

---

# ÅQVIST’S DYADIC DEONTIC LOGIC **E** IN HOL

CHRISTOPH BENZMÜLLER

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

ALI FARJAMI

*University of Luxembourg, Luxembourg*  
farjami110@gmail.com

XAVIER PARENT

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

---

## Abstract

We devise a shallow semantical embedding of Åqvist’s dyadic deontic logic **E** in classical higher-order logic. This embedding is encoded in Isabelle/HOL, which turns this system into a proof assistant for deontic logic reasoning. The experiments with this environment provide evidence that this logic *implementation* fruitfully enables interactive and automated reasoning at the meta-level and the object-level.

*Keywords:* Dyadic deontic logic **E**; Preference based semantics; Classical higher-order logic; Semantic embedding; Automated reasoning.

## 1 Introduction

Normative notions such as obligation and permission are the subject of deontic logic [18] and conditional obligations are addressed in so-called *dyadic deontic logic*. A landmark and historically important family of dyadic deontic logics has been proposed by B. Hansson [20]. These logics have been recast in the framework of possible

---

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 - MIREL - MIning and REasoning with Legal texts. Benzmüller has been funded by the Volkswagen Foundation under project CRAP — Consistent Rational Argumentation in Politics.

world semantics by Åqvist [3]. They come with a preference semantics, in which a binary preference relation ranks the possible worlds in terms of betterness. The framework was motivated by the well-known paradoxes of *contrary-to-duty* (CTD) reasoning like Chisholm [14]’s paradox. In this paper we focus on the class of all preference models, in which no specific properties (like reflexivity or transitivity) are required of the betterness relation. This class of models has a known axiomatic characterisation, given by Åqvist’s system **E** [24].

When applied as a meta-logical tool, *simple type theory* [15], aka classical Higher-Order Logic (HOL), can help to better understand semantical issues of embedded object logics. The syntax and semantics of HOL are well understood [7] and there exist automated proof tools for it; examples include Isabelle/HOL [22] and LEO-II [11].

In this paper we devise an *embedding* of **E** in HOL. This embedding utilizes the *shallow semantical embedding* approach that has been put forward by Benzmüller[6] as a pragmatical solution towards universal logic reasoning. This approach uses classical higher-order logic as (universal) meta-logic to specify, in a shallow way, the syntax and semantics of various object logics, in our case system **E**. The embedding has been encoded in Isabelle/HOL to enable syntactical and semantical experiments in deontic reasoning.

Benzmüller et al. [9] developed an analogous shallow semantical embedding for the dyadic deontic logic proposed by Carmo and Jones [13]. A core difference concerns the notion of semantics employed in both papers, which leads to different semantical embeddings. Instead of the semantics based on preference models as employed by Hansson [20] and Åqvist [3], a neighborhood semantics is employed by Carmo and Jones [13].

Deep semantical embeddings of non-classical logics have been studied in the related literature [17, 16]. The emphasis in these works typically is on interactive proofs of meta-logical properties. While meta-logical studies [8, 19] are also in reach for the methods presented here, our interest is in proof automation at object level, i.e., proof automation of Åqvist’s system **E**. In other words, we are interested in practical normative reasoning applications of system **E** in which a high degree of automation at object level is required. Moreover, we are interested not only in the “propositional” system **E**, but also in quantified extensions of it. For this, we plan to accordingly adapt the achievements of previous work [10, 4]. Making deep semantical embeddings scale for quantified non-classical logics, on the contrary, seems more challenging and less promising regarding proof automation.

The article is structured as follows. Sec. 2 describes system **E** and Sec. 3 introduces HOL. The semantical embedding of **E** in HOL is then devised and studied in Sec. 4. This section also shows the faithfulness (viz. soundness and completeness)

of the embedding. Sec. 5 discusses the implementation in Isabelle/HOL [22]. Sec. 6 concludes the paper.

## 2 Dyadic Deontic Logic **E**

The language of **E** is obtained by adding the following operators to the syntax of propositional logic:  $\Box$  (for necessity);  $\Diamond$  (for possibility); and  $\bigcirc(-/-)$  (for conditional obligation).  $\bigcirc(\psi/\phi)$  is read “If  $\phi$ , then  $\psi$  is obligatory”. The set of well-formed formulas is defined in the straightforward way. Iteration of the modal and deontic operators is permitted, and so are “mixed” formulas, e.g.,  $\bigcirc(q/p) \wedge p$ . We put  $\top =_{df} \neg q \vee q$ , for some propositional symbol  $q$ , and  $\perp =_{df} \neg \top$ . A preference model is a structure  $M = \langle W, \succeq, V \rangle$  where

- $W$  is a non-empty set of items called possible worlds;
- $\succeq \subseteq W \times W$  (intuitively,  $\succeq$  is a betterness or comparative goodness relation; “ $s \succeq t$ ” can be read as “world  $s$  is at least as good as world  $t$ ”);
- $V$  is a function assigning to each atomic sentence a set of worlds. (i.e  $V(q) \subseteq W$ ).

No specific properties (like reflexivity or transitivity) are required of the betterness relation.

Given a preference model  $M = \langle W, \succeq, V \rangle$  and a world  $s \in W$ , we define the satisfaction relation  $M, s \models \varphi$  (read as “world  $s$  satisfies  $\varphi$  in model  $M$ ”) by induction on the structure of  $\varphi$  as described below. Standard Deontic Logic (SDL) [18] is based on two classes of states: good/bad (or green/red). Preference models allow gradations between good and bad. The closer a world is to ideality, the better it is. Intuitively, the evaluation rule for the dyadic obligation operator puts  $\bigcirc(\psi/\phi)$  true whenever all the best  $\phi$ -worlds are  $\psi$ -worlds. Here best is defined in terms of optimality rather than maximality [24]. A  $\phi$ -world is optimal, if it is as least as good as any other  $\phi$ -world. We define  $V^M(\varphi) = \{s \in W \mid M, s \models \varphi\}$  and  $\text{opt}_{\succeq}(V(\varphi)) = \{s \in V(\varphi) \mid \forall t(t \models \varphi \rightarrow s \succeq t)\}$ . Whenever the model  $M$  is obvious from context, we write  $V(\varphi)$  instead of  $V^M(\varphi)$ .

$$\begin{aligned}
 M, s &\models p \text{ if and only if } s \in V(p) \\
 M, s &\models \neg\varphi \text{ if and only if } M, s \not\models \varphi \text{ (that is, not } M, s \models \varphi) \\
 M, s &\models \varphi \vee \psi \text{ if and only if } M, s \models \varphi \text{ or } M, s \models \psi \\
 M, s &\models \Box\varphi \text{ if and only if } V(\varphi) = W
 \end{aligned}$$

$$M, s \models \bigcirc(\psi/\varphi) \text{ if and only if } \text{opt}_{\succeq}(V(\varphi)) \subseteq V(\psi)$$

As usual, a formula  $\varphi$  is valid in a preference model  $M = \langle W, \succeq, V \rangle$  (notation:  $M \models \varphi$ ) if and only if, for all worlds  $s \in W$ ,  $M, s \models \varphi$ . A formula  $\varphi$  is valid (notation:  $\models \varphi$ ) if and only if it is valid in every preference model. The notions of semantic consequence and satisfiability in a model are defined as usual.

System **E** is defined by the following axioms and rules:

All truth functional tautologies	(PL)
S5-schemata for $\Box$ and $\Diamond$	(S5)
$\bigcirc(\psi_1 \rightarrow \psi_2/\varphi) \rightarrow (\bigcirc(\psi_1/\varphi) \rightarrow \bigcirc(\psi_2/\varphi))$	(COK)
$\bigcirc(\psi/\varphi) \rightarrow \Box \bigcirc(\psi/\varphi)$	(Abs)
$\Box\psi \rightarrow \bigcirc(\psi/\varphi)$	(Nec)
$\Box(\varphi_1 \leftrightarrow \varphi_2) \rightarrow (\bigcirc(\psi/\varphi_1) \leftrightarrow \bigcirc(\psi/\varphi_2))$	(Ext)
$\bigcirc(\varphi/\varphi)$	(Id)
$\bigcirc(\psi/\varphi_1 \wedge \varphi_2) \rightarrow \bigcirc(\varphi_2 \rightarrow \psi/\varphi_1)$	(Sh)
If $\vdash \varphi$ and $\vdash \varphi \rightarrow \psi$ then $\vdash \psi$	(MP)
If $\vdash \varphi$ then $\vdash \Box\varphi$	(N)

The notions of theoremhood, deducibility and consistency are defined as usual.

The following theorem tells us that system **E** is the weakest system that characterises preference models. It also tells us that the assumptions of reflexivity and totalness of  $\succeq$  do not modify the logic, in the sense that they do not add new validities (or theorems).

**Theorem 1.** *System **E** is sound and complete with respect to the class of all preference models. System **E** is also sound and complete with respect to the class of those in which  $\succeq$  is reflexive, and with respect to the class of those in which  $\succeq$  is total (for all  $s, t \in W$ ,  $s \succeq t$  or  $t \succeq s$ ).*

*Proof.* See Parent [24]. □

**E** comes first in a family of three systems. Consider the condition of limitedness, whose role is to rule out infinite chains of strictly better worlds. Formally: if  $V(\phi) \neq \emptyset$ , then  $\text{opt}_{\succeq}(V(\phi)) \neq \emptyset$ . Such a condition boosts the logic to system **F**, obtained by supplementing **E** with D\*:

$$\Diamond\phi \rightarrow (\bigcirc(\psi/\phi) \rightarrow P(\psi/\phi)) \tag{D*}$$

Similarly, the additional assumption of transitivity of  $\succeq$  boosts the logic to system **G**, obtained by supplementing **F** with Sp:

$$(P(\psi/\phi) \wedge \bigcirc((\psi \rightarrow \chi)/\phi) \rightarrow \bigcirc(\chi/(\phi \wedge \psi))) \quad (\text{Sp})$$

None of **F** and **G** will concern us in this paper.

### 3 Classical Higher-order Logic

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

#### 3.1 Syntax of HOL

To define 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  $(V_\alpha)_{\alpha \in T}$ .<sup>1</sup> We employ Church-style typing, where each term  $t_\alpha$  explicitly encodes its type information in subscript  $\alpha$ .

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  $\alpha$ .
- Every typed variable symbol  $X_\alpha \in V_\alpha$  is a HOL term of type  $\alpha$ .
- If  $s_{\alpha \rightarrow \beta}$  and  $t_\alpha$  are HOL terms of types  $\alpha \rightarrow \beta$  and  $\alpha$ , respectively, then  $(s_{\alpha \rightarrow \beta} t_\alpha)_\beta$ , called *application*, is an HOL term of type  $\beta$ .
- If  $X_\alpha \in V_\alpha$  is a typed variable symbol and  $s_\beta$  is an HOL term of type  $\beta$ , 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  $\lambda$ -calculus. In order to extend this base framework into logic HOL we simply ensure that the signature

---

<sup>1</sup>For example in Sec. 4 we assume constant symbol  $r$ , with type  $i \rightarrow i \rightarrow o$  as part of 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  $\alpha$ ) 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  $((\vee_{o \rightarrow o \rightarrow o} s_o) t_o)$ .

From the selected set of primitive connectives, other logical connectives can be introduced as abbreviations.<sup>2</sup> 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  $\exists X_\alpha s := \neg \forall X_\alpha \neg s$ .

Each occurrence of a variable in a term is either bound by a  $\lambda$  or free. We use  $free(s)$  to denote the set of variables with a free occurrence in  $s$ . We consider two terms to be *equal* if the terms are the same up to the names of bound variables, that is, we consider  $\alpha$ -conversion implicitly.

*Substitution* of a term  $s_\alpha$  for a variable  $X_\alpha$  in a term  $t_\beta$  is denoted by  $[s/X]t$ . Since we consider  $\alpha$ -conversion implicitly, we assume the bound variables of  $t$  to avoid variable capture.

Well-known operations and relations on HOL terms include  *$\beta\eta$ -normalization* and  *$\beta\eta$ -equality*, denoted by  $s =_{\beta\eta} t$ ,  *$\beta$ -reduction* and  *$\eta$ -reduction*. A  *$\beta$ -redex*  $(\lambda X s)t$   $\beta$ -reduces to  $[t/X]s$ . An  *$\eta$ -redex*  $\lambda X(sX)$ , where  $X \notin free(s)$ ,  $\eta$ -reduces to  $s$ . We write  $s =_\beta t$  to mean  $s$  can be converted to  $t$  by a series of  $\beta$ -reductions and expansions. Similarly,  $s =_{\beta\eta} t$  means  $s$  can be converted to  $t$  using both  $\beta$  and  $\eta$ .

---

<sup>2</sup>As demonstrated by Andrews [2], we could in fact start out with only primitive equality in the signature (for all types  $\alpha$ ) and introduce all other logical connectives as abbreviations based on it.

### 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 literature [7].

The semantics of choice for the remainder is Henkin semantics, i.e., we work with Henkin's general models [21]. 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  $D_\alpha$ , such that  $D_o = \{T, F\}$  (for truth and falsehood). The  $D_{\alpha \rightarrow \beta}$  are collections of functions mapping  $D_\alpha$  into  $D_\beta$ .

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  $D_\alpha$ , called the *denotation* of  $p_\alpha$ . The logical connectives  $\neg$ ,  $\vee$ ,  $\Pi$  and  $=$  are always given their expected, standard denotations:<sup>3</sup>

- $I(\neg_{o \rightarrow o}) = \text{not} \in D_{o \rightarrow o}$  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  $b$ .
- $I(\Pi_{(\alpha \rightarrow o) \rightarrow o}) = \text{all} \in D_{(\alpha \rightarrow o) \rightarrow o}$  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  $\alpha$ .

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  $X_\alpha$  to elements in  $D_\alpha$ .  $g[d/W]$  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  $s_\alpha$  on a model  $M = \langle D, I \rangle$  under assignment  $g$  is an element  $d \in D_\alpha$  defined in the following way:

---

<sup>3</sup>Since  $=_{\alpha \rightarrow \alpha \rightarrow o}$  (for all types  $\alpha$ ) 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 [1]: if such identity relations did not exist in the  $D_{\alpha \rightarrow \alpha \rightarrow o}$ , then Leibniz equality in Henkin semantics might not denote as intended.

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

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 a Henkin model  $M$  under assignment  $g$  if and only if  $\|s_o\|^{M,g} = T$ ; this is also expressed by writing that  $M, g \models^{HOL} s_o$ . An HOL formula  $s_o$  is called *valid* in  $M$ , which is expressed by writing that  $M \models^{HOL} s_o$ , if and only if  $M, g \models^{HOL} s_o$  for all assignments  $g$ . Moreover, a formula  $s_o$  is called *valid*, expressed by writing that  $\models^{HOL} s_o$ , if and only if  $s_o$  is valid in all Henkin models  $M$ .



## 4 Embedding **E** into HOL

### 4.1 Semantical Embedding

The formulas of **E** 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  $\tau$  in the remainder. The HOL signature is assumed to contain the constant symbol  $r_{i \rightarrow \tau}$ . Moreover, for each propositional symbol  $p^j$  of **E**, the HOL signature must contain the corresponding constant symbol  $p_\tau^j$ . 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 a formula  $\varphi$  of **E** into a formula  $[\varphi]$  of HOL of type  $\tau$ . The mapping is defined recursively:

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

$\neg_{\tau \rightarrow \tau}$ ,  $\vee_{\tau \rightarrow \tau \rightarrow \tau}$ ,  $\Box_{\tau \rightarrow \tau}$ ,  $\bigcirc_{\tau \rightarrow \tau \rightarrow \tau}$  abbreviate the following formulas of HOL:

$$\begin{aligned} \neg_{\tau \rightarrow \tau} &= \lambda A_\tau \lambda X_i \neg(A X) \\ \vee_{\tau \rightarrow \tau \rightarrow \tau} &= \lambda A_\tau \lambda B_\tau \lambda X_i (A X \vee B X) \\ \Box_{\tau \rightarrow \tau} &= \lambda A_\tau \lambda X_i \forall Y_i (A Y) \\ \bigcirc_{\tau \rightarrow \tau \rightarrow \tau} &= \lambda A_\tau \lambda B_\tau \lambda X_i \forall W_i ((\lambda V_i (A V \wedge (\forall Y_i (A Y \rightarrow r_{i \rightarrow \tau} V Y)))) W \rightarrow B W)^4 \end{aligned}$$

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  $vld_{\tau \rightarrow o} = \lambda A_\tau \forall S_i (A S)$ . With this definition, validity of a formula  $s$  in **E** corresponds to the validity of the formula  $(vld [s])$  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 preference models into Henkin models is employed.

<sup>4</sup>If  $\text{opt}_{\succeq}(A)$  is taken as an abbreviation for  $\lambda V_i (A V \wedge (\forall Y_i (A Y \rightarrow r_{i \rightarrow \tau} V Y)))$ , then this can be simplified to  $\bigcirc_{\tau \rightarrow \tau \rightarrow \tau} = \lambda A_\tau \lambda B_\tau \lambda X_i (\text{opt}_{\succeq}(A) \subseteq B)$ .

**Definition 1** (Preference model  $\Rightarrow$  Henkin model). *Let  $M = \langle W, \succeq, V \rangle$  be a preference model. Let  $p^1, \dots, p^m \in PV$ , for  $m \geq 1$  be propositional symbols and  $\lfloor p^j \rfloor = p^j_\tau$  for  $j = 1, \dots, m$ . A Henkin model  $H^M = \langle \{D_\alpha\}_{\alpha \in T}, I \rangle$  for  $M$  is defined as follows:  $D_i$  is chosen as the set of possible worlds  $W$  and all other sets  $D_{\alpha \rightarrow \beta}$  are chosen as (not necessarily full) sets of functions from  $D_\alpha$  to  $D_\beta$ . 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, in particular, it is required that  $D_\tau$  and  $D_{i \rightarrow \tau}$  contain the elements  $Ip^j_\tau$  and  $Ir_{i \rightarrow \tau}$ . Interpretation  $I$  is constructed as follows:*

1. For  $1 \leq i \leq m$ ,  $Ip^j_\tau \in D_\tau$  is chosen such that  $Ip^j_\tau(s) = T$  iff  $s \in V(p^j)$  in  $M$ .
2.  $Ir_{i \rightarrow \tau} \in D_{i \rightarrow \tau}$  is chosen such that  $Ir_{i \rightarrow \tau}(s, u) = T$  iff  $s \succeq u$  in  $M$ .

Since we assume that there are no other symbols (besides the  $r$ , the  $p^j$  and the primitive logical connectives) 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 preference model  $M$ . For all formulas  $\delta$  of  $\mathbf{E}$ , all assignments  $g$  and worlds  $s$  it holds:*

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

*Proof.* See appendix. □

**Lemma 2** (Henkin model  $\Rightarrow$  Preference model). *For every Henkin model  $H = \langle \{D_\alpha\}_{\alpha \in T}, I \rangle$  there exists a corresponding preference model  $M$ . Corresponding here means that for all formulas  $\delta$  of  $\mathbf{E}$  and for all assignments  $g$  and worlds  $s$ ,*

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

*Proof.* Suppose that  $H = \langle \{D_\alpha\}_{\alpha \in T}, I \rangle$  is a Henkin model. Without loss of generality, we can assume that the domains of  $H$  are denumerable [21]. We construct the corresponding preference model  $M$  as follows:

- $W = D_i$ .
- $s \succeq u$  for  $s, u \in W$  iff  $Ir_{i \rightarrow \tau}(s, u) = T$ .
- $s \in V(p^j_\tau)$  iff  $Ip^j_\tau(s) = T$  for all  $p^j$ .

Moreover, the above construction ensures that  $H$  is a Henkin model for  $M$ . Hence, Lemma 1 applies. This ensures that for all formulas  $\delta$  of **E**, for all assignments  $g$  and all worlds  $s$  we have  $\|\llbracket \delta \rrbracket S_i\|^{H,g[s/S_i]} = T$  if and only if  $M, s \models \delta$ .  $\square$

**Theorem 2** (Soundness and Completeness of the Embedding).

$$\models \varphi \text{ if and only if } \models^{HOL} vld \llbracket \varphi \rrbracket$$

*Proof.* (Soundness,  $\leftarrow$ ) The proof is by contraposition. Assume  $\not\models \varphi$ , i.e., there is a preference model  $M = \langle W, \succeq, V \rangle$ , and a world  $s \in W$ , such that  $M, s \not\models \varphi$ . By Lemma 1 for an arbitrary assignment  $g$  it holds that  $\|\llbracket \varphi \rrbracket S_i\|^{H^M,g[s/S_i]} = F$  in Henkin model  $H^M = \langle \{D_\alpha\}_{\alpha \in T}, I \rangle$ . Thus, by definition of  $\|\cdot\|$ , it holds that  $\|\forall S_i(\llbracket \varphi \rrbracket S_i)\|^{H^M,g} = \|vld \llbracket \varphi \rrbracket\|^{H^M,g} = F$ . Hence,  $H^M \not\models^{HOL} vld \llbracket \varphi \rrbracket$ . By definition  $\not\models^{HOL} vld \llbracket \varphi \rrbracket$ .

(Completeness,  $\rightarrow$ ) The proof is again by contraposition. Assume  $\not\models^{HOL} vld \llbracket \varphi \rrbracket$ , i.e., there is a Henkin model  $H = \langle \{D_\alpha\}_{\alpha \in T}, I \rangle$  and an assignment  $g$  such that  $\|vld \llbracket \varphi \rrbracket\|^{H,g} = F$ . By Lemma 2, there is a preference model  $M$  such that  $M \not\models \varphi$ . Hence,  $\not\models \varphi$ .  $\square$

*Remark:* In contrast to a deep logical embedding, in which the syntactical structure and the semantics of logic  $L$  would be formalized in full detail (using e.g., structural induction and recursion), only the core differences in the semantics of both system **E** and meta-logic HOL have been explicitly encoded in our shallow semantical embedding. In a certain sense we have thus shown, that system **E** can in fact be identified and handled as a natural fragment of HOL.

## 5 Implementation in Isabelle/HOL

The semantical embedding as devised in Sec. 4 has been implemented in the higher-order proof assistant Isabelle/HOL [22]. 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 a certain man go to help his neighbors;
2. It ought to be that if he goes he tells them he is coming;
3. If he does not go, he ought not to tell them he is coming;
4. He does not go.

```

1 theory DDLE imports Main
2 begin
3 typedecl i - <type for possible worlds>
4 type_synonym  $\tau$  = "(i $\Rightarrow$ bool)"
5 consts aw::i (* actual world *)
6
7 abbreviation(input) mtrue  :: " $\tau$ " ("T") where "T  $\equiv$   $\lambda w$ . True"
8 abbreviation(input) mfalse :: " $\tau$ " ("⊥") where "⊥  $\equiv$   $\lambda w$ . False"
9 abbreviation(input) mnnot  :: " $\tau \Rightarrow \tau$ " ("¬") [52]53 where "¬  $\equiv$   $\lambda w$ .  $\neg \psi(w)$ "
10 abbreviation(input) mand   :: " $\tau \Rightarrow \tau \Rightarrow \tau$ " ("infixr" ^) 51 where " $\varphi \wedge \psi \equiv \lambda w$ .  $\varphi(w) \wedge \psi(w)$ "
11 abbreviation(input) mor    :: " $\tau \Rightarrow \tau \Rightarrow \tau$ " ("infixr" v) 50 where " $\varphi \vee \psi \equiv \lambda w$ .  $\varphi(w) \vee \psi(w)$ "
12 abbreviation(input) mimp   :: " $\tau \Rightarrow \tau \Rightarrow \tau$ " ("infixr"  $\rightarrow$ ) 49 where " $\varphi \rightarrow \psi \equiv \lambda w$ .  $\varphi(w) \rightarrow \psi(w)$ "
13 abbreviation(input) mequ   :: " $\tau \Rightarrow \tau \Rightarrow \tau$ " ("infixr"  $\leftrightarrow$ ) 48 where " $\varphi \leftrightarrow \psi \equiv \lambda w$ .  $\varphi(w) \leftrightarrow \psi(w)$ "
14
15 abbreviation(input) mbox   :: " $\tau \Rightarrow \tau$ " ("□") where "□  $\equiv$   $\lambda \varphi w$ .  $\forall v$ .  $\varphi(v)$ "
16 consts r :: "i $\Rightarrow$  $\tau$ " (infixr "r" 70)
17 - <the betterness relation r, used in definition of  $\bigcirc$  <_> >
18 abbreviation(input) mopt   :: " $\tau \Rightarrow \tau$ " ("opt<_>")
19 where "opt< $\varphi$ >  $\equiv$  ( $\lambda v$ . (  $\varphi(v) \wedge (\forall x$ . ( $\varphi(x) \rightarrow v r x$  ) ) )"
20 abbreviation(input) msubset :: " $\tau \Rightarrow \tau \Rightarrow$ bool" (infix "⊆" 53)
21 where " $\varphi \subseteq \psi \equiv \forall x$ .  $\varphi x \rightarrow \psi x$ "
22 abbreviation(input) mcond  :: " $\tau \Rightarrow \tau \Rightarrow \tau$ " (" $\bigcirc$ <_>")
23 where " $\bigcirc \langle \psi | \varphi \rangle \equiv \lambda w$ . opt< $\varphi$ >  $\subseteq \psi$ "
24
25 abbreviation(input) valid  :: " $\tau \Rightarrow$ bool" ("|_") [8]109
26 where "|p|  $\equiv \forall w$ . p w"
27 definition cactual :: " $\tau \Rightarrow$ bool" ("|_|") [7]105
28 where "|p|_|  $\equiv p(aw)$ "
29
30 lemma True nitpick [satisfy, user_axioms, show_all, expect=genuine] oops
    
```

 Figure 1: Shallow semantical embedding of  $\mathbf{E}$  in Isabelle/HOL.

These statements can be given a consistent formalisation in System  $\mathbf{E}$ ; cf. Fig. 2. This is confirmed by the model finder Nitpick [12] integrated with Isabelle/HOL. Nitpick computes an intuitive, small model for the scenario consisting of one possible world  $i_1$ . The actual world is  $i_1$  also. Preference relation  $r$  is interpreted in this model as  $r = \emptyset$ . In the actual world the man doesn't go to help his neighbors and doesn't tell them that he is coming. That is,  $V(\neg go) = V(\neg tell) = \{i_1\}$ . Also, we have  $op(V(\top)) = \emptyset$ . So,  $i_1 \models \bigcirc(go/\top)$  by the evaluation rule for  $\bigcirc$ . Similarly,  $op(V(go)) = op(V(\neg go)) = \emptyset$  implies  $i_1 \models \bigcirc(tell/go)$  and  $i_1 \models \bigcirc(\neg tell/\neg go)$ .

```

31
32 section {* Chisholm Scenario *}
33
34 consts go :: "τ" tell :: "τ"
35
36 context (*Chisholm Scenario*)
37 assumes
38 ax1: "[ O<go|T> ]" (*It ought to be that a certain man go to help his neighbours.*) and
39
40 ax2: "[ O<tell|go > ]"(*It ought to be that if he goes he tells them he is coming.*) and
41
42 ax3: "[ O<-tell|-go> ]" (*If he does not go, he ought not to tell them he is coming.*) and
43
44 ax4 : "[¬go]₁" (*He does not go.*)
45
46 begin
47
48 lemma True nitpick [satisfy, user_axioms, show_all,expect=genuine] oops
49
50 end

```

Proof state  Auto update  Search:  100%

Nitpick found a model for card i = 1:

```

Constants:
aw = i₁
go = (λx. _)(i₁ := False)
(r) = (λx. _)(i₁ := (λx. _)(i₁ := False))
tell = (λx. _)(i₁ := False)

```

Output  Query  Sledgehammer  Symbols

M8,17 (1856/3855) (isabelle.isabelle.UTF-8-Isabelle) | n m r o U G | 30/695MB 2:58 PM

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

## 6 Conclusion

A shallow semantical embedding of Åqvist’s dyadic deontic logic **E** in classical higher-order logic has been presented and shown to be faithful (sound and complete). The works presented here and in Benz Müller et al. [9] provide the theoretical foundation for the implementation and automation of dyadic deontic logic within existing theorem provers and proof assistants for HOL. We do not provide new logics. Instead, we provide an empirical infrastructure for assessing practical aspects of ambitious, state-of-the-art deontic logics; this has not been done before.

We end this paper by listing a number of topics for future research. First, it would be worthwhile to study the shallow semantical embedding of the stronger systems **F** and **G** in HOL, and look at the three systems from the point of view of a semantics defining best in terms of maximality rather than optimality [23, 24]. Second, we could employ our implementation to systematically study some meta-logical properties of these systems within Isabelle/HOL. Third, it would be interesting to study the quantified extensions of system **E**, **F** and **G**. Previous work has focused on monadic modal logic and conditional logic [4, 5, 10]. Last, but not least, experiments could investigate whether the provided implementation already supports non-trivial applications in practical normative reasoning, or whether further improvements are required.

## Acknowledgements

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

## References

- [1] Andrews, P.B.: General models and extensionality. *Journal of Symbolic Logic* **37**(2), 395–397 (1972)
- [2] Andrews, P.B.: Church’s type theory. In: Zalta, E.N. editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, spring 2014 edition, (2014)
- [3] Åqvist, L.: Deontic logic. In: *Handbook of philosophical logic*, pp. 147–264. Springer, Dordrecht (2002)
- [4] Benz Müller, C.: Automating quantified conditional logics in HOL. In: Rossi, F. (eds.) *23rd International Joint Conference on Artificial Intelligence, IJCAI-13*, Beijing, China, pp. 746–753, AAAI Press, (2013)
- [5] Benz Müller, C.: Cut-elimination for quantified conditional logic. *Journal of Philosophical Logic*, **46**(3), 333–353, (2017)
- [6] Benz Müller, C.: Universal (Meta-)Logical Reasoning: Recent Successes. *Science of Computer Programming* (in print, preprint: <http://dx.doi.org/10.13140/RG.2.2.11039.61609/2>), (2018)
- [7] Benz Müller, C., Brown, C., Kohlhase, M.: Higher-order semantics and extensionality. *Journal of Symbolic Logic*, **69**(4), 1027–1088, (2004)
- [8] Benz Müller, C., Claus, M., Sultana, N.: Systematic verification of the modal logic cube in Isabelle/HOL. In: Kaliszzyk, C., Paskevich, A. (eds.) *PxTP 2015*, Berlin, Germany, EPTCS, vol. 186, pp. 24–41 (2015).
- [9] Benz Müller, C., Farjami, A., Parent., X.: A dyadic deontic logic in HOL. In: Broersen, J., Condoravdi, C., Nair, S., Pigozzi, G. (eds.) *Deontic Logic and Normative Systems* —

- 14th International Conference, DEON 2018, Utrecht, The Netherlands, 3-6 July, 2018, pp. 33–50, College Publications, UK, (2018)
- [10] Benzmüller, C., Paulson, L.C.: Quantified multimodal logics in simple type theory. *Logica Universalis (Special Issue on Multimodal Logics)*, **7**(1), 7–20, (2013)
- [11] Benzmüller, C., Sultana, N., Paulson, L. C., Theiß, F.: The higher-order prover LEO-II. *Journal of Automated Reasoning*, **55**(4), 389–404, (2015)
- [12] Blanchette, J.C., Nipkow, T.: Nitpick: A counterexample generator for higher-order logic based on a relational model finder. In: Kaufmann, M., Paulson, L. C. (eds.) *International Conference on Interactive Theorem Proving 2010*, LNCS, vol.6172 pp. 131–146, Springer, (2010)
- [13] Carmo, J. M. C. L. M., Jones, A. J. I.: Completeness and decidability results for a logic of contrary-to-duty conditionals. *Journal of Logic and Computation* **23**(3), 585–626 (2013)
- [14] Chisholm, R. M.: Contrary-to-duty imperatives and deontic logic. *Analysis*, **24**(2), 33–36 (1963)
- [15] Church, A.: A formulation of the simple theory of types. *Journal of Symbolic Logic*, **5**(2), 56–68, (1940)
- [16] Doczkal, C., Bard, J.: Completeness and decidability of converse PDL in the constructive type theory of Coq. In: Andronick, J., Felty, A. P. (eds.) *International Conference on Certified Programs and Proofs, CPP 2018, Los Angeles, USA, Proceedings of the 7th ACM SIGPLAN*, pp. 42–52, ACM, New York, USA (2018)
- [17] Doczkal, C., Smolka, G.: Completeness and decidability results for CTL in constructive type theory. *Journal of Automated Reasoning* **56**(32), 343–365 (2016)
- [18] Gabbay, D., Horty, J., Parent, X., van der Meyden, R., van der Torre, L.: *Handbook of deontic logic and normative systems. Volume 1*. College Publications, UK, (2013)
- [19] Kirchner, D., Benzmüller, C., Zalta, E.: Mechanizing principia logico-metaphysica in functional type theory. CoRR <https://arxiv.org/abs/1711.06542>, (2017)
- [20] Hansson, B.: An analysis of some deontic logics. *Nous*, 373–398 (1969)
- [21] Henkin, L.: Completeness in the theory of types. *Journal of Symbolic Logic*, **5**(2), 81–91, (1950)
- [22] Nipkow, T., Paulson, L.C., Wenzel, M.: Isabelle/HOL — A proof assistant for higher-Order logic, volume 2283 of *Lecture Notes in Computer Science*. Springer, (2002)
- [23] Parent, X.: Maximality vs optimality in dyadic deontic logic - Completeness results for systems in Hansson’s tradition. *Journal of Philosophical Logic*, **43**(6), 1101–1128 (2014)
- [24] Parent, X.: Completeness of Åqvist’s systems E and F. *The Review of Symbolic Logic*, **8**(1), 164–177 (2015)

## Appendix

### Proof for Lemma 1

In the proof we implicitly employ curring and uncuring, and we associate sets with their characteristic functions. Throughout the proof whenever possible we omit types in order to avoid making the notation too cumbersome. The proof of lemma 1 is by induction on the structure of  $\delta$ . We start with the case where  $\delta$  is  $p^j$ . We have

$$\begin{aligned}
 & \| [p^j] S \|^{H^M, g[s/S_i]} = T \\
 \Leftrightarrow & \| p^j_\tau S \|^{H^M, g[s/S_i]} = T \\
 \Leftrightarrow & Ip^j_\tau(s) = T \\
 \Leftrightarrow & s \in V(p^j) \quad (\text{by definition of } H^M) \\
 \Leftrightarrow & M, s \models p^j
 \end{aligned}$$

In the inductive cases we make use of the following **induction hypothesis**: *For sentences  $\delta'$  structurally smaller than  $\delta$  we have: For all assignments  $g$  and states  $s$ ,  $\| [\delta'] S \|^{H^M, g[s/S_i]} = T$  if and only if  $M, s \models \delta'$ .*

We consider each inductive case in turn:

(a)  $\delta = \varphi \vee \psi$ . In this case:

$$\begin{aligned}
 & \| [\varphi \vee \psi] S \|^{H^M, g[s/S_i]} = T \\
 \Leftrightarrow & \| ([\varphi] \vee_{\tau \rightarrow \tau \rightarrow \tau} [\psi]) S \|^{H^M, g[s/S_i]} = T \\
 \Leftrightarrow & \| ([\varphi] S) \vee ([\psi] S) \|^{H^M, g[s/S_i]} = T \quad (([\varphi] \vee_{\tau \rightarrow \tau \rightarrow \tau} [\psi]) S =_{\beta\eta} ([\varphi] S) \vee ([\psi] S)) \\
 \Leftrightarrow & \| [\varphi] S \|^{H^M, g[s/S_i]} = T \text{ or } \| [\psi] S \|^{H^M, g[s/S_i]} = T \\
 \Leftrightarrow & M, s \models \varphi \text{ or } M, s \models \psi \quad (\text{by induction hypothesis}) \\
 \Leftrightarrow & M, s \models \varphi \vee \psi
 \end{aligned}$$

(b)  $\delta = \neg\varphi$ . In this case:

$$\begin{aligned}
 & \| [\neg\varphi] S \|^{H^M, g[s/S_i]} = T \\
 \Leftrightarrow & \| (\neg_{\tau \rightarrow \tau} [\varphi]) S \|^{H^M, g[s/S_i]} = T \\
 \Leftrightarrow & \| \neg([\varphi] S) \|^{H^M, g[s/S_i]} = T \quad ((\neg_{\tau \rightarrow \tau} [\varphi]) S =_{\beta\eta} \neg([\varphi] S)) \\
 \Leftrightarrow & \| [\varphi] S \|^{H^M, g[s/S_i]} = F \\
 \Leftrightarrow & M, s \not\models \varphi \quad (\text{by induction hypothesis}) \\
 \Leftrightarrow & M, s \models \neg\varphi
 \end{aligned}$$

(c)  $\delta = \Box\varphi$ . We have the following chain of equivalences:

$$\begin{aligned}
 & \| [\Box\varphi] S \|^{H^M, g[s/S_i]} = T \\
 \Leftrightarrow & \| (\lambda X \forall Y ([\varphi] Y)) S \|^{H^M, g[s/S_i]} = T
 \end{aligned}$$



- $\Leftrightarrow \|\forall Y(\lfloor \varphi \rfloor Y)\|^{H^M, g[s/S_i]} = T$   
 $\Leftrightarrow$  For all  $a \in D_i$  we have  $\|\lfloor \varphi \rfloor Y\|^{H^M, g[s/S_i][a/Y_i]} = T$   
 $\Leftrightarrow$  For all  $a \in D_i$  we have  $\|\lfloor \varphi \rfloor Y\|^{H^M, g[a/Y_i]} = T$  ( $S \notin \text{free}(\lfloor \varphi \rfloor) = \emptyset$ )  
 $\Leftrightarrow$  For all  $a \in D_i$  we have  $M, a \models \varphi$  (by induction hypothesis)  
 $\Leftrightarrow M, s \models \Box \varphi$

(d)  $\delta = \bigcirc(\psi/\varphi)$ . We have the following chain of equivalences:

- $\|\lfloor \bigcirc(\psi/\varphi) \rfloor S\|^{H^M, g[s/S_i]} = T$   
 $\Leftrightarrow \|\lfloor (\lambda X \forall W ((\lambda V (\lfloor \varphi \rfloor V \wedge (\forall Y (\lfloor \varphi \rfloor Y \rightarrow r V Y)))) W \rightarrow \lfloor \psi \rfloor W)) S\|^{H^M, g[s/S_i]} = T$   
 $\Leftrightarrow \|\forall W ((\lambda V (\lfloor \varphi \rfloor V \wedge (\forall Y (\lfloor \varphi \rfloor Y \rightarrow r V Y)))) W \rightarrow \lfloor \psi \rfloor W)\|^{H^M, g[s/S_i]} = T$   
 $\Leftrightarrow$  For all  $u \in D_i$  we have:  
 $\|\lfloor (\lambda V (\lfloor \varphi \rfloor V \wedge (\forall Y (\lfloor \varphi \rfloor Y \rightarrow r V Y)))) W \rightarrow \lfloor \psi \rfloor W\|^{H^M, g[s/S_i][u/W_i]} = T$   
 $\Leftrightarrow$  For all  $u \in D_i$  we have:  
 If  $\|\lfloor (\lambda V (\lfloor \varphi \rfloor V \wedge (\forall Y (\lfloor \varphi \rfloor Y \rightarrow r V Y)))) W\|^{H^M, g[s/S_i][u/W_i]} = T$ ,  
 then  $\|\lfloor \psi \rfloor W\|^{H^M, g[s/S_i][u/W_i]} = T$   
 $\Leftrightarrow$  For all  $u \in D_i$  we have:  
 If  $\|\lfloor \varphi \rfloor W\|^{H^M, g[s/S_i][u/W_i]} = T$  and  
 $\|\forall Y (\lfloor \varphi \rfloor Y \rightarrow r W Y)\|^{H^M, g[s/S_i][u/W_i]} = T$ ,  
 then  $\|\lfloor \psi \rfloor V\|^{H^M, g[s/S_i][u/W_i]} = T$   
 $\Leftrightarrow$  For all  $u \in D_i$  we have:  
 If  $\|\lfloor \varphi \rfloor W\|^{H^M, g[s/S_i][u/W_i]} = T$  and  
 for all  $t \in D_i$  we have  $\|\forall Y (\lfloor \varphi \rfloor Y \rightarrow r W Y)\|^{H^M, g[s/S_i][u/W_i][t/Y_i]} = T$ ,  
 then  $\|\lfloor \psi \rfloor W\|^{H^M, g[s/S_i][u/W_i]} = T$   
 $\Leftrightarrow$  For all  $u \in D_i$  we have:  
 If  $\|\lfloor \varphi \rfloor W\|^{H^M, g[s/S_i][u/W_i]} = T$  and  
 for all  $t \in D_i$  we have  $\|\lfloor \varphi \rfloor Y\|^{H^M, g[s/S_i][u/W_i][t/Y_i]} = T$  implies  $Ir_{i \rightarrow \tau}(u, t) = T$ ,  
 then  $\|\lfloor \psi \rfloor W\|^{H^M, g[s/S_i][u/W_i]} = T$   
 $\Leftrightarrow$  For all  $u \in D_i$  we have:  
 If  $u \in V(\varphi)$  and  
 for all  $t \in D_i$  we have  $t \in V(\varphi)$  implies  $u \succeq t$ ,  
 then  $u \in V(\psi)$  (**see the justification \***)  
 $\Leftrightarrow \text{opt}_{\succeq}(V(\varphi)) \subseteq V(\psi)$   
 $\Leftrightarrow M, s \models \bigcirc(\psi/\varphi)$

**Justification \***: What we need to show is:  $\|\lfloor \varphi \rfloor\|^{H^M, g[s/S_i]}$  is identified with  $V(\varphi)$  (analogously  $\psi$ ). By induction hypothesis, for all assignments  $g$  and states  $s$ , we have  $\|\lfloor \varphi \rfloor S\|^{H^M, g[s/S_i]} = T$  if and only if  $M, s \models \varphi$ . Expanding the details of this

equivalence we have: for all assignments  $g$  and states  $s$

$$\begin{aligned}
 & s \in \llbracket [\varphi] \rrbracket^{H^M, g[s/S_i]} \quad (\text{functions to type } o \text{ are associated with sets}) \\
 \Leftrightarrow & \llbracket [\varphi] \rrbracket^{H^M, g[s/S_i]}(s) = T \\
 \Leftrightarrow & \llbracket [\varphi] \rrbracket^{H^M, g[s/S_i]} \llbracket S \rrbracket^{H^M, g[s/S_i]} = T \\
 \Leftrightarrow & \llbracket [\varphi] S \rrbracket^{H^M, g[s/S_i]} = T \\
 \Leftrightarrow & M, s \models \varphi \\
 \Leftrightarrow & s \in V(\varphi)
 \end{aligned}$$

Hence,  $s \in \llbracket [\varphi] \rrbracket^{H^M, g[s/S_i]}$  if and only if  $s \in V(\varphi)$ .

By extensionality we thus know that  $\llbracket [\varphi] \rrbracket^{H^M, g[s/S_i]}$  is identified with  $V(\varphi)$ . Moreover, since  $H^M$  obeys the Denotatpflicht we know that  $V(\varphi) \in D_\tau$ .