

Maximizing Return On Investment for Sustainable Operations through Smart Workload Migration

Georgia Christofidi*
IMDEA Software Institute
Universidad Politécnica de Madrid

Francisco Álvarez Terribas
Jesus Alberto Omaña Iglesias
Nicolas Kourtellis†
Telefónica Research

Thaleia Dimitra Doudali
IMDEA Software Institute

Abstract

The push for sustainable computing has led to an increased adoption of carbon-aware workload migration, where applications are shifted to time periods and regions with lower carbon intensity. Recently proposed solutions rely on the existence of infrastructures across multiple countries, assuming access to greener locations. However, this does not reflect the reality of national companies that only have on-premise IT infrastructure. For such companies to become more sustainable, they must inevitably increase their monetary expenses and expand their native or cloud infrastructure to other countries. Our analysis quantifies this trade-off between sustainability and cost: while migrating workloads to greener regions reduces carbon emissions by an order of magnitude, it also doubles deployment costs. Yet, we uncover a key insight: the application type and characteristics, such as latency and storage footprint, significantly impact carbon emissions and migration costs. This highlights the need for novel solutions that maximize returns on investments.

1 Overview

The growing push for sustainable computing has driven the design of new resource management and workload execution techniques [4, 6, 7, 9, 10], that shift workload execution in times and locations where the carbon intensity of electricity generation [3] allows for reduced overall carbon emissions. More specifically, spatial shifting moves tasks to regions with greener energy sources [10], such as Sweden or Nepal [3], while temporal shifting delays non-urgent workloads[9], such as batch jobs [6], until there is access to greener, renewable energy sources, such as solar energy. The above carbon-aware techniques are increasingly integrated with traditional resource management techniques, prioritizing the reduction of emissions at the potential expense of performance implications or increased operational costs [7].

However, the above carbon optimization strategies **are viable only for multinational** cloud providers and large companies with data centers that span multiple countries. National companies have infrastructure within few countries, typically one, whose carbon intensity and energy sources

may not be very ‘green’ [3]. Such companies can reduce their carbon emissions with spatial shifting by expanding their native or cloud infrastructure to greener countries, which incurs substantial monetary costs. Similarly, temporal shifting is not possible for many real-world, latency-sensitive applications that these companies host. Lastly, compliance with data privacy regulations, such as GDPR, often prohibits offloading data outside of regions, such as the EU.

The above limitations present a trade-off: while carbon-aware workload execution can reduce emissions, it is applicable only for a limited set of applications (e.g., long-lasting batch jobs [7]) and incurs additional cost for expanding infrastructure to greener countries. We then ask; **how can a national company operate with reduced carbon emissions, in return for minimal cost and uninterrupted user service and satisfaction?** In this work, we preliminarily quantify the carbon-cost trade-off that national companies face, when they rent or buy remote resources, in order to make application deployment greener. We consider a realistic scenario of microservice applications with user traffic that is first processed in a local data center before being redirected to a remote, lower-carbon region. Next, we present the details of our experimental characterization.

Motivational experiment. To analyze the trade-off between carbon efficiency and cost, we consider a company that deploys its entire cloud-edge infrastructure in Spain, a country with a moderate carbon intensity of 206 grams of CO₂ equivalent per kilowatt-hour (gCO₂eq/kWh) [3]. Recent works [6, 9] suggest offloading workloads to ‘greener’ locations, such as Sweden (20 gCO₂eq/kWh) [3]. Our study quantifies the differences in carbon footprint, deployment and migration cost (\$) between these two regions.

Performance of microservice applications. We evaluate two microservice-based applications from DeathStar-Bench [2, 5]: a social network (SN) application with 24 microservices and a media streaming (MS) application with 32 microservices. Each application runs on a separate node, with Kubernetes and Prometheus. We deploy a workload that sends 1,000 requests to each application at time steps that follow a Poisson distribution, emulating multiple concurrent users over 10 minutes. Table 1 summarizes the performance results of the 2 different microservice applications. We observe that the MS application has 2.89× higher average latency than the SN, while requiring less memory and similar

*Work done while Georgia Christofidi was at Telefonica Research. Currently she is a PhD student at the IMDEA Software Institute and Universidad Politécnica de Madrid.

† Nicolas Kourtellis is currently at Keysight.

App	AVG Lat	P95/P99 Lat	Storage	Memory
SN	9.49 ms	42.12/89.16 ms	0.85 GB	5.15 GB
MS	26.08 ms	70.65/127.52 ms	0.75 GB	3.71 GB

Table 1. Application performance regardless of location.

App(Location)	Carbon	Local + Move
SocialNet(ES)	72.72	0.0912 + 0.02
SocialNet(SWE)	7.06	0.0864 + 0.02
MediaStream(ES)	166.17	0.0456 + 0.02
MediaStream(SWE)	16.13	0.0432 + 0.02

Table 2. Comparison of carbon (mgCO₂eq) and local (\$/hr) operational cost per location, plus the cost of moving (\$) the app from ES to SWE.

storage. This performance difference arises because composing and uploading a movie review is more computationally demanding than creating a social media post.

Carbon footprint. Next, we capture the difference of running the applications in Spain and Sweden. Table 2 captures the total carbon emissions measured in milligrams of CO₂ equivalent (mgCO₂eq), computed by multiplying total energy consumption with country-specific carbon intensity [8]. We observe that running applications in Sweden reduces emissions by an order of magnitude, aligning with the difference in the carbon intensity of the countries. Finally, we see that the MS application generates more than 2× higher carbon emissions than the SN, due to its longer latency, making its migration to Sweden more impactful for sustainability.

Resource cost. To realize a greener deployment, a company located in Spain needs to expand its infrastructure with additional resources in Sweden, leading to additional resource costs. To quantify those, we use the Amazon EC2 On-Demand Pricing[1] and consider the t3.large instance for the social network application and the t3.medium instance for the media streaming application. Table 2 reports the on-demand hourly rate in the eu-south-2 region for Spain vs. the eu-north-1 region for Sweden. In addition, we calculate the cost of moving the application from Spain to Sweden, by multiplying the storage size with the data transfer cost per GB from Spain to Sweden in AWS[1]. We observe that deploying both applications in Sweden requires a cost similar to deploying them in Spain. Thus, if an application is moved within Europe, then double the budget is needed to acquire similar infrastructure in a different country, because the application needs to run on both locations, to meet performance requirements. Also, there is an additional cost associated with the actual migration, which is proportional to the application’s storage needs. Finally, we observe that the cost of deploying the media streaming application in Sweden is half of that of the social network one.

2 Lessons Learned & Discussion

Our findings highlight the trade-off between carbon emissions and cost for two real-world microservice applications,

media streaming and social network. For both of the applications that we studied, migration to a greener area has **an order of magnitude reduction in carbon emissions, while it doubles the cost of deployment.** In addition, we observe that the media streaming application generates twice the carbon emissions as the social network, and it has a smaller memory and storage footprint, which translates to reduced migration costs and overheads.

The above observations indicate that a greener deployment is not always affordable for a small or medium-sized national company, since it could double its operational cost. However, a **smart choice** of *which* application to migrate to greener locations can lead to a significant reduction in the absolute number of carbon emissions, in return for reduced migration overheads. Therefore, we conclude that extensive and insightful workload characterization is critical for expanding operations in greener regions by smart offloading of selected applications. Our future work will build upon this insight and design a new system solution that deploys online selective workload migration, and focuses on applications that maximize carbon savings while minimizing cost.

References

- [1] [n.d.]. AWS Pricing On-Demand. <https://aws.amazon.com/ec2/pricing/on-demand/>.
- [2] [n.d.]. DeathStar Benchmark. <https://github.com/delimitrou/DeathStarBench>.
- [3] [n.d.]. Electricity Maps. <https://app.electricitymaps.com/>.
- [4] Noman Bashir et al. 2024. The Sunk Carbon Fallacy: Rethinking Carbon Footprint Metrics for Effective Carbon-Aware Scheduling. In *Proceedings of the 2024 ACM Symposium on Cloud Computing*. 542–551.
- [5] Yu Gan et al. 2019. An Open-Source Benchmark Suite for Microservices and Their Hardware-Software Implications for Cloud & Edge Systems. In *Proceedings of the Twenty-Fourth International Conference on Architectural Support for Programming Languages and Operating Systems* (Providence, RI, USA) (ASPLOS ’19). Association for Computing Machinery, New York, NY, USA, 3–18. <https://doi.org/10.1145/3297858.3304013>
- [6] Walid A. Hanafy et al. 2023. CarbonScaler: Leveraging Cloud Workload Elasticity for Optimizing Carbon-Efficiency. *Proc. ACM Meas. Anal. Comput. Syst.* 7, 3, Article 57 (Dec. 2023), 28 pages. <https://doi.org/10.1145/3626788>
- [7] Walid A. Hanafy et al. 2024. Going Green for Less Green: Optimizing the Cost of Reducing Cloud Carbon Emissions. In *Proceedings of the 29th ACM International Conference on Architectural Support for Programming Languages and Operating Systems, Volume 3* (La Jolla, CA, USA) (ASPLOS ’24). Association for Computing Machinery, New York, NY, USA, 479–496. <https://doi.org/10.1145/3620666.3651374>
- [8] Loïc Lannelongue et al. 2021. Green algorithms: quantifying the carbon footprint of computation. *Advanced science* 8, 12 (2021), 2100707.
- [9] Adam Lechowicz et al. 2023. The Online Pause and Resume Problem: Optimal Algorithms and An Application to Carbon-Aware Load Shifting. *Proc. ACM Meas. Anal. Comput. Syst.* 7, 3, Article 45 (Dec. 2023), 32 pages. <https://doi.org/10.1145/3626776>
- [10] Thanathorn Sukprasert et al. 2024. On the Limitations of Carbon-Aware Temporal and Spatial Workload Shifting in the Cloud. In *Proceedings of the Nineteenth European Conference on Computer Systems* (Athens, Greece) (EuroSys ’24). Association for Computing Machinery, New York, NY, USA, 924–941. <https://doi.org/10.1145/3627703.3650079>