

Multi-Agent System Techniques for Autonomous Vehicle Navigation and Traffic Management*

Coordination, Control and Decision-Making in Intelligent Transportation Systems

Diana Brebeanu[†]

Cheriton School of Computer Science
University of Waterloo
Waterloo, Ontario, Canada
dbrebean@uwaterloo.ca

ABSTRACT

As autonomous vehicles continue to evolve and integrate into real-world transportation systems, the demand for intelligent, adaptive, and fail-safe decision-making processes becomes increasingly important. This survey paper reviews recent multi-agent system techniques developed to enhance autonomous navigation and traffic coordination, with a particular focus on two challenges: merging into a lane and ensuring safe decision-making under uncertainty. The first section explores the task of lane merging, a scenario that requires both strategic coordination among vehicles and precise motion control. Models based on game theory, specifically two-player static games, are used to determine optimal merging sequences, simulating how vehicles negotiate right-of-way based on mixed-strategy Nash equilibria. These strategic decisions are then coupled with bi-objective motion planning that leverages Pontryagin's Minimum Principle to minimize fuel consumption and travel time. A novel solution approach to traverse the Pareto front with a modified quicksort algorithm is introduced, allowing for efficient identification of optimal trade-offs between competing objectives. The second focus area centers on ensuring safe decision-making and addresses the challenges of operating under uncertain and adversarial conditions. Reinforcement learning is used to add robustness guarantees through the use of an adversarial training framework. The Robust Reinforcement Learning with Safety Guarantees (RRL-SG) model incorporates a learned adversary to simulate worst-case environmental perturbations and also adds a safety mask based on the Responsibility-Sensitive Safety (RSS) framework to eliminate high-risk actions. Furthermore, a Safe-Robust Markov Decision Process (MDP) formulation is introduced to iteratively refine policy and value functions using an actor-critic framework and robust policy improvement strategies. This dual analysis reveals important trade-offs and complementarities between coordination-centric and safety-centric methodologies. While merging strategies emphasize collaborative optimization in structured environments, robust RL approaches are designed to withstand uncertainty and maintain safety in unpredictable conditions. Together, these methodologies highlight the complex nature of autonomous navigation and demonstrate how integrating optimization, game theory, and machine learning can lead to more resilient and cooperative autonomous vehicle systems.

CCS CONCEPTS

- Computing methodologies
- Theory of Computation
- Information Systems
- Applied Computing

KEYWORDS

Autonomous vehicles, multi-agent systems, game theory, motion planning, reinforcement learning, decision-making under uncertainty, safe navigation, intelligent transportation systems, on-ramp merging, robust control

ACM Reference format:

Diana Brebeanu, 2025. Multi-Agent System Techniques for Autonomous Vehicle Navigation and Traffic Management: Coordination, Control and Decision-Making in Intelligent Transportation Systems. In *Proceedings of Waterloo ECE 493 (WATERLOO '25)*. ECE 493, Waterloo, ON, Canada, 2 pages. <https://doi.org/10.1145/1234567890>

1 Introduction

As autonomous vehicles move from concept to reality, the intersection of autonomous driving, traffic control and connected vehicle networks presents important opportunities for researchers to tackle the challenges of ensuring safe traffic management and control of vehicle responses to changes in their environment. These have a fundamental impact on public safety through the co-operation of different vehicles to reduce potential accidents as well as on transportation efficiency through optimizing traffic flow to reduce congestion and carbon emissions. Achieving these results requires robust solutions to challenges such as vehicle communication, real-time decision making for joint benefits as well as fail-safe systems that prioritize safety. In this survey paper, we will focus on techniques used for merging into a lane and ensuring safe decision-making processes centered around changing lanes.

1.1 Merging into a lane

In 2021, the paper “On-Ramp Merging Strategy for Connected and Automated Vehicles Based on Complete Information Static Game” by Min et al. [6] modelled the problem of merging into a lane by separating it into two concerns: merging

sequence and motion planning. Merging sequence refers to scheduling the order in which cars leave the merging ramp and join the main road while motion planning is concerned with ensuring vehicles move in a safe way to ensure avoidance of collisions while also optimizing fuel consumption and traffic efficiency. The work done in , "Cooperative Game Approach to Optimal Merging Sequence and on-Ramp Merging Control of Connected and Automated Vehicles" by Jing et al. [3] is also covered and focuses on modelling the merging sequence problem as a two-player complete information static game in which players simultaneously select their strategies. This is a competitive game in which the vehicles are competing for higher priority in the merging sequence. Each of the vehicles selects their strategy by calculating their mixed Nash strategy equilibrium.

For the motion planning part, it is formulated as a bi-objective optimization problem based on Pontryagin's Minimum Principle. In this approach, the parameters of the cost function are determined using a search method called varying scale-grid, where the search scale progressively narrows down to refine the solution. This method of parameter search is not commonly found in existing literature. The solution to this optimization problem is then identified by employing a quicksort algorithm to traverse the Pareto front, a set of solutions that are non-dominated and outperform other solutions in the entire solution space. This method of using the Pareto front for solution selection and the quicksort algorithm for optimization is also a novel contribution not typically seen in previous studies.

This two-pronged approach of using game theory for merging sequence and optimal control for motion planning demonstrates a well-integrated framework that mirrors the layered complexity of real-world merging behavior in connected and automated vehicle (CAV) systems. By decoupling the sequence decision from the physical control of vehicle movement, the researchers effectively reduced the computational complexity of the problem while preserving the interactions between strategic and dynamic elements. Notably, the mixed Nash equilibrium formulation captures the uncertainty and strategic interdependence between vehicles, aligning well with the decentralized nature of CAVs. Meanwhile, the bi-objective motion planning approach balances competing goals of safety and efficiency, and the use of Pareto front traversal with a quicksort-based search introduces a novel, computationally efficient method for identifying optimal trade-offs. Together, these methods reflect an advancement in both the modeling precision and the solution techniques for autonomous on-ramp merging, providing a foundation for future work to build more responsive and cooperative merging systems.

1.2 Ensuring Safe Decision-Making Processes

As previous sections focused on calculating decisions for vehicles in different traffic situations, it is also critical to ensure that these decisions consider the safety of the passengers with the highest importance. Reviewing the paper by He et al. [2], a new reinforcement learning technique is introduced to ensure collision safety through the development of an adversary model to simulate worst-case uncertainties. This is combined with an actor-critic algorithm to enable the agent to learn policies against adversarial perturbations to make trustworthy decisions and reduce collision likelihood. The framework for this is named RRL-SG (Robust Reinforcement Learning – Safety Guarantees).

Building on this foundation, the adversary model introduced by He et al. serves not just as a stress test but as a core training component to enhance the robustness of the decision-making process. By purposefully introducing perturbations that simulate rare or extreme driving conditions—such as sudden pedestrian crossings or unpredictable vehicle maneuvers—the model forces the reinforcement learning agent to account for edge-case scenarios that traditional training might overlook. This adversarial training ensures that the agent's learned policy is not only effective under nominal conditions but also resilient when exposed to unexpected hazards. The actor-critic framework further supports this by continuously adjusting both the policy and value estimations in response to these adversarial challenges, ultimately improving the system's ability to generalize and maintain safety across a wide range of traffic scenarios.

2 Related Works

2.1 Merging into a lane

In 2017, Kang and Rahka [4] proposed a decision-making process for calculating a merging sequence based on non-cooperative game theory. In this paper, they discuss two possible states that the vehicles could be in: original state and lagging state when compared to the car in the other lane. The game played between the two cars determines which car gets to keep its original state and which car will lag the other one. A similar approach is used in the Min et al. paper [6] to determine the order of vehicles.

In 2019, Jing et al. [3] modelled merging sequence as a multi-player game which was then decomposed into multiple two-player games. This technique was then used in the paper by Min et al [6] and combined with motion planning.

For the motion planning problem, two primary approaches have been explored: centralized and decentralized systems. Centralized systems involve a central controller that manages the traffic flow. For instance, Cao et al. (2015) [1] proposed a model where vehicles on the main road receive information from a central controller about vehicles on the ramp, allowing them to adjust their state to optimize the merging process. In contrast, decentralized systems rely on each vehicle making decisions based on the information it receives. Wang et al. (2018) [8] explored this approach in the context of on-ramp merging, where vehicles cooperate using Vehicle-to-Vehicle (V2V) and Vehicle-to-Infrastructure (V2I) communication to adapt their speeds and positions, thereby improving traffic flow during merging scenarios.

2.2 Ensuring Safe Decision-Making Processes

A traditional and most popular technique used in decision-making is based on rule-based systems such as finite-state machines (FSM). This is generally simple to implement but relies on the expert knowledge of specialists which makes it difficult to design rules for vehicles to follow during more complicated traffic situations.

An alternative to this approach is reinforcement learning. As reinforcement learning is based around an agent interacting with the environment to learn optimal behaviour (policies) for similar future interactions, it is a useful tool for solving complex sequential decision-making problems that autonomous vehicles

are tasked with. Ye et al. [9] published a paper outlining a reinforcement learning technique for determining policies for changing lanes and that these can be generalized to new environments. Moreover, reinforcement learning techniques have been used for determining optimal target speeds or speed patterns. However, most studies only consider at most one point of uncertainty which lowers the robustness of many of these approaches and their safety guarantees. The paper by He et al. [2] aims to consider policy robustness against multiple uncertainties.

3 Methodology

3.1 Merging into a lane

The problem of computing how cars on a merging ramp will merge into the main road is split into two parts. The first is the merging sequence problem, which is essentially a scheduling problem. This is concerned with determining the order of the cars while merging. This means comparing a car in the main road with a car on the merging ramp, determining which car will go in front of the other after having merged. The second part is motion planning which deals with the actual motion of vehicles to avoid collisions, optimize fuel consumption and enhance traffic flow.

3.1.1 Merging Sequence. The merging environment is modelled as a ramp that merges into a main road. This area is split into two sections: the control area before the merging ramp intersects with the main road (called the merging point) and the merging area where the ramp and main road unify. The control area is where the vehicles are meant to adjust their speeds according to a central controller before merging. In modelling the problem, several simplifying assumptions are made including that there are no time delays for sending signals between the centralized controller and the vehicles, no overtaking is allowed, and we are only looking at situations where there is one lane on the main road and one lane on the merging ramp. No overtaking means that no cars will speed up to pass a vehicle, they will only decide whether to slow down to lag. This will also help ensure safety.

The merging sequence is modelled as a two-player game. When a vehicle V_m in the merging lane enters the merging section, a game begins between itself and the first vehicle in the main road V_i which will determine who will be in the lead after merging. After this, another game will start between that the vehicle that lost the game (is now behind the leader) and the next vehicle in the other ramp.

Player 1	Player 2	
	Leader(λ)	Follower($1-\lambda$)
Leader(θ)	$(-c, -c)$	(m_1, m_2)
Follower($1-\theta$)	(n_1, n_2)	$(-c, -c)$

Figure 1: Payoff matrix for a two-player merging sequence game from [6]

The payoff matrix is shown in Figure 1 where c represents the cost of a collision between vehicles and is defined as described in (1)

$$c = \frac{1}{\| \chi_1(t_g) - \chi_2(t_g) \|} \quad (1)$$

The χ function describes the state of a player. The cost c is then greater if the states of each player are similar. m_1, m_2, n_1, n_2 are the payoffs of each player when they co-operate and are related to fuel consumption and travel time which will be explained in the motion planning section.

Solving for the mixed-strategy Nash equilibrium of the payoff matrix gives the following best response probability distributions in (2) for each player resulting in the mixed Nash equilibrium as $[(\theta^*, 1 - \theta^*), (\lambda^*, 1 - \lambda^*)]$

$$\theta^* = \frac{n_2 + c}{2c + m_2 + n_2} \text{ and } \lambda^* = \frac{m_1 + c}{2c + m_1 + n_1} \quad (2)$$

Because the equilibrium is a pure strategy, not a mixed one, if the player's states are similar such that $\chi_1(t_g) \rightarrow \chi_2(t_g)$ which causes $c \rightarrow \infty$, then $\theta^* \rightarrow \frac{1}{2}$ and $\lambda^* \rightarrow \frac{1}{2}$ which means both players could choose to be the follower or leader in the merging sequence. To avoid collisions, a convention is set that Player 1 will become the leader and Player 2 will become the follower.

3.1.2 Motion Planning. Building on top of the work by Jing et al [3], the cost function for each vehicle V_i is defined as

$$J_i = \frac{1}{2} \int_{t_i^0}^{t_i^f} w_1 a_i^2(t) dt + w_2 T_i \quad (3)$$

Where $a_i(t)$ represents the acceleration of vehicle V_i at moment t and T_i represents the travel time ($T_i = t_i^f - t_i^0$). The first term in the cost function relates to fuel consumption while the second term relates to travel time with the coefficients w_1 and w_2 representing their weights.

Motion planning is then formulated as a global optimization problem to find the minimum cost by minimizing over possible accelerations with added limits on minimum and maximum possible values of accelerations and velocity.

$$\min_{a_i} J_i \quad (4)$$

$$a_{\min} \leq a_i(t) \leq a_{\max}, 0 \leq v_i(t) \leq v_{\max}, \forall t \in \{t_i^0, t_i^f\}$$

Pontyagin's principle is then used to calculate the relationship between a vehicle's position, velocity and acceleration used during merging from this optimization problem dependent on the cost function. The payoffs from the matrix in Figure 1 are then derived from these values.

The two parameters w_1 and w_2 for the cost function are determined through running a varying-scale grid search to find an unbiased Pareto solution. This means a solution which is nondominated, any deviation would deteriorate at least one of the objectives: fuel consumption and travel time.

The solution is found placing initial limits l_1 and l_2 on w_1 and w_2 respectively and initializing search step lengths σ_1 and σ_2 along with scaling factors α_1 and α_2 . Then the fuel consumption and travel time in (5) and (6) respectively are calculated based on the values of w_1 and w_2 (substituting in W in the equation). Note that $f_i(t)$ calculates the fuel consumption based on the previously calculated velocity and acceleration of the car.

$$F = \sum_i^W \sum_j^N f_i(j) \quad (5)$$

$$T = \sum_i^W T_i \quad (6)$$

The Pareto front is then found using modified quicksort and the new values for F and T are then set to w_1 and w_2 . These are then used to calculate the new limits l_1 and l_2 and new search step lengths σ_1 and σ_2 . This process is then repeated until there are no longer significant changes in the step values.

3.2 Ensuring Safe Decision-Making Processes

To train an agent to anticipate multiple uncertainties, He et al. [2] train an adversarial agent online to model worst-case uncertainties through generating optimal perturbations of the observed states that a vehicle is in and the characteristics of the environment it interacts with. Furthermore, a safety mask is added to ensure collision safety which transforms the probability of selecting an unsafe decision to zero. This is then used with an adversarial robust actor algorithm which allows the agent to learn robust policies against these perturbations.

Pertaining to the agent's own characteristics, the state of an agent is designed to have 15 dimensions which include its own velocity, acceleration and lane index as well as the relative distance and velocity of the six nearest vehicles in the same or adjacent lanes. The action space is discrete and contains the following possible actions: changing lanes to the left, changing lanes to the right, maintaining the current state, accelerating at a fixed rate of 1.47 m/s^2 and decelerating at -2.00 m/s^2 .

3.1.1 Adversary Model. The optimal adversarial perturbations observed states, and environmental dynamics are represented as Δ_o^* and Δ_d^* in the form of probability distributions. The adversarial model takes in the state s of the agent and outputs these two perturbations. Δ_o^* is meant to add noise to the agent's observations such that it causes the worst-case confusion in the agent, expressed as maximizing the average variation distance on perturbed policies. Δ_d^* is then meant to minimize the expected return of the agent through causing the environment to behave in a way that causes worst-case rewards for the agent.

To quantify policy variation under adversarial observation noise, the model employs the Jensen–Shannon (JS) divergence, a symmetric and bounded variant of the Kullback–Leibler (KL) divergence. The JS divergence measures the difference between the original policy and the perturbed policy, encapsulated in an objective function J_o , which accounts for the change in action distributions due to observation perturbations.

To perturb the environmental dynamics, the agent's expected return is estimated by the Q function $Q_s^\pi(s)$ which takes in the state of the agent and assumes the action follows the policy π . This is then captured by the objective function

$$J_d(s, Q^\pi, \Delta_d) = \Delta_d Q^\pi(s) \quad (1)$$

These two objective functions are then combined and weighted to give (2)

$$J_\Delta(s, \pi, Q^\pi, \Delta) = (\alpha - 1)J_o(s, \pi, \Delta_o) + \alpha J_d(s, Q^\pi, \Delta_d) \quad (2)$$

This is then used in the following optimization problem to get the optimal parameters of the adversary model. The optimization has been simplified using the hyperbolic tangent and softmax functions on the perturbations

$$\theta \in \argmin_\theta E[J_\Delta(s, \pi, Q^\pi; \theta)] \quad (3)$$

3.1.2 Safety Mask. To ensure the collision safety of autonomous vehicles, a safety mask based on the Responsibility-Sensitive Safety (RSS) framework is developed using a jerk-bounded model. This model, derived from Intel, accounts for a realistic braking profile where a vehicle decelerates with a bounded jerk until a minimum deceleration is reached, followed by constant deceleration until a full stop. The minimum longitudinal safe distance D_{min}^{RSS} is computed accordingly considering vehicle speeds, accelerations, and braking dynamics.

This approach is then extended to lateral maneuvers such as lane changes, where a minimum lateral safety distance is introduced as a scaled version of D_{min}^{RSS} using a safety coefficient. The safety mask modifies decision-making probabilities by assigning negative infinite rewards to actions that violate these longitudinal or lateral safety distances. For instance, if the distance to a vehicle in the target lane is insufficient, lane-change and acceleration actions are suppressed. This masking technique combines reinforcement learning with rule-based filtering of unsafe maneuvers, enhancing decision safety without extensive retraining.

3.1.3 Adversarial Robust Actor-Critic Algorithm. A Markov Decision Process (MDP) is used as the basis for finding the optimal policy of the agent. In this case, the standard MDP is extended to model the behaviour of the agent under the adversarial perturbations and safety mask constraints. A new Safe-Robust MDP is defined as a max-min problem in (4) where T is the last time step and $\beta > 0$ is a trade-off coefficient:

$$\max_\pi \min_\Delta E[\sum_{t=0}^T \gamma^t r(s_t, a_t) + \beta J_\Delta(s, \pi, Q^\pi, \Delta)] \quad (4)$$

The reward function r was designed to represent driving safety, passenger comfort and travel efficiency and as opposed to Min et al. [6], fuel consumption was not taken into consideration. The reward function rewarded the agent for driving at high speeds while also penalizing the agent if any maneuvers caused a collision or performed risky movements such as high-speed lane changes.

Finding the optimal policy is done in two steps, safe-robust policy evaluation and robust policy improvement which are iterated until convergence.

3.1.3.1 Safe-Robust Policy Evaluation. The policy evaluation phase estimates the expected return of a fixed policy under environmental uncertainty. To account for possible adversarial perturbations, the standard Bellman backup operator $\Psi^{\pi, \Delta}$ is used to capture the effect of such perturbations on future value estimates.

$$\Psi^{\pi, \Delta} Q^\pi(s_t) = r_a(s_t, a_t) + \gamma \pi(s_{t+1}) Q^\pi(s_{t+1}) \quad (5)$$

This includes in an augmented reward that reflects both the reward and the influence of uncertainty and action-value function $Q^\pi(\cdot)$ which estimates the expected return based on the state and action of the agent when it follows policy π .

$$r_a = r(s_t, a_t) + \gamma \beta J_\Delta(\cdot) \quad (6)$$

To perform policy evaluation effectively, the authors employ two separate critic networks, each representing a parameterized action-value function. These networks are trained using a loss function that measures the discrepancy between predicted and target value estimates. The target is computed conservatively by taking the minimum value across both critic networks, which helps mitigate overestimation bias commonly observed in value-

based reinforcement learning. The critic networks are then updated by minimizing the squared Bellman error using stochastic gradient descent. To stabilize training further, the target networks are maintained using Polyak averaging, which softly updates target parameters as a weighted average of current and previous values. This approach enhances the safety and robustness of policy evaluation by explicitly modeling uncertainty and leveraging conservative value estimation strategies during training.

3.1.3.2 Robust Policy Improvement. This step is concerned with optimizing the policy π given the action-value function $Q^\pi(\cdot)$ under the adversarial perturbations. This is represented as a max-min game

$$\max_{\pi} \min_{\Delta} E[J(\pi, \Delta)] \quad (7)$$

Where $J(\cdot)$ corresponds to the objective function from (4) such that $J(\pi, \Delta) = \pi(s)Q^\pi(s) + \beta J_{\Delta}(s, \pi, Q^\pi, \Delta)$ and recalling that J_{Δ} comes from (2).

The optimal policy π^* and optimal adversarial perturbation Δ^* are then found by

1. Fixing an arbitrary policy π
2. Solving for Δ^* according to the optimization problem in (8)
3. Learning π^* with Δ^* using according to (9)

$$\Delta^* = \arg \min_{\Delta} E[J(\pi, \Delta)] \quad (8)$$

$$\pi^* = \arg \max_{\pi} E[J(\pi, \Delta^*)] \quad (9)$$

We can see that this represents a zero-sum game and based on the derived results, a convergence of the policy improvement is guaranteed.

4 Experimental Results

4.1 Merging into a lane

To explore the impact of varying parameters, the simulation focused on a control area with 10 vehicles, a control zone that was 40 meters long, and a 30-meter merging area. Vehicle speeds were randomly initialized with a normal distribution (mean of 15 m/s), and all vehicles were controlled to merge at a constant speed of 13.4 m/s. Constraints for acceleration and velocity were set within limits of $[-3, 3]$ m/s² and $[0, 30]$ m/s, respectively.

The initial step size for the search was then set at 0.05, with a scaling factor of 0.2. The search for the coefficients, w_1 and w_2 , was conducted across a grid range of $[0, 1]$, with surrogate objectives of fuel consumption and total travel time being equally weighted. After identifying potential Pareto-optimal solutions, the best solution was selected based on minimal bias, yielding coefficient values of $w_1 = 0.30$ and $w_2 = 0.80$. The search was refined by narrowing the range to $[0.25, 0.35]$ for w_1 and $[0.75, 0.85]$ for w_2 , with a finer precision of 0.01. The final optimal coefficient values were $w_1 = 0.32$ and $w_2 = 0.79$, which were selected after the search precision requirement was met.

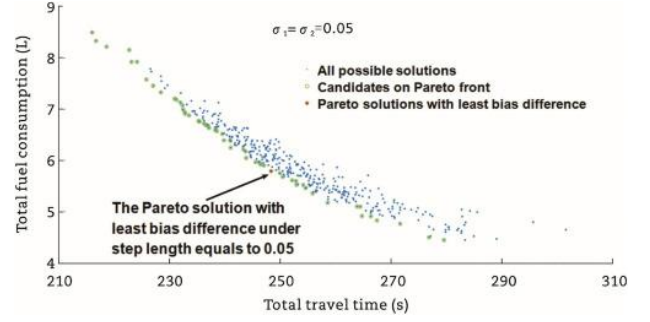


Figure 2: Simulation results from grid search under a 0.05 step length [2]

The experiments showed that all constraints were able to be satisfied, and the vehicles were all able to cooperatively merge. Moreover, the centralized controller ensured that vehicles passed the merging area in a coordinated manner, maintaining a constant distance of 30 meters between vehicles, ensuring safety. The vehicles also adhered to a consistent time headway, calculated as the length of the merging area divided by the merging speed (M/v_m), ensuring smooth and safe merging.

4.2 Ensuring Safe Decision-Making Processes

To evaluate the effectiveness and robustness of the proposed RRL-SG framework for safe decision-making, comparisons were conducted against several state-of-the-art RL baselines. These included dueling double deep Q-networks (D3QN) as a representative Q-learning algorithm, proximal policy optimization (PPO) for on-policy learning, soft actor-critic (SAC) for off-policy learning, and observation adversarial reinforcement learning (OARL) for robust RL.

Performance was assessed using a variety of metrics tailored to reflect both efficiency and safety in autonomous driving. The expected return was used as the primary indicator of overall policy performance. Additionally, average running speed and number of collisions were recorded to evaluate travel efficiency and traffic safety, respectively. To assess robustness, the extent of policy change under adversarial perturbations was measured; smaller changes indicated stronger resistance to attacks. For the on-ramp merging task, an additional metric—the merging success rate—was used, defined as the rate at which a vehicle successfully merged into the main lane without collisions.

Experiments were carried out in the Simulation of Urban MObility (SUMO) environment. Agents were trained and tested across two scenarios: a highway setting and an on-ramp merging scenario. In the highway scenario, traffic density was varied by altering the vehicle spawn probability P across low (0.06), normal (0.12), and high (0.24) levels. Each agent was trained under normal traffic density and tested under all three densities. For each algorithm, five independent training runs were performed using different random seeds. Evaluations were averaged over 100 test episodes, with ten episodes sampled per evaluation to account for stochastic variability in traffic flow.

To evaluate robustness, trained agents were tested under adversarial observational attacks generated by a learned adversary model. During these tests, the input states to the agent were

perturbed to simulate the impact of adversarial environmental noise, thus enabling a realistic assessment of policy stability and resilience. The highway scenario was used both for training and testing, while the on-ramp merging scenario was exclusively used for testing. The focus was on assessing the generalization of trained agents to this unseen task, specifically evaluating merging behavior in terms of safety and success under realistic traffic dynamics.

Quantitatively, in normal-density traffic without adversarial interference, RRL-SG achieved return improvements of approximately 22.31%, 7.22%, 10.34%, and 1.97% over D3QN, PPO, SAC, and OARL, respectively. These gains became even more pronounced in high-density environments, where returns improved by up to 78.63%, 47.41%, 25.45%, and 13.84%, respectively. Under adversarial conditions in high-density traffic, RRL-SG delivered dramatic return gains of 7669.57%, 2666.25%, 511.57%, and 8.99% over the same baselines.

Overall, the proposed framework demonstrates a significant improvements over state-of-the-art baselines across various metrics related to safety, robustness, and expected return. Compared to D3QN, PPO, SAC, and OARL, the agent consistently achieved superior performance in both adversarial and non-adversarial settings. These improvements are attributed to the integration of an RSS-based safety mask, which restricts the action space to safe regions, thereby mitigating the risk of collisions and enhancing learning efficiency by reducing unnecessary exploration.

5 Conclusion

As autonomous vehicles continue to integrate into existing traffic systems, the challenge of safely managing their interactions with traditional vehicles becomes increasingly critical. This survey paper highlights key techniques in addressing specific traffic scenarios, particularly lane merging and ensuring safe decision-making processes. By utilizing concepts such as game theory, optimization, and advanced communication systems, researchers have made significant strides in designing systems that not only improve safety but also enhance overall traffic efficiency. The methodologies discussed, such as the two-player game models for merging and the optimization of fuel consumption and travel time, provide valuable insights into how self-driving cars can interact with each other and their environment. In particular, the merging strategies illustrate a dual-layered approach: strategic planning through Nash equilibria to determine merging order, and dynamic motion planning using Pontryagin's Minimum Principle to minimize travel time and fuel consumption while avoiding collisions. The potential benefits of these techniques extend beyond just collision avoidance to include smoother traffic flow, reduced congestion, and a positive impact on environmental sustainability.

In contrast, the techniques for ensuring safe decision-making focus on robustness under uncertainty and the handling of adversarial conditions. The introduction of an adversarial reinforcement learning framework (RRL-SG) extends traditional policy learning by explicitly accounting for worst-case scenarios using adversary-generated perturbations. Additionally, safety is proactively embedded through the use of a safety mask based on the Responsibility-Sensitive Safety (RSS) framework, which guarantees collision avoidance by eliminating unsafe decisions from the agent's action space. This approach differs from merging

models by prioritizing policy robustness and system reliability in unpredictable environments over traffic flow optimization.

Together, these techniques underscore the complex nature of autonomous vehicle navigation. While merging solutions excel at structured interaction and efficiency under known conditions, robust decision-making frameworks ensure system resilience and safety in complex, uncertain environments. The potential benefits of both domains extend beyond collision avoidance to include smoother traffic flow, reduced congestion, and increased trustworthiness.

However, while these solutions show promise, they also highlight the complexity of achieving fully coordinated autonomous driving systems. The integration of real-time decision-making, vehicle-to-vehicle and vehicle-to-infrastructure communication, and fail-safe mechanisms remains an ongoing challenge. Future research will need to refine these techniques and explore new ones to ensure that autonomous vehicles can seamlessly interact within human-driven traffic systems, ultimately paving the way for safer, more efficient roads.

ACKNOWLEDGMENTS

I would like to thank Professor Zahedi for his valuable instruction on the fundamentals of game theory and multi-agent systems, as well as for his thoughtful responses to questions on these topics which aided in writing this paper.

REFERENCES

- [1] W. Cao, M. Mukai, T. Kawabe, H. Nishira, and N. Fujiki, "Cooperative vehicle path generation during merging using model predictive control with real-time optimization," *Control Engineering Practice*, vol. 34, pp. 98–105, Jan. 2015, doi: 10.1016/j.conengprac.2014.10.005.
- [2] X. He, W. Huang, C. Lv, "Toward Trustworthy Decision-Making for Autonomous Vehicles: A Robust Reinforcement Learning Approach with Safety Guarantees," *Engineering*, vol. 33, pp. 77–89, Feb. 2024, doi: 10.1016/j.eng.2023.10.005.
- [3] S. Jing, F. Hui, X. Zhao, J. Rios-Torres, and A. J. Khattak, "Cooperative Game Approach to Optimal Merging Sequence and on-Ramp Merging Control of Connected and Automated Vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 20, no. 11, pp. 4234–4244, Nov. 2019, doi: 10.1109/TITS.2019.2925871.
- [4] K. Kang and H. A. Rakha, "Game Theoretical Approach to Model Decision Making for Merging Maneuvers at Freeway On-Ramps," *Transportation Research Record*, vol. 2623, no. 1, pp. 19–28, 2017, doi: 10.3141/2623-03.
- [5] B. Liu, W. Han, E. Wang, S. Xiong, L. Wu, Q. Wang, "Multi-Agent Attention Double Actor-Critic Framework for Intelligent Traffic Light Control in Urban Scenarios with Hybrid Traffic," *IEEE Transactions on Mobile Computing*, vol. 23, no. 1, pp. 660–672, Jan. 2024, doi: 10.1109/tmc.2022.3233879.
- [6] H. Min, Y. Fang, X. Wu, G. Wu, X. Zhao, "On-Ramp Merging Strategy for Connected and Automated Vehicles Based on Complete Information Static Game," *Journal of Traffic and Transportation Engineering (English Edition)*, vol. 8, no. 4, pp. 582–595, Aug. 2021, doi: 10.1016/j.jtte.2021.07.003.
- [7] C. Spatharis, and K. Blekas, "Multiagent Reinforcement Learning for Autonomous Driving in Traffic Zones with Unsignalized Intersections," *Journal of Intelligent Transportation Systems*, vol. 28, no. 1, pp. 103–119, Aug. 14, 2022, doi: 10.1080/15472450.2022.2109416.
- [8] Z. Wang, G. Wu, and M. Barth, "Distributed Consensus-Based Cooperative Highway On-Ramp Merging Using V2X Communications," *SAE Technical Paper* 2018-01-1177, 2018.
- [9] F. Ye, P. Wang, C. -Y. Chan and J. Zhang, "Meta Reinforcement Learning-Based Lane Change Strategy for Autonomous Vehicles," 2021 IEEE Intelligent Vehicles Symposium (IV), Nagoya, Japan, 2021, pp. 223–230, doi: 10.1109/IV48863.2021.9575379.