

ENPM 667 - Project 1

A Technical Report on Deep Multi-agent Reinforcement Learning for Highway On-Ramp Merging in Mixed Traffic

(Dong Chen , Mohammad R. Hajidavalloo , Zhaojian Li , Kaian Chen ,
Yongqiang Wang , Longsheng Jiang³ , Yue Wang)

Kunj Golwala UID: 121118271

Shreya Kalyanaraman UID: 121166647



A. James Clark School of Engineering (Robotics)
University of Maryland, College Park, MD - 20742

Table of Contents

Abstract	i
1. Introduction	2
2. Background.....	3
2.1 Preliminaries of Reinforcement Learning.....	3
2.1.1. State-Action Value Function and Bellman Equation	3
2.1.2. Actor-Critic Framework.....	3
2.2. Multi-agent Reinforcement Learning (MARL)	4
3. Ramp merging as MARL	7
3.1 Problem Formulation	7
3.2 MA2C for CAVs	9
3.3 DNN Settings	10
4. Priority-based Safety Enhancement.....	11
4.1 Trajectory Propagation and Motion Primitives	12
4.2 Priority Assignment.....	12
4.3 Priority-based Safety Supervisor	13
4.4. Algorithm 1 Priority-Based Safety Supervisor.....	14
4.5. Algorithm 2 MARL for AVs with Safety Supervisor.....	16
5. Numerical Experiments	18
5.1 Reward Function Designs.....	19
5.2 Curriculum Learning for Hard Traffic Mode	20
5.3 Performance of the Priority-based Safety Supervisor.....	20
5.4. Comparison with Benchmark Algorithms.....	22
6. Conclusions and Discussions	24
References	25
APPENDIX	26

List of Figures

1. Flow chart of system and simulation setup	8
2. Variation of merging reward with distance x.....	9
3. Schematic representation of the proposed DNN architecture.....	11
4. Training curves for the n-step priority-based safety supervisor.....	21
5. Average speed during training for the n-step priority-based safety..... supervisor	22
6. Training curves comparison between the proposed MARL policy..... (baseline (bs) + $T_n = 8$) and state-of-the-art MARL benchmark	23

List of Tables

1. Testing performance comparison of collision rate.....23
and average speed (m/s) between n-step safety
supervisor based on the baseline (bs) method.

A Technical Review of Deep Multi-agent Reinforcement Learning for Highway On-Ramp Merging in Mixed Traffic

Dong Chen¹, Mohammad R. Hajidavalloo², Zhaojian Li³, Kaian Chen⁴, Yongqiang Wang⁵, Longsheng Jiang⁶, and Yue Wang⁷

ABSTRACT

The results of Dong Chen et al.'s "Deep Multi-agent Reinforcement Learning for Highway On-Ramp Merging in Mixed Traffic" investigation are presented and replicated in this work. The report addresses on-ramp merging in surroundings populated by human-driven vehicles (HDVs) and autonomous vehicles (AVs). Viewed as a multi-agent reinforcement learning (MARL) problem, a distributed architecture is used to encourage adaptive cooperation among autonomous cars under different traffic conditions. Important improvements cover parameter sharing and a priority-based safety supervisor to enable safe, effective interactions among agents, hence lowering crash rates and enhancing general traffic flow. Furthermore, curriculum learning is used to progressively raise traffic complexity so that autonomous cars may control more challenging environments as the learning process advances. This framework shows great safety improvement over current MARL methods, hence improving traffic efficiency. This paper attempts to reproduce and explain their techniques.

Keywords: Reinforcement Learning, Multi-agent Systems, Autonomous Vehicles



I. INTRODUCTION

Technological advancements like Tesla's Autopilot^[1] have rendered autonomous vehicles (AVs) more prevalent and sophisticated. These vehicles, equipped with sensors, machine learning algorithms, and decision-making frameworks, have the potential to significantly enhance road safety and improve traffic efficiency. Public road acceptance of AVs has, nevertheless, accompanied an increase in AV-related incidents, particularly in demanding traffic conditions. Sometimes AVs' inability to dynamically and securely adapt to changing driving conditions causes these incidents. The issue gets more crucial in mixed traffic conditions, as AVs have to live with and interact with human-driven vehicles (HDVs).

Highway on-ramp merging calls for AVs to negotiate a merging lane in coordination with cars in the through lane. This coordination consists of complex interactions between AVs modifying their speed to squeeze into gaps and HDVs may be changing their behavior to allow merging cars. In a perfect world, AVs and HDVs would work in perfect harmony; AVs would use these gaps to merge effectively while HDVs changed their speed or location to produce merging gaps. Achieving this degree of synchronizing is difficult in reality, especially in mixed traffic where the presence of human drivers adds variable and uncertainty in driving actions. Rule-based have also been studied and these methods rely on pre-defined heuristics or protocols, such as simple if-then rules (e.g., speed upon observed gaps). These techniques have complicated and dynamic characteristics of actual traffic that test them even if they operate well in controlled or basic conditions. Conversely, optimization-based methods—such as Model Predictive Control (MPC)—predict and maximize vehicle interactions across a limited time horizon using mathematical models. MPC seeks to identify an ideal series of activities to get desired outcomes, such as decreasing travel time or avoiding crashes, viewing model traffic as a dynamic system and These methods are less useful for real-time decision-making under mixed traffic conditions, though, because they depend on extremely precise models of traffic dynamics and entail large computing overhead.

More recently, data-driven methods leveraging reinforcement learning (RL) have shown promise as a fix for problems with highway on-ramp merging. Agents in the field of machine learning sometimes referred to as RL learn to make decisions by interacting with an environment and receiving feedback. Iteratively through trial and error, an RL agent finds a policy—a mapping from observable states to actions—that maximizes its cumulative reward over time. Although traditional RL methods show potential, they usually rely on single-agent situations, seeing other automobiles on the road as either static or dynamic environmental obstructions rather than as interactive players with their intentions and behaviors. These approaches thus miss the complexity of vehicle-to-vehicle interactions under mixed traffic conditions. This work extends single-agent RL to a multi-agent reinforcement learning (MARL) paradigm addressing its limits. Many agents in MARL, including autonomous cars, learn policies not in an isolated manner but rather collectively so that their activities complement one another. Within the framework of highway on-ramp merging, this method enables several AVs to learn cooperative merging techniques considering the presence and behaviors of nearby HDVs. Parameter sharing—where all AVs employ a shared policy and value network—allows consistent decision-making and scalability across many traffic conditions, hence defining a fundamental component of this MARL system. Furthermore, included in the framework are localized rewards that motivate AVs to maximize interactions with their nearby neighbors, thereby fostering cooperation while preserving computational economy.

A priority-based safety supervisor is proposed to improve the MARL framework's safety and learning efficiency even more. This system serves as protection during the training process since it controls the agents' behavior depending on traffic circumstances and predefined safety limits, therefore reducing collision chances. Safety first guarantees that, even in simulated surroundings, learning does not affect the physical well-being of vehicles or passengers. Using curriculum learning—a training approach modeled on human learning processes—is another essential element. Starting with basic, low-traffic surroundings and subsequently adding more difficult settings, training scenarios in curriculum learning are progressively made more complex. This steady increase in complexity lets AVs develop strong decision-making skills gradually, adjusting to ever-challenging merging conditions as training advances.

The suggested approach is compared to state-of-the-art MARL methods, and testing reveals notable increases in traffic flow and safety criteria. The framework solves the important difficulties of highway on-ramp merging in mixed traffic by allowing AVs to develop adaptive, cooperative merging behaviors while considering the uncertain character of HDVs. This paper emphasizes the possibilities of MARL methods, improved with safety monitoring and curriculum learning, to open the path for safer and more effective integration of autonomous vehicles into challenging traffic systems.

II. BACKGROUND

II.1. Preliminaries of Reinforcement Learning

For any discrete time interval t :

- The agent observes the current state $s_t \in S$, where S denotes the state space encompassing all feasible environmental configurations.
- The agent selects an action $a_t \in A$, where A represents all permitted actions.
- The surroundings provide a reward $r_t \in \mathbb{R}$, which measures the instantaneous gain (or loss) of performing action a_t in state s_t .
- The agent transitions to the next state s_{t+1} , determined by the external dynamics.

The agent aims to maximize the total reward R_t , defined as follows:

$$R_t = \sum_{k=0}^{\infty} \gamma^k r_{t+k}, \quad (1)$$

- r_{t+k} is the reward received k steps into the future.
- $\gamma \in (0, 1]$ is the discount factor, which determines the importance of future rewards relative to immediate rewards.

II.1.1. State-Action Value Function and Bellman Equation

A key concept in reinforcement learning is the **state-action value function**, also known as the Q -function. This function $Q^\pi(s_t, a_t)$ represents the expected cumulative reward starting from state s_t , taking action a_t , and then following the policy π thereafter:

$$Q^\pi(s_t, a_t) = \mathbb{E} \left[r_t + \gamma \mathbb{E}_{s_{t+1} \sim P(\cdot | s_t, a_t)} \left[\max_{a_{t+1}} Q^\pi(s_{t+1}, a_{t+1}) \right] \right]. \quad (2)$$

- $P(s_{t+1} | s_t, a_t)$ is the transition probability that determines the likelihood of transitioning to s_{t+1} given s_t and a_t .
- $\mathbb{E}[\cdot]$ represents the expected value, which accounts for the stochasticity of rewards and state transitions.

The state value function, $V^\pi(s_t)$, measures the expected cumulative reward starting from state s_t and following policy π :

$$V^\pi(s_t) = \mathbb{E}[R_t | s_t = s]. \quad (3)$$

The relationship between V^π and Q^π is:

$$V^\pi(s_t) = \mathbb{E}_{a_t \sim \pi} [Q^\pi(s_t, a_t)]. \quad (4)$$

II.1.2. Actor-Critic Framework

Modern RL algorithms often use the actor-critic ^[2]framework, which combines two neural networks:

1. Critic Network: This network, parameterized by ϕ , learns the value function $V_\phi^\pi(s_t)$, which estimates the value of a given state.
2. Actor Network: This network, parameterized by θ , learns the policy $\pi_\theta(a_t | s_t)$, which maps states to actions.

The actor network is trained to maximize the following objective function. The objective function for the actor network in reinforcement learning is given as:

$$J_\theta^\pi = \mathbb{E}_{\pi_\theta} [\log \pi_\theta(a_t | s_t) A_t], \quad (5)$$

- J_θ^π (Objective Function): This represents the *objective function* for the actor network. The goal of the actor is to maximize J_θ^π , which corresponds to improving the policy π_θ to choose actions that yield higher rewards.
- \mathbb{E}_{π_θ} (Expectation): This denotes the *expectation* over the policy $\pi_\theta(a_t | s_t)$. It indicates averaging over all possible actions a_t that the policy could take in a given state s_t , weighted by their probability under the policy.

- $\log \pi_\theta(a_t | s_t)$ (Log-Probability of the Action): This is the *logarithm of the probability* of taking action a_t in state s_t under the policy π_θ . The logarithm helps stabilize and smooth the optimization process, aiding in gradient calculations.
- $\pi_\theta(a_t | s_t)$ (Policy): This represents the *policy* parameterized by θ . It is the probability distribution over actions a_t given a state s_t . The reinforcement learning process aims to optimize this policy to maximize rewards.
- A_t (Advantage Function): The advantage function quantifies how much better or worse the chosen action a_t is compared to the average action the policy would take in the same state s_t . It is defined as:

$$A_t = Q^\pi(s_t, a_t) - V^\pi(s_t), \quad (6)$$

- $Q^\pi(s_t, a_t)$: The state-action value function, which represents the expected return for taking action a_t in the state s_t and following policy π thereafter.
- $V^\pi(s_t)$: The state value function, which represents the expected return from state s_t under the policy π .

$$J_\theta^\pi = \mathbb{E}_{\pi_\theta} [\log \pi_\theta(a_t | s_t) A_t], \quad (7)$$

- $A_t = Q^\pi(s_t, a_t) - V^\pi(s_t)$ is the advantage function, which quantifies how much better taking action a_t is compared to the average action at state s_t .

The critic network is trained to minimize the value function loss, which measures the error in estimating the value of states:

$$J^\phi = \min_{\phi} \mathbb{E}_D [(R_t + \gamma V_\phi(s_{t+1}) - V_\phi(s_t))^2], \quad (8)$$

- D is an experience replay buffer that stores previously encountered experiences for training [3].

By combining the actor and critic networks, the agent can learn both the optimal policy and an accurate value function simultaneously, leading to more efficient learning.

II.2. Multi-agent Reinforcement Learning (MARL)

Over the past decade, the field of reinforcement learning has made remarkable progress, successfully addressing various sequential decision-making problems within machine learning. Most of the RL applications, such as game playing, robotics, and autonomous driving, naturally fall under the category of multi-agent reinforcement learning (MARL). In MARL, we have multiple agents interacting with the environment and among themselves. Largely, MARL algorithms can be classified as cooperative, non-cooperative, and partially cooperative based on the nature of interactions between the agents.

In a cooperative setting, all the agents collectively work towards a shared goal, whereas in a non-cooperative setting, they compete against each other to reach their goals first. The partial setting is a mix of both cooperative and competitive agents. In this paper, the focus is set on a cooperative setting, where the agents work together to achieve a shared goal. The agents learn that the best solution is to coordinate with each other rather than act in isolation.

1. **Independent Q-Learning:** An independent MARL framework where each agent learns independently and simultaneously treats every other agent as part of the environment. The agent will ignore the presence of other agents and conduct its learning and updates its policies as if it were the only agent in the environment. The algorithm is fully scalable, but it suffers from non-stationarity and partial observability. Let's break down and understand the three terms better in terms of independent Q-learning. Scalability here means that the algorithm can handle a large number of agents without much impact on the efficiency. Adding more agents doesn't complicate the learning process because the independent learning process doesn't involve any communication between agents. Non-stationarity is a problem relevant in the context of decentralized learning algorithms, which can be better illustrated through the following illustration. Imagine two robots that are learning to navigate in a shared environment. Each robot updates its Q values based on its experiences, and as each robot learns, it affects the environment in ways that influence the other robot's learning process. Partial observability is a fundamental challenge in MARL, because the agent perceives only a limited part of the environment based on its observations made locally. The agent must make decisions without knowing what the other agents can see or knowing the full state.

2. **Off-policy MARL algorithm:** In this off-policy ^[5] MARL algorithm, collaboration is achieved by utilizing a centralized critic network that estimates the state-action value using global observations and actions. A centralized critic network has access to the global state, which allows the critic to judge how the agents' joint action affects the overall objective. It is used during training that handles non-stationary behavior, by stabilizing learning considering the joint impact of all agents' actions. This network is not used during deployment, because once trained, all the agents use their own local policy based on their local observations.
3. **Adaptive Communication and Reward Scaling:** These terms refer to techniques that improve the effectiveness of agent interactions during training. What adaptive communication means is the agent decides when, what and how to share information. Reward Scaling helps the agent focus on what really matters around them. It helps them silence the rewards from faraway neighbors and focus on those who directly impact them. By incorporating such techniques it is possible to reduce unnecessary communication.

The above examples work well with an increasing number of agents, but they might not be as efficient in a dynamic real-world environment. The reason behind this being, the MARL algorithms assume a fixed communication topology, where the connection remains constant between the agents. In real-world situations, the communication links may change, and the algorithm has to be retrained to handle the changes. There is an increasing need for algorithms that prioritize decentralized adaptability, transferable policies, and online learning to handle dynamic situations. In MARL, parameter sharing is widely used with homogeneous agents, which are agents with similar roles and perform the same functions. Single agents RL algorithms such as PPO (Proximal Policy Optimization) and ACKTR (Actor-Critic using Kronecker-Factored Trust Region) are extended to MARL with parameter sharing as MAPPO and MAACKTR. Let's learn how these algorithms work.

1. Proximal Policy Optimization ^[6]: PPO is a policy-gradient algorithm that learns a policy $\pi_\theta(a | s)$ which maps states (s) to actions (a) and the policy is updated to maximize the cumulative reward. It makes the algorithm more stable by preventing too many changes in a single update. It uses a clipped objective function to control how much a policy can change.

The final objective for Proximal Policy Optimization (PPO) is to maximize the clipped surrogate objective:

$$L^{CLIP}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right] \quad (9)$$

- \mathbb{E}_t : The expectation over time steps t .
- $r_t(\theta) = \frac{\pi_\theta(a_t | s_t)}{\pi_{\text{old}}(a_t | s_t)}$: The importance of sampling ratio, comparing the probability of action a_t under the new policy π_θ and the old policy π_{old} .
- \hat{A}_t : The advantage function, which measures how much better an action a_t is compared to the average action in state s_t .
- $\text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)$: A clipping function that restricts $r_t(\theta)$ to the range $[1 - \epsilon, 1 + \epsilon]$, where ϵ is a small hyper parameter (e.g., $\epsilon = 0.2$).
- \min : The minimum function ensures that the objective does not increase beyond the clipped range, thus preventing large policy updates.

If $r_t(\theta) \in [1 - \epsilon, 1 + \epsilon]$, the objective uses $r_t(\theta) \hat{A}_t$. If $r_t(\theta)$ goes outside this range, the clipped value $\text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t$ is used to suppress large updates.

MAPPO is extended to multi agent environments by using parameter sharing. Sharing the parameters allows MAPPO to handle varying number of agents without requiring changes to the architecture or training process. The core equation of MAPPO is the clipped objective function, but the θ represents shared parameters, meaning the policy and value functions are being shared across all agents. The value function is being updated to minimize the loss. MAPPO ensures sufficient multi-agent learning by sharing both the policy and the value function.

$$\mathcal{L}_{\text{value}}(\phi) = \mathbb{E}_t \left[\left(V_\phi(s_t) - G_t \right)^2 \right] \quad (10)$$

Where:

- $V_\phi(s_t)$: The predicted value of state s_t using shared value function parameters ϕ .
- G_t : Discounted return from time step t .

2. Actor-Critic using Kronecker-Factored Trust Region [7]: ACKTR is a method for training that uses an advanced optimization technique to learn faster and more precisely. It has an actor-critic system where the actor decides what action to take based on the current state and the critic evaluates how good the action is and provides feedback.

- Actor: Learns the policy $\pi(a | s)$, which defines the probability of taking an action a given the state s
- Critic: Learns the value function $V(s)$, which estimates how good it is to be in a certain state s .

ACKTR improves its learning process by using natural gradients, which are smarter when compared to regular or first-order gradients. While the first-order gradients treat all parameters equally and do not consider how each parameter affects the environment, natural gradients take the geometry of the environment using the Fisher information matrix. Natural gradients adjust the size and direction of improvement, treating all the parameters usually. The FIM measures how sensitive the outputs of the model are to changes in its parameters, but calculating it exactly is computationally expensive for large models. K-FAC makes it easier through the calculation of the Fisher matrix by assuming the Kronecker product of two smaller matrices.

$$\nabla_\theta^{\text{natural}} J(\theta) = F^{-1} \nabla_\theta J(\theta) \quad (11)$$

$$F \approx A \otimes G \quad (12)$$

$$\nabla_\theta^{\text{natural}} J(\theta) \approx (A^{-1} \otimes G^{-1}) \nabla_\theta J(\theta) \quad (13)$$

- $\nabla_\theta J(\theta)$: The standard gradient of the objective function $J(\theta)$ with respect to the parameters θ . It represents the steepest direction for improving the policy.
- F : The Fisher Information Matrix (FIM), which measures how sensitive the policy's output is to changes in the parameters. It captures the geometry of the parameter space.
- F^{-1} : The inverse of the Fisher Information Matrix. It scales the gradient to adjust for parameter sensitivity, ensuring efficient updates in different directions.
- $\nabla_\theta^{\text{natural}} J(\theta)$: The natural gradient, which is the scaled version of the standard gradient. It accounts for the geometry of the parameter space and provides more efficient updates.
- A : The covariance matrix of activations, calculated as $A = \mathbb{E}[aa^T]$, where a represents the activations.
- G : The covariance matrix of gradients, calculated as $G = \mathbb{E}[gg^T]$, where g represents the gradients of the loss with respect to the layer's outputs.
- \otimes : Kronecker product, a matrix operation that combines two matrices into a larger block matrix.

3. Multi-Agent Advantage Actor-Critic: The Parameter-Sharing Advantage Actor-Critic (MA2C) algorithm is a multi-agent reinforcement learning (MARL) approach that extends the traditional Actor-Critic (A2C) framework for cooperative multi-agent environments. In MA2C, all agents share the same policy and value networks (parameter sharing), which reduces computational overhead and ensures consistency among agents. Each agent uses its local observations and receives localized rewards, enabling decentralized decision-making. This makes MA2C scalable for environments with many agents, such as fleet management, where vehicles (agents) need to allocate resources efficiently. In the fleet management problem, MA2C optimizes fleet performance by maximizing efficiency (e.g., reducing wait times or fuel consumption) through cooperative behavior. Parameter sharing allows the algorithm to scale without requiring individual networks for each agent. Experiments confirm MA2C's strong performance due to its efficiency, scalability, and ability to handle decentralized tasks. In the paper, MA2C is used also as a benchmark algorithm to compare against the proposed method in the on-ramp merging problem.

The paper introduces a **novel on-policy Multi-Agent Reinforcement Learning (MARL)** algorithm tailored for the challenging *on-ramp merging problem* in autonomous driving. This algorithm is designed to address gaps in efficiency and safety by incorporating the following key features:

- **Action Masking:** Prevents invalid or unsafe actions (e.g., lane changes to non-existent lanes or exceeding speed limits), improving stability and reducing training noise.
- **Priority-Based Safety Supervisor:** Ensures safety-critical decisions by regulating actions based on traffic scenarios and system constraints.
- **Parameter Sharing:** All agents share a common policy network, which enhances scalability, ensures consistency and reduces computational costs.
- **Local Reward Shaping:** Rewards are shaped to reflect local interactions (e.g., with neighboring vehicles), enabling efficient credit assignment and reducing communication overhead.

The paper compares the performance of this proposed algorithm against **benchmarks like MA2C** (from Section III) in Section V, showcasing its improvements in safety, efficiency, and overall merging behavior. This novel approach combines scalability, safety, and real-time decision-making for autonomous vehicle coordination.

III. RAMP MERGING AS MARL

The on-ramp merging problem is formulated as a partially observable Markov decision process (POMDP) in this section. To solve the defined POMDP, the actor-critic-based MARL algorithm is presented, including action masking, effective reward function design, and parameter-sharing mechanisms.

III.1. Problem Formulation

This study models the on-ramp merging scenario in mixed traffic as a model-free multi-agent network. It is assumed that the graph $G = (V, E)$ represents the network, where each agent $i \in V$ interacts with its neighbors $N_i := \{j \mid e_{ij} \in E\}$ through edge connections e_{ij} . Let $S := \{x_i \mid x_i \in S\}$ and $A := \{x_i \mid x_i \in A\}$ signify the global state space and action space, respectively. The system's underlying dynamics can be described using the state transition distribution $P : S \times A \times S \rightarrow [0, 1]$. A decentralized MARL framework is adopted, where each agent i observes only a portion of the environment, such as surrounding vehicles. This design aligns with real-world traffic dynamics, where AVs can only perceive or interact with vehicles in their immediate vicinity. As a result, the overall system is effectively modeled as a POMDP M_G , which is defined by the tuple $(\{A_i, S_i, R_i\}_{i \in V}, T)$.

• **Action Space** Based on the designs described in the paper^[9], an agent i has a defined action space A_i that consists of high-level control instructions like turn left, turn right, cruise, acceleration, and deceleration. Lower-level controllers produce the necessary steering and throttle signals needed to navigate autonomous vehicles (AVs) once a high-level action is chosen. Figure (to be inserted) shows the general system configuration comprising these hierarchical control systems. The joint action space of all the AVs captures the overall action space for the system, implying,

$$A = A_1 \times A_2 \times \dots \times A_N. \quad (14)$$

• **State Space** The agent i is described as a matrix with dimensions $N_{N_i} \times W$ where N_{N_i} is the number of cars observed by the agent and W defines the number of attributes utilized to characterize every observed vehicle. These attributes comprise:

- **ispresent:** A binary variable denoting whether a vehicle is observable close to the ego vehicle.
- **x:** The observed vehicle's longitudinal relative location to the ego vehicle.
- **y:** The observed vehicle's lateral relative location to the ego vehicle.
- **v_x:** The observed vehicle's longitudinal velocity relative to the ego vehicle.
- **v_y:** The observed vehicle's lateral velocity relative to the ego vehicle.

Based on the agent's local sensing capabilities, it is expected that the ego vehicle can only observe "neighboring vehicles", defined as the N_{N_i} nearest cars located within a longitudinal distance of 150 meters from the ego vehicle. Setting $N_{N_i} = 5$ produced the highest performance in the case of the on-ramp merging situation shown in Fig. 1. The global state of the system is expressed as the Cartesian product of the states of every individual agent, obtained from:

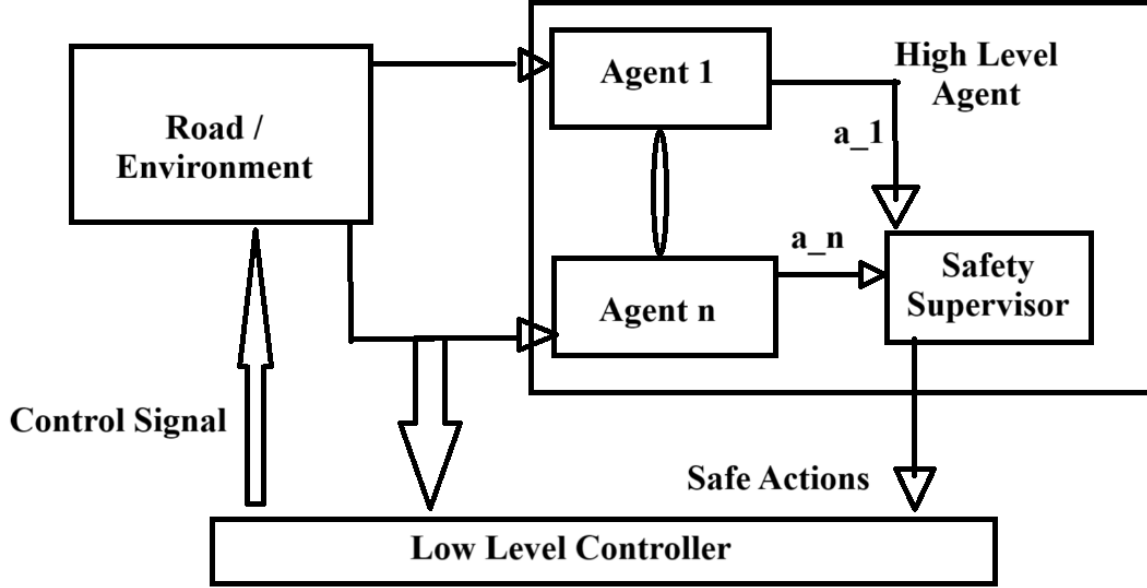


Figure 1. Flow chart of system and simulation setup

$$S = S_1 \times S_2 \times \dots \times S_N. \quad (15)$$

- **Reward Function**

Training reinforcement learning agents to show the intended behavior depends critically on the reward function R_i . The aim is to let agents negotiate the merging space effectively and safely. Agent i gets the reward at time t from:

$$r_{i,t} = w_c r_c + w_s r_s + w_h r_h + w_m r_m, \quad (16)$$

where w_c , w_s , w_h , and w_m are positive weighting factors assigned to the collision penalty (r_c), speed evaluation (r_s), headway time evaluation (r_h), and merging cost evaluation (r_m), respectively. Safety is the priority here, hence w_c gets a far higher value than the other weights. The reward components are defined as follows:

- Collision evaluation (r_c): Assigned a value of -1 in case of a collision and 0 otherwise.
- Speed evaluation (r_s): Defined as:

$$r_s = \min \left\{ \frac{v_t - v_{\min}}{v_{\max} - v_{\min}}, 1 \right\}, \quad (17)$$

where v_t is the current speed of the ego vehicle. Following recommendations from the US Department of Transportation (20–30 m/s) and the speed range observed in the Next Generation Simulation (NGSIM) dataset (minimum speed 6–8 m/s), the minimum and maximum speeds are set as $v_{\min} = 10$ m/s and $v_{\max} = 30$ m/s, respectively.

- Headway time evaluation (r_h): Calculated as:

$$r_h = \log \frac{d_{\text{headway}}}{t_h v_t}, \quad (18)$$

where d_{headway} is the distance between the current car and the vehicle in front, t_h is a preset safe headway period. The ego vehicle is penalized if t_h is violated, that is, if the headway time is too short, while else it is rewarded. This work sets t_h as advised in [10].

- **Merging cost (r_m):** This component is designed to penalize excessive waiting time on the merging lane, which can lead to traffic deadlocks. The merging cost is defined as:

$$r_m = -\exp\left(-\frac{(x-L)^2}{10L}\right), \quad (19)$$

where x represents the distance traveled by the ego vehicle on the ramp, and L denotes the total length of the ramp (refer to Fig. 1).

As shown in Fig. 3, the penalty increases as the ego vehicle gets closer to the merging endpoint. This formulation encourages efficient use of the ramp by minimizing delays and avoiding congestion.

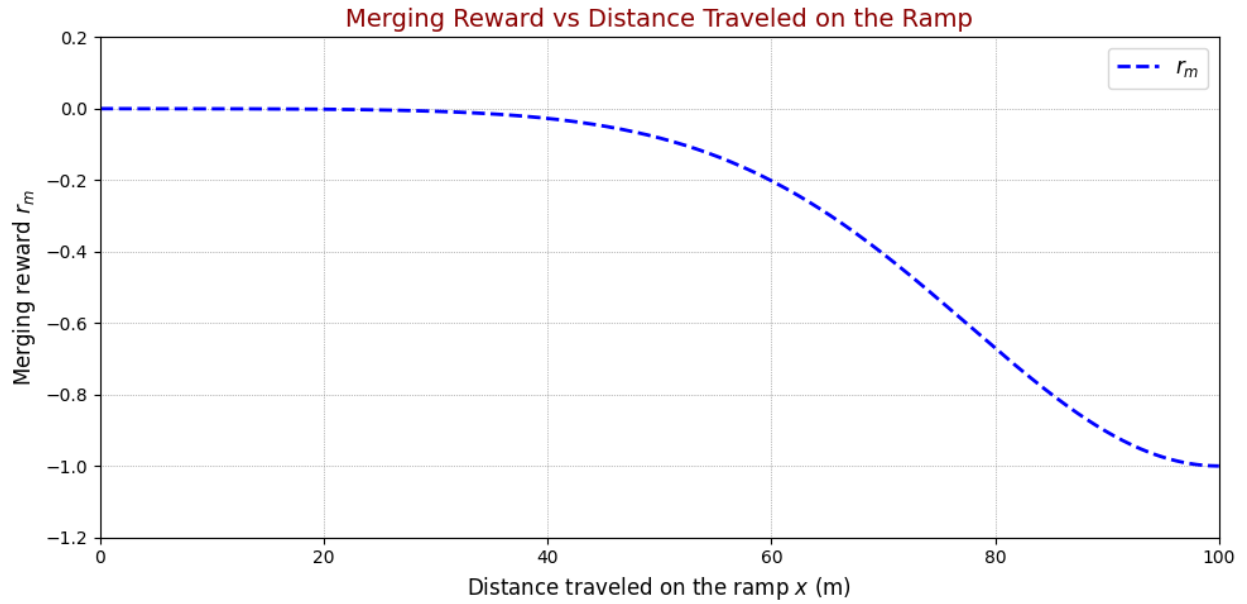


Figure 2. Variation of merging reward with distance x

III.2. MA2C for CAVs

Maximizing the worldwide reward $R_{g,t}$ defined as the total of individual agent rewards is the main goal of the cooperative MARL framework:

$$R_{g,t} = \sum_{i=1}^N r_{i,t}. \quad (20)$$

Under ideal circumstances, every agent would get the same average worldwide reward computed as:

$$r_{i,t} = \frac{1}{N} R_{g,t}, \quad \text{for } i = 1, 2, \dots, N. \quad (21)$$

This common reward system does, however, have major drawbacks. First, aggregating the global reward brings significant latency and increases communication overhead, which might be troublesome in systems with real-time restrictions, including autonomous vehicles (AVs). Second, the credit assignment issue results from a single global reward failing to represent the unique efforts of every vehicle adequately. This can restrict scalability to a small number of agents and compromise learning efficiency.

A localized incentive assignment method is proposed to address these issues. Under this approach, the reward of every ego vehicle is found only based on those of its surrounding neighbors. The agent i gets reward at time t as follows:

$$r_{i,t} = \frac{1}{|v_i|} \sum_{j \in v_i} r_{j,t}, \quad (22)$$

where $v_i = i \cup N_i$ represents the set of the ego vehicle i and its neighboring vehicles, and $|v_i|$ is the size of this set.

This localized reward system guarantees that agents concentrate mostly on interactions with nearby cars, therefore assigning rewards more relevant to the success or failure of a task. For tasks involving road-based interactions, this approach is especially successful since vehicles mostly affect their immediate surroundings.

The network backbone used in this study is depicted in Fig. 4, where the actor and critic networks share the same low-level representations. The policy loss and value function error are integrated into a single overall loss function, utilizing shared network parameters. The overall loss function is defined as:

$$J(\theta_i) = J^{\pi_{\theta_i}} - \beta_1 J^{V_{\phi_i}} + \beta_2 H(\pi_{\theta_i}(s_t)), \quad (23)$$

where β_1 and β_2 are weighting coefficients for the value function loss and the entropy regularization term, respectively. The entropy regularization term is expressed as:

$$H(\pi_{\theta_i}(s_t)) = \mathbb{E}_{\pi_{\theta_i}}[-\log(\pi_{\theta_i}(s_t))], \quad (24)$$

which encourages agents to explore new states, thereby improving learning efficiency.

The policy loss, $J^{\pi_{\theta_i}}$, is given by:

$$J^{\pi_{\theta_i}} = \mathbb{E}_{\pi_{\theta_i}} \left[\log(\pi_{\theta_i}(a_{t,i} | s_{t,i})) A_{t,i}^{\pi_{\theta_i}} \right], \quad (25)$$

where $A_{t,i}^{\pi_{\theta_i}}$ represents the advantage function, defined as:

$$A_{t,i}^{\pi_{\theta_i}} = r_{t,i} + \gamma V_{\phi_i}(s_{t+1,i}) - V_{\phi_i}(s_{t,i}). \quad (26)$$

The loss for updating the state value function, $J^{V_{\phi_i}}$, is:

$$J^{V_{\phi_i}} = \min_{\phi_i} \mathbb{E}_t \left[(r_{t,i} + \gamma V_{\phi_i}(s_{t+1,i}) - V_{\phi_i}(s_{t,i}))^2 \right]. \quad (27)$$

Separate experience replay buffers are maintained for each agent, while the shared policy network parameters are updated collectively across agents. This approach is especially effective in cooperative multi-agent reinforcement learning scenarios involving mixed traffic. Mini batches of sampled trajectories are employed to iteratively update the network parameters using Equation (23), ensuring reduced variance and enhanced training stability.

III.3. DNN Settings

Fig. 4 shows the deep neural network (DNN) architectural layout. To enhance scalability and robustness, the observation $s_{i,t}$ is divided into three sub-groups based on their physical characteristics: $s_{i,t}^1 \cup s_{i,t}^2 \cup s_{i,t}^3$ where $s_{i,t}^1$, $s_{i,t}^2$, and $s_{i,t}^3$ correspond to *ispresent states*, *position states*, and *speed states*, respectively. Every sub-group passes an independent fully connected (FC) layer, and the resultant encoded representations are concatenated into a single vector. Actor and critic networks share a 128-neuron FC layer derived from this concatenated vector. In the standard setting, the actor-network produces logits l_i that are passed through a Softmax layer to generate action probabilities: $\pi_{\theta_i}(s_i) = \text{softmax}([l_1, l_2, l_3, l_4, l_5])$ where $\pi_{\theta_i}(s_i)$ represents the probability distribution over actions. Next, actions from this distribution—that is, $a_i \sim \pi_{\theta_i}(s_i)$ are sampled. However, this sampling process introduces some challenges:

- **Invalid/unsafe actions:** Actions judged hazardous or invalid still carry non-zero probability. These dangerous behaviors could be chosen during training depending on the stochastic character of sampling, therefore causing possibly negative behaviors or system instability.
- **Policy updates with invalid actions:** Execution of undesirable actions affects and impedes policy training since the gathered experience is useless and misleading.

The proposed invalid action masking technique in [8] is used to help reduce these problems. This method "masks out" illegal behavior, therefore guaranteeing that only legitimate actions are gathered. This erroneous action mask, which is included in the actor network for strong policy training, is shown in Figure 3.

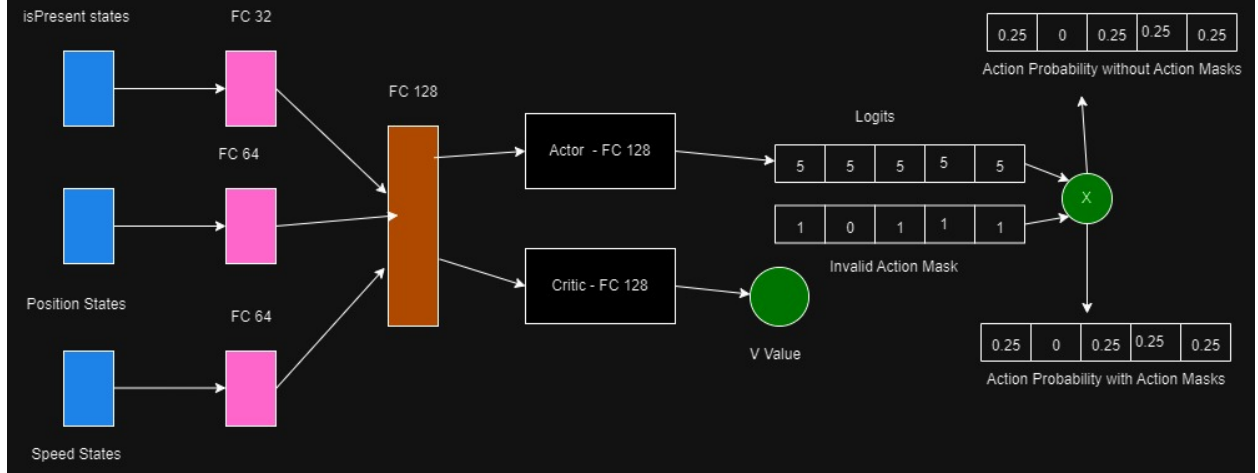


Figure 3. Schematic representation of the proposed DNN architecture

IV. PRIORITY-BASED SAFETY ENHANCEMENT

The rule-based action masking system outlined earlier is effective at preventing improper actions, such as invalid lane changes or unsafe acceleration. However, it does not fully address the risks of inter-vehicle or vehicle-obstacle collisions. To manage these risks in complex, dynamic, and congested mixed-traffic environments, a more comprehensive safety-enhancing approach is proposed. This system leverages vehicle dynamics and multi-step predictions over a defined *prediction horizon* (T_n) to ensure collision-free behavior. The primary objective of the safety supervisor is to predict potential collisions within the prediction horizon and adjust unsafe actions proactively. This involves modeling the behavior of both Human-Driven Vehicles (HDVs) and Autonomous Vehicles (AVs) to ensure safe and efficient navigation. **1. Human-Driven Vehicles (HDVs)** For HDVs, the safety supervisor uses the following models to predict behavior:

- **Intelligent Driver Model (IDM):** This model estimates longitudinal acceleration (a_i) of HDVs based on their current speed (v_i) and headway distance ($d_{headway}$) to the vehicle ahead:

$$a_i = a_{max} \left(1 - \left(\frac{v_i}{v_{des}} \right)^4 - \left(\frac{s^*(v_i, \Delta v)}{d_{headway}} \right)^2 \right), \quad (28)$$

where $s^*(v_i, \Delta v)$ represents the desired safe distance, which depends on speed, time headway (T), and braking deceleration:

$$s^*(v_i, \Delta v) = s_0 + v_i T + \frac{v_i \Delta v}{2\sqrt{a_{max}b}}. \quad (29)$$

- **MOBIL Lane-Change Model:** This model predicts when lane changes are safe and beneficial by evaluating the acceleration gain (a_{gain}) associated with the maneuver:

$$a_{gain} > a_{threshold}. \quad (30)$$

Lane changes occur only if they provide a significant benefit and are safe for all vehicles involved.

2. Autonomous Vehicles (AVs) For AVs, high-level decisions such as acceleration or lane changes are determined by a *Multi-Agent Reinforcement Learning (MARL) agent*. These high-level actions are then implemented through *low-level PID controllers*, which generate precise control signals. The PID controller computes these signals as:

$$u(t) = K_p e(t) + K_i \int e(t) dt + K_d \frac{de(t)}{dt}, \quad (31)$$

where $e(t)$ is the error between the desired and actual state, and K_p , K_i , and K_d are proportional, integral, and derivative gains, respectively. This ensures that the AVs execute decisions smoothly and reliably.

IV. 1. Trajectory Propagation and Motion Primitives

Vehicle trajectories are propagated using the *kinematic bicycle model*, which simulates vehicle motion based on velocity (v) and steering angle (δ):

$$\dot{x} = v \cos(\theta), \quad \dot{y} = v \sin(\theta), \quad \dot{\theta} = \frac{v}{L} \tan(\delta), \quad (32)$$

where (x, y) is the vehicle's position, θ is the orientation, L is the wheelbase, and δ is the steering angle. These predicted trajectories, referred to as *motion primitives*, allow the supervisor to evaluate future paths and ensure collision-free navigation.

The safety supervisor integrates predictive models for both HDVs and AVs and uses motion primitives to evaluate potential risks within the prediction horizon. The framework provides:

- Proactive collision avoidance by simulating future vehicle trajectories.
- Accurate modeling of HDV behavior using IDM and MOBIL.
- Efficient decision-making for AVs using MARL and PID controllers.

The proposed framework and simulation setup are illustrated in **Fig.1**. This system ensures safe, efficient, and reliable operation in high-density, mixed-traffic environments by effectively addressing both safety and efficiency.

IV. 2. Priority Assignment

Based on the joint motion primitives of all AVs, it is possible to forecast if a collision might arise inside the next T_n using the motion models for HDVs. This framework designs the safety-enhancing mechanism using the cooperative efforts of all AVs. While determining a joint safe action to prevent collisions is computationally difficult, spotting possible collisions from a given joint activity is rather easy. This results from the size of the action space, $|A_i|^N$ where N is the overall count of AVs. Particularly in cases with several AVs, the size of the action space expands exponentially as N increases, rapidly rendering the job computationally intractable. To address time constraints effectively, we propose a sequential, priority-based safety enhancement scheme that balances computational efficiency with real-time applicability. The core idea is to coordinate AVs in a sequential order, giving higher priority to AVs with smaller safety margins. For instance, AVs near the end of the merging lane or close to the defined safety boundary (e.g., with small time headway) are assigned higher priorities. The following rationales are considered for assigning priorities:

1. Vehicles on the merging lane are given higher priority compared to vehicles on the through lane, as merging is a time-critical process, especially near the lane end.
2. Vehicles closer to the merging lane's end are prioritized higher because they are more likely to encounter collisions or cause deadlocks^[11].
3. Vehicles with shorter time headway are assigned higher priority since they are at greater risk of colliding with preceding vehicles.

Based on these principles, the priority index p_i for the ego vehicle i is defined as:

$$p_i = \alpha_1 p_m + \alpha_2 p_d + \alpha_3 p_h + \sigma_i, \quad (33)$$

where α_1 , α_2 , and α_3 are positive weighting factors for the merging priority metric (p_m), the distance-to-end metric (p_d), and the time headway metric (p_h), respectively. The term $\sigma_i \sim \mathcal{N}(0, 0.01)$ is a small random variable added to avoid ties between two vehicles with the same priority index.

The metrics are defined as follows: **Merging priority metric (p_m):**

$$p_m = \begin{cases} 0.5, & \text{if on the merging lane,} \\ 0, & \text{otherwise.} \end{cases} \quad (34)$$

Distance-to-end priority metric (p_d):

$$p_d = \begin{cases} \frac{x}{L}, & \text{if on the merging lane,} \\ 0, & \text{otherwise.} \end{cases} \quad (35)$$

Here, x represents the distance the ego vehicle has traveled on the merging lane, and L is the total length of the ramp.

Time headway priority metric (p_h): The time headway priority is computed using the time headway definition:

$$p_h = -\log \left(\frac{d_{\text{headway}}}{t_h v_t} \right), \quad (36)$$

where d_{headway} is the headway distance, t_h is the predefined headway time, and v_t is the ego vehicle's speed.

IV. 3. Priority-based Safety Supervisor

The Priority-Based Safety Supervisor is an advanced system designed to ensure the safe operation of Autonomous Vehicles (AVs) in dynamic, mixed-traffic environments. It works by analyzing predicted vehicle motions, assigning priorities to AVs based on safety-critical factors, and adjusting their actions to prevent collisions. The supervisor processes vehicles sequentially, starting with those in the most critical situations, and ensures both computational efficiency and real-time applicability.

At each time step t , the supervisor generates a priority list P_t for all AVs, ranked by their priority scores. These scores are determined using metrics such as proximity to the merging lane endpoint and time headway to the vehicle ahead. The AV with the highest priority ($P_t[0]$) is selected first for a safety check. This involves analyzing its exploratory action, which is the proposed action generated by its action network, to determine if it could lead to a collision.

Collision prediction is based on the motions of both AVs and Human-Driven Vehicles (HDVs). The motions of HDVs are estimated using the *Intelligent Driver Model (IDM)* for acceleration and the *MOBIL model* for lane changes. For lower-priority AVs, their motions are assumed to remain the same as in the last time step. Using these predictions, the supervisor checks if the trajectories of any two vehicles come within a safety threshold over a prediction horizon (T_n). If no collision is detected, the exploratory action of $P_t[0]$ is deemed safe and approved.

If a collision is predicted, the exploratory action is marked as unsafe. The supervisor then evaluates alternative actions from the set of valid actions (A_{valid}) for the AV. Each action is assessed based on its safety margin ($d_{\text{sm},k}$), which is the minimum safe distance to surrounding vehicles. For lane-changing actions, the safety margin is calculated as the minimum distance to preceding and following vehicles in both the current and target lanes. For speed adjustments, the safety margin is the minimum headway distance to the vehicle ahead. The action with the highest safety margin is selected as the safest alternative, using the formula:

$$a'_t = \arg \max_{a_t \in A_{\text{valid}}} \left(\min_{k \in T_n} d_{\text{sm},k} \right) \quad (37)$$

After selecting a safe action, the trajectory of $P_t[0]$ is updated, and the vehicle is removed from the priority list. The process then moves to the next-highest-priority vehicle ($P_t[1]$), and the same safety-check procedure is repeated. For lower-priority AVs, the updated trajectories of higher-priority AVs are considered during the safety checks. This process continues until all vehicles in the priority list have been processed.

The system is designed to ensure collision-free motion by addressing unsafe exploratory actions and selecting the safest alternatives dynamically. It balances safety and efficiency by processing vehicles in order of priority and focusing on the most critical scenarios. For instance, in the case of lane changes, the supervisor evaluates the minimum distance to vehicles in both the current and target lanes, while for speed adjustments, it considers the headway distance to the vehicle ahead.

By employing this structured approach, the Priority-Based Safety Supervisor ensures that AVs can make safe and cooperative decisions in real-world traffic scenarios. It proactively detects and mitigates potential collisions, enhances computational efficiency, and enables real-time decision-making in complex environments, making it a robust solution for traffic safety management.

This algorithm ensures the safe operation of autonomous vehicles (AVs) in dynamic and complex traffic scenarios. It prioritizes vehicles based on their risk levels and iteratively ensures their actions are safe, modifying unsafe actions when necessary. The steps of the algorithm are as follows:

IV. 4. Algorithm 1 - Priority-Based Safety Supervisor

Algorithm 1 Priority-Based Safety Supervisor

Parameter : $L, \alpha_1, \alpha_2, \alpha_3, t_h, w, T_n$

Output : $a_i, i \in \mathcal{V}$

```

1 for  $i = 0$  to  $N$  do
2   Compute the priority scores according to Eq. (10) Rearrange ego vehicles to list  $P_t$  according to their priority
   scores
3 for  $j = 0$  to  $|P_t|$  do
4   Obtain the highest-priority vehicle  $P_t[0]$  Find its neighboring vehicles  $N_{P_t[0]}$  Predict trajectories  $\zeta_v, v \in P_t[0] \cup N_{P_t[0]}$ 
   for  $T_n$  time steps if trajectories are overlapped then
5     Replace the risky action as  $a_t \leftarrow a'_t$  according to Eq. (13) Replace the trajectory  $\zeta_{P_t[0]}$  with  $\zeta'_{P_t[0]}$ 
6   Remove  $P_t[0]$  from  $P_t$  Update  $P_t[i] \leftarrow P_t[i+1], i = 1, 2, \dots$ 

```

IV.4. 1 Input Parameters and Output

1. Parameters:

- L : Total length of the merging ramp.
- $\alpha_1, \alpha_2, \alpha_3$: Weighting factors for merging priority, distance-to-end priority, and time headway priority, respectively.
- t_h : Time headway threshold.
- w : Safety margin parameter.
- T_n : Time horizon for collision prediction.

2. Output:

- Safe actions a_i for each AV i in the system.

IV. 4. 2. Algorithm Steps

Step 1: Compute Priority Scores: For each vehicle i , calculate the priority score p_i using the formula:

$$p_i = \alpha_1 p_m + \alpha_2 p_d + \alpha_3 p_h + \sigma_i \quad (38)$$

- p_m : Merging priority.
- p_d : Distance-to-end priority.
- p_h : Time headway priority.
- σ_i : Small random value to avoid ties.

Step 2: Create a Priority List: Arrange vehicles in a priority list P_t in descending order of their priority scores. The highest-priority vehicle is placed at the top.

Step 3: Process Each Vehicle in Priority Order:

- Select the highest-priority vehicle $P_t[0]$.
- Identify its neighboring vehicles $N_{P_t[0]}$, which include other AVs and human-driven vehicles (HDVs).

Step 4: Predict Future Trajectories: Predict the trajectories ζ_v for T_n steps for the selected vehicle $P_t[0]$ and its neighboring vehicles $N_{P_t[0]}$.

Step 5: Check for Potential Collisions:

- If no trajectory overlap is detected, the exploratory action a_t is deemed safe and selected.
- If a trajectory overlap is detected, replace the exploratory action a_t with a safe action a'_t determined by:

$$a'_t = \arg \max_{a_t \in A_{\text{valid}}} \left(\min_{k \in T_n} d_{sm,k} \right) \quad (39)$$

- A_{valid} : Set of valid actions.
- $d_{sm,k}$: Safety margin at step k .

Step 6: Update Trajectory: Update the trajectory $\zeta_{P_i[0]}$ with the new safe trajectory $\zeta'_{P_i[0]}$.

Step 7: Remove Processed Vehicle:

- Remove $P_i[0]$ from the priority list.
- Update the priority list by moving the next highest-priority vehicle to the top.

Step 8: Repeat Until All Vehicles are Processed: Continue steps 3 to 7 until the priority list is empty.

Remark 1: Implementation with Vehicle-to-Infrastructure (V2I) Communication

The priority-based safety supervision scheme utilizes *Vehicle-to-Infrastructure (V2I)* communication to coordinate Autonomous Vehicles (AVs) and Human-Driven Vehicles (HDVs) near merging ramps. A central communication station monitors traffic conditions and calculates priority scores for AVs based on real-time observations. At each discrete time step t , the priority score for an AV i is computed as:

$$p_i = \alpha_1 p_m + \alpha_2 p_d + \alpha_3 p_h + \sigma_i, \quad (40)$$

p_m : Merging priority metric, defined as:

$$p_m = \begin{cases} 0.5, & \text{if on the merging lane,} \\ 0, & \text{otherwise.} \end{cases} \quad (41)$$

p_d : Distance-to-end metric, calculated as:

$$p_d = \begin{cases} \frac{x}{L}, & \text{if on the merging lane,} \\ 0, & \text{otherwise.} \end{cases} \quad (42)$$

Here, x is the distance traveled on the merging lane, and L is the total length of the ramp.

p_h : Time headway priority metric, computed as:

$$p_h = -\log \left(\frac{d_{\text{headway}}}{t_h v_t} \right), \quad (43)$$

- d_{headway} : Headway distance to the vehicle ahead.
- t_h : Predefined safe time headway.
- v_t : Vehicle speed.
- $\alpha_1, \alpha_2, \alpha_3$: Positive weighting factors that balance the impact of p_m , p_d , and p_h , respectively.
- σ_i : A small random variable ($\sigma_i \sim \mathcal{N}(0, 0.01)$) to prevent ties between vehicles with identical priority scores.

Once the priority scores are computed:

- AVs send their exploratory actions (proposed movements) to the central station.
- The station applies Algorithm 1 to process AVs sequentially based on their priority scores and generate safe actions.

The sequential nature of Algorithm 1 ensures computational efficiency. For instance, with a *prediction horizon* $T_n = 8$, the system processes actions in approximately 28.13 milliseconds per decision in complex traffic conditions (refer to Table II in Section V-C). This computational speed makes real-time implementation feasible, enabling updated control commands within one sampling step. Future optimizations will further enhance computational efficiency.

Remark 2: Importance of the Prediction Horizon T_n

The prediction horizon T_n is a critical hyperparameter that governs how far into the future the safety supervisor predicts potential collisions. The effects of small and large T_n are as follows:

Impact of Small T_n

- A small T_n makes the safety supervisor “short-sighted,” potentially failing to anticipate imminent collisions.
- This limitation could lead to no feasible solutions after a few steps.

Impact of Large T_n

- A large T_n accumulates uncertainties in human-driven vehicle (HDV) behavior due to noise and variations in actual vehicle motion.
- This results in overly conservative decisions, sacrificing performance to guarantee safety over extended horizons.

Optimal Value of T_n

To strike a balance, cross-validation experiments indicate that $T_n = 8$ or $T_n = 9$ provides the optimal trade-off. This value ensures:

- Adequate foresight to predict and resolve potential collisions.
- Computational efficiency for real-time decision-making.

IV. 5 . Algorithm 2 - MARL for AVs with Safety Supervisor

- **Parameters:** $\gamma, \eta, T, M, \beta_1, \beta_2$.
- **Outputs:** θ .

The algorithm begins by initializing the necessary parameters, such as the initial state of the system, the current time step, and an empty replay buffer to store experiences. The training process is divided into multiple epochs, each consisting of a fixed number of time steps. For each time step within an epoch, every agent (vehicle) collects information about its current state and determines its next action using a policy network enhanced with action masking to prevent invalid moves. These exploratory actions are then sent to the priority-based safety supervisor, which checks if they are safe. If an action is deemed unsafe, the safety supervisor replaces it with a predefined safer alternative based on a specific rule. The agent then executes the safe action, and the experience, including the state, action, reward, and resulting state, is stored in the replay buffer.

This process continues until the end of the time step sequence for the epoch. Afterward, the algorithm evaluates whether the episode is complete, either because the task has been finished successfully or because a collision or unsafe event occurred (signaled by the “DONE” flag). Once the episode is complete, the policy network of each agent is updated using the collected experiences stored in the buffer. The agents then reset to their initial states, and the next epoch begins. This cycle of exploration, safety validation, action execution, experience storage, and learning repeats for the specified number of epochs, allowing the agents to iteratively refine their policies and improve their performance over time.

The combination of action masking, safety supervision, and reinforcement learning ensures that the agents can explore the environment safely while progressively learning better strategies. The safety supervisor guarantees that risky actions are intercepted and replaced with safe ones, preventing accidents while still enabling the agents to learn from their experiences. This structured approach balances exploration and exploitation while maintaining the safety of the learning process.

Algorithm 2 MARL for AVs with Safety Supervisor

```
1: Initialize  $s_0, t \leftarrow 0, D \leftarrow \emptyset$ ;
2: for  $j = 0$  to  $M - 1$  do
3:   for  $t = 0$  to  $T - 1$  do
4:     for  $i \in \mathcal{V}$  do
5:       Observe  $s_i$ ;
6:       Update  $a_{i,t} \sim \pi_{\theta_i}(\cdot | s_i)$  with action masking;
7:     end for
8:     for  $i \in \mathcal{V}$  do
9:       Check the actions by Algorithm 1;
10:      if safe then
11:        Execute  $a_{i,t}$ ;
12:        Update  $D_i \leftarrow (s_{i,t}, a_{i,t}, r_{i,t}, v_{i,t})$ ;
13:      else
14:        Update  $a_{i,t} \leftarrow a'_{i,t}$  and execute  $a'_{i,t}$ ;
15:        Update  $D_i \leftarrow (s_{i,t}, a'_{i,t}, r_{i,t}, v_{i,t})$ ;
16:      end if
17:    end for
18:    Update  $t \leftarrow t + 1$ ;
19:    if DONE then
20:      for  $i \in \mathcal{V}$  do
21:        Update  $\theta_i \leftarrow \theta_i + \eta \nabla_{\theta_i} J(\theta_i)$ ;
22:      end for
23:    end if
24:    Initialize  $D_i \leftarrow \emptyset, i \in \mathcal{V}$ ;
25:  end for
26:  Update  $j \leftarrow j + 1$ ;
27: end for
28: Update  $s_0, t \leftarrow 0; = 0$ 
```

Step 1: Initialization: Initialize the initial state s_0 . Set the time $t \leftarrow 0$. Initialize the replay buffer $D \leftarrow \emptyset$, which will store experiences for learning.

Step 2: Outer Loop Over Epochs: The algorithm runs for M epochs. Each epoch represents a complete training cycle.

Step 3: Inner Loop Over Time Steps: For each epoch, the algorithm progresses through T time steps. At each time step:

Step 3.1: Agent Observation and Action Selection

- For each agent $i \in \mathcal{V}$, observe the current state s_i .
- Select an action $a_{i,t}$ using the policy $\pi_{\theta_i}(\cdot | s_i)$ with action masking.
- Action masking ensures that invalid or impractical actions are excluded from consideration.

Step 4: Safety Check: The priority-based safety supervisor checks the proposed actions $a_{i,t}$ for all agents:

- **If the action is safe:**
 - The agent executes $a_{i,t}$.
 - The experience tuple $(s_{i,t}, a_{i,t}, r_{i,t}, v_{i,t})$, which includes the state, action, reward, and next state, is stored in the replay buffer D_i .
- **If the action is unsafe:**
 - Replace the unsafe action $a_{i,t}$ with a safe action $a'_{i,t}$, which is executed instead.
 - Store the experience $(s_{i,t}, a'_{i,t}, r_{i,t}, v_{i,t})$ in D_i .

Step 5: Update Time Step: Increment the time t to proceed to the next step.

Step 6: Policy Update (When Episode Ends)

- Check if the episode has ended, which can occur if:
 - The task is completed, or
 - A collision happens (flagged by the *DONE* signal).
- If the episode ends, update the policy parameters θ_i of each agent using the collected experiences:

$$\theta_i \leftarrow \theta_i + \eta \nabla_{\theta_i} J(\theta_i),$$

where $J(\theta_i)$ is the objective function and η is the learning rate.

Step 7: Reset Buffers: After the episode ends, clear the replay buffers D_i for all agents to prepare for the next episode.

Step 8: Epoch Update: Once all time steps in an epoch are completed, increment the epoch counter j and begin the next epoch.

Step 9: Final Reset: After all epochs are completed, reset the initial state s_0 and the time t to their default values.

V. NUMERICAL EXPERIMENTS

This section evaluates the performance of the proposed Multi-Agent Reinforcement Learning (MARL) algorithm in a simulated road merging scenario. The evaluation focuses on two key aspects: training efficiency (how quickly and effectively the algorithm learns) and collision rate (its ability to avoid accidents during the merging process).

The road used for the simulation is 520 meters long, with a merging lane that begins at 320 meters and stretches for 100 meters. There are 12 evenly spaced spawn points along the main lane and the ramp lane between 0 and 220 meters. Vehicles that move beyond the end of the road are removed from the display, but their motion dynamics continue to be updated in the simulation. The scenario is tested under three levels of traffic density: easy, medium, and hard. In easy mode, 1-3 autonomous vehicles (AVs) and 1-3 human-driven vehicles (HDVs) are introduced. Medium mode includes 2-4 AVs and 2-4 HDVs, while hard mode features 4-6 AVs and 3-5 HDVs.

Each training episode begins with a randomized number of AVs and HDVs, which are assigned starting positions with slight variations (noise of ± 1.5 meters) around their spawn points. The vehicles' initial speeds are randomly selected between 25 and 27 meters per second. The AVs make decisions and take actions every 0.2 seconds, corresponding to a control frequency of 5 Hz. To simulate real-world conditions, HDVs' acceleration and steering inputs are perturbed with 5% random noise.

The MARL algorithm is trained over 2 million steps, corresponding to approximately 20,000 episodes, with each episode lasting up to 100 steps. Training is conducted with three different random seeds to ensure the robustness of the results, with the same seeds shared among all agents. The algorithm is evaluated every 200 training episodes by running three evaluation episodes, allowing periodic assessment of its performance and progress.

The reward function used for training is designed to encourage safe and efficient driving. Agents are rewarded based on their ability to avoid collisions (weight of 200), drive safely (weight of 1), maintain a safe distance from other vehicles (weight of 4), and merge successfully (weight of 4). The algorithm uses a discount factor (γ) of 0.99 to balance immediate and future rewards and a learning rate (η) of 5×10^{-4} to control the speed of updates to the agents' policies. Additional parameters, including priority coefficients ($\alpha_1, \alpha_2, \alpha_3$) set to 1 and loss function weights ($\beta_1 = 1$ and $\beta_2 = 0.01$), are used to balance learning priorities.

The simulation is built on a modified version of the highway-env simulator, an open-source tool designed for traffic scenarios. This simulator incorporates default settings from the Intelligent Driver Model (IDM) and MOBIL models, which handle vehicle dynamics and lane-changing behaviors. The experiments are conducted on an Ubuntu 18.04 server with an AMD 9820X processor and 64 GB of memory.

The proposed MARL algorithm is compared against a baseline version that does not include the safety supervisor introduced in this work. This comparison demonstrates the improvements in collision avoidance and overall performance brought by the safety supervisor. Additionally, a video demonstration of the training process is available online, providing a visual representation of the algorithm in action.

V. 1. Reward Function Designs

In this subsection, we evaluate the performance of the proposed Multi-Agent Reinforcement Learning (MARL) framework by analyzing two aspects: (1) the comparison between local and global reward designs, and (2) the impact of the safety penalty weight w_c in the reward function.

1. Local vs. Global Reward Designs

The proposed local reward design is compared to the global reward design used in previous works.

- **Local Reward:** Each agent receives a reward based on its own actions and outcomes.
- **Global Reward:** Each agent receives the average reward of all agents, defined as:

$$r_{i,t} = \frac{1}{N} \sum_{j=1}^N r_{j,t},$$

where N is the total number of agents, and $r_{j,t}$ is the reward of agent j at time t .

The evaluation results (shown in Fig. 8) demonstrate that the proposed local reward design outperforms the global reward design in terms of:

- **Higher Rewards:** Agents achieve better performance with local rewards.
- **Faster Convergence:** The learning process is faster with local rewards.

This difference is particularly evident in **Hard mode**, where the global reward design fails, achieving evaluation rewards below zero. The poor performance of global rewards in challenging scenarios is due to:

1. **Credit Assignment Issues:** Global rewards make it difficult to link an agent's actions to the overall outcome as the number of agents increases.
2. **Lack of Correlation:** The averaged global rewards do not align well with the actions of individual agents, reducing their effectiveness in learning.

In contrast, local rewards allow agents to optimize their own behaviors, resulting in better outcomes across all traffic scenarios, especially in complex conditions.

2. Impact of Safety Penalty Weight w_c

The effect of the safety penalty weight w_c in the reward function (Eq. 14) is studied using different values of w_c in the **Medium traffic mode**, while keeping other coefficients unchanged. The reward function is given as:

$$r_i = w_c \cdot \text{safety metric} + w_s \cdot \text{smoothness metric} + w_h \cdot \text{headway metric} + w_m \cdot \text{merging success metric}.$$

Effect of w_c :

- **High w_c :** Encourages safer behavior by penalizing unsafe actions more heavily. For example:
 - When $w_c \geq 100$, there are no collisions, indicating a strong emphasis on safety.
 - As w_c increases further, agents become overly conservative, reducing the overall traffic speed.
- **Low w_c :** Prioritizes traffic efficiency over safety, which can lead to collisions.

The testing results (shown in Table I) confirm that selecting $w_c \geq 100$ eliminates collisions. However, very high values of w_c slow down traffic flow as agents become excessively cautious.

Optimal Trade-Off: To balance safety and efficiency, w_c is set to 200 in the following experiments. This value provides a reasonable trade-off, ensuring a safe margin while maintaining traffic efficiency.

These findings emphasize the importance of reward function design in ensuring both safety and efficiency in multi-agent traffic scenarios.

V. 2. Curriculum Learning for Hard Traffic Mode

Curriculum Learning has been adopted to improve both the training speed and inference accuracy in high-traffic density scenarios, often referred to as the "Hard" traffic mode. The way people pick up difficult tasks—beginning with simpler ideas and working gradually toward more difficult ones—inspired this learning paradigm. Curriculum learning in this context refers to fine-tuning pre-trained models from simpler traffic modes—that is, Easy and Medium—for more difficult conditions, including the Hard traffic mode.

Curriculum learning lets us build on a basis laid in easier environments rather than directly from scratch for the Hard mode. This procedure greatly improves training effectiveness since the model gains from the knowledge already learned in previous phases. The model adjusts more successfully to the complexity of high-density traffic by spreading acquired representations and actions.

Safety-critical uses like autonomous driving especially benefit from this method. Direct training in a complicated environment might result in many "blind" explorations, in which the agent acts untrained and maybe dangerously. Such investigations can cause training inefficiencies or crashes. Starting from a well-trained baseline model, curriculum learning reduces the frequency of dangerous activities and increases the general safety and stability of the training process, therefore mitigating these risks.

The outcomes clearly show how successful curricula are. By greatly accelerating the convergence rate, pre-trained models from simpler traffic modes help the model to learn optimal policies more rapidly. Better final performance results from this method as well as from beginning training from scratch.

Moreover, the benefits of curriculum learning help to deal with traffic efficiently. This development emphasizes the capacity of the method to increase not only the operational effectiveness of the trained model but also the safety and stability of training.

Curriculum learning is implemented in the next tests for the Hard traffic mode considering these major advantages. This method guarantees that, in complicated, high-density traffic conditions, the model keeps improving while reducing training risks and boosting efficiency.

V. 3. Performance of the Priority-based Safety Supervisor

We assess in this subsection the performance of the suggested priority-based safety supervisor and its effects on traffic flow, training efficiency, and assessment incentives. As shown, the suggested safety supervisor exhibits better sample efficiency by faster convergence over all three traffic densities. Moreover, even in demanding Hard traffic conditions, the safety supervisor gets better assessment incentives. This development results from its capacity to substitute safe activities for dangerous ones, especially in the early stages of discovery. Early termination brought on by dangerous behavior is avoided by the safety supervisor, therefore improving the efficiency of learning.

The average vehicle speed during training shown in Figure 12 indicates traffic flow. The results unequivocally reveal that algorithms with the safety supervisor retain faster training speeds than the baseline technique, that is, without safety supervision. This underlines the two advantages of the suggested strategy: it increases traffic efficiency and training stability. Due to more frequent contacts in dense traffic, where lower speeds maintain safety by lowering accident risks, vehicle speeds decrease with increasing traffic density—26 m/s, 24 m/s, and 22 m/s for Easy, Medium, and Hard traffic modes, respectively.

The MARL algorithms for every traffic density are tested over three random seeds for thirty epochs following training. With a suitable forecast horizon ($T_n \geq 7$) the MARL framework removes collisions in all traffic modes, according to the results. By contrast, the baseline approach shows collision rates of 0.07 and 0.16 for Medium and Hard traffic categories respectively. Shorter prediction horizons $T_n = 0$ or $T_n = 3$ the safety supervisor finds difficult in demanding situations. For example, the limited horizon results in "short-sighted" decisions, which causes infeasibility after a few steps, hence the collision rate in the Hard mode stays at 0.03 even $T_n = 6$. On the other hand, too wide prediction horizons $T_n = 10, 12, 14$ could spread uncertainty and cause agents to adopt too cautious behavior to preserve safety resulting in higher accident rates and slower speeds.

The safety supervisor achieves best traffic efficiency with a suitable prediction horizon $T_n = 7, 8$. For the Easy mode, for instance, the optimum average speed is obtained with the safety supervisor ($T_n = 8$ at 27.72 m/s, vs 23.52 m/s for the baseline approach). Although wider prediction horizons raise inference time, the suggested framework is still possible for real-time implementation with enough processing capacity at the infrastructure.

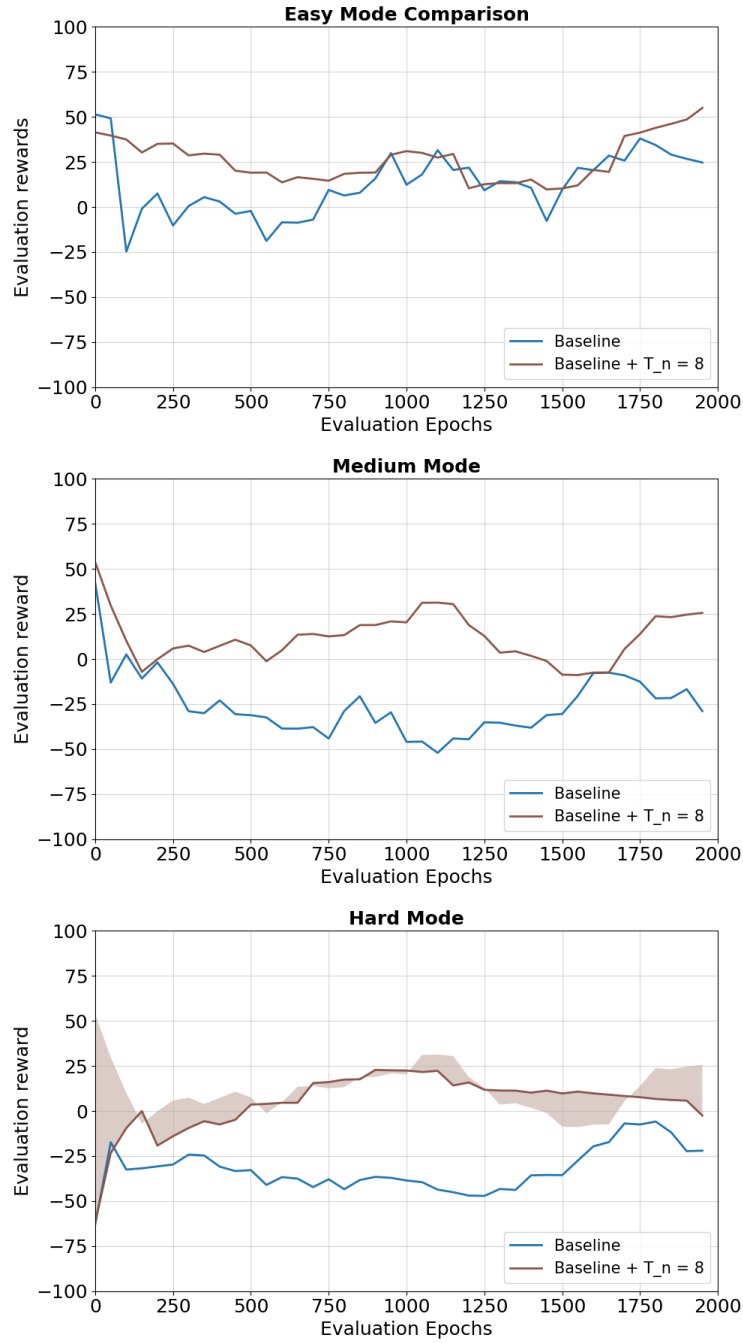


Figure 4. Training curves for the n-step priority-based safety supervisor.

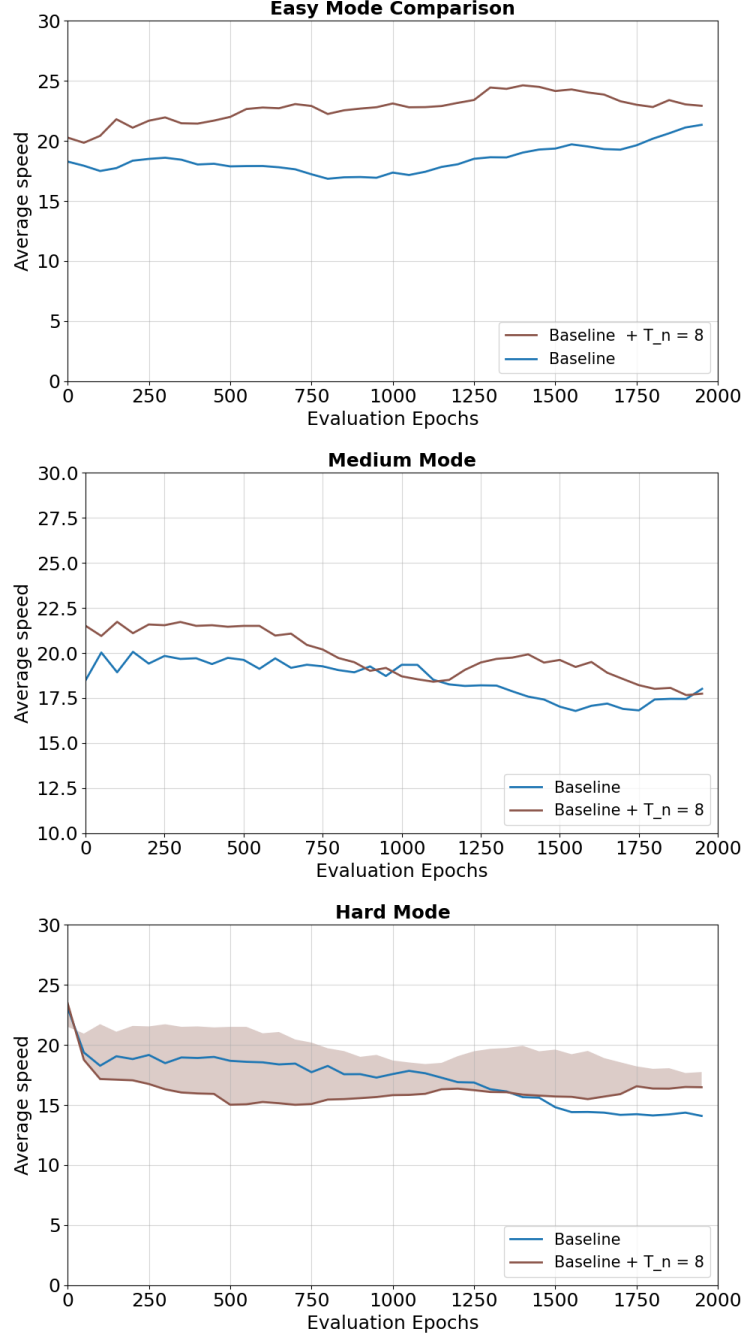


Figure 5. Average speed during training for the n -step priority-based safety supervisor

V. 4. Comparison with Benchmark Algorithms

In this subsection, we evaluate the proposed method against a state-of-the-art MARL benchmark, including MAACKTR, introduced in Section II. All MARL benchmarks rely on parameter sharing among agents to handle dynamic numbers of agents and utilize global rewards and discrete action spaces. The 8-step safety supervisor (baseline + $T_n = 8$) is applied across all traffic levels, as it offers a balanced trade-off between collision rates and prediction efficiency.

Figure 6 presents the evaluation results during training for the described MARL algorithm and MAACKTR algorithm. The proposed method (baseline + $T_n = 8$) consistently outperforms the benchmarks across all traffic levels. Notably, in the Hard traffic mode, the proposed method demonstrates significant improvements in sample efficiency and

Table 1. Testing performance comparison of collision rate and average speed (m/s) between n-step safety supervisor based on the baseline (bs) method.

Scenarios	Metrics	bs	bs + T _n = 8
Easy Mode	collision rate	0	0
	avg. speed	23.53	27.72
Medium Mode	collision rate	0.07	0
	avg. speed	20.30	24.08
Hard Mode	collision rate	0.16	0
	avg. speed	21.71	22.73

training performance over the benchmarks. Moreover, the proposed method achieves relatively higher training speeds, which enhances training efficiency.

After training, all algorithms were tested across three random seeds for 2000 epochs each, with average collision rates and vehicle speeds summarized in Table III. The results reveal that the proposed method effectively avoids collisions and achieves superior efficiency compared to other benchmark algorithms.

Furthermore, the study highlights the impact of system complexity as the number of vehicles increases. Higher traffic density exacerbates model mismatches, making it challenging to generate collision-free policies in such scenarios. These findings emphasize the robustness of model-free approaches in handling complex tasks like on-ramp merging under dense traffic conditions.

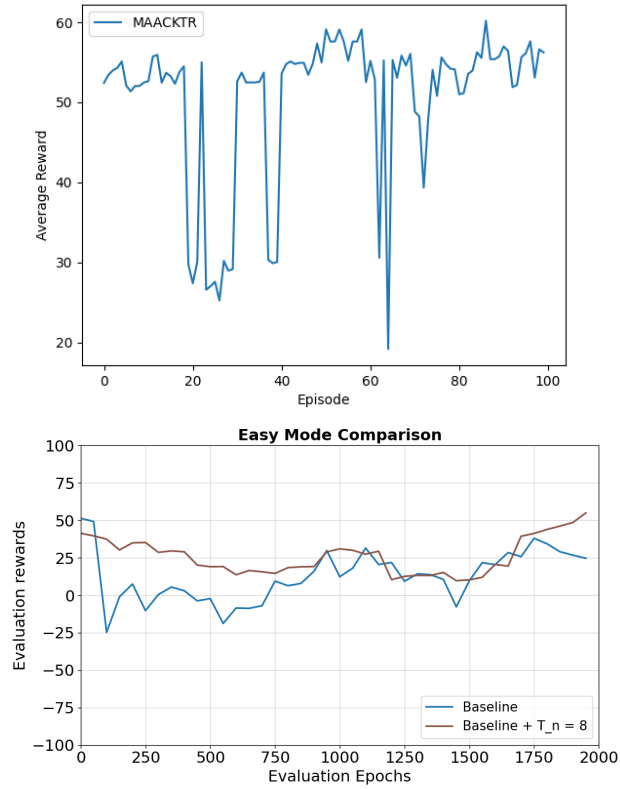


Figure 6. Training curves comparison between the proposed MARL policy (baseline (bs) + T_n = 8) and state-of-the-art MARL benchmark

VI. CONCLUSIONS AND DISCUSSIONS

From Figure 4, we observe that the baseline + $T_n = 8$ approach with 8 safety supervisors outperforms the MARL approach without the safety supervisor for all - easy, medium, and hard traffic density scenarios. This is especially seen in the high traffic density mode where the performance of the 8 safety supervisors mode is significantly better than the baseline method. This can also be verified from Figure 5, where the average speed plots vs epochs is plotted, where the safety supervisor shows its potential to achieve higher average speeds and hence ensure faster navigation and higher throughput through the traffic with minimal collision rate.

In Figure 6, the performance of a state-of-the-art algorithm- MAACKTR is compared with the designed algorithm in the paper. It can be observed from this plot, that the rewards for the MAACKTR algorithm vary abruptly with frequency oscillation. This may lead to undesirable behavior and movement of the ego vehicle as compared to the proposed algorithm which showcases reasonable stability. As the training is extended further to a higher number of epochs, the proposed algorithm is expected to return even better results.

This report replicates the work of Chen et al., implementing their deep multi-agent reinforcement learning approach for autonomous highway on-ramp merging in mixed traffic. By using MARL, curriculum learning, and a priority-based safety supervisor, the framework achieves significant improvements in safety and traffic flow, addressing the challenges of dynamic and mixed-traffic environments. Further experiments and simulations can be conducted to evaluate the robustness of the approach in different traffic conditions.

REFERENCES

- [1] “Future of driving,” <https://www.tesla.com/autopilot>, accessed: 2021-03-31.
- [2] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, “Asynchronous methods for deep reinforcement learning,” in International conference on machine learning, 2016, pp. 1928–1937.
- [3] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, “Playing atari with deep reinforcement learning,” arXiv preprint arXiv:1312.5602, 2013.
- [4] T. Chu, J. Wang, L. Codecà, 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.
- [5] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. P. Abbeel, and I. Mordatch, “Multi-agent actor-critic for mixed cooperative-competitive environments,” in Advances in neural information processing systems, 2017, pp. 6379–6390.
- [6] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, “Proximal policy optimization algorithms,” arXiv preprint arXiv:1707.06347, 2017.
- [7] Y. Wu, E. Mansimov, S. Liao, R. Grosse, and J. Ba, “Scalable trust region method for deep reinforcement learning using kronecker-factored approximation,” arXiv preprint arXiv:1708.05144, 2017.
- K. Lin, R. Zhao, Z. Xu, and J. Zhou, “Efficient large-scale fleet management via multi-agent deep reinforcement learning,” in Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery Data Mining, 2018, pp. 1774–1783.
- [9] N. Li, H. Chen, I. Kolmanovsky, and A. Girard, “An explicit decision tree approach for automated driving,” in Dynamic Systems and Control Conference, vol. 58271. American Society of Mechanical Engineers, 2017, p. V001T45A003.
- [10] T. Ayres, L. Li, D. Schleuning, and D. Young, “Preferred time-headway of highway drivers,” in ITSC 2001. 2001 IEEE Intelligent Transportation Systems. Proceedings (Cat. No. 01TH8585). IEEE, 2001, pp. 826–829.
- [11] M. Bouton, A. Nakhaei, K. Fujimura, and M. J. Kochenderfer, “Cooperation-aware reinforcement learning for merging in dense traffic,” in 2019 IEEE Intelligent Transportation Systems Conference (ITSC). IEEE, 2019, pp. 3441–3447.

APPENDIX

PYTHON CODE: MULTI-AGENT REINFORCEMENT LEARNING TRAINING AND EVALUATION

```
1 import torch as torch_lib
2 import os
3 import configparser
4 from torch import nn, optim
5 import numpy as np
6 from common.Agent import BaseAgent
7 from common.Model import PolicyNetwork, ValueNetwork, SharedActorCritic
8 from common.utils import calculate_entropy, convert_to_tensor, encode_action, VideoLogger
9
10 # Load configurations
11 config_path = 'configs/configs.ini'
12 config = configparser.ConfigParser()
13 config.read(config_path)
14 random_seed = config.getint('MODEL_CONFIG', 'torch_seed')
15 torch_lib.manual_seed(random_seed)
16 torch_lib.backends.cudnn.benchmark = False
17 torch_lib.backends.cudnn.deterministic = True
18 os.environ['PYTHONHASHSEED'] = str(random_seed)
19
20
21 class MARL(BaseAgent):
22     """
23     Multi-Agent Reinforcement Learning (MARL) using Advantage Actor-Critic (A2C)
24     """
25
26     def __init__(self, environment, state_size, action_size, **kwargs):
27         super(MARL, self).__init__(environment, state_size, action_size, **kwargs)
28
29         # Configuration Parameters
30         self.env = environment
31         self.state_dim = state_size
32         self.action_dim = action_size
33         self.rollout_steps = kwargs.get('rollout_steps', 10)
34         self.gamma = kwargs.get('discount_factor', 0.99)
35         self.entropy_weight = kwargs.get('entropy_weight', 0.01)
36         self.gradient_clipping = kwargs.get('max_gradient_norm', 0.5)
37         self.batch_limit = kwargs.get('batch_size', 64)
38         self.actor_lr = kwargs.get('actor_learning_rate', 0.0005)
39         self.critic_lr = kwargs.get('critic_learning_rate', 0.001)
40         self.shared_policy = kwargs.get('shared_policy', True)
41         self.use_gpu = torch_lib.cuda.is_available()
42
43         # Define Actor and Critic Networks
44         if self.shared_policy:
45             self.actor_critic = SharedActorCritic(self.state_dim, self.action_dim, kwargs.get('
46                 hidden_layer_size', 128))
47             self.optimizer = optim.Adam(self.actor_critic.parameters(), lr=self.actor_lr)
48         else:
49             self.policy_network = PolicyNetwork(self.state_dim, self.action_dim)
50             self.value_network = ValueNetwork(self.state_dim)
51             self.policy_optimizer = optim.Adam(self.policy_network.parameters(), lr=self.actor_lr)
52             self.value_optimizer = optim.Adam(self.value_network.parameters(), lr=self.critic_lr)
53
54         if self.use_gpu:
55             if self.shared_policy:
56                 self.actor_critic.cuda()
57             else:
58                 self.policy_network.cuda()
59                 self.value_network.cuda()
60
61         # Storage for tracking performance
```

```

61     self.training_rewards = []
62     self.episode_lengths = []
63     self.average_speeds = []
64
65     def collect_experience(self):
66         """
67         Collects experience by interacting with the environment over a series of steps.
68         """
69         states, actions, rewards, policies, masks = [], [], [], [], []
70         total_reward = 0
71         done = True
72
73         for _ in range(self.rollout_steps):
74             if done:
75                 state = self.env.reset()
76                 state_tensor = convert_to_tensor(state, self.use_gpu)
77
78                 # Get action and policy probabilities
79                 if self.shared_policy:
80                     policy_distribution, _ = self.actor_critic(state_tensor)
81                 else:
82                     policy_distribution = self.policy_network(state_tensor)
83
84                 action = torch_lib.multinomial(policy_distribution, 1).item()
85                 next_state, reward, done, _ = self.env.step(action)
86
87                 # Record the results
88                 states.append(state)
89                 actions.append(encode_action(action, self.action_dim))
90                 rewards.append(reward)
91                 policies.append(policy_distribution)
92                 masks.append(1 - int(done))
93
94                 state = next_state
95                 total_reward += reward
96
97         return states, actions, rewards, policies, masks, total_reward
98
99     def compute_loss(self, states, actions, rewards, policies, masks):
100         """
101         Computes the actor and critic losses based on collected experiences.
102         """
103         states_tensor = convert_to_tensor(states, self.use_gpu)
104         actions_tensor = convert_to_tensor(actions, self.use_gpu)
105         rewards_tensor = convert_to_tensor(rewards, self.use_gpu)
106         masks_tensor = convert_to_tensor(masks, self.use_gpu)
107
108         if self.shared_policy:
109             policy_distribution, value_estimate = self.actor_critic(states_tensor)
110         else:
111             policy_distribution = self.policy_network(states_tensor)
112             value_estimate = self.value_network(states_tensor)
113
114         log_probs = policy_distribution.log_prob(actions_tensor)
115         advantages = rewards_tensor - value_estimate
116         policy_loss = -(log_probs * advantages.detach()).mean()
117         value_loss = nn.MSELoss()(value_estimate, rewards_tensor)
118         entropy_loss = -calculate_entropy(policy_distribution) * self.entropy_weight
119
120         return policy_loss, value_loss, entropy_loss
121
122     def train_step(self, memory):
123         """
124         Trains the model for one step using sampled experiences.
125         """
126         experiences = memory.sample(self.batch_limit)
127         states, actions, rewards, policies, masks = experiences
128

```

```

129     policy_loss, value_loss, entropy_loss = self.compute_loss(states, actions, rewards,
130         policies, masks)
131     total_loss = policy_loss + value_loss + entropy_loss
132
133     self.optimizer.zero_grad()
134     total_loss.backward()
135     if self.gradient_clipping:
136         nn.utils.clip_grad_norm_(self.actor_critic.parameters(), self.gradient_clipping)
137     self.optimizer.step()
138
139 def evaluate(self, environment, episodes=10):
140     """
141     Evaluates the model over a number of episodes.
142     """
143     total_rewards = []
144     for _ in range(episodes):
145         state = environment.reset()
146         done = False
147         episode_reward = 0
148         while not done:
149             state_tensor = convert_to_tensor(state, self.use_gpu)
150             if self.shared_policy:
151                 action = self.actor_critic(state_tensor).argmax().item()
152             else:
153                 action = self.policy_network(state_tensor).argmax().item()
154             state, reward, done, _ = environment.step(action)
155             episode_reward += reward
156         total_rewards.append(episode_reward)
157     return np.mean(total_rewards)

```

Listing 1. Python Code for Proposed algorithm

SCRIPT FOR TRAINING THE PROPOSED MULTI-AGENT REINFORCEMENT LEARNING (MARL) MODEL IN A DYNAMIC ENVIRONMENT

```
1 from __future__ import print_function, division
2 from MARL import MARL
3 from common.utils import agg_double_list, copy_file, init_dir
4 from datetime import datetime
5
6 import argparse
7 import configparser
8 import sys
9 import os
10 import gym
11 import numpy as np
12 import matplotlib.pyplot as plt
13 import highway_env # Import the highway environment
14
15 # Add the highway environment path
16 sys.path.append("../highway-env")
17
18
19 def parse_args():
20     """
21     Parses command-line arguments for training or evaluation.
22     """
23     default_base_dir = "./results/"
24     default_config_dir = 'configs/configs.ini'
25
26     parser = argparse.ArgumentParser(description='Train or evaluate a policy on an RL environment using MARL.')
27
28     parser.add_argument('--base-dir', type=str, default=default_base_dir,
29                         help="Base directory to save results.")
30     parser.add_argument('--option', type=str, default='train',
31                         help="'train' to train the model, 'evaluate' to evaluate the model.")
32     parser.add_argument('--config-dir', type=str, default=default_config_dir,
33                         help="Path to the configuration file.")
34     parser.add_argument('--model-dir', type=str, default='',
35                         help="Path to a pretrained model directory.")
36     parser.add_argument('--evaluation-seeds', type=str, default=','.join([str(i) for i in range(0, 600, 20)]),
37                         help="Random seeds for evaluation, separated by commas.")
38
39     return parser.parse_args()
40
41
42 def train(args):
43     """
44     Train the MARL model using the configuration provided in the config file.
45     """
46     base_dir = args.base_dir
47     config_dir = args.config_dir
48     config = configparser.ConfigParser()
49     config.read(config_dir)
50
51     # Create an experiment folder
52     timestamp = datetime.now().strftime("%b-%d_%H_%M_%S")
53     output_dir = os.path.join(base_dir, timestamp)
54     dirs = init_dir(output_dir)
55     copy_file(dirs['configs'])
56
57     model_dir = args.model_dir if os.path.exists(args.model_dir) else dirs['models']
58
59     # Load model and training configuration
60     BATCH_SIZE = config.getint('MODEL_CONFIG', 'BATCH_SIZE')
61     MEMORY_CAPACITY = config.getint('MODEL_CONFIG', 'MEMORY_CAPACITY')
62     ROLL_OUT_N_STEPS = config.getint('MODEL_CONFIG', 'ROLL_OUT_N_STEPS')
```

```

63 reward_gamma = config.getfloat('MODEL_CONFIG', 'reward_gamma')
64 training_strategy = config.get('MODEL_CONFIG', 'training_strategy')
65 actor_hidden_size = config.getint('MODEL_CONFIG', 'actor_hidden_size')
66 critic_hidden_size = config.getint('MODEL_CONFIG', 'critic_hidden_size')
67 MAX_GRAD_NORM = config.getfloat('MODEL_CONFIG', 'MAX_GRAD_NORM')
68 ENTROPY_REG = config.getfloat('MODEL_CONFIG', 'ENTROPY_REG')
69 epsilon = config.getfloat('MODEL_CONFIG', 'epsilon')
70 alpha = config.getfloat('MODEL_CONFIG', 'alpha')
71 shared_network = config.getboolean('MODEL_CONFIG', 'shared_network')
72 reward_type = config.get('MODEL_CONFIG', 'reward_type')
73
74 # Training configurations
75 actor_lr = config.getfloat('TRAIN_CONFIG', 'actor_lr')
76 critic_lr = config.getfloat('TRAIN_CONFIG', 'critic_lr')
77 MAX_EPISODES = config.getint('TRAIN_CONFIG', 'MAX_EPISODES')
78 EPISODES_BEFORE_TRAIN = config.getint('TRAIN_CONFIG', 'EPISODES_BEFORE_TRAIN')
79 EVAL_INTERVAL = config.getint('TRAIN_CONFIG', 'EVAL_INTERVAL')
80 EVAL_EPISODES = config.getint('TRAIN_CONFIG', 'EVAL_EPISODES')
81 reward_scale = config.getfloat('TRAIN_CONFIG', 'reward_scale')
82
83 # Initialize environment
84 env = gym.make('merge-multi-agent-v0')
85 env_eval = gym.make('merge-multi-agent-v0')
86
87 for param in ['seed', 'simulation_frequency', 'duration', 'policy_frequency', '
88             COLLISION_REWARD',
89             'HIGH_SPEED_REWARD', 'HEADWAY_COST', 'HEADWAY_TIME', 'MERGING_LANE_COST', '
90             traffic_density',
91             'safety_guarantee', 'n_step']:
92     env.config[param] = config.get('ENV_CONFIG', param, fallback=None)
93
94 env_eval.config.update(env.config)
95
96 # Model initialization
97 state_dim = env.n_s
98 action_dim = env.n_a
99 marl = MARL(env, state_dim, action_dim,
100             memory_capacity=MEMORY_CAPACITY, roll_out_n_steps=ROLL_OUT_N_STEPS,
101             reward_gamma=reward_gamma, reward_scale=reward_scale,
102             actor_hidden_size=actor_hidden_size, critic_hidden_size=critic_hidden_size,
103             actor_lr=actor_lr, critic_lr=critic_lr,
104             entropy_reg=ENTROPY_REG, batch_size=BATCH_SIZE,
105             training_strategy=training_strategy, shared_network=shared_network,
106             reward_type=reward_type, max_grad_norm=MAX_GRAD_NORM)
107
108 # Load existing model if available
109 marl.load(model_dir, train_mode=True)
110
111 episodes, eval_rewards = [], []
112 best_eval_reward = -float('inf')
113
114 while marl.n_episodes < MAX_EPISODES:
115     marl.explore()
116     if marl.n_episodes >= EPISODES_BEFORE_TRAIN:
117         marl.train()
118
119 # Evaluate periodically
120 if marl.episode_done and (marl.n_episodes % EVAL_INTERVAL == 0):
121     rewards, *_ = marl.evaluation(env_eval, dirs['train_videos'], EVAL_EPISODES)
122     avg_reward = np.mean(rewards)
123     print(f"Episode {marl.n_episodes}, Avg. Reward: {avg_reward:.2f}")
124     episodes.append(marl.n_episodes)
125     eval_rewards.append(avg_reward)
126
127 if avg_reward > best_eval_reward:
128     marl.save(dirs['models'], marl.n_episodes)
129     best_eval_reward = avg_reward

```

```

129     np.save(os.path.join(output_dir, 'eval_rewards.npy'), np.array(eval_rewards))
130
131     # Save training progress and final model
132     marl.save(dirs['models'], MAX_EPISODES)
133     plt.figure()
134     plt.plot(eval_rewards)
135     plt.xlabel("Episodes")
136     plt.ylabel("Evaluation Rewards")
137     plt.savefig(os.path.join(output_dir, "training_rewards.png"))
138     plt.show()
139
140
141 def evaluate(args):
142     """
143     Evaluate a pretrained MARL model.
144     """
145     if not os.path.exists(args.model_dir):
146         raise FileNotFoundError("Pretrained model not found.")
147
148     config = configparser.ConfigParser()
149     config.read(os.path.join(args.model_dir, 'configs.ini'))
150
151     env = gym.make('merge-multi-agent-v0')
152     for param in ['seed', 'simulation_frequency', 'duration', 'policy_frequency', '
153                 COLLISION_REWARD',
154                 'HIGH_SPEED_REWARD', 'HEADWAY_COST', 'HEADWAY_TIME', 'MERGING_LANE_COST', '
155                 traffic_density',
156                 'safety_guarantee', 'n_step']:
157         env.config[param] = config.get('ENV_CONFIG', param, fallback=None)
158
159     marl = MARL(env, state_dim=env.n_s, action_dim=env.n_a)
160     marl.load(args.model_dir)
161
162     rewards, *_ = marl.evaluation(env, args.model_dir, eval_episodes=config.getint('TRAIN_CONFIG'
163                                         , 'EVAL_EPISODES'))
164     print(f"Average Rewards: {np.mean(rewards):.2f}")
165
166
167 if __name__ == "__main__":
168     args = parse_args()
169     if args.option == 'train':
170         train(args)
171     elif args.option == 'evaluate':
172         evaluate(args)
173     else:
174         raise ValueError("Invalid option. Use 'train' or 'evaluate'.")

```

Listing 2. Script for training the proposed Multi-Agent Reinforcement Learning (MARL) model in a dynamic environment

CONFIGURATION FILE

The following is the configuration file used in the project:

```
1 [MODEL_CONFIG]
2 epsilon = 1e-5;
3 alpha = 0.99;
4 reward_gamma = 0.99
5 MAX_GRAD_NORM = 5
6 ROLL_OUT_N_STEPS = 100
7 MEMORY_CAPACITY = 100
8 ; only use the latest ROLL_OUT_N_STEPS for training A2C
9 BATCH_SIZE = 100
10 ENTROPY_REG = 0.01
11 ; seeds for pytorch, 0, 2000, 2021
12 torch_seed = 2021
13
14 ; concurrent
15 training_strategy = concurrent
16 actor_hidden_size = 128
17 critic_hidden_size = 128
18 shared_network = True
19 action_masking = True
20 state_split = True
21 ; "greedy", "regionalR", "global_R"
22 reward_type = regionalR
23
24 [TRAIN_CONFIG]
25 MAX_EPISODES = 2000
26 EPISODES_BEFORE_TRAIN = 1
27 EVAL_EPISODES = 3
28 EVAL_INTERVAL = 50
29 reward_scale = 20.
30 actor_lr = 5e-4
31 critic_lr = 5e-4
32 test_seeds =
33     0,25,50,75,100,125,150,175,200,325,350,375,400,425,450,475,500,525,550,575
34
35 [ENV_CONFIG]
36 ; seed for the environment, 0, 2000, 2021
37 seed = 2021
38 simulation_frequency = 15
39 duration = 20
40 policy_frequency = 5
41 COLLISION_REWARD = 200
42 HIGH_SPEED_REWARD = 1
43 HEADWAY_COST = 4
44 HEADWAY_TIME = 1.2
45 MERGING_LANE_COST = 4
46 traffic_density = 3
47 safety_guarantee = True
48 n_step = 8
```

Listing 3. Configuration File: configs.ini

CUSTOM HIGHWAY MERGING ENVIRONMENT

```
1 import numpy as np
2 from gym.envs.registration import register
3 from typing import Tuple
4
5 from highway_env import utils
6 from highway_env.envs.common.abstract import AbstractEnv, MultiAgentWrapper
7 from highway_env.road.lane import LineType, StraightLane, SineLane
8 from highway_env.road.road import Road, RoadNetwork
9 from highway_env.vehicle.controller import ControlledVehicle, MDPVehicle
10 from highway_env.road.objects import Obstacle
11 from highway_env.vehicle.kinematics import Vehicle
12
13
14 class MergeEnvironment(AbstractEnv):
15     """
16     A simulation environment for highway merging scenarios.
17
18     The focus vehicle navigates a highway segment that includes a merge point, where vehicles
19     enter from an on-ramp. Rewards are based on maintaining high speeds, avoiding collisions,
20     and facilitating smooth merging interactions.
21     """
22     n_a = 5
23     n_s = 25
24
25     @classmethod
26     def default_config(cls) -> dict:
27         config = super().default_config()
28         config.update({
29             "observation": {
30                 "type": "Kinematics"},
31             "action": {
32                 "type": "DiscreteMetaAction",
33                 "longitudinal": True,
34                 "lateral": True},
35             "controlled_vehicles": 1,
36             "screen_width": 600,
37             "screen_height": 120,
38             "centering_position": [0.3, 0.5],
39             "scaling": 3,
40             "simulation_frequency": 15, # [Hz]
41             "duration": 20, # total duration
42             "policy_frequency": 5, # [Hz]
43             "reward_speed_range": [10, 30],
44             "COLLISION_REWARD": -200, # Penalty for collisions
45             "HIGH_SPEED_REWARD": 1, # Reward for maintaining high speed
46             "HEADWAY_COST": 4, # Cost for insufficient spacing
47             "HEADWAY_TIME": 1.2, # [s]
48             "MERGING_LANE_COST": 4, # Cost for staying in merging lane
49             "traffic_density": 1, # Levels of traffic density (easy/hard)
50         })
51         return config
52
53     def _reward(self, action: list) -> float:
54         """Calculate the overall reward for all controlled vehicles."""
55         return sum(self._vehicle_reward(action, vehicle) for vehicle in self.controlled_vehicles) \
56             / len(self.controlled_vehicles)
57
58     def _vehicle_reward(self, action: int, vehicle: Vehicle) -> float:
59         """
60         Compute the reward for an individual vehicle.
61         :param action: The action taken by the vehicle.
62         :param vehicle: The vehicle for which the reward is calculated.
63         :return: A reward value for the given state-action pair.
64         """
65         speed_ratio = utils.lmap(vehicle.speed, self.config["reward_speed_range"], [0, 1])
```

```

66
67     # Penalize for being in the merging lane
68     if vehicle.lane_index == ("b", "c", 1):
69         merging_penalty = -np.exp(-(vehicle.position[0] - sum(self.ends[:3])) ** 2 / (10 *
70             self.ends[2]))
71     else:
72         merging_penalty = 0
73
74     # Penalize for inadequate headway
75     headway_dist = self._compute_headway_distance(vehicle)
76     headway_penalty = np.log(headway_dist / (self.config["HEADWAY_TIME"] * vehicle.speed)) if
77         vehicle.speed > 0 else 0
78
79     # Compute total reward
80     reward = (
81         self.config["COLLISION_REWARD"] * (-1 * vehicle.crashed)
82         + (self.config["HIGH_SPEED_REWARD"] * np.clip(speed_ratio, 0, 1))
83         + self.config["MERGING_LANE_COST"] * merging_penalty
84         + self.config["HEADWAY_COST"] * (headway_penalty if headway_penalty < 0 else 0)
85     )
86     return reward
87
88 def _regional_reward(self):
89     """Calculate rewards for vehicles within a regional context."""
90     for vehicle in self.controlled_vehicles:
91         neighbor_vehicles = []
92
93         # Determine vehicles in the vicinity
94         if vehicle.lane_index in [("a", "b", 0), ("b", "c", 0), ("c", "d", 0)]:
95             v_front, v_rear = self.road.surrounding_vehicles(vehicle)
96             if self.road.network.side_lanes(vehicle.lane_index):
97                 v_side_front, v_side_rear = self.road.surrounding_vehicles(vehicle, self.road
98                     .network.side_lanes(vehicle.lane_index)[0])
99             else:
100                 v_side_front, v_side_rear = None, None
101         else:
102             v_side_front, v_side_rear = self.road.surrounding_vehicles(vehicle)
103             v_front, v_rear = None, None
104
105         # Collect all neighbor vehicles
106         for v in [v_front, v_side_front, vehicle, v_rear, v_side_rear]:
107             if isinstance(v, MDPVehicle) and v:
108                 neighbor_vehicles.append(v)
109
110         regional_reward = sum(v.local_reward for v in neighbor_vehicles)
111         vehicle.regional_reward = regional_reward / len(neighbor_vehicles)
112
113 def step(self, action: int) -> Tuple[np.ndarray, float, bool, dict]:
114     """Perform an environment step."""
115     obs, reward, done, info = super().step(action)
116     obs = np.array(obs).reshape((len(obs), -1))
117     info["agents_rewards"] = tuple(v.local_reward for v in self.controlled_vehicles)
118     return obs, reward, done, info
119
120 def _reset(self, num_CAV=0) -> None:
121     """Reset the environment for a new episode."""
122     self._make_road()
123     self._populate_road(num_CAV)
124     self.action_is_safe = True
125     self.T = int(self.config["duration"] * self.config["policy_frequency"])
126
127 def _make_road(self) -> None:
128     """Construct the road network."""
129     net = RoadNetwork()
130
131     c, s = LineType.CONTINUOUS_LINE, LineType.STRIPED
132     net.add_lane("a", "b", StraightLane([0, 0], [sum(self.ends[:2]), 0], line_types=[c, c]))

```

```

130     net.add_lane("b", "c", StraightLane([sum(self.ends[:2]), 0], [sum(self.ends[:3]), 0],
131         line_types=[c, s]))
132     net.add_lane("c", "d", StraightLane([sum(self.ends[:3]), 0], [sum(self.ends), 0],
133         line_types=[c, c]))
134
135     road = Road(network=net, np_random=self.np_random, record_history=self.config["
136         show_trajectories"])
137     self.road = road
138
139     def _populate_road(self, num_CAV: int) -> None:
140         """Populate the road with vehicles."""
141         road = self.road
142         self.controlled_vehicles = []
143
144         # Spawn vehicles
145         for _ in range(num_CAV):
146             vehicle = self.action_type.vehicle_class(road, road.network.get_lane(("a", "b", 0)).
147                 position(10, 0), speed=25)
148             self.controlled_vehicles.append(vehicle)
149             road.vehicles.append(vehicle)
150
151 register(
152     id='merge-v1',
153     entry_point='highway_env.envs:MergeEnvironment',
154 )
155
156 register(
157     id='merge-multi-agent-v0',
158     entry_point='highway_env.envs:MergeEnvironment',
159 )

```

Listing 4. Script for training the proposed Multi-Agent Reinforcement Learning (MARL) model in a dynamic environment