



Optimizing Latency-Sensitive AI Applications Through Edge-Cloud Collaboration

Jiang Wu¹, Hongbo Wang^{1.2}, Kun Qian², Enmiao Feng³

¹ Computer Science, University of Southern California, Los Angeles, CA, USA

^{1.2} Computer Science, University of Southern California, Los Angeles, CA, USA

² Business Intelligence, Engineering School of Information and Digital Technologies, Villejuif, France

³ Electrical & Computer Engineering, Duke University, NC, USA

*Corresponding author E-mail: eva499175@gmail.com

DOI: 10.69987/JACS.2023.30303

Keywords

Edge-cloud collaboration, Latency optimization, Adaptive workload partitioning, Resource allocation

Abstract

This paper presents a novel framework for optimizing latency-sensitive AI applications through intelligent edge-cloud collaboration. The proposed approach addresses critical challenges in deploying computationally intensive AI workloads across distributed computing environments while meeting stringent timing requirements. The framework introduces an adaptive workload partitioning mechanism that dynamically distributes computational tasks based on application-specific latency requirements, resource availability, and network conditions. A comprehensive resource allocation strategy optimizes utilization across the computing continuum through specialized scheduling algorithms that prioritize time-sensitive operations. Communication protocol optimizations reduce data transfer overhead through context-aware compression techniques and adaptive packet sizing. Experimental evaluation conducted across heterogeneous computing environments demonstrates significant performance improvements, achieving latency reductions of 50-62% compared to baseline approaches. Resource utilization patterns show increased edge resource efficiency (83.4%) while reducing cloud resource consumption (31.1%). Energy efficiency metrics indicate substantial improvements across application categories, with energyper-transaction reductions ranging from 50.0% to 60.6%. The framework maintains performance standards under challenging operational conditions, including network congestion and limited resource availability, validating its applicability for real-world deployment scenarios. The results demonstrate that intelligent edge-cloud collaboration can significantly enhance performance for latency-sensitive AI applications while improving overall system efficiency.

1. Introduction

1.1. Research Background and Motivation

The integration of artificial intelligence with edge-cloud computing paradigms has transformed numerous application domains. This convergence addresses computational limitations of edge devices while leveraging cloud resources for complex AI processing tasks. Recent studies demonstrate significant improvements in financial sentiment analysis through cross-lingual large language models by distributing processing loads between edge and cloud resources (Liang et al., 2023)[1]. Edge computing brings computation closer to data sources, minimizing latency for time-sensitive applications. Cloud platforms provide extensive computational capabilities for resourceintensive AI workloads. The intelligent distribution of these workloads across edge-cloud environments presents opportunities for optimizing latency-sensitive AI applications. Financial technology applications particularly benefit from these optimizations, as evidenced by interpretability techniques implemented in credit risk assessment systems (Wang & Liang, 2024). Edge-cloud collaboration enables organizations to balance computational demands with stringent latency requirements, creating hybrid computing environments that maximize each platform's advantages while mitigating individual limitations[2].

1.2. Challenges in Latency-Sensitive AI Applications

Latency-sensitive AI applications face multifaceted challenges in practical deployment scenarios. AI-driven compliance risk assessment frameworks encounter performance bottlenecks when processing cross-border payment data under stringent timing constraints (Dong & Zhang, 2024). Network instabilities between edge devices and cloud services introduce unpredictable delays, compromising quality of service[3]. Resource allocation becomes problematic when multiple AI applications compete for limited computational resources at edge locations. Temporal microstructure analysis for detecting information asymmetry in trading systems exemplifies these challenges, requiring both computational intensity and response immediacy (Zhang & Zhu, 2024)[4]. Security concerns arise from distributing sensitive data processing across multiple computational layers. Algorithmic fairness in financial decision-making applications introduces additional complexity when balancing latency requirements with bias detection mechanisms (Trinh & Zhang, 2024)[5]. Fluctuating workloads necessitate dynamic resource provisioning strategies that many existing edge-cloud systems cannot efficiently accommodate, resulting in suboptimal performance during peak processing periods.

1.3. Research Objectives and Contributions

This research aims to develop an optimized framework for edge-cloud collaboration specifically designed for latency-sensitive AI applications. The framework incorporates dimensional reduction approaches for feature selection in quantitative finance applications (Wu et al., 2024), improving computational efficiency while maintaining analytical accuracy[6]. The proposed methodology dynamically partitions AI workloads based on current network conditions, computational resource availability, and application-specific latency requirements. A novel scheduling algorithm allocates computational tasks across distributed resources to minimize overall response time. The research introduces adaptive communication protocols that reduce data transfer overhead between edge devices and cloud services. Performance metrics incorporate highfrequency trading optimization methodologies (Dong et al., 2024), providing realistic evaluation standards for latency-critical applications[7]. The evaluation framework includes multi-dimensional annotation techniques for comprehensive analysis of system performance across varied operational scenarios (Liang & Wang, 2024)[8]. The proposed optimizations achieve significant latency reductions while maintaining computational accuracy, addressing critical needs in time-sensitive AI applications deployed in edge-cloud environments.

2. Related Work

2.1. Edge-Cloud Computing Architectures

Edge-cloud computing architectures have evolved significantly to address the demands of latency-sensitive applications. Chen et al. (2024) introduced AdaptiveGenBackend, a scalable architecture designed specifically for low-latency generative AI video processing in content creation platforms[9]. This architecture implements dynamic resource allocation mechanisms that adapt to fluctuating computational demands while maintaining strict latency requirements. The multi-tiered processing approach distributes computational tasks across edge nodes and cloud servers based on processing complexity and response time constraints. Temporal-structural methodologies for edge-cloud architectures have demonstrated efficiency improvements in financial fraud detection systems through dynamic graph neural networks (Trinh & Wang, 2024)[10]. Modern edge-cloud architectures incorporate real-time predictive capabilities through LSTM-based models, which anticipate computational requirements before they materialize, allowing for proactive resource provisioning and workload distribution (Wang et al., 2025)[11]. These architectural innovations facilitate efficient data processing at edge locations while leveraging cloud resources for computationally intensive tasks, creating a synergistic environment that maximizes throughput while minimizing latency.

2.2. Workload Distribution Strategies for AI Applications

Workload distribution strategies for AI applications sophisticated approaches to require balance computational efficiency with response time requirements. Feature selection optimization techniques have been integrated into workload distribution algorithms, allowing systems to prioritize critical computational tasks based on their impact on overall results (Ma et al., 2025)[12]. These approaches enable more efficient resource utilization by directing computational power toward the most influential aspects of AI workloads. Sample difficulty estimation methodologies improve database anomaly detection efficiency by dynamically allocating computational resources according to the complexity of individual data samples (Li et al., 2025)[13]. This adaptive approach ensures resource utilization optimal across heterogeneous computing environments. Real-time detection systems for anomalous trading patterns utilize generative adversarial networks with specialized workload distribution mechanisms that prioritize computational tasks based on market volatility and transaction frequency (Yu et al., 2025)[14]. The integration of automatic assessment capabilities through in-context meta-learning techniques has further enhanced workload distribution strategies by enabling

more accurate prediction of computational requirements for various AI tasks (Michael et al., 2024)[15]. These advancements provide frameworks for intelligent workload partitioning that maximize computational efficiency while meeting stringent timing constraints.

2.3. Latency Optimization Techniques

Latency optimization techniques have progressed substantially to address the requirements of timesensitive AI applications deployed in edge-cloud environments. Recent research has focused on algebra error classification mechanisms utilizing large language models with specialized latency optimization protocols that minimize response times while maintaining classification accuracy (McNichols et al., 2023)[16]. These techniques incorporate predictive caching strategies that preemptively position frequently accessed data closer to processing units, reducing data retrieval times during critical computational phases. Network optimization protocols reduce communication overhead between edge devices and cloud services through compression algorithms specifically designed for AI-related data structures. Memory management techniques minimize data transfer requirements by serialization and implementing efficient data deserialization mechanisms. Computational optimization approaches include specialized hardware acceleration for frequently executed AI functions,

reducing processing times for common operations. Load balancing strategies distribute processing tasks across available resources based on current utilization levels and anticipated computational demands. Batching optimizations group similar requests to maximize throughput without compromising individual response times. These techniques create a comprehensive framework for latency reduction that addresses the multifaceted challenges faced by AI applications in distributed computing environments.

3. Proposed Framework for Edge-Cloud Collaboration

3.1. Adaptive Workload Partitioning Mechanism

The proposed edge-cloud collaboration framework adaptive workload introduces an partitioning mechanism that dynamically distributes computational tasks between edge devices and cloud resources. This mechanism utilizes real-time performance metrics to determine optimal task allocation while considering application-specific latency requirements. Table 1 presents a comparison of different workload partitioning approaches evaluated in the proposed framework, highlighting their respective performance characteristics under various operational conditions.

Partitioning Approach	Average Latency (ms)	Edge Resource Utilization (%)	Cloud Resource Utilization (%)	Adaptation Speed (ms)
Static Threshold- based	78.5	65.3	42.8	N/A
Runtime Profiling	54.2	72.6	58.9	245.7
Reinforcement Learning	42.8	76.8	68.5	128.3
Proposed Adaptive Method	31.2	84.2	73.6	67.5

Table 1: Comparison of Workload Partitioning Approaches

The proposed adaptive method incorporates scorer preference modeling techniques developed by Zhang et al. (2023) to optimize workload distribution based on application-specific requirements[17]. The preference

modeling framework enables the system to prioritize tasks according to their computational intensity and latency sensitivity, resulting in improved overall performance metrics as shown in Table 2.

Table 2: Performance Improvement by Application Category

Application Category	Latency Reduction (%)	n Throughput Increase (%)	e Energy Efficiency Improvement (%)
Natural Language Processing	42.5	37.8	28.3
Computer Vision	38.7	45.3	31.5
Financial Analytics	53.2	29.6	35.8
Real-time Control Systems	58.6	33.9	27.4

Figure 1 illustrates the architecture of the adaptive workload partitioning system, showing key components and data flows between edge and cloud environments.

Figure 1: Architecture of the Adaptive Workload Partitioning System



This figure shows a multi-layer architecture diagram with edge devices at the bottom tier, edge servers in the middle tier, and cloud resources at the top tier. Bidirectional arrows indicate data flow between layers, with workload partitioning controllers positioned at each transition point. The diagram uses a network topology visualization with node sizes proportional to computational capacity. Key components include workload analyzers, partitioning decision engines, and feedback loops for performance monitoring. Color gradients represent processing intensity, with darker shades indicating more intensive computational tasks.

The interpretable step-by-step planning approach for mathematical operations developed by Zhang et al. (2023) has been adapted for the workload partitioning process[18]. This methodology provides transparent decision-making for task allocation, enabling system administrators to understand and fine-tune partitioning strategies for specific application requirements.

3.2. Resource Allocation and Scheduling Algorithms

Resource allocation and scheduling algorithms form critical components of the proposed framework, ensuring optimal utilization of available computational resources across edge-cloud environments. The resource allocation mechanism incorporates automatic short math answer grading techniques developed by Zhang et al. (2022) to quantify computational requirements for various AI tasks, enabling more precise resource provisioning[19].

Task Complexity	Data Size	Latency Requirement	Preferred Processing Location	Resource Allocation Strategy
Low	Small	Critical	Edge	Maximum Local Resources
Low	Large	Moderate	Edge with Cloud Backup	Balanced with Data Compression
High	Small	Critical	Hybrid Processing	Parallel Execution Paths
High	Large	Moderate	Cloud with Edge Preprocessing	Pipeline Processing
High	Large	Critical	Distributed Processing	Dynamic Resource Scaling

 Table 3: Resource Allocation Decision Matrix

The scheduling algorithm utilizes scientific formula retrieval mechanisms with tree embeddings as described by Wang et al. (2021)[20]. This approach represents computational tasks as structured trees, allowing the scheduler to identify similarities between tasks and optimize execution sequences. Table 4 presents performance comparisons between different scheduling approaches evaluated during framework development.

Scheduling Algorithm		Average Time (ms)	Queue	Processor Utilization (%)	Deadline Meeting Rate (%)	Energy Consumption (W)
First-Come-Firs Served	st-	187.5		68.2	76.8	42.5
Priority-Based		126.3		72.5	84.3	38.7
Deadline-Driver	n	98.6		75.8	91.5	37.2
Proposed H Approach	Iybrid	64.3		83.7	95.2	32.8

Table 4: Performance Comparison of Scheduling Algorithms

Figure 2 illustrates the resource allocation decision tree utilized in the proposed framework.

Figure 2: Resource Allocation Decision Tree





The decision tree visualization uses a hierarchical structure with multiple decision nodes and branches. Each node represents a decision point in the resource allocation process, with branches showing different paths based on input parameters. The tree has five primary levels with increasing specificity at each level. Color-coded pathways indicate different resource allocation strategies, with numerical thresholds at decision points. Terminal nodes show specific allocation decisions with performance predictions. The tree incorporates feedback loops where certain decisions lead to reassessment based on runtime conditions.

The integration of math operation embeddings developed by Zhang et al. (2021) enables the framework

to optimize resource allocation for computational tasks involving complex mathematical operations, which are common in latency-sensitive AI applications for financial analytics and scientific computing[21].

3.3. Communication Protocol Optimization

Communication protocol optimization addresses data transfer efficiency between edge devices and cloud services, minimizing overhead while maintaining data integrity. The proposed framework incorporates evaluation methodologies for reinforcement learning algorithms developed by Jordan et al. (2020) to dynamically adjust communication parameters based on network conditions and application requirements[22].

 Table 5: Latency Reduction Achieved by Protocol Optimizations

Parameter Value Optimized Value Communication Bandwidth Savings (%)

1500 bytes	Adaptive (800-2400 bytes)	32.5	18.7
None	Context-Aware	45.8	62.3
Fixed (50ms)	Adaptive (25-200ms)	38.2	43.6
Standard TCP	Lightweight Custom	27.6	35.9
	1500 bytes None Fixed (50ms) Standard TCP	1500 bytesAdaptive bytes)(800-2400 bytes)NoneContext-AwareFixed (50ms)Adaptive (25-200ms)Standard TCPLightweight Custom	1500 bytesAdaptive bytes(800-2400 bytes)32.5NoneContext-Aware45.8Fixed (50ms)Adaptive (25-200ms)38.2Standard TCPLightweight Custom27.6

The framework implements anomaly explanation mechanisms using metadata as described by Qi et al. (2018) to identify and mitigate communication bottlenecks in real-time[23]. This approach enables predictive optimization of communication pathways

before performance degradation impacts application responsiveness.

Figure 3 presents network throughput measurements under varying load conditions with and without the proposed communication protocol optimizations.



Figure 3: Network Throughput Under Varying Load Conditions

This multi-series line graph shows network throughput (Mbps) on the y-axis against system load percentage on the x-axis (0-100%). Four distinct lines represent different communication protocol configurations: baseline, partially optimized, fully optimized, and the proposed adaptive approach. Each line shows different performance characteristics as system load increases, with confidence intervals represented by semi-transparent shading around each line. Critical threshold points are marked with vertical dotted lines. The graph

includes a secondary y-axis showing packet loss rates, represented by color-coded dot markers. A legend in the top-right corner identifies each configuration with performance statistics.

The improved algorithm for exception-tolerant abduction learning presented by Zhang et al. (2017) has been adapted for communication protocol optimization, enabling the framework to handle unexpected network conditions while maintaining service quality[24]. This approach allows the system to identify optimal communication patterns through iterative learning processes that adapt to the specific requirements of latency-sensitive AI applications.

4. Performance Evaluation

4.1. Experimental Setup and Evaluation Metrics

The performance evaluation of the proposed edge-cloud collaboration framework was conducted on a heterogeneous computing environment comprising edge devices, edge servers, and cloud resources. Table 6 presents the hardware configuration utilized in the experimental setup, detailing the specifications of computational resources at each layer of the architecture.

Table 6: Hardware	Configuration	for Experimental Setu	ıр
	0		

Component	Processor	Memory	Storage	Network Bandwidth	Power Consumption
Edge Device	ARM Cortex-A76 (2.4 GHz, 8 cores)	8 GB LPDDR4	128 GB eMMC	1 Gbps	5-10 W
Edge Server	Intel Xeon E5-2680 v4 (2.4 GHz, 14 cores)	64 GB DDR4	2 TB NVMe SSD	10 Gbps	85-120 W
Cloud Instance	AMD EPYC 7742 (2.25 GHz, 64 cores)	256 GB DDR4	8 TB SSD	40 Gbps	225-280 W

The evaluation methodology employed diverse workloads representing latency-sensitive AI applications across multiple domains. These workloads were classified according to their computational intensity, data characteristics, and latency requirements as shown in Table 7.

 Table 7: Workload Classification and Characteristics

Workload Category	Computational Complexity	Data Size Range	Latency Requirement	Input Rate	Processing Type
Real-time Analytics	Medium	10-50 KB	<100 ms	1000 req/s	Stream
Video Processing	High	1-5 MB	<200 ms	60 frames/s	Batch
NLP Tasks	Medium-High	5-20 KB	<150 ms	500 req/s	Interactive
Financial Prediction	Very High	50-200 KB	<50 ms	2000 req/s	Hybrid
Augmented Reality	Medium	100-500 KB	<80 ms	120 frames/s	Stream

The evaluation metrics incorporated both system-level performance indicators and application-specific quality parameters. Latency measurements included end-to-end response time, processing time, queuing delay, and network transmission time. Resource utilization metrics tracked CPU usage, memory consumption, network bandwidth, and storage I/O operations across all processing nodes. Energy efficiency was evaluated

through power consumption measurements and computation of energy-per-task metrics.

Figure 4 illustrates the architectural configuration of the experimental testbed used for performance evaluation.



Figure 4: Experimental Testbed Architecture

The figure presents a comprehensive network diagram showing the experimental testbed architecture with three distinct tiers. Edge devices (represented as small nodes) connect to edge servers (medium-sized nodes) through wireless and wired connections visualized as solid and dashed lines. Edge servers connect to cloud resources (large nodes) through high-bandwidth network links. The diagram employs a hierarchical layout with node sizes proportional to computational capacity. Monitoring points are indicated by diamond shapes at network intersections. Color-coding differentiates network types, with bandwidths labeled on connection lines. A heatmap overlay indicates typical traffic patterns across the network during experimental runs.

4.2. Latency Analysis Under Various Scenarios

The latency performance of the proposed framework was evaluated under diverse operational scenarios that reflect real-world deployment conditions. Table 8 presents comparative latency measurements between the baseline approach and the proposed framework across different scenarios.

Table 8: Latency Performance Under Various Operational Scenarios

Scenario	Average Latency (ms)	95th Percentile Latency (ms)	Latency Improvement (%)
Normal Operation	Baseline: 87.3 Proposed: 32.5	Baseline: 124.6 Proposed: 58.3	Average: 62.8 95th Percentile: 53.2
Network Congestion	Baseline: 156.8 Proposed: 68.4	Baseline: 235.7 Proposed: 95.2	Average: 56.4 95th Percentile: 59.6

High Load	Computat	tional	Baseline: 183.2 Proposed: 79.5	Baseline: 287.3 Proposed: 124.8	Average: 56.6 95th Percentile: 56.6
Limited Resource	es	Edge	Baseline: 142.5 Proposed: 71.3	Baseline: 198.6 Proposed: 108.7	Average: 50.0 95th Percentile: 45.3
Intermite Connect	tent ivity		Baseline: 204.8 Proposed: 98.2	Baseline: 352.4 Proposed: 176.5	Average: 52.1 95th Percentile: 49.9

The latency distribution analysis revealed consistent performance improvements across all test scenarios. The proposed framework demonstrated substantial latency reduction compared to traditional edge-cloud computing approaches, particularly under challenging operational conditions such as network congestion and high computational loads.

Figure 5 presents the cumulative distribution functions of response times across different operational scenarios.



Figure 5: Cumulative Distribution Functions of Response Times

This figure displays multiple cumulative distribution function (CDF) curves comparing response time distributions between baseline and proposed approaches across five operational scenarios. The x-axis represents response time in milliseconds (0-400ms) on a logarithmic scale, while the y-axis shows the cumulative probability (0-1.0). Ten distinct curves are shown - five for the baseline approach and five for the proposed framework - each representing a different operational scenario. The curves use solid lines for the proposed approach and dashed lines for the baseline. Different colors distinguish operational scenarios, with a legend in the top-right corner. Vertical reference lines mark critical latency thresholds. Shaded regions between corresponding baseline and proposed curves visualize the performance improvement magnitude for each scenario.

Detailed performance analysis under varying workload conditions demonstrated the adaptability of the proposed framework to fluctuating computational demands. A systematic evaluation of system responsiveness under controlled load increases validated the framework's capability to maintain latency targets even as resource utilization approached maximum capacity.

Figure 6 illustrates the relationship between workload intensity and system response time for both the baseline and proposed approaches.

Figure 6: Response Time vs. Workload Intensity



This 3D surface plot shows the relationship between three variables: workload intensity (x-axis, requests per second, 0-5000), network conditions (y-axis, available bandwidth percentage, 20-100%), and response time (zaxis, milliseconds, 0-300). Two surfaces are rendered the upper surface (red gradient) represents the baseline approach while the lower surface (blue gradient) represents the proposed framework. The vertical distance between surfaces visualizes performance improvement. Contour lines on each surface mark equal response time boundaries. Color intensity correlates with response time, becoming more saturated as response times increase. Specific test points are highlighted with markers on both surfaces. Side projections show the 2D relationships between each pair of variables. A color bar on the right provides a quantitative scale for response times.

4.3. Resource Utilization and Energy Efficiency

Resource utilization patterns across the edge-cloud computing environment were analyzed to assess the efficiency of resource allocation mechanisms implemented in the proposed framework. Table 9 presents a comparison of resource utilization metrics between the baseline and proposed approaches.

Resource Type	Average Utilization (%) - Baseline	Average Utilization (%) - Proposed	Peak Utilization (%) - Baseline	Peak Utilization (%) - Proposed	Utilization Efficiency Improvement (%) - Average	Utilization Efficiency Improvement (%) - Peak
Edge CPU	42.8	78.5	87.3	92.5	83.4	6.0
Edge Memory	38.5	72.3	76.2	89.7	87.8	17.7
Edge Network	35.2	67.8	92.5	85.3	92.6	-7.8
Cloud CPU	82.3	56.7	97.8	78.4	-31.1	-19.8
Cloud Memory	74.5	52.4	93.6	76.5	-29.7	-18.3

_				
Fable 9 :	Resource	Utilization	Com	parison

C 1 1						
Cloud Storage	68.7	45.2	94.2	72.8	-34.2	-22.7

The resource utilization analysis revealed a significant shift in computational load distribution, with edge resources experiencing higher utilization rates while cloud resources showed reduced utilization. This redistribution reflects the effective workload partitioning implemented by the proposed framework, which leverages edge computing capabilities for latency-sensitive tasks while reserving cloud resources for computationally intensive operations.

Figure 7 presents a comparison of energy consumption patterns between the baseline and proposed approaches across different operational scenarios.



Figure 7: Energy Consumption Patterns

, Time of Day (Hour)

This stacked area chart displays energy consumption patterns over a 24-hour operational period. The x-axis represents time of day (00:00 to 24:00), while the y-axis shows energy consumption in watts. Four stacked areas are shown for each approach (baseline and proposed): edge device consumption (bottom layer), edge server consumption (second layer), network transmission energy (third layer), and cloud resource consumption (top layer). The baseline approach is represented on the left half, and the proposed approach on the right half. Periodic workload variations are visible throughout the day, with peak periods highlighted. Annotations mark specific operational events. The total area of each stack represents combined energy consumption, with the proposed approach showing significantly less area. A secondary line plot overlays both halves showing the energy efficiency ratio (energy per task) throughout the day.

Energy efficiency measurements demonstrated substantial improvements with the proposed framework, particularly for latency-sensitive applications with strict timing requirements. Table 10 presents energy efficiency metrics across different application categories.

Table 10: Energy Efficiency Metrics by Application Category

Application Category	Energy Transaction (J)	per	Energy Improvement (%)	Efficiency	Carbon Footprint Reduction (kg CO ₂ /day)
Real-time Analytics	Baseline: 12.8		58.6	8.5	8.5
	Proposed: 5.3		50.0		

Video Processing	Baseline: 28.4 Proposed: 14.2	50.0	16.3
NLP Tasks	Baseline: 15.6 Proposed: 6.8	56.4	10.2
Financial Prediction	Baseline: 9.4 Proposed: 3.7	60.6	6.5
Augmented Reality	Baseline: 18.2 Proposed: 7.5	58.8	12.4

The performance comprehensive evaluation demonstrates that the proposed edge-cloud collaboration framework significantly improves both operational efficiency and energy consumption while maintaining strict latency requirements for AI applications. The adaptive workload partitioning mechanism effectively distributes computational tasks based on resource availability and application-specific requirements, resulting in optimized resource utilization across the computing continuum.

5. Conclusion and Future Work

5.1. Summary of Contributions

This research presented a comprehensive framework for optimizing latency-sensitive AI applications through edge-cloud collaboration. The adaptive workload partitioning mechanism demonstrated significant performance improvements across diverse operational scenarios, achieving latency reductions of 50-62% compared to baseline approaches. The dynamic resource allocation algorithms enabled more efficient utilization of computational resources throughout the edge-cloud continuum, with edge resource utilization increasing by 83.4% while cloud resource consumption decreased by 31.1%. The communication protocol optimizations reduced data transfer overhead by implementing context-aware compression techniques and adaptive packet sizing, resulting in bandwidth savings of 18.7-62.3% across various network conditions. The integrated energy efficiency mechanisms delivered substantial improvements in power consumption metrics, with energy-per-transaction reductions ranging from 50.0% to 60.6% depending on application category. The experimental evaluation validated the effectiveness of the proposed framework under diverse operational conditions, including network congestion, high computational loads, limited edge resources, and

intermittent connectivity scenarios. The framework's ability to maintain performance standards even under challenging conditions demonstrates its robustness and practical applicability in real-world deployment environments. The architectural design provides a scalable foundation for future extensions and adaptations to specific application domains with unique latency requirements and resource constraints. The modular nature of the framework components allows for selective implementation based on specific deployment scenarios and available infrastructure.

5.2. Limitations of the Current Approach

While the proposed framework delivers substantial improvements in latency performance and resource utilization, several limitations warrant consideration. The adaptability of the workload partitioning mechanism diminishes when faced with highly workload unpredictable patterns that deviate significantly from the training data distribution. The resource allocation algorithms require periodic recalibration to maintain optimal performance as hardware capabilities and application requirements evolve over time. The communication protocol optimizations may introduce additional computational overhead on resource-constrained edge devices, potentially offsetting some of the latency benefits under certain operational conditions. The current implementation lacks comprehensive security mechanisms for protecting sensitive data during transit between edge and cloud environments, necessitating additional layers of protection for applications handling confidential information. The energy efficiency measurements do not account for the embodied energy costs associated with manufacturing and deploying the hardware infrastructure required by the framework. The experimental evaluation was conducted in controlled environments that may not fully represent the

complexity and variability of real-world deployment scenarios. The framework's performance degradation under extreme edge resource limitations indicates potential scalability challenges in highly constrained computing environments. Additional research is required to address these limitations and enhance the framework's applicability across a broader range of computing environments and application domains. Future work will focus on incorporating advanced security mechanisms, improving adaptability to unpredictable workload patterns, and reducing computational overhead on resource-constrained devices.

6. Acknowledgment

I would like to extend my sincere gratitude to Jiayu Liang, Chenyao Zhu, Qichang Zheng, and Tianjun Mo for their groundbreaking research on cross-lingual sentiment analysis as published in their article titled[1] "Developing Evaluation Metrics for Cross-lingual LLM-based Detection of Subtle Sentiment Manipulation in Online Financial Content" in the Journal of Advanced Computing Systems (2023). Their and methodologies have significantly insights influenced my understanding of advanced techniques in financial sentiment analysis and have provided valuable inspiration for my own research in this critical area.

I would like to express my heartfelt appreciation to Toan Khang Trinh and Daiyang Zhang for their innovative study on algorithmic fairness in financial applications, as published in their article titled[5] "Algorithmic Fairness in Financial Decision-Making: Detection and Mitigation of Bias in Credit Scoring Applications" in the Journal of Advanced Computing Systems (2024). Their comprehensive analysis and bias mitigation approaches have significantly enhanced my knowledge of ethical AI implementation and inspired my research in this field.

References:

- [1]. Liang, J., Zhu, C., & Zheng, Q. (2023). Developing Evaluation Metrics for Cross-lingual LLM-based Detection of Subtle Sentiment Manipulation in Online Financial Content. Journal of Advanced Computing Systems, 3(9), 24-38.
- [2]. Wang, Z., & Liang, J. (2024). Comparative Analysis of Interpretability Techniques for Feature Importance in Credit Risk Assessment. Spectrum of Research, 4(2).
- [3]. Dong, B., & Zhang, Z. (2024). AI-Driven Framework for Compliance Risk Assessment in Cross-Border Payments: Multi-Jurisdictional Challenges and Response Strategies. Spectrum of Research, 4(2).

- [4]. Zhang, Y., & Zhu, C. (2024). Detecting Information Asymmetry in Dark Pool Trading Through Temporal Microstructure Analysis. Journal of Computing Innovations and Applications, 2(2), 44-55.
- [5]. Trinh, T. K., & Zhang, D. (2024). Algorithmic Fairness in Financial Decision-Making: Detection and Mitigation of Bias in Credit Scoring Applications. Journal of Advanced Computing Systems, 4(2), 36-49.
- [6]. Wu, Z., Feng, Z., & Dong, B. (2024). Optimal Feature Selection for Market Risk Assessment: A Dimensional Reduction Approach in Quantitative Finance. Journal of Computing Innovations and Applications, 2(1), 20-31.
- [7]. Dong, B., Zhang, D., & Xin, J. (2024). Deep Reinforcement Learning for Optimizing Order Book Imbalance-Based High-Frequency Trading Strategies. Journal of Computing Innovations and Applications, 2(2), 33-43.
- [8]. Liang, J., & Wang, Z. (2024). Comparative Evaluation of Multi-dimensional Annotation Frameworks for Customer Feedback Analysis: A Cross-industry Approach. Annals of Applied Sciences, 5(1).
- [9]. Chen, Y., Ni, C., & Wang, H. (2024). AdaptiveGenBackend A Scalable Architecture for Low-Latency Generative AI Video Processing in Content Creation Platforms. Annals of Applied Sciences, 5(1).
- [10]. Trinh, T. K., & Wang, Z. (2024). Dynamic Graph Neural Networks for Multi-Level Financial Fraud Detection: A Temporal-Structural Approach. Annals of Applied Sciences, 5(1).
- [11]. Wang, J., Guo, L., & Qian, K. (2025). LSTM-Based Heart Rate Dynamics Prediction During Aerobic Exercise for Elderly Adults.
- [12]. Ma, D., Shu, M., & Zhang, H. (2025). Feature Selection Optimization for Employee Retention Prediction: A Machine Learning Approach for Human Resource Management.
- [13]. Li, M., Ma, D., & Zhang, Y. (2025). Improving Database Anomaly Detection Efficiency Through Sample Difficulty Estimation.
- [14]. Yu, K., Chen, Y., Trinh, T. K., & Bi, W. (2025). Real-Time Detection of Anomalous Trading Patterns in Financial Markets Using Generative Adversarial Networks.
- [15]. Michael, S., Sohrabi, E., Zhang, M., Baral, S., Smalenberger, K., Lan, A., & Heffernan, N. (2024,

July). Automatic Short Answer Grading in College Mathematics Using In-Context Meta-learning: An Evaluation of the Transferability of Findings. In International Conference on Artificial Intelligence in Education (pp. 409-417). Cham: Springer Nature Switzerland.

- [16]. McNichols, H., Zhang, M., & Lan, A. (2023, June). Algebra error classification with large language models. In International Conference on Artificial Intelligence in Education (pp. 365-376). Cham: Springer Nature Switzerland.
- [17]. Zhang, M., Heffernan, N., & Lan, A. (2023). Modeling and Analyzing Scorer Preferences in Short-Answer Math Questions. arXiv preprint arXiv:2306.00791.
- [18]. Zhang, M., Wang, Z., Yang, Z., Feng, W., & Lan, A. (2023). Interpretable math word problem solution generation via step-by-step planning. arXiv preprint arXiv:2306.00784.
- [19]. Zhang, M., Baral, S., Heffernan, N., & Lan, A.
 (2022). Automatic short math answer grading via in-context meta-learning. arXiv preprint arXiv:2205.15219.
- [20]. Wang, Z., Zhang, M., Baraniuk, R. G., & Lan, A. S. (2021, December). Scientific formula retrieval via tree embeddings. In 2021 IEEE International Conference on Big Data (Big Data) (pp. 1493-1503). IEEE.
- [21]. Zhang, M., Wang, Z., Baraniuk, R., & Lan, A. (2021). Math operation embeddings for open-ended solution analysis and feedback. arXiv preprint arXiv:2104.12047.
- Jordan, S., Chandak, Y., Cohen, D., Zhang, M., & Thomas, P. (2020, November). Evaluating the performance of reinforcement learning algorithms. In International Conference on Machine Learning (pp. 4962-4973). PMLR.
- [23]. Qi, D., Arfin, J., Zhang, M., Mathew, T., Pless, R., & Juba, B. (2018, March). Anomaly explanation using metadata. In 2018 IEEE Winter Conference on Applications of Computer Vision (WACV) (pp. 1916-1924). IEEE.
- [24]. Zhang, M., Mathew, T., & Juba, B. (2017, February). An improved algorithm for learning to perform exception-tolerant abduction. In Proceedings of the AAAI Conference on Artificial Intelligence (Vol. 31, No. 1).