

Evaluating Lightweight Reinforcement Learning for Adaptive Traffic Signal Control:

A Comparative Study on TAPAS Cologne and Hangzhou Scenarios

Maryam Salem Alshehyari
Machin Learning
24010700
maryam.alshehyari@mbzuai.ac.ae

Bashayer Abdulla Alsereidi
Machin Learning
24010641
bashayer.alsereidi@mbzuai.ac.ae

Naema Mohamed Alkhzaimi
Machin Learning
24010724
naema.alkhzaimi@mbzuai.ac.ae

Project Repository:

<https://github.com/MaryamAlshehyari/Adaptive-RL-Traffic-Signal-Control>

Abstract—This work presents a comparative study of lightweight reinforcement learning (RL) models for adaptive traffic signal control using realistic urban simulation environments. As cities face increasing congestion and inefficiencies in static traffic light systems, data-driven control strategies offer promising alternatives. Leveraging the `sumo-rl` framework and two benchmark scenarios—TAPAS Cologne and Hangzhou 4x4 grid—we evaluate multiple RL techniques, including tabular Q-Learning, Deep Q-Networks (DQN), Quantile Regression DQN (QRDQN), and Proximal Policy Optimization (PPO). Our objective is to investigate the trade-offs between performance and computational simplicity, focusing on average waiting time, total vehicle arrivals, and stopped vehicle counts.

Each model was trained independently per scenario, with additional hybrid setups utilizing Q-table initialization and supervised pretraining. Experimental results demonstrate that PPO consistently outperforms other models, achieving the lowest waiting times and earliest completion across both environments. Tabular Q-Learning showed competitive results in smaller networks but lacked scalability. DQN variants were sensitive to initialization and struggled to generalize in more complex settings. QRDQN, while theoretically robust, underperformed in practice without careful tuning.

This study highlights the viability of lightweight and modular RL strategies for traffic signal control and provides practical insights into their deployment trade-offs. Our findings suggest that actor-critic methods, particularly PPO, offer the most balanced solution for real-time traffic optimization with scalable potential.

Index Terms—Reinforcement Learning, Traffic Signal Control, SUMO, Smart Transportation, Deep Q-Networks, Proximal Policy Optimization, Multi-Agent Systems, Internet of Things.

I. INTRODUCTION

Urbanization has led to a substantial surge in vehicular traffic, placing immense strain on urban infrastructure and existing traffic management systems [1]. Traditional traffic signal control methods—typically based on static timing plans or pre-configured actuated responses—are often unable to adapt to real-time traffic fluctuations. This lack of adaptability contributes directly to critical urban challenges such as traffic

congestion, increased travel delays, prolonged vehicle idling, and elevated emissions of greenhouse gases and pollutants [2].

Recent advancements in the Internet of Things (IoT) and machine learning technologies offer promising avenues to modernize traffic control systems. IoT-enabled infrastructures allow for the real-time collection of traffic data, while machine learning, and in particular reinforcement learning (RL), provides mechanisms for systems to learn optimal traffic signal policies dynamically [3]. Through continuous interaction with simulated environments, RL agents are capable of adjusting signal timings in response to real-time conditions without requiring pre-defined rules or expert supervision [4].

Several studies have demonstrated that deep RL methods, including Deep Q-Networks (DQN), Proximal Policy Optimization (PPO), and related architectures, can outperform traditional fixed-time or actuated signal control approaches in simulation environments [5]–[7]. However, these models often involve significant computational overhead, long training times, and hyperparameter sensitivity—factors that challenge their applicability in real-world deployments, particularly in resource-constrained settings [8].

This work aims to address these limitations by evaluating a spectrum of lightweight reinforcement learning models—including tabular Q-learning, DQN, QRDQN, and PPO—within realistic, large-scale urban traffic simulations using the `sumo-rl` framework. Focusing on two widely recognized benchmark scenarios—TAPAS Cologne and Hangzhou—we conduct a comparative study that emphasizes:

- **Resource-Efficient Modeling:** Assessing the performance of low-complexity RL agents under constrained computational settings.
- **Realistic Evaluation:** Benchmarking RL agents against fixed-time and actuated control baselines using key performance indicators (KPIs) such as average waiting time, queue length, and vehicle throughput.
- **Scalability and Feasibility:** Investigating the trade-offs between model simplicity, adaptability, and real-world deployability.

By prioritizing efficiency and feasibility over brute-force optimization, this study contributes to the growing field of intelligent traffic systems by offering insights into the capabilities of accessible, deployable RL-based control policies in realistic urban environments.

II. RELATED WORK

The integration of Internet of Things (IoT) technologies with traffic signal control has led to a variety of innovative approaches aimed at improving urban mobility. Early work by Smith and Brown [9] focused on collecting real-time traffic data such as vehicle counts and speeds, applying optimization techniques to adjust signal phases. Although their system showed improvements in small-scale simulations, it lacked scalability and adaptability, relying solely on traditional optimization without any learning capability.

Sharma et al. [10] adopted a machine learning approach, using pre-trained models (e.g., decision trees and support vector machines) to predict traffic flow and guide signal timing. While effective in short-term forecasting, their method lacked real-time adaptability and struggled to handle dynamic traffic conditions, particularly across multiple intersections.

Subsequent studies have turned to deep reinforcement learning (DRL) for more flexible, adaptive traffic control. Chen et al. [5], [7] implemented DRL frameworks like CoLight and PressLight, allowing agents to learn signal-switching policies through environment interactions. Their systems achieved significant improvements in average waiting time and intersection throughput. However, the training processes were computationally expensive, requiring powerful GPUs and long simulation times, with unstable performance during early training phases.

Shang et al. [6] advanced the field by integrating Graph Neural Networks (GNNs) with multi-agent RL architectures, modeling intersections as nodes in a graph to coordinate control policies. This approach improved global coordination across traffic networks but introduced complexity and high infrastructure demands, making real-world implementation challenging.

More recently, Shao et al. [11] proposed MoveLight, a movement-level DRL framework leveraging multi-head attention and fine-grained signal optimization. Their system was evaluated on realistic benchmark datasets, including TAPAS Cologne and Hangzhou, using SUMO and CityFlow simulators. Despite promising performance, MoveLight remains experimental, with sensitivity to hyperparameter tuning and no large-scale deployment to date.

Unlike prior studies that have prioritized maximizing performance through complex, resource-intensive deep reinforcement learning models, this work adopts a more pragmatic approach. Our objective is to investigate the viability of **lightweight, resource-efficient RL agents** under realistic constraints, which are more representative of typical deployment conditions in cities with limited infrastructure or compute capabilities.

Motivated by the success and popularity of benchmark datasets such as **TAPAS Cologne** and **Hangzhou**, which are commonly used in high-performance RL research, we intentionally applied these same scenarios as testbeds to examine how simpler models perform in the same environments. This provides a direct lens through which to compare algorithmic efficiency and trade-offs in performance, scalability, and learning complexity.

By focusing on models such as **Q-learning**, **DQN**, **QRDQN**, and **PPO**, this study contributes a **baseline-oriented, deployment-conscious perspective** to the field of adaptive traffic signal control—one that prioritizes interpretability, reproducibility, and accessibility over computational intensity.

III. SYSTEM MODEL

A. Overview

The proposed system is a modular reinforcement learning framework integrated with microscopic traffic simulation for adaptive traffic signal control. It combines SUMO (Simulation of Urban Mobility) with RL libraries such as Stable-Baselines3, PyTorch, and Ray RLlib to enable both centralized and decentralized agent training across diverse urban scenarios. The architecture supports tabular and deep reinforcement learning methods, allowing for flexible deployment and comparative experimentation. Figure 1 illustrates the modular architecture of the system, detailing the flow of traffic data through SUMO and the experimental reinforcement learning agents.

B. Core Components

Traffic Simulator: The SUMO engine serves as the foundation for environment dynamics, offering high-resolution simulation of urban traffic networks using XML-based topology (`.net.xml`) and routing (`.rou.xml`) files.

Environment Wrappers: Environment interactions are mediated via the `sumo-rl` Python API, with additional wrappers from `gymnasium` (for SB3), and `PettingZoo` (for RLlib) enabling compatibility with single-agent and multi-agent RL pipelines.

Observation and Action Handling: Observation vectors include phase encodings, lane queue lengths, and density metrics, while actions are defined as discrete traffic phase switches. These are encapsulated in modular handlers within `observations.py` and `traffic_signal.py`.

Reinforcement Learning Algorithms: The framework supports a diverse algorithmic range:

- **Q-Learning:** Tabular agents with epsilon-greedy exploration, implemented via `sumo-rl`.
- **DQN:** Deep Q-Networks using Stable-Baselines3 and custom PyTorch implementations.
- **DQN with Q-table Initialization:** SB3-based DQN initialized with tabular Q-values to enhance convergence.
- **DQN with Q-table Pretraining:** Supervised PyTorch training on prior Q-tables before standard RL.

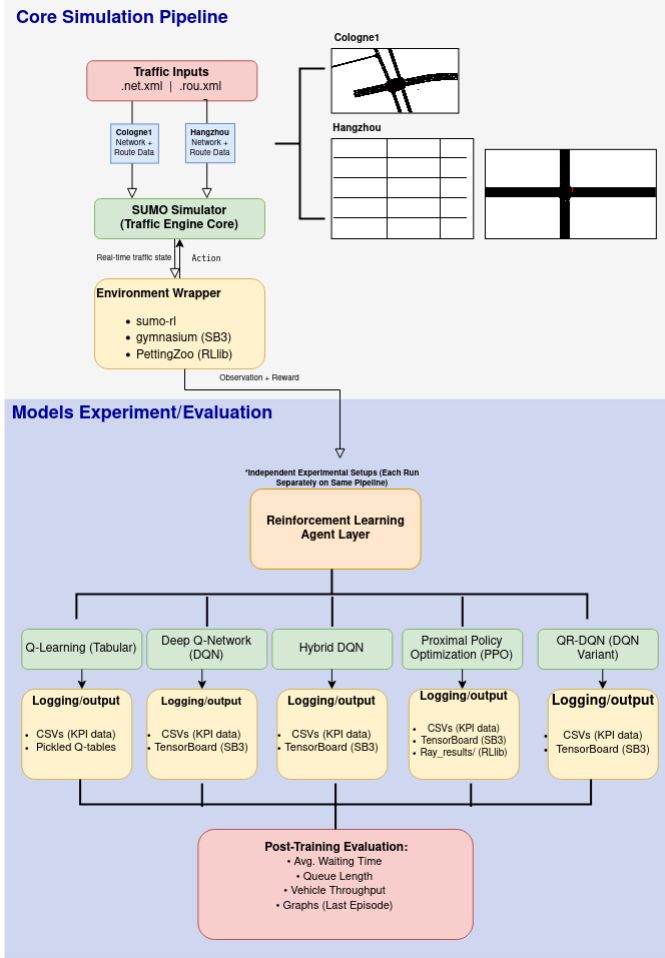


Fig. 1: System architecture showing simulation components, environment wrappers, independent RL agent setups, and post-training evaluation pipeline.

- **PPO:** Proximal Policy Optimization, evaluated using both SB3 (single-agent) and RLlib (multi-agent, distributed).
- **QRDQN:** A distributional reinforcement learning agent implemented using `sb3_contrib.QRDQN`, which learns a full quantile-based value distribution rather than scalar Q-values. This allows for more robust decision-making under uncertainty.

Multi-Agent Coordination: For multi-intersection setups, the architecture supports fully decentralized learning, assigning one agent per traffic light without shared parameters or coordination layers. Multi-agent compatibility is handled through SUMO’s native traffic light identifiers and `env.ts_ids`.

C. Traffic Scenarios and Datasets

Two real-world traffic scenarios were used to evaluate and compare the performance of multiple RL algorithms under diverse urban conditions:

- **TAPAS Cologne1:** This scenario is derived from the TAPAS Cologne dataset, based on microscopic traffic simulation of the Cologne metropolitan area in Germany.

For our experiments, a single 1x1 intersection was selected from the broader network to create a focused, controllable testbed (see Fig. 2). The dataset includes empirically generated vehicle flows that replicate real-world mobility behavior, including turning probabilities and realistic signal phase patterns.

- **Hangzhou 4x4 Grid:** The Hangzhou scenario models a dense 4x4 urban grid extracted from the Gudang District in Hangzhou, China. While structured as a grid for multi-agent control, the vehicle flows and road design are based on actual city traffic statistics. Each intersection is controlled independently, enabling fully decentralized training across 16 simultaneous agents (see Fig. 4).

The simulations were run using the SUMO engine, which computes real-time vehicle dynamics, traffic light behavior, and phase transitions. Input files include network definitions (‘.net.xml’) and route files (‘.rou.xml’), which provide the static layout and dynamic traffic inflows. These scenarios were chosen to test algorithm robustness across different levels of complexity and coordination.

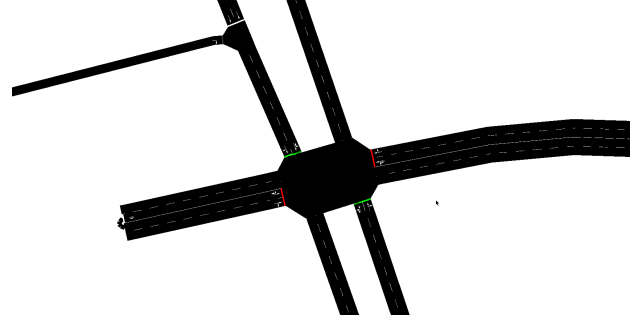


Fig. 2: TAPAS Cologne1 scenario — single real-world intersection cropped from the Cologne road network.

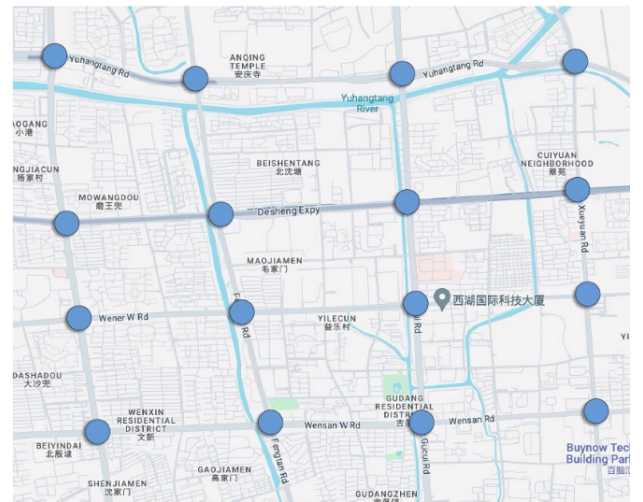


Fig. 3: Hangzhou 4x4 grid scenario — realistic multi-intersection urban environment with distributed control.

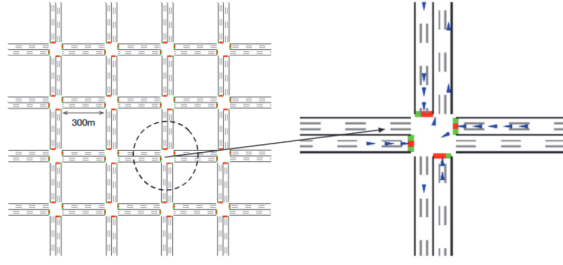


Fig. 4: Hangzhou 4x4 grid scenario — road network in simulator.

D. Execution and Logging

Training is conducted in episode-based or timestep-based modes, with episode lengths ranging from 3,000 to 100,000 simulation seconds. Output includes structured CSV logs, TensorBoard summaries, and saved model checkpoints. Q-tables are serialized using Pickle, and per-agent logs are organized by experiment type and timestamp.

E. Post-Training Evaluation

Performance metrics for all models were collected from the final episodes of their respective training processes. These final episodes represent the agent’s behavior after convergence and serve as the basis for comparison across models. Nonetheless, since exploration was already significantly decayed by the final training episodes, the collected results provide a reliable estimate of each model’s performance across the core metrics: average waiting time, number of stopped vehicles, total arrivals, and episode completion time.

IV. RESULTS AND EVALUATION

The performance of the proposed traffic signal control models was evaluated using three core Key Performance Indicators (KPIs):

- **System Mean Waiting Time (seconds):** The average time vehicles spent idling at intersections.
- **Agents Total Stopped:** The number of vehicles stopped across all traffic signals.
- **System Total Arrived Vehicles:** The cumulative count of vehicles that successfully exited the simulation.

Each metric was logged across simulation timesteps using SUMO output and visualized to assess both real-time dynamics and final convergence performance.

A. Cologne1 Results

To contextualize the performance of learning-based models, a fixed-time baseline was implemented using SUMO’s default static signal control. As shown in Table I, it failed to optimize flow, indicating that mere completion speed is not indicative of efficiency. By contrast, even the most basic RL setup (Q-learning) reduced average waiting time by over 99%, clearly demonstrating the value of adaptive control. This reinforces the critical limitation of static systems in dynamic urban

environments and justifies the use of intelligent signal control strategies.

Table I summarizes the performance of each model on the Cologne1 intersection.

TABLE I: Cologne1 – Final KPI Comparison

Model	Wait Time ↓	Stopped ↓	Arrived ↑	Finish Step ↓
Fixed-Time Baseline	1567.21	7839	322.25	Didn’t finish
Q-Learning	2.83	4.74	477.02	1755
DQN	287.38	4.67	479.59	5690
DQN + Q-table Init	6.98	12.52	352.50	650
PPO	0.75	5.64	440.08	665
QR-DQN	35.85	17.44	397.21	1135

(See Appendix, Fig. 7 for detailed results.)

Among all, **PPO** clearly demonstrated the best performance with the lowest average waiting time (0.75s) and fastest completion. While Q-Learning achieved slightly higher vehicle throughput, it required longer simulation steps. QR-DQN underperformed with elevated wait times and congestion, suggesting instability in the learning process.

B. Hangzhou Results

Table II compares final KPIs on the more complex Hangzhou 4x4 traffic grid.

TABLE II: Hangzhou – Final KPI Comparison

Model	Wait Time ↓	Stopped ↓	Arrived ↑	Finish Step ↓
Fixed-Time Baseline	353.66	25.15	1755.13	Didn’t finish
Q-Learning	0.0736	5.26	1775.30	4370
DQN (Custom Multi)	0.1069	6.52	1770.71	4475
DQN + Q-table Pretrain	0.1120	6.44	1772.49	4350
PPO (Ray RLlib)	0.0811	5.31	1771.59	4380

PPO again delivered the most balanced results, achieving the lowest waiting time and strong throughput without requiring aggressive pretraining. DQN with Q-table pretraining showed the fastest finish, but incurred a slight trade-off in overall flow efficiency.

C. PPO Training Stability and Convergence Analysis

To further support the evaluation of PPO’s performance, we present a compact visual overview of its training progression and final policy behavior for both test scenarios. Figures 5 and 6 respectively summarize PPO’s behavior on the Cologne1 and Hangzhou setups using paired subplots (training summary and final episode performance).

In the Cologne1 setup (Figure 5), PPO shows a steady reduction in mean waiting time and number of stopped vehicles over training episodes, with corresponding improvements in throughput. Despite occasional performance fluctuations, the model trends toward convergence. The final episode visualization reflects consistent, low-latency traffic flow and efficient clearance of queues, confirming that PPO has successfully generalized a performant signal policy.

In Hangzhou (Figure 6), PPO demonstrates faster convergence under the larger-scale, multi-agent setting. Key indicators stabilize early, and final episode metrics reflect sustained improvements: reduced waiting times, minimal vehicle queuing, and high exit rates. This suggests strong adaptability

and scalability of the actor-critic method even under higher complexity.

Together, these visualizations validate PPO’s superior convergence stability and real-time performance, supporting its selection as the top-performing model across both experimental domains.

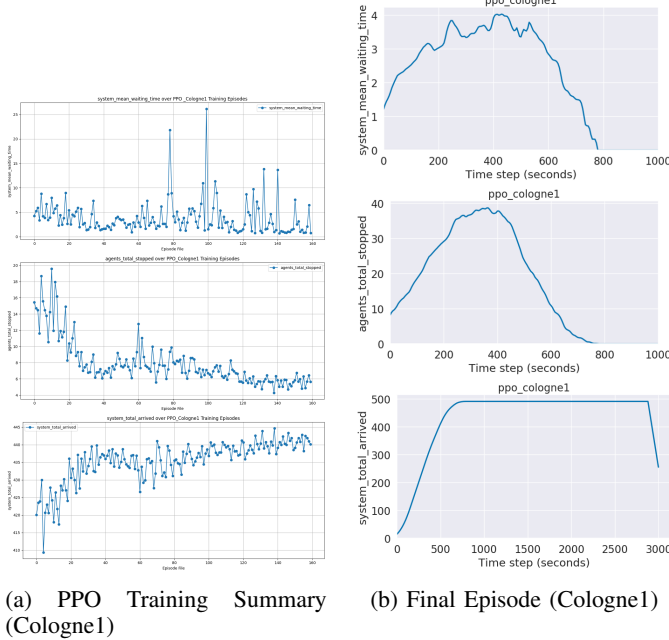


Fig. 5: PPO performance overview on Cologne1 scenario.

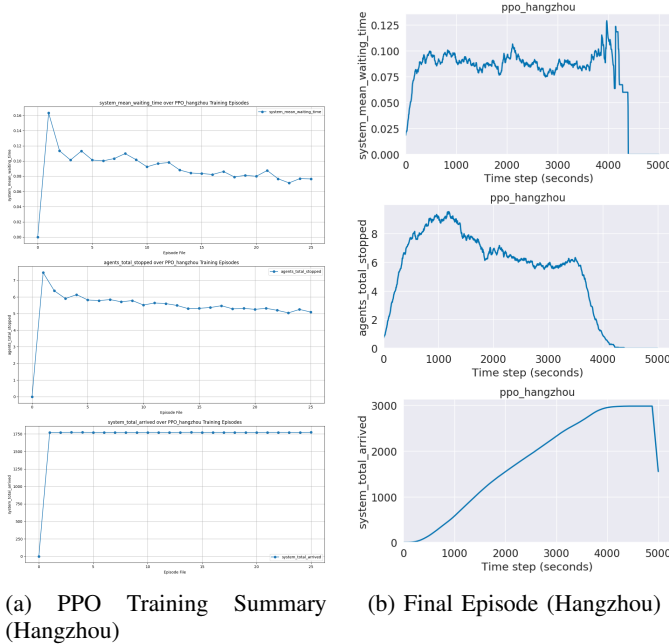


Fig. 6: PPO performance overview on Hangzhou scenario.

D. General Observations

Across both environments, PPO consistently outperformed other methods in terms of average wait time and system smoothness. While tabular Q-Learning delivered competitive results on smaller networks, it lacked the flexibility to scale efficiently to more complex grids. Deep Q-Learning variants exhibited mixed behavior — faster convergence with Q-table initialization in some cases, but unstable performance in others.

The evaluation confirms that actor-critic methods (PPO) not only offer superior performance but also demonstrate greater generalizability across network sizes. Visual inspection of the simulation curves further supports this conclusion, with PPO models exhibiting flatter, more stable trends over time.

V. CONCLUSION AND FUTURE WORK

This study investigated a range of reinforcement learning techniques for adaptive traffic signal control using realistic urban simulation environments. By leveraging both tabular and deep reinforcement learning methods—such as Q-Learning, Deep Q-Networks (DQN), Quantile Regression DQN (QRDQN), and Proximal Policy Optimization (PPO)—we explored the balance between performance, scalability, and deployment feasibility.

Our results demonstrate that actor-critic algorithms like PPO offer the most consistent and scalable performance across different network complexities. In the Cologne1 scenario, PPO achieved the lowest waiting time and fastest episode completion, outperforming both traditional Q-Learning and deep Q-based models. In the more complex Hangzhou 4x4 grid, PPO again led in wait time reduction, while DQN with Q-table pretraining showed competitive final-step efficiency.

While tabular Q-Learning proved surprisingly effective in simpler scenarios, its limitations became evident in larger topologies due to state-space explosion and lack of generalization. Deep Q-Networks performed variably—highlighting the challenges of reward shaping, exploration control, and hyperparameter tuning. QR-DQN offered modest gains but suffered from instability under longer simulation episodes.

In our opinion, the most promising direction is the integration of scalable actor-critic frameworks (e.g., PPO via Ray RLlib), which allow multi-agent deployment without requiring handcrafted exploration strategies or prior tabular policies. These models were not only easier to train but also more robust under unseen traffic patterns.

Challenges and Limitations

- **Training Time:** Large-scale simulations required long run-times, limiting extensive hyperparameter exploration.
- **Reward Engineering:** Some models were highly sensitive to reward definitions, particularly DQN variants.
- **Real-World Transferability:** While simulation results are encouraging, real-world deployment would require integration with sensor networks, real-time constraints, and safety protocols.

Future research may focus on:

- Extending the current models to support transfer learning across networks (e.g., train on Cologne, adapt to Hangzhou).
- Integrating real-time traffic sensor data (e.g., loop detectors, GPS probes) to enable deployment in live smart-city environments, beyond static simulation datasets.
- Exploring communication-aware multi-agent policies that allow traffic lights to coordinate their actions.
- Investigating transformer-based or graph neural network (GNN) policies to learn from the spatial topology of intersections.

In conclusion, this project validates the feasibility and impact of reinforcement learning for urban traffic management and offers a comparative benchmark that highlights trade-offs between interpretability, efficiency, and performance across algorithmic choices.

REFERENCES

- [1] D. Schrank, B. Eisele, T. Lomax, and J. Bak, “Urban mobility report,” Texas A&M Transportation Institute, 2019, accessed: 2023-10-15. [Online]. Available: <https://mobility.tamu.edu/umr/>
- [2] M. Barth and K. Boriboonsomsin, “Real-world carbon dioxide impacts of traffic congestion,” *Transportation Research Record*, vol. 2058, pp. 163–171, 2008.
- [3] A. Zanella, N. Bui, A. Castellani, L. Vangelista, and M. Zorzi, “Internet of things for smart cities,” *IEEE Internet of Things Journal*, vol. 1, no. 1, pp. 22–32, 2014.
- [4] L. Chen, K. Zhao, and H. Wang, “Optimizing traffic flow with reinforcement learning: A study on traffic light management,” in *IEEE Conference Proceedings*, 2023.
- [5] L. Chen, H. Zhang, and J. Wang, “Optimizing traffic flow using reinforcement learning: Colight and presslight,” in *Proceedings of the IEEE Conference on Intelligent Transportation Systems*, 2023.
- [6] Y. Shang, M. Liu, and K. Zhang, “Graph-based cooperation multi-agent rl for traffic signal control,” *IEEE Internet of Things Journal*, 2025.
- [7] L. Chen, H. Zhang, and J. Wang, “Optimizing traffic flow with reinforcement learning: A study on traffic light management,” in *Proceedings of the IEEE Conference on Intelligent Transportation Systems*, 2023.
- [8] J. Shang, S. Meng, and X. Zhou, “Graph-based cooperation multi-agent reinforcement learning for intelligent traffic signal control,” *IEEE Internet of Things Journal*, 2025.
- [9] J. Smith and L. Brown, “Iot-based traffic management with optimization techniques,” *Smart Cities*, vol. 5, no. 4, p. 66, 2022.
- [10] R. Sharma, A. Mehta, and S. Raghavan, “Iot-based traffic prediction and signal control for smart cities,” *ResearchGate Preprint*, 2021, available: <https://www.researchgate.net/publication>.
- [11] Q. Shao, L. Wu, and Z. Liu, “Movelight: Enhancing traffic signal control through movement-centric deep reinforcement learning,” *arXiv preprint arXiv:2407.17303*, 2024. [Online]. Available: <https://arxiv.org/abs/2407.17303>

APPENDIX A MODEL TRAINING CURVES

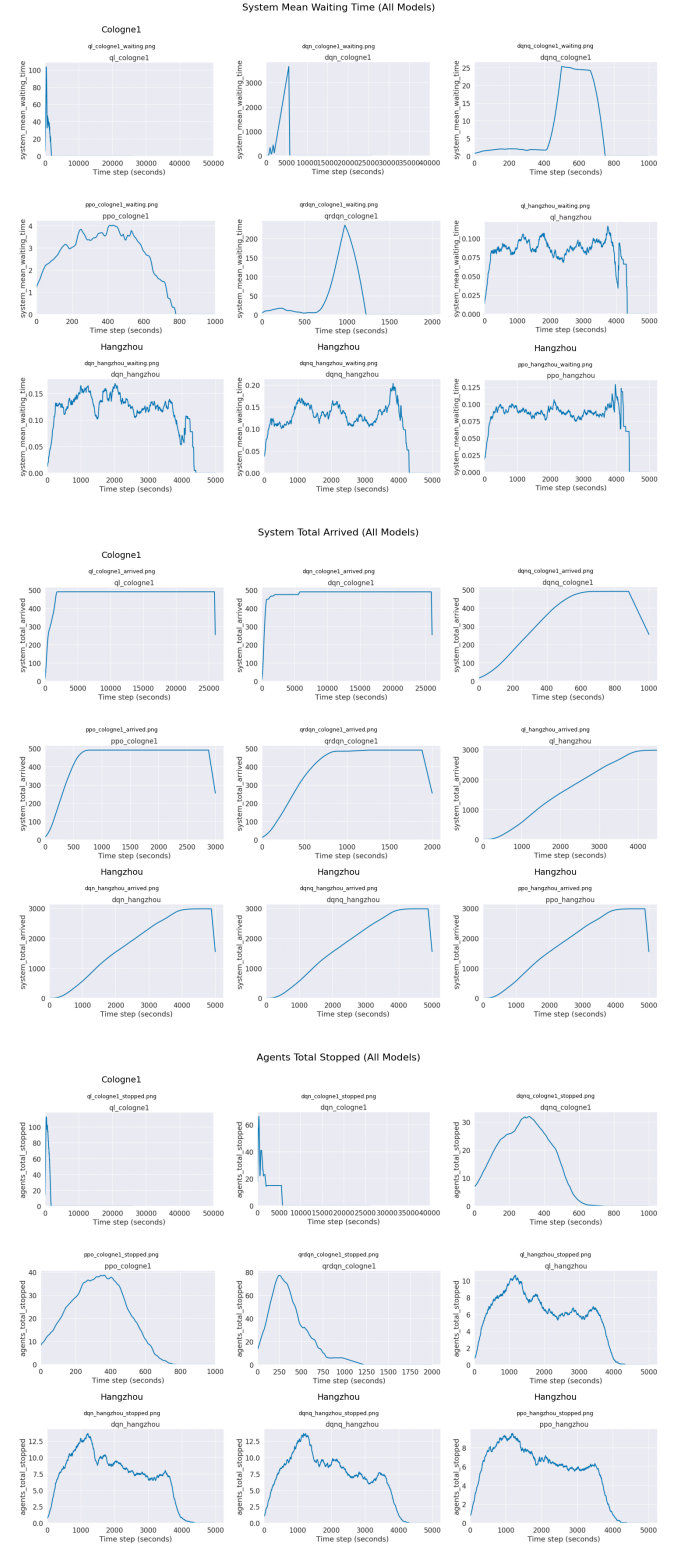


Fig. 7: Final episode KPIs comparison across all models: (a) waiting time, (b) total arrived vehicles, and (c) stopped agents.