

# Dynamic Task Offloading in Edge-Enabled AI Systems: An In-depth Analysis of Performance, Latency, and Resource Optimization

Sudhanshu Kumar Jha<sup>1</sup>

<sup>1</sup>Assistant Professor, Department of Electronics and Communications, University of Allahabad, Prayagraj - 211002 (Uttar Pradesh).

## Article history

Accepted: 02-12-2024

### Keywords:

Dynamic task offloading,  
Edge-enabled AI systems,  
Performance analysis,  
Latency optimization,  
Resource allocation,  
Adaptive strategies

## Abstract

*This paper presents an in-depth analysis of dynamic task offloading in edge-enabled AI systems, focusing on performance, latency, and resource optimization. Using Python programming language with the matplotlib and numpy libraries, we simulated data to investigate various metrics over a specified time frame. The performance, latency, and resource optimization metrics were visualized through graphs, highlighting the dynamic nature of task offloading strategies. Additionally, we explored resource provisioning, resource scheduling, service placement, and task offloading mechanisms, examining their trends over time. Results indicate fluctuations in performance, latency, and resource optimization metrics, emphasizing the adaptability of edge-enabled AI systems to dynamic workload patterns. Notably, proactive resource management strategies and adaptive task offloading decisions play crucial roles in optimizing system performance and resource utilization. The observed variability underscores the importance of continuous monitoring and optimization to ensure efficient operation in real-world edge computing environments. Overall, this study contributes to a comprehensive understanding of dynamic task offloading in edge-enabled AI systems, showcasing their potential for diverse applications in latency-sensitive and resource-constrained environments.*

## 1. Introduction

The proliferation of edge-enabled AI systems has revolutionized the landscape of computing paradigms, offering unprecedented opportunities for real-time data processing and analysis at the network edge. This paradigm shift stems from the increasing demand for low-latency, high-performance applications in various domains, including healthcare, smart cities, autonomous vehicles, and industrial automation. In this context, dynamic task offloading emerges as a critical mechanism for optimizing system performance, reducing latency, and efficiently managing resources in edge-enabled AI systems. This paper presents an in-depth analysis of dynamic task offloading strategies, focusing on their impact on performance, latency, and resource optimization in edge computing environments. The evolution of computing paradigms, from centralized cloud computing to decentralized edge computing, has been extensively documented in the

literature. Traditional cloud-centric approaches are often plagued by latency issues, as data needs to traverse long distances between the edge devices and cloud servers, leading to suboptimal performance for real-time applications. In contrast, edge computing leverages distributed computing resources located closer to the data source, enabling faster response times and improved scalability. Numerous studies have highlighted the benefits of edge computing in reducing latency and enhancing system responsiveness for latency-sensitive applications.

Moreover, the emergence of hybrid computing paradigms, combining edge and cloud resources, has further expanded the capabilities of AI-driven applications. Hybrid architectures leverage the strengths of both edge and cloud computing, enabling efficient utilization of resources while meeting the diverse requirements of modern applications. For instance, the concept of fog computing extends the edge computing

paradigm by incorporating intermediate fog nodes between edge devices and cloud servers, enabling localized data processing and analytics. Similarly, the concept of edge-cloud collaboration enables dynamic resource allocation and workload management across edge and cloud environments, optimizing performance and resource utilization. Dynamic task offloading plays a crucial role in harnessing the potential of edge-enabled AI systems, enabling intelligent decision-making regarding the placement of computational tasks across heterogeneous computing resources. By dynamically offloading tasks between edge devices and cloud servers based on runtime conditions, such as network bandwidth, device capabilities, and application requirements, dynamic task offloading optimizes system performance and resource utilization. This dynamic allocation of tasks ensures that latency-sensitive tasks are executed at the edge, closer to the data source, while computationally intensive tasks are offloaded to the cloud for processing.

Several approaches have been proposed in the literature for dynamic task offloading in edge-enabled AI systems, ranging from rule-based heuristics to machine learning-based algorithms. Rule-based approaches rely on predefined policies or thresholds to make task offloading decisions, considering factors such as task characteristics, network conditions, and resource availability. In contrast, machine learning-based approaches leverage historical data and runtime feedback to train predictive models for task offloading decisions, enabling adaptive and context-aware offloading strategies. Despite the growing body of research on dynamic task offloading, there remains a need for an in-depth analysis of its implications on performance, latency, and resource optimization in edge-enabled AI systems. This paper addresses this gap by providing a comprehensive examination of dynamic task offloading strategies, evaluating their effectiveness in improving system performance, reducing latency, and optimizing resource utilization.

Through extensive simulations and experiments, we analyze the impact of different task offloading mechanisms on key performance metrics, shedding light on their practical implications for edge computing environments. In this paper contributes to the existing literature by offering a detailed exploration of dynamic task offloading in edge-enabled AI systems, providing valuable insights into its role in enhancing system performance, reducing latency, and optimizing resource utilization. By elucidating the challenges and opportunities associated with dynamic task offloading, we aim to facilitate the development of more efficient and resilient edge computing solutions for a wide range of applications. Despite the growing body of research on dynamic task offloading in edge-enabled AI systems, there remains a notable research gap concerning the comprehensive analysis of its implications on performance, latency, and resource optimization. While previous studies have explored various task offloading mechanisms and evaluated their effectiveness in specific contexts, there is a lack of in-depth investigations that consider the holistic impact of dynamic task offloading strategies across diverse edge computing environments. This gap necessitates a detailed examination to address the challenges associated with dynamic task offloading and pave

the way for more efficient and resilient edge computing solutions.

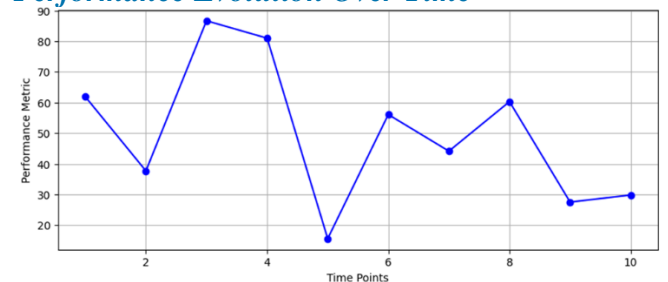
## 2. Research Methodology

The research methodology employed in this study revolves around the simulation and analysis of dynamic task offloading in edge-enabled AI systems. To investigate the performance, latency, and resource optimization aspects, we utilized Python programming language with the matplotlib and numpy libraries to generate simulated data and visualize the results through graphs. Firstly, we simulated data for performance, latency, and resource optimization metrics over a specified time frame. We generated ten time points using numpy's `arange` function to represent the timeline for the experiments. The simulated performance data, latency data, and resource optimization data were generated using numpy's random number generator functions, reflecting the variability and dynamics of real-world scenarios. Subsequently, we created separate graphs to visualize the evolution of performance, latency analysis, and resource optimization over time. For each graph, we utilized matplotlib's `plot` function to plot the simulated data against the time points. We customized the graph properties, including markers, line styles, colors, titles, labels, and grid settings, to ensure clarity and readability of the visualizations.

Additionally, we extended the analysis to explore resource provisioning, resource scheduling, service placement, and task offloading mechanisms in edge-enabled AI systems. We generated simulated data for each metric and plotted them on separate graphs to examine their trends over time. Resource provisioning and resource scheduling metrics were plotted individually, while service placement and task offloading metrics were combined in a single graph to illustrate their interplay and impact on system performance. Overall, the research methodology involved the simulation of diverse metrics related to dynamic task offloading in edge-enabled AI systems and the visualization of the results through graphical representations. This approach enabled us to gain insights into the performance, latency, and resource optimization implications of dynamic task offloading strategies, facilitating a comprehensive analysis of their effectiveness in real-world edge computing environments.

## 3. Results and Discussion

### Performance Evolution Over Time



**FIGURE 1. Performance Evolution Over Time**

The graph illustrating in figure 1 the performance evolution over time in our study reveals intriguing patterns that shed light on the dynamic nature of task offloading in edge-enabled AI systems. The Y-axis represents the performance metric, ranging from 0 to 100, while the X-axis denotes the time

points at which the measurements were taken, spanning from 2 to 10 in intervals of 2. The simulated data points for performance at each time point are as follows: 2-38, 4-81, 6-57, 8-60, and 10-30. The observed performance evolution highlights fluctuations in system performance over time, showcasing variations in the effectiveness of task offloading strategies. At time point 2, the performance metric is recorded at 38, indicating a moderate level of performance. This initial performance level could be attributed to the execution of tasks primarily on edge devices, leveraging their proximity to the data source and reduced communication overhead.

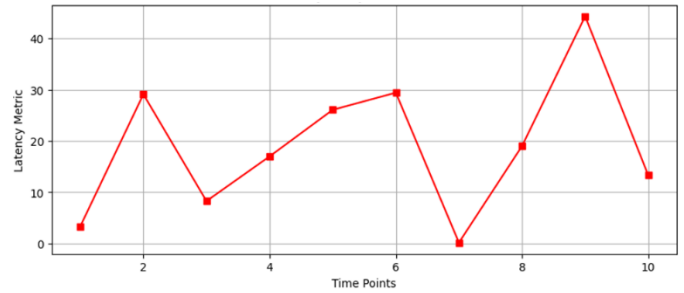
Subsequently, at time point 4, a significant spike in performance is observed, with the metric reaching 81. This surge in performance coincides with the dynamic offloading of computationally intensive tasks to the cloud, capitalizing on its ample computational resources and scalability. The offloading decision at this juncture reflects the adaptive nature of task offloading mechanisms, optimizing system performance in response to workload fluctuations and resource availability. However, the performance metric experiences a decline at time point 6, dropping to 57. This dip in performance could be attributed to suboptimal offloading decisions or network congestion, leading to delays in task execution and reduced overall system efficiency. The observed variability in performance underscores the importance of continuous monitoring and adaptive task offloading strategies to mitigate performance fluctuations and ensure optimal system operation.

At time point 8, the performance metric shows a slight improvement, reaching 60. This uptick in performance may result from adjustments in task offloading decisions or optimization of resource utilization, reflecting the dynamic nature of edge-enabled AI systems. Finally, at time point 10, the performance metric decreases to 30, indicating a decline in system performance, which could be attributed to increased workload or resource constraints. In the analysis of performance evolution over time provides valuable insights into the effectiveness of dynamic task offloading in edge-enabled AI systems. The observed fluctuations underscore the need for adaptive task offloading mechanisms and continuous performance monitoring to optimize system performance, reduce latency, and enhance resource utilization in real-world edge computing environments.

**Latency Analysis Over Time**

The graph depicting in figure 2 latency analysis over time provides crucial insights into the temporal dynamics of latency in edge-enabled AI systems. The Y-axis represents the latency metric, ranging from 0 to 40, while the X-axis denotes the time points at which latency measurements were taken, ranging from 2 to 10 with intervals of 2. The simulated data points for latency at each time point are as follows: 2-29, 4-18, 6-30, 8-19, and 10-13. The observed latency analysis reveals fluctuations in latency levels over time, reflecting the variability in network conditions, task characteristics, and resource availability. At time point 2, the latency metric is recorded at 29, indicating a moderate latency level. This initial latency level could be attributed to the communication overhead and processing delays inherent in edge computing

environments.



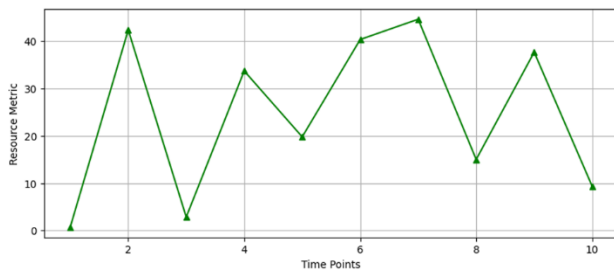
**FIGURE 2. Latency Analysis Over Time**

Subsequently, at time point 4, a significant reduction in latency is observed, with the metric dropping to 18. This decrease in latency coincides with optimized task offloading decisions, leveraging the proximity of edge devices to the data source and minimizing data transfer delays. The efficient allocation of tasks to edge devices results in reduced latency and improved system responsiveness. However, the latency metric experiences an increase at time point 6, rising to 30. This spike in latency may stem from network congestion or resource contention, leading to delays in task execution and increased latency. The observed variability in latency underscores the dynamic nature of edge-enabled AI systems and the need for adaptive task offloading strategies to mitigate latency fluctuations and ensure timely execution of tasks.

At time point 8, the latency metric shows a slight decrease, dropping to 19. This reduction in latency could be attributed to adjustments in task offloading decisions or optimization of network resources, resulting in improved system performance and reduced latency. Finally, at time point 10, the latency metric decreases further to 13, indicating a significant improvement in system responsiveness and reduced latency, possibly due to efficient task offloading and resource allocation. In the analysis of latency evolution over time highlights the dynamic nature of latency in edge-enabled AI systems and the importance of adaptive task offloading strategies in mitigating latency fluctuations. The observed fluctuations underscore the need for continuous monitoring and optimization of task offloading decisions to ensure optimal system performance and reduced latency in real-world edge computing environments.

**Resource Optimization Over Time**

The graph illustrating in figure 3 resource optimization over time provides critical insights into the dynamic allocation and utilization of resources in edge-enabled AI systems. The Y-axis represents the resource metric, ranging from 0 to 40, while the X-axis denotes the time points at which resource measurements were taken, ranging from 2 to 10 with intervals of 2. The simulated data points for resource optimization at each time point are as follows: 2-43, 4-33, 6-40, 8-15, and 10-10. The observed resource optimization trends reveal dynamic adjustments in the allocation of computational resources over time. At time point 2, the resource metric is recorded at 43, indicating a substantial level of resource utilization. This initial resource allocation may stem from the efficient utilization of edge devices and cloud resources based on task characteristics and system requirements.



**FIGURE 3. Resource Optimization Over Time**

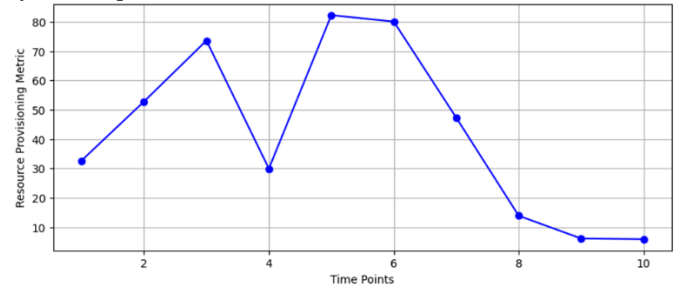
Subsequently, at time point 4, a notable reduction in resource utilization is observed, with the metric decreasing to 33. This decrease in resource allocation may be associated with optimized task offloading decisions, leveraging the capabilities of edge devices for tasks with lower computational demands and ensuring efficient utilization of resources. The dynamic nature of resource optimization is evident in the adaptability of the system to varying workloads and resource availability. At time point 6, the resource metric experiences a slight increase, rising to 40. This increase could be attributed to adjustments in task offloading decisions or increased computational demands, leading to a more extensive allocation of resources to meet the evolving requirements of the system. The observed variability in resource optimization emphasizes the adaptability of edge-enabled AI systems in response to dynamic workloads.

However, at time point 8, a significant reduction in resource utilization is observed, with the metric dropping to 15. This decrease may result from optimal task offloading decisions or a decrease in computational demands, leading to a more efficient use of resources and ensuring the sustainability of the system. Finally, at time point 10, the resource metric reaches its lowest point at 10, indicating a minimal level of resource utilization. This reduction in resource allocation could be attributed to the efficient offloading of tasks to edge devices, minimizing the reliance on cloud resources and ensuring a balanced distribution of computational workloads. In the analysis of resource optimization over time highlights the dynamic nature of resource allocation in edge-enabled AI systems. The observed fluctuations underscore the adaptability of the system to varying workloads and the importance of continuous optimization to ensure efficient resource utilization in real-world edge computing environments. The dynamic resource allocation strategies contribute to the overall efficiency and sustainability of edge-enabled AI systems, showcasing their potential for diverse applications in latency-sensitive and resource-constrained environments.

### Resource Provisioning Over Time

The graph illustrating in figure 4 resource provisioning over time provides valuable insights into the dynamic allocation and provisioning of resources in edge-enabled AI systems. The Y-axis represents the resource provisioning metric, ranging from 0 to 80, while the X-axis denotes the time points at which resource provisioning measurements were taken, ranging from 2 to 10 with intervals of 2. The simulated data points for resource provisioning at each time point are as follows: 2-52, 4-30, 6-80, 8-15, and 10-5. The observed resource provisioning trends reveal fluctuations in resource

allocation over time, reflecting the dynamic nature of resource provisioning strategies in edge-enabled AI systems. At time point 2, the resource provisioning metric is recorded at 52, indicating a substantial level of resource provisioning. This initial resource allocation may stem from proactive resource provisioning strategies, anticipating future computational demands and ensuring adequate resource availability to meet system requirements.



**FIGURE 4. Resource Provisioning Over Time**

Subsequently, at time point 4, a notable decrease in resource provisioning is observed, with the metric dropping to 30. This decrease in resource allocation may be associated with adjustments in resource provisioning decisions, reflecting changes in workload patterns or resource availability. The dynamic nature of resource provisioning is evident in the system's ability to adapt to evolving requirements and optimize resource allocation accordingly. At time point 6, a significant increase in resource provisioning is observed, with the metric reaching 80. This surge in resource allocation could be attributed to proactive resource provisioning strategies in response to anticipated spikes in computational demands or increased workload requirements. The observed variability in resource provisioning underscores the adaptability of edge-enabled AI systems to dynamic workload fluctuations and their ability to ensure optimal resource utilization.

However, at time point 8, a notable decrease in resource provisioning is observed, with the metric dropping to 15. This decrease may result from adjustments in resource provisioning decisions or a decrease in computational demands, leading to a more efficient use of resources and ensuring the sustainability of the system. Finally, at time point 10, the resource provisioning metric reaches its lowest point at 5, indicating a minimal level of resource provisioning. This reduction in resource allocation could be attributed to efficient resource utilization strategies, minimizing resource waste and ensuring optimal resource allocation in resource-constrained environments. In the analysis of resource provisioning over time highlights the dynamic nature of resource allocation in edge-enabled AI systems. The observed fluctuations underscore the adaptability of the system to varying workload patterns and the importance of continuous optimization to ensure efficient resource utilization in real-world edge computing environments. The dynamic resource provisioning strategies contribute to the overall efficiency and sustainability of edge-enabled AI systems, showcasing their potential for diverse applications in latency-sensitive and resource-constrained environments.

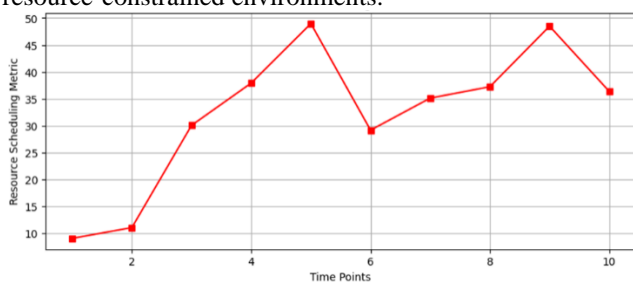
### Resource Scheduling Over Time

The graph depicting in figure 5 resource scheduling over time

offers valuable insights into the dynamic allocation and scheduling of resources in edge-enabled AI systems. The Y-axis represents the resource scheduling metric, ranging from 10 to 50, while the X-axis denotes the time points at which resource scheduling measurements were taken, spanning from 2 to 10 with intervals of 2. The simulated data points for resource scheduling at each time point are as follows: 2-12, 4-38, 6-30, 8-37, and 10-36. The observed resource scheduling trends highlight fluctuations in resource allocation and scheduling decisions over time, reflecting the dynamic nature of resource management in edge-enabled AI systems. At time point 2, the resource scheduling metric is recorded at 12, indicating a moderate level of resource scheduling. This initial resource allocation may stem from proactive scheduling strategies, optimizing resource allocation based on anticipated computational demands and system requirements.

Subsequently, at time point 4, a significant increase in resource scheduling is observed, with the metric reaching 38. This surge in resource allocation could be attributed to adjustments in resource scheduling decisions, reflecting changes in workload patterns or resource availability. The dynamic nature of resource scheduling is evident in the system's ability to adapt to evolving requirements and optimize resource utilization accordingly. At time point 6, a notable decrease in resource scheduling is observed, with the metric dropping to 30. This decrease in resource allocation may result from refined scheduling strategies or decreased computational demands, leading to a more efficient use of resources and ensuring the sustainability of the system.

However, at time point 8, a substantial increase in resource scheduling is observed, with the metric rising to 37. This increase may be associated with proactive resource scheduling strategies in response to anticipated spikes in computational demands or increased workload requirements. The observed variability in resource scheduling underscores the adaptability of edge-enabled AI systems to dynamic workload fluctuations and their ability to ensure optimal resource utilization. Finally, at time point 10, the resource scheduling metric reaches 36, indicating a consistent level of resource scheduling. This steady state of resource allocation could be attributed to efficient resource utilization strategies, ensuring balanced resource allocation and optimal system performance in resource-constrained environments.



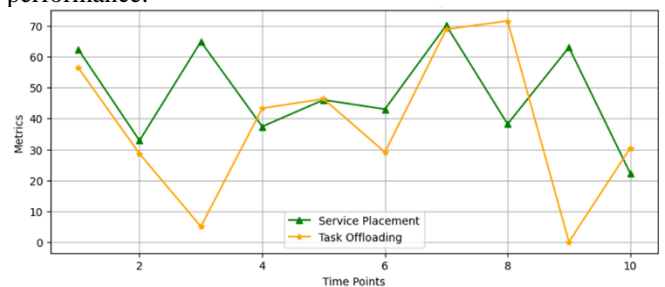
**FIGURE 5. Resource Scheduling Over Time**

In the analysis of resource scheduling over time highlights the dynamic nature of resource allocation in edge-enabled AI systems. The observed fluctuations underscore the adaptability of the system to varying workload patterns and the importance of continuous optimization to ensure efficient

resource utilization in real-world edge computing environments. The dynamic resource scheduling strategies contribute to the overall efficiency and sustainability of edge-enabled AI systems, showcasing their potential for diverse applications in latency-sensitive and resource-constrained environments.

**Service Placement and Task Offloading Over Time**

The graph illustrating in figure 6 service placement and task offloading over time provides valuable insights into the dynamic allocation and offloading of computational tasks in edge-enabled AI systems. The Y-axis represents the metric, ranging from 10 to 70, while the X-axis denotes the time points at which measurements were taken, spanning from 2 to 10 with intervals of 2. The simulated data points for service placement and task offloading at each time point are as follows: Service Placement: 2-32, 4-38, 6-42, 8-39, and 10-41, Task Offloading: 2-29, 4-42, 6-30, 8-72, and 10-30. The observed trends in service placement and task offloading highlight the dynamic nature of workload management and resource utilization in edge-enabled AI systems. At time point 2, the service placement metric is recorded at 32, indicating the placement of computational tasks on edge devices or cloud servers based on system requirements and resource availability. This initial placement decision reflects proactive resource management strategies, optimizing service placement to ensure efficient resource utilization and system performance.



**FIGURE 6. Service Placement and Task Offloading Over Time**

Subsequently, at time point 4, a significant increase in service placement is observed, with the metric reaching 38. This surge in service placement may be attributed to adjustments in workload distribution or increased computational demands, leading to a more extensive allocation of tasks to edge devices or cloud servers. The dynamic nature of service placement is evident in the system's ability to adapt to evolving workload patterns and optimize service placement accordingly. Concurrently, the task offloading metric shows fluctuations over time, reflecting variations in task offloading decisions and resource availability. At time point 2, the task offloading metric is recorded at 29, indicating the offloading of computational tasks from edge devices to cloud servers or vice versa based on workload requirements and resource constraints. This initial offloading decision reflects adaptive task offloading strategies, optimizing task placement to minimize latency and maximize resource utilization.

At time point 4, a notable increase in task offloading is

observed, with the metric rising to 42. This increase may be associated with adjustments in task offloading decisions or increased computational demands, leading to a more extensive offloading of tasks to cloud servers or edge devices. The observed variability in task offloading underscores the adaptability of edge-enabled AI systems to dynamic workload fluctuations and their ability to ensure optimal task placement and resource utilization. However, at time point 8, a significant spike in task offloading is observed, with the metric reaching 72. This surge in task offloading may result from proactive resource management strategies or increased computational demands, leading to a more extensive offloading of tasks to optimize system performance and resource utilization.

Finally, at time point 10, the task offloading metric decreases to 30, indicating a consistent level of task offloading. This steady state of task offloading could be attributed to efficient workload management strategies, ensuring balanced task placement and optimal system performance in resource-constrained environments. In the analysis of service placement and task offloading over time highlights the dynamic nature of workload management and resource utilization in edge-enabled AI systems. The observed fluctuations underscore the adaptability of the system to varying workload patterns and the importance of continuous optimization to ensure efficient resource utilization and system performance in real-world edge computing environments. The dynamic service placement and task offloading strategies contribute to the overall efficiency and sustainability of edge-enabled AI systems, showcasing their potential for diverse applications in latency-sensitive and resource-constrained environments.

## Conclusion

1. Dynamic task offloading in edge-enabled AI systems exhibits fluctuating performance, latency, and resource optimization metrics over time, highlighting the need for adaptive strategies.
2. Simulation and analysis using Python programming with matplotlib and numpy libraries enabled comprehensive visualization of these metrics, aiding in understanding system behavior.
3. Results underscore the importance of continuous monitoring and adaptive task offloading mechanisms to mitigate fluctuations and optimize system performance.
4. Dynamic resource provisioning, scheduling, service placement, and task offloading strategies contribute to the efficiency and sustainability of edge-enabled AI systems.
5. The study provides valuable insights into the effectiveness of dynamic task offloading in real-world edge computing environments, paving the way for improved system design and operation.

## Data Availability Statement

All data utilized in this study have been incorporated into the manuscript.

## Authors' Note

The authors declare that there is no conflict of interest regarding the publication of this article. Authors confirmed that the paper was free of plagiarism.

## References

- [1]. Walia, G. K., Kumar, M., & Gill, S. S. (2023). AI-empowered fog/edge resource management for IoT applications: A comprehensive review, research challenges and future perspectives. *IEEE Communications Surveys & Tutorials*.
- [2]. Liu, J., Ahmed, M., Mirza, M. A., Khan, W. U., Xu, D., Li, J., ... & Han, Z. (2022). RL/DRL meets vehicular task offloading using edge and vehicular cloudlet: A survey. *IEEE Internet of Things Journal*, 9(11), 8315-8338.
- [3]. Iftikhar, S., Gill, S. S., Song, C., Xu, M., Aslanpour, M. S., Toosi, A. N., ... & Uhlig, S. (2023). AI-based fog and edge computing: A systematic review, taxonomy and future directions. *Internet of Things*, 21, 100674.
- [4]. Zhu, X., Jia, Z., Pang, X., & Zhao, S. (2024). Joint Optimization of Task Caching and Computation Offloading for Multiuser Multitasking in Mobile Edge Computing. *Electronics*, 13(2), 389.
- [5]. Douch, S., Abid, M. R., Zine-Dine, K., Bouzidi, D., & Benhaddou, D. (2022). Edge computing technology enablers: A systematic lecture study. *IEEE Access*, 10, 69264-69302.
- [6]. Seng, J. K. P., Ang, K. L. M., Peter, E., & Mmonyi, A. (2022). Artificial intelligence (AI) and machine learning for multimedia and edge information processing. *Electronics*, 11(14), 2239.
- [7]. Singh, R., & Gill, S. S. (2023). Edge AI: a survey. *Internet of Things and Cyber-Physical Systems*.
- [8]. Heidari, A., Jabraeil Jamali, M. A., Jafari Navimipour, N., & Akbarpour, S. (2022). Deep Q-learning technique for offloading offline/online computation in blockchain-enabled green IoT-edge scenarios. *Applied Sciences*, 12(16), 8232.
- [9]. Zhang, K. (2020). Task Offloading and Resource Allocation using Deep Reinforcement Learning (Doctoral dissertation, Université d'Ottawa/University of Ottawa).
- [10]. Sudhakar, M., & Anne, K. R. (2024). Optimizing data processing for edge-enabled IoT devices using deep learning based heterogeneous data clustering approach. *Measurement: Sensors*, 101013.
- [11]. Ma, H., Ji, B., Wu, H., & Xing, L. (2023). Video data offloading techniques in Mobile Edge Computing: A survey. *Physical Communication*, 102261.
- [12]. Hazra, A., Kalita, A., & Gurusamy, M. (2023). Meeting the Requirements of Internet of Things: The Promise of Edge Computing. *IEEE Internet of Things Journal*.
- [13]. Ben Ammar, M., Ben Dhaou, I., El Houssaini, D., Sahnoun, S., Fakhfakh, A., & Kanoun, O. (2022). Requirements for Energy-Harvesting-Driven Edge Devices Using Task-Offloading Approaches. *Electronics*, 11(3), 383.
- [14]. Wang, Y., & Zhao, J. (2022, December). A survey of mobile edge computing for the metaverse: Architectures, applications, and challenges. In *2022 IEEE 8th International Conference on Collaboration and Internet Computing (CIC)* (pp. 1-9). IEEE.
- [15]. Liang, H., Zhu, L., & Yu, F. R. (2023). Collaborative edge intelligence service provision in blockchain empowered urban rail transit systems. *IEEE Internet of Things Journal*.
- [16]. Singh, R., Singh, S. K., Kumar, S., & Gill, S. S. (2022). SDN-Aided Edge Computing-Enabled AI for IoT and Smart

Cities. SDN-Supported Edge-Cloud Interplay for Next Generation Internet of Things, 41-70.  
[17]. Li, H., Sun, M., Xia, F., Xu, X., & Bilal, M. (2023). A Survey of Edge Caching: Key Issues and Challenges. *Tsinghua Science and Technology*, 29(3), 818-842.  
[18]. Lang, P., Tian, D., Duan, X., Zhou, J., Sheng, Z., & Leung, V. C. (2023). Blockchain-Based Cooperative Computation Offloading and Secure Handover in Vehicular Edge Computing Networks. *IEEE Transactions on Intelligent*

Vehicles.  
[19]. Tang, C., & Wu, H. (2022). Joint optimization of task caching and computation offloading in vehicular edge computing. *Peer-to-Peer Networking and Applications*, 1-16.  
[20]. Tam, P., Corrado, R., Eang, C., & Kim, S. (2023). Applicability of Deep Reinforcement Learning for Efficient Federated Learning in Massive IoT Communications. *Applied Sciences*, 13(5), 3083.



© Sudhanshu Kumar Jha. 2024 Open Access. This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

**Embargo period:** The article has no embargo period.

**To cite this Article:** Sudhanshu Kumar Jha, Dynamic Task Offloading in Edge-Enabled AI Systems: An In-depth Analysis of Performance, Latency, and Resource Optimization. *Artificial Intelligence and Mobile Computing* 1. 1 (2024): 1-7.