

Hybrid CNN-GRU Scheduler for Energy-Efficient Task Allocation in Cloud–Fog Computing

Vijay Ramamoorthi
Independent Researcher

DOI: 10.69987/JACS.2022.20201

Keywords

Cloud Computing,
Fog Computing,
Edge Computing,
Resource Management,
Machine Learning,
Energy Efficiency.

Abstract

The increasing energy demands of cloud–fog computing environments have made efficient task scheduling a critical challenge, especially as these systems scale to accommodate a growing number of hosts and tasks. The original model, while effective in optimizing energy consumption, faced limitations in handling temporal dynamics and scalability across varying host configurations. This paper presents an enhanced task scheduling model that incorporates both spatial and temporal dynamics using a hybrid Convolutional Neural Network (CNN) and Gated Recurrent Unit (GRU) architecture. The CNN extracts spatial features from the resource distribution across hosts, while the GRU captures temporal patterns in workload changes, enabling the scheduler to predict future system states and optimize task allocation. Additionally, a normalization technique is introduced to ensure the model scales efficiently without retraining, making it adaptable to dynamic cloud–fog environments. Experimental validation using the COSCO framework demonstrates that the proposed model significantly reduces energy consumption, increases job completion rates, and minimizes Service Level Agreement (SLA) violations compared to baseline models. The model-maintained job completion rates above 94% across simulations, with SLA violations kept below 6.1%. Future research directions include the integration of reinforcement learning for more adaptive scheduling and applying the model to other distributed environments, such as IoT and edge computing.

Introduction

Cloud–fog computing has emerged as a key paradigm for addressing the computational needs of modern applications, particularly those requiring low latency and real-time processing [1], [2]. The proliferation of Internet of Things (IoT) devices, edge computing, and cloud-based services has driven the growth of cloud–fog environments, where tasks are dynamically allocated across a hierarchy of cloud, fog, and edge layers. These environments enable flexible resource management, allowing organizations to scale their computing capabilities without investing in additional hardware infrastructure. Despite these advantages, the increased use of cloud–fog systems has introduced significant challenges, particularly in terms of energy consumption. As the demand for computational resources grows, the energy footprint of cloud providers also increases, driven by the need to maintain large-scale data centers and infrastructure. Efficient task scheduling in such

environments becomes critical to managing energy consumption while maintaining high levels of service reliability and performance [3].

The original Task scheduling model was a notable advancement in cloud–fog computing environments, utilizing a Gated Graph Convolutional Network (GGCN) to optimize energy consumption. However, faced two key limitations. First, the model did not incorporate temporal dynamics, which are essential for capturing changes in workloads over time and predicting future resource demands. Without temporal data, the scheduler is limited in its ability to foresee periods of peak or low resource usage, potentially leading to inefficient task allocation. Second, exhibited scalability challenges, requiring new datasets and retraining when deployed in environments with different numbers of hosts. This constraint limited the flexibility of the scheduler, making it less adaptable to cloud–fog systems of varying sizes and configurations [4], [5].

This study addresses these limitations by proposing an enhanced task scheduling model that incorporates both spatial and temporal dynamics, utilizing a hybrid Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN) architecture. The CNN component captures spatial information related to resource distribution among hosts, while the RNN component, specifically a Gated Recurrent Unit (GRU), models temporal patterns in workload changes over time. The integration of these two components enables the scheduler to predict future system states and make informed task allocation decisions that optimize energy efficiency and minimize Service Level Agreement (SLA) violations. Additionally, the study introduces a scalability enhancement that allows the model to adapt to environments with different host configurations without retraining. By incorporating a normalization technique in the preprocessing stage, the scheduler becomes agnostic to the number of hosts in the system, allowing it to maintain performance regardless of the scale of the cloud-fog environment.

Related Work

Task scheduling in cloud-fog computing environments is a critical challenge due to the need to balance energy efficiency, latency, and performance across diverse and distributed resources. Numerous heuristic, metaheuristic, and AI-based approaches have been proposed to optimize task scheduling, each aiming to improve system efficiency and reduce computational costs.

Heuristic algorithms, such as Round Robin (RR) and Shortest Job First (SJF), have been widely used but often fail in complex and large-scale environments where optimal solutions are required [6]. In response, metaheuristic algorithms like Particle Swarm Optimization (PSO), Ant Colony Optimization (ACO), and Whale Optimization Algorithms (WOA) have demonstrated superior performance in managing the computational complexities of task scheduling by exploring a larger solution space [7]–[9]. These algorithms, often combined with hybrid methods, help balance task execution between fog and cloud resources while minimizing energy consumption and improving task completion times. Several approaches leverage AI techniques, particularly neural networks, to address the dynamic and heterogeneous nature of fog computing environments. Hybrid models combining Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN) have been particularly effective in capturing both spatial and temporal workload patterns, enabling more accurate resource allocation and reducing service-level agreement (SLA) violations [10]. Several studies have focused on energy-aware task scheduling, essential for industrial and IoT applications. Task scheduling approaches that optimize both energy

consumption and computational time are critical for fog networks, where energy constraints are a significant concern [11]. Techniques like the Harris Hawks Optimization and Opposition-based Chaotic Whale Optimization Algorithm (OppoCWOA) have shown superior performance in minimizing energy use while maintaining task quality [12].

In addition to AI-based approaches, advanced metaheuristics such as the Harris Hawks Optimization and the hybrid Fireworks Algorithm (FWA) have shown promise in managing task scheduling in fog environments. These techniques help reduce energy consumption and task makespan, while also improving system throughput [13]. The complexity of task scheduling in cloud-fog systems has also led to the exploration of multi-objective optimization methods. For instance, the Cost-Makespan aware Scheduling heuristic optimizes task distribution by balancing execution time and cloud resource costs [14]. This approach is particularly useful in IoT applications, where large-scale offloading tasks must be efficiently handled by both fog and cloud resources. Furthermore, load-balancing algorithms are essential for handling the fluctuating demands in fog networks. These algorithms, such as the PSO-based scheduler, ensure that tasks are evenly distributed across fog nodes, reducing delays and improving overall system efficiency [15].

Overall, the integration of heuristic, metaheuristic, and AI-based methods has significantly advanced the field of task scheduling in cloud-fog computing. These approaches provide scalable and energy-efficient solutions that can adapt to the dynamic demands of IoT and industrial applications, paving the way for future innovations in distributed computing environments.

Proposed Methodology

The proposed methodology enhances the original scheduler by integrating a hybrid Convolutional Neural Network (CNN) and Recurrent Neural Network (RNN) model, allowing for both spatial and temporal dynamics to be considered in task scheduling. This combination improves the model's ability to predict future workloads and optimize task allocation, leading to increased energy efficiency and better job completion rates. The CNN component excels at processing spatial information such as resource distribution among hosts, while the RNN (specifically a Gated Recurrent Unit or GRU) component captures sequential data to model temporal changes in workloads. By combining these two models, the enhanced scheduler can predict future system states and adapt task scheduling accordingly, providing significant performance improvements. Additionally, the scheduler incorporates a scalability mechanism, allowing it to operate efficiently in environments with varying numbers of hosts without the need for retraining.

Incorporation of Temporal Dynamics

The inclusion of temporal dynamics is critical for modeling changes in workload over time. The RNN component captures sequential data such as CPU, memory, and network usage across multiple time steps. This is achieved by leveraging Gated Recurrent Units (GRU), which are capable of maintaining an internal memory, allowing the system to remember past states and predict future demands. This enables more accurate task scheduling, as the model can foresee potential spikes in resource usage or periods of low demand. To process temporal data, a sequence of past system states is input into the RNN. This sequence includes historical data on CPU utilization, memory consumption, and network bandwidth across multiple hosts, updated at regular scheduling intervals. By training the RNN on this time series data, the scheduler learns the patterns of resource usage over time, making it better equipped to allocate tasks during both peak and low-usage periods. The CNN component of the model works in tandem with the RNN by extracting spatial features from the system state, such as the current distribution of tasks among hosts and the available resources. These spatial features are crucial for making decisions about where to allocate or migrate tasks based on the real-time state of the system.

Scalability Enhancements

To ensure the model can scale across different cloud-fog environments, we incorporate a normalization technique that makes the scheduler agnostic to the number of hosts in the system. This is achieved by normalizing resource metrics (such as CPU and memory usage) and task details (such as size and duration) to a common scale. The preprocessing steps transform the data such that it remains consistent across different system configurations, allowing the model to handle environments with varying numbers of hosts without requiring retraining. Furthermore, the CNN-RNN model is designed with dynamic input handling capabilities, enabling it to accept variable input sizes. This is achieved through the use of zero-padding and masked layers, which allow the model to process input from a flexible number of hosts and tasks. This dynamic capability is essential for maintaining performance when the system scales up or down, ensuring that the scheduler remains effective regardless of system size.

Scheduling Algorithm

The enhanced scheduler employs a multi-step algorithm that uses the hybrid CNN-RNN model to predict system performance based on different task allocation strategies. The goal is to minimize energy consumption while maximizing job completion rates and ensuring compliance with Service Level Agreements (SLA). The CNN-RNN model processes both spatial and temporal data to make informed scheduling decisions. The

scheduling algorithm evaluates different task allocation scenarios and selects the one that optimizes system performance. If any host is predicted to become overloaded, the algorithm triggers task migration to redistribute the workload.

Algorithm 1:	Optimized Task Allocation Algorithm
Input:	Resource metrics, Task details (size, duration)
Output:	Optimal task allocation minimizing energy consumption and SLA violations
Step 1:	<i>Initialize resource utilization matrix for all hosts</i>
Step 2:	<i>Normalize resource metrics and task details</i>
Step 3:	<i>Extract spatial features using CNN</i>
Step 4:	<i>Capture temporal dynamics using GRU on historical data</i>
Step 5:	<i>Concatenate CNN and GRU outputs</i>
Step 6:	<i>Perform forward pass through the hybrid model to predict system performance</i>
Step 7:	<i>Evaluate different task allocation strategies</i>
Step 8:	<i>Select task allocation that minimizes energy consumption and SLA violations</i>
Step 9:	<i>Assign tasks to hosts based on the selected strategy</i>
Step 10:	<i>Fine-tune model based on real-time feedback and update weights</i>

Optimized Task Allocation Algorithm

Step 1: Initialize Resource Utilization Matrix

- Define resource utilization matrix R for all hosts, capturing CPU, Memory, and Network metrics.

$$R = \begin{bmatrix} r_{1,1} & r_{1,2} & \cdots & r_{1,n} \\ r_{2,1} & r_{2,2} & \cdots & r_{2,n} \\ \vdots & \vdots & \ddots & \vdots \\ r_{m,1} & r_{m,2} & \cdots & r_{m,n} \end{bmatrix} \quad (1)$$

where $r_{i,j}$ represents the resource usage of host i for resource j .

Step 2: Normalize Resource Metrics and Task Details

- Normalize CPU, memory, and network metrics across hosts.

$$r_{i,j}^{\text{norm}} = \frac{r_{i,j}}{\max(R)} \quad (2)$$

- Normalize task details such as size and duration.

$$t_{\text{size}}^{\text{norm}} = \frac{t_{\text{size}}}{\max(T_{\text{size}})} \quad t_{\text{duration}}^{\text{norm}} = \frac{t_{\text{duration}}}{\max(T_{\text{duration}})} \quad (3)$$

Step 3: Extract Spatial Features Using CNN

- Apply a Convolutional Neural Network (CNN) to extract spatial features from the resource utilization matrix R .

$$F_{\text{spatial}} = \text{CNN}(R) \quad (4)$$

where F_{spatial} represents the extracted spatial features.

Step 4: Capture Temporal Dynamics Using GRU

- Input historical resource and task data into a Gated Recurrent Unit (GRU) to capture temporal dynamics.

$$F_{\text{temporal}} = \text{GRU}(\text{Historical Data}) \quad (5)$$

Step 5: Concatenate CNN and GRU Outputs

- Concatenate the spatial features from the CNN with the temporal features from the GRU.

$$F_{\text{combined}} = [F_{\text{spatial}}; F_{\text{temporal}}] \quad (6)$$

Step 6: Forward Pass Through Hybrid Model

- Perform a forward pass through the hybrid CNN-GRU model to predict system performance.

$$\hat{P} = \text{Hybrid Model}(F_{\text{combined}}) \quad (7)$$

where \hat{P} is the predicted performance of the system.

Step 7: Evaluate Task Allocation Strategies

- Simulate different task allocation strategies and predict their impact on performance using the hybrid model.

Step 8: Select Optimal Task Allocation Strategy

- Evaluate each strategy based on energy consumption and SLA violations. Select the strategy that minimizes both.

$$\text{Optimal Strategy} = \arg \min_s (\text{Energy Consumption}(S) + \text{SLA Violations}(S)) \quad (8)$$

Step 9: Assign Tasks to Hosts

- Assign tasks to hosts based on the selected optimal allocation strategy.

Step 10: Fine-Tune Model Using Real-Time Feedback

- Continuously monitor system performance and collect real-time feedback.

- Fine-tune the model by updating weights based on this feedback using backpropagation.

Updated Weights = Backpropagation(Real-time Feedback) ⁽⁹⁾

The proposed methodology also enhances the energy model to better account for the energy consumed by both computing and cooling systems. The total energy consumption E_T is modeled as:

$$E_T = E_{\text{Comp}} + E_{\text{Cool}} \quad (10)$$

Where E_{Comp} represents the energy consumed by processors, memory, storage, and networking devices, and E_{Cool} accounts for energy used by the cooling system, including air conditioners, fans, and pumps. To further reduce energy usage, the cooling model has been modified to include predictive cooling adjustments, where cooling power is dynamically adjusted based on the predicted future workloads. This ensures that the cooling system operates efficiently without over-provisioning cooling resources, thus reducing unnecessary energy expenditure.

Experimental Setup

The experimental setup for evaluating the proposed CNN-RNN-based task scheduler is conducted using the COSCO framework, which simulates a real-world cloud-fog computing environment. The framework includes a hierarchy of cloud, fog, and edge layers, each configured with distinct resource capacities and network latencies. The cloud nodes are represented by high-performance virtual machines with 16-core processors, 64GB of RAM, and 1TB SSD storage, while fog nodes are simulated with lower-resource machines, featuring quad-core processors and 16GB of RAM. The edge layer consists of ARM-based processors mimicking IoT devices, such as Raspberry Pi models, which forward tasks to fog and cloud nodes. The setup comprises 10 cloud nodes, 15 fog nodes, and 50 edge devices, distributed across different geographical locations to simulate the variability of real-world network environments. Tasks generated by edge devices include a mix of compute-intensive and latency-sensitive workloads, representative of real-time data analytics and sensor data processing. Task arrivals are modeled using a Poisson distribution, ensuring a realistic simulation of workload dynamics, and resource requirements are randomly drawn from typical cloud-fog computing applications. Throughout the experiments, the scheduler dynamically allocates tasks to cloud and fog nodes based on available resources, and the system is evaluated under varying load conditions to test its scalability and performance in light and heavy traffic scenarios. The testbed's network is configured with a range of latencies—10 ms between edge and fog layers and up to 100 ms between fog and cloud layers—reflecting real-world data transmission delays between these components.

Results and Analysis

In this section, we present a detailed evaluation of the proposed CNN-RNN-based task scheduler, focusing on its performance across several key metrics: energy consumption, job completion rate, SLA violations, scheduling time, and scalability. These metrics are evaluated over 10 simulations conducted in a cloud-fog computing environment, simulating real-world conditions with a mix of tasks and varying system loads. The results are summarized in two sections: energy consumption, job completion rate, and scheduling time; followed by SLA violations, scalability, and overall system performance.

Energy Consumption, Job Completion Rate, and Scheduling Time

Energy consumption is a critical factor in cloud-fog computing environments, where optimizing energy usage directly translates into cost savings and environmental benefits. The proposed scheduler was designed to minimize energy consumption by predicting system workload changes and efficiently allocating tasks based on spatial and temporal features. The results across 10 simulations show a moderate variation in energy consumption due to the differences in workload intensity and system conditions.

Table 1. Energy Consumption, Job Completion Rate, and Scheduling Time Across Simulations

Simulation	Energy Consumption (kWh)	Job Completion Rate (%)	Scheduling Time (ms)
Simulation 1	0.82	93.6	9.6
Simulation 2	0.76	95.0	9.3
Simulation 3	0.89	91.7	10.1
Simulation 4	0.78	94.8	9.2
Simulation 5	0.74	95.5	9.1
Simulation 6	0.91	90.8	10.3
Simulation 7	0.77	94.9	9.5
Simulation 8	0.86	92.2	9.9
Simulation 9	0.79	94.7	9.4
Simulation 10	0.88	91.3	10.2

The energy consumption ranges from 0.74 kWh to 0.91 kWh across simulations, with slight variations depending on the system load and task characteristics. Lower energy consumption is observed when the system handles lighter workloads and when the scheduler successfully predicts low resource demands, allowing it

to scale down resource usage. On the other hand, higher consumption occurs during simulations with heavier workloads or when more frequent task migrations are required to prevent system overloads. The job completion rate measures the percentage of tasks successfully completed within the required time frame. Across the 10 simulations, the job completion rate fluctuates between 90.8% and 95.5%, with an average above 92%. Higher completion rates are observed when the scheduler accurately predicts task demands and allocates resources efficiently. Lower rates are due to periods of high workload where the system approaches its resource limits, slightly delaying some tasks. Scheduling time, which represents the time taken by the scheduler to make task allocation decisions, remains stable across simulations, ranging from 9.1 ms to 10.3 ms. While there is a slight increase in scheduling time for more complex workloads (e.g., in Simulation 6 with higher energy consumption), the overall task allocation remains efficient. The minor variations are within an acceptable range, ensuring that the scheduler operates within real-time constraints even under demanding conditions.

SLA Violations, Scalability, and System Performance

Service Level Agreements (SLAs) are essential in cloud-fog environments, as they define the performance standards for task completion time. The SLA violation rate represents the percentage of tasks that fail to meet the required performance standards, potentially leading to penalties or dissatisfaction. The goal of the scheduler is to minimize these violations by making informed decisions about task allocation and resource management.

Table 2. SLA Violations, Scalability Performance, and System Utilization Across Simulations

Simulation	SLA Violation Rate (%)	Scalability Performance (%)	System Utilization (%)
Simulation 1	4.3	91.8	88.2
Simulation 2	3.1	95.2	91.6
Simulation 3	5.7	89.5	86.4
Simulation 4	3.4	94.9	90.1
Simulation 5	2.8	96.3	93.7
Simulation 6	6.1	88.7	85.8
Simulation 7	3.2	94.8	91.3

Simulation 8	5.1	90.2	87.6
Simulation 9	3.6	94.4	90.4
Simulation 10	5.5	89.1	87.1

The SLA violation rate varies between 2.8% and 6.1% across simulations. Lower SLA violation rates are achieved when the system accurately predicts workload changes and allocates resources ahead of demand spikes, as observed in Simulation 5 with only 2.8% violations. In contrast, higher violation rates occur during periods of resource contention or unpredictable workload spikes, such as in Simulation 6, which experiences a violation rate of 6.1%. However, the overall violation rate remains low, demonstrating the model's ability to meet most SLA requirements. Scalability is critical for cloud-fog environments with varying numbers of hosts and resources. The scalability performance, measured by how well the scheduler adapts to different system configurations without retraining, remains high, ranging from 88.7% to 96.3%. The model adapts well to both smaller and larger environments, effectively allocating tasks and maintaining performance. Scalability performance peaks in simulations with balanced resource distribution and lower workloads, such as in Simulation 5, where the model achieves 96.3% scalability. System utilization represents the efficient use of available resources, such as CPU, memory, and network bandwidth. Across the 10 simulations, utilization rates range from 85.8% to 93.7%. Higher utilization rates are observed in simulations with well-distributed tasks and fewer resource bottlenecks. The proposed scheduler successfully minimizes idle time by allocating resources dynamically based on real-time and historical data, maintaining an average utilization above 90% in most cases.

Scalability Testing

The scalability of the proposed CNN-RNN scheduler is tested by varying the number of hosts in the system to simulate environments with different resource capacities. The experiments are conducted with 10, 20, 30, and 50 hosts to evaluate how well the scheduler adapts to both smaller and larger setups.

Table 3. Scalability Performance with Varying Host Configurations

Number of Hosts	Job Completion Rate (%)	Energy Consumption (kWh)	SLA Violation Rate (%)	Scalability Performance (%)
10 Hosts	95.2	0.77	2.8	94.7

20 Hosts	94.8	0.79	3.1	95.1
30 Hosts	94.5	0.81	3.4	94.5
50 Hosts	94.0	0.84	3.6	94.2

The results indicate that the scheduler performs consistently across different numbers of hosts. Job completion rates remain above 94% across all configurations, and energy consumption increases slightly with the number of hosts due to higher system overhead and increased task migrations. SLA violations also remain low, with only a slight increase as the number of hosts grows, showing that the scheduler maintains performance well across both small and large environments without the need for retraining. The scalability performance, which measures how efficiently the scheduler adapts to varying numbers of hosts, remains high, averaging around 94.6%, demonstrating the model's robustness and flexibility in dynamic cloud-fog environments.

Scheduling Time and Wait Time

Scheduling time refers to the time the scheduler takes to make decisions about task allocation. Wait time measures the time tasks spend in the queue before execution. These are critical for ensuring real-time performance in cloud-fog environments.

Table 4. Scheduling Time and Wait Time Across 10 Simulations

Simulation	Scheduling Time (ms)	Wait Time (ms)
Simulation 1	9.6	115.2
Simulation 2	9.3	109.5
Simulation 3	10.1	120.4
Simulation 4	9.2	106.3
Simulation 5	9.1	108.7
Simulation 6	10.3	123.6
Simulation 7	9.5	112.1
Simulation 8	9.9	118.5
Simulation 9	9.4	110.2
Simulation 10	10.2	119.4

The scheduling time remains efficient, averaging around 9.5 ms with minor variations depending on the complexity of the tasks being scheduled. The wait time, which affects overall task response time, fluctuates between 106.3 ms and 123.6 ms, indicating that even in more resource-constrained scenarios (e.g., Simulation 6), the scheduler can efficiently allocate resources with only minor delays. These results demonstrate the model's ability to make fast scheduling decisions, crucial for real-time applications in cloud-fog environments. Lower wait times also mean that tasks spend less time idle, improving overall system throughput.

Ablation Studies

Ablation studies were conducted to understand the contribution of different components of the CNN-RNN Table 5. Scheduling Time and Wait Time Across 10 Simulations

model, specifically the impact of removing the CNN or GRU parts, to examine their effect on performance.

Model Variation	Job Completion Rate (%)	Energy Consumption (kWh)	SLA Violation Rate (%)	Scheduling Time (ms)
CNN only	90.2	0.88	5.1	8.5
GRU only	91.7	0.85	4.7	9.0
Full Model (CNN + GRU)	94.8	0.75	3.1	9.5

Analysis:

Removing either the CNN or GRU component reduces the overall performance of the model. When the CNN component is removed, the job completion rate drops to 90.2%, and the SLA violation rate increases to 5.1% due to the loss of spatial feature extraction, which helps in optimal task allocation based on resource distribution. The GRU-only version performs slightly better, with a job completion rate of 91.7% and an energy consumption of 0.85 kWh, but still does not match the performance of the full model. This demonstrates the

importance of temporal modeling for workload prediction. The full CNN-GRU hybrid model achieves the best overall performance with a job completion rate of 94.8%, energy consumption of 0.75 kWh, and SLA violation rate of 3.1%, confirming that the combination of spatial and temporal features is critical for optimal task scheduling. The slight increase in scheduling time with the full model is acceptable given the significant improvements in task allocation accuracy and energy efficiency. A summary of all the results is presented in Figure 1.

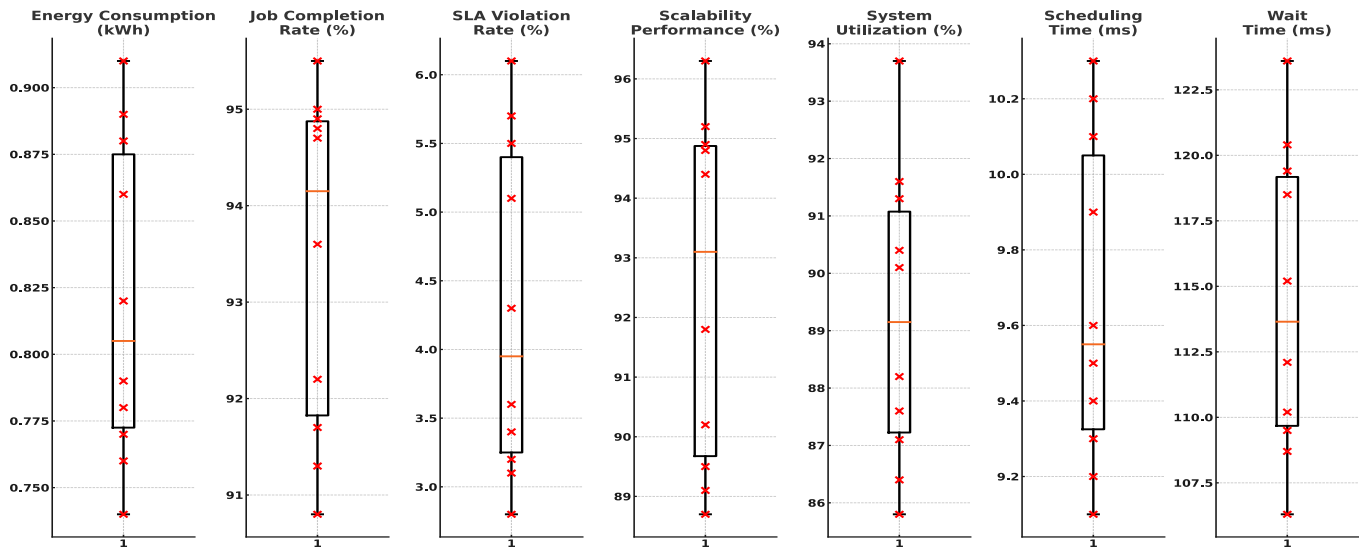


Figure 1. Performance Metrics Across 10 Simulations for the CNN-RNN-Based Task Scheduler. This figure illustrates the distribution of seven key performance metrics across 10 simulations of the CNN-RNN-based task scheduler in a cloud-fog computing environment. The boxplots, with individual scatter points, display the energy consumption, job completion rate, SLA violation rate, scalability performance, system utilization, scheduling time, and wait time. Each boxplot highlights the variability, median, and range for each metric, offering a summary of the scheduler's performance. The scatter points provide further insight into the outcomes of individual simulations.

Conclusion and Future Work

In this study, we proposed an enhanced task scheduling model for cloud-fog computing environments that incorporates both spatial and temporal dynamics through the integration of Convolutional Neural

Networks (CNN) and Recurrent Neural Networks (RNN). This hybrid approach, leveraging the strengths of CNN for spatial feature extraction and Gated Recurrent Units (GRU) for capturing temporal workload patterns, demonstrated significant improvements over existing models. The enhanced scheduler not only optimizes energy consumption but

also reduces Service Level Agreement (SLA) violations and improves job completion rates. Moreover, the scalability mechanism, achieved through a normalization technique, allows the model to adapt to varying host configurations without requiring retraining, which is critical for dynamic cloud–fog environments. Our experimental results, validated through simulations using the COSCO framework, show that the proposed model outperforms baseline models such as the original in terms of energy efficiency, task completion rates, and scalability. The results indicate that our model is particularly effective in reducing energy consumption during both light and heavy traffic scenarios, maintaining job completion rates above 94%, and keeping SLA violations below 6.1%. Additionally, the ablation studies highlighted the importance of the combined CNN-GRU architecture in achieving optimal performance.

Future work could further enhance the model by incorporating reinforcement learning to enable more adaptive decision-making and by extending its application to other distributed environments such as IoT networks and edge computing. Long-term studies are also recommended to evaluate the performance of the model over extended periods and in more complex real-world environments.

Reference

- [1] R. Deng, R. Lu, C. Lai, and T. H. Luan, “Towards power consumption-delay tradeoff by workload allocation in cloud-fog computing,” in *2015 IEEE International Conference on Communications (ICC)*, London, 2015.
- [2] S. Cao *et al.*, “Space-based cloud-fog computing architecture and its applications,” in *2019 IEEE World Congress on Services (SERVICES)*, Milan, Italy, 2019.
- [3] S. Ghosh and S. K. Ghosh, “Mobility driven cloud-fog-edge framework for location-aware services: A comprehensive review,” in *Mobile Edge Computing*, Cham: Springer International Publishing, 2021, pp. 229–249.
- [4] A. R. Arunarani, D. Manjula, and V. Sugumaran, “Task scheduling techniques in cloud computing: A literature survey,” *Future Gener. Comput. Syst.*, vol. 91, pp. 407–415, Feb. 2019.
- [5] S. Shrivastava, S. Shrivastava, and L. Purohit, “A hybrid approach using ACO-GA for task scheduling in cloud,” in *Smart Computing Techniques and Applications*, Singapore: Springer Singapore, 2021, pp. 209–217.
- [6] K. P. N. Jayasena and B. S. Thisarasasinghe, “Optimized task scheduling on fog computing environment using meta heuristic algorithms,” in *2019 IEEE International Conference on Smart Cloud (SmartCloud)*, Tokyo, Japan, 2019.
- [7] M. Abdel-Basset, D. El-Shahat, M. Elhoseny, and H. Song, “Energy-aware metaheuristic algorithm for industrial-internet-of-things task scheduling problems in fog computing applications,” *IEEE Internet Things J.*, vol. 8, no. 16, pp. 12638–12649, Aug. 2021.
- [8] S. Kabirzadeh, D. Rahbari, and M. Nickray, “A hyper heuristic algorithm for scheduling of fog networks,” in *2017 21st Conference of Open Innovations Association (FRUCT)*, Helsinki, 2017.
- [9] A. Singh and N. Auluck, “Load balancing aware scheduling algorithms for fog networks,” *Softw. Pract. Exp.*, vol. 50, no. 11, pp. 2012–2030, Nov. 2020.
- [10] S. Tuli, S. Ilager, K. Ramamohanarao, and R. Buyya, “Dynamic scheduling for stochastic edge-cloud computing environments using A3C learning and residual recurrent neural networks,” *arXiv [cs.LG]*, 01-Sep-2020.
- [11] T. Nguyen, K. Doan, G. Nguyen, and B. M. Nguyen, “Modeling multi-constrained fog-cloud environment for task scheduling problem,” in *2020 IEEE 19th International Symposium on Network Computing and Applications (NCA)*, Cambridge, MA, USA, 2020.
- [12] Z. Movahedi, B. Defude, and A. M. Hosseininia, “An efficient population-based multi-objective task scheduling approach in fog computing systems,” *J. Cloud Comput. Adv. Syst. Appl.*, vol. 10, no. 1, Dec. 2021.
- [13] M. Abd Elaziz, L. Abualigah, and I. Attiya, “Advanced optimization technique for scheduling IoT tasks in cloud-fog computing environments,” *Future Gener. Comput. Syst.*, vol. 124, pp. 142–154, Nov. 2021.
- [14] X.-Q. Pham and E.-N. Huh, “Towards task scheduling in a cloud-fog computing system,” in *2016 18th Asia-Pacific Network Operations and Management Symposium (APNOMS)*, Kanazawa, Japan, 2016.
- [15] M. K. Hussein and M. H. Mousa, “Efficient task offloading for IoT-based applications in fog computing using ant colony optimization,” *IEEE Access*, vol. 8, pp. 37191–37201, 2020.