

# A Survey of Reinforcement Learning-Based Motion Planning for Autonomous Driving: Lessons Learned from a Driving Task Perspective

Zhuoren Li, Guizhe Jin, Ran Yu, Zhiwen Chen, Wei Han, Nan Li, Lu Xiong, Bo Leng, Jia Hu, Ilya Kolmanovsky and Dimitar Filev

**Abstract**—Reinforcement learning (RL), with its ability to explore and optimize policies in complex, dynamic decision-making tasks, has emerged as a promising approach to addressing motion planning (MoP) challenges in autonomous driving (AD). Despite rapid advancements in RL and AD, a systematic description and interpretation of the RL design process tailored to diverse driving tasks remains underdeveloped. This survey provides a comprehensive review of RL-based MoP for AD, focusing on lessons from task-specific perspectives. We first outline the fundamentals of RL methodologies, and then survey their applications in MoP, analyzing scenario-specific features and task requirements to shed light on their influence on RL design choices. Building on this analysis, we summarize key design experiences, extract insights from various driving task applications, and provide guidance for future implementations. Additionally, we examine the frontier challenges in RL-based MoP, review recent efforts to address these challenges, and propose strategies for overcoming unresolved issues.

**Index Terms**—Reinforcement learning, autonomous driving, motion planning, survey.

## I. INTRODUCTION

**R**EINFORCEMENT learning (RL) is a machine learning paradigm that focuses on solving sequential decision-making and control challenges [1]. In contrast to supervised learning such as imitation learning (IL) [2]), where the agent directly learns a policy with labels of expert data, an RL agent generates its policy by interacting with the environment, and evaluating and iterating itself by statistically maximizing long-term rewards with its trial-and-error property [3]. The RL agent still learns a mapping between inputs and outputs rather than hidden patterns within the data. With RL methods surpassing human world champions in Go [4], Starcraft II [5], automobile racing [6], and drone racing [7], RL has been recognized as a promising approach for AD, especially for motion planning (MoP) [8], [9].

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible. (*corresponding author: Bo Leng, Jia Hu.*)

Zhuoren Li, Guizhe Jin, Ran Yu, Zhiwen Chen, Nan Li, Wei Han, Lu Xiong and Bo Leng are with the School of Automotive Studies, Tongji University, Shanghai 201804, China. (email: 1911055@tongji.edu.cn)

Jia Hu is with the Key Laboratory of Road and Traffic Engineering of the Ministry of Education, Tongji University, Shanghai 201804, China.

Ilya Kolmanovsky is with the Department of Aerospace Engineering, University of Michigan, Ann Arbor, MI 48109, USA.

Dimitar Filev is with the Hagler Institute for Advanced Study, Texas A&M University, College Station, TX 77840 USA.

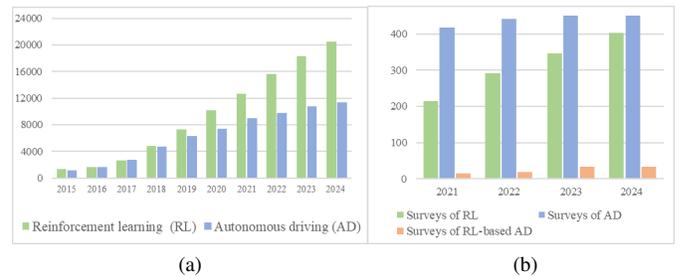


Fig. 1. Search result of Web of Science until 2024: (a) topic search for RL and AD. (b) topic search for surveys for RL, AD, and RL-based AD.

According to the search results from Web of Science (WOS), the number of research papers on the RL and AD topics has surged over the past decade, as shown in Fig. 1. In particular, owing to the complexity of interaction with the environment in different MoP problems [10], [11], RL has proven highly applicable to these tasks [12]. Recently, research on RL technologies applied to MoP has explored in a variety of driving tasks [13].

Most existing surveys focus on the overall technology of AD, or focus on specific functions such as localization, perception (especially object detection), communication, etc., with very few studies on MoP [14]. There are even fewer surveys summarizing RL-based MoP studies. Several references such as [3], [15] have reviewed some studies and applications of RL-based MoP for AD. Nevertheless, most of them focus on the perspective of categorizing RL methodologies, and do not clearly define the connection between RL and the specific driving tasks. Some surveys such as [16], [17] have tried to categorize and discuss RL-based MoP research according to driving scenarios, and provide insight into some state-of-the-art RL research from a problem-driven perspective. However, their summarization is incomplete, ignoring some rare driving tasks, such as parking and racing. Most importantly, they do not provide a detailed introduction to the scenario characteristics and task requirements corresponding to AD tasks, as well as their impact on RL model design.

Moreover, the limitations and challenges identified by most surveys, such as driving safety, policy robustness, sample efficiency, and scenario generalization, have been further explored in recent years. Despite the existence of several summaries of

advanced theoretical approaches to RL [18], [19], [20], to the best of our knowledge, there is no review that comprehensively summarizes the application of these state-of-the-art technologies to the field of MoP for AD. With the rapid development of RL-based AD technologies in both academia and industry, holistic and thorough review of recent investigations is needed.

This article analyzes and summarizes recent advanced work from a comprehensive driving task perspective (although owing to space limitations, we are unable to include some impressive RL-based MoP papers in this article). Our study aims to systematically answer the following questions: *How can RL be employed to formulate a MoP model for specific AD tasks? What are the generic design paradigms and customized adaptations of RL for various driving tasks? What are the advances addressing the current challenges for RL-based MoP?* The contributions of this article include the following:

- We outline the fundamentals of RL methodologies, and then focus on their applications in MoP for AD, where various driving tasks are systematically characterized to shed light on their influence on RL design.
- We summarize several developments in RL-based MoP for AD, extract insights from various driving task applications, and provide guidance for future implementations.
- The current challenges in RL applications to MoP for AD are discussed, and beyond pointing out challenges and future directions, a comprehensive review of recent exploratory efforts to address these issues with advanced methods is undertaken.

The structure of this article is shown in Fig. 2, and the remainder of it is organized as follows: Section II briefly introduces the basics of RL and RL-based MoP. Section III reviews research on RL-based MoP from a driving task perspective. Section IV discusses the lessons learned from RL-based MoP design for various driving tasks, and offers experiences and insights. Section V analyzes the current challenges in RL-based MoP and details exploratory efforts to apply advanced RL theories to address them, exploring outlooks and opportunities. Section VI concludes this article.

## II. BASICS OF RL AND RL-BASED MOP FOR AD

### A. Basic Theory and Algorithm of Reinforcement Learning

Perception, action, and goal are the three key elements of RL: After perceiving information about the environment state, the RL agent can take actions to influence the environment to achieve its goal. In RL, the agent is not concerned with how to act based on expert data, rather, iterates the policy by evaluating action performance through reward signals and improves its policy to achieve its goal. In general, the RL model can be formulated as a Markov Decision Process (MDP) [21] satisfying the Markov property: The future states depend only on the current state. Specifically, an MDP problem can be defined by a tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \gamma \rangle$ :

- $\mathcal{S}$  and  $\mathcal{A}$  denote the state and action spaces, respectively, i.e.  $s_t \in \mathcal{S}$  and  $a_t \in \mathcal{A}$ .
- $\mathcal{T} : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ ,  $\mathcal{T}(s_{t+1}, s_t, a_t)$  is the transition function from a current state-action pair  $(s_t, a_t)$  to a new state  $s_{t+1}$  at the next time step with probability

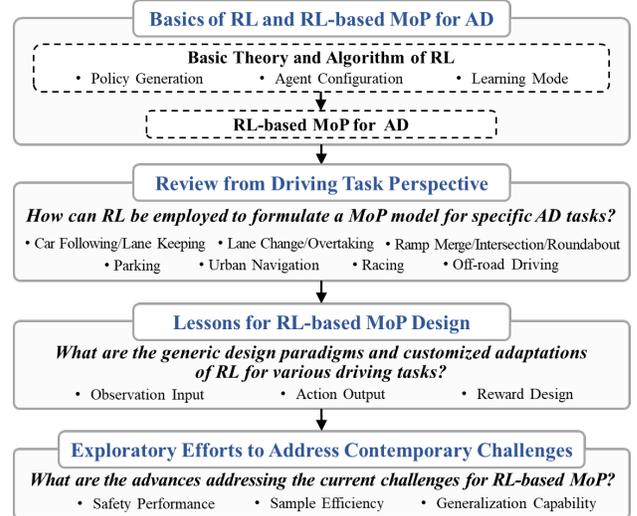


Fig. 2. The schematic of the survey structure of RL-based MoP for AD.

$P(s_{t+1} | s_t, a_t)$ , which is referred to as the environmental dynamics (system dynamics).

- $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  is the reward function used to evaluate the agent's performance.
- $\gamma \in [0, 1]$  denotes the discount factor for the present value of the future reward.

To describe not fully observable states, the MDP problem can be extended to a partially observable MDP (POMDP) [22]. For POMDPs, an observation space  $\mathcal{O}$ , an observation function  $\Omega(a_t, s_{t+1}, o_{t+1}) : \mathcal{S} \rightarrow \mathcal{O}$ , and the probability  $P(o_{t+1} | a_t, s_{t+1})$  of observing  $o_{t+1}$  after the agent executed at and reached  $s_{t+1}$ .

The policy  $\pi : (a_t | s_t)$ , maps the observed state  $s_t$  to a probability of an action  $a_t$ , which represents the driving maneuver in the AD driving task. The set of all possible policies is expressed by  $\Pi$ . The sequence  $\{s_0, a_0, s_1, a_1, \dots, s_t, a_t, \dots\}$  generated by the RL agent with the policy  $\pi$  is called trajectory or rollout. The solution objective of the MDP is to find the optimal policy  $\pi^*$  resulting in the highest expected discounted return over all possible trajectories, where  $h$  is the current timestep and  $H$  is the finite horizon (for an infinite horizon  $H$  is set to  $\infty$ ). Furthermore, the expectation of return following the policy  $\pi$  from a state  $s$  is defined as the value-function:

$$V_\pi(s) = \mathbb{E}[G_t | s_t = s] = \mathbb{E}\left[\sum_{t=h}^{h+H} \gamma^{t-h} \mathcal{R}_{t+1} | s_h = s\right] \quad (1)$$

where  $G_t$  means the total return for the current state  $s_t$ . Similarly, the action-value function, i.e., ‘‘Q-value function’’ is defined as:

$$\begin{aligned} Q_\pi(s_t, a_t) &= \mathbb{E}[G_t | s_t = s, a_t = a] \\ &= \mathbb{E}[\mathcal{R}_t + \gamma Q_\pi(s_{t+1}, \pi(a_{t+1}) | s_{t+1})] \end{aligned} \quad (2)$$

According to whether the state transition probability  $\mathcal{T}$  is known, RL methods can be classified into model-based and model-free. Typical model-based RL methods can utilize dynamic programming (DP) [23] to find the optimal policy with known environment dynamics. However, since the state transition function in many engineering applications (e.g.,

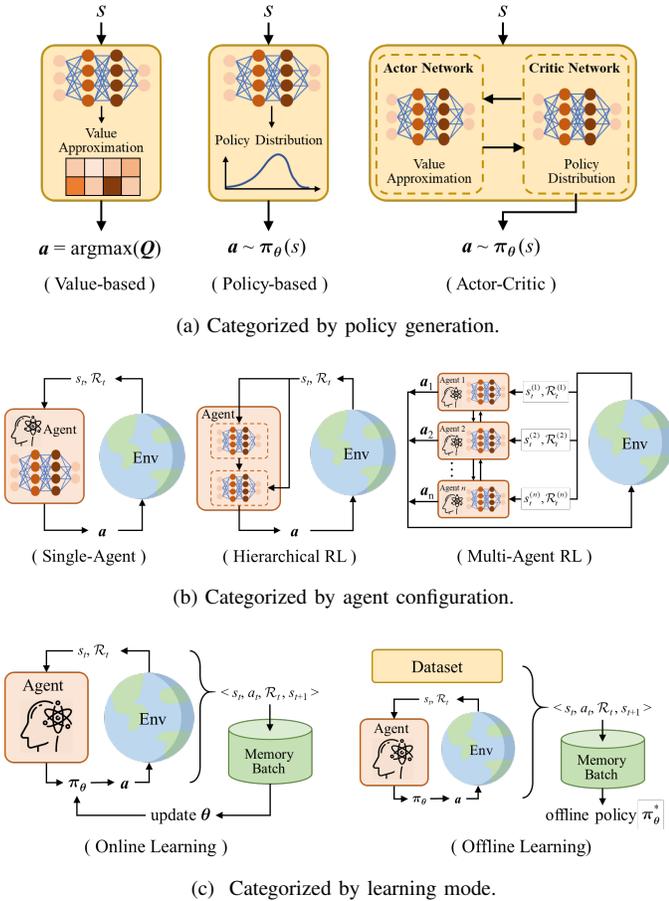


Fig. 3. RL methods with different categorization.

MoP for AD) is often unclear, it is a challenge to model the interaction between the agent and environment, limiting the application of model-based RL methods [24]. In contrast, model-free RL implicitly construct environment dynamics during learning, which are typically solved by Monte Carlo methods [25] and temporal difference (TD) methods [26].

When the states of the environment and the agent are high-dimensional or even infinite, it is impractical to store all Q-values. One widely used method is to use deep neural networks (DNNs) as a nonlinear Q-value function approximator over high-dimensional state spaces. Subsequently,  $\pi_\theta$  is denoted as the policy parameterized by the network parameter  $\theta$ , which aims to fit arbitrarily complex policy distribution functions.

In this article, RL algorithms are categorized by the difference in policy generation, agent configuration, and learning mode, as shown in Fig. 3, and we focus on model-free methods that are more applicable to MoP for AD.

### 1) Policy generation

a) *Value-based Methods*: These methods explicitly identify an optimal value function and learn the optimal policy from the value function. Q-learning is one of the most classic RL models. The optimal policy  $\pi^*$  of Q-learning aims to maximize the Q-value and can be defined as:

$$\arg \max_{\pi} Q_{\pi}(s_t, a_t) = \arg \max_{\pi} \mathbb{E} \left[ \sum_{t=h}^{h+H} \gamma^t \mathcal{R}(s_t, \pi(a_t|s_t)) \right] \quad (3)$$

The RL agent can update their policies by estimating Q-value as follows:

$$Q_{\pi}(s_t, a_t) \leftarrow Q_{\pi}(s_t, a_t) + \alpha \left[ \mathcal{R}_t + \gamma \max_{a_{t+1} \in \mathcal{A}} Q_{\pi}(s_{t+1}, a_{t+1}) - Q_{\pi}(s_t, a_t) \right] \quad (4)$$

where  $\alpha$  is the learning rate. With the use of DNNs, Q-learning has evolved into far-reaching algorithms represented by Deep Q-Network (DQN) [27], Double DQN (DDQN) [28], dueling DQN [29], and Dueling Double DQN (D3QN) [30]. In practice, the outstanding aspects of DQN are the experience replay and the design of the target network. The former breaks the correlation between experience samples and improves data utilization efficiency. The latter introduces a target network with parameters updated periodically during Q-network updates, thereby alleviating instability from rapid fluctuations in the Q-network. The loss function of the Q-network in the DQN can be expressed as:

$$l_t^Q(\theta) = \frac{1}{2} \left[ \mathcal{R}_t + \gamma \max_{a_{t+1} \in \mathcal{A}} Q'(s_{t+1}, a_{t+1}; \theta') - Q(s_t, a_t; \theta) \right]^2 \quad (5)$$

where  $Q'(\cdot; \theta')$  is the target network. Furthermore, the DDQN implements action selection and value evaluation with different Q networks, which reduces overestimation bias. The dueling DQN models the value function separately from the advantage function to improve the stability of the strategy. D3QN combines the techniques underlying the three algorithms above to obtain a more advanced value-based approach.

b) *Policy-based Methods*: Unlike value-based methods that indirectly obtain the policy by the optimal value function, policy-based methods directly iterate the parameters of the differentiable policy function. Such policy-based methods are more suitable for continuous control problems with infinite action sets. Specifically, the objective function for directly optimizing a stochastic policy function  $\pi_\theta$  is:

$$J(\theta) = \mathbb{E}_{\pi_\theta} [G_t | s_t = s]. \quad (6)$$

Policy gradient methods [25] use gradient descent to estimate the policy parameters that maximize the expected reward:

$$\nabla J(\theta) = \mathbb{E} \left[ \sum_{t=h}^{h+H} \gamma^t \mathcal{R}_t(s_t, \pi_\theta(a_t|s_t)) \nabla \log \pi_\theta(a_t|s_t) \right] \quad (7)$$

$$\theta \leftarrow \theta + \nabla J(\theta)$$

The value function still needs to be computed via the policy-based approach to update the policy. The REINFORCE algorithm [31] uses the Monte Carlo method to estimate  $Q_{\pi}(s_t, a_t)$ , but the estimation results exhibit large variance. In addition, the advantage function  $A_{\pi}(s_t, a_t) = Q_{\pi}(s_t, a_t) - V_{\pi}(s_t, a_t)$  [32] can be utilized to replace  $Q_{\pi}(s_t, a_t)$  to emphasize better actions. Note that policy-based methods can also use a deterministic policy (determining an action based on the state  $s$ , i.e.  $a = \mu_\theta(s)$ , which can be more efficient), rather than just a stochastic policy (selecting an action from a probability distribution,  $a \sim \pi_\theta(\cdot|s)$ ). In this case, the gradient of the objective function can be expressed as:

$$\nabla J(\mu_\theta) = \mathbb{E}_{\mu_\theta} [\nabla_{\theta} \mu_\theta(s) \nabla_a Q_{\mu_\theta}(s_{t+1}, a_t | a_t = \mu_\theta)] \quad (8)$$

The deterministic policy focuses only on exploitation during training and not on exploration. Therefore, the Deterministic Policy Gradient algorithm (DPG) [33] utilizes an off-policy approach to optimize the deterministic policy by sampling from the stochastic policy to ensure sufficient exploration.

*c) Actor-Critic Methods:* Actor-Critic methods are a special type of policy-based method that integrates techniques from value-based methods, where the actor is the policy function  $\pi_\theta$  generating actions to obtain the maximum return, and the critic is the value function  $V_{\pi_\theta}$  that estimates the actions. This coupled structure integrates the flexibility of policy optimization and the stability of value estimation. The Deep Deterministic Policy Gradient (DDPG) [34], Proximal Policy Optimization (PPO) [35], and Soft Actor-Critic (SAC) [36] algorithms are typical algorithms that utilize the actor-critic framework. In particular, the SAC algorithm maximizes the entropy of the actions while maximizing the expected return, thus encouraging exploration to obtain better performance. This has made it a popular paradigm in recent years [37].

## 2) Agent Configuration

In a single agent configuration, all interactions with the environment occur through a single agent. Specially, hierarchical RL (HRL) leverages hierarchical abstraction techniques [38] to decomposes an agent into multiple components, simplifying complex tasks by breaking them into subtasks learned by subagents. Not all subagents interact with the environment; typically, the actions from high-level subagents are concatenated into the state space of low-level subagents to provide context and guidance [8] while low-level subagents can control the entire agent. HRL is grounded in Semi MDP (SMDP), which includes the option selection policy  $\pi_{\mathcal{O}}(o_t|s_t)$  and the option internal policy  $\pi_o(a_t|s_t)$ . The high-level agent selects an option  $o_t$ , and then the low-level agent executes the policy  $\pi_o(a_t|s_t)$  corresponding to  $o_t$ , continuing until the option is interrupted [39]. Depending on whether the policies of the high-level and low-level agents are trained synchronously, HRL can be categorized into synchronous and asynchronous architectures. Methods with synchronous architectures are usually composed of a high-level policy providing coarse-grained subgoals, and a low-level policy to achieve fine-grained control [40]. Asynchronous HRL pre-trains multiple low-level policies for different tasks and trains the high-level policy to invoke them appropriately [41].

Multi-Agent RL (MARL) enables multiple agents to independently interact with the shared environment. Each agent has its own task, but its observations and rewards are influenced by the joint actions of all agents. Meanwhile, a single agent's long-term optimization objective also impacts the policy learning of other agents. Given the differences in observations among agents, the interaction process between agents and the environment is typically described by Markov Game (MG) [42], which is defined by an extension tuple  $\langle \mathcal{S}, N, \mathcal{A}^{(i)}_{i=1\sim N}, \mathcal{R}^{(i)}_{i=1\sim N}, \mathcal{T}, \gamma, \Omega, \mathcal{O}^{(i)}_{i=1\sim N} \rangle$ , where  $\mathcal{A}^{(i)}_{i=1\sim N}$  is the action sets for N agents,  $\mathcal{R}^{(i)}_{i=1\sim N}$  is the reward set and  $\mathcal{T} : \mathcal{S} \times \mathcal{A}^{(1)} \times \dots \times \mathcal{A}^{(i)} \times \dots \times \mathcal{A}^{(N)} \rightarrow [0, 1]$  is the transition function. Each agent receives a local observation  $\mathcal{O}^{(i)}$  by  $\Omega(\mathcal{S}, i)$ .

Relationships between agents can be categorized as coop-

erative, competitive, or mixed [43]. Additionally, the training process and action execution in MARL systems can generally be classified into two paradigms: centralized and decentralized. However, centralized execution requires real-time communication and a shared policy among all agents, which is difficult to implement in real-world systems [44]. Therefore, researchers often use two main architectures: i) Centralized Training Decentralized Execution (CTDE): During training, a central critic controls the global perspective and updates the policies of all agents based on their states and actions. ii) Decentralized Training Decentralized Execution (DTDE): Each agent is trained and operated independently, without the need to access global information. However, as the number of agents increases, the state space grows exponentially, making it challenging and slow to train a MARL system [45].

## 3) Learning Mode

Online RL allows the agent to freely interact with the environment and thus collect experience. The RL agent is required to collect sample data (trial-and-error experience) by itself in the training environment and relies on these data to update the policy. This allows the RL agent to discover an unknown optimal policy. However, online RL usually suffers from sample inefficiency in some tasks and places high demands on the fidelity of the training environment.

Offline RL is a framework dedicated to policy optimization from static, previously collected datasets, and it capitalizes on historical interaction data to derive optimal policy. In contrast to online RL, Offline RL relies solely on a pre-established dataset  $\mathcal{D}$ , thereby eliminating the need for ongoing exploration while mitigating associated risks. The core objective of offline RL is to minimize the Bellman error:

$$\begin{aligned} \nabla J(\theta) = & \mathbb{E}_{s_t, a_t, s_{t+1} \sim \mathcal{D}} [\mathcal{R}_t + \gamma \mathbb{E}_{a_{t+1} \sim \pi_{\text{off}}} [(Q^{\pi_\theta}(s_{t+1}, a_{t+1})) \\ & - Q^{\pi_\theta}(s_t, a_t)]^2] \end{aligned} \quad (9)$$

Achieving accurate error estimation requires alignment between the evaluation policy and the target policy. However, offline RL inherently aims to discover policies that outperform the original policy, which introduces an unavoidable distributional shift. This shift occurs when the state-action distribution under the learned policy diverges from that under the original policy, leading to inaccuracies in value estimation due to cumulative biases from sampling and function approximation.

To address this distributional shift, Offline RL methods are broadly divided into model-based and model-free methods. Model-based methods leverage learned dynamics models to estimate uncertainty and handle distributional discrepancies. Prominent examples include MOREl [46], MOPO [47], COMBO [48], etc. Model-free methods are further split into explicit and implicit regularization techniques. Explicit regularization methods, such as (Batch-Constrained Q-learning) BCQ [49], (Bootstrapping Error Accumulation Reduction) BEAR [50], Conservative Q-Learning (CQL) [51], etc., impose direct constraints on policy improvement to limit distributional divergence and encourage conservative policy update.

Additionally, the inability to interact with the environment to find more rewarding regions further restricts the performance of offline RL.

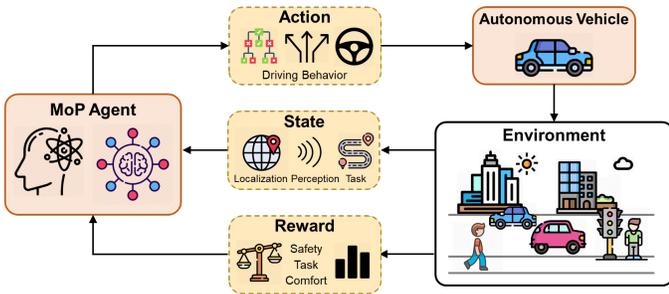


Fig. 4. RL algorithm applied to MoP for AD.

### B. RL-based Motion Planning for Autonomous Driving

MoP for AD generally refers to the planning process for generating feasible states and control sequences, and it is aimed at achieving safe and efficient movement. It generally requires a given route, or specified task to consider the evolution of the agent and environment dynamics [52]. A schematic of application of RL to MoP for AD is shown in Fig. 4, where the RL agent learns a driving policy from trial-and-error data. The ego vehicle (EV) states and environmental observations usually constitute the state space of the RL agent, and the action output by the RL agent is used for high-level behavioral-type decisions and for direct control of the vehicle maneuvering at a low-level.

For instance, value-based methods are widely used for behavioral planning in MoP [12], [53], [54]. The discrete action output of the value-Based RL fits well to supervisory control solutions where the higher level commands by the RL planner are implemented by the legacy motion control systems [55]. Meanwhile, policy-based methods can output the continuous control commands such as the steering angle and acceleration [56], [57], [58]. In recent years, the superior performance of the actor-critic methods has led to the direct learning of vehicle control commands becoming the mainstream direction in the current research [13], [59], [60]. Furthermore, it transpired that HRL motion planning has a similar algorithm architecture to the rule-based modular approach. Different sub-agent can be created to learn the policy for decision-making, trajectory planning, motion control tasks separately. Some works [61], [62], [63], [64] train the high-level policy to select discrete semantic decision actions, and then utilize a separate low-level policy to directly control the steering angle and acceleration, achieving more precise and flexible motion control while ensuring clear driving objectives. Several studies have used MARL to better model the interaction between vehicles and provide a global perspective on multi-vehicle control. CTDE methods are commonly used to generate multi-vehicle policies when collaborative tasks are involved, such as maintaining formation, and cooperative lane changing or merging [65], [66].

The core advantage of RL is its theoretical framework, which is focuses on optimizing decisions for long-term returns rather than merely imitating observed behavior. This capability enables RL to potentially outperform human drivers by uncovering innovative driving policies that extend beyond traditional rule-based models.

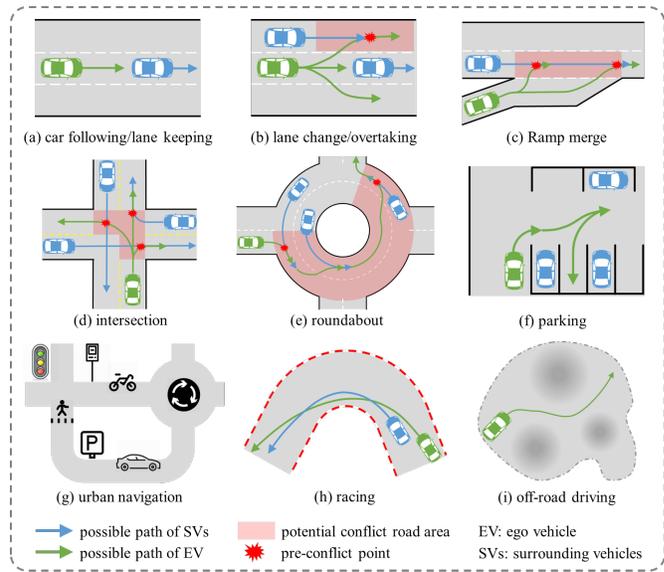


Fig. 5. Illustration of RL-based MoP for different driving tasks.

## III. A REVIEW FROM THE DRIVING TASK PERSPECTIVE

Most RL MoP studies in the AD field have focused on specific driving tasks, ranging from single tasks such as lane keeping or car following, to multi-task integrated urban navigation, etc. Different driving tasks and their application scenarios usually have their own unique characteristics, which have an enormous impact on the design of RL models. From the driving tasks perspective, this section describes the scenario characteristics and task requirements of different driving tasks. On this basis, we review the RL-based MoP literature under these tasks (as illustrated in Fig. 5), especially how they design an RL model for AD.

### A. Car Following/Lane Keeping

Car following (CF) and lane keeping (LK) are the two simple autonomous driving tasks for which early applications of RL approach have been explored. The former task aims to adjust the longitudinal speed to maintain a suitable speed between the EV and the front vehicle (FV), whereas the latter focuses on lateral distance control.

Zhu et al. [67] use the speed of EV, the speed difference from the FV, and the headway distance as the observed states, and then directly control acceleration using DDPG algorithm. Meanwhile, Time to Collision (TTC), Time Headway, and jerk are prioritized in the design of the reward function. Furthermore, Shi et al. [65] add the rear vehicle (RV) information into the state space, and correspondingly considers the safety reward and efficiency reward related to the RV. Chen et al. [54] further consider the cut-in maneuvers of vehicles from adjacent lanes. Specifically, a target acceleration it selected via DDQN by discretizing a continuous acceleration interval.

For the LK task, Kendall et al. [68] control the steering angle and target speed through the DDPG algorithm, with the state space containing the vehicle's speed and steering angle, as well as monocular camera images from the environment. Notably, they conducted real-world experiments on

a 250 meter section of road, using a modified Renault Twizy vehicle to learn the driving policy online. Moreover, Peng et al. [53] exploit D3QN to control quantized steering angle and acceleration values, promoting the EV follows the road centerline. Given the target path points, Tian et al. [69] add the lateral distance from the target path to the preview points in the observation space, and use two actor networks to control steering angle and vehicle speed, improving motion accuracy.

### B. Lane Change/Overtaking

Lane change (LC) is a common driving maneuver, and it causes large collision accidents [70]. The purpose of LC is to avoid collision and improve driving efficiency. On structured roads, lane changes are often accompanied by overtaking behavior, i.e. a continuous LC to obtain a faster driving speed. Some researchers divide overtaking maneuvers into three phases: moving to the target passing lane, overtaking another vehicle and then moving back to the original lane [71].

Many studies address such tasks through high-level behavior planning with physical feature inputs, such as continuous feature states of both the EV and surrounding vehicles (SVs) (e.g. surrounding six [72], or eight [12] vehicles, and vehicles within a certain range). Reference [13] use a discrete state grid of the surrounding environment as the input. Specifically, references [55], [62], [73], [74] use DQN and its improved algorithms to output three semantic actions—lane change to the left (LCL), lane change to the right (LCR), and lane keeping (LK)—focusing exclusively on lateral behaviors. Among them, reference [74] realizes the first application of RL lane-changing policy in the real world. Based on this, reference [72] further incorporates ac/deceleration in the action space, but do not fundamentally change the output form. These actions with low control granularity still limit the impact of the RL agent on vehicle’s maneuverability. More recently, an increasing number of studies have used the DDPG, PPO, SAC techniques, etc., to directly control the steering angle and acceleration [13], [60], [75]. Meanwhile, sensor data from LiDAR [76], camera [77], etc., are utilized as observation inputs to achieve direct mapping between perception and control commands. However, even small differences in adjacent inputs may cause significant fluctuations in the control commands output from the policy network.

Several studies proposed approaches to generating trajectory targets to indirectly control vehicles, aiming to balance flexibility and stability of lane change/overtaking behavior. For example, Yu et al. [71] select a trajectory from a given discrete trajectory set, which is then sent to a tracker module. Lu et al. [78] allow the agent to output a target point location as well as a desired vehicle speed, and then optimize the motion sequence for lane change or overtaking.

Despite differences in state and action space, most studies are consistent in reward design because of driving task characteristics. The safety reward is essential and is crucial and is typically associated with collision [64], relative distance [55], TTC [67], etc. Efficiency is also important and the efficiency reward is often dependent on the vehicle speed [73], the degree of task completion [79], etc. Other components of the reward function can represent comfort related to acceleration

and jerk [80], as well as adherence to traffic rules [68], e.g., overtaking on the left side. Furthermore, some studies design segmented rewards according to the overtaking phase to represent the goals of different phases [71].

### C. Ramp Merge/Intersection/Roundabout

1) *Ramp Merge*: the ramp merge task is typically prompted by the driving lanes of the EV and SVs will overlap in the future, resulting in a forced interaction between EV and SVs. The EV needs to adjust its speed before arriving at the lane merge point to find an acceptable gap between the SVs in the target lane. Therefore, the merge task requires the road geometry be used as an extra observation input compared to the lane changing/overtaking tasks [81], [82].

The simplest way is to learn longitudinal control, which ensures the EV drives to the right place at the appropriate time. Notable approaches include learning the ac/deceleration behavior [83] or the speed control command [84]. Some merge tasks allow the EV to complete the merge operation anywhere within a lane between a start and an end merge point [85]. Like lane changing/overtaking, RL solutions for merging can generate steering angle and acceleration commands to directly control the vehicle for more flexible merging maneuvers [59]. The merge task typically introduces an additional reward for reaching the target lane [86], while some studies also consider the driving distance or time to complete the merge [83], [87]. Due to the stronger interactions in merging scenarios, it is also crucial to devote attention to the collaborative behavior of SVs. Several studies employed game theory to model these interaction [88]. Recent research has advanced with MARL framework [89], [90] that provides each agent with strategies for the merging process. This approach learns the interaction characteristics among vehicles and helps them perform actions with a consistent optimal goal.

2) *Intersection*: Intersection scenario is similar to a ramp merge but with more lane conflicts (turn vs. straight, unsignalized intersections, etc.), complex road structures and road elements, and diverse driving behaviors. This makes intersection one of the most challenging tasks for AD on structured roads [17]. Early approaches focused on controlling vehicle longitudinal behavior through physical feature inputs, by adjusting vehicle acceleration/deceleration [91], or deciding whether to yield or assert the right of way (e.g., wait, pass, yield, take up, give up, etc.) [63], [92], [93], so as to pass through the intersection successfully. Subsequent research has sought to increase agent flexibility by directly controlling vehicle motion and exploiting the reward function reflecting tracking error, collisions, success rate and passage time [94], [95]. Some studies also incorporate additional rewards for violating traffic rules, such as crossing solid lines or running red lights [96], [97]. Other studies break down various scenarios into sub-tasks for training [98], or employ state machines [99] to manage these multiple tasks. However, handing semantic constraints for agents in these approaches remains challenging.

Due to the increasing number of scenario features and state observation inputs, recent research has attempted extracting feature information directly from raw sensor data. Ren et al. [95] utilize LiDAR point cloud data to extract features

of traffic participants, including vehicles, bicycles, and pedestrians. In [60], multi-view camera images are projected into the bird's-eye view (BEV) format to capture global scene features. References [97], [100] address sensor occlusion at intersections with roadside sensing information supplement. Similarly, several studies have explored the use of MARL to address driving through intersections. Antonio et al. [101] iteratively process the observations of the relative positions, speeds, and driving intentions and then individually control the desired speed of each vehicle. Zhao et al. [102] integrate the positional and speed information of each vehicle into a global state feature, and output a joint action based on each vehicle's desired speed.

3) *Roundabout*: the roundabout scenario can be viewed as a combination of two T-junctions and a circular multi-lane road. It involves both merge/intersection and lane change/overtaking tasks. At the entrance, the EV is required to perform a merging task similar to that at an intersection, while driving in the roundabout may involve lane changes and moving to the outside lane before exiting. This combination of multi-scenario features and multi-tasking creates a significant challenge for the MoP system, especially for environment encoding.

Zhang et al. [13] divide the state input into an environmental representation (ER) and a task representation (TR). The ER focuses the physical features of eight SVs, and the TR includes relative lane and exit distances. The action space represents macro-scale behavior (change lane or not) and mesoscale behaviors of desired acceleration and action time. The authors of [13] use MPC to generate pre-trained trajectories, which are embedded in the actor-critic network to improve the learning efficiency. Additionally, TR vectors are replicated in the environment encoding process to emphasize task success in the later training stages.

#### D. Parking

A parking scenario involves an unstructured environment in an urban area with partially regular roads and perpendicular, parallel or diagonal parking slots. Parking tasks have been widely studied, and automated parking technologies have been deployed in many produced vehicles. Current research tends to improve parking flexibility, i.e., reduce "D-R" gear shifting in unconventional or narrow parking spaces [103]. RL can be used for finding the optimal parking path.

Most studies directly control the vehicle's motion during the parking process, and use the position, velocity, and heading angle of the EV as necessary observation inputs, see, e.g., [58], [104]. For the reward design, safety and parking targets are necessary to encourage the EV to reach the target position and heading angle without collision. Additionally, parking is encouraged to be completed as quickly as possible to enhance efficiency [58], and smoothness of control commands is promoted to improve comfort [104]. [57] proposes a unified approach capable of coping with perpendicular, parallel or diagonal parking slots, with proximity sensor data as one of the observation inputs.

#### E. Urban Navigation

Urban scenario encompasses almost all above driving tasks. Unlike driving tasks in a single scenario, navigation in urban areas requires agents to be able to simultaneously understand the characteristics of different scenarios, including various combinations of merges, intersections, roundabouts, etc. Urban scenarios involve vehicles and pedestrians with different characteristics, and much more semantic traffic elements, which can lead to more complex interactions.

Most related studies use the E2E framework to address this task. For example, Reference [105] deals with a BEV rendering of the scene (map, routing, surrounding objects and previous ego states) that is compressed to a low dim latent space using a Variational Auto-encoders (VAE). The latent state is then fed to an RL controller.

Scenario generalization capabilities need to be specifically considered in navigation tasks, as long-tailed/out-of-distribution scenarios that are difficult to simulate with training data may appear in a city at any time. Anzalone et al. [106] present an E2E RL framework in the CARLA environment, utilizing the entire town map. The training process ranges from simple routing under speed constraints to a more complex phase involving randomized starting points and dynamic pedestrian scenarios. Zhan et al. [107] utilize transformers to aggregate a collection of variable-sized and unordered traffic participants into a state vector. They implement offline training and decouple it from downstream RL control to prevent overfitting and improve generalizability. In another study, Jin et al. [108] combine VAE with Generative Adversarial Networks (GAN) to encode input RGB images, thereby reducing the state dimension while also lowering the collision rate in adverse weather conditions. Hu et al. [109] propose a query-based design and connect perception, prediction and planning nodes in an integrated large parameterized E2E framework that incorporates full-stack driving tasks.

These approaches focus more on the generalization capability and comprehensive driving performance of RL-based MoP under different tasks, which usually requires a large network model and large amounts of training data.

#### F. Racing

Racing is a specialized, niche automotive activity, that typically takes place on a fixed, enclosed circuit, requiring intense competition with other vehicles. The object of the racing task is to drive as fast as possible. Similar to lane changing/overtaking, racing requires continually surpassing SVs to improve the position within the vehicle pack, but the lane changing process is not bound by clear lanes and only needs to be within the racetrack. Additionally, racing places greater emphasis on controlling the vehicle to follow the best route, especially around the corners.

The current way of controlling racing cars is dominated by the direct output of their acceleration and steering angle, but some studies also introduce high-level semantic behaviors for collaborative or adversarial purposes [41]. Typically, the observation inputs contain the position, speed, and heading angle of the EV, radar [110] or BEV images [111], and track progression information [6]. Moreover, some studies have

further considered information such as tire temperature and engine speed [41] to achieve better control at high speeds. In particular, the reward function should consider not only safety and efficiency, but also overtaking rewards [110] and even sportsmanship rules [6] due to the competitive nature of racing. Specifically, some studies use human presentation data to assist with policy updating or pre-training [41]. To obtain more robust overtaking strategies, some studies use curriculum learning for staged training [111]. In addition, [76] pursues to MARL to enhance agent consideration for high interaction and competitiveness. Notably, RL is now able to compete directly with and overtake top human racing players in a racing simulator [6].

### G. Off-road Driving

In the context of off-road driving task, distinct road boundaries and traffic signs commonly found in urban scenarios may be lacking. In such situations, it is more important to consider the terrain, irregular obstacles, and cartographic data.

Huang et al. [112] propose an RL-based 2.5D multi-objective path planning method. On the processed small-size 2.5D maps, a reward function that combines terrain, distance, and boundary information is designed to achieve multi-objective path planning that balances energy consumption and distance through DQN method. A multi-objective RL is proposed in [113] as a solution to the path planning problem of an unmanned mining truck in an irregular environment. The feasible path is obtained by extrapolating the steering angle output by RL within a kinematic model in the simulation. This technique is able to plan a path from the starting point to the target in less time than the hybrid A\* algorithm. Zhang et al. [114] combine RL with Dynamic Window Approach (DWA). The state space includes local elevation map, vehicle attitude, obstacle object features, and the target information, while the action space consists of the weight parameters of the evaluation function in the DWA and the time period.

In addition, RL in the off-road scenario has to accommodate more complex with vehicle dynamics. Wang et al. [115] propose a model-based RL algorithm that trains a probabilistic dynamic model to consider model-uncertainty, thereby improving the accuracy of system predictions. More specifically, they train a System Identification Transformer and an Adaptive Dynamics Model under a variety of simulated dynamics, improving robustness and adaptability.

## IV. LESSONS FOR RL-BASED MOP DESIGN

Most studies on RL-based MoP focus on specific driving tasks. Each driving task typically involves distinct scenario characteristics and task requirements, which significantly affect the RL agent design. The effective application of RL to a particular driving task requires careful consideration of several critical design elements, including the design of the observation input, the action output, the reward function, and the training environment. According to the review in Sec. III, these design components vary substantially across different driving tasks and algorithms. In addition, these manual designs strongly influence subsequent self-learning and policy

iterations. This section summarizes and analyzes the patterns of RL model design, aiming to extract lessons learned from various driving tasks and to provide clear guidelines for the application of RL-based MoP techniques for AD.

### A. Observation Input

#### 1) State Space Design

Unlike imitation learning, which constructs loss functions directly from expert data to establish input-output mappings, the RL agent collects feedback indirectly through interactions with the environment. This enables the RL agent to discover optimal solutions beyond expert data, but also makes establishing input-output mappings more difficult. Therefore, using low-dimensional feature data as model inputs simplifies the problem and accelerates the convergence process.

*a) Physical Features:* Physical features, which are highly abstract features processed by perception modules, are the most commonly used input information. They usually include the motion state of the EV (e.g., position, heading, speed, chassis states) and surrounding traffic participants (e.g., relative position, speed, distance).

*b) Sensor Input:* To reduce the loss of input features, sensor information (e.g., camera images, bird's-eye-view (BEV), LiDAR point cloud, etc.) is fed directly to the RL agent, aiming to achieve higher performance. In most cases, even with raw sensor inputs, abstract physical features remain essential. The reason is that it is difficult for RL agent to learn effective policy from high-dimensional and multi-source sensor information entirely on its own. In recent popular intersection/urban navigation tasks, the observation inputs are usually multimodal, including both physical features and multi-source sensor information. In some collaborative tasks, such as formation driving and cooperative merging, extra V2X communication information is considered.

*c) Auxiliary Representation:* Additionally, some auxiliary representational inputs may be useful to help the RL agent better understand the surrounding environment and the task requirements. For instance, a grid map [116] or risk potential field informing driving costs in [117]. For navigation in urban scenarios, the prior knowledge including road map, global route, and traffic rules is commonly used.

The detailed categories and descriptions of the observation inputs are provided in TABLE I.

#### 2) Multi-Model Observation Input

As driving tasks of interest to researchers become increasingly complex and RL algorithms continue to advance, research on RL-based MoP has evolved from using primarily physical features as observation inputs to integrating a broader range of information. Notably, the urban navigation tasks predominantly employs E2E architecture, which encompasses most independent driving tasks and has become the most prevalent paradigm for future research. However, E2E RL also results in a very redundant observation space and usually requires the integration of a large number of sensor inputs in addition to physical features. Furthermore, different driving tasks need to be supplemented with specific auxiliary representations, such as map information for urban navigation or terrain information for off-road driving. While these multi-source

TABLE I  
COMMON OBSERVATION INPUT CATEGORIES AND DESCRIPTIONS

Observation Input Category		Description in Detail
Physical Feature	EV State	ego position, head, speed, chassis states (e.g., acc., steering [64], yaw [66], engine speed [6], [41])
	Object State	feature information of surroundings from perception module such as relative position, speed, head, etc., [62], [85], [118]
Sensor Input	Camera	raw visual data from camera [53] [104], or from BEV [111], [78]
	LiDAR	point cloud data with spatial information [110], [57], [78]
Auxiliary Representation	Navigation Information	destination information, navigation route, etc. [105]
	Road Map	road structure topology information [110], [64]
	Grid Map	occupancy grid map [116], or elevation map describing surroundings [66]
	Traffic Rule	traffic light, speed limit, stop lines, etc. [106], [119]
	Risk Field	risk value of surroundings [75]
	History Info	historical trajectory information [73]
	Metrics	indirect metrics describing driving states such as TTC [120], [110]

observation inputs contain nearly all environmental features, directly feeding them into policy networks for learning may increase computational complexity and impede the efficient extraction of latent features. In addition, when these multi-source observation inputs are incomplete or perturbed, policy execution may be unstable.

### B. Action Output

The agent interacts with the environment by executing actions, updating its state accordingly. In the current RL theoretical framework, the form of action output directly determines the model type, and is generally classified into discrete and continuous action. In addition, there are some approaches that use indirect actions and hierarchical actions.

#### 2) Action Space Design

a) *Discrete Commands*: Discrete commands can represent high level decisions such as LCL, LCR and LK for the lateral vehicle behavior or a discrete set of acceleration/deceleration values in term of longitudinal dynamics. They can be viewed as prescribing vehicle semantic behavior. In some early studies, a set of discrete instructions was also used as an action space, thus the agent could to select one of discrete values, e.g., for the steering angle. Moreover, Yu et al. [71] allow the RL agent to select candidate trajectories from a real-time generated trajectory set.

b) *Continuous Commands*: However, by restricting to a discrete set of commands, one may not be able to obtain optimal solutions in various dynamic scenarios, especially at the vehicle control level. In addition, sudden changes in discrete commands can cause jerky oscillations in driving maneuvers. If an RL agent aims to directly control vehicle motion, it typically outputs continuous commands such as the steering angle and acceleration. Therefore, Policy-based and Actor-Critic methods are widely used for direct vehicle control; they can be applied to nearly all AD MoP tasks.

TABLE II  
COMMONLY USED ACTION OUTPUT FOR RL-BASED MoP

Action Output	Category	Description in Detail
Discrete Behavior	Lateral Behavior	semantic behavior e.g., LK, LCL, LCR [62], [118]
	Longitudinal Behavior	semantic behavior e.g., brake/acceleration/speed keeping, discrete acceleration set [53], [54], [121]
Continuous Control	Lateral Control	steering angle, steering wheel angle [59], [57]
	Longitudinal Control	acceleration, pedal degree [110], [58], target speed [64]
Indirect Command	Interaction Intention	take/give away, stop/go, etc. [84], [92]
	Maneuver Parameter	coefficient of polynomial curve [64], sample points, etc. [73] [78], distance for trajectory planning [122]
Hierarchical Action		semantic behavior + control command, etc. [123], output mixed action simultaneously [61], [62]

c) *Indirect Commands*: In addition, an indirect command can be an output to affect the vehicle’s motion. For instance, some studies indirectly plan continuous feasible trajectories by learning trajectory parameters (such as polynomial coefficients, objective function weights, etc.). Through parameterized action based RL framework, heterogeneous trajectory parameters can be generated synchronously, as in recent P-DQN [124], RL-TPA [122], and other approaches. Other works generate interactive actions to determine reference states for MoP, and then solve for optimal trajectories by, for example, constructing an optimal control problem, e.g., RL+MPC [125], etc.

d) *Hierarchical Actions*: HRL uses different networks to determine heterogeneous components of whole actions. The hierarchy can be serial (upper layer outputs discrete behaviors, lower layer outputs control commands), parallel (heterogeneous control commands are output simultaneously), or even hybrid [61], [120], [123]. Complex driving tasks can be simplified and split in this way but may suffer from sparse rewards. In particular, actions in different layers may be executed at different timescales, e.g., the upper layer’s lane change decision may be updated over a long timestep, whereas the lower layer’s steering commands may be outputted over a short timestep. In strictly HRL theory, a synchronously learned factor  $\beta$  is introduced to determine the update timing for upper level action [38], but this remains underexplored in MoP research. In addition, the optimal consistency of the hierarchically generated actions may be affected since the upper-layer network cannot fully access the policy information of the lower-layer network when generating actions. Note also that the parameterized action based RL can be further extended to hierarchical architectures [126].

#### 2) Control Granularity

Early research on RL for AD focused on decision-making for discrete behaviors. In recent years, direct control of the steering angle and acceleration has become a simple and popular choice; notably it facilitates determining continuous driving

actions through a much larger policy network. Thinking about control granularity cannot be ignored. At the same time, applying RL to trajectory-level actions can enrich the control granularity of RL-based MoP and improve the agent’s ability to focus on driving behaviors, control them more accurately and cope with complex, dynamic road environments. However, the loss function in RL is typically generated indirectly based on accumulated reward signals, rather than being directly derived from expert trajectory data, as in IL. Consequently, achieving convergence for high-dimensional waypoint actions is often difficult. Related research based on a parameterized action space [126] has emerged as a promising direction to facilitate action granularity design. Enabling the RL agent to output trajectory waypoints can further exploit the capabilities of RL. For instance, developers could evaluate the RL interaction with the environment objectively, which makes it easier to design a safety guarantee. Furthermore, this could greatly enhance the interpretability of RL MoP.

The commonly used action outputs for RL-based MoP are summarized in TABLE II.

### C. Reward Function

Reward design can significantly influence the performance of the RL agent as it directly informs the loss function required for network updating. Driving is a multi-attribute problem, and these attributes may include the time to reach destination, travel distance, collision, legal compliance, energy consumption, passenger experience, impacts on the traffic environment, etc. [127]. Defining a driving performance metric for an autonomous vehicle involves identifying various attributes and quantitatively describing them, and then combining them into a utility function. Current RL models for AD typically formulate the reward function as a weighted linear combination [16].

#### 1) Reward Attributes

Safety, efficiency, comfort, and traffic compliance are typical attributes considered when designing reward functions for AD. On this basis, different driving tasks can incorporate unique reward functions to better achieve task-specific goals.

*a) Safety:* Whether a collision occurs, including collision with vehicles or pedestrians, or out of road, is a direct measure of the safety. Some metrics that reflect the degree of potential danger are used to jointly describe the safety attribute, e.g., time to collision (TTC), and distance to SVs (DTV).

*b) Efficiency:* For efficiency, speed is a commonly used metric, e.g., driving at the desired speed or as fast as possible. Additionally, success, i.e. reaching the goal or completing the task, and the corresponding costs can be used as other important metrics for specific tasks. Examples include the success rate and time spent passing through an intersection, merging onto a main road, and completing a parking task.

*c) Comfort:* The reward for comfort is relatively straightforward to define and is usually correlated with the smoothness of the vehicle’s motion, including jerk, and lateral acceleration. Some studies use the variance of acceleration and the steering angle to further measure motion smoothness/comfort. In particular, racing and off-road tasks prioritize maneuverability over comfort.

TABLE III  
REWARD DESIGN FOR RL MOP APPROACHES

Reward	Category	Description in Detail
Safety	Collision	penalties for the occurrence of a collision [54], [58], [72]
	Potential Danger	indirect indicators e.g., TTC [67], THW [89], distance, risk value.
	Out of road	penalties for driving out of road [53]
Efficiency	Speed	reward for high speed or closing to desired speed, penalties for low speed. [54]
	Success	immediate reward when the goal/targets are successfully accomplished. [57], [62], [104]
	Success Cost	cost of reaching the goal, e.g., time spent, path distance, gear shifting number [111], [69]
Comfort	Smoothness	oscillation of states, e.g., jerk [64], frequent lane changes [123], lateral acceleration [78], variance of the steering angle and acceleration. [54], [58], [104]
Traffic compliance	Traffic Signal	follow the traffic light and traffic sign [96], [97]
	Lane Rules	overtaking on the left side or right side, driving into the correct lane before the intersection, etc.. [12], [128]
	Other Specificities	alternate right-of-way, road diversions, etc. [93]

*d) Traffic Compliance:* Traffic rules represent a highly complex set of guidelines, encompassing multiple semantic levels of understanding and evaluation [128]. Common traffic rule conformance includes adhering to traffic signals, staying in the correct lane, not speeding, etc. There are also specific rules for certain scenarios, such as alternate right-of-way and road diversions. A suitable generalization paradigm has yet to be established in the literature, since this attribute primarily appears in the urban navigation task and most studies approach it indirectly through multi-modal inputs and expert data labeling.

The common reward functions mentioned above from both actual attributes and shaping rewards are listed in Table III.

#### 2) Reward Utilization

Most RL related studies use weighted summation to combine different rewards. Knox et al. [127] discusses the calculation of weight factor limits, using crash, idle and success attributes as examples. However, human-manual weighting does not effectively harmonize the trade-offs and conflicts between multiple objectives. Parameter tuning methods such as GLIS [129] could be used to optimize the weight coefficient, which is an optional means. In addition, Inverse Reinforcement Learning methods are applied to learn the weight of each attribute [130] or the reward value [131] from expert experience.

Nevertheless, weighted tuning has a limited impact on improving achievable performance, and agents can still be skewed toward larger single-attribute rewards. The recent Multi-Critic approach excels in accommodating multiple objectives simultaneously [132]. Yuan et al. [133] decompose the value estimation based on a single reward function into decentralized estimation based on multiple reward functions through multiple Q-networks, which allows agents to better

balance multiple learning objectives. Moreover, [134] incorporates the context as an input to construct a reward machine to transform the reward functions for different tasks/scenarios, enhancing the adaptability to environmental changes.

### 3) Reward Shaping

When reward signals from objective attributes are sparse, it is a natural idea to encourage and indicate seemingly desirable maneuvers in reward functions, which is formalized as reward shaping [127]. For example, adding a reward for staying near the lane centerline can help a vehicle to quickly learn how to keep on track. However, combining this partially shaped rewards with existing safety rewards may lead the RL agent unexpectedly fall into a local optimum, such as persistently following SVs at a low speed, which is not actually the desired driving behavior. Common shaped rewards via one or more attributes include suggesting zero steering angle [105], increasing the separation distance with SVs [12], overtaking other vehicles [76], etc. While reward shaping improves the learning efficiency, it may reduce the achievable performance by subjectively changing the preference order of the reward function. As Russell and Norvig assert [135], *“It is better to design performance metrics according to what one actually wants to be achieved in the environment, rather than according to how one thinks the agent should behave”*. The survey [127] boils it down to a pithy description: *“specify how to measure outcomes, not how to achieve them.”* Despite its theoretical drawbacks, reward shaping remains effective in RL-based MoP methods as of this time. Until a better learning way emerge, reward shaping techniques, such as risk-aware shaping for safety [136] or directly traffic rule guidance, can enhancing driving performance to a certain extent. However, the potential negative consequences of each shaping operation need to be carefully considered. Designing effective reward functions remains an open problem, limiting the RL performance in MoP for AD as well as in other control tasks.

## V. EXPLORATORY EFFORTS TO ADDRESS CONTEMPORARY CHALLENGES

Although there have been many significant achievements in RL-based MoP, there are still many challenges in applying it to real-world AD systems. Owing to page limitations, we focus on three attributes that have the greatest impact on RL-based MoP for AD, i.e., safety performance, sample efficiency, and generalization capability. Other attributes, such as interpretability and ethics, are not discussed in this survey, and interested readers are referred to [15], [17], which address such topics. This section reviews recent exploratory efforts for these three frontier issues and proposes directions for future research. Since promoting sample efficiency and generalization capability share some common technical aspects, we distinguish them according to the primary motivation for using these techniques to enhance the performance of RL-based MoP.

### A. Safety Performance

Safety is a fundamental requirement for AD. However, the RL agent may sometimes prioritize maximizing the overall reward over ensuring safety, especially under conditions where

multiple objectives are considered. This can lead to unsafe or even disastrous behaviors, which is the most important hindrance to the application of RL to real-world AD [65]. Consequently, an increasing number of researchers have focused on the safety of RL-based MoP methods and have begun to explore the application of Safe RL.

Safe RL is often modeled as the Constraint MDP (CMDP) [137], which additionally minimizes safety-related cost  $C_\pi(s) = \mathbb{E}[\sum_{t=h}^{H+h} \gamma^{t-h} c_{t+1} | s_h = s]$  while maximizing cumulative reward expectations, where  $c_t$  is the safety-related cost value at timestep  $t$ . The objective of CMDP is to find a policy  $\pi_\theta \in \Pi_C$  to maximize the expected reward, where  $\Pi_C = \{\pi_\theta | C_\pi(s) \leq C_{thres}\}$  represents the safe policy set with a cost threshold  $C_{thres}$ . Safe RL applied in the MoP can usually be categorized as: i) Policy objective optimization: This method uses the cumulative cost values on the trajectories to search for safe policies, gradually converging to safe set. ii) Hard safety constraint: Stricter requirements on the safety of each step are imposed during training or testing through predefined constraints. This type of approach can further enhance safety, but is more conservative.

#### 1) Policy Objective Optimization

Constrained Policy Optimization (CPO) is frequently used to guide the generation of a safer driving policy [138]. Wen et al. [139] employed parallel CPO agents to collect sufficient safe and feasible experiences for policy updates, mitigating the driving risks of CPO failures in hazardous situations. Additionally, Lagrangian-based methods transform constrained safety optimization problems into unconstrained problems via Lagrange multipliers. In [64], a Lagrangian network adaptively adjusted penalties for constraint violations, while a feasible value network evaluates policy feasibility. Furthermore, inspired by the amygdala mechanism, Lv et al. [96] employ a fear model to recognize potential dangers and contingencies, aiming to maximize the expected return while adhering to the fear constraint.

Many researchers integrate Control Lyapunov Functions (CLFs) [140] or Control Barrier Functions (CBFs) [141] as constraints. In [142], a CLF based on the relative distance to obstacles is established, treating the collision probability as a risk factor in the critic with the policy gradient to improve safety. Udatha et al. [143] implement a distance-based probabilistic CBF, which is then converted into linear control constraints to ensure that policy updates adhere to safety requirements. Moreover, Yang et al. [144] take this further by learning a barrier function from collected unsafe and initial states, eliminating the need for prior knowledge.

Meanwhile, some works consider uncertainty in safety guiding the agent’s exploration. In [145], RL policy is updated only when its performance confidence exceeds the baseline, achieving safer behaviors. Zhang et al. [146] use variance from ensemble critic networks to encourage exploration and to determine when to switch from a Lagrangian-based approach to a rule-based approach. In [147], CVaR-based distributional critics facilitate the safety policy update, with the policy space adaptively expanding when actions near the boundary are identified as safe.

#### 2) Hard Safety Constraint

Setting driving rules or rule-based MoP as a safety filter is an intuitive way to enhance the policy's safety [148].

Gu et al. [149] propose a method that combines traditional MoP method with RL, where the safety buffers around obstacles constrain the RL output to collision-free path points. Wang et al. [150] develop a CBF that account for both longitudinal and lateral constraints, combined with predefined traffic rules, to ensure EV safety. Some filters are constructed based on conditional criteria. Reference [151] uses Linear Temporal Logic (LTL) based on prior safety rules to assess the current policy's safety, triggering a rule-based emergency response if the RL action is deemed unsafe. References [118] and [125] employ MPC-based longitudinal pre-planning to assess whether a safe and feasible acceleration can be generated. If unsafe, the vehicle remains in its original lane, and the masked unsafe decision is fed back to update the DQN's network.

Uncertainty can also be used in constraint design for safe RL, which typically includes aleatoric uncertainty and epistemic uncertainty [42]. Aleatoric uncertainty can be expressed as the risk from the scenario measured by the distribution of returns. It allows the RL agent to balance risk and efficiency after convergence and achieve performance similar to those trained in a risk-sensitive way [152]. Epistemic uncertainty usually arises from the insufficient scene training and can be represented by the variance of the ensemble network. In [153], the RL policy reverts to a rule-based policy if the uncertainty evaluated exceeds a safety threshold.

Predicted information can be leveraged to ensure safety over a long horizon. In [12], actions are mapped to trajectories, and the risk of each action is assessed based on the EV's own trajectory and the predicted trajectories of SVs, with high-risk actions being discarded and replaced by safer alternatives. The effectiveness of this approach is validated in real-world lane-change experiments with different vehicle speeds and gaps, significantly reducing the risky behavior of the RL agent. Moreover, Krasowski et al. [154] and Gu et al. [64] introduce the concept of a safe action set based on a prediction embedded framework, which is used to replace the actions of RL with safe alternatives in the event of a failure in the following vehicle strategy. Additionally, [155] and [156] construct optimization-based filters to guarantee that the agent remains safe at all times, while minimizing modifications to the RL policy. A related work [157] uses MPC as a filter to ensure that the agent always stays within a safe invariant set.

## B. Sample Efficiency

Owing to the interactive update paradigm of RL, a substantial number of samples are usually required to construct learning experiences with feedback rewards to generate feasible policies, which leads to sample efficiency problem. This problem is particularly evident in the AD MoP field because of the open interaction environment with large state space and hard-to-collect long-tail data, which results in slow driving policy convergence and lower-than-expected driving performance.

Since the complexity of interactions within the environment, it is challenging to objectively and effectively obtain reward signals in an AD task. In addition, many driving tasks exhibit

temporal correlations, which can further amplify the effects of delayed rewards. In addition, the RL agent must spend a considerable amount of time on constant trial-and-error in the massive exploration space. Besides, it is difficult to gain valuable experience to further improve driving policy performance in the late training stage. These factors contribute to the sample data cost.

To address these challenges, researchers have aimed to enable the RL agent to learn more driving experience from limited samples, thus improving the overall performance of MoP, and accelerating its deployment and application in AD.

### 1) Learning from Demonstration (LfD)

To facilitate faster learning of optimal driving by the RL agent, learning from demonstration (LfD) takes an inspiration from human learning styles. LfD can effectively handle initial exploration where the reward signal is too sparse or the exploration space is too large to be covered. A demonstration can usually be a priori rule models or expert data from human drivers or, alternatively, pre-trained policies.

*a) Learning from a rule-based planner:* The RL agent can be simply and directly guided through rule-based policy demonstration. Alighanbari et al. [158] generates switchable policy through the NMPC controller, and the experience generated by NMPC is used to guide the DDPG to speed up learning. Zhang et al. [13] design an optimization-based trajectory planner to offer the possible motion state data of the EV according to different decision parameter values. When recalled in RL, the related planning parameters are quickly obtained via Nelder-Mean search method. In [159], an expert system consisting of constrained iterative LQR and PID controllers is incorporated into RL training to improve sample efficiency in autonomous overtaking tasks. Similarly, Li et al. [12] design a formalized rule-based correction mechanism considering predicted risks, where multi-memory batches are set to store expert guidance experiences to further improve sample efficiency.

*b) Learning from Human-Guidance:* Combining human guidance with RL can be a promising way to alleviate the sample efficiency issue. A common tactic is to use demonstrations from human experts as a sampling experience for the RL agent. DQfD [160] incorporates expert demonstrations into the replay buffer with extra priority. Liu et al. [161] combine the objectives of reward maximization and expert imitation, and then sample the experiences from both the agent's self-exploration and the human demonstrations with an adaptive dynamic sampling ratio. Gao et al. [162] propose a unified Normalized Actor-Critic, where soft policy gradient formulations are used to reduce the Q-values of actions that were not observed from the demonstrations, thereby mitigating learning bias from low-quality demonstrations. In [125], human online interventions are triggered when an agent outputs unfavorable actions, which can limit unsafe exploration during training, and provide demonstrations in complex scenarios. Similarly, Wu et al. [163] establish an integrated framework including human/RL action switch mechanism, advantage-based prioritized experience, and human-intervention reward shaping. Its unique discriminatory ability for the quality of human guidance contributes to better learning performance.

*c) Learning with a pre-trained policy:* A near-optimal policy

extracted from offline demonstrations can be effectively used for online fine-tuning [164]. Huang et al. [165] distill human prior knowledge into imitative expert policy using Behavior Cloning (BC). Subsequently, a penalty term based on Kullback–Leibler (KL) divergence is added to the reward function, making it fast close to the expert policy in online learning. Shi et al. [117] employ DAGGER to train an IL agent for online RL initialization, which only requires a small amount of scene data to address the learning inefficiency under sparse rewards. In [166], Decision Transformer [167], which is an approach lying in between BC and offline RL, is used to extract a lightweight policy from large-scale offline guidance strategies during online interactions. It outperforms policy initialization via both BC and offline RL for safety-critical navigation and AD tasks.

### 2) Task Differentiation

It challenging to directly learn an effective driving policy in a complex MoP task. Decomposing the task into different parts is a feasible way [168]. Instead of learning to deal with the whole task directly starting from a complex environment, the agent learns the different sub-tasks in stages. The initial task guides the RL agent to perform better on the final task [169], reducing the learning complexity and improving convergence and sample efficiency.

*a) Curriculum Learning:* Curriculum Learning (CL) is a training technique that breaks down the learning process into tasks of increasing complexity. It enables incremental learning, which helps premature failure under high complexity and enhances learning efficiency. Traditionally, CL relies on manually defined stages, with task difficulties set by human experts. Shi et al. [170] design three-stage curriculum RL including adaptive cruise control, lane changing, and overtaking tasks with different reward functions. Anzalone et al. [106] progress from static to complex traffic and weather conditions through five stages, with data augmentation in the final phases to learn complex behaviors. Research on automatic curriculum generation has emerged to overcome the limitations of manual curriculum design. Banerjee et al. [111] used Bayesian optimization to automatically select curriculum through probabilistic inference on curriculum-reward functions. Niu et al. [171] decompose driving policy optimization into evaluation, scenario selection, and training. By dynamically estimating failure probabilities and resampling historical scenarios, this method provides real-time curriculum adaptation, improving learning robustness. Reference [172] introduce a task-driven labeled PAMDP based on LTL progression, which decomposes the training task at an abstract level and informs the RL agent of its current task progress. This technique enhances the exploration efficiency but so far limited to robotic grasping, with no application to complex MoP tasks in AD.

*b) Transfer Learning:* Transfer learning (TL) leverages knowledge reuse techniques [173] to utilize knowledge learned from related tasks. Originally, TL was intended to effectively transfer policies to new environments for better generalizability, which will be described in later Sec V.C later. At the same time, TL also contributes to sample efficiency, by allowing accelerating the learning process of new tasks with fewer samples. Specifically, given a set of source domain  $\mathcal{M}_S$  and

target domain  $\mathcal{M}_t$ , TL learns the optimal policy  $\pi^*$  and the target domain by utilizing exterior information from  $\mathcal{M}_S$  and interior information from  $\mathcal{M}_t$ . For example, Yan et al. [174] use the policy trained in the source domain as the initialization policy for the target domain, thereby improving the policy performance and convergence speed in the target domain. Shu et al. [121] improve the control performance and learning efficiency of the Dueling DQN through three transfer rules.

*c) Hierarchical Learning:* Hierarchical learning architecture (as discussed in Sec IV.B) decomposes the overall task at the action-output level, which can also enhance the feasibility of rapidly learning policy for complex tasks [172]. For MoP for AD, generating only the steering angle usually results in the vehicle deviating far from the lane centerline, because it is difficult for an RL agent to quickly distinguish between lane-changing and lane-following behaviors during the learning process [175]. High-level discrete semantic behaviors and low-level control commands can be combined well to achieve more precise and flexible motion control while ensuring clear driving objectives [123]. Xia et al. [176] also define high-level semantic behavior, and they couple them to low-level control actions, which computes fine-grained actions based on coarse-grained decisions that output them synchronously. The parameterized action space has achieved promising results in learning manipulation skills [177], but it has not been explored much in MoP for AD.

### 3) Promoting Exploration

Efficiently exploring the environment and gathering informative experiences is also important for accelerating learning toward the optimal policy [42]. Uncertainty-oriented exploration generally considers epistemic uncertainty and aleatoric uncertainty, similar to the safety considerations discussed in Sec. V.A. Lee et al. [178] leverage epistemic uncertainty to guide the policy in exploring unknown environments with high-uncertainty, allowing the RL agent to develop a more comprehensive understanding of the surroundings. Both types of uncertainty are considered in [179] within a single system to enhance the robustness of exploration against environmental noise. Intrinsic motivation-oriented exploration typically heuristically utilizes various types of reward-agnostic information to promote exploration. In the absence of an explicit reward signal, the RL agent can use intrinsic motivation to evaluate the quality of its actions. Ma et al. [180] use intrinsic curiosity to drive the agent to explore the environment in advance and collect experience. Curiosity here is represented as the error in predicting the outcome of the agent’s actions in its current state, i.e., the agent learns from the prediction error of forward dynamics. Wu et al. [181] use recurrent neural network to generate an intrinsic reward to encounter the RL agent to explore environment, improving exploration efficiency.

## C. Generalization Capability

As noted in Sec. III, most studies have been conducted in low-cost simulation environments tailored to single-task settings. However, task variability can lead to policy failure when applied across different environments. Furthermore, owing to the inherent incomplete limitations of the RL training

process, it has poor generalization ability in rare scenarios [182]. The ability for long-term multi-task learning is required to enhance the generalization ability of RL agents to the variety of ever-changing complex scenarios that real-world AD applications may face. The generalization ability includes both the policy that can be transferred to various driving tasks, and the robustness in response to perturbations during the execution [183]. Some cutting-edge techniques explored for the generalization capability of RL, but they have not yet been well applied in the field of MoP in AD.

### 1) Knowledge Transfer

a) *Transfer Learning*: Knowledge reuse in TL can improve generalization across different but related or similar tasks. Balakrishnan et al. [184] apply a domain randomization technique to generalize the policy learned in the simple WiseMove environment to the high-fidelity simulator WiseSim. Furthermore, Kevin et al. [185] combine domain adaptation and domain randomization techniques. By integrating virtual training with real-world data, they reduce the sim-to-real transfer gap in AD applications. Hieu et al. [186] pre-train on demonstration data using a combination of temporal difference and supervised loss, and then continuously update the policy by incorporating demonstration data with newly collected data, resulting in strong performance across different road conditions and weather conditions. Shoeleh et al. [187] propose a skill-based transfer learning and domain adaptation method, which helps the agent discover the state-action mapping that represents the relationship between the source and target tasks, thereby providing knowledge generalization across multiple tasks.

b) *Meta Learning*: Meta-RL aims to producing a broadly generalizable policy [188]. Generally, Meta-RL consists of an inner loop, which focuses on learning the specific task, and an outer loop, where the agent extracts knowledge from multiple tasks to improve its adaptability. In [189], following meta-training on the lane-change task under different traffic densities, the policy is able to safely handle meta-testing scenarios with high traffic densities. Deng et al. [190] utilize parallel unfolding and multi-task objectives, meanwhile they design a two-stage constraint adaptation strategy to achieve rapid adaptation to new tasks by reusing meta-training data.

c) *Continual Learning*: Continual RL (Cont-RL) aims to address the challenge of the ability limitation in handling new tasks without forgetting previously acquired knowledge, which enables continuous learning and adaptation to new environments [191]. Cont-RL requires an appropriate balance between the old and new tasks, with adequate generalizability to accommodate their distributional differences. Wei et al. [192] propose a shared feature extractor with an EWC loss to mitigate catastrophic forgetting and perform velocity control tasks across different environments. Cao et al. [193] introduce a disengagement-case imagination augment continual learning (DICL) method, which is capable of constructing imagination-based environments corresponding to disengagement cases and then improving driving policies within them.

### 2) Policy Stability

a) *Disturbance Robustness*: Robust RL focuses on learning policies that exhibit performance robustness against system external disturbances, such as model mismatch and environ-

mental perturbations. Typically, robust RL is modeled as a two-player zero-sum Markov game, where an adversarial agent trains alongside the ego agent to maximize disturbances, forcing the RL agent to develop a robust policy in response to these disturbances. He et al. [194] develop an adversarial agent that maximizes the Jensen-Shannon (JS) divergence between the policy and the original policy under observation disturbances. The RL agent incorporates the JS divergence as a constraint and use Lagrangian dual optimization to update its policy, thereby ensuring robustness to observation disturbances. Similarly in [86], the White-Box Adversarial Attack technique is employed to amplify the disturbance of each observation. Then, the policies of both the ego agent and the adversarial agent are weighted to output mixed actions, simulating the disturbances caused by environmental changes. Additionally, high-uncertainty RL policies can be replaced with more stable baseline policies [153], enabling timely adaptation to changes in the environment .

b) *Uncertainty Adaptation*: Uncertainty can also be leveraged to improve RL generalization performance. Epistemic uncertainty can reveal the test scenarios that are underrepresented in the training. Lutjens et al. [195] use MC-Dropout and Bootstrapping to achieve parallelized epistemic uncertainty estimation to promote more cautious actions in unknown environments, thereby improving policy generalization. In addition, Hoi et al. [196] propose a risk-conditioned distributional soft actor-critic method that learns risk-sensitive policies based on aleatoric uncertainty. It supports the adjustment of risk-level sensitivity without retraining, enabling safe generalization across various scenarios. Some works further attempt to leverage both epistemic and aleatoric uncertainty within a system to improve adaptability to the environment [197]. Hoel et al. [152] propose Ensemble Quantile Networks (EQN), where Bayesian estimates of epistemic uncertainty are obtained through model ensemble methods, which are used to select actions with lower risk in unknown environments. Meanwhile, aleatoric uncertainty is implicitly learned through quantile functions, balancing risk and time efficiency, thereby further enhancing the generalization capability.

### 3) Scenario Representation

Establishing an effective and compact feature representation from observations that supports subsequent policy reasoning is also a key to improving generalization performance.

Duan et al. [198] map surrounding vehicles' features into encoding vectors, which are then summed element-wise to create a representation set, ensuring that the policy network remains unaffected by vehicle arrangement and preventing strategy fluctuations. Similarly, an MLP encoding function and summation operator are constructed to handle traffic flows with participants of varying types and quantities [95]. The transformer is commonly employed to enhance scene understanding. G. Zhan et al. [107] construct a transformation module to extract observations from surrounding participants, the ego vehicle, and traffic lights, and use an Aggregation Module to combine these features into a fixed-dimensional representation, enabling the agent to adapt to dynamic environments with varying numbers of traffic participants.

LLMs present a promising avenue for improving scenario

understanding in MoP [199]. To release the agent from the burden of understanding the multi-modal data, LLMs can be used to extract meaningful feature representations and translate semantic or task information. Reference [200] model graph-structured reasoning through perception, prediction and planning question-answer pairs to mimic the human reasoning process, enabling the agent to correctly handle unseen deployment on Waymo after training only on NuScenes. Reference [201] argues that the absence of task-relevant representations may hinder the mapping of the network from state to reward. Based on this motivation, they employ LLMs to generate task-related state representations accompanied by intrinsic reward functions for RL, which apply to both continuous and discontinuous reward scenarios, improving the adaptability of RL to new tasks.

#### D. Open Challenges and Outlook

##### 1) Safety Consideration

**a) Safety Evaluation:** Both policy objective optimization and safety hard constraints require an evaluation of the risk of unsafety. Most methods rely on the assumption of a priori environment dynamics, such as simple safety rules, safety sets, state predictions, etc. [12], [64]. It is difficult to validate these assumptions until the autonomous vehicle is in action, which results in a model mismatch in the explicit knowledge of environmental dynamics in the open interactive environment. We believe that learning environment dynamics offers a promising solution to alleviate this problem by performing several step forward simulation with learned dynamics. The mismatch problem can be mitigated while retaining the inherent utility of the models and providing better generalizability. Moreover, further incorporating RL's learning uncertainty into the risk evaluation process allows for a more comprehensive identification of dangerous scenarios and safety vulnerabilities with limited priori knowledge.

**b) Trade-off between Safety and Rewards:** It is crucial to consider the trade-off between rewards and safety performance in an AD environment with complex interactions. On the one hand, if strict conservative cost functions are adopted, this may lead to poor reward utility. In contrast, open constraints and costs can lead to unsafe policy [202]. On the other hand, policy-objective optimization encourages the agent to converge on a highly rewarding policy that satisfies the constraints, but it lacks theoretical safety guarantees. Hard safety constraints can limit the exploration space and tend toward an overly conservative driving policy [19]. Research on effectively combining the strengths of the two approaches is urgently needed. Importantly, less conservative but equally meaningful safety assurances depicting practically acceptable assumptions [203] is likely to become a technical priority.

##### 2) Efficiently Learning

**a) Knowledge Integration:** Knowledge integration, i.e., the incorporation of external knowledge into the RL training, is an important direction for improve sample efficiency. Embedding human or prior knowledge into RL policy has been shown to be effective in improving sample efficiency. However, the performance of the current policy extracted from expert

demonstrations remains a concern. RL with human feedback (RLHF) [163], [204] is a popular approach, where reward models are learned from human-evaluation data. However, it becomes difficult to maintain when dealing with extensive training tasks. This area holds promise for further exploration. Future research should focus on how to efficiently extract quality human driving knowledge and effectively combine offline and online learning methods to integrate this prior knowledge into RL models.

**b) Efficient Exploration:** In environments with sparse, delayed rewards, several exploration methods have yielded promising results through uncertainty guidance or intrinsic motivation. However, most related research has not been applied to MoP tasks for AD. The currently dominant RL-based MoP paradigm integrates large multi-source observation inputs, and the difficulty of exploration increases as the size and complexity of the state-action space grow. This is accompanied by a substantial increase in the computational cost of learning uncertainty, while the intrinsic motivation for facilitating exploration becomes more challenging to construct. How efficient exploration can be achieved in a large and complex state-action space remains unclear.

On the basis of task differentiation, one promising way is to establish a universal approach to extract the hierarchical structure of different environments, such as the recently popular parameterized action space design. In addition, representation learning [205] has been leveraged in several recent studies that include improving policy performance in environments with image states and hybrid actions.

##### 3) Generalization

**a) Simulation Fidelity:** Deviations between simulations and the real world may lead to policy bias during implementation. Therefore, focusing on the development of more realistic and diverse simulation environments becomes critical. Many early works focused on accurately modeling the kinematic and dynamic behaviors of vehicles, creating simulation scenarios that do not adequately reflect the real world. In recent years, game engine-based simulators have provided more physically and visually realistic data. The collection of multi-modal data from the real world and then reconstructing the environment have been widely explored recently. Data-driven simulators are beginning to replace game engine-based simulators by generating synthetic data directly from real data, achieving high fidelity with low costs [206]. Despite the high data fidelity achieved through learning from real data, online interaction capabilities remain insufficient in current technologies. Combining training in low-fidelity simulators with validation in high-fidelity real data, or collecting data in realistic field testing to continuously improve the policy are seemingly viable technical routes. In general, Closed-loop interaction model generation still needs to be further explored.

**b) Reasoning Ability:** RL agents lack a basic understanding of the world and often rely on extensive trial-and-error to make rational decisions and understand the correlations between different factors. They may also face difficulties in recognizing and learning invariant mechanisms from limited environments and tasks [207]. Causal models, which study how relevant features of the world interact with each other,

formalize knowledge in a formative way and use invariance for effective knowledge transfer. Causal reasoning mechanisms that incorporate additional assumptions or prior knowledge to analyze and understand behaviors and their consequences enable agents to imagine and gain insights from scenarios that are missing from the collected data. Given the success of causal reasoning in various fields such as computer vision [186], causal RL has been recognized as an understudied but significant research direction with the potential to significantly improve the performance of RL-based MoP in generalization problems.

**c) Evaluation of Evolution:** For tasks with clear feature distinctions, such as overtaking, merging, and intersections, the generalization ability of RL-based MoP can be measured by comparing individual task metrics. However, real-world driving tasks are continuously changing and involve complex, coupled combinations of scenarios. Thus, the agent needs to flexibly and efficiently update and evolve to cope with edge cases encountered during driving. Having separate validation or testing phases may not effectively reflect the effectiveness of continuous evolution. Designing a broad and novel set of metrics to enhance generalization capabilities for RL agents is valuable, and at the same time a challenge that still needs to be explored more deeply by the research community.

**d) LLMs Enhancement:** The emergence of LLMs represents an important milestone in the field of natural language processing, and they have shown powerful capabilities in many real-world applications [208]. References [209], [210] exemplify integrating LLMs to enhance the interaction and generalization of AD systems. With extensive pre-trained knowledge and a high level of generalizability, the integration of LLMs and RL is considered a key development. LLMs could offer several capabilities to facilitate the generalization of RL-based MoP: i) enhancing multi-modal information understanding to provide predictions or suggestions from the context, thus reducing the need for RL agents to interact with a broad range of environments; ii) designing integrated rewards based on multi-disciplinary attributes and adaptively adjusting them based on scenario understanding, enhancing the multi-objective adaptability of learned policies; and iii) providing decision-level demonstration or guidance and further extracting driving knowledge for RL agents.

The application of LLMs to AD is still in its infancy, mainly using interface calls and model fine-tuning [210], [211], with no well-developed work yet combining LLMs deeply with RL-based MoP method. However, preliminary efforts show promise in achieving advanced functionality that facilitates the generalization capabilities of RL-based MoP. Importantly, the biases and hallucinations inherent in LLMs that may lead to distorted or inaccurate interpretations of multi-modal inputs [208], as well as the large computational costs and response times required by LLMs are currently significant challenges for applying LLMs in practice.

## VI. CONCLUSION

With its ability to explore and optimize policies in complex, dynamic decision-making tasks, reinforcement learning (RL)

has emerged as a promising approach for addressing motion planning (MoP) challenges in autonomous driving (AD). This survey provides a comprehensive review of RL-based MoP for AD, focusing on lessons learned from the driving task perspective. We outline the basic theory of RL methodologies, and then delve into their applications in MoP for diverse driving tasks. Scenario-specific features and task requirements are analyzed to illuminate their influence on RL design. On this basis, we summarize key experiences and extract insights for future implementations. Furthermore, we discuss three key frontier issues in RL-based MoP for AD, summarize how some representative emerging technologies are trying to solve them (especially over the past three years), and propose related open issues and future outlooks.

We observe that other approaches and technologies in the field of artificial intelligence are crucial for facilitating the development of RL-based MoP. Future research directions will explore the integration of these advanced methods with frontier issues to promote RL to build AD systems with a better understanding of the world.

## REFERENCES

- [1] R. S. Sutton, A. G. Barto, *et al.*, *Reinforcement learning: An introduction*, vol. 1. MIT press Cambridge, 1998.
- [2] B. Zheng, S. Verma, J. Zhou, I. W. Tsang, and F. Chen, "Imitation learning: Progress, taxonomies and challenges," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 35, no. 5, pp. 6322–6337, 2024.
- [3] B. R. Kiran, I. Sobh, V. Talpaert, P. Mannion, A. A. A. Sallab, S. Yogamani, and P. Pérez, "Deep reinforcement learning for autonomous driving: A survey," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 6, pp. 4909–4926, 2022.
- [4] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, *et al.*, "Mastering the game of go without human knowledge," *Nature*, vol. 550, no. 7676, pp. 354–359, 2017.
- [5] O. Vinyals *et al.*, "Grandmaster level in starcraft ii using multi-agent reinforcement learning," *Nature*, vol. 575, no. 7782, pp. 350–354, 2019.
- [6] P. R. Wurman, S. Barrett, K. Kawamoto, J. MacGlashan, K. Subramanian, T. J. Walsh, R. Capobianco, A. Devlic, F. Eckert, F. Fuchs, *et al.*, "Outracing champion gran turismo drivers with deep reinforcement learning," *Nature*, vol. 602, no. 7896, pp. 223–228, 2022.
- [7] E. Kaufmann, L. Bauersfeld, A. Loquerio, M. Müller, V. Koltun, and D. Scaramuzza, "Champion-level drone racing using deep reinforcement learning," *Nature*, vol. 620, no. 7976, pp. 982–987, 2023.
- [8] S. Teng, X. Hu, P. Deng, B. Li, Y. Li, Y. Ai, D. Yang, L. Li, Z. Xuanyuan, F. Zhu, and L. Chen, "Motion planning for autonomous driving: The state of the art and future perspectives," *IEEE Trans. Intell. Veh.*, vol. 8, no. 6, pp. 3692–3711, 2023.
- [9] A. Tampuu, T. Matiisen, M. Semikin, D. Fishman, and N. Muhammad, "A survey of end-to-end driving: Architectures and training methods," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 33, no. 4, pp. 1364–1384, 2022.
- [10] W. Schwarting, J. Alonso-Mora, and D. Rus, "Planning and decision-making for autonomous vehicles," *Annu. Rev. Control Robot. Auton. Syst.*, vol. 1, no. 1, pp. 187–210, 2018.
- [11] Z. Li, J. Hu, B. Leng, L. Xiong, and Z. Fu, "An integrated of decision making and motion planning framework for enhanced oscillation-free capability," *IEEE Trans. Intell. Transp. Syst.*, vol. 25, no. 6, pp. 5718–5732, 2024.
- [12] Z. Li, L. Xiong, B. Leng, P. Xu, and Z. Fu, "Safe reinforcement learning of lane change decision making with risk-fused constraint," in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, pp. 1313–1319, 2023.
- [13] Y. Zhang, B. Gao, L. Guo, H. Guo, and H. Chen, "Adaptive decision-making for automated vehicles under roundabout scenarios using optimization embedded reinforcement learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 32, no. 12, pp. 5526–5538, 2020.
- [14] L. Chen, Y. Li, C. Huang, B. Li, Y. Xing, D. Tian, L. Li, Z. Hu, X. Na, Z. Li, S. Teng, C. Lv, J. Wang, D. Cao, N. Zheng, and F.-Y. Wang, "Milestones in autonomous driving and intelligent vehicles: Survey of surveys," *IEEE Trans. Intell. Veh.*, vol. 8, no. 2, pp. 1046–1056, 2023.

- [15] R. Zhao, Y. Li, Y. Fan, F. Gao, M. Tsukada, and Z. Gao, "A survey on recent advancements in autonomous driving using deep reinforcement learning: Applications, challenges, and solutions," *IEEE Trans. Intell. Transp. Syst.*, pp. 1–34, 2024.
- [16] Z. Zhu and H. Zhao, "A survey of deep rl and il for autonomous driving policy learning," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 9, pp. 14043–14065, 2022.
- [17] J. Wu, C. Huang, H. Huang, C. Lv, Y. Wang, and F.-Y. Wang, "Recent advances in reinforcement learning-based autonomous driving behavior planning: A survey," *Transp. Res. Part C Emerg. Technol.*, vol. 164, p. 104654, 2024.
- [18] S. Gu, L. Yang, Y. Du, G. Chen, F. Walter, J. Wang, and A. Knoll, "A review of safe reinforcement learning: Methods, theories, and applications," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 46, no. 12, pp. 11216–11235, 2024.
- [19] W. Zhao, T. He, R. Chen, T. Wei, and C. Liu, "State-wise safe reinforcement learning: A survey," *arXiv:2302.03122*, 2023.
- [20] J. Xing, D. Wei, S. Zhou, T. Wang, Y. Huang, and H. Chen, "A comprehensive study on self-learning methods and implications to autonomous driving," *IEEE Trans. Neural Netw. Learn. Syst.*, 2024.
- [21] J. Duan, Y. Guan, S. E. Li, Y. Ren, Q. Sun, and B. Cheng, "Distributional soft actor-critic: Off-policy reinforcement learning for addressing value estimation errors," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 33, no. 11, pp. 6584–6598, 2022.
- [22] G. Shani, J. Pineau, and R. Kaplow, "A survey of point-based pomdp solvers," *Auton. Agents and Multi-Agent Syst.*, vol. 27, pp. 1–51, 2013.
- [23] R. A. Howard, *Dynamic Programming and Markov Processes*. Cambridge, MA, USA: MIT Press, 1960.
- [24] Q. Huang, "Model-based or model-free, a review of approaches in reinforcement learning," in *International Conference on Computing and Data Science (CDS)*, pp. 219–221, 2020.
- [25] R. S. Sutton, "Reinforcement learning: An introduction," *A Bradford Book*, 2018.
- [26] R. S. Sutton, "Learning to predict by the methods of temporal differences," *Machine learning*, vol. 3, pp. 9–44, 1988.
- [27] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [28] H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," *Proc. AAAI Conf. Artif. Intell.*, vol. 30, 2016.
- [29] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas, "Duelling network architectures for deep reinforcement learning," in *Proc. Int. Conf. Mach. Learn. (ICML)*, pp. 1995–2003, 2016.
- [30] M. Gök, "Dynamic path planning via dueling double deep q-network (d3qn) with prioritized experience replay," *Applied Soft Computing*, vol. 158, p. 111503, 2024.
- [31] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine learning*, vol. 8, pp. 229–256, 1992.
- [32] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, "High-dimensional continuous control using generalized advantage estimation," *arXiv:1506.02438*, 2015.
- [33] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller, "Deterministic policy gradient algorithms," in *Proc. Int. Conf. Mach. Learn. (ICML)*, pp. 387–395, 2014.
- [34] T. Lillicrap, "Continuous control with deep reinforcement learning," in *Proc. Int. Conf. Learn. Representations (ICLR)*, 2016.
- [35] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv:1707.06347*, 2017.
- [36] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *Proc. Int. Conf. Mach. Learn. (ICML)*, pp. 1861–1870, 2018.
- [37] X. Wang, S. Wang, X. Liang, D. Zhao, J. Huang, X. Xu, B. Dai, and Q. Miao, "Deep reinforcement learning: A survey," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 35, no. 4, pp. 5064–5078, 2024.
- [38] J. Lai, J. Wei, and X. Chen, "Overview of hierarchical reinforcement learning," *Comput. Eng. Appl.*, vol. 57, no. 3, pp. 72–79, 2021.
- [39] R. S. Sutton, D. Precup, and S. Singh, "Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning," *Artif. Intell.*, vol. 112, no. 1-2, pp. 181–211, 1999.
- [40] T. D. Kulkarni, K. Narasimhan, A. Saeedi, and J. Tenenbaum, "Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 29, 2016.
- [41] L. Chen, Y. He, W. Pan, F. R. Yu, and Z. Ming, "A novel generalized meta hierarchical reinforcement learning method for autonomous vehicles," *IEEE Network*, vol. 37, no. 4, pp. 230–236, 2023.
- [42] J. Hao et al., "Exploration in deep reinforcement learning: From single-agent to multiagent domain," *IEEE Trans. Neural Netw. Learn. Syst.*, 2023.
- [43] L. M. Schmidt, J. Brosig, A. Plinge, B. M. Eskofier, and C. Mutschler, "An introduction to multi-agent reinforcement learning and review of its application to autonomous mobility," in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, pp. 1342–1349, 2022.
- [44] R. Zhang, J. Hou, F. Walter, S. Gu, J. Guan, F. Röhrbein, Y. Du, P. Cai, G. Chen, and A. Knoll, "Multi-agent reinforcement learning for autonomous driving: A survey," *arXiv:2408.09675*, 2024.
- [45] R. Xu, W. Chen, H. Xiang, X. Xia, L. Liu, and J. Ma, "Model-agnostic multi-agent perception framework," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pp. 1471–1478, 2023.
- [46] R. Kidambi, A. Rajeswaran, P. Netrapalli, and T. Joachims, "Morel: Model-based offline reinforcement learning," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 33, pp. 21810–21823, 2020.
- [47] T. Yu, G. Thomas, L. Yu, S. Ermon, J. Y. Zou, S. Levine, C. Finn, and T. Ma, "Mopo: Model-based offline policy optimization," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 33, pp. 14129–14142, 2020.
- [48] T. Yu, A. Kumar, R. Rafailov, A. Rajeswaran, S. Levine, and C. Finn, "Combo: Conservative offline model-based policy optimization," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 34, pp. 28954–28967, 2021.
- [49] S. Fujimoto, D. Meger, and D. Precup, "Off-policy deep reinforcement learning without exploration," in *Proc. Int. Conf. Mach. Learn. (ICML)*, pp. 2052–2062, 2019.
- [50] A. Kumar, J. Fu, M. Soh, G. Tucker, and S. Levine, "Stabilizing off-policy q-learning via bootstrapping error reduction," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 32, 2019.
- [51] A. Kumar, A. Zhou, G. Tucker, and S. Levine, "Conservative q-learning for offline reinforcement learning," in *Proc. Adv. Neural Inf. Process. Syst.*, vol. 33, pp. 1179–1191, 2020.
- [52] L. Claussmann, M. Revilloud, D. Gruyer, and S. Glaser, "A review of motion planning for highway autonomous driving," *IEEE Trans. Intell. Transp. Syst.*, vol. 21, no. 5, pp. 1826–1848, 2020.
- [53] B. Peng, Q. Sun, S. E. Li, D. Kum, Y. Yin, J. Wei, and T. Gu, "End-to-end autonomous driving through dueling double deep q-network," *Automotive Innovation*, vol. 4, pp. 328–337, 2021.
- [54] Q. Chen, W. Zhao, L. Li, C. Wang, and F. Chen, "Es-dqn: A learning method for vehicle intelligent speed control strategy under uncertain cut-in scenario," *IEEE Trans. Veh. Technol.*, vol. 71, no. 3, pp. 2472–2484, 2022.
- [55] S. Nagesh Rao, H. E. Tseng, and D. Filev, "Autonomous highway driving using deep reinforcement learning," in *Proc. IEEE Int. Conf. Syst. Man Cybern. (SMC)*, pp. 2326–2331, 2019.
- [56] X. Zhang, L. Wu, H. Liu, Y. Wang, H. Li, and B. Xu, "High-speed ramp merging behavior decision for autonomous vehicles based on multi-agent reinforcement learning," *IEEE Internet of Things Journal*, 2023.
- [57] B. Sousa, T. Ribeiro, J. Coelho, G. Lopes, and A. F. Ribeiro, "Parallel, angular and perpendicular parking for self-driving cars using deep reinforcement learning," in *2022 IEEE International Conference on Autonomous Robot Systems and Competitions (ICARSC)*, pp. 40–46, 2022.
- [58] X. Tang, Y. Yang, T. Liu, X. Lin, K. Yang, and S. Li, "Path planning and tracking control for parking via soft actor-critic under non-ideal scenarios," *IEEE/CAA J. Autom. Sin.*, 2023.
- [59] J. Duan, Y. Kong, C. Jiao, Y. Guan, S. E. Li, C. Chen, B. Nie, and K. Li, "Distributional soft actor-critic for decision-making in on-ramp merge scenarios," *Automotive Innovation*, vol. 7, no. 3, pp. 403–417, 2024.
- [60] Z. Huang, J. Wu, and C. Lv, "Efficient deep reinforcement learning with imitative expert priors for autonomous driving," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 34, no. 10, pp. 7391–7403, 2022.
- [61] J. Peng, S. Zhang, Y. Zhou, and Z. Li, "An integrated model for autonomous speed and lane change decision-making based on deep reinforcement learning," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 11, pp. 21848–21860, 2022.
- [62] H. Lu, C. Lu, Y. Yu, G. Xiong, and J. Gong, "Autonomous overtaking for intelligent vehicles considering social preference based on hierarchical reinforcement learning," *Automotive Innovation*, vol. 5, no. 2, pp. 195–208, 2022.
- [63] Z. Qiao, J. Schneider, and J. M. Dolan, "Behavior planning at urban intersections through hierarchical reinforcement learning," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pp. 2667–2673, 2021.

- [64] Z. Gu, L. Gao, H. Ma, S. E. Li, S. Zheng, W. Jing, and J. Chen, "Safe-state enhancement method for autonomous driving via direct hierarchical reinforcement learning," *IEEE Trans. Intell. Transp. Syst.*, vol. 24, no. 9, pp. 9966–9983, 2023.
- [65] T. Shi, Y. Ai, O. ElSamadisy, and B. Abdulhai, "Bilateral deep reinforcement learning approach for better-than-human car-following," in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, pp. 3986–3992, 2022.
- [66] S. Chen, M. Wang, W. Song, Y. Yang, and M. Fu, "Multi-agent reinforcement learning-based decision making for twin-vehicles cooperative driving in stochastic dynamic highway environments," *IEEE Trans. Veh. Technol.*, vol. 72, no. 10, pp. 12615–12627, 2023.
- [67] M. Zhu, Y. Wang, Z. Pu, J. Hu, X. Wang, and R. Ke, "Safe, efficient, and comfortable velocity control based on reinforcement learning for autonomous driving," *Transp. Res. Part C Emerg. Technol.*, vol. 117, p. 102662, 2020.
- [68] A. Kendall *et al.*, "Learning to drive in a day," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pp. 8248–8254, 2019.
- [69] Y. Tian, X. Cao, K. Huang, C. Fei, Z. Zheng, and X. Ji, "Learning to drive like human beings: A method based on deep reinforcement learning," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 7, pp. 6357–6367, 2021.
- [70] H. Jula, E. B. Kosmatopoulos, and P. A. Ioannou, "Collision avoidance analysis for lane changing and merging," *IEEE Trans. veh. tech.*, vol. 49, no. 6, pp. 2295–2308, 2000.
- [71] Y. Yu, C. Lu, L. Yang, Z. Li, F. Hu, and J. Gong, "Hierarchical reinforcement learning combined with motion primitives for automated overtaking," in *Proc. IEEE Intell. Veh. Symposium (IV)*, pp. 1–6, 2020.
- [72] N. Li, D. W. Oyler, M. Zhang, Y. Yildiz, I. Kolmanovsky, and A. R. Girard, "Game theoretic modeling of driver and vehicle interactions for verification and validation of autonomous vehicle control systems," *IEEE Trans. Control Syst. Tech.*, vol. 26, no. 5, pp. 1782–1797, 2018.
- [73] S. Li, C. Wei, and Y. Wang, "Combining decision making and trajectory planning for lane changing using deep reinforcement learning," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 9, pp. 16110–16136, 2022.
- [74] M. U. Yavas, T. Kumbasar, and N. K. Ure, "A real-world reinforcement learning framework for safe and human-like tactical decision-making," *IEEE Trans. Intell. Transp. Syst.*, vol. 24, no. 11, pp. 11773–11784, 2023.
- [75] Z. Qi, T. Wang, J. Chen, D. Narang, Y. Wang, and H. Yang, "Learning-based path planning and predictive control for autonomous vehicles with low-cost positioning," *IEEE Trans. Intell. Veh.*, vol. 8, no. 2, pp. 1093–1104, 2023.
- [76] M. Kaushik, N. Singhanian, and K. M. Krishna, "Parameter sharing reinforcement learning architecture for multi agent driving," in *Proceedings of the 2019 4th International Conference on Advances in Robotics*, pp. 1–7, 2019.
- [77] T. Feng, X. Xu, X. Zhang, and X. Zhang, "An improved ddpq algorithm with barrier function for lane-change decision-making of intelligent vehicles," in *Proc. Artif. Intell.: First CAAI Int. Conf. (CICAI)*, pp. 127–139, 2021.
- [78] X. Lu, F. X. Fan, and T. Wang, "Action and trajectory planning for urban autonomous driving with hierarchical reinforcement learning," *arXiv:2306.15968*, 2023.
- [79] X. Wang and M. Althoff, "Safe reinforcement learning for automated vehicles via online reachability analysis," *IEEE Trans. Intell. Veh.*, pp. 1–15, 2023.
- [80] W. Huang, F. Braghin, and Z. Wang, "Learning to drive via apprenticeship learning and deep reinforcement learning," in *Proc. IEEE Int. Conf. Tools Artif. Intell. (ICTAI)*, pp. 1536–1540, 2019.
- [81] R. Valiente, B. Toghi, R. Pedarsani, and Y. P. Fallah, "Robustness and adaptability of reinforcement learning-based cooperative autonomous driving in mixed-autonomy traffic," *IEEE Open J. Intell. Transp. Syst.*, vol. 3, pp. 397–410, 2022.
- [82] S. Hwang, K. Lee, H. Jeon, and D. Kum, "Autonomous vehicle cut-in algorithm for lane-merging scenarios via policy-based reinforcement learning nested within finite-state machine," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 10, pp. 17594–17606, 2022.
- [83] L. Ye, Z. Wen, H. Zhang, and B. Ran, "A general drl-based framework using mode-selection tangent time projection for mixed on-ramp merging," in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, pp. 3207–3214, 2023.
- [84] T. Tram, I. Batkovic, M. Ali, and J. Sjöberg, "Learning when to drive in intersections by combining reinforcement learning and model predictive control," in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, pp. 3263–3268, 2019.
- [85] Y. Lin, J. McPhee, and N. L. Azad, "Anti-jerk on-ramp merging using deep reinforcement learning," in *Proc. IEEE Intell. Veh. Symposium (IV)*, pp. 7–14, 2020.
- [86] X. He, B. Lou, H. Yang, and C. Lv, "Robust decision making for autonomous vehicles at highway on-ramps: A constrained adversarial reinforcement learning approach," *IEEE Trans. Intell. Transp. Syst.*, vol. 24, no. 4, pp. 4103–4113, 2022.
- [87] C. Jiang, H. Liu, C. Qiu, S. Zhang, and W. Zhuang, "Ramp merging sequence and trajectory optimization for connected and autonomous vehicles using deep reinforcement learning," in *Proc. IEEE Int. Conf. Adv. Motion Control (AMC)*, pp. 1–7, 2024.
- [88] B. M. Albaba and Y. Yildiz, "Driver modeling through deep reinforcement learning and behavioral game theory," *IEEE Transactions on Control Systems Technology*, vol. 30, no. 2, pp. 885–892, 2021.
- [89] D. Chen, M. R. Hajidavalloo, Z. Li, K. Chen, Y. Wang, L. Jiang, and Y. Wang, "Deep multi-agent reinforcement learning for highway on-ramp merging in mixed traffic," *IEEE Trans. Intell. Transp. Syst.*, vol. 24, no. 11, pp. 11623–11638, 2023.
- [90] X. Zhang, L. Wu, H. Liu, Y. Wang, H. Li, and B. Xu, "High-speed ramp merging behavior decision for autonomous vehicles based on multiagent reinforcement learning," *IEEE Internet of Things Journal*, vol. 10, no. 24, pp. 22664–22672, 2023.
- [91] E. Leurent and J. Mercat, "Social attention for autonomous decision-making in dense traffic," *arXiv:1911.12250*, 2019.
- [92] W. Xiao, Y. Yang, X. Mu, Y. Xie, X. Tang, D. Cao, and T. Liu, "Decision-making for autonomous vehicles in random task scenarios at unsignalized intersection using deep reinforcement learning," *IEEE Trans. Veh. Technol.*, vol. 73, no. 6, pp. 7812–7825, 2024.
- [93] Z. Qiao, K. Muelling, J. Dolan, P. Palanisamy, and P. Mudalige, "Pomdp and hierarchical options mdp with continuous actions for autonomous driving at intersections," in *Proc. Intell. Transp. Syst. Conf. (ITSC)*, pp. 2377–2382, 2018.
- [94] J. Li, L. Sun, J. Chen, M. Tomizuka, and W. Zhan, "A safe hierarchical planning framework for complex driving scenarios based on reinforcement learning," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pp. 2660–2666, 2021.
- [95] Y. Ren, J. Jiang, G. Zhan, S. E. Li, C. Chen, K. Li, and J. Duan, "Self-learned intelligence for integrated decision and control of automated vehicles at signalized intersections," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 12, pp. 24145–24156, 2022.
- [96] X. He, J. Wu, Z. Huang, Z. Hu, J. Wang, A. Sangiovanni-Vincentelli, and C. Lv, "Fear-neuro-inspired reinforcement learning for safe autonomous driving," *IEEE Trans. Pattern Anal. Mach. Intell.*, 2023.
- [97] S. Wang, Z. Wang, R. Jiang, R. Yan, and L. Du, "Trajectory jerking suppression for mixed traffic flow at a signalized intersection: A trajectory prediction based deep reinforcement learning method," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 10, pp. 18989–19000, 2022.
- [98] Y. Liu, Y. Gao, Q. Zhang, D. Ding, and D. Zhao, "Multi-task safe reinforcement learning for navigating intersections in dense traffic," *J. Frankl. Inst.*, vol. 360, no. 17, pp. 13737–13760, 2023.
- [99] J. Jiang, Y. Ren, Y. Guan, S. E. Li, Y. Yin, D. Yu, and X. Jin, "Integrated decision and control at multi-lane intersections with mixed traffic flow," in *Journal of Physics: Conference Series*, vol. 2234, p. 012015, 2022.
- [100] Z. Bai, P. Hao, W. Shangguan, B. Cai, and M. J. Barth, "Hybrid reinforcement learning-based eco-driving strategy for connected and automated vehicles at signalized intersections," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 9, pp. 15850–15863, 2022.
- [101] G.-P. Antonio and C. Maria-Dolores, "Multi-agent deep reinforcement learning to manage connected autonomous vehicles at tomorrow's intersections," *IEEE Trans. Veh. Technol.*, vol. 71, no. 7, pp. 7033–7043, 2022.
- [102] R. Zhao, Y. Li, F. Gao, Z. Gao, and T. Zhang, "Multi-agent constrained policy optimization for conflict-free management of connected autonomous vehicles at unsignalized intersections," *IEEE Trans. Intell. Transp. Syst.*, vol. 25, no. 6, pp. 5374–5388, 2024.
- [103] J. Hu, Y. Feng, S. Li, H. Wang, J. So, and J. Zheng, "Mirroring the parking target: An optimal-control-based parking motion planner with strengthened parking reliability and faster parking completion," *IEEE Trans. Intell. Transp. Syst.*, vol. 25, no. 11, pp. 16157–16170, 2024.
- [104] J. Shi, K. Li, C. Piao, J. Gao, and L. Chen, "Model-based predictive control and reinforcement learning for planning vehicle-parking trajectories for vertical parking spaces," *Sensors*, vol. 23, no. 16, 2023.
- [105] J. Chen, B. Yuan, and M. Tomizuka, "Model-free deep reinforcement learning for urban autonomous driving," in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, pp. 2765–2771, 2019.
- [106] L. Anzalone, P. Barra, S. Barra, A. Castiglione, and M. Nappi, "An end-to-end curriculum learning approach for autonomous driving

- scenarios,” *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 10, pp. 19817–19826, 2022.
- [107] G. Zhan, Y. Jiang, S. E. Li, Y. Lyu, X. Zhang, and Y. Yin, “A transformation-aggregation framework for state representation of autonomous driving systems,” *IEEE Trans. Intell. Transp. Syst.*, vol. 25, no. 7, pp. 7311–7322, 2024.
- [108] Y.-L. Jin, Z.-Y. Ji, D. Zeng, and X.-P. Zhang, “Vwp:an efficient drl-based autonomous driving model,” *IEEE Transactions on Multimedia*, vol. 26, pp. 2096–2108, 2024.
- [109] Y. Hu, J. Yang, L. Chen, K. Li, C. Sima, X. Zhu, S. Chai, S. Du, T. Lin, W. Wang, L. Lu, X. Jia, Q. Liu, J. Dai, Y. Qiao, and H. Li, “Planning-oriented autonomous driving,” *arXiv:2212.10156*, 2023.
- [110] Y. Song, H. Lin, E. Kaufmann, P. Dürr, and D. Scaramuzza, “Autonomous overtaking in gran turismo sport using curriculum reinforcement learning,” in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pp. 9403–9409, 2021.
- [111] R. Banerjee, P. Ray, and M. Campbell, “Improving environment robustness of deep reinforcement learning approaches for autonomous racing using bayesian optimization-based curriculum learning,” *arXiv:2312.10557*, 2023.
- [112] S. Huang, X. Wu, and G. Huang, “Deep reinforcement learning-based multi-objective path planning on the off-road terrain environment for ground vehicles,” *arXiv:2305.13783*, 2023.
- [113] J. Zhao, Y. Wang, Y. Zhang, M. Wu, and R. Li, “A multi-objective deep reinforcement learning method for path planning in shovel loading scenario,” in *Proc. IEEE Int. Conf. Unmanned Syst.*, pp. 913–918, 2023.
- [114] Y. Zhang and C. Li, “On hierarchical path planning based on deep reinforcement learning in off- road environments,” in *Proc. Int. Conf. Autom., Robot. Applications (ICARA)*, pp. 461–465, 2024.
- [115] S. J. Wang, H. Zhu, and A. M. Johnson, “Pay attention to how you drive: Safe and adaptive model-based reinforcement learning for off-road driving,” in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pp. 16954–16960, 2024.
- [116] J. Shi, T. Zhang, Z. Zong, S. Chen, J. Xin, and N. Zheng, “Task-driven autonomous driving: Balanced strategies integrating curriculum reinforcement learning and residual policy,” *IEEE Robot. Autom. Lett.*, 2024.
- [117] J. Shi, T. Zhang, J. Zhan, S. Chen, J. Xin, and N. Zheng, “Efficient lane-changing behavior planning via reinforcement learning with imitation learning initialization,” in *Proc. IEEE Intell. Veh. Symposium (IV)*, pp. 1–8, 2023.
- [118] K. Yuan, Y. Huang, S. Yang, Z. Zhou, Y. Wang, D. Cao, and H. Chen, “Evolutionary decision-making and planning for autonomous driving based on safe and rational exploration and exploitation,” *Engineering*, vol. 33, pp. 108–120, 2024.
- [119] H. Tian, K. Reddy, Y. Feng, M. Qudus, Y. Demiris, and P. Angeloudis, “Enhancing autonomous vehicle training with language model integration and critical scenario generation,” *arXiv:2404.08570*, 2024.
- [120] K. B. Naveed, Z. Qiao, and J. M. Dolan, “Trajectory planning for autonomous vehicles using hierarchical reinforcement learning,” in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, pp. 601–606, 2021.
- [121] H. Shu, T. Liu, X. Mu, and D. Cao, “Driving tasks transfer using deep reinforcement learning for decision-making of autonomous vehicles in unsignalized intersection,” *IEEE Trans. Veh. Technol.*, vol. 71, no. 1, pp. 41–52, 2022.
- [122] G. Jin, Z. Li, B. Leng, W. Han, and L. Xiong, “Stability enhanced hierarchical reinforcement learning for autonomous driving with parameterized trajectory action,” in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, 2024.
- [123] L. Chen, Y. He, Q. Wang, W. Pan, and Z. Ming, “Joint optimization of sensing, decision-making and motion-controlling for autonomous vehicles: A deep reinforcement learning approach,” *IEEE Trans. Veh. Technol.*, vol. 71, no. 5, pp. 4642–4654, 2022.
- [124] J. Xiong, Q. Wang, Z. Yang, P. Sun, L. Han, Y. Zheng, H. Fu, T. Zhang, J. Liu, and H. Liu, “Parameterized deep q-networks learning: Reinforcement learning with discrete-continuous hybrid action space,” *arXiv:1810.06394*, 2018.
- [125] K. Yuan, Y. Huang, S. Yang, M. Wu, D. Cao, Q. Chen, and H. Chen, “Evolutionary decision-making and planning for autonomous driving: A hybrid augmented intelligence framework,” *IEEE Trans. Intell. Transp. Syst.*, vol. 25, no. 7, pp. 7339–7351, 2024.
- [126] Z. Li, G. Jin, R. Yu, B. Leng, and L. Xiong, “Interaction-aware deep reinforcement learning approach based on hybrid parameterized action space for autonomous driving,” in *Proc. SAE Intell. Connected Veh. Symposium (SAE ICVS)*, 2024.
- [127] W. B. Knox, A. Allievi, H. Banzhaf, F. Schmitt, and P. Stone, “Reward (mis)design for autonomous driving,” *Artificial Intelligence*, vol. 316, p. 103829, 2023.
- [128] G. Li, Y. Yang, S. Li, X. Qu, N. Lyu, and S. E. Li, “Decision making of autonomous vehicles in lane change scenarios: Deep reinforcement learning approaches with risk awareness,” *Transp. Res. Part C Emerg. Technol.*, vol. 134, p. 103452, 2022.
- [129] M. Zhu, D. Piga, and A. Bemporad, “C-GLISp: Preference-based global optimization under unknown constraints with applications to controller calibration,” *IEEE Trans. Control Syst. Tech.*, vol. 30, pp. 2176–2187, Sept. 2022.
- [130] Z. Wu, L. Sun, W. Zhan, C. Yang, and M. Tomizuka, “Efficient sampling-based maximum entropy inverse reinforcement learning with application to autonomous driving,” *IEEE Robot. Autom. Lett.*, vol. 5, no. 4, pp. 5355–5362, 2020.
- [131] X. Wen, S. Jian, and D. He, “Modeling the effects of autonomous vehicles on human driver car-following behaviors using inverse reinforcement learning,” *IEEE Trans. Intell. Transp. Syst.*, vol. 24, no. 12, pp. 13903–13915, 2023.
- [132] S. Mysore, G. Cheng, Y. Zhao, K. Saenko, and M. Wu, “Multi-critic actor learning: Teaching RL policies to act with style,” in *Proc. Int. Conf. Learn. Representations (ICLR)*, 2022.
- [133] W. Yuan, M. Yang, Y. He, C. Wang, and B. Wang, “Multi-reward architecture based reinforcement learning for highway driving policies,” in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, pp. 3810–3815, 2019.
- [134] R. T. Icarte, T. Q. Klassen, R. Valenzano, and S. A. McIlraith, “Reward machines: Exploiting reward function structure in reinforcement learning,” *J. Artif. Intell. Res.*, vol. 73, pp. 173–208, 2022.
- [135] S. J. Russell and P. Norvig, *Artificial intelligence: a modern approach*. Pearson, 2016.
- [136] L.-C. Wu, Z. Zhang, S. Haesaert, Z. Ma, and Z. Sun, “Risk-aware reward shaping of reinforcement learning agents for autonomous driving,” in *Proc. Annu. Conf. IEEE Ind. Electronics Soc.*, pp. 1–6, 2023.
- [137] H. Ma, C. Liu, S. E. Li, S. Zheng, W. Sun, and J. Chen, “Learn zero-constraint-violation safe policy in model-free constrained reinforcement learning,” *IEEE Trans. Neural Netw. Learn. Syst.*, vol. early access, 2024.
- [138] X. Chen, B. Xu, M. Hu, Y. Bian, Y. Li, and X. Xu, “Safe efficient policy optimization algorithm for unsignalized intersection navigation,” *IEEE/CAA J. Autom. Sin.*, vol. 11, no. 9, pp. 2011–2026, 2024.
- [139] L. Wen, J. Duan, S. E. Li, S. Xu, and H. Peng, “Safe reinforcement learning for autonomous vehicles through parallel constrained policy optimization,” in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, pp. 1–7, 2020.
- [140] T. J. Perkins and A. G. Barto, “Lyapunov design for safe reinforcement learning,” *Journal of Machine Learning Research*, vol. 3, no. Dec, pp. 803–832, 2002.
- [141] A. D. Ames, S. Coogan, M. Egerstedt, G. Notomista, K. Sreenath, and P. Tabuada, “Control barrier functions: Theory and applications,” in *2019 18th European control conference (ECC)*, pp. 3420–3431, 2019.
- [142] L. Zhang, R. Zhang, T. Wu, R. Weng, M. Han, and Y. Zhao, “Safe reinforcement learning with stability guarantee for motion planning of autonomous vehicles,” *IEEE Trans. Neural Netw. Learn. Sys.*, vol. 32, no. 12, pp. 5435–5444, 2021.
- [143] S. Udatha, Y. Lyu, and J. Dolan, “Reinforcement learning with probabilistically safe control barrier functions for ramp merging,” in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pp. 5625–5630, 2023.
- [144] Y. Yang, Y. Jiang, Y. Liu, J. Chen, and S. E. Li, “Model-free safe reinforcement learning through neural barrier certificate,” *IEEE Robot. Autom. Lett.*, vol. 8, no. 3, pp. 1295–1302, 2023.
- [145] S. Nagesh Rao *et al.*, “Robust ai driving strategy for autonomous vehicles,” in *AI-enabled Technologies for Autonomous and Connected Vehicles*, pp. 161–212, Springer, 2022.
- [146] Z. Zhang, Q. Liu, Y. Li, K. Lin, and L. Li, “Safe reinforcement learning in autonomous driving with epistemic uncertainty estimation,” *IEEE Trans. Intell. Transp. Syst.*, 2024.
- [147] K. Stachowicz and S. Levine, “Racer: Epistemic risk-sensitive rl enables fast driving with fewer crashes,” *arXiv:2405.04714*, 2024.
- [148] A. Baheri, S. Nagesh Rao, H. E. Tseng, I. Kolmanovsky, A. Girard, and D. Filev, “Deep reinforcement learning with enhanced safety for autonomous highway driving,” in *Proc. IEEE Intell. Veh. Symposium (IV)*, pp. 1550–1555, 2020.
- [149] S. Gu, G. Chen, L. Zhang, J. Hou, Y. Hu, and A. Knoll, “Constrained reinforcement learning for vehicle motion planning with topological reachability analysis,” *Robotics*, vol. 11, no. 4, p. 81, 2022.

- [150] X. Wang, "Ensuring safety of learning-based motion planners using control barrier functions," *IEEE Robot. Autom. Lett.*, vol. 7, no. 2, pp. 4773–4780, 2022.
- [151] B. Gangopadhyay, H. Soora, and P. Dasgupta, "Hierarchical program-triggered reinforcement learning agents for automated driving," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 8, pp. 10902–10911, 2021.
- [152] C.-J. Hoel, K. Wolff, and L. Laine, "Ensemble quantile networks: Uncertainty-aware reinforcement learning with applications in autonomous driving," *IEEE Trans. Intell. Transp. Syst.*, vol. 24, no. 6, pp. 6030–6041, 2023.
- [153] K. Yang, X. Tang, S. Qiu, S. Jin, Z. Wei, and H. Wang, "Towards robust decision-making for autonomous driving on highway," *IEEE Trans. Veh. Technol.*, vol. 72, no. 9, pp. 11251–11263, 2023.
- [154] H. Krasowski, Y. Zhang, and M. Althoff, "Safe reinforcement learning for urban driving using invariably safe braking sets," in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, pp. 2407–2414, 2022.
- [155] N. Li, Y. Li, and I. Kolmanovsky, "A unified safety protection and extension governor," *IEEE Transactions on Automatic Control*, 2024.
- [156] Y. Li, N. Li, H. E. Tseng, A. Girard, D. Filev, and I. Kolmanovsky, "Safe reinforcement learning using robust action governor," in *Learning for Dynamics and Control*, pp. 1093–1104, PMLR, 2021.
- [157] B. Tearle, K. P. Wabersich, A. Carron, and M. N. Zeilinger, "A predictive safety filter for learning-based racing control," *IEEE Robot. Autom. Lett.*, vol. 6, no. 4, pp. 7635–7642, 2021.
- [158] S. Alighanbari and N. L. Azad, "Deep reinforcement learning with nmpc assistance nash switching for urban autonomous driving," *IEEE Trans. Intell. Veh.*, vol. 8, no. 3, pp. 2604–2615, 2022.
- [159] J. Lu, G. Alcan, and V. Kyrki, "Integrating expert guidance for efficient learning of safe overtaking in autonomous driving using deep reinforcement learning," *arXiv:2308.09456*, 2023.
- [160] T. Hester, M. Vecerik, O. Pietquin, M. Lanctot, T. Schaul, B. Piot, D. Horgan, J. Quan, A. Sendonaris, I. Osband, *et al.*, "Deep q-learning from demonstrations," in *AAAI Conf. Artif. Intell.*, vol. 32, 2018.
- [161] H. Liu, Z. Huang, J. Wu, and C. Lv, "Improved deep reinforcement learning with expert demonstrations for urban autonomous driving," in *Proc. IEEE Intell. Veh. Symposium (IV)*, pp. 921–928, 2022.
- [162] Y. Gao, H. Xu, J. Lin, F. Yu, S. Levine, and T. Darrell, "Reinforcement learning from imperfect demonstrations," *arXiv:1802.05313*, 2018.
- [163] J. Wu, Z. Huang, W. Huang, and C. Lv, "Prioritized experience-based reinforcement learning with human guidance for autonomous driving," *IEEE Trans. Neural Netw. Learn. Sys.*, vol. 35, no. 1, pp. 855–869, 2024.
- [164] I. Uchendu, T. Xiao, Y. Lu, B. Zhu, M. Yan, J. Simon, M. Bennice, C. Fu, C. Ma, J. Jiao, *et al.*, "Jump-start reinforcement learning," in *Proc. Int. Conf. Mach. Learn. (ICML)*, pp. 34556–34583, 2023.
- [165] Z. Huang, J. Wu, and C. Lv, "Efficient deep reinforcement learning with imitative expert priors for autonomous driving," *IEEE Trans. Neural Netw. Learn. Sys.*, vol. 34, no. 10, pp. 7391–7403, 2023.
- [166] J. Li, X. Liu, B. Zhu, J. Jiao, M. Tomizuka, C. Tang, and W. Zhan, "Guided online distillation: Promoting safe reinforcement learning by offline demonstration," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pp. 7447–7454, 2024.
- [167] T. Ota, "Decision mamba: Reinforcement learning via sequence modeling with selective state spaces," *arXiv:2403.19925*, 2024.
- [168] M. Guo and M. Bürger, "Geometric task networks: Learning efficient and explainable skill coordination for object manipulation," *IEEE Transactions on Robotics*, vol. 38, no. 3, pp. 1723–1734, 2022.
- [169] X. Wang, Y. Chen, and W. Zhu, "A survey on curriculum learning," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 44, no. 9, pp. 4555–4576, 2022.
- [170] J. Shi, T. Zhang, Z. Zong, S. Chen, J. Xin, and N. Zheng, "Task-driven autonomous driving: Balanced strategies integrating curriculum reinforcement learning and residual policy," *IEEE Robot. Autom. Lett.*, vol. 9, no. 11, pp. 9454–9461, 2024.
- [171] H. Niu, Y. Xu, X. Jiang, and J. Hu, "Continual driving policy optimization with closed-loop individualized curricula," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pp. 6850–6857, 2024.
- [172] H. Wang, H. Zhang, L. Li, Z. Kan, and Y. Song, "Task-driven reinforcement learning with action primitives for long-horizon manipulation skills," *IEEE Trans. Cybern.*, vol. 54, no. 8, pp. 4513–4526, 2024.
- [173] F. L. Da Silva and A. H. R. Costa, "A survey on transfer learning for multiagent reinforcement learning systems," *J. Artif. Intell. Res.*, vol. 64, pp. 645–703, 2019.
- [174] Z. Yan and C. Wu, "Reinforcement learning for mixed autonomy intersections," in *Proc. IEEE Intell. Transp. Syst. Conf. (ITSC)*, pp. 2089–2094, 2021.
- [175] X. Xiao, B. Liu, G. Warnell, and P. Stone, "Motion planning and control for mobile robot navigation using machine learning: a survey," *Autonomous Robots*, vol. 46, no. 5, pp. 569–597, 2022.
- [176] Y. Xia, S. Liu, Q. Yu, L. Deng, Y. Zhang, H. Su, and K. Zheng, "Parameterized decision-making with multi-modality perception for autonomous driving," in *Proc. IEEE Int. Conf. Data Eng. (ICDE)*, pp. 4463–4476, 2024.
- [177] M. Dalal, D. Pathak, and R. R. Salakhutdinov, "Accelerating robotic reinforcement learning via parameterized action primitives," in *Proc. Adv. Neural Inf. Proces. Syst.*, vol. 34, pp. 21847–21859, 2021.
- [178] K. Lee, M. Laskin, A. Srinivas, and P. Abbeel, "Sunrise: A simple unified framework for ensemble learning in deep reinforcement learning," in *Proc. Int. Conf. Mach. Learn. (ICML)*, pp. 6131–6141, 2021.
- [179] T. Kanazawa, H. Wang, and C. Gupta, "Distributional actor-critic ensemble for uncertainty-aware continuous control," in *Proc. Int. Joint Conf. Neural Networks*, pp. 1–10, 2022.
- [180] Z. Ma, X. Liu, and Y. Huang, "Unsupervised reinforcement learning for multi-task autonomous driving: Expanding skills and cultivating curiosity," *IEEE Trans. Intell. Transp. Syst.*, vol. 25, no. 10, pp. 14209–14219, 2024.
- [181] Y. Wu, S. Liao, X. Liu, Z. Li, and R. Lu, "Deep reinforcement learning on autonomous driving policy with auxiliary critic network," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 34, no. 7, pp. 3680–3690, 2023.
- [182] Y. Lu *et al.*, "Imitation is not enough: Robustifying imitation with reinforcement learning for challenging driving scenarios," in *Proc. IEEE/RSJ Int. Conf. Intell. Rob. Syst. (IROS)*, pp. 7553–7560, 2023.
- [183] M. Everett, B. Lütjens, and J. P. How, "Certifiable robustness to adversarial state uncertainty in deep reinforcement learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 33, no. 9, pp. 4184–4198, 2022.
- [184] A. Balakrishnan, J. Lee, A. Gaurav, K. Czarnecki, and S. Sedwards, "Transfer reinforcement learning for autonomous driving: From wise-move to wisemim," *ACM Trans. Model. Comput. Simul. (TOMACS)*, vol. 31, no. 3, pp. 1–26, 2021.
- [185] K. L. Voogd, J. P. Allamaa, J. Alonso-Mora, and T. D. Son, "Reinforcement learning from simulation to real world autonomous driving using digital twin," *IFAC-PapersOnLine*, vol. 56, no. 2, pp. 1510–1515, 2023.
- [186] C. Gao, Y. Zheng, W. Wang, F. Feng, X. He, and Y. Li, "Causal inference in recommender systems: A survey and future directions," *ACM Transactions on Information Systems*, vol. 42, no. 4, pp. 1–32, 2024.
- [187] F. Shoeleh and M. Asadpour, "Skill based transfer learning with domain adaptation for continuous reinforcement learning domains," *Applied Intelligence*, vol. 50, no. 2, pp. 502–518, 2020.
- [188] J. Beck, R. Vuorio, E. Z. Liu, Z. Xiong, L. Zintgraf, C. Finn, and S. Whiteson, "A survey of meta-reinforcement learning," *arXiv:2301.08028*, 2023.
- [189] F. Ye, P. Wang, C.-Y. Chan, and J. Zhang, "Meta reinforcement learning-based lane change strategy for autonomous vehicles," in *Proc. IEEE Intell. Veh. Symposium (IV)*, pp. 223–230, 2021.
- [190] Q. Deng, R. Li, Q. Hu, Y. Zhao, and R. Li, "Context-aware meta-rl with two-stage constrained adaptation for urban driving," *IEEE Trans. Veh. Technol.*, 2023.
- [191] D. Abel, A. Barreto, B. Van Roy, D. Precup, H. P. van Hasselt, and S. Singh, "A definition of continual reinforcement learning," *Proc. Adv. Neural Inf. Proces. Syst.*, vol. 36, 2024.
- [192] D. Wei, J. Xing, S. Yang, Y. Lu, and Y. Huang, "Continual reinforcement learning for autonomous driving with application on velocity control under various environment," in *Proc. CAA Int. Conf. Veh. Control Intell. (CVCI)*, pp. 1–8, 2023.
- [193] Z. Cao *et al.*, "Autonomous driving policy continual learning with one-shot disengagement case," *IEEE Trans. Intell. Veh.*, vol. 8, no. 2, pp. 1380–1391, 2022.
- [194] X. He and C. Lv, "Towards safe autonomous driving: Decision making with observation-robust reinforcement learning," *Automotive Innovation*, vol. 6, no. 4, pp. 509–520, 2023.
- [195] B. Lütjens, M. Everett, and J. P. How, "Safe reinforcement learning with model uncertainty estimates," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pp. 8662–8668, 2019.
- [196] J. Choi, C. Dance, J.-E. Kim, S. Hwang, and K.-s. Park, "Risk-conditioned distributional soft actor-critic for risk-sensitive navigation," in *Proc. IEEE Int. Conf. Robot. Autom. (ICRA)*, pp. 8337–8344, 2021.
- [197] S. Li *et al.*, "Learning locomotion for quadruped robots via distributional ensemble actor-critic," *IEEE Robot. Autom. Lett.*, vol. 9, no. 2, pp. 1811–1818, 2024.

- [198] J. Duan, D. Yu, S. E. Li, W. Wang, Y. Ren, Z. Lin, and B. Cheng, "Fixed-dimensional and permutation invariant state representation of autonomous driving," *IEEE Trans. Intell. Transp. Syst.*, vol. 23, no. 7, pp. 9518–9528, 2021.
- [199] W. X. Zhao *et al.*, "A survey of large language models," *arXiv:2303.18223*, 2023.
- [200] C. Cui *et al.*, "A survey on multimodal large language models for autonomous driving," in *Proc. IEEE/CVF Winter Conf. Appl. Comput. Vision Workshops*, pp. 958–979, 2024.
- [201] B. Wang, Y. Qu, Y. Jiang, J. Shao, C. Liu, W. Yang, and X. Ji, "Llm-empowered state representation for reinforcement learning," *arXiv:2407.13237*, 2024.
- [202] W. Zhao, T. He, and C. Liu, "Probabilistic safeguard for reinforcement learning using safety index guided gaussian process models," in *Learning for Dynamics and Control Conference*, pp. 783–796, 2023.
- [203] K.-C. Hsu, H. Hu, and J. F. Fisac, "The safety filter: A unified view of safety-critical control in autonomous systems," *Annu. Rev. Control Robot. Auton. Syst.*, vol. 7, 2023.
- [204] G. Swamy, C. Dann, R. Kidambi, Z. S. Wu, and A. Agarwal, "A minimalist approach to reinforcement learning from human feedback," *arXiv:2401.04056*, 2024.
- [205] E. Kargar and V. Kyrki, "Increasing the efficiency of policy learning for autonomous vehicles by multi-task representation learning," *IEEE Trans. Intell. Veh.*, vol. 7, no. 3, pp. 701–710, 2022.
- [206] Z. Yang, Y. Chen, J. Wang, S. Manivasagam, W.-C. Ma, A. J. Yang, and R. Urtasun, "Unisim: A neural closed-loop sensor simulator," in *Proc. IEEE/CVF Conf. Comput. Vision and Pattern Recognit.*, pp. 1389–1399, 2023.
- [207] Z. Deng, J. Jiang, G. Long, and C. Zhang, "Causal reinforcement learning: A survey," *arXiv:2307.01452*, 2023.
- [208] Y. Cao and Y. Li, "Survey on large language model-enhanced reinforcement learning: Concept, taxonomy, and methods," *arXiv:2404.00282*, 2024.
- [209] A. Keysan, A. Look, E. Kosman, G. Gürsun, J. Wagner, Y. Yao, and B. Rakitsch, "Can you text what is happening? integrating pre-trained language encoders into trajectory prediction models for autonomous driving," *arXiv:2309.05282*, 2023.
- [210] L. Wen *et al.*, "Dilu: A knowledge-driven approach to autonomous driving with large language models," *arXiv:2309.16292*, 2023.
- [211] M. Yildirim, B. Dagda, and S. Fallah, "Highwayllm: Decision-making and navigation in highway driving with rl-informed language model," *arXiv:2405.13547*, 2024.