# Multi-Agent Reinforcement Learning for Autonomous Driving: A Survey

Ruiqi Zhang[1,2], Jing Hou[1], Florian Walter[3], *Member, IEEE*, Shangding Gu[4], Jiayi Guan[1], Florian Röhrbein[5], *Senior Member, IEEE*, Yali Du[6], Panpan Cai[7], Guang Chen[1,4,*], *Member, IEEE*, and Alois Knoll[4], *Fellow, IEEE*

*Abstract*—**Reinforcement Learning (RL) is a potent tool for sequential decision-making and has achieved performance surpassing human capabilities across many challenging real-world tasks. As the extension of RL in the multi-agent system domain, multi-agent RL (MARL) not only need to learn the control policy but also requires consideration regarding interactions with all other agents in the environment, mutual influences among different system components, and the distribution of computational resources. This augments the complexity of algorithmic design and poses higher requirements on computational resources. Simultaneously, simulators are crucial to obtain realistic data, which is the fundamentals of RL. In this paper, we first propose a series of metrics of simulators and summarize the features of existing benchmarks. Second, to ease comprehension, we recall the foundational knowledge and then synthesize the recently advanced studies of MARL-related autonomous driving and intelligent transportation systems. Specifically, we examine their environmental modeling, state representation, perception units, and algorithm design. Conclusively, we discuss open challenges as well as prospects and opportunities. We hope this paper can help the researchers integrate MARL technologies and trigger more insightful ideas toward the intelligent and autonomous driving.**

*Index Terms*—**Multi-agent reinforcement learning, autonomous driving, artificial intelligence**
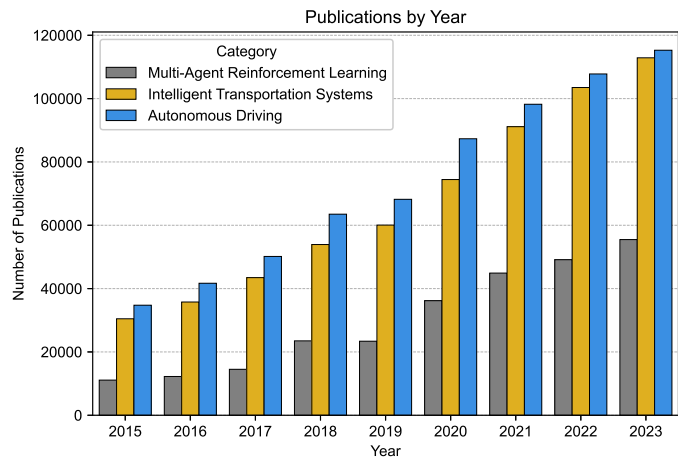
Fig. 1. The number of keywords *multi-agent reinforcement learning*, *autonomous driving* and *intelligent transportation systems* publications from 2015 to 2023 (from Dimension AI [14]). These three research topics are in rapid development and obtaining increasing attention from academia.

## I. INTRODUCTION

LARGE-SCALE autonomous driving systems have attracted tons of attention and millions of funding from industry, academia, and government in recent years [1], [2]. The motivation behind developing such a system is to replace human drivers with automated controllers. It can significantly reduce the time consumption and workload of driving, enhance the efficiency and safety of transportation systems, and promote economic development. Generally, to detect the vehicle states and generate reliable control policies, automated vehicles (AVs) should be equipped with massive electric units, like visual sensors including radars, light detection and ranging (LiDAR), RGB-Depth (RGB-D) cameras, event cameras, inertial measurement units (IMU), global positioning system (GPS) and so on [3]–[5]. A salient challenge in this topic is to build a robust and efficient algorithm that is capable of processing massive information and translating this data into

real-time operations. Early works divide this big issue into perception, planning, and control problems and solve them independently, known as modular autonomous driving.

On the other hand, as a powerful toolkit for sequential decision-making, reinforcement learning (RL) can optimize agent behavior models with the reward signal. As its evolution, deep RL combines the advantages of RL and deep neural networks, which enable to abstract complex observations and learn efficient feature representations [6]. In the past representative research, it has exhibited performances in domains such as board games [7], [8], video games [9], [10] and robotic control [11]–[13], where it rivaled or even surpassed human performances. For autonomous driving, RL brings end-to-end control into reality, which transitions directly from what the vehicle senses to what the vehicle should do, like human drivers. While RL has obtained many remarkable achievements on AVs, most of the related work has approached the issue from the perspective of individual vehicles, which leads to self-centric and possibly aggressive driving strategies, which may cause safety accidents and reduce the efficiency of transportation systems.

For real-world traffic systems, we typically define them as multi-agent systems (MAS) and aim to optimize the efficiency of the entire system rather than merely maximizing individual interests. In MAS, all agents make decisions and interact within a shared environment. This means that the states of each agent depend not only on its actions but also on the

[1]Ruiqi Zhang, Jing Hou, Jiayi Guan and Guang Chen are with Tongji University, Shanghai, China. [2]Ruiqi Zhang is with University of California, Berkeley, United States. [3]Florian Walter is with Machine Intelligence Lab, University of Technology Nuremberg, Germany. [4]Shangding Gu, Guang Chen and Alois Knoll are with Technical University of Munich, Germany. [5]Florian Röhrbein is with Chemnitz University of Technology, Germany. [6]Yali Du is with King's College London, United Kingdom. [7]Panpan Cai is with Shanghai Jiao Tong University, China. *Corresponding author: Guang Chen (mail to: guang@in.tum.de)
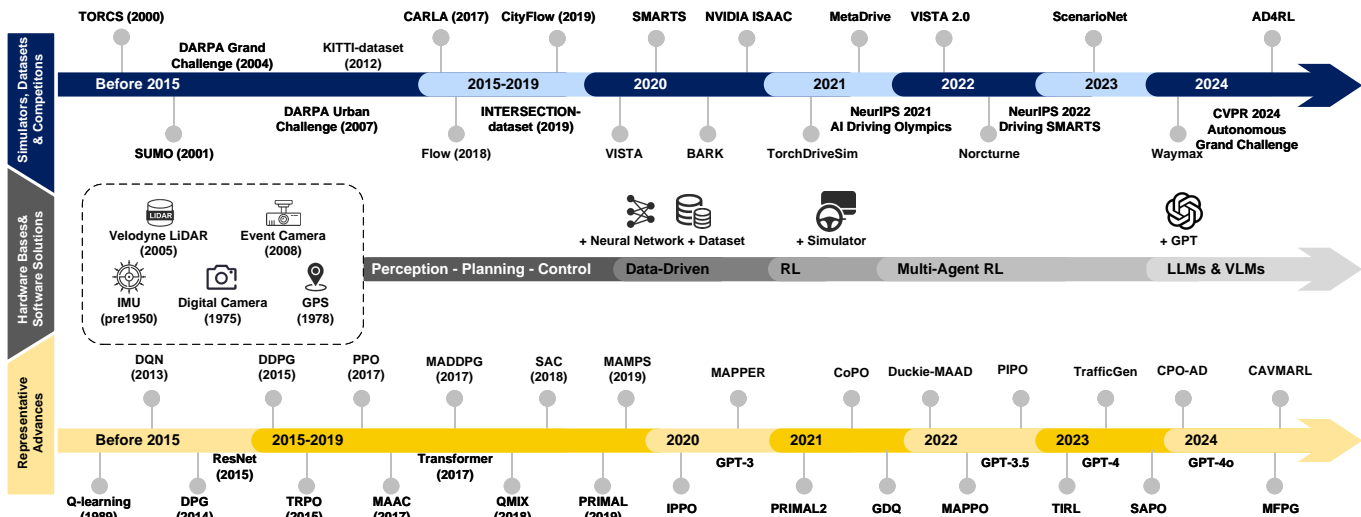
Fig. 2. Timeline of the evolution and representative studies of autonomous driving, deep learning and RL. Based on existing hardware, early research was conducted with a hierarchical scheme, i.e. perception-planning-control. Around 2014, with the rapid development of deep learning and the emergence of datasets, data-driven methods became the mainstream for a time. From 2015 to 2019, numerous RL algorithms appeared, and people realized the opportunities of end-to-end control through simulators. In 2017, MADDPG [24] introduced single-agent RL into multi-agent systems, leading to massive subsequent research on how to control large-scale autonomous vehicles through MARL. In 2020, ChatGPT was launched and made language models receive massive attention.

behaviors of others, making the environmental dynamics non-stationary and time-variant. Additionally, depending on the task setup, agents may either cooperate or compete with each other. In such complex scenarios, manually programming a priori actions is nearly impossible [15]. Thanks to significant advancements in multi-agent reinforcement learning (MARL), substantial breakthroughs have been achieved in traffic control [16], [17], energy distribution [18], [19], large-scale robotic control [20], [21], and economic modeling and prediction [22], [23]. Figure 1 illustrates the number of publications on these related research topics. Using the Dimensions database AI search [14], we searched for keywords including *multi-agent reinforcement learning*, *autonomous driving*, and *intelligent transportation*. The statistical results show that academia is highly interested in these issues, and the related research field is in a rapid growth phase. To accelerate further research and assist new researchers in getting started quickly, we have reviewed over 200 publications, open-source software, and code bases, and then systematically summarized existing achievements and the latest advancements in this paper.

Here, we note other recent reviews. In the milestone series [25]–[27], the authors summarized the blueprint from the history to the future and introduce the influential algorithms in autonomous driving briefly. There are also many surveys [28]–[30] that introduce the basic theory and application of RL and analyze the state-of-the-art (SoTA) algorithms for autonomous driving at their published time, but they mainly focus on single-agent learning. Authors of survey [31] first defined the hierarchical-structured autonomous driving system and limited their scope to local motion planning. They illustrated the kinetics of vehicle and demonstrated how sampling and search-based approaches work mathematically. However, they ignored the contribution of learning-based methods. In the recent survey of motion planning [2], researchers comprehensively investigated the pipeline and learning methods including deep

learning, inverse RL and imitation learning, and MARL. Similarly, a detailed overview covers the latest taxonomy and methodology in trajectory prediction [32]. There are some excellent reviews that summarized MARL approaches for AVs [1], [33], [34]. Nevertheless, researchers have made significant progress in theory and application, as well as in advanced robotic simulators in recent years. As a crucial component of online RL training, simulators determine the sim-to-real gap, that is, whether the policy learned by the agents can be easily transferred to physical robots. Consequently, to enable engineers and researchers to capture the latest updates and accelerate technological progress, we comprehensively summarize the technologies, challenges, and prospects in this field. Generally, the main contribution of this paper can be summarized as follows:

- We propose a series of criteria of benchmarks, analyzing and summarizing the features of advanced simulators, datasets and competitions of large-scale autonomous driving in detail.
- We categorize the state-of-the-art MARL methodologies, and review their technical improvements, insights and unsolved challenges in this field comprehensively.
- We capture the latest advances from relevant fields and delve into future directions of MARL-based autonomous driving from multiple perspectives and dimensions.
- We publish and maintain the GitHub repository[1] to continuously report and update the latest studies in MARL-based autonomous driving, intelligent transportation systems and other relevant areas.

In Fig. 2, we visualize the history of development in MARL, datasets, simulators, hardware and software for autonomous driving, and other related fields. In general, with the development of large-scale datasets and deep learning, autonomous driving has stepped from hierarchical control to the data-driven

[1][Online]. Available: https://github.com/ispc-lab/MARL4AD

era. With the emergence of advanced simulators, the RL-based method steps onto the stage, and then new techniques like large language models bring more opportunities. We will analyze them later in detail, and the rest of this article is organized as follows: In Section II, we first describe the metrics of benchmarks. We also analyze features of the most advanced autonomous driving simulators and datasets. In Section III, we recall the basic concepts, definitions, and open issues in RL and MARL. In Section IV, we exhaustively introduce the SoTA MARL algorithms for autonomous driving. Specifically, we analyze their state and action setup, methodological insights and applications. In Section V, we point out the existing challenges and give out the possible solutions. In Section VI, we capture the latest advances and propose promising directions toward safer and more intelligent autonomous driving.

## II. AUTONOMOUS DRIVING BENCHMARKS

RL is always data-hungry. Generally, it requires continuous interaction with the environment to obtain behavior trajectories, which facilitates more accurate value estimations from deep neural networks [35], [36]. However, due to the economic damage caused by uncertain exploration processes, we typically would not deploy our RL policies on real robots directly. Consequently, within the RL paradigm, data from real driving and high-fidelity simulators are ubiquitously adopted in the development of RL-based autonomous driving. In this section, we will introduce various data sources for large-scale MARL in autonomous driving and traffic systems.

### A. What is important for a good benchmark?

A benchmark involves the simulation of physical models, optical rendering, environment and interaction mechanisms, algorithms, and other complex tasks. For MARL-based autonomous driving, we identify the following crucial criteria. We list them out here and will analyse existing data resourses by these metrics.

*1) Realism and Fidelity:* High realism and fidelity ensure that the simulator can accurately replicate real-world driving conditions like weathers, lights, environmental dynamics, etc., and mitigate their distributional bias before real-world deployment. Deep learning models especially deep RL demand accurate data. In this case, the realism simulator ensures that its data is representative of real-world scenarios.

*2) Scalability:* Scalability ensures that simulators can handle the dynamic environments with numerous entities and variables, so that we can mimic the complexity of real-world scenarios. Scalable MARL algorithms can efficiently learn from the interaction among time-variant numbers of agents like the real traffic system.

*3) Diversity:* Diverse scenarios ensure that vehicles can be tested in a wide range of situations, including different traffic conditions, weather, and road structures. For deployment in the real world, it also contributes to the development of more robust autonomous driving systems and makes AVs more reliable in unpredictable real-world conditions. Simultaneously, diversity also implies that the agents could be heterogeneous, which brings more complex interactions and matches to the realities of actual traffic systems.

*4) Efficiency:* Time and computational resources are both significant concerns for MARL and autonomous driving [37], [38]. Lightweight simulators reduce computational consumption, and experiments work on cheaper and smaller hardware. Meanwhile, highly-parallelized simulator allows multiple environments to run concurrently and promotes the training process of MARL algorithms. Note that there is always a hard trade-off between fidelity and efficiency [39]. High fidelity requires complex computations to replicate real-world scenarios, especially for 3D visual information.

*5) Transferability:* Transferability requires the simulator to support various sensors technically and to replicate their characterisics. For autonomous driving tasks, there are significant differences in the data formats and frame rates of LiDAR, IMU, and cameras. It is essential for the simulator to maintain consistency in the sensor parameters of the intelligent agents with those of commercially available devices. Moreover, transferability is also presented in the simulator's compatibility with the vehicle's device in terms of programming language, communication protocols, and computing platforms.

*6) Features, Maintenance and Supports:* Reproducing the effectiveness of algorithm is always time-consuming. Therefore, it would be beneficial for developers to provide fundamental and verified baselines for testing, also with user-friendly application programming interfaces (APIs), annotations and tutorial documentation. These provisions would establish a fair and open comparison standard and improve the efficiency of subsequent developments. Furthermore, lasting maintenance is necessary. Hardware and software frequently discard old features and develop new ones, which can make data and code outdated. Therefore, continuous maintenance of datasets and code bases is an important task.

### B. Advanced Simulators

The selection of a simulator for MARL-based autonomous driving is a critical step. A good selection would save resources and enable the generation of vast amounts of diverse and high-quality training data, which is essential for effectively training MARL algorithms. Additionally, the simulator allows for parallelized testing and training and significantly accelerates the development process by reducing the time required to experiment with various scenarios and conditions. The advantages on extensive data generation and enhanced time efficiency make simulators indispensable for advancing autonomous driving technologies through MARL.

*1) The Open Racing Car Simulator: TORCS* [40] was first released in 2000. As a highly modular simulator for multi-agent racing, each race car offers low-level APIs to access partial vehicle states and provides visual information from multiple perspectives. After decades of evolution, it has accurate and editable vehicle dynamics, including the rotational inertia of different components, mechanical structures, tire dynamics, and a simplified aerodynamic model. The simulator supports discrete-time simulations with high frequency up to 500Hz, which allows for the development of complex, high-speed, and aggressive driving controllers on it [41]. Many representative competitions and RL-based research conduct

their research on TORCS [42]–[44]. Afterward, *Gym-TORCS* was released and aligned to OpenAI Gym [45], which provides unified APIs and a Python wrapper to facilitate the rapid development of RL-based controller [46]. However, it still lacks support for MARL and a paralleled environment. Later, *MADRaS* [47] filled this vacancy and offered both single-agent and multi-agent environments and interfaces, which could be used to test autonomous vehicle algorithms both heuristic and learning based on an inherently multi-agent setting.

*2) Simulation of Urban Mobility: SUMO* [48] has become a famous benchmark for the simulation of large road networks. It supports a wide range of scenarios and rich APIs so that users can customize traffic scenarios easily. Meanwhile, SUMO provides well-documented tutorials to assist users in implementing their simulations. In recent years, numerous studies have utilized this simulator to develop efficient and safe MARL algorithms for complex scenarios and obtain notable achievements [49]–[51]. To accelerate RL research, developers released a new simulator *Flow* [52] with interfaces to the distributed RL framework RLLIB [53] to achieve high-frequency traffic flow simulation. At the same time, it permits the integration with Amazon Web Services (AWS) elastic compute cloud and expands the variety of controllers, which brings higher flexibility and enables the training of large-scale RL policies. In *CityFlow* [54], developers improved their computational speed to 20 times higher than SUMO. Hence, the real-time simulation of city-level traffic networks becomes possible. CityFlow also expands interfaces for MARL algorithms and allows external data import, which means it can simulate accurate data and generate nearly authentic samples for policy optimization.

*3) Scalable Multi-Agent Reinforcement Learning Training School: SMARTS* [55] is one of the most advanced traffic simulators proposed by Huawei Noah's Ark Lab. It is established on the Social Agent Zoo platform and provides various heterogeneous agent assets for structured traffic flow. SMARTS also has Gym-standardized APIs and integrates broader MARL libraries like PyMARL [56], MALib [57] and RLLIB. Moreover, it supports SUMO as the background provider but optimizes its vehicle dynamics with a Bullet-based physical engine [58]. Implementing high-speed distributed computation introduces a bubble mechanism, which allows the elastic assignment of computational resources on local or remote machines. Furthermore, it offers a strong visualization toolkit through web streaming and allows developers to monitor the process from anywhere. Unlike other one-off works, SMARTS has established a substantial and stable community and promotes many promising works [59]–[62]. So far, its developers have maintained and constantly expanded the simulator's functionalities.

*4) MetaDrive: MetaDrive* [63] is one of the latest multi-agent system simulators based on Panda3D [64], possessing a broad asset library and allowing the import of external data. Beyond the given structured scenarios, developers can easily customize road map and traffic flows, and set the attributes of scene components via high-level APIs. Meanwhile, it establishes hierarchical management of assets by defining four types of manager classes, which facilitates developers in customizing agent mixtures, interactions and policy generalizability tests. Theoretically, MetaDrive can create an infinite variety of traffic scenarios. It employs state vectors as the agents' observations and provides abundant RL benchmarks, including model-free RL [65], [66], imitation learning [67], and offline RL [68]. Although it forsakes fine-grained visual information, it enables more rapid and efficient simulation and has triggered many insightful works [20], [69].

*5) CAR Learning to Act: CARLA* [70] is one of the state-of-the-art open-source 3D simulators for its realistic Unreal Engine 4-based dynamics simulation, lightweight optical rendering, and comprehensive technical support. For environmental information, it provides detailed scenes with various architectures, road configurations, and natural conditions, especially diverse weather settings, which are beneficial for the generalizability test of policy. Through its free asset library, it can simulate high-density traffic flow, pedestrians with different behaviors and traffic lights and signs, which makes it possible for automated agents to comprehend traffic regulations. Significantly, CARLA supports various sensors including LiDAR, RGB-D cameras, GPS, radar, and event cameras with editable characteristics and realistic noise. With rich C++ and Python-based APIs, researchers can freely define the attributes of agents and then analyze the collected data. Nowadays, CARLA not only makes substantial contributions to the MARL field but also continuously impacts computer vision research. A welcome trend is that its developers have established an official website with a good tutorial, blog, and user community for further updates and improvement. As a good supplement, based on CARLA, some researchers have developed a new benchmark named *MACAD* [71] with a faster implementation for MARL and integrated Gym-like APIs. Generally, CARLA significantly facilitates systematic research [72]–[74] for vision-based driving and reduces the simulation-to-reality (sim-to-real) gap, which is crucial for deploying the large-scale driving controller.

*6) Virtual Image Synthesis and Transformation for Autonomy: VISTA* [75] is a data-driven simulator and synthesizes time series of perceptual inputs from real-world. In contrast to physical simulators, VISTA aims to reconstruct that world and synthesize novel viewpoints within the environment via inputting real data of the physical world. Different sensing modalities, environments, dynamics, and tasks with varying complexity are supported. Meanwhile, it is highly modular, customizable, and extensible. Since the behavior trajectories are generated from real data, the sim-to-real gap is minimized, which is empirically validated by researchers in real-world experiments. In the subsequent *VISTA 2.0* [76], researchers integrate UNet architecture [77] to reproduce the dense output of LiDAR and apply temporal interpolation to estimate the events through RGB images, providing an additional two sensor types and data formats for AVs. Simultaneously, researchers also introduce a version that supports multi-agent interactions and validates basic scenarios involving multiple AVs [78]. However, its application in large-scale complex scenarios remains unexplored.

*7) NVIDIA ISAAC-Sim: ISAAC-Sim* [79] is established on the PhysX5 engine and supports GPU-based photorealism

TABLE I
A SUMMARY OF THE ADVANCED AUTONOMOUS DRIVING AND MARL SIMULATORS

| | Simulators | Year | Last Update | State Description | MARL Support | Link |
|---|---|---|---|---|---|---|
| TrafficFlow-Oriented | SUMO [48] | 2001 | 2024-07 | High-level states | ✔ | https://eclipse.dev/sumo/ |
| | Flow [52] | 2018 | 2019-08 | High-level states | ✔ | https://flow-project.github.io |
| | Highway-env [84] | 2018 | 2024-04 | Binary image, grid map, high-level states, etc. | ✔ | https://github.com/Farama-Foundation/HighwayEnv |
| | CityFlow [54] | 2019 | 2020-12 | High-level states | ✔ | https://cityflow-project.github.io/ |
| | BARK [85] | 2020 | 2022-06 | Roadmap, high-level states, etc | ✗ | https://github.com/bark-simulator/bark |
| | MADRaS [47] | 2020 | 2020-10 | Image stream, high-level states | ✔ | https://github.com/madras-simulator/MADRaS |
| | SMARTS [55] | 2020 | 2024-07 | BEV image, grid map, LiDAR, high-level states, etc. | ✔ | https://github.com/huawei-noah/SMARTS |
| | MetaDrive [63] | 2021 | 2024-06 | RGB camera, high-level states, etc. | ✔ | https://github.com/metadriverse/metadrive |
| | TBSim [86] | 2021 | 2023-10 | BEV image, etc. | ✗ | https://github.com/NVlabs/traffic-behavior-simulation |
| | TorchDriveSim [87] | 2021 | 2024-07 | BEV image, etc. | ✗ | https://github.com/inverted-ai/torchdrivesim |
| | Intersim [88] | 2022 | 2023-06 | Rasterized image, vectorized states | ✗ | https://github.com/Tsinghua-MARS-Lab/InterSim |
| | Nocturne [89] | 2022 | 2022-10 | High-level states | ✔ | https://github.com/facebookresearch/nocturne |
| | ScenarioNet [90] | 2024 | 2024-05 | RGB camera, high-level states, etc. (depends on external datasets) | ✔ | https://github.com/metadriverse/scenarionet |
| | Waymax [91] | 2024 | 2024-03 | Roadmap, Bounding boxes, high-level states, etc. | ✔ | https://github.com/waymo-research/waymax |
| Fidelity-Oriented | TORCS [40] | 2000 | 2020-02 | Image Stream | ✗ | https://sourceforge.net/projects/torcs/ |
| | Gym-TORCS [46] | 2017 | 2017-02 | Image stream, high-level states | ✗ | https://github.com/ugo-nama-kun/gym_torcs |
| | CARLA [70] | 2017 | 2024-07 | RGB-D & event camera, LiDAR, high-level states, etc. | ✔ | https://github.com/carla-simulator/carla |
| | MACAD [71] | 2020 | 2023-01 | RGB camera | ✔ | https://github.com/praveen-palanisamy/macad-gym |
| | ISAAC Sim [92] | 2020 | 2024-06 | RGB-D camera, LiDAR, high-level states, etc. | ✔ | https://developer.nvidia.com/isaac/sim |
| | Vista [75], [76] | 2020 | 2022-05 | RGB-D & event camera, LiDAR, etc. (depends on external datasets) | ✔ | https://github.com/vista-simulator/vista |

with real-time ray tracing and the MDL material definition for physically based rendering, which provides nearly real scenarios for robotic training and testing. Besides its rich asset library, it also permits developers to import mesh models, customized assets, and URDF format robot structure files, which significantly enhance the flexibility in robotic research. More importantly, it facilitates communication with other important tools like ROS [80], ROS2 [81], and Gazebo [82], simplifying the real-world deployment. In the latest version, human simulation is introduced and supports the simulation for available LiDAR brands like Ouster, Hesai, and Slamtec. Ray-traced (RTX) technology, can feedback more accurate LiDAR data w.r.t. various reflective materials or under different conditions. However, the high-fidelity optical and physical simulations increase the hardware load and computational cost so RTX-based GPU is required. Another option is the more lightweight *NVIDIA Isaac Lab* [83]. Compared to Isaac Sim, it concentrates on the training of RL policies and supports large-scale, high-throughput multi-GPU distributed computations.

In table I, we summarize the key features of these representative simulators. Overall, we categorize simulators into two major types: traffic flow-oriented and fidelity-oriented. Traffic flow-oriented simulators prioritize parallelization and compu-

tational efficiency, often at the expense of accurate dynamics modeling and rendering precision, making them suitable for large-scale vehicle simulations. In contrast, fidelity-oriented simulators offer better optical rendering and dynamic accuracy with a smaller sim-to-real gap, but they consume significantly more computational resources. Specifically, we give out their website and code sources and list the supported sensors and information and note that most simulators provide abstracted observations, such as vectorized lane direction or the collision distance to the vehicle ahead. Here, we collectively refer to these processed features and representations along with the vehicle's kinematic information (like speed, acceleration or orientation) as the *high-level states*. We also highlight the simulators that are still being maintained, and their updates are worth following. Researchers and practitioners should choose among these platforms based on their specific project needs, the complexity of the traffic scenarios, and the desired level of realism and interactivity in the simulations. We also notice that many unmentioned studies are showing fancy single-agent results on game engines [93]–[95] and we acknowledge their contributions, but they can hardly adapt to MARL paradigm so they will not be covered in our survey.

## C. Datasets

After decades of development, to effectively solve real-world driving problems, developers have collected tons of on-road data and established rich repositories with different sensors in various scenarios, such as KITTI [96], nuScenes [97], and Waymo Dataset [98]. However, for decision-making problems, these datasets do not record real-time actions like accelerating, braking, steering, or actuator outputs, so they are typically used only for visual tasks or multi-modal sensor fusion. For example, although the IMU and GPS history is given out in BDD-100K [99], the information is on a different domain from human driver actions, and further policy adaption and transferring would be difficult. Moreover, most datasets address the autonomy of a single vehicle rather than collaboration between multiple AVs. Collecting and aligning multi-vehicle data simultaneously is often expensive and difficult, so nowadays MARL paradigms are mostly verified in simulators. However, real vehicle decision data is still invaluable especially for imitation learning [67] and offline RL [68]. Compared to data obtained in simulators, although theoretically we cannot interact for infinite time like simulation, the data distribution is closer to real driving.

To provide diverse interactive data for sequential decision-making, the INTERACTION dataset [100] investigates driving habits across different cultural backgrounds and collects bird's-eye-view data via hovering drones with fine-grained annotations. Based on the nations and road structures, it compiles 11 sub-datasets including highways, roundabouts, intersections, merges, and unstructured roads. Additionally, the dataset includes rare collision data and aggressive driving behaviors. Afterwards, the latest research proposes a new benchmark AD4RL [101] based on Next-Generation Simulation (NGSim) US-101 dataset [102]. This benchmark provides 19 datasets from real-world human drivers after fine correction, value normalization, and alignment with partially observable Markov decision process. In other words, it can directly work as policy trajectories in the RL scheme. Additionally, it offers 7 popular offline RL algorithms applied in 3 realistic driving scenarios. A unified decision-making process model is also attached for verification across various scenarios on Flow simulator [52], which serves as a reference framework for algorithm design. Recently, the development and supplement of datasets for the offline MARL approach in large-scale autonomous driving and traffic systems remain a work in progress. We consider this an up-and-coming area and will discuss it later (see Section VI-B).

## D. Competitions

We notice that many recent competitions have been hosted during the international conference sessions, and here we address and appreciate the contribution of their organizers. These competitions provide a platform for researchers to showcase and compare their algorithms, drive the emergence of new techniques, and promote the application of MARL in real-world scenarios.

The earliest multiple vehicles involved in competition can be traced back to the DARPA's Urban Challenge [103]. Although this challenge involved controlling only one vehicle, it required real-time controller adjustment based on a dynamic environment. In the past five years, an increasing number of companies and institutions have recognized MARL's potential and organized academic competitions based on relevant simulators and datasets. At DAI 2020, a multi-vehicle control competition was organized using the SMARTS [55] simulator. Participants were required to develop a parameter-sharing multi-agent model to control a group of agents to accomplish short missions in ramp, double merge, T-junction, crossroads, and roundabout. Instead of testing in the simulation, the DuckieTown [104] AI Driving Olympics (AI-DO) competition on NeurIPS 2021 required participants to deploy real vehicles on scaled-down tracks. The vehicles had to navigate complex urban roads, avoiding pedestrians and other vehicles. The AI-Do competition marks a milestone in the practical deployment of embodied intelligence in complex MAS. At NeurIPS 2022, Huawei's Noah's Ark Lab again organized competitions based on the SMARTS simulator, featuring online and offline reinforcement learning tracks. OpenDriveLab consecutively hosted multi-track autonomous driving competitions on CVPR 2023 and 2024. In the latest Autonomous Grand Challenge, the organizers provided the offline end-to-end autonomous driving scale competition and the online CARLA competition. For the former one, participants were required to develop motion planning algorithms for complex scenarios using offline data from Motional nuPlan [105] dataset. The second track required the design of flexible policy learning methods in the CARLA [70] simulator. These competitions establish standardized benchmarks for evaluating different algorithms and offer a unified framework to compare performance and identify the most effective solutions.

## III. REINFORCEMENT LEARNING PRELIMINARIES

In this section, we will introduce the fundamental definitions and presentations of RL and MARL. We want to make sure the readers can comprehend the rest of this paper under a unified mathematical language.

### A. Deep Reinforcement Learning

Reinforcement Learning is a classical approach to sequential decision-making problems. We typically use the Markov decision process (MDP) to describe the decision-making procedure in RL paradigm, where the state distribution at the next time-step is only determined by the action and state of the current time-step and irrelevant to its history. Specifically, MDP can be denoted as a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{T}, \gamma)$, where $\mathcal{S}$ and $\mathcal{A}$ present the state space, action space and the stochastic observation space. $\mathcal{T} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to [0, 1]$ is the transition probability under a given state-action pair and includes the uncertainty of the system. $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to r$ indicates the reward function which issues immediate reward value $r$ at each time-step and $\gamma$ presents the discount factor in the optimization objective. For a behavior trajectory, Equation (1) accumulates the discounted reward value $r_i$ at each time-step and our objective is to find an optimal policy to maximize the total return $\mathcal{G}_t$.

$$\mathcal{G}_t = r_{t+1} + \gamma r_{t+2} + \cdots = \sum_{k=0} \gamma^k r_{t+k+1} \qquad (1)$$

$$\mathcal{V}(s_t) = \mathbb{E}_{s \sim \mathcal{S}} \left[ \sum\nolimits_{k=0} \gamma^k r_{t+k+1} \right] \quad (2)$$

$$= \mathbb{E}[r_t + \gamma V(s_{t+1})|s_t = s] \quad (3)$$

$$\mathcal{Q}(s_t, a_t) = \mathbb{E}_{s \sim \mathcal{S}, a \sim \pi_\theta} [r_t + \gamma Q(s_{t+1}, a_{t+1})|s_t = s, a_t = a] \quad (4)$$

However, due to the stochastic environmental dynamics, we usually cannot quantify the quality of a state or action directly. For this issue, a common solution is estimating the expected subsequent reward given a specific state using the Bellman equation, referred to as the state value. The Eq. (3) defines $\mathcal{V}$-function and evaluates the expected accumulative reward (i.e., return) of all subsequent states, which enables us to estimate state values through the dataset (or interactive trajectories) collected from simulation and real-world. Another form of the Bellman Equation uses state-action pair $(s_t, a_t)$ and estimates the state-action value with $\mathcal{Q}$-function as Eq. (4). Here, the objective we evaluate is the expected return with a given state and action $a_t \sim \pi_\theta(s_t)$ sampled from policy $\pi$ with learnable parameters $\theta$. Essentially, the policy is a distributional function of actions with a given state. In deep RL, the $\mathcal{V}$-function, $\mathcal{Q}$-function, and policy $\pi$ are represented by neural networks. In general, the Bellman equation presents the accumulative reward given a behavior trajectory, and establishes the connection with the ground truth and the value estimation, so that we can the error between them to learn a better state evaluation via deep learning.

Meanwhile, to learn the policy safely while avoiding robot damage, we collect data (i.e., behavior trajectories) in the simulator and identify the policy that maximizes the accumulative reward. Based on how the policy learns from the collected data, the mainstream deep RL methods can be categorized into *on-policy* and *off-policy* RL. On-policy RL learns the value of the policy being carried out by the agent, meaning the policy used to make decisions is the same as the policy being improved. The advantage of on-policy RL is that it directly updates the target policy with the latest data and is more robust under sparse reward function [36]. However, it requires more exploration during training as it must generate diverse experiences with the current policy. Contrarily, off-policy RL learns from different policies, which could be the old policy [65] or a separate exploratory policy. Hence, it's more sample efficient as it can use data from multiple policies [36].

Another taxonomy is from their prerequisite and assumptions. Ideally, when the information of the environment during MDP can be fully known, Dynamic Programming [106] can solve the Bellman Equations via policy iteration or value iteration [6]. If the environment dynamics are partially known or unknown, the Monte Carlo method [6] estimates the unbiased value by implementing possible trajectories from the current state and calculating the return as Equation (5).

$$\mathcal{V}(s_t) \leftarrow \mathcal{V}(s_t) + \alpha \left[ \mathcal{G}_t - \mathcal{V}(s_t) \right] \quad (5)$$

$$\mathcal{V}(s_t) \leftarrow \mathcal{V}(s_t) + \alpha \left[ \mathcal{G}_t + \lambda \mathcal{V}(s_{t+1}) - \mathcal{V}(s_t) \right] \quad (6)$$

Similarly, bootstrapping provides another estimation-based solution Temporal Difference (TD) method. It predicts the future state from current time-step and as Equation (6), we show the simplest method TD(0), where $\mathcal{G}_t + \lambda \mathcal{V}(s_{t+1}) - \mathcal{V}(s_t)$ is TD-error at time $t$. In this way, TD-error facilities the update of the objective function in either on-policy form like SARSA [107] or off-policy form like Q-learning [108]. Different from the value function-based methods above, another class of methods focuses on directly optimizing the parameters of the policy network, known as policy-based learning. In RL, the policy $\pi_\theta(a|s^*) = P(a|s^*, \theta)$ is essentially composed of a set of parameters $\theta$, which determines the probability of selecting an action under a given state $s^*$.

$$\mathcal{J}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum\nolimits_{\Psi_t \in \tau} \Psi_t \right] \quad (7)$$

$$\nabla_\theta \mathcal{J}(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} \left[ \sum\nolimits_{(s_t, a_t, \Psi_t) \in \tau} \log \left( \pi_\theta(a_t|s_t) \right) \Psi_t \right] \quad (8)$$

Accordingly, we aim to gain the optimal $\theta$, which maximizes our optimization objective in Equation (7), where $\Psi_t$ is a designable reward-related variable. In the simplest design, it is the reward itself. The roll-out trajectories sampled by policy $\pi_\theta$ are denoted as $\tau$. Under the mild conditions [109], we can obtain the gradient w.r.t. parameter $\theta$ via the key observation [110] as Equation (8). The REINFORCE [111] algorithm utilizes $\Psi_t = Q(s_t, a_t) - b$ with a learned coefficient $b$ to normalize Q-values. Similarly, the PPO algorithm [65] proposes introducing the advantage function to modify and optimize the local approximation of the real return. Compared to value-based methods, policy gradient approaches circumvent curse of dimensionality arising from discretization in continuous action spaces. Essentially, policy gradient methods generate a probability distribution over actions applicable to discrete and continuous action spaces.

The actor-critic architecture learns the policy and value function simultaneously. The actor (policy) makes decisions based on the current policy, and the critic (value function) evaluates the action selected by the actor. Actor-critic methods strike a balance between the high-bias value-based methods and high-variance policy-based methods and can efficiently operate in high-dimensional environments [66]. Due to space constraints, this paper will not undertake a deep dive into these categories and the implementation details. However, they are available in the reference if needed.

### B. Multi-Agent Deep Reinforcement Learning

According to the assumption and agent settings, we can introduce diverse probabilistic models to represent the tasks. As an extension of MDP, Markov Game (MG) presents the interaction process of a multi-agent system. MG is denoted a tuple $\left( \mathcal{N}, \mathcal{S}, \{\mathcal{A}^i\}, \{\mathcal{R}^i\}, \mathcal{P}, \gamma \right)$, where $\mathcal{N} = \{1, 2, \cdots, N\}$ denotes the set of interacting agents and $\mathcal{S}$ is the global state from all agents with $i \in \mathcal{N}$. Similarly, $\{\mathcal{A}^i\}$ and $\{\mathcal{R}^i\}$ are the set of individual action and reward. Note that $\mathcal{A}^i$ is the action space of the $i$-th agent so the joint action space is $\mathcal{A} := \mathcal{A}^1 \times \mathcal{A}^2 \times \cdots \times \mathcal{A}^N$. Like single-agent RL, $\mathcal{P}$ and $\gamma$ indicate the transition probability and discount factor. At time step $t$, each agent executes action $a_t^i \sim \pi^i$ according to the system state $s_t$, and then the system transits to the next state $s_{t+1}$ and obtain the reward $r_{t+1}^i := \mathcal{R}^i(s_{t+1}|s_t, a_t^1, \cdots, a_t^N)$. Hence, we can

obtain the value function for each agent like single-agent RL as Equation (9).

$$\mathcal{V}^i_{\pi^i,\pi^{-i}}(s_t) = \mathbb{E}_{a^i \sim \pi^i}\left[\sum_{k=0}\gamma^k r^i_{t+k+1}\right] \qquad (9)$$

$$\mathcal{V}^i_{\pi^i_*,\pi^{-i}_*}(s_t) \geq \mathcal{V}^i_{\pi^i,\pi^{-i}}(s_t), \forall i \in \mathcal{N} \qquad (10)$$

Ideally, we hope that the policy of each agent $\pi^i$ should give out the best response to the joint policy of other agents $\pi^{-i}$ in the MAS. The Nash equilibrium (NE) [112] provides a more specific description for this case. Mathematically, we denote the joint policy in NE case as $\pi^i_* = \{\pi^1_*, \pi^2_* \cdots, \pi^N_*\}$ and then it can be described via value inequality (10). Intuitively, the Nash Equilibrium indicates a situation with a joint policy, where no agent can change its policy to improve the performance unilaterally. However, the solutions representing NE in the state space are not unique in most cases [15]. For a given joint policy $\hat{\pi}$ and any other policy $\pi$, if it satisfies $\mathcal{V}^i_{\hat{\pi}}(s) \geq \mathcal{V}^i_\pi(s)$ for all states $s \in \mathcal{S}$, and there is at least one state $s^*$ making $\mathcal{V}^i_{\hat{\pi}}(s^*) > \mathcal{V}^i_\pi(s^*)$, the policy $\hat{\pi}$ is Pareto-optimal and it Pareto-dominates $\pi$. Accordingly, only a Nash Equilibrium without any other policy with greater value can be regarded as Pareto-optimal.

In autonomous driving, vehicles' observations are limited by the field of view (FoV) of the sensor and geographic location, so that they can only obtain incomplete information of environment. To this end, we typically modify the model by introducing the set $\mathcal{O}_i$ to denote observation space. Hence, We can define it as a decentralized partially observable MDP (Dec-POMDP) with representation in a seven-element tuple $\left(\mathcal{N}, \mathcal{S}, \{\mathcal{A}^i\}, \{\mathcal{O}^i\}, \{\mathcal{R}^i\}, \mathcal{P}, \gamma\right)$. In MARL, the intricacies of inter-agent relationships influence the design of reward functions and training schemes. These relationships are typically classified into competitive, cooperative, and mixed. In competitive settings, agents strive to maximize their objectives, often at the expense of other agents. In other words, reward functions are typically zero-sum, where one agent's gain is equivalent to another's loss. Contrarily, cooperative agents are required to work collectively towards a common goal. In this case, the rewards are shared and distributed among agents based on their contribution to achieving the collective objective. Mixed interactions involve elements of both cooperation and competition; therefore, rewards in mixed scenarios are more complex and require a trade-off between individual and group incentives.

### C. Learning Schemes

With the increasing number of agents, the complexities of state and action spaces would grow up exponentially, which causes the policy learning of MAS present a computational challenge. Generally, the entire process can be divided into two stages: Training and Testing. Training indicates the process where agents acquire data from interaction to obtain experience and update policies. After that, we evaluate the policy performance in the environment without any policy optimization, which is referred to as testing. Based on the classic categorization [113], the training process can be broadly classified into two paradigms: centralized and decentralized.



(a) An example of CTDE scheme
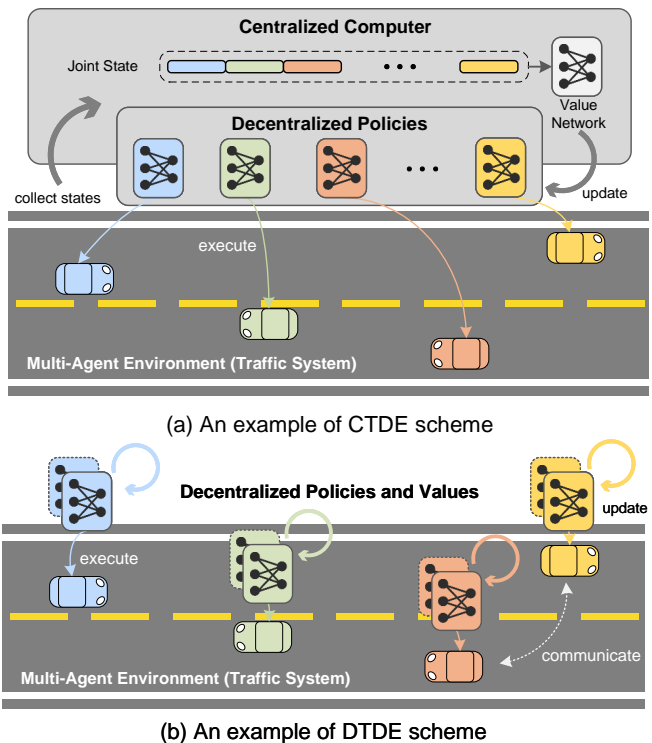


(b) An example of DTDE scheme

Fig. 3. The Centralized Training-Decentralized Execution (CTDE) scheme versus Decentralized Training and Execution (DTDE) scheme. In CTDE paradigm, it utilizes a central controller to concatenate all observations and distribute policies for all agents. In the DTDE paradigm, the agent computes the policy itself. According to hardware settings, they can communicate with each other and exchange information, or detect others' states by sensors.

In centralized training, agents update their policies with communication and information exchange, while additional information is removed during testing. In contrast, decentralized training is conducted in a distributed manner, where each agent independently performs and develops individual policies without extra information and communication among agents. The centralized execution asks all agents to follow the command from one joint policy with an unconstrained and instantaneous information exchange [15]. However, in real MAS, this strong assumption can hardly be achieved, and the agents are distributed various tasks and required to work independently. Therefore, we generally adopt decentralized execution, especially for autonomous driving cases. As shown in Figure 3, mainstream MARL algorithms can be categorized into two paradigms: centralized training with decentralized execution (CTDE) and decentralized training with decentralized execution (DTDE).

*1) CTDE Schemes:* Early centralized learning schemes are created to handle partial observable environments [114]. It assumes that a central controller (critic) collects the states and actions from all agents and updates the policies (actor) for them. Hence, the CTDE scheme with the Dec-POMDP model comes into being and has been applied in a deal of systems [115]–[117]. In this case, we could use single-agent RL algorithms and toolkits to simplify the problem. However, for MAS systems with heterogeneous agents or multiple task objectives, this scheme can hardly learn feasible policies to

meet various requirements.

*2) DTDE Schemes:* In the fully decentralized paradigm, agents are trained and operate independently without access to global information. Therefore, it's more suitable to handle the decision-making and planning tasks for large-scale robotic systems [20], [21]. Classical DTDE approaches were limited by non-stationary and low efficiency, while recent research has significantly alleviated it via online evolution [118], advanced value learning [119], [120] and other techniques. Meanwhile, to deal with insufficient information from other agents, the decentralized MAS is generally established with a communication network where agents can exchange message or observe and estimate others' states. In the AD context, agents are typically configured to communicate with other vehicles within the signal range via Bluetooth or WiFi or to estimate their states through visual sensors.

### D. Issues of MARL

*1) Non-Stationarity:* Non-stationarity is a significant issue for decentralized MAS, where agents interact within a shared environment and update their policies synchronously. Consequently, each vehicle doesn't know the full environment dynamics, and the next state doesn't depend solely on its action and current state so it would break the Markov assumption [15]. In the CTDE paradigm [24], the centralized critic has access to all agents' observations and actions. Since only the actor computes the policy and the critic component can be removed during testing, agents in the CTDE scheme have fully decentralized execution. For independent policy learners, a naive approach is to let agents either ignore the presence of others or proceed under the assumption that the behaviors of others are static [121]. In this context, agents are independent learners, which enable the conventional single-agent RL algorithms for policy learning and have been proven to achieve excellent results on various benchmarks [122]. However, in complex and stochastic environments, independent policy learning may result in sub-optimal performance or tend to exhibit a propensity for over-fitting to the policies of other agents, leading to a lack of generalizability in testing. To improve independent learning performance, researchers propose adopting different learning rates with shared rewards to achieve the optimal joint policy [123]. Another method involves the refinement of experience replay. Due to the non-stationarity of the environment, experience replay may store more irrelevant experiences to decentralized learning with the increasing time steps. Importance sampling corrections for stable experience replay is also a solution, which adjusts the weights between the prior and the new experience under different environment dynamics. This approach has been proven to enhance the performance of independent learners in complex gaming and robotic environments [124]–[126].

*2) Partial Observability:* In partially observable environments, agents do not have access to the full state of the environment. Instead, each agent receives only a local observation that provides incomplete or noisy information about the true state. To make effective decisions, agents must estimate or infer the underlying state of the environment from their partial observations. This often requires maintaining a a probability distribution over possible states, which adds computational complexity and uncertainty to the decision-making process. At the same time, in a partially observable setting, different agents might have different pieces of information about the environment. Coordinating actions effectively requires agents to share information or make decisions based on limited and potentially inconsistent knowledge.

For example, the policy learned with given sensors might require more work to find a workable latent representation and feasible policy from limited observations. To this end, a promising direction is using more complex architectures to enhance the representative capability of neural networks. Recent DRQN algorithm replaces the first layer of the DQN with a recurrent Long Short-Term Memory (LSTM) unit [127]. Combined with the concurrent experience replay trajectories mechanism, it mitigates the information limit on the hysteric Q-learning [128]. Inspired by this, researchers propose a weakly cooperative traffic model under traffic scenarios and apply an independent policy learning algorithm that utilizes a forgetting experience mechanism and a loose weight training mechanism to alleviate both partial observability and non-stationarity [129]. Besides, MADDPG algotithm [24] with recurrent actor-critic has been demonstrated to learn low-variance stable policies in partially observable environments with various constraints of agent-level communication [130].

*3) Credit Assignment:* The credit assignment (CA) problem is one of the crucial challenges in developing CTDE MARL. It involves determining how to allocate the global reward obtained through interaction with the environment [131]. The most widely used method in such tasks is the shared team reward, which may lead to the problem of lazy agents, where only partial agents in the system work hard. Essentially, it's a tricky sub-optimal policy. To address this, Value Decomposition is the most widely-used solution, which divides the joint value function into individual agent value functions and then selects the best action for each agent. A straightforward linear decomposition represents the total reward function as a sum of each agent's reward function. Recent research focuses on solving credit assignments for heterogeneous agents through nonlinear assignments using a learnable weighted network [131]. Another approach uses multi-agent policy gradient algorithms, which impose a monotonic constraint between the joint action value and individual policies [132]–[134]. The recent SoTA method can find the global-optimal policy via polarized policy gradient [135], shortening the distance between multiple non-optimal joint action values. However, this problem still exists in large-scale and complicated heterogeneous robotic control problems.

*4) Scalability:* The joint action space would increase exponentially with the increasing number of agents. Meanwhile, the number and density of agents are changing under specific real-world settings like the traffic flow of intersection, leading to variations in the dimensionality of observed information, which classic multi-layer perceptron (MLP) or convolutional neural network (CNN)-based [136] value networks struggle to deal with. Although we can apply the solutions such as LSTM [127], transformers [137], or dimensionality reduction

[21], centralized methods require substantial computational resources, memory, and bandwidth to calculate and distribute their strategies after receiving action state information from each agent. A possible solution to the curse of dimensionality is independent learning. However, as mentioned previously, this approach must produce consistent results in non-stationary environments and may be limited by partial observability. In the context of autonomous driving, each agent is typically set up to exchange information with specific other nearby agents. For this, researchers propose a distributed Q-learning [138], assuming that each agent knows only its local actions and rewards, but agents can send their Q-values to their nearest peers and update their Q-networks locally. A scalable actor-critic method is established on exponential decay property with average reward within dynamic environments, which can handle the scaling state-action space size of local adjacent agents [139]. In MAMBA [140], researchers utilize model-based RL to sustain a world model [141] for each agent during execution and generate efficient rollouts for training, removing the necessity of interaction with the environment.

*5) Communication Mechanism:* As mentioned before, the communication mechanisms help MARL overcome issues of non-stationarity and partial observability. Besides, a better design of communication mechanisms contributes to reducing the lowest required hardware and bandwidth while enhancing learning efficiency. For the MAS in the real world, the agents are set to communicate with the adjacent or nearby agent. The communication relations can be predefined, changed, or even learned and updated during the process [142], [143]. Meanwhile, we can introduce the proxy to facilitate the collection, processing, and distribution of the message [144]. The message could contain both the past and current observations and actions, and it also can be compressed and processed using sequence networks or autoencoders [145]. In some studies [146]–[148], the intended action or policy distributions can also be shared among agents and serve as additional input information for policy or value networks or both of them [149]. Note that some recent researches regard communication as a learnable process and focus on updating and adjusting communication protocols. Communication learning encompasses the communication policies and the content of messages [142]. It can be implemented through backpropagating gradients from the communicatees [150], [151] or through an independent reinforcement learning thread [152], [153].

## IV. STATE-OF-THE-ART METHODOLOGIES

This section will introduce recently the most advanced MARL methodologies for motion planning and control of multi-vehicle systems. We cannot encompass all the related studies, but select representative techniques in this survey are sourced from reports published in the most influential conferences and journals. Furthermore, we encourage the researchers to report more relevant works to our website.

### A. Centralized Multi-Agent RL

In the CTDE scheme, each vehicle has an independent policy network, and a core computer is set to merge and process the information from all vehicles. We first get the merged observation from all vehicles, evaluate the system state by a pre-defined global reward function, and then train the independent policies after credit assignment. PRIMAL [154] is a milestone work in centralized training for pathfinding. It assigns each agent an independent and fine-designed parameter-sharing actor-critic network and trains them with A3C [155] algorithm. In this work, researchers illustrate that independent policies lead to selfish behaviors, and a hand-crafted reward function with a safety penalty is a good solution. Additionally, there is a switch to allow agents to learn from interaction or expert demonstrations. The combination of reinforcement learning and imitation learning contributes to fast learning and alleviates the negative impact of selfish behaviors on the overall system. In this paper, a discrete grid world is defined, and the local state of each agent is set as the information of a $10\times10$ block with the unit vector directed toward the goal. To verify the feasibility in the real world, the authors also implement PRIMAL on AVs in a factory mockup.

In MADDPG [24], the authors propose the first generalizable CTDE algorithm based on deep deterministic policy gradient (DDPG) [156] with a toy multiple-particles environments. It provides an essential platform with easy vehicle dynamics to learn the continuous driving policies with continuous observation and action spaces under design-free scenarios and attracts many remarkable followers [21], [157]. Meanwhile, the combination of value function decomposition methods and CTDE scheme has achieved better scalability w.r.t. the number of agents and mitigates the impact of non-stationary on policy training, thereby improving performance in large-scale multi-agent systems [116], [158]. These methods have been verified in complex scenarios like unsignalized intersections in Highway-Env [84], [159]. Also, expert demonstration contributes to reducing the risk of converging to sub-optimal policies [159]. To verify the feasibility of deploying the CTDE approach in mapless navigation tasks, Global Dueling Q-learning (GDQ) [160] sets up an independent DDQN [161] for each turtlebot3 in the MPE [24] to train policies and estimate values. Additionally, they introduced a global value network that combines the outputs of the value networks of every agent to estimate the joint state value. This method has been proven to be more effective than normal value decomposition methods. Meanwhile, researchers also attempt to extend fundamental algorithms in single-agent RL such as PPO [65] or SAC [66] to multi-agent tasks and provide many significant baselines like MAAC [162] and MAPPO [163]. In particular, MAPPO has been verified comprehensively on massive benchmarks and has systematic guidance of hyperparameter selection and training. To overcome the sim-to-real gap and deploy MAPPO on real robots, developers train a policy in the Duckietown-Gym simulator for following waypoints on the ground. The MAPPO policy network adopts recurrent neural network [164] to recall the knowledge of the prior state and output the high-level target linear velocity and angular rate for each vehicle. Like most indoor navigation tasks, the optical track system captures the position and attitude of vehicles. With the linearized inverse kinetics, the executive low-level command of the vehicle can be obtained

after domain adaptation. This work reveals how to deploy the CTDE scheme on real robots, and the engineering experience is valuable for future studies.

## B. Independent Policy Optimization

Regarding practical deployment challenges such as communication, bandwidth, and system complexity, the fully decentralized system reduces communication overhead and bandwidth requirements by allowing agents to operate independently without constant coordination. Additionally, it is easier to deploy in environments with limited or unreliable communication infrastructure, lowers decision-making latency, and simplifies local computation for each agent. These factors make decentralized MARL a more practical and adaptable approach for real-world multi-agent applications. In recent years, Independent Policy Optimization (IPO) [165] has obtained increasing attention, and massive related approaches have been proposed. Concurrently, the complexity of the scenarios addressed in these studies and the scale of the agents involved have also been increasing synchronously, which reflects that decentralized learning matches the demands of large-scale autonomous driving in the real world more.

To solve the scalability issue in centralized schemes, MAP-PER [166] employs a decentralized actor-critic based on the A2C [155] algorithm. Firstly, the local observations of the occupancy map are represented as a 3-channel image containing static scenes, dynamic obstacles, and planned trajectory information from A* planner [167]. These 3-channel observations are abstracted into a latent vector via a CNN, along with waypoint information abstracted by an MLP, input into shared fully connected layers. Later, two independent MLPs output action probabilities and value estimates, respectively. Besides, MAPPER employs an extra evolutionary algorithm to eliminate bad policies during the optimization process. Compared with PRIMAL [154], MAPPER can learn faster and handle dynamic obstacles more effectively in large-scale scenarios. Another work scalability is G2RL [168], a grid map navigation method that can be used for any arbitrary number of agents. Similarly, it leverages A* to provide each agent with a global guiding path. Meanwhile, the local occupancy map is input into a local DDQN [161] planner to capture local observation and generate a corrective command to avoid dynamic obstacles. Since there is no need for communication between agents, this method does not require consideration of communication delays and can be extended to any scale.

As the successor to PRIMAL, PRIMAL$_2$ [169] retains the same hierarchical structure, i.e., an A* planner generating global paths and agent training guided by A3C and imitation learning. The key difference lies in PRIMAL$_2$'s fully decentralized training approach, which enhances its flexibility in handling structured and high-density complex scenarios. Like MAPPER, it adopts an $11 \times 11$ observation range and splits observations into multi-channel image inputs. The first 4 channels include static obstacles, the agent's own goal point, other agents' positions, and other agents' goal points. Channels 5-8 provide the local path from A* and the positions of other agents at three future timesteps within the observation range.
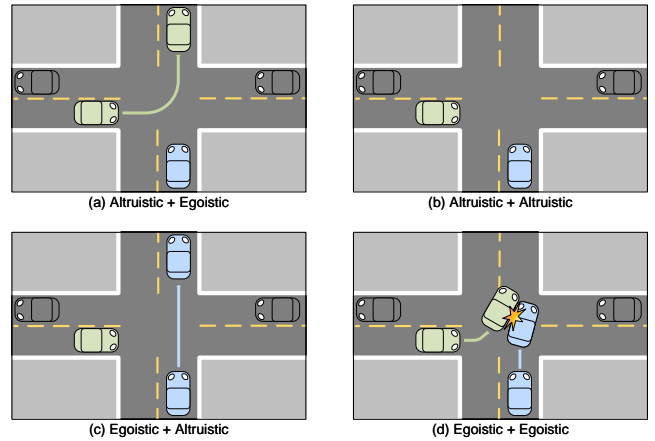


Fig. 4. Four examples at free intersection to show how social preference effects the behaviors of AVs. In (a) and (c), the combination (i.e. one egoistic + one altruistic) is healthy; in (b), two altruistic cars would wait for each other; in (d), two egoistic cars would crash into each other.

The final 3 channels offer the X and Y coordinate offsets of corridor exits and a boolean state indicating whether other agents are blocking the path. The finer observation inputs allow PRIMAL$_2$ to effectively address the agent deadlock issues in high-density, complex occupancy grids with shorter generated paths than its predecessor.

The aforementioned methods are developed for structured occupancy grids with discrete action spaces and applicable to automated ground vehicles in structured warehouses and freight terminals. While there are differences from real traffic systems, these methods remain inspirational for subsequent work. Other decentralized learning studies are conducted on more advanced continuous benchmarks [24], [63], [70]. For instance, in PIPO [21], researchers develop an end-to-end motion planning scheme using the permutation-invariant property of MAS with a graph neural network. They defined a progressively larger continuous scenario in MPE with various static obstacles. During training, the random permutation of observed states of other agents enhanced the feature representation of the actor-critic network. We note that there are numerous excellent and representative decentralized training schemes, but we categorize them under other subtopics and will elaborate on them in the following sections.

## C. Learning with Social Preference

Although independent policy learning is feasible in many tasks, it would lead to each agent being self-centric [20] when the interests of multiple agents conflict, the pure egoistic independent policy learning may fail. Therefore, an important issue is balancing the agents' egoism and altruism. In Fig. 4, we give a toy example to illustrate how social preference affects the agents' behaviors. If the agents cannot balance their altruistic and egoistic behaviors, these two would crash or get stopped by each other. Hence, social behaviors and preferences should be considered in policy learning [170]. To find a mathematical presentation of social preference, in the early work, researchers first propose to use a trigonometric

function like Eq. (12) to balance the individual and global reward, which inspires afterward studies [20], [62].

Afterward, as the representative work of driving with social preference, Coordinate Policy Optimization [20] (CoPO) draws inspiration from the self-driven particle systems in nature like fish and bird flocks and implements a decentralized MARL method with heterogeneous vehicle settings in MetaDrive [63]. It proposes hierarchical coordination for policy optimization to trade off the egoism and altruism of the policy. It utilizes the IPO for each vehicle within the traffic system to avoid the credit assignment problem. More specifically, by introducing a local coordination factor (LCF) in the training process, the agent seeks the optimal policy to maximize the averaged reward from all adjacent agents in its observable range. As Fig. 5, for the i-th agent $a_i$ in the system with time-invariant observed adjacent agents $N(i, t)$ in the range $d_n$, the agent should balance its ego reward $r_i$ and the average reward $r_i^N(t)$ as shown in Eq. (12).

$$r_i^N(t) = \frac{1}{|N(i,t)|} \sum_{j \in N(i,t)} r_j(t) \qquad (11)$$

$$r_i^C(t) = \cos(\phi) r_i(t) + \sin(\phi) r_i^N(t), \phi \in [-90°, 90°] \quad (12)$$

LCF $\phi$ indicates the specific altruism of the policy under a certain scenario, so the major difficulty of local coordination comes from selecting optimal LCF to maximize the global reward. Hence, in the global coordination, the author presents LCF as a Gaussian distribution $\phi \sim \mathcal{N}(\phi_\mu, \phi_\sigma)$ and the global objective is to find out the optimal policy with maximum accumulative rewards $J^G(\theta_1, \theta_2, \cdots) = \mathbb{E}_\tau \left[ \sum_{i \in N} \sum_{t=0}^T r_i(t) \right]$. To enable LCF to be learnable, the authors transfer the global objective into individual objective following with an easy factorization technique [116] and then derive the gradient from Eq. (13) to (15), where $\theta_i^o$ and $\theta_i^n$ denote the old new updated policy network of agent $i$.

$$\triangledown_\Phi J_i^G(\theta_i^n) = \triangledown_{\theta_i^n} J_i^G(\theta_i^n) \cdot \triangledown_\Phi \theta_i^n \qquad (13)$$

$$\triangledown_{\theta_i^n} J_i^G(\theta_i^n) = \mathbb{E} \left[ \triangledown_{\theta_i^n} \min(\rho A^G, \mathrm{clip}(\rho, 1 - \epsilon, 1 + \epsilon) A^G) \right] \qquad (14)$$

$$\triangledown_\Phi \theta_i^n = \alpha \cdot \mathbb{E} \left[ \triangledown_{\theta_i^o} \log \pi_{\theta_i^o(a_i|s)} \triangledown_\Phi A_{\Phi,i}^C \right] \qquad (15)$$

Here, $\Phi = [\phi_\mu, \phi_\sigma]$ denotes the mean and variance of LCF and is learnable, and $\alpha$ denotes the learning rate. In Eq. (14) and (15), $A^G$ and $A_{\Phi,i}^C$ is the advantage of global and locally coordinated reward respectively following the idea in [65]. Since Eq. (14) has no relevant term of $\Phi$, it can be regarded as a constant in the LCF objective. Hence, the objective of learning can be established as follows:

$$J^{LCF}(\Phi) = \mathbb{E} \left[ \triangledown_{\theta_i^n} J_i^G(\theta_i^n) \right] \cdot \left[ \triangledown_{\theta_i^o} \log \pi_{\theta_i^o(a_i|s)} A_{\Phi,i}^C \right] \quad (16)$$

Introducing a learnable factor to present social preference outperforms independent learning and mean-field policy optimization. It has also inspired many valuable discussions in this field about learning independent policy with sociological design and led to other excellent works. In subsequent research, people leverage more advanced language models like Transformer [137] to process social preference. For instance, in Social-Attention Policy Optimization (SAPO) [62], the
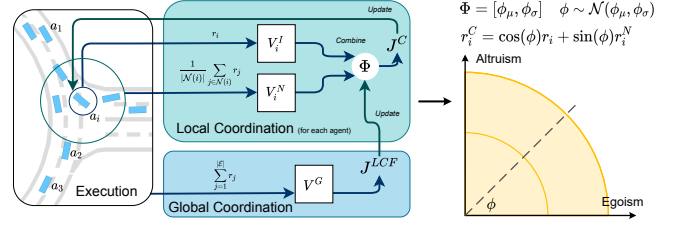


Fig. 5. The architecture of Coordinate Policy Optimization [20]. The bi-level training process enables to balance egoism and altruism.

authors introduce a multi-head attention layer to select the most interactive agent of the ego vehicle. In this work, the state of the ego vehicle is embedded into queries and compared to all the keys, and the interactive agent could be selected from the non-zero value index. After Gumbel Softmax [171] with gathered one-hot values, it obtains a compact permutation-invariant feature as the residual input of observation. Essentially, social coordination is a unique design of the reward function to balance collective and individual interests and align with the operational logic of real-world society.

### D. Safe and Trust-Worthy Learning

Safety is integral and the first priority to the deployment of autonomous driving systems, as they directly impact the reliability and people's lives of AVs. Recent RL researchers put massive efforts into ensuring the learned policy would never cause safety issues in the exploration process and after deployment. Specifically, inspired by [172], we categorize existing safety standards and methods in MARL into three types. First, soft safety guarantees involve designing safety penalty terms to reduce the probability of dangerous behavior. With fine-tuned rewards, the learning algorithm can be guided to prioritize safety alongside other performance metrics. However, although they have been proven to effectively improve safety performances in MAS, the limitation of soft guarantees is that they rely on the assumption that the reward function can accurately capture all safety aspects, which is often challenging in complex environments. The second is the probabilistic guarantees happening in the optimization process. For example, some recent MARL algorithms leverage the Lagrange constraints [21] or safety threshold during policy optimization process [173], [174]. Essentially, this improves policy gradient and helps avoid dangerous exploration behaviors. However, since the policy is still represented as a probability distribution, we cannot obtain a clear, explainable, and stable safety boundary for this method. Meanwhile, the vital safety constraints in real-world driving are instantaneous and deterministic [175]. For example, collision avoidance is a state-wise instantaneous constraint that only depends on the current state of the system rather than historical trajectories or random variables.

The third and safest approach is using hard safety boundaries to apply instantaneous strong corrections for agents' actions. For instance, researchers propose to learn the centralized shielding [184] from the joint action to correct any unsafe action in MAS at any risky time. Alternatively, combined with

TABLE II
RECENT MULTI-AGENT REINFORCEMENT LEARNING METHODOLOGY IN AUTONOMOUS DRIVING
AND INTELLIGENT TRANSPORTATION SYSTEM (RANKED BY YEAR)

| Research | Year | Training Scheme | Simulator | Maximum Agents | Vehicle Kinetics | Action Space Description | Real-World Experiment |
|---|---|---|---|---|---|---|---|
| MAMPS [176] | 2019 | Centralized | MPE | 4 | Omnidirection | Continuous, 2 dims | ✗ |
| PRIMAL [154] | 2019 | Centralized | Grid Map | 1024 | Grid model | Discrete, 5 dims | ✔ |
| MAPPER [166] | 2020 | Decentralized | Grid Map | 150 | Grid model | Discrete, 9 dims (add ↗, ↘, ↙, ↖ ) | ✗ |
| G2RL [168] | 2020 | Decentralized | Grid Map | 128 | Grid model | Discrete, 5 dims | ✗ |
| MACAD [71] | 2020 | Decentralized | MACAD | 3 | Bicycle model | Discrete, 9 dims (predefined action tuples) | ✗ |
| CoPO [20] | 2021 | Decentralized | MetaDrive | 40 | Full-vehicle model | Continuous, 2 dims | ✗ |
| GDQ [160] | 2021 | Centralized | MPE | 8 | Omnidirection | Continuous, 2 dims | ✔ |
| PRIMAL$_2$ [169] | 2021 | Decentralized | Grid Map | 2048 | Grid model | Discrete, 5 dims | ✔ |
| PIPO [21] | 2022 | Decentralized | MPE | 512 | Omnidirection | Continuous, 2 dims | ✗ |
| SSMAQL [177] | 2022 | Decentralized | Gym | 4 | Bicycle model | Discrete, 5 dims (high-level action) | ✗ |
| Duckie-MAAD [178] | 2022 | Centralized | Gym | 3 | Omnidirection | Discrete, 4 dims (high-level action) | ✔ |
| QMIXwD [159] | 2023 | Centralized | Highway-env | 4 | Bicycle model | Discrete, 3 dims (high-level action) | ✗ |
| TIRL [179] | 2023 | Decentralized | Grid Map | 64 | Grid model | Discrete, 5 dims | ✔ |
| SAPO [62] | 2023 | Decentralized | SMARTS | 4 | Bicycle model | Discrete, 3 dims (high-level action) | ✗ |
| CS-MADDPG [180] | 2024 | Centralized | Highway-Env | 2 | Bicycle model | Discrete, 5 dims (high-level action) | ✗ |
| CAVMARL [181] | 2024 | Decentralized | CARLA | 30 | Bicycle model | Discrete, 3 dims (high-level action) | ✗ |
| MFPG [182] | 2024 | Decentralized | ROS | 8 | Omnidirection | Continuous, 1 dim | ✔ |
| CPO-AD [183] | 2024 | Centralized | CARLA | 22 | Longitude only | Continuous, 1 dim | ✗ |

model predictive control, researchers propose a multi-agent model predictive shielding algorithm that provably guarantees safety for any policy learned from MARL [176]. However, due to the centralized setting, these methods cannot scale to the massive number of agents. Another widely-applied method for safety guarantee is control barrier function (CBF) [185], [186]. Assuming known explicit dynamics, early CBF-based safe controllers could ensure the safety of control strategies in simple tasks and environments [185]. However, real-world autonomous driving introduces complex nonlinear dynamics and intricate agent-level interactions, which lead the hand-crafted CBF to become infeasible [183]. Therefore, some studies have proposed incorporating additional neural networks to extract CBF for MAS and visualize it to interpret the safety boundaries [157], [181], [187]. For instance, MDBC [157] introduces neural barrier certificates for each agent and achieves scalable and super safe decentralized controllers for up to 1024 AVs and drones in particle environments. In the latest CAVMARL [181], researchers establish a safe action mapping via CBF-based quadratic programming. This controller leverages truncated Q-function to ensure the scalable joint state-action estimation under a centralized scheme and generates steering angle and acceleration with mathematically provable safety certificates.

## E. Methodological Summary

As shown in Table II, we collect representative works on MARL in outdoor autonomous driving, traffic system control, and structured scene transportation in the past five years. Meanwhile, we list their taxonomy, the maximum number of agents, the simulators, and whether real-world experiments are conducted. Here, we note that the action settings can be completely different even with the same simulation type. For example, in PRIMAL and PRIMAL$_2$, the agent's actions are set as $(\uparrow, \longrightarrow, \downarrow, \longleftarrow, *)$, representing four movements in the horizontal and vertical directions in a 2D grid map, along with staying in place. In contrast, MAPPER adds four additional diagonal movements ($\nearrow, \searrow, \swarrow, \nwarrow$ ) for the agents. Additionally, we find that many studies adopt predefined high-level action commands to simplify tasks. The policy network outputs discrete values that map to corresponding preset actions, and then a low-level controller takes the actions, generates commands, and sends them to the actuators. Two other specific examples are MFPG [182] and CPO-AD [183]. They preset a low-level unidirectional control mapping and only consider the movement of AVs in one direction.

Furthermore, we summarize three trends from past studies in this field. First, early research is limited by algorithm diversity and simulator performance and focuses more on centralized

MARL in grid maps. However, recent studies have discussed the potential of decentralized approaches with more complex continuous observation. Second, only a few studies conduct real-world experiments and use only discrete simulators and few agents. That is what the future works could improve. Third, the newer studies adopt more complicated designs and integrate more methods from other fields, such as data compression and machine vision.

## V. OPEN QUESTIONS AND CHALLENGES

In this section, we present the main challenges in MARL. Note that the problems faced by the CTDE and DTDE schemes are different, and though some feasible solutions have been proposed to solve these issues, they are still not unique and perfect. We hope that readers can become aware of their existence and properties in advance, thereby gaining a better understanding of the motivation and technical innovation from subsequent advanced methodologies.

### A. Multi-modal Information

Autonomous driving is a sequential decision-making process that leverages multi-modal information. Compared to MLPs in the original algorithm designs, recent research focuses more on designing more complicated neural network modules to learn better representations from multiple sensor information and to suit time series information.

*1) Multi-Sensor Fusion and Integration:* AVs acquire kinematic sensors like GPS and IMU, visual information from RGB-D cameras, LiDAR, and event cameras, and output executable commands to actuators. The multi-sensor information fusion enhances system safety redundancy, ensuring the vehicle can continue operating safely even if one sensor fails [188]. However, sensors with different principles possess various properties in real-world scenarios. For instance, like human drivers, RGB cameras capture the 3-channel information from reflected rays on the imaging plane and are cost-effective. However, RGB cameras are sensitive to light intensity and would overexpose in strong light or be covered by thermal noise in dark scenarios. Meanwhile, they are unsuitable for capturing high-speed objects for motion blur. Conversely, event cameras utilize parallelized pixels to capture brightness changes caused by high-speed moving objects, but they contain no color information [5]. LiDAR or depth camera can provide precise 3D spatial information. Also, information on kinematics and dynamics is significant for motion planning. Hence, it's challenging to merge these diverse data varying in formats, physical concepts, dimensions, and magnitudes. Early sensor fusion happens before the sensor information is input into the neural network. It concatenates the multi-modal information directly and relies on neural networks to extract helpful features [11], [20], [180]. This method is easy to implement but may result in the loss of geometric information and lead to more complex feature extraction. Contrarily, in the middle fusion paradigm, independent encoders extract key features from various modalities and squeeze them into latent vectors [21], [166], [169]. These feature vectors are concatenated along the feature dimension into a long vector
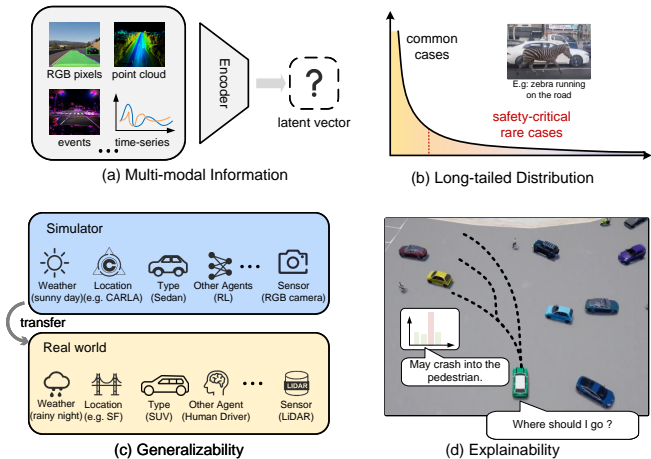


Fig. 6. Some examples of open questions and challenges. (a) AVs have to process multi-modal information from various sensors. (b) The long-tail distribution makes covering all scenarios in simulators or datasets difficult. (c) The learned policy should be generalizable to different environments and agents, and overcome the sim-to-real gap. (d) The learned policy should give a clear explanation for its action.

as the network input. Current research focuses on transferring the attention mechanism of language models to visual tasks and capturing contextual relationships between different modalities. Recent advances in sequential modeling and audio-visual fusion demonstrate that Transformer [137] is competent in modeling the information interaction for sequential or cross-modal data [189], [190].

*2) Learning Representations:* Massive approaches are involved in discussing the representation of the multi-modal information for vehicle sensor systems. Classic autonomous driving algorithms rely on HD maps, using diverse formulations and topological signs to represent map structures like segmented maps, vectorized centerlines, and landlines. For vehicle-centered maps, some studies directly input the sensor outputs like 3-channel pixel information or distances from LiDAR to network [11], [21], [191]. 2D grids or occupancy maps are a classical representation of vehicle-centered maps, which have been used in navigation in complex structured environments [154], [166], [169] and high-speed racing [192]. Alternatively, with depth cameras and LiDAR, the 3D scene can be reconstructed through SLAM [3] or Neural Radiance Field (NeRF) [193] and is represented in a 3D occupancy grid. However, high-resolution grids can lead to unacceptable computational consumption [30]. Recent research [194]–[196] proposes to unify multi-modal features in the shared bird's-eye view (BEV) representation space, which preserves both geometric and semantic information and suit for downstream tasks like navigation and obstacle avoidance. Although various representations offer potential designs for integration with MARL, determining which approach is the most effective remains an open and unresolved question. Additionally, balancing hardware costs with algorithmic performance and identifying which information is essential are crucial considerations for representation design.

## B. Robustness and Generalizability

Robustness and generalizability are critical factors for the effectiveness of MARL algorithms in autonomous driving. Robustness allows the vehicle to safely navigate a wide range of real-world conditions, while generalizability ensures the system can be applied broadly across different environments and scenarios. These qualities are fundamental to advancing autonomous driving technologies and achieving practical deployment in diverse and unpredictable real-world settings.

*1) Long-tailed Problem:* Autonomous vehicles are supposed to handle a broad spectrum of driving conditions, which include unexpected road obstructions, erratic behaviors from other drivers, unusual weather, and sudden mechanical failures. Numerous unpredictable corner cases lead to a gap in the vehicle's ability to respond effectively. As shown in Fig. 6, long-tailed problem [197] indicates a tricky data distribution of learning method, i.e., AVs' policies training can cover most common cases, but may fail to make safety-critical decisions in the rare scenarios that it has never seen before.

Hand-craft test bases [48], [70], [198], [199] and heuristic methods [200], [201] are leveraged to generate a broader range of traffic scenes. However, these generated scenarios demand significant human effort and domain expertise to create effective rules and often fail to capture the complexity of real-world traffic and structures accurately. Recently, researchers have been trying to achieve automatic traffic scene generation from real-world datasets to a unified representation and synthesize long realistic trajectories for RL [69], [90]. By exposing autonomous driving systems to these simulated environments, developers can ensure that the vehicles learn to manage unusual situations. This method complements real-world testing and helps build a more robust and adaptable system. Meanwhile, continual learning [202] and transfer learning [203] are practical tools to address long-tailed problems. These approaches enable MAS to learn from new experiences and update policies to incorporate novel situations so the vehicle can deal with everyday and rare scenarios throughout its operational life. In general, efficiently establishing realistic scenarios covering the long-tailed distribution of these rare and critical conditions is still a long-term challenge for autonomous driving.

*2) Sim-to-Real Gap:* Although it's useful for policy learning and pre-hoc testing, simulated environments often fail to capture real-world scenarios' full complexity and variability. Factors such as weather, road variations, and unpredictable behavior of other road components are difficult to replicate accurately in simulations [204]. This disparity can lead to autonomous driving systems performing well in simulated tests but failing to handle the nuances of real-world driving. First, simulators typically use idealized models of sensors and actuators, which do not account for the imperfections and limitations of real devices. This can lead to discrepancies in sensor readings and actuator responses, impacting the system's performance when deployed in real-world scenarios. Latency is another critical factor that differentiates simulations from real driving. In simulations, information transfer and processing times can be minimized or overlooked. However, in real-

world applications, delays in communication, sensor data processing, decision-making, and actuation can significantly affect the system's responsiveness [205]. Although we can arbitrarily define the frequency and delay of sensors and actuators in the advanced simulators [70], [79], it is challenging to model the fluctuations and disturbances of signal in the real world. Additionally, for reinforcement learning, predicting the agent's state after a collision helps design more comprehensive reward functions and avoid such scenarios. Modeling and simulating collisions remain a significant challenge involving complex mechanics, geometry, and material science. However, we notice some interesting discussion on differentiable collision processes in recent research [206].

The sim-to-real gap can be bridged by incorporating more detailed and diverse scenarios, realistic physics models, and accurate representations of sensor noise and failures. By creating more comprehensive and challenging simulations, developers can better prepare autonomous driving systems for the complexities of real-world operation. Meanwhile, through domain adaptation [207], we can improve the transferability of models trained in simulated environments to real applications. It involves adjusting a model trained in the source domain (simulation) to perform well in the target domain (real world) [208], [209]. This process aims to reduce the discrepancy between the two domains, allowing the model to generalize better to real-world conditions. Alternatively, we can integrate real-world data into the training and testing processes [69], [90], [210] so that it is possible to refine and validate simulation models and ensure a more close matching to real conditions.

## C. Safety Certificates

Safety is a paramount concern in the development of AV. Ensuring the reliability and robustness of these systems involves addressing several critical aspects, from extensive real-world testing to mitigating issues related to multi-agent communication and data transfer.

*1) Safety in Real World:* Industrial leaders emphasize the importance of cumulative test lengths to ensure the reliability of autonomous driving because real-world testing provides invaluable insights that cannot be fully captured in simulation. It allows for the evaluation of AVs in real traffic, weather, and road conditions and figures out the potential flaws that may not emerge during simulation. This type of testing is crucial for understanding how the system interacts with human drivers, pedestrians, and other road users, ensuring that it operates safely and effectively in a live environment. Although recent advancements [157], [186], [187] show the reliability of the hard safety constraints and the interpretation of barriers through learning methods, the long-tailed problem means we still cannot sample all risky scenarios and guarantee 100% safety with these approaches in the real world. Additionally, only a few research institutions and universities have access to specialized testing facilities for real-world traffic testing like MCity [211]. Although some studies include real-world validation and demonstrations, most are limited to down-scaled and simplified platforms [154], [160], [169]. Consequently, the current safety validation of MARL-related research remains

significantly inadequate. We believe this can be improved by accelerating the development of testing and research infrastructure and enhancing regional and international academic collaboration.

*2) Private Information Concern:* MARL-based traffic systems rely heavily on data exchange and communication between agents. In MAS, AVs need to frequently share data related to their environment, position, and intended actions to coordinate effectively. This data exchange, while essential for the system's functioning, can expose sensitive information. Personal data, such as location history, can be vulnerable to unauthorized access if adequate security measures are not in place. External attacks significantly threaten the privacy and security of MARL systems in autonomous driving. Attackers can exploit vulnerabilities in AVs' communication protocols and data storage systems. To alleviate these concerns, developing robust communication protocols that include authentication and verification mechanisms can help prevent unauthorized access and data tampering. Meanwhile, real-time anomaly detection systems [212] can help identify and mitigate potential attacks by monitoring for unusual patterns or behaviors in the data exchange. Decentralized data storage [213] can also enhance data security by providing a tamper-proof ledger of all data exchanges and transactions.

### D. Explainability

Explainability emphasizes the understandable and causable decision-making process, which ensures the actions are transparent and reasonable. However, nowadays, most learning-based methods still adopt the black-box deep neural network as the main component. As the primary shortcomings of black-box models, the lack of transparency poses significant safety concerns [214]. This opacity can lead to a lack of trust and confidence in the system, especially in safety-critical situations where understanding the reasoning behind a decision is crucial. For instance, if an AV makes an unexpected maneuver, it is essential to know whether the decision is based on valid reasoning or a flaw in the system.

Explainable methods like decision tree [215] or rule-based systems [216] can be used to establish the mathematically understandable controller. However, fine system-level design needs expert domain knowledge. State representation learning can create a low-dimensional and meaningful representation of the state space by processing high-dimensional raw observation data [217]. This approach captures the environmental variations influenced by the agent's actions, facilitating the extrapolation of explanations. Mainly, it is helpful in reinforcement learning for robotics and control, as it aids in understanding how the agent interprets observations and identifies what is relevant for learning to act effectively [218], [219]. Besides, applying post-hoc analysis methods to existing models can help interpret their decisions [220]. Techniques such as feature importance analysis, saliency map, and layer-wise relevance propagation can provide insights into which features and how these features influence specific decisions. Incorporating attention mechanisms in the model architecture can highlight which parts of the input data are most relevant to the decision-making

process [137], [221]–[223]. This can help in understanding the focus areas of the model during critical decision points. Another solution is using model-agnostic methods like local interpretable model-agnostic explanations (LIME) [214], [224] to provide the predictive explanations of any black-box model, which approximate the model's behavior locally to explain individual predictions.

## VI. FUTURE DIRECTIONS

In the last section, we briefly introduce the latest advancements and explain why we believe these directions are promising. We hope that this information will inspire researchers and lead to more outstanding research.

### A. Model-based MARL

Model-based RL has achieved significant progress in single AV driving. By incorporating additional neural networks, we can model complex nonlinear dynamics and state transition functions [225], [226]. In recent research, researchers implement real-world high-speed racing in complicated tracks via extra 4 networks for environmental dynamics and prediction of future state, observation, and reward [192]. However, there is no free lunch. While model-based reinforcement learning offers better performance and improves explainability, it also increases the requirement for computational resources. Typically, centralized approaches model the environment through joint actions and observations and get rid of the non-stationarity and partial observability. However, scalability remains a significant challenge, especially for heterogeneous MAS. Conversely, decentralized approaches are easier to scale but struggle to reach consensus in non-stationary dynamics and partially observable environments. Beyond the design of communication protocols, recent research is exploring the possibility of abstracted and simplified modeling, such as inferring and predicting the subsequent actions of other agents. Additionally, decentralized paradigms introduce more networks for learning models, so both the selection of model representations and the design of efficient network architectures are promising topics.

### B. Development of Offline Multi-Agent Datasets

Offline paradigms can improve the practicality and realism of reinforcement learning. Specifically, online trial-and-error learning can cause financial losses and social disruption in mission-critical systems. Training policies in simulations suffer from an inherent gap between simulated environments and real dynamics and limits the ability to leverage extensive previously collected datasets [68]. We observe that some recent studies are exploring how to use previously collected datasets to train single-vehicle offline reinforcement learning paradigms, paving the way for the expansion into MAS [227]. As mentioned before, although independent learning can achieve optimal individual policies, it follows a self-centric optimization and lacks the validation of social preference and altruism. From this assumption, it is necessary to redevelop and collect datasets at the system level rather than the individual level. Combining with data augmentation

approaches [228]–[230], these datasets can be exponentially expanded to cover most of the daily and long-tailed scenarios. Additionally, there are opportunities to combine offline MARL with other emerging methods, such as safe-RL [231] or meta-RL [232]. In general, offline MARL and system-level datasets are essential for advancing large-scale autonomous driving. From pre-collected data and integrating simulations, we hope developers can create safer, more efficient, and more reliable autonomous driving systems in the future.

### C. Human-in-the-Loop Learning

Human-in-the-loop (HITL) learning [233] can enhance the effectiveness and safety of autonomous driving systems. These methods incorporate human feedback into the learning process, integrate human expertise and intuition, and ensure that systems can handle complex and dynamic real-world scenarios more effectively by leveraging human judgment to guide the learning process. Specifically, HITL learning involves human operators providing feedback on the system's performance during training. This feedback can come in various forms, such as corrections to the vehicle's actions, suggestions for alternative routes, or evaluations of the system's decision-making process. By integrating human feedback, autonomous driving systems can better identify and avoid potential hazards and reduce the probability of accidents. Besides, around human demonstrations, we can design extra modules to ensure the effective learning of policy correction and safety improvement [234]–[236]. Reinforcement Learning with Human Feedback (RLHF) is a recently proposed powerful paradigm that transforms human feedback into guidance signals. Particularly, it eliminates the requirement for manual reward design, which plays a vital role in robotic control and language models. Besides modeling the reward, RLHF can help autonomous systems quickly adapt to new and unexpected driving scenarios by incorporating real-time human feedback. In the latest research [237], the first large-scale benchmark is proposed by integrating a series of RLHF baselines and over 1000 hybrid agents and expert driving trajectories in SMARTS [55] simulator. Meanwhile, its exciting results demonstrate that RLHF can follow the abstraction metrics and learn good human behavior. However, HITL means that it is inevitable that human driver participation will be required. How to minimize the need for human intervention, effectively extract key strategy improvements from these interventions, and prevent catastrophic forgetting will be important topics for future research in this field.

### D. Language Models for Autonomous Driving

As an eye-catching rising star, we have witnessed the contributions language models to the control communities [238]–[240]. Although they are not necessarily RL frameworks, we still discuss their possibilities and future in our paper. Nowadays, large language models (LLMs) [241]–[243] and visual language models (VLMs) [244], [245] show significant potential for autonomous driving. These models can process and understand both textual and visual information and trigger new approaches to enhance the capabilities of AVs. LLMs can contribute to autonomous driving by providing sophisticated

natural language understanding and generation support. These models can be used to interpret verbal instructions, understand traffic signs and signals, and process natural language inputs from passengers [246]–[248]. Although they show expressive common sense reasoning and understanding capabilities, they are mostly leveraged only in simple decision-making due to the inability of visual information. VLMs combine visual and linguistic data to create a better understanding of the environment [244], [249]. These models can interpret complex visual scenes and provide contextual understanding, which is crucial for making informed driving decisions. Moreover, they can be used to recognize and describe objects, predict the behavior of pedestrians and other vehicles, and understand the broader context of the driving environment [245], [250], [251]. Recently, we discovered an interesting discussion and imagination of language models under the MARL paradigm in the latest review [240]. For instance, language models could simplify HITL interventions from complex manual operations to voice inputs. Alternatively, it can be used for context distillation [252] and make the complex observations and communication information in MAS more compact. Thus, the distilled features would be more suitable for edge computing platforms. This is a rapidly growing and highly significant field that warrants continuous attention.

### REFERENCES

[1] T. T. Nguyen, N. D. Nguyen, and S. Nahavandi, "Deep reinforcement learning for multiagent systems: A review of challenges, solutions, and applications," *IEEE Transactions on Cybernetics*, vol. 50, no. 9, pp. 3826–3839, 2020.

[2] S. Teng, X. Hu, P. Deng, B. Li, Y. Li, Y. Ai, D. Yang, L. Li, Z. Xuanyuan, F. Zhu, and L. Chen, "Motion planning for autonomous driving: The state of the art and future perspectives," *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 6, pp. 3692–3711, 2023.

[3] J. Zhang and S. Singh, "Loam: Lidar odometry and mapping in real-time.," in *Robotics: Science and systems*, pp. 1–9, 2014.

[4] P. Liu, G. Chen, Z. Li, D. Clarke, Z. Liu, R. Zhang, and A. Knoll, "Neurodfd: Towards efficient driver face detection with neuromorphic vision sensor," in *2022 International Conference on Advanced Robotics and Mechatronics (ICARM)*, pp. 268–273, 2022.

[5] G. Chen, H. Cao, J. Conradt, H. Tang, F. Rohrbein, and A. Knoll, "Event-based neuromorphic vision for autonomous driving: A paradigm shift for bio-inspired visual sensing and perception," *IEEE Signal Processing Magazine*, vol. 37, no. 4, pp. 34–49, 2020.

[6] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.

[7] D. Silver, T. Hubert, J. Schrittwieser, I. Antonoglou, M. Lai, A. Guez, M. Lanctot, L. Sifre, D. Kumaran, T. Graepel, *et al.*, "A general reinforcement learning algorithm that masters chess, shogi, and go through self-play," *Science*, vol. 362, no. 6419, pp. 1140–1144, 2018.

[8] D. Silver, J. Schrittwieser, K. Simonyan, I. Antonoglou, A. Huang, A. Guez, T. Hubert, L. Baker, M. Lai, A. Bolton, *et al.*, "Mastering the game of go without human knowledge," *nature*, vol. 550, no. 7676, pp. 354–359, 2017.

[9] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, *et al.*, "Grandmaster level in starcraft ii using multi-agent reinforcement learning," *Nature*, vol. 575, no. 7782, pp. 350–354, 2019.

[10] F. Fuchs, Y. Song, E. Kaufmann, D. Scaramuzza, and P. Dürr, "Super-human performance in gran turismo sport using deep reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4257–4264, 2021.

[11] R. Zhang, J. Hou, G. Chen, Z. Li, J. Chen, and A. Knoll, "Residual policy learning facilitates efficient model-free autonomous racing," *IEEE Robotics and Automation Letters*, vol. 7, no. 4, pp. 11625–11632, 2022.

[12] Y. Song, A. Romero, M. Müller, V. Koltun, and D. Scaramuzza, "Reaching the limit in autonomous racing: Optimal control versus reinforcement learning," *Science Robotics*, vol. 8, no. 82, p. eadg1462, 2023.

[13] E. Kaufmann, L. Bauersfeld, A. Loquercio, M. Müller, V. Koltun, and D. Scaramuzza, "Champion-level drone racing using deep reinforcement learning," *Nature*, vol. 620, no. 7976, pp. 982–987, 2023.

[14] "Dimensions: Ai-powered research engine." [Online]. Available: https://www.dimensions.ai.

[15] S. Gronauer and K. Diepold, "Multi-agent deep reinforcement learning: a survey," *Artificial Intelligence Review*, pp. 1–49, 2022.

[16] M. Brittain and P. Wei, "Autonomous separation assurance in an high-density en route sector: A deep multi-agent reinforcement learning approach," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pp. 3256–3262, 2019.

[17] J. Yang, J. Zhang, and H. Wang, "Urban traffic control in software defined internet of things via a multi-agent deep reinforcement learning approach," *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 6, pp. 3742–3754, 2021.

[18] J. Perolat, J. Z. Leibo, V. Zambaldi, C. Beattie, K. Tuyls, and T. Grae-pel, "A multi-agent reinforcement learning model of common-pool resource appropriation," *Advances in neural information processing systems*, vol. 30, 2017.

[19] L. Ding, Z. Lin, X. Shi, and G. Yan, "Target-value-competition-based multi-agent deep reinforcement learning algorithm for distributed nonconvex economic dispatch," *IEEE Transactions on Power Systems*, vol. 38, no. 1, pp. 204–217, 2023.

[20] Z. Peng, Q. Li, K. M. Hui, C. Liu, and B. Zhou, "Learning to simulate self-driven particles system with coordinated policy optimization," *Advances in Neural Information Processing Systems*, vol. 34, pp. 10784–10797, 2021.

[21] R. Zhang, G. Chen, J. Hou, Z. Li, and A. Knoll, "Pipo: Policy optimization with permutation-invariant constraint for distributed multi-robot navigation," in *2022 IEEE International Conference on Multisensor Fusion and Integration for Intelligent Systems (MFI)*, pp. 1–7, 2022.

[22] D. Qiu, J. Wang, J. Wang, and G. Strbac, "Multi-agent reinforcement learning for automated peer-to-peer energy trading in double-side auction market.," in *IJCAI*, pp. 2913–2920, 2021.

[23] A. Shavandi and M. Khedmati, "A multi-agent deep reinforcement learning framework for algorithmic trading in financial markets," *Expert Systems with Applications*, vol. 208, p. 118124, 2022.

[24] R. Lowe, Y. WU, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," in *Advances in Neural Information Processing Systems*, vol. 30, 2017.

[25] L. Chen, Y. Li, C. Huang, B. Li, Y. Xing, D. Tian, L. Li, Z. Hu, X. Na, Z. Li, S. Teng, C. Lv, J. Wang, D. Cao, N. Zheng, and F.-Y. Wang, "Milestones in autonomous driving and intelligent vehicles: Survey of surveys," *IEEE Transactions on Intelligent Vehicles*, vol. 8, no. 2, pp. 1046–1056, 2023.

[26] L. Chen, Y. Li, C. Huang, Y. Xing, D. Tian, L. Li, Z. Hu, S. Teng, C. Lv, J. Wang, D. Cao, N. Zheng, and F.-Y. Wang, "Milestones in autonomous driving and intelligent vehicles—part i: Control, computing system design, communication, hd map, testing, and human behaviors," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 53, no. 9, pp. 5831–5847, 2023.

[27] L. Chen, S. Teng, B. Li, X. Na, Y. Li, Z. Li, J. Wang, D. Cao, N. Zheng, and F.-Y. Wang, "Milestones in autonomous driving and intelligent vehicles—part ii: Perception and planning," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, vol. 53, no. 10, pp. 6401–6415, 2023.

[28] S. Aradi, "Survey of deep reinforcement learning for motion planning of autonomous vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 2, pp. 740–759, 2020.

[29] B. R. Kiran, I. Sobh, V. Talpaert, P. Mannion, A. A. A. Sallab, S. Yo-gamani, and P. Pérez, "Deep reinforcement learning for autonomous driving: A survey," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 6, pp. 4909–4926, 2021.

[30] L. Chen, P. Wu, K. Chitta, B. Jaeger, A. Geiger, and H. Li, "End-to-end autonomous driving: Challenges and frontiers," *arXiv*, vol. 2306.16927, 2023.

[31] 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.

[32] Y. Huang, J. Du, Z. Yang, Z. Zhou, L. Zhang, and H. Chen, "A survey on trajectory-prediction methods for autonomous driving," *IEEE Transactions on Intelligent Vehicles*, vol. 7, no. 3, pp. 652–674, 2022.

[33] L. M. Schmidt, J. Brosig, A. Plinge, B. M. Eskofier, and C. Mutschler, "An introduction to multi-agent reinforcement learning and review of its application to autonomous mobility," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1342–1349, 2022.

[34] P. Yadav, A. Mishra, and S. Kim, "A comprehensive survey on multi-agent reinforcement learning for connected and automated vehicles," *Sensors*, vol. 23, no. 10, p. 4710, 2023.

[35] 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.

[36] X. Wang, S. Wang, X. Liang, D. Zhao, J. Huang, X. Xu, B. Dai, and Q. Miao, "Deep reinforcement learning: A survey," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–15, 2022.

[37] S. D. Pendleton, H. Andersen, X. Du, X. Shen, M. Meghjani, Y. H. Eng, D. Rus, and M. H. Ang, "Perception, planning, control, and coordination for autonomous vehicles," *Machines*, vol. 5, no. 1, 2017.

[38] J. Hou, G. Chen, R. Zhang, Z. Li, S. Gu, and C. Jiang, "Spreeze: High-throughput parallel reinforcement learning framework," *arXiv preprint arXiv:2312.06126*, 2023.

[39] J. Cameron, S. Myint, C. Kuo, A. Jain, H. Grip, P. Jayakumar, and J. Overholt, "Real-time and high-fidelity simulation environment for autonomous ground vehicle dynamics," in *Annual Ground Vehicle Systems Engineering And Technology Symposium (GVSETS) Symposium*, 2013.

[40] B. Wymann, E. Espié, C. Guionneau, C. Dimitrakakis, R. Coulom, and A. Sumner, "Torcs, the open racing car simulator," *Software available at http://torcs. sourceforge. net*, vol. 4, no. 6, p. 2, 2000.

[41] B. Wymann, C. Dimitrakakis, A. Sumner, E. Espié, and C. Guionneau, "Torcs: The open racing car simulator," 2015.

[42] D. Loiacono, P. L. Lanzi, J. Togelius, E. Onieva, D. A. Pelta, M. V. Butz, T. D. Lönneker, L. Cardamone, D. Perez, Y. Sáez, M. Preuss, and J. Quadflieg, "The 2009 simulated car racing championship," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 2, no. 2, pp. 131–147, 2010.

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

[44] D. Loiacono, L. Cardamone, and P. L. Lanzi, "Simulated car racing championship: Competition software manual," *arXiv preprint arXiv:1304.1672*, 2013.

[45] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, "Openai gym," *arXiv preprint arXiv:1606.01540*, 2016.

[46] U. Namakemono, "Gym torcs." [Online]. Available: https://github.com/ugo-nama-kun/gym_torcs.

[47] A. Santara, S. Rudra, S. A. Buridi, M. Kaushik, A. Naik, B. Kaul, and B. Ravindran, "Madras: Multi agent driving simulator," *Journal of Artificial Intelligence Research*, vol. 70, pp. 1517–1555, 2021.

[48] D. Krajzewicz, G. Hertkorn, C. Rössel, and P. Wagner, "Sumo (simulation of urban mobility)-an open-source traffic simulation," in *Proceedings of the 4th middle East Symposium on Simulation and Modelling (MESM20002)*, pp. 183–187, 2002.

[49] S. Liu, Y. Wang, X. Chen, Y. Fu, and X. Di, "Smart-eflo: An integrated sumo-gym framework for multi-agent reinforcement learning in electric fleet management problem," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 3026–3031, 2022.

[50] T. Chu, J. Wang, L. Codecà, and Z. Li, "Multi-agent deep reinforcement learning for large-scale traffic signal control," *IEEE Transactions on Intelligent Transportation Systems*, vol. 21, no. 3, pp. 1086–1095, 2020.

[51] B. Jiang, S. N. Givigi, and J.-A. Delamer, "A marl approach for optimizing positions of vanet aerial base-stations on a sparse highway," *IEEE Access*, vol. 9, pp. 133989–134004, 2021.

[52] C. Wu, A. R. Kreidieh, K. Parvate, E. Vinitsky, and A. M. Bayen, "Flow: A modular learning framework for mixed autonomy traffic," *IEEE Transactions on Robotics*, vol. 38, no. 2, pp. 1270–1286, 2022.

[53] E. Liang, R. Liaw, R. Nishihara, P. Moritz, R. Fox, K. Goldberg, J. Gonzalez, M. Jordan, and I. Stoica, "Rllib: Abstractions for distributed reinforcement learning," in *International Conference on Machine Learning*, pp. 3053–3062, PMLR, 2018.

[54] H. Zhang, S. Feng, C. Liu, Y. Ding, Y. Zhu, Z. Zhou, W. Zhang, Y. Yu, H. Jin, and Z. Li, "Cityflow: A multi-agent reinforcement learning environment for large scale city traffic scenario," in *The world wide web conference*, pp. 3620–3624, 2019.

[55] M. Zhou, J. Luo, J. Villella, Y. Yang, D. Rusu, J. Miao, W. Zhang, M. Alban, I. Fadakar, Z. Chen, *et al.*, "Smarts: Scalable multi-agent reinforcement learning training school for autonomous driving," *arXiv preprint arXiv:2010.09776*, 2020.

[56] M. Samvelyan, T. Rashid, C. Schroeder de Witt, G. Farquhar, N. Nardelli, T. G. Rudner, C.-M. Hung, P. H. Torr, J. Foerster, and S. Whiteson, "The starcraft multi-agent challenge," in *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 2186–2188, 2019.

[57] M. Zhou, Z. Wan, H. Wang, M. Wen, R. Wu, Y. Wen, Y. Yang, Y. Yu, J. Wang, and W. Zhang, "Malib: A parallel framework for population-based multi-agent reinforcement learning," *Journal of Machine Learning Research*, vol. 24, no. 150, pp. 1–12, 2023.

[58] E. Coumans and Y. Bai, "Pybullet, a python module for physics simulation for games, robotics and machine learning." http://pybullet.org, 2021.

[59] J. Guan, G. Chen, J. Huang, Z. Li, L. Xiong, J. Hou, and A. Knoll, "A discrete soft actor-critic decision-making strategy with sample filter for freeway autonomous driving," *IEEE Transactions on Vehicular Technology*, vol. 72, no. 2, pp. 2593–2598, 2023.

[60] H. Liu, Z. Huang, X. Mo, and C. Lv, "Augmenting reinforcement learning with transformer-based scene representation learning for decision-making of autonomous driving," *arXiv preprint arXiv:2208.12263*, 2022.

[61] E. Leurent, *Safe and efficient reinforcement learning for behavioural planning in autonomous driving*. PhD thesis, Université de Lille, 2020.

[62] Z. Dai, T. Zhou, K. Shao, D. H. Mguni, B. Wang, and H. Jianye, "Socially-attentive policy optimization in multi-agent self-driving system," in *Conference on Robot Learning*, pp. 946–955, PMLR, 2023.

[63] Q. Li, Z. Peng, L. Feng, Q. Zhang, Z. Xue, and B. Zhou, "Metadrive: Composing diverse driving scenarios for generalizable reinforcement learning," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 3, pp. 3461–3475, 2023.

[64] M. Goslin and M. Mine, "The panda3d graphics engine," *Computer*, vol. 37, no. 10, pp. 112–114, 2004.

[65] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.

[66] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, "Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor," in *International Conference on Machine Learning*, pp. 1861–1870, PMLR, 2018.

[67] 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.

[68] S. Levine, A. Kumar, G. Tucker, and J. Fu, "Offline reinforcement learning: Tutorial, review, and perspectives on open problems," *arXiv preprint arXiv:2005.01643*, 2020.

[69] L. Feng, Q. Li, Z. Peng, S. Tan, and B. Zhou, "Trafficgen: Learning to generate diverse and realistic traffic scenarios," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3567–3575, 2023.

[70] A. Dosovitskiy, G. Ros, F. Codevilla, A. Lopez, and V. Koltun, "CARLA: An open urban driving simulator," in *Proceedings of the 1st Annual Conference on Robot Learning*, vol. 78 of *Proceedings of Machine Learning Research*, pp. 1–16, PMLR, 13–15 Nov 2017.

[71] P. Palanisamy, "Multi-agent connected autonomous driving using deep reinforcement learning," in *2020 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–7, 2020.

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

[73] S. Lan, Z. Wang, E. Wei, A. K. Roy-Chowdhury, and Q. Zhu, "Collaborative multi-agent video fast-forwarding," *IEEE Transactions on Multimedia*, pp. 1–14, 2023.

[74] D. Ye, T. Zhu, C. Zhu, W. Zhou, and P. S. Yu, "Model-based self-advising for multi-agent learning," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 34, no. 10, pp. 7934–7945, 2023.

[75] A. Amini, I. Gilitschenski, J. Phillips, J. Moseyko, R. Banerjee, S. Karaman, and D. Rus, "Learning robust control policies for end-to-end autonomous driving from data-driven simulation," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 1143–1150, 2020.

[76] A. Amini, T.-H. Wang, I. Gilitschenski, W. Schwarting, Z. Liu, S. Han, S. Karaman, and D. Rus, "Vista 2.0: An open, data-driven simulator for multimodal sensing and policy learning for autonomous vehicles," in *2022 International Conference on Robotics and Automation (ICRA)*, pp. 2419–2426, 2022.

[77] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pp. 234–241, Springer, 2015.

[78] T.-H. Wang, A. Amini, W. Schwarting, I. Gilitschenski, S. Karaman, and D. Rus, "Learning interactive driving policies via data-driven simulation," in *2022 International Conference on Robotics and Automation (ICRA)*, pp. 7745–7752, 2022.

[79] "Isaac sim." [Online]. Available: https://developer.nvidia.com/isaac-sim, 2022.

[80] M. Quigley, K. Conley, B. Gerkey, J. Faust, T. Foote, J. Leibs, R. Wheeler, A. Y. Ng, *et al.*, "Ros: an open-source robot operating system," in *ICRA workshop on open source software*, p. 5, 2009.

[81] S. Macenski, T. Foote, B. Gerkey, C. Lalancette, and W. Woodall, "Robot operating system 2: Design, architecture, and uses in the wild," *Science Robotics*, vol. 7, no. 66, p. eabm6074, 2022.

[82] N. Koenig and A. Howard, "Design and use paradigms for gazebo, an open-source multi-robot simulator," in *2004 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS) (IEEE Cat. No.04CH37566)*, vol. 3, pp. 2149–2154 vol.3, 2004.

[83] "Isaac-lab." [Online]. Available: https://https://github.com/isaac-sim/IsaacLab.

[84] E. Leurent, "An environment for autonomous driving decision-making." [Online]. Available: https://github.com/eleurent/highway-env.

[85] J. Bernhard, K. Esterle, P. Hart, and T. Kessler, "Bark: Open behavior benchmarking in multi-agent environments," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 6201–6208, 2020.

[86] D. Xu, Y. Chen, B. Ivanovic, and M. Pavone, "Bits: Bi-level imitation for traffic simulation," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2929–2936, 2023.

[87] A. Ścibior, V. Lioutas, D. Reda, P. Bateni, and F. Wood, "Imagining the road ahead: Multi-agent trajectory prediction via differentiable simulation," in *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pp. 720–725, 2021.

[88] Q. Sun, X. Huang, B. C. Williams, and H. Zhao, "Intersim: Interactive traffic simulation via explicit relation modeling," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 11416–11423, 2022.

[89] E. Vinitsky, N. Lichtlé, X. Yang, B. Amos, and J. Foerster, "Nocturne: a scalable driving benchmark for bringing multi-agent learning one step closer to the real world," *Advances in Neural Information Processing Systems*, vol. 35, pp. 3962–3974, 2022.

[90] Q. Li, Z. M. Peng, L. Feng, Z. Liu, C. Duan, W. Mo, and B. Zhou, "Scenarionet: Open-source platform for large-scale traffic scenario simulation and modeling," in *Advances in Neural Information Processing Systems* (A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, eds.), vol. 36, pp. 3894–3920, 2024.

[91] C. Gulino, J. Fu, W. Luo, G. Tucker, E. Bronstein, Y. Lu, J. Harb, X. Pan, Y. Wang, X. Chen, *et al.*, "Waymax: An accelerated, data-driven simulator for large-scale autonomous driving research," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[92] V. Makoviychuk, L. Wawrzyniak, Y. Guo, M. Lu, K. Storey, M. Macklin, D. Hoeller, N. Rudin, A. Allshire, A. Handa, *et al.*, "Isaac gym: High performance gpu based physics simulation for robot learning," in *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.

[93] M. Martinez, C. Sitawarin, K. Finch, L. Meincke, A. Yablonski, and A. Kornhauser, "Beyond grand theft auto v for training, testing

and enhancing deep learning in self driving cars," *arXiv preprint arXiv:1712.01397*, 2017.

[94] "Deepdrive: A simulator that allows anyone with a pc to push the state-of-the-art in self-driving." [Online]. Available: https://github.com/deepdrive/deepdrive.

[95] S. Shah, D. Dey, C. Lovett, and A. Kapoor, "Airsim: High-fidelity visual and physical simulation for autonomous vehicles," in *Field and Service Robotics: Results of the 11th International Conference*, pp. 621–635, Springer, 2018.

[96] A. Geiger, P. Lenz, and R. Urtasun, "Are we ready for autonomous driving? the kitti vision benchmark suite," in *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3354–3361, 2012.

[97] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuscenes: A multimodal dataset for autonomous driving," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 11621–11631, 2020.

[98] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, *et al.*, "Scalability in perception for autonomous driving: Waymo open dataset," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2446–2454, 2020.

[99] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Madhavan, and T. Darrell, "Bdd100k: A diverse driving dataset for heterogeneous multitask learning," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2636–2645, 2020.

[100] W. Zhan, L. Sun, D. Wang, H. Shi, A. Clausse, M. Naumann, J. Kummerle, H. Konigshof, C. Stiller, A. de La Fortelle, *et al.*, "Interaction dataset: An international, adversarial and cooperative motion dataset in interactive driving scenarios with semantic maps," *arXiv preprint arXiv:1910.03088*, 2019.

[101] D. Lee, C. Eom, and M. Kwon, "Ad4rl: Autonomous driving benchmarks for offline reinforcement learning with value-based dataset," in *2024 IEEE International Conference on Robotics and Automation (ICRA)*, 2024.

[102] V. Alexiadis, J. Colyar, J. Halkias, R. Hranac, and G. McHale, "The next generation simulation program," *Institute of Transportation Engineers. ITE Journal*, vol. 74, no. 8, p. 22, 2004.

[103] M. Buehler, K. Iagnemma, and S. Singh, *The DARPA urban challenge: autonomous vehicles in city traffic*, vol. 56. Springer Science & Business Media, 2009.

[104] L. Paull, J. Tani, H. Ahn, J. Alonso-Mora, L. Carlone, M. Cap, Y. F. Chen, C. Choi, J. Dusek, Y. Fang, *et al.*, "Duckietown: an open, inexpensive and flexible platform for autonomy education and research," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1497–1504, IEEE, 2017.

[105] H. Caesar, J. Kabzan, K. S. Tan, W. K. Fong, E. Wolff, A. Lang, L. Fletcher, O. Beijbom, and S. Omari, "nuplan: A closed-loop ml-based planning benchmark for autonomous vehicles," *arXiv preprint arXiv:2106.11810*, 2021.

[106] R. A. Howard, *Dynamic programming and markov processes*. John Wiley, 1960.

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

[108] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, pp. 279–292, 1992.

[109] R. S. Sutton, D. McAllester, S. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," *Advances in neural information processing systems*, vol. 12, 1999.

[110] A. Ilyas, L. Engstrom, S. Santurkar, D. Tsipras, F. Janoos, L. Rudolph, and A. Madry, "A closer look at deep policy gradients," in *International Conference on Learning Representations*, 2020.

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

[112] J. F. Nash Jr, "Equilibrium points in n-person games," *Proceedings of the national academy of sciences*, vol. 36, no. 1, pp. 48–49, 1950.

[113] L. Steels, "The biology and technology of intelligent autonomous agents," *Robotics and Autonomous Systems*, vol. 1, no. 15, pp. 1–2, 1995.

[114] K. Zhang, Z. Yang, and T. Başar, "Multi-agent reinforcement learning: A selective overview of theories and algorithms," *Handbook of reinforcement learning and control*, pp. 321–384, 2021.

[115] E. A. Hansen, D. S. Bernstein, and S. Zilberstein, "Dynamic programming for partially observable stochastic games," in *AAAI*, vol. 4, pp. 709–715, 2004.

[116] T. Rashid, M. Samvelyan, C. S. De Witt, G. Farquhar, J. Foerster, and S. Whiteson, "Monotonic value function factorisation for deep multi-agent reinforcement learning," *The Journal of Machine Learning Research*, vol. 21, no. 1, pp. 7234–7284, 2020.

[117] J. Foerster, G. Farquhar, T. Afouras, N. Nardelli, and S. Whiteson, "Counterfactual multi-agent policy gradients," in *Proceedings of the AAAI conference on artificial intelligence*, 2018.

[118] M. Jaderberg, W. M. Czarnecki, I. Dunning, L. Marris, G. Lever, A. G. Castaneda, C. Beattie, N. C. Rabinowitz, A. S. Morcos, A. Ruderman, *et al.*, "Human-level performance in 3d multiplayer games with population-based reinforcement learning," *Science*, vol. 364, no. 6443, pp. 859–865, 2019.

[119] K. Zhang, Z. Yang, H. Liu, T. Zhang, and T. Basar, "Fully decentralized multi-agent reinforcement learning with networked agents," in *International Conference on Machine Learning*, pp. 5872–5881, PMLR, 2018.

[120] W. Mao, L. Yang, K. Zhang, and T. Basar, "On improving model-free algorithms for decentralized multi-agent reinforcement learning," in *International Conference on Machine Learning*, pp. 15007–15049, PMLR, 2022.

[121] M. Tan, "Multi-agent reinforcement learning: Independent vs. cooperative agents," in *International Conference on Machine Learning*, pp. 330–337, 1993.

[122] C. S. de Witt, T. Gupta, D. Makoviichuk, V. Makoviychuk, P. H. Torr, M. Sun, and S. Whiteson, "Is independent learning all you need in the starcraft multi-agent challenge?," *arXiv preprint arXiv:2011.09533*, 2020.

[123] L. Matignon, G. J. Laurent, and N. Le Fort-Piat, "Hysteretic q-learning: an algorithm for decentralized reinforcement learning in cooperative multi-agent teams," in *2007 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pp. 64–69, IEEE, 2007.

[124] J. Foerster, N. Nardelli, G. Farquhar, P. H. Torr, P. Kohli, and S. Whiteson, "Stabilising experience replay for deep multi-agent reinforcement learning," in *International Conference on Machine Learning*, pp. 1146–1155, PMLR, 2017.

[125] G. Papoudakis, F. Christianos, A. Rahman, and S. V. Albrecht, "Dealing with non-stationarity in multi-agent deep reinforcement learning," *arXiv preprint arXiv:1906.04737*, 2019.

[126] J. Wu, X. Sun, A. Zeng, S. Song, S. Rusinkiewicz, and T. Funkhouser, "Spatial intention maps for multi-agent mobile manipulation," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 8749–8756, IEEE, 2021.

[127] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[128] S. Omidshafiei, J. Pazis, C. Amato, J. P. How, and J. Vian, "Deep decentralized multi-task multi-agent reinforcement learning under partial observability," in *International Conference on Machine Learning*, pp. 2681–2690, PMLR, 2017.

[129] C. Zhang, S. Jin, W. Xue, X. Xie, S. Chen, and R. Chen, "Independent reinforcement learning for weakly cooperative multiagent traffic control problem," *IEEE Transactions on Vehicular Technology*, vol. 70, no. 8, pp. 7426–7436, 2021.

[130] R. E. Wang, M. Everett, and J. P. How, "R-maddpg for partially observable environments and limited communication," *arXiv preprint arXiv:2002.06684*, 2020.

[131] K. Jiang, W. Liu, Y. Wang, L. Dong, and C. Sun, "Credit assignment in heterogeneous multi-agent reinforcement learning for fully cooperative tasks," *Applied Intelligence*, vol. 53, no. 23, pp. 29205–29222, 2023.

[132] T. Zhang, Y. Li, C. Wang, G. Xie, and Z. Lu, "Fop: Factorizing optimal joint policy of maximum-entropy multi-agent reinforcement learning," in *International conference on machine learning*, pp. 12491–12500, PMLR, 2021.

[133] M. Zhou, Z. Liu, P. Sui, Y. Li, and Y. Y. Chung, "Learning implicit credit assignment for cooperative multi-agent reinforcement learning," *Advances in neural information processing systems*, vol. 33, pp. 11853–11864, 2020.

[134] B. Peng, T. Rashid, C. Schroeder de Witt, P.-A. Kamienny, P. Torr, W. Böhmer, and S. Whiteson, "Facmac: Factored multi-agent centralised policy gradients," *Advances in Neural Information Processing Systems*, vol. 34, pp. 12208–12221, 2021.

[135] W. Chen, W. Li, X. Liu, S. Yang, and Y. Gao, "Learning explicit credit assignment for cooperative multi-agent reinforcement learning via polarization policy gradient," in *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 11542–11550, 2023.

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

[137] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[138] S. Kar, J. M. F. Moura, and H. V. Poor, "$\mathcal{QD}$-learning: A collaborative distributed strategy for multi-agent reinforcement learning through consensus $+$ innovations," *IEEE Transactions on Signal Processing*, vol. 61, no. 7, pp. 1848–1862, 2013.

[139] G. Qu, Y. Lin, A. Wierman, and N. Li, "Scalable multi-agent reinforcement learning for networked systems with average reward," *Advances in Neural Information Processing Systems*, vol. 33, pp. 2074–2086, 2020.

[140] V. Egorov and A. Shpilman, "Scalable multi-agent model-based reinforcement learning," in *Autonomous Agents and Multiagent Systems*, pp. 381–390, 2022.

[141] D. Ha and J. Schmidhuber, "Recurrent world models facilitate policy evolution," in *Advances in Neural Information Processing Systems*, vol. 31, 2018.

[142] C. Zhu, M. Dastani, and S. Wang, "A survey of multi-agent deep reinforcement learning with communication," *Autonomous Agents and Multi-Agent Systems*, vol. 38, no. 1, p. 4, 2024.

[143] S. Bhattacharya, S. Kailas, S. Badyal, S. Gil, and D. Bertsekas, "Multiagent reinforcement learning: Rollout and policy iteration for pomdp with application to multirobot problems," *IEEE Transactions on Robotics*, vol. 40, pp. 2003–2023, 2024.

[144] Y. Liu, W. Wang, Y. Hu, J. Hao, X. Chen, and Y. Gao, "Multi-agent game abstraction via graph attention neural network," in *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 7211–7218, 2020.

[145] W. Kim, M. Cho, and Y. Sung, "Message-dropout: An efficient training method for multi-agent deep reinforcement learning," in *Proceedings of the AAAI conference on artificial intelligence*, pp. 6079–6086, 2019.

[146] J. Jiang and Z. Lu, "Learning attentional communication for multi-agent cooperation," *Advances in neural information processing systems*, vol. 31, 2018.

[147] T. Chu, S. Chinchali, and S. Katti, "Multi-agent reinforcement learning for networked system control," *arXiv preprint arXiv:2004.01339*, 2020.

[148] W. Kim, J. Park, and Y. Sung, "Communication in multi-agent reinforcement learning: Intention sharing," in *International Conference on Learning Representations*, 2020.

[149] Y. Niu, R. R. Paleja, and M. C. Gombolay, "Multi-agent graph-attention communication and teaming.," in *AAMAS*, pp. 964–973, 2021.

[150] S. Sukhbaatar, R. Fergus, *et al.*, "Learning multiagent communication with backpropagation," *Advances in neural information processing systems*, vol. 29, 2016.

[151] B. Freed, R. James, G. Sartoretti, and H. Choset, "Sparse discrete communication learning for multi-agent cooperation through backpropagation," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 7993–7998, 2020.

[152] J. Foerster, I. A. Assael, N. De Freitas, and S. Whiteson, "Learning to communicate with deep multi-agent reinforcement learning," *Advances in neural information processing systems*, vol. 29, 2016.

[153] N. Gupta, G. Srinivasaraghavan, S. Mohalik, N. Kumar, and M. E. Taylor, "Hammer: Multi-level coordination of reinforcement learning agents via learned messaging," *Neural Computing and Applications*, pp. 1–16, 2023.

[154] G. Sartoretti, J. Kerr, Y. Shi, G. Wagner, T. K. S. Kumar, S. Koenig, and H. Choset, "Primal: Pathfinding via reinforcement and imitation multi-agent learning," *IEEE Robotics and Automation Letters*, vol. 4, no. 3, pp. 2378–2385, 2019.

[155] 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*, pp. 1928–1937, PMLR, 2016.

[156] 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.

[157] Z. Qin, K. Zhang, Y. Chen, J. Chen, and C. Fan, "Learning safe multi-agent control with decentralized neural barrier certificates," *International Conference on Learning Representations*, 2021.

[158] P. Sunehag, G. Lever, A. Gruslys, W. M. Czarnecki, V. Zambaldi, M. Jaderberg, M. Lanctot, N. Sonnerat, J. Z. Leibo, K. Tuyls, *et al.*, "Value-decomposition networks for cooperative multi-agent learning based on team reward," in *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 2085–2087, 2018.

[159] C. Huang, J. Zhao, H. Zhou, H. Zhang, X. Zhang, and C. Ye, "Multi-agent decision-making at unsignalized intersections with reinforcement learning from demonstrations," in *2023 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1–6, 2023.

[160] E. Marchesini and A. Farinelli, "Centralizing state-values in dueling networks for multi-robot reinforcement learning mapless navigation," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 4583–4588, 2021.

[161] H. Van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," in *Proceedings of the AAAI conference on artificial intelligence*, 2016.

[162] R. Lowe, Y. I. Wu, A. Tamar, J. Harb, O. Pieter Abbeel, and I. Mordatch, "Multi-agent actor-critic for mixed cooperative-competitive environments," *Advances in neural information processing systems*, vol. 30, 2017.

[163] C. Yu, A. Velu, E. Vinitsky, J. Gao, Y. Wang, A. Bayen, and Y. Wu, "The surprising effectiveness of ppo in cooperative multi-agent games," *Advances in Neural Information Processing Systems*, vol. 35, pp. 24611–24624, 2022.

[164] W. Zaremba, I. Sutskever, and O. Vinyals, "Recurrent neural network regularization," *arXiv preprint arXiv:1409.2329*, 2014.

[165] C. Daskalakis, D. J. Foster, and N. Golowich, "Independent policy gradient methods for competitive reinforcement learning," *Advances in neural information processing systems*, vol. 33, pp. 5527–5540, 2020.

[166] Z. Liu, B. Chen, H. Zhou, G. Koushik, M. Hebert, and D. Zhao, "Mapper: Multi-agent path planning with evolutionary reinforcement learning in mixed dynamic environments," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 11748–11754, 2020.

[167] P. E. Hart, N. J. Nilsson, and B. Raphael, "A formal basis for the heuristic determination of minimum cost paths," *IEEE Transactions on Systems Science and Cybernetics*, vol. 4, no. 2, pp. 100–107, 1968.

[168] B. Wang, Z. Liu, Q. Li, and A. Prorok, "Mobile robot path planning in dynamic environments through globally guided reinforcement learning," *IEEE Robotics and Automation Letters*, vol. 5, no. 4, pp. 6932–6939, 2020.

[169] M. Damani, Z. Luo, E. Wenzel, and G. Sartoretti, "Primal$_2$: Pathfinding via reinforcement and imitation multi-agent learning - lifelong," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 2666–2673, 2021.

[170] W. Schwarting, A. Pierson, J. Alonso-Mora, S. Karaman, and D. Rus, "Social behavior for autonomous vehicles," *Proceedings of the National Academy of Sciences*, vol. 116, no. 50, pp. 24972–24978, 2019.

[171] E. Jang, S. Gu, and B. Poole, "Categorical reparameterization with gumbel-softmax," in *International Conference on Learning Representations*, 2016.

[172] L. Brunke, M. Greeff, A. W. Hall, Z. Yuan, S. Zhou, J. Panerati, and A. P. Schoellig, "Safe learning in robotics: From learning-based control to safe reinforcement learning," *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 5, pp. 411–444, 2022.

[173] S. Gu, J. G. Kuba, Y. Chen, Y. Du, L. Yang, A. Knoll, and Y. Yang, "Safe multi-agent reinforcement learning for multi-robot control," *Artificial Intelligence*, vol. 319, p. 103905, 2023.

[174] D. Ying, Y. Zhang, Y. Ding, A. Koppel, and J. Lavaei, "Scalable primal-dual actor-critic method for safe multi-agent rl with general utilities," *Advances in Neural Information Processing Systems*, vol. 36, 2024.

[175] W. Zhao, T. He, R. Chen, T. Wei, and C. Liu, "State-wise safe reinforcement learning: A survey," *arXiv preprint arXiv:2302.03122*, 2023.

[176] W. Zhang, O. Bastani, and V. Kumar, "Mamps: Safe multi-agent reinforcement learning via model predictive shielding," *arXiv preprint arXiv:1910.12639*, 2019.

[177] B. Toghi, R. Valiente, D. Sadigh, R. Pedarsani, and Y. P. Fallah, "Social coordination and altruism in autonomous driving," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 12, pp. 24791–24804, 2022.

[178] E. Candela, L. Parada, L. Marques, T.-A. Georgescu, Y. Demiris, and P. Angeloudis, "Transferring multi-agent reinforcement learning policies for autonomous driving using sim-to-real," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 8814–8820, 2022.

[179] L. Chen, Y. Wang, Z. Miao, Y. Mo, M. Feng, Z. Zhou, and H. Wang, "Transformer-based imitative reinforcement learning for multirobot path planning," *IEEE Transactions on Industrial Informatics*, vol. 19, no. 10, pp. 10233–10243, 2023.

[180] Z. Zheng and S. Gu, "Safe multi-agent reinforcement learning with bilevel optimization in autonomous driving," *arXiv preprint arXiv:2405.18209*, 2024.

[181] S. Han, S. Zhou, J. Wang, L. Pepin, C. Ding, J. Fu, and F. Miao, "A multi-agent reinforcement learning approach for safe and efficient behavior planning of connected autonomous vehicles," *IEEE Transactions on Intelligent Transportation Systems*, vol. 25, no. 5, pp. 3654–3670, 2024.

[182] B. Cai, C. Wei, and Z. Ji, "Deep reinforcement learning with multiple unrelated rewards for agv mapless navigation," *IEEE Transactions on Automation Science and Engineering*, pp. 1–18, 2024.

[183] R. Zhao, Y. Li, F. Gao, Z. Gao, and T. Zhang, "Multi-agent constrained policy optimization for conflict-free management of connected autonomous vehicles at unsignalized intersections," *IEEE Transactions on Intelligent Transportation Systems*, vol. 25, no. 6, pp. 5374–5388, 2024.

[184] I. ElSayed-Aly, S. Bharadwaj, C. Amato, R. Ehlers, U. Topcu, and L. Feng, "Safe multi-agent reinforcement learning via shielding," in *Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 483–491, 2021.

[185] A. D. Ames, J. W. Grizzle, and P. Tabuada, "Control barrier function based quadratic programs with application to adaptive cruise control," in *53rd IEEE Conference on Decision and Control*, pp. 6271–6278, 2014.

[186] M. Srinivasan and S. Coogan, "Control of mobile robots using barrier functions under temporal logic specifications," *IEEE Transactions on Robotics*, vol. 37, no. 2, pp. 363–374, 2021.

[187] C. Dawson, S. Gao, and C. Fan, "Safe control with learned certificates: A survey of neural lyapunov, barrier, and contraction methods for robotics and control," *IEEE Transactions on Robotics*, vol. 39, no. 3, pp. 1749–1767, 2023.

[188] A. Singh, "Transformer-based sensor fusion for autonomous driving: A survey," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 3312–3317, 2023.

[189] X. Bai, Z. Hu, X. Zhu, Q. Huang, Y. Chen, H. Fu, and C.-L. Tai, "Transfusion: Robust lidar-camera fusion for 3d object detection with transformers," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 1090–1099, 2022.

[190] X. Chen, T. Zhang, Y. Wang, Y. Wang, and H. Zhao, "Futr3d: A unified sensor fusion framework for 3d detection," in *proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 172–181, 2023.

[191] 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.

[192] A. Brunnbauer, L. Berducci, A. Brandstätter, M. Lechner, R. Hasani, D. Rus, and R. Grosu, "Latent imagination facilitates zero-shot transfer in autonomous racing," in *2022 International Conference on Robotics and Automation (ICRA)*, pp. 7513–7520, 2022.

[193] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," in *European Conference on Computer Vision*, pp. 405–421, 2020.

[194] T. Liang, H. Xie, K. Yu, Z. Xia, Z. Lin, Y. Wang, T. Tang, B. Wang, and Z. Tang, "Bevfusion: A simple and robust lidar-camera fusion framework," *Advances in Neural Information Processing Systems*, vol. 35, pp. 10421–10434, 2022.

[195] Z. Liu, H. Tang, A. Amini, X. Yang, H. Mao, D. L. Rus, and S. Han, "Bevfusion: Multi-task multi-sensor fusion with unified bird's-eye view representation," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2774–2781, 2023.

[196] T. Wang, F. Lu, Z. Zheng, G. Chen, and C. Jiang, "Rcdn: Towards robust camera-insensitivity collaborative perception via dynamic feature-based 3d neural modeling," *arXiv preprint arXiv:2405.16868*, 2024.

[197] Y. Zhang, B. Kang, B. Hooi, S. Yan, and J. Feng, "Deep long-tailed learning: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 9, pp. 10795–10816, 2023.

[198] S. Akhauri, L. Y. Zheng, and M. C. Lin, "Enhanced transfer learning for autonomous driving with systematic accident simulation," in *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 5986–5993, 2020.

[199] S. Suo, S. Regalado, S. Casas, and R. Urtasun, "Trafficsim: Learning to simulate realistic multi-agent behaviors," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10400–10409, 2021.

[200] A. Gambi, M. Mueller, and G. Fraser, "Automatically testing self-driving cars with search-based procedural content generation," in *Proceedings of the 28th ACM SIGSOFT International Symposium on Software Testing and Analysis*, pp. 318–328, 2019.

[201] F. Hauer, T. Schmidt, B. Holzmüller, and A. Pretschner, "Did we test all scenarios for automated and autonomous driving systems?," in *2019 IEEE Intelligent Transportation Systems Conference (ITSC)*, pp. 2950–2955, 2019.

[202] L. Wang, X. Zhang, H. Su, and J. Zhu, "A comprehensive survey of continual learning: Theory, method and application," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–20, 2024.

[203] F. Zhuang, Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He, "A comprehensive survey on transfer learning," *Proceedings of the IEEE*, vol. 109, no. 1, pp. 43–76, 2021.

[204] X. Hu, S. Li, T. Huang, B. Tang, R. Huai, and L. Chen, "How simulation helps autonomous driving: A survey of sim2real, digital twins, and parallel intelligence," *IEEE Transactions on Intelligent Vehicles*, vol. 9, no. 1, pp. 593–612, 2024.

[205] J.-M. Georg, J. Feiler, S. Hoffmann, and F. Diermeyer, "Sensor and actuator latency during teleoperation of automated vehicles," in *2020 IEEE Intelligent Vehicles Symposium (IV)*, pp. 760–766, 2020.

[206] L. Montaut, Q. L. Lidec, A. Bambade, V. Petrik, J. Sivic, and J. Carpentier, "Differentiable collision detection: a randomized smoothing approach," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 3240–3246, 2023.

[207] J. Li, Z. Yu, Z. Du, L. Zhu, and H. T. Shen, "A comprehensive survey on source-free domain adaptation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–22, 2024.

[208] K. Arndt, M. Hazara, A. Ghadirzadeh, and V. Kyrki, "Meta reinforcement learning for sim-to-real domain adaptation," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 2725–2731, 2020.

[209] Z. Liu, G. Chen, Z. Li, Y. Kang, S. Qu, and C. Jiang, "Psdc: A prototype-based shared-dummy classifier model for open-set domain adaptation," *IEEE Transactions on Cybernetics*, vol. 53, no. 11, pp. 7353–7366, 2023.

[210] S. Tan, K. Wong, S. Wang, S. Manivasagam, M. Ren, and R. Urtasun, "Scenegen: Learning to generate realistic traffic scenes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 892–901, 2021.

[211] U. Briefs, "Mcity grand opening," *Research Review*, 2015.

[212] A. R. Javed, M. Usman, S. U. Rehman, M. U. Khan, and M. S. Haghighi, "Anomaly detection in automated vehicles using multistage attention-based convolutional neural network," *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 7, pp. 4291–4300, 2021.

[213] P. K. Sharma and J. H. Park, "Blockchain based hybrid network architecture for the smart city," *Future Generation Computer Systems*, vol. 86, pp. 650–655, 2018.

[214] M. T. Ribeiro, S. Singh, and C. Guestrin, "" why should i trust you?" explaining the predictions of any classifier," in *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144, 2016.

[215] C. Brewitt, B. Gyevnar, S. Garcin, and S. V. Albrecht, "Grit: Fast, interpretable, and verifiable goal recognition with learned decision trees for autonomous driving," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 1023–1030, 2021.

[216] A. Likmeta, A. M. Metelli, A. Tirinzoni, R. Giol, M. Restelli, and D. Romano, "Combining reinforcement learning with rule-based controllers for transparent and general decision-making in autonomous driving," *Robotics and Autonomous Systems*, vol. 131, p. 103568, 2020.

[217] A. Heuillet, F. Couthouis, and N. Díaz-Rodríguez, "Explainability in deep reinforcement learning," *Knowledge-Based Systems*, vol. 214, p. 106685, 2021.

[218] A. Raffin, A. Hill, R. Traoré, T. Lesort, N. Díaz-Rodríguez, and D. Filliat, "Decoupling feature extraction from policy learning: assessing benefits of state representation learning in goal based robotics," in *ICLR Workshop on Structure and Priors in Reinforcement Learning*, 2019.

[219] R. Traoré, H. Caselles-Dupré, T. Lesort, T. Sun, G. Cai, D. Filliat, and N. Díaz-Rodríguez, "Discorl: Continual reinforcement learning via policy distillation," in *NeurIPS Workshop on Deep Reinforcement Learning*, 2019.

[220] A. Krajna, M. Brcic, T. Lipic, and J. Doncevic, "Explainability in reinforcement learning: perspective and position," *arXiv preprint arXiv:2203.11547*, 2022.

[221] J. Kim and J. Canny, "Interpretable learning for self-driving cars by visualizing causal attention," in *Proceedings of the IEEE international conference on computer vision*, pp. 2942–2950, 2017.

[222] J. Kim, A. Rohrbach, T. Darrell, J. Canny, and Z. Akata, "Textual explanations for self-driving vehicles," in *Proceedings of the European conference on computer vision (ECCV)*, pp. 563–578, 2018.

[223] K. Chitta, A. Prakash, B. Jaeger, Z. Yu, K. Renz, and A. Geiger, "Transfuser: Imitation with transformer-based sensor fusion for autonomous driving," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 45, no. 11, pp. 12878–12895, 2023.

[224] X. Zhao, W. Huang, X. Huang, V. Robu, and D. Flynn, "Baylime: Bayesian local interpretable model-agnostic explanations," in *Uncertainty in artificial intelligence*, pp. 887–896, PMLR, 2021.

[225] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi, "Dream to control: Learning behaviors by latent imagination," in *International Conference on Learning Representations*, 2019.

[226] D. Hafner, T. P. Lillicrap, M. Norouzi, and J. Ba, "Mastering atari with discrete world models," in *International Conference on Learning Representations*, 2020.

[227] X. Fang, Q. Zhang, Y. Gao, and D. Zhao, "Offline reinforcement learning for autonomous driving with real world driving data," in *2022 IEEE 25th International Conference on Intelligent Transportation Systems (ITSC)*, pp. 3417–3422, 2022.

[228] M. Reuse, M. Simon, and B. Sick, "About the ambiguity of data augmentation for 3d object detection in autonomous driving," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 979–987, 2021.

[229] X. Zhang, N. Tseng, A. Syed, R. Bhasin, and N. Jaipuria, "Simbar: Single image-based scene relighting for effective data augmentation for automated driving vision tasks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3718–3728, 2022.

[230] Z. Zheng, Y. Cheng, Z. Xin, Z. Yu, and B. Zheng, "Robust perception under adverse conditions for autonomous driving based on data augmentation," *IEEE Transactions on Intelligent Transportation Systems*, vol. 24, no. 12, pp. 13916–13929, 2023.

[231] H. Lin, W. Ding, Z. Liu, Y. Niu, J. Zhu, Y. Niu, and D. Zhao, "Safety-aware causal representation for trustworthy offline reinforcement learning in autonomous driving," *IEEE Robotics and Automation Letters*, vol. 9, no. 5, pp. 4639–4646, 2024.

[232] L. Meng, M. Wen, C. Le, X. Li, D. Xing, W. Zhang, Y. Wen, H. Zhang, J. Wang, Y. Yang, *et al.*, "Offline pre-trained multi-agent decision transformer," *Machine Intelligence Research*, vol. 20, no. 2, pp. 233–248, 2023.

[233] X. Wu, L. Xiao, Y. Sun, J. Zhang, T. Ma, and L. He, "A survey of human-in-the-loop for machine learning," *Future Generation Computer Systems*, vol. 135, pp. 364–381, 2022.

[234] J. Wu, Z. Huang, C. Huang, Z. Hu, P. Hang, Y. Xing, and C. Lv, "Human-in-the-loop deep reinforcement learning with application to autonomous driving," *arXiv preprint arXiv:2104.07246*, 2021.

[235] Z. Peng, Q. Li, C. Liu, and B. Zhou, "Safe driving via expert guided policy optimization," in *Conference on Robot Learning*, pp. 1554–1563, PMLR, 2022.

[236] Q. Li, Z. Peng, and B. Zhou, "Efficient learning of safe driving policy via human-ai copilot optimization," *arXiv preprint arXiv:2202.10341*, 2022.

[237] Y. Yuan, J. Hao, Y. Ma, Z. Dong, H. Liang, J. Liu, Z. Feng, K. Zhao, and Y. Zheng, "Uni-rlhf: Universal platform and benchmark suite for reinforcement learning with diverse human feedback," in *International Conference on Learning Representations*, 2024.

[238] L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch, "Decision transformer: Reinforcement learning via sequence modeling," *Advances in neural information processing systems*, vol. 34, pp. 15084–15097, 2021.

[239] M. Janner, Q. Li, and S. Levine, "Offline reinforcement learning as one big sequence modeling problem," *Advances in neural information processing systems*, vol. 34, pp. 1273–1286, 2021.

[240] C. Sun, S. Huang, and D. Pompili, "Llm-based multi-agent reinforcement learning: Current and future directions," *arXiv preprint arXiv:2405.11106*, 2024.

[241] J. D. M.-W. C. Kenton and L. K. Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," in *Proceedings of NAACL-HLT*, pp. 4171–4186, 2019.

[242] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar, *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.

[243] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altenschmidt, S. Altman, S. Anadkat, *et al.*, "Gpt-4 technical report," *arXiv preprint arXiv:2303.08774*, 2023.

[244] X. Zhou, M. Liu, B. L. Zagar, E. Yurtsever, and A. C. Knoll, "Vision language models in autonomous driving and intelligent transportation systems," *arXiv preprint arXiv:2310.14414*, 2023.

[245] L. Wen, X. Yang, D. Fu, X. Wang, P. Cai, X. Li, M. Tao, Y. Li, X. Linran, D. Shang, *et al.*, "On the road with gpt-4v (ision): Explorations of utilizing visual-language model as autonomous driving agent," in *ICLR 2024 Workshop on Large Language Model (LLM) Agents*, 2024.

[246] Z. Xu, Y. Zhang, E. Xie, Z. Zhao, Y. Guo, K. K. Wong, Z. Li, and H. Zhao, "Drivegpt4: Interpretable end-to-end autonomous driving via large language model," *arXiv preprint arXiv:2310.01412*, 2023.

[247] H. Sha, Y. Mu, Y. Jiang, L. Chen, C. Xu, P. Luo, S. E. Li, M. Tomizuka, W. Zhan, and M. Ding, "Languagempc: Large language models as decision makers for autonomous driving," *arXiv preprint arXiv:2310.03026*, 2023.

[248] C. Cui, Y. Ma, X. Cao, W. Ye, and Z. Wang, "Drive as you speak: Enabling human-like interaction with large language models in autonomous vehicles," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 902–909, 2024.

[249] Y. Huang, Y. Chen, and Z. Li, "Applications of large scale foundation models for autonomous driving," *arXiv preprint arXiv:2311.12144*, 2023.

[250] Y. Li, W. Zhang, K. Chen, Y. Liu, P. Li, R. Gao, L. Hong, M. Tian, X. Zhao, Z. Li, *et al.*, "Automated evaluation of large vision-language models on self-driving corner cases," *arXiv preprint arXiv:2404.10595*, 2024.

[251] C. Pan, B. Yaman, T. Nesti, A. Mallik, A. G. Allievi, S. Velipasalar, and L. Ren, "Vlp: Vision language planning for autonomous driving," *arXiv preprint arXiv:2401.05577*, 2024.

[252] C. Snell, D. Klein, and R. Zhong, "Learning by distilling context," *arXiv preprint arXiv:2209.15189*, 2022.