

# A Machine-Checked Correctness Proof of Normalization by Evaluation for Simply Typed Lambda Calculus

Author: András Kovács

Advisor: Ambrus Kaposi

Budapest, 2017

## Abstract

We implement and prove correct normalization by evaluation for simply typed lambda calculus, using the proof assistant Agda. The correctness proof consists of soundness, completeness and stability with respect to  $\beta\eta$ -conversion. The algorithm uses a Kripke model over the preordered set of order-preserving context embeddings. Completeness is given by a logical relation between the term model and the Kripke model. For soundness and stability, we first need to prove that evaluation, quoting and unquoting all commute with order-preserving embeddings, which amounts to showing that the codomain of the evaluation function can be refined to a presheaf model. Then, soundness is shown by a logical relation on semantic values which fall into the refined model. Stability is proved by induction on normal and neutral terms. Decidability of  $\beta\eta$ -conversion for terms also follows from correctness. The main advantage of the formalization is that normalization has a concise and structurally recursive specification and is presented separately from the correctness proofs. Also, the syntax of terms remains simple, making no use of explicit substitutions or closures. Overall, the technical overhead of the development is relatively light.

## Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
1.1	Overview . . . . .	2
1.2	Related Work . . . . .	3
<b>2</b>	<b>Metatheory</b>	<b>4</b>
2.1	Type Theory . . . . .	4
2.2	Agda . . . . .	5
2.2.1	Basic Constructions . . . . .	5
2.2.2	Formalization . . . . .	7
2.2.3	Standard Library . . . . .	7

<b>3</b>	<b>Syntax</b>	<b>9</b>
3.1	Core Syntax . . . . .	10
3.2	Context Embeddings and Substitutions . . . . .	11
3.2.1	Embeddings . . . . .	11
3.2.2	Substitutions . . . . .	12
3.3	Conversion . . . . .	14
<b>4</b>	<b>Normalization</b>	<b>15</b>
4.1	Normal terms . . . . .	15
4.2	Preliminaries: Models of STLC . . . . .	16
4.3	Implementation . . . . .	18
4.4	Kripke Models . . . . .	21
<b>5</b>	<b>Correctness</b>	<b>23</b>
5.1	Categorical Notions . . . . .	23
5.2	Laws for Embedding and Substitution . . . . .	25
5.2.1	Embeddings . . . . .	25
5.2.2	Substitutions . . . . .	26
5.3	Completeness . . . . .	28
5.4	Presheaf Refinement . . . . .	31
5.5	Soundness . . . . .	34
5.6	Stability . . . . .	39
<b>6</b>	<b>Variations</b>	<b>40</b>
6.1	Direct Presheaf Model . . . . .	40
6.2	More Efficient Evaluation . . . . .	43
<b>7</b>	<b>Discussion and Future Work</b>	<b>45</b>

# 1 Introduction

## 1.1 Overview

Normalization animates the syntax of typed  $\lambda$ -calculi, making them useful as programming languages. Also, in contrast to general terms, normal forms possess a more restricted structure which simplifies formal reasoning. Moreover, normalization allows one to decide convertibility of terms, which is required during type checking polymorphic and dependent type theories.<sup>1</sup>

This article presents an efficient and verified implementation of  $\beta\eta$ -normalization for simply typed  $\lambda$ -calculus (STLC) in Agda, an implementation of constructive type theory. In this setting, the syntax can be seen as an embedded language, and manipulating embedded terms is a limited form of metaprogramming. Normalization can be reused

---

<sup>1</sup>In particular, the types of System  $F_\omega$  largely correspond to simply typed  $\lambda$ -terms.

in the implementation of proof tactics, decision procedures or domain-specific languages, therefore efficiency is desirable. It is also preferable that the syntax remains as simple and self-contained as possible, both for pedagogical purposes and for the ease of reuse in other Agda developments. For this reason we choose an intrinsic syntax with de Bruijn indices and implicit substitutions.

We are not concerned with mapping the algorithm to lower-level abstract machines or hardware. Thus we are free to make use of the most convenient evaluation mechanism at hand: that of the metalanguage. This way, we can gloss over implementation details of efficient higher-order evaluation. This is the core idea of normalization by evaluation (NbE) from an operational point of view. Also, totality of evaluation in the metatheory is always implicitly assumed from the inside. This lets us implement normalization in a structurally recursive way, making its totality trivial. NbE also naturally supports  $\eta$ -normalization.

If we are to trust a particular normalization algorithm, we need to prove the following properties (borrowing terminology from [3]):

- *Completeness*: normalization maps a term to a normal form which is convertible to the term.
- *Soundness*: normalization maps convertible terms to equal normal forms.

Additionally, we require stability, which establishes that normalization is aligned with the specification of normal forms:

- *Stability*: normalization acts as the identity function on normal terms.

Section 2 describes the metatheory we work in. In Section 3, we formalize the syntax of STLC along with  $\beta\eta$ -conversion. Section 4 presents the normalization function, and Section 5 proves its correctness. Section 6 describes two alternative formalizations, one with more efficient evaluation and one with more concise correctness proofs. We discuss the results and possible future research in the final section.

The Agda formalization can be found at <https://github.com/AndrasKovacs/stlc-nbe>.

## 1.2 Related Work

Martin-Löf [27] first used NbE (though the term had not been coined yet) to show decidability of type checking for the 1975 version of his intuitionistic type theory. Berger and Schwichtenberg reintroduced NbE in 1991 [10] as a normalization method for lambda calculi. Abel’s work [1] gives a comprehensive overview of NbE, and also develops it for a range of type theories. The approach differs markedly from ours, as it uses extrinsic syntax and separate typed and untyped phases of evaluation.

Catarina Coquand’s work [18] is the most closely related to our development. She formally proves soundness and completeness of NbE for simple type theory. However, she considers a nameful syntax with explicit substitutions while we use intrinsic de Bruijn variables with implicit substitutions.

Altenkirch and Chapman [3] formalize a big-step normalization algorithm in Agda. This development uses an explicit substitution calculus and implements normalization using an environment machine with first-order closures. The algorithm works similarly to ours on the common syntactic fragment, but it is not structurally recursive and hence specifies evaluation as a inductive relation, using the Bove-Capretta method [11].

Altenkirch and Kaposi [5] use a glued presheaf model (“presheaf logical predicate”) for NbE for a minimal dependent type theory. This development provided the initial inspiration for the formalization in this article; it seemed plausible that “scaling down” dependently typed NbE to simple type theory would yield a formalization that is more compact than existing ones. This turned out to be the case; however the discussion of resulting development is relegated to Section 6.1, because using a Kripke model has the advantage of clean separation of the actual algorithm and its naturality proofs, which was deemed preferable to the more involved presheaf construction. Also, our work uses a simple presheaf model rather than the glued model used in *Ibid.* or in the work of Altenkirch, Hoffman and Streicher [4].

## 2 Metatheory

In this section the metatheory used for formalization is presented, first broadly, then its specific implementation in Agda.

### 2.1 Type Theory

The basic system we use is intensional Martin-Löf type theory (MLTT). For an overview and tutorial see the first chapter of the Homotopy Type Theory book [35]. However, there is no canonical definition of intensional type theory. There are numerous variations concerning type universes and the range of allowed type definitions. Our development uses the following features:

- *Strictly positive inductive definitions.* These are required for an intrinsically well-typed syntax. We do not require induction-recursion, induction-induction, higher induction or large inductive types.
- *Two predicative universes, named  $\mathbf{Set}_0$  and  $\mathbf{Set}_1$ , with large elimination into  $\mathbf{Set}_0$ .*  $\mathbf{Set}_0$  is the base universe in Agda, i. e. the universe of small types. Large elimination is needed for recursive definitions of semantic types, since inductive definitions for them would not be positive and hence would be illegal. We only need  $\mathbf{Set}_1$  to type the large eliminations. In Agda  $\mathbf{Set}$  is a synonym for  $\mathbf{Set}_0$ .
- *Function extensionality as an axiom.* This posits that two functions are equal if they take equal arguments to equal results. We could eschew it by using setoid reasoning with semantic equivalence, as C. Coquand does [18]. However, we only use function extensionality in correctness proofs, and the normalization function

does not refer to it or to any other postulate. We are content with *logical* as opposed to computational content for correctness proofs. This way, the numerous congruence lemmas needed for setoid reasoning can be skipped. Also, function extensionality has computational interpretation in cubical [14] and observational [7] type theories, to which this development could be plausibly ported in the future (when practical implementations of the mentioned theories become available).

In Agda 2.5.2, there is robust support for dependent pattern matching without Streicher’s K axiom [33]. We use the `--without-K` language option for all of our code in the main development. However, we do use axiom K in Section 6.1 for the direct presheaf model, but only for technical convenience, as it can be avoided by additional uniqueness proofs for equalities.

## 2.2 Agda

Agda is a dependently typed programming language and proof assistant. Its design and core features are described in Ulf Norell’s thesis [30]. The latest documentation is available online [21]. For a book-length introduction geared towards beginners, see [34]. We summarize here the core constructions and our notational conventions and liberties.

### 2.2.1 Basic Constructions

Two essential artifacts in Agda code are *inductive type definitions* and *function definitions*.

*Inductive definitions* introduce new types. Agda has *open type universes*, which means that programmers may freely add new types as long as they preserve logical consistency<sup>2</sup>. An inductive definition involves a type constructor declaration followed by zero or more term constructor declarations. For example, the type of natural numbers is defined the following way:

```
data ℕ : Set where
  zero : ℕ
  suc   : ℕ → ℕ
```

Inductive definitions may have *parameters* and *indices*. The former are implicitly quantified over the term constructors, but must be uniform in the constructors’ return types. The latter must be explicitly quantified in term constructors, but are allowed to vary. The definition for length-indexed vectors exhibits both:

```
data Vec (A : Set) : ℕ → Set where
  nil   : Vec A zero
  cons  : {n : ℕ} → A → Vec A n → Vec A (suc n)
```

---

<sup>2</sup>Agda checks *strict positivity* as a sufficient condition for consistency.

In a type constructor declaration, parameters are listed left to the colon, while indices are to the right. `(A : Set)` is a parameter, so it is implicitly quantified and is the same in the return types of `nil` and `cons`. In contrast, the length index is quantified in `cons`. We can use brackets to make parameters implicit; here `cons` has an implicit first parameter, and can be used like `cons zero nil` for a value of `Vec N (suc zero)`. Implicit arguments are filled in by Agda’s unification algorithm. Alternatively, the `∀` symbol can be used to leave the types of parameters implicit, and we could define `cons` as follows:

```
cons : ∀ {n} → A → Vec A n → Vec A (suc n)
```

Additionally, Agda has records, providing syntactic sugar and namespace management for iterated  $\Sigma$  types. For example, semigroups can be defined as elements of the following record type:

```
record Semigroup : Set1 where
  field
    S      : Set
    _@_    : M → M → M
    assoc  : ∀ m1 m2 m3 → (m1 @ (m2 @ m3)) ≡ ((m1 @ m2) @ m3)
```

Elements of `Semigroup` can be given as `record {S = ...; _@_ = ...; assoc = ...}`. Whenever we have some `(sg : Semigroup)` in context, we can bring all of its fields into the current namespace by `open Semigroup sg`. Also, given an `(sg : Semigroup)` in scope, we can use field names as projections, as in `(S sg)` or `(assoc sg)`. We will sometimes omit record instances and simply write `S` or `assoc` (or whatever record fields we are working with) when it is clear from where we are projecting from.

*Function definitions* are given by pattern matching. In Agda, pattern matching implements branching evaluation as it is usual in functional languages, but it also refines type indices based on particular selected cases. The main practical difference between parameters and indices is that we can gain information about indices by pattern matching, while we can’t infer anything about parameters just by having a parameterized term.

For example, in the following definition of concatenation for `Vec`, the result type is refined depending on whether the first argument is empty:

```
_+_ : N → N → N
zero + m = m
suc n + m = suc (n + m)

_++_ : ∀ {A n m} → Vec A n → Vec A m → Vec A (n + m)
nil   ++ ys = ys
cons x xs ++ ys = cons x (xs ++ ys)
```

In the `cons x xs ++ ys` case, the result type is refined to `Vec A (suc (n + m))`, allowing us to build a `suc`-long result on the right hand side with `cons`.

Function definitions must be total, hence pattern matching must be exhaustive and recursive calls must be structurally decreasing. Dependent pattern matching with structural recursion allows us to write the same proofs as inductive eliminators, but with a more convenient interface and less administrative burden; see e. g. [23] for a translation from patterns to eliminators.

As a notational convention, we shall take cues from the Idris [12] programming language and sometimes leave implicit parameters implicit even in type declarations of functions or constructors. In this style, declaring `_++_` looks as follows:

```
_++_ : Vec A n → Vec A m → Vec A (n + m)
```

This is not valid Agda, but we shall do this whenever the types and binding status of parameters are obvious. Sometimes we also omit type annotations from bindings when they can be inferred or not relevant for exposition.

### 2.2.2 Formalization

The project can be found at <https://github.com/AndrasKovacs/stlc-nbe>. Most Agda code listings in this article are also included in the formal development, although the versions here may include syntactic liberties and abbreviations. We indicate the names of relevant source files in the examples as Agda comments or in expository text.

The formalization involves a considerable amount of details. We omit a significant part of them from our explanations; in particular, we often only present the types of definitions and omit the implementations. It may be useful to read the source code and this article at the same time, and also have an interactive Agda environment at hand.

### 2.2.3 Standard Library

The formal development does not have any external dependencies. We need less than a hundred lines of code from the Agda standard library [22], but we choose to include it alongside the project in `Lib.agda`, for portability, and also make some changes.

Most of `Lib.agda` concerns equality and equational reasoning. Propositional equality is the same as in standard Agda:

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

We also postulate function extensionality, for functions with both implicit and explicit arguments (the “i” in `fexti` stands for “implicit”):

```
postulate
  fext  : (∀ x → f x ≡ g x) → f ≡ g
  fexti : (∀ x → f {x} ≡ g {x}) → (λ {x} → f {x}) ≡ (λ {x} → g {x})
```

Note that the types of `f` and `g` are left implicit above, in accordance with our notational liberties. In full Agda we write:

```
fext :
  ∀ {i j}{A : Set i}{B : A → Set j}{f g : (x : A) → B x}
  → ((x : A) → f x ≡ g x) → f ≡ g
```

Here the quantification noise is increased by universe indices. Universe polymorphism is only used in `Lib.agda`, but we will henceforth omit universe indices anyway.

Our style of equational reasoning deviates from the standard library. First, we use a more compact notation for symmetry and transitivity, inspired by the Homotopy Type Theory book [35]:

```
_■_ : x ≡ y → y ≡ z → x ≡ z -- transitivity
refl ■ refl = refl

_⁻¹ : x ≡ y → y ≡ x -- symmetry
refl ⁻¹ = refl
```

`_⁻¹` is postfix, so if we have  $(p : x \equiv y)$ , then  $(p^{-1})$  has type  $(y \equiv x)$ .

Second, for reasoning about congruences we mainly use two operations. The first (named `ap` in the HoTT book and `cong` in the standard library) expresses that functions respect equality:

```
_&_ : (f : A → B) → x ≡ y → f x ≡ f y
f & refl = refl
infixl 9 _&_
```

The second one expresses that *non-dependent function application* respects equality as well:

```
_⊗_ : f ≡ g → x ≡ y → f x ≡ g y
refl ⊗ refl = refl
infixl 8 _⊗_
```

Using these two operators, we can lift non-dependent functions with arbitrary arities to congruences. For example, if we have  $(p : a_0 \equiv a_1)$ ,  $(q : b_0 \equiv b_1)$  and  $(f : A \rightarrow B \rightarrow C)$ , then  $(f \& p \otimes q)$  has type  $(f a_0 b_0 \equiv f a_1 b_1)$ . Note that both operators associate left, and `_&_` binds more strongly. This is reminiscent of the lifting syntax with applicative functors [29] in Haskell programming.

Furthermore, we use coercion instead of transport (see HoTT book [35, pp. 72], also named `subst` in the Agda standard library):

```
coe : A ≡ B → A → B
coe refl a = a
```

`subst` can be recovered by `_&_` and `coe`:



```
subst : (P : A → Set) → x ≡ y → P x → P y
subst P eq = coe (P & eq)
```

The choice for `coe` is mainly stylistic and rather subjective. `coe (P & eq)` does not look worse than `subst P eq`, but `coe` by itself is used often and is more compact than `subst id`.

In Agda, sometimes an explicit “equational reasoning” syntax is used. For example, commutativity for addition may look like the following:

```
+ -comm : (m n : ℕ) → m + n ≡ n + m
+ -comm zero    n = + -right-identity n -1
+ -comm (suc m) n = begin
  suc m + n      ≡{ }
  suc (m + n)    ≡{ suc & + -comm m n }
  suc (n + m)    ≡{ + -suc n m -1 }
  n + suc m      QED
```

In formal development we do not use this style and instead keep the intermediate expressions implicit:

```
+ -comm (suc m) n = (suc & + -comm n m) ■ (+ -suc n m -1)
```

This is because in an interactive environment one can step through the intermediate points fairly easily, and writing them out adds significant visual clutter. However, we will use informal equational reasoning when it aids explanation.

We also borrow  $\Sigma$  (dependent sum),  $\tau$  (unit type),  $\perp$  (empty type) and  $\_u\_$  (disjoint union) from the standard library. Elements of  $\Sigma$  are constructed as  $(a, b)$  pairs, and the projections are `proj1` and `proj2`.  $A \times B$  is defined in terms of  $\Sigma$  as the non-dependent pair type.  $\tau$  has `tt` (“trivially true”) as sole element. Injection into disjoint unions is done with `inj1` and `inj2`. The *ex falso* principle is named `⊥-elim`, with type  $\{A : \text{Set}\} \rightarrow A \rightarrow A$ .

### 3 Syntax

In this section we present the syntax of simply typed lambda calculus - in our case intrinsic syntax with implicit substitutions and de Bruijn indices. We take term conversion to be part of the syntax and therefore include it in this section.<sup>3</sup> Since the definition of conversion refers to substitution and weakening, they are also presented here.

---

<sup>3</sup>We take the view that the “proper” definition of the syntax would be an initial object in some suitable category of models, e. g. the initial category with families. In this way, models include conversion as propositional equations, so the syntax is quotiented by conversion as well; see e. g. [13] and [6]. We do not use this definition because quotients are not yet natively supported in Agda, and we find value in developing a conventional presentation first.

### 3.1 Core Syntax

The complete definition is the following:

```
-- Syntax.agda
data Ty : Set where
  ι      : Ty
  _⇒_    : Ty → Ty → Ty

data Con : Set where
  •      : Con
  _,_    : Con → Ty → Con

data _∈_ (A : Ty) : Con → Set where
  vz : A ∈ (Γ , A)
  vs : A ∈ Γ → A ∈ (Γ , B)

data Tm Γ : Ty → Set where
  var : A ∈ Γ → Tm Γ A
  lam : Tm (Γ , A) B → Tm Γ (A ⇒ B)
  app : Tm Γ (A ⇒ B) → Tm Γ A → Tm Γ B
```

There is a base type  $\iota$  (Greek iota) and function types. Contexts are just lists of types.  $\_ \in \_$  is an inductively defined membership relation on contexts. It has the same structure as de Bruijn indices; for this reason they are named **vz** for “zero” and **vs** for “suc”. Terms are parameterized by a  $\Gamma$  context and indexed by a type. **var** picks an entry from the contexts, **lam** is abstraction while **app** is function application.

As an example, the identity function on the base type is represented as:

```
id-ι : Tm • (ι ⇒ ι)
id-ι = lam (var vz) -- with nameful syntax: λ (x : ι) . x
```

The presented syntax is *intrinsic*. In other words, only the well-typed terms are defined. An *extrinsic* definition would first define untyped *preterms* and separately a typing relation on preterms and contexts:

```
data Tm : Set where -- omitted constructors
data _⊢_∈_ (Γ : Con) → Tm → Ty → Set where -- omitted constructors
```

Here,  $\Gamma \vdash t \in A$  expresses that the preterm  $t$  is well-formed and has type  $A$  in context  $\Gamma$ .

There are advantages and disadvantages to both intrinsic and extrinsic definitions. With intrinsic syntax, well-formedness and type preservation come for free. On the other hand, proofs and computations which do not depend on types are easier to formulate with preterms. In the case of STLC, the type system is simple enough so that carrying around type information never becomes burdensome, so in this setting intrinsic definitions are all-around more convenient. For polymorphic systems such as System

$F$ , intrinsic typing introduces significant bureaucratic overhead, by tagging terms with type coercions; see Benton et al. [9] for approaches to dealing with them. For extrinsic syntax, powerful automation for reasoning about substitution is available in Coq, allowing spectacularly compact proofs [32]. It remains to be seen if it can be practically adapted to intrinsic syntax. For dependent type theories, intrinsic typing with quotient inductive-inductive definitions seems to be most promising, as it requires a relatively compact set of rules [6], in contrast to extrinsic approaches which suffer from a veritable explosion of well-formedness conditions and congruence rules.

## 3.2 Context Embeddings and Substitutions

In order to specify  $\beta$ -conversion, we need a notion of substitution. Also, the  $\eta$ -rule must refer to a notion of *weakening* or *embedding*: it mentions the “same” term under and over a  $\lambda$  binder, but the two occurrences cannot be definitionally the same, since they are in different contexts and thus have different types. We need to express the notion that whenever one has a term in a context, one can construct essentially the same term in a larger context. We use *order-preserving-embeddings* for this. Embeddings enable us to state the  $\eta$ -rule, but we also make use of them for defining substitutions.

### 3.2.1 Embeddings

We define order-preserving embeddings the following way:

```
-- Embedding.agda
data OPE : Con → Con → Set where
  •      : OPE • •
  drop   : OPE  $\Gamma$   $\Delta$  → OPE ( $\Gamma$  , A)  $\Delta$ 
  keep   : OPE  $\Gamma$   $\Delta$  → OPE ( $\Gamma$  , A) ( $\Delta$  , A)
```

Elements of  $\text{OPE } \Gamma \Delta$  express that  $\Delta$  can be obtained from  $\Gamma$  by dropping zero or more entries in order.

Some choices can be made here. C. Coquand [19] uses an explicit identity embedding constructor with type  $(\forall \{\Gamma\} \rightarrow \text{OPE } \Gamma \Gamma)$  instead of our  $\bullet$ . This slightly simplifies some proofs and slightly complicates others. It does not have propositionally unique identity embeddings<sup>4</sup>, but overall the choice between this definition and ours is fairly arbitrary.

An alternative solution is using *renamings*. Renamings are essentially substitutions containing only variables, alternatively definable as functions mapping variables of smaller contexts to variables of larger contexts [28, pp. 24]:

```
Ren : Con → Con → Set
Ren  $\Gamma$   $\Delta$  = {A : Ty} → A ∈  $\Delta$  → A ∈  $\Gamma$ 
```

---

<sup>4</sup>Since wrapping the identity embedding with any number of `keep`-s also yields identity embeddings.

However, renamings are able to express *permutations* of contexts as well, which we have no need for in this development. Thus, renamings are larger than what is strictly necessary, which may or may not result in technical difficulties. We aim to choose the smallest workable representation, as a general principle.

Embeddings have action on variables and terms, reconstructing them in larger contexts. We denote embedding by lowercase “e” subscripts:

$$\begin{aligned}
& \epsilon_e : \text{OPE } \Gamma \Delta \rightarrow A \in \Delta \rightarrow A \in \Gamma \\
& \epsilon_e \bullet \quad \quad \quad v \quad \quad = v \\
& \epsilon_e (\text{drop } \sigma) v \quad \quad = vs (\epsilon_e \sigma v) \\
& \epsilon_e (\text{keep } \sigma) vz \quad \quad = vz \\
& \epsilon_e (\text{keep } \sigma) (vs v) = vs (\epsilon_e \sigma v) \\
\\
& \text{Tm}_e : \text{OPE } \Gamma \Delta \rightarrow \text{Tm } \Delta A \rightarrow \text{Tm } \Gamma A \\
& \text{Tm}_e \sigma (\text{var } v) \quad = \text{var } (\epsilon_e \sigma v) \\
& \text{Tm}_e \sigma (\text{lam } t) \quad = \text{lam } (\text{Tm}_e (\text{keep } \sigma) t) \\
& \text{Tm}_e \sigma (\text{app } f a) = \text{app } (\text{Tm}_e \sigma f) (\text{Tm}_e \sigma a)
\end{aligned}$$

Identity embeddings keep every entry:

$$\begin{aligned}
& \text{id}_e : \forall \{\Gamma\} \rightarrow \text{OPE } \Gamma \Gamma \\
& \text{id}_e \{\bullet\} \quad \quad = \bullet \\
& \text{id}_e \{\Gamma, A\} = \text{keep } (\text{id}_e \{\Gamma\})
\end{aligned}$$

We define `wk` as shorthand for the embedding that only drops the topmost entry:

$$\begin{aligned}
& \text{wk} : \forall \{A \Gamma\} \rightarrow \text{OPE } (\Gamma, A) \Gamma \\
& \text{wk} = \text{drop id}_e
\end{aligned}$$

### 3.2.2 Substitutions

A salient feature of our syntax is the lack of explicit substitution, i. e. we define substitutions separately as lists of terms, and their action on terms is given as a recursive function. We choose implicit substitutions for two reasons:

1. It is the canonical presentation of STLC in textbooks [31, 24].
2. Lighter formalization. With implicit substitutions, the definition of conversion becomes significantly smaller. With explicit substitutions, we need congruence closure of  $\beta$  and  $\eta$  conversion on terms and substitutions, plus rules for the action of substitution on terms. With implicit substitutions, only terms need to be considered. For illustration, Altenkirch and Chapman’s definition [3] of conversion for explicit substitution calculus has twenty-four rules, while this development has only seven. It also follows that with our syntax, equational reasoning about substitutions involves definitional or propositional equalities, which are automatically respected by all constructions. There is minor complication involving semantic interpretation of substitutions, discussed in Section 5.5.

Strictly speaking, stating  $\beta$ -conversion only requires single substitutions. However, simultaneous substitution is easier to define and reason about. Thus, we shall define the former using the latter and identity substitutions. We have substitutions as lists of terms:

```
-- Substitution.agda
data Sub (Γ : Con) : Con → Set where
  •       : Sub Γ
  _,_    : Sub Γ Δ → Tm Γ A → Sub Γ (Δ , A)
```

Elements of  $\text{Sub } \Gamma \Delta$  contain a  $\text{Tm } \Gamma A$  for each  $A$  type in  $\Delta$ . In other words, an element of  $\text{Sub } \Gamma \Delta$  assigns to each variable in  $\Delta$  a term in  $\Gamma$ .

Again, there are some choices in defining **Sub**. Identity substitutions (which have the identity action on terms) could be represented as an explicit constructor, as well as **Sub** composition and weakenings. Our **Sub** can be seen as a normal form of substitutions: explicit composition and weakening forms tree-like structures, but we flatten them to just lists of terms. This eliminates redundancy, and allows us to define compositions and identities recursively, which provides us with more definitional equalities.

We could also give a functional definition (see [28, pp. 24]), similarly to how we did for renamings before:

```
Sub : Con → Con → Set
Sub Γ Δ = {A : Ty} → A ∈ Δ → Tm Γ A
```

However, this contains many definitionally distinct functions with the same action on terms, and so we dismiss it on grounds of unnecessary largeness.

Before we can define identity substitutions and the action of substitution on terms, we need an operation with type  $(\forall \{A \in \Gamma\} \rightarrow \text{Sub } \Gamma \Delta \rightarrow \text{Sub } (\Gamma , A) (\Delta , A))$ . This enables pushing substitutions under  $\lambda$ -binders when recursing on terms. Constructing this substitution requires that we embed all terms in the input substitution into the extended  $(\Gamma , A)$  context. We define this as right composition with an embedding:

```
_s ° e_ : Sub Δ Σ → OPE Γ Δ → Sub Γ Σ
•       s ° e δ = •
(σ , t) s ° e δ = (σ s ° e δ) , Tm_e δ t
```

The notation  $\_s \circ e\_$  expresses that we have a substitution on the left and an embedding on the right. Its type resembles that of ordinary function composition, if we consider both **Sub** and **OPE** as context morphisms. We will expand on the categorical interpretation of this in Section 5.2.

There is a canonical injection  $\ulcorner\_ \urcorner^{\text{OPE}}$  from **OPE** to **Sub**:

```
drop_s : Sub Γ Δ → Sub (Γ , A) Δ
drop_s σ = σ s ° e wk

keep_s : Sub Γ Δ → Sub (Γ , A) (Δ , A)
```

$\text{keep}_s \sigma = \text{drop}_s \sigma , \text{var } vz$

```

 $\Gamma \vdash_{\text{ope}} : \text{OPE } \Gamma \Delta \rightarrow \text{Sub } \Gamma \Delta$ 
 $\Gamma \bullet \vdash_{\text{ope}} = \bullet$ 
 $\Gamma \text{ drop } \sigma \vdash_{\text{ope}} = \text{drop}_s \Gamma \sigma \vdash_{\text{ope}}$ 
 $\Gamma \text{ keep } \sigma \vdash_{\text{ope}} = \text{keep}_s \Gamma \sigma \vdash_{\text{ope}}$ 

```

We note that  $\text{keep}_s$  is the needed operation for pushing substitutions under binders. Now the action of substitution on terms and variables can be defined:

```

 $\epsilon_s : \text{Sub } \Gamma \Delta \rightarrow A \in \Delta \rightarrow \text{Tm } \Gamma A$ 
 $\epsilon_s (\sigma , t) \text{ vz} = t$ 
 $\epsilon_s (\sigma , t) (\text{vs } v) = \epsilon_s \sigma v$ 

 $\text{Tm}_s : \text{Sub } \Gamma \Delta \rightarrow \text{Tm } \Delta A \rightarrow \text{Tm } \Gamma A$ 
 $\text{Tm}_s \sigma (\text{var } v) = \epsilon_s \sigma v$ 
 $\text{Tm}_s \sigma (\text{lam } t) = \text{lam } (\text{Tm}_s (\text{keep}_s \sigma) t)$ 
 $\text{Tm}_s \sigma (\text{app } f a) = \text{app } (\text{Tm}_s \sigma f) (\text{Tm}_s \sigma a)$ 

```

$\epsilon_s$  simply looks up a term from a substitution<sup>5</sup>, while  $\text{Tm}_s$  recurses into terms to substitute all variables. Identity substitutions can be defined as well:

```

 $\text{id}_s : \{\Gamma : \text{Con}\} \rightarrow \text{Sub } \Gamma \Gamma$ 
 $\text{id}_s \{\bullet\} = \bullet$ 
 $\text{id}_s \{\Gamma , A\} = \text{keep}_s \text{id}_s$ 

```

Single substitution with a  $(t : \text{Tm } \Gamma A)$  term is given by  $(\text{id}_s , t)$ . This assigns the  $t$  term to the zeroth de Bruijn variable and leaves all other variables unchanged.

### 3.3 Conversion

The conversion relation is given as:

```

-- Conversion.agda
data _~_ {Γ} : ∀ {A} → Tm Γ A → Tm Γ A → Set where
  η      : t ~ lam (app (Tm_e wk t) (var vz))
  β      : app (lam t) t' ~ Tm_s (id_s , t') t

```

---

<sup>5</sup>The immediate reason for not naming  $\epsilon_s$  simply “lookup” is that we will need several different lookup functions for various purposes. For any sizable formal development, naming schemes eventually emerge, for better or worse. The author has found that uniformly encoding salient information about types of operations in their names is worthwhile to do. It is also not easy to hit the right level of abstraction; too little and we repeat ourselves too much or neglect to include known structures, but too much abstraction may cause the project to be less accessible and often also less convenient to develop in the first place. Many of the operations and proofs here could be defined explicitly using categorical language, i. e. bundling together functors’ actions on objects, morphisms and category laws in records, then projecting out required components. Our choice is to keep the general level of abstraction low for this development.

```

lam    : t ~ t' → lam t ~ lam t'
app    : f ~ f' → a ~ a' → app f a ~ app f' a'

~refl  : t ~ t
_~-1_ : t ~ t' → t' ~ t
_~■_   : t ~ t' → t' ~ t'' → t ~ t''

```

$\eta$  and  $\beta$  are the actual conversion rules. We push  $t$  under a lambda with `(Tme wk)` in  $\eta$ , and use single substitution for  $\beta$ . The rest are rules for congruence (`lam`, `app`) and equivalence closure (`~refl`, `_~-1_`, `_~■_`). The syntax for equivalence rules is the same as for propositional equality, except for the  $\sim$  prefixes.

## 4 Normalization

In this section we specify normal forms and implement normalization. Then, we discuss Kripke models and how normalization can be expressed in terms of them.

### 4.1 Normal terms

Our definition of normal forms is entirely standard: they are either lambdas or *neutral* terms of base type, and neutral terms are variables applied to zero or more normal arguments:

```

-- NormalForm.agda
mutual
  data Nf (Γ : Con) : Ty → Set where
    ne  : Ne Γ ι → Nf Γ ι
    lam : Nf (Γ , A) B → Nf Γ (A ⇒ B)

  data Ne (Γ : Con) : Ty → Set where
    var : A ∈ Γ → Ne Γ A
    app : Ne Γ (A ⇒ B) → Nf Γ A → Ne Γ B

```

Agda allows us to overload `var`, `lam` and `app` for normal forms.  $\beta$ -normality is obvious, since only variables can be applied. Normal forms are also  $\eta$ -long: since `ne` only injects neutrals of base type, all normal terms of function type must in fact be lambdas. Only requiring  $\beta$ -normality would be as simple as having `(ne : ∀ {A} → Ne Γ A → Nf Γ A)`. Alternatively, normal forms can be given as

```

data Nf (Γ : Con) : (A : Ty) → Tm Γ A → Set

```

with the same construction as before except that it is a predicate on general terms, expressing their normality. This definition would be roughly as convenient as the one we use. We also need actions of context embeddings:

$Nf_e : OPE \ \Gamma \ \Delta \rightarrow Nf \ \Delta \ A \rightarrow Nf \ \Gamma \ A$   
 $Ne_e : OPE \ \Gamma \ \Delta \rightarrow Ne \ \Delta \ A \rightarrow Ne \ \Gamma \ A$

Normal and neutral terms can be injected back to terms. Also, injection commutes with embedding (i. e. injection is natural).

$\ulcorner \_ \urcorner Nf : Nf \ \Gamma \ A \rightarrow Tm \ \Gamma \ A$   
 $\ulcorner \_ \urcorner Ne : Ne \ \Gamma \ A \rightarrow Tm \ \Gamma \ A$

$\ulcorner \_ \urcorner Ne\text{-}nat : \forall \ \sigma \ n \rightarrow \ulcorner Ne_e \ \sigma \ n \urcorner Ne \equiv Tm_e \ \sigma \ \ulcorner n \urcorner Ne$   
 $\ulcorner \_ \urcorner Nf\text{-}nat : \forall \ \sigma \ n \rightarrow \ulcorner Nf_e \ \sigma \ n \urcorner Nf \equiv Tm_e \ \sigma \ \ulcorner n \urcorner Nf$

## 4.2 Preliminaries: Models of STLC

Clearly, STLC is a small fragment of Agda, so we should be able to interpret the syntax back to Agda types and constructions in a straightforward way. From a semantic viewpoint, the most straightforward interpretation of the syntax is called the *standard model*. From an operational viewpoint, the standard model is just a well-typed interpreter [8] for STLC as an embedded language. It is implemented as follows:

```
-- Misc/StdModel.agda
Tys : Ty → Set
Tys ⊥      = ⊥
Tys (A ⇒ B) = Tys A → Tys B

Cons : Con → Set
Cons •      = ⊤
Cons (Γ , A) = Cons Γ × Tys A

∈s : ∀ {Γ A} → A ∈ Γ → (Cons Γ → Tys A)
∈s vz      Γs = proj2 Γs
∈s (vs v) Γs = ∈s v (proj1 Γs)

Tms : ∀ {Γ A} → Tm Γ A → (Cons Γ → Tys A)
Tms (var v)   Γs = ∈s v Γs
Tms (lam t)   Γs = λ as → Tms t (Γs , as)
Tms (app f a) Γs = Tms f Γs (Tms a Γs)
```

Traditionally, notation for semantics uses double brackets, e. g. the interpretation of types would look like this:

$\llbracket \_ \rrbracket : Ty \rightarrow Set$   
 $\llbracket \bot \rrbracket = \bot$   
 $\llbracket A \Rightarrow B \rrbracket = \llbracket A \rrbracket \rightarrow \llbracket B \rrbracket$

We instead opt for notating semantic interpretations with superscripts on type constructors (contrast *subscripts*, which we use for denoting operations for embeddings and



substitutions). Syntactic overloading in Agda is not convenient for our purposes here, so we have to distinguish different constructions (terms or types, etc.) and different models in any case. Therefore it makes sense to just reuse the names of types in the syntax, and distinguish different models by different superscripts. Similarly, we annotate names of function arguments with superscripts, if they are elements of certain semantic sets. For example, for some model  $\mathbf{M}$  we use  $\Gamma^{\mathbf{M}}$ ,  $\Delta^{\mathbf{M}}$  etc. for elements of semantic contexts, and  $\mathbf{a}^{\mathbf{M}}$ ,  $\mathbf{b}^{\mathbf{M}}$  etc. for semantic terms.

As to the implementation of the model: we interpret the base type with the empty type, and functions as functions. Contexts are interpreted as lists of semantic values. Terms are interpreted as functions from semantic contexts to semantic types. Now, normalizing  $(\text{Tm}^s \ (\text{lam} \ (\text{var} \ \text{vz})) \ \text{tt})$  yields the  $(\lambda \ \mathbf{a}^s \rightarrow \mathbf{a}^s)$  Agda term, so indeed we interpret a syntactic identity function with a metatheoretic identity function.

The standard model already has a key feature of NbE: it uses metatheoretic functions for evaluation. However, the standard model does not allow one to go back to syntax from semantics. If we have an  $(\mathbf{f} : \mathbf{A}^s \rightarrow \mathbf{B}^s)$  semantic function, all we can do with it is to apply it to an argument and extract a  $\mathbf{B}^s$ . *Weak evaluation* of programs assumes that all variables are mapped to semantic values. This corresponds to the usual notion of program evaluation in common programming languages. In contrast, strong normalization requires reducing terms inside function bodies, under binders. Bound variables are not mapped to any semantic value, so they block evaluation. In STLC, variables block function application; in richer systems variables may block case analysis or elimination of inductively defined data as well. See [1, pp. 12] for an overview for various approaches for implementing computation under binders.

NbE adds additional structure to semantic types, internalizing computation with blocking variables, which allows one to get back to syntax, moreover, to a subset of syntax containing only normal terms.

Adding progressively more structure to models allows us to prove more properties about the syntax. The standard model yields a simple proof of consistency, namely that there is no term with base type in the empty context:

```
consistency : Tm • ı → ⊥
consistency t = Tms t tt
```

*Kripke models* contain slightly more structure than the standard model, allowing us to implement normalization, which is a *completeness* theorem from a logical viewpoint [20], as per the Curry-Howard correspondence. Adding yet more structure yields *presheaf models*, which enable correctness proofs for normalization as well. We summarize the computational and logical interpretations below. Be mindful though that the table picks out specifics things relevant to our interest, and there are many more computations and logical interpretations possible for each class of models.

Table 1: Overview of models

Model	Computation	Proof
Standard	Type-safe interpreter	Consistency
Kripke	Normalizer	Completeness
Presheaf	Normalizer	Proof-relevant completeness

### 4.3 Implementation

We denote the model for normalization with  $^n$  superscripts. The implementation follows a similar shape as the standard model. The key difference - from which others follow - is in the interpretation of functions:

```
-- Normalization.agda
Tyn : Ty → Con → Set
Tyn ⊢      Γ = Nf Γ ⊢
Tyn (A ⇒ B) Γ = ∀ {Δ} → OPE Δ Γ → Tyn A Δ → Tyn B Δ
```

Each type is mapped to a  $(\text{Con} \rightarrow \text{Set})$  predicate. Hence, semantic types are not just sets anymore, but families of sets indexed by contexts. The base type is mapped to the family of its normal forms. Functions are mapped to semantic functions which take semantic inputs to semantic outputs in *any context larger than  $\Gamma$* . The additional  $\text{OPE}$  parameter is the key to the algorithm. As it was noted in Section 4.2, semantic functions must be able to take blocking variables as inputs. Blocking variables must be necessarily “fresh” with respect to the  $\Gamma$  context of a semantic function. Here, we can apply a semantic function to suitable embedding into a larger context, allowing to subsequently apply it to semantic terms containing variables pointing into that context. In the simplest case (and the only case we will actually use), if we have

$$f^n : \forall \{ \Delta \} \rightarrow \text{OPE } \Delta \Gamma \rightarrow \text{Ty}^n A \Delta \rightarrow \text{Ty}^n B \Delta$$

then we also have

$$f^n (\text{wk } \{A\}) : \text{Ty}^n A (\Gamma, A) \rightarrow \text{Ty}^n B (\Gamma, A)$$

The interpretations of contexts and variables are analogous to those in the standard model: contexts become lists of semantic values and variables look up semantic values from semantic contexts. However, we define semantic contexts inductively rather than recursively, for technical convenience.

```
data Conn : Con → Con → Set where
  •      : Conn • Δ
  _,_    : Conn Γ Δ → Tyn A Δ → Conn (Γ , A) Δ

∈n : ∀ {Γ A} → A ∈ Γ → ∀ {Δ} → Conn Γ Δ → Tyn A Δ
∈n v z      (Γn , tn) = tn
∈n (vs v) (Γn , _ ) = ∈n v Γn
```

When interpreting terms, variables and applications are easy. In the latter case we just apply the recursive result to  $\text{id}_e$  - we can choose *not* to conjure up fresh variables if there is no need to. We run into a bit of a bump in the interpretation of lambdas:

```

TmN : ∀ {Γ A} → Tm Γ A → ∀ {Δ} → ConN Γ Δ → TyN A Δ
TmN (var v)   ΓN = ∈N v ΓN
TmN (lam t)   ΓN = λ σ aN → TmN t (? , aN)  -- ? marks the bump
TmN (app f a) ΓN = TmN f ΓN ide (TmN a ΓN)

```

Clearly, the idea is to evaluate the inner  $t$  term in  $(\Gamma^N, a^N)$  as we did in the standard model, but  $\Gamma^N$  does not have the right type. It has type  $\text{Con}^N \Gamma \Delta$ , but we need to return in some  $\Sigma$  extended context, with  $(\sigma : \text{OPE } \Sigma \Delta)$  witnessing the extension. So, what we need is an action of context embedding on semantic contexts, and on semantic terms as well, since the former are just lists of the latter. To implement it, we first need *composition* for context embeddings:

```

_◦_ : OPE Δ Σ → OPE Γ Δ → OPE Γ Σ
σ ◦_ • = σ
σ ◦_ drop δ = drop (σ ◦_ δ)
drop σ ◦_ keep δ = drop (σ ◦_ δ)
keep σ ◦_ keep δ = keep (σ ◦_ δ)

```

It witnesses transitivity for the embedding relation. Then, embeddings on semantic terms is easy enough:

```

TyNe : OPE Δ Γ → TyN A Γ → TyN A Δ
TyNe {ι}   σ tN = Nfe σ tN
TyNe {A ⇒ B} σ tN = λ δ aN → tN (σ ◦_ δ) aN

ConNe : OPE Σ Δ → ConN Γ Δ → ConN Γ Σ
ConNe σ • = •
ConNe σ (ΓN , tN) = ConNe σ ΓN , TyNe σ tN

```

For  $\text{Ty}^{N_e} \{ \iota \}$ , we use embedding of normal forms. For  $\text{Ty}^{N_e} \{ A \Rightarrow B \}$ , we compose the externally given  $\sigma$  embedding with the input embedding  $\delta$ .  $\text{Con}^{N_e}$  is just pointwise application of  $\text{Ty}^{N_e}$ .  $\text{Tm}^N$  can be now given as:

```

TmN : ∀ {Γ A} → Tm Γ A → ∀ {Δ} → ConN Γ Δ → TyN A Δ
TmN (var v)   ΓN = ∈N v ΓN
TmN (lam t)   ΓN = λ σ aN → TmN t (ConNe σ ΓN , aN)
TmN (app f a) ΓN = TmN f ΓN ide (TmN a ΓN)

```

We still need a quote function to transform semantic terms to normal terms. We also need to mutually define an unquote function which sends neutral terms to semantic terms, for reasons shortly illuminated. Then, normalization is given as evaluation followed by quotation, where evaluation uses a semantic context given by unquoting each variable in the context.

### mutual

$$\begin{aligned}
q^N &: \forall \{A \ \Gamma\} \rightarrow Ty^N \ A \ \Gamma \rightarrow Nf \ \Gamma \ A \\
q^N \ \{\iota\} &\quad t^N = t^N \\
q^N \ \{A \Rightarrow B\} \ t^N &= \text{lam} \ (q^N \ (t^N \ \text{wk} \ (u^N \ (\text{var} \ \text{vz})))) \\
\\
u^N &: \forall \{A \ \Gamma\} \rightarrow Ne \ \Gamma \ A \rightarrow Ty^N \ A \ \Gamma \\
u^N \ \{\iota\} &\quad n = ne \ n \\
u^N \ \{A \Rightarrow B\} \ n &= \lambda \ \sigma \ a^N \rightarrow u^N \ (\text{app} \ (Ne_e \ \sigma \ n) \ (q^N \ a^N)) \\
\\
u^{cN} &: \forall \ \{\Gamma\} \rightarrow Con^N \ \Gamma \ \Gamma \\
u^{cN} \ \{\bullet\} &= \bullet \\
u^{cN} \ \{\Gamma, A\} &= Con^N_e \ \text{wk} \ u^{cN} \ , \ u^N \ (\text{var} \ \text{vz}) \\
\\
nf &: \forall \ \{\Gamma \ A\} \rightarrow Tm \ \Gamma \ A \rightarrow Nf \ \Gamma \ A \\
nf \ t &= q^N \ (Tm^N \ t \ u^{cN})
\end{aligned}$$

At first, it can be difficult to build a mental model of the operation of the algorithm. On the highest level, normalization alternates between evaluation and quoting: a term is first evaluated, then if it has a function type, its semantic value is applied to a “blocking” input and we quote the result, which can be a semantic function again. This evaluate-quote alternation proceeds until the result is not a function but a base value. Note that although  $q^N$  does not directly call  $Tm^N$ , it does apply semantic functions, and this application results in an *internal* call to  $Tm^N$ , accordingly to the definition of semantic application.

In fact, this encapsulation of recursive  $Tm^N$  calls is a crucial detail which makes this whole definition structurally recursive and thus total. In the  $(Tm^N \ (\text{lam} \ t) \ \Gamma^N = \lambda \ \sigma \ a^N \rightarrow Tm^N \ t \ (Con^N_e \ \sigma \ \Gamma^N, a^N))$  clause, the recursive  $Tm^N$  call is evidently structural, and because it happens under a *metatheoretic* lambda, semantic functions can be applied to *any* value in any definition without compromising structurality.

To explain the previous scare quotes around “blocking”: in the  $(t^N \ \text{wk} \ (u^N \ (\text{var} \ \text{vz})))$  application, the  $(u^N \ (\text{var} \ \text{vz}))$  term plays the role of blocking input, but it is not quite entirely blocking.  $u^N$  performs an operation best described as *semantic  $\eta$ -expansion*: it acts as identity on neutral base terms, and from neutral function it produces semantic function which build up neutral applications from their inputs. Thus,  $(u^N \ \{\iota\} \ (\text{var} \ \text{vz}))$  simply reduces to  $(ne \ (\text{var} \ \text{vz}))$ , while  $(u^N \ \{\iota \Rightarrow i\} \ (\text{var} \ \text{vz}))$  reduces to  $(\lambda \ \sigma \ a^N \rightarrow ne \ (\text{app} \ (\text{var} \ (\text{E}_e \ \sigma \ \text{vz})) \ a^N))$ . In short,  $u^N$  returns semantic values which yield properly  $\eta$ -expanded normal forms when quoted. As convoluted this may seem, it is actually a very elegant solution.

## 4.4 Kripke Models

We show now how Kripke models relate to the above definition of normalization. Following the presentation of Altenkirch [2], Kripke models can be defined in Agda for the syntax of STLC as follows (omitting universe polymorphism):

```
-- Misc/Kripke.agda
record KripkeModel : Set1 where
  field
    W      : Set
    _≤_     : W → W → Set
    ≤refl   : ∀ {w} → w ≤ w
    _≤■_    : ∀ {w1 w2 w3} → w1 ≤ w2 → w2 ≤ w3 → w1 ≤ w3
    _||-ι_  : W → Set
    ι-mono  : ∀ {w1 w2} → w1 ≤ w2 → w1 ||-ι → w2 ||-ι
```

A Kripke model consists of a preordered set of “worlds”, denoted  $W$ , along with a forcing predicate for the base type, and a *monotonicity condition*  $\iota\text{-mono}$ . On a concrete level, we can obtain the definition of `KripkeModel` by noticing the details in our previous `N` model which can be abstracted out, and packing them into a record. Indeed, the interpretation of function types, contexts and terms can be derived from this amount of data. Agda allow us to include them in the record as well, as definitions which depend on the record fields:

```
record KripkeModel : Set1 where
  field
    ... -- the same as before

    _||-Ty_ : W → Ty → Set
    w ||-Ty ι      = w ||-ι
    w ||-Ty (A ⇒ B) = ∀ {w'} → w ≤ w' → w' ||-Ty A → w' ||-Ty B

    _||-Con_ : W → Con → Set
    w ||-Con •      = τ
    w ||-Con (Γ , A) = (w ||-Con Γ) × (w ||-Ty A)

    _||-_ : Con → Ty → Set
    Γ ||- A = ∀ w → w ||-Con Γ → w ||-Ty A

    Ty-mono : w ≤ w' → w ||-Ty A → w' ||-Ty A
    Ty-mono {ι}      σ p = ι-mono σ p
    Ty-mono {A ⇒ B} σ p δ q = p (σ ≤■ δ) q

    Con-mono : w ≤ w' → w ||-Con Γ → w' ||-Con Γ
    Con-mono {•}      σ p = p
    Con-mono {Γ , A} σ p = Con-mono σ (proj1 p) , Ty-mono {A} σ (proj2 p)
```

```

||-Tm : Tm  $\Gamma$  A  $\rightarrow$   $\Gamma$  ||- A
||-Tm (var vz)      w ( $\_$  , q) = q
||-Tm (var (vs v)) w (p , q) = ||-Tm (var v) w p
||-Tm (lam t)       w p  $\sigma$  a = ||-Tm t  $\_$  (Con-mono  $\sigma$  p , a)
||-Tm (app f a)     w p       = ||-Tm f w p  $\leq$ refl (||-Tm a w p)

```

With this, everything up to  $Tm^N$  in our previous implementation can be obtained by defining the corresponding model:

```

N : KripkeModel
N = record {
  W      = Con
;  $\_ \leq \_$  =  $\lambda \Gamma \Delta \rightarrow OPE \Delta \Gamma$ 
;  $\leq$ refl = ide
;  $\_ \leq \blacksquare \_$  =  $\_ \circ_e \_$ 
;  $\_ ||-\iota$  =  $\lambda \Gamma \rightarrow Nf \Gamma \iota$ 
;  $\iota$ -mono = Nfe }

```

This is an interesting exercise in abstraction, but what is the deeper significance? STLC can be viewed as a simple intuitionistic propositional logic, with a single atomic proposition  $\iota$ . What is a right notion of semantics, though? We could abstract out the interpretation of base types from the standard model as given in Section 4.2, and get another notion of models; would not that suffice? From a logician's point of view, the issue is that it does not give us completeness. Semantics in general is necessary because we use logic to talk about concepts about which we have underlying intuition, but we have no *a priori* reason to believe that syntactic constructions talk meaningfully about any intuitive concept.

Soundness and completeness together express that syntactic and semantic entailment are logically equivalent; the lack of completeness indicates that there is a mismatch, that semantics is larger than necessary. In Agda, we define soundness and completeness for Kripke models as follows<sup>6</sup>:

```

sound      = Tm  $\Gamma$  A  $\rightarrow$  ( $\forall$  M  $\rightarrow$  let open KripkeModel M in  $\Gamma$  ||- A)
complete = ( $\forall$  M  $\rightarrow$  let open KripkeModel M in  $\Gamma$  ||- A)  $\rightarrow$  Tm  $\Gamma$  A

```

The **let open** KripkeModel M incantation makes locally available all fields and definitions of M. We give here a brief informal proof for soundness and completeness.

**Theorem.** *STLC is sound and complete with respect to Kripke models.*

*Proof.* Soundness follows immediately from  $||-Tm$ . For completeness, we instantiate the hypothesis with the previously defined N model, then define quoting and unquoting the same way as in Section 4.3. Using them, we can produce a  $(Nf \Gamma A)$ , which can be canonically injected back to  $(Tm \Gamma A)$ .  $\square$

---

<sup>6</sup>Note that these are wholly different from soundness and completeness as defined for a normalization algorithm in Section 1.1.

With the standard set-theoretic semantics for first-order logic, the intuitive content is clear: formulas talk about truth and falsehood. What is an intuitive meaning of Kripke semantics, though? A possible answer is that it talks about *knowledge*.  $\mathbf{W}$  worlds can be understood as states of knowledge, and  $\_ \leq \_$  denotes increasing knowledge. Monotonicity expresses that if something is proven in some state of knowledge, it will never be invalidated, not matter how “wiser” we get.

Note that in the formal development the direct  $\mathbf{W}$  definitions are used instead of Kripke models. This is primarily to make the algorithm as transparent as possible to lay readers. In the `Normalization.agda` file, the core algorithm is laid out in about 50 lines, and together with the definition of syntax and embedding it is about a hundred lines.

## 5 Correctness

In this section, we prove completeness, soundness and stability for the algorithm.

### 5.1 Categorical Notions

First, we shall introduce a modest number of concepts from category theory. The prime motivation is that they allow us to compactly describe the contents of our proofs. They also make it easier to mentally keep track of proof obligations, since often a few words of categorical definitions can be unpacked into a dozen lemmas. When proofs align with category theory, that is usually also indicative that one has the right approach. That said, we shall keep category theory to a minimum and only introduce definitions that are directly applicable to our work.

A *category* is defined as follows.

```
-- Misc/Category.agda
record Category : Set1 where
  field
    Obj      : Set                -- "objects"
    morph    : Obj → Obj → Set  -- "morphisms"
    id       : morph I I
    _◦_      : morph J K → morph I J → morph I K
    idl      : f ◦ id ≡ f
    idr      : id ◦ f ≡ f
    ass      : f ◦ (g ◦ h) ≡ (f ◦ g) ◦ h
```

In short, categories have a set of objects and sets of morphisms between objects, with identity morphisms, composition, and associativity and identity laws for composition. Note that this is a naive definition which does not take into account size issues and

subtleties arising from the ambient type theory; for a more rigorous treatment see [35, pp. 319].

In written exposition,  $obj : C$  expresses that  $obj$  is an object of  $C$ , while  $f : C(A, B)$  means that  $f$  is a morphism of  $C$  from  $A$  to  $B$ .

Most importantly, the syntax of STLC itself forms a category, which is sometimes called the *syntactic category*. We denote it as **STLC**. Its objects are contexts, its morphisms are substitutions with the identity substitution as identity morphism. We have not yet defined composition for substitution, which provides composition for **STLC**, nor the composition laws. Providing these will be a significant part of our proof obligations about substitutions.

Another important category is the category of sets, denoted **Set**. In our setting, **Set** has as objects Agda types and functions as morphisms (again, ignoring subtleties).

*Functors* are structure-preserving mappings between categories:

```
record Functor (C D : Category) : Set1 where
  field
    Obj⇒      : Obj C → Obj D
    morph⇒    : morph I J → morph (Obj⇒ I) (Obj⇒ J)
    id⇒       : morph⇒ id ≡ id
    ◦⇒        : morph⇒ (f ◦ g) ≡ morph⇒ f ◦ morph⇒ g
```

We use  $F : C \rightarrow D$  to express that  $F$  is a functor from  $C$  to  $D$ . Functors map objects to objects and morphisms to morphisms. Also, they must map identities to identities and compositions to compositions. We will use the terms *action on objects* and *action on morphisms* to denote the first two, and *functor laws* for the latter two.

A *contravariant functor* flips morphisms, i. e. it has

```
morph⇒ : morph I J → morph (Obj⇒ J) (Obj⇒ I)
```

The composition law changes accordingly. Contravariant functors are often denoted as  $F : C^{op} \rightarrow D$  (meaning that  $F$ 's domain is the *opposite* category of  $C$ , where all morphisms are flipped). Regular non-flipping functors are sometimes called *covariant functors*.

A *presheaf* is a contravariant functor from some  $C$  category to **Set**, i. e. it maps objects to sets and morphisms to functions, with morphism arrows being flipped in the process.

*Natural transformations* are mappings between  $F, G : C \rightarrow D$  functors. They consist of a  $\phi$  family of morphisms for each  $I$  object of  $C$ , and a *naturality* condition, which expresses that  $\phi$  commutes with any lifted  $C$ -morphism.

```
record Nat {C D : Category} (F G : Functor C D) : Set1 where
  field
    Φ      : {I : Obj C} → morph (Obj⇒ F I) (Obj⇒ G I)
    nat    : {I J : Obj C} {f : morph I J} → Φ ◦ morph⇒ F f ≡ morph⇒ G f ◦ Φ
```



$F : C \rightarrow D$  functors themselves form a category, where functors are objects and natural transformations are morphisms.

This development mostly utilizes presheaves and natural transformations between presheaves. More concretely, our setting will be mostly  $PSh(\mathbf{OPE})$ , the category of presheaves on  $\mathbf{OPE}$ , where  $\mathbf{OPE}$  is the category of contexts and order-preserving context embeddings. We examine  $\mathbf{OPE}$  in detail in the next section.

## 5.2 Laws for Embedding and Substitution

Before we attempt to prove correctness we should pin down the laws of embedding and substitution. People who set out to write normalization proofs soon find that a jumble of twenty-odd substitution lemmas is required to make headway. The naive proving process works by repeatedly hitting roadblocks and reacting by adding more lemmas. A more prudent way is to characterize all of them beforehand in categorical terms, which lends confidence that a given set of lemmas is complete.

### 5.2.1 Embeddings

Contexts and order-preserving embeddings form a category, which we call  $\mathbf{OPE}$ . Objects are contexts, morphisms are elements of the  $\mathbf{OPE}$  type. Identity and composition are as previously defined in Section 3.2.1 and Section 4.3:

```
ide  : ∀ {Γ} → OPE Γ Γ
_◦e_ : OPE Δ Σ → OPE Γ Δ → OPE Γ Σ
```

The category laws can be proven with simple induction.

```
idle : ide ◦e σ ≡ σ           -- left identity
idre : σ ◦e ide ≡ σ           -- right identity
asse : (σ ◦e δ) ◦e ν ≡ σ ◦e (δ ◦e ν) -- associativity
```

Additionally, variables and terms of a given type are presheaves on  $\mathbf{OPE}$ .  $(\lambda \Gamma \rightarrow A \in \Gamma)$  maps each object of  $\mathbf{OPE}$  to an object of  $\mathbf{Set}$ , that is, an Agda type. Likewise does  $(\lambda \Gamma \rightarrow Tm \Gamma A)$ . They also have action on morphisms: these are precisely the embedding operations on variables and terms, defined in Section 3.2.1.

```
εe  : OPE Γ Δ → (A ∈ Δ → A ∈ Γ)
Tme : OPE Γ Δ → (Tm Δ A → Tm Γ A)
```

Note that the input is a morphism in  $\mathbf{OPE}$ , i. e. a context embedding, and the output is a morphism in  $\mathbf{Set}$ , which is an Agda function. Also note the contravariance: the order of  $\Gamma$  and  $\Delta$  is flipped in the output. Functor laws for variables and terms are as follows:

```
-- Embedding.agda
ε-ide : ∀ v → εe ide v ≡ v
```

```

 $\epsilon \circ_e$  :  $\forall \sigma \delta v \rightarrow \epsilon_e (\sigma \circ_e \delta) v \equiv \epsilon_e \delta (\epsilon_e \sigma v)$ 
 $\text{Tm-id}_e$  :  $\forall t \rightarrow \text{Tm}_e \text{id}_e t \equiv t$ 
 $\text{Tm} \circ_e$  :  $\forall \sigma \delta t \rightarrow \text{Tm}_e (\sigma \circ_e \delta) t \equiv \text{Tm}_e \delta (\text{Tm}_e \sigma t)$ 

```

Here, the functor laws are given in a pointwise form. The categorically precise type of  $\epsilon\text{-id}_e$  would be the following:

```

 $\epsilon\text{-id}_e$  :  $\epsilon_e \text{id}_e \equiv (\lambda v \rightarrow v)$ 

```

This expresses that  $\text{id}_e$  is mapped to the identity function in **Set**. Instead of using this definition, we apply both sides of the equation to a variable, and similarly transform the types for the other laws. The reason for this is that proving the categorically precise forms requires function extensionality which we would like to avoid when possible.

Normal and neutral terms are also presheaves on **OPE** in a similar manner:

```

-- NormalForm.agda
 $\text{Nf}_e$  :  $\text{OPE } \Gamma \Delta \rightarrow \text{Nf } \Delta A \rightarrow \text{Nf } \Gamma A$ 
 $\text{Ne}_e$  :  $\text{OPE } \Gamma \Delta \rightarrow \text{Ne } \Delta A \rightarrow \text{Ne } \Gamma A$ 
 $\text{Nf} \circ_e$  :  $\forall \sigma \delta t \rightarrow \text{Nf}_e (\sigma \circ_e \delta) t \equiv \text{Nf}_e \delta (\text{Nf}_e \sigma t)$ 
 $\text{Ne} \circ_e$  :  $\forall \sigma \delta t \rightarrow \text{Ne}_e (\sigma \circ_e \delta) t \equiv \text{Ne}_e \delta (\text{Ne}_e \sigma t)$ 
 $\text{Nf-id}_e$  :  $\forall t \rightarrow \text{Nf}_e \text{id}_e t \equiv t$ 
 $\text{Ne-id}_e$  :  $\forall t \rightarrow \text{Ne}_e \text{id}_e t \equiv t$ 

```

This is to be expected, since neutral and normal terms are just terms with restricted structure.

### 5.2.2 Substitutions

We previously introduced **STLC** as the syntactic category, containing contexts as objects and substitutions as morphisms. However, we then presented **OPE** as yet another category with contexts as objects. Why single out **STLC** as the category which characterizes syntax? Well, the conversion relation is an integral part of the syntax, and **OPE** does not contain the morphisms needed to define  $\beta$ -conversion. In contrast, **STLC** contains all substitutions as well as all embeddings, since embeddings can be viewed as restricted substitutions. **OPE** is a *wide subcategory* [17] of **STLC**: it contains all the objects of **STLC** but only some of its morphisms.

We review identity substitution here (previously defined in Section 3.2.2). Composition is elementwise substitution of terms in a substitution.

```

 $\text{id}_s$  :  $\forall \{\Gamma\} \rightarrow \text{Sub } \Gamma \Gamma$ 
 $\text{id}_s \{\bullet\} = \bullet$ 
 $\text{id}_s \{\Gamma, A\} = \text{keep}_s \text{id}_s$ 

 $\_ \circ_s \_$  :  $\text{Sub } \Delta \Sigma \rightarrow \text{Sub } \Gamma \Delta \rightarrow \text{Sub } \Gamma \Sigma$ 
 $\bullet \circ_s \delta = \bullet$ 
 $(\sigma, t) \circ_s \delta = (\sigma \circ_s \delta), \text{Tm}_s \delta t$ 

```

We also have substitution operations for terms and variables:

$$\begin{aligned} \epsilon_s &: \text{Sub } \Gamma \Delta \rightarrow A \in \Delta \rightarrow \text{Tm } \Gamma A \\ \text{Tm}_s &: \text{Sub } \Gamma \Delta \rightarrow \text{Tm } \Delta A \rightarrow \text{Tm } \Gamma A \end{aligned}$$

We aim to establish that **STLC** is a category, and that  $(\text{Tm } \_ A)$  and  $(A \in \_)$  are presheaves on it. This is significantly more complicated than what we have seen for **OPE**. Identity and composition for substitutions are actually defined using embeddings, thus proving properties of substitution involves a proving a variety of lemmas about their interaction with embedding.

The categorical structure of proof obligations is less clear here; although we have clear goals (category and presheaf laws), this author is unaware of a compact and abstract explanation of the process of building **STLC** around **OPE**. In [9], substitution laws are constructed the same way as here, but the authors of Ibid. also do not provide a semantic explanation. For contrast, see [5], where substitutions are given first as part of the syntax, and renamings (an alternative of embeddings) are defined later as a subcategory.

There are four different operations of composing embeddings and substitutions, since there are two choices for both arguments. We define the missing  $\_e \circ_s \_$  operation here:

$$\begin{aligned} \_e \circ_e \_ &: \text{OPE } \Delta \Sigma \rightarrow \text{OPE } \Gamma \Delta \rightarrow \text{OPE } \Gamma \Sigma \text{ -- defined in 4.3} \\ \_e \circ_s \_ &: \text{Sub } \Delta \Sigma \rightarrow \text{Sub } \Gamma \Delta \rightarrow \text{Sub } \Gamma \Sigma \text{ -- defined in this section} \\ \_s \circ_e \_ &: \text{Sub } \Delta \Sigma \rightarrow \text{OPE } \Gamma \Delta \rightarrow \text{Sub } \Gamma \Sigma \text{ -- defined in 3.2.2} \\ \_e \circ_s \_ &: \text{OPE } \Delta \Sigma \rightarrow \text{Sub } \Gamma \Delta \rightarrow \text{Sub } \Gamma \Sigma \\ \bullet \quad \_e \circ_s \delta &= \delta \\ \text{drop } \sigma \_e \circ_s (\delta, t) &= \sigma \_e \circ_s \delta \\ \text{keep } \sigma \_e \circ_s (\delta, t) &= (\sigma \_e \circ_s \delta), t \end{aligned}$$

There are variations of the category laws where occurrences of composition can be any of the four versions, and it appears that some (not all) of these laws have to be proven. Similarly, functor laws for terms and variables have multiple versions differing in the mapped composition. All in all, we need to prove the following theorems in order:

$$\begin{aligned} &\text{-- Substitution.agda} \\ \text{ass}_{see} &: \forall \sigma \delta v \rightarrow (\sigma \_s \circ_e \delta) \_s \circ_e v \equiv \sigma \_s \circ_e (\delta \_e \circ_e v) \\ \text{ass}_{ese} &: \forall \sigma \delta v \rightarrow (\sigma \_e \circ_s \delta) \_s \circ_e v \equiv \sigma \_e \circ_s (\delta \_s \circ_e v) \\ \\ \text{idl}_{es} &: \forall \sigma \rightarrow \text{id}_e \_e \circ_s \sigma \equiv \sigma \\ \text{idl}_{se} &: \forall \sigma \rightarrow \text{id}_s \_s \circ_e \sigma \equiv \ulcorner \sigma \urcorner_{ope} \\ \text{idr}_{es} &: \forall \sigma \rightarrow \sigma \_e \circ_s \text{id}_s \equiv \ulcorner \sigma \urcorner_{ope} \\ \\ \epsilon_{-e \circ_s} &: \forall \sigma \delta v \rightarrow \epsilon_s (\sigma \_e \circ_s \delta) v \equiv \epsilon_s \delta (\epsilon_e \sigma v) \\ \text{Tm}_{-e \circ_s} &: \forall \sigma \delta t \rightarrow \text{Tm}_s (\sigma \_e \circ_s \delta) t \equiv \text{Tm}_s \delta (\text{Tm}_e \sigma t) \\ \\ \epsilon_{-s \circ_e} &: \forall \sigma \delta v \rightarrow \epsilon_s (\sigma \_s \circ_e \delta) v \equiv \text{Tm}_e \delta (\epsilon_s \sigma v) \end{aligned}$$

```

Tm- $\circ_s \circ_e$  :  $\forall \sigma \delta t \rightarrow Tm_s (\sigma \circ_s \circ_e \delta) t \equiv Tm_e \delta (Tm_s \sigma t)$ 

ass $_{s \circ_e s}$  :  $\forall \sigma \delta v \rightarrow (\sigma \circ_s \circ_e \delta) \circ_s v \equiv \sigma \circ_s (\delta \circ_e \circ_s v)$ 
ass $_{s \circ_s e}$  :  $\forall \sigma \delta v \rightarrow (\sigma \circ_s \delta) \circ_s \circ_e v \equiv \sigma \circ_s (\delta \circ_s \circ_e v)$ 

-- Functor laws for  $(A \in \_)$  and  $(Tm \_ A)$ 
 $\epsilon \circ_s$  :  $\forall \sigma \delta v \rightarrow \epsilon_s (\sigma \circ_s \delta) v \equiv Tm_s \delta (\epsilon_s \sigma v)$ 
Tm- $\circ_s$  :  $\forall \sigma \delta t \rightarrow Tm_s (\sigma \circ_s \delta) t \equiv Tm_s \delta (Tm_s \sigma t)$ 
 $\epsilon$ -id $_s$  :  $\forall v \rightarrow \epsilon_s id_s v \equiv var v$ 
Tm-id $_s$  :  $\forall t \rightarrow Tm_s id_s t \equiv t$ 

-- Category laws for STLC
idr $_s$  :  $\forall \sigma \rightarrow \sigma \circ_s id_s \equiv \sigma$ 
idl $_s$  :  $\forall \sigma \rightarrow id_s \circ_s \sigma \equiv \sigma$ 
ass $_s$  :  $\forall \sigma \delta v \rightarrow (\sigma \circ_s \delta) \circ_s v \equiv \sigma \circ_s (\delta \circ_s v)$ 

```

The proofs are moderately interesting and involve straightforward induction and equational reasoning. They comprise about 110 lines of Agda.

### 5.3 Completeness

In this section we prove completeness of normalization. It expresses that the output of normalization is convertible to its input.

```

complete :  $\forall \{\Gamma A\} (t : Tm \Gamma A) \rightarrow t \sim \ulcorner nf\ t \urcorner \lrcorner Nf$ 

```

$nf$  refers to the function defined in Section 4.3. In the return type, the injection  $\ulcorner \_ \urcorner \lrcorner Nf$  converts  $(Nf \Gamma A)$  back to  $(Tm \Gamma A)$ .

Since  $nf$  is defined with a Kripke model, it is clear that proving properties about it have to involve similar model structures. The overall structure of the proof is very similar to normalization itself; we will have interpretation of types, contexts, terms, and quoting and unquoting. In general, proofs about properties of functions must have the same inductive “shape” as the functions in question. The main difference here is that types are interpreted as *relations*:

```

-- Completeness.agda
 $\approx$  :  $\forall \{A \Gamma\} \rightarrow Tm \Gamma A \rightarrow Ty^N A \Gamma \rightarrow Set$ 
 $\approx$  {1} :  $t\ t^N = t \sim \ulcorner t^N \urcorner \lrcorner Nf$ 
 $\approx$  {A  $\Rightarrow$  B}{ $\Gamma$ } :  $t\ t^N =$ 
   $\forall \{\Delta\} (\sigma : OPE \Delta \Gamma) \{a\ a^N\} \rightarrow a \approx a^N \rightarrow app (Tm_e \sigma t) a \approx t^N \sigma a^N$ 

```

$\approx$  is a *Kripke logical relation*. Here, “logical relation” simply means that it is a relation defined by induction on types. “Kripke” indicates that the interpretation of function types is generalized to work in extended contexts. The purpose of the generalization is the same as with evaluation: it provides us with a “fresh variable” supply, allowing us to quote semantic functions and go under binders.

$\_ \approx \_$  essentially extends conversion to relate syntactic and semantic terms. At the base type, it is precisely conversion; at function types, it expresses that syntactic and semantic application respect  $\_ \approx \_$ .

We extend  $\_ \approx \_$  to a relation between substitutions and semantic contexts:

```
data  $\_ \approx^c \_$  { $\Gamma$ } :  $\forall$  { $\Delta$ }  $\rightarrow$  Sub  $\Gamma$   $\Delta \rightarrow$  ConN  $\Delta$   $\Gamma \rightarrow$  Set where
  •  $\_ \approx^c \_ : \bullet \approx^c \bullet$ 
   $\_ , \_ : \sigma \approx^c \delta \rightarrow t \approx t' \rightarrow (\sigma , t) \approx^c (\delta , t')$ 
```

Substitutions are list of syntactic terms while semantic context are lists of semantic terms, so they can be elementwise related by  $\_ \approx \_$ .

Mirroring <sup>N</sup>, we need monotonicity for semantic terms and contexts:

```
 $\approx_e : \forall \sigma \rightarrow t \approx t^N \rightarrow Tm_e \sigma t \approx Ty^N_e \sigma t^N$ 
 $\approx^c_e : \forall \sigma \rightarrow \delta \approx^c \nu^N \rightarrow (\delta \circ_e \sigma) \approx^c Con^N_e \sigma \nu^N$ 
```

The only interesting case here is the base type:

```
 $\approx_e \{1\} \sigma t \approx t^N = ?$ 
```

Here, the goal has type  $(Tm_e \sigma t \sim \ulcorner Nf_e \sigma t^N \urcorner Nf)$ , and we have  $(t \approx t^N : t \sim \ulcorner t^N \urcorner Nf)$ . First, we witness  $(Tm_e \sigma t \sim Tm_e \sigma \ulcorner t^N \urcorner Nf)$  by showing that embedding respects conversion:

```
-- Embedding.agda
 $\sim_e : \forall \sigma \rightarrow t \sim t' \rightarrow Tm_e \sigma t \sim Tm_e \sigma t'$ 
```

This can be proven by induction on  $(t \sim t')$ , with the help of substitution laws defined in Section 5.2. Now,  $(\sim_e \sigma t \approx t^N : Tm_e \sigma t \sim Tm_e \sigma \ulcorner t^N \urcorner Nf)$ , and the type can be shown to be equal to the goal using  $\ulcorner Nf$ -nat.

The only needed lemma which has no direct counterpart in <sup>N</sup> is the following:

```
 $\_ \sim \_ \approx \_ : \forall \{A \Gamma\} \{t t'\} \{t^N : Ty^N A \Gamma\} \rightarrow t \sim t' \rightarrow t' \approx t^N \rightarrow t \approx t^N$ 
 $\_ \sim \_ \approx \_ \{1\} \quad p \ q = p \sim \_ \ q$ 
 $\_ \sim \_ \approx \_ \{A \Rightarrow B\} \quad p \ q = \lambda \sigma \ a \approx a^N \rightarrow app (\sim_e \sigma p) \sim_{refl} \_ \approx \_ \ q \sigma \ a \approx a^N$ 
```

Now, the interpretation of terms can be defined. This is often called the “fundamental theorem” of a logical relation interpretation.

```
 $Tm \approx : \forall (t : Tm \Gamma A) \{ \sigma \} \{ \delta^N : Con^N \Gamma \Delta \} \rightarrow \sigma \approx^c \delta^N \rightarrow Tm_s \sigma t \approx Tm^N t \delta^N$ 
```

The proof is by induction on the input term. Witnesses for variables are simply looked up from  $(\sigma \approx^c \delta^N)$ , the same way as with other models. If the input term is an application, we use the inductive hypotheses the same way as in  $Tm^N$ , although first we need to rewrite the goal type with a functor law:

```
 $Tm \approx (app \ f \ a) \ \{ \sigma \} \ \sigma \approx \delta^N$ 
rewrite  $Tm$ -ide  $(Tm_s \sigma \ f) \ ^{-1} = Tm \approx \ f \ \sigma \approx \delta^N \ id_e \ (Tm \approx \ a \ \sigma \approx \delta^N)$ 
```

Only the case for  $(lam \ t)$  is non-trivial:

$\text{Tm} \approx (\text{lam } t) \{ \sigma \} \{ \delta \} \sigma \approx \delta^N \vee \{ a \} \{ a^N \} a \approx a^N = ?$

The goal is

$\text{app } (\text{lam } (\text{Tm}_e (\text{keep } v) (\text{Tm}_s (\text{keep}_s \sigma) t))) a \approx \text{Tm}^N t (\text{Con}^N_e \vee \delta^N, a^N)$

By using the induction hypotheses, we get the following:

$\text{Tm} \approx t (\approx^c_e \vee \sigma \approx \delta^N, a \approx a^N)$   
 $: \text{Tm}_s ((\sigma \circ_e v), a) t \approx \text{Tm}^N t (\text{Con}^N_e \vee \delta^N, a^N)$

The left hand side of the  $\approx$  does not yet match the goal. We need to use  $\sim \blacksquare \approx$  to compose it with a conversion proof on the left. The  $\beta$  rule can be used, but the result type does not line up on the right hand side:

$\beta (\text{Tm}_e (\text{keep } v) (\text{Tm}_s (\text{keep}_s \sigma) t)) a$   
 $: \text{app } (\text{lam } (\text{Tm}_e (\text{keep } v) (\text{Tm}_s (\text{keep}_s \sigma) t))) a \sim$   
 $\text{Tm}_s (\text{id}_s, a) (\text{Tm}_e (\text{keep } v) (\text{Tm}_s (\text{keep}_s \sigma) t))$

Hence we need to coerce this type to the desired one, using embedding and substitution laws. We omit the equality proof here.<sup>7</sup>

$\text{Tm} \approx (\text{lam } t) \{ \sigma \} \{ \delta^N \} \sigma \approx \delta^N \vee \{ a \} \{ a^N \} a \approx a^N =$   
 $\text{coe eq } (\beta (\text{Tm}_e (\text{keep } v) (\text{Tm}_s (\text{keep}_s \sigma) t)) a) \text{ -- eq omitted}$   
 $\sim \blacksquare \approx$   
 $\text{Tm} \approx t (\approx^c_e \vee \sigma \approx \delta^N, a \approx a^N)$

Now, quoting, unquoting and completeness can be defined as follows:

**mutual**

$q \approx : \forall \{A\} \rightarrow t \approx t^N \rightarrow t \sim \ulcorner q^N t^N \urcorner Nf$   
 $q \approx \{ \iota \} \quad t \approx t^N = t \approx t^N$   
 $q \approx \{ A \Rightarrow B \} t \approx t^N = \eta \_ \sim \blacksquare \text{lam } (q \approx (t \approx t^N \text{wk } (u \approx (\text{var } vz))))$

$u \approx : \forall \{A\} (n : Ne \ulcorner A \urcorner) \rightarrow \ulcorner n \urcorner Ne \approx u^N n$   
 $u \approx \{ \iota \} \quad n = \sim \text{refl}$   
 $u \approx \{ A \Rightarrow B \} n \sigma \{ a \} \{ a^N \} a \approx a^N$   
 $\text{rewrite } \ulcorner \urcorner Ne \text{-nat } \sigma n^{-1} =$   
 $\text{app } \sim \text{refl } (q \approx a \approx a^N) \sim \blacksquare \approx u \approx (\text{app } (Ne_e \sigma n) (q^N a^N))$

$u^{c \approx} : \forall \{ \Gamma \} \rightarrow \text{id}_s \{ \Gamma \} \approx^c u^{c N}$   
 $u^{c \approx} \{ \bullet \} = \bullet$   
 $u^{c \approx} \{ \Gamma, A \} = \approx^c_e \text{wk } u^{c \approx}, u \approx (\text{var } vz)$

$\text{complete} : \forall \{ \Gamma A \} (t : \text{Tm } \Gamma A) \rightarrow t \sim \ulcorner nf t \urcorner Nf$   
 $\text{complete } t = \text{coe eq } (q \approx (\text{Tm} \approx t u^{c \approx})) \text{ -- eq omitted}$

<sup>7</sup>Generally, explicit reasoning about substitutions is omitted from textbooks as well, because it tends to be technical and uninteresting.

Since, `nf` was defined as evaluation followed by quotation, we define `complete` this way as well. Expanding the type of `complete`, we get:

$$\text{complete} : \forall \{\Gamma A\} (t : \text{Tm } \Gamma A) \rightarrow t \sim \ulcorner q^N (\text{Tm}^N t \text{ u}^{c^N}) \urcorner \text{Nf}$$

However,  $(q \approx (\text{Tm} \approx t \text{ u}^{c \approx}))$  has type  $(\text{Tm}_s \text{ id}_s t \sim \ulcorner q^N (\text{Tm}^N t \text{ u}^{c^N}) \urcorner \text{Nf})$ , hence we need an extra coercion to rewrite the result type along  $\text{Tm-id}_s$ . This completes the proof.  $\square$

Now, since a large part of the proof was just reiterating the structure of  $^N$ , would not it be better to implement normalization and completeness in one go? For instance, the following could be implemented:

$$\text{nf} : \forall \{\Gamma A\} (t : \text{Tm } \Gamma A) \rightarrow \Sigma (\text{Nf } \Gamma A) \lambda n \rightarrow t \sim \ulcorner n \urcorner \text{Nf}$$

However, this has a disadvantage: the proof terms threaded through add considerable noise to any subsequent proof about normalization. It is also contrary to our goal of implementing the algorithm as simply as possible.

## 5.4 Presheaf Refinement

For soundness, one might be tempted to try directly defining another logical relation, and proceeding similarly as in the previous section. That approach would fall short. The reason is that at some point one has to prove naturality for evaluation, quoting and unquoting, i. e. that they commute with embeddings.

$$\begin{aligned} \text{Tm}^N\text{-nat} & : \forall t \sigma \Gamma^N \rightarrow \text{Tm}^N t (\text{Con}^N_e \sigma \Gamma^N) \equiv \text{Ty}^N_e \sigma (\text{Tm}^N t \Gamma^N) \\ q^N\text{-nat} & : \forall \sigma t^N \rightarrow \text{Nf}_e \sigma (q^N t^N) \equiv q^N (\text{Ty}^N_e \sigma t^N) \\ u^N\text{-nat} & : \forall \sigma n \rightarrow \text{Ty}^N_e \sigma (u^N n) \equiv u^N (\text{Ne}_e \sigma n) \end{aligned}$$

Unfortunately, these are not actually provable, because  $\text{Ty}^N$  is too large: it contains “unnatural” semantic terms. Recall its definition:

$$\begin{aligned} \text{Ty}^N & : \text{Ty} \rightarrow \text{Con} \rightarrow \text{Set} \\ \text{Ty}^N \iota & \quad \Gamma = \text{Nf } \Gamma \iota \\ \text{Ty}^N (A \Rightarrow B) \Gamma & = \forall \{\Delta\} \rightarrow \text{OPE } \Delta \Gamma \rightarrow \text{Ty}^N A \Delta \rightarrow \text{Ty}^N B \Delta \end{aligned}$$

$(\text{Ty}^N \iota \Gamma)$  is always fine, but  $(\text{Ty}^N (A \Rightarrow B) \Gamma)$  contains ill-behaved values, for instance a semantic function which returns different values depending on the length of the input context  $\Delta$ . Such a function clearly does not commute with embedding. We would like to restrict semantic functions to a well-behaved subset.

A *presheaf model* of STLC interprets types and contexts as presheaves. In particular, it interprets function types as *presheaf exponentials* [15] whose definition expands to a semantic function together with the naturality condition we are looking for. There is a presheaf model over **OPE** which can be used for normalization and which is a straightforward refinement of the Kripke model  $^N$ . Recall that a presheaf is a contravariant set-valued functor, hence the interpretation of types has to include four components:

action on objects, action on morphisms and the two functor laws. Below, for the interpretation of types the action on objects is defined and we give the types of the other components:

```

-- action on objects
TyP : Ty → Con → Set
TyP ι      Γ = Nf Γ ι
TyP (A ⇒ B) Γ =
  Σ (∀ {Δ} → OPE Δ Γ → TyP A Δ → TyP B Δ) λ fP →
  ∀ {Δ Σ} (σ : OPE Δ Γ) (δ : OPE Σ Δ) aP
    → fP (σ ∘e δ) (TyPe δ aP) ≡ TyPe δ (fP σ aP)

-- action on morphisms
TyPe      : OPE Δ Γ → TyP A Γ → TyP A Δ

-- functor laws
TyP-ide : ∀ tP → TyPe ide tP ≡ tP
TyP-∘e  : ∀ tP σ δ → TyPe (σ ∘e δ) tP ≡ TyPe δ (TyPe σ tP)

```

The key change is in the interpretation of function types, where we have a  $\Sigma$  of a function and a naturality condition.

Note that the action on morphisms was already present in the Kripke model under the name of “monotonicity condition”. In fact, every presheaf model is a Kripke model (but not vice versa). In Kripke models, types are also interpreted as functors, but from a *preordered set* rather than a general category.

Preordered sets can be viewed as categories with at most a single morphism between any two objects [16]. If  $I \leq J$ , then there is a morphism between  $I$  and  $J$ , otherwise there is no such morphism. Identity and composition of morphisms comes from reflexivity and transitivity of preordering.

In a Kripke model, functor laws hold by default, since the uniqueness of morphisms implies that identities necessarily map to identities and compositions to compositions. However, our **OPE** as it is defined in Agda is actually not a preorder. It can be informally considered as such if we only care about the mere existence of an embedding between two contexts. Formally, we could use *propositional truncation* [35, pp. 117] on the sets of morphisms of **OPE** to get an actual preorder, but that is beyond the capabilities of Agda and the scope of this article.

In the presheaf model, terms are interpreted as morphisms from semantic contexts to semantic types, in other words, natural transformations. Concretely, this means that the evaluation function<sup>8</sup> should be mutually defined with the proof that it commutes with embeddings, i. e. its naturality. Similarly, quoting and unquoting are defined mutually with their naturality proofs.

In this section, we do not use the <sup>P</sup> model - it is described in Section 6.1 in more

---

<sup>8</sup>Recall that the evaluation function is the interpretation of terms.



detail. Instead, the previous definition of normalization using  $^N$  is kept unchanged, but we shall need additional machinery to prove the necessary naturalities. The solution is the following: first, we define a logical predicate on semantic values which restricts them to natural inhabitants. Then, we give all naturality proofs, using restricted semantic values when required, and also prove that the  $\text{Nm}^N$  evaluation function outputs natural values.

The naturality predicate is defined below. We use superscript  $^P$  notation for definitions related to this predicate. Here,  $^P$  means “presheaf refinement”, not “presheaf model”.

```
-- PresheafRefinement.agda
TyP : ∀ {Γ A} → TyN A Γ → Set
TyP {Γ} {ι}      tN = τ
TyP {Γ} {A ⇒ B} tN =
  ∀ {Δ} (σ : OPE Δ Γ) {aN : TyN A Δ}
  → TyP aN
  → (∀ {Σ} (δ : OPE Σ Δ) → tN (σ ∘e δ) (TyNe δ aN) ≡ TyNe δ (tN σ aN))
  × TyP (tN σ aN)
```

At base type every semantic value is natural. A semantic function is natural if when applied to a natural argument it commutes with embedding and yields a natural result.

We extend naturality to semantic contexts and prove that it is preserved by embedding:

```
TyPe : ∀ σ → TyP tN → TyP (TyNe σ tN)

data ConP : ∀ {Γ Δ} → ConN Γ Δ → Set where
  •      : ConP •
  _,_    : ConP ΓN → TyP tN → ConP (ΓN , tN)

ConPe : ∀ σ → ConP ΓN → ConP (ConNe σ ΓN)
```

Next are functor laws for types and contexts. These are provable without reference to naturality. However, function extensionality is required.

```
TyN-ide : ∀ {A} (tN : TyN A Γ) → TyNe ide tN ≡ tN
TyN-ide {A = ι}      tN = Nf-ide tN
TyN-ide {A = A ⇒ B} tN = fexti λ Δ → fext λ δ → tN & idle δ

TyN-∘e : ∀ {A} (tN : TyN A Σ) σ δ → TyNe (σ ∘e δ) tN ≡ TyNe δ (TyNe σ tN)
TyN-∘e {A = ι}      tN σ δ = Nf-∘e σ δ tN
TyN-∘e {A = A ⇒ B} tN σ δ = fexti λ Ξ → fext λ ν → tN & asse σ δ ν

-- proofs for Con are just elementwise extensions as always
ConN-ide : ∀ σN → ConNe ide σN ≡ σN
ConN-∘e : ∀ σ δ ν → ConNe (σ ∘e δ) ν ≡ ConNe δ (ConNe σ ν)
```

We proceed to naturality for evaluation, quoting and unquoting, which are provable this time, thanks to the additional  $\mathsf{P}$  assumptions:

$$\begin{aligned}
\mathsf{Tm}^{\mathsf{P}} & : \forall t \rightarrow \mathsf{Con}^{\mathsf{P}} \Gamma^{\mathsf{N}} \rightarrow \mathsf{Ty}^{\mathsf{P}} (\mathsf{Tm}^{\mathsf{N}} t \Gamma^{\mathsf{N}}) \\
\mathsf{Tm}^{\mathsf{N}}\text{-nat} & : \forall t \sigma \rightarrow \mathsf{Con}^{\mathsf{P}} \Gamma^{\mathsf{N}} \rightarrow \mathsf{Tm}^{\mathsf{N}} t (\mathsf{Con}^{\mathsf{N}_e} \sigma \Gamma^{\mathsf{N}}) \equiv \mathsf{Ty}^{\mathsf{N}_e} \sigma (\mathsf{Tm}^{\mathsf{N}} t \Gamma^{\mathsf{N}}) \\
\mathsf{u}^{\mathsf{P}} & : \forall n \rightarrow \mathsf{Ty}^{\mathsf{P}} (\mathsf{u}^{\mathsf{N}} n) \\
\mathsf{q}^{\mathsf{N}}\text{-nat} & : \forall \sigma t^{\mathsf{N}} \rightarrow \mathsf{Ty}^{\mathsf{P}} t^{\mathsf{N}} \rightarrow \mathsf{Nf}_e \sigma (\mathsf{q}^{\mathsf{N}} t^{\mathsf{N}}) \equiv \mathsf{q}^{\mathsf{N}} (\mathsf{Ty}^{\mathsf{N}_e} \sigma t^{\mathsf{N}}) \\
\mathsf{u}^{\mathsf{N}}\text{-nat} & : \forall \sigma n \rightarrow \mathsf{Ty}^{\mathsf{N}_e} \sigma (\mathsf{u}^{\mathsf{N}} n) \equiv \mathsf{u}^{\mathsf{N}} (\mathsf{Ne}_e \sigma n) \\
\mathsf{u}^{c\mathsf{P}} & : \forall \{\Gamma\} \rightarrow \mathsf{Con}^{\mathsf{P}} (\mathsf{u}^{c\mathsf{N}} \{\Gamma\})
\end{aligned}$$

The proofs are obtained using mutual induction and substitution laws.

## 5.5 Soundness

Soundness expresses that normalization maps convertible terms to the same normal form:

$$\mathsf{sound} : t \sim t' \rightarrow \mathsf{nf} t \equiv \mathsf{nf} t'$$

Unfolding  $\mathsf{nf}$  in the return type yields  $(\mathsf{q}^{\mathsf{N}} (\mathsf{Tm}^{\mathsf{N}} t) \equiv \mathsf{q}^{\mathsf{N}} (\mathsf{Tm}^{\mathsf{N}} t'))$ . It is clear that  $(\mathsf{Tm}^{\mathsf{N}} t)$  should be in some relation to  $(\mathsf{Tm}^{\mathsf{N}} t')$  such that we can prove that  $\mathsf{q}^{\mathsf{N}}$  takes related inputs to equal outputs.

The simplest option is to define that relation as propositional equality. Since we assume function extensionality, propositional equality on semantic values is equivalent to the following logical relation:

$$\begin{aligned}
\_\approx\_ & : \forall \{A \Gamma\} \rightarrow \mathsf{Ty}^{\mathsf{N}} A \Gamma \rightarrow \mathsf{Ty}^{\mathsf{N}} A \Gamma \rightarrow \mathsf{Set} \\
\_\approx\_ \{\iota\} & \quad \{\Gamma\} t t' = t \equiv t' \\
\_\approx\_ \{A \Rightarrow B\} \{\Gamma\} t t' & = \forall \{\Delta\} (\sigma : \mathsf{OPE} \Delta \Gamma) \{a a'\} \rightarrow a \approx a' \rightarrow t \sigma a \approx t' \sigma a'
\end{aligned}$$

In other words, equality of semantic values is just propositional equality at base types and pointwise equality for semantic functions. However, we are again foiled by requirements of naturality:  $\_\approx\_$  relates unnatural values as well, but some needed equalities can only be shown on natural domains. Hence, we amend the above definition:

$$\begin{aligned}
& \text{-- Soundness.agda} \\
\_\approx\_ & : \forall \{A \Gamma\} \rightarrow \mathsf{Ty}^{\mathsf{N}} A \Gamma \rightarrow \mathsf{Ty}^{\mathsf{N}} A \Gamma \rightarrow \mathsf{Set} \\
\_\approx\_ \{\iota\} & \quad \{\Gamma\} t t' = t \equiv t' \\
\_\approx\_ \{A \Rightarrow B\} \{\Gamma\} t t' & = \\
& \quad \forall \{\Delta\} (\sigma : \mathsf{OPE} \Delta \Gamma) \{a a'\} \rightarrow \mathsf{Ty}^{\mathsf{P}} a \rightarrow \mathsf{Ty}^{\mathsf{P}} a' \rightarrow a \approx a' \rightarrow t \sigma a \approx t' \sigma a'
\end{aligned}$$

Alternatively,  $\_\approx\_$  and  $\mathsf{Ty}^{\mathsf{P}}$  could be merged into a single *partial equivalence relation*. In this definition, two functions are related if both are natural and take related inputs to related outputs. The naturality predicate can be recovered by setting  $(\mathsf{Ty}^{\mathsf{P}} t)$  to be  $(t \approx t)$ . This alternative was explored in the formalization but appeared to be less convenient than separate naturality predicates.

To prove soundness, first we extend  $\underline{\approx}$  to semantic contexts and prove monotonicity, similarly as in the completeness proof:

```

data  $\underline{\approx}^c$  :  $\forall \{ \Gamma \Delta \} \rightarrow \text{Con}^N \Gamma \Delta \rightarrow \text{Con}^N \Gamma \Delta \rightarrow \text{Set}$  where
  •  $\bullet : \bullet \approx^c \bullet$ 
   $\underline{\approx}, \underline{\approx} : \sigma \approx^c \delta \rightarrow t \approx t' \rightarrow (\sigma, t) \approx^c (\delta, t')$ 

 $\approx_e : \forall \sigma \rightarrow t \approx t' \rightarrow \text{Ty}^{N_e} \sigma t \approx \text{Ty}^{N_e} \sigma t'$ 
 $\approx^c_e : \forall \sigma \rightarrow \delta \approx^c v \rightarrow \text{Con}^{N_e} \sigma \delta \approx^c \text{Con}^{N_e} \sigma v$ 

```

Next, we prove that evaluating the same term in related semantic contexts yields related semantic values.

```

--  $\epsilon \approx$  just looks up a witness, as before.
 $\epsilon \approx : (v : A \in \Gamma) \rightarrow \sigma \approx^c \delta \rightarrow \epsilon^N v \sigma \approx \epsilon^N v \delta$ 

 $\text{Tm} \approx : (t : \text{Tm } \Gamma A) \rightarrow \text{Con}^P \Gamma^N \rightarrow \text{Con}^P \Gamma^{N'} \rightarrow \Gamma^N \approx^c \Gamma^{N'} \rightarrow \text{Tm}^N t \Gamma^N \approx \text{Tm}^N t \Gamma^{N'}$ 
 $\text{Tm} \approx (\text{var } v) \quad \Gamma^P \Gamma^{P'} \sigma \approx \delta = \epsilon \approx v \sigma \approx \delta$ 
 $\text{Tm} \approx (\text{lam } t) \quad \Gamma^P \Gamma^{P'} \sigma \approx \delta =$ 
   $\lambda v a^P a^{P'} a^N \approx a^{N'}$ 
   $\rightarrow \text{Tm} \approx t (\text{Con}^{P_e} v \Gamma^P, a^P) (\text{Con}^{P_e} v \Gamma^{P'}, a^{P'}) (\approx^c_e v \sigma \approx \delta, a \approx a')$ 
 $\text{Tm} \approx (\text{app } f a) \Gamma^P \Gamma^{P'} \sigma \approx \delta =$ 
   $\text{Tm} \approx f \Gamma^P \Gamma^{P'} \sigma \approx \delta \text{ id}_e (\text{Tm}^P a \Gamma^P) (\text{Tm}^P a \Gamma^{P'}) (\text{Tm} \approx a \Gamma^P \Gamma^{P'} \sigma \approx \delta)$ 

```

The proof is entirely mechanical, although it is complicated by the need to carry along  $\text{Con}^P$  proofs, in order to produce  $(\text{Ty}^P a \Gamma)$  naturality proofs for inputs in the  $(\text{app } f a)$  case.

Next, we prove the more general lemma expressing that convertible terms are evaluated to related values:

```

 $\sim \sim : t \sim t' \rightarrow \text{Con}^P \Gamma^N \rightarrow \text{Con}^P \Gamma^{N'} \rightarrow \Gamma^N \approx^c \Gamma^{N'} \rightarrow \text{Tm}^N t \Gamma^N \approx \text{Tm}^N t' \Gamma^{N'}$ 

```

This is the bulk of the work. We proceed with induction on  $(t \sim t')$ . First, let us consider the constructors for the equivalence closure, namely reflexivity, symmetry and transitivity. The latter two can be proven in general for  $\underline{\approx}$  and  $\underline{\approx}^c$ :

```

 $\underline{\approx}^{-1} : t \approx t' \rightarrow t' \approx t$ 
 $\underline{\approx}^{c-1} : \sigma \approx^c \delta \rightarrow \delta \approx^c \sigma$ 
 $\underline{\approx} \blacksquare : t \approx t' \rightarrow t' \approx t'' \rightarrow t \approx t''$ 
 $\underline{\approx}^c \blacksquare : \sigma \approx^c \delta \rightarrow \delta \approx^c v \rightarrow \sigma \approx^c v$ 

```

Reflexivity is not generally provable, but we only need it on the image of  $\text{Tm}^N$ , which is already covered by  $\text{Tm} \approx$ . We can implement the cases for equivalence closure in  $\sim \sim$ :

```

 $\sim \sim \{ t = t \} \sim \text{refl } \Gamma^P \Gamma^{P'} \Gamma^N \approx \Gamma^{N'} =$ 
   $\text{Tm} \approx t \Gamma^P \Gamma^{P'} \Gamma^N \approx \Gamma^{N'}$ 
 $\sim \sim (t' \sim t \sim^{-1}) \Gamma^P \Gamma^{P'} \Gamma^N \approx \Gamma^{N'} =$ 
   $\sim \sim t' \sim t \Gamma^{P'} \Gamma^P (\Gamma^N \approx \Gamma^{N'} \approx^{c-1}) \approx^{-1}$ 

```

$$\begin{aligned}
& \approx (t \sim t' \sim_{\blacksquare} t' \sim t'') \Gamma^P \Gamma^{P'} \Gamma^N \approx \Gamma^{N'} = \\
& \approx t \sim t' \Gamma^P \Gamma^{P'} (\Gamma^N \approx \Gamma^{N'} \approx_{\mathbf{c} \blacksquare} (\Gamma^N \approx \Gamma^{N'} \approx_{\mathbf{c}^{-1}})) \approx_{\blacksquare} \approx t' \sim t'' \Gamma^P \Gamma^{P'} \Gamma^N \approx \Gamma^{N'}
\end{aligned}$$

The **app** and **lam** congruences are also straightforward:

$$\begin{aligned}
& \approx (\text{lam } t \sim t') \Gamma^P \Gamma^{P'} \Gamma^N \approx \Gamma^{N'} = \\
& \quad \lambda v \ a^P \ a^{P'} \ a \approx a' \rightarrow \\
& \quad \approx t \sim t' (\text{Con}^P_e v \ \Gamma^P, a^P) (\text{Con}^P_e v \ \Gamma^{P'}, a^{P'}) (\approx_{\mathbf{c}_e} v \ \Gamma^N \approx \Gamma^{N'}, a \approx a') \\
& \approx (\text{app } \{f = f\} \{f' = f'\} \{a = a\} \{a' = a'\} t \sim t' a \sim a') \Gamma^P \Gamma^{P'} \Gamma^N \approx \Gamma^{N'} = \\
& \approx t \sim t' \Gamma^P \Gamma^{P'} \Gamma^N \approx \Gamma^{N'} \text{id}_e (\text{Tm}^P a \ \Gamma^P) (\text{Tm}^P a' \ \Gamma^{P'}) (\approx a \sim a' \ \Gamma^P \Gamma^{P'} \Gamma^N \approx \Gamma^{N'})
\end{aligned}$$

The real work happens with the  $\beta$  and  $\eta$  rules. Let us see  $\eta$ :

$$\approx (\eta \ t) \ \{\Gamma^N\} \ \{\Gamma^{N'}\} \ \Gamma^P \ \Gamma^{P'} \ \Gamma^N \approx \Gamma^{N'} \ v \ \{a\} \ \{a'\} \ a^P \ a^{P'} \ a \approx a' = ?$$

The goal has the following type:

$$\text{Tm}^N \ t \ \Gamma^N \ v \ a \approx \text{Tm}^N \ (\text{Tm}_e \ \text{wk } t) \ (\text{Con}^N_e \ v \ \Gamma^{N'}, a') \ \text{id}_e \ a'$$

Making full use of the induction hypothesis gives us the following:

$$\text{Tm} \approx t \ \Gamma^P \ \Gamma^{P'} \ \Gamma^N \approx \Gamma^{N'} \ v \ a^P \ a^{P'} \ a \approx a' : \text{Tm}^N \ t \ \Gamma^N \ v \ a \approx \text{Tm}^N \ t \ \Gamma^{N'} \ v \ a'$$

Hence, we need to rewrite (or coerce) the right hand side of the result type with an equality.

$$\text{eq} : \text{Tm}^N \ (\text{Tm}_e \ \text{wk } t) \ (\text{Con}^N_e \ v \ \Gamma^{N'}, a') \ \text{id}_e \ a' \equiv \text{Tm}^N \ t \ \Gamma^{N'} \ v \ a'$$

This should give us a moment of pause. Obviously, we need to prove something about the action of  $\text{Tm}^N$  on *embedded terms* such as  $(\text{Tm}_e \ \text{wk } t)$ . However, we haven't even defined an interpretation for embeddings in any model. Embedding and substitution is implicit in our syntax, therefore we are not obliged to interpret them. Note that in the presheaf model in Section 5.4, we interpreted *types* as presheaves, and for each type we defined its action on embeddings, with embeddings considered as **OPE** morphisms, but that is not the same as interpreting embeddings themselves, considered as **STLC** morphisms.

We interpret embeddings as natural transformations between semantic contexts:

$$\begin{aligned}
& \text{OPE}^N : \text{OPE } \Gamma \ \Delta \rightarrow \forall \ \{\Sigma\} \rightarrow \text{Con}^N \ \Gamma \ \Sigma \rightarrow \text{Con}^N \ \Delta \ \Sigma \\
& \text{OPE}^N \bullet \quad \Gamma^N \quad = \delta^N \\
& \text{OPE}^N \ (\text{drop } \sigma) \ (\Gamma^N, \_) = \text{OPE}^N \ \sigma \ \Gamma^N \\
& \text{OPE}^N \ (\text{keep } \sigma) \ (\Gamma^N, t^N) = \text{OPE}^N \ \sigma \ \Gamma^N, t^N \\
& \text{OPE}^N\text{-nat} : \forall \ \sigma \ \delta \ \Gamma^N \rightarrow \text{OPE}^N \ \sigma \ (\text{Con}^N_e \ \delta \ \Gamma^N) \equiv \text{Con}^N_e \ \delta \ (\text{OPE}^N \ \sigma \ \Gamma^N) \\
& \text{-- proof omitted}
\end{aligned}$$

While an embedding encodes that one context can be obtained from another by dropping zero or more elements, a semantic embedding actually performs the dropping of elements on semantic contexts. For instance,  $(\text{OPE}^N \ \text{wk})$  acts as a “tail” function, dropping the topmost element.

Additionally, since the identity embedding is defined recursively (as opposed to as a constructor), we have to prove that it is interpreted as the identity natural transformation:

$$\text{OPEN-id}_e : \forall \Gamma^N \rightarrow \text{OPEN id}_e \Gamma^N \equiv \Gamma^N$$

Now, we can give the action of evaluation on embedded terms:

$$\begin{aligned} \epsilon_e^N &: \forall \sigma \ v \ \Gamma^N \rightarrow \epsilon^N (\epsilon_e \sigma \ v) \quad \Gamma^N \equiv \epsilon^N \ v \ (\text{OPEN} \sigma \ \Gamma^N) \\ \text{Tm}_e^N &: \forall \sigma \ t \ \Gamma^N \rightarrow \text{Tm}^N (\text{Tm}_e \sigma \ t) \ \Gamma^N \equiv \text{Tm}^N t \ (\text{OPEN} \sigma \ \Gamma^N) \end{aligned}$$

Note that these are essentially *computation* rules, but given as propositional equalities. If the syntax had explicit substitutions (and, by extension, explicit embeddings), these rules would be actual definitional equalities in the evaluation function. It is a trade-off: we forgo these definitional equalities in favor of others stemming from the recursive definition of embedding and substitution, and in favor of the greatly reduced number of conversion rules.

We can turn our attention back to the equality proof required for the  $\eta$  case. By informal equational reasoning:

$$\begin{aligned} &\text{Tm}^N (\text{Tm}_e \text{ wk } t) (\text{Con}_e^N \ v \ \Gamma^{N'}, a') \text{id}_e a' \equiv \\ &\quad \text{-- by } (\text{Tm}_e^N \text{ wk } t (\text{Con}_e^N \ v \ \Gamma^{N'}, a')) \\ &\text{Tm}^N t (\text{OPEN id}_e (\text{Con}_e^N \ v \ \Gamma^{N'})) \text{id}_e a' \equiv \\ &\quad \text{-- by } (\text{OPEN-id}_e (\text{Con}_e^N \ v \ \Gamma^{N'})) \\ &\text{Tm}^N t (\text{Con}_e^N \ v \ \Gamma^{N'}) \text{id}_e a' \equiv \\ &\quad \text{-- by } (\text{Tm}^N\text{-nat } t \ v \ \{\Gamma^{N'}\} \ \Gamma^{P'}) \\ &\text{Tm}^N t \ \Gamma^{N'} (v \circ_e \text{id}_e) a' \equiv \\ &\quad \text{-- by } (\text{idr}_e \ v) \\ &\text{Tm}^N t \ \Gamma^{N'} \ v \ a' \end{aligned}$$

In Agda, we can rewrite the goal type with the above equality, then return the induction hypothesis:

$$\begin{aligned} &\sim (\eta \ t) \ \{\Gamma^N\} \ \{\Gamma^{N'}\} \ \Gamma^P \ \Gamma^{P'} \ \Gamma^N \approx \Gamma^{N'} \ v \ \{a\} \ \{a'\} \ a^P \ a^{P'} \ a \approx a' \\ &\quad \text{rewrite} \\ &\quad \text{Tm}_e^N \text{ wk } t \ (\text{Con}_e^N \ v \ \Gamma^{N'}, a') \\ &\quad | \text{OPEN-id}_e (\text{Con}_e^N \ v \ \Gamma^{N'}) \\ &\quad | \text{Tm}^N\text{-nat } t \ v \ \{\Gamma^{N'}\} \ \Gamma^{P'} \\ &\quad | \text{idr}_e \ v \\ &= \text{Tm}^N t \ \Gamma^P \ \Gamma^{P'} \ \Gamma^N \approx \Gamma^{N'} \ v \ a^P \ a^{P'} \ a \approx a' \end{aligned}$$

Now, consider the  $\beta$  case:

$$\sim (\beta \ t \ t') \ \{\Gamma^N\} \ \{\Gamma^{N'}\} \ \Gamma^P \ \Gamma^{P'} \ \sigma \approx \delta = ?$$

The goal is  $(\text{Tm}^N t (\text{Con}_e^N \text{id}_e \Gamma^N, \text{Tm}^N t' \Gamma^N) \approx \text{Tm}^N (\text{Tm}_s (\text{id}_s, t') t) \Gamma^{N'})$ . Note that the right hand side is the evaluation of a substituted expression. Similarly as with embeddings, an interpretation must be given for substitutions as well as the action of evaluation on substituted expressions.

Substitutions are also interpreted as natural transformations, however, naturality only holds on natural inputs:

$$\begin{aligned} \text{Sub}^N &: \text{Sub } \Gamma \Delta \rightarrow \forall \{\Sigma\} \rightarrow \text{Con}^N \Gamma \Sigma \rightarrow \text{Con}^N \Delta \Sigma \\ \text{Sub}^N \bullet & \quad \Gamma^N = \bullet \\ \text{Sub}^N (\sigma, t) \Gamma^N &= \text{Sub}^N \sigma \Gamma^N, \text{Tm}^N t \Gamma^N \end{aligned}$$

$$\text{Sub}^N\text{-nat} : \forall \sigma \delta \Gamma^N \rightarrow \text{Con}^P \Gamma^N \rightarrow \text{Sub}^N \sigma (\text{Con}^N_e \delta \Gamma^N) \equiv \text{Con}^N_e \delta (\text{Sub}^N \sigma \Gamma^N)$$

$(\text{Sub}^N \sigma \Gamma^N)$  evaluates each term in  $\sigma$  in  $\Gamma^N$ , thereby converting a list of terms (substitution) into a list of values (semantic context). We also need two related lemmas:

$$\begin{aligned} \text{Sub}^N\text{-}_s \circ_e &: \forall \sigma \delta \Gamma^N \rightarrow \text{Sub}^N (\sigma \circ_e \delta) \Gamma^N \equiv \text{Sub}^N \sigma (\text{OPE}^N \delta \Gamma^N) \\ \text{Sub}^N\text{-id}_s &: \forall \Gamma^N \rightarrow \text{Sub}^N \text{id}_s \Gamma^N \equiv \Gamma^N \end{aligned}$$

There is no deep reason why we need  $\text{Sub}^N\text{-}_s \circ_e$  but not the other variants; only this lemma is needed in the rest of the soundness proof.

The action of evaluation on substituted terms is a bit more involved. Naturality must be assumed for the input semantic context, or else the proof does not go through - in fact, this is the reason why naturality predicates were introduced. This implies that the action of evaluation cannot be stated as propositional equality, only as  $\approx$ , the previously defined equality of semantic values.

$$\text{Tm}_s^N : \forall \sigma t \Gamma^N \rightarrow \text{Con}^P \Gamma^N \rightarrow \Gamma^N \approx^c \Gamma^N \rightarrow \text{Tm}^N (\text{Tm}_s \sigma t) \Gamma^N \approx \text{Tm}^N t (\text{Sub}^N \sigma \Gamma^N)$$

Now,  $\text{Tm}_s^N$  can be used to prove a  $\approx$  where the right hand side matches the goal type for the  $\beta$  case.

$$\begin{aligned} &(\text{Tm}_s^N (\text{id}_s, t') t \Gamma^{N'} \Gamma^{P'} ((\sigma \approx \delta \approx^{c-1}) \approx^{c \blacksquare} \sigma \approx \delta) \approx^{-1}) \\ &: \text{Tm}^N t (\text{Sub}^N \text{id}_s \Gamma^{N'}, \text{Tm}^N t' \Gamma^{N'}) \approx \text{Tm}^N (\text{Tm}_s (\text{id}_s, t') t) \Gamma^{N'} \end{aligned}$$

We can use transitivity of  $\approx$  and erase semantic  $\text{id}_e$  actions, after which the goal becomes

$$\text{Tm}^N t (\Gamma^N, \text{Tm}^N t' \Gamma^N) \approx \text{Tm}^N t (\Gamma^{N'}, \text{Tm}^N t' \Gamma^{N'})$$

Note that the evaluation context differs on the sides, but the evaluated term is the same. Thus,  $\text{Tm}^N$  can be used to prove the goal, which completes the proof of  $\approx$ .

$$\begin{aligned} &(\text{Tm}^N t (\Gamma^P, \text{Tm}^P t' \Gamma^P) (\Gamma^{P'}, \text{Tm}^P t' \Gamma^{P'}) (\sigma \approx \delta, (\text{Tm}^N t' \Gamma^P \Gamma^{P'} \sigma \approx \delta))) \\ &: \text{Tm}^N t (\Gamma^N, \text{Tm}^N t' \Gamma^N) \approx \text{Tm}^N t (\Gamma^{N'}, \text{Tm}^N t' \Gamma^{N'}) \end{aligned}$$

After this, quoting, unquoting and soundness is easily given.

**mutual**

$$\begin{aligned} q\approx &: \forall \{A \Gamma\} \{t t' : \text{Ty}^N A \Gamma\} \rightarrow t \approx t' \rightarrow q^N t \equiv q^N t' \\ q\approx \{\text{id}\} & \quad t \approx t' = t \approx t' \\ q\approx \{A \Rightarrow B\} & \quad t \approx t' = \text{lam } \& \ q\approx (t \approx t' (\text{wk } \{A\})) (u^P \_) (u^P \_) (u\approx \text{refl}) \\ \\ u\approx &: \forall \{A \Gamma\} \{n n' : \text{Ne } \Gamma A\} \rightarrow n \equiv n' \rightarrow u^N n \approx u^N n' \\ u\approx \{\text{id}\} & \quad p = \text{ne } \& \ p \end{aligned}$$

```

u ≈ {A ⇒ B} p = λ σ aP aP' a ≈ a' → u ≈ (app & (Nee σ & p) ⊗ q ≈ a ≈ a')

uc ≈ : ∀ {Γ} → uc N {Γ} ≈c uc N
uc ≈ {•} = •
uc ≈ {Γ , A} = ≈ce wk uc ≈ , u ≈ refl

sound : ∀ {Γ A} {t t' : Tm Γ A} → t ~ t' → nf t ≡ nf t'
sound t ~ t' = q ≈ (≈ t ~ t' uc P uc P uc ≈)

□

```

## 5.6 Stability

Stability expresses that normalization has no action on normal terms:

```
stable : (n : Ne Γ A) → nf Γ n `Nf ≡ n
```

It differs in several ways from the other two correctness properties. First, it is not quite as crucial; failing either soundness or completeness results in something else than a normalization algorithm, but instability is tolerable. A primary goal of normalization is to have a decision procedure for conversion, however, it is achievable without stability:

```

-- decidable Nf equality, given by simple induction
Nf≡? : (n n' : Nf Γ A) → (n ≡ n') ∪ (n ≡ n' → ⊥)

decidableConversion : (t t' : Tm Γ A) → (t ~ t') ∪ (t ~ t' → ⊥)
decidableConversion t t' with Nf≡? (nf t) (nf t')
... | inj₁ p =
  inj₁ (complete t ~■ coe ((λ x → Γ x `Nf ~ t') & p-1) (complete t' ~-1))
... | inj₂ p =
  inj₂ (λ q → p (sound q))

```

Rather, stability tells us whether a normalization algorithm matches up exactly with a definition of normal forms. In our case, if we leave normalization essentially unchanged, but do not mandate  $\eta$ -long normal forms as we currently do, there would be  $\eta$ -short normal forms which are  $\eta$ -expanded by normalization:

```

t : Tm (ι ⇒ ι) (ι ⇒ ι)
t = var vz

nf-η-expands : nf t ≡ (lam (app (var (vs vz)) (var vz)))
nf-η-expands = refl

```

Here, normalization “thinks” that normal forms are more restricted than they actually are. To regain stability, normalization would have to be modified so that it does not do any  $\eta$ -conversion and leaves neutral terms unchanged. That said, an unstable algorithm is still fine - it just does “more” than what is required. One could argue

that instability should be avoided because unstable algorithms perform unnecessary computation.

Stability is much simpler to prove than previous correctness properties. Simple mutual induction on normal and neutral terms suffices, making use of previously proved naturalities and straightforward equational reasoning.

```

stableE : (v : A ∈ Γ) → EN v ucN ≡ uN (var v)
stableE vz      = refl
stableE (vs v) =
  EN-nat v wk ucN
  ■ TyNe wk & stableE v
  ■ uN-nat wk (var v)
  ■ (λ x → uN (var (vs x))) & E-ide v

mutual
stable : (n : Nf Γ A) → nf ⊢ n ⊢ Nf ≡ n
stable (ne n) = stableNe n
stable (lam n) = lam & stable n

stableNe : (n : Ne Γ A) → TmN ⊢ n ⊢ Ne ucN ≡ uN n
stableNe (var v) = stableE v
stableNe (app f a) =
  (λ x → x ide (TmN ⊢ a ⊢ Nf ucN)) & stableNe f
  ■ (λ x → uN (app (Nee ide f) x)) & stable a
  ■ (λ x → uN (app x a)) & Ne-ide f

```

□

## 6 Variations

In this section we look at two alternate variants of the formalization with different design choices. The first one is more compact but includes a less transparent implementation of normalization. The second one contains a more efficient evaluation function at the cost of slightly larger proofs.

### 6.1 Direct Presheaf Model

In Section 5.4 the presheaf model of STLC was described, but not pursued to its full extent. Going for a full presheaf model lets us dispense with about 100 lines of Agda, which is a significant portion of the roughly 700 lines of the whole development. The formalization can be found at <https://github.com/AndrasKovacs/stlc-nbe/tree/master>.

The <sup>N</sup> model is changed to a presheaf model, therefore it includes all mutually defined



naturality proofs. The interpretation of types is now the same as  $\text{Ty}^p$  in Section 5.4, although here we also use  $^n$  superscripts for the presheaf model<sup>9</sup>.

A benefit of this approach is that the  $\text{Ty}^p$  logical predicate for naturality can be entirely dropped, since naturality is already baked into  $\text{Ty}^n$ . The drawback is that all functor laws and naturalities must be proven even when defining the evaluation function. We list here the types (only the types, for brevity) of the definitions required to implement normalization, along with brief descriptions of categorical interpretations. Below,  $(f : \text{PSh}(\text{OPE})(A, B))$  means that  $f$  is a natural transformation in  $\text{PSh}(\text{OPE})$  between  $A$  and  $B$  presheaves.

```
-- Presheaf.agda

-- Ty^n A : PSh(OPE)
-----
Ty^n   : Ty → Con → Set                -- action on objects
Ty^ne  : OPE Δ Γ → Ty^n A Γ → Ty^n A Δ -- action on morphisms

Ty^n-id_e : ∀ t^n → Ty^ne id_e t^n ≡ t^n
Ty^n-∘_e  : ∀ t^n σ δ → Ty^ne (σ ∘_e δ) t^n ≡ Ty^ne δ (Ty^ne σ t^n)

-- Con^n Γ : PSh(OPE)
-----
Con^n   : Con → Con → Set                -- action on objects
Con^ne  : OPE Δ Σ → Con^n Γ Σ → Con^n Γ Δ -- action on morphisms

Con^n-id_e : ∀ Γ^n → Con^ne id_e Γ^n ≡ Γ^n
Con^n-∘_e  : ∀ σ δ Γ^n → Con^ne (σ ∘_e δ) Γ^n ≡ Con^ne δ (Con^ne σ Γ^n)

-- ∈^N {Γ}{A} v : PSh(OPE)(Con^n Γ, Ty^n A)
-----
∈^N      : ∀ {Γ A} → A ∈ Γ → ∀ {Δ} → Con^n Γ Δ → Ty^n A Δ
∈^N-nat  : ∀ (v : A ∈ Γ) σ Γ^n → ∈^N v (Con^ne σ Γ^n) ≡ Ty^ne σ (∈^N v Γ^n)

-- Tm^n {Γ}{A} t : PSh(OPE)(Con^n Γ, Ty^n A)
-----
Tm^n     : ∀ {Γ A} → Tm Γ A → ∀ {Δ} → Con^n Γ Δ → Ty^n A Δ
Tm^n-nat : ∀ t σ Γ^n → Tm^n t (Con^ne σ Γ^n) ≡ Ty^ne σ (Tm^n t Γ^n)

-- q^n {A} : PSh(OPE)(Ty^n A, Nf _ A)
-----
q^n      : ∀ {A Γ} → Ty^n A Γ → Nf Γ A
q^n-nat  : ∀ σ t^n → Nf_e σ (q^n t^n) ≡ q^n (Ty^ne σ t^n)
```

---

<sup>9</sup>It stands for “normalization”, which it achieves, but by a different model than previously.

```
-- uN {A} : PSh(OPE)(Ne _ A, TyN A)
-----
uN      : ∀ {A Γ} → Ne Γ A → TyN A Γ
uN-nat  : ∀ σ n → TyNe σ (uN n) ≡ uN (Nee σ n)
```

```
-----
ucN    : ∀ {Γ} → ConN Γ Γ
nf      : ∀ {Γ A} → Tm Γ A → Nf Γ A
```

The proofs for completeness and stability are largely unchanged. For soundness, we again need to interpret embeddings and substitutions and define the action of evaluation on them.

```
-- Presheaf.agda
```

```
-- OPEN {Γ}{Δ} σ : PSh(OPE)(ConN Γ, ConN Δ)
-----
OPEN      : ∀ {Γ Δ} → OPE Γ Δ → ∀ {Σ} → ConN Γ Σ → ConN Δ Σ
OPEN-nat  : ∀ σ δ ΓN → OPEN σ (ConNe δ ΓN) ≡ ConNe δ (OPEN σ ΓN)

OPEN-ide : ∀ ΓN → OPEN ide ΓN ≡ ΓN
TmeN     : ∀ σ t ΓN → TmN (Tme σ t) ΓN ≡ TmN t (OPEN σ ΓN)
```

```
-- SubN {Γ}{Δ} σ : PSh(OPE)(ConN Γ, ConN Δ)
-----
SubN      : ∀ {Γ Δ} → Sub Γ Δ → ∀ {Σ} → ConN Γ Σ → ConN Δ Σ
SubN-nat  : ∀ σ δ ΓN → SubN σ (ConNe δ ΓN) ≡ ConNe δ (SubN σ ΓN)

SubN-ids : ∀ ΓN → SubN ids ΓN ≡ ΓN
SubN-soe : ∀ σ δ ΓN → SubN (σ soe δ) ΓN ≡ SubN σ (OPEN δ ΓN)
TmsN     : ∀ σ t ΓN → TmN (Tms σ t) ΓN ≡ TmN t (SubN σ ΓN)
```

The soundness proof is dramatically simpler this time. The reason is that propositional equality is already right for semantic values, so there is no need for any additional `_≈_` logical relation, and propositional equality is already an equivalence relation which is respected by all constructions in the metatheory. Hence, soundness is simply given as follows:

```
-- Soundness.agda
≈~ : ∀ {t t' : Tm Γ A} → t ~ t' → (σ : ConN Γ Δ) → TmN t σ ≡ TmN t' σ

sound : ∀ {Γ A}{t t' : Tm Γ A} → t ~ t' → nf t ≡ nf t'
sound t~t' = qN & ≈~ t~t' ucN
```

Here, the proof for `≈~` still makes use of `OPEN` and `SubN`, but it is much simpler than before and altogether takes 11 lines of Agda.

There is a slight disadvantage of the direct presheaf model: because normalization

is defined mutually with naturalities, and some naturality proofs refer to function extensionality, it is not anymore trivial that normalization always computes. What if normalization gets stuck on the postulated extensionality? It is easy to prove that normalization does compute, but this proof can only be done *external to Agda*. Internally to Agda we obviously have no means to reason about postulates or distinguish them from defined proofs.

## 6.2 More Efficient Evaluation

We describe a more efficient evaluation function in this section. The correctness proofs are also formalized for this version, requiring only minor modifications. The formalization can be found at <https://github.com/AndrasKovacs/stlc-nbe/blob/efficient-appN>.

Let us review the definition of  $\text{Tm}^N$ :

$$\begin{aligned} \text{Tm}^N &: \forall \{\Gamma \ A\} \rightarrow \text{Tm} \ \Gamma \ A \rightarrow \forall \{\Delta\} \rightarrow \text{Con}^N \ \Gamma \ \Delta \rightarrow \text{Ty}^N \ A \ \Delta \\ \text{Tm}^N \ (\text{var } v) \quad \Gamma^N &= \epsilon^N \ v \ \Gamma^N \\ \text{Tm}^N \ (\text{lam } t) \quad \Gamma^N &= \lambda \ \sigma \ a^N \rightarrow \text{Tm}^N \ t \ (\text{Con}^N_e \ \sigma \ \Gamma^N, \ a^N) \\ \text{Tm}^N \ (\text{app } f \ a) \ \Gamma^N &= \text{Tm}^N \ f \ \Gamma^N \ \text{id}_e \ (\text{Tm}^N \ a \ \Gamma^N) \end{aligned}$$

Notice the  $\text{id}_e$  in the  $(\text{app } f \ a)$  case. After each application to  $\text{id}_e$ , the  $\Gamma^N$  context is changed to  $(\text{Con}^N_e \ \text{id}_e \ \Gamma^N)$  during the evaluation of the internal  $t$  body of  $(\text{Tm}^N \ f \ \Gamma^N)$ . This means that every application has an overhead proportional to the size of the current context. Also, recall  $\text{Ty}^N_e$  for functions:

$$\text{Ty}^N_e \ \{A \Rightarrow B\} \ \sigma \ t^N = \lambda \ \delta \ a^N \rightarrow t^N \ (\sigma \circ_e \delta) \ a^N$$

$(\text{Ty}^N_e \ \text{id}_e \ t^N)$  creates a new closure around  $t^N$  which has essentially no effect, and these closures accumulate on the top of semantic functions in  $\Gamma^N$  after each  $(\text{app } f \ a)$  evaluation. This is certainly bad for performance. We would like evaluation to be essentially as efficient for weak reduction as the standard model (as defined in Section 4.2). This can be achieved by defining  $\text{id}_e$  and  $\text{Con}^N_e$  in a way such that  $(\text{Con}^N_e \ \text{id}_e \ \Gamma^N)$  immediately reduces to  $\Gamma^N$ .

First,  $\text{id}_e$  needs to be a primitive constructor:

```
-- Embedding.agda
data OPE : Con → Con → Set where
  id_e   : ∀ {Γ} → OPE Γ Γ
  drop   : ∀ {A Γ Δ} → OPE Γ Δ → OPE (Γ , A) Δ
  keep   : ∀ {A Γ Δ} → OPE Γ Δ → OPE (Γ , A) (Δ , A)
```

This slightly complicates the identity functor law for  $(\text{Tm} \_ A)$  and  $(A \in \_)$ , since identity embeddings are not unique anymore. We define a predicate on embeddings which picks out the identities:

```

Ide : ∀ {Γ} → OPE Γ Γ → Set
Ide ide      = τ
Ide (drop σ) = ⊥
Ide (keep σ) = Ide σ

```

Id<sub>e</sub> embeddings are id<sub>e</sub>-s wrapped in zero or more keep-s. We prove identity functor laws using Id<sub>e</sub>:

```

E-ide : ∀ (v : A ∈ Γ) {σ : OPE Γ Γ} {_ : Ide σ} → Ee σ v ≡ v
Tm-ide : ∀ (t : Tm Γ A) {σ : OPE Γ Γ} {_ : Ide σ} → Tme σ t ≡ t

```

Secondly, Con<sup>N</sup><sub>e</sub> is redefined as follows:

```

ConNe : OPE Σ Δ → ConN Γ Δ → ConN Γ Σ
ConNe ide      ΓN      = ΓN
ConNe (drop σ) •      = •
ConNe (drop σ) (ΓN , tN) = (ConNe (drop σ) ΓN) , TyNe (drop σ) tN
ConNe (keep σ) •      = •
ConNe (keep σ) (ΓN , tN) = (ConNe (keep σ) ΓN) , (TyNe (keep σ) tN)

```

We get the desired definitional equality, but in turn lose two equalities which are quite useful to have. We can resurrect them as propositional equalities:

```

-- PresheafRefinement.agda
ConNe -• : ∀ σ → ConNe σ • ≡ •
ConNe -, : ∀ σ ΓN tN → ConNe σ (ΓN , tN) ≡ (ConNe σ ΓN , TyNe σ tN)

```

We must manually rewrite along these equalities when necessary. It is not a significant burden though. Altogether, the changes described here result in about 20 additional lines in the formalization.

We eliminated the overhead on the evaluation of application, but what about quoting? It also contains semantic application:

```

qN : ∀ {A Γ} → TyN A Γ → Nf Γ A
qN {ι}      tN = tN
qN {A ⇒ B} tN = lam (qN (tN wk (uN (var vz))))

```

Here, the (t<sup>N</sup> wk) is essential, but it is far less of an overhead than the id<sub>e</sub>-s before. Notice that q<sup>N</sup> is *type-directed*, and for each semantic value it is called as many times as there are function arguments in the value's type. Of course, q<sup>N</sup> is applied to various sub-terms of a term during normalization (via u<sup>N</sup>), so the overall number of q<sup>N</sup> calls is likely greater than the number of arguments in a term's type. Still, it is only to be expected that *normalization* is more costly than standard interpretation, since it does more work. *Evaluation* on its own is no less efficient than standard interpretation, with the current optimized implementation.

## 7 Discussion and Future Work

The primary goal of this development was to formalize  $\beta\eta$ -normalization for STLC as simply and compactly as possible, with the sub-goal of implementing normalization itself in a transparent and efficient manner. Overall, these goals are achieved in the main formalization and its two alternate versions as well, with slightly varying trade-offs with respect to the desiderata. In particular, our development has the following advantages:

- Our syntax is conventional and matches the common presentation in pedagogical literature. Also, the used metatheory is conservative.
- The formalization is complete: no postulates (aside from function extensionality) or holes remain.
- The formalization is compact: not counting the standard library, the main version is roughly 700 lines, while the direct presheaf version is 600 lines.
- The normalization algorithm is structurally recursive, obviating the need for termination proofs.

As to possible extensions and improvements to this work, while it seems that the current approach is uniquely well-suited to STLC and its non-polymorphic extensions, there are complications and rough edges when we attempt to scale to more powerful type theories.

For System  $F$  with intrinsic syntax and implicit substitutions, this author has an unpublished mostly-complete formalization of an NbE algorithm [26], although there are no correctness proofs yet. It seems that System  $F_\omega$  is about as far as we can take intrinsic syntax with implicit substitutions. With full dependent types, intrinsic syntax only seems feasible with higher-inductive definitions and explicit substitutions. Implicit substitutions remain only feasible with extrinsic syntax.

Another possible improvement would be proof automation for substitution and embedding laws. A significant part of the development is essentially “boilerplate” related to this. For extrinsic syntax, there are very powerful automation libraries available in the proof assistant Coq [25, 32]. For intrinsic syntax, automatic generation of substitution lemmas seems more difficult, but even if we prove categorical lemmas by hand, we could still make use of automation for proving ad hoc term equations. This seems doable in Agda, although the automation facilities of Agda are not as nearly mature as that of Coq.

## References

- [1] Andreas Abel. “Normalization by Evaluation: Dependent Types and Impredicativity”. Habilitation thesis. PhD thesis. Ludwig-Maximilians-Universität München, 2013.

- [2] Thorsten Altenkirch. “Normalisation by completeness”. In: *Talk given at the Workshop on Normalization by Evaluation in Los Angeles*. 2009.
- [3] Thorsten Altenkirch and James Chapman. “Big-step normalisation”. In: *Journal of Functional Programming* 19.3-4 (2009), pp. 311–333.
- [4] Thorsten Altenkirch, Martin Hofmann, and Thomas Streicher. “Categorical reconstruction of a reduction free normalization proof”. In: *Category Theory and Computer Science*. Springer. 1995, pp. 182–199.
- [5] Thorsten Altenkirch and Ambrus Kaposi. “Normalisation by evaluation for dependent types”. In: *LIPICs-Leibniz International Proceedings in Informatics*. Vol. 52. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik. 2016.
- [6] Thorsten Altenkirch and Ambrus Kaposi. “Type theory in type theory using quotient inductive types”. In: *ACM SIGPLAN Notices*. Vol. 51. 1. ACM. 2016, pp. 18–29.
- [7] Thorsten Altenkirch, Conor McBride, and Wouter Swierstra. “Observational equality, now!” In: *Proceedings of the 2007 workshop on Programming languages meets program verification*. ACM. 2007, pp. 57–68.
- [8] Lennart Augustsson and Magnus Carlsson. *An exercise in dependent types: A well-typed interpreter*. 1999.
- [9] Nick Benton et al. “Strongly typed term representations in Coq”. In: *Journal of automated reasoning* 49.2 (2012), pp. 141–159.
- [10] Ulrich Berger and Helmut Schwichtenberg. “An inverse of the evaluation functional for typed lambda-calculus”. In: *Logic in Computer Science, 1991. LICS’91., Proceedings of Sixth Annual IEEE Symposium on*. IEEE. 1991, pp. 203–211.
- [11] Ana Bove and Venanzio Capretta. “Modelling general recursion in type theory”. In: *Mathematical Structures in Computer Science* 15.04 (2005), pp. 671–708.
- [12] Edwin Brady. “Idris, a general-purpose dependently typed programming language: Design and implementation”. In: *Journal of Functional Programming* 23.05 (2013), pp. 552–593.
- [13] Simon Castellan. *Dependent type theory as the initial category with families*. 2014. URL: <http://iso.mor.phis.me/archives/2011-2012/stage-2012-goteburg/report.pdf> (visited on 03/10/2017).
- [14] Cyril Cohen et al. “Cubical type theory: a constructive interpretation of the univalence axiom”. In: *arXiv preprint arXiv:1611.02108* (2016).
- [15] nLab contributors. *closed monoidal structure on presheaves*. 2017. URL: <https://ncatlab.org/nlab/show/closed+monoidal+structure+on+presheaves> (visited on 04/07/2017).
- [16] nLab contributors. *preorder*. 2017. URL: <https://ncatlab.org/nlab/show/preorder> (visited on 04/07/2017).
- [17] nLab contributors. *wide subcategory*. 2017. URL: <https://ncatlab.org/nlab/show/wide+subcategory> (visited on 04/07/2017).

- [18] Catarina Coquand. “A formalised proof of the soundness and completeness of a simply typed lambda-calculus with explicit substitutions”. In: *Higher-Order and Symbolic Computation* 15.1 (2002), pp. 57–90.
- [19] Catarina Coquand. “A proof of normalization for simply typed lambda calculus written in ALF”. In: *Proceedings of the Workshop on Types for Proofs and Programs*. Citeseer. 1992, pp. 85–92.
- [20] Thierry Coquand and Peter Dybjer. “Intuitionistic model constructions and normalization proofs”. In: *Mathematical Structures in Computer Science* 7.1 (1997), pp. 75–94.
- [21] Agda developers. *Agda documentation*. 2017. URL: <https://agda.readthedocs.io/en/v2.5.2/> (visited on 03/10/2017).
- [22] Agda developers. *Agda standard library*. 2017. URL: <https://github.com/agda/agda-stdlib> (visited on 03/10/2017).
- [23] Healfdene Goguen, Conor McBride, and James McKinna. “Eliminating dependent pattern matching”. In: *Algebra, Meaning, and Computation* (2006), pp. 521–540.
- [24] Robert Harper. *Practical foundations for programming languages*. Cambridge University Press, 2016.
- [25] Steven Keuchel, Stephanie Weirich, and Tom Schrijvers. “Needle & Knot: Binder Boilerplate Tied Up”. In: *Proceedings of the 25th European Symposium on Programming Languages and Systems - Volume 9632*. New York, NY, USA: Springer-Verlag New York, Inc., 2016, pp. 419–445. ISBN: 978-3-662-49497-4. DOI: 10.1007/978-3-662-49498-1\_17. URL: [http://dx.doi.org/10.1007/978-3-662-49498-1\\_17](http://dx.doi.org/10.1007/978-3-662-49498-1_17).
- [26] András Kovács. *Normalization by evaluation for intrinsic System F*. 2017. URL: <https://github.com/AndrasKovacs/sysF-NbE> (visited on 04/07/2017).
- [27] Per Martin-Löf. “An intuitionistic theory of types: Predicative part”. In: *Studies in Logic and the Foundations of Mathematics* 80 (1975), pp. 73–118.
- [28] Conor McBride. *Datatypes of Datatypes*. 2015. URL: <http://staff.mmc.ssfed.ru/~ulysses/Edu/SSGEP/conor/conor.pdf>.
- [29] Conor McBride and Ross Paterson. “Applicative programming with effects”. In: *Journal of functional programming* 18.01 (2008), pp. 1–13.
- [30] Ulf Norell. *Towards a practical programming language based on dependent type theory*. Vol. 32. Citeseer, 2007.
- [31] Benjamin C Pierce. *Types and programming languages*. MIT press, 2002.
- [32] Steven Schäfer, Tobias Tebbi, and Gert Smolka. “Autosubst: Reasoning with de Bruijn Terms and Parallel Substitutions”. In: *Interactive Theorem Proving - 6th International Conference, ITP 2015, Nanjing, China, August 24-27, 2015*. Ed. by Xingyuan Zhang and Christian Urban. LNAI. Springer-Verlag, Aug. 2015.
- [33] Thomas Streicher. *Investigations into intensional type theory*. <http://www.mathematik.tu-darmstadt.de/~streicher/HabilStreicher.pdf>, 1993.

- [34] Aaron Stump. *Verified Functional Programming in Agda*. New York, NY, USA: Association for Computing Machinery and Morgan & Claypool, 2016. ISBN: 978-1-97000-127-3.
- [35] The Univalent Foundations Program. *Homotopy Type Theory: Univalent Foundations of Mathematics*. Institute for Advanced Study: <https://homotopytypetheory.org/book>, 2013.