# A Survey of Deep RL and IL for Autonomous Driving Policy Learning

Zeyu Zhu, *Member, IEEE,* Huijing Zhao, *Senior Member, IEEE,*

*Abstract*—Autonomous driving (AD) agents generate driving policies based on online perception results, which are obtained at multiple levels of abstraction, e.g., behavior planning, motion planning and control. Driving policies are crucial to the realization of safe, efficient and harmonious driving behaviors, where AD agents still face substantial challenges in complex scenarios. Due to their successful application in fields such as robotics and video games, the use of deep reinforcement learning (DRL) and deep imitation learning (DIL) techniques to derive AD policies have witnessed vast research efforts in recent years. This paper is a comprehensive survey of this body of work, which is conducted at three levels: First, a taxonomy of the literature studies is constructed from the system perspective, among which five modes of integration of DRL/DIL models into an AD architecture are identified. Second, the formulations of DRL/DIL models for conducting specified AD tasks are comprehensively reviewed, where various designs on the model state and action spaces and the reinforcement learning rewards are covered. Finally, an in-depth review is conducted on how the critical issues of AD applications regarding driving safety, interaction with other traffic participants and uncertainty of the environment are addressed by the DRL/DIL models. To the best of our knowledge, this is the first survey to focus on AD policy learning using DRL/DIL, which is addressed simultaneously from the system, task-driven and problem-driven perspectives. We share and discuss findings, which may lead to the investigation of various topics in the future.

*Index Terms*—deep reinforcement learning, deep imitation learning, autonomous driving policy

## I. INTRODUCTION

**A**UTONOUMOUS driving (AD) has received extensive attention in recent decades [1–4] and could be a promising solution for improving road safety [5], traffic flow [6] and fuel economy [7], among other factors. A typical architecture of an AD system is illustrated in Fig. 1, which is composed of perception, planning and control modules. An AD agent generates **driving policies** based on online perception results, which are obtained at multiple levels of abstraction, e.g., behavior planning, motion planning and control. The earliest autonomous vehicles can be dated back to [8–10]. One milestone was the Defense Advanced Research Projects Agency (DARPA) Grand Challenges [11, 12]. Recent years have witnessed a huge boost in AD research, and many products and prototyping systems have been developed. Despite the fast development of the field, AD still faces substantial challenges in complex scenarios for the realization of safe, efficient

① Sequences of route points through road network ② Behavior decision: lane change, car follow, stop, etc.
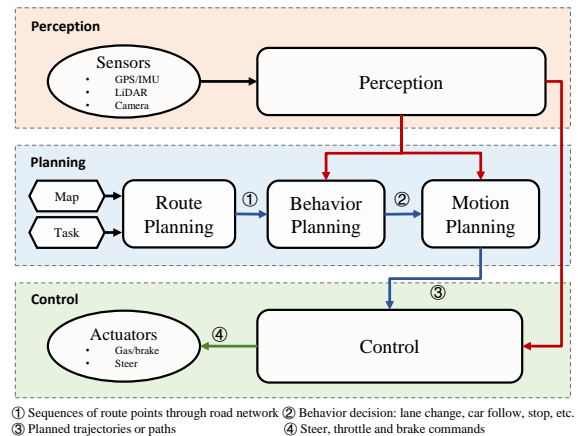③ Planned trajectories or paths ④ Steer, throttle and brake commands

Fig. 1. Architecture of autonomous driving systems. A general abstraction that is based on [13–16].

and harmonious driving behaviors [17, 18]. **Reinforcement learning** (**RL**) is a principled mathematical framework for solving sequential decision making problems [19–21]. **Imitation learning** (**IL**), which is closely related, refers to learning from expert demonstrations. However, the early methods of both were limited to relatively low-dimensional problems. The rise of **deep learning** (**DL**) techniques [22, 23] in recent years has provided powerful solutions to this problem through the appealing properties of **deep neural networks** (**DNNs**): function approximation and representation learning. DL techniques enable the scaling of RL/IL to previously intractable problems (e.g., high-dimensional state spaces), which have increased in popularity for complex locomotion [24], robotics [25] and autonomous driving [26–28] tasks. Unless otherwise stated, this survey focuses on **Deep RL** (**DRL**) and **Deep IL** (**DIL**).

A large variety of DRL/DIL models have been developed for learning AD policies, which are reviewed in this paper. Several surveys are relevant to this study. [13, 15] survey the motion planning and control methods of automated vehicles before the era of DL. [29–33] review general DRL/DIL methods without considering any particular applications. [4] addresses the deep learning techniques for AD with a focus on perception and control, while [34] addresses control only. [35] provides a taxonomy of AD tasks to which DRL models have been applied and highlights the key challenges. However, none of these studies answers the following questions:

*How can DRL/DIL models be integrated into AD systems from the perspective of system architecture? How can they be formulated to accomplish specified AD tasks? How can methods be designed that address the challenging issues of AD, such as safety, interaction with other traffic participants,*

*and uncertainty of the environment?*

This study seeks answers to the above questions. To the best of our knowledge, this is the first survey to focus on AD policy learning using DRL/DIL, which is addressed from the system, task-driven and problem-driven perspectives. Our contributions are threefold:

- A taxonomy of the literature is presented from the system perspective, from which five modes of integration of DRL/DIL models into an AD architecture are identified. The studies on each mode are reviewed, and the architectures are compared. It is found that the vast research efforts have focused mainly on exploring the potential of DRL/DIL in accomplishing AD tasks, while intensive studies are needed on the optimization of the architectures of DRL/DIL embedded systems toward real-world applications.

- The formulations of DRL/DIL models for accomplishing specified AD tasks are comprehensively reviewed, where various designs on the model state and action spaces and the reinforcement learning rewards are covered. It is found that these formulations rely heavily on empirical designs, which are brute-force approaches and lack theoretic support. Changing the designs or tuning the parameters could result in substantially different driving policies, which may pose large challenges to the AD system's stability and robustness in real-world deployment.

- The critical issues of AD applications that are addressed by the DRL/DIL models regarding driving safety, interaction with other traffic participants and uncertainty of the environment are comprehensively discussed. It is found that driving safety has been well studied. A typical strategy is to combine with traditional methods to ensure a DRL/DIL agent's functional safety; however, striking a balance between optimal policies and hard constraints remains non-trivial. The studies on the latter two issues are still highly preliminary, in which the problems have been addressed from divergent perspectives and the studies have not been conducted systematically.

The remainder of this paper is organized as follows: Sections II and III briefly introduce (D)RL and (D)IL, respectively. Section IV reviews the research on DRL/DIL in AD from a system architecture perspective. Section V reviews the task-driven methods and examines the formulations of DRL/DIL models for completing specified AD tasks. Section VI reviews problem-driven methods, in which specified autonomous vehicle problems and challenges are addressed. Section VII discusses the remaining challenges, and the conclusions of the survey are presented in Section VIII.

## II. PRELIMINARIES OF REINFORCEMENT LEARNING

### A. Problem Formulation

Reinforcement learning (RL) is a principled mathematical framework that is based upon the paradigm of trial-and-error learning, where an agent interacts with its environment through a trade-off between exploitation and exploration [36–38]. **Markov decision processes (MDPs)** are a mathematically idealized form of RL [21], which are represented as

$(\mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma)$, where $\mathcal{S}$ and $\mathcal{A}$ denote the sets of states and actions, respectively, and $\mathcal{P}(s_{t+1}|s_t, a_t) : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \to [0, 1]$ is the **transition/dynamics function** that maps state-action pairs onto a distribution of next-step states. The numerical immediate reward function $\mathcal{R}(s_t, a_t, s_{t+1}) : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ serves as a learning signal. A discount factor $\gamma \in [0, 1]$ determines the present value of future rewards (lower values encourage more myopic behaviors). MDPs' states satisfy the **Markov property** [21], namely, future states depend only on the immediately preceding states and actions. **Partially observable MDPs (POMDPs)** extend MDPs to problems in which access to fully observable Markov property states is not available. A POMDP has an observation set $\Omega$ and an observation function $\mathcal{O}$, where $\mathcal{O}(a_t, s_{t+1}, o_{t+1}) = p(o_{t+1}|a_t, s_{t+1})$ is the probability of observing $o_{t+1}$ given that the agent has executed action $a_t$ and reached state $s_{t+1}$ [39]. For the theory and algorithms of POMDPs, we refer readers to [39, 40].

The MDP agent selects an action $a_t \in \mathcal{A}$ at each time step $t$ based on the current state $s_t$, receives a numerical reward $r_{t+1}$ and visits a new state $s_{t+1}$. The generated sequence $\{s_0, a_0, r_1, s_1, a_1, r_2, ...\}$ is called a **rollout** or **trajectory**. The expected cumulative reward in the future, namely, the **expected discounted return** $G_t$ after time step $t$, is defined as [21]:

$$G_t \doteq r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + ... = \sum_{k=0}^{T} \gamma^k r_{t+k+1} \quad (1)$$

where $T$ is a finite value or $\infty$ for finite and infinite horizon problems, respectively. The **policy** $\pi(a|s)$ maps states to probabilities of selecting each possible action. The **value function** $v_\pi(s)$ is the expected return following $\pi$ from state $s$:

$$v_\pi(s) \doteq \mathbb{E}_\pi[G_t|s_t = s] \quad (2)$$

Similarly, the **action-value function** $q_\pi(s, a)$ is defined as:

$$q_\pi(s, a) \doteq \mathbb{E}_\pi[G_t|s_t = s, a_t = a] \quad (3)$$

which satisfies the recursive **Bellman equation** [41] :

$$q_\pi(s_t, a_t) = \mathbb{E}_{s_{t+1}}[r_{t+1} + \gamma q_\pi(s_{t+1}, \pi(s_{t+1}))] \quad (4)$$

The objective of RL is to identify the optimal policy that maximizes the expected return $\pi^* = \arg\max_\pi \mathbb{E}_\pi[G_t|s_t = s]$. The methods can be divided into three groups, as shown in Fig. 2 (a).

### B. Value-based Methods

To solve a reinforcement learning problem, one can identify an optimal action-value function and recover the optimal policy from the learned state-action values.

$$q_{\pi^*}(s, a) = \max_\pi q_\pi(s, a) \quad (5)$$

$$\pi^*(s) = \arg\max_a q_{\pi^*}(s, a) \quad (6)$$

Q-learning [42] is one of the most popular methods, which estimates Q values through temporal difference (**TD**):

$$q_\pi(s_t, a_t) \leftarrow q_\pi(s_t, a_t) + \alpha(Y - q_\pi(s_t, a_t)) \quad (7)$$
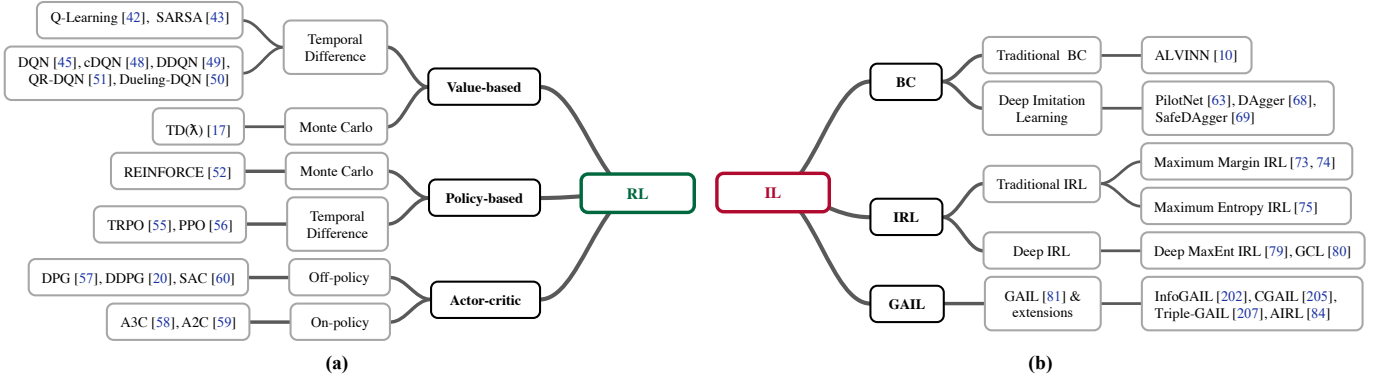
Fig. 2. A taxonomy of the general methods of reinforcement learning (RL) and imitation learning (IL)

where $Y = r_{t+1} + \gamma \max_{a_{t+1} \in \mathcal{A}} q_\pi(s_{t+1}, a_{t+1})$ is the temporal difference target and $\alpha$ is the learning rate. This can be regarded as a standard regression problem in which the error to be minimized is $Y - q_\pi(s_t, a_t)$. Q-learning is **off-policy** since it updates $q_\pi$ based on experiences that are not necessarily generated by the derived policy, while SARSA [43] is an **on-policy** algorithm that uses the derived policy to generate experiences. Another distinction is that SARSA uses target $Y = r_{t+1} + \gamma q_\pi(s_{t+1}, a_{t+1})$. In contrast to TD methods, Monte Carlo methods estimate the expected return through averaging the results of multiple rollouts, which can be applied to non-Markovian episodic environments. TD and Monte Carlo have been further combined in TD($\lambda$) [21].

The early methods [21, 42, 43] of this group rely on tabular representations. A major problem is the "curse of dimensionality" [44], namely, an increase in the number of state features would result in exponential growth of the number of state-action pairs that must be stored. Recent methods use DNNs to approximate a parameterized value function $q(s, a; \omega)$, of which Deep Q-networks (DQNs) [45] are the most representative, which learn the values by minimizing the following loss function:

$$J(\omega) = \mathbb{E}_t[(Y - q(s_t, a_t; \omega))^2] \qquad (8)$$

where $Y = r_{t+1} + \gamma \max q(s_{t+1}, a_{t+1}; \omega^-)$ is the target, $\omega^-$ denotes the parameters of the target network, and $\omega$ can be learnt based on the gradients.

$$\omega \leftarrow \omega - \alpha \mathbb{E}_t[(Y - q(s_t, a_t; \omega)) \nabla q(s_t, a_t; \omega)] \qquad (9)$$

The major contributions of DQN are the introduction of the target network and experience replay. To avoid rapidly fluctuating Q values and stabilize the training, the target network is fixed for a specified number of iterations during the update of the primary Q-network and subsequently updated to match the primary Q-network. Moreover, experience replay [46], which maintains a memory that stores transitions $(s_t, a_t, s_{t+1}, r_{t+1})$, increases the sample efficiency. A later study improves the uniform sample experience replay by introducing priority [47]. Continuous DQN (cDQN) [48] derives a continuous variant of DQN based on normalized advantage functions (NAFs). Double DQN (DDQN) [49] addresses the overestimation problem of DQN through the use of a double estimator. Dueling-DQN [50] introduces a dueling architecture in which

both the value function and associated advantage function are estimated. QR-DQN [51] utilizes distributional reinforcement learning to learn the full value distribution rather than only the expectation.

### C. Policy-based Methods

Alternatively, one can directly search and optimize a parameterized policy $\pi_\theta$ to maximize the expected return:

$$\max_\theta J(\theta) = \max_\theta v_{\pi_\theta}(s) = \max_\theta \mathbb{E}_{\pi_\theta}[G_t | s_t = s] \qquad (10)$$

where $\theta$ denotes the policy parameters, which can be optimized based on the policy gradient theorem [21]:

$$\begin{aligned} \nabla J(\theta) &\propto \sum_s \mu(s) \sum_a q_\pi(s, a) \nabla \pi_\theta(a|s) \\ &= \mathbb{E}_\pi[\sum_a q_\pi(s_t, a) \nabla \pi_\theta(a|s_t)] \\ &= \mathbb{E}_\pi[G_t \nabla \ln \pi_\theta(a_t|s_t)] \end{aligned} \qquad (11)$$

where $\mu(s)$ denotes the state visitation frequency. REINFORCE [52] is a straightforward Monte Carlo policy-based method that selects the discounted return $G_t$ following the policy $\pi_\theta$ to estimate the policy gradient in Eqn.11. The parameters are updated as follows [21]:

$$\theta \leftarrow \theta + \alpha G_t \nabla \ln \pi_\theta(a_t|s_t) \qquad (12)$$

This update intuitively increases the log probability of actions that lead to higher returns. Since empirical returns are used, the resulting gradients suffer from high variances. A common technique for reducing the variance and accelerating the learning is to replace $G_t$ in Eqn. 11 and 12 by $G_t - b(s_t)$ [21, 52], where $b(s_t)$ is a baseline. Alternatively, $G_t$ can be replaced by the advantage function [53, 54] $A_{\pi_\theta}(s, a) = q_{\pi_\theta(s,a)} - v_{\pi_\theta}(s)$.

One problem of policy-based methods is poor gradient updates may result in newly updated policies that deviate wildly from previous policies, which may decrease the policy performance. Trust region policy optimization (TRPO) [55] solves this problem through optimization of a surrogate objective function, which guarantees the monotonic improvement of policy performance. Each policy gradient update is constrained by using a quadratic approximation of the Kullback-Leibler (KL) divergence between policies. Proximal policy optimization (PPO) [56] improved upon TRPO through the application

of an adaptive penalty on the KL divergence and a heuristic clipped surrogate objective function. The requirement for only a first-order gradient also renders PPO easier to implement than TRPO.

### D. Actor-Critic Methods

Actor-critic methods have the advantages of both value-based and policy-based methods, where a neural network **actor** $\pi_\theta(a|s)$ selects actions and a neural network **critic** $q(s, a; \omega)$ or $v(s; \omega)$ estimates the values. The actor and critic are typically updated alternately according to Eqn. 11 and Eqn. 8, respectively. Deterministic policy gradient (DPG) [57] is an off-policy actor-critic algorithm that derives deterministic policies. Compared with stochastic policies, DPG only integrates over the state space and requires fewer samples in problems with large action spaces. Deep deterministic policy gradient (DDPG) [24] utilizes DNNs to operate on high-dimensional state spaces with experience replay and a separate actor-critic target network, which is similar to DQN. Exploitation of parallel computation is an alternative to experience replay. Asynchronous advantage actor-critic (A3C) [58] uses advantage estimates rather than discounted returns in the actor-critic framework and asynchronously updates policy and value networks on multiple parallel threads of the environment.

The parallel independent environments stabilize the learning process and enable more exploration. Advantage actor critic (A2C) [59], which is the synchronous version of A3C, uses a single agent for simplicity or waits for each agent to finish its experience to collect multiple trajectories. Soft actor critic (SAC) [60] benefits from the addition of an entropy term to the reward function to encourage better exploration.

## III. PRELIMINARIES OF IMITATION LEARNING

### A. Problem Formulation

Imitation learning possesses a simpler form and is based on learning from demonstrations (LfD) [61]. It is attractive for AD applications, where interaction with the real environment could be dangerous and vast amount of human driving data can be easily collected [62]. A demonstration dataset $\mathcal{D} = \{\xi_i\}_{i=0..N}$ contains a set of trajectories, where each trajectory $\xi_i = \{(s_t^i, a_t^i)\}_{t=0..T}$ is a sequence of state-action pairs, and action $a_t^i \in \mathcal{A}$ is taken by expert at state $s_t^i \in S$ under the guidance of an underlying policy $\pi_E$ of the expert [32]. Using the collected dataset $\mathcal{D}$, a common optimization-based strategy of imitation learning is to learn a policy $\pi^* : \mathcal{S} \to \mathcal{A}$ that mimics the expert's policy $\pi_E$ by satisfying [33]

$$\pi^* = \arg\min_\pi \mathbb{D}(\pi_E, \pi) \tag{13}$$

where $\mathbb{D}$ is a similarity measure between policies. The methods for solving the problem can be divided into three groups, as shown in Fig. 2 (b), which are reviewed below.

### B. Behavior Clone

Behavior clone (**BC**) formulates the problem as a supervised learning process with the objective of matching the learned policy $\pi_\theta$ to the expert policy $\pi_E$:

$$\min_\theta \mathbb{E}||\pi_\theta - \pi_E||_2 \tag{14}$$

which is typically realized by minimizing the L2-loss:

$$J(\theta) = \mathbb{E}_{(s,a)\sim\mathcal{D}}[(\pi_\theta(s) - a)^2] \tag{15}$$

Early research on imitation learning can be dated back to the ALVINN system [10], which used a 3-layer neural network to perform road following based on front camera images. In the most recent decade, deep imitation learning (DIL) has been conducted using DNNs as policy function approximators and has realized success in end-to-end AD systems [63–65]. BC performs well for states that are covered by the training distribution but generalizes poorly to new states due to compounding errors in the actions, which is also referred to as **covariate shift** [66, 67]. To overcome this problem, Ross et al. [68] proposed DAgger, which improves upon supervised learning by using a primary policy to collect training examples while running a reference policy simultaneously. In each iteration, states that are visited by the primary policy are also sent to the reference policy to output expert actions, and the newly generated demonstrations are aggregated into the training dataset. SafeDAgger [69] extends on DAgger by introducing a safe policy that learns to predict the error that is made by a primary policy without querying the reference policy. In addition to dataset aggregation, data augmentation techniques such as the addition of random noise to the expert action have also been commonly used in DIL [70, 71].

### C. Inverse Reinforcement Learning

The inverse reinforcement learning problem, which was first formulized in the study of Ng et al. [72], is to identify a reward function $r_\theta$ for which the expert behavior is optimal:

$$\max_\theta \mathbb{E}_{\pi_E}[G_t|r_\theta] - \mathbb{E}_\pi[G_t|r_\theta] \tag{16}$$

Early studies utilized linear function approximation of reward functions and identified the optimal reward via maximum margin approaches [73, 74]. By introducing the maximum entropy principle, Ziebart et al. [75] eliminated the reward ambiguity between demonstrations and expert policy where multiple rewards may explain the expert behavior. The reward function is learned through maximizing the posterior probability of observing expert trajectories:

$$J(\theta) = \mathbb{E}_{\xi_i\sim\mathcal{D}}[\log P(\xi_i|r_\theta)] \tag{17}$$

where the probability of a trajectory satisfies $P(\xi_i|r_\theta) \propto \exp(r_\theta(\xi_i))$. Several studies have extended the reward functions to nonlinear formulations through Gaussian processes [76] or boosting [77, 78]. However, the above methods generally operate on low-dimensional features. The use of rich and expressive function approximators, in the form of neural networks, has been proposed to learn reward functions directly on raw high-dimensional state representations [79, 80].

A problem that is encountered with IRL is that to evaluate the reward function, a forward reinforcement learning process must be conducted to obtain the corresponding policy, thereby rendering IRL inefficient and computationally expensive. Many early approaches require solving an MDP in the inner loop of each iterative optimization [20, 72, 74, 75]. Perfect knowledge of the system dynamics and an efficient

offline solver are needed in these methods, thereby limiting their applications in complex real-world scenarios such as robotic control. Finn et al. [80] proposed guided cost learning (GCL), which handles unknown dynamics in high-dimensional complex systems and learns complex neural network cost functions through an efficient sample-based approximation.

### D. Generative Adversarial Imitation Learning

Generative adversarial imitation learning (GAIL) [81] directly learns a policy from expert demonstrations while requiring neither the reward design in RL nor the expensive RL process in the inner loop of IRL. GAIL establishes an analogy between imitation learning and generative adversarial networks (GANs) [82]. The generator $\pi_\theta$ serves as a policy for imitating expert behavior by matching the state-action distribution of demonstrations, while the discriminator $D_\omega \in (0, 1)$ serves as a surrogate reward function for measuring the similarity between generated and demonstration samples. The objective function of GAIL is formulated in the following min-max form:

$$\min_{\pi_\theta} \max_{D_\omega} \mathbb{E}_{\pi_\theta}[\log D_\omega(s, a)] + \mathbb{E}_{\pi_E}[\log(1 - D_\omega(s, a))] - \lambda H(\pi_\theta) \tag{18}$$

where $H(\pi)$ is a regularization entropy term. The generator and the discriminator are updated with the following gradients:

$$\nabla_\theta J(\theta) = \mathbb{E}_\pi[\nabla_\theta \log \pi_\theta(a|s) Q(s, a)] - \lambda \nabla_\theta H(\pi_\theta)$$
$$\nabla_\omega J(\omega) = \mathbb{E}_\pi[\nabla_\omega \log D_\omega(s, a)] + \mathbb{E}_{\pi_E}[\nabla_\omega \log(1 - D_\omega(s, a))] \tag{19}$$

Finn et al. [83] mathematically proved the connection among GANs, IRL and energy-based models. Fu et al.[84] proposed adversarial inverse reinforcement learning (AIRL) based on an adversarial reward learning formulation, which can recover reward functions that are robust to dynamics changes.

## IV. ARCHITECTURES OF DRL/DIL INTEGRATED AD SYSTEMS

AD systems have been studied for decades [13–16], which are commonly composed of modular pipelines, as illustrated in Fig. 1. How can DRL/DIL models be integrated into an AD system and collaborate with other modules? This section reviews the literature from the system architecture perspective, from which five modes are identified, as illustrated in Fig. 3. A classification of the studies in each mode is presented in Table I, along with the exploited DRL/DIL methods, the upstream module for perception, the targeted AD tasks, and the advantages and disadvantages of the architectures, among other information. Below, we detail each mode of the architectures, which is followed by a comparison of the number of studies that correspond to each mode or to the use of DRL or DIL methods.

### A. Mode 1. DRL/DIL Integrated Control

Many studies have applied DRL/DIL to control, which can be abstracted as the architecture of Mode 1 and is illustrated in Fig. 3. Bojarski et al. [63] proposed a well-known end-to-end self-driving control framework. They trained a nine-layer CNN by supervised learning to learn the steering policy without explicit manual decomposition of sequential modules. However, their model only adapts to lane keeping. An alternative option is to feed traditional perception results into the DNN control module. Tang et al. [96] proposed the use of environmental rasterization encoding, along with the relative distance, velocity and acceleration, as input to a two-branch neural network, which was trained via proximal policy optimization.

Although Mode 1 features a simple structure and adapts to a large variety of learning methods, it is limited to specified tasks; thus, it has difficulty addressing scenarios in which multiple driving skills that are conditioned on various situations are needed. Moreover, bypassing and ignoring behavior planning or motion planning processes may weaken the model's interpretability and performance in complex tasks (e.g., urban navigation).

### B. Mode 2. Extension of Mode 1 with High-level Command

As illustrated in Fig. 3, Mode 2 extends Mode 1 by considering the high-level planning output. The control module is composed of either a general model for all behaviors [104, 110] or a series of models for distinct behaviors [26, 70, 105, 106, 108]. Chen et al. [110] projected detected environment vehicles and the routing onto a bird-view road map, which served as the input of a policy network. Liang et al. [26] built on conditional imitation learning (**CIL**) [70] and proposed the branched actor network, as illustrated in Fig. 3. These methods learn several control submodules for distinct behaviors. A gating control command from high-level route planning and behavior planning modules is responsible for the selection of the corresponding control submodule.

Although Mode 2 mitigates the problems that are encountered with Mode 1, it has its own limitations. A general model may be not sufficient for capturing diverse behaviors. However, learning a model for each behavior increases the demand for training data. Moreover, Mode 2 may be not as computationally efficient as Mode 1 since it requires high-level planning modules that are determined in advance to guide the control module.

### C. Mode 3. DRL/DIL Integrated Motion Planning

Mode 3 integrates DRL/DIL into the motion planning module, and its architecture is illustrated in Fig. 3. Utilizing the planning output (e.g., routes and driving decisions) from high-level modules, along with the current perception output, DNNs are trained to predict future trajectories or paths. DIL models [71, 111–113] are the mainstream choices for implementing this architecture. As illustrated in Fig. 3, Sun et al. [112] proposed training a neural network that imitates long-term MPC via Sampled-DAgger, where the policy network's input was from perception (obstacles, environment, and current states) and decision-making (driving decisions). Alternatively, Wulfmeier et al. [113] proposed projecting the LiDAR point cloud onto a grid map, which is sent to the DNN. The DNN is responsible for learning a cost function that guides the
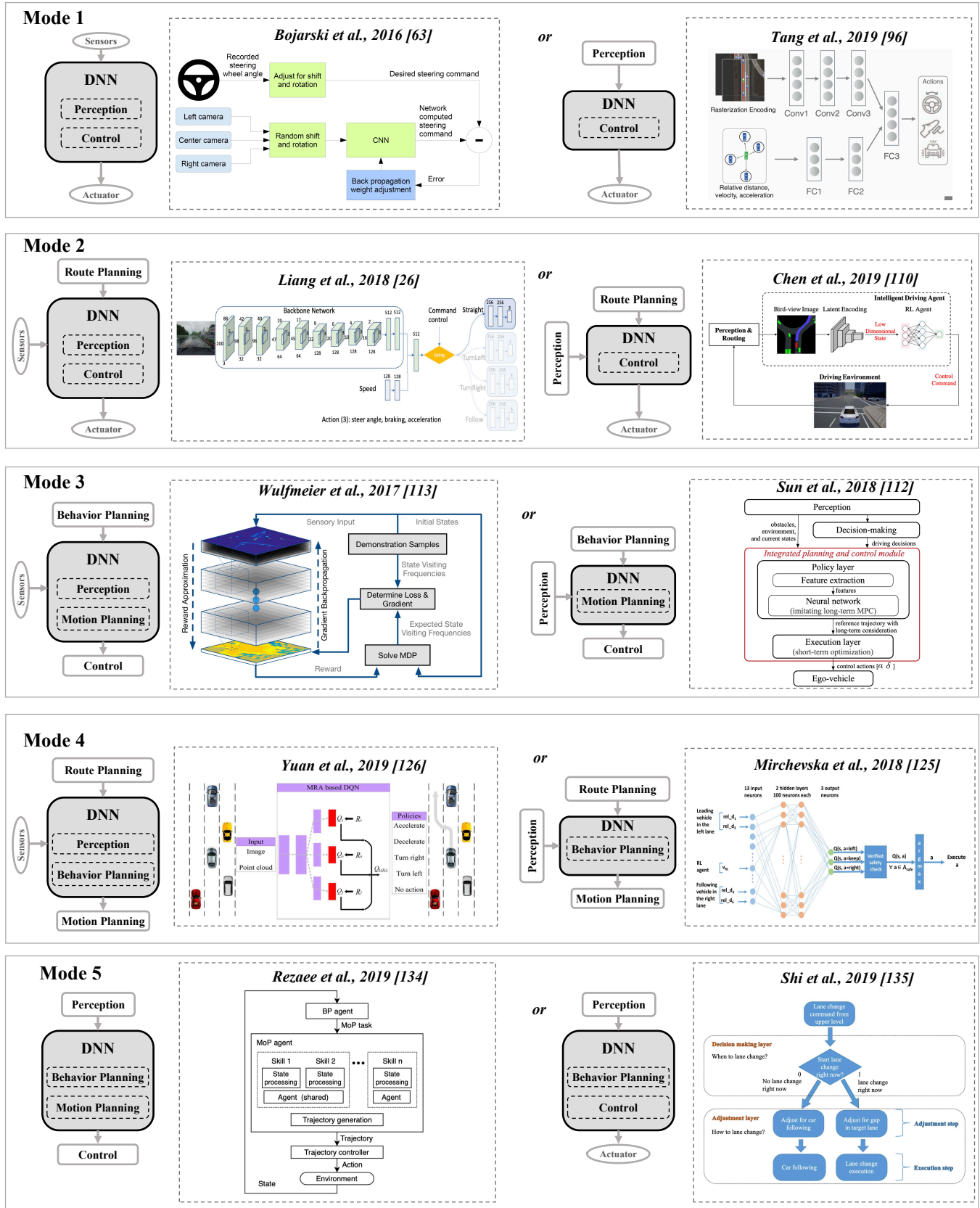
Fig. 3. Integration modes of DRL/DIL models into AD architectures

TABLE I
A TAXONOMY OF THE LITERATURE BY THE INTEGRATION MODES OF **DRL** OR **DIL** MODELS INTO AD ARCHITECTURES

| Architecture | Advantages | Disadvantages | Studies | Methods | Perception[1] | Tasks[2] | Remarks[3] |
|---|---|---|---|---|---|---|---|
| **Mode 1. DRL/DIL Integrated Control** | - Features a simple structure and adapts to various learning methods. | - Limited to specified tasks or skills. - Bypassing planning modules weakens the model's interpretability and capability. | [10], [85], [63], [64], [86], [87], [65], [88] | BC | D | road/lane following urban driving | safety: [69] |
| | | | [69], [89] | DAgger | D | road/lane following | |
| | | | [90], [91], [92], [93], [94], [95] | NQL DQN | T | lane changing traffic merging imminent events intersection | safety: [90] interaction: [95] |
| | | | [96], [97] | PPO | T | traffic merging urban driving | |
| | | | [28], [98] | DDPG | D | road/lane following imminent events | |
| | | | [99], [100], [101], [102], [103] | DDPG | T | road/lane following lane changing overtaking imminent events | |
| **Mode 2. Extension of Mode 1 with High-level Command** | - Considers both high-level planning and perception - Generates distinct control behaviors according to the high-level decisions | - Single model may not capture sufficiently diverse behaviors - Learning a model for each behavior increases the training cost and the demand for data | [104] | BC | D | urban driving | |
| | | | [105], [70], [106] | CIL | D | urban navigation | |
| | | | [107], [108] | UAIL | D | urban navigation | uncertainty: [107], [108] |
| | | | [109] | DDPG | T | parking | |
| | | | [26] | DDPG | D | urban navigation | |
| | | | [110] | DDQN/TD3/SAC | T | roundabout | |
| **Mode 3. DRL/DIL Integrated Motion Planning** | - Learn to imitate human trajectories - Efficient forward prediction process | - No guarantee on safety or feasibility | [111], [71] | BC | T | urban driving | safety: [71] |
| | | | [112] | DAgger | T | highway driving | |
| | | | [113] | MaxEnt DIRL | D | urban driving | |
| | | | [114] | DDQN/DQfD | T | valet parking | safety: [114] |
| | | | [115] | SAC | T | traffic merge | |
| **Mode 4. DRL/DIL Integrated Behavior Planning** | - The DNNs need only plan high-level behavioral actions. | - Simple and few actions limit the control precision and diversity of driving styles. - Complicated and too many actions increase the training cost. | [116] | AIRL | T | lane change | |
| | | | [117] | IRL | D | lane change | |
| | | | [118], [119], [120], [121], [122], [123], [124], [125] | DQN | T | lane change intersection | |
| | | | [126], [127] | DQN | D | lane change | |
| | | | [128], [129] | Q-Learning | T | lane changing | |
| | | | [130], [131] | DDQN | T | lane changing urban driving | |
| | | | [132] | DQfD | T | lane changing | safety: [132] |
| | | | [133] | QR-DQN | D | highway driving | |
| **Mode 5. DRL/DIL Integrated Hierarchical Planning and Control** | - Simultaneously plan at various levels of abstraction. | - Hierarchical DNNs increase the training cost and decrease the convergence speed. | [134], [135] | DQN | T | cruise control lane changing | |
| | | | [136] | Hierarchical Policy gradient | T | traffic light passing | |
| | | | [27] | DDPG | D | lane changing | |
| | | | [137] | DDPG | T | intersection | |
| | | | [138] | DDPG | T | urban driving | |

[1] Type of upstream perception module: (D)eep learning method/(T)raditional method
[2] For detailed information about AD tasks, see Table III
[3] Studies that also address safety VI-A, interaction VI-B and uncertainty VI-C problems are labelled.

motion planning. The control part in Mode 3 typically utilizes traditional control techniques such as PID [71] or MPC [112].

One major disadvantage of Mode 3 is that although DNN planned trajectories can imitate human trajectories, their safety and feasibility cannot be guaranteed.

### D. Mode 4. DRL/DIL Integrated Behavior Planning

Many studies focus on integrating DRL/DIL into the behavior planning module and deriving high-level driving policies. The corresponding architecture is presented in Fig. 3, where DNNs decide behavioral actions and the subsequent motion planning and control modules typically utilize traditional

methods. Many studies build upon DQN and its variants [118–126]. As illustrated in Fig. 3, Yuan et al. [126] decomposed the action space into five typical behavioral decisions on a highway. Compared with DRL, studies that employ DIL to learn high-level policies are limited. A recent study by Wang et al. [116] proposed the use of augmented adversarial inverse reinforcement learning (AIRL) to learn human-like decision-making on highways, where the action space consists of all possible combinations of lateral and longitudinal decisions.

In Mode 4, the design of behavioral actions is nontrivial, and one must balance the training cost and the diversity of driving

TABLE II
COMPARISON OF THE LITERATURE BY DRL/DIL INTEGRATION MODES

| Perception | Control (Modes 1&2) | | Motion Planning (Mode 3) | | Behavior Planning (Mode 4) | | Hierarchical P. & C. [1] (Mode 5) | |
|---|---|---|---|---|---|---|---|---|
| | DRL | DIL | DRL | DIL | DRL | DIL | DRL | DIL |
| **Traditional** | **15** | 0 | **2** | **3** | **13** | 1 | **5** | 0 |
| **DNN** | 3 | **16** | 0 | 1 | 3 | 1 | 1 | 0 |
| **Subtotal** | **18 (52.9%)** | 16 (47.1%) | 2 (33.3%) | **4 (66.7%)** | **16 (88.9%)** | 2 (11.1%) | **6 (100%)** | 0 (0%) |
| **Total** | **34 (53.1%)** | | 6 (9.4%) | | 18 (28.1%) | | 6 (9.4%) | |

[*] The values in this table are the numbers and percentages of papers in Table I that belong to the corresponding categories.
[1] Abbreviation for "Hierarchical Planning and Control"

TABLE III
A TAXONOMY OF THE LITERATURE BY SCENARIOS AND AD TASKS

| Scenario | AD Task | Description | Ref. DRL Methods | Ref. DIL Methods |
|---|---|---|---|---|
| **Urban** | **Intersection** | Learn to drive through intersections (while interacting and negotiating with other traffic participants). | DQN [90, 119–121] cDQN [138] DDPG [137, 138] | — |
| | **Roundabout** | Learn to drive through roundabouts (while interacting and negotiating with other traffic participants). | DDQN [110] TD3 [110] SAC [110] | Horizon GAIL[139] |
| | **Urban Navigation** | Learn to drive in urban environments to reach specified objectives by following global routes. | DDPG [26] | BC [104] CIL [70, 105, 106] UAIL [107, 108] |
| | **Urban Driving** | Learn to drive in urban environments without specified objectives. | DDQN [122] PPO [140] Policy gradient [136] | BC [65, 111] DIRL [113] |
| **Highway** | **Lane Change (LC)** | Learn to decide and perform lane changes. | DQN [91, 118] DDQN [141] DDPG [27] | Projection IRL [117] AIRL [116] |
| | **Lane Keep (LK)** | Learn to drive while maintaining the current lane. | DQN [102] DDPG [28] | BC [86, 87] SafeDAgger [69] |
| | **Cruise Control** | Learn a policy for adaptive cruise control. | NQL [92] DQN [134] Policy gradient [142] Actor-critic [143, 144] | — |
| | **Traffic Merging** | Learn to merge into highways. | DQN [93] PPO [96] | — |
| | **Highway Driving** | Learn a general policy for driving on a highway, which may include multiple behaviors such as LC and LK. | DQN [126, 145] QR-DQN [133] | GAIL [146] PS-GAIL [147] MaxEnt IRL [148] |
| **Others** | **Road Following** | Learn to simply follow one road. | — | BC [10, 63, 64, 85] DAgger [89] |
| | **Imminent Events** | Learn to avoid or mitigate imminent events such as collisions. | DQN [94] DDPG [98] | RAIL [149] |

styles. Simple and few behavioral actions limits the control precision and diversity of driving styles, whereas many and sophisticated actions increases the training cost.

### E. Mode 5. DRL/DIL Integrated Hierarchical Planning and Control

As illustrated in Fig. 3, Mode 5 blurs the lines between planning and control, where a single DNN outputs hierarchical actions [27] based on the parameterized action space [151] or hierarchical DNNs output actions at multiple levels [134–136]. Rezaee et al. [134] proposed an architecture, which is illustrated in Fig. 3, in which BP (behavior planning) is used to make high-level decisions regarding transitions between discrete states and MoP (motion planning) generates a target trajectory according to the decisions that are made via BP. Qiao et al. [137] built upon hierarchical MDP (HMDP) and realized their model through an options framework. In their implementation, the hierarchical options were modeled as high-level decisions (SlowForward or Forward). Based on high-level decisions and current observations, low-level control policies were derived.

Mode 5 simultaneously plans at multiple levels, and the low-level planning process considers high-level decisions. However, the use of hierarchical DNNs results in the increase of the training cost and potentially the decrease of the convergence speed since one poorly trained high-level DNN may mislead and disturb the learning process of low-level DNNs.
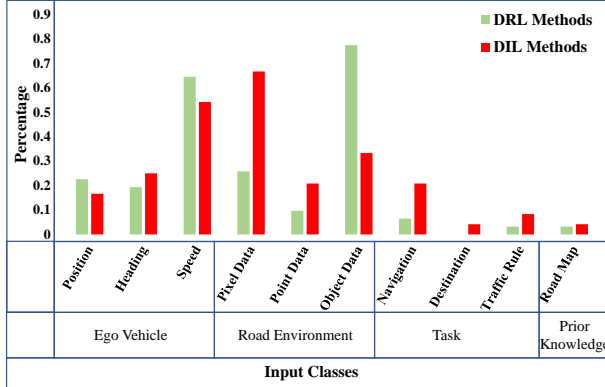
### F. Statistical Comparison

A statistical comparison among modes in terms of the number of studies is presented in Table II. Current studies on architectures are premature, unsystematic and unbalanced:

- Most studies focus on integrating DRL/DIL into control (Mode1&2), followed by behavior planning (Mode 4).
- For Mode 1&2, DRL methods mainly adopt traditional method-based perception, while DNN-based perception is preferred in DIL methods.
- DRL seems to be more popular for high-level decision making (Mode 4&5), while DIL is chosen more frequently for low-level control (Mode 1&2).
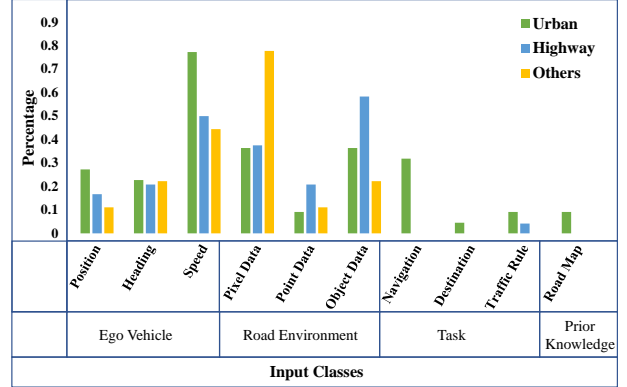
Future studies may address the imbalance problem and identify potential new architectures.

TABLE IV
INPUTS OF THE DRL/DIL METHODS FOR AD TASKS

| Information Source | Class | Inputs | Ref. | |
|---|---|---|---|---|
| | | | DRL Methods | DIL Methods |
| Ego Vehicle | Position Information | ego position | [90, 91, 93, 110, 120, 122, 136] | [111, 146, 147, 149] |
| | Heading Information | heading angle, orientation, steering, yaw, and yaw rate | [28, 91, 93, 119, 121, 138] | [64, 65, 139, 146, 147, 149] |
| | Speed Information | speed/velocity and acceleration | [26, 90, 91, 118–122, 136–138, 141] [28, 93, 94, 98, 102, 134, 144, 150] | [65, 70, 105–108, 117, 139] [89, 146–149] |
| Road Environment | Pixel Data | camera RGB images | [26–28, 126, 133, 145] | [65, 69, 70, 86, 87, 104–108] [10, 19, 63, 64, 85, 89] |
| | | semantically segmented images | [98] | — |
| | | 2D bird's-eye-view images | [96] | — |
| | Point Data | LiDAR sensor readings | [126, 133, 145] | [139, 146, 147, 149] |
| | | 2D LiDAR grid map | — | [113] |
| | Object Data | other road users' information: relative speed, position, and distance to ego | [110, 118–122, 138, 141] [92–94, 96, 134, 142–144] | [111, 116, 147, 149] |
| | | lane/road information: ego vehicle's distance to lane markings, road curvature, and lane width | [91, 93, 102, 118, 122, 134, 137, 141] | [116, 146, 147, 149] |
| Task | Navigation Information | navigational driving commands or planned routes | [26, 110] | [70, 104–106, 111] |
| | Destination Information | destination position, distance or angle to destination | — | [139] |
| | Traffic Rule Information | traffic lights' state, speed limit, and desired speed | [136] | [111, 148] |
| Prior Knowledge | Road Map | 2D top-down road map images | [110] | [111] |



Fig. 4. The percentages of the literature that uses certain data as input. (a) Comparison by DRL/DIL methods. (b) Comparison by scenarios.

## V. TASK-DRIVEN METHODS

DRL/DIL studies in AD can be categorized according to their application scenarios and targeted tasks, as listed in Table III. These task-driven studies use DRL/DIL to solve specified AD tasks, where the formulations of DRL/DIL can be decomposed into several key components: 1) *state space and input design*, 2) *action space and output design*, and 3) *reinforcement learning reward design*.

### A. State Space and Input Design

Table IV classifies commonly used inputs according to information source (ego vehicle/road environment/task/prior knowledge). A statistical comparison of the percentages of input classes that are used in DRL/DIL methods is presented in Fig. 4(a). Compared with the task and prior knowledge, the ego vehicle and road environment seem to be more popular information sources. In all input classes, the ego vehicle speed, pixel data (e.g., camera images) and object data (e.g., other road users' relative speeds and positions) are the most commonly used. Another significant difference between DRL and DIL is that DRL models prefer object data while DIL models prefer pixel data. The selection of low-dimensional object data rather than high-dimensional pixel data as input for DRL models renders the problem more tractable and accelerates the training procedure. Fig. 4(b) presents a statistical comparison of preferences for input classes in various scenarios. Input from the task (e.g., goal positions) and prior knowledge (e.g., road maps) are mostly used in urban scenarios. Point data (e.g., LiDAR sensor readings) and object data are more commonly used in highway scenarios.

Aside from deciding the choice of input, the application of a dynamic input size is an important problem since the number of cars or pedestrians in the ego's vicinity varies over time. Unfortunately, standard DRL/DIL methods rely on inputs of fixed size. The use of occupancy grid maps as inputs of CNNs is a practical solution [119, 121]. However, this solution imposes a trade-off between computational workload and expressiveness. Low-resolution grids decrease the computational burden at the cost of being imprecise in their representation of the environment, whereas for high-resolution grids, most computations may be redundant due to sparsity of the grid

TABLE V
ACTIONS OF THE DRL/DIL METHODS FOR AD TASKS

| Action Category | Subclass | Action Outputs | Ref | |
|---|---|---|---|---|
| | | | DRL Methods | DIL Methods |
| Behavior-level Actions | Acceleration Related | e.g., full brake, decelerate, continue, and accelerate | [119, 121, 122, 126, 133, 142, 145] | — |
| | Lane Change Related | e.g., keep, LLC, and RLC; choice of lane change gaps | [118, 122, 133, 141, 145] | [116] |
| | Turn Related | e.g., straight, left-turn, right-turn, and stop | [126] | [65, 117] |
| | Interaction Related | e.g., take/give way and follow a vehicle; wait/go | [120, 121] | — |
| Trajectory-level Actions | Planned Trajectory | future path 2D points | — | [111, 113] |
| Control-level Actions | Lateral | discrete steer angles | [19] | [10] |
| | | continuous steer | — | [63, 64, 85–87] |
| | | continuous angular speed or yaw acceleration | [91] | [65] |
| | Longitudinal | discrete acceleration values | [90, 94] | — |
| | | continuous acceleration | [92, 143, 144] | — |
| | | continuous brake and throttle | [150] | — |
| | Simultaneous Lateral & Longitudinal | continuous steer/turn-rate and speed/acceleration/throttle | [28, 93, 96, 110] | [70, 89, 104, 106, 107, 146, 147] |
| | | continuous steer, acceleration/throttle and brake | [26, 98, 102] | [105, 108] |
| | | continuous steer and binary brake decision | — | [69] |
| Hierarchical Actions | Behavior & Control | e.g., behavior-level pass/stop, control-level acceleration, and steer | [27, 136–138] | — |
| | Behavior & Trajectory | e.g., behavior-level maintenance, LLC, RLC and trajectory-level path points | [134] | — |

TABLE VI
REWARDS OF THE DRL METHODS FOR AD TASKS

| Category | Subclass | Description | Ref. |
|---|---|---|---|
| Safety | Avoid Collision | Impose penalties if a collision occurs | [26, 90, 105, 110, 116, 118–122, 137, 138] [69, 94, 96, 98, 126, 133, 143, 145, 149] |
| | Time to Collision | Impose penalties if the time to collision (TTC) is below a safe threshold | [118–120, 122, 144] |
| | Distance to Other Vehicles | Impose penalties if this distance is shorter than a safe threshold | [93, 134, 142, 144] |
| | Number of Lane Changes | Impose penalties if the number of lane changes is too large or reward a smaller number of lane changes | [118, 126, 133, 134, 141, 145] |
| | Out of Road | Impose penalties on driving out of road | [19, 27, 96, 102, 118] |
| Efficiency | Speed | Reward higher speed until the maximum speed limit is reached; Impose penalties if the speed is lower than the minimum speed limit | [26, 105, 110, 118, 119, 122, 136, 138, 140, 141] [27, 28, 93, 96, 102, 126, 133, 134, 145, 150] |
| | Success | Reward the agent if it finished the task successfully | [90, 96, 116, 118, 120, 121, 137, 138, 143] |
| | Number of Overtakes | Reward a higher number of overtakes for efficiency | [27, 126, 133, 145] |
| | Time | Impose a negative reward in each step to encourage the agent finish the task faster or penalize the agent if the task cannot be finished within a time threshold | [91, 110, 121, 137] |
| | Distance to the Destination | Provide a larger reward the closer the agent is to the destination | [105, 136, 137] |
| Comfort | Jerk | Impose penalties if the longitudinal or lateral control is too urgent | [91, 93, 96, 110, 120, 136, 138, 144, 149, 150] |
| Traffic Rules | Lane Mark Invasion | Impose penalties if the agent invades the lane marks | [26, 105, 149] |
| | Distance to the Lane Centerlines | Impose penalties if the agent deviates from the lane centerlines or routing baselines | [19, 27, 96, 110, 138, 140] |
| | Wrong Lane | Impose penalties if the agent is in the wrong lane, e.g., staying in the left-turn lane if the assigned route is straight | [122] |
| | Blocking Traffic | Impose penalties if the agent blocks the future paths of other vehicles that have the right of way | [137] |

maps. Furthermore, a grid map is still limited by its defined size, and agents outside this region is neglected. Alternatively, Everett et al. [152] proposed to leverage LSTMs' ability to encode a variable number of agents' observations. Huegle et al. [141] suggested the use of deep sets [153] as a flexible and permutation-invariant architecture to handle dynamic input. Dynamic input remains an open topic for future studies.

## B. Action Space and Output Design

A self-driving DRL/DIL agent can plan at different levels of abstraction, namely, low-level control, high-level behavioral planning and trajectory planning, or even at multiple levels simultaneously. According to this, Table V categorizes mainstream action spaces into four groups: behavior-level actions, trajectory-level actions, control-level actions and hierarchical actions. Behavior-level actions are usually designed according to specified tasks. For speed control, acceleration-related actions (e.g., full brake, decelerate, continue, and accelerate [119]) are commonly used. Lane change actions

(e.g., keep/LLC/LRC) and turn actions (e.g., turn left/right/go straight) are preferred in highway scenarios and urban scenarios, respectively. Trajectory-level actions refer to the planned or predicted trajectories/paths [111, 113], which are typically composed of future path 2D points. Control-level actions refer to low-level control commands (e.g., steer, acceleration, throttle, and brake), which are divided into three classes: lateral control, longitudinal control and simultaneous lateral and longitudinal control. Early studies focused mainly on discrete lateral [10, 19] or discrete longitudinal control [143], while continuous control was considered later [63, 64, 85, 91]. Continuous control is demonstrated in [102] to produce smoother trajectories than discrete control for lane keeping, which may make passengers feel more comfortable. Simultaneous lateral and longitudinal control has received wide attention, especially in urban scenarios [70, 105–108]. Recently, hierarchical actions have attracted more attention [27, 134, 136–138], which provide higher robustness and interpretability.
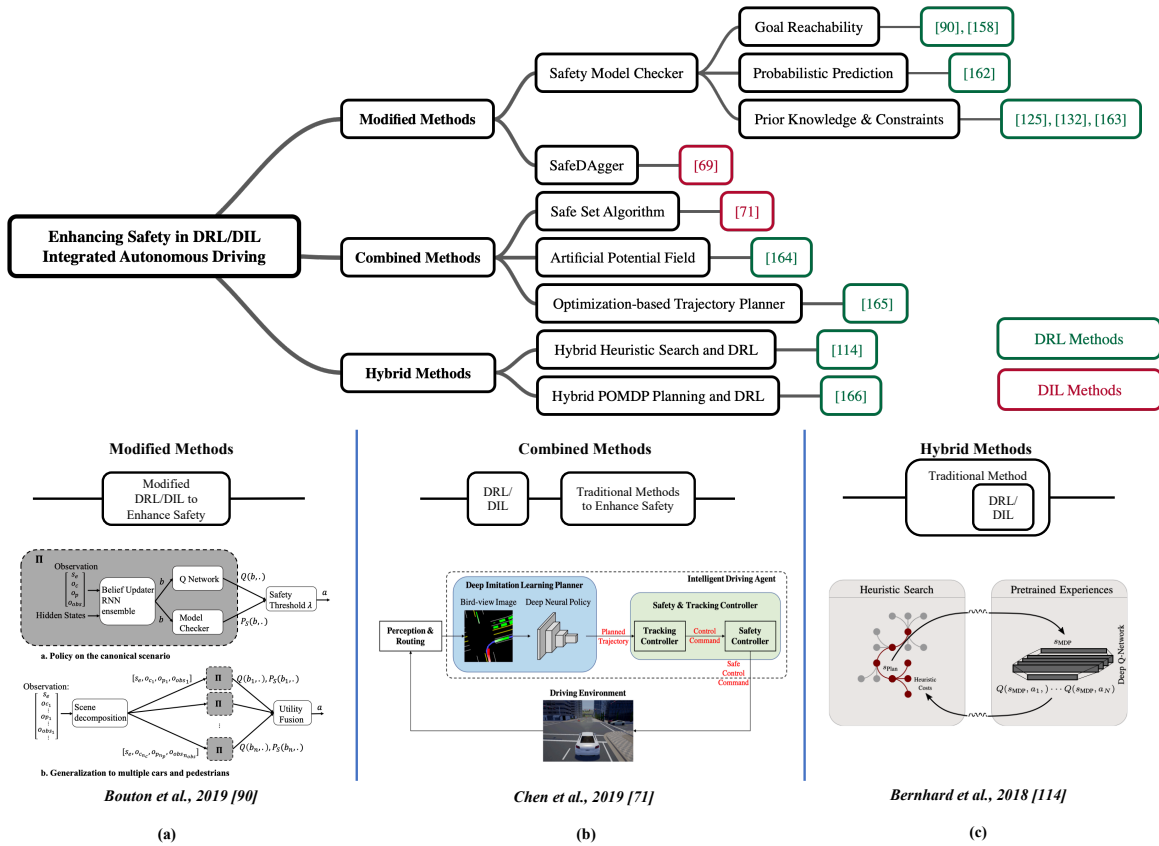
Fig. 5. A taxonomy of the literature on how driving safety is addressed by DRL/DIL models. (a)-(c) Three main methods with typical examples in literature.

## C. Reinforcement Learning Reward Design

A major problem that limits RL's real-world AD applications is the lack of underlying reward functions. Further, the ground truth reward, if it exists, may be multi-modal since human drivers change objectives according to the circumstances. To simplify the problem, current DRL models for AD tasks commonly formulate the reward function as a linear combination of factors, as presented in Fig. 6. A large proportion of studies consider safety and efficiency. Reward terms that are used in the literature are listed in Table VI. Collison and speed are the most common reward terms when considering safety and efficiency, respectively. However, empirically designed reward functions rely heavily on expert knowledge. It is difficult to balance rewards terms, which affects the trained policy performance. Recent studies
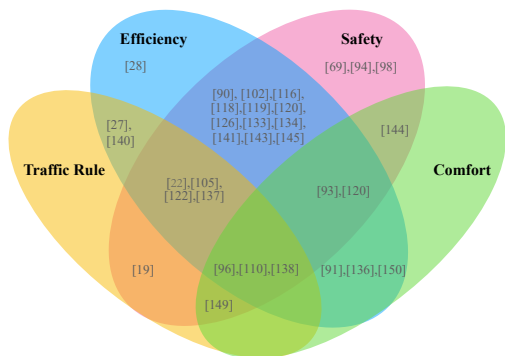
on predictive reward and multi-reward RL may inspire future investigation. Hayashi et al. [154] proposed a predictive reward function that is based on the prediction error of a deep predictive network that models the transition of the surrounding environment. Their hypothesis is that the movement of surrounding vehicles becomes unpredictable when the ego vehicle performs an unnatural driving behavior. Yuan et al. [126] decomposed a single reward function into multi-reward functions to better represent multi-dimensional driving policies through a branched version of Q networks.

## VI. PROBLEM-DRIVEN METHODS

AD application has special requirements on factors such as driving safety, interaction with other traffic participants and uncertainty of the environment. This section reviews the literature from the problem-driven perspective with the objectives of determining how these critical issues are addressed by the DRL/DIL models and identifying the challenges that remain.

### A. Safety-enhanced DRL/DIL for AD

Although DRL/DIL can learn driving policies for complex high-dimensional problems, they only guarantee optimality of the learned policies in a statistical sense. However, in safety-critical AD systems, one failure (e.g., collision) would cause catastrophe. Below, we review representative methods for enhancing safety of DRL/DIL in the AD literature. Fig. 5 categorizes the methods into three groups: (a) *modified methods*: methods that modify the original DRL/DIL algorithms, (b) *combined methods*: methods that combine DRL/DIL with



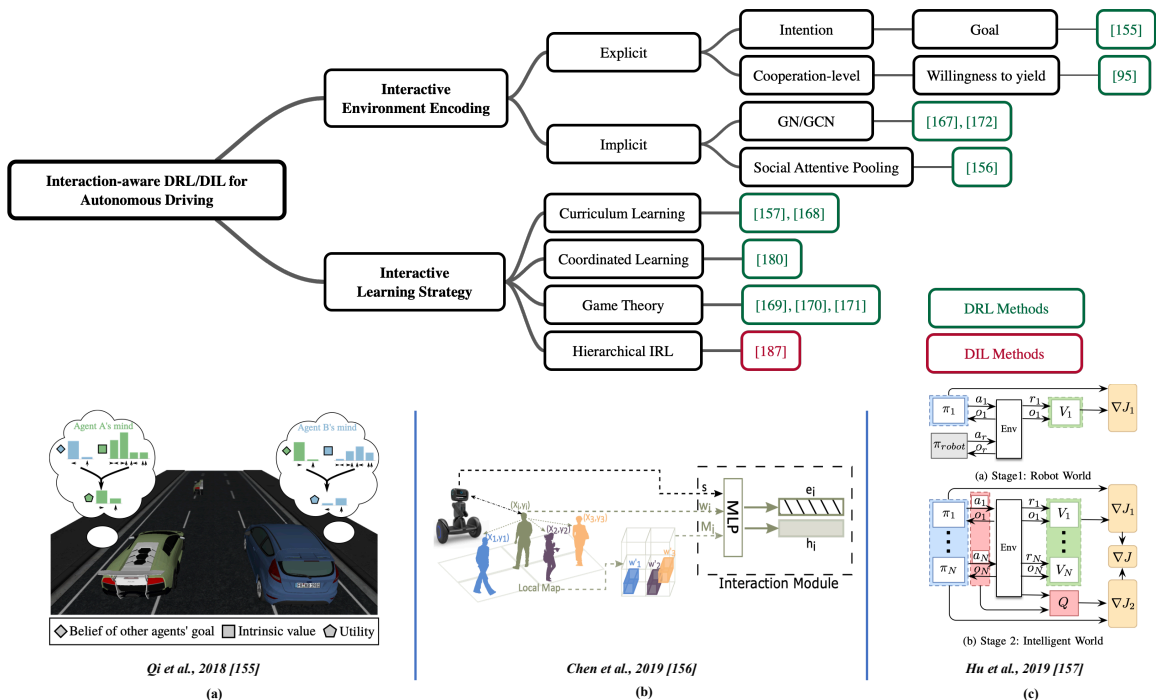Fig. 6. Reinforcement learning rewards for AD tasks.

Fig. 7. A taxonomy of the literature on interaction-aware DRL/DIL models for AD. (a) Qi et al. [155], (b) Chen et al. [156] are examples of explicit and implicit interactive environment encoding, respectively. (c) Hu et al. [157] is an example of interactive learning strategy.

traditional methods, and (c) *hybrid methods*: methods that integrate DRL/DIL into traditional methods.

*1) Modified Methods:* As illustrated in Fig. 5(a), modified methods modify the standard DRL/DIL algorithms to enhance safety, typically by constraining the exploration space [158–160]. A safety model checker is introduced to identify the set of actions that satisfy the safety constraints at each state. This can be realized through several approaches, such as goal reachability [90][161], probabilistic prediction [162] and prior knowledge & constraints [163]. Bouton et al. [90][161] use a probabilistic model checker, as illustrated in Fig. 5(a), to compute the probability of reaching the goal safely at each state-action pair. Then, safe actions are identified by applying a user-defined threshold on the probability. However, the proposed model checker requires a discretization of the state space and a full transition model. Alternatively, Isele et al. [162] proposed the use of probabilistic prediction to identify potentially dangerous actions that would cause collision, but the safety guarantee may not be sufficiently strong if the prediction is not accurate. Prior knowledge & constraints (e.g., lane changes are disallowed if they will lead to small time gaps) are also exploited [125, 132, 163]. For DIL, Zhang et al. [69] proposed SafeDAgger, in which a safety policy is learned to predict the error made by a primary policy without querying the reference policy. If the safety policy determines that it is unsafe to let the primary policy drive, the reference policy will take over. One drawback is that the quality of the learned policy may be limited by that of the reference policy.

*2) Combined Methods:* Various studies combine standard DRL/DIL with traditional rule-based methods to enhance safety. In contrast to the modified methods that are discussed above, combined methods don't modify the learning process of standard DRL/DIL. As presented in Fig. 5(b), Chen et al. [71]

proposed a framework in which DIL plans trajectories, while the rule-based tracking and safe set controller ensure safe control. Xiong et al. [164] proposed the linear combination of the control output from DDPG, artificial potential field and path tracking modules. According to Shalev-Shwartz et al. [165], hard constraints should be injected outside the learning framework. They decompose the double-merge problem into a composition of a learnable DRL policy and a trajectory planning module with non-learnable hard constraints. The learning part enables driving comfort, while the hard constraints guarantee safety.

*3) Hybrid Methods:* Hybrid methods integrate DRL/DIL into traditional heuristic search or POMDP planning methods. As presented in Fig. 5(c), Bernhard et al. [114] integrated experiences in the form of pretrained Q-values into Hybrid A* as heuristics, thereby overcoming the statistical failure rate of DRL while still benefitting computationally from the learned policy. However, the experiments are limited to stationary environments. Pusse et al. [166] presented a hybrid solution that combines DRL and approximate POMDP planning for collision-free autonomous navigation in simulated critical traffic scenarios, which benefits from advantages of both methods.

### B. Interaction-aware DRL/DIL for AD

Interaction is one of the intrinsic characteristics of traffic environments. An intelligent agent should reason beforehand about the behaviors of other traffic participants to passively react or actively adjust its own policy to cooperate or compete with other agents. This section reviews the interaction modeling methods and two groups of interaction-aware DRL/DIL methods for AD, as presented in Fig. 7. One group of methods focus on interactive environment encoding, while the other focus on interactive learning strategies.

*1) Interaction Modeling:* The simplest way to model interaction between multi-agents is to use a standard Markov decision process (MDP), where the other traffic participants are only treated as part of the environment [141, 156]. POMDP is another common interaction model [95, 155, 167], where the agent has limited sensing capabilities. A Markov game (MG) is also used for modeling interaction scenarios. According to whether agents have the same importance, methods can be categorized into three groups: 1) equal importance [157, 168, 169], 2) one vs. others [170], and 3) proactive-passive pair [171].

*2) Interactive Environment Encoding:* Interactive encoding of the environment is a popular research direction. As presented in Fig. 7, mainstream methods can be divided into two groups. One group of methods explicitly model other agents and utilize active reasoning about other agents in the algorithm workflow. POMDP is a common choice for these methods, where the intentions/cooperation levels of other agents are modeled as unobservable states that must be inferred. Qi et al. [155] proposed an intent-aware multi-agent planning framework, as presented in Fig. 7(a), which decouples intent prediction, high-level reasoning and low-level planning. The maintained belief regarding other agents' intents (objectives) was considered in the planning process. Bouton et al. [95] proposed a similar method that maintains a belief regarding the cooperation levels (e.g., the willingness to yield to the ego vehicle) of other drivers.

The other group of methods focuses on utilizing special neural network architectures to capture the interplay between agents by their relation or interaction representations. These methods are usually agnostic regarding the intentions of other agents. Jiang et al. [167] proposed graph convolutional reinforcement learning, in which the multi-agent environment is constructed as a graph. Agents are represented by nodes, and each node's corresponding neighbors are determined by distance or other metrics. Then, the latent features that are produced by graph convolutional layers are exploited to learn cooperation. Similarly, Huegle et al. [172] built upon graph neural networks [173] and proposed the deep scene architecture for learning complex interaction-aware scene representations. Inspired by social pooling [174, 175] and attention models [176, 177], Chen et al. [156] proposed a socially attentive DRL method for interaction-aware robot navigation through a crowd. As illustrated in Fig. 7(b), they extracted pairwise features of interaction between the robot and each human and captured the interactions among humans via local maps. A self-attention mechanism was subsequently used to infer the relative importance of neighboring humans and aggregate interaction features.

*3) Interactive Learning Strategy:* Various learning strategies have been used to learn interactive policies. Curriculum learning has been used to learn interactive policies [157][168], which can decouple complex problems into simpler problems. As presented in Fig. 7(c), Hu et al. [157] proposed an interaction-aware decision making approach that leverages curriculum learning. First, a decentralized critic for each agent is learned to generate distinct behaviors, where the agent does not react to other agents and only learns how to execute rational

actions to complete its own task. Second, a centralized critic is learned to enable agents to interact with each other to realize joint success and maintain smooth traffic. One limitation of these methods is that new models must be learned if the number of agents increases. Based on dynamic coordination graph (DCG) [179], Chao et al. [180] proposed a strategic learning solution for coordinating multiple autonomous vehicles in highways. DCG was utilized to explicitly model the continuously changing coordination dependencies among vehicles. Another group of interaction-aware methods use game theory[181]. Game theory has already been applied in robotics tasks such as robust control [182, 183] and motion planning [184, 185]. Recent years have also witnessed the increasing application of game theory in interaction-aware AD policy learning [169–171]. Li et al. [170] proposed the combination of hierarchical reasoning game theory (i.e. "level-$k$" reasoning [186]) and reinforcement learning. Level-$k$ reasoning is used to model intelligent vehicles' interactions in traffic, while RL evolves these interactions in a time-extended scenario. Ding et al. [171] introduced a proactive-passive game theoretical lane changing framework. The proactive vehicles learn to take actions to merge, while the passive vehicles learn to create merging space. Fisac et al. [169] proposed a novel game-theoretic real-time trajectory planning algorithm. The dynamic game is hierarchically decomposed into a long-horizon "strategic" game and a short-horizon "tactical" game. Furthermore, the long-horizon interaction game is solved to guide short-horizon planning, thereby implicitly extending the planning horizon and pushing the local trajectory optimization closer to global solutions. Apart from combining game theory and RL, solving an imitation learning problem under game-theoretic formalism is another approach. Sun et al. [187] proposed an interactive probabilistic prediction approach that was based on hierarchical inverse reinforcement learning (HIRL). They modeled the problem from the perspective of a two-agent game by explicitly considering the responses of one agent to the other. However, some of the current game-theoretic interaction-aware methods are limited by their two-vehicle settings and simulation experiments [169, 171, 187].

## C. Uncertainty-aware DRL/DIL for AD

Before deployment of a learned model, it is important to determine what it does not understand and estimate the uncertainty of the decision making output. As presented in Fig. 8, this section reviews the uncertainty-aware DRL/DIL methods for AD from three aspects: 1) *AD and deep learning uncertainty*, 2) *uncertainty estimation methods*, and 3) *multimodal driving behavior learning*.

*1) Autonomous Driving and Deep Learning Uncertainty:* Autonomous driving has inherent uncertainty, while deep learning methods have deep learning uncertainty, and the two can intersect. The AD uncertainty can be categorized as:

- **Traffic environment uncertainty** [107, 108, 178, 188–190]. Stochastic and dynamic interactions among agents with distinct behaviors lead to intrinsic irreducible randomness and uncertainty in a traffic environment.
- **Driving behavior uncertainty** [191]. Human driving behavior is multi-modal and stochastic (e.g., a driver can
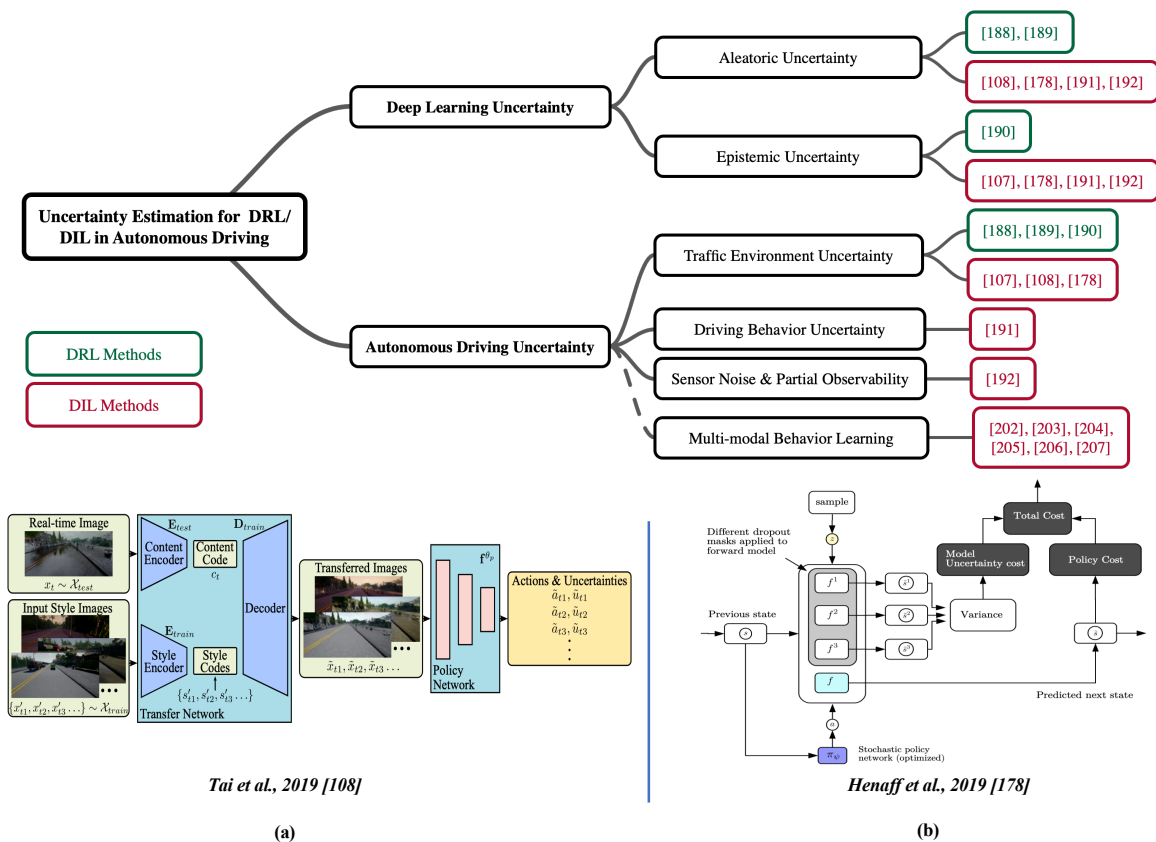
Fig. 8. A taxonomy of the literature on uncertainty-aware DRL/DIL models. (a) Tai et al. [108] addresses aleatoric uncertainty, while (b) Henaff et al. [178] addresses both aleatoric and epistemic uncertainties.

either make a left lane change or a right lane change when he comes up behind a van that is moving at a crawl).
- **Partial observability and sensor noise uncertainty** [192]. In real-world scenarios, the AD agent usually has limited partial observability (e.g., due to occlusion), and there is noise in the sensor observation.

Deriving from Bayesian deep learning approaches, Gal et al. [193] categorized deep learning uncertainty as **aleatoric/data** and **epistemic/model uncertainties**. Aleatoric uncertainty results from incomplete knowledge about the environment (e.g., partial observability and measurement noise), which can't be reduced through access to more or even unlimited data but can be explicitly modeled. In contrast, epistemic uncertainty originates from an insufficient dataset and measures what our model doesn't know, which can be eliminated with sufficient training data. We refer readers to [193, 194] for a deeper background on predictive uncertainty in deep neural networks. Although it is sometimes possible to use only aleatoric [108, 188, 189] or epistemic [107, 190] uncertainty to develop a reasonable model, the ideal approach would be to combine these two uncertainty estimates [178, 191, 192].

*2) Uncertainty Estimation Methods:* Aleatoric uncertainty is usually learned by using the heteroscedastic loss function [194]. The regression task and the loss are formulated as

$$[\tilde{\mathbf{y}}, \tilde{\sigma}] = \mathbf{f}^{\theta}(\mathbf{x}) \tag{20}$$

$$\mathcal{L}(\theta) = \frac{1}{2} \frac{\parallel \mathbf{y} - \tilde{\mathbf{y}} \parallel^2}{\tilde{\sigma}^2} + \frac{1}{2} \log \tilde{\sigma}^2 \tag{21}$$

where $\mathbf{x}$ denotes the input data and $\mathbf{y}$ and $\tilde{\mathbf{y}}$ denote the regression ground truth and the prediction output, respectively. $\theta$ denotes the model parameters, and $\tilde{\sigma}$ is another output of the model, which represents the standard variance of data $\mathbf{x}$ (the aleatoric uncertainty). The loss function can be interpreted as penalizing a large prediction error when the uncertainty is small and relaxing constraints on the prediction error when the uncertainty is large. In practice, the network predicts the log variance $\log \tilde{\sigma}^2$ [194]. Tai et al. [108] proposed an end-to-end real-to-sim visual navigation deployment pipeline, as illustrated in Fig. 8(a). An uncertainty-aware IL policy is trained with the heteroscedastic loss and outputs actions, along with associated uncertainties. A similar technique is proposed by Lee et al. [192].

The epistemic uncertainty is usually estimated via two popular methods: Monte Carlo (MC)-dropout [195, 196] and ensembles [197, 198]. These methods are similar in the sense that both apply probabilistic reasoning on the network weights. The variance of the model output serves as an estimate of the model uncertainty. However, multiple stochastic forward passes through dropout sampling may be time-consuming, while ensemble methods have higher training and storage costs. Kahn et al. [190] proposed an uncertainty-aware RL method that utilizes MC-dropout and bootstrapping [199], where the confidence for a specified obstacle is updated iteratively. Guided by the uncertainty cost, the agent behaves more carefully in unfamiliar scenarios in the early training phase. As presented in Fig. 8(b), Henaff et al. [178] proposed

training a driving policy by unrolling a learned dynamics model over multiple time steps while explicitly penalizing the original policy cost and an uncertainty cost that represents the divergence from the training dataset. Their method estimates both the aleatoric and epistemic uncertainties.

The uncertainty estimation methods that are presented above depend mainly on sampling. A novel uncertainty estimation method that utilizes a mixture density network (MDN) was proposed by Choi et al. [191] for learning from complex and noisy human demonstrations. Since an MDN outputs the parameters for constructing a Gaussian mixture model (GMM), the total variance of the GMM can be calculated analytically, and the acquisition of uncertainty requires only a single forward pass. Distributional reinforcement learning [51, 200, 201] offers another approach for modelling the uncertainty that is associated with actions. It models the RL return $R$ as a random variable that is subject to the probability distribution $Z(r|s,a)$ and the Q-value as the expected return $Q(s,a) = \mathbb{E}_{r \sim Z(r|s,a)}[r]$. In the AD domain, Wang et al. [188] applied distributional DDPG to an energy management strategy (EMS) problem as a case study to evaluate the effects of estimating the uncertainty that is associated with various actions at various states. Bernhard [189] et al. presented a two-step approach for risk-sensitive behavior generation that combined offline distribution reinforcement learning with online risk assessment, which increased safety in intersection crossing scenarios.

*3) Multi-Modal Driving Behavior Learning:* The inherent uncertainty in driving behavior results in multi-modal demonstrations. Many multi-modal imitation learning methods have been proposed. InfoGAIL [202] and Burn-InfoGAIL [203] infer latent/modal variables by maximizing the mutual information between latent variables and state-action pairs. VAE-GAIL [204] introduces a variational auto-encoder for inferring modal variables. However, due to the lack of labels in the demonstrations, these algorithms tend to distinguish latent labels without considering semantic information or the task context. Another direction focuses on labeled data in expert demonstrations. CGAIL [205] sends the modal labels directly to the generator and the discriminator. ACGAIL [206] introduces an auxiliary classifier for reconstructing the modal information, where the classifier cooperates with the discriminator to provide the adversarial loss to the generator. Nevertheless, the above methods mainly leverage random sampling of latent labels from a known prior distribution to distinguish multiple modalities. The trained models rely on manually specified labels to output actions; hence, they cannot select modes adaptively according to environmental scenarios. Recently, Fei et al. [207] proposed Triple-GAIL, which can learn adaptive skill selection and imitation jointly from expert demonstrations, and generated experiences by introducing an auxiliary skill selector.

## VII. DISCUSSION

Although DRL and DIL attract significant amounts of interest in AD research, they remain far from ready for real-world applications, and challenges are faced at the architecture, task and algorithm levels. However, solutions remain largely underexplored. In this section, we discuss these challenges along with future investigation. DRL and DIL also have their own technical challenges; we refer readers to comprehensive discussions in [29, 31, 33].

### A. System architecture

The success of modern AD systems depends on the meticulous design of architectures. The integration of DRL/DIL methods to collaborate with other modules and improve the performance of the system remains a substantial challenge. Studies have demonstrated various ways that DRL/DIL models could be integrated into an AD system. As illustrated Fig. 3, some studies propose new AD architectures, e.g., Modes 1&2, where an entire pipeline from the input of the sensor/perception to the output of the vehicle's actuators is covered. However, the traditional modules of sequential planning are missing in these new architectures, and driving policy is addressed at the control level only. Hence, these AD systems could adapt to only simple tasks, such as road following, that require neither the guidance of goal points nor switching of driving behaviors. The extension of these architectures to accomplish more complicated AD tasks remains a substantial challenge. Other studies utilize the traditional AD architectures, e.g., Modes 3&4, where DRL/DIL models are studied as substitutes for traditional modules to improve the performance in challenging scenarios. Mode 5 studies use both new and traditional architectures. Overall, the research effort until now has been focused more on exploring the potential of DRL/DIL in accomplishing AD tasks, whereas the design of the system architectures has yet to be intensively investigated.

### B. Formulation of driving tasks

Various DRL/DIL formulations have been established for accomplishing AD tasks. However, these formulations rely heavily on empirical designs. As reviewed in Section V, the state space and input data are designed case by case, and ad-hoc reward functions are usually adopted with hand-tuned coefficients that balance the costs regarding safety, efficiency, comfort, and traffic rules, among other factors. Such designs are very brute-force approaches, which lack both theoretical proof and in-depth investigations. Changing the designs or tuning the parameters could result in substantially different driving policies. However, in real-world deployment, more attention should be paid to the following questions: What design could realize the most optimal driving policy? Could such a design adapt to various scenes? How can the boundary conditions of designs be identified? To answer these questions, rigorous studies with comparative experiments are needed.

### C. Safe driving policy

AD applications have high requirements on safety, and guaranteeing the safety of a DRL/DIL integrated AD system is of substantial importance. Compared to traditional rule-based methods, DNN has been widely acknowledged as having poor interpretability. Its "black-box" nature renders difficult

the prediction of when the agent may fail to generate a safe policy. Deep models for real-world AD applications must address unseen or rarely seen scenarios, which is difficult for DL methods as they optimize objectives at the level of expectations over specified instances. To solve this problem, a general strategy is to combine traditional methods to ensure a DRL/DIL agent's functional safety. As reviewed in Fig. 5, various methods have been proposed in the literature, where the problems are usually formulated as compositions of learned policies with hard constraints [125, 132, 163]. However, balancing between the learned optimal policy and the safety guarantee by hard constraints is non-trivial and requires intensive investigation in the future.

### D. Interaction with traffic participants

The capability of human-like interaction is required of self-driving agents for sharing the roads with other traffic participants. As reviewed in Section VI-B, interaction-aware DRL/DIL is a rising topic, but the following problems remain: First, current studies attempt to solve the problem from various perspectives, and the systematic studies are needed. Second, few interaction-aware DIL methods are available, while interaction-aware DRL methods are limited to simplified scenarios that involve only a few agents. The combination of interaction-aware trajectory prediction methods [208–210] may be an open topic of potential value. Third, game theory and multi-agent reinforcement learning (**MARL**) [31] are highly correlated for interactive scenarios. MARL methods usually build on concepts of game theory (e.g., Markov games) to model the interaction process. Apart from DRL/DIL, methods are available for learning interactive policies through traditional game theory approaches, such as Nash equilibrium [211], level-$k$ reasoning [212] and game tree search [213]. Exploiting POMDP planning to learn interactive polices is also a trend [214, 215]. These methods have satisfactory interpretability but are limited to simplified or coarse discretizations of the agents' action space [213, 215]. Although the simplification reduces the computation burden, it tends to also lower the control precision. In the future, the combination of these methods and DRL/DIL may be promising.

### E. Uncertainty of the environment

Decision-making under uncertainty has been studied for decades [216, 217]. Nevertheless, modeling the uncertainty in DRL/DIL formulations remains challenging, especially under complex uncertain traffic environments. Several problems have been identified in current research: First, most uncertainty-aware methods follow the style of deep learning predictive uncertainty [193] without a deeper investigation. Is computing the predictive uncertainty of DNNs sufficient for AD tasks? Second, can the computed uncertainty be effectively utilized to realize a better decision making policy? Various methods incorporate the uncertainty cost into the global cost functions [178, 190], while other methods utilize uncertainty to generate risk-sensitive behavior [189]. Future efforts are needed to identify more promising applications. Third, human-driving behavior is uncertain, or multi-modal. However, DIL performs well for demonstrations from one expert rather than multiple experts [218]. A naive solution is to neglect the multi-modality and treat the demonstrations as if there is only one expert. The main side effect is that the model tends to learn an average policy rather than a multi-modal policy [202]. Thus, determining whether DRL/DIL learn effectively from noisy uncertain naturalistic driving data and generate multi-modal driving behavior according to various scenarios is meaningful.

### F. Validation and benchmarks

Validation and benchmarks are especially important for AD, but far from sufficient effort has been made regarding these aspects. First, comparison between DRL/DIL integrated architectures and traditional architectures is usually neglected in the literature, which is meaningful for identifying the quantitative performance gains and disadvantages of introducing DRL/DIL. Second, systematic comparison between DRL/DIL architectures is necessary. A technical barrier of the former two problems is the lack of a reasonable benchmark. High-fidelity simulators such as CARLA [105] may provide a virtual platform on which various architectures can be deployed and evaluated. Third, exhaustive validation of trained policies before deployment is of vital importance. However, validation is challenging. Real-world testing on vehicles has high costs in terms of time, finances and human labor and could be dangerous. Empirical validation through simulation can reduce the amount of required field testing and can be used as a first step for performance and safety evaluation. However, verification through simulation only ensures the performance in a statistical sense. Even small variations between the simulators and the real scenario can have drastic effects on the system behavior. Future studies are needed to identify practical, effective, low-risk and economical validation methods.

## VIII. CONCLUSIONS

In this study, a comprehensive survey is presented that focuses on autonomous driving policy learning using DRL/DIL, which is addressed simultaneously from the system, task-driven and problem-driven perspectives. The study is conducted at three levels: First, a taxonomy of the literature studies is presented from the system perspective, from which five modes of integration of DRL/DIL models into an AD architecture are identified. Second, the formulations of DRL/DIL models for accomplishing specified AD tasks are comprehensively reviewed, where various designs on the model state and action spaces and the reinforcement learning rewards are covered. Finally, an in-depth review is presented on how the critical issues of AD applications regarding driving safety, interaction with other traffic participants and uncertainty of the environment are addressed by the DRL/DIL models. The major findings are listed below, from which potential topics for future investigation are identified.

- DRL/DIL attract significant amounts of interest in AD research. However, literature studies in this scope have focused more on exploring the potential of DRL/DIL in accomplishing AD tasks, whereas the design of the system architectures remains to be intensively investigated.

- Many DRL/DIL models have been formulated for accomplishing AD tasks. However, these formulations rely heavily on empirical designs, which are brute-force approaches and lack both theoretical proof and in-depth investigations. In the real-world deployment of such models, substantial challenges in terms of stability and robustness may be encountered.

- Driving safety, which is the main issue in AD applications, has received the most attention in the literature. However, the studies on interaction with other traffic participants and the uncertainty of the environment remain highly preliminary, in which the problems have been addressed from divergent perspectives, and have not been conducted systematically.

## REFERENCES

[1] C. Urmson and W. Whittaker, "Self-driving cars and the urban challenge," *IEEE Intelligent Systems*, vol. 23, no. 2, pp. 66–68, 2008.

[2] S. Thrun, "Toward robotic cars," *Communications of the ACM*, vol. 53, no. 4, pp. 99–106, 2010.

[3] A. Eskandarian, *Handbook of intelligent vehicles*. Springer, 2012, vol. 2.

[4] S. M. Grigorescu, B. Trasnea, T. T. Cocias, and G. Macesanu, "A survey of deep learning techniques for autonomous driving," *J. Field Robotics*, vol. 37, no. 3, pp. 362–386, 2020.

[5] W. H. Organization *et al.*, "Global status report on road safety 2018: Summary," World Health Organization, Tech. Rep., 2018.

[6] A. Talebpour and H. S. Mahmassani, "Influence of connected and autonomous vehicles on traffic flow stability and throughput," *Transportation Research Part C: Emerging Technologies*, vol. 71, pp. 143–163, 2016.

[7] W. Payre, J. Cestac, and P. Delhomme, "Intention to use a fully automated car: Attitudes and a priori acceptability," *Transportation research part F: traffic psychology and behaviour*, vol. 27, pp. 252–263, 2014.

[8] E. D. Dickmanns and A. Zapp, "Autonomous high speed road vehicle guidance by computer vision," *IFAC Proceedings Volumes*, vol. 20, no. 5, pp. 221–226, 1987.

[9] C. Thorpe, M. H. Hebert, T. Kanade, and S. A. Shafer, "Vision and navigation for the carnegie-mellon navlab," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 10, no. 3, pp. 362–373, 1988.

[10] D. A. Pomerleau, "Alvinn: An autonomous land vehicle in a neural network," in *Advances in Neural Information Processing Systems*, 1989, pp. 305–313.

[11] S. Thrun, M. Montemerlo, H. Dahlkamp, D. Stavens, A. Aron, J. Diebel, P. Fong, J. Gale, M. Halpenny, G. Hoffmann *et al.*, "Stanley: The robot that won the darpa grand challenge," *Journal of field Robotics*, vol. 23, no. 9, pp. 661–692, 2006.

[12] M. Buehler, K. Iagnemma, and S. Singh, *The DARPA urban challenge: autonomous vehicles in city traffic*. springer, 2009, vol. 56.

[13] D. González, J. Pérez, V. Milanés, and F. Nashashibi, "A review of motion planning techniques for automated vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 17, no. 4, pp. 1135–1145, 2016.

[14] X. Li, Z. Sun, D. Cao, Z. He, and Q. Zhu, "Real-time trajectory planning for autonomous urban driving: Framework, algorithms, and verifications," *IEEE/ASME Transactions on Mechatronics*, vol. 21, no. 2, pp. 740–753, 2016.

[15] B. Paden, M. Čáp, S. Z. Yong, D. Yershov, and E. Frazzoli, "A survey of motion planning and control techniques for self-driving urban vehicles," *IEEE Transactions on Intelligent Vehicles*, vol. 1, no. 1, pp. 33–55, 2016.

[16] S. Ulbrich, A. Reschka, J. Rieken, S. Ernst, G. Bagschik, F. Dierkes, M. Nolte, and M. Maurer, "Towards a functional system architecture for automated vehicles," *arXiv preprint arXiv:1703.08557*, 2017.

[17] L. Li, K. Ota, and M. Dong, "Humanlike driving: Empirical decision-making system for autonomous vehicles," *IEEE Trans. Veh. Technol.*, vol. 67, no. 8, pp. 6814–6823, 2018. [Online]. Available: https://doi.org/10.1109/TVT.2018.2822762

[18] W. Schwarting, J. Alonso-Mora, and D. Rus, "Planning and decision-making for autonomous vehicles," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 1, 05 2018.

[19] G. Yu and I. K. Sethi, "Road-following with continuous learning," in *the Intelligent Vehicles' 95. Symposium*. IEEE, 1995, pp. 412–417.

[20] A. Y. Ng, A. Coates, M. Diel, V. Ganapathi, J. Schulte, B. Tse, E. Berger, and E. Liang, "Autonomous inverted helicopter flight via reinforcement learning," in *International Symposium on Experimental Robotics*, ser. Springer Tracts in Advanced Robotics, vol. 21. Springer, 2004, pp. 363–372.

[21] R. S. Sutton and A. G. Barto, *Reinforcement learning: an introduction*. MIT Press, 2018.

[22] Y. LeCun, Y. Bengio, and G. E. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436–444, 2015.

[23] I. J. Goodfellow, Y. Bengio, and A. C. Courville, *Deep Learning*, ser. Adaptive computation and machine learning. MIT Press, 2016.

[24] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra, "Continuous control with deep reinforcement learning," in *International Conference on Learning Representations*, 2016.

[25] Y. Zhu, R. Mottaghi, E. Kolve, J. J. Lim, A. Gupta, L. Fei-Fei, and A. Farhadi, "Target-driven visual navigation in indoor scenes using deep reinforcement learning," in *International Conference on Robotics and Automation*. IEEE, 2017, pp. 3357–3364.

[26] X. Liang, T. Wang, L. Yang, and E. Xing, "Cirl: Controllable imitative reinforcement learning for vision-based self-driving," in *the European Conference on Computer Vision (ECCV)*, 2018, pp. 584–599.

[27] Y. Chen, C. Dong, P. Palanisamy, P. Mudalige, K. Muelling, and J. M. Dolan, "Attention-based hierarchical deep reinforcement learning for lane change behaviors in autonomous driving," in *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 0–0.

[28] A. Kendall, J. Hawke, D. Janz, P. Mazur, D. Reda, J.-M. Allen, V.-D. Lam, A. Bewley, and A. Shah, "Learning to drive in a day," in *International Conference on Robotics and Automation*. IEEE, 2019, pp. 8248–8254.

[29] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "Deep reinforcement learning: A brief survey," *IEEE Signal Processing Magazine*, vol. 34, no. 6, pp. 26–38, 2017.

[30] Y. Li, "Deep reinforcement learning: An overview," *CoRR*, vol. abs/1701.07274, 2017.

[31] T. T. Nguyen, N. D. Nguyen, and S. Nahavandi, "Deep reinforcement learning for multi-agent systems: A review of challenges, solutions and applications," *CoRR*, vol. abs/1812.11794, 2018.

[32] A. Hussein, M. M. Gaber, E. Elyan, and C. Jayne, "Imitation learning: A survey of learning methods," *ACM Computing Surveys (CSUR)*, vol. 50, no. 2, pp. 1–35, 2017.

[33] T. Osa, J. Pajarinen, G. Neumann, J. A. Bagnell, P. Abbeel, and J. Peters, "An algorithmic perspective on imitation learning," *Found. Trends Robotics*, vol. 7, no. 1-2, pp. 1–179, 2018.

[34] S. Kuutti, R. Bowden, Y. Jin, P. Barber, and S. Fallah, "A survey of deep learning applications to autonomous vehicle control," *CoRR*, vol. abs/1912.10773, 2019.

[35] B. R. Kiran, I. Sobh, V. Talpaert, P. Mannion, A. A. A. Sallab, S. K. Yogamani, and P. Pérez, "Deep reinforcement learning for autonomous driving: A survey," *CoRR*, vol. abs/2002.00444, 2020.

[36] S. B. Thrun, "Efficient exploration in reinforcement learning," 1992.

[37] M. Coggan, "Exploration and exploitation in reinforcement learning," *Research supervised by Prof. Doina Precup, CRA-W DMP Project at McGill University*, 2004.

[38] Z. Hong, T. Shann, S. Su, Y. Chang, T. Fu, and C. Lee, "Diversity-driven exploration strategy for deep reinforcement learning," in *Advances in Neural Information Processing Systems*, 2018.

[39] G. Shani, J. Pineau, and R. Kaplow, "A survey of point-based pomdp solvers," *Autonomous Agents and Multi-Agent Systems*, vol. 27, no. 1, pp. 1–51, 2013.

[40] W. S. Lovejoy, "A survey of algorithmic methods for partially observed markov decision processes," *Annals of Operations Research*, vol. 28, no. 1, pp. 47–65, 1991.

[41] R. Bellman and R. Kalaba, "On the role of dynamic programming in statistical communication theory," *IRE Trans. Inf. Theory*, vol. 3, no. 3, pp. 197–203, 1957.

[42] C. J. C. H. Watkins and P. Dayan, "Technical note q-learning," *Mach. Learn.*, vol. 8, pp. 279–292, 1992.

[43] G. A. Rummery and M. Niranjan, *On-line Q-learning using connectionist systems*. University of Cambridge, Department of Engineering Cambridge, UK, 1994, vol. 37.

[44] R. Bellman, "Dynamic programming," *Science*, vol. 153, no. 3731, pp.

34–37, 1966.

[45] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. A. Riedmiller, A. Fidjeland, G. Ostrovski, and et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.

[46] L.-J. Lin, "Self-improving reactive agents based on reinforcement learning, planning and teaching," *Machine learning*, vol. 8, no. 3-4, pp. 293–321, 1992.

[47] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized experience replay," in *International Conference on Learning Representations*, 2016.

[48] S. Gu, T. P. Lillicrap, I. Sutskever, and S. Levine, "Continuous deep q-learning with model-based acceleration," in *International Conference on Machine Learning*, 2016.

[49] H. v. Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," in *the Thirtieth AAAI Conference on Artificial Intelligence*, 2016, pp. 2094–2100.

[50] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas, "Dueling network architectures for deep reinforcement learning," in *International Conference on Machine Learning*, 2016, pp. 1995–2003.

[51] W. Dabney, M. Rowland, M. G. Bellemare, and R. Munos, "Distributional reinforcement learning with quantile regression," in *AAAI Conference on Artificial Intelligence*, 2018, pp. 2892–2901.

[52] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Mach. Learn.*, vol. 8, pp. 229–256, 1992.

[53] L. C. Baird, "Reinforcement learning in continuous time: Advantage updating," in *IEEE International Conference on Neural Networks*, vol. 4. IEEE, 1994, pp. 2448–2453.

[54] J. Schulman, P. Moritz, S. Levine, M. I. Jordan, and P. Abbeel, "High-dimensional continuous control using generalized advantage estimation," in *International Conference on Learning Representations*, 2016.

[55] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International Conference on Machine Learning*, 2015, pp. 1889–1897.

[56] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *CoRR*, vol. abs/1707.06347, 2017.

[57] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. A. Riedmiller, "Deterministic policy gradient algorithms," in *International Conference on Machine Learning*, vol. 32, 2014, pp. 387–395.

[58] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *International Conference on Machine Learning*, 2016, pp. 1928–1937.

[59] J. Wang, Z. Kurth-Nelson, H. Soyer, J. Z. Leibo, D. Tirumala, R. Munos, C. Blundell, D. Kumaran, and M. M. Botvinick, "Learning to reinforcement learn," in *the 39th Annual Meeting of the Cognitive Science Society, CogSci 2017, London, UK, 16-29 July 2017*. cognitivesciencesociety.org, 2017.

[60] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, and S. Levine, "Soft actor-critic algorithms and applications," *CoRR*, vol. abs/1812.05905, 2018.

[61] B. D. Argall, S. Chernova, M. M. Veloso, and B. Browning, "A survey of robot learning from demonstration," *Robotics Auton. Syst.*, vol. 57, no. 5, pp. 469–483, 2009.

[62] H. John and C. James, "Ngsim interstate 80 freeway dataset," US Fedeal Highway Administration, FHWA-HRT-06-137, Washington, DC, USA, Tech. Rep., 2006.

[63] M. Bojarski, D. Del Testa, D. Dworakowski, B. Firner, B. Flepp, P. Goyal, L. D. Jackel, M. Monfort, U. Muller, J. Zhang *et al.*, "End to end learning for self-driving cars," *arXiv preprint arXiv:1604.07316*, 2016.

[64] M. Bojarski, P. Yeres, A. Choromanska, K. Choromanski, B. Firner, L. Jackel, and U. Muller, "Explaining how a deep neural network trained with end-to-end learning steers a car," *arXiv preprint arXiv:1704.07911*, 2017.

[65] H. Xu, Y. Gao, F. Yu, and T. Darrell, "End-to-end learning of driving models from large-scale video datasets," in *IEEE conference on Computer Vision and Pattern Recognition*, 2017, pp. 2174–2182.

[66] D. A. Pomerleau, "Efficient training of artificial neural networks for autonomous navigation," *Neural computation*, vol. 3, no. 1, pp. 88–97, 1991.

[67] S. Ross and D. Bagnell, "Efficient reductions for imitation learning," in *International Conference on Artificial Intelligence and Statistics*, ser. JMLR Proceedings, vol. 9. JMLR.org, 2010, pp. 661–668.

[68] S. Ross, G. J. Gordon, and D. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," in *International Conference on Artificial Intelligence and Statistics*, ser. JMLR Proceedings, vol. 15, 2011, pp. 627–635.

[69] J. Zhang and K. Cho, "Query-efficient imitation learning for end-to-end autonomous driving," *arXiv preprint arXiv:1605.06450*, 2016.

[70] F. Codevilla, M. Miiller, A. López, V. Koltun, and A. Dosovitskiy, "End-to-end driving via conditional imitation learning," in *International Conference on Robotics and Automation*. IEEE, 2018, pp. 1–9.

[71] J. Chen, B. Yuan, and M. Tomizuka, "Deep imitation learning for autonomous driving in generic urban scenarios with enhanced safety," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019, pp. 2884–2890.

[72] A. Y. Ng and S. J. Russell, "Algorithms for inverse reinforcement learning," in *International Conference on Machine Learning*, 2000, pp. 663–670.

[73] P. Abbeel and A. Y. Ng, "Apprenticeship learning via inverse reinforcement learning," in *International Conference on Machine Learning*, 2004, p. 1.

[74] N. D. Ratliff, J. A. Bagnell, and M. Zinkevich, "Maximum margin planning," in *International Conference on Machine Learning*, vol. 148, 2006, pp. 729–736.

[75] B. D. Ziebart, A. L. Maas, J. A. Bagnell, and A. K. Dey, "Maximum entropy inverse reinforcement learning," in *Aaai*, vol. 8, 2008, pp. 1433–1438.

[76] S. Levine, Z. Popovic, and V. Koltun, "Nonlinear inverse reinforcement learning with gaussian processes," in *Advances in Neural Information Processing Systems*, 2011, pp. 19–27.

[77] N. D. Ratliff, D. M. Bradley, J. A. Bagnell, and J. E. Chestnutt, "Boosting structured prediction for imitation learning," in *Advances in Neural Information Processing Systems*, 2006, pp. 1153–1160.

[78] N. D. Ratliff, D. Silver, and J. A. Bagnell, "Learning to search: Functional gradient techniques for imitation learning," *Auton. Robots*, vol. 27, no. 1, pp. 25–53, 2009.

[79] M. Wulfmeier, P. Ondruska, and I. Posner, "Maximum entropy deep inverse reinforcement learning," *arXiv preprint arXiv:1507.04888*, 2015.

[80] C. Finn, S. Levine, and P. Abbeel, "Guided cost learning: Deep inverse optimal control via policy optimization," in *International Conference on Machine Learning*, 2016, pp. 49–58.

[81] J. Ho and S. Ermon, "Generative adversarial imitation learning," in *Advances in Neural Information Processing Systems*, 2016, pp. 4565–4573.

[82] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems*, 2014, pp. 2672–2680.

[83] C. Finn, P. F. Christiano, P. Abbeel, and S. Levine, "A connection between generative adversarial networks, inverse reinforcement learning, and energy-based models," *CoRR*, vol. abs/1611.03852, 2016.

[84] J. Fu, K. Luo, and S. Levine, "Learning robust rewards with adversarial inverse reinforcement learning," *CoRR*, vol. abs/1710.11248, 2017.

[85] U. Muller, J. Ben, E. Cosatto, B. Flepp, and Y. L. Cun, "Off-road obstacle avoidance through end-to-end learning," in *Advances in neural information processing systems*, 2006, pp. 739–746.

[86] V. Rausch, A. Hansen, E. Solowjow, C. Liu, E. Kreuzer, and J. K. Hedrick, "Learning a deep neural net policy for end-to-end control of autonomous vehicles," in *2017 American Control Conference (ACC)*. IEEE, 2017, pp. 4914–4919.

[87] H. M. Eraqi, M. N. Moustafa, and J. Honer, "End-to-end deep learning for steering autonomous vehicles considering temporal dependencies," *arXiv preprint arXiv:1710.03804*, 2017.

[88] D. Wang, C. Devin, Q.-Z. Cai, F. Yu, and T. Darrell, "Deep object-centric policies for autonomous driving," in *International Conference on Robotics and Automation*. IEEE, 2019, pp. 8853–8859.

[89] Y. Pan, C.-A. Cheng, K. Saigol, K. Lee, X. Yan, E. Theodorou, and B. Boots, "Agile autonomous driving using end-to-end deep imitation learning," *arXiv preprint arXiv:1709.07174*, 2017.

[90] M. Bouton, A. Nakhaei, K. Fujimura, and M. J. Kochenderfer, "Safe reinforcement learning with scene decomposition for navigating complex urban environments," in *Intelligent Vehicles Symposium*. IEEE, 2019, pp. 1469–1476.

[91] P. Wang, C.-Y. Chan, and A. de La Fortelle, "A reinforcement learning based approach for automated lane change maneuvers," in *Intelligent Vehicles Symposium*. IEEE, 2018, pp. 1379–1384.

[92] X. Chen, Y. Zhai, C. Lu, J. Gong, and G. Wang, "A learning model for personalized adaptive cruise control," in *Intelligent Vehicles Symposium*. IEEE, 2017, pp. 379–384.
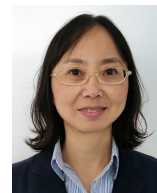
[93] P. Wang and C.-Y. Chan, "Formulation of deep reinforcement learning architecture toward autonomous driving for on-ramp merge," in *International Conference on Intelligent Transportation Systems*. IEEE, 2017, pp. 1–6.

[94] H. Chae, C. M. Kang, B. Kim, J. Kim, C. C. Chung, and J. W. Choi, "Autonomous braking system via deep reinforcement learning," in *International Conference on Intelligent Transportation Systems*. IEEE, 2017, pp. 1–6.

[95] M. Bouton, A. Nakhaei, K. Fujimura, and M. J. Kochenderfer, "Cooperation-aware reinforcement learning for merging in dense traffic," in *IEEE Intelligent Transportation Systems Conference*. IEEE, 2019, pp. 3441–3447.

[96] Y. Tang, "Towards learning multi-agent negotiations via self-play," in *IEEE International Conference on Computer Vision Workshops*, 2019, pp. 0–0.

[97] A. Folkers, M. Rick, and C. Büskens, "Controlling an autonomous vehicle with deep reinforcement learning," in *Intelligent Vehicles Symposium*. IEEE, 2019, pp. 2025–2031.

[98] H. Porav and P. Newman, "Imminent collision mitigation with reinforcement learning and vision," in *International Conference on Intelligent Transportation Systems*. IEEE, 2018, pp. 958–964.

[99] S. Wang, D. Jia, and X. Weng, "Deep reinforcement learning for autonomous driving," *arXiv preprint arXiv:1811.11329*, 2018.

[100] P. Wang, H. Li, and C.-Y. Chan, "Continuous control for automated lane change behavior based on deep deterministic policy gradient algorithm," in *Intelligent Vehicles Symposium*. IEEE, 2019, pp. 1454–1460.

[101] M. Kaushik, V. Prasad, K. M. Krishna, and B. Ravindran, "Overtaking maneuvers in simulated highway driving using deep reinforcement learning," in *Intelligent Vehicles Symposium*. IEEE, 2018, pp. 1885–1890.

[102] A. E. Sallab, M. Abdou, E. Perot, and S. Yogamani, "End-to-end deep reinforcement learning for lane keeping assist," *arXiv preprint arXiv:1612.04340*, 2016.

[103] R. Vasquez and B. Farooq, "Multi-objective autonomous braking system using naturalistic dataset," in *IEEE Intelligent Transportation Systems Conference*. IEEE, 2019, pp. 4348–4353.

[104] S. Hecker, D. Dai, and L. Van Gool, "End-to-end learning of driving models with surround-view cameras and route planners," in *the European Conference on Computer Vision (ECCV)*, 2018, pp. 435–453.

[105] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "Carla: An open urban driving simulator," *arXiv preprint arXiv:1711.03938*, 2017.

[106] M. Abdou, H. Kamal, S. El-Tantawy, A. Abdelkhalek, O. Adel, K. Hamdy, and M. Abaas, "End-to-end deep conditional imitation learning for autonomous driving," in *2019 31st International Conference on Microelectronics (ICM)*. IEEE, 2019, pp. 346–350.

[107] Y. Cui, D. Isele, S. Niekum, and K. Fujimura, "Uncertainty-aware data aggregation for deep imitation learning," in *International Conference on Robotics and Automation*. IEEE, 2019, pp. 761–767.

[108] L. Tai, P. Yun, Y. Chen, C. Liu, H. Ye, and M. Liu, "Visual-based autonomous driving deployment from a stochastic and uncertainty-aware perspective," *arXiv preprint arXiv:1903.00821*, 2019.

[109] M. Buechel and A. Knoll, "Deep reinforcement learning for predictive longitudinal control of automated vehicles," in *International Conference on Intelligent Transportation Systems*. IEEE, 2018, pp. 2391–2397.

[110] J. Chen, B. Yuan, and M. Tomizuka, "Model-free deep reinforcement learning for urban autonomous driving," in *IEEE Intelligent Transportation Systems Conference*. IEEE, 2019, pp. 2765–2771.

[111] M. Bansal, A. Krizhevsky, and A. Ogale, "Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst," *arXiv preprint arXiv:1812.03079*, 2018.

[112] L. Sun, C. Peng, W. Zhan, and M. Tomizuka, "A fast integrated planning and control framework for autonomous driving via imitation learning," in *Dynamic Systems and Control Conference*, vol. 51913. American Society of Mechanical Engineers, 2018, p. V003T37A012.

[113] M. Wulfmeier, D. Rao, D. Z. Wang, P. Ondruska, and I. Posner, "Large-scale cost function learning for path planning using deep inverse reinforcement learning," *The International Journal of Robotics Research*, vol. 36, no. 10, pp. 1073–1087, 2017.

[114] J. Bernhard, R. Gieselmann, K. Esterle, and A. Knol, "Experience-based heuristic search: Robust motion planning with deep q-learning," in *International Conference on Intelligent Transportation Systems*. IEEE, 2018, pp. 3175–3182.

[115] P. Hart, L. Rychly, and A. Knoll, "Lane-merging using policy-based reinforcement learning and post-optimization," in *IEEE Intelligent Transportation Systems Conference*. IEEE, 2019, pp. 3176–3181.

[116] P. Wang, D. Liu, J. Chen, H. Li, and C.-Y. Chan, "Human-like decision making for autonomous driving via adversarial inverse reinforcement learning," *arXiv*, pp. arXiv–1911, 2019.

[117] S. Sharifzadeh, I. Chiotellis, R. Triebel, and D. Cremers, "Learning to drive using inverse reinforcement learning and deep q-networks," *arXiv preprint arXiv:1612.03653*, 2016.

[118] A. Alizadeh, M. Moghadam, Y. Bicer, N. K. Ure, U. Yavas, and C. Kurtulus, "Automated lane change decision making using deep reinforcement learning in dynamic and uncertain highway environment," in *IEEE Intelligent Transportation Systems Conference*. IEEE, 2019, pp. 1399–1404.

[119] N. Deshpande and A. Spalanzani, "Deep reinforcement learning based vehicle navigation amongst pedestrians using a grid-based state representation," in *IEEE Intelligent Transportation Systems Conference*. IEEE, 2019, pp. 2081–2086.

[120] T. Tram, I. Batkovic, M. Ali, and J. Sjöberg, "Learning when to drive in intersections by combining reinforcement learning and model predictive control," in *International Conference on Intelligent Transportation Systems*. IEEE, 2019, pp. 3263–3268.

[121] D. Isele, R. Rahimi, A. Cosgun, K. Subramanian, and K. Fujimura, "Navigating occluded intersections with autonomous vehicles using deep reinforcement learning," in *International Conference on Robotics and Automation*. IEEE, 2018, pp. 2034–2039.

[122] C. Li and K. Czarnecki, "Urban driving with multi-objective deep reinforcement learning," in *International Conference on Autonomous Agents and MultiAgent Systems (AAMAS)*. International Foundation for Autonomous Agents and Multiagent Systems, 2019, pp. 359–367.

[123] M. P. Ronecker and Y. Zhu, "Deep q-network based decision making for autonomous driving," in *International Conference on Robotics and Automation Sciences*. IEEE, 2019, pp. 154–160.

[124] P. Wolf, K. Kurzer, T. Wingert, F. Kuhnt, and J. M. Zollner, "Adaptive behavior generation for autonomous driving using deep reinforcement learning with compact semantic states," in *Intelligent Vehicles Symposium*. IEEE, 2018, pp. 993–1000.

[125] B. Mirchevska, C. Pek, M. Werling, M. Althoff, and J. Boedecker, "High-level decision making for safe and reasonable autonomous lane changing using reinforcement learning," in *International Conference on Intelligent Transportation Systems*. IEEE, 2018, pp. 2156–2162.

[126] W. Yuan, M. Yang, Y. He, C. Wang, and B. Wang, "Multi-reward architecture based reinforcement learning for highway driving policies," in *International Conference on Intelligent Transportation Systems*. IEEE, 2019, pp. 3810–3815.

[127] J. Lee and J. W. Choi, "May i cut into your lane?: A policy network to learn interactive lane change behavior for autonomous driving," in *IEEE Intelligent Transportation Systems Conference*. IEEE, 2019, pp. 4342–4347.

[128] C. You, J. Lu, D. Filev, and P. Tsiotras, "Highway traffic modeling and decision making for autonomous vehicle using reinforcement learning," in *Intelligent Vehicles Symposium*. IEEE, 2018, pp. 1227–1232.

[129] L. Wang, F. Ye, Y. Wang, J. Guo, I. Papamichail, M. Papageorgiou, S. Hu, and L. Zhang, "A q-learning foresighted approach to ego-efficient lane changes of connected and automated vehicles on freeways," in *IEEE Intelligent Transportation Systems Conference*. IEEE, 2019, pp. 1385–1392.

[130] C.-J. Hoel, K. Wolff, and L. Laine, "Automated speed and lane change decision making using deep reinforcement learning," in *International Conference on Intelligent Transportation Systems*. IEEE, 2018, pp. 2148–2155.

[131] Y. Zhang, P. Sun, Y. Yin, L. Lin, and X. Wang, "Human-like autonomous vehicle speed control by deep reinforcement learning with double q-learning," in *Intelligent Vehicles Symposium*. IEEE, 2018, pp. 1251–1256.

[132] D. Liu, M. Brännstrom, A. Backhouse, and L. Svensson, "Learning faster to perform autonomous lane changes by constructing maneuvers from shielded semantic actions," in *IEEE Intelligent Transportation Systems Conference*. IEEE, 2019, pp. 1838–1844.

[133] K. Min, H. Kim, and K. Huh, "Deep distributional reinforcement learning based high-level driving policy determination," *IEEE Transactions on Intelligent Vehicles*, vol. 4, no. 3, pp. 416–424, 2019.

[134] K. Rezaee, P. Yadmellat, M. S. Nosrati, E. A. Abolfathi, M. Elmahgiubi, and J. Luo, "Multi-lane cruising using hierarchical planning and reinforcement learning," in *International Conference on Intelligent Transportation Systems*. IEEE, 2019, pp. 1800–1806.

[135] T. Shi, P. Wang, X. Cheng, C.-Y. Chan, and D. Huang, "Driving decision and control for automated lane change behavior based on deep reinforcement learning," in *IEEE Intelligent Transportation Systems*

*Conference*. IEEE, 2019, pp. 2895–2900.

[136] J. Chen, Z. Wang, and M. Tomizuka, "Deep hierarchical reinforcement learning for autonomous driving with distinct behaviors," in *Intelligent Vehicles Symposium*. IEEE, 2018, pp. 1239–1244.

[137] 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 *International Conference on Intelligent Transportation Systems*. IEEE, 2018, pp. 2377–2382.

[138] C. Paxton, V. Raman, G. D. Hager, and M. Kobilarov, "Combining neural networks and tree search for task and motion planning in challenging environments," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2017, pp. 6059–6066.

[139] F. Behbahani, K. Shiarlis, X. Chen, V. Kurin, S. Kasewa, C. Stirbu, J. Gomes, S. Paul, F. A. Oliehoek, J. Messias *et al.*, "Learning from demonstration in the wild," in *International Conference on Robotics and Automation*. IEEE, 2019, pp. 775–781.

[140] L. Chen, Y. Chen, X. Yao, Y. Shan, and L. Chen, "An adaptive path tracking controller based on reinforcement learning with urban driving application," in *Intelligent Vehicles Symposium*. IEEE, 2019, pp. 2411–2416.

[141] M. Huegle, G. Kalweit, B. Mirchevska, M. Werling, and J. Boedecker, "Dynamic input for deep reinforcement learning in autonomous driving," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019, pp. 7566–7573.

[142] C. Desjardins and B. Chaib-Draa, "Cooperative adaptive cruise control: A reinforcement learning approach," *IEEE Transactions on intelligent transportation systems*, vol. 12, no. 4, pp. 1248–1260, 2011.

[143] D. Zhao, B. Wang, and D. Liu, "A supervised actor–critic approach for adaptive cruise control," *Soft Computing*, vol. 17, no. 11, pp. 2089–2099, 2013.

[144] D. Zhao, Z. Xia, and Q. Zhang, "Model-free optimal control based intelligent cruise control with hardware-in-the-loop demonstration [research frontier]," *IEEE Computational Intelligence Magazine*, vol. 12, no. 2, pp. 56–69, 2017.

[145] K. Min, H. Kim, and K. Huh, "Deep q learning based high level driving policy determination," in *Intelligent Vehicles Symposium*. IEEE, 2018, pp. 226–231.

[146] A. Kuefler, J. Morton, T. Wheeler, and M. Kochenderfer, "Imitating driver behavior with generative adversarial networks," in *Intelligent Vehicles Symposium*. IEEE, 2017, pp. 204–211.

[147] R. P. Bhattacharyya, D. J. Phillips, B. Wulfe, J. Morton, A. Kuefler, and M. J. Kochenderfer, "Multi-agent imitation learning for driving simulation," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 1534–1539.

[148] M. Kuderer, S. Gulati, and W. Burgard, "Learning driving styles for autonomous vehicles from demonstration," in *International Conference on Robotics and Automation*. IEEE, 2015, pp. 2641–2646.

[149] R. P. Bhattacharyya, D. J. Phillips, C. Liu, J. K. Gupta, K. Driggs-Campbell, and M. J. Kochenderfer, "Simulating emergent properties of human driving behavior using multi-agent reward augmented imitation learning," in *International Conference on Robotics and Automation*. IEEE, 2019, pp. 789–795.

[150] Z. Huang, X. Xu, H. He, J. Tan, and Z. Sun, "Parameterized batch reinforcement learning for longitudinal control of autonomous land vehicles," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 49, no. 4, pp. 730–741, 2017.

[151] M. J. Hausknecht and P. Stone, "Deep reinforcement learning in parameterized action space," in *International Conference on Learning Representations*, 2016.

[152] M. Everett, Y. F. Chen, and J. P. How, "Motion planning among dynamic, decision-making agents with deep reinforcement learning," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 3052–3059.

[153] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Póczos, R. Salakhutdinov, and A. J. Smola, "Deep sets," in *Advances in Neural Information Processing Systems 30*, 2017.

[154] D. Hayashi, Y. Xu, T. Bando, and K. Takeda, "A predictive reward function for human-like driving based on a transition model of surrounding environment," in *International Conference on Robotics and Automation*. IEEE, 2019, pp. 7618–7624.

[155] S. Qi and S.-C. Zhu, "Intent-aware multi-agent reinforcement learning," in *International Conference on Robotics and Automation*. IEEE, 2018, pp. 7533–7540.

[156] C. Chen, Y. Liu, S. Kreiss, and A. Alahi, "Crowd-robot interaction: Crowd-aware robot navigation with attention-based deep reinforcement learning," in *International Conference on Robotics and Automation*. IEEE, 2019, pp. 6015–6022.

[157] Y. Hu, A. Nakhaei, M. Tomizuka, and K. Fujimura, "Interaction-aware decision making with adaptive strategies under merging scenarios," *arXiv preprint arXiv:1904.06025*, 2019.

[158] J. García, Fern, and o Fernández, "A comprehensive survey on safe reinforcement learning," *Journal of Machine Learning Research*, vol. 16, no. 42, pp. 1437–1480, 2015.

[159] N. Jansen, B. Könighofer, S. Junges, and R. Bloem, "Shielded decision-making in mdps," *CoRR*, vol. abs/1807.06096, 2018.

[160] N. Fulton and A. Platzer, "Safe reinforcement learning via formal methods: Toward safe control through proof and learning," in *the Thirty-Second AAAI Conference on Artificial Intelligence*, 2018, pp. 6485–6492.

[161] M. Bouton, J. Karlsson, A. Nakhaei, K. Fujimura, M. J. Kochenderfer, and J. Tumova, "Reinforcement learning with probabilistic guarantees for autonomous driving," *arXiv preprint arXiv:1904.07189*, 2019.

[162] D. Isele, A. Nakhaei, and K. Fujimura, "Safe reinforcement learning on autonomous vehicles," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 1–6.

[163] M. Mukadam, A. Cosgun, A. Nakhaei, and K. Fujimura, "Tactical decision making for lane changing with deep reinforcement learning," 2017.

[164] X. Xiong, J. Wang, F. Zhang, and K. Li, "Combining deep reinforcement learning and safety based control for autonomous driving," *arXiv preprint arXiv:1612.00147*, 2016.

[165] S. Shalev-Shwartz, S. Shammah, and A. Shashua, "Safe, multi-agent, reinforcement learning for autonomous driving," *CoRR*, vol. abs/1610.03295, 2016.

[166] F. Pusse and M. Klusch, "Hybrid online pomdp planning and deep reinforcement learning for safer self-driving cars," in *Intelligent Vehicles Symposium*. IEEE, 2019, pp. 1013–1020.

[167] J. Jiang, C. Dun, T. Huang, and Z. Lu, "Graph convolutional reinforcement learning," in *International Conference on Learning Representations*, 2020.

[168] A. Mohseni-Kabir, D. Isele, and K. Fujimura, "Interaction-aware multi-agent reinforcement learning for mobile agents with individual goals," in *International Conference on Robotics and Automation*. IEEE, 2019, pp. 3370–3376.

[169] J. F. Fisac, E. Bronstein, E. Stefansson, D. Sadigh, S. S. Sastry, and A. D. Dragan, "Hierarchical game-theoretic planning for autonomous vehicles," in *International Conference on Robotics and Automation*. IEEE, 2019, pp. 9590–9596.

[170] 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 Transactions on control systems technology*, vol. 26, no. 5, pp. 1782–1797, 2017.

[171] G. Ding, S. Aghli, C. Heckman, and L. Chen, "Game-theoretic cooperative lane changing using data-driven models," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 3640–3647.

[172] M. Huegle, G. Kalweit, M. Werling, and J. Boedecker, "Dynamic interaction-aware scene understanding for reinforcement learning in autonomous driving," *arXiv preprint arXiv:1909.13582*, 2019.

[173] T. N. Kipf and M. Welling, "Semi-supervised classification with graph convolutional networks," in *International Conference on Learning Representations*, 2017.

[174] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, F. Li, and S. Savarese, "Social LSTM: human trajectory prediction in crowded spaces," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 961–971.

[175] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, "Social GAN: socially acceptable trajectories with generative adversarial networks," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2255–2264.

[176] A. Vemula, K. Muelling, and J. Oh, "Social attention: Modeling attention in human crowds," in *International Conference on Robotics and Automation*. IEEE, 2018, pp. 1–7.

[177] Y. Hoshen, "VAIN: attentional multi-agent predictive modeling," in *Advances in Neural Information Processing Systems*, 2017, pp. 2701–2711.

[178] M. Henaff, A. Canziani, and Y. LeCun, "Model-predictive policy learning with uncertainty regularization for driving in dense traffic," in *International Conference on Learning Representations*, 2019.

[179] C. Guestrin, M. G. Lagoudakis, and R. Parr, "Coordinated reinforcement learning," in *International Conference on Machine Learning*, 2002, pp. 227–234.

[180] C. Yu, X. Wang, X. Xu, M. Zhang, H. Ge, J. Ren, L. Sun, B. Chen, and

G. Tan, "Distributed multiagent coordinated learning for autonomous driving in highways based on dynamic coordination graphs," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 2, pp. 735–748, 2019.

[181] K. Leyton-Brown and Y. Shoham, *Essentials of Game Theory: A Concise Multidisciplinary Introduction*, ser. Synthesis Lectures on Artificial Intelligence and Machine Learning. Morgan & Claypool Publishers, 2008.

[182] G. P. Papavassilopoulos and M. G. Safonov, "Robust control design via game theoretic methods," in *IEEE Conference on Decision and Control*, 1989, pp. 382–387 vol.1.

[183] U. Branch, S. Ganebnyi, S. Kumkov, V. Patsko, and S. Pyatko, "Robust control in game problems with linear dynamics," *International Journal of Mathematics, Game Theory and Algebra*, vol. 3, 01 2007.

[184] Hong Zhang, V. Kumar, and J. Ostrowski, "Motion planning with uncertainty," in *International Conference on Robotics and Automation*, vol. 1, 1998, pp. 638–643 vol.1.

[185] M. Zhu, M. Otte, P. Chaudhari, and E. Frazzoli, "Game theoretic controller synthesis for multi-robot motion planning part i: Trajectory based algorithms," in *International Conference on Robotics and Automation*, 2014, pp. 1646–1651.

[186] P. Wilson and D. Stahl, "On players' models of other players: Theory and experimental evidence," *Games and Economic Behavior*, vol. 10, pp. 218–254, 07 1995.

[187] L. Sun, W. Zhan, and M. Tomizuka, "Probabilistic prediction of interactive driving behavior via hierarchical inverse reinforcement learning," in *International Conference on Intelligent Transportation Systems*. IEEE, 2018, pp. 2111–2117.

[188] P. Wang, Y. Li, S. Shekhar, and W. F. Northrop, "Uncertainty estimation with distributional reinforcement learning for applications in intelligent transportation systems: A case study," in *IEEE Intelligent Transportation Systems Conference*. IEEE, 2019, pp. 3822–3827.

[189] J. Bernhard, S. Pollok, and A. Knoll, "Addressing inherent uncertainty: Risk-sensitive behavior generation for automated driving using distributional reinforcement learning," in *sium*. IEEE, 2019, pp. 2148–2155.

[190] G. Kahn, A. Villaflor, V. Pong, P. Abbeel, and S. Levine, "Uncertainty-aware reinforcement learning for collision avoidance," *arXiv preprint arXiv:1702.01182*, 2017.

[191] S. Choi, K. Lee, S. Lim, and S. Oh, "Uncertainty-aware learning from demonstration using mixture density networks with sampling-free variance modeling," in *International Conference on Robotics and Automation*. IEEE, 2018, pp. 6915–6922.

[192] K. Lee, K. Saigol, and E. A. Theodorou, "Safe end-to-end imitation learning for model predictive control," *arXiv preprint arXiv:1803.10231*, 2018.

[193] Y. Gal, "Uncertainty in deep learning," *University of Cambridge*, vol. 1, no. 3, 2016.

[194] A. Kendall and Y. Gal, "What uncertainties do we need in bayesian deep learning for computer vision?" in *Advances in neural information processing systems*, 2017, pp. 5574–5584.

[195] Y. Gal and Z. Ghahramani, "Bayesian convolutional neural networks with bernoulli approximate variational inference," *CoRR*, vol. abs/1506.02158, 2015.

[196] Y. Gal, R. Islam, and Z. Ghahramani, "Deep bayesian active learning with image data," in *International Conference on Machine Learning*, 2017, pp. 1183–1192.

[197] T. G. Dietterich, "Ensemble methods in machine learning," in *Multiple Classifier Systems, First International Workshop, MCS 2000, Cagliari, Italy, June 21-23, 2000, Proceedings*, ser. Lecture Notes in Computer Science, vol. 1857. Springer, 2000, pp. 1–15.

[198] B. Lakshminarayanan, A. Pritzel, and C. Blundell, "Simple and scalable predictive uncertainty estimation using deep ensembles," in *Advances in Neural Information Processing Systems 30*, 2017.

[199] B. Efron and R. Tibshirani, *An Introduction to the Bootstrap*. Springer, 1993.

[200] M. G. Bellemare, W. Dabney, and R. Munos, "A distributional perspective on reinforcement learning," in *International Conference on Machine Learning*, 2017, pp. 449–458.

[201] W. Dabney, G. Ostrovski, D. Silver, and R. Munos, "Implicit quantile networks for distributional reinforcement learning," in *International Conference on Machine Learning*, 2018, pp. 1104–1113.

[202] Y. Li, J. Song, and S. Ermon, "Infogail: Interpretable imitation learning from visual demonstrations," in *Advances in Neural Information Processing Systems*, 2017, pp. 3812–3822.

[203] A. Kuefler and M. J. Kochenderfer, "Burn-in demonstrations for multi-modal imitation learning," *arXiv preprint arXiv:1710.05090*, 2017.

[204] Z. Wang, J. S. Merel, S. E. Reed, N. de Freitas, G. Wayne, and N. Heess, "Robust imitation of diverse behaviors," in *Advances in Neural Information Processing Systems*, 2017, pp. 5320–5329.

[205] J. Merel, Y. Tassa, D. TB, S. Srinivasan, J. Lemmon, Z. Wang, G. Wayne, and N. Heess, "Learning human behaviors from motion capture by adversarial imitation," *arXiv preprint arXiv:1707.02201*, 2017.

[206] J. Lin and Z. Zhang, "Acgail: Imitation learning about multiple intentions with auxiliary classifier gans," in *Pacific Rim International Conference on Artificial Intelligence*. Springer, 2018, pp. 321–334.

[207] C. Fei, B. Wang, Y. Zhuang, Z. Zhang, J. Hao, H. Zhang, X. Ji, and W. Liu, "Triple-gail: A multi-modal imitation learning framework with generative adversarial nets," in *the Twenty-Ninth International Joint Conference on Artificial Intelligence (IJCAI)*, 2020, pp. 2929–2935.

[208] E. Schmerling, K. Leung, W. Vollprecht, and M. Pavone, "Multimodal probabilistic model-based planning for human-robot interaction," in *International Conference on Robotics and Automation*. IEEE, 2018, pp. 1–9.

[209] J. Li, H. Ma, and M. Tomizuka, "Interaction-aware multi-agent tracking and probabilistic behavior prediction via adversarial learning," in *International Conference on Robotics and Automation*. IEEE, 2019, pp. 6658–6664.

[210] H. Ma, J. Li, W. Zhan, and M. Tomizuka, "Wasserstein generative learning with kinematic constraints for probabilistic interactive driving behavior prediction," in *Intelligent Vehicles Symposium*. IEEE, 2019, pp. 2477–2483.

[211] A. Turnwald, D. Althoff, D. Wollherr, and M. Buss, "Understanding human avoidance behavior: interaction-aware decision making based on game theory," *International Journal of Social Robotics*, vol. 8, no. 2, pp. 331–351, 2016.

[212] R. Tian, S. Li, N. Li, I. Kolmanovsky, A. Girard, and Y. Yildiz, "Adaptive game-theoretic decision making for autonomous vehicle control at roundabouts," in *IEEE Conference on Decision and Control (CDC)*. IEEE, 2018, pp. 321–326.

[213] D. Isele, "Interactive decision making for autonomous vehicles in dense traffic," in *IEEE Intelligent Transportation Systems Conference*. IEEE, 2019, pp. 3981–3986.

[214] H. Bai, S. Cai, N. Ye, D. Hsu, and W. S. Lee, "Intention-aware online pomdp planning for autonomous driving in a crowd," in *International Conference on Robotics and Automation*. IEEE, 2015, pp. 454–460.

[215] C. Hubmann, J. Schulz, G. Xu, D. Althoff, and C. Stiller, "A belief state planner for interactive merge maneuvers in congested traffic," in *International Conference on Intelligent Transportation Systems*. IEEE, 2018, pp. 1617–1624.

[216] C. A. Holloway, *Decision making under uncertainty: models and choices*. Prentice-Hall Englewood Cliffs, NJ, 1979, vol. 8.

[217] M. J. Kochenderfer, *Decision making under uncertainty: theory and application*. MIT press, 2015.

[218] R. Camacho and D. Michie, "Behavioral cloning A correction," *AI Mag.*, vol. 16, no. 2, p. 92, 1995.

**Zeyu Zhu** received B.S. degree in computer science from Peking University, Beijing, China, in 2019, where he is currently pursuing the Ph.D. degree with the Key Laboratory of Machine Perception (MOE), Peking University. His research interests include intelligent vehicles, reinforcement learning, and machine learning.

**Huijing Zhao** received B.S. degree in computer science from Peking University in 1991. She obtained M.E. degree in 1996 and Ph.D. degree in 1999 in civil engineering from the University of Tokyo, Japan. From 1999 to 2007, she was a postdoctoral researcher and visiting associate professor at the Center for Space Information Science, University of Tokyo. In 2007, she joined Peking University as a tenure-track professor at the School of Electronics Engineering and Computer Science. She became an associate professor with tenure on 2013 and was promoted to full professor on 2020. She has research interest in several areas in connection with intelligent vehicle and mobile robot, such as machine perception, behavior learning and motion planning, and she has special interests on the studies through real world data collection.