Brondolin, R., Sardelli, T. and Santambrogio, M.D. (2018) 'Deep-mon: dynamic and energy efficient power monitoring for container-based infrastructures', *Proceedings of the IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, IEEE, pp.676–684.
- Columbus, L. (2017) 'State of cloud adoption and security', *Forbes*.
- Comerford, T. (2015) 'How data center operators can avoid energy price hikes this winter', *Data Center Knowledge*.
- Da Silva, V.G., Kirikova, M. and Alksnis, G. (2018) 'Containers for virtualization: an overview', *Applied Computer Systems*, Vol. 23, No. 1, pp.21–27.
- Danilak, R. (2017) *Why energy is a Big and Rapidly Growing Problem for Data Centers*, Frobes Technology Council.
- DIAMANTI (2019) *Container Adoption Benchmark Survey*, Technical Report.
- Estrada, Z.J., Stephens, Z., Pham, Kalbarczyk, Z. and Iyer, R.K. (2014) 'A performance evaluation of sequence alignment software in virtualized environments', *Proceedings of the 14th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing (CCGrid)*, IEEE, pp.730–737.
- Fargate AWS (2020) *AWS Fargate serverless compute for containers*. Available online at: <https://aws.amazon.com/en/fargate/>
- Garg, S.K. and Buyya, R. (2012) 'Green cloud computing and environmental sustainability', *Cloud computing and Distributed Systems*, pp.315–340.
- GRID5000 (2020) *GRID5000 large-scale and flexible testbed for experiment-driven research*. Available online at: <https://www.grid5000.fr/>
- Guitart, J. (2017) 'Toward sustainable data centers: a comprehensive energy management strategy', *Computing*, Vol. 99, No. 6, pp.597–615.
- Hammadi, A. and Mhamdi, L. (2014) 'A survey on architectures and energy efficiency in data center networks', *Computer Communications*, Vol. 40, pp.1–21.
- Hinz, M., Koslovski, G.P., Miers, C.C., Pilla, L.L. and Pillon, M.A. (2018) 'A cost model for IAAS clouds based on virtual machine energy consumption', *Journal of Grid Computing*, Vol. 16, No. 3, pp.493–512.
- Hu, J., Deng, J. and Wu, J. (2013) 'A green private cloud architecture with global collaboration', *Telecommunication Systems*, Vol. 52, No. 2, pp.1269–1279.
- Jain, A., Mishra, M., Peddoju, S.K. and Jain, N. (2013) 'Energy efficient computing- green cloud computing', *Proceedings of the International Conference on Energy Efficient Technologies for Sustainability*, pp.978–982.
- Kim, N.Y., Ryu, J.H., Kwon, B.W., Pan, Y. and Park, J.H. (2018) 'CF-cloudorch: container fog node-based cloud orchestration for IoT networks', *The Journal of Supercomputing*, Vol. 74, No. 12, pp.7024–7045.
- King, C.I. (2019) *Stress-ng: A Stress-Testing Swiss Army Knife*, CANONICAL.
- Kominos, C.G., Seyvet, N. and Vandikas, K. (2017) 'Bare-metal, virtual machines and containers in openstack', *Proceedings of the 20th Conference on Innovations in Clouds, Internet and Networks (ICIN)*, pp.36–43.
- Kubernetes (2020) *Kubernetes components*. Available online at: <https://kubernetes.io/docs/concepts/overview/components/>
- Kurpicz, M., Orgerie, A-C., Sobe, A. and Felber, P. (2018) 'Energy-proportional profiling and accounting in heterogeneous virtualized environments', *Sustainable Computing: Informatics and Systems*, Vol. 18, pp.175–185.
- Leitner, P., Cito, J. and Stöckli, E. (2016) 'Modelling and managing deployment costs of microservice-based cloud applications', *Proceedings of the 9th International Conference on Utility and Cloud Computing*, ACM, pp.165–174.
- Li, Z., Tesfatsion, S., Bastani, S., Ali-Eldin, A., Elmroth, E., Kihl, M. and Ranjan, R. (2017) 'A survey on modeling energy consumption of cloud applications: deconstruction, state of the art, and trade-off debates', *IEEE Transactions on Sustainable Computing*, Vol. 2, No. 3, pp.255–274.
- McCalpin, J.D. et al. (1995) 'Memory bandwidth and machine balance in current high performance computers', *IEEE Computer Society Technical Committee on Computer Architecture (TCCA) 'Newsletter'*, Vol. 2, pp.19–25.
- Mell, P. and Grance, T. et al. (2011) *The NIST Definition of Cloud Computing*, NIST Special Publication.
- Mentz, L.L., Loch, W. and Koslovski, G. (2020) 'Comparative experimental analysis of docker container networking drivers', *Proceedings of the 9th IEEE International Conference on Cloud Networking (IEEE CloudNet'20)*, IEEE, USA.
- Moore, S. (2020) *Gartner Forecasts Strong Revenue Growth for Global Container Management Software and Services through 2024*, Sydney, Australia. Available online at: <https://www.gartner.com/en/newsroom/press-releases/2020-06-25-gartner-forecasts-strong-revenue-growth-for-global-co>
- Morabito, R. (2015) 'Power consumption of virtualization technologies: an empirical investigation', *IEEE/ACM 8th International Conference on Utility and Cloud Computing (UCC)*, IEEE, pp.522–527.
- Mortimer, M. (2018) *iperf3 Documentation Release 0.1.10*.
- Pandikumar, S., Kabilan, S.P. and Amalraj, L. (2012) 'Article: green it: a study and analysis of environmental impact of social networks and search engines', *International Journal of Computer Applications*, Vol. 60, No. 6. Doi: 10.5120/9695-4135.
- Sharma, P., Chaufourmier, L., Shenoy, P. and Tay, Y.C. (2016) 'Containers and virtual machines at scale: a comparative study', *Proceedings of the 17th International Middleware Conference*, ACM, pp.1–13.
- Sharma, P., Pegus II, P., Irwin, D., Shenoy, P., Goodhue, J. and Culbert, J. (2017) 'Design and operational analysis of a green data center', *IEEE Internet Computing*, pp.1–1.

- Souppaya, M., Morello, J. and Scarfone, K. (2017) *Application Container Security Guide*, NIST Special Publication.
- Vohra, D. (2018) *Amazon Fargate Quick Start Guide: Learn how to use AWS Fargate to Run Containers with Ease*, Packt Publishing Ltd.
- Wu, C., Buyya, R. and Ramamohanarao, K. (2019) 'Cloud pricing models: taxonomy, survey, and interdisciplinary challenges', *ACM Computing Surveys (CSUR)*, Vol. 52, No. 6, pp.1–36.
- Zakarya, M. and Gillam, L. (2018) 'Managing energy, performance and cost in large scale heterogeneous datacenters using migrations', *Future Generation Computer Systems*, Vol. 93, pp.529–547.
- Zhang, J., Li, K., Guo, D., Qi, H., Yu, H., Jin, Y. and Sangaiah, A.K. (2018) 'Sustainable green data center: guaranteeing flow deadlines in chains of virtual network functions with mrouting', *Sustainable Computing: Informatics and Systems*, Vol. 19, pp.223–232.