

Multi-Agent Reinforcement Learning-based Cooperative Autonomous Driving in Smart Intersections

Taoyuan Yu, Kui Wang, Zongdian Li, Tao Yu, and Kei Sakaguchi

Abstract—Unsignalized intersections pose significant safety and efficiency challenges due to complex traffic flows. This paper proposes a novel roadside unit (RSU)-centric cooperative driving system leveraging global perception and vehicle-to-infrastructure (V2I) communication. The core of the system is an RSU-based decision-making module using a two-stage hybrid reinforcement learning (RL) framework. At first, policies are pre-trained offline using conservative Q-learning (CQL) combined with behavior cloning (BC) on collected dataset. Subsequently, these policies are fine-tuned in the simulation using multi-agent proximal policy optimization (MAPPO), aligned with a self-attention mechanism to effectively solve inter-agent dependencies. RSUs perform real-time inference based on the trained models to realize vehicle control via V2I communications. Extensive experiments in CARLA environment demonstrate high effectiveness of the proposed system, by: (i) achieving failure rates below 0.03% in coordinating three connected and autonomous vehicles (CAVs) through complex intersection scenarios, significantly outperforming the traditional Autoware control method, and (ii) exhibiting strong robustness across varying numbers of controlled agents and shows promising generalization capabilities on other maps.

I. INTRODUCTION

Intersection management is regarded as a bottleneck in the development of intelligent transportation systems (ITS), considering the complex and uncertain nature of urban intersections [1]. According to statistics from the Federal Highway Administration (FHWA) and the National Highway Traffic Safety Administration (NHTSA), intersection-related fatalities account for a substantial proportion of total traffic accident deaths, especially at unsignalized intersections, which reportedly accounted for 68% of such fatalities in 2024 [2], [3]. Unsignalized intersections have become accident hotspots due to blind spots and the lack of clear interaction rules between motor vehicles, non-motorized vehicles, and pedestrians. In such environments, connected and autonomous vehicles (CAVs) should possess highly effective perception, prediction, and coordinated decision-making capabilities to minimize conflicts and ensure safe and smooth driving [4].

With the popularization of autonomous vehicles (AVs), mixed traffic scenarios involving AVs and human-driven vehicles (HDVs) are becoming common, thereby introducing novel challenges to traffic participants. Vehicle-to-everything (V2X) communication technologies have emerged as a promising solution to improve roadway efficiency and safety [5], typically including vehicle-to-vehicle (V2V), vehicle-to-infrastructure (V2I), vehicle-to-pedestrian (V2P), and vehicle-to-network (V2N) links [6]. Among these, V2I communication facilitates real-time data exchange between

CAVs and smart roadside infrastructure (e.g., RSUs), laying a foundation for constructing cooperative driving systems [7].

Building upon V2I capabilities, the design of RSU-based cooperative systems has appeared as an attractive research topic for both academia and industry in recent years [8], [9]. Research on optimizing traffic flow at intersections has explored various methods. Traditional methods often rely on model-based optimization or game-theoretic frameworks to allocate right-of-way and manage traffic, achieving progress in reducing delays under certain conditions [10]–[12]. However, these methods typically lack the adaptability required to effectively handle the high complexity and uncertainty inherent in dynamic real-world traffic scenarios. To address these limitations, multi-agent reinforcement learning (MARL) has emerged as a promising alternative, with research investigating hierarchical structures or integrating perception modules to enhance coordination [13]–[15]. Nevertheless, many existing MARL applications in this domain tend to treat all vehicles uniformly, often overlooking the critical need for distinct policies tailored to specific driving roles (e.g., turning left, going straight, turning right) and failing to fully capture the complex interaction dynamics.

Further advancements aim to enhance MARL’s capabilities in complex environments. Self-attention mechanisms have been incorporated to dynamically model inter-agent dependencies, potentially improving generalization and decision efficiency [16], [17]. However, the effectiveness validation of adopting self-attention in MARL is still unexplored, particularly its adaptability to dynamically varying numbers of interacting vehicles within realistic transportation contexts like unsignalized intersections. Similarly, hybrid offline-online RL frameworks offer potential benefits by leveraging real-world data for safer and more efficient policy learning [18]–[21]. However, the practical implementation and demonstrated effectiveness of these hybrid approaches in coordinating multiple cooperative vehicles through complex and real-world intersection scenarios still require deeper investigation. Therefore, there exists a research gap in developing and validating an integrated MARL framework capable of robustly and efficiently coordinating vehicles with diverse intentions within the complex dynamics of intersections.

To address such challenges, we propose an innovative RSU-centric intelligent management system for unsignalized intersections. Utilizing bird-eye-view (BEV) perception from RSU-mounted LiDAR, the system employs a centralized MARL decision module featuring role-specific policy networks integrated with a self-attention mechanism, allowing for dynamic modeling of interactions between distinct

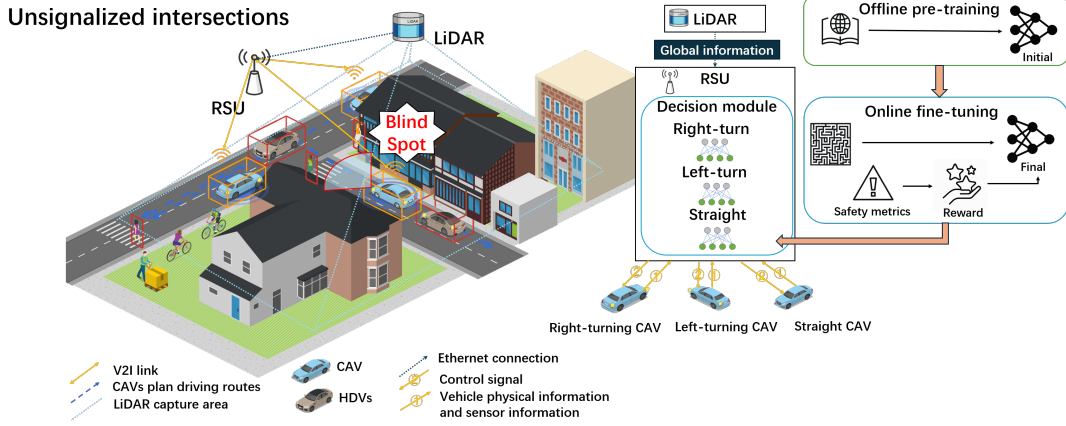


Fig. 1: High-level system design of the RSU-CAVs cooperative system

driving roles and flexible adaptation to varying numbers of vehicle participants. The policy networks are developed using a two-stage hybrid learning approach, involving offline pre-training on the collected dataset, followed by online fine-tuning in the simulation environment. The proposed system demonstrates significant advantages in adaptability, generalization, safety, and efficiency, while also reducing model deployment complexity and computational demands. The key contributions of this research are as follows:

- A novel hybrid RL framework combining offline pre-training and online fine-tuning techniques to enable cooperative driving for CAVs at unsignalized intersections.
- Development of personalized policy networks tailored to distinct driving roles (e.g., left-turn, straight, right-turn) at intersections.
- Integration of a self-attention mechanism into role-based MARL to enhance policy adaptability to varying vehicle numbers and model dynamic interactions.
- Demonstration of the model's generalization capability and rapid adaptability across diverse unsignalized intersection scenarios.

The remainder of the paper is structured as follows: Section II presents the overall architecture of the RSU-CAVs cooperative system. Section III describes the proposed algorithm in detail. Section IV demonstrates experiment results. Finally, Section V summarizes the paper and outlines future research directions.

II. RSU-CAVs COOPERATIVE SYSTEM

The overall system architecture of the proposed RSU-CAV cooperative framework is illustrated in Fig. 1. This system employs an RSU, equipped with sensors like LiDAR, for comprehensive intersection monitoring [22] and centralized decision-making for multiple CAVs at an unsignalized intersection. Leveraging the RSU's BEV perception, this centralized method overcomes the inherent limitations of individual vehicle perception, providing a global understanding of the traffic situation essential. This contrasts with individualistic methods, where each vehicle optimizes only its own goals, which can lead to competitive standoffs, inefficient gridlocks,

or unsafe maneuvers in unsignalized interactions. Instead, our framework enhances collective safety and overall traffic throughput by resolving conflicts harmoniously.

To effectively manage the intersection's complexities and uncertainties, the RSU utilizes adaptive decision-making policies developed through an RL-based method. This method begins by employing offline RL to instill foundational driving knowledge and essential interaction behaviors into the policy from collected datasets, establishing a robust and competent initial strategy. Subsequently, these foundational policies undergo targeted refinement using online RL within a simulation environment. This allows the policies to adapt to the intersection's unique dynamic characteristics and optimize performance for the required multi-vehicle cooperative tasks. Compared to learning entirely from scratch, this hybrid RL approach offers the advantages of accelerating learning convergence during the online refinement phase and ultimately yielding coordination policies that are more robust and effective in handling real-world complexities. Furthermore, deploying the computationally intensive analysis and multi-vehicle coordination logic onto the RSU also reduces the processing burden on individual CAVs. This allocation of tasks improves the system's real-time responsiveness and simplifies requirements for the CAVs [23].

As CAVs approach the intersection, they maintain continuous information exchange with the RSU. Based on real-time traffic data and the CAVs' approach trajectories, the RSU determines each vehicle's driving role (e.g., left-turn, straight, right-turn). Subsequently, leveraging its pre-loaded and role-based strategy networks within the centralized decision module, the RSU computes vehicle control signals, including throttle input, braking force, and steering angles. These control commands are transmitted in real-time to the corresponding CAVs through V2I communication, enabling direct command execution. Concurrently, the RSU continuously monitors the comprehensive real-time traffic conditions at the intersection, including the states and predicted movements of CAVs, observed HDVs, and pedestrians, alongside traffic flow smoothness, collision risks, and any abnormal situations. This continuous monitoring provides the neces-

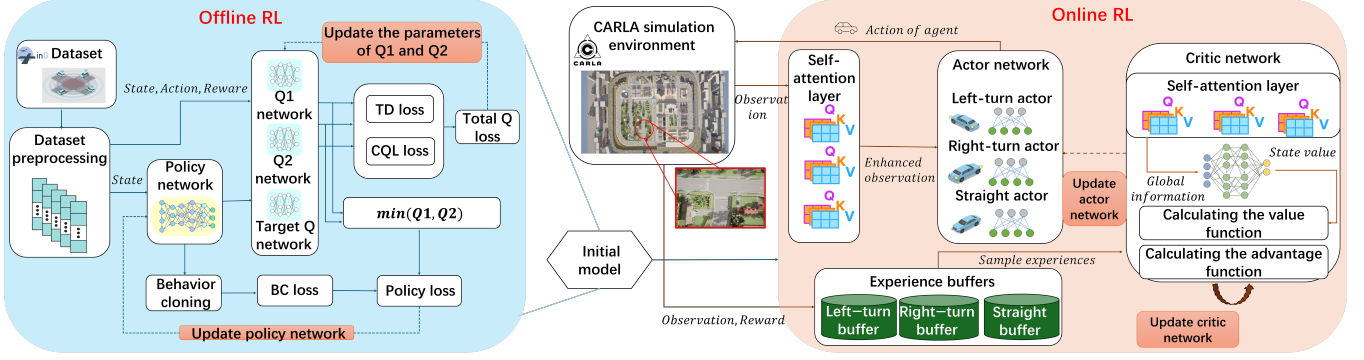


Fig. 2: Offline-online hybrid RL algorithm framework design

sary real-time inputs for the RSU's decision networks and facilitates ongoing performance evaluation.

III. HYBRID REINFORCEMENT LEARNING FRAMEWORK

As shown in Fig. 2, we propose a two-stage learning framework to develop effective cooperative driving strategies for complex unsignalized intersections. This approach first employs offline pre-training on the collected dataset, using offline RL to safely learn foundational driving skills and traffic priors. Subsequently, online fine-tuning within the CARLA simulator [24] allows agents to adapt to specific environment dynamics. This synergistic framework combines the safety and data efficiency of offline learning with the adaptability and performance optimization capabilities of online interaction. The specific methodologies for offline pre-training and online fine-tuning are discussed in this section.

A. Observation Space

At each time step t , the state space can be defined $\mathbf{s}(t)$, which encompasses all traffic participants monitored by RSU within the intersection. the RSU utilizes the global information to generate specific perspective information for each CAV, constructing an individual observation vector $\mathbf{o}(t)$. This construction process, denoted conceptually as $\mathbf{o}(t) = \mathcal{O}(\mathbf{s}(t))$, reflects simulated perception and processing limitations, meaning each $\mathbf{o}(t)$ is a partially observable and potentially noisy representation of the true state $\mathbf{s}(t)$. The observation $\mathbf{o}(t)$ is decomposed as:

$$\mathbf{o}(t) = [\mathbf{o}_{\text{core}}, \mathbf{o}_{\text{veh}}, \mathbf{o}_{\text{ped}}, \mathbf{o}_{\text{role}}, \mathbf{o}_{\text{ctx}}], \quad (1)$$

where \mathbf{o}_{core} denotes the ego vehicle features including speed magnitude, global position, heading angle, distance to intersection center, and junction occupancy status; \mathbf{o}_{veh} denotes relative positions and velocities of nearby vehicles within a predefined ranged; \mathbf{o}_{ped} denotes information on nearby pedestrians, including detection flags, distance, and relative angle; \mathbf{o}_{role} denotes the one-hot encoding of the agent's role (e.g., left-turn, straight, or right-turn); and \mathbf{o}_{ctx} denotes scenario-level identifiers used to curriculum learning.

B. Action Space

Throughout the learning framework, we define a unified two-dimensional continuous action space \mathcal{A} for the vehicle. It is structured as:

$$\mathbf{a}(t) = [\mathbf{a}_{\text{acc}}, \mathbf{a}_{\text{steer}}] \in \mathbb{R}^2 \quad (2)$$

where \mathbf{a}_{acc} denotes the longitudinal acceleration and $\mathbf{a}_{\text{steer}}$ denotes the steering angular velocity. It is important to note that during the offline pre-training phase, as ground-truth control signals are unavailable in the source data, the action $\mathbf{a}(t)$ in the dataset is estimated by analyzing the state transitions between consecutive changes in velocity and heading. In the online fine-tuning phase, the policy network directly outputs two-dimensional action.

C. Reward Function

To effectively guide the agent in learning desired cooperative driving behaviors during complex online interactions, we design a structured reward function $\mathcal{R}_{\text{online}}(\mathbf{s}(t), \mathbf{a}(t), \mathbf{s}(t+1))$ to translate high-level objectives into real-time feedback. The overall reward $r(t)$ is structured as:

$$r(t) = \sum w_k r_k(\mathbf{s}(t), \mathbf{a}(t), \mathbf{s}(t+1)), \quad (3)$$

where r_i denotes individual reward components and w_i denotes the corresponding weights. The reward terms include:

$$r_i \in \{r_{\text{safety}}, r_{\text{eff}}, r_{\text{comfort}}, r_{\text{task}}, r_{\text{yield}}, r_{\text{coop}}, r_{\text{penalty}}\} \quad (4)$$

where r_{safety} denotes the penalty for hazardous behavior based on metrics like minimum time-to-collision (TTC) and distance to nearby vehicles or pedestrians; r_{eff} denotes the efficiency that encourages maintaining a reasonable speed that is compatible with traffic flow; r_{comfort} denotes the penalty for large acceleration changes; r_{task} denotes the reward for all agents to cooperatively reach the navigation target; r_{yield} and r_{coop} denotes the rewards for promoting compliance with traffic rules and cooperation; and r_{penalty} denotes a severe penalty imposed on events like collisions or timeouts. Each term is scaled by its corresponding weight w_k , where w_{safety} and w_{penalty} are typically assigned larger values due to their critical safety implications.

D. Offline Pre-training: Networks and Algorithm

The primary goal of the offline pre-training phase is to provide a high-quality initialization for the subsequent online fine-tuning stage. In this stage, we train models independently for each driving role (left-turn, straight, right-turn) to incorporate role-specific prior knowledge. We first partition the InD dataset [25] based on vehicle intentions to create subsets $\mathcal{D}_{\text{role}}$.

For each subset, we apply an offline reinforcement learning algorithm that combines CQL with BC [26] [27]. The algorithm is implemented using an actor-critic framework to learn effectively from fixed datasets while mitigating distributional shift and imitating expert behavior.

To stabilize learning and reduce Q-value overestimation, the critic uses twin Q-networks $Q_{\theta_{i,1}}, Q_{\theta_{i,2}}$ with target networks. Each Q-network is trained with the following objective:

$$L_Q(\theta_{i,j}) = \mathbb{E}_{(\mathbf{o}, \mathbf{a}, r, \mathbf{o}') \sim \mathcal{D}_{\text{role}=i}} \left[\frac{1}{2} (Q_{\theta_{i,j}}(\mathbf{o}, \mathbf{a}) - y)^2 \right] + \alpha_{\text{CQL}} L_{\text{CQL-reg}}(\theta_{i,j}) \quad (5)$$

Here, $y = r + \gamma(1 - d) \min_j Q_{\theta'_{i,j}}(\mathbf{o}', \pi_{\phi_i}(\mathbf{o}'))$ denotes the TD target computed using target Q networks and the current policy.

The policy network π_{ϕ_i} is trained by minimizing the BC loss along with maximizing the expected conservative Q-value:

$$L_\pi(\phi_i) = \mathbb{E}_{\mathbf{o} \sim \mathcal{D}_{\text{role}=i}} \left[- \min_{j=1,2} Q_{\theta_{i,j}}(\mathbf{o}, \pi_{\phi_i}(\mathbf{o})) \right] + \lambda_{\text{BC}} \mathbb{E}_{(\mathbf{o}, \mathbf{a}) \sim \mathcal{D}_{\text{role}=i}} [\|\pi_{\phi_i}(\mathbf{o}) - \mathbf{a}\|^2] \quad (6)$$

where α_{CQL} and λ_{BC} denotes hyperparameters controlling the strength of CQL regularization and BC imitation, respectively.

Each role-specific actor $\pi_{\phi_{\text{role}}}$ and critic $Q_{\theta_{\text{role}}}$ are parameterized by multi-layer perceptrons (MLPs), which take normalized state inputs s_t sampled from $\mathcal{D}_{\text{role}}$. The output corresponds to action distribution parameters or state-action values.

Self-attention are not introduced at this stage to focus training on extracting robust role-based patterns using standard MLPs. Training stability is further improved by soft target updates and the adam optimizer. Successful offline training yields a set of pre-trained weights for the role-conditioned actor and critic networks, which are reused during the online fine-tuning phase to accelerate learning and enhance performance.

E. Online Fine-tuning: Networks and Algorithm

The online fine-tuning phase employs the MAPPO algorithm [28], selected for its effectiveness in multi-agent coordination. This stage builds upon the offline models by integrating role-specific actor networks ($\pi_{\phi_{\text{left}}}, \pi_{\phi_{\text{straight}}}, \pi_{\phi_{\text{right}}}$) with a shared critic network V_ψ .

To improve reasoning over dynamic environments, we augment both actor and critic networks with multi-head self-attention (MHSA). MHSA allows the model to jointly attend

to information from different representation subspaces at different positions. The core component is the scaled dot-product attention:

$$\text{A}(Q, K, V) = \text{softmax} \left(\frac{QK^\top}{\sqrt{d_k}} \right) V \quad (7)$$

where Q , K , and V denote the query, key, and value matrices, and d_k denotes the dimension of the keys. MHSA computes h attention 'heads' in parallel. For each head i , the input embedding E is linearly projected using learned weights W_i^Q, W_i^K, W_i^V to obtain the head's specific query, key, and value:

$$\mathbf{h}_i = \text{Attention}(EW_i^Q, EW_i^K, EW_i^V) \quad (8)$$

The outputs of the parallel heads are then concatenated and linearly projected using weights W^O to produce the final MHSA output:

$$M(E) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O \quad (9)$$

Here, E denotes the initial embedded observation features derived from the online observation \mathbf{o}_t via a learnable linear projection layer, and W_i^Q, W_i^K, W_i^V, W^O denote trainable weight matrices.

Online learning proceeds via an interact-learn loop. Agents generate trajectories:

$$\tau = \{(\mathbf{o}_t, \mathbf{a}_t, r_{t+1}, V_\psi(\mathbf{o}_t), \log \pi_{\phi_{\text{role}}}(\mathbf{a}_t | \mathbf{o}_t))\}_{t=0}^T \quad (10)$$

Advantage estimates and returns are computed using generalized advantage estimation (GAE):

$$\hat{A}_t^{\text{GAE}} = \sum_{l=0}^{T-t-1} (\gamma\lambda)^l \delta_{t+l}, \quad \delta_t = r_{t+1} + \gamma V_\psi(\mathbf{o}_{t+1}) - V_\psi(\mathbf{o}_t) \quad (11)$$

$$\hat{R}_t = \hat{A}_t^{\text{GAE}} + V_\psi(\mathbf{o}_t) \quad (12)$$

To enhance data efficiency, we adopt prioritized experience replay (PER). Each transition t is assigned a priority p_t proportional to its absolute TD error $|\delta_t|$, and sampled with probability $P(t) \propto p_t$. To correct the bias introduced by this non-uniform sampling, importance sampling (IS) weights are applied:

$$w_t = \left(\frac{1}{B \cdot P(t)} \right)^\beta \quad (13)$$

Here, B denotes the size of the replay buffer, and β denotes an exponent that controls the amount of importance sampling correction.

The shared critic is updated to minimize the weighted value loss:

$$L^{VF}(\psi) = \mathbb{E}_{t \sim \text{PER}} \left[w_t (V_\psi(\mathbf{o}_t) - \hat{R}_t)^2 \right] \quad (14)$$

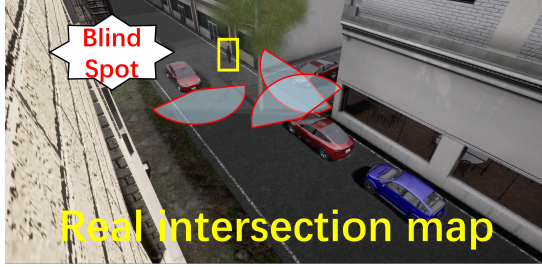
where \hat{R}_t denotes the target return calculated in Eq.(12).

Each role-specific actor $\pi_{\phi_{\text{role}}}$ is trained using the following weighted objective, which includes the PPO clipped surrogate loss and an entropy bonus $S[\cdot]$:

$$L^{\text{CLIP+S}}(\phi_{\text{role}}) = \mathbb{E}_{t \sim \text{PER}} \left[w_t \left(-L_t^{\text{CLIP}}(\phi_{\text{role}}) - c_2 \cdot S[\pi_{\phi_{\text{role}}}(\mathbf{o}_t)] \right) \right] \quad (15)$$



(a)



(b)

Fig. 3: Experimental scenario and generalization scenario settings (a) CARLA example map, (b) Real intersection map

The PPO surrogate loss L_t^{CLIP} is defined as:

$$L_t^{\text{CLIP}} = \min \left(r_t \hat{A}_t, \text{clip}(r_t, 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \quad (16)$$

where ϵ denotes the PPO clipping hyperparameter, and r_t denotes the probability ratio between the current policy and the old policy:

$$r_t = \frac{\pi_{\phi_{\text{role}}}(\mathbf{a}_t \mid \mathbf{o}_t)}{\pi_{\phi_{\text{old}}}(\mathbf{a}_t \mid \mathbf{o}_t)} \quad (17)$$

We further enhance training with adam optimizer and gradient clipping. These stabilizing techniques enable robust fine-tuning and effective adaptation to dynamic, multi-agent environments.

IV. EXPERIMENTS AND ANALYSIS

Experiments were conducted using the CARLA simulator paired with Unreal Engine in synchronous mode. The primary scenario involves one intersection within CARLA's Town03 map (shown in Fig. 3). In this simulation scenario, a total of 5 vehicles are present: a variable number (1 to 3) CAVs, designated "red", are controlled by the proposed system, while the remaining background vehicles, designated "blue", are managed by CARLA's Traffic Manager. Additionally, up to 3 pedestrians are randomly spawned on sidewalks and programmed to cross the road. A real intersection map based on the Institute of Science Tokyo campus was utilized for generalization testing. We assume the RSU possesses BEV perception and performs inference using the decision model obtained after online fine-tuning to determine driving strategies, subsequently sending control signals to the CAVs via simulated V2I communication.



Fig. 4: Offline pre-training results

A. Baselines and Evaluation Metrics

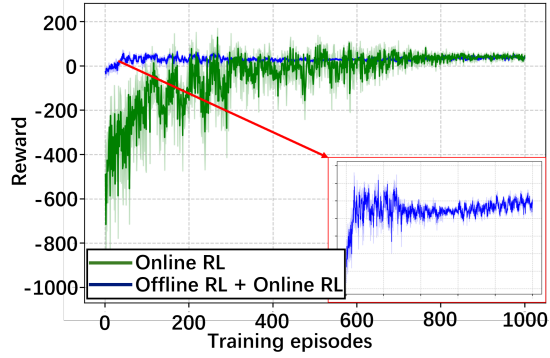
We used online-only MAPPO as part of an ablation study; we trained a MAPPO agent that undergoes no offline pre-training and starts training directly in CARLA from scratch (with the network structure and algorithm parameters being the same as the online phase of our proposed method), used to compare and evaluate the performance of our proposed algorithm. We will primarily compare them in terms of convergence speed. Additionally, we employ the open-source autonomous driving software stack, Autoware Universe [29] as a representative of traditional autonomous driving systems. In our experiments, Autoware is configured to control a single vehicle navigating the intersection within the identical scenario used for our single-agent RL tests (i.e., 1 Autoware-controlled vehicle, 4 background Traffic Manager vehicles, and 3 pedestrians). Its performance provides a benchmark for comparison against established methods.

B. Offline Pre-training Results

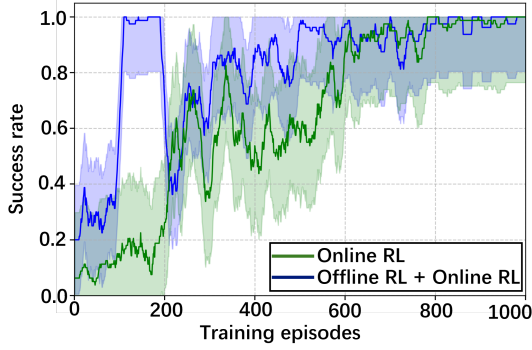
The purpose of the offline pre-training phase is to extract effective driving priors from the InD dataset, providing a high-quality model initialization for the online phase. As shown in Fig. 4, which displays the changes in the critic networks' Q1 and Q2 losses and the reward improvement metric during offline training, we can observe the learning progress. The loss values steadily converge as training progresses, indicating that the critic network effectively learned state-action value relationships from the offline data and that the training process possessed good stability. Furthermore, the reward improvement metric eventually stabilizes around 112%, exceeding the 100% baseline. This demonstrates that the policy learned via CQL combined with BC outperforms the average behavior present in the dataset in terms of optimizing the offline reward objective, successfully learning strategies beyond mere imitation, and providing a quality initialization basis for online fine-tuning.

C. Online Training Results

To validate the effectiveness of online fine-tuning and the value of offline pre-training, we compare the training progress of our proposed hybrid method against the online-only baseline. Fig. 5 presents the convergence curves for both



(a)



(b)

Fig. 5: Comparison of online training w/ and w/o offline RL, (a) reward, (b) success rate.

average episode reward and success rate during the external online training phase for both methods.

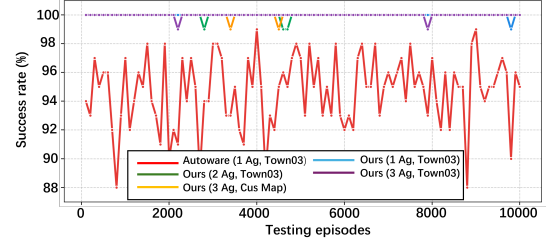
As shown in Fig. 5, our proposed hybrid method, leveraging the pre-trained model, exhibits a much higher initial average episode reward compared to the online-only method starting from scratch. Furthermore, the hybrid approach converges more rapidly to its final reward level, indicating that offline pre-training significantly accelerates the online learning process and contributes to achieving strong final performance. Additionally, the convergence trend for success rate mirrors that of the reward. Our proposed method reaches high success rates considerably faster, whereas the online-only method requires significantly more training episodes to approach similar levels of reliability. This comparison further demonstrates that offline pre-training markedly enhances both the efficiency of the online learning phase and the robustness of the ultimately learned policy.

D. Final Model Performance Evaluation and Generalization Analysis

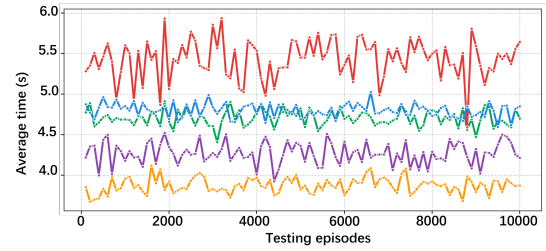
We evaluated the model through 10,000 performance test episodes on both the Town03 intersection and the custom map, comparing it against baselines. A summary of key performance indicator comparisons is shown in Fig. 6 and Tab. I. The proposed algorithmic model demonstrates high safety and reliability across all test scenarios on the Town03

TABLE I: Performance comparison summary

Method / Scenario	Failure rate (%)	Avg. time (s)
Ours (1 Agent, Town03)	0.01	5.52
Ours (2 Agent, Town03)	0.03	5.49
Ours (3 Agent, Town03)	0.02	5.25
Autoware (1 Agent, Town03)	5.31	5.77
Ours (3 Agent, Real Map)	0.02	5.15



(a)



(b)

Fig. 6: Final model performance evaluation results (a) success rate by testing episodes, (b) average travel time by testing episodes

map. When controlling a single vehicle, the failure rate was 0.01%. Compared to the 5.31% failure rate exhibited by the Autoware baseline in the identical single-vehicle scenario, our single-agent controller shows higher performance.

Notably, despite the significant increase in coordination complexity when scaling from single- to multi-vehicle scenarios, our system did not exhibit a marked decline in success rate. Specifically, failure rate was 0.03% in the two-vehicle coordination scenario and decreased to 0.02% in the three-vehicle coordination scenario. The combination of the RSU's BEV perspective and the self-attention mechanism contributes to this robustness, demonstrating our method's effectiveness in handling complex multi-agent cooperative tasks.

As shown in Fig. 6 and Tab. I, the traffic efficiency results indicate that our method also demonstrated strong performance. The average travel time in the single-vehicle scenario was 5.52 seconds, outperforming the 5.77 seconds recorded by the Autoware baseline. As the number of controlled vehicles increased, the average travel time showed a slight downward trend: 5.49 seconds for the two-vehicle scenario and 5.25 seconds for the three-vehicle scenario. This indicates that the multiple RL agents coordinated by our system formed a highly effective collaborative passage pattern. Their interactions proved more efficient than those with a larger number of background vehicles controlled by

the Traffic Manager.

Finally, in the generalization test, the three-vehicle model trained in Town03 was deployed on the real intersection map. The model achieved an extremely low failure rate of 0.02% and an average travel time of 5.15 seconds in this novel environment. This suggests that the expanded collective field of view inherent in the three-vehicle setup mitigated the impact of individual visual blind spots. Furthermore, as this map featured shorter traversal distance, it enabled our system to operate more smoothly and efficiently. This result strongly validates the excellent generalization capability of the learned policy, showcasing its adaptability to different environmental characteristics and providing a solid foundation for the practical application of the method.

V. CONCLUSION AND FUTURE WORKS

This research addresses the complex coordination challenges at unsignalized intersections by proposing an RSU-based centralized cooperative driving framework. The framework employs a two-stage method to train the decision model: offline pre-training initializes policies, followed by online fine-tuning in the simulation environment. Extensive experiments demonstrate the method's effectiveness, achieving high success rate and strong coordination robustness in scenarios with up to three controlled vehicles. Future primary work involves the proof-of-concept (PoC) experiments to fully validate the system effectiveness in the real world.

REFERENCES

- [1] K. F. Chu, A. Lam, and V. Li, "Traffic signal control using end-to-end off-policy deep reinforcement learning," *IEEE Transactions on Intelligent Transportation Systems*, vol. PP, pp. 1–12, 2021.
- [2] U.S. Department of Transportation, Federal Highway Administration. (2024, August) Mire 2.1: Model inventory of roadway elements. Report No. FHWA-SA-24-052. [Online]. Available: https://highways.dot.gov/sites/fhwa.dot.gov/files/2024-08/MIRE_2.1_FINAL_508v3.pdf
- [3] National Highway Traffic Safety Administration. (2024, April) Nhtsa estimates 39,345 traffic fatalities in 2024. U.S. Department of Transportation. [Online]. Available: <https://nhtsa.gov/press-releases/nhtsa-2023-traffic-fatalities-2024-estimates>
- [4] S. Chen, X. Hu, J. Zhao, R. Wang, and M. Qiao, "A review of decision-making and planning for autonomous vehicles in intersection environments," *World Electric Vehicle Journal*, vol. 15, no. 3, p. 99, 2024.
- [5] K. Wang, C. She, Z. Li, T. Yu, Y. Li, and K. Sakaguchi, "Roadside units assisted localized automated vehicle maneuvering: An offline reinforcement learning approach," in *2024 IEEE 27th International Conference on Intelligent Transportation Systems (ITSC)*, 2024, pp. 1709–1715.
- [6] Z. Li, K. Wang, T. Yu, and K. Sakaguchi, "Het-SDVN: SDN-based radio resource management of heterogeneous V2X for cooperative perception," *IEEE Access*, vol. 11, pp. 76 255–76 268, 2023.
- [7] D. Suo, B. Mo, J. Zhao, and S. E. Sarma, "Proof of travel for trust-based data validation in V2I communication," *IEEE Internet of Things Journal*, vol. 10, no. 11, pp. 9565–9584, 2023.
- [8] H. Alemayehu and A. Sargolzaei, "Testing and verification of connected and autonomous vehicles: A review," *Electronics*, vol. 14, no. 3, p. 600, 2025.
- [9] F. Sana, N. L. Azad, and K. Raahemiar, "Autonomous vehicle decision-making and control in complex and unconventional scenarios—a review," *Machines*, vol. 11, no. 7, p. 676, 2023.
- [10] Y. Zhu, Z. He, and G. Li, "A bi-hierarchical game-theoretic approach for network-wide traffic signal control using trip-based data," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 9, pp. 15 408–15 419, 2022.
- [11] M. Gallo, "Combined optimisation of traffic light control parameters and autonomous vehicle routes," *Smart Cities*, vol. 7, no. 3, pp. 1060–1088, 2024.
- [12] M. Ghadi, "A grid-based framework for managing autonomous vehicles' movement at intersections," *Periodica Polytechnica Transportation Engineering*, vol. 52, no. 3, pp. 235–245, 2024.
- [13] Y. Shi, H. Dong, C. He, Y. Chen, and Z. Song, "Mixed vehicle platoon forming: A multi-agent reinforcement learning approach," *IEEE Internet of Things Journal*, vol. PP, pp. 1–1, 01 2025.
- [14] H. Taghavifar, C. Hu, C. Wei, A. Mohammadzadeh, and C. Zhang, "Behaviorally-aware multi-agent rl with dynamic optimization for autonomous driving," *IEEE Transactions on Automation Science and Engineering*, vol. PP, pp. 1–1, 01 2025.
- [15] Y. Zhang, R. Hao, T. Zhang, X. Chang, Z. Xie, and Q. Zhang, "A trajectory optimization-based intersection coordination framework for cooperative autonomous vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 9, pp. 14 674–14 688, 2021.
- [16] S. Iqbal and F. Sha, "Actor-attention-critic for multi-agent reinforcement learning," *arXiv preprint arXiv:1810.02912*, 2019, iCML 2019 camera ready version.
- [17] R. Younas, H. M. Raza Ur Rehman, I. Lee, B.-W. On, S. Yi, and G. S. Choi, "Sa-marl: Novel self-attention-based multi-agent reinforcement learning with stochastic gradient descent," *IEEE Access*, vol. 13, pp. 35 674–35 687, 2025.
- [18] W.-C. Tseng, T.-H. J. Wang, Y.-C. Lin, and P. Isola, "Offline multi-agent reinforcement learning with knowledge distillation," in *Advances in Neural Information Processing Systems*, vol. 35. Curran Associates, Inc., 2022, pp. 226–237.
- [19] Z. Li, F. Nie, Q. Sun, F. Da, and H. Zhao, "Boosting offline reinforcement learning for autonomous driving with hierarchical latent skills," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024, pp. 18 362–18 369.
- [20] Z. Wang, G. Wu, and M. J. Barth, "Cooperative eco-driving at signalized intersections in a partially connected and automated vehicle environment," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 5, pp. 2029–2038, 2019.
- [21] Z. Wang, K. Han, and P. Tiwari, "Digital twin-assisted cooperative driving at non-signalized intersections," *IEEE Transactions on Intelligent Vehicles*, vol. 7, no. 2, pp. 198–209, 2021.
- [22] K. Wang, Z. Li, K. Nonomura, T. Yu, K. Sakaguchi, O. Hashash, and W. Saad, "Smart mobility digital twin based automated vehicle navigation system: A proof of concept," *IEEE Transactions on Intelligent Vehicles*, pp. 1–14, 2024.
- [23] K. Wang, T. Yu, Z. Li, K. Sakaguchi, O. Hashash, and W. Saad, "Digital twins for autonomous driving: A comprehensive implementation and demonstration," in *2024 International Conference on Information Networking (ICOIN)*, 2024, pp. 452–457.
- [24] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An open urban driving simulator," in *Proceedings of the 1st Annual Conference on Robot Learning*, 2017, pp. 1–16.
- [25] J. Bock, R. Krajewski, T. Moers, S. Runde, L. Vater, and L. Eckstein, "The ind dataset: A drone dataset of naturalistic road user trajectories at german intersections," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, 2020, pp. 1929–1934.
- [26] A. Kumar, A. Zhou, G. Tucker, and S. Levine, "Conservative q-learning for offline reinforcement learning," in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020, pp. 1179–1191.
- [27] D. A. Pomerleau, "Alvin: An autonomous land vehicle in a neural network," in *Advances in Neural Information Processing Systems*, vol. 1. Morgan-Kaufmann, 1988.
- [28] C. Yu, A. Velu, E. Vinitzky, J. Gao, Y. Wang, A. Bayen, and Y. Wu, "The surprising effectiveness of ppo in cooperative, multi-agent games," *arXiv preprint arXiv:2103.01955*, 2021, accepted at NeurIPS 2022 Datasets and Benchmarks. [Online]. Available: <https://arxiv.org/abs/2103.01955>
- [29] Autoware. [Online]. Available: <https://autoware.org/>