

Löb's Theorem in Coq: Taking stock of where we are

Janis Bailitis

Bachelor Project, Saarland University

Advisors: Dr. Yannick Forster and Dr. Dominik Kirst

Supervisor: Prof. Dr. Gert Smolka

22nd April 2024

This is a brief memo providing an overview of the current state of my Bachelor project, focusing on the goals achieved so far, and how the results relate to the work previously done at the Programming Systems Lab. Further, we aim to point out how our approach fits into the broader landscape of (provability) logic.

1 Preliminaries

At first, let us fix a notational convention: If we have a formula φ and write $\varphi(x_1, \dots, x_n)$, we indicate that φ has n free variables, unless stated otherwise.

We work in constructive type theory and build upon some key results obtained from formalising logic in type theory. We rely on results established by Kirst et al. [HK23],[KH23],[KP23]. Most importantly, we build upon a variant of Church's thesis for Robinson arithmetic, stating that any function $f : \mathbb{N} \rightarrow \mathbb{N}$ can be represented by a Σ_1 formula in \mathbb{Q} . More formally, we assume:

Axiom 1.1 (CT_Q, cf. [HK23]) For all $f : \mathbb{N} \rightarrow \mathbb{N}$ there exists a Σ_1 -formula $\varphi(x, y)$ such that for all $x : \mathbb{N}$, $\mathbb{Q} \vdash \forall y. \varphi(\bar{x}, y) \leftrightarrow y \equiv \overline{fx}$.

Later, this axiom was refined to account for partial functions $\mathbb{N} \dashrightarrow \mathbb{N}$. The initial variant of CT_Q can be derived from this statement.

Axiom 1.2 ((partial) CT_Q, cf. [KP23]) For all $f : \mathbb{N} \dashrightarrow \mathbb{N}$ there exists a Σ_1 -formula $\varphi(x, y)$ such that for all $x, y : \mathbb{N}$, $fx \downarrow y \leftrightarrow \mathbb{Q} \vdash \forall z. \varphi(\bar{x}, z) \leftrightarrow z \equiv \bar{y}$.

Kirst and Peters [KP23, section 5] use the term CT_Q to refer to Axiom 1.2, while we will refer to Axiom 1.1 by CT_Q as we do not use partial CT_Q directly. Using CT_Q, we can derive the **Weak Representability Theorem** playing a pivotal role in our project:

Fact 1.3 (Weak Representability Theorem) Suppose $P : \mathbb{N} \rightarrow \mathbb{P}$ is an enumerable predicate. We can find a Σ_1 -formula $\varphi(x)$ such that for all $x : \mathbb{N}$, $Px \leftrightarrow \mathbb{Q} \vdash \varphi(\bar{x})$.

Another result following from partial CT_Q we will use is the **Strong Separation Theorem**:

Fact 1.4 (Strong Separation Theorem) Let $P, R : \mathbb{N} \rightarrow \mathbb{P}$ be disjoint, semi-decidable predicates. There exists a Σ_1 -formula $\varphi(x)$ satisfying

$$\forall n. Pn \rightarrow Q \vdash \varphi(\bar{n}) \wedge Rn \rightarrow Q \vdash \dot{\neg}\varphi(\bar{n})$$

This gives rise to the following fact for decidable predicates:

Corollary 1.5 Let $P : \mathbb{N} \rightarrow \mathbb{P}$ be a decidable predicate. There exists a Σ_1 -formula $\varphi(x)$ such that

$$\forall n. Pn \leftrightarrow Q \vdash \varphi(\bar{n}) \wedge \neg Pn \leftrightarrow Q \vdash \dot{\neg}\varphi(\bar{n}).$$

Proof Since P is decidable, both P and $\lambda n. \neg Pn$ are semi-decidable. Further, both are clearly disjoint. By the previous fact, we obtain a Σ_1 -formula $\varphi(x)$ such that

$$\forall n. Pn \rightarrow Q \vdash \varphi(\bar{n}) \wedge \neg Pn \rightarrow Q \vdash \dot{\neg}\varphi(\bar{n}).$$

For each of the the converses $Q \vdash \varphi(\bar{n}) \rightarrow Pn$ and $Q \vdash \dot{\neg}\varphi(\bar{n}) \rightarrow \neg Pn$, we decide Pn . If the decision matches the conclusion of the implications, we are done. Otherwise, for the first converse, we know that $\neg Pn$ as well as $Q \vdash \varphi(\bar{n})$. Exploiting $\neg Pn \rightarrow Q \vdash \dot{\neg}\varphi(\bar{n})$ gives us $Q \vdash \dot{\neg}\varphi(\bar{n})$. Contradictory as Q is consistent using the standard model. Similarly, one handles the case Pn when proving $Q \vdash \dot{\neg}\varphi(\bar{n}) \rightarrow \neg Pn$. ■

Equipped with these notions, we started our exploration of Löb's theorem.

2 Löb's Theorem

The initial goal was to show Löb's Theorem, a result of provability logic stating that when we prove a formula in a sufficiently strong formal system, we can assume provability of this formula without additional cost. More formally, there is a modality, the box-modality (written \Box), expressing provability of a formula. Löb's theorem can then be stated as the following rule, where \mathcal{T} is an arbitrary theory:

Theorem 2.1 (Löb) Suppose φ is a sentence. In order to prove $\mathcal{T} \vdash \varphi$, it suffices to show $\mathcal{T} \vdash \Box\varphi \rightarrow \varphi$.

Three questions arose: Firstly, how can \Box be defined? Secondly – given a sensible definition of \Box – how can Löb's Theorem be shown? Lastly, which theory should we choose? We assumed that Peano arithmetic suffices, so we raised the question whether Robinson arithmetic is still enough.

The second question admits a fortunate answer. As long as \Box obeys four axioms (pointed out below), Löb's theorem can be derived. In particular, we can abstract away concrete implementation details of the modality's definition. The required axioms are as follows:

1. Necessitation: For all $\varphi : \mathbb{F}$, $\mathcal{T} \vdash \varphi$ implies $\mathcal{T} \vdash \Box\varphi$.
2. Modal fixed points: Whenever φ is a sentence, we can construct another sentence ψ such that $Q \vdash \psi \leftrightarrow (\Box\psi \rightarrow \varphi)$.
3. Internal necessitation: For all $\varphi : \mathbb{F}$, we can derive $\mathcal{T} \vdash \Box\varphi \rightarrow \Box\Box\varphi$.
4. Box distributivity: For all $\varphi, \psi : \mathbb{F}$, we have $\mathcal{T} \vdash \Box(\varphi \rightarrow \psi) \rightarrow \Box\varphi \rightarrow \Box\psi$.

For any operation \Box on formulas and theories \mathcal{T} obeying these rules, Löb's theorem can be derived; such an argument has been done by Boolos et al. [BBJ07, theorem 18.4], with a slightly different handling of modal fixed points. A first version of this proof is formalised in Coq in the file `coq/Loeb_formalisation_fol.v`. As of now, the formalisation needs some minor adjustments to account for the current formulation of our axioms, but it is a purely modal-logical argument not using any internals of formulas and deduction systems (apart from the inference rules).

3 Defining the Box modality

In search of a definition for \Box , we wanted to obtain a formula $\text{Prov}(x)$ characterising provability in our deduction system. At first, we worked in Robinson arithmetic. Intuitively, we asked for the property $\forall \tau : \mathbb{F}. Q \vdash \tau \leftrightarrow Q \vdash \text{Prov}(\tau)$. By defining $\Box \tau := \text{Prov}(\tau)$, we would have a candidate for the Box modality. Unfortunately, $\text{Prov}(\tau)$ is not syntactically correct, as we may only substitute in terms, but not entire formulas. Gödel numbering is the road to success in resolving this issue:

Fact 3.1 (Gödel numberings) There exist functions $\text{göd} : \mathbb{F} \rightarrow \mathbb{N}$ and $\text{göd}^{-1} : \mathbb{N} \rightarrow \mathbb{F}$ inverting each other, i.e. $\forall n. \text{göd}(\text{göd}^{-1}(n)) = n$ and $\forall \varphi. \text{göd}^{-1}(\text{göd}(\varphi)) = \varphi$. We call $\text{göd}(\varphi)$ the **Gödel number** of φ .

Remark 3.2 In the Coq mechanisation, we consistently use `g` for göd and `f` for göd^{-1} . In the future, we will work with coercions to allow for a more readable Coq development.

Definition 3.3 For a formula φ , we define $\ulcorner \varphi \urcorner := \overline{\text{göd}(\varphi)}$. This is the **Gödel numeral** or **quine quote** of φ .

The term ‘quine quote’ was used by Norrish [Nor18].

Having this machinery set up, we can derive the Box modality.

Lemma 3.4 There exists a Σ_1 -formula $\text{Prov}(x)$ satisfying the following equivalence: $\forall \tau. Q \vdash \tau \leftrightarrow Q \vdash \text{Prov}(\ulcorner \tau \urcorner)$.

Proof We apply the weak representability theorem. The predicate $\lambda \varphi. Q \vdash \varphi$ is enumerable by standard techniques, cf. [FKS19]. A routine argument concludes enumerability of $\lambda n. Q \vdash \text{göd}^{-1}(n)$. Applying this theorem, we obtain a Σ_1 -formula $\text{Prov}(x)$ such that for all $m : \mathbb{N}$

$$Q \vdash \text{göd}^{-1}(m) = (\lambda n. Q \vdash \text{göd}^{-1}(n))m \leftrightarrow Q \vdash \text{Prov}(\overline{m}).$$

If τ is any formula, we can plug in $\text{göd}(\tau)$ for m and arrive at

$$Q \vdash \tau = Q \vdash \text{göd}^{-1}(\text{göd}(\tau)) \leftrightarrow Q \vdash \varphi(\overline{\text{göd}(\tau)}) = Q \vdash \varphi(\ulcorner \tau \urcorner)$$

and are done. ■

Finally, we can define the Box modality.

Definition 3.5 (Box modality) For any formula τ , we set $\Box \tau := \text{Prov}(\ulcorner \tau \urcorner)$.

Immediately, we obtain the necessitation rule.

Fact 3.6 (Necessitation for Box) For all formulas τ , we have that $Q \vdash \tau$ implies $Q \vdash \Box \tau$.

Proof Unfold the definition of \Box . The remaining goal is a direct consequence of Lemma 3.4. ■

Next, we were interested in whether this definition is strong enough for the remaining three properties our modality should have. Kirst noted that we may need to switch to Peano arithmetic for Box distributivity and internal necessitation, while the existence of modal fixed points may be derivable from Q alone, which was our subsequent aim. Kirst suspected that, along the way, we need Gödel's diagonal lemma, hence we inspected this result in depth.

4 The diagonal lemma

In computability theory, we are often interested in a program's behaviour on its own Gödel number. In a similar fashion, we can analyse a formula's truth or provability when its own Gödel numeral is substituted in for the free variables. More specifically, we are interested in substituting in a formula's Gödel numeral for one of the free variables.

Definition 4.1 (Diagonalisation, cf. [Smi22]) Let $\varphi(x)$ be any formula (with potentially more than one free variable). We say that $\varphi(\ulcorner \varphi \urcorner)$ is the **diagonalisation** of φ . This gives rise to a function $\text{diag} : \mathbb{F} \rightarrow \mathbb{F}, \varphi \mapsto \varphi(\ulcorner \varphi \urcorner)$.

We can define a similar function on the level of Gödel numbers:

Definition 4.2 $\text{diag}_{\mathbb{N}} : \mathbb{N} \rightarrow \mathbb{N}, n \mapsto \text{göd}(\text{diag}(\text{göd}^{-1}(n)))$ sends a formula's Gödel number to the Gödel number of its diagonalisation.

Note that CT_{Q} applies to $\text{diag}_{\mathbb{N}}$; this is why we introduced this function at all.

We have a straightforward interchangeability result.

Lemma 4.3 Let φ be a formula. Then $\text{diag}_{\mathbb{N}}(\text{göd}(\varphi)) = \text{göd}(\text{diag}(\varphi))$.

Proof Immediate from the definitions and the identity $\text{göd}^{-1}(\text{göd}(\varphi)) = \varphi$. ■

Ultimately, our aim is to show the following result, the well known **diagonal lemma**.

Theorem 4.4 (Diagonal lemma for Q) Let $B(x)$ be a formula. There exists a sentence G satisfying $\text{Q} \vdash G \leftrightarrow B(\ulcorner G \urcorner)$.

Interestingly, diagonalisation is nowhere mentioned in the theorem statement. However, diagonal thinking plays a pivotal role in its proof. Before heading straight to showing the claim, let us discuss a preliminary result making our argument easier to digest.

Lemma 4.5 There exists a Σ_1 -formula $\text{dg}(x, y)$ such that for all $\psi : \mathbb{F}, \text{Q} \vdash \forall y. \text{dg}(\ulcorner \psi \urcorner, y) \leftrightarrow y \equiv \ulcorner \text{diag}(\psi) \urcorner$.

Proof With CT_{Q} , we find a Σ_1 -formula $\text{dg}(x, y)$ satisfying, for all $x : \mathbb{N}$,

$$\text{Q} \vdash \forall y. \text{dg}(\bar{x}, y) \leftrightarrow y \equiv \overline{\text{diag}_{\mathbb{N}}(x)}.$$

Suppose $\psi : \mathbb{F}$. We obtain, by setting $x = \text{göd}(\psi)$,

$$\text{Q} \vdash \forall y. \text{dg}(\ulcorner \psi \urcorner, y) \leftrightarrow y \equiv \overline{\text{diag}_{\mathbb{N}}(\text{göd}(\psi))}.$$

We finish by rewriting with Lemma 4.3. ■

Having all this in mind, we can discuss the proof of Theorem 4.4. Apart from getting the substitutions right, over which we abstract here, the argument itself is not very involved provided that we instantiate the existential quantifier with a sensible candidate. Indeed, finding this formula is a bit delicate and hence the true art in showing this claim. The proof idea is from [BBJ07].

Proof (of Theorem 4.4) First, we define G . We inspect the formula $F(x) := \dot{\exists}y.dg(x, y) \wedge B(y)$. We set $G := \text{diag}(F) = F(\ulcorner F \urcorner)$. Clearly, G is closed.

We are left to show $Q \vdash G \leftrightarrow B(\ulcorner G \urcorner)$. After unfolding the definitions of G and F , we arrive at

$$Q \vdash (\dot{\exists}y.dg(\ulcorner F \urcorner, y) \wedge B(y)) \leftrightarrow B(\ulcorner G \urcorner).$$

The direction from left to right is straightforward: We introduce our assumption, and destruct the existential as well as the conjunction. That is, we have, for some term y ,

$$Q, dg(\ulcorner F \urcorner, y), B(y) \vdash B(\ulcorner G \urcorner).$$

We apply Lemma 4.5 to our assumption $dg(\ulcorner F \urcorner, y)$ and obtain the proof state

$$Q, dg(\ulcorner F \urcorner, y), B(y), y \equiv \ulcorner \text{diag}(F) \urcorner \vdash B(\ulcorner G \urcorner).$$

Since $G = \text{diag}(F)$, we finish by rewriting on the meta level followed by a rewrite with $y \equiv \ulcorner G \urcorner$ in the deduction system.

For the converse, we have to show

$$Q, B(\ulcorner G \urcorner) \vdash \dot{\exists}y.dg(\ulcorner F \urcorner, y) \wedge B(y).$$

Instantiating the existential with $\ulcorner G \urcorner$ leaves us to prove

$$Q, B(\ulcorner G \urcorner) \vdash dg(\ulcorner F \urcorner, \ulcorner G \urcorner) \wedge B(\ulcorner G \urcorner).$$

By the assumption and conjunction rules as well as weakening, we are left to show

$$Q \vdash dg(\ulcorner F \urcorner, \ulcorner G \urcorner)$$

following from Lemma 4.5 as $G = \text{diag}(F)$. ■

5 Applications of the diagonal lemma

Having established the diagonal lemma, we can draw our attention to some important corollaries, including the existence of modal fixed points, contributing to the proof of Löb's theorem.

5.1 Modal fixed points

We will now clarify how the diagonal lemma contributes to the fixed point theorem. The high-level idea is, given a sentence φ , to apply the diagonal lemma to the formula $\text{Prov}(x) \rightarrow \varphi$. The result will then satisfy our requirements. The following proof accounts for all technical details.

Theorem 5.1 (Modal fixed points) For all sentences φ , there exists a sentence ψ such that $Q \vdash \psi \leftrightarrow (\Box \psi \rightarrow \varphi)$

Proof We use the diagonal lemma to obtain a sentence ψ satisfying $Q \vdash \psi \leftrightarrow (\text{Prov}(\ulcorner \psi \urcorner) \rightarrow \varphi)$. Applying the definition of Box gives us our goal $Q \vdash \psi \leftrightarrow (\Box \psi \rightarrow \varphi)$. ■

5.2 Tarski's indefinability result

In this and the subsequent section, we explore certain limits of our deductive system following from the diagonal lemma. From CT_Q , we concluded capability of Q to represent type-theoretic functions. In a similar fashion, we can ask whether there is a notion to 'represent' predicates $\mathbb{N} \rightarrow \mathbb{P}$ inside the deduction system. Indeed, there is, and one talks of a predicate being **definable**. The definition is taken from Boolos et al. [BBJ07].

Definition 5.2 (Definable) Let $P : \mathbb{N} \rightarrow \mathbb{P}$ be a predicate and \mathcal{T} a theory. The formula $\varphi(x)$ **defines** P in \mathcal{T} if, for all $n : \mathbb{N}$, we have both $Pn \rightarrow \mathcal{T} \vdash \varphi(\bar{n})$ and $\neg Pn \rightarrow \mathcal{T} \vdash \dot{\neg}\varphi(\bar{n})$.

Definition 5.3 (Definability) A predicate $P : \mathbb{N} \rightarrow \mathbb{P}$ is said to be **definable** in a theory \mathcal{T} if there exists a formula defining P in \mathcal{T} .

Remark 5.4 Every decidable predicate is definable in any extension of Q by Corollary 1.5 and weakening.

Next, we focussed on the predicate $\lambda n. \mathcal{T} \vdash \text{göd}^{-1}(n)$ for any consistent extension \mathcal{T} of Q . It turned out that it is not definable.

Lemma 5.5 (Indefinability lemma) Let \mathcal{T} be a consistent theory extending Q . $\lambda n. \mathcal{T} \vdash \text{göd}^{-1}(n)$ is not definable in \mathcal{T} .

Before proving this claim, we remind ourselves that we can use classical reasoning when proving distinguished claims, the so called **stable** claims (cf. [Smo24]):

Definition 5.6 Suppose $P : \mathbb{P}$. We say that P is **stable** if $\neg\neg P \rightarrow P$.

Fact 5.7 (Classical reasoning for stable claims) Suppose $P, Q : \mathbb{P}$. We have $\text{stable}(P) \rightarrow (Q \vee \neg Q \rightarrow P) \rightarrow P$.

We are now ready to show the indefinability lemma. The proof idea is due to Boolos et al. [BBJ07].

Proof (of Lemma 5.5) Suppose $P := \lambda n. \mathcal{T} \vdash \text{göd}^{-1}(n)$ was definable witnessed by $\varphi(x)$. Using the diagonal lemma, we find a sentence G such that $Q \vdash G \leftrightarrow \dot{\neg}\varphi(\ulcorner G \urcorner)$.

We have to show \perp , which is stable by a straightforward argument. Using Fact 5.7, we assume that $\mathcal{T} \vdash G \vee \mathcal{T} \not\vdash G$. Case analysis.

If $\mathcal{T} \vdash G$, there exists a list $A \subseteq \mathcal{T}$ witnessing this, i.e. $A \vdash G$. As $\varphi(x)$ defines P , and $P(\text{göd}(G))$ holds, we conclude $\mathcal{T} \vdash \varphi(\ulcorner G \urcorner)$, i.e. $B \vdash \varphi(\ulcorner G \urcorner)$ for some list $B \subseteq \mathcal{T}$. By weakening, we obtain $Q, A, B \vdash G \leftrightarrow \dot{\neg}\varphi(\ulcorner G \urcorner)$, $Q, A, B \vdash G$ as well as $Q, A, B \vdash \varphi(\ulcorner G \urcorner)$. We apply \mathcal{T} 's consistency to our goal and have to show $\mathcal{T} \vdash \perp$. Straightforward from the assumptions, as Q, A, B is contained in \mathcal{T} .

If we obtain $\mathcal{T} \not\vdash G$, we apply this assumption to our goal and have to show $\mathcal{T} \vdash G$. Now, since $\neg P(\text{göd}(G))$ holds, we conclude $\mathcal{T} \vdash \dot{\neg}\varphi(\ulcorner G \urcorner)$, i.e. $A \vdash \dot{\neg}\varphi(\ulcorner G \urcorner)$ for some list $A \subseteq \mathcal{T}$. Again, by weakening, have $Q, A \vdash G \leftrightarrow \dot{\neg}\varphi(\ulcorner G \urcorner)$ as well as $Q, A \vdash \dot{\neg}\varphi(\ulcorner G \urcorner)$. As our goal is $\mathcal{T} \vdash G$, we are done since Q, A is contained in \mathcal{T} . ■

While this result is somewhat abstract, it has interesting consequences. One of which is **Tarski's Theorem** stating that the theory of sentences correct for \mathbb{N} is not definable.¹

¹As of now, I am not sure how this 'correctness' is modelled in our set-up. In Boolos' et al. book, a formula is 'correct' if it is true in the his notion of standard interpretation, which most likely coincides with our [KH23, Definition 4] definition of \vDash with \mathbb{N} as domain type, picking the standard meta-level predicates for the predicate symbols. I am working on resolving this issue.

Corollary 5.8 (Tarski’s Theorem.) The theory $\lambda\varphi.\mathbb{N} \vDash \varphi$ is not definable.

Proof The theory in question is a consistent extension of Q and hence not definable by Lemma 5.5. ■

Finally, we can conclude this section with an argument that for any consistent extension \mathcal{T} of Q , it is not decidable whether a formula is derivable from \mathcal{T} or not. Most prominently, we cannot decide whether a formula is a theorem of Q .

Corollary 5.9 Let \mathcal{T} be a consistent extension of Q . The predicate $\lambda\varphi.\mathcal{T} \vdash \varphi$ is not decidable.

Proof Suppose, for contradiction, that $P := \lambda\varphi.\mathcal{T} \vdash \varphi$ was decidable. A routine argument derives decidability of $R := \lambda n.\mathcal{T} \vdash \text{göd}^{-1}(n)$. By Corollary 1.5, R is definable, whereas Lemma 5.5 asserts undefinability of R . Absurd. ■

In essence, the results derived in this section indicate certain boundaries our deduction system has when it comes to interaction with the meta level – even if we allow for arbitrary consistent extensions of Q . This why these (and some more) theorems are also known as ‘limitative theorems’, for instance in book by Boolos et al. [BBJ07].

5.3 Gödel’s first incompleteness theorem

Next, we will focus on another consequence of the diagonal lemma: The existence of a sentence that can neither be proven nor refuted in Q . That is, we will prove a variant of Gödel’s first incompleteness result.

Theorem 5.10 (Gödel’s first incompleteness theorem for Q) There exists a sentence G such that neither $Q \vdash G$ nor $Q \vdash \neg G$.

Before heading to the proof, we need to draw our attention to intuitionistic and classical reasoning. As we want our results to be as general as possible, we do not stick to intuitionistic (\vdash_i) or classical (\vdash_c) reasoning in our theorem statements. Unfortunately, the current proof needs classical reasoning within the deduction system once, so we could only show that $Q \not\vdash_c \neg G$. While this immediately refutes any intuitionistic derivation due to classical reasoning being strictly stronger, the proof-mode had trouble destructing assumptions when in classical mode, so we needed to find a way circumventing this. We found such a way, and the details can be found in our proof. We would like to emphasise our reliance on the following result:

Fact 5.11 (Σ_1 -conservativity) Let φ be a Σ_1 -sentence such that $Q \vdash_c \varphi$. Then $Q \vdash \varphi$.

As opposed to the library, this is a slightly different yet derivable formulation.

For the proof of Gödel’s result, we use Boolos’ et al. version of Σ_1 stating that we can write $\text{Prov}(x) = \exists y.\text{Prf}(x, y)$ for some Δ_1 -formula $\text{Prf}(x, y)$. We elaborate on this step in the discussion.

Definition 5.12 (ω -(in)consistency) Let \mathcal{T} be a theory. We say that \mathcal{T} is **ω -inconsistent** if it proves $\mathcal{T} \vdash \exists x.\tau(x)$ for formula some $\tau(x)$, but also admits a derivation of $\mathcal{T} \vdash \neg\tau(\bar{n})$ for all $n : \mathbb{N}$.

A theory not being ω -inconsistent is said to be **ω -consistent**.

Axiom 5.13 Q is ω -consistent.

Proof (of Theorem 5.10) By the diagonal lemma, we can find a sentence G such that $Q \vdash G \leftrightarrow \neg \text{Prov}(\ulcorner G \urcorner)$. We show that $Q \not\vdash G$ as well as $Q \not\vdash \neg G$.

Suppose $Q \vdash G$. By Lemma 3.4, conclude $Q \vdash \text{Prov}(\ulcorner G \urcorner)$. Furthermore, we obtain $Q \vdash \neg \text{Prov}(\ulcorner G \urcorner)$ from $Q \vdash G \leftrightarrow \neg \text{Prov}(\ulcorner G \urcorner)$. This is contradictory since Q is consistent using the standard model.

Now, assume $Q \vdash \neg G$. We show that Q is ω -inconsistent, contradicting Axiom 5.13. We pick $\text{Prf}(\ulcorner G \urcorner, y)$ as witness for ω -inconsistency.

At first, we show that $Q \vdash \exists x. \text{Prf}(\ulcorner G \urcorner, x)$. It suffices to verify $Q \vdash \text{Prov}(\ulcorner G \urcorner)$ by definition of Prf . By Fact 5.11, it is enough to show $Q \vdash_c \text{Prov}(\ulcorner G \urcorner)$. This fact is applicable since $\text{Prov}(\ulcorner G \urcorner)$ is a Σ_1 -sentence. By the contradiction rule, we are left to prove

$$Q, \neg \text{Prov}(\ulcorner G \urcorner) \vdash_c \perp,$$

for which

$$Q, \neg \text{Prov}(\ulcorner G \urcorner) \vdash \perp$$

suffices. We apply our assumption $Q \vdash \neg G$ and end up in the proof state

$$Q, \neg \text{Prov}(\ulcorner G \urcorner) \vdash G.$$

Due to $Q \vdash G \leftrightarrow \neg \text{Prov}(\ulcorner G \urcorner)$ obtained from the diagonal lemma, we are done.

Next, we verify that $Q \vdash \neg \text{Prf}(\ulcorner G \urcorner, \bar{n})$ for all $n : \mathbb{N}$. Now that $\text{Prf}(x, y)$ is Δ_1 and $\ulcorner G \urcorner$ as well as \bar{n} are closed as numerals, we get $Q \vdash \text{Prf}(\ulcorner G \urcorner, \bar{n})$ or $Q \vdash \neg \text{Prf}(\ulcorner G \urcorner, \bar{n})$. The latter case proves our goal. For the former case, we employ falsity elimination. With Q being consistent, we conclude $Q \not\vdash G$ from $Q \vdash \neg G$. From $Q \vdash \text{Prf}(\ulcorner G \urcorner, \bar{n})$, we derive $Q \vdash \exists x. \text{Prf}(\ulcorner G \urcorner, x)$, which is equivalent to $Q \vdash \text{Prov}(\ulcorner G \urcorner)$. By Lemma 3.4, we acquire a derivation of $Q \vdash G$, clashing with our assumption $Q \not\vdash G$. ■

Remark 5.14 Let \mathcal{T} be a consistent and ω -consistent extension of Q admitting a formula $\varphi(x)$ satisfying $\forall \tau : \mathbb{F}. \mathcal{T} \vdash \tau \leftrightarrow \mathcal{T} \vdash \varphi(\ulcorner \tau \urcorner)$. The same proof shows Gödel's incompleteness theorem for \mathcal{T} . We underline this claim's lack of verification.

If $\lambda \varphi. \mathcal{T} \vdash \varphi$ is enumerable, I postulate the existence of such a formula with a slightly generalised version of the weak representability theorem, but it is not fully certain whether this approach works.

Remark 5.15 In the proof above, we have only used the ω -consistency for the Σ_1 -formula $\exists x. \text{Prf}(\ulcorner G \urcorner, x)$. Using the lemma `Sigma1_witness` from the library, we can show ω -consistency of Q for such a special kind of formulas. In consequence, we can show Gödel's theorem for Q without relying on more than CT_Q being true.

In the proof, we used many assumptions requiring sophisticated arguments. Most prominently, we used that $\text{Prov}(x)$ is of the form $\exists y. \text{Prf}(x, y)$ and we relied on ω -consistency of Q on Σ_1 -formulas. Also, we used Σ_1 -conservativity, although this may be obsolete. Still, this proof is important since many textbooks, such as [Smi22] and [BBJ07] use a similar argumentation showing the claim.

We can, however, employ Fact 1.4 to make the proof significantly easier, which is what we will do next. At first, we need an adapted provability formula.

Lemma 5.16 There exists a Σ_1 -formula $\text{SProv}(x)$ satisfying

$$\forall \varphi : \mathbb{F}. Q \vdash \varphi \rightarrow Q \vdash \text{SProv}(\ulcorner \varphi \urcorner) \wedge Q \vdash \neg \varphi \rightarrow Q \vdash \neg \text{SProv}(\ulcorner \varphi \urcorner).$$

Proof We apply Fact 1.4. Since $\lambda\varphi.Q \vdash \varphi$ is enumerable by standard techniques (cf. [FKS19]), $\lambda\varphi.Q \vdash \dot{\neg}\varphi$ can be shown enumerable, too. We conclude that both $P := \lambda n.Q \vdash \text{göd}^{-1}(n)$ and $R := \lambda n.Q \vdash \dot{\neg}\text{göd}^{-1}(n)$ are also enumerable. As \mathbb{N} is discrete, we obtain that P and R are semi-decidable. By virtue of Q 's consistency, P and R are disjoint. We obtain a Σ_1 -formula $S\text{Prov}(x)$, such that

$$\forall n : \mathbb{N}. Pn \rightarrow Q \vdash S\text{Prov}(\bar{n}) \wedge Rn \rightarrow Q \vdash \dot{\neg}S\text{Prov}(\bar{n}).$$

Now, suppose $\varphi : \mathbb{F}$. Setting $n = \text{göd}(\varphi)$ in the result from above yields

$$P(\text{göd}(\varphi)) \rightarrow Q \vdash S\text{Prov}(\ulcorner\varphi\urcorner) \wedge R(\text{göd}(\varphi)) \rightarrow Q \vdash \dot{\neg}S\text{Prov}(\ulcorner\varphi\urcorner).$$

We are done since $P(\text{göd}(\varphi)) = Q \vdash \varphi$ and $R(\text{göd}(\varphi)) = Q \vdash \dot{\neg}\varphi$. ■

Remark 5.17 While we have $Q \vdash \text{Prov}(\ulcorner\varphi\urcorner) \rightarrow Q \vdash \varphi$ for any formula φ , the implication $Q \vdash S\text{Prov}(\ulcorner\varphi\urcorner) \rightarrow Q \vdash \varphi$ has not been inspected yet.

We are now ready to present a much shorter proof of Gödel's first incompleteness theorem.

Proof (of Theorem 5.10, alternative) We apply the diagonal lemma to $\dot{\neg}S\text{Prov}(x)$ and obtain a sentence G such that $Q \vdash G \leftrightarrow \dot{\neg}S\text{Prov}(\ulcorner G \urcorner)$. We will show: $Q \not\vdash G$ and $Q \not\vdash \dot{\neg}G$.

Suppose $Q \vdash G$. Our reasoning is the same as in our first proof. By Lemma 5.16, we obtain $Q \vdash S\text{Prov}(\ulcorner G \urcorner)$. From $Q \vdash G \leftrightarrow \dot{\neg}S\text{Prov}(\ulcorner G \urcorner)$ and $Q \vdash G$, we learn $Q \vdash \dot{\neg}S\text{Prov}(\ulcorner G \urcorner)$. Contradiction to Q 's consistency.

Next, let a derivation of $Q \vdash \dot{\neg}G$ be given. From Lemma 5.16, we get $Q \vdash \dot{\neg}S\text{Prov}(\ulcorner G \urcorner)$. Using $Q \vdash G \leftrightarrow \dot{\neg}S\text{Prov}(\ulcorner G \urcorner)$, we obtain $Q \vdash G$. Again, this clashes with consistency of Q . ■

6 Discussion

In some books, the diagonalisation of a formula is defined differently, for instance in Boolos' Burgess' and Jeffrey's book [BBJ07]. Instead of picking $\varphi(\ulcorner\varphi\urcorner)$ directly, the formula $\dot{\exists}x.x \equiv \ulcorner\varphi\urcorner \wedge \varphi(x)$ is chosen. Clearly, both are logically equivalent, but the latter imposes technical overhead since we need additional rewrites when using it in proofs. Still, the second formulation does not use any substitutions. However, we opted for the first formulation; this also used in the literature, for instance by Smith [Smi22].

Both authors also use a property similar to our CT_Q -assumption. Since they are not working in constructive type theory, they impose computability conditions on the underlying function f : Both authors arrive at a formulation stating that $\varphi(x, y)$ captures f in Q if $Q \vdash \dot{\forall}y.\varphi(\bar{x}, y) \leftrightarrow y \equiv \bar{z}$, provided that $f(x) = z$. While Smith requires f to be primitive recursive and therefore total, we can always replace z with $f(x)$ and obtain our notion of CT_Q . Boolos et al., on the contrary, allow f to be an arbitrary recursive function, so $f(x)$ does not need to be defined for all x . Note that this notion of representability asserts $f(x) = z \rightarrow Q \vdash \dot{\forall}y.\varphi(\bar{x}, y) \leftrightarrow y \equiv \bar{z}$ but does not enforce the converse $Q \vdash \dot{\forall}y.\varphi(\bar{x}, y) \leftrightarrow y \equiv \bar{z} \rightarrow f(x) = z$ to hold true. Our definition of partial CT_Q requires both directions. Still, this does not make any difference in our setting as we only work with total functions.

Lastly, CT_Q states that the representing formula is Σ_1 . This coincides with Smith's result, up to minor differences in the definition of Σ_1 : Instead of being a Δ_1 (cf. [KP23, definition 23]; this property is semantic) formula preceded by existential quantifiers, Smith's Σ_1 -formulas are those that

do not use quantification² (he calls them Δ_0 formulas; this property is purely syntactic) preceded by existentials. Boolos et al. are even stronger and claim that the representing formula is of the form $\exists x.\tau(x)$ with $\tau(x)$ being Δ_0 . Hermes and Kirst already showed a somewhat related theorem: Σ_1 -compression, cf. [HK23].

Any Δ_0 -sentence (that is, a Δ_0 formula with closed terms substituted in for the free variables) can be decided within Q (cf. [Smi22, theorem 23]), so any Δ_0 -formula is a Δ_1 -formula in our sense.

References

- [BBJ07] George S. Boolos, John P. Burgess and Richard C. Jeffrey. **Computability and Logic**. 5. edition. Cambridge University Press, 2007.
- [FKS19] Yannick Forster, Dominik Kirst and Gert Smolka. ‘On synthetic undecidability in coq, with an application to the entscheidungsproblem’. In: **Proceedings of the 8th ACM SIGPLAN International Conference on Certified Programs and Proofs, CPP 2019, Cascais, Portugal, January 14-15, 2019**. Ed. by Assia Mahboubi and Magnus O. Myreen. ACM, 2019, pp. 38–51. DOI: [10.1145/3293880.3294091](https://doi.org/10.1145/3293880.3294091). URL: <https://doi.org/10.1145/3293880.3294091>.
- [HK23] Marc Hermes and Dominik Kirst. ‘An Analysis of Tennenbaum’s Theorem in Constructive Type Theory’. In: **CoRR** abs/2302.14699 (2023). DOI: [10.48550/ARXIV.2302.14699](https://doi.org/10.48550/ARXIV.2302.14699). arXiv: [2302.14699](https://arxiv.org/abs/2302.14699). URL: <https://doi.org/10.48550/arXiv.2302.14699>.
- [KH23] Dominik Kirst and Marc Hermes. ‘Synthetic Undecidability and Incompleteness of First-Order Axiom Systems in Coq’. In: **J. Autom. Reason.** 67.1 (2023), p. 13. DOI: [10.1007/s10817-022-09647-X](https://doi.org/10.1007/s10817-022-09647-X). URL: <https://doi.org/10.1007/s10817-022-09647-X>.
- [KP23] Dominik Kirst and Benjamin Peters. ‘Gödel’s Theorem Without Tears - Essential Incompleteness in Synthetic Computability’. In: **31st EACSL Annual Conference on Computer Science Logic (CSL 2023)**. Ed. by Bartek Klin and Elaine Pimentel. Vol. 252. Leibniz International Proceedings in Informatics (LIPIcs). Dagstuhl, Germany: Schloss Dagstuhl – Leibniz-Zentrum für Informatik, 2023, 30:1–30:18. ISBN: 978-3-95977-264-8. DOI: [10.4230/LIPIcs.CSL.2023.30](https://doi.org/10.4230/LIPIcs.CSL.2023.30). URL: <https://drops-dev.dagstuhl.de/entities/document/10.4230/LIPIcs.CSL.2023.30>.
- [Nor18] Michael Norrish. **LSS 2018: Computability and Incompleteness. Gödel’s First Incompleteness Theorem, Tarski’s Indefinability of Truth**. 2018. URL: <https://comp.anu.edu.au/lss/lectures/2023/lss-computability-4.pdf>.
- [Smi22] Peter Smith. **Gödel Without (Too Many) Tears**. 2. edition. Logic Matters, Cambridge, 2022. URL: <https://www.logicmatters.net/resources/pdfs/GWT2edn.pdf>.
- [Smo24] Gert Smolka. **Modeling and Proving in Computational Type Theory Using the Coq Proof Assistant**. Textbook under construction. 2024. URL: <https://www.ps.uni-saarland.de/~smolka/drafts/mpctt.pdf>.

²Strictly speaking, bounded quantification is allowed, but as this can always be recast as sequence of conjunctions or disjunctions, we can omit quantifiers at all.