



AI-DRIVEN RESOURCE ALLOCATION FOR DELAY-SENSITIVE APPLICATIONS IN EDGE COMPUTING

Ramu S/O Kastori Nayak
Research Scholar

Dr. Shashi
Guide

Professor, Chaudhary Charansing University Meerut.

ABSTRACT

The rapid growth of delay-sensitive applications such as autonomous systems, augmented reality, and real-time analytics has intensified the demand for efficient resource management in edge computing environments. Traditional resource allocation techniques often struggle to meet stringent latency requirements due to dynamic workloads and limited edge resources. This paper proposes an AI-driven resource allocation framework that leverages machine learning algorithms to optimize the distribution of computational and network resources at the edge. By incorporating real-time data and predictive analytics, the proposed system dynamically adapts to varying workload patterns, minimizing latency and improving Quality of Service (QoS). Experimental evaluations demonstrate that the AI-based approach significantly outperforms conventional methods in terms of response time, resource utilization, and scalability. The results highlight the potential of artificial intelligence to enable intelligent, adaptive, and efficient resource management for next-generation delay-sensitive applications in edge computing.



KEYWORDS: *Edge Computing, Artificial Intelligence, Resource Allocation, Delay-Sensitive Applications, Latency Optimization, Machine Learning, Real-Time Systems, Quality of Service (QoS), Task Scheduling.*

INTRODUCTION

The proliferation of Internet of Things (IoT) devices, smart applications, and real-time services has significantly transformed modern computing paradigms. Applications such as autonomous vehicles, augmented and virtual reality, smart healthcare, and industrial automation demand ultra-low latency and high reliability. Traditional cloud computing architectures, while powerful, often fail to meet these stringent requirements due to network delays and centralized processing limitations. This has led to the emergence of edge computing, which brings computation and storage closer to end users, thereby reducing latency and improving responsiveness. Despite its advantages, edge computing introduces new challenges, particularly in resource allocation. Edge nodes typically have limited computational power, storage, and energy resources compared to centralized cloud data centers. Moreover, the dynamic and heterogeneous nature of workloads generated by delay-sensitive applications makes efficient resource management a complex task. Static or rule-based allocation strategies are often inadequate in such environments, as they cannot adapt to rapid changes in demand or network conditions. Artificial Intelligence (AI) has emerged as a promising solution to address these challenges. By leveraging machine learning techniques, AI-driven systems can analyze real-time data, predict workload patterns, and make intelligent decisions regarding resource distribution. Techniques such as

reinforcement learning, deep learning, and predictive analytics enable adaptive and proactive resource allocation, ensuring that delay-sensitive tasks are processed within strict latency constraints while maximizing resource utilization.

This work focuses on developing an AI-driven resource allocation framework tailored for delay-sensitive applications in edge computing environments. The proposed approach aims to dynamically allocate computational and network resources based on real-time system conditions and application requirements. By optimizing task scheduling and minimizing response time, the framework seeks to enhance overall system performance and Quality of Service (QoS). The remainder of this paper is organized as follows: Section II reviews related work in edge computing and AI-based resource management. Section III describes the proposed methodology and system architecture. Section IV presents experimental results and performance evaluation. Finally, Section V concludes the paper and outlines future research directions.

AIMS AND OBJECTIVES

Aim

The primary aim of this study is to design and develop an AI-driven resource allocation framework that efficiently manages computational and network resources in edge computing environments to support delay-sensitive applications with minimal latency and improved Quality of Service (QoS).

Objectives

1. To analyze the challenges associated with resource allocation in edge computing, particularly for delay-sensitive and real-time applications.
2. To design an intelligent resource allocation model using Artificial Intelligence techniques such as machine learning or reinforcement learning.
3. To develop a dynamic task scheduling mechanism that adapts to varying workloads and network conditions in real time.
4. To minimize latency and response time by optimizing the allocation of edge resources for critical applications.
5. To improve resource utilization efficiency across distributed edge nodes while avoiding overload and underutilization.

REVIEW OF LITERATURE

Edge computing has emerged as a promising paradigm to support delay-sensitive applications by bringing computational resources closer to end users. A comprehensive review by recent studies highlights that edge computing significantly reduces latency and improves Quality of Service (QoS), making it suitable for applications such as autonomous systems, augmented reality, and industrial IoT. However, efficient resource allocation and task scheduling remain complex due to the distributed and heterogeneous nature of edge environments. Several studies have focused on task scheduling and resource allocation challenges in edge computing. According to a state-of-the-art review, task scheduling in edge environments is a multi-objective optimization problem involving computation offloading, resource selection, and user mobility. These problems are often NP-hard, making it difficult to achieve optimal solutions in real time. Similarly, recent systematic reviews emphasize that dynamic workloads and limited resources further complicate scheduling decisions, necessitating adaptive and intelligent approaches. Early research primarily relied on traditional optimization and heuristic-based approaches. For instance, resource allocation schemes for augmented reality applications have been modeled as matching problems between tasks and computing resources to ensure timely execution and system efficiency. Other studies proposed delay-sensitive algorithms in container-based edge environments, focusing on minimizing end-to-end latency despite fluctuating workloads and complex system states. While these approaches provide useful solutions, they often lack adaptability to highly dynamic environments.

With the advancement of Artificial Intelligence, AI-driven techniques have gained significant attention in recent literature. Machine learning and deep learning methods enable predictive and adaptive resource allocation by analyzing real-time data and workload patterns. In particular, Deep Reinforcement Learning (DRL) has been widely used to model resource allocation as a Markov Decision Process, allowing systems to learn optimal policies under uncertain and dynamic conditions. Studies show that DRL-based approaches outperform traditional methods in terms of energy efficiency and latency reduction. Additionally, optimization frameworks and intelligent provisioning models have been proposed to address latency-sensitive requirements. For example, Lyapunov optimization-based techniques dynamically allocate resources while ensuring system stability and low latency, demonstrating improved performance compared to baseline methods. Recent works also explore application-aware resource allocation and service placement strategies in multi-access edge computing (MEC), considering heterogeneous QoS requirements and data dependencies.

RESEARCH METHODOLOGY

This research adopts a design-oriented and simulation-based methodology to develop an intelligent resource allocation framework for delay-sensitive applications in edge computing environments. The study focuses on integrating artificial intelligence techniques with edge infrastructure to address challenges such as latency minimization, efficient resource utilization, and dynamic workload management. The proposed system is modeled using a multi-layer architecture consisting of end devices, edge nodes, and cloud servers, where delay-sensitive tasks generated by user applications are offloaded to nearby edge nodes for faster processing while the cloud layer supports large-scale computation and backup operations. To emulate realistic scenarios, the study utilizes synthetic and real-world-inspired datasets that represent varying task arrival rates, computational requirements, and network conditions. Workloads are designed to reflect the behavior of latency-critical applications such as real-time analytics and interactive services. The core of the methodology lies in the implementation of an AI-driven resource allocation model, where machine learning techniques, particularly reinforcement learning and deep learning, are employed to enable adaptive decision-making. The system learns from the environment by observing parameters such as system load, network latency, and available resources, and accordingly determines optimal task offloading and resource distribution strategies. A reward mechanism is defined to guide the learning process, aiming to minimize response time while maximizing Quality of Service and resource efficiency.

The framework also incorporates a dynamic task scheduling mechanism that prioritizes delay-sensitive tasks based on urgency and system state. This ensures that critical applications receive timely processing even under fluctuating workloads and constrained resources. The proposed model is implemented using simulation tools and programming environments suitable for edge computing and machine learning, enabling controlled experimentation and performance testing. The evaluation of the system is carried out using key performance metrics including latency, throughput, resource utilization, energy efficiency, and Quality of Service. The results obtained from the AI-driven approach are compared with traditional and heuristic-based resource allocation methods to assess improvements in performance and adaptability. Analytical and graphical methods are used to interpret the results and validate the effectiveness of the proposed framework. Overall, this methodology provides a comprehensive approach to designing and evaluating an intelligent, scalable, and efficient resource allocation system tailored for delay-sensitive applications in edge computing.

STATEMENT OF THE PROBLEM

The rapid growth of delay-sensitive applications such as real-time analytics, autonomous systems, augmented reality, and smart healthcare has created a significant demand for ultra-low latency and high reliability in computing systems. Although edge computing has emerged as a promising solution by bringing computational resources closer to end users, it introduces critical challenges in efficient resource allocation due to limited computational capacity, dynamic workloads, and heterogeneous network conditions. Traditional resource allocation and task scheduling methods

are often static or heuristic-based, making them inadequate for handling the highly dynamic and unpredictable nature of edge environments. These approaches fail to adapt in real time to fluctuations in task demands, network latency, and resource availability, resulting in increased response times, inefficient utilization of resources, and degradation in Quality of Service (QoS). Furthermore, the coexistence of multiple delay-sensitive applications competing for limited edge resources exacerbates the complexity of decision-making.

Another major issue lies in balancing multiple conflicting objectives, such as minimizing latency, maximizing resource utilization, and reducing energy consumption, while ensuring fairness among users and applications. Existing solutions often address these objectives in isolation, leading to suboptimal system performance. Additionally, the lack of intelligent and predictive mechanisms limits the system's ability to proactively manage workloads and prevent congestion or resource bottlenecks. Therefore, there is a critical need for an adaptive and intelligent resource allocation framework that can dynamically respond to changing system conditions and efficiently manage edge resources for delay-sensitive applications. The problem addressed in this research is how to design and implement an AI-driven resource allocation mechanism that minimizes latency, improves QoS, and ensures optimal utilization of limited edge resources in a highly dynamic and distributed computing environment.

DISCUSSION

The findings of this study demonstrate that integrating Artificial Intelligence into resource allocation significantly enhances the performance of edge computing systems, particularly for delay-sensitive applications. The proposed AI-driven framework effectively addresses the limitations of traditional static and heuristic-based approaches by enabling dynamic, real-time decision-making. Through continuous learning and adaptation, the system responds efficiently to fluctuating workloads, varying network conditions, and resource constraints inherent in edge environments. One of the key observations is the substantial reduction in latency achieved by the AI-based model. By intelligently prioritizing delay-sensitive tasks and optimizing task offloading decisions, the framework ensures faster response times compared to conventional methods. This improvement is particularly critical for applications such as autonomous systems and real-time analytics, where even minor delays can significantly impact performance and user experience. Additionally, the use of predictive analytics allows the system to anticipate workload changes, further enhancing responsiveness and minimizing processing delays.

Another important outcome is the improvement in resource utilization. The AI-driven approach dynamically distributes workloads across available edge nodes, preventing both overloading and underutilization. This balanced allocation not only increases system efficiency but also contributes to better energy management. The ability to simultaneously optimize multiple objectives—such as latency, throughput, and energy efficiency—highlights the robustness of the proposed model. The discussion also reveals that reinforcement learning-based strategies are particularly effective in handling the complexity of edge environments. By modeling resource allocation as a sequential decision-making problem, the system learns optimal policies over time without requiring explicit programming for every possible scenario. However, the performance of such models depends on factors such as training time, quality of data, and the design of reward functions. Improper tuning may lead to slower convergence or suboptimal decisions in highly dynamic conditions.

CONCLUSION

This study highlights the critical role of intelligent resource management in supporting delay-sensitive applications within edge computing environments. The proposed AI-driven resource allocation framework demonstrates the ability to effectively overcome the limitations of traditional approaches by enabling dynamic, adaptive, and real-time decision-making. By leveraging machine learning techniques, particularly reinforcement learning and predictive analytics, the system efficiently allocates computational and network resources based on changing workload demands and network conditions. The results indicate significant improvements in key performance metrics, including

reduced latency, enhanced resource utilization, and improved Quality of Service (QoS). The framework's ability to prioritize time-critical tasks and optimize task offloading decisions ensures that delay-sensitive applications are processed with minimal response time. Additionally, the integration of predictive capabilities allows the system to proactively manage workloads, further increasing efficiency and system reliability.

Despite these advancements, challenges such as computational overhead, scalability, and security concerns remain important considerations for real-world implementation. Addressing these issues will be essential for the broader adoption of AI-driven solutions in edge computing. In conclusion, AI-driven resource allocation presents a powerful and promising approach for managing delay-sensitive applications in distributed edge environments. The proposed framework not only enhances system performance but also lays the foundation for future research in developing more robust, scalable, and intelligent edge computing systems.

REFERENCES

1. Shi, W., Cao, J., Zhang, Q., Li, Y., & Xu, L. (2016). Edge computing: Vision and challenges.
2. Mao, Y., You, C., Zhang, J., Huang, K., & Letaief, K. B. (2017). A survey on mobile edge computing: The communication perspective.
3. Mach, P., & Becvar, Z. (2017). Mobile edge computing: A survey on architecture and computation offloading.
4. Zhang, K., Mao, Y., Leng, S., He, Y., & Maharjan, S. (2018). Mobile-edge computing for vehicular networks: A promising network paradigm with predictive off-loading.
5. Sun, X., Ansari, N., & Qiu, X. (2019). Forecasting data-driven resource allocation in mobile edge computing.
6. Chen, X., Jiao, L., Li, W., & Fu, X. (2020). Efficient multi-user computation offloading for mobile-edge cloud computing.
7. Wang, S., Zhang, X., Zhang, Y., & Wang, L. (2020). A survey on reinforcement learning for resource management in edge computing.
8. Li, R., Zhao, Z., Sun, Q., & Zhou, S. (2021). Deep reinforcement learning for resource allocation in edge computing systems.
9. Huang, L., Bi, S., & Zhang, Y. J. (2021). Deep reinforcement learning for delay-sensitive computation offloading in mobile edge computing.
10. Zhang, Y., & Ansari, N. (2022). Edge intelligence and AI-driven resource management: A survey.