

# The Cedilleum Language Specification

## Syntax, Typing, Reduction, and Elaboration

Christopher Jenkins

October 19, 2018

## 1 Introduction

The document describes *Cedilleum*, a general-purpose dependently typed programming language with inductive datatypes. Unlike most languages of this description, the underlying theory of Cedilleum is *not* the Calculus of Inductive Constructions (CIC)[PM15]. Instead, Cedilleum is designed so that it may easily be translated to *Cedille Core* – a compact core theory in which induction is derivable for lambda-encoded datatypes – while still providing high-level features like pattern-matching and recursive definitions. That said, the formal specification of Cedilleum as a self-contained language has a lot in common with CIC – see in particular Section 8 of [Inr18], which served as the basic template for much of this document’s formal development.

### 1.1 Data-type Declarations

Before diving into the details, let us take a bird’s-eye view of the language by showing some simple example data-type definitions and functions over them.

```
-- Non-recursive
data Bool: ★ =
  | tt: Bool
  | ff: Bool
.

-- Recursive
data Nat: ★ =
  | zero: Nat
  | succ: Nat → Nat
.

-- Recursive, parameterized, indexed
data Vec (A: ★): Nat → ★ =
  | vnil : Vec zero
  | vcons: ∀ n: Nat. A → Vec n → Vec (succ n)
.
```

Figure 1: Definition of natural numbers and length-indexed lists

Figure 1 shows some definitions of inductive datatypes, and modulo differences in syntax should seem straightforward to programmers that have used languages like Agda, Idris, or Coq. The key things to note are:

- In constructor type signatures, recursive occurrences of the inductive data-type being defined (such as in `suc : Nat → Nat`) must be positive, *but not strictly positive*.

- In parameterized types (like `Vec` with parameter  $(A: \star)$ ) occurrences of the inductive type being defined are not written applied to its parameters.

For example, the constructor declaration `vnil : Vec zero` results in the term `vnil` having type  $\forall A: \star. \text{Vec } A \star$  (with  $\cdot$  denoting type application)

- In the constructor declaration `vcons :  $\forall n: \text{Nat}. A \rightarrow \text{Vec } n \rightarrow \text{Vec } (\text{succ } n)$` , the argument `n` is *computationally irrelevant* (also called *erased*). This is because it is introduced by the irrelevant dependent function former  $\forall$ , as opposed to the relevant function former  $\Pi$ . More will be said of this when we discuss the type system of Cedilleum.

## 1.2 Function Definitions

```
-- Non-recursive
ite :  $\forall X: \star. \text{Bool} \rightarrow X \rightarrow X \rightarrow X$ 
    =  $\Lambda X. \lambda b. \lambda \text{then}. \lambda \text{else}. \mu' b \{$ 
      | tt  $\mapsto$  then
      | ff  $\mapsto$  else
     $\}$ .

pred :  $\text{Nat} \rightarrow \text{Nat}$ 
    =  $\lambda n. \mu' n \{$ 
      | zero  $\mapsto$  n
      | succ n'  $\mapsto$  n'
     $\}$ 

-- Recursive
add :  $\text{Nat} \rightarrow \text{Nat} \rightarrow \text{Nat}$ 
    =  $\lambda n. \lambda m.$ 
       $\mu \text{ add-rec}. n @(\lambda x: \text{Nat}. \text{Nat}) \{$ 
        | zero  $\mapsto$  m
        | succ r  $\mapsto$  succ (add-rec r)
       $\}$ .
```

Figure 2: Functions over inductive datatypes

Figure 2 shows functions defined over inductive datatypes using pattern matching and recursion. The first difference to note between the definitions is that functions `ite` and `pred` perform “mere” pattern matching on their arguments by using  $\mu'$ , whereas `add` uses  $\mu$  which provides combined pattern-matching and fix-point recursion. In `add`,  $\mu$  binds `add-rec` as the name of the fixpoint function for recursion on `n`. From this alone the reader might expect that  $\mu'$  is merely syntactic sugar for the more verbose  $\mu$ , but actually the difference is a bit more subtle than this, as we will see shortly.

The first major departure of Cedilleum from other languages with inductive datatypes can be seen in the type of `add-rec`. Ordinarily, its type would be  $\Pi x: \text{Nat}. \text{Nat}$  (corresponding to the motive  $(\lambda x: \text{Nat}. \text{Nat})$ ) – but in Cedilleum, its type is  $\text{Nat}/\text{add-rec} \rightarrow \text{Nat}$  (where we read  $\text{Nat}/\text{add-rec}$  as a single identifier) and by extension for `add-rec r` to be well-typed, the `r` bound in the pattern `succ r` must have type  $\text{Nat}/\text{add-rec}$  and not the more usual  $\text{Nat}$ . Why? For recursive functions in Cedilleum, *termination is guaranteed by the type system* and not by a separate syntactic check for structurally-decreasing arguments. The type  $\text{Nat}/\text{add-rec}$  indicates the types of those terms which the function `add-rec` may legally take as arguments without risking non-termination, and within a case branch variables of this type are introduced in the place of recursive arguments to constructors<sup>1</sup> – that is, `r` in the pattern `succ r`. In the next subsection, we will see more of how to use “recursive occurrence” types like  $\text{Nat}/\text{add-rec}$ , including how to do further pattern matching on them and conversion back to the original type, but for now it suffices to consider them an artifice for type-guided termination.

```

-- Recursive, parameterized, indexed
vappend : ∀ A: *. ∀ m: Nat. ∀ n: Nat. Vec ·A m → Vec ·A n → Vec ·A (add m n)
= Λ A. Λ m. Λ n. λ xs. λ ys.
  μ vappend-rec. xs @(λ i: Nat. λ zs: Vec ·A i. Vec ·A (add i n)) {
    | vnil           ↦ ys
    | vcons -m' x xs' ↦ [ zs = vappend-rec -m' xs' ] - vcons -(add m' n) x zs
  }.

```

Figure 3: Dependent functions over inductive datatypes

Figure 3 shows the classic dependent function `vappend` over length-indexed lists `Vec`. Like `add`, it is defined by fixpoint recursion, here over the argument `xs`, and as with `add-rec` the fixpoint function `vappend-rec` has type  $\forall i: \text{Nat}. \Pi zs: \text{Vec}/\text{vappend-rec } i. \text{Vec} \cdot A (\text{add } i \text{ } n)$ , where the  $\forall$  indicates that `i` is an *irrelevant* argument. Note again the missing parameter `A` in the type `Vec/vappend-rec i` – this is not a typo, but rather an indication that `A` is “baked-in” to the type `Vec/vappend-rec`. Aside from this the two cases of `append` are mostly straightforward: in the `vnil` branch the expected type is `Vec ·A (add zero n)` which converts to `Vec ·A n`, so `ys` suffices; in the `vcons` branch we bind subdata `m': Nat`, `x: A`, and `xs': Vec/vappend-rec m'`, with `-m'` indicating that `m'` is bound *irrelevantly*, then we make a local binding `zs` by invoking `vappend` recursively on `m'` and `xs'` (where again the `-m'` indicates `m'` is an *irrelevant* argument) before producing a result whose type is convertible with the expected `Vec ·A (add (suc m') n)`.

### 1.3 Histomorphic Recursion

We now return to the strange types `Nat/add-rec` and `Vec/append-rec`. The reader may well ask whether they serve any purpose other than marking what recursive calls are legal in a function – and the answer is yes! Cedilleum uses types like these to support *histomorphic* or *course-of-values* recursion, which is to say recursion schemes that can dig arbitrarily deeply into the recursive occurrences of data to compute results. As a motivating example, Figure 4 shows how we can define division on natural numbers in Cedilleum.

## 2 Syntax

**Identifiers** Figure 5 gives the metavariables used in our grammar for identifiers. We consider all identifiers as coming from two distinct lexical “pools” – regular identifiers (consisting of identifiers *id* given for modules and definitions, term variables *u*, and type variables *X*) and kind identifiers  $\kappa$ . In Cedilleum source files (as in the parent language Cedille) kind variables should be literally prefixed with  $\kappa$  – the suffix can be any string that would by itself be a legal non-kind identifier. For example, `myDef` is a legal term / type variable and a legal name for a definition, whereas  $\kappa\text{myDef}$  is only legal as a kind definition.

**Untyped Terms** The grammar of pure (untyped) terms the untyped  $\lambda$ -calculus augmented with a primitives for combination fixed-point and pattern-matching definitions (and an auxiliary pattern-matching construct).

**Modules and Definitions** All Cedilleum source files start with production *mod*, which consists of a module declaration, a sequence of import statements which bring into scope definitions from other source files, and a sequence of *commands* defining terms, types, and kinds. As an illustration, consider the first few lines of a hypothetical `list.ced`:

```

module list .

import nat .

```

```

iterate : ∀ R: *. Nat → (R → R) → R → R
= Λ R. λ n. λ f. λ x.
  μ iterate-rec. n @(λ _: Nat. Nat) {
    | zero    ↦ r
    | succ r' ↦ f (iterate-rec r')
  }.

minus : Nat → Nat → Nat
= λ m. λ n.
  μ minus-rec. n @(λ _: Nat. Nat) {
    | zero    ↦ m
    | succ r  ↦ pred (minus-rec r)
  }.

lt : Nat → Nat → Bool
= λ m. λ n.
  μ' (minus m n) {
    | zero    ↦ ff
    | succ r  ↦ tt
  }.

divide : Nat → Nat → Nat
= λ m. λ n.
  μ divide-rec. m @(λ _: Nat. Nat) {
    | zero    ↦ zero
    | succ r  ↦
      ite (lt (succ (fromNat/divide-rec r)) n) zero
        ([ pred' = λ x: Nat/divide-rec. μ' x { | zero ↦ x | succ x' ↦ x' }
          ] - succ (divide-rec (iterate (pred n) pred' r)))
  }

```

Figure 4: Histomorphic recursion and division

Imports are handled first by consulting a global options files known to the Cedilleum compiler (on \*nix systems `~/cedille/options`) containing a search path of directories, and next (if that fails) by searching the directory containing the file being checked.

Term and type definitions are given with an identifier, a classifier (type or kind, resp.) to check the definition against, and the definition. For term definitions, giving classifier (i.e. the type) is optional. As an example, consider the definitions for the type of Church-encoded lists and two variants of the nil constructor, the first with a top-level type annotation and the second with annotations sprinkled on binders:

```

cList : * → *
= λ A : * . ∀ X : * . (A → X → X) → X → X .

cNil  : ∀ A : * . cList · A
= Λ A . Λ X . λ c . λ n . n .
cNil' = Λ A : * . Λ X : * . λ c : A → X → X . λ n : X . n .

```

Kind definitions are given without classifiers (all kinds have super-kind  $\square$ ), e.g.  $\kappa\text{func} = * \rightarrow *$

Inductive datatype definitions take a set of *parameters* (term and type variables which remain constant throughout the definition) well as a set of *indices* (term and type variables which *can* vary), followed by zero or more constructors. Each constructor begins with “|” (though the grammar can be relaxed so that the first of these is optional) and then an identifier and type is given. As an example, consider the following two definitions for lists and vectors (length-indexed lists).

$id$		identifiers for definitions
$u, c$		term variables
$X$		type variables
$\kappa$		kind variables
$x$	$::= id \mid u \mid X$	non-kind variables
$y$	$::= x \mid \kappa$	all variables

Figure 5: Identifiers

$p$	$::= u$	variables
	$\lambda u. p$	functions
	$p p$	applications
	$\mu u. p \{pcase^*\}$	fixed-point and pattern matching
	$\mu' p \{pcase^*\}$	simple pattern matching
$pcase$	$::= \mid u u^* \mapsto p$	

Figure 6: Untyped terms

```

data Bool : * =
  | tt : Bool
  | ff : Bool
  .

data Nat : * =
  | zero : Nat
  | suc : Nat → Nat
  .

data List (A : *) : * =
  | nil : List
  | cons : A → List → List
  .

data Vec (A : *) : Nat → * =
  | vnil : ∀ n : Nat . {n ≃ zero} ⇒ Vec n
  | vcons : ∀ n : Nat . ∀ m : Nat . A → Vec n → {m ≃ suc n} ⇒ Vec m
  .

```

**Types and Kinds** In Cedilleum, the expression language is stratified into three main “classes”: kinds, types, and terms. Kinds and types are listed in Figure 8 and terms are listed in Figure 9 along with some auxiliary grammatical categories. In both of these figures, the constructs forming expressions are listed from lowest to highest precedence – “abstractors” ( $\lambda \Lambda \Pi \forall$ ) bind most loosely and parentheses most tightly. Associativity is as-expected, with arrows ( $\rightarrow \Rightarrow$ ) and applications being left-associative and abstractors being right-associative.

The language of kinds and types is similar to that found in the Calculus of Implicit Constructions<sup>2</sup>. Kinds are formed by dependent and non-dependent products ( $\Pi$  and  $\rightarrow$ ) and a base kind for types which can classify terms ( $*$ ). Types are also formed by the usual (dependent and non-dependent) products ( $\Pi$  and  $\rightarrow$ ) and also *implicit* products ( $\forall$  and  $\Rightarrow$ ) which quantify over erased arguments (that is, arguments that

---

<sup>2</sup>Cite

<i>mod</i>	::=	<b>module</b> <i>id</i> . <i>imprt</i> * <i>cmd</i> *	module declarations
<i>imprt</i>	::=	<b>import</b> <i>id</i> .	module imports
<i>cmd</i>	::=	<i>defTermOrType</i> <i>defDataType</i> <i>defKind</i>	definitions
<i>defTermOrType</i>	::=	<i>id</i> <i>checkType</i> ? = <i>t</i> . <i>id</i> : <i>K</i> = <i>T</i> .	term definition type definition
<i>defKind</i>	::=	$\kappa = K$	kind definition
<i>defDataType</i>	::=	<b>data</b> <i>id param</i> * : <i>K</i> = <i>constr</i> * .	datatype definitions
<i>checkType</i>	::=	: <i>T</i>	annotation for term definition
<i>param</i>	::=	( <i>x</i> : <i>C</i> )	
<i>constr</i>	::=	<i>id</i> : <i>T</i>	

Figure 7: Modules and definitions

disappear at run-time).  $\Pi$ -products are only allowed to quantify over terms as all types occurring in terms are erased at run-time, but  $\forall$ -products can quantify over types *and* terms because terms can be erased. Meanwhile, non-dependent products ( $\rightarrow$  and  $\Rightarrow$ ) can only “quantify” over terms because non-dependent type quantification does not seem particularly useful. Besides these, Cedilleum features type-level functions and applications (with term and type arguments), and a primitive equality type for untyped terms. Last of all is the “hole” type ( $\bullet$ ) for writing partial type signatures or incomplete type applications. There are term-level holes as well, and together the two are intended to help facilitate “hole-driven development”: any hole automatically generates a type error and provides the user with useful contextual information.

We illustrate with another example: what follows is a module stub for **DepCast** defining dependent casts – intuitively, functions from  $a : A$  to  $B$   $a$  that are also equal<sup>3</sup> to identity – where the definitions **CastE** and **castE** are incomplete.

```
module DepCast .
```

```
CastE  $\triangleleft$   $\Pi$  A :  $\star$  . (A  $\rightarrow$   $\star$ )  $\rightarrow$   $\star$  =  $\bullet$  .
```

```
castE  $\triangleleft$   $\forall$  A :  $\star$  .  $\forall$  B : A  $\rightarrow$   $\star$  . CastE  $\cdot$  A  $\cdot$  B  $\Rightarrow$   $\Pi$  a : A . B a =  $\bullet$  .
```

**Annotated Terms** Terms can be explicit and implicit functions (resp. indicated by  $\lambda$  and  $\Lambda$ ) with optional classifiers for bound variables, let-bindings, applications  $t\ t'$ ,  $t\ -t'$ , and  $t\ \cdot T$  (resp. to another term, an erased term, or a type). In addition to this there are a number of useful operators for equational reasoning, type casting, providing annotations, and pattern matching. Each operator will be discussed in more detail in Section 4, but a few concrete programs in Cedilleum are given below merely to give a better idea of the syntax of the language.

```
isvnil :  $\forall$  A :  $\star$  .  $\forall$  n : Nat . Vec  $\cdot$  A n  $\rightarrow$  Bool
=  $\Lambda$  A .  $\Lambda$  n .  $\lambda$  xs .
   $\mu'$  xs @( $\Lambda$  n .  $\lambda$  xs . Bool)
  { | vnil -n -eq  $\mapsto$  tt
    | vcons -n -m x xs -eq  $\mapsto$  ff
  }
```

---

<sup>3</sup>Module erasure, discussed below

Sorts $S$	$::= \square$	sole super-kind
	$K$	kinds
Classifiers $C$	$::= K$	kinds
	$T$	types
Kinds $K$	$::= \Pi x : C . K$	explicit product
	$C \rightarrow K$	kind arrow
	$\star$	the kind of types that classify terms
	$(K)$	
Types $T$	$::= \Pi x : T . T$	explicit product
	$\forall x : C . T$	implicit product
	$\lambda x : C . T$	type-level function
	$T \Rightarrow T'$	arrow with erased domain
	$T \rightarrow T'$	normal arrow type
	$T \cdot T'$	application to another type
	$T t$	application to a term
	$\{ p \simeq p' \}$	untyped equality
	$(T)$	
	$X$	type variable
	$\bullet$	hole for incomplete types

Figure 8: Kinds and types

```

vlength :  $\forall A : \star . \forall n : \text{Nat} . \text{Vec} \cdot A \ n \rightarrow \text{Nat}$ 
=  $\Lambda A . \Lambda n . \lambda xs .$ 
   $\mu \text{len} . xs @ (\Lambda n . \lambda x . \text{Nat})$ 
  { | vnil -n -eq  $\mapsto$  zero
    | vcons -n -m x xs -eq  $\mapsto$  suc (len -n xs)
  }

```

### 3 Erasure

The definition of the erasure function given in Figure 10 takes the annotated terms from Figures 8 and 9 to the untyped terms of Figure 6. The last two equations indicate that any type or erased arguments in the the zero or more *vararg*'s of pattern-match case are indeed erased. The additional constructs introduced in the annotated term language such as  $\beta$ ,  $\phi$ , and  $\rho$ , are all erased to the language of pure terms.

### 4 Type System (sans Inductive Datatypes)

The inference rules for classifying expressions in Cedilleum are stratified into two judgments. Figure 12 gives the uni-directional rules for ensuring types are well-kinded and kinds are well-formed. Future versions of Cedilleum will allow for bidirectional checking for both typing *and* sorting, allowing for a unification of these two figures. Most of these rules are similar to what one would expect from the Calculus of Implicit Constructions, so we focus on the typing rules unique to Cedilleum.

The typing rule for  $\rho$  shows that  $\rho$  is a primitive for rewriting by an (untyped) equality. If  $t$  is an expression that synthesizes a proof that two terms  $t_1$  and  $t_2$  are equal, and  $t'$  is an expression synthesizing type  $[t_1/x] T$  (where, as per the footnote,  $t_1$  does not occur in  $T$ ), then we may essentially rewrite its type to

<sup>5</sup>Where we assume  $t$  does not occur anywhere in  $T$

<sup>5</sup>Where  $\mathbf{tt} = \lambda x. \lambda y. x$  and  $\mathbf{ff} = \lambda x. \lambda y. y$

Subjects $s$	$::= t$	term
	$T$	type
Terms $t$	$::= \lambda x \text{ class}^?.t$	normal abstraction
	$\Lambda x \text{ class}^?.t$	erased abstraction
	$[ \text{defTermOrType} ] - t$	let definitions
	$\rho \ t - t'$	equality elimination by rewriting
	$\phi \ t - t' \ \{t''\}$	type cast
	$\chi \ T - t$	check a term against a type
	$\delta - t$	ex falso quodlibet
	$\theta \ t \ t'^*$	elimination with a motive
	$t \ t'$	applications
	$t \ -t'$	application to an erased term
	$t \ .T$	application to a type
	$\beta \ \{t\}$	reflexivity of equality
	$\varsigma \ t$	symmetry of equality
	$\mu \ u \ . \ t \ \text{motive}^? \ \{\text{case}^*\}$	type-guarded pattern match and fixpoint
	$\mu' \ t \ \text{motive}^? \ \{\text{case}^*\}$	auxiliary pattern match
	$u$	term variable
	$(t)$	
	$\bullet$	hole for incomplete term
$\text{case}$	$::= \mid c \ \text{vararg}^* \mapsto t$	pattern-matching cases
$\text{vararg}$	$::= u$	normal constructor argument
	$-u$	erased constructor argument
	$\cdot X$	type constructor argument
$\text{class}$	$::= : C$	
$\text{motive}$	$::= @ T$	motive for induction

Figure 9: Annotated Terms

$[t_2/x] T$ . The rule for  $\beta$  is reflexivity for equality – it witnesses that a term is equal to itself, provided that the type of the equality is well-formed. The rule for  $\varsigma$  is symmetry for equality. Finally,  $\phi$  acts as a “casting” primitive: the rule for its use says that if some term  $t$  witnesses that two terms  $t_1$  and  $t_2$  are equal, and  $t_1$  has been judged to have type  $T$ , then intuitively  $t_2$  can also be judged to have type  $T$ . (This intuition is justified by the erasure rule for  $\phi$  – the expression erases to  $[t_2]$ ). The last rule involving equality is for  $\delta$ , which witnesses the logical principle *ex falso quodlibet* – if a certain impossible equation is proved (namely that the two Church-encoded booleans **tt** and **ff** are equal), then *any* type desired is inhabited.

The two remaining primitives are not essential to the theory but are useful additions for programmers. The rule for  $\chi$  allows the user to provide an explicit top-level annotation for a term, and  $\theta$  embodies “elimination with a motive”, using the expected type of an application to infer some type arguments. (TODO)

## 5 Inductive Datatypes

Before we can provide the typing rules for introduction and usage of inductive datatypes, some auxiliary definitions must be given. The syntax for these, and the structure of this entire section, borrows heavily from the conventions of the Coq documentation<sup>6</sup>. The author believes it is worthwhile to restate this development in terms of the Cedille type system, rather than merely pointing readers to the Coq documentation and

<sup>6</sup><https://coq.inria.fr/refman/language/cic.html#inductive-definitions>



$ x $	$=$	$x$
$ \star $	$=$	$\star$
$ \square $	$=$	$\square$
$ \beta \{t\} $	$=$	$ t $
$ \delta t $	$=$	$ t $
$ \chi T^? - t $	$=$	$ t $
$ \theta t t'^* $	$=$	$ t   t'^* $
$ \varsigma t $	$=$	$ t $
$ t t' $	$=$	$ t   t' $
$ t - t' $	$=$	$ t $
$ t \cdot T $	$=$	$ t $
$ \rho t - t' $	$=$	$ t' $
$ \forall x:C. C' $	$=$	$\forall x: C .  C' $
$ \Pi x:C. C' $	$=$	$\Pi x: C .  C' $
$ \lambda u:T. t $	$=$	$\lambda u.  t $
$ \lambda u. t $	$=$	$\lambda u.  t $
$ \lambda X:K. C $	$=$	$\lambda X: K .  C $
$ \Lambda x:C. t $	$=$	$ t $
$ \phi t - t' \{t''\} $	$=$	$ t'' $
$ [x = t : T] - t' $	$=$	$(\lambda x.  t' )  t $
$ [X = T : K] - t $	$=$	$ t $
$ \{t \simeq t'\} $	$=$	$\{ t  \simeq  t' \}$
$ \mu u. \cdot t \text{ motive}^? \{case^*\} $	$=$	$\mu u.  t  \{ case^* \}$
$ \mu' t \text{ motive}^? \{case^*\} $	$=$	$\mu'  t  \{ case^* \}$
$ id \text{ vararg}^* \mapsto t $	$=$	$id  vararg^*  \mapsto  t $
$ -u $	$=$	
$ \cdot T $	$=$	

Figure 10: Erasure for annotated terms

asking them to infer the differences between the two systems.

To begin with, the production *defDataType* gives the concrete syntax for datatype definitions, but it is not a very useful notation for representing one in the abstract syntax tree. In our typing rules we will instead use the notation  $\text{Ind}_M[p](\Gamma_I := \Sigma)$ , where

- $M$  is a meta-variable ranging over constant labels “C” and “A” (used to distinguish **c**oncrete and **a**bstracted inductive definitions – more on this below)
- $p$  is the number of **p**arameters of the inductive definition
- $\Gamma_I$  is a typing context binding *one* type variable  $I$ , the inductive type being defined
- $\Sigma$  is a typing context containing the  $n$  data constructors  $c_1, \dots, c_n$  of  $I$ .

For example, consider the **List** and **Vec** definitions from Section 2. These will be represented in the AST as

$$\text{Ind}_C[1](List : \star \rightarrow \star := \begin{array}{ll} nil & : \forall A : \star. List \cdot A \\ cons & : \forall A : \star. A \rightarrow List \cdot A \rightarrow List \cdot A \end{array})$$

Figure 11: Contexts

Typing contexts  $\Gamma ::= \emptyset \mid x:C, \Gamma \mid x=s:C, \Gamma$

$$\begin{array}{c}
\overline{\Gamma \vdash \star : \square} \\
\\
\frac{\Gamma \vdash C : S \quad \Gamma, y : C \vdash C' : S'}{\Gamma \vdash \Pi y : C. C' : S'} \quad \frac{\Gamma \vdash C : S \quad \Gamma, y : C \vdash C' : \star}{\Gamma \vdash \forall y : C. C' : \star} \\
\\
\frac{FV(p \ p') \subseteq \text{dom}(\Gamma)}{\Gamma \vdash \{p \simeq p'\} : \star} \quad \overline{\Gamma \vdash \kappa : \Gamma(\kappa)} \quad \overline{\Gamma \vdash X : \Gamma(X)} \\
\\
\frac{\Gamma \vdash \Pi x : C. K : \square \quad \Gamma, x : C \vdash T : K}{\Gamma \vdash \lambda x : C. T : \Pi x : C. K} \quad \frac{\Gamma \vdash T : \Pi x : K. K' \quad \Gamma \vdash T' : K}{\Gamma \vdash T \cdot T' : [T'/x]K'} \quad \frac{\Gamma \vdash T : \Pi x : T'. K \quad \Gamma \vdash_\downarrow t : T'}{\Gamma \vdash T t : [t/x]K}
\end{array}$$

Figure 12: Sort checking  $\boxed{\Gamma \vdash C : S}$ 

$$\begin{array}{c}
\overline{\Gamma \vdash_\delta u : \Gamma(u)} \quad \frac{\Gamma \vdash T : K \quad \Gamma, x : T \vdash_\delta t : T'}{\Gamma \vdash_\delta \lambda x : T. t : \Pi x : T. T'} \quad \frac{\Gamma, x : T \vdash_\downarrow t : T'}{\Gamma \vdash_\downarrow \lambda x. t : \Pi x : T. T'} \\
\\
\frac{\Gamma \vdash C : S \quad x \notin FV(|t|) \quad \Gamma, x : C \vdash_\delta t : T}{\Gamma \vdash_\delta \Lambda x : C. t : \forall x : C. T} \quad \frac{x \notin FV(|t|) \quad \Gamma, x : C \vdash_\delta t : T}{\Gamma \vdash_\downarrow \Lambda x. t : \forall x : C. T} \quad \frac{\Gamma \vdash_\uparrow t : \Pi x : T'. T \quad \Gamma \vdash_\downarrow t' : T'}{\Gamma \vdash_\delta t t' : [t'/x]T} \\
\\
\frac{\Gamma \vdash_\uparrow t : \forall X : K. T' \quad \Gamma \vdash T : K}{\Gamma \vdash_\delta t \cdot T : [T/X]T'} \quad \frac{\Gamma \vdash_\uparrow t : \forall x : T'. T \quad \Gamma \vdash_\downarrow t' : T'}{\Gamma \vdash_\delta t \cdot t' : [t'/x]T} \quad \frac{\Gamma \vdash_\uparrow t : T' \quad |T'| =_\beta |T|}{\Gamma \vdash_\downarrow t : T} \\
\\
\frac{\Gamma \vdash T : K \quad \Gamma \vdash_\downarrow t : T \quad \Gamma, id = t : T \vdash_\delta t' : T'}{\Gamma \vdash_\delta [id : T = t] \cdot t' : T'} \quad \frac{\Gamma \vdash_\uparrow t : T \quad \Gamma, id = t : T \vdash_\delta t' : T'}{\Gamma \vdash_\delta [id = t] \cdot t' : T'} \quad \frac{\Gamma \vdash_\uparrow t : \{t_1 \simeq t_2\} \quad \Gamma \vdash_\uparrow t' : [t_1/x] T}{\Gamma \vdash_\delta \rho t \cdot t' : [t_2/x] T} \quad 4 \\
\\
\frac{\Gamma \vdash K : \square \quad \Gamma \vdash T : K \quad \Gamma, id = T : K \vdash_\delta t' : T'}{\Gamma \vdash_\delta [id : K = T] \cdot t' : T'} \quad \frac{\Gamma \vdash \{t' \simeq t'\} : \star}{\Gamma \vdash_\downarrow \beta \{t\} : \{t' \simeq t'\}} \quad \frac{\Gamma \vdash_\delta t : \{t_1 \simeq t_2\}}{\Gamma \vdash_\delta \varsigma t : \{t_2 \simeq t_1\}} \\
\\
\frac{\Gamma \vdash_\downarrow t : \{|t_1| \simeq |t_2|\} \quad \Gamma \vdash_\delta t_1 : T}{\Gamma \vdash_\delta \phi t \cdot t_1 \{t_2\} : T} \quad \frac{\Gamma \vdash_\downarrow t : T}{\Gamma \vdash_\uparrow \chi T \cdot t : T} \quad \frac{\Gamma \vdash_\downarrow t : \{\mathbf{tt} \simeq \mathbf{ff}\}}{\Gamma \vdash_\downarrow \delta \cdot t : T} \quad 5 \\
\\
\frac{\Gamma \vdash_\uparrow t : ?? \quad \Gamma \vdash_\downarrow t'^* : ??}{\Gamma \vdash_\downarrow \theta t t'^* : T}
\end{array}$$

Figure 13: Type checking  $\boxed{\Gamma \vdash_\delta s : C}$  (sans inductive datatypes)

and

$$\text{Ind}_C[1](Vec : \star \rightarrow Nat \rightarrow \star := \begin{array}{ll} vnil & : \forall A : \star. \forall n : Nat. \{n \simeq zero\} \Rightarrow Vec \cdot A \ n \\ vcons & : \forall A : \star. \forall n : Nat. \forall m : Nat. A \rightarrow Vec \cdot A \ n \rightarrow \{m \simeq succ \ n\} \Rightarrow Vec \cdot A \ m \end{array})$$

All inductive types the user will define will be concrete inductive definitions, and have global scope. Abstracted definitions are automatically generated during fix-point pattern matching, and have local scope.

For an inductive datatype definition to be well-formed, it must satisfy the following conditions (each of which is explained in more detail in Subsections 5.1 and 5.2):

- The kind of  $I$  must be (at least) a *p-arity of kind*  $\star$ .
- The types of each  $id \in \Sigma$  must be *types of constructors of*  $I$

- The definition must satisfy the *non-strict* positivity condition.

Similarly, the notation in the grammar of Cedilleum  $\mu'$  and  $\mu$  for pattern matching is inconvenient, and we will represent them in the AST as resp.  $\mu'(t, P, t_{i=1..n})$  and  $\mu(x_{\text{rec}}, I', x_{\text{to}}, t, P, t_{i=1..n})$ . Translation from the form given in the grammar to this form is discussed in detail below, but is as expected. In particular, we enforce that patterns are exhaustive and non-overlapping, and that  $I'$  and  $x_{\text{to}}$  (which are not present in the grammar but an automatically generated identifier) are fresh w.r.t the global and local context. For example, consider the pattern-matches given in the code listings for `isvnil` and `vlength` above. These would be translated into the AST as

$$\mu'(xs, \Lambda n. \lambda x. \text{Bool}, \begin{array}{l} \Lambda n. \Lambda eq. tt \\ \Lambda n. \Lambda m. \lambda x. \lambda xs. \Lambda eq. ff \end{array})$$

and

$$\mu(\text{len}, \text{Vec}/\text{len}, \text{toVec}/\text{len}, xs, \Lambda n. \lambda x. \text{Nat}, \begin{array}{l} \Lambda n. \Lambda eq. \text{zero} \\ \Lambda n. \Lambda m. \lambda x. \lambda xs. \Lambda eq. \text{suc}(\text{len} - n \text{ xs}) \end{array})$$

In general, the generated name for  $I'$  and  $x_{\text{to}}$  that users will write in Cedilleum programs will be of the form “ $I/x_{\text{rec}}$ ” and “ $\text{to}I/x_{\text{rec}}$ ”.

For a pattern construct ( $\mu$  or  $\mu'$ ) in the AST to be well-formed, it must satisfy the following conditions (each of which is, again, explained in more detail in Subsections 5.3, 5.5, and 5.6):

- The motive  $P$  must be well-kinded
- $P$  must be a legal motive to be used in eliminating the inductive type  $I$  of the scrutinee  $t$
- Each branch  $t_i$  must have the type expected given the constructor  $c_i \in \Sigma$  and the motive  $P$ .

## 5.1 Auxiliary Definitions

**Contexts** To ease the notational burden, we will introduce some conventions for writing contexts within terms and types.

- We write  $\lambda \Gamma$ ,  $\Lambda \Gamma$ ,  $\forall \Gamma$ , and  $\Pi \Gamma$  to indicate some form of abstraction over each variable in  $\Gamma$ . For example, if  $\Gamma = x_1:T_1, x_2:T_2$  then  $\lambda \Gamma. t = \lambda x_1:T_1. \lambda x_2:T_2. t$ . Additionally, we will also write  $\prod \Gamma$  to indicate an arbitrary mixture of  $\Pi$  and  $\forall$  quantified variables. Note that *if  $\prod \Gamma$  occurs multiple times within a definition or inference rule, the intended interpretation is that all occurrences have the same mixture of  $\Pi$  and  $\forall$  quantifiers.*
- $\|\Gamma\|$  denotes the length of  $\Gamma$  (the number of variables it binds)
- We write  $s \Gamma$  to indicate the sequence of variable arguments in  $\Gamma$  given as arguments to  $s$ . Implicit in this notation is the removal of typing annotations from the variables  $\Gamma$  when these variables are given as arguments to  $s$ .

Since in Cedilleum there are three flavors of applications (to a type, to an erased term, and to an unerased term), we will only use this notion when the type or kind of  $s$  is known, which is sufficient to disambiguate the flavor of application intended for each particular binder in  $\Gamma$ . For example, if  $s$  has type  $\forall X:\star. \forall x:X. \Pi x':X. X$  and  $\Gamma = X:\star, x:X, x':X$  then  $s \Gamma = s \cdot X \cdot x \cdot x'$

- $\Delta$  and  $\Delta'$  are notations we will use for a specially designated contexts associating type variables with both global “concrete” and local “abstracted” inductive data-type declarations. The purpose of this latter sort of declaration is to enable type-guided termination of definitions using fixpoints (see Section 5.7) For example, given just the (global) data type declaration of  $\text{Vec}$ , we would have  $\Delta(\text{Vec}) = \text{Ind}_C[1](\Gamma_{\text{Vec}} := \Sigma)$ , where  $\Gamma_{\text{Vec}} = \text{Vec}:\star \rightarrow \text{Nat} \rightarrow \star$  and  $\Sigma$  binds data constructors  $\text{vnil}$  and  $\text{vcons}$  to the appropriate types.

**$p$ -arity** A kind  $K$  is a  $p$ -arity if it can be written as  $\Pi \Gamma. K'$  for some  $\Gamma$  and  $K'$ , where  $\|\Gamma\| = p$ . For an inductive definition  $\text{Ind}_M[p](\Gamma_I := \Sigma)$ , requiring that the kind  $\Gamma_I(I)$  is a  $p$ -arity of  $\star$  ensures that  $I$  *really does have*  $p$  parameters.

**Types of Constructors**  $T$  is a *type of a constructor of  $I$*  iff

- it is  $I s_1 \dots s_n$
- it can be written as  $\forall s:C. T$  or  $\Pi s:C. T$ , where (in either case)  $T$  is a type of a constructor of  $I$

**Positivity condition** The positivity condition is defined in two parts: the positivity condition of a type  $T$  of a constructor of  $I$ , and the positive occurrence of  $I$  in  $T$ . We say that a type  $T$  of a constructor of  $I$  satisfies the positivity condition when

- $T$  is  $I s_1 \dots s_n$  and  $I$  does not occur anywhere in  $s_1 \dots s_n$
- $T$  is  $\forall s:C. T'$  or  $\Pi s:C. T'$ ,  $T'$  satisfies the positivity condition for  $I$ , and  $I$  occurs *only* positively in  $C$

We say that  $I$  occurs only positively in  $T$  when

- $I$  does not occur in  $T$
- $T$  is of the form  $I s_1 \dots s_n$  and  $I$  does not occur in  $s_1 \dots s_n$
- $T$  is of the form  $\forall s:C. T'$  or  $\Pi s:C. T'$ ,  $I$  occurs only positively in  $T'$ , and  $I$  *does not* occur positively in  $C$

## 5.2 Well-formed inductive definitions

Let  $\Gamma_P$ ,  $\Gamma_I$ , and  $\Sigma$  be contexts such that  $\Gamma_I$  associates a single type-variable  $I$  to kind  $\Pi \Gamma_P. K$  and  $\Sigma$  associates term variables  $c_1 \dots c_n$  with corresponding types  $\forall \Gamma_P. T_1, \dots \forall \Gamma_P. T_n$ . Then the rule given in Figure 14 states when an inductive datatype definition may be introduced, provided that the following side conditions hold:

Figure 14: Introduction of inductive datatype

$$\frac{\emptyset \vdash \Gamma_I(I) : \square \quad \|\Gamma_P\| = p \quad (\Gamma_I, \Gamma_P \vdash T_i : \star)_{i=1..n}}{\text{Ind}_M[p](\Gamma_I := \Sigma) \text{ wf}}$$

- Names  $I$  and  $c_1 \dots c_n$  are distinct from any other inductive datatype type or constructor names, and distinct amongst themselves
- Each of  $T_1 \dots T_n$  is a type of constructor of  $I$  which satisfies the positivity condition for  $I$ . Furthermore, each occurrence of  $I$  in  $T_i$  is one which is applied to the parameters  $\Gamma_P$ .
- Identifiers  $I, c_1, \dots, c_n$  are fresh w.r.t the global context, and do not overlap with each other nor any identifiers in  $\Gamma_P$ .

When an inductive data-type has been defined using the *defDataType* production, it is understood that this always a concrete inductive type, and it (implicitly) adds to a global typing context the variable bindings in  $\Gamma_I$  and  $\Sigma$ . Similarly, when checking that the kind  $\Gamma_I(I)$  and type  $T_i$  are well-sorted and well-kinded, we assume an (implicit) global context of previous definitions.

Figure 15: Valid elimination kinds

$$\frac{}{\llbracket T : \star \mid T \rightarrow \star \rrbracket} \quad \frac{\llbracket T s : K \mid K' \rrbracket}{\llbracket T : \Pi s : C. K \mid \Pi s : C. K' \rrbracket}$$

### 5.3 Valid Elimination Kind

When type-checking a pattern match (either  $\mu$  or  $\mu'$ ), we need to know that the given motive  $P$  has a kind  $K$  for which elimination of a term with some inductive data-type  $I$  is permissible. We write this judgment as  $\llbracket T : K' \mid K \rrbracket$ , which should be read “the type  $T$  of kind  $K'$  can be eliminated through pattern-matching with a motive of kind  $K$ ”. This judgment is defined by the simple rules in Figure 15. For example, a valid elimination kind for the indexed type family  $Vec \cdot X$  (which has kind  $\Pi n : Nat. \star$ ) is  $\Pi n : Nat. \Pi x : Vec \cdot X n. \star$

### 5.4 Valid Branch Type

Another piece of kit we need is a way to ensure that, in a pattern-matching expression, a particular branch has the correct type given a particular constructor of an inductive data-type and a motive. We write  $\{\{c : T\}\}_I^P$  to indicate the type corresponding to the (possibly partially applied) constructor  $c$  of  $I$  and its type  $T$ . We abbreviate this notation to  $\{\{c\}\}^P$  when the inductive type variable  $I$ , and the type  $T$  of  $c$ , is known from the (meta-language) context.

$$\begin{aligned} \{\{c : I \bar{T} \bar{s}\}\}_I^P &= P \bar{s} c \\ \{\{c : \forall x : T'. T\}\}_I^P &= \forall x : T'. \{\{c \cdot x : T\}\}_I^P \\ \{\{c : \forall x : K. T\}\}_I^P &= \forall x : K. \{\{c \cdot x : T\}\}_I^P \\ \{\{c : \Pi x : T'. T\}\}_I^P &= \Pi x : T'. \{\{c x : T\}\}_I^P \end{aligned}$$

where we leave implicit the book-keeping required to separate the parameters  $\bar{T}$  from the indices  $\bar{s}$ .

The biggest difference between this definition and the similar one found in the Coq documentation is that types can have implicit and explicit quantifiers, so we must make sure that the types of branches have implicit / explicit quantifiers (and the subjects  $c$  have applications for types, implicit terms, and explicit terms), corresponding to those of the arguments to the data constructor for the pattern for the branch.

### 5.5 Well-formed Patterns

Figure 16: Well-formedness of a pattern

$$\frac{\Gamma \vdash P : K \quad \Sigma = c_1 : \forall \Gamma_P. T_1, \dots, c_n : \forall \Gamma_P. T_n \quad \|\bar{T}\| = \|\Gamma_P\| = p \quad \llbracket I \bar{T} : \Gamma(I) \mid K \rrbracket \quad (\Gamma, \Delta \vdash_\downarrow t_i : \{\{c_i \bar{T}\}\}^P)_{i=1..n}}{WF-Pat(\Gamma, \Delta, \text{Ind}_M[p](\Gamma_I := \Sigma), \bar{T}, \mu'(t, P, t_{i=1..n}))}$$

Figure 16 gives the rule for checking that a pattern  $\mu'(t, P, t_{i=1..n})$  is well-formed. We check that the motive  $P$  is well-kinded at kind  $K$ , that the given parameters  $\bar{T}$  match the expected number  $p$  from the inductive data-type declaration, that an inductive data-type  $I$  instantiated with the given parameters  $\bar{T}$  can be eliminated to a type of kind  $K$ , and that the given branches  $t_i$  account for each of the constructors  $c_i$  of  $\Sigma$  and have the required branch type  $\{\{c_i \bar{T}\}\}^P$  under the given local context  $\Gamma$  and context of inductive data-type declarations  $\Delta$ .

## 5.6 Generation of Abstracted Inductive Definitions

Cedilleum supports *histomorphic* recursion (that is, having access to all previous recursive values) where termination is ensured through typing. In order to make this possible, we need a mechanism for tracking the global definitions of *concrete* inductive data types as well the locally-introduced *abstract* inductive data type representing the recursive occurrences suitable for a fixpoint function to be called on.

If  $I$  is an inductive type such that  $\Delta(I) = \text{Ind}_C[p](\Gamma_I := \Sigma)$  and  $I'$  is a fresh type variable, then we define function  $\text{Hist}(\Delta, I, \bar{T}, I')$  producing an abstracted (well-formed) inductive definition  $\text{Ind}_A[0](\Gamma_{I'} := \Sigma')$ , where

- $\Gamma_{I'}(I') = \forall \Gamma_D. \star$  if  $\Gamma_I(I) = \forall \Gamma_P. \forall \Gamma_D. \star$  (and  $\|\Gamma_P\| = \|\bar{T}\| = p$ )

That is, the kind of  $I'$  is the same as the kind of  $I \bar{T}$

- $\Sigma' = c'_1 : \forall \Gamma_D. \prod_{\forall} \Gamma_{A'_1}. I' \Gamma_D, \dots, c'_n : \forall \Gamma_D. \prod_{\forall} \Gamma_{A'_n}. I \bar{T} \Gamma_D,$

when each of the concrete constructors  $c_i$  in  $\Sigma$  are associated with type  $\forall \Gamma_P. \forall \Gamma_D. \prod_{\forall} \Gamma_{A_i}. I \Gamma_P \Gamma_D$  and each  $\Gamma_{A'_i} = [\lambda \Gamma_P. I' / I, \bar{T} / \Gamma_P] \Gamma_{A_i}$ .

That is, trasforming the concrete constructors of the inductive datatype  $I$  to “abstracted” constructors involves replacing each recursive occurrence of  $I \Gamma_P$  with the fresh type variable  $I$ , and instantiating each of the parameters  $\Gamma_P$  with  $\bar{T}$ .

Users of Cedilleum will see “punning” of the concrete constructors  $c_i$  and abstracted constructors  $c'_i$ . In particular, when using fix-point pattern matching branch labels will be written with the constructors for the concrete inductive data-type, and the expected type of a branch given by the motive will pretty-print using the concrete constructors. In the inference rules, however, we will take more care to distinguish the abstract constructors (see Subsection 5.7).

## 5.7 Typing Rules

Figure 17: Use of an inductive datatype  $\text{Ind}_M[p](\Gamma_I := \Sigma)$

$$\begin{array}{c}
 \frac{\Gamma \vdash_{\uparrow} t : I \bar{T} \bar{s} \quad \text{WFPat}(\Gamma, \Delta, \Delta(I), \bar{T}, \mu'(t, P, t_{i=1..n}))}{\Gamma, \Delta \vdash_{\delta} \mu'(t, P, t_{i=1..n}) : P \bar{s} t} \\
 \\
 \frac{\begin{array}{l} \Gamma \vdash_{\uparrow} t : I \bar{T} \bar{s} \quad \Delta(I) = \text{Ind}_C[p](\Gamma_I := \Sigma) \quad \Gamma_I(I) = \prod \Gamma_P. \prod \Gamma_D. \star, \|\Gamma_P\| = p \quad \text{Hist}(\Delta, I, \bar{T}, I') = \text{Ind}_A[0](\Gamma_{I'} := \Sigma') \\ \Gamma' = \Gamma, \Gamma_{I'}, x_{\text{to}} = \lambda \Gamma_D. \lambda x. x : \forall \Gamma_D. I' \Gamma_D \rightarrow I \bar{T} \Gamma_D, x_{\text{rec}} : \forall \Gamma_D. \prod x : I' \Gamma_D. P \Gamma_D (x_{\text{to}} \Gamma_D x) \quad \Delta' = \Delta, \text{Hist}(\Delta, I, \bar{T}, I') \end{array}}{\frac{\text{WFPat}(\Gamma', \Delta', \Delta'(I'), \emptyset, \mu'(t, P, t_{i=1..n}))}{\Gamma, \Delta \vdash_{\delta} \mu(x_{\text{rec}}, I', x_{\text{to}}, t, P, t_{i=1..n}) : P \bar{s} t}}
 \end{array}$$

The first rule of Figure 17 is for typing simple pattern matching with  $\mu'$ . We need to know that the scrutinee  $t$  is well-typed at some inductive type  $I \bar{T} \bar{s}$ , where  $\bar{T}$  represents the parameters and  $\bar{s}$  the indices. Then we defer to the judgment  $\text{WFPat}$  to ensure that this pattern-matching expression is a valid elimination of  $t$  to type  $P$ .

The second rule is for typing pattern-matching with fix-points, and is significantly more involved. As above we check the scrutinee  $t$  has some inductive type  $I \bar{T} \bar{s}$ . We confirm that  $I$  is a *concrete* inductive data-type by looking up its definition in  $\Delta$ , and then generate the abstracted definition  $\text{Hist}(\Delta, I, \bar{T}, I')$  for some fresh  $I'$ . We then add to the local typing context  $\Gamma_{I'}$  (the new inductive type  $I'$  with its associated kind) and two new variables  $x_{\text{to}}$  and  $x_{\text{rec}}$ .

- $x_{\text{to}}$  is the *revealer*. It casts a term of an abstracted inductive data-type  $I' \Gamma_D$  to the concrete type  $I \bar{T} \Gamma_D$ . Crucially, it is an *identity* cast (the implicit quantification  $\Lambda \Gamma_D$  disappears after erasure). The intuition why this should be the case is that the abstracted type  $I'$  only serves to mark the recursive occurrences of  $I$  during pattern-matching to guarantee termination.
- $x_{\text{rec}}$  is the *recursor* (or the inductive hypothesis). Its result type  $P' \Gamma_D x$  utilizes  $x_{\text{to}}$  in  $P'$  to be well-typed, as the  $x$  in this expression has type  $I' \Gamma_D$ , but  $P$  expects an  $I \bar{T} \Gamma_D$ . Because  $x_{\text{to}}$  erases to the identity, uses of the  $x_{\text{rec}}$  will produce expressions whose types will not interfere with producing the needed result for a given branch (see the extended example – TODO).

With these definitions, we finish the rule by checking that the pattern is well-formed using the augmented local context  $\Gamma'$  and context of inductive data-type definitions  $\Delta'$ .

## References

- [Inr18] Inria. The Coq Documentation. <https://coq.inria.fr/refman/index.html>, 2018.
- [PM15] Christine Paulin-Mohring. Introduction to the calculus of inductive constructions, 2015.