

# Löb's Theorem

## A functional pearl of dependently typed quining

Anonymous

**Keywords** Agda, Löb, quine, self-reference

### Abstract

Löb's theorem states that to prove that a proposition is provable, it is sufficient to prove the proposition under the assumption that it is provable. The Curry-Howard isomorphism identifies formal proofs with abstract syntax trees of programs; Löb's theorem thus states that self-interpreters are impossible for total languages. We formalize a few variations of Löb's theorem in Agda using an inductive-inductive encoding of terms indexed over types. We verify the consistency of our formalizations relative to Agda by giving them semantics via interpretation functions.

*If P's answer is 'Bad!', Q will suddenly stop.  
But otherwise, Q will go back to the top,  
and start off again, looping endlessly back,  
till the universe dies and turns frozen and black.*

Excerpt from *Scooping the Loop Snooper: A proof that the Halting Problem is undecidable* (Pullum 2000))

### TODO

- cite Using Reflection to Explain and Enhance Type Theory?

### 1. Introduction

Löb's theorem has a variety of applications, from providing an induction rule for program semantics involving a “later” operator (Appel et al. 2007), to proving incompleteness of a logical theory as a trivial corollary, from acting as a no-go theorem for a large class of self-interpreters (TODO: mention  $F_\omega$ ?), to allowing robust cooperation in the Prisoner's Dilemma with Source Code (Barasz et al. 2014), and even in one case curing social anxiety (Yudkowsky 2014).

TODO: Talk about what's special about this paper earlier. Maybe here? Maybe a bit further down?

“What is Löb's theorem, this versatile tool with wondrous applications?” you may ask.

Consider the sentence “if this sentence is true, then you, dear reader, are the most awesome person in the world.” Suppose that this sentence is true. Then you, dear reader are the most awesome person in the world. Since this is exactly what the sentence asserts,

the sentence is true, and you, dear reader, are the most awesome person in the world. For those more comfortable with symbolic logic, we can let  $X$  be the statement “you, dear reader, are the most awesome person in the world”, and we can let  $A$  be the statement “if this sentence is true, then  $X$ ”. Since we have that  $A$  and  $A \rightarrow X$  are the same, if we assume  $A$ , we are also assuming  $A \rightarrow X$ , and hence we have  $X$ , and since assuming  $A$  yields  $X$ , we have that  $A \rightarrow X$ . What went wrong?<sup>1</sup>

It can be made quite clear that something is wrong; the more common form of this sentence is used to prove the existence of Santa Claus to logical children: considering the sentence “if this sentence is true, then Santa Claus exists”, we can prove that Santa Claus exists. By the same logic, though, we can prove that Santa Claus does not exist by considering the sentence “if this sentence is true, then Santa Claus does not exist.” Whether you consider it absurd that Santa Claus exist, or absurd that Santa Claus not exist, surely you will consider it absurd that Santa Claus both exist and not exist. This is known as Curry's paradox.

Have you figured out what went wrong?

The sentence that we have been considering is not a valid mathematical sentence. Ask yourself what makes it invalid, while we consider a similar sentence that is actually valid.

Now consider the sentence “if this sentence is provable, then you, dear reader, are the most awesome person in the world.” Fix a particular formalization of provability (for example, Peano Arithmetic, or Martin-Löf Type Theory). To prove that this sentence is true, suppose that it is provable. We must now show that you, dear reader, are the most awesome person in the world. *If provability implies truth*, then the sentence is true, and then you, dear reader, are the most awesome person in the world. Thus, if we can assume that provability implies truth, then we can prove that the sentence is true. This, in a nutshell, is Löb's theorem: to prove  $X$ , it suffices to prove that  $X$  is true whenever  $X$  is provable. Symbolically, this is

$$\Box(\Box X \rightarrow X) \rightarrow \Box X$$

where  $\Box X$  means “ $X$  is provable” (in our fixed formalization of provability).

Let us now return to the question we posed above: what went wrong with our original sentence? The answer is that self-reference with truth is impossible, and the clearest way I know to argue for this is via the Curry-Howard Isomorphism; in a particular technical sense, the problem is that self-reference with truth fails to terminate.

The Curry-Howard Isomorphism establishes an equivalence between types and propositions, between (well-typed, terminating, functional) programs and proofs. See Table 1 for some examples. Now we ask: what corresponds to a formalization of provability?

<sup>1</sup>Those unfamiliar with conditionals should note that the “if ... then ...” we use here is the logical “if”, where “if false then  $X$ ” is always true, and not the counter-factual “if”.

Logic	Programming	Set Theory
Proposition	Type	Set of Proofs
Proof	Program	Element
Implication ( $\rightarrow$ )	Function ( $\rightarrow$ )	Function
Conjunction ( $\wedge$ )	Pairing ( $.$ )	Cartesian Product ( $\times$ )
Disjunction ( $\vee$ )	Sum ( $+$ )	Disjoint Union ( $\sqcup$ )
Gödel codes	ASTs	—

**Table 1.** The Curry-Howard isomorphism between mathematical logic and functional programming

If a proof of  $P$  is a terminating functional program which is well-typed at the type corresponding to  $P$ , and to assert that  $P$  is provable is to assert that the type corresponding to  $P$  is inhabited, then an encoding of a proof is an encoding of a program. Although mathematicians typically use Gödel codes to encode propositions and proofs, a more natural choice of encoding programs will be abstract syntax trees. In particular, a valid syntactic proof of a given (syntactic) proposition corresponds to a well-typed syntax tree for an inhabitant of the corresponding syntactic type.

Unless otherwise specified, we will henceforth consider only well-typed, terminating programs; when we say “program”, the adjectives “well-typed” and “terminating” are implied.

Before diving into Löb’s theorem in detail, we’ll first visit a standard paradigm for formalizing the syntax of dependent type theory. (TODO: Move this?)

## 2. Quines

What is the computational equivalent of the sentence “If this sentence is provable, then  $X$ ”? It will be something of the form “ $??? \rightarrow X$ ”. As a warm-up, let’s look at a Python program that returns a string representation of this type.

To do this, we need a program that outputs its own source code. There are three genuinely distinct solutions, the first of which is degenerate, and the second of which is cheeky (or sassy?). These “cheating” solutions are:

- The empty program, which outputs nothing.
- The code `print(open(__file__, 'r').read())`, which relies on the Python interpreter to get the source code of the program.

Now we develop the standard solution. At a first gloss, it looks like:

```
(lambda T: '(' + T + ') -> X') "???"
```

Now we need to replace “ $???$ ” with the entirety of this program code. We use Python’s string escaping function (`repr`) and replacement syntax (`("foo %s bar" % "baz")` becomes `"foo baz bar"`):

```
(lambda T: '(' + T % repr(T) + ') -> X')
  ("(lambda T: '(' + T %% repr(T) + ') -> X')\n (%s)")
```

This is a slight modification on the standard way of programming a quine, a program that outputs its own source-code.

Suppose we have a function  $\square$  that takes in a string representation of a type, and returns the type of syntax trees of programs producing that type. Then our Löbian sentence would look something like (if  $\rightarrow$  were valid notation for function types in Python)

```
(lambda T:  $\square$  (T % repr(T)) -> X)
  ("(lambda T:  $\square$  (T %% repr(T)) -> X)\n (%s)")
```

Now, finally, we can see what goes wrong when we consider using “if this sentence is true” rather than “if this sentence is provable”.

Provability corresponds to syntax trees for programs; truth corresponds to execution of the program itself. Our pseudo-Python thus becomes

```
(lambda T: eval(T % repr(T)) -> X)
  ("(lambda T: eval(T %% repr(T)) -> X)\n (%s)")
```

This code never terminates! So, in a quite literal sense, the issue with our original sentence was that, if we tried to phrase it, we’d never finish.

Note well that the type  $(\square X \rightarrow X)$  is a type that takes syntax trees and evaluates them; it is the type of an interpreter. (TODO: maybe move this sentence?)

## 3. Abstract Syntax Trees for Dependent Type Theory

The idea of formalizing a type of syntax trees which only permits well-typed programs is common in the literature. (TODO: citations) For example, here is a very simple (and incomplete) formalization with  $\Pi$ , a unit type ( $\top$ ), an empty type ( $\perp$ ), and `lambda`. (TODO: FIXME: What’s the right level of simplicity?) TODO: mention convention of “?”

We will use some standard data type declarations, which are provided for completeness in Appendix A.

```
mutual
  infixl 2  $\triangleright$  _

data Context : Set where
   $\varepsilon$  : Context
   $\triangleright$  _ : ( $\Gamma$  : Context) -> Type  $\Gamma$  -> Context

data Type : Context -> Set where
  'T' :  $\forall \{ \Gamma \} \rightarrow$  Type  $\Gamma$ 
  ' $\perp$ ' :  $\forall \{ \Gamma \} \rightarrow$  Type  $\Gamma$ 
  'II' :  $\forall \{ \Gamma \}$ 
    -> (A : Type  $\Gamma$ ) -> Type ( $\Gamma \triangleright A$ ) -> Type  $\Gamma$ 
```

```
data Term : { $\Gamma$  : Context} -> Type  $\Gamma$  -> Set where
  'tt' :  $\forall \{ \Gamma \} \rightarrow$  Term { $\Gamma$ } 'T'
  ' $\lambda$ ' :  $\forall \{ \Gamma A B \} \rightarrow$  Term { $\Gamma \triangleright A$ } B -> Term ('II' A B)
```

An easy way to check consistency of a syntactic theory which is weaker than the theory of the ambient proof assistant is to define an interpretation function, also commonly known as an unquoter, or a denotation function, from the syntax into the universe of types. Here is an example of such a function:

```
mutual
  [ ]c : Context -> Set
  [  $\varepsilon$  ]c =  $\top$ 
  [  $\Gamma \triangleright T$  ]c =  $\Sigma$  [  $\Gamma$  ]c [ T ]T

  [ ]T :  $\forall \{ \Gamma \}$ 
    -> Type  $\Gamma$  -> [  $\Gamma$  ]c -> Set
  [ 'T' ]T [  $\Gamma$  ] =  $\top$ 
  [ ' $\perp$ ' ]T [  $\Gamma$  ] =  $\perp$ 
  [ 'II' A B ]T [  $\Gamma$  ] = ( $x$  : [ A ]T [  $\Gamma$  ]) -> [ B ]T ([  $\Gamma$  ], x)

  [ ]t :  $\forall \{ \Gamma T \}$ 
    -> Term { $\Gamma$ } T -> ([  $\Gamma$  ] : [  $\Gamma$  ]c) -> [ T ]T [  $\Gamma$  ]
  [ 'tt' ]t [  $\Gamma$  ] = tt
  [ ' $\lambda$  f' ]t [  $\Gamma$  ] x = [ f ]t ([  $\Gamma$  ], x)
```

TODO: Maybe mention something about the denotation function being “local”, i.e., not needing to do anything but the top-level case-analysis?

## 4. This Paper

In this paper, we make extensive use of this trick for validating models. We formalize the simplest syntax that supports Löb’s theorem and prove it sound relative to Agda in 12 lines of code; the understanding is that this syntax could be extended to support basically anything you might want. We then present an extended version of this solution, which supports enough operations that we can prove our syntax sound (consistent), incomplete, and nonempty. In a hundred lines of code, we prove Löb’s theorem under the assumption that we are given a quine; this is basically the well-typed functional version of the program that uses `open(__file__, 'r').read()`. Finally, we sketch our implementation of Löb’s theorem (code in an appendix) based on the assumption only that we can add a level of quotation to our syntax tree; this is the equivalent of letting the compiler implement `repr`, rather than implementing it ourselves. We close with an application to the prisoner’s dilemma, as well as some discussion about avenues for removing the hard-coded `repr`. **TODO: Ensure that this ordering is accurate**

## 5. Prior Work

**TODO: Use of Löb’s theorem in program logic as an induction principle? (TODO)**

**TODO: Brief mention of Lob’s theorem in Haskell / elsewhere / ? (TODO)**

## 6. Trivial Encoding

We begin with a language that supports almost nothing other than Löb’s theorem. We use  $\Box T$  to denote the type of Terms of whose syntactic type is  $T$ . We use  $\ulcorner T \urcorner$  to denote the syntactic type corresponding to the type of (syntactic) terms whose syntactic type is  $T$ . **TODO: This is probably unclear. Maybe mention repr?**

```
data Type : Set where
  '→' : Type → Type → Type
  '□' : Type → Type
```

```
data □ : Type → Set where
  Löb : ∀ {X} → □ ('□' X '→' X) → □ X
```

The only term supported by our term language is Löb’s theorem. We can prove this language consistent relative to Agda with an interpreter:

```
[_]ᵀ : Type → Set
[A '→' B]ᵀ = [A]ᵀ → [B]ᵀ
['□' T]ᵀ = □ T
```

```
[_]ᵀ : ∀ {T : Type} → □ T → [T]ᵀ
[Löb □ 'X' → X]ᵀ = [□ 'X' → X]ᵀ (Löb □ 'X' → X)
```

To interpret Löb’s theorem applied to the syntax for a compiler  $f$  ( $\Box 'X' \rightarrow X$  in the code above), we interpret  $f$ , and then apply this interpretation to the constructor Löb applied to  $f$ .

Finally, we tie it all together:

```
lōb : ∀ {X'} → □ ('□' X' '→' X') → [X']ᵀ
lōb f = [Löb f]ᵀ
```

This code is deceptively short, with all of the interesting work happening in the interpretation of Löb.

What have we actually proven, here? It may seem as though we’ve proven absolutely nothing, or it may seem as though we’ve proven that Löb’s theorem always holds. Neither of these is the case. The latter is ruled out, for example, by the existence of an self-interpreter for  $F_\omega$  (Brown and Palsberg 2016).<sup>2</sup>

<sup>2</sup>One may wonder how exactly the self-interpreter for  $F_\omega$  does not contradict this theorem. In private conversations with Matt Brown, we found that

We have proven the following. Suppose you have a formalization of type theory which has a syntax for types, and a syntax for terms indexed over those types. If there is a “local explanation” for the system being sound, i.e., an interpretation function where each rule does not need to know about the full list of constructors, then it is consistent to add a constructor for Löb’s theorem to your syntax. This means that it is impossible to contradict Löb’s theorem no matter what (consistent) constructors you add. We will see in the next section how this gives incompleteness, and discuss in later sections how to *prove Löb’s theorem*, rather than simply proving that it is consistent to assume.

## 7. Encoding with Soundness, Incompleteness, and Non-Emptiness

By augmenting our representation with top ( $\top$ ) and bottom ( $\perp$ ) types, and a unique inhabitant of  $\top$ , we can prove soundness, incompleteness, and non-emptiness.

```
data Type : Set where
  '→' : Type → Type → Type
  '□' : Type → Type
  '⊤' : Type
  '⊥' : Type
```

```
-- "□" is sometimes written as "Term"
data □ : Type → Set where
  Löb : ∀ {X} → □ ('□' X '→' X) → □ X
  'tt' : □ '⊤'
```

```
[_]ᵀ : Type → Set
[A '→' B]ᵀ = [A]ᵀ → [B]ᵀ
['□' T]ᵀ = □ T
['⊤']ᵀ = ⊤
['⊥']ᵀ = ⊥
```

```
[_]ᵀ : ∀ {T : Type} → □ T → [T]ᵀ
[Löb □ 'X' → X]ᵀ = [□ 'X' → X]ᵀ (Löb □ 'X' → X)
['tt']ᵀ = tt
```

```
¬_ : Set → Set
¬ T = T → ⊥
```

```
'¬' : Type → Type
'¬' T = T '→' '⊥'
```

```
lōb : ∀ {X'} → □ ('□' X' '→' X') → [X']ᵀ
lōb f = [Löb f]ᵀ
```

```
incompleteness : ¬ □ ('¬' ('□' '⊥'))
incompleteness = lōb
```

```
soundness : ¬ □ '⊥'
soundness x = [x]ᵀ
```

```
non-emptiness : □ '⊤'
non-emptiness = 'tt'
```

the  $F_\omega$  self-interpreter does not have a separate syntax for types, instead indexing its terms over types in the metalanguage. This means that the type of Löb’s theorem becomes either  $\Box (\Box X \rightarrow X) \rightarrow \Box X$ , which is not strictly positive, or  $\Box (X \rightarrow X) \rightarrow \Box X$ , which, on interpretation, must be filled with a general fixpoint operator. Such an operator is well-known to be inconsistent.

```
no-interpreters :  $\neg (\forall \{X'\} \rightarrow \Box (' \Box ' X' \rightarrow ' X'))$ 
no-interpreters interp = löb (interp { ' ' })
```

What is this incompleteness theorem? **TODO:** Incorporate this: Let's banish "truth". Sometimes it is useful to formalize a notion of provability. For example, you might want to show neither assuming  $T$  nor assuming  $\neg T$  yields a proof of contradiction. You cannot phrase this is  $\neg T \wedge \neg \neg T$ , for that is absurd. Instead, you want to say something like  $(\neg \Box T) \wedge \neg \Box (\neg T)$ , i.e., it would be absurd to have a proof object of either  $T$  or of  $\neg T$ . After a while, you might find that meta-programming in this formal syntax is nice, and you might want it to be able to formalize every proof, so that you can do all of your solving reflectively. If you're like me, you might even want to reason about the reflective tactics themselves in a reflective manner; you'd want to be able to add levels of quotation to quoted things to talk about such tactics. The incompleteness theorem says that this isn't possible. For any fixed language of syntactic proofs which is powerful enough to represent itself, there will always be some valid proofs that you cannot reflect into your syntax. In particular, you might be able to prove that your syntax has no proofs of  $\perp$  (by interpreting any such proof). But you'll be unable to quote that proof. This is what the incompleteness theorem that I stated says. (As I understand it, incompleteness, fundamentally, is a result about the limitations of formalizing provability.)

**TODO:** Does this code need any explanation? Maybe for no-interpreters?

## 8. Encoding with Quines

We now weaken our assumptions further. Rather than assuming Löb's theorem, we instead assume only a type-level quine in our representation. Recall that a *quine* is a program that outputs some function of its own source code. A *type-level quine at  $\phi$*  is program that outputs the result of evaluating the function  $\phi$  on the abstract syntax tree of its own type. Letting `Quine  $\phi$`  denote the constructor for a type-level quine at  $\phi$ , we have an isomorphism between `Quine  $\phi$`  and  $\phi \ulcorner \text{Quine } \phi \urcorner^\top$ , where  $\ulcorner \text{Quine } \phi \urcorner^\top$  is the abstract syntax tree for the source code of `Quine  $\phi$` . Note that we assume constructors for "adding a level of quotation", turning abstract syntax trees for programs of type  $T$  into abstract syntax trees for abstract syntax trees for programs of type  $T$ ; this corresponds to `repr`.

```
infixl 3 _ 'a_'
infixl 3 _w' 'a_'
infixl 3 _ ' '
infixl 2 _▷_
infixr 2 _ 'o'_
infixr 1 _ '→'_
```

We begin with an encoding of contexts and types, repeating from above the constructors of `'→'`, `'□'`, `'⊤'`, and `'⊥'`. We add to this a constructor for quines (`Quine`), and a constructor for syntax trees of types in the empty context (`'Typeε'`). Finally, rather than proving weakening and substitution as mutually recursive definitions, we take the easier but more verbose route of adding constructors that allow adding and substituting extra terms in the context. **TODO:** (McBride 2010) Note that `'□'` is now a function of the represented language, rather than a meta-level operator **TODO:** Does this need more explanation?

```
mutual
data Context : Set where
  ε : Context
  _▷_ : (Γ : Context) → Type Γ → Context

data Type : Context → Set where
```

```
'→' :  $\forall \{\Gamma\} \rightarrow \text{Type } \Gamma \rightarrow \text{Type } \Gamma \rightarrow \text{Type } \Gamma$ 
'⊤' :  $\forall \{\Gamma\} \rightarrow \text{Type } \Gamma$ 
'⊥' :  $\forall \{\Gamma\} \rightarrow \text{Type } \Gamma$ 
'Typeε' :  $\forall \{\Gamma\} \rightarrow \text{Type } \Gamma$ 
'□' :  $\forall \{\Gamma\} \rightarrow \text{Type } (\Gamma \triangleright \text{'Typeε'})$ 
Quine :  $\text{Type } (\varepsilon \triangleright \text{'Typeε'}) \rightarrow \text{Type } \varepsilon$ 
W :  $\forall \{\Gamma A\} \rightarrow \text{Type } \Gamma \rightarrow \text{Type } (\Gamma \triangleright A)$ 
W1 :  $\forall \{\Gamma A B\} \rightarrow \text{Type } (\Gamma \triangleright B) \rightarrow \text{Type } (\Gamma \triangleright A \triangleright (W B))$ 
_ ' ' :  $\forall \{\Gamma A\} \rightarrow \text{Type } (\Gamma \triangleright A) \rightarrow \text{Term } A \rightarrow \text{Type } \Gamma$ 
```

In addition to `'λ'` and `'tt'`, we now have the AST-equivalents of Python's `repr`, which we denote as  $\ulcorner \urcorner^\top$  for the type-level add-quote function **TODO: should this be called add-quote?**, and  $\ulcorner \urcorner^\top$  for the term-level add-quote function. We add constructors `quine→` and `quine←` that exhibit the isomorphism that defines our type-level quine constructor, though we elide a constructor declaring that these are inverses, as we find it unnecessary.

To construct the proof of Löb's theorem, we need a few other standard constructors, such as `'VAR0'`, which references a term in the context; `_ 'a_'`, which we use to denote function application; `_ 'o'_`, a function composition operator; and `'□' ' ' ' T ⊤'`, the variant of `'VAR0'` which adds an extra level of syntax-trees. We also include a number of constructors that handle weakening and substitution; this allows us to avoid both inductive-recursive definitions of weakening and substitution, and defining a judgmental equality or conversion relation.

```
data Term : {Γ : Context} → Type Γ → Set where
  'λ' :  $\forall \{\Gamma A B\} \rightarrow \text{Term } \{\Gamma \triangleright A\} (W B) \rightarrow \text{Term } (A \rightarrow B)$ 
  'tt' :  $\forall \{\Gamma\} \rightarrow \text{Term } \{\Gamma\} \text{'⊤'}$ 
   $\ulcorner \urcorner^\top$  :  $\forall \{\Gamma\} \rightarrow \text{Type } \varepsilon \rightarrow \text{Term } \{\Gamma\} \text{'Typeε'}$ 
   $\ulcorner \urcorner^\top$  :  $\forall \{\Gamma T\} \rightarrow \text{Term } \{\varepsilon\} T \rightarrow \text{Term } \{\Gamma\} (\ulcorner \urcorner^\top \text{' ' } T \urcorner^\top)$ 
  quine→ :  $\forall \{\phi\} \rightarrow \text{Term } \{\varepsilon\} (\text{Quine } \phi \rightarrow \phi \ulcorner \text{Quine } \phi \urcorner^\top)$ 
  quine← :  $\forall \{\phi\} \rightarrow \text{Term } \{\varepsilon\} (\phi \ulcorner \text{Quine } \phi \urcorner^\top \rightarrow \text{Quine } \phi)$ 
  'VAR0' :  $\forall \{\Gamma T\} \rightarrow \text{Term } \{\Gamma \triangleright T\} (W T)$ 
  _ 'a_' :  $\forall \{\Gamma A B\} \rightarrow \text{Term } \{\Gamma\} (A \rightarrow B) \rightarrow \text{Term } \{\Gamma\} A \rightarrow \text{Term } \{\Gamma\} B$ 
  _ 'o'_ :  $\forall \{\Gamma A B C\} \rightarrow \text{Term } \{\Gamma\} (B \rightarrow C) \rightarrow \text{Term } \{\Gamma\} (A \rightarrow B) \rightarrow \text{Term } \{\Gamma\} (A \rightarrow C)$ 
  '□' ' ' ' T ⊤' :  $\forall \{T\} \rightarrow \text{Term } \{\varepsilon \triangleright \ulcorner \urcorner^\top \text{' ' } T \urcorner^\top\} (W (\ulcorner \urcorner^\top \text{' ' } \ulcorner \urcorner^\top \text{' ' } T \urcorner^\top \urcorner^\top))$ 
  →SW1SV→W :  $\forall \{\Gamma T X A B\} \{x : \text{Term } \{\Gamma\} X\} \rightarrow \text{Term } (T \rightarrow (W_1 A \ulcorner \text{'VAR}_0' \rightarrow W B \urcorner^\top)) \rightarrow \text{Term } (T \rightarrow A \ulcorner x \rightarrow B \urcorner^\top)$ 
  ←SW1SV→W :  $\forall \{\Gamma T X A B\} \{x : \text{Term } \{\Gamma\} X\} \rightarrow \text{Term } ((W_1 A \ulcorner \text{'VAR}_0' \rightarrow W B \urcorner^\top) \ulcorner x \rightarrow T \urcorner^\top) \rightarrow \text{Term } ((A \ulcorner x \rightarrow B \urcorner^\top) \rightarrow T)$ 
```

```

w : ∀ {Γ A T} → Term A → Term {Γ ▷ T} (W A)
w → : ∀ {Γ A B X}
  → Term {Γ} (A '→' B)
  → Term {Γ ▷ X} (W A '→' W B)
-w''''_a : ∀ {A B T}
  → Term {ε ▷ T} (W ('□' '' Γ A '→' B ⊥T))
  → Term {ε ▷ T} (W ('□' '' Γ A ⊥T))
  → Term {ε ▷ T} (W ('□' '' Γ B ⊥T))

```

```

□ : Type ε → Set _
□ = Term {ε}

```

To verify the soundness of our syntax, we provide a model for it and an interpretation into that model. We call particular attention to the interpretation of  $\square$ , which is just  $\text{Term } \{\varepsilon\}$ ; to  $\text{Quine } \phi$ , which is the interpretation of  $\phi$  applied to  $\text{Quine } \phi$ ; and to the interpretations of the quine isomorphism functions, which are just the identity functions.

```

max-level : Level
max-level = |zero -- also works for higher levels

```

```
mutual
```

```

[ ]c : (Γ : Context) → Set (lsuc max-level)
[ ε ]c = ⊤
[ Γ ▷ T ]c = Σ [ Γ ]c [ T ]T

[ ]T : ∀ {Γ}
  → Type Γ → [ Γ ]c → Set max-level
[ A '→' B ]T [Γ] = [ A ]T [Γ] → [ B ]T [Γ]
[ '⊥' ]T [Γ] = ⊥
[ '⊥' ]T [Γ] = ⊥
[ 'Typeε' ]T [Γ] = Lifted (Type ε)
[ '□' ]T [Γ] = Lifted (Term {ε} (lower (Σ.snd [Γ])))
[ Quine φ ]T [Γ] = [ φ ]T ([Γ], lift (Quine φ))
[ W T ]T [Γ] = [ T ]T (Σ.fst [Γ])
[ W1 T ]T [Γ] = [ T ]T (Σ.fst (Σ.fst [Γ]), Σ.snd [Γ])
[ T '' x ]T [Γ] = [ T ]T ([Γ], [ x ]t [Γ])

```

```

[ ]t : ∀ {Γ T}
  → Term {Γ} T → ([Γ] : [ Γ ]c) → [ T ]T [Γ]
[ 'λ' f ]t [Γ] x = [ f ]t ([Γ], x)
[ 'tt' ]t [Γ] = tt
[ 'x ⊥T' ]t [Γ] = lift x
[ 'x ⊥T' ]t [Γ] = lift x
[ quine→ ]t [Γ] x = x
[ quine← ]t [Γ] x = x
[ 'VAR0' ]t [Γ] = Σ.snd [Γ]
[ g 'o' f ]t [Γ] x = [ g ]t [Γ] ([ f ]t [Γ] x)
[ f ''a x ]t [Γ] = [ f ]t [Γ] ([ x ]t [Γ])
[ 'ΓVAR0⊥t' ]t [Γ] = lift Γ lower (Σ.snd [Γ]) ⊥T
[ ←SW1SV→Wf ]t [Γ] = [ f ]t
[ →SW1SV→Wf ]t [Γ] = [ f ]t
[ w x ]t [Γ] = [ x ]t (Σ.fst [Γ])
[ w→f ]t [Γ] = [ f ]t (Σ.fst [Γ])
[ f w''''a x ]t [Γ]
  = lift (lower ([ f ]t [Γ]) ''a lower ([ x ]t [Γ]))

```

To prove Löb's theorem, we must create the sentence "if this sentence is provable, then X", and then provide an inhabitant of that type. We can define this sentence, which we call 'H', as the type-level quine at the function  $\lambda v. \square v \rightarrow 'X'$ . We can then convert back and forth between the types  $\square 'H'$  and  $\square 'H' \rightarrow 'X'$  with our quine isomorphism functions, and a bit

of quotation magic and function application gives us a term of type  $\square 'H' \rightarrow \square 'X'$ ; this corresponds to the inference of the provability of Santa Claus' existence from the assumption that the sentence is provable. We compose this with the assumption of Löb's theorem, that  $\square 'X' \rightarrow 'X'$ , to get a term of type  $\square 'H' \rightarrow 'X'$ , i.e., a term of type 'H'; this is the inference that when provability implies truth, Santa Claus exists, and hence that the sentence is provable. Finally, we apply this to its own quotation, obtaining a term of type  $\square 'X'$ , i.e., a proof that Santa Claus exists.

```

module inner ('X' : Type ε)
  (f' : Term {ε} ('□' '' Γ 'X' ⊥T '→' 'X'))
  where
    'H' : Type ε
    'H' = Quine (W1 '□' '' 'VAR0' '→' W 'X')

    'toH' : □ ('□' '' Γ 'H' ⊥T '→' 'X') '→' 'H'
    'toH' = ←SW1SV→W quine←

    'fromH' : □ ('H' '→' ('□' '' Γ 'H' ⊥T '→' 'X'))
    'fromH' = →SW1SV→W quine→

    '□H'→□X' : □ ('□' '' Γ 'H' ⊥T '→' '□' '' Γ 'X' ⊥T)
    '□H'→□X'
      = 'λ' (w Γ 'fromH' ⊥T
        w''''a 'VAR0'
        w''''a 'ΓVAR0⊥t)

    'h' : Term 'H'
    'h' = 'toH' ''a (f' 'o' '□H'→□X')

    Löb : □ 'X'
    Löb = 'fromH' ''a 'h' ''a Γ 'h' ⊥T

```

```

Löb : ∀ {X} → □ ('□' '' Γ X ⊥T '→' X) → □ X
Löb {X} f = inner.Löb X f

```

```

[ ] : Type ε → Set _
[ T ] = [ T ]T tt

```

```

'⊥' : ∀ {Γ} → Type Γ → Type Γ
'⊥' T = T '→' '⊥'

```

```

löp : ∀ {X} → □ ('□' '' Γ 'X' ⊥T '→' 'X') → [ 'X' ]
löp f = [ ]t (Löb f) tt

```

```

⊥_ : ∀ {ℓ m} → Set ℓ → Set (ℓ ⊔ m)
⊥_ {ℓ} {m} T = T → ⊥ {m}

```

As above, we can again prove soundness, incompleteness, and non-emptiness.

```

incompleteness : ¬ □ ('⊥' ('□' '' Γ '⊥' ⊥T))
incompleteness = löb

```

```

soundness : ¬ □ '⊥'
soundness x = [ x ]t tt

```

```

non-emptiness : Σ (Type ε) (λ T → □ T)
non-emptiness = '⊥', 'tt'

```



$B \text{ Says} \backslash A \text{ Says}$	Cooperate	Defect
Cooperate	(1 year, 1 year)	(0 years, 3 years)
Defect	(3 years, 0 years)	(2 years, 2 years)

**Table 2.** The payoff matrix for the prisoner’s dilemma; each cell contains (the years  $A$  spends in prison, the years  $B$  spends in prison).

## 9. Digression: Application of Quining to The Prisoner’s Dilemma

In this section, we use a slightly more enriched encoding of syntax; see Appendix B for details.

### 9.1 The Prisoner’s Dilemma

The Prisoner’s Dilemma is a classic problem in game theory. Two people have been arrested as suspects in a crime and are being held in solitary confinement, with no means of communication. The investigators offer each of them a plea bargain: a decreased sentence for ratting out the other person. Each suspect can then choose to either cooperate with the other suspect by remaining silent, or defect by ratting out the other suspect. The possible outcomes are summarized in Table 2.

Suspect  $A$  might reason thusly: “Suppose the other suspect cooperates with me. Then I’d get off with no prison time if I defected, while I’d have to spend a year in prison if I cooperate. Similarly, if the other suspect defects, then I’d get two years in prison for defecting, and three for cooperating. In all cases, I do better by defecting.” If suspect  $B$  reasons similarly, then both decide to defect, and both get two years in prison, despite the fact that both prefer the (Cooperate, Cooperate) outcome over the (Defect, Defect) outcome!

### 9.2 Adding Source Code

We have the intuition that if both suspects are good at reasoning, and both know that they’ll reason the same way, then they should be able to mutually cooperate. One way to formalize this is to talk about programs (rather than people) playing the prisoner’s dilemma, and to allow each program access to its own source code and its opponent’s source code (Barasz et al. 2014).

We have formalized this framework in Agda: we use ‘Bot’ to denote the type of programs that can play in such a prisoner’s dilemma; each one takes in source code for two ‘Bot’s and outputs a proposition which is true (a type which is inhabited) if and only if it cooperates with its opponent. Said another way, the output of each bot is a proposition describing the assertion that it cooperates with its opponent.

open lob

```
-- ‘Bot’ is defined as the fixed point of
-- ‘Bot’
-- ↔ (Term ‘Bot’ → Term ‘Bot’ → ‘Type’)
‘Bot’ : ∀ {Γ} → Type Γ
‘Bot’ {Γ}
= Quine (W1 ‘Term’ “ ‘VAR0’
‘→’ W1 ‘Term’ “ ‘VAR0’
‘→’ W ‘Type’ Γ))
```

To construct an executable bot, we could do a bounded search for proofs of this proposition; one useful method described in (Barasz et al. 2014) is to use Kripke frames. This computation is, however, beyond the scope of this paper.

The assertion that a bot  $b_1$  cooperates with a bot  $b_2$  is the result of interpreting the source code for the bot, and feeding the resulting function the source code for  $b_1$  and  $b_2$ .

```
-- N.B. “□” means “Term {ε}”, i.e., a term in
-- the empty context
_cooperates-with_ : □ ‘Bot’ → □ ‘Bot’ → Type ε
b1 cooperates-with b2 = lower (ll b1 tt (lift b1) (lift b2))
```

We now provide a convenience constructor for building bots, based on the definition of quines, and present four relatively simple bots: DefectBot, CooperateBot, FairBot, and PrudentBot.

```
make-bot : ∀ {Γ}
→ Term {Γ} > ‘□’ ‘Bot’ > W (‘□’ ‘Bot’)
(W (W (‘Type’ Γ)))
→ Term {Γ} ‘Bot’
make-bot t
= ←SW1SV→SW1SV→W
quine← “a ‘λ’ (→w (‘λ’ t))
```

```
‘DefectBot’ : □ ‘Bot’
‘CooperateBot’ : □ ‘Bot’
‘FairBot’ : □ ‘Bot’
‘PrudentBot’ : □ ‘Bot’
```

The first two bots are very simple: DefectBot never cooperates (the assertion that DefectBot cooperates is a contradiction), while CooperateBot always cooperates. We define these bots, and prove that DefectBot never cooperates and CooperateBot always cooperates.

```
‘DefectBot’ = make-bot (w (w ⊥ ⊥))
‘CooperateBot’ = make-bot (w (w ⊤ ⊤))
```

```
DB-defects : ∀ {b}
→ ¬ ll ‘DefectBot’ cooperates-with b ]
DB-defects {b} pf = pf
```

```
CB-cooperates : ∀ {b}
→ ll ‘CooperateBot’ cooperates-with b ]
CB-cooperates {b} = tt
```

We can do better than DefectBot, though, now that we have source code. FairBot cooperates with you if and only if it can find a proof that you cooperate with FairBot. By Löb’s theorem, to prove that FairBot cooperates with itself, it suffices to prove that if there is a proof that FairBot cooperates with itself, then FairBot does, in fact, cooperate with itself. This is obvious, though: FairBot decides whether or not to cooperate with itself by searching for a proof that it does, in fact, cooperate with itself.

To define FairBot, we first define what it means for the other bot to cooperate with some particular bot.

```
-- We can “evaluate” a bot to turn it into a
-- function accepting the source code of two
-- bots.
‘eval-bot’ : ∀ {Γ}
→ Term {Γ} (‘Bot’
‘→’ (‘□’ ‘Bot’ ‘→’ ‘□’ ‘Bot’ ‘→’ ‘Type’ Γ))
‘eval-bot’ = →SW1SV→SW1SV→W quine→
```

```
-- We can quote this, and get a function that
-- takes the source code for a bot, and
-- outputs the source code for a function that
-- takes (the source code for) that bot’s
-- opponent, and returns an assertion of
-- cooperation with that opponent
‘eval-bot’ : ∀ {Γ}
→ Term {Γ} (‘□’ ‘Bot’
```

```

    '→' '□' ({- other -} '□' 'Bot' '→' 'Type' Γ))
"eval-bot"
= 'λ' (w ⊢ "eval-bot" ⊢
  w "" "a" "VAR0"
  w "" "a" "Γ" "VAR0" ⊢)

-- The assertion "our opponent cooperates with
-- a bot b" is equivalent to the evaluation of
-- our opponent, applied to b. Most of the
-- noise in this statement is manipulation of
-- weakening and substitution.
"other-cooperates-with" : ∀ {Γ}
  → Term {Γ}
  ▷ '□' 'Bot'
  ▷ W ('□' 'Bot')
  (W (W ('□' 'Bot')) '→' W (W ('□' 'Type' Γ)))
"other-cooperates-with" {Γ}
= "eval-other"
  'o' w → (w (w → (w ('λ' "Γ" "VAR0" ⊢))))
where
  "eval-other"
  : Term {Γ ▷ '□' 'Bot' ▷ W ('□' 'Bot')}
  (W (W ('□' ('□' 'Bot' '→' 'Type' Γ)))
  "eval-other"
  = w → (w (w → (w "eval-bot"))) "a" "VAR0"

  "eval-other"
  : Term (W (W ('□' ('□' 'Bot'))
  '→' W (W ('□' 'Type' Γ)))
  "eval-other"
  = ww → (w → (w (w → (w "" "a")))) "a" "eval-other"

-- A bot gets its own source code as the first
-- argument (of two)
"self" : ∀ {Γ}
  → Term {Γ ▷ '□' 'Bot' ▷ W ('□' 'Bot')}
  (W (W ('□' 'Bot')))
"self" = w "VAR0"

-- A bot gets its opponent's source code as
-- the second argument (of two)
"other" : ∀ {Γ}
  → Term {Γ ▷ '□' 'Bot' ▷ W ('□' 'Bot')}
  (W (W ('□' 'Bot')))
"other" = "VAR0"

-- FairBot is the bot that cooperates iff its
-- opponent cooperates with it
"FairBot"
= make-bot ("□" ('other-cooperates-with' "a" "self"))

```

We now come to the final bot: PrudentBot. You do better in the prisoner's dilemma if you cooperate whenever that's required for mutual cooperation, and you defect whenever your opponent would cooperate even if you defected. PrudentBot formalizes an approximation to this intuition: PrudentBot cooperates with you if and only if it can prove that you cooperate with it, and it can prove that you defect against DefectBot (when it assumes that DefectBot does not cooperate with you).

PrudentBot defects against DefectBot. Since there is no proof of  $\perp$ , PrudentBot does not find a proof that DefectBot cooperates with it, and so it will not cooperate with DefectBot.

By Löb's theorem, PrudentBot cooperates with itself. Under the assumption that  $\perp$  is unprovable, PrudentBot can prove that it de-

fects against DefectBot. If we further assume that PrudentBot can find a proof that it cooperates with itself (which we are allowed to do when proving the hypothesis of Löb's theorem), then PrudentBot will, in fact, cooperate with itself. Hence, by Löb's theorem, we can prove that PrudentBot will cooperate with itself.

We leave the formalization of this proof to the reader, and present only the definition of PrudentBot.

```

-- Convenience notation for triply quoted
-- negation in a context with at least two
-- terms
ww "" "¬" : ∀ {Γ A B}
  → Term {Γ ▷ A ▷ B} (W (W ('□' ('Type' Γ))))
  → Term {Γ ▷ A ▷ B} (W (W ('□' ('Type' Γ))))
ww "" "¬" T = T ww "" "→" w (w ⊢ "⊥" ⊢)

-- PrudentBot cooperates if its opponent
-- cooperates with PrudentBot, and if, under
-- the assumption that ⊥ is unprovable (¬□⊥),
-- its opponent does not cooperate with
-- DefectBot
"PrudentBot"
= make-bot ("□"
  (( "other-cooperates-with" "a" "self")
  ww "" "×"
  (¬□⊥ ww "" "→" other-defects-against-DB)))
where
  other-defects-against-DB
  : Term {Γ ▷ '□' 'Bot' ▷ W ('□' 'Bot')}
  (W (W ('□' ('Type' _))))
  other-defects-against-DB
  = ww "" "¬"
  ('other-cooperates-with'
  "a" w (w ⊢ "DefectBot" ⊢))

  ¬□⊥ : ∀ {Γ A B}
    → Term {Γ ▷ A ▷ B} (W (W ('□' ('Type' Γ))))
  ¬□⊥ = w (w ⊢ "¬" ('□' '⊥') ⊢)

```

## 10. Encoding with Add-Quote Function

Now we return to our proving of Löb's theorem. Included in the artifact for this paper<sup>3</sup> is code that replaces the Quine constructor with simpler constructors. Because the lack of  $\beta$ -reduction in the syntax clouds the main points and makes the code rather verbose, we do not include the code in the paper, and instead describe the most interesting and central points.

Recall our Python quine from section 2:

```

(lambda T: □ (T % repr(T)) → X)
  ("(lambda T: □ (T %% repr(T)) → X)\n (%s)")

```

To translate this into Agda, we need to give a type to  $T$ . Clearly,  $T$  needs to be of type  $\text{Type } ???$  for some context  $???$ . Since we need to be able to substitute something into that context, we must have  $T : \text{Type } (\Gamma \triangleright ???)$ , i.e.,  $T$  must be a syntax tree for a type, with a hole in it.

What's the shape of the thing being substituted? Well, it's a syntax tree for a type with a hole in it. . . Uh-oh. Our quine's type, naïvely, is infinite!

We know of two ways to work around this. Classical mathematics, which uses Gödel codes instead of abstract syntax trees, uses an untyped representation of proofs. It's only later in the proof of

<sup>3</sup>In lob-build-quine.lagda.

Löb's theorem that a notion of a formula being “well-formed” is introduced.

Here, we describe an alternate approach. Rather than giving up types all-together, we can “box” the type of the hole, to hide it. Using `fst` and `snd` to denote projections from a  $\Sigma$  type, using  $\ulcorner A \urcorner$  to denote the abstract syntax tree for  $A$ ,<sup>4</sup> and using `%s` to denote the first variable in the context (written as ‘VAR<sub>0</sub>’ in previous formalizations above), we can write:

```
dummy : Type (ε ▷  $\ulcorner \Sigma \text{ Context Type} \urcorner$ )
repr :  $\Sigma \text{ Context Type} \rightarrow \text{Term } \{\varepsilon\} \ulcorner \Sigma \text{ Context Type} \urcorner$ 

cast-fst
  :  $\Sigma \text{ Context Type} \rightarrow \text{Type } (\varepsilon ▷ \ulcorner \Sigma \text{ Context Type} \urcorner)$ 
cast-fst (ε ▷  $\ulcorner \Sigma \text{ Context Type} \urcorner$ , T) = T
cast-fst (_, _) = dummy

LöbSentence : Type ε
LöbSentence
  = (λ (T :  $\Sigma \text{ Context Type}$ )
    → □ (cast-fst T % repr T) ‘→’ X)
    (ε ▷  $\ulcorner \Sigma \text{ Context Type} \urcorner$ 
    ,  $\ulcorner$  (λ (T :  $\Sigma \text{ Context Type}$ )
      → □ (cast-fst T % repr T) ‘→’ X)
      (%s)  $\urcorner$ )
```

In this pseudo-Agda code, `cast-fst` unboxes the sentence that it gets, and returns it if it is the right type. Since the sentence is, in fact, always the right type, what we do in the other cases doesn't matter.

Summing up, the key ingredients to this construction are:

- A type of syntactic terms indexed over a type of syntactic types (and contexts)
- Decidable equality on syntactic contexts at a particular point ( $\Sigma \text{ Context Type}$ ), with appropriate reduction on equal things
- $\Sigma$  types, projections, and appropriate reduction on their projections
- Function types
- A function `repr` which adds a level of quotation to any syntax tree
- Syntax trees for all of the above

In any formalization of dependent type theory with all of these ingredients, we can prove Löb's theorem.

## 11. Conclusion

What remains to be done is formalizing Martin-Löf type theory without assuming `repr` nor assuming a constructor for the type of syntax trees (‘Context’, ‘Type’, and ‘Term’ or ‘□’ in our formalizations), and instead support inductive types, and construct these things inductively, and by folding over the inductive definitions there-in.

**TODO:** Figure out what to say in the conclusion

**TODO:** Mention  $F_\omega$  again?

## A. Standard Data-Type Declarations

```
open import Agda.Primitive public
using (Level; _⊔_; lzero; lsuc)
```

```
infixl 1 _, _
```

<sup>4</sup>Note that  $\ulcorner \_ \urcorner$  would not be a function in the language, but a meta-level operation.

```
infixr 2 _×_
infixl 1 _≡_
```

```
record  $\top \{ \ell \} : \text{Set } \ell$  where
  constructor tt
```

```
data  $\perp \{ \ell \} : \text{Set } \ell$  where
```

```
record  $\Sigma \{ a p \} (A : \text{Set } a) (P : A \rightarrow \text{Set } p)$ 
  :  $\text{Set } (a \sqcup p)$ 
  where
  constructor _,_
  field
    fst : A
    snd : P fst
```

```
data Lifted  $\{ a b \} (A : \text{Set } a) : \text{Set } (b \sqcup a)$  where
  lift : A  $\rightarrow$  Lifted A
```

```
lower :  $\forall \{ a b A \} \rightarrow \text{Lifted } \{ a \} \{ b \} A \rightarrow A$ 
lower (lift x) = x
```

```
_×_ :  $\forall \{ \ell \ell' \} (A : \text{Set } \ell) (B : \text{Set } \ell') \rightarrow \text{Set } (\ell \sqcup \ell')$ 
A × B =  $\Sigma A (\lambda \_ \rightarrow B)$ 
```

```
data _≡_  $\{ \ell \} \{ A : \text{Set } \ell \} (x : A) : A \rightarrow \text{Set } \ell$  where
  refl : x ≡ x
```

```
sym :  $\{ A : \text{Set} \} \rightarrow \{ x : A \} \rightarrow \{ y : A \} \rightarrow x \equiv y \rightarrow y \equiv x$ 
sym refl = refl
```

```
trans :  $\{ A : \text{Set} \} \rightarrow \{ x y z : A \} \rightarrow x \equiv y \rightarrow y \equiv z \rightarrow x \equiv z$ 
trans refl refl = refl
```

```
transport :  $\forall \{ A : \text{Set} \} \{ x : A \} \{ y : A \} \rightarrow (P : A \rightarrow \text{Set})$ 
            $\rightarrow x \equiv y \rightarrow P x \rightarrow P y$ 
transport P refl v = v
```

## B. Encoding of Löb's Theorem for the Prisoner's Dilemma

```
module lob where
```

```
infixl 2  $\triangleright$ _
infixl 3  $\ulcorner$ _
infixr 1  $\urcorner$ _
infixr 1  $\ulcorner$ “→” $\urcorner$ _
infixr 1  $\ulcorner$ ww“→” $\urcorner$ _
infixl 3  $\ulcorner$ a $\urcorner$ _
infixl 3  $\ulcorner$ w“” $\urcorner$ a $\urcorner$ _
infixr 2  $\ulcorner$ o $\urcorner$ _
infixr 2  $\ulcorner$ × $\urcorner$ _
infixr 2  $\ulcorner$ “×” $\urcorner$ _
infixr 2  $\ulcorner$ w“×” $\urcorner$ _
```

```
mutual
```

```
data Context : Set where
  ε : Context
   $\triangleright$  : (Γ : Context)  $\rightarrow$  Type Γ  $\rightarrow$  Context
```

```
data Type : Context  $\rightarrow$  Set where
  ‘T’ :  $\forall \{ \Gamma \} \rightarrow \text{Type } \Gamma$ 
  ‘⊥’ :  $\forall \{ \Gamma \} \rightarrow \text{Type } \Gamma$ 
```



```

- '→' : ∀ {Γ} → Type Γ → Type Γ → Type Γ
- '×' : ∀ {Γ} → Type Γ → Type Γ → Type Γ
- 'Type' : ∀ Γ → Type Γ
- 'Term' : ∀ {Γ} → Type (Γ ▷ 'Type' Γ)
- Quine : ∀ {Γ} → Type (Γ ▷ 'Type' Γ) → Type Γ
- W : ∀ {Γ A} → Type Γ → Type (Γ ▷ A)
- W1 : ∀ {Γ A B}
  → Type (Γ ▷ B)
  → Type (Γ ▷ A ▷ W B)
- ' ' : ∀ {Γ A}
  → Type (Γ ▷ A)
  → Term A
  → Type Γ

```

data Term : {Γ : Context} → Type Γ → Set where

```

'tt' : ∀ {Γ} → Term {Γ} 'T'
'λ' : ∀ {Γ A B}
  → Term {Γ ▷ A} (W B)
  → Term (A '→' B)
'VAR0' : ∀ {Γ T} → Term {Γ ▷ T} (W T)
- '⊢' : ∀ {Γ}
  → Type Γ
  → Term {Γ} ('Type' Γ)
- '⊢' : ∀ {Γ T}
  → Term {Γ} T
  → Term {Γ} ('Term' '' ⊢ T ⊢T)
- '⊢VAR0' : ∀ {Γ T}
  → Term {Γ ▷ 'Term' '' ⊢ T ⊢T}
  (W ('Term' '' ⊢ 'Term' '' ⊢ T ⊢T ⊢T))
- '⊢VAR0' : ∀ {Γ}
  → Term {Γ ▷ 'Type' Γ}
  (W ('Term' '' ⊢ 'Type' Γ ⊢T))
- 'a' : ∀ {Γ A B}
  → Term {Γ} (A '→' B)
  → Term {Γ} A
  → Term {Γ} B
- '×' : ∀ {Γ}
  → Term {Γ} ('Type' Γ
    '→' 'Type' Γ
    '→' 'Type' Γ)
- quine→ : ∀ {Γ φ}
  → Term {Γ}
  (Quine φ '→' φ '' ⊢ Quine φ ⊢T)
- quine← : ∀ {Γ φ}
  → Term {Γ}
  (φ '' ⊢ Quine φ ⊢T '→' Quine φ)
- SW : ∀ {Γ X A} {a : Term A}
  → Term {Γ} (W X '' a)
  → Term X
- →SW1SV→W
  : ∀ {Γ T X A B} {x : Term X}
  → Term {Γ}
  (T '→' (W1 A '' 'VAR0' '→' W B) '' x)
  → Term {Γ}
  (T '→' A '' x '→' B)
- ←SW1SV→W
  : ∀ {Γ T X A B} {x : Term X}
  → Term {Γ}
  ((W1 A '' 'VAR0' '→' W B) '' x '→' T)
  → Term {Γ}
  ((A '' x '→' B) '→' T)

```

```

→SW1SV→SW1SV→W
  : ∀ {Γ T X A B} {x : Term X}
  → Term {Γ} (T '→' (W1 A '' 'VAR0'
    '→' W1 A '' 'VAR0'
    '→' W B) '' x)
  → Term {Γ} (T '→' A '' x '→' A '' x '→' B)
- ←SW1SV→SW1SV→W
  : ∀ {Γ T X A B} {x : Term X}
  → Term {Γ} ((W1 A '' 'VAR0'
    '→' W1 A '' 'VAR0'
    '→' W B) '' x
    '→' T)
  → Term {Γ} ((A '' x '→' A '' x '→' B) '→' T)
- w : ∀ {Γ A T}
  → Term {Γ} A
  → Term {Γ ▷ T} (W A)
- w→ : ∀ {Γ A B X}
  → Term {Γ ▷ X} (W (A '→' B))
  → Term {Γ ▷ X} (W (A '→' W B))
- →w : ∀ {Γ A B X}
  → Term {Γ ▷ X} (W (A '→' W B))
  → Term {Γ ▷ X} (W (A '→' B))
- ww→ : ∀ {Γ A B X Y}
  → Term {Γ ▷ X ▷ Y} (W (W (A '→' B)))
  → Term {Γ ▷ X ▷ Y} (W (W A) '→' W (W B))
- →ww : ∀ {Γ A B X Y}
  → Term {Γ ▷ X ▷ Y} (W (W A) '→' W (W B))
  → Term {Γ ▷ X ▷ Y} (W (W (A '→' B)))
- 'o' : ∀ {Γ A B C}
  → Term {Γ} (B '→' C)
  → Term {Γ} (A '→' B)
  → Term {Γ} (A '→' C)
- 'w'''a' : ∀ {Γ A B T}
  → Term {Γ ▷ T} (W ('Term' '' ⊢ A '→' B ⊢T))
  → Term {Γ ▷ T} (W ('Term' '' ⊢ A ⊢T))
  → Term {Γ ▷ T} (W ('Term' '' ⊢ B ⊢T))
- 'a' : ∀ {Γ A B}
  → Term {Γ} ('Term' '' ⊢ A '→' B ⊢T
    '→' 'Term' '' ⊢ A ⊢T
    '→' 'Term' '' ⊢ B ⊢T)
- '□' : ∀ {Γ A B}
  → Term {Γ ▷ A ▷ B}
  (W (W ('Term' '' ⊢ 'Type' Γ ⊢T)))
  → Term {Γ ▷ A ▷ B}
  (W (W ('Type' Γ)))
- '→' : ∀ {Γ}
  → Term {Γ} ('Type' Γ)
  → Term {Γ} ('Type' Γ)
  → Term {Γ} ('Type' Γ)
- 'ww'''→''' : ∀ {Γ A B}
  → Term {Γ ▷ A ▷ B}
  (W (W ('Term' '' ⊢ 'Type' Γ ⊢T)))
  → Term {Γ ▷ A ▷ B}
  (W (W ('Term' '' ⊢ 'Type' Γ ⊢T)))
  → Term {Γ ▷ A ▷ B}
  (W (W ('Term' '' ⊢ 'Type' Γ ⊢T)))
- 'ww'''×''' : ∀ {Γ A B}
  → Term {Γ ▷ A ▷ B}
  (W (W ('Term' '' ⊢ 'Type' Γ ⊢T)))
  → Term {Γ ▷ A ▷ B}
  (W (W ('Term' '' ⊢ 'Type' Γ ⊢T)))

```

```

→ Term {Γ ▷ A ▷ B}
(W (W ('Term' '' Γ 'Type' Γ ⊢T)))

□ : Type ε → Set _
□ = Term {ε}

'□' : ∀ {Γ} → Type Γ → Type Γ
'□' T = 'Term' '' Γ T ⊢T

_ "×" _ : ∀ {Γ}
→ Term {Γ} ('Type' Γ)
→ Term {Γ} ('Type' Γ)
→ Term {Γ} ('Type' Γ)
A "×" B = "×" ''a A ''a B

max-level : Level
max-level = lzero

mutual
[ ]c : (Γ : Context) → Set (lsuc max-level)
[ ε ]c = ⊤
[ Γ ▷ T ]c = Σ [ Γ ]c [ T ]T

[ ]T : {Γ : Context}
→ Type Γ
→ [ Γ ]c
→ Set max-level
[ W T ]T [Γ]
= [ T ]T (Σ.fst [Γ])
[ W1 T ]T [Γ]
= [ T ]T (Σ.fst (Σ.fst [Γ]), Σ.snd [Γ])
[ T '' x ]T [Γ] = [ T ]T ([Γ], [ x ]t [Γ])
[ 'Type' Γ ]T [Γ]
= Lifted (Type Γ)
[ 'Term' ]T [Γ]
= Lifted (Term (lower (Σ.snd [Γ])))
[ A '→' B ]T [Γ] = [ A ]T [Γ] → [ B ]T [Γ]
[ A '×' B ]T [Γ] = [ A ]T [Γ] × [ B ]T [Γ]
[ '⊥' ]T [Γ] = ⊥
[ '⊥' ]T [Γ] = ⊥
[ Quine φ ]T [Γ] = [ φ ]T ([Γ], (lift (Quine φ)))

[ ]t : ∀ {Γ : Context} {T : Type Γ}
→ Term T
→ ([Γ] : [ Γ ]c)
→ [ T ]T [Γ]
[ Γ x ⊢T ]t [Γ] = lift x
[ Γ x ⊢t ]t [Γ] = lift x
[ 'Γ'VAR0'⊢t' ]t [Γ]
= lift Γ lower (Σ.snd [Γ]) ⊢t
[ 'Γ'VAR0'⊢T' ]t [Γ]
= lift Γ lower (Σ.snd [Γ]) ⊢T
[ f ''a x ]t [Γ] = [ f ]t [Γ] ([ x ]t [Γ])
[ 'tt' ]t [Γ] = tt
[ quine→ {φ} ]t [Γ] x = x
[ quine← {φ} ]t [Γ] x = x
[ 'λ' f ]t [Γ] x = [ f ]t ([Γ], x)
[ 'VAR0' ]t [Γ] = Σ.snd [Γ]
[ SW t ]t = [ t ]t
[ ←SW1SV→W f ]t = [ f ]t

```

```

[ →SW1SV→W f ]t = [ f ]t
[ ←SW1SV→SW1SV→W f ]t = [ f ]t
[ →SW1SV→SW1SV→W f ]t = [ f ]t
[ w x ]t [Γ] = [ x ]t (Σ.fst [Γ])
[ w→f ]t [Γ] = [ f ]t [Γ]
[ →w f ]t [Γ] = [ f ]t [Γ]
[ ww→f ]t [Γ] = [ f ]t [Γ]
[ →ww f ]t [Γ] = [ f ]t [Γ]
[ "×" ]t [Γ] A B = lift (lower A '×' lower B)
[ g 'o' f ]t [Γ] x = [ g ]t [Γ] ([ f ]t [Γ] x)
[ f w ''a x ]t [Γ]
= lift (lower ([ f ]t [Γ]) ''a lower ([ x ]t [Γ]))
[ "''a' ]t [Γ] f x
= lift (lower f ''a lower x)
[ "□" {Γ} T ]t [Γ]
= lift ('Term' '' lower ([ T ]t [Γ]))
[ A "→" B ]t [Γ]
= lift
(lower ([ A ]t [Γ]) '→' lower ([ B ]t [Γ]))
[ A ww "→" B ]t [Γ]
= lift
(lower ([ A ]t [Γ]) "→" lower ([ B ]t [Γ]))
[ A ww "×" B ]t [Γ]
= lift
(lower ([ A ]t [Γ]) "×" lower ([ B ]t [Γ]))

```

```

module inner ('X' : Type ε)
('f' : Term {ε} ('□' 'X' '→' 'X'))
where
'H' : Type ε
'H' = Quine (W1 'Term' '' 'VAR0' '→' W 'X')

'toH' : □ (('□' 'H' '→' 'X') '→' 'H')
'toH' = ←SW1SV→W quine←

'fromH' : □ ('H' '→' ('□' 'H' '→' 'X'))
'fromH' = →SW1SV→W quine→

'□'H'→□'X' : □ ('□' 'H' '→' '□' 'X')
'□'H'→□'X'
= 'λ' (w Γ 'fromH' ⊢t
w ''a 'VAR0'
w ''a 'Γ'VAR0'⊢t)

'h' : Term 'H'
'h' = 'toH' ''a ('f' 'o' '□'H'→□'X')

Löb : □ 'X'
Löb = 'fromH' ''a 'h' ''a Γ 'h' ⊢t

```

```

Löb : ∀ {X}
→ Term {ε} ('□' X '→' X) → Term {ε} X
Löb {X} f = inner.Löb X f

```

```

[ ] : Type ε → Set _
[ T ] = [ T ]T tt

```

```

'⊢' _ : ∀ {Γ} → Type Γ → Type Γ
'⊢' T = T '→' '⊥'

```

```


$$\begin{aligned}
& \_w \text{“} \times \text{”} \_ : \forall \{ \Gamma X \} \\
& \rightarrow \text{Term} \{ \Gamma \triangleright X \} (W ('Type' \Gamma)) \\
& \rightarrow \text{Term} \{ \Gamma \triangleright X \} (W ('Type' \Gamma)) \\
& \rightarrow \text{Term} \{ \Gamma \triangleright X \} (W ('Type' \Gamma)) \\
A \text{“} \times \text{”} B = w \rightarrow (w \rightarrow (w \text{“} \times \text{”}) \text{“} \_a A \text{”} \_a B) \\
\\
\text{l\"ob} : \forall \{ 'X' \} \rightarrow \Box (' \Box ' X ' \rightarrow ' X') \rightarrow \llbracket 'X' \rrbracket \\
\text{l\"ob} f = \llbracket \text{L\"ob} f \rrbracket^t \text{tt} \\
\\
\lrcorner \_ : \forall \{ \ell \ell' \} \rightarrow \text{Set } \ell \rightarrow \text{Set } (\ell \sqcup \ell') \\
\lrcorner \_ \{ \ell \} \{ \ell' \} T = T \rightarrow \perp \{ \ell' \} \\
\\
\text{incompleteness} : \lrcorner \Box (' \lrcorner ' (' \Box ' \perp')) \\
\text{incompleteness} = \text{l\"ob} \\
\\
\text{soundness} : \lrcorner \Box ' \perp ' \\
\text{soundness } x = \llbracket x \rrbracket^t \text{tt} \\
\\
\text{non-emptiness} : \Sigma (\text{Type } \varepsilon) (\lambda T \rightarrow \Box T) \\
\text{non-emptiness} = \text{“} \top \text{”}, \text{“} \text{tt} \text{”}
\end{aligned}$$


```

## Acknowledgments

We would like to thank Benja Fallenstein for invaluable insight into the initial formalization of Löb’s theorem in type theory, Jack Gallagher for insight and ideas in simplifying and improving the formalization of self-representation, Patrick LaVictoire for facilitating the workshop where most of the initial ideas originated, Matt Brown for helping us discover the reason that the self-interpreter for  $F_\omega$  doesn’t contradict this formalization of Löb’s theorem, and Adam Chlipala for timely and helpful suggestions on the writing of this paper.

## References

- A. W. Appel, P.-A. Mellies, C. D. Richards, and J. Vouillon. A very modal model of a modern, major, general type system. *ACM SIGPLAN Notices*, 42(1):109–122, 2007. URL <https://www.cs.princeton.edu/~appel/papers/modalmodel.pdf>.
- M. Barasz, P. Christiano, B. Fallenstein, M. Herreshoff, P. LaVictoire, and E. Yudkowsky. Robust cooperation in the prisoner’s dilemma: Program equilibrium via provability logic. *ArXiv e-prints*, Jan 2014. URL <http://arxiv.org/pdf/1401.5577v1.pdf>.
- M. Brown and J. Palsberg. Breaking through the normalization barrier: A self-interpreter for f-omega. In *Proceedings of the 43rd Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, pages 5–17. ACM, 2016. doi: 10.1145/2837614.2837623. URL <http://compilers.cs.ucla.edu/pop116/pop116-full.pdf>.
- C. McBride. Outrageous but meaningful coincidences: Dependent type-safe syntax and evaluation. In *Proceedings of the 6th ACM SIGPLAN workshop on Generic programming*, pages 1–12. ACM, 2010. URL <https://personal.cis.strath.ac.uk/conor.mcbride/pub/DepRep/DepRep.pdf>.
- G. K. Pullum. Scooping the loop snooper, October 2000. URL <http://www.lel.ed.ac.uk/~gpullum/loopsnoop.html>.
- B. Yudkowsky. Lob’s theorem cured my social anxiety, February 2014. URL <http://agentyduck.blogspot.com/2014/02/lobs-theorem-cured-my-social-anxiety.html>.