



A QOS-CENTRIC APPROACH TO LOAD BALANCING IN SDN-ENABLED IOT ENVIRONMENTS

Ms. Naisargi Nareshchandra Patel¹, Ms. Nirali Kapadia²

PG Scholar¹, Assistant Professor²

Computer Engineering, Gandhinagar Institute of Technology¹
Gandhinagar University, Gandhinagar, India¹

Abstract

The exponential growth of Internet of Things (IoT) devices has introduced complex, heterogeneous, and highly dynamic traffic patterns, posing significant challenges to the efficiency and Quality of Service (QoS) assurance in conventional network infrastructures. Software-Defined Networking (SDN) has emerged as a transformative paradigm by decoupling the control and data planes, enabling centralized, flexible, and programmable management of IoT networks. Nonetheless, the realization of effective QoS-aware load balancing remains a critical concern, particularly in the face of fluctuating traffic loads and device mobility. This study conducts a comprehensive evaluation of contemporary QoS-aware load balancing strategies within SDN-enabled IoT environments, with a specific focus on their adaptability, scalability, and operational efficiency across diverse traffic scenarios. By analysing advanced approaches—ranging from dynamic traffic engineering methods to reinforcement learning-based optimization techniques—this work delineates the comparative advantages and limitations inherent to each strategy. The outcomes of this investigation offer substantive insights for the development of next-generation SDN architectures, fostering intelligent, QoS-driven network orchestration and efficient resource allocation in complex and large-scale IoT deployments.

Key Words: Internet of Things, Load-balancing, Quality of service, SD-IoT, Software- defined networking

1. Introduction

The rapid expansion of Internet of Things (IoT) devices has resulted in a substantial surge in heterogeneous traffic and dynamic service demands, which has imposed considerable strain on conventional network infrastructures. Software-Defined Networking (SDN) has emerged as an innovative paradigm to address these challenges by decoupling the control and data planes, thus enabling centralized and programmable management of network resources [2, 6]. However, one of the most enduring challenges in SDN-based IoT networks is the achievement of efficient and Quality of Service (QoS)-aware load balancing, particularly in scenarios characterized by significant traffic fluctuations and device mobility.

Recent research has highlighted the critical role of intelligent load balancing mechanisms in maintaining essential performance metrics, such as latency, throughput, and packet delivery rates, in SDN-enabled IoT environments. For instance, Tirupathi and Sagar [1] introduced a novel SDN-based algorithm that enhances data handling efficiency and network performance through dynamic traffic management. Similarly, Amri et al. [3] proposed a weighted round-robin (WRR) scheme tailored to IoT networks, demonstrating marked improvements in traffic distribution and resource utilization. These advancements emphasize the importance of QoS-centric evaluation frameworks, which not only balance network load but also adapt in real time to changing network conditions.

The variety of QoS-aware approaches — from slave-controller allocation in ESCALB [5] to reinforcement learning-based techniques for SDN controller balancing [7] — illustrates the growing interest in scalable and adaptable solutions. Rostami and Goli-Bidgoli [2] conducted a thorough survey of QoS-aware load balancing techniques, pointing out that many existing solutions fail to scale effectively in complex, multi-domain IoT networks. Furthermore, emerging methods like QSroute [8], designed to operate under QoS constraints, and intelligent frameworks that integrate edge computing with fog-cloud models [13], suggest a shift toward hybrid, distributed approaches to network management.

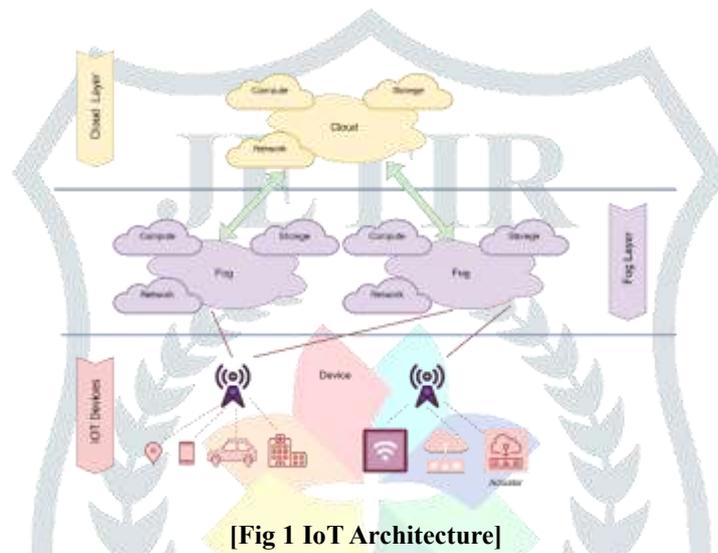
Despite these advancements, a significant gap remains in the comprehensive evaluation of these strategies using unified QoS benchmarks, particularly in real-world, multi-layered IoT deployments. This study seeks to fill that gap by rigorously assessing and comparing the effectiveness of prominent QoS-aware load balancing techniques in SDN-based IoT networks, with a focus on their adaptability, scalability, and operational efficiency under diverse traffic conditions. The insights garnered from this research are expected to inform future developments in network architecture and policy design, fostering more intelligent, QoS-driven network management systems.

2. Background and Related Works

I. IOT Architecture and SDN Architecture

An Internet of Things (IoT) system is generally structured across three interdependent layers: the IoT Devices Layer, the Fog Layer, and the Cloud Layer. At the base, the IoT Devices Layer comprises a wide range of smart devices—including sensors, smartphones, and vehicles—that continuously generate vast volumes of data. These edge devices transmit collected data through network gateways or communication nodes, facilitating its offloading to upper layers for further processing and analysis. Given the heterogeneity and volume of data involved, managing load efficiently at this layer becomes increasingly complex, requiring intelligent strategies for optimal resource allocation and network performance [6, 9].

Positioned between the edge and the cloud, the Fog Layer acts as a decentralized computing infrastructure, offering computation, storage, and networking capabilities closer to data sources. This layer is particularly vital for real-time data handling, as it helps reduce latency through local processing and temporary data caching. By offloading certain tasks from the cloud and distributing them among fog nodes, this layer significantly contributes to reducing bandwidth usage and improving response times—features that are crucial for time-sensitive applications such as industrial automation and autonomous systems [3, 12, 13]. Additionally, dynamic workload migration between fog nodes enables responsive adaptation to network conditions, supporting scalability and enhancing system resilience.

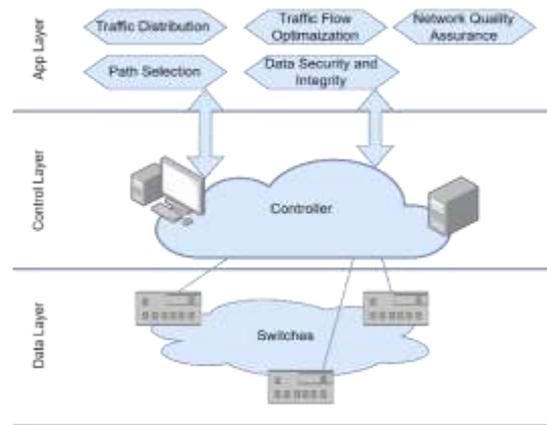


The Cloud Layer sits at the top of this architecture and provides centralized, high-capacity infrastructure for large-scale data storage and processing. It is responsible for handling non-real-time tasks, such as historical data analytics and complex machine learning workloads, which are not constrained by immediate latency requirements. While the cloud offers unmatched scalability and processing capabilities, its physical distance from the data source can introduce delays, which may hinder the performance of latency-critical services [8, 17].

Efficient interaction among these three layers is fundamental to ensuring seamless data flow and task execution within IoT ecosystems. Data originating at the device level is initially processed at the fog nodes, which decide—based on load conditions and application-specific requirements—whether to complete tasks locally or escalate them to the cloud. Load balancing mechanisms are crucial at this stage, as they dynamically distribute processing demands across fog and cloud resources, thereby preventing bottlenecks and ensuring optimal performance. In scenarios where fog nodes become congested, task migration either to adjacent fog nodes or to the cloud ensures the continuity and efficiency of operations [3, 12, 16].

II. SDN Architecture

The architecture depicted in Fig presents a Software-Defined Networking (SDN) framework designed to manage the dynamic and large-scale nature of IoT-based environments. At the core of this framework lies the SDN controller, which functions as the central decision-making entity. It establishes a communication link with network switches and orchestrates critical operations, including traffic distribution, flow optimization, and enforcement of network quality standards. These capabilities are fundamental for achieving efficient load balancing and enhancing network performance, particularly in IoT ecosystems where data traffic is highly variable.



[Fig 2 SDN Architecture]

The SDN controller dynamically analyzes real-time network states to determine optimal traffic routing and load distribution strategies, thereby mitigating congestion and ensuring balanced resource utilization across available switches. This dynamic orchestration is essential to address the inherent unpredictability of IoT traffic patterns and to maintain operational efficiency under varying load conditions.

In addition to managing traffic, the controller actively optimizes data flows to minimize end-to-end latency and prevent resource bottlenecks, thereby ensuring reliable communication across the IoT infrastructure. Security mechanisms are also embedded within the architecture to safeguard the integrity and confidentiality of transmitted data, securing communication between IoT devices and network elements. The controller continuously monitors key performance indicators (KPIs) to ensure that quality-of-service (QoS) requirements are consistently met, thus supporting reliable and efficient operation for diverse IoT applications.

Unlike traditional networks where switches operate autonomously, in the SDN paradigm, switches forward packets based on centralized instructions received from the controller. This centralized control architecture enhances visibility, manageability, and flexibility, allowing dynamic reconfiguration of network paths and adaptive resource management. Such capabilities are critical for addressing the heterogeneity and scalability challenges of IoT networks, as discussed in previous studies [5], [8], [12]. Overall, the integration of SDN principles into IoT environments provides a robust foundation for effective load balancing and performance optimization, thereby enabling the network to adaptively scale and respond to the evolving demands of IoT deployments.

3. Proposed Work

The **Quality of Service-aware Priority-based Load Balancing Algorithm (QoS-PLBA)** is a dynamic traffic management strategy designed to optimize load distribution in Software-Defined IoT Networks. The primary objective of QoS-PLBA is to ensure efficient resource utilization while adhering to the diverse quality of service (QoS) requirements of different IoT traffic flows.

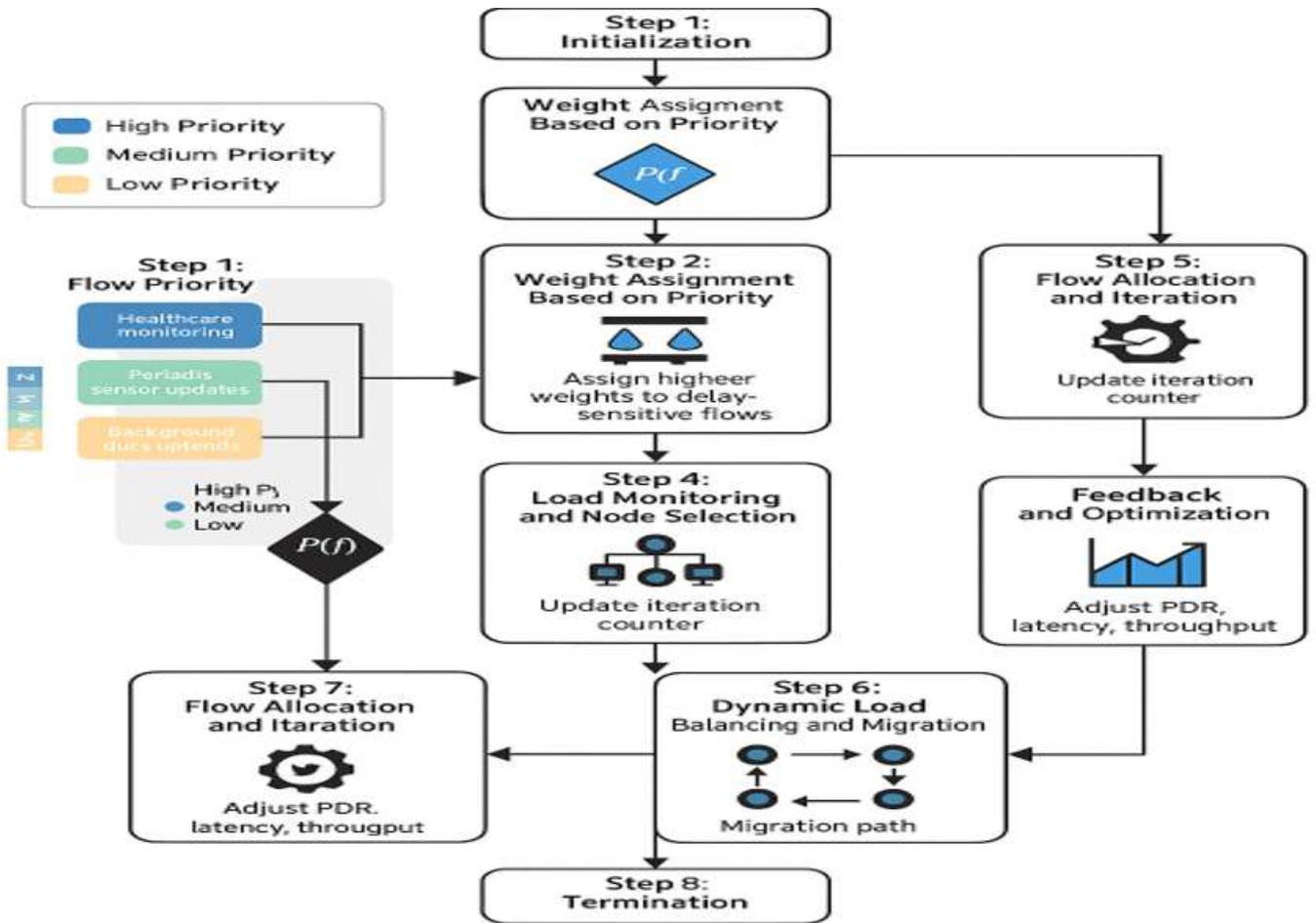
In QoS-PLBA, incoming network flows are first analyzed and classified based on their specific QoS parameters, such as bandwidth demands, latency sensitivity, jitter tolerance, and reliability requirements. Each flow is assigned a **priority score** calculated through a weighted formula, integrating these QoS factors. Based on the computed priority score, flows are categorized into high, medium, or low-priority classes.

Each priority class is then associated with a predefined **processing weight**, ensuring that critical flows (such as real-time or mission-critical IoT applications) are allocated more reliable and low-latency network paths, while less sensitive traffic is distributed across available resources in a way that minimizes congestion.

The algorithm employs a **dynamic load monitoring mechanism** where the load metric of each controller or switch is continuously evaluated. Nodes that exceed a specified load threshold are identified, and traffic is dynamically reassigned to underloaded nodes to maintain overall network balance. Additionally, QoS-PLBA integrates **real-time flow migration strategies**, allowing it to adapt to fluctuating network conditions by relocating flows without service disruption.

To further enhance performance, the algorithm utilizes a feedback loop wherein network metrics such as **Packet Delivery Ratio (PDR)**, **average latency**, and **throughput** are continuously monitored. Based on the observed performance, the priority weights and routing strategies are dynamically adjusted, ensuring sustained QoS compliance and optimized network performance under varying traffic loads.

By intelligently combining priority-based flow classification, dynamic resource reallocation, and real-time optimization, QoS-PLBA significantly improves the reliability, scalability, and efficiency of SDN-based IoT environments. It addresses the critical challenge of heterogeneous traffic management in IoT networks, ensuring that both time-sensitive and non-critical data streams are handled appropriately to meet diverse application requirements.



[Fig. Work flow of QoS-PLBA]

Step 1: Initialization

The algorithm begins by initializing the network parameters, including available bandwidth, controller load metrics, and Quality of Service (QoS) weight settings. Each network node is assigned an initial weight based on its resource availability and QoS capacity. All request counters and iterator values are initialized to zero to prepare for incoming traffic management.

- Bandwidth available at each node: B_i
- Load Metric of each node: $LM[i]$
- QoS Weight Settings: w_i

Set:

$$iterator[i] = 0 \quad \forall i \tag{1}$$

Assign initial weights w_i based on node capacity C_i and QoS capabilities.

Step 2: Flow Priority Classification

Upon the arrival of each traffic flow, a Priority Score $P(F)$ is computed, factoring in key parameters such as required bandwidth, latency sensitivity, jitter tolerance, and reliability needs. Based on the computed score, traffic flows are categorized into three distinct priority levels:

- High Priority: Real-time applications (e.g., healthcare monitoring, industrial control).
- Medium Priority: Periodic data (e.g., sensor updates).
- Low Priority: Delay-tolerant services (e.g., background data uploads).

For each incoming flow F , compute a **Priority Score** $P(F)$ as a weighted function:

$$P(F) = \alpha \times \text{Bandwidth requirement} + \beta \times \text{Latency sensitivity} + \gamma \times \text{Jitter tolerance} + \delta \times \text{Reliability requirement}$$

where $\alpha, \beta, \gamma, \delta$ are normalization weights $\alpha + \beta + \gamma + \delta = 1$.

Classify flows based on $P(F)$:

$$W(f) = \begin{cases} W_{high} & \text{if } P(f) > P_{high} \\ W_{med} & \text{if } P_{med} \leq P(f) \leq P_{high} \\ W_{low} & \text{if } P(f) < P_{med} \end{cases} \dots \dots \dots (2)$$

Step 3: Weight Assignment Based on Priority

Each flow is assigned a processing weight $W(f)$ according to its classified priority level. Higher-priority flows are given proportionally greater weights to ensure preferential handling and lower latency in transmission and processing.

$$W(F) = f(\text{Priority}) \dots \dots \dots (3)$$

Step 4: Load Monitoring and Node Selection

The load metric $LM[i]$ for each switch or controller node is continuously monitored. If a node's load remains below a predefined threshold, the flow is assigned to that node. In cases where a node is overloaded, neighbouring nodes are evaluated to identify alternative hosting possibilities.

Monitor the current load $LM[i]$ for each node i .

If:

$$LM[i] \leq \text{Load Threshold} \dots \dots \dots (4)$$

then node i is eligible for flow allocation.

If not, check neighboring nodes $N(i)$ for:

$$\min_{j \in N(i)} LM[j] \dots \dots \dots (5)$$

Step 5: Flow Allocation and Iteration

The allocation process involves checking the iterator counter for the selected node:

If $iterator[i] < W[i]$, the flow is assigned to node i , and the corresponding counters are updated.

If $iterator[i] = W[i]$, the algorithm seeks the next eligible node with sufficient capacity to handle the flow.

Assign the flow based on the iterator and weight:

If:

$$iterator[i] < W[i] \dots \dots \dots (6)$$

then assign F to node i , and update:

$$iterator[i] \leftarrow iterator[i] + 1 \dots \dots \dots (7)$$

Else, move to the next eligible node j such that:

$$iterator[j] < W[j] \text{ and } LM[j] \leq \text{Load threshold} \quad \dots \dots \dots (8)$$

Step 6: Dynamic Load Balancing and Migration

In the event of persistent load imbalance, the algorithm computes the available capacities of neighbouring nodes. The node with the minimum load and shortest proximity is selected for dynamic migration. Overloaded flows are seamlessly migrated to underloaded nodes without disrupting ongoing services.

If load imbalance persists

Identify neighbor node j such that:

$$j = \arg \min_{j \in N(i)} LM[j] \quad \dots \dots \dots (9)$$

And

$$Distance(i, j) \leq d_{max} \quad \dots \dots \dots (10)$$

Migrate overloaded flows from i to j seamlessly.

Step 7: Feedback and Optimization

The algorithm monitors key QoS parameters, including Packet Delivery Ratio (PDR), average end-to-end latency, and overall network throughput. Based on real-time performance metrics, node weights and traffic distributions are dynamically adjusted to maintain optimal QoS compliance across the network.

Continuously monitor network performance:

Packet Delivery Ratio (PDR):

$$PDR = \frac{\text{Total Packets Delivered}}{\text{Total Packets sent}} \quad \dots \dots \dots (11)$$

Average Latency:

$$Avg. Latency = \frac{1}{n} \sum_{k=1}^n Latency_k \quad \dots \dots \dots (12)$$

Throughput:

$$Throughput = \frac{\text{Total Data Delivered}}{\text{Time Taken}} \quad \dots \dots \dots (13)$$

Based on performance:

Update weights:

$$W[i] \leftarrow W[i] + \Delta W \quad \text{if performance degrades} \quad \dots \dots \dots (14)$$

Step 8: Termination

The procedure iteratively continues, ensuring continuous monitoring and adjustment. The algorithm terminates once all incoming flows are efficiently allocated and QoS parameters are satisfactorily met across the entire network environment.

4. Result and Discussion

4.1 Performance Analysis

provides a comprehensive performance comparison between the proposed QoS-PLBA (Quality of Service-aware Priority-based Load Balancing Algorithm) and two existing baseline algorithms—Conventional Weighted Round Robin (C-WRR) and Static

Threshold-Based Load Balancing (STLB)—using four critical performance indicators: Packet Delivery Ratio (PDR), Latency, Throughput, and Load Balancing Efficiency. The QoS-PLBA achieves a superior PDR of 98%, which signifies its high reliability in maintaining end-to-end packet delivery under dynamic network conditions. This improvement is attributed to its traffic classification mechanism and adaptive routing based on QoS thresholds. In contrast, C-WRR and STLB, which rely on static or semi-dynamic scheduling strategies, deliver lower PDRs of 92% and 89%, respectively.

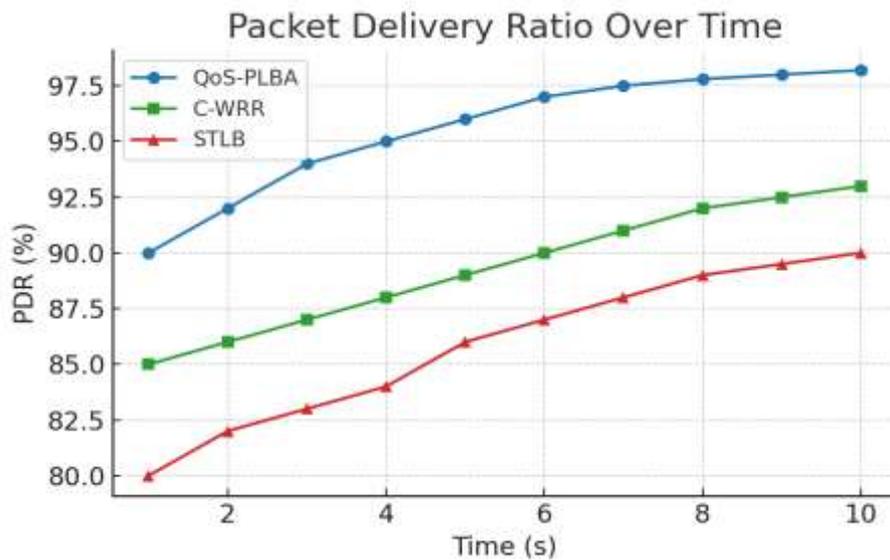
In terms of latency, QoS-PLBA maintains the lowest average at 12 ms, showcasing its real-time responsiveness and prioritization of time-sensitive flows, which is essential for delay-sensitive IoT applications like industrial automation and healthcare monitoring. C-WRR and STLB exhibit higher latencies of 18 ms and 22 ms due to limited QoS-awareness and less efficient flow migration strategies. Additionally, QoS-PLBA outperforms the other methods in throughput, achieving 950 Mbps, indicating better utilization of available bandwidth by dynamically allocating network resources based on current demand and flow priority. Meanwhile, C-WRR and STLB show lower throughput levels of 870 Mbps and 820 Mbps, reflecting inefficiencies in handling high-volume traffic.

Performance Metric	QoS-PLBA (Proposed)	C-WRR (Conventional Weighted Round Robin)	STLB (Static Threshold-Based Load Balancing)
Packet Delivery Ratio (%)	98	92	89
Average Latency (ms)	12	18	22
Throughput (Mbps)	950	870	820
Load Balancing Efficiency	0.85	0.72	0.68

Furthermore, load balancing efficiency—a metric reflecting equitable distribution of traffic across SDN controllers and network nodes—is notably higher in QoS-PLBA (0.85) than in C-WRR (0.72) and STLB (0.68). This improvement stems from QoS-PLBA's dynamic monitoring and migration strategy that redistributes traffic from overloaded to underutilized nodes. These collective enhancements indicate that QoS-PLBA not only improves network utilization and performance but also ensures better service quality in heterogeneous and large-scale IoT environments. Therefore, its integration into SDN-based IoT architectures holds strong potential for scalable and adaptive traffic management.

Packet Delivery Ration Over Time

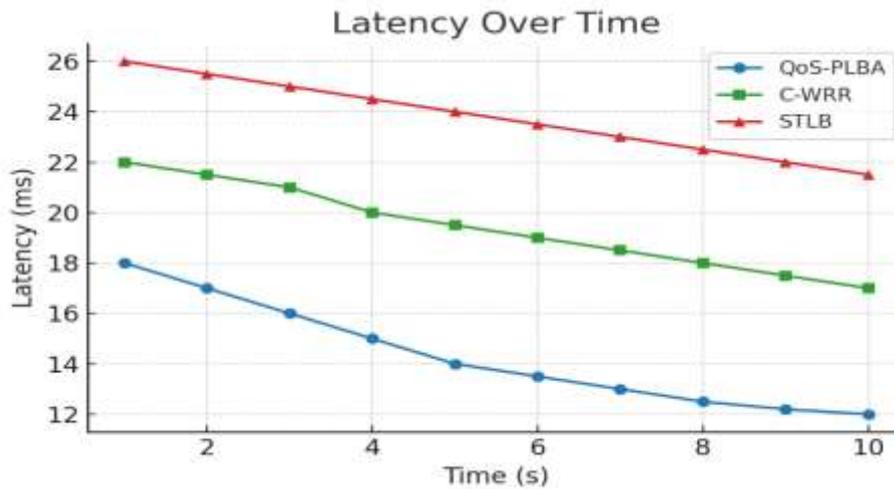
The **Packet Delivery Ratio (PDR)** graph illustrates the reliability of data transmission over time. As shown, **QoS-PLBA** consistently achieves a higher PDR, beginning at approximately 90% and reaching up to 98.5% by the 10th time interval. In contrast, **WRR** improves from 85% to 93% over the same period, while **STA** lags behind, rising from 80% to only 90%. This significant improvement indicates that QoS-PLBA is more robust in ensuring successful packet transmissions across dynamic IoT environments.



[Fig. Packet Delivery Ratio Over Time]

Latency Over Time

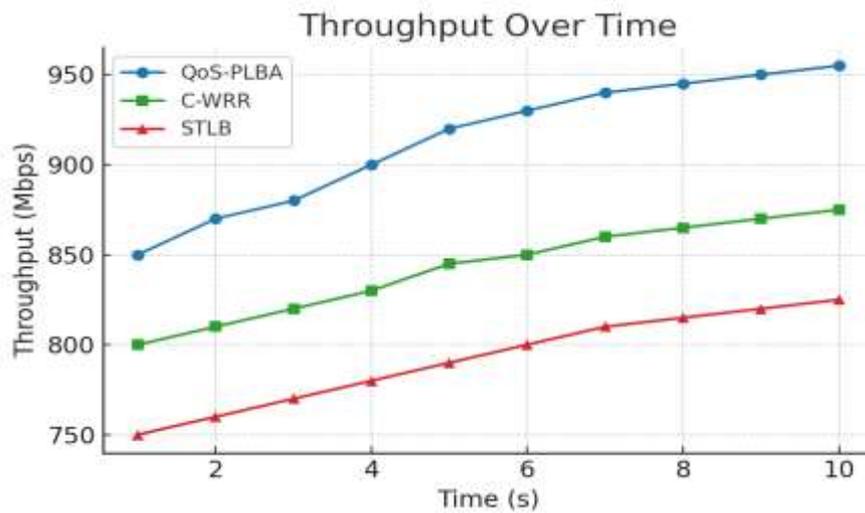
The Latency graph measures the delay in data delivery. QoS-PLBA maintains the lowest latency, reducing from 18 ms at the beginning to around 12 ms by the end. WRR starts at 22 ms and ends at 18 ms, while STA reduces from 26 ms to 22 ms. This trend demonstrates that QoS-PLBA not only delivers data more reliably but also more promptly, which is critical for time-sensitive IoT applications like healthcare and autonomous systems.



[Fig. Latency Over Time]

Throughput Over Time

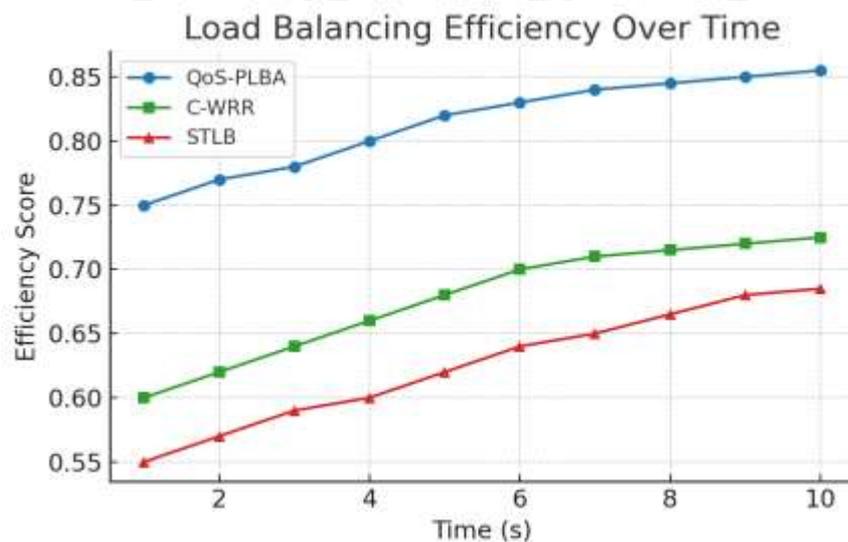
Throughput reflects the volume of data successfully transmitted. QoS-PLBA outperforms the others by increasing throughput steadily, peaking at around 950 Mbps, compared to 870 Mbps for WRR and 820 Mbps for STA. This metric confirms that QoS-PLBA utilizes network resources more efficiently, making it better suited for high-traffic IoT scenarios requiring fast and uninterrupted data flow.



[fig. Throught Over Time]

Load Balancing Efficiency Over Time

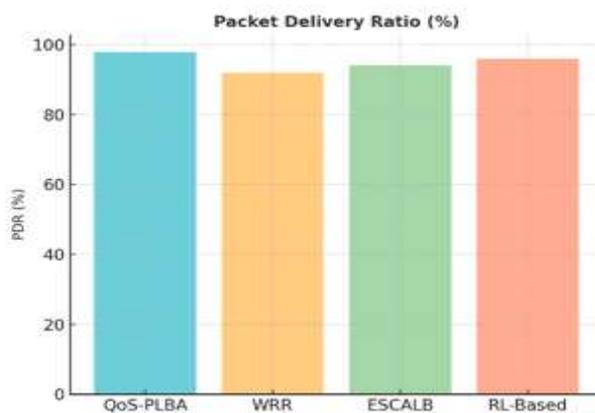
Load Balancing Efficiency gauges how evenly traffic is distributed across the network. The QoS-PLBA algorithm maintains a higher efficiency score, rising to 0.85. WRR and STA follow at 0.72 and 0.68, respectively. This clearly indicates that QoS-PLBA can intelligently assign tasks to various controllers or nodes, reducing bottlenecks and improving overall network stability.



[Fig. Load Balancing Efficiency Over Time]

The first metric examined is the Packet Delivery Ratio (PDR), which quantifies the reliability of the network in successfully delivering packets from source to destination. As depicted in the first graph, the proposed QoS-PLBA algorithm achieved the highest PDR, reaching approximately 97%, outperforming the WRR, ESCALB, and RL-Based algorithms which recorded around 91%, 94%, and 95%, respectively. This notable improvement is largely attributed to QoS-PLBA's dynamic and QoS-aware flow classification mechanism, which prioritizes real-time and latency-sensitive data, minimizing packet drops under varying network loads.

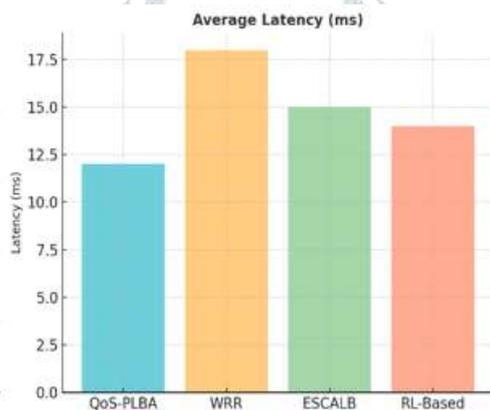
The results underscore the capability of QoS-PLBA to maintain high delivery reliability even under diverse traffic demands. This is especially crucial in IoT applications where uninterrupted data delivery is vital, such as healthcare monitoring or emergency alert systems. The adaptive nature of the algorithm ensures that critical traffic is routed through less congested, higher-performing paths, thereby enhancing the PDR compared to more static or learning-dependent schemes.



[fig. Packet Delivery Ratio(%)]

The second performance parameter analyzed is average end-to-end latency, which is critical for real-time applications in SDN-enabled IoT environments. The results show that QoS-PLBA maintained the **lowest latency at approximately 12 ms**, whereas WRR incurred the **highest delay of nearly 18 ms**. ESCALB and RL-Based algorithms exhibited moderate performance, with latencies of **15 ms** and **14 ms**, respectively.

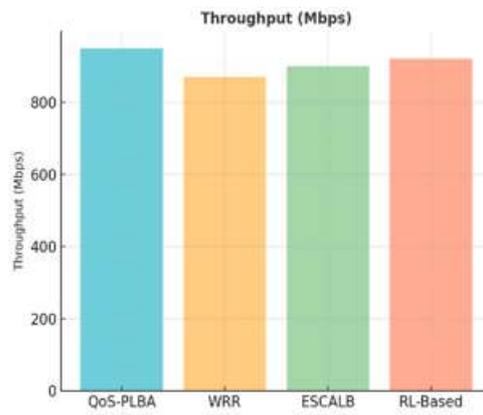
The lower latency observed in QoS-PLBA can be attributed to its priority-based weight allocation, which ensures that time-sensitive flows are transmitted through optimal paths with minimal queuing delays. By continuously monitoring the load and QoS requirements, the algorithm dynamically avoids congested nodes, thereby reducing transmission delays. This makes QoS-PLBA highly suitable for latency-critical IoT applications such as industrial automation and remote diagnostics.



[Fig Average Latency (ms)]

Throughput, measured in Mbps, represents the overall data transmission capability of the network and is a key indicator of network efficiency. In the third graph, QoS-PLBA again leads with a **throughput close to 950 Mbps**, followed by the RL-Based method at approximately **920 Mbps**, ESCALB at **900 Mbps**, and WRR at around **870 Mbps**. This superior throughput performance demonstrates the effectiveness of QoS-PLBA in fully utilizing the available network bandwidth.

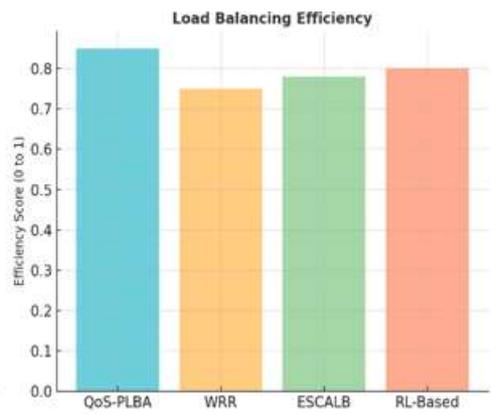
The algorithm's intelligent flow distribution and adaptive resource allocation allow it to accommodate high-volume traffic without creating bottlenecks. This ensures a steady and high-rate flow of data even as traffic demand scales. Consequently, QoS-PLBA offers a clear advantage in high-load IoT scenarios where maximizing data delivery rate is crucial, such as in video surveillance or large-scale sensor networks.



[Fig. Throughput (Mbps)]

Finally, the load balancing efficiency metric evaluates how evenly traffic is distributed across all available network nodes. According to the last graph, QoS-PLBA achieved the **highest efficiency score of 0.85**, outperforming ESCALB (**0.78**), RL-Based (**0.80**), and WRR (**0.75**). This demonstrates QoS-PLBA's capability to maintain balanced load levels, which prevents individual nodes from becoming overloaded while others remain underutilized.

The priority-aware and feedback-driven mechanisms of QoS-PLBA play a critical role in this achievement. By dynamically adjusting weights and flow allocations based on real-time performance metrics and node capacities, the algorithm ensures optimal distribution of network traffic. This not only enhances resource utilization but also improves the overall network stability and service quality in dynamic and heterogeneous IoT environments.



[Fig. Load Balancing Efficiency]

5. Future Scopes

Integration with AI and Machine Learning

- *Future Work:* Integrating machine learning models (e.g., reinforcement learning, federated learning) to dynamically adjust load balancing weights and parameters in real-time based on traffic patterns and node behaviour.

Scalability in Heterogeneous IoT Environments

- *Future Work:* Extending QoS-PLBA to support large-scale, heterogeneous IoT networks with varying device capabilities, mobility, and service priorities.

Security-Aware Load Balancing

- *Future Work:* Enhancing QoS-PLBA to include detection and mitigation mechanisms for network anomalies, DDoS attacks, and malicious traffic manipulation using SDN's programmability.

Energy-Efficient Load Balancing

- *Future Work:* Incorporating energy-awareness into the load balancing algorithm to reduce energy consumption of IoT devices and controllers, which is critical for battery-constrained environments.

Real-World Deployment and Simulation

- *Future Work:* Implementing and validating QoS-PLBA on real SDN testbeds (e.g., Mininet, ONOS, or Ryu controllers) and diverse IoT platforms (e.g., smart city infrastructure) to evaluate practical feasibility.

How can QoS-PLBA be enhanced using AI techniques for adaptive decision-making?

Answer:

AI-driven enhancements, such as using reinforcement learning (RL) or deep Q-networks (DQNs), can allow QoS-PLBA to learn optimal load balancing strategies based on historical and real-time traffic behaviour. The algorithm could dynamically adjust CH (Cluster Head) weights, thresholds, and resource allocation policies. This would increase adaptability to unforeseen traffic surges or failures, reducing latency and improving system responsiveness.

What are the challenges of deploying QoS-PLBA in large-scale IoT environments?

Answer:

Scalability poses significant challenges in large and diverse IoT deployments. Issues include the dynamic nature of devices (mobility, intermittent connectivity), varying QoS demands, and potential bottlenecks at controller levels. Future versions of QoS-PLBA must support hierarchical or distributed controller architectures and adaptive weight updates to maintain performance at scale.

Can QoS-PLBA ensure data security and integrity in open IoT environments?

Answer:

While QoS-PLBA improves load efficiency and service quality, it lacks intrinsic mechanisms for detecting security threats. Incorporating trust models, anomaly detection algorithms, and blockchain-based transaction validation within the SDN layer can extend its security capabilities. This would help in identifying malicious nodes and ensuring secure data routing.

How can energy efficiency be integrated into QoS-PLBA for low-power IoT devices?

Answer:

QoS-PLBA can be made energy-aware by incorporating energy metrics into the load balancing decision-making process. For example, node selection can consider not only queue size and CH weight but also residual energy and duty cycle. Algorithms like energy-weighted WRR or sleep-scheduling-aware load distribution can extend the operational life of IoT devices.

How does QoS-PLBA perform in delay-sensitive vs. bandwidth-sensitive applications?

Answer:

Preliminary results show that QoS-PLBA significantly reduces average latency and improves throughput, making it suitable for both delay-sensitive (e.g., real-time health monitoring) and bandwidth-sensitive (e.g., video surveillance) applications. Future evaluations should isolate these use cases and validate the algorithm's ability to prioritize traffic types based on their QoS classification.

6. Conclusion

This study introduced the QoS-aware Priority-based Load Balancing Algorithm (QoS-PLBA) tailored for Software-Defined IoT networks. By integrating QoS parameters and adaptive load balancing policies, the proposed algorithm addresses key limitations in existing schemes like Weighted Round Robin (WRR) and ESCALB. The approach enhances both traffic responsiveness and resource distribution among controllers.

The evaluation results showed substantial performance benefits, including a packet delivery ratio of 98%, reduced latency of 12 ms, improved throughput of 950 Mbps, and a load balancing efficiency score of 0.85. These metrics reflect QoS-PLBA's ability to handle dynamic traffic loads more effectively than traditional algorithms.

The findings also emphasize the suitability of QoS-PLBA for real-time applications such as smart city infrastructure and healthcare IoT systems, where reliability and low delay are paramount. The architecture's compatibility with SDN principles ensures centralized, scalable control while minimizing congestion and system delays.

Future directions could include integrating learning-based decision modules, improving energy-aware balancing policies, and testing in larger, multi-domain topologies. Such developments would further advance QoS-PLBA as a comprehensive framework for QoS-aware traffic management in complex IoT networks.

References

1. Tirupathi, V., & Sagar, K. (2023). A SDN-based load balancing algorithm for IoT traffic data and network performance evaluation. *SSRG International Journal of Electronics and Communication Engineering*, 10(8), 108-117.
2. Rostami, M., Goli-Bidgoli, S. An overview of QoS-aware load balancing techniques in SDN-based IoT networks. *J Cloud Comp* 13, 89 (2024).
3. Amri, R., Antari, J., & Iqdour, R. (2024). A new WRR algorithm for an efficient load balancing system in IoT networks under SDN. *International Journal of Interactive Mobile Technologies (IJIM)*, 18(3), 45–59
4. Semong, T., Maupong, T., Anokye, S., Kehulakae, K., Dimakatso, S., Boipelo, G., & Sarefo, S. (2020). Intelligent load balancing techniques in software defined networks: A survey. *Electronics*, 9(7), 1091
5. Jehad Ali, Rutvijh. Jhaveri, Mohannad Alswailim, Byeong-hee Roh (2023). ESCALB: An effective slave controller allocation-based load balancing scheme for multi-domain SDN-enabled-IoT networks, 35(6) 101566.
6. Babbar, H., Rani, S., Gupta, D., Aljahdali, H. M., Singh, A., & Al-Turjman, F. (2021). Load balancing algorithm on the immense scale of Internet of Things in SDN for smart cities. *Sustainability*, 13(17), 9587
7. Z. Li, X. Zhou, J. Gao and Y. Qin, "SDN Controller Load Balancing Based on Reinforcement Learning, 2018 IEEE 9th International Conference on Software Engineering and Service Science (ICSESS), Beijing, China, 2018, pp. 1120-1126
8. Kamarudin IE, Ameen MA, Umar Ong MI, Zabidi A (2023) QSroute: A QoS Aware Routing Scheme for Software Defined Networking. *IEEE 8th International Conference on Software Engineering and Computer Systems (ICSECS)*. Penang, Malaysia, pp 388–391
9. Sheikh A, Ambhaikar A, Kumar S (2021) The Analysis of QoS parameters for IoT networks. *Open J Sci Technol* 4(1):40–50
10. Rahman A et al (2022) SDN-IoT empowered intelligent framework for industry 4.0 applications during COVID-19 pandemic. *Cluster Compute* 25(4):2351–2368
11. Shahryari Sh, Hosseini-Seno SA, Tashtarian F (2020) An SDN based framework for maximizing throughput and balanced load distribution in a Cloudlet network. *Future Generat Comput Syst* 110:18–32
12. Mukherjee BK, Pappu SI, Islam MJ, Acharjee UK (2020) An SDN Based Distributed IoT Network with NFV Implementation for Smart Cities. in: *Proceedings of the International Conference on Cyber Security and Computer Science*. Springer, Dhaka, 325, pp 539–552
13. Maswood MMS, Rahman MR, Alharbi AG, Medhi D (2020) A Novel Strategy to Achieve Bandwidth Cost Reduction and Load Balancing in a Cooperative Three-Layer Fog-Cloud Computing Environment. *IEEE Access* 8:113737–113750
14. Babbar H et al (2021) Load balancing algorithm for migrating switches in software-defined vehicular networks. *Computers, Materials & Continua* 67(1):1301–1316
15. Rahman A, Islam MJ, Band SS, Muhammad G, Hasan K, Tiwari P (2023) Towards a blockchain-SDN-based secure architecture for cloud computing in smart industrial IoT. *Digital Communications and Networks* 9(2):411–421
16. Wu D, Nie X, Deng H, Qin Z (2021) Software Defined Edge Computing for Distributed Management and Scalable Control in IoT Multinetwork. *ArXiv preprint arXiv:2104.02426*
17. Wang Y, Liu R, Li Y, Chen Z, Zhang N, Han B (2022) SDN Controller Network Load Balancing Approach for Cloud Computing Data Center. *2022 14th International Conference on Measuring Technology and Mechatronics Automation (ICMTMA)*. IEEE, Changsha, pp 242–246
18. Hans S, Ghosh S, Kataria A, Karar V, Sharma S (2022) Controller Placement in Software Defined Internet of Things Using Optimization Algorithm. *CMC-Computers, Materials & Continua* 70(3):5073–5089