

AI Ethics Framework: A Quantum Perspective

By Walter Evans

Date: January 28, 2026

Abstract

This paper explores the intersection of artificial intelligence ethics and quantum cognition theory. We propose a novel framework for understanding ethical decision-making in AI systems through the lens of quantum probability theory. Our findings suggest that traditional binary ethical frameworks may be insufficient for capturing the nuanced nature of AI decision-making processes. This paper introduces a quantum-inspired model that embraces ambiguity and context-dependency in ethical judgments. Key findings include the application of quantum superposition as a model for ethical ambiguity, the exploration of entanglement effects in multi-agent AI systems, and the use of measurement problem analogies in ethical observation and auditing. We argue that this quantum perspective provides a more robust and psychologically realistic foundation for developing and governing ethical AI.

Keywords

Artificial Intelligence, AI Ethics, Quantum Cognition, Quantum Probability, Ethical Frameworks, Decision-Making, Multi-Agent Systems

1. Introduction

The rapid advancement **of** artificial intelligence (AI) has brought **to** the forefront a host **of** complex ethical challenges. As AI systems become increasingly autonomous **and** integrated **into** the fabric **of** society, the need **for** robust ethical frameworks **to** guide their development **and** deployment has become paramount [1]. Existing frameworks, largely derived **from** classical deontological **and** utilitarian ethics, often rely **on** binary logic **and** deterministic rules. These models, **while** valuable, struggle **to** account **for** the ambiguity, context-dependency, **and** inherent uncertainty that characterize many real-world ethical dilemmas [2]. The discrete, rule-based nature **of** these classical frameworks often fails **to** capture the more fluid **and** probabilistic nature **of** human moral intuition, a gap that becomes increasingly problematic **as** we strive **to** build AI that aligns **with** human values.

This paper introduces a novel approach **to** AI ethics **by** drawing upon the principles **of** quantum cognition theory. Quantum cognition **is** an emerging field that applies the mathematical formalism **of** quantum mechanics **to** model human judgment **and** decision-making, particularly **in** situations **where** classical probability theory fails **to** explain observed psychological phenomena [3]. By leveraging concepts such **as** superposition, entanglement, **and** measurement, quantum cognition offers a powerful **new** lens through which **to** understand **and** model the complexities **of** human thought. We argue that this quantum-inspired perspective can provide a more psychologically realistic **and** mathematically rigorous foundation **for** AI ethics.

We propose a framework **where** an AI's ethical state is not a single, definite value (e.g., "ethical" or "unethical") but rather a superposition of multiple potential ethical states, each with a certain probability. The act of auditing or observing the AI's behavior is analogous to a quantum measurement, which collapses this superposition into a single, definite outcome. This model allows for a more nuanced understanding of ethical ambiguity and the role of the observer in shaping ethical judgments. Furthermore, we explore the concept of entanglement to model the correlated ethical states of multi-agent AI systems, providing a new way to analyze systemic ethical risks and emergent behaviors.

2. Literature Review

2.1. Prevailing AI Ethics Frameworks

The discourse on AI ethics has produced a variety of frameworks and principles aimed at ensuring the responsible development of AI. A global convergence is emerging around five core principles: **transparency**, **justice and fairness**, **non-maleficence**, **responsibility**, and **privacy** [4]. Organizations such as UNESCO, the OECD, and IBM have all put forth guidelines that emphasize these principles in various forms [5] [6] [7]. For example, the principle of transparency calls for AI systems to be explainable and their decision-making processes to be understandable by humans. The principle of fairness seeks to prevent and mitigate bias in AI algorithms, ensuring that they do not perpetuate or amplify existing social inequalities.

While these principles provide essential high-level guidance, their practical implementation remains a significant challenge. Many of these frameworks are based on classical, deterministic logic, which assumes that ethical problems can be broken down into a set of clear rules and that a single "correct" answer can be derived. However, this approach often falls short when faced with the complexity and ambiguity of real-world ethical dilemmas, which often involve conflicting values and uncertain outcomes.

2.2. Fundamentals of Quantum Cognition

Quantum cognition **is** a research program that uses the mathematical principles **of** quantum theory **to** model human cognition [3]. It **is** important **to** note that quantum cognition does **not** claim that the brain **is** a quantum computer. Rather, it uses the *formalism* **of** quantum mechanics **as** a modeling tool, **in** much the same way that classical probability theory has been used **to** model rational choice. Quantum probability theory provides a different **set of rules for** calculating probabilities than classical probability theory, **and** these rules appear **to** be better suited **for** explaining a range **of** psychological phenomena, including violations **of** classical probability **in** decision-making [8].

Key concepts **from** quantum theory that are used **in** quantum cognition include:

- **Superposition:** **In** quantum mechanics, a particle can exist **in** a combination **of** multiple states at once. **In** quantum cognition, this **is** used **to** model the ambiguity **and** uncertainty **of** mental states. **For** example, before making a decision, a person may be **in** a superposition **of** multiple choices.
- **Entanglement:** This **is** a phenomenon **where** two **or** more quantum particles become linked **in** such a way that their fates are intertwined, no matter how far apart they are. **In** quantum cognition, this can model the correlated judgments **and** decisions **of** multiple individuals **or**, **as** we propose, multiple AI agents.
- **Measurement:** The act **of** measuring a quantum system collapses its superposition **into a single**, definite state. **In** quantum cognition, this **is** analogous **to** the act **of** making a decision **or** answering a question, which resolves the ambiguity **of** a mental state.

2.3. Quantum Approaches to Decision-Making and Ethics

There is a growing body of research that applies quantum cognition models to decision-making under uncertainty [9]. These models have been successful in explaining a variety of cognitive biases and paradoxical findings that are difficult to account for with classical models. For example, quantum models have been used to explain the order effects of questions in surveys and the conjunction fallacy, where people rate the probability of two events occurring together as higher than the probability of one of the events occurring alone.

More recently, researchers have begun to explore the application of quantum cognition to ethical decision-making [10]. It has been posited that quantum models can provide a more credible representation of human behavior in ethical dilemmas, which are often characterized by cognitive conflict and moral ambiguity [11]. These models can capture the way in which the act of deliberation can change a person's ethical preferences, and how different ethical principles can be in a state of superposition before a decision is made.

3. A Quantum-Inspired Framework for AI Ethics

Building on the foundations of quantum cognition, we propose a new framework for AI ethics that embraces the principles of superposition, entanglement, and measurement. This framework is not intended to replace existing ethical principles but rather to provide a new mathematical and conceptual language for their implementation and analysis.

3.1. Superposition and Ethical Ambiguity

In our framework, the ethical state of an AI system is represented by a state vector in a Hilbert space. Each basis vector of this space corresponds to a distinct ethical position or choice. Before a decision is made, the AI exists in a superposition of these basis states, with each state having a certain probability amplitude. For example, an autonomous vehicle facing a difficult choice in a potential accident scenario could be in a superposition of "swerve to hit obstacle A" and "continue to hit obstacle B".

The advantage of this approach is that it allows us to model the inherent ambiguity and conflict of ethical dilemmas in a mathematically precise way. Instead of forcing a binary choice at the outset, the superposition principle allows the AI to hold multiple, potentially contradictory, ethical considerations in a state of potentiality. The evolution of this superposition over time can be modeled using a Schrödinger-like equation, which can incorporate various factors such as new information, changing context, and the AI's own internal ethical principles.

3.2. Entanglement and Multi-Agent Systems

The concept of entanglement provides a powerful tool for modeling the ethical behavior of multi-agent AI systems. When two or more AI agents become entangled, their ethical states become correlated, even if they are not directly communicating. This means that a decision made by one agent can instantaneously affect the ethical state of the other agents in the system. For example, a team of autonomous drones on a search and rescue mission could have their ethical states entangled, leading to a coordinated and collectively ethical response to a complex situation.

This has significant implications for understanding systemic ethical risks. A group of interacting AI agents, such as in a high-frequency trading system, could become entangled in a way that leads to emergent, and potentially unethical, collective behavior that would not be predictable from the behavior of any single agent in isolation. Our framework provides a way to analyze and potentially mitigate these systemic risks by studying the entanglement properties of the multi-agent system.

3.3. Measurement and Ethical Observation

The act **of** auditing, judging, **or** even simply observing an AI's ethical decision is analogous to a quantum measurement. This measurement collapses the AI's superposition of ethical states into a single, definite outcome. This has a profound implication: the observer is not a passive bystander but an active participant in the creation of the ethical reality they are observing. The way in which a question is framed, the metrics that are used for an audit, and the context in which a decision is reviewed all act as a measurement basis that influences the outcome.

This perspective aligns **with** the growing recognition **in** the field **of** AI ethics that context **and** framing are crucial. It suggests that we cannot simply build an "ethical AI" **in** a vacuum **and then** expect it **to** behave ethically **in** all situations. Instead, we must design the entire socio-technical system, including the human observers **and** auditors, **in** a way that promotes ethical outcomes. Our framework provides a mathematical language **for** talking about the co-creation **of** ethical reality **by** humans **and** AI.

4. Implications and Future Research

The quantum-inspired framework for AI ethics proposed in this paper has several important implications for the development, governance, and auditing of AI systems. By providing a more nuanced and psychologically realistic model of ethical decision-making, it can help us to build AI that is more aligned with human values. For example, AI systems could be designed to maintain a state of ethical superposition for as long as possible, allowing for more flexible and context-sensitive responses to complex situations.

In terms of governance, our framework suggests that we need to move beyond simple rule-based approaches and develop more sophisticated methods for auditing and overseeing AI systems. This could involve the use of "quantum-like" simulations to explore the full range of an AI's potential ethical behaviors, rather than just testing it on a limited set of pre-defined scenarios. The concept of entanglement also highlights the need for a more systemic approach to AI governance, which takes into account the interactions between multiple AI agents and their environment.

Future research should focus on developing the mathematical and computational tools needed to implement this framework in practice. This will involve further work on quantum probability theory, as well as the development of new algorithms for simulating and analyzing quantum-like cognitive models. Empirical research is also needed to test the predictions of this framework and to explore its applicability to a wide range of AI systems and ethical dilemmas.

5. Conclusion

The challenge of building ethical AI is one of the most pressing issues of our time. The classical, deterministic frameworks that have dominated the field of AI ethics are increasingly showing their limitations in the face of the complexity and ambiguity of real-world ethical problems. This paper has proposed a new approach, inspired by the principles of quantum cognition, that embraces this ambiguity and provides a more robust foundation for AI ethics.

By modeling the ethical states of AI systems as superpositions, exploring the entanglement of multi-agent systems, and framing ethical observation as a form of measurement, our framework offers a new set of conceptual and mathematical tools for reasoning about AI ethics. While much work remains to be done, we believe that this quantum perspective has the potential to significantly advance our ability to create AI that is not only intelligent but also wise.

References

- [1] Floridi, L., & Cowls, J. (2018). AI4People—An ethical framework for a good AI society: Opportunities, risks, principles, and recommendations. *Minds and Machines*, 28(4), 689-707.
- [2] Whittaker, M., et al. (2018). *AI Now Report 2018*. AI Now Institute.
- [3] Busemeyer, J. R., & Bruza, P. D. (2012). *Quantum models of cognition and decision*. Cambridge University Press.
- [4] Jobin, A., Ienca, M., & Vayena, E. (2019). The global landscape of AI ethics guidelines. *Nature Machine Intelligence*, 1(9), 389-399.
- [5] UNESCO. (2021). *Recommendation on the Ethics of Artificial Intelligence*. UNESCO.
- [6] OECD. (2019). *Recommendation of the Council on Artificial Intelligence*. OECD/LEGAL/0449.
- [7] IBM. (n.d.). *AI Ethics*. Retrieved from <https://www.ibm.com/think/topics/ai-ethics>

- [8] Pothos, E. M., & Busemeyer, J. R. (2013). Can quantum probability provide a new direction for cognitive modeling?. *Behavioral and Brain Sciences*, 36(3), 255-274.
- [9] Khrennikov, A. (2018). Quantum probability in decision making from quantum information to quantum-like modeling in biology, economics, and psychology. *Scientific reports*, 8(1), 1-13.
- [10] Yilmaz, L. (2020). A quantum cognition model for simulating ethical dilemmas among multi-perspective agents. *Journal of Simulation*, 14(3), 195-209.
- [11] White, L., & Pothos, E. M. (2017). Quantum cognition models of ethical decision-making. In *Proceedings of the 15th International Conference on Cognitive Modeling*.