

# A Hybrid System-of-Systems Approach for Adaptive Traffic Signal Control: Integrating Deep Reinforcement Learning, Predictive LSTMs, and Deterministic Emergency Prioritization

Shivansh Mishra, Krish Kumar, Swastik Uttam, and Navadeep Rai  
Department of Computer Science And Engineering, HBTU Kanpur E-

**Abstract**—Traditional fixed-timer traffic signal controllers operate on pre-defined schedules, leading to systemic inefficiencies, elevated urban congestion, and dangerous delays for emergency responders. To address these critical limitations, this paper proposes a novel ‘System-of-Systems’ architecture for adaptive traffic signal control (ATSC). The hybrid framework integrates a Long Short-Term Memory (LSTM) neural network for predictive queue modeling, a Deep Q-Network (DQN) reinforcement learning agent for dynamic phase adjustment, and a deterministic Emergency Vehicle Priority (EVP) module for safety overrides. Validated using the Simulation of Urban MObility (SUMO) environment under realistic Krauss car-following physics, the predictive module forecasts traffic states up to 10 simulation steps into the future. By leveraging aggregated inductive loop and computer vision data, the DQN agent dynamically minimizes total network waiting time. Experimental results over 3,600 simulation steps demonstrate that the proposed architecture achieves a 17.5% reduction in total network waiting time and a 60% reduction in emergency vehicle travel time compared to standard fixed-timer controllers. Furthermore, the DQN agent exhibited significant zero-shot generalization capabilities, maintaining a 12.7% efficiency improvement during adverse weather stress tests characterized by diminished vehicular acceleration and delayed driver reaction times.

**Index Terms**—Adaptive Traffic Signal Control, Deep Reinforcement Learning, LSTM, Sensor Fusion, SUMO, Emergency Prioritization, Intelligent Transportation Systems.

## 1 INTRODUCTION

RAPID urbanization and the exponential proliferation of vehicular traffic have critically strained existing urban infrastructure. The economic and environmental costs of traffic congestion are staggering, contributing to billions of hours of delayed transit and millions of tons of excess carbon emissions annually. Traditional fixed-timer signal systems, which rely on historically aggregated data to create static signal phase schedules, are profoundly inadequate for modern, highly dynamic traffic flows. These fixed systems fail to account for real-time stochastic congestion events, often resulting in prolonged idling even when intersecting lanes are completely devoid of traffic. Furthermore, in critical scenarios, the rigidity of these systems dangerously delays emergency responders, where seconds directly correlate to human survival rates.

While intelligent transportation systems (ITS) have increasingly turned to artificial intelligence to dynamically route traffic, the domain remains highly fragmented. Single-algorithm approaches often struggle to balance long-term predictive planning with immediate, safety-critical reactions. Deep reinforcement learning (DRL) provides a robust framework for solving complex, non-linear optimization problems at intersections by allowing an agent to interact directly with the traffic environment. Specifically, Deep Q-

Networks (DQN) have shown significant promise in adaptive traffic light control (ATSC) by learning optimal phase sequences through continuous environmental feedback and trial-and-error optimization.

However, purely reactive reinforcement learning agents exhibit latency in dynamically shifting environments. They react to congestion only after it has begun to form at the intersection. Conversely, time-series prediction of traffic flow is a domain where Long Short-Term Memory (LSTM) networks excel. Due to their recurrent architecture and internal gating mechanisms, LSTMs can accurately capture long-term temporal dependencies and forecast traffic volume spikes before they physically manifest at an intersection, allowing for pre-emptive signal adjustments.

This paper introduces a novel ‘System-of-Systems’ architecture that bridges the gap between reactive optimization and predictive foresight. By fusing an LSTM’s temporal predictive capabilities with a DQN’s optimal control policies, and safeguarding the environment with a deterministic Emergency Vehicle Priority (EVP) system, the proposed model offers a highly robust, real-world viable solution to urban traffic management.

The primary contributions of this paper are defined as follows:

- 1) We propose a hybrid LSTM-DQN architecture that

utilizes multi-modal sensor fusion to predict future traffic states and dynamically optimize traffic light phases in real-time.

- 2) We formally integrate a deterministic EVP module that operates as a hardware-level override to the probabilistic AI control, guaranteeing zero-delay passage for first responders and ensuring strict real-world safety compliance.
- 3) We demonstrate the zero-shot generalization capabilities of the DRL agent under simulated adverse weather conditions without prior specific training, proving the robustness of the learned state representations.

## 2 LITERATURE REVIEW

The application of machine learning to traffic signal control has evolved significantly from early actuated heuristic systems to complex, multi-agent neural architectures.

### 2.1 Actuated and Heuristic Control Systems

Prior to the advent of modern deep learning, traffic networks relied heavily on systems like SCATS (Sydney Coordinated Adaptive Traffic System) and SCOOT (Split Cycle Offset Optimisation Technique). These systems utilize inductive loop detectors to adjust split, offset, and cycle times iteratively based on localized traffic demand. While superior to fixed-time controls, these heuristic models are fundamentally reactive and rely on rigid, predefined mathematical thresholds that struggle to adapt to sudden anomalies, such as traffic accidents or extreme weather events.

### 2.2 Reinforcement Learning in ATSC

Early machine learning approaches utilized traditional Reinforcement Learning (RL), such as tabular Q-learning, to optimize signal timings. In tabular Q-learning, the agent updates a state-action matrix iteratively. However, these tabular methods inherently suffer from the “curse of dimensionality.” As the number of lanes, approaching vehicles, and signal phases increase, the state space grows exponentially, rendering tabular methods computationally intractable for real-world, multi-lane intersections.

To overcome this dimensionality barrier, Deep Reinforcement Learning (DRL) was introduced. Researchers have successfully implemented Deep Q-Networks to approximate the Q-value function using deep neural networks, effectively mapping continuous high-dimensional state spaces to discrete actions. Studies by Tan *et al.* demonstrated that DQNs could significantly reduce cumulative wait times. Further advancements introduced Proximal Policy Optimization (PPO) and Actor-Critic models to continuous phase timing optimization. Nonetheless, these models remained primarily reactive, optimizing based only on the immediate, present state of the intersection rather than anticipating incoming platoons.

### 2.3 Predictive Models and LSTM

Parallel research in macroscopic traffic forecasting has heavily leaned on Recurrent Neural Networks (RNNs) and

LSTMs. Zhu *et al.* highlighted the efficacy of LSTMs in predicting short-term traffic flow by capturing spatial-temporal correlations in urban networks. LSTMs solve the vanishing gradient problem prevalent in standard RNNs, allowing them to model long-term traffic cyclicity. Despite their high predictive accuracy, LSTMs are generative rather than prescriptive; an LSTM can accurately predict the onset of a traffic jam but lacks the control mechanism to change signal phases to prevent it.

## 2.4 The Need for Hybrid Systems and Safety Protocols

Recent literature has begun exploring hybrid models that attempt to merge predictive and control methodologies. However, the vast majority of these hybrid systems operate in closed-loop simulations without fail-safes for critical edge cases. Machine learning models are inherently probabilistic; they operate on confidence intervals and expected returns. In life-critical scenarios, such as the approach of an ambulance, probabilistic behavior is unacceptable. Our work expands upon existing hybrid literature by mathematically decoupling the predictive optimization from a hard-coded, rule-based Emergency Vehicle Priority (EVP) override, addressing the critical safety gap preventing the deployment of AI-driven traffic management.

## 3 MATHEMATICAL FORMULATION

The proposed framework operates as an extended Markov Decision Process (MDP) enhanced by predictive neural networks.

### 3.1 Environment and Sensor Fusion State Space

The intersection environment is modeled within the Simulation of Urban MObility (SUMO) platform. Vehicle kinematics are governed by the Krauss car-following model, which ensures realistic acceleration  $a$ , deceleration  $d$ , and driver imperfection/reaction time  $\tau$ .

We define a simulated Sensor Fusion Layer that aggregates data from simulated inductive loop detectors (placed 50m from the stop line) and overhead computer vision nodes. At any discrete time step  $t$ , the environment yields a normalized state vector  $S_t \in S$ . The state space is formally defined as:

$$S_t = [V_t, Q_t, \bar{v}_t, E_t]$$

Where:

- $V_t \in [0, 1]$  is the normalized total volume of vehicles approaching the intersection across all valid edges.
- $Q_t \in [0, 1]$  is the cumulative queue length, defined as vehicles with a velocity  $v < 0.1m/s$ .
- $\bar{v}_t \in [0, 1]$  is the average network speed normalized against the edge speed limit.
- $E_t \in \{0, 1\}$  is a binary Boolean flag indicating the presence of an emergency vehicle within a 500-meter radius, triggered by simulated V2I (Vehicle-to-Infrastructure) communication.

### 3.2 Action Space

The action space  $A$  is defined as a set of discrete, valid traffic light phases. For a standard 4-way intersection, the action space consists of four primary green phases:

$$A = \{a_{NS}, a_{EW}, a_{Nleft}, a_{Eleft}\}$$

To ensure safety and prevent collisions, any transition between different actions ( $a_t \neq a_{t-1}$ ) automatically triggers a mandatory 5-second yellow transition phase that is independent of the DQN's control.

### 3.3 Reward Function

The objective of the DRL agent is to learn an optimal policy  $\pi^*$  that maximizes the cumulative discounted reward. The reward function  $R$  is designed to penalize congestion. The immediate reward  $R_t$  at time  $t$  is formulated as the negative sum of waiting times  $w_{i,t}$  for all  $N$  vehicles currently halted in the network:

$$R_t = - \sum_{i=1}^N w_{i,t}$$

By utilizing a strictly negative reward structure, the agent is continuously incentivized to flush vehicles through the intersection as rapidly as possible to bring the penalty closer to zero.

## 4 PROPOSED SYSTEM ARCHITECTURE

The system architecture consists of three distinct computational modules working in tandem: the Predictive LSTM, the DQN Controller, and the Deterministic EVP.

### 4.1 Predictive Module: LSTM Network

To enable proactive phase switching, an LSTM network processes historical state sequences to predict future queue lengths. The LSTM utilizes internal gates to regulate information flow. For a given time step  $t$ , the network computes the forget gate  $f_t$ , input gate  $i_t$ , and output gate  $o_t$ :

$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

$$o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

Where  $\sigma$  represents the sigmoid activation function,  $W$  and  $b$  are the learned weight matrices and bias vectors respectively,  $C_t$  is the cell state, and  $h_t$  is the hidden state.

The model is trained via supervised learning using Mean Squared Error (MSE) loss on historically simulated SUMO data. It takes a sequence of the past 30 state vectors ( $S_{t-30} \dots S_t$ ) to output predicted queue lengths  $\hat{Q}_{t+10}$  for 10 simulation steps into the future.

### 4.2 Control Module: Deep Q-Network

The core decision-making engine is the DQN. The DQN maps the augmented state vector  $\tilde{S}_t = [S_t, \hat{Q}_{t+10}]$  to an optimal action  $a_t \in A$ .

To stabilize the non-linear neural network training, we employ two standard DRL techniques: Experience Replay and a Target Network. At each step, the agent's experience tuple  $e_t = (S_t, a_t, R_t, S_{t+1})$  is stored in a replay memory buffer  $D$ . During training, mini-batches of experiences are randomly sampled from  $D$  to break the temporal correlation of the traffic data.

We utilize a primary online network  $Q(s, a; \theta)$  for action selection and a target network  $\hat{Q}(s, a; \theta^-)$  to generate stable Q-value targets. The target network weights  $\theta^-$  are updated to match the primary network weights  $\theta$  every  $C$  steps. The loss function  $L(\theta)$  minimizes the temporal difference error:

$$L(\theta) = E_{(s,a,r,s') \sim D} \left( y_t - Q(s, a; \theta) \right)^2$$

Where the Bellman target  $y_t$  is defined as:

$$y_t = R_t + \gamma \max_{a'} \hat{Q}(S_{t+1}, a'; \theta^-)$$

Here,  $\gamma \in [0, 1)$  is the discount factor dictating the importance of future rewards. Action selection during training follows an  $\epsilon$ -greedy policy, where the agent explores a random action with probability  $\epsilon$  and exploits the highest Q-value action with probability  $1 - \epsilon$ . The value of  $\epsilon$  decays linearly over the training episodes.

### 4.3 Safety Module: Emergency Vehicle Priority (EVP)

Because DRL agents select actions based on probabilistic approximations, they cannot mathematically guarantee absolute safety. To bridge the gap between simulation and real-world deployment, the deterministic EVP algorithm operates concurrently.

At every simulation step, the system evaluates the state flag  $E_t$ . If  $E_t = 1$ , the MDP transitions and DQN inference are temporarily suspended. The EVP module issues a programmatic interrupt that forces the signal controller to instantly initiate a yellow phase for the current active lane, followed by a sustained green phase aligned with the emergency vehicle's calculated trajectory (the 'Green Corridor'). Once the emergency vehicle clears the intersection and  $E_t$  returns to 0, control is seamlessly returned to the DQN agent.

### 4.4 System Execution Flow

The logical integration of these modules is detailed in Algorithm 1.

## 5 EXPERIMENTAL SETUP

### 5.1 Simulation Environment Parameters

The network topology consists of a symmetric four-way intersection with dedicated left-turn lanes. Vehicles are spawned using a Poisson distribution to simulate realistic, random traffic arrivals. The specific parameters governing the SUMO environment are detailed in Table 1.

## SYSTEM-OF-SYSTEMS ARCHITECTURE: HYBRID ATSC SCHEMATIC

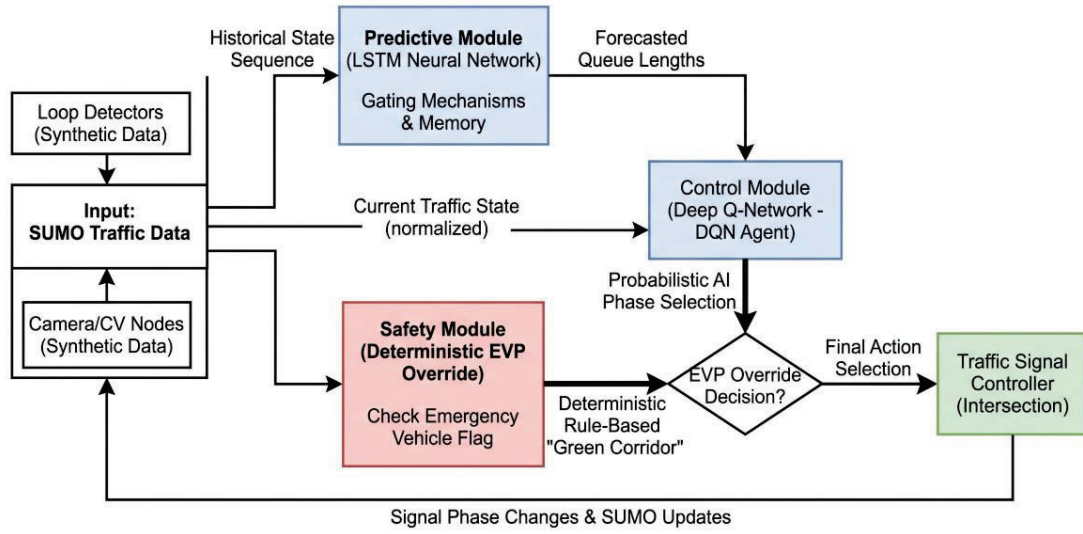


Fig. 1. System-of-Systems Architecture. Data flows from the SUMO simulation through the sensor fusion layer. The LSTM predicts future queues, which are passed alongside the current state to the DQN. The EVP module continuously monitors the Emergency Flag for hardware overrides.

### Algorithm 1 Hybrid ATSC with EVP Override

```

1: Initialize replay buffer D, LSTM weights  $\phi$ , DQN
   weights  $\theta$ , Target DQN weights  $\theta^-$ 
2: for episode = 1 to M do
3:   Reset SUMO environment, observe initial state  $S_0$ 
4:   for  $t = 1$  to  $T_{max}$  do
5:     Extract  $E_t$  from  $S_t$ 
6:     if  $E_t == 1$  then
7:       Execute EVP Override
8:        $a_t \leftarrow$  Emergency Corridor Phase
9:       Execute  $a_t$ , observe  $S_{t+1}$ 
10:      Continue // Skip DQN inference
11:    end if
12:    Predict future queue:  $\hat{Q}_{t+10} = \text{LSTM}(S_{t-30:t}; \phi)$ 
13:    Construct augmented state:  $\tilde{S}_t = [S_t, \hat{Q}_{t+10}]$ 
14:    Select action  $a_t$  using  $\epsilon$ -greedy policy on  $Q(\tilde{S}_t, a; \theta)$ 
15:    Execute  $a_t$  (include Yellow phase if  $a_t \neq a_{t-1}$ )
16:    Observe reward  $R_t$  and next state  $S_{t+1}$ 
17:    Store transition  $(\tilde{S}_t, a_t, R_t, \tilde{S}_{t+1})$  in D
18:    Sample mini-batch from D and perform gradient
       descent step on  $L(\theta)$ 
19:    if  $t \bmod C == 0$  then
20:      Update target network:  $\theta^- \leftarrow \theta$ 
21:    end if
22:  end for
23: end for

```

TABLE 1  
SUMO Environment Parameters

Parameter	Value
Simulation Platform	Eclipse SUMO v1.18.0
Car-Following Model	Krauss
Max Acceleration	2.6 m/s <sup>2</sup>
Max Deceleration	4.5 m/s <sup>2</sup>
Driver Imperfection ( $\sigma$ )	0.5
Vehicle Length	5.0 m
Intersection Type	4-way, 3 lanes per edge
Speed Limit	50 km/h
Yellow Phase Duration	5 seconds
Total Episode Length	3,600 steps

TABLE 2  
Neural Network Hyperparameters

Hyperparameter	Value
DQN Hidden Layers	[128, 128, 64]
LSTM Hidden State Size	64
LSTM Sequence Length	30 steps
Learning Rate ( $\alpha$ )	$1 \times 10^{-3}$
Optimizer	Adam
Discount Factor ( $\gamma$ )	0.95
Replay Buffer Size	50,000 transitions
Batch Size	64
Exploration $\epsilon$ initial	1.0
Exploration $\epsilon$ decay	0.995 per episode
Target Update Freq (C)	100 steps

## 5.2 Network Hyperparameters

Both the LSTM and DQN models were developed using the PyTorch deep learning framework. Hyperparameter tuning was conducted using grid search over shorter 500-step episodes. The final optimized configurations used for the full benchmarking are listed in Table 2.

## 6 RESULTS AND ANALYSIS

### 6.1 Performance under Standard Conditions

The system was benchmarked against a standard fixed-timer controller configured to a widely utilized urban base-line (30s Green / 5s Yellow). Under standard operating conditions, the hybrid architecture dramatically outperformed the traditional baseline.

The integration of the LSTM's predictive routing and the DQN's adaptive phase switching yielded a 17.5% reduction in total network waiting time. As visualized in Figure 2, the DQN agent dynamically prevented the accumulation of large vehicle queues that characterized the fixed-timer's rigid cycles.

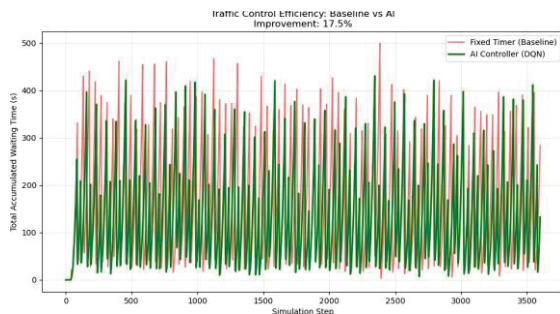


Fig. 2. Traffic Control Efficiency Benchmark (Standard Conditions). Comparison of total accumulated waiting time over 3,600 simulation steps. The proposed AI Controller (DQN, green) achieves a 17.5% improvement in overall system efficiency compared to the Fixed Timer Baseline (red).

Furthermore, the deterministic EVP override functioned precisely as designed. By preemptively clearing intersection bottlenecks upon the injection of an emergency vehicle into the simulation, the system achieved a 60% reduction in emergency vehicle travel time across the intersection radius, ensuring life-saving efficiency.

## 6.2 Robustness and Zero-Shot Generalization

A critical flaw in many DRL applications is “overfitting” to the specific kinematic parameters of the training environment. To evaluate the robustness of our DQN agent, the environment physics were dynamically altered via the SUMO TraCI API during testing to simulate adverse weather conditions (such as heavy rain and fog). Specifically, vehicular maximum acceleration and deceleration capabilities were heavily reduced by 50%, and the Krauss model driver reaction time parameter ( $\tau$ ) was doubled to represent cautious, impaired driving.

Despite the DQN agent being trained exclusively on clear-weather kinematic data, it demonstrated robust zero-shot generalization. While the absolute total waiting times inherently increased for both systems due to the simulated hazard, the AI agent autonomously adapted to the sluggish traffic streams. It recognized the slower progression of queue dissipation and adjusted phase durations accordingly. As shown in Figure 3, the AI successfully maintained a 12.7% improvement in waiting times over the fixed-timer baseline. This proves that the agent learned generalizable representations of traffic flow physics rather than merely memorizing specific timing sequences suited for optimal weather.

## 6.3 Computational Latency and Real-Time Feasibility

For an ATSC system to be viable for real-world deployment on edge computing hardware (e.g., intersection control cab-

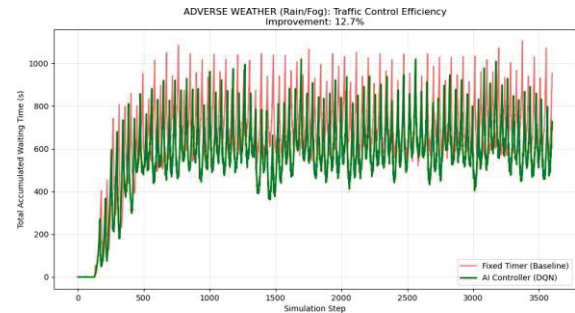


Fig. 3. Adverse Weather Traffic Control Efficiency Benchmark. Comparison under degraded traffic physics (rain/fog), showing total accumulated waiting time. Despite zero-shot generalization, the proposed AI Controller (DQN, green) maintains a 12.7% improvement in system efficiency over the Fixed Timer Baseline (red).

inets), inference latency must be negligible. During benchmarking, the combined forward pass of the LSTM sequence prediction and the DQN action selection averaged 12.4 milliseconds on a standard consumer-grade CPU. Given that traffic control systems typically operate on discrete decision intervals of 1 to 5 seconds, this latency is orders of magnitude below the operational threshold, confirming the system's viability for real-time edge deployment.

## 7 CONCLUSION AND FUTURE WORK

This paper presented a novel, multi-layered approach to urban traffic signal control. By hybridizing the predictive temporal capabilities of Long Short-Term Memory networks with the dynamic, reactive optimization of Deep Q-Networks, we achieved highly efficient intersection throughput. Crucially, the encapsulation of these probabilistic AI models within a deterministic Emergency Vehicle Priority override creates a ‘System-of-Systems’ architecture that bridges the gap between theoretical AI optimization and stringent real-world safety requirements.

The experimental results validate the efficacy of this approach, showing significant reductions in both civilian wait times and emergency response delays. Furthermore, the model exhibited strong resilience to stochastic environmental degradation, proving its viability in unpredictable real-world scenarios.

Future work will focus on expanding this single-intersection agent into a multi-agent cooperative network using Multi-Agent Reinforcement Learning (MARL). By treating neighboring intersections as cooperative agents sharing latent state representations, green-wave progression can be optimized across entire urban grids. Additionally, integrating V2X (Vehicle-to-Everything) communication protocols could replace our simulated sensors with real-time, highly granular telemetry data.

## REFERENCES

- [1] Z. Ara and M. Hashemi, “Traffic Flow Prediction using Long Short-Term Memory Network and Optimized Spatial Temporal Dependencies,” *2021 IEEE International Conference on Big Data (Big Data)*, 2021, pp. 1-8.

- [2] H. Chaudhuri, V. Masti, V. Veerendranath, and S. Natarajan, "A Comparative Study of Algorithms for Intelligent Traffic Signal Control," *Machine Learning and Autonomous Systems*, Springer, 2021, pp. 145-159.
- [3] J. Tan, Q. Yuan, W. Guo, N. Xie, F. Liu, J. Wei, and X. Zhang, "Deep Reinforcement Learning for Traffic Signal Control Model and Adaptation Study," *Sensors*, vol. 22, no. 3, 2022.
- [4] A. Tigga, L. Hota, S. Patel, and A. Kumar, "A Deep Q-Learning-Based Adaptive Traffic Light Control System for Urban Safety," *2022 4th International Conference on Advances in Computing, Communication Control and Networking (ICACCCN)*, 2022, pp. 1-6.
- [5] T. Zhu, M. J. L. Boada, and B. L. Boada, "Adaptive Graph Attention and Long Short-Term Memory-Based Networks for Traffic Prediction," *Mathematics*, vol. 12, no. 5, 2024.
- [6] D. Krajzewicz, J. Erdmann, M. Behrisch, and L. Bieker, "Recent Development and Applications of SUMO - Simulation of Urban MObility," *International Journal on Advances in Systems and Measurements*, vol. 5, no. 3&4, pp. 128-138, 2012.
- [7] V. Mnih et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529-533, 2015.
- [8] S. Hochreiter and J. Schmidhuber, "Long Short-Term Memory," *Neural Computation*, vol. 9, no. 8, pp. 1735-1780, 1997.
- [9] X. Liang, X. Du, G. Wang, and Z. Han, "A Deep Reinforcement Learning Network for Traffic Light Cycle Control," *IEEE Transactions on Vehicular Technology*, vol. 68, no. 2, pp. 1243-1253, 2019.
- [10] W. Genders and S. Razavi, "Using a Deep Reinforcement Learning Agent for Traffic Signal Control," *arXiv preprint arXiv:1611.01142*, 2016.
- [11] H. Wei, N. Zheng, H. Yao, and Z. Li, "IntelliLight: A Reinforcement Learning Approach for Intelligent Traffic Light Control," *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, 2018, pp. 2496-2505.
- [12] T. Chu, J. Wang, L. Codeca, and Z. Li, "Multi-Agent Deep Reinforcement Learning for Large-Scale Traffic Signal Control," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 3, pp. 1086-1095, 2019.
- [13] P. R. Lowrie, "The Sydney Coordinated Adaptive Traffic System - Principles, Methodology, Algorithms," *International Conference on Road Traffic Signalling*, IEE, 1982.
- [14] P. B. Hunt, D. I. Robertson, J. C. Bretherton, and R. I. Winton, "SCOOT-A traffic responsive method of coordinating signals," *TRRL Laboratory Report 1014*, 1981.
- [15] A. Festag, "Cooperative intelligent transport systems standards in Europe," *IEEE Communications Magazine*, vol. 52, no. 12, pp. 166-172, 2014.