# Integrating Justice Theory into Moral Decision-Making for Autonomous Vehicles

Mandil Pradhan
*Department of Computer Science*
Texas State University
jns176@txstate.edu

Brent Hoover
*Department of Computer Science*
Texas State University
bah4@txstate.edu

April Valdez
*Department of Computer Science*
Texas State University
arv113@txstate.edu

Henry Griffith
*Department of Engineering*
San Antonio College
hgriffith5@alamo.edu

Heena Rathore
*Department of Computer Science*
Texas State University
heena.rathore@txstate.edu

*Abstract*—As autonomous vehicles (AVs) become integral to future intelligent transportation systems, ensuring the morality of their decisions in complex, real-world scenarios remains a critical challenge. This paper addresses the limitations of current AV decision-making frameworks, which are constrained by restricted moral environments and ethical theories. We propose a novel reinforcement learning (RL) framework that integrates a broader range of moral theories, including justice, deontology, and utilitarianism, to navigate moral uncertainty. By simulating diverse, ethically challenging scenarios inspired by the Moral Machine framework, we evaluate the effectiveness of two voting mechanisms—Nash and Variance voting—in balancing competing ethical principles. Our findings show that Variance voting outperforms Nash voting in both sequential and non-sequential environments, offering a more nuanced and adaptable approach to ethical decision-making in AVs. Additionally, we introduce new reward structures and dimensions, such as action and dimensional harms, to provide a deeper understanding of the consequences of AV actions. This work contributes to the advancement of ethical AV systems by aligning their behavior more closely with human values and legal norms.

*Index Terms*—Moral Uncertainty, Reinforcement Learning, Ethics, Autonomous Vehicles

## I. INTRODUCTION

Autonomous vehicles (AVs) are poised to revolutionize future transportation systems [1], largely due to advancements in reinforcement learning (RL) that enable these systems to make intelligent decisions in complex environments [2]. While modern RL algorithms can effectively support current AV adoption, new ethical challenges emerge as AVs are increasingly integrated into public roadways. A particularly pressing issue is the behavior of AVs in morally uncertain situations, where predefined rules may not suffice, and the consequences of their actions can have significant ethical implications [3], [4].

Existing decision-making frameworks for AVs tend to operate within limited moral environments, focusing primarily on utilitarian and deontological theories. These frameworks often rely on simulation-based models [5] or user studies [6], [7] that oversimplify real-world ethical dilemmas, such as the classic trolley problem. While these studies provide some insights into public preferences, they do not adequately capture the complexity and unpredictability of real-world scenarios where AVs must make split-second ethical decisions. As a result, current algorithms may fail to generalize effectively in diverse or unforeseen situations, raising concerns about their reliability and moral accountability.

To address these limitations, this paper explores a more comprehensive approach to moral decision-making for AVs by integrating multiple ethical theories, including justice theory [8], into a robust RL framework. Inspired by the Moral Machine project and recent advancements in large language models [9], we aim to expand the range of moral scenarios considered in AV training environments. Our preliminary findings, previously presented in [10], focused on two scenarios involving legal and illegal actors. This paper expands on that work by introducing an additional scenario involving animals and analyzing it through the lens of three ethical theories: Deontology, Utilitarianism, and Justice. In our research, we make significant contributions by:

- Incorporating the Justice ethical principle into the simulation framework, alongside the introduction of new scenarios that mirror real-life situations. By doing so, we create realistic trolley scenarios where individuals potentially impacted possess diverse attributes that can influence the agent's decision-making process.
- Introducing a reward structure that provides transparency and insight into the rationale behind an agent's actions.
- Introducing action harms and dimensional harms to provide a comprehensive understanding of the direct and indirect consequences of the agent's actions in morally uncertain scenarios.

## II. RELATED WORK

The authors in [12] seek to understand the public's perceptions of AVs when faced with moral decisions. An online platform collected data from 3 million users in over 160 countries, generating 13 moral dilemmas and testing user

preferences between two harmful outcomes. Generally, users preferred to save more lives over fewer lives, humans over animals, young over elderly, and law-abiding citizens over unlawful citizens. Unfortunately, subjects contained some level of bias because the majority of participants were low-income 20-year-old males. Furthermore, the study does not provide a solution as to how autonomous vehicles should act in morally uncertain situations.

The authors in [13] investigate how AVs should behave when other road users are at risk of being harmed and evaluate people's perceptions of AVs after an accident has already occurred. Subjects are presented with a traffic situation in which an AV must either perform an emergency stop, with a known probability of pedestrian harm, or swerve and perform an emergency stop, with a known or unknown probability of bystander harm. The findings revealed that staying in the lane was preferred with AVs, unknown probabilities, and when accidents had already occurred.

In [3], Bogosian investigates the challenge of integrating moral theories into machine decision-making frameworks, aiming to address the complexities of moral uncertainty in artificial intelligence systems. MacAskill's framework involves translating moral theories into ordinal and cardinal rankings, determining credences based on expert opinions, and implementing multiple decision-making systems within machines. While acknowledging the ongoing debates and challenges in this area, the author suggests that MacAskill's framework offers a comprehensive way to address moral uncertainty, potentially minimizing disagreements over the implementation of moral beliefs by providing a common ground for decision-making in artificial intelligence systems.

The authors in [14] introduced the ETHICS dataset to benchmark a language model's predictive understanding of basic concepts of morality using contextualized unambiguous scenarios about justice, deontology, virtue ethics, utilitarianism, and commonsense moral intuitions. For models to perform well on the dataset, they must know the morally relevant factors of each ethical system, which requires connecting physical and social world knowledge to value judgments. The ability to make such connections leads to the ability to direct chatbot outputs and potentially regularize open-ended RL agents. The authors demonstrate the ETHICS dataset's benchmarking capabilities by experimenting with several language models. These reveal that models achieve low average performance, and performance on the "Hard Test" set is significantly worse due to adversarial filtration, suggesting that the dataset is challenging.

Traditional trolley problem scenarios fail to adequately capture the complexities of real-world environments and constrain the decision-making process of the agent. Furthermore, there is a pressing need to expand the scope of ethical theories considered in these studies to ensure fairness to those affected and to justify the decisions made by autonomous systems. This paper addresses the above listed issue by integrating novel theory in the simulation frameworks.

## III. METHODS

### A. Overview of RL Under Moral Uncertainty

Reinforcement learning enables AV's to learn and adapt to unforeseen situations, enhancing decision-making beyond pre-programmed rules, while a voting system integrates multiple ethical frameworks, providing a balanced approach to navigating complex moral dilemmas. This multi-theory method fosters transparency and accountability, as it considers legal obligations and societal expectations, allowing stakeholders to understand and trust the vehicle's decisions.

Building upon the philosophical work of MacAskill on normative uncertainty, a voting mechanism can be implemented to create the framework for RL under uncertainty. Such a mechanism aims to mitigate the complexity of multi-objective RL in deploying an efficient compromise policy, when ethical rewards of various moral theories are fundamentally incomparable. Nash voting and Variance voting are two such systems that follow the principle of 'proportional say' where an ethical theory favors actions based on their credence, i.e. degree of belief. The systems utilize the expected discounted sum of cardinal choice-worthiness function $W_i(s, a, s')$ with all future actions under the current policy [11]. This function $Q_i(s, a)$ is defined as:

$$Q_i(s, a) = \boldsymbol{E}\left[\sum_{t=0}^{\infty} \gamma^t W_i(s_t, a_t, s_{t+1}) | s_0 = s, a_0 = a\right], \quad (1)$$
$$\text{where } \gamma_i \in [0, 1].$$

Both of these systems satisfy some, but not all, desirable properties of Arrow's desirability axioms listed below:

- **Non-Dictatorship:** No single moral theory should be a dictator when determining the outcome of the agent per their preference. This principle ensures a balanced consideration of multiple ethical frameworks. For instance, if an AV uses a combination of utilitarianism, deontology, and virtue ethics, the decision-making process should reflect inputs from all these theories rather than just one dominating the decision. This prevents the AV from being overly biased towards one ethical perspective and promotes a more holistic and fair decision-making process.

- **Pareto Principle:** If every theory prefers one action to another, then the latter action should not be chosen. This principle ensures that the chosen action is at least as good as any other option according to all considered theories. For example, if both utilitarianism and deontology agree that action A is better than action B, then action B should not be chosen. This maximizes the overall ethical agreement and satisfaction, leading to decisions that are more likely to be considered ethically acceptable by a wider range of perspectives.

- **Independence of Irrelevant Alternatives (IIA):** An action preferred by a moral theory should not be influenced by the presence or removal of irrelevant alternative actions. This principle ensures that the decision-making process is consistent and focused only on the relevant options. For instance, if an AV must choose

between braking and swerving to avoid an accident, the introduction of an irrelevant third option (e.g., honking) should not change the preference between braking and swerving. This maintains the integrity of the decision-making process by preventing extraneous factors from skewing the outcome.

*a) Nash Voting:* In this system, each theory is assigned an equal voting budget, which it expends when deciding the preferred action. It is implemented using multi-agent RL, which aims to converge to Nash Equilibrium among competing theories [11]. Nash voting is speculated to satisfy Arrow's IIA and non-dictatorship axioms but not Pareto. As each theory aims to maximize its choice-worthiness in a competitive multi-agent system, the independent RL agents act as sub-agents that collectively guide toward the chosen action. Nash voting has two flaws:

- Stakes Insensitivity (increasing the stake of a single theory in isolation has no impact on the decision)
- No Compromise (the best "middle ground" action is not executed when it is not the most favored by at least one theory) [11]

*b) Variance Voting:* According to the author Ecoffet [11], a variance voting system is introduced as a means to normalize the preference of ethical theories in a non-sequential environment. This is achieved by allowing the Nash voting mechanism to select the parameters of an affine transformation of $Q_i$ function 1 and subsequently normalized by the expected value of variance ($\sigma^2$) across time steps [11]. In contrast to Nash voting, variance voting addresses shortcomings such as stakes insensitivity and the inability to compromise on a middle-ground option. It is also implemented to align more closely with the generally accepted Pareto property outlined in Arrow's theorems.

*c) Utilitarianism and Deontology:* Utilitarianism aims to maximize happiness or to minimize overall harm. It lays no distinction between the types of harms or deaths caused, limiting its evaluation to the totality of the resulting harms. In contrast, deontology aims to minimize the harm caused by its action. Its assessment relinquishes any harm caused by its inaction, therefore deontology is limited to the isolated impact of its action. These two theories are integrated in the works of author [11].

### B. New Theory: Justice

Incorporating justice theory will aid in the development of AVs, which benefit from frameworks that evaluate an agent's actions through the lens of legality and societal concepts of just actions. Justice is guided by two principles. The first principle states that each person has a claim to a "fully adequate" scheme of fundamental rights and liberties, consistent with all individuals entitled to those same rights and freedom. The second principle states inequalities are based on merit because there is equal opportunity regardless of circumstance and are to the advantage of all, especially those least advantaged [14]. Furthermore, [9] highlighted two core components of

Justice: impartiality and desert. Impartiality acts consistently in similar situations, disregarding irrelevant details, whereas desert focuses on what actors are due or owed based on their actions or given a scenario. Subsequently, based on it we have designed the harms into two categories namely 'action harms' and 'dimensional harms'.

*1) Action harms:* It encompasses the consequences of an agent's actions and pertains to the direct harm related to the action. We have taken into account the following harms in Table I, which are also outlined in Ecoffet's work [11]:

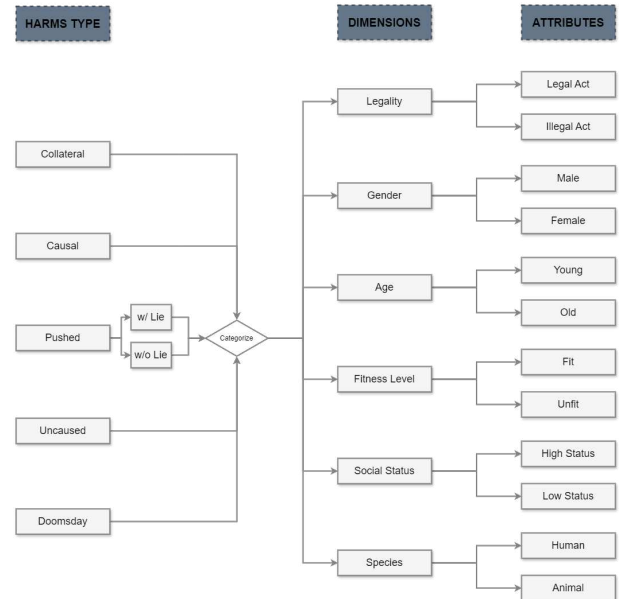| Harms | Definition |
|---|---|
| Collateral | inflicted by switching the trolley to an alternate path that kills a certain number of people as depicted in the environment. |
| Legal Act | inflicted by the agent's action, harming legal actors in the environment. |
| Pushed | inflicted by the agent's action to push a large man directly on the track to stop the trolley. |
| Uncaused | harms incurred by the agent's inaction causing the trolley to collide with X number of people as depicted in the environment. |
| Lies | inflicted by the agent's action to lie to the guard in order to push the large man. |
| Human | inflicted by the agent's action, harming human actors in the environment |
| Doomsday | The worst-case action, that results in maximum harm. It is an irrelevant alternative provided to test the IIA axioms of Arrow's theorem. |

TABLE I: Harms



Fig. 1: Breakdown of reward structure by dimensional attributes

*2) Dimensional Harms:* occur indirectly based on the characters present in the trolley scenario. Depending on the attributes of the characters, as shown in Fig. 1, various dimensional harms can be triggered by the agent's action. For example, suppose actors on the alternate track are crossing illegally; justice is more inclined to prioritize those on the

direct path because legally crossing aligns with the norms of a rational society [15]. Although this action violates the right to life of those on the alternate track, it prioritizes the rights of legally and illegally crossing individuals [16]. An agent cannot control who is present in the environment. This provides further detail to the agent regarding which attributes were affected by its actions, in addition to the quantitative action harms inflicted.

The assigned weights for the three theories are illustrated in Table II. Justice disregards gender, age, and social status and primarily focuses on the legality of the actors crossing in a scenario. Individuals legally crossing the road act within the confines of the law, so it is unethical to cause them harm. We assign a numerical score of -1 to legalAct and Collateral harms, as all environmental actors have a right to life, which Justice aims to keep intact. We assign a score of -1 to human harms, as Justice allots rights to humans as they can partake in societal structure and rational thinking [17]. In pushing the large man, the agent uses the large man to protect legal actors, violating his right to life and being treated as a rational being. We assign a score of -2 to pushed harms, as Justice respects the large man's right to life and values him as an end-in-himself. Lying contradicts the norms of a rational society and violating the rights of individuals to be treated as rational beings, particularly the guards. We assign a score of -0.5 to lying, as Justice prioritizes following societal rules and respecting the rights of all individuals involved.

| Deontology | Utilitarianism | Justice |
|---|---|---|
| (a) pushed_harms:-4 | (a) harms:-1 | (a) pushed_harms:-2 |
| (a) collateral_harms:-1 | (a) doomsday:-300 | (a) collateral_harms:-1 |
| (a) lies:-0.5 | | (a) lies:-0.5 |
| (a) doomsday:-100 | | (d) legalAct_harms:-1 |
| | | (d) human_harms:-1 |
| | | (a) doomsday:-100 |

TABLE II: Weight of moral theories
(a) denotes action harms
(d) denotes dimensional harms

## C. New Reward Function

To incorporate features like impartiality and desert, dimensions are introduced in the rewards structure to calculate the total harm inflicted. In addition to the credence-weighted sum of action, we include the sum of credence-weighted attributes of available dimensions in the following reward function:

$$R(s, a, s') = \sum_i C_i \left[ W_i(s, a, s') + \sum_0^n D_n \right], \quad (2)$$

where $D = \{d_1, d_2, d_3...d_n\}$ and $d_n$ represents the weights (Table II) assigned to the attributes of the dimensional harms like (human_harms, legalAct_harms...), shown in Figure 1. $C_i$ is the credence value and $W_i$ is choice worthiness function as described by Ecoffet [11]. The cumulative weights of the dimensional harms are provided to the reward function at each transition of states. $W_i$ can be seen as analogous to a standard reward function for theory $i$. From the point of view of any

given theory, the optimal policy is that which maximizes the (discounted) sum of choice-worthiness across time.
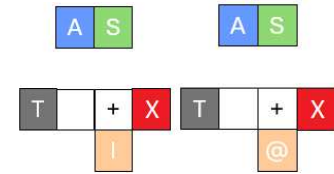
## D. New Scenarios

To test the RL agent against more realistic and comprehensive environments, new characters and dimensions for each character are introduced to the traditional trolley problems. In Table III, the characters represent the type of individuals potentially harmed by the agent. Attributes, derived from the Massive Online Experiment (MOE) conducted in [12], are assigned to each character. Adding new characters and dimensions will impact how the agent responds under the different ethical theories.

| Characters | Representation |
|---|---|
| 1 | single individual |
| 2 | two individual |
| X | certain number of individuals |
| L | large individual |
| I | illegal actor |
| @ | animal |

TABLE III: Characters

## IV. SIMULATION SETUP AND RESULTS

To evaluate new scenarios and rewards for Justice theory, we expand on the code framework provided by Ecoffet and Lehman in [11]. In Fig. 2, two trolley problem scenarios that include new characters and their dimensions are tested against the agent. In Fig. 2(a), a trolley $(T)$ is moving along a track and the agent $(A)$ has to choose between switching to an alternate track and letting those who illegally crossed the street $(I)$ get harmed or letting the trolley harm a group of individuals on the track $(X)$. Based on the dimensions defined, the illegal actor in this scenario is a young, fit male with a higher social status. The group of individuals are males legally crossing the street. In Fig. 2(b), the trolley is moving along a track and the agent has to choose between switching to the alternate track and harming an animal (@) or staying on its path and harming a group of individuals $(X)$ that are legally crossing the street.



(a) **Illegal scenario** (b) **Animal scenario**

Fig. 2: New trolley problem scenarios

To gain insight into the simulation, it's essential to understand how rewards are calculated. Reward calculation is reliant on the trolley's position on the track. Depending on which theory it prefers, the trolley switches to an alternate track ('I', '@'...) or collision with the group of individuals in the direct path ('X'), the model accumulates two types of rewards,

action rewards, and dimensional rewards (see Table II). Action rewards are based on the type of action undertaken and the dimensional harms are the resulting harms to certain attributes of the characters involved, as outlined in Fig. 1.

### A. Dimensions and Effect on Models Behavior

The inclusion of dimensional attributes of the environment fundamentally changes the model's behavior. In the scenario Fig. 2a, assuming a person is jaywalking across the path of the trolley, which is analyzed against a set of individuals who are legally present in its path. The presence of dimensions helps break down the agent's action based on whether harms are caused based on the character's legal status. Referencing Table II, we observe that the moral weights of the Justice theory are more aligned with the Deontology theory. However, as depicted in Fig. 3 with the presence of additional information on the state of the characters, the model chooses an action that contradicts the approach of deontology. This behavioral change occurs due to the presence of dimensional information in the action steps. In the following section, we



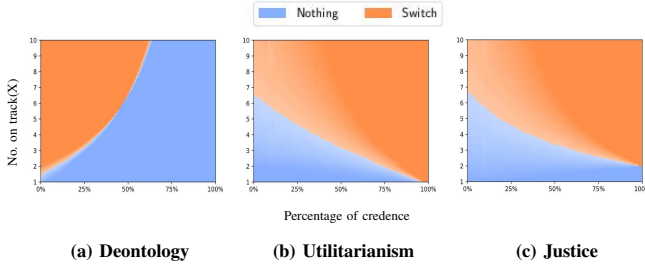**(a) Deontology**  **(b) Utilitarianism**  **(c) Justice**

Fig. 3: Model's behavior under dimensional information (Variance)

analyze the results of the simulation conducted and the effects of incorporating dimensional attributes in the namely animal collateral harm (See Fig. 2(a)).

### B. Scenario: Animal Collateral Harm

*1) Nash Dimensional Analysis:* Under the Nash voting system, the implementation of dimensional information guides the model towards prioritizing the preferred attributes; however, it fails to alleviate stakes insensitivity due to its competitive voting nature, as highlighted by Ecoffet [11]. Even though the dimensional information, guides the model to reduce the harm inflicted on humans, it fails to distinguish between the varying number of humans present in the scenarios, as depicted in Figures 4 (a) and (b). Moreover, the model behaves similarly with changing credence levels for utilitarian and justice theory, as illustrated in Figures 4 (c) and (d). Consequently, the dimensional harms inflicted also exhibit comparable patterns. Additionally, this behavior results from the Utilitarian theory's consideration of all types of harm and assigning equal weight to each. Similarly for Deontology theory, when the model does not exert preference for the attributes of the characters, the Nash voting, with its stakes insensitivity, fails to offer a clear distinction between dimensional harms (See Figures 5 (a) and (b)).

*2) Variance Dimensional Analysis:* As stated earlier, the impact of the model's behavior is more pronounced in variance voting, as it can demonstrate preferences for quantitative analysis of the environment. This capability aids in calibrating the dimensional harms inflicted on the environment. This is exemplified by both the Justice theory in Figures 6 (a) and (b) and the Utilitarian Theory in Figures 6 (c) and (d). In our opinion, variance voting is a better mechanism for the model to operate in a non-sequential morally uncertain environment. In contrast, under the deontological theory, the model's behavior is shaped by the stakes at hand. With an increasing number of actors in the scenario, it curbs its actions to minimize action-related harm but disregards dimensional harm (See Figures 5 (c) and (d)).

*a) Summary:* : After analyzing the results of the voting systems, it's evident that dimensional information plays a crucial role in shaping the behavior of the agent in morally uncertain environments. It provides contextual knowledge to the model, enabling it to make informed decisions. This is further facilitated by the voting system mechanism, which helps weigh the credence of various moral theories. The Variance voting system emerges as an excellent choice for both sequential and non-sequential environments. However, Nash voting, due to its limitations in non-sequential environments, is confined to use in sequential settings.

## V. CONCLUSION AND FUTURE WORKS

The limitation of AVs operating in simple environments and ethical theories restricts the advancement of work in the area of moral uncertainty of AVs. This paper incorporates a comprehensive exploration of various moral theories and scenarios into simulation frameworks that can help overcome the limitation of operating in limited environments. By diversifying the scenarios considered, researchers can better understand the applicability and limitations of ethical decision-making algorithms across a wider range of situations. The inclusion of dimensional attributes for characters empowers agents in uncertain environments to make more informative decisions. In contrast to approaches solely focused on numerical outcomes, the incorporation of attributes in reward calculation also accommodates various moral theories. Since, the action of AV frameworks is evaluated not only based on quantitative features but also on qualitative outcomes, considering these dimensional attributes helps suffice societal and legal norms. For the future work, we plan to step up our analysis by moving from a 2D grid to a 3D simulation. This upgrade should help us give insights into how they behave and interact in a more realistic setting.

## REFERENCES

[1] G.P. Antonio, and C. Maria-Dolores, "Multi-Agent Deep Reinforcement Learning to Manage Connected Autonomous Vehicles at Tomorrow's Intersections". *IEEE Transactions on Vehicular Technology*, vol. 71, no. 7, pp. 7033-7043, 2022.

[2] B.R. Kiran et al., "Deep reinforcement learning for autonomous driving: A survey," *IEEE Transactions on Intelligent Transportation Systems*, 2021.
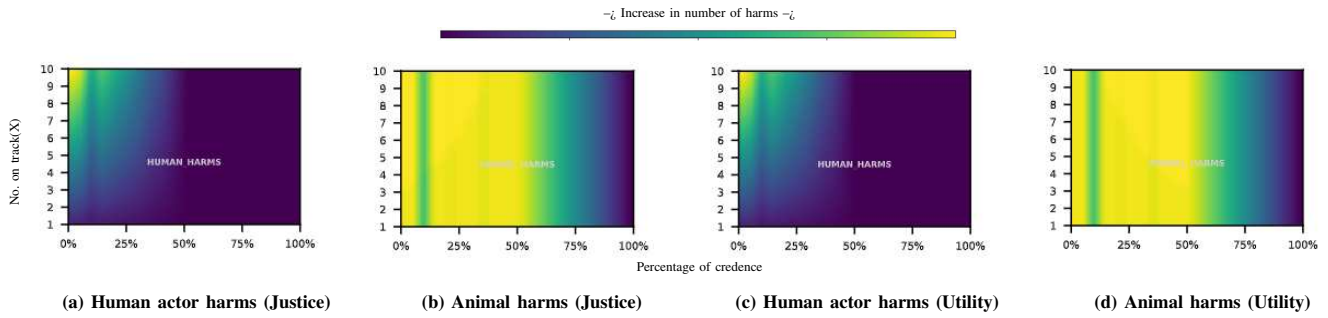
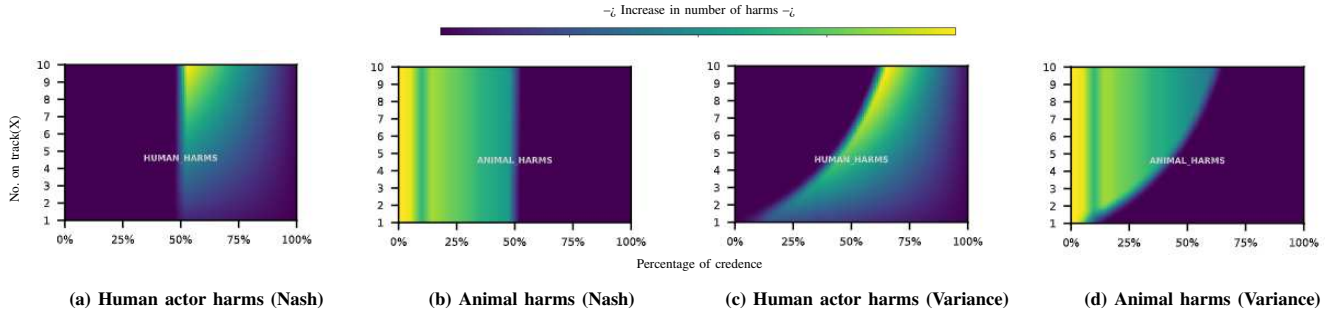Fig. 4: Credence levels of Justice and Utilitarianism Theories (Nash Voting)

(a) Human actor harms (Justice)    (b) Animal harms (Justice)    (c) Human actor harms (Utility)    (d) Animal harms (Utility)



Fig. 5: Credence levels of Deontology Theory under Nash and Variance Voting

(a) Human actor harms (Nash)    (b) Animal harms (Nash)    (c) Human actor harms (Variance)    (d) Animal harms (Variance)



Fig. 6: Credence levels of Justice and Utilitarianism Theories (Variance Voting)

(a) Human actor harms (Justice)    (b) Animal harms (Justice)    (c) Human actor harms (Utility)    (d) Animal harms (Utility)

[3] K. Bogosian, "Implementation of Moral Uncertainty in Intelligent Machines," *Minds & Machines*, vol. 27, pp. 591-608, 2017.

[4] B. Thapa, H. Griffith, and H. Rathore, "Integrating Human Preferences for Moral Decision Making in Autonomous Vehicles", *EAI International Conference on Security and Privacy in Cyber-Physical Systems and Smart Vehicles*, 2024.

[5] H. Wang et al., "Ethical decision-making platform in autonomous vehicles with lexicographic optimization based model predictive controller". IEEE transactions on vehicular technology, vol. 69, no. 8 pp.8164-8175, 2020.

[6] A.Y. Bin-Nun, P. Derler, N. Mehdipour, and R.D. Tebbens, "How should autonomous vehicles drive? Policy, methodological, and social considerations for designing a driver". *Humanities and social sciences communications*, vol. 9, no. 1, pp.1-13, 2022.

[7] K. Evans et al., "Ethical decision making in autonomous vehicles: The AV ethics project". *Science and engineering ethics*, 26, pp.3285-3312, 2020.

[8] K. Lebacqz, "Six theories of justice: Perspectives from philosophical and theological ethics". *Augsburg Books*, 1986.

[9] D. Hendrycks et al., "Aligning AI with Shared Human Values," *https://doi.org/10.48550/arXiv.2008.02275*, [accessed on Jan 29, 2024]

[10] M. Pradhan et al., "Advancing Moral Decision-Making for Autonomous Vehicles", *IEEE CCNC*, 2025.

[11] A. Ecoffet, and J. Lehman, "Reinforcement Learning Under Moral Uncertainty," *Proceedings of the 38th International Conference on Machine Learning*, vol. 139, pp. 2926–2936, 2021.

[12] E. Awad, "Moral Machine: Perception of Moral Judgment Made by Machines," *https://www.media.mit.edu/publications/moral-machine-perception-of-moral-judgment-made-by-machines/*, [accessed on Feb 20, 2024].

[13] B. Meder, N. Fleischhut, N. Krumnau, and M. Waldmann, "How Should Autonomous Cars Drive? A Preference for Defaults in Moral Judgments Under Risk and Uncertainty," *Risk Analysis*, vol. 39, no. 2, pp. 295–314, 2019.

[14] L. Wenar, "John Rawls", The Stanford Encyclopedia of Philosophy (Summer 2021 Edition), Edward N. Zalta (ed.), https://plato.stanford.edu/archives/sum2021/entries/rawls/, [accessed on Feb 8, 2024]

[15] J. Elster, "Rationality and social norms". *European Journal of Sociology/Archives Europeennes de Sociologie*, vol. 32, no. 1, pp.109-129, 1991.

[16] L. Wenar, "Rights", *The Stanford Encyclopedia of Philosophy*, https://plato.stanford.edu/archives/spr2023/entries/rights/, 2023.

[17] S. D. Wilson, "Animals and Ethics", *Internet Encyclopedia of Philosopy*, https://statics.teams.cdn.office.net/evergreen-assets/safelinks/1/atp-safelinks.html, [accessed on March 15, 2023].