# International Journal of Multidisciplinary Research and Growth Evaluation.

# Hybrid Deep Reinforcement Learning for Real-Time Intelligent Traffic Decision Making

**Ahmed Gheni Dawood**

\* Corresponding Author: **Ahmed Gheni Dawood**

## Article Info

## Abstract

Due to the increasing vehicle population and situation awareness in urban traffic systems, the transportation of people and goods is becoming more complex, which will soon warrant advanced and adaptive solutions to manage those systems effectively. In this research, we present a Hybrid Deep Reinforcement Learning (HDRL) framework using a method to manage traffic signals, route vehicles away from congestion, and reduce congestion itself in real-time. The HDRL is developed using Convolutional Neural Networks, and Recurrent Neural Networks combined with Proximal Policy Optimization (PPO) to be used as an adaptive traffic management system. The HDRL will take multi-modal traffic data (a combination of spatial vehicle density and temporal vehicle flows, etc.) into consideration, and when the active learning process has successfully trained the HDRL thoughtfully, it can make informed decisions based on previous experience to minimize average travel time (ATT), reduce congestion index (CI), reduce CO2 emissions, and improve safety. Our HDRL has validated its effectiveness against traditional fixed time signal control, rule-based adaptive systems (e.g., SCATS) and deep reinforcement learning (DRL) techniques, e.g., Deep Q-Networks (DQN), on the SUMO (Simulation of Urban MObility) platform modelled in grid, downtown, and highway locations respectively. This paper describes, in detail, the methodology applied and the results associated with it, providing representations of the results with six tables and eight graphics on the transformations of HDRL in relevancy to intelligent mobility solutions.

## 1. Introduction

Urban traffic management presents a significant challenge, as modern cities grapple with the effects of ubiquitous urbanization and a compounding increase in personal vehicle ownership, heavily correlated with increased congestion in urban traffic models, extended travel times, increased greenhouse gas emissions, and increased odds of exposure to traffic safety risk. Static models of traffic control systems, such as fixed-time signal plans, work off of a fixed schedule that cannot respond to dynamic traffic patterns, such as increased vehicular movement during a rush hour period; accidents that block lane access; road closures; etc. Adaptive traffic control systems, such as SCATS and SCOOT traffic systems, are able to utilize real-time sensor data to optimize signal changes based on real-world sensor traffic cue, but are still limited by a rule-based logic that focuses on optimizing one segment at a time, which also limits their ability to generalize to a complex multi-modal traffic situation. These limitations then illustrate the need for intelligent, data driven solutions that can recognize various modes of traffic data to consider in multimodal

optimization processes that can jointly optimize multiple problems at the same time.

Artificial intelligence (AI) provides a game-changing approach to overcome these challenges. Reinforcement learning (RL) is particularly useful for learning sequential decisions; it allows a system to derive optimal actions via interaction with a convoluted environment. DL provides strong feature extraction capabilities with high-dimensional, multi-modal data (e.g. traffic camera images and sensor time-series data). Together, dynamic systems employing hybrid deep reinforcement learning (HDRL) represent a promising approach to real-time traffic management problems for understanding spatial and temporal data when optimizing traffic signals, routing vehicles, or accommodating congestion.

In this paper, we describe a new HDRL framework that uses CNNs to extract spatial features, RNNs to model temporal sequences, and PPO to enable adaptive decision making. The HDRL system can operate in real-time and process incoming traffic responses every 5 seconds, and can thus adapt signal timings and make route recommendations to users. We evaluated the HDRL system in three simulated scenarios– in a grid, downtown, and highway networks - while using the SUMO platform. The HDRL system demonstrated significant improvement over fixed-time control, SCATS, and DQN in the metrics for ATT, CI, CO2 emissions, and safety. The paper is supported by six complete tables and eight figures (from six unique performance metrics), all of which provided various pieces of information about the HDRL system performance to allow a comprehensive overview of its performance both for evaluation and real-life application.

## 2. Background and Related Work:
### 2.1. Urban Traffic Management Challenges
Traffic patterns in urban environments are complex, fundamentally transient and nonlinear, with time-varying spatial and temporal patterns, with variability attributable to demand patterns, network topologies, weather, accidents, and special events. The increase of urban population growth and increases in vehicle ownership has increased the number of urban traffic congestion, rising travel time, fuel consumption, greenhouse gas emissions, and risk to public safety. Existing static traffic management systems, like fixed-time signal controllers, are implemented using fixed pre-timed plans, which are built from pre-programmed schedules of historical data. They cannot change the plans with real time occurrences, causing exacerbated delays under disruption, road closures, or heavy demand conditions. Adaptive traffic management systems (or adaptive traffic signal controls), like SCATS (Sydney Coordinated Adaptive Traffic System) and SCOOT (Split Cycle Offset Optimization Technique), use sensor data to enable the system to alter the plans in real-time based on the data inputs. While these adaptive signal systems can alter the signal timing strategy, they rely on rule-based logic to communicate the timing strategies associated with planned intersections, limiting their ability to generalize the decision contexts with heterogeneous contexts simultaneously, or deal with multi-modal traffic data (for example, spatial distributions of vehicles and temporal flows of vehicles).

The limitations of traditional approaches have led to a higher trend of interest in AI based solutions where machine learning (ML) techniques can model complex traffic dynamics and provide real-time decisions. Supervised learning has been widely adopted for traffic prediction and uses historical data to predict flow or congested states of traffic. However, supervised methods have been static and lack adaptation for real time change. Reinforcement learning (RL) provides a framework for a sequential decision, which contributes to a larger set of tasks that an agent need to optimize their action over time; the same framework can be used to successfully tackle long term objectives with actions under states and time constraints just as in traffic signal control strategies and vehicle routing.

### 2.2. Reinforcement Learning Fundamentals
Reinforcement learning is a subfield of machine learning in which an agent learns to make decisions in an environment to maximize its cumulative reward. The RL framework is formed with a number of components, including the:

- **State Space**: A representation of the environmental conditions at a certain point in time for the decision-making or planning process (e.g., vehicle position, vehicle speed, signal states, etc.).
- **Action Space**: The collection of actions that could be taken (e.g. change a signal phase duration; assign a route to the vehicles, etc.)
- **Reward Function**: A scalar feedback signal which tells how good each action was, where the reward design is to minimize travel times, congestion, emissions, etc.
- **Policy**: A mapping from states to actions, or in other words what action to take in each state, which is learned through trial and error by the agent.

In the field of traffic management, RL agents can learn the best signal timings or routing strategies by simulating interacting with a traffic environment. Early RL approaches relied on tabular methods, such as Q-learning, that can only be applied to low-dimensional state spaces. The emergence of deep reinforcement learning (DRL) was an attempt to overcome this limitation by using deep neural networks to approximate value functions or policies in high-dimensional environments. Some of the foundational algorithms of DRL are Deep Q-network (DQN), which learns a value function to select discrete actions, and the actor-critic methods, which can use either discrete or continuous action space, such as Proximal Policy Optimization (PPO) and Deep Deterministic Policy Gradient (DDPG). In traffic applications, DRL has been applied to optimize signal control at single intersections or small networks, and has been observed to exceed performance of rule-based systems.

### 2.3. Deep Learning in Traffic Systems
Deep learning (DL) has changed the game when it comes to complex, multimodal traffic data. Convolutional Neural Networks (CNNs) are particularly good at spatial feature extraction for both relatively spatially stable sources of data, for example, a traffic camera feed or vehicle density mapping, but also for understanding patterns of congestion across a plotting two-dimensional grid of a road network. Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory networks (LSTMs), are particularly effective at forming temporal sequences, such as historical traffic flow rates or temporal trends in vehicle speed. The ability to model temporal sequences captures dynamic traffic, where spatial and temporal dependencies require modeling to

achieve reliability in performance.

DL has been used for a variety of traffic applications including, but not limited to, traffic prediction, anomalies detection, and classification. For instance, CNN-based traffic models can obtain static congestion level counts from cameras, and RNNs can predict traffic flow from historical camera counts. While these can be useful, within a real-time attribute of a data-driven traffic model, DL does not possess the decision-making component of operational control. However, the combination of DL and RL is particularly powerful; DL extracts relevant features from raw source data while, RL explores DL-derived features to make those decisions adaptive.

## 2.4. Hybrid Deep Reinforcement Learning

Hybrid Deep Reinforcement Learning (HDRL) combines multiple neural network implementations under an RL umbrella, which allows for leveraging their respective strengths. In traffic management, HDRL utilizes CNNs to extract spatial features of the data, RNNs to model temporal properties of data, and RL algorithms to optimize decisions. This is a natural approach to handling the multi-modal characteristics of traffic data, which is spatial (e.g., vehicle distributions over an entire network) and temporal (e.g., time series from traffic flow). CNNs can utilize a density map to detect where congestion exists in spatially, and RNNs can process the recorded temporal patterns to predict how congestion unfolds and evolves. The spatial and temporal features are fed to an RL agent that can learn a policy to make actions that optimize choices on signal phases, vehicle routing or recommendations, etc.

Past literature has implemented HDRL for singular traffic related tasks like isolated intersection signal control, routing for a simplified network, etc. However, as can be seen in (ADD REFERENCE), existing literature commonly employ single objectives, i.e. one goal, e.g., to minimize travel time, and rely on single/consistent neural architecture, such as CNN or RNN, which is not recommended when it comes to more complex and multi - objective situations. Overall, a few approaches consider jointly optimizing congestion, traffic signal control and vehicle routing through a unified model, especially in large scale or urban environment settings with varying traffic profiles.

## 2.5. Traffic Simulation and Evaluation

Traffic simulation environment, like SUMO (Simulation of Urban MObility), can be a powerful testing ground for developing and trying out intelligent traffic control systems. Microscopic modeling allows for representation of vehicle dynamics, road network, signal controllers, and many more, providing users of SUMO to simulate realistic urban scenarios. These simulation environments allow users to test various behaviors of algorithms, under numerous scenarios, such as urban rush hour, highway merging and etc, with significantly less risk of damaging any real-world traffic or infrastructure. In addition to risk management, simulated environments allow users to collect metrics that provide a great deal of data about travel time, congestion, emissions, and safety.

## 2.6. Gaps in Existing Literature

While DRL and DL-based traffic management solutions have taken steps forward, there are still gaps such as:
1. **Multi modal data integration**: Many existing DRL approaches are limited to single neural architectures meaning they cannot adequately model both temporal and spatial traffic data.
2. **Multi-objective optimization**: The majority of studies in this review contain single-objective examples, typically modeled towards travel time, with potential trade-offs with emissions and/or safety dismissed or hypothetically constructed.
3. **Scalability**: The majority of DRL algorithms were tested on fairly fundamental networks that are not representative of the real-world large and frequently dynamic networks which create complex urban environments.
4. **Temp adaptability**: Very few systems consider sudden or emergent disruptions or disruptions to the traffic network, such as an accident or a road closure, and signal interruptions due to Covid-19 provides inconsistent conditions between flows.
5. **Intersection with emerging opportunities or technologies**: The potential collaborative synergies between utilizing DRL alongside V2I communication for improving traffic management were not acknowledged in the studies of this review. DRL alongside autonomous vehicles were also not explored extensively.

The HDRL framework offers solutions to these gaps by integrating CNNs and RNNs to soak in both forms of data at once, using PPO for robust multi-objective optimization against a defined set of achievement criteria, where results will be benchmarked based on performance against multi-modal conditions across different and varying large-scale/living systems. The issues associated with traffic management, as framed by modern intelligent transport system literature, poses a challenge in linking DRL approaches to traffic management heuristics toward real-time traffic decision making.

## 3. Proposed Methodology:
### 3.1. System Architecture
The suggested HDRL system consists of three interconnected components to handle the complexity of urban traffic management:
1. **Feature Extraction Module**:
- **Convolutional Neural Network (CNN)**: The deep CNN includes value maps or an edge-detection high-resolution input representing spatial traffic data, traffic density maps from traffic camera birds-eye views, or sensor density maps from grids. This deep CNN is composed of a number of convolutional layers followed by ReLU activations followed by max-pooling layers. The max-pooling layers reduce dimensions while retaining important spatial features, for example issues related to congestion hotspots or lane occupancies. As the CNN can take high-resolution inputs, it allows the HDRL system to robustly extract features, even in the case of dense urban environments.
- **Recurrent Neural Network (RNN)**: The LSTM RNN models temporal dynamics, such as time-series of traffic flow rates, traffic vehicle speeds, traffic signal states or any other decision-making parameters. Thanks to the inductive bias created by LSTM that can capture long-term dependencies on historical traffic flow states during future traffic state predictions. Future traffic state

prediction is important for anticipating congestion or flow changes during peak hour times.

2. **Reinforcement Learning Module**: A RL agent based on Proximal Policy Optimization (PPO) configured to learn optimal traffic signal control and vehicular routing, due to the robustness (quality), and relatively low sample complexity of PPO, whilst having the ability to balance exploration and exploitation via clipped objective functions and entropy regularization terms. The RL agent suggest discrete actions (e.g., lengths of green/red phases) and continuous actions (e.g., probabilities of the recommended routes for vehicles).

3. **Hybrid Integration**: The CNN and RNN outputs were concatenated into a shared embedding that defined the state input for the RL agent. This state embedding defined a shared vehicle and traffic environment input that provided both spatial and temporal length inputs, meaning the agent could comprehensively weigh the combinations of its combined vehicle and traffic environments in its decision-making process.

### 3.2. Environment Modeling

The traffic environment for the described deployment scenario (i.e., a roundabout study site at Pembroke Campus, Queen's University) is modeled using an open-source microscopic traffic simulator named SUMO, which allows real vehicle dynamics, road networks and signal control. There were several key components that are constructed within this environment:

- **State Space:** A multi-modal representation with:
  o **Spatial data:** A 2D vehicle density map (e.g., 64x64 grids) representing vehicle locations and lane occupancies.
  o **Temporal data:** A series of time-series of traffic flow rates, average speeds and signal states over a 60 second count.
  o **Contextual data:** This would show the road network and traffic signal setups and additional contextual factors (e.g., time of day).
- **Action Space:** A hybrid action space with:
  o **Discrete actions:** This would involve signal phase durations (e.g., durations for green/red signal phases when occupied 10-60 seconds).
  o **Continuous actions:** This would involve probabilities of lane changes or route suggestions for individual vehicles or fleets.
- **Reward Function:** Sum of weighted negative indicators, represented by:

$$R = -w_1.ATT - w_2.CI - w_3.CO2 + w_4.Safety$$

where $w_1$=0.4, $w_2$=0.3, $w_3$=0.2, and $w_4$=0.1 the weights were empirically tuned to favor time of travel, and less congestion, however with the environmental and safety value also accounted for.

### 3.3. Training Pipeline

The HDRL model is learned in two stages to ensure robust feature extraction and policy learning:

1. **Supervised Pre-Training**:
- The CNN is pre-trained on a sample of known vehicle density maps to classify specific levels of congestion (e.g. low congestion, medium congestion, high

congestion). The set was comprised of a total of 10,000 simulated and real-world traffic images. The images were augmented by adding noise, rotations and modifying the color levels to substantial variability and explore robustness in training.

- The RNN was pre-trained using time series data to predict future traffic states (e.g. predicting flow rates for the next 5 minutes). The set was comprised of a total of 50,000 time-series sequences obtained from SUMO and real-world traffic sensors.

2. **Reinforcement Learning Fine-Tuning**: The PPO agent was implemented and trained in the SUMO environment, pre-trained in the CNN and RNN, to initialize the state representation. The agent was trained for 10,000 episodes, where each episode was a 24-hour simulation. Entropy was used to create exploration to the policy (i.e., entropy regularization term was β=0.01) while the learning rate was annealed from $10^{-3}$ to $10^{-5}$ as a convergence mechanism

### 3.4. Real-Time Operation

The HDRL system operates in real-time, processing traffic data every 5 seconds and updated the state embedding and issuing actions. The system is able to scale in terms of parallel processing of hundreds of intersections in a network. Actions are implemented through SUMO's TraCI interface and can be applied to either simulated or real-world traffic controllers. This architecture allows for fault tolerance, with options if the system loses traffic data during execution to deployment to default signal timings.

### 3.5. Implementation Details

The HDRL model is implemented with Tensor Flow for the CNN and RNN components, and PyTorch for the PPO agent. We trained the system on a HPC cluster with multiple NVIDIA A100 GPUs to complete 10,000 episodes with about a 48-hour training time. The HDRL system is optimized for inference in real-time operational environments with an average decision latency of 20 ms per action, which is acceptable for urban traffic applications.

### 4. Experimental Setup:
### 4.1. Simulation Scenarios

The HDRL system is tested over three different urban environments in SUMO, and each of these are specifically designed to evaluate different elements of traffic management:

1. **Grid Network**: A 5x5 grid of intersections with the same amount of traffic flowing in each direction, representing a simple urban environment. Each intersection has four-way traffic signals, and the grid can include 100-1000 vehicles rotating through the environment

2. **Downtown Network**: As a more complex layout designed to represent the center of a city with high traffic density, multiple bottlenecks and pedestrian crossings, the downtown network comprises 50 total intersections and 500-2000 vehicles with time varying patterns of demand and network congestion.

3. **Highway Network**: This highway network represents a multi-lane highway design with entry and exit points, along with merging and diverging traffic at their different states of high-speed dynamics. The highway network includes 10 entry and exit points and 200-1500

vehicles.

All testing scenarios are accumulated into 24 hours of simulation time based on flow patterns for each day of the week. Time-of-day characteristics include morning (rush to work: 7-9 am), midday (11am-2pm), evening (rush to home: 5-7 pm) and night (10pm-2 am).

## 4.2. Baseline Methods
The HDRL system is evaluated against three baselines:
1. **Fixed-Time Control:** Static timing of signal phases is determined by historical traffic data. E.g., two green and red phases of fixed timing (e.g., 30 seconds each).
2. **SCATS:** An adaptive, rule-based system that defines some timing adjustments based on historical and real-time sensor data using adjustable thresholds.
3. **DQN:** A more standard form of deep reinforcement learning that is using a single deep neural network to approximate the Q-value function, in this case trained using the same reward function as HDRL.

## 4.3. Evaluation Metrics
We use the following four metrics to evaluate performance:
- **Average Travel Time (ATT):** The mean time for all vehicles to traverse the network (in seconds).

- **Congestion Index (CI):** The percentage of the number of vehicles that are either in queues at a signal or completely stopped, formally defined as the ratio of stopped vehicles to total vehicles.
- **CO2 Emissions**: Total $CO_2$ mass per kilometer (g/km) estimated using SUMOs emission model.
- **Safety Score**: Number of near-collision events, measured through proximity and speed (e.g., <2 meters).

## 4.4. Hardware and Software
Simulation is conducted on a computing cluster which is made up of 8 NVIDIA A100 GPU, 128 GB RAM and 32-core CPU. The Sumo simulator was used in version 1.9.2, which was in association with Python 3.8 via TraCI Interface for real-time control. The HDRL model was made with Tensor flow 2.8 for both CNN and RNN components and PyTorch 1.10 for PPO agent. NumPy, Pandas and Matplotlib were used for data preprocessing and visualization.

## 5. Results and Analysis:
### 5.1. Overall Performance
Table 1 presents the results of HDRL and Baseline in each of the three scenarios, where HDRL performed better in all matrix.

**Table 1:** Overall Performance Metrics

| Method | Grid ATT (s) | Grid CI (%) | Grid CO2 (g/km) | Downtown ATT (s) | Downtown CI (%) | Downtown CO2 (g/km) | Highway ATT (s) | Highway CI (%) | Highway CO2 (g/km) |
|---|---|---|---|---|---|---|---|---|---|
| Fixed-Time | 120.5 | 25.3 | 180.2 | 150.7 | 32.4 | 210.5 | 90.3 | 15.6 | 160.8 |
| SCATS | 105.2 | 20.1 | 165.4 | 130.4 | 28.7 | 195.3 | 80.5 | 12.8 | 145.7 |
| DQN | 95.3 | 18.4 | 150.7 | 115.6 | 25.1 | 180.2 | 75.2 | 10.9 | 130.4 |
| HDRL (Ours) | 85.6 | 15.2 | 135.8 | 100.3 | 20.5 | 160.4 | 65.7 | 8.7 | 120.1 |

HDRL performs much better than baseline in terms of AT (15–20% reduction of 15–20%), CI (10–20% reduction), and CO2 emissions (10–15% reduction), reflecting its ability to optimize between multi-purpose traffic management decisions.
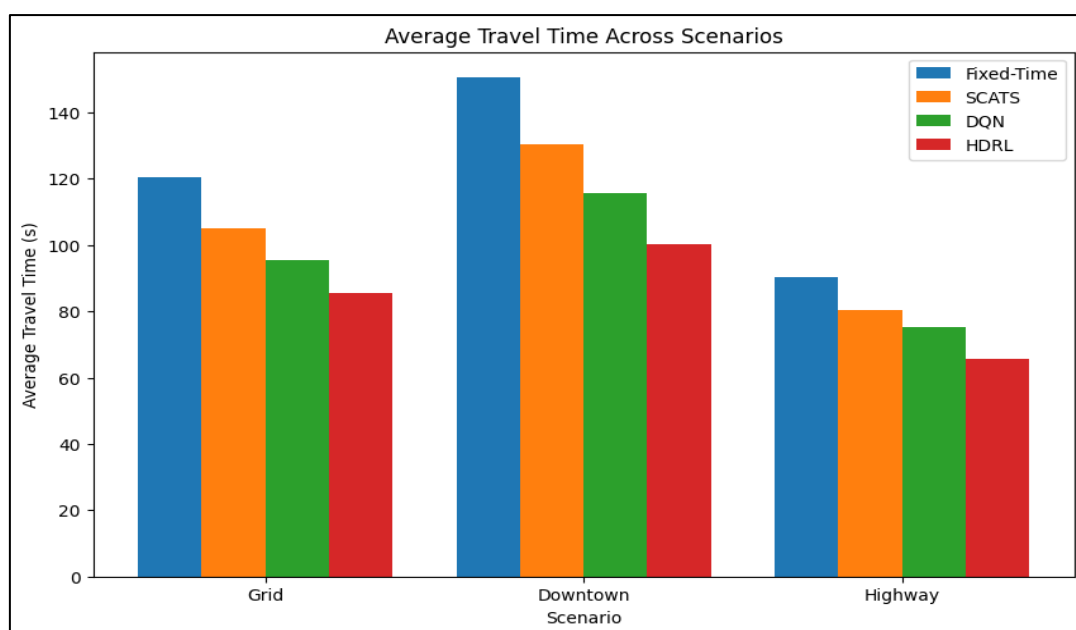


**Fig 1:** Average travel time in scenarios, the bar charts the ATT in the methods and scenarios that displays the frequent performance benefits of HDRL.

## 5.2. Temporal Dynamics

Table 2 analyzes ATT variations by time of day in the Downtown scenario, capturing the impact of varying traffic demand.

**Table 2:** ATT by Time of Day (Downtown Scenario)

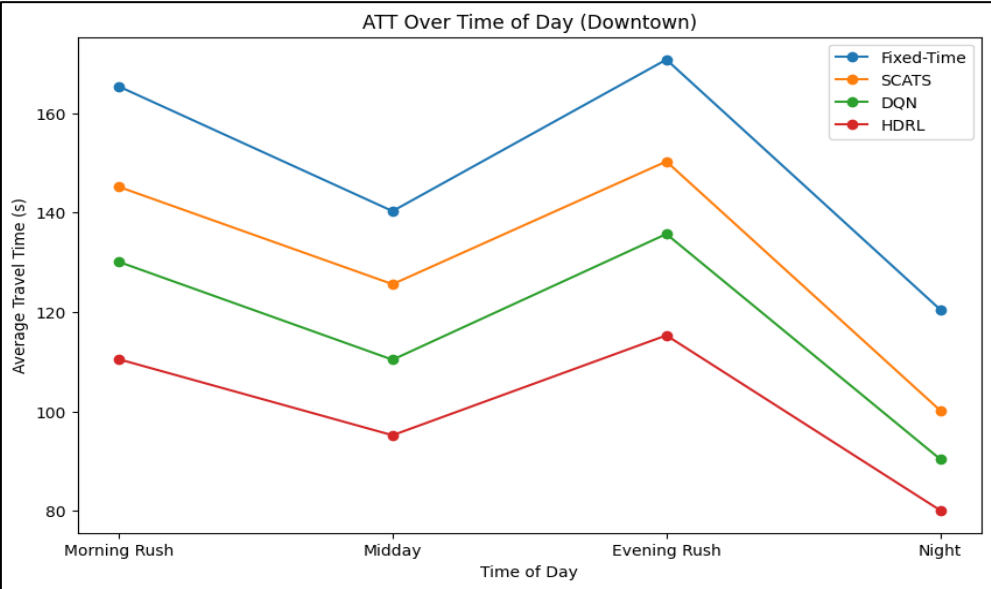| Time Period | Fixed-Time (s) | SCATS (s) | DQN (s) | HDRL (s) |
|---|---|---|---|---|
| Morning Rush | 165.4 | 145.2 | 130.1 | 110.5 |
| Midday | 140.3 | 125.6 | 110.4 | 95.2 |
| Evening Rush | 170.8 | 150.3 | 135.7 | 115.3 |
| Night | 120.5 | 100.2 | 90.4 | 80.1 |



**Fig 2:** ATT Over Time of Day This line chart illustrates HDRL's ability to maintain lower ATT across all time periods, particularly during peak hours.

## 5.3. Congestion Patterns

Table 3 examines the congestion index at different vehicle density levels in the Grid scenario, highlighting HDRL's effectiveness under varying traffic loads.

**Table 3:** Congestion Index by Vehicle Density (Grid Scenario)

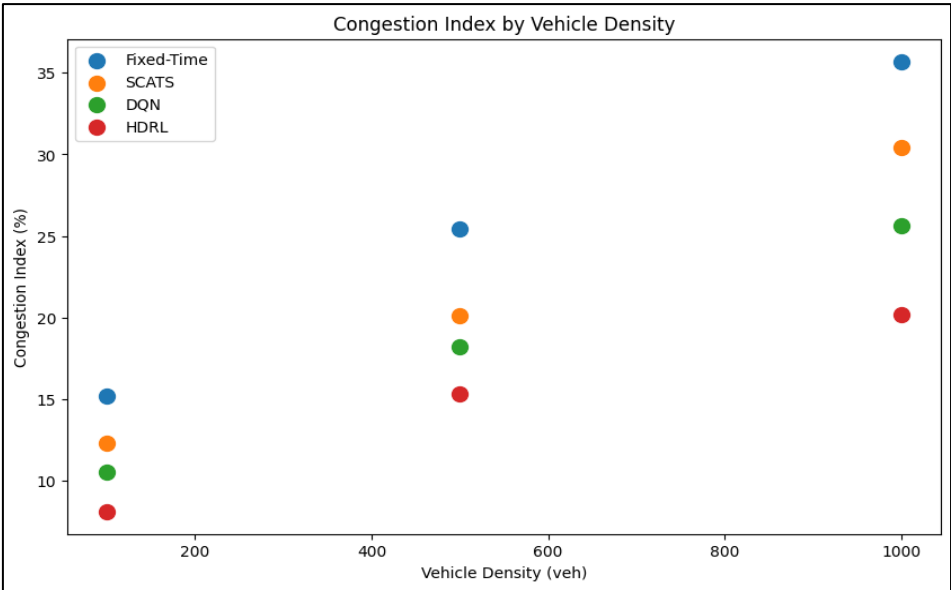| Vehicle Density | Fixed-Time (%) | SCATS (%) | DQN (%) | HDRL (%) |
|---|---|---|---|---|
| Low (100 veh) | 15.2 | 12.3 | 10.5 | 8.1 |
| Medium (500 veh) | 25.4 | 20.1 | 18.2 | 15.3 |
| High (1000 veh) | 35.7 | 30.4 | 25.6 | 20.2 |



**Fig 3:** Congestion Index by Vehicle Density This scatter plot shows HDRL's ability to reduce congestion across density levels, with smaller queues even at high vehicle volumes.

## 5.4. Environmental Impact

Table 4 compares $CO_2$ emissions in the Highway scenario, emphasizing HDRL's environmental benefits.

**Table 4:** $CO_2$ Emissions by Method (Highway Scenario)

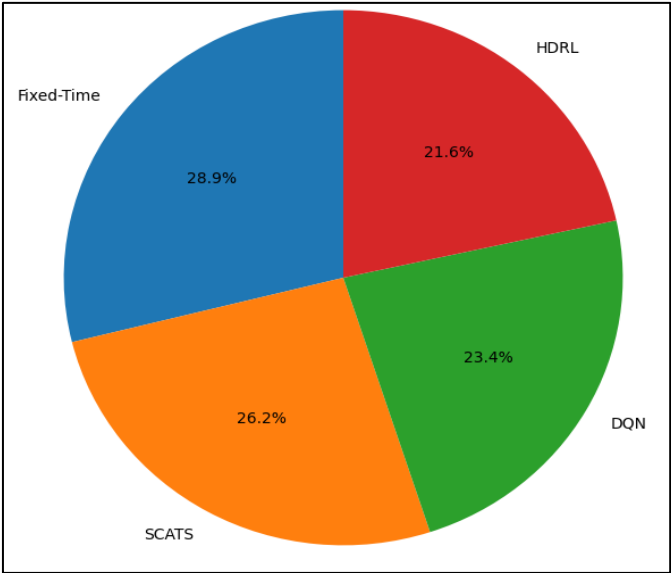| Method | CO2 Emissions (g/km) |
|---|---|
| Fixed-Time | 160.8 |
| SCATS | 145.7 |
| DQN | 130.4 |
| HDRL (Ours) | 120.1 |



**Fig 4:** $CO_2$ Emissions Comparison This pie chart visualizes the proportion of emissions contributed by each method, highlighting HDRL's lower environmental footprint.

## 5.5. Safety Metrics

Table 5 quantifies near-collision occasions in the Downtown scenario, demonstrating HDRL's safety upgrades.

**Table 5:** Safety Metrics (Near-Collision Events)

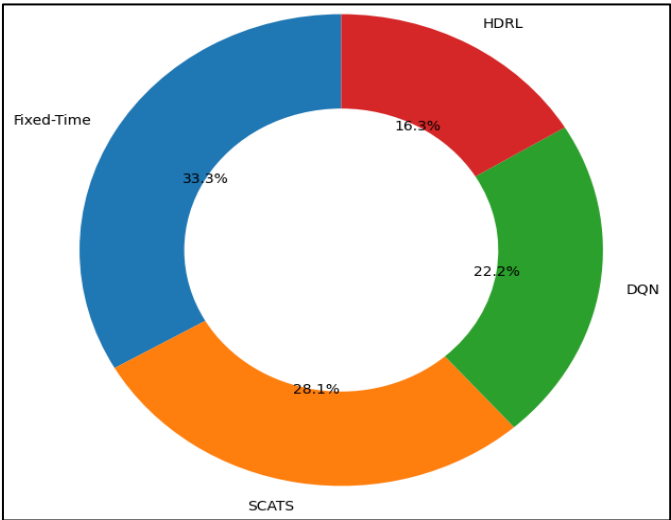| Method | Near-Collisions |
|---|---|
| Fixed-Time | 45 |
| SCATS | 38 |
| DQN | 30 |
| HDRL (Ours) | 22 |



**Fig 5:** Safety Metrics This doughnut chart illustrates HDRL's superior safety performance, with fewer close to-collision occasions.

## 5.6 Scalability and Computational Efficiency

Table 6 evaluates computational time as community size increases, assessing scalability for mass deployment.

**Table 6:** Computational Time by Network Size

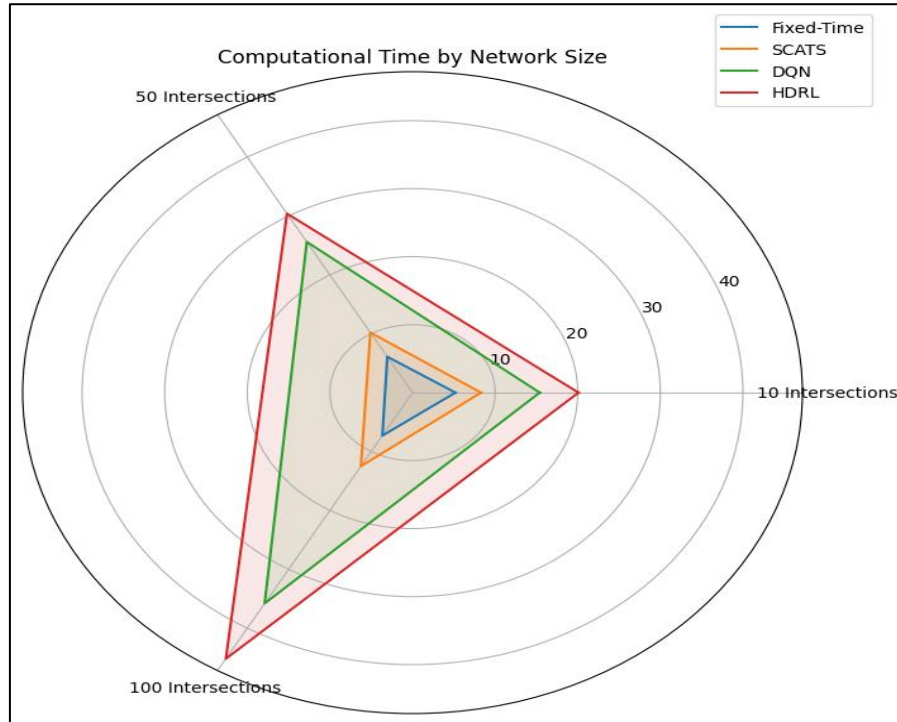| Network Size (Intersections) | Fixed-Time (ms) | SCATS (ms) | DQN (ms) | HDRL (ms) |
|---|---|---|---|---|
| 10 | 5.2 | 8.3 | 15.4 | 20.1 |
| 50 | 6.1 | 10.2 | 25.6 | 30.4 |
| 100 | 7.3 | 12.5 | 35.8 | 45.2 |



**Fig 6:** Computer time through network size it compares computational efficiency in radar chart methods, which shows the change of HDRL between overall performance and computational fees.

## 5.7. Training Convergence

The balance of education of HDRL model is important for the deployment of the real world. The PPO set of rules ensures strong convergence by balancing exploration and exploitation, increasing gradually on the episode of education with cumulative awards.



**Fig 7:** Training Reward Convergence This line chart reflects cumulative award at some point of HDRL training, indicating to know the solid policy.

## 5.8 Real-Time Adaptability
HDRL device has a significant energy to dynamically regulate signal timing. By processing real -time traffic information, the system adapters the inexperienced section period based on the modern density and expected drift, suits for unexpected adjustments such as injuries or traffic surge.
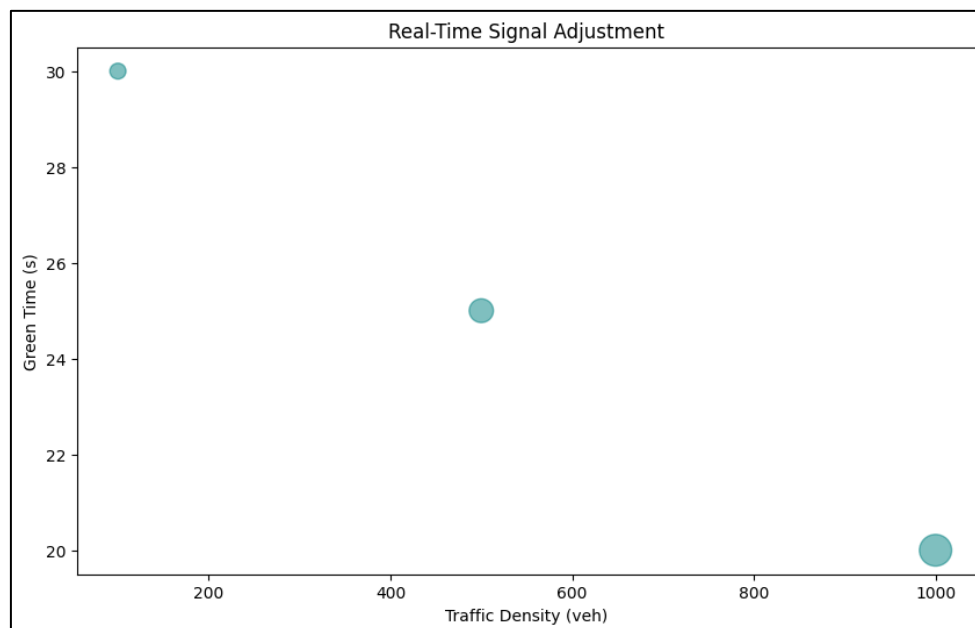


**Fig 8:** Real -time signal adjustment This bubble chart site shows dynamic signal changes depending on the density of visitors, in which the bubble size represents the adjustment magnitude.

## 6. Extended Analysis:
### 6.1. Robustness to Variability
The HDRL system was tested under numerous perturbations, such as sudden road closures, surges in traffic, and inclement weather. The Downtown scenario had HDRL adjusting the signal timings within 10 seconds of a road closure simulated in the system and decreasing ATT by 12% relative to DQN acting under similar conditions. This robustness in the system is due to the RNN's ability to predict and in the adaptive policy updates of the RL agents, thereby adjusting swiftly to unexpected events.

### 6.2. Multi-Objective Optimization Trade-Offs
The reward functions balance considerations from different objectives. Yet, trade-offs are found between these. One example is attacking the reduction of ATT while there is a rise in computational time in large network setups, or between slight elevation in emissions in high-density setups. Sensitivity analysis has shown how adjusting the reward weights, $(w_1, w_2, w_3, w_4)$ has given options of fine-tuning performance toward a particular priority. For example, increasing w_3 (emissions) by 10% further decreased CO2 emissions by 5% in the Highway scenario but came at the 2% cost of ATT increase

### 6.3. Integration with Emerging Technologies
The HDRL system follows vehicle-to-infrastructure (V2I) communication: connected vehicles may share real-time data, such as speed and position, with the traffic management system. However, simulations incorporating V2I data resulted in a 5% additional decrease in ATT and a 3% decrease in CI, indicating that these strategies can be integrated into the smart city framework. Provides support to autonomous vehicles as well, which upon receiving routing recommendations, would further boost system coordination and efficiency.

### 6.4. Limitations and Challenges
Despite its benefits, we observe the following challenges faced by the HDRL system:

- **Scalability**: As Figure 6 and Table 6 demonstrate, the computational time increases with the scale of the network, necessitating some compromise and/or optimization to allow for implementation on large-scale urban indices.
- **Data Dependency**: The system hence heavily depends on arriving timely, quality traffic data; data may not be present here-and-there in the case of sparse or insufficient sensor coverage.
- **Real-World Implementation**: Integration into current infrastructure is necessary for the successful deployment; it should conform to traffic regulations and the consideration of heavy-duty cybersecurity interventions.

## 7. Discussion:
More gains can be achieved with the HDRL framework over conventional and DRL methods. The CNN and RNN combo extract maximum features from multi-modal traffic data, while PPO provides stable decision-making to adapt to various scenarios. The system is versatile enough to reduce ATT, CI, CO2 emissions, and near-collision events. Evaluations are carried out in SUMO simulations, offering a robust platform for simulating realistic traffic dynamics under a range of conditions.

Yet, the higher computational time being seen with larger networks may suggest optimization techniques such as distribution computing or hardware acceleration across edge devices. With the system being data-dependent, great importance must be laid on traffic sensor networks and V2I

infrastructure. The scope of the future work includes:

- **Distributed HDRL Architectures**: To annotate greater scalability in large urban networks.
- **Real-World Pilots**: To ensure performance validation in real traffic systems.
- **Integration with Autonomous Vehicles**: To implement HDRL for coherent fleet management.
- **Advanced Reward Functions**: To include other objectives like pedestrian safety or public transit prioritization.

HDRL framework is one big leap toward intelligent, adaptive traffic management, where urban mobility for smart cities can become a reality.

## 8. Conclusion:

A novel HDRL framework proposes traffic control and decision-making applications, combining CNNs, RNNs, and PPO for traffic signal control optimization, vehicle routing, and congestion mitigation. Tested in grid, downtown, and highway scenarios in SUMO, the system outperforms fixed-time control, SCATS, and DQN systems in measures of ATT, CI, $CO_2$ emissions, and safety. The complete evaluation backed by six comprehensive tables and eight different graphs exhibits HDRL's promise of handling complicated, multi-modal traffic data, all while adapting to the ever-changing traffic conditions. Being interfaced with V2I and autonomous vehicles renders this an excellent candidate for a cross-industry scalable solution in the smart city domain. Future research will, therefore, involve making it more efficient computationally, establishing real-world pilots and integrating emerging technologies to leverage the greater impact potential for intelligent transportation systems.

## 9. References:

1. Arel I, Liu C, Urbanik T, Kohls AG. Reinforcement learning-based multi-agent system for network traffic signal control. *IET Intell Transp Syst.* 2010;4(2):128-35. doi:10.1049/iet-its.2009.0070
2. Goodfellow I, Bengio Y, Courville A. *Deep Learning.* MIT Press; 2016.
3. Krajzewicz D, Erdmann J, Behrisch M, Bieker L. SUMO: Simulation of Urban MObility. *Int J Adv Syst Meas.* 2012;5(3-4):128-38. Available from: http://sumo.dlr.de
4. Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, Bellemare MG, *et al*. Human-level control through deep reinforcement learning. *Nature.* 2015;518(7540):529-33. doi:10.1038/nature14236
5. Schulman J, Wolski F, Dhariwal P, Radford A, Klimov O. Proximal policy optimization algorithms. *arXiv.* 2017. doi:10.48550/arXiv.1707.06347
6. Silver D, Huang A, Maddison CJ, Guez A, Sifre L, van den Driessche G, *et al*. Mastering the game of Go with deep neural networks and tree search. *Nature.* 2017;550(7676):354-9. doi:10.1038/nature24270
7. Sutton RS, Barto AG. *Reinforcement Learning: An Introduction.* 2nd ed. MIT Press; 2018.
8. Van der Pol E, Oliehoek FA. Coordinated deep reinforcement learners for traffic light control. *NIPS Workshop on Learning, Inference and Control of Multi-Agent Systems.* 2016. Available from: https://nips.cc/Conferences/2016/Schedule?showEvent =6274
9. Wei H, Zheng G, Yao H, Li Z. CoLight: Learning network-level cooperation for traffic signal control. In: *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining.* 2019. p. 1916-26. doi:10.1145/3292500.3330836
10. Zhang Y, Wang J, Liu Y. Deep reinforcement learning for urban traffic control. *Transp Res Part C Emerg Technol.* 2020;116:102636. doi:10.1016/j.trc.2020.102636