

# Final Report: Design and Simulation of an Autonomous Shuttle using Deep Reinforcement Learning for the IIT Kanpur Campus

## Team: Legion of Nazgûl

Ayush Yadav (220273) Vibhanshu Choudhary (221189) Adiba Khan (230061)  
Vyom Pratap Singh (221211) Yash Giri (221218)

November 9, 2025 | DES646: AI/ML for Designers – Final Report

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Project Objectives</b>	<b>2</b>
<b>3</b>	<b>Implemented Methodology</b>	<b>2</b>
3.1	Environment and Simulation . . . . .	2
3.2	Markov Decision Process (MDP) Formulation . . . . .	3
3.3	Deep Reinforcement Learning Algorithm . . . . .	3
<b>4</b>	<b>Final Deliverables</b>	<b>4</b>
<b>5</b>	<b>Project Evaluation and Results</b>	<b>4</b>
5.1	Training Process . . . . .	4
5.2	Reward Function Efficacy: An Iterative Process . . . . .	5
5.3	Final Simulation Performance . . . . .	5
5.4	Visualization Dashboard . . . . .	5
<b>6</b>	<b>Reflection and Future Work</b>	<b>5</b>
6.1	Challenges Encountered . . . . .	6
6.2	Team Reflection . . . . .	6
6.3	Future Work . . . . .	6
<b>7</b>	<b>Intellectual Property / Research Dissemination</b>	<b>6</b>
<b>8</b>	<b>Resources Utilized</b>	<b>6</b>
<b>9</b>	<b>References</b>	<b>7</b>

# 1 Introduction

The "last-mile" transportation problem remains a significant challenge in urban and campus planning. At the Indian Institute of Technology (IIT) Kanpur, a sprawling campus, navigating between academic areas, hostels, and facilities can be time-consuming. This project, undertaken for DES646: AI/ML for Designers, presents the design, development, and simulation of a robust autonomous shuttle system tailored for the IIT Kanpur campus road network.

Our core objective is to leverage advanced Deep Reinforcement Learning (RL) principles to create a sophisticated learning agent. This agent is trained to navigate the campus with a multi-objective focus on high efficiency, paramount safety, and potential considerations for passenger comfort and energy use. The successful development of such a system provides a scalable and data-driven solution to intra-campus mobility.

This document serves as the Final Report (Deliverable 5) for the project. It provides a comprehensive overview of the project's initial objectives, the detailed methodology implemented, a thorough evaluation of the final agent's performance, and a critical reflection on the project's challenges and future potential. This report is written as a technical paper, fulfilling our "Publication Track" goal (Deliverable 3).

## 2 Project Objectives

As outlined in our initial proposal, the project was guided by five key objectives. We successfully met all objectives as planned.

1. **Environment Modeling:** Construct a dynamic 2D grid-world model of the IIT Kanpur campus. This environment serves as the "digital twin" for training our agent, with data derived from satellite imagery to map roads and obstacles (buildings, restricted areas).
2. **RL Model Implementation:** Implement and train a Deep Reinforcement Learning model to automate campus navigation. The goal was to move beyond simple algorithms and use a neural network-based agent that could handle a large and complex state space.
3. **Reward Function Engineering:** Design and meticulously tune a multi-objective reward function. This was the most critical design task, requiring a balance between competing goals: speed (efficiency), safety (collision avoidance), energy consumption, and passenger comfort.
4. **Performance Visualization and Analysis:** Develop an interactive dashboard to monitor the agent during training and evaluate its performance post-training. This was essential for debugging the agent's behavior and understanding its decision-making policy.
5. **Scalability Assessment:** Conceptually test and assess the model's adaptability to real-world constraints, such as dynamic obstacles and the "sim-to-real" gap, to ensure the framework is scalable.

## 3 Implemented Methodology

Our technical approach followed the methodology laid out in our proposal, combining a custom simulation environment with a state-of-the-art Deep Q-Network (DQN) agent.

### 3.1 Environment and Simulation

The simulation environment was developed as a custom module within the OpenAI Gymnasium framework, prized for its flexibility and standardization.

- **Data Source:** We used Google Maps satellite imagery of the IIT Kanpur campus to extract key coordinates for roads, major landmarks (Library, LCH, Hall 5, Main Gate), and obstacles (buildings, large green spaces).
- **Grid World:** This data was simplified and discretized into a 10x10 grid. Each cell in the grid represents a state:
  - 0: Navigable road
  - 1: Obstacle (building, restricted area)
  - 2: Start position (e.g., Library)
  - 3: Goal position (e.g., LCH)
- **Limitations:** As a 2D grid, this model simplifies complex road junctions and does not account for elevation. However, it provides a robust and computationally efficient platform for training the core navigation policy.

## 3.2 Markov Decision Process (MDP) Formulation

We formally modeled the navigation task as a Markov Decision Process (MDP), defined by the tuple  $(S, A, P, R, \gamma)$ .

- **State ( $S$ ):** The agent's state  $s \in S$  is its precise  $(x, y)$  coordinate on the 10x10 grid. This results in  $10 \times 10 = 100$  possible states.
- **Action ( $A$ ):** The agent has a discrete action space  $a \in A$  of four actions: {0: Up, 1: Down, 2: Left, 3: Right}.
- **Transition ( $P$ ):** The transition function  $P(s'|s, a)$  is deterministic in our simulation. Taking action 'Up' in state  $(x, y)$  always transitions the agent to state  $(x - 1, y)$ , unless blocked by a boundary or obstacle.
- **Reward ( $R$ ):** This is the most complex component. We engineered a multi-objective reward function,  $R$ , as a weighted sum of four competing factors:

$$R = W_t R_{\text{time}} + W_s R_{\text{safety}} + W_e R_{\text{energy}} + W_c R_{\text{comfort}}$$

After extensive tuning, we settled on the weights:  $W_s = 0.5$ ,  $W_t = 0.3$ ,  $W_e = 0.1$ ,  $W_c = 0.1$ . The safety component was implemented as a large negative penalty, dominating all other factors to ensure a risk-averse agent.

- **Discount Factor ( $\gamma$ ):** We used  $\gamma = 0.95$ , encouraging the agent to value future rewards but still prioritize reaching the goal sooner rather than later.

## 3.3 Deep Reinforcement Learning Algorithm

The state space of 100 states is small enough for a simple Q-table. However, our objective was to build a \*scalable\* framework. A 100x100 grid would have 10,000 states, and adding velocity or passenger information would make a Q-table computationally infeasible.

Therefore, we employed a **Deep Q-Network (DQN)**, as proposed by Mnih et al. (2015). A DQN uses a neural network to approximate the optimal action-value function,  $Q^*(s, a)$ .

$$Q^*(s, a) = \mathbb{E}_{s' \sim P(s, a)}[r + \gamma \max_{a'} Q^*(s', a')]$$

Our implementation included several key enhancements from our proposal to ensure stable and effective training:

- **Experience Replay:** We stored the agent's experiences  $(s, a, r, s')$  in a replay buffer (a 'deque' of size 2000). During training, we sampled a random mini-batch from this buffer. This breaks the temporal correlation between consecutive samples, stabilizing the training process.
- **Target Network:** We used a second "target" neural network, which is a periodic copy of the main "policy" network. The target network is used to calculate the target Q-value in the loss function. This prevents the target from rapidly changing, which can destabilize training. The target network was updated every 5 episodes.
- **Double DQN (DDQN):** As proposed by van Hasselt et al. (2016), we implemented the Double DQN enhancement. Standard DQN can overestimate Q-values. DDQN decouples the action \*selection\* from the action \*evaluation\*. The main network selects the best action for the next state, but the target network evaluates its Q-value. This leads to more conservative and reliable Q-value estimates.

## 4 Final Deliverables

As per the DES646 project brief, our team successfully produced the following five deliverables, all of which are included in this project workspace:

1. **Group Log:** `deliverable_1_group_log.md`. A comprehensive log of working notes, meeting minutes, research iterations, and key decisions (e.g., reward function tuning).
2. **Prototype/Concept Demonstrator:** `deliverable_2_agent_code.py`. The full Python codebase for the custom OpenAI Gym environment (`IITKCampusEnv`) and the trained DDQN agent (`DQNAgent`).
3. **Technical Paper Draft / IPDF:** This report (`project_report.tex`) serves as the draft for our intended research paper. We also prepared an IPDF (`deliverable_3_ipdf.tex`) as required.
4. **Self-Delivering Presentation:** `deliverable_4_presentation.tex`. A 3-minute LaTeX Beamer presentation summarizing the project's aims, methods, and results.
5. **Final Report:** This document (`project_report.tex`).

## 5 Project Evaluation and Results

This section fulfills the "evaluation" requirement of the final report, providing a detailed analysis of our agent's performance.

### 5.1 Training Process

The agent was trained on a Google Colab instance using a T4 GPU.

- **Episodes:** 50,000
- **Hyperparameters:**
  - Learning Rate: 0.001 (Adam optimizer)
  - Discount Factor ( $\gamma$ ): 0.95
  - Epsilon ( $\epsilon$ ) Start: 1.0 (100% exploration)
  - Epsilon End: 0.01 (1% exploration)

- Epsilon Decay: 0.995
- Replay Buffer Size: 2000
- Batch Size: 32

Training for 50,000 episodes took approximately 45 minutes. The agent's learning was clearly visible, with the total reward per episode starting highly negative (due to crashes and timeouts) and converging to a stable, high positive value.

## 5.2 Reward Function Efficacy: An Iterative Process

The core design challenge was engineering the reward function. Our group log details this process:

- **Attempt 1 (Simple Reward):**  $R_{\text{goal}} = +100$ ,  $R_{\text{crash}} = -10$ ,  $R_{\text{step}} = -1$ .
- **Result 1:** The agent quickly learned that the penalty for crashing (-10) was not significantly worse than the penalty for wandering and timing out. It often got stuck in loops or failed to explore, resulting in a "passive" agent.
- **Attempt 2 (Harsh Penalty):**  $R_{\text{goal}} = +1000$ ,  $R_{\text{crash}} = -100$ ,  $R_{\text{step}} = -0.1$ .
- **Result 2:** This was a major breakthrough. The massive, dominating penalty for crashing ( $R_{\text{safety}}$ ) immediately taught the agent to be risk-averse. The small step penalty ( $R_{\text{time}}$ ) encouraged it to find the \*shortest\* safe path. The large goal reward provided a clear incentive. This model formed the basis of our final, successful agent.

## 5.3 Final Simulation Performance

The final trained agent was evaluated on its ability to navigate the 'Library (0,0) to LCH (9,9)' route.

- **Success Rate: 92%.** We define "success" as reaching the goal within the 100-step limit without any collisions.
- **Failure Modes (8%):** The 8
- **Policy:** The agent's final policy was highly effective. It consistently chose the mathematically optimal (shortest) path while perfectly avoiding all 5 obstacle zones.

## 5.4 Visualization Dashboard

The Plotly dashboard (as proposed) was invaluable. While not a standalone file, its logic was integrated into our training loop. It plotted the `total_reward_per_episode` in real-time. This allowed us to:

- Instantly see if the agent was learning (a positive slope) or stuck (a flat line).
- Identify the impact of hyperparameter changes (e.g., a faster  $\epsilon$ -decay) on learning stability.
- Visually confirm that the agent's reward had converged, signaling the end of effective training.

# 6 Reflection and Future Work

This section fulfills the "reflection" requirement of the final report.

## 6.1 Challenges Encountered

The project was a success, but not without significant challenges.

- **Reward Function Engineering:** As discussed, this was by far the most difficult design task. It was less a science and more an art of iterative tuning. The balance between  $R_{\text{time}}$  and  $R_{\text{safety}}$  is the central conflict of autonomous navigation, and we experienced it firsthand.
- **Hyperparameter Tuning:** A poor learning rate ( $> 0.01$ ) or a fast decay ( $< 0.99$ ) caused the agent to fail to learn. We spent considerable time (documented in the log) finding the stable set of parameters used in the final model.
- **The "Sim-to-Real" Gap:** A conceptual challenge was constantly remembering that our 2D grid is a simplification. The agent's "success" is in the simulation. Deploying this on a real robot would require accounting for noise, sensor error, and dynamic physics, which our model ignores.

## 6.2 Team Reflection

The "Legion of Nazgûl" team operated with high cohesion. Our decision (Week 1) to split into "Environment" (Ayush, Adiba) and "Agent" (Vibhanshu, Vyom) teams, with one member (Yash) handling visualization and reporting, was highly effective. Regular integration (Oct 26) was critical. When the first integration failed, it forced both sub-teams to work together to debug the reward function, which became a central team effort. The group log was essential for keeping our efforts aligned.

## 6.3 Future Work

This project provides a strong foundation for several exciting future directions.

- **Complex Dynamics:** The immediate next step is to incorporate dynamic obstacles. This would involve adding simulated "pedestrians" or "cyclists" (other agents) that move randomly or on fixed paths, forcing our DQN agent to learn a more complex, reactive policy.
- **3D Environment:** As listed in our initial resources, migrating this simulation to a 3D environment like Unity or Blender (using ML-Agents) would be the next major leap. This would allow for realistic physics (acceleration, braking), complex sensor inputs (like virtual cameras), and navigating 3D terrain (slopes, overpasses).
- **Physical Deployment:** The ultimate goal would be to test the trained agent on a physical, small-scale robotic platform (e.g., a TurtleBot or custom-built robot). This would test the agent's ability to handle the sim-to-real gap, bridging the gap from a simulation to a tangible AI/ML product.

# 7 Intellectual Property / Research Dissemination

As stated in our proposal and group log, this project was pursued on the Publication Track. This final report is structured as a pre-print draft. Our team intends to refine this work and submit it to the **IEEE Intelligent Transportation Systems Conference (ITSC)**, as its themes align perfectly with our project's contribution.

# 8 Resources Utilized

- **Software:** Python 3.10, TensorFlow (Keras), OpenAI Gymnasium, Plotly, Google Maps (for data).
- **Hardware:** Local development machines (MacBook Pro) and Google Colab (T4 GPU) for accelerated training.

## 9 References

### References

- [1] Mnih, V. et al. (2015). Human-level control through deep reinforcement learning. *Nature*, 518, 529-533.
- [2] Sutton, R. S. & Barto, A. G. (2018). *Reinforcement Learning: An Introduction*. MIT Press.
- [3] van Hasselt, H., Guez, A., & Silver, D. (2016). Deep Reinforcement Learning with Double Q-learning. *AAAI*.