

Löb's Theorem

A functional pearl of dependently typed quining

Jason Gross

MIT CSAIL
jgross@mit.edu

Name2 Name3

Affiliation2/3
Email2/3

Categories and Subject Descriptors CR-number [subcategory]:
third-level

General Terms Agda, Lob, quine, self-reference

Keywords Agda, Lob, quine, self-reference

Abstract

This is the text of the abstract.

*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* (Pullum 2000))

TODO

- cite Using Reflection to Explain and Enhance Type Theory?

1. Introduction

Löb's theorem has a variety of applications, from proving incompleteness of a logical theory as a trivial corollary, to acting as a no-go theorem for a large class of self-interpreters (TODO: mention F_ω ?), from allowing robust cooperation in the Prisoner's Dilemma with Source Code (Barasz et al. 2014), to 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 $A \rightarrow B$

are the same, if we assume A , we are also assuming $A \rightarrow B$, and hence we have B , and since assuming A yields B , we have that $A \rightarrow B$. 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.

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? 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 math-

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 counterfactual "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

ematicians 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 program `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: □ (T % repr(T)) -> X)
  ("(lambda T: □ (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 Π , 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 _▷_

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

data Type : Context -> Set where
  'T' : ∀ {Γ} -> Type Γ
  '⊥' : ∀ {Γ} -> Type Γ
  'II' : ∀ {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 {Γ ▷ 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
  [ε]ᶜ = ⊤
  [Γ ▷ T]ᶜ = Σ [Γ]ᶜ [T]ᵀ

  [_]ᵀ : ∀ {Γ} -> Type Γ -> [Γ]ᶜ -> Set
  ['T']ᵀ [Γ] = ⊤
  ['⊥']ᵀ [Γ] = ⊥
  ['II' 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)
```

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 sup-

port 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:

```
löp : ∀ {X} → □ (□ X → X) → [X]ᵀ
löp 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 a self-interpreter for F_ω (Brown and Palsberg 2016).²

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

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-Emptyness

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

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

```
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 → '⊥'
```

```
löp : ∀ {X} → □ (□ X → X) → [X]ᵀ
löp f = [Löb f]ᵀ
```

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

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

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

```
no-interpreters : ¬ (∀ {X} → □ (□ X → X))
no-interpreters interp = löp (interp ['⊥'])
```

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 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 $\phi \vdash \text{Quine } \phi$, where $\vdash \text{Quine } \phi$ 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`.

```
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. 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
  '→' : ∀ {Γ} → Type Γ → Type Γ → Type Γ
  '⊢' : ∀ {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))
  _'o'_ : ∀ {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 \vdash_{repr} for the type-level add-quote function **TODO: should this be called add-quote?**, and \vdash_{repr} 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 `'VAR0'tt`, 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' : ∀ {Γ A B} → Term {Γ ⊢ A} (W B) → Term (A '→' B)
```

```
'tt' : ∀ {Γ} → Term {Γ} '⊢'
Γ ⊢tt : ∀ {Γ} → Type ε → Term {Γ} 'Typeε'
Γ ⊢tt : ∀ {Γ T} → Term {ε} T → Term {Γ} ('□' '→' T ⊢tt)
quine→ : ∀ {φ} → Term {ε} (Quine φ '→' φ '→' Quine φ ⊢tt)
quine← : ∀ {φ} → Term {ε} (φ '→' Quine φ ⊢tt '→' Quine φ)
'VAR0' : ∀ {Γ T} → Term {Γ ⊢ T} (W T)
_ 'a_ : ∀ {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)
'Γ'VAR0'tt : ∀ {T}
  → Term {ε ⊢ '□' '→' T ⊢tt} (W ('□' '→' '□' '→' T ⊢tt ⊢tt))
→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)
w : ∀ {Γ A T} → Term {Γ} A → Term {Γ ⊢ T} (W A)
w→ : ∀ {Γ A B X}
  → Term {Γ} (A '→' B)
  → Term {Γ ⊢ X} (W A '→' W B)
_ wtt _ : ∀ {A B T}
  → Term {ε ⊢ T} (W ('□' '→' A '→' B ⊢tt))
  → Term {ε ⊢ T} (W ('□' '→' A ⊢tt))
  → Term {ε ⊢ T} (W ('□' '→' B ⊢tt))
```

```
□ : 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 `'□'`, 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 any higher level
```

```
mutual
[ ]c : (Γ : Context) → Set (|suc 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.proj2 [Γ]))
[ Quine φ ]T [Γ] = [ φ ]T ([Γ], lift (Quine φ))
[ W T ]T [Γ] = [ T ]T (Σ.proj1 [Γ])
[ W1 T ]T [Γ] = [ T ]T ((Σ.proj1 (Σ.proj1 [Γ])), (Σ.proj2 [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
[ Γ x ⊥ ]t [Γ] = lift x
[ Γ x ⊥ ]t [Γ] = lift x
[ quine → ]t [Γ] x = x
[ quine ← ]t [Γ] x = x
[ 'VAR0' ]t [Γ] = Σ.proj2 [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 [Γ] = lift Γ lower (Σ.proj2 [Γ]) ⊥
[ ←SW1SV → W f ]t = [ f ]t
[ →SW1SV → W f ]t = [ f ]t
[ w x ]t [Γ] = [ x ]t (Σ.proj1 [Γ])
[ w → f ]t [Γ] = [ f ]t (Σ.proj1 [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 and 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 $\Box 'H' \rightarrow 'X'$; 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' ⊥ '→' 'X'))
  where
    'H' : Type ε
    'H' = Quine (W1 '□' 'Γ' 'VAR0' '→' W 'X')

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

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

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

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

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

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

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

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

    löb : ∀ {X} → □ ('□' 'Γ' 'X' ⊥ '→' 'X') → [ 'X' ]

```

```

löb 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' '⊥' ⊥))
incompleteness = löb

```

```

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

```

```

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

TODO: Explain Prisoner's dilemma

(Barasz et al. 2014)

open lob

```

-- a bot takes in the source code for itself,
-- for another bot, and spits out the assertion
-- that it cooperates with this bot

```

```

'Bot' : ∀ {Γ} → Type Γ

```

```

'Bot' {Γ}
  = Quine (W1 'Term' 'Γ' 'VAR0'
    '→' W1 'Term' 'Γ' 'VAR0'
    '→' W ('Type' Γ))

```

```

_cooperates-with_ : □ 'Bot' → □ 'Bot' → Type ε

```

```

b1 cooperates-with b2 = lower ([ b1 ]t tt (lift b1) (lift b2))

```

```

'eval-bot' : ∀ {Γ}

```

```

  → Term {Γ} ('Bot' '→' ('□' 'Bot' '→' '□' 'Bot' '→' 'Type' Γ))

```

```

'eval-bot' = →SW1SV → SW1SV → W quine→

```

```

'eval-bot' : ∀ {Γ}

```

```

  → Term {Γ} ('□' 'Bot'

```

```

    '→' '□' ({- other -} '□' 'Bot' '→' 'Type' Γ))

```

```

'eval-bot' = 'λ' (w Γ 'eval-bot' ⊥ w 'a' 'VAR0' w 'a' 'Γ' 'VAR0' ⊥)

```

```

'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')

```

```

'self' : ∀ {Γ}
  → Term {Γ ▷ '□' 'Bot' ▷ W ('□' 'Bot')}
    (W (W ('□' 'Bot')))
'self' = w 'VAR0'

'other' : ∀ {Γ}
  → Term {Γ ▷ '□' 'Bot' ▷ W ('□' 'Bot')}
    (W (W ('□' 'Bot')))
'other' = 'VAR0'

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))

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 'Γ' '⊥' '⊃' '⊣')

'DefectBot' : □ 'Bot'
'CooperateBot' : □ 'Bot'
'FairBot' : □ 'Bot'
'PrudentBot' : □ 'Bot'

'DefectBot' = make-bot (w (w 'Γ' '⊥' '⊃'))
'CooperateBot' = make-bot (w (w 'Γ' '⊣' '⊃'))
'FairBot' = make-bot ('□' ('other-cooperates-with' 'a' 'self'))
'PrudentBot'
  = make-bot ('□'
    (φ0 ww'''_ ×'''
      (¬□⊥ ww'''_ →''' other-defects-against-DefectBot)))
where
  φ0 : ∀ {Γ}
    → Term {Γ ▷ '□' 'Bot' ▷ W ('□' 'Bot')}
      (W (W ('□' ('Type' Γ))))
    → Term {Γ} 'Bot'
  φ0 = 'other-cooperates-with' 'a' 'self'

other-defects-against-DefectBot
  : Term {Γ ▷ '□' 'Bot' ▷ W ('□' 'Bot')}
    (W (W ('□' ('Type' Γ))))
other-defects-against-DefectBot
  = ww'''_
    ('other-cooperates-with' 'a' w (w 'Γ' 'DefectBot' '⊣'))

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

```

10. Encoding with Add-Quote Function

- (appendix) - Discuss whiteboard phrasing of sentence with sigmas
- It remains to show that we can construct
 - Discuss whiteboard phrasing of untyped sentence
 - Given:
 - X
 - □ = Term
 - f : □ 'X' → X
 - define y : X

- Suppose we have a type $H \cong \text{Term } \ulcorner H \rightarrow X \urcorner$, and we have
- toH : $\text{Term } \ulcorner H \rightarrow X \urcorner \rightarrow H$
- fromH : $H \rightarrow \text{Term } \ulcorner H \rightarrow X \urcorner$
- quote : $H \rightarrow \text{Term } \ulcorner H \urcorner$
-
- Then we can define
- $y = (\lambda h : H. f (\text{subst } (\text{quote } h) h) (\text{toH } \ulcorner h : H. f (\text{subst } (\text{quote } h) h) \urcorner}))$

11. Removing add-quote and actually tying the knot (future work 1)

- Temporary outline section to be moved
-

How do we construct the Curry–Howard analogue of the Löbian sentence? A quine is a program that outputs its own source code (). We will say that a *type-theoretic quine* is a program that outputs its own (well-typed) abstract syntax tree. Generalizing this slightly, we can consider programs that output an arbitrary function of their own syntax trees.

- TODO: Examples of double quotation, single quotation, etc.
- Given any function ϕ from doubly-quoted syntactic types to singly-quoted syntactic types, and given an operator $\ulcorner _ \urcorner$ which adds an extra level of quotation, we can define the type of a *quine* at ϕ to be a (syntactic) type "Quine ϕ " which is isomorphic to " ϕ ($\ulcorner \text{Quine } \phi \urcorner$)".
- What's wrong is that self-reference with truth is impossible. In a particular technical sense, it doesn't terminate. Solution: Provability
- Quining / self-referential provability sentence and provability implies truth
- Curry–Howard, quines, abstract syntax trees (This is an interpreter!)

A. Standard Data-Type Declarations

```

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

infixl 1 _→_
infixr 2 _×_
infixl 1 _≡_

record T {ℓ} : Set ℓ where
  constructor tt

data ⊥ {ℓ} : Set ℓ where

record Σ {a p} (A : Set a) (P : A → Set p) : Set (a ⊔ p) where
  constructor _,_
  field
    proj1 : A
    proj2 : P proj1

data Lifted {a b} (A : Set a) : Set (b ⊔ a) where
  lift : A → Lifted A

lower : ∀ {a b A} → Lifted {a} {b} A → A
lower (lift x) = x

_×_ : ∀ {ℓ ℓ'} (A : Set ℓ) (B : Set ℓ') → Set (ℓ ⊔ ℓ')
A × B = Σ A (λ _ → B)

data _≡_ {ℓ} {A : Set ℓ} (x : A) : A → Set ℓ where
  refl : x ≡ x

```


Encoding of Löb's Theorem for the Prisoner's Dilemma

```

data Type : Context → Set where
  W : ∀ {Γ A} → Type Γ → Type (Γ ▷ A)
  W₁ : ∀ {Γ A B} → Type (Γ ▷ B) → Type (Γ ▷ A)
  " " : ∀ {Γ A} → Type (Γ ▷ A) → Term {Γ} A
  'Type' : ∀ Γ → Type Γ
  'Term' : ∀ {Γ} → Type (Γ ▷ 'Type' Γ)
  '→' : ∀ {Γ} → Type Γ → Type Γ → Type Γ
  '×' : ∀ {Γ} → Type Γ → Type Γ → Type Γ
  'Quine' : ∀ {Γ} → Type (Γ ▷ 'Type' Γ) → Type Γ
  'T' : ∀ {Γ} → Type Γ
  '⊥' : ∀ {Γ} → Type Γ

```

$$\begin{aligned}
& \rightarrow \text{Term } \{\Gamma\} ((A \text{ ' } x \rightarrow B) \text{ ' } \rightarrow T) \\
& \rightarrow \text{SW}_1 \text{SV} \rightarrow \text{SW}_2 \text{SV} \rightarrow W : \forall \{\Gamma T X A B\} \{x : \text{Term } X\} \\
& \rightarrow \text{Term } \{\Gamma\} (T \text{ ' } \rightarrow (W_1 A \text{ ' } \text{'VAR}_0 \text{ ' } \rightarrow W_1 A \text{ ' } \text{'VAR}_0 \text{ ' } \rightarrow W) \\
& \rightarrow \text{Term } \{\Gamma\} (T \text{ ' } \rightarrow A \text{ ' } x \rightarrow A \text{ ' } x \rightarrow B) \\
& \leftarrow \text{SW}_1 \text{SV} \rightarrow \text{SW}_2 \text{SV} \rightarrow W : \forall \{\Gamma T X A B\} \{x : \text{Term } X\} \\
& \rightarrow \text{Term } \{\Gamma\} ((W_1 A \text{ ' } \text{'VAR}_0 \text{ ' } \rightarrow W_1 A \text{ ' } \text{'VAR}_0 \text{ ' } \rightarrow W B) \text{ ' } x \\
& \rightarrow \text{Term } \{\Gamma\} ((A \text{ ' } x \rightarrow A \text{ ' } x \rightarrow B) \text{ ' } \rightarrow T) \\
w : \forall \{\Gamma A T\} \rightarrow \text{Term } \{\Gamma\} A \rightarrow \text{Term } \{\Gamma \triangleright T\} (W A) \\
w \rightarrow : \forall \{\Gamma A B X\} \rightarrow \text{Term } \{\Gamma \triangleright X\} (W (A \text{ ' } \rightarrow B)) \rightarrow \text{Term } \{\Gamma \triangleright X\} \\
\rightarrow w : \forall \{\Gamma A B X\} \rightarrow \text{Term } \{\Gamma \triangleright X\} (W A \text{ ' } \rightarrow W B) \rightarrow \text{Term } \{\Gamma \triangleright X\} \\
ww \rightarrow : \forall \{\Gamma A B X Y\} \rightarrow \text{Term } \{\Gamma \triangleright X \triangleright Y\} (W (W (A \text{ ' } \rightarrow B))) \rightarrow \\
\rightarrow ww : \forall \{\Gamma A B X Y\} \rightarrow \text{Term } \{\Gamma \triangleright X \triangleright Y\} (W (W A) \text{ ' } \rightarrow W (W B) \\
\text{ ' } \circ \text{ ' } : \forall \{\Gamma A B C\} \rightarrow \text{Term } \{\Gamma\} (B \text{ ' } \rightarrow C) \rightarrow \text{Term } \{\Gamma\} (A \text{ ' } \rightarrow C) \\
\text{ ' } \text{w''}_a : \forall \{\Gamma A B T\} \rightarrow \text{Term } \{\Gamma \triangleright T\} (W ('Term' \text{ ' ' } \Gamma A \text{ ' } \rightarrow B \text{ ' } \rightarrow T) \\
\text{ ' } \text{w''}_a : \forall \{\Gamma A B\} \rightarrow \text{Term } \{\Gamma\} ('Term' \text{ ' ' } \Gamma A \text{ ' } \rightarrow B \text{ ' } \rightarrow \text{'Term' ' ' } \Gamma A \\
\text{ ' } \text{w''}_a : \forall \{\Gamma A B T\} \rightarrow \text{Term } \{\Gamma \triangleright T\} ('Type' (\Gamma \triangleright T) (W (A \text{ ' } \rightarrow B) \text{ ' } \rightarrow T) \\
\text{ ' } \text{w''}_a : \forall \{\Gamma A B\} \rightarrow \text{Term } \{\Gamma\} ('Type' (\Gamma \triangleright A) (W (A \text{ ' } \rightarrow B) \text{ ' } \rightarrow A) \\
\text{ ' } \text{w''}_a : \forall \{\Gamma A\} \rightarrow \text{Term } \{\Gamma \triangleright A\} ('Type' (\Gamma \triangleright A) (W (A \text{ ' } \rightarrow B) \text{ ' } \rightarrow A) \\
\text{ ' } \text{w''}_a : \forall \{\Gamma\} \rightarrow \text{Term } \{\Gamma\} ('Type' \Gamma) \rightarrow \text{Term } \{\Gamma\} ('Type' \Gamma) \\
\text{ ' } \text{w''}_a : \forall \{\Gamma A B\} \rightarrow \text{Term } \{\Gamma \triangleright A \triangleright B\} (W (W ('Term' \text{ ' ' } \Gamma A \text{ ' } \rightarrow B) \text{ ' } \rightarrow T) \\
\text{ ' } \text{w''}_a : \forall \{\Gamma A B\} \rightarrow \text{Term } \{\Gamma \triangleright A \triangleright B\} (W (W ('Term' \text{ ' ' } \Gamma A \text{ ' } \rightarrow B) \text{ ' } \rightarrow T) \\
\text{ ' } \text{w''}_a : \forall \{\Gamma A B\} \rightarrow \text{Term } \{\Gamma \triangleright A \triangleright B\} (W (W ('Term' \text{ ' ' } \Gamma A \text{ ' } \rightarrow B) \text{ ' } \rightarrow T)
\end{aligned}$$
$$\begin{aligned} & \llbracket \text{if } \epsilon \text{ then } A \text{ else } B \rrbracket^c : (\Gamma : \text{Context}) \rightarrow \text{Set} \text{ (Isuc max-level)} \\ & \llbracket \epsilon \rrbracket^c = \top \\ & \llbracket \Gamma \triangleright T \rrbracket^c = \Sigma \llbracket \Gamma \rrbracket^c \llbracket T \rrbracket^\top \end{aligned}$$

2016/2/29

$$\begin{aligned} \text{Löb} &: \forall \{X\} \rightarrow \text{Term } \{\varepsilon\} \text{ ('}\square\text{' } X \text{'}\rightarrow\text{' } X) \rightarrow \text{Term } \{\varepsilon\} X \\ \text{Löb } \{X\} f &= \text{inner.Löb } X f \end{aligned}$$

2016/2/29

9

```

module well-typed-syntax-helpers where
  open well-typed-syntax

```

```

infixl 3 _'a_
infixr 1 _'→'_
infixl 3 _'t'_
infixl 3 _'t'_1_
infixl 3 _'t'_2_
infixr 2 _'o'_

```

```

_ '→' _ : ∀ {Γ} → Type Γ → Type Γ → Type Γ
_ '→' _ {Γ} A B = _ '→' _ {Γ} A (W {Γ} {A} B)

```

```

_ 'a' _ : ∀ {Γ A B} → Term {Γ} (A '→' B) → Term A → Term B
_ 'a' _ {Γ} {A} {B} f x = SW ( _ 'a' _ {Γ} {A} {W B} f x )

```

```

_ 't' _ : ∀ {Γ A} {B : Type (Γ ▷ A)} → (b : Term {Γ ▷ A} B) → (a : Term {Γ ▷ A} B)
b 't' a = 'λ' b 'a a

```

```

substType-tProd : ∀ {Γ T A B} {a : Term {Γ} T} →
  Term {Γ} ((A '→' B) 'a)
  → Term {Γ} ( _ '→' _ {Γ} (A 'a) (B '1 a) )
substType-tProd {Γ} {T} {A} {B} {a} x = SW ((WSV (w x)) 't' a)

```

```

SV = substType-tProd

```

```

'λ' •' : ∀ {Γ A B} → Term {Γ ▷ A} (W B) -> Term (A '→' B)
'λ' •' f = 'λ' f

```

```

un'λ' : ∀ {Γ A B} → Term (A '→' B) → Term {Γ ▷ A} B
un'λ' f = SW1V (WV (w f) 'a 'VAR0')

```

```

weakenProd : ∀ {Γ A B C} →
  Term {Γ} (A '→' B)
  → Term {Γ ▷ C} (W A '→' W1 B)
weakenProd {Γ} {A} {B} {C} x = WV (w x)
wV = weakenProd

```

```

w1 : ∀ {Γ A B C} → Term {Γ ▷ B} C → Term {Γ ▷ A ▷ W {Γ} {A} B} (W1 {Γ} {A} {B = B} C)
w1 x = un'λ' (WV (w ('λ' x)))

```

```

_ 't'_1_ : ∀ {Γ A B C} → (c : Term {Γ ▷ A ▷ B} C) → (a : Term {Γ} A)
f 't'_1 x = un'λ' (SV ('λ' ('λ' f) 'a x))
_ 't'_2_ : ∀ {Γ A B C D} → (c : Term {Γ ▷ A ▷ B ▷ C} D) → (a : Term {Γ} A)
f 't'_2 x = un'λ' (S1V (un'λ' (SV ('λ' ('λ' ('λ' f) 'a x))))

```

```

S10W' : ∀ {Γ C T A} {a : Term {Γ} C} {b : Term {Γ} (T 'a)} → Term {Γ} (A '→' (S10W (T 'a)))
S10W' = S10W-1

```

```

S10W-weakenType : ∀ {Γ T A} {B : Type (Γ ▷ A)}
  → {a : Term {Γ} A}
  → {b : Term {Γ} (B 'a)}
  → Term {Γ} (W (W T) '1 a 'b)
  → Term {Γ} T
S10W-weakenType x = SW (S10W x)

```

```

S10WW = S10W-weakenType

```

```

S210W-weakenType : ∀ {Γ A B C T}
  {a : Term {Γ} A}
  {b : Term {Γ} (B 'a)}

```

```

{c : Term {Γ} (C '1 a 'b)} →
  Term {Γ} (W (W T) '2 a '1 b 'c)
  → Term {Γ} (T 'a)
S210W-weakenType x = S10W (S210W x)

```

```

S210WW = S210W-weakenType

```

```

W101W : ∀ {Γ A B C T} → Term {Γ ▷ A ▷ B ▷ W (W C)} (W1 (W1 (W T)))
W101W = W1010

```

```

WV-nd : ∀ {Γ A B C} →
  Term {Γ ▷ C} (W (A '→' B))
  → Term {Γ ▷ C} (W A '→' W B)
WV-nd x = 'λ' •' (W10 (SW1V (WV (w (WV x)) 'a 'VAR0)))

```

```

weakenProd-nd : ∀ {Γ A B C} →
  Term {Γ} (A '→' B)
  → Term {Γ ▷ C} (W A '→' W B)
weakenProd-nd {Γ} {A} {B} {C} x = WV-nd (w x)

```

```

WV-nd-tProd-nd : ∀ {Γ A B C D} →
  Term {Γ ▷ D} (W (A '→' B '→' C))
  → Term {Γ ▷ D} (W A '→' W B '→' W C)
WV-nd-tProd-nd x = 'λ' (WV-1 ('λ' (W101W (SW1V (wV (WV (WWW x))))

```

```

weakenProd-nd-Prod-nd : ∀ {Γ A B C D} →
  Term (A '→' B '→' C)
  → Term {Γ ▷ D} (W A '→' W B '→' W C)
weakenProd-nd-Prod-nd {Γ} {A} {B} {C} {D} x = WV-nd-tProd-nd (w x)
w → → = weakenProd-nd-Prod-nd

```

```

W1S1W' : ∀ {Γ A T'' T' T} {a : Term {Γ} A}
  → Term {Γ ▷ T'' ▷ W (T'' 'a)} (W1 (W (T'' 'a)))
  → Term {Γ ▷ T'' ▷ W (T'' 'a)} (W1 (W (T'' 'a)))
W1S1W' = weakenType1-substType-weakenType1-inv

```

```

substType-weakenType1-inv : ∀ {Γ A T' T}
  {a : Term {Γ} A}
  {b : Term {Γ} (B 'a)}
  → Term {Γ ▷ T'} (W (T' 'a))
  → Term {Γ ▷ T'} (W (T' 'a))
substType-weakenType1-inv {a = a} x = S1W1 (W1S1W' (w1 x) 't'_1 a)

```

```

S1W' = S1W-weakenType1-inv

```

```

_ 'o' _ : ∀ {Γ A B C}
  → Term {Γ} (A '→' B)
  → Term {Γ} (B '→' C)
  → Term {Γ} (A '→' C)
g 'o' f = 'λ' (w → f 'a (w → g 'a 'VAR0'))

```

```

WS00W1 : ∀ {Γ T' B A} {b : Term {Γ} B} {a : Term {Γ ▷ B} (W A)} {T}
  → Term {Γ ▷ T'} (W (W1 T'' 'a 'b))
  → Term {Γ ▷ T'} (W (T'' (SW (a 't' b))))
WS00W1 = weakenType-substType-substType-weakenType1

```

```

WS00W1' : ∀ {Γ T' B A} {b : Term {Γ} B} {a : Term {Γ ▷ B} (W A)} {T}
  → Term {Γ ▷ T'} (W (T'' (SW (a 't' b))))

```

open well-typed-syntax-context-helpers

' ε ': Term { ε } 'Context'
' ε ' = $\ulcorner \varepsilon \urcorner^c$

$\ulcorner \square \urcorner : \text{Type} \ (\varepsilon \triangleright \ulcorner \text{Type} \urcorner \text{ '' } \ulcorner \varepsilon \urcorner)$
 $\ulcorner \square \urcorner = \ulcorner \text{Term} \urcorner \text{ '' }_1 \ulcorner \varepsilon \urcorner$

```
module well-typed-syntax-eq-dec where
  open well-typed-syntax
```

$$\text{context-pick-if} : \forall \{\ell\} \{P : \text{Context} \rightarrow \text{Set } \ell\} \\ \{\Gamma : \text{Context}\} \\ (\text{dummy} : P (\varepsilon \triangleright \text{'}\Sigma \text{' 'Context' 'Type'})) \\ (\text{val} : P \Gamma) \rightarrow$$

$P(\varepsilon \triangleright \text{'}\Sigma\text{' 'Context' 'Type'})$
 $\text{context-pick-if } \{P = P\} \cdot \{\varepsilon \triangleright \text{'}\Sigma\text{' 'Context' 'Type'}\} \text{ dummy val} = \text{val}$
 $\text{context-pick-if } \{P = P\} \cdot \{\Gamma\} \text{ dummy val} = \text{dummy}$

$$\begin{aligned} \text{context-pick-if-refl} &: \forall \{ \ell \ P \ dummy \ val \} \rightarrow \\ &\quad \text{context-pick-if} \{ \ell \} \{ P \} \{ \varepsilon \triangleright ' \Sigma ' \text{ 'Context' 'Type' } \} \ dummy \ val \equiv val \\ \text{context-pick-if-refl} \{ P = P \} &= \text{refl} \end{aligned}$$

```
module well-typed-quoted-syntax where
  open well-typed-syntax
  open well-typed-syntax-helpers public
  open well-typed-quoted-syntax-defs public
  open well-typed-syntax-context-helpers public
  open well-typed-syntax-eq-dec public
```

infixr 2 _“o”_

$$\begin{aligned} \text{quote-sigma} &: (\Gamma v : \Sigma \text{ Context Type}) \rightarrow \text{Term } \{\varepsilon\} \text{ ('}\Sigma \text{' 'Context' 'Type'} \\ \text{quote-sigma } (\Gamma, v) &= \text{'existT' } \ulcorner \Gamma \urcorner^c \ulcorner v \urcorner^T \end{aligned}$$
$$\begin{aligned}
- \text{“}\circ\text{”} &: \forall \{A \ B \ C\} \\
&\rightarrow \square \left(\text{“}\square\text{”} \text{ “} (C \text{ “}\rightarrow\text{”} B) \right) \\
&\rightarrow \square \left(\text{“}\square\text{”} \text{ “} (A \text{ “}\rightarrow\text{”} C) \right) \\
&\rightarrow \square \left(\text{“}\square\text{”} \text{ “} (A \text{ “}\rightarrow\text{”} B) \right) \\
g \text{ “}\circ\text{”} f &= (\text{“}f\text{comp-nd”} \text{ “}a \text{ “}f\text{”} \text{ “}a \text{ “}g\text{”})
\end{aligned}$$
[illegible]

```

▷  $T_1^1(W_1, B, a)$ 
module well-typed-syntax-pre-interpretor where
  open well-typed-syntax
  open well-typed-syntax-helpers

```

```
max-level : Level
max-level = |suc |zero
```

```

module inner
  (context-pick-if' : ∀ ℓ (P : Context → Set ℓ)
    (Γ : Context)
    (dummy : P (ε ▷ 'Σ' 'Context' 'Type'))
    (val : P Γ) →
    P (ε ▷ 'Σ' 'Context' 'Type'))

```

```

(context-pick-if-refl' : ∀ ℓ P dummy val →
  context-pick-if' ℓ P (ε ▷ 'Σ' 'Context' 'Type') dummy val ≡ val)
where

context-pick-if : ∀ {ℓ} {P : Context → Set ℓ}
  {Γ : Context}
  (dummy : P (ε ▷ 'Σ' 'Context' 'Type'))
  (val : P Γ) →
P (ε ▷ 'Σ' 'Context' 'Type')
context-pick-if {P = P} dummy val = context-pick-if' _ P _ dummy val
context-pick-if-refl : ∀ {ℓ P dummy val} →
  context-pick-if {ℓ} {P} {ε ▷ 'Σ' 'Context' 'Type'} dummy val ≡ val
context-pick-if-refl {P = P} = context-pick-if-refl' _ P _ _

private
dummy : Type ε
dummy = 'Context'

cast-helper : ∀ {X T A} {x : Term X} → A ≡ T → Term {ε} (T "x →" A)
cast-helper refl = 'λ' 'VAR0'

cast'-proof : ∀ {T} → Term {ε} (context-pick-if {P = Type} (W dummy) T
  "→" T "existT' Γ ε ▷ 'Σ' 'Context' 'Type' Γ T T")
cast'-proof {T} = cast-helper {'Σ' 'Context' 'Type'}
  {context-pick-if {P = Type} {ε ▷ 'Σ' 'Context' 'Type'} (W dummy) T
  {T} (sym (context-pick-if-refl {P = Type} {dummy = W dummy}))}

cast-proof : ∀ {T} → Term {ε} (T "existT' Γ ε ▷ 'Σ' 'Context' 'Type' Γ T T
  "→" context-pick-if {P = Type} (W dummy) T "existT' Γ ε ▷ 'Σ' 'Context'
  "Type" T)
cast-proof {T} = cast-helper {'Σ' 'Context' 'Type'} {T}
  {context-pick-if {P = Type} {ε ▷ 'Σ' 'Context' 'Type'} (W dummy) T
  (context-pick-if-refl {P = Type} {dummy = W dummy})}

'idfun' : ∀ {T} → Term {ε} (T "→" T)
'idfun' = 'λ' 'VAR0'

mutual
[ ]c : (Γ : Context) → Set (lsuc max-level)
[ ]T : {Γ : Context} → Type Γ → [ ]c Γ → Set max-level

[ ]c ε = ⊤
[ ]c (Γ ▷ T) = Σ ([ ]c Γ) (λ Γ' → [ ]T T Γ')

[ ]T (T1 "x") [Γ] = [ ]T T1 ([Γ], [ ]t x [Γ])
[ ]T (T2 "1 a") ([Γ], A↓) = [ ]T T2 ([Γ], [ ]t a [Γ], A↓)
[ ]T (T3 "2 a") ([Γ], A↓, B↓) = [ ]T T3 ([Γ], [ ]t a [Γ], A↓, B↓)
[ ]T (T3 "3 a") ([Γ], A↓, B↓, C↓) = [ ]T T3 ([Γ], [ ]t a [Γ], A↓, B↓, C↓)
[ ]T (W T1) ([Γ], _) = [ ]T T1 [Γ]
[ ]T (W1 T2) ([Γ], A↓, B↓) = [ ]T T2 ([Γ], B↓)
[ ]T (W2 T3) ([Γ], A↓, B↓, C↓) = [ ]T T3 ([Γ], B↓, C↓)
[ ]T (T "→" T1) [Γ] = (T↓ : [ ]T T [Γ]) → [ ]T T1 ([Γ], T↓)
[ ]T 'Context' [Γ] = Lifted Context
[ ]T 'Type' ([Γ], T↓) = Lifted (Type (lower T↓))
[ ]T 'Term' ([Γ], T↓, t↓) = Lifted (Term (lower t↓))
[ ]T ('Σ' T T1) [Γ] = Σ ([ ]T T [Γ]) (λ T↓ → [ ]T T1 ([Γ], T↓))

[ ]t : ∀ {Γ : Context} {T : Type Γ} → Term T → ([Γ] : [ ]c Γ) → [ ]t T [Γ]
[ ]t (w t) ([Γ], A↓) = [ ]t t [Γ]
[ ]t ('λ' t) [Γ] T↓ = [ ]t t ([Γ], T↓)
[ ]t (t "a t1) [Γ] = [ ]t t [Γ] ([ ]t t1 [Γ])

[ ]t 'VAR0' ([Γ], A↓) = A↓
[ ]t (Γ Γc) [Γ] = lift Γ
[ ]t (Γ T TT) [Γ] = lift T
[ ]t (Γ t Tt) [Γ] = lift t
[ ]t 'quote-term' [Γ] (lift T↓) = lift Γ T↓ Tt
[ ]t ('quote-sigma' {Γ0} {Γ1}) [Γ] (lift Γ, lift T) = lift ('existT' {Γ0} {Γ1} T)
[ ]t 'cast' [Γ] T↓ = lift (context-pick-if
  {P = Type}
  {lower (Σ.proj1 T↓)})
  (W dummy)
  (lower (Σ.proj2 T↓)))
[ ]t (SW t) [Γ] = [ ]t t [Γ]
[ ]t (WS∀ t) [Γ] T↓ = [ ]t t [Γ] T↓
[ ]t (SW1V t) [Γ] = [ ]t t [Γ]
[ ]t (W∀ t) [Γ] T↓ = [ ]t t [Γ] T↓
[ ]t (W∀-1 t) [Γ] T↓ = [ ]t t [Γ] T↓
[ ]t (WW∀ t) [Γ] T↓ = [ ]t t [Γ] T↓
[ ]t (S1∀ t) [Γ] T↓ = [ ]t t [Γ] T↓
[ ]t (S10W-1 t) [Γ] = [ ]t t [Γ]
[ ]t (S10W t) [Γ] = [ ]t t [Γ]
[ ]t (WWS10W t) [Γ] = [ ]t t [Γ]
[ ]t (WS210Wε t) [Γ] = [ ]t t [Γ]
[ ]t (S210W t) [Γ] = [ ]t t [Γ]
[ ]t (W10 t) [Γ] = [ ]t t [Γ]
[ ]t (W10-inv t) [Γ] = [ ]t t [Γ]
[ ]t (W1010 t) [Γ] = [ ]t t [Γ]
[ ]t (S1W1 t) [Γ] = [ ]t t [Γ]
[ ]t (weakenType1-substType-weakenType1-inv t) [Γ] = [ ]t t [Γ]
[ ]t (weakenType1-substType-weakenType1 t) [Γ] = [ ]t t [Γ]
[ ]t (weakenType-substType-substType-weakenType1 t) [Γ] = [ ]t t [Γ]
[ ]t (weakenType-substType-substType-weakenType1-inv t) [Γ] = [ ]t t [Γ]
[ ]t (substType-W10 t) [Γ] = [ ]t t [Γ]
[ ]t (WS210W1 t) [Γ] = [ ]t t [Γ]
[ ]t (substType1-substType-tProd t) [Γ] T↓ = [ ]t t [Γ] T↓
[ ]t (substType2-substType-substType-W10-weakenType t) [Γ] = [ ]t t [Γ]
[ ]t (S10W2-weakenType t) [Γ] = [ ]t t [Γ]
[ ]t (weakenType-W10 t) [Γ] = [ ]t t [Γ]
[ ]t (beta-under-subst t) [Γ] = [ ]t t [Γ]
[ ]t 'proj1' [Γ] (x, p) = x
[ ]t 'proj2' ([Γ], (x, p)) = p
[ ]t ('existT' x p) [Γ] = [ ]t x [Γ], [ ]t p [Γ]
[ ]t (f "x") [Γ] = lift (lower ([ ]t f [Γ]) "lower ([ ]t x [Γ]))
[ ]t (f w "x") [Γ] = lift (lower ([ ]t f [Γ]) "lower ([ ]t x [Γ]))
[ ]t (f "→" x) [Γ] = lift (lower ([ ]t f [Γ]) "→" lower ([ ]t x [Γ])
[ ]t (f w "→" x) [Γ] = lift (lower ([ ]t f [Γ]) "→" lower ([ ]t x [Γ])
[ ]t (f w "→" x) [Γ] = lift (lower ([ ]t f [Γ]) "→" lower ([ ]t x [Γ])
[ ]t (w → x) [Γ] A↓ = [ ]t x (Σ.proj1 [Γ]) A↓
[ ]t ("→" → w "→") [Γ] T↓ = T↓
[ ]t 'tApp-nd' [Γ] f↓ x↓ = lift (SW (lower f↓ "a lower x↓))
[ ]t Γ ← Γ' [Γ] T↓ = T↓
[ ]t Γ → Γ' [Γ] T↓ = T↓
[ ]t ('fcomp-nd' {A} {B} {C}) [Γ] g↓ f↓ = lift ('o' {ε} (lower
  ('Γ' {B} {A} {b}) [Γ] = lift ('λ' {ε} ('VAR0' {ε} {__} {ε}
  ('Γ' {B} {A} {b}) [Γ] = lift ('λ' {ε} ('VAR0' {ε} {__} {ε}
[ ]t ('cast-refl' {T}) [Γ] = lift (cast-proof {T})
[ ]t ('cast-refl' {T}) [Γ] = lift (cast-proof {T})
[ ]t ('Γ' → {T} {B} {b} {c} {v}) [Γ] = lift ('idfun' {__} {ε}
[ ]t ('s ← {T} {B} {b} {c} {v}) [Γ] = lift ('idfun' {__} {ε}

module well-typed-syntax-interpreter where
open well-typed-syntax

```

open well-typed-syntax-eq-dec

max-level : Level

max-level = well-typed-syntax-pre-interpreter.max-level

$\llbracket _ \rrbracket^c : (\Gamma : \text{Context}) \rightarrow \text{Set} \text{ (lsuc max-level)}$
 $\llbracket _ \rrbracket^c = \text{well-typed-syntax-pre-interpreter.inner}.\llbracket _ \rrbracket^c$
 $(\lambda \ell P \Gamma \text{ dummy val} \rightarrow \text{context-pick-if } \{P = P\} \text{ dummy val})$
 $(\lambda \ell P \text{ dummy val} \rightarrow \text{context-pick-if-refl } \{P = P\} \{ \text{dummy} \})$

$\llbracket _ \rrbracket^T : \{\Gamma : \text{Context}\} \rightarrow \text{Type } \Gamma \rightarrow \llbracket _ \rrbracket^c \Gamma \rightarrow \text{Set max-level}$
 $\llbracket _ \rrbracket^T = \text{well-typed-syntax-pre-interpreter.inner}.\llbracket _ \rrbracket^T$
 $(\lambda \ell P \Gamma \text{ dummy val} \rightarrow \text{context-pick-if } \{P = P\} \text{ dummy val})$
 $(\lambda \ell P \text{ dummy val} \rightarrow \text{context-pick-if-refl } \{P = P\} \{ \text{dummy} \})$

$\llbracket _ \rrbracket^t : \forall \{\Gamma : \text{Context}\} \{T : \text{Type } \Gamma\} \rightarrow \text{Term } T \rightarrow (\llbracket \Gamma \rrbracket : \llbracket _ \rrbracket^c \Gamma) \rightarrow \llbracket _ \rrbracket^t T$
 $\llbracket _ \rrbracket^t = \text{well-typed-syntax-pre-interpreter.inner}.\llbracket _ \rrbracket^t$
 $(\lambda \ell P \Gamma \text{ dummy val} \rightarrow \text{context-pick-if } \{P = P\} \text{ dummy val})$
 $(\lambda \ell P \text{ dummy val} \rightarrow \text{context-pick-if-refl } \{P = P\} \{ \text{dummy} \})$

module well-typed-syntax-interpreter-full where

open well-typed-syntax

open well-typed-syntax-interpreter

$\text{Context}\varepsilon\Downarrow : \llbracket _ \rrbracket^c \varepsilon$
 $\text{Context}\varepsilon\Downarrow = \text{tt}$

$\text{Type}\varepsilon\Downarrow : \text{Type } \varepsilon \rightarrow \text{Set max-level}$
 $\text{Type}\varepsilon\Downarrow T = \llbracket _ \rrbracket^T T \text{ Context}\varepsilon\Downarrow$

$\text{Term}\varepsilon\Downarrow : \{T : \text{Type } \varepsilon\} \rightarrow \text{Term } T \rightarrow \text{Type}\varepsilon\Downarrow T$
 $\text{Term}\varepsilon\Downarrow t = \llbracket _ \rrbracket^t t \text{ Context}\varepsilon\Downarrow$

$\text{Type}\varepsilon\Downarrow\Downarrow : \forall \{A\} \rightarrow \text{Type } (\varepsilon \triangleright A) \rightarrow \text{Type}\varepsilon\Downarrow A \rightarrow \text{Set max-level}$
 $\text{Type}\varepsilon\Downarrow\Downarrow T A\Downarrow = \llbracket _ \rrbracket^T T (\text{Context}\varepsilon\Downarrow, A\Downarrow)$

$\text{Term}\varepsilon\Downarrow\Downarrow : \forall \{A\} \rightarrow \{T : \text{Type } (\varepsilon \triangleright A)\} \rightarrow \text{Term } T \rightarrow (x : \text{Type}\varepsilon\Downarrow A) \rightarrow \text{Type}\varepsilon\Downarrow\Downarrow T x$
 $\text{Term}\varepsilon\Downarrow\Downarrow t x = \llbracket _ \rrbracket^t t (\text{Context}\varepsilon\Downarrow, x)$

module löb where

open well-typed-syntax

open well-typed-quoted-syntax

open well-typed-syntax-interpreter-full

module inner ('X' : Type ε) ('f' : Term {ε ▷ ('□' "Γ 'X' ⊃^T)}) (W 'X') where
 X : Set _
 X = Typeε↓ 'X'

$f'' : (x : \text{Type}\varepsilon\Downarrow ('□' "Γ 'X' ⊃^T)) \rightarrow \text{Type}\varepsilon\Downarrow\Downarrow \{ '□' "Γ 'X' ⊃^T \} (W 'X') x$
 $f'' = \text{Term}\varepsilon\Downarrow\Downarrow f'$

dummy : Type ε
 dummy = 'Context'

$\text{cast} : (\Gamma v : \Sigma \text{Context Type}) \rightarrow \text{Type } (\varepsilon \triangleright \text{'Σ' 'Context' 'Type'})$
 $\text{cast } (\Gamma, v) = \text{context-pick-if } \{P = \text{Type}\} \{ \Gamma \} (W \text{ dummy}) v$

Hf : (h : Σ Context Type) → Type ε
 Hf h = (cast h " quote-sigma h '→' 'X')

qh : Term { (ε ▷ 'Σ' 'Context' 'Type') } (W ('Type' " 'ε'))

qh = f' w''' x

where

$f' : \text{Term } (W ('Type' "Γ ε ▷ \text{'Σ' 'Context' 'Type'} \text{ } ^\triangleright))$
 $f' = w \rightarrow \text{'cast' ' ' 'a' 'VAR}_0 \text{'}$

$x : \text{Term } (W ('Term' "Γ ε ^\triangleright \text{'Σ' 'Context' 'Type'} \text{ } ^\triangleright))$
 $x = (w \rightarrow \text{'quote-sigma' ' ' 'a' 'VAR}_0 \text{'})$

$h2 : \text{Type } (\varepsilon \triangleright \text{'Σ' 'Context' 'Type'})$
 $h2 = (W_1 \text{'□' " (qh w "→" w Γ 'X' ^\triangleright)})$

$h : \Sigma \text{Context Type}$
 $h = ((\varepsilon \triangleright \text{'Σ' 'Context' 'Type'}) , h2)$

H0 : Type ε

H0 = Hf h

H : Set

H = Term {ε} H0

$\text{'H0'} : \Box ('Type' "Γ ε ^\triangleright)$
 $\text{'H0'} = \Gamma H0 \text{ } ^\triangleright T$

$\text{'H'} : \text{Type } \varepsilon$
 $\text{'H'} = \text{'□' " 'H0'}$

$H0' : \text{Type } \varepsilon$
 $H0' = \text{'H' '→' 'X'}$

$H' : \text{Set}$
 $H' = \text{Term } \{ \varepsilon \} H0'$

$\text{'H0''} : \Box ('Type' "Γ ε ^\triangleright)$
 $\text{'H0''} = \Gamma H0' \text{ } ^\triangleright T$

$\text{'H''} : \text{Type } \varepsilon$
 $\text{'H''} = \text{'□' " 'H0''}$

toH-helper-helper : ∀ {k} → h2 ≡ k

→ □ (h2 " quote-sigma h '→' '□' "Γ h2 " quote-sigma h '→' 'X' ^\triangleright)
 → □ (k " quote-sigma h '→' '□' "Γ k " quote-sigma h '→' 'X' ^\triangleright)

toH-helper-helper p x = transport (λ k → □ (k " quote-sigma h '→' 'X' ^\triangleright)) x

toH-helper : □ (cast h " quote-sigma h '→' 'H')

toH-helper = toH-helper-helper

{k = context-pick-if {P = Type} {ε ▷ 'Σ' 'Context' 'Type'} (W du
 (sym (context-pick-if-refl {P = Type} {W dummy} {h2})))
 (SooW1 → (('→' → w "→" 'o' "fcomp-nd" 'a' ('s←←' 'o' "cas

$\text{'toH'} : \Box ('H' '→' 'H')$

$\text{'toH'} = \Gamma \rightarrow \Gamma \text{'o' "fcomp-nd" 'a' } (\Gamma \rightarrow \Gamma \text{'a' } \Gamma \text{toH-helper } ^\triangleright) \text{'o' } \Gamma \leftarrow \text{'a'}$

toH : H' → H

toH h' = toH-helper 'o' h'

fromH-helper-helper : ∀ {k} → h2 ≡ k

→ □ ('□' "Γ h2 " quote-sigma h '→' 'X' ^\triangleright → h2 " quote-sigm
 → □ ('□' "Γ k " quote-sigma h '→' 'X' ^\triangleright → k " quote-sigma h

fromH-helper-helper p x = transport (λ k → □ ('□' "Γ k " quote-sigma h '→' 'X' ^\triangleright)) x

fromH-helper : □ ('H' '→' cast h " quote-sigma h)

```

fromH-helper = fromH-helper-helper
{k = context-pick-if {P = Type} {ε ▷ 'Σ' 'Context' 'Type'}} (W dummy) h2}
(sym (context-pick-if-refl {P = Type} {W dummy} {h2}))
(S00W1' ← (Γ → 'Γ' 'o' 'fcomp-nd' 'a' (Γ → 'Γ' 'a' Γ 'λ' 'VAR0' 't' 'o' 'cast-refl' 'o' 's →') 'o' w' → → → →)))

'fromH' : □ ('H' → 'H')
'fromH' = Γ → 'Γ' 'o' 'fcomp-nd' 'a' (Γ → 'Γ' 'a' Γ fromH-helper 't') 'o' Γ ← 'Γ'

fromH : H → H'
fromH h' = fromH-helper 'o' h'

lob : □ 'X'
lob = fromH h' 'a' Γ h' 't'
where
  f' : Term {ε ▷ '□' 'H0'} (W ('□' ' (Γ '□' ' 'H0' 't' →) 'X' 't'))
  f' = Conv0 {'H0'} {'X'} (SW10 (w ∀ 'fromH' 'a' 'VAR0'))

  x : Term {ε ▷ '□' 'H0'} (W ('□' ' Γ 'H' 't'))
  x = w → 'quote-term' 'a' 'VAR0'

  h' : H
  h' = toH ('λ' (w → ('λ' 'f') 'a' (w → 'tApp-nd' 'a' f' 'a' x)))

lob : { 'X' : Type ε } → □ (('□' ' Γ 'X' 't') → 'X') → □ 'X'
lob { 'X' } 'f' = inner.lob 'X' (un'λ' 'f')

```

Acknowledgments

(Adam Chlipala, Matt Brown)
Acknowledgments, if needed.

References

- 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>.
- 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>.