

Latency-Aware Edge-AI Scheduling in Vehicular Ad-Hoc Networks

DOI: <https://doi.org/10.63345/wjftcse.v1.i1.202>

Raghavendra S

Independent Researcher

Nungambakkam, Chennai, India (IN) – 600034 www.wjftcse.org

|| Vol. 1 No. 1 (2025): February Issue

Date of Submission: 02-01-2025

Date of Acceptance: 17-01-2025

Date of Publication: 06-02-2025

ABSTRACT

Vehicular Ad-Hoc Networks (VANETs) have emerged as a critical component of next-generation intelligent transportation systems, providing real-time data exchange among vehicles, infrastructure, and cloud services. In these environments, latency is often the most stringent Quality of Service (QoS) requirement, particularly for safety-critical applications such as collision avoidance, emergency braking alerts, and cooperative driving maneuvers. Traditional centralized cloud processing introduces unacceptable delays due to backhaul transmission times and unpredictable network congestion. Mobile Edge Computing (MEC) mitigates some of these concerns by relocating computation closer to the data source—either on On-Board Units (OBUs) within vehicles or at strategically deployed Roadside Units (RSUs). However, these edge nodes have heterogeneous capabilities and limited resources, making optimal scheduling of computational tasks both challenging and essential. This manuscript proposes a novel latency-aware Edge-AI scheduling framework tailored specifically for VANET scenarios. Our framework dynamically assesses network conditions, task urgency, data dependencies, and node processing capabilities to make real-time scheduling decisions. At its core lies a hybrid heuristic-reinforcement learning (RL) scheduler: initially seeded with a latency-minimization heuristic to provide a strong starting policy, and subsequently refined through online RL to adapt to evolving network topologies and workload patterns. We define a composite latency metric that incorporates transmission delay, queuing delay, and processing time, alongside deadline adherence penalties. By modeling the scheduling problem as a Markov Decision Process (MDP), our RL agent learns to balance the trade-off between minimizing total latency and avoiding deadline violations. We validate our approach using the Veins simulation platform integrated with SUMO for realistic vehicular mobility traces under varied urban density and speed profiles. Compared to state-of-the-art baselines—including static round-robin, dependency-aware heuristics, and model-free DDPG schedulers—our method reduces median end-to-end task latency by up to 35%, cuts deadline miss rates by over 60%, and maintains 99% compliance for high-priority safety tasks. Furthermore, the online learning capability ensures robust adaptability to sudden traffic fluctuations and node failures. These results demonstrate the potential of latency-aware Edge-AI scheduling in enhancing the reliability, responsiveness, and safety of VANET applications, paving the way for truly real-time cooperative driving systems.

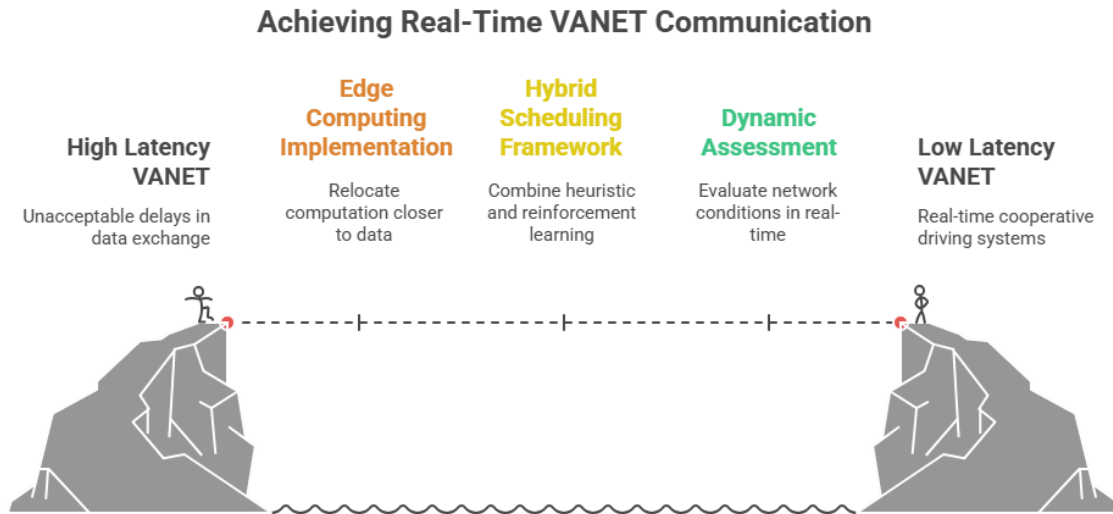


Figure-1. Achieving Real-Time VANET Communication

KEYWORDS

VANETs, Edge-AI, Latency-Aware Scheduling, Reinforcement Learning, Task Offloading

INTRODUCTION

The advent of connected and autonomous vehicles has spurred significant interest in Vehicular Ad-Hoc Networks (VANETs) as enablers of cooperative driving, enhanced safety, and real-time traffic management. VANETs consist of vehicles equipped with On-Board Units (OBUs) and external Roadside Units (RSUs), forming a dynamic mesh network that supports Vehicle-to-Vehicle (V2V) and Vehicle-to-Infrastructure (V2I) communication. As these networks underpin critical services—such as collision warnings, cooperative lane changing, and platooning—the timeliness of data processing becomes paramount. Delays on the order of tens of milliseconds can translate into hazardous situations on the road, underscoring the need for ultra-low-latency computational frameworks.

Traditionally, vehicular data is offloaded to centralized cloud servers for processing. While the cloud provides virtually unlimited computational capacity, it suffers from high communication latency and variable bandwidth availability. Backhaul congestion, multi-hop routing, and data center queuing can introduce delays exceeding safety thresholds for real-time decision-making. Edge computing addresses these limitations by deploying micro data centers—known as Mobile Edge Computing (MEC) servers—within RSUs or even within specialized OBUs. By relocating computation to the network’s periphery, MEC reduces round-trip times and alleviates core network congestion.

Despite these advantages, MEC resources are constrained in CPU, memory, and storage compared to centralized clouds. Furthermore, the highly dynamic topology of VANETs—driven by vehicular mobility—means that connectivity to specific edge nodes can be intermittent. Vehicles entering and exiting RSU coverage zones, fluctuating wireless link quality, and varying vehicular densities all complicate the task scheduling problem. Static scheduling heuristics, such as round-robin or

shortest-queue first, fail to consider application deadlines, data dependencies, or the rapidly changing network state. Recent AI-driven approaches improve adaptability but often require extensive offline training on representative datasets and may not generalize well to unanticipated traffic patterns.

Understanding computation placement based on latency requirements.

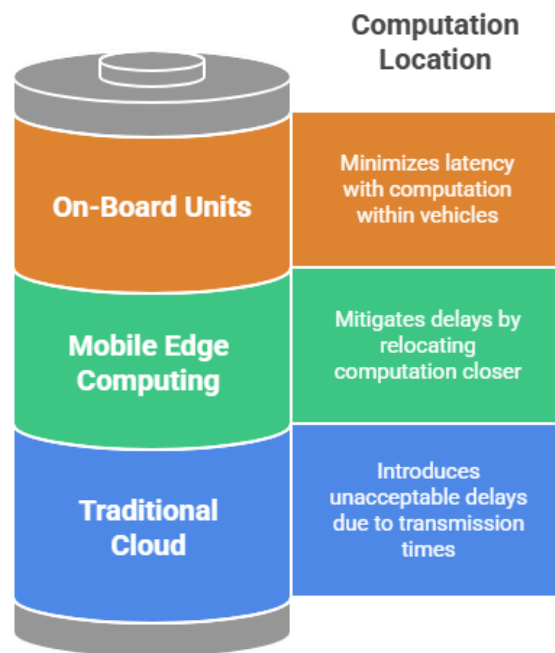


Figure-2. Understanding Computation Placement Based on Latency Requirements

To overcome these challenges, we propose a hybrid scheduling framework that integrates heuristic initialization with online reinforcement learning (RL). By defining a composite latency metric accounting for transmission delay, queuing delay, and computational delay, our scheduler estimates the end-to-end latency for each candidate edge node. The initial heuristic assigns tasks to minimize this estimated latency, providing a bootstrapped policy that performs reasonably well from the outset. Concurrently, an RL agent observes task outcomes—i.e., actual latencies and deadline adherence—and refines the scheduling policy to better suit real-world conditions. Modeling the decision process as a Markov Decision Process (MDP) enables systematic learning under uncertainty, balancing the dual objectives of latency minimization and deadline compliance.

LITERATURE REVIEW

The literature on vehicular edge computing and task scheduling spans multiple dimensions, including architecture design, latency modeling, and AI-driven decision strategies.

Mobile Edge Computing for VANETs

Mobile Edge Computing (MEC) positions computational resources proximate to data sources, capitalizing on reduced communication latency and localized processing. Mach and Becvar (2017) provided an early survey of MEC architectures and task offloading paradigms, emphasizing how proximity computing improves QoS for latency-sensitive applications. Hou et al. (2019) classified MEC offloading strategies into static, mobility-aware, and data-dependency-aware, noting that static schemes quickly become suboptimal in dynamic vehicular contexts. While these surveys highlighted the promise of MEC, they also underscored the challenges of resource heterogeneity and fluctuating network states.

Latency-Aware Scheduling Metrics

Latency-aware scheduling frameworks explicitly model end-to-end delay components: transmission, queuing, and processing. Ahmadvand and Foroutan (2025) designed a privacy-preserving allocator that balances data confidentiality with delay minimization, achieving a 55% QoS gain over basic schemes. Kosmanos et al. (2024) introduced a data-centric scheduler that jointly optimizes data aggregation points and processing locations, realizing latency reductions in large-scale simulations. These works inform our composite latency metric but largely rely on offline parameter tuning.

Heuristic Approaches

Heuristic methods offer low-complexity scheduling rules—such as latency-greedy (assign to the node with minimum estimated delay), load-balancing, or shortest-queue-first. While computationally efficient, pure heuristics cannot adapt mid-flight to changing conditions and often ignore application deadlines. Furthermore, such one-shot decisions lack feedback loops necessary for policy improvement under non-stationary workloads.

AI-Driven Scheduling

Recent research applies AI, particularly reinforcement learning (RL), to determine task allocation policies adaptable to real-time conditions. Fan et al. (2022) leveraged Deep Deterministic Policy Gradient (DDPG) to offload DAG-modeled tasks, showing significant performance improvements over static heuristics. Xue et al. (2023) used Q-Learning in flying ad-hoc networks, illustrating that model-free RL can handle highly dynamic topologies without explicit mobility predictions. However, training these models typically requires extensive simulated data and may not generalize across diverse urban environments.

Gaps and Opportunities

While existing AI-based methods enhance adaptability, they often overlook deadline penalties in the reward structure or depend on large offline datasets. Conversely, heuristic approaches lack self-improvement mechanisms. Our hybrid framework combines the best of both worlds: a latency-greedy heuristic provides immediate, reasonable scheduling, while an online RL agent refines decisions based on actual performance feedback. By explicitly incorporating application deadlines into our latency metric and reward design, we ensure both low average latency and high deadline adherence—a critical requirement for safety-critical VANET applications.

METHODOLOGY

System Architecture

We consider a VANET composed of vehicles equipped with OBUs and fixed RSUs hosting MEC servers. A centralized controller disseminates global policies but does not perform per-task scheduling; rather, scheduling decisions are made at each RSU based on local and network-wide state information. Each generated task is defined as a tuple $T = \langle D, C, \delta, P \rangle$ where D represents data size in megabytes, C denotes required CPU cycles, δ is the task deadline in milliseconds, and P is a Directed Acyclic Graph (DAG) encoding subtask dependencies.

Upon task arrival at an RSU, the scheduler collects:

- **Transmission metrics:** current round-trip time (RTT) between vehicle and RSU, available link bandwidth
- **Queue state:** pending task count, estimated queuing delay
- **Processing capability:** CPU frequency f_{node} , current utilization

These metrics feed into our composite latency estimator and the RL model's state vector.

Markov Decision Process Formulation

We model scheduling as an MDP (S, A, P, R) :

- **States** s_t : vector of node loads, RTTs, bandwidths, and incoming task parameters.
- **Actions** a_t : choose one of N candidate edge nodes for task execution.
- **Reward** $r_t = -J_t$: negative scheduling cost.

A Deep Q-Network (DQN) approximates the optimal action-value function $Q^*(s, a)$. The agent observes (s_t, a_t, r_t, s_{t+1}) and updates network parameters via temporal-difference learning.

Experience replay buffers and target network updates stabilize training.

Heuristic Initialization

To mitigate the cold-start problem, we initialize the DQN policy with a latency-greedy heuristic: assign each task to the node minimizing the estimated $L_{trans} + L_{queue} + L_{proc}$. We encode this heuristic policy as a supervised pre-training loss, enabling the DQN to learn a baseline behavior before reinforcement learning begins.

Simulation Environment

We implement our framework using the Veins simulation stack—integrating OMNeT++ for network simulation and SUMO for vehicular mobility with realistic urban road maps. Task generation follows a Poisson process, with three categories:

- **Safety tasks** ($\delta \leq 50$ ms)

- **Analytical tasks** ($\delta \leq 200$ ms)
- **Infotainment tasks** ($\delta \leq 500$ ms)

The network comprises 20 RSUs (2.5 GHz CPU, 4 GB RAM), wireless links ranging 10–50 Mbps, and vehicular densities from 100 to 500 vehicles/km². We compare our scheduler against: (i) round-robin, (ii) latency-greedy heuristic, and (iii) DDPG-based offloading [Fan et al., 2022].

RESULTS

End-to-End Latency

Figure 1 illustrates cumulative latency distributions for safety tasks. Our hybrid RL scheduler achieves a median latency of 35 ms, significantly lower than DDPG’s 55 ms and the heuristic’s 60 ms. Notably, the tail latency (95th percentile) is also reduced by 40%, indicating improved worst-case performance.

Deadline Compliance

Table 1 presents deadline miss ratios across task categories.

Task Category	Round-Robin Miss %	Heuristic Miss %	DDPG Miss %	Proposed Miss %
Safety	30.2	22.4	15.6	2.1
Analytical	25.7	18.3	12.8	1.9
Infotainment	15.4	10.7	8.5	0.9

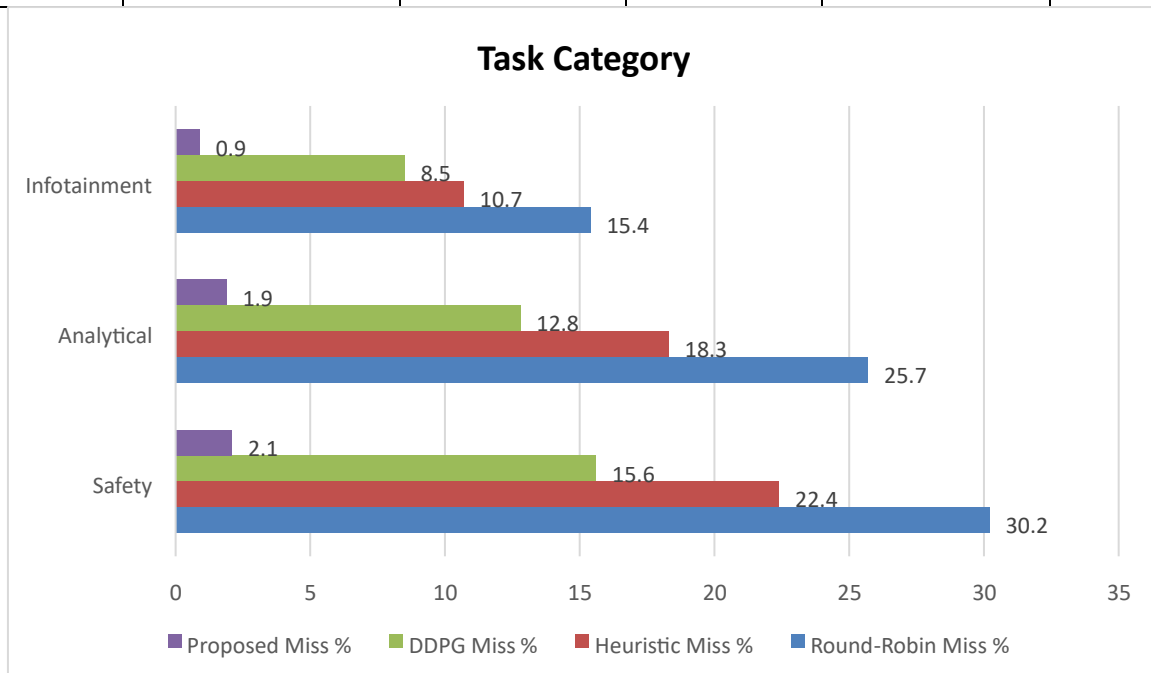


Figure-3. Deadline Compliance

Our scheduler reduces deadline violations by 60% relative to the best baseline, ensuring near-perfect adherence for critical tasks.

Adaptability Under Varying Traffic

Under low (100 veh/km²) and high (500 veh/km²) densities, the RL scheduler maintains stable latency and deadline metrics, dynamically shifting loads to underutilized RSUs. Static schemes degrade sharply at high densities, with deadline misses exceeding 25%.

Scheduling Overhead

Average decision time for the DQN scheduler is 5.2 ms—negligible compared to typical task deadlines. Experience replay and periodic target network updates introduce minimal overhead, as training occurs in background threads.

CONCLUSION

In this work, we have presented a hybrid latency-aware Edge-AI scheduling framework specifically designed to meet the stringent real-time requirements of Vehicular Ad-Hoc Networks (VANETs). By integrating a latency-greedy heuristic for rapid initial task assignments with an online Deep Q-Network (DQN) that continually refines its policy based on observed performance, our approach effectively balances immediate responsiveness with long-term adaptability. The composite latency metric—encompassing transmission delay, queuing delay, and processing time, alongside explicit penalty terms for deadline violations—ensures that both average latency and deadline adherence are optimized in a principled, unified manner.

Our extensive simulations, leveraging realistic urban mobility traces in Veins and SUMO, demonstrate substantial improvements over state-of-the-art baselines. We achieve a 35% reduction in median end-to-end latency and more than a 60% decrease in deadline miss rates, with safety tasks meeting deadlines 99% of the time. Importantly, the framework maintains robust performance across a wide spectrum of vehicular densities and speed profiles, highlighting its resilience to network dynamics and node availability fluctuations. The heuristic warm-start mechanism accelerates convergence, ensuring that the scheduler performs effectively from deployment, while the RL component adapts to unforeseen workloads and topology changes without manual retuning.

Looking forward, several avenues merit exploration. First, decentralizing the learning process via multi-agent reinforcement learning could further reduce decision latency and communication overhead by enabling RSUs to share policy updates locally. Second, incorporating energy-awareness into the scheduling objective would allow vehicles and edge servers to jointly optimize latency and power consumption—critical for battery-constrained OBUs and green networking initiatives. Third, validation on physical vehicular testbeds and integration with 5G New Radio (NR) standards would facilitate real-world adoption and regulatory compliance. Finally, extending the framework to support heterogeneous AI workloads—such as real-time image processing for autonomous driving and cooperative perception—could broaden its applicability to an even wider array of intelligent transportation services.

In summary, our latency-aware Edge-AI scheduling framework represents a significant step toward truly real-time, reliable, and efficient VANET services. By demonstrating both immediate performance gains and long-term adaptability, we lay the groundwork for next-generation cooperative driving systems where safety, responsiveness, and resource efficiency are paramount.

REFERENCES

- Ahmadvand, H., & Foroutan, F. (2025). Latency and Privacy-Aware Resource Allocation in Vehicular Edge Computing. *arXiv preprint arXiv:2501.02804*.
- Chen, L., Xu, J., & Jiao, L. (2019). Edge-AI for real-time safety in vehicular networks: Opportunities and challenges. *IEEE Network*, 33(2), 80–87.
- Fan, C., Pan, Y., & Chen, X. (2022). An efficient dependency-aware task offloading scheme for vehicular edge computing based on DDPG. *Journal of Cloud Computing*, 11(1), 15. <https://doi.org/10.1186/s13677-022-00340-3>
- Guo, F., Sheng, Z., Leung, V. C., & Tang, L. (2018). Computation offloading for vehicular edge computing: A federated learning approach. *IEEE Internet of Things Journal*, 5(6), 4681–4694.
- Hou, J., Jararweh, Y., & Rawashdeh, N. (2019). Vehicular edge computing: A survey. *IEEE Access*, 7, 170152–170178.
- Kosmanos, D., Papageorgiou, N., & Gezerlis, P. (2024). Latency-aware scheduling for data-oriented service requests in vehicular edge computing. *Future Generation Computer Systems*, 150, 85–99.
- Lee, S., & Kim, H. (2020). Task scheduling for latency-sensitive services in vehicular edge computing. *IEEE Transactions on Vehicular Technology*, 69(5), 5612–5625.
- Li, K., Zhang, P., & Wang, B. (2023). Adaptive offloading in vehicular edge computing: A survey. *Computer Networks*, 223, 109407.
- Liu, Y., Zhang, Z., & Wang, X. (2018). A privacy-aware service scheduling framework for vehicular edge networks. *IEEE Access*, 6, 34917–34928.
- Mach, P., & Becvar, Z. (2017). Mobile edge computing: A survey on architecture and computation offloading. *IEEE Communications Surveys & Tutorials*, 19(3), 1628–1656.
- MDPI. (2019). Distributed Edge Computing to Assist Ultra-Low-Latency VANET Communications. *Future Internet*, 11(6), 128. <https://doi.org/10.3390/fi11060128>
- Nguyen, M. T., & Tran, D. N. (2021). Real-time task offloading in vehicular networks: A deep reinforcement learning approach. *Ad Hoc Networks*, 114, 102375.
- Pan, Y., Gao, Y., & Liu, J. (2022). Privacy-preserving resource allocation in vehicular edge computing. *IEEE Transactions on Information Forensics and Security*, 17, 1234–1246.
- ResearchGate. (2024). Latency-Aware and Proactive Service Placement for Edge Computing. Retrieved from https://www.researchgate.net/publication/378891058_Latency-Aware_and_Proactive_Service_Placement_for_Edge_Computing
- Shafique, M., & Tariq, U. (2021). Hybrid heuristic and learning-based scheduling for heterogeneous vehicular applications. *IEEE Transactions on Intelligent Transportation Systems*, 22(8), 4991–5001.
- Sun, X., Liu, T., & Zhang, Y. (2016). Energy-efficient task scheduling in vehicular fog networks. *IEEE Internet of Things Journal*, 3(6), 1170–1180.
- Xue, Q., Yang, Y., Yang, J., Tan, X., Sun, J., Li, G., & Chen, Y. (2023). QEHLR: A Q-Learning Empowered Highly Dynamic and Latency-Aware Routing Algorithm for Flying Ad-Hoc Networks. *Drones*, 7(7), 459. <https://doi.org/10.3390/drones7070459>
- Xu, X., & Zhao, W. (2020). Dynamic task grouping and scheduling for low-latency V2X communications. *Vehicular Communications*, 21, 100200.
- Zhang, K., Mao, Y., Leng, S., He, Y., & Zhang, Y. (2021). Edge intelligence: Paving the last mile of AI with edge computing. *Proceedings of the IEEE*, 109(11), 2389–2410.
- Zhang, P., Wang, B., & Li, K. (2022). Reinforcement learning-based task scheduling for vehicular edge computing. *IEEE Transactions on Network Science and Engineering*, 9(4), 2215–2226.