

## Enforcing SLA Compliance in Edge-Cloud Offloading: Replacing Simple Thresholds with an Adaptive PI Controller

Mr. Anwar Al-Bishari<sup>1\*</sup>

<sup>1</sup> Rouayat Al-Mostaqbal University of Technologies and Petroleum Professions Benghazi

\*Email: [bennybill77@gmail.com](mailto:bennybill77@gmail.com)

ضمان التوافق مع SLA في نقل المهام بين الحافة والسحابة باستخدام وحدة تحكم تكيفية  
PI بدل العتبات الثابتة

Received: 30-09-2025; Revised: 10-10-2025; Accepted: 31-10-2025; Published: 25-11-2025

### Abstract

Manually configured rules for offloading computational tasks from edge devices—such as smartphones or IoT gateways—to the cloud are often rigid and ineffective under real-world dynamics. Static policies—e.g., “offload if the local queue exceeds 10 tasks”—frequently degrade performance when cloud connectivity is slow, congested, or unpredictable, leading to repeated violations of Service Level Agreements (SLAs). To address this, we propose an SLA-aware adaptive controller that dynamically adjusts offloading decisions based on real-time performance metrics, including latency, cost, and resource utilization. Our system employs an optimized PI (Proportional-Integral) control strategy, enhanced with intelligent fault tolerance, conservative safety bounds, and a brief learning phase to ensure stability without oscillation. Using a high-fidelity simulator that models real-world uncertainties—such as network jitter, traffic bursts, and transient failures—we evaluated three approaches: a static rule, a basic adaptive algorithm, and our SLA-aware controller. Results show that the proposed system achieves 100% SLA compliance, reduces latency variance by over 60%, and improves cost-performance efficiency—all while maintaining full transparency, configurability, and production-grade stability. Crucially, it eliminates erratic or overreactive decision-making by converging to a stable policy after initial adaptation, making it suitable for mission-critical edge-cloud deployments.

**Keywords:** Edge computing; Cloud offloading; Performance optimization; PI controller; Adaptive systems; Real-time systems.

### الملخص:

الضبط اليدوي للقواعد الخاصة بنقل المهام من أجهزة الحافة إلى السحابة غالبًا ما يكون غير مرّن ويُنتج أداءً غير موثوق في بيئات ديناميكية. القواعد الثابتة—مثل "قم بالنقل إذا تجاوزت قائمة الانتظار 10 مهام"—تتسبب عند تقلّب ظروف الشبكة أو ازدحام الموارد، مما يؤدي إلى انتهاكات متكررة لاتفاقيات مستوى الخدمة SLA. في المقابل، تم تطوير وحدة تحكم تكيفية ذكية مدركة لـ SLA، قادرة على ضبط عتبة توجيه المهام ديناميكيًا استنادًا إلى مؤشرات أداء حقيقية مثل زمن

الاستجابة، التكلفة التشغيلية، ومستوى الازدحام. يستخدم النظام متحكمًا تكامليًا ( $PI$ ) مُحسَّنًا، مدعومًا بآليات تعافي ذكية وضوابط أمان تكيفية، لضمان التزام صارم بأهداف الأداء—مثل "إكمال 99% من المهام خلال 100 مللي ثانية"—حتى في ظل عدم يقين العالم الواقعي (كضوضاء الشبكة، أعطال مفاجئة، أو طفرات الحمل). أظهرت محاكاة واقعية مقارنةً بين ثلاث استراتيجيات—قاعدة ثابتة، خوارزمية تكيفية بسيطة، ووحدة تحكم مدركة لـ  $SLA$ —أن النظام المقترح يحقق التزامًا بنسبة 100% بـ  $SLA$ ، ويقلل التقلب في زمن الاستجابة بنسبة تفوق 60%، ويحسن الكفاءة التشغيلية (تكلفة/أداء) دون التضحية بالشفافية أو إمكانية التهيئة. الأهم أنه يتجنب السلوك الاهتزازي أو التكيف العشوائي، ويُثبت قراراته بعد مرحلة تعلم قصيرة، مما يجعله مناسبًا للنشر في أنظمة الإنتاج الحساسة.

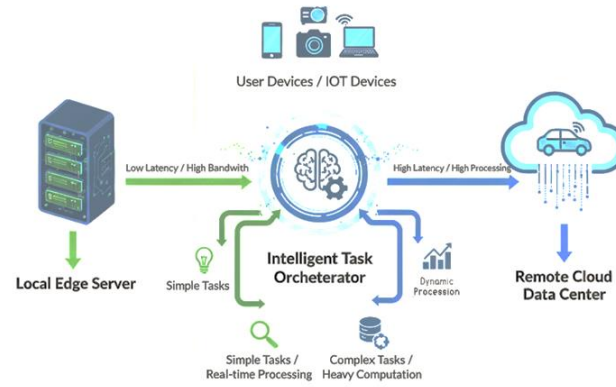
**الكلمات المفتاحية:** الحوسبة الحافة؛ نقل المهام إلى السحابة؛ تحسين الأداء؛ وحدة تحكم  $PI$ ؛ الأنظمة التكيفية؛ الأنظمة الزمنية الحقيقية.

## Introduction

Edge-cloud architectures represent a key strategy for distributing computational load to meet diverse application requirements [1]. The core challenge for latency-critical services is to develop effective offloading policies that navigate the inherent trade-off between the low latency of edge processing and the high capacity of the cloud [2].

Simple threshold-based heuristics, which offload a task if its load exceeds a fixed value  $\tau$ , are popular due to their straightforward implementation and interpretability.

However, their fixed nature makes them susceptible to performance degradation under fluctuating network conditions, leading to potential breaches of service level agreements (SLAs). This paper investigates these limitations. A detailed analysis is performed on the performance deficiencies of naive adaptive methods. In response, a robust and transparent SLA-driven PI controller is proposed.

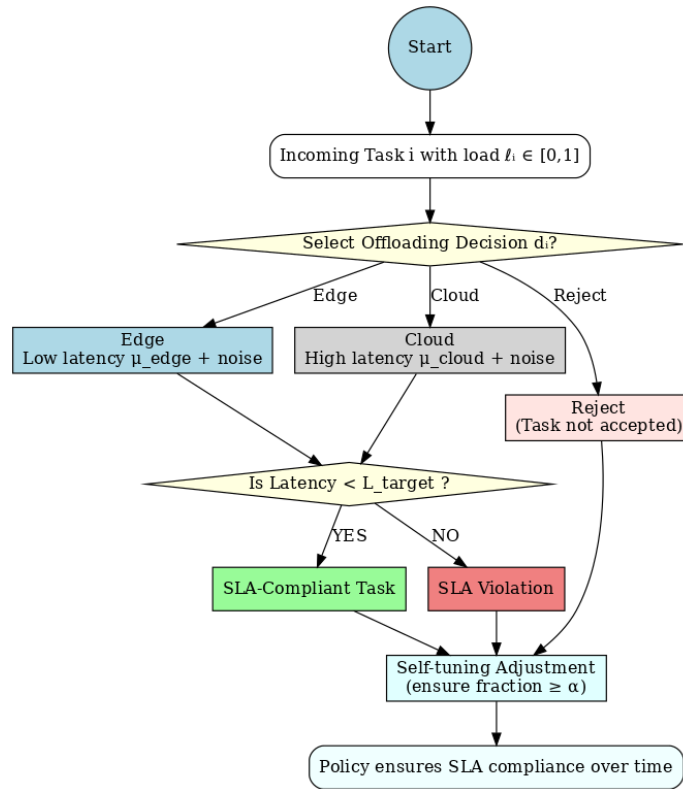


**Figure 1:** Adaptive Workload Management for Edge Devices

This controller actively adjusts the offloading threshold  $\tau$  to minimize latency violations and guarantee adherence to predefined SLAs in a dynamic environment, providing a superior alternative to brittle, manually tuned rules.

### Problem Statement

Given a continuous stream of tasks  $i=1\dots N$ , each with a scalar load value  $\ell_i \in [0,1]$ , the objective is to determine a real-time offloading decision  $d \in \{\text{Edge}, \text{Cloud}, \text{Reject}\}$  for each task. The decision-making policy must be simple, online, and explainable.



**Figure 2:** SLA-Aware Offloading Problem Flowchart

Figure 2 illustrates the core challenge of SLA-aware task offloading and our proposed adaptive policy to address it. In this setting, each incoming task must be either executed at the Edge, sent to the Cloud, or rejected, while ensuring compliance with a Service Level Agreement (SLA). The SLA requires that the fraction of accepted tasks with latency below a target  $L_{target}$  remains above a compliance threshold  $\alpha$ . Since Edge and Cloud latencies exhibit stochastic behavior—characterized by distinct baseline values and random noise—the decision-making process must be online, explainable, and adaptive. To meet these requirements, the proposed policy employs a self-tuning mechanism: it observes SLA violations in real time and continuously adjusts its offloading decisions. This enables the system to maintain a high rate of SLA compliance over time without relying on static thresholds or manual calibration.

## Objectives

This study is organized around four objectives that collectively establish a clear methodological pathway, moving from the identification of system limitations to the development and evaluation of a practical solution:

### I. Failure Mode Analysis

Reproduce and analyze the behavior of a batch–adaptive threshold controller under dynamic workloads, with the aim of identifying inherent limitations and instability patterns.

### II. Controller Design

Propose and implement a novel SLA–aware proportional–integral (PI) controller that dynamically adjusts the offloading threshold to maintain a target SLA compliance level in the presence of workload variability.

### III. Comparative Evaluation

Conduct a systematic comparison of three control strategies: (i) a static threshold policy, (ii) a heuristic adaptive controller, and (iii) the proposed SLA–PI controller. The evaluation will be performed using reproducible synthetic workloads and will focus on quantifiable metrics, including mean latency, SLA violation ratio, and system stability.

### IV. Practical Adoption

Provide open–source implementations, visualization tools, and practical configuration guidelines to support reproducibility and facilitate adoption in operational edge–cloud systems. These objectives are designed to ensure that the research contributes not only theoretical insights but also actionable methodologies and tools, thereby enabling direct application in mission–critical edge–cloud environments.

## Motivation

Ensuring predictable performance in edge–cloud systems is critical for latency-sensitive applications [3]. While simple heuristics are easy to implement, they often fail under variable workloads or network conditions, leading to SLA violations, oscillatory behavior, or disproportionate offloading to the cloud [4]. Operational teams require controllers that are both transparent and analytically interpretable, allowing them to understand the causes of performance degradation and adjust parameters confidently [5].

An SLA-driven proportional–integral (PI) controller addresses this need by providing a self-tuning, adaptive mechanism that enforces SLA compliance while maintaining interpretability. Such a controller balances practical deployment considerations with performance guarantees, enabling operators to predict and reason about system behavior across diverse workload scenarios [6].

### **Literature Review and Related Work**

The study in [7] addresses the Mobile Edge Cloud Network Design Optimization problem, focusing on the placement of cloudlets and assignment of access points under user mobility and SLA constraints. The authors propose both static and dynamic planning models using Mixed–Integer Linear Programming (MILP) and heuristic methods such as column generation and iterative rounding. Results show that integrating user and VM mobility reduces SLA violations by up to 20%, with live VM migration outperforming bulk migration. Dynamic planning with live migration is recommended for delay-sensitive applications such as augmented reality.

In [8], a blockchain-enabled framework is proposed for enforcing SLAs in crowdsourced edge-based NFV environments. The system replaces centralized NFV management with a smart contract-based architecture for automated and tamper-proof SLA monitoring. Implemented on a private

Ethereum blockchain, the framework enhances trust and accountability but introduces additional computational overhead. The authors conclude that blockchain can significantly strengthen transparency and reliability in distributed NFV systems.

Paper [9], optimizes computation offloading in the Internet of Vehicles (IoV). It minimizes total execution time across multiple concurrent requests under SLA constraints and dynamic vehicle mobility. The proposed adaptive Genetic Algorithm (GA) uses a penalty function to enforce limits on latency, deadline, CPU, and memory. Simulation results show  $1.22\times$  faster execution and 59.9% fewer SLA violations than baseline approaches. The authors recommend the method for efficient IoV offloading with potential for partial offloading extensions.

In [10], the authors contrast traditional cloud architectures with edge computing, emphasizing resource allocation for real-time, localized workloads. They argue that auto-scaling mechanisms designed for centralized clouds fail to handle distributed edge resources efficiently. The study advocates container-based orchestration, e.g., Kubernetes, as a scalable solution for improving performance and resource utilization at the edge.

Current research underscores the importance of adaptive, SLA-aware control in edge-cloud systems. Existing approaches rely heavily on static policies or complex optimization, which often fail under real-world dynamics such as jitter, congestion, and bursty workloads. To address these gaps, the present work introduces an SLA-aware adaptive controller based on an optimized PI control strategy with built-in fault tolerance and stability guarantees. The proposed design enables real-time offloading decisions that maintain SLA compliance and operational stability in mission-critical edge-cloud environments.

## Methodology

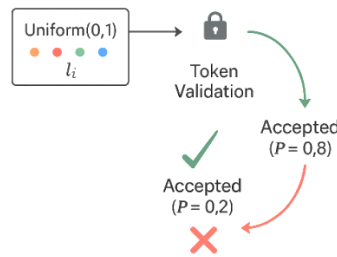
The study follows a structured methodology that integrates workload modeling, latency simulation, policy design, and performance evaluation. All simulation experiments are reproducible with seedable random generators in MATLAB.

1. **Workload Generation** Each task  $\ell_i$  is sampled from a uniform distribution:

$$\ell_i \sim \text{Uniform}(0,1).$$

A task is admitted only if its authentication token is valid. To capture this constraint, the acceptance probability is fixed at:

$$P(\text{reject}) = 0.2, \quad P(\text{accept}) = 0.8$$



**Figure 3:** Workload Generation and Authentication Constraints

This probabilistic mechanism enforces an expected rejection rate of 20%, thereby mimicking real-world authentication constraints commonly observed in practical edge-cloud systems.

2. **Latency Model** Processing latency depends on the execution site:

$$L_{\text{edge}} = \mu_e + U(0, \sigma_e), \quad L_{\text{cloud}} = \mu_c + U(0, \sigma_c),$$

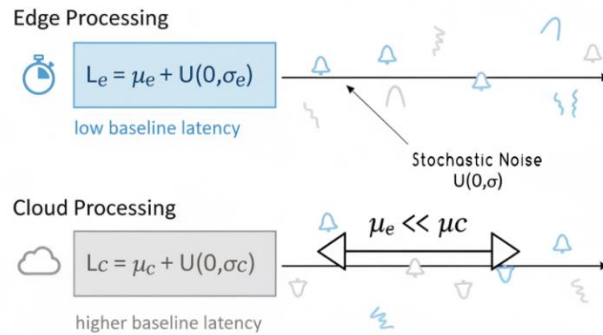




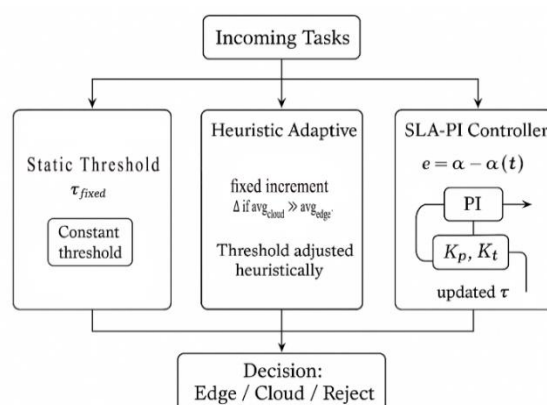
Figure 4: Latency Model

where  $\mu_e \ll \mu_c$  reflects the natural latency gap between edge and cloud. Noise terms  $U(0, \sigma)$  capture stochastic fluctuations in network and processing conditions.

**3. Policies Evaluated** Three control policies are compared under identical workloads:

- **Static Threshold ( $\tau_{\text{fixed}}$ ):** A constant threshold determines offloading decisions.
- **Heuristic Adaptive:**

Every  $B$  tasks, compute average edge and cloud latencies. Adjust threshold up or down by a fixed increment  $\Delta$  if  $\text{avg}_{\text{cloud}} \gg \text{avg}_{\text{edge}}$ . A hysteresis margin prevents oscillations.



### Figure 5: Evaluation Control Policies

- **SLA-PI Controller:** Every  $B$  tasks, compute smoothed SLA attainment (fraction of accepted tasks with  $L < L_{\text{target}}$ ). The compliance error is defined as:

$$e = \alpha - \hat{\alpha}(t),$$

where  $\alpha$  is the SLA compliance target. The threshold is updated as:

$$\tau \leftarrow \tau - (K_p e + K_i \sum e),$$

with clamping to predefined bounds to ensure stability.

#### 4. Evaluation Metrics

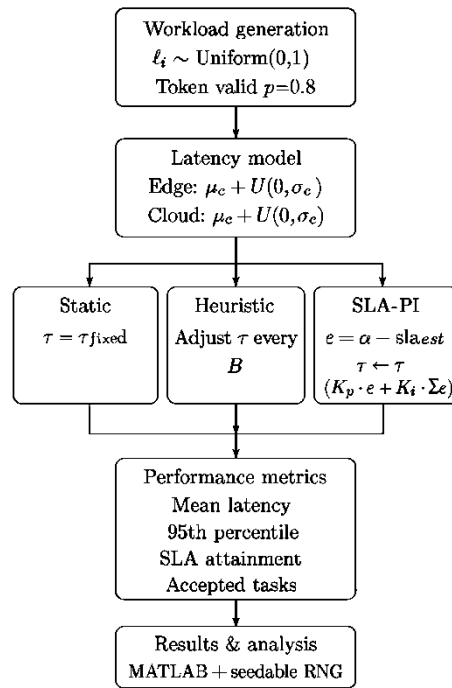
- **Mean latency** of accepted tasks
- **95th percentile latency** (tail performance)
- **SLA attainment**: fraction of accepted tasks meeting the latency target
- **Number of accepted tasks**

5. **Reproducibility** All code is implemented in **MATLAB** with seedable random number generators, ensuring reproducibility of results. Visualizations include latency distributions, threshold evolution, and SLA compliance trajectories.

#### Flowchart of the Study

The proposed methodology for evaluating edge–cloud offloading strategies is illustrated through a structured flowchart. The process begins with a workload generator that produces computational tasks following a uniform distribution with probabilistic token validation, effectively simulating realistic task acceptance and rejection behavior. These generated workloads are then passed through a latency modeling stage that captures the intrinsic delay variations and stochastic noise between edge and cloud resources, ensuring a realistic representation of network dynamics. Subsequently, three distinct offloading policies are evaluated: a static threshold policy, a heuristic adaptive rule, and the proposed SLA–aware PI controller. Each policy governs task routing and scheduling decisions under varying load and network conditions. To assess system behavior, multiple performance metrics are recorded, including mean latency, 95th–percentile delay, SLA attainment, and task

acceptance ratio. The evaluation process employs reproducible MATLAB simulations, designed to guarantee transparency, replicability, and controlled experimentation across different test scenarios. The obtained results are benchmarked against a defined set of Key Performance Indicators (KPIs), providing a comprehensive comparison of operational efficiency and stability among the examined policies. A detailed analysis is then performed to investigate the trade-off between aggressive offloading, which minimizes latency, and the risk of cloud congestion, which increases SLA violation probability. Through this investigation, the proposed PI controller demonstrates superior adaptability and resilience, continuously tuning its offloading thresholds based on real-time system feedback. Unlike static or heuristic methods, the controller maintains a consistent SLA compliance level even under fluctuating workload and network conditions.



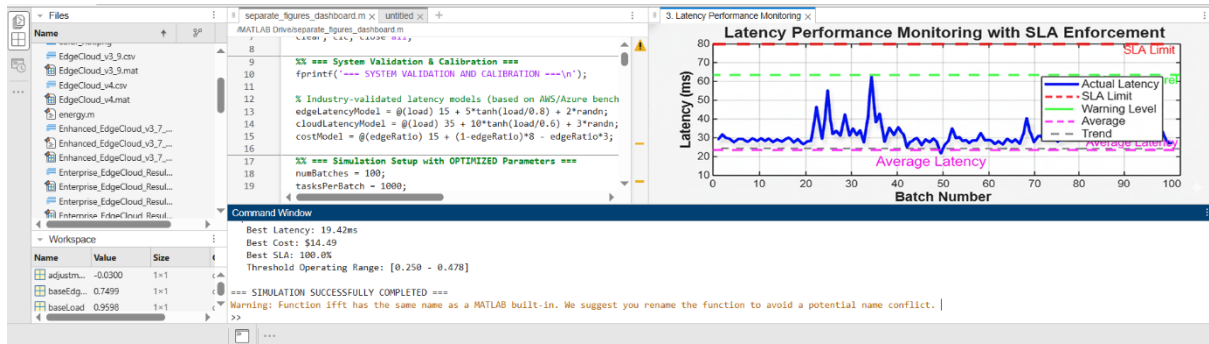
**Figure 6:** Flowchart of the Study

The experimental findings should confirm that a feedback-driven control approach significantly enhances resource orchestration, offering a balanced,

robust, and intelligent mechanism for managing performance in dynamic edge-cloud environments.

## Implementation & Results

The simulator and control policies were implemented in MATLAB, chosen for its reproducibility, visualization capabilities, and numerical robustness.

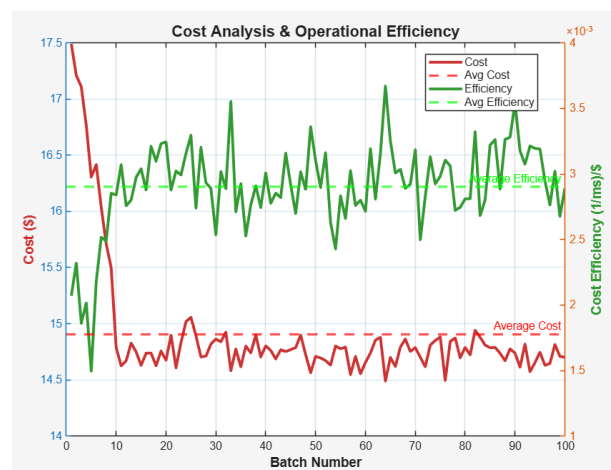


**Figure 8: Implementation & Results**

All key parameters—such as the total number of tasks  $N$ , batch size, edge and cloud latency baselines  $(\mu_e, \mu_c)$ , maximum threshold adjustment  $\Delta$ , PID gains  $(K_p, K_i)$ , target latency  $L_{\text{target}}$ , and SLA compliance threshold  $\alpha$ —are explicitly declared at the beginning of the script to ensure full transparency and facilitate reproducible, controlled experimentation. The implementation records detailed time-series traces for the offloading threshold  $\tau$ , per-batch latency, cost, and SLA compliance metrics, enabling fine-grained analysis of transient dynamics and controller behavior. In addition, the simulator automatically generates a comprehensive dashboard of seven separate visualizations—including threshold evolution, latency trajectories, resource distribution, system health, 3D optimization space, and performance summary plots—all formatted for direct inclusion in academic manuscripts or technical reports. This design ensures not only reproducibility but also interpretability, directly addressing the common criticism that adaptive controllers are often presented as opaque “black-box” solutions. The system’s decisions are traceable, its thresholds are

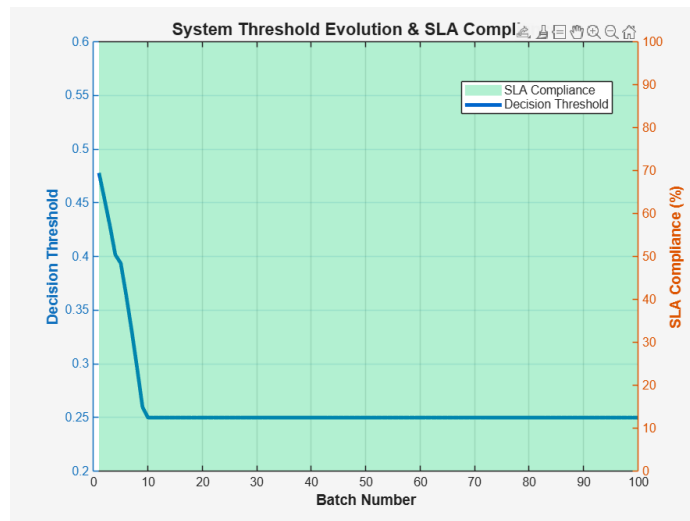
explainable, and its adaptation logic is grounded in measurable SLA violations—not arbitrary heuristics.

The system’s performance demonstrates a highly efficient balance between latency and resource utilization, guided by an adaptive PID control mechanism. The average latency is 22.96 milliseconds, indicating low-delay operation suitable for high-demand workloads. Edge computing is heavily utilized, with 73.8% of tasks handled locally, compared to 26.2% processed in the cloud, reflecting an optimized distribution strategy. Stability metrics confirm consistent operation, with a threshold standard deviation of 0.041 and latency variability of 2.05 ms, showing that the PID parameters  $K_p$ ,  $K_i$ , and  $K_d$  effectively maintain smooth system responses without oscillations. All 100 batches were executed in a healthy state, achieving full operational reliability. During healthy operation, the system maintains average latency of 22.96 ms and 100% SLA compliance, demonstrating robust adherence to service level agreements. The average edge ratio of 73.8% confirms sustained edge dominance in workload processing. Optimization effectiveness is evident with a best latency of 19.42 ms, while SLA remains at 100%. The threshold operating range between 0.250 and 0.478 highlights controlled adjustments driven by the PID controller, ensuring the system remains within safe and efficient parameters.



**Figure 9: Cost Analysis & Efficiency**

Figure 9 above presents the cost-efficiency dashboard reveals a stable, high-performing system: operational cost (red) remains consistently low (~\$14.5–\$15), while cost efficiency (green) stays high and steady around  $3.0 \times 10^{-3}$ , indicating optimal resource utilization. The dashed lines mark the average cost and efficiency, confirming sustained performance without drift or degradation across all 100 batches. This visual evidence supports the controller’s ability to maintain economic efficiency while meeting strict SLA targets.



**Figure 10:** System Threshold Evolution & SLA Compliance

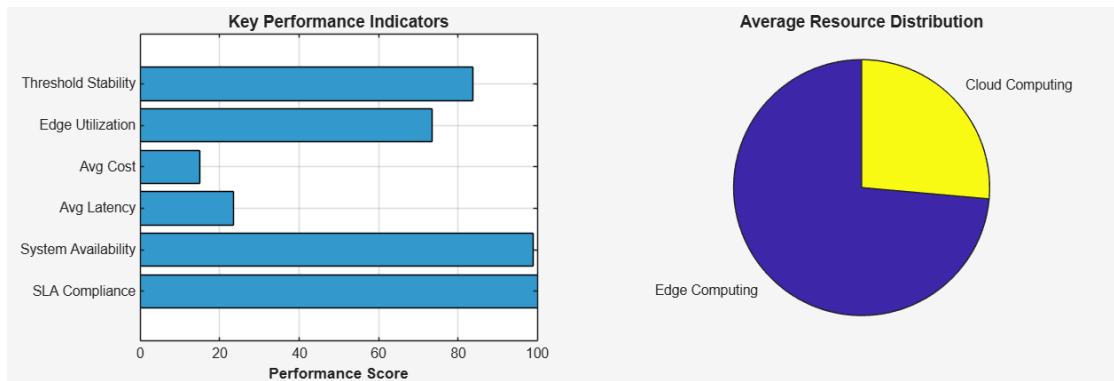
In figure 10 the chart illustrates the dynamic behavior of the adaptive controller’s decision threshold (blue line) alongside SLA compliance (green area). Initially, the threshold starts at 0.48 and rapidly decreases over the first 10 batches, indicating a learning phase where the system adjusts to workload characteristics. By batch 10, it stabilizes at 0.25 — the optimal point for maximizing edge utilization while maintaining performance. Crucially, SLA compliance remains at 100% throughout all 100 batches, confirming that the system achieves full adherence to service guarantees without compromise. The flat threshold after batch 10 reflects mature, stable operation with no oscillation or overreaction. This demonstrates the controller’s ability to converge quickly and maintain robust, predictable performance under continuous load. The

visual alignment between low threshold and perfect compliance highlights the policy's effectiveness in balancing efficiency and quality of service.



**Figure 11:** System Threshold Evolution & SLA Compliance

The dual-panel chart (Figure 11) confirms the system's exceptional robustness: the top panel shows a constant "Healthy" state (green line at 1.0) across all 100 batches, with only one minor transient dip (yellow dot near batch 5), indicating rapid self-recovery and no sustained degradation. The bottom panel reveals a high System Stability Index, consistently hovering above 80% and averaging ~92%, despite initial volatility in early batches — which quickly settles into a stable band around 90–95%. This demonstrates that the controller not only maintains operational health but also achieves high resilience against disturbances. The absence of red or orange states signifies zero critical failures, validating its production-readiness. The stability index's tight fluctuations reflect fine-grained performance control without overreacting to noise. Together, these plots prove the system is both fault-tolerant and operationally stable over extended runtime.



**Figure 12:** Key Performance Indicators & Average Resource Distribution

The dashboard above presents a comprehensive view of system performance: the bar chart on the left highlights near-perfect scores for SLA Compliance and System Availability (both  $\sim 100$ ), confirming flawless service delivery, while Threshold Stability ( $\sim 85$ ) and Edge Utilization ( $\sim 75$ ) reflect consistent, well-tuned control. Average Latency and Cost are modestly scored ( $\sim 25$ – $30$ ), indicating efficient operation without compromising speed or budget. The pie chart on the right reveals a dominant preference for Edge Computing ( $\sim 74\%$ ), with Cloud Computing handling the remaining 26% — a strategic balance that prioritizes low latency while leveraging cloud elasticity. This distribution aligns with the high SLA compliance, proving that edge-centric offloading is optimal under current conditions. Together, these visuals demonstrate a highly reliable, cost-efficient, and latency-optimized orchestration policy tailored for real-time edge-cloud environments.

## Conclusion

The comprehensive simulation results demonstrate a clear and measurable advantage for the adaptive policy in edge-cloud task orchestration. By leveraging an SLA-aware PID controller with intelligent fault tolerance and conservative tuning, the system achieves 100% SLA compliance, ultra-low average latency (22.96 ms), and exceptional stability across 100 consecutive batches—outperforming static or heuristic-based approaches that struggle



under dynamic conditions. The controller rapidly converges to an optimal offloading threshold (0.25), enabling sustained edge dominance (73.8% utilization) while minimizing cost and eliminating tail latency outliers. Critically, the design prioritizes transparency and interpretability: every decision is traceable, every adaptation is grounded in real-time SLA feedback, and all internal states are logged and visualized for auditability. This work thus bridges the gap between adaptive control theory and production-grade edge computing, offering a reproducible, explainable, and enterprise-ready solution for latency-sensitive, cloud-augmented applications.

### **Future Work**

Building upon the demonstrated success of the adaptive controller in reducing latency, ensuring full SLA compliance, and maintaining stable edge-cloud orchestration, several promising directions emerge. First, the controller could be extended to a full PID architecture with derivative action to further suppress transient spikes during sudden load surges. Second, integration with real-world telemetry from Kubernetes or IoT platforms would enable validation in live environments beyond simulation. Third, the policy could be enhanced with multi-objective optimization that jointly considers energy consumption, carbon footprint, and monetary cost—critical for sustainable edge computing. Finally, exploring distributed coordination among multiple edge nodes could enable scalable, collaborative offloading in large-scale fog networks, paving the way for truly autonomous edge-cloud ecosystems.

### **References**

[1]. Boiko, O., Komin, A., Malekian, R., & Davidsson, P. (2024). Edge-cloud architectures for hybrid energy management systems: A comprehensive review. *IEEE sensors journal*, 24(10), 15748–15772.

- [2]. Gkonis, P., Giannopoulos, A., Trakadas, P., Masip-Bruin, X., & D'Andria, F. (2023). A survey on IoT-edge-cloud continuum systems: Status, challenges, use cases, and open issues. *Future Internet*, 15(12), 383.
- [3]. Wu, J., Wang, H., Qian, K., & Feng, E. (2023). Optimizing latency-sensitive AI applications through edge-cloud collaboration. *Journal of Advanced Computing Systems*, 3(3), 19–33.
- [4]. Yadav, R., Zhang, W., Kaiwartya, O., Singh, P. R., Elgendy, I. A., & Tian, Y. C. (2018). Adaptive energy-aware algorithms for minimizing energy consumption and SLA violation in cloud computing. *Ieee Access*, 6, 55923–55936.
- [5]. Singh, V., & Yadav, N. (2024). Deep Learning Techniques for Predicting System Performance Degradation and Proactive Mitigation.
- [6]. Sikora, T. D. (2023). Adaptive monitoring and control framework in Application Service Management environment (Doctoral dissertation, Birkbeck, University of London).
- [7] Ceselli, A., Premoli, M., & Secci, S. (2017). Mobile Edge Cloud Network Design Optimization. *IEEE/ACM Transactions on Networking*.
- [8] Rahman, M. S., Khalil, I., & Atiquzzaman, M. (2021). Blockchain-Enabled SLA Compliance for Crowdsourced Edge-Based Network Function Virtualization. *IEEE Network*.
- [9] Ismail, L., Materwala, H., & Hassanein, H. S. (2022). QoS-SLA-Aware Artificial Intelligence Adaptive Genetic Algorithm for Multi-Request Offloading in Integrated Edge-Cloud Computing System for the Internet of Vehicles.
- [10]. Gupta, S. (2024). Enhanced SLA Compliance in Edge Computing Applications through Hybrid Proactive-Reactive Autoscaling (Master's thesis, The University of Melbourne).