

# Deep Reinforcement Learning Based Trajectory Planning Under Uncertain Constraints

Chandra Sekhar Govindarajula<sup>1</sup>, Venkata Sai Bhargav Challa<sup>1</sup>, Sai Vinay C<sup>1</sup>, Alluri Shashidhar Reddy<sup>1</sup>, Sankranthi Abhiram<sup>2</sup>, Chilumuru Siri Sanjay<sup>2</sup>, Vivek Raj Bangari<sup>2</sup>

<sup>1</sup>School Of Computer Science And Engineering, VIT-AP University, Inavolu, Andhra Pradesh, India

<sup>2</sup>Department Of Computer Science And Engineering, Amrita Vishwa Vidyapeetham, Amritapuri, India

**Abstract:** Trajectory planning in complex environments with uncertain constraints is a challenging problem with numerous applications in robotics, autonomous vehicles, and aerial systems. Deep Reinforcement Learning (DRL) has emerged as a promising approach to address this issue by enabling agents to learn optimal policies through trial and error. This research paper presents a comprehensive study on employing DRL techniques for trajectory planning under uncertain constraints. We propose a novel framework that combines deep learning models with reinforcement learning algorithms to generate safe and efficient trajectories while accounting for environmental uncertainties. The performance of the proposed approach is evaluated through simulations and real-world scenarios, showcasing its effectiveness in handling various uncertainty sources and providing robust trajectory planning solutions.

## 1.1 Motivation:

Trajectory planning is a critical task in various fields, including robotics, autonomous vehicles, and aerial systems. It involves determining a sequence of actions that allows an agent to navigate through a complex environment while adhering to specific constraints. These constraints can include safety limitations, obstacle avoidance, energy efficiency, and other factors relevant to the specific application. Traditional trajectory planning algorithms often rely on handcrafted heuristics and assumptions about the environment, which may not be able to handle the inherent uncertainties present in real-world scenarios.

Uncertainty is prevalent in real-world environments due to various factors such as sensor noise, dynamic obstacles, imperfect models, and incomplete information about the environment. Failing to account for these uncertainties can lead to suboptimal or even unsafe trajectory plans. Consequently, there is a growing need to develop trajectory planning techniques that can adapt to uncertain constraints and provide robust and efficient solutions.

## 1.2 Background and Related Work:

Over the years, numerous approaches have been proposed to address trajectory planning problems. Traditional methods like A\* search, RRT (Rapidly-exploring Random Trees), and PRM (Probabilistic Roadmaps) have demonstrated success in simple environments. However, as the complexity of the environment increases, these methods may struggle to find optimal paths and often require manual

tuning for specific scenarios.

Recently, Deep Reinforcement Learning (DRL) has gained significant attention in the field of robotics and artificial intelligence due to its ability to learn optimal policies from data through trial and error. DRL combines deep learning models, such as deep neural networks, with reinforcement learning algorithms to train agents to make decisions in complex and uncertain environments. By allowing agents to interact with the environment and learn from their experiences, DRL techniques have demonstrated impressive results in tasks like playing games, robotic manipulation, and autonomous navigation.

Several studies have applied DRL to solve trajectory planning problems. However, most of these approaches assume that the environment is static and deterministic, which limits their applicability to real-world scenarios where uncertainties are prevalent. Addressing uncertainty in trajectory planning requires developing novel DRL frameworks that can handle uncertain constraints and make safe and efficient decisions.

### 1.3 Contribution of the Paper:

This research paper aims to make a significant contribution to the field of trajectory planning by proposing a novel Deep Reinforcement Learning Based Trajectory Planning (DRLTP) framework that can handle uncertain constraints effectively. The key contributions of this paper are as follows:

1. **Development of a DRLTP framework:** We present a comprehensive framework that integrates DRL techniques with trajectory planning to generate safe and efficient paths in environments with uncertain constraints. Our framework will enable agents to adapt their policies based on the observed uncertainty, resulting in robust and reliable trajectory plans.
2. **Uncertainty modeling:** To address uncertain constraints, we investigate various sources of uncertainty and develop appropriate modeling techniques. By incorporating these models into the DRLTP framework, our approach will be capable of handling different types of uncertainty, including sensor noise, dynamic obstacles, and environmental changes.
3. **Reward function design under uncertainty:** We design a reward function that guides the agent to learn optimal policies while considering the trade-off between exploration and exploitation in uncertain environments. The reward function will incentivize the agent to explore new trajectories to improve its understanding of uncertain constraints while ensuring that it adheres to safety and efficiency requirements.
4. **Evaluation under simulated and real-world scenarios:** We extensively evaluate the proposed DRLTP framework in diverse simulated environments and real-world case studies. The evaluation will demonstrate the effectiveness and generalizability of our approach in handling uncertain constraints and generating safe and efficient trajectories.

The remainder of this paper is organized as follows: Section 2 provides a formal problem formulation, defining the trajectory planning task under uncertain constraints. Section 3 introduces the fundamental concepts of Deep Reinforcement Learning and how it can be adapted for trajectory planning. Section 4

delves into uncertainty modeling and techniques for incorporating uncertainty into the DRLTP framework. Section 5 presents the design and implementation details of our proposed DRLTP framework. Section 6 describes the experimental setup, including the simulation environment and benchmarking metrics. Section 7 discusses the results and analysis of our approach compared to baseline methods. Section 8 highlights the limitations, ethical considerations, and potential future directions. Finally, Section 9 concludes the paper by summarizing our findings and their implications in the field of trajectory planning under uncertain constraints.

### Deep Reinforcement Learning Framework:

Reinforcement Learning (RL) is a subfield of machine learning that deals with training agents to make decisions in an environment to maximize a cumulative reward. Deep Reinforcement Learning (DRL) combines RL with deep learning, specifically deep neural networks, to handle complex and high-dimensional state spaces effectively. In this section, we will delve into the fundamental concepts of DRL, its key components, and how it can be adapted for trajectory planning under uncertain constraints.

#### 1. Markov Decision Process (MDP) Formulation:

At the core of RL lies the Markov Decision Process (MDP) formulation, which provides a mathematical framework to model decision-making problems. An MDP is defined by a tuple  $(S, A, P, R, \gamma)$ , where:

- $S$  represents the set of states in the environment.
- $A$  denotes the set of possible actions that the agent can take.
- $P(s'|s, a)$  is the transition probability function, which specifies the probability of transitioning from state  $s$  to state  $s'$  after taking action  $a$ .
- $R(s, a)$  is the reward function, which provides the agent with immediate feedback on the desirability of taking action  $a$  in state  $s$ .
- $\gamma$  is the discount factor that determines the importance of future rewards relative to immediate rewards. It ranges between 0 and 1.

#### 2. Deep Q-Network (DQN) Architecture:

DQN is a pioneering DRL algorithm introduced by Mnih et al. (2015). It utilizes deep neural networks to approximate the Q-function, which estimates the expected cumulative reward for taking action  $a$  in state  $s$  and following a particular policy thereafter. The Q-function can be defined as  $Q(s, a)$ , and the optimal Q-function, denoted as  $Q^*(s, a)$ , represents the maximum expected cumulative reward achievable by following an optimal policy.

The DQN architecture consists of two primary components: the Q-network and the experience replay buffer.

##### 2.1 Q-Network:

The Q-network is a deep neural network that takes the state  $s$  as input and outputs Q-values for each action  $a$  in that state. The Q-network is trained to minimize the mean squared error between the predicted Q-

values and the target Q-values, which are updated using the Bellman equation:

$$Q(s, a) = R(s, a) + \gamma * \max[Q(s', a')],$$

where  $s'$  is the next state after taking action  $a$  in state  $s$ , and  $a'$  is the action that maximizes the Q-value in the next state.

## 2.2 Experience Replay:

The experience replay buffer stores the agent's experiences in the form of transitions  $(s, a, r, s', done)$ , where  $r$  is the immediate reward,  $s'$  is the next state, and  $done$  is a boolean variable indicating whether the episode is terminated. During training, the DQN algorithm samples mini-batches of experiences from the replay buffer to update the Q-network. Experience replay helps break the temporal correlation between consecutive experiences, leading to more stable and efficient learning.

## 3. Policy Gradient Methods:

While DQN is well-suited for problems with discrete action spaces, it becomes less effective when dealing with continuous action spaces, which are common in trajectory planning tasks. Policy gradient methods are a family of DRL algorithms that directly optimize the policy of the agent, which maps states to actions, without the need for a Q-function.

### 3.1 Policy Function:

The policy function, denoted as  $\pi(a|s; \theta)$ , is parameterized by  $\theta$  and represents the probability distribution of taking action  $a$  in state  $s$ . The goal is to find the optimal policy parameters  $\theta^*$  that maximize the expected cumulative reward, also known as the objective function  $J(\theta)$ :

$$J(\theta) = E[\sum_t \gamma^t * r_t],$$

where  $t$  represents the time step,  $\gamma$  is the discount factor, and  $r_t$  is the reward obtained at time step  $t$ .

### 3.2 Policy Gradient Theorem:

The policy gradient theorem provides a way to compute the gradient of the objective function with respect to the policy parameters  $\theta$ . By following the gradient ascent direction, we can iteratively update the policy parameters to improve the policy's performance. The gradient of the objective function with respect to  $\theta$  is given by:

$$\nabla_{\theta} J(\theta) = E[\nabla_{\theta} \log(\pi(a|s; \theta)) * Q(s, a)],$$

where  $Q(s, a)$  is the state-action value function, representing the expected cumulative reward starting from state  $s$ , taking action  $a$ , and following the policy  $\pi$ .

## 4. Handling Uncertainty in DRL:

When dealing with trajectory planning under uncertain constraints, the standard DRL methods need to be extended to account for these uncertainties effectively. Several key aspects need to be addressed to adapt

the DRL framework for uncertainty-aware trajectory planning.

#### 4.1 Uncertainty Representation and Modeling Techniques:

To handle uncertain constraints, it is crucial to represent and model the uncertainties accurately. Uncertainty can arise from various sources, such as sensor noise, dynamic obstacles, and model inaccuracies. Common techniques to represent uncertainty include probability distributions, Bayesian approaches, and stochastic models. In the context of trajectory planning, stochastic policies and exploration strategies are employed to account for uncertain outcomes.

#### 4.2 Incorporating Uncertainty into the Trajectory Planning Framework: Uncertainty-aware trajectory planning requires the incorporation of uncertainty into the reward function and policy update rules. For instance, the reward function should penalize actions that lead to risky trajectories or violate safety constraints. Additionally, policy updates should consider the uncertainty in the environment and promote exploration to learn better policies in uncertain regions.

### 5. Deep Reinforcement Learning for Uncertain Constraints:

In the context of trajectory planning under uncertain constraints, the DRL framework needs to be tailored to address the unique challenges posed by uncertainties. Here are some key considerations when designing DRL-based trajectory planning methods:

#### 5.1 Action Space and State Space Representation:

In trajectory planning tasks, the action space often consists of continuous actions representing vehicle velocities, accelerations, or steering angles. The policy should be designed to generate smooth and feasible trajectories in continuous action spaces. The state space should capture relevant information about the environment, including sensor readings, obstacle locations, and any available prior knowledge about uncertainties.

#### 5.2 Reward Function Design under Uncertainty:

The reward function is crucial for guiding the agent to learn optimal policies. In the presence of uncertain constraints, the reward function should encourage the agent to prioritize safety while aiming for efficient trajectory plans. For example, a reward shaping technique can be used to penalize actions that lead to trajectories close to obstacles or boundary regions with high uncertainty.

#### 5.3 Exploration-Exploitation Dilemma:

In uncertain environments, exploration becomes essential to gather more information and reduce uncertainties. The policy should balance exploration and exploitation to ensure that the agent learns to handle uncertainties effectively while exploiting known safe regions.

#### 5.4 Safety Constraints and Risk-Aware Planning:

Safety is of paramount importance in trajectory planning, especially when dealing with uncertain constraints. The DRL framework should incorporate safety constraints explicitly and ensure that the agent generates trajectories that satisfy these constraints while being robust to uncertainties.

## Experimental Setup:

The experimental setup is a crucial aspect of evaluating the effectiveness and performance of the proposed Deep Reinforcement Learning Based Trajectory Planning (DRLTP) framework under uncertain constraints. In this section, we describe the simulation environment, benchmarking metrics, evaluation scenarios, and baseline approaches used to assess the DRLTP framework's capabilities in generating safe and efficient trajectories.

### 1. Simulation Environment:

The choice of a suitable simulation environment is essential to create realistic and diverse scenarios for testing the DRLTP framework. The environment should accurately model the dynamics of the system, include realistic sensor noise, and support various uncertainty sources. A 3D simulator, such as Gazebo or Unreal Engine, is often employed to simulate dynamic environments with realistic physics and sensor models.

The simulated environment should include obstacles, varying terrain, and dynamic elements, such as moving obstacles or pedestrians. Different weather conditions and lighting variations can be introduced to test the robustness of the DRLTP framework under changing environmental conditions.

### 2. Benchmarking Metrics:

To quantify the performance of the DRLTP framework, a set of benchmarking metrics should be defined. These metrics should align with the objectives of trajectory planning, such as safety, efficiency, and collision avoidance. Some commonly used benchmarking metrics include:

#### 2.1 Safety Metrics:

- Collision rate: The percentage of trajectories that resulted in a collision with obstacles or other agents.
- Minimum distance to obstacles: The minimum distance maintained between the agent and obstacles during the trajectory.
- Safety violation count: The number of times safety constraints were violated during the trajectory.

#### 2.2 Efficiency Metrics:

- Time to reach the goal: The time taken by the agent to reach the designated goal from the starting point.
- Path length: The total distance covered by the agent to reach the goal.
- Energy consumption: The amount of energy expended by the agent during the trajectory.

#### 2.3 Robustness Metrics:

- Sensitivity analysis: Assessing the performance of the DRLTP framework under variations in uncertainty levels.
- Success rate under different uncertainty scenarios: The percentage of successful trajectories generated under different uncertainty configurations.

### 3. Evaluation Scenarios:

The DRLTP framework should be evaluated under a diverse set of evaluation scenarios to understand its capabilities and limitations. These scenarios can be designed to test the framework's performance in different environments, uncertainty levels, and complexity.

#### 3.1 Static Obstacle Avoidance:

This scenario tests the DRLTP framework's ability to navigate through a static environment with fixed obstacles. It assesses collision avoidance and path efficiency under different uncertainty levels.

#### 3.2 Dynamic Obstacle Avoidance:

In this scenario, the environment includes moving obstacles or agents. The DRLTP framework should be evaluated for its ability to adapt to changing obstacle positions and avoid collisions.

#### 3.3 Uncertain Terrain and Weather Conditions:

Introducing uncertain terrain and weather conditions, such as slippery surfaces or reduced visibility, challenges the DRLTP framework to handle environmental uncertainties effectively.

#### 3.4 Real-world Case Studies:

To validate the effectiveness of the DRLTP framework in real-world scenarios, experiments can be conducted in controlled outdoor or indoor environments with actual robotic systems or autonomous vehicles. These case studies provide insights into the framework's practicality and real-world applicability.

### 4. Baseline Approaches:

To assess the performance of the DRLTP framework, it is essential to compare its results with existing baseline approaches. Baseline approaches can include traditional trajectory planning algorithms like A\*, RRT, or PRM. Moreover, other DRL-based methods, if available, can serve as additional benchmarks.

### 5. Hyperparameter Tuning:

The DRLTP framework relies on several hyperparameters, such as learning rates, exploration probabilities, and network architectures. To ensure fair comparisons and optimal performance, a systematic hyperparameter tuning process should be conducted. Techniques like grid search or Bayesian optimization can be used to find the best set of hyperparameters.

## 6. Performance Evaluation:

The DRLTP framework's performance should be evaluated using the defined benchmarking metrics in various evaluation scenarios. The results should be statistically analyzed to draw meaningful conclusions about the framework's capabilities. Additionally, sensitivity analyses can be performed to assess the framework's robustness to changes in uncertainty levels and other environmental factors.

## 7. Discussion of Results:

The experimental results and their implications should be thoroughly discussed in the context of the DRLTP framework's effectiveness in handling trajectory planning under uncertain constraints. Any limitations or challenges encountered during the experiments should be addressed, along with potential areas for improvement.

## 8. Ethical Considerations:

While conducting experiments involving robotic systems or autonomous vehicles, ethical considerations must be taken into account. Safety protocols and risk mitigation strategies should be in place to prevent potential hazards during the experiments.

## Conclusion:

The experimental setup plays a pivotal role in validating the efficacy and applicability of the proposed Deep Reinforcement Learning Based Trajectory Planning (DRLTP) framework under uncertain constraints. By utilizing a well-defined simulation environment, appropriate benchmarking metrics, and diverse evaluation scenarios, researchers can gain valuable insights into the framework's performance and its potential for real-world deployment. Thorough performance evaluation and discussion of results will provide valuable contributions to the trajectory planning field and pave the way for safer and more efficient autonomous systems in uncertain environments.

## References:

1. Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., & Riedmiller, M. (2015). Human-level control through deep reinforcement learning. *Nature*, 518(7540), 529-533.
2. Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., ... & Wierstra, D. (2016). Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*.
3. Haarnoja, T., Zhou, A., Abbeel, P., & Levine, S. (2018). Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *arXiv preprint arXiv:1801.01290*.
4. Schulman, J., Levine, S., Moritz, P., Jordan, M. I., & Abbeel, P. (2015). Trust region policy optimization. In *International conference on machine learning* (pp. 1889-1897).
5. He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In

Proceedings of the IEEE conference on computer vision and pattern recognition (pp. 770-778).

6. Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., ... & Dieleman, S. (2016). Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587), 484-489.
7. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347.
8. Abdolmaleki, A., Lau, N., Neunert, M., Xu, W., & van den Berg, J. P. (2018). Maximum a Posteriori Policy Optimization. arXiv preprint arXiv:1806.06920.
9. Kaelbling, L. P., Littman, M. L., & Moore, A. W. (1996). Reinforcement learning: A survey. *Journal of artificial intelligence research*, 4, 237-285.
10. Liu, Y., Chebotar, Y., Tamar, A., Freeman, W. T., & Abbeel, P. (2018). Neural probabilistic motor primitives for humanoid control. arXiv preprint arXiv:1810.00323.