Navigating the Unimaginable: Decentralizing the Autonomous Vehicle's Ethical Compass

A White Paper by The Fijishi Center for Autonomous Systems Ethics.



Disclaimer: The following is intended for information purposes only, and may not be incorporated into any contract. It is not a commitment to deliver any material, code, or functionality, and should not be relied upon in making purchasing decisions. The development, release, and timing of any features or functionality described for Fijishi's products remains at the sole discretion of Fijishi.

Index

Executive Summary	Page 3
Introduction: The Unavoidable Dilemma	Page 3
The Challenge of the "Unimaginable": Why Algorithms Fail	Page 3
The Flawed Compass: Limitations of Centralized Ethics	Page 4
Decentralizing the Dilemma: A New Paradigm	Page 4
Fijishi Aeterna: A Framework for Adaptive Ethics	Page 5
Benefits of Decentralization	Page 5
Conclusion	Page 6

Executive Summary

Autonomous vehicles (AVs) promise unprecedented safety and efficiency, yet they confront a profound ethical void in rare, unavoidable accident scenarios. Faced with split-second decisions involving potential harm, such as choosing between protecting a passenger and a pedestrian, AVs currently lack a universally accepted, preprogrammed "right" answer. This white paper argues that relying solely on fixed algorithms for such complex moral judgments is insufficient and fails to reflect the dynamic nature of human ethics. We propose a decentralized approach, introducing the conceptual framework where organizations leverage platforms like Fijishi and AI systems like Aeterna to train AV ethical decision-making not through rigid rules, but by learning from diverse human inputs. This paradigm shift aims to cultivate an ethical compass within AI that is a continuously adapting consensus, mirroring the evolving moral landscape of society.

1. Introduction: The Unavoidable Dilemma

The advent of autonomous vehicles marks a significant leap forward in transportation technology. With their potential to eliminate human error – the cause of the vast majority of road accidents – AVs hold the promise of saving millions of lives annually. However, this promise is shadowed by a difficult truth: while AVs can drastically *reduce* accidents, they cannot eliminate them entirely. In exceedingly rare, unavoidable situations, an AV may be forced into a scenario where any potential action results in harm.

Consider the classic "trolley problem" transposed onto the road: an AV's path is obstructed, and the only options are to swerve and potentially harm pedestrians, or remain on course and potentially harm the passenger(s). These are not engineering challenges; they are profound ethical quandaries with no easy answers. Current AV development largely focuses on *avoiding* such scenarios through advanced sensing and predictive algorithms. But the "unimaginable" – the moment where an unavoidable choice must be made – remains a critical, unresolved hurdle.

The public rightfully demands to know how AVs will be programmed to act in these life-or-death situations. Who decides the hierarchy of lives? Is a passenger's life inherently prioritized over a pedestrian's? What if the pedestrian is a child, or the passenger is a family? Pre-programming a rigid hierarchy feels cold, calculating, and ultimately, ethically precarious because there is no global consensus on how to weigh these competing values.

2. The Challenge of the "Unimaginable": Why Algorithms Fail

Traditional AI and programming excel at executing defined rules and optimizing for clear objectives. However, ethical decision-making in complex, novel situations defies simple rule-based logic. Human ethics are not a static set of axioms; they are fluid, context-dependent, and influenced by a myriad of cultural, social, and personal values.

Attempting to hardcode ethical responses for every conceivable edge case in autonomous driving is an impossible task. The sheer number of variables – the number and age of potential victims, their relationship to the vehicle occupants, the nature of the potential harm, the legal and societal context – creates a combinatorial explosion of scenarios that cannot be exhaustively catalogued and prescribed.

Moreover, even if such a catalogue were possible, deciding on the "correct" rule for each scenario would require a centralized authority to impose its ethical framework on all others. This raises serious questions about representation, bias, and public acceptance. Whose ethics get programmed into the car? The engineer's? The company's? The regulator's? A centralized approach risks embedding the biases and values of a select few into machines making life-and-death decisions for everyone. This top-down imposition of ethics creates a trust deficit, as the public has no assurance that the AV's moral compass aligns with societal values.

The core problem is that the "moral compass" needed by autonomous vehicles in these critical moments is not a fixed, pre-determined algorithm. It needs to be something more nuanced, more representative, and capable of adapting to the complexities and diversity of human ethical reasoning.

3. The Flawed Compass: Limitations of Centralized Ethics

Current approaches to AV ethics often involve expert panels or internal company guidelines defining the rules of engagement for unavoidable accidents. While well-intentioned, these centralized methods suffer from inherent limitations. A small group, no matter how well-meaning, cannot possibly represent the full spectrum of ethical perspectives within a diverse society. Their decisions may be influenced by their own cultural backgrounds, values, and even corporate pressures.

Furthermore, ethics are not static. Societal values evolve over time, influenced by cultural shifts, technological advancements, and collective learning. A centralized, fixed ethical algorithm programmed today could become outdated or even unacceptable tomorrow. Relying on a rigid, pre-programmed moral code is akin to navigating with an ancient map in a rapidly changing landscape. It fails to account for the dynamic nature of human morality and the need for continuous adaptation.

The lack of transparency and public input in developing these centralized ethical frameworks further erodes trust. Without understanding *why* an AV might be programmed to act in a certain way in a crisis, the public is left to speculate and fear. This opaqueness hinders adoption and raises significant questions about accountability when an unavoidable accident occurs.

4. Decentralizing the Dilemma: A New Paradigm

Addressing the AV ethical challenge requires a fundamental shift in approach. Instead of attempting to define and impose a fixed ethical code, we must focus on building AI systems capable of *learning* and *adapting* their ethical framework from the source: diverse human input. This is the core principle behind decentralizing the dilemma. Ethics are, at their heart, a product of human interaction, consensus, and societal norms. Therefore, the training data for an AV's ethical compass should come directly from the collective wisdom and values of the people it serves. This decentralized approach acknowledges that there is no single "right" answer but rather a spectrum of ethically defensible positions, and the AV's behaviour in a crisis should ideally reflect a consensus derived from this diversity.

5. Fijishi Aeterna: A Framework for Adaptive Ethics

Imagine a framework where organizations developing AVs can move beyond preprogrammed rules and instead train their AI systems to understand and apply ethical reasoning derived from a broad cross-section of humanity. This is where hypothetical concepts like Fijishi and Aeterna come into play.

- Envision **Aeterna** as a platform or a set of tools that facilitates the collection and aggregation of diverse human ethical input. This could involve presenting a wide range of hypothetical accident scenarios to people from different demographics, cultures, and backgrounds. The platform captures not just their chosen action in a scenario (e.g., swerve or stay), but also their reasoning, their priorities, and the values that inform their decision. This data could be collected through interactive simulations, surveys, deliberative forums, or even by analysing how humans discuss and resolve ethical dilemmas in various contexts. Fijishi's role is to structure, anonymize, and synthesize this rich, qualitative ethical data into a format usable for AI training.
- Position Aeterna as an advanced AI training system specifically designed to learn ethical frameworks from the decentralized data provided by Fijishi. Unlike traditional machine learning that optimizes for performance metrics (like speed or safety under normal conditions), Aeterna is trained to understand the *patterns* and *priorities* embedded in the diverse human responses to ethical dilemmas. It doesn't learn rule A applies in situation X; it learns that, across a broad range of human input, there's a tendency to prioritize value Y over value Z in situations with characteristics P and Q. Aeterna learns *how* to weigh competing values and apply ethical principles based on the collective human consensus it has been trained on. Its "moral compass" is not a static algorithm but a continuously adapting model that refines its understanding as it processes new, diverse ethical input.

In this framework, AV developers use **Aeterna** to gather the ethical perspectives relevant to their target markets and then use **Aeterna** to train their specific AV AI models. The AI doesn't receive a rulebook; it receives a deeply nuanced understanding of the ethical landscape derived from thousands or millions of human inputs.

6. Benefits of Decentralization

This decentralized approach offers several significant advantages:

• Increased Societal Acceptance: When the public knows that the ethical decisions of AVs are informed by a broad range of human values, not just a

few programmers, trust is likely to increase. This transparency fosters greater acceptance and reduces fear surrounding the "unimaginable."

- **Reduced Bias:** By drawing from diverse inputs, the ethical framework learned by Aeterna is less likely to be skewed by the biases of a small, homogeneous group. It aims to reflect a more representative cross-section of societal values.
- **Robust and Representative Ethics:** The resulting ethical compass is more likely to be robust and applicable to a wider range of unforeseen scenarios because it is based on a deeper understanding of underlying ethical principles rather than brittle, scenario-specific rules.
- Adaptability: As societal values evolve, new data can be fed into Fijishi, and Aeterna can be retrained, allowing the AV's ethical compass to adapt over time. This ensures that autonomous vehicles remain aligned with contemporary ethical norms.
- Enhanced Accountability: While still complex, the decentralized model shifts the focus of accountability from a specific programmer's rule to the process of gathering and interpreting diverse human ethical input. This necessitates transparency in the data collection and training methodologies.

7. Conclusion

Autonomous vehicles present society with a profound ethical challenge – navigating unavoidable accident scenarios where there is no universally agreed-upon "right" answer. Relying on fixed, pre-programmed algorithms is an insufficient and potentially problematic approach that fails to capture the complexity and dynamic nature of human ethics.

The path forward lies in decentralizing the dilemma. By building systems that allow AI to learn ethical frameworks from the diverse inputs of the very people they will interact with, we can cultivate an ethical compass that is not a rigid algorithm but a continuously adapting consensus. Conceptual frameworks like **Aeterna** offer a glimpse into how organizations can facilitate this process, moving beyond imposed rules to a model of learned, representative ethics.

Embracing this decentralized paradigm is not just a technical challenge; it is a societal imperative. It requires collaboration between technologists, ethicists, social scientists, and the public. Only by grounding the AV's moral compass in the collective wisdom and evolving values of humanity can we build trust, ensure accountability, and navigate the "unimaginable" with greater confidence and ethical integrity.

This document is provided for information purposes only. This document is not warranted to be errorfree, nor subject to any other warranties or conditions, whether expressed orally or implied in law, including implied warranties and conditions of merchantability or fitness for a particular purpose. We specifically disclaim any liability with respect to this document and no contractual obligations are formed either directly or indirectly by this document. This document may not be reproduced or transmitted in any form or by any means, electronic or mechanical, for any purpose, without our prior written permission. To know more, please visit www.fijishi.com

©2025 Fijishi, and/or its affiliates. All rights reserved.