

# A Comprehensive Review of Reinforcement Learning for Autonomous Driving in the CARLA Simulator

Elahe Delavari<sup>1†</sup>, Feeza Khan Khanzada<sup>1†</sup>, Jaerock Kwon<sup>1\*</sup>

<sup>1\*</sup>Electrical and Computer Engineering, University Of Michigan-Dearborn, 4901 Evergreen Road, Dearborn, 48128, MI, USA.

\*Corresponding author(s). E-mail(s): [jrkwon@umich.edu](mailto:jrkwon@umich.edu);  
Contributing authors: [elahed@umich.edu](mailto:elahed@umich.edu); [feezakk@umich.edu](mailto:feezakk@umich.edu);

<sup>†</sup>These authors contributed equally to this work.

## Abstract

Autonomous-driving research has recently embraced deep Reinforcement Learning (RL) as a promising framework for data-driven decision making, yet a clear picture of how these algorithms are currently employed, benchmarked and evaluated is still missing. This survey fills that gap by systematically analysing around 100 peer-reviewed papers that train, test or validate RL policies inside the open-source CARLA simulator. We first categorize the literature by algorithmic family—model-free, model-based, hierarchical, and hybrid—and quantify their prevalence, highlighting that more than 80% of existing studies still rely on model-free methods such as DQN, PPO and SAC. Next, we explain the diverse state, action and reward formulations adopted across works, illustrating how choices of sensor modality (RGB, LiDAR, BEV, semantic maps, and carla kinematics states), control abstraction (discrete vs. continuous) and reward shaping are used across various literature. We also consolidate the evaluation landscape by listing the most common metrics (success rate, collision rate, lane deviation, driving score) and the towns, scenarios and traffic configurations used in CARLA benchmarks. Persistent challenges including sparse rewards, sim-to-real transfer, safety guarantees and limited behaviour diversity are distilled into a set of open research questions, and promising directions such as model-based RL, meta-learning and richer multi-agent simulations are outlined. By providing a unified taxonomy, quantitative statistics and a critical discussion of limitations, this review aims to serve both as a reference for newcomers and as a roadmap for advancing RL-based autonomous driving toward real-world deployment.

**Keywords:** Reinforcement Learning, Autonomous Driving, CARLA Simulator

## 1 Introduction

Autonomous driving represents one of the most transformative technological frontiers of the 21st century, with the potential to revolutionize transportation by enhancing safety, efficiency, and accessibility. By reducing human error—widely recognized as the leading cause of traffic accidents [1]—Autonomous Vehicles (AVs) promise safer roads, reduced congestion, and improved mobility. However, the development of robust and reliable AV systems remains a formidable challenge, primarily due to the complexity of real-world driving environments. These environments present dynamic traffic patterns, intricate multi-agent interactions, and rare but critical long-tail scenarios that traditional rule-based approaches struggle to handle effectively.

Reinforcement Learning (RL) has emerged as a powerful paradigm for training autonomous driving policies, offering an adaptive framework where agents learn by interacting with their environment and refining their decision-making strategies over time [2]. Unlike conventional supervised learning approaches that require extensive labeled datasets, RL enables AVs to optimize actions dynamically based on rewards, making it particularly suited for sequential decision-making tasks such as lane-keeping, obstacle avoidance, intersection handling, and route planning. Despite its potential, deploying RL-based autonomous systems in the real world poses several challenges, including safety risks, computational demands, and the need for active environmental interaction for training.

To address these challenges, simulation platforms have become an indispensable tool for RL research in autonomous driving. Among these, CARLA (Car Learning to Act) [3] has established itself as a premier open-source simulator designed specifically for AV research. Built on Unreal Engine, CARLA provides high-fidelity urban and highway driving environments, supporting diverse road layouts, dynamic traffic conditions, various weather scenarios, and an extensive suite of sensors (e.g., RGB cameras, LiDAR, depth sensors, and semantic segmentation). This flexibility enables safe, controlled, and scalable RL training, allowing researchers to benchmark and refine driving policies before real-world deployment.

Over the years, CARLA has become a cornerstone for RL-driven autonomous driving research, fostering advancements across multiple dimensions. Researchers have leveraged a wide array of RL techniques [4–20], including model-free and model-based approaches for handling variety of driving tasks. Furthermore, multimodal RL frameworks have been explored to incorporate rich sensor data, improving decision-making under uncertainty. While these methods have significantly advanced the field, challenges such as sparse reward signals, generalization across diverse environments, and the sim-to-real transfer gap remain open problems that demand further investigation.

This paper provides a comprehensive review of RL research within the CARLA simulator, analyzing the evolution of RL techniques and their applications in autonomous driving. The key contributions of this review are as follows:

- **Survey of RL Formulations** – We explore various RL approaches employed in CARLA, including model-free methods such as the Deep Q-Network (DQN), Proximal Policy Optimization (PPO), and Soft Actor-Critic (SAC), Deep Deterministic Policy Gradient (DDPG); model-based techniques; hierarchical RL; and hybrid strategies that integrate multiple learning paradigms.
- **Comparison of State and Action Spaces** – We examine how different sensor modalities (camera, LiDAR, multi-modal fusion) and control abstractions (low-level throttle/steering vs. high-level discrete maneuvers) impact policy learning and performance.
- **Reward Engineering and Terminal condition in RL for Autonomous Driving** – We highlight the critical role of reward function design, analyzing how different formulations balance safety, comfort, and efficiency in various driving tasks.
- **Evaluation metrics and Scenarios for RL Policies** – We review the evaluation methodologies used in CARLA, including key performance metrics such as collision rates, lane deviation, speed regulation, and success rates in navigation tasks along with the different scenarios used for training and testing.
- **Challenges and Future Directions** – We discuss persistent challenges in RL-driven autonomous driving often mentioned in the literature, including generalization to unseen environments, robustness in high-density traffic scenarios, and interpretability of learned strategies. Additionally, we outline promising avenues for future research to bridge these gaps.

By synthesizing recent advancements and identifying key challenges, this review aims to serve as a valuable resource for researchers and practitioners in the field. As simulation platforms like CARLA continue to evolve, they will play an increasingly critical role in the pursuit of safe and reliable autonomous driving systems. Through this work, we hope to provide clarity on the state of RL in CARLA, fostering further innovation and accelerating progress toward real-world deployment. This review focuses on single-agent RL approaches in the CARLA simulator. While Multi-Agent RL (MARL) is a relevant and growing area of research—particularly for urban traffic scenarios and coordination among multiple vehicles—it is beyond the scope of this review. we consider around 100 papers from IEEE and ACM.

## 2 CARLA simulator

CARLA is an open-source simulator explicitly designed to support research in autonomous driving and ADAS. Built on the Unreal Engine 4 <sup>1</sup>, CARLA provides a photorealistic 3D environment with configurable urban and highway scenes that include diverse road layouts, dynamic weather, and day-night cycles. Crucially, it exposes a flexible Python/C++ API for scene manipulation, traffic generation, and sensor configuration, enabling researchers to script complex scenarios and integrate with ROS [21], Autoware [22], or Apollo [22] stacks.

A key strength of CARLA lies in its extensible sensor suite. Out of the box, users can deploy RGB/depth cameras, LiDAR, radar, GPS, IMU, and semantic segmentation sensors—each with fully customizable intrinsic and extrinsic parameters. Custom

---

<sup>1</sup><https://www.unrealengine.com/>

sensors may also be defined via JSON, and ground-truth data (e.g., 2D/3D bounding boxes, semantic labels) can be obtained directly from the environment. This level of fidelity allows for rigorous evaluation of perception and sensor-fusion algorithms under controlled yet realistic conditions.

Underpinning CARLA’s realism is its client–server architecture. The server handles all physics, rendering, and world updates—leveraging GPU acceleration for accurate vehicle dynamics (including wheel friction, suspension, and center-of-mass modeling) and environment effects—while lightweight clients manage agent logic and scenario control. A built-in recorder facilitates seamless capture and replay of both sensor streams and control commands, supporting reproducible experiments and accelerated headless simulations.

In the comprehensive survey by Kaur et al. [23], CARLA and LGSVL [24] emerge as the state-of-the-art open-source simulators for end-to-end autonomous-vehicle testing, outperforming MATLAB/Simulink, CarSim, PreScan, and Gazebo across key requirements such as high-fidelity 3D environments, sensor variety, and scenario management. Unlike Gazebo—which demands manual creation of models and scene geometry via XML—CARLA supplies urban maps and traffic infrastructure (traffic lights, stop signs, lane markings) automatically from OpenDRIVE files, thus drastically reducing setup overhead. Moreover, CARLA’s seamless integration with major autonomy stacks via native ROS bridges contrasts with the external tooling required by LGSVL and others.

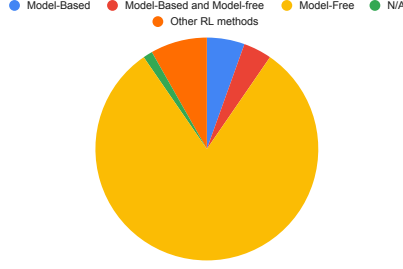
CARLA’s open-source ecosystem further distinguishes it from proprietary platforms (e.g., Waymo’s CarCraft, Uber’s DataViz). Researchers can inspect, extend, or optimize core modules—ranging from physics engines to AI agents—and contribute back improvements. This collaborative model accelerates feature development (e.g., the recent addition of an RSS-based safety assurance module) and ensures that CARLA remains at the forefront of sim-to-real research challenges.

### 3 Reinforcement Learning Methods

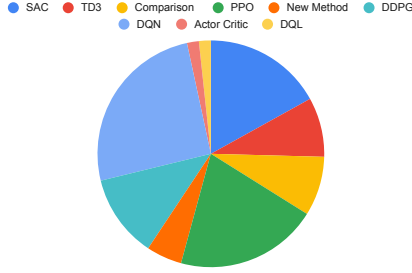
RL methods used in CARLA vary widely in their learning paradigms, exploration strategies, and suitability for different types of control tasks. Broadly, they can be categorized into model-free and model-based approaches, with several hybrid and alternative techniques (e.g., distributional, Bayesian, and hierarchical RL) gaining traction in recent years. Model-free methods like DQN, PPO, and SAC have become standard choices for end-to-end driving due to their empirical robustness, while model-based approaches aim for higher sample efficiency by learning environment dynamics. This section reviews the RL algorithms applied in CARLA-based driving research, analyzing their design, use cases, and key contributions.

Despite the potential for improved sample efficiency, MBRL remains underexplored in CARLA, as shown in Fig. 1. The majority of studies continue to rely on model-free algorithms due to their relative ease of implementation and empirical robustness.

As shown in Fig. 1, model-free approaches dominate the CARLA literature, accounting for over 80 % of all surveyed works, while model-based and hybrid methods remain relatively rare. Fig. 2 drills into that model-free slice—highlighting the shares



**Fig. 1** Distribution of RL model types used in CARLA-based studies. Model-free approaches dominate the literature (over 80%), while model-based and hybrid methods are significantly less explored.



**Fig. 2** Distribution of different model free RL.

of DQN, PPO, SAC, DDPG, etc. To complement these visual summaries, Table 1 provides a comprehensive listing of representative CARLA publications, organized by algorithmic family and method variant. This table serves as a roadmap for the subsections that follow, where we review each category in detail.

### 3.1 Model-Free Reinforcement Learning

Model-free RL (MFRL) methods learn optimal policies directly from experience without explicitly modeling the environment’s transition dynamics. These approaches are particularly well-suited to high-dimensional and complex environments like CARLA, where learning accurate dynamics models can be challenging. Model-free methods are typically divided into value-based approaches (e.g., DQN), which estimate action-value functions, and policy-based or actor-critic methods (e.g., PPO, DDPG, SAC), which directly optimize the policy with or without value function guidance. Although model-free methods often require large amounts of training data, they are widely used for end-to-end driving tasks in the CARLA simulator, as shown in Fig. 1.

#### 3.1.1 Actor Critic

Actor-critic methods combine value-based and policy-based RL by maintaining two main components: an *actor* and a *critic*. The actor is a parameterized policy, often denoted as  $\pi_{\theta}(a | s)$ , which directly maps states to actions or action probabilities. The

critic is a value function approximator, such as  $V_\omega(s)$  or  $Q_\omega(s, a)$ , used to evaluate and guide updates for the actor.

In the standard RL setting, an agent interacts with the environment through states  $s_t \in \mathcal{S}$ , actions  $a_t \in \mathcal{A}$ , and rewards  $r_t \in \mathbb{R}$ . At each timestep  $t$ , the agent observes a state  $s_t$ , selects an action  $a_t$  according to its policy, and receives a reward  $r_t$  along with the next state  $s_{t+1}$ .

In a typical actor-critic setting, the critic parameters  $\omega$  are updated by minimizing a Temporal-Difference (TD) error. For example, if the critic is a value function  $V_\omega$ , the update target for a given transition  $(s_t, a_t, r_t, s_{t+1})$  might be

$$y_t = r_t + \gamma V_\omega(s_{t+1}) \quad (1)$$

and the critic’s objective becomes

$$\mathcal{L}(\omega) = (V_\omega(s_t) - y_t)^2, \quad (2)$$

where  $\gamma$  is the discount factor.

The actor parameters  $\theta$  are updated by following the policy gradient, which often uses an advantage estimate to reduce variance. An example update rule is

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{t=1}^N \nabla_\theta \log \pi_\theta(a_t | s_t) A_t, \quad (3)$$

where  $A_t$  is an advantage function. One simple choice for the advantage is the TD error

$$A_t = r_t + \gamma V_\omega(s_{t+1}) - V_\omega(s_t). \quad (4)$$

Thus, the actor adjusts its policy in directions suggested by the critic’s estimated advantage of taking action  $a_t$  in state  $s_t$ . By updating the critic to better approximate the value function and updating the actor to maximize returns according to the critic’s feedback, actor-critic methods can learn efficiently in environments with large or continuous action spaces. The learned critic provides lower-variance gradient estimates than pure policy gradient methods, and the actor allows the policy to be optimized continuously rather than relying solely on value-based approaches. In the self-driving domain, Chronis *et al.* [25] and Udatha *et al.* [26] have employed actor–critic variants using separate networks for policy and value estimation.

### 3.1.2 Deep Deterministic Policy Gradient

The DDPG algorithm [27] is a model-free, off-policy actor-critic method designed for environments with continuous action spaces. It extends the Deterministic Policy Gradient (DPG) framework [28] by integrating deep neural networks and stabilization techniques introduced in the DQN [29] algorithm. Due to its ability to directly output real-valued actions, DDPG has been widely used in autonomous driving tasks in simulators like CARLA.

DDPG employs two neural networks:

- An **actor** network  $\mu(s|\theta^\mu)$  that deterministically maps a state  $s$  to a continuous action  $a$ .

- A **critic** network  $Q(s, a|\theta^Q)$  that estimates the expected return (Q-value) of taking action  $a$  in state  $s$ .

The critic is trained to minimize the loss between predicted Q-values and target Q-values, defined using the Bellman equation:

$$y_i = r_i + \gamma Q' \left( s_{i+1}, \mu' \left( s_{i+1} | \theta^{\mu'} \right) \middle| \theta^{Q'} \right), \quad (5)$$

where:

- $r_i$  is the reward at timestep  $i$ ,
- $s_{i+1}$  is the next state,
- $Q'$  and  $\mu'$  are target networks for the critic and actor respectively,
- $\gamma \in [0, 1]$  is the discount factor,
- $\theta^{Q'}$ ,  $\theta^{\mu'}$  are the parameters of the target networks.

These target networks are updated using a soft update rule:

$$\theta' \leftarrow \tau \theta + (1 - \tau) \theta', \quad (6)$$

where  $\tau \ll 1$  is a small constant (e.g.,  $\tau = 0.001$ ). The actor is updated by maximizing the critic's estimate of the expected return. Using the chain rule, the policy gradient is computed as:

$$\nabla_{\theta^\mu} J \approx \mathbb{E}_{s \sim \mathcal{D}} \left[ \nabla_a Q(s, a | \theta^Q) \big|_{a=\mu(s)} \nabla_{\theta^\mu} \mu(s | \theta^\mu) \right], \quad (7)$$

where  $\mathcal{D}$  is the replay buffer distribution. To improve stability and exploration, DDPG includes the following components:

- **Replay buffer**: Stores past transitions  $(s_t, a_t, r_t, s_{t+1})$  for off-policy training.
- **Target networks**: Slowly updated copies of the actor and critic to stabilize training.
- **Ornstein–Uhlenbeck process**: Temporally correlated noise added to actions for exploration in physical environments with inertia.
- **Batch normalization**: Applied to input and hidden layers to mitigate covariate shift and scale sensitivity.

Although DDPG was initially successful across various continuous control benchmarks, it is known to suffer from Q-function overestimation and instability during training. Several studies have explored and extended DDPG to address a variety of challenges in autonomous control. Wu et al. [4] applied DDPG to solve high-dimensional control problems, leveraging its suitability for continuous domains. Goel et al. [5] and Youssef et al. [6] adopted standard DDPG setups for adaptive and general AV control tasks. Peng et al. [30] proposed a hybrid Imitation learning method fusing visual information with the additional steering angle calculated by Pure-Pursuit (PP) called IPP-RL framework that combines imitation learning with DDPG, using steering signals from a pure pursuit controller to guide the learning process. The study

by Pérez-Gil et al. [31] presents a comparative analysis of DDPG and DQN, evaluating their control performance in autonomous driving tasks. To enhance learning efficiency, Fu et al. [32] combined DDPG with an Information Bottleneck mechanism, aiming to reduce redundancy in state representation and improve sample efficiency. Additionally, Chen et al. [33] improved DDPG by incorporating self-attention models for visual perception, prioritized experience replay for more effective sampling, and Ornstein-Uhlenbeck noise to encourage better exploration. Zhang et al. [34] proposed a self-learning lane-keeping algorithm using DDPG. Tsai et al. [35] applied DDPG and RDPG (Recurrent Deterministic Policy Gradient) showing that RDPG outperformed DDPG in generalization to unseen scenarios. Doe et al. [36] adopt a canonical actor-critic formulation in which the actor outputs continuous weighting parameters while the critic approximates the state-action value function. They stabilise learning through both an experience-replay buffer and slowly updated target networks. Building on the same foundation, Li et al. [37] also use separate actor and critic networks together with replay and target mechanisms, emphasising the algorithm’s off-policy nature to optimise a deterministic policy in real-valued action spaces. Ahmed et al. [38] likewise employ DDPG, tailoring the actor to generate continuous steering, throttle and brake commands while the critic guides policy updates via Q-value estimates.

### 3.1.3 Deep Q-Learning

Deep Q-Learning (DQL) extends classical Q-learning by using a deep neural network as a function approximator for the action-value function. Instead of storing a Q-table for all state-action pairs, a parameterized neural network  $Q_\theta(s, a)$  is trained to estimate the expected return. Given a transition  $(s, a, r, s')$ , the *target* for this network can be written as:

$$y = r + \gamma \max_{a'} Q_{\theta^-}(s', a'), \quad (8)$$

where  $\gamma$  is the discount factor, and  $Q_{\theta^-}$  is a *target network* whose parameters  $\theta^-$  are periodically updated from the main network  $\theta$ . The loss function for each mini-batch is

$$\mathcal{L}(\theta) = (Q_\theta(s, a) - y)^2. \quad (9)$$

A key mechanism in DQL is the experience replay buffer, which stores past transitions. At each training step, a random mini-batch is sampled from this buffer to break correlation between consecutive samples and stabilize training.

By decoupling the action-value estimation from target computation through a slowly updated target network, and by training on random samples from the replay buffer, DQL mitigates instabilities that arise from training a neural network with correlated data and non-stationary targets. This approach has been highly successful in discrete action domains, particularly in Atari game environments, where it surpasses human-level performance for many titles. In AV domain, Cheng et al. [39] utilize such a DQL architecture to achieve effective and stable learning in their longitudinal control setting.



### 3.1.4 Proximal Policy Optimization

PPO [40] is a model-free, on-policy policy gradient method that achieves stable and efficient policy updates through a first-order approximation to Trust Region Policy Optimization (TRPO). PPO has gained popularity for its simplicity, generality, and strong empirical performance across both discrete and continuous control tasks, making it a common choice in CARLA-based RL research. PPO maintains a stochastic policy  $\pi_\theta(a|s)$ , which is optimized to maximize expected returns by improving the likelihood of advantageous actions. The vanilla policy gradient objective is given by:

$$L^{\text{PG}}(\theta) = \mathbb{E}_t \left[ \log \pi_\theta(a_t|s_t) \hat{A}_t \right], \quad (10)$$

where  $\hat{A}_t$  is an estimator of the advantage function at timestep  $t$ . While this objective works in principle, performing multiple updates using the same batch of data often leads to instability due to large policy changes.

To address this, PPO introduces a clipped surrogate objective that restricts policy updates to a conservative range. Defining the probability ratio

$$r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)}, \quad (11)$$

the clipped objective becomes:

$$L^{\text{CLIP}}(\theta) = \mathbb{E}_t \left[ \min \left( r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right], \quad (12)$$

where  $\epsilon$  is a small constant (typically 0.1 or 0.2) controlling the trust region size. This objective penalizes updates that move the new policy too far from the old one, ensuring stable learning while still enabling improvement when  $\hat{A}_t$  is positive. The complete loss function includes terms for value function fitting and entropy regularization:

$$L^{\text{PPO}}(\theta) = \mathbb{E}_t \left[ L_t^{\text{CLIP}}(\theta) - c_1 L_t^{\text{VF}}(\theta) + c_2 S[\pi_\theta](s_t) \right], \quad (13)$$

where:

- $L_t^{\text{VF}} = (V_\theta(s_t) - V_t^{\text{target}})^2$  is the value function loss,
- $S[\pi_\theta](s_t)$  is the entropy of the policy to encourage exploration,
- $c_1, c_2$  are weighting coefficients.

Each PPO iteration involves collecting rollouts from the environment using the current policy  $\pi_{\theta_{\text{old}}}$ , computing advantages (e.g., via Generalized Advantage Estimation), and performing several epochs of minibatch updates to optimize the objective. PPO avoids the complexity of second-order optimization (as in TRPO) and is easy to integrate with deep architectures.

PPO’s training stability and simplicity have made it the dominant MFRL algorithm in autonomous driving research, leading to its use in a wide range of studies

that build on its core strengths. Building on these advantages, Carton et al. [7] implemented PPO directly for autonomous control tasks, while Mohammed et al. [8] used PPO within the RLlib framework to leverage distributed and scalable training. To improve learning stability and exploration, Deng et al. [41] integrated improved exploration strategies, while a follow-up study [42] employed PPO with clipped objectives and Generalized Advantage Estimation (GAE) with the context-aware state as input. Zhao et al. [43] proposed a fully end-to-end autonomous driving pipeline based on PPO. Wu et al. [44] also employed PPO for control policy training.

The “Roach” expert agent developed by Zhang et al. [45] was trained using PPO, demonstrating its effectiveness for expert-level policy imitation. Advanced perception mechanisms have been integrated with PPO in several studies. For instance, Agarwal et al. [46] incorporated pretrained BEV-SemSeg autoencoders, while Trumpp et al. [47] employed a Recurrent DriveNet with LSTM layers for temporal state propagation. Xing et al. [48] addressed domain generalization by combining PPO with a Cycle-Consistent Variational Autoencoder (VAE), enabling zero-shot policy transfer by disentangling domain-specific and general features. Anzalone et al. [49] explored curriculum learning, applying PPO across a five-stage training progression—from simplified conditions to dense urban scenarios. Finally, Silva et al. [50] trained three distinct PPO-based sub-policies (lane following, turning left, turning right) under a high-level planner, using data augmentation techniques to enhance robustness and generalization.

Albilani et al. [51] embed a PPO variant called GPPO inside a hierarchical Option-Critic framework: expert demonstrations generated by Answer-Set Programming guide low-level option policies, while a high-level controller selects among those options, yielding a two-tier architecture that accelerates learning and improves interpretability. Martínez-Gómez et al. [52] adopt a more conventional single-layer PPO agent whose shared actor-critic network outputs continuous speed references for the ego vehicle, capitalising on PPO’s clipped-surrogate objective to ensure stable updates in dense-traffic scenarios. Jin et al. [53] enhance the baseline with generalised advantage estimation and prioritised experience replay, producing “PPOE,” which speeds convergence and reduces collision frequency in their highway benchmark. Finally, Shi et al. [54] modify PPO for a hybrid high-/low-level decomposition: a convolutional actor-critic network processes deformed occupancy grids and traffic-rule vectors, a hybrid reward couples behavioural and motion objectives, and a small imitation-learning warm-start improves early performance.

### 3.1.5 Trust Region Policy Optimization

TRPO is a policy gradient method that ensures monotonic policy improvement by constraining the step size during updates, making it well-suited for tasks requiring stable learning. Although less commonly used than PPO due to its computational complexity, TRPO has been successfully applied in autonomous driving scenarios that demand precise control. Gutierrez-Moreno et al. [55] conducted a comparative study evaluating TRPO alongside PPO, DQN, and A2C. In their experiments targeting merge scenarios, TRPO achieved the highest success rate (92.9%) and was ultimately

chosen for deployment. This result highlights TRPO’s strength in producing reliable policies for complex decision-making tasks under safety-critical constraints.

### 3.1.6 Deep Q-Network

The Deep Q-Network (DQN) algorithm [29] was the first deep RL method to achieve human-level performance on a wide range of tasks using only high-dimensional raw visual input. It combines Q-learning with deep convolutional neural networks to approximate the action-value function  $Q(s, a)$ , enabling end-to-end learning of control policies from pixels. DQN estimates the optimal action-value function:

$$Q^*(s, a) = \max_{\pi} \mathbb{E} \left[ r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \cdots \right. \\ \left. \mid s_t = s, a_t = a, \pi \right], \quad (14)$$

which represents the maximum expected return after taking action  $a$  in state  $s$ , and thereafter following policy  $\pi$ . The action-value function is approximated by a neural network  $Q(s, a; \theta)$  with parameters  $\theta$ , trained to minimize the temporal difference (TD) error using the loss:

$$L_i(\theta_i) = \mathbb{E}_{(s, a, r, s')} \left[ (y_i - Q(s, a; \theta_i))^2 \right], \quad (15)$$

where the target is defined as:

$$y_i = r + \gamma \max_{a'} Q(s', a'; \theta_i^-), \quad (16)$$

and  $\theta_i^-$  are the parameters of a separate *target network* that is held fixed for several iterations and then periodically updated to match  $\theta_i$ .

To stabilize learning, DQN introduces two key techniques:

- **Experience replay:** Transitions  $(s_t, a_t, r_t, s_{t+1})$  are stored in a buffer and sampled uniformly to decorrelate training data and improve data efficiency.
- **Target network:** A separate target Q-network  $Q(s, a; \theta^-)$  is used to compute the TD target  $y_i$ , reducing the risk of divergence during updates.

DQN uses an  $\epsilon$ -greedy policy for exploration and is typically implemented with a convolutional Q-network that maps image-based states to Q-values over discrete actions.

DQN and its variants—despite being designed for discrete action spaces—have been widely applied in autonomous driving (notably in CARLA) by discretizing continuous controls such as steering and throttle. Notably, Zhang et al. [9] implemented a basic DQN with a convolutional neural network (CNN) for image processing. Yang et al. [10] significantly enhanced DQN performance using guided training (G-DDQN), state representation networks (GR-DDQN), dueling architectures (GRSD-DDQN), and safety rule constraints (GRS-DDQN). Bai et al. [56] proposed FEN-DQN, which

integrates a Feature Extraction Network composed of ResNet-50 and LSTM to generate affordance representations such as distance and angle. Muhammed et al. [57] introduced a Double 3-Step C51 DQN combining Categorical DQN with n-step returns and double Q-learning to reduce value overestimation. Elallid et al. [58] used a standard DQN to control an AV and avoid collisions.

Two studies by Chekroun et al. [59] and Toromanoff et al. [60] employed Rainbow-IQN Ape-X, a variant of DQN that incorporates expert demonstration integration, distributional RL, and prioritized experience replay in distributed settings. Safety-critical decision-making involving braking and driving modules was explored by Marouane et al. [61] using DQN. Other studies applied DQN to specific driving scenarios: Muhtadin et al. [62] focused on roundabouts, Elallid et al. [63] addressed complex traffic conditions, and May et al. [64] used CARLA-based simulations for self-driving control. Temporal dynamics were modeled using DQN with CNN-LSTM integration by Ahmed et al. [65], while Clemmons et al. [66] implemented a DQN-based controller with a non-visual state space and reward shaping to guide an ego vehicle through dynamic, multi-vehicle environments in CARLA. Li et al. [67] introduced three Q-learning variants—T-DQN, T-DDQN, and T-D3QN—for handling intersections. Finally, Liu et al. [68] employed Double Deep Q-Learning (DDQN) to scale RL to large state spaces in complex urban settings. Weng et al. [69] proposed a DQN-based car-following framework that uses RGB-D image inputs to guide control in autonomous driving. Hishmeh et al. [70] proposed a short-term planning framework for autonomous driving using the Dueling Deep Q-Network (Dueling DQN) algorithm. Dagdanov et al. [71] proposed DeFIX, a hybrid IL-RL framework where a DQN agent is trained on failure scenarios of an IL policy. The method uses a policy classifier to switch between the IL agent and specialized RL agents trained on mini-scenarios extracted from infractions, improving safety and performance in challenging CARLA urban tasks.

Deshpande et al. [72] propose a thresholded lexicographic Deep Q-learning architecture in which two independent Double DQN agents are trained separately for the safety and speed objectives. At execution time, their outputs are merged lexicographically: the safety agent’s action is accepted outright unless its estimated Q-value falls below a predefined threshold, in which case the speed agent’s recommendation is used. Maintaining distinct replay buffers and networks for the two objectives ensures that safety remains the higher-priority criterion without sacrificing the sample efficiency of off-policy Q-learning. An explainability-oriented study [73] augments a baseline DQN with Bayesian deep learning via Monte-Carlo Dropout. Multiple stochastic forward passes yield both the mean and standard deviation of each action’s Q-value; the agent then selects the action whose expected return is maximal and whose epistemic uncertainty lies below a user-defined threshold. This simple uncertainty filter provides a principled mechanism for trading off performance against risk, rendering the learned policy more transparent and trustworthy in safety-critical settings.

### 3.1.7 Soft Actor Critic

Soft Actor-Critic (SAC) [74] is a state-of-the-art, off-policy actor-critic algorithm that aims to maximize not only the expected cumulative reward but also the *entropy* of the

policy. By including an entropy term in the objective, SAC encourages more stochastic policies, which helps achieve better exploration in continuous control tasks and improves robustness during training. Unlike standard RL objectives that solely maximize expected reward, SAC incorporates an entropy term to encourage exploration. The overall objective for SAC can be written as:

$$J(\pi) = \mathbb{E}_{\substack{s_t \sim \rho_\pi \\ a_t \sim \pi(\cdot | s_t)}} \left[ \sum_{t=0}^{\infty} \gamma^t (r(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot | s_t))) \right], \quad (17)$$

where  $\rho_\pi$  denotes the state distribution under the policy  $\pi$ ,  $r(s_t, a_t)$  is the reward at time  $t$ ,  $\gamma$  is the discount factor, and  $\alpha$  is a temperature parameter that determines the trade-off between reward maximization and policy entropy. The term  $\mathcal{H}(\pi(\cdot | s_t))$  represents the entropy of the policy at state  $s_t$ .

SAC maintains two separate action-value (Q) networks,  $Q_{\theta_1}$  and  $Q_{\theta_2}$ , to mitigate overestimation bias. Both networks are trained to minimize a soft Bellman residual:

$$J_Q(\theta_i) = \mathbb{E}_{(s_t, a_t, r_t, s_{t+1}) \sim \mathcal{D}} \left[ \frac{1}{2} \left( Q_{\theta_i}(s_t, a_t) - y^{\text{soft}}(r_t, s_{t+1}, a_{t+1}) \right)^2 \right], \quad i = 1, 2 \quad (18)$$

where  $\mathcal{D}$  is a replay buffer of past transitions and

$$y^{\text{soft}}(r_t, s_{t+1}, a_{t+1}) = r_t + \gamma \left( \min_{j \in \{1, 2\}} Q_{\bar{\theta}_j}(s_{t+1}, a_{t+1}) - \alpha \log \pi_\phi(a_{t+1} | s_{t+1}) \right). \quad (19)$$

Here,  $\bar{\theta}_j$  indicates the parameters of a target network for the corresponding Q-function, and  $a_{t+1}$  is sampled from the current policy  $\pi_\phi(\cdot | s_{t+1})$ . The policy  $\pi_\phi(a_t | s_t)$  is typically parameterized as a Gaussian distribution whose mean and covariance are outputs of a neural network. The actor (policy) is updated by minimizing:

$$J_\pi(\phi) = \mathbb{E}_{s_t \sim \mathcal{D}} \left[ \mathbb{E}_{a_t \sim \pi_\phi(\cdot | s_t)} \left( \alpha \log \pi_\phi(a_t | s_t) - \min_{i \in \{1, 2\}} Q_{\theta_i}(s_t, a_t) \right) \right]. \quad (20)$$

Intuitively, this objective encourages the policy to select actions that maximize the action-value function, while also penalizing low-entropy policies through the  $\alpha \log \pi_\phi$  term. The temperature parameter  $\alpha$  balances the reward maximization term against the policy entropy term. Instead of fixing  $\alpha$ , SAC often employs an *automatic tuning* mechanism, which optimizes  $\alpha$  by solving:

$$J(\alpha) = \mathbb{E}_{a_t \sim \pi_\phi(\cdot | s_t)} \left[ -\alpha (\log \pi_\phi(a_t | s_t) + \bar{\mathcal{H}}) \right], \quad (21)$$

where  $\bar{\mathcal{H}}$  is a target entropy, typically chosen based on the action space dimension.

A wide range of studies have employed or extended SAC to enhance decision-making and control for autonomous driving. For instance, hierarchical extensions of SAC have been investigated, with DDPG variants included for baseline comparisons [11]. Wei et al. [12] introduced a continual RL framework based on SAC—referred to as EM-SAC—that aims to improve both adaptability and continual learning capabilities in autonomous driving decision and control systems.

Additionally, imitation learning methods, such as Behavior Cloning (BC) [75], GAIL [76], InfoGAIL [77], and MAGAIL [78], have been integrated with traffic simulators to train social vehicles. In this context, Zhu et al. [13] introduced the RITA traffic environment and demonstrated how an ego vehicle can be trained directly using SAC or GAIL, subsequently benchmarking these RL policies within the RITA-generated traffic environment. Other researchers have leveraged SAC specifically as a model-free, off-policy RL algorithm to address online decision-making tasks [79]. To boost data efficiency, Wu et al. [80] augmented SAC with Expert Demonstrations Augmentation (EDA) and Mixed Priority Sampling (MPS). Aghdasian et al. [81] further examined sensor-, image-, and fusion-based observation modalities under SAC-driven policies.

Beyond these approaches, hybrid or augmented architectures for SAC have also been proposed. For example, Huang et al. [82] developed Simoun, built on top of SAC and enriched with a dual-path network architecture (capturing motion and appearance) as well as a consistency-guided curiosity module. Focused Experience Replay (FER), introduced as a method to improve stability and learning speed, has been combined with SAC (principally SACv2) and shown to be adaptable to DDPG and TD3 [83]. Other lines of work further illustrate the versatility of SAC, whether integrated into custom reward functions [84] or applied to specialized driving scenarios such as autonomous emergency vehicle operation [85]. Han and Yilmaz [86] proposed a hybrid control strategy, SIRL, that combines a sparse expert policy with an SAC-based reinforcement learner using Bayesian Controller Fusion to enable safe exploration and accelerated training in urban driving scenarios. Igoe et al. [87] introduced Multi-Alpha Soft Actor-Critic (MAS), which treats the entropy coefficient as a learnable distribution to mitigate the stochastic bias inherent in SAC. Evaluated in the CARLA simulator, MAS demonstrated more stable and less overly cautious driving behaviors compared to standard SAC, while preserving sample efficiency and policy robustness. Gupta et al. [88] introduces NavSAC, a tailored Soft Actor-Critic agent within HyLEAR, encodes a  $400 \times 400$  intention image plus reward, speed, and prior action, then branches into  $V_\psi$ ,  $Q_\theta$ , and policy  $\pi_\phi$  over {Accelerate, Maintain, Decelerate}. Off-policy training with entropy-regularised SAC losses uses a replay buffer periodically infused with planner demonstrations, fusing imitation and exploration for faster, stable learning.

### 3.1.8 Twin Delayed Deep Deterministic Policy Gradient

Twin Delayed Deep Deterministic Policy Gradient (TD3) [89] is designed for continuous control and addresses the overestimation bias of Q-values often observed in

actor-critic methods like DDPG. TD3 tackles the Q-value overestimation that plagues DDPG by (1) training two independent critics and backing up the minimum of their estimates, (2) delaying actor updates so the policy sees more accurate values, and (3) adding clipped Gaussian noise to target actions to smooth the critics’ learning targets. These three tweaks combine to yield much stabler training and higher returns on standard continuous-control benchmarks

To mitigate overestimation, TD3 takes the *minimum* of the two target critics when computing the TD target. First, add small clipped noise to the target action:

$$\tilde{a} = \pi_{\phi'}(s') + \epsilon, \quad \epsilon \sim \text{clip}(\mathcal{N}(0, \sigma), -c, c). \quad (22)$$

The target value for a sampled transition  $(s, a, r, s')$  is then

$$y = r + \gamma \min(Q_{\theta'_1}(s', \tilde{a}), Q_{\theta'_2}(s', \tilde{a})). \quad (23)$$

Each critic network  $Q_{\theta_i}$  is updated by minimizing the mean-squared error between its predicted Q-value and the target  $y$ :

$$L(\theta_i) = \frac{1}{N} \sum_{(s, a, r, s') \sim \mathcal{D}} (Q_{\theta_i}(s, a) - y)^2, \quad (24)$$

where  $N$  is the minibatch size. Both critics are updated at every time step.

The actor network  $\pi_\phi$  is updated less frequently (e.g., once for every two updates of the critics). The gradient for the actor is computed by:

$$\nabla_\phi J \approx \frac{1}{N} \sum_{s \sim \mathcal{D}} \nabla_a Q_{\theta_1}(s, a) \Big|_{a=\pi_\phi(s)} \nabla_\phi \pi_\phi(s). \quad (25)$$

After each set of updates, the target networks are softly updated as:

$$\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i, \quad \phi' \leftarrow \tau \phi + (1 - \tau) \phi', \quad (26)$$

where  $\tau \ll 1$  is the target update rate (e.g.,  $\tau = 0.005$ ).

Several recent studies have enhanced and applied the TD3 algorithm to address safety, exploration, and performance in autonomous driving tasks. Jia et al. [90] proposed an improved TD3 variant that enforces safety constraints via a Lagrangian formulation, ensuring that the control policy respects prescribed safety limits. In another work, Xu et al. [91] introduced Pri-TD3, which builds upon the TD3 framework with Prioritized Experience Replay (PER) to mitigate Q-value overestimation and reduce policy variance more effectively than simpler methods such as DDPG.

Liao et al. [92] employed TD3 for real-time, safe, and comfortable lateral control, focusing on collision avoidance strategies under realistic driving conditions. Furthermore, other researchers have explored noise injection into the actor outputs during training to promote more robust exploration and improve policy convergence [93]. In the domain of intersection navigation, TD3 has been combined with a Local Goal Velocity (LGV) mechanism and further augmented with PER to refine the learning process and expedite the development of effective driving policies [94]. Li et al. [95] extend TD3 with a local-attention safety architecture in which the actor network is

augmented by an ego-attention module. This module dynamically identifies socially critical surrounding vehicles, enabling the policy to focus its continuous control outputs on the most influential interactions. To further promote collision avoidance, the authors design a custom safety-reward that penalises proximity violations while still rewarding efficient progress. The resulting algorithm retains TD3’s twin critic networks and delayed policy updates, yet the attention mechanism and safety-aware shaping together yield a policy that is both risk-averse and time-efficient in dense traffic scenarios.

### 3.1.9 New Methods

There are certain new RL-based methods introduced in the literature. For instance, khalil et al. [96] leverages SAC’s entropy-regularized objective to encourage exploration. SAC trains a policy (actor) and value function (critic) jointly in an off-policy manner by sampling from a replay buffer and updating parameters using Q-value targets .

In a different vein, Qi et al. [97] introduces a new framework merges MFRL with human cognitive processes to enhance decision-making. Specifically, the framework replaces conventional neural network decision-making modules with a cognitive model that follows the “Perception–Memory–Reasoning–Action” cycle, thereby infusing human-like reasoning loops into the learning process. To ensure transparency and interpretability, predefined production rules govern the decision-making steps of this cognitive model .

Couto et al. [98] proposed hGAIL, a hierarchical imitation learning framework that integrates Generative Adversarial Imitation Learning (GAIL) with Proximal Policy Optimization (PPO) for policy training. This model-free, on-policy approach enables more structured policy learning by combining adversarial imitation with hierarchical decision-making.

Jaafra et al. [99] proposed a context-aware autonomous driving framework that leverages meta-RL for rapid policy adaptation in dynamic urban environments. Their method combines a gradient-based meta-learning approach—using a Model-Agnostic Meta-Learning (MAML) strategy—to train a meta-policy capable of quickly adapting to new driving scenarios with minimal additional data.

Baheri et al [100] depart from value- and policy-gradient paradigms by adopting an evolutionary RL approach based on the Covariance Matrix Adaptation Evolution Strategy. CMA–ES maintains a population of candidate controllers whose parameter vectors are sampled from a multivariate normal distribution; after roll-outs, the distribution’s mean and covariance are iteratively updated to favour higher-performing individuals. By foregoing gradient information altogether, this black-box optimisation scheme sidesteps issues of non-differentiability and sparse rewards, enabling direct search in the controller’s parameter space while gradually adapting exploration to the covariance structure of successful policies.

Moreover, distributed off-policy RL has been employed in several recent approaches. These methods incorporate elements such as clipped double Q-learning from TD3, as well as the maximum entropy concept from SAC, to refine learning stability and exploration [101].



### 3.1.10 Comparative Studies

Some studies consider different methods in RL and conducted comparative studies. For instance, a variety of test scenarios have been constructed in recent studies to evaluate learned models rigorously and to provide standardized benchmark results under well-defined evaluation metrics [102]. One such learning framework partitions the overall driving task into distinct subtasks—lateral control (lane keeping and lane changing), longitudinal control, and decision-making—to address the multi-faceted nature of autonomous driving. Continuous-action subtasks (lane keeping, lane changing, and longitudinal control) are trained using an Actor-Critic approach (specifically, DDPG), while discrete decision-making tasks employ a DQN. To balance exploration and exploitation, multi-head actor networks introduce adaptive noise in the action selection mechanism [103].

Comparative evaluations of D3QN, A2C, and PPO and their variants have also been reported in intersection scenarios, where decision-making complexities often increase [104]. In another line of research, a child’s ride-on toy car was redesigned to serve as a self-driving platform; this setup was employed to compare the performance of three RL models (DQN, DDQN, and D3DQN) under real-world constraints [105]. Further expanding on efficiency concerns, Frauenknecht et al. [106] investigated data-efficient RL approaches for vehicle control, highlighting the importance of reducing sample complexity and computational overhead without sacrificing policy performance. They compare SAC, PETS-MPPI, MBPO, REDQ for their work.

## 3.2 Model-Based Reinforcement Learning

Model-based RL (MBRL) uses a learned or predefined model of the environment to guide policy optimization. Instead of relying exclusively on repeated interaction with the real environment, the agent maintains an approximate transition and reward function, often denoted by

$$\hat{f}_\theta(s, a) \approx s', \quad \hat{r}_\theta(s, a) \approx r, \quad (27)$$

where  $\theta$  are the model parameters,  $s \in \mathcal{S}$  is the current state,  $a \in \mathcal{A}$  is the action,  $r$  is the immediate reward, and  $s'$  is the next state. These approximations can be learned by minimizing prediction errors across a dataset of real transitions collected from environment interactions:

$$\mathcal{L}(\theta) = \sum_{(s, a, s', r) \in \mathcal{D}} \left\{ \|\hat{f}_\theta(s, a) - s'\|^2 + \|\hat{r}_\theta(s, a) - r\|^2 \right\}. \quad (28)$$

Once the model is learned, it can be used to *simulate* experience without requiring costly real-world rollouts. The agent can perform imaginary rollouts by repeatedly applying the learned transition and reward functions:

$$s_{t+1}^{(sim)} = \hat{f}_\theta(s_t^{(sim)}, a_t^{(sim)}), \quad r_{t+1}^{(sim)} = \hat{r}_\theta(s_t^{(sim)}, a_t^{(sim)}), \quad (29)$$

These simulated transitions are then used to evaluate or improve the policy. For instance, if the agent uses a parameterized policy  $\pi_\phi(a \mid s)$ , it can be improved by maximizing expected returns in simulation. In a policy-gradient framework, this might involve updating  $\phi$  via

$$\nabla_\phi J \approx \frac{1}{N} \sum \nabla_\phi \log \pi_\phi(a_t \mid s_t) G_t, \quad (30)$$

where  $G_t$  is an estimate of the return from simulated or real data. A popular strategy is to alternate between collecting real data and using the learned model for planning or imagination-based updates. This idea, originally introduced in the *Dyna* framework [107], typically follows these steps: gather transitions  $(s, a, s', r)$  from the real environment, update the model  $\hat{f}_\theta, \hat{r}_\theta$  using this data, and then generate additional “imagined” transitions by rolling out the updated model to refine the policy or value function.

Although MBRL can dramatically improve sample efficiency, it also introduces potential model bias: if the learned model is inaccurate in certain state-action regions, planning on this incorrect model can lead to suboptimal or destabilizing actions. Addressing this challenge often involves estimating uncertainty (e.g., with ensemble models) or restricting simulation rollouts to shorter horizons where the model remains more reliable. In practice, successful methods like PETS (Probabilistic Ensembles with Trajectory Sampling) [108], MBPO (Model-Based Policy Optimization) [109], or Dreamer [110] combine these ideas to balance the benefits of high sample efficiency with robust policy learning.

MBRL techniques have recently garnered increasing attention for autonomous driving applications, as they aim to leverage internal world models to improve sample efficiency and reduce the need for risky real-world exploration. For instance, DreamerV3 [111] has been utilized to learn a latent representation of the driving environment, with the policy trained within this learned latent space for enhanced efficiency [14]. Expanding on the Dreamer paradigm, an approach called Iso-Dream++ decouples controllable and noncontrollable dynamics in the latent space, incorporating min-max variance constraints to avert training collapse and a sparse dependency mechanism to capture indirect influences between these decoupled states. This architecture generates latent imaginations from isolated controllable and noncontrollable factors, which guide the policy learning process [15].

Deep Imitative Reinforcement Learning (DIRL) represents another variant that combines imitation learning (IL) with a MBRL framework, using a predictive world model—Reveries-net—to simulate intricate interactions and refine driving policies in an offline fashion [16]. In a further illustration of the benefits of MBRL, a “world-on-rails” methodology leverages pre-recorded trajectories and treats the environment as non-responsive to ego-vehicle actions. This simplifies the modeling process, enabling the use of a forward model to predict vehicle dynamics while action-value functions are computed via tabular Bellman updates. The resulting action-value functions then supervise a visuomotor policy through distillation, effectively sidestepping the need for direct exploration or the risk of executing unsafe policies [17]. By integrating world modeling with dense offline supervision, such approaches demonstrate that MBRL can be both sample efficient and robust for autonomous driving tasks.

### 3.3 Hybrid RL (Combination of MBRL and MFRL)

Recent advancements have explored synergistic integrations of MBRL and MFRL to enhance both safety and performance in autonomous driving systems. One approach leverages a model-free constraint—captured via distributional RL in the form of a cost critic network to quantify risk—alongside a model-based High-Order Control Barrier Function (HOCBF) to enforce vehicle-dynamics safety constraints. The proposed framework employs Lagrangian optimization to update the policy under these constraints and integrates imitation learning during early training phases for improved stability [18].

Another example, termed Deductive Reinforcement Learning (DeRL), merges DDPG with a model-based “Deduction Reasoner” (DR). The DR utilizes a learned environment model to project future states and rewards, providing a “self-assessment value” that guides policy improvements. Additionally, a Semantic Encoder Module (SEM) extracts low-dimensional, robust representations, thereby reinforcing both interpretability and efficacy [19].

In parallel, research has also combined a DQN agent with the classical, model-based A\* path planner. Here, the A\* algorithm supplies waypoints, and the DRL agent learns to follow these waypoints while actively avoiding collisions. This hybrid solution underscores the importance of integrating established planning algorithms with learning-based modules to bolster reliability and safety [20].

### 3.4 Alternative RL Methods

While methods like DQN, PPO, and SAC are frequently used for autonomous driving in CARLA, a number of less-explored but promising RL approaches have also been applied. These methods offer alternative perspectives on decision-making, uncertainty handling, and policy structure. Below, we provide a brief overview of these methods, citing their foundational work as well as their application in the CARLA simulator.

#### 3.4.1 Bayesian Reinforcement Learning

Bayesian RL (BRL) maintains a posterior distribution over environment dynamics or value functions, allowing for uncertainty-aware exploration. A comprehensive survey of BRL techniques was presented by Ghavamzadeh et al. [112]. In the CARLA context, BRL was used by Gharaee et al. [113], where Gaussian Mixture Models (GMMs) were employed to cluster perceptual states and model action probabilities. This approach combined GMM-based modeling with temporal difference learning, enabling adaptive policy updates based on observation similarity.

#### 3.4.2 Hierarchical Reinforcement Learning

Hierarchical Reinforcement Learning (HRL) decomposes complex tasks into hierarchies of subpolicies. The options framework introduced by Sutton et al. [114] laid the groundwork for HRL by enabling temporal abstractions in RL. In the work by Li et al. [115], the authors introduced a hierarchical framework that integrates Imitation Learning (IL) and RL. IL is used to manage low-level control tasks, such as trajectory

following, while RL is responsible for high-level decision-making, including navigation and collision avoidance. The framework first learns from expert demonstrations and then refines its policy through RL, enabling more robust behavior in complex and dynamic driving environments. Separately, Wang et al. [116] adopted a two-level HRL structure designed for collision avoidance involving dynamic pedestrians. The upper layer selects from a set of Decision Primitives (DPs) based on a Predictive Risk Map (PRM), while the lower layer executes the chosen primitive through detailed control over steering, acceleration, and braking. This design allowed the agent to make safe and context-sensitive decisions in highly dynamic environments. He et al. [117] proposed a HRL framework for autonomous driving, where the upper decision-making layer selects motion primitives using a Double DQN, and the lower execution layer applies DDPG to realize continuous control. This architecture allows for skill reuse across scenarios and improves generalization by decomposing high-level decisions from low-level controls. Ben Naveed et al. [118] proposed a HRL framework for trajectory planning in autonomous driving, combining a high-level policy for maneuver selection (lane follow/wait vs. lane change) with low-level planners for trajectory generation. A PID controller was used to track waypoints instead of directly applying low-level controls, ensuring smooth and safe execution. Gangopadhyay et al. [119] introduced a Hierarchical Program-Triggered RL (HPRL) framework that combines symbolic programming with modular RL agents for autonomous driving. Each agent is trained for a specific maneuver (e.g., lane change, turn, straight drive) using DQN or DDPG, while a structured program governs task execution and enforces safety via embedded assertions. This design improves interpretability and verifiability by isolating learning to simple sub-tasks and validating the high-level controller with formal methods.

### 3.4.3 Maximum Entropy Reinforcement Learning

Maximum Entropy Reinforcement Learning (MaxEnt RL) encourages policies that maximize both reward and entropy, promoting exploration and robustness. The principle was originally introduced in the context of inverse RL by Ziebart [120], and later extended to RL through energy-based policies by Haarnoja et al. [121]. Khalil et al. [122] proposed a Sequential Latent MaxEnt RL framework for use in CARLA. This approach combined latent-space RL with motion prediction via a probabilistic graphical model, enhancing decision-making in dynamic driving scenarios.

### 3.4.4 Distributional Reinforcement Learning.

Distributional RL learns the full distribution of returns  $Z(s, a)$  rather than only their expectation  $Q(s, a)$ . This approach was first introduced by Bellemare et al. [123] and further refined through the Fully Parameterized Quantile Function (FQF) by Yang et al. [124]. Chen et al. [125] applied FQF in CARLA to learn separate policies for path and speed planning, allowing the model to better handle uncertainty due to occlusions and unpredictable pedestrian behavior.

**Table 1** Key CARLA publications by RL method.

Category / Method	Key CARLA Papers / Variants
<b>MFRL</b>	
Actor-Critic	[25, 26]
DQN	[9, 10, 56–73]
DQL	[39, 128]
PPO	[7, 8, 41–54]
TRPO	[55]
DDPG	[4–6, 30–38]
SAC	[11–13, 79–88]
TD3	[90–95]
<b>MBRL</b>	[14–17]
<b>Hybrid RL</b>	
Model-Based + Model-Free	[18–20]
<b>Alternative / Specialized RL Methods</b>	
Bayesian RL	[113]
Hierarchical RL	[115–119]
Distributional RL	[125]
Maximum Entropy RL	[122]
Skill-Based Hierarchical Offline RL	[127]
New / Evolutionary Methods	[96–101]
<b>Comparative Studies</b>	[102–106]

### 3.4.5 Skill-Based Hierarchical Offline Reinforcement Learning

Skill-based RL allows high-level policies to select from learned low-level behaviors. In offline settings, this approach has been explored by Sharma et al. [126], who proposed latent skill discovery in multi-task settings. Li et al. [127] adapted this concept to CARLA by training a high-level planner using a two-branch Variational Autoencoder (VAE) that captured both discrete and continuous driving skills, while the low-level controller executed steering and throttle commands accordingly.

## 4 State Space Representations

In RL for autonomous driving, the design of the observation space plays a critical role in shaping the agent’s perception and decision-making capabilities. CARLA provides a rich set of sensor modalities—such as RGB cameras, LiDAR, semantic segmentation, and direct vehicle state information—enabling diverse representations of the environment. Researchers often choose or combine these modalities based on task complexity, computational constraints, and the desired level of abstraction. This section surveys commonly used observation spaces in CARLA-based RL studies, grouped by their sensor types and fusion strategies, and discusses how each representation affects policy learning and generalization.

## 4.1 Front Facing RGB Images

Several studies adopt a single-stream visual input as the core observation. For instance, Elallid et al. [58] use a front-facing grayscale image, initially sized at  $640 \times 480$  but cropped and resized to  $192 \times 256$  to reduce irrelevant background. This preprocessed image is then flattened to feed into a Deep Q-Network (DQN). Zhang et al. [9] also rely heavily on pixel observations, with lane-deviation and collision indicators influencing only the reward function rather than the agent’s direct observation space. Liu et al. [68], uses end-to-end vision-based learning that processes standard front-facing RGB images with lane invasion and collision detector sensor. Meanwhile, Ahmad et al. [65] fuses an RGB image ( $400 \times 400 \times 3$ ) with vehicle speed and the angle from the road center, thereby adding minimal numeric inputs to anchor visual perception. Although purely vision-based designs can suffice for lane-following tasks, many researchers find that incorporating at least some numeric data—such as speed or steering angle—improves training stability and convergence rates. Xing et al. [48] uses a raw RGB front facing image for Latent Unified State Representation (LUSR) that separates domain-general features (e.g., vehicle dynamics) from domain-specific variations like weather, thus making the agent’s policy more transferrable. In a model-based setting, Pan et al. [15] decouples the state into controllable (ego-vehicle) and noncontrollable (other vehicles) factors from raw RGB front facing camera frame, enhancing planning accuracy.

Other works emphasize the incorporation of kinematic or command data to provide richer contextual cues. Carton et al. [7] employ a front-facing RGB camera at around  $192 \times 64$  resolution and augment it with vehicle speed and steering angle. Mohammed et al. [8] propose fusing a sequence of front-facing RGB frames, extracted via a pre-trained CNN, with current speed and route speed limits, either in early convolutional layers or at the final fully connected layers. By integrating control signals into the state, Wu et al. [44] provide the agent with steering, throttle, and brake values that are concatenated with a 64-dimensional feature embedding from a front-facing RGB feed ( $160 \times 80 \times 3$ ). Similarly, Zhao et al. [43] compress RGB frame using a Variational Autoencoder (VAE) and combine these latent features with vehicle speed, throttle, orientation, distance from lane center, waypoint and previous steering angle, yielding a compact yet information-rich input. Peng et al. [30] go a step further by including a Pure Pursuit (PP) steering angle in addition to speed and high-level route commands, demonstrating how geometric planning cues can guide visual learning to achieve smoother control policies.

Dimensionality reduction strategies frequently appear to manage the high computational demands of processing raw images. Savari et al. [79] employ an auto-encoder to generate a latent representation from the front camera view before feeding it into the Soft Actor-Critic (SAC) algorithm. Deductive RL [19] takes an alternative approach by giving the actor a reduced observation (monocular front-facing camera + forward speed), while the critic receives a privileged state from CARLA (e.g., distances to pedestrians or vehicles). This allows the actor to learn from a realistic sensor setup, but the critic’s extra information stabilizes training. Meanwhile, Zhao et al. [129] leverage

a Bird’s Eye View (BEV) semantic map only during training as an auxiliary supervision signal, relying solely on the front-facing ResNet50-processed RGB images plus high-level navigational commands at inference, thus ensuring real-time viability.

Temporal stacking of frames is also commonly leveraged for better motion awareness. Huang et al. [82] Stacks a sequence of recent frames, raw image frames processed by a dual-path CNN: one path learns appearance features per frame, the other learns motion features from frame-to-frame changes. Elallid et al. [93] provides the agent with four consecutive grayscale images ( $84 \times 84$ ) to capture short-term environment changes critical for intersection decisions. Elallid et al. [63] similarly stacks four RGB frames (downsampled to  $144 \times 144$  grayscale) and appends a 2D goal vector representing the target location. In racing contexts, Cai et al. [16] uses four-frame stacks of  $96 \times 96 \times 3$  inputs plus vehicle speed, exploiting temporal cues to execute rapid maneuvers. Li et al. [67] processes two consecutive monocular images (captured at times  $t$  and  $t - 2$ ) to estimate ego-vehicle motion in intersections. In a somewhat different direction, Chen et al. [33] encodes the front camera ( $64 \times 64 \times 3$ ) with a self-attention network and fuses additional sensor data (speed, acceleration, position), shifting focus toward efficient multi-modal integration. Finally, Li et al. [115] demonstrates how feature extraction with ResNet, combined with vehicle speed and high-level navigation commands (e.g., Turn Right, Turn Left, Go Straight), allows for hierarchical policy learning that separates strategic decision-making from direct motor control. Meanwhile, Youssef et al. [6] also utilize front-facing images for feature extraction but place special emphasis on high-level route instructions (turn left, turn right, or go straight), integrated through fully connected layers or the final actor network.

## 4.2 Bird Eye View Images

A growing number of studies have explored bird’s-eye view (BEV) representations to provide a structured, top-down depiction of the driving environment in CARLA. Such views typically display the ego vehicle’s surroundings—including lanes, vehicles, and traffic signals—in a spatially coherent format, facilitating more interpretable learning of road geometry and multi-agent interactions. For instance, Wang et al. [102] generate BEV images at resolutions of  $64 \times 64$ , placing the ego vehicle at the bottom center of the frame, with color-coding for ego-vehicle (blue), lanes (red) and surrounding vehicles (green). This strategy provides the RL agent with an easily parsable scene layout, improving its capacity for lane-changing and collision avoidance. Similarly, Deng et al. [41] enhance decision-making by integrating BEV inputs with a 6-dimensional measurement vector (speed, acceleration, steering angle, gear state) and a historical trajectory encoded via Gated Recurrent Units (GRU). This approach leverages short-term motion context to address partial observability issues, thereby promoting more consistent control in dynamic traffic scenarios.

In many BEV-based methods, the representation is augmented with additional high-level or kinematic data to capture the ego vehicle’s pose and local objectives. Deng et al. [42] build upon their earlier work by blending BEV images with real-time vehicle measurements and historical trajectories. They employ recurrent architectures (GRU, LSTM, MLP, or Transformers) to encode past observations, actions, and rewards,



consolidating them into a latent vector that is then fed to the RL policy. Historical data improves situational awareness in partially observable environments, leading to smoother handling of dynamic road users. LiDAR sensors can also reinforce the BEV perspective by supplying real-time object detection and depth cues. For example, Manikandan et al. [105] combine LiDAR-based obstacle distances with camera perception and kinematic signals (velocity, steering angle, lane offset) to better detect and localize objects in all directions. Fused RGB-LiDAR point clouds or RGB-D images likewise help the agent maintain a detailed map of surrounding vehicles, enabling more informed decisions under complex traffic conditions.

Several other works implement specialized BEV semantic segmentation or multi-channel images to highlight distinct road features. In Li et al. [14], the state space includes semantic masks, bounding boxes, and traffic signal states, while Zhang et al. [45] encodes drivable areas, lane boundaries, and pedestrians in separate segmentation channels. Extended approach, such as Agarwal et al. [46] combine BEV segmentation (vehicles, pedestrians, traffic lights) with waypoints from a global planner, along with kinematic data (speed, distance to goal, steering angle) to produce a richly detailed representation. Trumpp et al. [47] uses four stacked bev images into LSTM-based modules to capture temporal dependencies, further refining predictions and behavior modeling.

To handle highly dynamic scenarios, some methods integrate motion prediction or incorporate specialized LiDAR-based maps into the BEV. Khalil et al. [122] uses two BEV images: one for LiDAR data (ground and above-ground point clouds with route waypoints) and another for predicted trajectories of surrounding vehicles extracted from MotionNet [130]. Likewise, Li et al. [127] retains a BEV semantic segmentation image plus a measurement vector, while Tong et al. [85] adopts a  $256 \times 256$  BEV image covering a 35m radius around the ego vehicle. This latter work encodes an ego-vehicle state (speed, angular gap to route, lateral distance, and junction indicator), ensuring that emergency-vehicle maneuvers are sensitive to both local geometry and mission-critical efficiency. Altogether, these BEV-centric frameworks demonstrate the benefits of structured, top-down map representations, especially when complemented by numeric vehicle states, temporal models, or motion prediction modules. They enable more interpretable and robust policy learning in multi-agent and urban driving environments, where spatial awareness and accurate trajectory planning are essential.

### 4.3 Semantic Segmentic Images

Semantic segmentation has been widely adopted as a critical component of state representations in RL-based autonomous driving frameworks, often to encode meaningful spatial and categorical cues for decision-making. For instance, Wang et al. [101] integrates drivable area segmentation alongside vehicle speed, providing the RL agent with a clear delineation of navigable surfaces. In [59] Chekroun et al. used three RGB camera streams (center, left, right) are processed by pretrained encoders to extract segmentation features, traffic light status, and intersection presence over short temporal windows (stacking the frames), enabling the agent to capture both semantic and temporal information. Likewise, the pipeline in [60] Toromanoff et al. employs



a ResNet-18 pretrained on semantic segmentation data to encode visual frames and includes auxiliary signals such as traffic light states, intersection recognition, and lane positioning. These aggregated features are then stacked over time, furnishing the RL module with a richer temporal perspective.

In a more sensor-fused approach, the work in [61] Marouane et al. combines minimum distance to obstacles (obtained from depth and segmentation images), collision indicators, orientation angles, lateral distance metrics, and vehicle velocity, encapsulating both semantic and geometric factors critical to safe navigation. Going further, Gharaee et al. [113] adapted a bayesian perspective, where segmentation maps are segmented into six spatial regions (three vertical, two horizontal), and histograms of class distributions across five semantic categories— road, road-line, off-road, static objects, dynamic objects —are normalized and concatenated to form a compact yet informative state vector.

Other works concentrate on tailoring segmentation outputs to specific driving tasks. In [62] Muhtadin et al. used grayscale segmentation images (or advanced binary masks) filter out non-drivable terrain to simplify roundabout navigation policies. Finally, Silva et al. [50] processes front-facing semantic segmentation frames at  $84 \times 168$  pixels and stacks four consecutive frames to account for temporal context. These semantic inputs are merged with vehicle measurements—such as throttle, steering angle, speed, and heading—to form a unified representation that is passed through a CNN-based encoder (inspired by IMPALA), enabling robust lane-keeping and intersection handling. Collectively, these works underscore the versatility of semantic segmentation in RL-based driving, as it provides structured, high-level features that enhance situational awareness and facilitate safer and more efficient decision-making.

#### 4.4 Hybrid State Space

Several recent studies adopt *hybrid* state representations that fuse high-level perception and low-level vehicle data to tackle the complexity of autonomous driving tasks. For instance, Wang et al. [11] propose a hierarchical approach that splits decision-making between a high-level agent, which selects maneuvers (e.g., lane following, turns), and a low-level agent controlling steering, speed, and throttle. Their observation space combines semantic segmentation (ENet) for drivable areas, object detection (YOLOv4) for vehicles and pedestrians, and vehicle kinematics, leading to more informed decisions in urban settings. Khalil et al. [96] also employ a fused representation, relying on LiCaNext outputs to project LiDAR-based object detection and motion fields onto a BEV. By including temporal history and previous control inputs, they maintain critical environmental context while reducing the dimensionality of raw sensor data. Bai et al. [56] focus on a Feature Extraction Network (FEN) that processes camera images into semantic affordances—such as center distance and hazard indicators—while also including the agent’s speed, traffic light state, and leading vehicle parameters, helping a DQN handle the complexity of dynamic traffic. Muhammed et al. [57] merge depth and semantic segmentation images to provide pixel-wise distance-to-objects and class labels, processed by a CNN for enhanced spatial understanding and obstacle avoidance. Aghdasian et al. [81] similarly combine high-level visual features from a residual CNN with low-level kinematic data (e.g., velocity, orientation,

lane-center distance), creating a fused latent space that balances scene perception and precise motion control. In a different vein, the work P´erez-Gil et al. [31] integrates either raw RGB images or waypoints with a driving feature vector, which includes vehicle speed and lane center information, highlighting the importance of combining visual cues with essential geometric features.

Additional hybrid strategies leverage latent representations or domain adaptation to ensure more robust policy learning. For example, Chen et al. [17] fuses multi-camera (three) inputs to produce high-level navigation commands, Anzalone et al. [49] stitches multi-camera RGB images over time and combines them with road, vehicle, and navigational features for more efficient training in dynamic environments. Transformer-based multimodal fusion is seen in [32], where Fu et al. combined camera images, LiDAR bird’s-eye pseudo-images, and ego-vehicle dynamics to enable attention-based decision-making. Similarly, May et al. [64] merges data from three RGB cameras, radar, GPS, IMU, and collision sensors into a single representation. The integration of model-based path planners with deep RL is explored by Yurtsever et al. [20], combining CNN-encoded camera data with vehicle states and waypoint distances. Finally, Jo et al. [84] processes camera and LiDAR bird’s-eye views into a low-dimensional latent space alongside vehicle kinematics, underscoring how compressed multimodal features can guide effective policy learning. Collectively, these works illustrate that hybrid state representations—spanning raw imagery, semantic cues, object detection outputs, vehicle kinematics, and driver or route-level attributes—consistently improve situational awareness and decision-making performance in RL-based autonomous driving systems.

#### 4.5 CARLA Kinematics only State Space

Below is a consolidated description of a typical *CARLA kinematic-based observation* for autonomous driving tasks in RL, integrating approaches from multiple recent works. Although individual studies vary in the precise composition and dimensionality of the observation vector, the following elements commonly appear.

##### 1. Ego-Vehicle Kinematic State:

- **Position and Orientation:** A core component is the ego vehicle’s global or local-frame position (often denoted  $x, y$ ) and heading angle (yaw). Many works encode yaw as  $\cos(\text{yaw})$  and  $\sin(\text{yaw})$  to avoid discontinuities [10, 94]. Some papers also track roll, pitch, or higher-order motion states (e.g., yaw rate, pitch rate) to capture full 3D pose dynamics [106]. Others include an explicit *heading error* relative to a desired path or lane center instead of absolute yaw [5, 97].
- **Velocity and Acceleration:** The longitudinal speed  $v$  (m/s) is almost always present. Many studies also include lateral velocity  $v_y$ , longitudinal acceleration  $a_x$ , or even higher derivatives like  $\dot{x}, \dot{y}, \ddot{x}, \ddot{y}$  in local coordinates to capture the ego vehicle’s instantaneous motion [12, 90, 106, 125]. Yaw rate  $\dot{\gamma}$  can similarly be added [4, 80].
- **Lateral/Longitudinal Deviations:** Many lane-keeping or lane-changing tasks explicitly track the ego vehicle’s lateral offset from a reference lane boundary [18, 92], cross-track error [5], or distance to a goal position [12, 80].

## 2. Surrounding Traffic and Obstacles

- **Relative Positions and Velocities:** To handle multi-vehicle traffic, numerous works include distances or relative coordinates of up to two or more neighboring vehicles [4, 10, 90]. Relative velocity components ( $\Delta v_x, \Delta v_y$ ) are frequently used [90], and some studies extend this to angle or heading differences [92].
- **Distance to Obstacles or Road Boundaries:** In maneuvers such as lane changing or parking, the agent may need direct measures of distance to lane edges and obstacles on each side [80, 92]. In unsignalized intersection or pedestrian-avoidance tasks, explicit collision-risk terms or booleans (e.g., “obstacle ahead”) are added [97, 116].
- **Use of Graph or Waypoint Structures:** Some recent research vectorizes map information (lane segments, HD map nodes) and adjacency relationships (predecessor/successor/left/right) to better capture road topology [18, 94]. Others represent upcoming waypoints directly as coordinates in the ego frame [131].

## 3. Additional Task-Specific Features

- **Lane Indicators and Discrete Events:** Discrete lane IDs, turn signal statuses, or lane-change flags are occasionally encoded to help the network differentiate maneuver intentions [10, 90].
- **Driver-Centric Metrics:** In driver-assistance contexts, some works include “driver stress” or “excitement” measures as direct numeric inputs, reflecting a more human-centric policy [25].
- **Road Speed Limits and Traffic Signals:** For tasks involving traffic rules, the current speed limit or traffic signal states (green/red) may appear in the observation vector [25, 83].
- **Minimal vs. Extended Observation Spaces:** Simpler methods might only use  $(x, y, \text{yaw}, v)$  [66, 97], while more advanced approaches incorporate up to 20–30 features or even on-board sensor readings (lidar/radar arrays) processed into bounding distances [80, 92].

## 4. Typical Dimensionality and Implementation Details

- **Observation Dimension:** Ranges from as few as 4–6 scalars (for minimal steering or speed control tasks) [5, 97] up to 20+ dimensional vectors when including multiple vehicles, and lane offsets. Some wrappers for CARLA yield 6–10 continuous values by default [83].
- **Local vs. Global Coordinates:** Many prefer an *ego-centric* frame, ensuring the ego vehicle is always at the origin with heading zero [94, 125]. This can simplify learning because the agent’s network does not need to adapt to arbitrary global positions.
- **Stacking Observations:** Combining consecutive frames (e.g., last four states) is widely adopted to capture short-term dynamics and reduce partial observability [103, 104, 125].

In sum, CARLA kinematic-based observations typically revolve around the ego vehicle’s position, orientation, and velocities, often augmented by lane-relative or

goal-relative errors, and surrounding-vehicle states. These standardized kinematic features provide a strong foundation for policy learning in lane keeping, lane changing, intersection negotiation, and various maneuvering scenarios in urban or highway settings [4, 5, 10, 12, 18, 25, 66, 80, 83, 90, 92, 94, 97, 103, 104, 116, 125, 131]. Depending on the task, this base observation can be expanded with additional environment, driver, or sensor data, creating a flexible framework for RL in autonomous driving.

## 5 Action Space

The choice of action space is a fundamental design decision in RL for autonomous driving, as it directly shapes the agent’s ability to explore, learn, and execute precise maneuvers. In CARLA-based driving tasks, action spaces have been defined in discrete, and continuous forms to balance trade-offs between sample efficiency, control fidelity, and interpretability. We considered a new category for hierarchical studies as their action space included two different levels. In this section, we first review discrete formulations—ranging from coarse, high-level primitives to fine-grained binning strategies—before examining continuous designs that offer smooth control via normalized or custom scaled outputs. We then discuss hierarchical studies and how they considered their action space.

### 5.1 Discrete Action Space

Action spaces in RL for autonomous driving using CARLA can be broadly categorized into discrete, continuous, or hybrid. In autonomous driving applications using RL, the design of the discrete action space plays a critical role in balancing learning complexity and control precision. Across the literature, the discrete action space can be categorized along several dimensions. These categories, along with representative examples, are summarized in Table 2.

1. Control Dimension:

- Lateral Control: Decisions are focused on lateral maneuvers (e.g., lane changes or steering adjustments).
- Longitudinal Control: The focus is on speed-related decisions (e.g., throttle and braking).
- Combined Control: Both lateral and longitudinal controls are considered together.

2. Granularity:

- Coarse: Only a few discrete actions (often fewer than 10) are defined. These actions are mostly high-level actions like accelerate, decelerate, turn left, turn right, that are mapped to small combination of the steering and acceleration.
- Moderate: An intermediate number of actions (typically in the range of 10–30).
- Fine: A detailed discretization that divides the control range into many bins (e.g., 27, 108, 190 actions).

3. Mapping Strategy\*:

- **Direct Mapping:** Discrete actions are mapped directly to specific low-level control commands (e.g., steering and throttle values), which are applied to the environment without additional processing.
- **Indirect Mapping:** Discrete actions are translated into control signals by an external controller or planning module (e.g., a PID).

This categorization illustrates the diversity of design choices in discrete action space formulations for autonomous driving. A majority of work uses flat action space with direct mapping with different granularity. Some like [54, 61, 63, 128] use coarse action space to simplify the decision-making, while others [58–60] use fine-grained discretizations (up to 190 actions) to achieve better control precision. A small portion of studies incorporate indirect mapping, where the action chosen by the RL agent is given to a PID controller for execution [10, 71]. Finally, hierarchical approaches—such as [115, 118, 119] employ a multi-layered decision process, where a high-level policy selects abstract maneuvers (e.g., "turn left") that are refined into low-level controls. These strategies help address the complexity of continuous control while offering modularity and interpretability.

**Table 2** Taxonomy of Discrete Action Space Designs in RL-Based Autonomous Driving

Control Dim.	Granularity	Mapping Strategy	References
Lateral only	Coarse-Moderate	Direct	[68, 102]
Longitudinal only	Coarse-Fine	Direct	[39, 55–57, 67, 72]
Combined	Coarse	Direct	[54, 61–64, 66, 69, 70, 73, 88, 105, 113, 116, 128]
Combined	Fine-Moderate	Direct	[9, 14, 17, 20, 31, 58–60, 65]
Combined	Coarse-Moderate	Indirect (PID)	[10, 25, 71]

## 5.2 Continuous Action Space

In RL for autonomous driving, continuous action spaces are widely used to provide fine-grained control over vehicle motion. These spaces enable smooth steering and acceleration adjustments, which are essential for realistic and safe driving behavior, but they also introduce challenges in terms of training stability and interpretability.

---

\*In this study we categorized the mapping strategy of the papers that has not explicitly mentioned the low level controller as direct.

Table 3 summarizes representative studies and categorizes continuous action space designs across several key dimensions.

1. Control Dimensionality:

- Lateral Only: The agent controls only the steering angle, with speed managed externally or held constant.
- Longitudinal Only: The policy controls speed via throttle, or braking.
- Combined Control: The most common approach, where the agent jointly predicts steering and speed commands. This may include:
  - Combined (negative braking): The Steering range is  $[-1, 1]$  and Throttle and brake values in  $[-1, 1]$ , with negative values ( $[-1, 0]$ ) representing braking.
  - Combined (no brake): The Steering range is  $[-1, 1]$  and Throttle is bounded in  $[0, 1]$ , with braking implicitly handled by reducing throttle.
  - Combined (separate brake): The Steering range is  $[-1, 1]$  and Throttle is bounded in  $[0, 1]$ , Brake is represented as a separate continuous output in range  $[0, 1]$ .

2. Action Range and Normalization:

- Normalized Range: Many studies use normalized outputs (e.g.,  $[-1, 1]$ ) and apply activation functions like `tanh` to bound the values.
- Dedicated or Tuned Ranges: A few studies customize control ranges, for example scaling steering from  $[-0.5, 0.5]$  to  $[-40^\circ, 40^\circ]$  [46].

3. Mapping Strategy:

- Direct Mapping: The majority of methods directly apply the policy’s output to the simulator’s low-level control interface.
- Indirect Mapping: A subset of works uses the agent’s output to specify control targets (e.g., speed or acceleration), which are then achieved using external control modules such as PID controllers.

As shown in Table 3, continuous action space design is diverse, yet several patterns emerge. Most approaches adopt a combined control structure where throttle and steering are handled simultaneously—either as normalized values in  $[-1, 1]$  [42, 80], bounded distributions like Beta [45], or custom physical units [97]. Several studies simplify the control interface by merging throttle and brake into a single signal, using negative values to represent braking [45, 49], while others maintain separate channels for acceleration and braking to increase control expressiveness and realism [35, 93]. In longitudinal-only setups, throttle or acceleration is learned by the agent, while steering is handled via traditional controllers like PID or Stanley, enabling modular architectures [47, 94]. Interestingly, several works employ hybrid pipelines where the RL policy outputs high-level targets (e.g., desired speed or acceleration), which are then executed by a PID controller [87, 94]. A few recent studies introduce predictive control over future steps [129] or encode control commands using task-specific scalings or semantic units [91], showing a trend toward physically interpretable action spaces. Finally, lateral-only designs [5, 92] demonstrate that decoupling steering from speed

control remains useful for focused lane-keeping tasks or simplified evaluations. Overall, the diversity in continuous action space design reflects ongoing efforts to balance precision, safety, and learnability in real-world driving scenarios.

**Table 3** Taxonomy of Continuous Action Space Designs in RL-Based Autonomous Driving

Control Type	Action Range	Mapping	Reference
Longitudinal only	Dedicated range	Direct	[12, 47, 95, 104]
Longitudinal only	Dedicated range	Indirect	[26, 52, 94]
Lateral only	Dedicated range	Direct	[92]
Lateral only	Dedicated range	Indirect	[5]
Combined	Not Mentioned	Not Mentioned	[4, 7, 8, 13, 30, 32, 33, 43, 48, 53, 82–84, 90, 96, 99, 117, 129]
Combined (negative braking)	Steering: $[-1, 1]$ , Throttle and Brake: $[-1, 1]$	Direct	[6, 34, 37, 41, 42, 45, 49, 80, 98, 100, 101, 106]
Combined (negative braking)	Steering: $[-1, 1]$ , Throttle and Brake: $[-1, 1]$	Indirect	[85, 87]
Combined (no brake)	Steering: $[-1, 1]$ , Throttle: $[0, 1]$	Direct	[16, 31, 44, 50, 79, 81]
Combined (separate brake)	Steering: $[-1, 1]$ , Acceleration: $[0, 1]$ , Braking: $[0, 1]$	Direct	[19, 26, 35, 86, 93]
Combined (specific range)	Dedicated range	Direct	[15, 91, 97]
Combined (specific range)	Dedicated range	Indirect	[46]

### 5.3 Hierarchical Studies

Across these seven works mentioned in Table 4, a two-tier hierarchy is the common thread: a small, discrete set of high-level maneuvers drives exploration and safety,

**Table 4** Two-tier hierarchical action space designs

High-Level Actions	Low-Level Actions / Controller	Reference
Follow lane; Go straight; Turn left; Turn right	Imitation-learning policies mapping observations to continuous steering $[-1,1]$ , throttle $[0,1]$ and brake $[0,1]$	[115]
Lane follow / wait; Lane change	PID-tracked waypoints selected as sub-goals	[118]
Maneuver-specific agents: Straight-drive; Turn-left/right; Lane-change	Dedicated discrete action set of steering and throttle for each maneuver	[119]
Options via Option-Critic: Lane following; Intersection handling; Obstacle avoidance	Atomic maneuvers (accelerate; brake; turn left; turn right) via PPO; Answer Set Programming (ASP) rules override	[51]
Five lateral offsets (Frenet frame waypoints) followed by PID controller.	Six discrete brake/throttle settings	[125]
Road-following; Left turn; Right turn	Continuous throttle/brake/steer sampled from bounded Beta distributions in range $[0,1]$ . These raw commands are then passed through PID controllers.	[11]
Maintain lane; Change left; Change right	DDPG actor-critic outputs continuous steering $[-1,1]$ or accel/brake $[-1,1]$	[103]

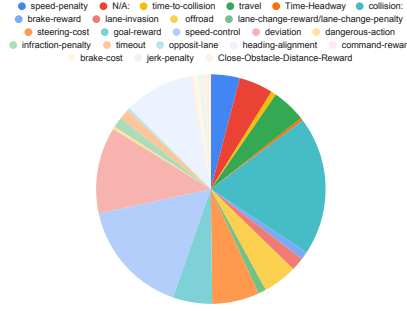
while low-level continuous controllers (imitation-learning policies, PID loops, or actor-critic networks) deliver smooth vehicle commands. The high-level action vocabularies vary—from simple lane-following vs. lane-change choices to richer option sets in Option-Critic architectures—but always aim to reduce policy complexity. At the low level, controllers are typically PID-based or learned via DDPG/PPO, striking a balance between precision and learnability. Together, these designs show that decoupling decision planning from continuous actuation not only improves sample efficiency but also enhances interpretability and provable safety in autonomous-driving RL systems.

## 6 Reward

Reward design is a fundamental component of RL in autonomous driving, as it directly shapes the agent’s behavior. A wide variety of reward terms have been proposed in the literature, reflecting the diverse priorities of different driving tasks—such as safety, goal completion, comfort, and rule compliance. Fig. 3 shows the distribution of reward terms across papers. Notably, collision penalties, speed control, and deviation from the lane center are among the most frequently used rewards, reflecting the emphasis on safety and lane-keeping behavior. Other commonly used rewards include goal-based



rewards, heading alignment, and steering smoothness (e.g., jerk and steering cost), which guide the agent toward efficient and comfortable driving.



**Fig. 3** Distribution of reward types used in RL for autonomous driving tasks.

Table 5 categorizes the different reward types along with their purpose and density, and cites key works that have employed them. This categorization aims to help future researchers choose or design reward functions suited to their task.

The definition of the most popular reward terms are mentioned in Table 6. Several studies introduce a diverse set of specialized reward components beyond the core categories to capture high-level commands, safety distances, comfort, efficiency, and rule compliance. *Command rewards* incentivize correct execution of navigational instructions—such as turning at intersections—in hierarchical frameworks [115]. *Opposite-lane penalties* discourage entering or remaining in lanes of oncoming traffic [113]. *Time-headway rewards* promote maintaining a safe time gap behind lead vehicles to avoid tailgating [12], while *time-to-collision penalties* penalize states with high imminent-collision risk [12, 90]. *Dangerous-action penalties* discourage deviations from expert human behavior [79]. *Collision-free-driving bonuses* provide a small positive reward for each collision-free timestep [73]. *Risk-reward terms* encode the time difference between the ego vehicle’s departure and a social vehicle’s arrival at a conflict area [95], while *time-differential rewards* align maneuvers temporally with available merge slots [52]. *Waypoint-following rewards* grant a bonus upon reaching designated waypoints [128]. *Failure penalties* impose costs for not reaching the goal, and *near-miss penalties* penalize moderate negative reward when pedestrians enter a predefined safety buffer [88]. *Fast-completion rewards* encourage rapid task execution [117]. *Steering or throttle jitter penalties* discourage abrupt control changes to ensure smoother, more comfortable driving [35]. *Safe-dist-to-lead-vehicle terms* maintain a safe following distance from the car ahead [35, 118], while *regular-timestep penalties and progress rewards* balance efficiency with purposeful movement [118]. *Unsmoothness penalties* discourage maneuvers that conflict with prior velocity profiles, promoting smoother transitions [118]. *Dist-to-lead-vehicle metrics* quantify the longitudinal gap to the lead vehicle [119], and *dist-to-nearest-obj rewards* promote safe clearance from nearby obstacles [119]. Finally, *incorrect-lane* and *incorrect-steer penalties* discourage misaligned lane occupancy and excessive steering during maneuvers

**Table 5** Categorization of Reward Functions Used in RL for Autonomous Driving

Reward Term	Purpose	Dense/Sparse	References
<b>Collision</b>	Safety	Sparse	[4–12, 15, 16, 18–20, 26, 30–33, 35, 38, 42–47, 49, 51–58, 61–73, 80, 85, 88, 91–94, 96, 99, 101–104, 106, 113, 115–119, 122, 125, 128]
<b>Speed Control</b>	Efficiency	Dense	[4, 6, 7, 10–12, 14, 16–20, 30, 32–34, 38, 41, 43, 45–47, 49–57, 59–61, 64, 65, 71, 72, 79, 81, 82, 84, 85, 87, 88, 91, 93, 95, 96, 99–103, 113, 116, 117, 122, 125]
<b>Deviation</b>	Safety	Dense	[101], [14], [5–7, 16–19, 31–35, 44–47, 49, 50, 59, 60, 62, 65, 66, 81, 82, 84–87, 91, 92, 103, 106, 119, 122, 125]
<b>Goal Reward</b>	Task Completion	Sparse	[20, 26, 54–57, 63, 64, 66, 80, 82, 86, 88, 90, 93–95, 100, 106, 115–119, 125, 128]
<b>Heading Alignment</b>	Safety	Dense	[7–9, 16–19, 32, 33, 44, 45, 49–51, 53, 59, 60, 62, 65, 66, 69, 84–86, 90, 96, 101, 103, 105, 119, 125]
<b>Steering Cost</b>	Comfort	Dense	[6, 10, 11, 14, 15, 18, 19, 30, 32, 41, 42, 45, 47, 53, 70, 88, 91, 92, 96, 106, 122]
<b>Jerk Penalty</b>	Comfort	Dense	[12, 54, 70, 104]
<b>Brake Cost/Reward</b>	Comfort	Dense	[17, 44, 56, 61, 91]
<b>Lane Invasion</b>	Safety	Dense/Sparse	[31, 46, 63, 68, 101]
<b>Offroad</b>	Safety	Dense	[6, 7, 11, 70, 79, 85, 93, 100, 116]
<b>Timeout</b>	Efficiency	Sparse	[43, 80, 85, 94]
<b>Infraction</b>	Safety	Dense	[42], [35, 45–47, 54, 71]
<b>Travel</b>	Efficiency/Progress	Dense	[14], [11, 12, 42, 63, 64, 66, 80, 91, 93, 96, 99, 104, 118]
<b>Lane Change</b>	Safety/Efficiency	Sparse	[4, 102, 103]
<b>Close-Obstacle Distance</b>	Safety	Dense	[33, 65, 105]
<b>Speed Penalty</b>	Safety	Dense	[8, 11, 42, 43, 57, 67, 70, 70, 71, 71, 80, 85, 91, 104, 122]
<b>Others</b>	N/A	N/A	[12, 35, 52, 73, 79, 88, 90, 95, 99, 113, 115, 117–119, 128]

[119], and *drive-side-walk penalties* ensure adherence to drivable space by penalizing overlap with sidewalks [99].

## 7 Terminal Condition

In RL, terminal conditions define when an episode ends. In CARLA-based driving tasks, terminal conditions are essential for shaping episode boundaries and determining when an agent has completed a task successfully or failed due to violations. Common

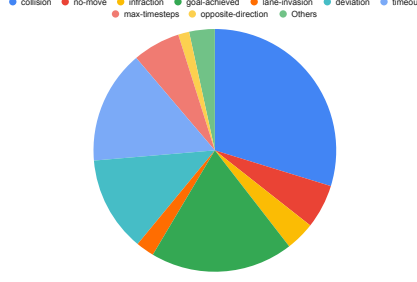
**Table 6** Definitions of Specialized Reward Terms

Term	Definition
Collision	Penalizes any collision with other vehicles, pedestrians, or objects to ensure safe driving.
Speed Control	Rewards the agent for maintaining a desired speed or penalizes deviation from a target speed range. Used to balance safety and efficiency.
Deviation	Penalizes the lateral offset from the lane center or planned trajectory to promote lane-keeping behavior.
Goal Reward	Provides a positive reward when the agent reaches its destination or finishing the task at hand. Encourages task completion.
Heading Alignment	Rewards alignment between the vehicle’s heading and the direction of the lane or target waypoint.
Steering Cost	Penalizes abrupt or large steering commands to encourage smoother control and comfort.
Jerk	Penalizes rapid changes in acceleration or steering, promoting smoother and more comfortable driving.
Brake Cost/Reward	Penalizes unnecessary or hard braking (cost) or rewards proper and timely use of braking (reward), aiming to reduce discomfort and unsafe stops.
Lane Invasion	Penalizes crossing over lane boundaries without intention, reflecting rule violation or unsafe maneuvering.
Offroad	Penalizes driving off the drivable road area, promoting adherence to lane and road limits.
Timeout	Penalizes failure to reach the goal within a time limit, encouraging efficient navigation.
Infraction	Penalizes the agent for violating traffic rules like running red lights, ignoring stop signs, etc.
Travel	Encourages consistent forward movement or progress through the environment based on the distance to the goal location.
Lane Change	Rewards safe and necessary lane changes or penalizes unnecessary or unsafe ones.
Close-Obstacle Distance	Rewards the agent for maintaining a safe distance from nearby obstacles while navigating close environments.
Speed Penalty	Penalizes speeding or exceeding safe velocity thresholds. Often complements speed control rewards. It can instead penalize the car not moving at all.

terminal triggers include collisions, lane invasions, reaching a goal location, going off-road, timeouts, or violating traffic rules. Properly chosen terminal conditions ensure that the learning process focuses on meaningful behaviors while penalizing unsafe or inefficient actions. They also influence the agent’s return estimates and training stability, particularly in sparse-reward settings. Table 7 summarized the most common terminal conditions used in the literature along with references. Fig. 4 illustrates the distribution of terminal conditions.

**Table 7** Common Termination Conditions in Autonomous Driving Scenarios (with example citations).

Condition	Fatal	Description	References
Collision	Y	The ego vehicle’s bounding box intersects with another object (e.g., another vehicle, pedestrian, barrier). Triggers include rear-end or side-impact collisions.	[4–12, 14, 20, 26, 30, 35, 38, 39, 41–45, 52, 54, 56–58, 60, 62–66, 69–73, 80–82, 85, 88, 90–93, 95, 96, 99–102, 104, 105, 122, 125, 128, 131]
No-move	N	The ego vehicle remains stationary or moves below a minimal threshold (speed or distance) for a set duration or number of timesteps. Common examples include failing to move for over 30–60 seconds.	[7, 11, 41, 43, 45, 51, 53, 60, 80, 101, 105]
Infraction	Y	The ego vehicle violates traffic rules (e.g., running red lights, ignoring stop signs, illegal turns). Sometimes includes driving on the wrong side of the road.	[14, 41, 42, 45, 46, 54, 60, 98]
Goal-achieved	N	The ego vehicle successfully reaches its target destination or completes the route (e.g., arriving at a waypoint or parking spot).	[4, 6, 8, 11, 14, 20, 26, 30, 37, 38, 42, 46, 51–54, 56–58, 63, 64, 66, 69, 72, 80, 81, 88, 90–93, 95, 99, 100, 104, 119, 125, 128]
Lane-invasion	Y	The ego vehicle crosses lane markings or invades adjacent lanes in an unintended manner. Sometimes used interchangeably with deviation if crossing lane boundaries is disallowed.	[18, 48, 98, 106, 131]
Deviation	Y	The ego vehicle strays from the intended path or lane centre beyond a tolerance (e.g., 2–3 m), or drives off-road entirely.	[4, 7, 11, 12, 18, 30, 37, 41–45, 48, 51, 53, 56, 57, 60, 79–81, 84, 85, 90–92, 96, 102, 103]
Timeout	N	The scenario’s simulation clock expires (e.g., 8 s to cross an intersection, 300 s to complete a route) without collision or reaching the goal.	[6–8, 10, 18, 26, 37–39, 46, 52, 54, 62, 70, 71, 73, 80, 85, 88, 91, 95–97, 100, 102, 104, 125, 128]
Max-timesteps	N	A hard limit on the total number of simulator steps (e.g., 800, 1000, or 3000 steps). Once reached, the episode terminates automatically.	[9, 15, 25, 30, 43, 48, 56, 57, 63, 82, 93, 106, 122]
Opposite-direction	Y	The ego vehicle orients and/or drives 180° against the valid heading, often leading to traveling into oncoming traffic.	[30, 66]
Others	N/A	Various custom or less common triggers, such as <i>excess-speed</i> (driving well above limits), <i>over-pass-destination</i> (traveling beyond the target), <i>stuck/blocking</i> (blocked in traffic or making no progress), <i>excessive-low-reward</i> (accumulating large negative rewards), <i>no-improvement</i> (no training progress for many episodes), <i>following-failure</i> , etc.	[11, 12, 25, 43, 66]

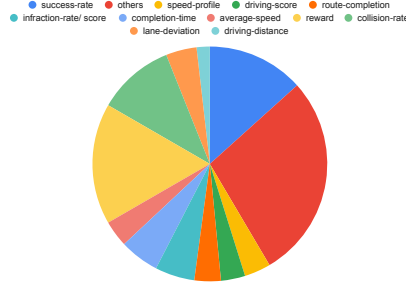


**Fig. 4** Distribution of terminal conditions used in RL for autonomous driving tasks.

## 8 Evaluation Metric

Evaluation metrics are critical for benchmarking RL agents in autonomous driving, offering quantitative insights into performance across diverse tasks and environments. The reviewed literature reveals a broad spectrum of metrics that reflect priorities such as safety, task completion, control stability, behavioral realism, and training efficiency. These metrics are often aligned with specific objectives—whether ensuring collision-free navigation, adherence to traffic rules, or replicating human-like behavior in maneuvers such as lane-following and intersection handling.

Based on our analysis, we categorize the metrics into four primary groups: (1) task completion and task progress (e.g., success rate, route completion, driving distance); (2) safety and comfort (e.g., collision rate, infraction score, lane deviation); (3) Efficiency (e.g., completion time, average speed); and (4) overall-preformance and overall driving quality (e.g., reward, driving score). Fig. 5 provides an overview of the distribution of evaluation metrics across the surveyed literature, while Table 8 summarizes the most commonly used metrics along with representative studies that adopt them.



**Fig. 5** Distribution of evaluation metrics used in RL for autonomous driving tasks.

While foundational metrics like success rate, collision rate, and reward dominate the literature, many studies introduce specialized metrics tailored to specific needs. For example, [101] assess lateral and longitudinal smoothness via the standard deviation of steering and speed, respectively. Other works incorporate human-likeness scores

**Table 8** Common Evaluation Metrics Used in RL for Autonomous Driving

Metric	Aspect	Description	References
Success Rate	Task-completion	Measures the percentage of episodes or routes completed without failure.	[4, 6, 7, 11, 12, 16–19, 26, 30, 34, 38, 41, 42, 45, 46, 51, 55, 56, 58, 60, 61, 69, 72, 80, 85, 86, 92, 94–96, 98, 100–102, 104, 113, 115–118, 125, 127]
Collision Rate	Safety	Tracks the frequency of collisions per episode or over the evaluation period.	[4, 5, 8–10, 26, 33, 39, 42, 49, 52–54, 56–58, 61, 63, 65, 69, 81, 82, 85, 86, 88, 90, 93, 95–98, 104, 105, 118, 125]
Route Completion	Task-progress	Indicates the percentage of the planned route successfully followed by the agent.	[8, 14, 17, 25, 32, 41, 45, 53, 59, 71, 98, 127]
Infraction Rate / Score	Safety	Captures violations such as off-road driving, collisions, or running red lights.	[8, 14, 17, 19, 32, 41, 42, 45–47, 53, 56, 59, 60, 71, 86, 98, 127]
Completion Time	Efficiency	Measures how long it takes the agent to finish a task or reach a goal.	[4, 9, 16, 31, 39, 43, 48, 52, 54, 61, 63, 72, 82, 88, 92, 93, 116, 131]
Average Speed	Efficiency	Reflects the mean speed of the ego vehicle. Sometimes analyzed with variance to evaluate consistency.	[10, 16, 47, 49–51, 62, 72, 87, 96, 102, 125]
Reward / Return	Overall performance	Average reward obtained per episode during evaluation. It reflects policy performance with respect to the task’s reward function.	[5, 6, 9, 15, 18, 20, 33, 35, 37, 38, 41–44, 48–51, 53, 55–58, 63–66, 69, 70, 73, 79, 80, 82, 83, 86, 87, 91–93, 96, 97, 99, 103–106, 115, 117–119, 122, 125, 127, 128]
Driving Score	Overall driving quality	A composite metric combining route completion, infraction penalties, and sometimes efficiency.	[14, 17, 32, 41, 45, 53, 59, 71, 82, 98, 127]
Speed Profile	Efficiency/Comfort	A speed profile is a time-based or distance-based representation of how a vehicle’s speed changes along its route. It gives a detailed picture of the vehicle’s acceleration and deceleration during a driving episode.	[10, 39, 44, 53, 56, 57, 64, 67, 84, 87, 101, 116]
Lane Deviation	Safety/Comfort	Lane deviation measures how far a vehicle drifts or stays away from the center of its lane while driving. It is a critical metric for evaluating how well an AV (or human driver) maintains lane discipline.	[4, 5, 8, 43, 44, 50, 51, 53, 61, 81, 84, 90, 96, 106]
Driving Distance	Task-progress/Efficiency	Driving distance refers to the total length of the path that a vehicle actually travels during a driving episode or simulation.	[50, 51, 53, 64, 65, 72]
Others	N/A	There are some papers that used other metrics including but not limited to <i>mean-route-error</i> : Average deviation from planned route (trajectory error), <i>inference-time</i> : Time taken by model to make a decision, <i>heading-deviation</i> : Difference between vehicle heading and road direction, <i>no-collision</i> : Boolean or count of episodes without collision, <i>average-episode-length</i> : Average number of steps per episode, <i>convergence</i> : Metric to track if training is stabilizing, <i>lane-invasion-rate</i> : Frequency of lane marking violations, <i>training-time</i> : Total training duration, episode length	[31, 37, 131], [18, 73, 88], [51, 84, 103], [38, 52, 72, 113], [8, 62, 101, 102], [33, 53, 55, 68], [65, 85, 118], [55, 88, 128], [43, 56, 57].

[39], command execution accuracy [68], or critical driving comfort measures [67]. Additionally, training-focused evaluations such as action-value convergence, episode length trends, and convergence diagnostics are increasingly used to track learning stability and robustness [91].

In total, we identified 95 unique evaluation metrics across the surveyed literature. Although not all are widely adopted, many provide valuable insights for specific deployment scenarios, particularly in safety-critical or human-in-the-loop contexts. To maintain clarity, Table 8 highlights metrics used in at least five studies (except the others category that has the least number of 3).

**Table 9** Scenario categories vs. CARLA towns. A checkmark (✓) indicates the scenario has been studied in that town. “Custom” = non-official or heavily modified maps.

Scenario Type	Town01	Town02	Town03	Town04	Town05	Town06	Town07	Town08	Town09	Town10	Custom
Urban driving / route completion	✓	✓	✓	✓	✓	✓	✓			✓	
Intersections / T-junctions	✓		✓			✓	✓				✓
Roundabouts			✓								
Highway / lane-change				✓	✓	✓					✓
Parking / partial control					✓	✓					
Racing / high-speed tracks											✓
Domain adaptation / multi-weather	✓	✓	✓	✓	✓	✓	✓				
Custom envs / special setups											✓

## 9 Towns and Scenario

In CARLA-based RL research, towns refer to different simulated environments (e.g., Town01 to Town10), each with unique layouts, road structures, and environmental features. Scenarios, on the other hand, define specific driving tasks or challenges within these towns, such as lane following, intersection handling, or collision avoidance. While CARLA includes some predefined maps and features, most scenarios are manually created by researchers to focus on particular behaviors or challenges. Some studies adopt standardized benchmarks like CoRL 2017, NoCrash, or the CARLA Leaderboard, which provide consistent evaluation settings. Others design custom routes and traffic situations to test specific hypotheses, evaluate generalization to unseen roads, or simulate rare edge cases. The choice of town and scenario plays a key role in shaping task difficulty and determining how well a trained policy performs in new or complex environments. Table 9 provides a consolidated overview of the scenario types examined in the literature and their corresponding CARLA towns. While CARLA offers ten official towns, the majority of studies focus on a limited subset, particularly Town01 and Town02, which feature relatively simple urban layouts suited for route completion and intersection handling. More complex designs such as roundabouts (Town03) and multi-lane highways (Town04–Town06) are less frequently explored, though they offer richer testing grounds for lane-change and merging maneuvers. Parking and partial-control experiments are concentrated in Town05–Town06 due to the inclusion of parking lots, whereas specialized tasks such as high-speed racing are typically conducted in custom-designed environments outside the default towns. This distribution highlights both the historical availability of certain maps and a tendency in the community to prioritize simpler, more accessible settings over more challenging layouts.

### 9.1 Urban Driving in Official CARLA Towns

A large subset of papers focuses on general urban driving and route-completion tasks using the default CARLA towns (Town01 through Town10) under varying levels of traffic density, weather, and route complexity. Several works concentrate on Town01 and Town02 for basic or intermediate urban tasks, such as lane keeping and navigation

with occasional turns. For instance, [31, 131] both use Town01 with predefined routes (sometimes with dynamic waypoints but no obstacles). Similarly, [50, 113, 115] employ Town01 and Town02 to study how RL agents adapt from simpler to more complex road structures, including traffic lights and pedestrian traffic. Other works ([6, 19, 46, 113, 127]) similarly center on Town01/Town02 benchmark (e.g., CoRL2017 or NoCrash benchmarks).

Town03 is a frequent choice for more complex urban settings, with denser traffic or trickier intersections. Papers such as [18, 32, 45, 122] train or evaluate on Town03 with heavy traffic (up to 100 vehicles or more), often testing advanced methods like motion-prediction-based RL. Similarly, [33, 63, 65] also adopt Town03 for navigating crossroads, roundabouts, and congested roads. Some authors ([14, 59, 101]) use multiple towns or LeaderBoard routes encompassing Town03 to rigorously validate their policies.

Meanwhile, Town04 appears in works such as [68, 82, 90, 96, 103], often featuring highways, urban blocks, or both. Several papers ([7, 12, 41, 42]) also highlight Town04 in meta-learning or continual-learning contexts, sampling diverse traffic and weather conditions. Town05 is popular for route-following or intermediate urban tasks ([20, 60, 66, 80]), and Town06 and Town07 often appear in more specialized scenarios (e.g., [17, 39, 43, 45, 48, 49]). Town10 appears as well, such as in [64], which tests dynamic weather conditions and city-style roads.

Within these general urban-driving papers, some focus on dense traffic or advanced collision-avoidance behaviors, like [47] (Town02 with 30 vehicles and 50 pedestrians), [85] (emergency vehicles in heavy traffic), and [94] (scaled-up intersections). Others emphasize particular training regimens, such as curriculum learning ([49]), offline RL ([127]), or an imitation+RL hybrid ([45], [30]). In all, these works collectively address the challenge of teaching an autonomous agent to navigate city roads with traffic lights, pedestrians, and dynamic vehicle populations.

## 9.2 Intersection-Focused and T-Junction Scenarios

A separate set of studies zooms in on intersection handling, particularly T-junctions or unsignalized intersections. For instance, [93] evaluates navigation in a T-intersection with progressively increasing pedestrian/vehicle density (up to 450 vehicles/pedestrians). Town03 also appears frequently for unprotected or unsignalized intersections ([104]), testing the agent’s ability to yield and merge properly. In [94], training focuses on a single left-turn intersection, with testing on four additional intersections of increasing complexity, while [67] designs “complex left-turn across path” collision-prone maneuvers. Similarly, [56] addresses a simpler single intersection in Town01, emphasizing red-light compliance and safe car-following behaviors. Other works such as [63] also highlight intersection-dense environments, scaling traffic to hundreds of vehicles and pedestrians. These intersection-focused studies often delve into safety-critical decision-making—when to brake, how to sequence turns, and how to handle partial occlusions (e.g. [125] with pedestrians in the blind-spot.)



### 9.3 Roundabouts and Specialized Urban Maneuvers

Several papers concentrate on roundabouts, an environment that combines intersection-like decision-making with circular traffic flow. In Town03, [62] contrasts a four-way roundabout (with intersections) to a simpler U-turn loop, while [45] and [17] also include roundabouts in their urban routes (alongside highways and city blocks). Roundabouts force the RL agent to negotiate merges, exits, and near-constant lateral control challenges.

Additionally, some works target narrower urban driving maneuvers beyond intersections or roundabouts. For example, [16] addresses a specialized high-speed track within an urban-ish environment, albeit more akin to a racing scenario. [5] and [105] highlight speed-change maneuvers in an urban setting, examining how RL policies handle varying velocity constraints (e.g., from 8 m/s to 16 m/s). Meanwhile, studies like [25] define multiple routes in Town03, some with sharp 45–180° turns, to stress-test an end-to-end policy’s turning ability in dense urban grids.

### 9.4 Highway and Lane-Change Scenarios

Another broad category focuses on highway driving or lane-change maneuvers in multi-lane roads. Certain papers combine official towns with highway segments—for example, [60] trains partly in Town05 and tests on a highway environment to see if the RL agent generalizes from urban to highway. In [15] sets up a highway scenario in Town04 with 30 vehicles in day/night conditions, while [97] uses a dual-lane highway with fluent vs. congested traffic. [10] presents a five-lane road in Town06, replicating ramp entries and variable traffic speeds.

Lane-changing is explicitly studied in [102] (a custom three-lane highway), [4] (a 2-lane straight road with discretionary lane changes), and [92] (two-lane roads with static/dynamic obstacles). These works typically measure collisions, comfortable lateral maneuvers, and safe merges, often restricting the environment to minimal intersections or signals. Such highway-focused research underscores how RL must handle higher speeds and more frequent lane merges than in typical urban blocks.

### 9.5 Parking and Partial-Task Scenarios

A smaller cluster of studies focuses on parking maneuvers or partial-control tasks (e.g., only longitudinal control). For instance, [80] tackles vertical and parallel parking in Town05, showcasing how RL can learn precise maneuvering into parking bays.

Considering the partial task scenarios [39] in Town06 considers only throttle and brake with a single lead vehicle, while steering remains static—this isolates speed-control strategies for collision avoidance. Such partial or specialized tasks often serve to simplify the action space so the RL agent can master one dimension of driving (speed or lateral alignment) before moving to more complex scenarios.

### 9.6 Racing or High-Speed Specialized Tracks

A distinct set of references moves beyond ordinary urban or highway driving into racing or closed-course scenarios. [16] implements a 1 km custom racetrack within

CARLA, containing tight turns and obstacles in “easy” or “hard” variations. Some works ([57, 106]) also mention custom closed-circuit setups, though they may not always be labeled “racing.” These tasks emphasize high-speed cornering, minimal margin for lane deviation, and sometimes transfer to real-world scaled vehicles. By contrast, typical urban RL tasks rarely push the car to racing speeds or require strictly apexing corners.

## 9.7 Domain Adaptation and Multi-Weather Training

Many papers investigate weather variations and domain adaptation—teaching agents to generalize from sunny noon to rain, fog, dusk, or even new towns. Works like [49] explicitly train in “soft” weather presets (e.g. clear, cloudy noon, etc) and test in “hard” conditions (e.g. heavy rain, sunset, etc). [30, 59, 96, 127] also vary meteorological settings (e.g., HardRainNoon, WetNoon, ClearSunset) to measure how robustly a learned policy transfers. Meanwhile, [48] trains the model in Town07 on late evening and tests in clear noon or hard rain, illustrating domain shift in both lighting and precipitation.

Separately, meta-learning or continual learning approaches ([7, 12, 41, 42]) treat each new weather, traffic pattern, or town layout as a “task.” The RL agent samples from multiple tasks at training time, improving its capacity to adapt quickly to unseen conditions (e.g., new towns or more severe weather). Some studies [41, 42] incorporate Town01–Town04 as meta-training environments, then measure performance on Town05–Town06.

## 9.8 Custom Environments and Specialized Configurations

Finally, a notable subset of studies forgo official CARLA towns entirely or modify them heavily. They create custom worlds for tighter experimental control, specialized tasks, or real-to-sim synergy. For example, [102] builds a three-lane highway map, [92] a two-lane 420 m road, and [125] a scenario with an occluding parked van to test pedestrian-blind-spot crossing. [9] deliberately uses a minimal “Town” to isolate reward function variations, while [79] sets up a simpler lane-keeping environment for online RL. Some papers ([13], [44]) mention using CARLA generally but do not specify which map or scenario, often focusing instead on algorithmic innovations (e.g., new RL frameworks or 3D perception modules).

Racing or track-based custom worlds appear in [16] and [106], the latter modeling a real test track for industrial validation. Others, like [81] and [57], mention creating diverse random routes or specialized setups, sometimes spanning multiple towns or custom spawns, to avoid overfitting. These custom environments underscore the flexibility of CARLA for both academic research and practical testing scenarios beyond the standard Town0X maps.

## 10 Limitation

Despite recent progress, RL-based autonomous driving in CARLA continues to face several persistent challenges that hinder scalability, safety, and real-world deployment.

These limitations span fundamental issues such as sample efficiency, generalization, safety guarantees, perception constraints, scenario realism, and long-horizon planning, among others. Table 10 provides a structured summary of the major limitation categories, while the following subsections discuss each in detail, supported by representative studies from the literature.

## 10.1 Sample Efficiency and Data Cost

Many recent autonomous-driving frameworks exhibit limited sample efficiency and impose high data costs. A common assumption that every expert demonstration is perfectly optimal leads to assigning a uniform, high reward; in datasets such as the CARLA autopilot logs, approximately 10% of the recorded actions are actually poor, injecting reward noise that hinders policy learning [59]. Models with dual-path architectures that improve observation interpretability, such as SIMOUN, still learn slowly because they must rely on finite datasets and sparse reward signals [82]. Other end-to-end approaches likewise struggle with low sample efficiency, demanding thousands of interaction episodes before achieving acceptable performance [91]. Overall, inadequate sample efficiency and weak generalization constrain RL agents in complex urban traffic, where rare corner cases are difficult to encounter within limited training tasks; meta-RL augmented with safety constraints or adversarial learning has been proposed as one route to more efficient and safer driving, but its benefits remain to be fully validated [41].

## 10.2 Generalization / Transfer and Sim-to-Real Gap

Despite impressive in-simulation results, contemporary autonomous-driving agents continue to face substantial challenges in generalisation, transfer, and sim-to-real deployment. ROACH [45], for example, performs strongly in CARLA yet struggles to transfer because of visual, dynamical, and interaction mismatches, and because simulators still model surrounding traffic with limited realism. A pronounced distribution gap between the environments encountered during training and those met during testing further degrades policy reliability [94]. Policies trained in a single town or under a single weather regime often fail when exposed to new cities or conditions, and many studies omit any true sim-to-real evaluation [98]. End-to-end conditional imitation learning, which maps images directly to control commands, finds it difficult to extract the correlations needed for robust decisions and cannot guarantee real-time performance in dense urban traffic [129]. Even safety-aware deep RL variants still display systematic errors—such as consistently wide turns—and, given the safety stakes, remain unsuitable for real-world deployment until success rates approach 100% [6]. Network-based longitudinal controllers are expected to benefit from additional human demonstrations, suggesting that richer data could improve real-world adaptability [39]. Current online safety-learning methods are unable to compute next-state estimates for dangerous actions in CARLA and must reuse the current state; future work proposes synthesising those transitions with GANs to assess generalisability [79]. Techniques that rely on segmentation to improve generalisation separate perception from

**Table 10** Summary of Limitations in RL-based Autonomous Driving with CARLA

Limitation Category	Cate-	Key Issues and Challenges	Representative References
Sample and Data Cost	Efficiency	Low sample efficiency; reliance on large datasets and thousands of episodes. Expert demonstrations assumed optimal despite noise (e.g., ~10% poor actions in autopilot logs). Sparse rewards and reliance on finite data hinder learning.	[41, 59, 82, 91]
Generalization / Sim-to-Real Gap		Strong performance in-simulation but poor transfer due to visual, dynamic, and traffic realism gaps. Policies trained in a single town/weather regime fail under new conditions. Limited sim-to-real evaluations. Safety-aware RL still exhibits systematic errors (e.g., wide turns). Scenario-specific training hampers transferability.	[6, 7, 39, 45, 51, 79, 94, 98, 102, 118, 129]
Safety and Risk-awareness	Guarantees	Lack of explicit safe-RL mechanisms; unstable across traffic situations without fine-tuned hyperparameters. Safety often delegated to external shields. Simulations omit actuation errors and rare corner cases. Constrained RL improves safety but reduces comfort. Safety-utility trade-off remains underexplored.	[10, 12, 32, 42, 85, 119]
Sparse and Shaped Rewards	Poorly	Rewards often sparse, inconsistent, or overly coarse. Aggregating multiple objectives into single scalars hides trade-offs. Few studies analyze macroscopic traffic efficiency. Oversimplified rewards can trap agents in local optima. Lack of systematic reward design analysis.	[4, 9, 47, 59, 99, 115]
Perception and Sensor Fusion	Limits	Over-reliance on CARLA’s waypoint oracle or RGB-only inputs. Limited use of LiDAR, RADAR, GPS, or inertial data. Simulation setups often differ from real sensor suites (e.g., lack of rear cameras, ignored depth/LiDAR channels). Poor robustness to adverse weather or unseen towns.	[8, 11, 16, 18, 31, 33, 53, 55, 62, 68, 69, 93, 105]
Scenario and Interaction	Realism	Simplified scenarios: fixed throttle, static obstacles, discretized action spaces, narrow geographic coverage. Missing pedestrians, bicycles, or complex multi-agent interactions. Evaluations often under fixed weather and limited maps. Robustness to unseen layouts rarely tested.	[5, 20, 25, 26, 33, 50, 65, 67, 68, 71, 90, 104, 122, 125]
Hierarchical / Long-Horizon Planning		Current agents are short-horizon and reactive, lacking explicit route planning or long-term reasoning. Few integrate trajectory prediction or multi-path planning. Limited foresight prevents goal-directed, mission-level control.	[47, 68, 70, 88, 94]
Hyper-Parameter Sensitivity and Training Stability		Performance highly sensitive to hyper-parameters, batch sizes, and replay strategies. On-policy methods (e.g., PPO) data-hungry; off-policy methods fragile without tuning. Limited GPU budgets and short training times worsen instability. CNN backbones not optimized for urban scenes.	[19, 45, 46, 62, 81, 83, 96, 99]
Insufficient World Validation	Real-	Overwhelming reliance on simulation; most methods defer real-world testing. Model-free RL impractical due to data demands. Prototype systems lack real sensor data. Controllers trained in simulation fail in real-world layouts. Real-world scalability and robustness remain unverified.	[18–20, 43, 46, 55, 60, 64, 73, 80, 85, 92, 106, 131]
Traffic-Rule and Social-Norm Compliance	Com-	Few frameworks enforce explicit traffic rules or social norms. Agents may generate illegal manoeuvres or ignore communicative cues (e.g., signals, brake lights). Road user models remain homogeneous, neglecting cultural/social diversity.	[53, 72, 95, 103, 122]
Performance Gap		Despite progress, RL policies underperform expert humans and SOTA baselines. Poor handling of dense traffic and non-standard intersections. Weak lateral control on curves. Policies lack rationality, interpretability, and trustworthiness.	[16, 32, 37, 49, 94]
Limited Diversity	Behavior	Agents exhibit narrow manoeuvre sets (lane-keeping, intersections, traffic lights). Rarely handle lane changes, roundabouts, or cooperative merges. Heavy reliance on scripted autopilot and homogeneous actors. Limited scenario templates and lack of heterogeneous traffic agents.	[30, 35, 38, 44, 46, 49, 50, 52, 55, 58, 64, 65, 69, 93, 97, 103, 122, 129]

control and may underuse shared representations; forthcoming work will explore alternative architectures and stronger data augmentation [7]. Existing benchmark suites lack low-speed traffic and therefore fail to expose agents to an important subset of urban scenarios; planned extensions will broaden test cases and include model-based

RL baselines [102]. Finally, guided RL benefits from explicit rules, yet too many handcrafted constraints risk overfitting and impairing generalisation; evaluations on the more demanding CARLA Leaderboard are intended to clarify this trade-off [51]. Scenario-specific policy training further hampers transfer. Manoeuvres learned for a given junction or speed profile rarely integrate into a unified behavioural graph, so skills do not carry over to unseen traffic compositions or road morphologies [118]

### 10.3 Safety guarantees and Risk-awareness

Many recent reinforcement-learning approaches for autonomous driving provide only limited assurances of safe operation. Most do not incorporate explicit safe-RL mechanisms, so there is no formal guarantee that the learned policy will avoid hazardous manoeuvres once deployed [32]. Empirical studies confirm that agents can remain unstable across a range of traffic situations unless their reward structure or hyperparameters are finely tuned [85]. Some frameworks delegate responsibility for collision avoidance to an external symbolic “safety shield,” effectively decoupling safety from learning [119]. Simulation evaluations tend to be overly idealised: they neglect actuation errors and therefore overestimate real-world safety margins, while coarse decision- and control-update intervals further mask timing-critical failures [10]. Because RL agents encounter only a fraction of long-tail corner cases during finite training, unsafe behaviour can still emerge under rare conditions; proposed remedies include embedding safety constraints directly into the objective or adding adversarial curricula to expose risky scenarios earlier [42]. Constrained variants such as EM-SAC demonstrate that stricter limits on deceleration reduce risk but at the cost of ride comfort, illustrating a broader trade-off between safety and utility that remains insufficiently quantified [12]. Overall, providing verifiable safety guarantees and robust risk awareness is an open challenge that future work must address before RL-based policies can be trusted in safety-critical autonomous-driving applications.

### 10.4 Sparse / Poorly Shaped Rewards

A pervasive limitation of current autonomous-driving research is the continued reliance on sparse or poorly shaped reward signals. In imitation-reinforcement hybrids, identical low-level manoeuvres receive a high reward when executed by the demonstrator but a low reward when generated during exploration, creating an internal inconsistency that can inflate the estimated value of demonstration actions and destabilise learning [59]. More broadly, existing studies seldom analyse how reward design influences macroscopic traffic efficiency [47]. Frameworks that embed hierarchical decision layers often retain rudimentary reward definitions whose coarse granularity restricts both flexibility and sample efficiency [115]. Even when multiple objectives are acknowledged, they are typically aggregated into a single scalar, obscuring trade-offs that—in real deployments—must be handled separately; a truly multi-dimensional reward formulation remains to be implemented [99]. Experimental evaluations are likewise confined to a small set of handcrafted reward variants, limiting generality [9]. Finally, oversimplified rewards can trap agents in local optima [4]. Collectively, these shortcomings

highlight the need for better-shaped, task-decomposed, and safety-aware reward functions, as well as for systematic studies of how reward design affects both microscopic behaviour and network-level traffic outcomes.

## 10.5 Perception limits and sensor fusion

Most contemporary RL pipelines for autonomous driving continue to operate with impoverished perceptual inputs and minimal sensor fusion, undermining their robustness and real-world fidelity. Numerous studies build directly on CARLA’s internal waypoint oracle, thereby bypassing the need to learn spatial reasoning and limiting adaptability once high-level guidance is absent; future deployments must replace this crutch with camera- and LiDAR-driven perception that matches a production vehicle’s sensor stack [31]. Vision modules themselves are frequently restricted to raw RGB frames, which are sample-inefficient and brittle; semantic renderings would provide more task-relevant abstractions and yield safer policies [16]. Several frameworks omit complementary modalities such as LiDAR, RADAR, GPS, or inertial data, despite their proven value for localisation and obstacle detection in adverse conditions [18, 33, 53, 55, 62, 68, 93]. Even when depth is available, other essential channels—for instance LiDAR in RGB-D setups—are ignored [69], and real-vehicle prototypes sometimes lack rear-facing cameras or exhibit mechanical steering dead-zones that constrain manoeuvrability [105]. Limited perception using semantic segmentation exposes models to poorly handled weather shifts, and degraded success rates in unseen towns and rain scenarios [11]. The over-idealised sensor suites used in simulation depart markedly from on-road hardware, as acknowledged by recent work that calls for richer multi-modal inputs, promoting more modular, encoder-based vision backbones [8]. Until such multi-sensor fusion and noise-aware perception become standard, the generality and safety of RL-based driving policies will remain fundamentally constrained.

## 10.6 Scenario Realism and Multi-Agent Interaction

Most RL evaluations for autonomous driving still rely on highly idealised scenarios that diverge from real-world traffic. Some studies fix longitudinal control to a constant throttle, preventing the agent from modulating acceleration in response to context [68]. Collision handling is likewise simplified: agents are tuned to skirt static props such as poles or fences but receive no training on moving vehicles or pedestrians, leaving dynamic-obstacle avoidance essentially untested [50]. At the control layer, many algorithms discretise steering and acceleration commands; coarse, fixed action sets reduce manoeuvre smoothness and precision which is not the case with real world driving which relies on continuous action space [20, 65, 67]. Geographic coverage is narrow as well: the DeFIX policy, for example, is validated only on CARLA’s Town05, so its performance in unseen maps remains unknown [71]. Entire pipelines are built and tested in virtual worlds devoid of live pedestrians or cross-traffic, with future work merely proposing to add additional cameras or more demanding scenes [25]. Controllers evaluated on empty roads achieve high success rates that quickly erode once dynamic obstacles or poor visibility are introduced, motivating plans to pair the controller with an external motion planner for ramp merges, overtakes, or low-visibility driving [5].

Other studies validate their method on only pedestrian-interaction scenarios; it should further be evaluated on more complex situations [125]. Most decision-making modules ignore bicycles, scooters, or pedestrians altogether—an omission acknowledged as a major avenue for future research on multi-agent intersections [104]. Likewise, frameworks tuned for low-speed dual-lane scenarios have yet to prove transferable to broader road networks or on-road trials [90]. Work on highway merging is typically limited to a simpler lane merging scenarios; extending these models to multiple interacting agents and standardizing benchmark datasets is recognised as essential for fair comparison and uncertainty modelling [26]. In some studies experimental protocols themselves reveal further weaknesses: mostly evaluations assume fixed weather conditions, leaving algorithmic resilience to rain, fog or low-illumination settings unverified [122]. Similarly, although strong results are often reported on a restricted subset of CARLA scenarios, robustness to unseen layouts is rarely demonstrated [33].

Together, these observations reveal a pronounced gap between the tidy, single-agent, discretised scenarios prevalent in current research and the dense, multi-actor, continuous-control conditions faced on public roads; closing this gap will require richer multi-agent benchmarks, continuous-action policies, and systematic evaluations that span diverse maps, traffic mixes, and environmental conditions.

## 10.7 Hierarchical / Long-Horizon Planning

Current RL agents for autonomous driving are still limited to short-horizon, reactive control and offer only rudimentary support for hierarchical or long-range planning. Some studies acknowledge that incorporating trajectory-prediction modules could substantially improve foresight but leave such components to future work [47]. Many controllers can avoid collisions and lane invasions yet cannot steer toward an explicit destination, because they operate without a route planner or goal-directed policy layer [68]. Even when the state representation encodes local road geometry, it omits navigation cues for instance intended future trajectories, preventing the policy from reasoning beyond the immediate scene [94]. Some systems can perform only simple, short-term manoeuvre selection and cannot sequence actions over extended horizons [70]. Some frameworks like HyLEAR [88] favours comfortable but often longer routes, and ongoing work is adding multi-path planning, prediction of neighbouring vehicles, and noise-aware perception precisely to extend its look-ahead capability. Overall, the absence of integrated, long-horizon planning modules remains a significant barrier to deploying RL-based driving stacks in real-world, goal-oriented missions.

## 10.8 Hyper-Parameter Sensitivity and Training Stability

A persistent obstacle to reliable RL based autonomous-driving is the pronounced sensitivity of both data efficiency and final performance to algorithmic hyper-parameters, network capacity, and training budgets. On-policy schemes such as PPO consume millions of interactions while discarding past experience, making them data-hungry and fragile in rare or edge-case situations; shifting to off-policy algorithms has been proposed as a remedy [46]. The absence of meta-RL modules limits cross-task adaptability [19]. Even within off-policy replay buffers, performance can hinge on finely tuned



parameters: Focused Experience Replay requires careful selection of the sampling-spread coefficient  $\beta_1$ ; poor choices induce premature convergence or loss of behavioural diversity [83]. Maximum-entropy variants show similar brittleness—larger batch sizes markedly reduce post-turn collisions and lane departures [96]. Constraints on hardware and training time further undermine stability: limited GPUs preclude use of richer perception branches, while truncated training runs leave DQN agents well short of convergence [62, 96]. Existing CNN backbones are rarely optimised for urban scenes; Neural Architecture Search (NAS) is flagged as future work to mitigate this mismatch [99]. Researchers also report vanishing gradients and slow convergence during feature extraction [81], and even expert distillation models plateau below leaderboard ceilings—suggesting that larger or more expressive networks are still needed [45]. Collectively, these findings highlight the fragility of current pipelines to hyper-parameter sensitivity and resource limitations, underscoring the need for systematic sensitivity studies, automated tuning, and more compute-efficient training algorithms to achieve robust, transferable performance.

## 10.9 Insufficient Real-World Validation

To date, the overwhelming majority of deep-RL driving studies remain confined to simulation, leaving their real-world generalisation essentially unverified. Several groups explicitly defer prototype deployment to future work: ROS-based controllers are good choice for transfer to an NVIDIA-powered test vehicle [131], affordance-encoded policies intend to retrain on real camera feeds before road world deployment [60], and both deductive RL, knowledge-enhanced, and decision making frameworks acknowledge that only simulated benchmarks have been considered so far [18, 19, 55]. Simulation-only evaluation also characterises work on urban driving policies, intersection navigation, emergency-vehicle manoeuvring, and hybrid RL-planning pipelines, all of which report no physical experiments [20, 46, 85]. Even studies that build custom campus maps or golf-cart prototypes have yet to collect any real sensor data [64]. The impracticality of deploying Model Free policies like SAC and PPO in the real world is apparent as they require large amounts of training data and they might still deviate from the intended trajectory within short amount of time [43, 106]. Parking controllers trained for 200 simulated episodes might fail in richer layouts, underscoring the gap between simulated and real parking environment [80]. Likewise, control studies considering the driving styles still rely on hand-tuned rewards that must be recalibrated with human driving data before field tests [92]. Finally, uncertainty-aware decision modules double computational latency and have so far been exercised only in high-fidelity simulators [73]. Across all these efforts, real-world scalability, robustness to sensor noise, and compliance with traffic regulations remain open challenges; comprehensive road trials are needed to validate simulated gains and to expose failure modes that only emerge in uncontrolled environments.

## 10.10 Traffic-rule and social-norm compliance

Current decision-making frameworks for autonomous driving lack rigorous mechanisms for traffic-rule and social-norm adherence. Most RL pipelines optimise for task success



without an explicit rule-checking layer, leaving the legality of generated manoeuvres uncertain [53, 72, 103]. Some studies likewise neglect communicative cues—such as turn indicators and brake lights—that are vital for inferring the intentions of surrounding vehicles [122]. Behavioural models of other road users remain overly homogeneous, focusing on kinematic states while ignoring cultural differences, driver habits and situational social preferences [95]. Accordingly, future research needs to embed formal traffic regulations directly into optimisation objectives or enforce them through constraint-satisfaction layers, enrich training environments with heterogeneous, rule-following agents to achieve robust social-norm compliance and legally safe behaviour [53, 72].

### 10.11 Performance Gap

Despite recent advances, the decision-making and control policies reported in the current literature remain noticeably inferior to both experienced human drivers and the best-performing machine baselines. Deep Imitative Reinforcement Learning (DIRL), for example, offers markedly better data efficiency and robustness, yet its ultimate driving proficiency still falls short of expert human performance, suggesting that policy expressiveness and learning capacity are insufficient for the complexities of real-world traffic [16]. In comparisons with SOTA such as CIRL [132], CAL [133] and CIRLS [134], some proposed methods underscoring persistent deficiencies in learning efficiency, generalisation and sample usage [49]. When exposed to denser traffic or non-standard intersections, trained agents exhibit erratic interaction handling, reflecting limited capacity to model the negotiation dynamics of multi-agent driving [94]. At the control level, lateral tracking accuracy deteriorates on curved segments, indicating an incomplete internalisation of vehicle dynamics that would likely be accentuated on a physical platform [37]. Finally, while several frameworks optimise for efficiency, they seldom incorporate explicit notions of decision rationality, interpretability or passenger trust, leaving the qualitative soundness of the chosen manoeuvres largely unexamined [32]. Addressing these shortcomings—by benchmarking against diverse weather and traffic scenarios, tightening integration of vehicle-dynamics knowledge, and enforcing interpretable decision criteria—remains essential for closing the gap between current autonomous policies, state-of-the-art baselines and human expertise.

### 10.12 Limited Behavior Diversity (Different Manuvers)

Existing deep-RL and imitation-learning pipelines exhibit a narrow behavioural repertoire that falls well short of the manoeuvre diversity required for dependable urban autonomy. Most studies confine training and evaluation to elementary tasks—lane keeping, traffic-light compliance or intersection handling—thereby sidestepping more interactive operations such as negotiated lane changes, roundabout circulation or multi-agent gap selection in heterogeneous traffic streams [46, 55]. Heavy reliance on manually crafted production rules further constrains behavioural coverage: encoding human cognition by hand is labour-intensive, difficult to scale and ill-suited to novel situations that deviate from the rule base [97]. Even when consistency is

achieved across several CARLA towns, overall performance still lags behind on finer-grained tasks, underscoring poor manoeuvre generalisation [49]. The experimental landscape is likewise limited. Many papers validate only one or two scenario templates—for instances, T-intersections, or five-lane highways—without pedestrians, irregular geometries or conflicting traffic patterns [35, 52, 65, 69, 93, 103, 129]. Training often depends on CARLA’s rule-abiding autopilot or a small, behaviourally homogeneous actor set, depriving agents of exposure to jaywalking pedestrians, reckless drivers or dense multi-class flows [64, 122]. Although incremental improvements are reported when additional episodes or temporal abstractions are introduced, authors themselves acknowledge the need for vastly richer scenario catalogues, higher-fidelity multi-agent interactions and more powerful high-level decision modules before reliable coverage of urban manoeuvre diversity can be claimed [30, 38, 44, 50, 58]. Addressing these shortcomings will require large-scale, heterogeneous simulation suites that blend rule-breaking agents with stochastic pedestrians, formal integration of lateral and longitudinal decision layers, and benchmark protocols that span the full gamut of urban driving behaviours—from complex merges to cooperative gap-sharing—in order to deliver policies that generalise beyond narrowly scripted tasks.

## 11 conclusion

We have presented a comprehensive survey of RL research in the CARLA simulator, covering nearly one hundred peer-reviewed studies spanning algorithms, state and action spaces, reward design, terminal conditions, and evaluation metrics. By organizing these contributions along common axes, we have provided newcomers with a clear map of the CARLA landscape and offered seasoned researchers a critical synthesis of where the field stands today.

Three prominent trends emerge from our review. First, CARLA-based work remains dominated by model-free approaches—over 80% of studies rely on variants of PPO, SAC, DQN, DDPG, or TD3—while model-based and hybrid methods are comparatively under-explored. Second, although bird’s-eye-view and semantically segmented images afford richer spatial structure, most works still employ kinematic-based or front-camera RGB state representations. Third, while a handful of standardized metrics (e.g., collision rate, lane deviation, success rate) recur across studies, many papers introduce bespoke evaluation measures tailored to specific tasks or scenarios, which can hinder direct comparison.

Despite these advances, several critical challenges remain. Sample efficiency continues to be a bottleneck, as state-of-the-art algorithms demand extensive trial-and-error even in simulation, limiting their real-world applicability. Insufficient real-world validation has left most CARLA-trained policies untested outside the simulator, obscuring their performance under true driving conditions. Generalization across towns, weather conditions, and sensor modalities remains elusive, and the sim-to-real gap persists despite early efforts in domain randomization. Safety guarantees and risk-sensitive decision making are seldom integrated into core learning pipelines, while sparse or poorly shaped reward functions further destabilize training. Perception studies often rely on a single modality—typically RGB imagery—neglecting the complementary

strengths of LiDAR, RADAR, and GPS. Finally, traffic-rule compliance and social-norm adherence are rarely embedded in learned policies, and most evaluations use simplified single-agent scenarios that fail to capture the complexity of real urban environments.

Addressing these gaps will require concerted efforts on multiple fronts. First, model-based and hybrid approaches—where learned dynamics models guide planning—hold promise for dramatically improving sample efficiency and accelerating convergence. Second, integrating formal safety constraints, uncertainty estimation, and risk-aware objectives into RL algorithms will be essential for certifiable autonomy. Third, a push toward standardized benchmarks, including multi-town, multi-weather scenario suites with unified evaluation metrics, will enable fair comparisons and reproducible progress. Finally, bridging the sim-to-real divide will necessitate richer domain-adaptation techniques, adversarial curricula, and real-vehicle pilot deployments to uncover and mitigate failure modes that only emerge outside the simulator.

## **Declarations**

### **Funding**

This research received no external funding.

### **Conflicts of interest/Competing interests**

The authors declare that they have no conflicts of interest or competing interests.

### **Availability of data and material (data transparency)**

Not applicable.

### **Code availability (software application or custom code)**

Not applicable.

### **Authors' contributions**

Elahe Delavari and Feeza Khan Khanzada contributed equally to this work. Both authors read and extracted information from the reviewed papers. Elahe Delavari prepared some of the tables and carried out the final edits and modifications for submission. Prof. Jaerock Kwon verified the correctness of the information and suggested revisions. All authors read and approved the final manuscript.

### **Acknowledgements**

Not applicable.

## References

- [1] National Highway Traffic Safety Administration: Critical reasons for crashes investigated in the national motor vehicle crash causation survey. Technical Report DOT HS 812 115, U.S. Department of Transportation (2015). Accessed: 2025-04-09. <https://crashstats.nhtsa.dot.gov/Api/Public/ViewPublication/812115>
- [2] Sallab, A.E., Abdou, M., Perot, E., Yogamani, S.: Deep reinforcement learning framework for autonomous driving. *Electronic Imaging* **29**(19), 70–76 (2017) <https://doi.org/10.2352/issn.2470-1173.2017.19.avm-023>
- [3] Dosovitskiy, A., Ros, G., Codevilla, F., Lopez, A., Koltun, V.: CARLA: An Open Urban Driving Simulator. In: *Proceedings of the 1st Annual Conference on Robot Learning*, pp. 1–16. PMLR, ??? (2017). ISSN: 2640-3498. <https://proceedings.mlr.press/v78/dosovitskiy17a.html> Accessed 2024-07-11
- [4] Wu, Y., Yin, Z., Yu, J., Zhang, M.: Lane Change Decision-Making through Deep Reinforcement Learning with Driver’s Inputs. In: *2022 IEEE 7th International Conference on Intelligent Transportation Engineering (ICITE)*, pp. 314–319 (2022). <https://doi.org/10.1109/ICITE56321.2022.10101421> . <https://ieeexplore.ieee.org/document/10101421/?arnumber=10101421> Accessed 2025-01-10
- [5] Goel, A., Chauhan, S.: Adaptive Look-ahead distance for Pure Pursuit Controller with Deep Reinforcement Learning Techniques. In: *Advances in Robotics - 5th International Conference of The Robotics Society*, pp. 1–5. ACM, Kanpur India (2021). <https://doi.org/10.1145/3478586.3478600> . <https://dl.acm.org/doi/10.1145/3478586.3478600> Accessed 2025-01-10
- [6] Youssef, F., Houda, B.: Deep reinforcement learning with external control: self-driving car application. In: *Proceedings of the 4th International Conference on Smart City Applications*, pp. 1–7. ACM, Casablanca Morocco (2019). <https://doi.org/10.1145/3368756.3369038> . <https://dl.acm.org/doi/10.1145/3368756.3369038> Accessed 2025-01-10
- [7] Carton, F., Filliat, D., Rabarisoa, J., Pham, Q.C.: Using Semantic Information to Improve Generalization of Reinforcement Learning Policies for Autonomous Driving. In: *2021 IEEE Winter Conference on Applications of Computer Vision Workshops (WACVW)*, pp. 144–151. IEEE, Waikola, HI, USA (2021). <https://doi.org/10.1109/WACVW52041.2021.00020> . <https://ieeexplore.ieee.org/document/9407810/> Accessed 2025-01-10
- [8] Mohammed, S., Argun, A., Ascheid, G.: A Unified Approach to Autonomous Driving in a High-Fidelity Simulator Using Vision-Based Reinforcement Learning. In: *2024 IEEE/SICE International Symposium on System Integration (SII)*, pp. 1093–1098 (2024).

- <https://doi.org/10.1109/SII58957.2024.10417385> . ISSN: 2474-2325.  
<https://ieeexplore.ieee.org/document/10417385/?arnumber=10417385> Accessed 2025-01-10
- [9] Zhang, M., Nakamoto, Y.: Toward ensuring better learning performance in reinforcement learning. In: 2022 Tenth International Symposium on Computing and Networking Workshops (CANDARW), pp. 134–139 (2022). <https://doi.org/10.1109/CANDARW57323.2022.00037> . ISSN: 2832-1324. <https://ieeexplore.ieee.org/document/10062717/?arnumber=10062717> Accessed 2025-01-10
- [10] Yang, Z., Pei, X., Xu, J., Zhang, X., Xi, W.: Decision-making in Autonomous Driving by Reinforcement Learning Combined with Planning & control. In: 2022 6th CAA International Conference on Vehicular Control and Intelligence (CVCI), pp. 1–6 (2022). <https://doi.org/10.1109/CVCI56766.2022.9964691> . <https://ieeexplore.ieee.org/document/9964691/?arnumber=9964691> Accessed 2025-01-10
- [11] Wang, J., Sun, H., Zhu, C.: Vision-Based Autonomous Driving: A Hierarchical Reinforcement Learning Approach. *IEEE Transactions on Vehicular Technology* **72**(9), 11213–11226 (2023) <https://doi.org/10.1109/TVT.2023.3266940> . Conference Name: IEEE Transactions on Vehicular Technology. Accessed 2025-01-10
- [12] Wei, D., Xing, J., Yang, S., Lu, Y., Huang, Y.: Continual Reinforcement Learning for Autonomous Driving with Application on Velocity Control under Various Environment. In: 2023 7th CAA International Conference on Vehicular Control and Intelligence (CVCI), pp. 1–8 (2023). <https://doi.org/10.1109/CVCI59596.2023.10397147> . <https://ieeexplore.ieee.org/document/10397147/?arnumber=10397147> Accessed 2025-01-10
- [13] Zhu, Z., Zhang, S., Zhuang, Y., Liu, Y., Liu, M., Gong, Z., Kai, S., Gu, Q., Wang, B., Cheng, S., Wang, X., Hao, J., Yu, Y.: RITA: Boost Driving Simulators with Realistic Interactive Traffic Flow. In: Proceedings of the Fifth International Conference on Distributed Artificial Intelligence. DAI '23, pp. 1–10. Association for Computing Machinery, New York, NY, USA (2023). <https://doi.org/10.1145/3627676.3627681> . <https://doi.org/10.1145/3627676.3627681> Accessed 2025-01-10
- [14] Li, Q., Jia, X., Wang, S., Yan, J.: Think2Drive: Efficient Reinforcement Learning by Thinking in Latent World Model for Quasi-Realistic Autonomous Driving (in CARLA-v2). arXiv. Comment: Accepted by ECCV 2024 (2024). <https://doi.org/10.48550/arXiv.2402.16720>
- [15] Pan, M., Zhu, X., Zheng, Y., Wang, Y., Yang, X.: Model-Based Reinforcement Learning with Isolated Imaginations. arXiv. Comment: arXiv admin note:

text overlap with arXiv:2205.13817 (2023). <https://doi.org/10.48550/arXiv.2303.14889>

- [16] Cai, P., Wang, H., Huang, H., Liu, Y., Liu, M.: Vision-Based Autonomous Car Racing Using Deep Imitative Reinforcement Learning. arXiv. Comment: 8 pages, 8 figures. IEEE Robotics and Automation Letters (RA-L) & IROS 2021 (2021). <https://doi.org/10.48550/arXiv.2107.08325>
- [17] Chen, D., Koltun, V., Krähenbühl, P.: Learning to Drive from a World on Rails. arXiv. Comment: Paper published in ICCV 2021(Oral); Code and data available at: [https://dotchen.github.io/world\\_on\\_rails/](https://dotchen.github.io/world_on_rails/) (2021). <https://doi.org/10.48550/arXiv.2105.00636>
- [18] Wang, C., Zhou, S., Wang, L., Lu, Z., Wu, C., Wen, X., Shou, G.: Autonomous Driving via Knowledge-Enhanced Safe Reinforcement Learning. IEEE Transactions on Intelligent Vehicles, 1–14 (2024) <https://doi.org/10.1109/TIV.2024.3412916>
- [19] Huang, C., Zhang, R., Ouyang, M., Wei, P., Lin, J., Su, J., Lin, L.: Deductive Reinforcement Learning for Visual Autonomous Urban Driving Navigation. IEEE Transactions on Neural Networks and Learning Systems **32**(12), 5379–5391 (2021) <https://doi.org/10.1109/TNNLS.2021.3109284>
- [20] Yurtsever, E., Capito, L., Redmill, K., Ozguner, U.: Integrating Deep Reinforcement Learning with Model-based Path Planners for Automated Driving. arXiv. Comment: 6 pages, 5 figures. Accepted for IEEE Intelligent Vehicles Symposium 2020 (2020). <https://doi.org/10.48550/arXiv.2002.00434>
- [21] Quigley, M., Conley, K., Gerkey, B., Faust, J., Foote, T., Leibs, J., Wheeler, R., Ng, A.Y., *et al.*: Ros: an open-source robot operating system. In: ICRA Workshop on Open Source Software, vol. 3, p. 5 (2009). Kobe
- [22] Autoware Foundation: Autoware: Open-source software for autonomous driving. <https://github.com/autowarefoundation/autoware>. Accessed: June 2025
- [23] Kaur, P., Taghavi, S., Tian, Z., Shi, W.: A Survey on Simulators for Testing Self-Driving Cars (2021). <https://arxiv.org/abs/2101.05337>
- [24] Rong, G., Shin, B.H., Tabatabaee, H., Lu, Q., Lemke, S., Možeiko, M., Boise, E., Uhm, G., Gerow, M., Mehta, S., Agafonov, E., Kim, T.H., Sterner, E., Ushiroda, K., Reyes, M., Zelenkovsky, D., Kim, S.: LGSVL Simulator: A High Fidelity Simulator for Autonomous Driving (2020). <https://arxiv.org/abs/2005.03778>
- [25] Chronis, C., Sardianos, C., Varlamis, I., Michail, D., Tserpes, K.: A driving profile recommender system for autonomous driving using sensor data and reinforcement learning. In: 25th Pan-Hellenic Conference on Informatics, pp. 33–38. ACM, Volos Greece (2021). <https://doi.org/10.1145/3503823.3503830>

<https://dl.acm.org/doi/10.1145/3503823.3503830> Accessed 2025-01-10

- [26] Udatha, S., Lyu, Y., Dolan, J.: Reinforcement learning with probabilistically safe control barrier functions for ramp merging. In: 2023 IEEE International Conference on Robotics and Automation (ICRA), pp. 5625–5630. <https://doi.org/10.1109/ICRA48891.2023.10161418> . <https://ieeexplore.ieee.org/document/10161418> Accessed 2025-05-06
- [27] Lillicrap, T.P., Hunt, J.J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., Wierstra, D.: Continuous control with deep reinforcement learning. arXiv:1509.02971 [cs] (2019). <https://doi.org/10.48550/arXiv.1509.02971> . <http://arxiv.org/abs/1509.02971> Accessed 2025-03-27
- [28] Silver, D., Lever, G., Heess, N., Degris, T., Wierstra, D., Riedmiller, M.: Deterministic Policy Gradient Algorithms. In: Proceedings of the 31st International Conference on Machine Learning, pp. 387–395. PMLR, ??? (2014). ISSN: 1938-7228. <https://proceedings.mlr.press/v32/silver14.html> Accessed 2025-03-27
- [29] Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A.A., Veness, J., Bellemare, M.G., Graves, A., Riedmiller, M., Fidjeland, A.K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., Hassabis, D.: Human-level control through deep reinforcement learning. *Nature* **518**(7540), 529–533 (2015) <https://doi.org/10.1038/nature14236> . Publisher: Nature Publishing Group. Accessed 2025-03-27
- [30] Peng, M., Gong, Z., Sun, C., Chen, L., Cao, D.: Imitative reinforcement learning fusing vision and pure pursuit for self-driving. In: 2020 IEEE International Conference on Robotics and Automation (ICRA), pp. 3298–3304 (2020). IEEE
- [31] Pérez-Gil, Ó., Barea, R., López-Guillén, E., Bergasa, L.M., Gómez-Huélamo, C., Gutiérrez, R., Díaz-Díaz, A.: Deep reinforcement learning based control for Autonomous Vehicles in CARLA. *Multimedia Tools and Applications* **81**(3), 3553–3576 (2022) <https://doi.org/10.1007/s11042-021-11437-3>
- [32] Fu, W., Li, Y., Ye, Z., Liu, Q.: Decision Making for Autonomous Driving Via Multimodal Transformer and Deep Reinforcement Learning. In: 2022 IEEE International Conference on Real-time Computing and Robotics (RCAR), pp. 481–486 (2022). <https://doi.org/10.1109/RCAR54675.2022.9872180>
- [33] Chen, M., Li, Y., Liu, Q., Lv, S., Xu, Y., Liu, Y.: Towards Autonomous Driving Decision by Combining Self-attention and Deep Reinforcement Learning. In: 2021 IEEE International Conference on Real-time Computing and Robotics (RCAR), pp. 1110–1115 (2021). <https://doi.org/10.1109/RCAR52367.2021.9517610>
- [34] Zhang, B., Xu, C., Su, Y., Xu, J.: A self-learning lane keeping algorithm. In: ICMLCA 2021; 2nd International Conference on Machine Learning and

Computer Application, pp. 1–7 (2021)

- [35] Tsai, J., Chang, Y.-T., Chuang, P.-H., You, Z.: An Autonomous Vehicle-Following Technique for Self-Driving Cars Based on the Semantic Segmentation Technique. In: 2023 IEEE International Symposium on Robotic and Sensors Environments (ROSE), pp. 1–7 (2023). <https://doi.org/10.1109/ROSE60297.2023.10410810>
- [36] Doe, D.M., Chen, D., Han, K., Wang, H., Xie, J., Han, Z.: DSORL: Data source optimization with reinforcement learning scheme for vehicular named data networks **24**(10), 11225–11237 <https://doi.org/10.1109/TITS.2023.3292033> . Accessed 2025-05-06
- [37] Li, L., Jiang, W., Shi, M., Wu, T.: Dynamic target following control for autonomous vehicles with deep reinforcement learning. In: 2022 International Conference on Advanced Robotics and Mechatronics (ICARM), pp. 386–391. <https://doi.org/10.1109/ICARM54641.2022.9959167> . <https://ieeexplore.ieee.org/document/9959167> Accessed 2025-05-06
- [38] Ahmed, M., Abobakr, A., Lim, C.P., Nahavandi, S.: Policy-based reinforcement learning for training autonomous driving agents in urban areas with affordance learning **23**(8), 12562–12571 <https://doi.org/10.1109/TITS.2021.3115235> . Accessed 2025-05-06
- [39] Cheng, S., Wang, N., Chen, F., Pipe, T.: Longitudinal Driving Skills Transfer from Driver to Smart Vehicle. In: 2021 26th International Conference on Automation and Computing (ICAC), pp. 1–6 (2021). <https://doi.org/10.23919/ICAC50006.2021.9594177> . <https://ieeexplore.ieee.org/document/9594177/?arnumber=9594177> Accessed 2025-01-10
- [40] Schulman, J., Wolski, F., Dhariwal, P., Radford, A., Klimov, O.: Proximal Policy Optimization Algorithms. arXiv. arXiv:1707.06347 [cs] (2017). <https://doi.org/10.48550/arXiv.1707.06347> . <http://arxiv.org/abs/1707.06347> Accessed 2025-03-27
- [41] Deng, Q., Zhao, Y., Li, R., Hu, Q., Liu, T., Li, R.: Context - Enhanced Meta-Reinforcement Learning with Data-Reused Adaptation for Urban Autonomous Driving. In: 2023 International Joint Conference on Neural Networks (IJCNN), pp. 1–8 (2023). <https://doi.org/10.1109/IJCNN54540.2023.10191187> . ISSN: 2161-4407. <https://ieeexplore.ieee.org/document/10191187/?arnumber=10191187> Accessed 2025-01-10
- [42] Deng, Q., Li, R., Hu, Q., Zhao, Y., Li, R.: Context-Aware Meta-RL With Two-Stage Constrained Adaptation for Urban Driving. IEEE Transactions on Vehicular Technology **73**(2), 1567–1581 (2024) <https://doi.org/10.1109/TVT.>



- [43] Zhao, J., Zhao, Y., Li, W., Zeng, C.: End-to-End Autonomous Driving Algorithm Based on PPO and Its Implementation. In: 2024 IEEE 13th Data Driven Control and Learning Systems Conference (DDCLS), pp. 1852–1858 (2024). <https://doi.org/10.1109/DDCLS61622.2024.10606596> . ISSN: 2767-9861. <https://ieeexplore.ieee.org/document/10606596/?arnumber=10606596> Accessed 2025-01-10
- [44] Wu, Y., Yuan, X.: Proximal Policy Optimization-based Reinforcement Learning for End-to-end Autonomous Driving. In: 2023 38th Youth Academic Annual Conference of Chinese Association of Automation (YAC), pp. 844–849 (2023). <https://doi.org/10.1109/YAC59482.2023.10401381> . ISSN: 2837-8601. <https://ieeexplore.ieee.org/document/10401381/?arnumber=10401381> Accessed 2025-01-10
- [45] Zhang, Z., Liniger, A., Dai, D., Yu, F., Gool, L.V.: End-to-End Urban Driving by Imitating a Reinforcement Learning Coach. arXiv. Comment: Published at ICCV 2021 (2021). <https://doi.org/10.48550/arXiv.2108.08265>
- [46] Agarwal, T., Arora, H., Schneider, J.: Learning Urban Driving Policies using Deep Reinforcement Learning. In: 2021 IEEE International Intelligent Transportation Systems Conference (ITSC), pp. 607–614 (2021). <https://doi.org/10.1109/ITSC48978.2021.9564412>
- [47] Trumpp, R., Büchner, M., Valada, A., Caccamo, M.: Efficient Learning of Urban Driving Policies Using Bird’s-Eye-View State Representations. arXiv. Comment: IEEE International Conference on Intelligent Transportation Systems 2023 (2023). <https://doi.org/10.48550/arXiv.2305.19904>
- [48] Xing, J., Nagata, T., Chen, K., Zou, X., Neftci, E., Krichmar, J.L.: Domain Adaptation In Reinforcement Learning Via Latent Unified State Representation. arXiv. Comment: Accepted by AAAI 2021; Fixed a typo in equation 3 (2021). <https://doi.org/10.48550/arXiv.2102.05714>
- [49] Anzalone, L., Barra, S., Nappi, M.: Reinforced Curriculum Learning For Autonomous Driving In Carla. In: 2021 IEEE International Conference on Image Processing (ICIP), pp. 3318–3322 (2021). <https://doi.org/10.1109/ICIP42928.2021.9506673>
- [50] Silva, V.A.S., Grassi, V.: Addressing Lane Keeping and Intersections using Deep Conditional Reinforcement Learning. In: 2021 Latin American Robotics Symposium (LARS), 2021 Brazilian Symposium on Robotics (SBR), and 2021 Workshop on Robotics in Education (WRE), pp. 330–335 (2021). <https://doi.org/10.1109/LARS/SBR/WRE54079.2021.9605436>

- [51] Albilani, M., Bouzeghoub, A.: Guided hierarchical reinforcement learning for safe urban driving. In: 2023 IEEE 35th International Conference on Tools with Artificial Intelligence (ICTAI), pp. 746–753. <https://doi.org/10.1109/ICTAI59109.2023.00115> . ISSN: 2375-0197. <https://ieeexplore.ieee.org/document/10356414> Accessed 2025-05-06
- [52] Martínez Gómez, L.M., Daza, I.G., Sotelo Vázquez, M.Á.: Temporal based deep reinforcement learning for crowded lane merging maneuvers. In: 2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC), pp. 2764–2769. <https://doi.org/10.1109/ITSC57777.2023.10422486> . ISSN: 2153-0017. <https://ieeexplore.ieee.org/document/10422486> Accessed 2025-05-06
- [53] Jin, Y.-L., Ji, Z.-Y., Zeng, D., Zhang, X.-P.: VWP:an efficient DRL-based autonomous driving model **26**, 2096–2108 <https://doi.org/10.1109/TMM.2022.3177942> . Accessed 2025-05-06
- [54] Shi, J., Zhang, T., Zhan, J., Chen, S., Xin, J., Zheng, N.: Efficient lane-changing behavior planning via reinforcement learning with imitation learning initialization. In: 2023 IEEE Intelligent Vehicles Symposium (IV), pp. 1–8. <https://doi.org/10.1109/IV55152.2023.10186577> . ISSN: 2642-7214. <https://ieeexplore.ieee.org/document/10186577> Accessed 2025-05-06
- [55] Gutiérrez-Moreno, R., Barea, R., López-Guillén, E., Arango, F., Revenga, P., Bergasa, L.M.: Decision Making for Autonomous Driving Stack: Shortening the Gap from Simulation to Real-World Implementations. In: 2024 IEEE Intelligent Vehicles Symposium (IV), pp. 3107–3113 (2024). <https://doi.org/10.1109/IV55156.2024.10588560>
- [56] Bai, Y., Du, J., Zhang, Y., Huang, Y.: FEN-DQN: An End-to-End Autonomous Driving Framework Based on Reinforcement Learning with Explicit Affordance. In: 2023 7th CAA International Conference on Vehicular Control and Intelligence (CVCI), pp. 1–6 (2023). <https://doi.org/10.1109/CVCI59596.2023.10397338> . <https://ieeexplore.ieee.org/document/10397338/?arnumber=10397338> Accessed 2025-01-10
- [57] Muhammed, A., Essam, H., Alber, B., Samuel, K., Muhammed, H., Wagdy, M., Khaled, N., Fawzy, H., Tarek, A., AbdelSalam, M., El-Kharashi, M.W.: Developing AI Agent with Functional Mockup Units for Car Autonomous Navigation. In: 2021 28th IEEE International Conference on Electronics, Circuits, and Systems (ICECS), pp. 1–5 (2021). <https://doi.org/10.1109/ICECS53924.2021.9665639> . <https://ieeexplore.ieee.org/document/9665639/?arnumber=9665639> Accessed 2025-01-10
- [58] Elallid, B.B., Benamar, N., Mrani, N., Rachidi, T.: Dqn-based reinforcement learning for vehicle control of autonomous vehicles interacting with pedestrians. In: 2022 International Conference on Innovation and Intelligence for Informatics, Computing, and Technologies (3ICT), pp. 489–493 (2022). IEEE

- [59] Chekroun, R., Toromanoff, M., Hornauer, S., Moutarde, F.: GRI: General Reinforced Imitation and Its Application to Vision-Based Autonomous Driving. arXiv (2022). <https://doi.org/10.48550/arXiv.2111.08575>
- [60] Toromanoff, M., Wirbel, E., Moutarde, F.: End-to-End Model-Free Reinforcement Learning for Urban Driving Using Implicit Affordances. arXiv. Comment: Accepted at main conference of CVPR 2020 (2020). <https://doi.org/10.48550/arXiv.1911.10868>
- [61] Marouane, C., Saad, B.: Safe Navigation Based on Deep Q-Network Algorithm Using an Improved Control Architecture. In: 2024 2nd International Conference on Electrical Engineering and Automatic Control (ICEEAC), pp. 1–6 (2024). <https://doi.org/10.1109/ICEEAC61226.2024.10576248>
- [62] Muhtadin, Meliaz, M.R., Dikairono, R., Purnama, I.K.E., Purnomo, M.H.: Deep Reinforcement Learning Control Strategy at Roundabout for i-CAR Autonomous Car. In: 2023 International Seminar on Intelligent Technology and Its Applications (ISITIA), pp. 473–478 (2023). <https://doi.org/10.1109/ISITIA59021.2023.10221077>
- [63] Elallid, B.B., Bagaa, M., Benamar, N., Mrani, N.: A Reinforcement Learning Based Approach for Controlling Autonomous Vehicles in Complex Scenarios. In: 2023 International Wireless Communications and Mobile Computing (IWCMC), pp. 1358–1364 (2023). <https://doi.org/10.1109/IWCMC58020.2023.10182377>
- [64] May, J., Poudel, S., Hamdan, S., Poudel, K., Vargas, J.: Using the CARLA Simulator to Train A Deep Q Self-Driving Car to Control a Real-World Counterpart on A College Campus
- [65] Ahmed, M., Lim, C.P., Nahavandi, S.: A Deep Q-Network Reinforcement Learning-Based Model for Autonomous Driving. In: 2021 IEEE International Conference on Systems, Man, and Cybernetics (SMC), pp. 739–744 (2021). <https://doi.org/10.1109/SMC52423.2021.9658892>
- [66] Clemmons, J., Jin, Y.-F.: Reinforcement Learning-Based Guidance of Autonomous Vehicles. In: 2023 24th International Symposium on Quality Electronic Design (ISQED), pp. 1–6 (2023). <https://doi.org/10.1109/ISQED57927.2023.10129362>
- [67] Li, G., Lin, S., Li, S., Qu, X.: Learning Automated Driving in Complex Intersection Scenarios Based on Camera Sensors: A Deep Reinforcement Learning Approach. IEEE Sensors Journal **22**(5), 4687–4696 (2022) <https://doi.org/10.1109/JSEN.2022.3146307>
- [68] Liu, Z., Hu, J., Song, T., Huang, Z.: A Methodology Based on Deep Reinforcement Learning to Autonomous Driving with Double Q-Learning. In: 2021 7th International Conference on Computer and Communications (ICCC), pp.

1266–1271 (2021). <https://doi.org/10.1109/ICCC54389.2021.9674600>

- [69] Friji, H., Ghazzai, H., Besbes, H., Massoud, Y.: A DQN-Based Autonomous Car-Following Framework Using RGB-D Frames. In: 2020 IEEE Global Conference on Artificial Intelligence and Internet of Things (GCAIoT), pp. 1–6 (2020). <https://doi.org/10.1109/GCAIoT51063.2020.9345899>
- [70] Hishmeh, L., Awad, F.: Deer in The Headlights: Short Term Planning via Reinforcement Learning Algorithms for Autonomous Vehicles. In: 2020 11th International Conference on Information and Communication Systems (ICICS), pp. 255–260 (2020). <https://doi.org/10.1109/ICICS49469.2020.239561>
- [71] [2210.16567] DeFIX: Detecting and Fixing Failure Scenarios with Reinforcement Learning in Imitation Learning Based Autonomous Driving. <https://arxiv.org/abs/2210.16567>
- [72] Deshpande, N., Vaufreydaz, D., Spalanzani, A.: Navigation In Urban Environments Amongst Pedestrians Using Multi-Objective Deep Reinforcement Learning. arXiv. <https://doi.org/10.48550/arXiv.2110.05205> . <http://arxiv.org/abs/2110.05205> Accessed 2025-05-06
- [73] Explainable Artificial Intelligence (XAI) Approach for Reinforcement Learning Systems | Proceedings of the 39th ACM/SIGAPP Symposium on Applied Computing. <https://dl.acm.org/doi/10.1145/3605098.3635992> Accessed 2025-05-06
- [74] Haarnoja, T., Zhou, A., Abbeel, P., Levine, S.: Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor. arXiv. <https://doi.org/10.48550/arXiv.1801.01290> . <http://arxiv.org/abs/1801.01290> Accessed 2025-07-07
- [75] Pomerleau, D.A.: Efficient training of artificial neural networks for autonomous navigation. *Neural computation* **3**(1), 88–97 (1991)
- [76] Ho, J., Ermon, S.: Generative adversarial imitation learning. *Advances in neural information processing systems* **29** (2016)
- [77] Li, Y., Song, J., Ermon, S.: Infogail: Interpretable imitation learning from visual demonstrations. *Advances in neural information processing systems* **30** (2017)
- [78] Song, J., Ren, H., Sadigh, D., Ermon, S.: Multi-agent generative adversarial imitation learning. *Advances in neural information processing systems* **31** (2018)
- [79] Savari, M., Choe, Y.: Online Virtual Training in Soft Actor-Critic for Autonomous Driving. In: 2021 International Joint Conference on Neural Networks (IJCNN), pp. 1–8 (2021). <https://doi.org/10.1109/IJCNN52387.2021.9533791> . ISSN: 2161-4407. <https://ieeexplore.ieee.org/document/9533791/?arnumber=9533791> Accessed

2025-01-10

- [80] Wu, Y., Wang, L., Lu, X., Wu, Y., Zhang, H.: Reinforcement Learning-based Autonomous Parking with Expert Demonstrations. In: 2023 7th CAA International Conference on Vehicular Control and Intelligence (CVCI), pp. 1–6 (2023). <https://doi.org/10.1109/CVCI59596.2023.10397136>. <https://ieeexplore.ieee.org/document/10397136/?arnumber=10397136> Accessed 2025-01-10
- [81] Aghdasian, A.J., Ardakani, A.H., Aqabakee, K., Abdollahi, F.: Autonomous Driving using Residual Sensor Fusion and Deep Reinforcement Learning. In: 2023 11th RSI International Conference on Robotics and Mechatronics (ICRoM), pp. 265–270 (2023). <https://doi.org/10.1109/ICRoM60803.2023.10412516>. ISSN: 2572-6889. <https://ieeexplore.ieee.org/document/10412516/?arnumber=10412516> Accessed 2025-01-10
- [82] Huang, Y., Peng, P., Zhao, Y., Zhai, Y., Xu, H., Tian, Y.: Simoun: Synergizing Interactive Motion-appearance Understanding for Vision-based Reinforcement Learning. In: 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 176–185 (2023). <https://doi.org/10.1109/ICCV51070.2023.00023>. ISSN: 2380-7504. <https://ieeexplore.ieee.org/document/10377012/?arnumber=10377012> Accessed 2025-01-10
- [83] Kong, S.-H., Nahrendra, I.M.A., Paek, D.-H.: Enhanced Off-Policy Reinforcement Learning With Focused Experience Replay. *IEEE Access* **9**, 93152–93164 (2021) <https://doi.org/10.1109/ACCESS.2021.3085142>
- [84] Jo, S.-B., Kim, P.-S., Jeong, H.-Y.: An Integrated Reward Function of End-to-End Deep Reinforcement Learning for the Longitudinal and Lateral Control of Autonomous Vehicles. In: 2022 IEEE 95th Vehicular Technology Conference: (VTC2022-Spring), pp. 1–5 (2022). <https://doi.org/10.1109/VTC2022-Spring54318.2022.9860646>
- [85] Tong, Z., Ayala, A., Sandoval, E.B., Cruz, F.: Urban Autonomous Driving of Emergency Vehicles with Reinforcement Learning. In: 2023 IEEE Latin American Conference on Computational Intelligence (LA-CCI), pp. 1–6 (2023). <https://doi.org/10.1109/LA-CCI58595.2023.10409469>
- [86] Han, Y., Yilmaz, A.: Learning to Drive Using Sparse Imitation Reinforcement Learning. *arXiv* (2022). <https://doi.org/10.48550/arXiv.2205.12128>
- [87] Igoe, C., Pande, S., Venkatraman, S., Schneider, J.: Multi-Alpha Soft Actor-Critic: Overcoming Stochastic Biases in Maximum Entropy Reinforcement Learning. In: 2023 IEEE International Conference on Robotics and Automation (ICRA), pp. 7162–7168 (2023). <https://doi.org/10.1109/ICRA48891.2023>

- [88] Gupta, D., Klusch, M.: HyLEAR: Hybrid deep reinforcement learning and planning for safe and comfortable automated driving. In: 2023 IEEE Intelligent Vehicles Symposium (IV), pp. 1–8. <https://doi.org/10.1109/IV55152.2023.10186781> . ISSN: 2642-7214. <https://ieeexplore.ieee.org/document/10186781> Accessed 2025-05-06
- [89] Fujimoto, S., Hoof, H., Meger, D.: Addressing function approximation error in actor-critic methods. In: Proceedings of the 35th International Conference on Machine Learning, pp. 1587–1596. PMLR. ISSN: 2640-3498. <https://proceedings.mlr.press/v80/fujimoto18a.html> Accessed 2025-07-07
- [90] Jia, R., Lu, S., Xie, W., Dai, L., Wu, Q., Zhang, J.: Safe Reinforcement Learning for Autonomous Lane-Changing: Considering the Collision Time and Lane Boundary Constraints. In: 2023 China Automation Congress (CAC), pp. 3484–3489 (2023). <https://doi.org/10.1109/CAC59555.2023.10450529> . ISSN: 2688-0938. <https://ieeexplore.ieee.org/document/10450529/?arnumber=10450529> Accessed 2025-01-10
- [91] Xu, T., Wang, Q., Yang, Y., Meng, J., Wang, Y.: End-to-End Autonomous Driving Decision-Making Solution Based on Pri-TD3. In: 2023 5th International Conference on Robotics, Intelligent Control and Artificial Intelligence (RICAI), pp. 525–529 (2023). <https://doi.org/10.1109/RICAI60863.2023.10489418> . <https://ieeexplore.ieee.org/document/10489418/?arnumber=10489418> Accessed 2025-01-10
- [92] Liao, Y., Yu, G., Chen, P., Zhou, B., Li, H.: Lateral Motion Control for Obstacle Avoidance in Autonomous Driving Based on Deep Reinforcement Learning. In: 2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC), pp. 5229–5234 (2023). <https://doi.org/10.1109/ITSC57777.2023.10422057> . ISSN: 2153-0017. <https://ieeexplore.ieee.org/document/10422057/?arnumber=10422057> Accessed 2025-01-10
- [93] Elallid, B.B., El Alaoui, H., Benamar, N.: Deep reinforcement learning for autonomous vehicle intersection navigation. In: 2023 International Conference on Innovation and Intelligence for Informatics, Computing, and Technologies (3ICT), pp. 308–313 (2023). IEEE
- [94] Liu, Y., Zhang, Q., Gao, Y., Zhao, D.: Deep-Reinforcement-Learning-Based Driving Policy at Intersections Utilizing Lane Graph Networks. IEEE Transactions on Cognitive and Developmental Systems **16**(5), 1759–1774 (2024) <https://doi.org/10.1109/TCDS.2024.3384269>
- [95] Li, S., Peng, K., Hui, F., Li, Z., Wei, C., Wang, W.: A decision-making approach for complex unsignalized intersection by deep reinforcement learning **73**(11),

- [96] Khalil, Y.H., Mouftah, H.T.: Exploiting Multi-Modal Fusion for Urban Autonomous Driving Using Latent Deep Reinforcement Learning. *IEEE Transactions on Vehicular Technology* **72**(3), 2921–2935 (2023) <https://doi.org/10.1109/TVT.2022.3217299> . Conference Name: IEEE Transactions on Vehicular Technology. Accessed 2025-01-10
- [97] Qi, H., Hou, E., Liu, G., Ye, P.: Cognitive Reinforcement Learning For Autonomous Driving. In: 2023 IEEE 3rd International Conference on Digital Twins and Parallel Intelligence (DTPI), pp. 1–5 (2023). <https://doi.org/10.1109/DTPI59677.2023.10365456>
- [98] Couto, G.C.K., Antonelo, E.A.: Hierarchical generative adversarial imitation learning with mid-level input generation for autonomous driving on urban environments. *IEEE Transactions on Intelligent Vehicles* (2024)
- [99] Jaafra, Y., Deruyver, A., Laurent, J.L., Naceur, M.S.: Context-Aware Autonomous Driving Using Meta-Reinforcement Learning. In: 2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA), pp. 450–455 (2019). <https://doi.org/10.1109/ICMLA.2019.00084>
- [100] Baheri, A., Kolmanovsky, I., Girard, A., Tseng, H.E., Filev, D.: Vision-based autonomous driving: A model learning approach. In: 2020 American Control Conference (ACC), pp. 2520–2525. <https://doi.org/10.23919/ACC45564.2020.9147510> . ISSN: 2378-5861. <https://ieeexplore.ieee.org/document/9147510> Accessed 2025-05-06
- [101] Wang, G., Niu, H., Zhu, D., Hu, J., Zhan, X., Zhou, G.: A Versatile and Efficient Reinforcement Learning Framework for Autonomous Driving. *arXiv*. Comment: 8 pages, 6 figures (2022). <https://doi.org/10.48550/arXiv.2110.11573>
- [102] Wang, J., Zhang, Q., Zhao, D.: Benchmarking Lane-changing Decision-making for Deep Reinforcement Learning. In: 2021 7th International Conference on Robotics and Artificial Intelligence, pp. 26–32. ACM, Guangzhou China (2021). <https://doi.org/10.1145/3505688.3505693> . <https://dl.acm.org/doi/10.1145/3505688.3505693> Accessed 2025-01-10
- [103] Cai, Y., Yang, S., Wang, H., Teng, C., Chen, L.: A Decision Control Method for Autonomous Driving Based on Multi-Task Reinforcement Learning. *IEEE Access* **9**, 154553–154562 (2021) <https://doi.org/10.1109/ACCESS.2021.3126796>
- [104] Xu, S.-Y., Chen, X.-M., Wang, Z.-J., Hu, Y.-H., Han, X.-T.: Decision-Making Models for Autonomous Vehicles at Unsignalized Intersections Based on Deep Reinforcement Learning. In: 2022 International Conference on Advanced Robotics and Mechatronics (ICARM), pp.



- 672–677 (2022). <https://doi.org/10.1109/ICARM54641.2022.9959664> .  
<https://ieeexplore.ieee.org/document/9959664/?arnumber=9959664> Accessed 2025-01-10
- [105] Manikandan, N.S., Kaliyaperumal, G., Wang, Y.: Ad Hoc-Obstacle Avoidance-Based Navigation System Using Deep Reinforcement Learning for Self-Driving Vehicles. *IEEE Access* **11**, 92285–92297 (2023) <https://doi.org/10.1109/ACCESS.2023.3297661> . Conference Name: IEEE Access. Accessed 2025-01-10
  - [106] Frauenknecht, B., Ehlgen, T., Trimpe, S.: Data-efficient Deep Reinforcement Learning for Vehicle Trajectory Control. *arXiv*. *arXiv:2311.18393 [cs]* (2023). <https://doi.org/10.48550/arXiv.2311.18393> . <http://arxiv.org/abs/2311.18393> Accessed 2025-01-10
  - [107] Sutton, R.S.: Dyna, an integrated architecture for learning, planning, and reacting **2**(4), 160–163 <https://doi.org/10.1145/122344.122377> . Accessed 2025-07-07
  - [108] Chua, K., Calandra, R., McAllister, R., Levine, S.: Deep Reinforcement Learning in a Handful of Trials using Probabilistic Dynamics Models. *arXiv*. <https://doi.org/10.48550/arXiv.1805.12114> . <http://arxiv.org/abs/1805.12114> Accessed 2025-07-07
  - [109] Janner, M., Fu, J., Zhang, M., Levine, S.: When to Trust Your Model: Model-Based Policy Optimization. *arXiv*. <https://doi.org/10.48550/arXiv.1906.08253> . <http://arxiv.org/abs/1906.08253> Accessed 2025-07-07
  - [110] Hafner, D., Lillicrap, T., Fischer, I., Villegas, R., Ha, D., Lee, H., Davidson, J.: Learning latent dynamics for planning from pixels. In: *Proceedings of the 36th International Conference on Machine Learning*, pp. 2555–2565. PMLR. ISSN: 2640-3498. <https://proceedings.mlr.press/v97/hafner19a.html> Accessed 2025-07-07
  - [111] Hafner, D., Pasukonis, J., Ba, J., Lillicrap, T.: Mastering Diverse Domains through World Models. *arXiv*. <https://doi.org/10.48550/arXiv.2301.04104> . <http://arxiv.org/abs/2301.04104> Accessed 2025-07-07
  - [112] Ghavamzadeh, M., Mannor, S., Pineau, J., Tamar, A.: Bayesian reinforcement learning: A survey. *Foundations and Trends in Machine Learning* **8**(5-6), 359–483 (2015)
  - [113] Gharaee, Z., Holmquist, K., He, L., Felsberg, M.: A Bayesian Approach to Reinforcement Learning of Vision-Based Vehicular Control. *arXiv* (2021). <https://doi.org/10.48550/arXiv.2104.03807>
  - [114] Sutton, R.S., Precup, D., Singh, S.: Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. In: *Artificial Intelligence*, vol. 112, pp. 181–211 (1999)



- [115] Li, Z.: A Hierarchical Autonomous Driving Framework Combining Reinforcement Learning and Imitation Learning. In: 2021 International Conference on Computer Engineering and Application (ICCEA), pp. 395–400 (2021). <https://doi.org/10.1109/ICCEA53728.2021.00084>
- [116] Wang, H., Lu, C., Gong, J.: Hierarchical Reinforcement Learning with Successor Representation for Intelligent Vehicle Collision Avoidance of Dynamic Pedestrian. In: 2024 9th International Conference on Computer and Communication Systems (ICCCS), pp. 575–581 (2024). <https://doi.org/10.1109/ICCCS61882.2024.10603229>
- [117] He, X., Yang, L., Lu, C., Li, Z., Gong, J.: Adaptive Decision Making at the Intersection for Autonomous Vehicles Based on Skill Discovery. arXiv. Comment: Accepted by IEEE ITSC 2022 (2022). <https://doi.org/10.48550/arXiv.2207.11724>
- [118] Naveed, K.B., Qiao, Z., Dolan, J.M.: Trajectory Planning for Autonomous Vehicles Using Hierarchical Reinforcement Learning. arXiv. Comment: 7 pages, 5 figures (2020). <https://doi.org/10.48550/arXiv.2011.04752>
- [119] Gangopadhyay, B., Soora, H., Dasgupta, P.: Hierarchical Program-Triggered Reinforcement Learning Agents For Automated Driving. IEEE Transactions on Intelligent Transportation Systems **23**(8), 10902–10911 (2022) <https://doi.org/10.1109/TITS.2021.3096998> [arXiv:2103.13861](https://arxiv.org/abs/2103.13861) [cs]. Comment: The paper is under consideration in Transactions on Intelligent Transportation Systems
- [120] Ziebart, B.D.: Modeling purposeful adaptive behavior with the principle of maximum causal entropy. PhD thesis, Carnegie Mellon University (2010)
- [121] Haarnoja, T., Tang, H., Abbeel, P., Levine, S.: Reinforcement learning with deep energy-based policies. In: Proceedings of the 34th International Conference on Machine Learning, pp. 1352–1361 (2017)
- [122] Khalil, Y.H., Mouftah, H.T.: Integration of Motion Prediction with End-to-end Latent RL for Self-Driving Vehicles. In: 2021 International Wireless Communications and Mobile Computing (IWCMC), pp. 1111–1116 (2021). <https://doi.org/10.1109/IWCMC51323.2021.9498581>
- [123] Bellemare, M.G., Dabney, W., Munos, R.: A distributional perspective on reinforcement learning. In: International Conference on Machine Learning, pp. 449–458 (2017)
- [124] Yang, Z., Zhang, X., Meng, D., Wen, Y., Zhang, B., Xu, Z.: Fully parameterized quantile function for distributional reinforcement learning. In: Advances in Neural Information Processing Systems (2019)
- [125] Chen, X., Yang, Y., Xu, S., Fu, S., Yang, D.: Motion Planning

- for Autonomous Vehicles in Uncertain Environments Using Hierarchical Distributional Reinforcement Learning. In: 2024 36th Chinese Control and Decision Conference (CCDC), pp. 1844–1851 (2024). <https://doi.org/10.1109/CCDC62350.2024.10587952> . ISSN: 1948-9447. <https://ieeexplore.ieee.org/document/10587952/?arnumber=10587952> Accessed 2025-01-10
- [126] Sharma, A., Jain, A., Finn, C., Levine, S.: Emergent latent skills in multi-task reinforcement learning. In: Conference on Robot Learning, pp. 973–991 (2020)
- [127] Li, Z., Nie, F., Sun, Q., Da, F., Zhao, H.: Boosting Offline Reinforcement Learning for Autonomous Driving with Hierarchical Latent Skills. arXiv (2023). <https://doi.org/10.48550/arXiv.2309.13614>
- [128] Hussonnois, M., Jun, J.-Y.: End-to-end autonomous driving using the apex algorithm in carla simulation environment. In: 2022 Thirteenth International Conference on Ubiquitous and Future Networks (ICUFN), pp. 18–23. <https://doi.org/10.1109/ICUFN55119.2022.9829674> . ISSN: 2165-8536. <https://ieeexplore.ieee.org/document/9829674> Accessed 2025-05-06
- [129] Zhao, X., Zhu, X., Su, Y.: A Real-time and Robust Visual-based Autonomous Driving Model. In: 2021 6th International Conference on Computational Intelligence and Applications (ICCIA), pp. 257–261 (2021). <https://doi.org/10.1109/ICCIA52886.2021.00057> . <https://ieeexplore.ieee.org/document/9644099/?arnumber=9644099> Accessed 2025-01-10
- [130] Wu, P., Chen, S., Metaxas, D.: MotionNet: Joint Perception and Motion Prediction for Autonomous Driving Based on Bird’s Eye View Maps. arXiv. <https://doi.org/10.48550/arXiv.2003.06754> . <http://arxiv.org/abs/2003.06754> Accessed 2025-07-07
- [131] Pérez-Gill, Ó., Barea, R., López-Guillén, E., Bergasa, L.M., Gómez-Huélamo, C., Gutiérrez, R., Díaz, A.: Deep Reinforcement Learning based control algorithms: Training and validation using the ROS Framework in CARLA Simulator for Self-Driving applications. In: 2021 IEEE Intelligent Vehicles Symposium (IV), pp. 1268–1273 (2021). <https://doi.org/10.1109/IV48863.2021.9575616>
- [132] Liang, X., Wang, T., Yang, L., Xing, E.: Cirl: Controllable imitative reinforcement learning for vision-based self-driving. In: Proceedings of the European Conference on Computer Vision (ECCV) (2018)
- [133] Sauer, A., Savinov, N., Geiger, A.: Conditional Affordance Learning for Driving in Urban Environments (2018). <https://arxiv.org/abs/1806.06498>
- [134] Codevilla, F., Santana, E., Lopez, A.M., Gaidon, A.: Exploring the limitations of behavior cloning for autonomous driving. In: Proceedings of the IEEE/CVF

