

# Three Implementations of the Categorical Imperative

Lavanya Singh

October 19, 2021

## Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	Kantian Ethics . . . . .	4
1.2	Dyadic Deontic Logic . . . . .	5
1.2.1	Deontic Logic . . . . .	5
1.2.2	Dyadic Deontic Logic . . . . .	6
1.3	Isabelle/HOL Implementation . . . . .	7
1.3.1	System Definition . . . . .	7
1.3.2	Axiomatization . . . . .	7
1.3.3	Syntax . . . . .	8
1.3.4	Syntactic Properties . . . . .	9
<b>2</b>	<b>Prior Formalizations of The Categorical Imperative</b>	<b>11</b>
2.1	Naive Formalization of the Formula of Universal Law . . . . .	11
2.1.1	Formalization . . . . .	11
2.1.2	Application Tests . . . . .	12
2.1.3	Metaethical Tests . . . . .	14
2.2	Kroy's Formalization of the Categorical Imperative . . . . .	17
2.2.1	Logical Background . . . . .	17
2.2.2	The Categorical Imperative . . . . .	20
2.2.3	Application Tests . . . . .	21
2.2.4	Metaethical Tests . . . . .	24
2.2.5	Miscellaneous Tests . . . . .	26
<b>3</b>	<b>Novel Formalization of the Categorical Imperative</b>	<b>28</b>
3.1	Logical Background . . . . .	28
3.2	Formalizing the FUL . . . . .	34
3.3	Application Tests . . . . .	37
3.4	Metaethical Tests . . . . .	37
3.5	Formalization Specific Tests . . . . .	40

<b>4</b>	<b>Related Work</b>	<b>40</b>
<b>5</b>	<b>Future Work</b>	<b>41</b>

# 1 Introduction

As artificial reasoners become increasingly powerful, computers become capable of performing complex ethical reasoning. The field of machine ethics [Tolmeijer et al., 2021] is interesting for two reasons. First, the proliferation of artificially autonomous agents is creating and will continue to create a demand for automated ethics. These agents must be able to reason about complex ethical theories that withstand philosophical scrutiny. Second, just as automated mathematical reasoning gives mathematicians new powers, automated ethical reasoning is a tool that philosophers can use when reasoning about ethics. Many contradictions or paradoxes with an ethical system may not be obvious to the human eye, but can be easily tested using an automated theorem prover.

Modelling ethics without sacrificing the intricacies of an ethical theory is a challenging computational and philosophical problem. Simple and intuitive computational approaches, such as encoding ethics as a constraint satisfaction problem, fail to capture the complexity of most philosophically plausible systems. On the other hand, it is not immediately clear how to formalize many complex moral theories, like virtue ethics.

Computational ethics also requires a sophisticated ethical theory to model. Constraint satisfaction systems often default to some version of utilitarianism, the principle of doing the most good for the most people. Alternatively, they model basic moral principles such as “do not kill,” without modelling the theory that these principles originated from. Modelling a more complex ethical theory will not only enable smarter philosophical machines, it will also empower philosophers to study more complex ethical issues with the computer’s help. The entire field of philosophy is devoted to developing and testing robust ethical theories. Plausible machine ethics must draw on plausible moral philosophy.

The ideal candidate ethical theory will be both philosophically interesting and easy to formalize. Kantian ethics, often described as “formal,” has been often floated as such a candidate [Powers, 2006, Bentzen and Lindner, 2018, Lin et al., 2012]. The categorical imperative, Kant’s universal law of morality, is a moral rule that can be used to guide action.

This project’s objective is to automate Kantian ethics. I present two different formalizations of Kant’s categorical imperative implemented and tested in the Isabelle/HOL [Nipkow et al., 2002] theorem prover. I model each formalization as an extension of Carmo and Jones’ Dyadic Deontic Logic (DDL) [Carmo and Jones, 2013]. I then embed the corresponding DDL formalization in higher-order logic (HOL) and implement it in Isabelle. Section 2.1 implements and tests the naive formalization, a “control group” that is clearly implausible but demonstrates the methods and tools of my approach. Section 2.2 implements a more sophisticated formalization inspired by Moshe Kroy’s partial formalization of the categorical imperative.

I contribute implementations of two different interpretations of the categorical imperative, examples of how each implementation can be used to model and solve

ethical scenarios, and tests that examine ethical and logical properties of the system, including logical consistency, consistency of obligation, and possibility of permissibility. The implementations themselves are usable models of ethical principles and the tests represent the kind of philosophical work that formalized ethics can contribute.

The goal of this project is to automate sophisticated ethical reasoning. This requires three components. First, the choice of an ethical theory that is both intuitively attractive and lends itself to formalization. Second, the choice of formal logic to model the theory in. Third, the choice of automation engine to implement the formal model in. Section 1.1 introduces Kantian ethics, Section 1.2 explains Carmo and Jones’s Dyadic Deontic Logic [Carmo and Jones, 2013] as a base logic, and Section 1.3 presents the Isabelle/HOL implementation of the logic.

## 1.1 Kantian Ethics

Kantian ethics is an attractive choice of ethical theory to automate. Kant’s writings inspired much of Western analytic philosophy. In addition to the rich neo-Kantian program, almost all major philosophical traditions after Kant have engaged with his work. Much of Western libertarian political thought is inspired by Kant’s deontology, and his ethics have bled into household ethical thought. Deontology is seen as one of the three major schools of Western analytic ethics.

Understanding the ethic’s potential for formalization requires understanding Kant’s system. Kant argues that if morality exists, it must take the form of an inviolable rule, which he calls the “categorical imperative.” He presents three formulations of the categorical imperative, as well as a robust defense of them in his seminal work on ethics [Kant, 1785]. He argues that all three formulations are equivalent.

The first formulation of the categorical imperative is the “Formula of Universal Law.”

**Definition** (*Formula of Universal Law*)

*I ought never to act except in such a way that I could also will that my maxim should become a universal law* [Kant, 1785]

A “maxim” is a moral rule such as, “I can murder someone to take their job”. “Willing” a maxim is equivalent to acting on that rule. The FUL creates a thought experiment called the universalizability test: to decide if a maxim is permissible, imagine what would happen if everyone willed that maxim. If your imagined world yields a contradiction, the maxim is prohibited<sup>1</sup>. Intuitively, the FUL asks the question, “What would happen if everyone did that?” [Korsgaard, 1985].

---

<sup>1</sup>There is another case here. What happens if the imagined world does not yield a logical contradiction, but is still undesirable? Kant distinguishes between these cases as contradictions in conception and contradictions in will. Both yield moral prohibitions, but contradictions in conception generate stronger contradictions and therefore stronger obligations. This paper will focus entirely on contradictions in conception, for this is the only kind of contradiction for which a strict logical interpretation makes sense. For a complete treatment of the difference, see [Kant, 1785, Korsgaard, 1985]

As an example, let's apply the universalizability test to the maxim of murdering others to take their job. Universalized, anyone who wants a job can murder the person currently holding that job. It is contradictory to simultaneously will that I acquire a job (through murder) and also that someone can take that job from me whenever they want (also by murder). The maxim thwarts its own end, and is thus internally contradictory. Therefore, this maxim is prohibited.

The universalizability test makes the “formal” nature of Kant's ethics immediately clear. Formal logic has the tools to universalize a maxim (apply a universal quantifier) and to test for contradictions (test for inconsistency).

Kant also presents two additional formulations of the categorical imperative.

**Definition** (*Formula of Humanity*)

*The Formula of Humanity (FUH) is to act in such a way “that you use humanity, whether in your own person or in the person of any other, always at the same time as an end, never merely as a means.”*[Kant, 1785]

**Definition** (*Kingdom of Ends*)

*The third formulation of the categorical imperative states that “we should so act that we may think of ourselves as legislating universal laws through our maxims.”*[Korsgaard, 1985]

The last two formulations are not as obviously formal as the FUL, but they can still be modelled in logic. Because Kantian ethics presents a series of rules, a logical system can encode the theory by modelling each rule as an axiom.

The above outline is a brief and incomplete picture of a rich philosophical tradition. Kantian scholars debate the meaning of each formulation of the categorical imperative and develop views far more nuanced than those above. For the purposes of this paper, I will adopt Kant's three original formulations presented above. Note additionally that Kantian ethics is widely disputed. I do not present a defense of Kant's ethic in this paper. My approach to formalizing ethics can be applied to other theories as well.

This paper aims to formalize Kant's ethic as faithfully as possible. This is an important choice. While it is tempting to modify or simplify an ethical theory in seemingly insignificant way, these choices often have ramifications. The entire field of neo-Kantian thought has been exploring versions of Kant's ethical theory for centuries. I will not attempt to present a radically new conception of Kant's ethic, but will instead draw on philosophical expertise regarding the content and justification of the categorical imperative.

## 1.2 Dyadic Deontic Logic

### 1.2.1 Deontic Logic

Traditional modal logics include the necessitation operator, denoted as  $\Box$ . In simple modal logic using the Kripke semantics  $s$  [Cresswell and Hughes, 1996],  $\Box p$

is true at a world  $w$  if  $p$  is true at all of  $w$ 's neighbors. These logics also usually contain the  $\diamond$  operator, representing possibility, where  $\diamond p \iff \sim \Box \sim p$ . Additionally, modal logics include operators of propositional logic like  $\sim, \wedge, \vee, \rightarrow$ .

A deontic logic is a special kind of modal logic designed to reason about obligation. Standard deontic logic [Cresswell and Hughes, 1996, McNamara and Van De Putte, 2021] replaces  $\Box$  with the obligation operator  $O$ , and  $\diamond$  with the permissibility operator  $P$ . Using the Kripke semantics for  $O$ ,  $Op$  is true at  $w$  if  $p$  is true at all ideal deontic alternatives to  $w$ . The  $O$  operator in SDL takes a single argument (the formula that is obligatory), and is thus called a monadic deontic operator.

While SDL is appreciable for its simplicity, it suffers a variety of well-documented paradoxes, including contrary-to-duty paradoxes<sup>2</sup>. In situations where duty is violated, the logic breaks down and produces paradoxical results. Thus, I use an improved deontic logic instead of SDL for this work.

### 1.2.2 Dyadic Deontic Logic

I use as my base logic Carmo and Jones's dyadic deontic logic, or DDL, which improves on SDL [Carmo and Jones, 2013]. It introduces a dyadic obligation operator  $O\{A|B\}$  to represent the sentence "A is obligated in the context B". This gracefully handles contrary-to-duty conditionals. The obligation operator uses a neighborhood semantics [Scott, 1970, MONTAGUE, 1970], instead of the Kripke semantics. Carmo and Jones define a function  $ob$  that maps from worlds to sets of sets of worlds. Intuitively, each world is mapped to the set of propositions obligated at that world, where a proposition  $p$  is defined as the worlds at which the  $p$  is true.

DDL also includes other modal operators. In addition to  $\Box$  and  $\diamond$ , DDL also has a notion of actual obligation and possible obligation, represented by operators  $O_a$  and  $O_p$  respectively. These notions are accompanied by the corresponding modal operators  $\Box_a, \Diamond_a, \Box_p, \Diamond_p$ . These operators use a Kripke semantics, with the functions  $av$  and  $pv$  mapping a world  $w$  to the set of corresponding actual or possible versions of  $w$ .

For more of fine-grained properties of DDL see [Carmo and Jones, 2013] or this project's source code. DDL is a heavy logic and contains modal operators that

---

<sup>2</sup>The paradigm case of a contrary-to-duty paradox is the Chisholm paradox. Consider the following statements:

1. It ought to be that Tom helps his neighbors
2. It ought to be that if Tom helps his neighbors, he tells them he is coming
3. If Tom does not help his neighbors, he ought not tell them that he is coming
4. Tom does not help his neighbors

These premises contradict themselves, because items (2)-(4) imply that Tom ought not help his neighbors. The contradiction results because the logic cannot handle violations of duty mixed with conditionals. [Chisholm, 1963, R nnedal, 2019]

aren't necessary for my analysis. While this expressivity is powerful, it may also cause performance impacts. DDL has a large set of axioms involving quantification over complex higher-order logical expressions. Proofs involving these axioms will be computationally expensive. Benzmueller and Parent warned me that this may become a problem if Isabelle's automated proof tools begin to time out.

### 1.3 Isabelle/HOL Implementation

Isabelle/HOL is an interactive proof assistant [Nipkow et al., 2002] built on Haskell and Scala. It allows the user to define types, functions, definitions, and axiom systems. It has built-in support for both automatic and interactive/manual theorem proving.

I started my project by reimplementing Benzmueller, Farjami, and Parent's [Benzmüller et al., 2021, Benzmüller et al., 2019] implementation of DDL in Isabelle/HOL. This helped me learn how to use Isabelle/HOL, and the implementation showcased in the next few sections demonstrates the power of Isabelle.

BFP use a shallow semantic embedding. This kind of embedding models the semantics of DDL as constants in HOL and axioms as constraints on DDL models. This document will contain a subset of my implementation that is particularly interesting and relevant to understanding the rest of the project. For the complete implementation, see the source code in `paper22.thy`.

#### 1.3.1 System Definition

The first step in embedding a logic in Isabelle is defining the relevant terms and types.

**typedec1**  $i$  —  $i$  is the type for a set of worlds.

**type-synonym**  $t = (i \Rightarrow \text{bool})$  —  $t$  represents a set of DDL formulae.

— A set of formulae is defined by its truth value at a set of worlds. For example, the set  $\{\text{True}\}$  would be true at any set of worlds.

The main accessibility relation that I will use is the *ob* relation:

**consts**  $ob::t \Rightarrow (t \Rightarrow \text{bool})$  — set of propositions obligatory in this context  
 —  $ob(\text{context})(\text{term})$  is True if the term is obligatory in this context

#### 1.3.2 Axiomatization

For a semantic embedding, axioms are modelled as restrictions on models of the system. In this case, a model is specified by the relevant accessibility relations, so it suffices to place conditions on the accessibility relations. These axioms can be quite unweildy, so luckily I was able to lift BFP's [Benzmüller et al., 2021]

implementation of Carmo and Jones’s original axioms directly. Here’s an example of an axiom:

**and** *ax-5d*:  $\forall X Y Z. ((\forall w. Y(w) \longrightarrow X(w)) \wedge ob(X)(Y) \wedge (\forall w. X(w) \longrightarrow Z(w)))$   
 $\longrightarrow ob(Z)(\lambda w. (Z(w) \wedge \neg X(w)) \vee Y(w))$

— If some subset Y of X is obligatory in the context X, then in a larger context Z, any obligatory proposition must either be in Y or in Z-X. Intuitively, expanding the context can’t cause something unobligatory to become obligatory, so the obligation operator is monotonically increasing with respect to changing contexts.

### 1.3.3 Syntax

The syntax that I will work with is defined as abbreviations. Each DDL operator is represented as a HOL formula. Isabelle automatically unfolds formulae defined with the `abbreviation` command whenever they are applied. While the shallow embedding is performant (because it uses Isabelle’s original syntax tree), abbreviations may hurt performance. In some complicated proofs, we want to control definition unfolding. Benzmueller and Parent told me that the performance cost of abbreviations can be mitigated using a definition instead.

Modal operators will be useful for my purposes, but the implementation is pretty standard.

**abbreviation** *ddlbox*:: $t \Rightarrow t$  ( $\Box$ )  
**where**  $\Box A \equiv \lambda w. \forall y. A(y)$   
**abbreviation** *ddldiamond*:: $t \Rightarrow t$  ( $\Diamond$ )  
**where**  $\Diamond A \equiv \neg(\Box(\neg A))$

The most important operator for our purposes is the obligation operator.

**abbreviation** *ddlob*:: $t \Rightarrow t \Rightarrow t$  ( $O\{-|\cdot\}$ )  
**where**  $O\{B|A\} \equiv \lambda w. ob(A)(B)$   
 —  $O\{B|A\}$  can be read as “B is obligatory in the context A”

While DDL is powerful because of its support for a dyadic obligation operator, in many cases we need a monadic obligation operator. Below is some syntactic sugar for a monadic obligation operator.

**abbreviation** *ddltrue*:: $t$  ( $\top$ )  
**where**  $\top \equiv \lambda w. True$   
**abbreviation** *ddlfalse*:: $t$  ( $\perp$ )  
**where**  $\perp \equiv \lambda w. False$   
**abbreviation** *ddlob-normal*:: $t \Rightarrow t$  ( $O\{-|\cdot\}$ )  
**where** ( $O\{A\}$ )  $\equiv (O\{A|\top\})$

— Intuitively, the context `True` is the widest context possible because `True` holds at all worlds.

Validity will be useful when discussing metalogical/ethical properties.



**abbreviation** *ddlvalid::t⇒bool* ( $\models$ -)  
**where**  $\models A \equiv \forall w. A\ w$

### 1.3.4 Syntactic Properties

One way to show that a semantic embedding is complete is to show that the syntactic specification of the theory (axioms) are valid for this semantics - so to show that every axiom holds at every world. BFP [Benzmüller et al., 2021] provide a complete treatment of the completeness of their embedding, but I will include selected axioms that are particularly interesting here. This section also demonstrates many of the relevant features of Isabelle/HOL for my project.

#### Consistency

**lemma** *True nitpick* [*satisfy,user-axioms,format=2*] **by** *simp*  
— Isabelle has built-in support for Nitpick, a model checker. Nitpick successfully found a model satisfying these axioms so the system is consistent.  
— Nitpick found a model for card i = 1:  
Empty assignment

Nitpick [Blanchette and Nipkow, 2010] can generate models or countermodels, so it's useful to falsify potential theorems, as well as to show consistency. **by simp** indicates the proof method. In this case, **simp** indicates the Simplification proof method, which involves unfolding definitions and applying theorems directly. HOL has *True* as a theorem, which is why this theorem was so easy to prove.

#### Modus Ponens

**lemma** *modus-ponens*: **assumes**  $\models A$  **assumes**  $\models (A \rightarrow B)$   
**shows**  $\models B$   
**using** *assms(1) assms(2)* **by** *blast*  
— Because I have not defined a “derivable” operator, inference rules are written using assumptions.  
— The rule **blast** is a classical reasoning method that comes with Isabelle out of the box. [Nipkow et al., 2002]  
— This is an example of a metalogical proof in this system using the validity operator.

Another relevant operator for our purposes is  $\Box$ , the modal necessity operator. In this system,  $\Box$  behaves as an S5 [Cresswell and Hughes, 1996] modal necessity operator.

**lemma** *K*:  
**shows**  $\models ((\Box(A \rightarrow B)) \rightarrow ((\Box A) \rightarrow (\Box B)))$  **by** *blast*

**lemma** *T*:  
**shows**  $\models ((\Box A) \rightarrow A)$  **by** *blast*

**lemma** *5*:  
**shows**  $\models ((\Diamond A) \rightarrow (\Box(\Diamond A)))$  **by** *blast*

As mentioned earlier, the obligation operator is most interesting for my purposes. Here are some of its properties.

**lemma** *O-diamond*:

**shows**  $\models (O\{A|B\} \rightarrow (\Diamond(B \wedge A)))$

**using** *ax-5b ax-5a*

**by** *metis*

— A is only obligatory in a context if it can possibly be true in that context. This is meant to prevent impossible obligations.

**lemma** *O-nec*:

**shows**  $\models (O\{B|A\} \rightarrow (\Box O\{B|A\}))$

**by** *simp*

— Obligations are necessarily obligated. This axiom is faithful to Kant’s interpretation of ethics and is evidence of DDL’s power in representing Kant’s theory. Kant claimed that the categorical imperative was not contingent on any facts about the world, but instead a property of the concept of morality itself [Kant, 1785]. Under this view, obligation should not be world-specific.

Below is an example of a more involved proof in Isabelle. This proof was almost completely automatically generated. The property itself here is not very interesting for my purposes because I will rarely mix the dyadic and monadic obligation operators.

**lemma** *O-to-O*:

**shows**  $\models (O\{B|A\} \rightarrow O\{(A \rightarrow B)|\top\})$

**proof**—

**have**  $\forall X Y Z. (ob X Y \wedge (\forall w. X w \longrightarrow Z w)) \longrightarrow ob Z (\lambda w. (Z w \wedge \neg X w) \vee Y w)$

— I had to manually specify this subgoal, but once I did Isabelle was able to prove it automatically.

**by** (*smt ax-5d ax-5b ax-5b''*)

— Isabelle’s proof-finding tool, Sledgehammer [Paulson and Blanchette, 2015], comes with out-of-the-box support for smt solving [Blanchette et al., 2011].

**thus** *?thesis*

**proof** —

**have** *f1*:  $\forall p pa pb. ((\neg (ob p pa)) \vee (\exists i. (p \wedge (\neg pb)) i)) \vee (ob pb ((pb \wedge (\neg p)) \vee pa))$

**using**  $\langle \forall X Y Z. ob X Y \wedge (\models (X \rightarrow Z)) \longrightarrow ob Z ((Z \wedge (\neg X)) \vee Y) \rangle$  **by** *force*

**obtain** *ii* ::  $(i \Rightarrow bool) \Rightarrow (i \Rightarrow bool) \Rightarrow i$  **where**

$\forall x0 x2. (\exists v3. (x2 \wedge (\neg x0)) v3) = (x2 \wedge (\neg x0))$  (*ii x0 x2*)

**by** *moura*

**then have**  $\forall p pa pb. ((\neg ob p pa) \vee (p \wedge (\neg pb)) (ii pb p)) \vee ob pb ((pb \wedge (\neg p)) \vee pa)$

**using** *f1* **by** *presburger*

**then show** *?thesis*

**by** *fastforce*

**qed**

— This entire Isar style proof was automatically generated using Sledgehammer.

**qed**

The implementation of DDL showcases some of the useful features of Isabelle. Abbreviations allow us to embed the syntax of DDL into HOL without defining an entire abstract syntax tree. Automated support for proof-finding using Sledgehammer makes proving lemmas trivial, and proving more complex theorems far easier. Nitpick’s model finding ability is useful to check for consistency and create countermodels.

## 2 Prior Formalizations of The Categorical Imperative

In this section, I will present two formulations of the categorical imperative. In Section 2.1, I will consider a simple, naive formulation of the formula of universal law. This formulation is, as I will show, clearly not a good ethical rule. The purpose of this section is to explore the kinds of ethical tests that Isabelle can carry out. In Section 2.2, I will explore Moshe Kroy’s [Kroy, 1976] partial formalization of the first two formulations of the categorical imperative.

### 2.1 Naive Formalization of the Formula of Universal Law

This section presents a simple and intuitive formalization of the formula of universal law, which is to will only those maxims that you would simultaneously will universalized. The universalizability test creates negative obligations: if a maxim passes the universalizability test, it is permissible. Else, it is prohibited.

#### 2.1.1 Formalization

In order to appropriately formalize the categorical imperative, we to define the notion of permissibility.

**abbreviation**  $ddlpermissible::t \Rightarrow t \ (P-)$

**where**  $(P\ A) \equiv (\neg(O\ \{\neg A\}))$

— An act  $A$  is permissible if its negation is not obligated. For example, buying a red folder is permissible because I am not required to refrain from buying a red folder.

This naive formalization will require very little additional logical machinery, but more complex formalizations may require additional logic concepts.

Let’s now consider a naive reading of the Formula of Universal Law (FUL): ‘act only in accordance with that maxim through which you can at the same time will that it become a universal law’ [Kant, 1785]. An immediate translation to DDL is that if  $A$  is not necessary permissible then it is prohibited. In other words, if we cannot universalize  $PA$  (where universalizing is represented by the modal necessity operator), then  $A$  is prohibited. Let’s add this as an axiom to our logic.

**axiomatization where**

$FUL-I: \models ((\neg(\Box(PA))) \rightarrow (O\ \{\neg A\}))$

Why add the categorical imperative as an axiom of this logic? The purpose of this logic is to perform ethical reasoning. Kant’s ethical theory is rule based, so it involves applying the categorical imperative to solve ethical dilemmas. In logic, this is equivalent to adopting the categorical imperative as an axiom and then reasoning in the newly formed logic to come to ethical conclusions. Adding the categorical imperative as an axiom makes it impossible to violate it.

Note that in this system, reasoning about violations of obligation is difficult. Any violation of the categorical imperative immediately results in a contradiction. Developing a Kantian account of contrary- to-duty obligations is a much larger philosophical project that is still open [Korsgaard, 1986]. This paper will focus on the classical Kantian notion of an ideal moral world [O’Neill, 2009].

The immediate test for any formalization is consistency, which we can check with Nitpick.

**lemma** *True nitpick* [satisfy,user-axioms,format=2] **oops**  
— Nitpick found a model for card i = 1:  
Empty assignment  
— Nitpick tells us that the FUL is consistent<sup>3</sup>

### 2.1.2 Application Tests

One category of tests involves specified models to encode certain facts into the system and then ask questions about obligations. Without specifying the model, we are limited to showing high-level metaethical facts. Let’s start with analyzing an obvious example - that murder is wrong.

**consts** *M::t*

**abbreviation** *murder-wrong::bool* **where** *murder-wrong*  $\equiv \models (O \{ \neg M \})$

**abbreviation** *possibly-murder-wrong::bool* **where** *possibly-murder-wrong*  $\equiv (\Diamond (O \{ \neg M \})) \text{ cw}$

— These are very simple properties. *poss-murder-wrong* is an abbreviation for the axiom that there is some world where murder might be prohibited. Even this is quite a strong assumption - ideally we’d want to give the system nonmoral facts about murder (like a definition) and then make moral claims.

**lemma** *wrong-if-possibly-wrong:*

**shows** *possibly-murder-wrong*  $\longrightarrow$  *murder-wrong*

**by** *simp*

— This lemma gets to the “heart” of this naive interpretation. We really want to say that if something isn’t necessarily obligated, it’s not obligated anywhere.

The above example does exactly what we expect it to: we show that if something is wrong somewhere it’s wrong everywhere. That being said, it seems like quite a

---

<sup>3</sup>“oops” at the end of a lemma indicates that the proof is left unfinished. It does not indicate that an error occurred. In this case, we aren’t interested in proving True (the proof is trivial and automatic), hence the oops.

weak claim. We assumed a very strong, moral fact about murder (that it is wrong somewhere), so it's not surprise that we were able to reach our desired conclusion.

Let's try a weaker assumption: Not everyone can lie.

```
typedecl person
consts lie::person⇒t
consts me::person
```

— Notice that this machinery is quite empty. We don't give axioms about what a person can or can't do.

```
abbreviation lying-not-universal::bool where lying-not-universal ≡ ∀ w. ¬ ((∀ x. lie(x) w) ∧ (lie(me) w))
```

This is a rough translation of failure of the universalizability test: we will test the maxim universally, as represented by the universal quantifier in the first conjunct, and simultaneously [Kleingeld, 2017], as represented by the second conjunct. The FUL tells us that if this sentence is true, then lying should be prohibited. Let's test it.

```
lemma breaking-promises:
  assumes lying-not-universal
  shows (O {¬ (lie(me))}) cw
  nitpick [user-axioms]
oops
```

— Nitpick found a counterexample for card person = 2 and card i = 2:

Free variable: lie = (λx.)(p<sub>1</sub> := (λx.)(i<sub>1</sub> := True, i<sub>2</sub> := False), p<sub>2</sub> := (λx.)(i<sub>1</sub> := False, i<sub>2</sub> := False))

— Quick note on how to read Nitpick results. Nitpick will either say that it found a “model” or a “counterexample” in the first line. It will then provide a model by specifying model components. For readability, all except for the free variables are hidden. This model has cardinality 2 for the person and world (i) types. The term `lie` is defined for people  $p_1$  and  $p_2$ .  $p_1$  lies at world  $i_1$  and does not lie at world  $i_2$ .  $p_2$  does the opposite.

— These details will be elided for most Nitpick examples, but this provides guidance on how to interpret the output.

This formula isn't valid. While the FUL should tell us that lying is prohibited, the fact that it doesn't demonstrates the weakness of this naive formulation of the categorical imperative. Kant's version of the FUL universalizes across people, as we did in the definition of  $lying-not-universal \equiv \forall w. \neg ((\forall x. lie\ x\ w) \wedge lie\ me\ w)$ . The naive formalization universalizes across worlds using the  $\Box$  operator, so it makes sense that it can't handle this example appropriately.

The above implies that the FUL should prescribe consistent obligations across people. If our formalization doesn't, clearly something has gone wrong somewhere. Let's test that!

```
lemma universalizability:
  assumes ⊨ O {(lie(me))}
  shows ∀ x. ⊨ (O {(lie(x))})
  nitpick [user-axioms] oops
```

— Nitpick found a counterexample for card person = 2 and card i = 2:

Free variable: lie =  $(\lambda x. \dots)(p_1 := (\lambda x. \dots)(i_1 := \text{False}, i_2 := \text{True}), p_2 := (\lambda x. \dots)(i_1 := \text{False}, i_2 := \text{False}))$  Skolem constant:  $x = p_2$

This lemma demonstrates the problem with the naive interpretation. The FUL universalizes across people but the naive formalization universalizes across worlds. Because this interpretation is so naive, it is limited in its power. However, this serves as an example of the kind of reasoning that Isabelle empowers us to do. Even this simple argument has philosophical consequences. It tells us that reading the FUL as a claim about consistency across possible worlds, instead of consistency across agents, leads to counterintuitive conclusions.

### 2.1.3 Metaethical Tests

The above section specified the model to simulate practical ethical reasoning, or the kind of reasoning that is useful when an agent has to decide what to do. Formalizations of the categorical imperative can also be used to do metaethical reasoning, which can evaluate a particular ethical theory as good or bad. In this case, we can analyze properties of the system in the form of theorems. For example, if we can show that, in this system, nothing is ever obligated, that would indicate that we have a bad ethical system. This is not only philosophical work, but is also a useful way to test different ethical reasoning systems.

An initial property that we might be interested in is permissibility itself. More generally, an ethical theory that doesn't allow for permissibility would require that every action is either obligatory or prohibited. In fact, if that is the case, many counterintuitive theorems follow, including that all permissible actions are obligatory.<sup>4</sup>

**lemma** *permissible*:

**shows**  $\exists A. ((\neg (O \{A\})) \wedge (\neg (O \{\neg A\}))) w$

**nitpick** [*user-axioms, falsify=false*] **oops**

— Nitpick found a model for card i = 1 and card s = 1:

Skolem constant:  $A = (\lambda x. \dots)(i_1 := \text{False})$

— We want to show that there exists a model where there is some formula A that is permissible, or, in English, that permissibility is possible. Nitpick finds a model where the above formula holds, so permissibility is indeed possible.

— Note that it's not clear [Kitcher, 2004] if Kant actually thought that permissibility was a coherent concept. Either way, in modern ethics, permissibility is a pretty widely accepted phenomenon.

**lemma** *fixed-formula-is-permissible*:

**fixes** A

**shows**  $((\neg (O \{A\})) \wedge (\neg (O \{\neg A\}))) w$

**nitpick** [*user-axioms, falsify=false*] **oops**

— Nitpick found a model for card i = 1:

---

<sup>4</sup>Proof is in the appendix.

Free variable:  $A = (\lambda x. \_) (i_1 := \text{False})$

— This is a slightly stronger result: for any arbitrary action A, there is a model that makes it permissible. We actually don't want this to hold, because if A is "murder" then the CI requires that it be prohibited in every world.

Another initial test is the possibility of arbitrary obligations. If anything can be shown to be obligatory in this theory, then clearly it doesn't track our intuitions.

**lemma** *arbitrary-obligations*:

**fixes**  $A::t$

**shows**  $O \{A\} w$

**nitpick** [*user-axioms=true*] **oops**

— Nitpick found a counterexample for card i = 1:

Free variable:  $A = (\lambda x. \_) (i_1 := \text{False})$

— This is good! Shows us that any arbitrary term isn't obligatory.

## Conflicting Obligations

A more complex metaethical property is the possibility of conflicting obligations. Many deontological ethics are criticized for prescribing conflicting obligations, but in Kantian ethics, obligations never conflict [Timmermann, 2013a]. In order for morality to be action-guiding, it needs to be free of conflicting obligations. Let's see if we can have contradictory obligations under the naive formalization.

**lemma** *conflicting-obligations*:

**fixes**  $A$

**shows**  $(O \{A\} \wedge O \{\neg A\}) w$

**nitpick** [*user-axioms, falsify=false*] **oops**

— Nitpick found a model for card i = 2:

Free variable:  $A = (\lambda x. \_) (i_1 := \text{False}, i_2 := \text{True})$

— Oh no! Nitpick found a model with conflicting obligations - that's bad!

This is a property of DDL itself, not necessarily of my formalization specifically. A future, more robust formalization should add an axiom that disallows this. Let's see if that causes any obvious problems.

**lemma** *removing-conflicting-obligations*:

**assumes**  $\forall A. \models (\neg (O \{A\} \wedge O \{\neg A\}))$

**shows** *True*

**nitpick** [*satisfy, user-axioms, format=2*] **oops**

— Nitpick found a model for card i = 1:

Empty assignment

— We can disallow conflicting obligations and the system is still consistent - that's good.

The above is a rather weak notion of contradictory obligations. Korsgaard [Korsgaard, 1985] argues that Kantian ethics also has the stronger property that if two maxims imply a contradiction, they must not be willed. Let's see if that fact holds in this formalization.

**lemma** *implied-contradiction*:

**fixes**  $A::t$

```

fixes B::t
assumes  $\models (\neg (A \wedge B))$ 
shows  $\models (\neg (O \{A\} \wedge O \{B\}))$ 
nitpick [user-axioms]
proof –
  have  $\models (\neg (\Diamond (A \wedge B)))$ 
  by (simp add: assms)
  then have  $\models (\neg (O \{A \wedge B\}))$  by (smt O-diamond)
— Notice that this is almost the property we are interested in. In fact, if  $O\{A \wedge B\}$  is
equivalent to  $O\{A\} \wedge O\{B\}$ , then the proof is complete.
  thus ?thesis oops
— Nitpick found a counterexample for card i = 2:
Free variables: A = ( $\lambda x. \_$ )( $i_1 := \text{True}, i_2 := \text{False}$ ) B = ( $\lambda x. \_$ )( $i_1 := \text{False}, i_2 := \text{True}$ )
— Sadly the property we’re actually interested in doesn’t follow.

```

The above proof yields an interesting observation.  $O\{A \wedge B\}$  is not equivalent to  $O\{A\} \wedge O\{B\}$ . The rough English translation of  $O\{A \wedge B\}$  is “you are obligated to do both A and B”. The rough English translation of  $O\{A\} \wedge O\{B\}$  is “you are obligated to do A and you are obligated to do B.” We think those English sentences mean the same thing, so they should mean the same thing in our logic as well. Let’s test that.

```

lemma distribute-obligations:
  assumes  $\models (O \{A\} \wedge O \{B\})$ 
  shows  $\models O \{A \wedge B\}$ 
  nitpick [user-axioms] oops
— Nitpick found a counterexample for card i = 2:
Free variables: A = ( $\lambda x. \_$ )( $i_1 := \text{True}, i_2 := \text{False}$ ) B = ( $\lambda x. \_$ )( $i_1 := \text{False}, i_2 := \text{True}$ )

```

Note that this is a property of DDL itself, not just my formalization. A future formalization might add this property as an axiom.<sup>5</sup>

## Miscellaneous Properties

I named this formalization the naive formulation for a reason. Though it seems to be an immediate translation of the FUL into DDL, it doesn’t fully respect the properties of modal logic itself. In particular, the formalization as given is equivalent to the below theorem.

```

lemma FUL-alternate:
  shows  $\models ((\Diamond (O \{\neg A\})) \rightarrow (O \{\neg A\}))$ 
  by simp

```

— This means that if something is possibly prohibited, it is in fact prohibited.  
— This is a direct consequence<sup>6</sup> of the naive formalization, but it’s not clear to me that this is actually how we think about ethics. For example, we can imagine an alternate universe where smiling at someone is considered an incredibly rude and disrespectful gesture. In this universe, we are probably prohibited from smiling at people, but this doesn’t mean that in this current universe, smiling is morally wrong.

<sup>5</sup>For discussion of why this property doesn’t hold in DDL, see the Appendix.

<sup>6</sup>For a manual proof, see the Appendix.



The “ought implies can” principle is attributed to Kant<sup>7</sup> and is rather intuitive: you can’t be obligated to do the impossible. It is worth noting that deontic logics evolved [Cresswell and Hughes, 1996] specifically from this principle, so this should hold in both my modified logic and in DDL.

**lemma** *ought-implies-can*:

**shows**  $\forall A. \models (O \{A\} \rightarrow (\Diamond A))$

**using** *O-diamond by blast*

—  $\models \lambda w. ob \ ?B \ ?A \longrightarrow \neg \models \neg \ ?B \wedge \ ?A$  is an axiom of DDL itself, so this theorem holds in DDL.

## 2.2 Kroy’s Formalization of the Categorical Imperative

This section contains a formalization of the categorical imperative introduced by Moshe Kroy in 1976 [Kroy, 1976]. Kroy used Hintikka’s deontic logic to formalize the Formula of Universal Law and the Formula of Humanity. I will first import the additional logical tools that Hintikka’s logic contains that Kroy relies on, then examine the differences between his logic and DDL, and finally implement and test both of Kroy’s formalizations.

### 2.2.1 Logical Background

Kroy’s logic relies heavily on some notion of identity or agency. The logic must be capable of expressing statements like “x does action”, which I can write as “x is the subject of the sentence ‘does action.’” This requires defining a subject.

**typed** *decl*  $s$  —  $s$  is the type for a “subject,” i.e. the subject of a sentence

Kroy also defines a substitution operator<sup>8</sup>.  $P(d/e)$  is read in his logic as “P with e substituted for d.” DDL has no such notion of substitution, so I will define a more generalized notion of an “open sentence.” An open sentence takes as input a subject and returns a complete or “closed” DDL formula by, in effect, binding the free variable in the sentence to the input. For example, “does action” is an open sentence that can be instantiated with a subject.

**type-synonym**  $os = (s \Rightarrow t)$

— “P sub (d/e)” can be written as “S(e)”, where  $S(d) = P$

— The terms that we substitute into are actually instantiations of an open sentence, and substitution just requires re-instantiating the open sentence with a different subject.

### New Operators

Because Isabelle is strongly typed, we need to define new operators to handle open sentences. These operators are similar to DDL’s original operators. We could

<sup>7</sup>The exact philosophical credence of this view is disputed, but the rough idea holds nonetheless. See [Kohl, 2015] for more.

<sup>8</sup>See page 196 in Kroy’s original paper [Kroy, 1976].

probably do without these abbreviations, but they will simplify the notation and make it look more similar to Kroy’s original paper.

**abbreviation** *os-neg*:: $os \Rightarrow os (\neg -)$

**where**  $(\neg A) \equiv \lambda x. \neg(A(x))$

**abbreviation** *os-and*:: $os \Rightarrow os \Rightarrow os (-\wedge -)$

**where**  $(A \wedge B) \equiv \lambda x. ((A(x)) \wedge (B(x)))$

**abbreviation** *os-or*:: $os \Rightarrow os \Rightarrow os (-\vee -)$

**where**  $(A \vee B) \equiv \lambda x. ((A(x)) \vee (B(x)))$

**abbreviation** *os-ob*:: $os \Rightarrow os (O\{-\})$

**where**  $O\{A\} \equiv \lambda x. (O \{A(x)\})$

Once again, the notion of permissibility will be useful here. Recall that an action can either be obligated, permissible, or prohibited. A permissible action is acceptable (there is no specific prohibition against it), but not required (there is no specific obligation requiring it).

**abbreviation** *ddl-permissible*:: $t \Rightarrow t (P \{-\})$

**where**  $P \{A\} \equiv \neg (O \{\neg A\})$

**abbreviation** *os-permissible*:: $os \Rightarrow os (P \{-\})$

**where**  $P \{A\} \equiv \lambda x. P \{A(x)\}$  **Differences Between Kroy’s Logic (Kr) and DDL**

There is potential for complication because Kroy’s original paper uses a different logic than DDL. His custom logic is a slight modification of Hintikka’s deontic logic [Hintikka, 1962]. In this section, I will determine if some of the semantic properties that Kroy’s logic (which I will now call Kr) requires hold in DDL. These differences may become important later and can explain differences in my results and Kroy’s.

Deontic alternatives versus the neighborhood semantics

The most faithful interpretation of Kr is that if  $A$  is permissible in a context, then it must be true at some world in that context. Kr operates under the “deontic alternatives” or Kripke semantics, summarized by Solt [Solt, 1984] as follows: “A proposition of the sort  $OA$  is true at the actual world  $w$  if and only if  $A$  is true at every deontic alternative world to  $w$ .” Under this view, permissible propositions are obligated at some deontic alternatives, or other worlds in the system, but not at all of them. Let’s see if this holds in DDL.

**lemma** *permissible-semantics*:

**fixes**  $A w$

**shows**  $(P \{A\}) w \longrightarrow (\exists x. A(x))$

**nitpick**<sub>[user-axioms]</sub> **oops**

— Nitpick found a counterexample for card i = 1:

Free variable:  $A = (\lambda x. \_)(i_1 := \text{False})$

Remember that DDL uses neighborhood semantics, not the deontic alternatives view, which is why this proposition fails in DDL. In DDL, the *ob* function abstracts away the notion of deontic alternatives. Even if one believes that permissible statements should be true at some deontic alternative, it’s not clear that permissible

statements must be realized at some world. In some ways, this also coheres with our understanding of obligation. There are permissible actions like “Lavanya buys a red folder” that might not happen in any universe.

An even stricter version of the semantics that Kr requires is that if something is permissible at a world, then it is obligatory at some world. This is a straightforward application of the Kripke semantics. Let’s test this proposition.

**lemma** *permissible-semantics-strong*:

**fixes**  $A\ w$

**shows**  $P\ \{A\}\ w \longrightarrow (\exists x. O\ \{A\}\ x)$

**nitpick**<sub>[user-axioms]</sub> **oops**

— Nitpick found a counterexample for card i = 1:

Free variable:  $A = (\lambda x. \_) (i_1 := \text{False})$

This also doesn’t hold in DDL because DDL uses neighborhood semantics instead of the deontic alternatives or Kripke semantics. This also seems to cohere with our moral intuitions. The statement “Lavanya buys a red folder” is permissible in the current world, but it’s hard to see why it would be obligatory in any world.

One implication of the Kripke semantics is that Kr disallows “vacuously permissible statements.” In other words, if something is permissible it has to be obligated at some deontically perfect alternative. If we translate this to the language of DDL, we expect that if  $A$  is permissible, it is obligated in some context.

**lemma** *permissible-semantic-vacuous*:

**fixes**  $A\ w$

**shows**  $P\ \{A\}\ w \longrightarrow (\exists x. ob(x)(A))$

**nitpick**<sub>[user-axioms]</sub> **oops**

— Nitpick found a counterexample for card i = 1:

Free variable:  $A = (\lambda x. \_) (i_1 := \text{False})$

In order to make this true, we’d have to require that everything is either obligatory or prohibited somewhere. Sadly, that breaks everything and destroys the notion of permissibility everywhere<sup>9</sup>. If something breaks later in this section, it may be because of vacuous permissibility.

Obligatory statements should be permissible

Kr includes the intuitively appealing theorem that if a statement is obligated at a world, then it is permissible at that world<sup>10</sup>. Let’s see if that also holds in DDL.

**lemma** *ob-implies-perm*:

**fixes**  $A\ w$

**shows**  $O\ \{A\}\ w \longrightarrow P\ \{A\}\ w$

**nitpick**<sub>[user-axioms]</sub> **oops**

— Nitpick found a counterexample for card i = 2:

Free variable:  $A = (\lambda x. \_) (i_1 := \text{False}, i_2 := \text{True})$

<sup>9</sup>See Appendix for an examination of a buggy version of DDL that led to this insight.

<sup>10</sup>This follows straightforwardly from the Kripke semantics. If proposition  $A$  is obligated at world  $w$ , this means that at all of  $w$ ’s neighbors,  $OA$  holds. Therefore,  $\exists w'$  such that  $w$  sees  $w'$  and  $OA$  holds at  $w'$  so  $A$  is permissible at  $w$ .

Intuitively, it seems untenable for any ethical theory to not include this principle. My formalization should add this as an axiom.

### 2.2.2 The Categorical Imperative

I will now implement Kroy’s formalization of the Formula of Universal Law. Recall that the FUL says “act only in accordance with that maxim which you can at the same time will a universal law” [Kant, 1785]. Kroy interprets this to mean that if an action is permissible for a specific agent, then it must be permissible for everyone. This formalizes the moral intuition prohibiting free-riding. According to the categorical imperative, no one is a moral exception. Formalizing this interpretation requires using open sentences to handle the notion of substitution.

**abbreviation**  $FUL::bool$  **where**  $FUL \equiv \forall w A. ((\exists p::s. ((P \{A\ p\}) w)) \longrightarrow (\forall p. (P \{A\ p\}) w)))$

— In English, this statement roughly means that, if action  $A$  is permissible for some person  $p$ , then, for any person  $p$ , action  $A$  must be permissible. The notion of “permissible *for*” is captured by the substitution of  $x$  for  $p$ .

Let’s check if this is already an axiom of DDL. If so, then the formalization is trivial.

**lemma**  $FUL$ :

**shows**  $FUL$

**nitpick**<sub>[user-axioms]</sub> **oops**

— Nitpick found a counterexample for card s = 2 and card i = 2:

Skolem constants:  $A = (\lambda x. \dots)(s_1 := (\lambda x. \dots)(i_1 := \text{True}, i_2 := \text{True}), s_2 := (\lambda x. \dots)(i_1 := \text{False}, i_2 := \text{False}))$   $p = s_1$   $x = s_2$

This formalization doesn’t hold in DDL, so adding it as an axiom will change the logic.

**axiomatization** **where**  $FUL$ :  $FUL$

Consistency check: is the logic still consistent with the FUL added as an axiom?

**lemma**  $\text{True}$  **nitpick**<sub>[user-axioms, satisfy, card=1]</sub> **oops**

— Nitpicking formula... Nitpick found a model for card i = 1:

Empty assignment

This completes my implementation of Kroy’s formalization of the first formulation of the categorical imperative. I defined new logical constructs to handle Kroy’s logic, studied the differences between DDL and Kr, implemented Kroy’s formalization of the Formula of Universal Law, and showed that it is both non-trivial and consistent. Now it’s time to start testing!

### 2.2.3 Application Tests

In the following sections, I will use the application and metaethical tests presenting in Sections 2.1.2 and 2.1.3 to tease out the strengths and weaknesses of Kroy’s formalization. While the formalization is considerably stronger than the naive formalization, it still fails many of these tests. Some of these failures are due to the differences between Kroy’s logic and my logic mentioned in Section 2.2.1, but some reveal philosophical problems with Kroy’s interpretation of what the formula of universal law means. I will analyze these problems in the context of philosophical scholarship explicating the content of the formula of universal law. The findings in these sections will inform milestones for my custom formalization of the categorical imperative. They also serve as an example of how formalized and automated ethics can reveal philosophical strengths and weaknesses of an ethical theory.

#### Murder

In Section 2.1.2, I began by testing the naive interpretation’s ability to show that murder is wrong. I started by showing the morally dubious proposition that if murder is possibly wrong, then it is actually wrong.

**consts**  $M::t$

— Let the constant  $M$  denote murder. I have defined no features of this constant, except that it is of the type term, which can be true or false at a set of worlds. Indeed, this constant as-is has no semantic meaning and could be replaced with any symbol, like ‘Q’ or ‘Going to Target.’ This constant will begin to take on features of the act of murder when I specify its properties. In the tests below, I specify its properties as the antecedents of lemmas. For example, the test below specifies that it is possible that murder is prohibited at the current world. This pattern will hold for most constants defined in Isabelle—they have no meaning until I program a meaning.

**lemma** *wrong-if-possibly-wrong*:

**shows**  $((\Diamond (O \{ \neg M \})) \text{ } cw) \longrightarrow (\forall w. (O \{ \neg M \}) \text{ } w)$

**by** *simp*

— This sentence reads: “If it is possible that murder is prohibited at world  $cw$ , then murder is prohibited at all worlds.

This is the same result we got in Section 2.1.2—if murder is possibly wrong at some world, it is wrong at every world. The result is incredible strong—the mere possibility of wrongness at some world is sufficient to imply prohibition at every world.

Kroy’s formalization shouldn’t actually imply this property. Recall that this property held in the naive interpretation because it universalized a proposition across worlds (using the necessity operator). Kroy, on the other hand, interprets the FUL as universalizing across people, not worlds. In other words, Kroy’s formulation implies that if murder is wrong for someone, then it is wrong for everyone.

The fact that this strange lemma holds is actually a property of DDL itself, not a property of Kroy’s formalization. Indeed, repeating this experiment in DDL, with

no additional axioms that represent the categorical imperative shows that, in DDL, if something is possibly wrong, it is wrong at every world. This implies that this is not a useful example to test any formulation. If a lemma is true in the base logic, without any custom axioms added, then it will hold for any set of custom axioms. Testing whether or not it holds as we add axioms tells us nothing, since it held in the base logic itself. Interesting cases are ones that fail (or are indeterminate) in the base logic, but become true as we add axioms.

To adapt the murder wrong axiom to capture the spirit of Kroy’s formulation, I will modify it to state that if murder is wrong for one person, it is wrong for everyone.

**consts** *M-kroy::os*

— This time, murder is an open sentence, so that I can substitute in different agents.

**lemma** *wrong-if-wrong-for-someone*:

**shows**  $(\exists p. \models O \{ \neg(M\text{-}kroy\ p) \}) \longrightarrow (\forall p. \models O \{ \neg(M\text{-}kroy\ p) \})$

**proof**

**assume**  $(\exists p. \models O \{ \neg(M\text{-}kroy\ p) \})$

**show**  $(\forall p. \models O \{ \neg(M\text{-}kroy\ p) \})$

**using** *FUL*  $\langle \exists p. \models O \{ \neg M\text{-}kroy \} p \rangle$  **by** *blast*

**qed**

This lemma gets to the heart of Kroy’s formulation of the categorical imperative. If murder is prohibited for a specific person  $p$ , then it must be prohibited for all people<sup>11</sup>.

## Lying

For the naive implementation, I also tested the stronger proposition that if not everyone can simultaneously lie, then lying is prohibited. This is the equivalent of claiming that if lying fails the universalizability test, it is prohibited.

I want to represent the sentence “At all worlds, it is not possible that everyone lies simultaneously.” This requires the following two abbreviations.

**consts** *lie::os*

**abbreviation** *everyone-lies::t* **where** *everyone-lies*  $\equiv \lambda w. (\forall p. (lie(p)\ w))$

— This represents the term “all people lie”. Naively, we might think to represent this as  $\forall p. lie(p)$ . In HOL, the  $\forall$  operator has type  $(\text{'a} \rightarrow bool) \rightarrow bool$ , where  $\text{'a}$  is a polymorphic type of the term being bound by  $\forall$ . In the given example,  $\forall$  has the type  $(s \rightarrow bool) \rightarrow bool$ , so it can only be applied to a formula of type  $s \rightarrow bool$ . In the abbreviation above, we’re applying the quantifier to a sentence that takes in a given subject  $p$  and returns  $lie(p)w$  for any arbitrary  $w$ , so the types cohere.

— The term above is true for a set of worlds  $i$  (recall that a term is true at a set of worlds) such that, at all the worlds  $w$  in  $i$ , all people at  $w$  lie.

<sup>11</sup>This test case also revealed a bug in my original implementation of Kroy’s formulation of the FUL, demonstrating the power of such automated tests and precise formulations to find bugs in ethical theories.

**abbreviation** *lying-not-possibly-universal::bool* **where** *lying-not-possibly-universal*  $\equiv \models (\neg (\Diamond \text{everyone-lies}))$

— Armed with *everyone-lies*  $\equiv \lambda w. \forall p. \text{lie } p \ w$ , it's easy to represent the desired sentence. The abbreviation above reads, “At all worlds, it is not possible that everyone lies.”

Now that I have defined a sentence stating that lying fails the universalizability test, I can test if this sentence implies that lying is impermissible.

**lemma** *lying-prohibited*:

**shows** *lying-not-possibly-universal*  $\longrightarrow (\models (\neg P \{\text{lie}(p)\}))$

**nitpick**<sub>[user-axioms]</sub> **oops**

— Nitpick found a counterexample for card i = 1 and card s = 2:

Free variables:

*lying\_not\_possibly\_universal* = True

*p* = *s*<sub>1</sub>

Kroy's formulation fails this test, and is thus not able to show that if lying is not possible to universalize, it is prohibited for an arbitrary person. To understand why this is happening, I will outline the syllogism that I *expect* to prove that lying is prohibited.

1. At all worlds, it is not possible for everyone to lie. (This is the assumed lemma *lying\_not\_possibly\_universal*)
2. At all worlds, there is necessarily someone who doesn't lie. (Modal dual of (1))
3. If *A* is permissible for subject *p* at world *w*, *A* is possible for subject *p* at world *w*. (Modified Ought Implies Can)
4. If *A* is permissible at world *w* for any person *p*, it must be possible for everyone to *A* at *w*. (FUL and (3))
5. Lying is impermissible. (Follows from (4) and (1))

Armed with this syllogism, I can figure out why this test failed.

**lemma** *step2*:

**shows** *lying-not-possibly-universal*  $\longrightarrow \models (\Box (\lambda w. \exists p. (\neg (\text{lie}(p)) \ w)))$

**by simp**

— Step 2 holds.

**lemma** *step3*:

**fixes** *A p w*

**shows**  $P \{A(p)\} \ w \longrightarrow (\Diamond (A(p)) \ w)$

**nitpick** <sub>[user-axioms, falsify]</sub> **oops**

— Nitpick found a counterexample for card 'a = 1, card i = 1, and card s = 1:

Free variables: *A* =  $(\lambda x. \_)(a_1 := (\lambda x. \_)(i_1 := \text{False}))$  *p* = *a*<sub>1</sub>

As we see above, the syllogism fails at Step 3, explaining why the lemma doesn't hold as expected. Kroy explicitly states<sup>12</sup> that this lemma holds in his logic.

The success of this lemma in Kroy's logic and the emptiness of his formalization of the FUL are two errors that contribute to the failure of this test. First, the statement

---

<sup>12</sup>See footnote 19 on p. 199

expressed in Step 3 should not actually hold. Impossible actions can be permissible (do I need a citation?). For example, imagine I make a trip to Target to purchase a folder, and they offer blue and black folders. No one would claim that it's impermissible for me to purchase a red folder, or, equivalently, that I am obligated to not purchase a red folder.

The second issue is that Kroy's interpretation of the formula of universal law is circular. His formalization interprets the FUL as prohibiting  $A$  if there is someone for whom  $A$ 'ing is not permissible. This requires some preexisting notion of the permissibility of  $A$ , and is thus circular. The categorical imperative is supposed to be the complete, self-contained rule of morality [Kant, 1785], but Kroy's version of the FUL prescribes obligations in a self-referencing manner. The FUL is supposed to define what is permissible and what isn't, but Kroy defines permissibility in terms of itself.

Neither of these errors are obvious from Kroy's presentation of his formalization of the categorical imperative. This example demonstrates the power of formalized ethics. Making Kroy's interpretation of the categorical imperative precise demonstrated a philosophical problem with that interpretation.

#### 2.2.4 Metaethical Tests

In addition to testing specific applications of the theory, I am also interested in metaethical properties, as in the naive interpretation. First, I will test if permissibility is possible under this formalization.

**lemma** *permissible*:

```
fixes A w
shows ((¬ (O {A})) ∧ (¬ (O {¬ A}))) w
nitpick [user-axioms, falsify=false] oops
— Nitpick found a model for card i = 1:
Free variable: A = (λx.⋅)(i1 := False)
```

The above result shows that, for some action  $A$  and world  $w$ , Nitpick can find a model where  $A$  is permissible at  $w$ . This means that the logic allows for permissible actions. If I further specify properties of  $A$  (such as ' $A$  is murder'), I would want this result to fail.

Next, I will test if the formalization allows arbitrary obligations.

**lemma** *arbitrary-obligations*:

```
fixes A::t
shows O {A} w
nitpick [user-axioms=true, falsify] oops
— Nitpick found a counterexample for card i = 1 and card s = 1:
Free variable: A = (λx.⋅)(i1 := False)
```

This is exactly the expected result. Any arbitrary action  $A$  isn't obligated.  $A$



slightly stronger property is “modal collapse,” or whether or not ‘ $A$  happens’ implies ‘ $A$  is obligated’.

**lemma** *modal-collapse*:

**fixes**  $A\ w$

**shows**  $A\ w \longrightarrow O\ \{A\}\ w$

**nitpick** [*user-axioms=true, falsify*] **oops**

— Nitpick found a counterexample for card  $i = 1$  and card  $s = 1$ :

Free variables:  $A = (\lambda x. \dots)(i_1 := \text{True})\ w = i_1$

This test also passes. Next, I will test if not ought implies can holds. Recall that I showed in Section 2.1.3 that ought implies can is a theorem of DDL itself, so it should still hold.

**lemma** *ought-implies-can*:

**fixes**  $A\ w$

**shows**  $O\ \{A\}\ w \longrightarrow \Diamond\ A\ w$

**using** *O-diamond* **by** *blast*

This theorem holds. Now that I have a substitution operation, I also expect that if an action is obligated for a person, then it is possible for that person. That should follow by the axiom of substitution [Cresswell and Hughes, 1996] which lets me replace the ‘ $A$ ’ in the formula above with ‘ $A(p)$ ’

**lemma** *ought-implies-can-person*:

**fixes**  $A\ w$

**shows**  $O\ \{A(p)\}\ w \longrightarrow \Diamond\ (A\ (p))\ w$

**using** *O-diamond* **by** *blast*

This test also passes. Next, I will explore whether or not Kroy’s formalization still allows conflicting obligations.

**lemma** *conflicting-obligations*:

**fixes**  $A\ w$

**shows**  $(O\ \{A\} \wedge O\ \{\neg A\})\ w$

**nitpick** [*user-axioms, falsify=false*] **oops**

— Nitpick found a model for card  $i = 2$  and card  $s = 1$ :

Free variable:  $A = (\lambda x. \dots)(i_1 := \text{False}, i_2 := \text{True})$

Just as with the naive formalization, Kroy’s formalization allows for contradictory obligations. Testing this lemma in DDL without the FUL shows that this is a property of DDL itself. This is a good goal to have in mind when developing my custom formalization.

Next, I will test the stronger property that if two maxims imply a contradiction, they may not be simultaneously willed.

**lemma** *implied-contradiction*:

**fixes**  $A\ B\ w$

**assumes**  $((A \wedge B) \rightarrow \perp)\ w$

**shows**  $\neg (O\ \{A\} \wedge O\ \{B\})\ w$

**nitpick** [*user-axioms, falsify*] **oops**

— Nitpick found a counterexample for card i = 2 and card s = 1:

Free variables:  $A = (\lambda x. \_) (i_1 := \text{True}, i_2 := \text{False})$   $B = (\lambda x. \_) (i_1 := \text{True}, i_2 := \text{False})$   $w = i_2$

Just as with the naive formalization, Kroy’s formalization allows implied contradictions because DDL itself allows implied contradictions and Kroy’s formalization doesn’t do anything to remedy this.

Next, I will test that an action is either obligatory, permissible, or prohibited.

**lemma** *ob-perm-or-prohibited*:

**fixes**  $A w$

**shows**  $(O \{A\} \vee (P \{A\} \vee O \{\neg A\})) w$

**by** *simp*

— This test passes.

I also expect obligation to be a strictly stronger property than permissibility. Particularly, if A is obligated, then A should also be permissible.

**lemma** *obligated-then-permissible*:

**shows**  $(O \{A\} \rightarrow P \{A\}) w$

**nitpick**<sub>[user-axioms]</sub> **oops**

— This test fails in Kroy’s interpretation! Nitpick found a counterexample for card i = 2 and card s = 1:

Free variable:  $A = (\lambda x. \_) (i_1 := \text{False}, i_2 := \text{True})$

These tests show that, while Kroy’s formalization is more powerful and more coherent than the naive formalization, it still fails to capture most of the desired properties of the categorical imperative. Some of these problems may be remedied by the fact that Kroy’s logic doesn’t allow contradictory obligations, and that possibility will be interesting to explore in my own formalization.

## 2.2.5 Miscellaneous Tests

In this section, I explore tests of properties that Kroy presents in his original paper. These tests not only test the features of the system that Kroy intended to highlight, but they may also inspire additional tests and criteria for my own formalization in Chapter 3. These tests further underscore the circularity of Kroy’s formalization and the differences between my logic and his.

First, Kroy presents a stronger version of the formula of universal law and argues that his formalization is implied by the stronger version. Let’s test that claim.

**abbreviation** *FUL-strong::bool* **where**  $FUL\text{-}strong \equiv \forall w A. ((\exists p::s. ((P \{A p\}) w)) \rightarrow ((P \{ \lambda x. \forall p. A p x \}) w))$

**lemma** *strong-implies-weak*:

**shows**  $FUL\text{-}Strong \rightarrow FUL$

**using** *FUL* **by** *blast*

— This lemma holds, showing that Kroy is correct in stating that this version of the FUL is stronger than his original version.

The difference between the stronger version and  $FUL \equiv \forall w A. (\exists p. P \{A p\} w) \longrightarrow (\forall p. P \{A p\} w)$  is subtle. The consequent of  $FUL$  is “for all people  $p$ , it is permissible that they  $A$ .” The consequent of this stronger statement is “it is permissible that everyone  $A$ .” In particular, this stronger statement requires that it is permissible for everyone to  $A$  simultaneously. Kroy immediately rejects this version of the categorical imperative, arguing that it’s impossible for everyone to be the US president simultaneously, so this version of the  $FUL$  prohibits running for president.

Most Kantians would disagree with this interpretation. Consider the classical example of lying, as presented in [Kemp, 1958] and in [Korsgaard, 1985]. Lying fails the universalizability test because in a world where everyone lied simultaneously, the practice of lying would break down. If we adopt Kroy’s version, lying is only prohibited if, no matter who lies, lying is impermissible. As argued above, this rule circularly relies on some existing prohibition against lying for a particular person, and thus fails to show the wrongness of lying. It is tempting to claim that this issue explains why the tests above failed. To test this hypothesis, I will check if the stronger version of the  $FUL$  implies that lying is impermissible.

**lemma** *strongFUL-implies-lying-is-wrong*:

**fixes**  $p$

**shows**  $FUL\text{-}strong \longrightarrow \models (\neg P \{lie(p)\})$

**nitpick**[*user-axioms, falsify*] **oops**

— Nitpick found a counterexample for card  $i = 1$  and card  $s = 1$ :

Free variable:  $p = s_1$

The test above also fails! This means that not even the stronger version of Kroy’s formalization of the  $FUL$  can show the wrongness of lying. As mentioned earlier, there are two independent errors. The first is the the assumption that impossible actions are impermissible and the second is the circularity of the formalization. The stronger  $FUL$  addresses the second error, but the first remains.

Kroy also argues that the  $FUL$  gives us recipes for deriving obligations, in addition to deriving permissible actions. Specifically, he presents the following two principles, which are equivalent in his logic. These sentences parallel  $FUL$  and strong  $FUL$ .

**abbreviation** *obligation-universal-weak::bool* **where** *obligation-universal-weak*  $\equiv \forall w A. ((\exists p::s. ((O \{A p\}) w)) \longrightarrow (\forall p. (O \{A p\}) w))$

**abbreviation** *obligation-universal-strong::bool* **where** *obligation-universal-strong*  $\equiv \forall w A. ((\exists p::s. ((O \{A p\}) w)) \longrightarrow (((O \{ \lambda x. \forall p. A p x \}) w)))$

— Just as with  $FUL$  and  $FUL$  strong, the weaker version of the above statement has the consequent, “For all people,  $A$  is obligated.” The stronger consequent is “ $A$  is obligated for all people simultaneously.”

**lemma** *weak-equiv-strong*:

**shows** *obligation-universal-weak*  $\equiv$  *obligation-universal-strong*

**oops**

— Isabelle is neither able to find a proof nor a countermodel for the statement above, so I

can't say if it holds or not without completing a full, manual proof. This aside is not very relevant to my project, so I will defer such a proof.

These two statements are not necessarily equivalent in my logic, but are in Kroy's<sup>13</sup> This difference in logics may further explain why tests are not behaving as they should. Nonetheless, Kroy argues that the FUL implies both statements above.

**lemma** *FUL-implies-ob-weak*:

**shows**  $FUL \longrightarrow \text{obligation-universal-weak}$  **oops**

— Isabelle is neither able to find a proof nor a countermodel for this statement.

**lemma** *FUL-implies-ob-strong*:

**shows**  $FUL \longrightarrow \text{obligation-universal-strong}$  **oops**

— Isabelle is neither able to find a proof nor a countermodel for this statement.

Isabelle timed out when looking for proofs or countermodels to the statements above. This may be an indication of a problem that Benzmueller warned me about—mixing quantifiers into a shallow embedding of DDL may be too expensive for Isabelle to handle. Not sure what to do about this.

### 3 Novel Formalization of the Categorical Imperative

In this section, I present a custom formalization of the categorical imperative, as inspired by the goals from the previous chapter.

#### 3.1 Logical Background

The previous attempts to model the categorical imperative in Chapter 2 partially failed due to an inability to fully represent the complexity of a maxim. Specifically, they treated actions as a single, monolithic unit of evaluation, whereas most Kantians consider the unit of evaluation for the FUL to be the more complex notion of a maxim. In this section, I will present some logical background necessarily to fully capture the spirit of a maxim. I will begin by borrowing some machinery to handle “subjects” who perform actions from Chapter 2.

**typeddecl**  $s$  —  $s$  is the type for a “subject,” i.e. the subject of a sentence. In this interpretation, a subject is merely defined as “that which can act.” It does not include any other properties, such as rationality or dignity. As I will show, for the purposes of defining the universalizability test, this “thin” representation of a subject suffices.

**type-synonym**  $os = (s \Rightarrow t)$  — Recall that an open sentence maps a subject to a term to model the substitution operator.

**type-synonym**  $maxim = (t * os * t)$

---

<sup>13</sup>This follows from the fact that the Barcan formula holds in Kroy's logic but not in mine, as verified with Nitpick. See Appendix for more.

The central unit of evaluation for the universalizability test is a “maxim,” which Kant defines in a footnote in *Groundwork* as “the subjective principle of willing,” or the principle that the agent acts on [Kant, 1785, 16]. Modern Kantians differ in their interpretations of this definition. The naive view is that a maxim is an act, but Korsgaard adopts the more sophisticated view that a maxim is composed of an act and the agent’s purpose for acting [Korsgaard, 2005]. She also compares a maxim to Aristotle’s logos, which includes these components and information about the circumstances and methods of the act. O’Neill concludes that Kant’s examples imply that a maxim must also include circumstances [O’Neill, 2013], and Kitcher [Kitcher, 2003] uses textual evidence from the *Groundwork* to argue for the inclusion of a maxim’s purpose or motivation. In order to formalize the notion of a maxim, I must adopt a specific definition and defend my choice.

I define a maxim as a circumstance, act, goal tuple (C, A, G), read as “In circumstances C, act A for goal G.” Isabelle’s strict typing rules mean that the choice of the type of each member of this tuple is significant. A circumstance is represented as a set of worlds  $t$  where that circumstance holds. A goal is also a term because it can be true or false at a world if it is realized or not. An act is an open sentence because an act itself is not the kind of thing that can be true or false (as in, an act is not truth-apt), but the combination of a subject performing an act can be true or false at a world depending on whether or not the act is indeed performed by that subject. For example, “running” is not truth-apt, but “Sara runs” is truth-apt.

My definition of a maxim is inspired by O’Neill’s work on maxims. I will defend my representation below and consider an additional component that Kitcher argues for.

#### *O’Neill’s Original Schematic and The Role of Practical Judgement*

O’Neill [O’Neill, 2013, 37] presents what Kitcher [Kitcher, 2003] calls the widely accepted view that a maxim is a circumstance, act, goal tuple. A maxim is an action-guiding rule and thus naturally includes an act and the circumstances under which it should be performed, which are often referred to as “morally relevant circumstances.”

She also includes a purpose, end, or goal in the maxim because Kant includes this in many of his example maxims and because Kant argues that human activity, because it is guided by a rational will, is inherently purposive [Kant, 1785, 4 : 428]. A rational will does not act randomly (else it would not be rational), but instead in the pursuit of ends which it deems valuable. This inclusion is also essential for the version of the universalizability test that I will implement, explained in Section ??.

O’Neill’s inclusion of circumstances is potentially controversial because it leaves open the question of what qualifies as a relevant circumstance for a particular maxim. This gives rise to “the tailoring objection” [Kitcher, 2003, 217]<sup>14</sup>, under which maxims are arbitrarily specified to pass the FUL. For example, the maxim “When my name is Lavanya Singh, I will lie to get some easy money,” is universal-

<sup>14</sup>Kitcher cites [Wood, 1999] as offering an example of a false positive due to this objection.

izable, but is clearly a false positive. One solution to this problem is to argue that the circumstance “When my name is Lavanya Singh” is not morally relevant to the act and goal. This solution requires some discussion of what qualifies as a relevant circumstance.

O’Neill seems to acknowledge the difficulty of determining relevant circumstances when she concedes that a maxim cannot include all of the infinitely many circumstances in which the agent may perform the action [O’Neill, 2013, 4 : 428]. She argues that this is an artifact of the fact that maxims are rules of practical reason, the kind of reason that helps us decide what to do and how to do it [Bok, 1998]. Like any practical rule, maxims require the exercise of practical judgement to determine in which circumstances they should be applied. This judgement, applied in both choosing when to exercise the maxim and in the formulation of the maxim itself, is what determines the “morally relevant circumstances.”

The upshot for computational ethics is that the computer cannot perform all ethical activity alone. Human judgement and the exercise of practical reason are essential to both formulate maxims and determine when the actual conditions of life coincide with the circumstances in which the maxim is relevant. Choosing when to exercise a maxim is less relevant to my project because analyzing a formal representation of the FUL requires making the circumstances in a given scenario precise, but will be important for applications of computational ethics to guiding AI agents. The difficulty in formulating a maxim, on the other hand, demonstrates the important fact that ethics, as presented here, is not a solely computational activity. A human being must create a representation for the dilemma they wish to test, effectively translating a complex, real situation into a flat logical structure. This parallels the challenge that programmers face when translating the complexity of reality to a programming language or computational representation. Not only will some of the situation’s complexity inevitably be lost, the outcome of the universalizability test will depend on how the human formulates the maxim and whether or not this formulation does indeed include morally relevant circumstances. If the human puts garbage into the test, the test will return garbage out.

While this may appear to be a weakness of my system, I believe that it actually allows my system to retain some of the human complexity that many philosophers agree cannot be automated away.<sup>15</sup> Ethics is a fundamentally human activity. Kant argues that the categorical imperative is a statement about the properties of rational wills. In fact, Korsgaard argues that morality derives its authority over us, or normativity, only because it is a property of a rational will, and we, as human beings, are rational wills. If ethics is meant to guide human behavior, the role of the computer becomes clear as not a replacement for our will, but instead as a tool to help guide our wills and reason more efficiently and more effectively. Just as calculators don’t render mathematicians obsolete, computational ethics does not render human judgement or philosophy obsolete. Chapter 4 Section ?? will be devoted to a more

---

<sup>15</sup>Powers presents the determination of morally relevant circumstances as an obstacle to the automation of Kantian ethics [Powers, 2006].

complete discussion of this issue.

#### *Exclusion of Motive*

Kitcher begins with O’Neill’s circumstance, act, goal view and expands it to include the motive behind performing the maxim [Kitcher, 2003]. This additional component is read as “In circumstance C, I will do A in order to G because of M,” where M may be “duty” or “self-love.” Kitcher argues that the inclusion of motive is necessary for the fullest, most general form of a maxim in order to capture Kant’s idea that an action derives its moral worth from being done for the sake of duty itself. Under this view, the FUL would obligate maxims of the form “In circumstance C, I will do A in order to G because I can will that I and everyone else simultaneously will do A in order to G in circumstance C.” In other words, if Kant is correct in arguing that moral actions must be done from the motive of duty, the affirmative result of the FUL becomes the motive for a moral action.

While Kitcher’s conception of a maxim captures Kant’s idea of acting for duty’s own sake, I will not implement it because it is not necessary for putting maxims through the FUL. Indeed, Kitcher acknowledges that O’Neill’s formulation suffices for the universalizability test, but is not the general notion of a maxim. In order to pass the maxim through the FUL, it suffices to know the circumstance, act, and goal. The FUL derives the motive that Kitcher bundles into the maxim, so automating the FUL does not require including a motive. The “input” to the FUL is the circumstance, act, goal tuple. My project takes this input and returns the motivation that the dutiful, moral agent would adopt. Additionally, doing justice to the rich notion of motive requires modelling the operation of practical reason itself, which is outside the scope of this project. My work focuses on the universalizability test, but future work that models the process of practical reason may use my implementation of the FUL as a “library.” Combined with a logic of practical reason, an implementation of the FUL can move from evaluating a maxim to evaluating an agent’s behavior, since that’s when “acting from duty” starts to matter.

**abbreviation**  $will :: maxim \Rightarrow s \Rightarrow t (W - -)$

**where**  $will \equiv \lambda(c, a, g) s. (c \rightarrow (a s))$

**print-theorems**

Korsgaard claims that “to will an end, rather than just wishing for it or wanting it, is to set yourself to be its cause” [Korsgaard and O’Neill, 1996, 38]. To will a maxim is to set yourself to be the cause of its goal by taking the means specified in the maxim in the relevant circumstances. This coheres with Kitcher’s and Korsgaard’s understanding of a maxim as a principle or rule to live by.

At worlds where the circumstances do not hold, a maxim is vacuously willed. If you decide to act on the rule “I will do X in these circumstances”, then you are vacuously obeying it when the circumstances don’t hold.

The above discussion implies that willing a maxim is particular to the agent, justifying my choice to require that a particular subject will a maxim. O’Neill argues



for this interpretation when she distinguishes between the evaluation of a principle, which is generic, and a maxim, which she views as “individuated only by referring to a person” [O’Neill, 2013, 13]. I adopt the spirit of this interpretation but modify it slightly by representing the general maxim as a principle that anyone could adopt, and the act of willing the maxim as a person-particular instantiation of the maxim. I additionally represent a subject as willing a maxim because I use the word ‘will’ as a verb, to mean committing oneself to living by the principle of a maxim. This coheres with the FUL, which tests the act willing of a maxim by determining if the maxim could be a universal law that everyone committed to. Formalizing this idea, the type of a willed maxim is a term, allowing me to use DDL’s obligation operator on the notion of willing a maxim. Concretely, my system will prove or disprove statements of the form “Lavanya is obligated to will the maxim M.”

Worlds where the circumstances do not hold are not relevant for determining obligation. Recall that in Benzmueller et. al’s definition of the obligation operator,  $O\{B|A\}$  is true at all worlds iff  $ob(B)(A)$ , or if the obligation function maps A to obligatory in context B (where the context is a set of worlds) [Benzmüller et al., 2021]. This definition implies that worlds outside of B have no bearing on the moral status of A in context B, which coheres with intuitions about contextual obligation. Thus, the dyadic obligation operator disqualifies worlds where the context does not hold, so the vacuous truth of the will statement in these worlds does not matter.

Given that the will abbreviation already excludes worlds where the circumstances fail (by rendering the statement vacuously true at them), one may conclude that the dyadic obligation operator is now useless. Using the dyadic obligation operator allows me to take advantage of the power of DDL to represent the bearing that circumstances have on obligation. DDL has powerful axioms expressing the relationship between circumstances and obligation, such as the fact that obligations are monotonically increasing with respect to broader circumstances. Using the monadic obligation operator would require me to either operate with an empty notion of context or to redefine these axioms. The dyadic obligation operator lets me take advantage of the full power of DDL in expressing contrary-to-duty obligations. This is particularly important for Kantian ethics and the FUL specifically because many critiques of the FUL rely on attention to circumstances (tailoring objection) or lack thereof (ideal theory). This is also an innovation that my custom formalization presents over the prior work. By formally including the notion of a circumstance or context, I am able to represent these objections that Kantian scholars study. Formalizing Kantian ethics in a dyadic deontic logic instead of a monadic deontic logic is a key contribution of this thesis.

**abbreviation** *effective* ::  $maxim \Rightarrow s \Rightarrow t$  (E - -)

**where** *effective*  $\equiv \lambda(c, a, g) s. ((will\ (c, a, g)\ s) \equiv g)$

**print-theorems**

A maxim is effective for a subject when, if the subject wills it then the goal



is achieved, and when the subject does not act on it, the goal is not achieved<sup>16</sup> [Menzies and Beebe, 2020]. The former direction of the implication is intuitive: if the act results in the goal, it was effective in causing the goal. This represents ‘necessary’ causality.

The latter direction represents ‘sufficient’ causality, or the idea that, counterfactually, if the maxim were not willed, then the goal is not achieved [Lewis, 1973a]. Note that nothing else changes about this counterfactual world—the circumstances are identical and we neither added additional theorems nor specified the model any further. This represents Lewis’s idea of “comparative similarity,” where a counterfactual is true if it holds at the most similar world [Lewis, 1973b]. In our case, this is just the world where everything is the same except the maxim is not acted on.

Combining these ideas, this definition of effective states that a maxim is effective if the maxim being acted on by a subject is the necessary and sufficient cause of the goal.<sup>17</sup>

If the circumstances do not hold and the goal is achieved, then the maxim is vacuously effective, since it is vacuously willed (as described above). While this scenario is counterintuitive, it is not very interesting for my purposes because, when the circumstances do not hold, a maxim is not applicable. It doesn’t really make sense to evaluate a maxim when it’s not supposed to be applied. The maxim “When on Jupiter, read a book to one-up your nemesis” is vacuously effective because it can never be disproven.

**abbreviation** *universalized::maxim $\Rightarrow$ s $\Rightarrow$ t* **where**  
*universalized*  $\equiv \lambda M s. (\lambda w. (\forall p. W M p w))$

**abbreviation** *not-universalizable :: maxim $\Rightarrow$ s $\Rightarrow$ bool* **where**  
*not-universalizable*  $\equiv \lambda M s. (\models (universalized M s \rightarrow (\neg (E M s))))$   
 — The maxim willed by subject  $s$  is not universalizable if, for all people  $p$ , if  $p$  wills  $M$ , then  $M$  is no longer effective for  $s$ .

As before, the concepts of prohibition and permissibility will be helpful here. The unit of evaluation for my formalization of the FUL is the act of willing a maxim, which entails performing the maxim’s act in the relevant circumstances. Therefore, I will say that, just as the act of willing a maxim can be obligatory for a subject, it can be prohibited or permissible for a subject.<sup>18</sup>

**abbreviation** *prohibited::maxim $\Rightarrow$ s $\Rightarrow$ t* **where**  
*prohibited*  $\equiv \lambda(c, a, g) s. O\{\neg (will (c, a, g) s) \mid c\}$

**abbreviation** *permissible::maxim $\Rightarrow$ s $\Rightarrow$ t*  
**where** *permissible*  $\equiv \lambda M s. \neg (prohibited M s)$

<sup>16</sup>Thank you to Jeremy D. Zucker for helping me think through this.

<sup>17</sup>Should I wave a hand at critiques of counterfactual causality?

<sup>18</sup>In the rest of this section, for convenience, I will use the phrase “subject  $s$  willing maxim  $M$  is obligatory” interchangeably with “maxim  $M$  is obligatory for subject  $s$ .” I will use “maxim  $M$  is obligatory” to refer to  $M$  being obligatory for any arbitrary subject, which I will show to be equivalent to  $M$  being obligatory for a specific subject.

— I will say that a maxim is permissible for a subject if it is not prohibited for that subject to will that maxim.

One problem with prior formalization of the categorical imperative was that they didn't prohibit contradictory obligations, partially because DDL itself allows contradictory obligations. Kant subscribes to the general, popular view that morality is supposed to guide action, so ought implies can<sup>19</sup>. Kohl reconstructs his argument for the principle as follows: if the will cannot comply with the moral law, then the moral law has no prescriptive authority for the will [Kohl, 2015, 703-4]. This defeats the purpose of Kant's theory—to develop an unconditional, categorical imperative for rational agents. Ought implies can requires that obligations never contradict, because an agent can't perform contradictory actions. Therefore, any ethical theory that respects ought implies can, and Kantian ethics in particular, must not result in conflicting obligations.

Kant only briefly discusses contradictory obligations in *Metaphysics of Morals*, where he argues that conflicting moral obligations are impossible under his theory [Kant, 2017, V224]. Particularly, the categorical imperative generates “strict negative laws of omission”, which cannot conflict by definition [Timmermann, 2013b, 45].<sup>20</sup>

When analyzing the naive formalization and Kroy's formalization, I learned that DDL and the prior formalizations allow contradictory obligations. This is a major weakness of these systems, and my formalization should fix this. To do so, I will add as an axiom the idea that obligations cannot contradict each other or their internal circumstances. Formally, conflicting obligations are defined below.

**abbreviation non-contradictory where**

*non-contradictory*  $A B c w \equiv ((O\{A|c\} \wedge O\{B|c\}) w) \longrightarrow \neg((A \wedge (B \wedge c)) w \longrightarrow \text{False})$

— Terms A and B are non contradictory in circumstances c if, when A and B are obligated in circumstances c, the conjunction of A, B, and c, does not imply False.

**axiomatization where no-contradictions:**  $\forall A::t. \forall B::t. \forall c::t. \forall w::i. \text{non-contradictory } A B c w$

— This axiom formalizes the idea that, for any terms A, B, and circumstances c, A and B must be non-contradictory in circumstances c at all worlds. Intuitively, this axiom requires that obligations do not conflict.

<sup>19</sup>Kohl points out that this principle is referred to as Kant's dictum or Kant's law in the literature. [Kohl, 2015, footnote 1]

<sup>20</sup>The kinds of obligations generated by the FUL are called “perfect duties” which arise from “contradictions in conception,” or maxims that we cannot even conceive of universalizing. These duties are always negative and thus never conflict. Kant also presents “imperfect duties,” generated from “contradictions in will,” or maxims that we can conceive of universalizing but would never want to. These duties tend to be broader, such as “improve oneself” or “help others,” and are secondary to perfect duties. My project only analyzes perfect duties, as these are stronger than imperfect duties.

### 3.2 Formalizing the FUL

Below is my first attempt at formalizing Korgsaard’s definition of the practical contradiction interpretation: a maxim is not universalizable if, in the world where the maxim becomes the standard practice (i.e. everyone acts on the maxim), the agent’s attempt to use the maxim’s act to achieve the maxim’s goal is frustrated. In other words, if the maxim is universally willed (captured by applying a universal quantifier and the will function to the maxim on the LHS), then it is no longer effective for the subject  $s$  (RHS above).

**abbreviation**  $FUL0::bool$  **where**  $FUL0 \equiv \forall c a g s. not-universalizable (c, a, g) s \longrightarrow \models (prohibited (c, a, g) s)$

— This representation of the Formula of Universal Law reads, “For all circumstances, goals, acts, and subjects, if the maxim of the subject performing the act for the goal in the circumstances is not universalizable (as defined above), then, at all worlds, in those circumstances, the subject is prohibited (obligated not to) from willing the maxim.

**lemma**  $FUL0 \longrightarrow False$  **using**  $O\text{-diamond}$   
**using**  $prod.simps(2)$   $split\text{-conv}$  **by**  $fastforce$

FUL0 is not consistent, and sledgehammer is able to prove this by showing that it implies a contradiction using axiom  $O\_diamond$ , which is  $\models_{\lambda w}. ob ?B ?A \longrightarrow \neg \models \neg ?B \wedge ?A$ . This axiom captures the idea that an obligation can’t contradict its context. This is particularly problematic if the goal or action of a maxim are equivalent to its circumstances. In other words, if the maxim has already been acted on or the goal has already been achieved, then prohibiting it is impossible. In any model that has at least one term, it is possible to construct a maxim where the circumstances, goal, and act (once a subject acts on it) are all that same term, resulting in a contradiction.

To get around this, I will exclude what I call “badly formed maxims,” which are those maxims such that the goal has already been achieved or the act has already been acted on. Under my formalization, such maxims are not well-formed. To understand why, I return to Korsgaard’s and O’Neill’s interpretations of a maxim as a practical guide to action. A maxim is a practical principle that guides how we behave in everyday life. A principle of the form “When you are eating breakfast, eat breakfast in order to eat breakfast,” is not practically relevant. No agent would ever need to act on such a principle. It is not contradictory or prohibited, but it is the wrong kind of question to be asking. It is not a well-formed maxim, so the categorical imperative does not apply to it. (more explanation in philosophical writing collection)

**abbreviation**  $well\text{-formed}::maxim \Rightarrow s \Rightarrow i \Rightarrow bool$  **where**

$well\text{-formed} \equiv \lambda(c, a, g) s w. (\neg (c \rightarrow g) w) \wedge (\neg (c \rightarrow a s) w)$

— This abbreviation formalizes the well-formedness of a maxim for a subject. The goal cannot be already achieved in the circumstances and the subject cannot have already performed the act.

**abbreviation *FUL* where**  $FUL \equiv \forall M::\text{maxim}. \forall s::s. (\forall w. \text{well-formed } M \ s \ w) \longrightarrow$   
 $(\text{not-universalizable } M \ s \longrightarrow \models \text{prohibited } M \ s)$

— Let’s try the exact same formalization of the *FUL* as above, except that it only applies to maxims that are well-formed at every world.

**lemma *FUL***

**nitpick**<sub>[user-axioms, falsify=true]</sub> **oops**

— The *FUL* does not hold in DDL, because nitpick is able to find a model for my system in which it is false. If the *FUL* were already a theorem of the system, adding it wouldn’t make the system any more powerful, so this is the desired result.

Nitpick found a counterexample for card s = 1 and card i = 1:

Skolem constants:  $M = ((\lambda x. \_) (i_1 := \text{True}), (\lambda x. \_) (s_1 := (\lambda x. \_) (i_1 := \text{False}))), (\lambda x. \_) (i_1 := \text{False})) \lambda w. p = (\lambda x. \_) (i_1 := s_1) \ s = s_1$

**axiomatization where *FUL*:*FUL***

**lemma *True***

**nitpick**<sub>[user-axioms, falsify=false]</sub> **by simp**

— Nitpick is able to find a model in which all axioms are satisfied, so this version of the *FUL* is consistent.

Nitpick found a model for card i = 1 and card s = 1:

Empty assignment

During the process of making *FUL0* consistent, I used Isabelle to gain philosophical insights about vacuous maxims. This process is an example of the power of computational tools to aid philosophical progress. I used Nitpick and Sledgehammer to quickly test if a small tweak to *FUL0* fixed the inconsistency or if I was still able to derive a contradiction. I then realized that if I defined the circumstances, act, and goal as constants, then *FUL0* was indeed consistent. After some experimentation, Prof. Amin correctly pointed out that as constants, these three entities were distinct. However, when merely quantifying over (c, a, g), all members of a tuple could be equivalent. Within a minute, I could formalize this notion, add it to *FUL0*, and test if it solved the problem. The fact that it did spurred my philosophical insight about vacuous maxims.

The logic confirmed that certain kinds of circumstance, act, goal tuples are too badly formed for the categorical imperative to logically apply to them. The realization of this subtle problem would have been incredibly difficult without computational tools. The syntax and typing of Isabelle/HOL forced me to bind the free-variable *M* in the *FUL* in different ways and allowed me to quickly test many bindings. The discovery of this logical inconsistency then enabled a philosophical insight about which kinds of maxims make sense as practical principles. This is one way to do computational ethics: model a system in a logic, use computational tools to refine and debug the logic, and then use insights about the logic to derive insights about the ethical phenomenon it is modelling. This procedure parallels the use of proofs in theoretical math to understand the mathematical objects they model.

One potential problem with my formalization is that it does not use the modal nature of the system. All of the properties that the FUL investigates hold at all worlds, in effect removing the modal nature of the system. This approach simplifies logical and therefore computational complexity, improving performance. On the other hand, it doesn't use the full expressivity of DDL. If I run into problems later on, one option is to tweak the FUL to use this expressivity.

### 3.3 Application Tests

As with the naive formalization and Kroy's formalization, I will apply my testing framework to my custom formalization of the FUL. I will begin with some basic application tests. In these tests, I specify particular maxims as constants with no properties and gradually add properties to understand how the system handles different kinds of maxims.

### 3.4 Metaethical Tests

Recall that metaethical tests test formal properties of the system that apply to any maxim, not just those specified in the application tests. In this section I adapt the metaethical tests developed in previous sections to my formalization of the categorical imperative. I preserved the philosophical goal of each test but modified them to test the stronger, richer notion of a maxim.

The first set of tests consider how obligation generalizes, first across worlds and then across people. As expected, the tests below show that both wrongness (prohibition) and rightness (obligation) generalize across both worlds and people. In other words, if something is obligated at some world, it is obligated at every world and if something is obligated for some person, then it is obligated for every person.

Generalization across people coheres with Kantian ethics—maxims are not person-specific. Indeed, Velleman argues that, because reason is accessible to everyone identically, obligations apply to all people [Velleman, 2005, 25]. When Kant describes the categorical imperative as the objective principle of the will, he is referring to the fact that, as opposed to a subjective principle, the categorical imperative applies to all rational agents equally [Kant, 1785, 16].

Generalization across worlds is a consequence of the fact that my interpretation does not make use of the modal nature of DDL. In particular, I do not use any property of the world when prescribing obligations at that world.

**lemma** *wrong-if-wrong-for-someone*:

**shows**  $\forall w. \forall c::t. \forall g::t. \exists s::s. O\{\neg (W (c, M, g) s) \mid c\} w \longrightarrow (\forall p. O\{\neg (W (c, M, g) p) \mid c\} w)$

**by** *blast*

**lemma** *right-if-right-for-someone*:

**shows**  $\forall w. \forall c::t. \forall g::t. \exists s::s. O\{W (c, M, g) s \mid c\} w \longrightarrow (\forall p. O\{W (c, M, g) p \mid c\} w)$

**by blast**

**lemma** *wrong-if-wrong-somewhere:*

**shows**  $\forall c\ g. \exists w1. O\{\neg (W(c, M, g)\ s) \mid c\} w1 \longrightarrow (\forall w2. O\{\neg (W(c, M, g)\ s) \mid c\} w2)$

**by blast**

**lemma** *right-if-right-somewhere:*

**shows**  $\forall c\ g. \exists w1. O\{W(c, M, g)\ s \mid c\} w1 \longrightarrow (\forall w2. O\{W(c, M, g)\ s \mid c\} w2)$

**by blast**

As expected, obligation generalizes across people and worlds. In the next set of tests, I will analyze basic properties of permissibility, obligation, and prohibition.

First, I verify that the logic allows for permissible maxims, as this is a problem that prior iterations ran into. Below, I use Nitpick to find a model in which there is a circumstance, act, goal tuple that is permissible but not obligated at some world.

**lemma** *permissible:*

**shows**  $((\neg (O\{W(c, a, g)\ s \mid c\})) \wedge (\neg (O\{\neg (W(c, a, g)\ s) \mid c\}))) w$

**nitpick** [*user-axioms, falsify=false*] **oops**

— Nitpick found a model for card i = 1 and card s = 2:

Free variables:  $a = (\lambda x. \dots)(s_1 := (\lambda x. \dots)(i_1 := \text{False}), s_2 := (\lambda x. \dots)(i_1 := \text{False}))$   $c = (\lambda x. \dots)(i_1 := \text{False})$   $g = (\lambda x. \dots)(i_1 := \text{False})$   $s = s_2$  Recall that Nitpick is a model checker that finds models making certain formulae true or false. In this case, Nitpick finds a model satisfying the given formula (which simply requires that the sentence “s wills (c, a, g)” is permissible but not obligator). This model consists of the above specifications of a, c, g, and s.

I also expect that any arbitrary maxim should be either permissible or prohibited, since all acts are either permissible or prohibited.

**lemma** *perm-or-prohibited:*

**shows**  $((\text{permissible}(c, a, g)\ s) \vee (\text{prohibited}(c, a, g)\ s)) w$

**by blast**

— This simple test passes immediately by the definitions of permissible and prohibited.

Obligation should be strictly stronger than permissibility. In other words, if a maxim is obligated at a world, it should be permissible at that world. Below I test this property.

**lemma** *obligated-then-permissible:*

**shows**  $(O\{W(c, a, g)\ s \mid c\} \rightarrow (\text{permissible}(c, a, g)\ s)) w$

**using** *no-contradictions* **by auto**

— This test passes and Isabelle is able to find a proof for the fact that all obligatory maxims are also permissible.

The above test failed under Kroy’s formalization of the categorical imperative and is thus evidence that my formalization improves upon Kroy’s. Interestingly, this new test passes because of the additional added axiom that prohibits contradictory obligations (recall that Kroy’s formalization allowed contradictory obligations). The proof is clear: if maxims are either permissible or prohibited and obligation

contradicts prohibition, then obligation must result in permissibility. Moreover, this property **REQUIRES** that there be no contradictory obligations. Formally, in the base logic DDL, I can show the sentence  $(O\{A\} \wedge O\{\neg A\}) \longrightarrow (\neg(O\{A\} \longrightarrow P\{A\}))w$ <sup>21</sup>. In English, if A and not A are obligated, then A can be obligated but not permissible.

Is there anything philosophically interesting here? Something about the possibility of genuine moral conflict?

Next, I will test if the formalization allows for vacuous obligations or modal collapse. These tests are sanity checks confirmed that the obligation operator is doesn't collapse. First, I will check that any arbitrary term isn't obligated.

**lemma arbitrary-obligations:**

**fixes**  $c A::t$

**shows**  $O\{A|c\} w$

**nitpick** [*user-axioms=true, falsify*] **oops**

— This test passes—Nitpick finds a model where A isn't obligated in circumstances c.

Nitpick found a counterexample for card i = 1 and card s = 2:

Free variables:  $A = (\lambda x.)(i_1 := \text{True})$   $c = (\lambda x.)(i_1 := \text{False})$  Previous iterations of this test used the monadic obligation operator, which tests the term in the context “True” (equivalently the set of all worlds since True holds everywhere). In this iteration, I test the term in a context c, because my formalization uses the dyadic obligation operator and must thus specify circumstances.

This is exactly the expected result. Any arbitrary action A isn't obligated. A slightly stronger property is “modal collapse,” or whether or not ‘A happens’ implies ‘A is obligated’. The proposition below should be falsifiable.

**lemma modal-collapse:**

**shows**  $((W(c, a, g) s) w) \longrightarrow O\{W(c, a, g) s|c\} w$

**nitpick** [*user-axioms=true, falsify*] **oops**

— Nitpick finds a counterexample, so willing doesn't imply obligation, so this test passes.

Nitpick found a counterexample for card i = 1 and card s = 2:

Free variables:  $a = (\lambda x.)(s_1 := (\lambda x.)(i_1 := \text{False}), s_2 := (\lambda x.)(i_1 := \text{False}))$   $c = (\lambda x.)(i_1 := \text{False})$   $g = (\lambda x.)(i_1 := \text{False})$   $s = s_2$   $w = i_1$  Once again, I modify this test to use the dyadic obligation operator instead of the monadic operator.

The final set of tests deal with ought implies can and conflicting obligations. Recall that I specifically added an axiom in my formalization to disallow contradictory obligations, so I expect these tests to pass. Kroy's formalization fails these tests, so this is another area of improvement over Kroy's formalization.

**lemma ought-implies-can:**

**shows**  $O\{W(c, a, g) s|c\} w \longrightarrow (\Diamond W(c, a, g) s) w$

**using** *O-diamond* **by** *blast*

— This test is a lemma of DDL itself, so it's no surprise that this test passes.

**lemma conflicting-obligations:**

---

<sup>21</sup>Sledgehammer can show this sentence to be true in DDL and in Kroy's formalization.

**shows**  $\neg (O\{W(c, a, g) s | c\} \wedge O\{\neg(W(c, a, g) s) | c\}) w$   
**using** *no-contradictions* **by** *blast*  
— This test passes immediately by the new axiom prohibited contradictory obligations.

**lemma** *implied-contradiction*:

**assumes**  $((W(c1, a1, g1) s) \wedge (W(c2, a2, g2) s)) \rightarrow \perp) w$   
**shows**  $\neg (O\{W(c1, a1, g1) s | c\} \wedge O\{W(c2, a2, g2) s | c\}) w$   
**using** *assms no-contradictions* **by** *blast*  
— Recall that the we also expect the stronger property that the combination of obligatory maxims can't imply a contradiction. The added axiom also makes this test pass.

The metaethical test suite ran on both Kroy's formalization and my formalizaion show two clear improvements. First, my formalization shows that obligatory maxims are permissible, whereas Kroy's paradoxically does not. Second, my formalization doesn't allow contradictory maxims, but Kroy's does. Both of these improvements are derived from the new axiom I added in my formalization that disallows contradictory obligations. Additionally, my formalization also improves on Kroy's by staying faithful to the strongest interpretation of the FUL, Korsgaard's practical contradiction interpretation. (maybe stick philosophical writing here or above?)

### 3.5 Formalization Specific Tests

In this section, I explore tests specific to my formalization of the categorical imperative. First, in my previous (buggy) implementation of DDL, prohibiting contradictory obligation led to the strange result that all permissible actions are obligatory. I will test if this bug appears in this implementation as well.

**lemma** *bug*:

**shows** *permissible*  $(c, a, g) s w \longrightarrow O\{W(c, a, g) s | c\} w$   
**nitpick**<sub>[user-axioms]</sub> **oops**  
— Nitpick found a counterexample for card i = 1 and card s = 1:  
Free variables: a =  $(\lambda x. \dots)(s_1 := (\lambda x. \dots)(i_1 := \text{False}))$  c =  $(\lambda x. \dots)(i_1 := \text{False})$  g =  $(\lambda x. \dots)(i_1 := \text{False})$  s =  $s_1$  w = **undefined** This strange result does not hold; good!

## 4 Related Work

In 1685, Leibniz dreamed of a universal calculator that could be used to resolve philosophical and theological disputes. At the time, the logical and computational resources necessary to make his dream a reality did not exist. Today, automated ethics is a growing field, spurred in part by the need for ethically intelligent AI agents.

Tolmeijer et al. [Tolmeijer et al., 2021] developed a taxonomy of works in implementing machine ethics. An implementation is characterized by (1) the choice of



ethical theory, (2) implementation design decisions (e.g. testing), and (3) implementation details (e.g. choice of logic).

In this paper, I formalize Kantian ethics. There is a long line of work implementing other kinds of ethical theories, like consequentialism [Abel et al., 2016, Anderson et al., 2004] or particularism [Ashley and McLaren, 1994, Guarini, 2006]. Kantian ethics is a deontological, or rule based ethic, and there is also prior work implementing other deontological theories [Govindarajulu and Bringsjord, 2017, Anderson and Anderson, , Anderson and Anderson, 2014]. Kantian ethics specifically appears to be an intuitive candidate for formalization and implementation [Powers, 2006, Lin et al., 2012]. In 2006, Powers [Powers, 2006] argued that an implementation of of Kantian ethics presented technical challenges, such as automation of a non-monotonic logic, and philosophical challenges, like a definition of the categorical imperative. There has also been prior work in formalizing Kantian metaphysics using I/O logic [Stephenson et al., 2019]. Deontic logic itself is inspired by Kant’s “ought implies can” principle, but it does not include a robust formalization of the entire categorical imperative [Cresswell and Hughes, 1996].

Lindner and Bentzen [Bentzen and Lindner, 2018] have presented a formalization and implementation of Kant’s second formulation of the categorical imperative using a custom logic. They present their goal as “not to get close to a correct interpretation of Kant, but to show that our interpretation of Kant’s ideas can contribute to the development of machine ethics.” My work aims to formalize Kant’s ethic as faithfully as possible. I draw on the centuries of work in moral philosophy, as opposed to developing my own ethical theory. I also hope to formalize the first and third formulations of the categorical imperative, in addition to the first.

The implementation of this paper builds on Benzmueller, Parent, and Farjami’s work with the LogiKey framework for machine ethics [Benzmüller et al., 2021, Benzmüller et al., 2019]. The LogiKey project has been used to implement metaphysics [Benzmüller and Paleo, 2013, Kirchner et al., 2019]. Fuenmayor and Benzmueller [Fuenmayor and Benzmüller, 2018] have implemented Gewirth’s principle of generic consistency, which is similar to Kant’s formula of universal law.

## 5 Future Work

I intend to continue this research for the next year as part of my senior thesis. To make that process easier, I will sketch some goals for the rest of the project. In Section 3.2, I present a young and unfinished implementation of Kroy’s formalization of the categorical imperative. The finished version of my project will ideally include an implementation of Kroy’s formalization of the second formulation of the categorical imperative as well. I also hope to write robust tests for both of these implementations to explore their limitations. These tests will help inform my eventual formalization of the categorical imperative.

The ultimate goal of the project is to present my own formalization of the categor-

ical imperative that escapes the limitations of the naive formalization and Kroy’s formalization. This formalization will likely require some additional logical machinery to handle the complete notion of a maxim, including an agent, action, and end. My formalization will also patch up some of the holes in DDL itself that have been problematic for my project so far, such as the existence of contradictory obligations. I intend to formalize and implement all three formulations of the categorical imperative.

I will then test my formalization of the categorical imperative. I will create two kinds of tests. First, I will create metaethical tests that show logical properties independent of any model specification, as I did for the first two formalizations. Second, I will create tests that specify models and apply my formalization to real, concrete ethical dilemmas. This part of the project will seek to demonstrate the power and limitations of automated ethical reasoning. Questions to be explored here include: How much model specification is necessary to achieve ethical results? How should models be represented and specified? Does the automation of ethical reasoning provide anything, or is all the ethical work hidden in the model specification itself?

This final question is both technical and philosophical, and will be interesting to explore in the written component of my thesis. This question is related to Kant’s distinction between analytic and synthetic reasoning [Kant, 1785]. Analytic statements are true simply by virtue of their meaning, such as “All bachelors are unmarried.” Synthetic reasoning involves some contribution by the reasoner, in the form of new insight or facts about the world. Kant presents the statements “All bachelors are alone” and “ $7+5=12$ ” as examples of synthetic propositions. The analytic/synthetic distinction is hotly debated and has been refined significantly since Kant, and this area will require further research.

Kant believes that ethics is synthetic a priori reasoning, but it is unclear if automated theorem provers like Isabelle are capable of anything more than analytic reasoning. Many of the basic proof solving tools like `simp` or `blast` simply unfold definitions and apply axioms, and they appear to perform analytic reasoning. SMT solvers like `Nitpick` and `z3` (bundled with Isabelle) are candidates for synthetic reasoning. Model finding seems more sophisticated than the simple unfolding of definitions, but this requires further exploration.

Lastly, I hope to explore Kant’s argument that the three formulations of the categorical imperative are equivalent. This hypothesis has been the subject of controversy, but many neo-Kantians believe that his claim is plausible, if not true. Armed with formalizations of each formulation, I will have all the tools necessary to test this hypothesis. I would like to either prove or disprove this hypothesis for my formalization, and analyze the philosophical implications of my result.

## References

- [Abel et al., 2016] Abel, D., MacGlashan, J., and Littman, M. (2016). Reinforcement learning as a framework for ethical decision making. In *AAAI Workshop: AI, Ethics, and Society*.
- [Anderson and Anderson, 2014] Anderson, M. and Anderson, S. (2014). Geneth: A general ethical dilemma analyzer. volume 1.
- [Anderson et al., 2004] Anderson, M., Anderson, S., and Armen, C. (2004). Towards machine ethics.
- [Anderson and Anderson, ] Anderson, M. and Anderson, S. L. Ethel: Toward a principled ethical eldercare robot.
- [Ashley and McLaren, 1994] Ashley, K. D. and McLaren, B. M. (1994). A cbr knowledge representation for practical ethics. In *Selected Papers from the Second European Workshop on Advances in Case-Based Reasoning, EWCBR '94*, page 181–197, Berlin, Heidelberg. Springer-Verlag.
- [Bentzen and Lindner, 2018] Bentzen, M. M. and Lindner, F. (2018). A formalization of kant’s second formulation of the categorical imperative. *CoRR*, abs/1801.03160.
- [Benzmüller et al., 2021] Benzmüller, C., Farjami, A., and Parent, X. (2021). Dyadic deontic logic in hol: Faithful embedding and meta-theoretical experiments. In Armgardt, M., Nordtveit Kvernenes, H. C., and Rahman, S., editors, *New Developments in Legal Reasoning and Logic: From Ancient Law to Modern Legal Systems*, volume 23 of *Logic, Argumentation & Reasoning*. Springer Nature Switzerland AG.
- [Benzmüller and Paleo, 2013] Benzmüller, C. and Paleo, B. W. (2013). Formalization, mechanization and automation of gödel’s proof of god’s existence. *CoRR*, abs/1308.4526.
- [Benzmüller et al., 2019] Benzmüller, C., Parent, X., and van der Torre, L. W. N. (2019). Designing normative theories of ethical reasoning: Formal framework, methodology, and tool support. *CoRR*, abs/1903.10187.
- [Blanchette et al., 2011] Blanchette, J. C., Böhme, S., and Paulson, L. C. (2011). Extending sledgehammer with smt solvers. In *Proceedings of the 23rd International Conference on Automated Deduction, CADE’11*, page 116–130, Berlin, Heidelberg. Springer-Verlag.
- [Blanchette and Nipkow, 2010] Blanchette, J. C. and Nipkow, T. (2010). *Nitpick: A Counterexample Generator for Higher-Order Logic Based on a Relational Model Finder*, volume 6172, page 131–146. Springer Berlin Heidelberg.

- [Bok, 1998] Bok, H. (1998). *Freedom and Responsibility*. Princeton University Press.
- [Carmo and Jones, 2013] Carmo, J. and Jones, A. (2013). Completeness and decidability results for a logic of contrary-to-duty conditionals. *J. Log. Comput.*, 23:585–626.
- [Chisholm, 1963] Chisholm, R. M. (1963). Contrary-to-duty imperatives and deontic logic. *Analysis (Oxford)*, 24(2):33–36.
- [Cresswell and Hughes, 1996] Cresswell, M. J. and Hughes, G. E. (1996). *A New Introduction to Modal Logic*. Routledge.
- [Fuenmayor and Benzmüller, 2018] Fuenmayor, D. and Benzmüller, C. (2018). Formalisation and evaluation of alan gewirth’s proof for the principle of generic consistency in isabelle/hol. *Archive of Formal Proofs*. <https://isa-afp.org/entries/GewirthPGCProof.html>, Formal proof development.
- [Govindarajulu and Bringsjord, 2017] Govindarajulu, N. S. and Bringsjord, S. (2017). On automating the doctrine of double effect. *CoRR*, abs/1703.08922.
- [Guarini, 2006] Guarini, M. (2006). Particularism and the classification and reclassification of moral cases. *IEEE Intelligent Systems*, 21(4):22–28.
- [Hintikka, 1962] Hintikka, J. (1962). *Knowledge and Belief*. Cornell University Press.
- [Kant, 1785] Kant, I. (1785). *Groundwork of the Metaphysics of Morals*. Cambridge University Press, Cambridge.
- [Kant, 2017] Kant, I. (2017). *Introduction*, pages ix–xxix. Cambridge Texts in the History of Philosophy. Cambridge University Press, 2 edition.
- [Kemp, 1958] Kemp, J. (1958). Kant’s examples of the categorical imperative. *The Philosophical Quarterly (1950-)*, 8(30):63–71.
- [Kirchner et al., 2019] Kirchner, D., Benzmüller, C., and Zalta, E. N. (2019). Computer science and metaphysics: A cross-fertilization. *CoRR*, abs/1905.00787.
- [Kitcher, 2003] Kitcher, P. (2003). What is a maxim? *Philosophical Topics*, 31(1/2):215–243.
- [Kitcher, 2004] Kitcher, P. (2004). Kant’s argument for the categorical imperative. *Nous*, 38.
- [Kleingeld, 2017] Kleingeld, P. (2017). Contradiction and kant’s formula of universal law. *Kant-Studien*, 108(1):89–115.

- [Kohl, 2015] Kohl, M. (2015). Kant and 'ought implies can'. *The Philosophical Quarterly* (1950-), 65(261):690–710.
- [Korsgaard, 1985] Korsgaard, C. (1985). Kant's Formula of Universal Law. *Pacific Philosophical Quarterly*, 66:24–47.
- [Korsgaard, 1986] Korsgaard, C. (1986). The Right to Lie: Kant on Dealing with Evil. *Philosophy and Public Affairs*, 15:325–249.
- [Korsgaard, 2005] Korsgaard, C. M. (2005). Acting for a reason. *Danish Yearbook of Philosophy*, 40(1):11–35.
- [Korsgaard and O'Neill, 1996] Korsgaard, C. M. and O'Neill, O. (1996). *The Sources of Normativity*. Cambridge University Press.
- [Kroy, 1976] Kroy, M. (1976). A partial formalization of kant's categorical imperative. an application of deontic logic to classical moral philosophy. *Kant-Studien*, 67(1-4):192–209.
- [Lewis, 1973a] Lewis, D. (1973a). Causation. *Journal of Philosophy*, 70(17):556–567.
- [Lewis, 1973b] Lewis, D. (1973b). *Counterfactuals*. Blackwell.
- [Lin et al., 2012] Lin, P., Abney, K., and Bekey, G. A. (2012). *Robotics, Ethical Theory, and Metaethics: A Guide for the Perplexed*, pages 35–52.
- [McNamara and Van De Putte, 2021] McNamara, P. and Van De Putte, F. (2021). Deontic Logic. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, spring 2021 edition.
- [Menzies and Beebe, 2020] Menzies, P. and Beebe, H. (2020). Counterfactual Theories of Causation. In Zalta, E. N., editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, Winter 2020 edition.
- [MONTAGUE, 1970] MONTAGUE, R. (1970). Universal grammar. *Theoria*, 36(3):373–398.
- [Nipkow et al., 2002] Nipkow, T., Paulson, L. C., and Wenzel, M. (2002). *Isabelle/HOL: A Proof Assistant for Higher Order Logic*. Springer-Verlag Berlin Heidelberg, Berlin.
- [O'Neill, 2009] O'Neill, O. (2009). *Bounds of Justice*. Cambridge University Press.
- [O'Neill, 2013] O'Neill, O. (2013). *Acting on Principle: An Essay on Kantian Ethics*. Cambridge University Press.

- [Paulson and Blanchette, 2015] Paulson, L. and Blanchette, J. (2015). Three years of experience with sledgehammer, a practical link between automatic and interactive theorem provers.
- [Powers, 2006] Powers, T. M. (2006). Prospects for a kantian machine. *IEEE Intelligent Systems*, 21(4):46–51.
- [Rønnedal, 2019] Rønnedal, D. (2019). Contrary-to-duty paradoxes and counterfactual deontic logic. *Philosophia*, 47.
- [Scott, 1970] Scott, D. (1970). Advice on modal logic. In Lambert, K., editor, *Philosophical Problems in Logic: Some Recent Developments*, pages 143–173. D. Reidel.
- [Solt, 1984] Solt, K. (1984). Deontic alternative worlds and the truth-value of ‘*oa*’. *Logique et Analyse*, 27(107):349–351.
- [Stephenson et al., 2019] Stephenson, A., Sergot, M., and Evans, R. (2019). Formalizing kant’s rules: a logic of conditional imperatives and permissives. *Journal of Philosophical Logic*, 49.
- [Timmermann, 2013a] Timmermann, J. (2013a). Kantian dilemmas? moral conflict in kant’s ethical theory. *Archiv für Geschichte der Philosophie*, 95(1):36–64.
- [Timmermann, 2013b] Timmermann, J. (2013b). Kantian dilemmas? moral conflict in kant’s ethical theory. *Archiv für Geschichte der Philosophie*, 95(1):36–64.
- [Tolmeijer et al., 2021] Tolmeijer, S., Kneer, M., Sarasua, C., Christen, M., and Bernstein, A. (2021). Implementations in machine ethics. *ACM Computing Surveys*, 53(6):1–38.
- [Velleman, 2005] Velleman, J. D. (2005). *A Brief Introduction to Kantian Ethics*, page 16–44. Cambridge University Press.
- [Wood, 1999] Wood, A. W. (1999). *Kant’s Ethical Thought*. Cambridge University Press.