© 2014 The Authors and IOS Press.

This article is published online with Open Access by IOS Press and distributed under the terms

of the Creative Commons Attribution Non-Commercial License.

doi:10.3233/978-1-61499-419-0-93

Automating Gödel's Ontological Proof of God's Existence with Higher-order Automated Theorem Provers

Christoph Benzmüller¹ and Bruno Woltzenlogel Paleo²

Abstract. Kurt Gödel's ontological argument for God's existence has been formalized and automated on a computer with higher-order automated theorem provers. From Gödel's premises, the computer proved: necessarily, there exists God. On the other hand, the theorem provers have also confirmed prominent criticism on Gödel's ontological argument, and they found some new results about it.

The background theory of the work presented here offers a novel perspective towards a *computational theoretical philosophy*.

1 INTRODUCTION

Kurt Gödel proposed an argumentation formalism to prove the existence of God [23, 30]. Attempts to prove the existence (or non-existence) of God by means of abstract, ontological arguments are an old tradition in western philosophy. Before Gödel, several prominent philosophers, including St. Anselm of Canterbury, Descartes and Leibniz, have presented similar arguments. Moreover, there is an impressive body of recent and ongoing work (cf. [31, 19, 18] and the references therein). Ontological arguments, for or against the existence of God, illustrate well an essential aspect of metaphysics: some (necessary) facts for our existing world are deduced by purely a priori, analytical means from some abstract definitions and axioms.

What motivated Gödel as a logician was the question, whether it is possible to deduce the existence of God from a small number of foundational (but debatable) axioms and definitions, with a mathematically precise, formal argumentation chain in a well defined logic.

In theoretical philosophy, formal logical confrontations with such ontological arguments had been so far (mainly) limited to paper and pen. Up to now, the use of computers was prevented, because the logics of the available theorem proving systems were not expressive enough to formalize the abstract concepts adequately. Gödel's proof uses, for example, a complex higher-order modal logic (HOML) to handle concepts such as *possibility* and *necessity* and to support quantification over individuals and properties.

Current works [10, 9] of the first author and Paulson illustrate that many expressive logics, including quantified (multi-)modal logics, can be embedded into the classical higher-order logic (HOL), which can thus be seen as a universal logic [6]. For this universal logic, efficient automated theorem provers have been developed in recent years, and these systems were now employed in our work.

Gödel defines God (see Fig. 1) as a being who possesses all *positive* properties. He does not extensively discuss what positive properties are, but instead he states a few reasonable (but debatable) ax-

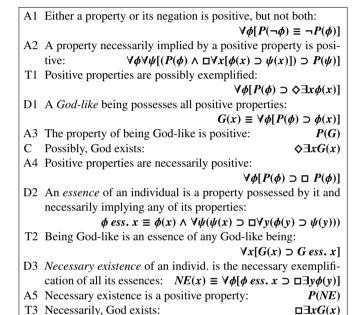


Figure 1. Scott's version of Gödel's ontological argument [30].

ioms that they should satisfy. Various slightly different versions of axioms and definitions have been considered by Gödel and by several philosophers who commented on his proof (cf. [31, 3, 2, 19, 1, 18]).

The overall idea of Gödel's proof is in the tradition of Anselm's argument, who defined God as some entity of which nothing greater can be conceived. Anselm argued that existence in the actual world would make such an assumed being even greater; hence, by definition God must exist. Gödel's ontological argument is clearly related to this reasoning pattern. However, it also tries to fix some fundamental weaknesses in Anselm's work. For example, Gödel explicitly proves that God's existence is possible, which has been a basic assumption of Anselm. Because of this, Anselm's argument has been criticized as incomplete by Leibniz. Leibniz instead claimed that the assumption should be derivable from the definition of God as a perfect being and from the notion of perfection. Gödel's proof addresses this critique, and it also addresses the critique of others, including Kant's objection that existence should not be treated as a predicate. On the other hand, Gödel's work still leaves room for criticism, in particular, his axioms are so strong that they imply modal collapse, that is, a situation where contingent truths and necessary truths coincide. More information on the philosophical debate on Gödel's proof is provided in [31].

¹ Freie Universität Berlin, Germany, email: c.benzmueller@fu-berlin.de; this author has been supported by the German National Research Foundation (DFG) under grants BE 2501/9-1 and BE 2501/11-1.

² Technical University Vienna, Austria, email: bruno@logic.at

We have analyzed Dana Scott's version of Gödel's proof [30] (cf. Fig. 1) for the first-time with an unprecedented degree of detail and formality with the help of higher-order automated theorem provers (HOL-ATPs).³ The following has been done (and in this order): (i) a detailed natural deduction proof; (ii) a formalization in TPTP THF syntax [33]; (iii) an automatic verification of the consistency of the axioms and definitions with Nitpick [16]; (iv) an automatic demonstration of the theorems with the provers LEO-II [11] and Satallax [17]; (v) a step-by-step formalization using the Coq proof assistant [15]; (vi) a formalization using the Isabelle proof assistant [26], where the theorems (and some additional lemmata) have been automated with the Isabelle tools Sledgehammer and Metis. Subsequently, we have studied additional consequences of Gödel's axioms, including modal collapse and monotheism, and we have investigated variations of the proof, for example, by switching from constant domain semantics to varying domain semantics.

In this paper we focus on the core aspect of our work related to AI: proof automation with HOL-ATPs (cf. aspects (ii)–(iv) above). The particular contributions of this paper are as follows: In Sec. 2 we present an elegant embedding of HOML [21, 25] in HOL [5, 7]. This background theory extends and adapts the work as presented in [9, 10]. In Sec. 3, we present details on the encoding of this embedding and of Gödel's argument in the concrete THF syntax [33] for HOL, and we report on the experiments we have conducted with HOL-ATPs. The main findings of these experiments are summarized in Sec. 4. Related and future work is addressed in Sec. 5, and the paper is concluded in Sec. 6. None of the above aspects have been addressed (at least not in depth) in any of our other existing (short and partly non-reviewed) publications on the subject [12, 13, 14, 34].

2 THEORY FRAMEWORK

An embedding of quantified modal logic (with first-order and propositional quantifiers) in HOL has been presented in [10]. The theory below extends this work: quantifiers for all types are now supported, and nested uninterpreted predicate and function symbols of arbitrary types are allowed as opposed to allowing top-level uninterpreted predicate symbols over individual variables only.

2.1 Higher-order modal logic

A notion of HOML is introduced that extends HOL with a modal operator \square . An appropriate notion of semantics for HOML is obtained by adapting Henkin semantics for HOL (cf. [24] and [21]). The presentation in this section is adapted from [25] and [5].

Def. 1 The set T of simple types is freely generated from the set of basic types $\{o, \mu\}$ (o stands for Booleans and μ for individuals) using the function type constructor \rightarrow . We may avoid parentheses, and $\alpha \rightarrow \alpha \rightarrow \alpha$ then stands for $(\alpha \rightarrow (\alpha \rightarrow \alpha))$, that is, function types associate to the right.

Def. 2 The grammar for HOML is:

$$\begin{split} s,t & ::= & p_{\alpha} \mid X_{\alpha} \mid (\lambda X_{\alpha^{\blacksquare}} s_{\beta})_{\alpha \to \beta} \mid (s_{\alpha \to \beta} t_{\alpha})_{\beta} \mid (\neg_{o \to o} s_{o})_{o} \mid \\ & ((\vee_{o \to o \to o} s_{o}) t_{o})_{o} \mid (\forall_{(\alpha \to o) \to o} (\lambda X_{\alpha^{\blacksquare}} s_{o}))_{o} \mid (\Box_{o \to o} s_{o})_{o} \end{split}$$

where $\alpha, \beta \in T$. p_{α} denotes typed constants and X_{α} typed variables (distinct from p_{α}). Complex typed terms are constructed via abstraction and application. The type of each term is given as a subscript. Terms s_o of type o are called formulas. The logical connectives of choice are $\neg_{o \to o}$, $\lor_{o \to o \to o}$, $\lor_{(\alpha \to o) \to o}$ (for $\alpha \in T$), and $\Box_{o \to o}$. Type subscripts may be dropped if irrelevant or obvious. Similarly, parentheses may be avoided. Binder notation $\forall X_{\alpha}$ so is used as shorthand for $\forall_{(\alpha \to o) \to o}(\lambda X_{\alpha}$ so, and infix notation $s \lor t$ is employed instead of $((\lor s) t)$. From the above connectives, other logical connectives, such as \top , \bot , \land , \supset , \equiv , \exists , and \diamondsuit , can be defined in the usual way.

Def. 3 Substitution of a term A_{α} for a variable X_{α} in a term B_{β} is denoted by [A/X]B. Since we consider α -conversion implicitly, we assume the bound variables of B avoid variable capture.

Def. 4 Two common relations on terms are given by β -reduction and η -reduction. A β -redex has the form $(\lambda X_\bullet s)t$ and β -reduces to [t/X]s. An η -redex has the form $(\lambda X_\bullet sX)$ where variable X is not free in s; it η -reduces to s. We write $s=_{\beta}t$ to mean s can be converted to t by a series of β -reductions and expansions. Similarly, $s=_{\beta\eta}t$ means s can be converted to t using both β and η . For each $s_{\alpha} \in HOML$ there is a unique β -normal form and a unique $\beta\eta$ -normal form.

Def. 5 A frame D is a collection $\{D_{\alpha}\}_{{\alpha} \in T}$ of nonempty sets D_{α} , such that $D_{\sigma} = \{T, F\}$ (for truth and falsehood). The $D_{{\alpha} \to {\beta}}$ are collections of functions mapping D_{σ} into $D_{{\beta}}$.

Def. 6 A variable assignment g maps variables X_{α} to elements in D_{α} . g[d/W] denotes the assignment that is identical to g, except for variable W, which is now mapped to d.

Def. 7 A model for HOML is a quadruple $M = \langle W, R, D, \{I_w\}_{w \in W} \rangle$, where W is a set of worlds (or states), R is an accessibility relation between the worlds in W, D is a frame, and for each $w \in W$, $\{I_w\}_{w \in W}$ is a family of typed interpretation functions mapping constant symbols p_α to appropriate elements of D_α , called the denotation of p_α in world w (the logical connectives \neg , \lor , \forall , and \Box are always given the standard denotations, see below). Moreover, it is assumed that the domains $D_{\alpha \to \alpha \to o}$ contain the respective identity relations on objects of type α (to overcome the extensionality issue discussed in [4]).

Def. 8 The value $||s_{\alpha}||^{M,g,w}$ of a HOML term s_{α} on a model $M = \langle W, R, D, \{I_w\}_{w \in W} \rangle$ in a world $w \in W$ under variable assignment g is an element $d \in D_{\alpha}$ defined in the following way:

- 1. $||p_{\alpha}||^{M,g,w} = I_w(p_{\alpha})$ and $||X_{\alpha}||^{M,g,w} = g(X_{\alpha})$
- 2. $||(s_{\alpha \to \beta} t_{\alpha})_{\beta}||^{M,g,w} = ||s_{\alpha \to \beta}||^{M,g,w} (||t_{\alpha}||^{M,g,w})$
- 3. $\|(\lambda X_{\alpha^{\bullet}} s_{\beta})_{\alpha \to \beta}\|^{M,g,w} = the function f from D_{\alpha} to D_{\beta} such that <math>f(d) = \|s_{\beta}\|^{M,g[d/X_{\alpha}],w}$ for all $d \in D_{\alpha}$
- 4. $\|(\neg_{o \to o} s_o)_o\|^{M,g,w} = T \text{ iff } \|s_o\|^{M,g,w} = F$
- 5. $\|((\vee_{o \to o \to o} s_o) t_o)_o\|^{M,g,w} = T \text{ iff } \|s_o\|^{M,g,w} = T \text{ or } \|t_o\|^{M,g,w} = T$
- 6. $\|(\forall_{(\alpha \to o) \to o}(\lambda X_{\alpha^{\bullet}} s_o))_o\|^{M,g,w} = T \text{ iff for all } d \in D_{\alpha} \text{ we have } \|s_o\|^{M,g[d/X_{\alpha}],w} = T$
- 7. $\|(\Box_{o \to o} s_o)_o\|^{M,g,w} = T$ iff for all $v \in W$ with wRv we have $\|s_o\|^{M,g,v} = T$

Def. 9 A model $M = \langle W, R, D, \{I_w\}_{w \in W} \rangle$ is called a standard model iff for all $\alpha, \beta \in T$ we have $D_{\alpha \to \beta} = \{f \mid f : D_\alpha \longrightarrow D_\beta\}$. In a Henkin model function spaces are not necessarily full. Instead it is only required that $D_{\alpha \to \beta} \subseteq \{f \mid f : D_\alpha \longrightarrow D_\beta\}$ (for all $\alpha, \beta \in T$) and that the valuation function $\|\cdot\|^{Mg,w}$ from above is total (i.e., every term denotes). Any standard model is obviously also a Henkin model. We consider Henkin models in the remainder.

³ All sources of our formalization are publicly available at https://github.com/FormalTheology/GoedelGod. Our work has attracted major public interest, and leading media institutions worldwide have reported on it; some exemplary links to respective media reports and interviews are available at the above URL (see 'Press' subfolder).

Def. 10 A formula s_o is true in model M for world w under assignment g iff $||s_o||^{M,g,w} = T$; this is also denoted as $M, g, w \models s_o$. A formula s_o is called valid in M iff M, g, $w \models s_o$ for all $w \in W$ and all assignments g. Finally, a formula so is called valid, which we denote $by \models s_o$, iff s_o is valid for all M. Moreover, we write $\Gamma \models \Delta$ (for sets of formulas Γ and Δ) iff there is a model $M = \langle W, R, D, \{I_w\}_{w \in W} \rangle$, an assignment g, and a world $w \in W$, such that $M, g, w \models s_o$ for all $s_o \in \Gamma$ and $M, g, w \models t_o$ for at least one $t_o \in \Delta$.

The above definitions introduce higher-order modal logic K. In order to obtain logics KB and S5 respective conditions on accessibility relation R are postulated: R is a symmetric relation in logic KB, and it is an equivalence relation in logic S5. If these restriction apply, we use the notations \models^{KB} and \models^{S5} . Gödel's argument has been developed and studied in the context of logic S5 (and logic S5 has subsequently been criticized). However, the HOL-ATPs discovered (cf. Sec. 4) that logic KB is sufficient.

An important issue for quantified modal logics is whether constant domain or varying domain semantics is considered. The theory above introduces constant domains. Terms (other than those of Boolean type) are modeled as rigid, that is, their denotation is fixed for all worlds. An adaptation to varying or cumulative domains is straightforward (cf. [20]). Moreover, non-rigid terms could be modeled; that is, terms whose denotation may switch from world to world. The respective assumptions of Gödel are not obvious to us.

Classical higher-order logic 2.2

HOL is easily obtained from HOML by removing the modal operator \Box from the grammar, and by dropping the set of possible worlds W and the accessibility relation R from the definition of a model. Nevertheless, we explicitly state the most relevant definitions for the particular notion of HOL as employed in this paper. One reason is that we do want to carefully distinguish the HOL and HOML languages in the remainder (we use boldface fonts for HOL and standard fonts for HOML). There is also a subtle, but harmless, difference in the HOL language as employed here in comparison to the standard presentation: here three base types are employed, whereas usually only two base types are considered. The third base type plays a crucial role in our embedding of HOML in HOL.

Def. 11 The set **T** of simple types freely generated from a set of basic types $\{o, \mu, \iota\}$ using the function type constructor \rightarrow . o is the type of Booleans, μ is the type of individuals, and type ι is employed as the type of possible worlds below. As before we may avoid parentheses.

Def. 12 The grammar for higher-order logic HOL is:

$$s, t \quad ::= \quad p_{\alpha} \mid X_{\alpha} \mid (\lambda X_{\alpha^{\blacksquare}} s_{\beta})_{\alpha \to \beta} \mid (s_{\alpha \to \beta} t_{\alpha})_{\beta} \mid \neg_{o \to o} s_{o} \mid ((\vee_{o \to o \to o} s_{o}) t_{o}) \mid \forall_{(\alpha \to o) \to o} (\lambda X_{\alpha^{\blacksquare}} s_{o})$$

where $\alpha, \beta \in T$. The text from Def. 2 analogously applies, except that we do not consider the modal connectives \square and \diamondsuit .

The definitions for substitution (Def. 3), β - and η -reduction (Def. 4), frame (Def. 5), and assignment (Def. 6) remain unchanged.

Def. 13 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_{α} to appropriate elements of D_{α} , called the denotation of p_{α} (the logical connectives \neg , \lor , and \forall are always given the standard denotations, see below). Moreover, we assume that the domains $D_{\alpha \to \alpha \to \rho}$ contain the respective identity relations.

Def. 14 The value $||s_{\alpha}||^{M,g}$ of a HOL term s_{α} on a model $M = \langle D, I \rangle$ under assignment g is an element $d \in D_{\alpha}$ defined in the following

- 1. $||p_{\alpha}||^{M,g} = I(p_{\alpha})$ and $||X_{\alpha}||^{M,g} = g(X_{\alpha})$
- 2. $\|(s_{\alpha \to \beta} t_{\alpha})_{\beta}\|^{M,g} = \|s_{\alpha \to \beta}\|^{M,g} (\|t_{\alpha}\|^{M,g})$ 3. $\|(\lambda X_{\alpha} \cdot s_{\beta})_{\alpha \to \beta}\|^{M,g} = \text{the function } f \text{ from } D_{\alpha} \text{ to } D_{\beta} \text{ such that }$ $f(d) = ||s_{\beta}||^{M,g[d/X_{\alpha}]}$ for all $d \in D_{\alpha}$
- $\|(\neg_{o\rightarrow o} s_o)_o\|^{M,g} = T \text{ iff } \|s_o\|^{M,g} = F$
- 5. $\|((\mathsf{V}_{o\to o\to o} s_o) t_o)_o\|^{M,g} = T \text{ iff } \|s_o\|^{M,g} = T \text{ or } \|t_o\|^{M,g} = T$
- 6. $\|(\forall_{(\alpha \to o) \to o}(\lambda X_{\alpha \bullet} s_o))_o\|^{M,g} = T \text{ iff for all } d \in D_\alpha \text{ we have}$

The definition for standard and Henkin models (Def. 9), and for truth in a model, validity, etc. (Def. 10) are adapted in the obvious way, and we use the notation $M, g \models s_o, \models s_o$, and $\Gamma \models \Delta$. As for HOML, we assume Henkin semantics in the remainder.

2.3 **HOML** as a fragment of **HOL**

The encoding of HOML in HOL is simple: we identify HOML formulas of type o with certain HOL formulas of type $\iota \to o$. The HOL type $\iota \to o$ is abbreviated as σ in the remainder. More generally, we define for each HOML type $\alpha \in T$ the associated raised HOL type $\lceil \alpha \rceil$ as follows: $\lceil \mu \rceil = \mu$, $\lceil \sigma \rceil = \sigma = \iota \rightarrow \sigma$, and $\lceil \alpha \rightarrow \beta \rceil = \lceil \alpha \rceil \rightarrow \lceil \beta \rceil$. Hence, all HOML terms are rigid, except for those of type o.

Def. 15 HOML terms s_{α} are associated with type-raised HOL terms $\lceil s_{\alpha} \rceil$ in the following way:

$$\lceil p_{\alpha} \rceil = p_{\lceil \alpha \rceil}$$

$$\lceil X_{\alpha} \rceil = X_{\lceil \alpha \rceil}$$

$$\lceil (s_{\alpha \to \beta} t_{\alpha}) \rceil = (\lceil s_{\alpha \to \beta} \rceil \lceil t_{\alpha} \rceil)$$

$$\lceil (\lambda X_{\alpha} \cdot s_{\beta}) \rceil = (\lambda \lceil X_{\alpha} \rceil \cdot \lceil s_{\beta} \rceil)$$

$$\lceil (\neg s_{\alpha \to \sigma} s_{\alpha}) \rceil = (\dot{\neg} s_{\alpha \to \sigma} \lceil s_{\alpha} \rceil)$$

$$\lceil ((\forall_{\alpha \to \sigma \to \sigma} s_{\alpha}) t_{\alpha}) \rceil = ((\dot{\lor} s_{\alpha \to \sigma \to \sigma} \lceil s_{\alpha} \rceil) \lceil t_{\alpha} \rceil)$$

$$\lceil ((\forall_{(\alpha \to \sigma) \to \sigma} (\lambda X_{\alpha} \cdot s_{\beta}) \rceil = (\dot{\lor} s_{\alpha \to \sigma} s_{\alpha} \rceil)$$

$$\lceil ((\exists_{\alpha \to \sigma} s_{\alpha}) \rceil = (\dot{\Box} s_{\alpha \to \sigma} \lceil s_{\alpha} \rceil)$$

 $\dot{\neg}$, $\dot{\lor}$, $\dot{\lor}$, and $\dot{\Box}$ are the type-raised modal HOL connectives associated with the corresponding modal HOML connectives. They are defined as follows (where $\mathbf{r}_{t\to t\to 0}$ is a new constant symbol in HOL associated with the accessibility relation R of HOML):

$$\dot{\neg}_{\sigma \to \sigma} = \lambda s_{\sigma^{\blacksquare}} \lambda W_{\iota^{\blacksquare}} \neg (s \ W)$$

$$\dot{\lor}_{\sigma \to \sigma \to \sigma} = \lambda s_{\sigma^{\blacksquare}} \lambda t_{\sigma^{\blacksquare}} \lambda W_{\iota^{\blacksquare}} s \ W \lor t \ W$$

$$\dot{\lor}_{(\alpha \to \sigma) \to \sigma} = \lambda s_{\alpha \to \sigma^{\blacksquare}} \lambda W_{\iota^{\blacksquare}} \lor X_{\alpha^{\blacksquare}} s \ X \ W$$

$$\dot{\Box}_{\sigma \to \sigma} = \lambda s_{\sigma^{\blacksquare}} \lambda W_{\iota^{\blacksquare}} \lor V_{\iota^{\blacksquare}} \neg (r_{t \to t \to \sigma} \ W \ V) \lor s \ V$$

As before, we write $\dot{\forall} X_{\alpha} \cdot s_{\sigma}$ as shorthand for $\dot{\forall}_{(\alpha \to \sigma) \to \sigma} (\lambda X_{\alpha} \cdot s_{\sigma})$. Further operators, such as $\dot{\top}$, $\dot{\bot}$, $\dot{\land}$, $\dot{\supset}$, $\dot{\equiv}$, $\dot{\diamondsuit}$, and $\dot{\exists}$ ($\dot{\exists}X_{\alpha}$ s_{\sigma} is used as shorthand for $\dot{\exists}_{(\alpha \to \sigma) \to \sigma}(\lambda X_{\alpha \bullet} s_{\sigma})$) can now be easily defined.⁴ The above equations can be treated as abbreviations in HOL theorem provers. Alternatively, they can be stated as axioms where = is either Leibniz equality or primitive equality (if additionally provided in the HOL grammar, as is the case for most modern HOL provers).

⁴ We could introduce further modal operators, such as the difference modality **D**, the global modality **E**, nominals with !, and the @ operator (cf. [10]).

As a consequence of the above embedding we can express HOML proof problems elegantly in the type-raised syntax of HOL. Using rewriting or definition expanding, we can reduce these representations to corresponding statements containing only the basic HOL connectives $\neg_{\sigma \to \sigma}$, $\bigvee_{\sigma \to \sigma \to \sigma}$, and $\bigvee_{(\alpha \to \sigma) \to \sigma}$.

Ex. 1 The HOML formula $\Box \exists P_{\mu \to \sigma} P a_{\mu}$ is associated with the type raised HOL formula $\dot{\Box} \dot{\exists} P_{\mu \to \sigma} P a_{\mu}$, which rewrites into the following $\beta \eta$ -normal HOL term of type σ

$$\lambda W_{\iota} \forall V_{\iota} \neg (r W V) \lor \neg \forall P_{\mu \to \sigma} \neg (P a_{\mu} V)$$

Next, we define validity of type-raised modal HOL propositions s_{σ} in the obvious way: s_{σ} is valid iff for all possible worlds w_t we have $w_t \in s_{\sigma}$, that is, iff (s_{σ}, w_t) holds.

Def. 16 Validity is modeled as an abbreviation for the following λ -term: $valid = \lambda s_{\iota \to o^{\bullet}} \forall W_{\iota^{\bullet}} s W$ (alternatively, we could define validity simply as $\forall_{(\iota \to o) \to o}$). Instead of $valid s_{\sigma}$ we also use the notation $[s_{\sigma}]$.

Ex. 2 We analyze whether the type-raised modal HOL formula $\dot{\Box} \dot{\exists} P_{\mu \to \sigma^*}(P \ a_{\mu})$ is valid or not. For this, we formalize the HOL proof problem $[\dot{\Box} \dot{\exists} P_{\mu \to \sigma^*}(P \ a_{\mu})]$, which expands into $\forall W_{\iota^*} \forall V_{\iota^*} \neg (r \ W \ V) \lor \neg \forall P_{\mu \to \sigma^*} \neg (P \ a_{\mu} \ V)$. It is easy to check that this term is valid in Henkin semantics: put $P = \lambda X_{\mu^*} \lambda Y_{\iota^*} \neg T$.

Theorem 1 (Soundness and Completeness) For all HOML formulas s_o we have:

$$\models s_o \quad iff \quad \models [\lceil s_o \rceil]$$

Proof sketch: The proof adapts the ideas presented in [10]. By contraposition it is sufficient to show $otin s_o$ iff $otin [[s_o]]^{M,g,w}$ (for some HOML model M, assignment g, and w) iff $otin [[w]^{M,g,w}] = [w]^{M,g}$ (for some HOL model M and assignment g) iff $otin [[s_o]] = [w]^{M,g[w/W]}$ (for some M, g, and w). We easily get the proof by choosing the obvious correspondences between D and D, W and D_t , U and U, U and U and U.

From Theorem 1 we get the following corollaries:

$$\models^{KB} s_o$$
 iff (symmetric $r_{t \to t \to 0}$) $\models [[s_o]]$
 $\models^{SS} s_o$ iff (equiv-rel $r_{t \to t \to 0}$) $\models [[s_o]]$

where symmetric and equiv-rel are defined in an obvious way.

Constant domain quantification is addressed above. Techniques for handling varying domain and cumulative domain quantification in the embedding of first-order modal logics in HOL have been outlined in [8]. These techniques, which have also been adapted for the theory above, cannot be presented here for space limitations.

Note that also non-rigid terms can easily be modeled by type-raising. For example, a non-rigid HOML constant symbol kingOfFrance $_{\mu}$ would be mapped to a type-raised (and thus world-depended) HOL constant symbol **kingOfFrance** $_{\iota \to \mu}$.

3 EXPERIMENTS

The above embedding has been encoded in the concrete THF0 syntax [33] for HOL; cf. the files Quantified_K/_KB/_S5.ax⁵ available

at https://github.com/FormalTheology/GoedelGod/tree/master/Formalizations/THF (all files mentioned below are provided under this URL). The definition for quantifier $\dot{\mathbf{Y}}_{((\mu \to \sigma) \to \sigma) \to \sigma}$, for example, is given as 6

Subsequently the axioms, definitions, and theorems from Fig. 1 and some further, related problems have been encoded in THF0. Then the THF0 compliant HOL-ATPs LEO-II [11], Satallax [17], and Nitpick [16] have been employed to automate the proof problems. LEO-II, which internally cooperates with the first-order prover E [29], was used exclusively in the initial round of experiments, that is, it was the first prover to automate Gödel's ontological argument.

Theorem T1 from Fig. 1, for example, is formalized as

This encodes the HOL formula

$$[\dot{\forall}\phi_{u\to\sigma} p_{(u\to\sigma)\to\sigma}\phi \dot{\supset} \dot{\Diamond}\dot{\exists}X_{\mu}\phi X]$$

v in the THF0 encoding stands for valid and p corresponds to the uppercase P, for 'positive', from Fig. 1. The respective encodings and the results of a series of recent experiments with LEO-II (version 1.6.2), Satallax (version 2.7), and Nitpick (version 2013) are provided in Fig. 2. The first row marked with T1, for example, shows that theorem T1 follows from axioms A2 and A1 (where only the ⊃-direction is needed); LEO-II and Satallax confirm this in 0.1 second. The experiments have been carried out w.r.t. the logics K and/or KB, and w.r.t. constant (const) and varying (vary) domain semantics for the domains of individuals. The exact dependencies (available axioms and definitions) are displayed for each single problem. The results of the prover calls are given in seconds. '-' means timeout. 'THM', 'CSA', 'SAT', and 'UNS' are the reported result statuses; they stand for 'Theorem', 'CounterSatisfiable', 'Satisfiable', and 'Unsatisfiable', respectively. The experiments can be easily reproduced: all relevant files have been provided at the above URL. For example, the two THF0 problem files associated with the first table row for T1 are T1_K_const_min.p and T1_K_vary_min.p, and those associated with the second row for T1 are T1_K_const_max.p and T1_K_vary_max.p, respectively. Moreover, a simple shell script call_tptp.sh is provided, which can be used to make remote calls to LEO-II, Satallax, and Nitpick installed at Sutcliffe's SystemOnTPTP infrastructure [32] at the University of Miami. The experiments used standard 2.80GHz computers with 1GB memory remotely located in Miami.

⁵ The formalization in these files slightly varies from the above theory w.r.t. technical details. For example, a generic \Box -operator is introduced that can be instantiated for different accessibility relations as e.g. required for multimodal logic applications (cf. [10]). Moreover, since THF0 does not support polymorphism, copies of the $\dot{\mathbf{V}}_{(\alpha \to \sigma) \to \sigma}$ and $\dot{\mathbf{J}}_{(\alpha \to \sigma) \to \sigma}$ connectives are provided only for the selected types $(\mu \to \sigma) \to \sigma$ and $((\mu \to \sigma) \to \sigma) \to \sigma$ as precisely required in Gödels's proof. The Isabelle version [13] and the Coq version of the encoding instead provide respective polymorphic definitions.

^{6 \$}i, \$o, and mu represent the HOL base types i, o, and μ. \$i>\$o encodes a function (predicate) type. Function application is represented by @, and for universal quantification, existential quantification and λ-abstraction the symbols!, ? and ^ are employed. ¬, V, Λ, and ⊃ are written as ~, |, &, and =>, respectively. The type-raised modal connectives are called mforall_*, mexists_*, mnot, mor, mand, mimplies, etc.

	HOL encoding	dependencies	logic	status	LEO-II const/vary	Satallax const/vary	Nitpick const/vary
A1	$[\dot{\forall} \phi_{\mu \to \sigma^{\blacksquare}} p_{(\mu \to \sigma) \to \sigma} (\lambda X_{\mu^{\blacksquare}} \dot{\neg} (\phi X)) \stackrel{.}{=} \dot{\neg} (p\phi)]$						
A2	$[\dot{\forall}\phi_{\mu o\sigma},\dot{\forall}\dot{\psi}_{\mu o\sigma},(p_{(\mu o\sigma) o\sigma}\phi\dot{\wedge}\dot{\Box}\dot{\forall}X_{\mu},(\phi X)]$	$(\downarrow \psi X)) \supset p \psi]$					
T1	$[\dot{\forall} \phi_{\mu \to \sigma} p_{(\mu \to \sigma) \to \sigma} \phi \supset \dot{\Diamond} \exists X_{\mu} \phi X]$	$A1(\supset), A2$	K	THM	0.1/0.1	0.0/0.0	—/—
	,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	A1, A2	K	THM	0.1/0.1	0.0/5.2	/
D1	$g_{\mu o \sigma} = \lambda X_{\mu} \dot{\forall} \phi_{\mu o \sigma} p_{(\mu o \sigma) o \sigma} \phi \dot{\supset} \phi X$						
A3	$[p_{(\mu o\sigma) o\sigma}g_{\mu o\sigma}]$						
C	$[\diamondsuit \exists X_{\mu^{\bullet}} g_{\mu o \sigma} X]$	T1, D1, A3	K	THM	0.0/0.0	0.0/0.0	—/— —/—
		A1, A2, D1, A3	K	THM	0.0/0.0	5.2/31.3	—/—
A4	$[\dot{\forall} \phi_{\mu o \sigma^{\blacksquare}} p_{(\mu o \sigma) o \sigma} \phi \ \dot{\supset} \ \dot{\Box} p \phi]$						
D2	$\operatorname{ess}_{(\mu \to \sigma) \to \mu \to \sigma} = \lambda \phi_{\mu \to \sigma} \lambda X_{\mu} \phi X \dot{\wedge} \dot{\forall} \psi_{\mu \to \sigma}$	$(\psi X \supset \Box \dot{\forall} Y_{\mu} \ (\phi Y \supset \psi Y))$					
T2	$[\forall X_{\mu} \ g_{\mu \to \sigma} X \supset (\operatorname{ess}_{(\mu \to \sigma) \to \mu \to \sigma} gX)]$	A1, D1, A4, D2	K	THM	19.1/18.3	0.0/0.0	—/— —/—
	, -, , , -	A1, A2, D1, A3, A4, D2	K	THM	12.9/14.0	0.0/0.0	—/—
D3	$NE_{\mu \to \sigma} = \lambda X_{\mu} \dot{\forall} \phi_{\mu \to \sigma} (ess \phi X \dot{\supset} \dot{\Box} \dot{\exists} Y_{\mu} \phi Y)$						
A5	$[p_{(\mu o \sigma) o \sigma} \mathrm{NE}_{\mu o \sigma}]$						
T3	$[\dot{\Box}\dot{\exists}X_{\mu}$, $g_{\mu o\sigma}X]$	D1, C, T2, D3, A5	K	CSA	—/—	—/—	3.8/6.2
		A1, A2, D1, A3, A4, D2, D3, A5	K	CSA	—/—	—/—	8.2/7.5
		D1, C, T2, D3, A5	KB	THM	0.0/0.1	0.1/5.3	—/—
		A1, A2, D1, A3, A4, D2, D3, A5	KB	THM	—/—	—/—	—/—
MC	$[s_{\sigma} \supset \dot{\Box} s_{\sigma}]$	D2, T2, T3	KB	THM	17.9/—	3.3/3.2	—/—
	[-00]	A1, A2, D1, A3, A4, D2, D3, A5	KB	THM	_/_	_/_	<u>_</u> /
FG	$[\dot{\forall}\phi_{\mu o\sigma}]\dot{\forall}X_{\mu}$ $(g_{\mu o\sigma}X\dot{\supset}(\dot{\lnot}(p_{(\mu o\sigma) o\sigma}\phi)\dot{\supset}\cdot$	$\dot{\neg}(\phi X)))$] A1, D1	KB	THM	16.5/—	0.0/0.0	_/_
	$\mu \vee \mu \rightarrow 0 \qquad \mu \vee \mu \rightarrow 0 \qquad - \vee \psi \rightarrow 0 \rightarrow$	A1, A2, D1, A3, A4, D2, D3, A5	KB	THM	12.8/15.1	0.0/5.4	_/_
MT	$[\dot{\forall} X_{\mu} \ \dot{\forall} Y_{\mu} \ (g_{\mu \to \sigma} X \ \dot{\supset} \ (g_{\mu \to \sigma} Y \ \dot{\supset} \ X \ \dot{=} \ Y))]$	D1.FG	KB	THM	_/_	0.0/3.3	_/_
	7/1	A1, A2, D1, A3, A4, D2, D3, A5	KB	THM	_/	_/	_/_
CO	∅ (no goal, check for consistency)	A1, A2, D1, A3, A4, D2, D3, A5	KB	SAT	—/—	—/—	7.3/7.4
D2'	$\operatorname{ess}_{(\mu \to \sigma) \to \mu \to \sigma} = \lambda \phi_{\mu \to \sigma} \lambda X_{\mu} \dot{\forall} \psi_{\mu \to \sigma} (\psi X)$	$ \dot{\Box} \dot{\forall} Y_{\mu} (\phi Y \dot{\supset} \psi Y)) $			•	-	•
CO'	\emptyset (no goal, check for consistency)	$A1(\supset), A2, D2', D3, A5$	KB	UNS	7.5/7.8	—/—	—/—
		A1, A2, D1, A3, A4, D2', D3, A5	KB	UNS	/	_/_	<u>_</u> /

Figure 2. HOL encodings and experiment results for Scott's version of Gödel's ontological argument from Fig. 1.

4 MAIN FINDINGS

Several interesting and partly novel findings have been contributed by the HOL-ATPs, including:

- 1. The axioms and definitions from Fig. 1 are consistent (cf. CO in Fig. 2)
- 2. Logic K is sufficient for proving T1, C and T2.
- 3. For proving the final theorem T3, logic KB is sufficient (and for K a countermodel is reported). This is highly relevant since several philosophers have criticized Gödel's argument for the use of logic S5. This criticism is thus provably pointless.
- 4. Only for T3 the HOL-ATPs still fail to produce a proof directly from the axioms; thus, T3 remains an interesting benchmark problem; T1, C, and T2 are rather trivial for HOL-ATPs.
- 5. Gödel's original version of the proof [23], which omits conjunct $\phi(x)$ in the definition of *essence* (cf. D2'), seems inconsistent (cf. the failed consistency check for CO' in Fig. 2). As far as we are aware of, this is a new result.
- 6. Gödel's axioms imply what is called the modal collapse (cf. MC in Fig. 2) φ ⊃ □φ, that is, contingent truth implies necessary truth (which can even be interpreted as an argument against free will; cf. [31]). MC is probably the most fundamental criticism put forward against Gödel's argument.
- 7. For proving T1, only the >-direction of A1 is needed. How-

- ever, the ⊂-direction of A1 is required for proving T2. Some philosophers (e.g. [3]) try to avoid MC by eluding/replacing the ⊃-direction of A1.
- 8. Gödel's axioms imply a 'flawless God', that is, an entity that can only have 'positive' properties (cf. FG in Fig. 2). However, a comment by Gödel in [23] explains that 'positive' is to be interpreted in a moral aesthetic sense only.
- 9. Another implication of Gödel's axioms is monotheism (see MT in Fig. 2). MT can easily be proved by Satallax from FG and D1. It remains non-trivial to prove it directly from Gödel's axioms.
- 10. All of the above findings hold for both constant domain semantics and varying domain semantics (for the domain of individuals).

The above findings, in particular (10), well illustrate that the theory framework from Sec. 2 has a great potential towards a flexible support system for *computational theoretical philosophy*. In fact, Gödel's ontological argument has been verified and even automated not only for one particular setting of logic parameters, but these logic parameters have been varied and the validity of the argument has been reconfirmed (or falsified, cf. D2' and CO') for the modified setting. Moreover, our framework is not restricted to a particular theorem proving system, but has been fruitfully employed with some of the most prominent automated and interactive theorem provers available to date.

5 RELATED AND FUTURE WORK

We are pioneering the computer-supported automation of modern versions of the ontological argument. There are two related papers [27, 28]. Both focus on the comparably simpler argument by Anselm. [27] encodes (a variant) of Anselm's argument in first-order logic and employs the theorem prover PROVER9 in experiments; this work has been criticized in [22]. The work in [28], which has evolved in parallel to ours, interactively verifies Anselm's argument in the higher-order proof assistant PVS. Note in particular, that both formalizations do not achieve the close correspondence between the original formulations and the formal encodings that can be found in our approach.

A particular strength of our universal logic framework is that it can be easily adapted for logic variations and even supports flexible combinations of logics (cf. [6]). In ongoing and future work we will therefore investigate further logic parameters for Gödel's argument, including varying domains at higher types and non-rigid terms. We plan to make the entire landscape of results available to the interested communities. This is relevant, since philosophers are sometimes imprecise about the very details of the logics they employ.

6 CONCLUSION

While computers can now calculate, play games, translate, plan, learn and classify data much better than we humans do, tasks involving philosophical and theological inquiries have remained mostly untouched by our technological progress up to now. Due to the abstract and sophisticated types of reasoning they require, they can be considered a challenging frontier for automated reasoning.

We accepted this challenge and decided to tackle, with automated reasoning techniques, a philosophical problem that is almost 1000 years old: the ontological argument for God's existence, firstly proposed by St. Anselm of Canterbury and greatly improved by Descartes, Leibniz, Gödel and many others throughout the centuries. So far, there was no AI system capable of dealing with such complex problems. We created a prototypical infrastructure extending widely used systems such as LEO-II, Satallax, and Nitpick (and Isabelle and Coq) to allow them to cope with modalities; and using the extended systems we were able to automatically reconstruct and verify Gödel's argument, as well as discover new facts and confirm controversial claims about it. This is a landmark result, with media repercussion in a global scale, and yet it is only a glimpse of what can be achieved by combining computer science, philosophy and theology.

Our work, in this sense, offers new perspectives for a computational theoretical philosophy. The critical discussion of the underlying concepts, definitions and axioms remains a human responsibility, but the computer can assist in building and checking rigorously correct logical arguments. In case of logico-philosophical disputes, the computer can check the disputing arguments and partially fulfill Leibniz' dictum: Calculemus — Let us calculate!

ACKNOWLEDGEMENTS

We thank Alexander Steen, Max Wisniewski, and the anonymous reviewers for their comments and suggestions.

REFERENCES

- R.M. Adams, 'Introductory note to *1970', in Kurt Gödel: Collected Works Vol. 3: Unpubl. Essays and Letters, Oxford Univ. Press, (1995).
- [2] A.C. Anderson and M. Gettings, 'Gödel ontological proof revisited', in Gödel'96: Logical Foundations of Mathematics, Computer Science, and Physics: Lecture Notes in Logic 6, 167–172, Springer, (1996).

- [3] C.A. Anderson, 'Some emendations of Gödel's ontological proof', *Faith and Philosophy*, **7**(3), (1990).
- [4] P.B. Andrews, 'General models and extensionality', *Journal of Symbolic Logic*, 37(2), 395–397, (1972).
- [5] P.B. Andrews, 'Church's type theory', in *The Stanford Encyclopedia of Philosophy*, ed., E.N. Zalta, spring 2014 edn., (2014).
- [6] C. Benzmüller, 'HOL based universal reasoning', in *Handbook of the 4th World Congress and School on Universal Logic*, ed., J.Y. Beziau et al., pp. 232–233, Rio de Janeiro, Brazil, (2013).
- [7] C. Benzmüller and D. Miller, 'Automation of higher-order logic', in Handbook of the History of Logic, Volume 9 — Logic and Computation, Elsevier, (2014). Forthcoming; preliminary version available at http://christoph-benzmueller.de/papers/B5.pdf.
- [8] C. Benzmüller, J. Otten, and Th. Raths, 'Implementing and evaluating provers for first-order modal logics', in *Proc. of the 20th European Conference on Artificial Intelligence (ECAI)*, pp. 163–168, (2012).
- [9] C. Benzmüller and L.C. Paulson, 'Exploring properties of normal multimodal logics in simple type theory with LEO-II', in *Festschrift in Honor of Peter B. Andrews on His 70th Birthday*, ed., C. Benzmüller et al., 386–406, College Publications, (2008).
- [10] C. Benzmüller and L.C. Paulson, 'Quantified multimodal logics in simple type theory', Logica Universalis, 7(1), 7–20, (2013).
- [11] C. Benzmüller, F. Theiss, L. Paulson, and A. Fietzke, 'LEO-II a cooperative automatic theorem prover for higher-order logic', in *Proc. of IJCAR* 2008, number 5195 in LNAI, pp. 162–170. Springer, (2008).
- [12] C. Benzmüller and B. Woltzenlogel-Paleo, 'Formalization, mechanization and automation of Gödel's proof of God's existence', arXiv:1308.4526, (2013).
- [13] C. Benzmüller and B. Woltzenlogel-Paleo, 'Gödel's God in Is-abelle/HOL', Archive of Formal Proofs, (2013).
- [14] C. Benzmüller and B. Woltzenlogel-Paleo, 'Gödel's God on the computer', in *Proceedings of the 10th International Workshop on the Implementation of Logics*, EPiC Series. EasyChair, (2013). Invited abstract.
- [15] Y. Bertot and P. Casteran, Interactive Theorem Proving and Program Development, Springer, 2004.
- [16] J.C. Blanchette and T. Nipkow, 'Nitpick: A counterexample generator for higher-order logic based on a relational model finder', in *Proc. of ITP 2010*, number 6172 in LNCS, pp. 131–146. Springer, (2010).
- [17] C.E. Brown, 'Satallax: An automated higher-order prover', in *Proc. of IJCAR 2012*, number 7364 in LNAI, pp. 111 117. Springer, (2012).
- [18] R. Corazzon. Contemporary bibliography on ontological arguments: http://www.ontology.co/biblio/ontological-proof-contemporary-biblio.htm.
- [19] M. Fitting, Types, Tableaux and Gödel's God, Kluwer, 2002.
- [20] M. Fitting and R.L. Mendelsohn, First-Order Modal Logic, volume 277 of Synthese Library, Kluwer, 1998.
- [21] D. Gallin, Intensional and Higher-Order Modal Logic, North-Holland, 1975.
- [22] P. Garbacz, 'PROVER9's simplifications explained away', Australasian Journal of Philosophy, 90(3), 585–592, (2012).
- [23] K. Gödel, Appx.A: Notes in Kurt Gödel's Hand, 144–145. In [31], 2004.
- [24] L. Henkin, 'Completeness in the theory of types', *Journal of Symbolic Logic*, 15(2), 81–91, (1950).
- [25] R. Muskens, 'Higher Order Modal Logic', in *Handbook of Modal Logic*, ed., P Blackburn et al., 621–653, Elsevier, Dordrecht, (2006).
- [26] T. Nipkow, L.C. Paulson, and M. Wenzel, Isabelle/HOL: A Proof Assistant for Higher-Order Logic, number 2283 in LNCS, Springer, 2002.
- [27] P.E. Oppenheimera and E.N. Zalta, 'A computationally-discovered simplification of the ontological argument', *Australasian Journal of Philos*ophy, 89(2), 333–349, (2011).
- [28] J. Rushby, 'The ontological argument in PVS', in Proc. of CAV Workshop "Fun With Formal Methods", St. Petersburg, Russia, (2013).
- [29] S. Schulz, 'E a brainiac theorem prover', *AI Communications*, **15**(2), 111–126, (2002).
- [30] D. Scott, Appx.B: Notes in Dana Scott's Hand, 145–146. In [31], 2004.
- [31] J.H. Sobel, Logic and Theism: Arguments for and Against Beliefs in God, Cambridge U. Press, 2004.
- [32] G. Sutcliffe, 'The TPTP problem library and associated infrastructure', Journal of Automated Reasoning, 43(4), 337–362, (2009).
- [33] G. Sutcliffe and C. Benzmüller, 'Automated reasoning in higher-order logic using the TPTP THF infrastructure.', *Journal of Formalized Rea*soning, 3(1), 1–27, (2010).
- [34] B. Woltzenlogel-Paleo and C. Benzmüller, 'Automated verification and reconstruction of Gödel's proof of God's existence', OCG J., (2013).