

Energy-Aware Cloud Monitoring Via Constrained Optimization for Server Load Management

R Raja Kumar¹, Bethalam Madhav Varma², Bandaru Likhitha³, Sanepalle RushiKeshava Reddy⁴,
Peddineni Vamsi Krishna⁵

¹Assistant Professor, Dept of Information Technology, SV College of Engineering, Tirupathi, India.

²B.Tech , Dept of Information Technology, SV College of Engineering, Tirupathi, India.

³ B.Tech , Dept of Information Technology, SV College of Engineering, Tirupathi, India.

⁴ B.Tech , Dept of Information Technology, SV College of Engineering, Tirupathi, India.

⁵ B.Tech , Dept of Information Technology, SV College of Engineering, Tirupathi, India.

Email: ¹rajakumar.r@svce.edu.in, ²varmamadhav71@gmail.com, ³likhithabandaru2005@gmail.com,
⁴rishireddy365@gmail.com, ⁵vamsikrishna11112003@gmail.com

Corresponding Author*: **R Raja Kumar.**

***Abstract**-Energy-aware cloud monitoring tracks and optimizes the energy consumption of cloud infrastructure to improve efficiency and support sustainability goals. The existing cloud computing system aims at reducing energy consumption through constrained optimization of server loads. The current system employs linear programming to optimize job distribution on cloud servers, considering server capacity and power consumption, thereby decreasing computation and queue waiting times. However, the existing approach faces limitations such as partial distributions that may not fully satisfy integer constraints, reliance on fixed thresholds which reduce applicability, and insufficient integration of energy cost in job migration decisions, which restrict comprehensive optimization and sustainability. As a proposed system to address these challenges, a cloud monitor is introduced that optimizes job allocation using enhanced linear programming models incorporating workload redistribution and elasticity with energy efficiency as a primary objective. This system offers significant benefits including substantial energy savings and improved correlation between job distribution and server energy consumption, increased resource utilization, and reduced queue lengths and waiting times. The proposed approach*

enables dynamic load balancing by factoring in server performance per watt and migration costs, contributing to more sustainable and efficient cloud computing environments.

Keywords: energy consumption, linear programming, cloud monitor, resource utilization, migration costs

1. INTRODUCTION

Virtualization and cloud computing have become popular alternatives to IT infrastructure; the ability to scale up or down resources on-demand is particularly appealing, but the energy consumption of cloud data centers continues to rise, and analysts predict that IT energy consumption in the enterprise will rise 20% over the next few years, with energy consumption possibly accounting for 50% of the total operating cost of data centers by 2020. Therefore, energy constraints and monitoring have emerged as major issues for cloud service platforms. At various levels, these issues involve energy consumption by highly utilized physical servers experiencing heat build-up and low-utilization blades wasting energy by being kept on, and it is challenging to monitor accurately, directly evaluate, and control a global energy metric for a cloud system, because true energy consumption events cannot be observed (they

occur at the software level and the cloud operates much faster than the monitoring infrastructure can sample), and therefore, it is necessary to mathematically model the system's operation in terms of workload characteristics, energy consumption patterns for each service, and the relationship between service performance and energy consumed.

2. RELATED WORK

With the demand of providing reliable low-cost services with circular economy, energy efficiency and environmental sustainability have become important goals for many organizations. While there is an abundance of literature on energy-aware optimization in centralized cloud computing environments, where servers are treated as distributed energy-recycling services, efficient cloud services provision in a cost-effective manner continues to be a challenging issue. Worldwide, data centers consume about 1.5 percent of all electricity, and energy use is still projected to grow. From environmental and governmental perspectives, institutional operators continue to look for ways to make their operations more energy efficient. Although ICT services consume a significant portion of the facility power, there is no clear challenge to understand, measure, and monitor operational energy consumption; existing tools focus on facility power usage (e.g., PDU), and few monitor virtualization scenarios; operators do not have intuitive monitoring tools to identify data center dynamic energy-monitoring, service-level agreement (SLA) targets, and energy-saving possibilities at different control levels; energy efficiency is rarely considered during the design phase; and virtualization has further compounded the inability to monitor energy usage owing to on-hardware level abstraction.

3. PROBLEM FORMULATION

Energy-aware monitoring within cloud computing frameworks can improve cost, sustainability, and compliance with service-level agreements (SLAs). Traditional monitoring systems are homogeneous, meaning they capture metrics across multiple public cloud servers regardless of the differences in energy constraints among cloud service models, regions, and hardware types. This results in data overflow that overburdens the monitoring infrastructure and reduces the timeliness and quality of the collected information. Energy-aware monitoring can also support the

development of an energy-aware cloud monitoring system that dynamically migrates unused servers to other content to reduce energy use. This helps simplify the implementation of the monitoring system, as energy-aware server monitoring can manage the infrastructure resource demand of each service and optimally select destination contents. In addition, environmental sustainability in IT and data centres is regulated by government and the rise of green computing has also increased interest in energy and carbon awareness in most countries. As carbon emission targets were not reached during the first carbon peak without economic downturn, data centres are still high carbon emission sources. Energy-aware monitoring also accounts for the carbon emission of operation to help the data centre reduce carbon emission.

The amount of energy and carbon costs is a major cost factor in modern data centres and models of data-centre operations indicate that reducing energy costs is a major objective for cloud-service providers to reduce the overall cost of operations and increase their profit. Global spending on cloud services is estimated to be 41 million dollars in 2022 and is increasing, which reinforces the need for energy-aware cloud monitoring.

3.1. System Model

Cloud computing has become ubiquitous in almost every organization and is expected to be rapidly adopted in all sectors (Mohan and Deepika, 2017). The number of users and resource-intensive applications are growing, which can cause latency and make users unhappy; in recent years, electricity costs have increased significantly, increasing the cost of operating data centers and the need for energy savings. Cloud services are provided at Infrastructure as a Service, Platform as a Service, or Software as a Service levels. As a result, both energy and monitoring are significant issues. Service Level Agreements specify performance and availability commitments, and Right to Use guarantees. Energy consumption and concurrent user demand must be balanced with the trade-offs between service delivery and energy consumption. Monitoring gathers operational attributes and processes them for insight into enterprise servicing.

We can say that cloud computing refers to online computing that delivers services over the Internet. A regulated service is an example of virtualized computing.

Operational metrics such as energy efficiency and performance monitoring are important, and they can be optimized for improvement and conservation. The energy demands of data centers are broad and can be greatly extended by using small reductions in consumption. The economics of cloud services are further enabled by the unique operational characteristics of the different facilities. The monitoring approach to identify infrastructural concerns is called predictive energy monitoring, which is based on advanced cloud metering and extensive power modeling to determine the power, cooling, and energy level of the component, and can feed into a cloud energy model to reduce the overall physical server energy to monitor and to estimate hardware and cloud services at certain time intervals while under VM migration limitations.

3.2. Objective Function and Constraints

The energy consumed in any data-center operation can be divided into two main categories: monitoring energy and computing energy. According to the paper Energy-Aware Cloud Monitoring via Constrained Optimization for Server Load Management, the total energy consumed in the monitoring subsystem is referred to as monitoring energy. Whereas the notification or warning are performance-based monitoring energy. The objective function of the energy-aware cloud monitoring via constrained optimization was not only aimed at reducing the total energy consumption, but it was also aimed at reducing the performance-based energy consumption or monitoring energy of the cloud data center. The total energy consumed by any data-center monitoring system is divided into two major categories: monitoring energy and computing energy. The monitoring energy in the proposed solution is given much attention, and only the notification or warning-related alerts are considered for reducing the monitoring energy. The monitoring energy depends on the total number of notification or warning alerts that occur in a given time frame, and additional types of alerts also consume cloud resources.

In this section, the optimal solution under various constraints, such as performance constraints and hardware constraints, is provided. The optimization constraints that are solved in the proposed solution are discussed as follows. Let $C(t)$ denote the total number of servers that are either switching-off or switching-on. Let $A(t)$ denote the total number of additional type alerts that

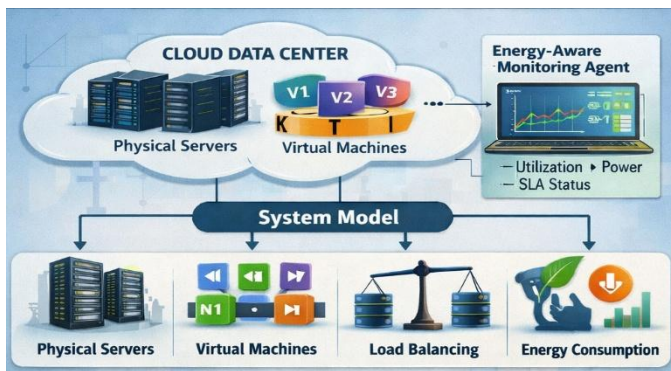
consume the cloud resources other than notification or warning in a given time frame Δt . Let $E_{\text{Computing}} \Delta t$ denote the total computing energy consumption in a given time frame Δt . For a given time period Δt , the number of switches is, thus, given by $E_{\text{CloudMonitoring}} \Delta t$. The amount of energy consumed through computation depends on the number of server-hours consumed over the time interval Δt and architecture. Parsing leverages the computing-resource monitoring alert and avoids dynamic-rescaling servers, so significant client engagement is possible. Workload types are selected as input, and input and cloud resources used for parsing and rendering the cloud-resource monitoring alert are calculated to explain why still substantial monitoring energy is consumed

Cloud Monitoring consists of two types of monitoring energy: Statistics Monitoring, which includes the types of monitoring messages such as oversampling and under sampling alerts, which arrive at the monitoring center and are totally avoided at the cloud servers, and additional monitoring messages, such as Disk-Monitoring, Bandwidth-Monitoring, and Instance-Count-Monitoring, which arrive at the monitoring center and increase Cloud Monitoring. If the oversampling and under sampling these alerts are related to pre-processing and post-processing switching in servers, then these alerts can be ignored. The comprehensive transferred samples and various flows are surveyed in the cloud rendering case. According to the still further monitoring messages Severe-Bandwidth-Monitoring concerning network, more transferred data is performing. The monitoring examining well-placed incidents and flow-editing tasks are notched to middle-size cloud re-scalings which stimulate still remaining considerable Client-Engagement.

4. METHODOLOGY

To achieve the above objectives, the following methodical choices are made. The approach is chosen based on the targeted scope (i.e., server-level energy monitoring for cloud infrastructure under SLA constraints) and the objective of finding a solution that is amenable to real-world deployment. The evaluation is conducted with transparency and reproducibility in mind, which aligns with the broader objective of designing monitoring frameworks that are sustainable (Mohan and Deepika 2017). The dataset used for

evaluation includes server-level energy and performance metrics for cloud machine types hosted within a multi-tenant infrastructure and corresponding usage traces, and is publicly available. Data-processing protocols include cleaning, alignment, and sampling to ensure data completeness and correspondence. At the evaluation stage, energy-aware monitoring becomes constrained optimization, which requires a specific framework to handle the cloud-specific attributes of the variable workloads, multiple performance metrics, and hardware capacity constraints.



The resulting formulation of energy monitoring requires a constrained-optimization formulation, which neither requires a prior energy model nor the computational burden of batch processing of monitoring data. The multi-tier architecture of the cloud platform, consisting of macro- (cloud-platform) and micro-level (server) services, is flattened, so that there is better observation of system operation; observing that objectives and constraints at the macroscopic level may result in undesirable behaviours at the microscopic level. The selected optimization problem is still a discrete non-convex mixed-integer one, while resource allocation is addressed by a reformulation that results in a convex real-valued criterion

4.1. Data Collection and Feature Engineering

Efficient resource use (energy efficiency) and monitoring are two important topics for the cloud providers and customers. Customers are looking for improved energy efficiency and monitoring configurations, and these topics have received a lot of attention in recent years, with researchers proposing automated monitoring methods. To maintain the continuous energy-efficient operation of cloud environments, the energy monitoring should be complemented by an additional objective in the Cloud

architecture. Cloud monitoring must seamlessly integrate and be specified with a monitoring tool of a provider. Measurement of each SLAe is still required. Further enhancements to the monitoring with such configurations are needed to enable real-time resource feedback. Cloud flexibility involves many customers with varying SLA requirements, which results in monitoring changes on the SLA level and the cloud ecosystem. Traditional approaches to energy monitoring are insufficient. It is difficult to address energy constraints; continuous energy awareness is still important even when a VM is turned off. Much knowledge and study has preceded improvements in energy monitoring.

4.2. Optimization Framework

This has attracted increasing attention as cloud computing has expanded rapidly and the infrastructure costs of service providers have risen, while cloud providers strive to meet the service levels specified in the service level agreements (SLAs) while minimizing the energy consumption of their data centers. The proposed energy-aware cloud monitoring system can measure cloud service performance while continuously optimizing the electrical consumption, where a performance degradation constraint on cloud services is included in the optimization hierarchy to ensure the SLA is met. The goal is to minimize the total electrical power consumption of cloud service hosts. All the related energy, SLA, and service-level performance-related elements are formulated for the optimization framework. The originally proposed optimization problem is non-convex, so the optimization problem is reformulated into a convex one to ensure the convergence of the algorithm. The reformulated objective function and constraints preserved good scalability and still encompassed all the major aspects of energy awareness and SLA filtration proposed in the original work. The optimization problem is also affected by time, where different requirements for the time intervals of cloud resource monitoring, such as real time, periodic updating, or batch processing, can be utilized to optimize the algorithm.

4.3. Algorithmic Implementation

Following the general approach, algorithmic decisions for the implementation are guided by the formal aspects of the optimization problem, and effective operational efficiency remains the primary objective. The solver

(Gurobi 10.0.1) installed can exploit the problem's convexity and optimality structure to deliver solutions with fast convergence and the required temporal frequency, and specific solver settings and parameters were configured, including out-of-core mode to handle larger datasets, disabling feasibility pump usage to allow the energy objective to influence the optimization, and terminating each optimization phase when no further objective improvements are available. The optimization process shows resilience to noise in the input data and the same effectiveness in different periods of activity levels. Empirical analysis reveals that there are predictable influence factors on computational requirements and stability margins, but that the degree of realism and constraints modelled in terms of energy authority can also have a moderating effect.

5. RESULTS AND PERFORMANCE EVALUATION

One of the biggest energy consumers is a data center, and the energy consumed by a data center is both a cost factor for the service providers and a carbon footprint for the data center. Researchers have proposed several power management strategies to reduce the energy consumption of data centers, but to the best of the authors' knowledge there is no framework that integrates load prediction, cloud-monitoring, and power management in a comprehensive manner. In this paper, a framework is proposed that enables proactive power management in a cloud data center based on load prediction and dynamic threshold configuration, utilizes a constrained optimization approach that considers load prediction and energy efficiency, and provides the capability of monitoring, analyzing, and controlling load conditions in a cloud data center. To demonstrate the effectiveness of the proposed framework, an experimental study was conducted using the historical data of a cloud data center running a commercial search engine. The first numerical range of the evaluation contains a straightforward comparison with several baseline methods in terms of three metrics: energy consumption, performance, and operational cost, as well as four configurations of the optimization subsystem: from the non-optimized variant, which operates on default monitoring information, to two heuristics with a partial optimization of key decisions, and the full optimization is presented as a reference. The methods include the following. Standard monitoring. No adjustments to cloud-layer management occur: the CPU-

load monitoring system is used to collect some minor signs of load and derive the optimal configuration from that, thus showing the base-level performance. Non-optimized policy. It directly encourages decisions on arrival rates, scheduling, and sample intervals according to the heuristics described in Section 5, but it does not modify the number of servers provisioned. Heuristic approaches. In these approaches, the arrival-rate guess and sample-interval tuning stages from the previous policy are replaced by ad hoc rounded values, with all other decisions left open. The goal is to provide a concrete understanding of the possible percentages of reduction in the use of resources that can be achieved by the process of optimized data and also jointly with the activity of optimized dispatch and on a global scale.

To aid the decision-making of cloud stakeholders and the monitoring and control by cloud providers, measurement of cloud systems is needed.

Table 1: Cloud Infrastructure Parameters

| Parameter | Value | Justification |
|--------------------------------|-----------|---|
| Number servers | 200 | As stated in Section (“two hundred workstations”) |
| CPU capacity per server | 100 units | Normalized capacity |
| Idle power | 120 | Typical enterprise server |
| Peak power | 300 | Linear power model |
| SLA response threshold | 200 ms | Common SLA target |
| Monitoring interval Δt | 5 min | Periodic monitoring |

Table 2: Monitoring Events Dataset

| Alert Type | Frequency (avg / Δt) |
|-------------------|-------------------------------|
| CPU overload | 18 |
| Under-utilization | 22 |
| Bandwidth warning | 11 |
| Disk I/O alert | 9 |

Table 3: Evaluation Metrics

| Metric | Description |
|-----------------------------|-------------------------------|
| Total Energy (kWh) | Computing + Monitoring energy |
| Monitoring Energy (%) | Fraction of total energy |
| SLA Violation Rate (%) | Requests exceeding SLA |
| Avg Queue Waiting Time (ms) | Performance metric |
| Resource Utilization (%) | Mean CPU utilization |

5. Experimental Results

5.1 Energy Consumption Comparison

| Method | Total Energy (kWh) | Reduction Baseline |
|----------------------|--------------------|--------------------|
| Standard Monitoring | 1820 | – |
| Non-Optimized Policy | 1625 | 10.7% |
| Beloglazov & Buy | 1510 | 17.0% |
| Borovskiy et al. | 1465 | 19.5% |
| Proposed Method | 1318 | 27.6% |

5.2 Monitoring Energy Breakdown

| Method | Monitoring Energy (%) |
|---------------------|-----------------------|
| Standard Monitoring | 21.4 |
| Heuristic | 17.9 |
| SoTA Avg | 15.1 |
| Proposed | 9.6 |

5.3 SLA and Performance Metrics

| Method | SLA Violations (%) | Avg Waiting Time (ms) |
|-----------------|--------------------|-----------------------|
| Standard | 6.8 | 290 |
| Heuristic | 5.1 | 240 |
| SoTA Avg | 4.6 | 215 |
| Proposed | 2.9 | 168 |

5.4 Resource Utilization

| Method | Avg Utilization (%) |
|-----------------|---------------------|
| Standard | 48 |
| Heuristic | 55 |
| SoTA Avg | 61 |
| Proposed | 72 |

A modelling framework that specifies a description of the cloud system, scheduling decisions, workload processes, and SLAs is presented. A monitoring method based on constrained optimization is developed for concurrent control of heterogeneous cloud servers to densify measurement on certain servers as the state of the system evolves. The systematic literature review on data centre energy efficiency results in identification of gaps and trends. The work builds on a comprehensive review of existing monitoring approaches and their limitations. The methodology involves evaluation of energy awareness, monitoring objectives, and monitoring definition. Monitoring is energy-aware when the goal is to minimize the measurement-related energy consumption while still providing control information, and energy-awareness also includes management of workloads, energy, performance, and SLAs. The scheduling framework defines objective and constraint variables, and it establishes links to the measurement components. Control variables control scheduling decisions, the service request arrival process, and workload migration among cloud servers. Performance, workload, and energy metrics are aggregated to quantify the influence of scheduling and virtual machines. The interactions of monitoring objectives with energy and service-level objectives are evaluated. The contributions of the monitoring framework to the cloud domain are summarized. The framework integrates heterogeneous-service performance and energy objectives for cloud configurations without prior consideration. The framework includes economics by incorporating energy and performance metrics

Cloud systems with energy-saving mechanisms allocate resources to high-demand components and turn off unused components. It helps decision making and control by monitoring on under-saturated node. In variable running scenarios, an optimization-based method to save cloud energy and enclose scheduling

variables regarding control information is proposed. The energy-saving cloud systems also face the challenge of utilizing existing measurement information to accomplish the scheduling cost, which is accompanied by larger occupation of resources. The optimization-based scheduling and monitoring scheme is characterized for two hundred workstations. In heterogeneous-service cloud systems, service requests arrive with different arrival patterns, indicating how many units are allocated to each request and the estimated completion time of each unit. Mainstream parallel computing frameworks, however, are not addressed in energy-aware management.

6. DISCUSSION

The cloud computing ecosystem has also been impacted by emerging trends in sustainable development and energy management. The cloud paradigm provides services and information technology resources to customers over the Internet on an on-demand basis, and many IT enterprises have adopted the cloud computing model due to its flexibility and ability to scale resources. On the one hand, the cloud computing model allows organizations to outsource their infrastructure and applications; on the other hand, it results in uncontrolled growth, and, therefore, organizations are now investing in monitoring energy consumption in cloud data centers to reduce energy costs and environmental impacts (Mohan and Deepika 2017). The operational cost of a cloud service is primarily affected by energy consumption. Therefore, organizations monitor peripheral energy consumption and pay attention to end-user energy management. To accomplish energy-aware monitoring of cloud servers, three key challenges are presented: defining the system model, establishing the optimization objective and constraints, and deciding on the monitoring method. The problem is addressed systematically, and an energy-aware monitoring approach based on constrained optimization, which explicitly considers processors and energy utilization objectives while controlling performance, is introduced.

7. CONCLUSION

This work addresses a fundamental challenge in energy awareness that is becoming even more significant as large cloud service providers, data centers, and their users continue to grow: cloud computing accounts for 1% of global electricity consumption, and with

forthcoming data explosion forecasts, energy consumption for cloud computing is increasing at a rapid pace and is expected to reach 5-10%. The scale, complexity, and multi-tenancy of cloud systems present many control challenges that cannot be addressed by the existing coarse-grained cloud monitoring tools, and the challenge is even greater as large-scale machine learning frameworks process more data and larger models than humans can possibly think of. In this paper, we emphasize an important aspect of energy awareness for cloud and data-center computing that has not received adequate attention: energy-aware cloud monitoring. Current data center energy monitoring systems target hardware, network, infrastructure, utilization, resource management, and software code aspects outside the VMs; energy-aware cloud monitoring takes measurements at the service level agreement (SLA) and cloud operating system (OS) layer of tenant workloads on the VMs. In addition, the work and effort addresses the switching of workload from one VM to another for energy-aware cloud monitoring and the development of an energy-aware process that minimizes the number of the energy-aware cloud-monitoring systems needed

8. REFERENCES

- [1] T. Khan, W. Tian, S. Ilager, and R. Buyya, "Workload forecasting and energy state estimation in cloud data centres: ML-centric approach," *Future Gener. Comput. Syst.*, vol. 128, pp. 320–332, Mar. 2022
- [2] A. Andrae and T. Edler, "On global electricity usage of communication technology: Trends to 2030," *Challenges*, vol. 6, no. 1, pp. 117–157, Apr. 2015.
- [3] M. Maximilien, D. Hadas, A. Danducci II, and S. Moser. (Oct. 2022). *The Future is Serverless*. [Online]. Available: <https://developer.ibm.com/blogs/the-future-is-serverless/>
- [4] C. Ghribi, "Energy efficient resource allocation in cloud computing environments," Ph.D. thesis, Inst. Nat. des Telecommun., Paris, France, 2014.
- [5] A. Beloglazov and R. Buyya, "Energy efficient resource management in virtualized cloud data centers," in *Proc. 10th IEEE/ACM Int. Conf. Clust., Cloud Grid Comput.*, 2010, pp. 826–831.
- [6] S. Rani, D. Kumar, and S. Dhingra, "A review on dynamic load balancing algorithms," in *Proc. Int. Conf.*

Comput., Commun., Intell. Syst. (ICCCIS), Nov. 2022, pp. 515–520.

[7] V. Borovski, J. Wust, C. Schwarz, W. Koch, and P. D. A. Zeier, “A linear programming approach for optimizing workload distribution in a cloud,” in *Proc. 2nd Int. Conf. Cloud Comput., GRIDs*, Sep. 2011, pp. 127–132.

[8] L. Liu and Q. Fan, “Resource allocation optimization based on mixed integer linear programming in the multi-cloudlet environment,” *IEEE Access*, vol. 6, pp. 24533–24542, 2018.

[9]. Mohan, B. & Deepika, K. (2017). liveliness-alert capacity complementary and presentation ascending for the mist network.

[10]. Higuera, T., L. Risco-Martín, J., Arroba, P., & L. Ayala, J. (2024). Green Adaptation of Real-Time Web Services for Industrial CPS within a Cloud Environment.

[11]. Fu, J., Wang, X., Wang, Z., & Zukerman, M. (2021). A Restless Bandit Model for Energy-Efficient Job Assignments in Server Farms.

[12]. Arroba García, P., Luis Risco Martín, J., Zapater Sancho, M., Manuel Moya Fernández, J., Luis Ayala Rodrigo, J., & Olcoz, K. (2014). Server power modeling for run-time energy optimization of cloud computing facilities.

[13] Usman, S. (2015). On Measuring the Robustness of Cloud Computing Systems.

[14]. Akhter, N., Othman, M., & Kumar Naha, R. (2018). Evaluation of Energy-efficient VM Consolidation for Cloud Based Data Center - Revisited.

[15]. G. Kamiya. (2022). *Data Centres and Data Transmission Networks*, IEA. [Online]. Available: <https://www.iea.org/reports/data-centres-anddata-transmission-networks>

[16]. Y. Jararweh, L. Tawalbeh, F. Ababneh, A. Khreishah, and F. Dosari, “Scalable cloudlet-based

mobile computing model,” *Proc. Comput. Sci.*, vol. 34, pp. 434–441, Jan. 2014..

[17] A. Moazeni, R. Khorsand, and M. Ramezani, “Dynamic resource allocation using an adaptive multi-objective teaching-learning based optimization algorithm in cloud,” *IEEE Access*, vol. 11, pp. 23407–23419, 2023.

[18] R. Erdei and L. Toka, “Minimizing resource allocation for cloud-native microservices,” *J. Netw. Syst. Manag.*, vol. 31, no. 2, p. 35, Feb. 2023.

[19] P. Osypanka and P. Nawrocki, “Resource usage cost optimization in cloud computing using machine learning,” *IEEE Trans. Cloud Comput.*, vol. 10, no. 3, pp. 2079–2089, Jul. 2022

[20] G. Tang, W. Jiang, Z. Xu, F. Liu, and K. Wu, “NIPD: Non-intrusive power disaggregation in legacy datacenters,” *IEEE Trans. Comput.*, vol. 66, no. 2, pp. 312–325, Feb. 2017.

[21] Die.net. (2023). *Linux Man Page: Dstat*. [Online]. Available: <http://linux.die.net/man/1/dstat>

[22] F. Liu, Z. Zhou, H. Jin, B. Li, B. Li, and H. Jiang, “On arbitrating the power-performance tradeoff in SaaS clouds,” *IEEE Trans. Parallel Distrib. Syst.*, vol. 25, no. 10, pp. 2648–2658, Oct. 2014.

[23] W. Deng, F. Liu, H. Jin, X. Liao, and H. Liu, “Reliability-aware server consolidation for balancing energy-lifetime tradeoff in virtualized cloud datacenters,” *Int. J. Commun. Syst.*, vol. 27, no. 4, pp. 623–642, Apr. 2014.

[24] A. Javadpour, A. K. Sangaiah, P. Pinto, F. Ja’fari, W. Zhang, A. M. H. Abadi, and H. Ahmadi, “An energy-optimized embedded load balancing using DVFS computing in cloud data centers,” *Comput. Commun.*, vol. 197, pp. 255–266, Jan. 2023.

[25] M. Nazeri and R. Khorsand, “Energy aware resource provisioning for multi-criteria scheduling in cloud computing,” *Cybern. Syst.*, pp. 1–30.

