
A DYADIC DEONTIC LOGIC IN HOL

CHRISTOPH BENZMÜLLER

University of Luxembourg, Luxembourg, and Freie Universität Berlin, Germany
c.benzmueller@gmail.com

ALI FARJAMI

University of Luxembourg, Luxembourg
farjami110@gmail.com

XAVIER PARENT

University of Luxembourg, Luxembourg
xavier.parent@uni.lu

Abstract

A shallow semantical embedding of a dyadic deontic logic by Carmo and Jones in classical higher-order logic is presented. This embedding is proven sound and complete, that is, faithful.

The work presented here provides the theoretical foundation for the implementation and automation of dyadic deontic logic within off-the-shelf higher-order theorem provers and proof assistants.

Keywords: Logic of CTD conditionals by Carmo and Jones; Classical higher-order logic; Semantic embedding; Automated reasoning

1 Introduction

Dyadic deontic logic is the logic for reasoning with dyadic obligations ("it ought to be the case that ... if it is the case that ..."). A particular dyadic deontic logic, tailored to so-called contrary-to-duty (CTD) conditionals, has been proposed by Carmo and Jones [13]. We shall refer to it as DDL in the remainder. DDL comes with a neighborhood semantics and a weakly complete axiomatization over the class of finite models. The framework is immune to the well-known CTD paradoxes, like

This work has been supported by the European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No 690974.

Chisholm’s paradox [14, 19], and other related puzzles. However, the question of how to mechanise and automate reasoning tasks in DDL has not been studied yet.

This article addresses this challenge. We essentially devise a faithful semantical embedding of DDL in classical higher-order logic (HOL). The latter logic thereby serves as an universal meta-logic. Analogous to successful, recent work in the area of computational metaphysics (cf. [6] and the references therein), the key motivation is to mechanise and automate DDL on the computer by reusing existing theorem proving technology for meta-logic HOL. The embedding of DDL in HOL as devised in this article enables just this.

Meta-logic HOL [4], as employed in this article, was originally devised by Church [17], and further developed by Henkin [18] and Andrews [1, 3, 2]. It bases both terms and formulas on simply typed λ -terms. The use of the λ -calculus has some major advantages. For example, λ -abstractions over formulas allow the explicit naming of sets and predicates, something that is achieved in set theory via the comprehension axioms. Another advantage is, that the complex rules for quantifier instantiation at first-order and higher-order types is completely explained via the rules of λ -conversion (the so-called rules of α -, β -, and η -conversion) which were proposed earlier by Church [15, 16]. These two advantages are exploited in our embedding of DDL in HOL.

Different notions of semantics for HOL have been thoroughly studied in the literature [7, 20]. In this article we assume HOL with Henkin semantics and choice (cf. the detailed description by Benzmüller et. al. [7]). For this notion of HOL, which does not suffer from Gödel’s incompleteness results, several sound and complete theorem provers have been developed in the past decades [9]. We propose to reuse these theorem provers for the mechanisation and automation of DDL. The semantical embedding as devised in this article provides both the theoretical foundation for the approach and the practical bridging technology that is enabling DDL applications within existing HOL theorem provers.

The article is structured as follows: Section 2 outlines DDL and Sec. 3 introduces HOL. The semantical embedding of DDL in HOL is then devised and studied in Sec. 4. This section also addresses soundness and completeness, but due to space restrictions the proofs can only be sketched here; for details we refer to [8]. Section 5 discusses the implementation and automation of the embedding in Isabelle/HOL [21] and Sec. 6 concludes the paper.

2 The Dyadic Deontic Logic of Carmo and Jones

This section provides a concise introduction of DDL, the dyadic deontic logic proposed by Carmo and Jones. Definitions as required for the remainder are presented. For further details we refer to the literature [13, 12].

To define the formulas of DDL we start with a countable set P of propositional symbols, and we choose \neg and \vee as the only primitive connectives.

The set of *DDL formulas* is given as the smallest set of formulas obeying the following conditions:

- Each $p^j \in P$ is an (atomic) DDL formula.
- Given two arbitrary DDL formulas φ and ψ , then
 - $\neg\varphi$ — *classical negation*,
 - $\varphi \vee \psi$ — *classical disjunction*,
 - $\bigcirc(\psi/\varphi)$ — *dyadic deontic obligation: “it ought to be ψ , given φ ”*,
 - $\Box\varphi$ — *in all worlds*,
 - $\Box_a\varphi$ — *in all actual versions of the current world*,
 - $\Box_p\varphi$ — *in all potential versions of the current world*,
 - $\bigcirc_a\varphi$ — *monadic deontic operator for actual obligation*, and
 - $\bigcirc_p\varphi$ — *monadic deontic operator for primary obligation*

are also DDL formulas.

Further logical connectives can be defined as usual: $\varphi \wedge \psi := \neg(\neg\varphi \vee \neg\psi)$, $\varphi \rightarrow \psi := \neg\varphi \vee \psi$, $\varphi \leftrightarrow \psi := (\varphi \rightarrow \psi) \wedge (\psi \rightarrow \varphi)$, $\Diamond\varphi := \neg\Box\neg\varphi$, $\Diamond_a\varphi := \neg\Box_a\neg\varphi$, $\Diamond_p\varphi := \neg\Box_p\neg\varphi$, $\top := \neg q^j \vee q^j$, for some propositional symbol q^j , $\perp := \neg\top$, and $\bigcirc\varphi := \bigcirc(\varphi/\top)$.

A DDL *model* is a structure $M = \langle S, av, pv, ob, V \rangle$, where S is a non empty set of items called possible worlds, V is a function assigning a set of worlds to each atomic formula, that is, $V(p^j) \subseteq S$. $av: S \rightarrow \wp(S)$, where $\wp(S)$ is the power set of S , is a function mapping worlds to sets of worlds such that $av(s) \neq \emptyset$. $av(s)$ is the set of actual versions of the world s . $pv: S \rightarrow \wp(S)$ is another, similar mapping such that $av(s) \subseteq pv(s)$ and $s \in pv(s)$. $pv(s)$ is the set of potential versions of the world s . $ob: \wp(S) \rightarrow \wp(\wp(S))$ is a function mapping sets of worlds to sets of sets of worlds. $ob(\bar{X})$ is the set of propositions that are obligatory in context $\bar{X} \subseteq S$. The following conditions hold for ob (where $\bar{X}, \bar{Y}, \bar{Z}$ designate arbitrary subsets of S):

1. $\emptyset \notin ob(\bar{X})$.
2. If $\bar{Y} \cap \bar{X} = \bar{Z} \cap \bar{X}$, then $\bar{Y} \in ob(\bar{X})$ if and only if $\bar{Z} \in ob(\bar{X})$.

3. Let $\bar{\beta} \subseteq ob(\bar{X})$ and $\bar{\beta} \neq \emptyset$. If $(\cap\bar{\beta}) \cap \bar{X} \neq \emptyset$ (where $\cap\bar{\beta} = \{s \in S \mid \text{for all } \bar{Z} \in \bar{\beta} \text{ we have } s \in \bar{Z}\}$), then $(\cap\bar{\beta}) \in ob(\bar{X})$.
4. If $\bar{Y} \subseteq \bar{X}$ and $\bar{Y} \in ob(\bar{X})$ and $\bar{X} \subseteq \bar{Z}$, then $(\bar{Z} \setminus \bar{X}) \cup \bar{Y} \in ob(\bar{Z})$.
5. If $\bar{Y} \subseteq \bar{X}$ and $\bar{Z} \in ob(\bar{X})$ and $\bar{Y} \cap \bar{Z} \neq \emptyset$, then $\bar{Z} \in ob(\bar{Y})$.

Satisfiability of a formula φ for a model $M = \langle S, av, pv, ob, V \rangle$ and a world $s \in S$ is expressed by writing that $M, s \models \varphi$ and we define $V^M(\varphi) = \{s \in S \mid M, s \models \varphi\}$. In order to simplify the presentation, whenever the model M is obvious from context, we write $V(\varphi)$ instead of $V^M(\varphi)$. Moreover, we often use “iff” as shorthand for “if and only if”.

$M, s \models p^j$	iff	$s \in V(p^j)$
$M, s \models \neg\varphi$	iff	$M, s \not\models \varphi$ (that is, not $M, s \models \varphi$)
$M, s \models \varphi \vee \psi$	iff	$M, s \models \varphi$ or $M, s \models \psi$
$M, s \models \Box\varphi$	iff	$V(\varphi) = S$
$M, s \models \Box_a\varphi$	iff	$av(s) \subseteq V(\varphi)$
$M, s \models \Box_p\varphi$	iff	$pv(s) \subseteq V(\varphi)$
$M, s \models \bigcirc(\psi/\varphi)$	iff	$V(\psi) \in ob(V(\varphi))$
$M, s \models \bigcirc_a\varphi$	iff	$V(\varphi) \in ob(av(s))$ and $av(s) \cap V(\neg\varphi) \neq \emptyset$
$M, s \models \bigcirc_p\varphi$	iff	$V(\varphi) \in ob(pv(s))$ and $pv(s) \cap V(\neg\varphi) \neq \emptyset$

Our evaluation rule for $\bigcirc(_/_)$ is a simplified version of the one used by Carmo and Jones. Given the constraints placed on ob , the two rules are equivalent (cf. [5, result II-2-2]).

As usual, a DDL formula φ is *valid in a DDL model* $M = \langle S, av, pv, ob, V \rangle$, i.e. $M \models^{DDL} \varphi$, if and only if for all worlds $s \in S$ we have $M, s \models \varphi$. A formula φ is *valid*, denoted $\models^{DDL} \varphi$, if and only if it is valid in every DDL model.

3 Classical Higher-order Logic

In this section we introduce classical higher-order logic (HOL). The presentation, which has partly been adapted from [5], is rather detailed in order to keep the article sufficiently self-contained.

3.1 Syntax of HOL

For defining the syntax of HOL, we first introduce the set T of *simple types*. We assume that T is freely generated from a set of *basic types* $BT \supseteq \{o, i\}$ using the

function type constructor \rightarrow . Type o denotes the (bivalent) set of Booleans, and i a non-empty set of individuals.

For the definition of HOL, we start out with a family of denumerable sets of typed constant symbols $(C_\alpha)_{\alpha \in T}$, called the HOL *signature*, and a family of denumerable sets of typed variable symbols $(V_\alpha)_{\alpha \in T}$.¹ We employ Church-style typing, where each term t_α explicitly encodes its type information in subscript α .

The *language of HOL* is given as the smallest set of terms obeying the following conditions.

- Every typed constant symbol $c_\alpha \in C_\alpha$ is a HOL term of type α .
- Every typed variable symbol $X_\alpha \in V_\alpha$ is a HOL term of type α .
- If $s_{\alpha \rightarrow \beta}$ and t_α are HOL terms of types $\alpha \rightarrow \beta$ and α , respectively, then $(s_{\alpha \rightarrow \beta} t_\alpha)_\beta$, called *application*, is an HOL term of type β .
- If $X_\alpha \in V_\alpha$ is a typed variable symbol and s_β is an HOL term of type β , then $(\lambda X_\alpha s_\beta)_{\alpha \rightarrow \beta}$, called *abstraction*, is an HOL term of type $\alpha \rightarrow \beta$.

The above definition encompasses the simply typed λ -calculus. In order to extend this base framework into logic HOL we simply ensure that the signature $(C_\alpha)_{\alpha \in T}$ provides a sufficient selection of primitive logical connectives. Without loss of generality, we here assume the following *primitive logical connectives* to be part of the signature: $\neg_{o \rightarrow o} \in C_{o \rightarrow o}$, $\vee_{o \rightarrow o \rightarrow o} \in C_{o \rightarrow o \rightarrow o}$, $\Pi_{(\alpha \rightarrow o) \rightarrow o} \in C_{(\alpha \rightarrow o) \rightarrow o}$ and $=_{\alpha \rightarrow \alpha \rightarrow \alpha} \in C_{\alpha \rightarrow \alpha \rightarrow \alpha}$, abbreviated as $=^\alpha$. The symbols $\Pi_{(\alpha \rightarrow o) \rightarrow o}$ and $=_{\alpha \rightarrow \alpha \rightarrow \alpha}$ are generally assumed for each type $\alpha \in T$. The denotation of the primitive logical connectives is fixed below according to their intended meaning. *Binder notation* $\forall X_\alpha s_o$ is used as an abbreviation for $\Pi_{(\alpha \rightarrow o) \rightarrow o} \lambda X_\alpha s_o$. Universal quantification in HOL is thus modeled with the help of the logical constants $\Pi_{(\alpha \rightarrow o) \rightarrow o}$ to be used in combination with lambda-abstraction. That is, the only binding mechanism provided in HOL is lambda-abstraction.

HOL is a logic of terms in the sense that the *formulas of HOL* are given as the terms of type o . In addition to the primitive logical connectives selected above, we could assume *choice operators* $\epsilon_{(\alpha \rightarrow o) \rightarrow \alpha} \in C_{(\alpha \rightarrow o) \rightarrow \alpha}$ (for each type α) in the signature. We are not pursuing this here.

Type information as well as brackets may be omitted if obvious from the context, and we may also use infix notation to improve readability. For example, we may write $(s \vee t)$ instead of $((\vee_{o \rightarrow o \rightarrow o} s_o) t_o)$.

¹For example in Section 4 we will assume constant symbols *av*, *pv* and *ob* with types $i \rightarrow i \rightarrow o$, $i \rightarrow i \rightarrow o$ and $(i \rightarrow o) \rightarrow (i \rightarrow o) \rightarrow o$ as part of the signature.

From the selected set of primitive connectives, other logical connectives can be introduced as abbreviations.² For example, we may define $s \wedge t := \neg(\neg s \vee \neg t)$, $s \rightarrow t := \neg s \vee t$, $s \longleftrightarrow t := (s \rightarrow t) \wedge (t \rightarrow s)$, $\top := (\lambda X_i X) = (\lambda X_i X)$, $\perp := \neg \top$ and $\exists X_\alpha s := \neg \forall X_\alpha \neg s$.

The notions of *free variables*, α -*conversion*, $\beta\eta$ -*equality* (denoted as $=_{\beta\eta}$) and *substitution* of a term s_α for a variable X_α in a term t_β (denoted as $[s/X]t$) are defined as usual.

3.2 Semantics of HOL

The semantics of HOL is well understood and thoroughly documented. The introduction provided next focuses on the aspects as needed for this article. For more details we refer to the previously mentioned literature [7].

The semantics of choice for the remainder is Henkin semantics, i.e., we work with Henkin's general models [18]. Henkin models (and standard models) are introduced next. We start out with introducing frame structures.

A *frame* D is a collection $\{D_\alpha\}_{\alpha \in \mathbb{T}}$ of nonempty sets D_α , such that $D_o = \{T, F\}$ (for truth and falsehood). The $D_{\alpha \rightarrow \beta}$ are collections of functions mapping D_α into D_β .

A *model* for HOL is a tuple $M = \langle D, I \rangle$, where D is a frame, and I is a family of typed interpretation functions mapping constant symbols $p_\alpha \in C_\alpha$ to appropriate elements of D_α , called the *denotation of p_α* . The logical connectives \neg , \vee , Π and $=$ are always given their expected, standard denotations:³

- $I(\neg_{o \rightarrow o}) = \text{not} \in D_{o \rightarrow o}$ such that $\text{not}(T) = F$ and $\text{not}(F) = T$.
- $I(\vee_{o \rightarrow o \rightarrow o}) = \text{or} \in D_{o \rightarrow o \rightarrow o}$ such that $\text{or}(a, b) = T$ iff $(a = T \text{ or } b = T)$.
- $I(=_{\alpha \rightarrow \alpha \rightarrow o}) = \text{id} \in D_{\alpha \rightarrow \alpha \rightarrow o}$ such that for all $a, b \in D_\alpha$, $\text{id}(a, b) = T$ iff a is identical to b .

²As demonstrated by Andrews [4], we could in fact start out with only primitive equality in the signature (for all types α) and introduce all other logical connectives as abbreviations based on it. Alternatively, we could remove primitive equality from the above signature, since equality can be defined in HOL from these other logical connectives by exploiting Leibniz' principle, expressing that two objects are equal if they share the same properties. *Leibniz equality* \doteq^α at type α is thus defined as $s_\alpha \doteq^\alpha t_\alpha := \forall P_{\alpha \rightarrow o}(Ps \longleftrightarrow Pt)$. The motivation for the redundant signature as selected here is to stay close to the the choices taken in implemented theorem provers such as LEO-II and Leo-III and also to theory paper [7], which is recommended for further details.

³Since $=_{\alpha \rightarrow \alpha \rightarrow o}$ (for all types α) is in the signature, it is ensured that the domains $D_{\alpha \rightarrow \alpha \rightarrow o}$ contain the respective identity relations. This addresses an issue discovered by Andrews [2]: if such identity relations did not existing in the $D_{\alpha \rightarrow \alpha \rightarrow o}$, then Leibniz equality in Henkin semantics might not denote as intended.

- $I(\Pi_{(\alpha \rightarrow o) \rightarrow o}) = \text{all} \in D_{(\alpha \rightarrow o) \rightarrow o}$ such that for all $s \in D_{\alpha \rightarrow o}$, $\text{all}(s) = T$ iff $s(a) = T$ for all $a \in D_\alpha$; i.e., s is the set of all objects of type α .

Variable assignments are a technical aid for the subsequent definition of an interpretation function $\|\cdot\|^{M,g}$ for HOL terms. This interpretation function is parametric over a model M and a variable assignment g .

A *variable assignment* g maps variables X_α to elements in D_α . $g[d/W]$ denotes the assignment that is identical to g , except for variable W , which is now mapped to d .

The *denotation* $\|s_\alpha\|^{M,g}$ of an HOL term s_α on a model $M = \langle D, I \rangle$ under assignment g is an element $d \in D_\alpha$ defined in the following way:

$$\begin{aligned} \|p_\alpha\|^{M,g} &= I(p_\alpha) \\ \|X_\alpha\|^{M,g} &= g(X_\alpha) \\ \|(s_{\alpha \rightarrow \beta} t_\alpha)_\beta\|^{M,g} &= \|s_{\alpha \rightarrow \beta}\|^{M,g}(\|t_\alpha\|^{M,g}) \\ \|(\lambda X_\alpha s_\beta)_{\alpha \rightarrow \beta}\|^{M,g} &= \text{the function } f \text{ from } D_\alpha \text{ to } D_\beta \text{ such that} \\ &\quad f(d) = \|s_\beta\|^{M,g[d/X_\alpha]} \text{ for all } d \in D_\alpha \end{aligned}$$

A model $M = \langle D, I \rangle$ is called a *standard model* if and only if for all $\alpha, \beta \in T$ we have $D_{\alpha \rightarrow \beta} = \{f \mid f : D_\alpha \rightarrow D_\beta\}$. In a *Henkin model (general model)* function spaces are not necessarily full. Instead it is only required that for all $\alpha, \beta \in T$, $D_{\alpha \rightarrow \beta} \subseteq \{f \mid f : D_\alpha \rightarrow D_\beta\}$. However, it is required that the valuation function $\|\cdot\|^{M,g}$ from above is total, so that every term denotes. Note that this requirement, which is called *Denotatpflicht*, ensures that the function domains $D_{\alpha \rightarrow \beta}$ never become too sparse, that is, the denotations of the lambda-abstractions as devised above are always contained in them.

Corollary 1. *For any Henkin model $M = \langle D, I \rangle$ and variable assignment g :*

1. $\|(\neg_{o \rightarrow o} s_o)_o\|^{M,g} = T$ iff $\|s_o\|^{M,g} = F$.
2. $\|((\vee_{o \rightarrow o \rightarrow o} s_o) t_o)_o\|^{M,g} = T$ iff $\|s_o\|^{M,g} = T$ or $\|t_o\|^{M,g} = T$.
3. $\|((\wedge_{o \rightarrow o \rightarrow o} s_o) t_o)_o\|^{M,g} = T$ iff $\|s_o\|^{M,g} = T$ and $\|t_o\|^{M,g} = T$.
4. $\|((\rightarrow_{o \rightarrow o \rightarrow o} s_o) t_o)_o\|^{M,g} = T$ iff (if $\|s_o\|^{M,g} = T$ then $\|t_o\|^{M,g} = T$).
5. $\|((\leftarrow_{o \rightarrow o \rightarrow o} s_o) t_o)_o\|^{M,g} = T$ iff ($\|s_o\|^{M,g} = T$ iff $\|t_o\|^{M,g} = T$).
6. $\|\top\|^{M,g} = T$.
7. $\|\perp\|^{M,g} = F$.

8. $\|(\forall X_\alpha s_o)_o\|^{M,g} = T$ iff for all $d \in D_\alpha$ we have $\|s_o\|^{M,g[d/X_\alpha]} = T$.
9. $\|(\exists X_\alpha s_o)_o\|^{M,g} = T$ iff there exists $d \in D_\alpha$ such that $\|s_o\|^{M,g[d/X_\alpha]} = T$.

Proof. We leave the proof as an exercise to the reader. □

An HOL formula s_o is *true* in an Henkin model M under assignment g if and only if $\|s_o\|^{M,g} = T$; this is also expressed by writing that $M, g \models^{\text{HOL}} s_o$. An HOL formula s_o is called *valid* in M , which is expressed by writing that $M \models^{\text{HOL}} s_o$, if and only if $M, g \models^{\text{HOL}} s_o$ for all assignments g . Moreover, a formula s_o is called *valid*, expressed by writing that $\models^{\text{HOL}} s_o$, if and only if s_o is valid in all Henkin models M . Finally, we define $\Sigma \models^{\text{HOL}} s_o$ for a set of HOL formulas Σ if and only if $M \models^{\text{HOL}} s_o$ for all Henkin models M with $M \models^{\text{HOL}} t_o$ for all $t_o \in \Sigma$.

Note that any standard model is obviously also a Henkin model. Hence, validity of a HOL formula s_o for all Henkin models, implies validity of s_o for all standard models.

4 Modeling DDL as a Fragment of HOL

This section, the core contribution of this article, presents a shallow semantical embedding of DDL in HOL and proves its soundness and completeness. In contrast to a deep logical embedding, where the syntax and semantics of logic L would be formalized in full detail (using structural induction and recursion), only the core differences in the semantics of both DDL and meta-logic HOL are explicitly encoded here.

4.1 Semantical Embedding

DDL formulas are identified in our semantical embedding with certain HOL terms (predicates) of type $i \rightarrow o$. They can be applied to terms of type i , which are assumed to denote possible worlds. That is, the HOL type i is now identified with a (non-empty) set of worlds. Type $i \rightarrow o$ is abbreviated as τ in the remainder. The HOL signature is assumed to contain the constant symbols $av_{i \rightarrow \tau}$, $pv_{i \rightarrow \tau}$ and $ob_{\tau \rightarrow \tau \rightarrow o}$. Moreover, for each propositional symbol p^i of DDL, the HOL signature must contain the corresponding constant symbol p_τ^i . Without loss of generality, we assume that besides those symbols and the primitive logical connectives of HOL, no other constant symbols are given in the signature of HOL.

The mapping $[\cdot]$ translates DDL formulas s into HOL terms $[s]$ of type τ . The mapping is recursively⁴ defined:

$$\begin{aligned}
[p^j] &= p_\tau^j \\
[\neg s] &= \neg_\tau [s] \\
[s \vee t] &= \vee_{\tau \rightarrow \tau \rightarrow \tau} [s] [t] \\
[\Box s] &= \Box_{\tau \rightarrow \tau} [s] \\
[\bigcirc(t/s)] &= \bigcirc_{\tau \rightarrow \tau \rightarrow \tau} [s] [t] \\
[\Box_a s] &= \Box_{\tau \rightarrow \tau}^a [s] \\
[\Box_p s] &= \Box_{\tau \rightarrow \tau}^p [s] \\
[\bigcirc_a s] &= \bigcirc_{\tau \rightarrow \tau}^a [s] \\
[\bigcirc_p s] &= \bigcirc_{\tau \rightarrow \tau}^p [s]
\end{aligned}$$

$\neg_{\tau \rightarrow \tau}$, $\vee_{\tau \rightarrow \tau \rightarrow \tau}$, $\Box_{\tau \rightarrow \tau}$, $\bigcirc_{\tau \rightarrow \tau \rightarrow \tau}$, $\Box_{\tau \rightarrow \tau}^a$, $\Box_{\tau \rightarrow \tau}^p$, $\bigcirc_{\tau \rightarrow \tau}^a$ and $\bigcirc_{\tau \rightarrow \tau}^p$ thereby abbreviate the following HOL terms:

$$\begin{aligned}
\neg_{\tau \rightarrow \tau} &= \lambda A_\tau \lambda X_i \neg(A X) \\
\vee_{\tau \rightarrow \tau \rightarrow \tau} &= \lambda A_\tau \lambda B_\tau \lambda X_i (A X \vee B X) \\
\Box_{\tau \rightarrow \tau} &= \lambda A_\tau \lambda X_i \forall Y_i (A Y) \\
\bigcirc_{\tau \rightarrow \tau \rightarrow \tau} &= \lambda A_\tau \lambda B_\tau \lambda X_i (\text{ob } A B) \\
\Box_{\tau \rightarrow \tau}^a &= \lambda A_\tau \lambda X_i \forall Y_i (\neg(\text{av } X Y) \vee A Y) \\
\Box_{\tau \rightarrow \tau}^p &= \lambda A_\tau \lambda X_i \forall Y_i (\neg(\text{pv } X Y) \vee A Y) \\
\bigcirc_{\tau \rightarrow \tau}^a &= \lambda A_\tau \lambda X_i ((\text{ob } (\text{av } X) A) \wedge \exists Y_i (\text{av } X Y \wedge \neg(A Y))) \\
\bigcirc_{\tau \rightarrow \tau}^p &= \lambda A_\tau \lambda X_i ((\text{ob } (\text{pv } X) A) \wedge \exists Y_i (\text{pv } X Y \wedge \neg(A Y)))
\end{aligned}$$

Analyzing the truth of a translated formula $[s]$ in a world represented by term w_i corresponds to evaluating the application $([s] w_i)$. In line with previous work [10], we define $\text{vld}_{\tau \rightarrow o} = \lambda A_\tau \forall S_i (A S)$. With this definition, validity of a DDL formula s in DDL corresponds to the validity of formula $(\text{vld } [s])$ in HOL, and vice versa.

4.2 Soundness and Completeness

To prove the soundness and completeness, that is, faithfulness, of the above embedding, a mapping from DDL models into Henkin models is employed.

Definition 1 (Henkin model H^M for DDL model M). *For any DDL model $M = \langle S, \text{av}, \text{pv}, \text{ob}, V \rangle$, we define a corresponding Henkin model H^M . Thus, let a DDL model $M = \langle S, \text{av}, \text{pv}, \text{ob}, V \rangle$ be given. Moreover, assume that $p^j \in P$, for $j \geq 1$, are*

⁴A recursive definition is actually not needed in practice. By inspecting the equations below it should become clear that only the abbreviations for the logical connectives of DDL are required in combination with a type-lifting for the propositional constant symbols; cf. also Fig. 1.

the only propositional symbols of DDL. Remember that our embedding requires the corresponding signature of HOL to provide constant symbols p_τ^j such that $\lfloor p^j \rfloor = p_\tau^j$ for $j = 1, \dots, m$.

A Henkin model $H^M = \langle \{D_\alpha\}_{\alpha \in T}, I \rangle$ for M is now defined as follows: D_i is chosen as the set of possible worlds S ; all other sets $D_{\alpha \rightarrow \beta}$ are chosen as (not necessarily full) sets of functions from D_α to D_β . For all $D_{\alpha \rightarrow \beta}$ the rule that every term $t_{\alpha \rightarrow \beta}$ must have a denotation in $D_{\alpha \rightarrow \beta}$ must be obeyed (Denotatpflicht). In particular, it is required that D_τ , $D_{i \rightarrow \tau}$ and $D_{\tau \rightarrow \tau \rightarrow o}$ contain the elements Ip_τ^j , $Iav_{i \rightarrow \tau}$, $Ipv_{i \rightarrow \tau}$ and $Iob_{\tau \rightarrow \tau \rightarrow o}$. The interpretation function I of H^M is defined as follows:

1. For $j = 1, \dots, m$, $Ip_\tau^j \in D_\tau$ is chosen such that $Ip_\tau^j(s) = T$ iff $s \in V(p^j)$ in M .
2. $Iav_{i \rightarrow \tau} \in D_{i \rightarrow \tau}$ is chosen such that $Iav_{i \rightarrow \tau}(s, u) = T$ iff $u \in av(s)$ in M .
3. $Ipv_{i \rightarrow \tau} \in D_{i \rightarrow \tau}$ is chosen such that $Ipv_{i \rightarrow \tau}(s, u) = T$ iff $u \in pv(s)$ in M .
4. $Iob_{\tau \rightarrow \tau \rightarrow o} \in D_{\tau \rightarrow \tau \rightarrow o}$ is such that $Iob_{\tau \rightarrow \tau \rightarrow o}(\bar{X}, \bar{Y}) = T$ iff $\bar{Y} \in ob(\bar{X})$ in M .
5. For the logical connectives \neg , \vee , Π and $=$ of HOL the interpretation function I is defined as usual (see the previous section).

Since we assume that there are no other symbols (besides the p^i , av , pv , ob and \neg , \vee , Π , and $=$) in the signature of HOL, I is a total function. Moreover, the above construction guarantees that H^M is a Henkin model: $\langle D, I \rangle$ is a frame, and the choice of I in combination with the Denotatpflicht ensures that for arbitrary assignments g , $\|\cdot\|^{H^M, g}$ is a total evaluation function.

Lemma 1. Let H^M be a Henkin model for a DDL model M . In H^M we have for all $s \in D_i$ and all $\bar{X}, \bar{Y}, \bar{Z} \in D_\tau$ (cf. the conditions on DDL models as stated on page 3).⁵

- (av) $Iav_{i \rightarrow \tau}(s) \neq \emptyset$.
- (pv1) $Iav_{i \rightarrow \tau}(s) \subseteq Ipv_{i \rightarrow \tau}(s)$.
- (pv2) $s \in Ipv_{i \rightarrow \tau}(s)$.
- (ob1) $\emptyset \notin Iob_{\tau \rightarrow \tau \rightarrow o}(\bar{X})$.
- (ob2) If $\bar{Y} \cap \bar{X} = \bar{Z} \cap \bar{X}$, then $(\bar{Y} \in Iob_{\tau \rightarrow \tau \rightarrow o}(\bar{X}))$ iff $(\bar{Z} \in Iob_{\tau \rightarrow \tau \rightarrow o}(\bar{X}))$.
- (ob3) Let $\bar{\beta} \subseteq Iob_{\tau \rightarrow \tau \rightarrow o}(\bar{X})$ and $\bar{\beta} \neq \emptyset$.
If $(\cap \bar{\beta}) \cap \bar{X} \neq \emptyset$, where $\cap \bar{\beta} = \{s \in S \mid \text{for all } \bar{Z} \in \bar{\beta} \text{ we have } s \in \bar{Z}\}$,
then $(\cap \bar{\beta}) \in Iob_{\tau \rightarrow \tau \rightarrow o}(\bar{X})$.

⁵In the proof in [8] we implicitly employ currying and uncurrying, and we associate sets with their characteristic functions. This analogously applies to the remainder of this article.

- (ob4) If $\bar{Y} \subseteq \bar{X}$ and $\bar{Y} \in \text{Iob}_{\tau \rightarrow \tau \rightarrow o}(\bar{X})$ and $\bar{X} \subseteq \bar{Z}$,
then $(\bar{Z} \setminus \bar{X}) \cup \bar{Y} \in \text{Iob}_{\tau \rightarrow \tau \rightarrow o}(\bar{Z})$.
- (ob5) If $\bar{Y} \subseteq \bar{X}$ and $\bar{Z} \in \text{Iob}_{\tau \rightarrow \tau \rightarrow o}(\bar{X})$ and $\bar{Y} \cap \bar{Z} \neq \emptyset$,
then $\bar{Z} \in \text{Iob}_{\tau \rightarrow \tau \rightarrow o}(\bar{Y})$.

Proof. Each statement follows by construction of H^M for M . □

Lemma 2. Let $H^M = \langle \{D_\alpha\}_{\alpha \in T}, I \rangle$ be a Henkin model for a DDL model M . We have $H^M \models^{\text{HOL}} \Sigma$ for all $\Sigma \in \{AV, PV1, PV2, OB1, \dots, OB5\}$, where

- AV is $\forall W_i \exists V_i (av_{i \rightarrow \tau} W_i V_i)$
- PV1 is $\forall W_i \forall V_i (av_{i \rightarrow \tau} W_i V_i \rightarrow pv_{i \rightarrow \tau} W_i V_i)$
- PV2 is $\forall W_i (pv_{i \rightarrow \tau} W_i W_i)$
- OB1 is $\forall X_\tau \neg ob_{\tau \rightarrow \tau \rightarrow o} X_\tau (\lambda X_\tau \perp)$
- OB2 is $\forall X_\tau Y_\tau Z_\tau ((\forall W_i ((Y_\tau W_i \wedge X_\tau W_i) \longleftrightarrow (Z_\tau W_i \wedge X_\tau W_i)))$
 $\rightarrow (ob_{\tau \rightarrow \tau \rightarrow o} X_\tau Y_\tau \longleftrightarrow ob_{\tau \rightarrow \tau \rightarrow o} X_\tau Z_\tau))$
- OB3 is $\forall \beta_{\tau \rightarrow \tau \rightarrow o} \forall X_\tau$
 $(((\forall Z_\tau (\beta_{\tau \rightarrow \tau \rightarrow o} Z_\tau \rightarrow ob_{\tau \rightarrow \tau \rightarrow o} X_\tau Z_\tau)) \wedge \exists Z_\tau (\beta_{\tau \rightarrow \tau \rightarrow o} Z_\tau))$
 $\rightarrow ((\exists Y_i (((\lambda W_i \forall Z_\tau (\beta_{\tau \rightarrow \tau \rightarrow o} Z_\tau \rightarrow Z_\tau W_i)) Y_i) \wedge X_\tau Y_i))$
 $\rightarrow ob_{\tau \rightarrow \tau \rightarrow o} X_\tau (\lambda W_i \forall Z_\tau (\beta_{\tau \rightarrow \tau \rightarrow o} Z_\tau \rightarrow Z_\tau W_i))))$
- OB4 is $\forall X_\tau Y_\tau Z_\tau$
 $((\forall W_i (Y_\tau W_i \rightarrow X_\tau W_i) \wedge ob_{\tau \rightarrow \tau \rightarrow o} X_\tau Y_\tau \wedge \forall X_\tau (X_\tau W_i \rightarrow Z_\tau W_i))$
 $\rightarrow ob_{\tau \rightarrow \tau \rightarrow o} Z_\tau (\lambda W_i ((Z_\tau W_i \wedge \neg X_\tau W_i) \vee Y_\tau W_i)))$
- OB5 is $\forall X_\tau Y_\tau Z_\tau$
 $((\forall W_i (Y_\tau W_i \rightarrow X_\tau W_i) \wedge ob_{\tau \rightarrow \tau \rightarrow o} X_\tau Z_\tau \wedge \exists W_i (Y_\tau W_i \wedge Z_\tau W_i))$
 $\rightarrow ob_{\tau \rightarrow \tau \rightarrow o} Y_\tau Z_\tau)$

Proof. By construction of H^M for M in combination with Lemma 1. □

Lemma 3. Let H^M be a Henkin model for a DDL model M . For all DDL formulas δ , arbitrary variable assignments g and worlds s it holds:

$$M, s \models \delta \text{ if and only if } \|\llbracket \delta \rrbracket S_i\|^{H^M, g[s/S_i]} = T$$

Proof. By induction on the structure of δ . □

Lemma 4. For every Henkin model $H = \langle \{D_\alpha\}_{\alpha \in T}, I \rangle$ such that $H \models^{\text{HOL}} \Sigma$ for all $\Sigma \in \{AV, PV1, PV2, OB1, \dots, OB5\}$, there exists a corresponding DDL model M . Corresponding means that for all DDL formulas δ and for all assignments g and worlds s , $\|\llbracket \delta \rrbracket S_i\|^{H, g[s/S_i]} = T$ if and only if $M, s \models \delta$.

Proof. Suppose that $H = \langle \{D_\alpha\}_{\alpha \in T}, I \rangle$ is a Henkin model such that $H \models^{\text{HOL}} \Sigma$ for all $\Sigma \in \{\text{AV}, \text{PV1}, \text{PV2}, \text{OB1}, \dots, \text{OB5}\}$. Without loss of generality, we can assume that the domains of H are denumerable [18]. We construct the corresponding DDL model M as follows:

1. $S = D_i$,
2. $u \in av(s)$ for $s, u \in S$ iff $Iav_{i \rightarrow \tau}(s, u) = T$,
3. $u \in pv(s)$ for $s, u \in S$ iff $Ipv_{i \rightarrow \tau}(s, u) = T$,
4. $\bar{Y} \in ob(\bar{X})$ for $\bar{X}, \bar{Y} \in D_i \rightarrow D_o$ iff $Iob_{\tau \rightarrow \tau \rightarrow o}(\bar{X}, \bar{Y}) = T$, and
5. $s \in V(p^j)$ iff $Ipp_\tau^j(s) = T$.

Since $H \models^{\text{HOL}} \Sigma$ for all $\Sigma \in \{\text{AV}, \text{PV1}, \text{PV2}, \text{OB1}, \dots, \text{OB5}\}$, it is straightforward (but tedious) to verify that av , pv and ob satisfy the conditions as required for a DDL model.

Moreover, the above construction ensures that H is a Henkin model H^M for DDL model M . Hence, Lemma 3 applies. This ensures that for all DDL formulas δ , for all assignment g and all worlds s we have $\|\llbracket \delta \rrbracket S_i\|^{H^M, g[s/S_i]} = T$ if and only if $M, s \models \delta$. \square

Theorem 1 (Soundness and Completeness of the Embedding).

$$\models^{DDL} \varphi \text{ if and only if } \{\text{AV}, \text{PV1}, \text{PV2}, \text{OB1}, \dots, \text{OB5}\} \models^{HOL} \text{vld } \llbracket \varphi \rrbracket$$

Proof. (Soundness, \leftarrow) The proof is by contraposition. Assume $\not\models^{DDL} \varphi$, that is, there is a DDL model $M = \langle S, av, pv, ob, V \rangle$, and world $s \in S$, such that $M, s \not\models \varphi$. Now let H^M be a Henkin model for DDL model M . By Lemma 3, for an arbitrary assignment g , it holds that $\|\llbracket \varphi \rrbracket S_i\|^{H^M, g[s/S_i]} = F$. Thus, by definition of $\|\cdot\|$, it holds that $\|\forall S_i(\llbracket \varphi \rrbracket S)\|^{H^M, g} = \|\text{vld } \llbracket \varphi \rrbracket\|^{H^M, g} = F$. Hence, $H^M \not\models^{\text{HOL}} \text{vld } \llbracket \varphi \rrbracket$. Furthermore, $H^M \models^{\text{HOL}} \Sigma$ for all $\Sigma \in \{\text{AV}, \text{PV1}, \text{PV2}, \text{OB1}, \dots, \text{OB5}\}$ by Lemma 2. Thus, $\{\text{AV}, \text{PV1}, \text{PV2}, \text{OB1}, \dots, \text{OB5}\} \not\models^{\text{HOL}} \text{vld } \llbracket \varphi \rrbracket$.

(Completeness, \rightarrow) The proof is again by contraposition. Assume $\{\text{AV}, \text{PV1}, \text{PV2}, \text{OB1}, \dots, \text{OB5}\} \not\models^{\text{HOL}} \text{vld } \llbracket \varphi \rrbracket$, that is, there is a Henkin model $H = \langle \{D_\alpha\}_{\alpha \in T}, I \rangle$ such that $H \models^{\text{HOL}} \Sigma$ for all $\Sigma \in \{\text{AV}, \text{PV1}, \text{PV2}, \text{OB1}, \dots, \text{OB5}\}$, but $\|\text{vld } \llbracket \varphi \rrbracket\|^{H, g} = F$ for some assignment g . By Lemma 4, there is a DDL model M such that $M \not\models \varphi$. Hence, $\not\models^{DDL} \varphi$. \square

Each DDL reasoning problem thus represents a particular HOL problem. The embedding presented in this section, which is based on simple abbreviations, tells us how the two logics are connected.

5 Implementation in Isabelle/HOL

The semantical embedding as devised in Sec. 4 has been implemented in the higher-order proof assistant Isabelle/HOL [21]. Figure 1 displays the respective encoding. Figure 2 applies this encoding to Chisholm’s paradox (cf. [14]), which involves the following four statements:

1. It ought to be that Jones goes to assist his neighbors;
2. It ought to be that if Jones goes, then he tells them he is coming;
3. If Jones doesn’t go, then he ought not tell them he is coming;
4. Jones doesn’t go.

These statements can be given a consistent formalisation in DDL see Fig. 2. This is confirmed by the model finder Nitpick [11] integrated with Isabelle/HOL. Nitpick computes an intuitive, small model for the scenario consisting of two possible worlds i_1 and i_2 . Function ob is interpreted in this model as follows:

$$\begin{aligned} ob(\{i_1, i_2\}) &= \{\{i_1, i_2\}, \{i_1\}\} \\ ob(\{i_1\}) &= \{\{i_1, i_2\}, \{i_1\}\} \\ ob(\{i_2\}) &= \{\{i_1, i_2\}, \{i_2\}\} \\ ob(\emptyset) &= \emptyset \end{aligned}$$

The designated current world in the given model is i_2 , in which Jones doesn’t go to assist his neighbors and doesn’t tell them that he is coming. In the other possible world i_1 , Jones is going to assist them and he also tells them that he his coming. That is, $V(go) = V(tell) = \{i_1\}$. Also, we have $\{i_1\} \in ob(\{i_1, i_2\})$. So, $i_2 \models \bigcirc go$ by the evaluation rule for \bigcirc . Similarly, $\{i_1\} \in ob(\{i_1\})$ implies $i_2 \models \bigcirc(tell/go)$, and $\{i_2\} \in ob(\{i_2\})$ implies $i_2 \models \bigcirc(\neg tell/\neg go)$.

For further experiments, focusing on the automation of meta-theoretic aspects of DDL, we refer to [8, Fig. 2 and Fig. 3].

6 Conclusion

A shallow semantical embedding of Carmo and Jones’s logic of contrary-to-duty conditionals in classical higher-order logic has been presented, and shown to be faithful (sound and complete). This theory work has meanwhile been implemented in the proof assistant Isabelle/HOL. This implementation constitutes the first theorem

prover for the logic by Carmo and Jones that is available to date. The foundational theory for this implementation has been laid in this article.

There is much room for future work. First, experiments could investigate whether the provided implementation already supports non-trivial applications in practical normative reasoning, or whether further emendations and improvements are required. Second, the introduced framework could also be used to systematically analyse the properties of Carmo and Jones’s dyadic deontic logic within Isabelle/HOL. Third, analogous to previous work in modal logic [10], the provided framework could be extended to study and support first-order and higher-order variants of the framework.

Acknowledgements

We thank the anonymous reviewers for their valuable feedback and comments.

References

- [1] P.B. Andrews. Resolution in type theory. *Journal of Symbolic Logic*, 36(3):414–432, 1971.
- [2] P.B. Andrews. General models and extensionality. *Journal of Symbolic Logic*, 37(2):395–397, 1972.
- [3] P.B. Andrews. General models, descriptions, and choice in type theory. *Journal of Symbolic Logic*, 37(2):385–394, 1972.
- [4] P.B. Andrews. Church’s type theory. In E.N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, spring 2014 edition, 2014.
- [5] C. Benzmüller. Cut-elimination for quantified conditional logic. *Journal of Philosophical Logic*, 46(3):333–353, 2017.
- [6] C. Benzmüller. Recent successes with a meta-logical approach to universal logical reasoning (extended abstract). In S.A. da Costa Cavalheiro and J.L. Fiadeiro, editors, *Formal Methods: Foundations and Applications - 20th Brazilian Symposium, SBMF 2017, Recife, Brazil, November 29 - December 1, 2017, Proceedings*, volume 10623 of *Lecture Notes in Computer Science*, pages 7–11. Springer, 2017.
- [7] C. Benzmüller, C. Brown, and M. Kohlhasse. Higher-order semantics and extensionality. *Journal of Symbolic Logic*, 69(4):1027–1088, 2004.
- [8] C. Benzmüller, A. Farjami, and X. Parent. Faithful semantical embedding of a dyadic deontic logic in HOL. CoRR, <https://arxiv.org/abs/1802.08454>, 2018.
- [9] C. Benzmüller and D. Miller. Automation of higher-order logic. In D.M. Gabbay, J.H. Siekmann, and J. Woods, editors, *Handbook of the History of Logic, Volume 9 — Computational Logic*, pages 215–254. North Holland, Elsevier, 2014.

- [10] C. Benzmüller and L.C. Paulson. Quantified multimodal logics in simple type theory. *Logica Universalis (Special Issue on Multimodal Logics)*, 7(1):7–20, 2013.
- [11] J.C. Blanchette and T. Nipkow. Nitpick: A counterexample generator for higher-order logic based on a relational model finder. In *ITP 2010*, number 6172 in Lecture Notes in Computer Science, pages 131–146. Springer, 2010.
- [12] J. Carmo and A.J.I. Jones. Deontic logic and contrary-to-duties. In D. M. Gabbay and F. Guenther, editors, *Handbook of Philosophical Logic: Volume 8*, pages 265–343. Springer Netherlands, Dordrecht, 2002.
- [13] J. Carmo and A.J.I. Jones. Completeness and decidability results for a logic of contrary-to-duty conditionals. *J. Log. Comput.*, 23(3):585–626, 2013.
- [14] R.M. Chisholm. Contrary-to-duty imperatives and deontic logic. *Analysis*, 24:33–36, 1963.
- [15] A. Church. A set of postulates for the foundation of logic. *Annals of Mathematics*, 33(3):346–366, 1932.
- [16] A. Church. An unsolvable problem of elementary number theory. *American Journal of Mathematics*, 58(2):354–363, 1936.
- [17] A. Church. A formulation of the simple theory of types. *Journal of Symbolic Logic*, 5(2):56–68, 1940.
- [18] L. Henkin. Completeness in the theory of types. *Journal of Symbolic Logic*, 15(2):81–91, 1950.
- [19] P. McNamara. Deontic logic. In E.N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, winter 2014 edition, 2014.
- [20] R. Muskens. Intensional models for the theory of types. *Journal of Symbolic Logic*, 75(1):98–118, 2007.
- [21] T. Nipkow, L.C. Paulson, and M. Wenzel. *Isabelle/HOL — A Proof Assistant for Higher-Order Logic*, volume 2283 of *Lecture Notes in Computer Science*. Springer, 2002.

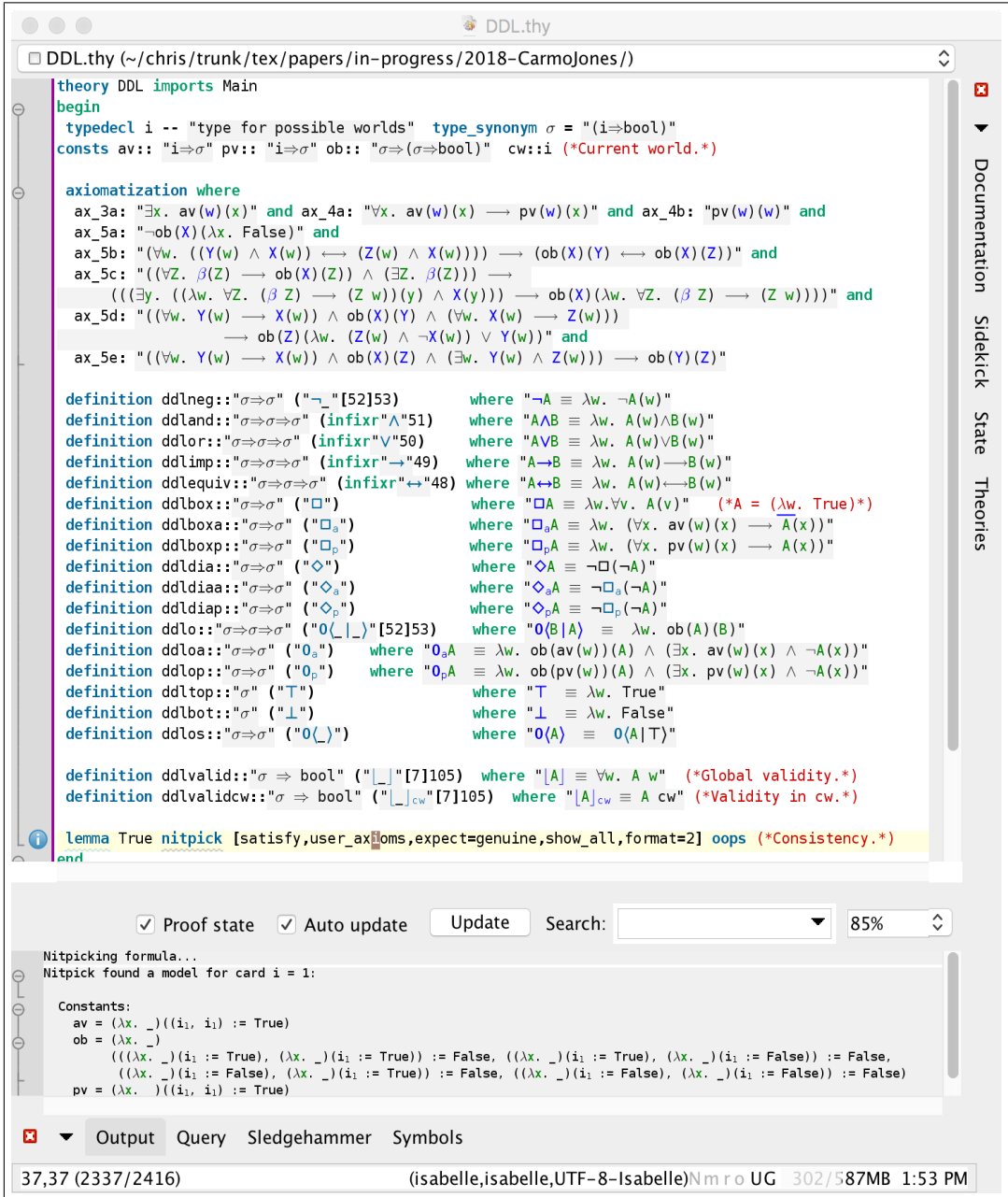


Figure 1: Shallow semantical embedding of DDL in Isabelle/HOL.

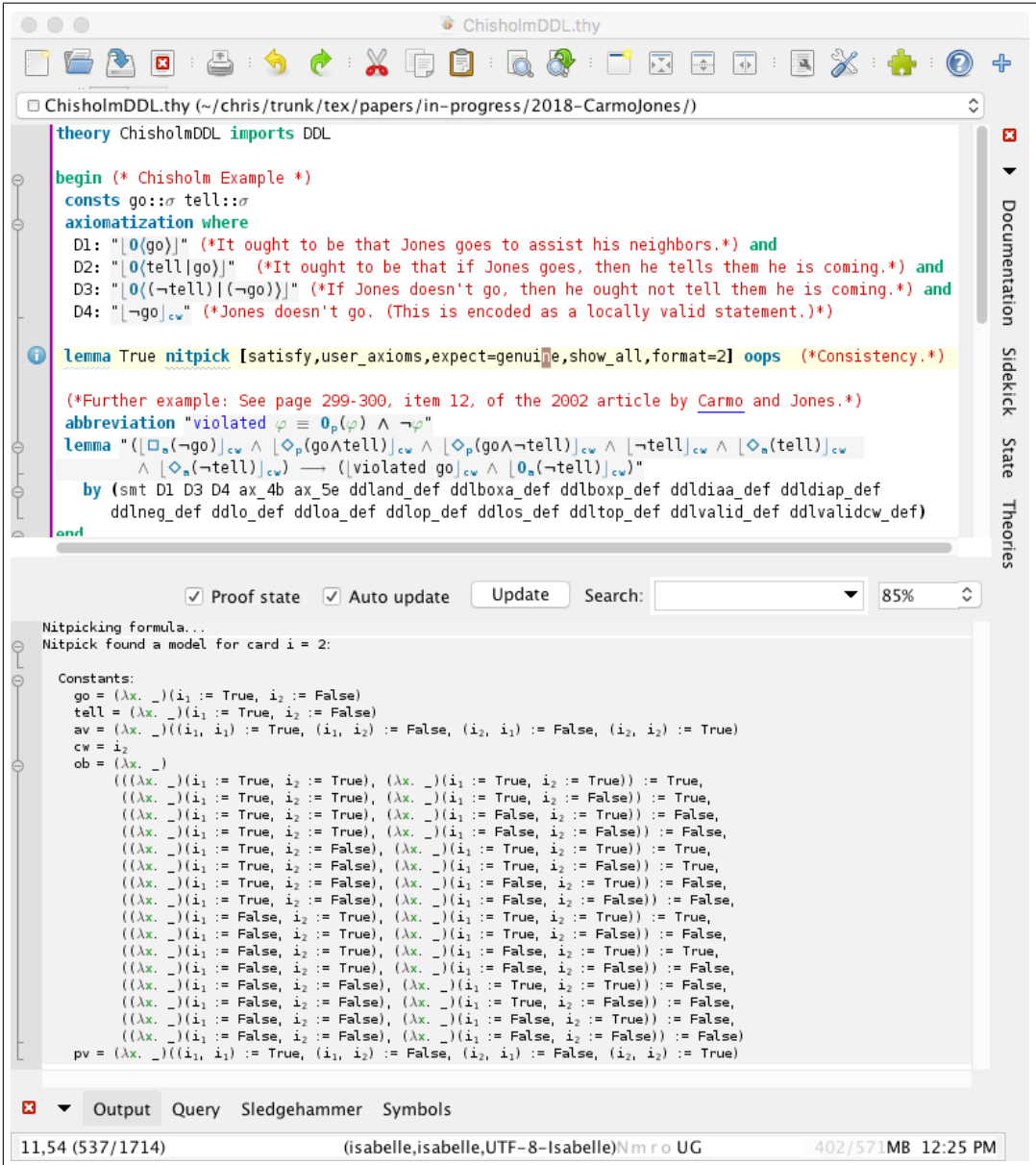


Figure 2: The Chisholm paradox scenario encoded in DDL (the shallow semantical embedding of DDL in Isabelle/HOL as displayed in Fig. 1 is imported here). Nitpick confirms consistency the encoded statements.