

# Autonomous Vehicle - Ethics Paper

James Derrod

jderrod@oxy.edu

Occidental College

## 1 Introduction

In the rapidly evolving field of artificial intelligence, the integration of deep reinforcement learning (DRL) with autonomous systems presents both groundbreaking potential and a multitude of ethical challenges.

For my comps project, I plan on developing a DRL algorithm to train an agent to successfully drive a car within a Unity environment. Unity is one of the largest, publicly available platforms for video game development, and allows for the use of the MLAgents toolkit, as well as rendering capabilities. The MLAgents toolkit is useful for integrating machine learning agents into Unity, facilitating the training, evaluation, and visualization of an agent's learning process and operational efficiency.

The problem of teaching an agent to drive a car using DRL in Unity can be framed as a Markov decision process (MDP), which is a mathematical framework for modeling decision making in situations where outcomes are partially random and partially under the control of a decision maker. Basic reinforcement learning is modeled as a Markov decision process, defined as:

- a set of environment and agent states,  $S$ ;
- a set of actions,  $A$ , the agent can take;
- A transition function to describe the probability of transition (at time  $t$ ) from state  $s$  to  $s'$  under action  $a$  [4].

In the context of my project, the components of the MDP are:

- **States ( $S$ ):**
  - A state represents the current condition of the environment as perceived by the vehicle/agent.
  - This could include:
    - \* Vehicle's position, orientation, speed, and acceleration.
    - \* Road conditions, number of lanes, size of lanes, traffic lights
    - \* Position and velocity of observed nearby vehicles and other dynamic obstacles.
    - \* Positions of non-moving observed obstacles such as barriers, lights, or signs.

- **Actions ( $A$ ):**

- Actions are the set of possible moves the agent can take at any state.
- This could include:
  - \* Accelerating, decelerating, or maintaining current speed. Defined by value in range  $[-1,1]$ .
  - \* Turning right, turning left, maintaining current direction. Defined by value in range  $[-1,1]$ .

- **Transition Function ( $P$ ):**

- This function will define the probability of moving from one state to another, given a particular action.
- **Ex.** If the agent decides to accelerate, this function will provide the probabilities of resultant states, such as:
  - \* Successfully overtaking another vehicle
  - \* An unwanted collision with an obstacle due to excessive speed.

Additionally, we want to define a reward function to shape the actions the agent will take given the goal of the agent.

- **Rewards ( $R$ ):**

- Rewards given for maintaining safe distance, staying within speed limits, obeying traffic rules, staying on road, moving towards goal.
- Negative rewards (penalties) given for collisions, near misses, lane/traffic violations.

Finally, we need to define a discount factor ( $\gamma$ ) which determines the present value of future rewards with a range of  $[0,1]$ . A higher discount factor would lead to an increased value for future rewards. A discount factor closer to 1 would encourage strategies that prioritize long-term benefits, such as maintaining a safe distance behind a vehicle to prevent a potential collision rather than prioritizing immediate gains like speeding to pass a slow vehicle.

With these elements defined, the problem of training an autonomous vehicle using DRL involves implementing a learning policy that maximizes cumulative future rewards.

This will be done through three steps:

1. **Initialization:** Start with random guesses for the policy, mapping from states to chosen actions.
2. **Experience Collection:** The agent will interact with the potentially dynamic environment, collecting and documenting experiences in terms of states, actions, rewards, and next states.
3. **Learning:** Implementing an algorithm such as Deep Q-Networks (DQN) to update the policy based on experiences collected. These updates aim to gradually converge towards an optimal policy that maximizes the cumulative reward.

My project aims to explore the capabilities of DRL within a controlled simulation, and push the boundaries of what autonomous systems can learn and achieve. However, the introduction of such sophisticated systems into widespread use brings with it a plethora of ethical considerations that must be thoroughly examined. This paper will address these issues, which include but are not limited to:

- Data bias
- Accessibility
- Power distribution
- Potential for misuse

I argue that completely addressing all ethical concerns associated with this technology is unattainable due to inherent limitations and complexities in the current technological landscape. By exploring each ethical concern in detail, this paper aims to highlight the significant challenges that and inevitable shortcomings of ethical management of deploying reinforcement learning in autonomous navigation systems.

## 2 Data Bias

Data bias occurs when an algorithm produces systematic errors due to inaccuracies in the data input. In machine learning, this often arises from non-representative or incomplete training data.

In autonomous driving systems, bias can lead to decisions that are not only incorrect, but potentially dangerous[2].

### 2.1 Sources of Data Bias in this Project

- **Simulated Environment Bias:** A simulated scenario will not fully account for the diversity of real world driving conditions.
- **Sampling Bias:** Selection of specific scenarios might disproportionately represent certain situations. A simulation will not accurately represent the variability in pedestrian behavior, complexities of different traffic patterns, or levels of visibility.

- **Human Input Bias:** "Safe Driving" may have a different definition for whoever is developing the rewards/algorithm, influenced by their own driving experience.

### 2.2 Impacts of Data Bias

- **Performance Limitations:** A model trained in biased conditions may excel in simulated tests but perform poorly in real world applications. This is often due to the model encountering situations it has not encountered in training and therefore has not learned how to handle those situations.
- **Safety Concerns:** A model trained predominately in daytime conditions might fail to appropriately recognize pedestrians at night, leading to higher risk of accidents.
- **Social Implications:** Systemic bias may emerge if, for example, a model performed well in geographically typical scenarios of the developed world, but poorly in less developed scenarios.

### 2.3 Challenges of Mitigating Data Bias

- **Complexity of Real World Representation:** Practical limitations must be acknowledged in simulating every possible driving condition due to technological constraints and the increase in complexity it would entail.
- **Bias Identification and Correction:** Biases are often not recognized until specific failures occur, prompting a return to the development phase to integrate new scenarios or re-balance the dataset.

## 3 Accessibility

Accessibility in autonomous vehicle technologies refers to the design of systems that are usable by people with a wide range of abilities and disabilities. In the context of deep learning and autonomous driving, ensuring accessibility is challenging and often overlooked in system design[1].

### 3.1 Sources of Accessibility Issues in this Project

- **Interface Design:** The user interface of the autonomous system may not be designed to accommodate users with varying physical, cognitive, and sensory abilities.
- **Technological Exclusivity:** Advanced technologies often require interactions that may exclude non-technical users or those with certain disabilities.
- **Economic Barriers:** The high cost of development and implementation of accessible technologies can

lead to prioritization decisions that exclude certain user groups.

### 3.2 Impacts of Poor Accessibility

- **User Exclusion:** If the system is not accessible, it may exclude people with disabilities from utilizing autonomous vehicles, potentially increasing societal inequalities.
- **Safety Risks:** Inaccessible emergency interfaces or controls can pose higher safety risks to users with disabilities.

### 3.3 Challenges of Enhancing Accessibility

- **Design Complexity:** Creating universally accessible systems increases the complexity of the design and testing phases.
- **Cost Implications:** The financial burden of implementing comprehensive accessibility features can be significant.
- **Continuous Adaptation:** Accessibility needs change as technologies and societal norms evolve, requiring ongoing updates and refinements.

## 4 Power Distribution

Power distribution in the context of autonomous vehicle technology refers to how decision-making capabilities and operational control are allocated among different stakeholders, including developers, users, regulatory bodies, and potentially malicious actors. The centralization or decentralization of power raises significant ethical concerns, particularly regarding autonomy, surveillance, and control [AVPowerInequity].

### 4.1 Sources of Power Imbalance in this Project

- **Developer Control:** Developers and corporations often retain significant control over autonomous technology, including proprietary software and hardware, data collection, and algorithm updates. This centralization of control can limit users' autonomy and potentially lead to dependencies that are difficult to break.
- **User Autonomy:** While autonomous vehicles promise enhanced mobility, they also risk diminishing the user's active decision-making role. Dependence on automated systems may erode individuals' skills and decision-making capacities, potentially leading to a passivity that can be exploited.
- **Regulatory Influence:** Governments and regulatory bodies possess the power to dictate terms of technology deployment, privacy standards, and safety regula-

tions. However, their influence is often reactive and lagging behind technological advances, which may prevent timely and effective governance.

### 4.2 Impacts of Power Distribution

- **Surveillance and Privacy:** Centralized control over autonomous vehicles could facilitate unprecedented levels of surveillance by developers or governments, especially as these vehicles are equipped to collect extensive data about users and their environments.
- **Manipulation and Control:** The concentration of power with developers or regulatory bodies could lead to scenarios where vehicle behavior is remotely manipulated for commercial or law enforcement purposes, raising acute ethical and privacy concerns.
- **Market Dominance:** The ability of large corporations to dominate the market with advanced technologies can stifle competition and innovation, potentially leading to monopolistic behaviors and further concentrating power.

### 4.3 Challenges of Addressing Power Imbalance

- **Technological Opacity:** The complexity and proprietary nature of autonomous vehicle technologies make it difficult for users and regulators to fully understand and effectively oversee these systems. This opacity can shield significant power imbalances from public scrutiny and regulatory oversight.
- **Evolving Technologies:** Rapid advancements in technology continually shift power dynamics, often outpacing the development of corresponding ethical norms and regulatory frameworks. This lag complicates efforts to establish fair and equitable power distributions.
- **Global Disparities:** Different countries and regions may have varying capabilities to adopt and regulate new technologies, leading to global disparities in power and control over autonomous vehicle technology. This can exacerbate inequalities and hinder comprehensive governance.

## 5 Potential for Misuse

The development of autonomous vehicles using deep reinforcement learning opens up several avenues for misuse that can have significant societal impacts. Whether by malicious actors or through unintended consequences, the potential for misuse is a critical ethical concern that complicates the deployment of this technology[3].

## 5.1 Sources of Potential Misuse

- **Malicious Exploits:** As autonomous vehicles are controlled by software that interacts with a complex environment, they are vulnerable to hacking and other cyber attacks. These exploits could lead to unauthorized control over vehicle behavior, endangering public safety.
- **Surveillance Abuse:** Given the data-intensive nature of autonomous driving technologies, there is a risk that this data could be used for intrusive surveillance. This could be perpetrated not just by state actors but also private companies or criminals, using data to track individuals without their consent.
- **Software Manipulation:** The software that governs autonomous vehicles can be altered to behave in unethical ways, such as prioritizing the safety of the vehicle over pedestrians in violation of ethical norms, or creating situations that could lead to insurance fraud.

## 5.2 Impacts of Misuse

- **Public Safety Risks:** Malicious control or failure modes introduced by exploits could result in accidents, injuries, or worse. The potential for widespread harm is significantly elevated with the proliferation of autonomous vehicles.
- **Privacy Violations:** Unauthorized surveillance and data misuse could lead to massive breaches of privacy, affecting countless individuals. These breaches would not only violate individual rights but could also undermine public trust in autonomous technologies.
- **Legal and Ethical Challenges:** Misuse of autonomous vehicle technology raises complex legal issues, such as liability in the case of accidents caused by manipulated software, and ethical issues regarding surveillance and data privacy.

## 5.3 Challenges of Mitigating Misuse

- **Security Measures:** Implementing robust security measures to protect against hacking and unauthorized access remains a major technical challenge. The complexity of autonomous systems makes them difficult to secure completely.
- **Regulatory Oversight:** There is a need for comprehensive regulatory frameworks that can keep pace with technological advancements and address the many ways in which misuse can occur. Developing these frameworks is complicated by the global nature of technology deployment and the varying capabilities of regulatory bodies.
- **Technological Sophistication:** As autonomous technologies become more sophisticated, so too do the

methods of misuse. Keeping ahead of malicious actors requires continuous innovation and vigilance in technology design and policy enforcement.

## 6 Conclusion

This project applies deep reinforcement learning (DRL) to develop autonomous vehicles within a Unity environment, bringing to light several ethical challenges. These include data bias, accessibility limitations, power distribution issues, and the potential for misuse. These problems highlight the inherent limitations of using simulated environments and centralized control systems in this technological field.

Despite the advancements demonstrated, the reasons presented argue that it is not possible to fully address all ethical concerns with the current technology, both within this project as well as the field of autonomous vehicles as a whole. The fast pace at which technology evolves makes it even harder to manage these issues effectively.

This highlights the importance of continuous ethical review and the need to adapt our methods as technology changes, and stresses the need for ongoing improvements to ensure that technological developments are in line with ethical standards.

## References

- [1] Alberto Dianin Elisa Ravazzoli, Georg Hauger. *Implications of Autonomous Vehicles for Accessibility and Transport Equity: A Framework Based on Literature*. 2021. URL: <https://www.mdpi.com/2071-1050/13/8/4448>.
- [2] David Danks, Alex John London. *Algorithm Bias in Autonomous Systems*. 2017. URL: <https://www.cmu.edu/dietrich/philosophy/docs/london/IJCAI17-AlgorithmicBias-Distrib.pdf>.
- [3] Sunniva F. Meyer Rune Elvik, Espen Johnsson. *Risk analysis for forecasting cyberattacks against connected and autonomous vehicles*. 2021. URL: <https://link.springer.com/article/10.1007/s12198-021-00236-4>.
- [4] Wikipedia. *Markov Decision Process*. 2024. URL: [https://en.wikipedia.org/wiki/Markov\\_decision\\_process](https://en.wikipedia.org/wiki/Markov_decision_process).