

An Object-Logic Explanation for the Inconsistency in Gödel’s Ontological Theory (Extended Abstract)

Christoph Benz Müller^{1*} and Bruno Woltzenlogel Paleo²

¹ Freie Universität Berlin, Berlin, Germany
c.benzmuel1er@fu-berlin.de

² Australian National University, Canberra, Australia
bruno.wp@gmail.com

Abstract. This paper discusses the inconsistency in Gödel’s ontological argument. Despite the popularity of Gödel’s argument, this inconsistency remained unnoticed until 2013, when it was detected automatically by the higher-order theorem prover LEO-II. Complementing the meta-logic explanation for the inconsistency available in our IJCAI 2016 paper [6], we present here a new purely object-logic explanation that does not rely on semantic argumentation.

1 Introduction

Kurt Gödel’s ontological argument for the existence of God [9, 14] is amongst the most discussed formal proofs in modern literature. A rich body of publications – including very recent ones – present, discuss, assess, criticize, modify and improve Gödel’s original work (see e.g. Sobel [15] and Oppy [12] and the references therein).

Scott’s version of Gödel’s argument was automatically reconstructed by higher-order automated theorem provers [4] and its correctness was verified step-by-step in the Coq proof assistant [5]. To bridge the gap between higher-order logics (HOL; cf. [1] and the references therein), as used by these systems, and higher-order *modal* logics (HOML; cf. [10] and the references therein), on which the ontological argument is based, the logic embedding approach [2, 4] was used.

However, Gödel’s original axioms, as used in his manuscript [9], are inconsistent. This fact has remained unnoticed to philosophers until 2013, when LEO-II [3] found a surprising refutation of the axioms. In [6] we extracted from LEO-II’s machine-oriented refutation an informal and human-oriented intuitive explanation for the inconsistency, and we reconstructed and verified it in the ISABELLE proof assistant. But that explanation relied on reasoning at the *meta-logic* (HOL) level, which was only possible because of the embedding. Here we complement that work with a purely *object-logic* (HOML) explanation, and we compare and formalize both explanations in the Coq proof assistant.

Applications of (first-order) theorem proving technology in metaphysics were first reported by Fitelson, Oppenheimer and Zalta [8, 11]. Later on, Rushby [13] used the PVS proof assistant. Common to both works is a significant amount of proof-hand-coding work as well as their focus on a non-modal formalization of St. Anselm’s simpler and older ontological argument.

* This work was supported by the German Research Foundation DFG grant BE2501/9-1,2

```

theory Scott_SSU imports QML_SSU
begin
consts P :: "(μ⇒σ)⇒σ"
axiomatization where
A1a: "[|∀φ. P(φ) → ¬P(φ)|]" and
A1b: "[|∀φ. ¬P(φ) → P(φ)|]" and
A2: "[|∀φ ψ. P(φ) ∧ □(∀x. φ(x) → ψ(x)) → P(ψ)|]"
definition G where
"G(x) = (∀φ. P(φ) → φ(x))"
axiomatization where
A3: "[|P(G)|]" and
A4: "[|∀φ. P(φ) → □(P(φ))|]"
definition ess (infixr "ess" 85) where
"φ ess x = φ(x) ∧ (∀ψ. ψ(x) → □(∀y. φ(y) → ψ(y)))"
definition NE where
"NE(x) = (∀φ. φ ess x → □(∃ φ))"
axiomatization where
A5: "[|P(NE)|]"

theorem T3: "[|□ (∃ G)|]" -- {* LE0-11 proves T3 in 2,5sec *}
sledgehammer [provers = remote_leo2]
by (metis (lifting, full_types)
A1a A1b A2 A3 A4 A5 G_def NE_def ess_def)

lemma True nitpick [satisfy,user_axioms,expect=genuine] oops
-- {* Consistency is confirmed by Nitpick *}

theorem T2: "[|∀x. G(x) → G ess x|]"
sledgehammer [provers = remote_leo2]
by (metis A1b A4 G_def ess_def)

lemma MC: "[|∀φ. φ → □(φ)|]" -- {* Modal Collapse *}
sledgehammer [provers = remote_satallax, timeout=600]
by (meson T2 T3 ess_def)
end

theory GoedelGodWithoutConjunctInEss_K imports QML
begin
consts P :: "(μ⇒σ)⇒σ"
definition ess (infixr "ess" 85) where
"φ ess x = (∀ψ. ψ(x) → □(∀y. φ(y) → ψ(y)))"
definition NE where
"NE x = (∀φ. φ ess x → □(∃ φ))"
definition EmptyProperty ("∅") where
"∅ = (λx.λw. False)"

axiomatization where
A1a: "[|∀φ. P(φ) → ¬P(φ)|]" and
A2: "[|∀φ. ∀ψ. (P(φ) ∧ □(∀x. φ(x) → ψ(x))) → P(ψ)|]"

theorem T1: "[|∀φ. P(φ) → □(∃(φ))|]"
by (metis A1a A2)

lemma L1: "[|∀x. (∅ ess x)|]" (* Empty Essence Lemma *)
by (metis EmptyProperty_def ess_def)

axiomatization where
A5: "[|P(NE)|]"

lemma False (* Inconsistency *)
by (metis EmptyProperty_def A5 L1 NE_def T1)
end

```

Fig. 1: Scott’s consistent axioms (left) and proof of the inconsistency of (a subset of) Gödel’s original axioms (right)

2 An Essential Difference in the Definitions of Essence

Gödel’s manuscript can be considered a translation of Leibniz’s ideas on the argument into modern modal logic. Gödel discussed his manuscript with Scott, who shared a slightly different version with a larger public. Scott’s version of the axioms and definitions, formalized in ISABELLE, is shown in Fig. 1. The main difference to Gödel’s version is an extra conjunct in the definition of *essence* (*ess*). For Scott, an essential property of an individual must be possessed by him/her. For Gödel, this is not required.

Gödel’s omission has been considered inessential and merely an oversight by many. For more than four decades, its serious consequences remained unnoticed, despite numerous analyses and criticisms of the argument. However, as explained here, the extra conjunct is in fact crucial. Without it, Gödel’s original axioms are inconsistent. With it, Scott’s axioms are consistent (cf. Fig. 1, where the model finder NITPICK [7] confirms consistency). In personal communication, Dana Scott confirmed that he was unaware that Gödel’s original axioms were inconsistent.

3 Automating HOML in HOL

In our experiments in this branch of metaphysics we utilize an embedding of HOMLs, such as **K**, **KB** and **S5** with various domain conditions (possibilist and actualist quantification), in HOL. More precisely, formulas in HOML are *lifted*, i.e., converted into predicates over worlds, which are themselves explicitly represented as terms. The logical constants of HOML are translated to HOL terms in such a way that, for instance, $\Box\phi$

and $\diamond\varphi$ (relative to a current world w_0) are mapped, respectively, to the HOL formulas $\forall w.(rw_0w) \rightarrow (\varphi w)$ and $\exists w.(rw_0w) \wedge (\varphi w)$. This form of embedding is precisely the well-known standard translation, which is here intra-logically realized — and extended for quantifiers — in HOL by stating a set of equations defining the logical constants. The resulting logic is the HOML **K** with rigid terms and constant domains (possibilist quantifiers). Other logics (e.g. **KB**, **S5**) are embedded by adding axioms that restrict the accessibility relation r . Varying domains and actualist quantifiers can be simulated by using an existence predicate to guard the quantifiers.

4 Intuitive Explanations for the Inconsistency

In the typical workflow during an attempt to prove a conjecture with a theorem prover, it is customary to check the consistency of the axioms first. For if the axioms are inconsistent, anything (including the conjecture) would be trivially derivable in classical logic (*ex falso quodlibet*). Surprisingly, when this routine check was performed on Gödel’s axioms [4], the LEO-II prover claimed that the axioms were inconsistent. Unfortunately, the refutation generated by LEO-II was barely human-readable. The refutation was based on machine-oriented inference rules (a higher-order resolution calculus [3]), and the text file had 153 lines (with an average of 184 characters per line) and used a machine-oriented syntax (TPTP THF [16]).

Although LEO-II’s resolution refutation is not easy to read for humans, it did contain relevant hints to the importance of the empty property³ $\lambda x.\perp$ (also denoted \emptyset , as in HOL it is customary to think of unary predicates as sets)⁴. Based on this hint, we conceived the following informal explanation for the inconsistency of Gödel’s axioms (reproduced without change from [6]):

1. From Gödel’s definition of essence ($\phi \text{ ess } x \leftrightarrow \forall \psi(\psi(x) \rightarrow \Box \forall y(\phi(y) \rightarrow \psi(y)))$) it follows that the empty property (or self-difference) is an essence of every individual
(Empty Essence Lemma): $\forall x (\emptyset \text{ ess } x)$
2. From axiom A5 (‘necessary existence’ is a positive property: $P(NE)$) and theorem T1 (*Positive properties are possibly exemplified*: $\forall \phi[P(\phi) \rightarrow \Diamond \exists x\phi(x)]$), it follows that NE is possibly exemplified:
 $\Diamond \exists x[NE(x)]$
3. Expanding the definition of ‘necessary existence’ ($NE(x) \equiv \forall \phi[\phi \text{ ess } x \rightarrow \Box \exists y\phi(y)]$), the following is obtained:
 $\Diamond \exists x[\forall \phi[\phi \text{ ess } x \rightarrow \Box \exists y[\varphi(y)]]]$
4. The sentence above holds for all φ and thus, in particular, for the empty property (or self-difference):
 $\Diamond \exists x[\emptyset \text{ ess } x \rightarrow \Box \exists y[\emptyset(y)]]$
5. By the Empty Essence Lemma, the antecedent of the implication above is valid. Therefore, the sentence above entails:
 $\Diamond \exists x[\Box \exists y[\emptyset(y)]]$
6. By definition of \emptyset :
 $\Diamond \exists x[\Box \perp]$
7. As the existential quantifier is binding no variable within its scope, the sentence is equi-valid with:
 $\Diamond \Box \perp$

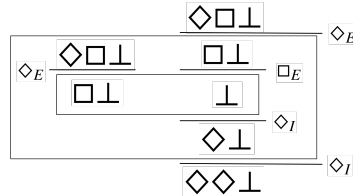
³ Note that the terms for the empty property ($\lambda x.\perp$) and for the property of self-difference ($\lambda x.x \neq x$) have identical denotations in the logic setting with functional and Boolean extensionality assumed here. For the proof to go through, it is irrelevant which property is used.

⁴ An additional lambda abstraction occurs in the empty property in LEO-II’s proof (and in the reconstruction in ISABELLE) because the embedding approach lifts the boolean type o to $\iota \rightarrow o$.

8. To see that the sentence above is contradictory, we may reason semantically, thinking of possible worlds. If w_0 is the arbitrary current world, the \diamond operator forces the existence of a world w accessible from w_0 such that $\Box\perp$ is true in w . But $\Box\perp$ can only be true in w , if there is no world w' accessible from w . In logics⁵ with a reflexive or symmetric accessibility relation (e.g. **KB**), it is easy to see that there must be a world w' accessible from w : either w' itself, in case of a reflexive relation, or w_0 , in case of a symmetric relation. In fact, even in **K**, with no accessibility condition, there must be a world w' accessible from w . The reason is that $\diamond\Box\perp$ should be *valid* (true in all worlds). Therefore, it is true in w as well, where the existence of an accessible world w' is forced by the \diamond operator. As a model for $\diamond\Box\perp$ (which is a consequence of Gödel's axioms) cannot be built, Gödel's axioms are inconsistent.

If we were to convert the informal proof above to a formal proof, the semantic reasoning in step 8 would require a leap to the meta-logic (HOL), in order to expand the definitions of modal operators and reason directly about possible worlds. The alternative proof below avoids this leap and remains purely within the object logic (HOML **K**):

- 8* . We must derive \perp from $\diamond\Box\perp$. In order to derive \perp , it suffices to show that there exists a derivable proposition such that its negation is also derivable. We choose $\diamond\Box\perp$ as the candidate proposition, and hence we must show that:
- (a) $\neg\diamond\Box\perp$ is *derivable*: this proposition is equi-valid to $\Box\Box\top$, which is trivially derivable from \top by two applications of the necessitation inference rule.
 - (b) $\diamond\Box\perp$ is *derivable*: and indeed, it can be derived (using a recently developed natural deduction calculus for modal logic **K** [5]) as follows:



An interesting and unusual feature of the derivation shown above is that the leftmost \diamond_E (diamond elimination) inference derives a formula ($\Box\perp$) that is never used as a premise. This is necessary because of the *eigen-box* condition, which requires that every box must be accessed by exactly one *strong* modal rule. The purpose of the *strong* \diamond_E inference is merely to create and access the innermost box that is needed by the *weak* \Box_E and \diamond_I inferences inside the outermost box.

The proofs above have been formalized and verified step-by-step in Coq. The complete proofs can be found in <https://github.com/FormalTheology/GoedelGod>. The following Coq script shows the formalization of step 8 of the meta-logic proof.

⁵ Interestingly, the refutation automatically generated by LEO-II uses a symmetric accessibility relation, and thus requires the modal logic **KB**. The informal, human-constructed refutations described here, on the other hand, requires only the weaker modal logic **K**. In our experiments LEO-II (like all other HOL provers) was still too weak to automatically prove the inconsistency already in logic **K**. Hence, this remains an open problem for automated theorem provers.

```

Lemma dia_box_false_to_false_meta: [(dia (box mFalse))] -> [mFalse].
Proof. intro H. intro w.
destruct (H w) as [w0 [R0 H0]]. destruct (H w0) as [w1 [R1 H1]].
box_elim H0 w1 HF. unfold mFalse in HF. destruct HF as [p [HF1 HF2]].
contradiction. Qed.

```

The other Coq scripts below show the formalization of step 8* of the object-logic proof.

```

Lemma mimplies_to_mnot: [mforall p:o, (p m-> mFalse) m-> (m~ p)].
Proof. mv. intro p. intro H. intro H0.
destruct (H H0) as [p0 [H1 H2]]. apply H2. exact H1. Qed.

Lemma dia_not_not_box: [ mforall p, (dia (m~ p)) m-> (m~ (box p)) ].
Proof. mv. intro p. intro H1. intro H2.
dia_e H1. apply H. box_e H2 H3. exact H3. Qed.

Lemma dia_box_false_to_false_object: [(dia (box mFalse))] -> [mFalse].
Proof. intro H. intro w. exists (dia (dia mFalse)).
split.
  dia_e (H w). dia_e (H w0). dia_i w0. dia_i w1. box_e H0 H3. exact H3.

  apply box_not_not_dia. box_i. apply box_not_not_dia. box_i.
  apply mimplies_to_mnot. intro H4. exact H4. Qed.

```

Due to a deliberate and disciplined use of only the simplest (and non-automatic) Coq tactics, there is a straightforward correspondence between the tactics used in the scripts above and the inference rules of the modal natural deduction calculus [5]. Therefore, the lengths of the proof scripts (in number of tactic applications) can serve as estimations for the lengths of the corresponding natural deduction proof. It is noticeable that the meta-logic proof is significantly shorter than the pure object-logic proof. An in-depth analysis reveals that the reasoning about the possible worlds semantics in the meta-logic proof acts as a short-cut: when it becomes impossible (in step 8) to build the third world w' (because \perp would have to hold in it, and thus w' would be contradictory), a contradiction at the HOL level can be immediately derived, completing the proof. In the object-level proof, on the other hand, such a contradiction has to be found in the arbitrary initial world w . This requires not only additional tedious logical inferences (cf. the proofs of the lemmas `mimplies_to_mnot` and `dia_not_not_box`), but also a non-trivial guessing of the contradictory proposition $\diamond\diamond\perp$, whose purpose is precisely to carry over the contradiction from w' back to w .

5 Conclusion

The axioms and definitions in Gödel's manuscript are inconsistent (even in the weakest modal logic **K**); this was detected automatically by the prover LEO-II. In our previous work [6], we presented a human-readable and intuitive meta-logic explanation for the inconsistency, and we formalized and semi-automatically reconstructed it in the ISABELLE proof assistant. Here this work was extended with an object-level explanation, and both explanations were formalized step-by-step in Coq proof assistant. A comparison of the formal Coq proofs of both explanations revealed that the meta-logic reasoning is more powerful, because it enables shortcuts and, therefore, requires fewer inferences

and guesses. We conjecture that this is not accidental, but rather a fundamental reason why the embedding approach is effective in practice.

It is kind of entertaining that our work reveals a mistake in Gödel’s manuscript and at same time further substantiates Gödel’s belief that “there is a scientific (exact) philosophy and theology, which deals with concepts of the highest abstractness; and this is also most highly fruitful for science.” [17][p. 316]. Indeed, through the investigation of Gödel’s mistake, we have been led to an interesting little conjecture in automated reasoning and proof theory (the global axiom $\diamond\Box\perp$ is inconsistent).

References

1. P.B. Andrews. Church’s type theory. In E.N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Spring 2014 edition, 2014.
2. C. Benzmüller and L.C. Paulson. Quantified multimodal logics in simple type theory. *Logica Universalis*, 7(1):7–20, 2013.
3. C. Benzmüller, L.C. Paulson, N. Sultana, and F. Thei. The higher-order prover LEO-II. *Journal of Automated Reasoning*, 55(4):389–404, 2015.
4. C. Benzmüller and B. Woltzenlogel-Paleo. Automating Gödel’s ontological proof of God’s existence with higher-order automated theorem provers. In Torsten Schaub, Gerhard Friedrich, and Barry O’Sullivan, editors, *ECAI 2014*, volume 263 of *Frontiers in Artificial Intelligence and Applications*, pages 93 – 98. IOS Press, 2014.
5. C. Benzmüller and B. Woltzenlogel Paleo. Interacting with modal logics in the coq proof assistant. In Lev D. Beklemishev and Daniil V. Musatov, editors, *Computer Science - Theory and Applications - 10th International Computer Science Symposium in Russia, CSR 2015, Listvyanka, Russia, July 13-17, 2015, Proceedings*, volume 9139 of *LNCS*, pages 398–411. Springer, 2015.
6. C. Benzmüller and B. Woltzenlogel Paleo. The inconsistency in Gödel’s ontological argument: A success story for AI in metaphysics. In *IJCAI 2016*, 2016.
7. J.C. Blanchette and T. Nipkow. Nitpick: A counterexample generator for higher-order logic based on a relational model finder. In *ITP 2010*, number 6172 in *LNCS*, pages 131–146. Springer, 2010.
8. B. Fitelson and E. N. Zalta. Steps toward a computational metaphysics. *J. Philosophical Logic*, 36(2):227–247, 2007.
9. K. Gödel. *Appx. A: Notes in Kurt Gödel’s Hand*, pages 144–145. In Sobel [15], 1970.
10. R. Muskens. Higher Order Modal Logic. In P. Blackburn et al., editor, *Handbook of Modal Logic*, Studies in Logic and Practical Reasoning, pages 621–653. Elsevier, Dordrecht, 2006.
11. P.E. Oppenheimer and E.N. Zalta. A computationally-discovered simplification of the ontological argument. *Australasian J. of Philosophy*, 89(2):333–349, 2011.
12. G. Oppy. Ontological arguments. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Spring 2015 edition, 2015.
13. J. Rushby. The ontological argument in PVS. In *Proc. of CAV Workshop “Fun With Formal Methods”*, St. Petersburg, Russia, 2013.
14. D. Scott. *Appx. B: Notes in Dana Scott’s Hand*, pages 145–146. In Sobel [15], 1972.
15. J.H. Sobel. *Logic and Theism: Arguments for and Against Beliefs in God*. Cambridge U. Press, 2004.
16. G. Sutcliffe and C. Benzmüller. Automated reasoning in higher-order logic using the TPTP THF infrastructure. *J. of Formalized Reasoning*, 3(1):1–27, 2010.
17. H. Wang. *A Logical Journey: From Gödel to Philosophy*. MIT Press, 1996.