

Löb's Theorem

A functional pearl of dependently typed quining

Jason Gross

MIT CSAIL
jgross@mit.edu

Jack Gallagher

MIRI
jack@gallabytes.com

Benya Fallenstein

MIRI
benya@intelligence.org

Keywords Agda, Löb's theorem, quine, self-reference, type theory

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 implies, for total languages which validate it, that self-interpreters are impossible. 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.

1. Introduction

*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)

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, 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).

In this paper, after introducing the content of Löb's theorem, we will present in Agda three formalizations of type-theoretic languages and prove Löb's theorem in and about these languages: one that shows the theorem is admissible as an axiom in a wide range of situations, one which proves Löb's theorem by assuming as an axiom the existence of quines (programs which output their own source code), and one which constructs the proof under even weaker assumptions; see section 5 for details.

“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 are the most awesome person in the world. Since this is exactly what the sentence asserts, the sentence is true, and you 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 . Thus since assuming A yields X , we have that $A \rightarrow X$. What went wrong?¹

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.

The problem is that the phrase “this sentence is true” is not a valid mathematical assertion; no language can encode a truth predicate for itself (Tarski 1936). However, some languages *can* encode assertions about *provability* (Gödel 1931). In section 2, we will dig into the difference between truth and provability from a computational perspective. We will present an argument for the indefinability of truth and for the definability of provability, which we hope will prove enlightening when we get to the formalization of Löb's theorem itself.

Now consider the sentence “if this sentence is provable, then Santa Claus exists.” Fix a formal system powerful enough to talk about which of its sentences are provable (for example, Peano Arithmetic, Martin-Löf Type Theory, or Gödel-Löb Modal Logic), and fix a formalization of provability in that system. To prove that our sentence is true, we suppose that it is provable. We must now show that Santa Claus exists. *If provability implies truth*, then the sentence is true, and thus Santa Claus exists. Hence, 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. If we let $\Box X$ denote the assertion “ X is provable,” then, symbolically, Löb's theorem becomes:

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

¹Those 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	—
$\Box X \rightarrow X$	Interpreters	—
(In)completeness	Halting problem	—

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

2. Quines and the Curry-Howard Isomorphism

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 we 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? 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 is abstract syntax trees (ASTs). 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.

Note well that the type $(\Box X \rightarrow X)$ is a type that takes syntax trees and evaluates them; it is the type of an interpreter or an unquoter.

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

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 outputs its own source code.

There are three genuinely distinct solutions, the first of which is degenerate, and the second of which is cheeky. These 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.
- A program with a “template” which contains a copy of the source code of all of the program except for the template itself, leaving a hole where the template should be. The program then substitutes a quoted copy of the template into the hole in the template itself. In code, we can use Python’s `repr` to get a quoted copy of the template, and we do substitution using Python’s replacement syntax: for example, `("foo %s bar" % "baz")` becomes `"foo baz bar"`. Our third solution, in code, is thus:

```
T = 'T = %s\nprint(T % repr(T))'
print(T % repr(T))
```

The functional equivalent, which does not use assignment, and which we will be using later on in this paper, is:

```
(lambda T: T % repr(T))
(' (lambda T: T % repr(T))\n (%s) ')
```

We can use this technique, known as quining (Hofstadter 1979; Kleene 1952), to describe self-referential programs.

Suppose Python had a function `□` that took a quoted representation of a Martin-Löf type (as a Python string), and returned a Python object representing the Martin-Löf type of ASTs of Martin-Löf programs inhabiting that type. Now consider the program

```
φ = (lambda T: □(T % repr(T)))
(' (lambda T: □(T % repr(T))\n (%s) ')
```

The variable φ evaluates to the type of ASTs of programs inhabiting the type corresponding to $T \% \text{repr}(T)$, where T is `' (lambda T: □(T % repr(T))\n (%s) '`. What Martin-Löf type does this string, $T \% \text{repr}(T)$, represent? It represents $\Box(T \% \text{repr}(T))$, of course. Hence φ is the type of syntax trees of programs that produce proofs of φ —in other words, φ is a Henkin sentence.

Taking it one step further, assume Python has a function $\Pi(a, b)$ which takes two Python representations of Martin-Löf types and produces the Martin-Löf type $(a \rightarrow b)$ of functions from a to b . If we also assume that these functions exist in the term language of string representations of Martin-Löf types, we can consider the function

```
def Lob(X):
    T = '(lambda T: Π(□(T % repr(T)), X))(%s)'
    φ = Π(□(T % repr(T)), X)
    return φ
```

What does `Lob(X)` return? It returns the type φ of abstract syntax trees of programs producing proofs that “if φ is provable, then X .” Concretely, `Lob(⊥)` returns the type of programs which prove Martin-Löf type theory consistent, `Lob(SantaClaus)` returns the variant of the Santa Claus sentence that says “if this sentence is provable, then Santa Claus exists.”

Let us now try producing the true Santa Claus sentence, the one that says “If this sentence is true, Santa Claus exists.” We need a function `Eval` which takes a string representing a Martin-Löf program, and evaluates it to produce a term. Consider the Python program

```
def Tarski(X):
    T = '(lambda T: Π(Eval(T % repr(T)), X))(%s)'
    φ = Π(Eval(T % repr(T)), X)
    return φ
```

Running `Eval(T % repr(T))` tries to produce a term that is the type of functions from `Eval(T % repr(T))` to X . Note that φ is itself the type of functions from `Eval(T % repr(T))` to X . If `Eval(T % repr(T))` could produce a term of type φ , then φ would evaluate to the type $\varphi \rightarrow X$, giving us a bona fide Santa Claus sentence. However, `Eval(T % repr(T))` attempts to produce the type of functions from `Eval(T % repr(T))` to X by evaluating `Eval(T % repr(T))`. This throws the function `Tarski` into an infinite loop which never terminates. (Indeed, choosing $X = \perp$ it’s trivial to show that there’s no way to write `Eval` such that `Tarski` halts, unless Martin-Löf type theory is inconsistent.)

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 (McBride 2010; Chapman 2009; Danielsson 2007). For example, here is a very simple (and incomplete) formalization with dependent function types (Π), a unit type (\top), an empty type (\perp), and functions (λ).

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

```
mutual
  infixl 2 _▷_

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

  data Type : Context → Set where
    'T' : ∀ {Γ} → Type Γ
    '⊥' : ∀ {Γ} → Type Γ
    'Π' : ∀ {Γ}
      → (A : Type Γ) → Type (Γ ▷ A) → Type Γ

  data Term : {Γ : Context} → Type Γ → Set where
    'tt' : ∀ {Γ} → Term {Γ} 'T'
    'λ' : ∀ {Γ A B} → Term B → Term {Γ} ('Π' 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. This function gives a semantic model to the syntax. Here is an example of such a function:

```
mutual
  [ ]c : Context → Set
  [ ε ]c = ⊤
  [ Γ ▷ T ]c = Σ [ Γ ]c [ T ]T

  [ ]T : ∀ {Γ}
    → Type Γ → [ Γ ]c → Set
  [ 'T' ]T Γ↓ = ⊤
  [ '⊥' ]T Γ↓ = ⊥
  [ 'Π' A B ]T Γ↓ = (x : [ A ]T Γ↓) → [ B ]T (Γ↓, x)

  [ ]t : ∀ {Γ T}
    → Term {Γ} T → (Γ↓ : [ Γ ]c) → [ T ]T Γ↓
  [ 'tt' ]t Γ↓ = tt
  [ 'λ' f ]t Γ↓ x = [ f ]t (Γ↓, x)
```

Note that this interpretation function has an essential property that we will call *locality*: the interpretation of any given constructor does not require doing case analysis on any of its arguments. By contrast, one could imagine an interpretation function that interpreted function types differently depending on their domain and codomain; for example, one might interpret $(\perp \rightarrow A)$ as \top , or one might interpret an equality type differently at each type, as in Observational Type Theory (Altenkirch et al. 2007).

4. This Paper

In this paper, we make extensive use of this trick for validating models. In section 6, 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 in section 7 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 in section 8 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()`. After taking a digression for an application of Löb’s theorem to the prisoner’s dilemma in section 9, we sketch in section 10 our implementation of Löb’s theorem (code in the supplemental material) based on only the assumption 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 in section 11 with some discussion about avenues for removing the hard-coded `repr`.

5. Prior Work

There exist a number of implementations or formalizations of various flavors of Löb’s theorem in the literature. Appel et al. use Löb’s theorem as an induction rule for program logics in Coq (Appel et al. 2007). Piponi formalizes a rule with the same shape as Löb’s theorem in Haskell, and uses it for, among other things, spreadsheet evaluation (Piponi 2006). Simmons and Toninho formalize a constructive provability logic in Agda, and prove Löb’s theorem within that logic (Simmons and Toninho 2012).

Gödel’s incompleteness theorems, easy corollaries to Löb’s theorem, have been formally verified numerous times (Shankar 1986, 1997; O’Connor 2005; Paulson 2015).

To our knowledge, our twelve line proof is the shortest self-contained formally verified proof of the admissibility of Löb’s theorem to date. We are not aware of other formally verified proofs of Löb’s theorem which interpret the modal \Box operator as an inductively defined type of syntax trees of proofs of a given theorem, as we do in this formalization, though presumably the modal \Box operator Simmons and Toninho could be interpreted as such syntax trees. Finally, we are not aware of other work which uses the trick of talking about a local interpretation function (as described at the end of section 3) to talk about consistent extensions to classes of encodings of type theory.

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 $\Box' T$ to denote the syntactic type corresponding to the type of (syntactic) terms whose syntactic type is T . For example, the type of a `repr` which operated on syntax trees would be $\Box T \rightarrow \Box (\Box' T)$.

```
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:

```

löb : ∀ {X'} → □ ('□' X' '→' X') → [ X' ]T
löbf = [ Löbf ]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_ω (Brown and Palsberg 2016).²

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

'¬' : Type → Type
'¬' T = T '→' '⊥'

löb : ∀ {X'} → □ ('□' X' '→' X') → [ X' ]T
löbf = [ Löbf ]t

-- There is no syntactic proof of absurdity

```

²One may wonder how exactly the self-interpreter for F_ω does not contradict this theorem. In private conversations with Matt Brown, we found that the F_ω 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 $\square (\square X \rightarrow X) \rightarrow \square X$, which is not strictly positive, or $\square (X \rightarrow X) \rightarrow \square X$, which, on interpretation, must be filled with a general fixpoint operator. Such an operator is well-known to be inconsistent.

```

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

```

```

-- but it would be absurd to have a syntactic
-- proof of that fact
incompleteness : ¬ □ ('¬' ('□' '⊥'))
incompleteness = löb

-- However, there are syntactic proofs of some
-- things (namely ⊤)
non-emptiness : □ '⊤'
non-emptiness = 'tt'

-- There are no syntactic interpreters, things
-- which, at any type, evaluate code at that
-- type to produce its result.
no-interpreters : ¬ (∀ {X'} → □ ('□' X' '→' X'))
no-interpreters interp = löb (interp { '⊥' })

```

What is this incompleteness theorem? Gödel's incompleteness theorem is typically interpreted as "there exist true but unprovable statements." In intuitionistic logic, this is hardly surprising. A more accurate rendition of the theorem in Agda might be "there exist true but inadmissible statements," i.e., there are statements which are provable meta-theoretically, but which lead to (meta-theoretically-provable) inconsistency if assumed at the object level.

This may seem a bit esoteric, so let's back up a bit, and make it more concrete. Let's begin by banishing "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 $(\neg \square T) \wedge \neg \square (\neg T)$, 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 us, 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, then, is this: your reflective system, no matter how powerful, cannot formalize every proof. 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 stated above says. Incompleteness, fundamentally, is a result about the limitations of formalizing provability.

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 ϕ* is program that outputs the result of evaluating the function ϕ on the abstract syntax tree of its own type. Letting $\text{Quine } \phi$ denote the constructor for a type-level quine at ϕ , we have an isomorphism between $\text{Quine } \phi$ and $\phi \ulcorner \text{Quine } \phi \urcorner^T$, where $\ulcorner \text{Quine } \phi \urcorner^T$ is the abstract syntax tree for the source code of $\text{Quine } \phi$. 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`.

We begin with an encoding of contexts and types, repeating from above the constructors of ' \rightarrow ', ' \square ', ' \top ', and ' \perp '. We add to

this a constructor for quines (Quine), and a constructor for syntax trees of types in the empty context ('Types'). 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. Note that ' \square ' is now a function of the represented language, rather than a meta-level operator.

Note that we use the infix operator $_ \text{' ' } _$ to denote substitution.

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

data Type : Context → Set where
  '→' _ : ∀ {Γ} → Type Γ → Type Γ → Type Γ
  'T' : ∀ {Γ} → Type Γ
  '⊥' : ∀ {Γ} → Type Γ
  'Typeε' : ∀ {Γ} → Type Γ
  '□' : ∀ {Γ} → Type (Γ ▷ 'Typeε')
  Quine : Type (ε ▷ 'Typeε') → Type ε
  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
```

In addition to ' λ ' and 'tt', we now have the AST-equivalents of Python's `repr`, which we denote as $\ulcorner _ \urcorner^T$ for the type-level add-quote function, and $\ulcorner _ \urcorner^t$ 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 found it unnecessary.

To construct the proof of Löb's theorem, we need a few other standard constructors, such as ' VAR_0 ', which references a term in the context; $_ \text{' ' } _$, which we use to denote function application; $_ \text{' o ' } _$, a function composition operator; and ' $\ulcorner \text{VAR}_0 \urcorner^t$ ', the variant of ' VAR_0 ' 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 avoid defining a judgmental equality or conversion relation.

```
data Term : {Γ : Context} → Type Γ → Set where
  'λ' : ∀ {Γ A B}
    → Term {Γ ▷ A} (W B) → Term (A '→' B)
  'tt' : ∀ {Γ}
    → Term {Γ} 'T'
  '⌊' _⌋T : ∀ {Γ} -- type-level repr
    → Type ε
    → Term {Γ} 'Typeε'
  '⌊' _⌋t : ∀ {Γ T} -- term-level repr
    → Term {ε} T
    → Term {Γ} ('□' '' Γ TT)
  quine→ : ∀ {φ}
    → Term {ε} (Quine φ '→' φ '' Γ Quine φT)
  quine← : ∀ {φ}
    → Term {ε} (φ '' Γ Quine φT '→' Quine φ)
  -- The constructors below here are for
  -- variables, weakening, and substitution
  'VAR0' : ∀ {Γ T}
    → Term {Γ ▷ T} (W T)
  _''_ : ∀ {Γ A B}
    → Term {Γ} (A '→' B)
```

```
→ Term {Γ} A
→ Term {Γ} B
_ 'o' _ : ∀ {Γ A B C}
  → Term {Γ} (B '→' C)
  → Term {Γ} (A '→' B)
  → Term {Γ} (A '→' C)
'⌊' VAR0 ⌋t : ∀ {T}
  → Term {ε ▷ '□' '' Γ TT}
  (W ('□' '' Γ '□' '' Γ TT TT))
→ SW1 SV → W : ∀ {Γ T X A B} {x : Term {Γ} X}
  → Term (T '→' (W1 A '' 'VAR0' '→' W B) '' x)
  → Term (T '→' A '' x '→' B)
← SW1 SV → W : ∀ {Γ T X A B} {x : Term {Γ} X}
  → Term ((W1 A '' 'VAR0' '→' W B) '' x '→' T)
  → Term ((A '' x '→' B) '→' 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 B T}
  → Term {ε ▷ T} (W ('□' '' Γ A '→' BT))
  → Term {ε ▷ T} (W ('□' '' Γ AT))
  → Term {ε ▷ T} (W ('□' '' Γ BT))
```

```
□ : 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 ϕ 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 (|suc max-level)
[ ε ]c = T
[ Γ ▷ T ]c = Σ [ Γ ]c [ T ]T

[ ]T : ∀ {Γ}
  → Type Γ → [ Γ ]c → Set max-level
[ 'Typeε' ]T Γ⌋ = Lifted (Type ε)
[ '□' ]T Γ⌋ = Lifted (Term {ε} (lower (snd Γ⌋)))
[ Quine φ ]T Γ⌋ = [ φ ]T (Γ⌋, lift (Quine φ))
-- The rest of the type-level interpretations
-- are the obvious ones, if a bit obscured by
-- carrying around the context.
[ A '→' B ]T Γ⌋ = [ A ]T Γ⌋ → [ B ]T Γ⌋
[ 'T' ]T Γ⌋ = T
[ '⊥' ]T Γ⌋ = ⊥
[ W T ]T Γ⌋ = [ T ]T (fst Γ⌋)
[ W1 T ]T Γ⌋ = [ T ]T (fst (fst Γ⌋), snd Γ⌋)
[ T '' x ]T Γ⌋ = [ T ]T (Γ⌋, [ x ]t Γ⌋)

[ ]t : ∀ {Γ T}
  → Term {Γ} T → (Γ⌋ : [ Γ ]c) → [ T ]T Γ⌋
[ '⌊' x ⌋T ]t Γ⌋ = lift x
[ '⌊' x ⌋t ]t Γ⌋ = lift x
[ quine→ ]t Γ⌋ x = x
[ quine← ]t Γ⌋ x = x
```

-- The rest of the term-level interpretations
 -- are the obvious ones, if a bit obscured by
 -- carrying around the context.

```

[λ' f] Γ ⊢ x = [f] (Γ ⊢, x)
[tt] Γ ⊢ = tt
[VAR₀] Γ ⊢ = snd Γ ⊢
[Γ' VAR₀ Γ ⊢] Γ ⊢ = lift Γ lower (snd Γ ⊢) Γ ⊢
[g 'o' f] Γ ⊢ x = [g] Γ ⊢ ([f] Γ ⊢ x)
[f 'a' x] Γ ⊢ = [f] Γ ⊢ ([x] Γ ⊢)
[←SW₁SV→W f] Γ ⊢ = [f] Γ ⊢
[→SW₁SV→W f] Γ ⊢ = [f] Γ ⊢
[w x] Γ ⊢ = [x] Γ ⊢ (fst Γ ⊢)
[w → f] Γ ⊢ = [f] Γ ⊢ (fst Γ ⊢)
[f w 'a' x] Γ ⊢
  = lift (lower ([f] Γ ⊢) 'a' lower ([x] 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 $\lambda v. \Box v \rightarrow 'X'$. We can then convert back and forth between the types $\Box 'H'$ and $\Box 'H' \rightarrow 'X'$ with our quine isomorphism functions, and a bit of quotation magic and function application gives us a term of type $\Box 'H' \rightarrow \Box 'X'$; 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 $\Box 'X' \rightarrow 'X'$, to get a term of type $\Box 'H' \rightarrow 'X'$, 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 $\Box 'X'$, i.e., a proof that Santa Claus exists.

```

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

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

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

    '□' 'H' → □ 'X' : □ ('□' 'Γ' 'H' Γ ⊢ → '□' 'Γ' 'X' Γ ⊢)
    '□' 'H' → □ 'X'
      = 'λ' (w Γ 'fromH' Γ ⊢
        w 'a' 'VAR₀'
        w 'a' 'Γ' 'VAR₀' Γ ⊢)

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

    Löb : □ 'X'
    Löb = 'fromH' 'a' 'h' 'a' 'h' Γ ⊢

  Löb : ∀ {X} → □ ('□' 'Γ' X Γ ⊢ → X) → □ X
  Löb {X} f = inner.Löb X f

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

  '¬' _ : ∀ {Γ} → Type Γ → Type Γ

```

A Says \ B Says	Cooperate	Defect
	(1 year, 1 year) (3 years, 0 years)	(0 years, 3 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).

$'\neg' T = T '→' '\perp'$

$\text{Löb} : \forall \{X\} \rightarrow \Box ('□' 'Γ' 'X' Γ ⊢ → 'X') \rightarrow \Box X$
 $\text{Löb} f = [] \Gamma \text{tt}$

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

$\text{incompleteness} : \neg \Box ('□' 'Γ' '\perp' \Gamma \text{tt})$
 $\text{incompleteness} = \text{Löb}$

$\text{soundness} : \neg \Box '\perp'$
 $\text{soundness} x = [x] \Gamma \text{tt}$

$\text{non-emptiness} : \Sigma (\text{Type } \varepsilon) (\lambda T \rightarrow \Box T)$
 $\text{non-emptiness} = 'T', 'tt'$

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'
--  $\leftrightarrow$  (Term 'Bot'  $\rightarrow$  Term 'Bot'  $\rightarrow$  'Type')
'Bot' :  $\forall \{ \Gamma \} \rightarrow$  Type  $\Gamma$ 
'Bot' {  $\Gamma$  }
  = Quine (W1 'Term' " 'VAR0'
    '→' W1 'Term' " 'VAR0'
    '→' W ('Type'  $\Gamma$ ))
```

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 { $\varepsilon$ }", i.e., a term in
-- the empty context
```

```
_cooperates-with_ : □ 'Bot'  $\rightarrow$  □ 'Bot'  $\rightarrow$  Type  $\varepsilon$ 
b1 cooperates-with b2 = lower (ll b1)t tt (lift b1) (lift b2))
```

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

```
make-bot :  $\forall \{ \Gamma \}$ 
   $\rightarrow$  Term {  $\Gamma \triangleright$  □ 'Bot'  $\triangleright$  W (□ 'Bot') }
    (W (W ('Type'  $\Gamma$ )))
   $\rightarrow$  Term {  $\Gamma$  } 'Bot'
make-bot t
  =  $\leftarrow$  SW1SV  $\rightarrow$  SW1SV  $\rightarrow$  W
  quine  $\leftarrow$  "a" 'λ' ( $\rightarrow$  w ('λ' t))
```

```
'DefectBot' : □ 'Bot'
'CooperateBot' : □ 'Bot'
'FairBot' : □ '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  $\Gamma$  '⊥'  $\neg$  T))
'CooperateBot' = make-bot (w (w  $\Gamma$  '⊤'  $\neg$  T))
```

```
DB-defects :  $\forall \{ b \}$ 
   $\rightarrow \neg$  ll 'DefectBot' cooperates-with b ll
DB-defects {b} pf = pf
```

```
CB-cooperates :  $\forall \{ b \}$ 
   $\rightarrow$  ll 'CooperateBot' cooperates-with b ll
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' :  $\forall \{ \Gamma \}$ 
   $\rightarrow$  Term {  $\Gamma$  } ('Bot'
    '→' (□ 'Bot' '→' '□' 'Bot' '→' 'Type'  $\Gamma$ ))
'eval-bot' =  $\rightarrow$  SW1SV  $\rightarrow$  SW1SV  $\rightarrow$  W quine  $\rightarrow$ 
```

```
-- 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' :  $\forall \{ \Gamma \}$ 
   $\rightarrow$  Term {  $\Gamma$  } ('□' 'Bot'
    '→' '□' ({- other -} '□' 'Bot' '→' 'Type'  $\Gamma$ ))
'eval-bot'
  = 'λ' (w  $\Gamma$  'eval-bot'  $\neg$  t
    w'''a 'VAR0'
    w'''a 'Γ' 'VAR0'  $\neg$  t)
```

```
-- 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' :  $\forall \{ \Gamma \}$ 
   $\rightarrow$  Term {  $\Gamma$ 
     $\triangleright$  □ 'Bot'
     $\triangleright$  W (□ 'Bot') }
    (W (W (□ 'Bot')) '→' W (W (□ ('Type'  $\Gamma$ ))))
'other-cooperates-with' {  $\Gamma$  }
  = 'eval-other'
  'o' w  $\rightarrow$  (w (w  $\rightarrow$  (w ('λ' 'Γ' 'VAR0'  $\neg$  t))))
```

```
where
```

```
'eval-other'
  : Term {  $\Gamma \triangleright$  □ 'Bot'  $\triangleright$  W (□ 'Bot') }
    (W (W (□ ('□' 'Bot' '→' 'Type'  $\Gamma$ ))))
'eval-other'
  = w  $\rightarrow$  (w (w  $\rightarrow$  (w "eval-bot")))
```

```
'eval-other'
  : Term (W (W (□ ('□' 'Bot'))
    '→' W (W (□ ('Type'  $\Gamma$ ))))
'eval-other'
  = ww  $\rightarrow$  (w  $\rightarrow$  (w (w  $\rightarrow$  (w "a")))
```

```
-- A bot gets its own source code as the first
-- argument (of two)
```

```
'self' :  $\forall \{ \Gamma \}$ 
   $\rightarrow$  Term {  $\Gamma \triangleright$  □ 'Bot'  $\triangleright$  W (□ 'Bot') }
    (W (W (□ 'Bot'))
'eval-bot' = w 'VAR0'
```

```
-- A bot gets its opponent's source code as
-- the second argument (of two)
```

```
'other' :  $\forall \{ \Gamma \}$ 
   $\rightarrow$  Term {  $\Gamma \triangleright$  □ 'Bot'  $\triangleright$  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 leave the proof that this formalization of FairBot cooperates with itself as an exercise for the reader. In Appendix C, we present an alternative formalization with a simple proof that FairBot cooperates with itself, but with no general definition of the type of bots; we relegate this code to an appendix so as to not confuse the reader by introducing a different way of handling contexts and weakening in the middle of this paper.

10. Encoding with an Add-Quote Function

Now we return to our proving of Löb’s theorem. Included in the artifact for this paper³ is code that replaces the Quine constructor with simpler constructors. Because the lack of β -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 Table 2:

```
(lambda T:  $\Pi(\Box(T \% \text{repr}(T)), X)$ )
(' (lambda T:  $\Pi(\Box(T \% \text{repr}(T)), X)$ )\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. What shape does that hole have? The shape is that of a syntax tree 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 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 Σ type, using $\ulcorner A \urcorner$ to denote the abstract syntax tree for A ,⁴ and using $\%s$ to denote the first variable in the context (written as ‘ VAR_0 ’ in previous formalizations above), we can write:

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

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

LöbSentence : Type  $\varepsilon$ 
LöbSentence
  = ( $\lambda$  ( $T : \Sigma \text{ Context Type}$ )
     $\rightarrow \Box$  ( $\text{cast-fst } T \% \text{repr } T$ ) ‘ $\rightarrow$ ’  $X$ )
    ( $\varepsilon \triangleright \ulcorner \Sigma \text{ Context Type } \urcorner$ 
    ,  $\ulcorner \lambda$  ( $T : \Sigma \text{ Context Type}$ )
       $\rightarrow \Box$  ( $\text{cast-fst } T \% \text{repr } T$ ) ‘ $\rightarrow$ ’  $X$ )
      (%s)  $\urcorner$ 
```

³In lob-build-quine.lagda.

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

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 (in particular, at $\Sigma \text{ Context Type}$), with appropriate reduction on equal things
- Σ 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 and without assuming a constructor for the type of syntax trees (‘ Context ’, ‘ Type ’, and ‘ Term ’ or ‘ \Box ’ in our formalizations). We would instead support inductive types, and construct these operators as inductive types and as folds over inductive types.

If you take away only three things from this paper, take away these:

1. There will always be some true things which are not possible to say, no matter how good you are at talking in type theory about type theory.
2. Giving meaning to syntax in a way that doesn’t use cases inside cases allows you to talk about when it’s okay to add new syntax.
3. If believing in something is enough to make it true, then it already is. Dream big.

A. Standard Data-Type Declarations

```
open import Agda.Primitive public
using (Level;  $\_ \sqsubseteq \_$ ; lzero; lsuc)
```

```
infixl 1  $\_ \rightarrow \_$ 
infixr 2  $\_ \times \_$ 
infixl 1  $\_ \equiv \_$ 
```

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

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

```
 $\neg \_ : \forall \{ \ell \ell' \} \rightarrow \text{Set } \ell \rightarrow \text{Set } (\ell \sqcup \ell')$ 
 $\neg \_ \{ \ell \} \{ \ell' \} T = T \rightarrow \perp \{ \ell' \}$ 
```

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


$$\begin{aligned}
& \text{'tt'} : \forall \{\Gamma\} \rightarrow \text{Term } \{\Gamma\} \text{ 'T'} \\
& \text{'}\lambda\text{' } : \forall \{\Gamma A B\} \\
& \quad \rightarrow \text{Term } \{\Gamma \triangleright A\} (W B) \\
& \quad \rightarrow \text{Term } (A \text{ '}\rightarrow\text{' } B) \\
& \text{'VAR}_0\text{' } : \forall \{\Gamma T\} \rightarrow \text{Term } \{\Gamma \triangleright T\} (W T) \\
& \text{'}\neg^T\text{' } : \forall \{\Gamma\} \\
& \quad \rightarrow \text{Type } \Gamma \\
& \quad \rightarrow \text{Term } \{\Gamma\} (\text{'Type' } \Gamma) \\
& \text{'}\neg^t\text{' } : \forall \{\Gamma T\} \\
& \quad \rightarrow \text{Term } \{\Gamma\} T \\
& \quad \rightarrow \text{Term } \{\Gamma\} (\text{'Term' } '' \neg^T T \neg^T) \\
& \text{'}\neg\text{'VAR}_0\text{'}\neg^t\text{' } : \forall \{\Gamma T\} \\
& \quad \rightarrow \text{Term } \{\Gamma \triangleright \text{'Term' } '' \neg^T T \neg^T\} \\
& \quad \quad (W (\text{'Term' } '' \neg^T \text{'Term' } '' \neg^T T \neg^T \neg^T)) \\
& \text{'}\neg\text{'VAR}_0\text{'}\neg^T\text{' } : \forall \{\Gamma\} \\
& \quad \rightarrow \text{Term } \{\Gamma \triangleright \text{'Type' } \Gamma\} \\
& \quad \quad (W (\text{'Term' } '' \neg^T \text{'Type' } \Gamma \neg^T)) \\
& \text{'}\neg\text{'a_}\text{' } : \forall \{\Gamma A B\} \\
& \quad \rightarrow \text{Term } \{\Gamma\} (A \text{ '}\rightarrow\text{' } B) \\
& \quad \rightarrow \text{Term } \{\Gamma\} A \\
& \quad \rightarrow \text{Term } \{\Gamma\} B \\
& \text{'}\times'''\text{' } : \forall \{\Gamma\} \\
& \quad \rightarrow \text{Term } \{\Gamma\} (\text{'Type' } \Gamma \\
& \quad \quad \text{'}\rightarrow\text{' } \text{'Type' } \Gamma \\
& \quad \quad \text{'}\rightarrow\text{' } \text{'Type' } \Gamma) \\
& \text{quine}\rightarrow\text{' } : \forall \{\Gamma \varphi\} \\
& \quad \rightarrow \text{Term } \{\Gamma\} \\
& \quad \quad (\text{Quine } \varphi \text{ '}\rightarrow\text{' } \varphi '' \neg \text{Quine } \varphi \neg^T) \\
& \text{quine}\leftarrow\text{' } : \forall \{\Gamma \varphi\} \\
& \quad \rightarrow \text{Term } \{\Gamma\} \\
& \quad \quad (\varphi '' \neg \text{Quine } \varphi \neg^T \text{'}\rightarrow\text{' } \text{Quine } \varphi) \\
& \text{SW} : \forall \{\Gamma X A\} \{a : \text{Term } A\} \\
& \quad \rightarrow \text{Term } \{\Gamma\} (W X '' a) \\
& \quad \rightarrow \text{Term } X \\
& \rightarrow \text{SW}_1 \text{SV} \rightarrow W \\
& : \forall \{\Gamma T X A B\} \{x : \text{Term } X\} \\
& \quad \rightarrow \text{Term } \{\Gamma\} \\
& \quad \quad (T \text{'}\rightarrow\text{' } (W_1 A '' \text{'VAR}_0\text{' } \text{'}\rightarrow\text{' } W B) '' x) \\
& \quad \rightarrow \text{Term } \{\Gamma\} \\
& \quad \quad (T \text{'}\rightarrow\text{' } A '' x \text{'}\rightarrow\text{' } B) \\
& \leftarrow \text{SW}_1 \text{SV} \rightarrow W \\
& : \forall \{\Gamma T X A B\} \{x : \text{Term } X\} \\
& \quad \rightarrow \text{Term } \{\Gamma\} \\
& \quad \quad ((W_1 A '' \text{'VAR}_0\text{' } \text{'}\rightarrow\text{' } W B) '' x \text{'}\rightarrow\text{' } T) \\
& \quad \rightarrow \text{Term } \{\Gamma\} \\
& \quad \quad ((A '' x \text{'}\rightarrow\text{' } B) \text{'}\rightarrow\text{' } T) \\
& \rightarrow \text{SW}_1 \text{SV} \rightarrow \text{SW}_1 \text{SV} \rightarrow W \\
& : \forall \{\Gamma T X A B\} \{x : \text{Term } X\} \\
& \quad \rightarrow \text{Term } \{\Gamma\} (T \text{'}\rightarrow\text{' } (W_1 A '' \text{'VAR}_0\text{' } \text{'}\rightarrow\text{' } W_1 A '' \text{'VAR}_0\text{' } \text{'}\rightarrow\text{' } W_1 A '' \text{'VAR}_0\text{'}))
\end{aligned}$$

```

    '→' 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)
- "→" " " _ : ∀ {Γ A B}
  → Term {Γ ▷ A ▷ B}
    (W (W ('Term' " " 'Type' Γ "T)))
  → Term {Γ ▷ A ▷ B}
    (W (W ('Type' Γ "T)))
  → Term {Γ ▷ A ▷ B}
    (W (W ('Term' " " 'Type' Γ "T)))
- "×" " " _ : ∀ {Γ A B}
  → Term {Γ ▷ A ▷ B}
    (W (W ('Term' " " 'Type' Γ "T)))
  → Term {Γ ▷ A ▷ B}
    (W (W ('Type' Γ "T)))
  → Term {Γ ▷ A ▷ B}
    (W (W ('Term' " " 'Type' Γ "T)))
  → Term {Γ ▷ A ▷ B}
    (W (W ('Type' Γ "T)))
  → Term {Γ ▷ A ▷ B}
    (W (W ('Term' " " 'Type' Γ "T)))

```

□ : Type ε → Set _

□ = Term {ε}

'□' : ∀ {Γ} → Type Γ → Type Γ

'□' T = 'Term' " " 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' ]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])

```

data $\in (T : \text{Type}) : \text{Context} \rightarrow \text{Set}$ where

First we want our “variable zero”, which lets us pick off the “top” element of the context.

$\text{top} : \forall \{\Gamma\} \rightarrow T \in (T :: \Gamma)$

Then we want a way to extend variables to work in larger contexts.

$\text{pop} : \forall \{\Gamma S\} \rightarrow T \in \Gamma \rightarrow T \in (S :: \Gamma)$

And, finally, we are ready to define the term language for our extended STLC.

$\text{data } _ \vdash _ (\Gamma : \text{Context}) : \text{Type} \rightarrow \text{Set where}$

The next few constructors are fairly standard. Before anything else, we want to be able to lift bindings into terms.

$\text{var} : \forall \{T\} \rightarrow T \in \Gamma \rightarrow \Gamma \vdash T$

Then the intro rules for all of our easier datatypes.

$\langle _ \rangle : \Gamma \vdash 'T'$

$_ \vdash _ : \forall \{A B\} \rightarrow \Gamma \vdash A \rightarrow \Gamma \vdash B \rightarrow \Gamma \vdash A ' \times ' B$

$_ \vdash _ \text{-elim} : \forall \{A\} \rightarrow \Gamma \vdash ' _ ' \rightarrow \Gamma \vdash A$

$\pi_1 : \forall \{A B\} \rightarrow \Gamma \vdash A ' \times ' B \rightarrow \Gamma \vdash A$

$\pi_2 : \forall \{A B\} \rightarrow \Gamma \vdash A ' \times ' B \rightarrow \Gamma \vdash B$

$' \lambda ' : \forall \{A B\} \rightarrow (A :: \Gamma) \vdash B \rightarrow \Gamma \vdash (A ' \rightarrow ' B)$

$_ \vdash _ \text{'a'} : \forall \{A B\} \rightarrow \Gamma \vdash (A ' \rightarrow ' B) \rightarrow \Gamma \vdash A \rightarrow \Gamma \vdash B$

At this point things become more delicate. To properly capture Gödel–Löb modal logic, abbreviated as GL, we want our theory to validate the rules

1. $\vdash A \rightarrow \vdash \Box A$
2. $\vdash \Box A \rightarrow \vdash \Box \Box A$

However, it should *not* validate $\vdash A \rightarrow \vdash \Box A$. If we only had the unary \Box operator we would run into difficulty later. Crucially, we couldn’t add the rule $\Gamma \vdash A \rightarrow \Gamma \vdash \Box A$, since this would let us prove $A \rightarrow \Box A$.

We will use Gödel quotes to denote the constructor corresponding to rule 1:

$\ulcorner _ \urcorner : \forall \{\Delta A\} \rightarrow \Delta \vdash A \rightarrow \Gamma \vdash (\Delta ' \vdash ' A)$

Similarly, we will write the rule validating $\Box A \rightarrow \Box \Box A$ as repr .

$\text{repr} : \forall \{\Delta A\} \rightarrow \Gamma \vdash (\Delta ' \vdash ' A) \rightarrow \Gamma \vdash (\Delta ' \vdash ' (\Delta ' \vdash ' A))$

We would like to be able to apply functions under \Box , and for this we introduce the so-called “distribution” rule. In GL, it takes the form $\vdash \Box (A \rightarrow B) \rightarrow \vdash \Box A \rightarrow \vdash \Box B$. For us it is not much more complicated.

$\text{dist} : \forall \{\Delta A B\}$
 $\rightarrow \Gamma \vdash (\Delta ' \vdash ' (A \rightarrow B))$
 $\rightarrow \Gamma \vdash (\Delta ' \vdash ' A)$
 $\rightarrow \Gamma \vdash (\Delta ' \vdash ' B)$

And, finally, we include the Löbian axiom.

$\text{Löb} : \forall \{\Delta A\}$
 $\rightarrow \Gamma \vdash (\Delta ' \vdash ' ((\Delta ' \vdash ' A) \rightarrow A))$
 $\rightarrow \Gamma \vdash (\Delta ' \vdash ' A)$

From these constructors we can prove the simpler form of the Löb rule.

$\text{Löb} : \forall \{\Gamma A\} \rightarrow \Gamma \vdash ((\Gamma ' \vdash ' A) \rightarrow A) \rightarrow \Gamma \vdash A$
 $\text{Löb } t = t \text{'a'} \text{Löb } \ulcorner t \urcorner$

Of course, because we are using DeBruijn indices, before we can do too much we’ll need to give an account of lifting. Thankfully, unlike when we were dealing with dependent type theory, we can define these computationally, and get for free all the congruences we had to add as axioms before.

Our definition of weakening is unremarkable, and sufficiently simple that Agsy, Agda’s automatic proof-finder, was able to fill in all of the code; we include it in the artifact and elide all but the type signature from the paper.

$\text{lift-tm} : \forall \{\Gamma A\} T \Delta \rightarrow (\Delta ++ \Gamma) \vdash A \rightarrow (\Delta ++ (T :: \Gamma)) \vdash A$
 Weakening is a special case of lift-tm .

$\text{wk} : \forall \{\Gamma A B\} \rightarrow \Gamma \vdash A \rightarrow (B :: \Gamma) \vdash A$

$\text{wk} = \text{lift-tm } _ \varepsilon$

Finally, we define function composition for our internal language.

$\text{infixl } 10 _ \circ ' _$
 $_ \circ ' _ : \forall \{\Gamma A B C\}$
 $\rightarrow \Gamma \vdash (B ' \rightarrow ' C)$
 $\rightarrow \Gamma \vdash (A ' \rightarrow ' B)$
 $\rightarrow \Gamma \vdash (A ' \rightarrow ' C)$
 $f \circ ' g = ' \lambda ' (\text{wk } f \text{'a'} (\text{wk } g \text{'a'} \text{var top}))$

Now we are ready to prove that FairBot cooperates with itself. Sadly, our type system isn’t expressive enough to give a general type of bots, but we can still prove things about the interactions of particular bots if we substitute their types by hand. For example, we can state the desired theorem (that FairBot cooperates with itself) as:

$\text{distf} : \forall \{\Gamma \Delta A B\}$
 $\rightarrow \Gamma \vdash (\Delta ' \vdash ' A \rightarrow B)$
 $\rightarrow \Gamma \vdash (\Delta ' \vdash ' A) \rightarrow (\Delta ' \vdash ' B)$
 $\text{distf } bf = ' \lambda ' (\text{dist } (\text{wk } bf) (\text{var top}))$

$\text{evf} : \forall \{\Gamma \Delta A\}$
 $\rightarrow \Gamma \vdash (\Delta ' \vdash ' A) \rightarrow (\Delta ' \vdash ' (\Delta ' \vdash ' A))$
 $\text{evf} = ' \lambda ' (\text{repr } (\text{var top}))$

$\text{fb-fb-cooperate} : \forall \{\Gamma A B\}$
 $\rightarrow \Gamma \vdash (\Gamma ' \vdash ' A) \rightarrow B$
 $\rightarrow \Gamma \vdash (\Gamma ' \vdash ' B) \rightarrow A$
 $\rightarrow \Gamma \vdash (A \times B)$

$\text{fb-fb-cooperate } a b$
 $= \text{Löb } (b \circ ' \text{distf } \ulcorner a \urcorner \circ ' \text{evf})$
 $, \text{Löb } (a \circ ' \text{distf } \ulcorner b \urcorner \circ ' \text{evf})$

We can also state the theorem in a more familiar form with a couple abbreviations

$\ulcorner _ \urcorner = _ \text{'a'} \text{'a'} \text{'a'}$
 $\Box = _ \vdash _ \varepsilon$

$\text{fb-fb-cooperate}' : \forall \{A B\}$
 $\rightarrow \Box (\ulcorner A \urcorner \rightarrow B)$
 $\rightarrow \Box (\ulcorner B \urcorner \rightarrow A)$
 $\rightarrow \Box (A \times B)$

$\text{fb-fb-cooperate}' = \text{fb-fb-cooperate}$

In the file `fair-bot-self-cooperates.lagda` in the artifact, we show all the meta-theoretic properties we had before: soundness, inhabitedness, and incompleteness.

Acknowledgments

We would like to thank 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_ω doesn’t contradict this formalization of Löb’s theorem, and Nate Soares and Adam Chlipala for their timely and invaluable suggestions on the writing of this paper.

References

- T. Altenkirch, C. McBride, and W. Swierstra. Observational equality, now! In *Proceedings of the 2007 Workshop on Programming Languages Meets Program Verification*, PLPV ’07, pages 57–68, New York, NY, USA, 2007. ACM. ISBN 978-1-59593-677-6. doi: 10.1145/1292597.1292608. URL <http://www.cs.nott.ac.uk/~pszttxa/publ/obseqnow.pdf>.
- 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>.
- J. Chapman. Type theory should eat itself. *Electronic Notes in Theoretical Computer Science*, 228:21–36, 2009. ISSN 1571-0661. doi: <http://dx.doi.org/10.1016/j.entcs.2008.12.114>. URL <http://www.sciencedirect.com/science/article/pii/S157106610800577X>. Proceedings of the International Workshop on Logical Frameworks and Metalanguages: Theory and Practice (LFMTP 2008).
- N. A. Danielsson. A formalisation of a dependently typed language as an inductive-recursive family. In T. Altenkirch and C. McBride, editors, *Types for Proofs and Programs*, volume 4502 of *Lecture Notes in Computer Science*, chapter A Formalisation of a Dependently Typed Language as an Inductive-Recursive Family, pages 93–109. Springer Berlin Heidelberg, Berlin, Heidelberg, 2007. ISBN 978-3-540-74464-1. doi: 10.1007/978-3-540-74464-1_7. URL http://dx.doi.org/10.1007/978-3-540-74464-1_7.
- K. Gödel. Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I. *Monatshefte für mathematik und physik*, 38(1):173–198, 1931. URL <http://www.w-k-essler.de/pdfs/goedel.pdf>.
- D. R. Hofstadter. *Gödel, Escher, Bach: An Eternal Golden Braid*. Vintage, 1979. ISBN 978-0394745022.
- S. C. Kleene. *Introduction to Metamathematics*. Wolters-Noordhoff, 1952. ISBN 0-7204-2103-9.
- 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>.
- R. O’Connor. Essential incompleteness of arithmetic verified by coq. *CoRR*, abs/cs/0505034, 2005. URL <http://arxiv.org/abs/cs/0505034>.
- L. C. Paulson. A mechanised proof of Gödel’s incompleteness theorems using Nominal Isabelle. *Journal of Automated Reasoning*, 55(1):1–37, 2015. URL <https://www.cl.cam.ac.uk/~lp15/papers/Formath/Goedel-ar.pdf>.
- D. Piponi. From löb’s theorem to spreadsheet evaluation, November 2006. URL <http://blog.sigfpe.com/2006/11/from-l-theorem-to-spreadsheet.html>.
- G. K. Pullum. Scooping the loop snooper, October 2000. URL <http://www.lel.ed.ac.uk/~gpullum/loopsnoop.html>.
- N. Shankar. *Proof-checking Metamathematics (Theorem-proving)*. PhD thesis, The University of Texas at Austin, 1986. AAI8717580.
- N. Shankar. *Metamathematics, Machines and Gödel’s Proof*. Cambridge University Press, 1997.
- R. J. Simmons and B. Toninho. Constructive provability logic. *CoRR*, abs/1205.6402, May 2012. URL <http://arxiv.org/abs/1205.6402>.
- A. Tarski. Der Wahrheitsbegriff in den formalisierten Sprachen. *Studia Philosophica*, 1:261–405, 1936. URL <http://www.w-k-essler.de/pdfs/Tarski.pdf>.
- B. Yudkowsky. Löb’s theorem cured my social anxiety, February 2014. URL <http://agentyduck.blogspot.com/2014/02/lobs-theorem-cured-my-social-anxiety.html>.