Secondary Publication



Parent, Xavier; Benzmüller, Christoph

Normative Conditional Reasoning as a Fragment of HOL

Date of secondary publication: 26.09.2023 Submitted Version (Preprint), Article Persistent identifier: urn:nbn:de:bvb:473-irb-910431

Primary publication

Parent, Xavier; Benzmüller, Christoph (2023): Normative Conditional Reasoning as a Fragment of HOL. Online: arXiv, S. 1-22, doi: 10.48550/arxiv.2308.10686.

Legal Notice

This work is protected by copyright and/or the indication of a licence. You are free to use this work in any way permitted by the copyright and/or the licence that applies to your usage. For other uses, you must obtain permission from the rights-holder(s).

This document is made available under a Creative Commons license.



The license information is available online: https://creativecommons.org/licenses/by/4.0/legalcode

Normative Conditional Reasoning as a Fragment of HOL

Xavier Parent^a and Christoph Benzmüller^{b c}

^a Technische Universität Wien, Favoritenstrasse 9, A-1040 Wien, Austria ^b Universität Bamberg, An der Weberei 5, 96047 Bamberg, Germany ^c Freie Universität Berlin, Arnimallee 7, 14195 Berlin, Germany

ARTICLE HISTORY

Compiled August 23, 2023

ABSTRACT

We report some results regarding the mechanization of normative (preferencebased) conditional reasoning. Our focus is on Åqvist's system **E** for conditional obligation (and its extensions). Our mechanization is achieved via a shallow semantical embedding in Isabelle/HOL. We consider two possible uses of the framework. The first one is as a tool for meta-reasoning about the considered logic. We employ it for the automated verification of deontic correspondences (broadly conceived) and related matters, analogous to what has been previously achieved for the modal logic cube. The second use is as a tool for assessing ethical arguments. We provide a computer encoding of a well-known paradox in population ethics, Parfit's repugnant conclusion. Whether the presented encoding increases or decreases the attractiveness and persuasiveness of the repugnant conclusion is a question we would like to pass on to philosophy and ethics.

KEYWORDS

Conditional obligation; betterness; Isabelle/HOL; automated theorem proving; population ethics; mere addition/repugnant conclusion paradox

1. Introduction

We report some results regarding the mechanization of normative (preference-based) conditional reasoning. Our focus is on Åqvist's system **E** for conditional obligation (and its extensions). Our mechanization is achieved via a shallow semantical embedding in Isabelle/HOL adapting the methods used by Benzmüller et al. (2015). To look at Standard Deontic Logic (SDL) and extensions (Chellas, 1980; Parent and van der Torre, 2021) would not be very interesting. First, no new insights would be gained, since SDL is a normal modal logic of type KD, which is already covered by the prior work of Benzmüller et al. Second, SDL is vulnerable to the well-known deontic paradoxes, including in particular Chisholm's paradox of contrary-to-duty obligation, see Parent and van der Torre (2021) for details. We thus focus here on Dyadic Deontic Logics (DDL) with a preference-based semantics, which originate from the works of Hansson (1969) and Lewis (1973). In these works, one uses an "intensional" conditional to represent conditional obligation sentences that is weaker than the one obtained us-

CONTACT X. Parent. Email: x.parent.xavier@gmail.com; C. Benzmüller: christoph.benzmueller@unibamberg.de

ing material implication. The semantics generalizes that of SDL: the SDL-ish binary classification of states into good/bad is relaxed to allow for grades of ideality (best, second-best, ...). In this framework, which is particularly popular in deontic logic (cf. the overview chapter by Parent (2021) in the second volume of the *Handbook of Deontic Logic*), a preference relation \succeq ranks the possible worlds in terms of comparative goodness or betterness.¹ The conditional obligation of ψ , given φ (notation: $\bigcirc(\psi/\varphi)$) is evaluated as true if the best φ -worlds are all ψ -worlds. Like in modal logic, different properties of the betterness relation yield different systems.

In this paper, our emphasis is on two possible uses of the mechanized tool. First, we employ it as a tool for meta-reasoning about the considered logics. So far the correspondences between properties and modal axioms have been established "with pen and paper". This raises the question of how much of these correspondences can be automatically explored by modern theorem-proving technology. The automatic verification of correspondences can be done for the modal cube (Benzmüller et al., 2015). We want to understand if it can also be done in DDL. As explained by Benzmüller et al. (2015) we believe that "automation facilities could be very useful for the exploration of the meta-theory of other logics, for example, conditional logics, since the overall methodology is obviously transferable to other logics of interest". Here we follow up on that suggestion, building on further prior results from Benzmüller et al. (2019), where the weakest available system (called \mathbf{E}) has faithfully been embedded in Higher-Order Logic (HOL). In the present paper we consider extensions of E. We look at connections or correspondences between axioms and semantic conditions as "extracted" by relevant soundness and completeness theorems. Thus, we take "correspondence" in the same (broad) sense as Hughes and Cresswell, who write:

"D, T, K4, KB [are] produced by adding a single axiom to K and [...] in each case the system turns out to be characterized by [sound and complete wrt] the class of models in which [the accessibility relation] R satisfies a certain condition. When such a situation obtains—i.e. when a system K+ α is characterized by the class of all models in which R satisfies a certain condition—we shall [...] say [...] that the wff α itself is characterized by that condition, or that the condition *corresponds* [their italics] to α ." (Hughes and Cresswell, 1984, p. 41)

The second use we consider for our mechanized system is as a tool for assessing ethical arguments in philosophical debates. As an illustration, we look at one of the well-known impossibility theorems in population ethics, the so-called "repugnant conclusion" due to Parfit (1984). We provide a computer encoding of the repugnant conclusion in order to make it amenable to formal analysis and computer-assisted experiments. We believe that the formalisation has the potential to further stimulate the philosophical debate on the repugnant conclusion, since the simplifications achieved are indeed quite far-reaching. In particular, our formalization indicates the possibility of a new take on the scenario. It has been suggested that, since the transitivity of "better than" is presupposed in the impossibility theorem, one can avoid such a result by simply giving up on transitivity. The solution is sometimes dismissed on the ground that it is too radical. The formalization reveals a less radical variant solution, which consists in weakening transitivity rather than giving it up wholesale. However, not all candidate weakenings of transitivity will do. For instance, a-cyclicity works fine, but not the interval order condition.

The paper is organized as follows. Section 2 recalls system \mathbf{E} and its extensions. Section 3 shows the embedding of \mathbf{E} in Isabelle/HOL. Section 4 studies the corre-

¹For $i \succeq j$, read "*i* is at least as good as *j*".

spondence between the properties of the betterness relation and the axioms. Section 5 discusses the repugnant conclusion. Section 6 concludes.²

2. System E

We describe the semantics and proof theory of system E and its extensions. This one introduces the primitive symbol $\bigcirc (_/_)$ for "it is obligatory that ... given that ...", from which symbol $P(_/_)$ for "it is permitted that ... given that ..." is defined. The language also has \Box and \diamondsuit .

2.1. Semantics

We start with the main ingredients of the semantics. A preference model is a structure $M = (W, \succeq, V)$, where W is a non-empty set of possible worlds, \succeq is a preference relation ranking elements of W in terms of betterness or comparative goodness, and V is a function assigning to each propositional letter a subset of W (intuitively, the subset of those worlds where the propositional letter is true). $a \succeq b$ may be read "a is at least as good as b". \succ is the strict counterpart of \succeq , defined by $a \succ b$ (a is strictly better than b) iff $a \succeq b$ and $b \succeq a$. \approx is the equal goodness relation, defined by $a \approx b$ (a and b are equally good) iff $a \succeq b$ and $b \succeq a$.

The truth conditions for modal and deontic formulas read:

- $M, a \vDash \Box \varphi$ iff $\forall b \in W$ we have $M, b \vDash \varphi$
- $M, a \vDash \bigcirc (\psi/\varphi)$ iff $\forall b \in \text{best}(\varphi)$ we have $M, b \vDash \psi$

When no confusion can arise, we omit the reference to M and simply write $a \models \varphi$. Intuitively, $\bigcirc(\psi/\varphi)$ is true if the best φ -worlds are all ψ -worlds. There is variation among authors regarding the formal definition of "best". It is sometimes cast in terms of maximality (we call this the max rule) and some other times cast in terms of optimality (we call this the opt rule). An φ -world a is maximal if it is not (strictly) worse than any other φ -world. It is optimal if it is at least as good as any φ -world. The two notions coincide only when "gaps" (incomparabilities) in the ranking are ruled out. Formally:

Max rule	Opt rule
$best(\varphi) = max(\varphi)$	$best(\varphi) = opt(\varphi)$

where

$$a \in \max(\varphi) \Leftrightarrow a \models \varphi \& \neg \exists b \ (b \models \varphi \& b \succ a)$$
$$a \in \operatorname{opt}(\varphi) \Leftrightarrow a \models \varphi \& \forall b \ (b \models \varphi \rightarrow a \succeq b)$$

The relevant properties of \succeq are (universal quantification over worlds is left implicit):

- Reflexivity: $a \succeq a$;
- Transitivity: if $a \succeq b$ and $b \succeq c$, then $a \succeq c$;
- totalness or (strong) connectedness: $a \succeq b$ or $b \succeq a$ (or both);

²The theory file is available for downloading at http://logikey.org under sub-repository "/Deontic-Logics/cube-dll/" (files "cube.thy" and "mere_addition.thy").

- Interval order: \succeq is reflexive and Ferrers (if $a \succeq b$ and $c \succeq d$, then $a \succeq d$ or $c \succeq b$);
- A-cyclicity: \succ contains no cycles of the form $a_1 \succ a_2 \succ ... \succ a_n \succ a_1$.

A-cyclicity and the interval order condition are discussed in Parent (2022). They are two weakenings of the assumption of transitivity. In particular, the interval order condition makes room for the idea of non-transitive equal goodness relation due to discrimination thresholds. These are cases where $a \approx b$ and $b \approx c$ but $a \not\approx c$ (see Luce 1956).

Lewis' limit assumption is meant to rule out sets of worlds without a "limit" (viz. a best element). Its exact formulation varies among authors. It exists in (at least) the following four versions, where best $\in \{\max, opt\}$

$$\underline{\text{Limitedness}} \\
 \text{If } \exists x \text{ s.t. } x \models \varphi \text{ then } \text{best}(\varphi) \neq \emptyset \\
 \underline{\text{Smoothness}} \text{ (or stopperedness)} \\
 \text{If } x \models \varphi, \text{ then: either } x \in \text{best}(\varphi) \text{ or } \exists y \text{ s.t. } y \succ x \& y \in \text{best}(\varphi) \\
 \text{ (SM)}$$

A betterness relation \succeq will be called "opt-limited" or "max-limited" depending on whether (LIM) holds with respect to opt or max. Similarly, it will be called "opt-smooth" or "max-smooth" depending on whether (SM) holds with respect to opt or max. For pointers to the literature, and the relationships between these versions of the limit assumption, see Parent (2014).

The above semantics may be viewed as a special case of the selection function semantics favored by Stalnaker and generalized by Chellas (1975). The preference relation is replaced with a selection function f from formulas to subsets of W, such that, for all φ , $f(\varphi) \subseteq W$. Intuitively, $f(\varphi)$ outputs all the best φ -worlds. The evaluation rule for the dyadic obligation operator is thus given as: $\bigcirc(\psi/\varphi)$ holds when $f(\varphi) \subseteq ||\psi||$, where $||\psi||$ is the set of ψ -worlds. It is known that when suitable constraints are put on the selection function, the two semantics validate exactly the same set of formulas– cf. Parent (2015) for details.³ The correspondence between constraints put on the selection function and modal axioms have been verified by automated means by Benzmüller et al. (2012). A comparison between this prior study and ours is left as a topic for future research.

2.2. Systems

The relevant systems are shown in Fig. 1. A line between two systems indicates that the system to the left is strictly included in the system to the right.



³One can go one step further, and make the selection function semantics an instance of a more general semantics equipped with a neighborhood function, like in traditional modal logic (cf. Chellas (1975)).

All contain the classical propositional calculus; they then add the following schemata:

• For **E** (the naming follows Parent (2021)):

S5-schemata for
$$\Box$$
 (S5)
 $\bigcirc (\psi \to \xi/\varphi) \to (\bigcirc (\psi/\varphi) \to \bigcirc (\xi/\varphi)$ (COK)
 $\bigcirc (\psi/\varphi) \to \Box \bigcirc (\psi/\varphi)$ (Abs)
 $\Box \varphi \to \bigcirc (\varphi/\psi)$ (Nec)

$$\Box(\varphi \leftrightarrow \psi) \to (\bigcirc(\xi/\varphi) \leftrightarrow \bigcirc(\xi/\psi))$$
(Ext)

$$\bigcirc (\varphi/\varphi)$$
 (Id)

$$\bigcirc (\xi/\varphi \land \psi) \to \bigcirc (\psi \to \xi/\varphi)$$
 (Sh)

If
$$\vdash \varphi$$
 then $\vdash \Box \varphi$ (N)

• For **F**: axioms of **E** plus

$$\Diamond \varphi \to (\bigcirc (\psi/\varphi) \to P(\psi/\varphi))$$
 (D^{*})

• For \mathbf{F} +(CM): axioms of \mathbf{F} plus

$$(\bigcirc(\psi/\varphi)\land\bigcirc(\xi/\varphi))\to\bigcirc(\xi/\varphi\land\psi)$$
 (CM)

• For \mathbf{F} +(DR): axioms of \mathbf{F} plus

$$\bigcirc (\xi/\varphi \lor \psi) \to (\bigcirc (\xi/\varphi) \lor \bigcirc (\xi/\psi))$$
(DR)

 $\bullet\,$ For ${\bf G}:$ axioms of ${\bf F}$ plus:

$$(P(\psi/\varphi) \land \bigcirc (\psi \to \xi/\varphi)) \to \bigcirc (\xi/\varphi \land \psi)$$
 (Sp)

We give an intuitive explanation for these axioms. COK is the conditional analogue of the familiar distribution axiom K. Abs is the absoluteness axiom of Lewis (1973), and reflects the fact that the ranking is not world-relative. Nec is the deontic counterpart of the familiar necessitation rule. Ext permits the replacement of necessarily equivalent formulas in the antecedent of deontic conditionals. Id is the deontic analogue of the identity principle. D^{*} rules out the possibility of conflicts between obligations, for a "consistent" context A. CM and DR correspond to the principle of cautious monotony and disjunctive rationality from the non-monotonic logic literature. CM tells us that complying with an obligation does not modify the other obligations arising in the same context. DR tells us that if a disjunction of states of affairs triggers an obligation, then at least one disjunct triggers this obligation. Due to Spohn, Sp is equivalent with the principle of rational monotony; $\bigcirc(\psi \to \xi/\varphi)$ is changed into $\bigcirc(\xi/\varphi)$. The principle says that realizing a permission does not modify the other obligations arising in the same context.

We give below the main soundness and completeness theorems. Those stated in Th. 2.1 hold under both the opt rule and the max rule. It is understood that limitedness is cast in terms of opt when the opt rule is applied, and in terms of max when the max rule is applied. The same holds for smoothness.

Theorem 2.1 (Soundness and completeness, Parent2021; 2022). (i) **E** is sound and complete w.r.t. the class of all preference models; (ii) **F** is sound and complete w.r.t. the class of preference models in which \succeq is limited; (iii) **F**+CM is sound and complete w.r.t. the class of preference models in which \succeq is smooth; (iv) **F**+DR is (weakly) sound and complete w.r.t. the class of (finite) preference models in which \succeq meets the interval order condition.

Theorem 2.2 (Soundness and completeness, Parent2014). (i) Under the opt-rule **G** is sound and complete w.r.t. the class of preference models in which \succeq is limited and transitive; (ii) under the max-rule, **G** is sound and complete w.r.t. the class of preference models in which \succeq is limited, transitive and total.

For more background on these systems, see Parent (2021) and the references therein.

2.3. Correspondences

Table 1 shows some of the known "correspondences" between semantic properties and formulas that can be extracted from Th. 2.1 and Th. 2.2. The leftmost column shows the properties of \succeq . The two middle columns show the corresponding modal axioms, the first column for the max rule, and the second one for the opt rule. As before it is understood that smoothness (resp. limitedness) is defined for max in the max column, and for opt in the opt column. The rightmost column gives the paper where the completeness theorem is established. The symbol \times indicates that the property (or pair of properties) is known not to correspond to any axiom, in the sense that the property does not modify the set of valid formulas. On the fifth (resp. seventh) line the parenthesis "(+smoothness)" (resp. "(+limitedness)") indicates that smoothness (resp. limitedness) is assumed in the background.⁴

Property	Formula (max)	Formula (opt)	Reference
reflexivity	×	×	Parent (2015)
totalness	×	×	Parent (2015)
limitedness	D^{\star}	D^{\star}	Parent (2015)
smoothness	CM	CM	Parent (2014)
transitivity (+smoothness)	×	Sp	Parent $(2014,2)$
transitivity+totalness	Sp	×	Parent (2014)
interval order (+limitedness)	DR	DR	Parent (2022)

Table 1.: Some correspondence

3. System E in Isabelle/HOL

Our modelling of System E in Isabelle/HOL reuses and adapts prior work (Benzmüller et al., 2019) and it instantiates and applies the LogiKEy methodology (Benzmüller et al., 2020), which supports plurality at different modelling layers.

3.1. LogiKEy

Classical higher-order logic (HOL) is fixed in the LogiKEy methodology and infrastructure (Benzmüller et al., 2020) as a *universal meta-logic* (Benzmüller, 2019) at the base layer (L0), on top of which a plurality of (combinations of) object logics can

⁴Even though smoothness does not play any apparent role in the validation of the axiom, the completeness result is for a class of models satisfying this property.

```
1 theory DDLcube imports Main
  2
3 begin
  4
5
      (*** We introduce Aqvist's system E from the 2019 IfColog paper ***)
  6
  7 typedecl i (* Possible worlds *)
  <sup>8</sup> type_synonym \sigma = "(i\Rightarrowbool)"
  s_{i} type_synonym \alpha = "i\Rightarrow \sigma" (* Type of <u>betterness</u> relation between worlds *)
 10 type_synonym \tau = "\sigma \Rightarrow \sigma"
11
12
13 consts aw::i (* Actual world *)
<sup>14</sup> abbreviation etrue :: "\sigma" ("\top") where "\top \equiv \lambda w. True"
abbreviation efalse :: "\sigma" ("\perp") where "\perp \equiv \lambdaw. False"
abbreviation enot :: "\sigma \Rightarrow \sigma" ("\neg_"[52]53) where "\neg \varphi \equiv \lambdaw. \neg \varphi(w)"
abbreviation eand :: "\sigma \Rightarrow \sigma \Rightarrow \sigma" (infixr" \wedge "51) where "\varphi \land \psi \equiv \lambda w. \varphi(w) \land \psi(w)"
abbreviation eor :: "\sigma \Rightarrow \sigma \Rightarrow \sigma" (infixr"\lor"50) where "\varphi \lor \psi \equiv \lambda w. \varphi(w) \lor \psi(w)"
abbreviation eimpf :: "\sigma \Rightarrow \sigma \Rightarrow \sigma" (infixr"\lor"50) where "\varphi \lor \psi \equiv \lambda w. \varphi(w) \lor \psi(w)"
abbreviation eimpf :: "\sigma \Rightarrow \sigma \Rightarrow \sigma" (infixr"\multimap"49) where "\varphi \rightarrow \psi \equiv \lambda w. \varphi(w) \rightarrow \psi(w)"
abbreviation eimpb :: "\sigma \Rightarrow \sigma \Rightarrow \sigma" (infixr"\leftarrow"49) where "\varphi \leftarrow \psi \equiv \lambda w. \psi(w) \rightarrow \varphi(w)"
21 abbreviation eequ :: "\sigma \Rightarrow \sigma \Rightarrow \sigma" (infixr" \leftrightarrow "48) where "\varphi \leftrightarrow \psi \equiv \lambda w. \varphi(w) \leftrightarrow \psi(w)"
22
23 abbreviation ebox :: "\sigma \Rightarrow \sigma" ("\Box") where "\Box \equiv \lambda \varphi w. \forall v. \varphi(v)"
24 abbreviation ddediomond :: "\sigma \Rightarrow \sigma" ("\diamond") where "\diamond \varphi \equiv \lambda w. \exists v. \varphi(v)"
25
<sub>26</sub> abbreviation evalid :: "\sigma \Rightarrow bool" ("| |"[8]109) (* Global validity *)
        where "|p| \equiv \forall w. p w"
27
28 abbreviation ecjactual :: "\sigma \Rightarrow bool" (" | 1"[7]105) (* Local validity — in world aw *)
        where ||p|| = p(aw)
29
31 consts r :: "a" (infixr "r" 70) (* Betterness relation *)
```

Figure 2.: Basic semantical ingredients; propositional and modal connectives

become encoded (layer L1). In the case of this paper, we encode extensions of System **E** at layer L1 in order to assess them. Employing these object logics notions of layer L1 we can then articulate a variety of logic-based domain-specific languages, theories and ontologies at the next layer (L2), thus enabling the modelling and automated assessment of different application scenarios (layer L3). Note that the assessment studies conducted in this paper at layer L3 do not require any further knowledge to be provided at layer L2; hence layer L2 modellings do not play a role in this paper.

LogiKEy significantly benefits from the availability of theorem provers for HOL, such as Isabelle/HOL which internally provides powerful automated reasoning tools such as *Sledgehammer* (Blanchette et al., 2013; Blanchette et al., 2016) and *Nitpick* (Blanchette and Nipkow, 2010). The automated theorem proving systems integrated via *Sledgehammer* include higher-order ATP systems, first-order ATP systems, and SMT (satisfiability modulo theories) solvers, and many of these systems in turn use efficient SAT solver technology internally. Indeed, proof automation with *Sledgehammer mer* and (counter-)model finding with *Nitpick* were invaluable in supporting our exploratory modeling approach at various levels. These tools were very responsive in automatically proving (*Sledgehammer*), disproving (*Nitpick*), or showing consistency by providing a model (*Nitpick*). In the first case, references to the required axioms and lemmas were returned (which can be seen as a kind of abduction), and in the case of models and counter-models they often proved to be very readable and intuitive. In this section and subsequent ones, we highlight some explicit use cases of *Sledgehammer* and *Nitpick*. They have been similarly applied at all levels as mentioned before.

```
36 abbreviation eopt :: "\sigma \Rightarrow \sigma" ("opt<_>") (* opt rule*)
      where "opt<\varphi> \equiv (\lambdav. ((\varphi)(v) \wedge (\forallx. ((\varphi)(x) \rightarrow v r x))) )"
37
38 abbreviation econdopt :: "\sigma \Rightarrow \sigma \Rightarrow \sigma" ("\odot < |_>")
      where "\odot < \psi | \varphi > \equiv \lambda w. \text{ opt} < \varphi > \subseteq \psi"
39
40 abbreviation eperm :: "\sigma \Rightarrow \sigma \Rightarrow \sigma" ("\mathcal{P} < | >")
      where "\mathcal{P} < \psi | \varphi > \equiv \neg \odot < \neg \psi | \varphi >"
41
42
43 abbreviation emax :: "\sigma \Rightarrow \sigma" ("max< >")
                                                                                      (* Max rule *)
44 where "max\langle \varphi \rangle \equiv (\lambda v. ((\varphi)(v) \land (\forall x. ((\varphi)(x) \longrightarrow (x r v \longrightarrow v r x)))))"

45 abbreviation econd :: "\sigma \Rightarrow \sigma \Rightarrow \sigma" ("\bigcirc < | >")

46 where "\bigcirc \langle \psi | \varphi \rangle \equiv \lambda w. \max \langle \varphi \rangle \subseteq \psi"
47 abbreviation euncobl :: "\sigma \Rightarrow \sigma" ("O<_>")
       where "O<\varphi> \equiv O<\varphi|T>"
48
49 abbreviation ddeperm :: "\sigma \Rightarrow \sigma \Rightarrow \sigma" ("P< | >")
       where "P<\psi | \varphi > \equiv \neg \bigcirc < \neg \psi | \varphi >"
51
52
    (* Settings for model finder Nitpick *)
53
54
55 nitpick_params [user_axioms, show_all, expect=genuine]
56
57
    (*** First consistency check ***)
58
59
60
    lemma True
       nitpick [satisfy] (* model found *)
61
62
       oops
63
    (*** The max-rule and opt-rule don't coincide ***)
64
65
    lemma "\odot < \psi | \varphi > \equiv \bigcirc < \psi | \varphi >"
       nitpick [card i=1] (* counterexample found for card i=1 *)
67
68
       oops
```

Figure 3.: Truth-conditions

3.2. Faithful embedding of system E

In the work of Benzmüller et al. (2019), it is shown that the embedding of **E** in Isabelle/HOL is faithful, in the sense that a formula φ in the language of **E** is valid in the class PREF of all preference models if and only if the HOL translation of φ (notation: $|\varphi|$) is valid in the class of Henkin models of HOL.

Theorem 3.1 (Faithfulness of the embedding).

 $\models^{\text{PREF}} \varphi \text{ if and only if } \models^{HOL} |\varphi|$

Remember that the establishment of such a result is our main success criterium at layer L1 in the LogiKEy methodology.

This first two screenshots show the encoding of \mathbf{E} in Isabelle/HOL. Fig. 2 shows the basic ingredients in the preferential model, and describes how the propositional and alethic modal connectives are handled. The betterness relation \succeq is encoded as a binary relational constant r (l. 31). In Fig. 3, the notions of optimality and maximality are encoded. Different pairs of modal operators (obligation, permission) are introduced to distinguish between the two types of truth-conditions. The model finder *Nitpick* is able to verify the consistency of the formalization (l. 60) and to verify the non-equivalence between the two types of truth-conditions (l. 66). *Sledgehammer* is able to show the validity of the axioms of \mathbf{E} . By presenting a suitable counter-model, *Nitpick* is able to show the invalidity of the axioms pertaining to the stronger systems.

3.3. Properties

The encoding of the properties of the betterness relation are shown in Figs. 4 and 5. On l. 99-104 of Fig. 4, one sees the different versions of Lewis' limit assumption. The

```
87 (**********
88 Properties
    ************************
89
90
    (* The standard properties of the betterness relation *)
91
92
93 abbreviation "reflexivity \equiv (\forall x. x r x)"
abbreviation "transitivity \equiv (\forall x \ y \ z. (x r y \land y r z) \longrightarrow x r z)"
 abbreviation "totalness \equiv (\forall x \ y, (x \ r \ y \lor y \ r \ x))"
96
97 (* 4 versions of Lewis's limit assumption *)
98
abbreviation "mlimitedness \equiv (\forall \varphi. (\exists x. (\varphi)x) \longrightarrow (\exists x. max < \varphi > x))"
100 abbreviation "msmoothness =
        (\forall \varphi \ x. ((\varphi)x \longrightarrow (\max \langle \varphi \rangle x \lor (\exists y. (y \mathbf{r} x \land \neg (x \mathbf{r} y) \land \max \langle \varphi \rangle y)))))"
101
abbreviation "olimitedness \equiv (\forall \varphi. (\exists x. (\varphi)x) \rightarrow (\exists x. opt < \varphi > x))"
abbreviation "osmoothness \equiv
      (\forall \varphi \ x. \ ((\varphi)x \longrightarrow (opt < \varphi > x \lor (\exists y. (y \ r \ x \land \neg (x \ r \ y) \land opt < \varphi > y)))))"
104
```

Figure 4.: Standard properties

property in Fig. 5 is the interval order condition. This one is usually described as the combination of totalness with the Ferrers condition as shown on l. 137. *Sledgehammer* is able to confirm a fact that has been generally overlooked in the literature, namely that totalness can be replaced by the simpler condition of reflexivity (l. 139-141). More

```
144 (* Interval order (reflexivity + Ferrers) *)
145
146 abbreviation Ferrers
147 where "Ferrers ≡ (∀x y z u. (x r u ∧ y r z) → (x r z ∨ y r u))"
148
149 theorem T2:
150 assumes Ferrers and reflexivity (*fact overlooked in the literature*)
151 shows totalness
152 by (simp add: assms(1) assms(2)) (* proof found *)
```

Figure 5.: Interval order

weakenings of transitivity Parent (2022) are encoded in the theory file. For simplicity's sake we give the example of quasi-transitivity and a-cyclicity. The encoding of the second is shown in Fig. 6. The formal definition of this notion reads: if $a \succ^* b$, where \succ^* is the transitive closure of \succ , then $b \not\succeq a$. In Isabelle/HOL, the transitive closure of a relation can be defined in a few lines, shown in Fig. 7.

```
139 (* A-cyclicity: cycles of strict <u>betterness</u> are ruled out*)
140
141 abbreviation loopfree
142 where "loopfree \equiv \forall x \ y. \ tcr_strict \ x \ y \longrightarrow (y \ r \ x \longrightarrow x \ r \ y)"
```

```
Figure 6.: A-cyclicity
```

Quasi-transitivity requires the strict betterness relation be transitive. Its encoding is shown in Fig. 8.

```
106 (* Weaker forms of transitivity--they require the notion of
107 transitive closure*)
108
definition transitive :: "\alpha \Rightarrow bool"
      where "transitive Rel \equiv \forall x \ y \ z. Rel x y \land Rel y z \longrightarrow Rel x z"
110
and definition sub rel :: "\alpha \Rightarrow \alpha \Rightarrow bool"
     where "sub_rel Rel1 Rel2 \equiv \forallu v. Rel1 u v \longrightarrow Rel2 u v"
-112
<sup>113</sup> definition assfactor::"\alpha \Rightarrow \alpha"
      where "assfactor Rel \equiv \lambda u v. Rel u v \land \negRel v u "
-114
115
    (* In HOL the transitive closure of a relation can be defined in a single line; here
116
        we apply the construction to \underline{betterness} relation \mathbf{r} and for its strict
117
        variant (\lambda u v. u r v \wedge \neg v r u) *)
118
119 definition tcr
      where "tcr \equiv \lambda x y. \forall Q. transitive Q \longrightarrow (sub rel r Q \longrightarrow Q \times y)"
-120
121
122 definition tcr strict
123 where "tcr_strict \equiv \lambda x y. \forall Q. transitive Q \longrightarrow (sub_rel (\lambda u v. u r v \land \neg v r u) Q \longrightarrow Q x y)"
```

Figure 7.: Transitive closure

Figure 8.: Quasi-transitivity

4. Correspondences

4.1. Max rule

Here we check known correspondences between modal axioms under the max rule.

First, *Nitpick* is able to confirm that the formula is not valid unless the matching property is assumed. Figs. 9 and 12 show that, when the relevant property is not assumed, counter-models for D^* , CM, DR and Sp are found by *Nitpick*.

```
189 (* Max-Limitedness corresponds to D *)
190
191 Lemma "|\Diamond \varphi \rightarrow (\bigcirc \langle \psi | \varphi \rangle \rightarrow \mathsf{P} \langle \psi | \varphi \rangle)|"
        nitpick [card i=3] (* counterexample found for card i=3 *)
192
193
        oops
194
195
     lemma " [ (\bigcirc <\psi | \varphi > \land \bigcirc <\chi | \varphi >) \rightarrow \bigcirc <\chi | \varphi \land \psi > ] "
        nitpick [card i=3] (* counterexample found *)
196
197
        oops
198
     \mathsf{lemma} \ "|\bigcirc <\chi | (\varphi \lor \psi) > \to ((\bigcirc <\chi | \varphi >) \lor (\bigcirc <\chi | \psi >)) | "
199
        nitpick [card i=3] (* counterexample found *)
200
201
        oops
```

Figure 9.: D^{*}, CM and DR invalid in general

In Figs. 10, 11 and 12, it is confirmed that if the property is assumed, then the axiom is validated. Thus, the implications having the form "property \Rightarrow axiom" are all verified; Fig. 10 shows it for limitedness and smoothness, Fig. 11 for the interval order condition, and Fig. 12 for the combination of transitivity and totalness. But the converse implications are all falsified by *Nitpick*. We will come back to this point later on.

```
203 theorem T8:
        assumes mlimitedness
204
        shows "D*": "[\Diamond \varphi \rightarrow \bigcirc \langle \psi | \varphi \rangle \rightarrow \mathsf{P} \langle \psi | \varphi \rangle]"
-205
206
        sledgehammer
        using assms
207
        oops
-208
209
210 lemma
        assumes "D*": "[\Diamond \varphi \rightarrow \neg (\bigcirc \prec \psi | \varphi \land \bigcirc \prec \neg \psi | \varphi \rangle)]"
211
-212
        shows mlimitedness
        nitpick [card i=3] (* counterexample found *)
213
-214
        oops
215
     (* Smoothness corresponds to cautious monotony *)
216
217
218 theorem T9:
        assumes msmoothness
219
-220
        shows CM: "[(\bigcirc <\psi | \varphi > \land \bigcirc <\chi | \varphi >) \rightarrow \bigcirc <\chi | \varphi \land \psi >]"
        using assms by force
-221
222
223 lemma
        assumes CM: "[(\bigcirc <\psi | \varphi > \land \bigcirc <\chi | \varphi >) \rightarrow \bigcirc <\chi | \varphi \land \psi >]"
224
-225
        shows msmoothness
        nitpick [card i=3]
                                         (* counterexample found *)
226
-227
       oops
```

Figure 10.: Limit assumption

```
230 lemma
231
        assumes reflexivity
         shows DR: [\bigcirc \langle \chi | \varphi \lor \psi \rangle \rightarrow (\bigcirc \langle \chi | \varphi \lor \lor \bigcirc \langle \chi | \psi \rangle)]
-232
        nitpick [card i=3] (* counterexample found *)
233
        oops
-234
235
236 theorem T10:
        assumes reflexivity and Ferrers
237
        shows DR: [\bigcirc <\chi | (\varphi \lor \psi) > \rightarrow (\bigcirc <\chi | \varphi > \lor \bigcirc <\chi | \psi >)]
-238
        by (metis assms)
-239
240
241 lemma
        assumes DR: [\bigcirc \langle \chi | \varphi \lor \psi \rangle \rightarrow (\bigcirc \langle \chi | \varphi \lor \lor \bigcirc \langle \chi | \psi \rangle)]
242
        shows reflexivity
-243
        nitpick [card i=1] (* counterexample found *)
244
        oops
245
246
247 lemma
        assumes DR: [\bigcirc \langle \chi | \varphi \lor \psi \rangle \rightarrow (\bigcirc \langle \chi | \varphi \lor \lor \bigcirc \langle \chi | \psi \rangle)]
248
        shows Ferrers
-249
        nitpick [card i=2] (* counterexample found *)
250
-251
        oops
```

Figure 11.: Interval order

4.2. Opt rule

The outcomes of our experimentation are the same as for the max rule except for one small change. Transitivity no longer needs totalness to validate Sp. This one only needs transitivity. Besides, the assumption of transitivity of the betterness relation gives us a principle of transitivity for a weak preference operator over formula, defined by $\varphi \geq \psi$ iff $P(\varphi/\varphi \lor \psi)$. This is shown in Fig. 13.

```
253 (*Transitivity and totalness corresponds to the Spohn axiom (Sp)*)
254
     lemma
255
        assumes transitivity
256
        shows Sp: "[( P < \psi | \varphi > \land \bigcirc < (\psi \rightarrow \chi) | \varphi >) \rightarrow \bigcirc < \chi | (\varphi \land \psi) >]"
-257
        nitpick [card i=3] (* counterexample found *)
258
-259
        oops
260
261 lemma
262
        assumes totalness
        shows Sp: "[( P < \psi | \varphi > \land \bigcirc <(\psi \rightarrow \chi) | \varphi >) \rightarrow \bigcirc <\chi | (\varphi \land \psi) >]"
-263
        nitpick [card i=3] (* counterexample *)
264
265
        oops
266
267 theorem T11:
268
        assumes transitivity and totalness
        shows Sp: "[( P < \psi | \varphi > \Lambda \bigcirc <(\psi \to \chi) | \varphi >) \rightarrow \bigcirc <\chi | (\varphi \land \psi) >]"
-269
        by (metis assms)
-270
271
272 theorem T12:
        assumes transitivity and totalness
273
        shows transit: "|( P < \varphi | \varphi \lor \psi > \land P < \psi | \psi \lor \chi >) \rightarrow P < \varphi | (\varphi \lor \chi) >|"
-274
        by (metis assms(1) assms(2))
-275
276
277 lemma
         \text{assumes} \quad \text{Sp: } "\lfloor ( \mathsf{P}\!\!<\!\!\psi | \varphi \!\!> \land \bigcirc \!\!<\!\!(\psi \!\!\rightarrow\!\!\chi) | \varphi \!\!>) \rightarrow \bigcirc \!\!<\!\!\chi | (\varphi \!\land\! \psi) \!\!> \rfloor " 
278
        shows totalness
-279
        nitpick [card i=1] (* counterexample found *)
280
-281
        oops
282
283 Lemma
        assumes Sp: "[( P < \psi | \varphi > \Lambda \bigcirc <(\psi \rightarrow \chi) | \varphi >) \rightarrow \bigcirc <\chi | (\varphi \land \psi) >]"
284
        shows transitivity
-285
        nitpick [card i=2] (* counterexample found *)
286
-287
        oops
```

Figure 12.: Transitivity and totalness (max)

```
321 (*transitivity*)
322
323 theorem T15:
        assumes transitivity
324
        shows Sp': "[( \mathcal{P} < \psi | \varphi > \land \odot < (\psi \to \chi) | \varphi >) \to \odot < \chi | (\varphi \land \psi) >]"
325
        by (metis assms)
326
327
328 theorem T16:
        assumes transitivity
329
        shows Trans: "[( \mathcal{P} < \varphi | \varphi \lor \psi > \land \mathcal{P} < \psi | \psi \lor \xi > ) \rightarrow \mathcal{P} < \varphi | \varphi \lor \xi >]"
330
        by (metis assms)
331
332
333 lemma
        assumes Sp: "[( \mathcal{P}<\psi | \varphi > \land \odot <(\psi \rightarrow \chi) | \varphi >) \rightarrow \odot <\chi | (\varphi \land \psi) >]"
334
        assumes Trans: "[( \mathcal{P} < \varphi | \varphi \lor \psi > \land \mathcal{P} < \psi | \psi \lor \xi >) \rightarrow \mathcal{P} < \varphi | \varphi \lor \xi >]"
335
        shows transitivity
336
        nitpick [card i=2] (* counterexample found *)
337
        oops
338
```

Figure 13.: Transitivity (opt)

4.3. Inclusion

In the work of Benzmüller et al. (2015), proper inclusion between systems in the modal cube are verified by looking at the model constraints of their respective axiomatizations. Because of the lack of full equivalence between modal axiom and property of the relation, we cannot do the same, at least not yet. Nor can we show equivalence between systems when restraining the number of worlds.

4.4. The $\exists \forall$ truth-conditions (Lewis)

Variant evaluation rules have been proposed for the conditional in order to handle some of the problems encountered with the usual pattern of evaluation in terms of best. We take the example of Lewis (1973)'s evaluation rule. In order to avoid commitment to the limit assumption, Lewis suggests that $\bigcirc(\psi/\varphi)$ should be true whenever there is no φ -world or there is a $\varphi \wedge \psi$ -world which starts a (possibly infinite) sequence of increasingly better $\varphi \wedge \psi$ -worlds. Formally:

$$a \models \bigcirc (\psi/\varphi) \text{ iff } \neg \exists b \ (b \models \varphi) \text{ or} \\ \exists b \ (b \models \varphi \land \psi \& \forall c \ (c \succeq b \Rightarrow c \models \varphi \rightarrow \psi)) \tag{(\exists \forall)}$$

We shall refer to the statement appearing at the right-hand-side of "iff" as the $\exists \forall$ rule. The encoding is shown in Fig.14.

abbreviation lewcond :: "
$$\sigma \Rightarrow \sigma \Rightarrow \sigma$$
" (" $\circ < [>>$ ")
where " $\circ < \psi | \varphi > \equiv \lambda v$. ($\neg (\exists x. (\varphi)(x)) \lor$
1 ($\exists x. ((\varphi)(x) \land (\psi)(x) \land (\forall y. ((y r x) \longrightarrow (\varphi)(y) \longrightarrow (\psi)(y)))))$ "
abbreviation lewperm :: " $\sigma \Rightarrow \sigma \Rightarrow \sigma$ " (" $\int < [>>$ ")
where " $\int < \psi | \varphi > \equiv \neg \circ < \neg \psi | \varphi >$ "
4
lemma True nitpick [satisfy,user_axioms,expect=genuine]
96_ oops

Figure 14.: $\exists \forall$ rule

Isabelle/HOL is able to verify in what sense the standard account in terms of best requires the limit assumption. The law "from $\Diamond \varphi$, $\bigcirc (\psi/\varphi)$ and $\bigcirc (\neg \psi/\varphi)$ infer $\bigcirc (\chi/\varphi)$ " is valid. This is known as the principle of "deontic explosion", often called DEX. It says that, in the presence of a conflict of duties (unless it is triggered by an "inconsistent" state of affairs) everything becomes obligatory. This has led most authors to make the limitedness assumption in order to validate D*, and hence make DEX harmless: the set { $\Diamond \varphi$, $\bigcirc (\psi/\varphi)$, $\bigcirc (\neg \psi/\varphi)$ } is not satisfiable. This is shown in Fig. 15. On l. 371, the validity of DEX is established under the max rule. On l. 374, DEX is falsified under the $\exists \forall$ rule.

Isabelle/HOL is also able to verify that when all the standard properties of the betterness relation are assumed, then the three evaluation rules collapse. This is shown in Fig. 15. L. 380-384 show the equivalence between the $\exists \forall$ rule and the opt rule, and l. 386-390 show the equivalence between the $\exists \forall$ rule and the max rule.

Questions of correspondence between properties and modal axioms are still under investigation. There are two extra complications. First, a completeness result is available for the strongest system **G** only: it is complete with respect to the class of models in which \succeq is transitive and total (and hence reflexive). Second, only two properties seem to have an import, but the matching between them and the axioms is not oneto-one: one property validates more than one axiom, sometimes in combination with the other property. This is shown in Table 2. The left column gives the axiom. The right column shows the property (or pair of properties) required to validate this one.

```
367 (***********
368 Relationship Lewis rule and max/opt rule
      ***************
369
370
371 (* deontic explosion-max rule *)
P_{372} \text{ theorem DEX: } "\lfloor (\Diamond \varphi \land \bigcirc \prec \psi | \varphi > \land \bigcirc \prec \neg \psi | \varphi >) \rightarrow \bigcirc \prec \chi | \varphi > \rfloor "
        by blast
373
374
\begin{array}{l} & \text{375} \text{ (* no-} \underline{\text{deontic}} \text{ explosion-} \underline{\text{lewis}} \text{ rule *)} \\ & \text{376} \text{ lemma DEX: } "[( \Diamond \varphi \land \circ \prec \psi | \varphi \succ \land \circ \prec \neg \psi | \varphi \succ) \rightarrow \circ \prec \chi | \varphi \succ]" \end{array}
         nitpick [card i=2] (* counterexample found*)
377
         oops
378
379
380 theorem T18:
         assumes mlimitedness and transitivity and totalness
381
         shows \| \circ \psi \| \varphi \to \phi \circ \psi \| \varphi \|
382
         sledgehammer
383
         by (smt (z3) assms)
384
385
386 theorem T19:
         assumes mlimitedness and transitivity and totalness
387
         shows "\lfloor \circ < \psi \mid \varphi > \leftrightarrow \bigcirc < \psi \mid \varphi > \rfloor"
388
         sledgehammer
389
         by (smt (z3) assms)
390
```

Figure 15.: Deontic explosion (DEX)

Axiom of \mathbf{G}	Property (or pair of properties) of \succeq
(D*)	totalness
(Sp)	transitivity

(Sp)	transitivity
(COK)	transitivity and totalness
(CM)	transitivity and totalness

Table 2.: Axioms and properties under the $\exists \forall$ rule-from Parent (2021)

```
476 (*axioms of E holding irrespective of the properties of r*)
477
478 theorem Abs: | \circ \langle \psi | \varphi \rangle \rightarrow \Box \circ \langle \psi | \varphi \rangle |
         by blast
479
480
<sup>481</sup> theorem Nec: "[\Box \psi \rightarrow \circ \langle \psi | \varphi \rangle]"
-482
         by blast
483
<sup>484</sup> theorem Ext: "[\Box(\varphi_1 \leftrightarrow \varphi_2) \rightarrow (\circ < \psi | \varphi_1 > \leftrightarrow \circ < \psi | \varphi_2 >)]"
-485
         by simp
486
_{487} theorem Id: "| \circ < \varphi | \varphi > |"
        by auto
488
 489
490 theorem Sh: [\circ \langle \psi | \varphi_1 \land \varphi_2 \rangle \rightarrow \circ \langle \langle \varphi_2 \rightarrow \psi \rangle | \varphi_1 \rangle]^*
          hy hlact
```

Figure 16.: Axioms independent of the properties $(\exists \forall \text{ rule})$

Figure 17.: Transitivity and totalness alone $(\exists \forall \text{ rule})$

```
\begin{array}{l} \text{lemma Sp: } & \|\left(\int \langle \psi | \varphi \rangle \wedge \circ \langle (\psi \rightarrow \chi) | \varphi \rangle \right) \rightarrow \circ \langle \chi | (\varphi \wedge \psi) \rangle \right]^{"} \\ & \text{nitpick } (* \underline{\text{countermodel}} *) \\ & \text{oops} \\ & \text{acl} \\ \text{acl} \\ \text{assumes } \\ \text{shows Sp: } & \| (\int \langle \psi | \varphi \rangle \wedge \circ \langle (\psi \rightarrow \chi) | \varphi \rangle) \rightarrow \circ \langle \chi | (\varphi \wedge \psi) \rangle \right]^{"} \\ & \text{acl} \\ & \text{sledgehammer } (* proof found*) \\ & \text{oops} \end{array}
```

Figure 18.: Transitivity and totalness alone $(\exists \forall \text{ rule, ct'd})$

```
<sup>2</sup>415 lemma COK: "| \circ \langle (\psi_1 \rightarrow \psi_2) | \varphi \rangle \rightarrow (\circ \langle \psi_1 | \varphi \rangle \rightarrow \circ \langle \psi_2 | \varphi \rangle) |"
416
          nitpick [card i=2] (* counterexample found *)
          oops
417
418
419 lemma
          assumes transitivity
420
          shows COK: [\circ<(\psi_1 \rightarrow \psi_2) | \varphi \rightarrow (\circ<\psi_1 | \varphi \rightarrow \circ<\psi_2 | \varphi \rightarrow)]"
nitpick [card i=2] (* counterexample found *)
421
422
          oops
423
424
425 lemma
          assumes totalness
426
          shows COK: (\circ < (\psi_1 \rightarrow \psi_2) | \varphi > \rightarrow (\circ < \psi_1 | \varphi > \rightarrow \circ < \psi_2 | \varphi >) ]
427
          nitpick [card i=3] (* counterexample found *)
428
429
          oops
 430
431 theorem T22:
          assumes transitivity and totalness
432
       shows COK: [\circ<(\psi_1 \rightarrow \psi_2) | \varphi \rightarrow (\circ<\psi_1 | \varphi \rightarrow \circ<\psi_2 | \varphi \rightarrow)]" by (smt (z3) assms)
433
434
435
```

Figure 19.: Transitivity and totalness together $(\exists \forall \text{ rule})$

```
437 Lemma CM: "[(\circ < \psi | \varphi > \land \circ < \chi | \varphi >) \rightarrow \circ < \chi | \varphi \land \psi >]"
        nitpick [card i=2] (* counterexample found *)
438
        oops
439
446
     lemma
441
        assumes transitivity
442
        shows CM: "[(\circ < \psi | \varphi > \land \circ < \chi | \varphi >) \rightarrow \circ < \chi | \varphi \land \psi >]"
443
        nitpick [card i=2] (* counterexample found *)
444
445
        oops
446
     lemma
447
        assumes totalness
448
        shows CM: "|(\circ < \psi | \varphi > \land \circ < \chi | \varphi >) \rightarrow \circ < \chi | \varphi \land \psi >|"
449
        nitpick [card i=3] (* counterexample found *)
456
451
        oops
452
453 theorem T23:
454
        assumes transitivity and totalness
        shows CM'': "[(\circ < \psi | \varphi > \land \circ < \chi | \varphi >) \rightarrow \circ < \chi | \varphi \land \psi >]"
455
        by (metis assms)
456
```

Figure 20.: Transitivity and totalness together $(\exists \forall \text{ rule,ct'd})$

In Fig. 16, *Sledgehammer* shows the validity of the axioms of **E** holding independently of the properties assumed of the betterness relation. In Figs. 17 and 18, *Sledgehammer* confirms that the D^{*} axiom and the Sp axiom call for totalnesss and transitivity, respectively. Similarly, Figs. 19 and 20 show that COK and CM call for *both* transitivity and totalness.

4.5. Discussion

To conclude, with regards to correspondence, the situation for conditional (deontic) logic is still slightly different from the one for traditional modal logic. In the latter setting, the full equivalence between the property of the relation and the modal formula is verified by automated means. In the former setting only the direction "property \Rightarrow axiom" is verified by automated means. To be more precise, what is verified is the fact that, if the property holds, then the axiom holds. What is not confirmed is the converse statement, that if the axiom holds then the property holds. This asymmetry deserves to be discussed.

First, it is usual to distinguish between validity on a frame and validity in a model based on a frame. A frame is a pair $\mathcal{F} = (W, R)$, with W a set of worlds and R the accessibility relation. A model based on $\mathcal{F} = (W, R)$ is the triplet $\mathcal{M} = (W, R, V)$ obtained by adding a specific valuation V, or a specific assignment of truth-values to propositional letters at worlds. For a formula to be valid on a frame \mathcal{F} , it must be valid in all models based on \mathcal{F} . In other words, it must be true for every assignment to the propositional letters. We have worked at the level of models. But in so-called correspondence theory, see e.g. van Benthem (2001), the link between formulas and properties is in general studied at the level of frames themselves. One shows that \mathcal{F} meets a given condition iff formula A is valid on \mathcal{F} . In a recent extension of the semantical embedding approach for public announcement logic PAL, cf. Benzmüller and Reiche (2022), an explicit dependency on the concrete evaluation domain has been modeled. It remains future work to study whether this idea can be further extended and adapted to also support a notion of validity for frames as needed here.

Second, the most we got is that a given property is a sufficient condition for the validity of the axiom, but not a necessary one. For instance, to disprove the implication

```
Free variables:
  \chi = (\lambda x. _)(i_1 := False, i_2 := True, i_3 := False)

\varphi = (\lambda x. _)(i_1 := False, i_2 := False, i_3 := True)
   \phi = (\lambda \mathbf{x})
               )(i1 := False, i2 := False, i3 := False)
Skolem constants:
  \varphi = (\lambda x. \_)(i_1 := True, i_2 := True, i_3 := True) x = i_3
  x = i_2
  \lambda y. x = (\lambda x. _)(i_1 := i_3, i_2 := i_1, i_3 := i_3)
Constant:
  (r) =
    (a) Model for HOL
                              i_1: \neg \varphi, \neg \psi, \neg \chi
                        (\downarrow i_2: \neg \varphi, \neg \psi, \chi)
   (b) Preferential model. An arrow from i_1 to i_2
  means i_1 \succeq i_2. No arrow from i_2 to i_1 means
  i_2 \not\succeq i_1
```

Figure 21.: A non-smooth model validating CM (max)

"CM \Rightarrow m-smoothness" under the max rule (Fig. 10), *Nitpick* exhibits a model in which CM holds and m-smoothness falsified. This model is shown in Fig. 21. The corresponding preferential model is also shown below. Smoothness is falsified, because it contains an infinite loop of strict betterness, making the smoothness condition fail for, e.g., $\varphi \lor \neg \varphi$. But CM (vacuously) holds, because the two conjuncts appearing in the antecedent of the axiom are both false. Indeed, i_3 is a maximal φ -world, and it falsifies ψ and χ . This shows that m-smoothness is not a necessary condition for the axiom to hold.

It is interesting to remark that *Nitpick* always presents a finite standard model. We leave it as a topic for future research to investigate if the crucial distinction between standard and non-standard models for HOL which, according to Andrews (2002), sheds so much light on the mysteries associated with the incompleteness theorems, has a bearing on the issue at hand.

Another open problem concerns the possibility of verifying "negative" results. As shown in Table 1, under the max rule transitivity alone does not correspond to any axiom. Also under both the max rule and the opt rule neither reflexivity nor totalness correspond to an axiom. Finally, under the $\exists \forall$ rule the limit assumption has no import. All this has been established with pen and paper. It would be worth exploring the question as to whether and how this problem could be tackled in Isabelle/HOL.

5. The repugnant conclusion

In this section, we illustrate another possible use of the logical machinery. It consists in employing it for the computer-aided assessment of ethical arguments in philosophy. We use the repugnant conclusion discussed by Parfit (1984). We provide a computer encoding of his argument for the repugnant conclusion in order to make it amenable to formal analysis and computer-assisted experiments. Isabelle/HOL is able to confirm the viability of a possible solution to the paradox, advocated by Temkin (1987) among others. It consists in resolving the paradox by abandoning the assumption of transitivity of "better than". Isabelle/HOL is also able to confirm the availability of a variant solution, which consists in weakening transitivity rather than abandonning it wholesale. Furthermore, one can show that not all the conceivable weakenings of transitivity will do.

The repugnant conclusion reads:

"For any perfectly equal population with very high positive welfare, there is a population with very low positive welfare which is better, other things being equal" (Parfit, 1984, p. 523).

The target is "total utilitarianism", according to which the best outcome is given by the total of well-being in it (Parfit, 1984, p. 387). This view implies that any loss in the quality of lives in a population can be compensated for by a sufficient gain in the quantity of a population. Fig. 22 illustrates the repugnant conclusion. The blocks correspond to two populations, A and Z. The width of each block represents the number of people in the corresponding population, the height represents their quality of life. All the lives in the above diagram have lives worth living. People's quality of life is much lower in Z than in A but, since there are many more people in Z, there is a greater quantity of welfare in Z as compared to A. Consequently, although the people in A lead very good lives and the people in Z have lives only barely worth living, Z is nevertheless better than A according to classical utilitarianism.



Figure 22.: Repugnant conclusion

It has been argued by e.g. Temkin (1987) that the repugnant conclusion can be blocked, by just dropping the assumption of the transitivity of "better than". This is best explained by considering a smaller version of the paradox, called the mere addition paradox. The repugnant conclusion is generated by iteration of the reasoning underlying the mere addition paradox.

The mere addition paradox is shown in Fig. 23. In population A, everybody enjoys a very high quality of life. In population A^+ there is one group of people as large as the group in A and with the same high quality of life. But A^+ also contains a number of people with a somewhat lower quality of life. In Parfit's terminology A^+ is generated from A by "mere addition". Population B has the same number of people as A^+ , their lives are worth living and at an average welfare level slightly above the average in A^+ , but lower than the average in A. The link with the repugnant conclusion is that by reiterating this structure (scenario B^+ and C, C^+ etc.), we end up with a population Z in which all lives have a very low positive welfare.



Figure 23.: Mere addition paradox

The following statements are all plausible:

- (P0) A is strictly better than B: A > B. Otherwise, in the original scenario, by parity of reasoning or consistency (scenario B^+ and C, C^+ etc.) one would have to deny that A is better than Z.
- (P1) A^+ is at least as good as $A: A^+ \ge A$. Justification: A^+ is not worse than (and hence at least as good as) A; the addition of lives worth living (the + people) cannot make a population worse.
- (P2) *B* is strictly better than A^+ : $B > A^+$. Justification: A^+ and *B* have the same size; the average welfare level in *B* is slightly above the average in A^+ , and the distribution is uniform across members. So *B* is better in regard to both average welfare (and thus also total welfare) and equality.

The relations \geq and > appearing in (P0)-(P2) apply to propositional formulas. It is usual to take $\varphi \geq \psi$ as a shorthand of $P(\varphi/\varphi \lor \psi)$, and $\varphi > \psi$ as a shorthand of $\varphi \geq \psi$ and $\psi \not\geq \varphi$. (Cf. Lewis (1973)). This is shown in Table 3.

Definiendum	Definiens	Reading
$\varphi \geq \psi$	$P(arphi ee arphi ee \psi)$	φ permitted, if $\varphi \lor \psi$
$\varphi > \psi$	$P(\varphi/\varphi \lor \psi) \land \bigcirc (\neg \psi/\varphi \lor \psi)$	φ permitted and ψ forbidden, if $\varphi \lor \psi$

Table 3.: Preference on formulas

Fig. 24 shows the encoding of P0-P2 in terms of obligation.

```
6
7 consts a::σ aplus::σ b::σ
8
9 (*the mere addition scenario*)
10
11 axiomatization where
12 (* A is striclty better than B*)
13 P0: "[(¬⊙<¬a|aVb>∧⊙<¬b|aVb>)]" and
14 (* Aplus is at least as good as A*)
15 P1: "[¬⊙<¬aplus|aVaplus>]" and
16 (* B is strictly better than Aplus*)
17 P2: "[(¬⊙<¬b|aplusVb> ∧ ⊙<¬aplus|aplusVb>)]"
18
```

Figure 24.: Encoding of the mere addition scenario

Fig. 25 shows some sample queries run on the scenario. On l. 19, the assumption of the transitivity of the betterness relation (on possible worlds) is assumed. *Sledgehammer* shows the inconsistency of (P0)-(P2). On l. 23, the assumption of transitivity is

not assumed. *Nitpick* shows the satisfiability of (P0)-(P2). The model generated by *Nitpick* is shown in Fig. 26.

```
19 (* Sledgehammer finds P0-P2 inconsistent given
20 transitivity of the <u>betterness</u> relation in the models*)
21
22 theorem TransIncons:
    assumes transitivity
23
    shows False
-24
25
    sledgehammer
    by (metis P0 P1 P2 assms)
-26
27
28
  (* Nitpick shows consistency in the absence of transitivity*)
29
30
  lemma true
31
32
    nitpick [satisfy, card i=3]
                                   (*model found*)
-33
    oops
```

Figure 25.: Sample queries on (P0)-(P2)

```
Nitpick found a model for card i = 3:

Constants:

(r) =

(\lambdax. _)

((i<sub>1</sub>, i<sub>1</sub>) := True, (i<sub>1</sub>, i<sub>2</sub>) := True,

(i<sub>1</sub>, i<sub>3</sub>) := False, (i<sub>2</sub>, i<sub>1</sub>) := False,

(i<sub>2</sub>, i<sub>2</sub>) := True, (i<sub>2</sub>, i<sub>3</sub>) := True,

(i<sub>3</sub>, i<sub>1</sub>) := True, (i<sub>3</sub>, i<sub>2</sub>) := True,

(i<sub>3</sub>, i<sub>3</sub>) := True)

a = (\lambdax. _)(i<sub>1</sub> := False, i<sub>2</sub> := False, i<sub>3</sub> := True)

aplus = (\lambdax. _)(i<sub>1</sub> := True, i<sub>2</sub> := False, i<sub>3</sub> := False)

b = (\lambdax. _)(i<sub>1</sub> := True, i<sub>2</sub> := False, i<sub>3</sub> := False)
```

Figure 26.: A non-transitive model satisfying (P0)-(P2)

```
43 (* Nitpick shows consistency if transitivity is weakened into <u>acyclicity</u> or quasi-transitivity*)
44
45 theorem Acyclcons:
    assumes loopfree
46
    shows true
47
    nitpick [show all,satisfy,card=3] (* model found for card i=3 *)
48
49
    oops
50
51 theorem Quasicons2:
52
    assumes Quasitransit
    shows true
-53
54
    nitpick [show_all,satisfy,card=4]
55
    oops
```

Figure 27.: A-cyclicity and quasi-transitivity



Figure 28.: Interval order

Nitpick is also able to confirm that the mere-addition paradox is avoided if transitivity ity is not rejected wholesale, but weakened into a-cyclicity or quasi-transitivity. This point has in general been overlooked in the literature. On the other hand, *Sledgehammer* can verify that this solution does not work for the interval order condition, which represents another candidate weakening of transitivity. The verifications are shown in Figs. 27 and 28.

6. Conclusion

Utilizing the LogiKEy methodology and framework we have developed mechanizations of extensions of Åqvist's preference-based system **E** for conditional obligation. We have illustrated the use of the resulting tool for (i) meta-logical studies and for (ii) object-level application studies in normative reasoning. Novel contributions, partly contributed by the automated reasoning tools in Isabelle/HOL, include the automated verification of the correspondence between semantic properties and modal axioms, and the formalization and mechanization of Parfit's argument for the repugnant conclusion. This one reveals the possibility of a take on the scenario usually under-appreciated in the literature, which consists in weakening transitivity rather than reject it wholesale. Future work includes the handling of the full equivalence between properties and formulas, the formalization of (and comparison with) other solutions to the repugnant conclusion, and the analysis of other variant paradoxes discussed in the literature.

Funding

The authors acknowledge TU Wien Bibliothek for financial support through its Open Access Funding Programme. Dr. X. Parent was funded in whole, or in part, by the Austrian Science Fund (FWF) [M3240 N, ANCoR project]. For the purpose of open access, the author has applied a CC BY public copyright licence to any Author Accepted Manuscript version arising from this submission.

References

Andrews, P. (2002). An Introduction to Mathematical Logic and Type Theory. Springer.

- Benzmüller, C. (2019). Universal (meta-)logical reasoning: Recent successes. Science of Computer Programming, 172:48–62.
- Benzmüller, C., Claus, M., and Sultana, N. (2015). Systematic verification of the modal logic cube in Isabelle/HOL. In Kaliszyk, C. and Paskevich, A., editors, *PxTP 2015*, volume 186, pages 27–41, Berlin, Germany. EPTCS.

- Benzmüller, C., Farjami, A., and Parent, X. (2019). Åqvist's dyadic deontic logic E in HOL. Journal of Applied Logics – IfCoLoG Journal of Logics and their Applications (Special Issue: Reasoning for Legal AI), 6(5):733–755.
- Benzmüller, C., Gabbay, D., Genovese, V., and Rispoli, D. (2012). Embedding and automating conditional logics in classical higher-order logic. Annals of Mathematics and Artificial Intelligence, 66(1-4):257–271.
- Benzmüller, C., Parent, X., and van der Torre, L. (2020). Designing normative theories for ethical and legal reasoning: Logikey framework, methodology, and tool support. Artificial Intelligence, 287:103348.
- Benzmüller, C. and Reiche, S. (2022). Automating public announcement logic with relativized common knowledge as a fragment of HOL in LogiKEy. *Journal of Logic and Computation*.
- Blanchette, J. C., Böhme, S., and Paulson, L. C. (2013). Extending Sledgehammer with SMT solvers. *Journal of Automated Reasoning*, 51(1):109–128.
- Blanchette, J. C., Kaliszyk, C., Paulson, L. C., and Urban, J. (2016). Hammering towards QED. Journal of Formalized Reasoning, 9(1):101–148.
- Blanchette, J. C. and Nipkow, T. (2010). Nitpick: A counterexample generator for higher-order logic based on a relational model finder. In Kaufmann, M. and Paulson, L. C., editors, *ITP* 2010, volume 6172 of *LNCS*, pages 131–146. Springer.

Chellas, B. (1975). Basic conditional logic. Journal of Philosophical Logic, 4(2):pp. 133–153.

- Chellas, B. (1980). Modal Logic. Cambridge University Press, Cambridge.
- Hansson, B. (1969). An analysis of some deontic logics. *Noûs*, 3(4):373–398.
- Hughes, G. E. and Cresswell, M. J. (1984). A companion to modal logic. Methuen, London.
- Lewis, D. (1973). Counterfactuals. Blackwell, Oxford.
- Luce, R. (1956). Semiorders and a theory of utility discrimination. *Econometrica*, 24:178–191.
- Parent, X. (2014). Maximality vs. optimality in dyadic deontic logic. Journal of Philosophical Logic, 43(6):1101–1128.
- Parent, X. (2015). Completeness of Aqvist's systems E and F. Review of Symbolic Logic, 8(1):164–177.
- Parent, X. (2021). Preference semantics for dyadic deontic logic: a survey of results. In Gabbay, D., Horty, J., Parent, X., van der Meyden, R., and van der Torre, L., editors, *Handbook* of Deontic Logic and Normative Systems, pages 1–70. College Publications, London. UK. Volume 2.
- Parent, X. (2022). On some weakened forms of transitivity in the logic of norms. In Giordanni, L. and Casini, G., editors, NMR 2022: 20th International Workshop on Non-Monotonic Reasoning. CEUR-WS. Extended abstract.
- Parent, X. and van der Torre, L. (2021). Introduction to Deontic Logic and Normative Systems. College Publications, London.
- Parfit, D. (1984). Reasons and Persons. Oxford University Press.
- Temkin, L. S. (1987). Intransitivity and the mere addition paradox. *Philosophy and Public Affairs*, 16(2):138–187.
- Van Benthem, J. (2001). Correspondence theory. In Gabbay, D. M. and Guenthner, F., editors, Handbook of Philosophical Logic, pages 325–408. Springer Netherlands, Dordrecht.