

1. Introduction (1p)

(a) Problem

(b) Contributions:

- i. In Section 4.a, I make a philosophical argument for why Kantian ethics is the most natural of the three major ethical traditions (deontology, virtue ethics, utilitarianism) to formalize.
- ii. In Section 4.b, I present a formalization of the practical contradiction interpretation of the Formula of Universal Law in Dyadic Deontic Logic.
- iii. In Section 4.c, I present an implementation of the above formalization in the Isabelle/HOL theorem prover. My implementation includes axioms and definitions such that my system, when given an appropriately represented input, can prove that the input is permissible, obligatory, or prohibited. It can also return a list of facts used in the proof and, in some cases, an Isar-style human readable proof. Because my system is equipped with a sophisticated definition of a maxim, I show that it improves on prior work by performing reasoning nuanced enough to distinguish between a joke and a lie.
- iv. In Section 4.d, I present a testing framework that can generally evaluate how faithful an implementation of automated Kantian ethics is. My framework includes meta-ethical tests and application tests inspired by philosophical literature. My testing framework shows that my formalization substantially improves on prior work and can be generalized to any implementation of automated Kantian ethics.
- v. In Section 4.e, I demonstrate new ethical insights discovered using my system and argue that computational methods like the one presented in this paper can help philosophers address ethical problems.

2. The Problem (1p)

- (a) In order for AI agents to navigate the world responsibly, they need to perform ethical reasoning.
- (b) This ethical reasoning should draw on existing work in philosophy and should be explainable.
- (c) Maybe put something about computational ethics here?

3. My Idea (2p)

- (a) Overview of the Top-Down, Theorem Prover approach
- (b) Testing Framework (introduce the idea, metaethical tests, application tests)
- (c) How it would fit into (a) an AI agent (b) philosophers' workflow

4. The Details (5p)

- (a) Why Kantian Ethics (1-1.5p)
- (b) Isabelle/HOL implementation (basic logical background and structure)
- (c) Features: robust definition of a maxim, ability to distinguish between lies and jokes, reusability of testing framework (maybe include the table showing which tests my implementation passes vs prior implementations?)
- (d) Philosophical insights discovered along the way (vacuous maxims example)
- (e) Limitations (needs an input parser and a common sense database)

5. Related Work (1-2p)

- (a) Talk about bottom-up vs top-down approaches
- (b) I improve prior work by (1) staying faithful to philosophical literature (2) building an explainable system

6. Conclusion (0.5p)