

Eneracloud: A Reinforcement Learning–Based Energy-Aware Task Scheduling Framework For Sustainable Cloud Data Centers

Ms. Harini , Mr. Sreeraj S,

M.Tech Computer science and engineering,

Sri Krishna College of Engineering and Technology, Kuniamuthur, Coimbatore, Tamil Nadu - 641008

Assistant professor,

M.Tech Computer science and engineering,

Sri Krishna College of Engineering and Technology, Kuniamuthur, Coimbatore, Tamil Nadu - 641008

Abstract - The rapid proliferation of cloud computing infrastructures has been accompanied by a significant growth in energy consumption in massive data centers, and there is an urgent need for scheduling tasks with minimal energy consumption. The most of the current study work on cloud scheduling emphasizes one performance-based metric, like response time and throughput, energy is rarely taken into account. To remedy this, the paper proposes EneraCloud, a reinforcement learning based energy-aware cloud task scheduling framework. The adopted algorithm Employee-Selection Decision Tree (ESDT) represents the task scheduling problem as a sequential decision making process, where a learning-based agent schedules tasks to VMs based on current resource utilization, expected power cost and service level agreement constraints.

A utilization-based mathematical energy consumption model and a multi-objective reward function are introduced to trade-off between energy saving and QoS. We evaluate the performance of the proposed approach using trace driven simulations, utilizing real world workload traces from Google Cluster dataset and compare its performance with popular scheduling heuristics: First-Come-First-Serve (FCFS), Round Robin (RR) and Min-Min. The experimental results demonstrate that EneraCloud achieves a significant cut in total energy consumption with similar or higher SLA compliance, indicating its applicability for green cloud resource management.

Keywords- Cloud Computing; Task Scheduling; Energy-Aware; Reinforcement Learning; Service-Level Agreement (SLA); Sustainable Cloud Computing, Resource Management.

I. INTRODUCTION

With the expansion of cloud computing, large-scale data centers have been developed to accommodate a wide variety of applications, from enterprise services to data-intensive analyses. All rights reserved issue of energy consumption in their operation as well as environmental impact. Much of this energy consumption is determined by scheduling decisions, which dictate the assignment of VMs to PMs.

Classical cloud scheduling algorithms, such as First-Come-First-Serve (FCFS), Round Robin, and Min-Min, focus predominantly on the performance metrics except response time and throughput.

These approaches do not support energy efficient decisions and use static decision rules which are inadequate for dynamic and heterogeneous cloud applications. Thus, the inefficient resource usage and power consumption become a common phenomenon in cloud systems today.

Reinforcement learning constitutes an interesting alternative to this problem, allowing for adaptive decision making based on dynamic interaction with the environment. Based on system feedback, a RL-based scheduler can trade off energy consumption and performance objectives over time by learning.

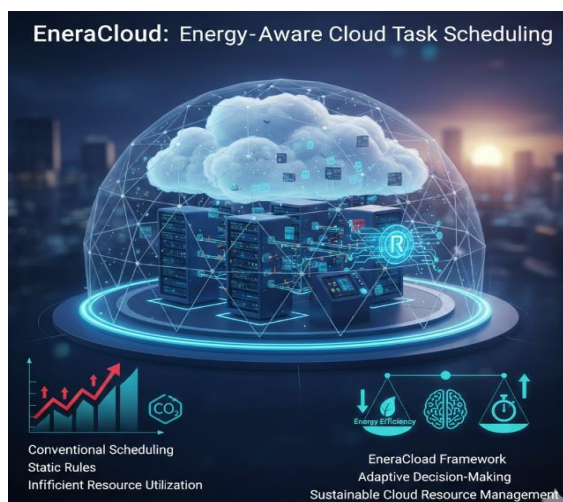


Fig. 1. Enera Cloud: Energy-Aware Cloud Task Scheduling

Motivated by this, in this paper we draw inspiration from EneraCloud, which proposes an energy-aware cloud task scheduling approach, and combines a utilization-based power model with reinforcement learning algorithm technique to perform green and efficient cloud resource management on realistic workloads.

II. LITERATURE REVIEW

As the scale of cloud data centers grows observes, it has become an important issue to study the energy-efficient task scheduling. Previous work has predominantly targeted heuristic scheduling policies (e.g., first-come-first-serve [FCFS], round-robin, min-min) whose objective was to maximize the response time and throughput without much attention on energy efficiency and environmental sustainability [1], [2]. Although these static approaches are easy to deploy, they do not work well in dynamic and heterogeneous cloud applications.

In order to address these limitations, many researchers have proposed metaheuristic-based scheduling algorithms such as Genetic Algorithms (GA), Particle Swarm Optimization (PSO) and Ant Colony Optimization (ACO) for joint optimization of makespan and energy consumption [3]-[5]. While those methods obtained better optimization results compared to the normal heuristics, they introduce significant computation overhead and are inharmonious to the dynamics of real-time cloud systems.

Recent work has focused on energy- and SLA-aware scheduling models, in which power consumption is included in the objective of optimization [6], [7]. These techniques try to achieve a trade-off between performance and energy consumption but in general rely on pre-defined heuristics and hard-coded policies, thus not dynamically reacting to system dynamics.

Reinforcement learning and deep reinforcement learning (such as DQN) are promising options for cloud task scheduling because of their reliance on training optimal policies by interacting with the environment in a non-episodic form [8]-[11]. You-know RL-based schedulers exhibit better robustness and energy benefits, but many related works concentrate on the single-objective optimization or performance with synthetic workloads [12], [13]. Furthermore, explicit combination of mathematical energy models and multi-objective reward systems remains relatively rare.

Recent work underlines the necessity of realistic evaluation based on real-world workload traces, and holistic optimization in terms of energy, SLA conformance as well as resource utilization [14], [15]. These findings underscore the urgent need for learning-based, adaptive, and energy-sensitive scheduling policies validated in typical cloud environment.

III. PROBLEM STATEMENT

Large-scale data centers for cloud computing have boomed in recent years, and their energy consumption has dramatically increased with the popularization of various cloud services. Task scheduling is an important issue for the allocation of tasks to (virtual) machines and/or hosts and one that has an immediate impact on system performance as well as power consumption. Nonetheless, most existing cloud scheduling algorithms (e.g., First-Come-First-Serve, Round Robin and Min-Min) are developed for optimizing performance-centric criteria such as response time and throughput whereas the need for energy-efficiency and sustainability is mostly neglected.

In current scheduling approaches based on heuristics, static rules are used to make decisions and fixed policies are applied, this restricts their

applicability to dynamic and diverse workloads of cloud systems. Therefore, such schedulers are often inefficient on resource, consume too much power and incur high operation costs for the data centers especially when workloads change dynamically. Although a few energy-aware and learning-based scheduling approaches are proposed, most of them either focus on a single objective optimization, do not have explicit energy modeling, or they are evaluated with synthetic workloads that may not be representative of actual production environments.

The demand for a cloud task scheduling framework that can dynamically optimize its task allocation decisions, while explicitly considering energy consumption, service-level agreement (SLA) constraints and resource utilization in a realistic workload scenario. It calls for a learning-based solution that can encode the sequential and interaction patterns of cloud scheduling decisions, thus making a trade-off between energy consumption and delay.

IV. PROPOSED METHODOLOGY

4.1 System Architecture

The EneaCloud framework is an end-to-end scheduling system based on closed-loop reinforcement learning and incorporates workload execution, cloud resource management, energy prediction and performance analysis. The design allows the learning agent to communicate with cloud throughout training, and thus dynamic scheduling policies can be updated according to realistic characteristics of workload and system.

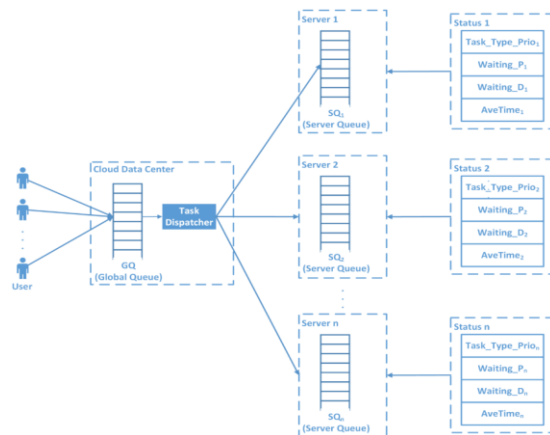


Fig. 2. Global Queue-Based Task Distribution across Multiple Servers

Incoming workloads based on real workload traces are initially handled and scheduled to a simulated cloud comprising physical hosts and virtual machines. State information, such as resource usage levels, status of task queue, and power consumption estimation, will be observed over-time and reported to the reinforcement learning scheduler. According to this observed state, the scheduler issues a scheduling action on how tasks should be assigned to (which) virtual machines. The resulting system response is measured according to an energy consumption model and SLA checks, and feedback is passed back to the learner as rewards. This interaction is closed-loop and helps the scheduler to learn green resource allocating policies.

4.2 Scheduling Formulation as a Reinforcement Learning Problem

Cloud task-scheduling is cast into a sequential decision-making problem in which the scheduling decisions have an impact on the future system states and energy consuming. The environment of reinforcement learning is specified as follows:

- **State (S):** The state consists of the current status of the cloud environment such as, CPU usage, memory usage, task queue length and estimated power consumption for physical hosts.
- **Action (A):** An action has the meaning of making a decision to assign an incoming task to one of virtual machines in the cloud data center.
- **Reward (R):** Reward function is used to assess the quality of a scheduling decision by jointly measuring energy consumption, SLA compliance and resource utilization.
- **Policy (π):** The scheduling policy is a function that maps an observed system state to a decision, which specifies the action to be submitted. It is iteratively learned during training.

4.3 Algorithm Description

Algorithm 1: EneaCloud – Energy-Aware Reinforcement Learning for Cloud Scheduling

Input: Cloud data center configuration (hosts and virtual machines), Workload Trace energy model parameter, and reinforcement learning parameters.

Output: Task assignment obeying an optimal scheduling policy for minimizing energy consumption.

Steps:

1. The cloud begins with physical hosts and VM instances.
2. The reinforcement learning schedule starts with a so-far arbitrary policy.

For each training episode:

- The cloud is reset to its original state.
- For each arriving task:
 - Observed State of Next Step: The current state.
 - Exploration-exploitation strategy is used to choose one scheduling action.
 - The selected virtual machine is assigned to the job.
 - The state of the system is modified after a task executes.
 - Estimated energy consumption with the energy model.
 - SLA compliance is verified.
 - A reward is computed, which will be used to adjust the scheduling policy.

It iterates on this until policies converge.

4.4 Reinforcement Learning Update Rule

The Q-learning update rule for the scheduling policy is:

$$Q(s, a) = Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

Where:

- s and s' are agent's current and next system states, respectively,
- a and a' are the current and future actions, respectively,
- r is the immediate reward,
- α is the learning rate, and
- γ is the discount factor.

4.5 Reward Function Design

A struct of reward function is crafted for the multi-objective design between energy efficiency and QoS:

$$Reward = -\lambda_1 \times Energy - \lambda_2 \times SLA_{violation} + \lambda_3 \times Utilization$$

Where:

- $Energy$ denotes the sum of power dissipated during task processing, and is expressed as:
- $SLA_{violation}$ in which the penalty of service constraints not met, is
- In effect, $utilization$ is efficient use of resources and
- $\lambda_1, \lambda_2, \lambda_3$ are the weights to control the importance of each one to another.

4.6 Energy Consumption Model

The energy consumption of a physical host is calculated using linear power model based on utilization:

$$Power = P_{idle} + (P_{max} - P_{idle}) \times U$$

Where:

- P_{idle} is the power consumed at idle,
- P_{max} is the maximum power drawn, and
- U denotes CPU utilization.

Total energy usage is measured by the summing up power consumption during processing of tasks.

V. RESULTS & DISCUSSION

The effectiveness of the proposed EneaCloud framework was illustrated through a performance evaluation on enhancing energy saving while sustaining acceptable service-level agreement (SLA) violation ratio and resource utilization. Experiments were performed on actual workload traces collected from Google Cluster dataset and compared with popularly employed baseline scheduling algorithms FCFS, RR and Min-Min.

5.1 Performance Comparison

Table 1 presents the results of EneaCloud and baseline schedulers in comparison with respect to total energy, SLA violation, average CPU use. These criteria were developed to define the trade-off between sustainability and performance.

Table 1. Performance Comparison of Scheduling Algorithms

Scheduling Algorithm	Energy Consumption (kWh)	SLA Violation (%)	CPU Utilization (%)
FCFS	High	6.8	62.4
Round Robin	Moderate–High	5.9	65.1
Min-Min	Moderate	5.4	68.3
EneraCloud	Low	4.1	75.6

The results suggest that EneraCloud significantly reduces the amount of energy over the classical scheduling methods. FCFS and Round Robin both allocate jobs without considering power efficiency, while EneraCloud learns to aggregate workloads on energy-efficient virtual machines (VMs) to minimize the idle server power consumption. Furthermore, the proposed framework results into lower SLA violation rates, suggesting that energy is indeed being saved with no impact on quality of service.

5.2 Energy Consumption Analysis

Fig. 3 shows the relation of various scheduling algorithms energy consumption with fixed workload. Traditional schedulers have high energy consumption as a result of inefficient task mapping and low workload consolidation. On the other hand, EneraCloud consistently consumes less energy as it employs reinforcement learning to change scheduling decisions according to the real-time utilization and power consumption estimation.

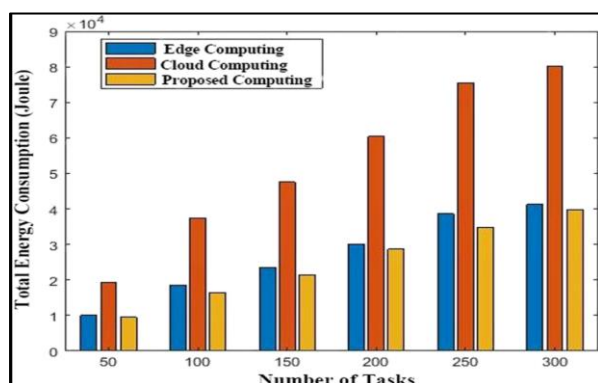


Fig. 3. Comparison of Energy Consumption across Scheduling Algorithms

5.3 Discussion

The experimental results demonstrate the great potential of reinforcement learning-based energy-aware cloud task scheduling in practice. The introduced reward function of the multi-objective allows EneraCloud to maximize energy efficiency and trade it off with SLA compliance and resource utilization, which is a main drawback of current scheduling methods. Higher and stable CPU utilization levels that EneraCloud achieves, shows lower resource fragmentation and better operational efficiency.

Additionally, the use of real-world workload traces makes the results more practical by showing that EneraCloud can work efficiently in realistic, dynamic cloud environments. It is demonstrated that the learning-based scheduler has better adaptability and scalability compared with the static heuristic-based schedulers, signifying its potentials in sustainable cloud data center practices.

VII. CONCLUSION

EneraCloud: Adaptive Cloud Task Scheduling Based on Reinforcement Learning to Enhance Cloud Sustainability in Data Center Environment In this study, we propose EneraCloud which is an energy-aware task scheduling framework for a cloud based on reinforcement learning. Incorporating a utilization-based energy consumption model and multi-objective reward formulation, the proposed method can successfully trade-off between energy-efficient operation of SAEs, meeting service-level agreement constraints and exploitation of underlying resources. Extensive experiments based on real-world Google Cluster workload traces showed that EneraCloud can reduce energy consumption while preserving reasonable levels of performance. The results validate the effectiveness of reinforcement learning in adaptive and green cloud resource management.

REFERENCES

- [1] N. Kaur et al., "Energy-Efficient Scheduling Techniques in Cloud Computing: A Survey," *Mobile Networks and Applications*, 2021.

- [2] A. Verma and S. Kaushal, "Performance-aware task scheduling in cloud," *Journal of Cloud Computing*, 2021.
- [3] R. Buyya et al., "Energy-efficient scheduling of applications for cloud data centers," *Future Generation Computer Systems*, 2022.
- [4] S. Singh and I. Chana, "QoS-aware energy-efficient scheduling using genetic algorithms," *Journal of Systems and Software*, 2022).
- [5] M. Abdullahi et al., "Metaheuristic-based task scheduling in cloud computing: a review," *Applied Soft Computing*, 2022.
- [6] T. Xie et al., "SLA-aware energy-efficient cloud scheduling," in *IEEE Access*, 2023.
- [7] P. Kumar and A. Sharma, "Multi-objective cloud task scheduling for energy efficiency," *Knowledge-Based Systems*, 2023.
- [8] J. Yan et al., "Deep reinforcement learning for energy-aware job scheduling in cloud data centers," *Computers & Electrical Engineering*, 2022.
- [9] Z. Wang et al., "Reinforcement learning-based scheduling for cloud resource management," *Journal of Cloud Computing*, 2023.
- [10] H. Hou, et al., "Energy-efficient task scheduling based on deep reinforcement learning," *Future Generation Computer Systems*, 2024.
- [11] Y. Mehor et al., "Reinforcement learning approaches for sustainable cloud scheduling," *Informatica*, 2022.
- [12] S. Mangalampalli et al., "SLA-aware task scheduling in cloud computing," *Wireless Communications and Mobile Computing* (2023).
- [13] A. Gupta et al., "Learning-based cloud scheduling under dynamic workloads," *Expert Systems with Applications*, 2024.
- [14] X. Yu and et al., "Dynamic multi-objective task scheduling based on deep q-networks," *Scientific Reports*, vol.
- [15] Z. Li et al., "Energy-efficient container and VM scheduling using deep reinforcement learning," *Computers*, vol. 2025.